# Facebook's disaggregated storage and compute for Map/Reduce

## Yun Jin

# Intro

# Facebook's Monthly Active Users

Grew by ~1.1B since 2010 monthly active users
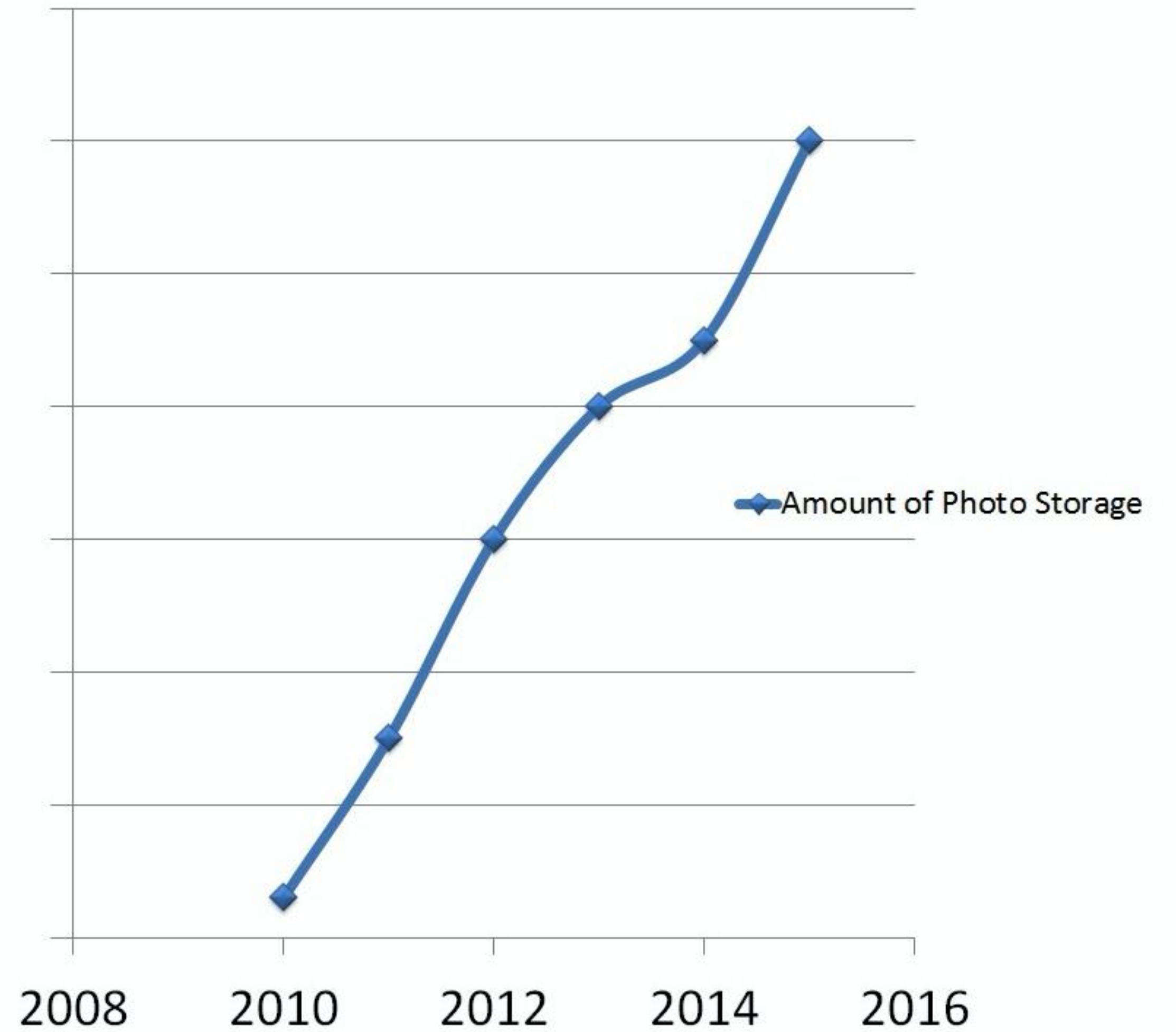


360 video
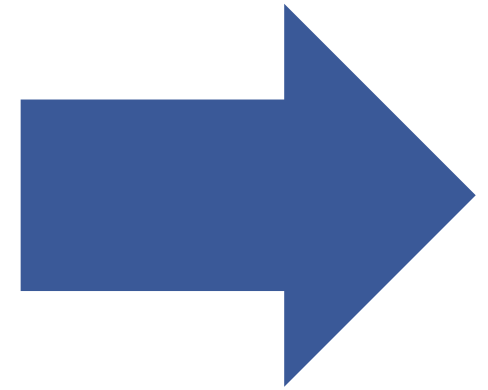
Videos vs pictures

Growth in Messenger

Advance of the Mobile

Era of the the Desktop

© Statista 2016

# Amount of Photo Storage

# So what is an Exabyte?

- 1 Exabyte == 1000 Petabytes
- 1 Petabyte == 1000 Terabytes
- 1 Exabyte = ~250,000 4 TB drives

- **250k drives stacked flat >30 times taller than Seattle Space Needle**

X 30 times!
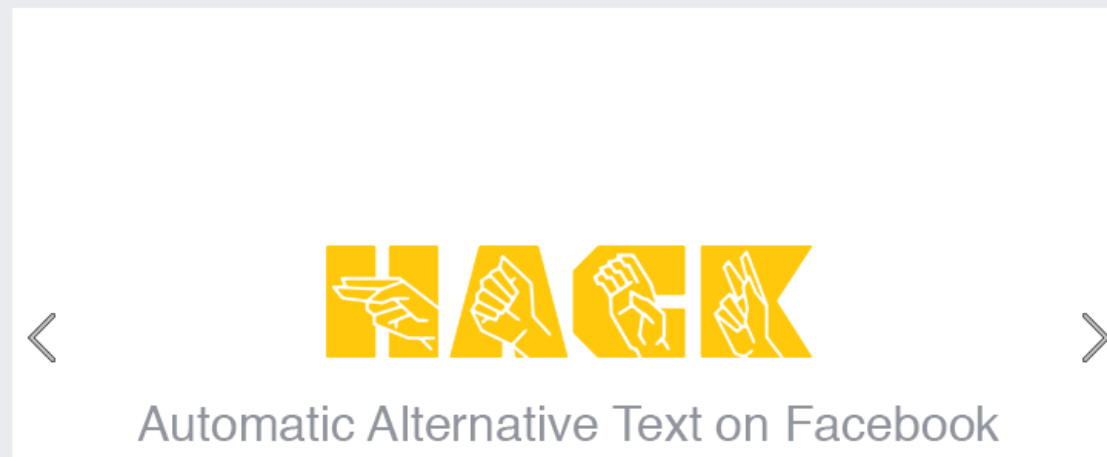
# Data Scientists
# are DW customers

# Data Scientist and DW



Research at Facebook

Home | Our Research | Academic Programs | Publications | Blog | People

## Core Data Science

Home | Publications | People | Blog

HACK

Automatic Alternative Text on Facebook

newsroom

Home | News | Products | Company Info | Directory | Media Gallery | Investor Re

Like | Share

April 4, 2016

## Using Artificial Intelligence to Help Blind People 'See' Facebook

# Data Scientist and DW



Research at Facebook

Home    Our Research    Academic Programs    Publications    Blog    People

## Core Data Science

Home    Publications    People    Blog

Automatic Alter

Research at Facebook

Home    Our Research    Academic Programs    Publications    Blog    People

## NFL Fan Friendships on Facebook

Blog

- Arizona Cardinals
- Atlanta Falcons
- Baltimore Ravens
- Buffalo Bills
- Carolina Panthers
- Chicago Bears
- Cincinnati Bengals
- Cleveland Browns
- Dallas Cowboys
- Denver Broncos
- Detroit Lions
- Green Bay Packers
- Houston Texans
- Indianapolis Colts
- Jacksonville Jaguars
- Kansas City Chiefs
- Miami Dolphins
- Minnesota Vikings
- New England Patriots
- New Orleans Saints
- New York Giants
- New York Jets
- Oakland Raiders
- Philadelphia Eagles
- Pittsburgh Steelers
- Saint Louis Rams
- San Diego Chargers
- San Francisco 49ers
- Seattle Seahawks
- Tampa Bay Buccaneers
- Tennessee Titans
- Washington Redskins

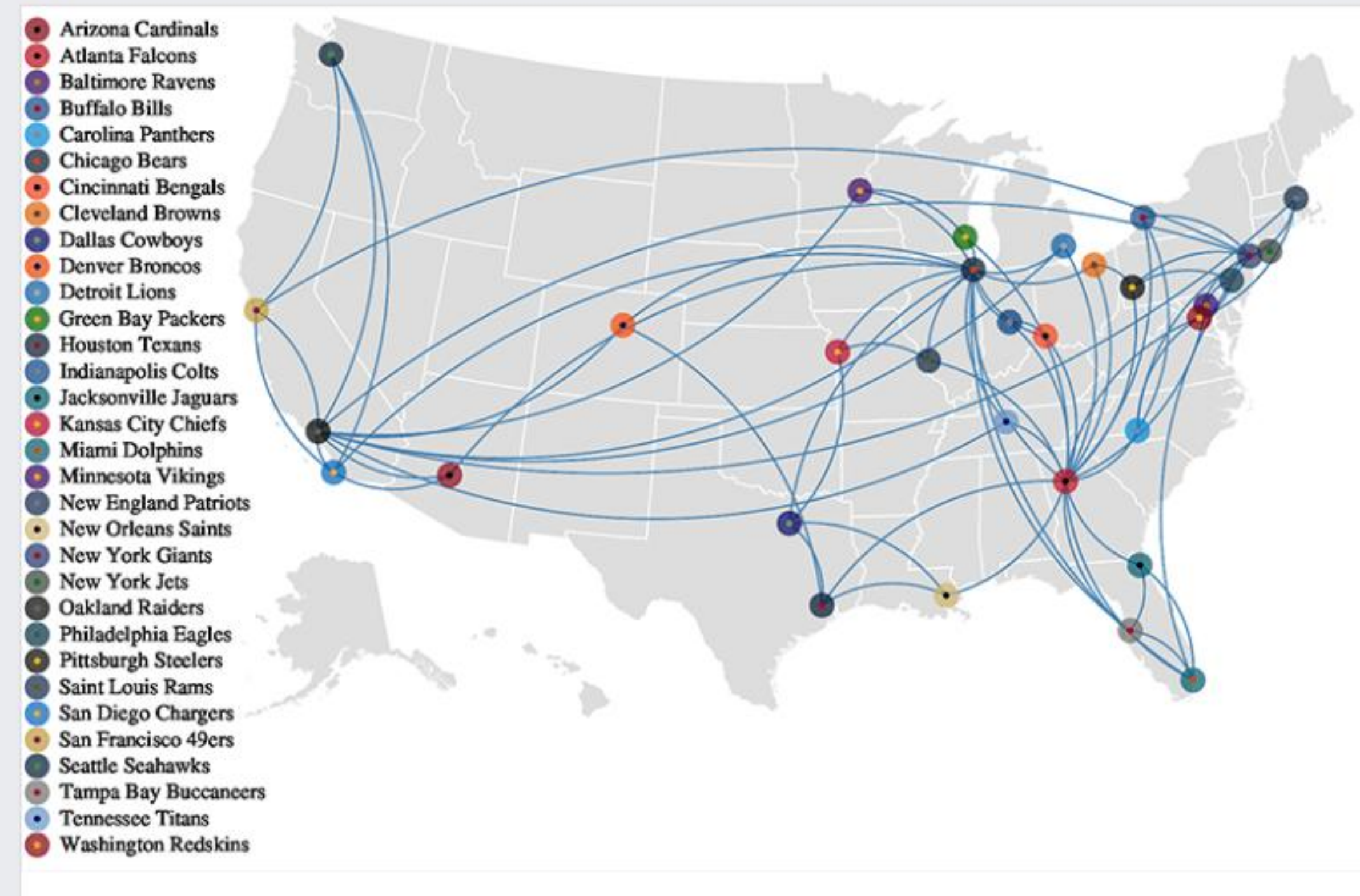newsroom    Like    Share

Home    News    Products    Company Info    Directory    Media Gallery    Investor Re

April 4, 2016

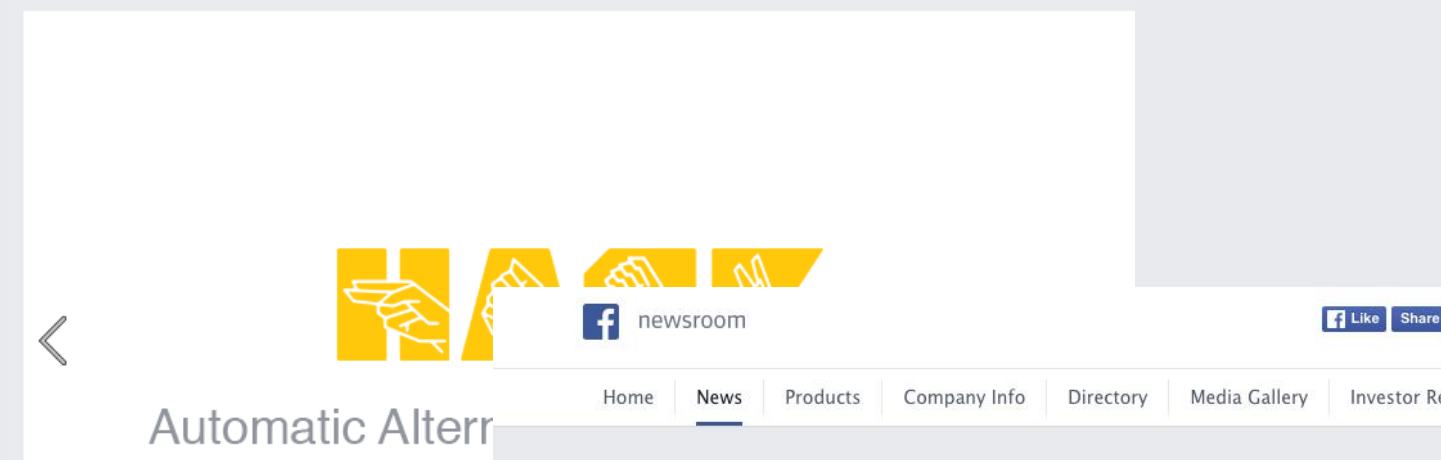## Using Artificial Intelligence to Help Blind People 'See' Facebook

# Data Scientist and DW

# Data Science@FB: dynamic, growing and hard to predict

# Data Warehouse at Facebook

# Data Warehouse at Facebook

# DW @FB: Hive on Hadoop (HDFS)

# DW @FB: Hive on Hadoop (HDFS)

DW @FB: Hive on Hadoop (HDFS)

Permanent and Temp Storage

# Dynamic demand vs. static resource allocation



| | Specification | Quantity |
|---|---|---|
| Storage | 4 TB SAS HDD | 15 |
| CPU | Intel Xeon 20 core | 2 |
| Network | 10 Gbps | 1 |

# Efficiency at hyper scale

| Specification | | Quantity |
|---|---|---|
| Storage | 4 TB SAS HDD | 15 |
| CPU | Intel Xeon 20 core | 2 |
| Network | 10 Gbps | 1 |

# Elasticity
## Splitting DC for Storage and Compute

| Specification | | Quantity |
|---|---|---|
| Storage | 4 TB SAS HDD | 15 |
| CPU | Intel Xeon 20 core | 2 |
| Network | 10 Gbps | 1 |

Compute

Storage

# Elasticity
## Splitting DC for Storage and Compute

| Specification | | Quantity |
|---|---|---|
| Storage | 4 TB SAS HDD | 15 |
| CPU | Intel Xeon 20 core | 2 |
| Network | 10 Gbps | 1 |

Compute

Storage

# Elasticity
## Splitting DC for Storage and Compute

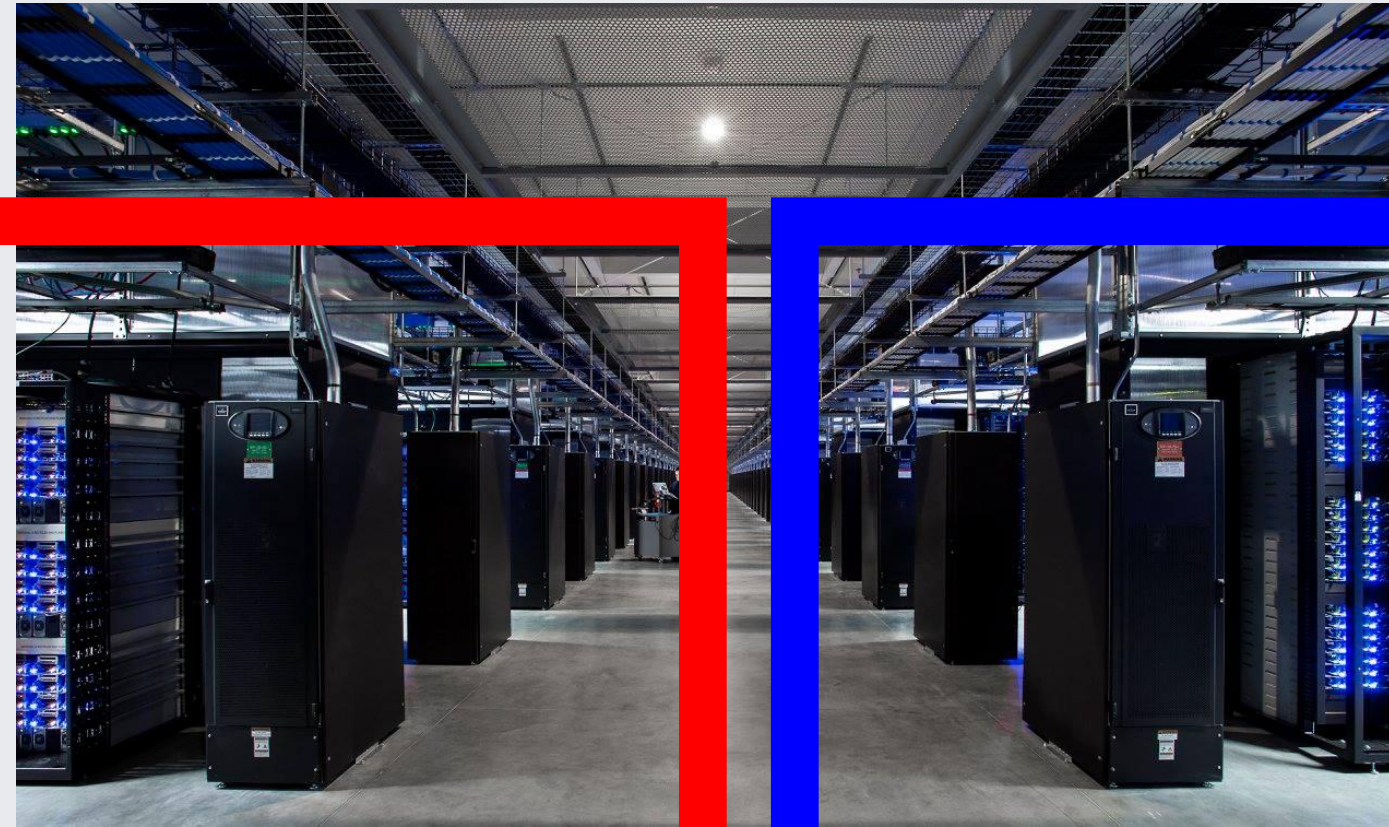| | Specification | Quantity |
|---|---|---|
| Storage | 4 TB SAS HDD | 15 |
| CPU | Intel Xeon 20 core | 2 |
| Network | 10 Gbps | 1 |

Compute

Storage

# Elasticity



**Compute**

- Add more compute if/when needed
- Upgrade to latest Intel's CPUs
- Replace compute every 3 years or sooner

**Storage**

- Keep HDDs longer than 3 years
- Grow storage capacity independent from compute
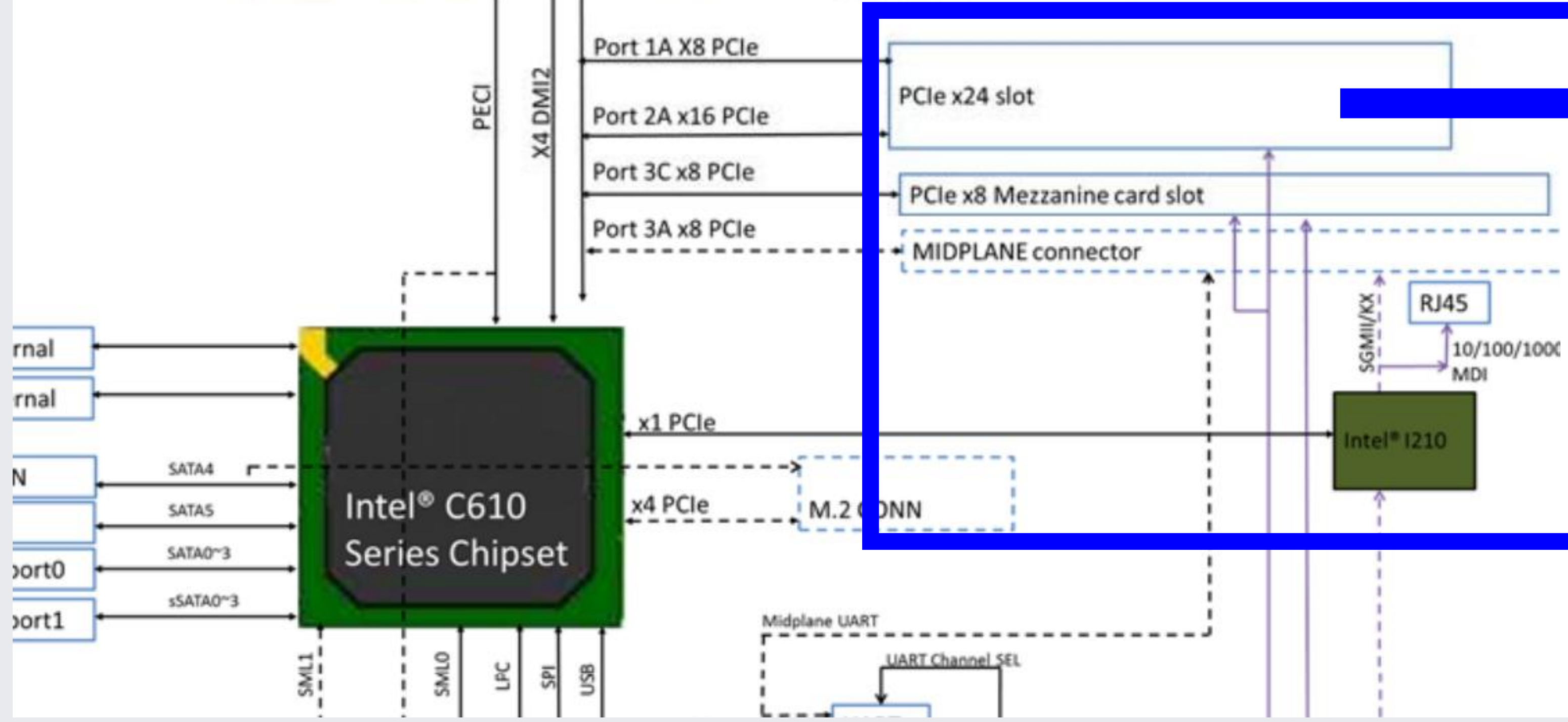- Migrate to 8-10 TB HDDs

**facebook**

# Storage and Compute separation
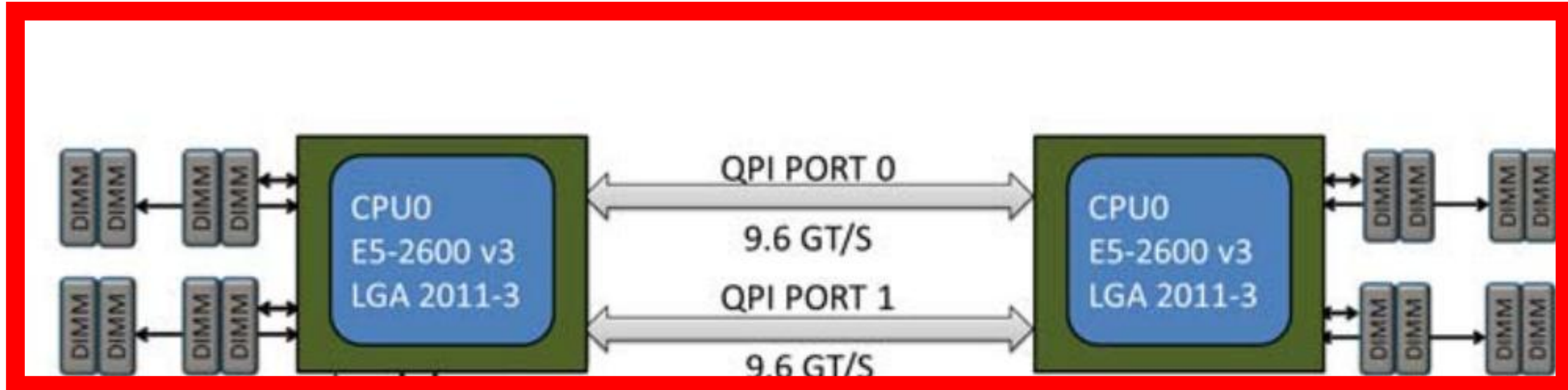
## Is that possible?

# Compute and storage: local vs. remote



Compute

Storage

Open Compute Project • Intel Motherboard • v3.1

# Last ~25 years of HDD evolution



Seagate 94171-327
(300MB)

iPhone 5    16 GB

| Specs | Value |
|---|---:|
| Form | 3.5" |
| Platters | 5 |
| Heads | 9 |
| Capacity | 300MB |
| Interface | SCSI |
| Seek time | 17ms |
| Data transfer rate | 1 MB/sec |

# Last ~25 years of HDD evolution

| Manufacturer | Capacity | Transfer speed (MB/sec) | Time to read all data | Year |
|---|---|---|---|---|
| Seagate | 300MB | 1 | 5 mins | 1990 |
| IBM | 10GB | 12 | 13 mins | 1998 |
| Seagate | 750GB | 72 | 3 hours | 2006 |
| Hitachi | **1TB** | 85 | **3.2 hours** | 2007 |
| WD/Seagate | **4TB** | 100 | **11 hours** | 2012 |
| Seagate | **8TB** | 120 | **18 hours** | 2014 |

# Last ~25 years of Ethernet



**Evolution of Network Bandwidth**

Legend:
- Ethernet (green)
- Token Ring* (red)
- MOST* (orange)
- Fibre Channel (teal)

Ethernet data points: Legacy 3Mbs (1980), 10 Mbps (1990), 100 Mbps (1996), 1 Gbps (2000), 10 Gbps (2004), 40 Gbps (2010), 100 Gbps (2012)

Token Ring*: 4 Mbps (1987), 16 Mbps (1990)

Fibre Channel: 800 Mbps (2000), 4 Gbps (2001), 8 Gbps (2008)

MOST*: 25 Mbps (2000), 50 Mbps, 150 Mbps (2011)

*Shared network architecture

# Last ~25 years of HDD vs. Ethernet



**Evolution of Network Bandwidth**

Legend:
- Ethernet
- Token Ring*
- MOST*
- Fibre Channel

*Shared network architecture

© 2013 Broadcom Corporation. All rights reserved

| Manufacturer | Capacity | Transfer speed (MB/sec) | Time to read all data | Year |
|---|---|---|---|---|
| Seagate | 300MB | 1 | 5 mins | 1990 |
| IBM | 10GB | 12 | 13 mins | 1998 |
| Seagate | 750GB | 72 | 3 hours | 2006 |
| Hitachi | 1TB | 85 | 3.2 hours | 2007 |
| WD/Seagate | 4TB | 100 | 11 hours | 2012 |
| Seagate | 8TB | 120 | 18 hours | 2014 |

## Ethernet/HDD speed ratio
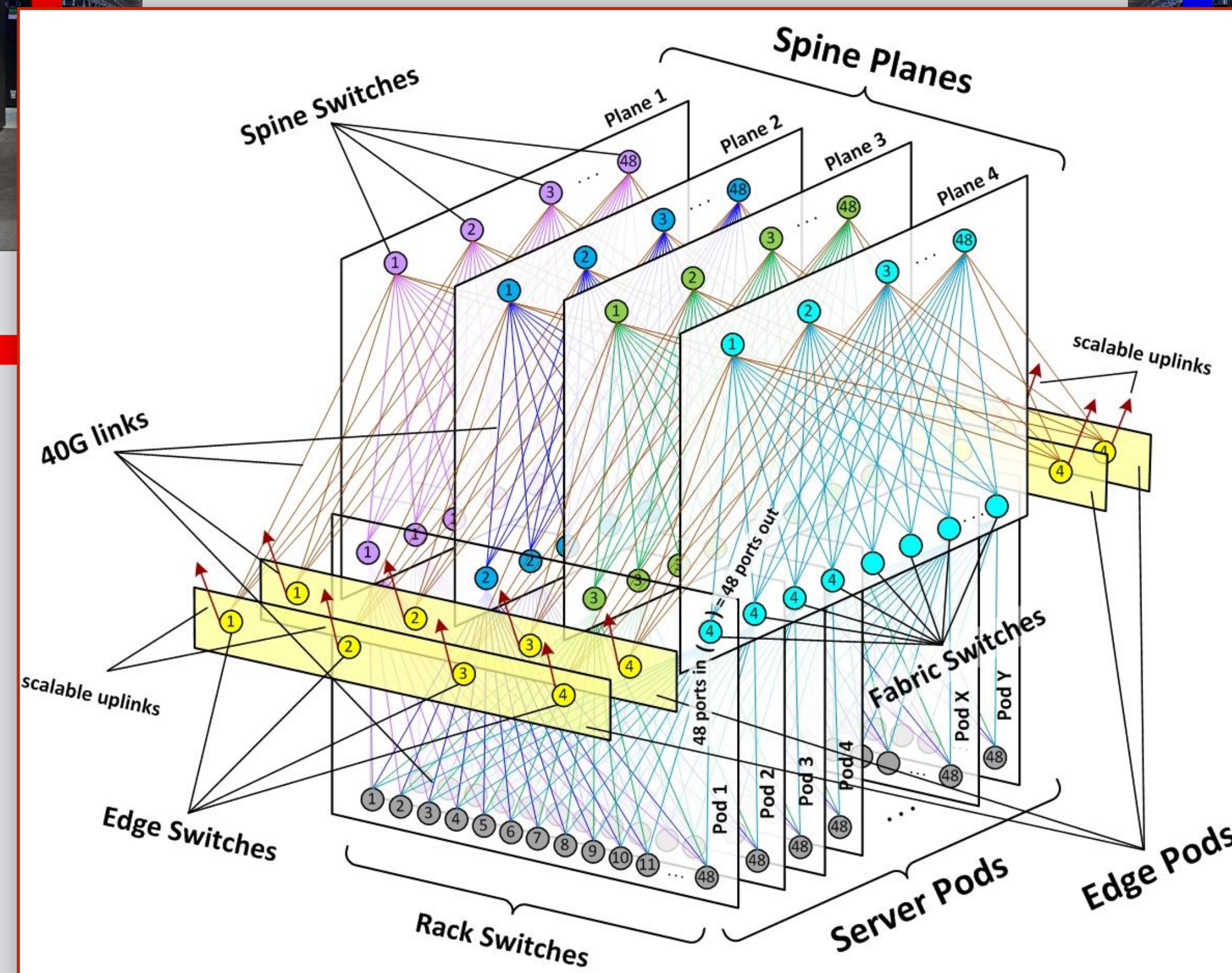


Ethernet vs HDD speed ratio

# When Ethernet is faster than HDDs
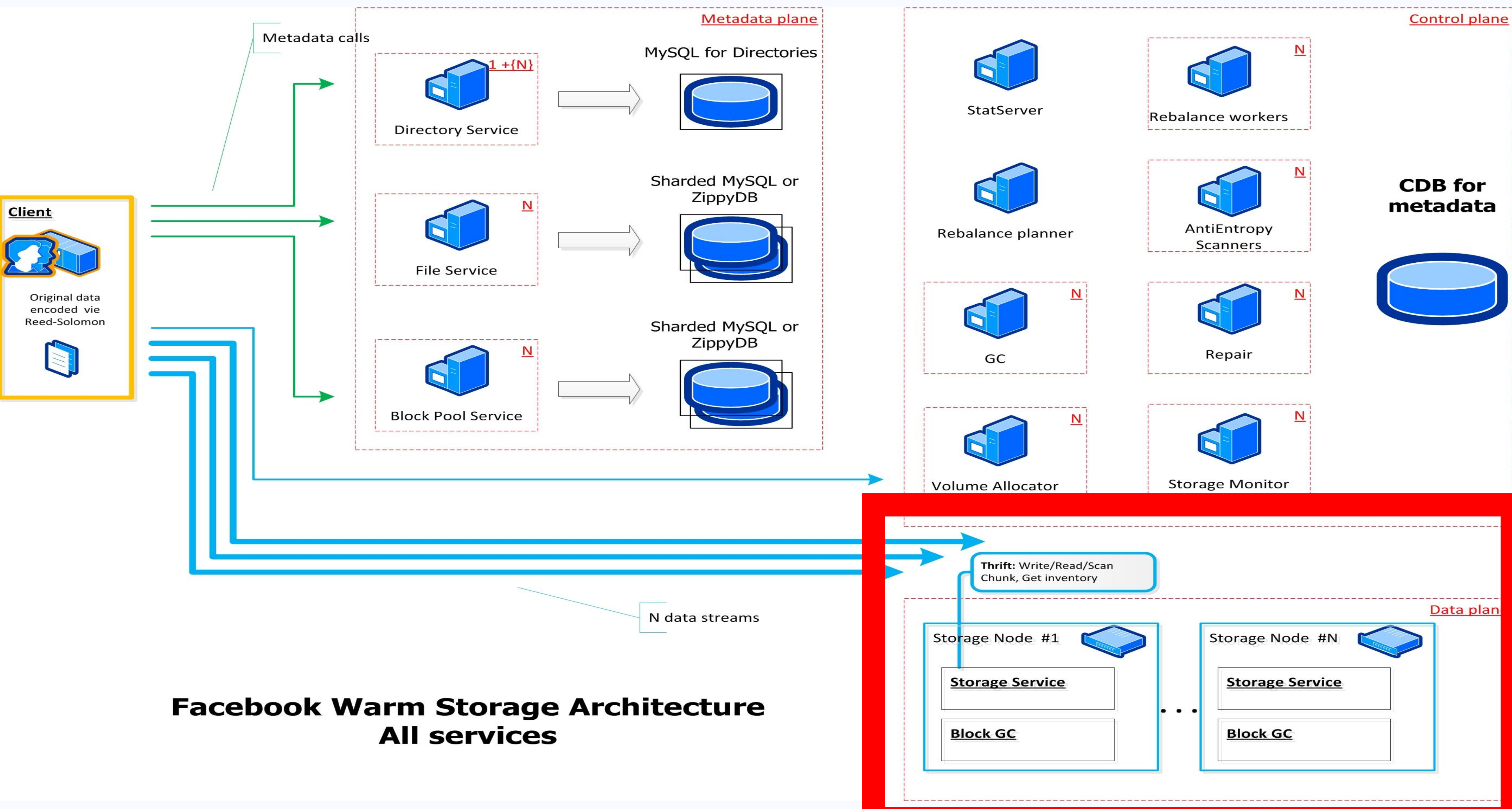# And Data Center network is 'flat'

Compute
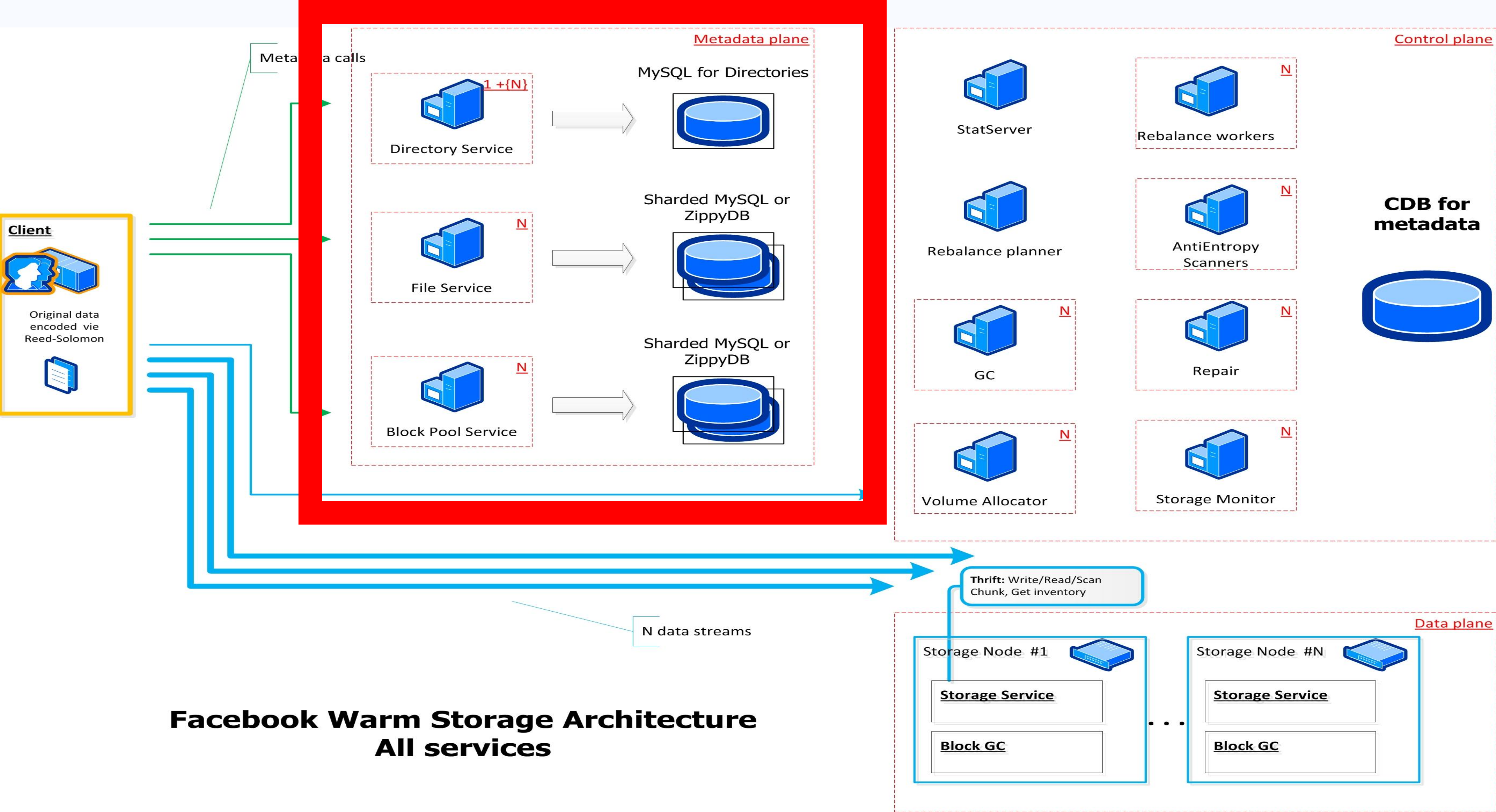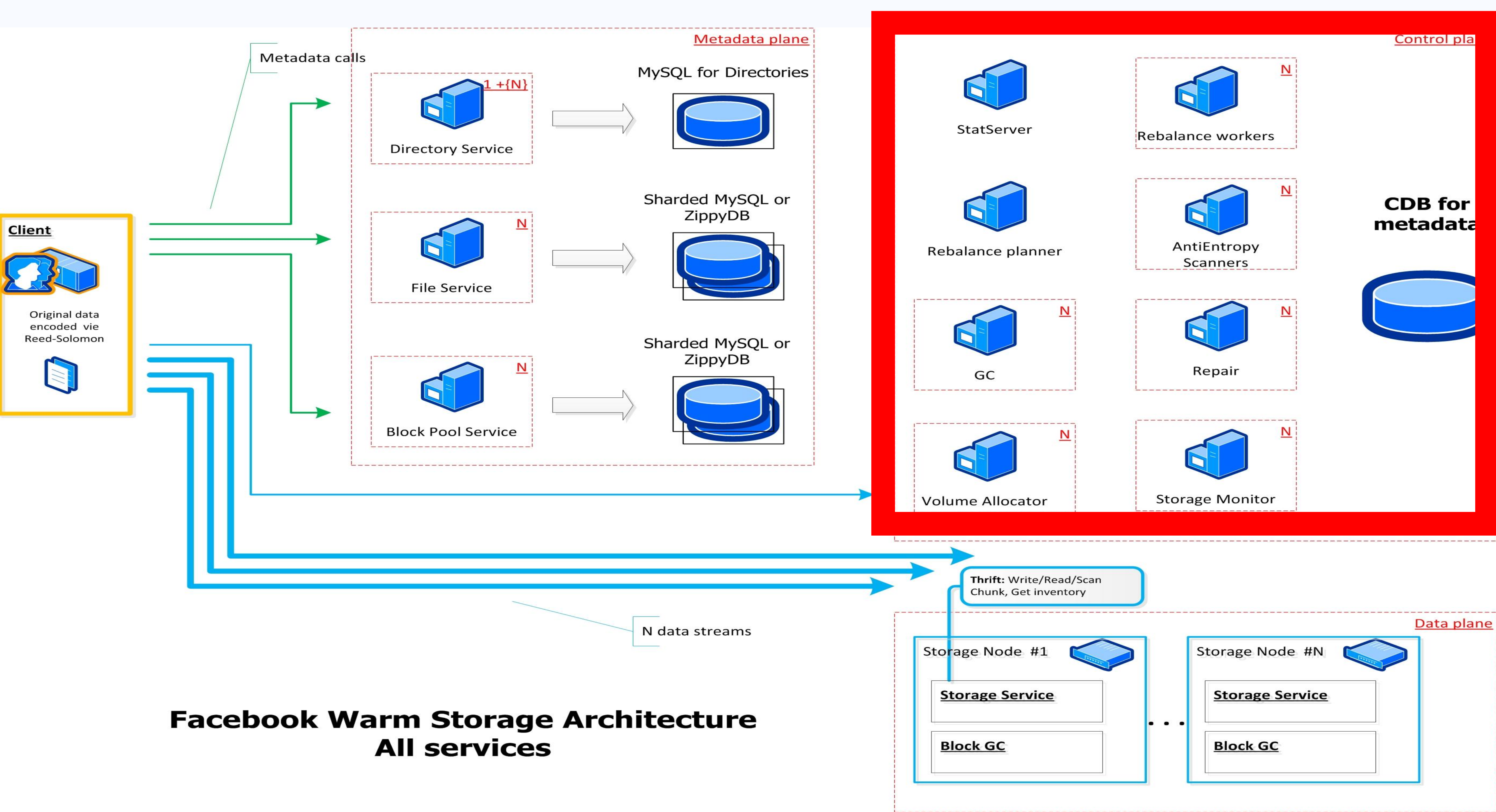
Storage

Facebook Warm Storage Architecture
All services

**Facebook Warm Storage Architecture All services**
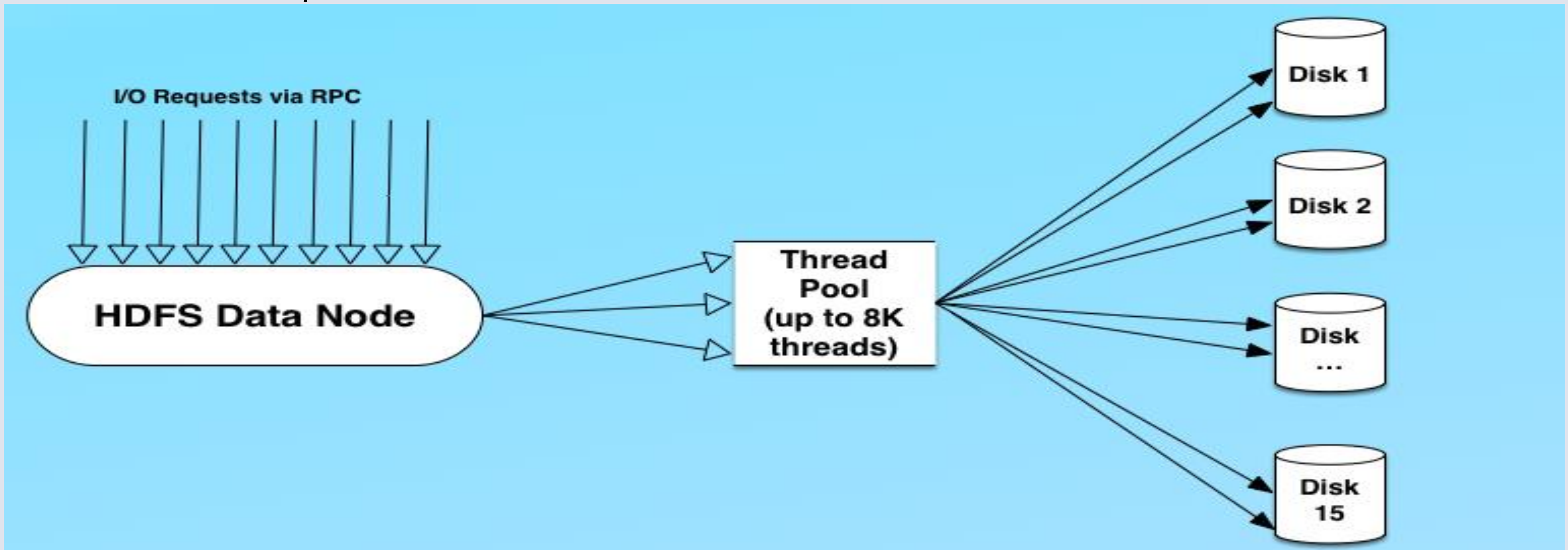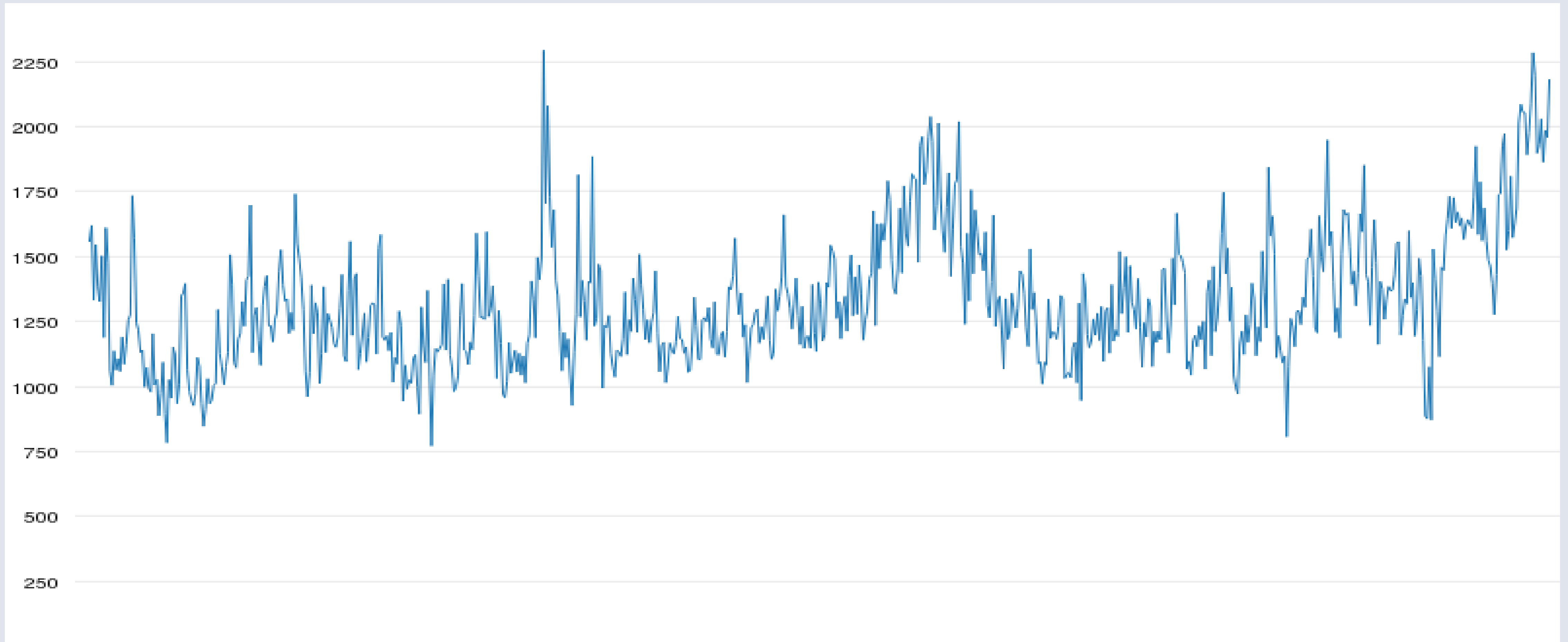
**Facebook Warm Storage Architecture
All services**

# Storage Service IO Model
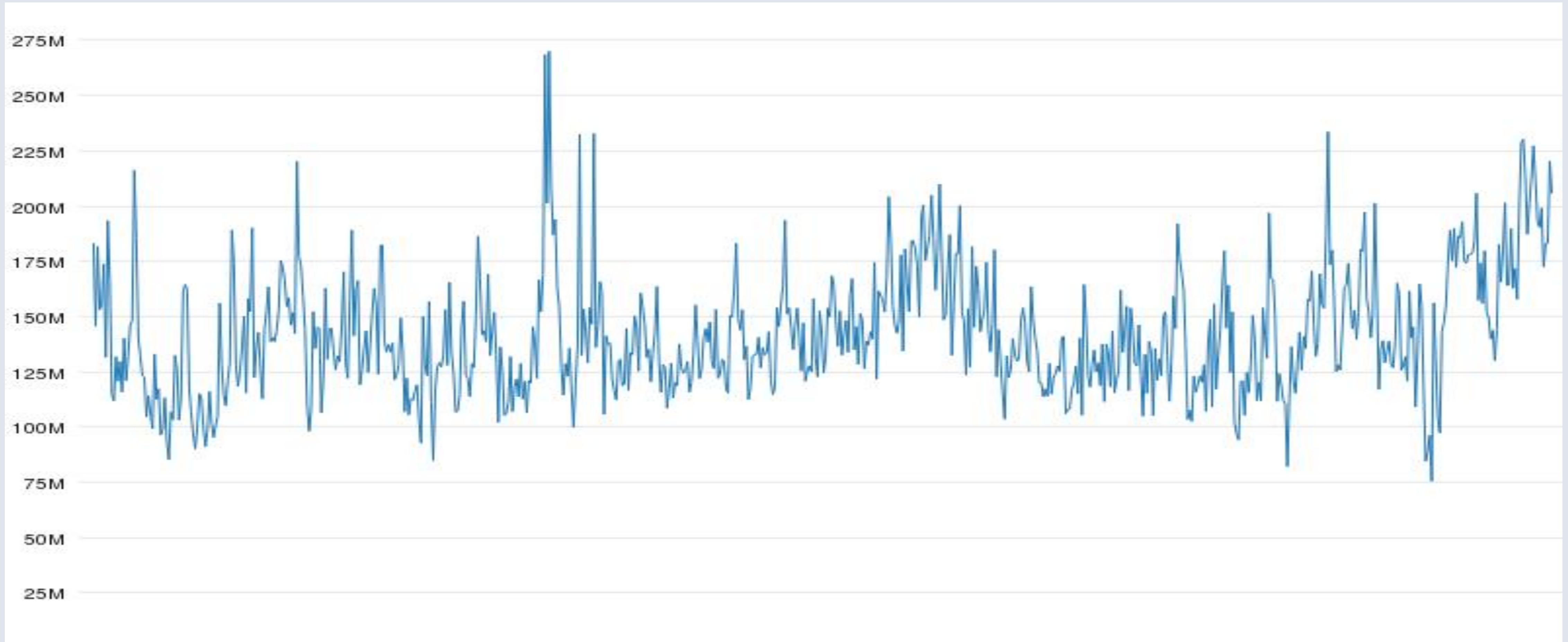
# HDFS Data Node I/O Model

- HDFS Disk Thrashing
  - 4k – 8K I/O threads

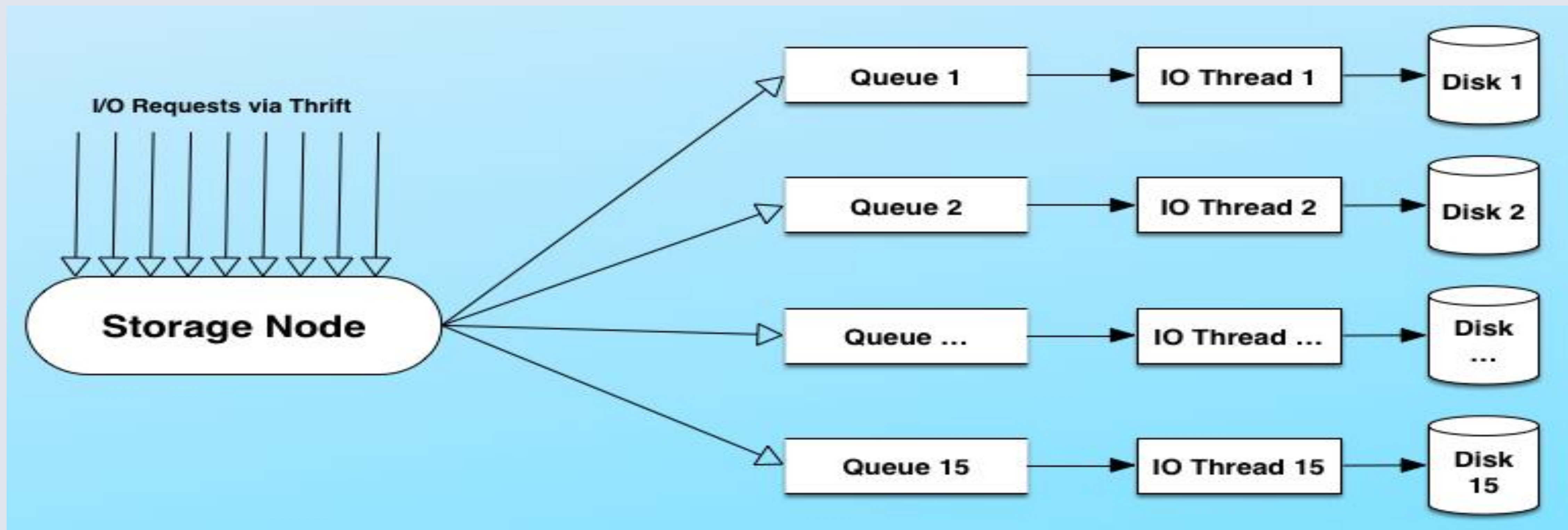# HDFS IOPS/Node – very large cluster

# HDFS (MB/s per Node – 15 disks )
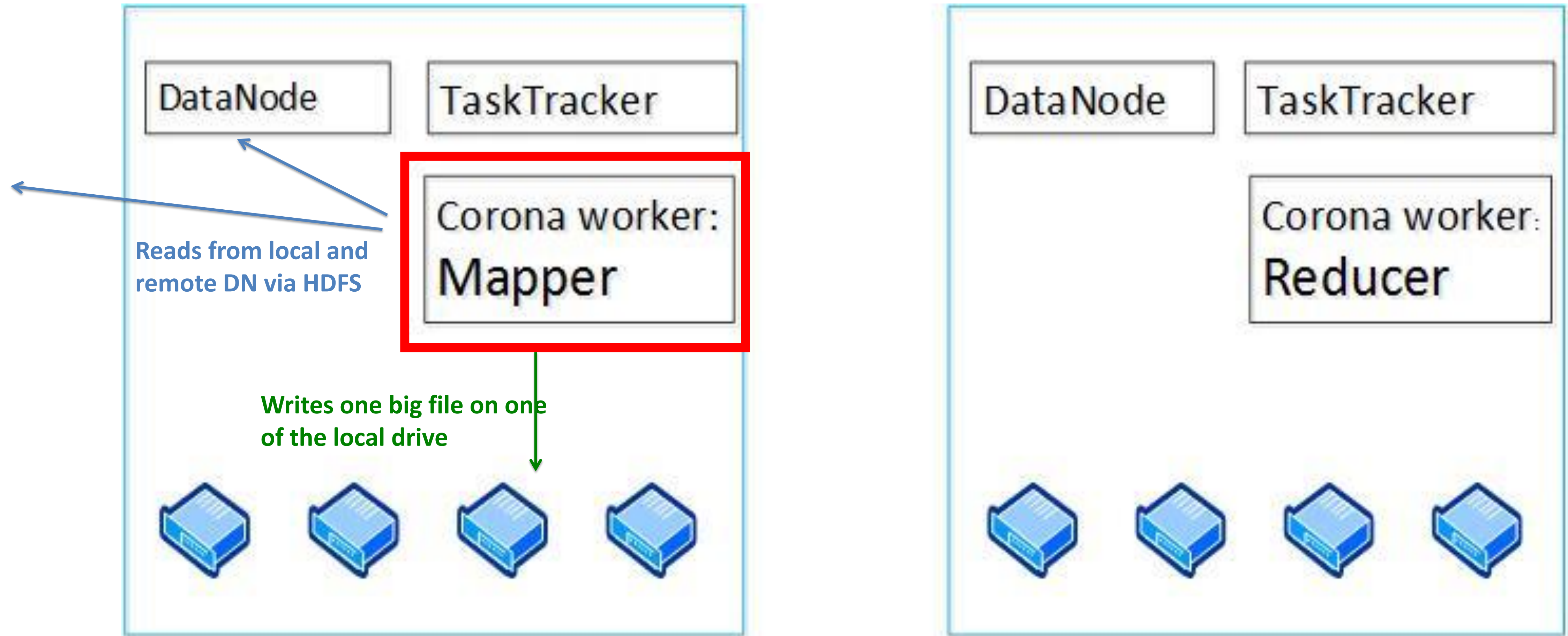
# Storage Node Single Thread I/O Model

- 1 thread per disk
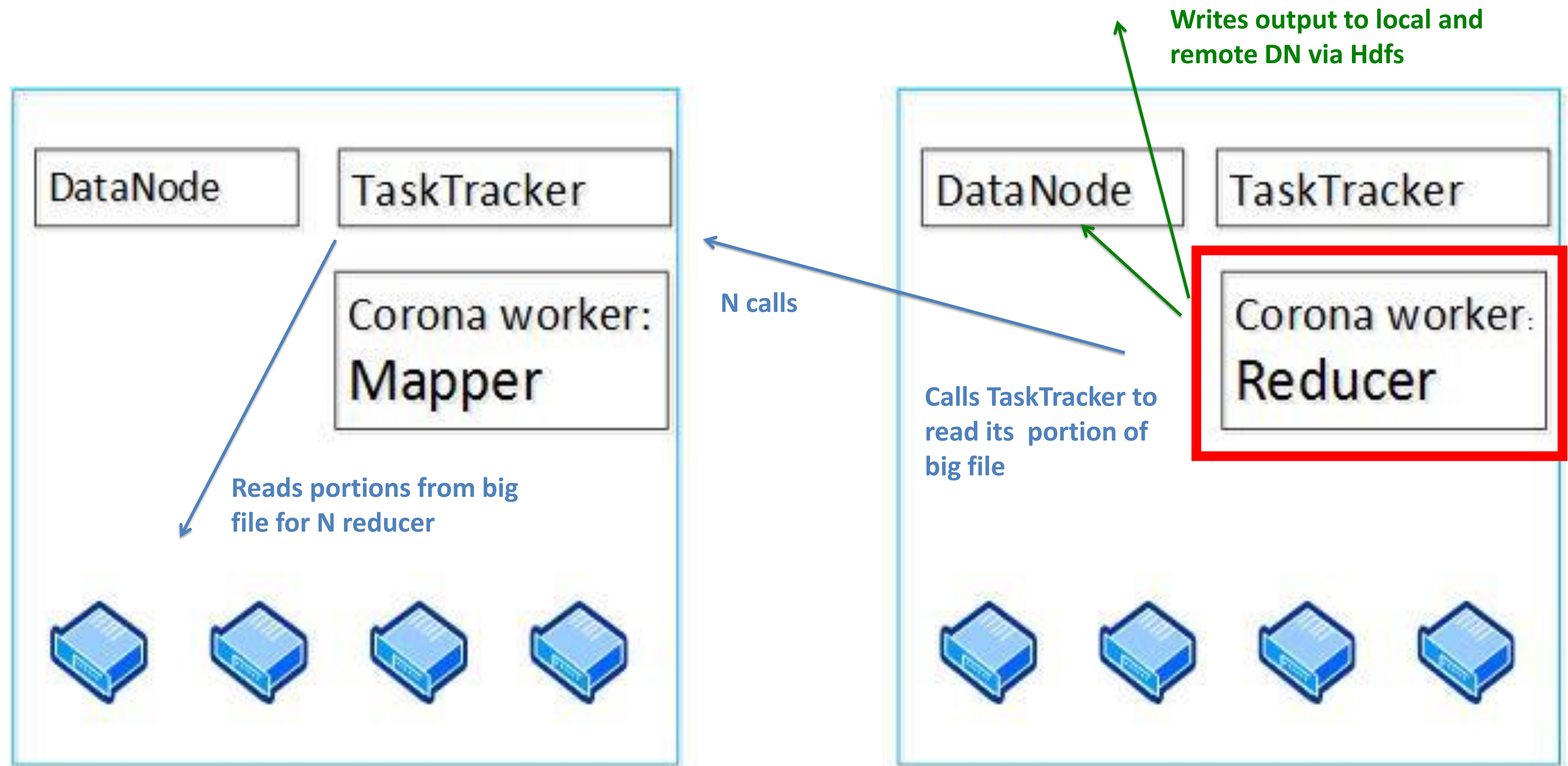- I/O operations are pushed to a priority queue

# facebook

# TempFS and Dis Aggregated storage

# Co-located HDFS Map Reduce

# Co-located HDFS Map Reduce

**Writes output to local and remote DN via Hdfs**

DataNode    TaskTracker

Corona worker:
**Mapper**

**N calls**

**Reads portions from big file for N reducer**

DataNode    TaskTracker

Corona worker:
**Reducer**

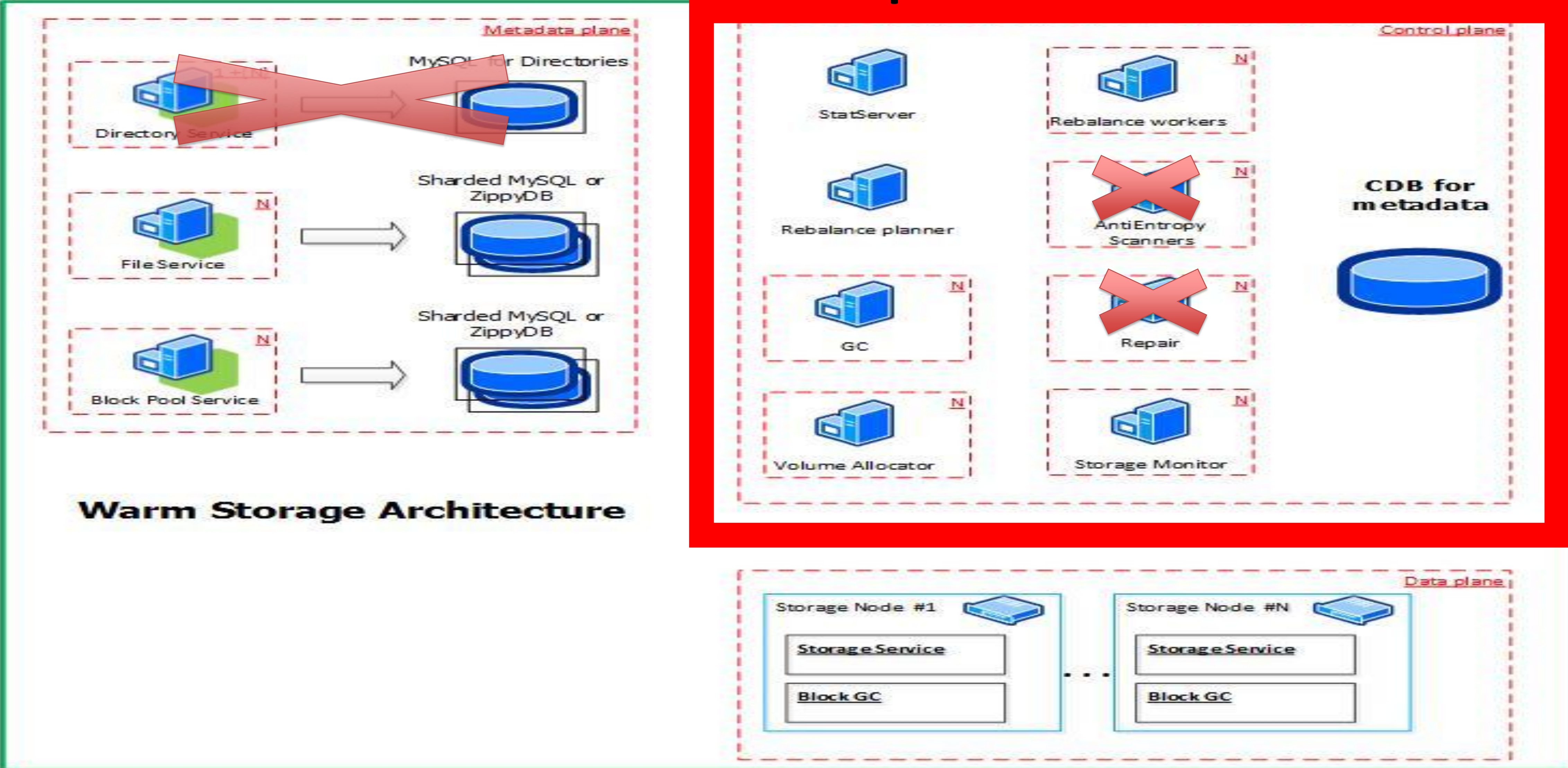**Calls TaskTracker to read its portion of big file**

# Disaggregated Map Reduce

# TempFS Architecture

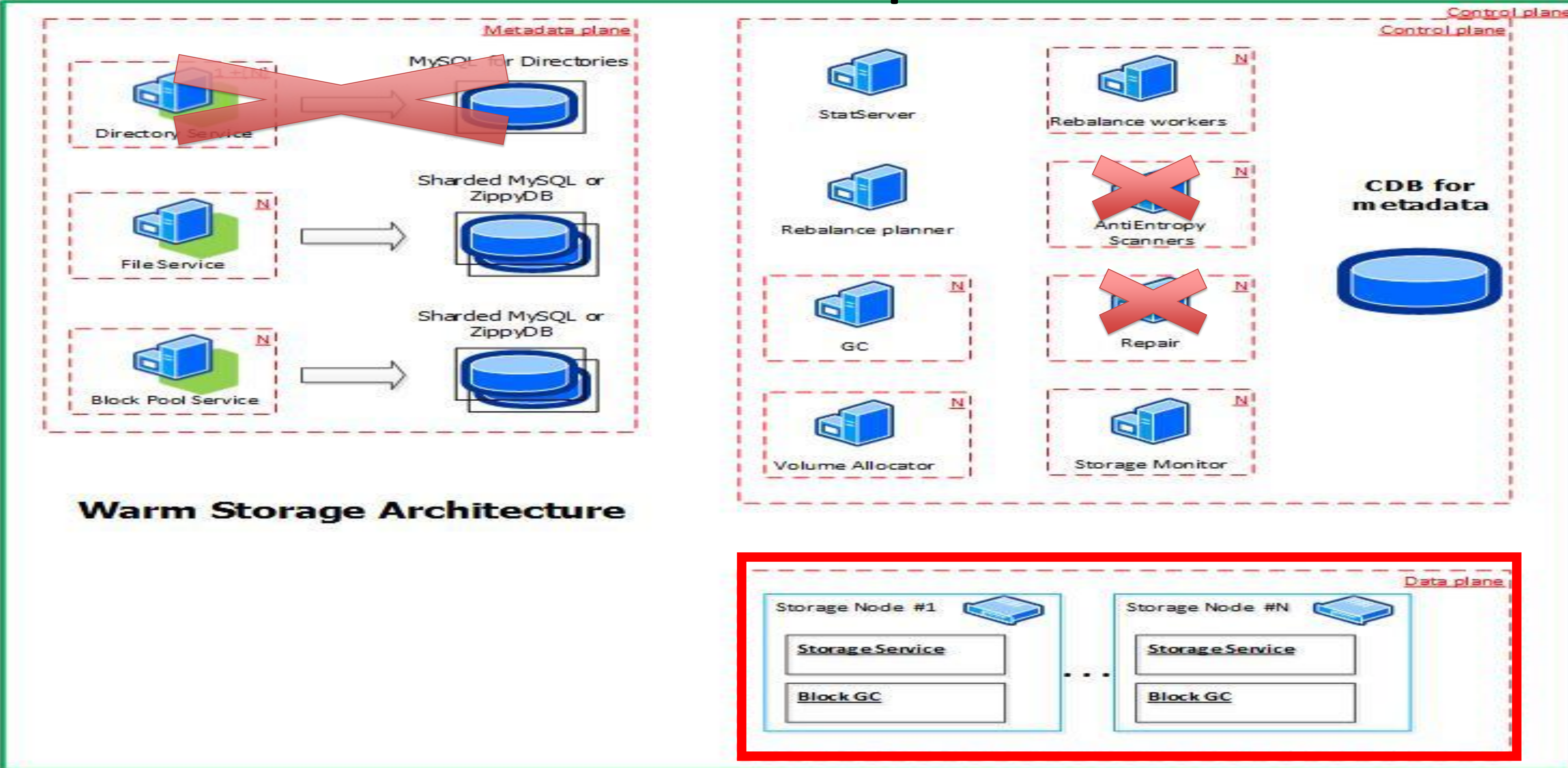# Changes in Warm Storage Architecture for TempFS

# Changes in Warm Storage Architecture for TempFS

# Changes in Warm Storage Architecture for TempFS

# Future Work

- Streaming protocol
- Optimizing corona for Disagg

**facebook**

# Conclusions

# •Conclusions

1. The end is close for Moore's and Kryder's  laws
2. But networking is still improving
3. Efficiency at hyper scale is hard

4. Storage and Compute separation gives better choices and helps with efficiency