



UNIVERSITAT POLITÈCNICA DE CATALUNYA  
BARCELONATECH

Facultat d'Informàtica de Barcelona



UNIVERSITAT POLITÈCNICA DE CATALUNYA  
(UPC) - BARCELONATECH

FACULTAT D'INFORMÀTICA DE BARCELONA (FIB)

MASTER IN INNOVATION AND RESEARCH IN  
INFORMATICS (MIRI)

HIGH PERFORMANCE COMPUTING (HPC)

---

# Enabling Analytic and HPC Workflows with COMPSs

---

FINAL MASTER THESIS (FMT)

2016-2017 | AUTUMN SEMESTER

*Author:*

Cristián RAMÓN-CORTÉS  
VILARRODONA  
(cristian.ramoncortes@bsc.es)

*Supervisor:*

Dra. Rosa M. BADIA SALA  
(rosa.m.badia@bsc.es)

COMPUTER ARCHITECTURE DEPARTMENT (DAC)



**Barcelona  
Supercomputing  
Center**

*Centro Nacional de Supercomputación*



"En esto, descubrieron treinta o cuarenta molinos de viento que hay en aquel campo, y así como don Quijote los vio, dijo a su escudero:

-La ventura va guiando nuestras cosas mejor de lo que acertáramos a desear; porque ves allí, amigo Sancho Panza, donde se descubren treinta, o pocos más, desaforados gigantes, con quien pienso hacer batalla y quitarles a todos las vidas, con cuyos despojos comenzaremos a enriquecer; que ésta es buena guerra, y es gran servicio de Dios quitar tan mala simiente de sobre la faz de la tierra.

-¿Qué gigantes? -dijo Sancho Panza.

-Aquéllos que allí ves -respondió su amo- de los brazos largos, que los suelen tener algunos de casi dos leguas.

-Mire vuestra merced -respondió Sancho- que aquéllos que allí se parecen no son gigantes, sino molinos de viento, y lo que en ellos parecen brazos son las aspas, que, volteadas del viento, hacen andar la piedra del molino.

-Bien parece -respondió don Quijote- que no estás cursado en esto de las aventuras: ellos son gigantes...que yo voy a entrar con ellos en fiera y desigual batalla.

Y diciendo esto, dio de espuelas a su caballo Rocinante, sin atender a las voces que su escudero Sancho le daba, advirtiéndole que, sin duda alguna, eran molinos de viento, y no gigantes..."

---

**Miguel de Cervantes Saavedra,**  
*Don Quijote de la Mancha*



## Dedication

Facing a challenging work needs self-efforts as well as the patience of the people around us, especially from our peers.

To my loving mother and father, Dolors and Joan, whose love, encouragement and gentle prodding guided me to such a success. I hope that this work will complete the dream they had for me many years ago when they chose to give me the best education they could.

Special thanks to my sweet sister, Marta, whose affection and support keeps me always up.

Last but not least, I cannot forget Laura, who did more than her share around the house as I was locked in the computer room. Without her unconditional love and constant encouragement, this would not have been possible.

Wholeheartedly,  
Cristián Ramón-Cortés Vilarrodona



## Declaration of Authorship

I hereby declare that, except where specific reference is made to the work of others, this Master's thesis has been composed by me and it is based on my own work. None of the contents of this dissertation have been previously published nor submitted, in whole or in part, to any other examination in this or any other university.

Signed:

---

Date:

---





## Acknowledgements

I gratefully thank my supervisor Rosa M. Badia Sala for all her assistance during my career at the *Barcelona Supercomputing Center (BSC-CNS)* and for giving me the opportunity to collaborate on this project.

I would also like to thank all my colleagues, current and former members of the *Workflows and Distributed Computing* team from the *Barcelona Supercomputing Center (BSC)* for their useful comments, remarks and engagement through the learning process of this Master's thesis: Jorge Ejarque, Francesc Lordan, Francisco Javier Conejero, Raul Sirvent, Daniele Lezzi, Pol Alvarez, Ramon Amela, Sandra Corella, Albert Serven, Adrià Aguilà and Sergio Rodríguez.

Special thanks to Kim Serradell Maronda for giving me the opportunity to work with the *NMMB* application and guiding me through its internals.



UNIVERSITAT POLITÈCNICA DE CATALUNYA (UPC) - BARCELONATECH

Facultat d'Informàtica de Barcelona (FIB)

MASTER IN INNOVATION AND RESEARCH IN INFORMATICS (MIRI)

High Performance Computing (HPC)

## *Abstract*

### **Enabling Analytic and HPC Workflows with COMPSs**

by Cristián RAMÓN-CORTÉS VILARRODONA

In the recent joint venture between High-Performance Computing (HPC) and Big-Data (BD) Ecosystems towards the Exascale Computing, the scientific community has realized that powerful programming models and high-level abstraction tools are a must. Within this context, the Barcelona Supercomputing Center (BSC) is developing the COMP Superscalar (COMPSs) programming model, whose main objective is to develop applications in a sequential way, while the Runtime System handles the inherent parallelism of the application and abstracts the programmer from the different underlying infrastructures. The parallelism is achieved by defining an application Interface that allows COMPSs to detect methods that operate on a set of parameters (called tasks), and execute them distributedly and transparently.

This Master Thesis aims to enhance COMPSs, adapting it to the needs of the Big-Data Ecosystems, by supporting Analytic and HPC workflows. To this end, we propose a straightforward integration with the execution of binaries, and *MPI* and *OmpSs* applications. Although the COMPSs programming model is kept untouched, we extend the COMPSs Annotations and some of the COMPSs internals such as the task schedulers and the worker executors.

To support our contribution, we have ported to COMPSs two real use cases. On the one hand, *NMMB BSC-Dust*, a workflow to predict the atmospheric life cycle of the desert dust and, on the other hand, *Guidance*, an integrated solution for Genome and Phenome association analysis.

**Keywords:** HPC, Distributed Computing, Workflows, COMPSs, PyCOMPSs



# Contents

<b>Dedication</b>	<b>v</b>
<b>Declaration of Authorship</b>	<b>vii</b>
<b>Acknowledgements</b>	<b>ix</b>
<b>Abstract</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Context . . . . .	2
1.3 Objectives . . . . .	2
1.3.1 Detailed Objectives . . . . .	2
1.4 Document Structure . . . . .	2
<b>2 State of the art</b>	<b>5</b>
2.1 Distributed libraries . . . . .	5
2.1.1 MPI . . . . .	5
2.1.2 Sockets . . . . .	8
2.2 Workflow Frameworks . . . . .	11
2.2.1 Frameworks with explicit workflows' definition . . . . .	11
2.2.1.1 Taverna . . . . .	11
2.2.1.2 Fireworks . . . . .	12
2.2.1.3 Kepler . . . . .	14
2.2.1.4 Galaxy . . . . .	15
2.2.2 Frameworks with implicit workflows' definition . . . . .	16
2.2.2.1 MapReduce . . . . .	16
2.2.2.2 Spark . . . . .	18
2.2.2.3 Swift . . . . .	20
2.2.2.4 COMP Superscalar (COMPSs) . . . . .	22
<b>3 COMPSs overview</b>	<b>23</b>
3.1 Programming Model . . . . .	24
3.2 Runtime System . . . . .	27
3.3 Task Workflow . . . . .	29
<b>4 Tools and methodology</b>	<b>31</b>
4.1 Tools . . . . .	31
4.2 Methodology . . . . .	31
4.2.1 Scientific method design . . . . .	31
4.2.2 Development strategy . . . . .	32
4.2.3 Validation strategy . . . . .	32

<b>5</b>	<b>Implementation</b>	<b>33</b>
5.1	Programming model annotations . . . . .	33
5.1.1	New task annotations . . . . .	33
5.1.2	Environment variables as annotations . . . . .	35
5.1.3	Versioning task annotation . . . . .	38
5.1.4	SchedulerHints task annotation . . . . .	40
5.1.5	New stream parameter annotation . . . . .	41
5.1.6	New prefix parameter annotation . . . . .	42
5.2	Scheduling modifications . . . . .	45
5.2.1	Treatment of non-native tasks . . . . .	45
5.2.2	Multi-node execution actions . . . . .	45
5.2.3	Treatment of SchedulerHints . . . . .	48
5.3	Worker enhancements . . . . .	48
5.3.1	Invokers . . . . .	49
5.3.2	External executors enhancement . . . . .	50
<b>6</b>	<b>Results and evaluation</b>	<b>53</b>
6.1	Proofs of concept . . . . .	53
6.1.1	BLAST . . . . .	53
6.1.1.1	Application description . . . . .	53
6.1.1.2	Purpose . . . . .	54
6.1.1.3	Evaluation . . . . .	54
6.1.2	Matmul . . . . .	58
6.1.2.1	Application description . . . . .	58
6.1.2.2	Purpose . . . . .	59
6.1.2.3	Hybrid COMPSs and MPI Matmul implementation . . . . .	59
6.1.2.4	Evaluation . . . . .	61
6.2	Use cases . . . . .	64
6.2.1	NMMB/BSC-Dust . . . . .	64
6.2.1.1	Application description . . . . .	64
6.2.1.2	Purpose . . . . .	64
6.2.1.3	NMMB/BSC-Dust implementation with COMPSs . . . . .	65
6.2.1.4	Evaluation . . . . .	66
6.2.2	GUIDANCE . . . . .	70
6.2.2.1	Application description . . . . .	70
6.2.2.2	Purpose . . . . .	71
6.2.2.3	GUIDANCE implementation with COMPSs . . . . .	71
6.2.2.4	Evaluation . . . . .	73
<b>7</b>	<b>Conclusions and Future work</b>	<b>75</b>
	<b>Bibliography</b>	<b>77</b>
	<b>Appendices</b>	<b>81</b>
<b>A</b>	<b>Blast: complete code</b>	<b>83</b>
A.1	Blast.java . . . . .	83
A.2	BlastItf.java . . . . .	86
A.3	BlastImpl.java . . . . .	87
A.4	BINARY.java . . . . .	88

<b>B Matmul: complete code</b>	<b>89</b>
B.1 Matmul.java . . . . .	89
B.2 MatmulItf.java . . . . .	91
B.3 MatmulImpl.java . . . . .	92
B.4 MPI.java . . . . .	93
B.5 Matmul.c . . . . .	93
B.6 Block.java . . . . .	95
<b>C NMMB/BSC-Dust: code highlights</b>	<b>99</b>
C.1 Nmmb.java . . . . .	99
C.2 NmmbItf.java . . . . .	112
<b>D GUIDANCE: code highlights</b>	<b>121</b>
D.1 Guidance.java . . . . .	121
D.2 GuidanceItf.java . . . . .	153





# List of Figures

2.1	MPI Hello example in C . . . . .	6
2.2	MPI Hello execution example . . . . .	7
2.3	MPI Hello diagram of execution . . . . .	7
2.4	Main code of Java Hello example with Sockets . . . . .	8
2.5	Master process code of Java Hello example with Sockets . . . . .	9
2.6	Slave process code of Java Hello example with Sockets . . . . .	10
2.7	Execution example of Hello with Sockets . . . . .	11
2.8	BLAST design example using Taverna . . . . .	12
2.9	FireWorks components . . . . .	12
2.10	Workflows' components in FireWorks . . . . .	13
2.11	Example of Workflow using FireWorks . . . . .	13
2.12	Lotka-Volterra workflow example using Kepler . . . . .	14
2.13	Galaxy graphical web-based platform to define Workflows . . . . .	15
2.14	Galaxy graphical web-based platform to execute Workflows . . . . .	16
2.15	Wordcount example on top of <i>Hadoop</i> . . . . .	17
2.16	Execution example of Wordcount using MapReduce . . . . .	18
2.17	Spark's Components . . . . .	19
2.18	Wordcount example in Java using Spark . . . . .	20
2.19	Swift programming language . . . . .	21
2.20	Swift simulation workflow example . . . . .	21
2.21	Swift simulation code example . . . . .	22
3.1	COMPSs overview . . . . .	23
3.2	Increment main class . . . . .	24
3.3	Increment helper methods class . . . . .	25
3.4	Increment Interface . . . . .	25
3.5	Sequential execution example of Increment . . . . .	26
3.6	COMPSs execution example of Increment . . . . .	26
3.7	COMPSs structure . . . . .	27
3.8	COMPSs Runtime overview . . . . .	28
3.9	COMPSs task execution workflow . . . . .	29
5.1	Binary annotation . . . . .	34
5.2	Complete Binary annotation . . . . .	34
5.3	MPI annotation . . . . .	34
5.4	Complete MPI annotation . . . . .	34
5.5	OmpSs annotation . . . . .	35
5.6	Complete OmpSs annotation . . . . .	35
5.7	Wordcount Interface . . . . .	35
5.8	Wordcount executions with different constraint values . . . . .	36
5.9	Wordcount Interface with environment variables . . . . .	36
5.10	Wordcount executions with environment variables as constraints . . . . .	36
5.11	Example with an Interface with all the available environment variables . . . . .	37

5.12	Example of complex environment variables on the <i>workingDir</i> field . . . . .	38
5.13	Example of previous versioning main code . . . . .	38
5.14	Example of previous versioning implementation 1 code . . . . .	38
5.15	Example of previous versioning implementation 2 code . . . . .	39
5.16	Example of the Interface of previous versioning . . . . .	39
5.17	Example of the Interface of previous versioning with constraints . . . . .	39
5.18	Example of the new Annotation Interface . . . . .	40
5.19	Extended example of the new Annotation Interface . . . . .	40
5.20	Example of an Interface with SchedulerHints . . . . .	41
5.21	Example of the different return types of the non-native tasks . . . . .	41
5.22	Example of the different stream annotations for non-native tasks . . . . .	42
5.23	Binary Tasks example for joint prefixes . . . . .	43
5.24	Main code example for joint prefixes . . . . .	43
5.25	Interface example of an application with prefixes . . . . .	44
5.26	Example of the main code calls to tasks with prefixes . . . . .	44
5.27	Example of the command executed inside each task using prefixes . . . . .	45
5.28	Example of a single node task flow . . . . .	46
5.29	Example of a multi-node task flow . . . . .	47
5.30	New structure of the COMPSs Worker Executors . . . . .	50
5.31	Execution time versus number of <i>ProcessBuilders</i> or Pipes . . . . .	51
6.1	Execution arguments of the COMPSs BLAST application . . . . .	53
6.2	Example of BLAST execution with $N = 8$ . . . . .	54
6.3	COMPSs BLAST application: new <i>align</i> task implementation . . . . .	54
6.4	COMPSs BLAST application: old <i>align</i> task implementation . . . . .	55
6.5	COMPSs BLAST application: new <i>align</i> task call . . . . .	55
6.6	COMPSs BLAST application: old <i>align</i> task call . . . . .	56
6.7	COMPSs BLAST application: new <i>align</i> 's interface annotation . . . . .	56
6.8	COMPSs BLAST application: old <i>align</i> 's interface annotation . . . . .	56
6.9	Multiplication of a Matrix divided in blocks . . . . .	58
6.10	Task execution graph of a <i>Matmul</i> example . . . . .	59
6.11	Hybrid COMPSs and MPI block layers . . . . .	60
6.12	Main multiplication loop of the Hybrid <i>Matmul</i> . . . . .	60
6.13	<i>multiplyAccumulative</i> 's interface annotation for the Hybrid COMPSs and MPI <i>Matmul</i> . . . . .	61
6.14	<i>Matmul</i> Strong Scaling analysis . . . . .	62
6.15	<i>Matmul</i> Weak Scaling analysis . . . . .	62
6.16	Example of four hour average AOD from NMMB/BSC-Dust . . . . .	64
6.17	NMMB/BSC-Dust step workflow . . . . .	65
6.18	NCVIEW plots of FIS 3D variable for both implementations . . . . .	66
6.19	NCVIEW plot of PS 3D variable for both implementations . . . . .	66
6.20	NCVIEW plot of SLP 3D variable for both implementations . . . . .	67
6.21	Tasks graph of NMMB/BSC-Dust with COMPSs . . . . .	68
6.22	Paraver task view of the NMMB/BSC-Dust execution with COMPSs . . . . .	68
6.23	GUIDANCE's schematic representation of the typical complete Genome and Phenome association analysis . . . . .	70
6.24	Example of previous GUIDANCE binary task implementation . . . . .	71
6.25	GUIDANCE's partial task graph . . . . .	72

# List of Tables

3.1	Useful arguments for the <i>runcompss</i> command . . . . .	27
5.1	Available fields for the <i>SchedulerHints</i> annotation . . . . .	40
5.2	Available stream types with their valid directions and execution behaviour .	42
6.1	Execution parameters of the BLAST application . . . . .	57
6.2	Task times . . . . .	57
6.3	NMMB/BSC-Dust code summary . . . . .	65
6.4	Execution times of the different NMMB/BSC-Dust implementations for the simulation of 1 day of global domain with 64 cores . . . . .	67
6.5	GUIDANCE code summary . . . . .	73



# List of Abbreviations

<b>API</b>	<b>Application Programming Interface</b>
<b>COMPSs</b>	<b>COMP Superscalar</b>
<b>CPU</b>	<b>Central Processing Unit</b>
<b>GUI</b>	<b>Graphical User Interface</b>
<b>I/O</b>	<b>Input / Output</b>
<b>MPI</b>	<b>Message Passing Interface</b>
<b>NEMS</b>	<b>NOAA Environmental Modeling System</b>
<b>NMMB</b>	<b>Nonhydrostatic Multiscale Model on the B-grid</b>
<b>NOAA</b>	<b>National Oceanic and Atmospheric Administration</b>
<b>OS</b>	<b>Operating System</b>
<b>PyCOMPSs</b>	<b>Python binding for COMP Superscalar</b>
<b>RAM</b>	<b>Random Access Memory</b>
<b>SSH</b>	<b>Secure SHell</b>
<b>VM</b>	<b>Virtual Machine</b>
<b>WSDL</b>	<b>Web Services Description Language</b>



# Glossary

## **CPU**

The Central Processing Unit (CPU) is the part of the computer that contains all the elements required to execute the instructions of software programs. Its main components are the main memory, the Processing Unit (PU) and the Control Unit (CU). Modern computers use multi-core processors, which are a single chip containing one or more cores.

## **Core**

A core is an individual processor that actually executes program instructions. Current single chip CPUs contain many cores and are referred as multi-processor or multi-cores.

## **MareNostrum III**

MareNostrum III is the most powerful supercomputer in Spain with 3056 nodes (48.896 processors), 115.5 TB of main memory and a peak performance of 1.1 Petaflops. It is hosted by the Barcelona Supercomputing Center (BSC).

## **Node**

A compute node refers to a single system within a cluster of many systems.

## **Scratch Space**

Supercomputers generally have what is called scratch space: disk space available for temporary use and only accessible from the compute nodes.

## **SSH**

A protocol to securely connect to a remote computer. This connection is generally for a command line interface, but it is possible to use GUI programs through SSH.

## **Environment Variable**

In Linux systems, each process has an execution environment. This environment can be inherited from the user session environment and can be extended with specific process variables. An Environment Variable is a value stored in the process environment that can affect its execution.

## **API**

An Application Programming Interface (API) is a set of methods and functions that are used by another software to produce an abstraction layer.

## **Operating System**

A system software that manages the hardware and provides services for computer software.

**Framework**

Framework stands for a set of standardized concepts, practices or criterias used to face a given problem. Specifically, it defines a set of programs, libraries, languages, and programming models used jointly in a project.

**Workflow**

A workflow is composed of tasks and dependencies between tasks. Workflows are commonly represented as graphs, with the nodes beeing tasks and the arrows representing the dependencies. Somehow, tasks must represent an action that must be done (i.e. the execution of a binary), and the dependencies must represent the requirements that must be satisfied to be able to execute the task (i.e. the machine availability or the required data).

**Binary**

A file containing a list of machine code instructions to perform a list of tasks.

**WSDL**

Web Services Description Language (WSDL) is an Extensible Markup Language (XML) used to describe web services.

**Graphical User Interface**

The GUI is a software that graphically interacts with the user of a computer to ease the data manipulation.



# Chapter 1

## Introduction

### 1.1 Motivation

Several years ago the industry required the research community to shift from sequential computing to parallel computing. The extreme increase in the computational loads and the decrease of the acceptable system's response time forced the community to "think in parallel". Although writing parallel codes is not as easy as writing sequential programs, we can assume that the actual computing resources of a multi-core processor are already able to exploit the inherent instruction-level parallelism of an application in a completely transparent way for programmers.

However, the next generation of applications were requiring more high-performance computing resources than those that a single computing resource could offer. Thus, it was necessary to use many machines communicated through networks to achieve a larger amount of computing capabilities. That is what we know today as *Distributed Computing*. Distributed applications exploit the *task-level* parallelism and are even more difficult to handle for programmers. One of the major issues that arise from both parallel and distributed programming is that writing in parallel is not as easy as writing sequential programs and, more often than expected, people that can develop useful and complete end-user applications is not capable of writing efficient parallel code (and the other way around). We can somehow consider that scientists interpreting results do not care about the computational load or distribution (how results are computed), but care of the results quality, the time spent to retrieve the results and the robustness of the system. In this line of "*writing programs that scale with increasing number of cores should be as easy as writing programs for sequential computers*" [3] several programming models have arisen to help the programmer handle the underlying infrastructure.

Nowadays, the scientific community not only wants to work in parallel and distributed systems but also needs to handle a large amount of data. In this sense, *Big-Data (BD) Ecosystems* appeared some years ago; allowing the community to store, check, retrieve and transform enormous amounts of data in acceptable response times. For this purpose, several programming models have also arisen but are completely different to the ones used by the *High Performance Computing* community.

In the race towards the Exascale Computing [2], the scientific community has realized that unifying *High Performance Computing (HPC)* platforms and *Big-Data (BD) Ecosystems* is a **must**. Nowadays, these two ecosystems differ significantly at hardware and software level but "*programming models and tools are perhaps the biggest point of divergence between the scientific-computing and big-data ecosystems*" [50]. In this respect, "*there is a clear desire on the part of a number of domain and computer scientists to see a convergence between the two high-level types of computing systems, software, and applications: Big Data Analytics (BDA) and High Performance*

*Computing (HPC)*" [22]. We can then conclude that the joint venture between HPC and BDA ecosystems opens a great opportunity for programming models that is worth to investigate.

## 1.2 Context

The current project is conducted as the *Final Master Thesis* in the *Master of Innovation and Research in Informatics - High Performance Computing (MIRI - HPC)* offered by the *Universitat Politècnica de Catalunya (UPC)* [16] and has been funded by the *Barcelona Supercomputing Center (BSC)* [7].

The project has been developed as a research support engineer in the *Workflows and Distributed Computing* group of the *Computer Science* department at the *BSC*. The main goal of this group is to ease the development of applications and services for distributed infrastructures (such as Clusters, Grids and Clouds) through the *COMP Superscalar (COMPSs)* [19] programming model. The group also covers Big-Data and HPC applications, energy constrained environments, GPU offloading and execution in mobile devices.

## 1.3 Objectives

To take profit of the abovementioned joint venture between HPC and BDA ecosystems this Final Master Thesis **enhances the COMPSs Framework to support Analytic and HPC workflows**. In next subsection we provide detailed information about the specific objectives that have been considered to successfully achieve this contribution.

### 1.3.1 Detailed Objectives

The following points summarize the main goals of the project:

**O1** Conduct a survey among the Analytic and HPC workflows to determine the most common used technologies, compare them to the solutions proposed by other state of the art Workflow Managers and determine the most adequate and feasible features for COMPSs.

**O2** Extend the current COMPSs syntax to support the selected features in a non-invasive way; keeping the backward compatibility. More specifically, extend the COMPSs annotations to support binary, MPI and OmpSs tasks.

**O3** Enhance the current COMPSs Runtime by implementing the required mechanisms to schedule and execute the new type of tasks. At scheduling level, this objective includes the treatment of multi-node tasks and the inclusion of scheduler hints for a certain type of tasks. At execution level, the COMPSs workers must transparently execute the different kind of tasks.

**O4** Demonstrate that the proposal is able to execute Analytic and HPC Workflows by executing and analysing real use-cases.

## 1.4 Document Structure

The rest of the document is organized as follows. Chapter 2 gives an overview of the current state of the art of HPC and Analytic workflows, describing the most used technologies and comparing them to what is supported in the most common Workflow Managers. Chapter 3 introduces the COMPSs programming model, describing its main features before our

---

contribution. Chapter 4 introduces describes the tools and the development methodology followed during all the project. Chapter 5 accurately describes our contribution to provide the reader a closer look of the solutions that we have chosen to face each specific challenge. Chapter 6 evaluates the obtained results and the performance of our proposal. Chapter 7 provides a brief summary of the thesis and, finally, in Chapter 8 we discuss the conclusions and state the guidelines for the future work.



## Chapter 2

# State of the art

Scientific programs from different fields (such as bioinformatics, biomechanics, earth sciences or engineering simulations) that are either computationally intensive or require a huge amount of storage capabilities are known as *e-Science* applications. In the past decade, the computational and storage requirements of these *e-Science* applications have grown enough to do not fit in a single machine. This fact has forced *e-Science* applications to move from parallel computing to distributed computing because the computing resources within a multi-core processor are limited by the architecture and the technology of the system and were already fully exploited. However, in the distributed computing paradigm, the only room for improvements is to exploit the *task-based* parallelism between nodes.

Considering that writing applications for parallel or distributed environments is much more difficult than writing sequential programs due to the concurrency issues, the IT community has developed several libraries and frameworks to ease the development of applications (and consequently, the development of *e-Science* applications). In an attempt to classify all the tools provided by the community that can be used for the design of *e-Science* applications, we have divided them into two groups: *Distributed Libraries* and *Workflow Frameworks*. Next, Section 2.1 defines and provides examples of some state of the art *Distributed Libraries* and, Section, 2.2 defines and introducing some state of the art *Workflow Frameworks*.

### 2.1 Distributed libraries

On the one hand, distributed computing can be achieved by programming the low-level communication between processes. Although using distributed libraries can lead to excellent performance results, it is a tedious work for the programmer to develop efficient code for a complex application using such libraries. Next subsections provide some examples of these libraries.

#### 2.1.1 MPI

The Message-Passing Interface (MPI) standard “*includes point-to-point message-passing, collective communications, group and communicator concepts, process topologies, environmental management, process creation and management, one-sided communications, extended collective operations, external interfaces, I/O, some miscellaneous topics, and a profiling interface*” [38]. The MPI standard includes bindings for C and Fortran, and its goal is “*to develop a widely used standard for writing message-passing programs*”[38]. It has several different implementations but the most well-known are *OpenMPI*[54] and *IMPI*[30]).

For the end user, using MPI requires handling explicitly the spawn of the processes, the code executed by each process and the communication between them. All this management is done by using API calls, and thus, the application code must be compiled and executed with the MPI compiler of the specific MPI implementation. The main advantage is that the users have full control of all the processes and the communication between them which, for

experienced users, leads to high efficient codes. However, for inexperienced users, an efficient code can become unreadable and handle many processes can become a tedious work. Moreover, when porting a sequential application to an MPI application, the users must explicitly distribute the data between the processes and retrieve back the results (which can lead to load imbalance or inefficient communications). Additionally, another inconvenience of MPI is that, once the application's execution has started, the number of processes cannot be changed dynamically, limiting the maleability of the applications.

Figure 2.1 shows an example of an MPI application written in C. The code spawns a given number of processes, sets up one process as coordinator and the rest as slaves that send back a message to the coordinator saying that they are ready to work. As seen in the figure, each process has a unique identifier, and the communication between them is done by using the MPI\_Send or MPI\_Receive API calls (obviously, more complex programs will require more complex API calls).

```

1 #include <assert.h>
2 #include <stdio.h>
3 #include <string.h>
4 #include <mpi.h>
5
6 int main(int argc, char **argv) {
7     char buf[256];
8     int my_rank, num_procs;
9     /* Initialize the infrastructure necessary for communication */
10    MPI_Init(&argc, &argv);
11    /* Identify this process */
12    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
13    /* Find out how many total processes are active */
14    MPI_Comm_size(MPI_COMM_WORLD, &num_procs);
15    /* Until this point, all programs have been doing exactly the same.
16       Here, we check the rank to distinguish the roles of the programs */
17    if (my_rank == 0) {
18        int other_rank;
19        printf("We have %i processes.\n", num_procs);
20        /* Send messages to all other processes */
21        for (other_rank = 1; other_rank < num_procs; other_rank++) {
22            sprintf(buf, "Hello %i!", other_rank);
23            MPI_Send(buf, sizeof(buf), MPI_CHAR, other_rank,
24                    0, MPI_COMM_WORLD);
25        }
26        /* Receive messages from all other process */
27        for (other_rank = 1; other_rank < num_procs; other_rank++) {
28            MPI_Recv(buf, sizeof(buf), MPI_CHAR, other_rank,
29                    0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
30            printf("%s\n", buf);
31        }
32    } else {
33        /* Receive message from process #0 */
34        MPI_Recv(buf, sizeof(buf), MPI_CHAR, 0,
35                0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
36        assert(memcmp(buf, "Hello ", 6) == 0),
37
38        /* Send message to process #0 */
39        sprintf(buf, "Process %i reporting for duty.", my_rank);
40        MPI_Send(buf, sizeof(buf), MPI_CHAR, 0,
41                0, MPI_COMM_WORLD);
42    }
43    /* Tear down the communication infrastructure */
44    MPI_Finalize();
45    return 0;
46 }

```

FIGURE 2.1: MPI Hello example in C.  
Source: Wikipedia, Message Passing Interface

Figure 2.2 shows the command line result and Figure 2.3 shows a diagram of executing the aforementioned code with 4 processes. Notice that the spawn time of the processes and the communication between them is not always done at the same time and thus, the diagram is only one of the possible execution diagrams of the same code. For instance, Process 0 will always send messages to Processes 1, 2 and 3 in the same order and will receive the messages back in the same order, but Processes 1, 2 and 3 can receive the message in different orders and send the reply in different orders. This issue is one of the hardest things to overcome when developing applications with MPI because blocking processes in a receive call can lead to significant overheads. For instance, in our diagram, Processes 2 and 3 have sent all their data, but Process 0 does not receive the message until the data from Process 1 is received.

```

1 $ mpicc example.c
2 $ mpirun -n 4 ./a.out
3 We have 4 processes.
4 Process 1 reporting for duty.
5 Process 2 reporting for duty.
6 Process 3 reporting for duty.

```

FIGURE 2.2: MPI Hello execution example

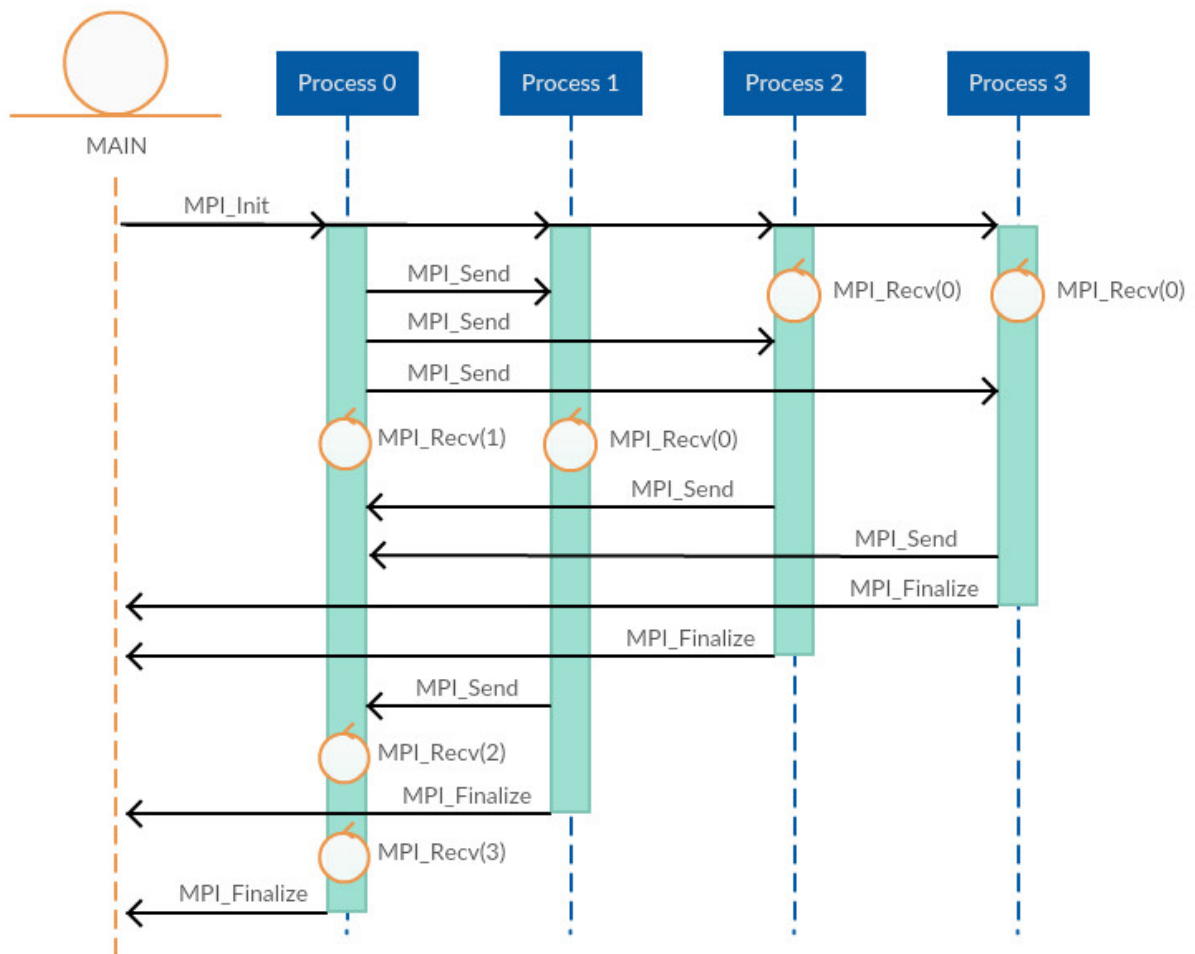


FIGURE 2.3: MPI Hello diagram of execution

## 2.1.2 Sockets

A Socket[42] is an endpoint on a machine to send and receive data and can be used either to communicate processes within the same machine or through the network. Each process refers to a socket by means of a *Socket Descriptor* and then establishes a connection through the socket creating a *Channel*. The underlying Operating System of the machine provides a socket API, but applications normally use a higher-level API implementation that depends on its language (for example, `sys/socket.h` for C [26], socket library for Python [49] or `java.net.Socket` for Java [40]).

Like the previous case, for the end user, using Sockets only enables the communication between different processes. Thus, it requires handling explicitly the spawn of the processes, the code executed by each process and the communication between them. Since the application will also use an API, the source code must be developed using this underlying technology and compiled accordingly. However, contrary to MPI, it does not require a specific compiler.

The next figures show the same example than the MPI case but written in Java and using Java Sockets. Figure 2.4 contains the main code to initialize and spawn all the processes.

```
1 package helloSocket;
2
3 public class HelloSocket {
4
5     private static final int BASE_SOCKET_PORT = 47_000;
6     private static int numThreads;
7     private static int[] socketPorts;
8
9     public static void main(String args[]) {
10         // Get arguments
11         if (args.length < 1) {
12             System.err.print("[ERROR] Invalid usage");
13             System.err.print("      Usage: ./main <numThreads>");
14             System.exit(1);
15         }
16         numThreads = Integer.valueOf(args[0]);
17         // Create the array of usable ports
18         socketPorts = new int[numThreads];
19         for (int i = 0; i < numThreads; ++i) {
20             socketPorts[i] = BASE_SOCKET_PORT + i;
21         }
22
23         Thread[] procs = new Thread[numThreads];
24         // Spawn the master process
25         procs[0] = new MasterProcess(numThreads, socketPorts);
26         procs[0].start();
27         // Spawn the slave processes
28         for (int i = 1; i < numThreads; ++i) {
29             procs[i] = new SlaveProcess(i, socketPorts[i]);
30             procs[i].start();
31         }
32
33         // Join the processes
34         for (int i = 0; i < numThreads; ++i) {
35             try {
36                 procs[i].join();
37             } catch (InterruptedException ie) {
38                 Thread.currentThread().interrupt();
39             }
40         }
41     }
42 }
```

FIGURE 2.4: Main code of Java Hello example with Sockets



Figure 2.5 contains the code for the main process (the coordinator).

```
package helloSocket;
2
import java.io.BufferedReader;
4 import java.io.BufferedWriter;
import java.io.IOException;
6 import java.io.InputStreamReader;
import java.io.OutputStreamWriter;
8 import java.net.Socket;

10
public class MasterProcess extends Thread {
12
    private static final String SERVER_MACHINE = null; // hostname
14
    private final int numThreads;
16 private final int[] socketPorts;

18
    public MasterProcess(int numThreads, int[] socketPorts) {
20         super();
        this.numThreads = numThreads;
22         this.socketPorts = socketPorts;
    }

24
    @Override
26 public void run() {
        System.out.println("We have " + this.numThreads + " processes");
28
        /* Send messages to all other processes */
30 for (int i = 1; i < this.numThreads; ++i) {
        try (Socket socketClient = new Socket(SERVER_MACHINE, this.socketPorts[i])) {
32             BufferedWriter writer = new BufferedWriter(new OutputStreamWriter(
                socketClient.getOutputStream()));
34             writer.write("Hello " + i);
            writer.flush();
36         } catch (IOException e) {
            System.err.println("[ERROR] Cannot send message to client "
38                 + this.socketPorts[i]);
            e.printStackTrace();
40         }
    }

42
    /* Receive messages from all other process */
44 for (int i = 1; i < this.numThreads; ++i) {
        try (Socket socketClient = new Socket(SERVER_MACHINE, this.socketPorts[i])) {
46             BufferedReader reader = new BufferedReader(new InputStreamReader(
                socketClient.getInputStream()));
48             String msg = null;
            while ((msg = reader.readLine()) != null) {
50                 System.out.println(msg);
            }
52         } catch (IOException e) {
            System.err.println("[ERROR] Cannot receive message from client "
54                 + this.socketPorts[i]);
            e.printStackTrace();
56         }
    }
58 }
60 }
```

FIGURE 2.5: Master process code of Java Hello example with Sockets

Figure 2.6 contains the code for the slave processes.

```

package helloSocket;
2
import java.io.BufferedReader;
import java.io.BufferedWriter;
4 import java.io.IOException;
import java.io.InputStreamReader;
import java.io.OutputStreamWriter;
6 import java.net.ServerSocket;
import java.net.Socket;
8
10
12 public class SlaveProcess extends Thread {
    private final int id;
    private final int socketPort;
14
16
18     public SlaveProcess(int id, int socketPort) {
        super();
        this.id = id;
        this.socketPort = socketPort;
20     }
22
24     @Override
    public void run() {
26         try (ServerSocket mySocket = new ServerSocket(this.socketPort)) {
            /* Receive message from process #0 */
28             try (Socket connectionSocket = mySocket.accept()) {
                BufferedReader reader = new BufferedReader(new InputStreamReader(
30                     connectionSocket.getInputStream()));
                String msg;
32                 while ((msg = reader.readLine()) != null) {
                    assert (msg.equals("Hello" + this.id));
34                 }
            } catch (IOException e) {
36                 System.err.println("[ERROR] Thread " + this.id
                    + " cannot read from process 0 socket");
38                 throw e;
            }
40
            /* Send message to process #0 */
42             try (Socket connectionSocket = mySocket.accept()) {
                BufferedWriter writer = new BufferedWriter(new OutputStreamWriter(
44                     connectionSocket.getOutputStream()));
                writer.write("Process " + this.id + " reporting for duty.");
                writer.flush();
46             } catch (IOException e) {
48                 System.err.println("[ERROR] Thread " + this.id
                    + " cannot write to process 0 socket");
50                 throw e;
            }
52         } catch (IOException e) {
            System.err.println("[ERROR] Thread " + this.id + " cannot open socket");
            e.printStackTrace();
54         }
56     }
58 }

```

FIGURE 2.6: Slave process code of Java Hello example with Sockets

Notice that, like in the previous case, since the master process reads the messages in an ordered manner, the execution output is always the same (see Figure 2.7) even if the process spawn order and the messages' arrival can vary between executions.

```
$ javac helloSocket/*
2 $ jar cf hello.jar helloSocket/
$ java -cp helloSocket.jar helloSocket.HelloSocket 4
4 We have 4 processes
Process 1 reporting for duty.
6 Process 2 reporting for duty.
Process 3 reporting for duty.
```

FIGURE 2.7: Execution example of Hello with Sockets

## 2.2 Workflow Frameworks

On the other hand, distributed computing can also be achieved by using a workflow framework. Frameworks are designed to encapsulate a programming model, a runtime or even a programming language that eases the development of distributed applications. Although all the frameworks hide the communication mechanisms between the different processes, they can be classified regarding its task definition. Thus, subsection 2.2.2 provides several examples of *frameworks with implicit workflows' definition* and subsection 2.2.1 provides several other examples of *frameworks with explicit workflows' definition*.

### 2.2.1 Frameworks with explicit workflows' definition

Frameworks with explicit workflows' definition allow the users to design the full application workflow using a receipt file or a visualization tool. The main advantage is that the users can specifically control the dependencies between the different stages and, thus, they have a clear overview of how their application is executed by the framework. However, its main advantage is also its main disadvantage, because designing complex workflows can result in a tedious work.

#### 2.2.1.1 Taverna

Taverna [44] is a Workflow Management System that includes a set of tools to design and execute scientific workflows. The Taverna Suite consists of three components: Taverna Engine (used for enacting workflows), Taverna Workbench (the desktop client application) and the Taverna Server (to execute the remote workflows). It has been recently moved to the Apache Incubator project, and it is being used in a wide variety of domains such as bioinformatics, biodiversity, chemistry, astronomy, data mining, digitalisation and image analysis.

Although Taverna has a command line interface, its most powerful tool is the visual workflow design. This tool allows the users to graphically define workflows by constructing a diagram with inputs, outputs, actions and dependencies between actions. Each action defines its input and output ports so that users can incrementally build a diagram by linking the input and output ports of the different actions or by adding static content values to the input ports. Notice that the actions can be of any type of service since Taverna provides a set of internally configured services but allows the users to add any service by providing its WSDL address. Moreover, when the workflow is finished, Taverna can validate and run it; allowing the users to check the execution status and the partial results.

Figure 2.8 shows a Taverna diagram example of a BLAST application. The workflow has four outputs generated by different actions and does not have any input since the values are obtained from static content values (i.e. program value or database value). The light blue,

green and magenta actions are services; being the green one the invocation to BLAST with program, database and query sequence.

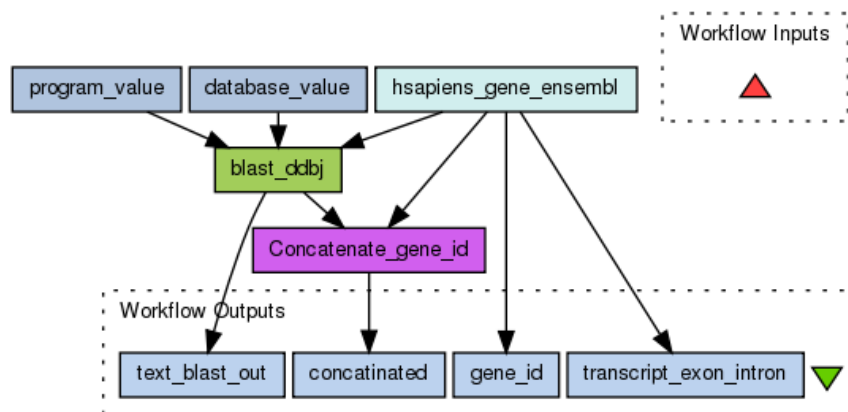


FIGURE 2.8: BLAST design example using Taverna  
Source: myExperiment[51]. Taverna workflows

As previously demonstrated, Taverna is a powerful tool to design workflows for inexperienced users since they do not have to deal directly with parallelism issues. However, this framework is only oriented to applications that can be defined as a pipeline of services and thus, its execution relies on the service availability. In this sense, the community has done a hard work in making available several services and predefined workflows that can be used as nested workflows at the end user applications.

### 2.2.1.2 Fireworks

FireWorks [32] is an open source Framework to define, manage and execute workflows. Workflows are defined using Python, JSON [46] or YAML [47], stored using MongoDB [37] and can be monitored through a web interface. The workflows' execution can be automated over arbitrary computing resources, and the framework provides fault-tolerant mechanisms and multiple execution modes (to run on different underlying infrastructures such as multi-core machines or clusters managed by queues). FireWorks includes two components (see Figure 2.9): LaunchPad (a server that manages the workflows) and one or more FireWorkers (a worker to run the jobs).

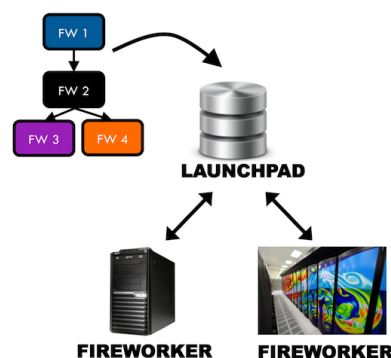


FIGURE 2.9: FireWorks components. Source: FireWorks[32]

Workflows in FireWorks have three main components (see Figure 2.10). Firstly, the users can define a FireTask which represents an atomic computing job. FireTasks must execute a single shell script or a Python function. Secondly, users can define a FireWork, which contains all the information needed to bootstrap the job execution. For example, a FireWork may contain a list of FireTasks to execute sequentially and its parameters. Finally, a Workflow is a set of Fireworks with dependencies between them.

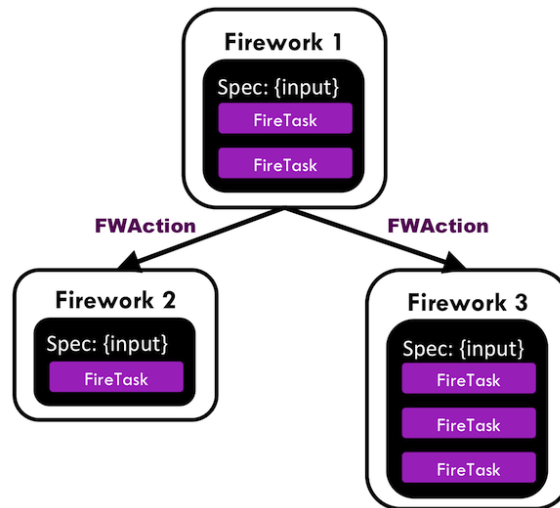


FIGURE 2.10: Workflows' components in FireWorks  
Source: FireWorks[32]

Next, Figure 2.11 shows a Python example using FireWorks. At left, we provide the code to create the workflow and, at right, a diagram of the obtained workflow. Notice that the FireTasks and FireWorks are easily created but that the Workflow must be explicitly constructed which, for complex workflows, can be difficult to understand. Thus, although the users do not need to lead with parallelism directly, they must define the workflow explicitly.

```

from fireworks import Firework, Workflow, FWorker, LaunchPad, ScriptTask
from fireworks.core.rocket_launcher import rapidfire

# set up the LaunchPad and reset it
launchpad = LaunchPad()
launchpad.reset('', require_password=False)

# define four individual FireWorks used in the Workflow
task1 = ScriptTask.from_str('echo "Ingrid is the CEO."')
task2 = ScriptTask.from_str('echo "Jill is a manager."')
task3 = ScriptTask.from_str('echo "Jack is a manager."')
task4 = ScriptTask.from_str('echo "Kip is an intern."')

fw1 = Firework(task1)
fw2 = Firework(task2)
fw3 = Firework(task3)
fw4 = Firework(task4)

# assemble Workflow from FireWorks and their connections by id
workflow = Workflow([fw1, fw2, fw3, fw4], {fw1: [fw2, fw3], fw2: [fw4], fw3: [fw4]})

# store workflow and launch it locally
launchpad.add_wf(workflow)
rapidfire(launchpad, FWorker())

```

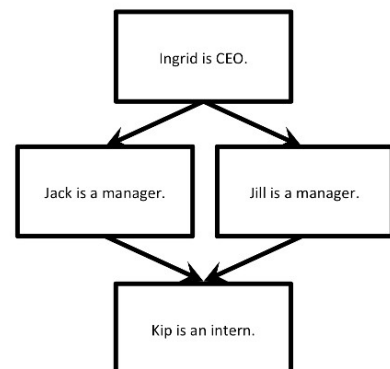


FIGURE 2.11: Example of Workflow using FireWorks  
Source: FireWorks[32]

### 2.2.1.3 Kepler

Kepler [5] is designed to create, execute and share models and analyses of scientific and engineering workflow applications. It can easily merge R [24] scripts, compiled C code or facilitate the execution of models remotely. The users use a Graphical User Interface to select and connect analytical components and data sources to define a scientific workflow.

Scientific Workflows in Kepler consist of customizable components, relations, and ports. On the one hand, the components can be either a director (to control the execution of a workflow), an actor (to execute the instructions given by the director) or a parameter (to add a configurable value). On the other hand, relations and ports are used to facilitate communication between the different components. Figure 2.12 shows a Lotka-Volterra workflow defined in the Kepler Interface where the components, the relations, and the ports are identified with callouts.

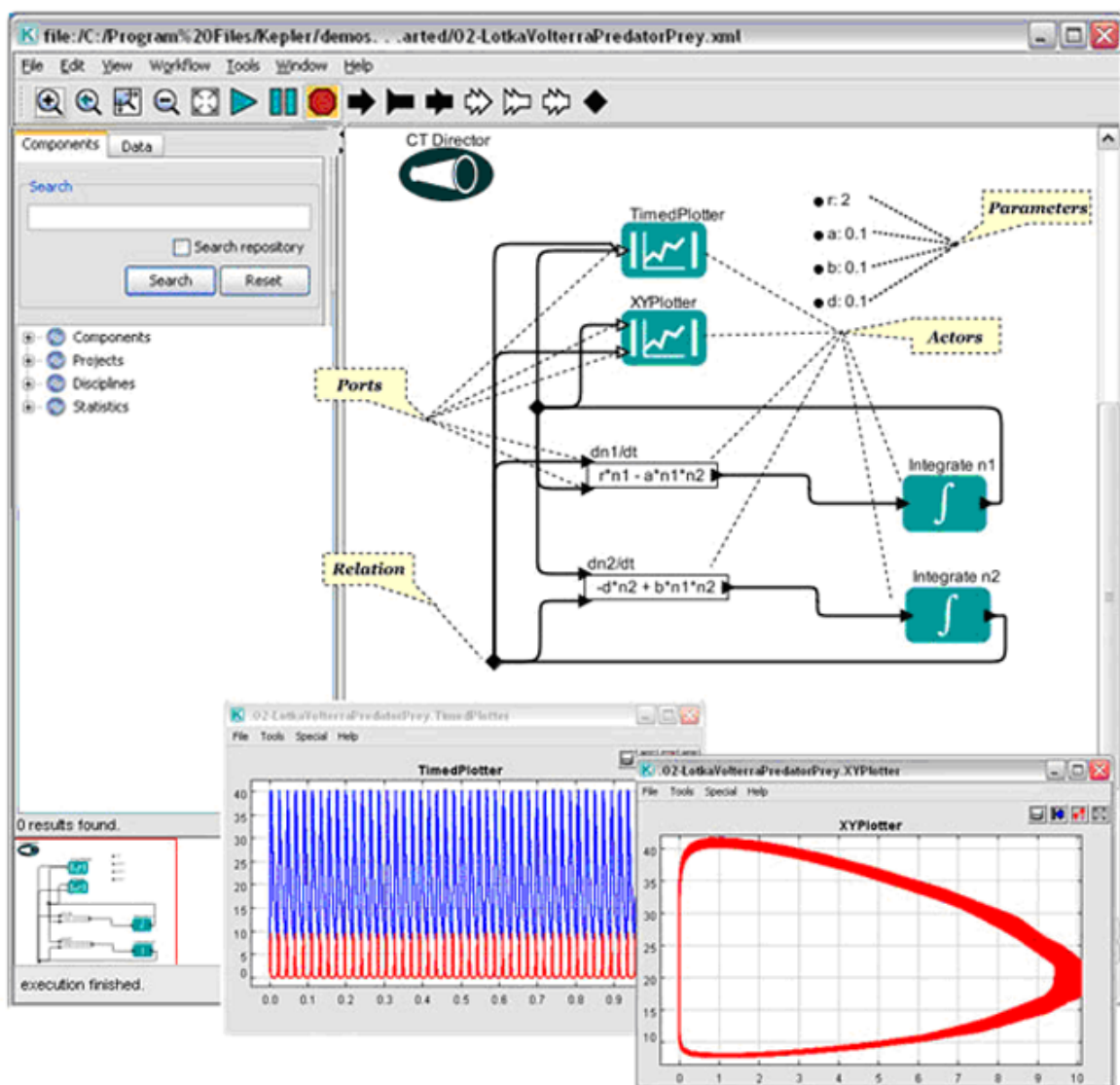


FIGURE 2.12: Lotka-Volterra workflow example using Kepler  
Source: Kepler examples

Although the Kepler Interface remembers to Taverna, this Framework does not restrict the users to execute services but rather allows them to orchestrate any type of execution.

However, like Taverna, Kepler is designed for users with little experience in parallel computing and all its efforts are intended to make the workflow design easier.

#### 2.2.1.4 Galaxy

Galaxy [28] is a web-based platform for biomedical research. It is designed for users without programming experience, providing a graphical interface to easily specify parameters and run tools and workflows. One of the greatest advantages of Galaxy is to share, publish, access and reproduce any analysis of the other users in an interactive web-based framework.

Galaxy workflows are graphically defined through its web-based platform (see Figure 2.13) by combining the execution of different tools. Several tools are provided for every user (such as retrieving data, calculating statistics or performing complex genome operations) but the users can also define custom tools. To build a workflow, the users define data dependencies between the tools' execution.

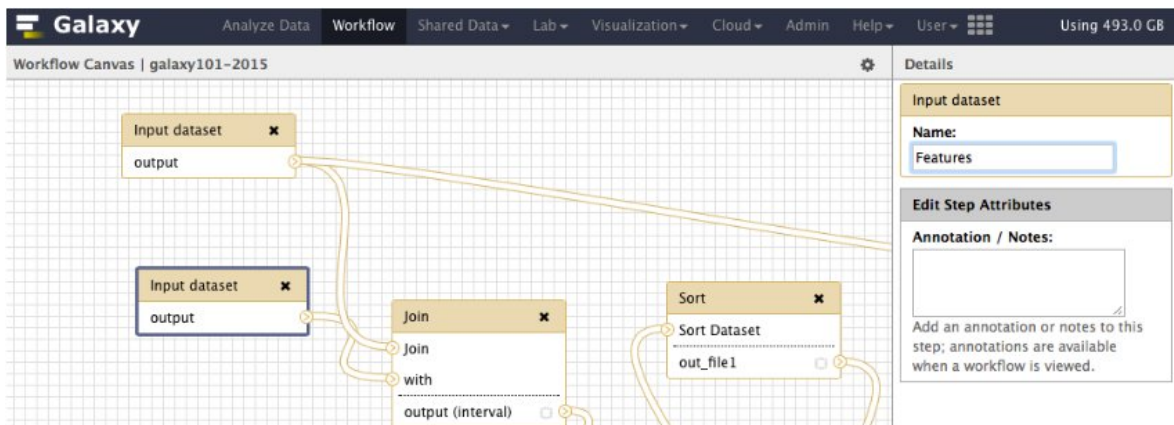


FIGURE 2.13: Galaxy graphical web-based platform to define Workflows  
Source: Galaxy Project tutorials

Galaxy also provides a section to execute the workflow. On execution time, the platform retrieves information of the execution of each tool, providing a live monitoring of the whole workflow (see Figure 2.14). Once the application has finished, the users can also analyze the output data.

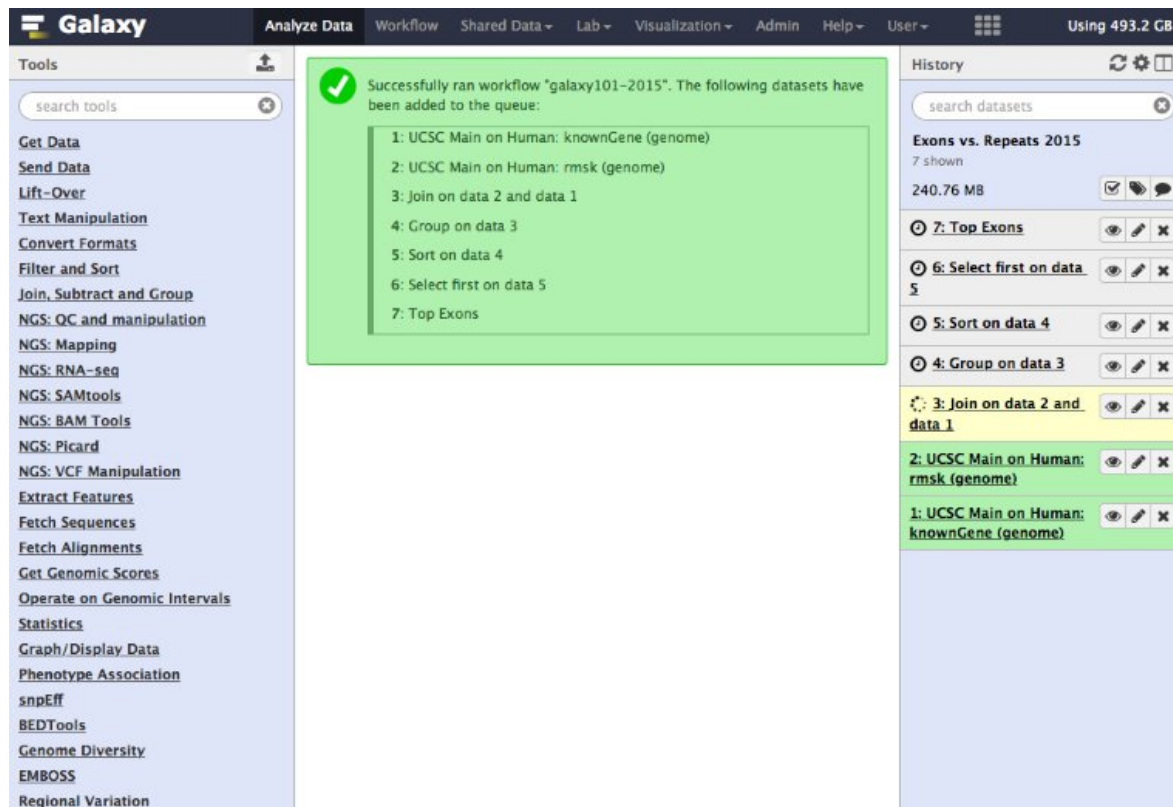


FIGURE 2.14: Galaxy graphical web-based platform to execute Workflows  
Source: Galaxy Project tutorials

Notice that Galaxy is built for specific biomedical research tools, rather than a general purpose platform. Although it provides a very clean, easy and fast way to build, execute and analyse workflows, it is only designed for biomedical applications.

## 2.2.2 Frameworks with implicit workflows' definition

Frameworks with implicit workflows' definition use different technologies to define the parallelism automatically. Users develop applications in an almost sequential manner and they do not need to handle explicitly the tasks spawned by a method call (that are distributed among the available resources). Consequently, the different frameworks differ on this "almost sequential manner" to develop the application. The main advantage of these frameworks is that the complexity of distributed programming is almost reduced to zero. However, the main disadvantage is that the users do not know beforehand how the framework will execute their application (for example, how many tasks will be created in a specific call).

### 2.2.2.1 MapReduce

MapReduce [21] is a programming model (and an underlying runtime) for processing and generating large datasets. The underlying runtime can automatically handle the parallelism; managing the inter-machine communication, taking advantage of the data locality and providing fault-tolerant mechanisms (recovering from partial failures of servers and storages by rescheduling the task jobs). In this sense, the users are unaware of the parallelism.



However, although it is a powerful and simple programming model, it is only suitable for applications that can be expressed in terms of *mappers* and *reducers*. This means that the applications must be designed with a *map* and a *reduce* function so that the MapReduce framework can only exploit the inherent parallelism.

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Job;
8 import org.apache.hadoop.mapreduce.Mapper;
9 import org.apache.hadoop.mapreduce.Reducer;
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
11 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
12
13 public class WordCount {
14
15     public static class TokenizerMapper extends Mapper<Object,Text,Text,IntWritable> {
16         private final static IntWritable one = new IntWritable(1);
17         private Text word = new Text();
18
19         public void map(Object key, Text value, Context context)
20             throws IOException, InterruptedException {
21
22             StringTokenizer itr = new StringTokenizer(value.toString());
23             while (itr.hasMoreTokens()) {
24                 word.set(itr.nextToken());
25                 context.write(word, one);
26             }
27         }
28     }
29
30     public static class IntSumReducer extends Reducer<Text,IntWritable,Text,IntWritable> {
31         private IntWritable result = new IntWritable();
32
33         public void reduce(Text key, Iterable<IntWritable> values, Context context)
34             throws IOException, InterruptedException {
35
36             int sum = 0;
37             for (IntWritable val : values) {
38                 sum += val.get();
39             }
40             result.set(sum);
41             context.write(key, result);
42         }
43     }
44
45     public static void main(String[] args) throws Exception {
46         Configuration conf = new Configuration();
47         Job job = Job.getInstance(conf, "word count");
48         job.setJarByClass(WordCount.class);
49         job.setMapperClass(TokenizerMapper.class);
50         job.setCombinerClass(IntSumReducer.class);
51         job.setReducerClass(IntSumReducer.class);
52         job.setOutputKeyClass(Text.class);
53         job.setOutputValueClass(IntWritable.class);
54         FileInputFormat.addInputPath(job, new Path(args[0]));
55         FileOutputFormat.setOutputPath(job, new Path(args[1]));
56         System.exit(job.waitForCompletion(true) ? 0 : 1);
57     }
58 }
59 }
```

FIGURE 2.15: Wordcount example on top of *Hadoop*  
Source: Hadoop Map Reduce tutorial

Figure 2.15 shows the implementation of a Wordcount application using the *Hadoop* [43] implementation of the MapReduce programming model. The Wordcount application counts the number of occurrences of each word in a given text, and it is a commonly used example for MapReduce because it is conceptually easy to define the map and the reduce functions. Notice that, in opposite to the previous frameworks, the user defines the *map*, *reduce* and the main code in a sequential-like way; without directly dealing with parallelism and without specifying the execution workflow.

When executing the previous MapReduce application (see Figure 2.16), the runtime processes the input, divides it into smaller sub-problems and distributes each subproblem to one of the available worker nodes. A worker can then repeat this step, leading to a multi-level tree structure where the intermediate nodes merge their worker results and send them back to their master node, or simply compute the result of the subproblem. In any case, the last master node merges all the intermediate values, collecting and combining the answers of all the sub-problems to obtain the final output.

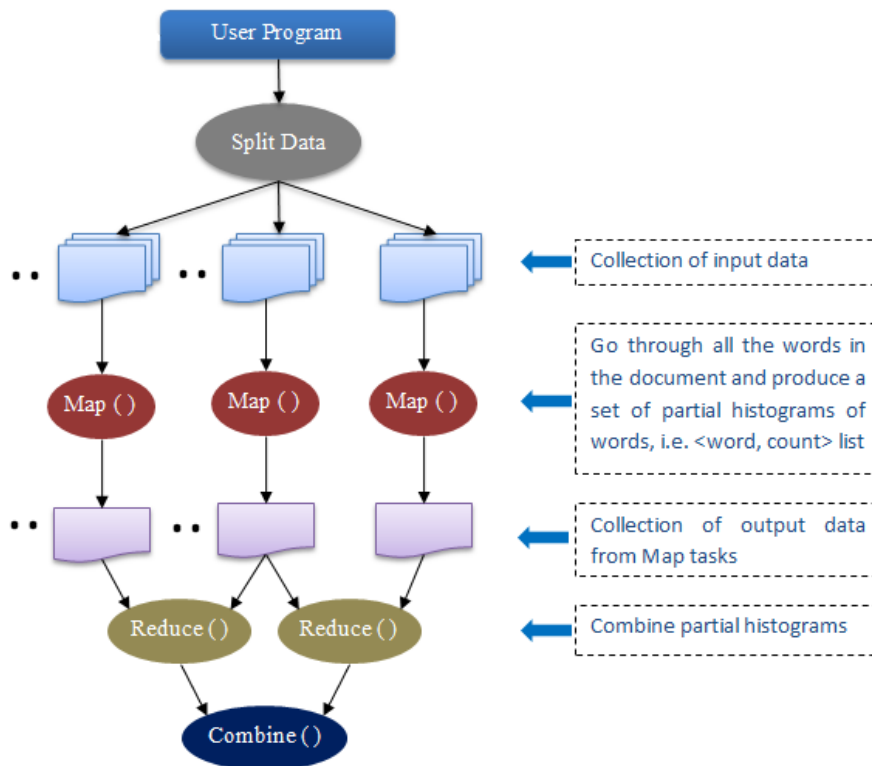


FIGURE 2.16: Execution example of Wordcount using MapReduce  
Source: SALSA Group, PTI Indiana University, Wordcount User Guide

As a final note, we must highlight that the MapReduce programming model is very popular since it has been implemented in several programming languages (including C++, Python, and Java) and implementations like *Hadoop* are widely used.

### 2.2.2.2 Spark

Apache Spark [52] is a fast and general-purpose large-scale data processing engine that supports general execution graphs, and that can run on top of Hadoop, Mesos, standalone

or in the cloud (see Figure 2.17). It extends the MapReduce programming model to efficiently support other types of computations (such as batch applications, interactive algorithms, SQL queries or streaming) and by allowing in-memory processing. These features make easy and inexpensive for users to combine different workloads and define complex pipelines. Moreover, it has been designed to be highly accessible, providing a simple API in Java, Scala [17], Python, R and SQL, and supporting a rich set of high-level tools including Spark SQL for SQL and structured data processing, MLib for machine learning, GraphX for graph processing, and Spark Streaming.

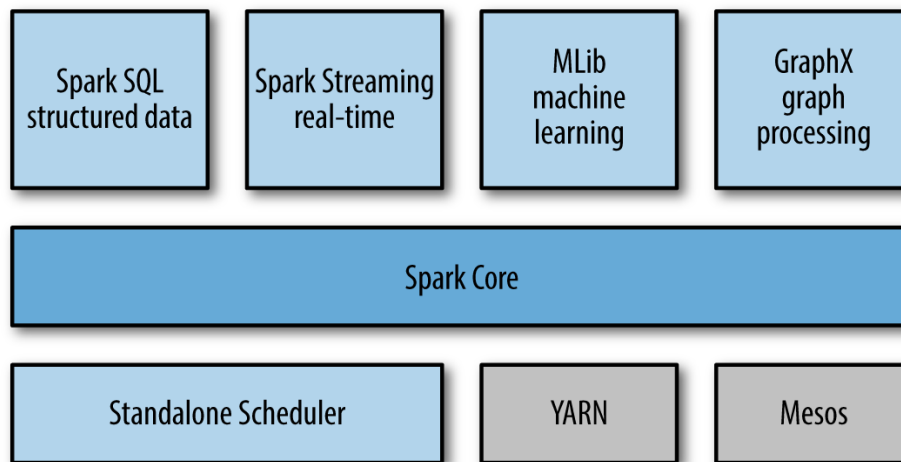


FIGURE 2.17: Spark's Components  
Source: Learning Spark [35]

*“At a high level, every Spark application consists of a driver program that runs the user’s main function and executes various parallel operations on a cluster. The main abstraction Spark provides is a Resilient Distributed Dataset (RDD), which is a collection of elements partitioned across the nodes of the cluster that can be operated on in parallel”*[52]. For the RDD creation, Spark relies on its underlying file system.

Figure 2.18 provides the Wordcount implementation in Java using the Spark API. Although the code can be written in a more compact way, we have kept it unrolled for the sake of clarity and to be easily mapped to the functions defined by the MapReduce programming model. Thus, as in the MapReduce example shown before, the functions provided are the Map and Reduce functions, and the main code is written in a sequential-like way. Notice that the code becomes significantly more readable than in the previous example and, as expected, the user does not explicitly define the workflow. In fact, when using Spark to develop an application, the user is unaware of the underlying infrastructure and the workflow diagram.

We must highlight that Apache Spark is widely used nowadays and it is becoming one of the most powerful Big Data engines. Although it still relies on the MapReduce programming model (making difficult the implementation of sophisticated algorithms which depend on the state of the processes running on other nodes), Spark provides an efficient and complete framework that separates the user from the common distributed programming issues such as data dependency analysis, task scheduling, fault-tolerant mechanisms and workers communication.

```

1 import java.util.Arrays;
import org.apache.spark.SparkConf;
3 import org.apache.spark.api.java.JavaPairRDD;
import org.apache.spark.api.java.JavaRDD;
5 import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.api.java.function.FlatMapFunction;
7 import org.apache.spark.api.java.function.Function2;
import org.apache.spark.api.java.function.PairFunction;
9 import scala.Tuple2;

11 public class WordCount {

13     private static final FlatMapFunction<String, String> SPLIT_WORDS
        = new FlatMapFunction<>() {
15         public Iterable<String> call(String s) {
            return Arrays.asList(s.split(" "));
17         }
        };

19     private static final PairFunction<String, String, Integer> MAP_WORDS
21     = new PairFunction<>() {
        public Tuple2<String, Integer> call(String s) {
23         return new Tuple2<String, Integer>(s, 1);
        }
25     };

27     private static final Function2<Integer, Integer, Integer> REDUCE_WORDS
        = new Function2<>() {
29         public Integer call(Integer a, Integer b) {
            return a + b;
31         }
        };

33     public static void main(String[] args) throws Exception {
35         String inputFile = args[0];
        String outputFile = args[1];
37         // Create a Java Spark Context
        SparkConf conf = new SparkConf().setAppName("wordCount");
39         JavaSparkContext sc = new JavaSparkContext(conf);
        // Load input data
41         JavaRDD<String> input = sc.textFile(inputFile);
        // Split
43         JavaRDD<String> words = input.flatMap(SPLIT_WORDS);
        // Transform into words and count
45         JavaPairRDD<String, Integer> pairs = words.mapToPair(MAP_WORDS);
        JavaPairRDD<String, Integer> counter = pairs.reduceByKey(REDUCE_WORDS);
47         // Save the word count
        counter.saveAsTextFile(outputFile);
49     }
}

```

FIGURE 2.18: Wordcount example in Java using Spark

### 2.2.2.3 Swift

Swift [53] is a fast and easy scripting language for executing distributed and parallel applications. The programming language is structured like a Shell script and builds a data-flow oriented workflow. The Swift programs run many copies of ordinary applications as soon as their inputs are available. The main advantage is that its workflows can be fast and easily defined, and can run over any underlying architecture (multicore computers, clusters, clouds or grids). The users must consider that Swift it is not designed to process large collections of data but to orchestrate programs that do that processing. In this sense, it is intended to be a pure Workflow Manager: handling the execution of such programs on remote sites, staging input and output files from these sites and choosing these sites.

In opposition to the previous cases, Swift is a programming language itself built on top of Java. Consequently, it provides data types, mapped data types, mapped functions, conventional expressions, structured data, loops and data flow instructions so that users can define their workflows (see Figure 2.19). Additionally, Swift uses a configuration file to specify the underlying infrastructure.

<ul style="list-style-type: none"> <li>▪ <b>Data types</b></li> </ul> <pre>string s = "hello world"; int i = 4; int A[];</pre> <ul style="list-style-type: none"> <li>▪ <b>Mapped data types</b></li> </ul> <pre>type image; image file1&lt;"snapshot.jpg"&gt;;</pre> <ul style="list-style-type: none"> <li>▪ <b>Mapped functions</b></li> </ul> <pre>app (file o) myapp(file f, int i) { mysim "-s" i @f @o; }</pre> <ul style="list-style-type: none"> <li>▪ <b>Conventional expressions</b></li> </ul> <pre>if (x == 3) {     y = x+2;     s = @strcat("y: ", y); }</pre>	<ul style="list-style-type: none"> <li>▪ <b>Structured data</b></li> </ul> <pre>image A[]&lt;array_mapper...&gt;;</pre> <ul style="list-style-type: none"> <li>▪ <b>Loops</b></li> </ul> <pre>foreach f,i in A {     B[i] = convert(A[i]); }</pre> <ul style="list-style-type: none"> <li>▪ <b>Data flow</b></li> </ul> <pre>analyze(B[0], B[1]); analyze(B[2], B[3]);</pre>
-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

FIGURE 2.19: Swift programming language. Source: Swift [53]

To exemplify how Swift works, Figure 2.20 shows the simulation workflow of the Swift code shown in Figure 2.21. Notice that the code first defines the two applications that the workflow will invoke: the simulation and the analysis. For both applications, it also defines their input and output parameters, and their execution instructions (in both cases, the execution is an invocation of a bash script to launch the simulation or the analysis). Next, it retrieves the workflow arguments and define the intermediate variables needed to execute the workflow. In a third step, it spawns *nsim* simulations; storing its result into separated files. Finally, it launches an *analyze* application that uses the output files of all the previously launched simulations.

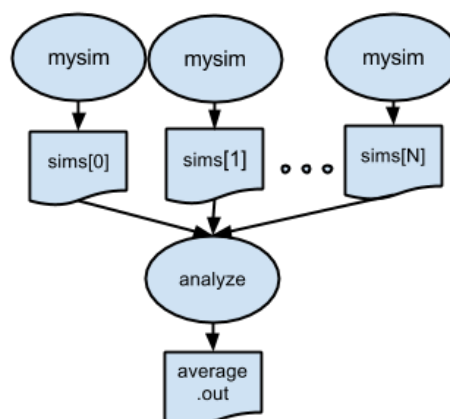


FIGURE 2.20: Swift simulation workflow example. Source: Swift tutorial[53]

```

type file;
2
app (file out, file log)
4
simulation (int sim_steps, int sim_range, int sim_values, file sim_script)
6
{
  bash @sim_script "--timesteps" sim_steps "--range" sim_range
 8
    "--nvalues" sim_values stdout=@out stderr=@log;
}
10
app (file out, file log) analyze (file s[], file stat_script)
12
{
  bash @stat_script filenames(s) stdout=@out stderr=@log;
14
}

16
int nsim = toInt(arg("nsim", "10"));
int steps = toInt(arg("steps", "1"));
18
int range = toInt(arg("range", "100"));
int values = toInt(arg("values", "5"));
20

file sims[];
22
file simulate_script <"simulate.sh">;
file stats_script <"stats.sh">;
24

foreach i in [0:nsim-1] {
26
  file simout <single_file_mapper; file=strcat("output/sim_",i,".out")>;
  file simlog <single_file_mapper; file=strcat("output/sim_",i,".log")>;
28
  (simout,simlog) = simulation(steps,range,values,simulate_script);
  sims[i] = simout;
30
}

32
file stats_out<"output/average.out">;
file stats_log<"output/average.log">;
34
(stats_out, stats_log) = analyze(sims,stats_script);

```

FIGURE 2.21: Swift simulation code example

We must highlight that Swift builds the workflow considering the data dependencies between applications and exploits the maximum parallelism between applications. Thus, users must define the workflow regarding data dependencies between applications but do not need to build the workflow structure explicitly. Moreover, the applications can be any type of executable, like a shell does, and do not need to follow any convention (like MapReduce obligates users to define *Map* and *Reduce* functions).

#### 2.2.2.4 COMP Superscalar (COMPSs)

Since this project is based in COMP Superscalar (COMPSs) [19] we have decided to provide a more in-depth view of it in Chapter 3.

## Chapter 3

# COMPSs overview

COMP Superscalar (COMPSs) [4] [33] is a task-based programming model that belongs to the family of Frameworks with implicit workflows. COMPSs applications consist of three parts: the application's code developed in a totally sequential manner, an application interface where the programmers specify which functions can be remotely executed (*tasks*) and a configuration file that describes the underlying infrastructure. With these three components, the COMPSs Runtime system exploits the inherent parallelism of the application at execution time by detecting the task calls and the data dependencies between them.

COMPSs natively supports Java applications but also provides bindings for Python (Py-COMPSs [23]) and C/C++. Furthermore, COMPSs allows applications to be executed on top of different infrastructures (such as multi-core machines, grids, clouds or containers) without modifying a single line of the application's code (see Figure 3.1). It also has fault-tolerant mechanisms for partial failures (with job resubmission and reschedule when task or resources fail), has a live monitoring tool through a built-in web interface, supports instrumentation using the Extrae [8] tool to generate post-mortem traces that can be analysed with Paraver [15], has an Eclipse IDE, and has pluggable cloud connectors and task schedulers.

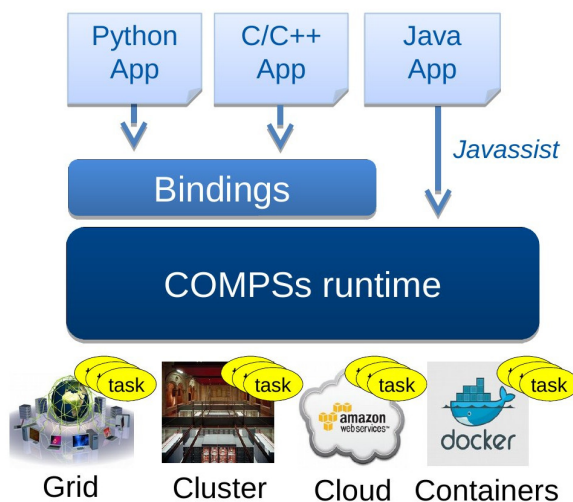


FIGURE 3.1: COMPSs overview

Additionally, the COMPSs model has three key characteristics:

- **Sequential Programming:** The users do not need to deal with any parallelization and distribution issue such as thread creation, synchronization, data distribution, messaging or fault-tolerance. COMPSs programmers only select which methods must be considered as tasks and the COMPSs Runtime spawns them asynchronously on a set of resources instead of executing them locally and sequentially.

- **Infrastructure Agnostic:** COMPSs model abstracts the application from the underlying infrastructure. Hence, COMPSs applications do not include any platform related detail such as deployment or resource management. This feature makes applications portable between infrastructures with different characteristics.
- **No APIs:** When using COMPSs native language, Java, the model does not require any special API, pragma or construct in the program. Since COMPSs instruments the application's code at execution time to detect the tasks, everything can be developed in the standard Java syntax and libraries.

### 3.1 Programming Model

The COMPSs programming model mainly involves choosing the right tasks for a sequential application. To illustrate how COMPSs applications are built from a sequential code, Figure 3.2 provides the main code and Figure 3.3 provides the helper methods of a sequential Increment application in Java. The Increment is a didactic application that takes  $N$  counters, initializes them to a random value and increments them by  $U$  units.

```

1 package increment;
3 import java.util.Random;
5
7 public class Increment {
9     private static final int MAX_INT = 10;
11
12     public static void main(String[] args) {
13         // Retrieve arguments
14         if (args.length != 2) {
15             System.err.println("[ERROR] Invalid number of arguments");
16             System.err.println("    Usage: increment.Increment <N> <U>");
17             System.exit(1);
18         }
19         int N = Integer.parseInt(args[0]);
20         int U = Integer.parseInt(args[1]);
21
22         // Initialize counters
23         Integer[] counters = new Integer[N];
24         for (int i = 0; i < N; ++i) {
25             counters[i] = new Random().nextInt(MAX_INT);
26             System.out.println("[LOG] Initial Counter " + i + " value is " + counters[i]);
27         }
28
29         // Increment U units each counter
30         for (int i = 0; i < U; ++i) {
31             for (int j = 0; j < N; ++j) {
32                 counters[j] = IncrementImpl.increment(counters[j]);
33             }
34         }
35
36         // Show final counter values
37         for (int i = 0; i < N; ++i) {
38             System.out.println("[LOG] Final Counter " + i + " value is " + counters[i]);
39         }
40     }
41 }

```

FIGURE 3.2: Increment main class



```

1 package increment;
3 public class IncrementImpl {
5     public static Integer increment(Integer counter) {
6         // Return the increased counter
7         return counter + 1;
8     }
9 }

```

FIGURE 3.3: Increment helper methods class

After implementing the sequential code, the users must create an *Interface* to transform their sequential application into a COMPSs application. COMPSs requires this *Interface* to *annotate* which methods must be considered *tasks* that can be remotely executed.

COMPSs Interfaces must be defined inside a file with the same name than the main class of the users' application but with the "Itf" suffix (for instance, in the previous example, the Interface must be stored in the *IncrementItf.java* file). Regarding its content, the interface must contain one entry per task, which is annotated with the `@Method` annotation<sup>1</sup>. For each *Method* annotation the users must also provide the *declaring class* and the parameters description. On the one hand, the declaring class of a function is the class containing the implementation of the task and its required to link the task to the method implementation. On the other hand, the parameters description is indicated by adding a `@Parameter` annotation to every task parameter. This annotation is required to build the task dependency graph since COMPSs uses data-flow graphs. The mandatory contents of the *Parameter* annotation are the *Type* (that must refer to any Java basic type, a string, an object or a file) and the *Direction* (where the only valid values are *IN*, *OUT* and *INOUT*).

```

1 package increment;
2 import integratedtoolkit.types.annotations.Parameter;
3 import integratedtoolkit.types.annotations.parameter.Direction;
4 import integratedtoolkit.types.annotations.parameter.Type;
5 import integratedtoolkit.types.annotations.task.Method;
6
7
8 public interface IncrementItf {
9
10     @Method(declaringClass = "increment.IncrementImpl")
11     Integer increment(
12         @Parameter(type = Type.OBJECT, direction = Direction.IN) Integer counter
13     );
14 }
15
16 }

```

FIGURE 3.4: Increment Interface

Figure 3.4 shows the Interface for the previous Increment application example. Notice that it only contains one task declaration, *increment*, with defined inside the *IncrementImpl* class and with two parameters. The first parameter is the return value of the function which has type *Integer* and, by default, has direction *OUT*. The second parameter is the *counter*

<sup>1</sup>There are more complex COMPSs annotations that are beyond the scope of this section. However, for any interested reader, the COMPSs annotations are described in-depth in the *COMPSs User Guide: Application Development* [19]

argument which has type *Integer* and direction IN because the function requires the input value of the parameter to increase it but does not modify it.

Since the COMPSs annotations do not interfere with the applications' code, all COMPSs applications can be sequentially executed. To do so, Figure 3.5 compiles the previous code and executes the application with  $N = 2$  counters that must be increased by  $U = 3$ .

```

$ javac increment/*
2 $ jar cf increment.jar increment/
$ java -cp increment.jar increment.Increment 2 3
4 [LOG] Initial Counter 0 value is 7
  [LOG] Initial Counter 1 value is 1
6 [LOG] Final Counter 0 value is 10
  [LOG] Final Counter 1 value is 4

```

FIGURE 3.5: Sequential execution example of Increment

On the other hand, the code can also be executed with COMPSs without recompiling the application's code. To do so, the users must invoke the *runcompss* command instead of the traditional *java* command. When done, the COMPSs Runtime will be setup and the application will be distributedly executed. Figure 3.6 provides the execution output of the Increment application and the task graph generated by its execution.

```

$ runcompss -g increment.Increment 2 3
[ INFO] Using default execution type: compss
[ INFO] Using default location for project file
[ INFO] Using default location for resources file
[ INFO] Using default language: java

----- Executing increment.Increment -----

WARNING: IT Properties file is null. Setting default values
[(790)  API] - Starting COMPSs Runtime v2.0
[LOG] Initial Counter 0 value is 5
[LOG] Initial Counter 1 value is 2
[LOG] Final Counter 0 value is 8
[LOG] Final Counter 1 value is 5
[(5090) API] - Execution Finished

```

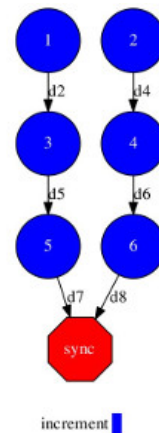


FIGURE 3.6: COMPSs execution example of Increment

The *runcompss* command has several command line arguments (that are fully detailed when executing *runcompss -help*) but Table 3.1 provides a short description of those we have found to be more useful.

Argument	Description
-d	Enables the debug mode
-g	Enables the final graph generation
-m	Enables the monitor tool
-t	Enables the tracing tool
--summary	Provides a task summary at the end of the execution
--lang=<str>	Enables the Python and C/C++ bindings
--project=<str>	Sets an specific project configuration file
--resources=<str>	Sets an specific resources configuration file
--classpath=<str>	Adds an specific classpath to the execution environment

TABLE 3.1: Useful arguments for the *runcomps* command

## 3.2 Runtime System

To abstract applications from the underlying infrastructure, COMPSs relies on its Runtime System to spawn a master process on the machine where the application is running and a worker process per available resource (see Figure 3.7). These processes are communicated through the network (using different communication adaptors) and can send messages to each other to orchestrate the distributed execution of the application.

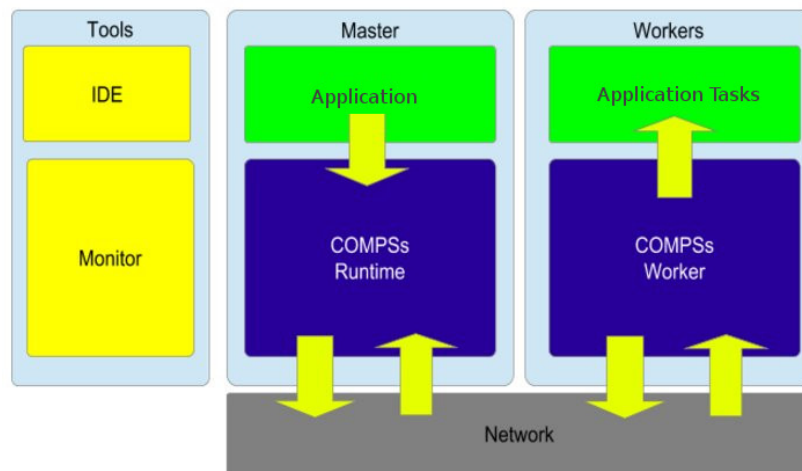


FIGURE 3.7: COMPSs structure

Once a Java application starts, the COMPSs Runtime [34] triggers a custom Java Class-Loader that uses Javassist [18] to instrument the application's main class. The instrumentation modifies the original code by inserting the necessary calls to the COMPSs API to generate tasks, handle data dependencies and add data synchronizations. To achieve the same purpose on Python applications, the Python Binding (PyCOMPSs) parses the decorators of the main code and adds the necessary calls to the COMPSs API. In the case of C/C++ applications, COMPSs also requires an Interface file that is used when compiling the application to generate stubs for the main code, add the required COMPSs API calls, and generate the

code for the tasks execution at the workers. In any case, as shown on the top of the Figure 3.8, the interaction between the application and the COMPSs Runtime is always made through the COMPSs API.

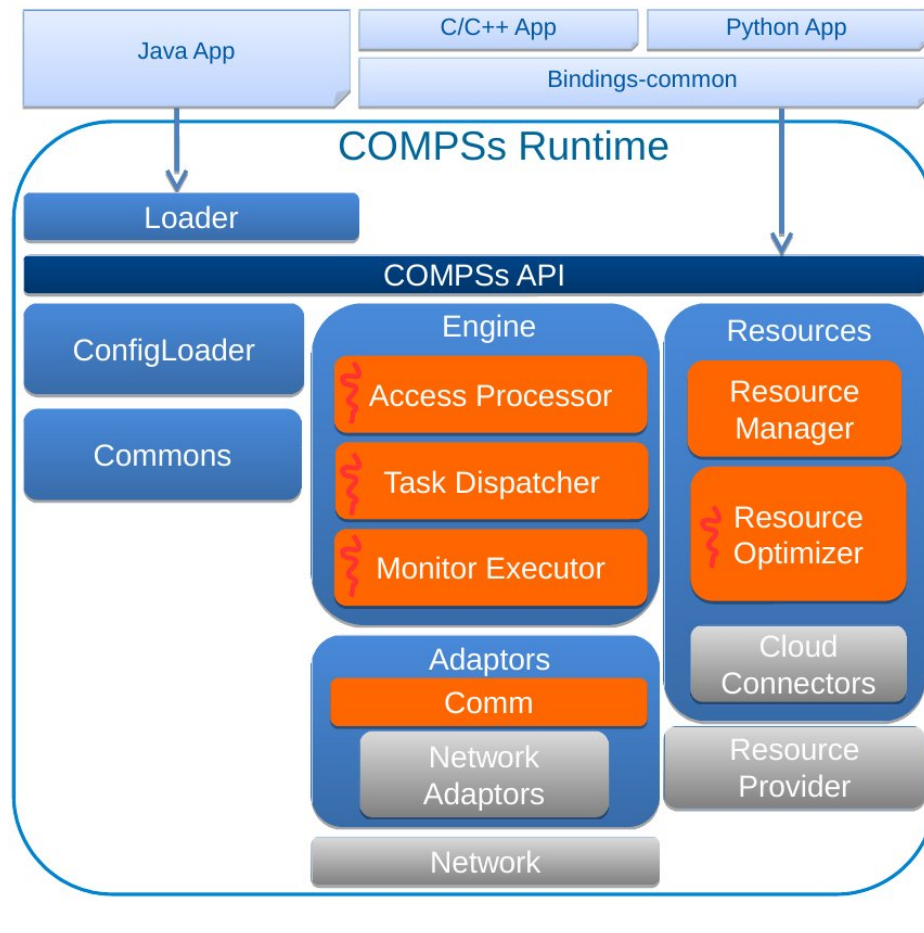


FIGURE 3.8: COMPSs Runtime overview

More in-depth, the COMPSs Runtime has five main components:

- **Commons** Contains the common structures used by all the Runtime components
- **ConfigLoader** Loads the project and the resources configuration files, the command line arguments, and the JVM configuration parameters.
- **Engine** Contains the submodules to handle the task detection, the data dependencies, and the task scheduling. More specifically, the *Access Processor* watches for the data accesses so that the Runtime can build the data dependencies between tasks, the *Task Dispatcher* controls the task life-cycle and the *Monitor Executor* controls the monitor structures for real-time and post-mortem monitoring.
- **Resources** Handles all the available resources in the underlying infrastructure. This component creates, destroys and monitors the state of all the available resources. Since COMPSs supports elasticity through cloud connectors, this component contains a *Resource Optimizer* subcomponent that takes care of creating and destroying resources.
- **Adaptors** Contains the different communication adaptors implementations. This layer is used to communicate the COMPSs Master and the COMPSs Workers and abstracts the rest of the Runtime from the different network adaptors.

### 3.3 Task Workflow

To clarify how COMPSs works when executing an application Figure 3.9 describes the task life-cycle. From the application's main code the COMPSs API registers the different tasks. Considering the registered tasks, COMPSs builds a task graph based on the data dependencies. This graph is then submitted to the *Task Dispatcher* that schedules the data-free tasks when possible. This means that a task is only scheduled when it is data-free, and there are enough free resources to execute it (each task can have different constraints, and thus, it is not scheduled if there is not a resource that satisfies the requirements).

Eventually, a task can be scheduled and, then, it is submitted to execution. This step includes the job creation, the transfer of the input data, the job transfer to the selected resource, the real task execution on the worker and the output retrieval from the worker back to the master. If any of these steps fail, COMPSs provides fault-tolerant mechanisms for partial failures.

Once the task has finished, COMPSs stores the monitoring data of the task, synchronizes any data required by the application, releases the data-dependant tasks so that they can be scheduled, and deletes the task.

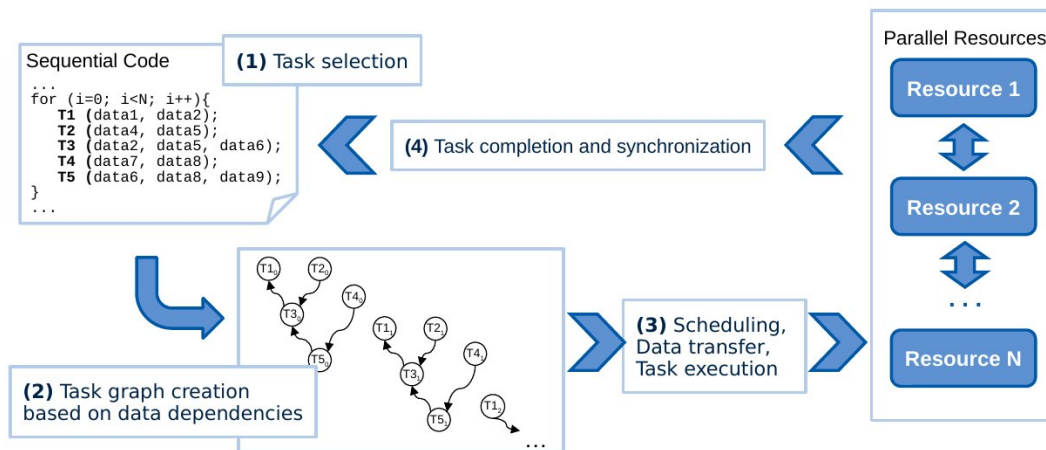


FIGURE 3.9: COMPSs task execution workflow



## Chapter 4

# Tools and methodology

On the one hand, this chapter lists the tools used to enhance COMPSs to enable analytic and HPC Workflows. Considering that this project is based on the extension of an already existing tool, COMPSs, most of the other tools used for its extension are already preset. On the other hand, Section 4.2, reports the development methodology of the project.

### 4.1 Tools

As stated, the principal tool used has been COMPSs as it was before this project (see Chapter 3). Since the COMPSs Runtime is developed in Java[39] language, the biggest part of the development of this project has been performed in Java language using the Eclipse[36] IDE and the Apache Maven[45] Software Project Management. The implementation has been completed with BASH[25] scripts to handle the communication between worker processes (see Section 5.3 for extended details).

For the results and the evaluation section, we have used Extrae[8] and Paraver[15] to validate the parallel executions and GnuPlot[27] to illustrate simulation results.

### 4.2 Methodology

We have selected a Design Research Method as the scientific method design, combined with a Test Driven Development strategy, always bearing in mind that the main goal of this project is to enhance COMPSs to enable analytic and HPC Workflows. Thus, we have selected conceptual use cases to theorize the COMPSs extensions (Section 4.2.1), carried them out following the development strategy (Section 4.2.2) and validated with real use cases (Section 4.2.3).

#### 4.2.1 Scientific method design

In a first step, during the Relevance Cycle, we have analyzed analytic and HPC applications that can potentially be ported to COMPSs. In this sense, we have concluded that the COMPSs Programming Model lacks an easy way to execute binaries, MPI[38] and OmpSs[14] applications; which prevents some users from using COMPSs.

In a second step, during the Design Cycle, we have carefully selected two real use cases that represent the needs of the scientific community and which allow us to evaluate the requirements that must be fulfilled. The first application, *NMMB BSC-Dust* [12], is a model to predict the atmospheric life cycle of the desert dust for regional and global domains. The second application, *GUIDANCE* [9], is an integrated solution for Genome and Phenome association analysis.

However, the requirements of these applications turned out to be too complicated to be faced in a single development phase, which led us to find out simpler use cases and evaluate these during the validation phase. Consequently, we have faced the problem gradually

by dividing the requirements into two different categories. In the first category, we have considered the execution of simple binaries. To drain the standard binary execution requirements, we have analyzed a COMPSs version of the BLAST[6] application. To enhance COMPSs, we have proceeded gradually by defining incremental tests and finally achieving all the BLAST requirements. In the second category, we have considered a hybrid COMPSs and MPI version of the Matrix Multiplication, and we have enhanced COMPSs accordingly. As in the previous case, to ensure that all the requirements are fulfilled, we have performed incremental tests to finally execute this application.

Moreover, at the end of this second step, we have evaluated the real applications. We must highlight that we have gone through a validation process to evolve from the simpler use cases to the real ones.

Finally, during the Rigor Cycle, we have deployed the project development on a production environment, and we have largely documented the new artifacts to share the knowledge with the community.

### 4.2.2 Development strategy

The development has been based on the aforementioned proof of concept applications (BLAST and Matmul), and we have followed a Test Driven Development strategy. Although we have applied this methodology for the first time during this project, it turned out to be very robust, flexible and appropriate for the case.

More in-depth, we have defined several easy-to-run tests for each simple use case. Once identified the application requirements, we have implemented the features needed inside the COMPSs Runtime in a top-down strategy. The development has been performed in such a way, rather than following a bottom-up strategy, because the chosen use cases strongly define the requirements from the COMPSs API (method and parameter annotations) and are more flexible with regard to the COMPSs internals.

In this sense, from the first simple use case (BLAST), we have first enhanced the COMPSs Annotations to extend the method and the parameter annotation. Next, we have associated this information to the tasks and propagated it through the COMPSs Runtime up to the COMPSs Worker. Once all the information has been transferred to the COMPSs Worker, we have enhanced the task treatment (through new invokers) and, finally, we have implemented the binary executors.

From the second use case (Matmul), we have first extended our new annotation with an MPI method annotation. Then we have extended the COMPSs Schedulers to support multi-node tasks and, finally, we have extended the new invokers and executors to run MPI binaries.

We would like to highlight that the Test Driven Development has helped us significantly in making this process iterative without turning it into a nightmare. During the development, both use cases have been tested several times on the local machine to evaluate the development.

### 4.2.3 Validation strategy

Even if we have tested locally simple applications, an extensive validation has been performed to move from the simple use cases to the real applications. This validation process includes porting all the applications to the MareNostrum III [10] Supercomputer to execute and to perform an in-depth analysis of both types, the simple and the real of use cases (see Chapter 6 for more details).



## Chapter 5

# Implementation

### 5.1 Programming model annotations

As shown in Chapter 3, the COMPSs Programming model defines annotations that must be added to the sequential code in order to run the applications in parallel. These annotations can be split into two groups:

- **Method Annotations** Annotations added to the sequential code methods to detect them as tasks and potentially execute them in parallel.
- **Parameter Annotations** Annotations added to the parameters of an annotated method to handle data dependencies and transfers.

Consequently, we have first extended the Method Annotations to detect Binary, MPI and OmpSs methods (from now on, *non-native* methods). Next, we have allowed the users to define annotations with environment variables. In this step, we have taken advantage of all these modifications to change also the annotations referring to the versioning and the scheduler hints. Finally, we have modified the parameter annotations to support specific binary needs.

#### 5.1.1 New task annotations

As previously stated, we have decided to firstly extend the method annotations to recognize the different type of tasks. The task annotations are used by the COMPSs Runtime to detect the methods that the user wants to execute as tasks. The annotations are defined in a separated file (known as *interface*) that the COMPSs Loader uses to detect the tasks while instrumenting the application code. In this sense, the interface defines the signature of the method that the COMPSs Loader must transform into a task. Since a Java signature must contain the fully qualified name (that is, the declaring class of the method, the method name, and the parameter types), the COMPSs annotations provide a declaring class for each *Method* annotation.

To extend the annotations to support non-native tasks we have firstly modified the COMPSs Loader to consider as potential task any method defined inside the *binary.BINARY*, *mpi.MPI* and *omps.OMPSS* classes. The methods found in these classes are then cross-validated with the tasks defined in the interface and selected as tasks if needed. This design decision is motivated by the fact that we consider that the annotation of non-native tasks must only refer to the real execution and, thus, we want to avoid a *declaringClass* field in the new annotation.

Secondly, we have extended the method annotations themselves to recognize the different type of tasks. For each non-native method that we are willing to support (execution of binaries, MPI binaries, and OmpSs binaries), we have created a new annotation that extends the information of the current *Method* annotation with the specificities of each type.

In the Binary case, we have created the `@Binary` annotation that **must** define the associated binary file name through the `binary` field (see Figure 5.1). Moreover, this annotation can optionally define the `workingDir` field (to set up the working directory of the binary execution) and, as any `Method` annotation does, the `priority` and `constraints` fields (see Figure 5.2).

```
1 @Binary(binary = "binary")
   void binaryTask();
```

FIGURE 5.1: Binary annotation using only mandatory fields

```
2 @Binary(binary = "binary",
3         workingDir = "/tmp/",
4         priority = true,
   constraints = @Constraints(computingUnits = "1"))
   void binaryTask();
```

FIGURE 5.2: Binary annotation using all available fields

In the MPI case, we have created the `@MPI` annotation that **must** define the MPI command to be run (also known as MPI Runner) through the `mpiRunner` field, the associated binary file name through the `binary` field, and the number of computing nodes to be reserved for the MPI execution through the `computingNodes` field (see Figure 5.3). Moreover, this annotation can optionally define the `workingDir` field (to set up the working directory of the binary execution) and, as any `Method` annotation does, the `priority` and `constraints` fields. Figure 5.4 shows a complete example of this annotation. Notice that the total number of requested Computing Units is  $2 \text{ nodes} \cdot 4 \frac{\text{CUs}}{\text{node}} = 8 \text{ total CUs}$ .

```
1 @MPI(mpiRunner = "mpirun",
2     binary = "mpiBinary",
3     computingNodes = "2")
   void mpiTask();
```

FIGURE 5.3: MPI annotation using only mandatory fields

```
2 @MPI(mpiRunner = "mpirun",
3     binary = "mpiBinary",
4     computingNodes = "2",
5     workingDir = "/tmp/",
6     priority = true,
   constraints = @Constraints(computingUnits = "4"))
   void mpiTask();
```

FIGURE 5.4: MPI annotation using all available fields

In the OmpSs case, we have created the `@OmpSs` annotation that **must** define the associated binary through the `binary` field (see Figure 5.5). Moreover, this annotation can optionally define the `workingDir` field (to set up the working directory of the binary execution) and, as any `Method` annotation does, the `priority` and `constraints` fields (see Figure 5.6).

```
1 @OmpSs(binary = "ompssBinary")
   void ompssTask();
```

FIGURE 5.5: OmpSs annotation using only mandatory fields

```
2 @OmpSs(binary = "ompssBinary",
3         priority = true,
4         workingDir = "/tmp/",
5         constraints = @Constraints(computingUnits = "4"))
   void ompssTask();
```

FIGURE 5.6: OmpSs annotation using all available fields

Only considering the modifications listed in this section, the current COMPSs annotations (for `Method` and `Service` tasks) were left intact. However, for the sake of clarity and to enhance the COMPSs capabilities, we have decided to modify the existant annotations to include environment variables as annotations' values, a clearer versioning annotation and new annotations for the upcoming schedulers.

### 5.1.2 Environment variables as annotations

A large variety of applications use different constraint values depending on the execution even if the code of the application remains the same. This is mostly because users want to adapt the task needs to the different data sizes and computational load. However, the current COMPSs version only allowed users to define different constraints for tasks by re-compiling all the application with a different interface. Figure 5.7 shows the Interface of a clear example of this case: a Wordcount application that always defines two tasks (`wordCount` and `mergeResults`) with the same code. The execution load of the `wordCount` task is highly related to the size of the data to be treated. Thus, users want to specify a different memory requirement depending on the input data size (because the file must be totally loaded in memory). However, to change the value of the `memorySize` constraint, they must recompile the whole application.

```
1 public interface WordcountItf {
2     @Method(declaringClass = "WordcountImpl")
3     public HashMap<String, Integer> mergeResults(
4         @Parameter HashMap<String, Integer> m1,
5         @Parameter HashMap<String, Integer> m2
6     );
7     @Method(declaringClass = "WordcountImpl")
8     @Constraints(memorySize = "4.0")
9     public HashMap<String, Integer> wordCount(
10        @Parameter(type = Type.FILE, direction = Direction.IN) String filePath
11    );
12 }
```

FIGURE 5.7: Wordcount Interface

Figure 5.8 shows how a user must run two different executions of this application, one with a big file (for instance, needing 16 Gb of memory) and one with a small file (needing 1 Gb of memory).

```

# Modify the application Interface
2 $ vi WordcountItf.java # Change the constraint value to 1.0
# Clean and build the application
4 $ mvn clean package
# Run the application with COMPSs
6 $ runcomps wordcount.Wordcount /path/to/small/file
# Modify the application Interface
8 $ vi WordcountItf.java # Change the constraint value from 1.0 to 16.0
# Clean and build the application
10 $ mvn clean package
# Run the application with COMPSs
12 $ runcomps wordcount.Wordcount /path/to/big/file

```

FIGURE 5.8: Wordcount executions with different constraint values

In our new approach, the annotations' values can be specified through environment variables that are resolved on the COMPSs Master node in execution time. Following the previous example, Figure 5.9 shows the same application using an environment variable to define the memory size annotation. Figure 5.10 illustrates how the user can easily execute the application several times with different input data files with only redefining the environment variable (we must highlight, that the users no longer need to recompile the whole application between executions).

```

public interface WordcountItf {
2   @Method(declaringClass = "WordcountImpl")
   public HashMap<String, Integer> mergeResults(
4     @Parameter HashMap<String, Integer> m1,
     @Parameter HashMap<String, Integer> m2
6   );
   @Method(declaringClass = "WordcountImpl")
8   @Constraints(memorySize = "${MEMORY_SIZE}")
   public HashMap<String, Integer> wordCount(
10    @Parameter(type = Type.FILE, direction = Direction.IN) String filePath
   );
12 }

```

FIGURE 5.9: Wordcount Interface with environment variables

```

# Clean and build the application
2 $ mvn clean package
# Modify the environment variable to be loaded by the Interface
4 export MEMORY_SIZE=1.0
# Run the application with COMPSs
6 $ runcomps wordcount.Wordcount /path/to/small/file
# Modify the environment variable to be loaded by the Interface
8 export MEMORY_SIZE=16.0
# Run the application with COMPSs
10 $ runcomps wordcount.Wordcount /path/to/big/file

```

FIGURE 5.10: Wordcount executions with environment variables as constraints

This behavior can be used with any already existing annotation. Next, Figure 5.11 provides a complete example of all the variables that can be resolved at execution time from environment variables (environment variables are detected by the use of the `$` at the beginning of the field value). Notice that the only values that we are preventing the users from defining with environment variables are the *declaringClass* field inside the *Method* annotation and the parameter annotations.

```

2 public class MainItf {
3
4     @Method(declaringClass = "MainImpl",
5             isModifier = "${isModifier}",
6             priority = "${isPrioritary}")
7     @Binary(binary = "${binary}",
8             workingDir = "${wd}",
9             priority = "${isPrioritary}")
10    @MPI(mpiRunner = "${mpiRunner}",
11         binary = "${mpiBinary}",
12         computingNodes = ${computingNodes},
13         workingDir = "${wd}",
14         priority = "${isPrioritary}")
15    @OmpSs(binary = "${ompssBinary}",
16           workingDir = "${wd}",
17           priority = "${isPrioritary}")
18    @Constraints(computingUnits = "${computingUnits}",
19                processorName = "${processorName}",
20                processorSpeed = "${processorSpeed}",
21                processorArchitecture = "${processorArchitecture}",
22                processorPropertyName = "${processorPropertyName}",
23                processorPropertyValue = "${processorPropertyValue}",
24                memorySize = "${memorySize}",
25                memoryType = "${memoryType}",
26                storageSize = "${storageSize}",
27                storageType = "${storageType}",
28                operatingSystemType = "${operatingSystemType}",
29                operatingSystemDistribution = "${operatingSystemDistribution}",
30                operatingSystemVersion = "${operatingSystemVersion}",
31                appSoftware = "${appSoftware}",
32                hostQueues = "${hostQueues}",
33                wallClockLimit = "${wallClockLimit}")
34    @SchedulerHints(isDistributed = "${isDistributed}",
35                   isReplicated = "${isReplicated}")
36    int task(
37        @Parameter(type = Type.STRING, direction = Direction.IN) String message
38    );
39 }

```

FIGURE 5.11: Example with an Interface with all the available environment variables

We must highlight that this modification seems to have low sides effects, but it is rather the other way around. For the end-users, these modifications imply that all the annotations' values are of type *string* rather than a specific type for each of them. For the COMPSs Runtime, it means that all the annotation values must be resolved and transformed to its real value (integer, float, string, etc.). Notice that while this transformation is simple when the final value is an integer or a float (because the constraint value can only be an integer/float value or an environment variable), it is not that simple when the final value is a string (because the constraint value can have zero, one or more than one environment variables). For instance, Figure 5.12 shows many different possibilities that the user can specify when providing the *workingDir* constraint.

```

workingDir = "${wd}"
2 workingDir = "$wd"
workingDir = "/path/${wd}"
4 workingDir = "${wd}/path"
workingDir = "/path/$wd/path"
6 workingDir = "/path/$wd1/path/${wd2}/path"

```

FIGURE 5.12: Example of complex environment variables on the *workingDir* field

Finally, when using environment variables as constraints, the Runtime will only raise an exception if the environment variable is not defined (*null*), empty or it cannot be parsed to its real type.

### 5.1.3 Versioning task annotation

Versioning is a mechanism that COMPSs provides for defining several implementations of the same task and deciding, at execution time, which is the best implementation to run.

In the previous COMPSs version, the user defines a task in the interface and defines more than one *declaringClass* field for the different implementations. Figures 5.13, 5.14 and 5.15 show how a task *sayHello* is called by the main code and two different implementations of this task with the same name but in separated files. Notice that the callee in the main code refers to the *Impl1* but in execution time the COMPSs Runtime will choose to execute any of the implementations (either *Impl1* or *Impl2*). That means that using the callee *Impl1.sayHello()* or *Impl2.sayHello()* will make no difference at execution time.

```

public class Hello {
2
   public static void main(String[] args) {
4       System.out.println("Hello World (from main code)");
       // Launch task
6       // We could use Impl2.sayHello() indistinctly
       Impl1.sayHello();
8   }
}

```

FIGURE 5.13: Example of previous versioning main code: Hello.java

```

1 public class Impl1 {
   public static void sayHello() {
3       System.out.println("Hello World (from Implementation 1)");
   }
5 }

```

FIGURE 5.14: Example of previous versioning Impl1.java

```

1 public class Impl2 {
2     public static void sayHello() {
3         System.out.println("Hello World (from Implementation 2)");
4     }
5 }

```

FIGURE 5.15: Example of previous versioning Impl2.java

Following with the example, Figure 5.16 defines its interface. The task *sayHello* is annotated with two different implementations declared in classes *Impl1* and *Impl2* respectively.

```

1 public interface HelloItf {
2
3     @Method(declaringClass = {"Impl1", "Impl2"})
4     void sayHello();
5 }

```

FIGURE 5.16: Example of the Interface of previous versioning: HelloItf.java

The user has also available the *MultiConstraints* annotation to add the needed requirements for each implementation. The programming model supports a global constraint annotation for the common constraints of all the implementations and a *MultiConstraint* annotation for the specific constraints of each implementation. Notice that the number of entries within the *MultiConstraint* annotation must be the same as the number of entries within the *declaringClass* field. Figure 5.17 redefines the Interface of the previous example to add a global constraint of 1 computing unit to both implementations and a specific memory constraint for each implementation. Obviously, this example can be extended to any constraint field.

```

1 public interface HelloItf {
2
3     @Method(declaringClass = {"Impl1", "Impl2"})
4     @Constraints(computingUnits = 1)
5     @MultiConstraints({
6         @Constraints(memorySize = 1),
7         @Constraints(memorySize = 2)})
8     void sayHello();
9 }

```

FIGURE 5.17: Example of the Interface of previous versioning with constraints: HelloItf.java

Extending this model to support non-native tasks is a problem since the different implementations are enclosed within the *Method* annotation (and the new *Binary*, *MPI* and *OmpSs* annotations are defined outside the *Method* annotation). Moreover, to check that the number of entries of the *MultiConstraints* annotation is the same than the number of declared implementations becomes expensive. Hence, our new implementation only allows to declare one implementation per method annotation, moves the *Constraints* annotation inside each task annotation to define the specific implementation requirements and allows the user to define one single *Constraint* clause to define the common requirements. Figure 5.18 shows how the previous example will be specified with the new COMPSs Annotations.

```

1 public interface NewHelloItf {
3     @Method(declaringClass = "Impl1",
4             constraints = @Constraints(memorySize = 1))
5     @Method(declaringClass = "Impl2",
6             constraints = @Constraints(memorySize = 2))
7     @Constraints(computingUnits = 1)
8     void sayHello();
9 }

```

FIGURE 5.18: Example of the new Annotation Interface: NewHelloItf.java

Notice that this new annotation syntax allows us to add the non-native tasks as different implementations easily. For example, Figure 5.19 adds a new binary implementation to the previous example (now, the *sayHello* method has three different implementations).

```

1 public interface NewHelloItf {
3     @Method(declaringClass = "Impl1",
4             constraints = @Constraints(memorySize = 1))
5     @Method(declaringClass = "Impl2",
6             constraints = @Constraints(memorySize = 2))
7     @Binary(binary = "${BINARY}")
8     @Constraints(computingUnits = 1)
9     void sayHello();
}

```

FIGURE 5.19: Extended example of the new Annotation Interface: NewHelloItf.java

#### 5.1.4 SchedulerHints task annotation

Since the previous modifications were already making the new COMPSs Annotations not backward compatible, we have profited to add a specific annotation for the user to define hints for the COMPSs Scheduler. In the current state, the new *SchedulerHints* annotation only supports the two fields shown in Table 5.1

Field	Type	Description
isDistributed	boolean	Forces the scheduler to schedule the instances of the task in a round-robin manner between nodes
isReplicated	boolean	Replicates the task execution in all the available nodes

TABLE 5.1: Available fields for the *SchedulerHints* annotation

We are conscient that these annotations may break, somehow, the programming model because users must be unaware of the underlying infrastructure and because the *isReplicated* field makes the parallel behavior different to the sequential behavior. However, these options are very useful for advanced users (for example, in cases where an initialization task must be executed in all the computational nodes) and may be, in the near future, very useful for scheduler enhancements (such as annotating map-reduce tasks).



Figure 5.20 illustrates how to use these annotations in the interface.

```

public interface MainItf {
2
    @Method(declaringClass = "MainImpl")
    @SchedulerHints(isReplicated = true)
    void replicatedTask();
6
    @Method(declaringClass = "MainImpl")
    @SchedulerHints(isDistributed = true)
    void distributedTask();
8
10 }

```

FIGURE 5.20: Example of an Interface with SchedulerHints: MainItf.java

### 5.1.5 New stream parameter annotation

Once the task annotations for non-native tasks was defined, we had to implement a way to communicate the Java application with the non-native tasks' execution. When executing standalone binaries, MPI processes or OmpSs processes the exit value of the processes is used as the return value. Thus, we have decided that the COMPSs non-native tasks must use the exit value of their internal binary as the return value of the task. In this sense, we have allowed the users to capture this value by defining the return type of the non-native task as an *int* (for implicit synchronization), as an *Integer* (for post-access synchronization) or to forget it (declaring the function as *void*). Figure 5.21 shows an Interface example of the three return types.

```

public interface MainItf {
2
    @Binary(binary = "${BINARY}")
    int binaryTask1();
4
    @Binary(binary = "${BINARY}")
    Integer binaryTask2();
6
    @Binary(binary = "${BINARY}")
    void binaryTask3();
8
10 }

```

FIGURE 5.21: Example of the different return types of the non-native tasks

However, the users not only need the process exit value to work with this kind of applications but need to set the Standard Input (*stdin*) and capture the Standard Output (*stdout*) and Error (*stderr*). For this purpose, we have created a new parameter annotation, *stream*, that allows the users to set some parameters as I/O streams for the non-native tasks. Stream parameters are not passed directly to the binary command but rather they are set as *stdin*, *stdout* or *stderr* of the binary process. Since this kind of redirection is restricted to files in *LINUX* Operating Systems, we have decided to keep the same restrictions to the annotation. Consequently, **all stream parameters must be files**.

Figure 5.22 shows the Interface of an application with two tasks that have a normal parameter (the first one, that will be sent directly to the binary execution), a file parameter to be used as *stdin* of the process, a file parameter to be used as *stdout* and a last file parameter to be used as *stderr*. The difference between *task1* and *task2* in this example is that the first

task will overwrite the *fileOut* and *fileErr* content (since the files are opened in *write* mode), and the second task will append the *fileOut* and *fileErr* content at the end of the file (since the files are opened in *append* mode).

```

1 public interface StreamItf {
3     @Binary(binary = "${BINARY}")
4     Integer task1(
5         @Parameter(type = Type.STRING, direction = Direction.IN) String normalParameter,
6         @Parameter(type = Type.FILE, direction = Direction.IN, stream = Stream.STDIN)
7         String fileIn,
8         @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
9         String fileOut,
10        @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
11        String fileErr
12    );
13
14    @Binary(binary = "${BINARY}")
15    Integer task2(
16        @Parameter(type = Type.STRING, direction = Direction.IN) String normalParameter,
17        @Parameter(type = Type.FILE, direction = Direction.IN, stream = Stream.STDIN)
18        String fileIn,
19        @Parameter(type = Type.FILE, direction = Direction.INOUT, stream = Stream.STDOUT)
20        String fileOut,
21        @Parameter(type = Type.FILE, direction = Direction.INOUT, stream = Stream.STDERR)
22        String fileErr
23    );
24 }

```

FIGURE 5.22: Example of the different stream annotations for non-native tasks

To summarize the last information retrieved from this example, Table 5.2 show the available modes for each stream type.

Type	Stream	Direction	Description
FILE	Stream.STDIN	Direction.IN	Sets the process <i>stdIn</i> . The file is opened in <i>read</i> mode
FILE	Stream.STDOUT	Direction.OUT	Sets the process <i>stdOut</i> . The file is opened in <i>write</i> mode
FILE	Stream.STDOUT	Direction.INOUT	Sets the process <i>stdOut</i> . The file is opened in <i>append</i> mode
FILE	Stream.STDERR	Direction.OUT	Sets the process <i>stdErr</i> . The file is opened in <i>write</i> mode
FILE	Stream.STDERR	Direction.INOUT	Sets the process <i>stdErr</i> . The file is opened in <i>append</i> mode

TABLE 5.2: Available stream types with their valid directions and execution behaviour

### 5.1.6 New prefix parameter annotation

Before going in-depth in this section, we must explain that COMPSs builds the task data dependencies graph taking into account the parameters annotated in the application Interface. Analyzing several binaries, we have found out that a non-negligible part of them use prefixes for each parameter. The prefixes used by binaries can be divided into two types:

- **Separated Prefix** A prefix that is written separately before the parameter value. This type of prefixes are of the form:

*./binary - param1 value - -param2 value - k value*

In fact, there is not a strong need that the parameter prefix starts with a dash but its the common behavior for Linux binaries.

- **Joint Prefix** A prefix that is written with the parameter value without beeing separated or with a separation character that it is not an empty space. This types of prefix vary a lot but are of the form:

*./binary - pValue - q = Value - -r = Value s = Value*

The separated prefixes do not represent a problem for the COMPSs programming model since they can be defined as a standalone string parameter that is finally passed to the binary. However, the joint prefixes do represent a problem for COMPSs since the users must prepend the prefix to the parameter, breaking the data dependencies between the tasks. For the sake of clarity, consider the two tasks shown in Figure 5.23 and the main code shown in Figure 5.24. Since the second task requires a joint prefix, when calling it from the main code the users must modify its value and prepend the prefix to the *fileName* variable. This string modification causes a synchronization in the appliaction's main code instead of creating a data dependency between the two tasks.

```

public interface ExampleItf {
2
    // Must execute: ./tmp/bin1 fileName
    @Binary(binary = "/tmp/bin1")
    Integer task1(
6        @Parameter(type = Type.FILE, direction = Direction.INOUT) String fileName
    );
8
    // Must execute: ./tmp/bin2 --file=fileName
    @Binary(binary = "/tmp/bin2")
    Integer task2(
12        @Parameter(type = Type.FILE, direction = Direction.INOUT) String fileName
    );
14 }

```

FIGURE 5.23: Binary Tasks example for joint prefixes

```

public static void main(String[] args) {
2    String fileName = "/tmp/file";
    BINARY.task1(fileName);
4    BINARY.task2("--file=" + fileName);
}

```

FIGURE 5.24: Main code example for joint prefixes

Consequently, for this second type of prefixes, we have created a new parameter annotation *prefix* that allows the users to define the prefix separately to the parameter value and its prepended to the parameter value just before the binary execution. This modification allows COMPSs to handle the data dependencies between parameters (since prefixes are immutable strings that do not define data dependencies) and allows the binaries to receive the parameter prefixes and its value together as a single parameter.

```

1 public interface StreamItf {
3     @Binary(binary = "binaryExample")
4     void task1(
5         @Parameter(type = Type.STRING, direction = Direction.IN)    String pPrefix,
6         @Parameter(type = Type.FILE,    direction = Direction.IN)    String fileIn,
7         @Parameter(type = Type.STRING, direction = Direction.IN)    String qPrefix,
8         @Parameter(type = Type.FILE,    direction = Direction.INOUT) String fileInOut,
9         @Parameter(type = Type.STRING, direction = Direction.IN)    String kPrefix,
10        @Parameter(type = Type.INT,     direction = Direction.IN)     int k
11    );
13    @Binary(binary = "binaryExample")
14    void task2(
15        @Parameter(type = Type.FILE,    direction = Direction.IN,    prefix = "-p=")
16        String fileIn,
17        @Parameter(type = Type.FILE,    direction = Direction.INOUT, prefix = "--q=")
18        String fileInOut,
19        @Parameter(type = Type.INT,     direction = Direction.IN,    prefix = "k")
20        int k
21    );
23    @Binary(binary = "binaryExample")
24    void task3(
25        @Parameter(type = Type.STRING, direction = Direction.IN)
26        String pPrefix,
27        @Parameter(type = Type.FILE,    direction = Direction.IN)
28        String fileIn,
29        @Parameter(type = Type.FILE,    direction = Direction.INOUT, prefix = "--q=")
30        String fileInOut,
31        @Parameter(type = Type.INT,     direction = Direction.IN,    prefix = "k")
32        int k
33    );
34 }

```

FIGURE 5.25: Interface example of an application with prefixes

Figure 5.25 shows the Interface of an application with three binary tasks, Figure 5.26 shows the main code calls to these tasks and Figure 5.27 shows the final binary command that is executed in the task. Notice that the first task, *task1*, only uses separated prefixes; the second task, *task2*, uses only joint prefixes and the third task, *task3*, is a hybrid example of both separated and joint prefixes.

```

1 public class BinaryPrefixesExample {
2
3     public static void main(String[] args)
4     {
5         String file1 = "file1.in"
6         String file2 = "file2.inout"
7         int kValue = 10;
8
9         // Launch task 1
10        task1("-p", file1, "--q", file2, "k", kValue);
11
12        // Launch task 2
13        task2(file1, file2, kValue);
14
15        // Launch task 3
16        task3("-p", file1, file2, kValue);
17    }
18 }

```

FIGURE 5.26: Example of the main code calls to tasks with prefixes

```
1 # TASK 1
  ./binaryExample -p file1.in --q file2.inout k 10
3
4 # TASK 2
5 ./binaryExample -p=file1.in --q=file2.inout k10
6
7 # TASK 3
  ./binaryExample -p file1.in --q=file2.inout k10
```

FIGURE 5.27: Example of the command executed inside each task using prefixes

## 5.2 Scheduling modifications

Before explaining the scheduler modifications, we must define how COMPSs handles the task creation, scheduling and execution. The COMPSs Runtime instruments the application's main code looking for invocations to the methods defined as tasks in the application interface. When these methods are detected, COMPSs creates a task that is submitted to the Task Analyzer component and substitutes the method call by an *executeTask()* call. When the Task Analyzer receives a new task, it computes its data dependencies and submits it to the Task Scheduler. Next, the Task Scheduler creates an Execution Action associated with the task and adds it to the execution queue. Eventually, the Execution Action will be scheduled and launched (this mechanism requires the task to be data-free and to have enough free resources to fulfill the task constraints). When the Execution Action is launched, a Job is created to monitor the task execution. This job includes the transfer of the job definition and all the input data to the target COMPSs Worker, the real task execution in the worker and the transfer of the output data back to the COMPSs Master. Once the job is completed, its data dependent Execution Actions are released (if any), and the job is destroyed (or, depending on the debug level, stored for post-mortem analysis).

### 5.2.1 Treatment of non-native tasks

Non-native tasks only represent a new way of executing tasks in the COMPSs Worker. Although the data structure that contains its information (*BinaryImplementation*, *MPIImplementation*, *OmpSsImplementation*) is quite different than the one storing Method / Service tasks (*MethodImplementation* and *ServiceImplementation* respectively), all of them inherit a common super-structure (*Implementation*) that allows the COMPSs Master to treat any type of task in the same way.

Consequently, the Scheduler component is also independent of the task execution and, thus, no modification has been added to enable the execution of non-native tasks. However, we must emphasize that although the execution of MPI tasks itself has not caused any modification (because it also extends from the same interface), the fact of using more than one computational node did.

### 5.2.2 Multi-node execution actions

That being said, the first design choice to enable multi-node task executions has been to associate several execution actions to the same task. This mapping allows the Scheduler to treat the data-dependencies and the resource consumption as it was done before. However, the execution actions associated with the same task must have different behaviors during

the execution phase because only one of the actions must really launch the job.

Consequently, we have extended the *ExecutionAction* in a *MultiNodeExecutionAction* class that is only used when a task requires more than one computing node (otherwise the previous *ExecutionAction* implementation is used). When the task scheduler receives a new multi-node task (*ExecuteTaskRequest*) it creates a new *MultiNodeGroup* instance and  $N$  *MultiNodeExecutionAction* instances (being  $N$  the number of nodes requested by the task). The *MultiNodeGroup* instance is shared among all the actions assigned to the same task execution, and it handles the actions' id within the group. More in-depth, when the *MultiNodeExecutionActions* are created the *MultiNodeGroup* assigns a nullable identifier to all of them. Once the actions are scheduled and launched, the *MultiNodeGroup* assigns a unique valid identifier between 1 and  $N$ . This action identifier is used during the action execution to act as an execution slave node (when the assigned identifier is different to 1) or to act as an execution master node (when the assigned identifier is 1). When the *MultiNodeExecutionAction* is identified as a slave, it no longer triggers a job execution, but rather reserves the requested resources and waits for its master action to complete. When the *MultiNodeExecutionAction* is identified as a master, it retrieves all the hostnames of its slave actions (for the MPI command) and behaves as a normal *ExecutionAction* (launches a job to monitor the input data transfers, the real task execution on the node and the output data transfers).

On the one hand, Figure 5.28 shows an example of the normal process. A task  $T1$  requiring 1 node (normal task) is submitted to the scheduler through the *ExecuteTaskRequest* request. The request is then processed and an *ExecutionAction* is created as it was done before. Eventually, the action is scheduled, launched and finally executed, creating a new job that will monitor the task execution in the target node.

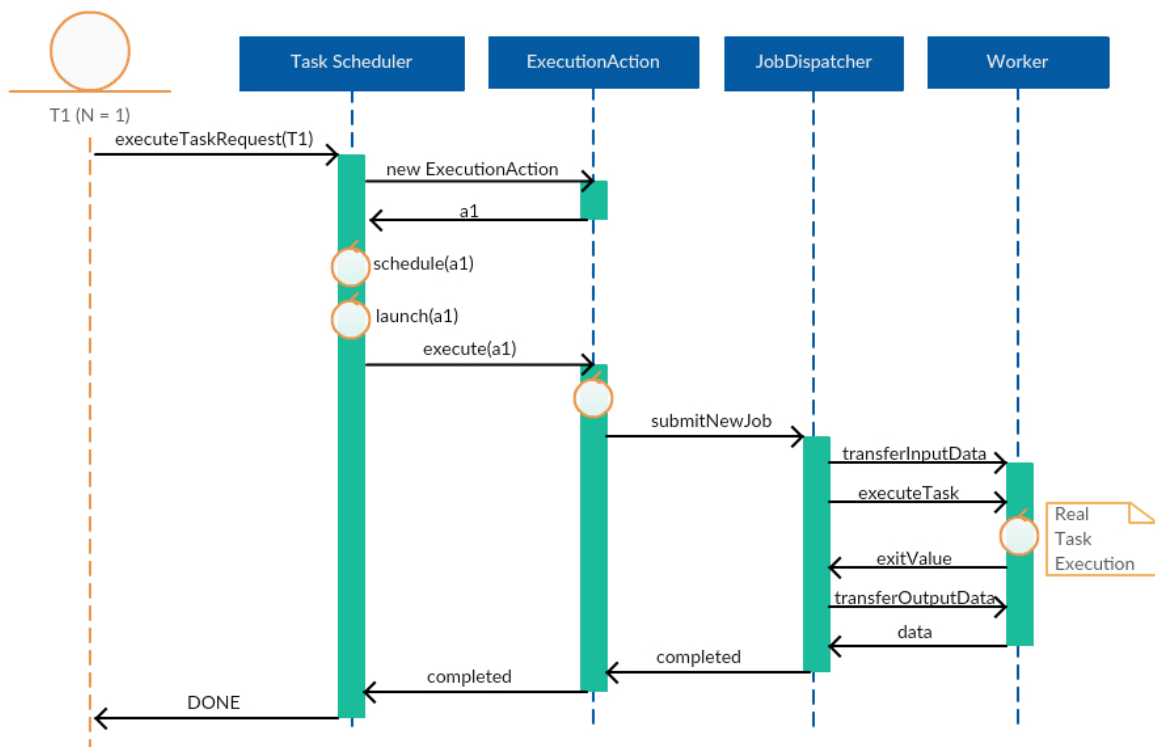


FIGURE 5.28: Example of a single node task flow

On the other hand, Figure 5.29 shows an example of the Multi-Node process. A task  $T2$  requiring 3 nodes (multi-node task) is submitted to the scheduler through the same *ExecutionTaskRequest* request. The request is then processed: a new action group (lets say  $g1$ ) is created (a new instance of the *MultiNodeGroup*) and 3 *MultiNodeAction* instances (lets say  $a1$ ,  $a2$  and  $a3$ ) are created. The action group  $g1$  is shared among all the three actions and assigns a nullable action identifier to all of them.

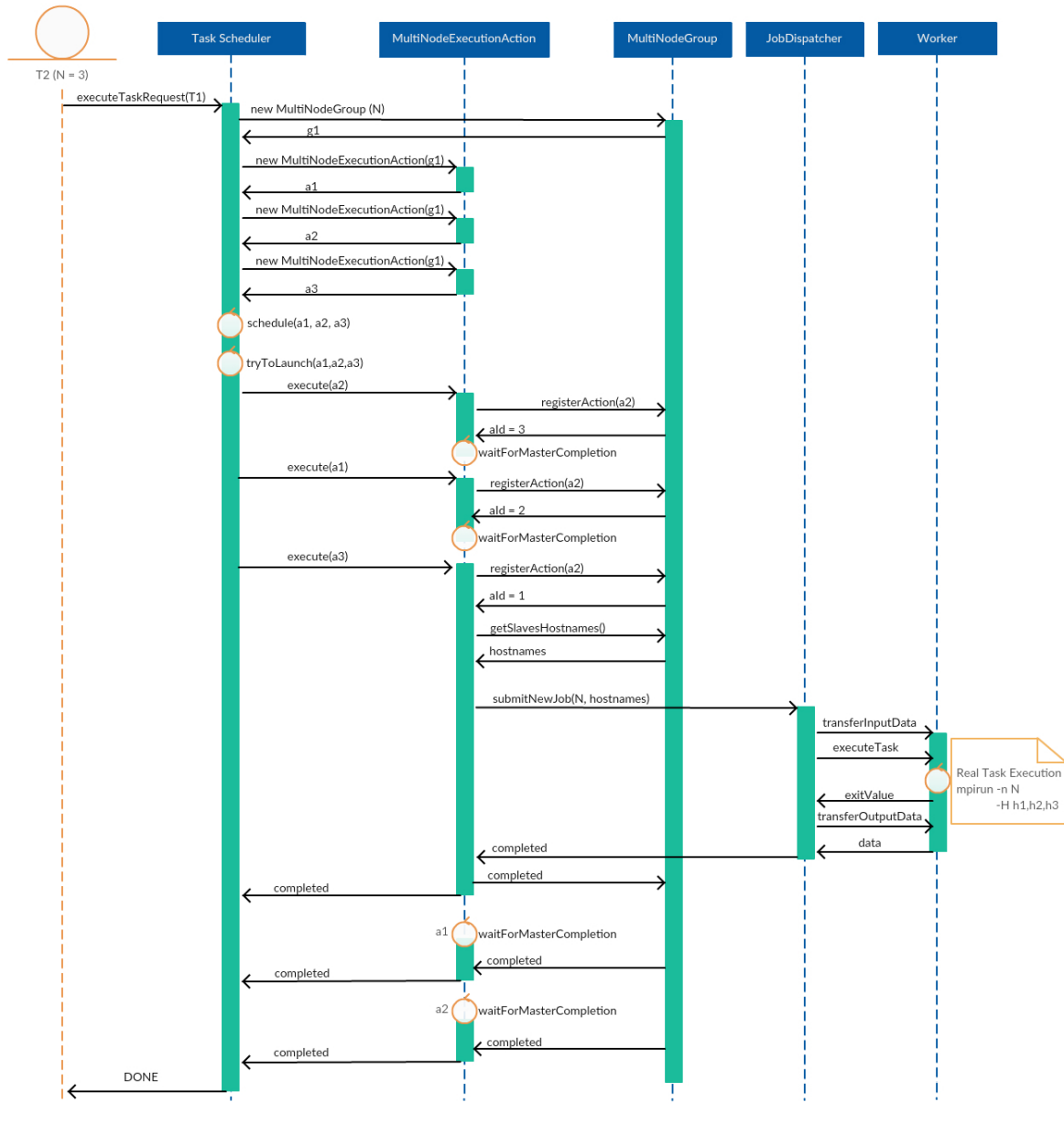


FIGURE 5.29: Example of a multi-node task flow

Eventually,  $a2$  is scheduled, launched and finally executed. On the execution phase, the action asks for an action identifier and the action group  $g1$  assigns it an `actionId = 3` (because the group size is 3 and no action has previously requested an identifier). Since the action identifier classifies  $a2$  as a slave action, the execute phase only reserves the task constraints and waits for the master action completion.

Eventually,  $a1$  is also scheduled, launched and finally executed. Following the same process than the previous action,  $a1$  is granted with `actionId = 2` (because the group size is 3 and

only one action has previously requested an identifier). Since the action  $a1$  is also classified as a slave, it reserves the task constraints and waits for the master action completion.

Finally,  $a3$  will also be scheduled, launched and finally executed. In this case, the action group assigns it an  $actionId = 1$ . Since it is the last action, it is now identified as master and during its execution phase it retrieves the hostnames of the resources assigned to all the actions inside the  $g1$  group (lets say,  $h1$  for  $a1$  and  $h2$  for  $a2$ ) and launches the execution job. The job will be then executed (lets say that the host assigned to this action  $a3$  is  $h3$ ) monitoring the input data transfers, performing the real task execution (for example, calling the MPI command inside the host  $h3$  with 3 nodes  $h1$ ,  $h2$  and  $h3$ ) and retrieving back the output data from  $h3$ .

Once the job is completed, the action  $a3$  is marked as completed (freeing all the reserved resources) and, then, it triggers its completion to all the slave actions registered in the group  $g1$ . Consequently,  $a1$  and  $a2$  are also marked as completed (and its resources are also freed). When all the actions within the group are marked as completed, the task is registered as *DONE* and follows the usual process: frees its data dependent tasks and it is stored for post-mortem analysis.

### 5.2.3 Treatment of SchedulerHints

Currently there are only two *SchedulerHints* available in COMPSs annotation: *isDistributed* and *isReplicated*. Both annotations are attached to the task definition and are treated when the *ExecuteTaskRequest* request is served.

When the *isDistributed* annotation is enabled, the request checks how many tasks of the same type have been already scheduled to each available worker. Then, it chooses the worker that has executed **less** tasks of the same type and forces the Execution Action to be scheduled to the selected worker. Notice that the computational cost of distributing a task in a Round Robin manner among the available workers is proportional to the number of available workers. To maintain consistency, when a resource is chosen as the target worker of an Execution Action during the *schedule* phase, a task counter in the target worker is increased. In this sense, notice that the memory cost is increased by Equation 5.1 since each available worker stores a list of counters of size equal to the number of different tasks registered in the Application Interface.

$$numWorkers \cdot numTypesTasks \cdot sizeof(int) \quad (5.1)$$

When the *isReplicated* annotation is enabled, the request creates one Execution Action per available worker and forces the actions to be executed in the selected workers. The task is then considered as *DONE* when all its Execution Actions are marked as *completed*. For this purpose, the task stores an execution counter initialized to the number of available workers that only releases the task (and its data dependent actions) when reaching zero.

As a final note, we highlight that the Scheduling Hints are evaluated during the scheduling phase. That means that the workers considered in both cases must be functional when the task scheduling is being treated (not when the task is really executed).

## 5.3 Worker enhancements

The Communication layer abstracts the Master node from the specific Communication Adaptors and thus, from the underlying infrastructure. However, the worker processes spawned by this layer are dependent on each Adaptor implementation. Currently, COMPSs supports the *NIO* and the *GAT* Communication Adaptors.



On the one hand, *GAT* Adaptor is built on top on the Java Grid Application Toolkit (JavaGAT) [48] which relies on the SSH connection between nodes. During the application execution, the Runtime spawns a new worker process per task execution. More specifically, when a task must be executed, the *GAT* Communication Adaptor creates a *GAT* Job, sends the job and the required data through SSH to the worker's resource, starts the worker process, executes the task itself, closes the worker process, and retrieves the job status, the job's log files, and the required output data. The *worker.sh* script orchestrates all the processes and launches a language dependant script for the real task execution (*GATWorker.java* for Java, *worker.py* for Python and *Worker* for C/C++). Although the implementation suffers from some performance overheads (because the overhead of spawning a new process on each task execution becomes non-negligible for small duration tasks), it provides a high connectivity interface since it only requires the SSH port to be opened.

On the other hand, *NIO* Adaptor is a more sophisticated implementation based on Java New I/O (NIO) library [41]. This adaptor spawns a persistent Java Worker Process per resource, rather than one per task execution, and the communication between Master and Workers is then made through Sockets. Hence, this Adaptor provides better performance than the *GAT* Adaptor but requires extra open ports between the available resources. Furthermore, the Worker processes persist during the full execution of the application, what also lets us have an object cache per worker, data communications between workers (rather than handling all the data in the Master resource) and thread binding mechanisms to map threads to specific cores of the machine. Finally, for the task execution, each worker has several *Executor* threads that can execute natively Java applications, or Python and C/C++ applications using a *ProcessBuilder*.

### 5.3.1 Invokers

To enable the execution of non-native tasks for any Communication Adaptor we have implemented a *GenericInvoker* class that provides an API for executing standard, MPI, and OmpSs binaries. This API is built on top of a *BinaryRunner* class that spawns, runs and monitors the execution of any binary command.

More specifically, the *BinaryRunner* class has two methods. Firstly, *createCMDParametersFromValues* serializes the received parameters to construct the binary arguments. This method is also in charge of processing the *Stream* annotations and redirecting the StdIn, StdOut, and StdErr when required. Secondly, *executeCMD* executes the received binary command (with all its parameters), monitors its execution and, finally, returns the exit value of the process.

On the other hand, the *GenericInvoker* class provide three functions: *invokeBinaryMethod*, *invokeMPIMethod*, and *invokeOmpSsMethod* to invoke respectively standard, MPI and OmpSs binaries. The three methods receive the binary path and the argument values, construct and execute the command by calling the *BinaryRunner* functions and return the exit value of the binary execution.

In a second step we have adapted each of the Communication Adaptor (*GAT* and *NIO*) to call this *GenericInvoker* when needed. In both cases, we have substituted the normal task execution by a *switch-case* that selects the required invoker considering the task's implementation type.

### 5.3.2 External executors enhancement

Currently, the PyCOMPSs binding is gaining relevance because Data Analytic Workflows can be easily designed in Python and thus, translated to PyCOMPSs. During the development of this project, we have found out that the COMPSs Workers loose performance because of the creation of the *ProcessBuilder* used to execute Python, C, and C++ tasks.

To solve this problem, we have completely redesigned the way the COMPSs Worker launches a Python, C or C++ task. First of all, during the Worker initialization, we use the *ProcessBuilder* to launch a single BASH script that creates  $N$  input pipes,  $N$  output pipes and  $N$  processes (being  $N$  the number of available *Executors* at the given Worker). These BASH processes persist until the Worker is stopped and use the pipes to communicate with the Java *Executors*. In fact, each Java *Executor* stores a pair of pipes so that the *Executors* are mapped one to one (see Figure 5.30).

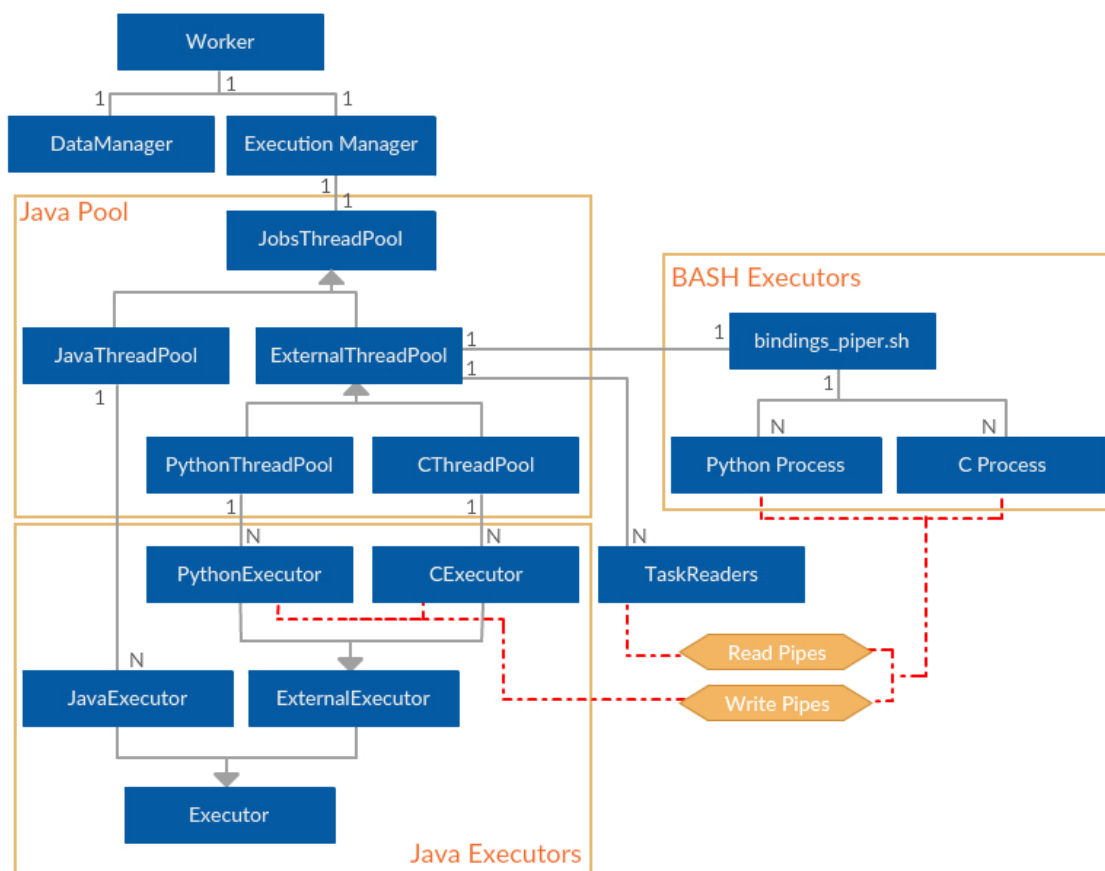


FIGURE 5.30: New structure of the COMPSs Worker Executors

When a task arrives at the Worker, it is processed by a Java *Executor* thread. Assuming that the task is a Python, C or C++ task, the Java *Executor* sends a command through its assigned input pipe to the BASH *Executor* with the task definition and waits for a completion message on the assigned output pipe. The BASH *Executor* then receives the task command, executes the task and sends its result back to the Java *Executor* through the output pipe. Notice that both executor threads are never active at the same time since, when one Java is processing the task the BASH process is listening, and when the BASH process is executing the task the Java process is listening.

At the end of the application execution, when the Worker is stopped, all the Java *Executors* send a *QUIT* command through the input pipe to kill its assigned BASH processes and then exit. To double check the shutdown process, the initial BASH script (spawned with the *ProcessBuilder* at the Worker creation) is killed with a bash *TRAP* to kill any remaining pipe or process.

Finally, Figure 5.31 shows the execution time in the y-axis versus the number of *Executors* ( $N$ ) in the x-axis for an implementation with *ProcessBuilders* (blue) and an implementation with pipes (red for the total time including the spawn and the destruction of the first BASH script, and yellow for the task execution time). We must highlight that the implementation using pipes speeds-up significantly the pre and post actions that must be done in every task execution, getting us to a lower overhead when using the COMPSs Bindings.

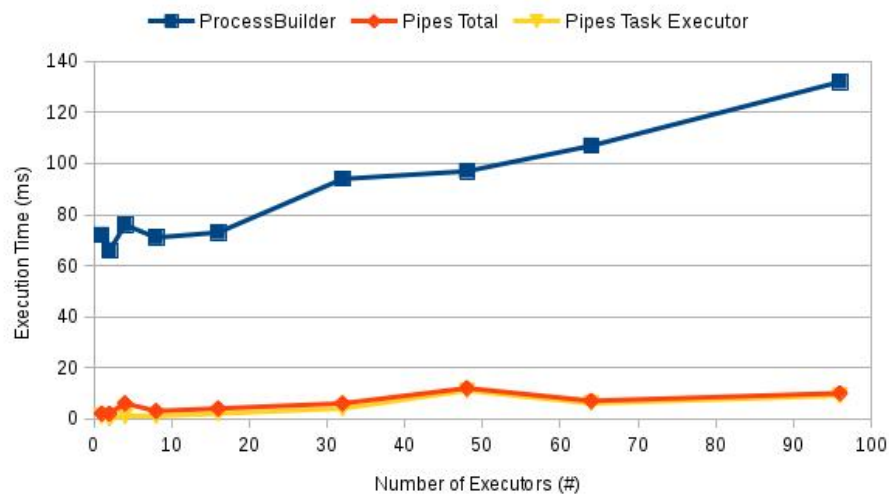


FIGURE 5.31: Execution time versus number of *ProcessBuilders* or Pipes



## Chapter 6

# Results and evaluation

### 6.1 Proofs of concept

#### 6.1.1 BLAST

##### 6.1.1.1 Application description

Basic Local Alignment Search Tool (BLAST) [6] is an algorithm to find regions of similarity between primary biological sequences. The program compares a nucleotide or protein (known as *query*) to a sequences' database and identifies sequences that resemble the *query* sequence above a certain threshold.

The COMPSs implementation of BLAST splits the *query* sequence on smaller *fragments*, comparing each fragment against the database and merging up the obtained results. Using different execution arguments, the users can select the *query* sequence, the number of fragments and the target database (Figure 6.1 provides the complete list of execution arguments).

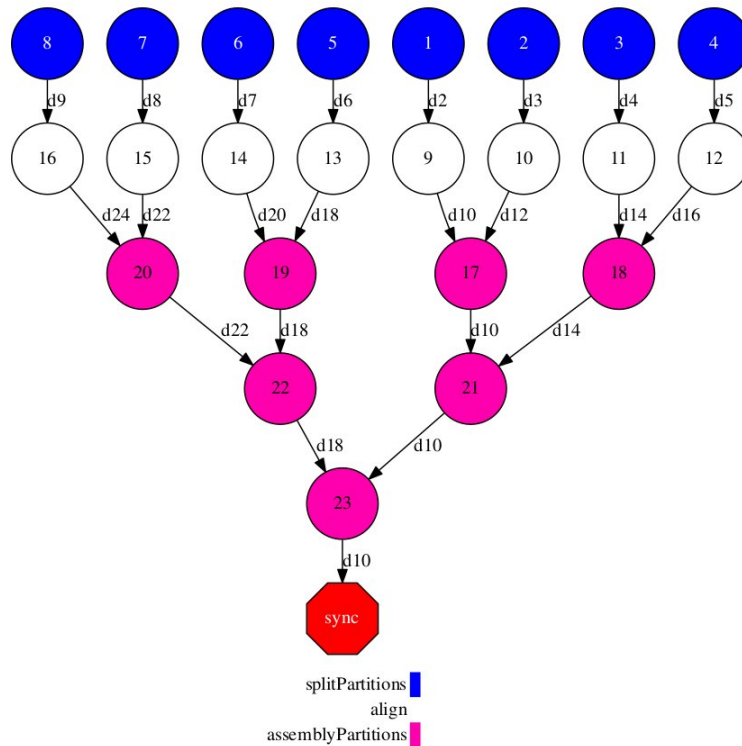
```
blast.BLAST <debug> <database> <query> <nFrgs> <tmpDir> <outputFile> <blastBinaryArgs>
2  - debug           Enables and disables the debug mode
   - database       The sequences' database path to compare with
4  - query          The path where the query sequence is stored
   - nFrgs         The number of fragments used to divide the sequence (N)
6  - tmpDir        A temporal directory for intermediate results
   - outputFile    The final output file with the sequence matches
8  - blastBinaryArgs The extra arguments to pass to the BLAST binary
```

FIGURE 6.1: Execution arguments of the COMPSs BLAST application

Regarding the code, it is divided in three main blocks:

- **Split:** Splits the *query* sequence in **N** fragments
- **Alignment:** Compares each fragment of the *query* sequence against the database invoking the BLAST binary
- **Assembly:** Merges all the intermediate files into a single file to produce the final result

The BLAST interface provides three types of tasks, one per block. Firstly, the input file is splitted by using **N** *splitPartitions* tasks (being **N** the number of fragments). Next, each fragment is processed using the *align* task. Finally, the partial results are merged by using **N-1** *assemblyPartitions* tasks. Figure 6.2 shows an example with **N = 8**.

FIGURE 6.2: Example of BLAST execution with  $N = 8$ 

### 6.1.1.2 Purpose

This proof of concept must demonstrate that our implementation in COMPSs eases the execution of binary files without degrading its performance. Moreover, the users must have enough flexibility to execute the binaries with different parameters, to synchronize (or not) the exit value, and to capture (or not) the *StdOut* and *StdErr*.

### 6.1.1.3 Evaluation

When porting the BLAST implementation with COMPSs to the new *Binary* annotation, we have left intact the execution parameters and the application behavior (blocks and number of tasks spawned per block). The only modifications that we have introduced are:

- ***align* task implementation:** A new dummy implementation of the *align* task is defined in the *binary.BINARY* file (see Figure 6.3). This implementation replaces the old one defined in the *BLASTImpl* file (see Figure 6.4), which is totally removed.

```

public static Integer align(String pFlag, String pMode, String dFlag, String database,
2     String iFlag, String partitionFile, String oFlag,
3     String partitionOutput, String extraCMDArgs) {
4
5     return -1;
6 }

```

FIGURE 6.3: COMPSs BLAST application: new *align* task implementation

```

public static void align(String databasePath, String partitionFile, String partitionOutput,
2     String blastBinary, String extraCMDArgs) throws BlastException {

4     Long startAlignment = System.currentTimeMillis();
    String cmd = blastBinary
6         + " -p blastx"
          + " -d " + databasePath
8         + " -i " + partitionFile
          + " -o " + partitionOutput
10        + " " + extraCMDArgs;
    Process p = Runtime.getRuntime().exec(cmd);
12    int exitValue = p.waitFor();
    Long alignmentTime = (System.currentTimeMillis() - startAlignment) / 1_000;
14    System.out.println("Alignment time " + alignmentTime + " s");

16    if (exitValue != 0) {
        try (BufferedInputStream bisErr = new BufferedInputStream(p.getErrorStream());
18             BufferedOutputStream bosErr = new BufferedOutputStream(
                new FileOutputStream(partitionFile + ".err"));) {
20            byte[] b = new byte[1_024];
            int read;
22            while ((read = bisErr.read(b)) >= 0) {
                bosErr.write(b, 0, read);
24            }
            bisErr.close();
26            bosErr.close();
        } catch (IOException ioe) {
28            String msg = "ERROR: Cannot retrieve CMD error content";
            System.err.println(msg);
30            throw new BlastException(msg, ioe);
        }

32        String msg = "Error executing Blast job, exit value is: " + exitValue;
34        System.err.println(msg);
        throw new BlastException(msg);
36    }
}

```

---

FIGURE 6.4: COMPSs BLAST application: old *align* task implementation

- ***align* task call:** Adapt the main code to the new *align* task call. The new task no longer wraps the binary execution, but rather spawns the *blast* binary directly. Thus, we must add the binary arguments explicitly and recover the exit value of the task. Figure 6.5 shows the new task call loop and Figure 6.6 shows the old one.

```

1 private static void alignSequences() throws BLASTException {
    final String pFlag = "-p";
3    final String pMode = "blastx";
    final String dFlag = "-d";
5    final String iFlag = "-i";
    final String oFlag = "-o";
7    int numAligns = BLAST.partialInputs.size();
    Integer[] exitValues = new Integer[numAligns];
9    for (int i = 0; i < numAligns; i++) {
        exitValues[i] = BINARY.align(pFlag, pMode, dFlag, BLAST.databasePath, iFlag,
11        BLAST.partInputs.get(i), oFlag, BLAST.partOutputs.get(i), BLAST.cmdArgs);
    }
13 }

```

---

FIGURE 6.5: COMPSs BLAST application: new *align* task call

```

1 private static void alignSequences() throws BLASTException {
2     for (int i = 0; i < partialInputs.size(); i++) {
3         BLASTImpl.align(databaseName, partialInputs.get(i), partialOutputs.get(i),
4             blastBinary, commandArgs);
5     }
6 }

```

FIGURE 6.6: COMPSs BLAST application: old *align* task call

- ***align* interface' annotation:** We delete the previous *align* method and add a new task with the Binary annotation. Figure 6.7 shows the new annotation and 6.8 shows the previous one.

```

@Binary(binary = "${BLAST_BINARY}")
2 Integer align(
3     @Parameter(type = Type.STRING, direction = Direction.IN) String pFlag,
4     @Parameter(type = Type.STRING, direction = Direction.IN) String pMode,
5     @Parameter(type = Type.STRING, direction = Direction.IN) String dFlag,
6     @Parameter(type = Type.STRING, direction = Direction.IN) String database,
7     @Parameter(type = Type.STRING, direction = Direction.IN) String iFlag,
8     @Parameter(type = Type.FILE, direction = Direction.IN) String partitionFile,
9     @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
10    @Parameter(type = Type.FILE, direction = Direction.OUT) String partitionOutput,
11    @Parameter(type = Type.STRING, direction = Direction.IN) String extraCMDArgs
12 );

```

FIGURE 6.7: COMPSs BLAST application: new *align*'s interface annotation

```

@Method(declaringClass = "blast.BLASTImpl")
2 void align(
3     @Parameter(type = Type.STRING, direction = Direction.IN) String databasePath,
4     @Parameter(type = Type.FILE, direction = Direction.IN) String partitionFile,
5     @Parameter(type = Type.FILE, direction = Direction.OUT) String partitionOutput,
6     @Parameter(type = Type.STRING, direction = Direction.IN) String blastBinary,
7     @Parameter(type = Type.STRING, direction = Direction.IN) String extraCMDArgs
8 );

```

FIGURE 6.8: COMPSs BLAST application: old *align*'s interface annotation

Although the previous figures highlight the major modifications, Appendix A provides the full code of the BLAST implementation using the new COMPSs annotations. Notice that the annotation and the task call are more complex than the previous version (since there are more parameters), but the task's implementation is totally suppressed. This means, the users no longer need to create a new *ProcessBuilder* from the commands, spawn a process, wait for its completion, retrieve the exit value, and read the process' output and error streams. Thus, we consider that these annotations provide a significant advantage regarding programmability.

Concerning performance, we have measured the execution time of the *align* task using both implementations. On the one hand, for the *old* implementation, we have measured the



total task time inside the Runtime and the time spent on the real binary execution in the task code inside the application. On the other hand, for the *new* implementation, we have measured the total task time and the time spent on the real binary execution instrumenting the Runtime.

The experiments have been run on the *MinoTauro* [11] machine that is a heterogeneous cluster hosted at the Barcelona Supercomputing Center (BSC) that has 61 Bull B505 blades with 2 Intel E5649 (6-Core) processors at 2.53 GHz, 24 GB of Main memory, 250 GB SSD and 2 Infiniband QDR (40 Gbit each), and 39 Bullx R421-E4 servers with 2 Intel Xeon E5-2630 v3 8-core processors at 2.4 GHz, 128 Gb of main memory, 120 Gb SSD and 1 PCIe 3.0 x8 8GT/s. For the experimentation, we have only requested two nodes (one COMPSs Master and one COMPSs Worker) of the first type, and we have only launched one task at a time (1 maximum task per node) to avoid distortion of the task execution time caused by the execution of other processes. Moreover, we have launched 10 BLAST executions with the parameters shown in Table 6.1.

Parameter	Value
<b>debug</b>	disabled
<b>database</b>	swissprot
<b>query</b>	sargasso.fasta (43 Kb)
<b>nFrag</b>	8
<b>tmpDir</b>	/scratch/tmp/
<b>Output file</b>	output.txt
<b>Binary Arguments</b>	None

TABLE 6.1: Execution parameters of the BLAST application

Notice that the total number of executed tasks is 80, since the *query* sequence is divided in 8 fragments (see Equation 6.1).

$$1 \frac{\text{align task}}{\text{fragment}} \cdot 8 \frac{\text{fragment}}{\text{execution}} \cdot 10 \text{ execution} = 80 \text{ align task} \quad (6.1)$$

Table 6.2 shows the minimum, the maximum and the mean execution times of the *align* task when running **10 times** the *BLAST* application under the aforementioned conditions. In terms of absolute values, the difference between the two implementations is negligible for the total task execution time (less than 1%) and for the binary execution time (less than 0.5%). Regarding the overhead, the old implementation adds 16.25 ms on average to the binary execution and the new implementation adds 15.50 ms on average. Hence, as expected, the overhead introduced by the new implementation can be considered the same than the old implementation, and there is no performance degradation when using the new Binary annotation.

	NEW		OLD	
	Total task	Binary	Total task	Binary
<b>Minimum</b>	9141	9119	9196	9164
<b>Mean</b>	15224	15208	15027	15010
<b>Maximum</b>	31204	31178	30463	30435

TABLE 6.2: Task times

## 6.1.2 Matmul

### 6.1.2.1 Application description

The Matrix Multiplication (*Matmul*) is a common operation in diverse fields such as numerical methods, earth science, industrial simulations, machine learning or bioinformatics. In fact, this operation has become so common that most of the underlying libraries (i.e. Intel Math Kernel Library - MKL [29]) provide an easy-to-use API call to an optimized implementation.

The common parallel approach to the Matrix Multiplication is to divide the matrix in smaller matrixes, called *blocks* (see Figure 6.9). When multiplying two block matrixes (i.e. *A* and *B*), the result matrix (i.e. *C*) is calculated block by block. Notice that the *C* blocks can be calculated independently and thus, in parallel.

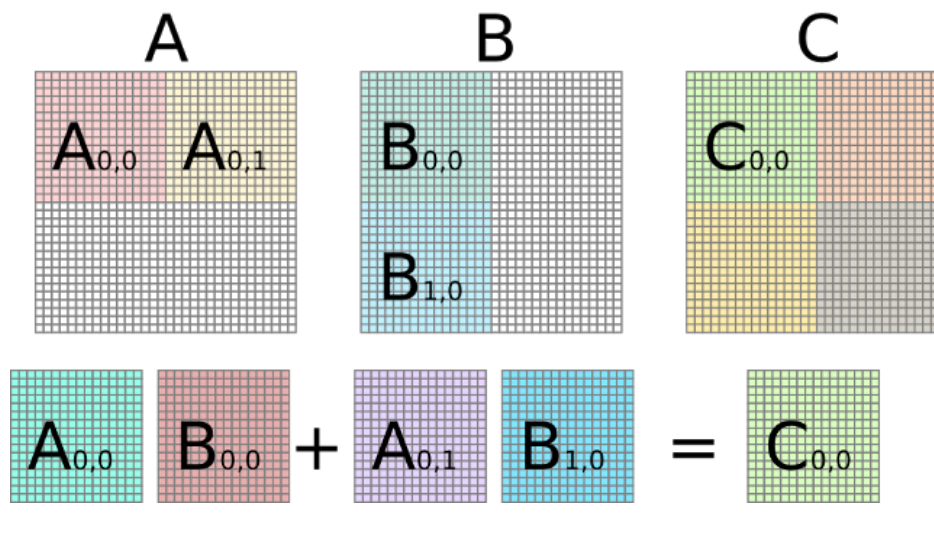


FIGURE 6.9: Multiplication of a Matrix divided in blocks  
Source : Dongrui She, GPU Assignment 5KK70

The COMPSs implementation of *Matmul* constructs two float matrixes, *A* and *B*, with *MSIZE* blocks of size *BFSIZE*. Both parameters are used passed through the application arguments. Notice that the matrixes are squared, and its real size is stated in Equation 6.2.

$$(MSIZE \cdot BFSIZE) \times (MSIZE \cdot BFSIZE) \text{ floats} \quad (6.2)$$

Consequently, the result matrix, *C*, is also squared and has the same size.

Regarding the code, it is divided in two main blocks:

- **Initialization:** Initialization of each block of matrixes *A* and *B* to random float values using files to store each block.
- **Multiplication Loop:** Multiplies and accumulates the two matrixes block by block using the IKJ-algorithm which has been proven to be the best algorithm in terms of performance [31] without using external libraries.

The initialization of each block is performed inside the *initializeBlock* method, and the multiplication and accumulation is done in the *multiplyAccumulative* method. Both methods are annotated in the *Matmul* Interface so that COMPSs spawn them as tasks. Notice

that since the multiplication and accumulation of each block of  $C$  is independent, the implementation spawns  $MSIZE \cdot MSIZE$  chains (number of total blocks of  $C$ ) of depth  $MSIZE$  (number of accumulations of a block of  $C$ , this is, the length of a row of  $A$  and the length of a column of  $B$ ). Figure 6.10 shows the task execution graph of a *Matmul* of  $MSIZE = 2$ ,  $BSIZE = 128$  where the tasks before the barrier point are the initialization of the matrixes  $A$ ,  $B$  and  $C$ , and the task-chains are the previously explained multiplications.

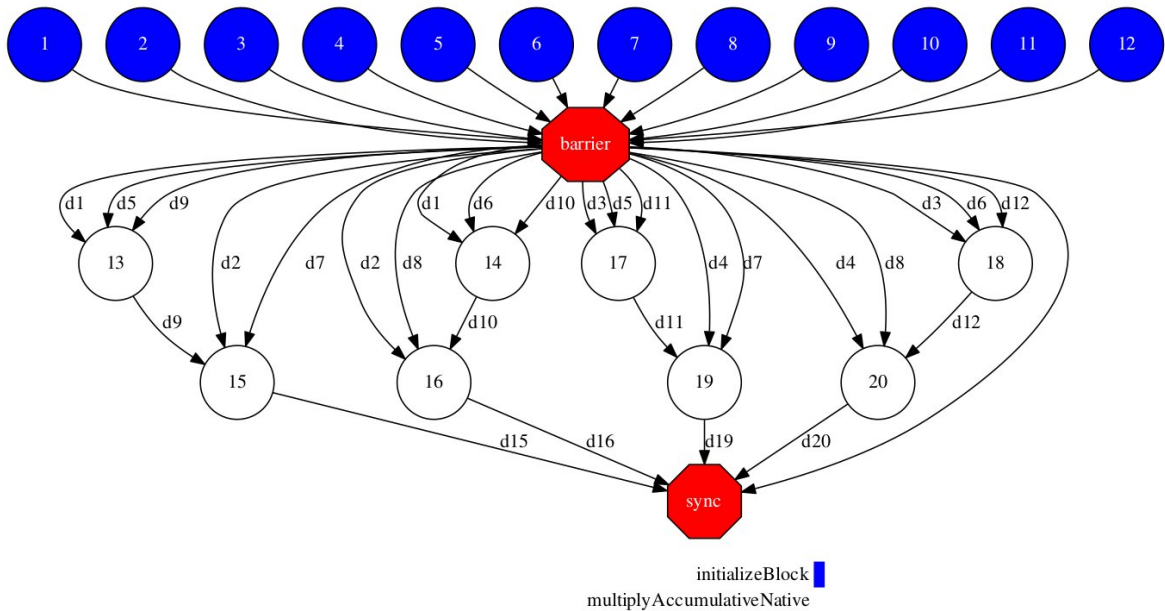


FIGURE 6.10: Task execution graph of a *Matmul* example with  $MSIZE = 2$

### 6.1.2.2 Purpose

This proof of concept must demonstrate that our enhancement of the COMPSs Runtime eases the execution of MPI tasks by obtaining a hybrid COMPSs-MPI version of the *Matmul*. This new version must be capable of orchestrating several MPI executions in more than one computational node without introducing a significant overhead to the MPI execution. Moreover, the users must have enough flexibility to execute the MPI binaries with different parameters, to synchronize (or not) the exit value, and to capture (or not) the *StdOut* and *StdErr*.

### 6.1.2.3 Hybrid COMPSs and MPI Matmul implementation

The idea behind building a Hybrid COMPSs and MPI *Matmul* implementation is to divide the matrix into bigger blocks and multiply them in parallel using MPI. To do so, we add an extra MPI layer under COMPSs that multiplies one block of matrixes  $A$  and  $B$  using  $P$  MPI processes, and accumulates the result on one block of  $C$  (see Figure 6.11). Since, for this MPI layer, each COMPSs block is treated as a single matrix, the MPI implementation is independent of the COMPSs implementation. In this sense, we have chosen the standard MPI *Matmul* implementation that divides the  $A$  matrix into line-blocks and broadcasts the entire  $B$  matrix.

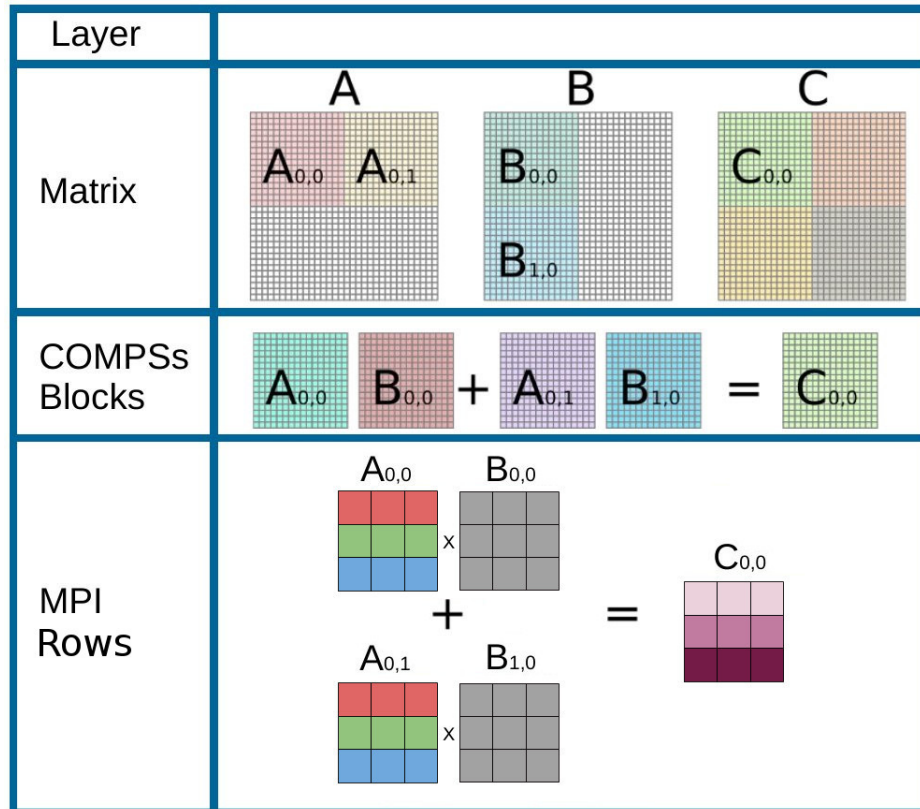


FIGURE 6.11: Hybrid COMPSs and MPI block layers

For this purpose, we have substituted in the previous implementation the *multiplyAccumulative* task by an MPI task that invokes any of the MPI implementation of the *Matmul*. Figure 6.12 shows the main multiplication loop invoking the MPI task and retrieving the exit value.

```

private static void computeMultiplication() {
2   Integer[][][] exitValues = new Integer[MSIZE][MSIZE][MSIZE];

4   // Launch tasks
   for (int i = 0; i < MSIZE; ++i) {
6       for (int k = 0; k < MSIZE; ++k) {
           for (int j = 0; j < MSIZE; ++j) {
8               exitValues[i][j][k] = MPI.multiplyAccumulative(BSIZE, AfileNames[i][k],
                   BfileNames[k][j], CfileNames[i][j]);
10            }
12        }
   }
}

```

FIGURE 6.12: Main multiplication loop of the Hybrid *Matmul*

Figure 6.13 shows the interface annotation of the MPI *multiplyAccumulative* task. Notice that the binary path, the number of computing nodes, and the number of computing units per node are retrieved from environment variables to easily parametrize the MPI binary path and the number of MPI processes during the experimentation phase. Moreover, the total number of processes available for the MPI *Matmul* is  $P = CN \cdot CUS$ . Regarding the

task parameters, we use the MPI matrix size (this is, the block size for the COMPSs layer) and the file path for the contents of the matrixes  $A$ ,  $B$  and  $C$  (which are, for the COMPSs layer, the matrix blocks).

```

1 @MPI(binary = "${MATMUL_BINARY}",
    mpiRunner = "mpirun",
3     computingNodes = "${CN}")
@Constraints(computingUnits = "${CUS}")
5 Integer multiplyAccumulative(
    @Parameter() int bsize,
7     @Parameter(type = Type.FILE, direction = Direction.IN) String aIn,
    @Parameter(type = Type.FILE, direction = Direction.IN) String bIn,
9     @Parameter(type = Type.FILE, direction = Direction.INOUT) String cOut
);

```

FIGURE 6.13: *multiplyAccumulative*'s interface annotation for the Hybrid COMPSs and MPI *Matmul*

Appendix B contains the complete code of this Hybrid COMPSs and MPI *Matmul* implementation.

#### 6.1.2.4 Evaluation

To evaluate the new Matrix Multiplication implementation, we have run all the experiments on the *Nord III* [13] cluster, hosted at the Barcelona Supercomputing Center (BSC) composed of one compute rack from the MareNostrum III [10]. This supercomputer has 1344 Intel SandyBridge EP-2670 cores at 2.6 GHz, 10.5 TB of main memory, a peak performance of 28 GigaFlops and InfiniBand interconnection.

We have compared the standalone MPI implementation and the Hybrid (COMPSs and MPI) implementation for a different number of processes and a fixed matrix size to evaluate the overhead caused by including a new COMPSs layer on top of the MPI. Regarding the number of processes, we have chosen multiples of 16 to fill the computational nodes completely. Concerning the matrix size, it has been fixed to 16384 long float elements because it is the biggest matrix that can fit inside the memory of a single node (uses up to 29 Gb of main memory).

On the one hand, we have performed a set of experiments with a single MPI execution but with an increasing number of processes. Figure 6.14 shows the strong scaling analysis of both, the MPI and the Hybrid implementations. The left y-axis represents the execution time (in seconds), the right y-axis represents the percent overhead of the Hybrid implementation with respect to the MPI implementation, and the x-axis represents the number of processes. Notice that none of the implementations scales because the execution time includes the reading, the computation and the writing times, and, in our implementations, the reading and the writing phases are performed only by the master process. Although there are more sophisticated implementations that initialize and store the data in parallel, this is beyond the scope of this section since we are only willing to evaluate the overhead of including a new layer on top of MPI. In this sense, notice that the overhead of wrapping the MPI *Matmul* implementation with COMPSs is below the 5% for all the cases.

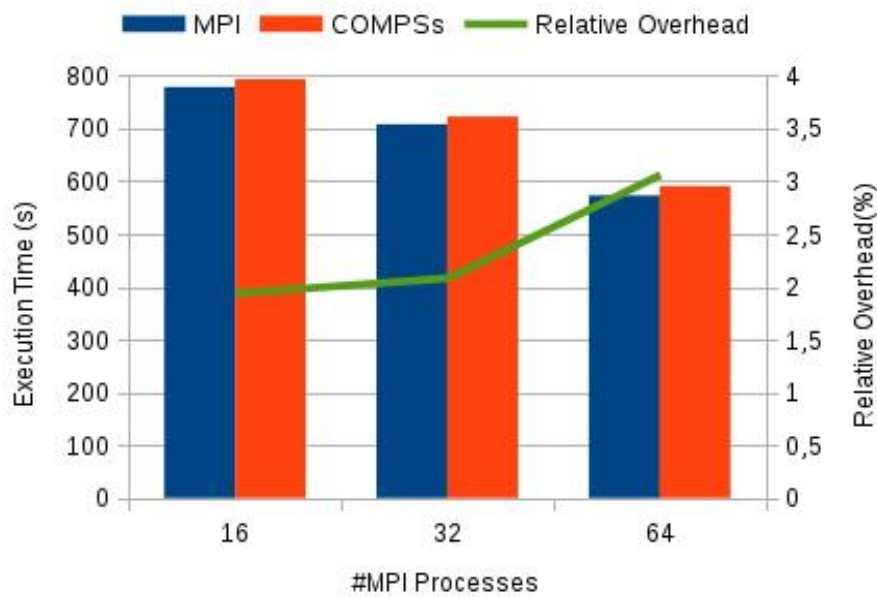


FIGURE 6.14: *Matmul* Strong Scaling analysis

On the other hand, we have performed a set of experiments with multiple MPI executions with the same amount of MPI threads. We have fixed the MPI threads to 16 to use one node completely, but as future work, it could be interesting to perform a more in-depth experimentation with an MPI thread size that does not fit properly to the node slots.

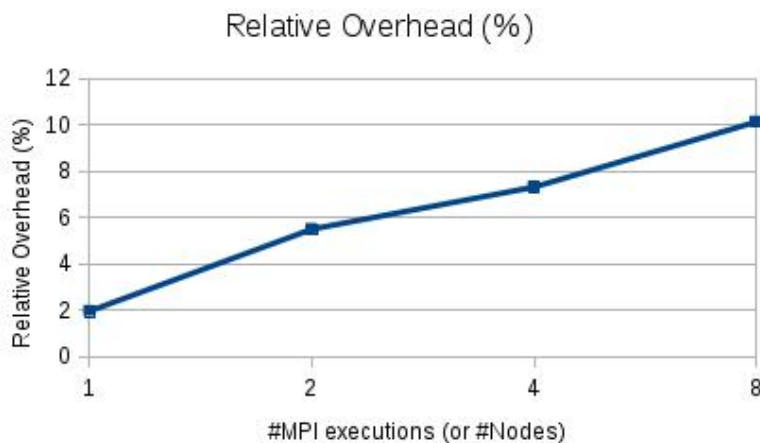


FIGURE 6.15: *Matmul* Weak Scaling analysis

Figure 6.15 shows the weak scaling analysis of the Hybrid implementation. The y-axis represents the percent overhead with respect to the standalone MPI implementation with 16 processes and matrix size 16384 long float elements, and the x-axis represents the number of simultaneous MPI executions of 16 threads each (or, what is equivalent, the number of nodes). Notice that, when using COMPSs to orchestrate several MPI executions, the overhead remains negligible (less than 8%) for up to 4 MPI executions (each of them using a complete node). However, when orchestrating 8 or more MPI executions, the overhead of

---

scheduling the tasks and managing the available resources becomes significant and requires a high-load computation in the MPI layer.

To conclude this section, the first set of experiments have demonstrated that including a new COMPSs layer on top of the MPI layer does not cause a significant overhead. Moreover, in the second set of experiments, we have shown that our enhancement of the COMPSs Runtime is capable of orchestrating up to 4 different MPI executions of 16 threads each without a significant overhead.

## 6.2 Use cases

### 6.2.1 NMMB/BSC-Dust

#### 6.2.1.1 Application description

NMMB/BSC-Dust model [12] is pluggable component of the Non-hydrostatic Multiscale Model (NMMB) designed and developed by the Earth Sciences Department of the Barcelona Supercomputing Center (BSC) in collaboration with NOAA/National Centers for Environmental Prediction (NCEP), NASA Goddard Institute for Space Studies and the International Research Institute for Climate and Society (IRI). This model provides short to medium-range weather and dust forecasts for regional and global domains (Figure 6.16 shows a simulation example of four-hour average AOD of the North African Domain on 16 May 2006).

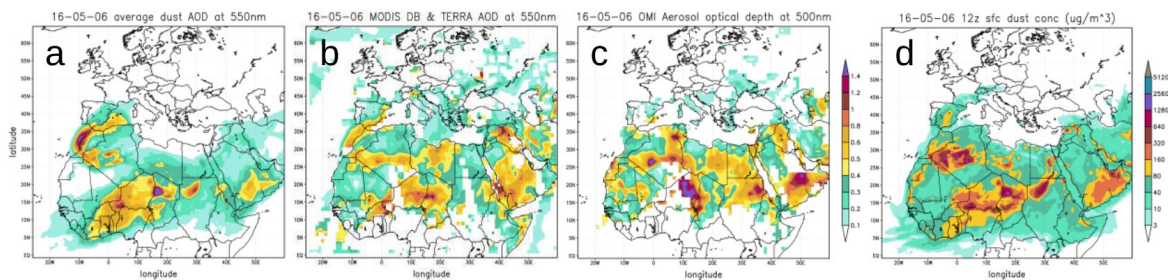


FIGURE 6.16: Example of four hour average AOD from NMMB/BSC-Dust  
Source: NMMB/BSC-Dust model [12]

NMMB/BSC-Dust model appears in the frame of quantifying the global dust emission from arid and semi-arid areas and its spatial distribution. Far from being useless, “*dust significantly affects the global and regional energy balance by absorbing and scattering shortwave and longwave radiation, dust transported by winds modifies atmospheric heating rates, temperature and stability, influences the hydrological cycle and impacts the human health*” [1]. In fact, in the past decades, several models have been developed to reproduce the dust cycle and estimate its influence on the climate system.

More in-depth, the NMMB/BSC-Dust model defines a complex analytic workflow in a *BASH* script of 778 lines that handles the execution of an initialization step (called *FIXED*) and a main loop that is executed for each timestep of the simulation period. The main loop is defined by three internal steps: the *VARIABLE* pre-process, the *NMMB/BSC-Dust model simulation* and the post-process step. Figure 6.17 depicts the general NMMB/BSC-Dust workflow. Although each of these steps spawns a large number of binary calls (the biggest part of them written in *Fortran 77* and *Fortran 90*), the only one that actually runs in parallel is the MPI simulation inside the *UMO\_MODEL* step. Furthermore, the script is already prepared to support different workflow options by redefining some variables inserted at its beginning (i.e. enabling and disabling a specific step, defining a global or a regional domain, changing the domain size, or changing the model data folder or parameters).

#### 6.2.1.2 Purpose

The purpose is to evaluate the contributions of this thesis by porting a real data science application to COMPSs so that it can benefit of the COMPSs programming model abstractions. The porting must transform the current NMMB/BSC-Dust implementation (a complex sequential workflow defined in *BASH*) into an easy, parallel, and portable code without any loss of performance.



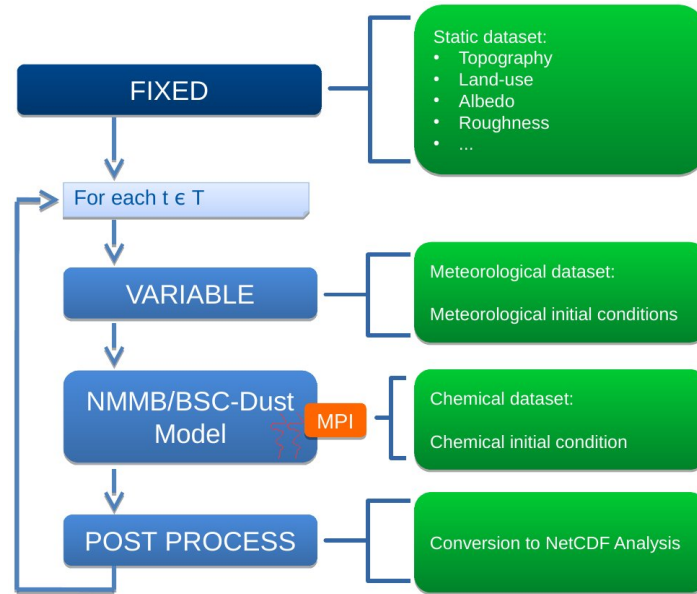


FIGURE 6.17: NMMB/BSC-Dust step workflow

### 6.2.1.3 NMMB/BSC-Dust implementation with COMPSs

From our point of view, having a complex workflow defined in such a long *BASH* script lacks reliability. On the one hand, because users need to modify the code to change the simulation parameters and, on the other hand, because maintaining this script becomes hard for unexperienced users.

Consequently, we have ported the NMMB/BSC-Dust workflow to a Java application. The new workflow has a main class to define the main step workflow (*Nmmb.java*), the application interface for COMPSs (*NmmbItf.java*), the binary and MPI tasks dummy definitions (*BINARY.java* and *MPI.java*), two classes to load and handle the simulation parameters (*NMMBConfigManager.java* and *NMMBParameters.java*), a class to handle the environment variables (*NMMBEnvironment.java*), a class to store all the constant values (*NMMBConstants.java*) and five util classes to encapsulate the file management, the logger printers and the *BASH* and *Fortran* executors. Figure 6.3 shows the code summary for the implementation and Appendix C contains the main file and the interface of the application.

Language	Files	Blank	Comment	Code
Fortran 90	23	394	2806	7581
Fortran 77	8	182	3568	6518
Java	18	558	887	2688
Bourne Shell	12	155	117	669
Maven	1	10	16	162
Bourne Again Shell	3	19	14	64
XML	1	0	0	13
<b>SUM</b>	<b>66</b>	<b>1318</b>	<b>7408</b>	<b>17695</b>

TABLE 6.3: NMMB/BSC-Dust code summary

This implementation offers three principal advantages. First, the execution parameters

are loaded from a configuration file that can be modified without recompiling the application (and without risking of involuntary modifications on the workflow behavior). Secondly, as it is built on top of COMPSs, it can rely on the power of the Runtime to abstract from the underlying infrastructure. In this sense, it can benefit from the *runcomps* and *enqueue\_comps* commands to run over different infrastructures without modifying the application code. Finally, the binaries inside each step have been parallelized so that the available resources are exploited as much as possible.

#### 6.2.1.4 Evaluation

To evaluate the new NMMB/BSC-Dust implementation we have run the experiments on the *Nord III* [13] cluster, hosted at the Barcelona Supercomputing Center (BSC) composed of one compute rack from the MareNostrum III [10]. This supercomputer has 1344 Intel SandyBridge EP-2670 cores at 2.6 GHz, 10.5 TB of main memory, a peak performance of 28 GigaFlops and InfiniBand interconnection. For the experimentation, we have run several simulations of 1 day (a single loop iteration) of a global domain with 64 cores available for the MPI simulation (4 nodes). In the case of the COMPSs implementation, we have requested 4 nodes for computation (COMPSs Workers) and an extra node for the COMPSs Master processes.

First, we have checked that the final result produced by both implementations is the same. The NMMB/BSC-Dust application produces a forecast binary file (*CASE.nc*) that can be visualized with the NCVIEW [20] visual browser for *netCDF* format files. As expected, in both cases, the file produced weights 210 Mb, and its visualization results are the same. Figures 6.18, 6.19 and 6.20 show three different views of the forecast result with both implementations (at right, the original result and, at left, with the new COMPSs implementation) to demonstrate that the visualizations are the same.

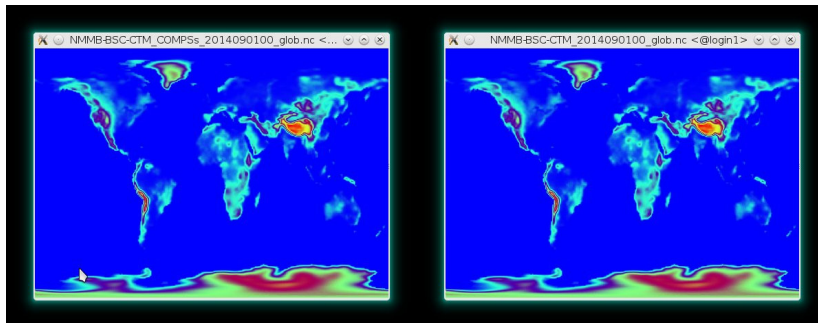


FIGURE 6.18: NCVIEW plot of FIS 3D variable for both implementations

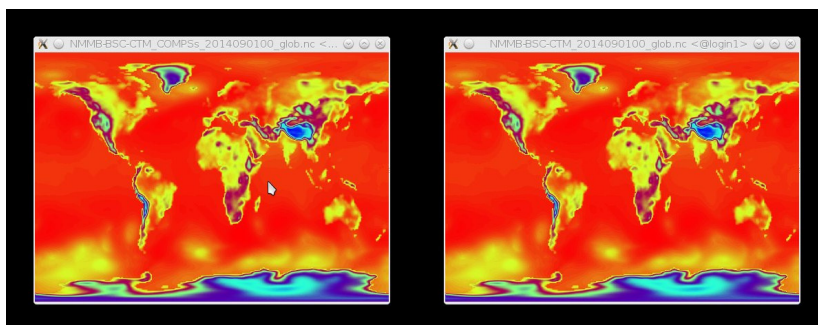


FIGURE 6.19: NCVIEW plot of PS 3D variable for both implementations

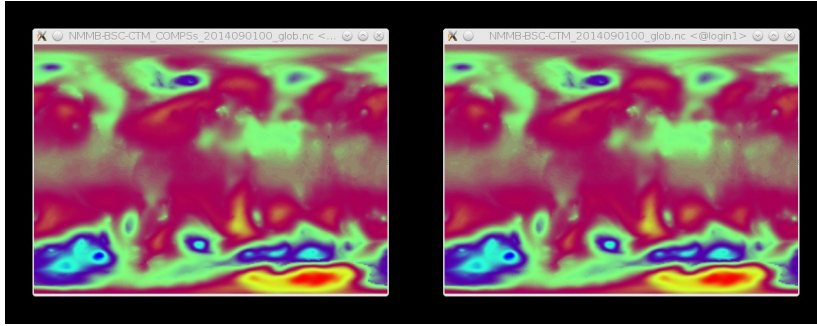


FIGURE 6.20: NCVIEW plot of SLP 3D variable for both implementations

Next, In terms of performance, we have measured the execution time of each step and the total execution time for both implementations. On the one hand, for the *old* implementation we have measured the time inserting *time* calls on the main BASH script. On the other hand, for the implementation with COMPSs, we have measured the time adding system calls on the main Java class. Table 6.4 summarizes the execution times for each implementation and the relative speed-ups. Notice that the implementation with COMPSs has an overall speed-up of  $s = 1.45$  because the fixed, the variable and the post-process steps are splitted in tasks that can run in parallel. Although we will proceed to a more in-depth study, we can advance that this speed-up is limited by the data dependencies between the tasks inside each step. Moreover, the execution time of the Model Simulation step remains the same, allowing us to ensure that the current implementation of the *MPI* annotations do not introduce a significant overhead.

Steps	Execution Times (s)		Speed-up (u)
	Previous Impl.	COMPSs Impl.	
<b>Fixed</b>	290	117	2.48
<b>Variable</b>	26	19	1.37
<b>Model simulation</b>	244	242	1.01
<b>Post process</b>	38	34	1.12
<b>Total</b>	<b>601</b>	<b>413</b>	<b>1.45</b>

TABLE 6.4: Execution times of the different NMMB/BSC-Dust implementations for the simulation of 1 day of global domain with 64 cores

Finally, we have proceeded to analyze in-depth the COMPSs implementation. Figure 6.21 shows the task graph of the NMMB/BSC-Dust when executed with COMPSs. For the sake of clarity, we have highlighted each of the different steps.

Notice that the *fortranCompile* tasks open a significant parallelism at the beginning of the FIXED and the VARIABLE phases, but the rest of the tasks define a quite complex graph with lots of dependencies. In fact, the current implementation defines 37 different tasks and executes 58 tasks for a single iteration of the main loop (31 for the fixed step, 24 for the variable step, 1 for the model simulation and 2 for the post-process). Furthermore, the tasks have at least 3 parameters, at most 47 parameters (*allprep* task) and 7 parameters in average.

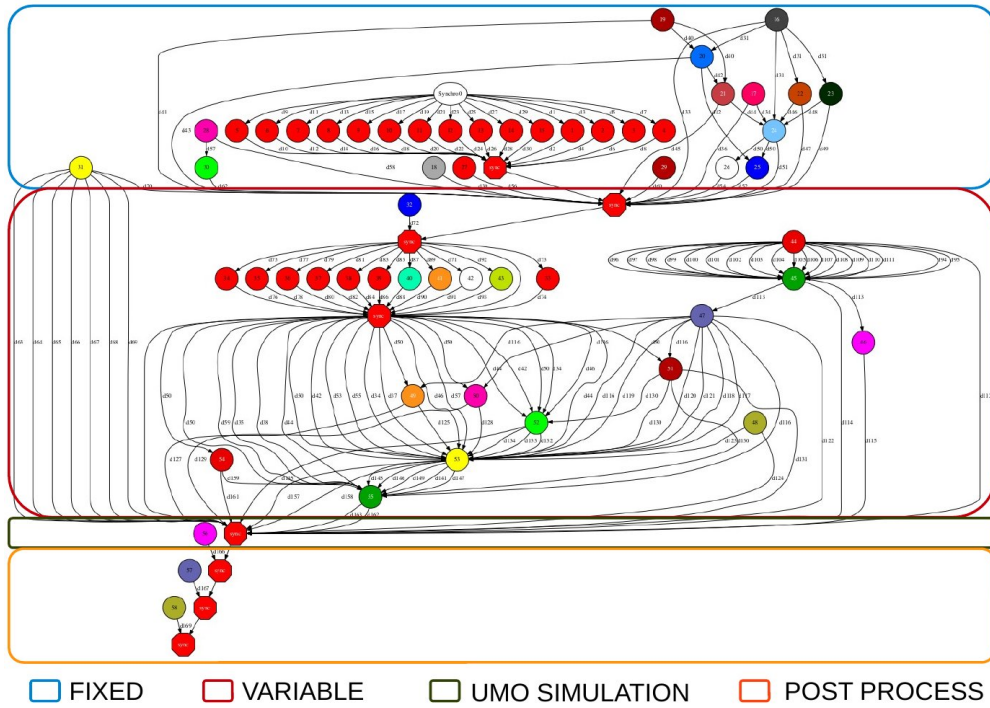


FIGURE 6.21: Tasks graph of NMMB/BSC-Dust with COMPS

Thanks to the integration of COMPSs with the Extrae tool we have also obtained post-mortem trace of the application’s execution. Figure 6.22 shows the Paraver view of the task trace of the NMMB/BSC-Dust implementation with COMPSs. For the sake of clarity, we have highlighted each of the steps at the bottom of the image.

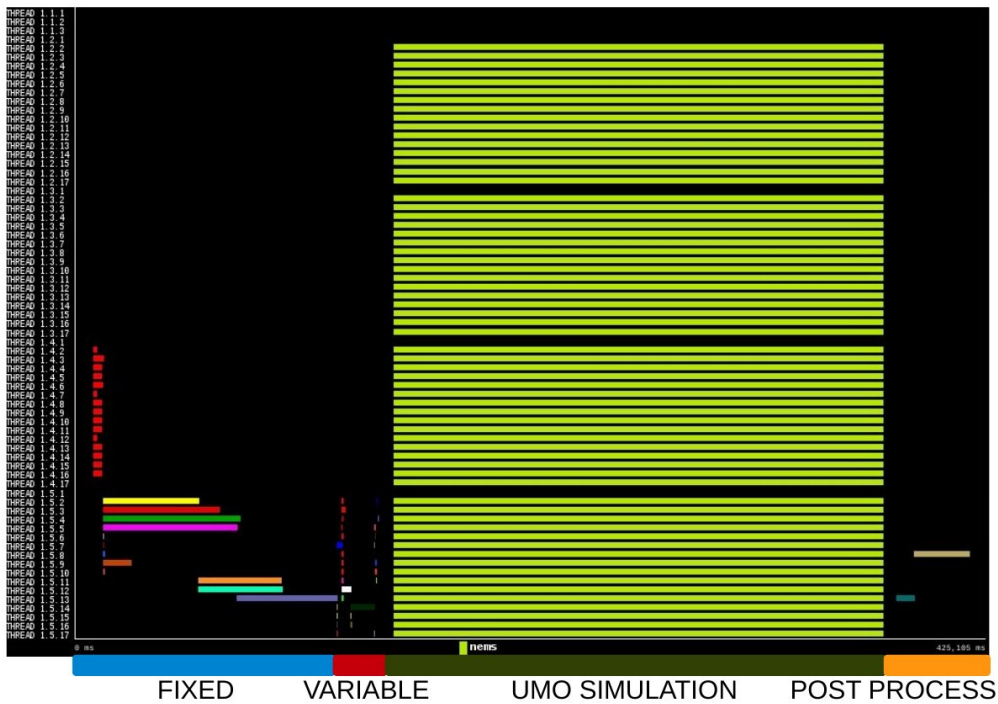


FIGURE 6.22: Paraver task view of the NMMB/BSC-Dust execution

Notice that the first node is reserved for the COMPSs Master and that the first thread of each node is reserved for the COMPSs Worker. For the rest of the threads, the application does not fill all the available resources during the *fixed*, *variable* and *post-process* steps because there are not enough tasks to run in parallel. In fact, the first compilation tasks (in red) are run in the fourth resource (threads 1.4.2 - 1.4.17) because it is the first available worker in the execution. However, the rest of the *fixed* and *variable* tasks are executed in the fifth resource (threads 1.5.2 - 1.5.17) because the first task has been arbitrarily scheduled to this resource and the rest of them is scheduled to the same resource because of data locality (notice that, as shown in the task graph, there is a synchronization point rather than a data dependency between the first compilation tasks and the rest of the tasks). Regarding the Model Simulation, the scheduler has reserved all the available resources because it requires 64 cores. On the *post-process* step, the two tasks are also scheduled to the fifth resource because of data locality.

As a final note, this implementation is a prototype of the NMMB/BSC-Dust running with COMPSs. For next implementations, since there are free resources during the *fixed*, the *variable* and the *post-process* steps, it would be worth to avoid synchronizations between them and try to overlap the execution of tasks of different steps. It will also be interesting to use dynamic resource management to release resources during these steps and acquire them again during the execution of the Model Simulation (since the requirements of this MPI step are the ones that obligate the application to request so many resources).

## 6.2.2 GUIDANCE

### 6.2.2.1 Application description

GUIDANCE [9] is a framework for large-scale genome and phenome-wide association studies on parallel computing platforms developed by the Computational Genomics Group at the Barcelona Supercomputing Center (BSC). This integrated framework provides an easy solution to perform Genome and Phenome association analysis, allowing the users to perform all the steps in a single execution or in a modular way with optional user intervention. The current GUIDANCE's implementation is already based on COMPSs to make the application integrable to multiple parallel and distributed platforms, and to ensure the efficient usage of the computing resources.

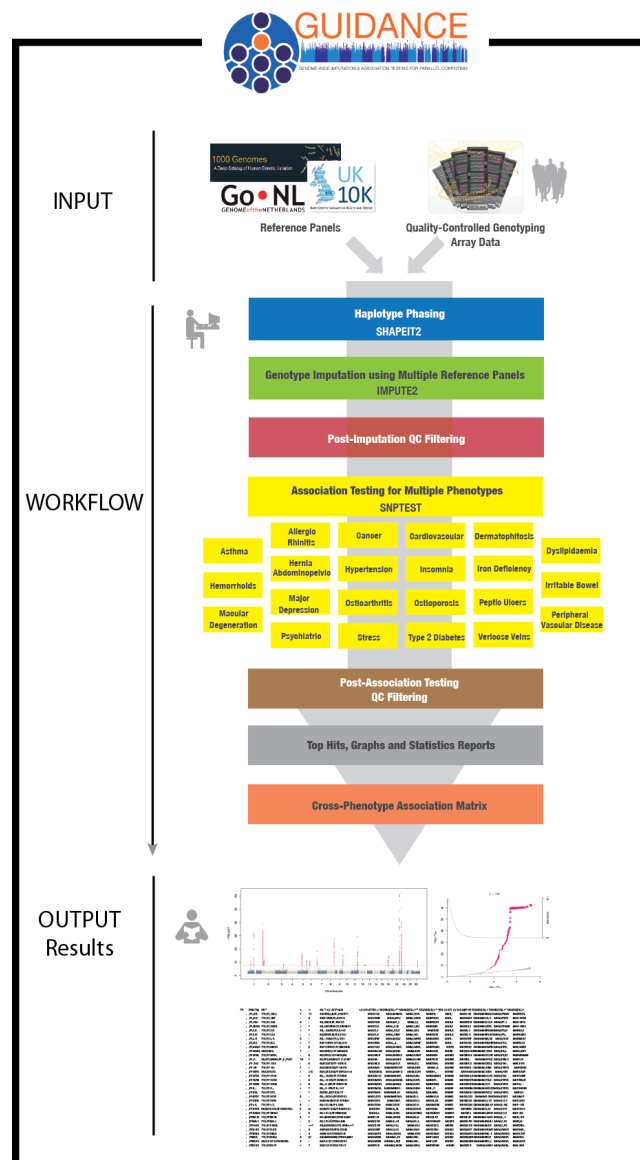


FIGURE 6.23: GUIDANCE's schematic representation of the typical complete Genome and Phenome association analysis. Source: [9]

Figure 6.23 provides a schematic representation of the steps performed by the typical complete Genome and Phenome association analysis. The application's workflow starts

retrieving Quality Controlled genetic data. Next, it goes through phasing and imputation using multiple panels and performs an association test considering multiple phenotypes (*SNPTEST*). Finally, GUIDANCE provides summary statistics and graphical representations of the results.

### 6.2.2.2 Purpose

The purpose is to evaluate the advantages of this thesis' developments with respect to the previous COMPSs model by upgrading a real data science application that already uses COMPSs. During the past year, the previous GUIDANCE implementation has become a huge workflow that can be hardly managed, mainly because of the invocation of many binary files. The upgrade must simplify the previous implementation without any loose of performance so that future engineers can keep improving the GUIDANCE framework in an easy way.

### 6.2.2.3 GUIDANCE implementation with COMPSs

GUIDANCE's main code is implemented in Java and defines a task workflow of both binary invocations and native methods. The main file (*Guidance.java*) contains the workflow description and a wrapper function for each task defined inside the *GuidanceImpl.java* file. All the tasks are duly annotated inside the *GuidanceItf.java* file.

```
public static void binaryTaskExample(String[] cmdParams) throws Exception {
2     /**** SOME PRE-PROCESS (if any) ****/
3
4     // Command construction
5     String cmd = BINARY + FLAGS + cmdParams;
6     ProcessBuilder pb = new ProcessBuilder(cmd.split(" "));
7     // Process spawn
8     Process p = pb.start();
9     // Handling the streams so that dead lock situation never occurs
10    byte[] b = new byte[1024];
11    int read;
12    BufferedInputStream bisInp = new BufferedInputStream(p.getInputStream());
13    BufferedOutputStream bosInp = new BufferedOutputStream(
14    new FileOutputStream(outputFile + ".stdout"));
15    while ((read = bisInp.read(b)) >= 0) {
16        bosInp.write(b, 0, read);
17    }
18    BufferedInputStream bisErr = new BufferedInputStream(p.getErrorStream());
19    BufferedOutputStream bosErr = new BufferedOutputStream(
20    new FileOutputStream(errorFile + ".stderr"));
21    while ((read = bisErr.read(b)) >= 0) {
22        bosErr.write(b, 0, read);
23    }
24    bisErr.close();
25    bisInp.close();
26    bosErr.close();
27    bosInp.close();
28    // Check the proper ending of the process
29    int exitValue = p.waitFor();
30    if (exitValue != 0) {
31        throw new Exception("[ERROR] Task error, exit value is: " + exitValue);
32    }
33
34    /**** SOME POST-PROCESS (if any) ****/
35 }
```

FIGURE 6.24: Example of previous GUIDANCE binary task implementation

The previous implementation defines 24 different tasks; where 9 of them spawn a binary execution. Considering only these 9 tasks that spawn a binary execution, all of them have a common structure (see Figure 6.24): they construct the binary command, spawn a *Process-Builder*, launch the binary execution, capture the command output and error, and retrieve the process exit value. In some specific cases, there is also some pre-process to parse or decompress the given arguments, and some post-process to rename or compress the output files. Notice that this structure is repeated several times among the code and gets the code dirty for the final user.

Regarding the execution, GUIDANCE's smallest test lasts 9780 seconds (less than 3 hours) and spawns 1010 tasks. The test size cannot be smaller because the test must cover all the possibilities within the workflow. As depicted in the cropped task graph shown in Figure 6.25, the GUIDANCE implementation has an initialization phase with phew parallelism, an intermediate phase with a scatter/gather structure, and a final phase to analyze the results.

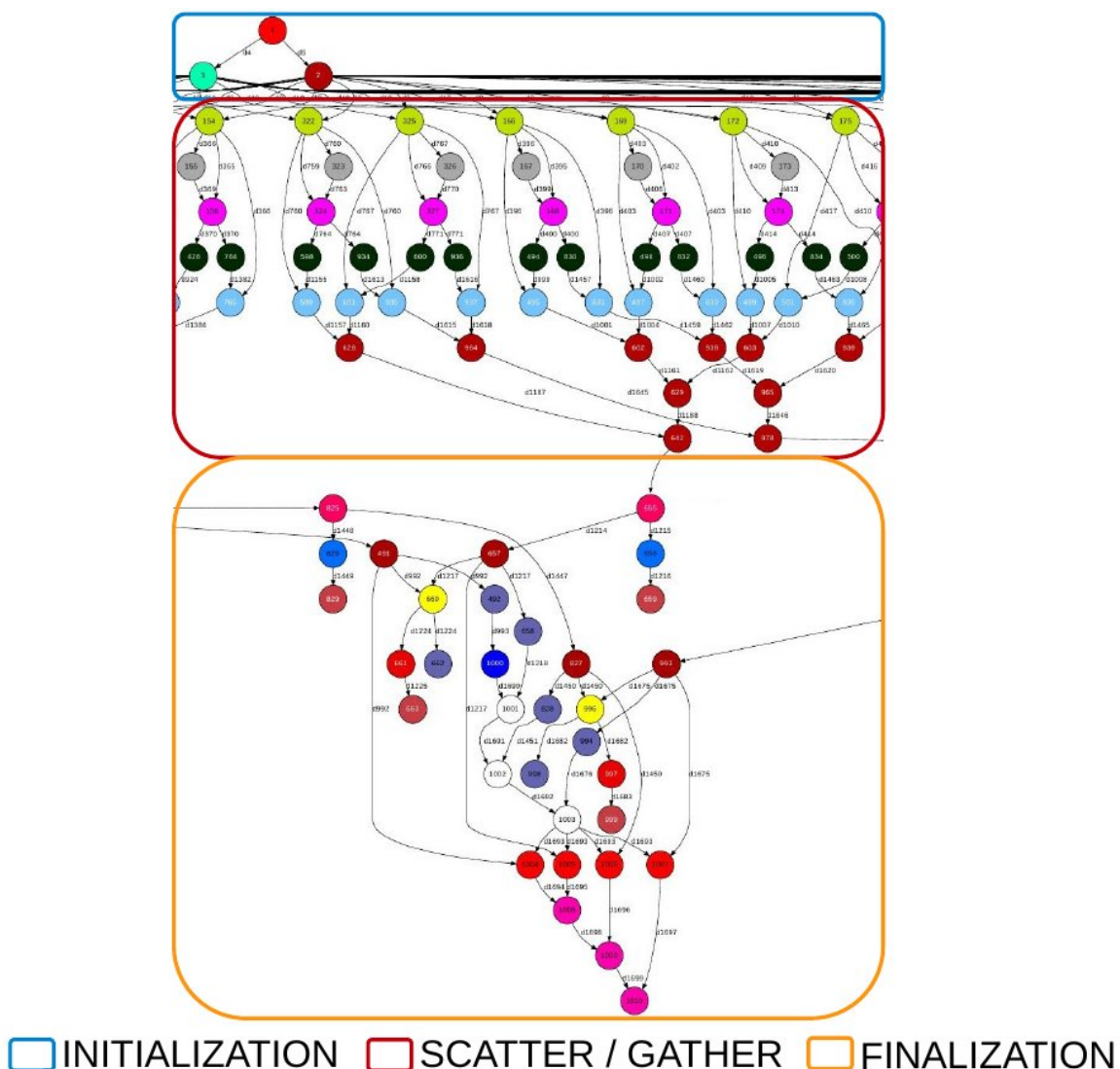


FIGURE 6.25: GUIDANCE's partial task graph



### 6.2.2.4 Evaluation

Regarding programmability, the new annotations allow the users to remove all the *ProcessBuilder* management required to launch binaries. In the GUIDANCE case, 9 tasks from the 24 available have been completely removed and substituted by a task with a binary annotation. Since some of these tasks require pre or post steps (such as compressing and decompressing files), we have added 4 additional pre/post tasks. Appendix D contains the main file and the interface of the GUIDANCE application.

Taking all into consideration, as shown in Figure 6.5 the tasks' code (*GuidanceImpl*) has been reduced by three and the application's main code (*Guidance*) has been reduced by a 20 %. On the other hand, the application's annotation interface has grown since the binary tasks have more parameters (the previous implementation had an average of 6.5 parameters per task, and the new implementation has an average of 9.0 parameters per task). The overall GUIDANCE's code has been reduced a 23.28 %, mostly due to the removal of the *ProcessBuilder* management.

File	Implementation	Num. Files	Blank	Comment	Code
Guidance	NEW	1	327	749	1296
GuidanceImpl	NEW	1	193	445	1074
BINARY	NEW	1	20	210	55
GuidanceItf	NEW	1	32	0	374
TOTAL	NEW	43	1862	4108	7125
Guidance	OLD	1	457	531	1637
GuidanceImpl	OLD	1	1066	1053	3488
GuidanceItf	OLD	1	54	41	352
TOTAL	OLD	14	2619	2891	9288

TABLE 6.5: GUIDANCE code summary

From our point of view, this new implementation is way easier to maintain than the previous one since the programmers do not need to handle with the binary invocations directly (but rather rely on the COMPSs *@Binary* annotation). Moreover, we have included an abstraction layer on the application's main code so that the programmers can add new steps to the workflow in an easy templated way.

Finally, due to time constraints and the shutdown of the MareNostrum III Supercomputer this application has only been evaluated in terms of programmability. As future work, we plan to validate the implementation with big runs in the new MareNostrum IV supercomputer. However, since previous cases did not show up any loss of performance, we do not expect any improvement nor deterioration of the execution times.



## Chapter 7

# Conclusions and Future work

This Master Thesis provides a first adaptation of the COMPSs programming model to the needs of the Big-Data Ecosystems. Concerning the COMPSs programming model, this thesis provides a new set of task annotations to easily integrate Java workflows with the execution of binaries, and *MPI* and *OmpSs* applications. To provide full support to this integration, we have also implemented two new parameter annotations to add string prefixes and Linux stream redirections. Taking advantage of the new annotations, we have profited to redesign the versioning annotations, to extend the task annotations to support environment variables, and added new annotations for scheduler hints.

Regarding the COMPSs internals, this project has extended the schedulers to support multi-node actions in a transparent and non-blocking fashion. Furthermore, we have redesigned the COMPSs Worker Executors to support the execution of non-native tasks and to reduce the overhead of executing bindings' tasks (Python, C, and C++ tasks).

The previous implementations have been validated against two proof of concept applications and two real uses cases. On the one side, the BLAST and the Matrix Multiplication applications have demonstrated that the new annotations reduce the code complexity while maintaining the performance. On the other side, the NMMB/BSC-Dust and the GUIDANCE applications have shown that the new COMPSs features adapt to the Big-Data Ecosystem requirements; providing a simpler, portable and efficient implementation.

During the development of this project, the COMPSs Bindings (for Python, C, and C++ languages) have increasingly gained importance to provide a comfortable integration with the data science workflows. Hence, as future work, we plan to integrate the new COMPSs annotations with the COMPSs bindings. Moreover, since the community has reacted positively to the overhead's reduction that the new COMPSs Worker Executors provide when executing binding tasks, we plan to design persistent Python, C and C++ workers that can communicate through pipes with the current Java persistent worker and reduce even more this overhead.

Regarding the COMPSs Runtime itself, as future work, we also plan to implement pluggable scheduler politics so that users can select the suitable scheduler for their applications. In this sense, the new schedulers could profit from the new scheduler hints annotation to have more information about the tasks and take more complex decisions.

Finally, concerning the real uses cases, on the one hand, we want to perform executions of the NMMB/BSC-Dust application in other supercomputers to take advantage of the portability that the new COMPSs implementation provides. On the other hand, we want to do a performance analysis of the new GUIDANCE implementation to validate and compare it against the previous implementation.



# Bibliography

- [1] Hausteijn et al. Perez et al. "Atmospheric dust modeling from meso to global scales with the online NMMB/BSC-Dust model – Part 2: Experimental campaigns in Northern Africa". In: *Atmospheric Chemistry and Physics Journal (ACP)* 12 (Mar. 2012), pp. 2933–2958. URL: <http://www.atmos-chem-phys.net/12/2933/2012/>.
- [2] Jack Dongarra et al. "The international Exascale Software Project roadmap". In: *International Journal of High Performance Computing Applications* 25.1 (Feb. 2011), pp. 3–60. URL: <http://hpc.sagepub.com/content/25/1/3>.
- [3] Krste Asanovic et al. "A view of the Parallel Computing Landscape". In: *Communications of the ACM* 52.10 (Oct. 2009), pp. 56–67. URL: <http://dl.acm.org/citation.cfm?id=1562783>.
- [4] R. M. Badia et al. "COMP superscalar, an interoperable programming framework". In: *SoftwareX* 3 (Dec. 2015), pp. 32–36. URL: <http://dx.doi.org/10.1016/j.softx.2015.10.004>.
- [5] Derik et al. Barseghian. "Workflows and extensions to the Kepler scientific workflow system to support environmental sensor data access and analysis". In: *Ecological Informatics* 5 (2010), pp. 42–50. URL: <http://dx.doi.org/10.1016/j.ecoinf.2009.08.008>.
- [6] National Center for Biotechnology Information (NCBI). *Blast*. URL: <https://blast.ncbi.nlm.nih.gov/Blast.cgi>.
- [7] Barcelona Supercomputing Center (BSC). *Barcelona Supercomputing Center (BSC)*. URL: <http://www.bsc.es>.
- [8] Barcelona Supercomputing Center (BSC). *Extræ Tool*. URL: <https://tools.bsc.es/extrae>.
- [9] Barcelona Supercomputing Center (BSC). *Guidance*. URL: <http://cg.bsc.es/guidance/>.
- [10] Barcelona Supercomputing Center (BSC). *MareNostrum 3 User Guide*. URL: <https://www.bsc.es/support/MareNostrum3-ug.pdf>.
- [11] Barcelona Supercomputing Center (BSC). *Minotauro*. URL: <https://www.bsc.es/innovation-and-services/supercomputers-and-facilities/minotauro>.
- [12] Barcelona Supercomputing Center (BSC). *NMMB BSC-Dust*. URL: [http://www.bsc.es/ESS/nmmb\\_bsc-dust](http://www.bsc.es/ESS/nmmb_bsc-dust).
- [13] Barcelona Supercomputing Center (BSC). *Nord III*. URL: <https://www.bsc.es/user-support/nord3.php>.
- [14] Barcelona Supercomputing Center (BSC). *OmpSs*. URL: <https://pm.bsc.es/ompss>.
- [15] Barcelona Supercomputing Center (BSC). *Paraver Tool*. URL: <https://tools.bsc.es/paraver>.
- [16] Universitat Politècnica de Catalunya (UPC). *Universitat Politècnica de Catalunya (UPC)*. URL: <http://www.upc.es>.

- [17] Scala Center. *Scala Programming Language*. URL: <https://www.scala-lang.org/>.
- [18] Shigeru Chiba. "Load-time Structural Reflection in Java". In: *ECOOP 2000 - Object-Oriented Programming 1850* (May 2000), pp. 313–336. URL: [http://dx.doi.org/10.1007/3-540-45102-1\\_16](http://dx.doi.org/10.1007/3-540-45102-1_16).
- [19] COMP Superscalar (COMPSs). *COMP Superscalar (COMPSs)*. URL: <http://compss.bsc.es>.
- [20] Scripps Institution of Oceanography David W. Pierce. *NCView*. URL: [http://meteora.ucsd.edu/~pierce/ncview\\_home\\_page.html](http://meteora.ucsd.edu/~pierce/ncview_home_page.html).
- [21] Jeffrey Dean and Sanjay Ghemawat. "MapReduce: simplified data processing on large clusters". In: *Communications of the ACM* 51.1 (Jan. 2008), pp. 107–113. URL: <http://dl.acm.org/citation.cfm?id=1327492>.
- [22] Ewa Deelman. "Big Data Analytics and High Performance Computing Convergence Through Workflows and Virtualization". In: *Big Data and Extreme-Scale Computing* (2016). URL: <http://www.exascale.org/bdec/sites/www.exascale.org/bdec/files/whitepapers/deelman-bdec2016.pdf>.
- [23] Jesús Labarta et al. Enric Tejedor Rosa M. Badia. "PyCOMPSs: Parallel computational workflows in Python". In: *The International Journal of High Performance Computing Applications (IJHPCA)* 31 (2017), pp. 66–82. URL: <http://dx.doi.org/10.1177/1094342015594678>.
- [24] R Foundation. *R Programming Language*. URL: <https://www.r-project.org/>.
- [25] GNU. *Bash*. URL: <https://www.gnu.org/software/bash/>.
- [26] GNU. *C Socket implementation*. URL: [https://www.gnu.org/software/libc/manual/html\\_node/Sockets.html](https://www.gnu.org/software/libc/manual/html_node/Sockets.html).
- [27] GNU. *GNU Plot*. URL: <http://www.gnuplot.info/>.
- [28] Galaxy Community HUB. *Galaxy Project*. URL: <https://galaxyproject.org/>.
- [29] Intel. *Intel Math Kernel Library (Intel MKL)*. URL: <https://software.intel.com/en-us/intel-mkl>.
- [30] Intel. *Intel MPI implementation*. URL: <https://software.intel.com/en-us/intel-mpi-library>.
- [31] F. G. Gustavson J. J. Dongarra and A. Karp. "Implementing linear algebra algorithms for dense matrices on a vector pipeline machine". In: *SIAM Review* 26.1 (Jan. 1984), pp. 91–112. URL: <http://dx.doi.org/10.1137/1026003>.
- [32] Anubhav et al. Jain. "FireWorks: a dynamic workflow system designed for high-throughput applications". In: *Concurrency and Computation: Practice and Experience* 27.17 (2015), pp. 5037–5059. URL: <http://dx.doi.org/10.1002/cpe.3505>.
- [33] Rosa M. Badia Javier Conejero Sandra Corella and Jesus Labarta. "Task-based programming in COMPSs to converge from HPC to big data". In: *The International Journal of High Performance Computing Applications* (Apr. 2017), p. 1094342017701278. URL: <https://doi.org/10.1177/1094342017701278>.
- [34] R. M. Badia et al. Lordan F. "ServiceSs: an interoperable programming framework for the Cloud". In: *Journal of Grid Computing* 12.1 (Mar. 2014), pp. 67–91. URL: <https://digital.csic.es/handle/10261/132141>.
- [35] Andy Konwinski Matei Zaharia Patrick Wendell and Holden Karau. *Learning Spark*. 2015.
- [36] Eclipse Members. *Eclipse IDE*. URL: <https://eclipse.org/>.

- [37] MongoDB. *MongoDB*. URL: <https://www.mongodb.com>.
- [38] MPI: *A Message-Passing Interface Standard*. June 2015. URL: <http://mpi-forum.org/docs/>.
- [39] Oracle. *Java Programming Language*. URL: <https://www.oracle.com/es/java/index.html>.
- [40] Oracle. *Java Socket implementation*. URL: <https://docs.oracle.com/javase/7/docs/api/java/net/Socket.html>.
- [41] Oracle. *JavaNIO*. URL: <http://www.oracle.com/technetwork/articles/javase/nio-139333.html>.
- [42] Oracle. *Socket Definition*. URL: <https://docs.oracle.com/javase/tutorial/networking/sockets/definition.html>.
- [43] Apache Organization. *Apache Hadoop*. URL: <http://hadoop.apache.org/>.
- [44] Apache Organization. *Taverna*. URL: <http://www.taverna.org.uk/>.
- [45] Apache Maven Organization. *Maven*. URL: <https://maven.apache.org/>.
- [46] JSON Organization. *JSON*. URL: <http://www.json.org/>.
- [47] Yaml Organization. *YAML*. URL: <http://yaml.org/>.
- [48] IBIS Project. *JavaGAT*. URL: <http://www.cs.vu.nl/ibis/javagat.html>.
- [49] Python. *Python Socket implementation*. URL: <https://docs.python.org/3/howto/sockets.html>.
- [50] Daniel A. Reed and Jack Dongarra. "Exascale Computing and Big Data". In: *Communications of the ACM* 58.7 (July 2015), pp. 56–68. URL: <http://cacm.acm.org/magazines/2015/7/188732-exascale-computing-and-big-data>.
- [51] Manchester Joint team from the universities of Southampton and Oxford in the UK. *myExperiment*. URL: <https://www.myexperiment.org/>.
- [52] Apache Spark. *Apache Spark*. URL: <http://spark.apache.org/>.
- [53] Swift. *Swift Programming Language*. URL: <http://swift-lang.org/>.
- [54] Open MPI Development Team. *Open MPI implementation*. URL: <https://www.open-mpi.org/>.





# Appendices



## Appendix A

# Blast: complete code

### A.1 Blast.java

```

1 package blast;
3 import java.io.BufferedReader;
import java.io.File;
5 import java.io.FileReader;
import java.io.IOException;
7
import java.util.UUID;
9 import java.util.ArrayList;
import java.util.LinkedList;
11 import java.util.List;
13 import binary.BINARY;
15 import blast.BlastImpl;
import blast.exceptions.BlastException;
17 import blast.utils.FilesManagement;
19 /**
 * BLAST implementation
21 *
 */
23 public class Blast {
25     private static final String ENV_BLAST_BINARY = "BLAST_BINARY";
27     private static boolean debug;
private static String databasePath;
29 private static String inputFileName;
private static int numFragments;
31 private static String tmpDir;
private static String outputFileName;
33 private static String commandArgs;
35     private static List<String> partialOutputs = null;
private static List<String> partialInputs = null;
37
/**
39 * MAIN CODE
 *
41 * @param args
 * @throws BlastException
43 */
public static void main(String[] args) throws BlastException {
45     // Parse application parameters
parseArgs(args);
47
// -----
49 // Start execution
try {
51     // Split sequence input file
splitSequenceFile();
53     // Submit tasks

```

```

        alignSequences();
55     // Assembly process
        String lastMerge=assembleSequences();
57     // Move result to expected output file
        FilesManagement.copyResult(lastMerge,Blast.outputFileName, Blast.debug);
59 } catch (BlastException be) {
        throw be;
61 } finally {
        // Clean up partial results
63     cleanUp();
        }
65 }

67 /**
68  * Parses the input parameters and initializes the class attributes
69  *
70  * @param args
71  */
private static void parseArgs(String[] args) {
73     Blast.debug = Boolean.parseBoolean(args[0]);
74     Blast.databasePath = args[1];
75     Blast.inputFileName = args[2];
76     Blast.numFragments = Integer.parseInt(args[3]);
77     Blast.tmpDir = args[4];
78     Blast.outputFileName = args[5];
79
80     Blast.commandArgs = "";
81     for (int i = 6; i < args.length; i++) {
82         Blast.commandArgs += args[i] + " ";
83     }
84 }

85 /**
86  * Splits the input sequence file into separated files
87  *
88  * @throws BlastException
89  */
private static void splitSequenceFile() throws BlastException {
91     System.out.println("Split sequence file");
92
93     // Read number of different sequences
94     int nsequences = 0;
95     try (BufferedReader bf = new BufferedReader(
96         new FileReader(Blast.inputFileName))) {
97         String line = null;
98         while ((line = bf.readLine()) != null) {
99             if (line.contains(">")) {
100                 nsequences++;
101             }
102         }
103     } catch (IOException ioe) {
104         String msg = "ERROR: Cannot read input file " + Blast.inputFileName;
105         System.err.print(msg);
106         throw new BlastException(msg, ioe);
107     }
108
109     System.out.println("- The total number of sequences is: " + nsequences);
110
111     // Calculate seqs per fragment and needed files
112     int seqsPerFragment = (int) Math
113         .round(((double) nsequences / (double) Blast.numFragments));
114     Blast.partialInputs = new ArrayList<>(Blast.numFragments);
115     Blast.partialOutputs = new ArrayList<>(Blast.numFragments);
116
117     if (Blast.debug) {
118         System.out.println(
119             "- The total number of sequences of a fragment is: "
120             + seqsPerFragment);
121         System.out
122             .println("\n- Splitting sequences among fragment files...");
123     }
124 }
125

```

```

127 // Split into files
128 for (int frag = 0; frag < Blast.numFragments; ++frag) {
129     // Creating fragment
130     UUID index = UUID.randomUUID();
131     String partitionFile = Blast.tmpDir + "seqFile" + index + ".sqf";
132     String partitionOutput = Blast.tmpDir + "resFile" + index
133         + ".result.txt";
134
135     // Initialize partial fragment information
136     BlastImpl.splitPartitions(Blast.inputFileName, partitionFile,
137         Blast.numFragments, frag);
138
139     // Store fileNames
140     Blast.partialInputs.add(partitionFile);
141     Blast.partialOutputs.add(partitionOutput);
142 }
143
144 if (Blast.debug) {
145     System.out.println("Input Files are: ");
146     for (int i = 0; i < Blast.partialInputs.size(); ++i) {
147         System.out.println("  - " + Blast.partialInputs.get(i));
148     }
149     System.out.println("Output Files are: ");
150     for (int i = 0; i < Blast.partialOutputs.size(); ++i) {
151         System.out.println("  - " + Blast.partialOutputs.get(i));
152     }
153 }
154
155 /**
156  * Aligns each fragment
157  *
158  */
159 private static void alignSequences() throws BlastException {
160     System.out.println("");
161     System.out.println("Aligning Sequences:");
162
163     final String pFlag = "-p";
164     final String pMode = "blastx";
165     final String dFlag = "-d";
166     final String iFlag = "-i";
167     final String oFlag = "-o";
168     int numAligns = Blast.partialInputs.size();
169     Integer[] exitValues = new Integer[numAligns];
170     for (int i = 0; i < numAligns; i++) {
171         if (Blast.commandArgs != null && !Blast.commandArgs.isEmpty()) {
172             exitValues[i] = BINARY.align(pFlag, pMode, dFlag,
173                 Blast.databasePath, iFlag, Blast.partialInputs.get(i),
174                 oFlag, Blast.partialOutputs.get(i), Blast.commandArgs);
175         } else {
176             exitValues[i] = BINARY.align(pFlag, pMode, dFlag,
177                 Blast.databasePath, iFlag, Blast.partialInputs.get(i),
178                 oFlag, Blast.partialOutputs.get(i));
179         }
180     }
181
182     if (Blast.debug) {
183         System.out.println("");
184         System.out.println(" - Number of fragments to assemble -> "
185             + Blast.partialOutputs.size());
186     }
187 }
188
189 /**
190  * Creates reduce tasks
191  *
192  * @return fileName of last reduce
193  */
194 private static String assembleSequences() {
195     // MERGE-REDUCE
196     LinkedList<Integer> q = new LinkedList<>();
197     for (int i = 0; i < Blast.partialOutputs.size(); i++) {

```

```

    q.add(i);
199 }

201 int x = 0;
    while (!q.isEmpty()) {
203     x = q.poll();
        if (!q.isEmpty()) {
205         int y = q.poll();

207         if (debug) {
            System.out.println(
209             " - Merging files -> " + Blast.partialOutputs.get(x)
                + " and " + Blast.partialOutputs.get(y));
211         }
            BlastImpl.assemblyPartitions(Blast.partialOutputs.get(x),
213             Blast.partialOutputs.get(y));
                q.add(x);
215         }
        }
217
    return Blast.partialOutputs.get(0);
219 }

221 /**
    * Cleans up the intermediate files
223     *
    */
225 private static void cleanUp() {
    // Cleaning intermediate sequence input files
227     for (int i = 0; i < Blast.partialInputs.size(); i++) {
        File fSeq = new File(Blast.partialInputs.get(i));
229         fSeq.delete();
        }
231
    for (int i = 0; i < Blast.partialOutputs.size(); i++) {
233         File fres = new File(Blast.partialOutputs.get(i));
            fres.delete();
235     }
    }
237 }

```

## A.2 BlastItf.java

```

package blast;
2
import integratedtoolkit.types.annotations.Parameter;
4 import integratedtoolkit.types.annotations.parameter.Direction;
import integratedtoolkit.types.annotations.parameter.Type;
6 import integratedtoolkit.types.annotations.task.Binary;
import integratedtoolkit.types.annotations.task.Method;
8
public interface BlastItf {
10
    @Method(declaringClass = "blast.BlastImpl")
12     void splitPartitions(
        @Parameter(type = Type.FILE, direction = Direction.IN) String inputFileName,
14         @Parameter(type = Type.FILE, direction = Direction.OUT) String partitionFile,
            @Parameter() int nFrag, @Parameter() int myFrag
16     );

18     @Binary(binary = "${BLAST_BINARY}")
        Integer align(
20         @Parameter(type = Type.STRING, direction = Direction.IN) String pFlag,
            @Parameter(type = Type.STRING, direction = Direction.IN) String pMode,
22             @Parameter(type = Type.STRING, direction = Direction.IN) String dFlag,
                @Parameter(type = Type.STRING, direction = Direction.IN) String database,
24             @Parameter(type = Type.STRING, direction = Direction.IN) String iFlag,

```

```

    @Parameter(type = Type.FILE, direction = Direction.IN) String partitionFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String partitionOutput,
    @Parameter(type = Type.STRING, direction = Direction.IN) String extraCMDArgs
);
26
30
@Binary(binary = "${BLAST_BINARY}")
Integer align(
    @Parameter(type = Type.STRING, direction = Direction.IN) String pFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String pMode,
    @Parameter(type = Type.STRING, direction = Direction.IN) String dFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String database,
    @Parameter(type = Type.STRING, direction = Direction.IN) String iFlag,
    @Parameter(type = Type.FILE, direction = Direction.IN) String partitionFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String partitionOutput
);
34
36
38
40
42
@Method(declaringClass = "blast.BlastImpl")
44 void assemblyPartitions(
    @Parameter(type = Type.FILE, direction = Direction.INOUT) String partialFileA,
    @Parameter(type = Type.FILE, direction = Direction.IN) String partialFileB
);
46
48
}

```

## A.3 BlastImpl.java

```

1 package blast;
3 import java.io.BufferedReader;
import java.io.BufferedWriter;
5 import java.io.FileReader;
import java.io.FileWriter;
7 import java.io.IOException;
9 import blast.exceptions.BlastException;
11 public class BlastImpl {
13 /**
    * Splits the input file into the fragment myFrag
    *
    * @param inputFileName
    * @param partitionFile
    * @param nFragments
    * @param myFrag
    * @throws BlastException
    */
    public static void splitPartitions(String inputFileName,
    String partitionFile, int nFragments, int myFrag)
    throws BlastException {
    int frag = 0;
    try (BufferedReader bf = new BufferedReader(
    new FileReader(inputFileName));
    BufferedWriter bw = new BufferedWriter(
    new FileWriter(partitionFile, true))) {
    String line = null;
    while ((line = bf.readLine()) != null) {
    if (line.contains(">")) {
    frag++;
    }
    if (frag % nFragments == myFrag) {
    bw.write(line);
    bw.newLine();
    }
    }
    }
}

```

```

41     } catch (IOException ioe) {
42         String msg = "ERROR: Cannot read input file " + inputFileA;
43         System.err.print(msg);
44         throw new BlastException(msg, ioe);
45     }
46 }
47
48 /**
49  * Assembles the two partial files
50  *
51  * @param partialFileA
52  * @param partialFileB
53  */
54 public static void assemblyPartitions(String partialFileA,
55     String partialFileB) {
56     System.out.println("Assembling partial outputs -> " + partialFileA
57         + " to " + partialFileB);
58
59     String line = null;
60     boolean append = true;
61     try (BufferedWriter bw = new BufferedWriter(
62         new FileWriter(partialFileA, append));
63         BufferedReader bfB = new BufferedReader(
64             new FileReader(partialFileB))) {
65
66         while ((line = bfB.readLine()) != null) {
67             bw.write(line);
68             bw.newLine();
69         }
70     } catch (IOException ioe) {
71         System.err.println("ERROR: Exception assembling partitions");
72         ioe.printStackTrace();
73     }
74 }
75 }

```

## A.4 BINARY.java

```

package binary;
2
3 /**
4  * Dummy class for BINARY tasks
5  *
6  */
7 public class BINARY {
8
9     public static Integer align(String pFlag, String pMode, String dFlag,
10         String database, String iFlag, String partitionFile, String oFlag,
11         String partitionOutput, String extraCMDArgs) {
12
13         return -1;
14     }
15
16     public static Integer align(String pFlag, String pMode, String dFlag,
17         String database, String iFlag, String partitionFile, String oFlag,
18         String partitionOutput) {
19
20         return -1;
21     }
22 }

```



## Appendix B

# Matmul: complete code

### B.1 Matmul.java

```

1 package matmul.files;
3 import java.io.IOException;
5 import integratedtoolkit.api.COMPSs;
import mpi.MPI;
7
/**
9  * MATMUL Implementation
10 *
11 */
public class Matmul {
13
14     private static int TYPE;
15     private static int MSIZE;
16     private static int BSIZE;
17
18     private static String[][] AfileNames;
19     private static String[][] BfileNames;
20     private static String[][] CfileNames;
21
22     /**
23      * MAIN CODE
24      *
25      * @param args
26      * @throws Exception
27      */
28     public static void main(String[] args) throws Exception {
29         // Check and get parameters
30         if (args.length != 3) {
31             System.out.println(
32                 "Usage: matmul.files.Matmul <type> <MSize> <BSize>");
33             throw new Exception("[ERROR] Incorrect number of parameters");
34         }
35         TYPE = Integer.parseInt(args[0]);
36         MSIZE = Integer.parseInt(args[1]);
37         BSIZE = Integer.parseInt(args[2]);
38
39         // Initialize matrices
40         System.out.println("[LOG] TYPE parameter value = " + TYPE);
41         System.out.println("[LOG] MSIZE parameter value = " + MSIZE);
42         System.out.println("[LOG] BSIZE parameter value = " + BSIZE);
43         initializeVariables();
44         initializeMatrix(AfileNames, true);
45         initializeMatrix(BfileNames, true);
46         initializeMatrix(CfileNames, false);
47
48         // Wait for initialization
49         COMPSs.barrier();
50
51         // Compute matrix multiplication C = A x B
52         long startTime = System.currentTimeMillis();
53         computeMultiplication();

```

```

55     long estimatedTime = System.currentTimeMillis() - startTime;
56     System.out.println("[TIME] EXECUTION TIME = " + estimatedTime);
57
58     // End
59     System.out.println("[LOG] Main program finished.");
60 }
61
62 /**
63  * Initializes the filenames
64  *
65  */
66 private static void initializeVariables() {
67     AfileNames = new String[MSIZE][MSIZE];
68     BfileNames = new String[MSIZE][MSIZE];
69     CfileNames = new String[MSIZE][MSIZE];
70     for (int i = 0; i < MSIZE; i++) {
71         for (int j = 0; j < MSIZE; j++) {
72             AfileNames[i][j] = "A." + i + "." + j;
73             BfileNames[i][j] = "B." + i + "." + j;
74             CfileNames[i][j] = "C." + i + "." + j;
75         }
76     }
77 }
78
79 /**
80  * Initializes each fileName with random values or 0s
81  *
82  * @param fileNames
83  * @param initRand
84  * @throws IOException
85  */
86 private static void initializeMatrix(String[][] fileNames, boolean initRand)
87     throws IOException {
88     for (int i = 0; i < MSIZE; ++i) {
89         for (int j = 0; j < MSIZE; ++j) {
90             MatmulImpl.initializeBlock(fileNames[i][j], BSIZE, initRand);
91         }
92     }
93 }
94
95 /**
96  * Main loop of matrix multiplication
97  *
98  */
99 private static void computeMultiplication() {
100     System.out.println("[LOG] Computing result");
101     Integer[][][] exitValues = new Integer[MSIZE][MSIZE][MSIZE];
102
103     // Launch tasks
104     for (int i = 0; i < MSIZE; ++i) {
105         for (int k = 0; k < MSIZE; ++k) {
106             for (int j = 0; j < MSIZE; ++j) {
107                 switch (TYPE) {
108                     case 1:
109                         exitValues[i][j][k] = MatmulImpl
110                             .multiplyAccumulativeNative(BSIZE,
111                                 AfileNames[i][k], BfileNames[k][j],
112                                 CfileNames[i][j]);
113                     case 2:
114                         exitValues[i][j][k] = MPI.multiplyAccumulativeMPI(
115                             BSIZE, AfileNames[i][k], BfileNames[k][j],
116                             CfileNames[i][j]);
117                     default:
118                         System.err.println("[ERROR] Invalid type");
119                         System.exit(1);
120                         break;
121                 }
122             }
123         }
124     }
125 }

```

```

127 // Sync: Wait loop
128 for (int i = 0; i < MSIZE; i++) {
129     for (int j = 0; j < MSIZE; j++) {
130         for (int k = 0; k < MSIZE; k++) {
131             if (exitValues[i][j][k] != 0) {
132                 System.err.println(
133                     "[ERROR] Some task failed with exitValue "
134                     + exitValues[i][j][k]);
135             }
136         }
137     }
138 }
139 }
141 }

```

## B.2 MatmulItf.java

```

1 package matmul.files;
2
3 import integratedtoolkit.types.annotations.Constraints;
4 import integratedtoolkit.types.annotations.Parameter;
5 import integratedtoolkit.types.annotations.parameter.Direction;
6 import integratedtoolkit.types.annotations.parameter.Type;
7 import integratedtoolkit.types.annotations.task.MPI;
8 import integratedtoolkit.types.annotations.task.Method;
9
10
11 public interface MatmulItf {
12
13     @Method(declaringClass = "matmul.files.MatmulImpl")
14     @Constraints(computingUnits = "1")
15     void initializeBlock(
16         @Parameter(type = Type.FILE, direction = Direction.OUT) String filename,
17         @Parameter() int BSIZE,
18         @Parameter() boolean initRand
19     );
20
21     @Method(declaringClass = "matmul.files.MatmulImpl")
22     @Constraints(computingUnits = "1")
23     Integer multiplyAccumulativeNative(
24         @Parameter() int bsize,
25         @Parameter(type = Type.FILE, direction = Direction.IN) String aIn,
26         @Parameter(type = Type.FILE, direction = Direction.IN) String bIn,
27         @Parameter(type = Type.FILE, direction = Direction.INOUT) String cOut
28     );
29
30     @MPI(binary = "${MATMUL_BINARY}",
31         mpiRunner = "mpirun",
32         computingNodes = "${NODES_PER_MPI_TASK}")
33     @Constraints(computingUnits = "${CUS_PER_MPI_TASK}")
34     Integer multiplyAccumulativeMPI(
35         @Parameter() int bsize,
36         @Parameter(type = Type.FILE, direction = Direction.IN) String aIn,
37         @Parameter(type = Type.FILE, direction = Direction.IN) String bIn,
38         @Parameter(type = Type.FILE, direction = Direction.INOUT) String cOut
39     );
40 }

```

### B.3 MatmulImpl.java

```

1 package matmul.files;
3 import java.io.FileOutputStream;
import java.io.IOException;
5
7 /**
 * Implementation of Task methods
 *
9 */
public class MatmulImpl {
11
12     private static final byte[] NEW_LINE = "\n".getBytes();
13
14     /**
15      * Initializes a block of size @BSIZE and stores it to @filename
16      *
17      * @param filename
18      * @param BSIZE
19      * @param initRand
20      * @throws IOException
21      */
22     public static void initializeBlock(String filename, int BSIZE,
23         boolean initRand) throws IOException {
24         try (FileOutputStream fos = new FileOutputStream(filename)) {
25             for (int iblock = 0; iblock < BSIZE; ++iblock) {
26                 for (int jblock = 0; jblock < BSIZE; ++jblock) {
27                     double value = (double) 0.0;
28                     if (initRand) {
29                         value = (double) (Math.random() * 10.0);
30                     }
31                     fos.write(String.valueOf(value).getBytes());
32                     fos.write(" ".getBytes());
33                 }
34                 fos.write(NEW_LINE);
35             }
36             fos.write(NEW_LINE);
37         } catch (IOException e) {
38             throw new IOException("[ERROR] Error initializing matrix", e);
39         }
40     }
41
42     /**
43      * Multiplies blocks @aFile and @bFile, accumulates result on block @cFile
44      * and stores it to disk
45      *
46      * @param BSIZE
47      * @param aFile
48      * @param bFile
49      * @param cFile
50      * @return
51      */
52     public static Integer multiplyAccumulativeNative(int BSIZE, String aFile,
53         String bFile, String cFile) {
54         Block a = new Block(BSIZE, aFile);
55         Block b = new Block(BSIZE, bFile);
56         Block c = new Block(BSIZE, cFile);
57
58         c.multiplyAccum(a, b);
59         c.blockToDisk(cFile);
60
61         return 0;
62     }
63 }

```

## B.4 MPI.java

```

1 package mpi;
2
3 /**
4  * Dummy class for MPI tasks
5  *
6  */
7 public class MPI {
8
9     /**
10    * Dummy implementation to call the matmul MPI binary
11    *
12    * @param bSize
13    * @param aIn
14    * @param bIn
15    * @param cOut
16    * @return
17    */
18    public static Integer multiplyAccumulativeMPI(int bSize, String aIn,
19        String bIn, String cOut) {
20        return -1;
21    }
22 }

```

## B.5 Matmul.c

```

1 #include "mpi.h"
2 #include <stdio.h>
3 #include <stdlib.h>
4
5 #define MASTER 0          /* taskid of first task */
6 #define FROM_MASTER 1    /* setting a message type */
7 #define FROM_WORKER 2   /* setting a message type */
8
9 void fillMatrix(const char* fileName, int matrixSize, double* mat) {
10     int i, j;
11
12     FILE *file;
13     file = fopen(fileName, "r");
14
15     // printf(" - Open file %s with size %d.\n", fileName, matrixSize);
16     for(i = 0; i < matrixSize; i++) {
17         for(j = 0; j < matrixSize; j++) {
18             if (!fscanf(file, "%lf", &mat[i*matrixSize + j])) {
19                 break;
20             }
21             // printf("%lf\n", mat[i*matrixSize + j]);
22         }
23     }
24     fclose(file);
25 }
26
27 void storeMatrix(const char* fileName, int matrixSize, const double* mat) {
28     int i, j;
29
30     FILE *file;
31     file = fopen(fileName, "w");
32
33     // printf(" - Store matrix in file %s with size %d.\n", fileName,
34     // matrixSize);
35     for(i = 0; i < matrixSize; i++) {
36         for(j = 0; j < matrixSize; j++) {

```

```

    fprintf(file, "%lf ", mat[i*matrixSize + j]);
39     // printf("%lf\n", mat[i*matrixSize + j]);
    }
41     fprintf(file, "\n");
    }
43     fclose(file);
}
45
int main (int argc, char *argv[]) {
47     /***** Initialize *****/
    if (argc != 5) {
49         // printf("Incorrect usage: matmul <matrixSize> <Ain> <Bin> <Cout>. Quitting...\n")
        ;
        exit(1);
51     }
    int matrixSize = atoi(argv[1]);    // Number of rows/columns in matrixes
53     char* ain = argv[2];            // FileName of Ain
    char* bin = argv[3];            // FileName of Bin
55     char* cout = argv[4];          // FileName of Cout

57     double a[matrixSize*matrixSize]; // Matrix A to be multiplied
    double b[matrixSize*matrixSize]; // Matrix B to be multiplied
59     double c[matrixSize*matrixSize]; // Result matrix C

61     // Initialize MPI env
    int mpiProcs;                    // Number of MPI Nodes
63     int taskid;                     // Task identifier
    MPI_Status status;               // Status for MPI communications
65     MPI_Request send_request;      // For async calls

67     MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &taskid);
69     MPI_Comm_size(MPI_COMM_WORLD, &mpiProcs);

71     // Misc variables
    int i, j, k, dest, rows;
73

    /***** master task *****/
75     if (taskid == MASTER) {
        // printf("Matmul with %d MPI nodes.\n", mpiProcs);
77

        // Initialize arrays
79         // printf("Initialize matrixes.\n");
        fillMatrix(ain, matrixSize, a);
81         fillMatrix(bin, matrixSize, b);
        fillMatrix(cout, matrixSize, c);
83

        // Send matrix data to the worker tasks
85         // printf("Send matrixes to workers.\n");
        int averow = matrixSize/mpiProcs;
87         int extra = matrixSize%mpiProcs;
        int offset = 0;
89         int mtype = FROM_MASTER;
        for (dest = 0; dest < mpiProcs; dest++) {
91             rows = (dest < extra) ? averow+1 : averow;
            MPI_Isend(&offset, 1, MPI_INT, dest, mtype, MPI_COMM_WORLD,
93             &send_request);
            MPI_Isend(&rows, 1, MPI_INT, dest, mtype, MPI_COMM_WORLD,
95             &send_request);
            MPI_Isend(&a[offset*matrixSize], rows*matrixSize, MPI_DOUBLE, dest,
97             mtype, MPI_COMM_WORLD, &send_request);
            MPI_Isend(&b, matrixSize*matrixSize, MPI_DOUBLE, dest, mtype,
99             MPI_COMM_WORLD, &send_request);
            MPI_Isend(&c, rows*matrixSize, MPI_DOUBLE, dest, mtype,
101             MPI_COMM_WORLD, &send_request);
            offset = offset + rows;
103         }
    }
105

    /***** worker task *****/
107     // Receive matrix
    // printf("Receive IN matrixes on process %d.\n", taskid);

```

```

109  int offset = 0;
110  int mtype = FROM_MASTER;
111  MPI_Recv(&offset, 1, MPI_INT, MASTER, mtype, MPI_COMM_WORLD, &status);
112  MPI_Recv(&rows, 1, MPI_INT, MASTER, mtype, MPI_COMM_WORLD, &status);
113  MPI_Recv(&a, rows*matrixSize, MPI_DOUBLE, MASTER, mtype, MPI_COMM_WORLD,
114  &status);
115  MPI_Recv(&b, matrixSize*matrixSize, MPI_DOUBLE, MASTER, mtype,
116  MPI_COMM_WORLD, &status);
117  MPI_Recv(&c[offset*matrixSize], rows*matrixSize, MPI_DOUBLE, MASTER, mtype,
118  MPI_COMM_WORLD, &status);
119
120  // Perform multiply accumulative
121  // printf("Perform multiply accumulative on process %d.\n", taskid);
122  for (i = 0; i < rows; i++) {
123      for (k = 0; k < matrixSize; k++) {
124          for (j = 0; j < matrixSize; j++) {
125              c[i*matrixSize + j] = c[i*matrixSize + j]
126              + a[i*matrixSize + k]*b[k*matrixSize + j];
127          }
128      }
129  }
130
131  // Send back result to master
132  // printf("Send result back to master on process %d.\n", taskid);
133  mtype = FROM_WORKER;
134  MPI_Isend(&offset, 1, MPI_INT, MASTER, mtype, MPI_COMM_WORLD,
135  &send_request);
136  MPI_Isend(&rows, 1, MPI_INT, MASTER, mtype, MPI_COMM_WORLD, &send_request);
137  MPI_Isend(&c, rows*matrixSize, MPI_DOUBLE, MASTER, mtype, MPI_COMM_WORLD,
138  &send_request);
139
140  /***** master task *****/
141  if (taskid == MASTER) {
142      // Receive results from worker tasks
143      // printf("Receive results.\n");
144      mtype = FROM_WORKER;
145      for (i = 0; i < mpiProcs; i++) {
146          MPI_Recv(&offset, 1, MPI_INT, i, mtype, MPI_COMM_WORLD, &status);
147          MPI_Recv(&rows, 1, MPI_INT, i, mtype, MPI_COMM_WORLD, &status);
148          MPI_Recv(&c[offset*matrixSize], rows*matrixSize, MPI_DOUBLE, i,
149  mtype, MPI_COMM_WORLD, &status);
150          // printf(" - Received results from task %d\n", i);
151      }
152
153      // Print result
154      //printMatrix(matrixSize, c);
155
156      // Store result to file
157      // printf("Store result matrix.\n");
158      storeMatrix(cout, matrixSize, c);
159
160      printf ("Done.\n");
161  }
162
163  /***** Finalize *****/
164  MPI_Finalize();
165  }

```

## B.6 Block.java

```

package matmul.files;
2
import java.io.BufferedReader;
4 import java.io.FileNotFoundException;
import java.io.FileOutputStream;
6 import java.io.FileReader;
import java.io.IOException;

```

```

8 import java.util.StringTokenizer;
10 public class Block {
12     private static final byte[] NEW_LINE = "\n".getBytes();
14     private final int bSize;
15     private final double[][] data;
16
17     /**
18      * Initialization of a block of size @bSize from file @filename
19      *
20      * @param bSize
21      * @param filename
22      */
23     public Block(int bSize, String filename) {
24         this.bSize = bSize;
25         this.data = new double[this.bSize][this.bSize];
26
27         // Retrieve values from file
28         try (BufferedReader br = new BufferedReader(new FileReader(filename))) {
29             StringTokenizer tokens;
30             String nextLine;
31
32             for (int i = 0; i < this.bSize; i++) {
33                 nextLine = br.readLine();
34                 tokens = new StringTokenizer(nextLine);
35                 for (int j = 0; j < bSize && tokens.hasMoreTokens(); j++) {
36                     this.data[i][j] = Double.parseDouble(tokens.nextToken());
37                 }
38             }
39         } catch (FileNotFoundException fnfe) {
40             fnfe.printStackTrace();
41             System.exit(-1);
42         } catch (IOException ioe) {
43             ioe.printStackTrace();
44             System.exit(-1);
45         }
46     }
47
48     /**
49      * Dumps the content of the block to file @filename
50      *
51      * @param filename
52      */
53     public void blockToDisk(String filename) {
54         try (FileOutputStream fos = new FileOutputStream(filename)) {
55             for (int i = 0; i < this.bSize; i++) {
56                 for (int j = 0; j < this.bSize; j++) {
57                     String str = String.valueOf(this.data[i][j]) + " ";
58                     fos.write(str.getBytes());
59                 }
60                 fos.write(NEW_LINE);
61             }
62         } catch (FileNotFoundException fnfe) {
63             fnfe.printStackTrace();
64             System.exit(-1);
65         } catch (IOException ioe) {
66             ioe.printStackTrace();
67             System.exit(-1);
68         }
69     }
70
71     /**
72      * Multiplies blocks @a and @b and accumulates result on the current
73      * block
74      *
75      * @param a
76      * @param b
77      */
78     public void multiplyAccum(Block a, Block b) {

```



```
80     for (int i = 0; i < this.bSize; ++i) {  
81         for (int k = 0; k < this.bSize; ++k) {  
82             for (int j = 0; j < this.bSize; ++j) {  
83                 this.data[i][j] += a.data[i][k] * b.data[k][j];  
84             }  
85         }  
86     }  
87 }  
88 }
```



## Appendix C

# NMMB/BSC-Dust: code highlights

### C.1 Nmmb.java

```

1 package nmmb;
3 import java.io.File;
  import java.util.Date;
5
  import org.apache.commons.configuration.ConfigurationException;
7 import org.apache.logging.log4j.LogManager;
  import org.apache.logging.log4j.Logger;
9
  import binary.BINARY;
11 import mpi.MPI;
  import nmmb.configuration.NMMBConfigManager;
13 import nmmb.configuration.NMMBConstants;
  import nmmb.configuration.NMMBEnvironment;
15 import nmmb.configuration.NMMBParameters;
  import nmmb.exceptions.MainExecutionException;
17 import nmmb.exceptions.TaskExecutionException;
  import nmmb.loggers.LoggerNames;
19 import nmmb.utils.FileManagement;
  import nmmb.utils.FortranWrapper;
21 import nmmb.utils.MessagePrinter;
23
  public class Nmmb {
25     // Loggers
  private static final Logger LOGGER_MAIN = LogManager.getLogger(LoggerNames.NMMB_MAIN);
27 private static final Logger LOGGER_FIXED = LogManager
    .getLogger(LoggerNames.NMMB_FIXED);
29 private static final Logger LOGGER_VARIABLE = LogManager
    .getLogger(LoggerNames.NMMB_VARIABLE);
31 private static final Logger LOGGER_UMO_MODEL = LogManager
    .getLogger(LoggerNames.NMMB_UMO_MODEL);
33 private static final Logger LOGGER_POST = LogManager.getLogger(LoggerNames.NMMB_POST);
35
  /**
  * Prints the usage
  *
  */
39 private static void usage() {
  LOGGER_MAIN.info("Invalid parameters for nmmb.Nmmb");
41 LOGGER_MAIN.info("  Usage: nmmb.Nmmb <configFilePath>");
  }
43
  /**
  * FIXED STEP
  */
47 private static void doFixed(NMMBParameters nmmbParams) throws TaskExecutionException {
  LOGGER_FIXED.info("Enter fixed process");
49
  // Prepare execution
51 nmmbParams.prepareFixedExecution();
  MessagePrinter fixedMP = new MessagePrinter(LOGGER_FIXED);
53

```

```

55 // Build the fortran executables
56 if (nmmBParams.isCompileBinaries()) {
57     Integer[] compilationEvs = new Integer[FortranWrapper.FIXED_FORTRAN_F90_FILES.length
58         + FortranWrapper.FIXED_FORTRAN_F_FILES.length];
59     int i = 0;
60     fixedMP.printInfoMsg("Building fixed executables");
61     for (String fortranFile : FortranWrapper.FIXED_FORTRAN_F90_FILES) {
62         String executable = NMMBEnvironment.FIX + fortranFile
63             + FortranWrapper.SUFFIX_EXE;
64         String src = NMMBEnvironment.FIX + fortranFile
65             + FortranWrapper.SUFFIX_F90_SRC;
66
67         compilationEvs[i++] = BINARY.fortranCompiler(FortranWrapper.MC_FLAG,
68             FortranWrapper.SHARED_FLAG, FortranWrapper.CONVERT_PREFIX,
69             FortranWrapper.CONVERT_VALUE, FortranWrapper.TRACEBACK_FLAG,
70             FortranWrapper.ASSUME_PREFIX, FortranWrapper.ASSUME_VALUE,
71             FortranWrapper.OPT_FLAG, FortranWrapper.FPMODEL_PREFIX,
72             FortranWrapper.FPMODEL_VALUE, FortranWrapper.STACK_FLAG,
73             FortranWrapper.OFLAG, executable, src);
74     }
75     for (String fortranFile : FortranWrapper.FIXED_FORTRAN_F_FILES) {
76         String executable = NMMBEnvironment.FIX + fortranFile
77             + FortranWrapper.SUFFIX_EXE;
78         String src = NMMBEnvironment.FIX + fortranFile
79             + FortranWrapper.SUFFIX_F_SRC;
80         compilationEvs[i++] = BINARY.fortranCompiler(FortranWrapper.MC_FLAG,
81             FortranWrapper.SHARED_FLAG, FortranWrapper.CONVERT_PREFIX,
82             FortranWrapper.CONVERT_VALUE, FortranWrapper.TRACEBACK_FLAG,
83             FortranWrapper.ASSUME_PREFIX, FortranWrapper.ASSUME_VALUE,
84             FortranWrapper.OPT_FLAG, FortranWrapper.FPMODEL_PREFIX,
85             FortranWrapper.FPMODEL_VALUE, FortranWrapper.STACK_FLAG,
86             FortranWrapper.OFLAG, executable, src);
87     }
88     // Sync master to wait for compilation
89     for (i = 0; i < compilationEvs.length; ++i) {
90         LOGGER_FIXED.debug("Compilation of " + i + " binary ended with status "
91             + compilationEvs[i]);
92         if (compilationEvs[i] != 0) {
93             throw new TaskExecutionException(
94                 "[ERROR] Error compiling binary " + i);
95         }
96     }
97     fixedMP.printInfoMsg("Finished building fixed executables");
98 }
99
100 // Begin binary calls
101 fixedMP.printHeaderMsg("BEGIN");
102
103 final int NUM_BINARIES = 16;
104 Integer[] fixedBinariesEvs = new Integer[NUM_BINARIES];
105 int i = 0;
106
107 fixedMP.printInfoMsg("Generate DEM height and sea mask files");
108 String topoDir = NMMBEnvironment.GEODATA_DIR + "topolkmDEM" + File.separator;
109 String seamaskDEM = NMMBEnvironment.OUTPUT + "seamaskDEM";
110 String heightDEM = NMMBEnvironment.OUTPUT + "heightDEM";
111 fixedBinariesEvs[i++] = BINARY.smmount(topoDir, seamaskDEM, heightDEM);
112
113 fixedMP.printInfoMsg("Generate landuse file");
114 String landuseDataDir = NMMBEnvironment.GEODATA_DIR + "landuse_30s"
115     + File.separator;
116 String landuse = NMMBEnvironment.OUTPUT + "landuse";
117 String kount_landuse = NMMBEnvironment.OUTPUT + "kount_landuse";
118 fixedBinariesEvs[i++] = BINARY.landuse(landuseDataDir, landuse, kount_landuse);
119
120 fixedMP.printInfoMsg("Generate landusenew file");
121 String landusenew = NMMBEnvironment.OUTPUT + "landusenew";
122 String kount_landusenew = NMMBEnvironment.OUTPUT + "kount_landusenew";
123 fixedBinariesEvs[i++] = BINARY.landusenew(NMMBEnvironment.GTOPO_DIR, landusenew,
124     kount_landusenew);
125
126 fixedMP.printInfoMsg("Generate mountains");

```

```
127 String topo30sDir = NMMBEnvironment.GEODATA_DIR + "topo_30s" + File.separator;
String heightmean = NMMBEnvironment.OUTPUT + "heightmean";
fixedBinariesEvs[i++] = BINARY.topo(topo30sDir, heightmean);
129
fixedMP.printInfoMsg("Generate standard deviation of topography height");
131 String stdh = NMMBEnvironment.OUTPUT + "stdh";
fixedBinariesEvs[i++] = BINARY.stdh(heightmean, seamaskDEM, topo30sDir, stdh);
133
fixedMP.printInfoMsg("Generate envelope mountains");
135 String height = NMMBEnvironment.OUTPUT + "height";
fixedBinariesEvs[i++] = BINARY.envelope(heightmean, stdh, height);
137
fixedMP.printInfoMsg("Generate top soil type file");
139 String soiltypeDir = NMMBEnvironment.GEODATA_DIR + "soiltype_top_30s"
+ File.separator;
141 String topsoiltype = NMMBEnvironment.OUTPUT + "topsoiltype";
fixedBinariesEvs[i++] = BINARY.topsoiltype(seamaskDEM, soiltypeDir, topsoiltype);
143
fixedMP.printInfoMsg("Generate bottom soil type file");
145 String soiltypePath = NMMBEnvironment.GEODATA_DIR + "soiltype_bot_30s"
+ File.separator;
147 String botsoiltype = NMMBEnvironment.OUTPUT + "botsoiltype";
fixedBinariesEvs[i++] = BINARY.botsoiltype(seamaskDEM, soiltypePath, botsoiltype);
149
fixedMP.printInfoMsg("Generate sea mask and reprocess mountains");
151 String seamask = NMMBEnvironment.OUTPUT + "seamask";
fixedBinariesEvs[i++] = BINARY.toposeamask(seamaskDEM, seamask, height, landuse,
153 topsoiltype, botsoiltype);
155
fixedMP.printInfoMsg("Reprocess standard deviation of topography height");
fixedBinariesEvs[i++] = BINARY.stdhtopo(seamask, stdh);
157
fixedMP.printInfoMsg("Generate deep soil temperature");
159 String soiltempPath = NMMBEnvironment.GEODATA_DIR + "soiltemp_1deg"
+ File.separator;
161 String deeptemperature = NMMBEnvironment.OUTPUT + "deeptemperature";
fixedBinariesEvs[i++] = BINARY.deeptemperature(seamask, soiltempPath,
163 deeptemperature);
165
fixedMP.printInfoMsg("Generate maximum snow albedo");
String maxsnowalbDir = NMMBEnvironment.GEODATA_DIR + "maxsnowalb"
+ File.separator;
167 String snowalbedo = NMMBEnvironment.OUTPUT + "snowalbedo";
169 fixedBinariesEvs[i++] = BINARY.snowalbedo(maxsnowalbDir, snowalbedo);
171
fixedMP.printInfoMsg("Generate vertical coordinate");
String dsg = NMMBEnvironment.OUTPUT + "dsg";
173 fixedBinariesEvs[i++] = BINARY.vcgenerator(dsg);
175
fixedMP.printInfoMsg("Generate highres roughness length for africa and asia");
String roughnessDir = NMMBEnvironment.GEODATA_DIR + "roughness_025s"
+ File.separator;
177 String roughness = NMMBEnvironment.OUTPUT + "roughness";
179 fixedBinariesEvs[i++] = BINARY.roughness(roughnessDir, roughness);
181
fixedMP.printInfoMsg("Generate co2 files");
String co2_data_dir = NMMBEnvironment.GEODATA_DIR + "co2data" + File.separator;
183 String co2_trans = NMMBEnvironment.OUTPUT + "co2_trans";
fixedBinariesEvs[i++] = BINARY.gfdlco2(dsg, co2_data_dir, co2_trans);
185
fixedMP.printInfoMsg(
187 "Generate lookup tables for aerosol scavenging collection efficiencies");
String lookup_aerosol2_rh00 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh00";
189 String lookup_aerosol2_rh50 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh50";
String lookup_aerosol2_rh70 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh70";
191 String lookup_aerosol2_rh80 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh80";
String lookup_aerosol2_rh90 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh90";
193 String lookup_aerosol2_rh95 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh95";
String lookup_aerosol2_rh99 = NMMBEnvironment.OUTPUT + "lookup_aerosol2.dat.rh99";
195 fixedBinariesEvs[i++] = BINARY.run_aerosol(nmmbParams.isCompileBinaries(),
nmmbParams.isCleanBinaries(), lookup_aerosol2_rh00, lookup_aerosol2_rh50,
197 lookup_aerosol2_rh70, lookup_aerosol2_rh80, lookup_aerosol2_rh90,
```

```

        lookup_aerosol2_rh95, lookup_aerosol2_rh99);
199
// Wait for binaries completion and check exit value
201 for (i = 0; i < fixedBinariesEvs.length; ++i) {
    LOGGER_FIXED.debug("Execution of " + i + " binary ended with status "
203         + fixedBinariesEvs[i]);
    if (fixedBinariesEvs[i] != 0) {
205         throw new TaskExecutionException("[ERROR] Error executing binary " + i);
    }
207 }

// Clean Up binaries
209 if (nmmBParams.isCleanBinaries()) {
    fixedMP.printInfoMsg("Clean up executables");
211 for (String fortranFile : FortranWrapper.FIXED_FORTRAN_F90_FILES) {
    String executable = NMMBEnvironment.FIX + fortranFile
213         + FortranWrapper.SUFFIX_EXE;
    File f = new File(executable);
215     if (f.exists()) {
        f.delete();
217     }
    }
219 for (String fortranFile : FortranWrapper.FIXED_FORTRAN_F_FILES) {
    String executable = NMMBEnvironment.FIX + fortranFile
221         + FortranWrapper.SUFFIX_EXE;
    File f = new File(executable);
223     if (f.exists()) {
        f.delete();
225     }
    }
227 }

/* End *****/
229 fixedMP.printHeaderMsg("END");
231

LOGGER_FIXED.info("Fixed process finished");
233 }

/*
235 * VARIABLE STEP
237 */
239 private static void compileVariable() throws TaskExecutionException {
    // Build the fortran objects
241 Integer[] depCompilationEvs = new Integer[FortranWrapper.VARIABLE_FORTRAN_F90_DEP_FILES
    .length];
    int objectIndex = 0;
243 for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F90_DEP_FILES) {
    String moduleDir = NMMBEnvironment.VRB;
245     String object = NMMBEnvironment.VRB + fortranFile
        + FortranWrapper.SUFFIX_OBJECT;
247     String src = NMMBEnvironment.VRB + fortranFile
        + FortranWrapper.SUFFIX_F90_SRC;
249

    depCompilationEvs[objectIndex++] = BINARY.fortranCompileObject(
251         FortranWrapper.MC_FLAG, FortranWrapper.SHARED_FLAG,
        FortranWrapper.CONVERT_PREFIX, FortranWrapper.CONVERT_VALUE,
253         FortranWrapper.TRACEBACK_FLAG, FortranWrapper.ASSUME_PREFIX,
        FortranWrapper.ASSUME_VALUE, FortranWrapper.OPT_FLAG,
255         FortranWrapper.FPMODEL_PREFIX, FortranWrapper.FPMODEL_VALUE,
        FortranWrapper.STACK_FLAG, FortranWrapper.CFLAG, src,
257         FortranWrapper.OFLAG, object, FortranWrapper.MODULE_FLAG, moduleDir);
    }
259 // Sync to check compilation status (dependency with task object is
    // also respected if this sync is erased)
261 for (int i = 0; i < depCompilationEvs.length; ++i) {
    LOGGER_VARIABLE.debug("Compilation of " + i
263         + " dependant binary ended with status " + depCompilationEvs[i]);
    if (depCompilationEvs[i] != 0) {
265         throw new TaskExecutionException("[ERROR] Error compiling binary " + i);
    }
267 }

```

```

269 // Build the fortran executables
Integer[] compilationEvs = new Integer[FortranWrapper.VARIABLE_FORTRAN_F90_FILES.length
271 + FortranWrapper.VARIABLE_FORTRAN_F_FILES.length
+ FortranWrapper.VARIABLE_GFORTRAN_F_FILES.length
273 + FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_W3.length
+ FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_DEPS.length + 1];
275
int executableIndex = 0;
277 for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F90_FILES) {
String executable = NMMBEnvironment.VRB + fortranFile
279 + FortranWrapper.SUFFIX_EXE;
String src = NMMBEnvironment.VRB + fortranFile
281 + FortranWrapper.SUFFIX_F90_SRC;

283 compilationEvs[executableIndex++] = BINARY.fortranCompiler(
FortranWrapper.MC_FLAG, FortranWrapper.SHARED_FLAG,
285 FortranWrapper.CONVERT_PREFIX, FortranWrapper.CONVERT_VALUE,
FortranWrapper.TRACEBACK_FLAG, FortranWrapper.ASSUME_PREFIX,
287 FortranWrapper.ASSUME_VALUE, FortranWrapper.OPT_FLAG,
FortranWrapper.FPMODEL_PREFIX, FortranWrapper.FPMODEL_VALUE,
289 FortranWrapper.STACK_FLAG, FortranWrapper.OFLAG, executable, src);
}
291 for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES) {
String executable = NMMBEnvironment.VRB + fortranFile
293 + FortranWrapper.SUFFIX_EXE;
String src = NMMBEnvironment.VRB + fortranFile + FortranWrapper.SUFFIX_F_SRC;
295 compilationEvs[executableIndex++] = BINARY.fortranCompiler(
FortranWrapper.MC_FLAG, FortranWrapper.SHARED_FLAG,
297 FortranWrapper.CONVERT_PREFIX, FortranWrapper.CONVERT_VALUE,
FortranWrapper.TRACEBACK_FLAG, FortranWrapper.ASSUME_PREFIX,
299 FortranWrapper.ASSUME_VALUE, FortranWrapper.OPT_FLAG,
FortranWrapper.FPMODEL_PREFIX, FortranWrapper.FPMODEL_VALUE,
301 FortranWrapper.STACK_FLAG, FortranWrapper.OFLAG, executable, src);
}
303
305 for (String fortranFile : FortranWrapper.VARIABLE_GFORTRAN_F_FILES) {
String executable = NMMBEnvironment.VRB + fortranFile
+ FortranWrapper.SUFFIX_EXE;
307 String src = NMMBEnvironment.VRB + fortranFile + FortranWrapper.SUFFIX_F_SRC;
compilationEvs[executableIndex++] = BINARY.gfortranCompiler(
309 FortranWrapper.BIG_O_FLAG, src, FortranWrapper.OFLAG, executable);
}
311
313 for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_W3) {
String executable = NMMBEnvironment.VRB + fortranFile
+ FortranWrapper.SUFFIX_EXE;
315 String src = NMMBEnvironment.VRB + fortranFile + FortranWrapper.SUFFIX_F_SRC;
String w3LibFlag = "-L" + NMMBEnvironment.UMO_LIBS
+ FortranWrapper.W3_LIB_DIR;
317 String bacioLibFlag = "-L" + NMMBEnvironment.UMO_LIBS
+ FortranWrapper.BACIO_LIB_DIR;
319 compilationEvs[executableIndex++] = BINARY.fortranCompilerWithW3(
FortranWrapper.MC_FLAG, FortranWrapper.SHARED_FLAG,
321 FortranWrapper.CONVERT_PREFIX, FortranWrapper.CONVERT_VALUE,
FortranWrapper.TRACEBACK_FLAG, FortranWrapper.ASSUME_PREFIX,
323 FortranWrapper.ASSUME_VALUE, FortranWrapper.OPT_FLAG,
FortranWrapper.FPMODEL_PREFIX, FortranWrapper.FPMODEL_VALUE,
325 FortranWrapper.STACK_FLAG, FortranWrapper.OFLAG, executable, src,
w3LibFlag, bacioLibFlag, FortranWrapper.W3_FLAG,
327 FortranWrapper.BACIO_FLAG);
}
329
331 for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_DEPS) {
String executable = NMMBEnvironment.VRB + fortranFile
+ FortranWrapper.SUFFIX_EXE;
333 String src = NMMBEnvironment.VRB + fortranFile + FortranWrapper.SUFFIX_F_SRC;
String object = NMMBEnvironment.VRB + FortranWrapper.MODULE_FLT
+ FortranWrapper.SUFFIX_OBJECT;
335 compilationEvs[executableIndex++] = BINARY.fortranCompileWithObject(
FortranWrapper.MC_FLAG, FortranWrapper.SHARED_FLAG,
337 FortranWrapper.CONVERT_PREFIX, FortranWrapper.CONVERT_VALUE,
FortranWrapper.TRACEBACK_FLAG, FortranWrapper.ASSUME_PREFIX,
339 FortranWrapper.ASSUME_VALUE, FortranWrapper.OPT_FLAG,

```

```

341         FortranWrapper.FPMODEL_PREFIX, FortranWrapper.FPMODEL_VALUE,
342         FortranWrapper.STACK_FLAG, FortranWrapper.OFLAG, executable, src,
343         object);
344     }
345     String source = NMMBEnvironment.VRB + FortranWrapper.READ_PAUL_SOURCE
346         + FortranWrapper.SUFFIX_F_SRC;
347     String executable = NMMBEnvironment.VRB + FortranWrapper.READ_PAUL_SOURCE
348         + FortranWrapper.SUFFIX_EXE;
349     compilationEvs[executableIndex++] = BINARY.compileReadPaulSource(source,
350         executable);
351
352     // Sync master to wait for compilation
353     for (int i = 0; i < compilationEvs.length; ++i) {
354         LOGGER_VARIABLE.debug("Compilation of " + i + " binary ended with status "
355             + compilationEvs[i]);
356         if (compilationEvs[i] != 0) {
357             throw new TaskExecutionException("[ERROR] Error compiling binary " + i);
358         }
359     }
360 }
361
362 private static void cleanUpVariableExe() {
363     for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F90_DEP_FILES) {
364         String executable = NMMBEnvironment.VRB + fortranFile
365             + FortranWrapper.SUFFIX_EXE;
366         File f = new File(executable);
367         if (f.exists()) {
368             f.delete();
369         }
370     }
371     for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F90_FILES) {
372         String executable = NMMBEnvironment.VRB + fortranFile
373             + FortranWrapper.SUFFIX_EXE;
374         File f = new File(executable);
375         if (f.exists()) {
376             f.delete();
377         }
378     }
379     for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES) {
380         String executable = NMMBEnvironment.VRB + fortranFile
381             + FortranWrapper.SUFFIX_EXE;
382         File f = new File(executable);
383         if (f.exists()) {
384             f.delete();
385         }
386     }
387     for (String fortranFile : FortranWrapper.VARIABLE_GFORTRAN_F_FILES) {
388         String executable = NMMBEnvironment.VRB + fortranFile
389             + FortranWrapper.SUFFIX_EXE;
390         File f = new File(executable);
391         if (f.exists()) {
392             f.delete();
393         }
394     }
395     for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_W3) {
396         String executable = NMMBEnvironment.VRB + fortranFile
397             + FortranWrapper.SUFFIX_EXE;
398         File f = new File(executable);
399         if (f.exists()) {
400             f.delete();
401         }
402     }
403     for (String fortranFile : FortranWrapper.VARIABLE_FORTRAN_F_FILES_WITH_DEPS) {
404         String executable = NMMBEnvironment.VRB + fortranFile
405             + FortranWrapper.SUFFIX_EXE;
406         File f = new File(executable);
407         if (f.exists()) {
408             f.delete();
409         }
410     }
411     String readPaulSource = NMMBEnvironment.VRB + FortranWrapper.READ_PAUL_SOURCE
412         + FortranWrapper.SUFFIX_EXE;

```



```

413     File f = new File(readPaulSource);
         if (f.exists()) {
415         f.delete();
         }
417     }

419     private static void doVariable(NMMBParameters nmmbParams, Date currentDate)
         throws TaskExecutionException {
421         LOGGER_VARIABLE.info("Enter variable process");

423         // Prepare execution
         nmmbParams.prepareVariableExecution(currentDate);
425         MessagePrinter variableMP = new MessagePrinter(LOGGER_VARIABLE);

427         /* Compile *****/
         if (nmmbParams.isCompileBinaries()) {
429             variableMP.printInfoMsg("Building variable executables");
                 compileVariable();
431             variableMP.printInfoMsg("Finished building variable executables");
         }

433         // Set variables for binary calls
435         variableMP.printHeaderMsg("BEGIN");

437         String CW = NMMBEnvironment.OUTPUT + "00_CW.dump";
         String ICEC = NMMBEnvironment.OUTPUT + "00_ICEC.dump";
439         String SH = NMMBEnvironment.OUTPUT + "00_SH.dump";
         String SOILT2 = NMMBEnvironment.OUTPUT + "00_SOILT2.dump";
441         String SOILT4 = NMMBEnvironment.OUTPUT + "00_SOILT4.dump";
         String SOILW2 = NMMBEnvironment.OUTPUT + "00_SOILW2.dump";
443         String SOILW4 = NMMBEnvironment.OUTPUT + "00_SOILW4.dump";
         String TT = NMMBEnvironment.OUTPUT + "00_TT.dump";
445         String VV = NMMBEnvironment.OUTPUT + "00_VV.dump";
         String HH = NMMBEnvironment.OUTPUT + "00_HH.dump";
447         String PRMSL = NMMBEnvironment.OUTPUT + "00_PRMSL.dump";
         String SOILT1 = NMMBEnvironment.OUTPUT + "00_SOILT1.dump";
449         String SOILT3 = NMMBEnvironment.OUTPUT + "00_SOILT3.dump";
         String SOILW1 = NMMBEnvironment.OUTPUT + "00_SOILW1.dump";
451         String SOILW3 = NMMBEnvironment.OUTPUT + "00_SOILW3.dump";
         String SST_TS = NMMBEnvironment.OUTPUT + "00_SST_TS.dump";
453         String UU = NMMBEnvironment.OUTPUT + "00_UU.dump";
         String WEASD = NMMBEnvironment.OUTPUT + "00_WEASD.dump";

455         String GFS_file = NMMBEnvironment.OUTPUT + "131140000.gfs";

457         String deco = NMMBEnvironment.VRB;

459         String llspl000 = NMMBEnvironment.OUTPUT + "llspl.000";
         String outtmp = NMMBEnvironment.OUTPUT + "llstmp";
         String outmst = NMMBEnvironment.OUTPUT + "llsmst";
463         String outsst = NMMBEnvironment.OUTPUT + "llgsst";
         String outsno = NMMBEnvironment.OUTPUT + "llgsno";
465         String outcic = NMMBEnvironment.OUTPUT + "llgcic";

467         String llgsst05 = NMMBEnvironment.OUTPUT + "llgsst05";
         String sstfileinPath = NMMBEnvironment.OUTPUT + "sst2dvar_grb_0.5";

469         String seamask = NMMBEnvironment.OUTPUT + "seamask";
         String albedo = NMMBEnvironment.OUTPUT + "albedo";
         String albedobase = NMMBEnvironment.OUTPUT + "albedobase";
473         String albedomnth = NMMBEnvironment.GEODATA_DIR + "albedo" + File.separator
             + "albedomnth";

475         String albedorrtn = NMMBEnvironment.OUTPUT + "albedorrtn";
         String albedorrtnldegDir = NMMBEnvironment.GEODATA_DIR + "albedo_rrtnldeg"
             + File.separator;

477         String vegfrac = NMMBEnvironment.OUTPUT + "vegfrac";
         String vegfracmnth = NMMBEnvironment.GEODATA_DIR + "vegfrac" + File.separator
             + "vegfracmnth";

483         String landuse = NMMBEnvironment.OUTPUT + "landuse";

```

```

485 String topsoiltype = NMMBEnvironment.OUTPUT + "topsoiltype";
486 String height = NMMBEnvironment.OUTPUT + "height";
487 String stdh = NMMBEnvironment.OUTPUT + "stdh";
488 String z0base = NMMBEnvironment.OUTPUT + "z0base";
489 String z0 = NMMBEnvironment.OUTPUT + "z0";
490 String ustar = NMMBEnvironment.OUTPUT + "ustar";
491
492 String sst05 = NMMBEnvironment.OUTPUT + "sst05";
493 String deeptemperature = NMMBEnvironment.OUTPUT + "deeptemperature";
494 String snowalbedo = NMMBEnvironment.OUTPUT + "snowalbedo";
495 String landusenew = NMMBEnvironment.OUTPUT + "landusenew";
496 String llgsst = NMMBEnvironment.OUTPUT + "llgsst";
497 String llgsno = NMMBEnvironment.OUTPUT + "llgsno";
498 String llgcic = NMMBEnvironment.OUTPUT + "llgcic";
499 String llsmst = NMMBEnvironment.OUTPUT + "llsmst";
500 String llstmp = NMMBEnvironment.OUTPUT + "llstmp";
501 String albedorrtmcorr = NMMBEnvironment.OUTPUT + "albedorrtmcorr";
502 String dzsoil = NMMBEnvironment.OUTPUT + "dzsoil";
503 String tskin = NMMBEnvironment.OUTPUT + "tskin";
504 String sst = NMMBEnvironment.OUTPUT + "sst";
505 String snow = NMMBEnvironment.OUTPUT + "snow";
506 String snowheight = NMMBEnvironment.OUTPUT + "snowheight";
507 String cice = NMMBEnvironment.OUTPUT + "cice";
508 String seamaskcorr = NMMBEnvironment.OUTPUT + "seamaskcorr";
509 String landusecorr = NMMBEnvironment.OUTPUT + "landusecorr";
510 String landusenewcorr = NMMBEnvironment.OUTPUT + "landusenewcorr";
511 String topsoiltypecorr = NMMBEnvironment.OUTPUT + "topsoiltypecorr";
512 String vegfraccorr = NMMBEnvironment.OUTPUT + "vegfraccorr";
513 String z0corr = NMMBEnvironment.OUTPUT + "z0corr";
514 String z0basecorr = NMMBEnvironment.OUTPUT + "z0basecorr";
515 String emissivity = NMMBEnvironment.OUTPUT + "emissivity";
516 String canopywater = NMMBEnvironment.OUTPUT + "canopywater";
517 String frozenprecratio = NMMBEnvironment.OUTPUT + "frozenprecratio";
518 String smst = NMMBEnvironment.OUTPUT + "smst";
519 String sh2o = NMMBEnvironment.OUTPUT + "sh2o";
520 String stmp = NMMBEnvironment.OUTPUT + "stmp";
521 String dsq = NMMBEnvironment.OUTPUT + "dsq";
522 String fcst = NMMBEnvironment.OUTPUT + "fcst";
523 String fcstDir = NMMBEnvironment.OUTPUT + "fcst";
524 String bocoPrefix = NMMBEnvironment.OUTPUT + "boco.";
525 String llsplPrefix = NMMBEnvironment.OUTPUT + "llspl.";
526
527 String source = NMMBEnvironment.OUTPUT + "source";
528 String sourceNETCDF = NMMBEnvironment.OUTPUT + "source.nc";
529 String sourceNCIncludeDir = NMMBEnvironment.VRB_INCLUDE_DIR;
530
531 String soildust = NMMBEnvironment.OUTPUT + "soildust";
532 String kount_landuse = NMMBEnvironment.OUTPUT + "kount_landuse";
533 String kount_landusenew = NMMBEnvironment.OUTPUT + "kount_landusenew";
534 String roughness = NMMBEnvironment.OUTPUT + "roughness";
535
536 // Begin binary calls
537 final int NUM_BINARIES = 12;
538 Integer[] variableBinariesEvs = new Integer[NUM_BINARIES];
539 int binaryIndex = 0;
540
541 variableMP.printInfoMsg("degrib gfs global data");
542 variableBinariesEvs[binaryIndex++] = BINARY.degribgfs_generic_05(CW, ICEC, SH,
543     SOILT2, SOILT4, SOILW2, SOILW4, TT, VV, HH, PRMSL, SOILT1, SOILT3, SOILW1,
544     SOILW3, SST_TS, UU, WEASD);
545
546 variableMP.printInfoMsg("GFS 2 Model");
547 variableBinariesEvs[binaryIndex++] = BINARY.gfs2model_rrtm(CW, ICEC, SH, SOILT2,
548     SOILT4, SOILW2, SOILW4, TT, VV, HH, PRMSL, SOILT1, SOILT3, SOILW1, SOILW3,
549     SST_TS, UU, WEASD, GFS_file);
550
551 variableMP.printInfoMsg("INC RRTM");
552 variableBinariesEvs[binaryIndex++] = BINARY.inc_rrtm(GFS_file, deco);
553
554 variableMP.printInfoMsg("CNV RRTM");
555 variableBinariesEvs[binaryIndex++] = BINARY.cnv_rrtm(GFS_file, llspl000, outtmp,
556     outmst, outsst, outsno, outcic);

```

```

557     variableMP.printInfoMsg("Degrib 0.5 deg sst");
559     variableBinariesEvs[binaryIndex++] = BINARY.degribsst(llgsst05, sstfileinPath);

561     variableMP.printInfoMsg("Prepare climatological albedo");
563     variableBinariesEvs[binaryIndex++] = BINARY.albedo(llspl000, seamask, albedo,
        albedobase, albedomnth);

565     variableMP.printInfoMsg("Prepare rrtm climatological albedos");
567     variableBinariesEvs[binaryIndex++] = BINARY.albedorrtm(llspl000, seamask,
        albedorrtm, albedorrtmldegDir);

569     variableMP.printInfoMsg("Prepare climatological vegetation fraction");
571     variableBinariesEvs[binaryIndex++] = BINARY.vegfrac(llspl000, seamask, vegfrac,
        vegfracmnth);

573     variableMP.printInfoMsg("Prepare z0 and initial ustar");
575     variableBinariesEvs[binaryIndex++] = BINARY.z0vegfrac(seamask, landuse,
        topsoiltype, height, stdh, vegfrac, z0base, z0, ustar);

577     variableMP.printInfoMsg("Interpolate to model grid and execute allprep (fcst)");
579     variableBinariesEvs[binaryIndex++] = BINARY.allprep(llspl000, llgsst05, sst05,
        height, seamask, stdh, deeptemperature, snowalbedo, z0, z0base, landuse,
581     landusenew, topsoiltype, vegfrac, albedorrtm, llgsst, llgsno, llgcic,
        llsmst, llstmp, albedorrtmcorr, dzsoil, tskin, sst, snow, snowheight,
583     cice, seamaskcorr, landusecorr, landusenewcorr, topsoiltypecorr,
        vegfraccorr, z0corr, z0basecorr, emissivity, canopywater, frozenprecratio,
585     smst, sh2o, stmp, dsg, fcst, albedo, ustar, fcstDir, bocoPrefix,
        llsplPrefix);

587     variableMP.printInfoMsg("Prepare the dust related variable (soildust)");
589     variableBinariesEvs[binaryIndex++] = BINARY.readpaulsource(seamask, source,
        sourceNETCDF, sourceNCIncludeDir);

591     variableMP.printInfoMsg("Dust Start");
593     variableBinariesEvs[binaryIndex++] = BINARY.dust_start(llspl000, soildust, snow,
        topsoiltypecorr, landusecorr, landusenewcorr, kount_landuse,
        kount_landusenew, vegfrac, height, seamask, source, z0corr, roughness);
595

597     // Wait for binaries completion and check exit value
599     for (int i = 0; i < variableBinariesEvs.length; ++i) {
        LOGGER_VARIABLE.debug("Execution of " + i + " binary ended with status "
601         + variableBinariesEvs[i]);
        if (variableBinariesEvs[i] != 0) {
            throw new TaskExecutionException("[ERROR] Error executing binary " + i);
603         }
        }

605     variableMP.printHeaderMsg("END");

607     // Clean Up binaries
609     if (nmmbParams.isCleanBinaries()) {
        variableMP.printInfoMsg("Clean up executables");
        cleanUpVariableExe();
611     }

613     // Post execution
615     String folderOutputCase = NMMBEnvironment.OUTNMMB + nmmbParams.getCase()
        + File.separator;
        nmmbParams.postVariableExecution(folderOutputCase);

617     LOGGER_VARIABLE.info("Variable process finished");
619 }

621 private static void copyFilesFromPreprocess(NMMBParameters nmmbParams)
        throws TaskExecutionException {
623     // Clean specific files
        final String[] outputFiles = new String[] { "isop.dat", "meteo-data.dat",
625         "chemic-reg", "main_input_filename", "main_input_filename2", "GWD.bin",
        "configure_file", "co2_trans", "ETAMPNEW_AERO", "ETAMPNEW_DATA" };
627     for (String file : outputFiles) {
        String filePath = NMMBEnvironment.UMO_OUT + file;

```

```

629     if (!FileManagement.deleteFile(filePath)) {
630         LOGGER_UMO_MODEL.debug(
631             "Cannot erase previous " + file + " because it doesn't exist.");
632     }
633 }

634 // Clean regular expr files
635 File folder = new File(NMMBEnvironment.UMO_OUT);
636 for (File file : folder.listFiles()) {
637     if ((file.getName().startsWith("lai") && file.getName().endsWith(".dat"))
638         || (file.getName().startsWith("pftp_")
639             && file.getName().endsWith(".dat"))
640         || (file.getName().startsWith("PET")
641             && file.getName().endsWith(".txt"))
642         || (file.getName().startsWith("PET")
643             && file.getName().endsWith("File"))
644         || (file.getName().startsWith("boco."))
645         || (file.getName().startsWith("boco_chem."))
646         || (file.getName().startsWith("nmm_b_history."))
647         || (file.getName().startsWith("tr"))
648         || (file.getName().startsWith("RRT"))
649         || (file.getName().endsWith(".TBL"))
650         || (file.getName().startsWith("fcstdone."))
651         || (file.getName().startsWith("restartdone."))
652         || (file.getName().startsWith("nmm_b_rst_"))
653         || (file.getName().startsWith("nmm_b_hst_"))) {
654
655         if (!FileManagement.deleteFile(file)) {
656             LOGGER_UMO_MODEL.debug("Cannot erase previous " + file.getName()
657                 + " because it doesn't exist.");
658         }
659     }
660 }
661 }
662 }
663
664 // Prepare UMO model files
665 String outputFolderPath = NMMBEnvironment.OUTPUT;
666 File outputFolder = new File(outputFolderPath);
667 for (File file : outputFolder.listFiles()) {
668     if (file.getName().startsWith("boco.")
669         || file.getName().startsWith("boco_chem.")) {
670         // Copy file
671         String targetPath = NMMBEnvironment.UMO_OUT + file.getName();
672         if (!FileManagement.copyFile(file.getAbsolutePath(), targetPath)) {
673             throw new TaskExecutionException("[ERROR] Error copying file from "
674                 + file.getName() + " to " + targetPath);
675         }
676     }
677 }
678
679 String chemicRegSrc = NMMBEnvironment.OUTPUT + "chemic-reg";
680 String chemicRegTarget = NMMBEnvironment.UMO_OUT + "chemic-reg";
681 if (!FileManagement.copyFile(chemicRegSrc, chemicRegTarget)) {
682     // TODO: Really no error when file does not exist?
683     LOGGER_UMO_MODEL.debug("Cannot copy file from " + chemicRegSrc + " to "
684         + chemicRegTarget + ". Skipping...");
685 }
686
687 String gwdSrc = NMMBEnvironment.OUTPUT + "GWD.bin";
688 String gwdTarget = NMMBEnvironment.UMO_OUT + "GWD.bin";
689 if (!FileManagement.copyFile(gwdSrc, gwdTarget)) {
690     // TODO: Really no error when file does not exist?
691     LOGGER_UMO_MODEL.debug("Cannot copy file from " + gwdSrc + " to " + gwdTarget
692         + ". Skipping...");
693 }
694
695 String inputDomain1Src = NMMBEnvironment.OUTPUT + "fcst";
696 String inputDomain1Target = NMMBEnvironment.UMO_OUT + "input_domain_01";
697 if (!FileManagement.copyFile(inputDomain1Src, inputDomain1Target)) {
698     throw new TaskExecutionException("[ERROR] Error copying file from "
699         + inputDomain1Src + " to " + inputDomain1Target);
700 }

```

```
701 String inputDomain2Src = NMMBEnvironment.OUTPUT + "soildust";
702 String inputDomain2Target = NMMBEnvironment.UMO_OUT + "main_input_filename2";
703 if (!FileManagement.copyFile(inputDomain2Src, inputDomain2Target)) {
704     throw new TaskExecutionException("[ERROR] Error copying file from "
705         + inputDomain2Src + " to " + inputDomain2Target);
706 }
707
708 // Copy aerosols scavenging coeff
709 String lookupAerosol2RH00Src = NMMBEnvironment.OUTPUT
710     + "lookup_aerosol2.dat.rh00";
711 String lookupAerosol2RH00Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH00";
712 if (!FileManagement.copyFile(lookupAerosol2RH00Src, lookupAerosol2RH00Target)) {
713     throw new TaskExecutionException("[ERROR] Error copying file from "
714         + lookupAerosol2RH00Src + " to " + lookupAerosol2RH00Target);
715 }
716
717 String lookupAerosol2RH50Src = NMMBEnvironment.OUTPUT
718     + "lookup_aerosol2.dat.rh50";
719 String lookupAerosol2RH50Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH50";
720 if (!FileManagement.copyFile(lookupAerosol2RH50Src, lookupAerosol2RH50Target)) {
721     throw new TaskExecutionException("[ERROR] Error copying file from "
722         + lookupAerosol2RH50Src + " to " + lookupAerosol2RH50Target);
723 }
724
725 String lookupAerosol2RH70Src = NMMBEnvironment.OUTPUT
726     + "lookup_aerosol2.dat.rh70";
727 String lookupAerosol2RH70Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH70";
728 if (!FileManagement.copyFile(lookupAerosol2RH70Src, lookupAerosol2RH70Target)) {
729     throw new TaskExecutionException("[ERROR] Error copying file from "
730         + lookupAerosol2RH70Src + " to " + lookupAerosol2RH70Target);
731 }
732
733 String lookupAerosol2RH80Src = NMMBEnvironment.OUTPUT
734     + "lookup_aerosol2.dat.rh80";
735 String lookupAerosol2RH80Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH80";
736 if (!FileManagement.copyFile(lookupAerosol2RH80Src, lookupAerosol2RH80Target)) {
737     throw new TaskExecutionException("[ERROR] Error copying file from "
738         + lookupAerosol2RH80Src + " to " + lookupAerosol2RH80Target);
739 }
740
741 String lookupAerosol2RH90Src = NMMBEnvironment.OUTPUT
742     + "lookup_aerosol2.dat.rh90";
743 String lookupAerosol2RH90Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH90";
744 if (!FileManagement.copyFile(lookupAerosol2RH90Src, lookupAerosol2RH90Target)) {
745     throw new TaskExecutionException("[ERROR] Error copying file from "
746         + lookupAerosol2RH90Src + " to " + lookupAerosol2RH90Target);
747 }
748
749 String lookupAerosol2RH95Src = NMMBEnvironment.OUTPUT
750     + "lookup_aerosol2.dat.rh95";
751 String lookupAerosol2RH95Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH95";
752 if (!FileManagement.copyFile(lookupAerosol2RH95Src, lookupAerosol2RH95Target)) {
753     throw new TaskExecutionException("[ERROR] Error copying file from "
754         + lookupAerosol2RH95Src + " to " + lookupAerosol2RH95Target);
755 }
756
757 String lookupAerosol2RH99Src = NMMBEnvironment.OUTPUT
758     + "lookup_aerosol2.dat.rh99";
759 String lookupAerosol2RH99Target = NMMBEnvironment.UMO_OUT + "ETAMPNEW_AERO_RH99";
760 if (!FileManagement.copyFile(lookupAerosol2RH99Src, lookupAerosol2RH99Target)) {
761     throw new TaskExecutionException("[ERROR] Error copying file from "
762         + lookupAerosol2RH99Src + " to " + lookupAerosol2RH99Target);
763 }
764
765 // Copy coupling previous day (if required)
766 if (nmmbParams.getCoupleDustInit()) {
767     String historySrc = NMMBEnvironment.OUTNMMB + nmmbParams.getCase()
768         + File.separator + "history_INIT.hhh";
769     String historyTarget = NMMBEnvironment.UMO_OUT + "history_INIT.hhh";
770     if (!FileManagement.copyFile(historySrc, historyTarget)) {
771         throw new TaskExecutionException("[ERROR] Error copying file from "
```

```

773         + historySrc + " to " + historyTarget);
774     }
775 }
776 }
777
778 /**
779  * Performs the UMO Model simulation step
780  *
781  */
782 private static void doUMOModel(NMMBParameters nmmbParams, Date currentDate)
783     throws TaskExecutionException {
784     LOGGER_UMO_MODEL.info("Enter UMO Model process");
785
786     // Prepare execution
787     copyFilesFromPreprocess(nmmbParams);
788     nmmbParams.prepareUMOModelExecution(currentDate);
789     MessagePrinter umoModelMP = new MessagePrinter(LOGGER_UMO_MODEL);
790
791     // Begin MPI call
792     umoModelMP.printHeaderMsg("BEGIN");
793     umoModelMP.printInfoMsg("Executing nmmb_esmf.x UMO-NMMb-DUST-RRTM model");
794
795     String stdoutFile = NMMBEnvironment.UMO_OUT + "nems.out";
796     String stderrFile = NMMBEnvironment.UMO_OUT + "nems.err";
797     Integer nemsEV = MPI.nems(stdoutFile, stderrFile);
798
799     LOGGER_UMO_MODEL.debug("Execution of mpirun NEMS ended with status " + nemsEV);
800     if (nemsEV != 0) {
801         throw new TaskExecutionException("[ERROR] Error executing mpirun nems");
802     }
803     umoModelMP
804         .printInfoMsg("Finished Executing nmmb_esmf.x UMO-NMMb-DUST-RRTM model");
805
806     // Post execution
807     nmmbParams.postUMOModelExecution(currentDate);
808
809     umoModelMP.printHeaderMsg("END");
810
811     LOGGER_UMO_MODEL.info("UMO Model process finished");
812 }
813
814 /**
815  * Performs the POST step
816  *
817  */
818 private static void doPost(NMMBParameters nmmbParams, Date currentDate)
819     throws TaskExecutionException {
820     // Define model output folder by case and date
821     String currentDateSTR = NMMBConstants.STR_TO_DATE.format(currentDate);
822     String hourSTR = (nmmbParams.getHour() < 10)
823         ? "0" + String.valueOf(nmmbParams.getHour())
824         : String.valueOf(nmmbParams.getHour());
825     String folderOutput = NMMBEnvironment.OUTNMMB + nmmbParams.getCase()
826         + File.separator + currentDateSTR + hourSTR + File.separator;
827
828     LOGGER_POST.info("Postproc_carbono process for DAY: " + currentDateSTR);
829
830     // Prepare execution
831     nmmbParams.preparePostProcessExecution(currentDate);
832     MessagePrinter postProcMP = new MessagePrinter(LOGGER_POST);
833
834     // Deploy files and compile binaries if needed
835     Integer evCompile = BINARY.preparePost(nmmbParams.isCompileBinaries(),
836         folderOutput);
837     if (evCompile != 0) {
838         throw new TaskExecutionException("[ERROR] Error preparing post process");
839     }
840
841     // Begin POST call
842     postProcMP.printHeaderMsg("BEGIN");
843
844     String domainSTR = (nmmbParams.getDomain()) ? "glob" : "reg";

```

```

845 String dateHour = currentDateSTR + hourSTR;
846 String pout = folderOutput + "new_pout_*.nc";
847 String CTM = folderOutput + "NMMB-BSC-CTM_" + dateHour + "_" + domainSTR + ".nc";

849 Integer ev = BINARY.executePostprocAuth(folderOutput, pout, CTM);
850 LOGGER_POST.debug("Execution of NCRAT ended with status " + ev);
851 if (ev != 0) {
852     throw new TaskExecutionException("[ERROR] Error executing post process");
853 }

855 // Post execution
856 nmmbParams.cleanPostProcessExecution(folderOutput);

857 postProcMP.printHeaderMsg("END");

859 LOGGER_POST.info("Post process finished");
861 }

863 /**
864  * Main NMMB Workflow
865  *
866  * @param args
867  *         args[0] : Configuration file path
868  *
869  * @throws MainExecutionException
870  */
871 public static void main(String[] args) throws MainExecutionException {
872     LOGGER_MAIN.info("Starting NMMB application");

873     // Check and get arguments
874     if (args.length != 1) {
875         usage();
876         System.exit(1);
877     }
878     String configurationFile = args[0];

881     // Load configuration
882     NMMBConfigManager nmmbConfigManager = null;
883     try {
884         LOGGER_MAIN.info("Loading NMMB Configuration file " + configurationFile);
885         nmmbConfigManager = new NMMBConfigManager(configurationFile);
886         LOGGER_MAIN.info("Configuration file loaded");
887     } catch (ConfigurationException ce) {
888         LOGGER_MAIN.error("[ERROR] Cannot load configuration file: "
889             + configurationFile + ". Aborting...", ce);
890         throw new MainExecutionException(ce);
891     }

893     // Compute the execution variables
894     NMMBParameters nmmbParams = new NMMBParameters(nmmbConfigManager);

895     // Prepare the execution
896     nmmbParams.prepareExecution();

897     // Fixed process (do before main time looping)
898     if (nmmbParams.doFixed()) {
899         try {
900             doFixed(nmmbParams);
901         } catch (TaskExecutionException tee) {
902             LOGGER_FIXED.error("[ERROR] Task exception on fixed phase. Aborting...",
903                 tee);
904             throw new MainExecutionException(tee);
905         }
906     }

909     // Start main time loop
910     Date currentDate = nmmbParams.getStartDate();
911     while (!currentDate.after(nmmbParams.getEndDate())) {
912         String currentDateSTR = NMMBConstants.STR_TO_DATE.format(currentDate);
913         LOGGER_MAIN.info(currentDateSTR + " simulation started");

915         // Create output folders if needed

```

```

917     nmmbParams.createOutputFolders(currentDate);
919
920     // Vrbl process
921     if (nmmbParams.doVariable()) {
922         try {
923             doVariable(nmmbParams, currentDate);
924         } catch (TaskExecutionException tee) {
925             LOGGER_VARIABLE
926                 .error("[ERROR] Task exception on variable phase at date "
927                     + currentDateSTR + ". Aborting...", tee);
928             throw new MainExecutionException(tee);
929         }
930     }
931
932     // UMO model run
933     if (nmmbParams.doUmoModel()) {
934         try {
935             doUMOModel(nmmbParams, currentDate);
936         } catch (TaskExecutionException tee) {
937             LOGGER_UMO_MODEL
938                 .error("[ERROR] Task exception on UMO Model phase at date "
939                     + currentDateSTR + ". Aborting...", tee);
940             throw new MainExecutionException(tee);
941         }
942     }
943
944     // Post process
945     if (nmmbParams.doPost()) {
946         try {
947             doPost(nmmbParams, currentDate);
948         } catch (TaskExecutionException tee) {
949             LOGGER_POST.error("[ERROR] Task exception on Post phase at date "
950                 + currentDateSTR + ". Aborting...", tee);
951             throw new MainExecutionException(tee);
952         }
953     }
954
955     LOGGER_MAIN.info(currentDateSTR + " simulation finished");
956
957     // Getting next simulation day
958     currentDate = Date.from(currentDate.toInstant()
959         .plusSeconds(NMMBConstants.ONE_DAY_IN_SECONDS));
960 }
961 }
962 }
963 }

```

## C.2 NmmbItf.java

```

1 package nmmb;
2
3 import integratedtoolkit.types.annotations.Constraints;
4 import integratedtoolkit.types.annotations.Parameter;
5 import integratedtoolkit.types.annotations.parameter.Direction;
6 import integratedtoolkit.types.annotations.parameter.Stream;
7 import integratedtoolkit.types.annotations.parameter.Type;
8 import integratedtoolkit.types.annotations.task.Binary;
9 import integratedtoolkit.types.annotations.task.MPI;
10 import nmmb.configuration.NMMBConstants;
11 import nmmb.utils.BinaryWrapper;
12 import nmmb.utils.FortranWrapper;
13
14
15 public interface NmmbItf {
16
17     // COMPILE STEP
18     @Binary(binary = FortranWrapper.FC)

```



```

19 Integer fortranCompiler(
    @Parameter(type = Type.STRING, direction = Direction.IN) String mcFlag,
21    @Parameter(type = Type.STRING, direction = Direction.IN) String sharedFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String covertPrefix,
23    @Parameter(type = Type.STRING, direction = Direction.IN) String convertValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String tracebackFlag,
25    @Parameter(type = Type.STRING, direction = Direction.IN) String assumePrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String assumeValue,
27    @Parameter(type = Type.STRING, direction = Direction.IN) String optFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelPrefix,
29    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String stackFlag,
31    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String executable,
33    @Parameter(type = Type.FILE, direction = Direction.IN) String source
);

35
@Binary(binary = FortranWrapper.FC)
Integer fortranCompileObject(
    @Parameter(type = Type.STRING, direction = Direction.IN) String mcFlag,
39    @Parameter(type = Type.STRING, direction = Direction.IN) String sharedFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String covertPrefix,
41    @Parameter(type = Type.STRING, direction = Direction.IN) String convertValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String tracebackFlag,
43    @Parameter(type = Type.STRING, direction = Direction.IN) String assumePrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String assumeValue,
45    @Parameter(type = Type.STRING, direction = Direction.IN) String optFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelPrefix,
47    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String stackFlag,
49    @Parameter(type = Type.STRING, direction = Direction.IN) String cFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String source,
51    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String object,
53    @Parameter(type = Type.STRING, direction = Direction.IN) String moduleFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String moduleDir
);

55
@Binary(binary = FortranWrapper.FC)
Integer fortranCompileWithObject(
59    @Parameter(type = Type.STRING, direction = Direction.IN) String mcFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String sharedFlag,
61    @Parameter(type = Type.STRING, direction = Direction.IN) String covertPrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String convertValue,
63    @Parameter(type = Type.STRING, direction = Direction.IN) String tracebackFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String assumePrefix,
65    @Parameter(type = Type.STRING, direction = Direction.IN) String assumeValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String optFlag,
67    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelPrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelValue,
69    @Parameter(type = Type.STRING, direction = Direction.IN) String stackFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
71    @Parameter(type = Type.STRING, direction = Direction.IN) String executable,
    @Parameter(type = Type.STRING, direction = Direction.IN) String source,
73    @Parameter(type = Type.FILE, direction = Direction.IN) String object
);

75
@Binary(binary = FortranWrapper.FC)
Integer fortranCompilerWithW3(
77    @Parameter(type = Type.STRING, direction = Direction.IN) String mcFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String sharedFlag,
79    @Parameter(type = Type.STRING, direction = Direction.IN) String covertPrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String convertValue,
81    @Parameter(type = Type.STRING, direction = Direction.IN) String tracebackFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String assumePrefix,
83    @Parameter(type = Type.STRING, direction = Direction.IN) String assumeValue,
    @Parameter(type = Type.STRING, direction = Direction.IN) String optFlag,
85    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelPrefix,
    @Parameter(type = Type.STRING, direction = Direction.IN) String fpmodelValue,
87    @Parameter(type = Type.STRING, direction = Direction.IN) String stackFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
89    @Parameter(type = Type.STRING, direction = Direction.IN) String executable,

```

```

91     @Parameter(type = Type.FILE, direction = Direction.IN) String source,
     @Parameter(type = Type.STRING, direction = Direction.IN) String w3libDir,
93     @Parameter(type = Type.STRING, direction = Direction.IN) String bacioLibDir,
     @Parameter(type = Type.STRING, direction = Direction.IN) String w3Lib,
95     @Parameter(type = Type.STRING, direction = Direction.IN) String bacioLib
     );
97
     @Binary(binary = FortranWrapper.GFC)
99     Integer gfortranCompiler(
     @Parameter(type = Type.STRING, direction = Direction.IN) String bigOFlag,
101     @Parameter(type = Type.FILE, direction = Direction.IN) String source,
     @Parameter(type = Type.STRING, direction = Direction.IN) String oFlag,
103     @Parameter(type = Type.STRING, direction = Direction.IN) String executable
     );
105
     @Binary(binary = BinaryWrapper.COMPILE_READ_PAUL_SOURCE)
107     Integer compileReadPaulSource(
     @Parameter(type = Type.FILE, direction = Direction.IN) String source,
109     @Parameter(type = Type.STRING, direction = Direction.IN) String executable
     );
111
     // FIXED STEP
113     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.SMMOUNT + FortranWrapper.SUFFIX_EXE)
115     Integer smmount(
     @Parameter(type = Type.STRING, direction = Direction.IN) String topoDir,
117     @Parameter(type = Type.FILE, direction = Direction.OUT) String seamaskDEM,
     @Parameter(type = Type.FILE, direction = Direction.OUT) String heightDEM
119     );
121
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.LANDUSE + FortranWrapper.SUFFIX_EXE)
123     Integer landuse(
     @Parameter(type = Type.STRING, direction = Direction.IN) String landuseDir,
125     @Parameter(type = Type.FILE, direction = Direction.OUT) String landuse,
     @Parameter(type = Type.FILE, direction = Direction.OUT) String kount_landuse
127     );
129
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.LANDUSENEW + FortranWrapper.SUFFIX_EXE)
131     Integer landusenew(
     @Parameter(type = Type.STRING, direction = Direction.IN) String gtopDir,
133     @Parameter(type = Type.FILE, direction = Direction.OUT) String landusenew,
     @Parameter(type = Type.FILE, direction = Direction.OUT) String kount_landusenew
135     );
137
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.TOPO + FortranWrapper.SUFFIX_EXE)
139     Integer topo(
     @Parameter(type = Type.STRING, direction = Direction.IN) String topoDir,
141     @Parameter(type = Type.FILE, direction = Direction.OUT) String heightmean
     );
143
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.STDH + FortranWrapper.SUFFIX_EXE)
145     Integer stdh(
     @Parameter(type = Type.FILE, direction = Direction.IN) String seamaskDEM,
     @Parameter(type = Type.FILE, direction = Direction.IN) String heightmean,
149     @Parameter(type = Type.STRING, direction = Direction.IN) String topoDir,
     @Parameter(type = Type.FILE, direction = Direction.OUT) String stdh
151     );
153
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.ENVELOPE + FortranWrapper.SUFFIX_EXE)
155     Integer envelope(
     @Parameter(type = Type.FILE, direction = Direction.IN) String heightmean,
157     @Parameter(type = Type.FILE, direction = Direction.IN) String stdh,
     @Parameter(type = Type.FILE, direction = Direction.OUT) String height
159     );
161
     @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
     + FortranWrapper.TOPSOILTYPE + FortranWrapper.SUFFIX_EXE)

```

```

163 Integer topsoiltype(
    @Parameter(type = Type.FILE, direction = Direction.IN) String seamaskDEM,
165     @Parameter(type = Type.STRING, direction = Direction.IN) String soiltypeDir,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String topsoiltype
167 );

169 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.BOTSOILTYPE + FortranWrapper.SUFFIX_EXE)
171 Integer botsoiltype(
    @Parameter(type = Type.FILE, direction = Direction.IN) String seamaskDEM,
173     @Parameter(type = Type.STRING, direction = Direction.IN) String soiltypePath,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String botsoiltype
175 );

177 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.TOPOSEAMASK + FortranWrapper.SUFFIX_EXE)
179 Integer toposeamask(
    @Parameter(type = Type.FILE, direction = Direction.IN) String seamaskDEM,
181     @Parameter(type = Type.FILE, direction = Direction.OUT) String seamask,
    @Parameter(type = Type.FILE, direction = Direction.INOUT) String height,
183     @Parameter(type = Type.FILE, direction = Direction.INOUT) String landuse,
    @Parameter(type = Type.FILE, direction = Direction.INOUT) String topsoiltype,
185     @Parameter(type = Type.FILE, direction = Direction.INOUT) String botsoiltype
    );

187
189 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.STDHTOPO + FortranWrapper.SUFFIX_EXE)
Integer stdhtopo(
191     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
    @Parameter(type = Type.FILE, direction = Direction.INOUT) String stdh
193 );

195 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.DEEPTEMPERATURE + FortranWrapper.SUFFIX_EXE)
197 Integer deeptemperature(
    @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
199     @Parameter(type = Type.STRING, direction = Direction.IN) String soiltempPath,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String deeptemperature
201 );

203 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.SNOWALBEDO + FortranWrapper.SUFFIX_EXE)
205 Integer snowalbedo(
    @Parameter(type = Type.STRING, direction = Direction.IN) String maxsnowalbDir,
207     @Parameter(type = Type.FILE, direction = Direction.OUT) String snowalbedo
    );

209
211 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.VCGENERATOR + FortranWrapper.SUFFIX_EXE)
Integer vcgenerator(
213     @Parameter(type = Type.FILE, direction = Direction.OUT) String dsg
    );

215
217 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.ROUGHNESS + FortranWrapper.SUFFIX_EXE)
Integer roughness(
219     @Parameter(type = Type.STRING, direction = Direction.IN) String roughnessDir,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String roughness
221 );

223 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/"
    + FortranWrapper.GFDLCO2 + FortranWrapper.SUFFIX_EXE)
225 Integer gfdlco2(
    @Parameter(type = Type.FILE, direction = Direction.IN) String dsg,
227     @Parameter(type = Type.STRING, direction = Direction.IN) String co2_data_dir,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String co2_trans
229 );

231 @Binary(binary = "${" + NMMBConstants.ENV_NAME_FIX + "}/lookup_tables/"
    + BinaryWrapper.RUN_AEROSOL)
233 Integer run_aerosol(
    @Parameter() boolean mustCompile,

```

```

235     @Parameter() boolean mustClean,
236     @Parameter(type = Type.FILE, direction = Direction.OUT)
237     String lookup_aerosol2_rh00,
238     @Parameter(type = Type.FILE, direction = Direction.OUT)
239     String lookup_aerosol2_rh50,
240     @Parameter(type = Type.FILE, direction = Direction.OUT)
241     String lookup_aerosol2_rh70,
242     @Parameter(type = Type.FILE, direction = Direction.OUT)
243     String lookup_aerosol2_rh80,
244     @Parameter(type = Type.FILE, direction = Direction.OUT)
245     String lookup_aerosol2_rh90,
246     @Parameter(type = Type.FILE, direction = Direction.OUT)
247     String lookup_aerosol2_rh95,
248     @Parameter(type = Type.FILE, direction = Direction.OUT)
249     String lookup_aerosol2_rh99
250 );
251
252
253 // VARIABLE STEP
254 @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
255         + BinaryWrapper.DEGRIB_GFS_GENERIC)
256 Integer degribgfs_generic_05(
257     @Parameter(type = Type.FILE, direction = Direction.OUT) String CW,
258     @Parameter(type = Type.FILE, direction = Direction.OUT) String ICEC,
259     @Parameter(type = Type.FILE, direction = Direction.OUT) String SH,
260     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILT2,
261     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILT4,
262     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILW2,
263     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILW4,
264     @Parameter(type = Type.FILE, direction = Direction.OUT) String TT,
265     @Parameter(type = Type.FILE, direction = Direction.OUT) String VV,
266     @Parameter(type = Type.FILE, direction = Direction.OUT) String HH,
267     @Parameter(type = Type.FILE, direction = Direction.OUT) String PRMSL,
268     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILT1,
269     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILT3,
270     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILW1,
271     @Parameter(type = Type.FILE, direction = Direction.OUT) String SOILW3,
272     @Parameter(type = Type.FILE, direction = Direction.OUT) String SST_TS,
273     @Parameter(type = Type.FILE, direction = Direction.OUT) String UU,
274     @Parameter(type = Type.FILE, direction = Direction.OUT) String WEASD
275 );
276
277 @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
278         + FortranWrapper.GFS2MODEL + FortranWrapper.SUFFIX_EXE)
279 Integer gfs2model_rrtm(
280     @Parameter(type = Type.FILE, direction = Direction.IN) String CW,
281     @Parameter(type = Type.FILE, direction = Direction.IN) String ICEC,
282     @Parameter(type = Type.FILE, direction = Direction.IN) String SH,
283     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILT2,
284     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILT4,
285     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILW2,
286     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILW4,
287     @Parameter(type = Type.FILE, direction = Direction.IN) String TT,
288     @Parameter(type = Type.FILE, direction = Direction.IN) String VV,
289     @Parameter(type = Type.FILE, direction = Direction.IN) String HH,
290     @Parameter(type = Type.FILE, direction = Direction.IN) String PRMSL,
291     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILT1,
292     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILT3,
293     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILW1,
294     @Parameter(type = Type.FILE, direction = Direction.IN) String SOILW3,
295     @Parameter(type = Type.FILE, direction = Direction.IN) String SST_TS,
296     @Parameter(type = Type.FILE, direction = Direction.IN) String UU,
297     @Parameter(type = Type.FILE, direction = Direction.IN) String WEASD,
298     @Parameter(type = Type.FILE, direction = Direction.OUT) String GFS_file
299 );
300
301 @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
302         + FortranWrapper.INC_RRTM + FortranWrapper.SUFFIX_EXE)
303 Integer inc_rrtm(
304     @Parameter(type = Type.FILE, direction = Direction.IN) String GFS_file,
305     @Parameter(type = Type.STRING, direction = Direction.IN) String deco
306 );

```

```

307     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
309           + FortranWrapper.CNV_RRTM + FortranWrapper.SUFFIX_EXE)
Integer cnv_rrtm(
311     @Parameter(type = Type.FILE, direction = Direction.IN) String GFS_file,
313     @Parameter(type = Type.FILE, direction = Direction.OUT) String llspl000,
315     @Parameter(type = Type.FILE, direction = Direction.OUT) String outtmp,
317     @Parameter(type = Type.FILE, direction = Direction.OUT) String outmst,
319     @Parameter(type = Type.FILE, direction = Direction.OUT) String outsst,
321     @Parameter(type = Type.FILE, direction = Direction.OUT) String outsno,
323     @Parameter(type = Type.FILE, direction = Direction.OUT) String outcic
325 );

327     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
329           + FortranWrapper.DEGRIB_SST + FortranWrapper.SUFFIX_EXE)
Integer degribsst(
331     @Parameter(type = Type.FILE, direction = Direction.OUT) String llgsst05,
333     @Parameter(type = Type.STRING, direction = Direction.IN) String sstfileinPath
335 );

337     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
339           + FortranWrapper.ALBEDO + FortranWrapper.SUFFIX_EXE)
Integer albedo(
341     @Parameter(type = Type.FILE, direction = Direction.IN) String llspl000,
343     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
345     @Parameter(type = Type.FILE, direction = Direction.OUT) String albedo,
347     @Parameter(type = Type.FILE, direction = Direction.OUT) String albedobase,
349     @Parameter(type = Type.STRING, direction = Direction.IN) String albedomnth
351 );

353     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
355           + FortranWrapper.ALBEDO_RRTM_1DEG + FortranWrapper.SUFFIX_EXE)
Integer albedorrtm(
357     @Parameter(type = Type.FILE, direction = Direction.IN) String llspl000,
359     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
361     @Parameter(type = Type.FILE, direction = Direction.OUT) String albedorrtm,
363     @Parameter(type = Type.STRING, direction = Direction.IN) String albedorrtm1degDir
365 );

367     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
369           + FortranWrapper.VEG_FRAC + FortranWrapper.SUFFIX_EXE)
Integer vegfrac(
371     @Parameter(type = Type.FILE, direction = Direction.IN) String llspl000,
373     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
375     @Parameter(type = Type.FILE, direction = Direction.OUT) String vegfrac,
377     @Parameter(type = Type.STRING, direction = Direction.IN) String vegfracmnth
379 );

381     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
383           + FortranWrapper.Z0_VEGUSTAR + FortranWrapper.SUFFIX_EXE)
Integer z0vegfrac(
385     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
387     @Parameter(type = Type.FILE, direction = Direction.IN) String landuse,
389     @Parameter(type = Type.FILE, direction = Direction.IN) String topsoiltype,
391     @Parameter(type = Type.FILE, direction = Direction.IN) String height,
393     @Parameter(type = Type.FILE, direction = Direction.IN) String stdh,
395     @Parameter(type = Type.FILE, direction = Direction.IN) String vegfrac,
397     @Parameter(type = Type.FILE, direction = Direction.OUT) String z0base,
399     @Parameter(type = Type.FILE, direction = Direction.OUT) String z0,
401     @Parameter(type = Type.FILE, direction = Direction.OUT) String ular
403 );

405     @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
407           + FortranWrapper.ALLPREP_RRTM + FortranWrapper.SUFFIX_EXE)
Integer allprep(
409     @Parameter(type = Type.FILE, direction = Direction.IN) String llspl000,
411     @Parameter(type = Type.FILE, direction = Direction.IN) String llgsst05,
413     @Parameter(type = Type.FILE, direction = Direction.OUT) String sst05,
415     @Parameter(type = Type.FILE, direction = Direction.IN) String height,
417     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
419     @Parameter(type = Type.FILE, direction = Direction.IN) String stdh,
421     @Parameter(type = Type.FILE, direction = Direction.INOUT) String deeptemperature,

```

```

379     @Parameter(type = Type.FILE, direction = Direction.INOUT) String snowalbedo,
380     @Parameter(type = Type.FILE, direction = Direction.IN) String z0,
381     @Parameter(type = Type.FILE, direction = Direction.IN) String z0base,
382     @Parameter(type = Type.FILE, direction = Direction.IN) String landuse,
383     @Parameter(type = Type.FILE, direction = Direction.IN) String landusenew,
384     @Parameter(type = Type.FILE, direction = Direction.IN) String topsoiltype,
385     @Parameter(type = Type.FILE, direction = Direction.IN) String vegfrac,
386     @Parameter(type = Type.FILE, direction = Direction.IN) String albedorrtm,
387     @Parameter(type = Type.FILE, direction = Direction.IN) String llgsst,
388     @Parameter(type = Type.FILE, direction = Direction.IN) String llgsno,
389     @Parameter(type = Type.FILE, direction = Direction.IN) String llgcic,
390     @Parameter(type = Type.FILE, direction = Direction.IN) String llmst,
391     @Parameter(type = Type.FILE, direction = Direction.IN) String llstmp,
392     @Parameter(type = Type.FILE, direction = Direction.OUT) String albedorrtmcorr,
393     @Parameter(type = Type.FILE, direction = Direction.OUT) String dzsoil,
394     @Parameter(type = Type.FILE, direction = Direction.OUT) String tskin,
395     @Parameter(type = Type.FILE, direction = Direction.OUT) String sst,
396     @Parameter(type = Type.FILE, direction = Direction.OUT) String snow,
397     @Parameter(type = Type.FILE, direction = Direction.OUT) String snowheight,
398     @Parameter(type = Type.FILE, direction = Direction.OUT) String cice,
399     @Parameter(type = Type.FILE, direction = Direction.OUT) String seamaskcorr,
400     @Parameter(type = Type.FILE, direction = Direction.OUT) String landusecorr,
401     @Parameter(type = Type.FILE, direction = Direction.OUT) String landusenewcorr,
402     @Parameter(type = Type.FILE, direction = Direction.OUT) String topsoiltypecorr,
403     @Parameter(type = Type.FILE, direction = Direction.OUT) String vegfraccorr,
404     @Parameter(type = Type.FILE, direction = Direction.OUT) String z0corr,
405     @Parameter(type = Type.FILE, direction = Direction.OUT) String z0basecorr,
406     @Parameter(type = Type.FILE, direction = Direction.OUT) String emissivity,
407     @Parameter(type = Type.FILE, direction = Direction.OUT) String canopywater,
408     @Parameter(type = Type.FILE, direction = Direction.OUT) String frozenprecratio,
409     @Parameter(type = Type.FILE, direction = Direction.OUT) String smst,
410     @Parameter(type = Type.FILE, direction = Direction.OUT) String sh2o,
411     @Parameter(type = Type.FILE, direction = Direction.OUT) String stmp,
412     @Parameter(type = Type.FILE, direction = Direction.IN) String dsq,
413     @Parameter(type = Type.FILE, direction = Direction.OUT) String fcst,
414     @Parameter(type = Type.FILE, direction = Direction.IN) String albedo,
415     @Parameter(type = Type.FILE, direction = Direction.IN) String ustar,
416     @Parameter(type = Type.STRING, direction = Direction.IN) String fcstDir,
417     @Parameter(type = Type.STRING, direction = Direction.IN) String bocoPrefix,
418     @Parameter(type = Type.STRING, direction = Direction.IN) String llsplPrefix
419 );

420
421
422 @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
423         + FortranWrapper.READ_PAUL_SOURCE + FortranWrapper.SUFFIX_EXE)
424 Integer readpaulsource(
425     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
426     @Parameter(type = Type.FILE, direction = Direction.OUT) String source,
427     @Parameter(type = Type.FILE, direction = Direction.OUT) String sourceNETCDF,
428     @Parameter(type = Type.STRING, direction = Direction.IN) String sourceNCIncludeDir
429 );

430
431 @Binary(binary = "${" + NMMBConstants.ENV_NAME_VRB + "}/"
432         + FortranWrapper.DUST_START + FortranWrapper.SUFFIX_EXE)
433 Integer dust_start(
434     @Parameter(type = Type.FILE, direction = Direction.IN) String llspl000,
435     @Parameter(type = Type.FILE, direction = Direction.OUT) String soildust,
436     @Parameter(type = Type.FILE, direction = Direction.IN) String snow,
437     @Parameter(type = Type.FILE, direction = Direction.IN) String topsoiltypecorr,
438     @Parameter(type = Type.FILE, direction = Direction.IN) String landusecorr,
439     @Parameter(type = Type.FILE, direction = Direction.IN) String landusenewcorr,
440     @Parameter(type = Type.FILE, direction = Direction.IN) String kount_landuse,
441     @Parameter(type = Type.FILE, direction = Direction.IN) String kount_landusenew,
442     @Parameter(type = Type.FILE, direction = Direction.IN) String vegfrac,
443     @Parameter(type = Type.FILE, direction = Direction.IN) String height,
444     @Parameter(type = Type.FILE, direction = Direction.IN) String seamask,
445     @Parameter(type = Type.FILE, direction = Direction.IN) String source,
446     @Parameter(type = Type.FILE, direction = Direction.IN) String z0corr,
447     @Parameter(type = Type.FILE, direction = Direction.IN) String roughness
448 );

449 // UMO MODEL STEP

```

```
451 @MPI(mpiRunner = "mpirun",
      binary = "${" + NMMBConstants.ENV_NAME_UMO_PATH + "}/MODEL/exe/"
453       + FortranWrapper.NEMS + FortranWrapper.SUFFIX_EXE,
      workingDir = "${" + NMMBConstants.ENV_NAME_UMO_OUT + "}/",
455       computingNodes = "${NEMS_NODES}")
      @Constraints(computingUnits = "${NEMS_CUS_PER_NODE}")
457       Integer nems (
          @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
459         String stdoutFile,
          @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
461         String stderrFile
      );
463
      // POST PROC STEP
465 @Binary(binary = "${" + NMMBConstants.ENV_NAME_POST_CARBONO + "}/"
          + BinaryWrapper.PREPARE_POSTPROC_AUTH)
467 Integer preparePost (
          @Parameter() boolean mustCompile,
469         @Parameter(type = Type.STRING, direction = Direction.IN) String folderOutput
      );
471
      @Binary(binary = "${" + NMMBConstants.ENV_NAME_POST_CARBONO + "}/"
473             + BinaryWrapper.EXEC_POSTPROC_AUTH)
      Integer executePostprocAuth (
475         @Parameter(type = Type.STRING, direction = Direction.IN) String folderOutput,
          @Parameter(type = Type.STRING, direction = Direction.IN) String sourcesPath,
477         @Parameter(type = Type.FILE, direction = Direction.OUT) String destFile
      );
479
  }
```





## Appendix D

# GUIDANCE: code highlights

### D.1 Guidance.java

```

package guidance;
2
import java.io.BufferedWriter;
4 import java.io.File;
import java.io.FileWriter;
6 import java.io.IOException;
import java.util.ArrayList;
8 import java.util.List;
import java.util.Map;
10
import binary.BINARY;
12
import java.text.SimpleDateFormat;
14 import java.util.Date;
16
import guidance.GuidanceImpl;
import guidance.exceptions.EnvironmentVariableException;
18 import guidance.exceptions.GuidanceTaskException;
import guidance.files.AssocFiles;
20 import guidance.files.CombinedPanelsFiles;
import guidance.files.CommonFiles;
22 import guidance.files.ImputationFiles;
import guidance.files.MergeFiles;
24 import guidance.files.PhenomeAnalysisFiles;
import guidance.files.ResultsFiles;
26 import guidance.utils.ChromoInfo;
import guidance.utils.Environment;
28 import guidance.utils.Environment.ENV_VARS;
import guidance.utils.ParseCmdLine;
30
32 /**
 * The higher-level class of Guidance
34 */
public class Guidance {
36
    private static final String PACKAGE_VERSION = "guidance_0.9.6";
38
    private static final long B_TO_GB = 1_048_576;
40
    private static final double PVA_THRESHOLD = 5e-8;
42 private static final String PVA_THRS = Double.toString(PVA_THRESHOLD);
44
    /**
46  * The main of Guidance begins here
    *
48  * @param args
    * @throws EnvironmentVariableException
50  * @throws Exception
    */
52 public static void main(String args[]) throws EnvironmentVariableException,
    IOException {

```

```

54 // Verify that all environment variables have been defined correctly
55 Environment.verify();
56
57 // Print information of Guidance version
58 printGuidancePackageVersion();
59
60 // Get the input arguments
61 ParseCmdLine parsingArgs = new ParseCmdLine(args);
62
63 // Verify and print the status of each stage
64 printStagesStatus(parsingArgs);
65 System.out.println("\n[Guidance] Verified stages status.");
66
67 // Get the file name where the list of commands is going to be saved
68 String listOfStagesFileName = parsingArgs.getListOfStagesFile();
69
70 // Verify whether the file exists or not.
71 File listOfStages = new File(listOfStagesFileName);
72 if (!listOfStages.exists()) {
73     System.out.println("\n[Guidance] File to store the tasks list: "
74         + listOfStagesFileName);
75 } else {
76     System.out.println("\n[Guidance] File to store the tasks list (overwritten): "
77         + listOfStagesFileName);
78 }
79 listOfStages.createNewFile();
80
81 ArrayList<String> listOfCommands = new ArrayList<>();
82 String datestring = new SimpleDateFormat("yyyy/MM/dd HH:mm:ss").format(new Date());
83 listOfCommands.add("#####");
84 listOfCommands.add("# List of tasks executed by Guidance workflow");
85 listOfCommands.add("# Date: " + datestring);
86 listOfCommands.add("# Parameters of the execution: ");
87 listOfCommands.add("#####");
88
89 // Get the names of the reference panel to be used in the execution
90 ArrayList<String> rpanelTypes = new ArrayList<>(parsingArgs.getRpanelTypes());
91 String outDir = parsingArgs.getOutDir();
92 ChromoInfo generalChromoInfo = new ChromoInfo();
93
94 System.out.println("[Guidance] We print testTypes information");
95 int numberOfTestTypes = parsingArgs.getNumberOfTestTypeName();
96 for (int kk = 0; kk < numberOfTestTypes; kk++) {
97     String tmp_test_type = parsingArgs.getTestTypeName(kk);
98     String tmp_responseVar = parsingArgs.getResponseVar(kk);
99     String tmp_covariables = parsingArgs.getCovariables(kk);
100     System.out.println("[Guidance] " + tmp_test_type + " = "
101         + tmp_responseVar + ";" + tmp_covariables);
102 }
103
104 // Main code of the workflow
105 // Now, we have to ask whether the GWAS process is mixed.
106 // printEnvVariables();
107 doMixed(parsingArgs, outDir, rpanelTypes, generalChromoInfo, listOfCommands);
108
109 // Finally, we print the commands in the output file defined for this.
110 try (BufferedWriter writer = new BufferedWriter(new FileWriter(listOfStages))) {
111     for (String str : listOfCommands) {
112         writer.write(str);
113         writer.newLine();
114         writer.newLine();
115     }
116
117     writer.flush();
118 } catch (IOException ioe) {
119     throw ioe;
120 }
121 }
122
123 /**
124 * Method that performs the complete workflow when "mixed" type of GWAS is chosen

```

```

126  *
127  * @param parsingArgs
128  * @param outDir
129  * @param rpanelTypes
130  * @param generalChromoInfo
131  * @param listOfCommands
132  */
133 private static void doMixed(ParseCmdLine parsingArgs, String outDir,
134     List<String> rpanelTypes, ChromoInfo generalChromoInfo,
135     ArrayList<String> listOfCommands) {
136
137     // Create some general objects
138     int startChr = parsingArgs.getStart();
139     int endChr = parsingArgs.getEnd();
140
141     String exclCgatFlag = parsingArgs.getExclCgatSnp();
142     String exclSVFlag = parsingArgs.getExclSVSnp();
143     int chunkSize = parsingArgs.getChunkSize();
144     String inputFormat = null;
145
146     // Create the names for the common files
147     CommonFiles commonFilesInfo = new CommonFiles(parsingArgs, outDir);
148
149     // Create the whole directory structure.
150     createDirStructure(parsingArgs, outDir, rpanelTypes, startChr, endChr);
151
152     // Create the names for mixed files
153     ImputationFiles imputationFilesInfo = new ImputationFiles(parsingArgs,
154         generalChromoInfo, outDir, rpanelTypes);
155
156     // Create the names for Association files
157     AssocFiles assocFilesInfo = new AssocFiles(parsingArgs, generalChromoInfo,
158         outDir, rpanelTypes);
159
160     // Create the names for Merge files
161     MergeFiles mergeFilesInfo = new MergeFiles(parsingArgs, generalChromoInfo,
162         outDir, rpanelTypes);
163
164     // Create the names for Results Files. Take into account this class it to
165     // generate file name for results
166     // of each testType and each rpanelType. The results for combined panels or
167     // phenome analysis are created
168     // by other class.
169     ResultsFiles resultsFilesInfo = new ResultsFiles(parsingArgs, outDir, rpanelTypes);
170     CombinedPanelsFiles combinedPanelsFilesInfo = new CombinedPanelsFiles(parsingArgs,
171         outDir, rpanelTypes);
172     PhenomeAnalysisFiles phenomeAnalysisFilesInfo = new PhenomeAnalysisFiles(
173         parsingArgs, outDir, rpanelTypes);
174
175     // Create all file names used in the workflow.
176     // Also, define which of them will be temporal or permanent.
177     setFinalStatusForCommonFiles(parsingArgs, commonFilesInfo);
178     setFinalStatusForImputationFiles(parsingArgs, imputationFilesInfo,
179         generalChromoInfo, rpanelTypes);
180     setFinalStatusForAssocFiles(parsingArgs, assocFilesInfo, generalChromoInfo,
181         rpanelTypes);
182
183     // boolean refPanelCombine = parsingArgs.getRefPanelCombine();
184     // int numOfChromosToProcess = endChr - startChr + 1;
185     for (int i = startChr; i <= endChr; i++) {
186         // We get the output pairs file name for mixed
187         String mixedPairsFile = commonFilesInfo.getPairsFile(i);
188         String mixedSampleFile = commonFilesInfo.getSampleFile(i);
189
190         String theChromo = Integer.toString(i);
191         String mixedGenFile = null;
192         String bedFile = null;
193         String bimFile = null;
194         String famFile = null;
195         String mixedBedFile = null;
196         String mixedBimFile = null;
197         String mixedFamFile = null;

```

```

198 String mixedBedToBedLogFile = null;
200
201 String gmapFile = parsingArgs.getGmapDir() + File.separator
202     + parsingArgs.getGmapFileName(i);
203 String mixedShapeitHapsFile = commonFilesInfo.getShapeitHapsFile(i);
204 String mixedShapeitSampleFile = commonFilesInfo.getShapeitSampleFile(i);
205 String mixedShapeitLogFile = commonFilesInfo.getShapeitLogFile(i);
206 String mixedExcludedSnpsFile = commonFilesInfo.getExcludedSnpsFile(i);
207
208 String mixedFilteredHaplotypesFile
209     = commonFilesInfo.getFilteredHaplotypesFile(i);
210 String mixedFilteredHaplotypesSampleFile
211     = commonFilesInfo.getFilteredHaplotypesSampleFile(i);
212 String mixedFilteredHaplotypesLogFile
213     = commonFilesInfo.getFilteredHaplotypesLogFile(i);
214 String mixedFilteredHaplotypesVcfFile
215     = commonFilesInfo.getFilteredHaplotypesVcfFile(i);
216 String mixedFilteredListOfSnpsFile
217     = commonFilesInfo.getListOfSnpsFile(i);
218
219 // Check if the input are in GEN o BED format
220 inputFormat = parsingArgs.getInputFormat();
221 if (inputFormat.equals("BED")) {
222     bedFile = commonFilesInfo.getBedFile();
223     bimFile = commonFilesInfo.getBimFile();
224     famFile = commonFilesInfo.getFamFile();
225
226     mixedBedFile = commonFilesInfo.getByChrBedFile(i);
227     mixedBimFile = commonFilesInfo.getByChrBimFile(i);
228     mixedFamFile = commonFilesInfo.getByChrFamFile(i);
229     mixedBedToBedLogFile = commonFilesInfo.getBedToBedLogFile(i);
230
231     doConvertFromBedToBed(parsingArgs, listOfCommands, bedFile, bimFile,
232         famFile, theChromo, mixedBedFile, mixedBimFile,
233         mixedFamFile, mixedBedToBedLogFile);
234
235     // Create the RsId list of SNPs that are AT, TA, CG, or GC
236     // In that case, because inputType is BED we pass the newBimFile
237     doCreateRsIdList(parsingArgs, listOfCommands, mixedBimFile, exclCgatFlag,
238         mixedPairsFile, inputFormat);
239
240     doPhasingBed(parsingArgs, listOfCommands, theChromo, mixedBedFile,
241         mixedBimFile, mixedFamFile, gmapFile,
242         mixedExcludedSnpsFile, mixedShapeitHapsFile,
243         mixedShapeitSampleFile, mixedShapeitLogFile,
244         mixedFilteredHaplotypesFile, mixedFilteredHaplotypesSampleFile,
245         mixedFilteredHaplotypesLogFile, mixedFilteredHaplotypesVcfFile,
246         mixedFilteredListOfSnpsFile, exclCgatFlag, exclSVflag);
247 } else { // The inputFormat is GEN
248     mixedGenFile = commonFilesInfo.getGenFile(i);
249     // Task
250     doCreateRsIdList(parsingArgs, listOfCommands, mixedGenFile, exclCgatFlag,
251         mixedPairsFile, inputFormat);
252
253     doPhasing(parsingArgs, listOfCommands, theChromo, mixedGenFile,
254         mixedSampleFile, gmapFile, mixedExcludedSnpsFile,
255         mixedShapeitHapsFile, mixedShapeitSampleFile,
256         mixedShapeitLogFile, mixedFilteredHaplotypesFile,
257         mixedFilteredHaplotypesSampleFile, mixedFilteredHaplotypesLogFile,
258         mixedFilteredHaplotypesVcfFile, mixedFilteredListOfSnpsFile,
259         exclCgatFlag, exclSVflag);
260
261 } // End of inputFormat.equals("GEN")
262
263 // *****
264 // * Now perform the imputation tasks
265 // * over the data for the different reference
266 // * panels that the user wants to include. The
267 // * list of the panels to be used are stored in
268 // * the List array rpanelTypes
269 // *****
270 for (int kk = 0; kk < rpanelTypes.size(); kk++) {

```

```

270         // String rpanelDir = parsingArgs.getRpanelDir(kk);
271         int maxSize = generalChromoInfo.getMaxSize(i);
272         int minSize = generalChromoInfo.getMinSize(i);
273
274         int lim1 = minSize;
275         int lim2 = lim1 + chunkSize - 1;
276         int j = 0;
277         for (j = minSize; j < maxSize; j += chunkSize) {
278             makeImputationPerChunk(parsingArgs, i, lim1, lim2, kk, gmapFile,
279                 imputationFilesInfo, commonFilesInfo, listOfCommands);
280             lim1 = lim1 + chunkSize;
281             lim2 = lim2 + chunkSize;
282         }
283     } // End for(kk=0; kk<rpanelTypes.size(); kk++)
284 } // End for(i=startChr;i <= endChr; i++)
285
286 // *****
287 // * COMPSs API Call to wait for all tasks *
288 // * comment out following line to include the synchronization point *
289 // *****
290 // COMPSs.barrier();
291
292 // Now we continue with the association
293 int numberOfTestTypes = parsingArgs.getNumberOfTestTypeName();
294 for (int tt = 0; tt < numberOfTestTypes; tt++) {
295     for (int kk = 0; kk < rpanelTypes.size(); kk++) {
296         for (int i = startChr; i <= endChr; i++) {
297             int maxSize = generalChromoInfo.getMaxSize(i);
298             int minSize = generalChromoInfo.getMinSize(i);
299
300             int lim1 = minSize;
301             int lim2 = lim1 + chunkSize - 1;
302             for (int j = minSize; j < maxSize; j += chunkSize) {
303                 makeAssociationPerChunk(parsingArgs, tt, kk, i, lim1, lim2,
304                     imputationFilesInfo, commonFilesInfo, assocFilesInfo,
305                     listOfCommands);
306
307                 lim1 = lim1 + chunkSize;
308                 lim2 = lim2 + chunkSize;
309             }
310
311             // Now we perform the merge of chunks for each chromosome
312             makeMergeOfChunks(parsingArgs, listOfCommands, tt, kk, i, minSize,
313                 maxSize, chunkSize, assocFilesInfo, mergeFilesInfo);
314
315             // We have to filterByAll the last merged File of each chromomose
316             String theLastReducedFile
317                 = mergeFilesInfo.getTheLastReducedFile(tt, kk, i);
318             String filteredByAllFile
319                 = mergeFilesInfo.getFilteredByAllFile(tt, kk, i);
320             String condensedFile
321                 = mergeFilesInfo.getCondensedFile(tt, kk, i);
322
323             doFilterByAll(parsingArgs, listOfCommands, theLastReducedFile,
324                 filteredByAllFile, condensedFile);
325         } // End for(i=startChr;i <= endChr; i++)
326
327         // Now we have to joint the condensedFiles of each chromosome. There is
328         // not problem if chr23 is being processed, because
329         // the format of condensedFiles is the same for all chromosome.
330         makeJointCondensedFiles(parsingArgs, listOfCommands, tt, kk, startChr,
331             endChr, mergeFilesInfo);
332
333         // Now we have to joint the filteredByAllFiles of each chromosome.
334         // Here there is an additional complexity due to the chr23,
335         // because the file format changes if chr23 is being processed. This
336         // situation is taken into account inside the next function.
337         String rpanelName = rpanelTypes.get(kk);
338
339         makeJointFilteredByAllFiles(parsingArgs, listOfCommands, tt, kk,
340             rpanelName, startChr, endChr, mergeFilesInfo);

```

```

342     String lastCondensedFile = mergeFilesInfo.getFinalCondensedFile(tt, kk);
344
345     String lastFilteredByAllFile
346         = mergeFilesInfo.getFinalFilteredByAllFile(tt, kk);
347
348     // Lets extract top-hits results from the lastFilteredByAllFile
349     String topHitsResults = resultsFilesInfo.getTopHitsFile(tt, kk);
350
351     String filteredByAllXFile = null;
352     if (endChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
353         filteredByAllXFile
354             = mergeFilesInfo.getAdditionalFilteredByAllXFile(tt, kk, 0);
355     } else {
356         filteredByAllXFile = lastFilteredByAllFile;
357     }
358
359     doGenerateTopHits(parsingArgs, listOfCommands, lastFilteredByAllFile,
360                     filteredByAllXFile, topHitsResults, PVA_THRS);
361
362     String qqPlotPdfFile = resultsFilesInfo.getQqPlotPdfFile(tt, kk);
363     String qqPlotTiffFile = resultsFilesInfo.getQqPlotTiffFile(tt, kk);
364
365     String manhattanPlotPdfFile
366         = resultsFilesInfo.getManhattanPdfFile(tt, kk);
367     String manhattanPlotTiffFile
368         = resultsFilesInfo.getManhattanTiffFile(tt, kk);
369
370     String correctedPvaluesFile
371         = resultsFilesInfo.getCorrectedPvaluesFile(tt, kk);
372
373     doGenerateQQManhattanPlots(parsingArgs, listOfCommands, lastCondensedFile,
374                               qqPlotPdfFile, manhattanPlotPdfFile,
375                               qqPlotTiffFile, manhattanPlotTiffFile, correctedPvaluesFile);
376
377     } // End for(kk=0; kk<rpanelTypes.size(); kk++)
378
379     // Now we continue with the combining of the results of the different
380     // reference panels.
381     // It is done if the refPanelCombine flag is true.
382     makeCombinePanels(parsingArgs, mergeFilesInfo, combinedPanelsFilesInfo,
383                     rpanelTypes, tt, listOfCommands);
384
385     } // End for(tt=0; tt<numberOfTestTypes;tt++)
386
387     makePhenotypeAnalysis(parsingArgs, mergeFilesInfo, resultsFilesInfo,
388                          phenomeAnalysisFilesInfo, rpanelTypes, listOfCommands);
389
390     System.out.println("\n[Guidance] All tasks are in execution, please wait...");
391 }
392
393 /**
394  * Method that generates all the tasks for imputation and association
395  *
396  * @param parsingArgs
397  * @param chrNumber
398  * @param lim1
399  * @param lim2
400  * @param panelIndex
401  * @param gmapFile
402  * @param imputationFilesInfo
403  * @param commonFilesInfo
404  * @param listOfCommands
405  */
406 private static void makeImputationPerChunk(ParseCmdLine parsingArgs, int chrNumber,
407     int lim1, int lim2, int panelIndex, String gmapFile,
408     ImputationFiles imputationFilesInfo, CommonFiles commonFilesInfo,
409     ArrayList<String> listOfCommands) {
410
411     int chunkSize = parsingArgs.getChunkSize();
412     double infoThreshold = parsingArgs.getInfoThreshold();
413
414     String knownHapFileName = parsingArgs.getRpanelHapFileName(panelIndex, chrNumber);

```

```

414 String rpanelDir = parsingArgs.getRpanelDir(panelIndex);
    String knownHapFile = rpanelDir + File.separator + knownHapFileName;
416
    String lim1S = Integer.toString(lim1);
418 String lim2S = Integer.toString(lim2);
    String chrS = Integer.toString(chrNumber);
420 String infoThresholdS = Double.toString(infoThreshold);
    String imputationTool = parsingArgs.getImputationTool();
422
    if (imputationTool.equals("impute")) {
424         String legendFileName=parsingArgs.getRpanelLegFileName(panelIndex, chrNumber);
            String legendFile = rpanelDir + File.separator + legendFileName;
426
            String mixedShapeitHapsFile = commonFilesInfo.getShapeitHapsFile(chrNumber);
428 String mixedShapeitSampleFile=commonFilesInfo.getShapeitSampleFile(chrNumber);
            String mixedPairsFile = commonFilesInfo.getPairsFile(chrNumber);
430 String mixedImputeFile = imputationFilesInfo.getImputedFile(panelIndex,
                chrNumber, lim1, lim2, chunkSize);
432 String mixedImputeFileInfo =imputationFilesInfo.getImputedInfoFile(panelIndex,
                chrNumber, lim1, lim2, chunkSize);
434 String mixedImputeFileSummary = imputationFilesInfo.getImputedSummaryFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
436 String mixedImputeFileWarnings = imputationFilesInfo.getImputedWarningsFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
438
            String mixedFilteredFile = imputationFilesInfo.getFilteredFile(panelIndex,
440 chrNumber, lim1, lim2, chunkSize);
            String mixedFilteredLogFile = imputationFilesInfo.getFilteredLogFile(
442 panelIndex, chrNumber, lim1, lim2, chunkSize);
444
            // We create the list of rsId that are greater than or equal to the
            // infoThreshold value
446 String mixedFilteredRsIdFile = imputationFilesInfo.getFilteredRsIdFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
448
            doImputationWithImpute(parsingArgs, listOfCommands, chrS, gmapFile,
450 knownHapFile, legendFile, mixedShapeitHapsFile,
                mixedShapeitSampleFile, lim1S, lim2S, mixedPairsFile,
452 mixedImputeFile, mixedImputeFileInfo, mixedImputeFileSummary,
                mixedImputeFileWarnings);
454
            doFilterByInfo(parsingArgs, listOfCommands, mixedImputeFileInfo,
456 mixedFilteredRsIdFile, infoThresholdS);
458
            doQctools(parsingArgs, listOfCommands, mixedImputeFile,mixedFilteredRsIdFile,
                mixedFilteredLogFile);
460 } else if (imputationTool.equals("minimac")) {
            String mixedFilteredHaplotypesFile
462 = commonFilesInfo.getFilteredHaplotypesFile(chrNumber);
            String mixedFilteredHaplotypesSampleFile
464 = commonFilesInfo.getFilteredHaplotypesSampleFile(chrNumber);
            String mixedFilteredListOfSnpsFile
466 = commonFilesInfo.getListOfSnpsFile(chrNumber);
468
            String mixedImputedMMFileName = imputationFilesInfo.getImputedMMFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
470 String mixedImputedMMInfoFile = imputationFilesInfo.getImputedMMInfoFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
472 String mixedImputedMMERateFile = imputationFilesInfo.getImputedMMERateFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
474 String mixedImputedMMRecFile = imputationFilesInfo.getImputedMMRecFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
476 String mixedImputedMMDoseFile = imputationFilesInfo.getImputedMMDoseFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
478 String mixedImputedMMLogFile = imputationFilesInfo.getImputedMMLogFile(
                panelIndex, chrNumber, lim1, lim2, chunkSize);
480
            doImputationWithMinimac(parsingArgs, listOfCommands, knownHapFile,
482 mixedFilteredHaplotypesFile, mixedFilteredHaplotypesSampleFile,
                mixedFilteredListOfSnpsFile, mixedImputedMMFileName,
484 mixedImputedMMInfoFile, mixedImputedMMERateFile,
                mixedImputedMMRecFile, mixedImputedMMDoseFile, mixedImputedMMLogFile,

```

```

486         chrS, lim1S, lim2S);
    } else {
488         System.err.println("\t[makeImputationPerChunk]: Error, the imputation tool "
+ imputationTool + " is not allowed...");
490         System.exit(1);
    }
492 }

494 /**
495  * Method that generates all the tasks for association
496  *
497  * @param parsingArgs
498  * @param testTypeIndex
499  * @param panelIndex
500  * @param chrNumber
501  * @param lim1
502  * @param lim2
503  * @param imputationFilesInfo
504  * @param commonFilesInfo
505  * @param assocFilesInfo
506  * @param listOfCommands
507  */
508 private static void makeAssociationPerChunk(ParseCmdLine parsingArgs,
    int testTypeIndex, int panelIndex, int chrNumber, int lim1,
510     int lim2, ImputationFiles imputationFilesInfo, CommonFiles commonFilesInfo,
    AssocFiles assocFilesInfo, ArrayList<String> listOfCommands) {
512
    int chunkSize = parsingArgs.getChunkSize();
514     double mafThreshold = parsingArgs.getMafThreshold();
    double infoThreshold = parsingArgs.getInfoThreshold();
516     double hweCohortThreshold = parsingArgs.getHweCohortThreshold();
    double hweCasesThreshold = parsingArgs.getHweCasesThreshold();
518     double hweControlsThreshold = parsingArgs.getHweControlsThreshold();

520     String responseVar = parsingArgs.getResponseVar(testTypeIndex);
    String covariables = parsingArgs.getCovariables(testTypeIndex);
522
    String chrS = Integer.toString(chrNumber);
524
    String mafThresholdS = Double.toString(mafThreshold);
    String infoThresholdS = Double.toString(infoThreshold);
526     String hweCohortThresholdS = Double.toString(hweCohortThreshold);
    String hweCasesThresholdS = Double.toString(hweCasesThreshold);
528     String hweControlsThresholdS = Double.toString(hweControlsThreshold);
530
    String snptestOutFile = assocFilesInfo.getSnptestOutFile(testTypeIndex,
532     panelIndex, chrNumber, lim1, lim2, chunkSize);
    String snptestLogFile = assocFilesInfo.getSnptestLogFile(testTypeIndex,
534     panelIndex, chrNumber, lim1, lim2, chunkSize);
536
    String mixedSampleFile = commonFilesInfo.getSampleFile(chrNumber);
    String mixedImputeFileInfo = imputationFilesInfo.getImputedInfoFile(panelIndex,
538     chrNumber, lim1, lim2, chunkSize);
540
    String mixedFilteredFile = imputationFilesInfo.getFilteredFile(panelIndex,
    chrNumber, lim1, lim2, chunkSize);
542
    doSnptest(parsingArgs, listOfCommands, chrS, mixedFilteredFile, mixedSampleFile,
544     snptestOutFile, snptestLogFile, responseVar, covariables);
546
    String summaryFile = assocFilesInfo.getSummaryFile(testTypeIndex, panelIndex,
    chrNumber, lim1, lim2, chunkSize);
548
    doCollectSummary(parsingArgs, listOfCommands, chrS, mixedImputeFileInfo,
550     snptestOutFile, summaryFile, mafThresholdS, infoThresholdS,
    hweCohortThresholdS, hweCasesThresholdS, hweControlsThresholdS);
552 }

554 /**
555  * Method to perform the merging fo chunks for each chromosome
556  *
557  * @param parsingArgs

```



```

558 * @param listOfCommands
559 * @param ttIndex
560 * @param rpanelIndex
561 * @param chr
562 * @param minSize
563 * @param maxSize
564 * @param chunkSize
565 * @param assocFilesInfo
566 * @param mergeFilesInfo
567 */
568 private static void makeMergeOfChunks(ParseCmdLine parsingArgs,
569     ArrayList<String> listOfCommands, int ttIndex, int rpanelIndex, int chr,
570     int minSize, int maxSize, int chunkSize, AssocFiles assocFilesInfo,
571     MergeFiles mergeFilesInfo) {
572
573     String theChromo = Integer.toString(chr);
574     int lim1 = minSize;
575     int lim2 = lim1 + chunkSize - 1;
576
577     int numberOfChunks = maxSize / chunkSize;
578     int module = maxSize % chunkSize;
579     if (module != 0) {
580         numberOfChunks++;
581     }
582     int indexA = 0;
583     int indexC = 0;
584     String reducedA = null;
585     String reducedB = null;
586     String reducedC = null;
587     for (int processedChunks = 0; processedChunks < 2 * numberOfChunks - 2;
588         processedChunks = processedChunks + 2) {
589         if (processedChunks < numberOfChunks) {
590             reducedA = assocFilesInfo.getSummaryFile(ttIndex, rpanelIndex, chr,
591                 lim1, lim2, chunkSize);
592             lim1 = lim1 + chunkSize;
593             lim2 = lim2 + chunkSize;
594         } else {
595             reducedA=mergeFilesInfo.getReducedFile(ttIndex, rpanelIndex, chr, indexA);
596             indexA++;
597         }
598
599         if (processedChunks < numberOfChunks - 1) {
600             reducedB = assocFilesInfo.getSummaryFile(ttIndex, rpanelIndex, chr,
601                 lim1, lim2, chunkSize);
602             lim1 = lim1 + chunkSize;
603             lim2 = lim2 + chunkSize;
604         } else {
605             reducedB=mergeFilesInfo.getReducedFile(ttIndex, rpanelIndex, chr, indexA);
606             indexA++;
607         }
608         reducedC = mergeFilesInfo.getReducedFile(ttIndex, rpanelIndex, chr, indexC);
609         doMergeTwoChunks(parsingArgs, listOfCommands, reducedA, reducedB, reducedC,
610             theChromo);
611         indexC++;
612     } // End of for of chunks
613 }
614
615 /**
616 * Method to perform the joint of condensed files of each rpanel
617 *
618 * @param parsingArgs
619 * @param listOfCommands
620 * @param ttIndex
621 * @param rpanelIndex
622 * @param startChr
623 * @param endChr
624 * @param mergeFilesInfo
625 */
626 private static void makeJointCondensedFiles(ParseCmdLine parsingArgs,
627     ArrayList<String> listOfCommands, int ttIndex, int rpanelIndex,
628     int startChr, int endChr, MergeFiles mergeFilesInfo) {

```

```

630     int indexA = 0;
631     int indexC = 0;
632     String condensedA = null;
633     String condensedB = null;
634     String condensedC = null;
635     int numberOfChrs = endChr - startChr + 1;
636
637     if (numberOfChrs == 1) { // There is only one chr to process.
638         condensedA = mergeFilesInfo.getCondensedFile(ttIndex, rpanelIndex, startChr);
639         condensedB = condensedA;
640
641         condensedC
642             = mergeFilesInfo.getAdditionalCondensedFile(ttIndex, rpanelIndex, indexC);
643         doJointCondenseFiles(parsingArgs, listOfCommands, condensedA, condensedB,
644             condensedC);
645     } else {
646         for (int processedCondensed = 0; processedCondensed < 2 * numberOfChrs - 2;
647             processedCondensed = processedCondensed + 2) {
648             if (processedCondensed < numberOfChrs) {
649                 int i = startChr + processedCondensed;
650                 condensedA = mergeFilesInfo.getCondensedFile(ttIndex, rpanelIndex, i);
651             } else {
652                 condensedA = mergeFilesInfo.getAdditionalCondensedFile(ttIndex,
653                     rpanelIndex, indexA);
654                 indexA++;
655             }
656
657             if (processedCondensed < numberOfChrs - 1) {
658                 int i = startChr + processedCondensed + 1;
659                 condensedB = mergeFilesInfo.getCondensedFile(ttIndex, rpanelIndex, i);
660             } else {
661                 condensedB = mergeFilesInfo.getAdditionalCondensedFile(ttIndex,
662                     rpanelIndex, indexA);
663                 indexA++;
664             }
665             condensedC = mergeFilesInfo.getAdditionalCondensedFile(ttIndex,
666                 rpanelIndex, indexC);
667
668             doJointCondenseFiles(parsingArgs, listOfCommands, condensedA,
669                 condensedB, condensedC);
670             indexC++;
671         } // End for
672     }
673 }
674
675 /**
676  * Method to perform the joint of filteredByAll files of each rpanel
677  *
678  * @param parsingArgs
679  * @param listOfCommands
680  * @param ttIndex
681  * @param rpanelIndex
682  * @param rpanelName
683  * @param startChr
684  * @param endChr
685  * @param mergeFilesInfo
686  */
687 private static void makeJointFilteredByAllFiles(ParseCmdLine parsingArgs,
688     ArrayList<String> listOfCommands, int ttIndex, int rpanelIndex,
689     String rpanelName, int startChr, int endChr, MergeFiles mergeFilesInfo) {
690
691     int endChrNormal = endChr;
692     if (startChr < endChr) {
693         if (endChr != ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
694             endChrNormal = endChr;
695         } else {
696             endChrNormal = endChr - 1;
697         }
698     }
699
700     int indexA = 0;
701     int indexC = 0;

```

```

702     int processedFiltered = 0;
       String rpanelFlag = "NO";
704     String filteredA = null;
       String filteredB = null;
706     String filteredC = null;
       int numberOfChrs = endChrNormal - startChr + 1;
708
       if (numberOfChrs == 1) { // There is only one chr to process.
710         rpanelFlag = "YES";
           filteredA=mergeFilesInfo.getFilteredByAllFile(ttIndex, rpanelIndex, startChr);
712         filteredB = filteredA;
714
           filteredC = mergeFilesInfo.getAdditionalFilteredByAllFile(ttIndex,
               rpanelIndex, indexC);
716         if (startChr != ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
           doJointFilteredByAllFiles (parsingArgs, listOfCommands, filteredA,
718             filteredB, filteredC, rpanelName, rpanelFlag);
           }
720     } else {
       for (processedFiltered = 0; processedFiltered < 2 * numberOfChrs - 2;
722         processedFiltered = processedFiltered + 2) {
           if (processedFiltered < numberOfChrs) {
724             int i = startChr + processedFiltered;
               filteredA = mergeFilesInfo.getFilteredByAllFile(ttIndex,
726                 rpanelIndex, i);
           } else {
728             filteredA = mergeFilesInfo.getAdditionalFilteredByAllFile(ttIndex,
               rpanelIndex, indexA);
730             indexA++;
           }
732
           if (processedFiltered < numberOfChrs - 1) {
734             int i = startChr + processedFiltered + 1;
               filteredB=mergeFilesInfo.getFilteredByAllFile(ttIndex, rpanelIndex, i);
736             } else {
               filteredB = mergeFilesInfo.getAdditionalFilteredByAllFile(ttIndex,
738                 rpanelIndex, indexA);
               indexA++;
740             rpanelFlag = "YES";
           }
742         filteredC = mergeFilesInfo.getAdditionalFilteredByAllFile(ttIndex,
           rpanelIndex, indexC);
744
           doJointFilteredByAllFiles (parsingArgs, listOfCommands, filteredA,
746             filteredB, filteredC, rpanelName, rpanelFlag);
           indexC++;
748         } // End for
       }
750
       // Now we process the chr 23 if this is defined in this execution
752     if (endChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
       rpanelFlag = "YES";
754     filteredA = mergeFilesInfo.getFilteredByAllFile(ttIndex, rpanelIndex, endChr);
       filteredC = mergeFilesInfo.getAdditionalFilteredByAllXFile(ttIndex,
756     rpanelIndex, 0);
758
       doJointFilteredByAllFiles (parsingArgs, listOfCommands, filteredA,
           filteredA, filteredC, rpanelName, rpanelFlag);
760     }
762 }
764 /**
       * Method that performs the last part of the workflow corresponding to the merging,
766     * combining and drawing of results
       *
768     * @param parsingArgs
       * @param mergeFilesInfo
770     * @param combinedPanelsFilesInfo
       * @param rpanelTypes
772     * @param ttIndex
       * @param listOfCommands

```

```

774     */
775     private static void makeCombinePanels(ParseCmdLine parsingArgs,
776         MergeFiles mergeFilesInfo, CombinedPanelsFiles combinedPanelsFilesInfo,
777         List<String> rpanelTypes, int ttIndex, ArrayList<String> listOfCommands) {
778
779         int startChr = parsingArgs.getStart();
780         int endChr = parsingArgs.getEnd();
781         int endChrNormal = endChr;
782         if (endChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
783             endChrNormal = ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES - 1;
784         }
785
786         String startChrS = Integer.toString(startChr);
787         String endChrS = Integer.toString(endChr);
788         String endChrNormalS = Integer.toString(endChrNormal);
789
790         String filteredPanelA = null;
791         String filteredPanelB = null;
792         String filteredPanelC = null;
793
794         String filteredXPanelA = null;
795         String filteredXPanelB = null;
796         String filteredXPanelC = null;
797         boolean refPanelCombine = parsingArgs.getRefPanelCombine();
798
799         String combinedCondensedFile = null;
800
801         // We have to ask if we have to combine results from different panels or not.
802         if (refPanelCombine == true) {
803             int startingIndex = 0;
804             filteredPanelA = mergeFilesInfo.getFinalFilteredByAllFile(ttIndex,
805                 startingIndex);
806
807             if (endChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
808                 filteredXPanelA = mergeFilesInfo.getAdditionalFilteredByAllXFile(ttIndex,
809                     startingIndex, 0);
810             }
811
812             for (int kk = 1; kk < rpanelTypes.size(); kk++) {
813                 filteredPanelB = mergeFilesInfo.getFinalFilteredByAllFile(ttIndex, kk);
814                 filteredPanelC = combinedPanelsFilesInfo.getCombinedFilteredByAllFile(
815                     ttIndex, kk - 1);
816
817                 if (startChr != ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
818                     doCombinePanelsComplex(parsingArgs, listOfCommands, filteredPanelA,
819                         filteredPanelB, filteredPanelC, startChrS, endChrNormalS);
820                 }
821                 filteredPanelA = filteredPanelC;
822
823                 if (endChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
824                     filteredXPanelB = mergeFilesInfo.getAdditionalFilteredByAllXFile(
825                         ttIndex, kk, 0);
826                     filteredXPanelC
827                         = combinedPanelsFilesInfo.getCombinedFilteredByAllXFile(
828                             ttIndex, kk - 1);
829
830                     doCombinePanelsComplex(parsingArgs, listOfCommands, filteredXPanelA,
831                         filteredXPanelB, filteredXPanelC, endChrS, endChrS);
832                     filteredXPanelA = filteredXPanelC;
833                 }
834             }
835
836             // Now we have to generate the combinedCondensedFile. We do it, by using the
837             // last filteredPanelC and the last filteredXPanelC
838             if (endChr != ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES) {
839                 filteredXPanelC = filteredPanelC;
840             }
841
842             if ((startChr == ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES)) {
843                 filteredPanelC = filteredXPanelC;
844             }

```

```

846     combinedCondensedFile
      = combinedPanelsFilesInfo.getCombinedCondensedFile(ttIndex);
848
849     doCombineCondensedFiles(parsingArgs, listOfCommands, filteredPanelC,
850         filteredXPanelC, combinedCondensedFile);
851
852     // Finally, we create topHits, QQ and Manhattan plots.
853     String topHitsCombinedResults
854         = combinedPanelsFilesInfo.getTopHitsFile(ttIndex);
855
856     doGenerateTopHits(parsingArgs, listOfCommands, filteredPanelC,
857         filteredXPanelC, topHitsCombinedResults, PVA_THRS);
858
859     String combinedQqPlotPdfFile
860         = combinedPanelsFilesInfo.getQqPlotPdfFile(ttIndex);
861     String combinedQqPlotTiffFile
862         = combinedPanelsFilesInfo.getQqPlotTiffFile(ttIndex);
863     String combinedManhattanPlotPdfFile
864         = combinedPanelsFilesInfo.getManhattanPdfFile(ttIndex);
865     String combinedManhattanPlotTiffFile
866         = combinedPanelsFilesInfo.getManhattanTiffFile(ttIndex);
867     String combinedCorrectedPvaluesFile
868         = combinedPanelsFilesInfo.getCorrectedPvaluesFile(ttIndex);
869
870     doGenerateQQManhattanPlots(parsingArgs, listOfCommands, combinedCondensedFile,
871         combinedQqPlotPdfFile, combinedManhattanPlotPdfFile,
872         combinedQqPlotTiffFile, combinedManhattanPlotTiffFile,
873         combinedCorrectedPvaluesFile);
874
875 }
876
877 /**
878  * Method that performs the last part of the workflow corresponding to the phenome
879  * Analysis, combining the results of each phenotype analysis
880  *
881  * @param parsingArgs
882  * @param mergeFilesInfo
883  * @param resultsFilesInfo
884  * @param phenomeAnalysisFilesInfo
885  * @param rpanelTypes
886  * @param listOfCommands
887  */
888 private static void makePhenotypeAnalysis(ParseCmdLine parsingArgs,
889     MergeFiles mergeFilesInfo, ResultsFiles resultsFilesInfo,
890     PhenomeAnalysisFiles phenomeAnalysisFilesInfo, List<String> rpanelTypes,
891     ArrayList<String> listOfCommands) {
892
893     int endChr = parsingArgs.getEnd();
894     String endChrS = Integer.toString(endChr);
895
896     String topHitsFile = null;
897     String filteredByAllFile = null;
898     String filteredByAllXFile = null;
899     // String condensedFile = null;
900     String phenomeFileA = null;
901     String phenomeFileB = null;
902     String phenomeFileC = null;
903
904     int numberOfTestTypes = parsingArgs.getNumberOfTestTypeName();
905     int numberOfRpanelsTypes = rpanelTypes.size();
906
907     // int maxPhenoIndex = numberOfTestTypes * numberOfRpanelsTypes - 1;
908     int phenoIndex = 0;
909
910     int ttIndex = 0;
911     int rpIndex = 0;
912     String ttName = parsingArgs.getTestTypeName(ttIndex);
913     String rpName = rpanelTypes.get(rpIndex);
914
915     // Lets create the header for the initPhenoMatrix task
916     ttIndex = 0;

```

```

918     rpIndex = 0;
919     ttName = parsingArgs.getTestTypeName(ttIndex);
920     rpName = rpanelTypes.get(rpIndex);
921     topHitsFile = resultsFilesInfo.getTopHitsFile(ttIndex, rpIndex);
922     phenomeFileA = phenomeAnalysisFilesInfo.getPhenotypeIntermediateFile(phenoIndex);
923     phenoIndex++;
924
925     doInitPhenoMatrix(parsingArgs, listOfCommands, topHitsFile, ttName, rpName,
926         phenomeFileA);
927     for (ttIndex = 0; ttIndex < numberOfTestTypes; ttIndex++) {
928         int startRp = 0;
929         if (ttIndex == 0) {
930             startRp = 1;
931         }
932         ttName = parsingArgs.getTestTypeName(ttIndex);
933         for (rpIndex = startRp; rpIndex < numberOfRpanelsTypes; rpIndex++) {
934             rpName = rpanelTypes.get(rpIndex);
935
936             topHitsFile = resultsFilesInfo.getTopHitsFile(ttIndex, rpIndex);
937             phenomeFileB
938                 = phenomeAnalysisFilesInfo.getPhenotypeIntermediateFile(phenoIndex);
939
940             doAddToPhenoMatrix(parsingArgs, listOfCommands, phenomeFileA, topHitsFile,
941                 ttName, rpName, phenomeFileB);
942
943             phenomeFileA = phenomeFileB;
944             phenoIndex++;
945         }
946     }
947
948     phenoIndex = 0;
949     // Lets do the fillinout of the phenomeAnalysis.
950     for (ttIndex = 0; ttIndex < numberOfTestTypes; ttIndex++) {
951         ttName = parsingArgs.getTestTypeName(ttIndex);
952
953         for (rpIndex = 0; rpIndex < numberOfRpanelsTypes; rpIndex++) {
954             rpName = rpanelTypes.get(rpIndex);
955
956             filteredByAllFile = mergeFilesInfo.getFinalFilteredByAllFile(ttIndex,
957                 rpIndex);
958             if (endChrS.equals("23")) {
959                 filteredByAllXFile
960                     = mergeFilesInfo.getAdditionalFilteredByAllXFile(ttIndex, rpIndex, 0);
961             } else {
962                 filteredByAllXFile = mergeFilesInfo.getFinalFilteredByAllFile(ttIndex,
963                     rpIndex);
964             }
965
966             phenomeFileB = phenomeAnalysisFilesInfo.getPhenotypeFile(phenoIndex);
967
968             doFilloutPhenoMatrix(parsingArgs, listOfCommands, phenomeFileA,
969                 filteredByAllFile, filteredByAllXFile, endChrS, ttName,
970                 rpName, phenomeFileB);
971
972             phenoIndex++;
973         }
974     }
975
976     // Last round to generate final results
977     phenoIndex = 0;
978     phenomeFileA = phenomeAnalysisFilesInfo.getPhenotypeFile(phenoIndex);
979
980     phenoIndex++;
981
982     // Lets do the finalization of the phenomeAnalysis.
983     for (ttIndex = 0; ttIndex < numberOfTestTypes; ttIndex++) {
984         int startRp = 0;
985         if (ttIndex == 0) {
986             startRp = 1;
987         }
988
989         ttName = parsingArgs.getTestTypeName(ttIndex);

```

```

990         for (rpIndex = startRp; rpIndex < numberOfRpanelsTypes; rpIndex++) {
991             rpName = rpanelTypes.get(rpIndex);
992             phenomeFileB = phenomeAnalysisFilesInfo.getPhenotypeFile(phenoIndex);
993             phenomeFileC = phenomeAnalysisFilesInfo.getPhenotypeFinalFile(phenoIndex);
994             doFinalizePhenoMatrix(parsingArgs, listOfCommands, phenomeFileA,
995                 phenomeFileB, ttName, rpName, phenomeFileC);
996             phenomeFileA = phenomeFileC;
997             phenoIndex++;
998         }
999     }
1000 }
1001
1002 /**
1003  * Method the wraps the execution of convertFromBedToBed tasks and store the command
1004  * in the listOfCommands
1005  *
1006  * @param parsingArgs
1007  * @param listOfCommands
1008  * @param bedFile
1009  * @param bimFile
1010  * @param famFile
1011  * @param theChromo
1012  * @param mixedBedFile
1013  * @param mixedBimFile
1014  * @param mixedFamFile
1015  * @param mixedBedToBedLogFile
1016  */
1017 private static void doConvertFromBedToBed(ParseCmdLine parsingArgs,
1018     ArrayList<String> listOfCommands, String bedFile, String bimFile,
1019     String famFile, String theChromo, String mixedBedFile, String mixedBimFile,
1020     String mixedFamFile, String mixedBedToBedLogFile) {
1021
1022     if (parsingArgs.getStageStatus("convertFromBedToBed") == 1) {
1023         // Prepare the parameter values
1024         String noWebFlag = "--noweb";
1025         String bedPrefix = "--bed";
1026         String bimPrefix = "--bim";
1027         String famPrefix = "--fam";
1028         String chromoPrefix = "--chr";
1029         String recodeFlag = "--recode";
1030         String outPrefix = "--out";
1031         String makeBedFlag = "--make-bed";
1032         String stdout = mixedBedFile + ".stdout";
1033         String stderr = mixedBedFile + ".stderr";
1034
1035         // Store the command on the list
1036         String cmdToStore = ENV_VARS.PLINK_BINARY.value() + " " + noWebFlag + " "
1037             + bedPrefix + " " + bedFile + " " + bimPrefix + " "
1038             + bimFile + " " + famPrefix + " " + famFile + " " + chromoPrefix
1039             + " " + theChromo + " " + recodeFlag + " " + outPrefix
1040             + " " + mixedBedFile + " " + makeBedFlag + " #" + mixedBedFile
1041             + " #" + mixedBimFile + " #" + mixedFamFile + " #"
1042             + mixedBedToBedLogFile + " stdout " + stdout + " stderr " + stderr;
1043         listOfCommands.add(cmdToStore);
1044
1045         // Launch the task
1046         // We do not capture the ExitValue because the task will fail and
1047         // we do not want to synchronize
1048         BINARY.convertFromBedToBed(noWebFlag, bedPrefix, bedFile, bimPrefix, bimFile,
1049             famPrefix, famFile, chromoPrefix, theChromo, recodeFlag, outPrefix,
1050             mixedBedFile, makeBedFlag, mixedBimFile, mixedFamFile,
1051             mixedBedToBedLogFile, stdout, stderr);
1052     }
1053 }
1054
1055 /**
1056  * Method the wraps the execution of createRsIdList task and store the command in
1057  * the listOfCommands
1058  *
1059  * @param parsingArgs
1060  * @param listOfCommands

```

```

1062     * @param mixedBimOrGenFile
1063     * @param exclCgatFlag
1064     * @param mixedPairsFile
1065     * @param inputFormat
1066     */
1067     private static void doCreateRsIdList(ParseCmdLine parsingArgs, ArrayList<String>
1068         listOfCommands, String mixedBimOrGenFile,
1069         String exclCgatFlag, String mixedPairsFile, String inputFormat) {
1070
1071         if (parsingArgs.getStageStatus("createRsIdList") == 1) {
1072             // Store the command on the list
1073             String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1074                 + "java createRsIdList " + mixedBimOrGenFile + " "
1075                 + exclCgatFlag + " " + mixedPairsFile + " " + inputFormat;
1076             listOfCommands.add(cmdToStore);
1077
1078             // Launch the task
1079             try {
1080                 GuidanceImpl.createRsIdList(mixedBimOrGenFile, exclCgatFlag,
1081                     mixedPairsFile, inputFormat);
1082             } catch (GuidanceTaskException gte) {
1083                 // COMPSs: This code is never reached since exception causes task to fail
1084                 System.err.println("[Guidance] Exception createRsIdList task");
1085                 System.err.println(gte.getMessage());
1086             }
1087         }
1088     }
1089
1090     /**
1091     * Method that wraps the execution of phasing task with bed input formats and store
1092     * the command in the listOfCommands
1093     *
1094     * @param parsingArgs
1095     * @param listOfCommands
1096     * @param theChromo
1097     * @param bedFile
1098     * @param bimFile
1099     * @param famFile
1100     * @param gmapFile
1101     * @param excludedSnpsFile
1102     * @param shapeitHapsFile
1103     * @param shapeitSampleFile
1104     * @param shapeitLogFile
1105     * @param filteredHaplotypesFile
1106     * @param filteredHaplotypesSampleFile
1107     * @param filteredHaplotypesLogFile
1108     * @param filteredHaplotypesVcfFile
1109     * @param listOfSnpsFile
1110     * @param exclCgatFlag
1111     * @param exclSVFlag
1112     */
1113     private static void doPhasingBed(ParseCmdLine parsingArgs, ArrayList<String>
1114         listOfCommands, String theChromo, String bedFile, String bimFile,
1115         String famFile, String gmapFile, String excludedSnpsFile,
1116         String shapeitHapsFile, String shapeitSampleFile, String shapeitLogFile,
1117         String filteredHaplotypesFile, String filteredHaplotypesSampleFile,
1118         String filteredHaplotypesLogFile, String filteredHaplotypesVcfFile,
1119         String listOfSnpsFile, String exclCgatFlag, String exclSVFlag) {
1120
1121         if (parsingArgs.getStageStatus("phasingBed") == 1) {
1122             // Prepare the parameter values
1123             String inputBedPrefix = "--input-bed";
1124             String inputMapPrefix = "--input-map";
1125             String chrXPrefix = ((theChromo.equals("23")) ? "--chrX : ");
1126             String outputMaxPrefix = " --output-max ";
1127             String threadPrefix = "--thread";
1128             String effectiveSizePrefix = "--effective-size";
1129             String outputLogPrefix = "--output-log";
1130             int numThreads = 16;
1131             int effectiveSize = 20_000;
1132             String shapeitStdOut = shapeitHapsFile + ".stdout";
1133             String shapeitStdErr = shapeitHapsFile + ".stderr";

```



```

1134     // Store the command on the list
1136     String cmdToStore = ENV_VARS.SHAPEIT_BINARY.value() + " " + inputBedPrefix
        + " " + bedFile + " " + bimFile + " " + famFile + " "
1138         + inputMapPrefix + " " + gmapFile + " " + chrXPrefix + " "
        + outputMaxPrefix + " " + shapeitHapsFile + " "
1140         + shapeitSampleFile + " " + threadPrefix + " " + numThreads + " "
        + effectiveSizePrefix + " " + effectiveSize + " "
1142         + outputLogPrefix + " " + shapeitLogFile + " stdout " + shapeitStdOut
        + " stderr " + shapeitStdErr;
1144
        listOfCommands.add(cmdToStore);
1146
        // Launch the task
1148     // We do not capture the ExitValue because the task will fail and
        // we do not want to synchronize
1150     BINARY.phasingBed(inputBedPrefix, bedFile, bimFile, famFile, inputMapPrefix,
        gmapFile, chrXPrefix, outputMaxPrefix, shapeitHapsFile,
1152     shapeitSampleFile, threadPrefix, numThreads, effectiveSizePrefix,
        effectiveSize, outputLogPrefix, shapeitLogFile, shapeitStdOut,
1154     shapeitStdErr);
    }
1156
    if (parsingArgs.getStageStatus("createListOfExcludedSnps") == 1) {
1158     // Store the command on the list
        String cmdToStore = ENV_VARS.JAVA_BINARY.value() + " "
1160         + "createListOfExcludedSnps.jar" + " " + shapeitHapsFile + " "
        + excludedSnpsFile + " " + exclCgatFlag + " " + exclSVFlag;
1162     listOfCommands.add(cmdToStore);

        // Launch the task
1164     try {
1166         GuidanceImpl.createListOfExcludedSnps(shapeitHapsFile, excludedSnpsFile,
            exclCgatFlag, exclSVFlag);
1168     } catch (GuidanceTaskException gte) {
        // COMPSs: This code is never reached since exception causes task to fail
1170     System.err.println("[Guidance] Exception createListOfExcludedSnps task");
        System.err.println(gte.getMessage());
1172     }
    }
1174

    if (parsingArgs.getStageStatus("filterHaplotypes") == 1) {
1176     // Prepare the parameter values
        String convertFlag = "--convert";
1178     String inputHapsPrefix = "--input-haps";
        String excludeSnpPrefix = "--exclude-snp";
1180     String outputHapsPrefix = "--output-haps";
        String outputLogPrefix = "--output-log";
1182     String outputVcfPrefix = "--output-vcf";

        // Store the command on the list
1184     String cmdToStore = ENV_VARS.SHAPEIT_BINARY.value() + " " + convertFlag + " "
        + inputHapsPrefix + " " + shapeitHapsFile + " " + shapeitSampleFile
1186         + " " + excludeSnpPrefix + " " + excludedSnpsFile + " "
        + outputHapsPrefix + " " + filteredHaplotypesFile + " "
1188         + filteredHaplotypesSampleFile + " " + outputLogPrefix + " "
        + filteredHaplotypesLogFile + " " + outputVcfPrefix
1190         + filteredHaplotypesVcfFile;
        listOfCommands.add(cmdToStore);
1192

        // Launch the task
1194     try {
1196         BINARY.filterHaplotypes(convertFlag, inputHapsPrefix, shapeitHapsFile,
            shapeitSampleFile, excludeSnpPrefix, excludedSnpsFile,
1198         outputHapsPrefix, filteredHaplotypesFile,
            filteredHaplotypesSampleFile, outputLogPrefix,
1200         filteredHaplotypesLogFile, outputVcfPrefix,
            filteredHaplotypesVcfFile);

1202         GuidanceImpl.postFilterHaplotypes(filteredHaplotypesFile, listOfSnpsFile);
1204     } catch (GuidanceTaskException gte) {
        // COMPSs: This code is never reached since exception causes task to fail

```

```

1206         System.err.println("[Guidance] Exception createListOfExcludedSnps task");
1207         System.err.println(gte.getMessage());
1208     }
1209 }
1210 }
1211
1212 /**
1213  * Method that wraps the execution of phasing task with gen input formats and store
1214  * the command in the listOfCommands
1215  *
1216  * @param parsingArgs
1217  * @param listOfCommands
1218  * @param theChromo
1219  * @param genFile
1220  * @param sampleFile
1221  * @param gmapFile
1222  * @param excludedSnpsFile
1223  * @param shapeitHapsFile
1224  * @param shapeitSampleFile
1225  * @param shapeitLogFile
1226  * @param filteredHaplotypesFile
1227  * @param filteredHaplotypesSampleFile
1228  * @param filteredHaplotypesLogFile
1229  * @param filteredHaplotypesVcfFile
1230  * @param listOfSnpsFile
1231  * @param exclCgatFlag
1232  * @param exclSVFlag
1233  */
1234 private static void doPhasing(ParseCmdLine parsingArgs, ArrayList<String>
1235     listOfCommands, String theChromo, String genFile, String sampleFile,
1236     String gmapFile, String excludedSnpsFile, String shapeitHapsFile,
1237     String shapeitSampleFile, String shapeitLogFile, String filteredHaplotypesFile,
1238     String filteredHaplotypesSampleFile, String filteredHaplotypesLogFile,
1239     String filteredHaplotypesVcfFile, String listOfSnpsFile, String exclCgatFlag,
1240     String exclSVFlag) {
1241
1242     if (parsingArgs.getStageStatus("phasing") == 1) {
1243         // Prepare the parameter values
1244         String inputGenPrefix = "--input-gen";
1245         String inputMapPrefix = "--input-map";
1246         String chromoFlag = ((theChromo.equals("23")) ? "--chrX" : "");
1247         String outputPrefix = " --output-max ";
1248         String threadPrefix = "--thread";
1249         String effectiveSizePrefix = "--effective-size";
1250         String outputLogPrefix = "--output-log";
1251         int numThreads = 16;
1252         int effectiveSize = 20_000;
1253         String shapeitStdOut = shapeitHapsFile + ".stdout";
1254         String shapeitStdErr = shapeitHapsFile + ".stderr";
1255
1256         // Store the command on the list
1257         String cmdToStore = ENV_VARS.SHAPEIT_BINARY.value() + " " + inputGenPrefix
1258             + " " + genFile + " " + sampleFile + " "
1259             + inputMapPrefix + " " + gmapFile + " " + chromoFlag + " "
1260             + outputPrefix + " " + shapeitHapsFile + " "
1261             + shapeitSampleFile + " " + threadPrefix + " " + numThreads + " "
1262             + effectiveSizePrefix + " " + effectiveSize + " "
1263             + outputLogPrefix + " " + shapeitLogFile + " stdout " + shapeitStdOut
1264             + " stderr " + shapeitStdErr;
1265         listOfCommands.add(cmdToStore);
1266
1267         // Launch the task
1268         BINARY.phasing(inputGenPrefix, genFile, sampleFile, inputMapPrefix, gmapFile,
1269             chromoFlag, outputPrefix, shapeitHapsFile, shapeitSampleFile,
1270             threadPrefix, numThreads, effectiveSizePrefix, effectiveSize,
1271             outputLogPrefix, shapeitLogFile, shapeitStdOut, shapeitStdErr);
1272     }
1273
1274     if (parsingArgs.getStageStatus("createListOfExcludedSnps") == 1) {
1275         // Store the command on the list
1276         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator

```

```

1278         + "java createListOfExcludedSnps.jar " + shapeitHapsFile + " "
1279           + excludedSnpsFile + " " + exclCgatFlag + " " + exclSVFlag;
1280
1281     listOfCommands.add(cmdToStore);
1282
1283     // Launch the task
1284     try {
1285         GuidanceImpl.createListOfExcludedSnps(shapeitHapsFile, excludedSnpsFile,
1286             exclCgatFlag, exclSVFlag);
1287     } catch (GuidanceTaskException gte) {
1288         // COMPSs: This code is never reached since exception causes task to fail
1289         System.err.println("[Guidance] Exception createListOfExcludedSnps task");
1290         System.err.println(gte.getMessage());
1291     }
1292 }
1293
1294 if (parsingArgs.getStageStatus("filterHaplotypes") == 1) {
1295     // Prepare the parameter values
1296     String convertFlag = "-convert";
1297     String inputHapsPrefix = "--input-haps";
1298     String excludeSnpPrefix = "--exclude-snp";
1299     String outputHapsPrefix = "--output-haps";
1300     String outputLogPrefix = "--output-log";
1301     String outputVcfPrefix = "--output-vcf";
1302
1303     // Store the command on the list
1304     String cmdToStore = ENV_VARS.SHAPEIT_BINARY.value() + " " + convertFlag + " "
1305         + inputHapsPrefix + " " + shapeitHapsFile + " " + shapeitSampleFile
1306         + " " + excludeSnpPrefix + " " + excludedSnpsFile + " "
1307         + outputHapsPrefix + " " + filteredHaplotypesFile + " "
1308         + filteredHaplotypesSampleFile + " " + outputLogPrefix + " "
1309         + filteredHaplotypesLogFile + " " + outputVcfPrefix
1310         + filteredHaplotypesVcfFile;
1311     listOfCommands.add(cmdToStore);
1312
1313     // Launch the task
1314     try {
1315         BINARY.filterHaplotypes(convertFlag, inputHapsPrefix, shapeitHapsFile,
1316             shapeitSampleFile, excludeSnpPrefix, excludedSnpsFile,
1317             outputHapsPrefix, filteredHaplotypesFile,
1318             filteredHaplotypesSampleFile, outputLogPrefix,
1319             filteredHaplotypesLogFile, outputVcfPrefix,
1320             filteredHaplotypesVcfFile);
1321         GuidanceImpl.postFilterHaplotypes(filteredHaplotypesFile, listOfSnpsFile);
1322     } catch (GuidanceTaskException gte) {
1323         // COMPSs: This code is never reached since exception causes task to fail
1324         System.err.println("[Guidance] Exception createListOfExcludedSnps task");
1325         System.err.println(gte.getMessage());
1326     }
1327 }
1328 }
1329
1330 /**
1331  * Method that wraps the execution of impute task and store the command in the
1332  * listOfCommands
1333  *
1334  * @param parsingArgs
1335  * @param listOfCommands
1336  * @param chrS
1337  * @param gmapFile
1338  * @param knownHapFile
1339  * @param legendFile
1340  * @param shapeitHapsFile
1341  * @param shapeitSampleFile
1342  * @param lim1S
1343  * @param lim2S
1344  * @param pairsFile
1345  * @param imputeFile
1346  * @param imputeFileInfo
1347  * @param imputeFileSummary
1348  * @param imputeFileWarnings
1349  */

```

```

1350 private static void doImputationWithImpute(ParseCmdLine parsingArgs, ArrayList<String>
1351     listOfCommands, String chrS, String gmapFile, String knownHapFile,
1352     String legendFile, String shapeitHapsFile, String shapeitSampleFile,
1353     String lim1S, String lim2S, String pairsFile, String imputeFile,
1354     String imputeFileInfo, String imputeFileSummary, String imputeFileWarnings) {
1355
1356     if (parsingArgs.getStageStatus("imputeWithImpute") == 1) {
1357         // Prepare the parameter values
1358         String prephasedFlag = "use_prephased_g";
1359         String mPrefix = "-m";
1360         String hPrefix = "-h";
1361         String lPrefix = "-l";
1362         String knownHapsPrefix = "-known_haps_g";
1363         String samplePrefix = "-sample_g";
1364         String intPrefix = "-int";
1365         String chromoFlag = ((chrS.equals("23")) ? "-chrX" : "");
1366         String exludeSnpsPrefix = "-exclude_snps_g";
1367         String imputeExcludedFlag = "-impute_excluded";
1368         String nePrefix = "-Ne";
1369         int ne = 20_000;
1370         String oPrefix = "-o";
1371         String iPrefix = "-i";
1372         String rPrefix = "-r";
1373         String wPrefix = "-w";
1374         String noSampleQCFlag = "-no_sample_qc_info";
1375         String oGzFlag = "-o_gz";
1376         String stdOutFile = imputeFile + ".stdout";
1377         String stdErrorFile = imputeFile + ".stderr";
1378
1379         // Store the command on the list
1380         String cmdToStore = ENV_VARS.IMPUTE2_BINARY.value() + " " + prephasedFlag
1381             + " " + mPrefix + " " + gmapFile + " " + hPrefix + " " + knownHapFile
1382             + " " + lPrefix + " " + legendFile + " " + knownHapsPrefix + " "
1383             + shapeitHapsFile + " " + samplePrefix + " " + shapeitSampleFile + " "
1384             + intPrefix + " " + lim1S + " " + lim2S + " " + chromoFlag + " "
1385             + exludeSnpsPrefix + " " + pairsFile + " " + imputeExcludedFlag + " "
1386             + nePrefix + " " + ne + " " + oPrefix + " " + imputeFile + " "
1387             + iPrefix + " " + imputeFileInfo + " " + rPrefix + " "
1388             + imputeFileSummary + wPrefix + " " + imputeFileWarnings + " "
1389             + noSampleQCFlag + " " + oGzFlag + " stdout " + stdOutFile
1390             + " stderr " + stdErrorFile;
1391
1392         listOfCommands.add(cmdToStore);
1393
1394         // Launch task
1395         BINARY.imputeWithImpute(prephasedFlag, mPrefix, gmapFile, hPrefix,
1396             knownHapFile, lPrefix, legendFile, knownHapsPrefix,
1397             shapeitHapsFile, samplePrefix, shapeitSampleFile, intPrefix,
1398             lim1S, lim2S, chromoFlag, exludeSnpsPrefix, pairsFile,
1399             imputeExcludedFlag, nePrefix, ne, oPrefix, imputeFile, iPrefix,
1400             imputeFileInfo, rPrefix, imputeFileSummary, wPrefix,
1401             imputeFileWarnings, noSampleQCFlag, oGzFlag, stdOutFile,
1402             stdErrorFile);
1403     }
1404 }
1405
1406 /**
1407  * Method that wraps the execution of Imputation task with Minimac2 and store the
1408  * command in the listOfCommands
1409  *
1410  *
1411  * @param parsingArgs
1412  * @param listOfCommands
1413  * @param knownHapFile
1414  * @param filteredHapsFile
1415  * @param filteredSampleFile
1416  * @param filteredListOfSnpsFile
1417  * @param imputedMMFileName
1418  * @param imputedMMInfoFile
1419  * @param imputedMMERateFile
1420  * @param imputedMMRecFile
1421  * @param imputedMMDoseFile

```

```

1422 * @param imputedMMLogFile
1423 * @param chrS
1424 * @param lim1S
1425 * @param lim2S
1426 */
1427 private static void doImputationWithMinimac(ParseCmdLine parsingArgs,
1428     ArrayList<String> listOfCommands, String knownHapFile,
1429     String filteredHapsFile, String filteredSampleFile,
1430     String filteredListOfSnpsFile, String imputedMMFileName,
1431     String imputedMMInfoFile, String imputedMMERateFile, String imputedMMRecFile,
1432     String imputedMMDoseFile, String imputedMMLogFile,
1433     String chrS, String lim1S, String lim2S) {
1434
1435     if (parsingArgs.getStageStatus("imputeWithMinimac") == 1) {
1436         // Prepare the parameter values
1437         String vcfReferenceFlag = "--vcfReference";
1438         String refHapsPrefix = "--refHaps";
1439         String snpsPrefix = "--snps";
1440         String shapeHapsPrefix = "--shape_haps";
1441         String samplePrefix = "--sample";
1442         String vcfStartPrefix = "--vcfstart";
1443         String vcfEndPrefix = "--vcfend";
1444         String chromoFlag = chrS.equals(MAX_NUMBER_OF_CHROMOSOMES_STR) ? "--chr" : "";
1445         String vcfWindowPrefix = "--vcfwindow";
1446         int vcfWindow = 250_000;
1447         String roundsPrefix = "--rounds";
1448         int rounds = 5;
1449         String statesPrefix = "--states";
1450         int states = 200;
1451         String prefixPrefix = "--prefix";
1452         String gzipFlag = "--gzip";
1453         String stderrFile = imputedMMFileName + ".stderr";
1454
1455         // Store the command on the list
1456         String cmdToStore = ENV_VARS.MINIMAC_BINARY.value() + " " + vcfReferenceFlag
1457             + " " + refHapsPrefix + " " + knownHapFile + " " + snpsPrefix + " "
1458             + filteredListOfSnpsFile + " " + shapeHapsPrefix + " "
1459             + filteredHapsFile + " " + samplePrefix + " " + filteredSampleFile
1460             + " " + vcfStartPrefix + " " + lim1S + " " + vcfEndPrefix + " "
1461             + lim2S + " " + chromoFlag + " " + chrS + " " + vcfWindowPrefix + " "
1462             + vcfWindow + " " + roundsPrefix + " " + rounds + " " + statesPrefix
1463             + " " + states + " " + prefixPrefix + " " + imputedMMFileName + " "
1464             + gzipFlag + " # " + imputedMMInfoFile + " # " + imputedMMERateFile
1465             + " # " + imputedMMRecFile + " # " + imputedMMDoseFile + " stdout "
1466             + imputedMMLogFile + " stderr " + stderrFile;
1467         listOfCommands.add(cmdToStore);
1468
1469         // Launch task
1470         BINARY.imputeWithMinimac(vcfReferenceFlag, refHapsPrefix, knownHapFile,
1471             snpsPrefix, filteredListOfSnpsFile, shapeHapsPrefix,
1472             filteredHapsFile, samplePrefix, filteredSampleFile, vcfStartPrefix,
1473             lim1S, vcfEndPrefix, lim2S, chromoFlag, chrS, vcfWindowPrefix,
1474             vcfWindow, roundsPrefix, rounds, statesPrefix, states, prefixPrefix,
1475             imputedMMFileName, gzipFlag, imputedMMInfoFile, imputedMMERateFile,
1476             imputedMMRecFile, imputedMMDoseFile, imputedMMLogFile, stderrFile);
1477     }
1478 }
1479
1480 /**
1481 * Method that wraps the execution of filterByInfo task and store the command in
1482 * the listOfCommands
1483 *
1484 * @param parsingArgs
1485 * @param listOfCommands
1486 * @param imputeFileInfo
1487 * @param filteredRsIdFile
1488 * @param infoThresholdS
1489 */
1490 private static void doFilterByInfo(ParseCmdLine parsingArgs, ArrayList<String>
1491     listOfCommands, String imputeFileInfo, String filteredRsIdFile,
1492     String infoThresholdS) {

```

```

1494     if (parsingArgs.getStageStatus("filterByInfo") == 1) {
1495         // We create the list of rsId that are greater or equal to the infoThreshold
1496         // Store the command on the list
1497         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1498             + "java filterByInfo " + imputeFileInfo + " " + filteredRsIdFile
1499             + " " + infoThresholdS;
1500         listOfCommands.add(cmdToStore);
1501
1502         // Launch the task
1503         try {
1504             GuidanceImpl.filterByInfo(imputeFileInfo, filteredRsIdFile,
1505                 infoThresholdS);
1506         } catch (GuidanceTaskException gte) {
1507             // COMPSs: This code is never reached since exception causes task to fail
1508             System.err.println("[Guidance] Exception filterByInfo task for controls");
1509             System.err.println(gte.getMessage());
1510         }
1511     }
1512 }
1513
1514 /**
1515  * Method that wraps the execution of qctoolS task and store the command in the
1516  * listOfCommands
1517  *
1518  * @param parsingArgs
1519  * @param listOfCommands
1520  * @param imputeFile
1521  * @param filteredRsIdFile
1522  * @param filteredFile
1523  * @param filteredLogFile
1524  */
1525 private static void doQctoolS(ParseCmdLine parsingArgs, ArrayList<String>
1526     listOfCommands, String imputeFile, String filteredRsIdFile,
1527     String filteredFile, String filteredLogFile) {
1528
1529     double mafThreshold = parsingArgs.getMafThreshold();
1530     String mafThresholdS = Double.toString(mafThreshold);
1531
1532     if (parsingArgs.getStageStatus("qctools") == 1) {
1533         // Prepare the parameter values
1534         String gPrefix = "-g";
1535         String ogPrefix = "-og";
1536         String inclRSIdsPrefix = "-incl-rsids";
1537         String omitChromosomeFlag = "-omit-chromosome";
1538         String forceFlag = "-force";
1539         String logPrefix = "-log";
1540         String mafPrefix = "-maf";
1541         String oneFlag = "1";
1542         String stdoutFile = filteredFile + ".stdout";
1543         String stderrFile = filteredFile + ".stderr";
1544         String filteredFileGz = filteredFile + ".gz";
1545
1546         // Store the command on the list
1547         String cmdToStore = ENV_VARS.QCTOOL_BINARY.value() + " " + gPrefix + " "
1548             + imputeFile + " " + ogPrefix + " " + filteredFile + " "
1549             + inclRSIdsPrefix + " " + filteredRsIdFile + " " + omitChromosomeFlag
1550             + " " + forceFlag + " " + logPrefix + " " + filteredLogFile + " "
1551             + mafPrefix + " " + mafThresholdS + " " + oneFlag + " stdout "
1552             + stdoutFile + " stderr " + stderrFile;
1553         listOfCommands.add(cmdToStore);
1554
1555         // Launch the task
1556         BINARY.qctoolS(gPrefix, imputeFile, ogPrefix, filteredFile, inclRSIdsPrefix,
1557             filteredRsIdFile, omitChromosomeFlag, forceFlag, logPrefix,
1558             filteredLogFile, mafPrefix, mafThresholdS, oneFlag, stdoutFile,
1559             stderrFile);
1560
1561         // Launch the post zip task
1562         try {
1563             GuidanceImpl.postSnptest(filteredFile, filteredFileGz);
1564         } catch (GuidanceTaskException gte) {
1565             // COMPSs: This code is never reached since exception causes task to fail

```

```

1566         System.err.println("[Guidance] Exception on zip post SNP Test files");
1567         System.err.println(gte.getMessage());
1568     }
1569 }
1570
1571 /**
1572  * Method that wraps the execution of snptest task and store the command in the
1573  * listOfCommands
1574  *
1575  * @param parsingArgs
1576  * @param listOfCommands
1577  * @param chrS
1578  * @param mergedGenFile
1579  * @param mergedSampleFile
1580  * @param snptestOutFile
1581  * @param snptestLogFile
1582  * @param responseVar
1583  * @param covariables
1584  */
1585 private static void doSnptest(ParseCmdLine parsingArgs, ArrayList<String>
1586     listOfCommands, String chrS, String mergedGenFile, String mergedSampleFile,
1587     String snptestOutFile, String snptestLogFile, String responseVar,
1588     String covariables) {
1589
1590     if (parsingArgs.getStageStatus("snptest") == 1) {
1591         // Prepare the parameter values
1592         String covarsSTR = covariables.replace(',', ' ');
1593         String dataPrefix = "-data";
1594         String oPrefix = "-0";
1595         String phenoPrefix = "-pheno";
1596         String covarsFlag = covariables.equals("none") ? "" : "-cov_names " + covarsSTR;
1597         String hweFlag = "-hwe";
1598         String logPrefix = "-log";
1599         final String chromo23Flags
1600             = "-method newml -assume_chromosome X -stratify_on sex -frequentist 1";
1601         final String chromoOtherFlags = "-method em -frequentist 1 2 3 4 5";
1602         String chromoFlag = chrS.equals(ChromoInfo.MAX_NUMBER_OF_CHROMOSOMES_STR)
1603             ? chromo23Flags : chromoOtherFlags;
1604         String stdoutFile = snptestLogFile + ".stdout";
1605         String stderrFile = snptestLogFile + ".stderr";
1606         String snpTestOutFileGz = snptestOutFile + ".gz";
1607
1608         // Store the command on the list
1609         String cmdToStore = ENV_VARS.SNPTEST_BINARY.value() + dataPrefix + " "
1610             + mergedGenFile + " " + mergedSampleFile + " " + oPrefix
1611             + " " + snptestOutFile + " " + phenoPrefix + " " + responseVar
1612             + " " + covarsFlag + " " + hweFlag + " " + logPrefix
1613             + " " + snptestLogFile + " " + chromoFlag + " " + stdoutFile
1614             + " " + stderrFile;
1615         listOfCommands.add(cmdToStore);
1616
1617         // Launch the task
1618         BINARY.snptest(dataPrefix, mergedGenFile, mergedSampleFile, oPrefix,
1619             snptestOutFile, phenoPrefix, responseVar, covarsFlag, hweFlag,
1620             logPrefix, snptestLogFile, chromoFlag, stdoutFile, stderrFile);
1621
1622         // Launch the post zip task
1623         try {
1624             GuidanceImpl.postSnptest(snptestOutFile, snpTestOutFileGz);
1625         } catch (GuidanceTaskException gte) {
1626             // COMPSS: This code is never reached since exception causes task to fail
1627             System.err.println("[Guidance] Exception on zip post SNP Test files");
1628             System.err.println(gte.getMessage());
1629         }
1630     }
1631 }
1632
1633 /**
1634  * Method that wraps the execution of collectSummary task and store the command in
1635  * the listOfCommands
1636  */

```

```

1638 * @param parsingArgs
1640 * @param listOfCommands
1642 * @param chrS
1644 * @param imputeFileInfo
1646 * @param snptestOutFile
1648 * @param summaryFile
1650 * @param mafThresholdS
1652 * @param infoThresholdS
1654 * @param hweCohortThresholdS
1656 * @param hweCasesThresholdS
1658 * @param hweControlsThresholdS
1660 */
1662 private static void doCollectSummary(ParseCmdLine parsingArgs, ArrayList<String>
1664 listOfCommands, String chrS, String imputeFileInfo, String snptestOutFile,
1666 String summaryFile, String mafThresholdS, String infoThresholdS,
1668 String hweCohortThresholdS, String hweCasesThresholdS,
1670 String hweControlsThresholdS) {
1672
1674     if (parsingArgs.getStageStatus("collectSummary") == 1) {
1676         // Submitting the collect_summary task per this chunk
1678
1680         // Store the command on the list
1682         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1684             + "java collectSummary " + chrS + " " + imputeFileInfo + " "
1686             + snptestOutFile + " " + summaryFile + " " + mafThresholdS + " "
1688             + infoThresholdS + " " + hweCohortThresholdS + " "
1690             + hweCasesThresholdS + " " + hweControlsThresholdS;
1692
1694         listOfCommands.add(cmdToStore);
1696
1698         // Launch the task
1700         try {
1702             GuidanceImpl.collectSummary(chrS, imputeFileInfo, snptestOutFile,
1704                 summaryFile, mafThresholdS, infoThresholdS,
1706                 hweCohortThresholdS, hweCasesThresholdS, hweControlsThresholdS);
1708         } catch (GuidanceTaskException gte) {
1710             // COMPSS: This code is never reached since exception causes task to fail
1712             System.err.println("[Guidance] Exception collectSummary task");
1714             System.err.println(gte.getMessage());
1716         }
1718     }
1720 }
1722
1724 /**
1726 * Method that wraps the jointCondensedFiles task and store the command in
1728 * the listOfCommands
1730 *
1732 * @param parsingArgs
1734 * @param listOfCommands
1736 * @param condensedA
1738 * @param condensedB
1740 * @param condensedC
1742 */
1744 private static void doJointCondenseFiles(ParseCmdLine parsingArgs,
1746 ArrayList<String> listOfCommands, String condensedA,
1748 String condensedB, String condensedC) {
1750
1752     if (parsingArgs.getStageStatus("jointCondensedFiles") == 1) {
1754         // Store the command on the list
1756         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1758             + "java jointCondensedFiles " + condensedA + " " + condensedB
1760             + " " + condensedC;
1762         listOfCommands.add(cmdToStore);
1764
1766         // Launch the task
1768         try {
1770             GuidanceImpl.jointCondensedFiles(condensedA, condensedB, condensedC);
1772         } catch (GuidanceTaskException gte) {
1774             // COMPSS: This code is never reached since exception causes task to fail
1776             System.err.println("[Guidance] Exception jointCondensedFiles task");
1778             System.err.println(gte.getMessage());
1780         }
1782     }
1784 }

```



```

1710     }
1711 }
1712
1713 /**
1714  * Method that wraps the doJointFilteredByAllFiles task and store the command in
1715  * the listOfCommands
1716  *
1717  * @param parsingArgs
1718  * @param listOfCommands
1719  * @param filteredByAllA
1720  * @param filteredByAllB
1721  * @param filteredByAllC
1722  * @param rpanelName
1723  * @param rpanelFlag
1724  */
1725 private static void doJointFilteredByAllFiles(ParseCmdLine parsingArgs,
1726     ArrayList<String> listOfCommands, String filteredByAllA,
1727     String filteredByAllB, String filteredByAllC, String rpanelName,
1728     String rpanelFlag) {
1729
1730     if (parsingArgs.getStageStatus("jointFilteredByAllFiles") == 1) {
1731         // Store the command on the list
1732         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1733             + "java jointFilteredByAllFiles " + filteredByAllA + " "
1734             + filteredByAllB + " " + filteredByAllC + " " + rpanelName
1735             + " " + rpanelFlag;
1736         listOfCommands.add(cmdToStore);
1737
1738         // Launch the task
1739         try {
1740             GuidanceImpl.jointFilteredByAllFiles(filteredByAllA, filteredByAllB,
1741                 filteredByAllC, rpanelName, rpanelFlag);
1742         } catch (GuidanceTaskException gte) {
1743             // COMPSs: This code is never reached since exception causes task to fail
1744             System.err.println("[Guidance] Exception getBestSnps task");
1745             System.err.println("The error message here is " + gte.getMessage());
1746         }
1747     }
1748 }
1749
1750 /**
1751  * Method that wraps the generateTopHits task and store the command in
1752  * the listOfCommands
1753  *
1754  * @param parsingArgs
1755  * @param listOfCommands
1756  * @param filteredFile
1757  * @param filteredXFile
1758  * @param topHitsResults
1759  * @param pvaThrS
1760  */
1761 private static void doGenerateTopHits(ParseCmdLine parsingArgs,
1762     ArrayList<String> listOfCommands, String filteredFile,
1763     String filteredXFile, String topHitsResults, String pvaThrS) {
1764
1765     if (parsingArgs.getStageStatus("generateTopHits") == 1) {
1766         // Store the command on the list
1767         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1768             + "java generateTopHits " + filteredFile + " " + filteredXFile
1769             + " " + topHitsResults + " " + pvaThrS;
1770         listOfCommands.add(cmdToStore);
1771
1772         // Launch the task
1773         try {
1774             GuidanceImpl.generateTopHitsAll(filteredFile, filteredXFile,
1775                 topHitsResults, pvaThrS);
1776         } catch (GuidanceTaskException gte) {
1777             // COMPSs: This code is never reached since exception causes task to fail
1778             System.err.println("[Guidance] Exception generateTopHits task");
1779             System.err.println(gte.getMessage());
1780         }
1781     }
1782 }

```

```

1782     }
1783
1784     /**
1785      * Method that wraps the generateQQManhattanPlots task and store the command in
1786      * the listOfCommands
1787      *
1788      * @param parsingArgs
1789      * @param listOfCommands
1790      * @param condensedFile
1791      * @param qqPlotFile
1792      * @param manhattanPlotFile
1793      * @param qqPlotTiffFile
1794      * @param manhattanPlotTiffFile
1795      * @param correctedPvaluesFile
1796      */
1797     private static void doGenerateQQManhattanPlots(ParseCmdLine parsingArgs,
1798         ArrayList<String> listOfCommands, String condensedFile,
1799         String qqPlotFile, String manhattanPlotFile, String qqPlotTiffFile,
1800         String manhattanPlotTiffFile, String correctedPvaluesFile) {
1801
1802         if (parsingArgs.getStageStatus("generateQQManhattanPlots") == 1) {
1803             // Prepare the parameter values
1804             String stdoutFile = correctedPvaluesFile + ".stdout";
1805             String stderrFile = correctedPvaluesFile + ".stderr";
1806
1807             // Store the command on the list
1808             String cmdToStore = ENV_VARS.R_SCRIPT_BIN_DIR.value() + File.separator
1809                 + "Rscript " + ENV_VARS.R_SCRIPT_DIR.value() + File.separator
1810                 + "qqplot_manhattan.R " + condensedFile + " " + qqPlotFile + " "
1811                 + manhattanPlotFile + " " + qqPlotTiffFile + " "
1812                 + manhattanPlotTiffFile + " " + correctedPvaluesFile + " stdout "
1813                 + stdoutFile + " stderr " + stderrFile;
1814             listOfCommands.add(cmdToStore);
1815
1816             // Launch the unzip task
1817             String condensedFileUnzip = condensedFile + ".templ";
1818             try {
1819                 GuidanceImpl.preGenerateQQManhattanPlots(condensedFile,
1820                     condensedFileUnzip);
1821             } catch (GuidanceTaskException gte) {
1822                 // COMPSs: This code is never reached since exception causes task to fail
1823                 System.err.println("[Guidance] Exception QQ Manhattan Plot files");
1824                 System.err.println(gte.getMessage());
1825             }
1826
1827             // Launch the task
1828             BINARY.generateQQManhattanPlots(condensedFileUnzip, qqPlotFile,
1829                 manhattanPlotFile, qqPlotTiffFile, manhattanPlotTiffFile,
1830                 correctedPvaluesFile, stdoutFile, stderrFile);
1831         }
1832     }
1833
1834     /**
1835      * Method that wraps the combinePanelsComplex task and store the command in the
1836      * listOfCommands
1837      *
1838      * @param parsingArgs
1839      * @param listOfCommands
1840      * @param resultsPanelA
1841      * @param resultsPanelB
1842      * @param lastResultFile
1843      * @param startChrS
1844      * @param endChrS
1845      */
1846     private static void doCombinePanelsComplex(ParseCmdLine parsingArgs,
1847         ArrayList<String> listOfCommands, String resultsPanelA, String resultsPanelB,
1848         String lastResultFile, String startChrS, String endChrS) {
1849
1850         if (parsingArgs.getStageStatus("combinePanelsComplex") == 1) {
1851             // Store the command on the list
1852             String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1853                 + "java combinePanelsComplex " + resultsPanelA + " " + resultsPanelB

```

```

1854         + " " + lastResultFile + " " + startChrS + " " + endChrS;
        listOfCommands.add(cmdToStore);
1856
        // Launch the task
1858     try {
        GuidanceImpl.combinePanelsComplex(resultsPanelA, resultsPanelB,
1860         lastResultFile, startChrS, endChrS);
    } catch (GuidanceTaskException gte) {
1862         // COMPSs: This code is never reached since exception causes task to fail
        System.err.println("[Guidance] Exception combinePanelsComplex task");
1864         System.err.println(gte.getMessage());
    }
1866 }
1868
1869 /**
1870  * Method that wraps the doCombineCondensedFiles task and store the command in
1871  * the listOfCommands
1872  *
1873  * @param parsingArgs
1874  * @param listOfCommands
1875  * @param filteredA
1876  * @param filteredX
1877  * @param combinedCondensedFile
1878  */
1879 private static void doCombineCondensedFiles(ParseCmdLine parsingArgs,
1880     ArrayList<String> listOfCommands, String filteredA,
1881     String filteredX, String combinedCondensedFile) {
1882
1883     double mafThreshold = parsingArgs.getMafThreshold();
1884     double infoThreshold = parsingArgs.getInfoThreshold();
1885     double hweCohortThreshold = parsingArgs.getHweCohortThreshold();
1886     double hweCasesThreshold = parsingArgs.getHweCasesThreshold();
1887     double hweControlsThreshold = parsingArgs.getHweControlsThreshold();
1888
1889     String mafThresholdS = Double.toString(mafThreshold);
1890     String infoThresholdS = Double.toString(infoThreshold);
1891     String hweCohortThresholdS = Double.toString(hweCohortThreshold);
1892     String hweCasesThresholdS = Double.toString(hweCasesThreshold);
1893     String hweControlsThresholdS = Double.toString(hweControlsThreshold);
1894
1895     if (parsingArgs.getStageStatus("combineCondensedFiles") == 1) {
1896         // Store the command on the list
1897         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1898             + "java combineCondensedFiles " + filteredA + " " + filteredX
1899             + " " + combinedCondensedFile + " " + mafThresholdS + " "
1900             + infoThresholdS + " " + hweCohortThresholdS + " "
1901             + hweCasesThresholdS + " " + hweControlsThresholdS;
1902         listOfCommands.add(cmdToStore);
1903
1904         // Launch the task
1905     try {
1906         GuidanceImpl.combineCondensedFiles(filteredA, filteredX,
1907             combinedCondensedFile, mafThresholdS, infoThresholdS,
1908             hweCohortThresholdS, hweCasesThresholdS, hweControlsThresholdS);
1909     } catch (GuidanceTaskException gte) {
1910         // COMPSs: This code is never reached since exception causes task to fail
1911         System.err.println("[Guidance] Exception combineCondensedFile task");
1912         System.err.println(gte.getMessage());
1913     }
1914 }
1915
1916 /**
1917  * Method that wraps the mergeTwoChunks task and store the command in
1918  * the listOfCommands
1919  *
1920  * @param parsingArgs
1921  * @param listOfCommands
1922  * @param reduceA
1923  * @param reduceB
1924  * @param reduceC

```

```

1926     * @param theChromo
1927     */
1928     private static void doMergeTwoChunks(ParseCmdLine parsingArgs,
1929         ArrayList<String> listOfCommands, String reduceA, String reduceB,
1930         String reduceC, String theChromo) {
1931
1932         if (parsingArgs.getStageStatus("mergeTwoChunks") == 1) {
1933             // Store the command on the list
1934             String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1935                 + "java mergeTwoChunks " + reduceA + " " + reduceB + " "
1936                 + reduceC + " " + theChromo;
1937             listOfCommands.add(cmdToStore);
1938
1939             // Launch the task
1940             try {
1941                 GuidanceImpl.mergeTwoChunks(reduceA, reduceB, reduceC, theChromo);
1942             } catch (GuidanceTaskException gte) {
1943                 // COMPSs: This code is never reached since exception causes task to fail
1944                 System.err.println("[Guidance] Exception mergeTwoChunks task");
1945                 System.err.println(gte.getMessage());
1946             }
1947         }
1948     }
1949
1950     /**
1951     * Method that wraps the filterByAll task and store the command in the listOfCommands
1952     *
1953     * @param parsingArgs
1954     * @param listOfCommands
1955     * @param inputFile
1956     * @param outputFile
1957     * @param outputCondensedFile
1958     */
1959     private static void doFilterByAll(ParseCmdLine parsingArgs,
1960         ArrayList<String> listOfCommands, String inputFile, String outputFile,
1961         String outputCondensedFile) {
1962
1963         double mafThreshold = parsingArgs.getMafThreshold();
1964         double infoThreshold = parsingArgs.getInfoThreshold();
1965         double hweCohortThreshold = parsingArgs.getHweCohortThreshold();
1966         double hweCasesThreshold = parsingArgs.getHweCasesThreshold();
1967         double hweControlsThreshold = parsingArgs.getHweControlsThreshold();
1968
1969         String mafThresholdS = Double.toString(mafThreshold);
1970         String infoThresholdS = Double.toString(infoThreshold);
1971         String hweCohortThresholdS = Double.toString(hweCohortThreshold);
1972         String hweCasesThresholdS = Double.toString(hweCasesThreshold);
1973         String hweControlsThresholdS = Double.toString(hweControlsThreshold);
1974
1975         // Task
1976         if (parsingArgs.getStageStatus("filterByAll") == 1) {
1977             // Store the command on the list
1978             String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
1979                 + "java filterByAll " + inputFile + " " + outputFile + " "
1980                 + outputCondensedFile + " " + mafThresholdS + " " + infoThresholdS
1981                 + " " + hweCohortThresholdS + " " + hweCasesThresholdS + " "
1982                 + hweControlsThresholdS;
1983             listOfCommands.add(cmdToStore);
1984
1985             // Launch the task
1986             try {
1987                 GuidanceImpl.filterByAll(inputFile, outputFile, outputCondensedFile,
1988                     mafThresholdS, infoThresholdS, hweCohortThresholdS,
1989                     hweCasesThresholdS, hweControlsThresholdS);
1990             } catch (GuidanceTaskException gte) {
1991                 // COMPSs: This code is never reached since exception causes task to fail
1992                 System.err.println("[Guidance] Exception filterByAll task");
1993                 System.err.println(gte.getMessage());
1994             }
1995         }
1996     }

```

```
1998  /**
2000   * Method that wraps the initPhenoMatrix task and store the command in the
2002   * listOfCommands
2004   *
2006   * @param parsingArgs
2008   * @param listOfCommands
2010   * @param topHitsFile
2012   * @param ttName
2014   * @param rpName
2016   * @param phenomeFile
2018   */
2020 private static void doInitPhenoMatrix(ParseCmdLine parsingArgs,
2022   ArrayList<String> listOfCommands, String topHitsFile, String ttName,
2024   String rpName, String phenomeFile) {
2026
2028   if (parsingArgs.getStageStatus("initPhenoMatrix") == 1) {
2030     // Store the command on the list
2032     String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
2034     + "java initPhenoMatrix " + topHitsFile + " " + ttName + " "
2036     + rpName + " " + phenomeFile;
2038     listOfCommands.add(cmdToStore);
2040
2042     // Launch the task
2044     try {
2046       GuidanceImpl.initPhenoMatrix(topHitsFile, ttName, rpName, phenomeFile);
2048     } catch (GuidanceTaskException gte) {
2050       // COMPSs: This code is never reached since exception causes task to fail
2052       System.err.println("[Guidance] Exception initPhenoMatrix task");
2054       System.err.println(gte.getMessage());
2056     }
2058   }
2060 }
2062
2064 /**
2066   * Method that wraps the addToPhenoMatrix task and store the command in the
2068   * listOfCommands
2070   *
2072   * @param parsingArgs
2074   * @param listOfCommands
2076   * @param phenomeFileA
2078   * @param topHitsFile
2080   * @param ttName
2082   * @param rpName
2084   * @param phenomeFileB
2086   */
2088 private static void doAddToPhenoMatrix(ParseCmdLine parsingArgs,
2090   ArrayList<String> listOfCommands, String phenomeFileA,
2092   String topHitsFile, String ttName, String rpName, String phenomeFileB) {
2094
2096   if (parsingArgs.getStageStatus("addToPhenoMatrix") == 1) {
2098     // Store the command on the list
2100     String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
2102     + "java addToPhenoMatrix " + phenomeFileA + " " + topHitsFile
2104     + " " + ttName + " " + rpName + " " + phenomeFileB;
2106     listOfCommands.add(cmdToStore);
2108
2110     // Launch the task
2112     try {
2114       GuidanceImpl.addToPhenoMatrix(phenomeFileA, topHitsFile, ttName,
2116       rpName, phenomeFileB);
2118     } catch (GuidanceTaskException gte) {
2120       // COMPSs: This code is never reached since exception causes task to fail
2122       System.err.println("[Guidance] Exception addToPhenoMatrix task");
2124       System.err.println(gte.getMessage());
2126     }
2128   }
2130 }
2132
2134 /**
2136   * Method that wraps the filloutPhenoMatrix task and store the command in
2138   * the listOfCommands
2140   *
2142   */
```

```

2070 * @param parsingArgs
2071 * @param listOfCommands
2072 * @param phenomeFileA
2073 * @param filteredByAllFile
2074 * @param filteredByAllXFile
2075 * @param endChrS
2076 * @param ttName
2077 * @param rpName
2078 * @param phenomeFileB
2079 */
2080 private static void doFilloutPhenoMatrix(ParseCmdLine parsingArgs,
2081     ArrayList<String> listOfCommands, String phenomeFileA,
2082     String filteredByAllFile, String filteredByAllXFile, String endChrS,
2083     String ttName, String rpName, String phenomeFileB) {
2084
2085     if (parsingArgs.getStageStatus("filloutPhenoMatrix") == 1) {
2086         // Store the command on the list
2087         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
2088             + "java filloutPhenoMatrix " + phenomeFileA + " " + filteredByAllFile
2089             + " " + filteredByAllXFile + " " + endChrS + " " + ttName + " "
2090             + rpName + " " + phenomeFileB;
2091         listOfCommands.add(cmdToStore);
2092
2093         // Launch the task
2094         try {
2095             GuidanceImpl.filloutPhenoMatrix(phenomeFileA, filteredByAllFile,
2096                 filteredByAllXFile, endChrS, ttName, rpName, phenomeFileB);
2097         } catch (GuidanceTaskException gte) {
2098             // COMPSs: This code is never reached since exception causes task to fail
2099             System.err.println("[Guidance] Exception filloutPhenoMatrix task");
2100             System.err.println(gte.getMessage());
2101         }
2102     }
2103 }
2104
2105 /**
2106 * Method that wraps the finalizePhenoMatrix task and store the command in the
2107 * listOfCommands
2108 *
2109 * @param parsingArgs
2110 * @param listOfCommands
2111 * @param phenomeFileA
2112 * @param phenomeFileB
2113 * @param ttName
2114 * @param rpName
2115 * @param phenomeFileC
2116 */
2117 private static void doFinalizePhenoMatrix(ParseCmdLine parsingArgs,
2118     ArrayList<String> listOfCommands, String phenomeFileA,
2119     String phenomeFileB, String ttName, String rpName, String phenomeFileC) {
2120
2121     if (parsingArgs.getStageStatus("finalizePhenoMatrix") == 1) {
2122         // Store the command on the list
2123         String cmdToStore = ENV_VARS.JAVA_HOME.value() + File.separator
2124             + "java finalizePhenoMatrix " + phenomeFileA + " "
2125             + phenomeFileB + " " + ttName + " " + rpName + " " + phenomeFileC;
2126         listOfCommands.add(cmdToStore);
2127
2128         // Launch the task
2129         try {
2130             GuidanceImpl.finalizePhenoMatrix(phenomeFileA, phenomeFileB, ttName,
2131                 rpName, phenomeFileC);
2132         } catch (GuidanceTaskException gte) {
2133             // COMPSs: This code is never reached since exception causes task to fail
2134             System.err.println("[Guidance] Exception finalizePhenoMatrix task");
2135             System.err.println(gte.getMessage());
2136         }
2137     }
2138 }
2139
2140 /**
2141 * Method to set the final status of each file that has been generated in the

```

```

2142 * execution So far, the three main status are: uncompressed, compressed, and deleted.
2143 * The default status is: uncompressed. That is the initial status assigned
2144 * to each file
2145 *
2146 * @param parsingArgs
2147 * @param commonFilesInfo
2148 */
2149 private static void setFinalStatusForCommonFiles(ParseCmdLine parsingArgs,
2150 CommonFiles commonFilesInfo) {
2151 // String removeTempFiles = parsingArgs.getRemoveTemporalFiles();
2152 String compressFiles = parsingArgs.getCompressFiles();
2153 String finalStatus = null;
2154 if (compressFiles.equals("NO")) {
2155     finalStatus = "uncompressed";
2156 } else {
2157     finalStatus = "compressed";
2158 }
2159
2160 int startChr = parsingArgs.getStart();
2161 int endChr = parsingArgs.getEnd();
2162 for (int i = startChr; i <= endChr; i++) {
2163     commonFilesInfo.setPairsFileFinalStatus(i, finalStatus);
2164     commonFilesInfo.setShapeitHapsFileFinalStatus(i, finalStatus);
2165     commonFilesInfo.setShapeitSampleFileFinalStatus(i, finalStatus);
2166 }
2167 }
2168
2169 /**
2170 * Method to set the final status of each file that has been generated in the
2171 * execution So far, the three main status are: uncompressed, compressed, and deleted.
2172 * The default status is: uncompressed. That is the initial status assigned
2173 * to each file
2174 *
2175 * @param parsingArgs
2176 * @param imputationFilesInfo
2177 * @param generalChromoInfo
2178 * @param refPanels
2179 */
2180 private static void setFinalStatusForImputationFiles(ParseCmdLine parsingArgs,
2181 ImputationFiles imputationFilesInfo,
2182 ChromoInfo generalChromoInfo, List<String> refPanels) {
2183
2184 int startChr = parsingArgs.getStart();
2185 int endChr = parsingArgs.getEnd();
2186 int chunkSize = parsingArgs.getChunkSize();
2187 String imputationTool = parsingArgs.getImputationTool();
2188
2189 // String removeTempFiles = parsingArgs.getRemoveTemporalFiles();
2190 String compressFiles = parsingArgs.getCompressFiles();
2191 String finalStatus = null;
2192 if (compressFiles.equals("NO")) {
2193     finalStatus = "uncompressed";
2194 } else {
2195     finalStatus = "compressed";
2196 }
2197
2198 if (imputationTool.equals("impute")) {
2199     for (int j = 0; j < refPanels.size(); j++) {
2200         // String rPanel = refPanels.get(j);
2201         for (int chromo = startChr; chromo <= endChr; ++chromo) {
2202             int lim1 = 1;
2203             int lim2 = lim1 + chunkSize - 1;
2204             int numberOfChunks = generalChromoInfo.getMaxSize(chromo)/chunkSize;
2205             int module = generalChromoInfo.getMaxSize(chromo) % chunkSize;
2206             if (module != 0)
2207                 numberOfChunks++;
2208
2209             for (int k = 0; k < numberOfChunks; k++) {
2210                 imputationFilesInfo.setImputedFileFinalStatus(j, chromo, lim1,
2211                     lim2, chunkSize, finalStatus);
2212                 imputationFilesInfo.setFilteredFileFinalStatus(j, chromo, lim1,
2213                     lim2, chunkSize, finalStatus);

```

```

2214         imputationFilesInfo.setImputedInfoFileFinalStatus(j, chromo, lim1,
2215             lim2, chunkSize, finalStatus);
2216
2217         lim1 = lim1 + chunkSize;
2218         lim2 = lim2 + chunkSize;
2219     }
2220 }
2221 }
2222 } else if (imputationTool.equals("minimac")) {
2223     for (int j = 0; j < refPanels.size(); j++) {
2224         // String rPanel = refPanels.get(j);
2225         for (int chromo = startChr; chromo <= endChr; ++chromo) {
2226             int lim1 = 1;
2227             int lim2 = lim1 + chunkSize - 1;
2228             int numberOfChunks = generalChromoInfo.getMaxSize(chromo) / chunkSize;
2229             int module = generalChromoInfo.getMaxSize(chromo) % chunkSize;
2230             if (module != 0)
2231                 numberOfChunks++;
2232
2233             for (int k = 0; k < numberOfChunks; k++) {
2234                 imputationFilesInfo.setImputedMMInfoFileFinalStatus(j, chromo,
2235                     lim1, lim2, chunkSize, finalStatus);
2236                 lim1 = lim1 + chunkSize;
2237                 lim2 = lim2 + chunkSize;
2238             }
2239         }
2240     }
2241 }
2242 }
2243 }
2244
2245 /**
2246  * Method to set the final status of each file that has been generated in the
2247  * execution. So far, the three main status are: uncompressed, compressed, and deleted
2248  * The default status is: uncompressed. That is the initial status assigned
2249  * to each file
2250  *
2251  * @param parsingArgs
2252  * @param assocFilesInfo
2253  * @param generalChromoInfo
2254  * @param refPanels
2255  */
2256 private static void setFinalStatusForAssocFiles(ParseCmdLine parsingArgs,
2257     AssocFiles assocFilesInfo, ChromoInfo generalChromoInfo,
2258     List<String> refPanels) {
2259
2260     int startChr = parsingArgs.getStart();
2261     int endChr = parsingArgs.getEnd();
2262     int chunkSize = parsingArgs.getChunkSize();
2263
2264     // String removeTempFiles = parsingArgs.getRemoveTemporalFiles();
2265     String compressFiles = parsingArgs.getCompressFiles();
2266     String finalStatus = null;
2267     if (compressFiles.equals("NO")) {
2268         finalStatus = "uncompressed";
2269     } else {
2270         finalStatus = "compressed";
2271     }
2272
2273     // Now we continue with the association
2274     int numberOfTestTypes = parsingArgs.getNumberOfTestTypeName();
2275     for (int tt = 0; tt < numberOfTestTypes; tt++) {
2276         for (int j = 0; j < refPanels.size(); j++) {
2277             // String rPanel = refPanels.get(j);
2278             for (int chromo = startChr; chromo <= endChr; ++chromo) {
2279                 int maxSize = generalChromoInfo.getMaxSize(chromo);
2280                 int total_chunks = maxSize / chunkSize;
2281                 int lim1 = 1;
2282                 int lim2 = lim1 + chunkSize - 1;
2283
2284                 for (int k = 0; k < total_chunks; k++) {
2285                     assocFilesInfo.setSnpTestOutFileFinalStatus(tt, j, chromo,

```



```

2286         lim1, lim2, chunkSize, finalStatus);
2287         assocFilesInfo.setSummaryFileFinalStatus(tt, j, chromo,
2288             lim1, lim2, chunkSize, finalStatus);
2289
2290         lim1 = lim1 + chunkSize;
2291         lim2 = lim2 + chunkSize;
2292     }
2293 }
2294 }
2295 }
2296 }
2297 }
2298 }

```

## D.2 GuidanceItf.java

```

package guidance;
2
import guidance.utils.Environment;
4 import integratedtoolkit.types.annotations.Constraints;
import integratedtoolkit.types.annotations.task.Binary;
6 import integratedtoolkit.types.annotations.task.Method;
import integratedtoolkit.types.annotations.Parameter;
8 import integratedtoolkit.types.annotations.parameter.Direction;
import integratedtoolkit.types.annotations.parameter.Stream;
10 import integratedtoolkit.types.annotations.parameter.Type;
12
public interface GuidanceItf {
14
    @Binary(binary = "$" + Environment.EV_PLINKBINARY )
16    @Constraints(computingUnits = "1", memorySize = "1.0")
    Integer convertFromBedToBed(
18        @Parameter(type = Type.STRING, direction = Direction.IN) String webFlag,
        @Parameter(type = Type.STRING, direction = Direction.IN) String bedPrefix,
20        @Parameter(type = Type.FILE, direction = Direction.IN) String bedFile,
        @Parameter(type = Type.STRING, direction = Direction.IN) String bimPrefix,
22        @Parameter(type = Type.FILE, direction = Direction.IN) String bimFile,
        @Parameter(type = Type.STRING, direction = Direction.IN) String famPrefix,
24        @Parameter(type = Type.FILE, direction = Direction.IN) String famFile,
        @Parameter(type = Type.STRING, direction = Direction.IN) String chromoPrefix,
26        @Parameter(type = Type.STRING, direction = Direction.IN) String chromo,
        @Parameter(type = Type.STRING, direction = Direction.IN) String recodeFlag,
28        @Parameter(type = Type.STRING, direction = Direction.IN) String outPrefix,
        @Parameter(type = Type.FILE, direction = Direction.OUT)String newBedFile,
30        @Parameter(type = Type.STRING, direction = Direction.IN) String makeBedFlag,
        @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
32        String newBimFile,
        @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
34        String newFamFile,
        @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
36        String newLogFile,
        @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
38        String stdoutOutputFile,
        @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
40        String stderrOutputFile
    );
42
    @Binary(binary = "$" + Environment.EV_SHAPEITBINARY)
44    @Constraints(computingUnits="16", memorySize = "20.0f")
    Integer phasingBed(
46        @Parameter(type = Type.STRING, direction = Direction.IN) String inputBedPrefix,
        @Parameter(type = Type.FILE, direction = Direction.IN) String bedFile,
48        @Parameter(type = Type.FILE, direction = Direction.IN) String bumFile,
        @Parameter(type = Type.FILE, direction = Direction.IN) String famFile,
50        @Parameter(type = Type.STRING, direction = Direction.IN) String inputMapPrefix,
        @Parameter(type = Type.FILE, direction = Direction.IN) String gmapFile,
52        @Parameter(type = Type.STRING, direction = Direction.IN) String chromoFlag,

```

```

54     @Parameter(type = Type.STRING, direction = Direction.IN) String outputPrefix,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitHapsFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitSampleFile,
56     @Parameter(type = Type.STRING, direction = Direction.IN) String threadPrefix,
    @Parameter(type = Type.INT, direction = Direction.IN) int numThreads,
58     @Parameter(type = Type.STRING, direction = Direction.IN)
    String effectiveSizePrefix,
60     @Parameter(type = Type.INT, direction = Direction.IN) int effectiveSize,
    @Parameter(type = Type.STRING, direction = Direction.IN) String outputLogPrefix,
62     @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitLogFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
64     String stdoutOutputFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
66     String stderrOutputFile
    );
68
    @Binary(binary = "$" + Environment.EV_SHAPEITBINARY)
70     @Constraints(computingUnits="16", memorySize = "20.0f")
    Integer phasing(
72         @Parameter(type = Type.STRING, direction = Direction.IN) String inputGenPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String inputGenFile,
74         @Parameter(type = Type.FILE, direction = Direction.IN) String inputSampleFile,
            @Parameter(type = Type.STRING, direction = Direction.IN) String inputMapPrefix,
76         @Parameter(type = Type.FILE, direction = Direction.IN) String gmapFile,
            @Parameter(type = Type.STRING, direction = Direction.IN) String chromoFlag,
78         @Parameter(type = Type.STRING, direction = Direction.IN) String outputPrefix,
            @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitHapsFile,
80         @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitSampleFile,
            @Parameter(type = Type.STRING, direction = Direction.IN) String threadPrefix,
82         @Parameter(type = Type.INT, direction = Direction.IN) int numThreads,
            @Parameter(type = Type.STRING, direction = Direction.IN)
84         String effectiveSizePrefix,
            @Parameter(type = Type.INT, direction = Direction.IN) int effectiveSize,
86         @Parameter(type = Type.STRING, direction = Direction.IN) String outputLogPrefix,
            @Parameter(type = Type.FILE, direction = Direction.OUT) String shapeitLogFile,
88         @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
            String stdoutOutputFile,
90         @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
            String stderrOutputFile
92     );

94     @Binary(binary = "$" + Environment.EV_IMPUTE2BINARY)
    @Constraints(computingUnits="1", memorySize = "14.0f")
96     Integer imputeWithImpute(
98         @Parameter(type = Type.STRING, direction = Direction.IN) String prephasedFlag,
            @Parameter(type = Type.STRING, direction = Direction.IN) String mPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String gmapFile,
100         @Parameter(type = Type.STRING, direction = Direction.IN) String hPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String knownHapFile,
102         @Parameter(type = Type.STRING, direction = Direction.IN) String lPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String legendFile,
104         @Parameter(type = Type.STRING, direction = Direction.IN) String knownHapsPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String shapeitHapsFile,
106         @Parameter(type = Type.STRING, direction = Direction.IN) String samplePrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String shapeitSampleFile,
108         @Parameter(type = Type.STRING, direction = Direction.IN) String intPrefix,
            @Parameter(type = Type.STRING, direction = Direction.IN) String lim1S,
110         @Parameter(type = Type.STRING, direction = Direction.IN) String lim2S,
            @Parameter(type = Type.STRING, direction = Direction.IN) String chromoFlag,
112         @Parameter(type = Type.STRING, direction = Direction.IN) String excludeSnpsPrefix,
            @Parameter(type = Type.FILE, direction = Direction.IN) String pairsFile,
114         @Parameter(type = Type.STRING, direction = Direction.IN)
            String imputeExcludedFlag,
116         @Parameter(type = Type.STRING, direction = Direction.IN) String nePrefix,
            @Parameter(type = Type.INT, direction = Direction.IN) int ne,
118         @Parameter(type = Type.STRING, direction = Direction.IN) String oPrefix,
            @Parameter(type = Type.FILE, direction = Direction.OUT) String imputeFile,
120         @Parameter(type = Type.STRING, direction = Direction.IN) String iPrefix,
            @Parameter(type = Type.FILE, direction = Direction.OUT) String imputeFileInfo,
122         @Parameter(type = Type.STRING, direction = Direction.IN) String rPrefix,
            @Parameter(type = Type.FILE, direction = Direction.OUT) String imputeFileSummary,
124         @Parameter(type = Type.STRING, direction = Direction.IN) String wPrefix,

```

```

126     @Parameter(type = Type.FILE, direction = Direction.OUT) String imputeFileWarnings,
    @Parameter(type = Type.STRING, direction = Direction.IN) String noSampleQCFlag,
    @Parameter(type = Type.STRING, direction = Direction.IN) String oGzFlag,
128     @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
    String stdoutFile,
130     @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
    String stderrFile
132 );

134 @Binary(binary = "$" + Environment.EV_MINIMACBINARY)
@Constraints(computingUnits="1", memorySize = "3.0f")
136 Integer imputeWithMinimac(
    @Parameter(type = Type.STRING, direction = Direction.IN) String vcfReferenceFlag,
138     @Parameter(type = Type.STRING, direction = Direction.IN) String refHapsPrefix,
    @Parameter(type = Type.FILE, direction = Direction.IN) String knownHapFile,
140     @Parameter(type = Type.STRING, direction = Direction.IN) String snpsPrefix,
    @Parameter(type = Type.FILE, direction = Direction.IN)
142     String filteredListOfSnpsFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String shapeHapsPrefix,
144     @Parameter(type = Type.FILE, direction = Direction.IN) String filteredHapsFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String samplePrefix,
146     @Parameter(type = Type.FILE, direction = Direction.IN) String filteredSampleFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String vcfStartPrefix,
148     @Parameter(type = Type.STRING, direction = Direction.IN) String lim1S,
    @Parameter(type = Type.STRING, direction = Direction.IN) String vcfEndPrefix,
150     @Parameter(type = Type.STRING, direction = Direction.IN) String lim2S,
    @Parameter(type = Type.STRING, direction = Direction.IN) String chromoFlag,
152     @Parameter(type = Type.STRING, direction = Direction.IN) String chromo,
    @Parameter(type = Type.STRING, direction = Direction.IN) String vcfWindowPrefix,
154     @Parameter(type = Type.INT, direction = Direction.IN) int vcfWindow,
    @Parameter(type = Type.STRING, direction = Direction.IN) String roundsPrefix,
156     @Parameter(type = Type.INT, direction = Direction.IN) int rounds,
    @Parameter(type = Type.STRING, direction = Direction.IN) String statesPrefix,
158     @Parameter(type = Type.INT, direction = Direction.IN) int states,
    @Parameter(type = Type.STRING, direction = Direction.IN) String prefixPrefix,
160     @Parameter(type = Type.STRING, direction = Direction.IN) String imputedMMFileName,
    @Parameter(type = Type.STRING, direction = Direction.IN) String gzipFlag,
162     @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
    String imputedMMInfoFile,
164     @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
    String imputedMMRateFile,
166     @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
    String imputedMMRecFile,
168     @Parameter(type = Type.FILE, direction = Direction.OUT, prefix = "#")
    String imputedMMDoseFile,
170     @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
    String stdoutFile,
172     @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
    String stderrFile
174 );

176 @Method(declaringClass = "guidance.GuidanceImpl")
@Constraints(computingUnits="1", memorySize = "14.0f")
178 void preGenerateQQManhattanPlots(
    @Parameter(type = Type.FILE, direction = Direction.IN) String condensedFile,
180     @Parameter(type = Type.FILE, direction = Direction.OUT) String condensedFileUnzip
    );

182
184 @Binary(binary = "$" + Environment.EV_RSCRIPTBINDIR
    + "/Rscript" + Environment.EV_RSCRIPTDIR + "/qqplot_manhattan.R")
@Constraints(computingUnits="1", memorySize = "14.0f")
186 Integer generateQQManhattanPlots(
    @Parameter(type = Type.FILE, direction = Direction.IN) String inputFile,
188     @Parameter(type = Type.FILE, direction = Direction.OUT) String qqPlotFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String manhattanPlotFile,
190     @Parameter(type = Type.FILE, direction = Direction.OUT) String qqPlotTiffFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT)
192     String manhattanPlotTiffFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT)
194     String correctedPvaluesFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
196     String stdoutFile,

```

```

198     @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
        String stderrFile
199     );
200
201     @Binary(binary = "$" + Environment.EV_SNPTESTBINARY)
202     @Constraints(computingUnits="1", memorySize = "2.0f")
        Integer snptest (
203         @Parameter(type = Type.STRING, direction = Direction.IN) String dataPrefix,
204         @Parameter(type = Type.FILE, direction = Direction.IN) String mergedGenFile,
205         @Parameter(type = Type.FILE, direction = Direction.IN)
206             String mergedSampleFile,
207         @Parameter(type = Type.STRING, direction = Direction.IN) String oPrefix,
208         @Parameter(type = Type.FILE, direction = Direction.OUT) String snpTestOutFile,
209         @Parameter(type = Type.STRING, direction = Direction.IN) String phenoPrefix,
210         @Parameter(type = Type.STRING, direction = Direction.IN) String responseVar,
211         @Parameter(type = Type.STRING, direction = Direction.IN) String covarsFlag,
212         @Parameter(type = Type.STRING, direction = Direction.IN) String hweFlag,
213         @Parameter(type = Type.STRING, direction = Direction.IN) String logPrefix,
214         @Parameter(type = Type.FILE, direction = Direction.OUT) String snpTestLogFile,
215         @Parameter(type = Type.STRING, direction = Direction.IN) String chromoFlag,
216         @Parameter(type = Type.FILE, direction = Direction.OUT,
217             stream = Stream.STDOUT) String stdoutFile,
218         @Parameter(type = Type.FILE, direction = Direction.OUT,
219             stream = Stream.STDERR) String stderrFile
220     );
221
222     @Method(declaringClass = "guidance.GuidanceImpl")
223     @Constraints(computingUnits="1", memorySize = "2.0f")
        void postSnptest (
224         @Parameter(type = Type.FILE, direction = Direction.IN) String snpTestOutFile,
225         @Parameter(type = Type.FILE, direction = Direction.OUT) String snpTestOutFileZip
226     );
227
228
229     @Binary(binary = "$" + Environment.EV_QCTOOLBINARY)
230     @Constraints(computingUnits="1", memorySize = "2.0f")
        Integer qctools (
231         @Parameter(type = Type.STRING, direction = Direction.IN) String gPrefix,
232         @Parameter(type = Type.FILE, direction = Direction.IN) String imputeFile,
233         @Parameter(type = Type.STRING, direction = Direction.IN) String ogPrefix,
234         @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredFile,
235         @Parameter(type = Type.STRING, direction = Direction.IN) String inclRSIDsPrefix,
236         @Parameter(type = Type.FILE, direction = Direction.IN) String inclusionRsIdFile,
237         @Parameter(type = Type.STRING, direction = Direction.IN)
238             String omitChromosomeFlag,
239         @Parameter(type = Type.STRING, direction = Direction.IN) String forceFlag,
240         @Parameter(type = Type.STRING, direction = Direction.IN) String logPrefix,
241         @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredLogFile,
242         @Parameter(type = Type.STRING, direction = Direction.IN) String mafPrefix,
243         @Parameter(type = Type.STRING, direction = Direction.IN) String mafThresholdS,
244         @Parameter(type = Type.STRING, direction = Direction.IN) String oneFlag,
245         @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDOUT)
246             String stdoutFile,
247         @Parameter(type = Type.FILE, direction = Direction.OUT, stream = Stream.STDERR)
248             String stderrFile
249     );
250
251
252     @Method(declaringClass = "guidance.GuidanceImpl")
253     @Constraints(computingUnits="1", memorySize = "2.0f")
        void postQCTool (
254         @Parameter(type = Type.FILE, direction = Direction.IN) String filteredFile,
255         @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredFileGz
256     );
257
258
259     @Binary(binary = "$" + Environment.EV_SHAPEITBINARY)
260     @Constraints(computingUnits="1", memorySize = "4.0f")
        void filterHaplotypes (
261         @Parameter(type = Type.STRING, direction = Direction.IN) String convertFlag,
262         @Parameter(type = Type.STRING, direction = Direction.IN) String inputHapsPrefix,
263         @Parameter(type = Type.FILE, direction = Direction.IN) String hapsFile,
264         @Parameter(type = Type.FILE, direction = Direction.IN) String sampleFile,
265         @Parameter(type = Type.STRING, direction = Direction.IN) String excludeSnpPrefix,
266         @Parameter(type = Type.FILE, direction = Direction.IN) String excludedSnpsFile,

```

```

270     @Parameter(type = Type.STRING, direction = Direction.IN) String outputHapsPrefix,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredHapsFile,
272     @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredSampleFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String outputLogPrefix,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredLogFile,
274     @Parameter(type = Type.STRING, direction = Direction.IN) String outputVcfPrefix,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredHapsVcfFile
276 );
278
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "4.0f")
280 void postFilterHaplotypes (
    @Parameter(type = Type.FILE, direction = Direction.IN) String filteredHapsFile,
282     @Parameter(type = Type.FILE, direction = Direction.OUT) String listOfSnpsFile
    );
284
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "1.0f")
286 void createRsIdList (
    @Parameter(type = Type.FILE, direction = Direction.IN) String genFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String exclCgatFlag,
290     @Parameter(type = Type.FILE, direction = Direction.OUT) String pairsFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String inputFormat
292 );
294
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "4.0f")
296 void createListOfExcludedSnps (
    @Parameter(type = Type.FILE, direction = Direction.IN) String shapeitHapsFile,
298     @Parameter(type = Type.FILE, direction = Direction.OUT) String excludedSnpsFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String exclCgatFlag,
300     @Parameter(type = Type.STRING, direction = Direction.IN) String exclSVFlag
    );
302
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "6.0f")
304 void filterByAll (
    @Parameter(type = Type.FILE, direction = Direction.IN) String inputFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String outputFile,
308     @Parameter(type = Type.FILE, direction = Direction.OUT) String outputCondensedFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String mafThresholds,
310     @Parameter(type = Type.STRING, direction = Direction.IN) String infoThresholds,
    @Parameter(type = Type.STRING, direction = Direction.IN)
312     String hweCohortThresholds,
    @Parameter(type = Type.STRING, direction = Direction.IN) String hweCasesThreshold,
314     @Parameter(type = Type.STRING, direction = Direction.IN)
    String mafControlsThreshold
316 );
318
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "6.0f")
320 void filterByInfo (
    @Parameter(type = Type.FILE, direction = Direction.IN) String imputeFileInfo,
322     @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String threshold
324 );
326
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "6.0f")
328 void jointFilteredByAllFiles (
    @Parameter(type = Type.FILE, direction = Direction.IN) String filteredByAllA,
330     @Parameter(type = Type.FILE, direction = Direction.IN) String filteredByAllB,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String filteredByAllC,
332     @Parameter(type = Type.STRING, direction = Direction.IN) String rpanelName,
    @Parameter(type = Type.STRING, direction = Direction.IN) String rpanelFlag
334 );
336
    @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "6.0f")
338 void jointCondensedFiles (
    @Parameter(type = Type.FILE, direction = Direction.IN) String inputAFile,
340     @Parameter(type = Type.FILE, direction = Direction.IN) String inputBFile,

```

```

342     @Parameter(type = Type.FILE, direction = Direction.OUT) String outputFile
    );

344     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
346     void collectSummary(
        @Parameter(type = Type.STRING, direction = Direction.IN) String chr,
348         @Parameter(type = Type.FILE, direction = Direction.IN) String firstImputeFileInfo,
        @Parameter(type = Type.FILE, direction = Direction.IN) String snptestOutFile,
350         @Parameter(type = Type.FILE, direction = Direction.OUT) String reduceFile,
        @Parameter(type = Type.STRING, direction = Direction.IN) String mafThresholdS,
352         @Parameter(type = Type.STRING, direction = Direction.IN) String infoThresholdS,
        @Parameter(type = Type.STRING, direction = Direction.IN)
354         String hweCohortThresholdS,
        @Parameter(type = Type.STRING, direction = Direction.IN) String hweCasesThresholdS,
356         @Parameter(type = Type.STRING, direction = Direction.IN)
        String hweControlsThresholdS
358     );

360     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
362     void initPhenoMatrix(
        @Parameter(type = Type.FILE, direction = Direction.IN) String topHitsFile,
364         @Parameter(type = Type.STRING, direction = Direction.IN) String ttName,
        @Parameter(type = Type.STRING, direction = Direction.IN) String rpName,
366         @Parameter(type = Type.FILE, direction = Direction.OUT) String phenomeFile
    );

368     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
370     void addToPhenoMatrix(
372         @Parameter(type = Type.FILE, direction = Direction.IN) String phenomeAFile,
        @Parameter(type = Type.FILE, direction = Direction.IN) String topHitsFile,
374         @Parameter(type = Type.STRING, direction = Direction.IN) String ttName,
        @Parameter(type = Type.STRING, direction = Direction.IN) String rpName,
376         @Parameter(type = Type.FILE, direction = Direction.OUT) String phenomeBFile
    );

378     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
380     void filloutPhenoMatrix(
382         @Parameter(type = Type.FILE, direction = Direction.IN) String phenomeAFile,
        @Parameter(type = Type.FILE, direction = Direction.IN) String filteredByAllFile,
384         @Parameter(type = Type.FILE, direction = Direction.IN) String filteredByAllXFile,
        @Parameter(type = Type.STRING, direction = Direction.IN) String thereisX,
386         @Parameter(type = Type.STRING, direction = Direction.IN) String ttName,
        @Parameter(type = Type.STRING, direction = Direction.IN) String rpName,
388         @Parameter(type = Type.FILE, direction = Direction.OUT) String phenomeBFile
    );

390     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
392     void finalizePhenoMatrix(
394         @Parameter(type = Type.FILE, direction = Direction.IN) String phenomeAFile,
        @Parameter(type = Type.FILE, direction = Direction.IN) String phenomeBFile,
396         @Parameter(type = Type.STRING, direction = Direction.IN) String ttName,
        @Parameter(type = Type.STRING, direction = Direction.IN) String rpName,
398         @Parameter(type = Type.FILE, direction = Direction.OUT) String phenomeCFile
    );

400     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
402     void mergeTwoChunks(
404         @Parameter(type = Type.FILE, direction = Direction.IN) String reduceFileA,
        @Parameter(type = Type.FILE, direction = Direction.IN) String reduceFileB,
406         @Parameter(type = Type.FILE, direction = Direction.OUT) String reduceFileC,
        @Parameter(type = Type.STRING, direction = Direction.IN) String chrS
408     );

410     @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "20.0f")
412     void combinePanelsComplex(

```

```
414     @Parameter(type = Type.FILE, direction = Direction.IN) String resultsFileA,
    @Parameter(type = Type.FILE, direction = Direction.IN) String resultsFileB,
416     @Parameter(type = Type.FILE, direction = Direction.OUT) String resultsFileC,
    @Parameter(type = Type.STRING, direction = Direction.IN) String chromoStart,
    @Parameter(type = Type.STRING, direction = Direction.IN) String chromoEnd
418 );

420 @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "6.0f")
422 void combineCondensedFiles(
    @Parameter(type = Type.FILE, direction = Direction.IN) String filteredA,
424     @Parameter(type = Type.FILE, direction = Direction.IN) String filteredX,
    @Parameter(type = Type.FILE, direction = Direction.OUT)
426     String combinedCondensedFile,
    @Parameter(type = Type.STRING, direction = Direction.IN) String mafThresholdS,
428     @Parameter(type = Type.STRING, direction = Direction.IN) String infoThresholdS,
    @Parameter(type = Type.STRING, direction = Direction.IN)
430     String hweCohortThresholdS,
    @Parameter(type = Type.STRING, direction = Direction.IN) String hweCasesThreshold,
432     @Parameter(type = Type.STRING, direction = Direction.IN)
    String mafControlsThreshold
434 );

436 @Method(declaringClass = "guidance.GuidanceImpl")
    @Constraints(computingUnits="1", memorySize = "8.0f")
438 void generateTopHitsAll(
    @Parameter(type = Type.FILE, direction = Direction.IN) String resultsAFile,
440     @Parameter(type = Type.FILE, direction = Direction.IN) String resultsBFile,
    @Parameter(type = Type.FILE, direction = Direction.OUT) String outputTopHitsFile,
442     @Parameter(type = Type.STRING, direction = Direction.IN) String pvaThreshold
    );
444 }
```