



Deep Learning I/O: Benchmark and Profiling

Hariharan Devarajan (hariharandev1@llnl.gov)



This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 (LLNL-PRES-861305)

Deep Learning to accelerate science

Variety of

applications for

DL

- Diverse scientific domains utilize deep learning to solve scientific problems.
 - accelerate time-to-results.
 - extract patterns of large data.
- Data is the core of all deep learning methods.

Understanding the I/O behavior of scientific DL applications is crucial for scientific discovery.



Data Management in DL frameworks

Data

Management

Popular DL frameworks provide methods for tuning the input pipeline.

Major phases for data management

- Data Ingestion
 - different data sources
 - o data representations
- Data Processing
 - Parallelism and pipeline through worker assignment.
- Data Communication
 - RDMA,
 - GPU-Direct, and
 - NVLink

Input pipeline in TensorFlow with optimization example.



https://towardsdatascience.com/building-efficient-data-pipelines-using-tensorflow-8f647f03b4ce

Traditional Vs Deep Learning Stack

I/O tasks	Data Phase	Process Management	Language
• Coordinated Vs Independent	 Bulk Synchronous Vs Asynchronous background I/O 	 MPI Vs. Fork Spawn Threading DDP Horovod 	 C/C++/Fortran Vs. Python/C/C++

Observation

I/O research for DL applications require an **understanding** of data access behavior of modern scientific applications.

There is a need of an I/O benchmark and profiler representing I/O for deep learning workloads



DLIO Benchmark

A Data-Centric Benchmark for Scientific DL Applications

<u>Code</u>



https://github.com/argonnelcf/dlio_benchmark.git





SCAN ME

https://github.com/hariharandevarajan/dlio-profiler.git

A Dataflow Tracer for large-scale Deep Learning workloads

DLIO Profiler

DLIO Benchmark

A Data-Centric Benchmark for Scientific DL Applications

High-Level Design

- Component Design
- Reader allows different data formats
- Checkpoint allows different checkpointing implementations
- Data loaders implement different input pipeline.

All Components support plugin design to add custom implementations.



Top Level





Workflow



Dataset



Reader



Train



Checkpoint



Supported Workloads



Current known adopters



DLIO Benchmark Scaling



Shared File applications



Observation

- Splitting data across multiple files achieves better I/O performance than using one large shared file.
- Additionally, only TFRecord have I/O optimization for input pipeline which is not present in scientific data formats such as HDF5.

<u>Code</u>



SCAN ME

https://github.com/hariharandevarajan/dlio-profiler.git

DLIO Profiler

A Data-Centric Benchmark for Scientific DL Applications

Challenges of Current tools

- Need for handling dynamic processes
- Multi-level profiling
 - Python
 - I/O Workers
 - C/C++ Code
 - System Calls
- Metadata Tagging



DLIO Profiler (High-Level Design)



Unet3D (Perfetto UI by Google)

🏫 Perfetto 🛛 😑	Q Search or type '>' for commands or ':' for SQL mode									
Navigation ^		00:00:00	00:00:02	00:00:04	00:00:06	00:00:08	00:00:10	00:00:12	00:00:14	00:00:16
🗁 Open trace file	19779d20:35t41 848 067 000	00:00:00 000 000 000 000 000 000 000 00	00:00:02 000 000 000	00:00:04 000 000 000	00:00:06 000 000 000	00:00:08 000 000 000	00:00:10 000 000 000	00:00:12 000 000 000	00:00:14 000 000 000	00:00:16 000 000 000
Dpen with legacy UI	× <i>≡</i>									
 Record new trace 	A Process 0									
Current Trace	Thread 213592	РуТ	DLIOBenchmark.initialize	P.	<module< th=""><th>».iter</th><th>DLIOBenchmark.run DLIOBenchmark.train <module>.yield <m TorchFramework.co Torchi</m </module></th><th>odule>.yield <module>.yie Framework.co TorchFramework</module></th><th>ld <module>.yield Tor .co TorchFramework.co</module></th><th>chFramework.co</th></module<>	».iter	DLIOBenchmark.run DLIOBenchmark.train <module>.yield <m TorchFramework.co Torchi</m </module>	odule>.yield <module>.yie Framework.co TorchFramework</module>	ld <module>.yield Tor .co TorchFramework.co</module>	chFramework.co
Show timeline	Thread 014647				TorchDa NPZRead	Torc TorchDatas TorchData NPZR NPZReader NPZReade	a Torch TorchDa T Torch r VPZRe NPZRead N NPZR	Datase Torc TorchDa Tor eader.r NPZR NPZRead NPZ	chD TorchDat Torch Torch (Rea NPZReade NPZRe NPZR	Torc Torc e NPZR NPZR
🛃 Download	Inread 214647				NPZRead.	NPZ NPZReader NPZReade	NPZR NPZRea N NPZR	eader NPZ NPZRea NPZ	Rea NPZRead NPZR NPZRe	e NPZR NPZ
🥃 Query (SQL)					Torc To	orchData TorchData TorchD P7Reader NP7Reader NP7Rea	Ti rchDat Torc T Torc	TorchDatas TorchD Torc To NP7Reader NP7Rea NP7R N	orchD To TorchDatase Torch 27Rea NP NP7Reader NP7/	hDatase To Reader r NP
🖄 Viz	Thread 214648				NPZR. NF	PZReade NPZReade NPZRe	NF Reade NPZ N NPZR 1	NPZReader NPZRe NPZ NF	ZRea NP NPZReader NPZR	Reader NP
C Metrics										

MuMMI (AI-driven MD workflow)

22

Scheduler Allocation Detai Nodes: 32 Processes: 22949 Thread allocations acr Compute: 0 I/O: 2 Events Recorded: 4M	ils ross node	es (inclu	des dynam	ically	created	l threads
Description of Dataset Used └── Files: 1276825						
Benavior of Application Split of Time in appli Total Time: 42815. Overall I/0: 372.4 Metrics by function Function Split of Time in appli Function Image: state of the stat	ication .186 sec 411 sec ount M 33K N 34K 94K 7K NA 33K	n 25 IA as thes NA	si mean e operatio 11MB	ze mediar ns do n NA	i 75 ot have trans 6KB	max bytes sfered 588MB
- write 13	31K 1	2KB	148KB	2KB	9KB	3MB
mkdir 16 opendir 1k Listat64 61 fcntl 61 rmdir 61	N 13 13 13	IA as thes	e operatio	ns do n	ot have tran	bytes sfered





(c) High-level Summary.

DLIO Profiler Performance

	Score-P	Darshan	Recorder	DLIO Profiler (40 workers)
# Events Captured (out of 1.1M)	68,752	189	1,389	1.1 M
Overhead	13%	23%	13%	7%
Load Time for 1M events	3.56 min	1.6 min	13 min	62 sec
Load Time for 10M events	34.6 min	17 min	2.2 hr	1.3 min
Load Time for 100M events (hr)	6.1 hr	3.3 hr	>12 hr	3.4 min
Trace size: 1M events	59 MB	9 MB	20 MB	8.4 MB
Trace size: 10M events	610 MB	76 MB	169 MB	56 MB
Trace size: 100M events	5 GB	638 MB	1.3 GB	556 MB

Conclusions

A list of all observations

We demonstrated the need for a new benchmark and tool for Al workloads Described the modular design of DLIO benchmark and 2 the supported workloads. Described the DLIO Profiler design along with analysis 3 friendly trace format Demonstrated some visualizations possible using the 4 DLIO Profiler trace.



Benchmark



Thank you

hariharandev1@llnl.gov



Profiler



This material is based upon work supported by the National Science Foundation under Grant no. OCI-1835764 and CSR-1814872. Also, this work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. This material is based upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research under the DOE Early Career Research Program (LLNL-PRES-861305). Finally, this research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357.