

Getting Started with TensorFlow on GPUs^o



Magnus Hyttsten

@MagnusHyttsten



Agenda





An Awkward Social Experiment (that I'm afraid you will be part of...)



GPU TECHNOLOGY CONFERENCE

SILICON VALLEY -



GPU TECHNOLOGY CONFERENCE

ROCKS!









CAT DOG









T





TensorFlow 2.0 Alpha is out







Easy

Simplified APIs. Focused on Keras and eager execution

Powerful

Flexibility and performance. Power to do cutting edge research and scale to > 1 exaflops

Scalable

Tested at Google-scale. Deploy everywhere







tf.data (Dataset) tf.feature_column

(Transfer Learning)

High-level APIs

Perform Distributed Training (talk @1pm)

E.g. V100



tf.data (Dataset) tf.feature_column

(Transfer Learning)

High-level APIs

Perform Distributed Training (talk @1pm)

E.g. V100

Built to Distribute and Scale



Premade Estimators

LinearRegressor(...) LinearClassifier(...) DNNRegressor(...) DNNClassifier(...) DNNLinearCombinedRegressor(...) DNNLinearCombinedClassifier(...) BaselineRegressor(...) BaselineClassifier(...) BoostedTreeRegressor(...) BoostedTreeClassifier(...)

Train locally

estimator =

estimator.train (input_fn= ...)
estimator.evaluate(input_fn= ...)
estimator.predict (input_fn= ...)

Premade Estimators

Datasets

Datasets

Premade Estimator - Wide & Deep

```
wide_columns = [
    tf.feature_column.bucketized_column(
        'age',=[18, 27, 40, 65])]
deep_columns = [
    tf.feature_column.numeric_column('visits'),
    tf.feature_column.numeric_column('clicks')]
```

```
tf.estimator.DNNLinearCombinedClassifier(
    linear_feature_columns=wide_columns,
    dnn_feature_columns=deep_columns,
    dnn_hidden_units=[100, 75, 50, 25])
```



tf.data (Dataset) tf.feature_column

(Transfer Learning)



Perform Distributed Training

E.g. V100

Custom Models

tf.keras

tf.keras.layers

model = tf.keras.models.Sequential([tf.keras.layers.Flatten(), tf.keras.layers.Dense(512, activation='relu'), tf.keras.layers.Dropout(0.2), tf.keras.layers.Dense(10, activation='softmax') 1) model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy']) (dataset, epochs=5) model.fit Datasets model.evaluate(dataset) model.predict (dataset)

TensorFlow Datasets

- 30+ available
- Add your own

- audio
 - o "nsynth"
- image
 - o "celeb_a"
 - o "cifar10"
 - o **"coco2014"**
 - "diabetic_retinopathy_detection"
 - "imagenet2012"
 - o "mnist"
 - "open_images_v4"

- structured
 - o "titanic"
- text
 - "imdb_reviews"
 - o **"lm1b"**
 - "squad"

- translate
 - "wmt_translate_ende"
 - "wmt_translate_enfr"
- video
 - "bair_robot_pushing_small"
 - "moving_mnist"
 - "starcraft_video"

TensorFlow Summary

- Datasets (tf.data) for the input pipeline
 - a. TensorFlow Datasets is great
 - b. tf.feature_columns are cool too
- Premade Estimators
- Keras Models (tf.keras)

The V-100



And why is it so good @ Machine Learning???



- High-Level We look at only parts of the power of GPUs
- Simple Overview More optimal designs exist
- Reduced Scope Only considering fully-connected layers, etc

Strengths of V100



- Built for Massively Parallel Computations
- Specific hardware / software to manage Deep Learning Workloads (Tensor Cores, mixed-precision execution, etc)

Strengths of V100



- Built for Massively Parallel Computations
- Specific hardware / software to manage Deep Learning Workloads (Tensor Cores, mixed-precision execution, etc)

Tesla SXM V100
 5376 cores (FP32)

My Questions Around the GPU

What are we going to do with 5376 FP32 cores?

The Unsatisfactory Answer

What are we going to do with 5376 FP32 cores?

"Execute things in parallel"!



What are we going to do with 5376 FP32 cores? "Execute things in parallel"!

Yes, but how can we exactly do that for ML Workloads?



What are we going to do with 5376 FP32 cores? "Execute things in parallel"!

Yes, but how can we exactly do that for ML Workloads? "Hey, that's your job - That's why we're here listening"!



What are we going to do with 5376 FP32 cores? "Execute things in parallel"!

Yes, but how can we exactly do that for ML Workloads? "Hey, that's your job - That's why we're here listening"!

Alright, let me try to talk about that then





- We may have a huge number of layers
- Each layer can have huge number of neurons
- --> There may be 100s millions or even billions * and + ops

All knobs are W values that we need to tune So that given a certain input, they generate the correct output



"Matrix Multiplication is

EATING (the computing resources of) **THE WORLD**"

$$h_{i_{j}} = [X_{0}, X_{1}, X_{2, ...}] * [W_{0}, W_{1}, W_{2, ...}]$$

$$h_{i_{j}} = X_0 W_0 + X_1 W_1 + X_2 W_2 + \dots$$

Matmul

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values W = [0.1, 0.1, ..., 0.1] # Then we need to have 256 weight values $h_{0.0} = X + W = [1+0.1 + 2+0.1 + ... + 256+0.1] == 32389.6$

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

 $h_{0,0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$



X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values
W = [0.1, 0.1, ..., 0.1] # Then we need to have 256 weight values
V = [0.1, 0.1, ..., 0.1] # Then we need to have 256 weight values





0.1

2

٠

*

 $1 \times 0.1 = 0.1$

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

 $h_{0.0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$



X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values \angle

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

 $h_{0,0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

h_{0.0} = X ★ W # [1★0.1 + 2★0.1 + ... + 256★0.1] == 32389.6

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

h_{0.0} = X ★ W # [1★0.1 + 2★0.1 + ... + 256★0.1] == 32389.6

GPU Execution

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values
W = [0.1, 0.1, ..., 0.1] # Then we need to have 256 weight values

 $h_{0.0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

 $h_{0,0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

h_{0,0} = X * W # [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values μ

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

 $h_{0,0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

X	W	X1_mul_vector
	0.1	1*0.1 = 0.1

2 0.1 $2 \times 0.1 = 0.2$

٠

 \cdot \cdot \cdot 6_{1} 0.1_{1} $256 \times 0.1 = 25.6$

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

 $h_{0,0} = X * W \# [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

X	W	X1_mul_vector	Multi-threaded
1	0.1	1*0.1 = 0.1	Execution (256 Threads)
2	0.1	2*0.1 = 0.2	
•	* .	•	
•	•	•	
•	•	•	
256	0.1	256*0.1 = 25.6	

GPU - #1 What about Summation?

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values W = [0.1, 0.1, ..., 0.1] # Then we need to have 256 weight values $h_{0,0} = X * W # [1*0.1 + 2*0.1 + ... + 256*0.1] == 32389.6$

X	W	X1_mul_vector	Multi-threaded
1	0.1	1*0.1 = 0.1	Execution (256 Threads)
2	0.1	2*0.1 = 0.2	
•	* .	•	
•	•	•	
•	•	•	
256	0.1	256*0.1 = 25.6	

GPU - #2 Summary Step

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values

 $W = [0.1, 0.1, \ldots, 0.1]$ # Then we need to have 256 weight values

h_{0.0} = X ★ W # [1★0.1 + 2★0.1 + ... + 256★0.1] == 32389.6

GPU - #2 Summary Step

X = [1.0, 2.0, ..., 256.0] # Let's say we have 256 input values
W = [0.1, 0.1, ..., 0.1] _____# Then we need to have 256 weight values

h_{0.0} = X ★ W # [1★0.1 + 2★0.1 + ... + 256★0.1] == 32389.6

Comparing - Order of Magnitude (sequences)

Many Knobs to Tune

But the type of calculation we perform is very suited for GPUs

	_							And a state of the local division of the loc	and the second second		CO LONG CO. LONG	State of the local division of	10000	- 100 C	CLL JONE VIEW	Co. SHE WARE	and the second		a set of the set of the	and the second se		1. Contract of the second sec second second sec	
		E) - IQ			- 👰	~ @ ,		. 🔘	1		۷		Ø	•	• 4		🥏 ≤			0	1	2
																							I
				il				I	11		Ĭ					j				Î			
																			F	1 E			
				211161	1111	111		1111	11111										T.	1111	11111		
									E	E		E	E	11.0	12 13	E	E	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	18 8 8	ł		1818	12 8 1
7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24		2	3	4	5	

Summary

- GPUs == Many Threads == Great for ML Workloads
- And now you know how this works
- Fortunately, you don't need to worry about implementation details

E multi-core CPU

T

\bullet \bigcirc multi-core **GPU** CPU

Work needed: NONE

(just use a GPU build)

Beyond That

Use Distribution Strategy API

There's a talk for that (@ 1pm)

tensorflow.org/learn

TensorFlow Courses

coursera.org/learn/introduction-tensorflow

udacity.com/tensorflow

Distribution Strategies

tensorflow.org/alpha/guide/distribute_strategy

@MagnusHyttsten