# RAMCloud: A Low-Latency Datacenter Storage System

## Ankita Kejriwal

## Stanford University

**(Joint work with Diego Ongaro, Ryan Stutsman, Steve Rumble,**

**Mendel Rosenblum and John Ousterhout)**

# What if you had...

## … a Storage System that provides:

- **Scale**
  - Data size: 10 PB
  - Accessible by 100,000 nodes (10 Million cores)

- **Uniform fast random access time to all data**
  - 100 B read: 2 µs RPC
  - 100 B write: 5 µs RPC
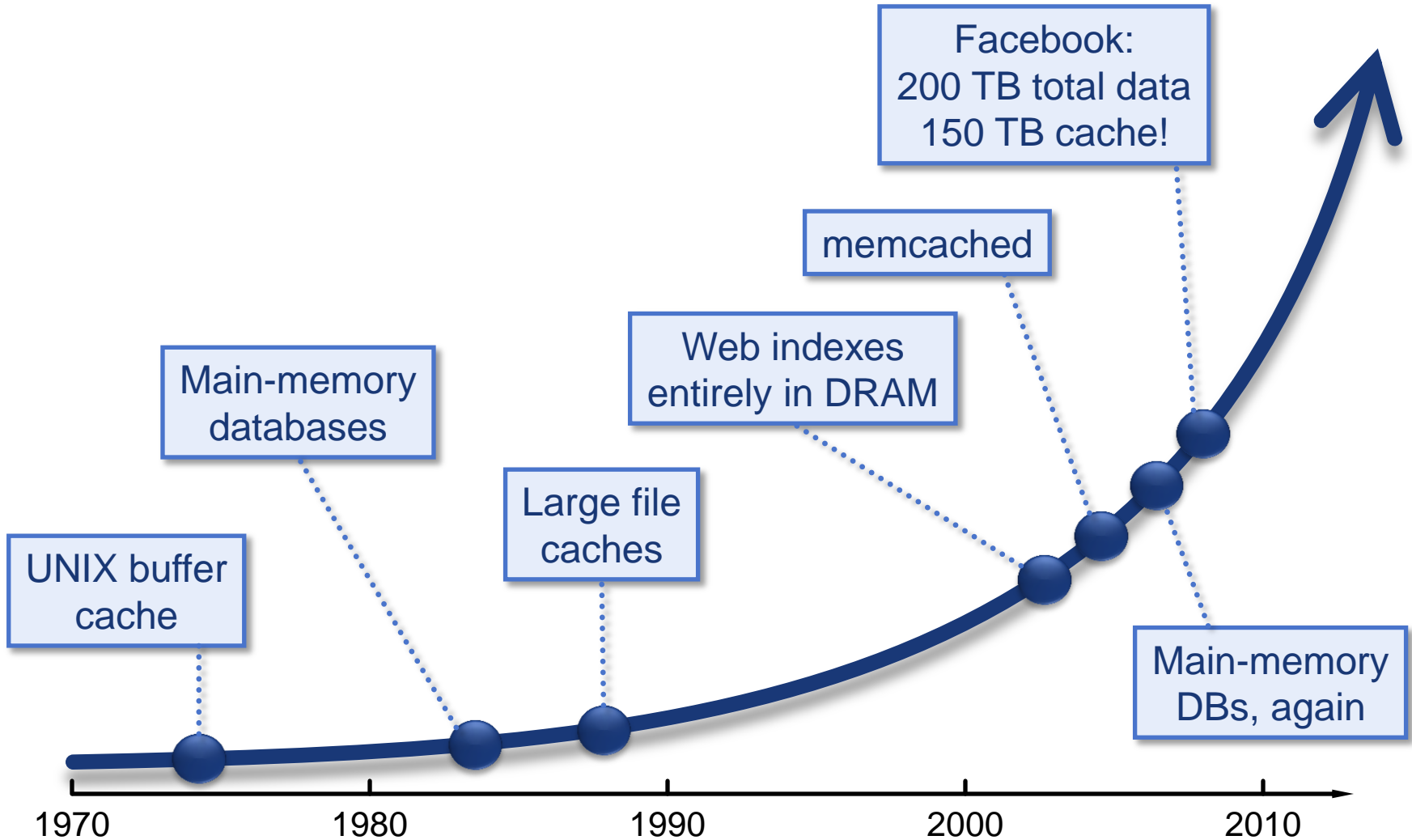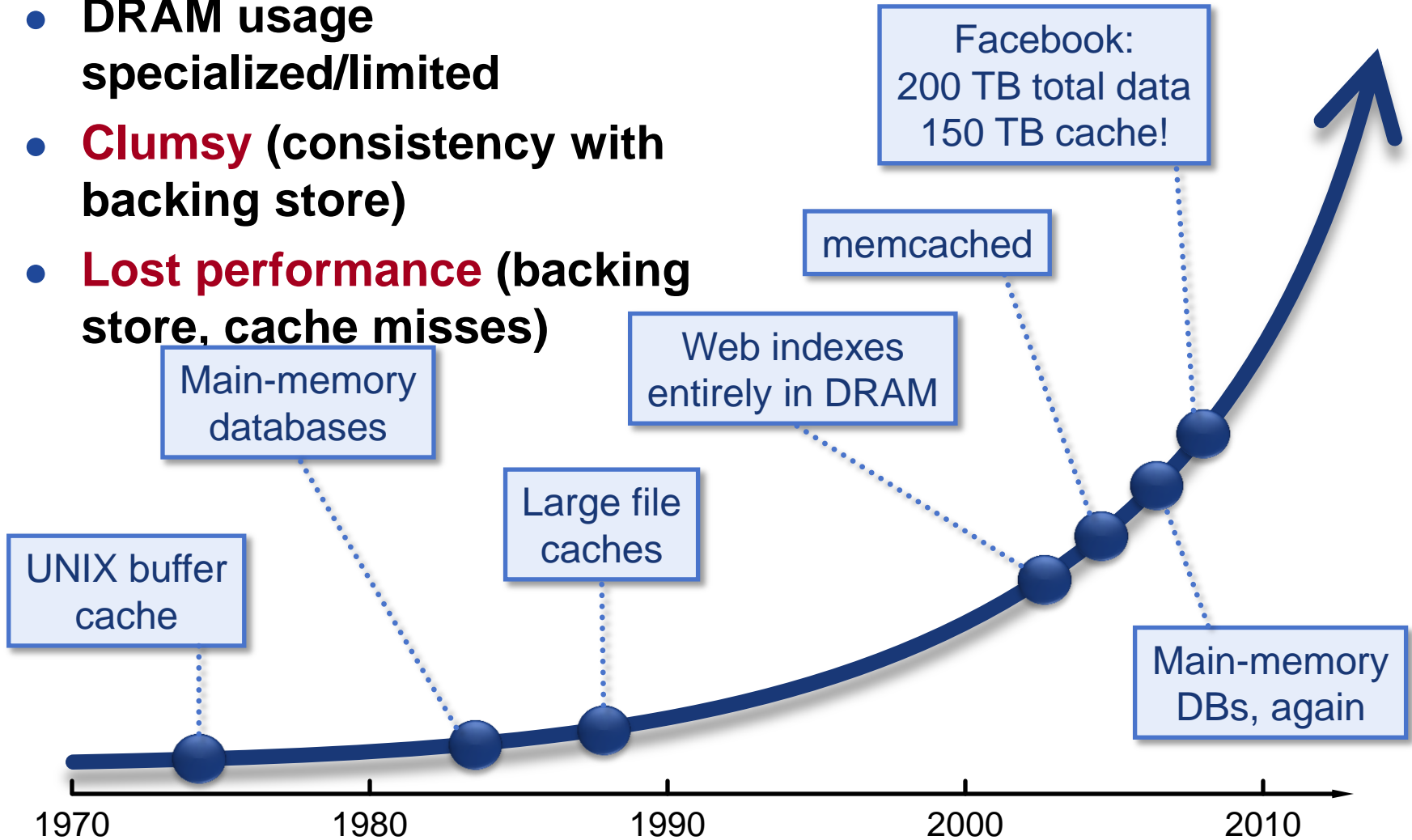
- **Durable and available**

# RAMCloud

- **General-purpose storage system**

- **All data always in DRAM**

- **Scale: 1000 – 10000 servers, 1 PB data**

- **Performance goals:**
  - High throughput: 1M ops/sec/server
  - Low-latency access: 5-10μs RPC

- **Durable and available**

- **Potential impact: enable new class of applications**
  - Primary motivation: Web sphere
  - Maybe HPC?

# DRAM in Storage Systems



Facebook:
200 TB total data
150 TB cache!

memcached

Web indexes
entirely in DRAM

Main-memory
databases

Large file
caches

UNIX buffer
cache

Main-memory
DBs, again

1970    1980    1990    2000    2010

# DRAM in Storage Systems
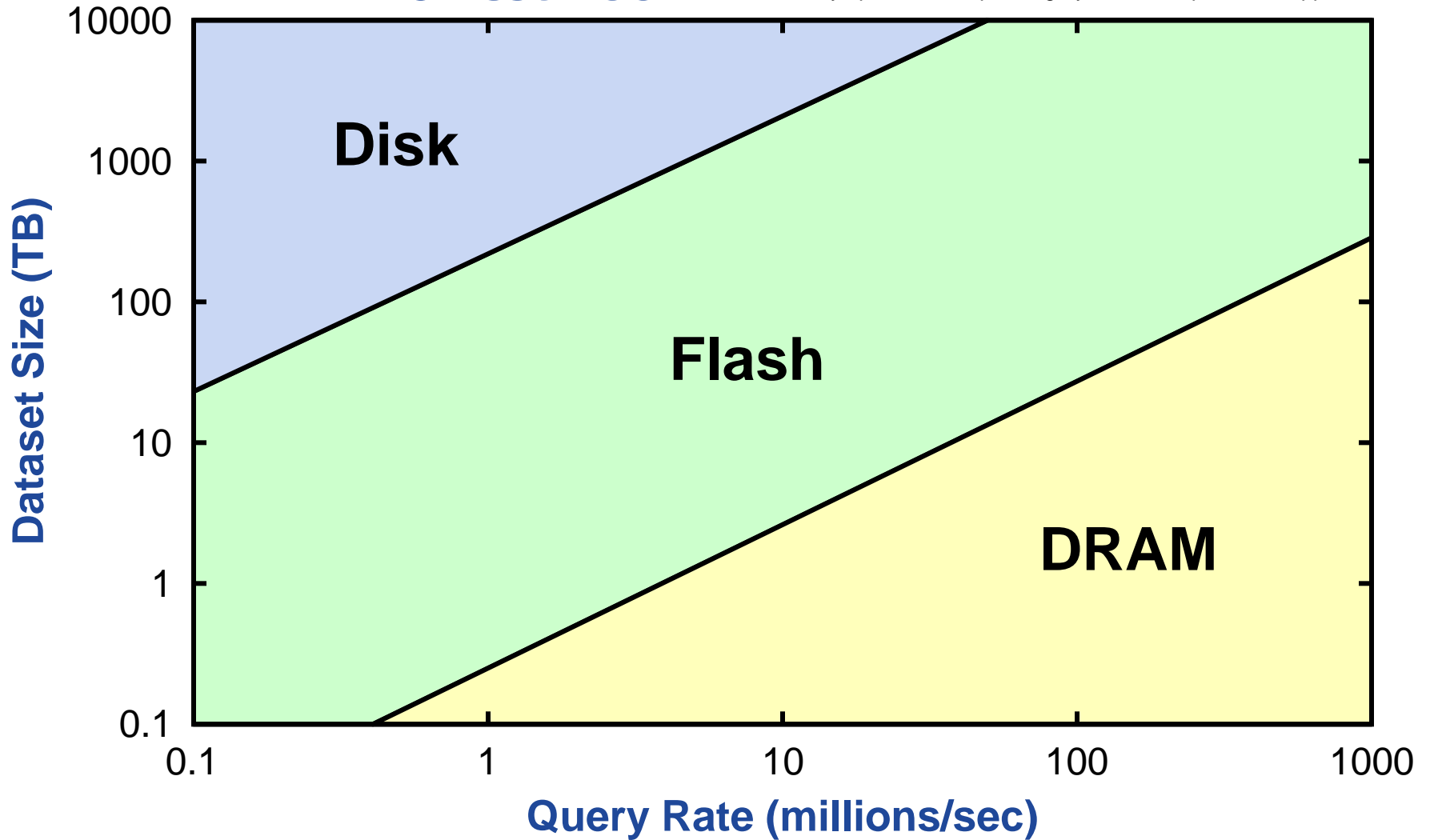
- **DRAM usage specialized/limited**
- **Clumsy** (consistency with backing store)
- **Lost performance** (backing store, cache misses)

Facebook:
200 TB total data
150 TB cache!

memcached

Web indexes entirely in DRAM

Main-memory databases

Large file caches

UNIX buffer cache

Main-memory DBs, again
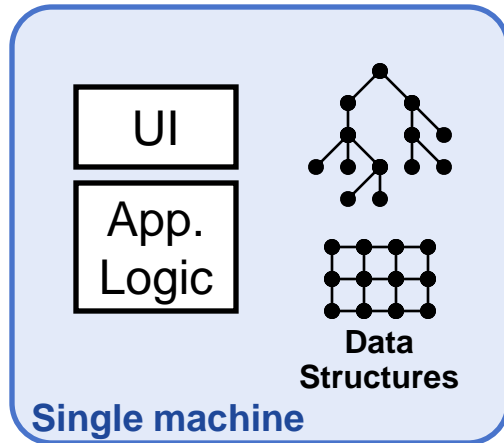
1970    1980    1990    2000    2010

# DRAM is cheaper!

**Lowest TCO**

from "Andersen et al., "FAWN: A Fast Array of Wimpy Nodes",
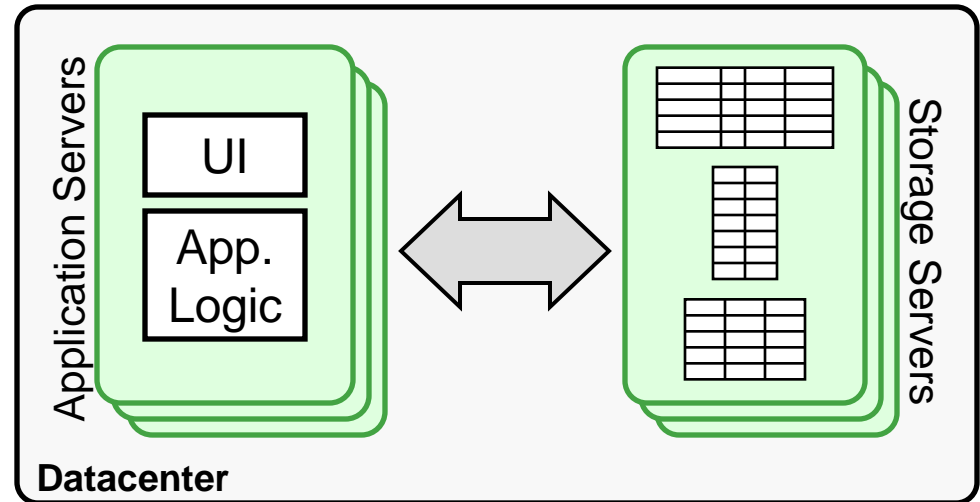Proc. 22nd Symposium on Operating System Principles, 2009, pp. 1-14.



**Dataset Size (TB)** (y-axis: 0.1, 1, 10, 100, 1000, 10000)

Regions: **Disk**, **Flash**, **DRAM**

**Query Rate (millions/sec)** (x-axis: 0.1, 1, 10, 100, 1000)

# Why Does Latency Matter?

**Traditional Application**
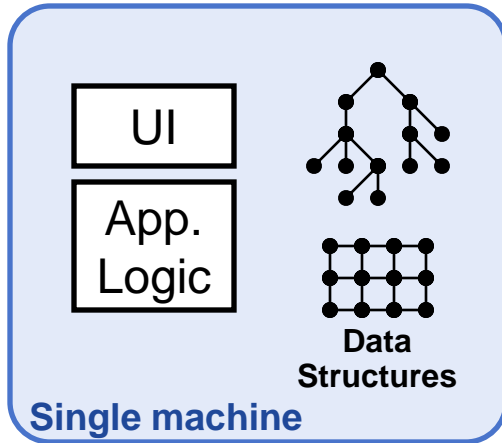
**Web Application**



**<< 1µs latency**

**0.5-10ms latency**

- **Large-scale apps struggle with high latency**
  - Random access data rate has not scaled!
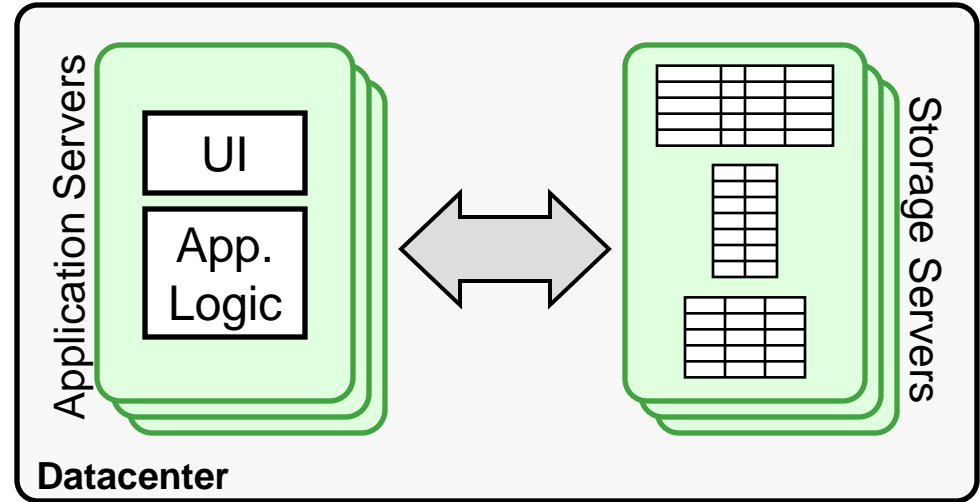  - Facebook: can only make 100-150 internal requests per page

# RAMCloud Goal: Scale and Latency

**Traditional Application**



**Single machine**

**Data Structures**

**<< 1μs latency**

**Web Application**



Application Servers
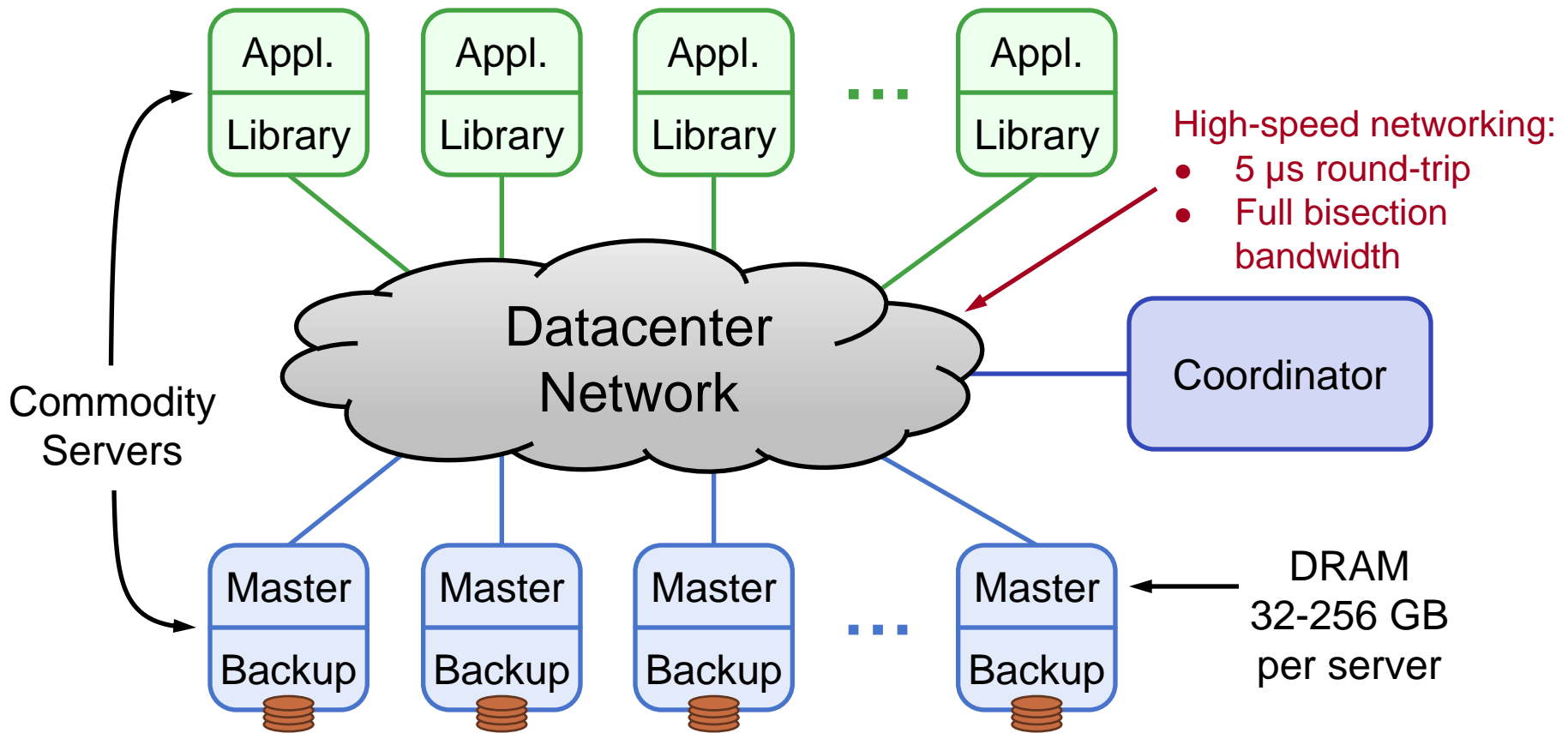
Storage Servers

**Datacenter**

**~~0.5-10ms~~ latency**
**5-10μs**

- **Enable new class of applications**

# RAMCloud Architecture

**1000 – 100,000 Application Servers**



High-speed networking:
- 5 µs round-trip
- Full bisection bandwidth

Coordinator

Datacenter Network

Commodity Servers

Appl. | Library
Appl. | Library
Appl. | Library
... | Appl. | Library

Master | Backup
Master | Backup
Master | Backup
... | Master | Backup

DRAM 32-256 GB per server
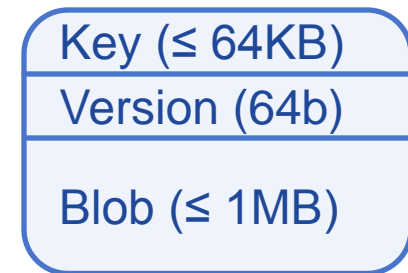
**1000 – 10,000 Storage Servers**

# Data Model: Key-Value Store

## Tables

key | value
key | value
key | value

key | value
key | value
key | value
key | value

key | value
key | value
key | value
key | value

Object

| Key (≤ 64KB) |
| Version (64b) |
| Blob (≤ 1MB) |

```
read(tableId, key)
    => blob, version

write(tableId, key, blob)
    => version

cwrite(tableId, key, blob, version)
    => version

delete(tableId, key)

enumerate(tableId)
```

Richer model in the future:
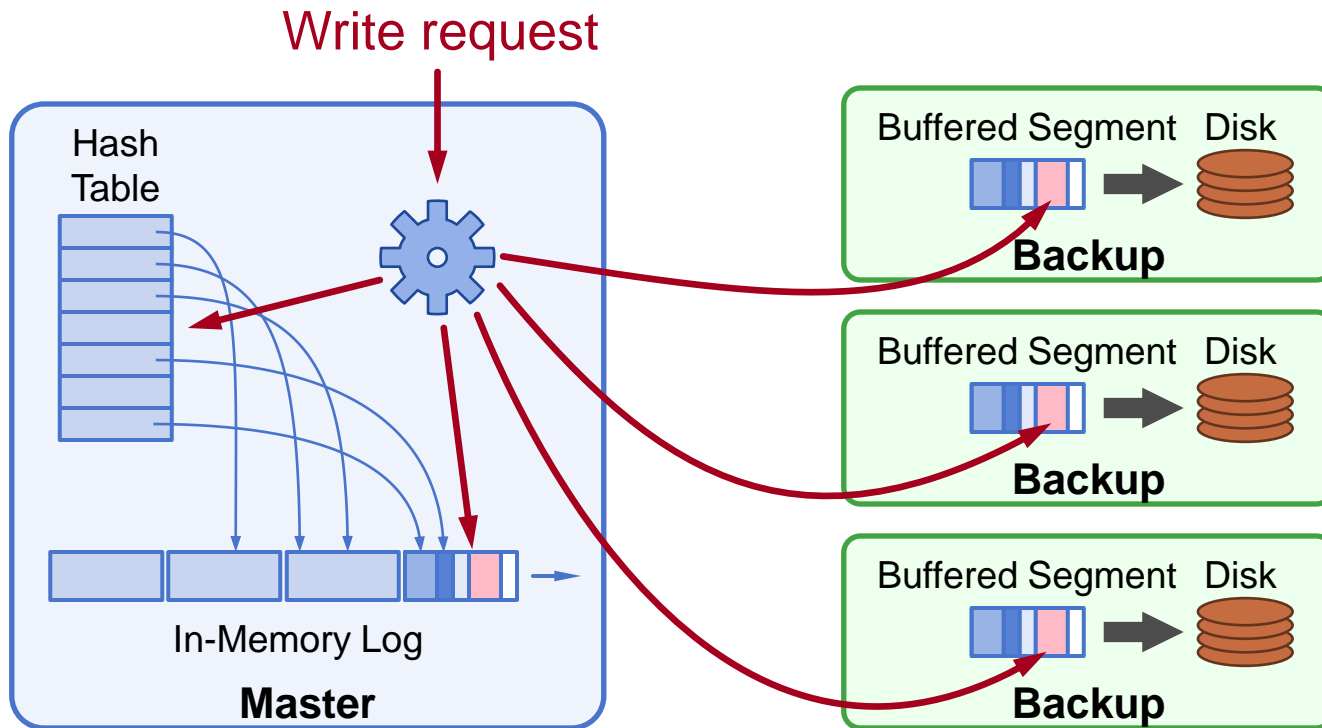- Indexes?
- Transactions?
- Graphs?

# Durability and Availability

- **Goals:**
  - No impact on performance
  - Minimum cost, energy

- **Keep replicas in DRAM of other servers?**
  - 3x system cost, energy
  - Still have to handle power failures

- **RAMCloud approach:**
  - 1 copy in DRAM
  - Backup copies on disk/flash: <span style="color:darkred">durability ~ free!</span>

- **Issues to resolve:**
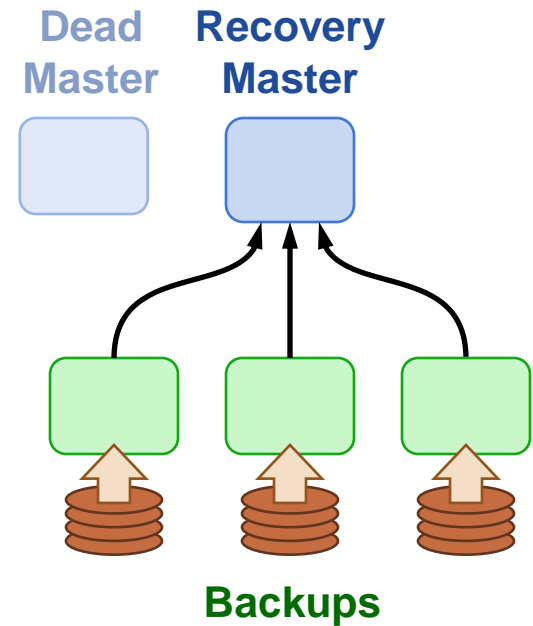  - Synchronous disk I/O's during writes??
  - Data unavailable after crashes??

# Buffered Logging

Write request



- **No disk I/O during write requests**

- **Log-structured: backup disks and master's memory**
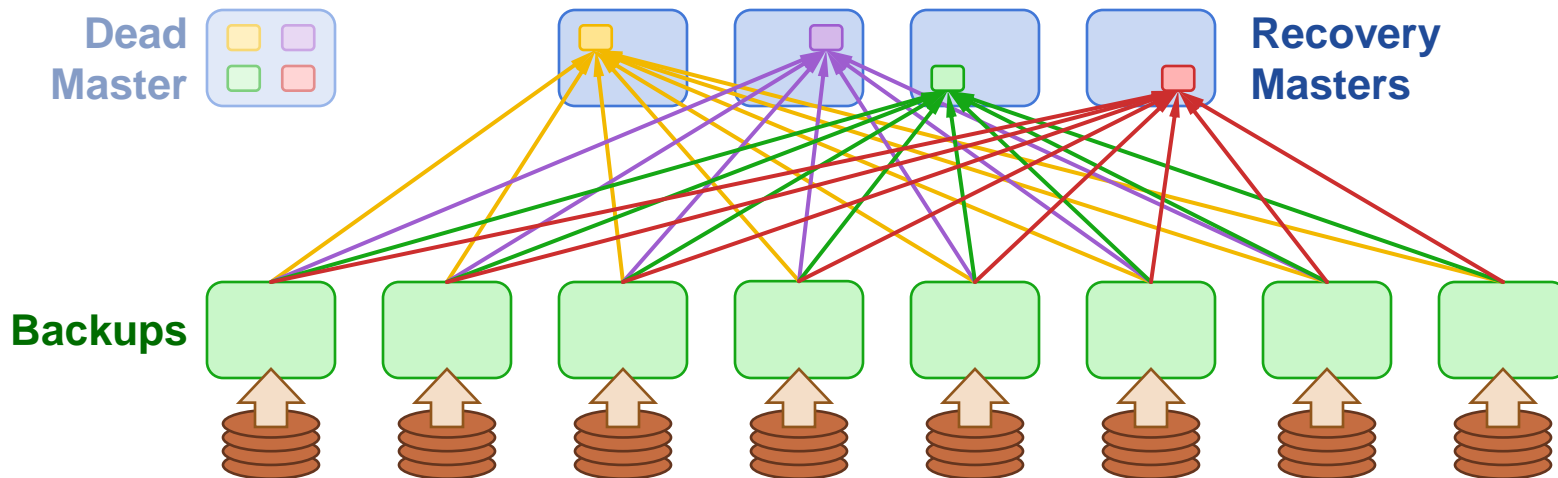
- **Log cleaning**

# Crash Recovery

- **Server crashes:**
  - Must replay log to reconstruct data

- **Crash recovery:**
  - Choose recovery master
  - Backup reads log info from disk
  - Transfers logs to recovery master
  - Recovery master replays log

- **Meanwhile, data is unavailable**

- **RAMCloud approach: fast crash recovery**
  - 1-2 seconds for 100 GB of data
  - Use system scale to get around bottlenecks

**Dead Master**  **Recovery Master**

**Backups**

# Fast Crash Recovery

- **Scatter backup data across backups**

- **Divide each master's data into partitions**
  - Recover each partition on a separate recovery master
  - Each backup divides its log data among recovery masters

# RAMCloud Project Status

- **Goal: build production-quality implementation**

- **Nearing 1.0-level release**

- **Current test cluster:**
    - 80 servers, 2 TB data
    - High speed Infiniband networking
    - Performance:
        - 100 B read: 5.3 µs RPC
        - 100 B write: 15 µs RPC

- **Interested in finding applications for RAMCloud**

# Is RAMCloud right for HPC apps?

**Properties of RAMCloud relevant to application developers:**

- **Durability and availability**

- **Key-value store**

- **Commodity hardware**

- **Read / write access latency**

- **Random access to small objects**

# Conclusion

- **General-purpose storage system**

- **All data always in DRAM**

- **Designed for:**

  - Scale: 1000 – 10000 servers, 1 PB data

  - Performance: 5-10μs RPC

- **Durable and available**

# Questions

- **Is RAMCloud appropriate for HPC Applications?**
  - Durability and availability
  - Key-value store
  - Commodity hardware
  - Read / write access latency
  - Random access to small objects

- **One thing that we could change to make RAMCloud interesting to you!**

# Thank you!