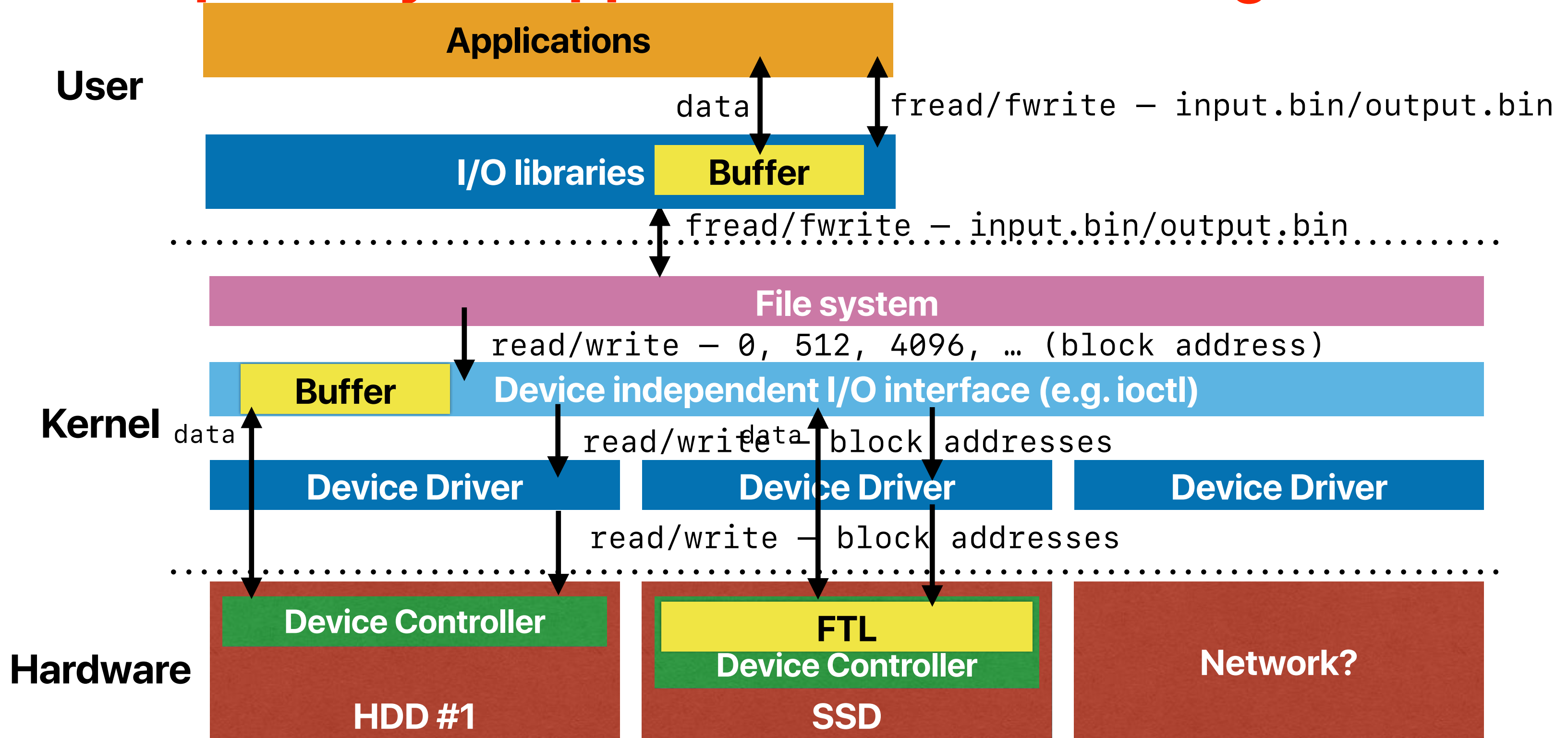


# **File systems over the network**

Hung-Wei Tseng

# Recap: How your application reaches storage device



# Recap: File systems on a computer

- Unix File System
  - Hierarchical directory structure
  - File — metadata (inode) + data
  - Everything is files
- BSD Fast File System — optimize for reads
  - Cylinder group — Layout data carefully with device characteristics, replicated metadata
  - Larger block size & fragments to fix the drawback
  - A few other new features
- Sprite Log-structured File System — optimize for small random writes
  - Computers cache a lot — reads are no more the dominating traffic
  - Aggregates small writes into large sequential writes to the disk
  - Invalidate older copies to support recovery

# Recap: Extent file systems — ext2, ext3, ext4

- Basically optimizations over FFS + Extent + Journaling (write-ahead logs)
- Extent — consecutive disk blocks
- A file in ext file systems — a list of extents
- Journal
  - Write-ahead logs — performs writes as in LFS
  - Apply the log to the target location when appropriate
- Block group
  - Modern H.D.Ds do not have the concept of "cylinders"
  - They label neighboring sectors with consecutive block addresses
  - Does not work for SSDs given the internal log-structured management of block addresses

# Recap: flash SSDs, NVM-based SSDs

- Asymmetric read/write behavior/performance
- Wear-out faster than traditional magnetic disks
- Another layer of indirection is introduced
  - Intensify log-on-log issues
  - We need to revise the file system design

# The introduction of virtual file system interface

**User-space**

Applications, user-space libraries

.....open, ..close, ..read, ..write, .....

**Virtual File System**

open, close, read, write, ...

**File system #1 (e.g. ext4)**

**File system #2 (e.g. f2fs)**

read/write - 0, 512, 4096, ... (block address)

**Device independent I/O interface (e.g. ioctl)**

data

**Device Driver**

read/write

**Device Driver**

data, block addresses

read/write - block addresses

**Device Controller**

**FTL**

**Device Controller**

**HDD #1**

**SSD**

**Hardware**

# Outline

- NFS
- Google file system

# Network File System



# The introduction of virtual file system interface

**User-space**

Applications, user-space libraries

.....open, ..close, ..read, ..write, .....

**Virtual File System**

open, close, read, write, ...

open, close, read, write, ...

**File system #1 (e.g. ext4)**

**File system #2 (e.g. f2fs)**

**File system #3 — NFS**

read/write — 0, 512, 4096, ... (block address)

open, close, read, write, ...

**Kernel**

**Device independent I/O interface (e.g. ioctl)**

**Network Stack**

data

read/write — block addresses

data

**Device Driver**

**Device Driver**

**Network Device Driver**

read/write — block addresses

**Device Controller**

**FTL**

**Device Controller**

**Device Controller**

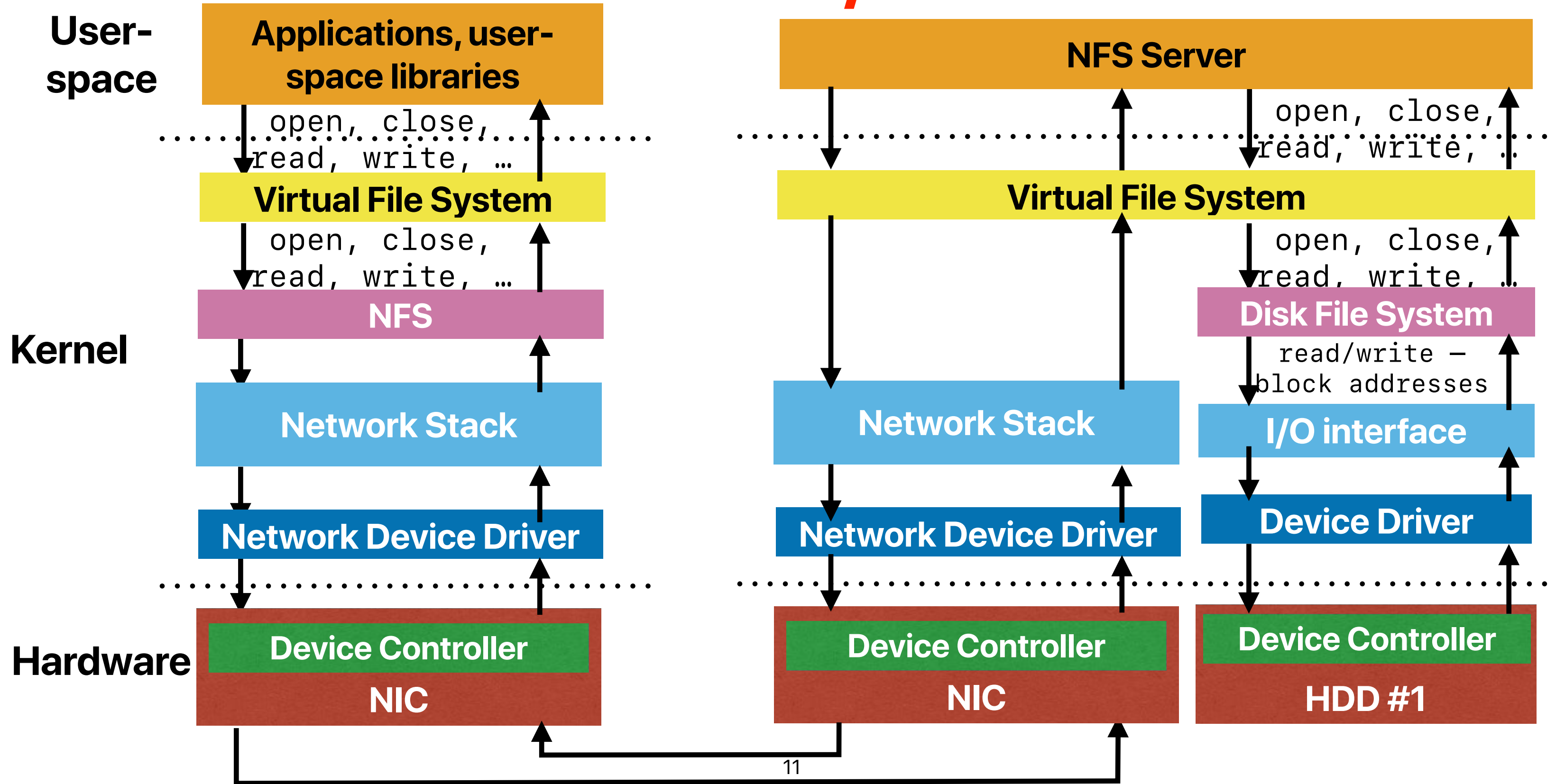
**HDD #1**

**SSD**

**NIC**

**Hardware**

# NFS Client/Server



# How does NFS handle a file?

- The client gives it's file system a tuple to describe data
  - Volume: Identify which server contains the file — represented by the mount point in UNIX
  - inode: Where in the server
  - generation number: version number of the file
- The local file system forwards the requests to the server
- The server response the client with file system attributes as local disks

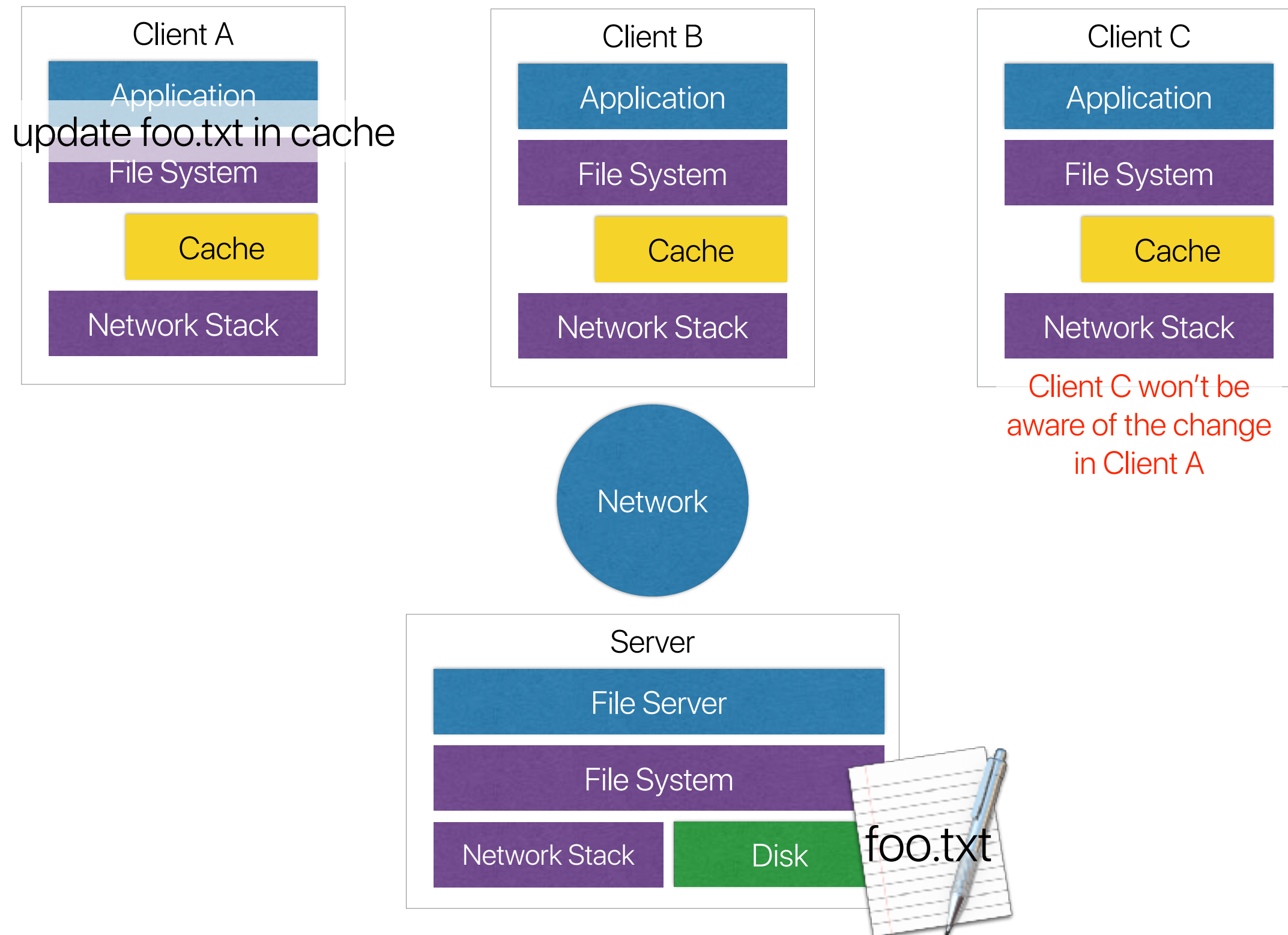
# Caching

- NFS operations are expensive
  - Lots of network round-trips
  - NFS server is a user-space daemon
- With caching on the clients
  - Only the first reference needs network communication
  - Later requests can be satisfied in local memory

# Idempotent operations

- Given the same input, always give the same output regardless how many times the operation is employed
- You only need to retry the same operation if it failed

# Think about this



# Solution

- Flush-on-close: flush all write buffer contents when close the file
  - Later open operations will get the latest content
- Force-getattr:
  - Open a file requires getattr from server to check timestamps
  - attribute cache to remedy the performance

# **The Google File System**

**Sanjay Ghemawat, Howard Gobioff, and**

**Shun-Tak Leung**

**Google**



# Why we care about GFS

- Conventional file systems do not fit the demand of data centers
- Workloads in data centers are different from conventional computers
  - Storage based on inexpensive disks that fail frequently
  - Many large files in contrast to small files for personal data
  - Primarily reading streams of data
  - Sequential writes appending to the end of existing files
  - Must support multiple concurrent operations
  - Bandwidth is more critical than latency

# Data-center workloads for GFS

- Google Search (Web Search for a Planet: The Google Cluster Architecture, IEEE Micro, vol. 23, 2003)
- MapReduce (MapReduce: Simplified Data Processing on Large Clusters, OSDI 2004)
  - Large-scale machine learning problems
  - Extraction of user data for popular queries
  - Extraction of properties of web pages for new experiments and products
  - Large-scale graph computations
- BigTable (Bigtable: A Distributed Storage System for Structured Data, OSDI 2006)
  - Google analytics
  - Google earth
  - Personalized search

# What GFS proposes?

- Maintaining the same interface
  - The same function calls
  - The same hierarchical directory/files
- Files are decomposed into large chunks (e.g. 64MB) with replicas
- Hierarchical namespace implemented with flat structure
- Master/chunkservers/clients

# Latency Numbers Every Programmer Should Know

Operations	Latency (ns)	Latency (us)	Latency (ms)	
L1 cache reference	0.5 ns			~ 1 CPU cycle
Branch mispredict	5 ns			
L2 cache reference	7 ns			14x L1 cache
Mutex lock/unlock	25 ns			
Main memory reference	100 ns			20x L2 cache, 200x L1 cache
Compress 1K bytes with Zippy	3,000 ns	3 us		
Send 1K bytes over 1 Gbps network	10,000 ns	10 us		
Read 4K randomly from SSD*	150,000 ns	150 us		~1GB/sec SSD
Read 1 MB sequentially from memory	250,000 ns	250 us		
Round trip within same datacenter	500,000 ns	500 us		
Read 1 MB sequentially from SSD*	1,000,000 ns	1,000 us	1 ms	~1GB/sec SSD, 4X memory
Read 512B from disk	10,000,000 ns	10,000 us	10 ms	20x datacenter roundtrip
Read 1 MB sequentially from disk	20,000,000 ns	20,000 us	20 ms	80x memory, 20X SSD
Send packet CA-Netherlands-CA	150,000,000 ns	150,000 us	150 ms	

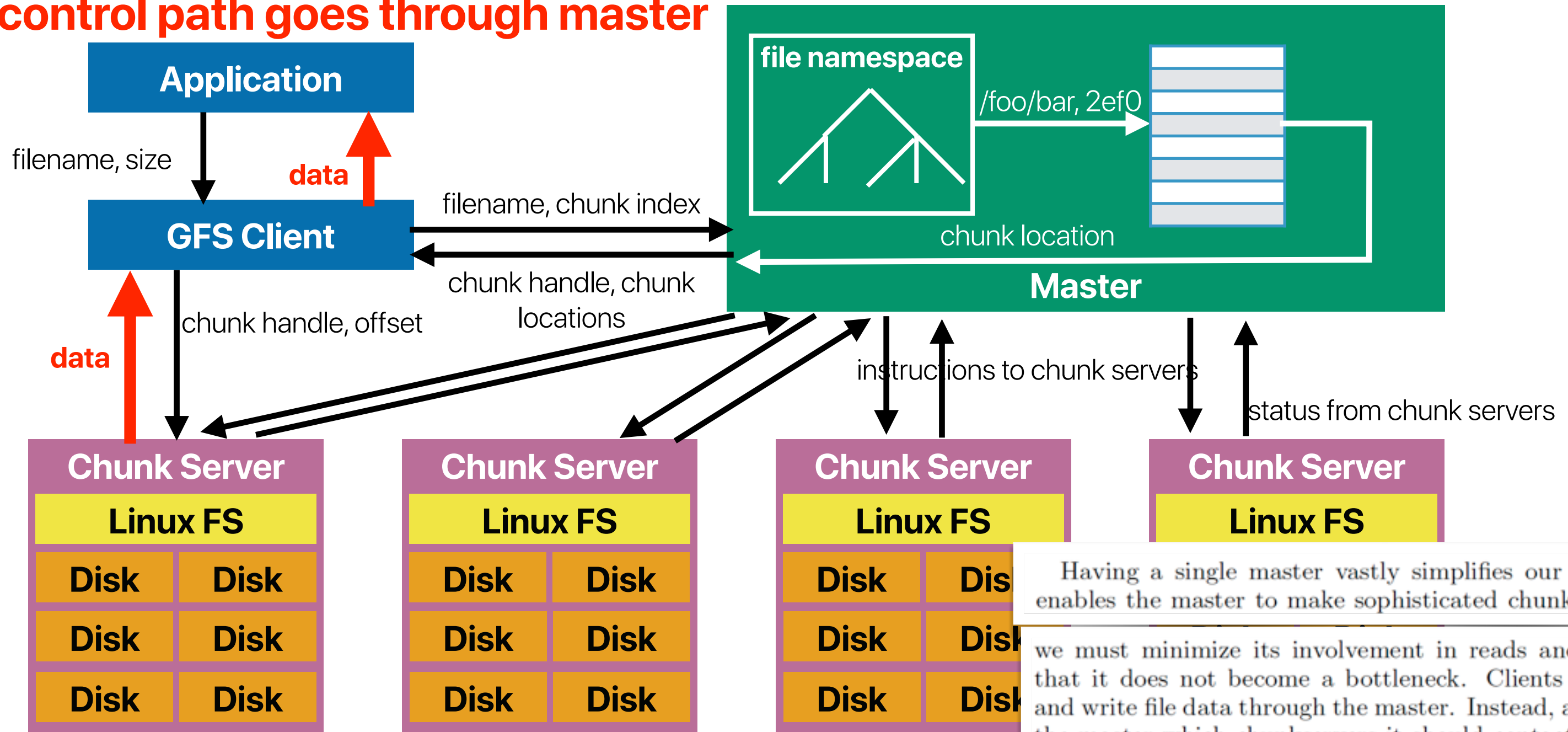
# Flat file system structure

- Directories are illusions
- Namespace maintained like a hash table

Unlike many traditional file systems, GFS does not have a per-directory data structure that lists all the files in that directory. Nor does it support aliases for the same file or directory (i.e, hard or symbolic links in Unix terms). GFS logically represents its namespace as a lookup table mapping full pathnames to metadata. With prefix compression, this

# GFS Architecture

decoupled data and control paths —  
only control path goes through master



load balancing, replicas among chunkservers

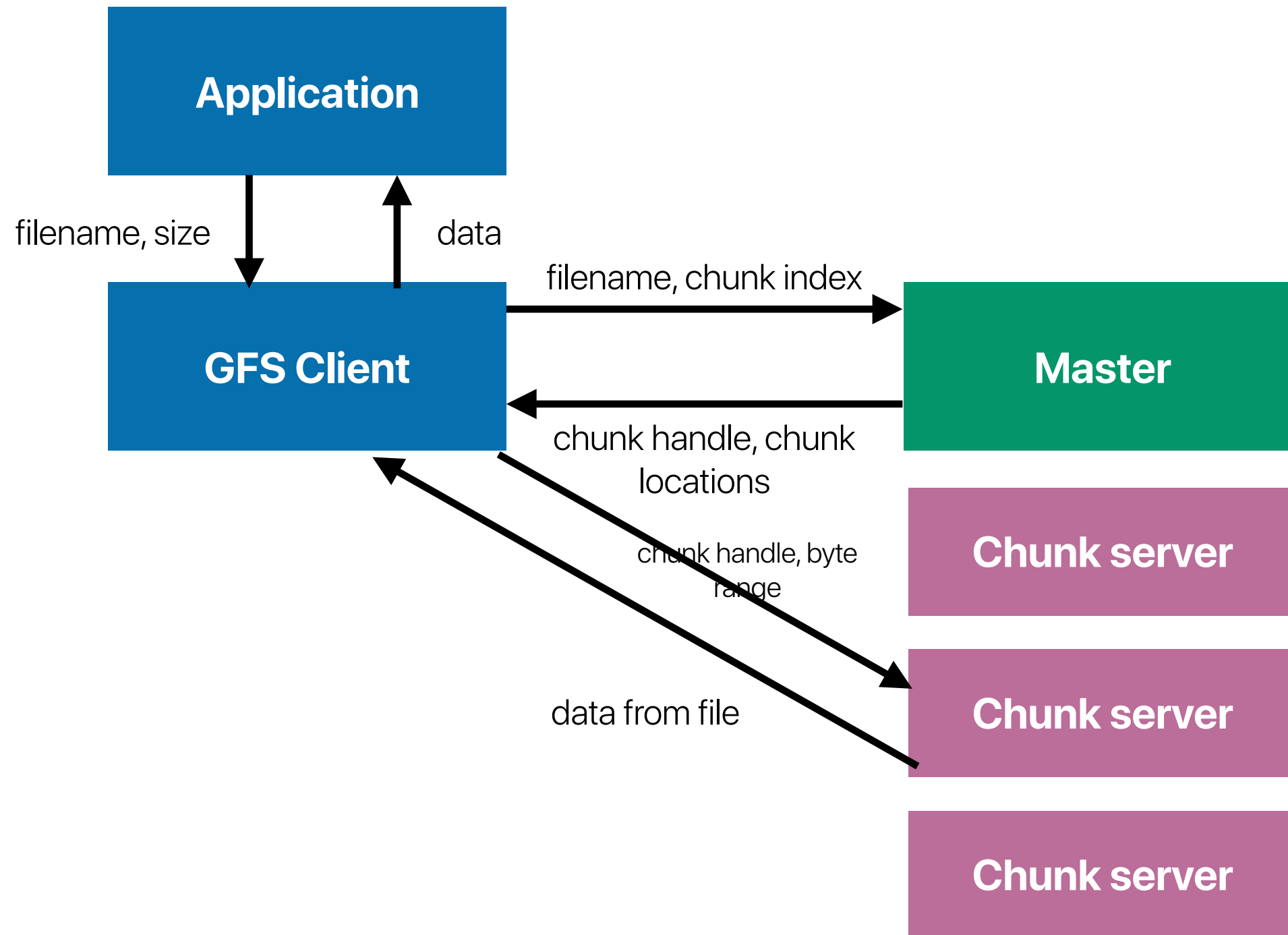
Having a single master vastly simplifies our design and enables the master to make sophisticated chunk placement

we must minimize its involvement in reads and writes so that it does not become a bottleneck. Clients never read and write file data through the master. Instead, a client asks the master which chunkservers it should contact. It caches this information for a limited time and interacts with the chunkservers directly for many subsequent operations.

# Distributed architecture

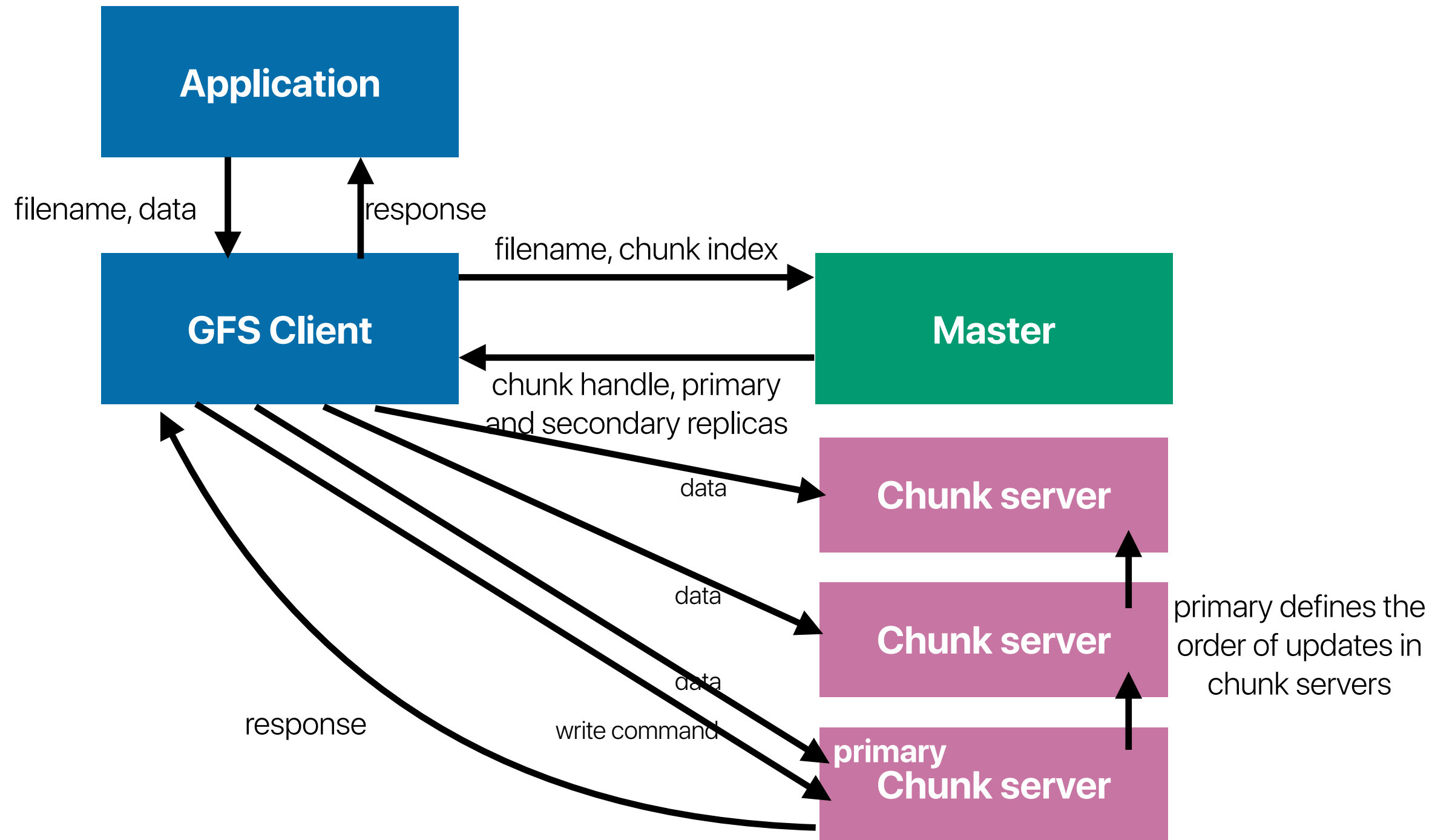
- Single master
  - maintains file system metadata including namespace, mapping, access control and chunk locations.
  - controls system wide activities including garbage collection and chunk migration.
- Chunkserver
  - stores data chunks
  - chunks are replicated to improve reliability (3 replicas)
- Client
  - APIs to interact with applications
  - interacts with masters for control operations
  - interacts with chunkservers for accessing data
  - Can run on chunkservers

# Reading data in GFS





# Writing data in GFS



# GFS: Relaxed Consistency model

- Distributed, simple, efficient
- Filename/metadata updates/creates are atomic
- Consistency modes

	Write — write to a specific offset	Append — write to the end of a file
Serial success	Defined	Defined with interspersed with inconsistent
Concurrent success	Consistent but undefined	
Failure	inconsistent	

- Consistent: all replicas have the same value
- Defined: replica reflects the mutation, consistent
- Applications need to deal with inconsistent cases themselves

# Real world, industry experience

- Linux problems (section 7)
  - Linux driver issues — disks do not report their capabilities honestly
  - The cost of fsync — proportion to file size rather than updated chunk size
  - Single reader-writer lock for mmap
  - Due to the open-source nature of Linux, they can fix it and contribute to the rest of the community
- **GFS is not open-sourced**

system behavior. When appropriate, we improve the kernel and share the changes with the open source community.

# Single master design

- GFS claims this will not be a bottleneck
- In-memory data structure for fast access
- Only involved in metadata operations — decoupled data/control paths
- Client cache
- What if the master server fails?

# The evolution of GFS

- Mentioned in "Spanner: Google's Globally-Distributed Database", OSDI 2012 — "tablet's state is stored in set of B-tree-like files and a write-ahead log, all on a distributed file system called Colossus (the successor to the Google File System)"
- Single master

proportionate increase in the amount of metadata the master had to maintain. Also, operations such as scanning the metadata to look for recoveries all scaled linearly with the volume of data. So the amount of work required of the master grew substantially. The amount of storage needed to retain all that information grew as well.

In addition, this proved to be a bottleneck for the clients, even though the clients issue few metadata operations themselves—for example, a client talks to the master whenever it does an open. When you have thousands of clients all talking to the master at the same time, given that the master is capable of doing only a few thousand operations a second, the average client isn't able to command all that many operations per second. Also bear in mind that there are applications such as MapReduce, where you might suddenly have a thousand tasks, each wanting to open a number of files. Obviously, it would take a long time to handle all those requests, and the master would be under a fair amount of duress.

acmqueue

Case Study  
GFS: Evolution on Fast-forward

A discussion between Kirk McKusick and Sean Quinlan about the origin and evolution of the Google File System.

**MCKUSICK** And historically you've had one cell per data center, right?

**QUINLAN** That was initially the goal, but it didn't work out like that to a large extent—partly because of the limitations of the single-master design and partly because isolation proved to be difficult. As a consequence, people generally ended up with more than one cell per data center. We also ended up doing what we call a "multi-cell" approach, which basically made it possible to put multiple GFS masters on top of a pool of chunkservers. That way, the chunkservers could be configured to have, say, eight GFS masters assigned to them, and that would give you at least one pool of underlying storage—with multiple master heads on it, if you will. Then the application was responsible for partitioning data across those different cells.



# The evolution of GFS

- Support for smaller chunk size — gmail

**QUINLAN** The distributed master certainly allows you to grow file counts, in line with the number of machines you're willing to throw at it. That certainly helps.

One of the appeals of the distributed multimaster model is that if you scale everything up by two orders of magnitude, then getting down to a 1-MB average file size is going to be a lot different from having a 64-MB average file size. If you end up going below 1 MB, then you're also going to run into other issues that you really need to be careful about. For example, if you end up having to read 10,000 10-KB files, you're going to be doing a lot more seeking than if you're just reading 100 1-MB files.

My gut feeling is that if you design for an average 1-MB file size, then that should provide for a much larger class of things than does a design that assumes a 64-MB average file size. Ideally, you would like to imagine a system that goes all the way down to much smaller file sizes, but 1 MB seems a reasonable compromise in our environment.

**MCKUSICK** What have you been doing to design GFS to work with 1-MB files?

**QUINLAN** We haven't been doing anything with the existing GFS design. Our distributed master system that will provide for 1-MB files is essentially a whole new design. That way, we can aim for something on the order of 100 million files per master. You can also have hundreds of masters.

# Lots of other interesting topics

- snapshots
- namespace locking
- replica placement
- create, re-replication, re-balancing
- garbage collection
- stable replica detection
- data integrity
- diagnostic tools: logs are your friends

# Do they achieve their goals?

- Storage based on inexpensive disks that fail frequently — replication, distributed storage
- Many large files in contrast to small files for personal data — large chunk size
- Primarily reading streams of data — large chunk size
- Sequential writes appending to the end of existing files — large chunk size
- Must support multiple concurrent operations — flat structure
- Bandwidth is more critical than latency — large chunk size



# Why we care about GFS

- Conventional file systems do not fit the demand of data centers
- Workloads in data centers are different from conventional computers
  - Storage based on inexpensive disks that fail frequently
    - MapReduce is fault tolerant
  - Many large files in contrast to small files for personal data
    - MapReduce aims at processing large amount of data once
  - Primarily reading streams of data
    - MapReduce reads chunks of large files
  - Sequential writes appending to the end of existing files
    - Output file keep growing as workers keep writing
  - Must support multiple concurrent operations
    - MapReduce has thousands of workers simultaneously
  - Bandwidth is more critical than latency
    - MapReduce only wants to finish tasks within "reasonable" amount of time

# What's missing in GFS?

- GFS only supports consistency models
- Scalability — single master
- Only efficient in dealing with large data
- No geo-redundancy

## Human Error Investigated in California Blackout's Spread to Six Million

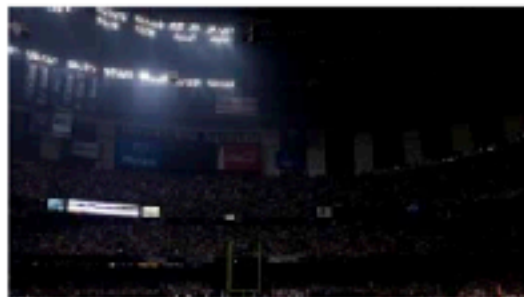
### 'Abnormality' caused power



Pat Yasinskas  
ESPN Staff Writer

NEW ORLEANS -- As he ran 108 yards for the longest kick in Super Bowl history, [Jacoby Jones](#) saw only one thing.

"Daylight," Jones said after helping the [Baltimore Ravens](#) defeat the [San Francisco 49ers](#). "Follow any avenue and it



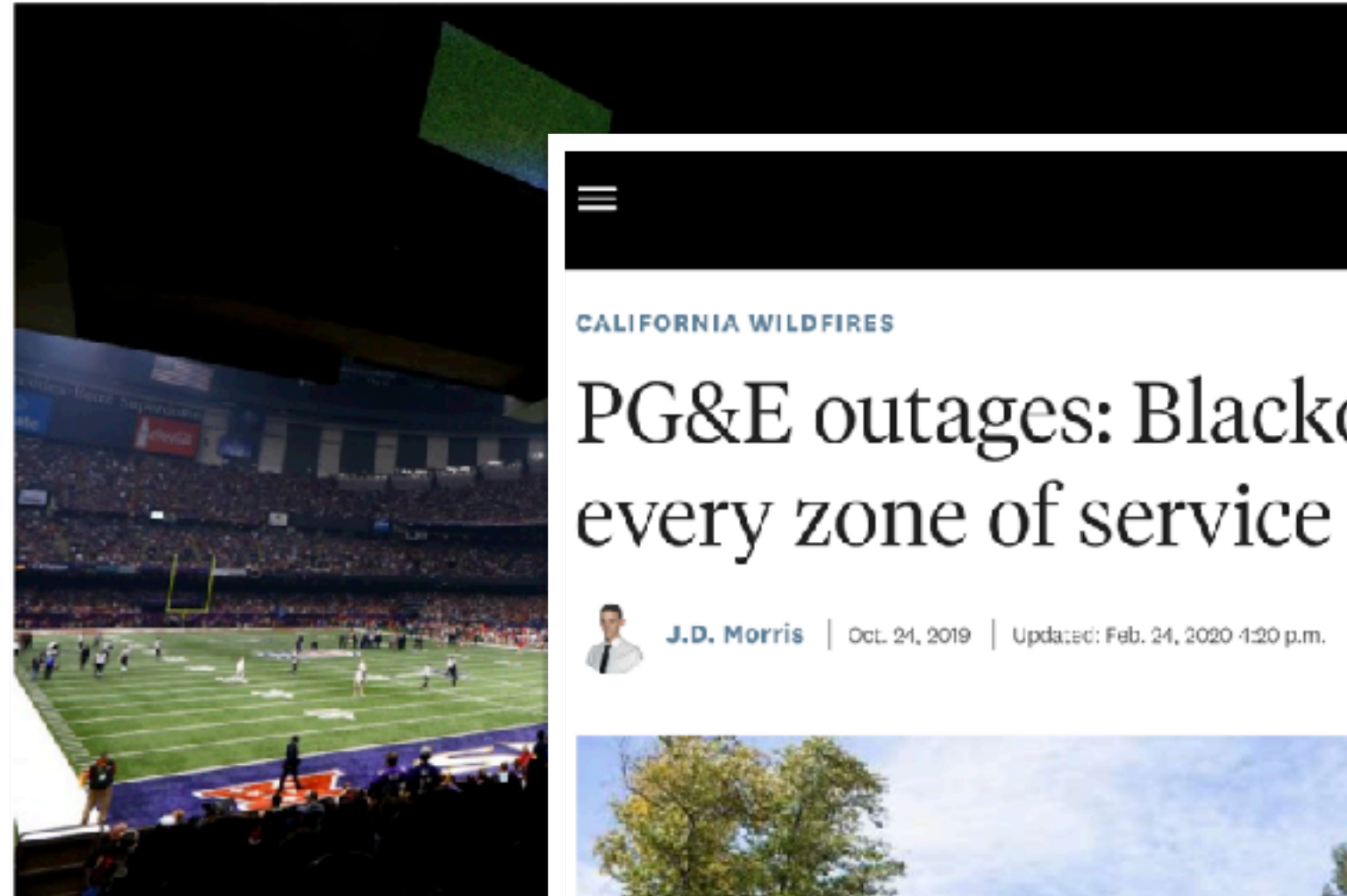
The majority of lights went out in the Superdome during the Super Bowl, causing a 34-minute delay. AP Photo/Mercio Sanchez

The irony was that the power outage, which lasted for a period of darkness, was a result of a human error.

Minutes after the game's kickoff, the lights in the Mercedes-Benz Superdome went out, and the game was stopped for 34 minutes, and the television broadcast was interrupted.

television broadcast was interrupted.

## Did Beyonce cause the Super Bowl blackout?



Fans look on to the field after a sudden power outage during the Super Bowl in New Orleans. (Getty Images)

By MEREDITH BLAKE FEB. 4, 2013 | 12 AM

Beyonce shut it down during the

Only a few minutes after the pop

### CALIFORNIA WILDFIRES

## PG&E outages: Blackouts could hit nearly every zone of service area by Sunday



J.D. Morris | Oct. 24, 2019 | Updated: Feb. 24, 2020 4:20 p.m.



# **Windows Azure Storage: A Highly Available Cloud Storage Service with Strong Consistency**

**Brad Calder, Ju Wang, Aaron Ogus, Niranjana Nilakantan, Arild Skjolsvold, Sam McKelvie, Yikang Xu, Shashwat Srivastav, Jiesheng Wu, Huseyin Simitci, Jaidev Haridas, Chakravarthy Uddaraju, Hemal Khatri, Andrew Edwards, Vaman Bedekar, Shane Mainali, Rafay Abbasi, Arpit Agarwal, Mian Fahim ul Haq, Muhammad Ikram ul Haq, Deepali Bhardwaj, Sowmya Dayanand, Anitha Adusumilli, Marvin McNett, Sriram Sankaran, Kavitha Manivannan, Leonidas Rigas**  
**Microsoft**



# Data center workloads for WAS

		%Requests	%Capacity	%Ingress	%Egress
All	Blob	17.9	70.31	48.28	66.17
	Table	46.88	29.68	49.61	33.07
	Queue	35.22	0.01	2.11	0.76
Bing	Blob	0.46	60.45	16.73	29.11
	Table	98.48	39.55	83.14	70.79
	Queue	1.06	0	0.13	0.1
XBox GameSaves	Blob	99.68	99.99	99.84	99.88
	Table	0.32	0.01	0.16	0.12
	Queue	0	0	0	0
XBox Telemetry	Blob	26.78	19.57	50.25	11.26
	Table	44.98	80.43	49.25	88.29
	Queue	28.24	0	0.5	0.45
Zune	Blob	94.64	99.9	98.22	96.21
	Table	5.36	0.1	1.78	3.79
	Queue	0	0	0	0

# Why Windows Azure Storage

- A cloud service platform for social network search, video streaming, XBOX gaming, records management, and etc. in M\$.
  - Must tolerate many different data abstractions: blobs, tables and queues
  - Data types:
    - Blob(Binary Large Objects) storage: pictures, excel files, HTML files, virtual hard disks (VHDs), big data such as logs, database backups -- pretty much anything.
- Table: database tables
- Queue: store and retrieve messages. Queue messages can be up to 64 KB in size, and a queue can contain millions of messages. Queues are generally used to store lists of messages to be processed asynchronously.

Large

Large

Small

# Why Windows Azure Storage (cont.)

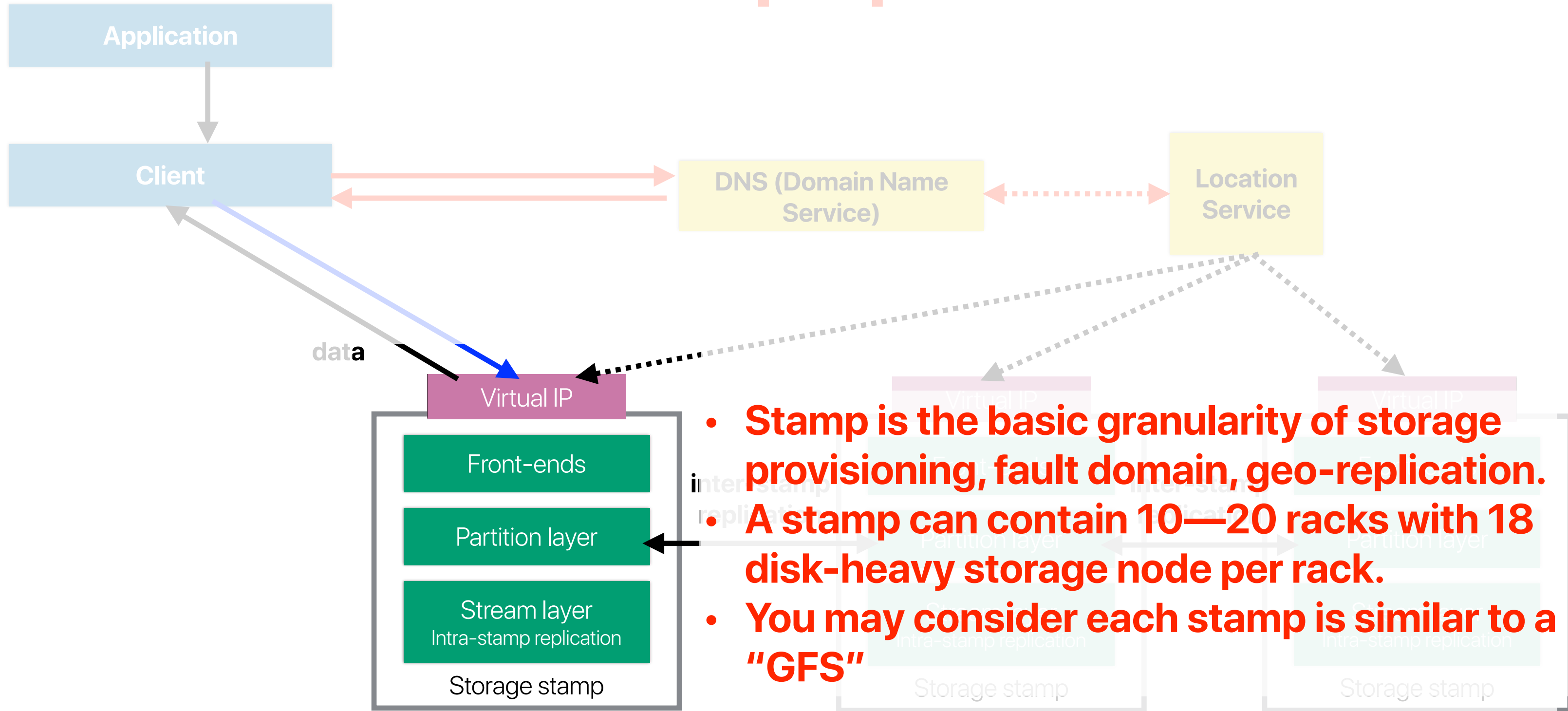
- Learning from feedbacks in existing cloud storage
  - Strong consistency
  - Global and scalable namespace/storage
  - Disaster recovery
  - Multi-tenancy and cost of storage

All problems in computer science can be solved by another level of  
indirection

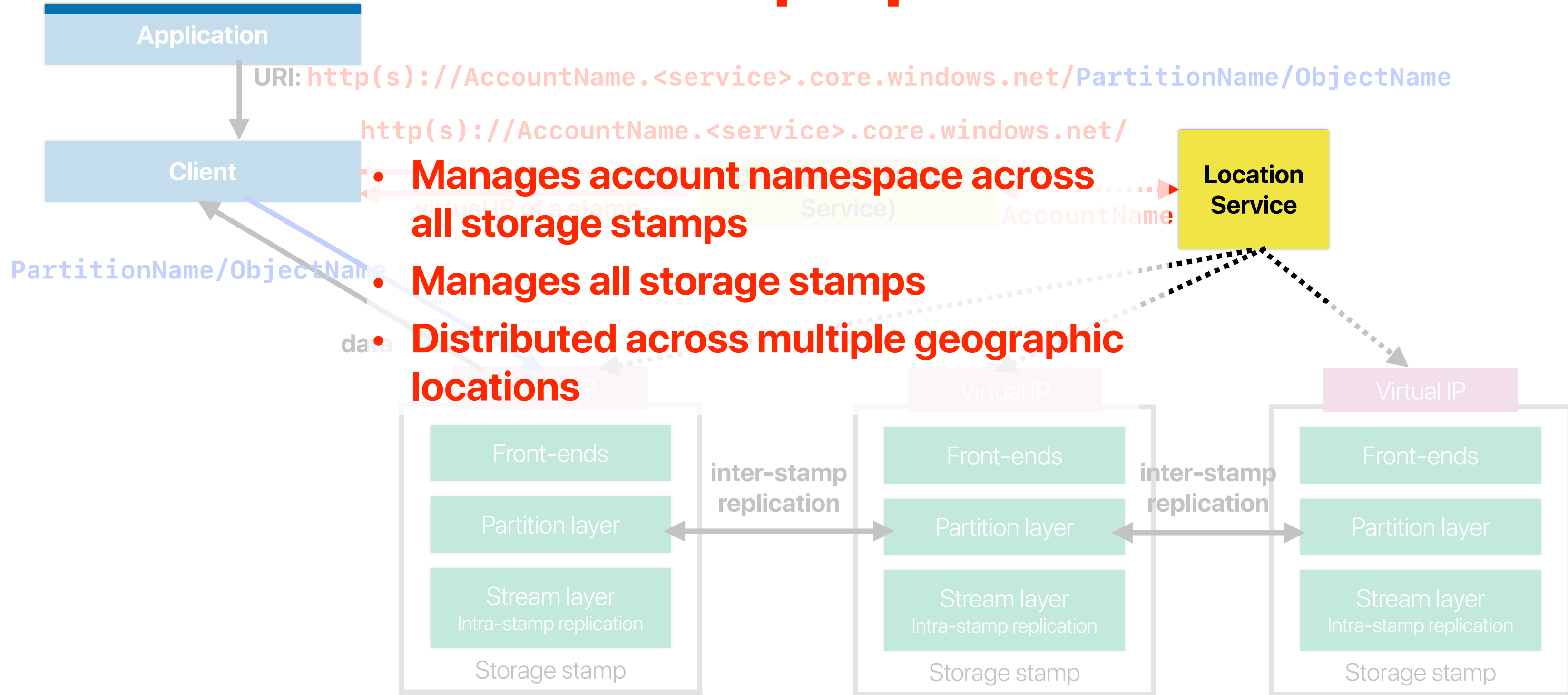
*–David Wheeler*



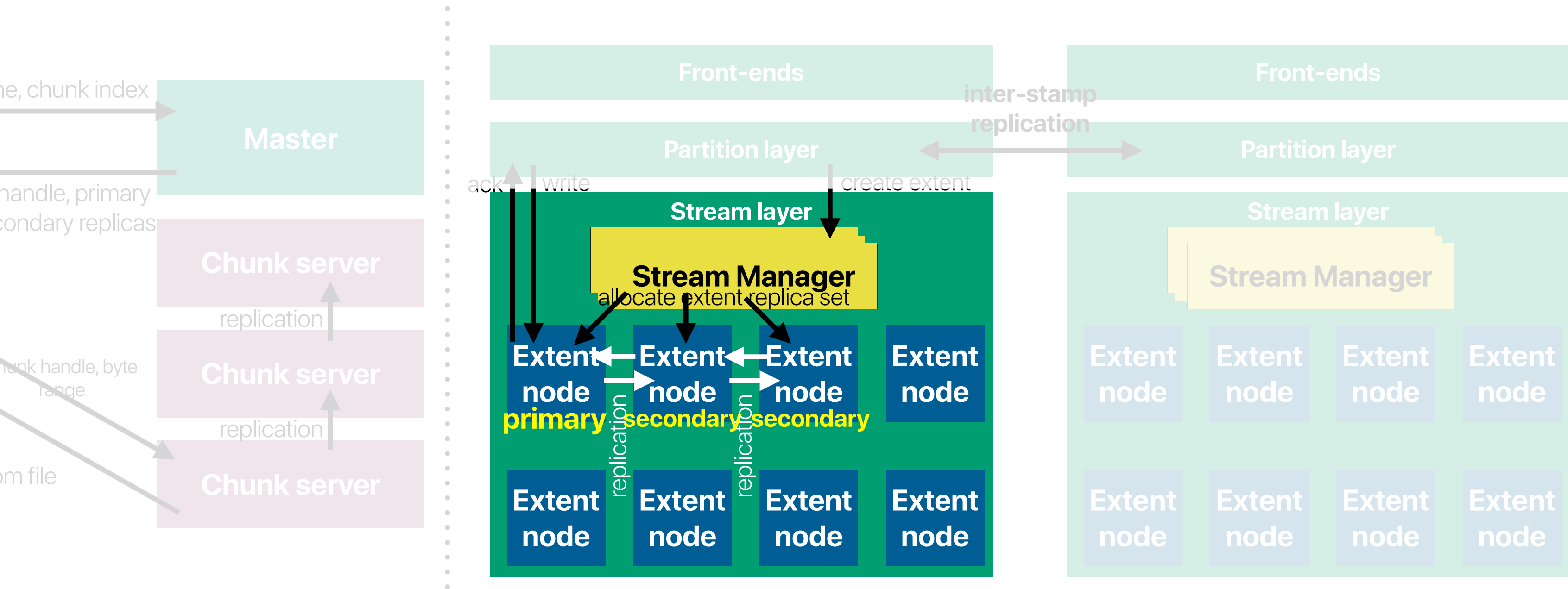
# What WAS proposes?



# What WAS proposes?



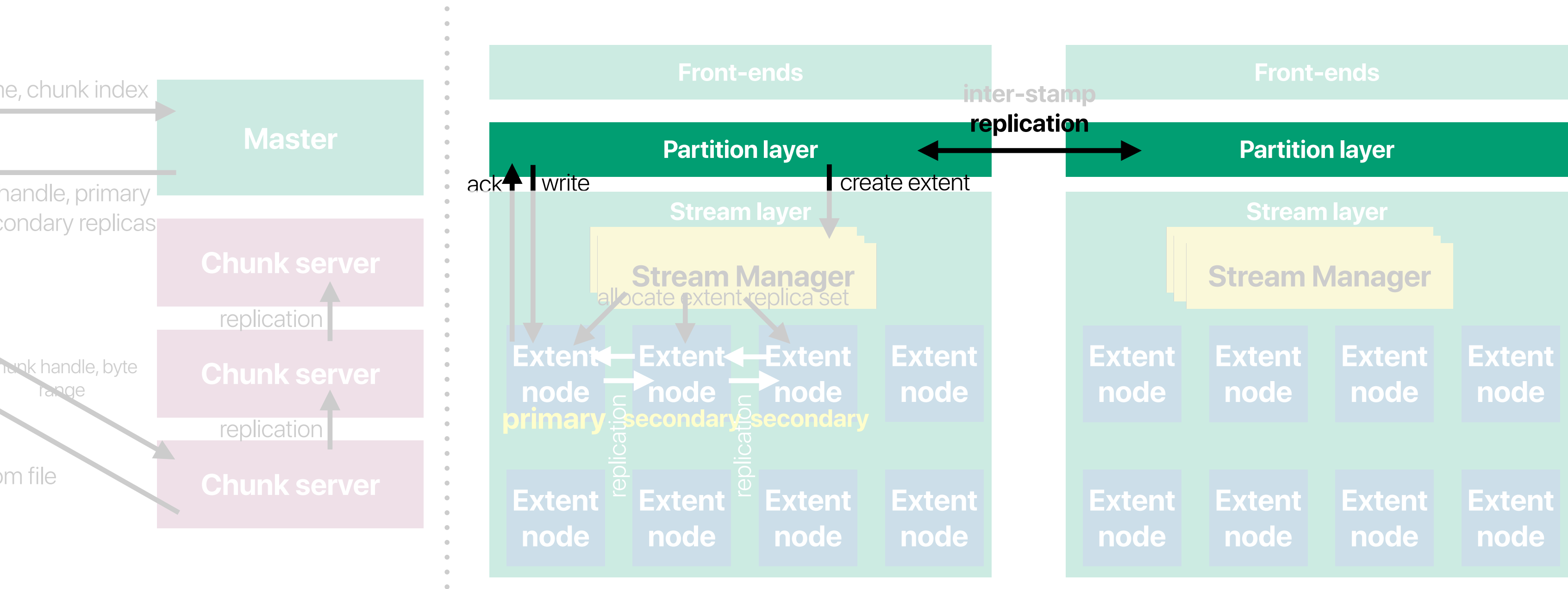
# GFS v.s. stamp in WAS



# Write failure

- Consider the case where 1 of 3 nodes handling a write fails and the current extent is sealed at latest commit boundary (end of extent) — that data will be on failed node
- new extent created
- SM chooses **three** new replicas to store extents
- client retries via new primary among the three new replicas
- failed node, upon restart, will coord w/ SM to synchronize its extent to the commit length decided upon

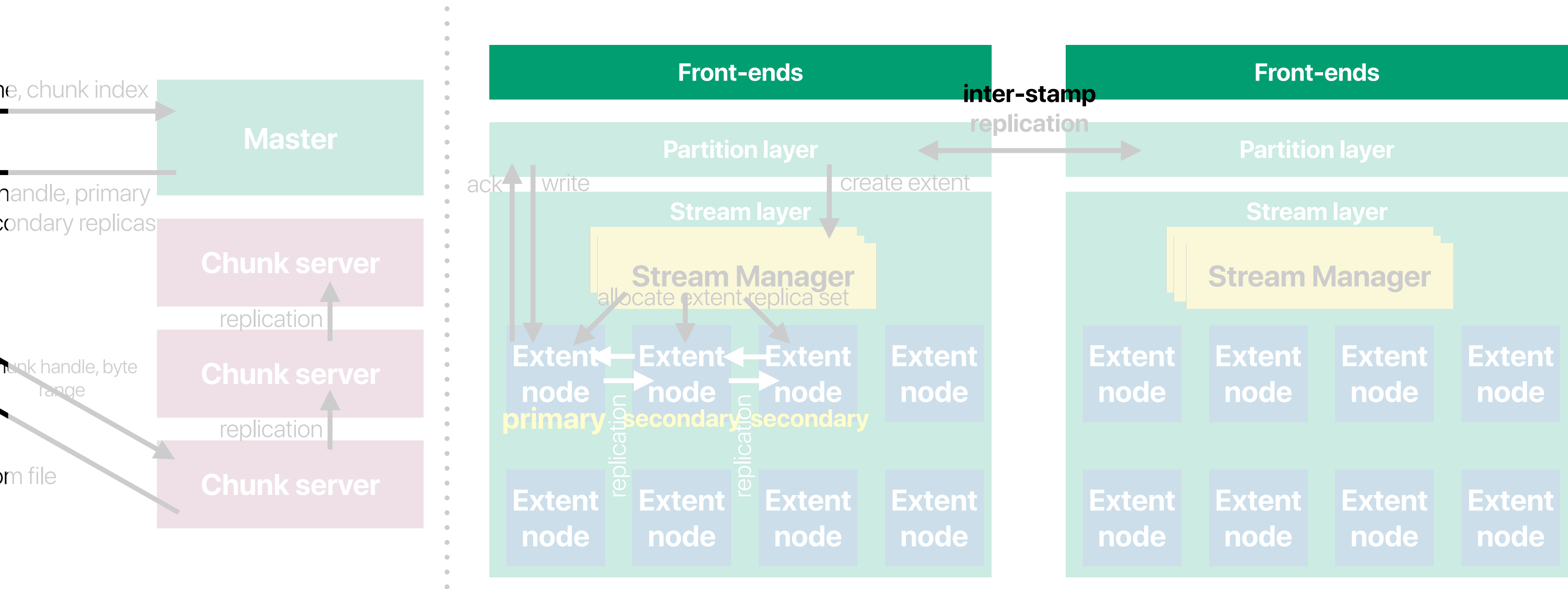
# GFS v.s. stamp in WAS



# Partition layer

- Managing high-level data abstractions
- Providing scalable object namespaces
- Providing transaction ordering and strong consistency for objects
- Storing object data on top of the stream layer
- Cache object data to reduce disk I/O

# GFS v.s. stamp in WAS

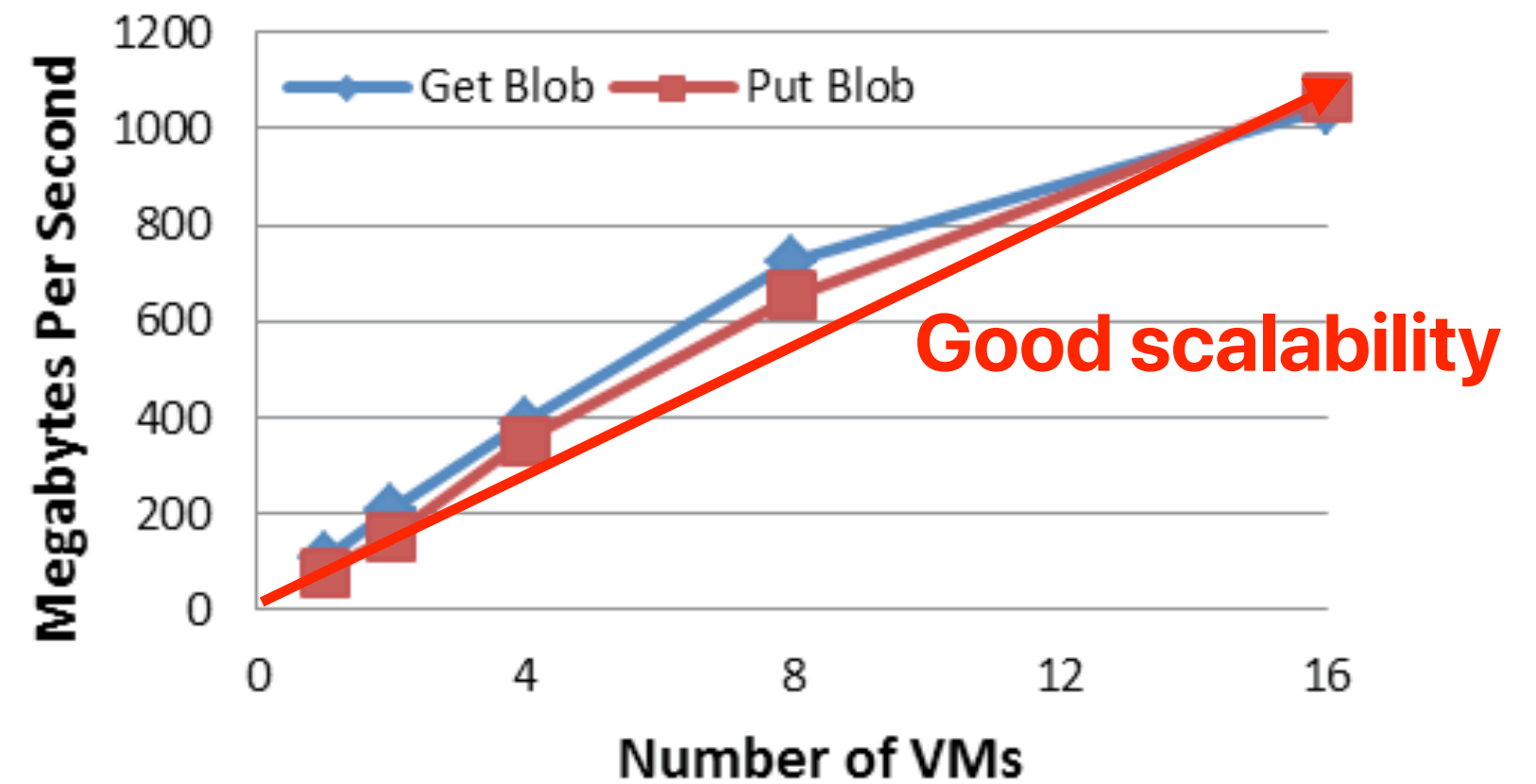
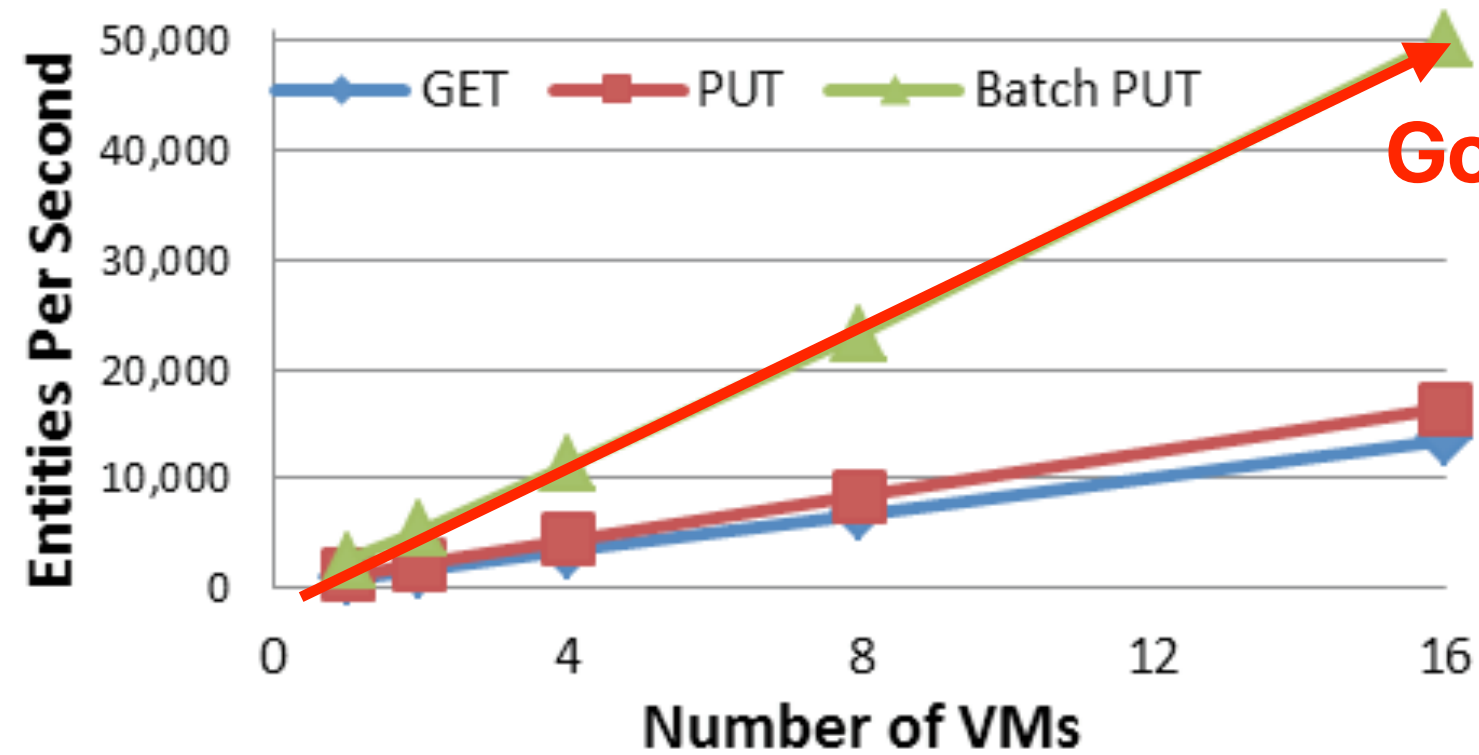


# Front-end layer

- A set of **stateless** servers taking incoming requests
  - Think about the benefits of stateless in NFS
- Keep partition maps to forward the request to the right server
  - A stamp can contain 10—20 racks with 18 disk-heavy storage node per rack
- Stream large objects directly from the stream layer and cache frequently accessed data for efficiency



# Are they doing well?



# GFS v.s. WAS

	GFS (OSDI 2003)	WAS (SOSP 2011)
File organizations	file chunk block	stream extent record
System architecture	master chunkserver	stream manager extent nodes
Data updates		append only updates
Consistency models	relaxed consistency	strong consistency
Data formats	files	multiple types of objects
Replications	intra-cluster replication	geo-replication
Usage of nodes	chunk server can perform both	separate computation and storage

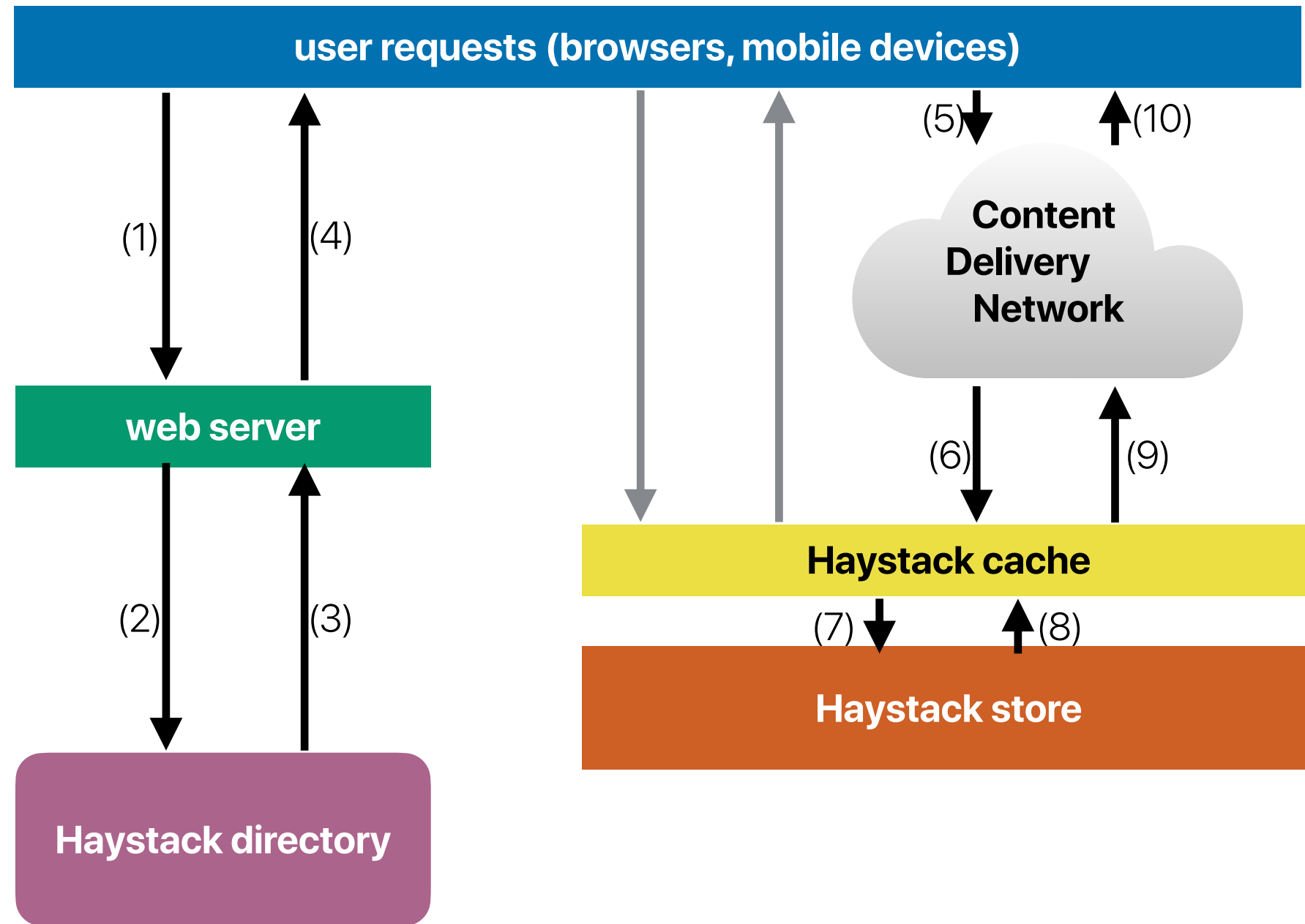
# **f4: Facebook's Warm BLOB Storage System**

**Subramanian Muralidhar, Wyatt Lloyd, Sabyasachi Roy, Cory Hill,  
Ernest Lin, Weiwen Liu, Satadru Pan, Shiva Shankar,  
Viswanath Sivakumar, Linpeng Tang, and Sanjeev Kumar.**

# The original NFS-based FB storage

- Within a data center with high-speed network, the round-trip latency of network accesses is not really a big deal
- However, the amount of metadata, especially directory metadata, is huge — cannot be cached
- As a result, each file access still requires ~ 10 inode/data requests from disks/network nodes — kill performance

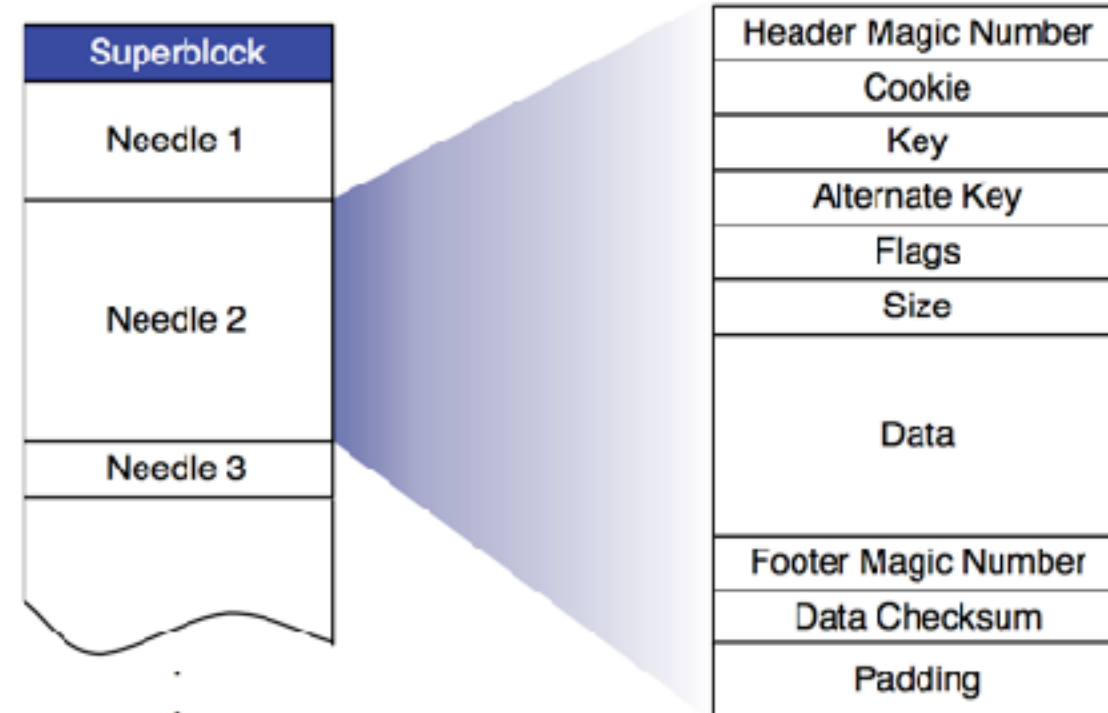
# Haystack



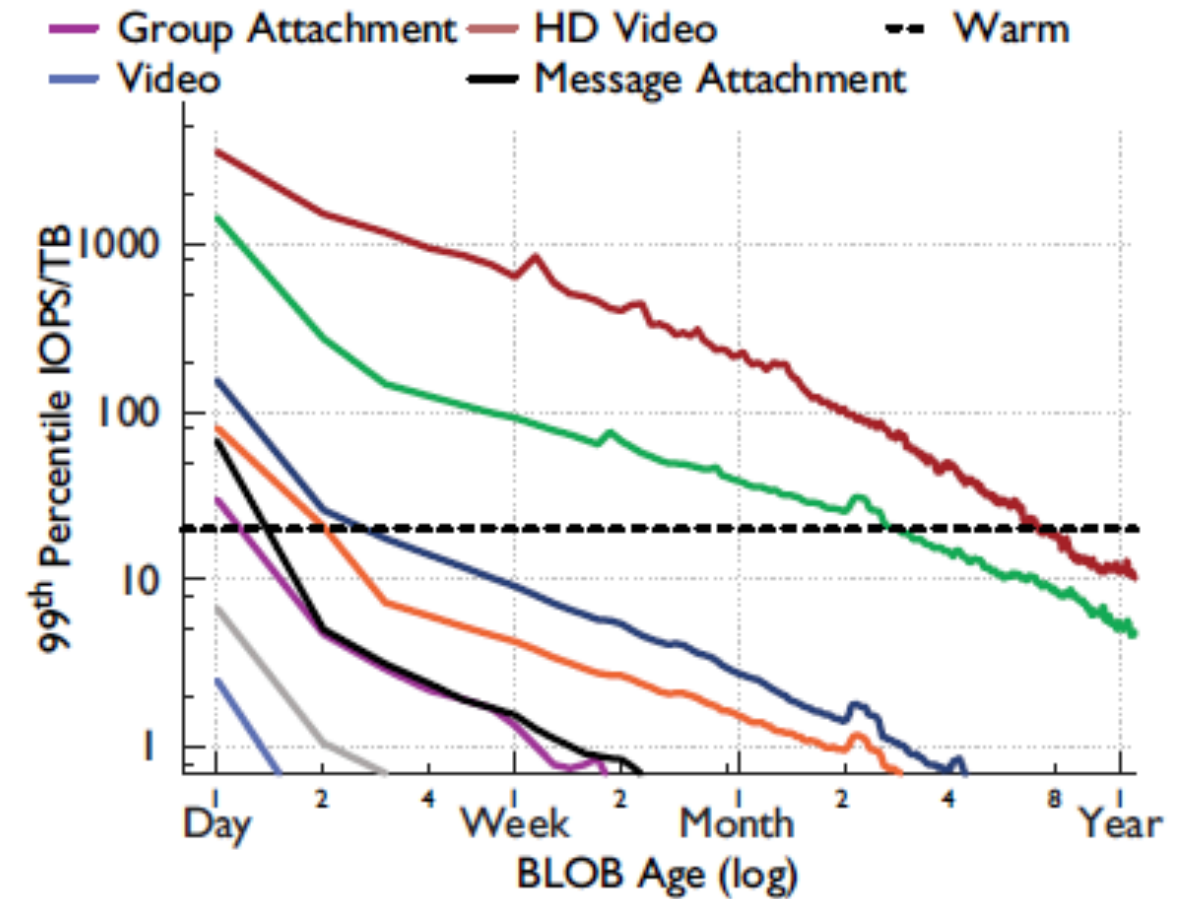
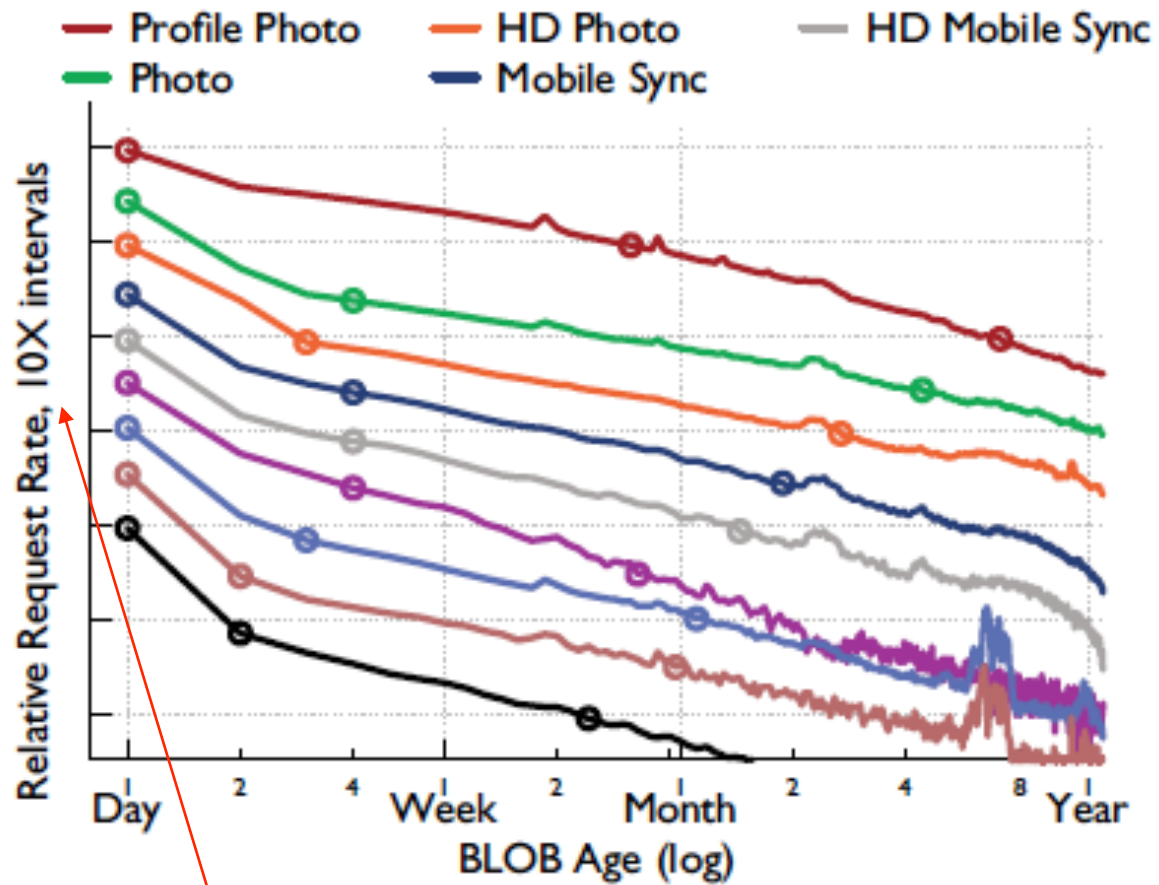
`http://<CDN>/<Cache>/<Machine ID>/<Logical volume,Photo>`

# Haystack

- Each storage unit provides 10TB of usable space, using RAID-6 — 20% redundancy for parity bits
  - Each storage split into 100 physical volumes (100GB)
  - Physical volumes on different machines grouped into logical volumes
  - A photo saved to a logical volume is written to all corresponding physical volumes — 3 replicas
- Each volume is actually just a large file
  - Needle represents a photo
  - Each needle is identified through the offset
  - Sealed (the same as WAS) once the file reaches 100GB



# "Temperature" of data



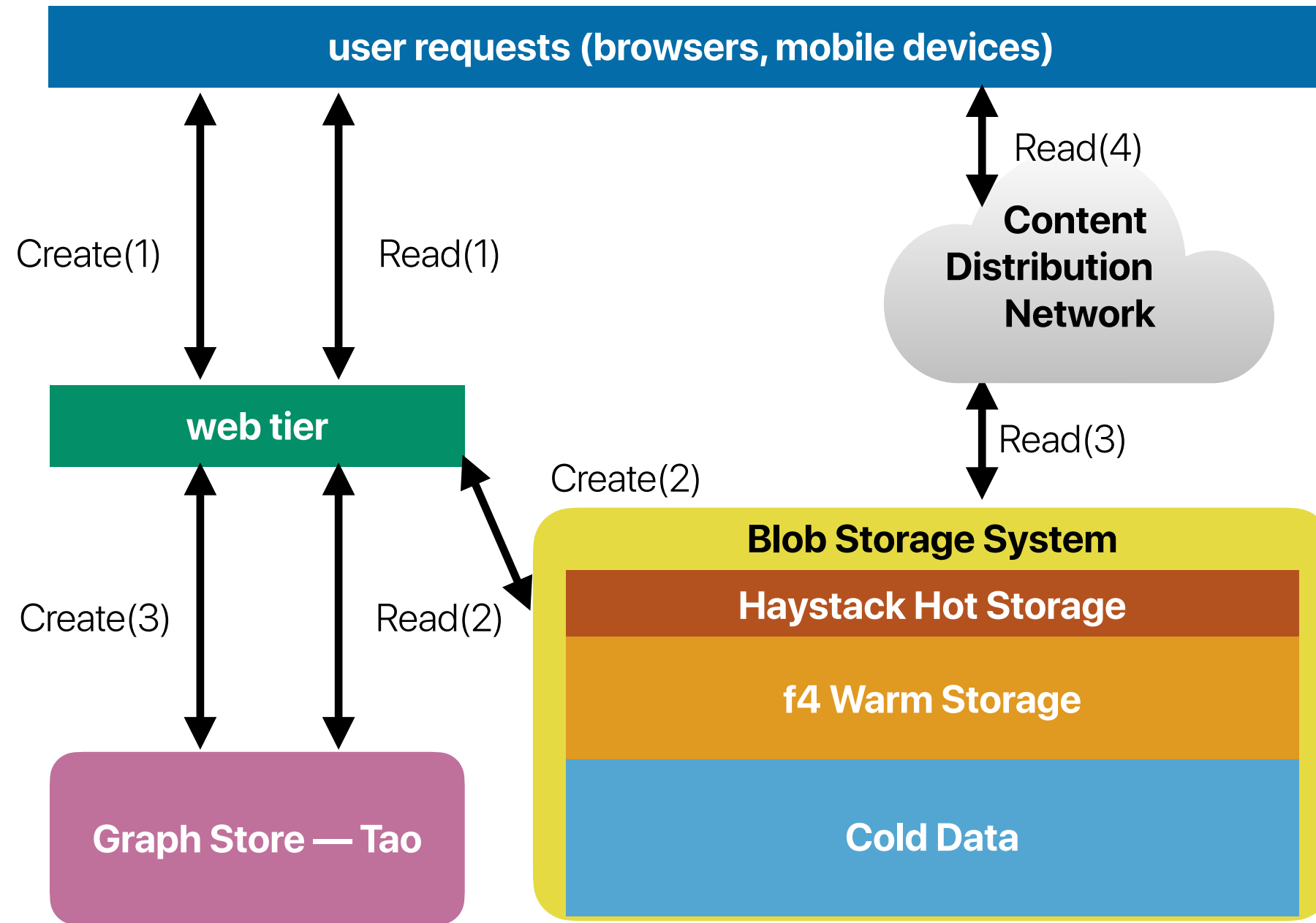
log scale — not encouraged to graph like this if you're writing a technical document or scientific paper

# "Temperature" of data

	Hot	Warm	Cold
Access Frequency	Most frequent	Less frequent	Rare
Pattern	Created often, delete often	Not so frequently read Not so frequently deleted Maybe read-only	Long-term storage, usually takes hours to retrieve
Size		65PB in 2014 and growing rapidly	

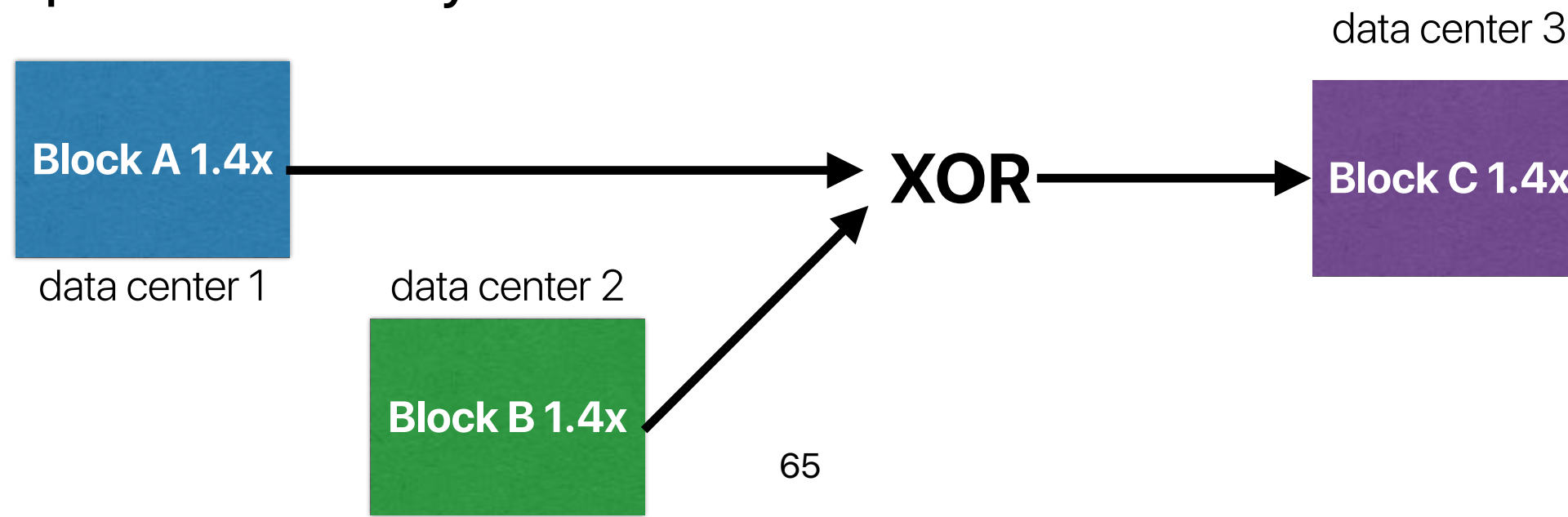


# Facebook storage architecture



# Storage efficiency

- Reed-Solomon erasure coding
  - Strips: 10GB data + 4GB parity — 1.4x space efficiency
  - One volume contains 10 strips
- XOR Geo-replication
  - Use XOR to reduce overhead further (e.g., Azure makes full copies)
  - Block A in DC1 + block B in DC2 -> parity block P in DC3
  - Any two blocks can be used to generate the third
  - 1.5x space efficiency
- $1.4 * 1.5 = 2.1x$  space efficiency in total



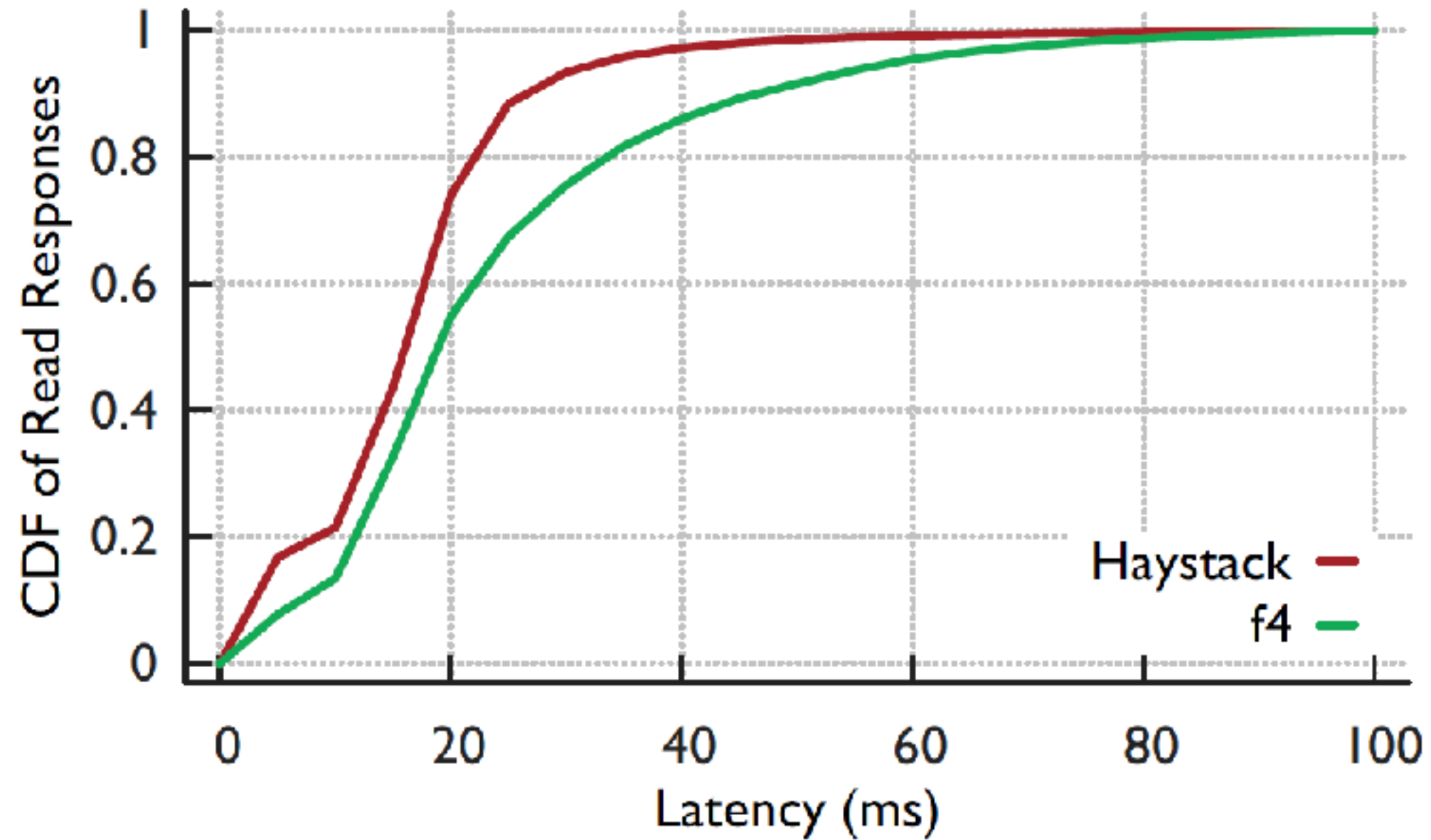
# Fault tolerance

- 1%-2% HDD fail in a year
  - replicate data across multiple disks
  - Use erasure coding for storage efficiency
    - $n$  blocks  $\rightarrow n + k$  blocks, can tolerate  $k$  simultaneous failures
    - higher cost for recovering data when there is a failure
- Host failures (periodically)
  - replicate coded blocks on different hosts
- Rack failures (multiple times/year)
  - replicate coded blocks on different racks
- Datacenter failures (rare, but catastrophic)
  - replicate blocks across data centers
  - use XOR to reduce overhead further (e.g., Azure makes full copies)
    - block A in DC1 + block B in DC2  $\rightarrow$  parity block P in DC3
    - any two blocks can be used to generate the third
- Index files
  - use normal triple replication (tiny, little benefit in coding them)

# What happens if fault occurs?

- Drive fails
  - Reconstruct blocks on another drive
  - Heavy disk, Network, CPU operation
  - one in background
- During failure, may need to reconstruct data online
  - rebuilder node reads BLOB from data + parity, reconstructs
  - only reads + reconstructs the BLOB (40KB), not the entire block (1GB)

# Performance of f4



# Cells

- Each cell contains 14 racks of 15 hosts, each host contains 30 4TB H.D.Ds.
- A unit of acquisition, deployment
- Storage for a set of volumes
- Similar to the idea of stamps

Make common case fast,  
make rare case correct