DATA LOADING: the Next Frontier in Scale-out Deep Learning

Emily Watkins



DEEP LEARNING EXECUTION FLOW ON PAPER



















1 TRAINING ITERATION





1 TRAINING ITERATION





DATASETS IMPACT THROUGHPUT



Throughput (Img/Sec)



Impact of data loading

1. Overview of input pipelines & impact of data format.

2. How does data load time fit into overall training throughput?

3. Options for improving throughput based on training datasets.



Impact of data loading

- **1.** Overview of input pipelines & impact of data format.
- 2. How does data load time fit into overall training throughput?
- **3. Options for improving throughput based on training datasets.**

















+										
I I GPU I Fan	Name Temp	Perf	Persis Pwr:Us	tence-M age/Cap	Bus-Id	Memo	Disp.A pry-Usage	-+-	Volatile GPU-Util	Uncorr. ECC Compute M.
1 0 1 N/A	Tesla 48C	V100- P0	-SXM2 296W	0ff / 300W	0000000 15756M	0:06: iB /	:00.0 Off 16130MiB		90%	0 Default
1 N/A	Tesla 49C	V100 P0	-SXM2 238W	0ff / 300W	0000000 15756M	0:07 iB /	:00.0 Off 16130MiB	1	98%	0 Default
1 Z 1 N/A	Tesla 52C	V100- P0	-SXM2 273W	0ff / 300W	0000000 15756M	0:0A: iB /	:00.0 Off 16130MiB	1	98%	0 Default
3 N/A	Tesla 47C	V100- P0	-SXM2 260W	0ff / 300W	0000000 15756M	0:0B iB /	:00.0 Off 16130MiB	I	98%	0 Default
i 4 i N/A	Tesla 53C	V100 P0	-SXM2 283W	0ff / 300W	0000000 15756M	0:85 iB /	:00.0 Off 16130MiB	1	98%	0 Default
I 5 I N/A	Tesla 49C	V100- P0	-SXM2 82W	0ff / 300W	0000000 15756M	0:86 iB /	:00.0 Off 16130MiB	I	91%	0 Default
1 6 1 N/A	Tesla 51C	V100- P0	-SXM2 78W	Off / 300W	0000000 15756M	0:89 iB /	:00.0 Off 16130MiB		99%	0 Default
7 N/A	Tesla 45C	V100 P0	- SXM2 69W	0ff / 300W	0000000 15756M	0:8A: iB /	:00.0 Off 16130MiB	1	99%	0 Default
+										+
I Proc	esses:	PID	Туре	Process	name					GPU Memory I Usage I
0	21	.859	C	python						15745MiB
1	21	860	C	python						15745MLB
1 2	21	861	C	python						15745M1B
3	21	862	C C	python						15745M1B
4	21	864	ć	python						15745M1B
	21	965	c	python						15745MLB
	21	866	ć	python						15745MLB
+			······	pychon						+



GPU UTILIZATION IS NOT EVERYTHING.

- NOT THE FULL PICTURE

- NOT GRANULAR ENOUGH







BATCH TIMING

training logs:	batch,	ms
batch_results:	: 917 ,	1040
batch_results:	918 ,	43
batch_results:	919 ,	89
batch_results:	920 ,	29
batch_results:	921 ,	1025
batch_results:	922 ,	37
batch_results:	923 ,	23
batch_results:	924 ,	90
batch_results:	925 ,	1053
batch_results:	926 ,	31
batch_results:	927 ,	21



BATCH TIMING ODDITIES

training logs:	batch,	ms
batch results	: 917,	1040
batch results	: 918,	43
batch_results	: 919,	89
batch_results	: 920,	29
batch_results	: 921,	<mark>1025</mark>
batch_results	: 922,	37
batch_results	: 923,	23
batch_results	: 924,	90
batch_results	: 925,	<mark>1053</mark>
batch_results	: 926,	31
batch_results	: 927,	21



WHAT'S "GOOD" PERFORMANCE?





WHAT'S "GOOD" PERFORMANCE?





TESTING WITH SYNTHETIC DATA





BREAK DOWN THE PROBLEM SPACE

REAL VS. SYNTHETIC DATA

Linear Performance Synthetic Performance	28,200 images/s
	27,200 images/s





TESTING WITH SYNTHETIC DATA







3x3 conv, 128

3x3 conv, 128

¥-----

24 © 2019 PURE STORAGE INC.

UNDER THE HOOD





CAT



BOAT

CAT



- 2. Associate labels
- 3. Shuffle
- 4. Read, crop, distort
- 5. Convert to tensor
- 6. Copy to GPU





ISSUE IDENTIFIED: INPUT PIPELINE

IMAGENET PERFORMANCE BENCHMARK, RESNET-50





WE BALANCED THIS PIPELINE





UNDER THE HOOD





CAT



BOAT

CAT





DOG



6. Copy to GPU

1. Enumerate

3. Shuffle

2. Associate labels



4. Read, crop, distort 5. Convert to tensor



! DATASETS FOR BENCHMARK JOBS





! DATASETS FOR BENCHMARK JOBS





TENSORS: LIMITED PROCESSING

tensor([[[0.5022, [0.4166, [0.4851,	0.6049, 0.6221, 0.6049,	0.7077, 0.6392, 0.5878,	· · · , · · · ,	0.7077, 0.9646, 1.4098,	0.6392, 0.9303, 1.3242,	0.7762], 0.9303], 1.2043],
[1.7694, [1.4098, [1.7865,	1.6324, 1.2214, 1.5297,	1.4612, 0.9988, 1.1529,	· · · , · · · ,	1.1015, 1.1015, 1.4612,	0.9988, 1.0844, 1.5468,	0.7762], 1.2728], 1.7694]],
[[-0.2675, [-0.4076, [-0.3550,	-0.1099, -0.1625, -0.2150,	0.0651, -0.0749, -0.1625,	· · · , · · · ,	0.2227, 0.5203, 0.8880,	0.1352, 0.4328, 0.8004,	0.2402], 0.4153], 0.7129],
[1.3431, [0.7829, [1.0455,	1.1506, 0.5553, 0.7479,	0.9580, 0.3277, 0.3803,	· · · , · · · ,	0.2577, 0.3277, 0.7654,	0.1527, 0.3102, 0.8529,	-0.0749], 0.5028], 1.0630]],
[[0.6879, [0.5136, [0.5834,	0.7925, 0.7402, 0.6879,	0.9145, 0.7402, 0.6705,	· · · , · · · ,	0.7751, 1.0714, 1.4374,	0.6705, 1.0017, 1.3502,	0.7751], 1.0017], 1.2805],
[2.0300, [1.4548, [1.6465,	1.8905, 1.2631, 1.3851,	1.7163, 1.0365, 1.0017,	· · · , · · · ,	1.1062, 1.1237, 1.4200,	1.0365, 1.1411, 1.5071,	0.8099], 1.2980], 1.6814]]])











PNG IMAGES - NICE LINEAR SCALE





PNG IMAGES – 50% THROUGHPUT

EVEN WHEN "TUNED"



PNG images: utilizes data pipeline to load & transform SYNTHETIC: data generated on GPUs



IMAGE REPRESENTATION





IMAGE FORMAT



Higher image complexity = Longer load time



IMAGE FORMAT & IMAGE SIZE



Larger image size = Longer load time



Impact of data loading

1. Overview of input pipelines & impact of data format.

2. How does data load time fit into overall training throughput?

3. Options for improving throughput based on training datasets.



LIFE OF A BATCH





	tensor([[[0.5022, [0.4166, [0.4851,	0.6049, 0.707 0.6221, 0.639 0.6049, 0.587	7,, 2,, 8,,	0.7077, 0.9646, 1.4098,	0.6392, 0.9303, 1.3242,	0.7762], 0.9303], 1.2043],	
P and	[1.7694, [1.4098, [1.7865,	1.6324, 1.461 1.2214, 0.998 1.5297, 1.152	2,, 18,, 19,,	1.1015, 1.1015, 1.4612,	0.9988, 1.0844, 1.5468,	0.7762], 1.2728], 1.7694]],	
	[[-0.2675, [-0.4076, [-0.3550,	-0.1099, 0.065 -0.1625, -0.074 -0.2150, -0.162	i1,, 9,, 5,,	0.2227, 0.5203, 0.8880,	0.1352, 0.4328, 0.8004,	0.2402], 0.4153], 0.7129],	
1	[1.3431, [0.7829, [1.0455,	1.1506, 0.958 0.5553, 0.327 0.7479, 0.380	0,, 7,, 3,,	0.2577, 0.3277, 0.7654,	0.1527, 0.3102, 0.8529,	-0.0749], 0.5028], 1.0630]],	
	[[0.6879, [0.5136, [0.5834,	0.7925, 0.914 0.7402, 0.740 0.6879, 0.670	5,, 12,, 15,,	0.7751, 1.0714, 1.4374,	0.6705, 1.0017, 1.3502,	0.7751], 1.0017], 1.2805],	
	[2.0300, [1.4548, [1.6465.	1.8905, 1.716 1.2631, 1.036 1.3851, 1.001	i3,, i5,, .7,,	1.1062, 1.1237, 1.4200,	1.0365, 1.1411, 1.5071,	0.8099], 1.2980], 1.6814]]])	





LIFE OF A BATCH





WORKER DEPENDENCY: BATCH 0

1 CPU worker, 1 GPU





WORKER DEPENDENCY: BATCH 1





WASTED TIME ADDS UP OVER EPOCHS





WASTED TIME ADDS UP OVER EPOCHS,

ESPECIALLY WITH LARGER DATA LOAD TIME





46

CONCURRENT WORKERS

SINGLE GPU







EXAMPLE WITH A REAL RATIO







EXAMPLE WITH A REAL RATIO

MULTI-GPU







IMAGE LOAD TIME V. GPU TIME

224x224 TENSOR, 8 GPUs, 16 WORKERS







IMAGE LOAD TIME V. GPU TIME

768x768 PNG, 8 GPUs, 16 WORKERS

time







IMAGE LOAD TIME V. GPU TIME

768x768 JPEG, 8 GPUs, 16 WORKERS

time







Impact of data loading

1. Overview of input pipelines & impact of data format.

2. How does data load time fit into overall training throughput?

3. Options for improving throughput based on training datasets.



Optimization Suggestions

- **1. Precaching:** saving a pre-transformed tensor version of your dataset can minimize data load time during training (at the cost of on-the-fly distortions).
- **2. Workload updates:** some workflows were designed around non-GPU compute environments. Investigate utilizing newer image loading libraries.
- **3. Data load on GPU:** Nvidia DALI can be used to offload CPU work to GPUs, especially if workflows have been updated to use GPU-friendly data loading libraries.



Conclusions

- **1.** Data format can significantly impact training throughput.
- 2. There is no one-size-fits-all input pipeline.
- 3. It's critical to have have a mental model for your input pipeline and a methodology for testing its performance.



QUESTIONS

@dataemilyw



