# Predicting and optimizing the performance of HDF5 applications

Donghe Kang

kang.1002@osu.edu
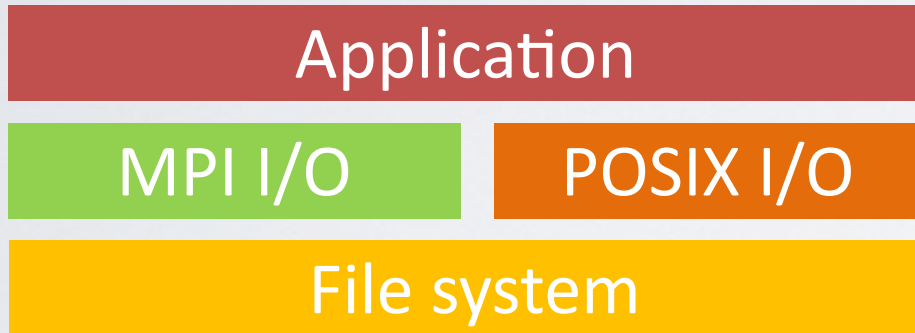
THE OHIO STATE
UNIVERSITY

# Software stack
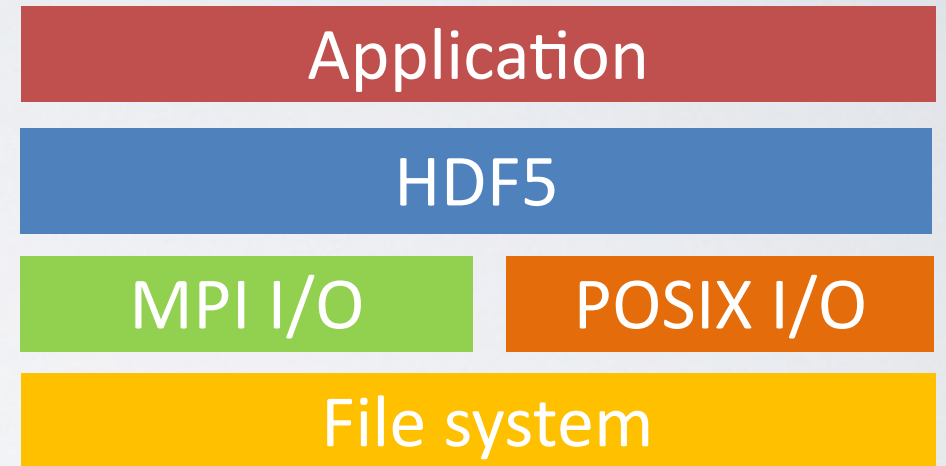
- **Previous modeling work**



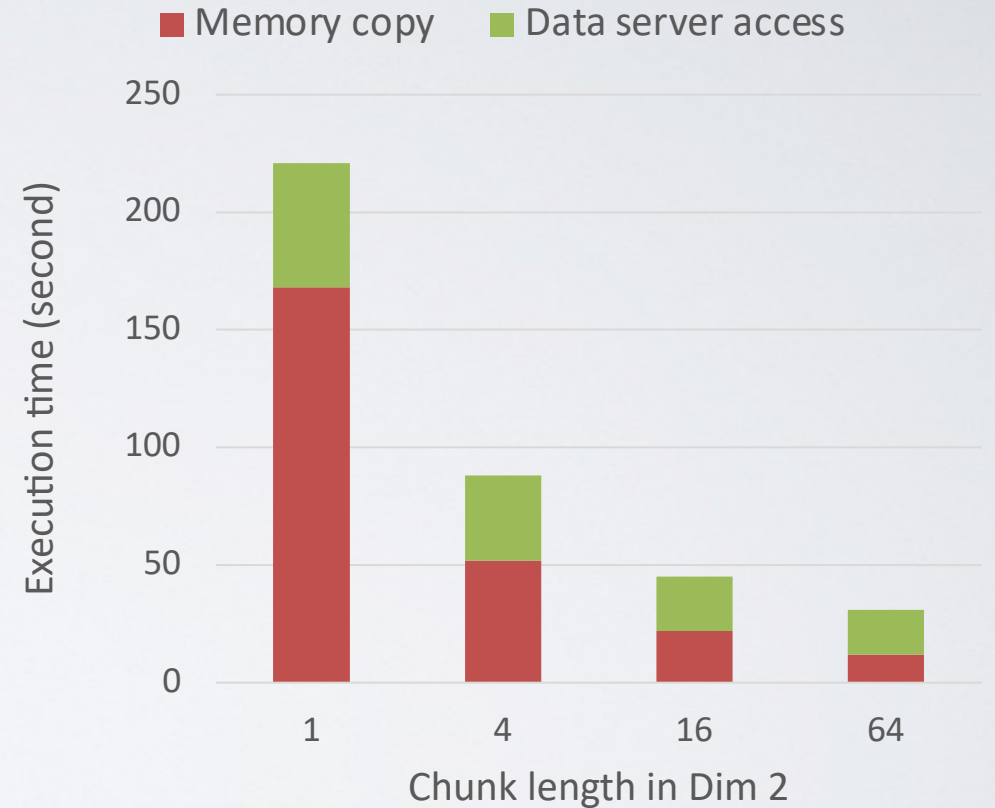Applications directly access data. Previous works build models to predict the I/O time.

- **This work**



How to predict the end-to-end performance of the array I/O requests?

# Motivation

- ## Read an entire array



Dimension 2

64

Dimension 1

36M

Load

Chunked on-disk storage layout

In-memory serial buffer



Memory copy     Data server access

Execution time (second)

250

200

150

100

50

0

1          4          16          64

Chunk length in Dim 2

3

# Read path in HDF5

```
┌─────────────────────────────────┐
│  Selected region and data buffer│
└─────────────────────────────────┘
               │
┌─────────────────────────────────┐
│   Locate chunks in the region   │
└─────────────────────────────────┘
               │
         ◇ Next chunk ◇ ──N──► ┌──────────┐
               │ Y             │ Complete │
               │               └──────────┘
         ◇ Hit in the ◇ ──Y──┐
         ◇ library cache ◇    │
               │ N            │
         ◇ Hit in the file ◇ ──Y──► ┌──────────────────┐
         ◇ system cache ◇            │  Read from file  │
               │ N                   │  system cache    │
                                     └──────────────────┘
         ┌──────────────────────┐
         │ Read from data servers│
         └──────────────────────┘
               │
┌─────────────────────────────────┐
│  Memcpy data from the library   │
│     cache to data buffer        │
└─────────────────────────────────┘
```

The cache read and memory copy operations can spend up to 90% and 87% of the execution time in our experiments.
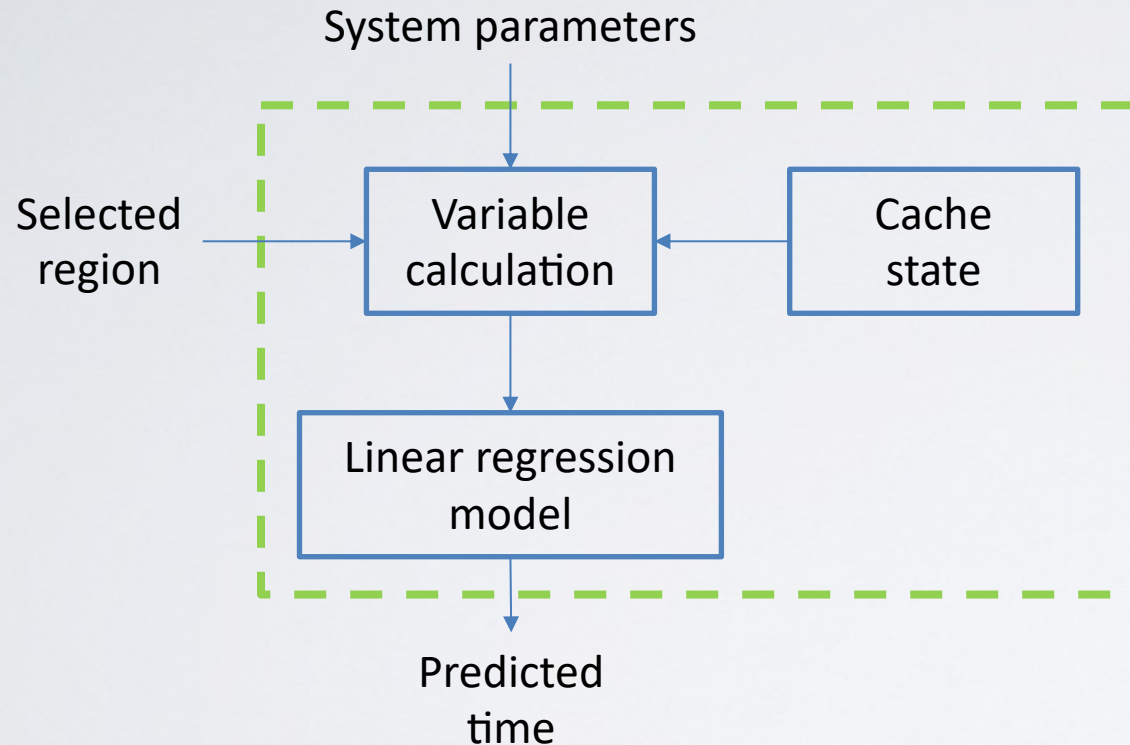
An end-to-end model should cover all the steps in the I/O path.

End-to-end time =

> cache read time +

> data server access time +

> memory copy time

4

# Solution overview

System parameters

Selected region → Variable calculation ← Cache state

Variable calculation → Linear regression model → Predicted time
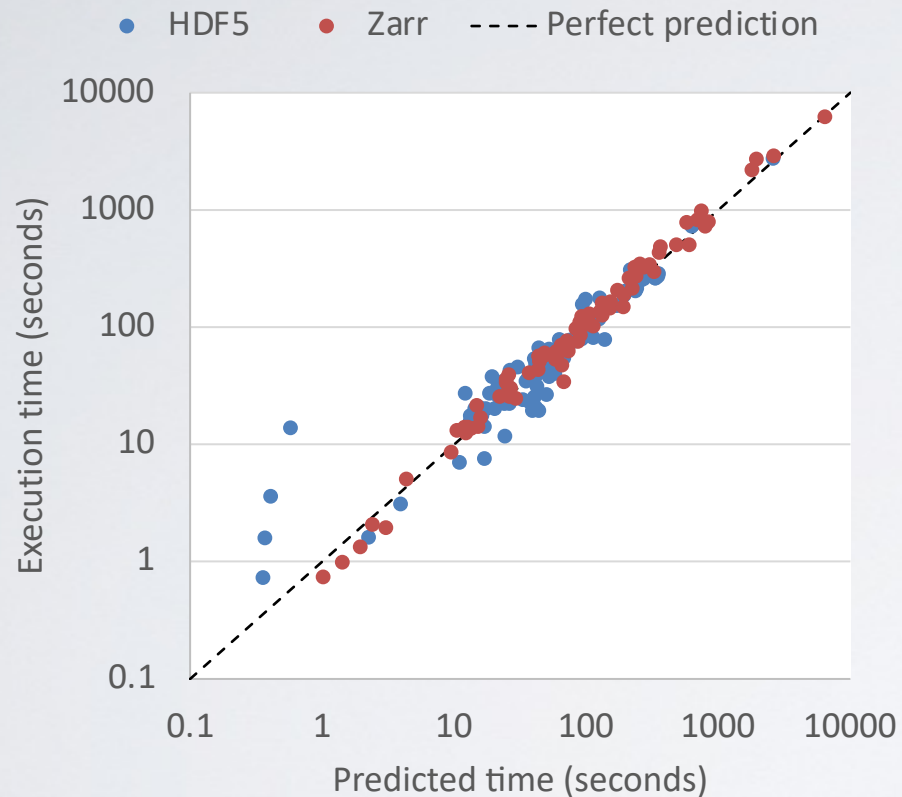
- The system parameters
  - the chunk shape
  - storage format (HDF5 or Zarr)
  - stripe count and size in the file system
  - cache size

- The solution contains three components
  - The cache state component maintains the status of the library cache and the file system cache
  - The variable calculation component computes the variables in the model based on the selected region
  - The model component predicts the cache read time, data server access time and memory copy time

# Model accuracy

Real and predicted time

## Experiment variables

| Variable | Value range |
|---|---|
| Array shape | $36M \times 64, 2.4K \times 2K \times 3K$ |
| Selected region | entire array, 1K rows, 1 column, 1K boxes |
| Chunk length in Dim 1 | 70K - 4.5M |
| Chunk length in Dim 2 | 1 - 64 |
| Number of processes | 1 - 64 |
| Stripe count | 1 - 64 |
| File system | Lustre, GPFS |

- Takeaways:
  - The RMSE of the model is 0.29
  - The model correctly predicts the fastest library between HDF5 and Zarr 94% of the time.

# Small array challenge

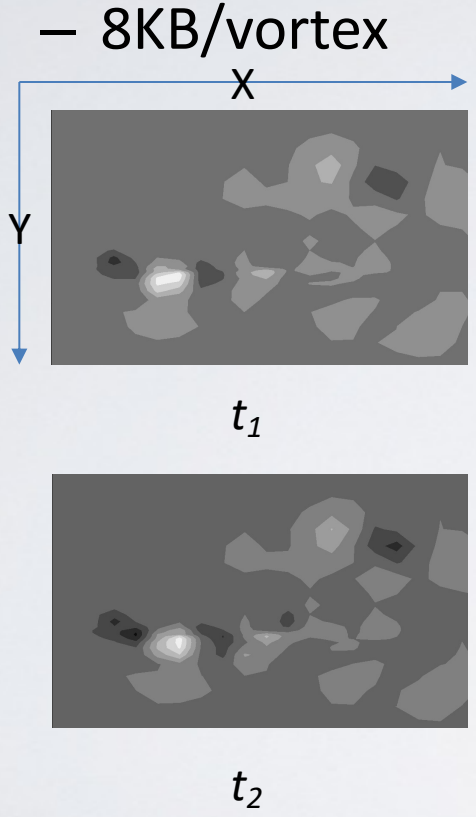- Supernova detection



Segment
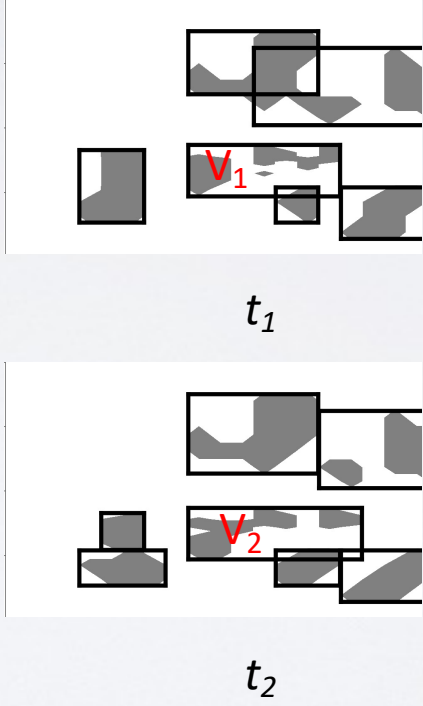
CNN

Real or bogus?

21✕21 pixels/image

# Small array challenge

- Vortices prediction
  - 8KB/vortex



X

Y

$t_1$

$t_2$

Identify vortices →

$V_1$

$t_1$

$V_2$

$t_2$
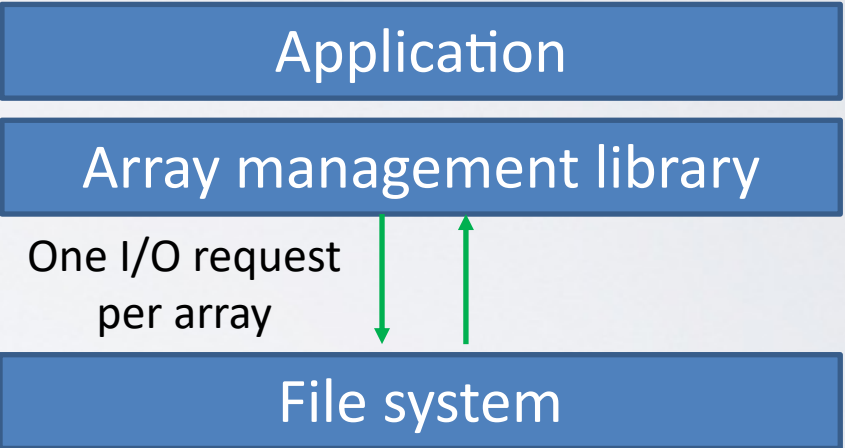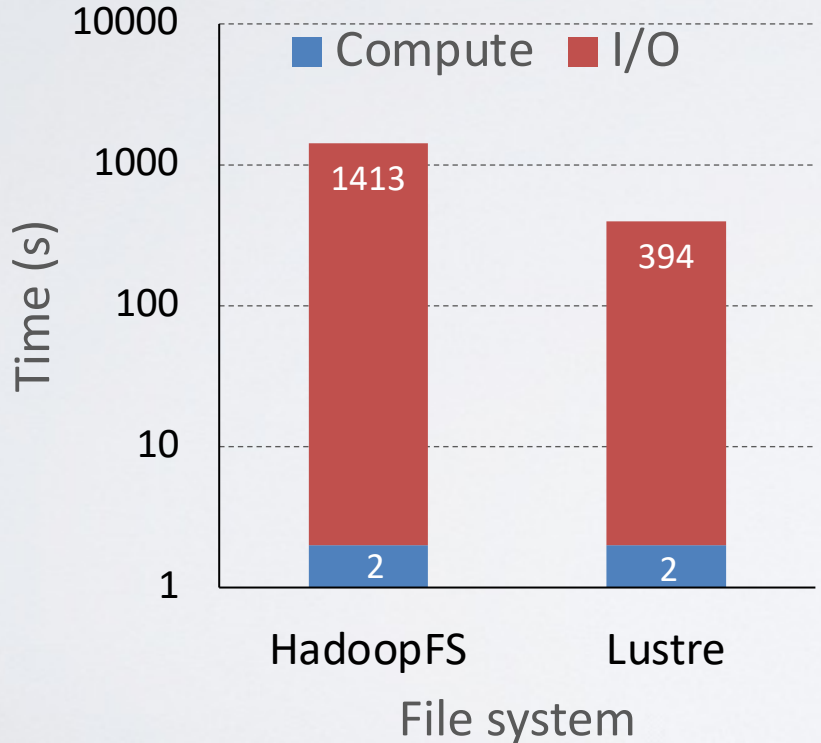
3D query →

$V_1 = V_2$

# Small array challenge

- One I/O per array
  - I/O takes 200✕ to 700✕ longer than computation



Bar chart — Time (s) vs File system (HadoopFS, Lustre). Legend: Compute (blue), I/O (red). HadoopFS: I/O 1413, Compute 2. Lustre: I/O 394, Compute 2.



Application
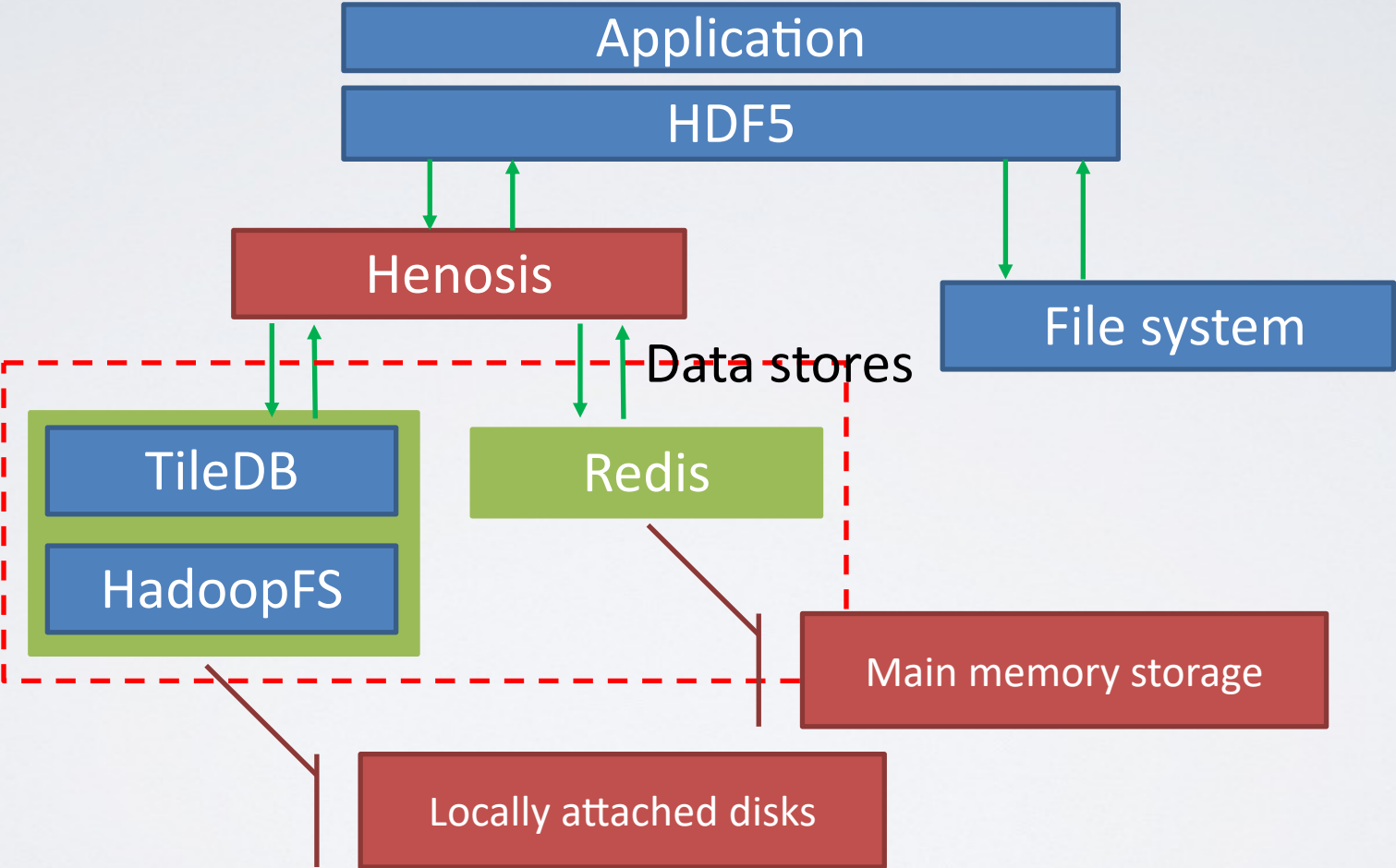Array management library
One I/O request per array
File system

9

# Small array challenge

Too many small I/O requests!
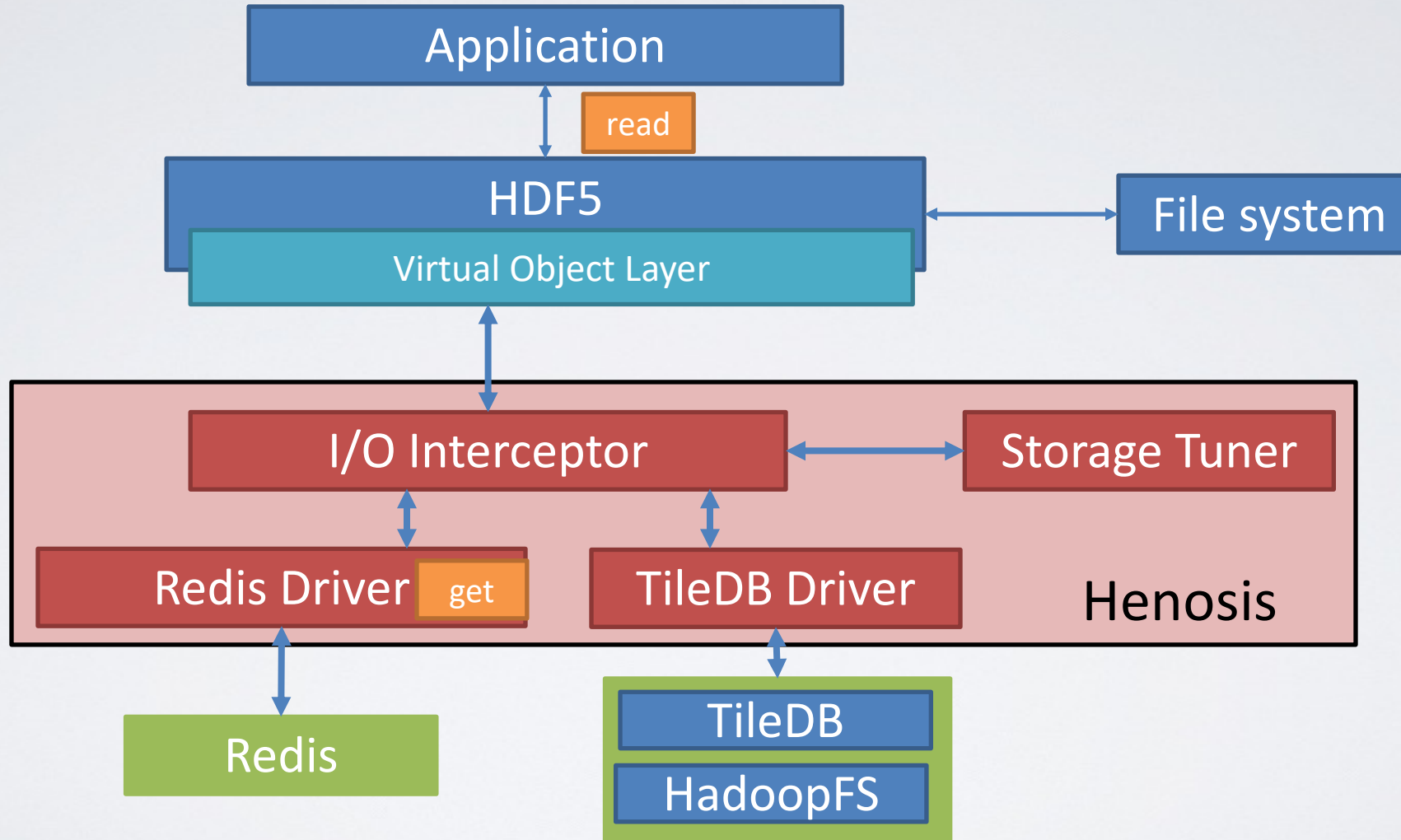
# HDF5 on heterogeneous data stores

# Goals

- Store arrays on heterogeneous data stores
  - Without modifying applications
- Accelerate small I/O requests
  - Placement → improve the performance of one request
  - Consolidation → reduce the number of requests
- Automatically decide the array storage layout
  - Which data store should an array be placed in?
  - How do we store small arrays in chunks?

# Goals

- **Store arrays on heterogeneous data stores**
  - **Without modifying applications**
- Accelerate small I/O requests
  - Placement → improve the performance of one request
  - Consolidation → reduce the number of requests
- Automatically decide the array storage layout
  - Which data store should an array be placed in?
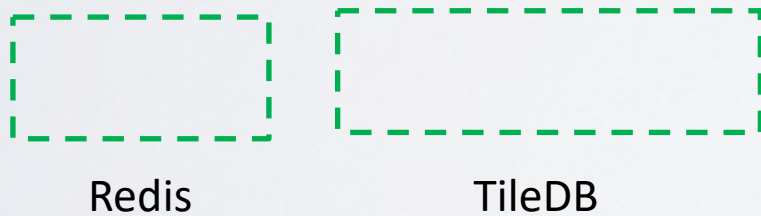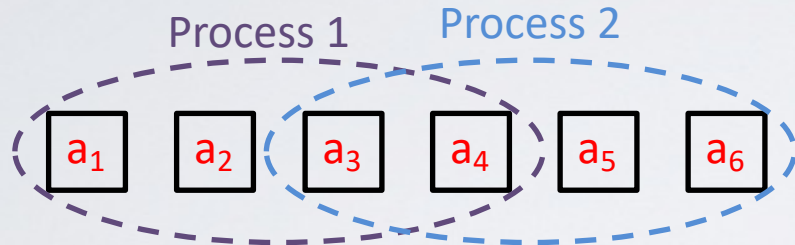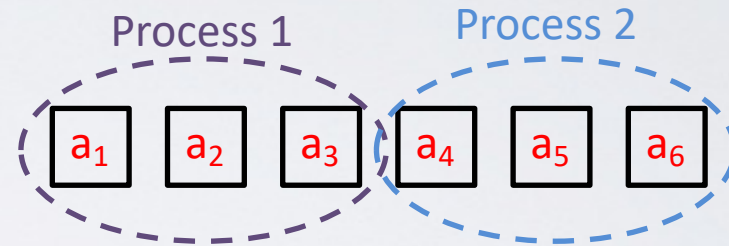  - How do we store small arrays in chunks?

# System architecture

# Goals

- Store arrays on heterogeneous data stores
  - Without modifying applications
- **Accelerate small I/O requests**
  - **Placement → improve the performance of one request**
  - **Consolidation → reduce the number of requests**
- Automatically decide the array storage layout
  - Which data store should an array be placed in?
  - How do we store small arrays in chunks?

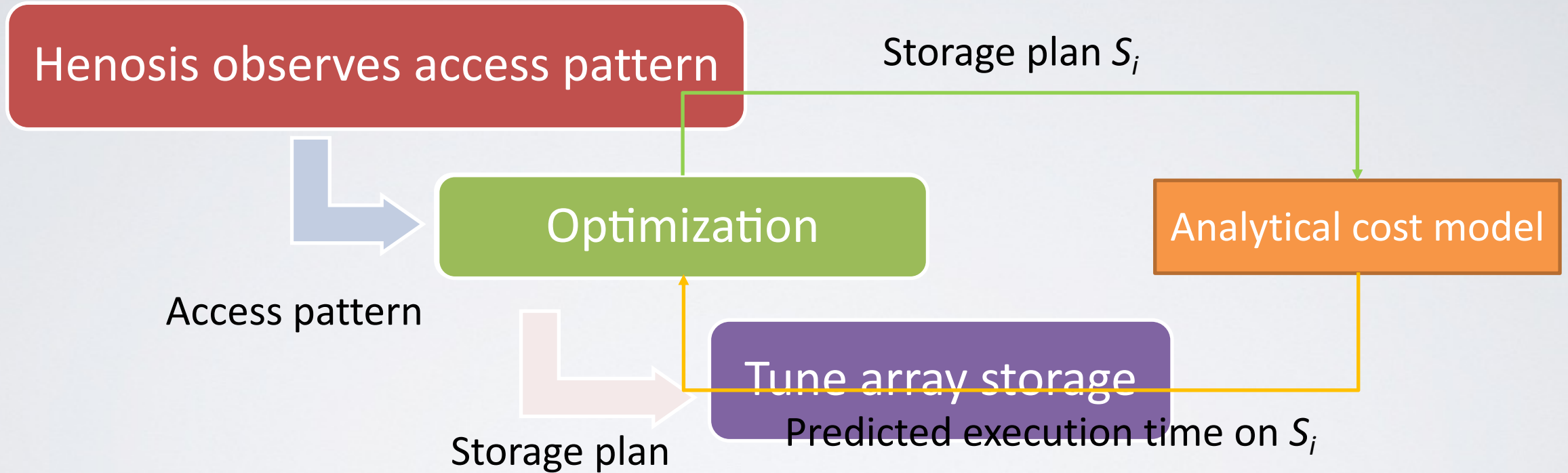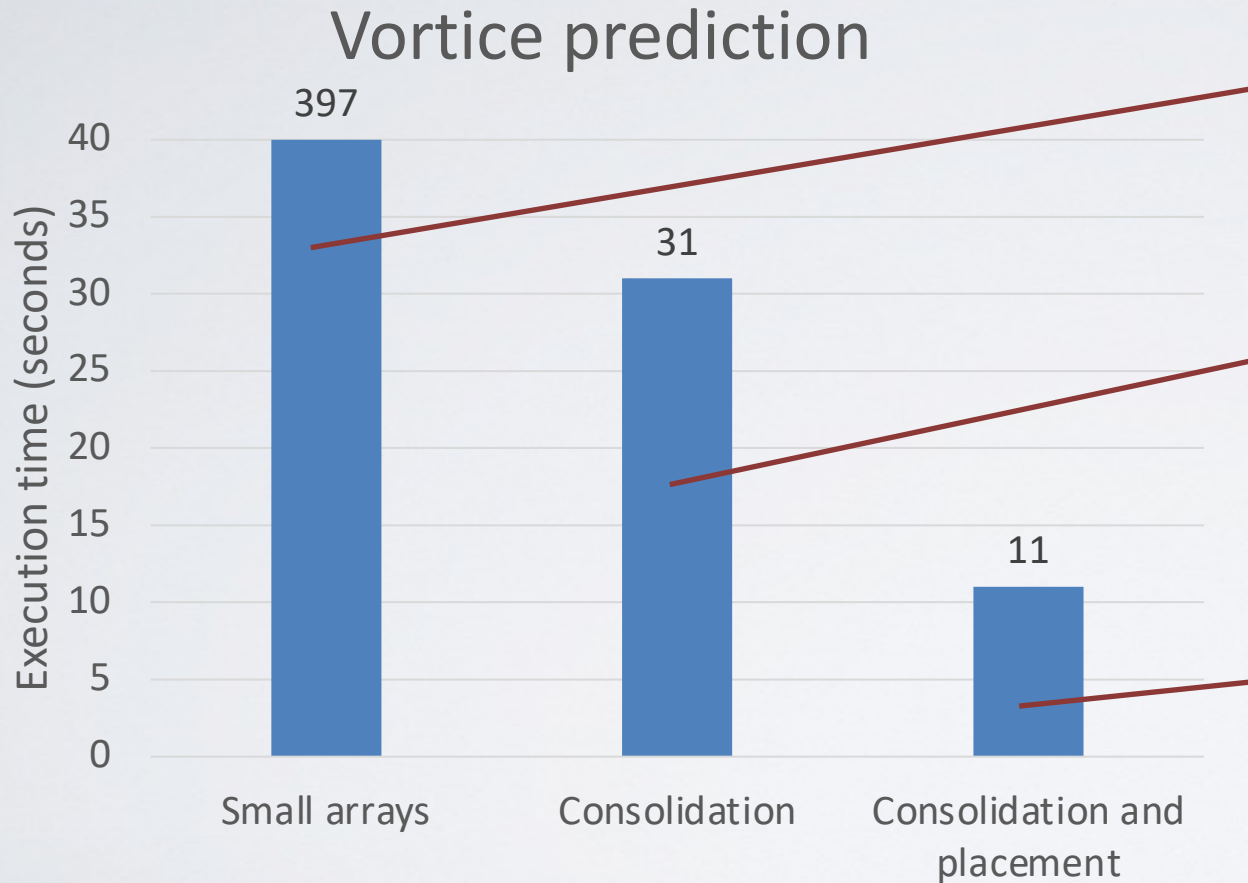# I/O acceleration techniques

- Placement

- Consolidation



Process 1  Process 2

$a_1$ $a_2$ $a_3$ $a_4$ $a_5$ $a_6$

Redis    TileDB

Process 1  Process 2

$a_1$ $a_2$ $a_3$ $a_4$ $a_5$ $a_6$

Chunk 1  Chunk 2

16

# Goals

- Store arrays on heterogeneous data stores
  - Without modifying applications
- Accelerate small I/O requests
  - Placement $\rightarrow$ improve the performance of one request
  - Consolidation $\rightarrow$ reduce the number of requests
- **Automatically decide the array storage layout**
  - **Which data store should an array be placed in?**
  - **How do we store small arrays in chunks?**

# Optimization workflow

Henosis observes access pattern

Storage plan $S_i$

Optimization

Analytical cost model

Access pattern

Tune array storage

Storage plan

Predicted execution time on $S_i$

# Optimization impact

## Vortice prediction

Execution time (seconds)

397

31

11

Small arrays   Consolidation   Consolidation and placement

Directly read many small arrays from HadoopFS

Consolidate small arrays with same access pattern in a chunk. Read in fewer I/O requests.

Place hot arrays and co-accessed arrays in Redis. Read less data from HadoopFS.

# Conclusions

- Reading from cache and memory copy to transform layouts can spend 90% of the execution time
  - An end-to-end model to cover the entire I/O path

- Applications spend 99% of the time to read small arrays
  - VOL transparently forward requests to two data stores
  - Placement and consolidation reduce the number of I/O requests to the slow data store
  - An analytical cost model helps to decide the storage layout

Personal webpage: https://web.cse.ohio-state.edu/~kang.1002/