# Improving Data Locality for Analytics Jobs on Kubernetes Using Alluxio

*Adit Madan | Alluxio*
*Gene Pang | Alluxio*

ALLUXIO

# Outline

Alluxio Overview
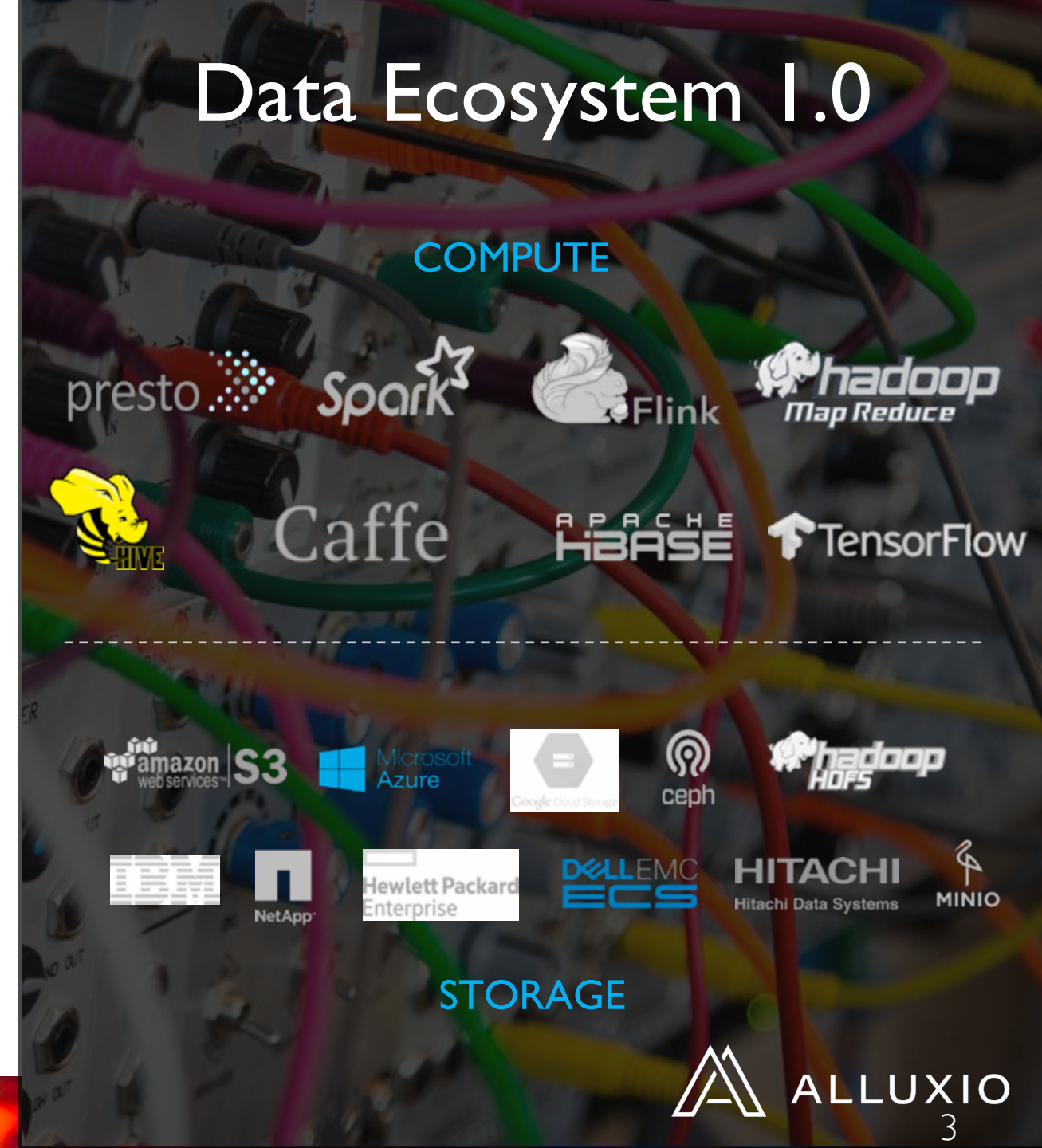
Data locality with Spark and Alluxio

Kubernetes Overview

Spark and Alluxio in Kubernetes

Alluxio Innovations for Structured Data

ALLUXIO

# Data Ecosystem - *Beta*

**COMPUTE**

hadoop Map Reduce

hadoop HDFS

**STORAGE**

# Data Ecosystem 1.0

COMPUTE

presto   Spark   Flink   hadoop Map Reduce

HIVE   Caffe   APACHE HBASE   TensorFlow

amazon web services | S3   Microsoft Azure   Google Cloud Storage   ceph   hadoop HDFS

IBM   NetApp   Hewlett Packard Enterprise   DELL EMC ECS   HITACHI Hitachi Data Systems   MINIO

STORAGE

ALLUXIO

ALLUXIO

3

# Data stack journey and innovation paths

**Co-located**

Co-located
compute & HDFS
on the same cluster

| MR / Hive |
|-----------|
| HDFS |

- Typically compute-bound clusters over 100% capacity
- Compute & I/O need to be scaled together even when not needed

**Disaggregated**

Disaggregated
compute & HDFS
on the same cluster

| Hive |
|------|
| HDFS |

- Compute & I/O can be scaled independently but I/O still needed on HDFS which is expensive

**HDFS for Hybrid Cloud**

Burst HDFS data in
the cloud,
public or private

**Support more frameworks**

Support Presto, Spark
and other computes
without app changes

**Transition to Object store**

Enable & accelerate
big data on
object stores

ALLUXIO

# Independent scaling of compute & storage

| Java File API | HDFS Interface | S3 Interface | POSIX Interface | REST API |

**ALLUXIO** Data Orchestration for the Cloud

| HDFS Driver | Swift Driver | S3 Driver | NFS Driver |

# Alluxio Data Orchestration for the Cloud

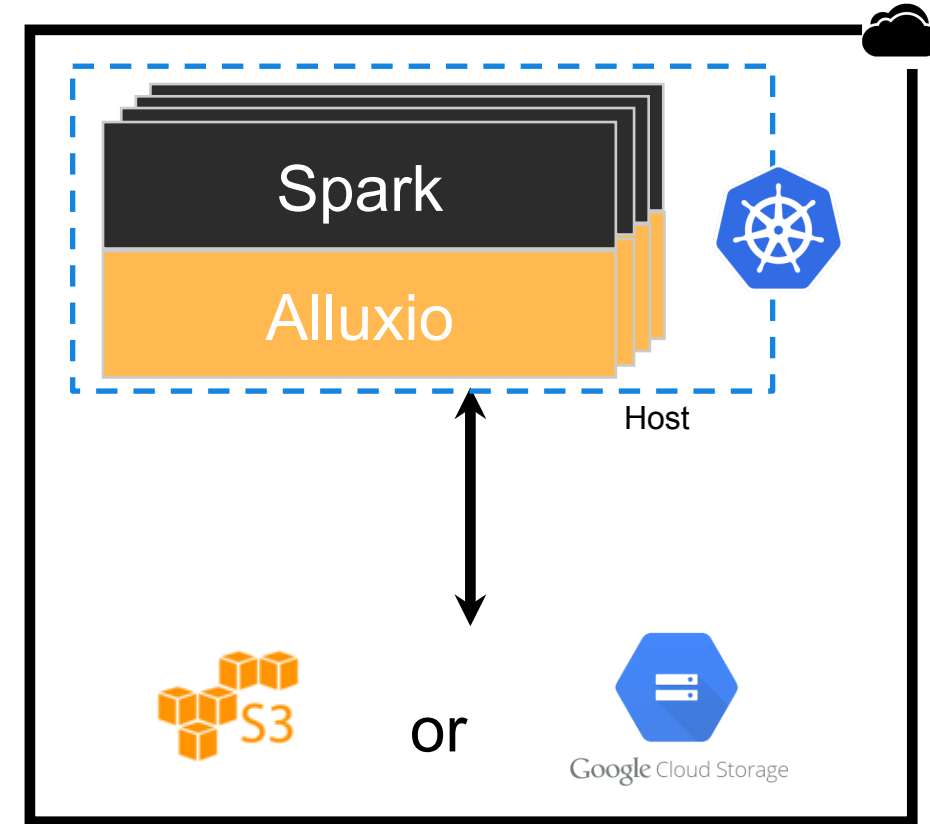| Global Namespace | Intelligent Caching | Data Management | Structured Data Catalog | Data Transformation |

# Why Alluxio in K8s?

Elastic Data for Elastic Compute

- Improve Data locality
  - Big data analytics or ML
  - Cache data close to compute
- Enable high-speed data sharing across jobs
  - A staging storage layer
- Unification of persistent storage
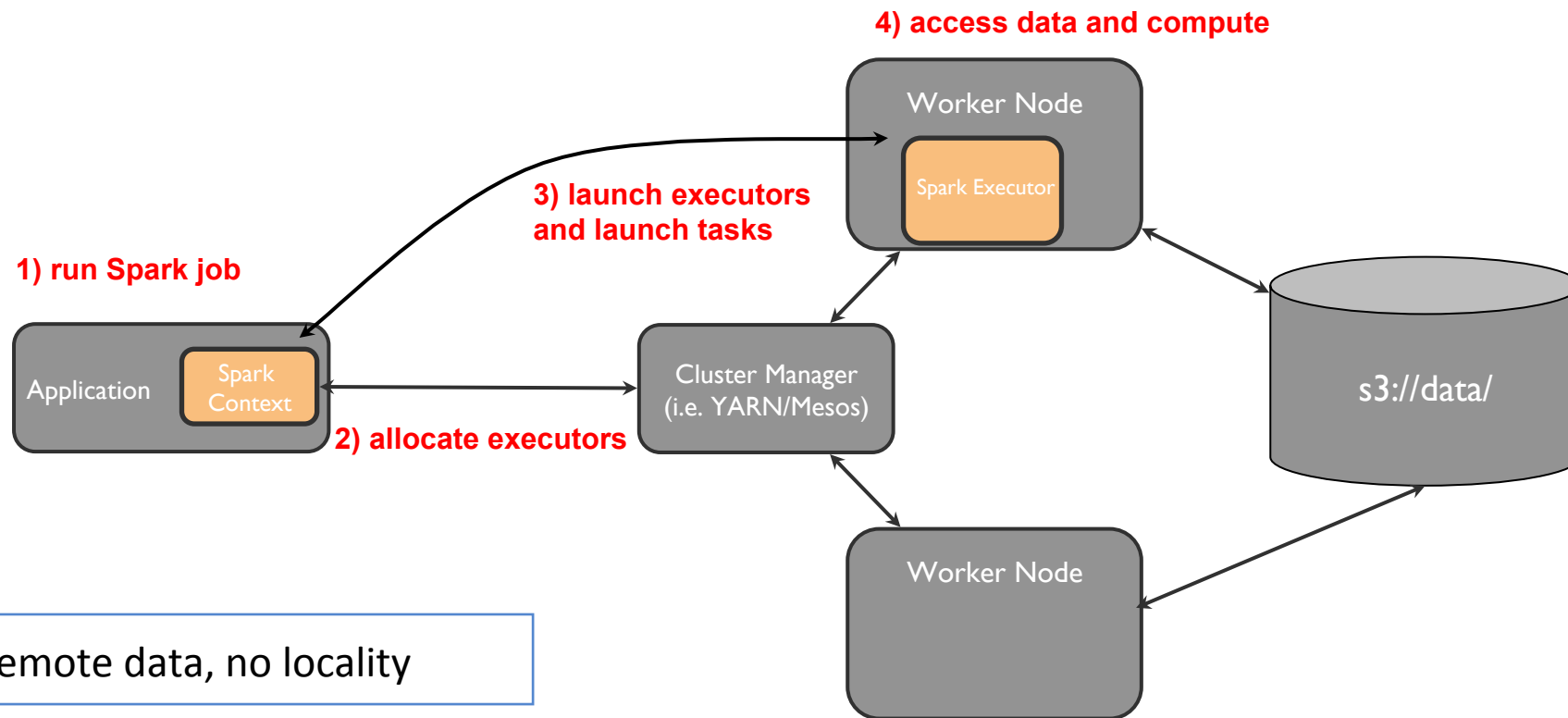  - Data abstraction across different storage

Read more at https://www.alluxio.io/blog/kubernetes-alluxio-and-the-disaggregated-analytics-stack/
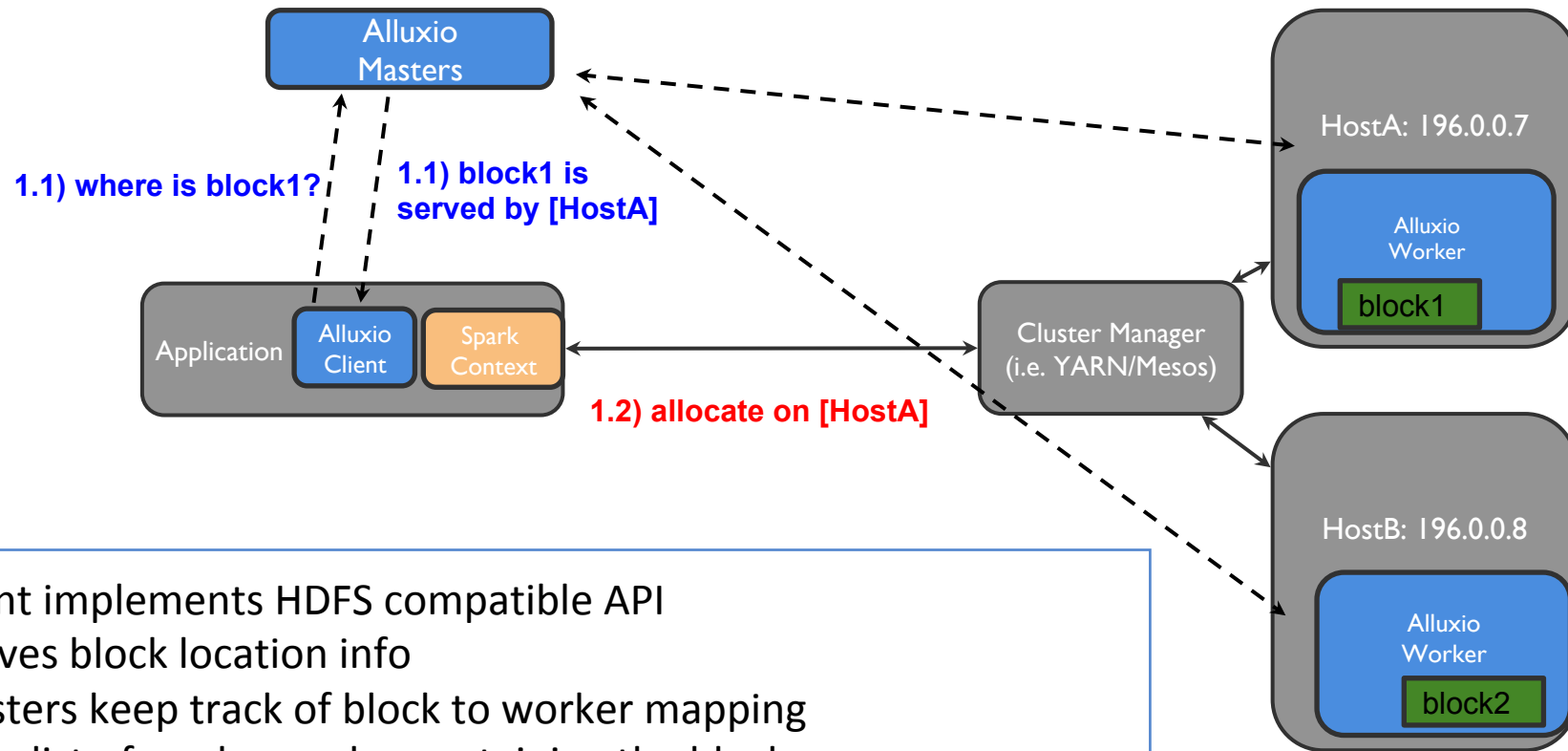
# Spark + Alluxio Data Locality

Without K8s

# Spark Workflow (w/o Alluxio)

**4) access data and compute**

Worker Node

Spark Executor

**3) launch executors and launch tasks**

**1) run Spark job**

Application

Spark Context

**2) allocate executors**

Cluster Manager (i.e. YARN/Mesos)
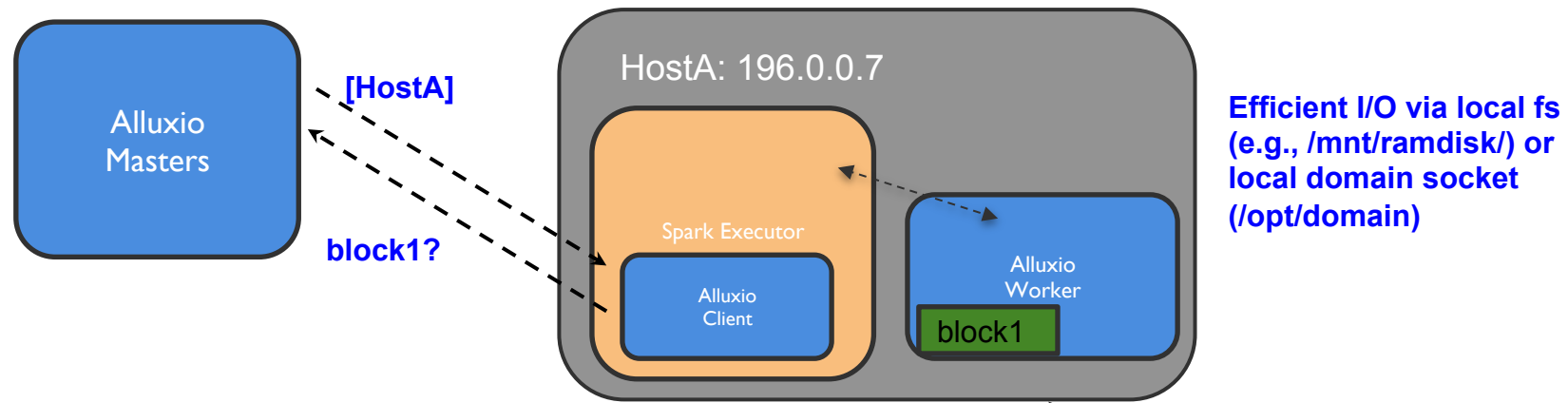
s3://data/

Worker Node

Takeaway: Remote data, no locality

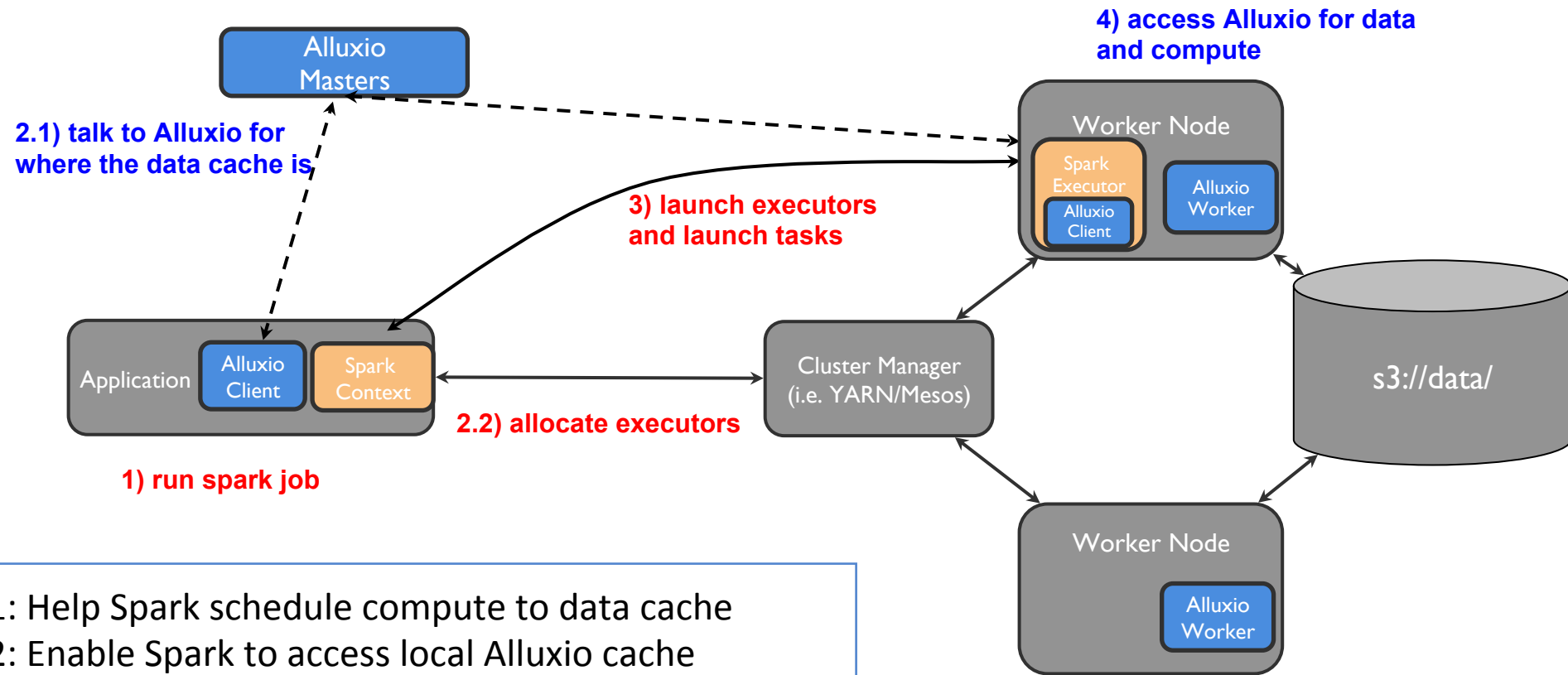# Step I: Schedule Compute to Data Location



Alluxio client implements HDFS compatible API
    Retrieves block location info
Alluxio masters keep track of block to worker mapping
    Serve a list of worker nodes containing the block

# Step 2: Detect+Exchange Data w/ local Worker

Alluxio Masters

[HostA]

block1?

HostA: 196.0.0.7

Spark Executor

Alluxio Client

Alluxio Worker

block1

Efficient I/O via local fs (e.g., /mnt/ramdisk/) or local domain socket (/opt/domain)

Spark Executor finds local Alluxio Worker
    Mechanism: Hostname comparison
Spark Executor talks to local Alluxio Worker
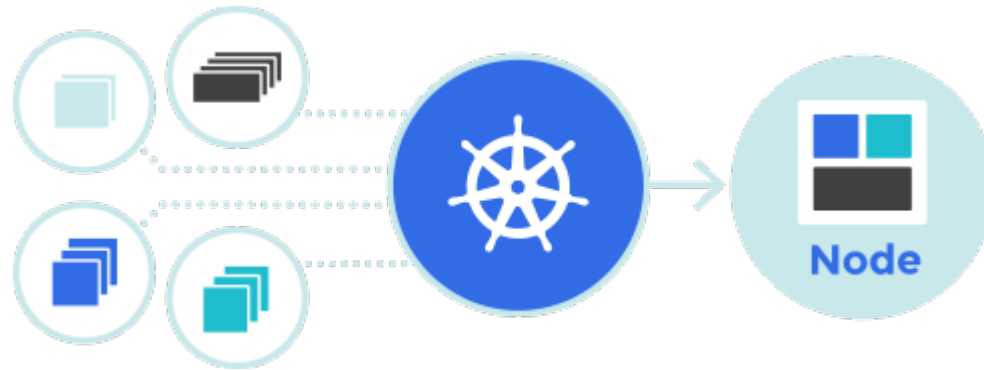    Mechanism: Either short-circuit access (via local FS) or local domain socket

HostB: 196.0.0.8

Alluxio Worker

s3://data/

# Recap: Spark Architecture w/ Alluxio



**4) access Alluxio for data and compute**

**2.1) talk to Alluxio for where the data cache is**

Alluxio Masters

Worker Node
- Spark Executor
- Alluxio Client
- Alluxio Worker

**3) launch executors and launch tasks**

Application
- Alluxio Client
- Spark Context

Cluster Manager (i.e. YARN/Mesos)

**2.2) allocate executors**

**1) run spark job**

s3://data/

Worker Node
- Alluxio Worker

Step 1: Help Spark schedule compute to data cache
Step 2: Enable Spark to access local Alluxio cache

# Kubernetes Overview

# Kubernetes (K8s) is…

"an open-source container-orchestration system for automating application deployment, scaling, and management."

# Container Orchestration

- Platform agnostic cluster management

- Service discovery and load balancing

- Storage orchestration

- Horizontal scaling

- Self-monitoring and self-healing

- Automated rollouts and rollbacks

# Key K8s Terms

- Node
  - A VM or physical machine
- Container
  - Container = Image once running on Docker Engine
- Pod
  - Schedulable unit of one or more containers running together
- Controller
  - Controls the desired state such as copies of a Pod
- DaemonSet
  - A Controller that ensures each Node has only one such Pod
- Persistent Volume
  - A storage resource with lifecycle independent of Pods
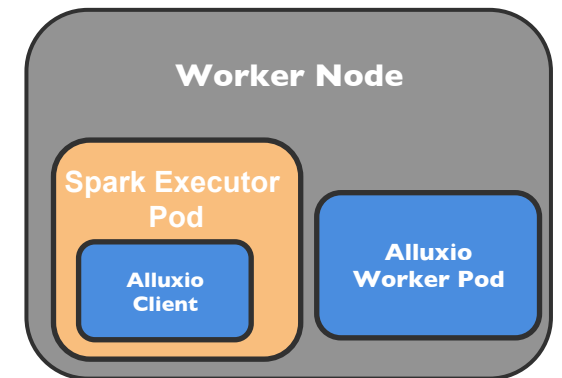
# Spark + Alluxio Data Locality

In K8s environment

# Spark on K8s Architecture

- Spark 2.3 added native K8s support

- *spark-submit* talks to API Server to launch Driver

- Spark Driver launches Executor Pods

- When the application completes,

  - Executor Pods terminate and are cleaned up

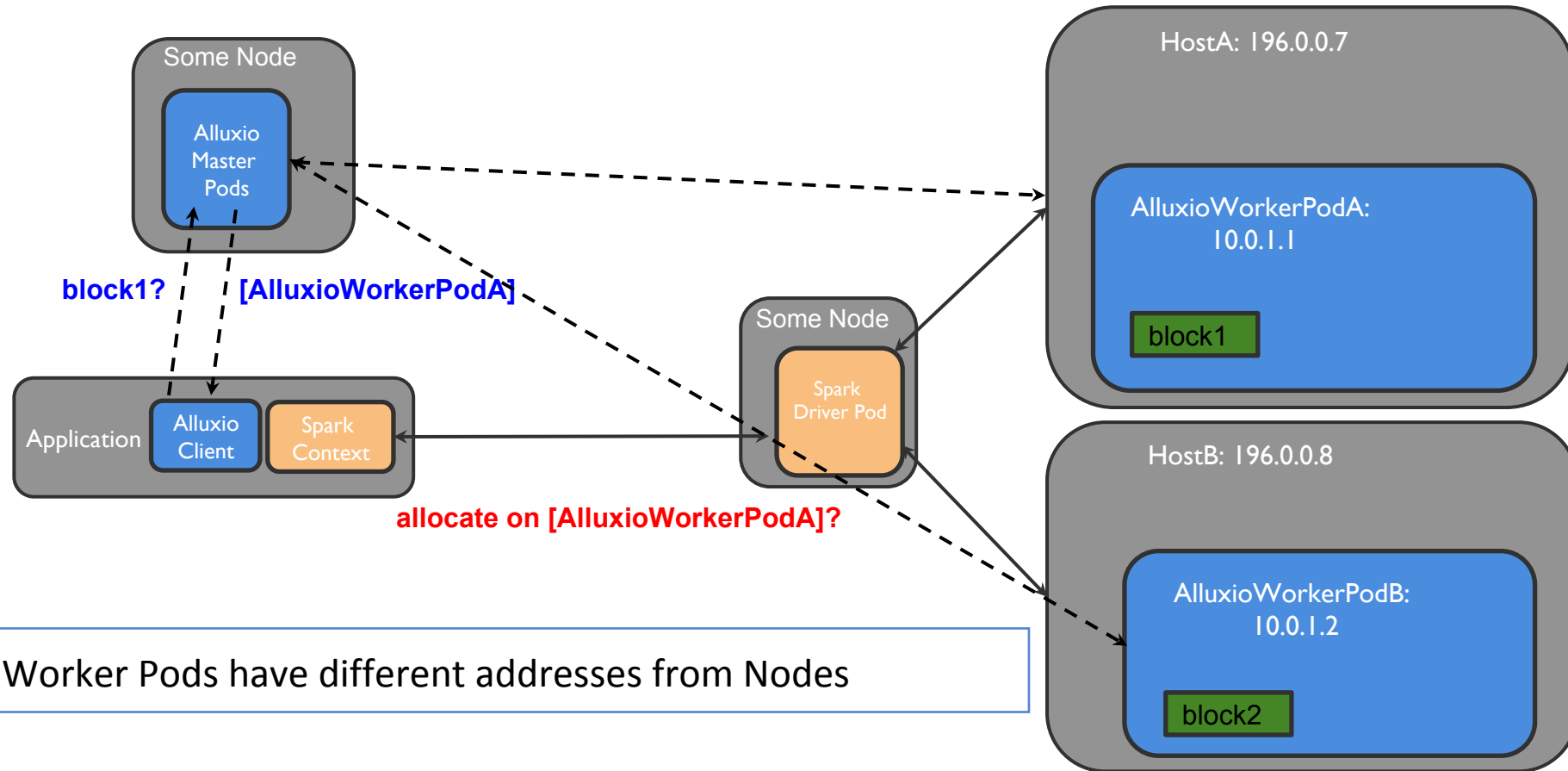  - Driver Pod persists logs and remains in "COMPLETED" state

# Deployment Model: Co-location

- Co-locate Spark Executor Pod w/ Alluxio Worker Pod
- Lifecycle
  - Spark Executors are ephemeral
  - Alluxio Workers persist across all Spark jobs
- Deployment order:
  - Deploy Alluxio cluster first (masters+workers)
  - An Alluxio Worker on each Node, by DaemonSet
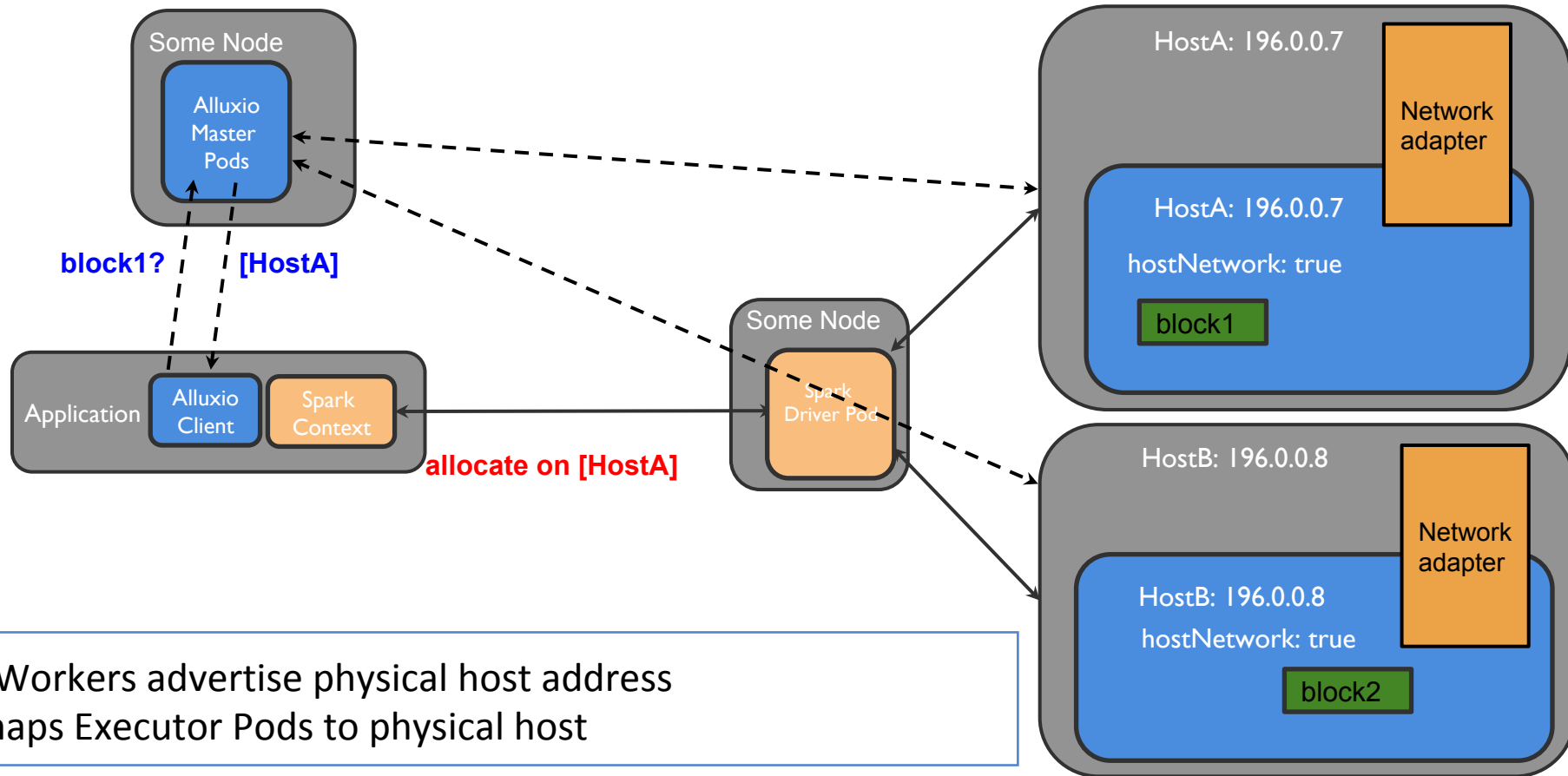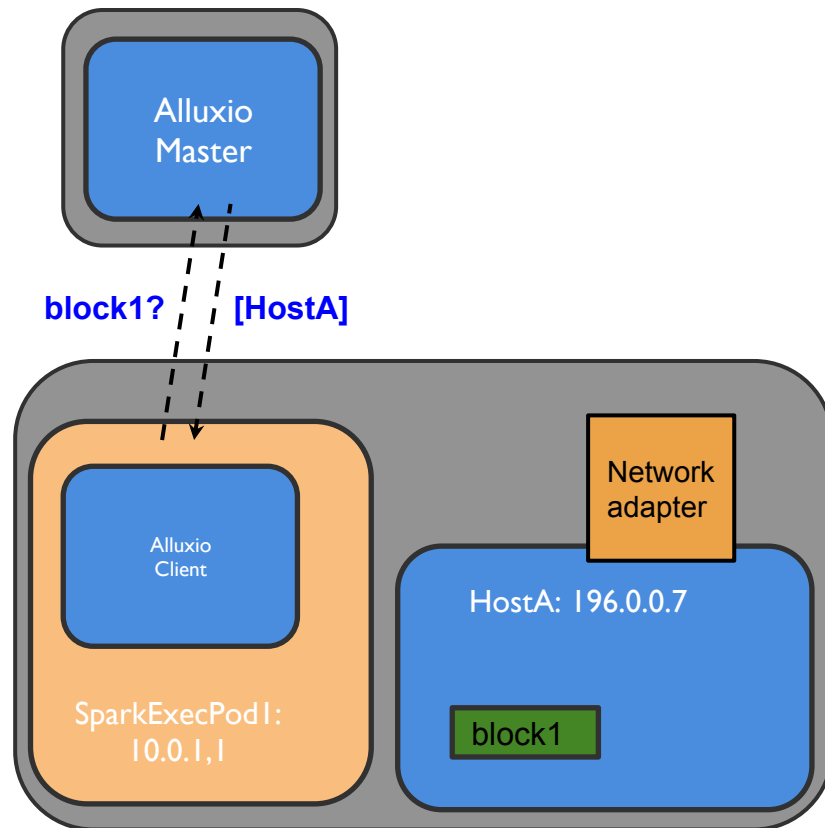  - *spark-submit* launches Spark Driver + Executors

**Worker Node**

**Spark Executor Pod**

**Alluxio Client**

**Alluxio Worker Pod**

# Challenge 1: Executor Allocation to Workers



Problem: Worker Pods have different addresses from Nodes

# Solution: Alluxio Workers w/ *hostNetwork*



Alluxio Workers advertise physical host address
Spark maps Executor Pods to physical host

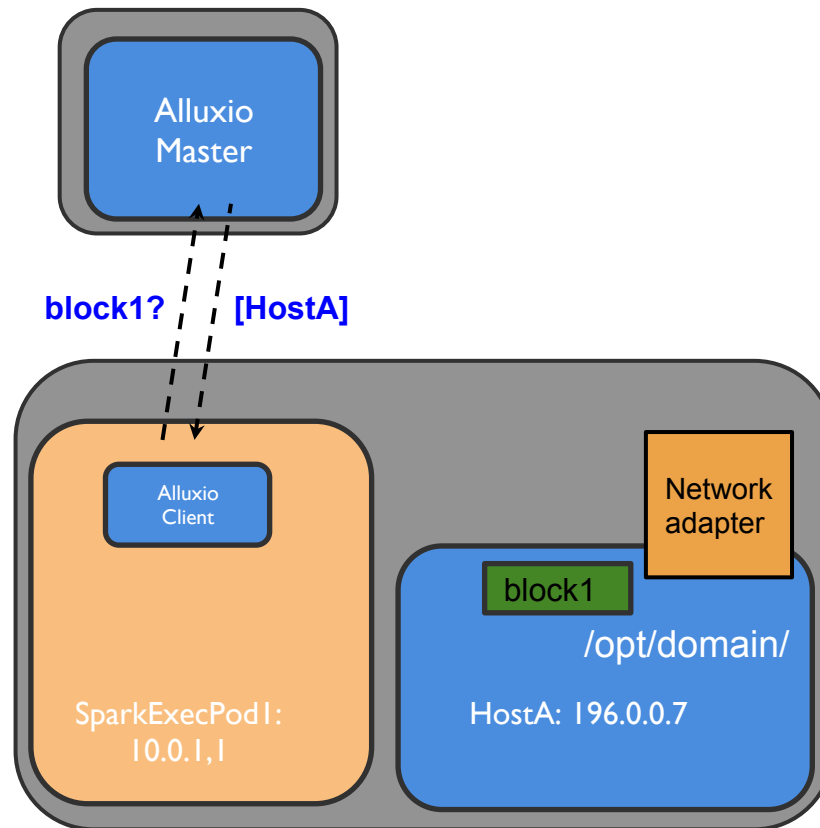# Challenge 2: Identify host-local Alluxio Pod



Problem:
Spark Executor Pod has different hostName
     Hostname HostA != SparkExecPod1
     Local Alluxio Worker not identified!
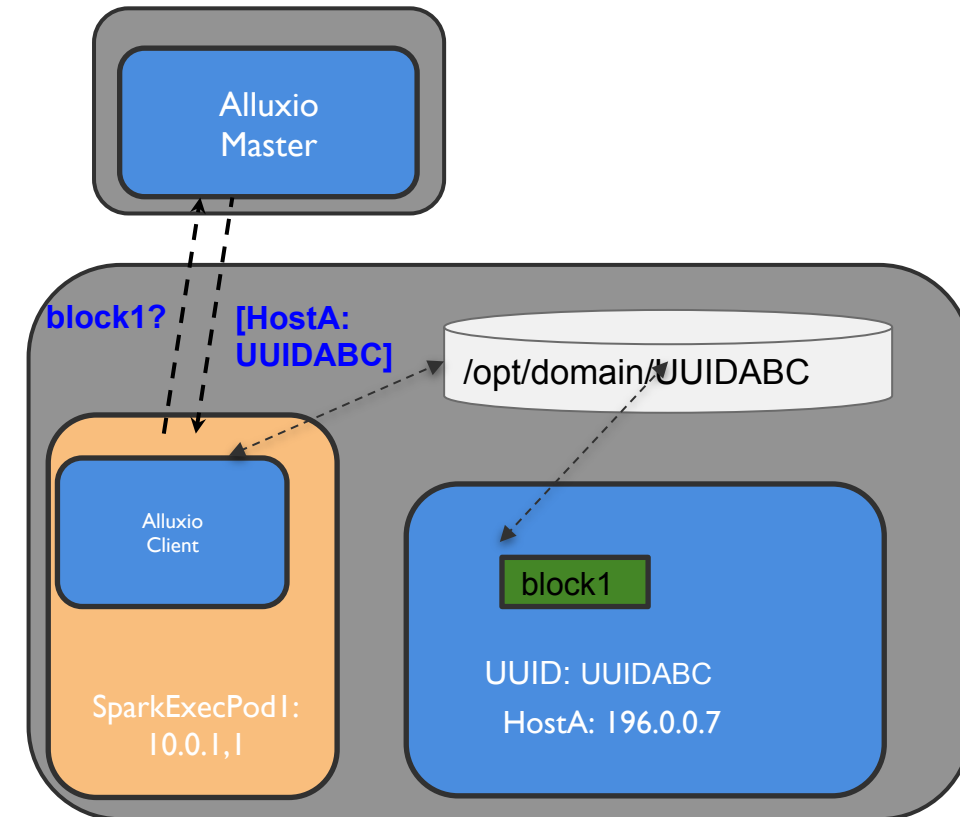
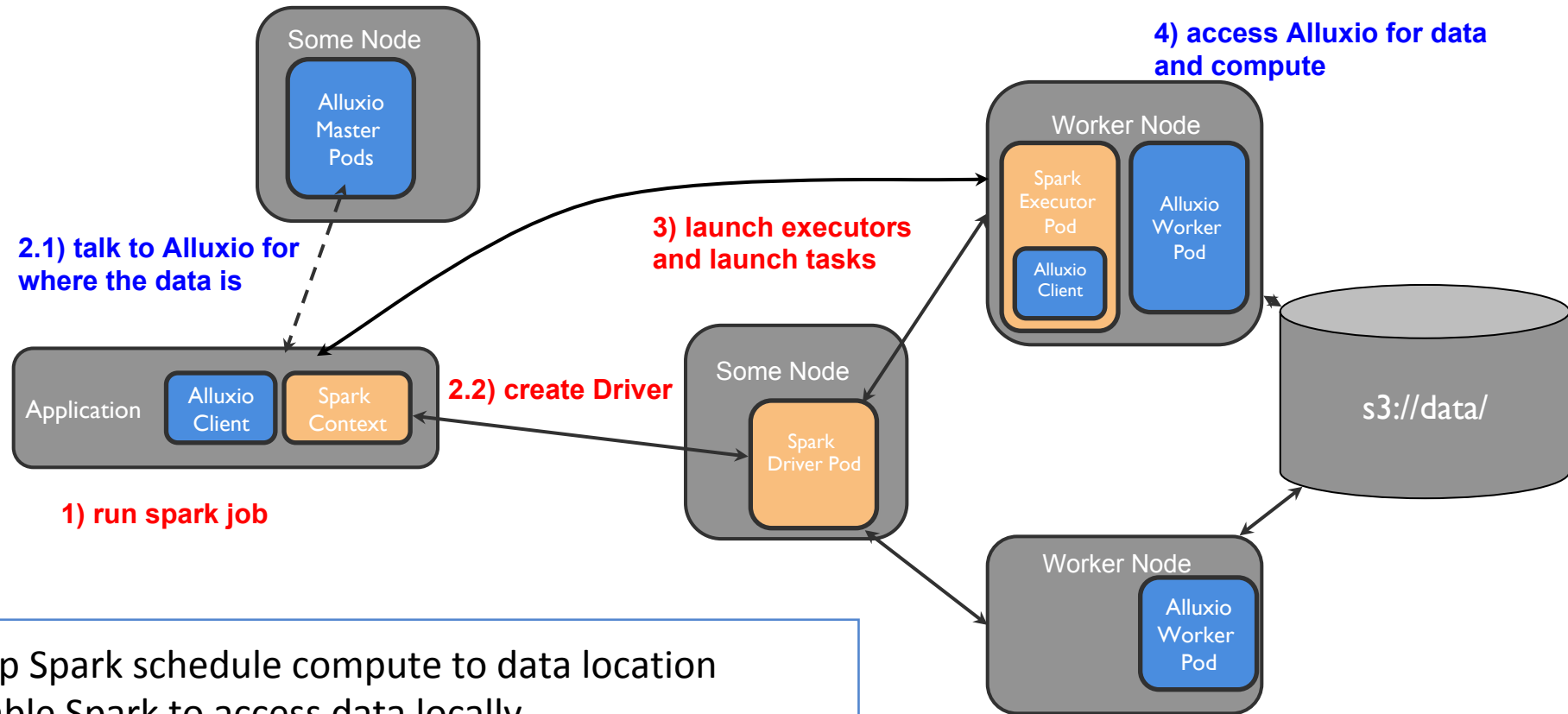# Challenge 3: Executor fails to find domain socket



Problem:
- Pods don't share the File System
- Domain socket */opt/domain* is in Alluxio Worker Pod

Alluxio Master

block1?    [HostA]

Alluxio Client

Network adapter

block1

/opt/domain/

SparkExecPod1: 10.0.1,1

HostA: 196.0.0.7

# Solution: Share a *hostPath* Volume b/w Pods

- Each Alluxio Worker has a UUID
- Share domain socket by a *hostPath* Volume
- Alluxio Client finds local worker's domain socket by finding file matching Worker UUID
- Worker domain socket path
  - */opt/domain/d -> /opt/domain/UUIDABC*
- Mount hostPath Volume to Spark Executor
  - Enabled in Spark 2.4

Alluxio Master

block1?   [HostA: UUIDABC]

/opt/domain/UUIDABC

Alluxio Client

block1

UUID: UUIDABC
HostA: 196.0.0.7

SparkExecPod1: 10.0.1,1

# Recap: Spark + Alluxio Architecture on K8s



Step 1: Help Spark schedule compute to data location
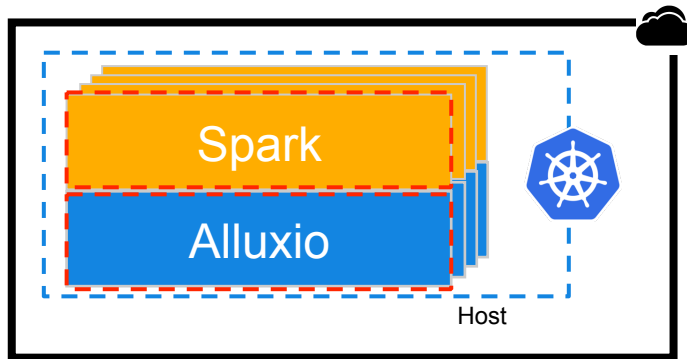Step 2: Enable Spark to access data locally

# Limitations and Future Work

- Enterprise environments may restrict *hostNetwork* and *hostPath*

- Alluxio workers need *hostNetwork*

  - Plan: Support container network translation

- The domain socket file requires a *hostPath* volume

  - Plan: Using Local Persistent Volumes

- Feedback/collaboration are welcome!
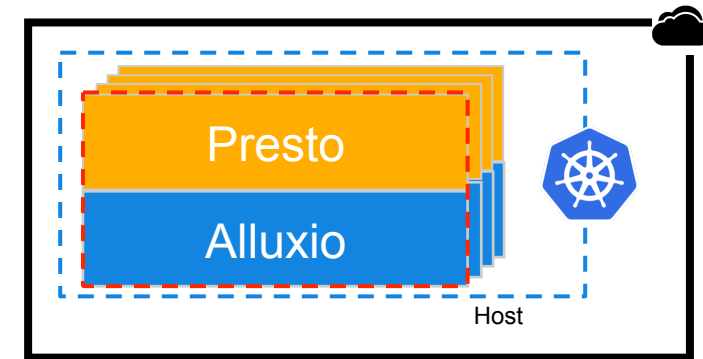
# Alternate Deployment Options

# Deploying Alluxio in K8s

**Spark**

**Alluxio**

Host

**Presto**

**Alluxio**

Host

Legend:     Pod

Alluxio and Compute in different pods on the same host

When do you use this?
- Compute, like Spark, is short running and ephemeral
- Alluxio data orchestration & access layer is long running and used across many jobs

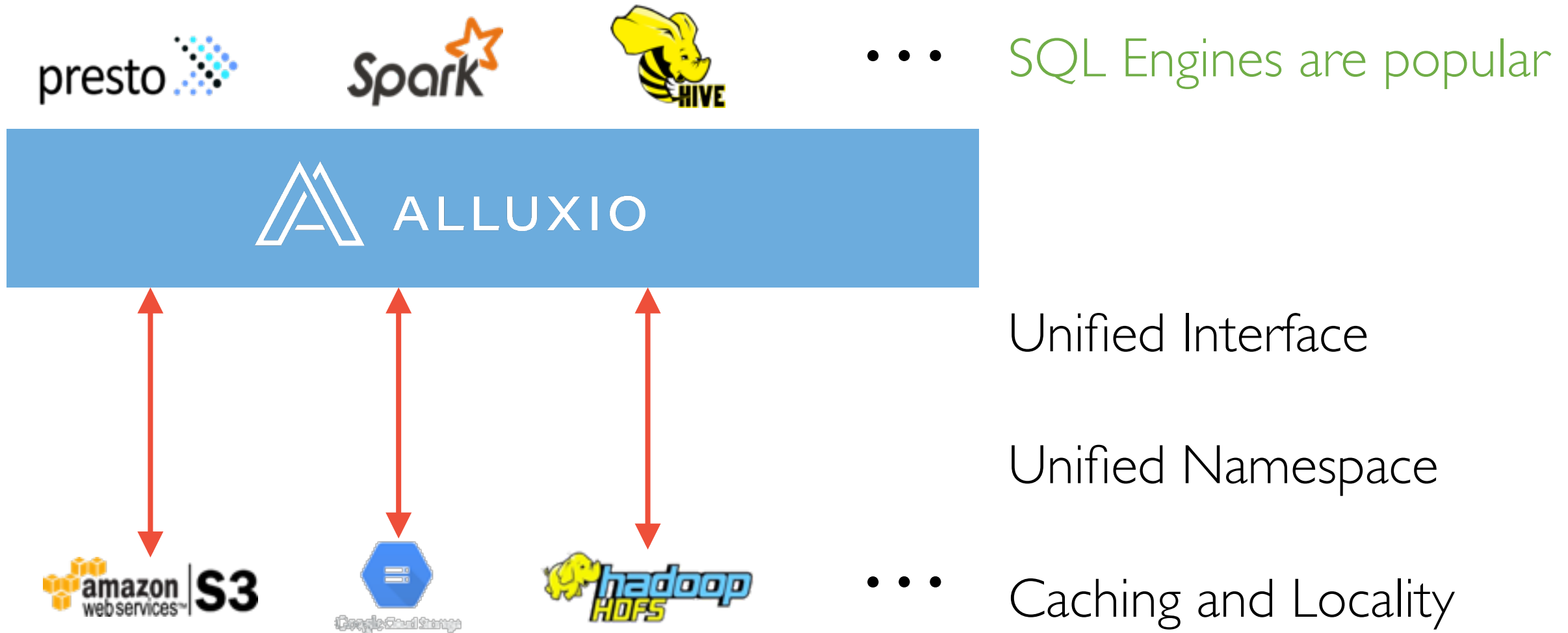Alluxio and Compute framework in the same pod

When do you use this?
- Compute, like Presto, is long running
- Data tier with Alluxio needs to be scaled along with compute tier

# Alluxio
# Structured Data Management

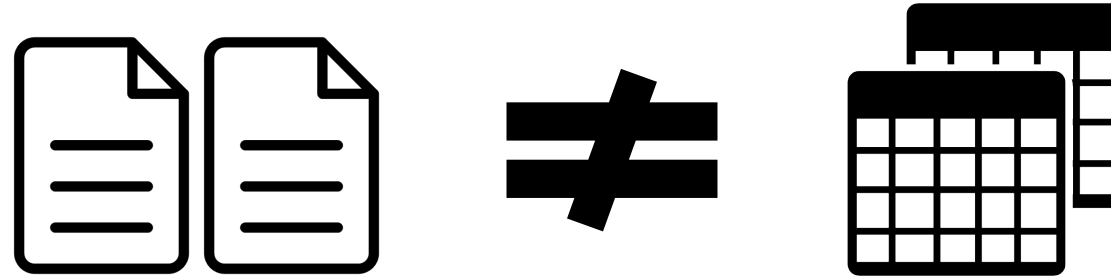## Innovations for Structured Data

# Common Alluxio Use Cases



SQL Engines are popular

Unified Interface

Unified Namespace

Caching and Locality

# Storage Systems                                    SQL Frameworks

Files/Objects                                                Tables

Directories                                                Schemas

Raw Bytes                                              Rows/Columns

**Impedance Mismatch**

**Storage Optimized**                          **Compute Optimized**

**Further Expand Benefits!**

# Benefits of Alluxio Data Orchestration

Caching

Unified Interface/Namespace

Schema-Aware Optimizations

Compute-Optimized Formats

Physical Data Independence
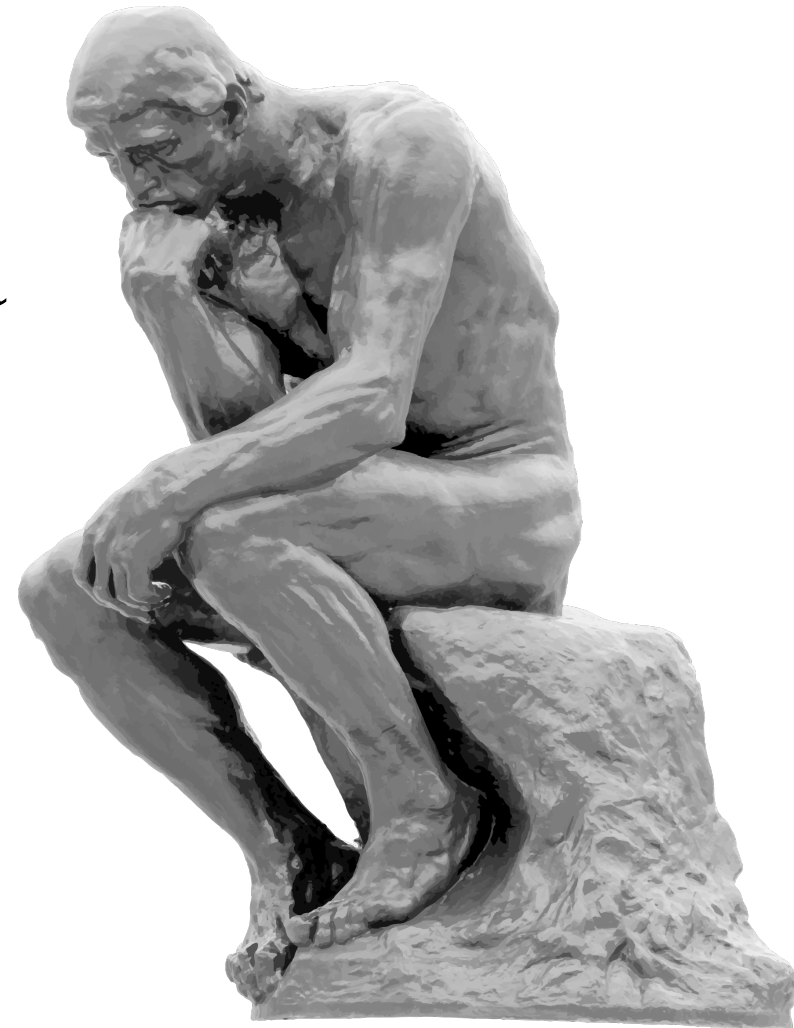
Storage Systems

SQL Frameworks

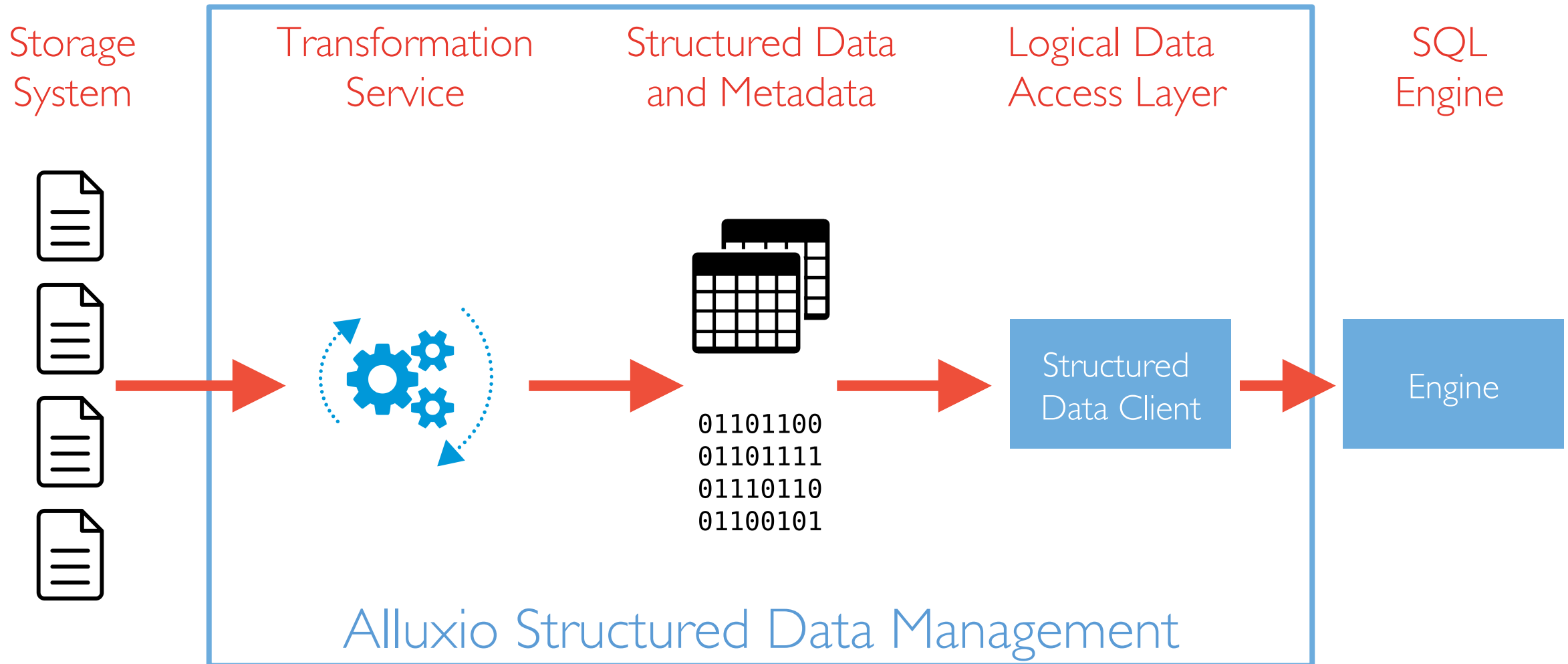# High-Level Philosophy

**Provide Structured Data APIs**

**Focus** on how frameworks interact with data

**Cache Logical Data Access**

**Focus** on caching what frameworks want

# Alluxio Structured Data Management



Storage System → Transformation Service → Structured Data and Metadata → Logical Data Access Layer → SQL Engine

01101100
01101111
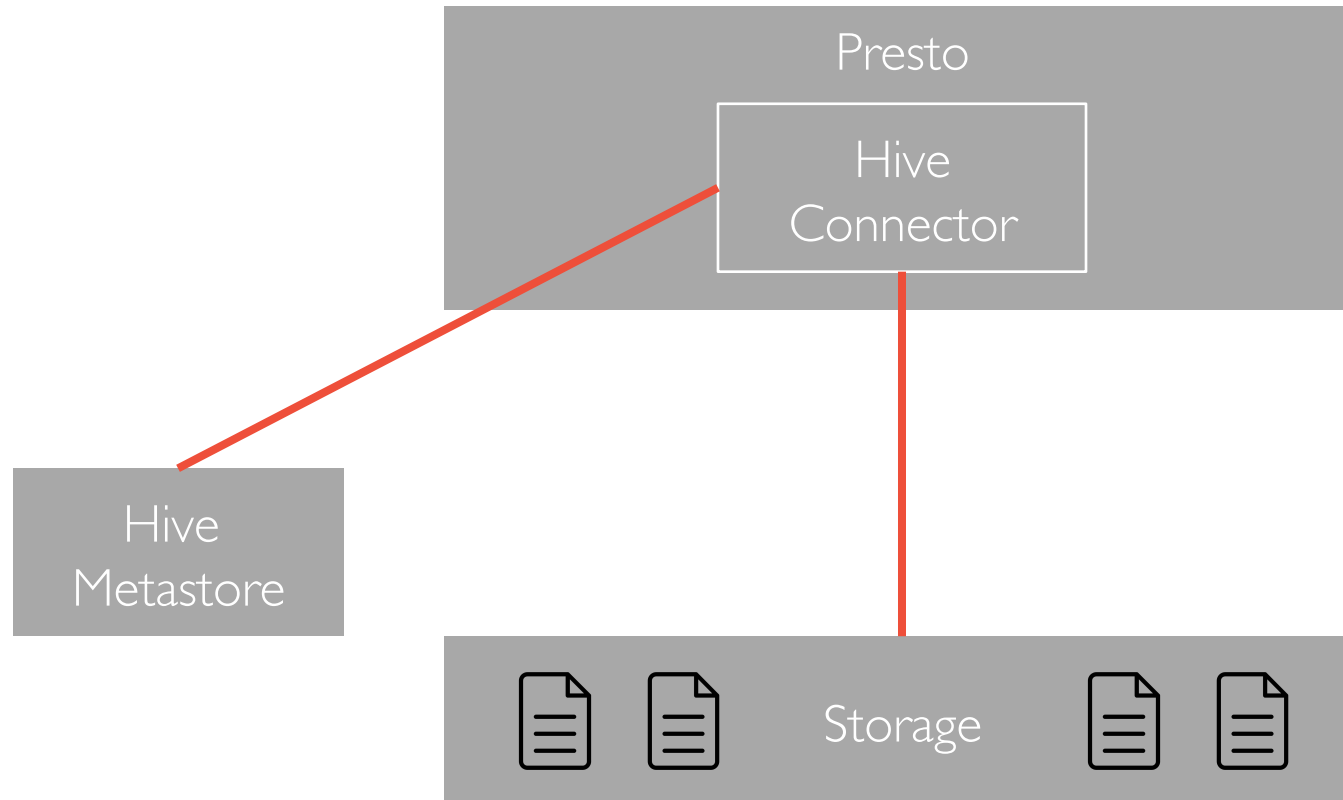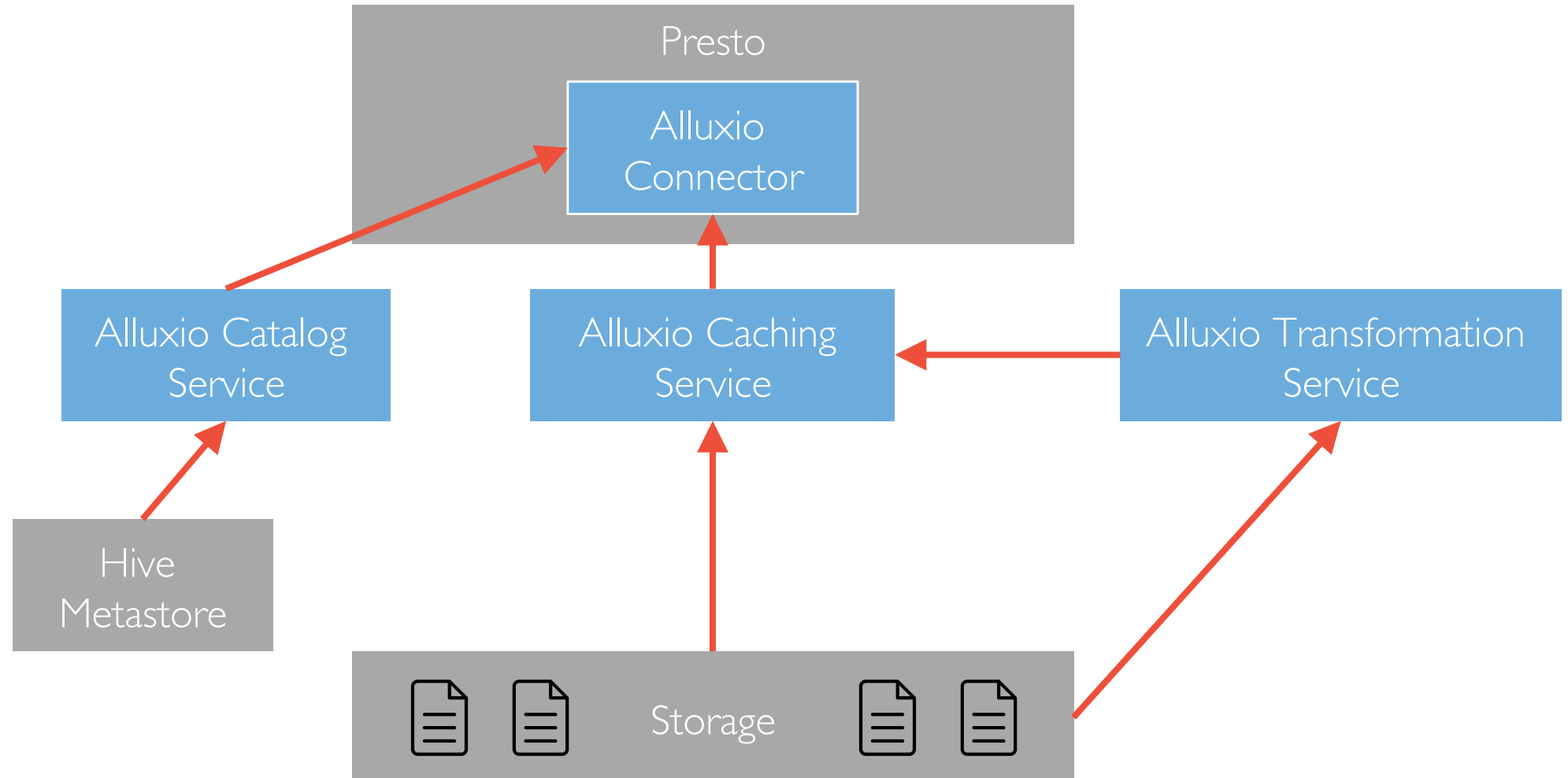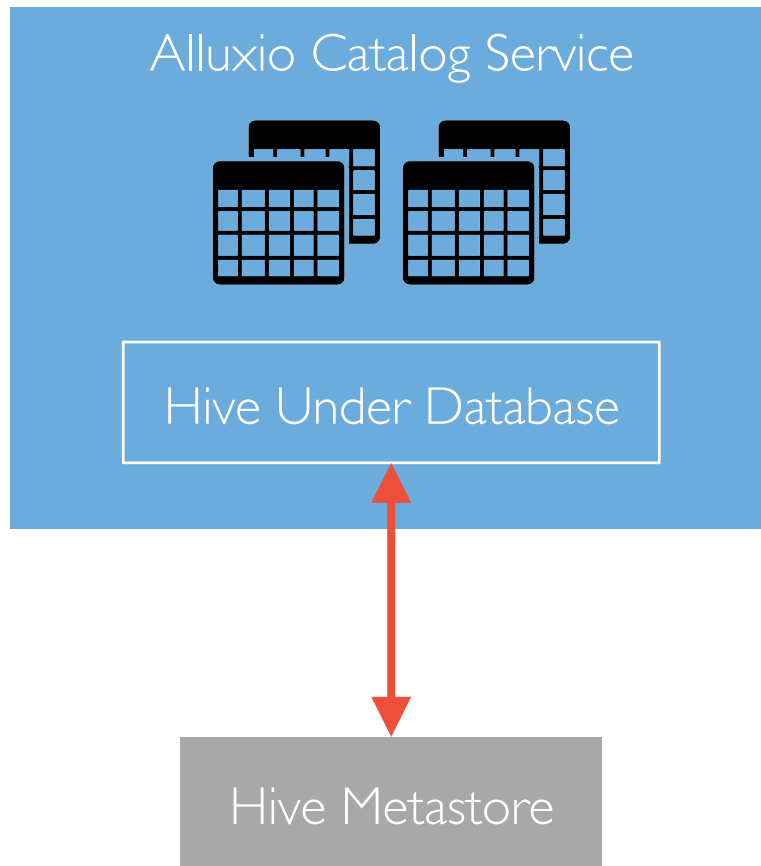01110110
01100101

Structured Data Client

Engine

Alluxio Structured Data Management

# Developer Preview in Alluxio 2.1

# Target Environment

# Alluxio Structured Data Management

# Alluxio Catalog Service

Alluxio Catalog Service

Hive Under Database

Hive Metastore

## Functionality

Manages metadata for structured data

Abstracts other database catalogs as Under Database (UDB)

## Benefits

Schema-aware optimizations

Simple deployment

# Alluxio Presto Connector

Tighter integration with Presto
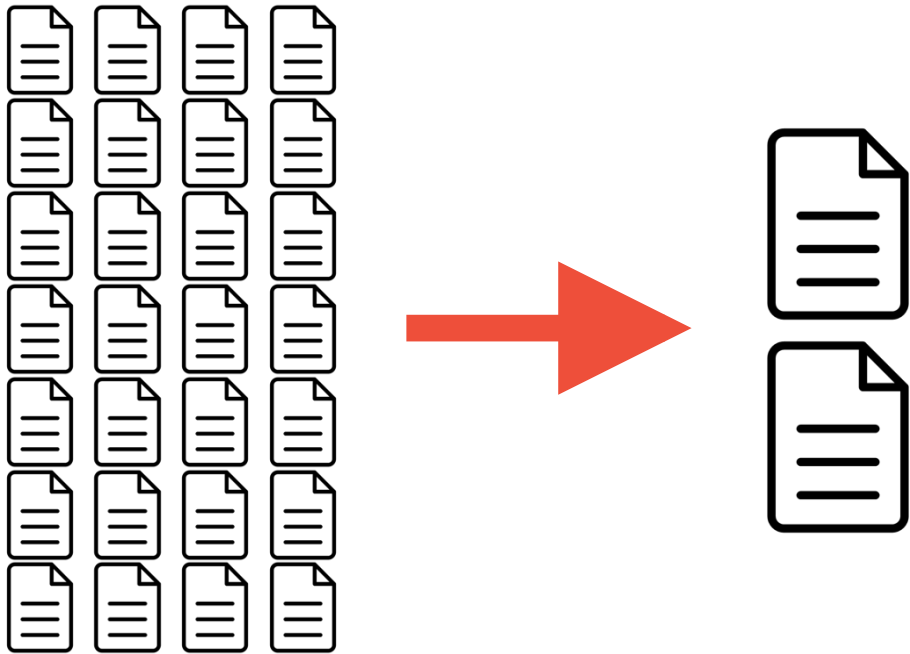
New plugin based on the Presto Hive connector

Available in Alluxio 2.1 distribution
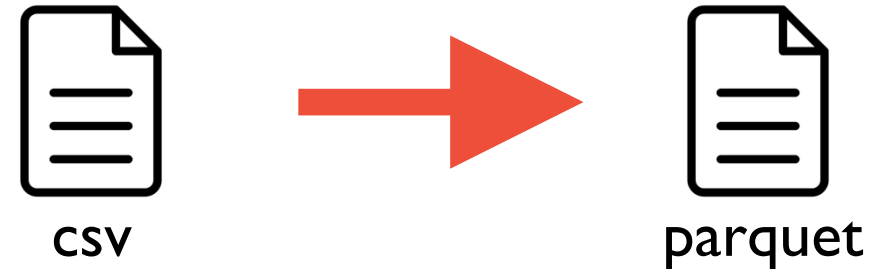
In Progress: Merging connector into Presto codebase

ALLUXIO

# Transformation Service

Transform data to be <span style="color:green">compute-optimized</span> independent from <span style="color:red">storage-optimized</span> format

**Coalesce**            **Format Conversion**



csv → parquet

ALLUXIO

# Demo!

# Demo

2 isolated AWS 10-node clusters

      Presto + Hive Metastore + S3 Data

      Presto + Alluxio + Hive Metastore + S3 Data

TPCDS sample dataset on S3

    ~10,000 CSV files

# Demo Summary

Attached existing Hive database into Alluxio Catalog

Alluxio Catalog served table metadata for Presto

Transformed store_sales by coalescing and converting CSV to Parquet

Presto Without Alluxio

20s

Alluxio Transformations

7s

Alluxio Transformations With Caching

3s

# Future Work

User community feedback/collaboration is important!

Future projects

     New UDB implementations (AWS Glue)

     More conversion formats (json)

     DDL/DML workloads (CREATE TABLE, INSERT, etc.)

     New Client APIs for structured data (Arrow)

# Developer Preview Available in Alluxio 2.1

Try it out!

     [Documentation](#)

Provide feedback

     Feature requests and issues in Github [Alluxio/alluxio](#)

## Thank You!