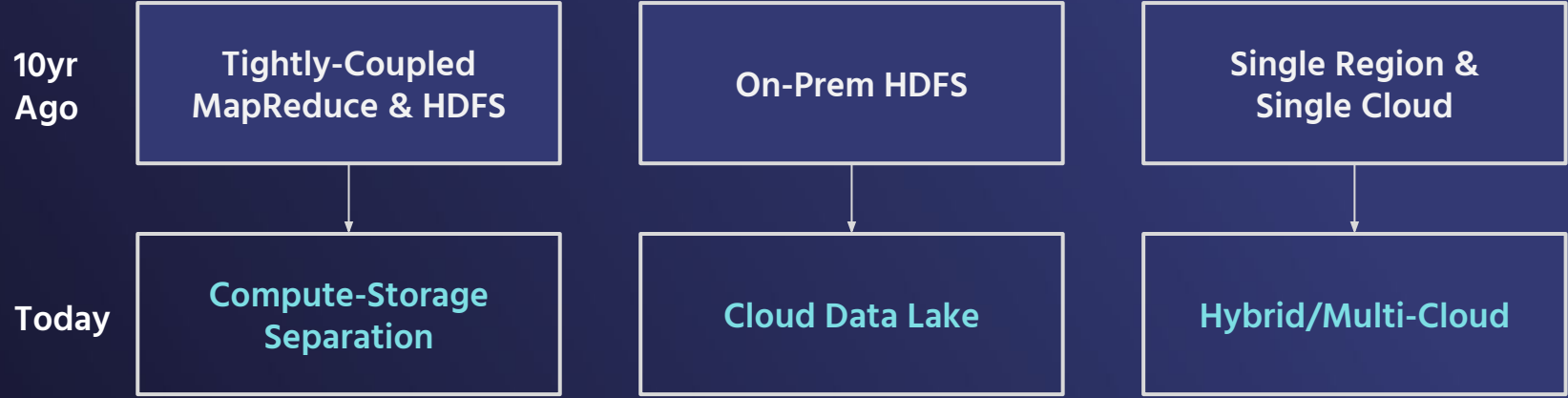# Tackling I/O Challenges in Modern Data Lakes

Hope Wang, Developer Advocate @ Alluxio
(hope.wang@alluxio.com)

# Why Should You Care About I/O?

# The Evolution of the Modern Data Stack

**10yr Ago**

| Tightly-Coupled MapReduce & HDFS | On-Prem HDFS | Single Region & Single Cloud |

**Today**

| Compute-Storage Separation | Cloud Data Lake | Hybrid/Multi-Cloud |

**More Elastic, Cheaper, More Scalable**

# The Evolution of the Modern Data Stack

**Today**

| Compute-Storage Separation | Cloud Data Lake | Hybrid/Multi-Cloud |
|---|---|---|

**Data is Remote from Compute; Locality is Missing**

**I/O Challenges**

# I/O Challenges

## Performance

- Analytics SQL: High query latency because of retrieving remote data
- Model Training: Training is slow because of loading remote data in each epoch (LISTing lots of small files is particularly slow)

## Cost

- LIST/GET/PUT operation costs add up quickly
- Cross-region data transfer (egress) fees
- GPU cycles are wasted waiting for data

## Reliability

- Job failures
- Amazon S3 errors:

```
503 Slow Down
503 Service Unavailable
```



amazon S3

AmazonS3Exception: Internal Error (Service: Amazon S3; Status Code: 500; Error Code: 500 Internal Error; Request ID: A4DBBEXAMPLE2C4D)

AmazonS3Exception: Slow Down (Service: Amazon S3; Status Code: 503; Error Code: 503 Slow Down; Request ID: A4DBBEXAMPLE2C4D)

# 10%

of your data is hot data

# 10%

of your data is hot data

→

Add a
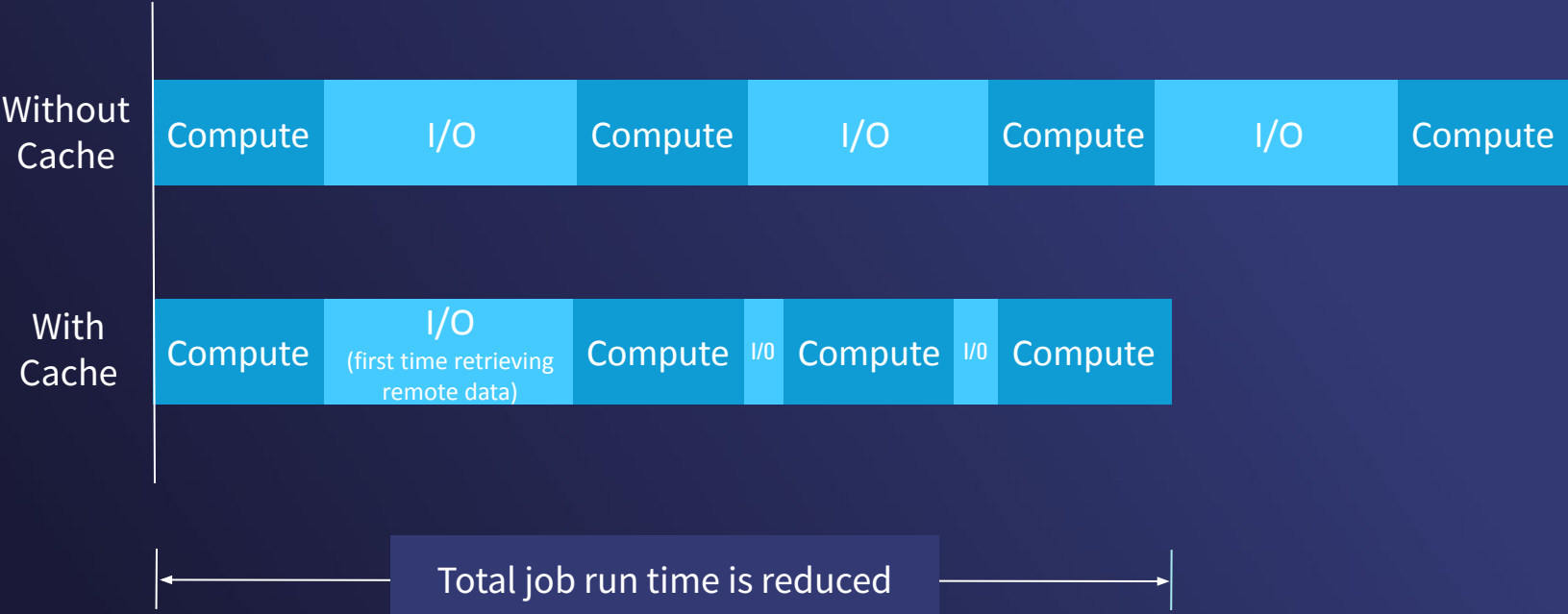
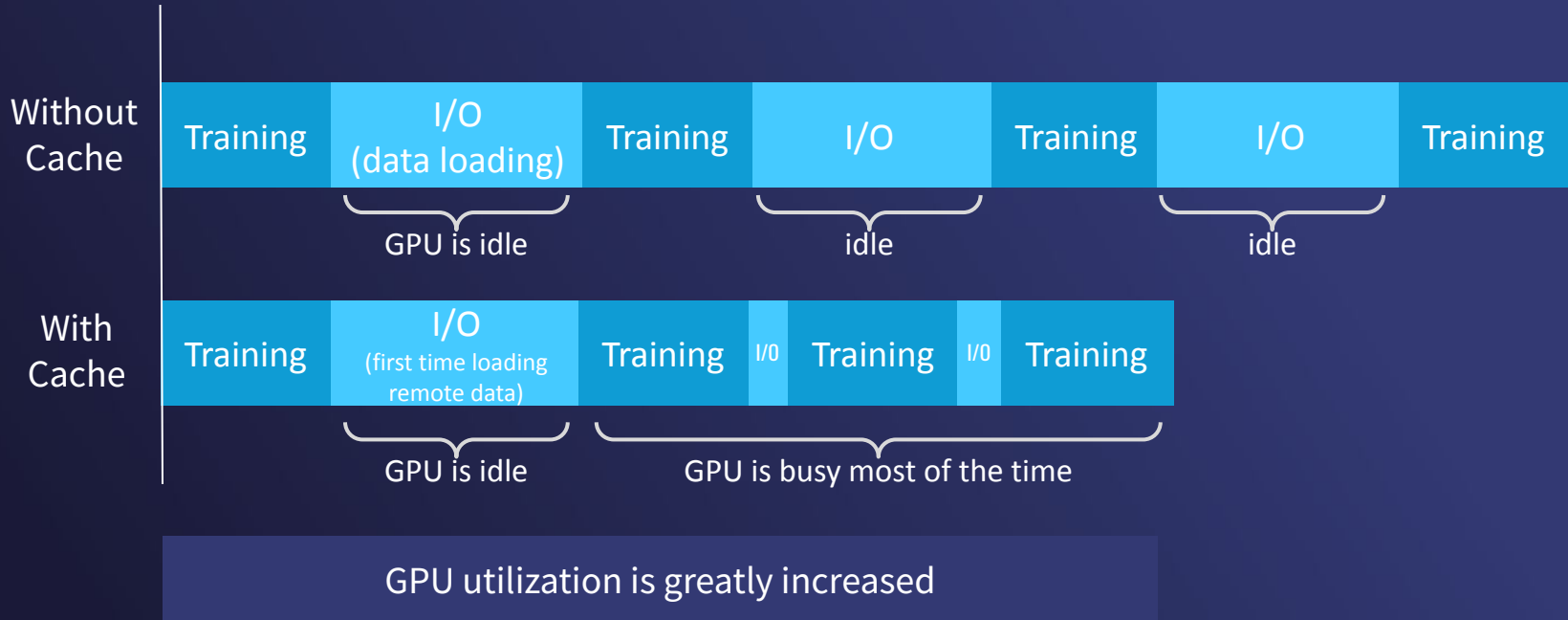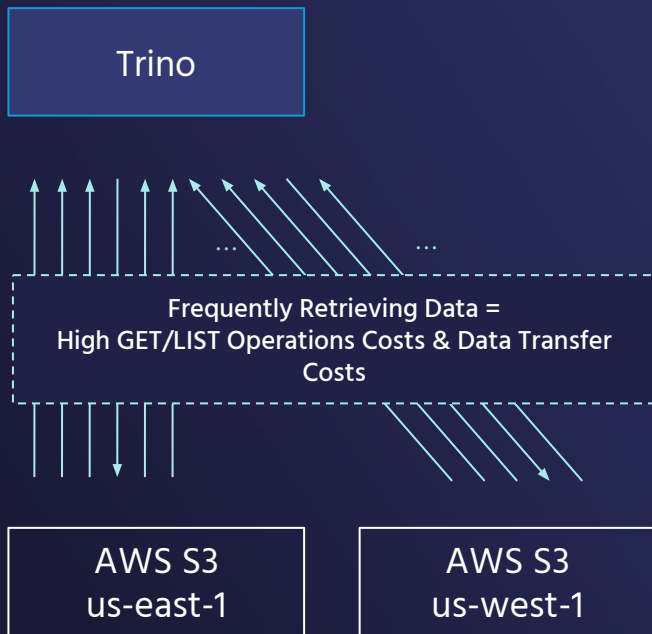## **Data Caching Layer**

between compute & storage

*Source: Alluxio*

# Reduce Latency

# Increase GPU Utilization

| | | | |
|---|---|---|---|
| **Without Cache** | Training | I/O (data loading) | Training | I/O | Training | I/O | Training |

GPU is idle | idle | idle

| | | | |
|---|---|---|---|
| **With Cache** | Training | I/O (first time loading remote data) | Training | I/O | Training | I/O | Training |

GPU is idle | GPU is busy most of the time

GPU utilization is greatly increased

# Reduce Cloud Storage Cost

## Without Cache

Trino

Frequently Retrieving Data =
High GET/LIST Operations Costs & Data Transfer
Costs

...                    ...

AWS S3
us-east-1

AWS S3
us-west-1

## With Cache

Trino

Fast Access with
Hot Data Cached

Data Cache

...                    ...

Only Retrieve Data When Necessary =
Lower S3 Costs

AWS S3
us-east-1

AWS S3
us-west-1

# Improve Reliability

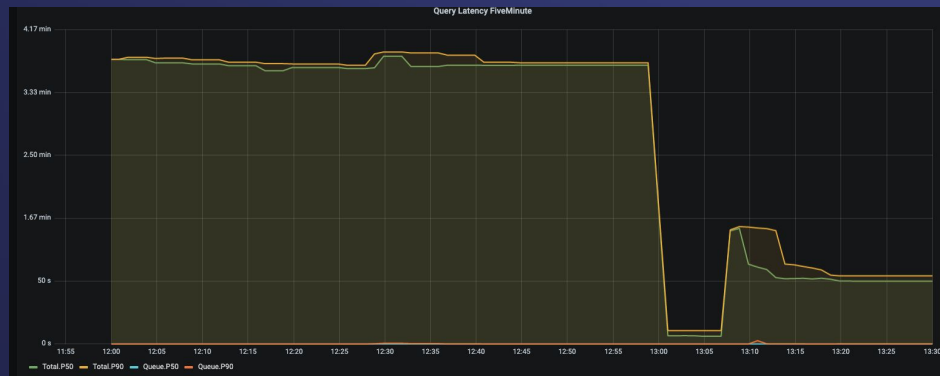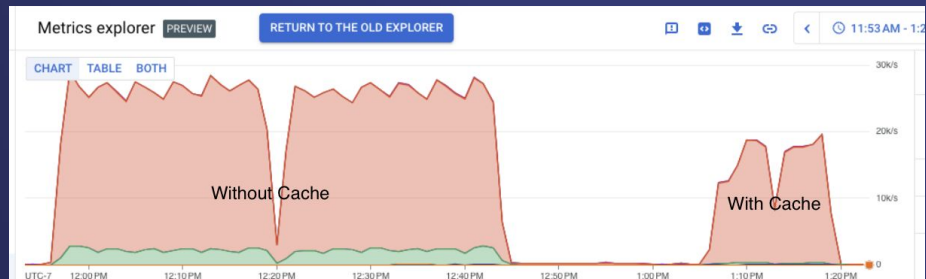Prevent Job Failures like "503 Service Unavailable"...


Prevent
Network
Congestion


Relieve
Overloaded
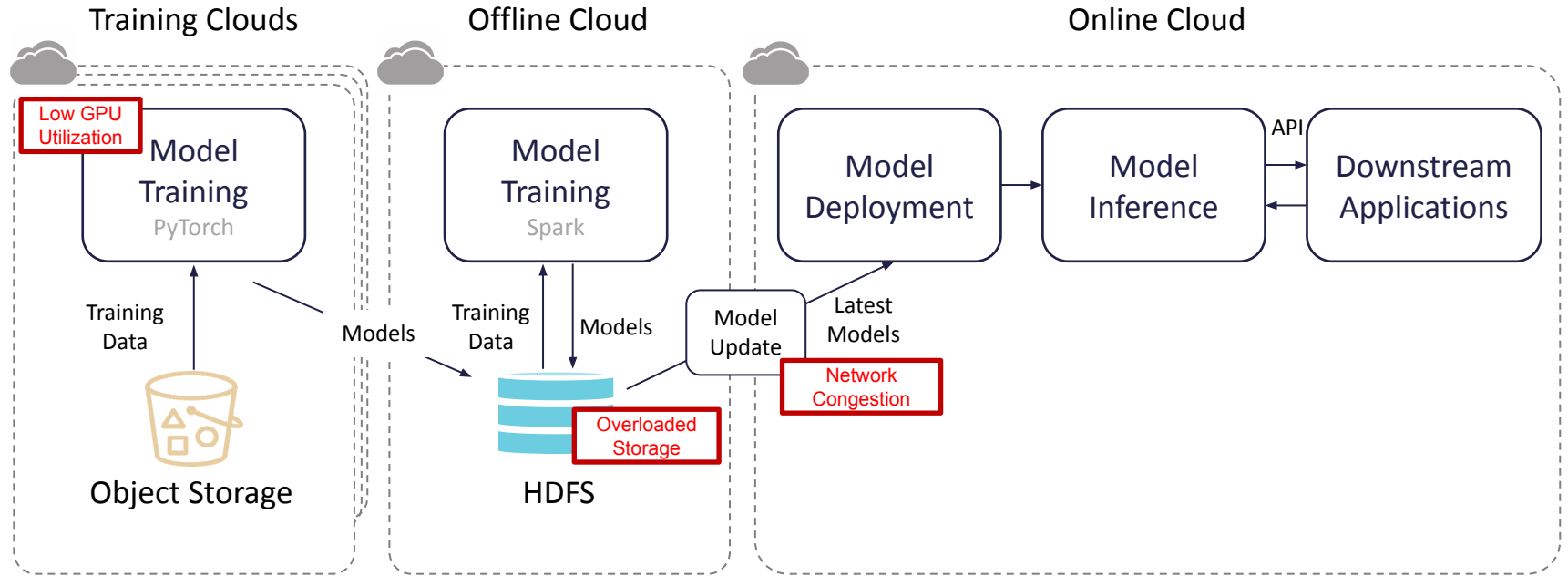Storage

# Data Caching for Presto @ Uber

- Uber deployed Presto-Alluxio local cache on cloud for initial evaluation
- Cost reduction & performance improvement after adopting Alluxio
  - 💰: >80% reduction of # of read requests to GCS (saving $$$$)
  - 🚀: 228 s ➡ 50 s reduction of P90 query latency (faster queries)
- This evaluation has greatly impacted Uber's decision on cloud migration
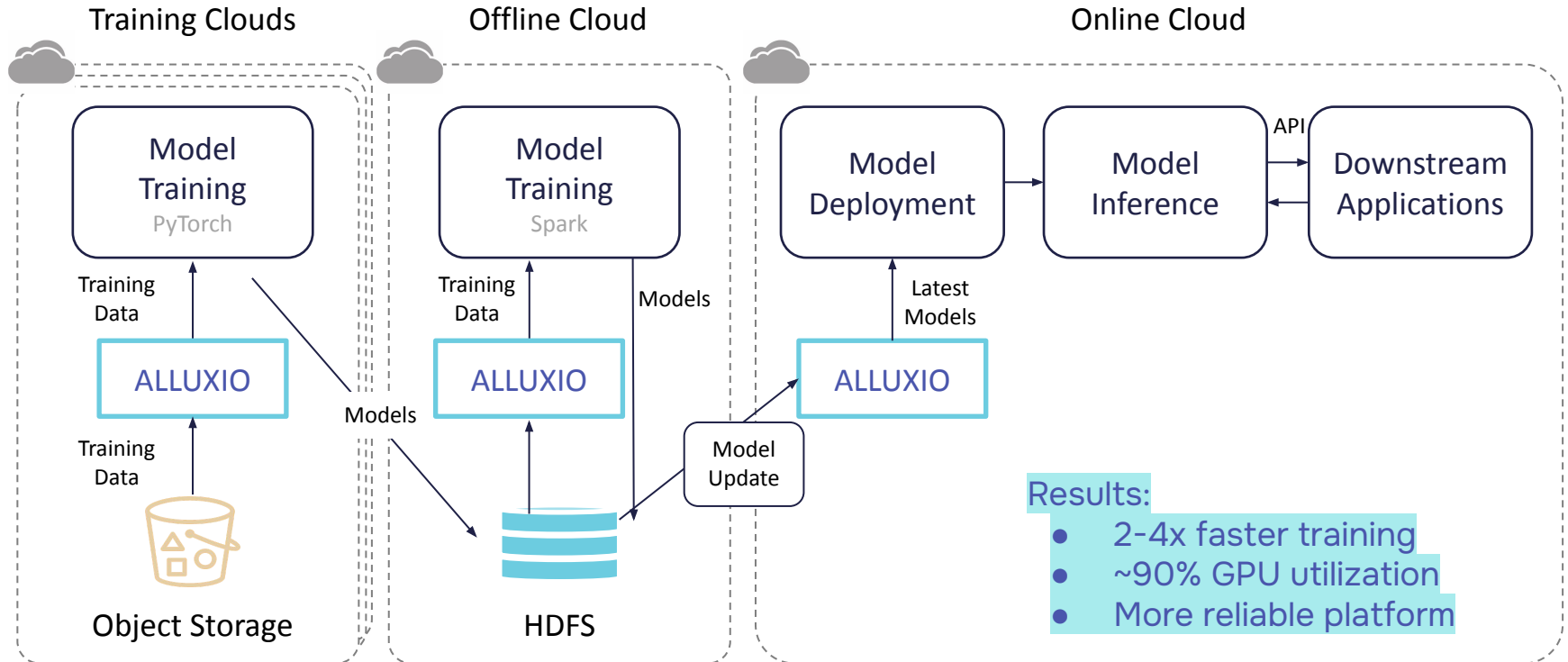
*Source: Uber+Alluxio Talk @ Community Over Code NA 2023*
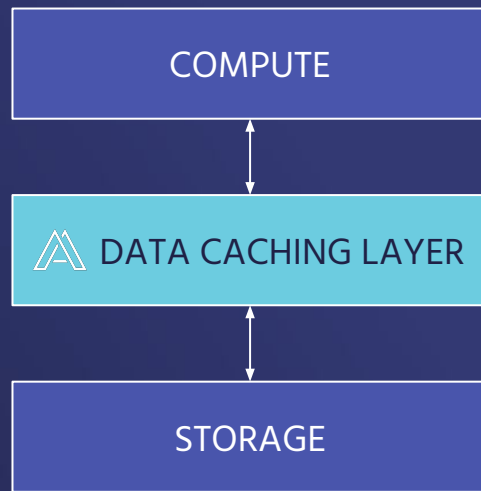
# Data Caching Across ML Data Pipeline

## With Cache



Training Clouds

Offline Cloud

Online Cloud

Model Training
PyTorch

Model Training
Spark

Model Deployment

Model Inference

Downstream Applications

API

Training Data

Training Data

Latest Models

ALLUXIO

ALLUXIO

ALLUXIO

Training Data

Models

Models

Model Update

Object Storage

HDFS

Results:
- 2-4x faster training
- ~90% GPU utilization
- More reliable platform

# Key Takeaways

- The evolution of modern data stack poses challenges for **data locality**
- You should care about I/O in data lake because it greatly impacts the **performance, cost** & **reliability** of your data platform
- Having a data caching layer **between compute and storage** can solve the I/O challenges
- You can use cache for both **analytics and AI workloads**

```
┌─────────────────────────────┐
│          COMPUTE            │
└─────────────────────────────┘
               ↕
┌─────────────────────────────┐
│   △  DATA CACHING LAYER      │
└─────────────────────────────┘
               ↕
┌─────────────────────────────┐
│          STORAGE            │
└─────────────────────────────┘
```

# Thanks!
# Let's Grab a Coffee!

**Do you have any questions?**

✉️ Email me: hope.wang@alluxio.com

in Connect with me: https://www.linkedin.com/in/hopechong/

Ping me on Alluxio Slack: https://alluxio.io/slack

Scan the QR code for a **Linktree** including great learning resources, exciting meetups & a community of data & AI infra experts!