

# Facebook's Tectonic Filesystem: Efficiency from Exascale

**Satadru Pan**<sup>1</sup>, Theano Stavrinou<sup>1,2</sup>, Yunqiao Zhang<sup>1</sup>, Atul Sikaria<sup>1</sup>, Pavel Zakharov<sup>1</sup>,  
Abhinav Sharma<sup>1</sup>, Shiva Shankar P<sup>1</sup>, Mike Shuey<sup>1</sup>, Richard Wareing<sup>1</sup>, Monika  
Gangapuram<sup>1</sup>, Guanglei Cao<sup>1</sup>, Christian Preseau<sup>1</sup>, Pratap Singh<sup>1</sup>, Kestutis  
Patiejunas<sup>1</sup>, JR Tipton<sup>1</sup>, Ethan Katz-Bassett<sup>3</sup>, and Wyatt Lloyd<sup>2</sup>

<sup>1</sup>Facebook, Inc., <sup>2</sup>Princeton University, <sup>3</sup>Columbia University

# Exabyte-Scale Storage Use Cases at FB

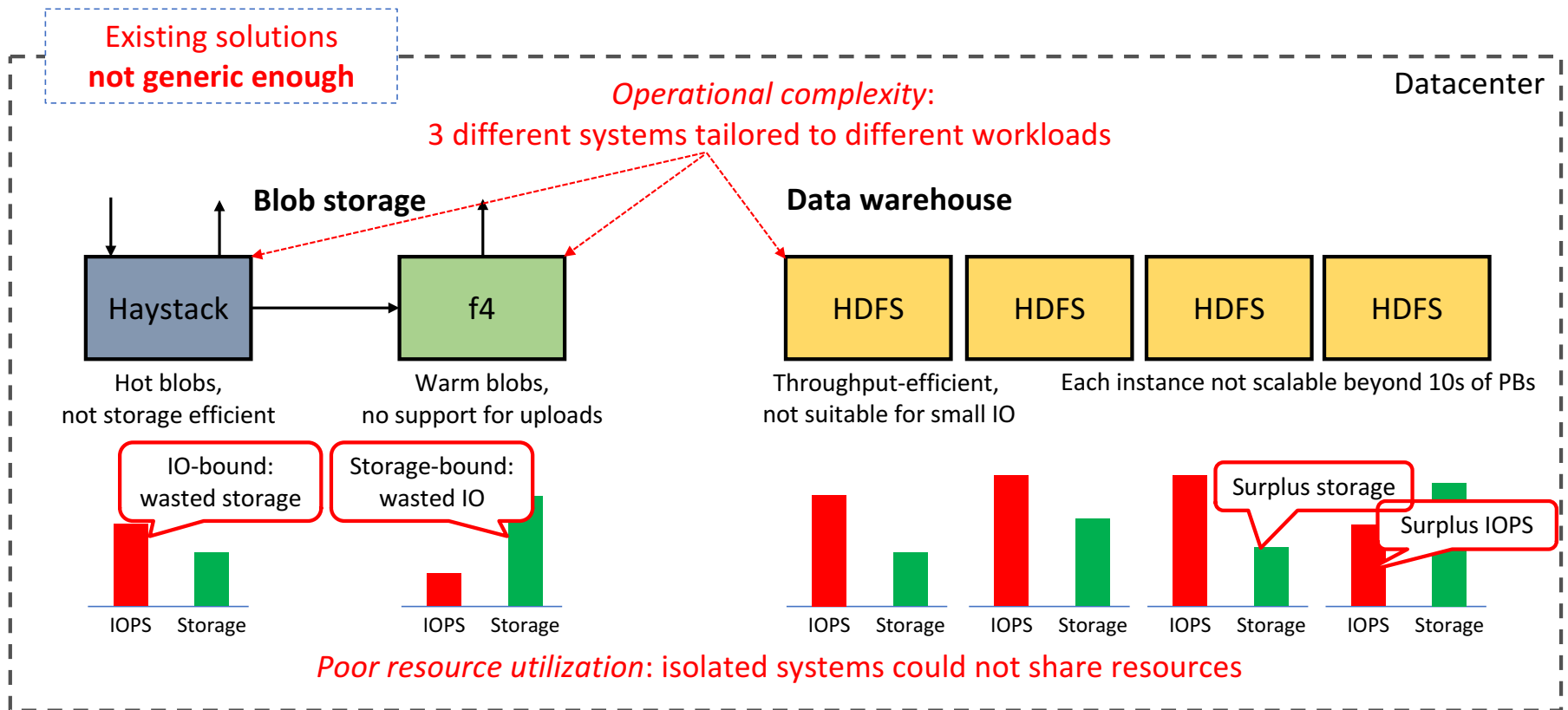
## **Blob storage**

- Photos and videos in Facebook, Messenger attachments
- Exabytes of data
- Several KBs to several MBs in size
- Latency sensitive

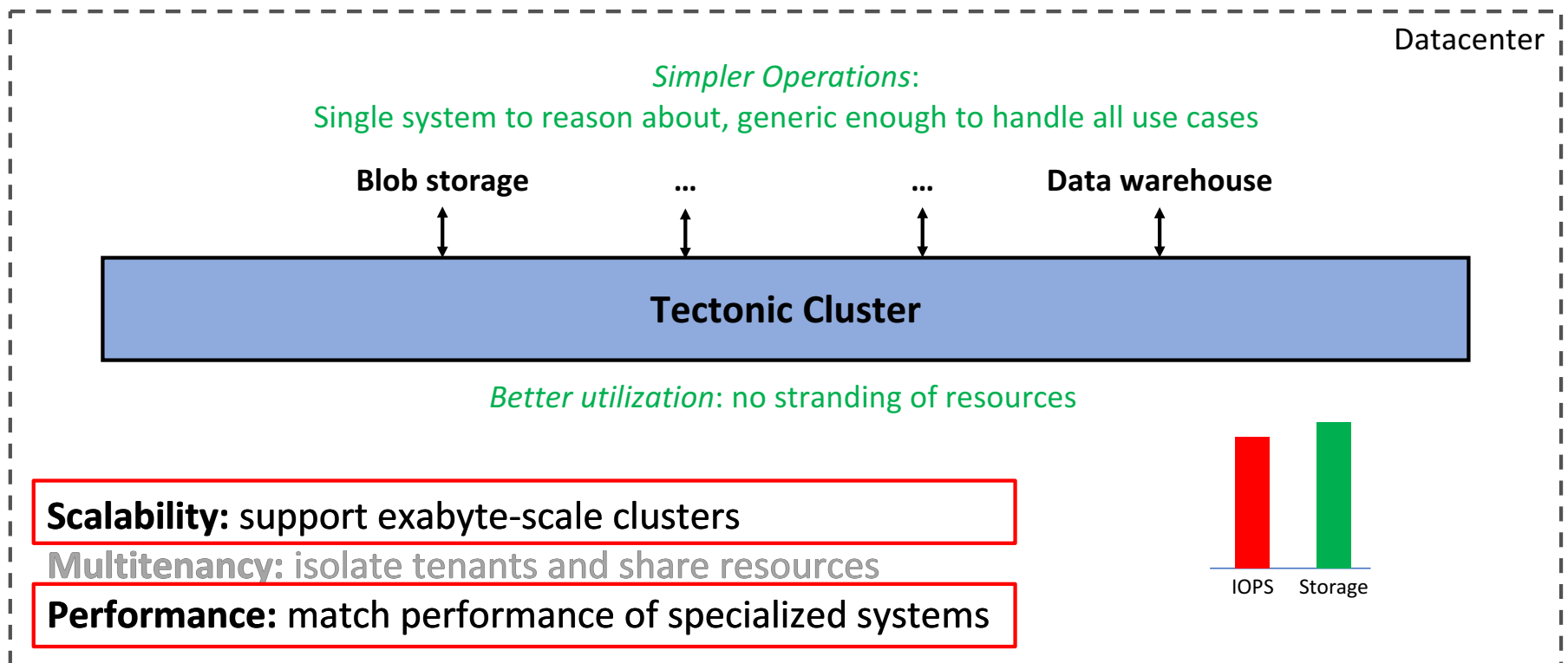
## **Data warehouse**

- Hive tables for data analytics, machine learning
- Exabytes of data
- Reads are order of multiple MBs, writes are 10s of MBs
- Throughput sensitive

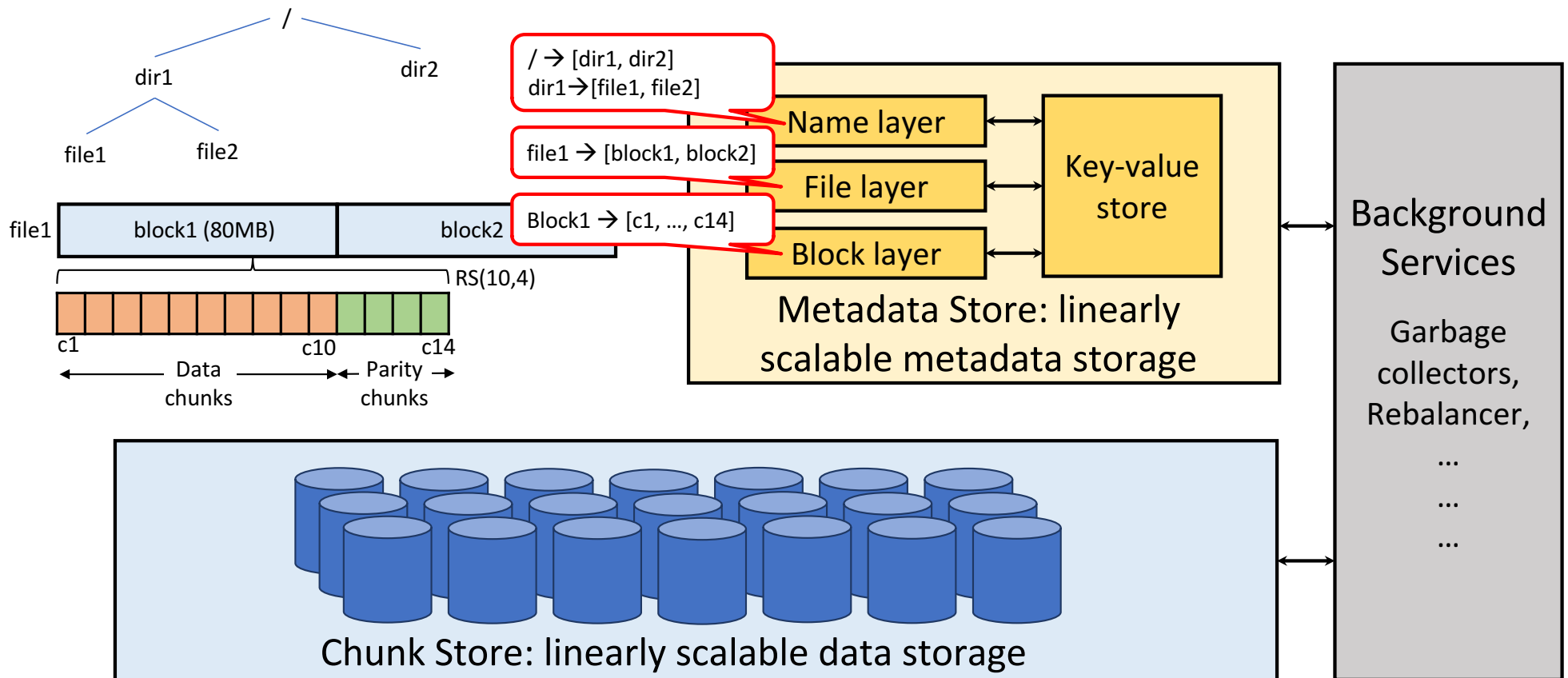
# Storage Infrastructure Before Tectonic



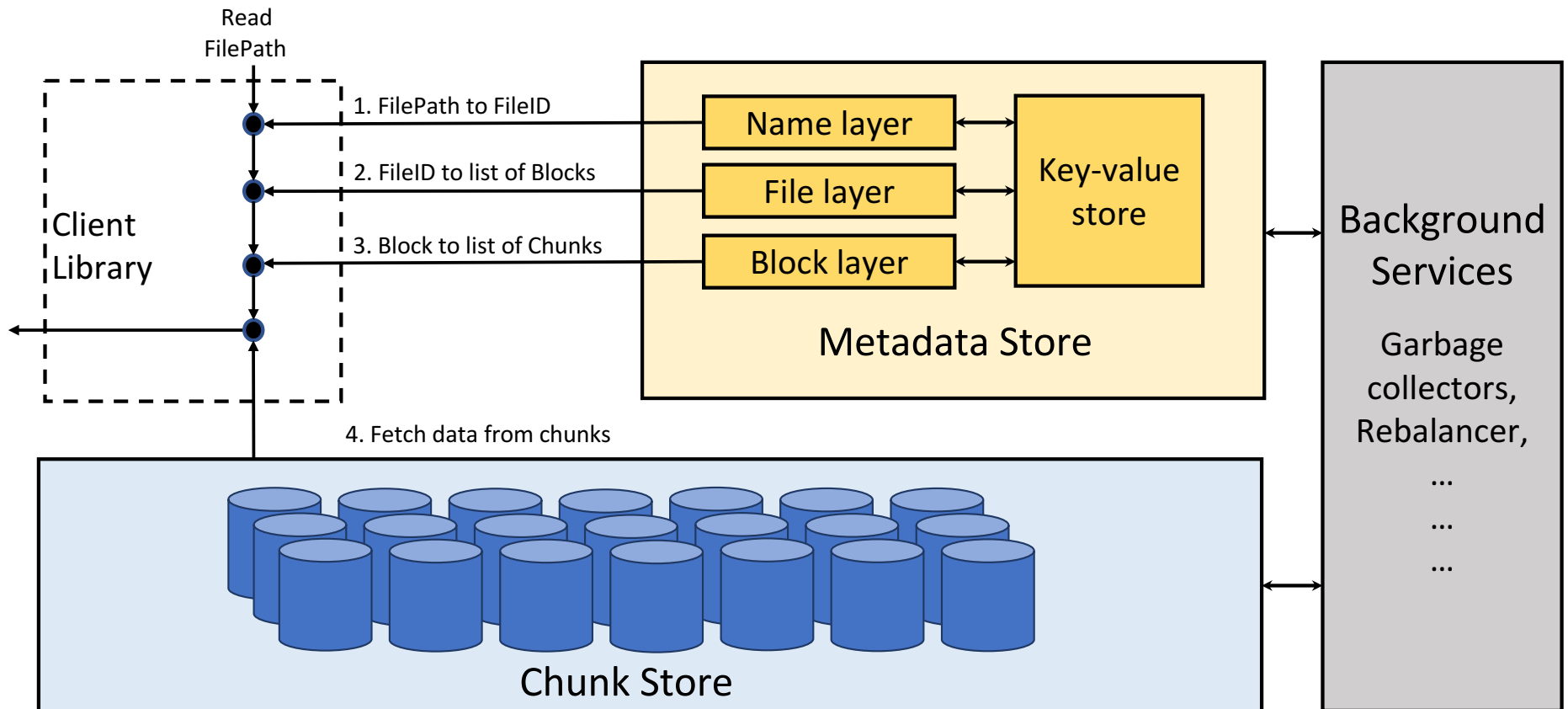
# Tectonic Overview



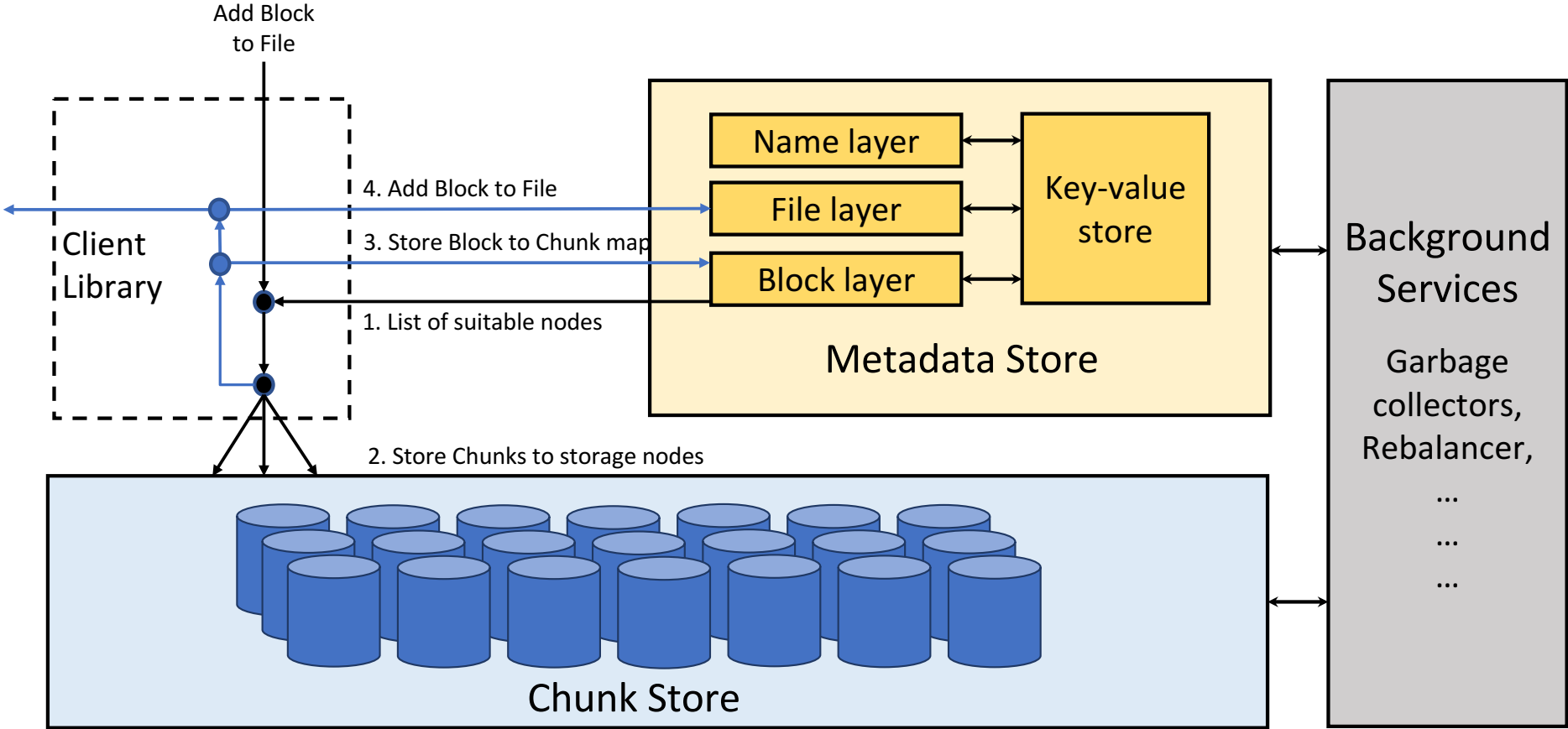
# Scalability: Support Exabyte Scale Clusters



# Scalability: Support Exabyte Scale Clusters



# Scalability: Support Exabyte Scale Clusters



# Performance: Match Specialized Systems

- Specialized storage systems optimize for the specific access pattern and performance requirements
- Tectonic uses *tenant-specific optimizations* to match the performance of specialized systems
- Optimizations are enabled by the Client Library, which runs in application binary
- Client library allows flexible and varying composition of Tectonic operations, which can be configured according to the needs of the tenant



# Tenant-specific Optimizations: Appends

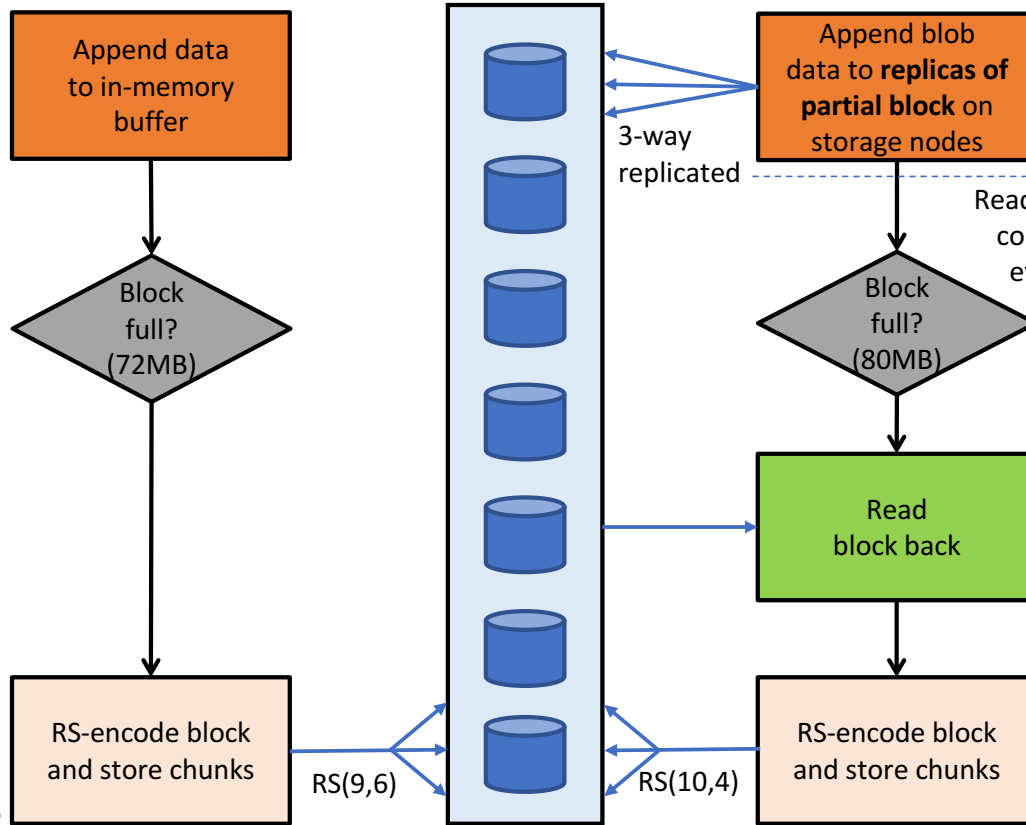
## Data warehouse

Files are large  
(100s of MBs)

Files are read  
after the creator  
closes the file

Minimize bytes  
written to store file  
to improve overall  
throughput

Read-after-write  
consistency  
only after file close



## Blob storage

Blob sizes are small  
(100s of KBs)

Blobs appended to log  
structured file

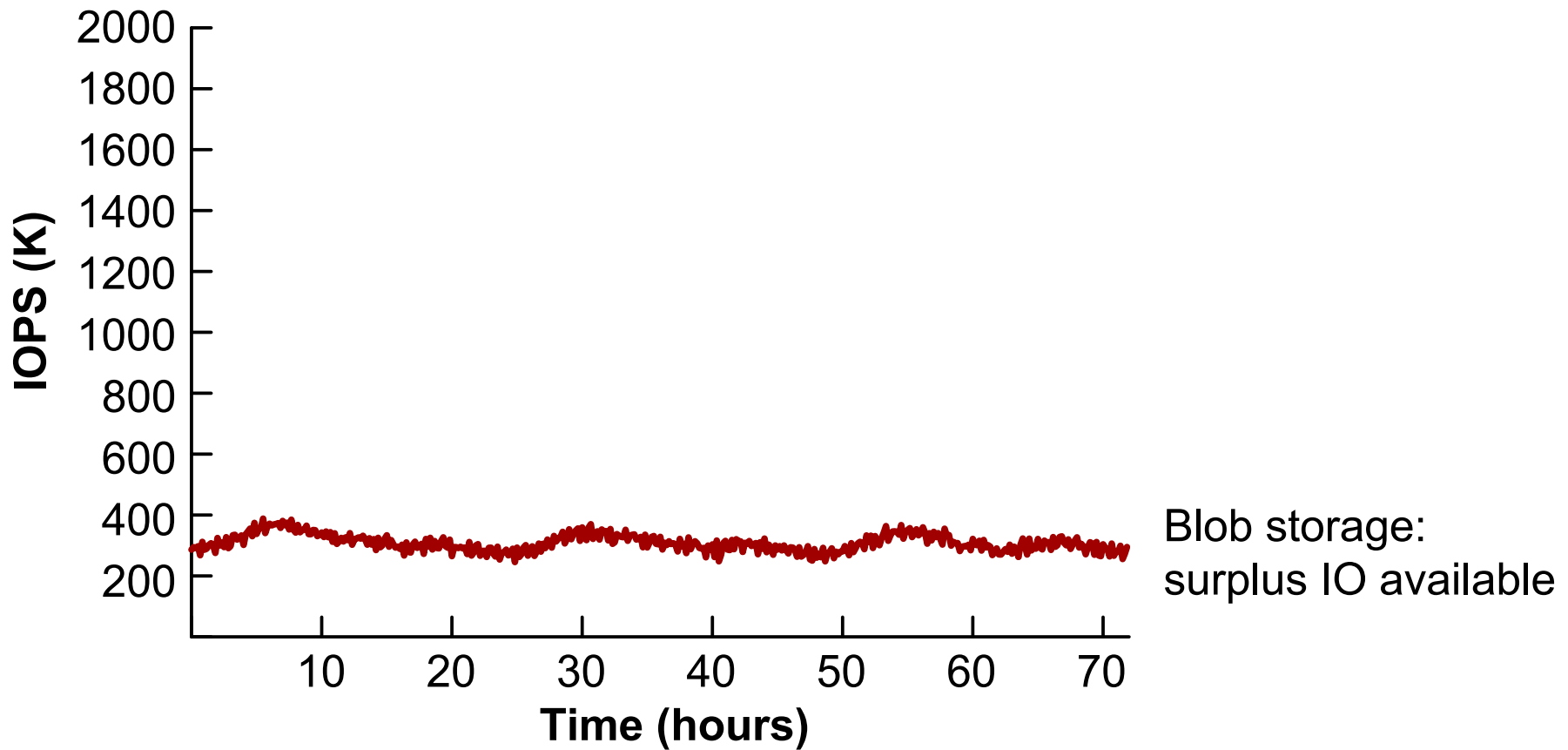
Blobs need to be persisted  
before acknowledging  
upload

Minimize latency  
for blob uploads,  
Later optimize  
storage

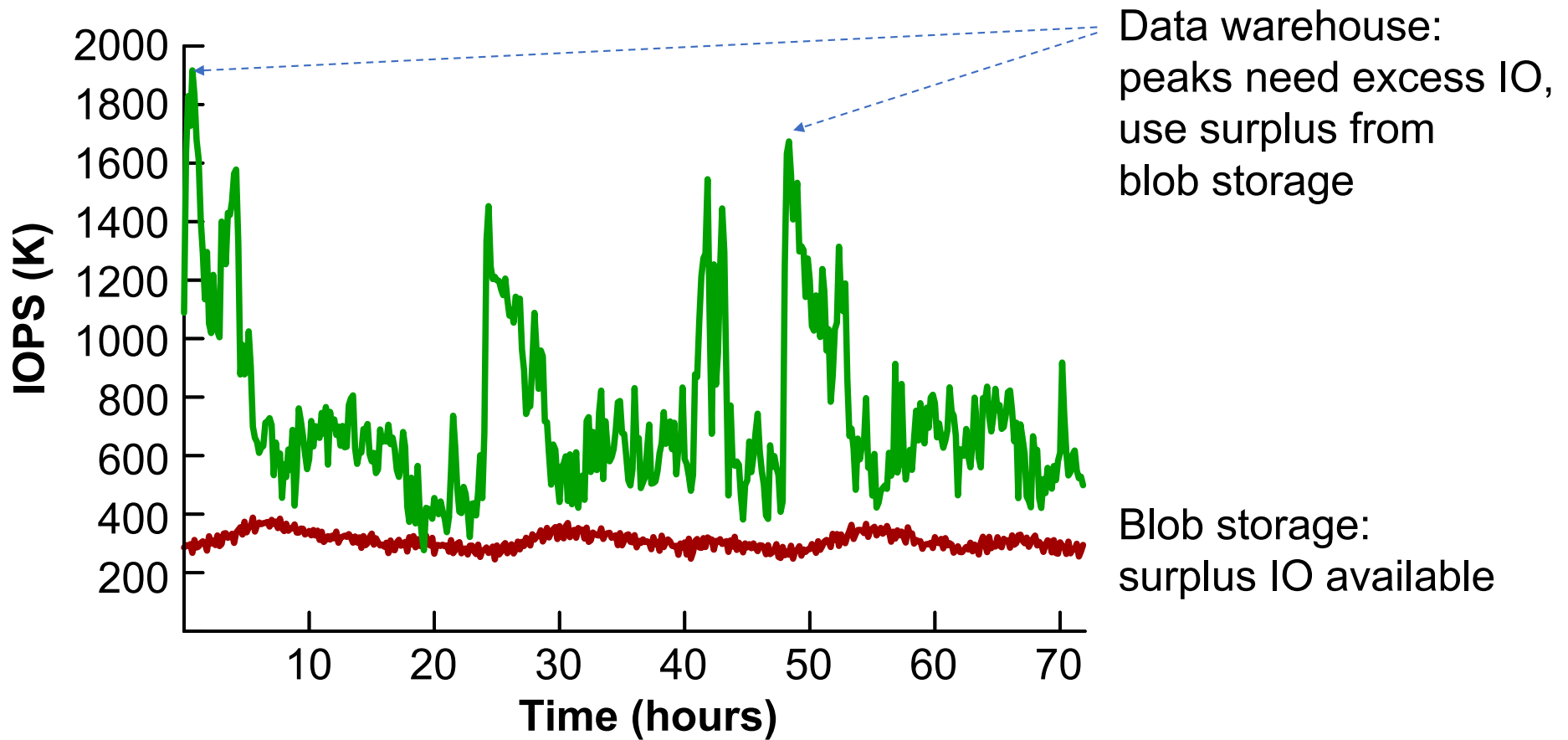
# Results

- Tectonic clusters are ~10x the size of HDFS clusters, which simplifies production operations
- Blob storage latency in Tectonic comparable to Haystack
- In a multitenant cluster, data warehouse uses surplus IO from blob storage to serve its peaks

# Efficiency From Storage Consolidation



# Efficiency From Storage Consolidation



# Tectonic Provides Datacenter-Scale Storage

- Replaced previous constellation of specialized storage systems
  - Simpler operations
  - Better resource utilization
- Tectonic's design addresses the key challenges:
  - Scalability: disaggregated linearly scalable components
  - Performance: tenant-specific optimizations via client library
  - ...

**Thank You**