# TACHYON

# A Reliable Memory-Centric Distributed Storage System

**Haoyuan (HY) Li**, Ali Ghodsi, Matei Zaharia, Scott Shenker, Ion Stoica

September, 2014, ISTC @ Portland

http://www.istc-cc.cmu.edu/

amplab
UC Berkeley

Intel Science & Technology
Center for Cloud Computing
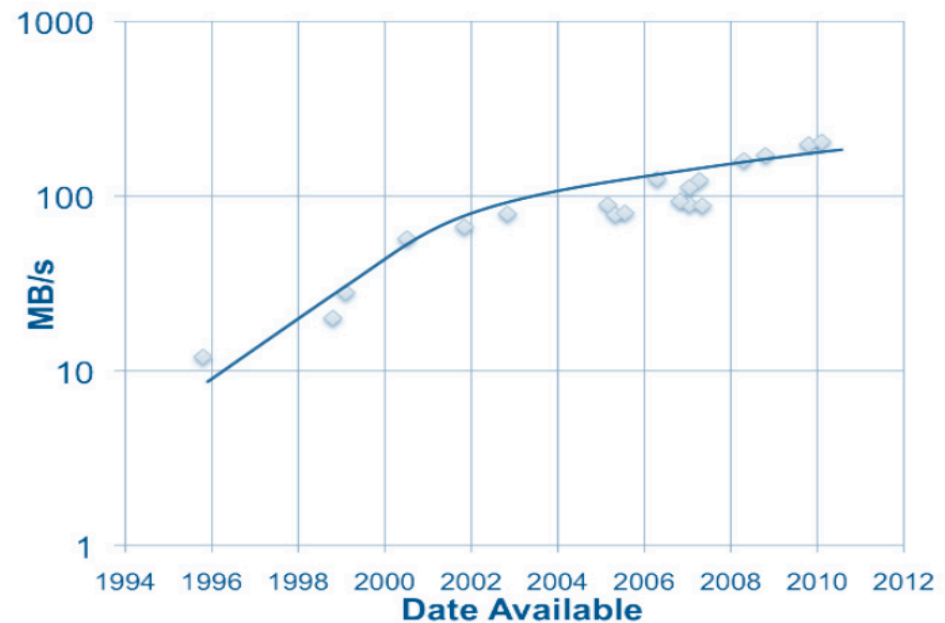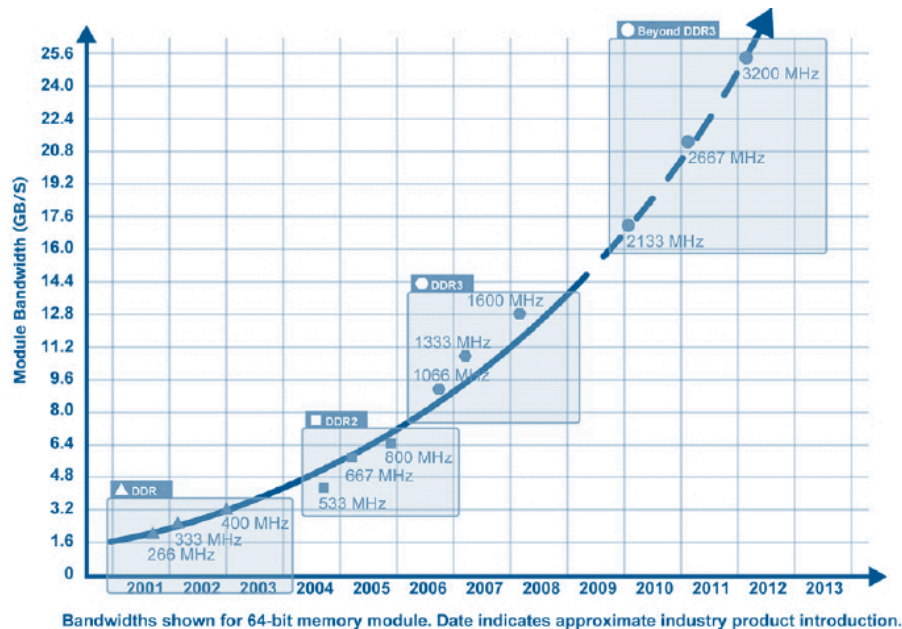
# Outline

- Overview
  - Feature 1: Memory Centric Storage Architecture
  - Feature 2: Lineage in Storage


- Challenges


- Open Source


- Future

# Outline

- **Overview**
  - Feature 1: Memory Centric Storage Architecture
  - Feature 2: Lineage in Storage


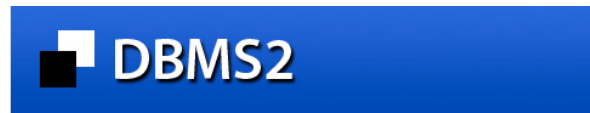- Challenges


- Open Source


- Future

# Memory is King

- RAM throughput increasing exponentially
- Disk throughput increasing slowly



Bandwidths shown for 64-bit memory module. Date indicates approximate industry product introduction.

Memory-locality key to interactive response time

# Realized by many…

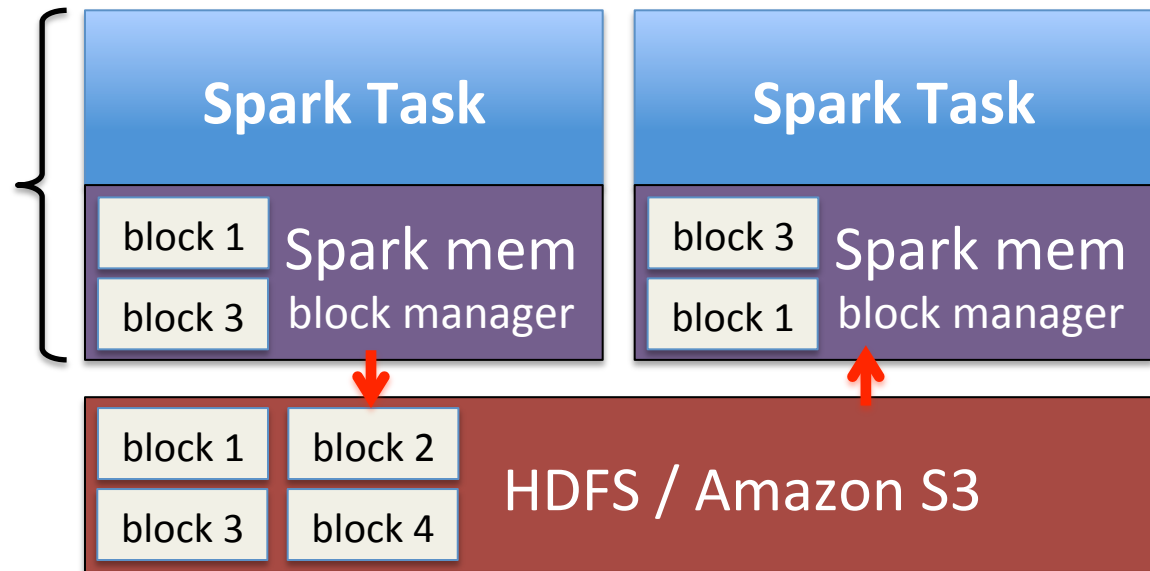- Frameworks already leverage memory

# Problem solved?

# An Example: Spark

- Fast in-memory data processing framework
  - Keep one in-memory copy inside JVM
  - Track lineage of operations used to derive data
  - Upon failure, use lineage to recompute data

Lineage Tracking

# Issue 1

## *Data Sharing is the bottleneck in analytics pipeline: Slow writes to disk*
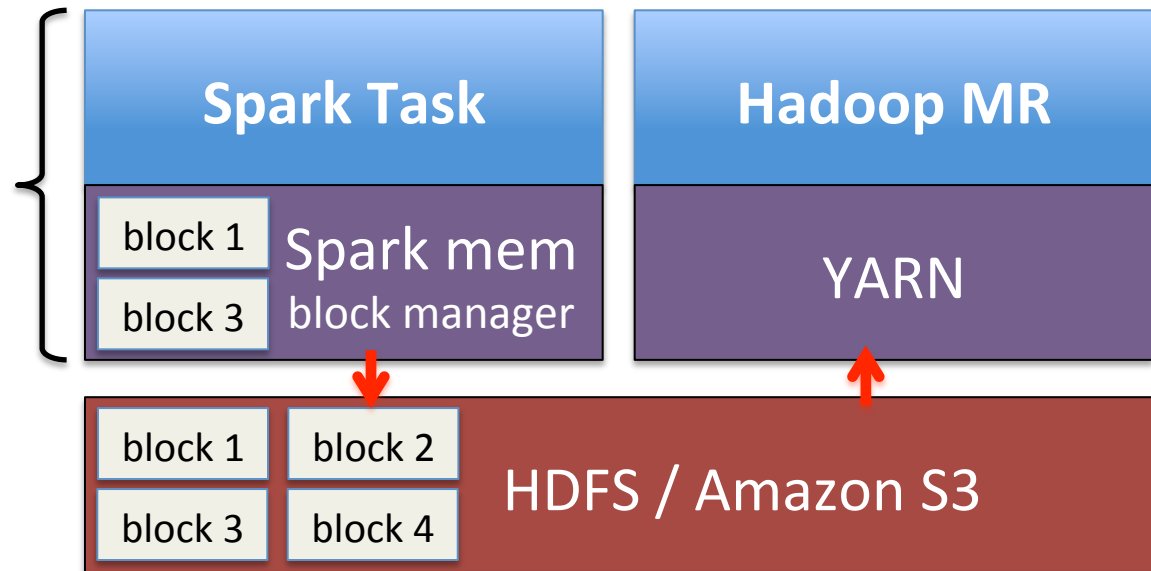
storage engine & execution engine same process (slow writes)

# Issue 1

## *Data Sharing is the bottleneck in analytics pipeline: Slow writes to disk*

storage engine &
execution engine
same process
(slow writes)

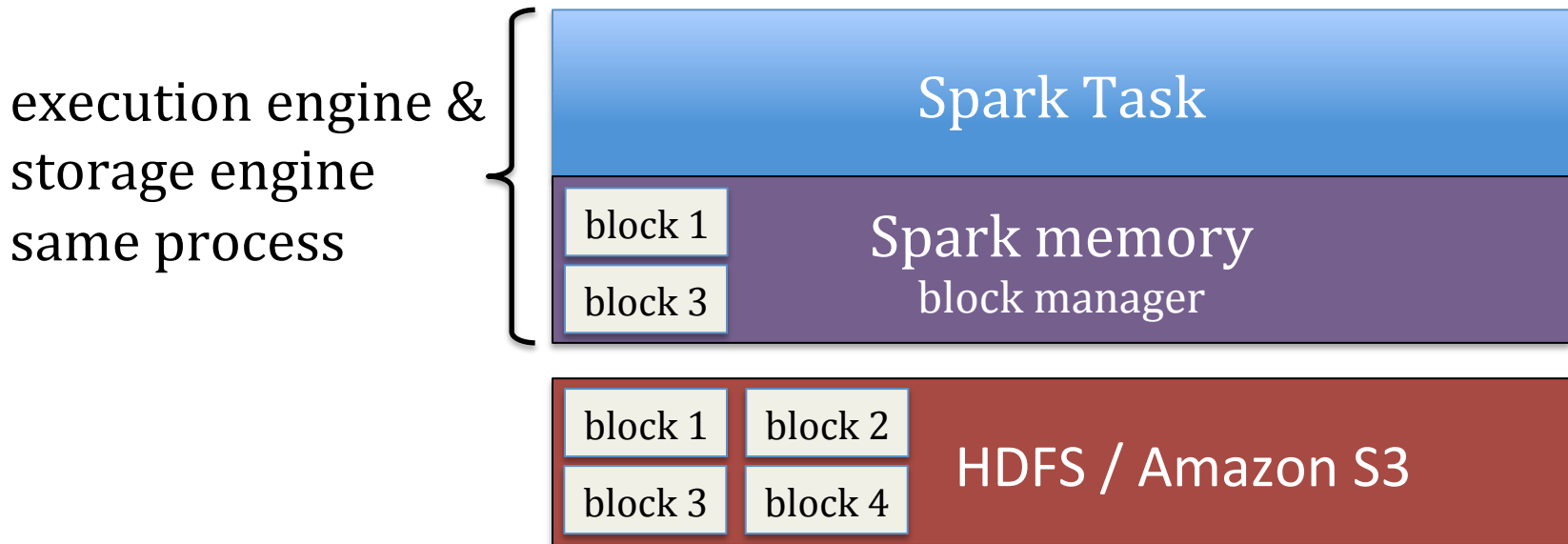| Spark Task | Hadoop MR |
|---|---|
| block 1 / block 3  Spark mem block manager | YARN |

| | HDFS / Amazon S3 |
|---|---|
| block 1 / block 2 / block 3 / block 4 | |

# Issue 2

## *Cache loss when process crashes.*

execution engine &
storage engine
same process

| Spark Task |
|---|
| block 1, block 3 — Spark memory block manager |

| block 1, block 2, block 3, block 4 — HDFS / Amazon S3 |

# Issue 2

## *Cache loss when process crashes.*

execution engine &
storage engine
same process

| crash |
| --- |

block 1
block 3 **Spark memory** block manager

block 1  block 2
block 3  block 4 HDFS / Amazon S3

# Issue 2

## *Cache loss when process crashes.*

execution engine &
storage engine
same process

crash

| block 1 | block 2 | HDFS / Amazon S3 |
|---------|---------|------------------|
| block 3 | block 4 | |

# Issue 3

## *In-memory Data Duplication &*
## *Java Garbage Collection*

execution engine &
storage engine
same process
(duplication & GC)

| Spark Task | Spark Task |
|---|---|
| block 1 block 3 Spark mem block manager | block 3 block 1 Spark mem block manager |

| | |
|---|---|
| block 1 block 2 block 3 block 4 | HDFS / Amazon S3 |

# Tachyon

***Reliable*** data sharing at **<span style="color:red">*memory-speed*</span>**
***within and across*** cluster frameworks/jobs

# Solution Overview

**Basic idea**

- Feature 1: <span style="color:red">memory-centric</span> storage architecture

- Feature 2: push <span style="color:red">lineage</span> down to storage layer

**Facts**

- One data copy in memory
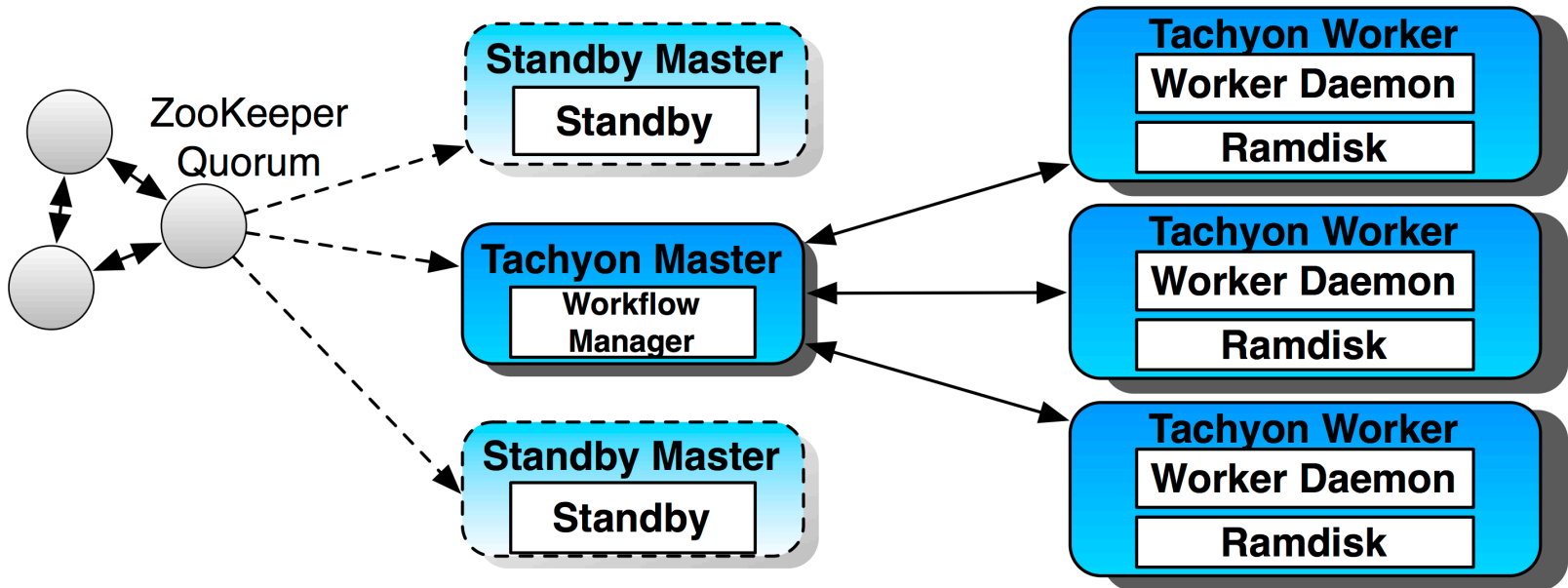
- Recomputation for fault-tolerance

# Stack

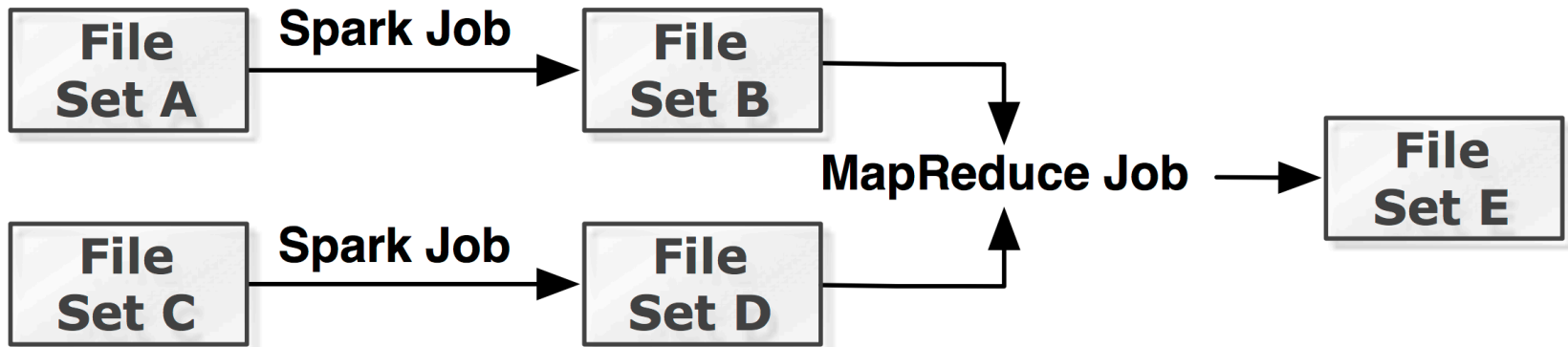Computation Frameworks
(Spark, MapReduce, Impala, H2O, ...)

Tachyon

Existing Storage Systems
(HDFS, S3, GlusterFS, ...)

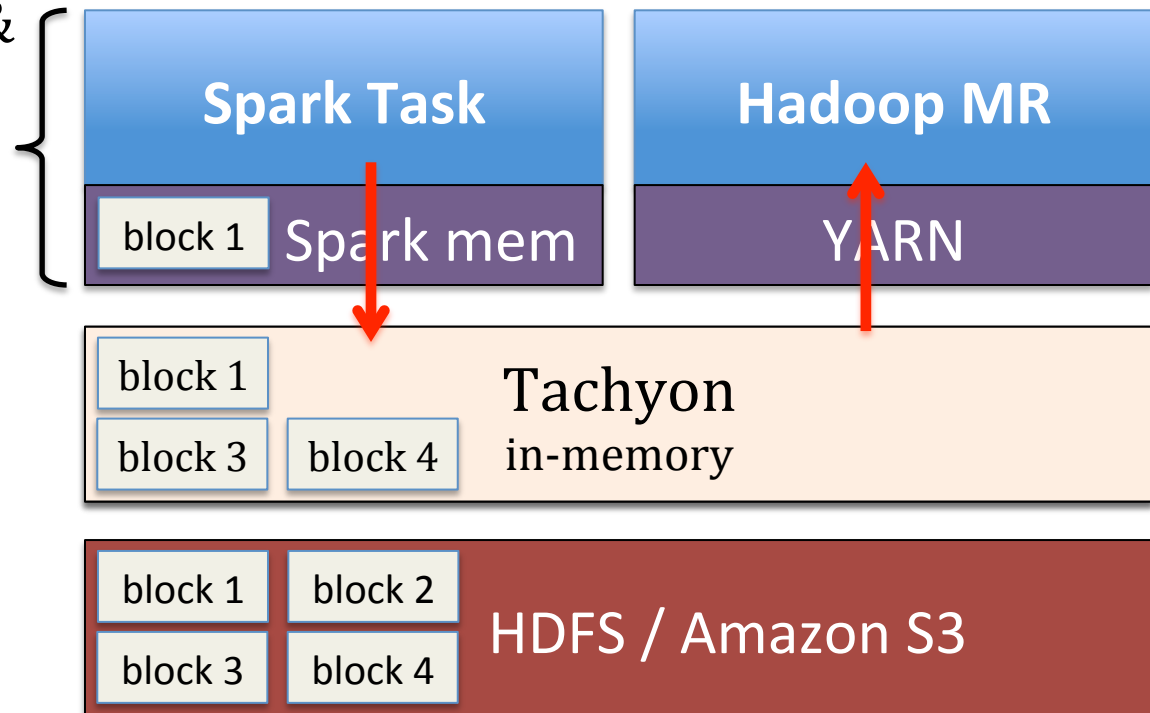# Memory-Centric Storage Architecture
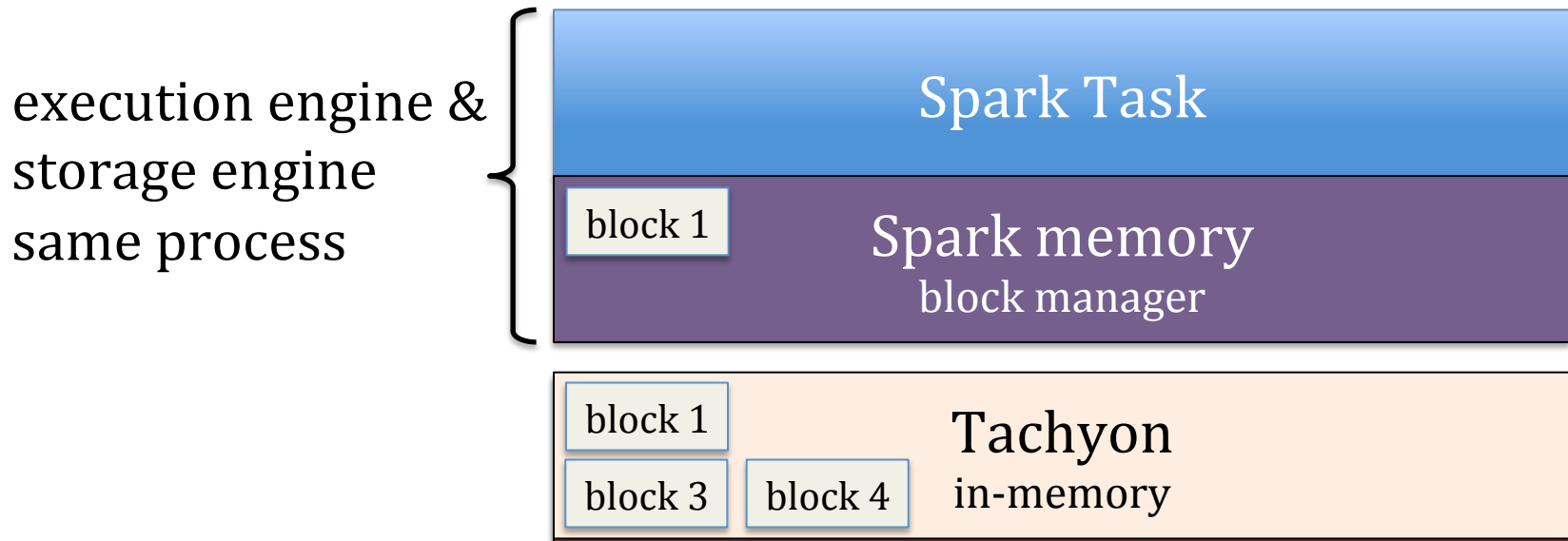
# Lineage in Storage

# Issue 1 revisited

## *Memory-speed data sharing among jobs in different frameworks*

execution engine &
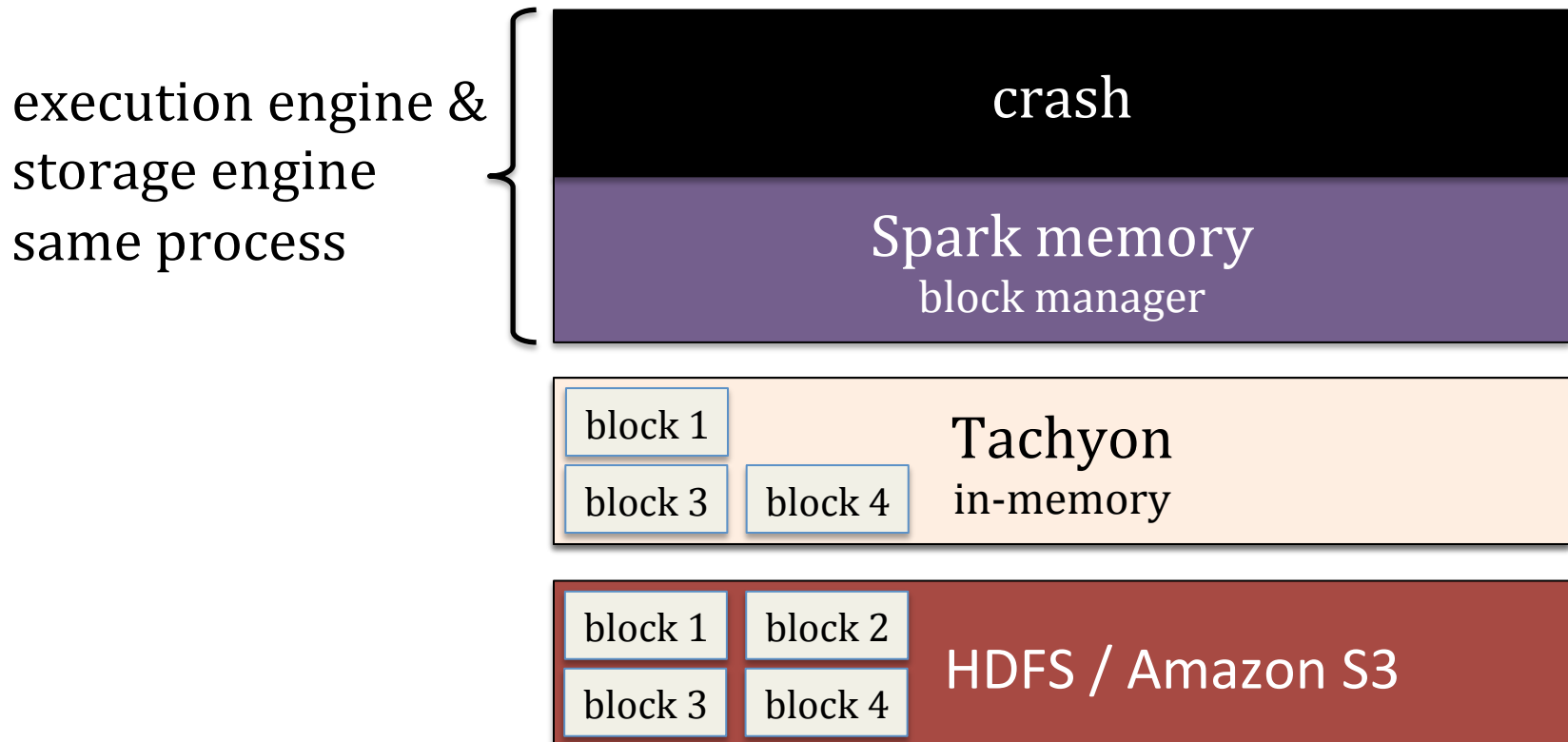storage engine
same process
(fast writes)

| Spark Task | Hadoop MR |
|------------|-----------|
| block 1  Spark mem | YARN |

**Tachyon** in-memory
- block 1
- block 3
- block 4

**HDFS / Amazon S3**
- block 1
- block 2
- block 3
- block 4

# Issue 2 revisited

## *Keep in-memory data safe, even when a job crashes.*

execution engine & storage engine same process

| Spark Task |
|---|
| block 1  Spark memory  block manager |

| block 1  Tachyon  in-memory |
|---|
| block 3  block 4 |

# Issue 2 revisited

*Keep in-memory data safe,
even when a job crashes.*

execution engine &
storage engine
same process

| crash |
| :---: |
| **Spark memory**<br>block manager |

| block 1 | | Tachyon |
| :---: | :---: | :---: |
| block 3 | block 4 | in-memory |

| block 1 | block 2 | HDFS / Amazon S3 |
| :---: | :---: | :---: |
| block 3 | block 4 | |

# Issue 2 revisited

## *Keep in-memory data safe, even when a job crashes.*

execution engine & storage engine same process

crash

block 1 | block 3 | block 4 | Tachyon in-memory

block 1 | block 2 | block 3 | block 4 | HDFS / Amazon S3

# Issue 3 revisited

## *No in-memory data duplication, much less GC*

execution engine &
storage engine
same process
(no duplication & GC)

| Spark Task | Spark Task |
|---|---|
| Spark mem | Spark mem |

Tachyon
in-memory

block 1

block 3 | block 4

HDFS / Amazon S3

block 1 | block 2

block 3 | block 4

# Outline

- Overview
  - Feature 1: Memory Centric Storage Architecture
  - Feature 2: Lineage in Storage

- **Challenges**

- Open Source

- Future

# Question 1: How long to get missing data back?

Lineage enables Asynchronous Checkpointing

# Edge Algorithm

- Checkpoint leaves
- Checkpoint hot files
- Bounded Recovery Cost

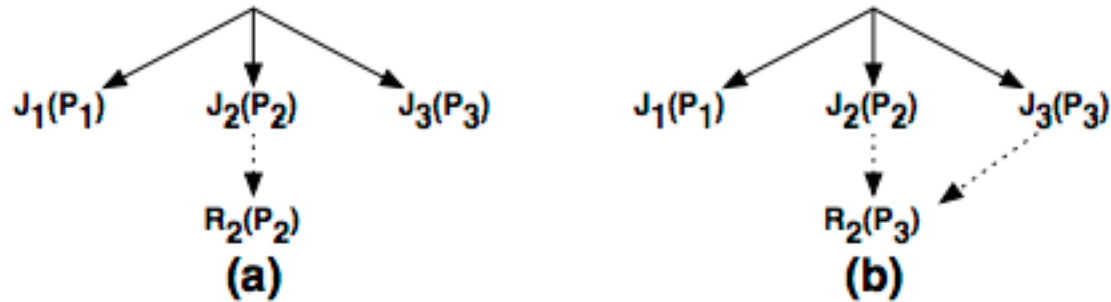# Question 2: How to allocate recomputation resource?



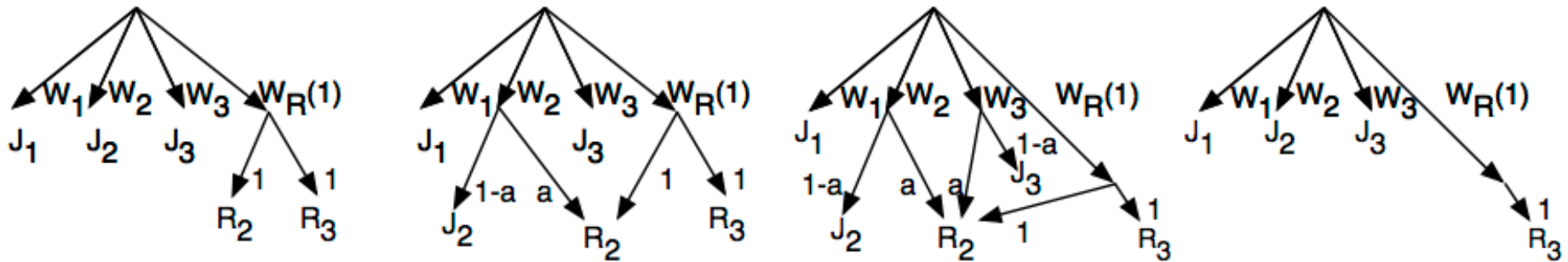Would recomputation slow down my high priority jobs?
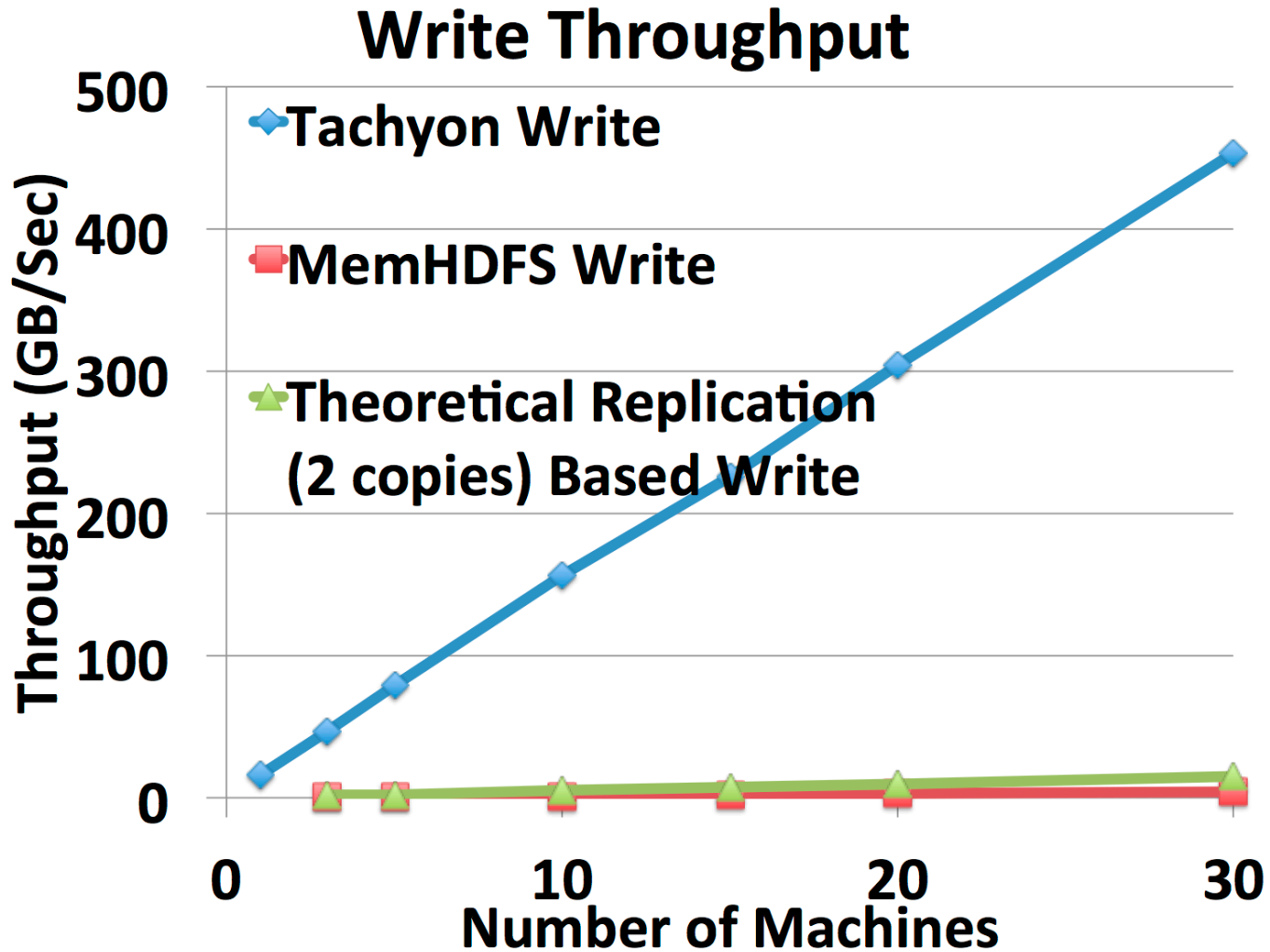Priority Inversion?

# Recomputation Resource Allocation

- Priority Based Scheduler



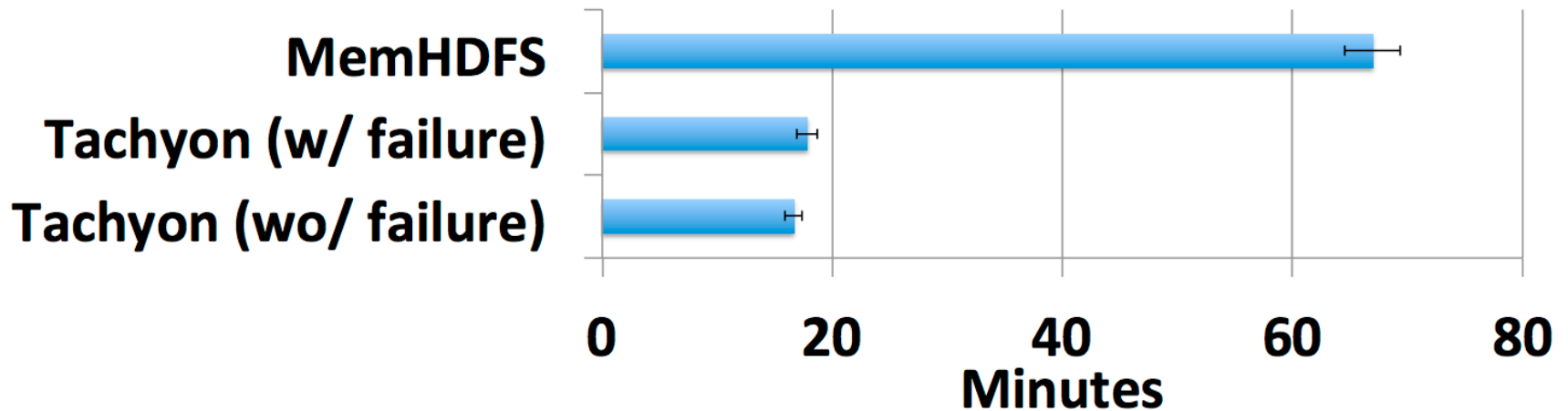- Fair Sharing Based Scheduler

# Comparison with in Memory HDFS



**Write Throughput**

Throughput (GB/Sec)

- Tachyon Write
- MemHDFS Write
- Theoretical Replication (2 copies) Based Write

Number of Machines
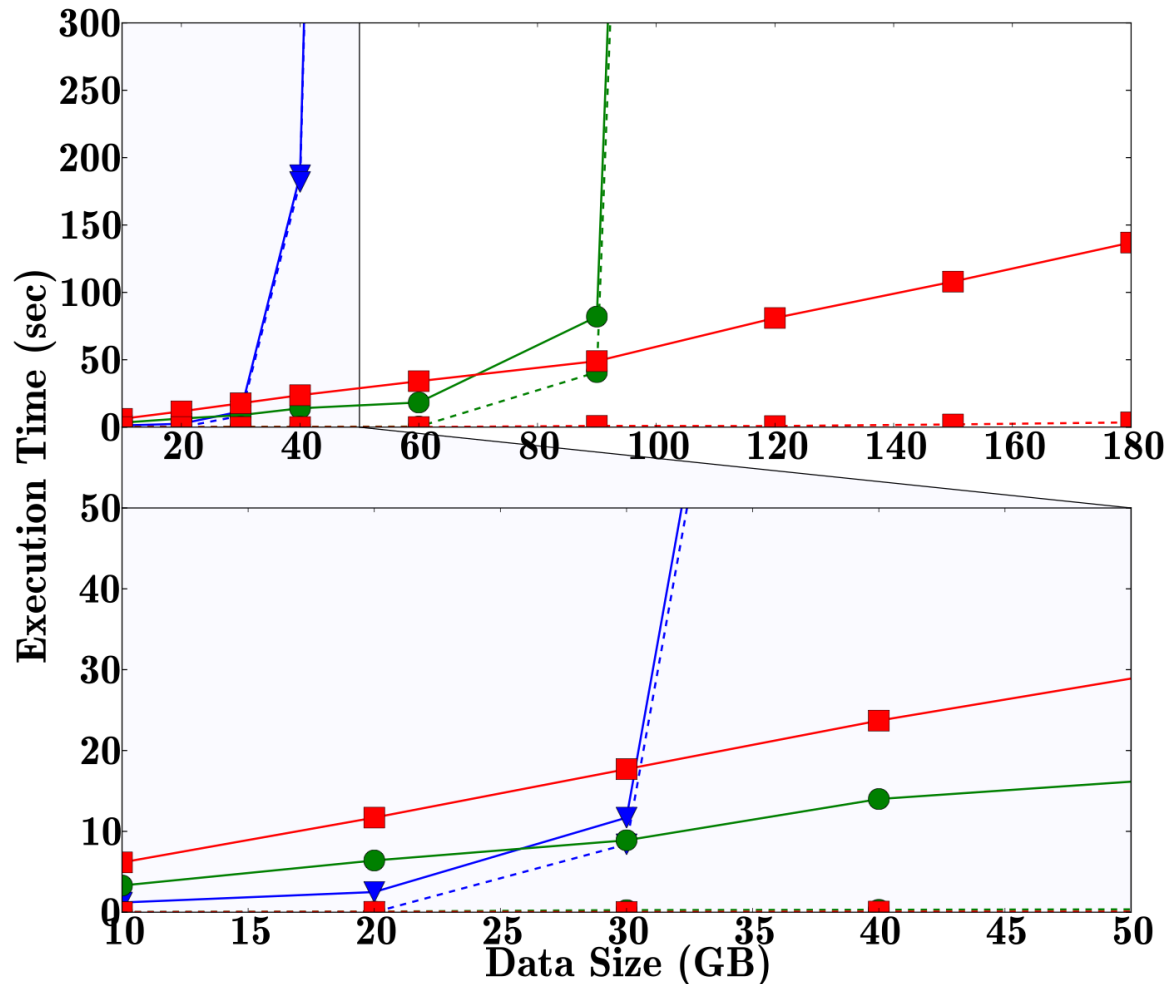
# Workflow Improvement



Performance comparison for realistic workflow. The workflow ran 4x faster on Tachyon than on MemHDFS. In case of node failure, applications in Tachyon still finishes 3.8x faster.
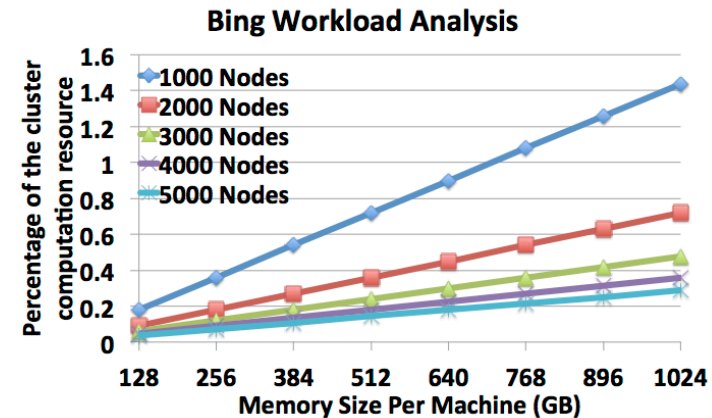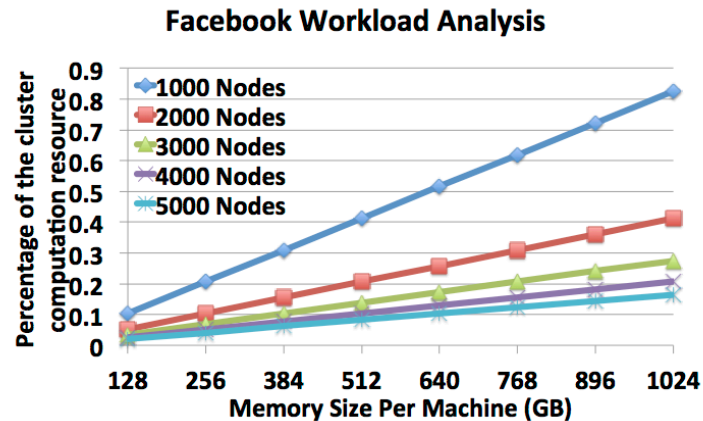
# Further Improve Spark's Performance

**Legend:**
- ▼ Spark Cache (Deserialized)
- ▼ Spark Cache GC (Deserialized)
- ● Spark Cache (Serialized)
- ● Spark Cache GC (Serialized)
- ■ TFS
- ■ TFS GC

**Grep Program**

# Recomputation Resource Consumption



**Trace Summary**



Facebook Workload Analysis



Bing Workload Analysis

# Outline

- Overview
  - Feature 1: Memory Centric Storage Architecture
  - Feature 2: Lineage in Storage

- Challenges

- **Open Source**

- Future
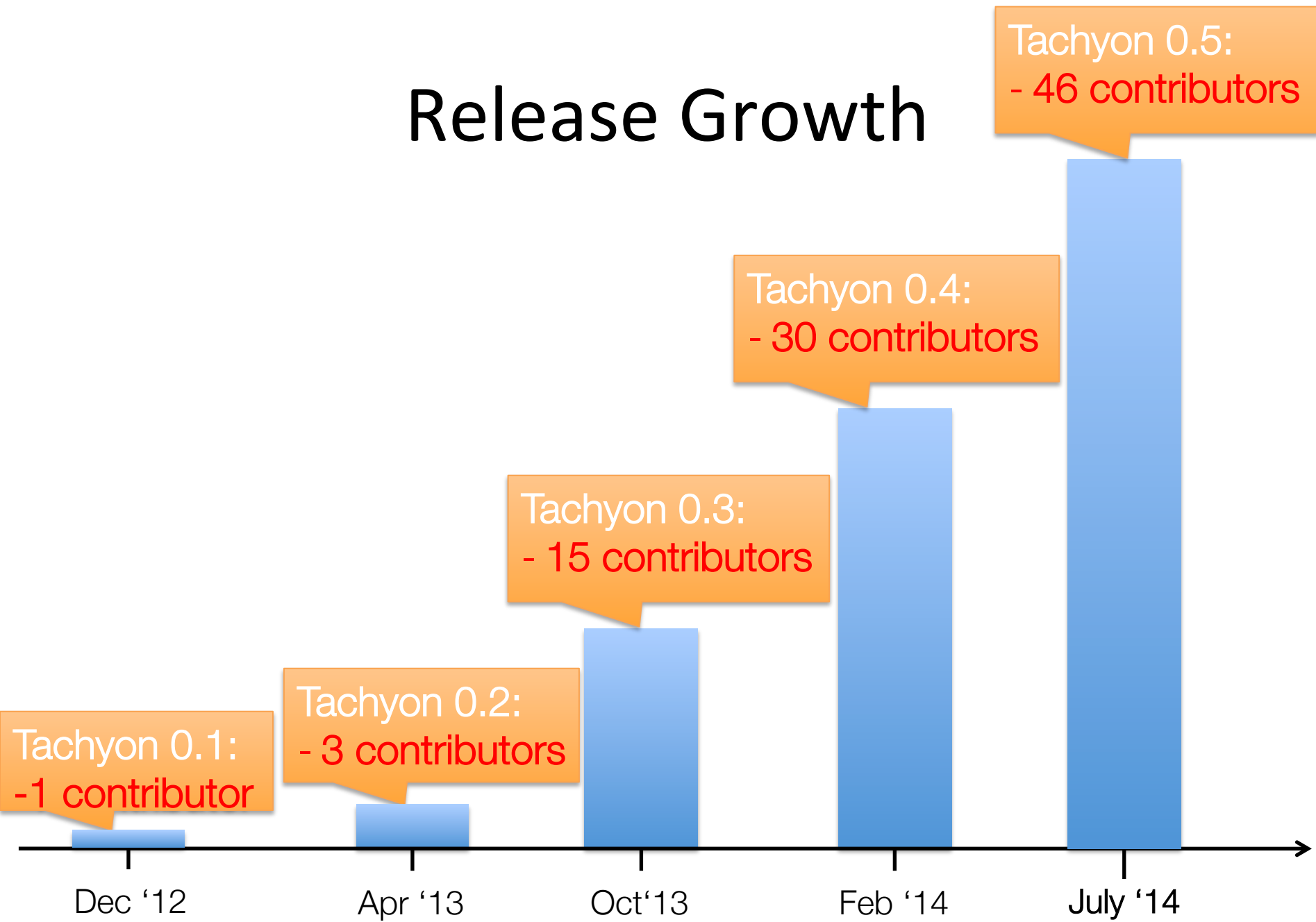
# TACHYON Open Source Status
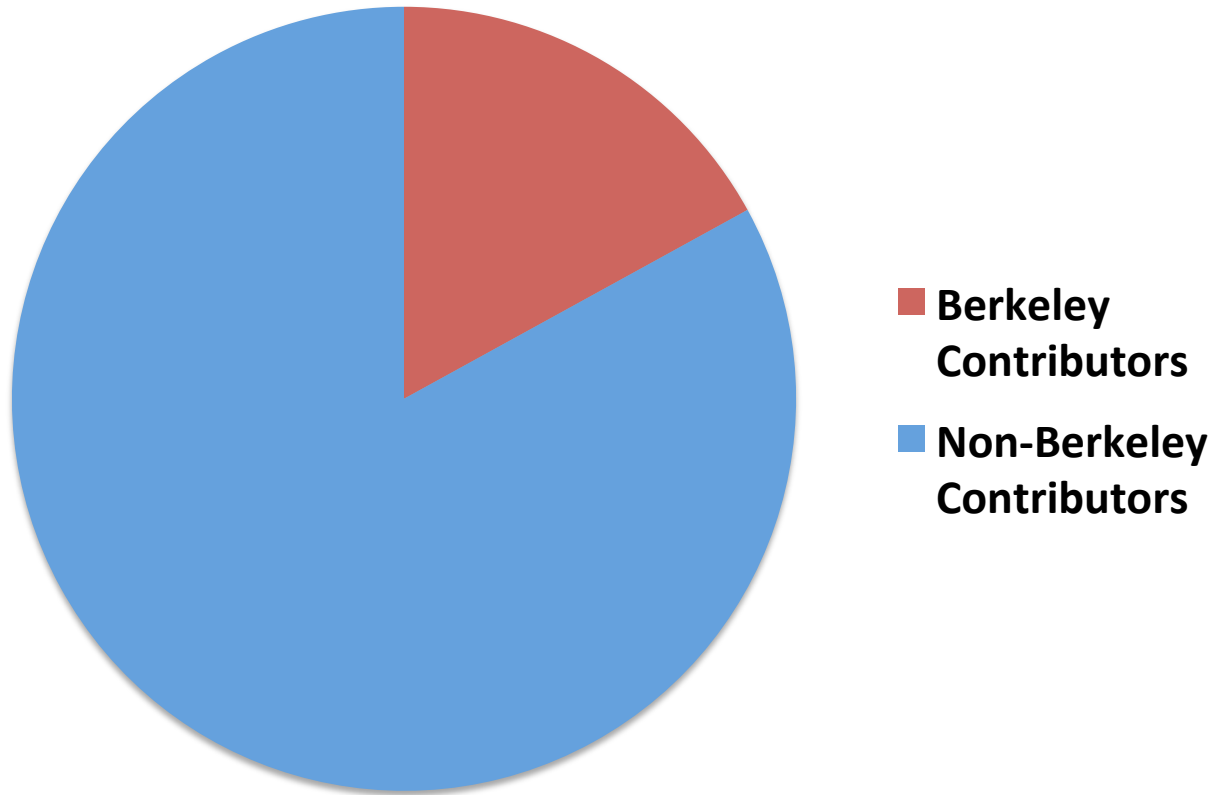
- Apache License 2.0, Version 0.5.0 (July 2014)

  PUBLIC  amplab / tachyon        ★ Unstar    859

- Deployed at tens of companies

- 15+ Companies Contributing

- No code change for Spark and MapReduce applications.

# Release Growth

Tachyon 0.5:
- 46 contributors

Tachyon 0.4:
- 30 contributors

Tachyon 0.3:
- 15 contributors

Tachyon 0.2:
- 3 contributors

Tachyon 0.1:
-1 contributor

Dec '12          Apr '13          Oct'13          Feb '14          July '14

# Open Community



- ■ **Berkeley Contributors**
- ■ **Non-Berkeley Contributors**

# Thanks to our Code Contributors!

| | | | |
|---|---|---|---|
| Aaron Davidson | David Capwell | Lukasz Jastrzebski | Rong Gu |
| Achal Soni | David Zhu | Manu Goyal | Sean Zhong |
| Ali Ghodsi | Du Li | Mark Hamstra | Seonghwan Moon |
| Andrew Ash | Fei Wang | Mingfei Shi | Shivaram Venkataraman |
| Anurag Khandelwal | Gerald Zhang | Mubarak Seyed | Srinivas Parayya |
| Aslan Bekirov | Grace Huang | Nick Lanham | Tao Wang |
| Bill Zhao | Haoyuan Li | Orcun Simsek | Thu Kyaw |
| Brad Childs | Henry Saputra | Pengfei Xuan | Timothy St. Clair |
| Calvin Jia | Hobin Yoon | Qianhao Dong | Vamsi Chitters |
| Chao Chen | Huamin Chen | Qifan Pu | Xi Liu |
| Cheng Chang | Jey Kottalam | Raymond Liu | Xiang Zhong |
| Cheng Hao | Joseph Tang | Reynold Xin | Xiaomin Zhang |
| Colin Patrick McCabe | Juan Zhou | Robert Metzger | Zhao Zhang |

# Tachyon
# is in Fedora 20

**Thanks to Redhat!**

# Commercially supported by  Atigeo™ and running in dozens of their customers' clusters
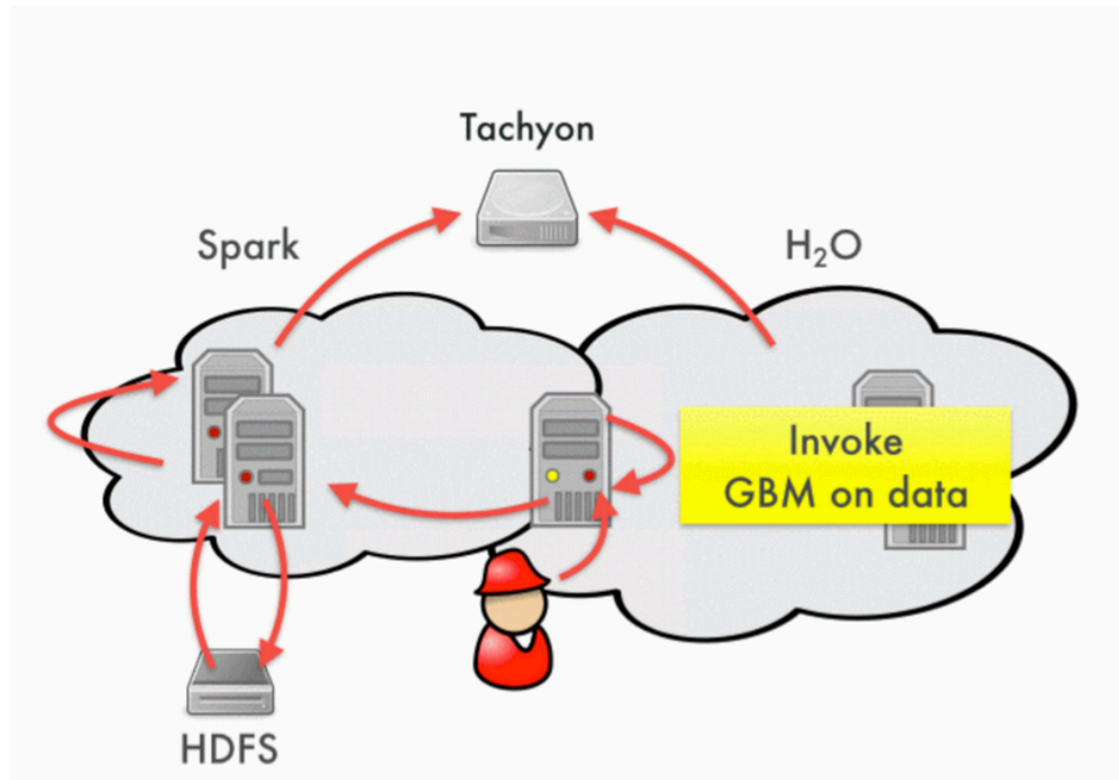
# Tachyon is the Default Off-Heap Storage Solution for Spark

# TACHYON in Spark H₂O

Today, data gets parsed and exchanged between Spark and H2O via Tachyon. Users can interactively query big data both via SQL and ML from within the same context.

# Reaching wider communities:
# e.g. GlusterFS



43

# Under Filesystem Choices
## (Big Data, Cloud, HPC, Enterprise)

# Outline

- Overview
  - Feature 1: Memory Centric Storage Architecture
  - Feature 2: Lineage in Storage


- Challenges


- Open Source


- **Future**

# Short Term Roadmap (0.6 Release)

- Ceph Integration (Ceph Community)

- Hierarchical Local Storage (Intel)

- Performance Improvement (Yahoo)

- Multi-tenancy (AMPLab)

- Mesos Integration (Mesos Community)

- ***Many more*** from AMPLab and Industry Collaborators.

# Features

- **Memory Centric Storage Architecture**
- **Lineage in Storage (alpha)**
- Hierarchical Local Storage
- Data Serving
- Scalable metadata management
- Different hardware
- More…
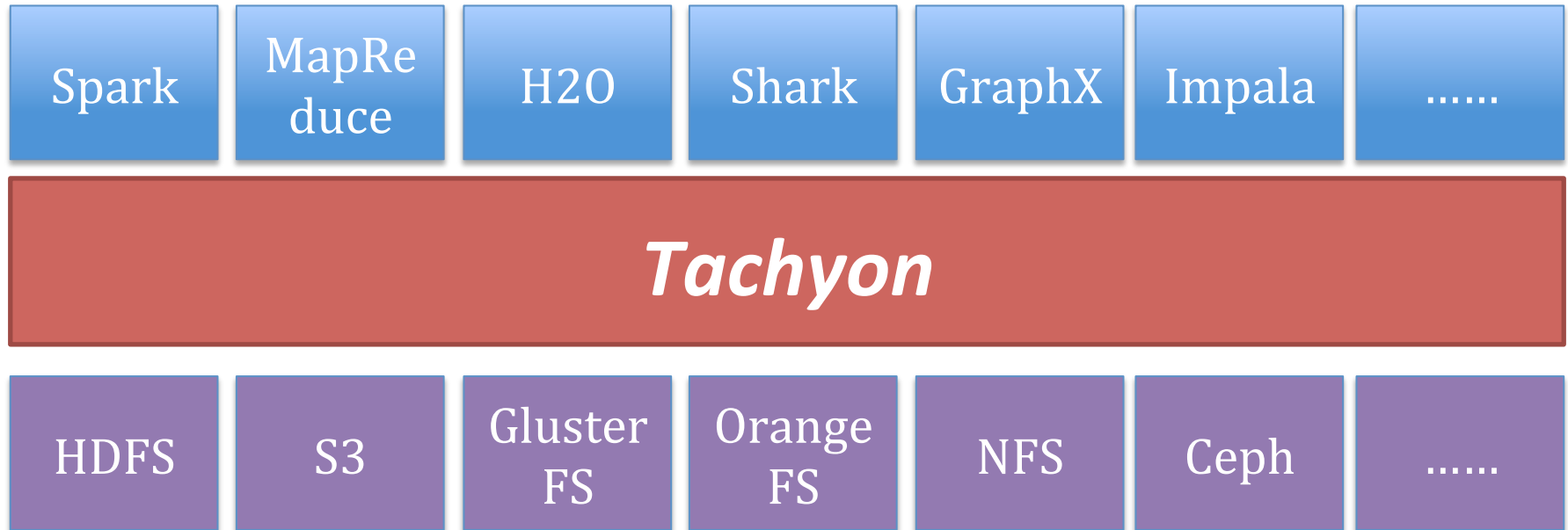- Your Requirements?

# Data Serving: An Example

Data Analytics Pipeline:

Query the results of batch jobs.

# What do we need?

Sequential I/O   **+**   Random Access!

# Tachyon Goal?

# Better Assist Other Components



| Spark | MapReduce | H2O | Shark | GraphX | Impala | ...... |

**Tachyon**

| HDFS | S3 | Gluster FS | Orange FS | NFS | Ceph | ...... |

**Welcome Collaboration!**

# *Thanks! Questions?*

- *More Information:*
  - *Website: [http://tachyon-project.org](http://tachyon-project.org)*
  - *Github: [https://github.com/amplab/tachyon](https://github.com/amplab/tachyon)*
  - *Meetup: [http://www.meetup.com/Tachyon](http://www.meetup.com/Tachyon)*
- *Email: [haoyuan@cs.berkeley.edu](mailto:haoyuan@cs.berkeley.edu)*