# Aeromancer: A Workflow Manager for Large-Scale MapReduce-Based Scientific Workflows

Presented by Sarunya Pumma
Supervisors: Dr. Wu-chun Feng, Dr. Mark Gardner, and Dr. Hao Wang

# Outline

- Introduction & Motivation
  - Cloud computing and its benefits
  - SeqInCloud and scalability analysis
- Aeromancer
  - Usage scenario
  - Workflow and architecture
  - Discussion
- Experiments & Initial Results

# Introduction & Motivation

What is Cloud?

# What is Cloud?

# How is Cloud beneficial for Us?

# How is Cloud beneficial for Us?



Next-generation sequencing (NGS) produces a huge amount of DNA sequence data

{ 2k sequencers produce 15 PB genetic data/year }

Cost of sequencing a human-size genome decreases over time

{ 
2011 → $95M
2012 → $6.5k
2014 → $1k
}

# How is Cloud beneficial for Us?

# Big Data

# How is Cloud beneficial for Us?

Cloud can accelerate our ability
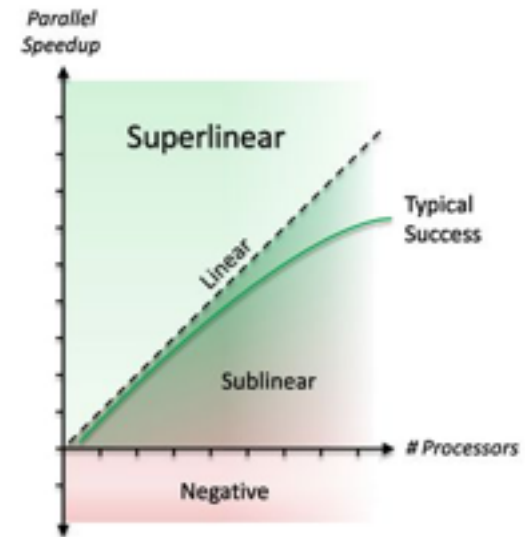to identify the cure for cancer

# Big Data Challenge

Provisioning compute and storage resources to analyze the exponentially increasing genomic data
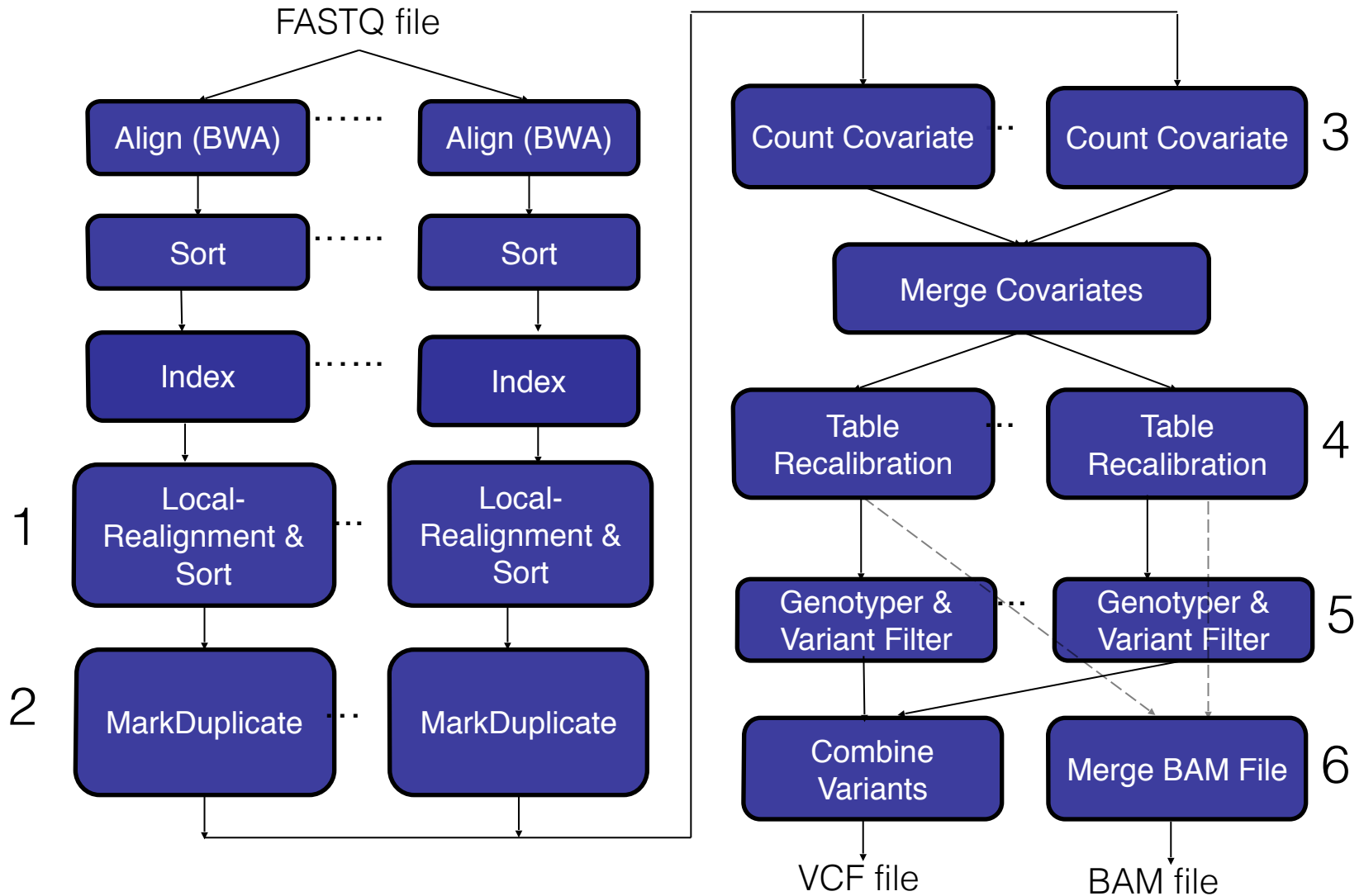
# SeqInCloud

- It is implemented based on GATK (Genome Analysis Toolkit)

- All six stages have to be executed in order

- Each stage is implemented as a small separated Hadoop program

  - It can be broken down into small jobs and run on multiple computers/processors

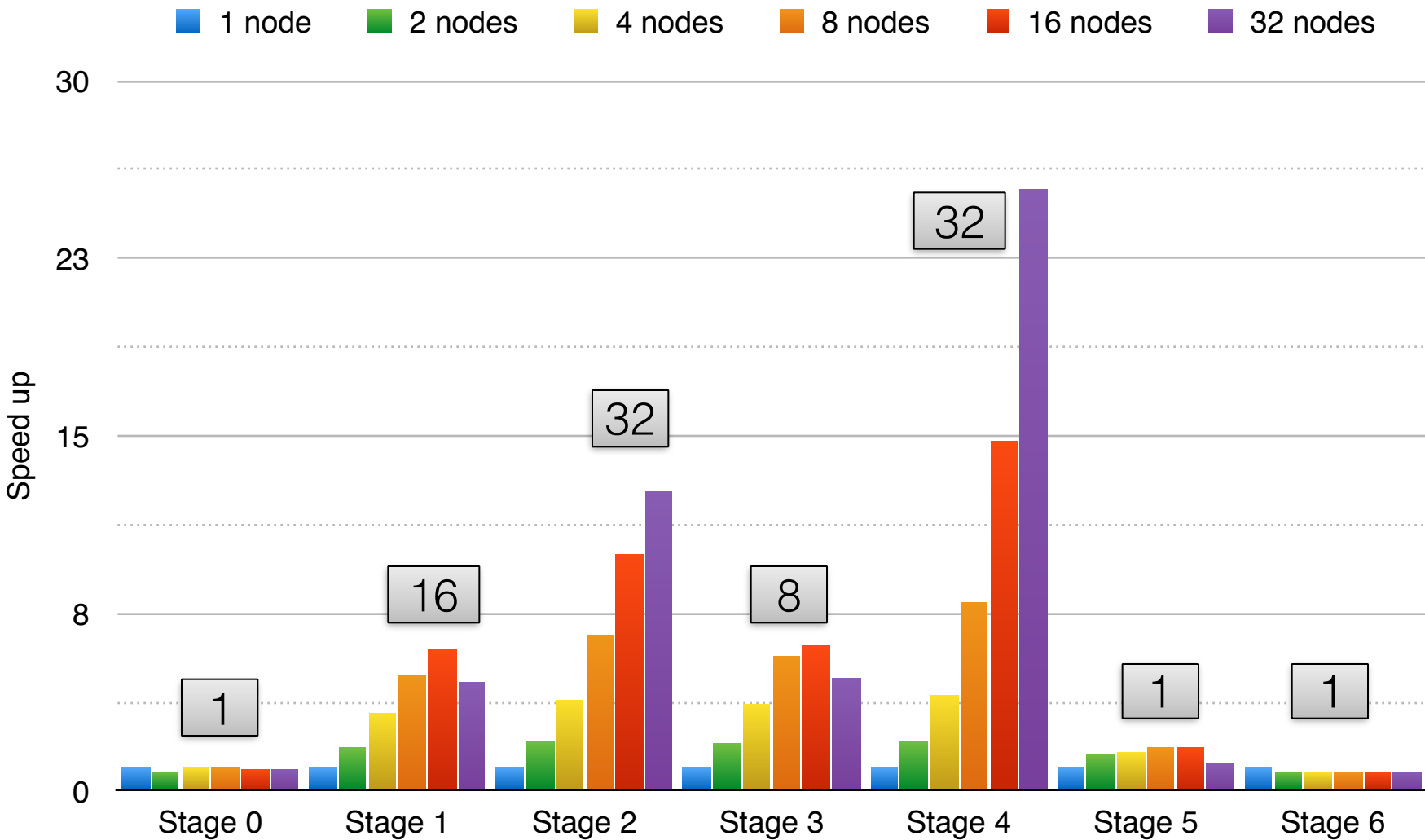  - We expect a **linear speed-up** on every stage
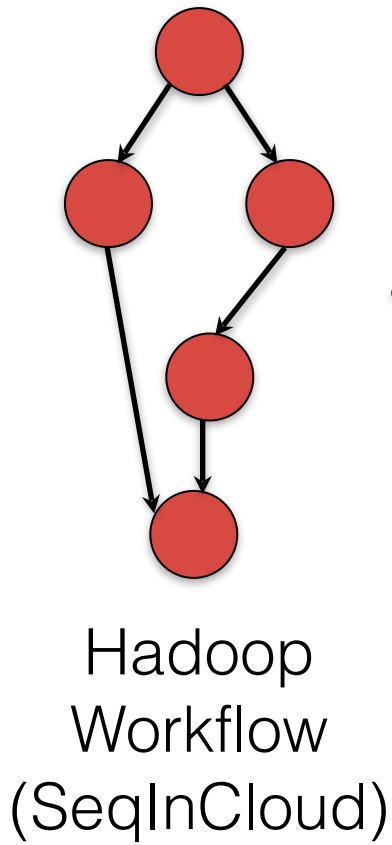
# SeqInCloud

## Genome variant analysis pipeline

# SeqInCloud Scalability Analysis

- To run SeqInCloud efficiently on Hadoop

    - Number of compute nodes must be varied based on the scalability of each stage

- Challenges

    1. How many nodes do we need for each stage?

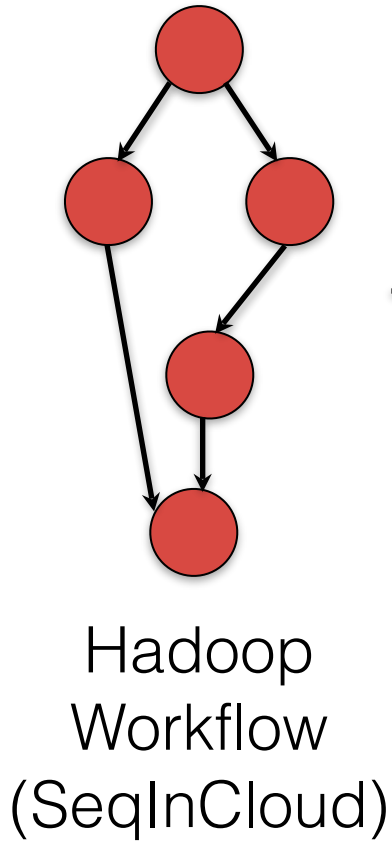    2. How to rapidly adjust a number of nodes with negligible overhead?

An automatic Hadoop workflow execution is needed!

Hadoop
Workflow
(SeqInCloud)

Hadoop
Workflow
Manager

Cloud
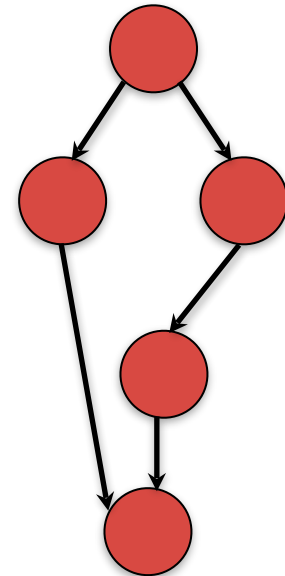
Hadoop
Workflow
(SeqInCloud)

Aeromancer

Cloud

# Aeromancer

- Platform for executing the scientific Hadoop workflows
    - Based on Hadoop
        - MapReduce
        - HDFS
    - Build on top of CloudGene
    - Workflow (DAG)
        - Node/Stage
        - Edge (represent dependency between two stages)

# Aeromancer Usage Scenario

1

Add a new application to Aeromancer

# Aeromancer Usage Scenario



YAML file

# Aeromancer Usage Scenario

**1** Add a new application to Aeromancer

**2** Submit a new job

# Aeromancer Usage Scenario



User Portal

# Aeromancer Usage Scenario

1 Add a new application to Aeromancer

2 Submit a new job

3 Wait for the execution to complete
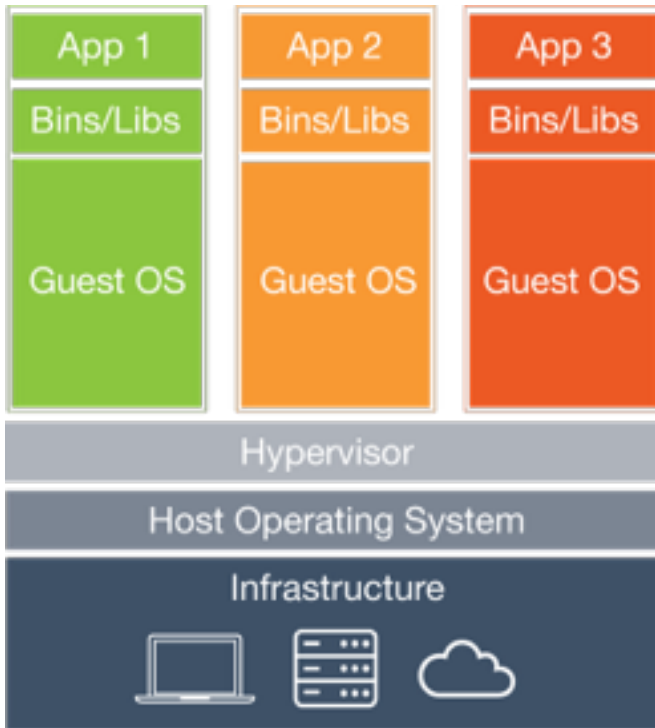
# Aeromancer Workflow & Architecture
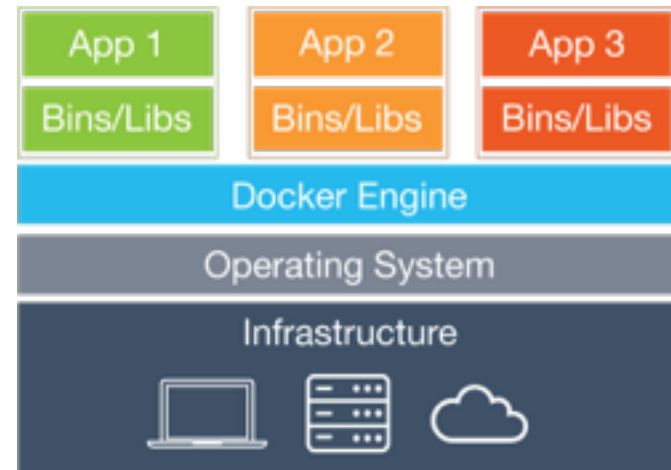
# Aeromancer Discussion

- Currently, Aeromancer worker nodes are **virtual machines** in the cloud

- Pro

  - Number of worker nodes are flexibly adjustable upon demand

- Con

  - The overhead of creating/deleting worker nodes is high

  - Execution speed drops significantly

**Containerization** may be able to
improve Aeromancer's performance

# Virtualization

# Containerization



| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Libs | Bins/Libs | Bins/Libs |
| Guest OS | Guest OS | Guest OS |
| Hypervisor | | |
| Host Operating System | | |
| Infrastructure | | |

| App 1 | App 2 | App 3 |
|-------|-------|-------|
| Bins/Libs | Bins/Libs | Bins/Libs |
| Docker Engine | | |
| Operating System | | |
| Infrastructure | | |

Reference: https://www.docker.com/whatisdocker
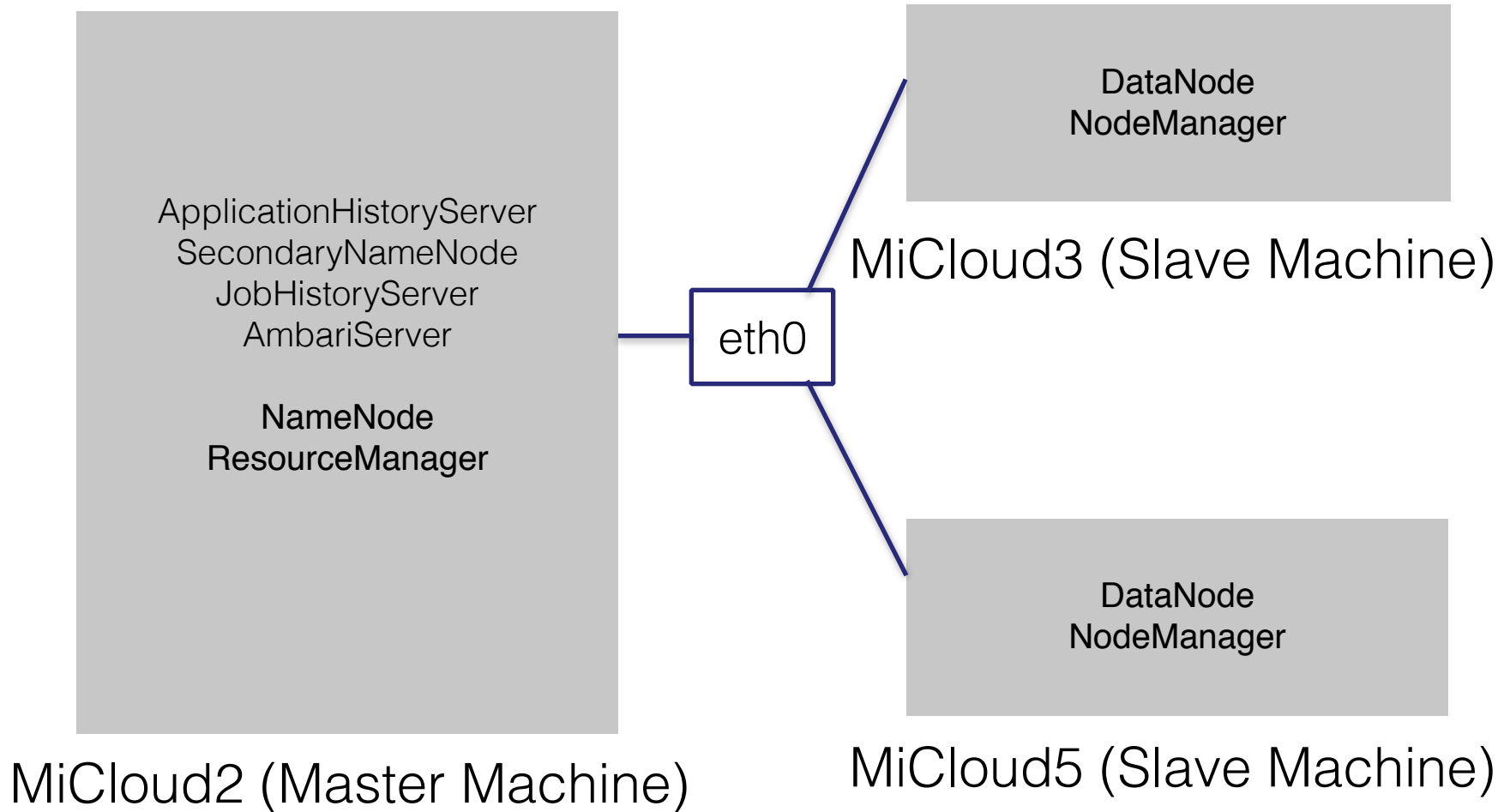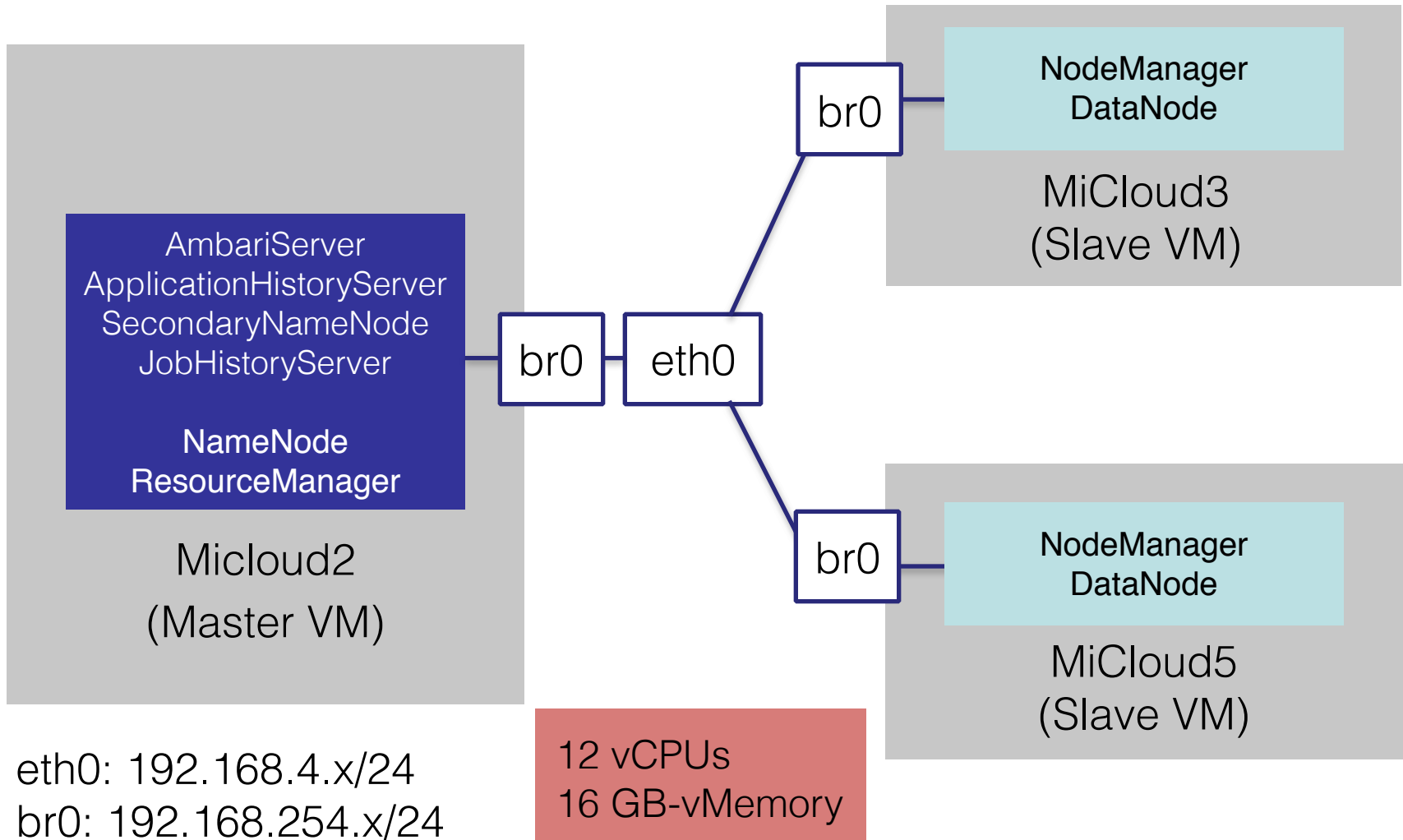
# Experiment

- Platforms
    - Bare metal machines
    - KVM
    - Docker
- Benchmarks (HiBench)
    - Sort
    - WordCount
    - TeraSort
- Performance
    - Speed (execution time)
    - HDFS throughput

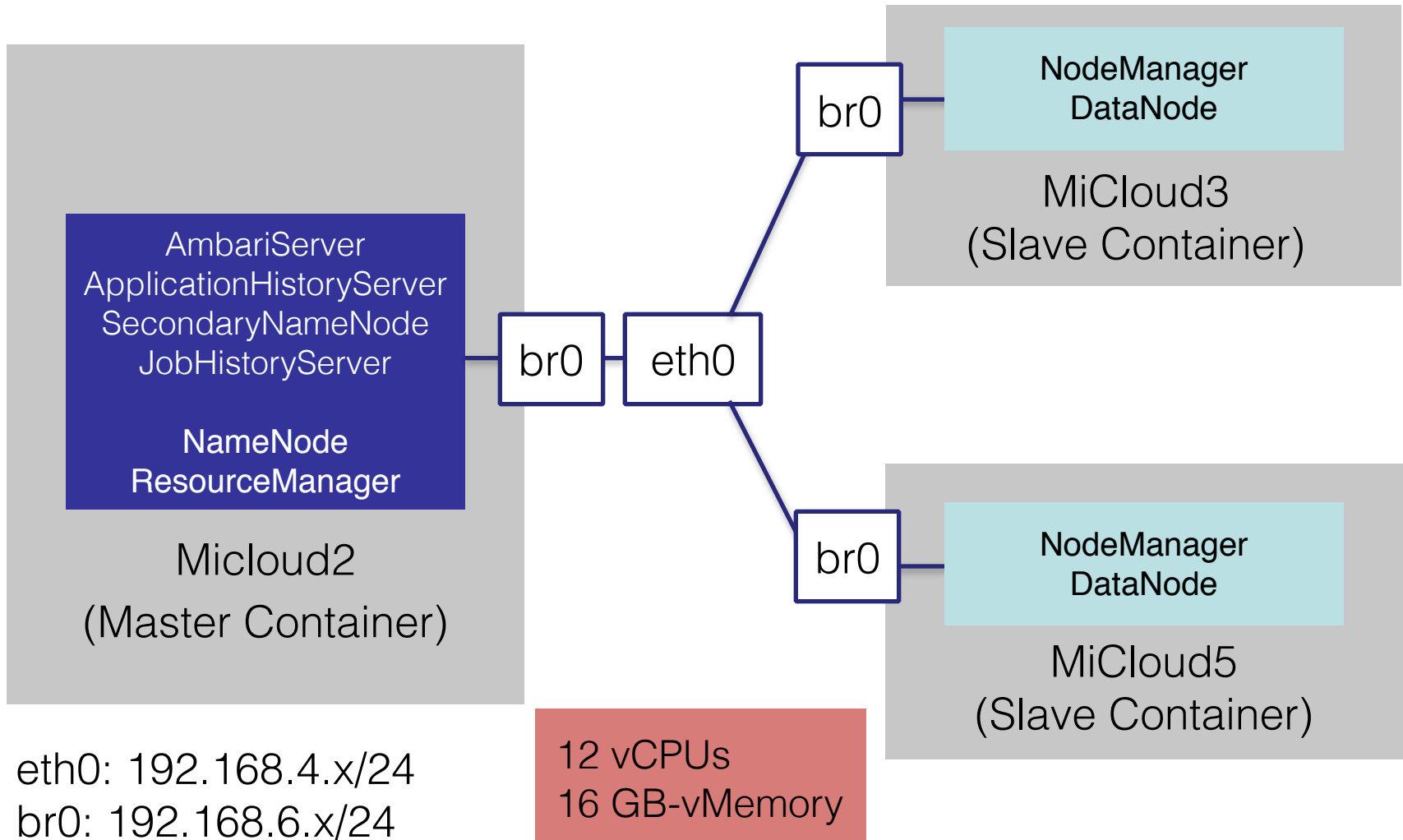- Nutch Indexing
- PageRank
- EnhancedDFSIO

# Experiment:
# Hadoop on Bare Metal Environment

# Experiment:
# Hadoop on Virtual Machine Environment

# Experiment:
# Hadoop on Docker

# Initial Results: Slow Down



Legend: KVM (blue), Docker (green)

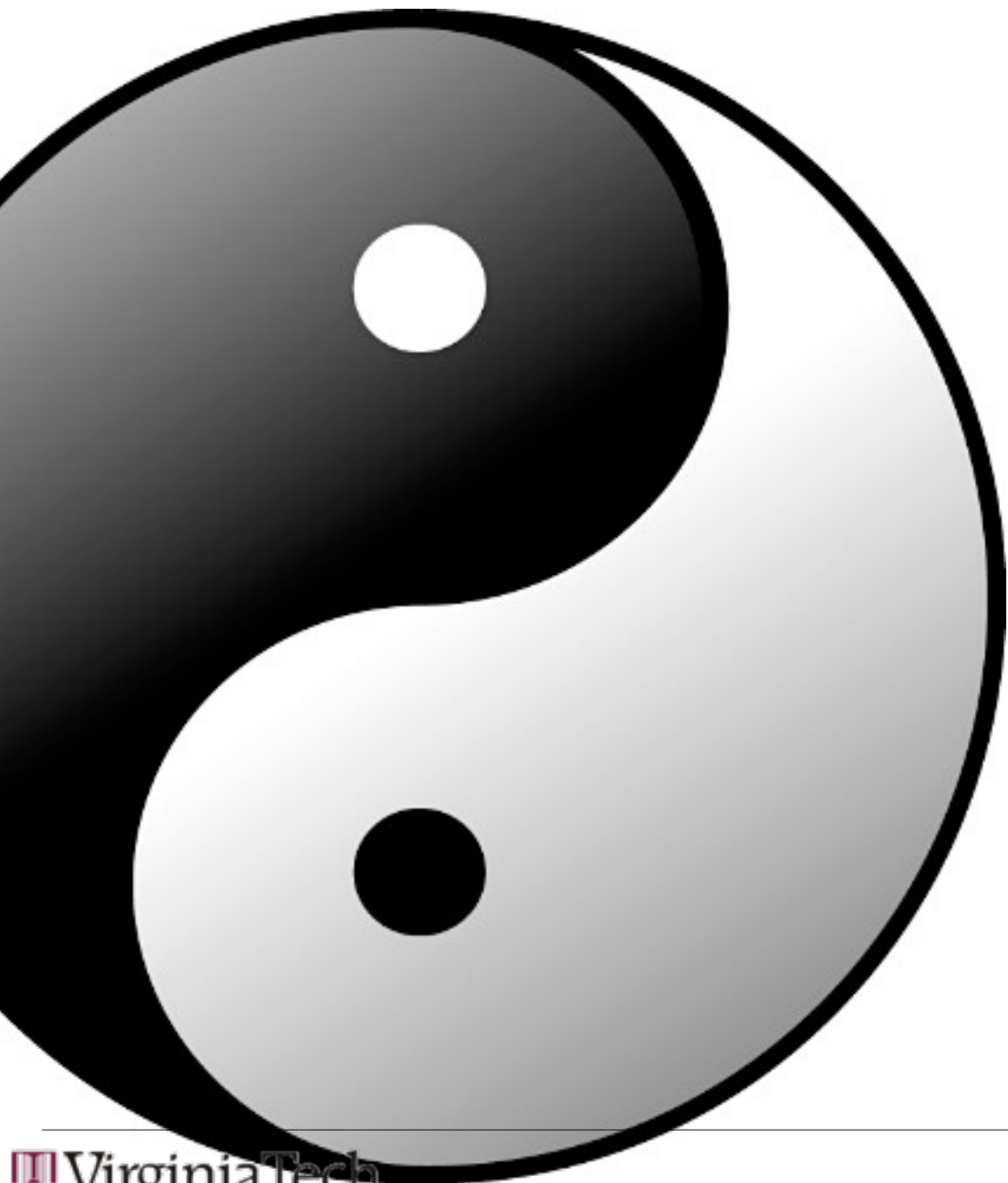| Benchmark | KVM | Docker |
|---|---|---|
| Sort | 62% | 95.4% |
| Wordcount | 92% | 98.6% |
| DFSIOE | 14% | 59.5% |
| NutchIndexing | 88% | 97.4% |
| TeraSort | 9% | 56.8% |
| PageRank | 83% | 94.8% |

# Future Work

- Automatic resources provisioning

  - Determining the number of computing nodes needed for each stage

    - Based on the scalability of each stage

    - Based on user's constraints (budget/expected response time)

- Automatic task scheduling

  - Determining the optimal execution location for each stage (either client or cloud)

Thank you

Virginia Tech
*Invent the Future*

SyNeRG
synergy.cs.vt.edu