



Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool, Peter Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff Young, Ryan Sepassi, Blake Hechtman, Google Al

## **Massive Model-Parallelism** (and/or data-parallelism) Made Easier

- Describe the overall computation in a TensorFlow-like language with named tensor-dimensions.
- Describe your physical cores as n-dimensional array of processors by specifying a "mesh".
- Describe which tensor-dimensions should be split across which dimensions of the mesh of processors.
- You're Done Mesh-TensorFlow compiles your graph into Single-Program-Multiple-Data (SPMD) TensorFlow code plus collective communication primitives.

## Motivation

- Ability to train **GIANT** multi-billion/trillion-parameter models which do not fit on one processor.
- Ability to process giant examples (spatial/temporal splitting of images/video, etc.)
- Low latency by splitting computation for one example across multiple processors.
- Traditional MIMD approaches to model-parallelism are: Tricky to design
- Create giant cumbersome graphs
- Are prone to bottlenecks

# Mesh-TensorFlow: Deep Learning for SuperComputers



allreduce across one mesh dimension.

Integrated on Google Cloud TPU along with examples like the



• Transformer models working in MeshTF

layout\_model\_parallel="vocab:m0,heads:m0,d\_ff:m0" layout\_dp\_mp="batch:m0,vocab:m1,heads:m1,d\_ff:m1"

Table 3: Transformer Machine-Translation Results.  $d_{model} = 1024, d_k = d_v = 128$ 

f	heads	$d_k, d_v$	Parameters	WMT14 EN-DE	WMT14 EN-FR	
			(Billions)	BLEU	BLEU	
18	4	128	0.15	25.5	41.8	
96	8	128	0.24	26.5	42.5	
92	16	128	0.42	27.1	43.3	
84	32	128	0.77	27.5	43.5	
68	64	128	1.48	27.5	43.8	
36	128	128	2.89	26.7	43.9	
96	16	64	0.21	28.4	41.8	[21]

Table 2: Transformer-Decoder Language Models:  $d_{model} = 1024, d_k = d_v = 256$ 

heads	Parameters	Billion-Word Benchmark	Wikipedia
	(Billions)	Word-Perplexity	Subword-Perplexity
4	0.14	35.0	8.74
8	0.22	31.7	8.03
16	0.37	28.9	7.44
32	0.67	26.8	6.99
64	1.28	25.1	6.55
128	2.48	24.1	6.24
256	4.90	24.0(23.5)	6.01
JN [20]	6.5	28.0	
nsemble [17]		26.1	
le (different methods)[17]	> 100	23.7	

• Trained Transformer models with up to 5B parameters on

(6PFLOPS/possible 11.5 PFLOP/s)

### Status

• Code is Open-Source on github - please contribute. https://github.com/tensorflow/mesh/tree/master/mesh\_tensor

Implementations to produce SPMD code for TPU or MIMD



