# ImageNet is the new MNIST

## Chris Ying

Research SWE @ Google Brain
g.co/brain

on behalf of **many** people across Google

# Goal: "Interactive ML supercomputing"

- Hardware
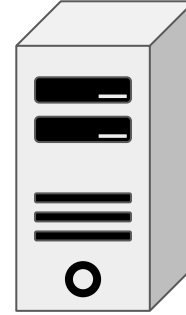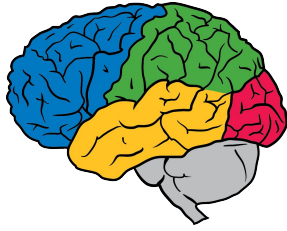  - Cloud TPUs
  - TPU pods

- Software
  - TensorFlow Datasets, Layers, and Estimator APIs (open-source)
  - XLA compiler (open-source) with TPU backend

- Research
  - Understanding of generalization gap
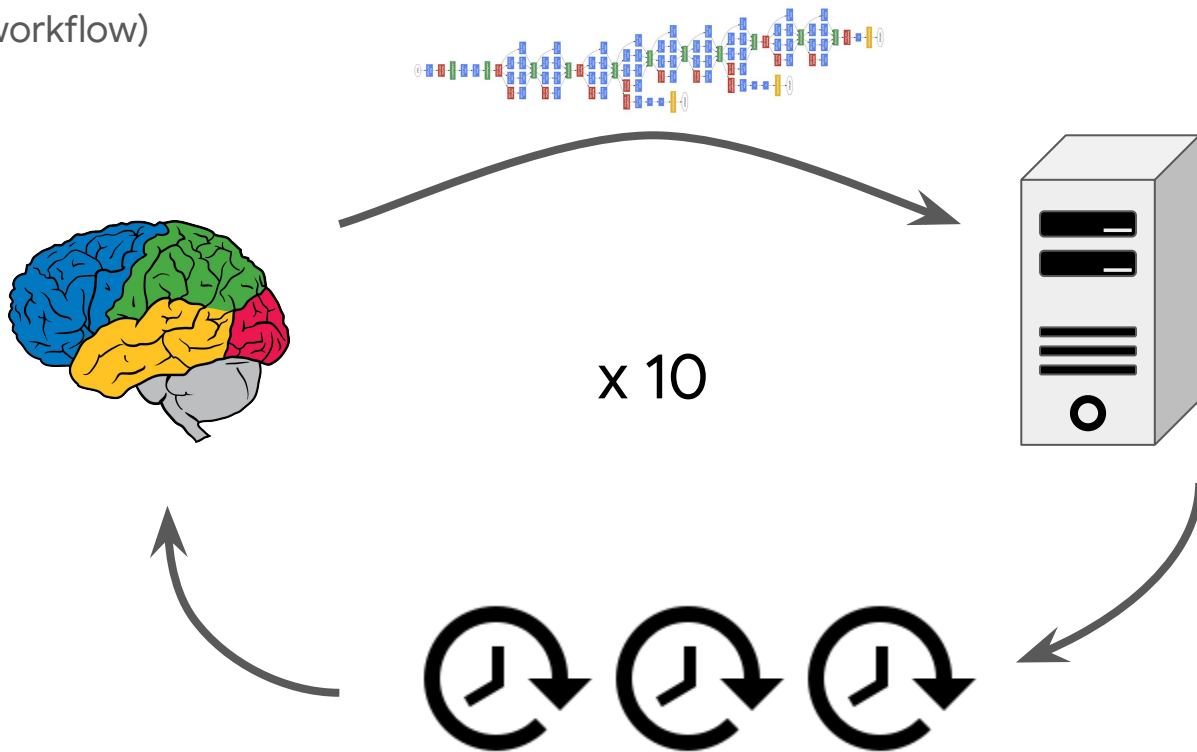  - Large-batch training advances

# Motivation

(classical workflow)
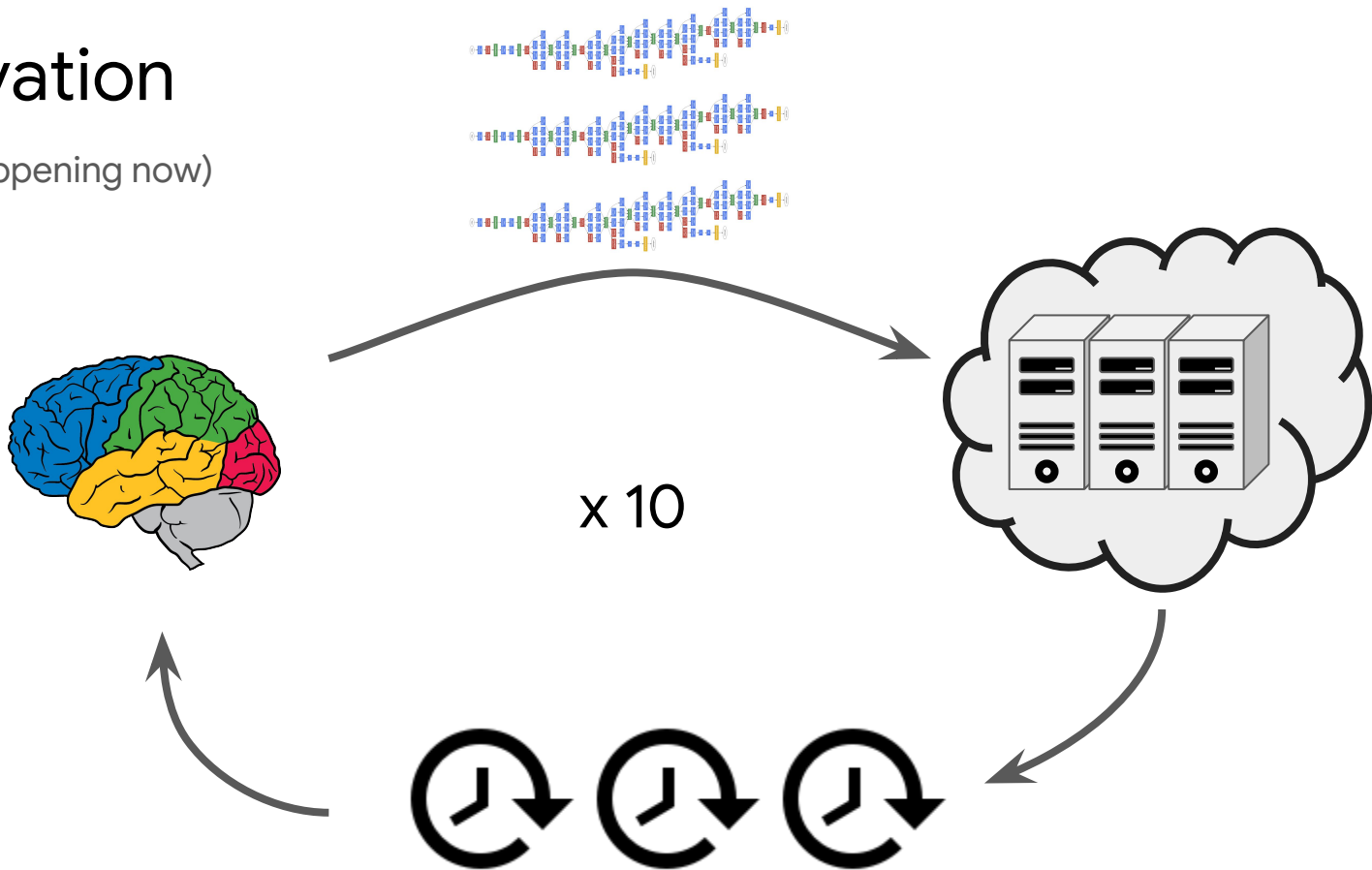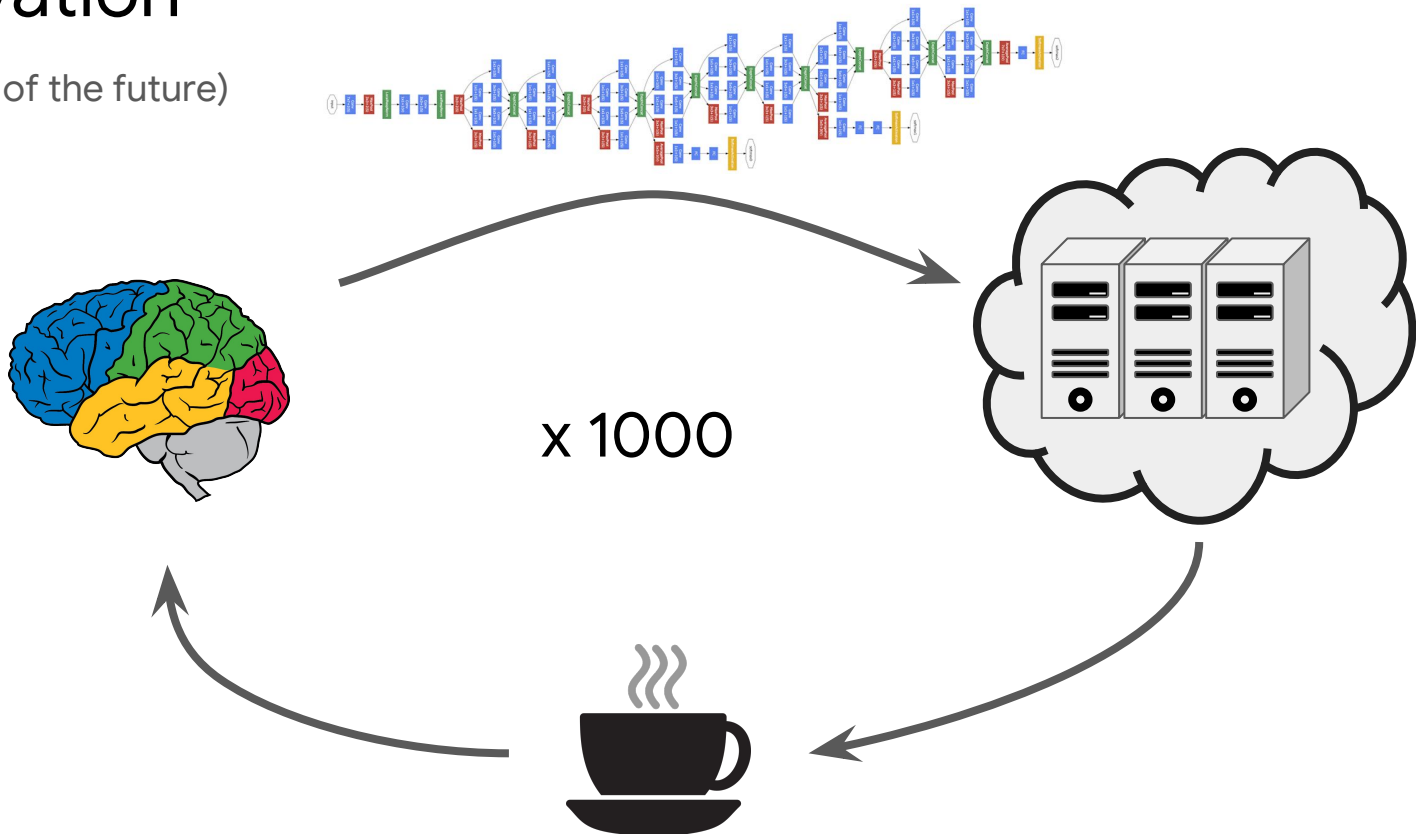
# Motivation

(classical workflow)



x 10

# Motivation

(what's happening now)

x 10

# Motivation

(our vision of the future)

x 1000

# ImageNet is the new MNIST



**MNIST**: 60,000 B&W images



**ImageNet**: 1,281,167 color images

# Motivating results

ResNet-50-v2 on ImageNet

| # of TPU devices | Batch size | Time to 90 epochs | Accuracy |
|---|---|---|---|
| 1 | 256 | 23 hours 22 minutes | **76.6%** |
| 4 | 1024 | 5 hours 48 minutes | 76.3% |
| 16 | 4096 | 1 hour 30 minutes | 76.5% |
| 32 | 8192 | **45 minutes** | 76.1% |
| 64 | 16384 | 22 minutes | 75.0% |

Only change between different runs is batch size (linearly scale LR) and hardware, <u>no model changes or hyperparameter re-tuning</u>!

# Cloud TPU



180 TFLOPS of computation, 64 GB of HBM memory, 2400 GB/s mem BW

Cloud TPU

# TPUv2 Chip



- 45 TFLOPS
- 16 GB of HBM
- 600 GB/s mem BW
- Vector unit: float32
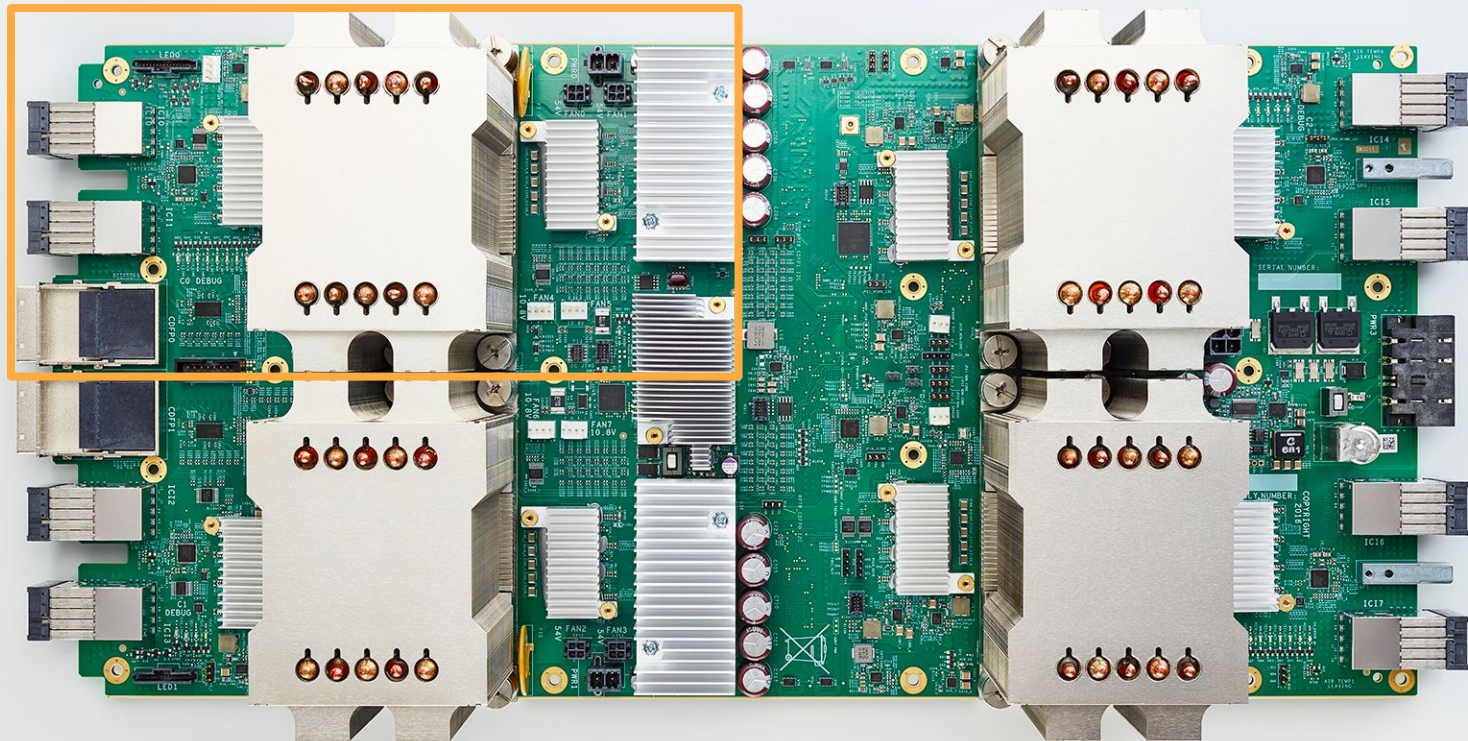- Scalar unit: float32
- Matrix unit (MXU): float32 input/output, reduced precision multiplication



HBM
8 GB

core

vector unit

scalar unit

MXU
128x128

core

vector unit

scalar unit

MXU
128x128

HBM
8 GB

# TPUv2 Chip

**Matrix Unit**

128x128 systolic array
float32 results*

* 10 GB of HBM
* 600 GB/s mem BW
* Scalar unit: 32b float
* MXU: 32b float accumulation but reduced precision for multipliers
* 45 TFLOPS

* reduced precision multiplication

HBM
8 GB

core

| vector unit | scalar unit |

MXU
128x128

core

| vector unit | scalar unit |

MXU
128x128

HBM
8 GB

# Matrix Unit Systolic Array

Computing y = Wx

Toy example: 3x3 systolic array
W = 3x3 matrix
batch_size(x) = 3

# Matrix Unit Systolic Array

## Computing y = Wx

with W = 3x3, batch_size(x) = 3

inputs

weights

$X_{33}$

$X_{32}$  $X_{23}$

$X_{31}$  $X_{22}$  $X_{13}$

$X_{21}$  $X_{12}$

Matrix Unit (MXU)

$W_{11}$ $X_{11}$  $W_{12}$  $W_{13}$

$W_{21}$  $W_{22}$  $W_{23}$

$W_{31}$  $W_{32}$  $W_{33}$

accumulation

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

$X_{33}$

$X_{32}$  $X_{23}$

$X_{31}$  $X_{22}$  $X_{13}$

weights

Matrix Unit (MXU)

| $W_{11}$ $X_{21}$ | $W_{12}X_{12}$ + $W_{11}X_{11}$ | $W_{13}$ |
| $W_{21}$ $X_{11}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}$ |

accumulation

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

$X_{33}$

$X_{32}$  $X_{23}$

weights

Matrix Unit (MXU)

| $W_{11}$ $X_{31}$ | $W_{12}X_{22}$ + $W_{11}X_{21}$ | $W_{13}X_{13}$ + ... |
| $W_{21}$ $X_{21}$ | $W_{22}X_{12}$ + $W_{21}X_{11}$ | $W_{23}$ |
| $W_{31}$ $X_{11}$ | $W_{32}$ | $W_{33}$ |

accumulation

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

$X_{33}$

outputs

Matrix Unit (MXU)

weights

| | | |
|---|---|---|
| $W_{11}$ | $W_{12}X_{32}$ + $W_{11}X_{31}$ | $W_{13}X_{23}$ + ... |
| $W_{21}$ $X_{31}$ | $W_{22}X_{22}$ + $W_{21}X_{21}$ | $W_{23}X_{13}$ + ... |
| $W_{31}$ $X_{21}$ | $W_{32}X_{12}$ + $W_{31}X_{11}$ | $W_{33}$ |

$Y_{11} = W_{11}X_{11} + W_{12}X_{12} + W_{13}X_{13}$

accumulation

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

weights

Matrix Unit (MXU)

| $W_{11}$ | $W_{12}$ | $W_{13}X_{33}$ $+$ $...$ |
| $W_{21}$ | $W_{22}X_{32}$ $+$ $W_{21}X_{31}$ | $W_{23}X_{23}$ $+$ $...$ |
| $W_{31}$ $X_{31}$ | $W_{32}X_{22}$ $+$ $W_{31}X_{21}$ | $W_{33}X_{13}$ $+$ $...$ |

accumulation

outputs

$Y_{21} = W_{11}X_{21} + W_{12}X_{22} + W_{13}X_{23}$

$Y_{11} = W_{11}X_{11} + W_{12}X_{12} + W_{13}X_{13}$

$Y_{12} = W_{21}X_{11} + W_{22}X_{12} + W_{23}X_{13}$

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

weights

outputs

**Matrix Unit (MXU)**

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
| $W_{21}$ | $W_{22}$ | $W_{23}X_{33}$ + ... |
| $W_{31}$ | $W_{32}X_{32}$ + $W_{31}X_{31}$ | $W_{33}X_{23}$ + ... |

$Y_{31} = W_{11}X_{31} + W_{12}X_{32} + W_{13}X_{33}$

$Y_{21} = W_{11}X_{21} + W_{12}X_{22} + W_{13}X_{23}$

$Y_{11} = W_{11}X_{11} + W$

$Y_{22} = W_{21}X_{21} + W_{22}X_{22} + W_{23}X_{23}$

$Y_{12} = W_{21}X_{11} + W_{22}X_{12} + W_{23}X_{13}$

$Y_{13} = W_{31}X_{11} + W_{32}X_{12} + W_{33}X_{13}$

accumulation

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

**inputs**

**weights**

**Matrix Unit (MXU)**

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
| $W_{21}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}X_{33}$ + ... |

**accumulation**

**outputs**

$Y_{31} = W_{11}X_{11} + W_{12}X_{12} + W_{13}X_{13}$

$Y_{21} = W_{11}X_{21} + W$

$Y_{32} = W_{21}X_{31} + W_{22}X_{32} + W_{23}X_{33}$

$Y_{22} = W_{21}X_{21} + W_{22}X_{22} + W_{23}X_{23}$

$Y_{12} = W_{21}X_{11} + W$

$Y_{23} = W_{31}X_{21} + W_{32}X_{22} + W_{33}X_{23}$

$Y_{13} = W_{31}X_{11} + W_{32}X_{12} + W_{33}X_{13}$

# Matrix Unit Systolic Array

## Computing y = Wx
with W = 3x3, batch_size(x) = 3

inputs

weights

outputs

Matrix Unit (MXU)

| $W_{11}$ | $W_{12}$ | $W_{13}$ |
| $W_{21}$ | $W_{22}$ | $W_{23}$ |
| $W_{31}$ | $W_{32}$ | $W_{33}$ |

accumulation

$Y_{31} = W_{11}X_{11} + W$

$Y_{32} = W_{21}X_{31} + W_{22}X_{32} + W_{23}X_{33}$

$Y_{22} = W_{21}X_{21} + W$

$Y_{33} = W_{31}X_{31} + W_{32}X_{32} + W_{33}X_{33}$

$Y_{23} = W_{31}X_{21} + W_{32}X_{22} + W_{33}X_{23}$

$Y_{13} = W_{31}X_{11} + W$

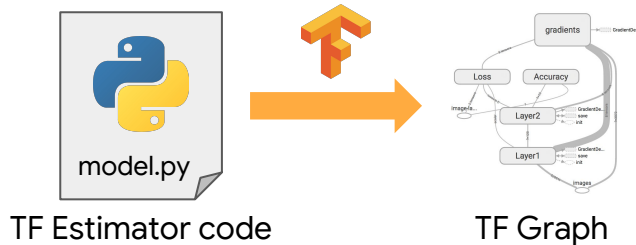Cloud TPU Pod
64 Cloud TPUs in 2-D toroidal mesh
11.5 petaFLOPS
4 terabytes of HBM memory

# Accelerated Linear Algebra (XLA)

- JIT / AOT compiler for linear algebra
- Targets multiple backends, e.g. CPUs, GPUs, and TPUs
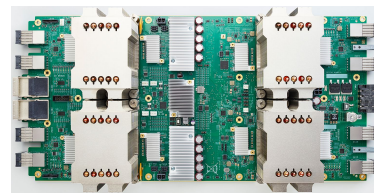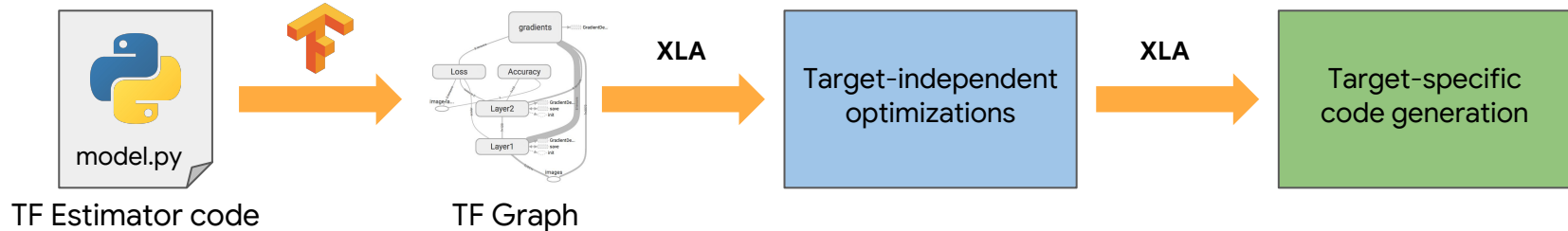- Compiler, runtime, and accelerator-specific optimizer

The life of a neural network:



TF Estimator code                    TF Graph

# Accelerated Linear Algebra (XLA)

- JIT / AOT compiler for linear algebra
- Targets multiple backends, e.g. CPUs, GPUs, and TPUs
- Compiler, runtime, and accelerator-specific optimizer

The life of a neural network:



TF Estimator code · TF Graph · **XLA** · Target-independent optimizations · **XLA** · Target-specific code generation
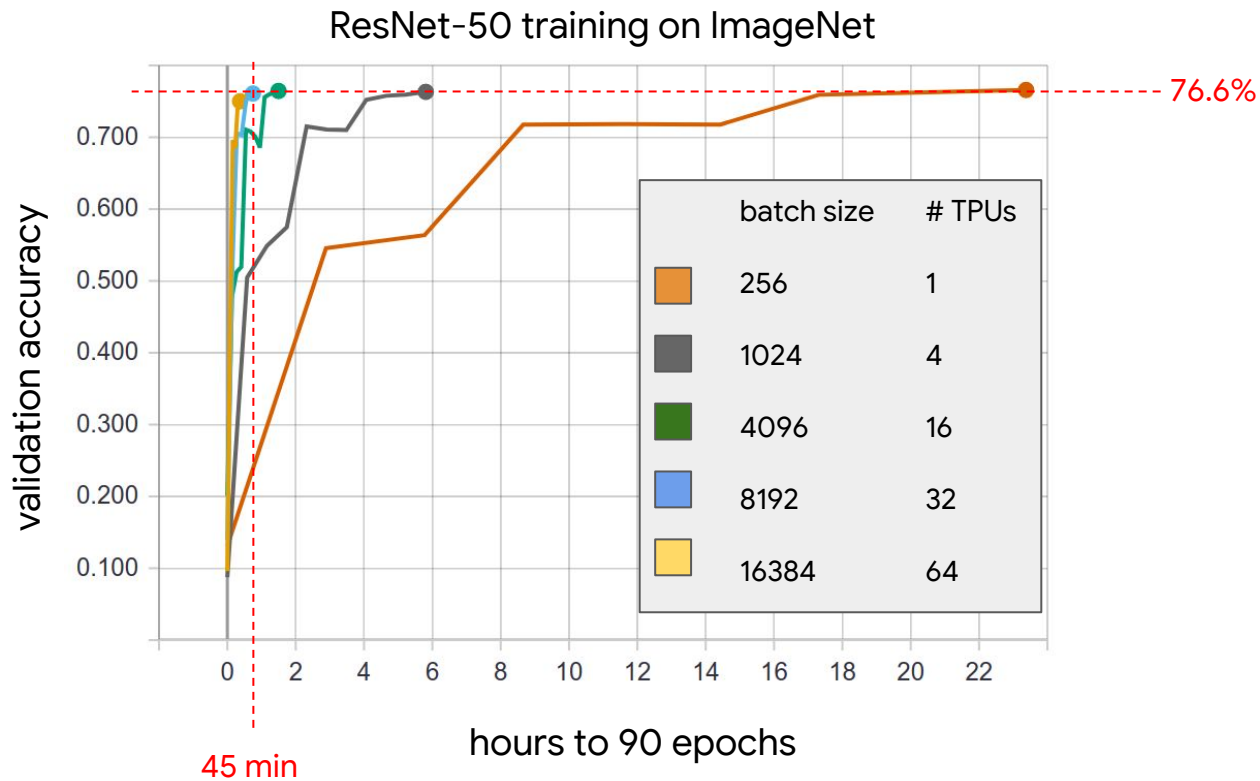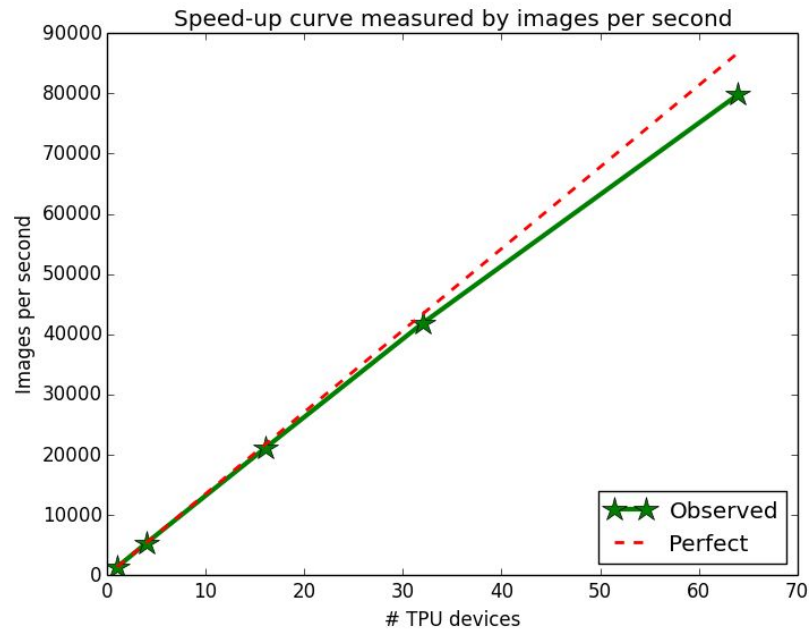
# Large batch training

- Understanding generalization gap (2016 N. Keskar et. al., 2017 E. Hoffer et. al.)

- Relationship of batch size and noise scale (2018 S. Smith et. al.)

- Learning rate scaling and schedule (2017 P. Goyal et. al.)

- New optimizers
    - K-FAC*: approximate Fisher information matrix (2015 J. Martens)
    - Neumann*:  approximate inverse Hessian (2018 S. Krishnan et. al.)
    - LARS: per-layer learning rate (2018 Y. You et. al.)

* stick around after this talk to hear more about these!

# Experiments



ResNet-50 training on ImageNet

| batch size | # TPUs |
|---|---|
| 256 | 1 |
| 1024 | 4 |
| 4096 | 16 |
| 8192 | 32 |
| 16384 | 64 |

validation accuracy

hours to 90 epochs

45 min

76.6%

# Experiments



Speed-up curve measured by images per second

# Experiments

| # of TPU devices | Batch size | Time to 90 epochs | Accuracy |
|---|---|---|---|
| 32 | 8192 | 44.9 minutes | **76.1%** |
| 64 | 8192 | 29.8 minutes | 75.7% |
| 64 | 16384 | 22.3 minutes | 75.0% |
| 64 | 65536 | 17.5 minutes | 65.4% |
| 64 | 8192 → 16384[1] | **29.5 minutes** | **76.1%** |

Only change between different runs is batch size (linearly scale LR) and hardware,
no model changes or hyperparameter re-tuning!
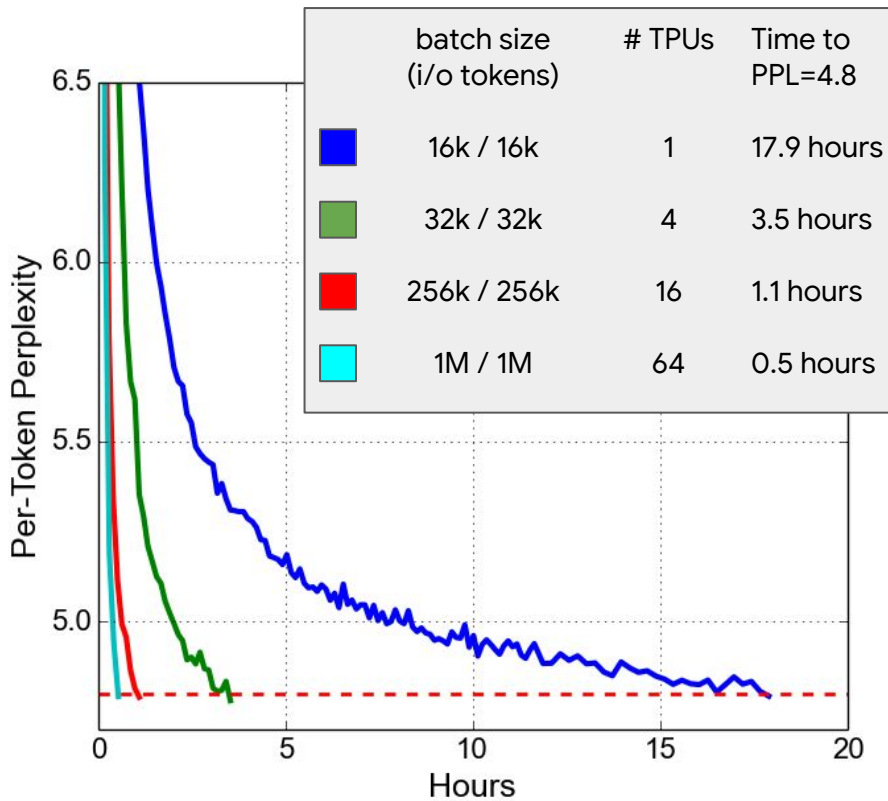
[1] Don't Decay the Learning Rate, Increase the Batch Size  (2018 S. Smith  et. al)

# More than just ImageNet

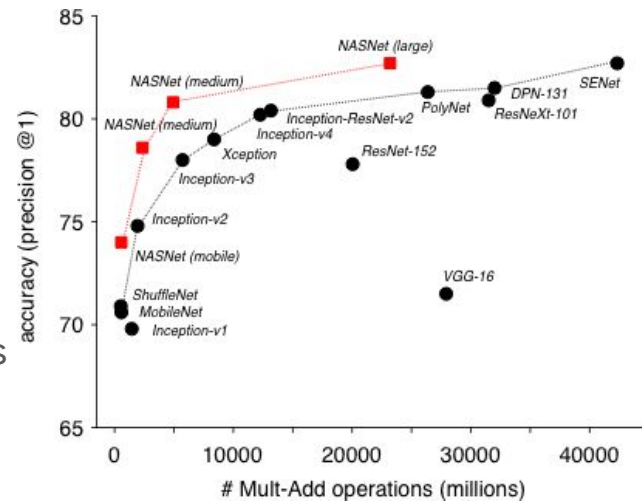Transformer model from "Attention is All You Need" (2017 A. Vaswani et. al.)

WMT'14 English–German translation task

Adam optimizer - same learning rate schedule across configurations



| | batch size (i/o tokens) | # TPUs | Time to PPL=4.8 |
|---|---|---|---|
| 🟦 | 16k / 16k | 1 | 17.9 hours |
| 🟩 | 32k / 32k | 4 | 3.5 hours |
| 🟥 | 256k / 256k | 16 | 1.1 hours |
| 🟦 | 1M / 1M | 64 | 0.5 hours |

# Implications



- Faster training enables neural architecture search
  - Reinforcement learning architectures beat existing models in accuracy and cost [1]
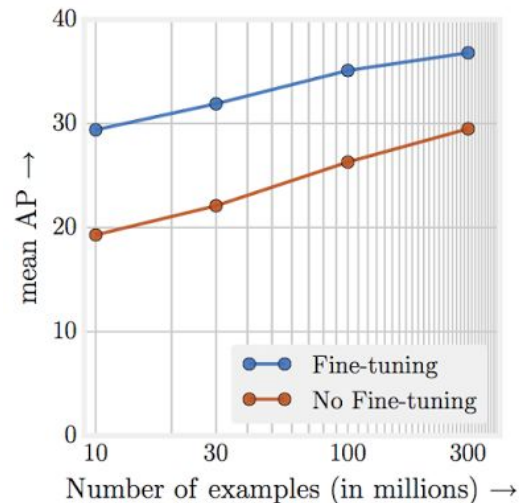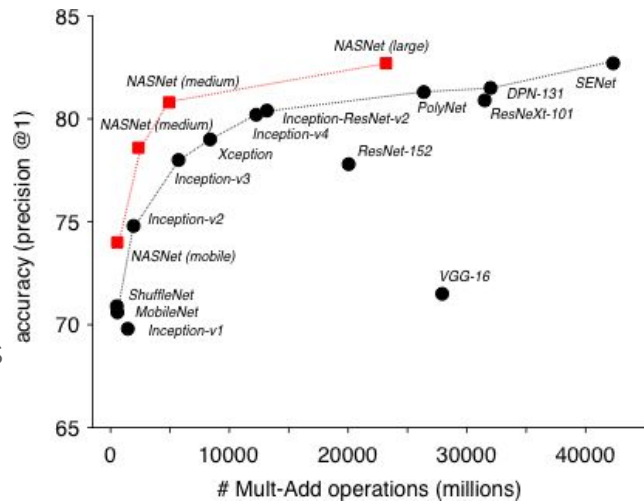
[1] Learning Transferable Architectures for Scalable Image Recognition (2017 B. Zoph et. al)

# Implications

- Faster training enables neural architecture search
  - Reinforcement learning architectures beat existing models in accuracy and cost [1]



- What's the "new ImageNet"?
  - Full ImageNet (14M), Open Images (9M), YouTube-8M
  - Performance increases logarithmically with data [2]



[1] Learning Transferable Architectures for Scalable Image Recognition (2017 B. Zoph et. al)

[2] Revisiting Unreasonable Effectiveness of Data in Deep Learning Era (2017 C. Sun et. al)

Pieter-jan

Brennan

Sam

Jonathan

Sameer

Zak

Quoc

Bjarke

Noam

Naveen

Chris

# Thank you!

chrisying@google.com

g.co/brain
g.co/tpusignup