# Tachyon: Reliable File Sharing at Memory-speed Across Cluster Frameworks

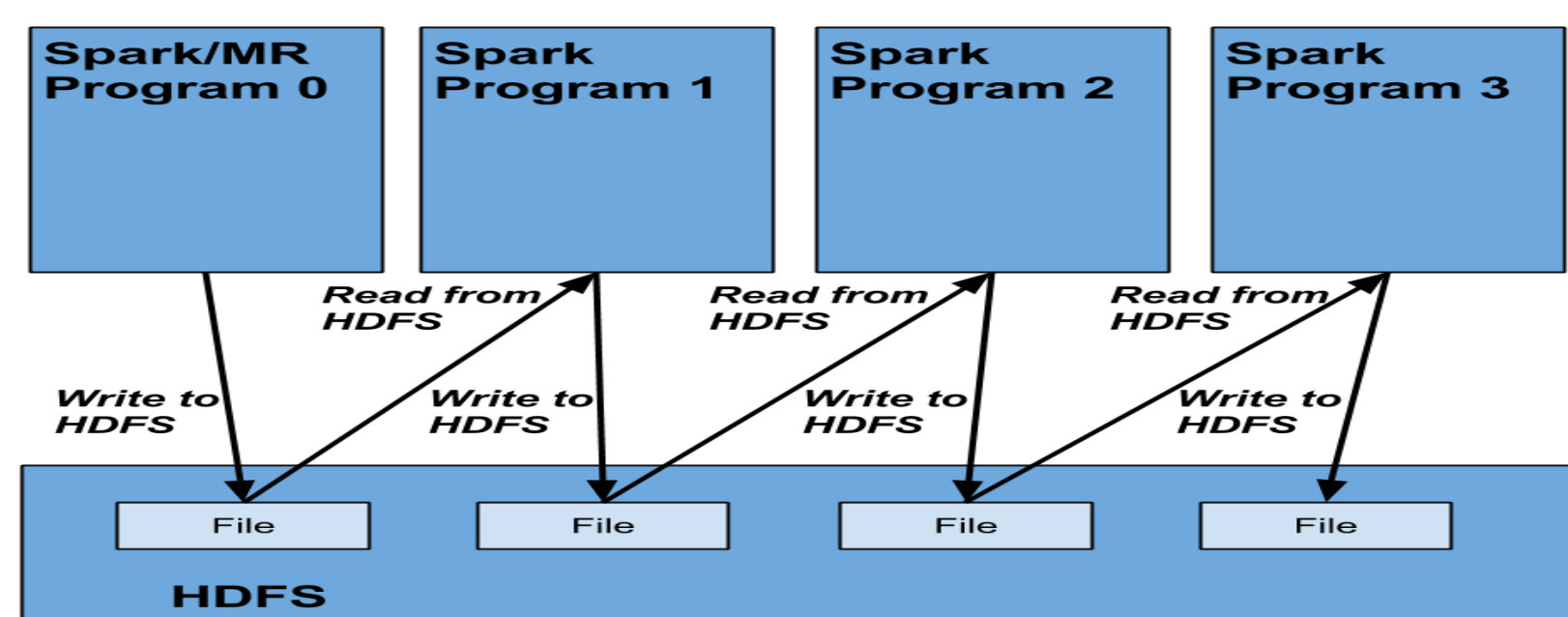Haoyuan Li, Ali Ghodsi, Matei Zaharia, Eric Baldeschwieler, Scott Shenker, Ion Stoica (UC Berkeley)

## Motivation

- Storage throughput has long been the bottleneck for I/O intensive applications.
- Hardware trend suggests existing solutions would not improve significantly.
- This overhead starts to dominate the running times of many applications.
- In data analytical workloads, files are commonly immutable while tasks are stateless and deterministic
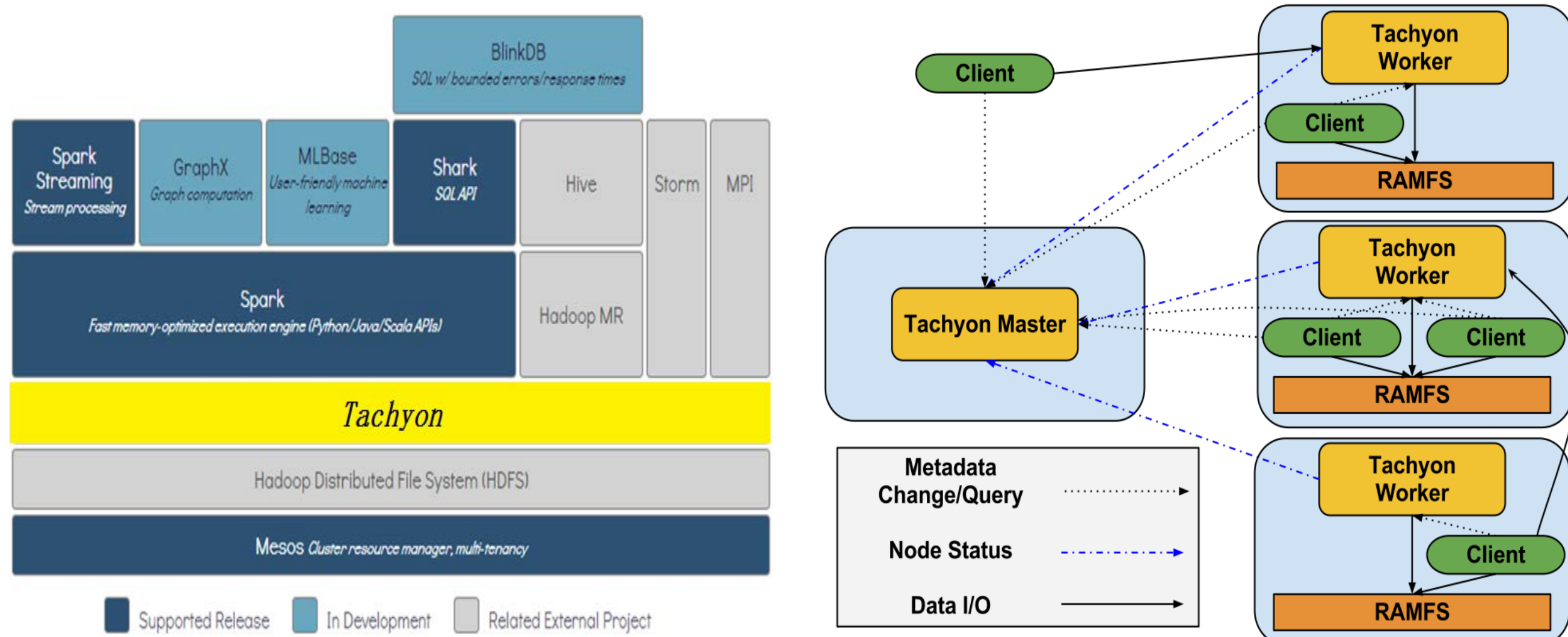
## Existing Systems

- Read: HDFS reads from disk. In memory K/V stores are optimized for small value lookup, but not large sequential read.
- Write: Replication based fault tolerant is the bottleneck.

## Example



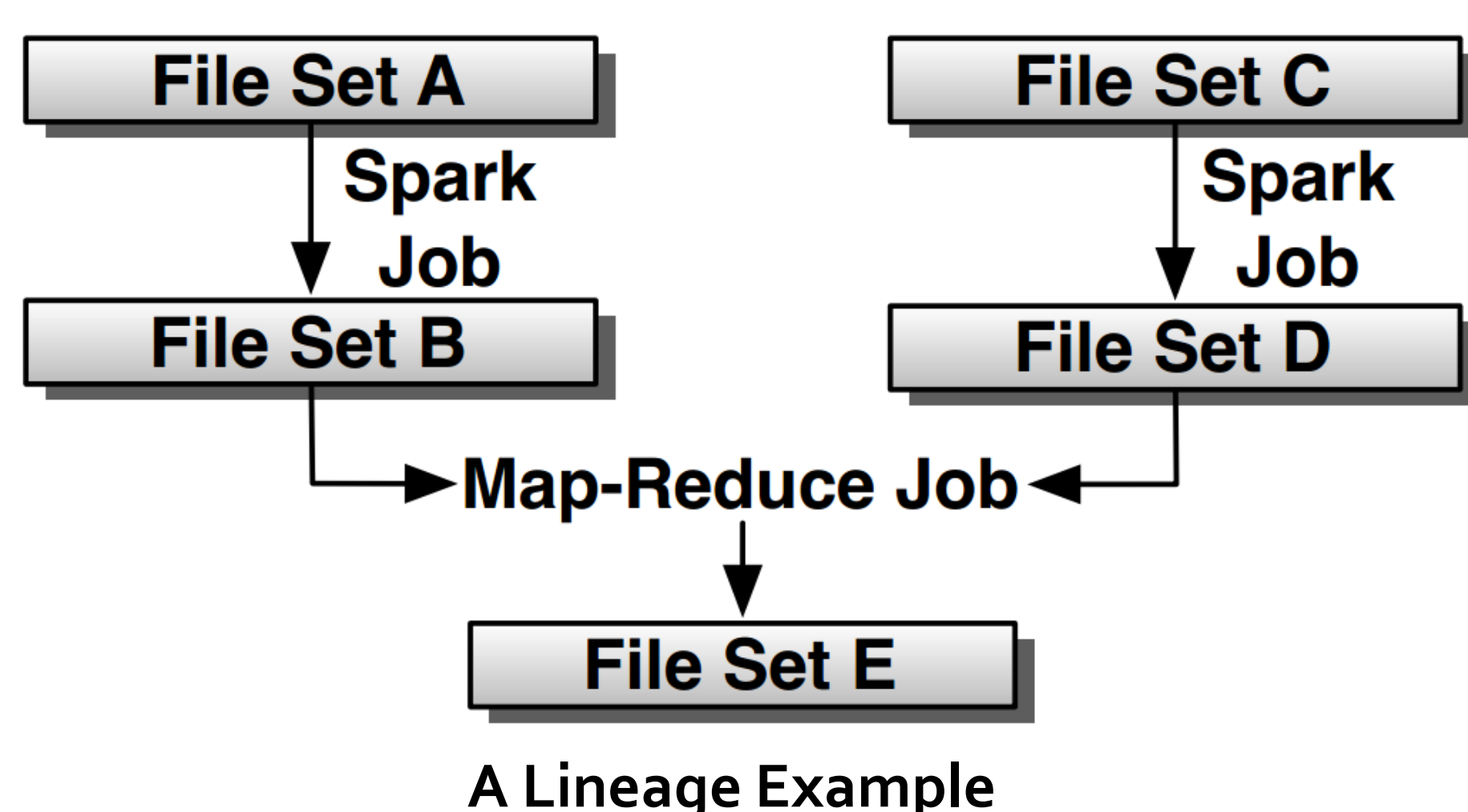**Goal: reliable file sharing at memory speed across cluster frameworks**

## Architecture



## Challenges

*How to achieve reliable file sharing without replication?*

## Recomputation Based Fault Recovery



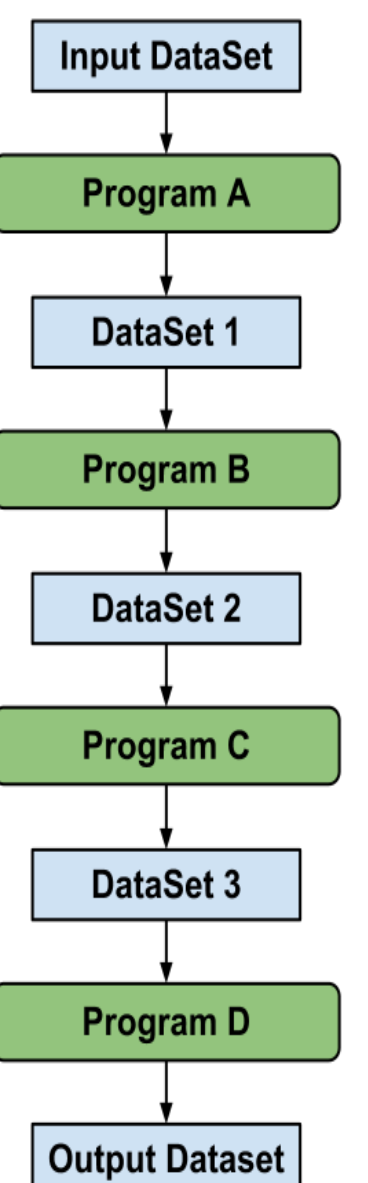A Lineage Example

## Generic Lineage API

Lineage Required Information:
- Input, Output Files
- Binary Program
- Configuration: various configuration formats
- Dependency Type: wide, narrow
- Tachyon provides API to capture above information
- Spark(200-line patch), MapReduce(200-line patch)
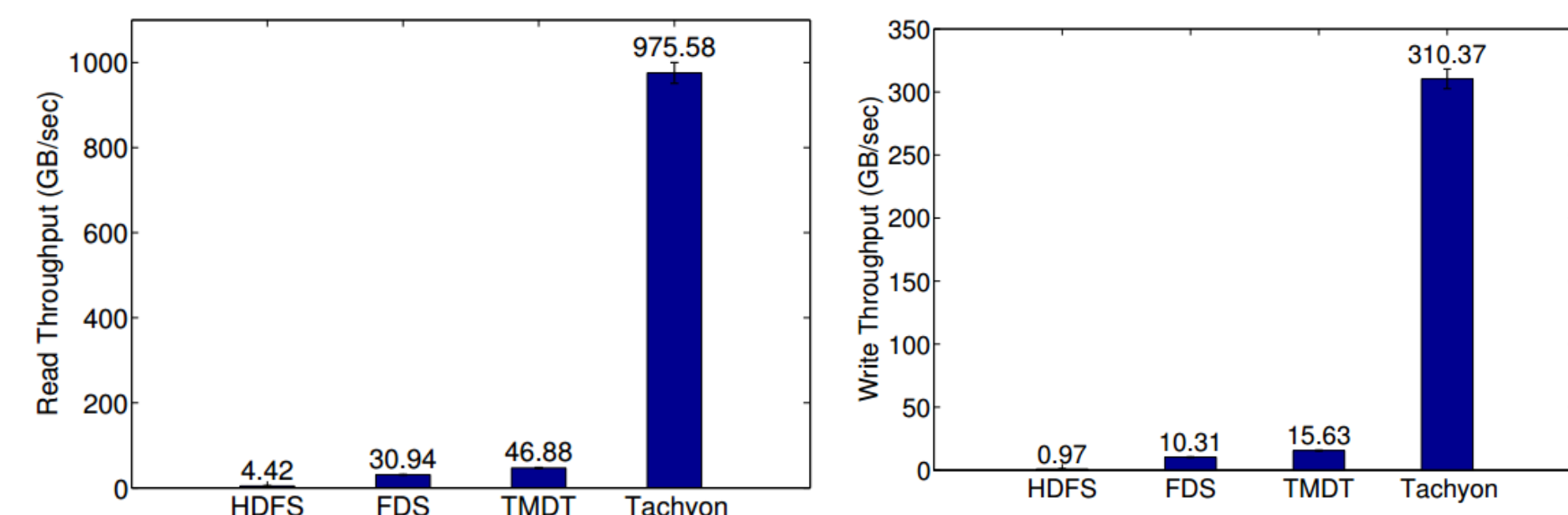- No burden on application programmers
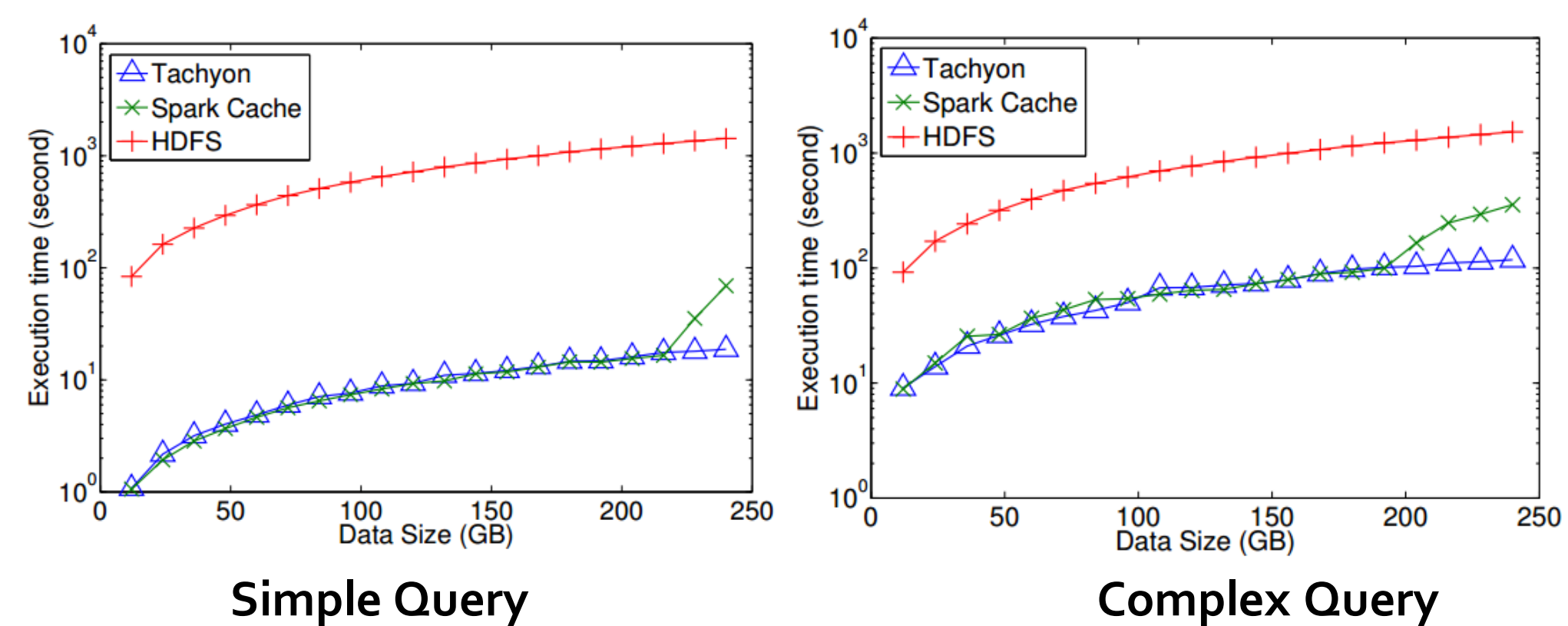
## Recomputation Cost

Checkpoint in the background
- Naïve: checkpoint based on creation order.
- Better: Snapshot checkpointing.
- Even under Naïve algorithm, performance is better than HDFS under failures.
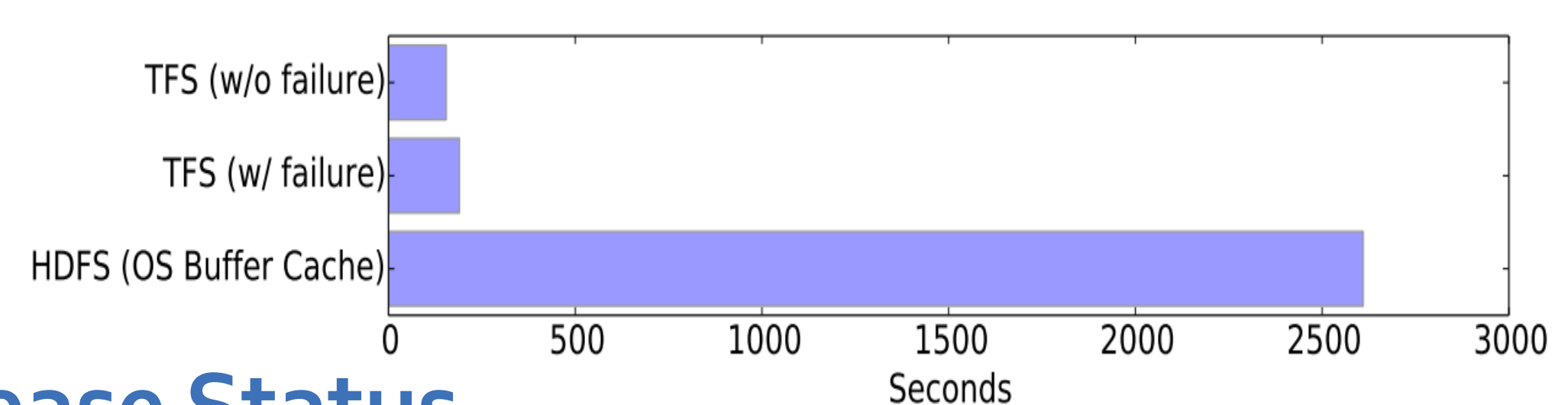- Under snapshot checkpointing, the recovery time is bounded.



## Read/Write Throughput



## Real Workload



Simple Query    Complex Query

## Realistic Workflow



## Release Status

- Developer Preview v0.3.0 (October 2013):
  - First read of files cached in memory
  - Write synchronously to UnderlayerFS (no lineage)
  - Spark and MapReduce can use it without changing any code:
- Current Features: Java-like file API; Compatible with Hadoop; Master fault tolerance, Native support for raw table; Pluggable underlayer file system; CLI; Web user interface; Whitelist, Pinlist.

## Conclusion

- Tachyon pushes lineage into file system layer to enable memory throughput read and write