| 76,000,000 | 80,000+ | 13,000+ | 2,400+ |
|:---:|:---:|:---:|:---:|
| downloads | commits | pull requests | contributors |

# TensorFlow
## Ecosystem

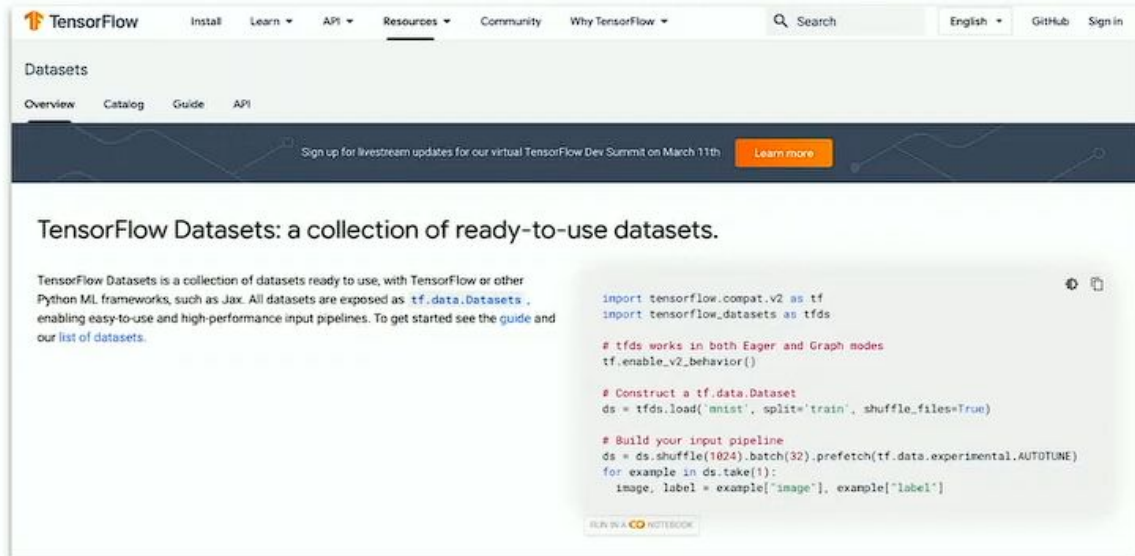| From research | To production | Deployed everywhere |

Empowering Responsible AI

Powered by the community

# Eager at the core and simple, performant data input pipelines in 2.x

eager execution, supporting numpy arrays

tf.data

TensorFlow Datasets

# Add-ons and extensions to the TensorFlow ecosystem

- TF Probability
- TF Graphics
- Mesh TensorFlow
- TF Model Garden
- TF Agents
- TF Text
- Swift for TensorFlow
- Sonnet
- JAX
- Neural Structured Learning
- TF Quantum
- …and more on tensorflow.org!

# TensorFlow 2.1 supports Cloud TPUs

# Challenge

**How do I use it?**

**Is it safe?**

**Is it fair?**

**Is it the latest version?**

tfhub.dev

# TensorFlow Hub

A comprehensive collection of models

Image

Text

Video

Audio

# Ready to use

Pre-trained models ready for transfer learning on your own datasets and deployable anywhere you want

**TensorFlow**
Extended

**TensorFlow**
.JS

**TensorFlow**
Lite

**Coral**

**Problem domain**

Image feature vector

**Architecture**

MobileNet V2

**Publisher**

Google

**Dataset**

ImageNet (ILSVRC-2012-CLS)

Format: TF2.0 Saved Model

Fine tunable: Yes

License: Apache-2.0 ⧉

Last updated: 2020-02-20

# Model formats

| Saved Model | .JS (v1, default) | .JS (v2, default) | .JS (v3, default) |

## Want to use this model?

To use this model, take a look at the example code, or at our user guide.

You can also try out the associated Colab.

**Copy URL to clipboard**

**Download Model**

**Open Colab Notebook**

Asset size: 2.62MB

# TF2 SavedModel

This is a SavedModel in TensorFlow 2 format. Using it requires TensorFlow 2 (or 1.15) and TensorFlow Hub 0.5.0 or newer.

# Overview

MobileNet V2 is a family of neural network architectures for efficient on-device image classification and related tasks, originally published by

- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, Liang-Chieh Chen: "Inverted Residuals and Linear Bottlenecks: Mobile Networks for Classification, Detection and Segmentation", 2018.

Mobilenets come in various sizes controlled by a multiplier for the depth (number of features) in the convolutional layers. They can also be trained for various sizes of input images to control inference speed.

# Powered by the community

Built, trained, and deployed already by the TensorFlow community

- DeepMind
- Google
- Microsoft AI for Earth
- NVIDIA
- The Metropolitan Museum of Art
- Global Biodiversity Information Facility
- Kaggle
- And more...

# Kaggle competition now supports 2.x

No setup

No provisioning

Datasets in optimal formats

TPUs and GPUs provided at no cost to users

# What is Kaggle?

- Community
- Competitions
- Datasets
- Notebooks
- Courses

> 4 million registered users

# Kaggle hosts machine learning competitions

# TF 2.1 + Kaggle == easy acceleration

- No setup
- No provisioning
- Datasets in optimal formats
- TPUs and GPUs provided at no cost to users

# Check out Tensorflow on Kaggle!

- Check out the code above
  - https://www.kaggle.com/philculliton/a-simple-tf-2-1-notebook
- Compete in the Flower Classification competition
  - https://www.kaggle.com/c/flower-classification-with-tpus
- Try TPUs with TF 2.1 in Notebooks
  - https://www.kaggle.com/notebooks

# Better experimentation with TensorFlow

## TensorBoard.dev

Upload and share your ML experiments with anyone



## Performance Profiler

Available in TensorBoard, Profiler provides overview of model performance and better debugging guidance

TF is speed!

# Life of a performance engineer

Capture Profile → Aggregate and analyze results → Optimize model → Perf OK 🙂

(Most boring and time consuming)

Repeat

# ML productivity

```
┌──────────────┐    ┌──────────────┐    ┌──────────────┐
│              │    │              │    │              │   Perf OK
│   Capture    │ >> │  Automated   │ >> │   Optimize   │ >>  🙂
│   Profile    │    │ performance  │    │    model     │
│              │    │  guidance!   │    │              │
└──────────────┘    └──────────────┘    └──────────────┘
      ▲                                        │
      │   (Used extensively inside Google to tune products)   Repeat
      └──────────────────────────────────────┘
```

CAPTURE PROFILE

Runs (29)

train/2020-03-02_14:46:29

Tools (5)

overview_page

Hosts (1)

localhost

## Performance Summary

**Average Step Time**
*lower is better*
*(σ = 66.7 ms)*                                    1805.1 ms

- **All Others Time**
  *(σ = 57.4 ms)*                                  717.5 ms

- **Compilation Time**
  *(σ = 0.0 ms)*                                   0.0 ms

- **Output Time**
  *(σ = 0.0 ms)*                                   0.0 ms

- **Input Time**
  *(σ = 0.1 ms)*                                   0.3 ms

- **Kernel Launch Time**
  *(σ = 21.6 ms)*                                  327.3 ms

- **Host Compute Time**
  *(σ = 2.3 ms)*                                   6.7 ms

- **Device to Device Time**
  *(σ = 0.0 ms)*                                   0.0 ms

- **Device Compute Time**
  *(σ = 25.8 ms)*                                  753.3 ms

**Device Compute Precisions**
*(out of Total Device Time)*
- 16-bit: 46.3
- 32-bit: 53.7

## Step-time Graph

**Step Time (in milliseconds)**

Legend:
- All others
- Compilation
- Output
- Input
- Kernel launch
- Host compute
- Device to device
- Device compute

*Step Number*

## Recommendation for Next Step

- Your program is NOT input-bound because only 0.0% of the total step time sampled is waiting for input. Therefore, you should focus on reducing other time.

- 18.1 % of the total step time sampled is spent on Kernel Launch.

- 39.7 % of the total step time sampled is spent on All Others time.

TensorBoard    SCALARS    GRAPHS    PROFILE                          INACTIVE

CAPTURE PROFILE

**ON HOST: TOTAL SELF-TIME (GROUPED BY TYPE)**
*(in microseconds) of a TensorFlow operation*

**ON HOST: TOTAL SELF-TIME**
*(in microseconds) of a TensorFlow operation*

Runs (29)

train/2020-03-02_14:46:29

- CollectiveReduce
- SoftmaxCrossEntropyWi...
- Dataset
- OneHot
- ParseExampleV2
- Other

41.6%  30.9%  21.9%

- Iterator::Model::Shard::P...
- Iterator::Model::Shard::P...
- Iterator::Model::Shard::P...
- replica_7/model/tf_op_la...
- model/tf_op_layer_Softm...
- replica_4/model/tf_op_la...
- replica_6/model/tf_op_la...
- replica_2/model/tf_op_la...
- replica_3/model/tf_op_la...

13.2%  9.7%

1/4

Tools (5)

tensorflow_stats

Hosts (1)

localhost

## TensorFlow operations

Host/device          Type          Operation

Note: To avoid sluggishness, only the 1000 most time-consuming operations out of the total 36340 operations are shown in the table below.

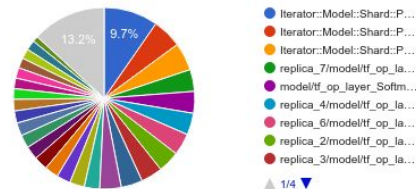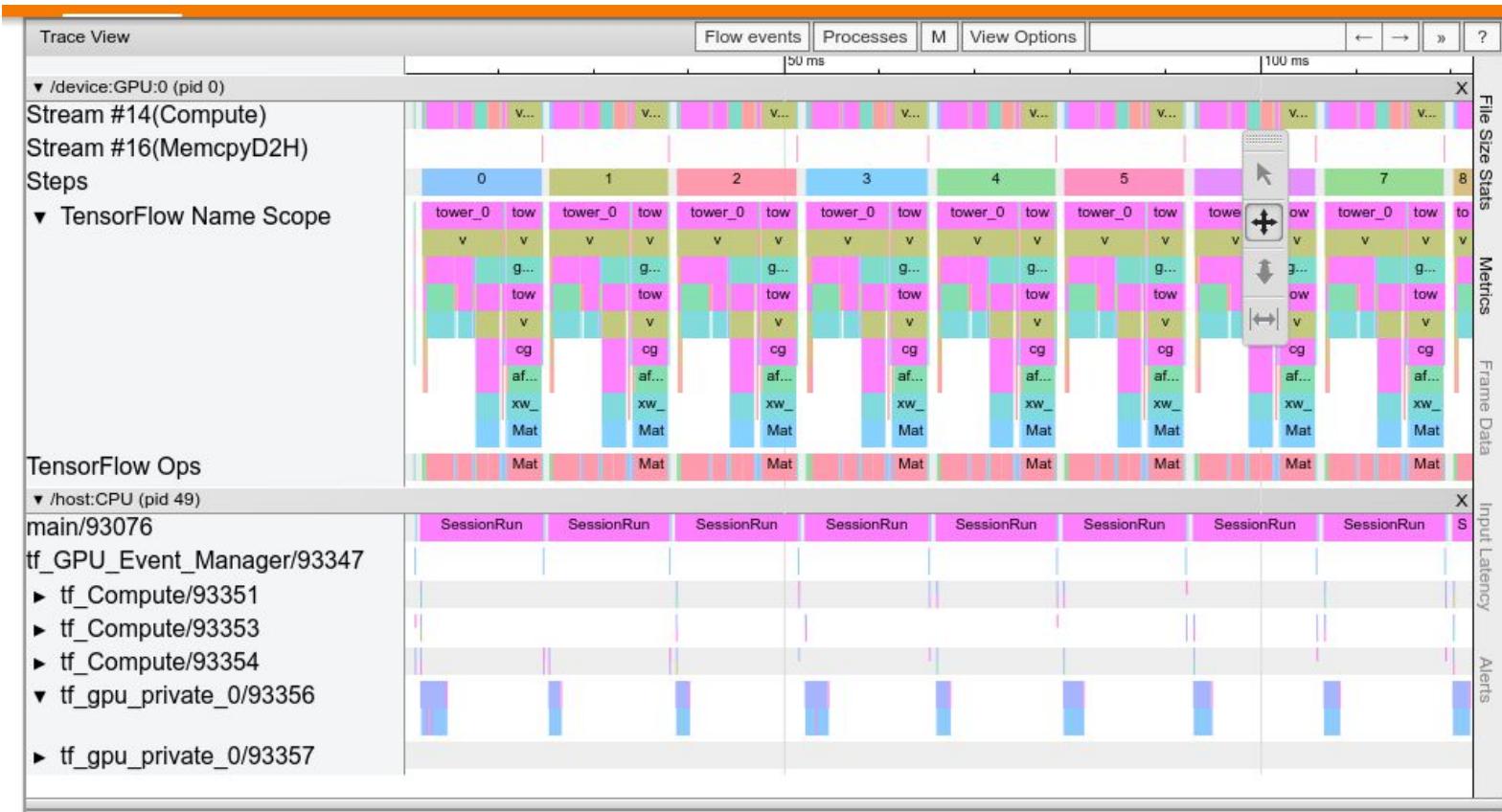| Rank | Host/device | Type | Operation | #Occurrences | Total time (us) | Avg. time (us) | Total self-time (us) | Avg. self-time (us) | Total self-time on Device (%) | Cumulative total-self time on Device (%) | Total self-time on Host (%) | Cumulative total-self time on Host (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 43 | Device | MatMul | gradient_tape/replica_3/model/transformer_v2/Transformer/decode/embedding_shared_weights_1/presoftmax_linear/MatMul/MatMul | 11 | 149,200 | 13,564 | 149,200 | 13,564 | 0.3% | 66.3% | 0% | 0% |
| 44 | Device | MatMul | gradient_tape/replica_6/model/transformer_v2/Transformer/decode/embedding_shared_weights_1/presoftmax_linear/MatMul/MatMul | 11 | 149,199 | 13,564 | 149,199 | 13,564 | 0.3% | 66.5% | 0% | 0% |
| 45 | Device | MatMul | gradient_tape/replica_5/model/transformer_v2/Transformer/decode/embedding_shared_weights_1/presoftmax_linear/MatMul/MatMul | 11 | 149,179 | 13,562 | 149,179 | 13,562 | 0.3% | 66.8% | 0% | 0% |
| 46 | Device | MatMul | gradient_tape/model/transformer_v2/Transformer/decode/embedding_shared_weights_1/presoftmax_linear/MatMul/MatMul | 11 | 149,163 | 13,560 | 149,163 | 13,560 | 0.3% | 67.1% | 0% | 0% |

# TF 2 Profiler Tool Set

Overview Page   Input Pipeline Analyzer   TensorFlow Stats   Trace Viewer   +   4 GPU/TPU expert tools

**Available Today on TensorBoard**

# Next Steps

- Tutorial:
  https://www.tensorflow.org/tensorboard/tensorboard_profiling_keras
- Guide: https://tensorflow.org/guide/profiler
- Github: https://github.com/tensorflow/profiler
- Two related talks in this afternoon:
  - "Scaling TensorFlow data processing with tf.data"
  - "Scaling TensorFlow 2 models to multi-worker GPUs"

Thank you!

TensorFlow
DEV SUMMIT 2020