



## Managing Datastore Locality at Scale with Akkio

Muthukaruppan Annamalai, Kaushik Ravichandran, Harish Srinivas, Igor Zinkov,  
Luning Pan, Tony Savor, David Nagle and Michael P. Hines

FACEBOOK, UNIVERSITY OF TORONTO\*

1



# Motivation

ALTOONA, IA

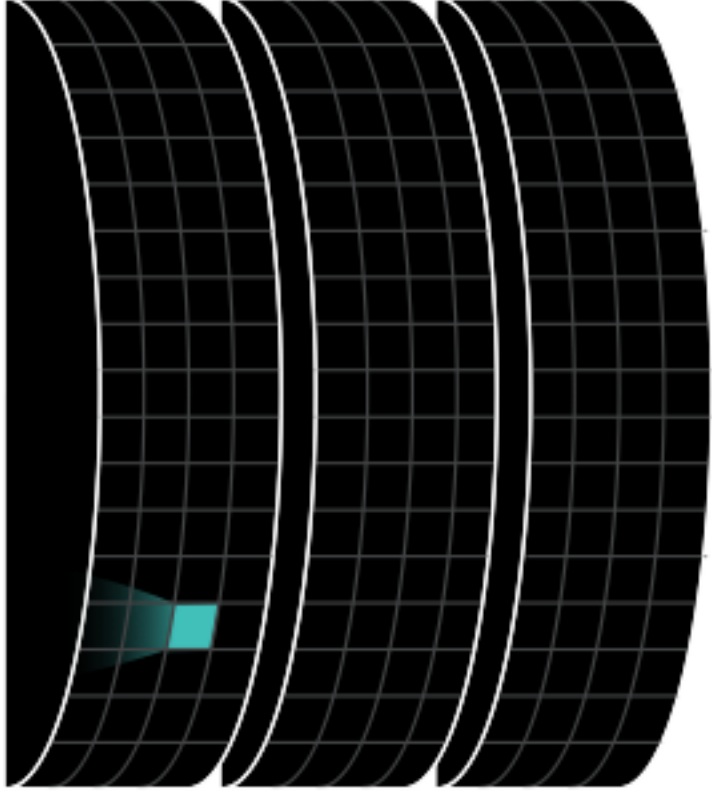
PRINEVILLE, OR



# Motivation



$\mu$ -shard



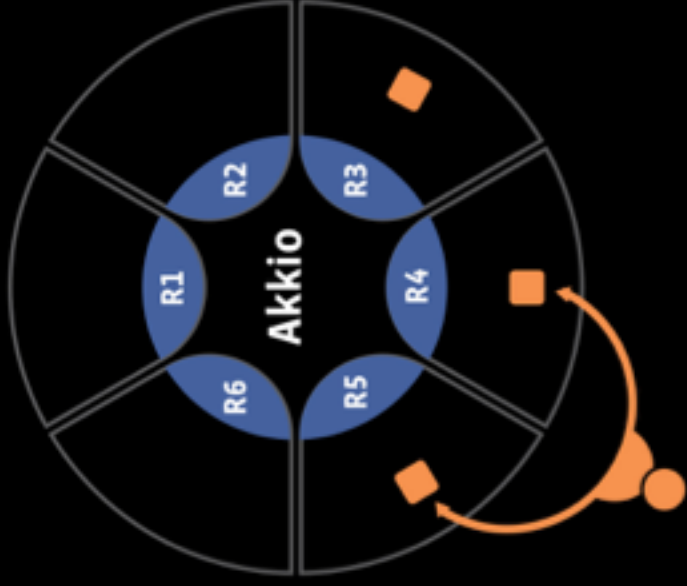
**Saving Space**





**Saving WAN  
Bandwidth**





Strong consistency  
is not expensive

Better serves  
workloads with low  
read-write ratios

|



## Flexible

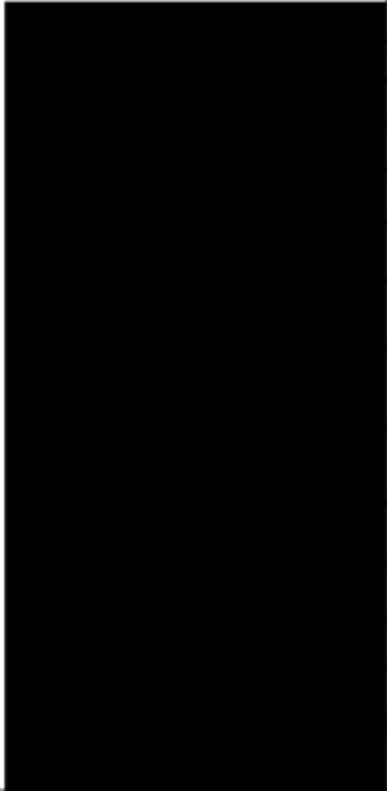
Supports five different  
durable storage systems

## Scalable

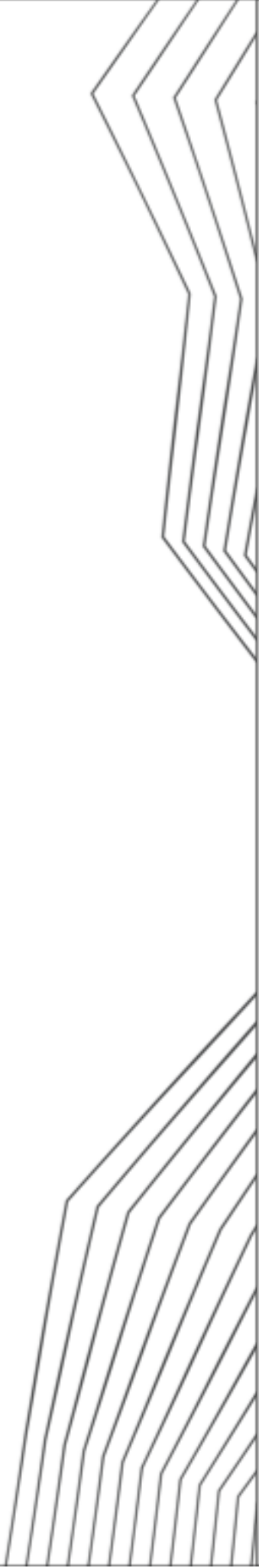
100's of millions of data access  
requests per second

OF CROSS DC TRAFFIC DAILY

OF DATA STORED

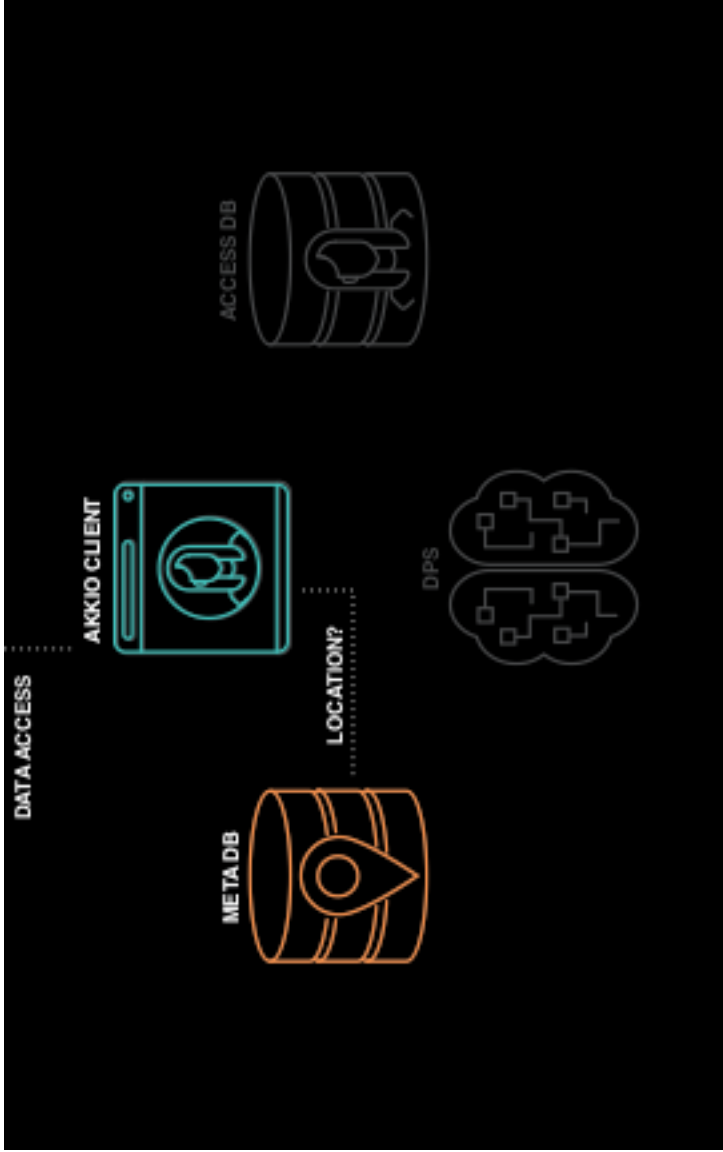


2



WEST COAST DB

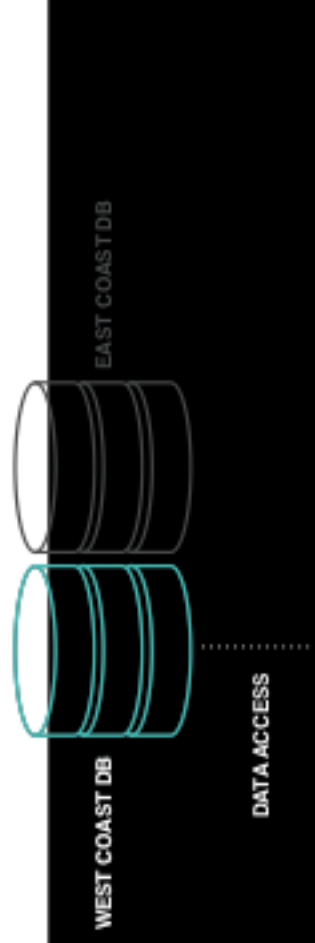
EAST COAST DB

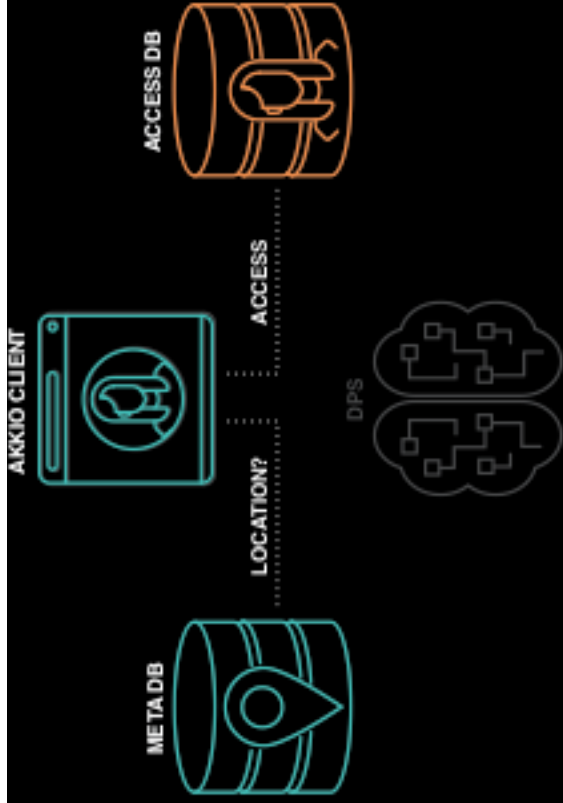


## Metadata DB

- Stores  $\mu$ -shard -> placement mapping
- A cache in front to absorb most reads

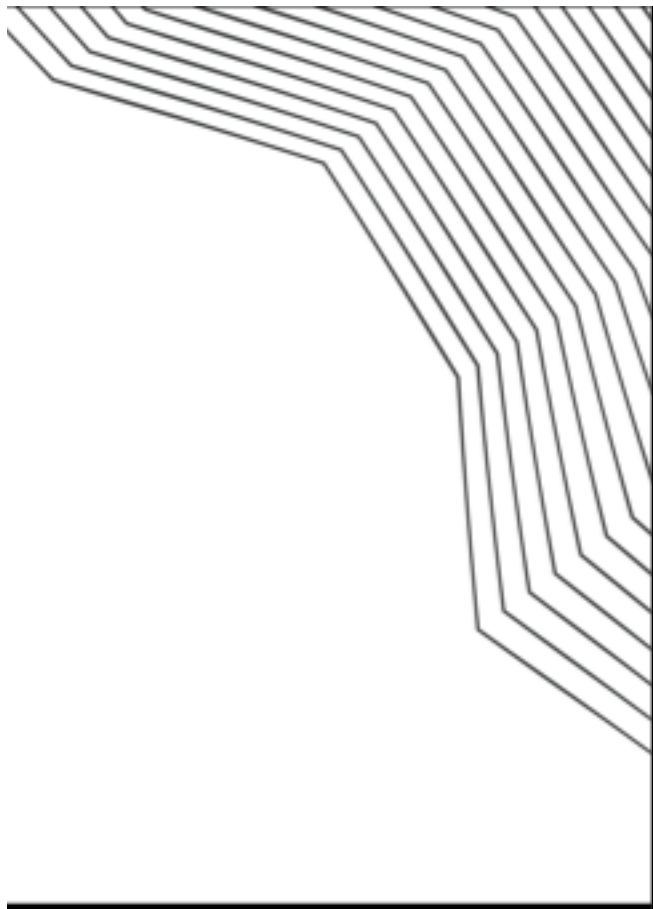
- Stale lookups are prevented through means like per- $\mu$ -shard Access Control Lists (ACLs)
- Typical footprint is ~300GB





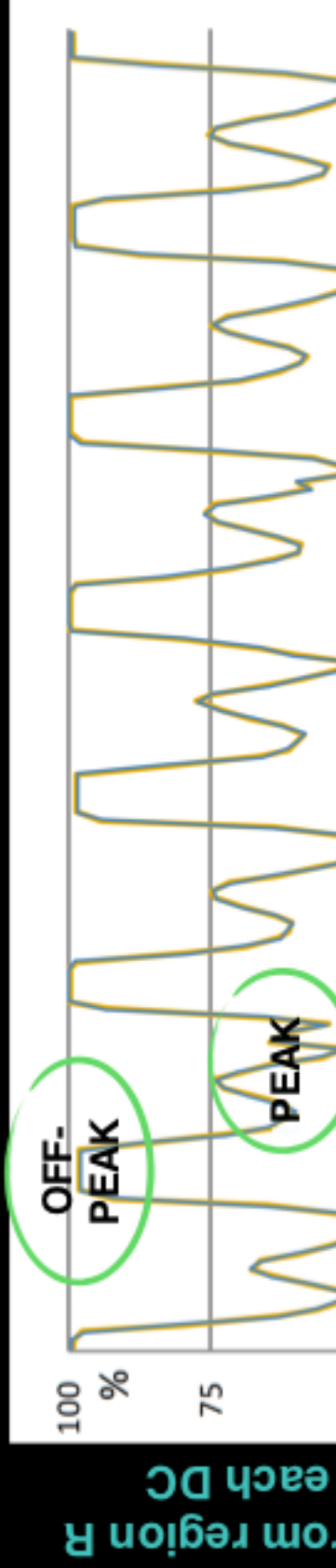
# Access DB

- Fundamentally a time-windowed, counter service
- One counter per- $\mu$ -shard, datacenter> tuple

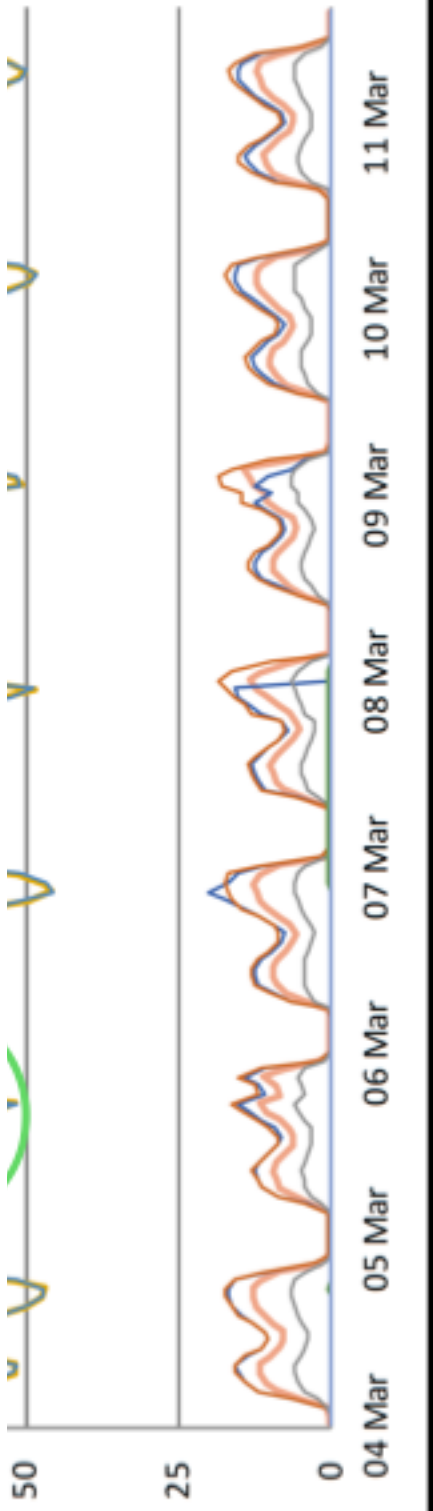


- Typical footprint is ~300GB
- Similar services can share access data
- Routing layer is a complete black box

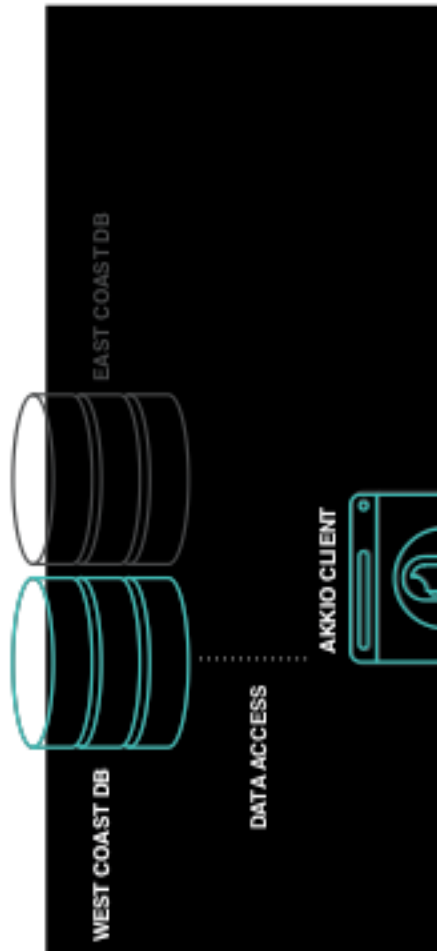
## Routing Layer Load Balancing? No problem!



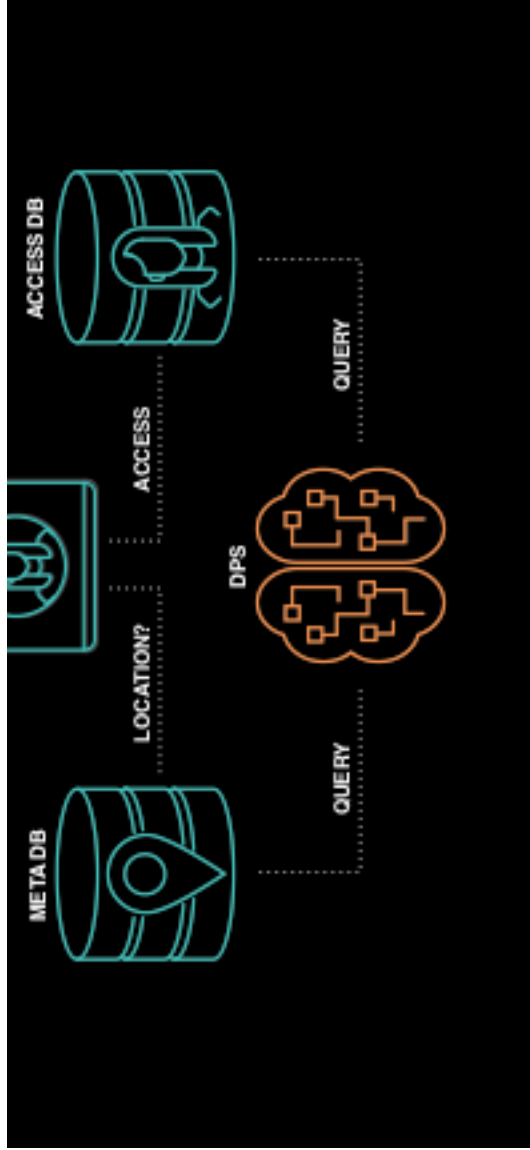
% of requests fr  
serviced by



Time







## Data Placement Service

- Data Placement Service knows list of possible placements and associated data centers
- Receives hints on sub-optimally placed  $\mu$ -shards from Akkio client
- New  $\mu$ -shards are provisioned through the Data Placement Service

- Runs a set of policies to determine migration eligibility
- Identifies an optimal location using a configurable scoring policy

## Default Scoring Policy

- Eliminate any placement possibilities with hot data centers
- Generate per-data center score
  - Number of times a  $\mu$ -shard was accessed from a data center in the last N days
    - N is configurable

- Generate per-placement score
  - Sum of corresponding data center scores
- Pick placement with highest score
  - If tied, break tie using resource usage information

## Fault Tolerance

- Migrations are serialized through a per- $\mu$ -shard lock on the Metadata DB
- A work item in the Metadata DB is created to track any work
  - Atomically created and deleted with the lock
  - Used to recover migrations
  - Monotonically increasing epoch to fence writes

# Migration

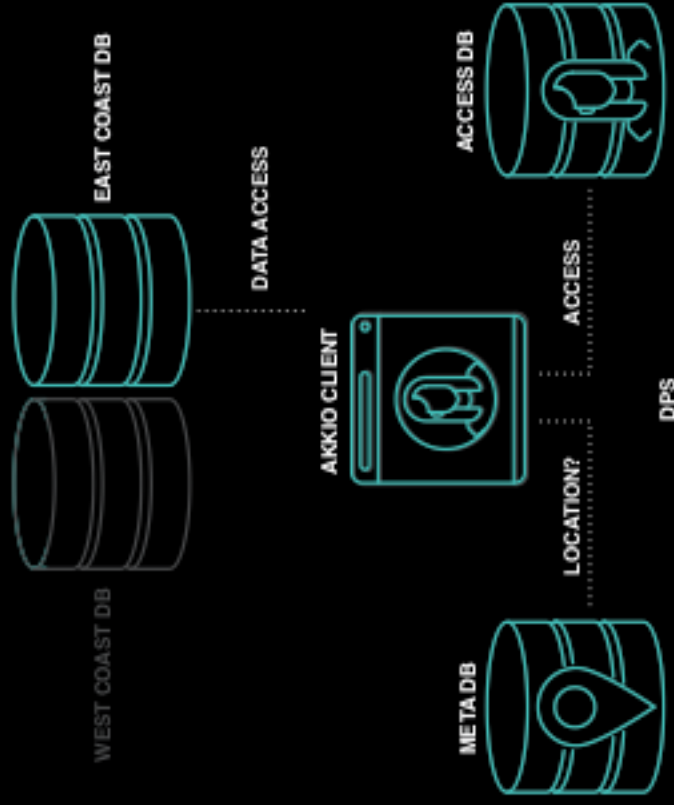
- Every backend has a custom Data Placement Service plug-in to enable reads/writes to backend
- Configurable migration strategy
  - Depends on capabilities of underlying store
    - Native timestamps
    - Fine-grained ACLs
    - Read-modify-write transactions

## Simple Migration Strategy

1. [Metadata DB] Atomically grab per- $\mu$ -shard lock & create work item
2. [Source] Set per- $\mu$ -shard ACL to R/O, drain writes, read data
3. [Destination] Write data, set per- $\mu$ -shard ACL to R/O
4. [Metadata DB] Update placement
5. [Source] Delete data & per- $\mu$ -shard ACL

6. [Destination] Set per- $\mu$ -shard ACL to R/W

7. [Metadata DB] Atomically delete lock & work item





QUERY

QUERY

3

## Use Case

40%

Storage Capacity  
with

75%

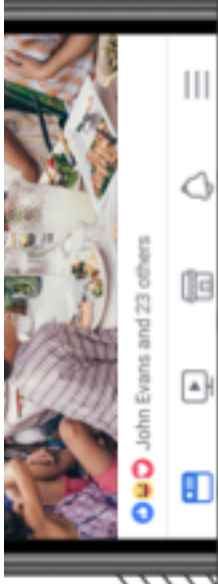
Reduction in  
WAN BW

68%

Savings in jobs  
Latencies







Typically its  
cheaper to move  
compute

Data does not need  
to be present  
everywhere

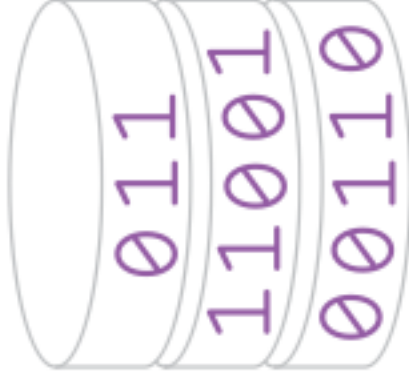
Savings on WAN  
bandwidth and  
latency wins



0001101000110  
1101010101101  
00100101001111



**HOT DATA**



**WARM DATA**



**COLD DATA**

A decorative border surrounds the page, featuring a repeating geometric pattern of nested, slightly offset lines that create a sense of depth and movement. The pattern is composed of multiple parallel lines that are not perfectly straight, giving it a hand-drawn or organic feel.

**Locality at fine-granularity works!**



**THANK YOU**

**POSTER BOOTH #9**

