# Rockhopper: a true HPC system with cloud concepts

Richard Knepper, Barbara Hallock, Craig Stewart, Matt Link, Matt Jacobs

rich@iu.edu, bahalloc@iu.edu, stewart@iu.edu, mrlink@iu.edu, mjacobs@penguincomputing.com

## ABSTRACT

A number of services for scientific computing based on cloud resources have recently drawn significant attention in both research and infrastructure provider communities. Most cloud resources currently available lack true high performance characteristics, such as high-speed interconnects or storage. Researchers studying cloud systems have pointed out that many cloud services do not provide service level agreements that may meet the needs of the research community. Furthermore, the lack of location information provided to the user and the shared nature of the systems use may create risk for users of the system, in the instance that their data is moved to an unknown location with an unknown level of security.

Indiana University and Penguin Computing have partnered to create a system, Rockhopper, which addresses many of these issues. This system is a true high performance resource, with on-demand allocations and control and tracking of jobs, situated at Indiana University's high-security datacenter facility. Rockhopper allows researchers to flexibly conduct their work under a number of use cases while also serving as an extension of cyberinfrastructure that scales from the researcher's local environment all the way up through large national resources.

We describe the architecture and ideas behind the creation of the system, present a use case for campus bridging, and provide a typical example of system usage. In a comparison of Rockhopper to a cloud-based system, we run the Trinity RNA-seq software against a number of datasets on both the Rockhopper system and on Amazon's EC2 service.

## ARCHITECTURE GOALS

The Rockhopper system is the result of a joint effort on the part of Indiana University and Penguin Computing, with the goal of providing true HPC capabilities to users at federally-funded research and development centers (FFRDC's) and educational institutions. Indiana University's Pervasive Technology Institute was charged with providing a flexible system for research computational needs. This system would need to allow researchers to get access to accounts and compute resources quickly, while providing the same security as IU's other HIPAA-compliant resources, with an environment that would be similar enough to IU's existing systems that transition costs to the new system would be low.

Rockhopper provides cycles for researchers at FFRDC's and educational institutions at a price competitive with cloud offerings such as Amazon EC2 or Windows Azure. The system is specifically oriented at users who need quick access to secure, flexible resources, such as: users with grant money that would otherwise be spent on a locally-administrated cluster; users who are in need of quick processing capability in order to meet paper deadlines; and users with overspent allocations on national resources who need to continue analyses. Effort was also made to ensure that owners of data would retain sole rights to the information they put on the system.

## HYBRID MODEL FOR CLOUD-LIKE HPC

Rockhopper constitutes a significant resource that creates a hybrid HPC/Cloud system. The hardware, scheduling, storage, interconnect, and software-as well as the support at IU and Penguin-are that of a traditional HPC resource, while the ability to receive access to the hardware quickly and and the ease of integrating into different types of workflow more closely resemble cloud offerings. This hybrid system makes it possible to conduct significant analyses and support larger workflows with the only condition of use being that the researcher at an FFRDC or educational institution.

In order to realize the full potential of this hybrid model, Rockhopper will continue to work best when it is possible to integrate with other resources and with a broad set of commonly-used research software, in order to meet the needs of the most researchers. Rockhopper also can take advantage of software for data management, including the Global Federated Filesystem (GFFS) component of Genesis II software and Globus Online. These technologies facilitate the use of the Rockhopper resource with data at researchers' institutions without resorting to cumbersome transfer mechanisms.

The ability of Rockhopper to participate in analyses across multiple systems and for users to easily make use of the Rockhopper cyberinfrastructure fulfills campus bridging efforts pursued by Indiana University and the XSEDE project, which seeks to make computational resources appear and work as if they were proximal to the researcher making use of them. IU and Penguin Computing are pursuing additional changes to improve Rockhopper's ability to engage in campus bridging efforts.

## AMAZON EC2 INSTANCE DETAILS

| Instance Type | vCPU | Memory | Storage | Price Per Hour |
|---|---|---|---|---|
| m2.4xlarge | 8 | 68.4GB | 2 x 840GB | $1.64 |
| cr1.8xlarge | 32 | 244GB | 2 x 120 GB (SSD) | $2.40 |

## ROCKHOPPER SYSTEM CONFIGURATION

| System Configuration | Aggregate Information | Per Node Information |
|---|---|---|
| Machine Type | High-performance computing Usage on demand Penguin Computing 1804 MPP cluster | 4x2.1 GHz 12-core AMD Opteron Altus 6172 processors |
| Operating System | CentOS 5 | |
| Memory Model | Distributed | |
| Processor Cores | 528 | 48 |
| Memory | 1.4TB | 128GB |
| Processing Capability | 5242 gigaflops | 403 gigaflops |

## USE CASE: TRINITY RNA-SEQ ASSEMBLIES ON ROCKHOPPER AND AMAZON EC2

The Trinity software for de novo reconstruction of transcriptomes from RNA-seq data3 is used to process large volumes of RNA sequence data. Trinity is composed of three separate programs which perform specific tasks, and runs best when using strand-specific data. This bioinformatics software can utilize a significant amount of memory per core and multiple cores on a single system. This software is installed on Rockhopper and another resource at IU dedicated to genome assembly, the Mason cluster 4. The standard use case for utilization of this software is to capture RNA sequences and process single large files in serial on a single node. In order to gauge the capabilities of Rockhopper for this task, staff at IU ran a typical set of Trinity analyses on the Rockhopper system and measured the number of core hours required by the analyses. Then the same datasets were transferred to a comparable Amazon EC2 Instance and analyzed with Trinity, recording the time required to complete the analysis.

## RESULTS COMPARISON: TRINITY RUNS ON ROCKHOPPER VS. EC2

| Data Set | Total Sequence Length | Core Hours | Penguin Pricing |
|---|---|---|---|
| 500MB | 19869875 | 66.44 | $6.03 |
| 1GB | 28555620 | 100.53 | $9.15 |
| 3GB | 48647454 | 245.53 | $22.40 |
| 5GB | 60051778 | 366.83 | $33.51 |
| 7GB | 69089235 | 488.48 | $44.66 |

Base costs for Rockhopper are $0.09/core hour and $0.10/fractional GB-month. These costs are calculated based on the assumption that the researcher would store data for one full month during analysis.

Base costs per Trinity run on EC2 do not include the cost of time setting up Trinity to run on the EC2 instance. As Amazon bills EC2 usage based on how long the instance is active, costs for analysis may end up being much higher than indicated, particularly if multiple runs must be completed.

| Data Set | Total Sequence Length | Elapsed Time | Amazon Pricing |
|---|---|---|---|
| 500MB | 18466805 | 1:44:02 | $2.84 |
| 1GB | 24264084 | 2:40:48 | $4.40 |
| 3GB | 34240171 | 5:33:16 | $9.11 |
| 5GB | 37885510 | 8:39:34 | $14.20 |
| 7GB | 41805558 | 10:25:19 | $17.09 |

## A NOTE ON EC2 INSTANCES

Amazon bills EC2 usage based on the amount of time the instance is open, rather than active; the above-indicated pricing does not include idle time or the time incurred setting up Trinity prior to running the analysis.

## CONCLUSIONS

Both cloud systems and traditional HPC systems have different advantages and disadvantages in support of scientific research. The Rockhopper system was created in order to alleviate concerns about security, data ownership, and location of analyses, while providing a true cluster-on-demand service that is as flexible as common cloud offerings, at a competitive price to commercial offerings. This system is capable of participating in complex workflows and working with portal and science gateways.

In comparison to Amazon EC2 as a resource for carrying out scientific analyses, Rockhopper provides a close competitor to EC2, representing a less expensive option for running analyses in a flexible and simple fashion. If an Amazon EC2 instance is configured and managed very vigilantly, it is possible to compete on price based on the analyses that we ran, but this will require a significant amount of preparation and management in order to realize cost savings compared to Rockhopper.

## Supported by:

NSF · IU · PENGUIN COMPUTING · XSEDE Extreme Science and Engineering Discovery Environment