# A Reliable Memory-Centric Distributed Storage System

**Haoyuan Li, Tachyon Nexus**
haoyuan@tachyonnexus.com
September 22, 2015 @ SDC 2015

# TACHYON
## NEXUS

- Team consists of Tachyon creators, top contributors, people from UC Berkeley, Google, CMU, VMware, Stanford, Facebook, etc.

- $7.5 million Series A from Andreessen Horowitz

- Committed to Tachyon Open Source

# TACHYON
## N E X U S

**WE'RE HIRING!**

# Outline

- Overview
  - Motivation
  - Tachyon Architecture
  - Using Tachyon
- Open Source
  - Status
  - Production Use Cases
- Roadmap

**TACHYON**
N E X U S

# Outline

- **Overview**
  - Motivation
  - Tachyon Architecture
  - Using Tachyon
- Open Source
  - Status
  - Production Use Cases
- Roadmap

**TACHYON**
N E X U S

# Tachyon: Born in UC Berkeley AMPLab



Cluster manager

Parallel computation framework

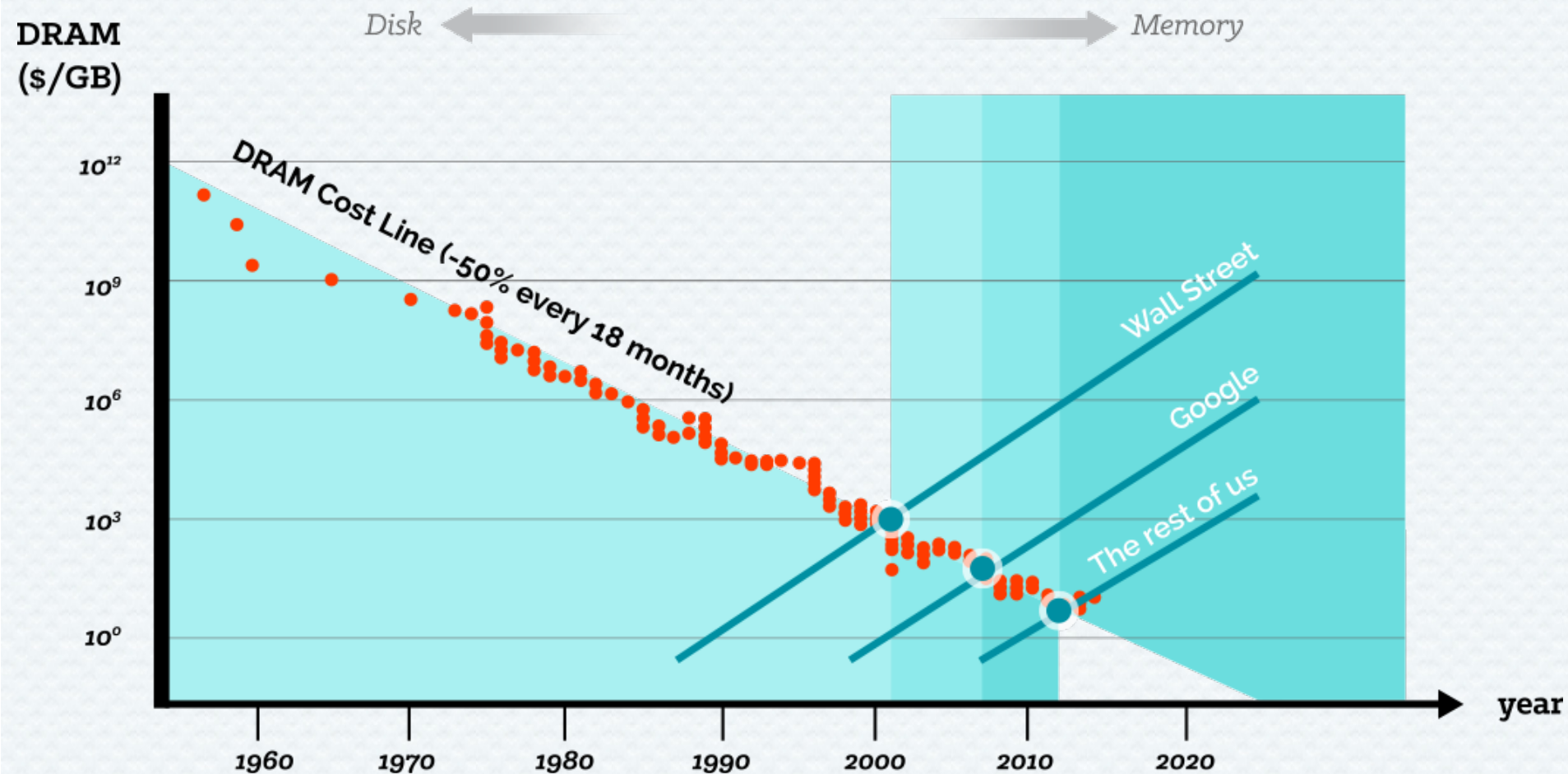Reliable, distributed memory-centric storage system

# **Why Tachyon?**

TACHYON
N E X U S

# Memory is Fast

- 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 times

# Memory is Cheaper



source: jcmit.com

**TACHYON NEXUS**

9

# Realized by many…

# Is the Problem Solved?

**TACHYON**
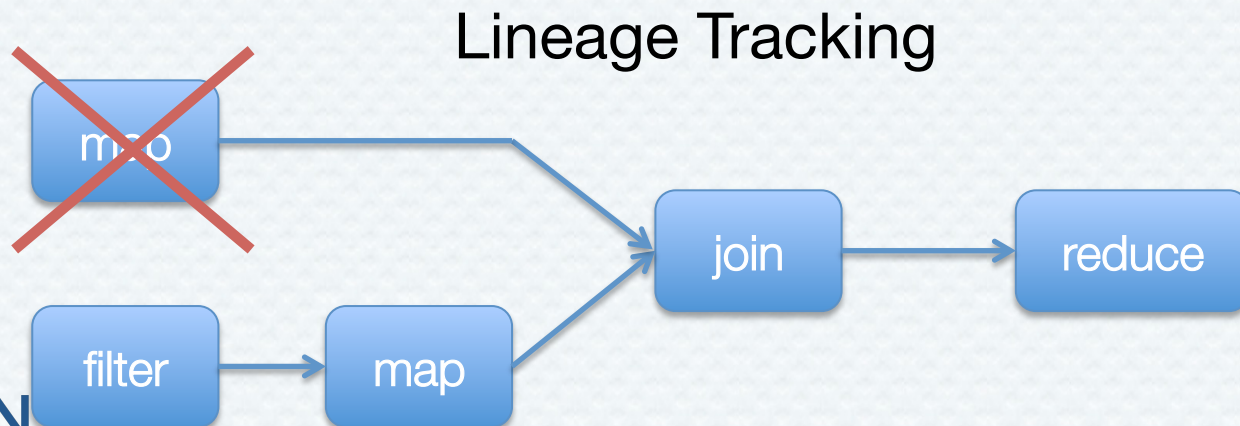N E X U S

# Missing a Solution for the Storage Layer
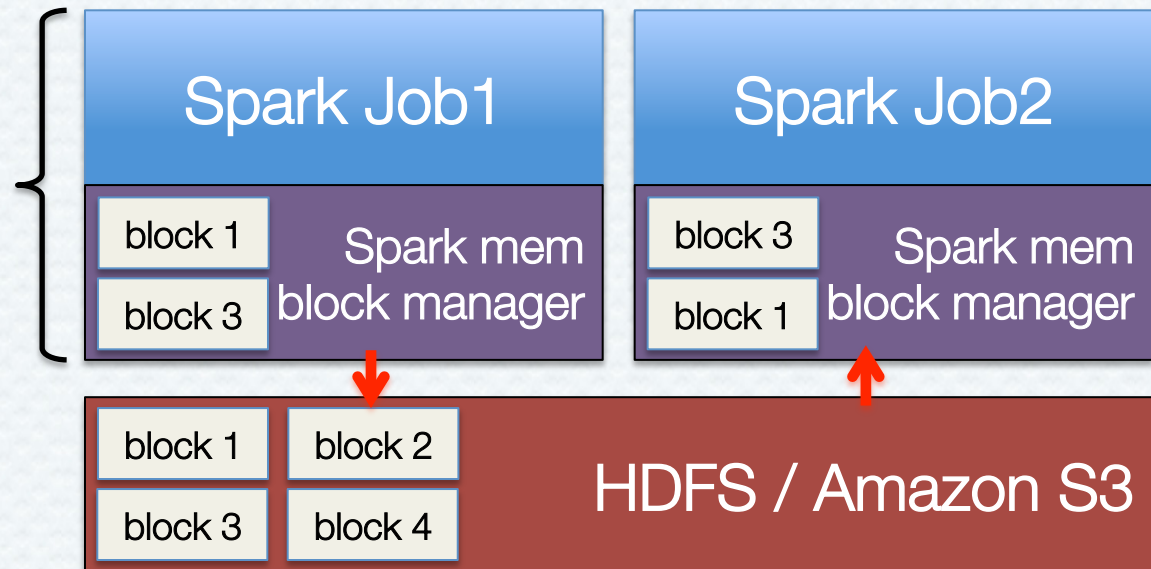
# An Example: Spark

- 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
<span style="color:red">(slow writes)</span>

| Spark Job1 | Spark Job2 |
|---|---|
| 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
N E X U S

# Issue 1

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

storage engine &
execution engine
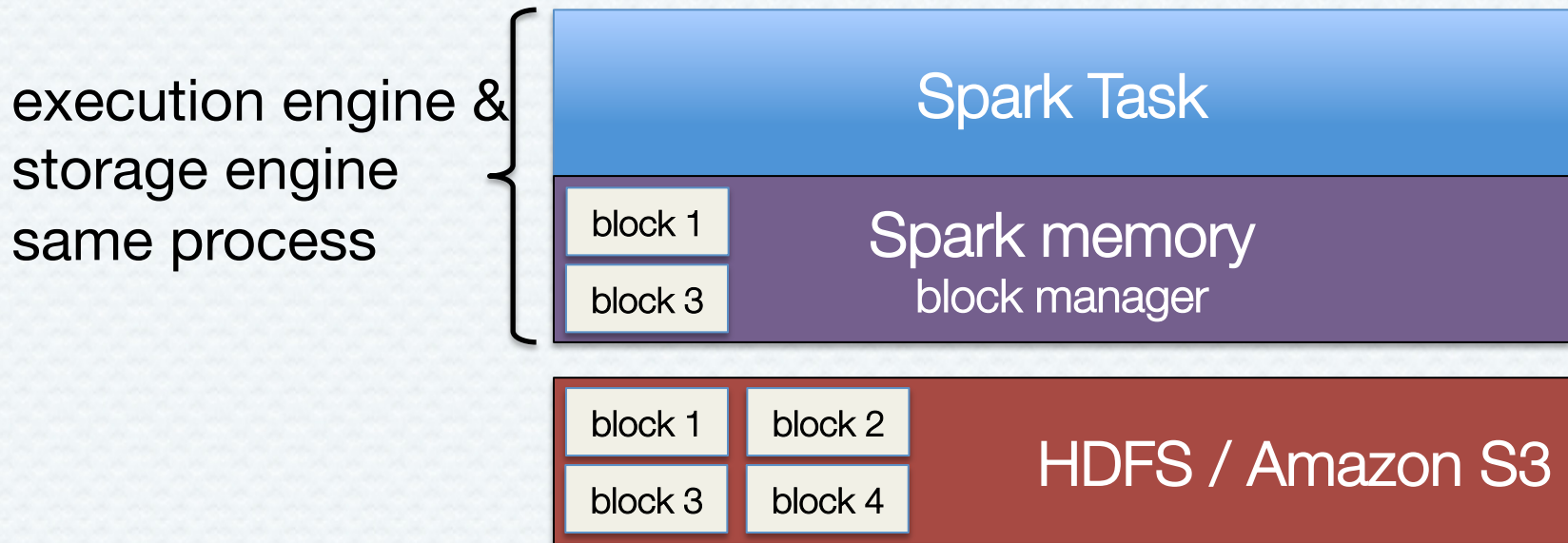same process
<span style="color:red">(slow writes)</span>



Spark Job

Hadoop MR Job

block 1

block 3

Spark mem
block manager

YARN

block 1 | block 2

block 3 | block 4

HDFS / Amazon S3

TACHYON
NEXUS

# Issue 2

## *Cache loss when process crashes*

execution engine &
storage engine
same process

| Spark Task |
| --- |

| block 1 | Spark memory |
| block 3 | block manager |

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

TACHYON
N E X U S

# Issue 2

## *Cache loss when process crashes*

execution engine & storage engine same process

crash

| block 1 | Spark memory |
|---------|--------------|
| block 3 | block manager |

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

**TACHYON**
N E X U S

17

# Issue 2

## *Cache loss when process crashes*

execution engine &
storage engine
same process

crash

block 1   block 2

block 3   block 4

HDFS / Amazon S3

**TACHYON**
N E X U S

# Issue 3

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

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

| Spark Task1 | Spark Task2 |
|---|---|
| block 1 | block 3 |
| block 3 Spark mem block manager | block 1 Spark mem block manager |

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

# Tachyon

*Reliable* data sharing at

*memory-speed* *within and across* cluster frameworks/jobs
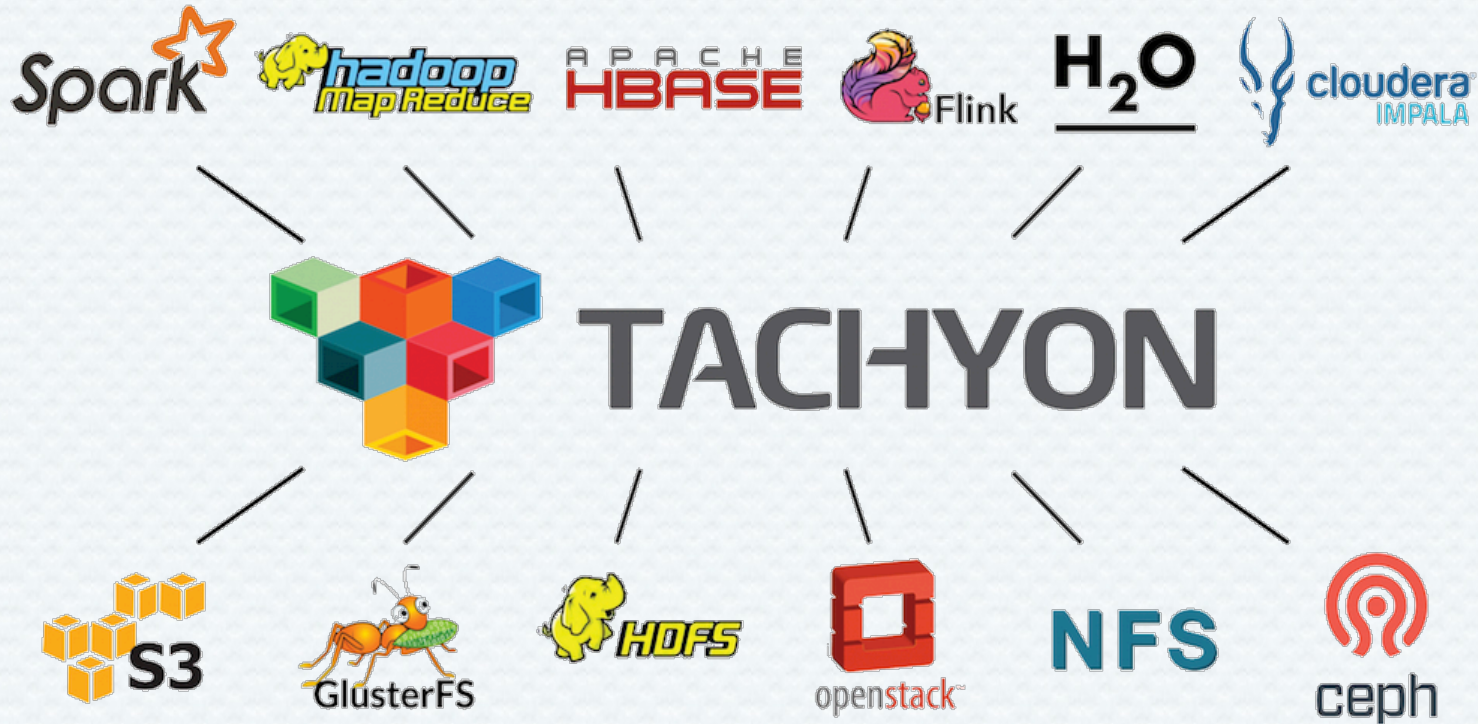
TACHYON
NEXUS

# Technical Overview

**Ideas**

- A **memory-centric** storage architecture
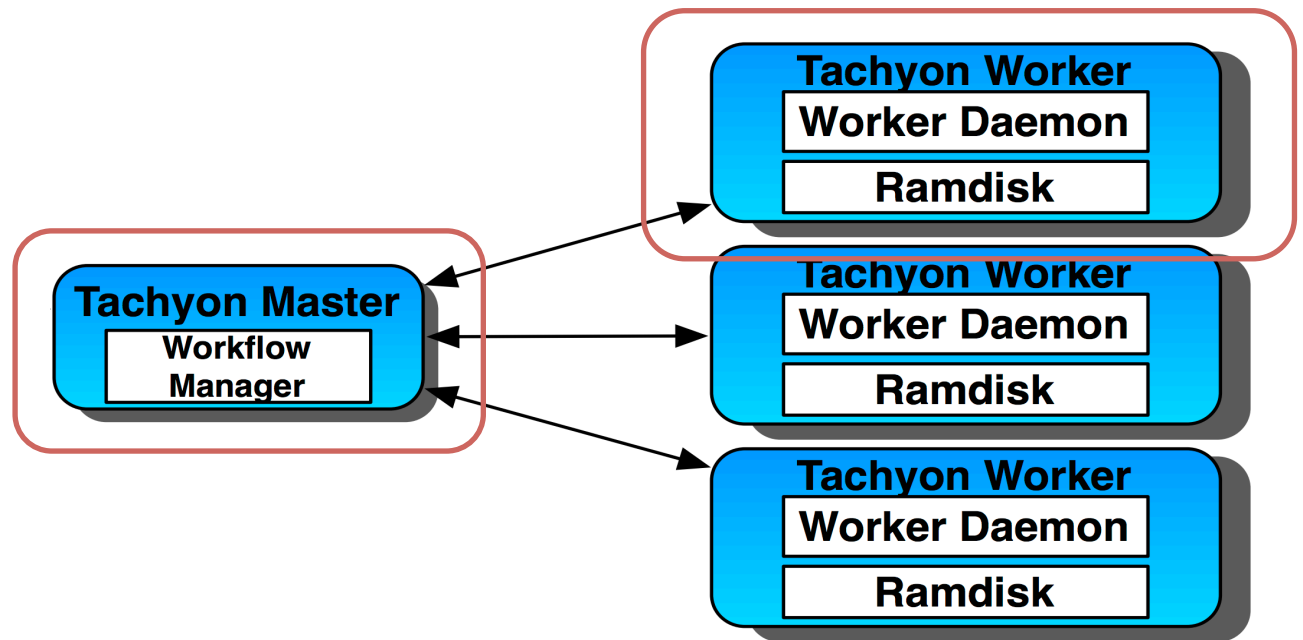- Push **lineage** down to storage layer
- Manage **tiered** storage

**Facts**

- One data copy in memory
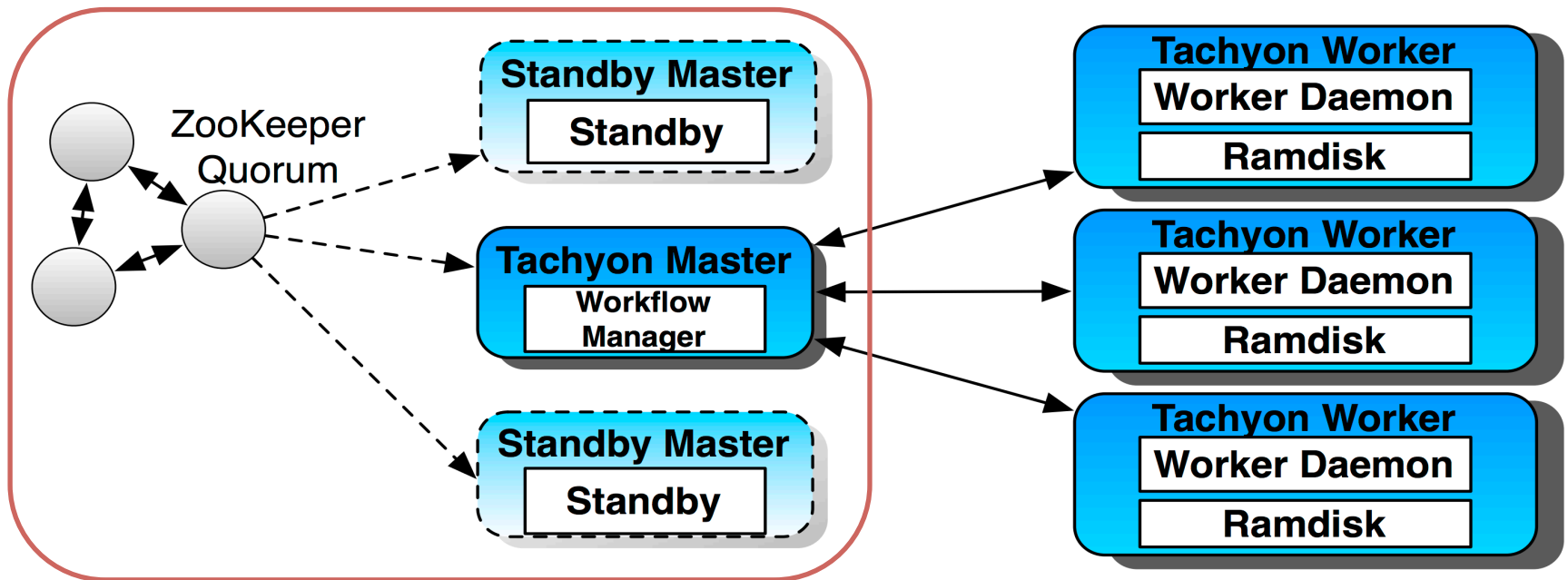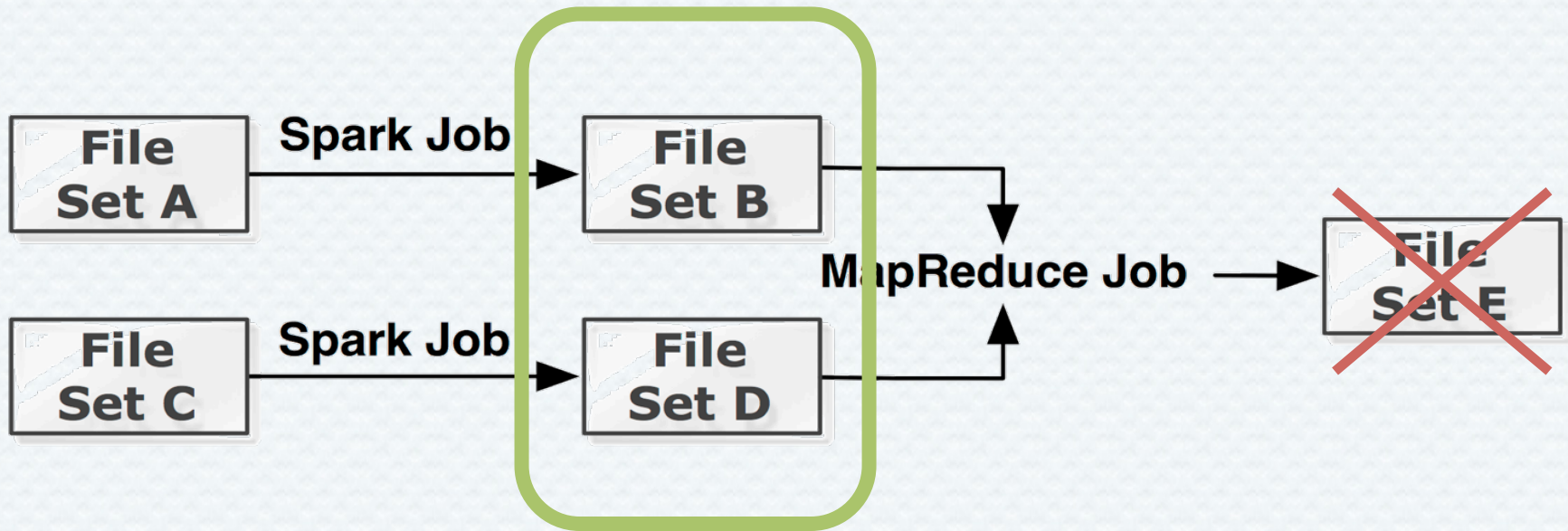- Re-computation for fault-tolerance

TACHYON
N E X U S

# Eco-System

# Tachyon Memory-Centric Architecture
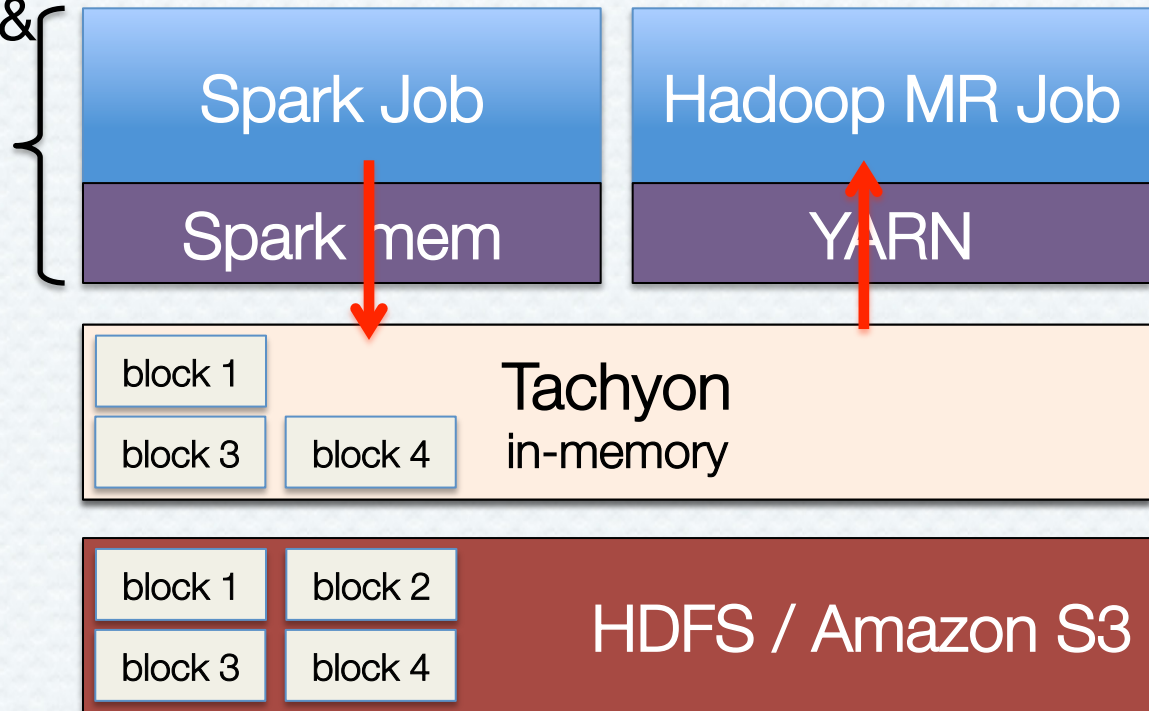
# Tachyon Memory-Centric Architecture

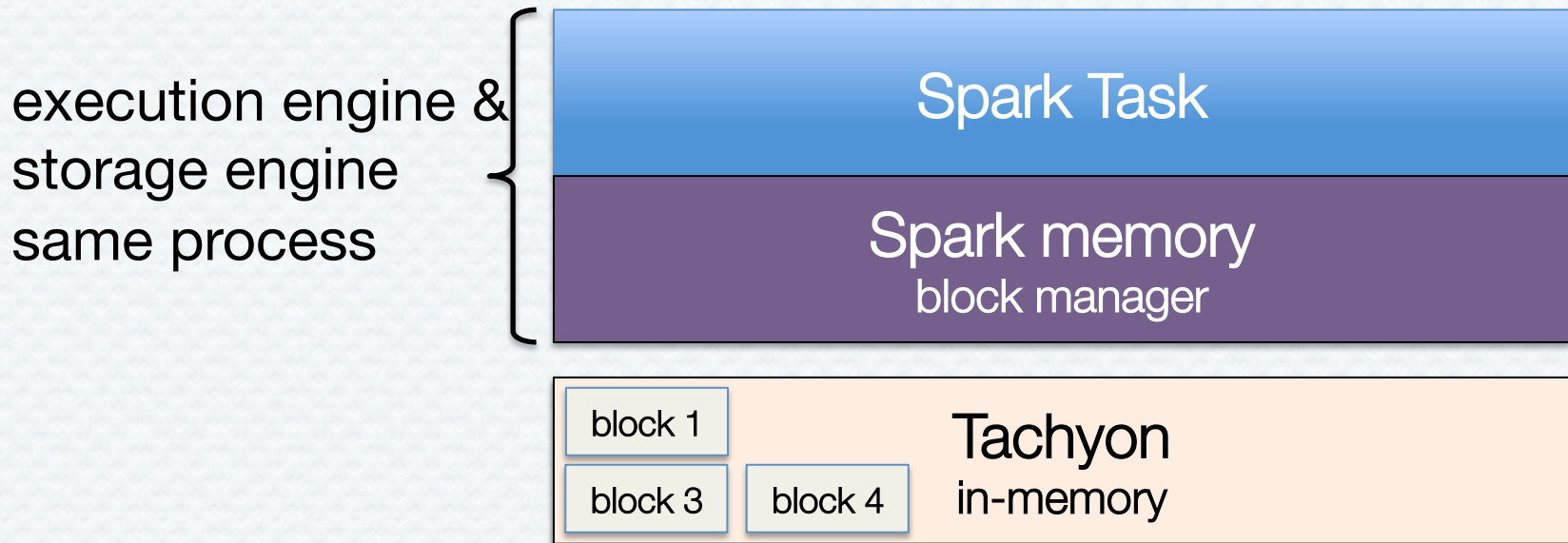# Lineage in Tachyon

# Issue 1 revisited

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

execution engine & storage engine same process (fast writes)

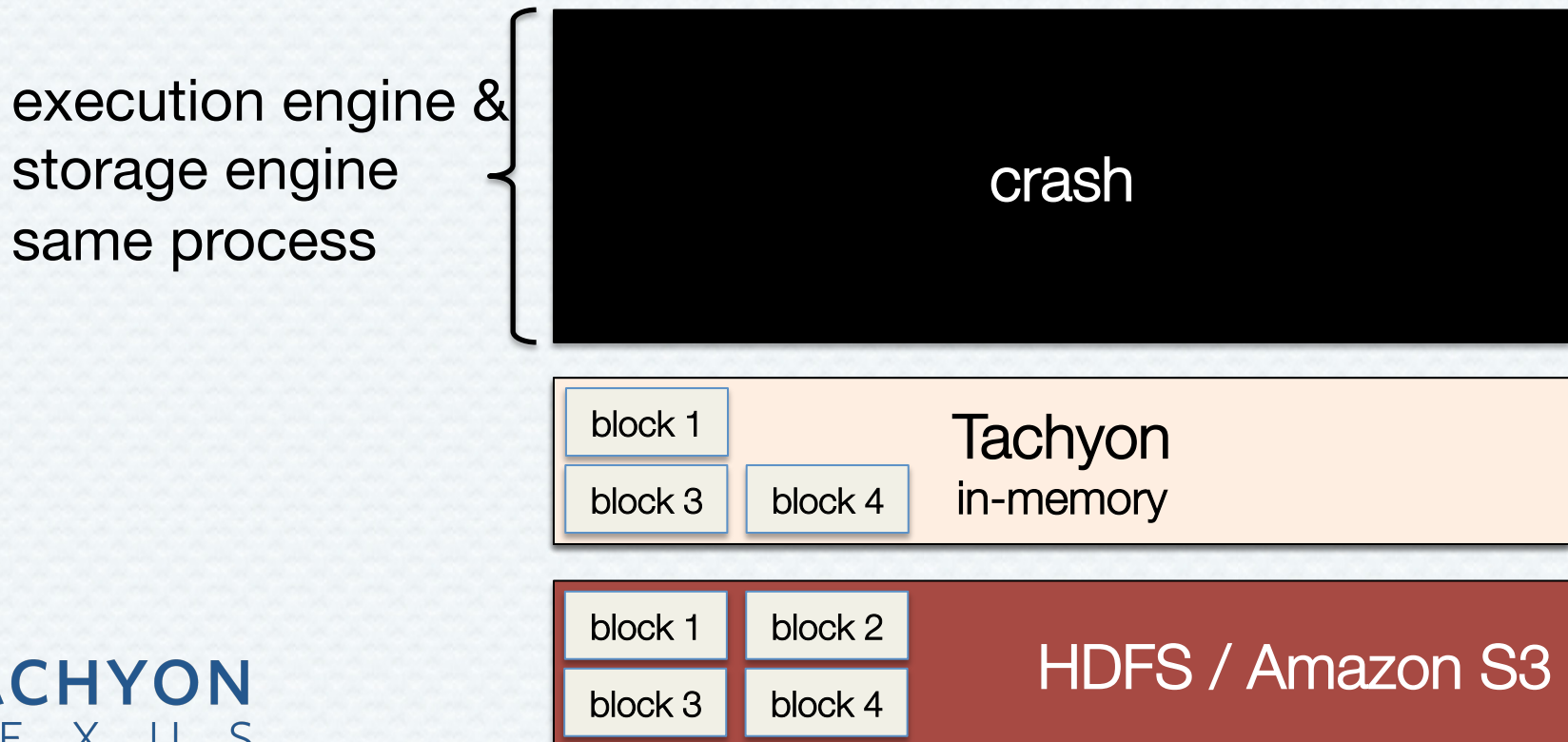# Issue 2 revisited
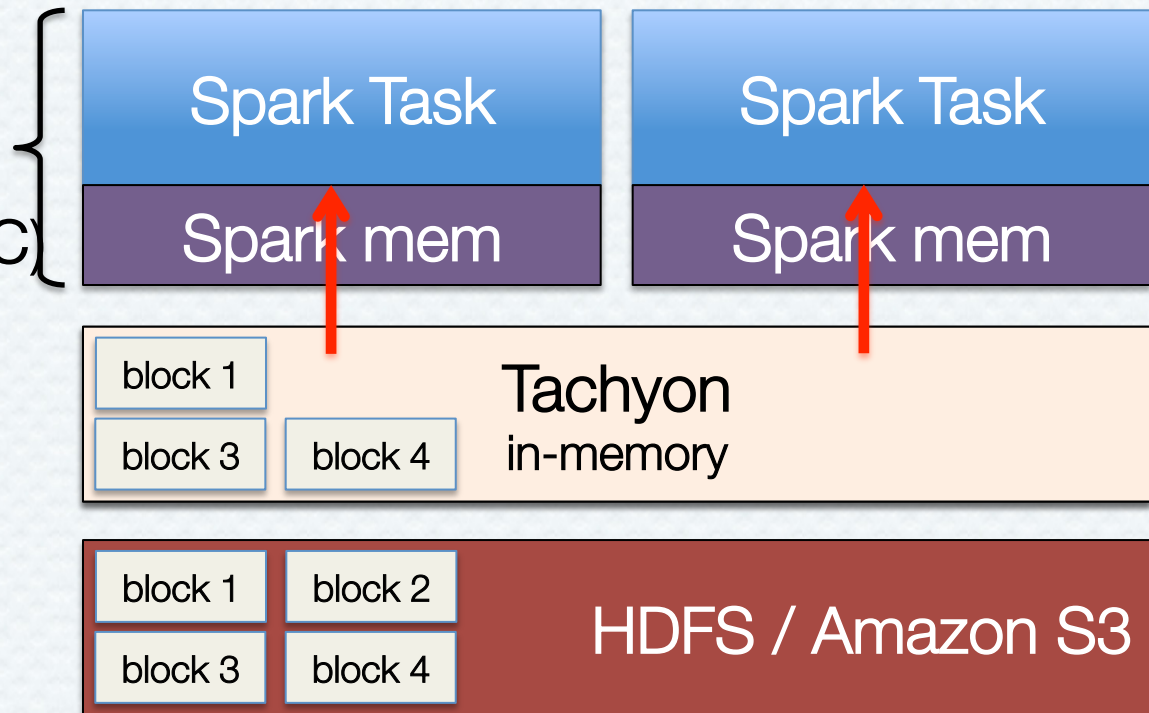
*Keep in-memory data safe, even when a job crashes.*

execution engine &
storage engine
same process

| Spark Task |
| --- |
| 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

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

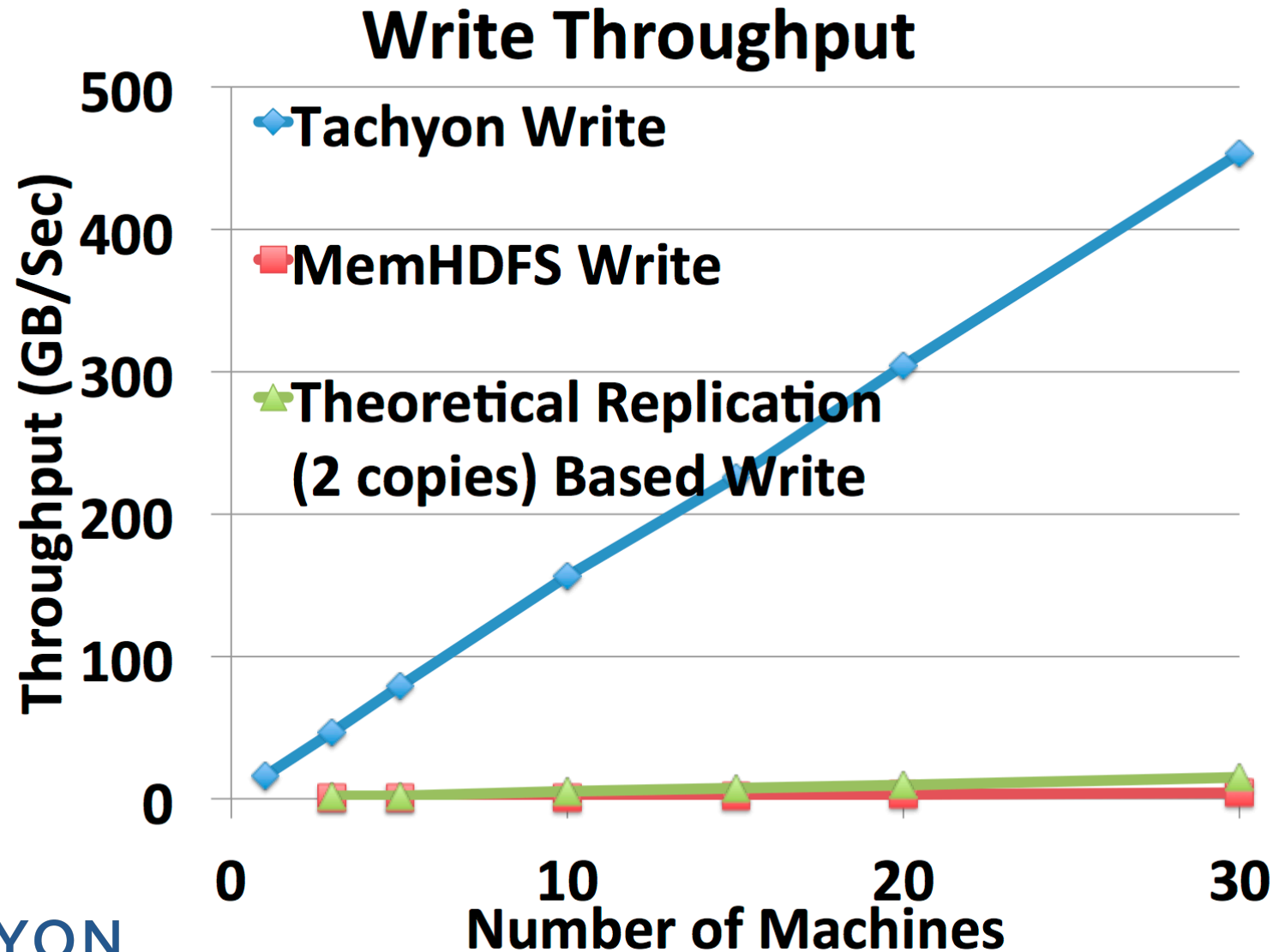| block 1 | block 2 | HDFS / Amazon S3 |
| block 3 | block 4 | |

TACHYON
N E X U S

# 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 |

block 1

block 3    block 4

Tachyon
in-memory

block 1    block 2

block 3    block 4

HDFS / Amazon S3

TACHYON
N E X U S

29

# Comparison with In-Memory HDFS

# Outline

- Overview
  - Motivation
  - Tachyon Architecture
  - Using Tachyon
- **Open Source**
  - Status
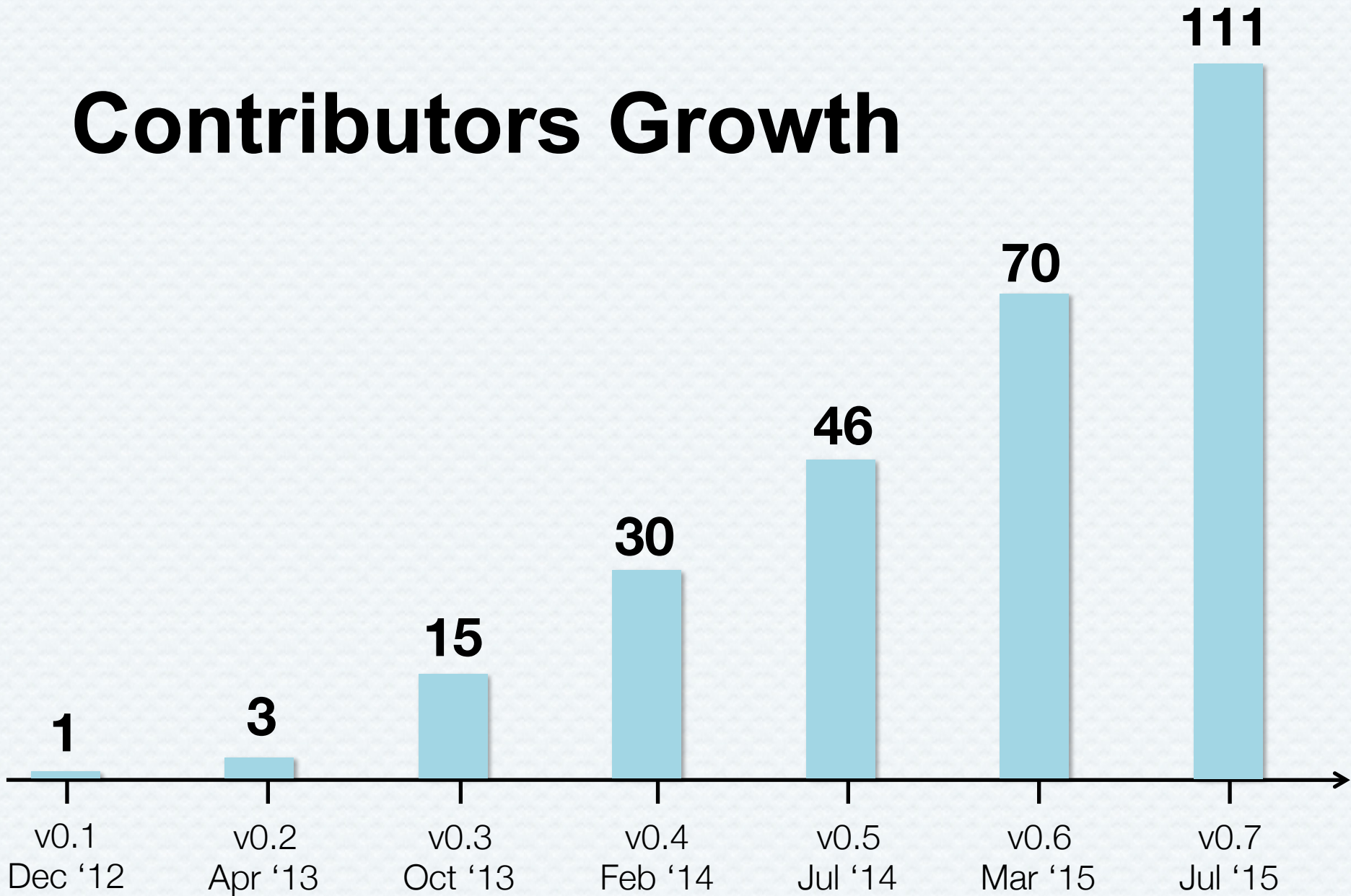  - Production Use Cases
- Roadmap

**TACHYON**
NEXUS

# Open Source Status

- Started at UC Berkeley AMPLab in Summer 2012

- Apache License 2.0, Version 0.7.1 (August 2015)

PUBLIC amplab / tachyon ★ Unstar 1,569
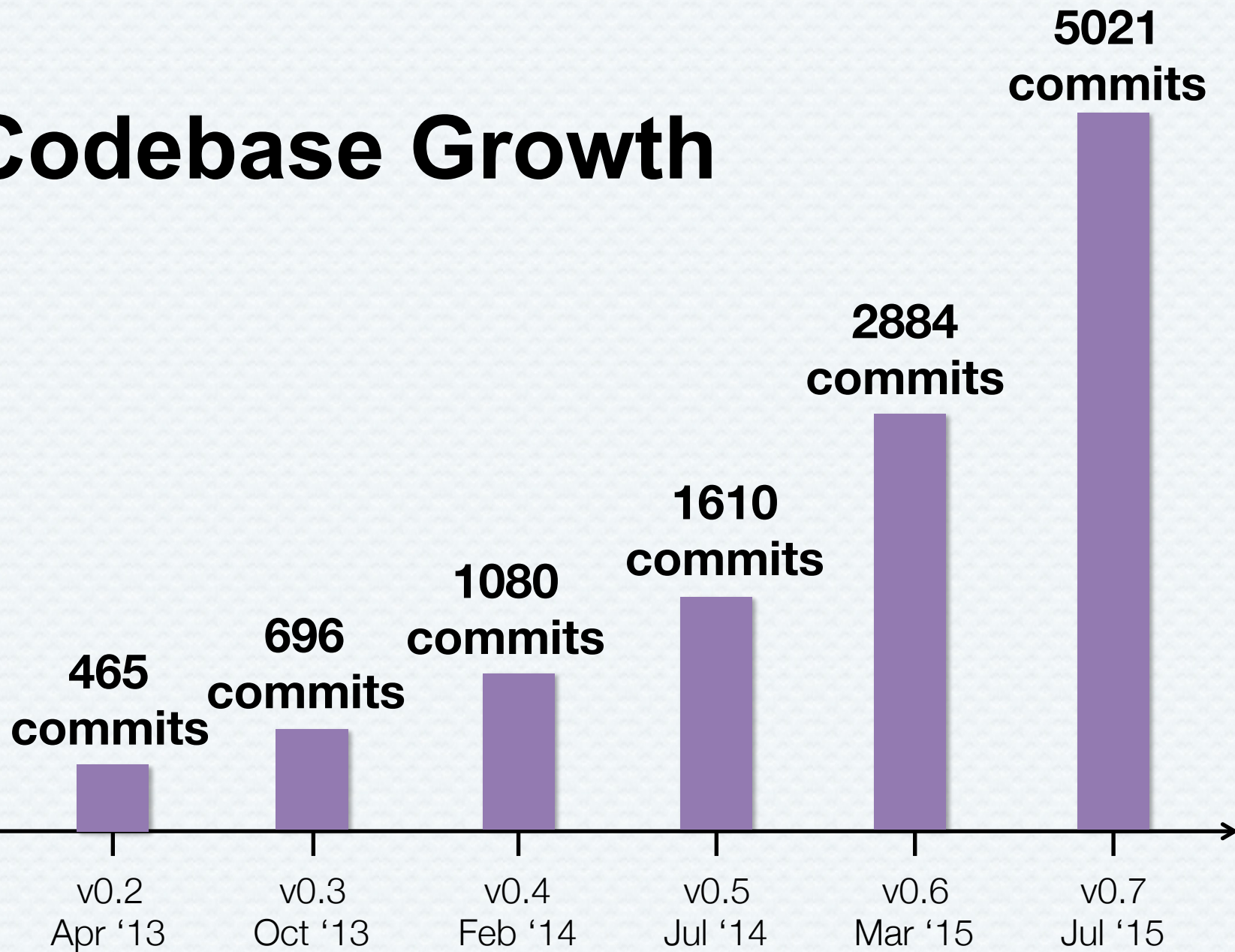
- Deployed at > 50 companies (July 2014)

- 30+ Companies Contributing

**TACHYON**
N E X U S

# Codebase Growth

5021 commits

2884 commits

1610 commits

1080 commits

696 commits

465 commits

v0.2 Apr '13

v0.3 Oct '13

v0.4 Feb '14

v0.5 Jul '14

v0.6 Mar '15

v0.7 Jul '15

TACHYON NEXUS

34

# Thanks to Our Contributors!

# Reported Tachyon Usage



**Tachyon is the** **-Heap** **lution**

Pivotal

**EMC²**

**GIGAOM**

**ZDNet**

databas
TRENDS AND APPLICAT

The Future Architecture of a Data Lake:
In-memory Data Exchange Platform

**IBM Research**

8.17.2015

Tachyon for ultra-fast Big Data processing

*Editor's note: This article is by cloud analytics infrastructure expert Gil Vernik, IBM Research-Haifa.*

Today's massive growth in data sets means that storage is increasingly becoming a critical bottleneck for system workloads. My storage team in Haifa, Israel wants to analyze and understand these massive volumes of data, and we need to store them somewhere reliable. Although disk space is an option, it's too slow to carry out fast Big Data processing. In-memory computing, which keeps the data in a server's RAM for fast access and processing, offers a good solution for processing Big Data workloads – but it's limited and expensive.

Enter Tachyon, a memory-centric distributed storage system that offers processing at memory-speed and reliable storage. Its software works with servers in clusters so there's plenty of room for storage, and a unique proprietary feature eliminates the need for replication to ensure fault tolerance. Now, we've connected Tachyon to Swift so it can work effortlessly with Swift and SoftLayer. The result? Tachyon is even more flexible and efficient.

IBM RESEARCH HOMEPAGE

**IBM Research**

IBM RESEARCH ON TWITTER

Follow us on Twitter

BLOG ARCHIVE

▼ 2015 (44)
  ▼ August (6)
    IBM's New Polymers Acclaimed for Use in 3-D Printi...
    Tachyon for ultra-fast Big Data

erns
cover. act.

rns    Contact    Q

‹ Previous    Next ›

ric file system

didn't know Tachyon, you could
an *only* move faster than the
speed of light. *This* Tachyon is part of the Berkeley Data Analytics Stack (**BDAS**), which

Excha

TACHYON
NEXUS

36

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

# Use Case: Baidu

- Framework: SparkSQL
- Under Storage: Baidu's File System
- Storage Media: MEM + HDD
- 100+ nodes deployment
- 1PB+ managed space
- 30x Performance Improvement

**TACHYON**
N E X U S

# Use Case: a SAAS Company

- Framework: Impala

- Under Storage: S3

- Storage Media: MEM + SSD

- 15x Performance Improvement

**TACHYON**
NEXUS

# Use Case: an Oil Company

- Framework: Spark

- Under Storage: GlusterFS

- Storage Media: MEM only

- Analyzing data in traditional storage

**TACHYON**
N E X U S

# Use Case: a SAAS Company

- Framework: Spark

- Under Storage: S3

- Storage Media: SSD only

- Elastic Tachyon deployment

**TACHYON**
N E X U S

# Outline

- Overview
  - Motivation
  - Tachyon Architecture
  - Using Tachyon
- Open Source
  - Status
  - Production Use Cases
- **Roadmap**

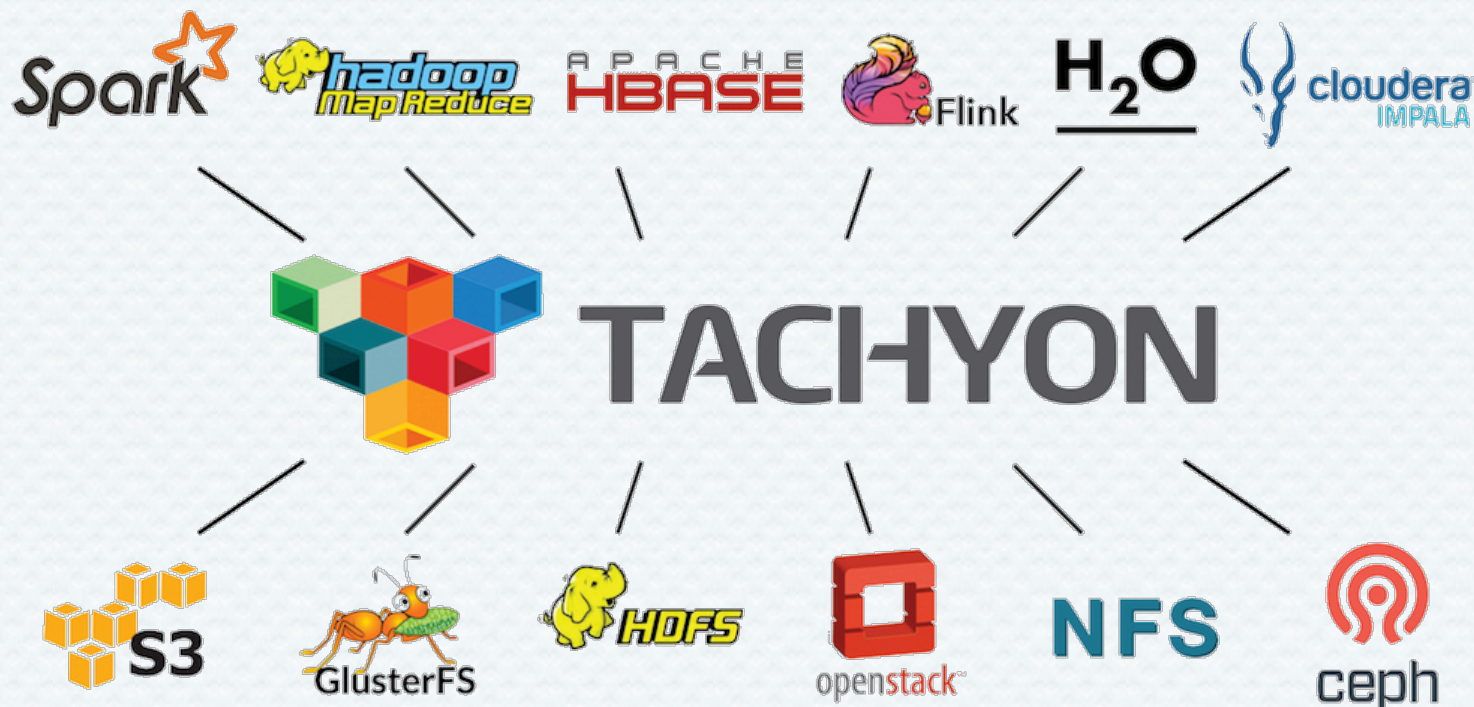**TACHYON**
N E X U S

# New Features

- Lineage in Storage (alpha)
- Tiered Storage (alpha)

**TACHYON**
NEXUS

# New Features

- Lineage in Storage (alpha)
- Tiered Storage (alpha)
- Data Serving
- Support for New Hardware
- …
- Your New Feature!

TACHYON
NEXUS

# Tachyon's Goal?

TACHYON
NEXUS

# Distributed Memory-Centric Storage: Better Assist Other Components



**Welcome Collaboration!**

JIRA New Contributor Tasks

# TACHYON

- Website: http://tachyon-project.org

- Github: https://github.com/amplab/tachyon

- Meetup: http://www.meetup.com/Tachyon

- News Letter Subscription: http://goo.gl/mwB2sX

- Email: haoyuan@tachyonnexus.com

TACHYON
NEXUS