

BIOCLOUD: Resource Provisioning Framework for Bioinformatics Applications in a Multi-Cloud Environment

Izzet F. Senturk¹, Bala Krishnan², Kamer Kaya^{1,3}, Qutaibah Malluhi², Umit V. Catalyurek¹

¹Dept. Biomedical Informatics, The Ohio State University, Columbus, OH, 43230

²KINDI Center for Computing Research, Qatar University, Doha, Qatar

³Computer Science & Engineering, Sabanci University, Istanbul, Turkey

1- Introduction

Next Generation Sequencing (NGS) has paved the way to extract genetic information from biological systems in a massively parallel fashion. It also made the technology pervasive by decreasing the sequencing cost significantly. The availability of NGS in a wider scale due to decreased sequencing costs has shifted the challenge from obtaining sequencing data to analyzing it which requires addressing scalability issues. Especially for small-scale labs, the lack of dedicated hardware that is required to complete the analysis within a reasonable amount of time drives the efforts to cloud computing.

2- Motivation

Even though the Cloud helps the researchers to avoid the high ownership and maintenance costs of such dedicated hardware, the complexity of configuring the analysis services that can dynamically scale on-demand can be a barrier for non-computer savvy bio-researchers. BIOCLOUD provides a ready-to-use and efficient analysis environment to the bio-researchers so that they can focus solely on the analysis, not on the complexities of the computational environment.

BIOCLOUD can handle peak loads by provisioning additional resources and configuring them in such a way that the resource utilization can be maximized. The core competency of BIOCLOUD is the ability of managing resources from multiple cloud providers and local clusters so that these resources are utilized simultaneously in a seamless and efficient manner. There are projects that provide Cloud support for NGS applications, such as Globus Genomics [1], but they are designed to work with single Cloud provider.

Cloud providers tend to develop custom APIs which complicates the interoperability of the tools to access their services. These custom APIs also make it difficult to migrate from one cloud provider to another or to use multiple providers simultaneously. BIOCLOUD saves the researchers the trouble and enables a hybrid multi-cloud model without any interoperability complications.

3- The BIOCLoud Framework

BIOCLOUD provides a single entry point to conduct bioinformatics analysis on a multi-cloud environment. BIOCLOUD leverages DeltaCloud [2] for resource provisioning and Galaxy [3] for abstract workflow

development and execution. Through its user-friendly web-interface, researchers can define their resources and specify their time and budget constraints. Based on this information, BIOCLoud designates a schedule for the submitted workflow and determines the target resources for execution. It also ensures the resource availability before executing a particular workflow step. With the features it provides, BIOCLoud can schedule different workflow steps on different cloud environments while scaling the resources on-demand.

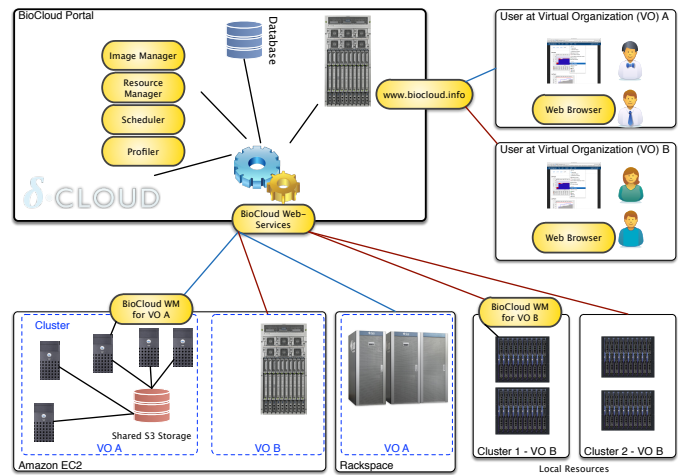


Fig. 1: Layered architecture of BIOCLoud

BIOCLOUD (Fig. 1) consists of two main components: BIOCLoud PORTAL and BIOCLoud WORKFLOW MANAGERS (WM). The portal is BIOCLoud's single access point of entry for all users (virtual organizations) where they can manage their resources. On the other hand, WMs are the management components of the system specific to a particular virtual organization (VO).

BIOCLOUD offers a loosely coupled architecture in which some of the decisions (e.g., when to dispatch a workflow step and where to run this step) are delegated to BIOCLoud PORTAL so that the scheduling logic is separated from the core workflow system. This provides the flexibility of updating the scheduling algorithm without requiring a

