

# Achieving optimal speed on TF 2

Zebroid Meeting 2020-07-21

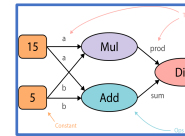
Andres Potapczynski

# Topics for today

## Relevant questions

- A What are the major changes from TF1 to TF2?
- B How does TF2 compare against TF1 and PyTorch?
- C What are the most important TF APIs for performance?

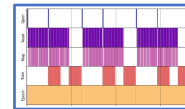
## APIs / Tools analyzed



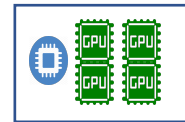
Computational Graph  
(tf.function)



Input Pipeline  
(tf.data)



Performance profiler  
(tf.profiler)



Distributed training  
(tf.distribute)



Preamble



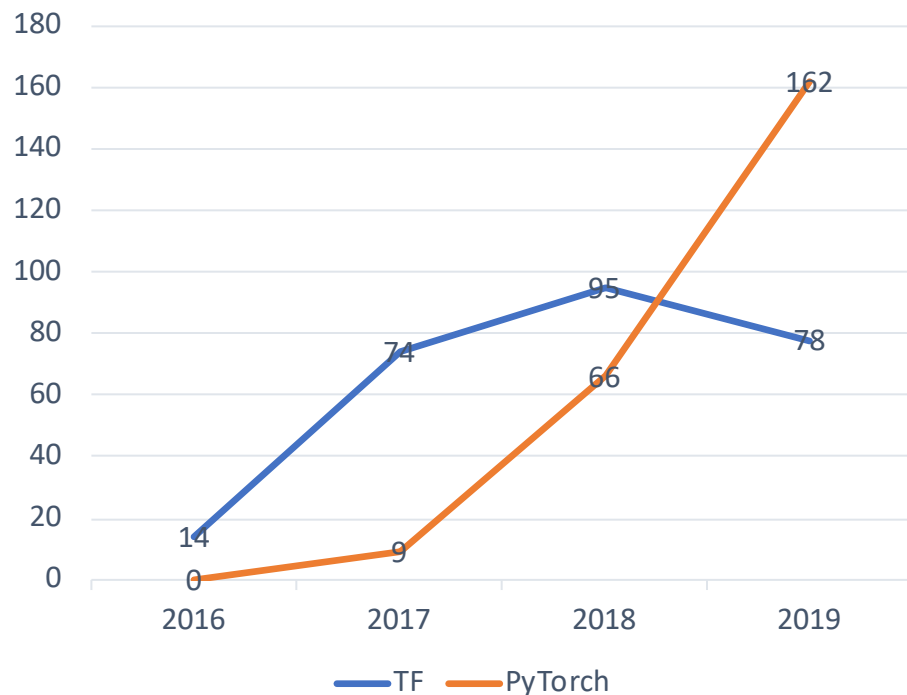
Focus

A

# Similar to PyTorch, TF 2 is easier to debug and more pythonic

## PyTorch has been gaining popularity in the community

# of NeurIPs papers with implementations on a given framework\*



## From (TF 1)\*\*

- Using graphs by default
- Defining nodes on a graph
- Constructing a graph with a fixed type in mind
- Using control dependencies to manage the execution
- Initializing variables manually
- Writing non-intuitive control flow statements and functions

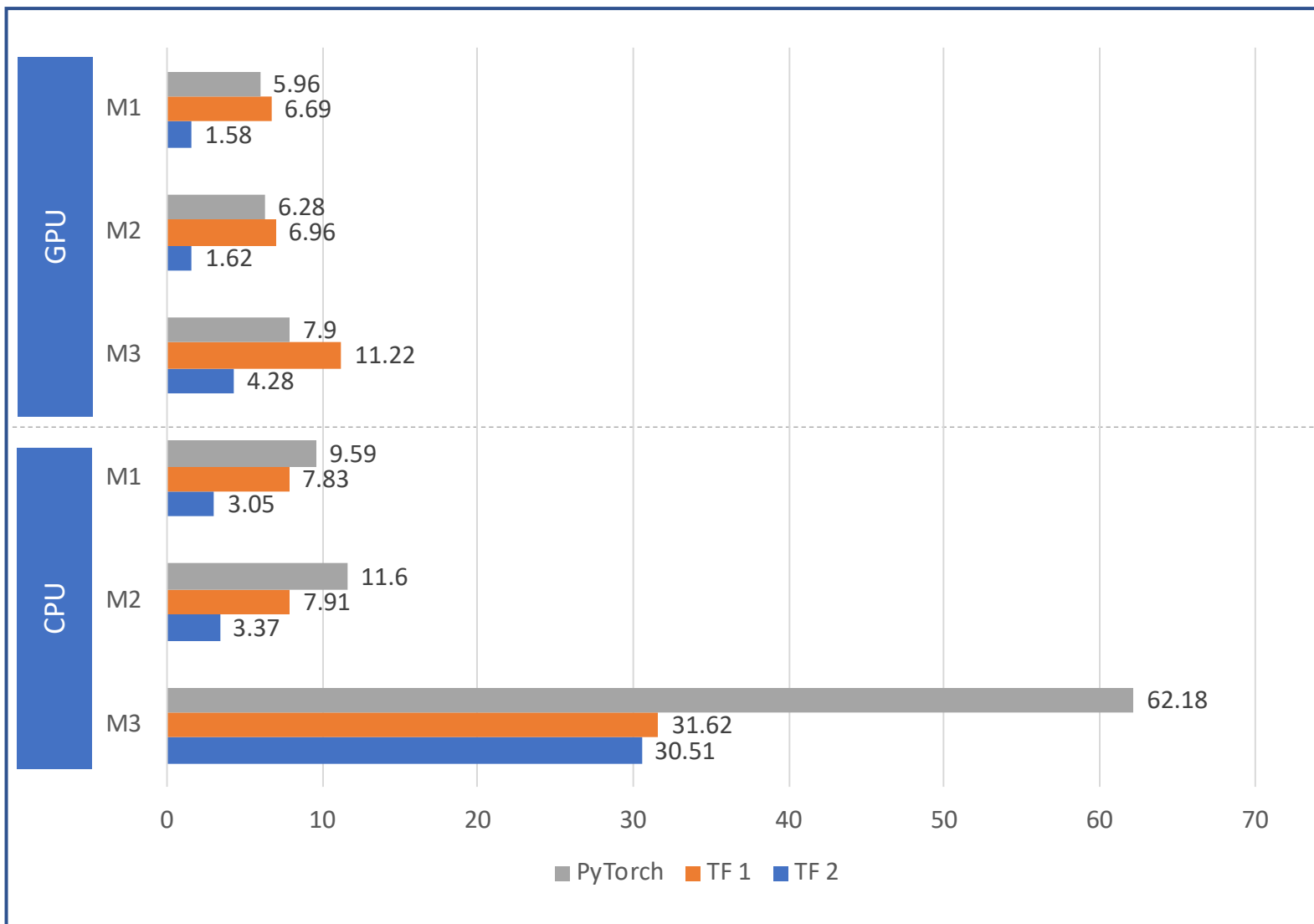
## To (TF 2)\*\*

- Running code on eager execution (enables per step debugging)
- Selecting operations to run under `tf.function`
- Having functions that "adapt" to different types
- Writing code imperatively to control order of execution
- Letting the variables be automatically initialized
- Using python syntax to define control flows (and letting Autograph do the translation)

\* <http://horace.io/pytorch-vs-tensorflow/>, \*\* `tf.function` and Autograph (TF Dev Summit 2019) talk

B

# In principle, TF 2 should not be slower than TF 1 (VAEs)



## Models on MNIST

### M1 VAE (dense linear architecture)

- Encoder: Dense(784) - Dense(400) – Dense(64)
- Decoder: Dense(64) – Dense(400) – Dense(784)

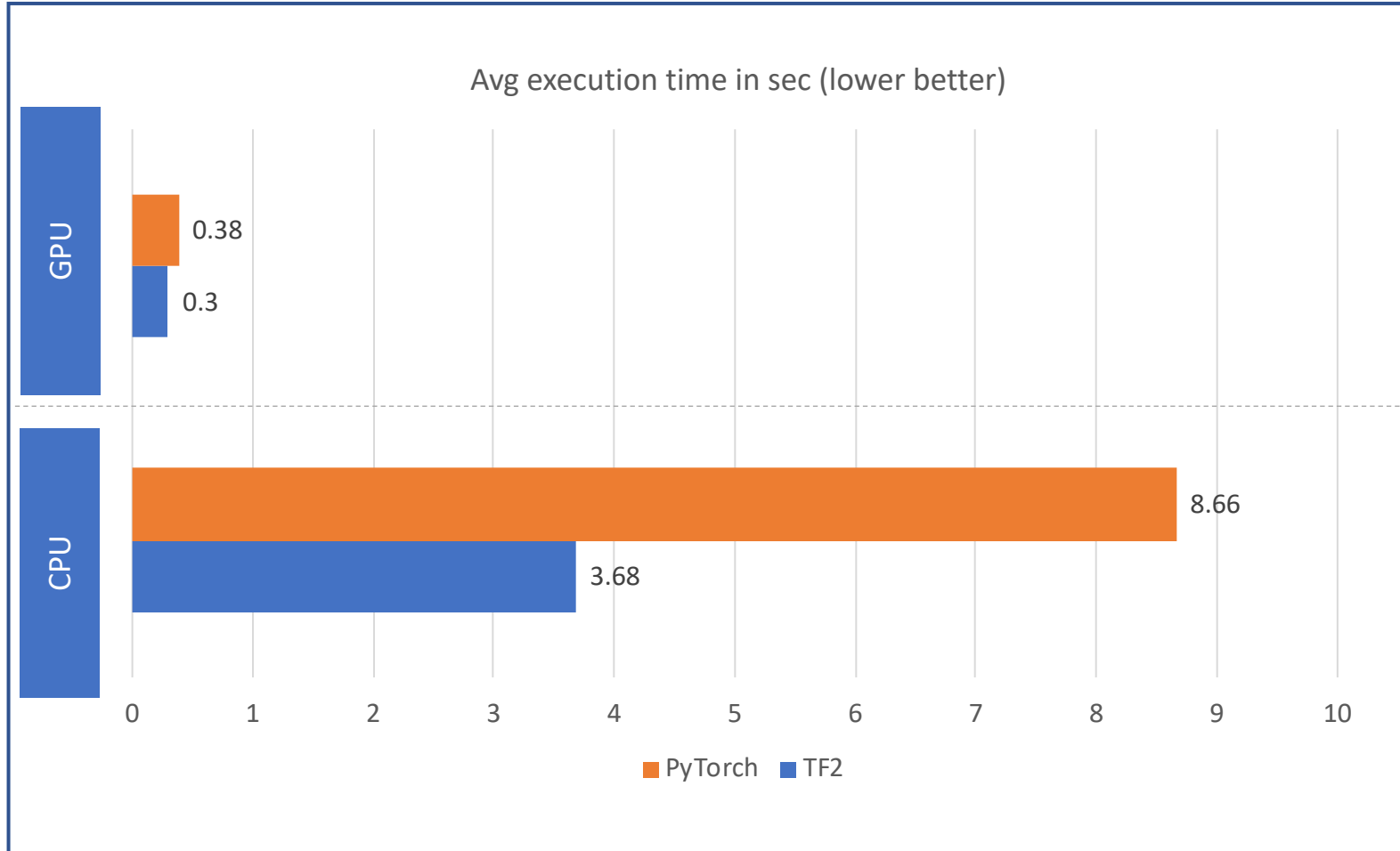
### M2 VAE (dense nonlinear architecture)

- Encoder: Dense(784) - Dense(400) – Dense(240) – Dense(64)
- Decoder: Dense(64) – Dense(240) – Dense(400) – Dense(784)

### M3 VAE (convolutional architecture)

- Encoder: Conv2D(32, 3, 2) - Conv2D(64, 3, 2) – Dense(64).
- Decoder: Dense(1568) - Conv2DT(64, 3, 2) - Conv2DT(32, 3, 2) - Conv2DT(1, 3, 1)

## B And in some cases faster than PyTorch



### Models

#### Generator

- 4 Conv2DTranspose layers with Batch Normalization and ReLU activations
- 1 Conv2DTranpose layer with Tanh at the end

#### Discriminator

- 4 Conv2D layers with Batch Normalization and LeakyReLU activations
- 1 Conv2D layer with Sigmoid at the end

# c TF 2 encapsulates several APIs for different tasks

NOT EXHAUSTIVE

|                       |   |                             |  |
|-----------------------|---|-----------------------------|--|
| High-level DL APIs    | <ul style="list-style-type: none"><li>• tf.keras</li><li>• tf.estimator</li></ul>   | Visualization               | <ul style="list-style-type: none"><li>• tf.summary</li><li>• tf.profile</li></ul>  |
| Low-level DL APIs     | <ul style="list-style-type: none"><li>• tf.nn</li><li>• tf.losses</li><li>• tf.metrics</li><li>• tf.optimizers</li><li>• tf.train</li><li>• tf.initializers</li></ul> | Deployment and optimization | <ul style="list-style-type: none"><li>• tf.distribute</li><li>• tf.saved_model</li><li>• tf.autograph</li><li>• tf.lite</li><li>• tf.quantization</li><li>• tf.tpu</li></ul> |
| Autodiff              | <ul style="list-style-type: none"><li>• tf.GradientTape</li><li>• tf.gradients</li></ul>  | Special data structures     | <ul style="list-style-type: none"><li>• tf.lookup</li><li>• tf.ragged</li><li>• tf.nest</li><li>• tf.sparse</li></ul>  |
| I/O and Preprocessing | <ul style="list-style-type: none"><li>• tf.data</li><li>• tf.feature_column</li><li>• tf.audio</li><li>• tf.image</li><li>• tf.io</li><li>• tf.queue</li></ul>        | Mathematics                 | <ul style="list-style-type: none"><li>• tf.math</li><li>• tf.linalg</li><li>• tf.signal</li><li>• tf.random</li><li>• tf.bitwise</li></ul>                                   |

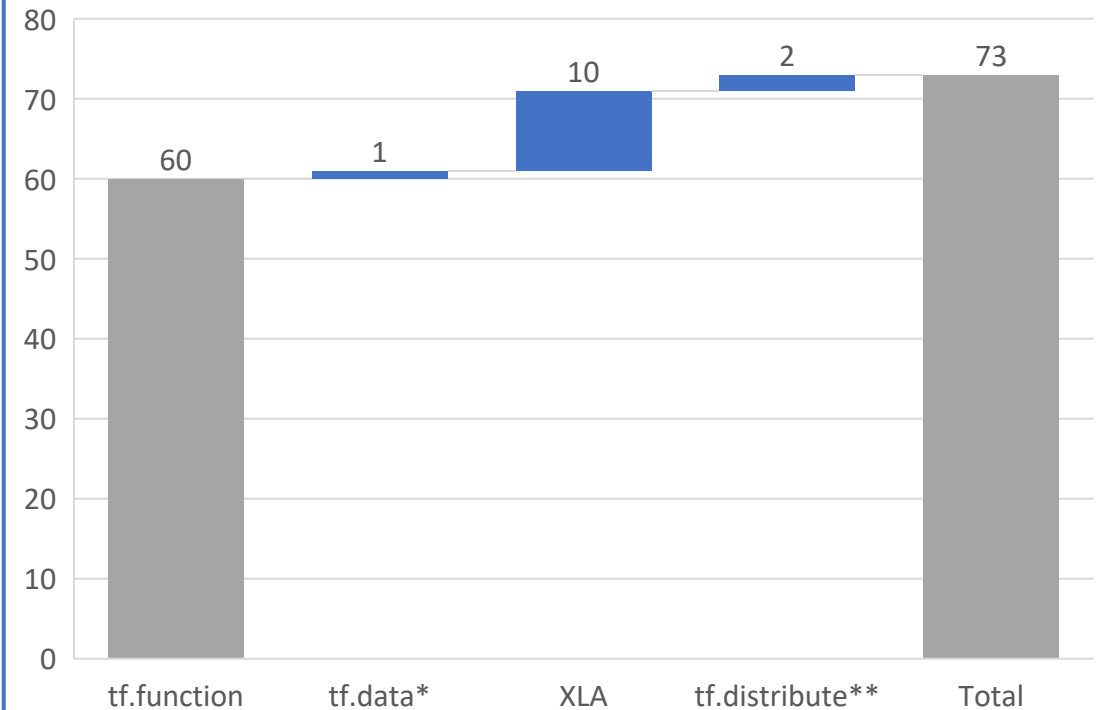
# c We focus on the APIs that drive performance

Covered

| Activity  | API / Tool   |
|---|--|
| <ul style="list-style-type: none"><li>• Ensure code is graph executable (optimize graph with XLA)</li></ul> | <ul style="list-style-type: none"><li>• <code>tf.function</code></li></ul>   |
| <ul style="list-style-type: none"><li>• Preprocess data efficiently</li></ul>                               | <ul style="list-style-type: none"><li>• <code>tf.data</code></li></ul>   |
| <ul style="list-style-type: none"><li>• Profile bottlenecks</li></ul>                                       | <ul style="list-style-type: none"><li>• <code>tf.profiler</code></li></ul>   |
| <ul style="list-style-type: none"><li>• Distribute across GPUs</li></ul>                                    | <ul style="list-style-type: none"><li>• <code>tf.distribute</code></li></ul>   |
| <ul style="list-style-type: none"><li>• Accelerate computation by reducing precision</li></ul>              | <ul style="list-style-type: none"><li>• <code>tf.keras.mixed_precision</code> / <code>tf.keras.quantize</code></li></ul> |
| <ul style="list-style-type: none"><li>• Shrink your models</li></ul>  | <ul style="list-style-type: none"><li>• <code>tf.keras.prune</code></li></ul>  |

## Even small models can benefit from all the optimizations

M3 VAE on MNIST, % reduction in execution time over baseline



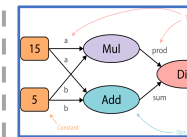
\* Data optimization consists of caching and prefetching, \*\* Even for small models, `tf.distribute` yields some improvement despite its overhead. Distributed across 2 Tesla P4 GPUs

# Topics for today

## Relevant questions

- A What are the major changes from TF1 to TF2?
- B How does TF2 compare against TF1 and PyTorch?
- C What are the most important TF APIs for performance?

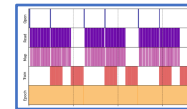
## APIs / Tools analyzed



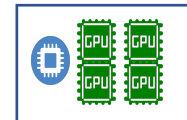
Computational Graph  
(tf.function)



Input Pipeline  
(tf.data)



Performance profiler  
(tf.profiler)



Distributed training  
(tf.distribute)

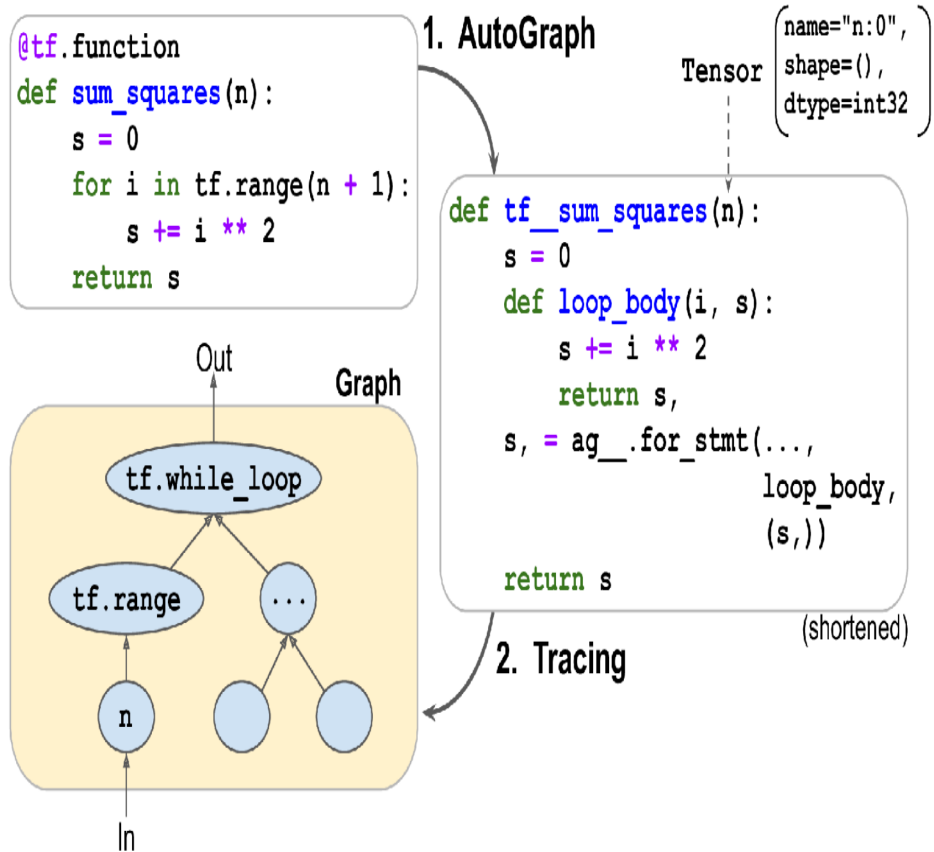


Next



# tf.function requires some rules to compile properly

## Translation of python code to graph by decoration



## Common mistakes

- Writing all the code and then trying to decorate it
- Calling external libraries (like NumPy)
- Multiple function retracing from feeding python objects
- Returning non tensor objects
- Creating variables on every call
- Using the tensor's shape value (specially batch size)
- Putting decorators on all functions

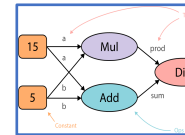
Lack of examples and documentation might be hampering the adoption of TF 2!

# Topics for today

## Relevant questions

- A What are the major changes from TF1 to TF2?
- B How does TF2 compare against TF1 and PyTorch?
- C What are the most important TF APIs for performance?

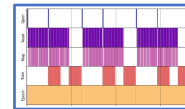
## APIs / Tools analyzed



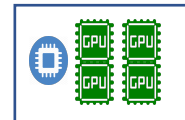
Computational Graph  
(tf.function)



Input Pipeline  
(tf.data)



Performance profiler  
(tf.profiler)



Distributed training  
(tf.distribute)



Next

# tf.data simplifies the input pipeline

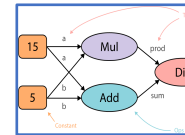
| Standard code example |  |
|-----------------------|--|
| Extract               | <pre>files = tf.data.Dataset.list_files(file_pattern) ds = tf.data.TFRecordDataset(files)</pre>  |
| Transform             | <pre>ds = ds.shuffle(buffer_size) ds = ds.repeat(epochs) ds = ds.map(pre_fn, num_parallel_calls=tf.data.experimental.AUTOTUNE) ds = ds.batch(batch_size) ds = ds.cache() ds = ds.prefetch(tf.data.experimental.AUTOTUNE)</pre> |
| Load                  | <pre>iterator = ds.make_one_shot_iterator() x = iterator.get_next()</pre>  |

# Topics for today

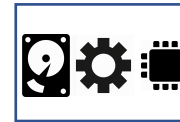
## Relevant questions

- A What are the major changes from TF1 to TF2?
- B How does TF2 compare against TF1 and PyTorch?
- C What are the most important TF APIs for performance?

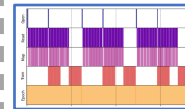
## APIs / Tools analyzed



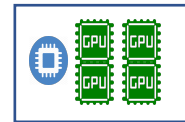
Computational Graph  
(tf.function)



Input Pipeline  
(tf.data)



Performance profiler  
(tf.profiler)



Distributed training  
(tf.distribute)



Next

# tf.profiler helps your model execute faster

NOT EXHAUSTIVE

EXAMPLE

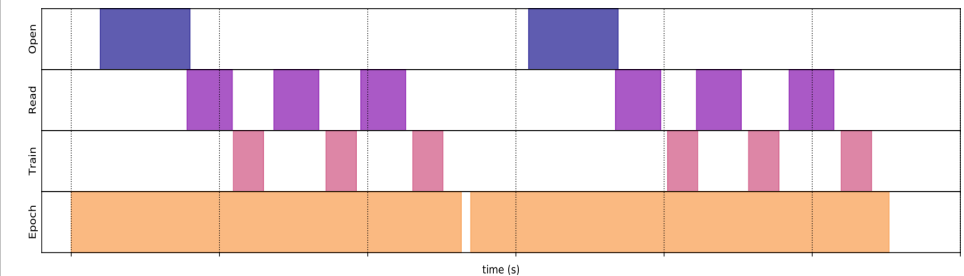
## Benefits (tools)

- Shows input bottlenecks (input pipeline analyzer)
- Guides how to speed up parts of training by exhibiting the cost of their ops (TF Stats)
- Breaks down ops per device and exhibits if they are waiting for input

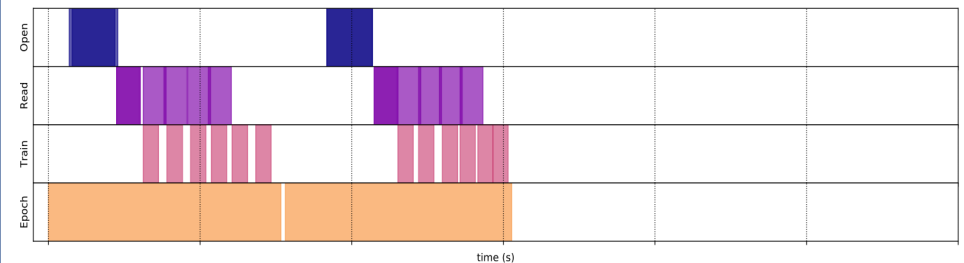
## Best Practices\*

- Use parallel calls on map, prefetch and cache input
- Utilize the devices more by increasing batch size
- Calculate metrics every few steps and reduce callbacks
- Reduce precision fp16 and make dimension divisible by 8
- Send data to multiple devices in parallel
- Use `tf.name_scope` to identify most costly ops in python construct

## Inadequate trace



## Optimized trace

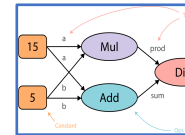


# Topics for today

## Relevant questions

- A What are the major changes from TF1 to TF2?
- B How does TF2 compare against TF1 and PyTorch?
- C What are the most important TF APIs for performance?

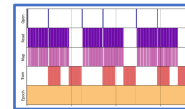
## APIs / Tools analyzed



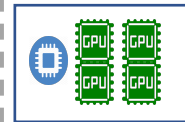
Computational Graph  
(tf.function)



Input Pipeline  
(tf.data)



Performance profiler  
(tf.profiler)

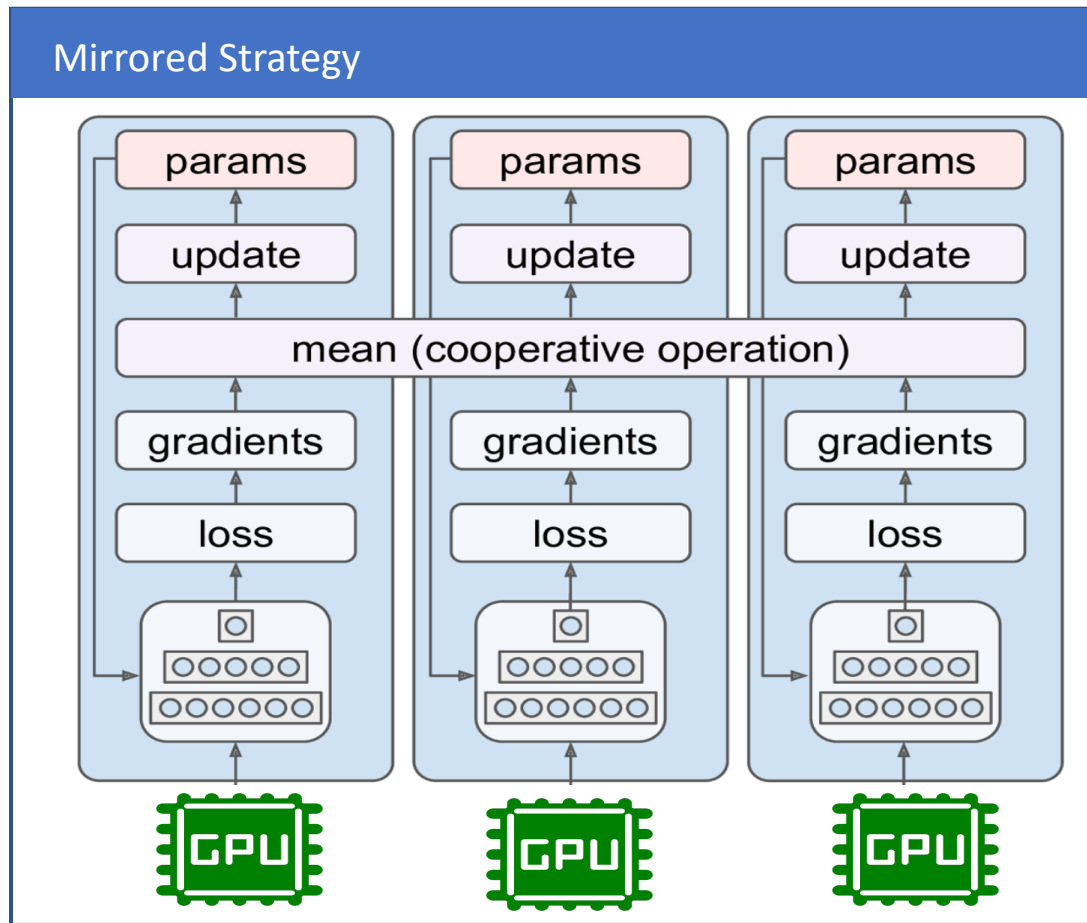


Distributed training  
(tf.distribute)



Next

# tf.distribute is easy to use since TF 2 is “strategy aware”



## Code changes

- Change data pipeline by distributing it
- Create model and optimizer under the strategy scope
- Define loss on a per example basis, create distributed train step function and aggregate loss and gradients using global batch size