

Spring 2013

Large-scale Molecular Dynamics Simulation with Forward Flux Sampling on Hadoop

Pengfei Xuan

Yueli Zheng

Sapna Sarupria

Amy Apon

Follow this and additional works at: https://tigerprints.clemson.edu/grads_symposium

Recommended Citation

Xuan, Pengfei; Zheng, Yueli; Sarupria, Sapna; and Apon, Amy, "Large-scale Molecular Dynamics Simulation with Forward Flux Sampling on Hadoop" (2013). *Graduate Research and Discovery Symposium (GRADS)*. 59.
https://tigerprints.clemson.edu/grads_symposium/59

This Poster is brought to you for free and open access by the Research and Innovation Month at TigerPrints. It has been accepted for inclusion in Graduate Research and Discovery Symposium (GRADS) by an authorized administrator of TigerPrints. For more information, please contact kokeefe@clemson.edu.

ABSTRACT

Simulating rare events is extremely difficult and requires massive computational resources and complex data processing workflow, which is determined by the nature of stochastic systems. To help computational scientists discover hard scientific problems in this area, we built a large-scale molecular dynamics simulation framework integrated with forward flux sampling (FFS) technique on Hadoop ecosystem. In this project, we port the customized FFS workflow to underlying MapReduce-based computing pipeline by using dataflow-driven design pattern and Gromacs application. The early works show that our framework is able to provide a scalable, fault-tolerance and efficient rare events simulation environment over varieties of computing infrastructures, while preserving the flexibility of the original scientific application.

BACKGROUND

Rare events are fluctuation-driven processes occurring infrequently such as the nucleation of crystals or protein folding events. It exists widely as kinds of different phenomena, such as earthquakes, financial crashes, telecommunication network failure, protein conformational change, activated chemical reactions and so on. It is notoriously difficult to simulate because the very few events are observed in the traditional simulation. The traditional simulation methods such as “brute-force” simulation is highly inefficient in the rare event simulation, because the waiting time between events is hundreds and thousands of times longer than the time scale of the event itself in these processes [1]. In this project, we use FFS [2] sampling method to improve the efficiency of rare events simulation.

STUDY CASE

We are trying to use FFS to simulate the nucleation and growth of gas hydrates to obtain detailed molecular-level picture of hydrate formation kinetics [3]. In the simulation, to perform FFS, Bash Script concurrently call GROMACS programs to generate a large volume of intermediate files (more than 100 thousands) with TB-level data size, and then calculate threshold to filter parts of results. The massive MD simulation gives a big overhead on the this approach.

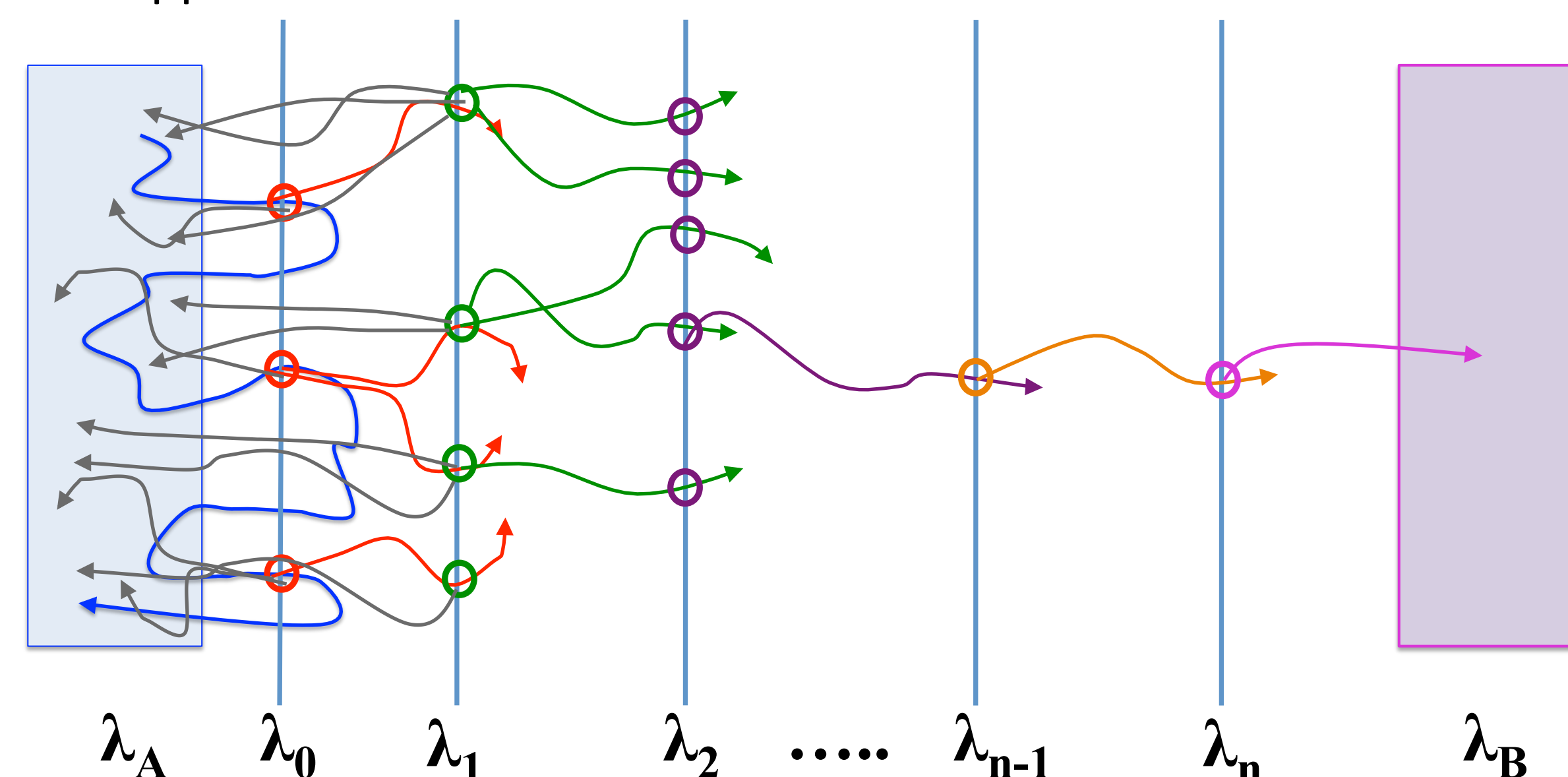


Figure 1. The first (A) and second (B) stages of the FFS method. The distribution of points at the interfaces depends on the history of the paths, as illustrated by the dashed lines in (b) [3].

SYSTEM ARCHITECTURE AND IMPLEMENTATION

Hadoop-based computing environment is designed to support the high-level dataflow-driven implementation by using various FFS techniques and Gromacs scientific application.

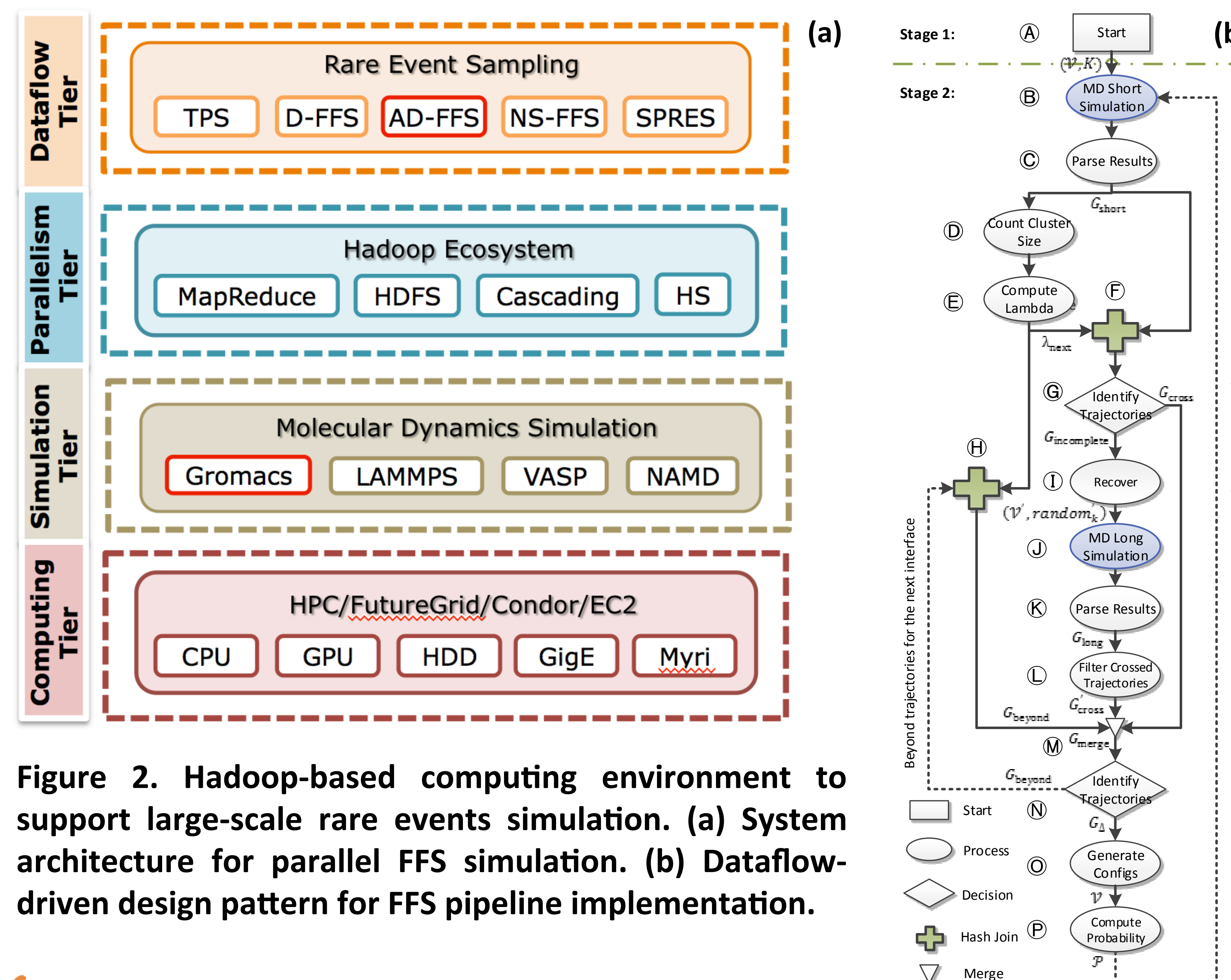


Figure 2. Hadoop-based computing environment to support large-scale rare events simulation. (a) System architecture for parallel FFS simulation. (b) Dataflow-driven design pattern for FFS pipeline implementation.

RUNTIME ENVIRONMENT

- Package FFS pipeline to JAR
- Upload Gromacs application and its libraries to Hadoop Distributed Cache
- Workflow engine loads FFS pipeline and assigns each MapReduce jobs to JobTracker
- Each Job is split into map tasks and reduce tasks in parallel and distributed over the entire HPC cluster
- TaskTracker in each computing node manages the execution of Gromacs and collects results from each short MD simulation
- The user or task in computing cluster can poll a global view of the status for all jobs and intermediate simulation results
- Application can directly makes use of underlying hardware features to speedup the performance (e.g. SSD, GPU and MIC coprocessor)

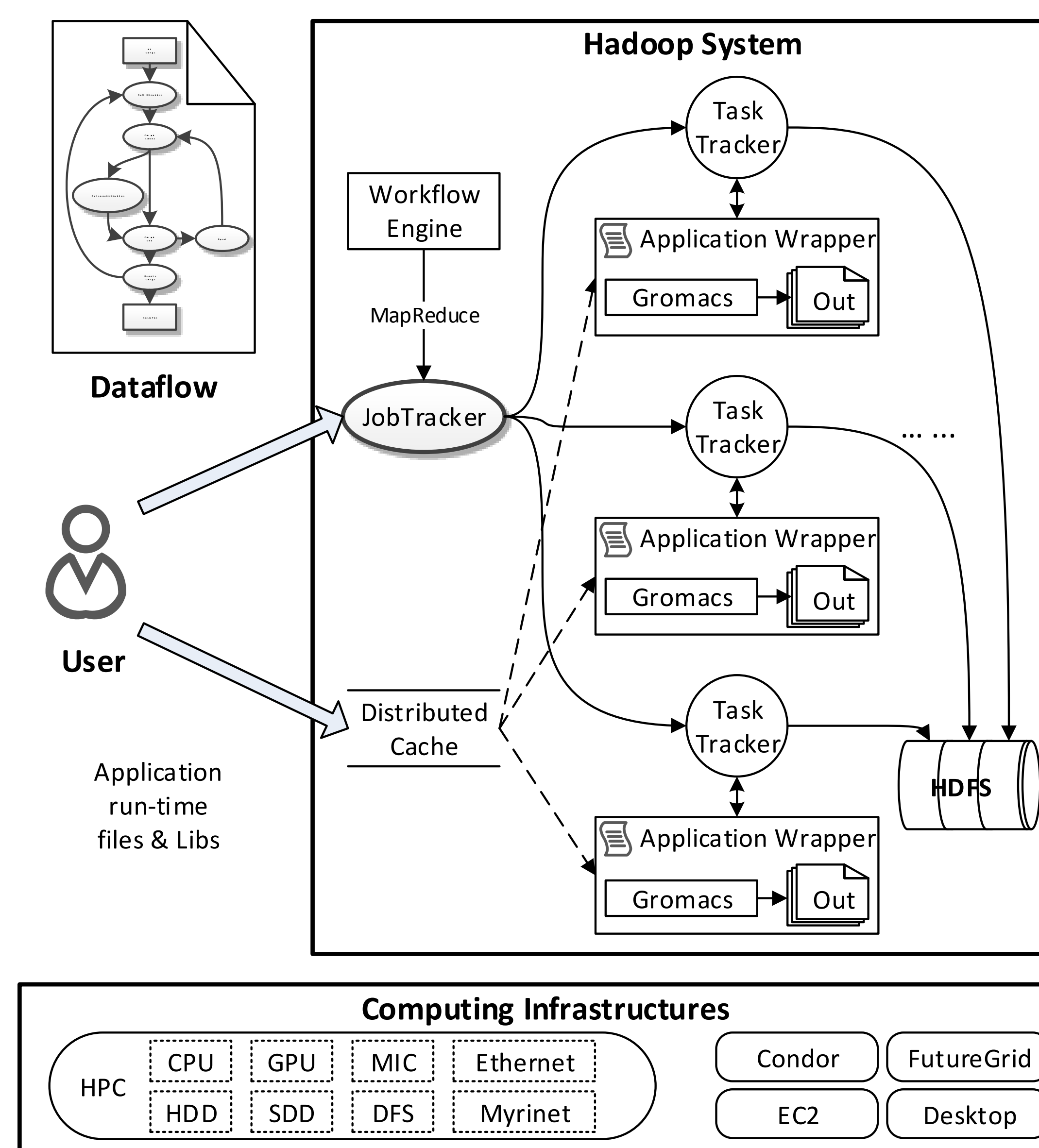


Figure 3. Parallel FFS on various computing infrastructures

RESULT

The data sets in Table 1 are obtained from the FFS implementation by using 64 computing nodes, which lists the cluster size at different interfaces and probabilities of the successful trajectories over all trajectories from current interfaces to next interfaces. With the cluster sizes at interfaces increasing from 70 to 380, the changing of the probability at every of forwarding interfaces can be observed. The original Bash Script-based implementation took more than 3 days to simulate only one branch for each move forward. Our new Hadoop-based implementation decreases simulation time for each move forward with using all branches from days-level to hours-level.

Table 1. FFS simulation result: list of cluster size at different interfaces and probabilities of the successful trajectories over all trajectories from previous interfaces to each different interface with simulation time.

Interface	Cluster Size	Probability	Running Time
$\lambda_A - \lambda_0$	70 - 201	0.010106383	31 min
$\lambda_0 - \lambda_1$	201 - 231	0.006554099	57 min
$\lambda_1 - \lambda_2$	231 - 259	0.005833333	72 min
$\lambda_2 - \lambda_3$	259 - 269	0.016380344	79 min
$\lambda_3 - \lambda_4$	269 - 285	0.005459424	75 min
$\lambda_4 - \lambda_5$	285 - 295	0.010316140	84 min
$\lambda_5 - \lambda_6$	295 - 308	0.006701596	82 min
$\lambda_6 - \lambda_7$	308 - 317	0.009507481	90 min
$\lambda_7 - \lambda_8$	317 - 324	0.009321255	118 min
$\lambda_8 - \lambda_9$	324 - 336	0.003568515	113 min
$\lambda_9 - \lambda_{10}$	336 - 342	0.002876391	113 min
$\lambda_{10} - \lambda_{11}$	342 - 347	0.003836410	110 min
$\lambda_{11} - \lambda_{12}$	347 - 357	0.005967250	115 min
$\lambda_{12} - \lambda_{13}$	357 - 363	0.004715673	118 min
$\lambda_{13} - \lambda_{14}$	363 - 368	0.003007519	132 min
$\lambda_{14} - \lambda_{15}$	368 - 374	0.006819289	170 min
$\lambda_{15} - \lambda_{16}$	374 - 380	0.008266819	110 min

REFERENCES

1. M. A. Rohrdanz, W. Zheng, and C. Clementi, “Discovering Mountain Passes via Torchlight: Methods for the Definition of Reaction Coordinates and Pathways in Complex Macromolecular Reactions,” Annual Review of Physical Chemistry, vol. 64, no. 1, 2013.
2. Allen, R.J., P.B. Warren, and P.R. Ten Wolde, *Sampling rare switching events in biochemical networks*. Physical review letters, 2005. 94(1): p. 18104.
3. Sarupria, S. and P.G. Debenedetti, *Molecular Dynamics Study of Carbon Dioxide Hydrate Dissociation*. The Journal of Physical Chemistry A, 2011.

ACKNOWLEDGEMENT

- Supported by NSF awards CNS-1228312 and CI-1212680CCIT
- OrangeFS Testbed and CITI Palmetto HPC resources