



Hewlett Packard
Enterprise

MLPerf Workload I/O Patterns: Is Hybrid I/O the Right Choice?

Rajeev Mishra

April, 2025

Agenda

What is MLPerf Storage

Workload Options in MLPerf

Hybrid I/O

Workload - data loaders

I/O Profile Deep Dive

Workload Performance: Inference



MLPerf Storage



1. Measures the performance of a storage system for AI Training



2. Supports A100 and H100



3. Data loader used to read the data



Workload and data loader

1. Unet3D/Pytorch

Primarily used for medical image segmentation and 3D object detection.

2. Resnet50/Tensorflow

Widely used for image classification and feature extraction tasks.

3. Cosmoflow/Tensorflow

Analyzing large-scale structure data and simulations



Hybrid I/O

- Combines the functionalities of buffered I/O and direct I/O.
- Following tunes which I/O to perform
 - `hybrid_io_read_threshold_bytes`
 - `hybrid_io_write_threshold_bytes`
- **Preliminary Assumptions**
 - Random reads
 - Application caching
 - Cache Thrashing
- GDS Direct I/O



I/O Profile

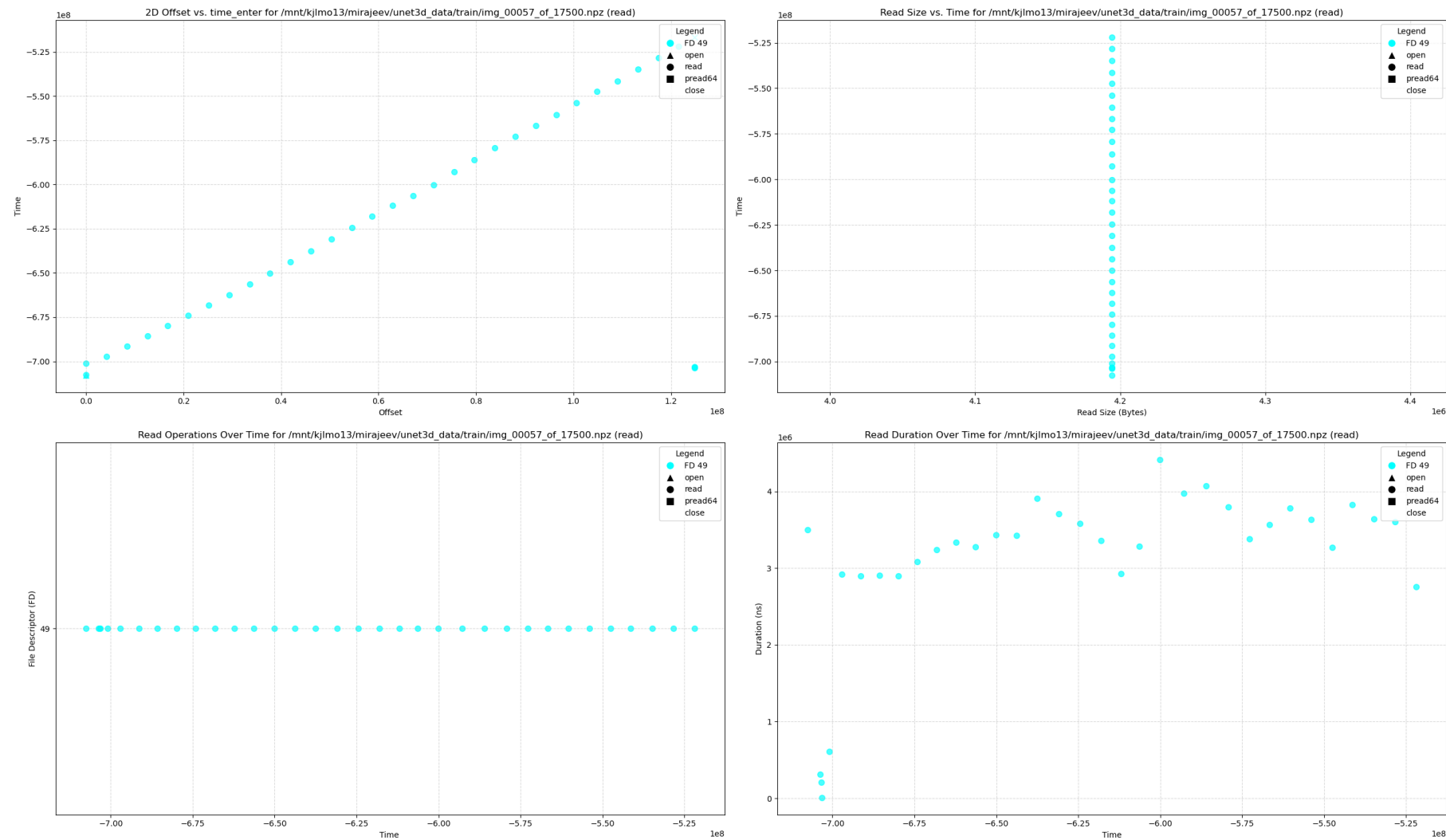


Unet3D

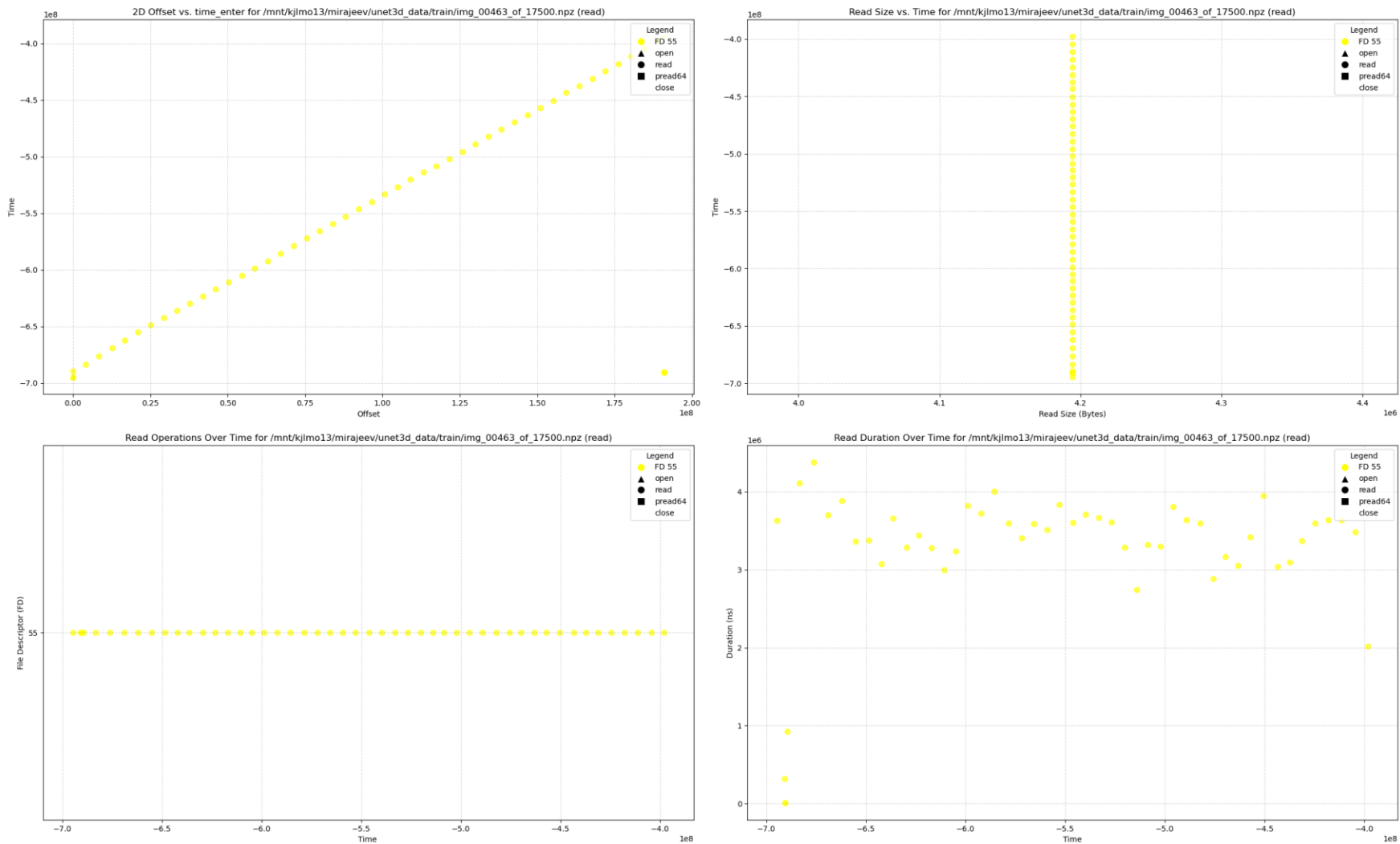
Bpf trace profile



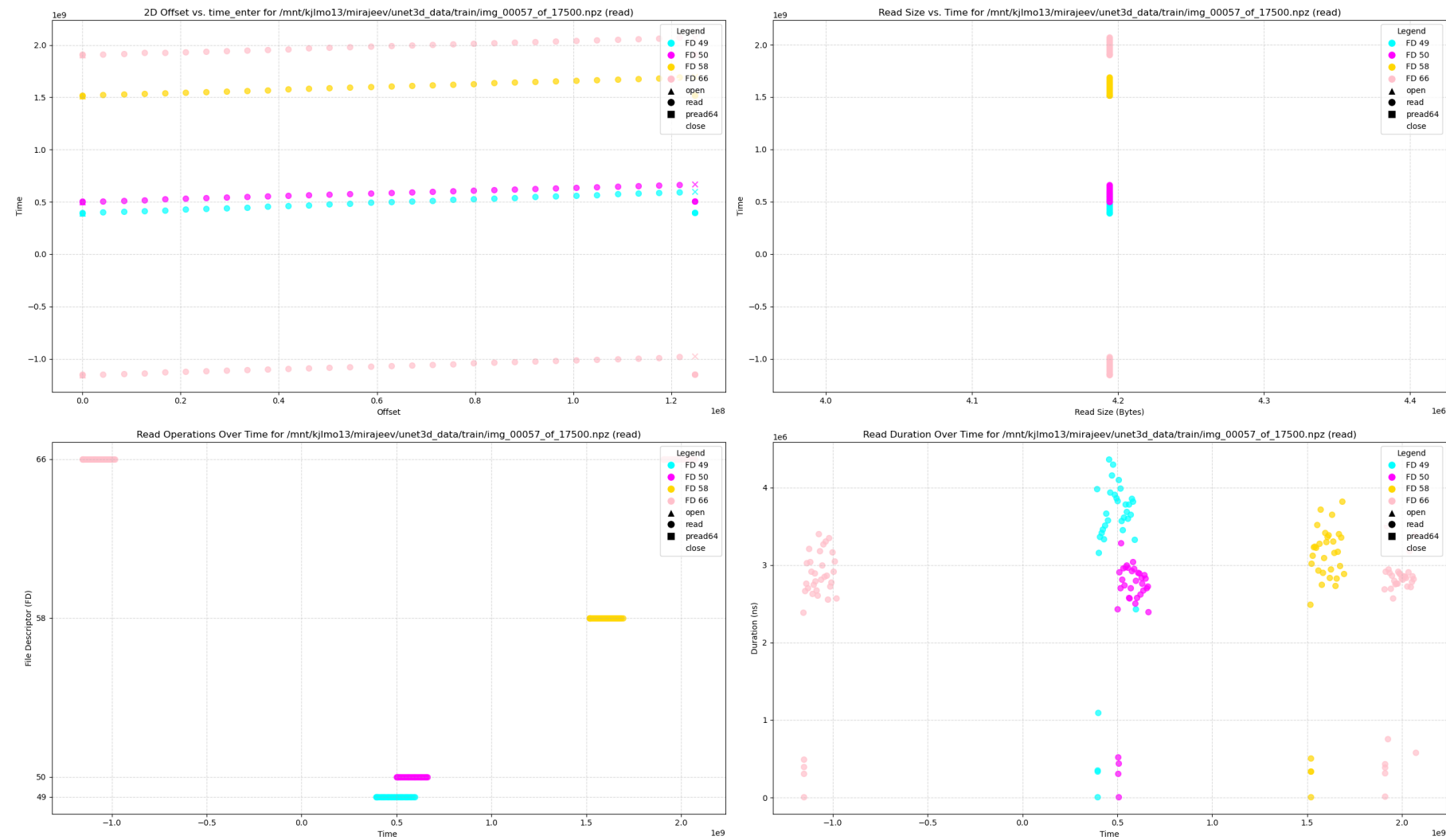
Unet3D Epoch 1 plots



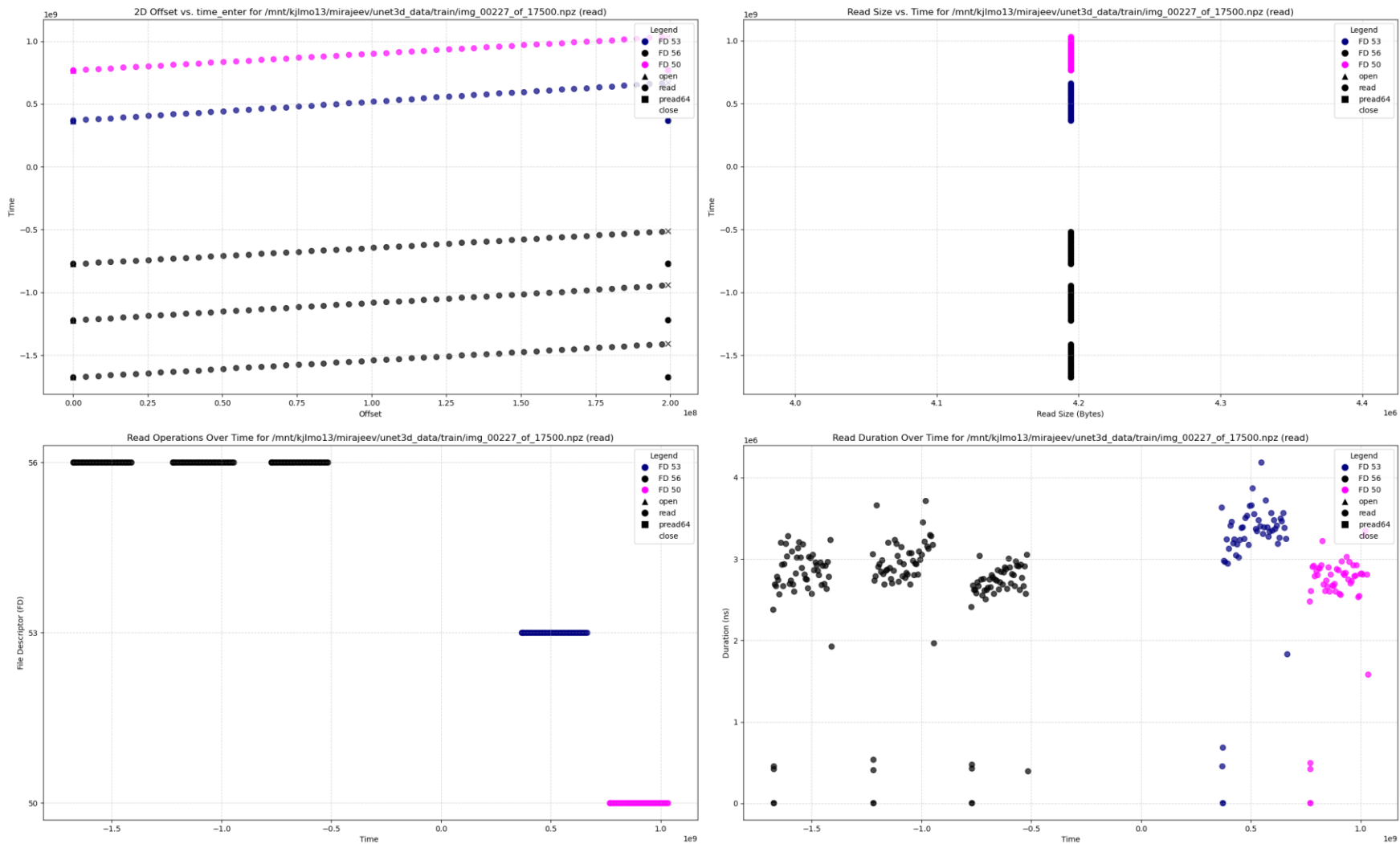
Unet3D Epoch 1 continues



Unet3D Epoch 5 plots



Unet3D Epoch 5 continue



Darshan

Bpf trace profile



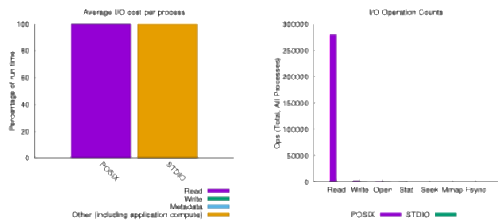
ResNet Darshan

python3 (3/21/2025)

1 of 3

jobid: 1960735 uid: 121685592 nprocs: 1 runtime: 380.1681 seconds

I/O performance estimate (at the POSIX layer): transferred 69821.1 MiB at 128.33 MiB/s
I/O performance estimate (at the STDIO layer): transferred 0.0 MiB at 25.51 MiB/s



Most Common Access Sizes (POSIX or MPI-IO)			File Count Summary (estimated by POSIX I/O access offsets)			
	access size	count	type	number of files	avg. size	max size
POSIX	262144	279040	total opened	515	136MiB	137MiB
	178	822	read-only files	512	137MiB	137MiB
	124573	512	write-only files	1	282KiB	282KiB
			read/write files	0	0	0
	179	435	created files	1	282KiB	282KiB

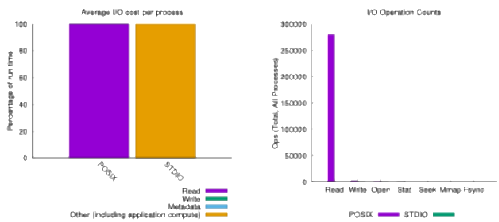
```
python3 dlio benchmark/dlio benchmark/main.py --config-path=/home/mirajeev/storage/storage-conf workload=resnet50 h100
++ workload.workflow.generate.data=False ++ workload.workflow.train=True
++ workload.dataset.data.folder=/mnt/kgimo13/mirajeev/resnet50_data ++ workload.workflow.profling=False
++ workload.profling.profler=None ++ hydra.output.subdir=configs ++ hydra.run.dir=/mnt/kgimo13/mirajeev/resnet50_h100
```

python3 (3/21/2025)

1 of 3

jobid: 1960736 uid: 121685592 nprocs: 1 runtime: 380.1677 seconds

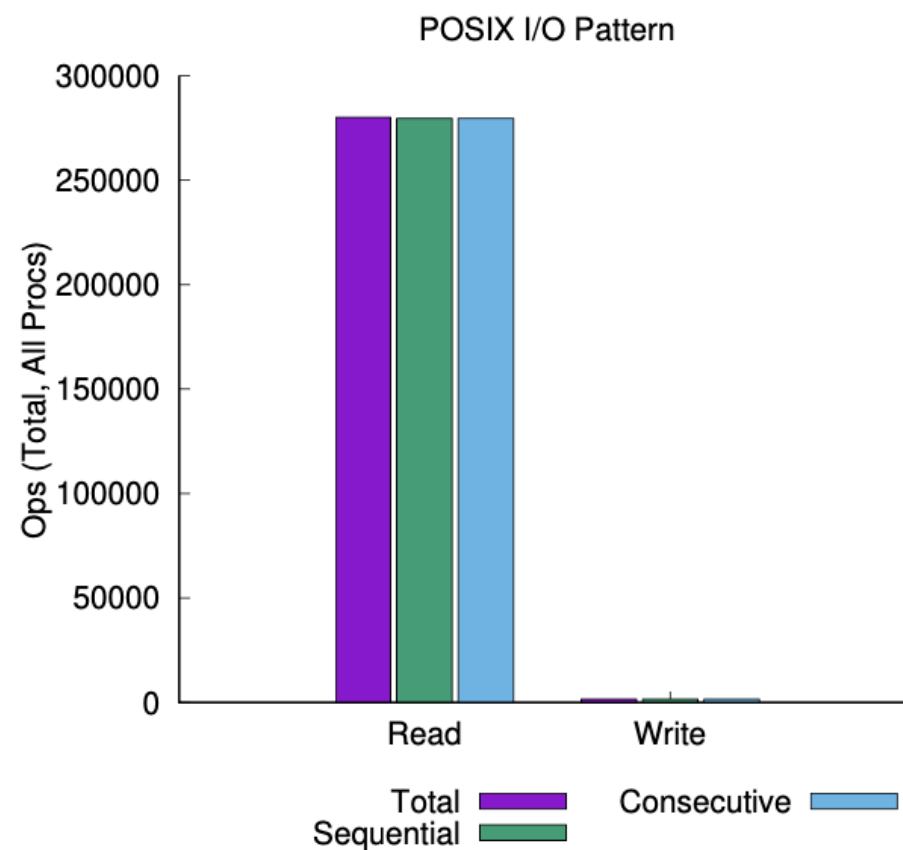
I/O performance estimate (at the POSIX layer): transferred 69821.1 MiB at 110.94 MiB/s
I/O performance estimate (at the STDIO layer): transferred 0.0 MiB at 27.94 MiB/s



Most Common Access Sizes (POSIX or MPI-IO)			File Count Summary (estimated by POSIX I/O access offsets)			
	access size	count	type	number of files	avg. size	max size
POSIX	262144	279040	total opened	515	136MiB	137MiB
	178	802	read-only files	511	137MiB	137MiB
	124573	512	write-only files	1	279KiB	279KiB
			read/write files	0	0	0
	179	449	created files	1	279KiB	279KiB

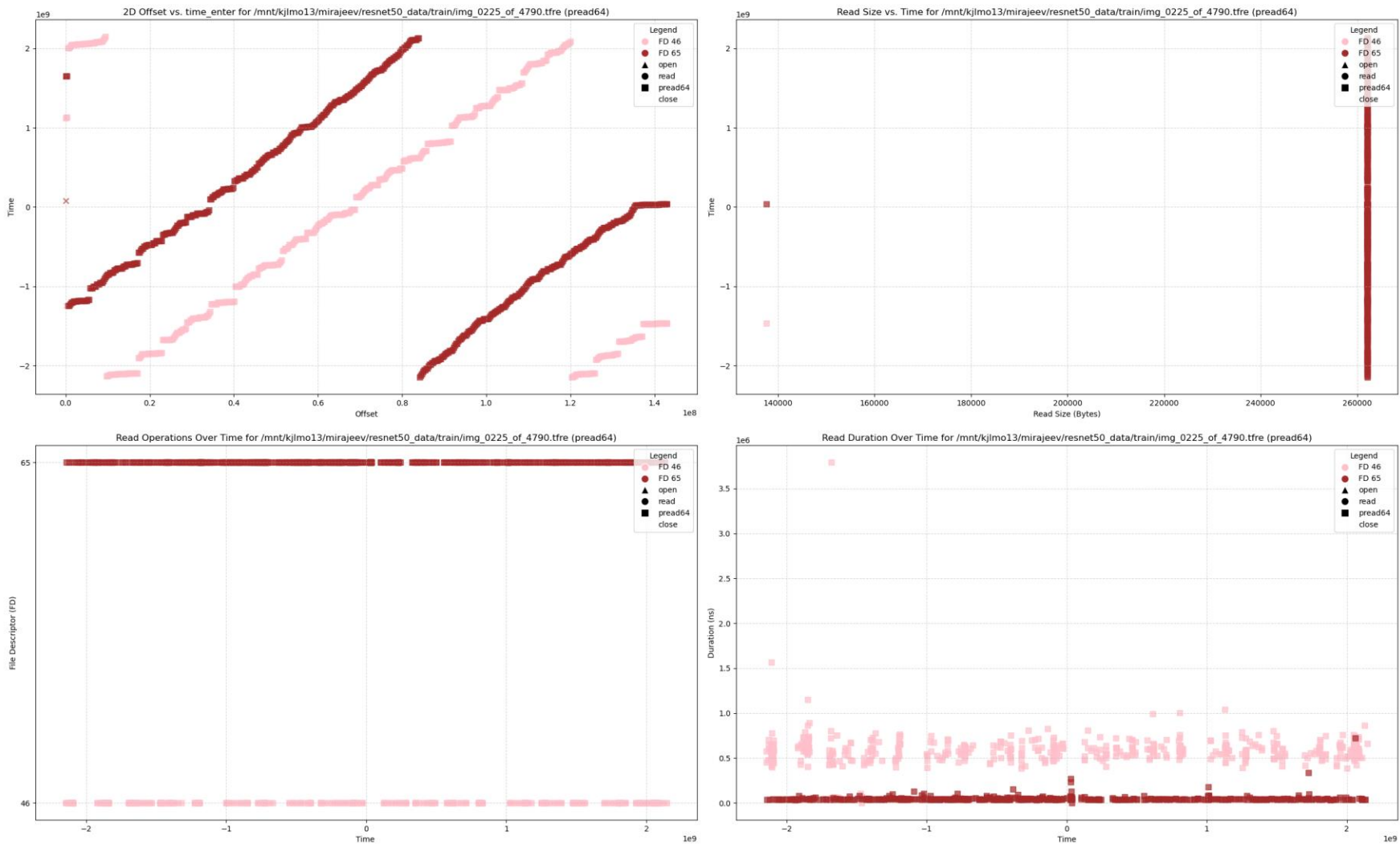
```
python3 dlio benchmark/dlio benchmark/main.py --config-path=/home/mirajeev/storage/storage-conf workload=resnet50 h100
++ workload.workflow.generate.data=False ++ workload.workflow.train=True
++ workload.dataset.data.folder=/mnt/kgimo13/mirajeev/resnet50_data ++ workload.workflow.profling=False
++ workload.profling.profler=None ++ hydra.output.subdir=configs ++ hydra.run.dir=/mnt/kgimo13/mirajeev/resnet50_h100
```



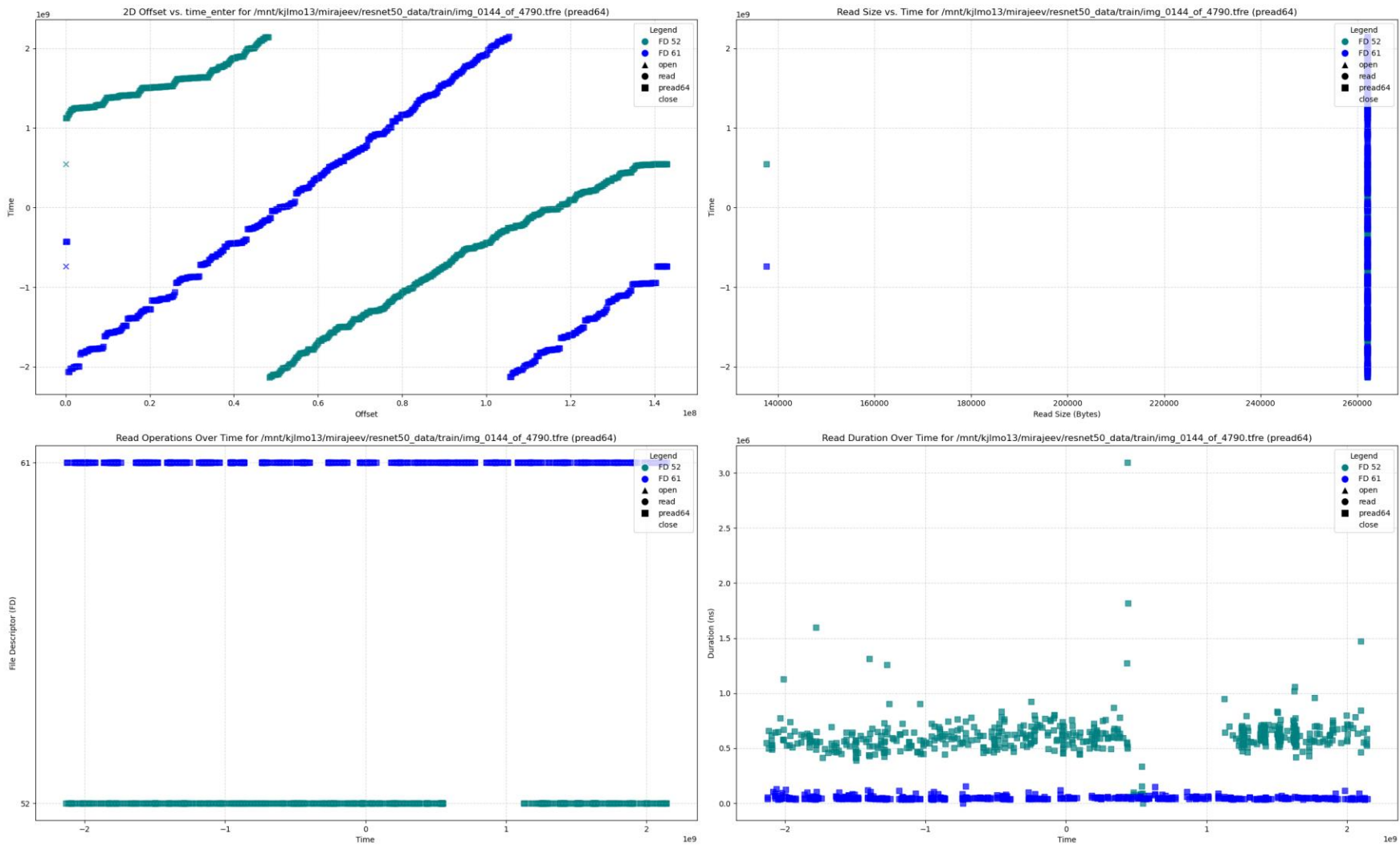


sequential: An I/O op issued at an offset greater than where the previous I/O op ended.

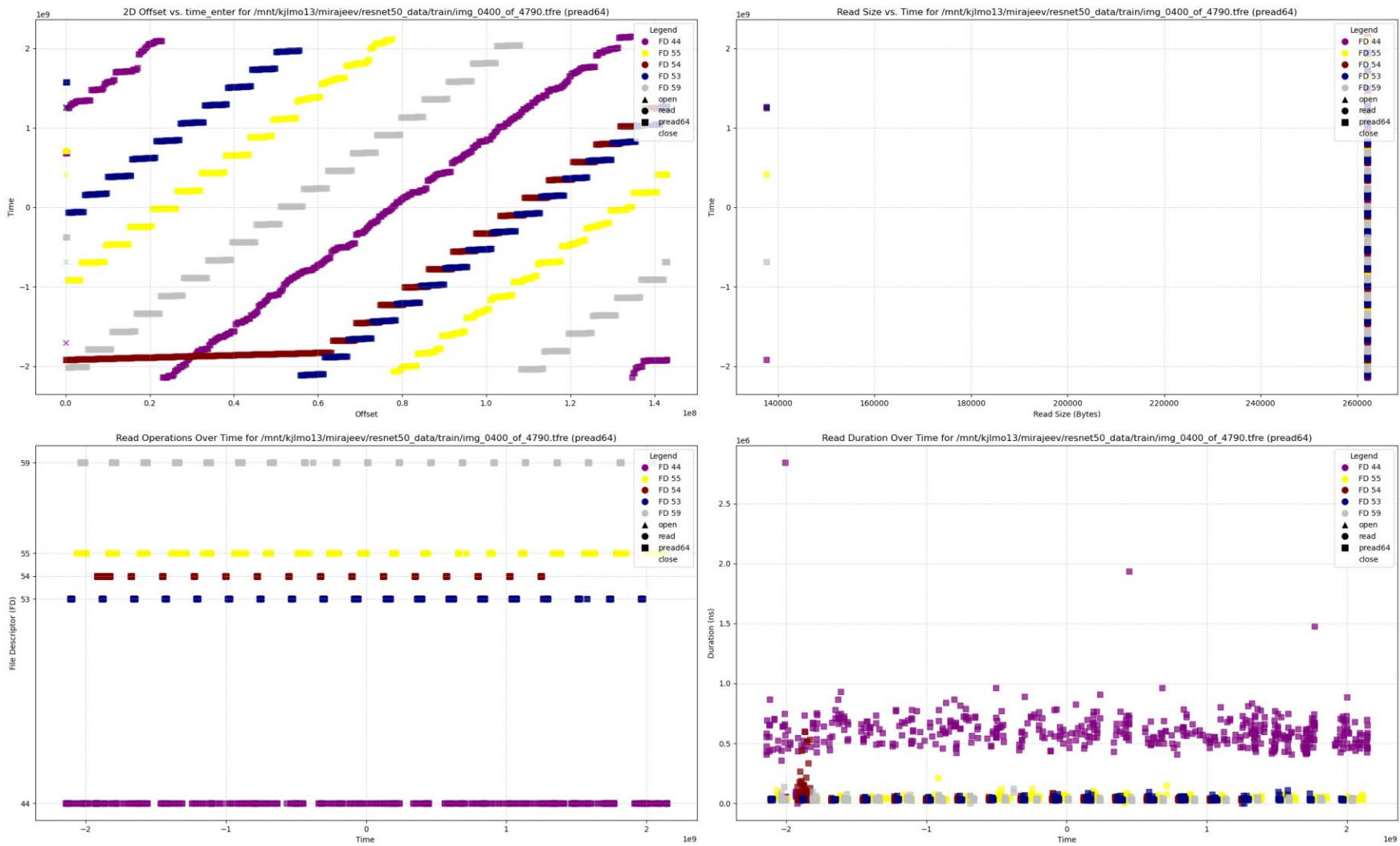
ResNet Epoch 1



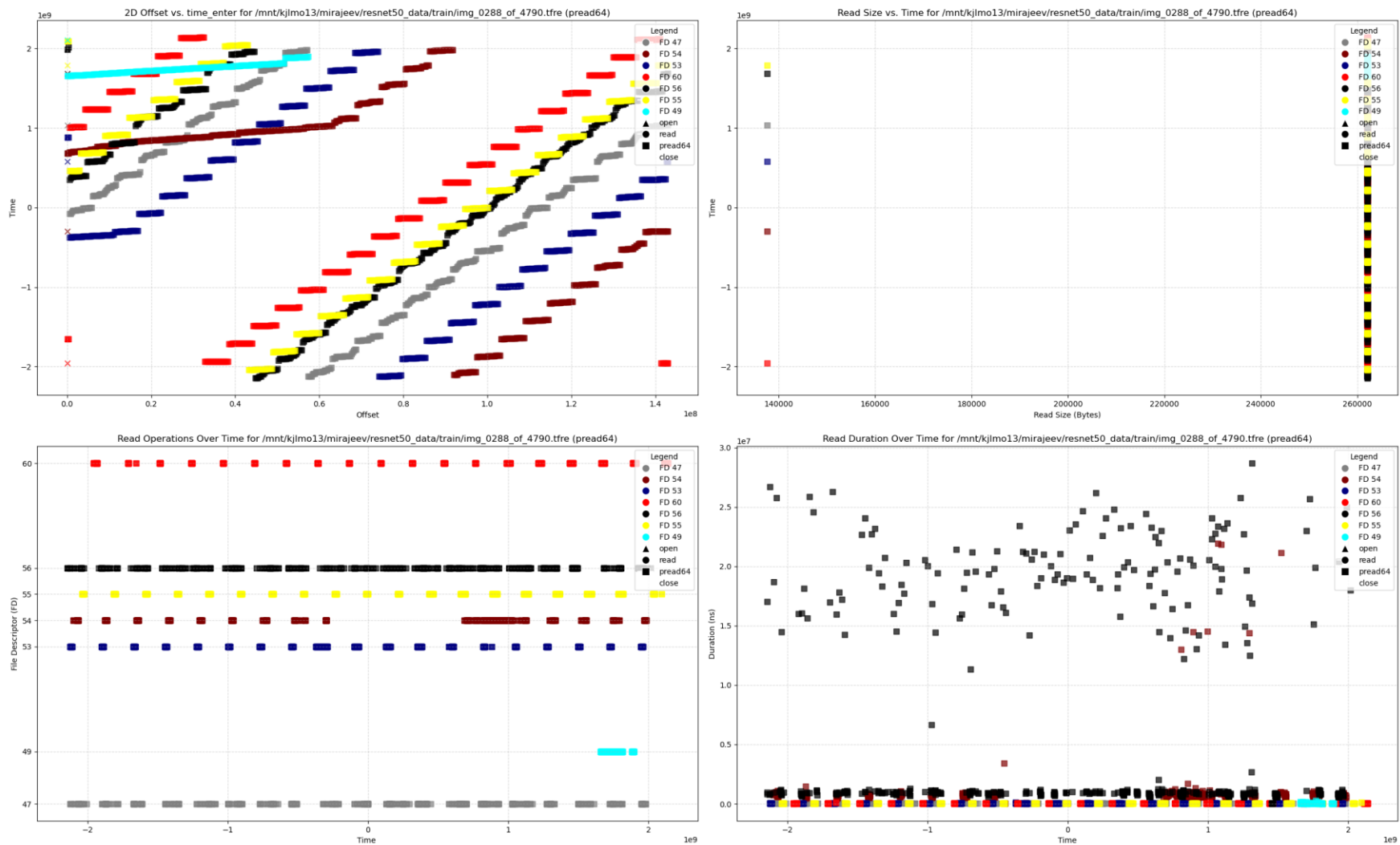
ResNet Epoch1 continue



ResNet Epoch5



ResNet Epoch 5 continue



Darshan

Bpftrace

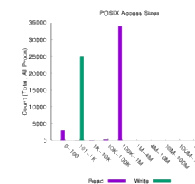
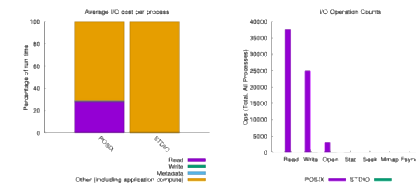


Page 10 of 10

2 of 4

jobid: 1962246	uid: 121685592	nprocs: 1	runtime: 133.1322 seconds
----------------	----------------	-----------	---------------------------

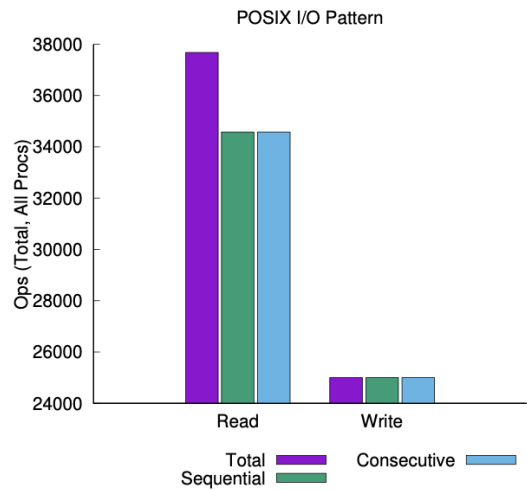
I/O performance *estimate* (at the POSIX layer): transferred **8376.6 MiB** at **217.84 MiB/s**
I/O performance *estimate* (at the STDIO layer): transferred **0.0 MiB** at **20.95 MiB/s**



Most Common Access Sizes (POSIX or MPI-IO)			File Count Summary (estimated by POSIX I/O access offsets)			
	access size	count	type	number of files	avg. size	max size
POSIX	262144	31464	total opened	1024	2.7MiB	4.3MiB
	179	10927	read-only files	1021	2.7MiB	2.9MiB
	178	10735	write-only files	1	4.3MiB	4.3MiB
	177	2567	read/write files	0	0	0
			created files	1	4.3MiB	4.3MiB

```
python3 dllo benchmark/dllo benchmark/main.py --config-path=/home/mirajeev/storage/storage-conf workload=cosmoflow h100
++ workload/workflow generate_data=False ++ workload/workflow train=True
++ workload/dataset folder=/mnt/kg1mo13/mirajeev/cosmoflow data ++ workload/workflow profiling=False
++ workload/profiling profiler=None ++ hydra.output subdir=configs ++ hydra.run.dir=/mnt/kg1mo13/mirajeev/cosmoflow h100
```

I/O Read Pattern



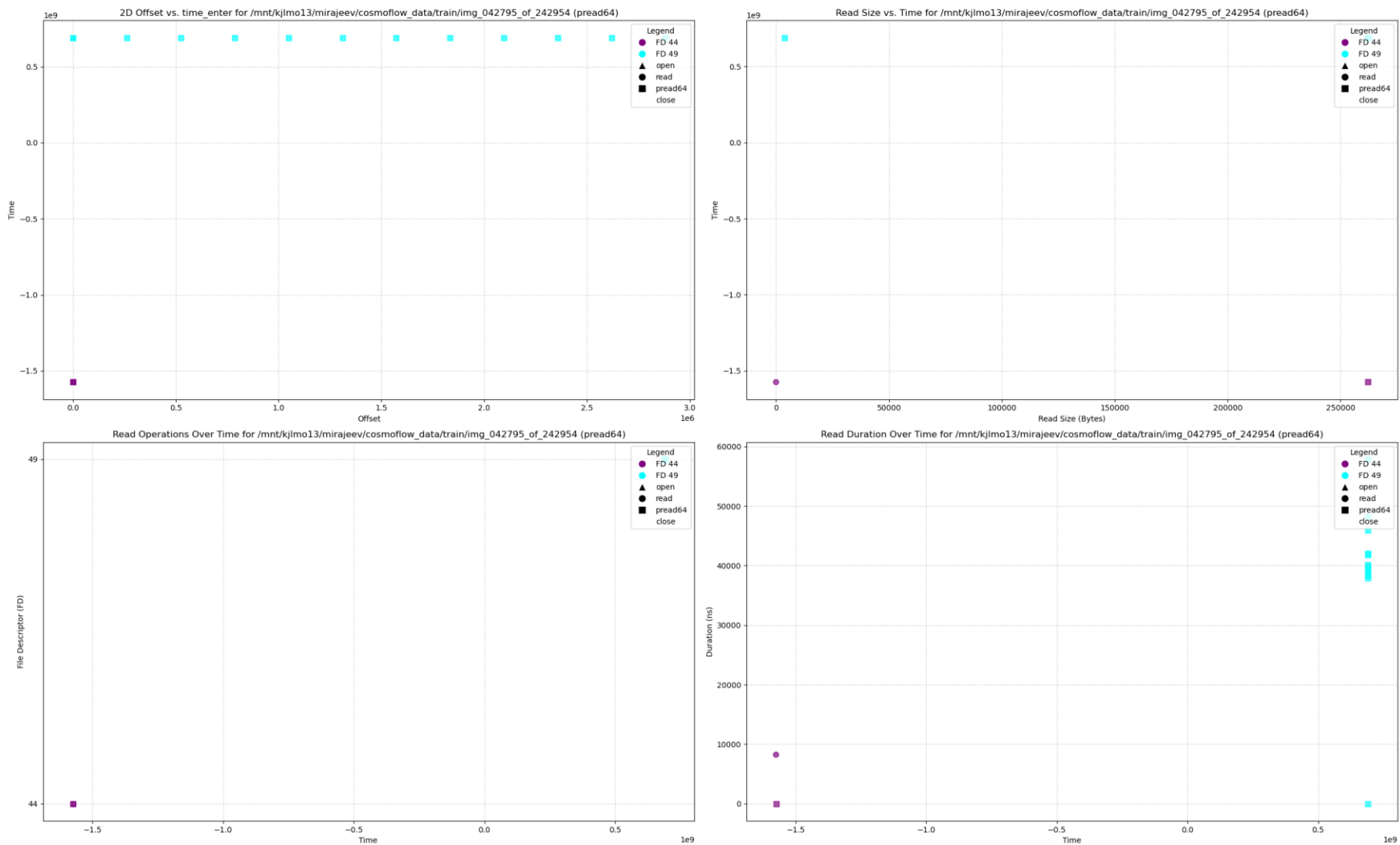
sequential: An I/O op issued at an offset greater than where the previous I/O op ended.
consecutive: An I/O op issued at the offset immediately following the end of the previous I/O op.

Variance in Shared Files (POSIX and STDIO)

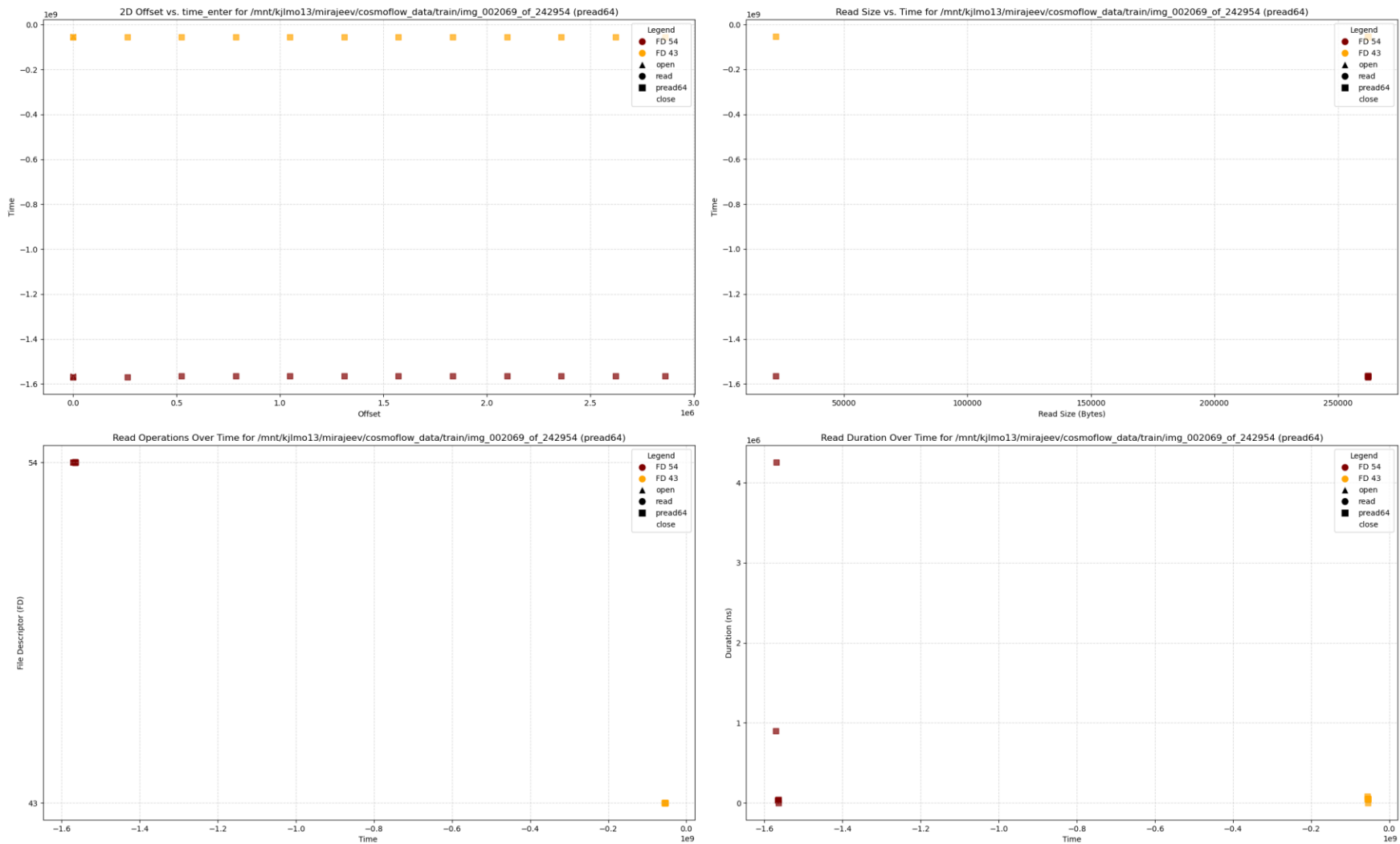
File Suffix	Processes	Fastest			Slowest			σ	
		Rank	Time	Bytes	Rank	Time	Bytes	Time	Bytes



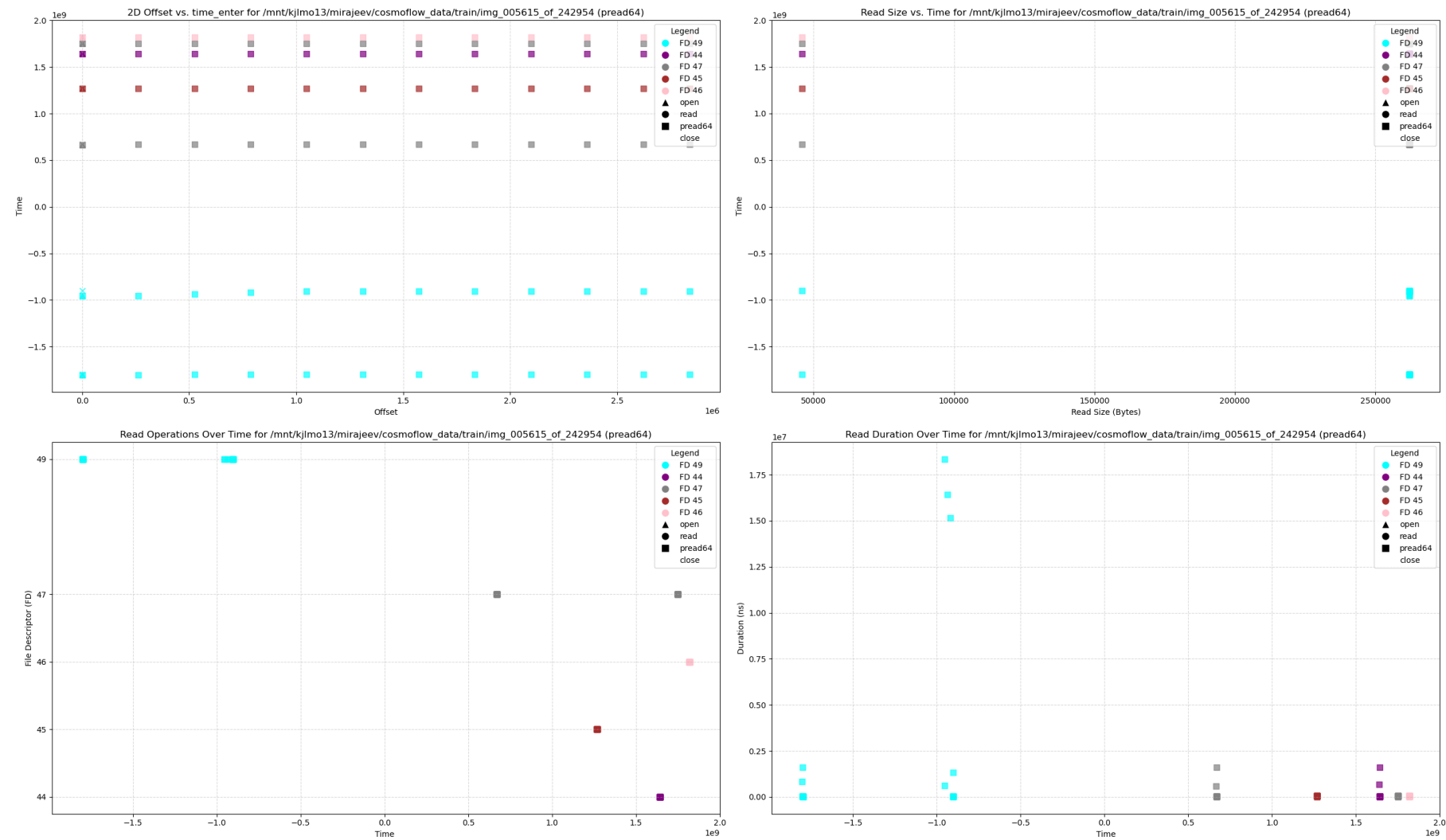
CosmoFlow Epoch1



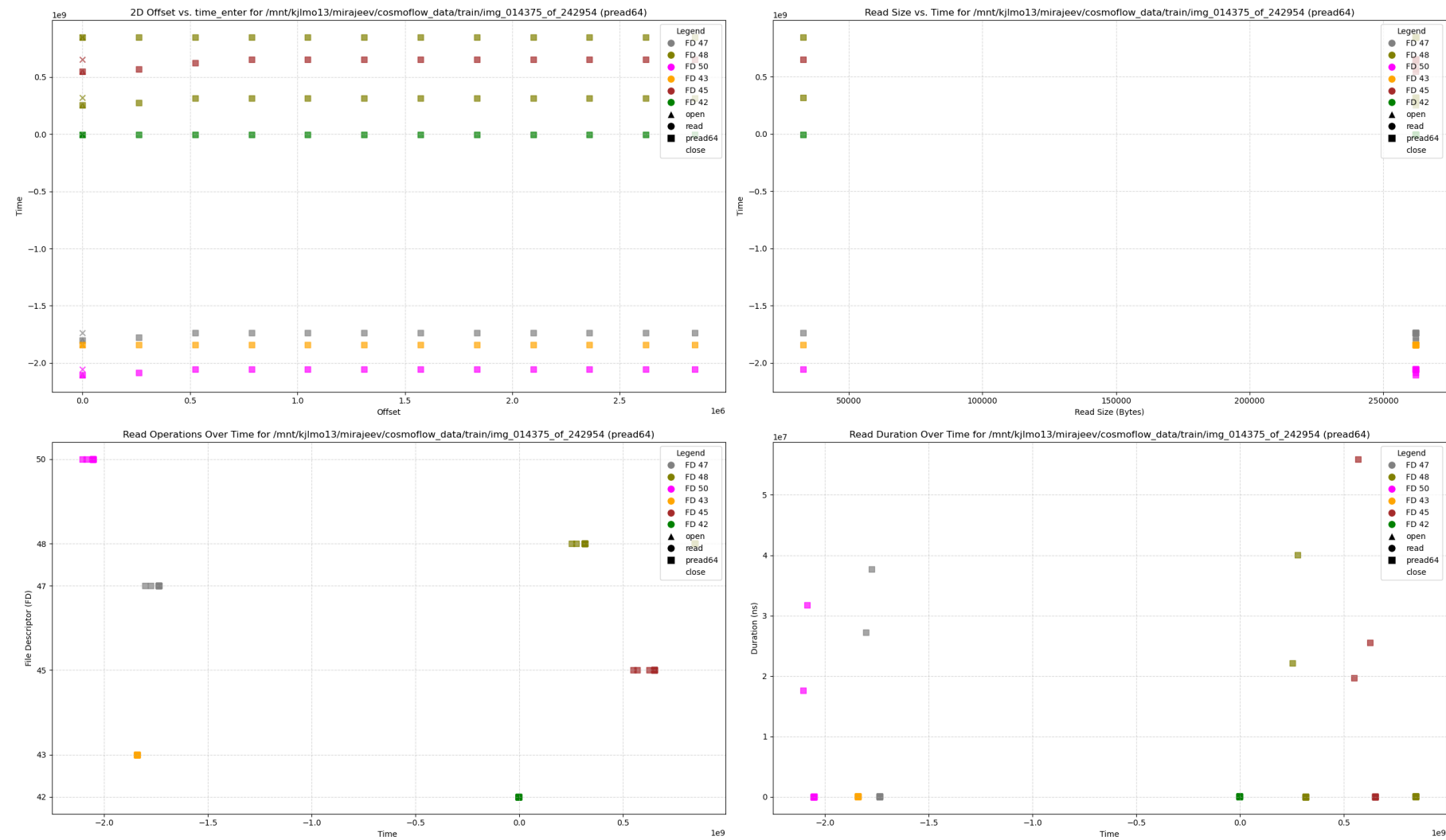
CosmoFlow Epoch1 continue



CosmoFlow Epoch5



CosmoFlow Epoch5 continue



IO Profile Inferences

UNet 3D

Read size: 4MB

Sequential & consecutive

I/O intensive (larger read size compared to others)



IO Profile Inferences Continue

ResNet50

Read size: 256K

Sequential & consecutive

CosmoFlow

Read size: 256K

Sequential & consecutive

IO Burst



Unet3d hybrid off

[METRIC] =====
[METRIC] Number of Simulated Accelerators: 2
[METRIC] Training Accelerator Utilization [AU] (%): 91.6906 (0.0000)
[METRIC] Training Throughput (samples/second): 35.4613 (0.0000)
[METRIC] Training I/O Throughput (MB/second): 4957.8110 (0.0000)
[METRIC] train_au_meet_expectation: success



Unet3d hybrid on

[INFO] Averaged metric over all epochs

[METRIC] =====

[METRIC] Number of Simulated Accelerators: 2

[METRIC] Training Accelerator Utilization [AU] (%): 34.9955 (0.0000)

[METRIC] Training Throughput (samples/second): 14.2703 (0.0000)

[METRIC] Training I/O Throughput (MB/second): 1995.1266 (0.0000)

[METRIC] train_au_meet_expectation: fail



Resnet hybrid off

[METRIC] =====
[METRIC] Number of Simulated Accelerators: 4
[METRIC] Training Accelerator Utilization [AU] (%): 92.8829 (0.0000)
[METRIC] Training Throughput (samples/second): 6632.4851 (0.0000)
[METRIC] Training I/O Throughput (MB/second): 725.2514 (0.0000)
[METRIC] train_au_meet_expectation: success
[METRIC]



Resnet hybrid on

[METRIC] =====
[METRIC] Number of Simulated Accelerators: 4
[METRIC] Training Accelerator Utilization [AU] (%): 92.6233 (0.0000)
[METRIC] Training Throughput (samples/second): 6613.6918 (0.0000)
[METRIC] Training I/O Throughput (MB/second): 723.1964 (0.0000)
[METRIC] train_au_meet_expectation: success
[METRIC] =====



Cosmoflow hybrid off

[METRIC] =====
[METRIC] Number of Simulated Accelerators: 2
[METRIC] Training Accelerator Utilization [AU] (%): 88.6635 (0.1147)
[METRIC] Training Throughput (samples/second): 504.6917 (2.4233)
[METRIC] Training I/O Throughput (MB/second): 1361.3829 (6.5368)
[METRIC] train_au_meet_expectation: success



Cosmoflow hybrid on

[METRIC] =====
[METRIC] Number of Simulated Accelerators: 2
[METRIC] Training Accelerator Utilization [AU] (%): 88.4566 (0.1159)
[METRIC] Training Throughput (samples/second): 503.4359 (1.6254)
[METRIC] Training I/O Throughput (MB/second): 1357.9954 (4.3844)
[METRIC] train_au_meet_expectation: success



Tuning

- Readahead will mostly help as the workload is sequential
- Testing was conducted on non-FOLIO systems
- Observed multiple misses in read ahead
- Size of read is same



Moving Forward

- Read ahead optimization
- Can fast path be improved for AI workload
- Will FOLIO changes help
- IO Burst how we can handle it efficiently



Inference

- Hybrid I/O (DIO/BIO) performs similarly on CosmoFlow and ResNet.
- BIO performance can be further improved with optimized readahead tuning.



Thank You



Rajeev Mishra
rajeevm@hpe.com

