

Blizzard: Fast, Cloud-scale Block Storage for Cloud-oblivious Applications

Microsoft®
Research



James Mickens, Jeremy Elson,
Edmund B. Nightingale,
Darren Gehring, Krishna
Nareddy



Asim Kadav
Vijay Chidamaram



Bin Fan



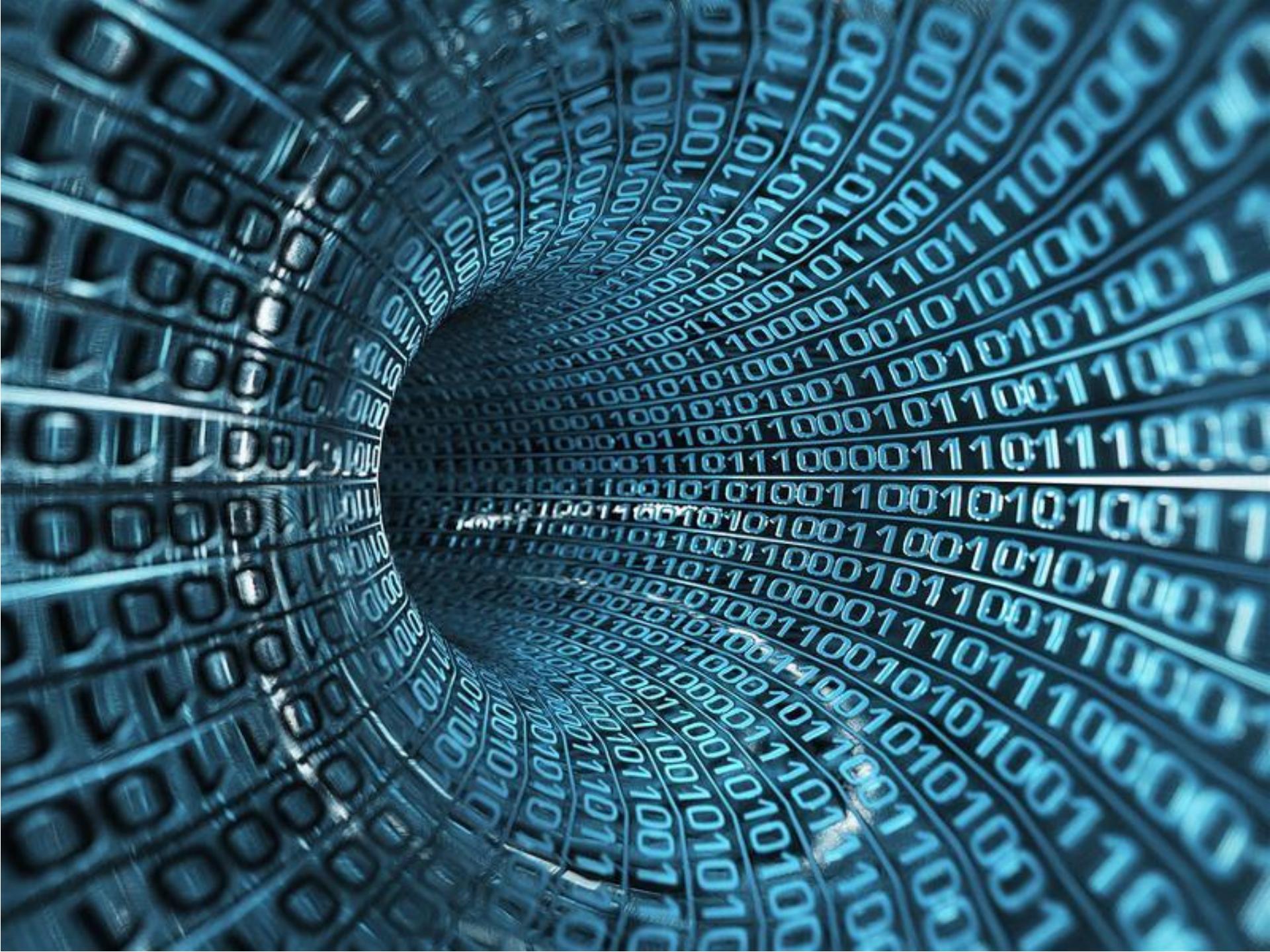
JOHNS HOPKINS



Osama Khan

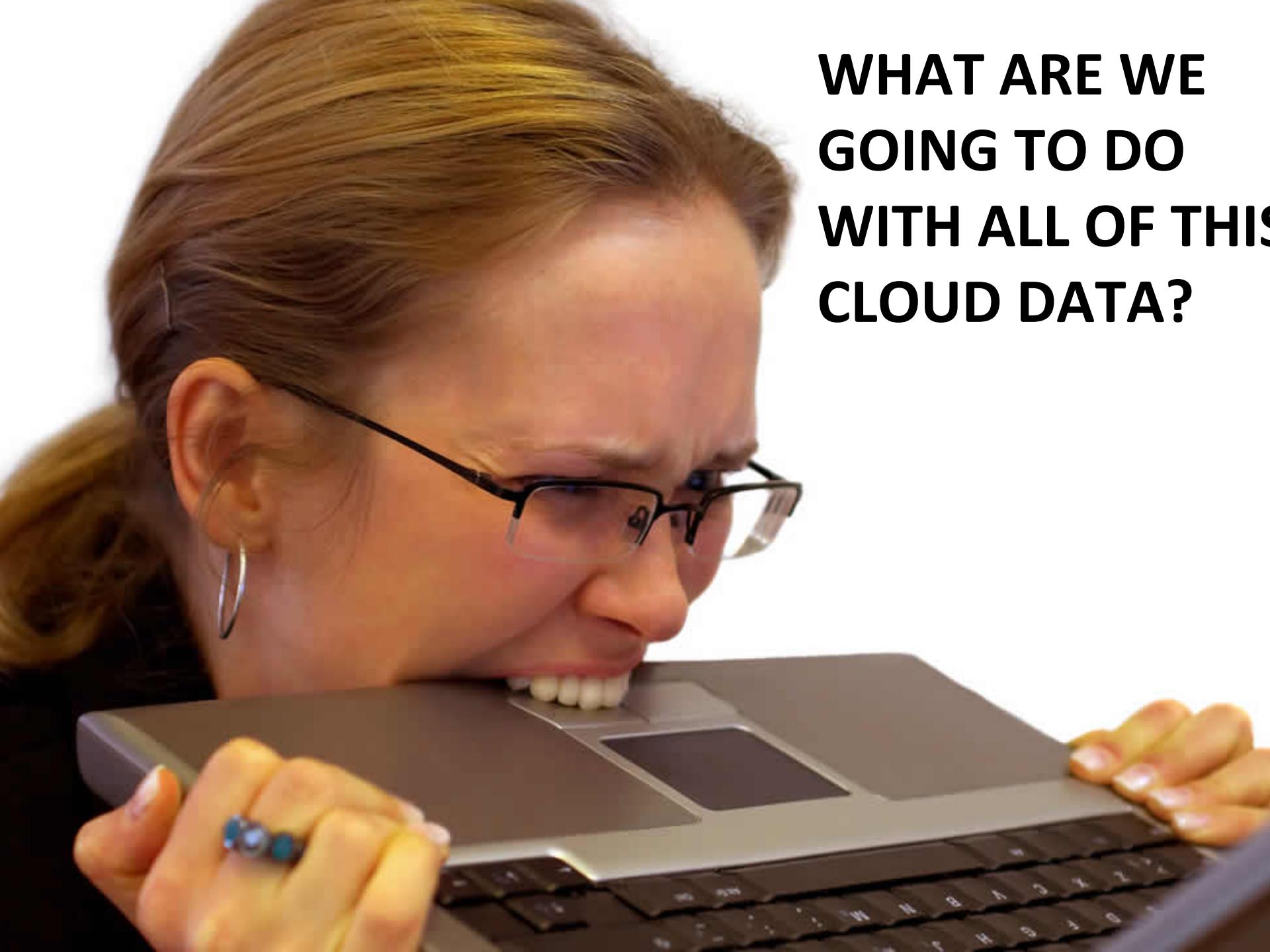


IT'S, UH, THE FUTURE

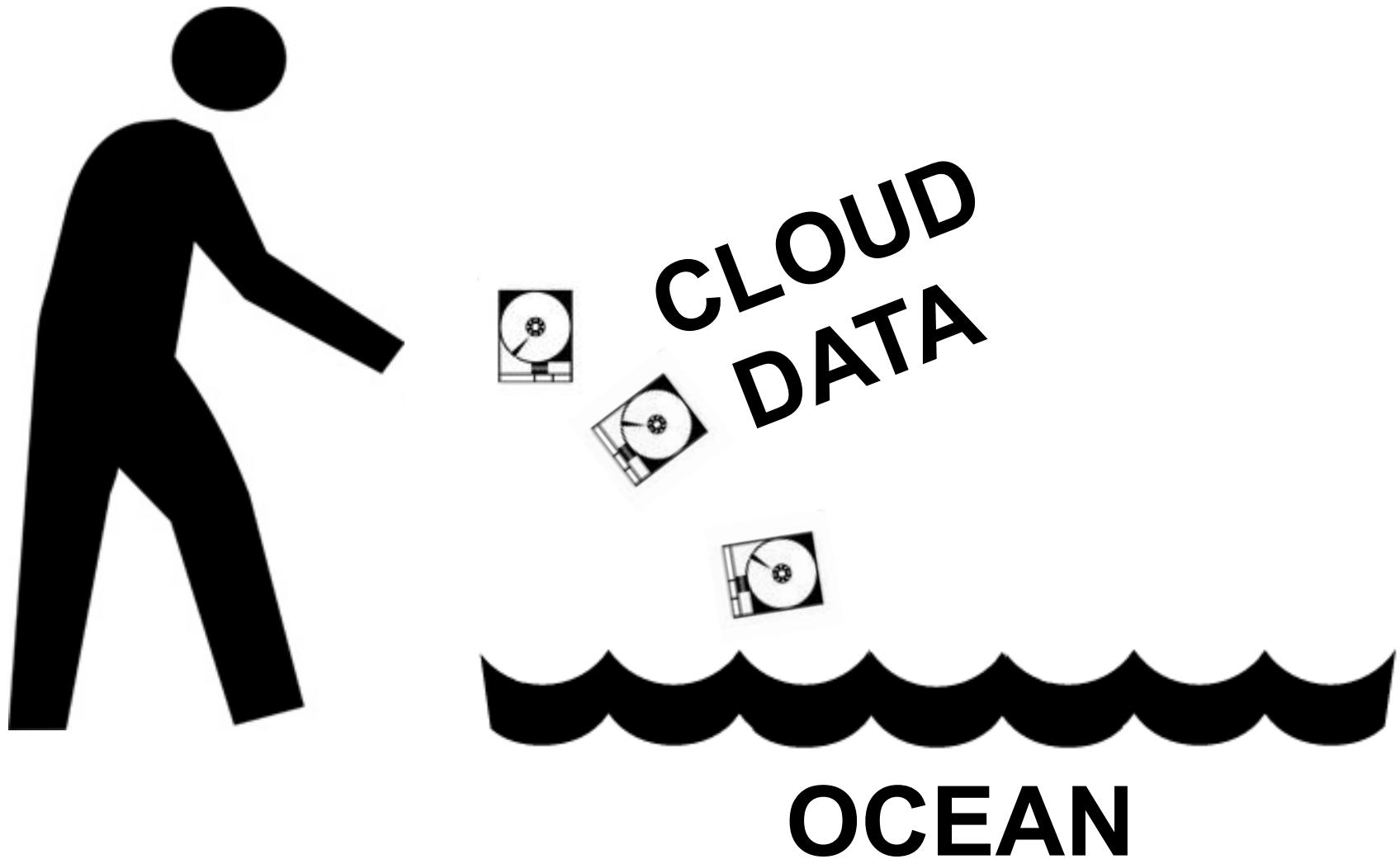






A close-up photograph of a woman with long brown hair pulled back, wearing black-rimmed glasses and a hoop earring. She is looking down at a silver laptop with a serious, focused expression. Her hands are visible on the keyboard and trackpad. The background is plain white.

**WHAT ARE WE
GOING TO DO
WITH ALL OF THIS
CLOUD DATA?**



My Goal

- Take unmodified POSIX/Win32 applications . . .
- Run those applications in the cloud . . .
- On the same hardware used to run big-data apps . . .
- . . . and give them cloud-scale IO performance!

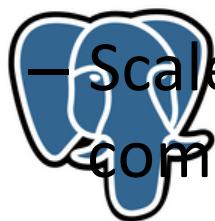
Postgre**SQL**



My Goal

- Take unmodified POSIX/Win32 applications . . .
- Run those applications in the cloud . . .
- On the same hardware used to run big-data apps . . .
- . . . and give them cloud-scale IO performance!

PostgreSQL
– Throughput > 1000 MB/s

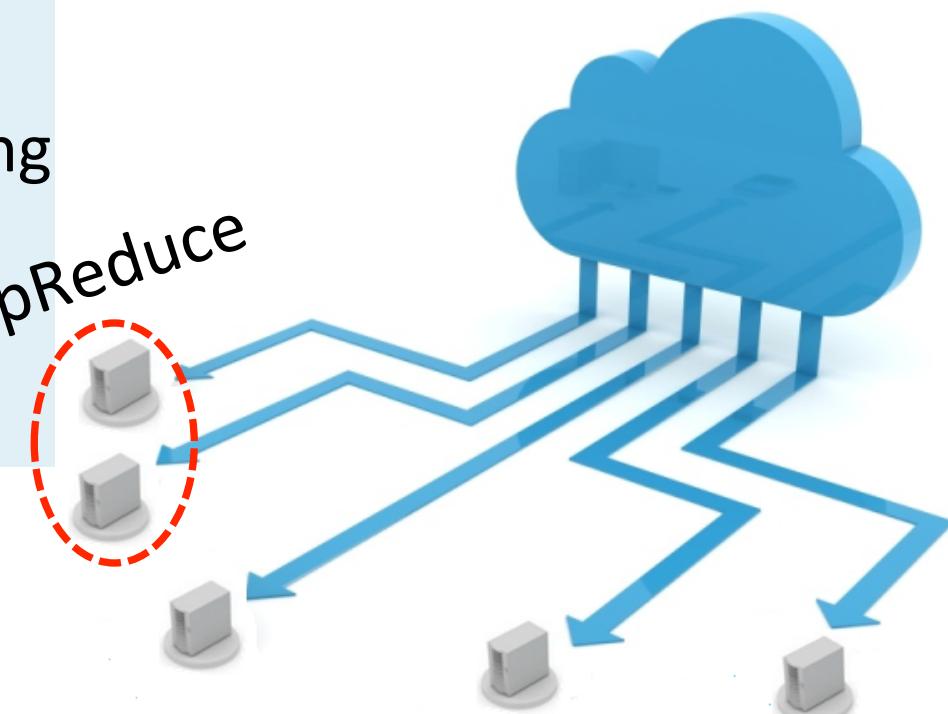


eclipse
– Scale-out architecture using commodity parts

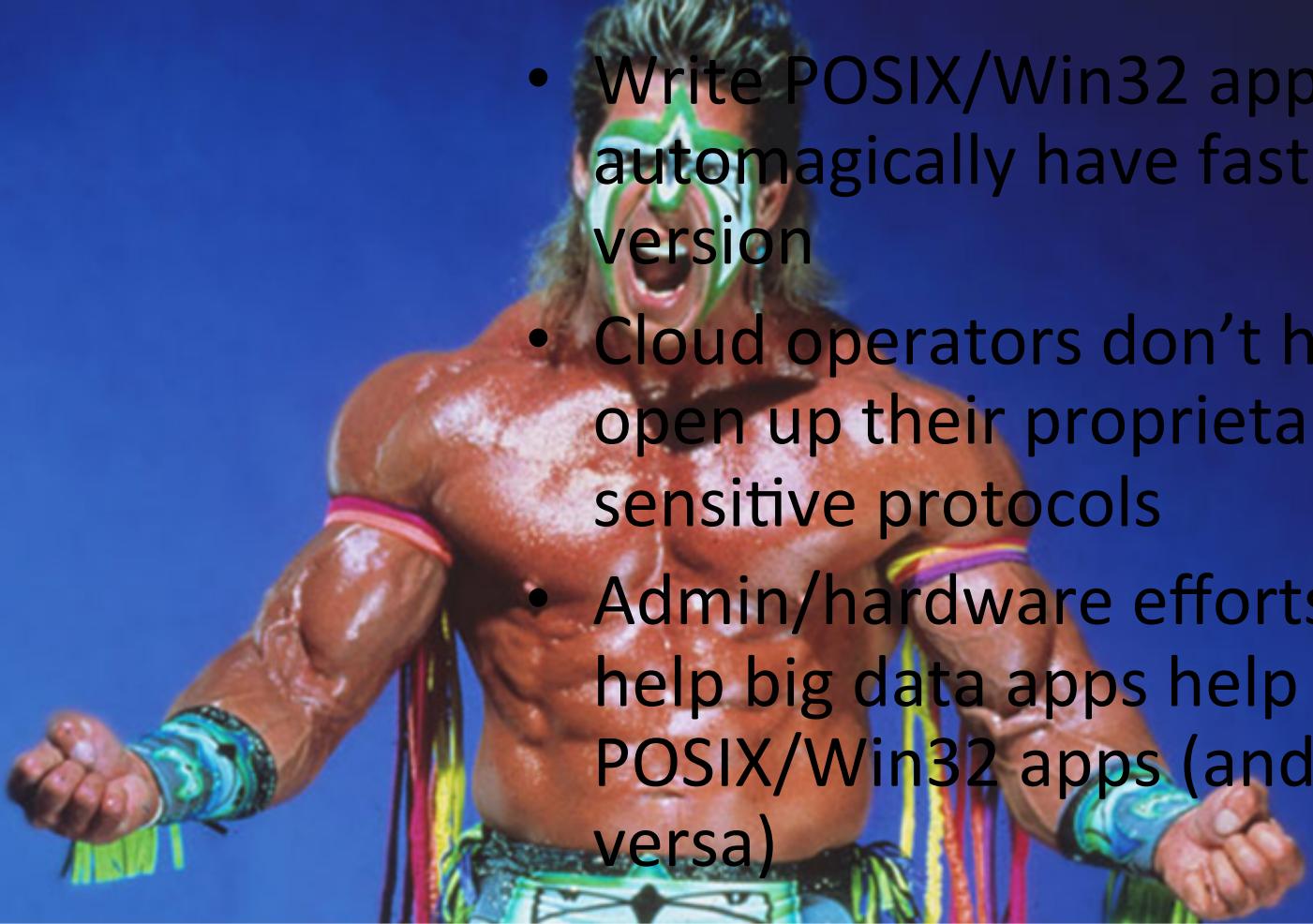
– Transparent failure recovery

Microsoft®
Exchange 2013

MapReduce



Why Do I Want To Do This?

- 
- Write POSIX/Win32 app once, automagically have fast cloud version
 - Cloud operators don't have to open up their proprietary or sensitive protocols
 - Admin/hardware efforts that help big data apps help POSIX/Win32 apps (and vice versa)

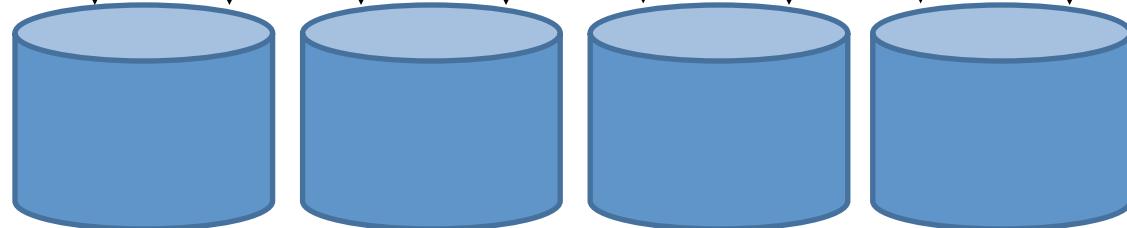
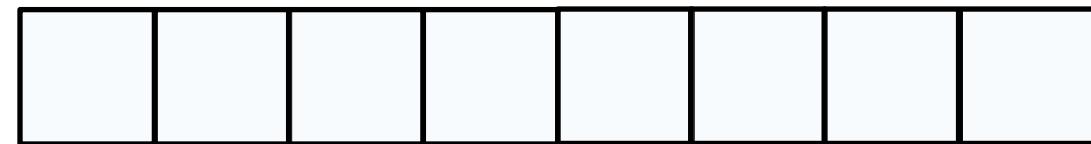
BECAUSE THIS WOULD BE AMAZING

Our Solution: Blizzard

PostgreSQL

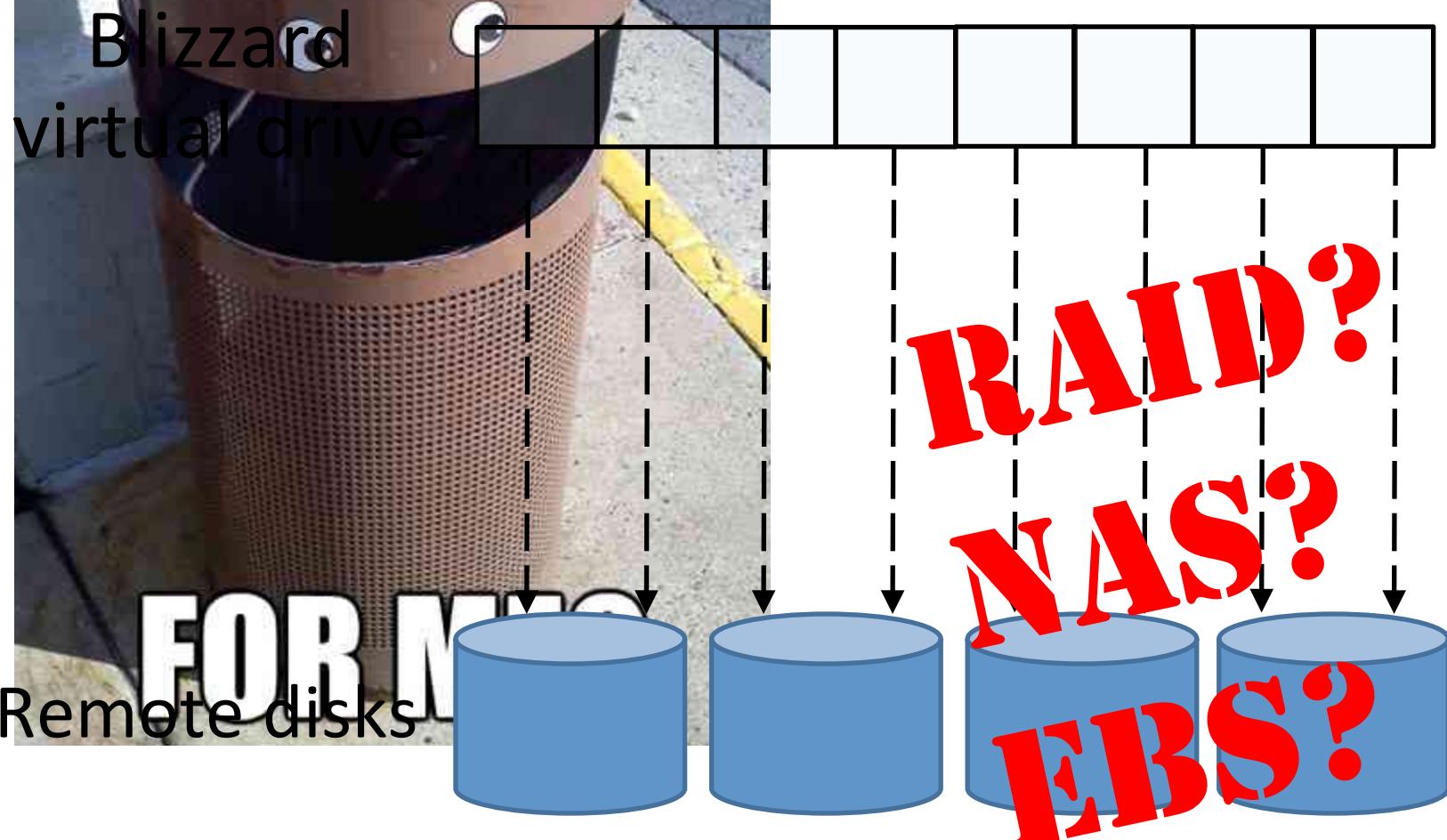


Blizzard
virtual drive



Remote disks

Our Solution: Blizzard



A close-up, profile view of a man's face, focusing on his right side. He has short, light-colored hair and is wearing a blue and white striped shirt. The lighting is dramatic, casting deep shadows on the left side of his face (the viewer's perspective) and highlighting the right side. His gaze is directed towards the camera. The background is solid black.

**THE PROBLEM
IS YOU**

LISTEN



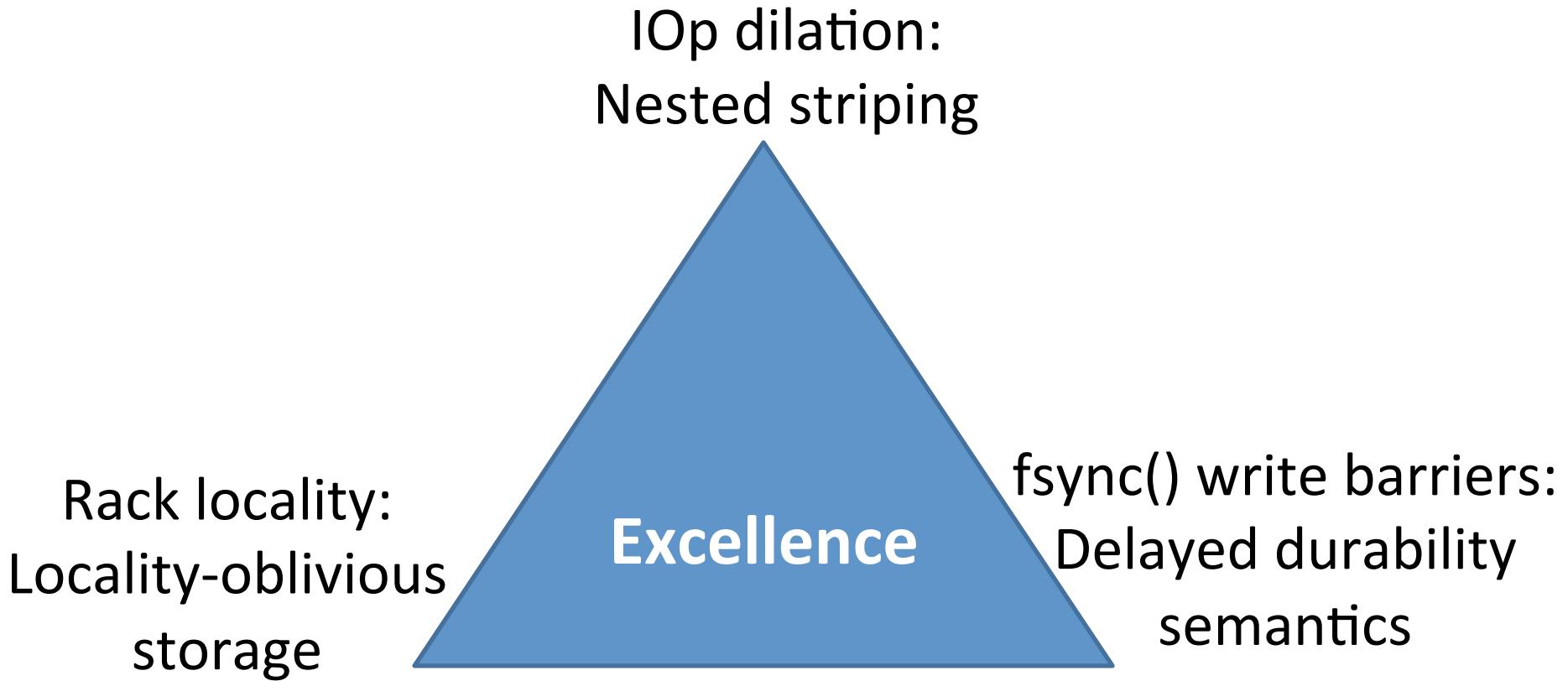
The naïve approach for implementing virtual disks does not **maximize spindle parallelism** for POSIX/Win32 applications which **frequently issue fsync() operations** to maintain consistency.

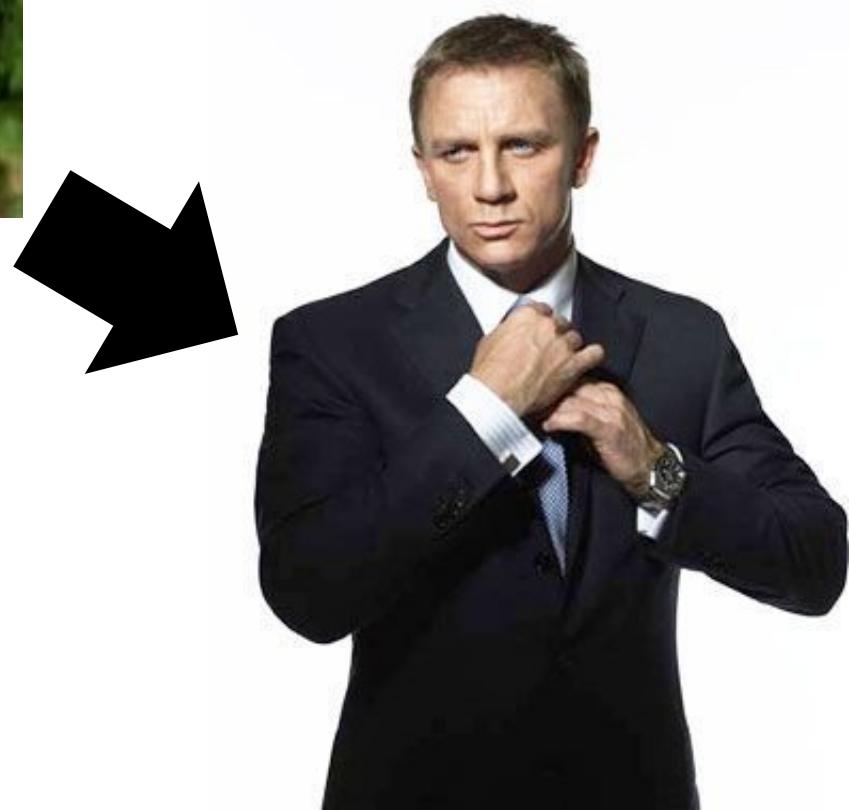
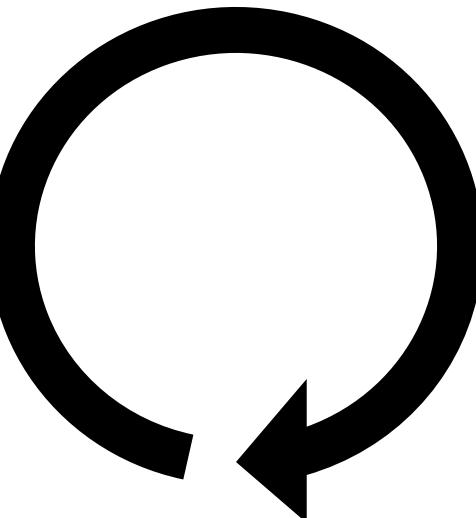
LISTEN



The naïve approach for implementing virtual disks does not **maximize spindle parallelism** for POSIX/Win32 applications which **frequently issue fsync() operations** to maintain consistency.

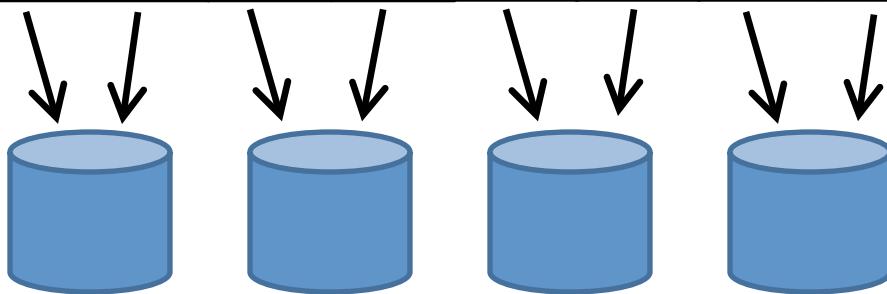
LISTEN



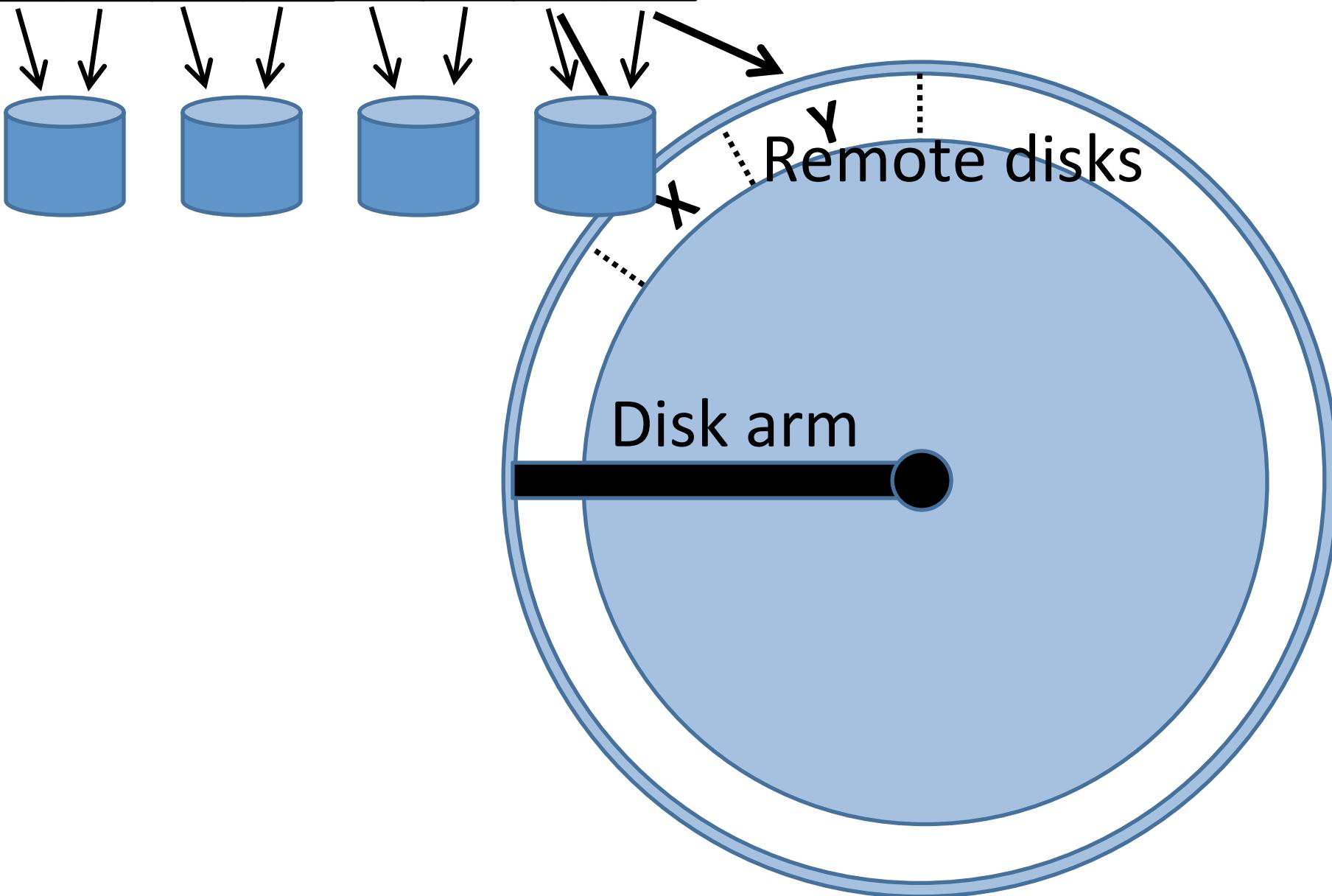


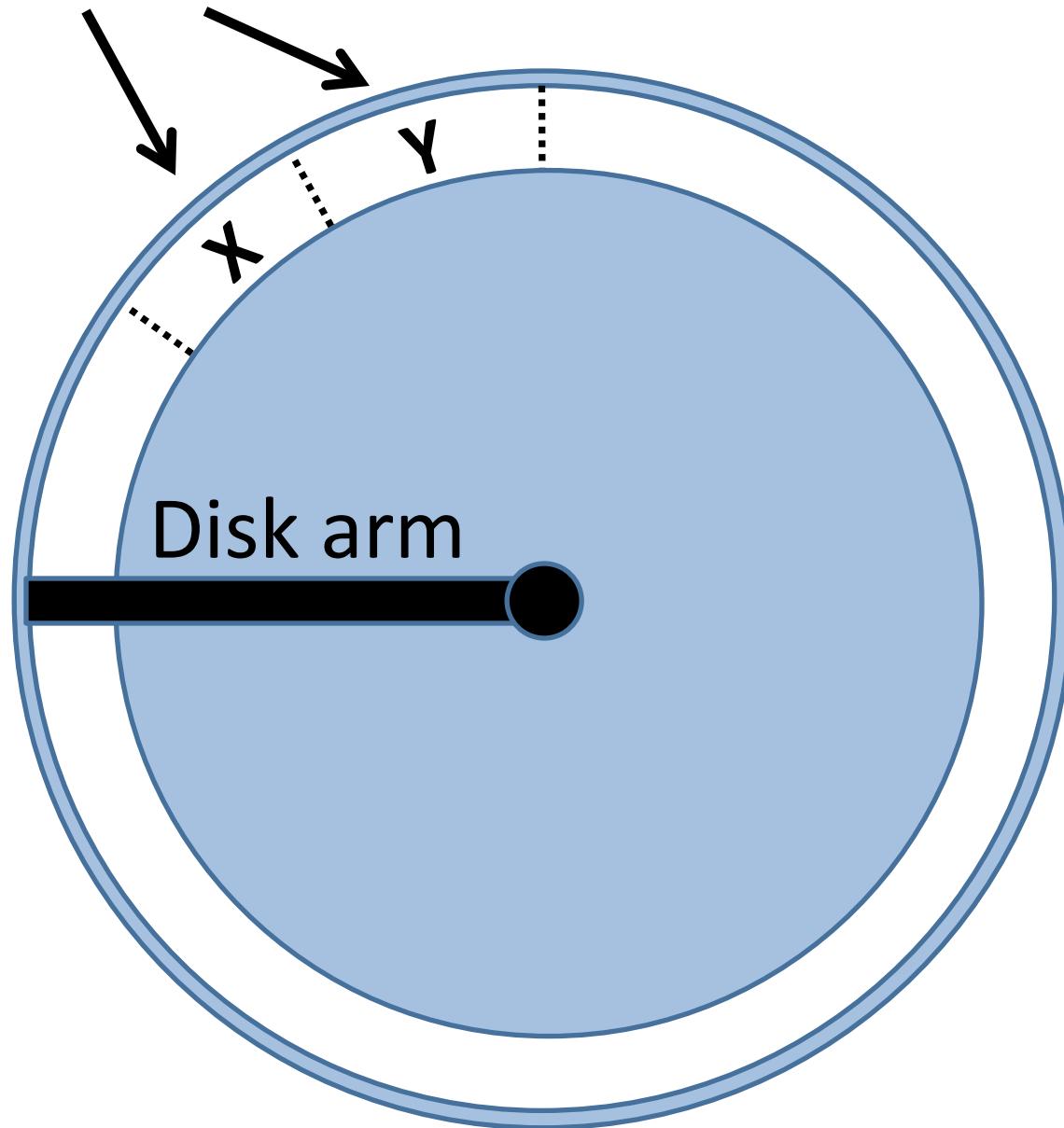
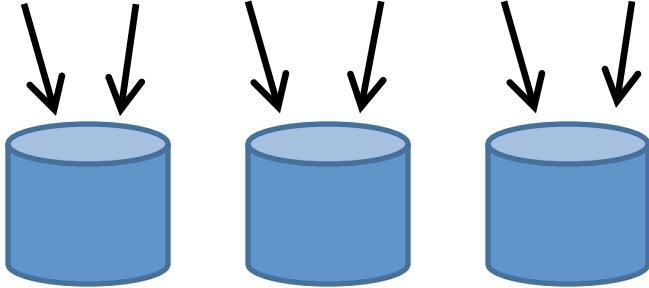


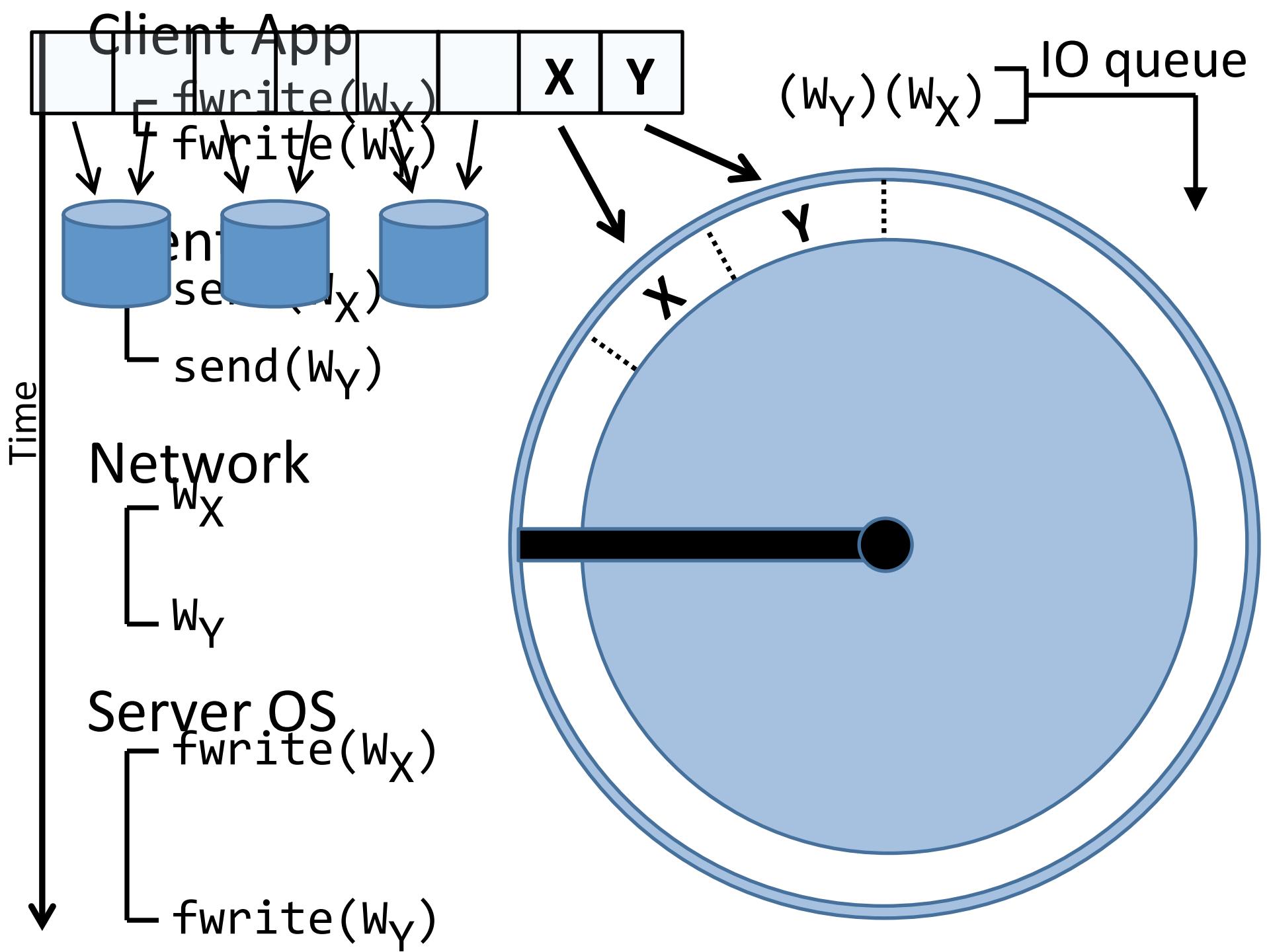
Virtual disk

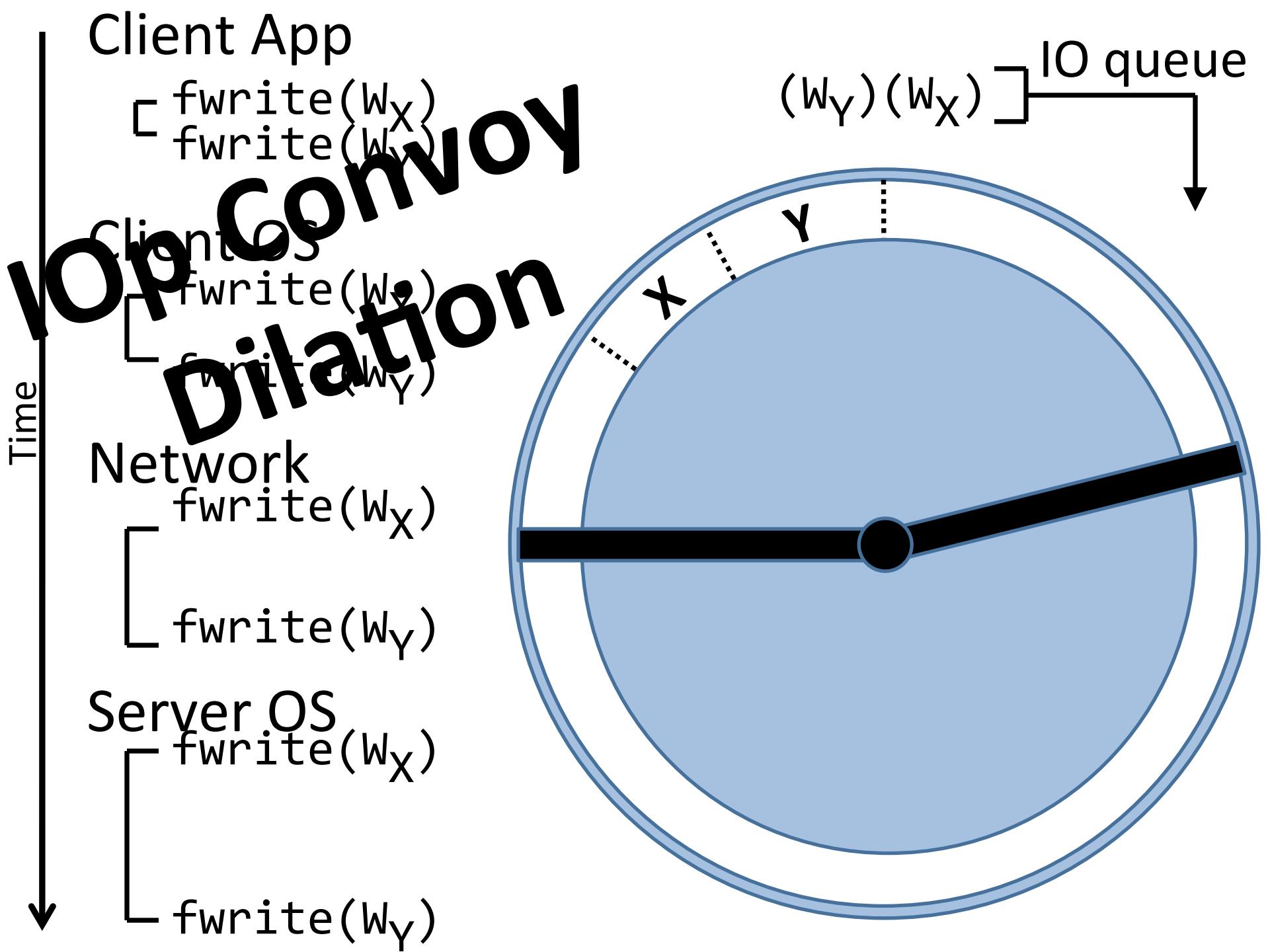


Remote disks





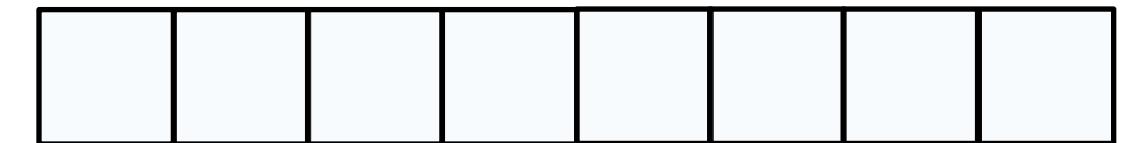




Fixing IOp Convoy Dilation

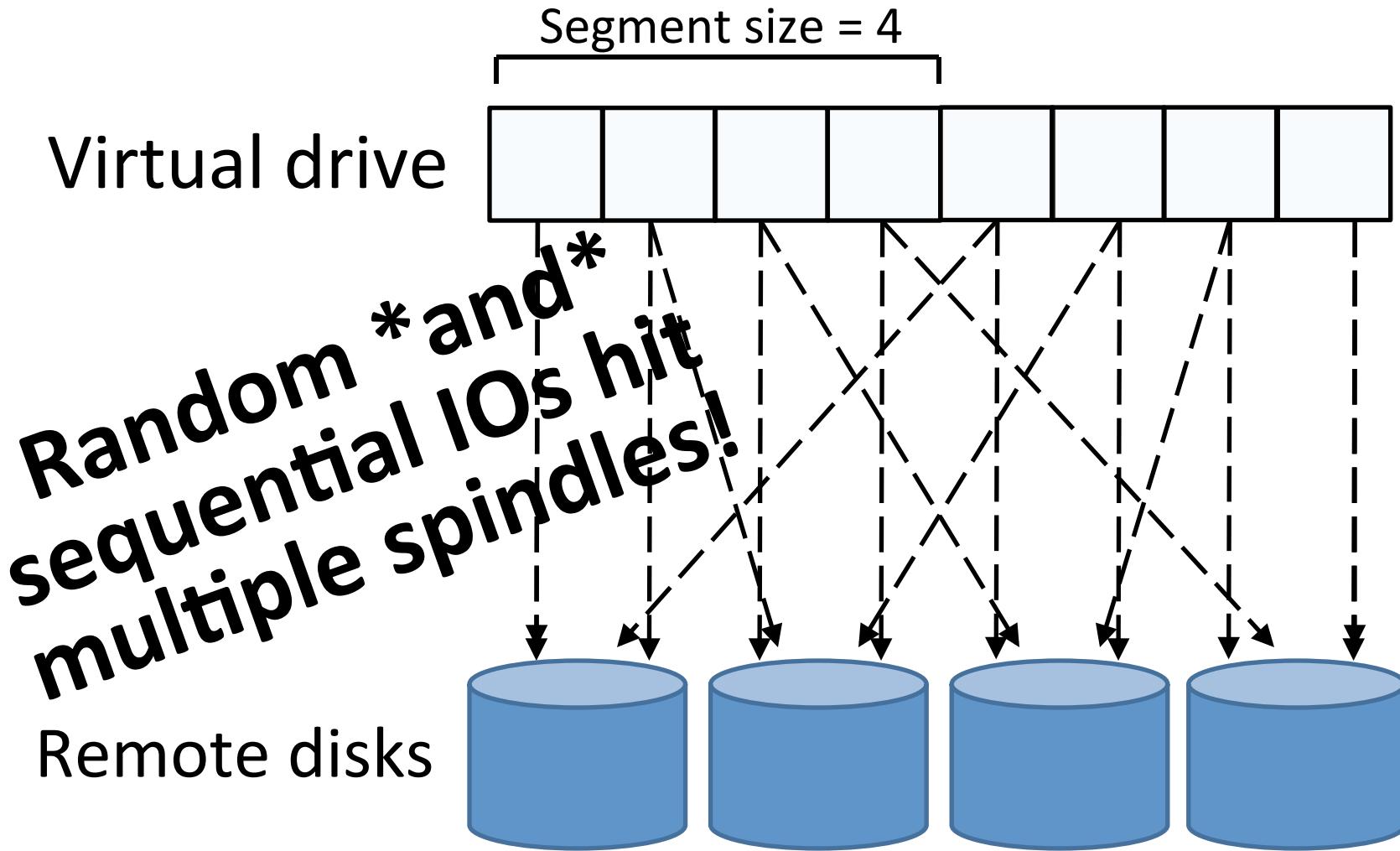
Virtual drive

Segment size = 4



Remote disks

Fixing IOp Convoy Dilation



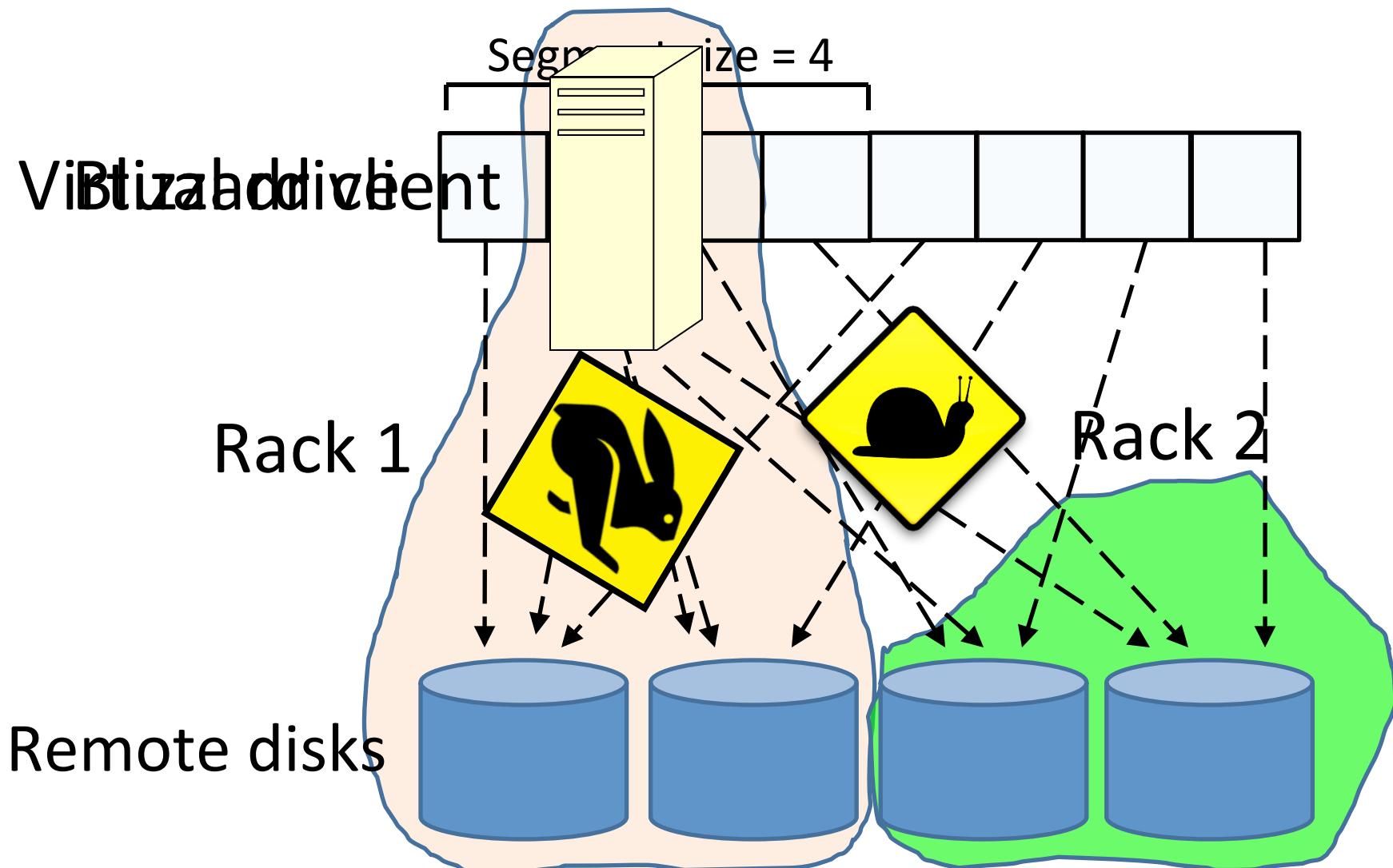
A close-up profile of Muhammad Ali's face, looking upwards and to the right. He has dark skin, a prominent nose, and a serious expression. His hair is dark and curly. The background is blurred.

**TOTAL
VICTORY**

Rack Locality



Rack Locality



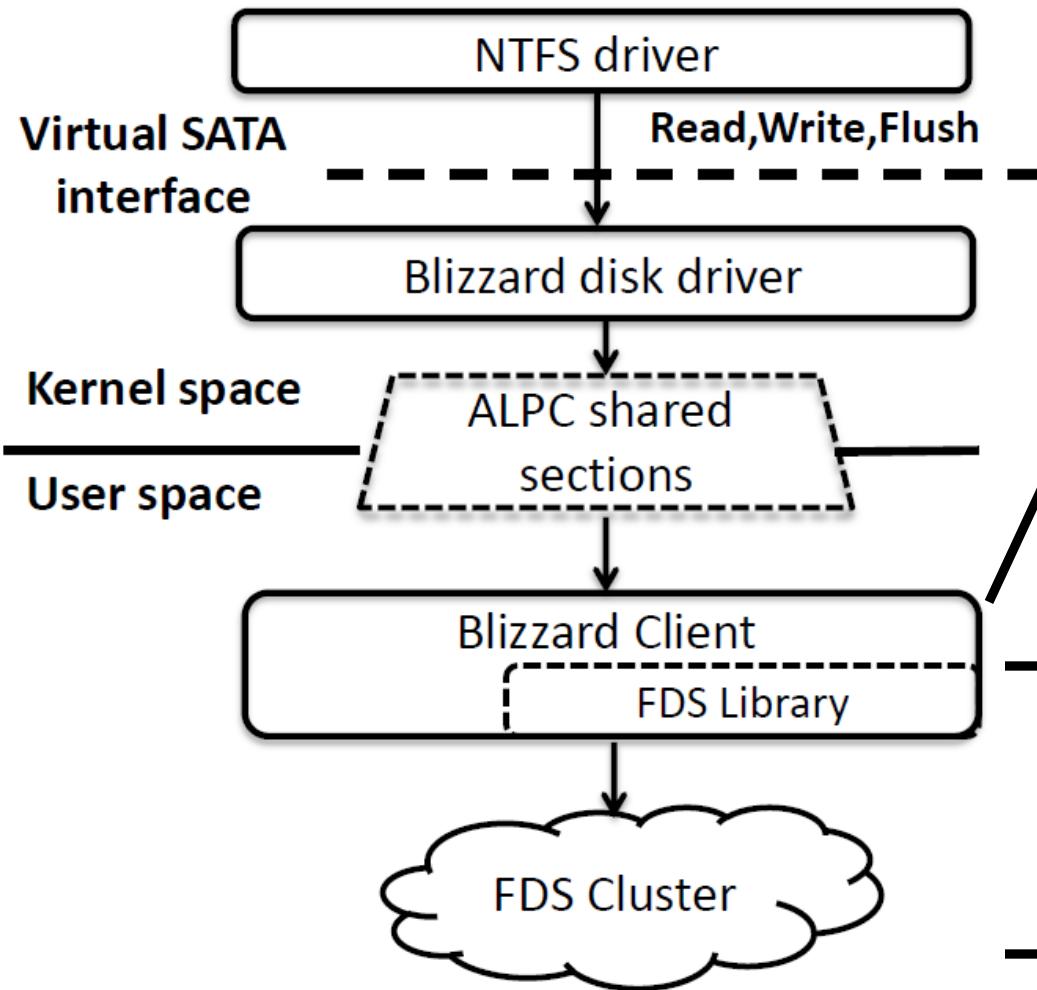


**AIN'T NOBODY
GOT TIME FOR
THAT**

FDS To The Rescue (OSDI 2012)

- Hardware architecture
 - Full bisection bandwidth network (no oversubscription)
 - Allocate each disk enough network bandwidth to drive disk at full sequential speed (1 disk \approx 128 MB/s \approx 1Gbps)
- Result: locality-oblivious storage
 - Any client can access any disk as fast as local
 - **Enables aggressive striping**

Blizzard as FDS Client



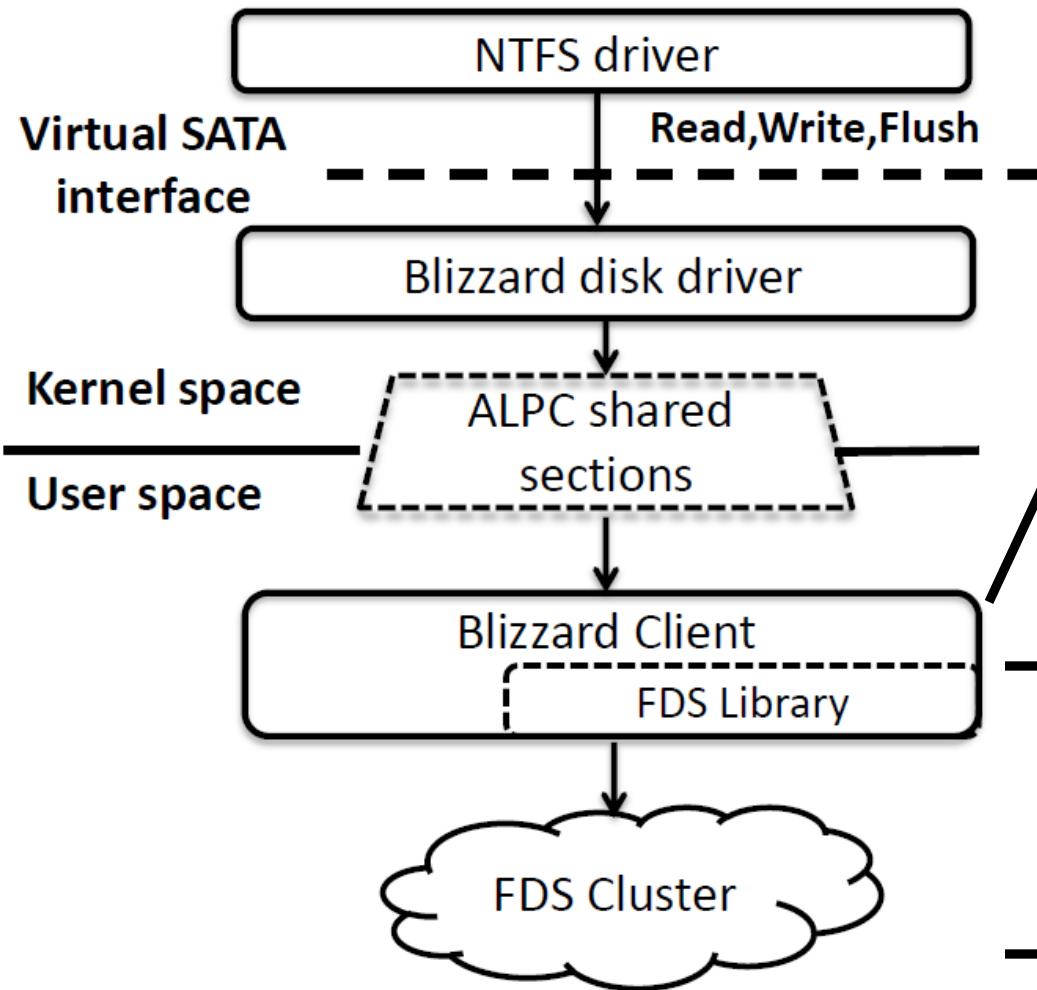
Blizzard client handles:

- Nested striping
- Delayed durability semantics

FDS provides:

- Locality-oblivious storage hardware
- Server-side failure recovery
- RTS/CTS to avoid edge congestion

Blizzard as FDS Client



Blizzard client handles:

- Nested striping
- Delayed durability semantics

FDS provides:

- Locality-oblivious storage hardware
- Server-side failure recovery
- RTS/CTS to avoid edge congestion



**POSIX/Win32
apps**

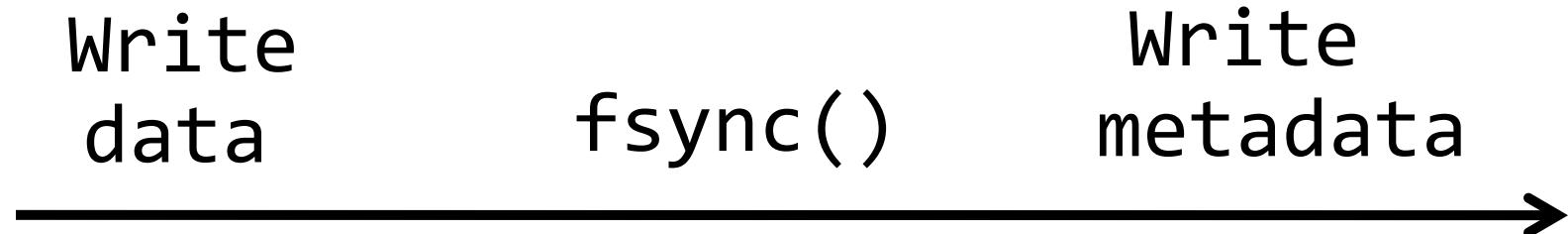
**Traditional
big-data apps**



ARE WE THERE YET?

The problem with fsync()

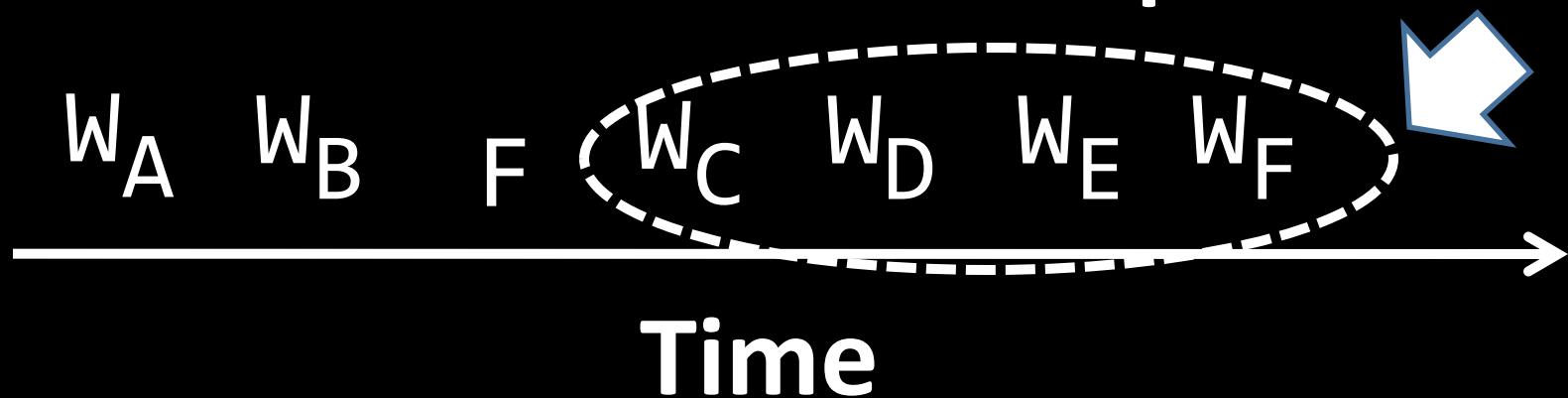
- Used by POSIX/Win32 file systems and applications to implement crash consistency
 - Disk only returns from fsync() when all prior writes have become durable
 - Ex: ensure data is written before metadata





WRITE BARRIERS
RUIN
BIRTHDAYS

Stalled operations
limit parallelism!

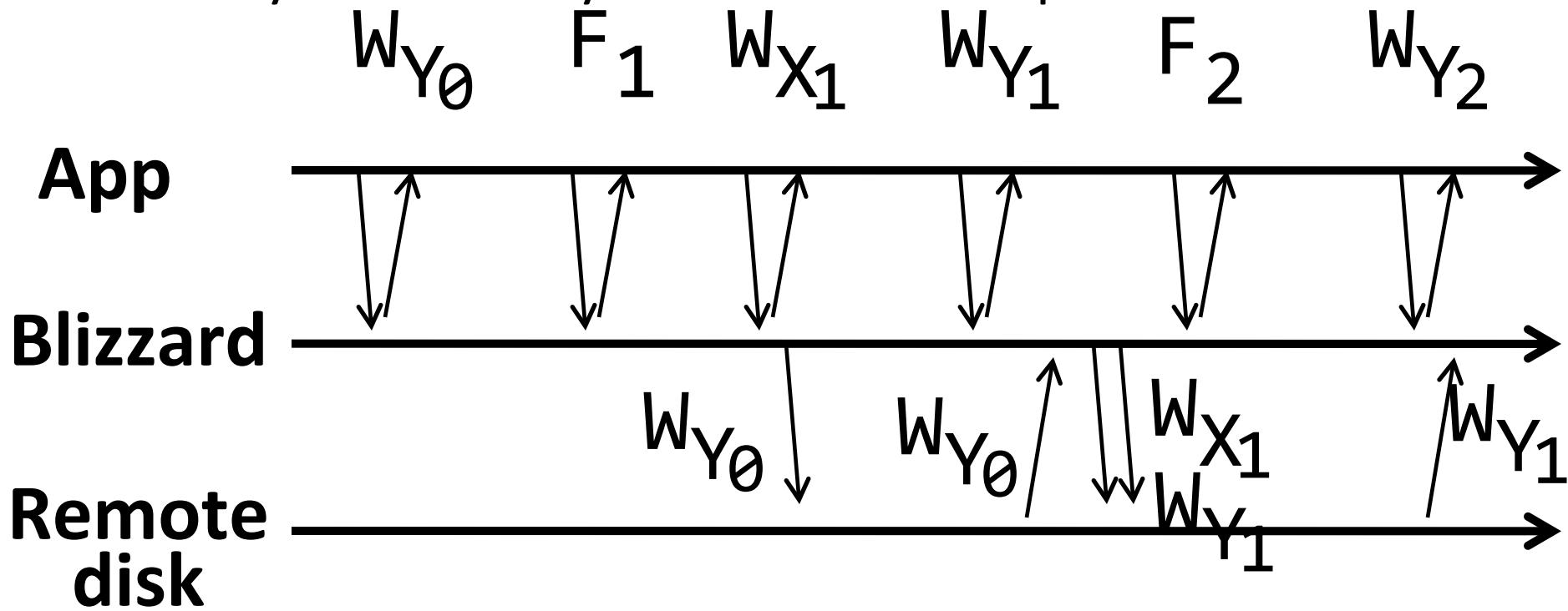


Delayed Durability

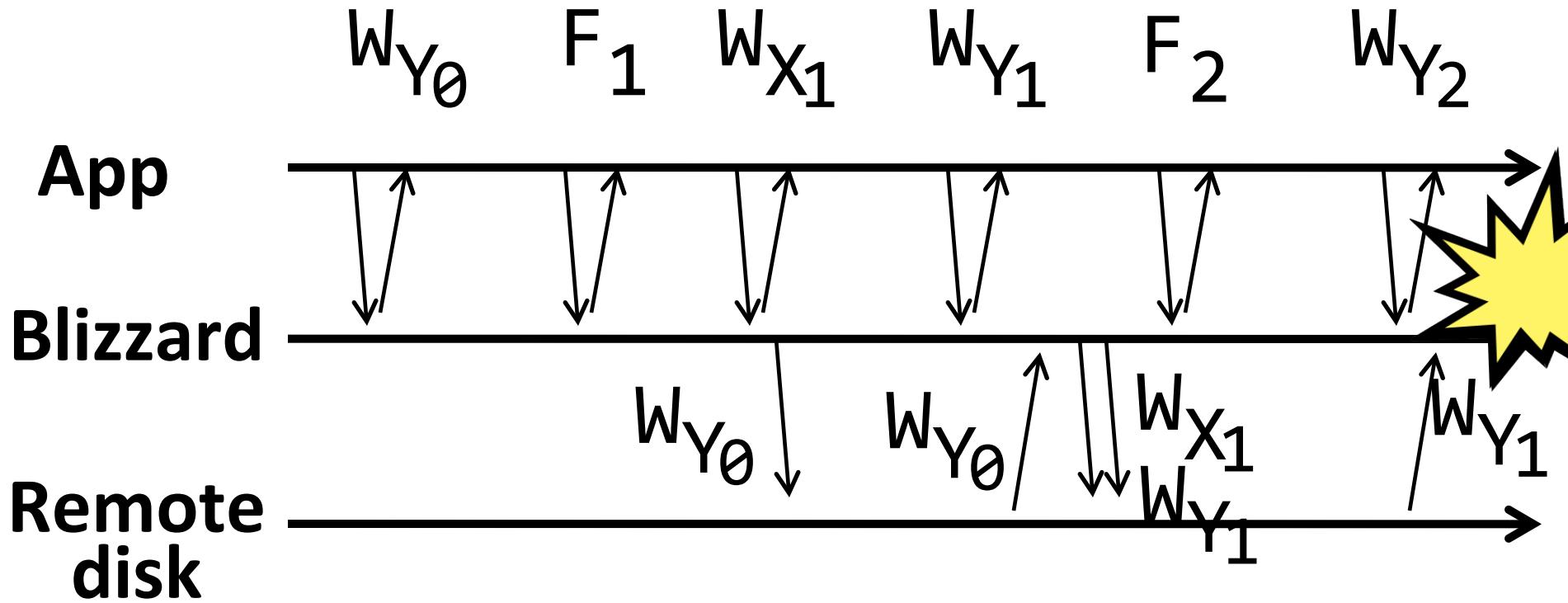
- Decouple durability from ordering
- Acknowledge `fsync()` immediately . . .
 - . . . but increment flush epoch
 - Tag writes with their epoch number,
asynchronously retire writes in epoch order

Delayed Durability

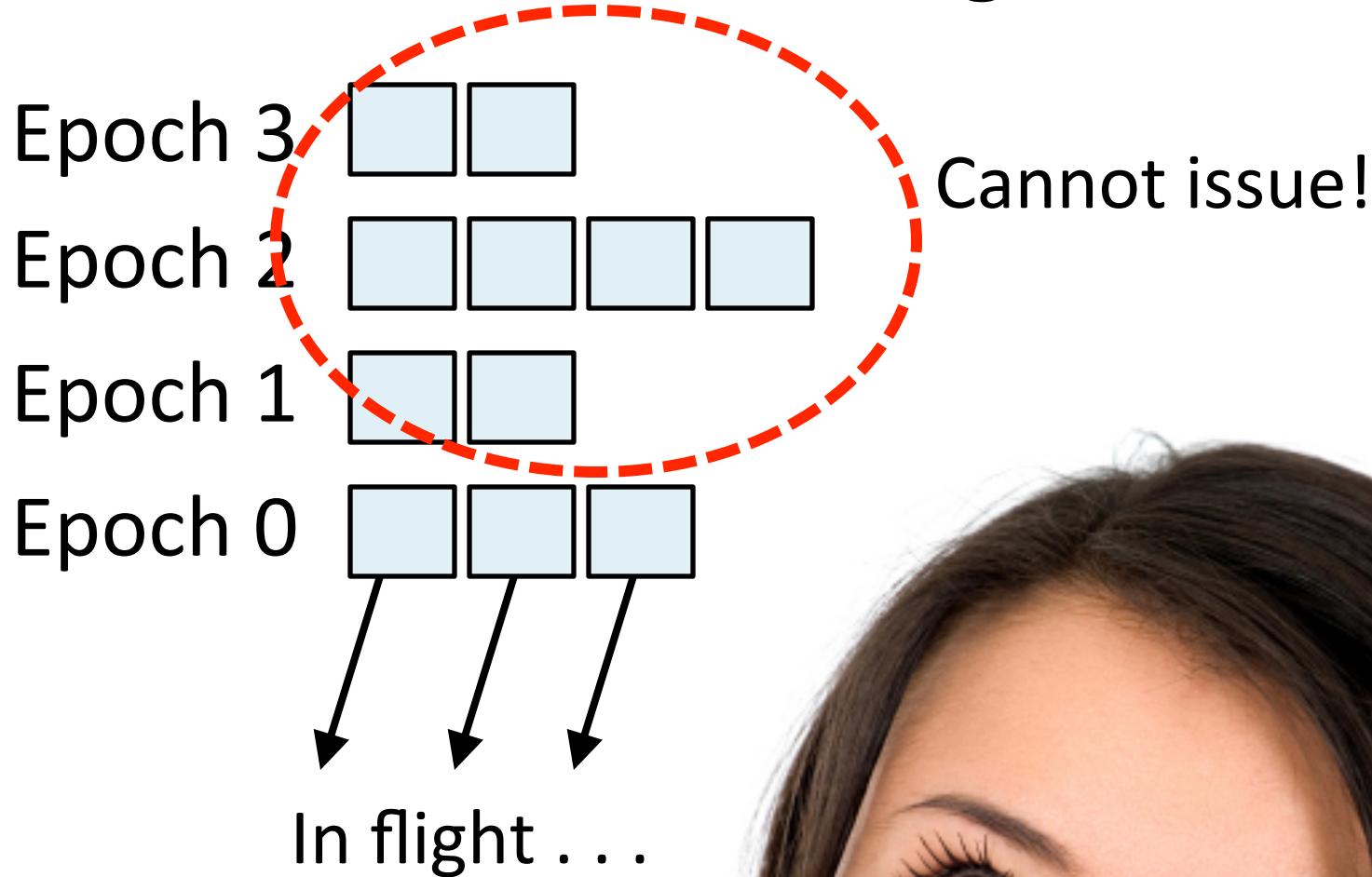
- Decouple durability from ordering
- Acknowledge `fsync()` immediately . . .
 - . . . but increment flush epoch
 - Tag writes with their epoch number,
asynchronously retire writes in epoch order



- All writes are acknowledged . . .
- . . . but only ~~Y₀~~ and ~~Y₁~~ are durable!
- ~~Decouples prefix consistency from ordering~~
 - = All epochs up to N-1 are durable
 - = Acknowledge + `Sync()` immediately . . .
 - = Some, all, or no writes from epoch N are durable
 - = but tag writes with their epoch number
 - = No writes from later epochs are durable
- Prefix consistency good enough for most apps, provides much better performance!



Isn't Blizzard buffering a lot of data?

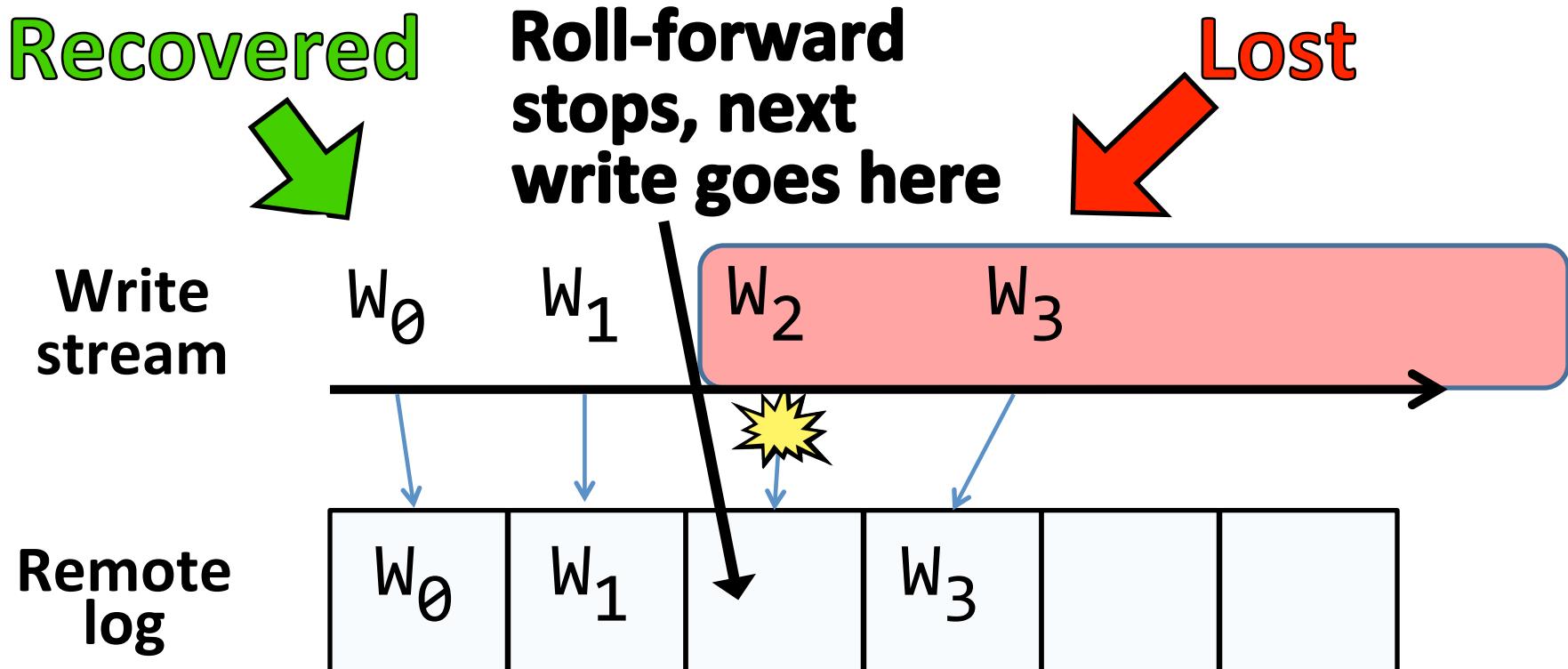


THESE ARE NO
SEED BASED
EXTRACTS
NECESSARY
FOR
EXTRACTION?
TECHNIQUES

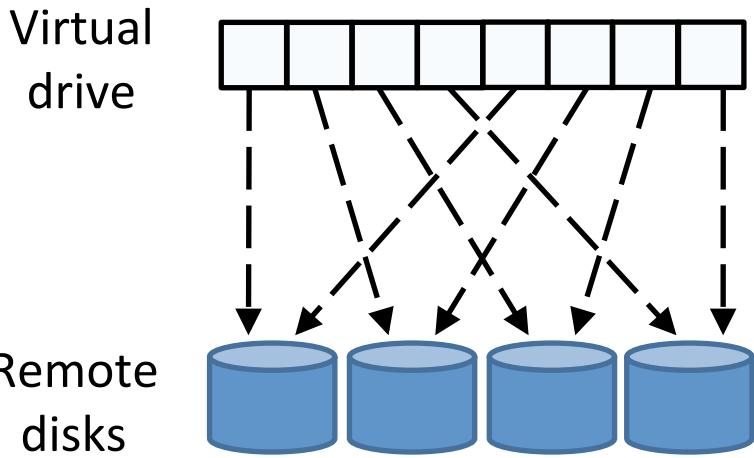


Log-based Writes

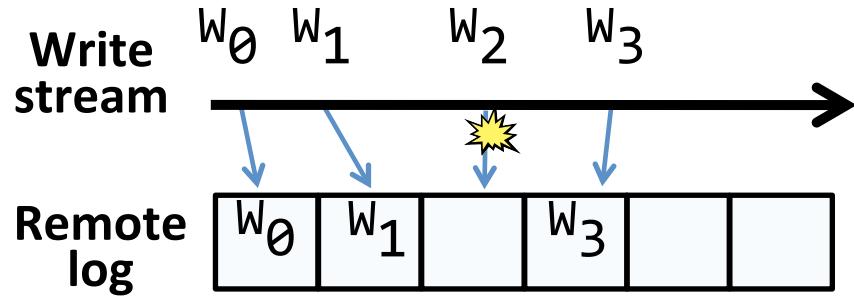
- Treat backing storage as a distributed log
 - Issue writes to log **immediately** and **in order**
 - On failure, roll forward from last checkpoint and stop when you find torn write, unallocated log block with old epoch number



Summary of Blizzard's Design



FDS



- Problem: IOp Dilation
- Solution: Nested striping

- Problem: Rack locality constrains parallelism
- Solution: Full-bisection networks, match disk and network bandwidth

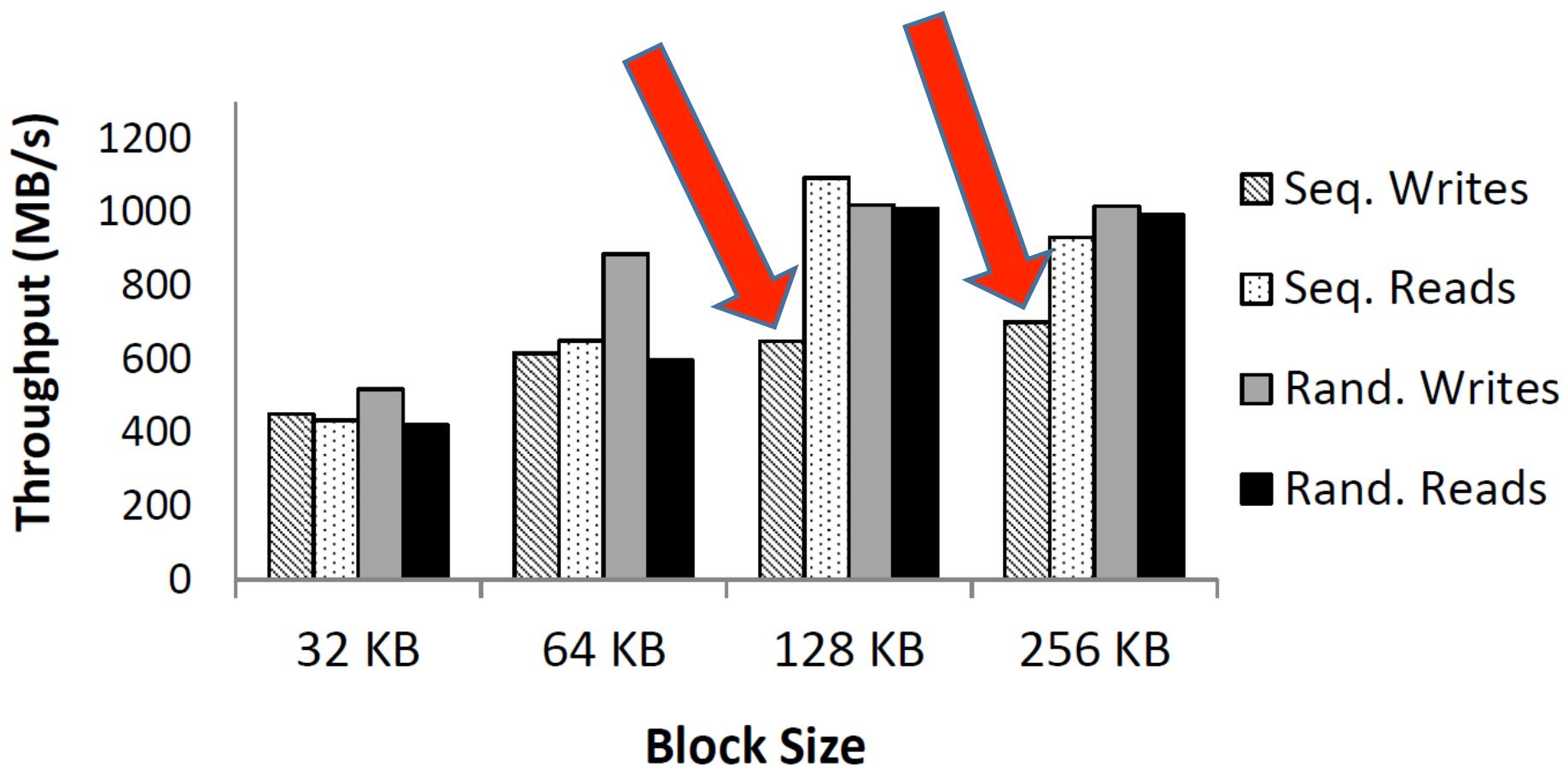
- Problem: Evil fsync()s
- Solution: Delayed durability

PROOF OF YOUR EXCELLENCE

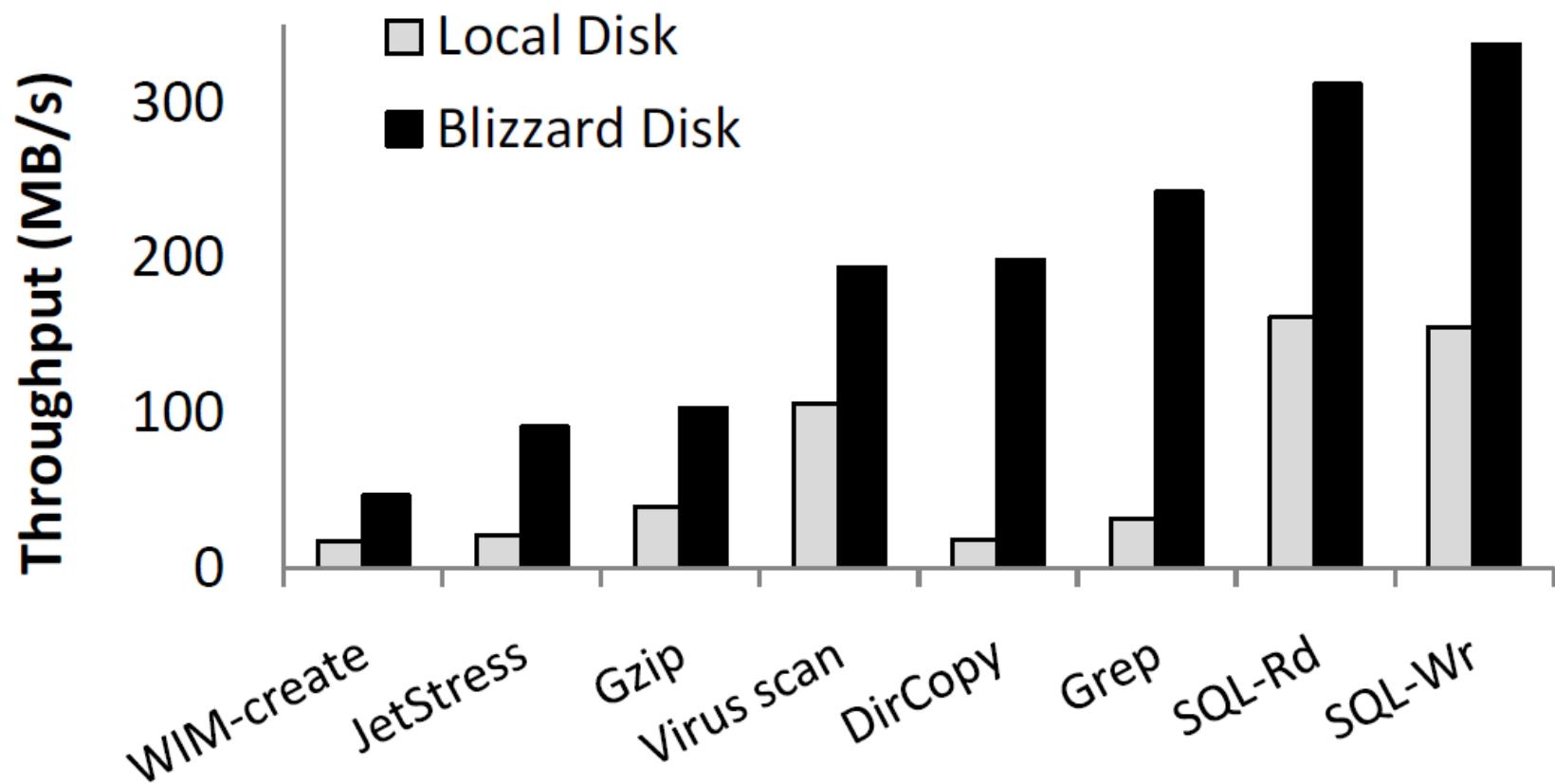


I DO NOT SEE IT

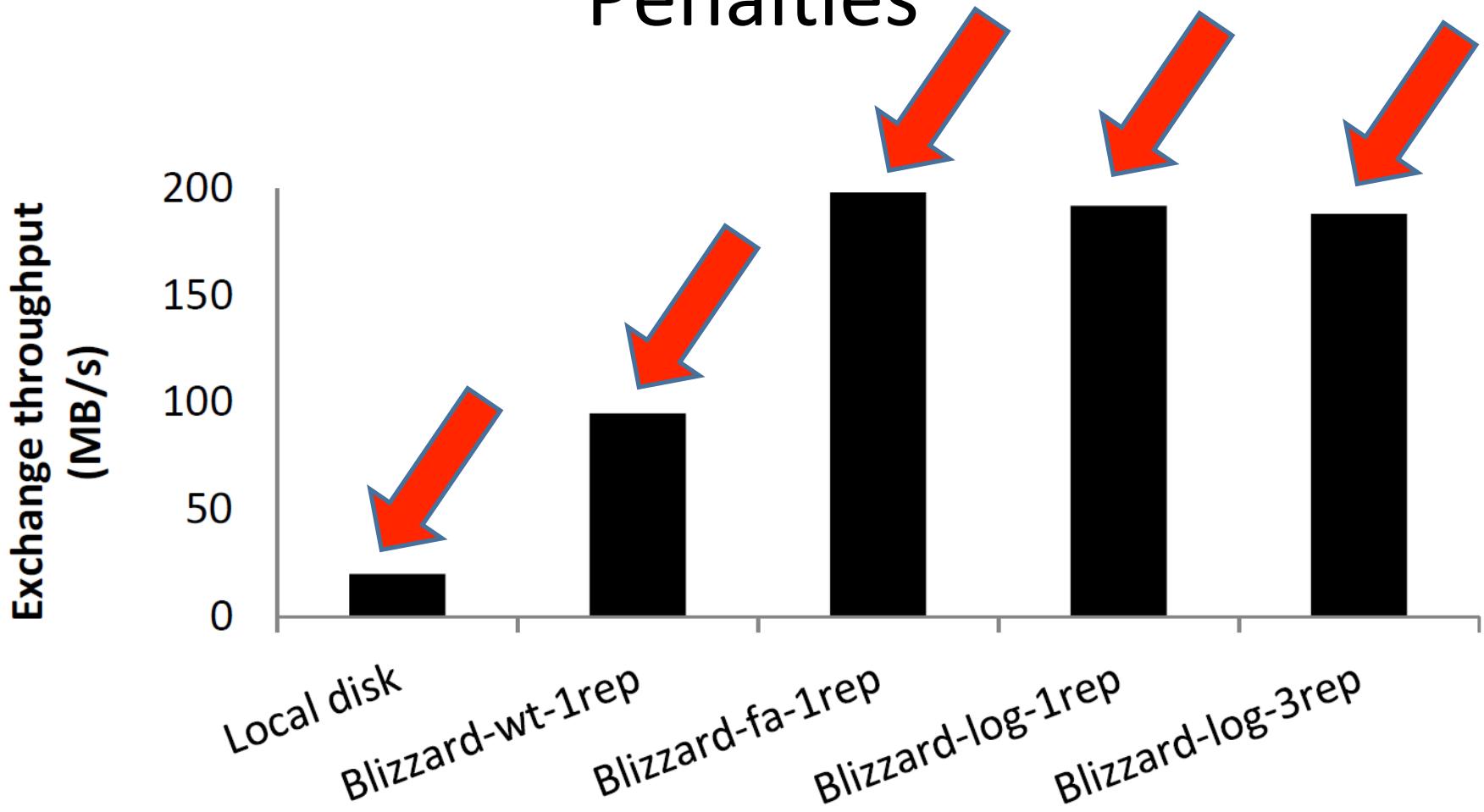
Throughput Microbenchmark



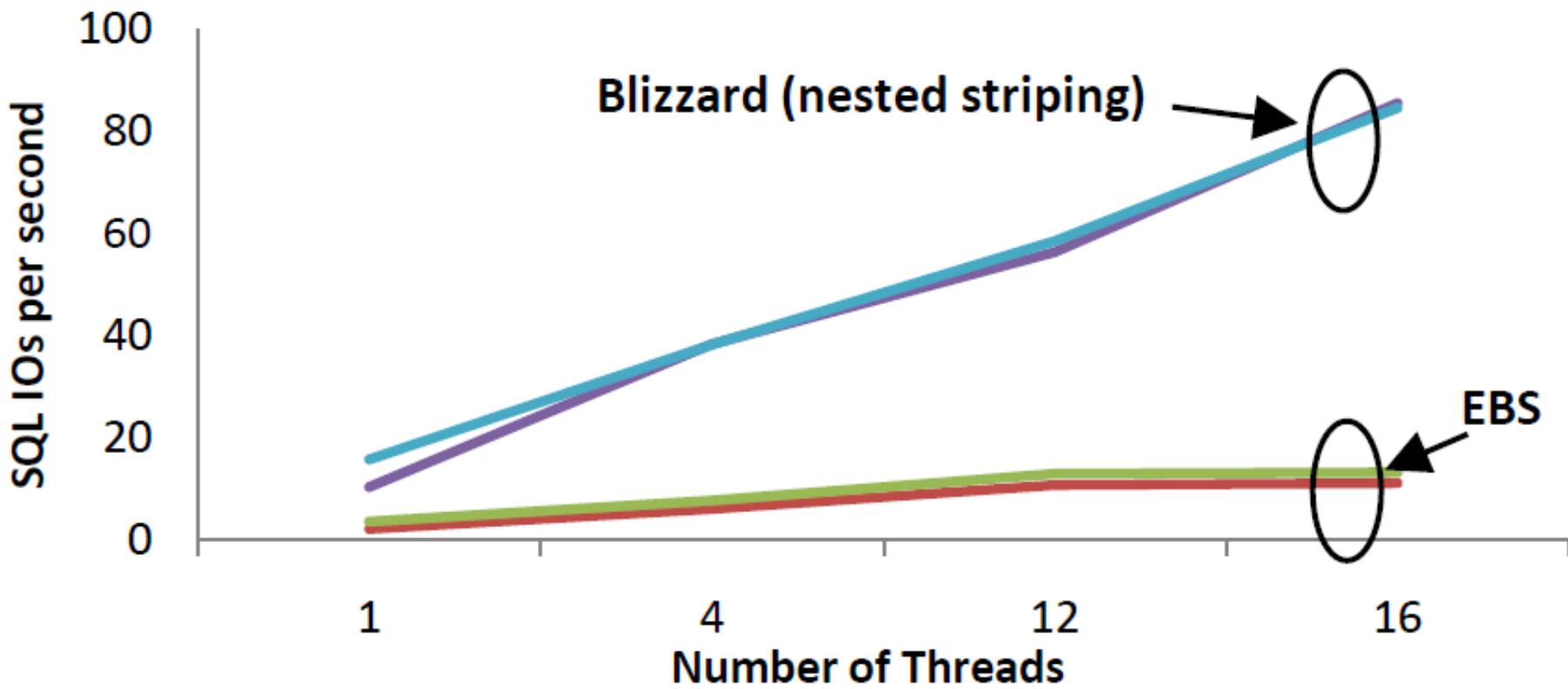
Application Macrobenchmarks



Delayed Durability: Hiding Replication Penalties



Blizzard vs EBS: Write IOps and Read IOps



Related Work

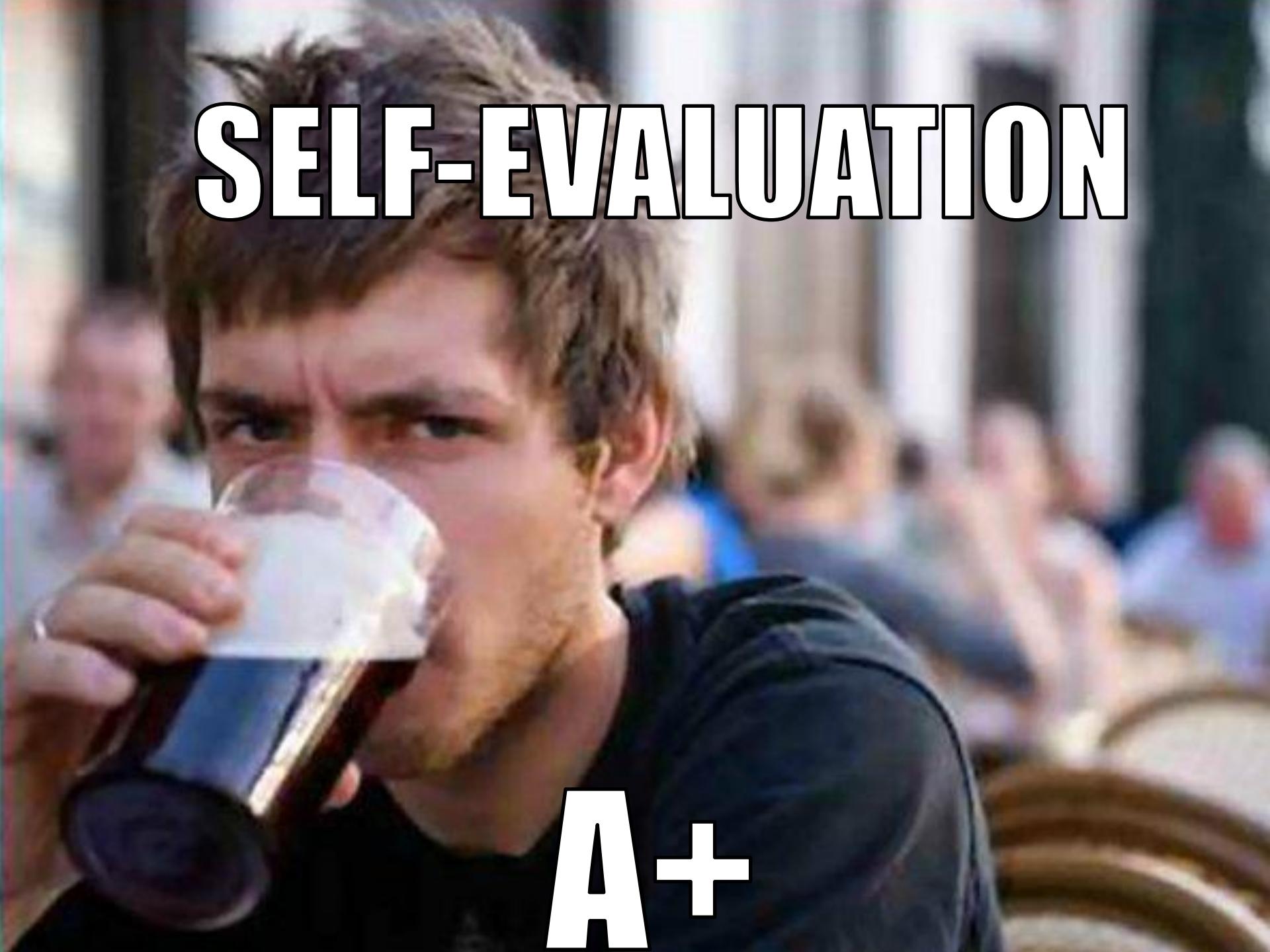
Virtual drives

- iSCSI, AoE
- Petal
- EBS, Azure Drive
- Salus

File systems

- BlueSky, pNFS
- OptFS
- BPFS

SELF-EVALUATION

A close-up photograph of a man with short brown hair, wearing a dark t-shirt. He is holding a large glass mug filled with beer and is in the middle of taking a drink. His gaze is directed straight at the camera with a serious, focused expression. The background is blurred, showing what appears to be a social gathering or a bar.

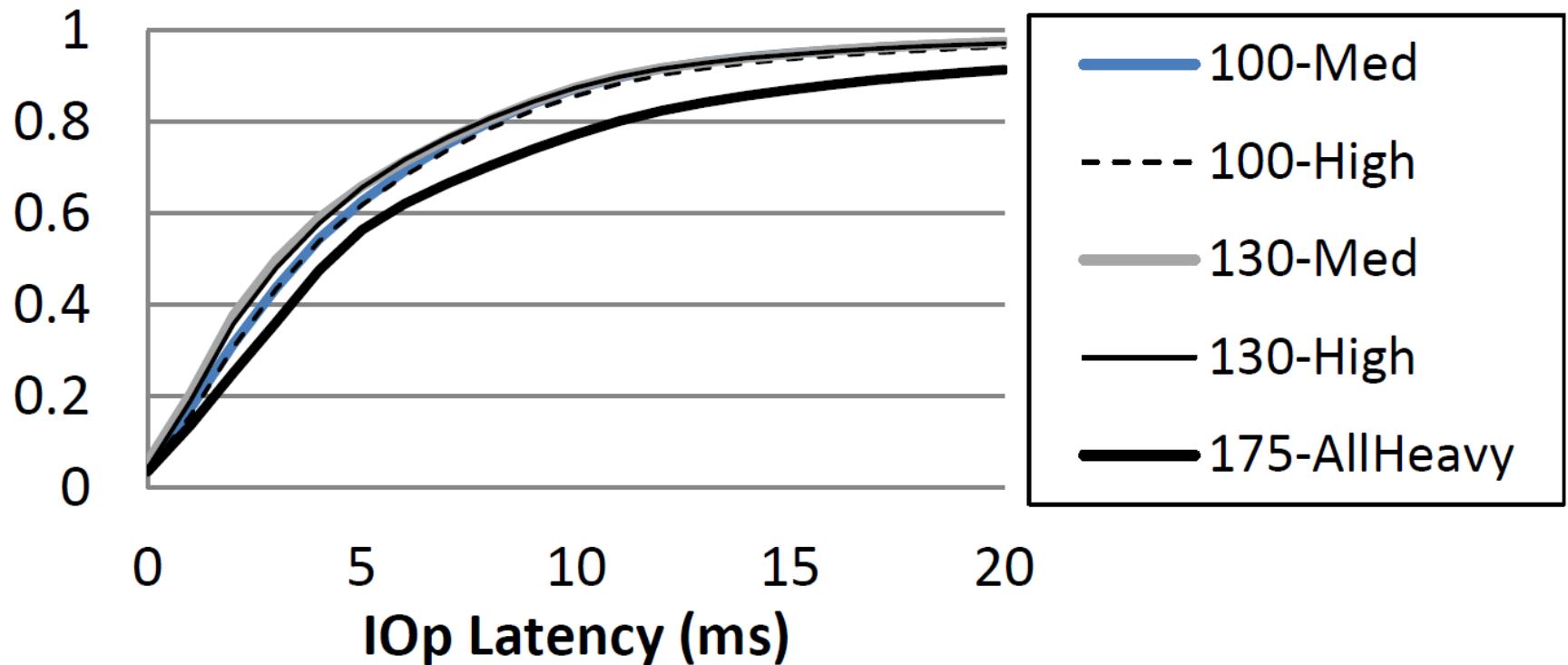
A+

Conclusions

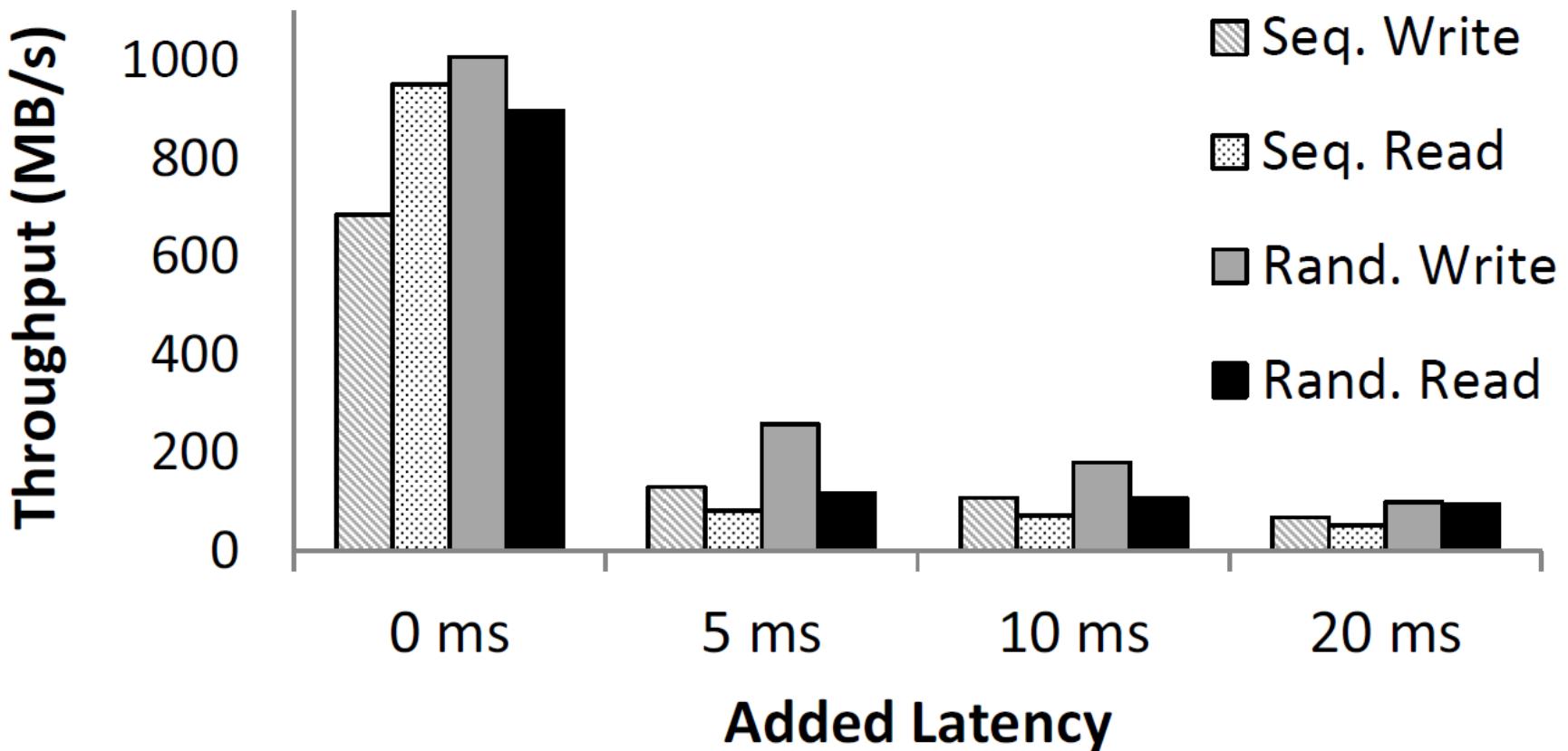
- Unmodified POSIX/Win32 apps can have cloud-scale IO!
 - Nested striping
 - FDS-style hardware substrate
 - Delayed durability semantics
- Raw perf: 1000+ MB/s
- 2x–10x app-level speedups



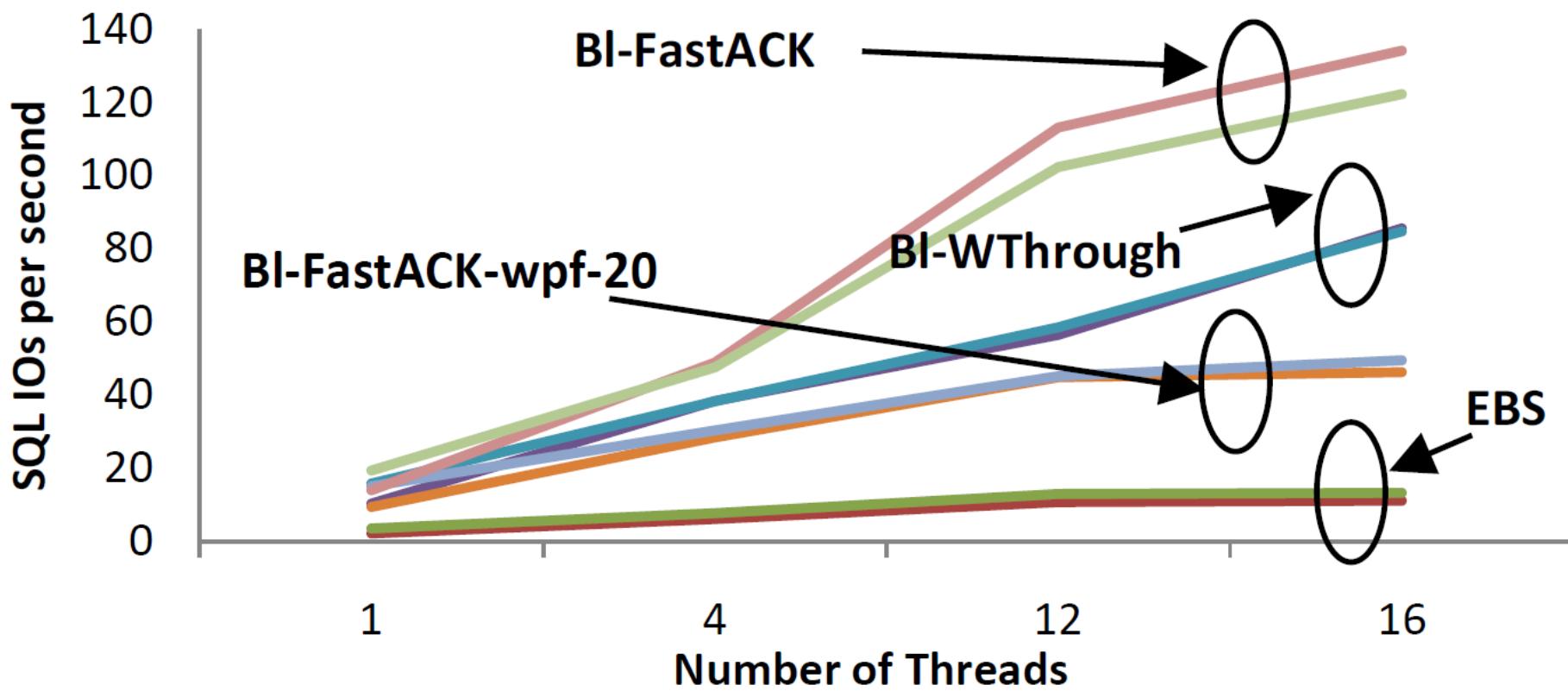
IOp Latency



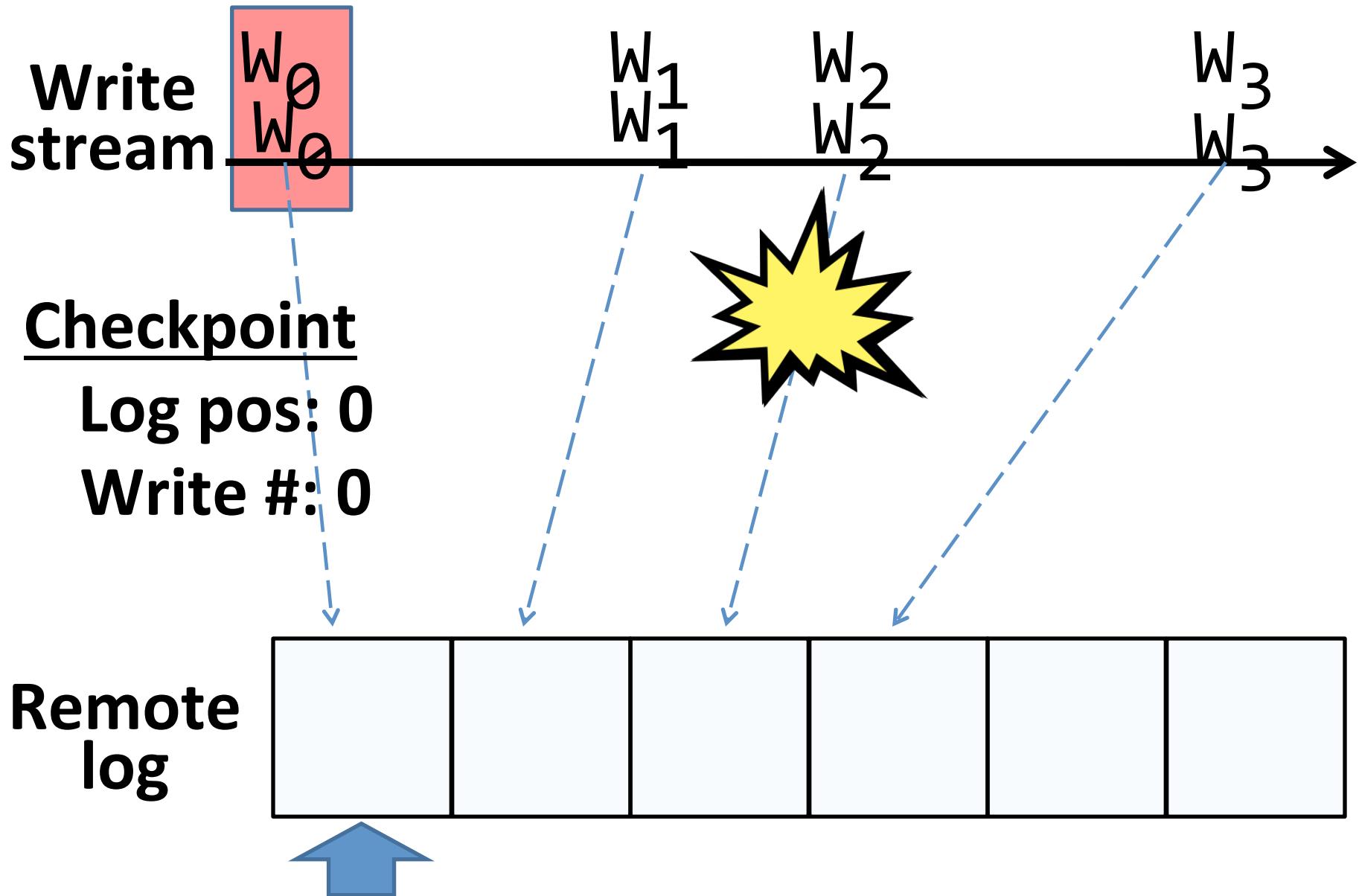
RTT Sensitivity



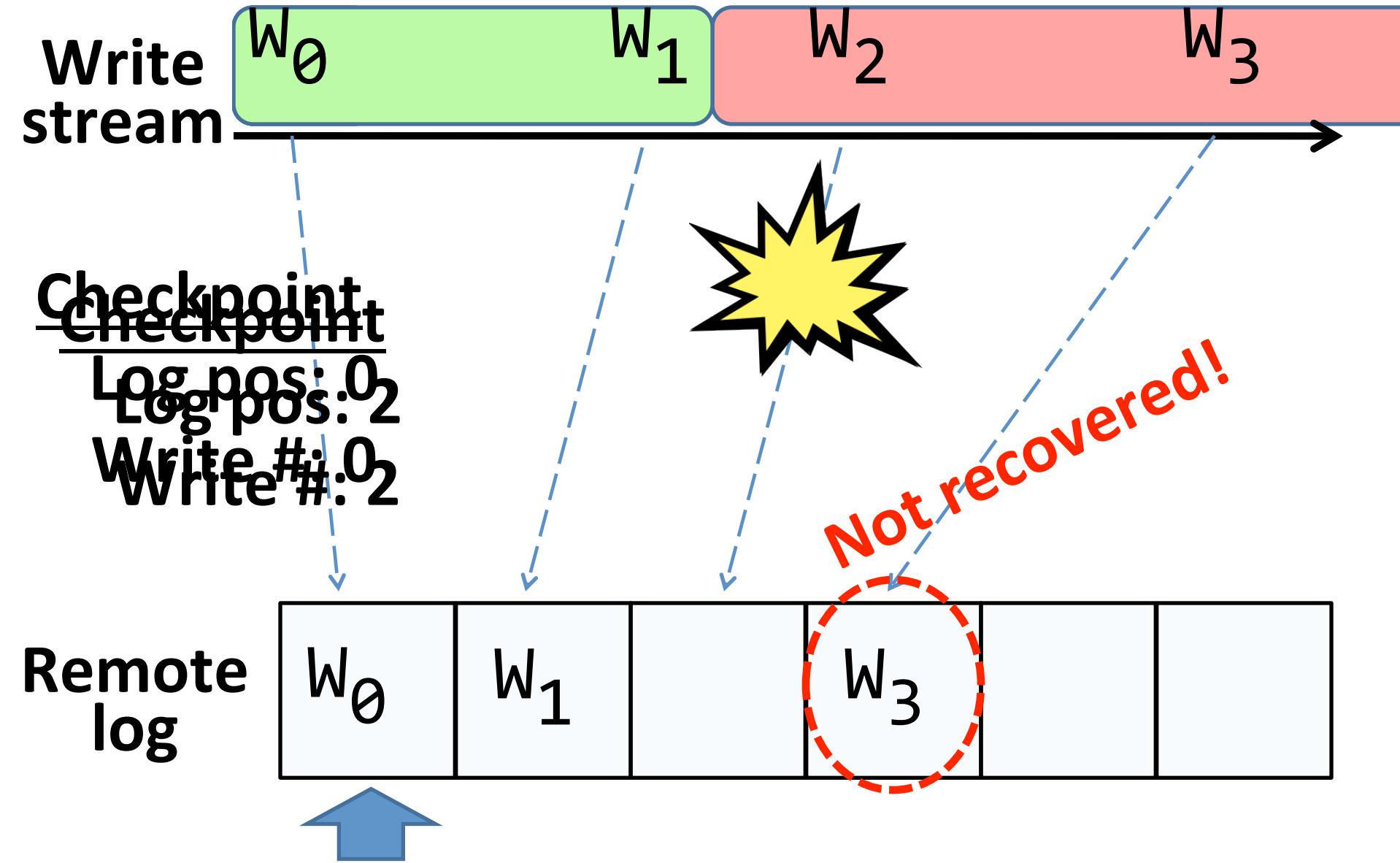
Blizzard vs. EBS



Log-Based Writes



Recovery



Additional Details

- Blizzard maps each virtual block to backing physical block in the log
 - Allocation map included in checkpoints
- To avoid IOp dilation, use random permutation to determine next log position!

Prior example: 0, 1, 2, 3, 4, ...

What Blizzard

really does: $X_{n+1} = (aX_n + c) \text{ mod } m$

(Checkpoint a, c, m, X_n)