



OpenStack at NDSU - Design Document

Justin Anderson
John Nagel
Cesar Ramirez
John Rensberger

Table of Contents

[Table of Contents](#)

[1. Overview](#)

[1.1 Introduction](#)

[1.2 Scope](#)

[2. High-Level Design](#)

[2.1 High-Level Component Design](#)

[2.2 Required Services](#)

[2.3 Instance Requirements](#)

[2.4 User Deployment](#)

[3. Non-Functional Design](#)

[4. Definitions and Acronyms](#)

[5. Document Review](#)

1. Overview

1.1 Introduction

John Deere and RDO are collecting beets harvest data. Dr. Ludwig wishes to store this data in a distributed file system for highly availability. Dr. Denton wants to later analyze and process this data using a Hadoop cluster. More data will be received on a weekly basis and it needs to be added to the storage system. The system administrator needs to receive reports of the current state of the system while Hadoop is running.

An important component of this project is being able to deploy a Hadoop cluster within OpenStack. Hadoop is an open source software framework for running applications that scale over large clusters and huge amounts of data. The main computational paradigm of Hadoop is MapReduce, which will break (map) bigger problems into smaller tasks that get distributed and then the result gets reassembled (reduce).

1.2 Scope

This project will use OpenStack to deploy a Hadoop cluster. This cluster is going to be easy to deploy and destroy once finished. Storage should also be capable of holding a multitude of terabytes and it should be able to be expanded as needed. The requirements of the Hadoop cluster should also be easy to configure via Heat (or any orchestration service).

2. High-Level Design

2.1 High-Level Component Design

Component	Description
Glance - Image Service	Glance provides a repository of images that can be retrieved in order to create new VMs.
Heat - Orchestration	Heat is a service to orchestrate multiple composite cloud applications using the AWS CloudFormation template format, through both an OpenStack-native REST API and a CloudFormation-compatible Query API.

2.2 Required Services

Component	Related Requirements
All	Gigabit wired networking between all components
Nova	Hardware capable of virtualization
Glance	Storage for images with Hadoop software installed and configured
Horizon	Web front-end for managing OpenStack
Quantum	Hadoop nodes will run in an isolated network
Cinder	Sufficient storage for VMs and data
Oz	Oz is a set classes and scripts that allow Heat to do automatic installations of various guest operations.
libguestfs	Is a set of tools used by Oz to access and modify VM disk images

2.3 Instance Requirements

- Linux images with Hadoop installed and configured
 - At the time of build this will be the latest stable version
 - This will run 1 master node and n slave nodes.
 - The master node runs the NameNode and JobTracker daemons
 - The slave nodes which run the DataNode and TaskTracker daemons
- Scripts to automate the deployment of Hadoop clusters using euca2ools
 - These will take care of provisioning instances and any network setup that needs to be done
 - No manual configuration or web-based provisioning should be required
 - These scripts will be modeled after existing EC2 deployment scripts available in the current Hadoop distribution

2.4 User Deployment

1. The user should be able to launch a Hadoop cluster using the provided scripts, or, ideally, using a template provided for OpenStack Heat.
2. The user should be able to launch Hadoop jobs from the master node using SSH.

3. Non-Functional Design

- Must use the CS department's Active Directory for authentication
- Permissions must be given on a need-to-use basis
- Data must be stored in a distributed and redundant storage
- All nodes within a single cluster should run in an isolated network
- Data should be able to be added directly and with ease
- Must provide documentation for future use and maintenance

4. Definitions and Acronyms

Term	Definition
Distributed file system	Any file system that allows access to files from multiple hosts sharing via a computer network.
Redundancy	Data will be written in more than one place. This allows for contingency in case of disk failure. Redundancy can be done at the level of bytes, files or filesystems.
Orchestration	Automated arrangement, coordination, and management of complex computer systems, middleware, and services.
Cluster	A set of connected computers that work together so that in many respects they can be viewed as a single system.
JEOS	JEOS is an acronym for Just Enough Operating System. It is an operating system that contains the bare, minimal components to run.
euca2ools	Euca2ools are command line tools for interacting with Amazon Web Services (AWS) and other AWS-compatible web services, such as Eucalyptus and OpenStack.

5. Document Review

Project Manager

Name:	John Rensberger
Date:	
Signature:	

Mentor

Name:	Michael Fork, Lance Bragstad, Mathew Odden, Adam Reznechek
Date:	
Signature:	