



C2: LANL-Seagate's Early Prototype for Near-Data, SQL-Like Query Processing

Qing Zheng, Scientist, Los Alamos National Laboratory

8/17/2023

LA-UR-23-29456



Managed by Triad National Security, LLC, for the U.S. Department of Energy's NNSA.

Overview

Problem

Scientific analytics increasingly bottlenecked on large data transfers

Goal

Evaluate how/how much in-drive analytics can help

C2

Our early prototype for assessing gains

Results

Sizeable speedups even when data transfer is not the primary bottleneck

Background: Scientific Datasets

Resemble tables with rows and columns

- Rows: records
- Columns: attributes

Traditional HPC data formats: HDF5, NetCDF (self describing, offset-based query interface)

We are also looking at leveraging industrial data formats (such as Apache Parquet, ORC, Avro) and analytics stacks to enable richer query types beyond offsets (e.g.: SQL)

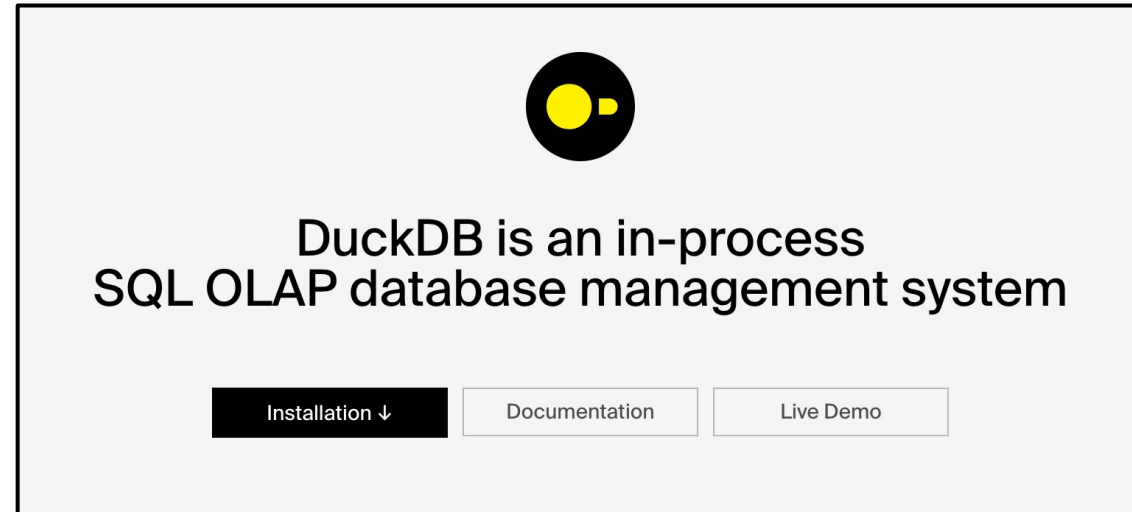
An Example

Store data in Parquet

- Self describing
- Columnar data model
 - Row Groups
 - Column Chunks

Run Queries using DuckDB

- Supports SQL
- Understands Parquet



```
SELECT * FROM 'test.parquet'  
WHERE X>Y
```

HPC Simulation Sciences Workflow

Simulation Phase

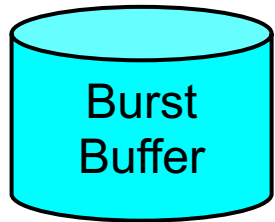
- User submits jobs
- Jobs run on compute nodes
- Jobs generate data (e.g.: in Parquet)
- Data is written to backend storage
- Storage likely tiered

Analytics Phase

- User runs queries against their data (e.g.: DuckDB)
- **A query may select only a tiny amount of data from a large dataset**
- But the reader program may still have to read the entire dataset from storage nodes

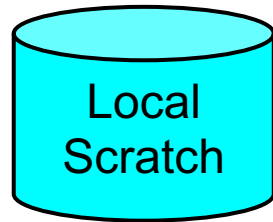
Can we return only data that is selected by a query?

Time to Read Back 1PB of Data



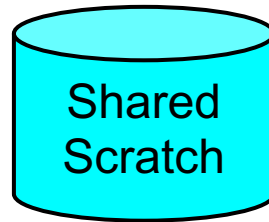
3.2TB/s

312s



1.2TB/s

14min



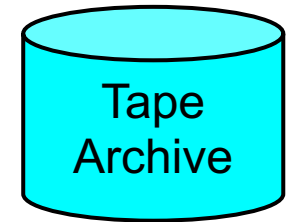
300GB/s

56min



100GB/s

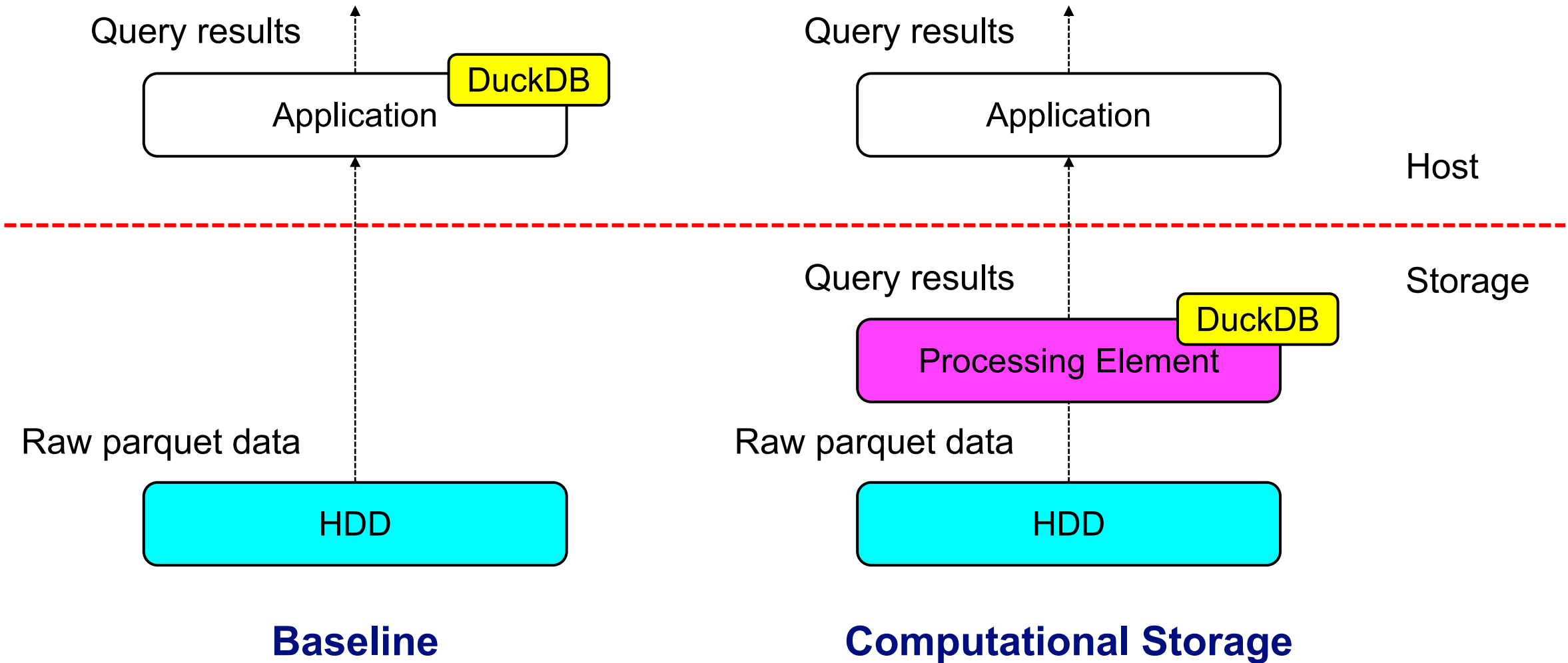
2.8hr



10GB/s

28hr

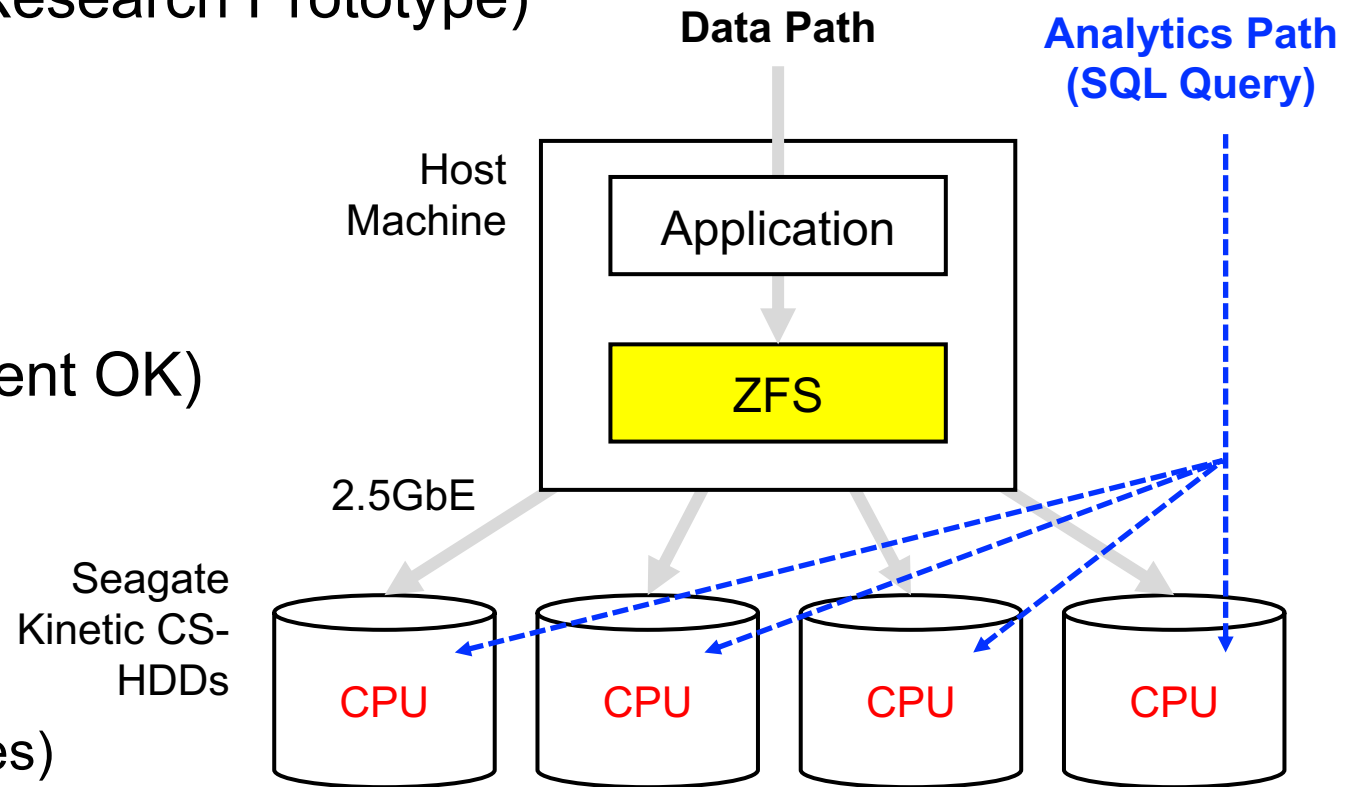
Why Computational Storage Might Help



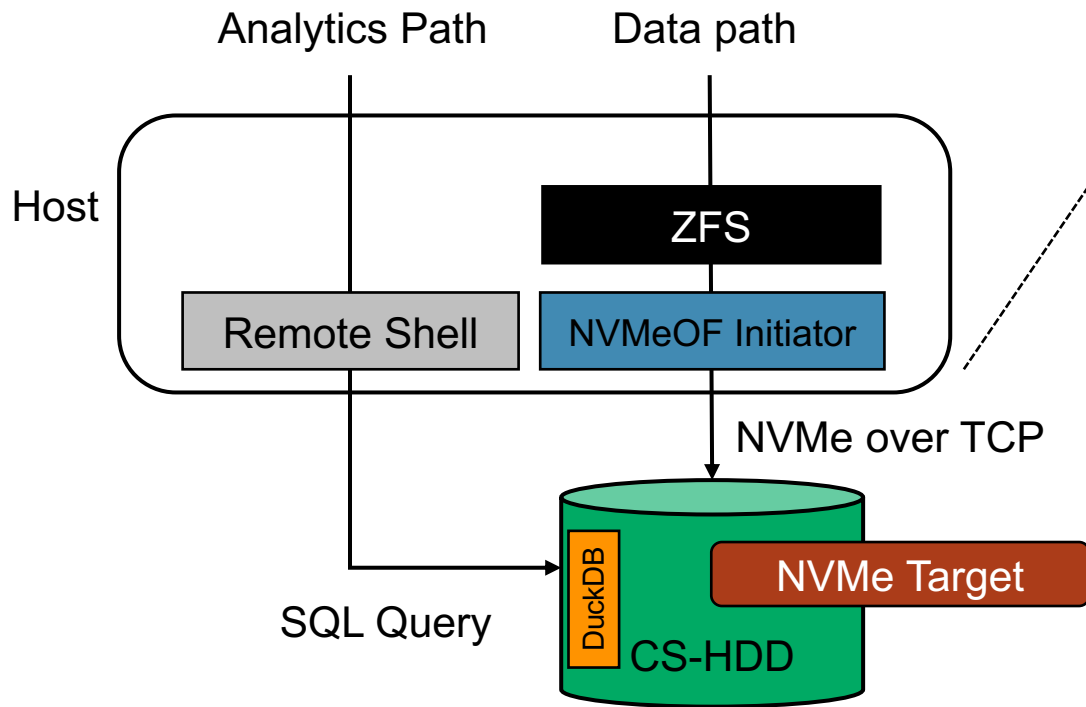
C2: In-Drive Query Processing

Our first near-data analytics prototype for cool storage tiers

- **Disk:** Kinetic CS-HDDs (Seagate's Research Prototype)
 - CPU: 2x ARM Cortex-A53 cores
 - RAM: 1GB
 - OS: Ubuntu Linux (C++ development OK)
 - Network: 2x 2.5GbE
- **Host Filesystem:** ZFS
 - Data protection: **RAID** (1, 2, or 3 parities)



A Close Look



Two Challenges

Drives have no knowledge of FS file-to-block mapping

- Solution: LibZDB (allow querying ZFS for mapping information)

A data row may be split over multiple drives

- Data alignment control

Evaluation

3 Scenarios

A) Host network is a bottleneck

- Can in-drive analytics improve performance?

B) Host CPU is a bottleneck

- Can in-drive analytics improve performance?

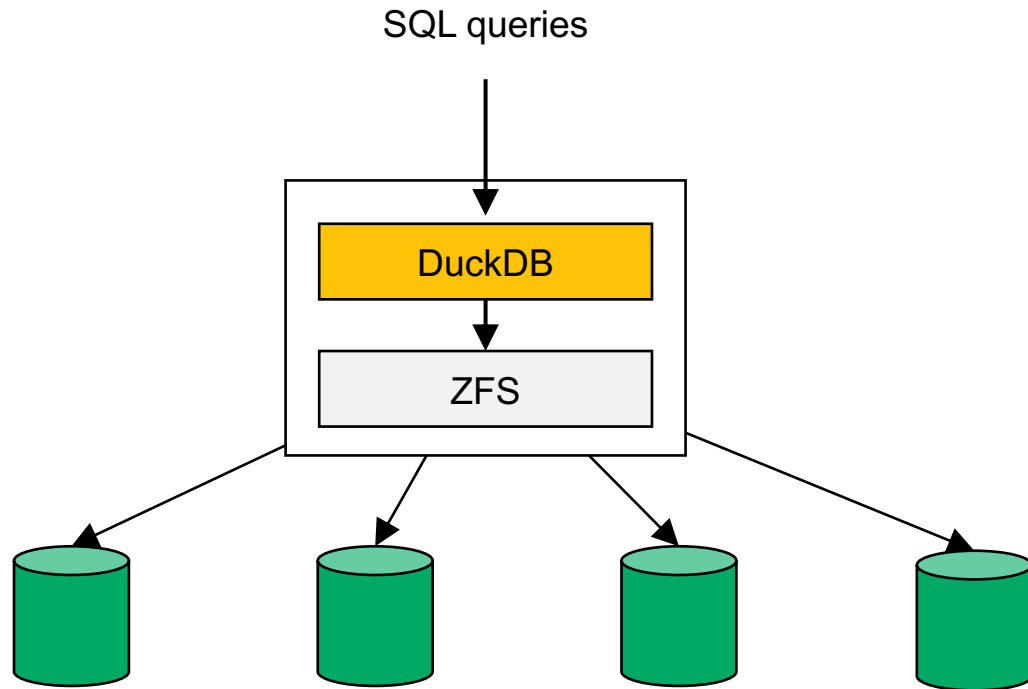
C) Host has abundant CPU & network

- Can in-drive analytics continue to improve performance?

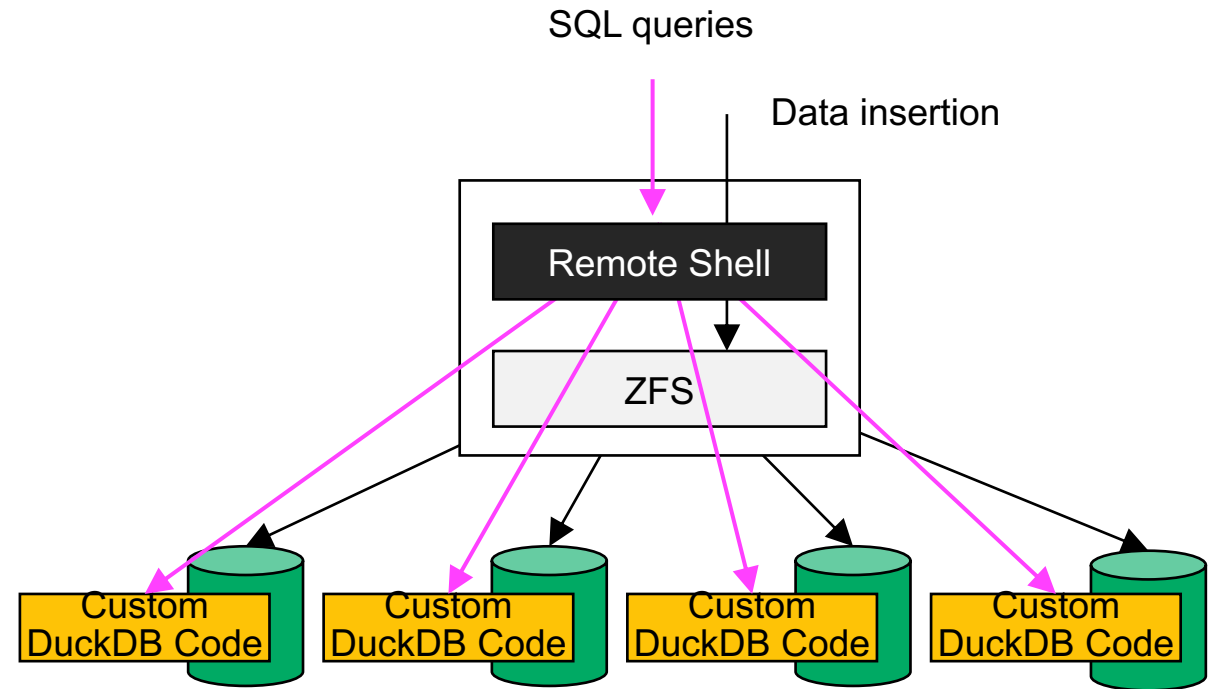
Experiment Setup

- 1 ZFS host (32 AMD CPU cores)
- **38 CS-HDDs**
 - 2x 16+3 **RAID** Pools
- **50GB** dataset from a real particle simulation
 - 2 billion rows (in Parquet fmt)
 - Columns: ID, x, y, z, ke
- **DuckDB queries**
 - SELECT * WHERE ke>X

Baseline vs. C2 Kinetic Runs



Baseline

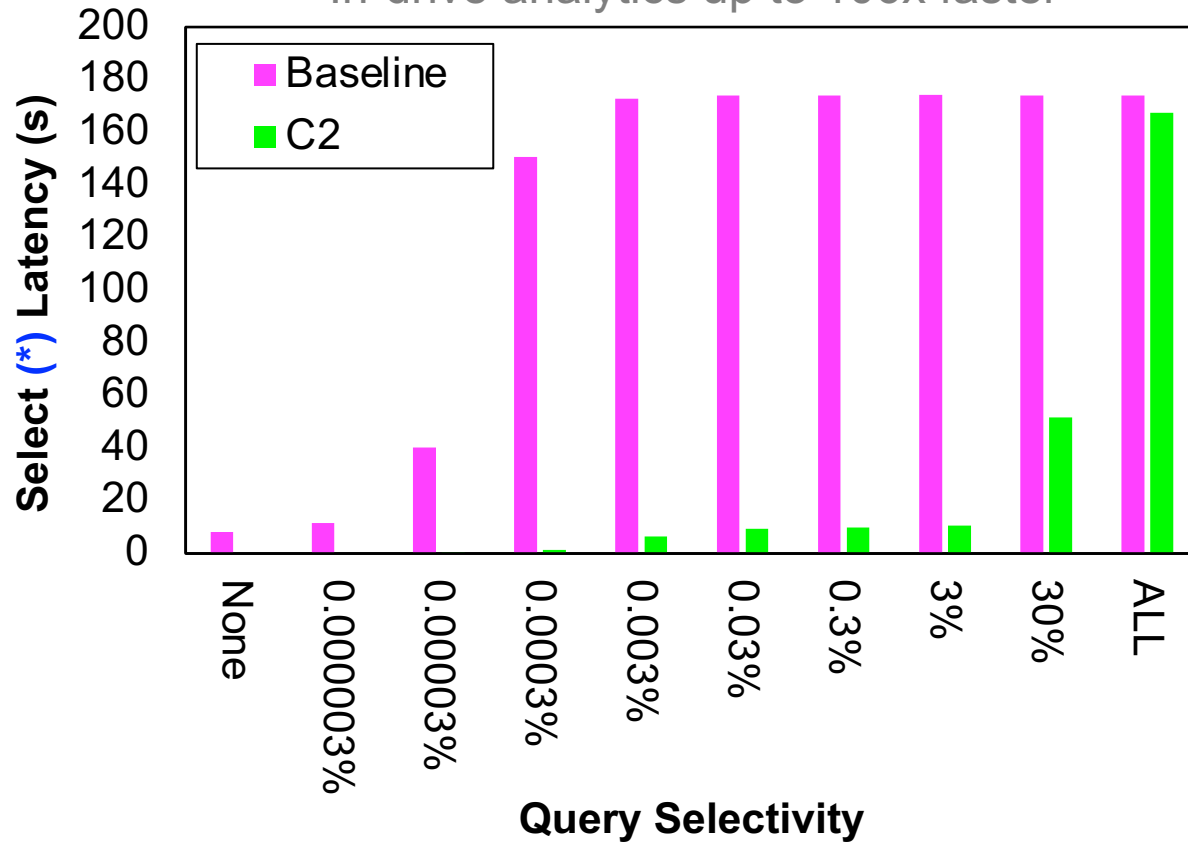


C2

Result

Case #1: Host network was the bottleneck

In-drive analytics up to 106x faster

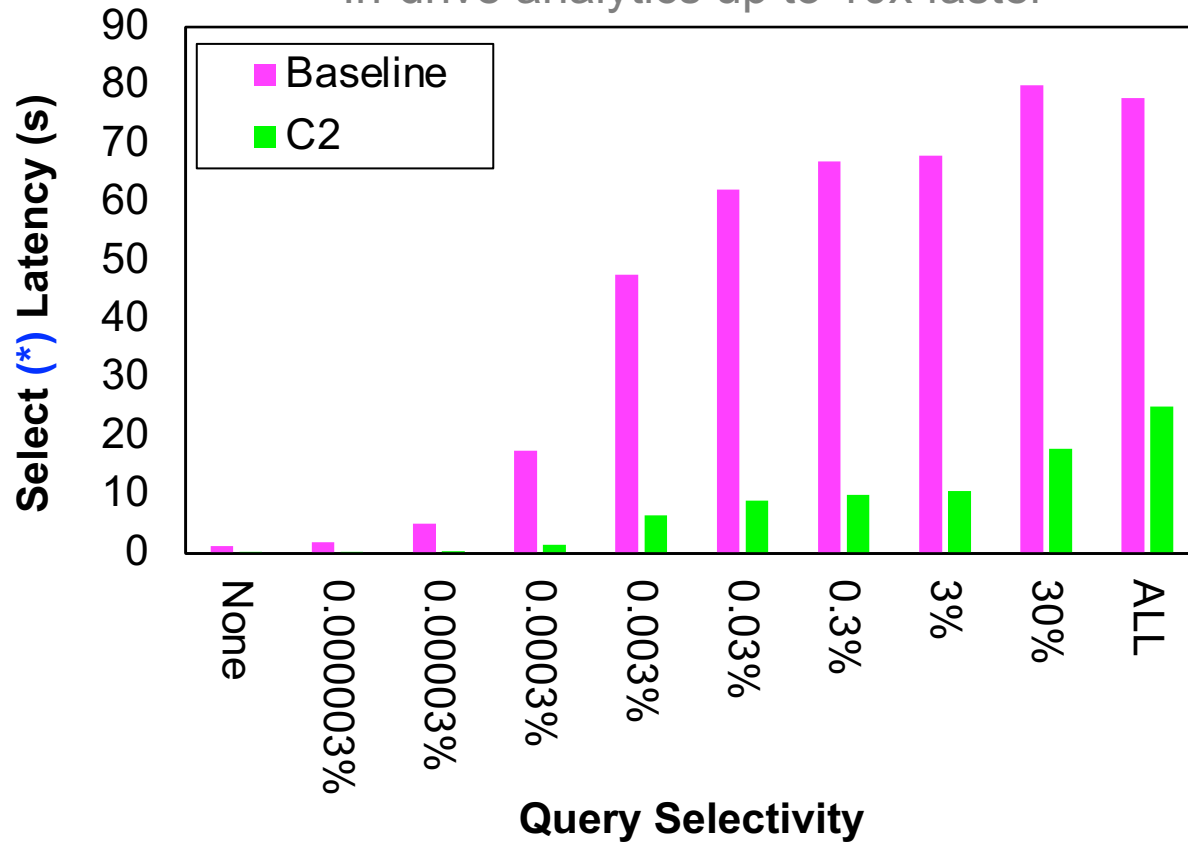


In-drive analytics allow sending less data over the network

Result

Case #2: Host CPU was the bottleneck

In-drive analytics up to 10x faster

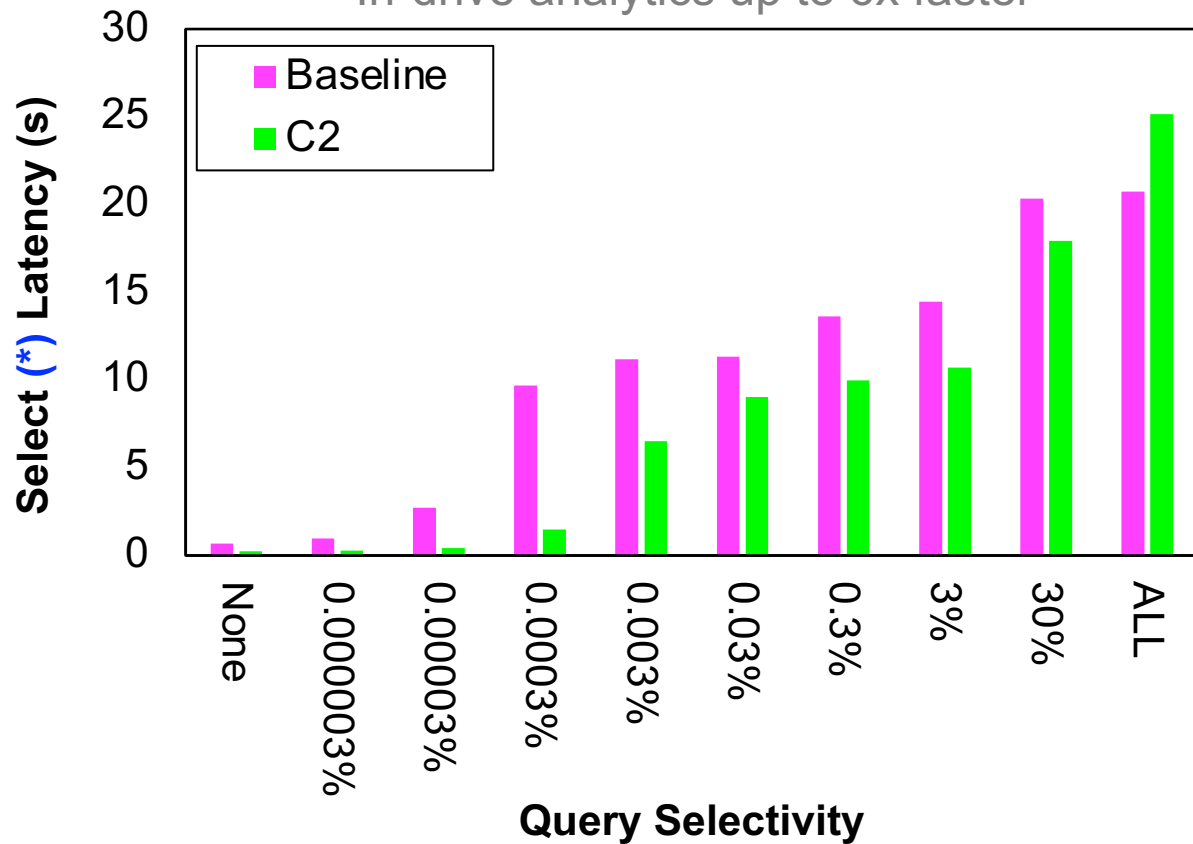


In-drive analytics allow massively parallel computing across drives

Result

Case #3: Host had abundant network & CPU

In-drive analytics up to 6x faster



In-drive analytics allow more fully utilizing disk bandwidth

Conclusion

Computational storage provides new ways of accelerating data-intensive applications

In-drive data management schemes matter (O_DIRECT, clustered index)

Layer violation: “cheating” one filesystem may be possible; cheating multiple layers of filesystems is hard (FS internal load balancing, fail over, compression, concurrency control)

Future directions: Block-based acceleration to object-based acceleration

Acknowledgement

Jason Lee (jasonlee@lanl.gov)

Brian Atkinson (batkinson@lanl.gov)

Jarrett Crews (jarrett@lanl.gov)

David Bonnie (dbonnie@lanl.gov)

Dominic Manno (dmanno@lanl.gov)

Gary Grider (ggrider@lanl.gov)

Philip Kufeldt (philip.kufeldt@seagate.com)

Evan Burgess
(evan.burgess@seagate.com)

Ivan Rodriguez
(ivan.rodriguez@seagate.com)

David Allen (david.j.allen@seagate.com)

John Bent (john.bent@seagate.com)

Bradley Settlemyer
(bsettlemyer@nvidia.com)

Thank you!