

github.com/Emma926/paradnn







A Systematic Methodology for Analysis of Deep Learning Hardware and Software Platforms

Yu (Emma) Wang, Gu-Yeon Wei, David Brooks Harvard University

Contact: ywang03@g.harvard.edu 3/3/2020



HARVARD John A. Paulson School of Engineering and Applied Sciences

Acknowledgement

Frank Chen, Glenn Holloway, Dan Janni, Peter Mattson, Lifeng Nai, David Patterson, Francesco Pontiggia, Parthasarathy Ranganathan, Vijay Reddi, Brennan Saeta, Zak Stone, Anitha Vijayakumar, Shibo Wang, Qiumin Xu, Doe Hyun Yoon, Cliff Young



HARVARD John A. Paulson School of Engineering and Applied Sciences



Challenges with ML Benchmarking

- Diversity in deep learning models used
 - Problem Domains, Models, Datasets
- Pace of field
 - State-of-the-art models evolve every few months
- Varying evaluation metrics
 - Accuracy, Time to train, Latency of inference
- Multi-disciplinary field
 - Algorithms, Systems, Hardware, ML Software Stacks



State of the art: MLPerf 0.6

Area	Benchmark	Dataset	Model	Reference Implementation
Vision	Image classification	ImageNet	ResNet-50	TensorFlow
	Object detection	COCO 2017	Mask R-CNN	Pytorch
	Object detection	COCO 2017	SSD-ResNet34	Pytorch
Language/	Translation	WMT Eng-Germ	Transformer	TensorFlow
Addio	Speech recognition	WMT Eng-Germ	GNMT	PyTorch
Commerce	Recommendation	MovieLens-20M	NCF	PyTorch
Action	Reinforcement Learning	Go	Mini-go	TensorFlow



State of the art: MLPerf 0.6

Area	Benchmark			Dataset		Model		Reference	
				\frown		Implementation			
Vision	1	Image classification		ImageNet		ResNet-50		TensorFlow	
		Object detection		COCO 2017	/	Mask R-CNN		Pytorch	
		Object detection	Τ	COCO 2017	V	SSD-ResNet34	X	Pytorch	
Language/		Translation	Τ	WMT Eng-Germ		Transformer		TensorFlow	
Audio		Speech recognition	Ι	WMT Eng-Germ	Λ	GNMT	X	PyTorch	
Commerce		Recommendation		MovieLens-20M		NCF	Λ	PyTorch	
Action		Reinforcement Learning		Go		Mini-go		TensorFlow	
		•	\checkmark		\checkmark		\checkmark		



Our Methodology





Our Methodology











ParaDnn vs MLPerf

ParaDnn

- Avoid drawing conclusions based on several arbitrary models
- Generate thousands of parameterized, end-to-end models
- Prepare hardware designs for future models
- Complement the use of existing real-world models, i.e. MLPerf



- Good for studying accuracy or convergence with real datasets
- Represent the specific models some people care about





ParaDnn Canonical Models



CNNs: Residual, Bottleneck



RNNs: RNN, LSTM, GRU





Models





Models



- ParaDnn covers a larger range than the real models
 - from 10k to ~1 billion parameters



Analysis Enabled by ParaDnn

- Roofline analysis of TPU v2
- Homogenous Platform Comparison: TPU v2 vs v3
- Heterogeneous Platform Comparison: TPU vs GPU

The Roofline Model









The Roofline Model





The Roofline Model





Transformer





FC Models

ParaDnn sweeps a large range of models, from memory-bound to compute-bound.



FC Models





FC Models





Memory-bound



TPU v2 vs v3?





















Architecture of TPU v2 vs v3





180 TFLOPS / Board



Figure is from https://cloud.google.com/tpu/docs/system-architecture

Google's Choice of TPU v3













31





Memory-bound: 1.5x speedup

32





Memory-bound, but benefit from 2x memory capacity:

Google's Choice of TPU v3









ParaDnn provides diverse set of operations, and shows different operations are sensitive to different system component upgrades.



TPU vs GPU?



Hardware Platforms

			Mem	Mem	Mem Bdw	Peak
Platform	Unit	Version	Туре	(GB)	(GB/s)	FLOPS
GPU	1	V100				
(DGX-1)	Pkg	(SXM2)	HBM2	16	900	125T
	1 Board					
TPU	(8 cores)	v2	HBM	8	2400	180T



Hardware Platforms

			Mem	Mem	Mem Bdw	Peak
Platform	Unit	Version	Туре	(GB)	(GB/s)	FLOPS
GPU	1	V100				
(DGX-1)	Pkg	(SXM2)	HBM2	16	900	125T
	1 Board					
TPU	(8 cores)	v2	HBM	8	2400	180T

300 GB/s per core

FC and CNN



FC and CNN





Hardware Platforms

			Mem	Mem	Mem Bdw	Peak
Platform	Unit	Version	Туре	(GB)	(GB/s)	FLOPS
GPU	1	V100				
(DGX-1)	Pkg	(SXM2)	HBM2	16	900	125T
	1 Board					
TPU	(8 cores)	v2	HBM	8	2400	180T
					300 GB/s	
					per core	



FC TPU/GPU Speedup colored with Batch Size





FC TPU/GPU Speedup colored with **Batch Size**





FC TPU/GPU Speedup colored with **Batch Size**





FC TPU/GPU Speedup colored with **Node Size**





Hardware Platforms

						1.44x
			Mem	Mem	Mem Bdw	Peak
Platform	Unit	Version	Туре	(GB)	(GB/s)	FLOPS
GPU	1	V100				
(DGX-1)	Pkg	(SXM2)	HBM2	16	900	125T
	1 Board					
TPU	(8 cores)	v2	HBM	8	2400	180T
					300 GB/s	
					per core	



CNN TPU/GPU Speedup colored with **Batch Size**





CNN TPU/GPU Speedup colored with **Batch Size**





CNN TPU/GPU Speedup colored with **Batch Size**





CNN TPU/GPU Speedup colored with **Filters**





Conclusion

- Parameterized methodology: ParaDnn + a set of analysis methods
- Single platform analysis: TPU v2
- Homogenous platform comparison: TPU v2 vs v3
- Heterogeneous platform comparison: TPU vs GPU



Limitations of this Work

- Does not include:
 - Inference
 - Multi-node system: multi-GPU, or TPU pods
 - Accuracy, convergence
 - Cloud overhead
- Tractability
 - Limit the range of hyperparameters and datasets
 - Small batch sizes (<16) and large batch sizes (> 2k) are not studied
 - Synthetic datasets do not include data infeed overhead
 - Iterations of TPU loop is 100. Larger numbers can slightly increase the performance.



Available: github.com/Emma926/paradnn

Questions?









HARVARD John A. Paulson School of Engineering and Applied Sciences