

Large Scale Deep Learning with TensorFlow

Jeff Dean
Google Brain Team
g.co/brain

In collaboration with **many** other people at Google

Background

Google Brain project started in 2011, with a focus on pushing state-of-the-art in neural networks. Initial emphasis:

- use large datasets, and
- large amounts of computation



to push boundaries of what is possible in perception and language understanding



Overview

- Cover our experience from past ~5 years
 - **Research:** speech, images, video, robotics, language understanding, NLP, translation, optimization algorithms, unsupervised learning, ...
 - **Production:** deployed systems for advertising, search, GMail, Photos, Maps, YouTube, speech recognition, image analysis, user prediction, ...
- Focus on neural nets, but many techniques more broadly applicable



What is the Google Brain Team?

- Research team focused on long term artificial intelligence research
 - Mix of computer systems and machine learning research expertise
 - Pure ML research, and research in context of emerging ML application areas:
 - robotics, language understanding, healthcare, ...

[g.co/brain](https://www.google.com/brain)



We Disseminate Our Work in Many Ways

- By publishing our work
 - See papers at research.google.com/pubs/BrainTeam.html
- By releasing TensorFlow, our core machine learning research system, as an open-source project
- By releasing implementations of our research models in TensorFlow
- By collaborating with product teams at Google to get our research into real products

What Do We Really Want?

- Build artificial intelligence algorithms and systems that learn from experience
- Use those to solve difficult problems that benefit humanity



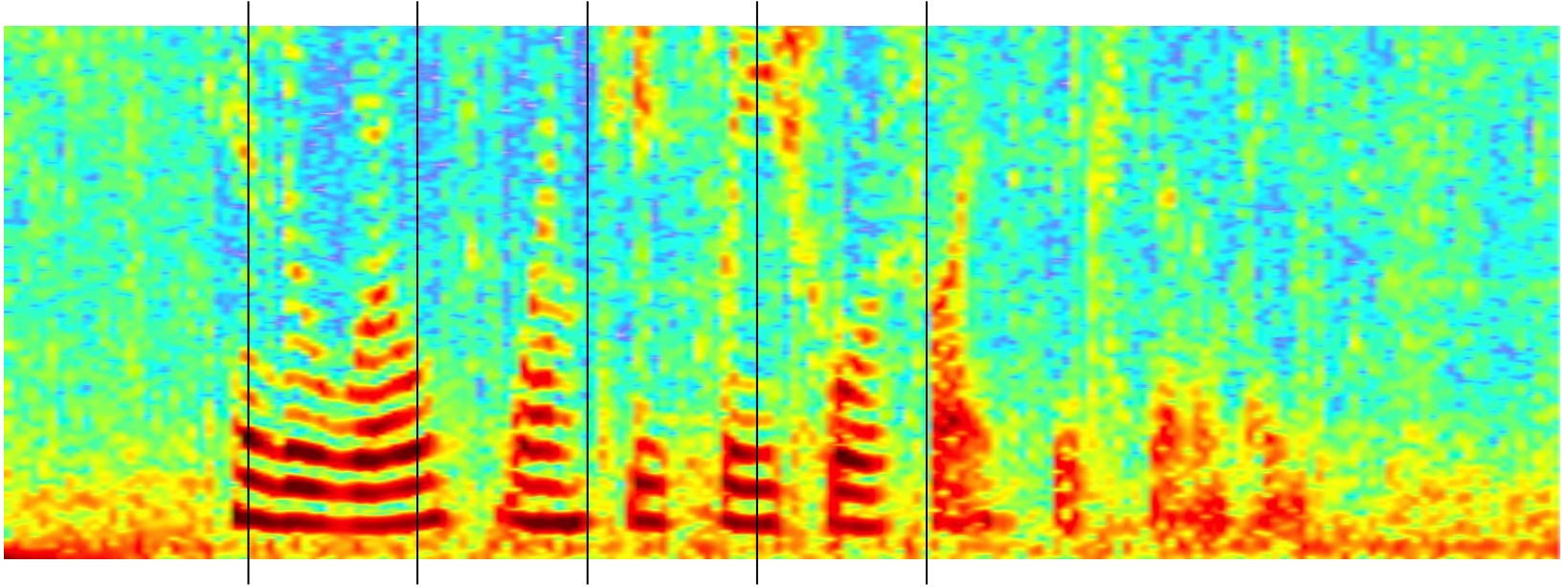
What do I mean by understanding?



What do I mean by understanding?



What do I mean by understanding?



What do I mean by understanding?

Query

[car parts for sale]

What do I mean by understanding?

Query

[car parts for sale]

Document 1

... car parking available for a small fee.
... parts of our floor model inventory for sale.

Document 2

Selling all kinds of automobile and pickup truck parts, engines, and transmissions.

Example Needs of the Future

- *Which of these eye images shows symptoms of diabetic retinopathy?*

- *Find me all rooftops in North America*

- *Describe this video in Spanish*

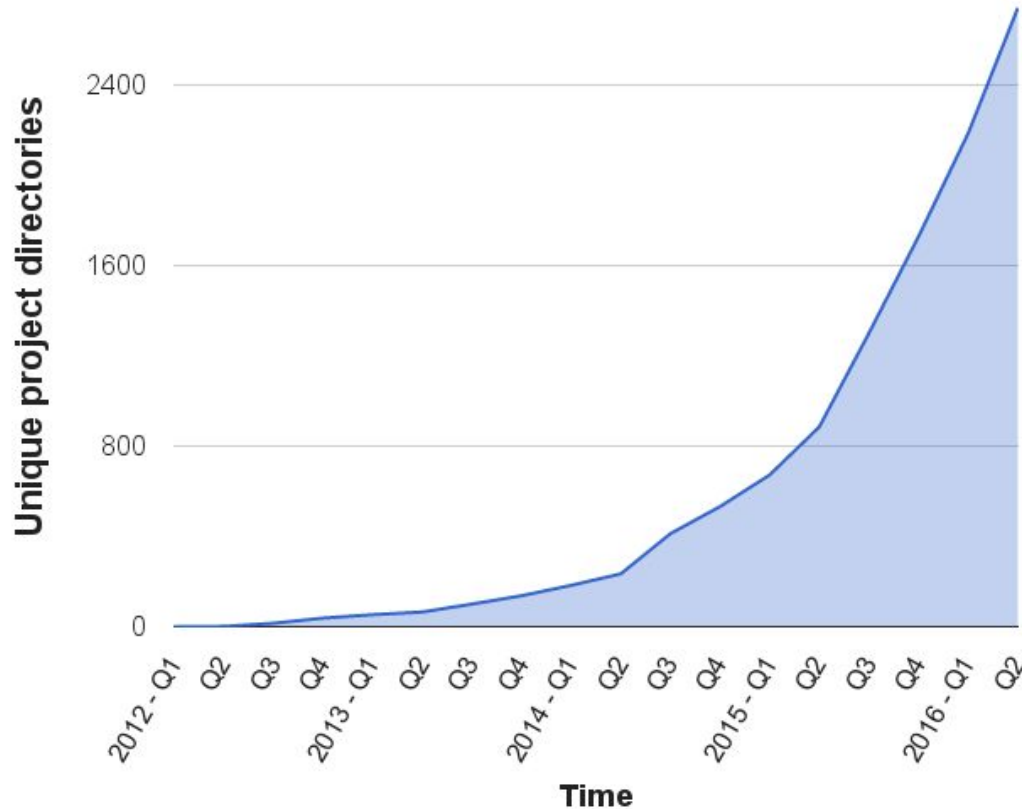
- *Find me all documents relevant to reinforcement learning for robotics and summarize them in German*

- *Find a free time for everyone in the Smart Calendar project to meet and set up a videoconference*

- *Robot, please fetch me a cup of tea from the snack kitchen*

Growing Use of Deep Learning at Google

of directories containing model description files



Across many products/areas:

- Android
- Apps
- drug discovery
- Gmail
- Image understanding
- Maps
- Natural language understanding
- Photos
- Robotics research
- Speech
- Translation
- YouTube
- ... many others ...



Overview

- Discuss *TensorFlow*, an open source machine learning system
 - Our primary research and production system
 - Show real examples
 - Explain what's happening underneath the covers



Two Generations of Distributed ML Systems

1st generation - DistBelief (Dean *et al.*, NIPS 2012)

- Scalable, good for production, but not very flexible for research

2nd generation - TensorFlow (see [tensorflow.org](https://www.tensorflow.org) and whitepaper 2015, [tensorflow.org/whitepaper2015.pdf](https://www.tensorflow.org/whitepaper2015.pdf))

- Scalable, good for production, but also flexible for variety of research uses
- Portable across range of platforms
- Open source w/ Apache 2.0 license



Important Property of Neural Networks

Results get better with

**more data +
bigger models +
more computation**



Need Both Large Datasets & Large, Powerful Models

“Scaling Recurrent Neural Network Language Models”, Williams et al.

2015

arxiv.org/pdf/1502.00512v1.pdf

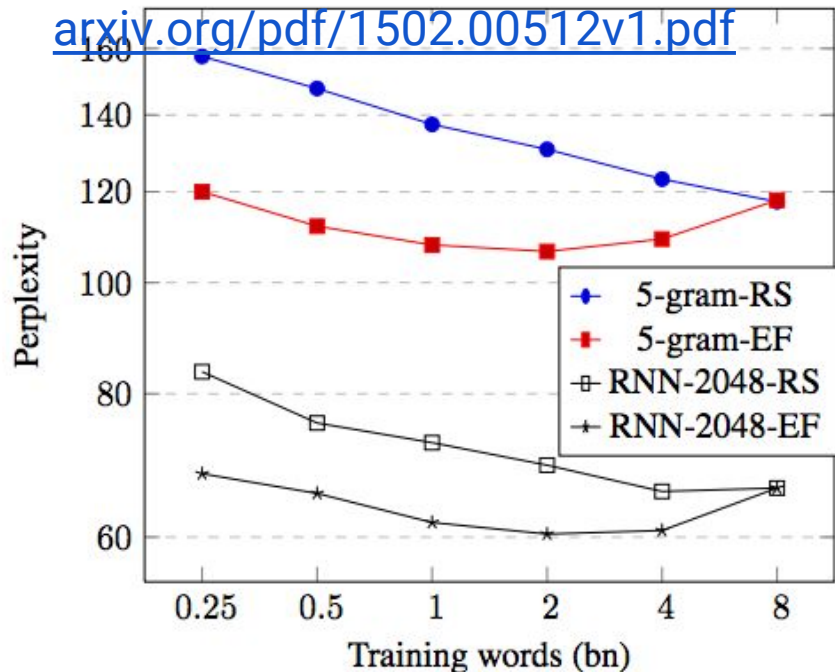


Fig. 2. Performance of 5-grams against $nstate$ 2048 RNNs with increasing training data size. We test on Randomly Selected (RS) splits and Entropy Filtered (EF) splits of the 8bn corpus.

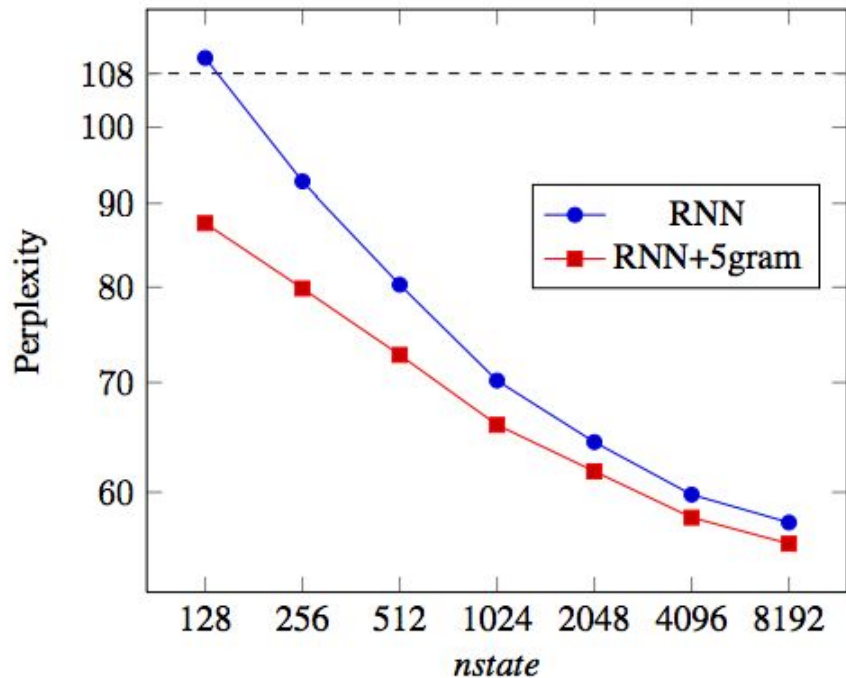


Fig. 3. Scaling $nstate$ trained on 1bn words of the entropy filtered 8bn corpus. Dashed line is the 5-gram baseline.

Large Datasets + Powerful Models

- Combination works incredibly well
- Poses interesting systems problems, though:
 - Need lots of computation
 - Want to train and do experiments quickly
 - Large-scale parallelism using distributed systems really only way to do this at very large scale
 - Also want to easily express machine learning ideas



Basics of Deep Learning

- Unsupervised cat
- Speech
- Vision
- General trend is towards more complex models:
 - Embeddings of various kinds
 - Generative models
 - Layered LSTMs
 - Attention



Learning from Unlabeled Images



- Train on 10 million images (YouTube)
- 1000 machines (16,000 cores) for 1 week.
- 1.15 billion parameters



Learning from Unlabeled Images



Top 48 stimuli from the test set



Optimal stimulus
by numerical optimization



Learning from Unlabeled Images



Top 48 stimuli from the test set



Optimal stimulus
by numerical optimization



Adding Supervision



Top stimuli for selected neurons.



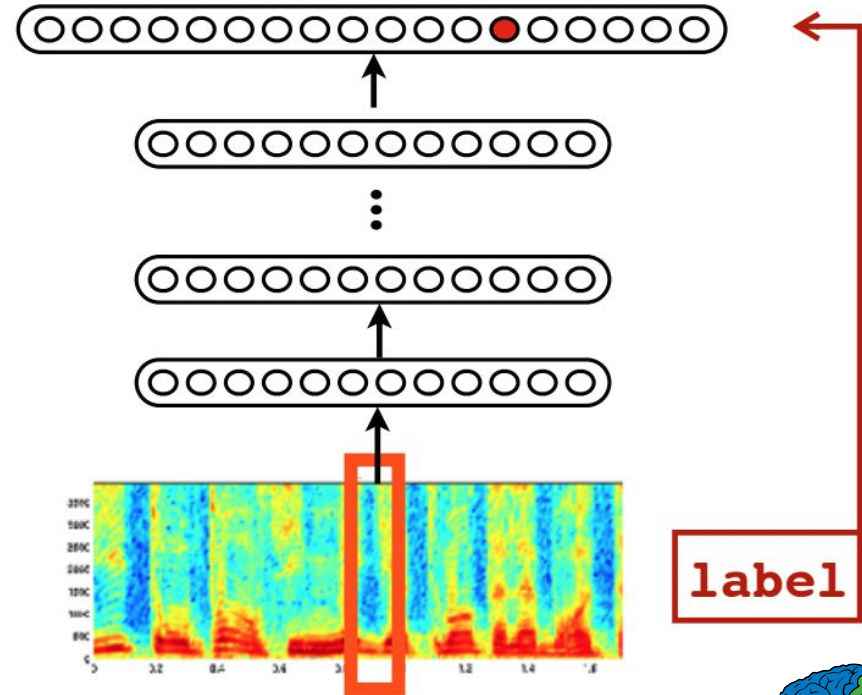
Speech: Feedforward Acoustic Models

Model speech frame-by-frame,
independently

Simple fully-connected networks

Deep Neural Networks for Acoustic Modeling in Speech Recognition

Hinton *et al.* IEEE Signal
Processing Magazine, 2012



CLDNNs

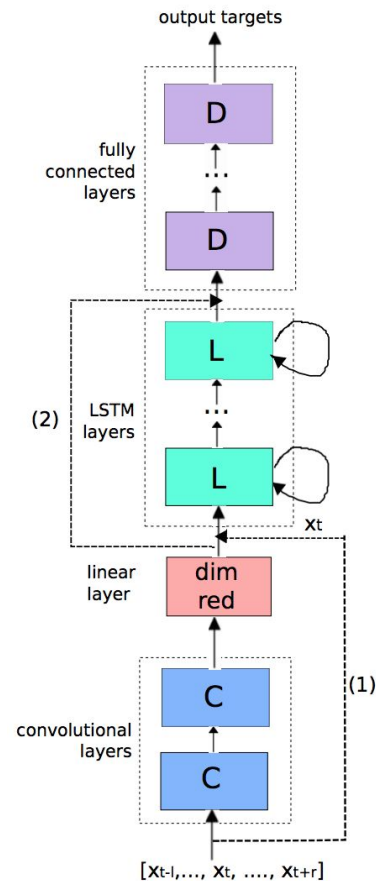
Model frequency invariance using 1D convolutions

Model time dynamics using an LSTM

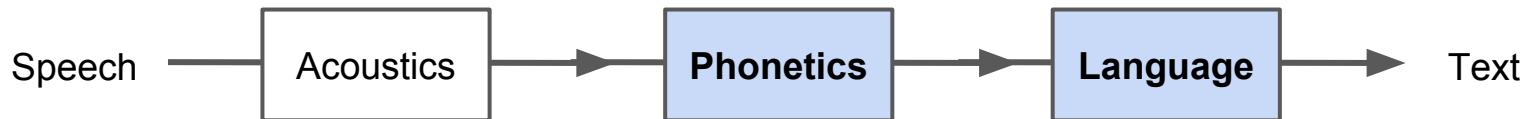
Use fully connected layers on top to add depth

**Convolutional, Long Short-Term Memory,
Fully Connected Deep Neural Networks**

Sainath *et al.* ICASSP'15



Trend: LSTMs end-to-end!



Train recurrent models that also incorporate **Lexical** and **Language Modeling**:

Fast and Accurate Recurrent Neural Network

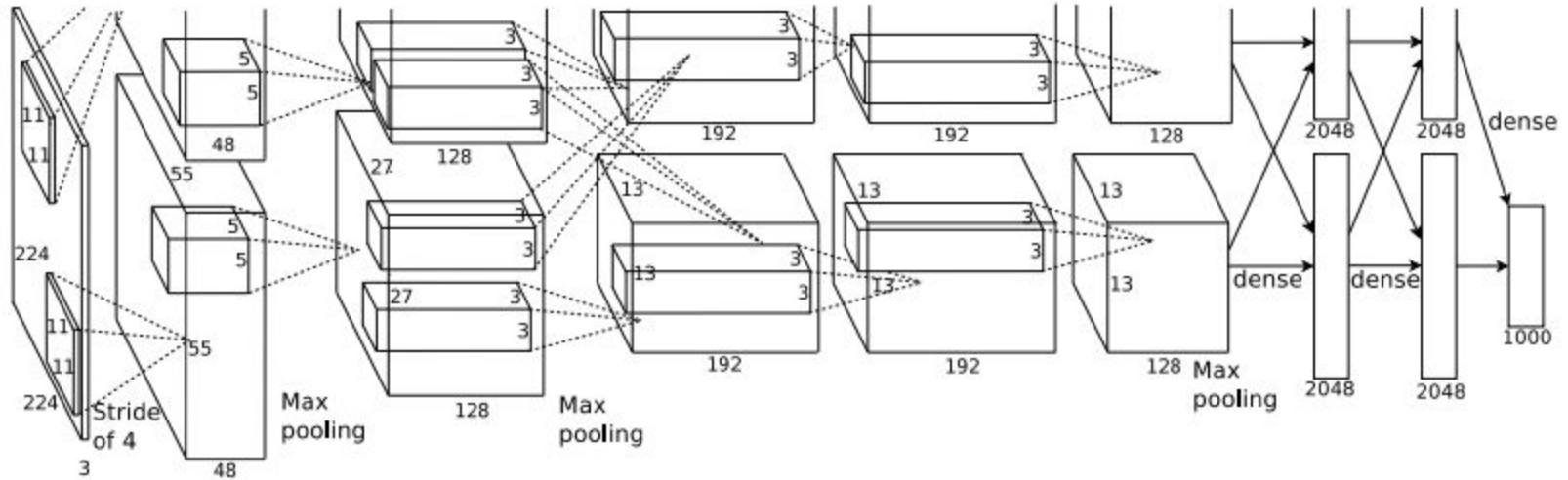
Acoustic Models for Speech Recognition, H. Sak *et al.* 2015

Deep Speech: Scaling up end-to-end speech recognition, A. Hannun *et al.* 2014

Listen, Attend and Spell, W. Chan *et al.* 2015



CNNs for Vision: AlexNet



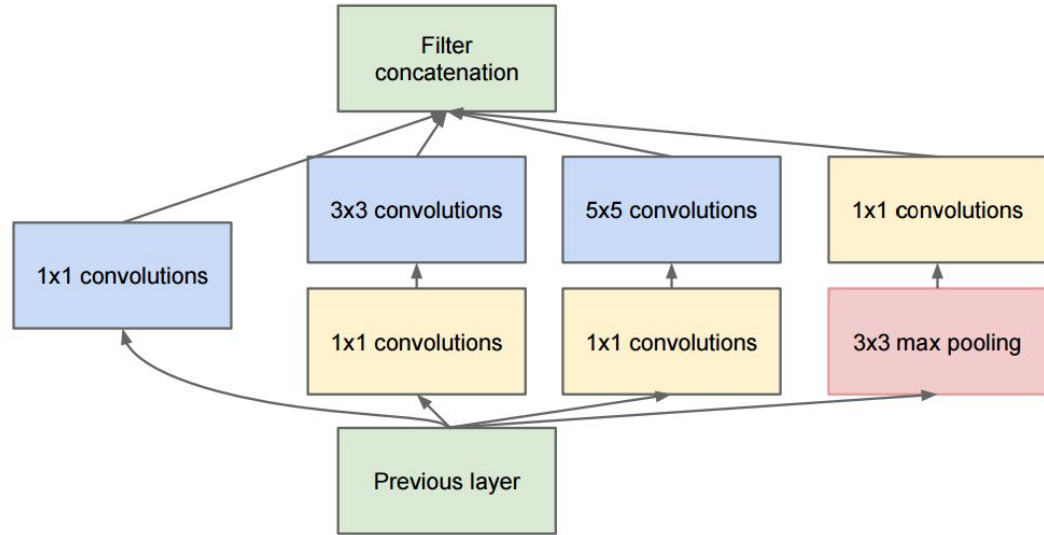
ImageNet Classification with Deep Convolutional Neural Networks

Krizhevsky, Sutskever and Hinton, NIPS 2012

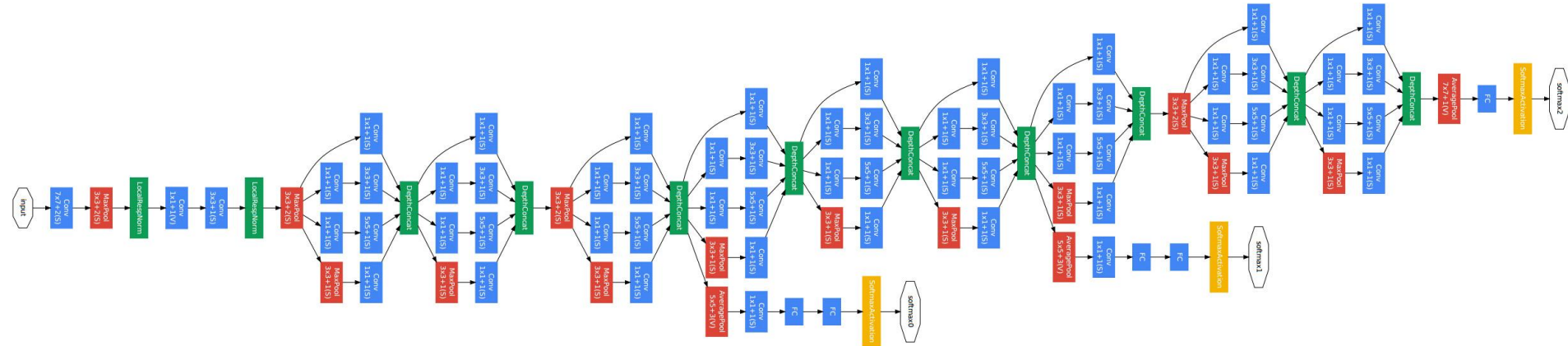


The Inception Architecture (GoogLeNet, 2015)

Basic module, which is then replicated many times



The Inception Architecture (GoogLeNet, 2015)



Going Deeper with Convolutions

Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich

ArXiv 2014, CVPR 2015



Inception-v3 (December 2015)

Rethinking the Inception Architecture for Computer Vision

Christian Szegedy

Google Inc.

szegedy@google.com

Vincent Vanhoucke

vanhoucke@google.com

Sergey Ioffe

sioffe@google.com

Jonathon Shlens

shlens@google.com

Zbigniew Wojna

University College London

zbigniewwojna@gmail.com

<http://arxiv.org/abs/1512.00567>



Rapid Progress in Image Recognition

ImageNet
challenge
classification
task

Team	Year	Place	Error (top-5)	Params
XRCE (pre-neural-net explosion)	2011	1st	25.8%	
Supervision (AlexNet)	2012	1st	16.4%	60M
Clarifai	2013	1st	11.7%	65M
MSRA	2014	3rd	7.35%	
VGG	2014	2nd	7.32%	180M
GoogLeNet (Inception)	2014	1st	6.66%	5M
Andrej Karpathy (human)	2014	N/A	5.1%	100 trillion?
BN-Inception (Arxiv)	2015	N/A	4.9%	13M
Inception-v3 (Arxiv)	2015	N/A	3.46%	25M

Models with small number of parameters fit easily in a mobile app (8-bit fixed point)



What do you want in a machine learning system?

- **Ease of expression:** for lots of crazy ML ideas/algorithms
- **Scalability:** can run experiments quickly
- **Portability:** can run on wide variety of platforms
- **Reproducibility:** easy to share and reproduce research
- **Production readiness:** go from research to real products





<http://tensorflow.org/>

and

<https://github.com/tensorflow/tensorflow>

Open, standard software for
general machine learning

Great for Deep Learning in
particular

First released Nov 2015

Apache 2.0 license

TensorFlow:

Large-Scale Machine Learning on Heterogeneous Distributed Systems

(Preliminary White Paper, November 9, 2015)

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng
Google Research*

Abstract

TensorFlow [1] is an interface for expressing machine learning algorithms, and an implementation for executing such algorithms. A computation expressed using TensorFlow can be executed with little or no change on a wide variety of heterogeneous systems, ranging from mobile devices such as phones

sequence prediction [47], move selection for Go [34], pedestrian detection [2], reinforcement learning [38], and other areas [17, 5]. In addition, often in close collaboration with the Google Brain team, more than 50 teams at Google and other Alphabet companies have deployed deep neural networks using DistBelief in a wide variety

<http://tensorflow.org/whitepaper2015.pdf>

TensorFlow: A system for large-scale machine learning

Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng

Google Brain

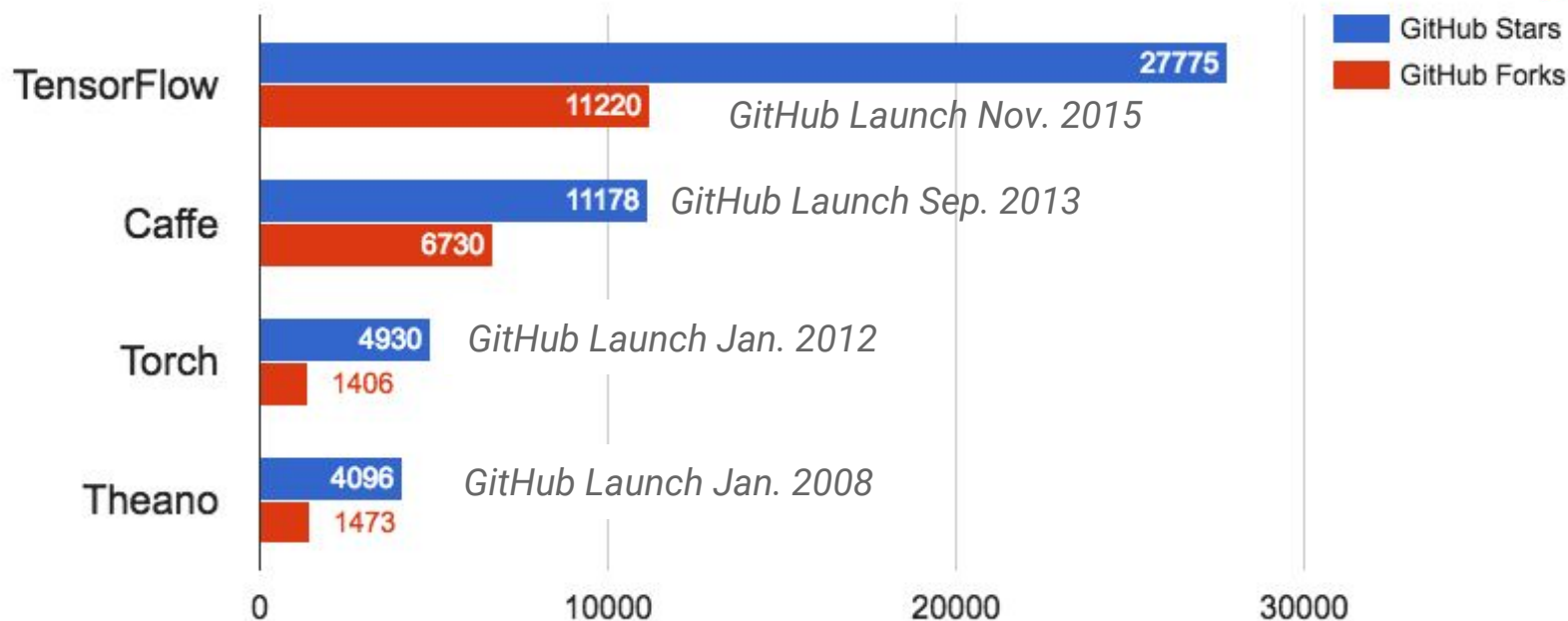
Preprint: arxiv.org/abs/1605.08695

Updated version to appear in OSDI 2016

Strong External Adoption



Adoption of Deep Learning Tools on GitHub

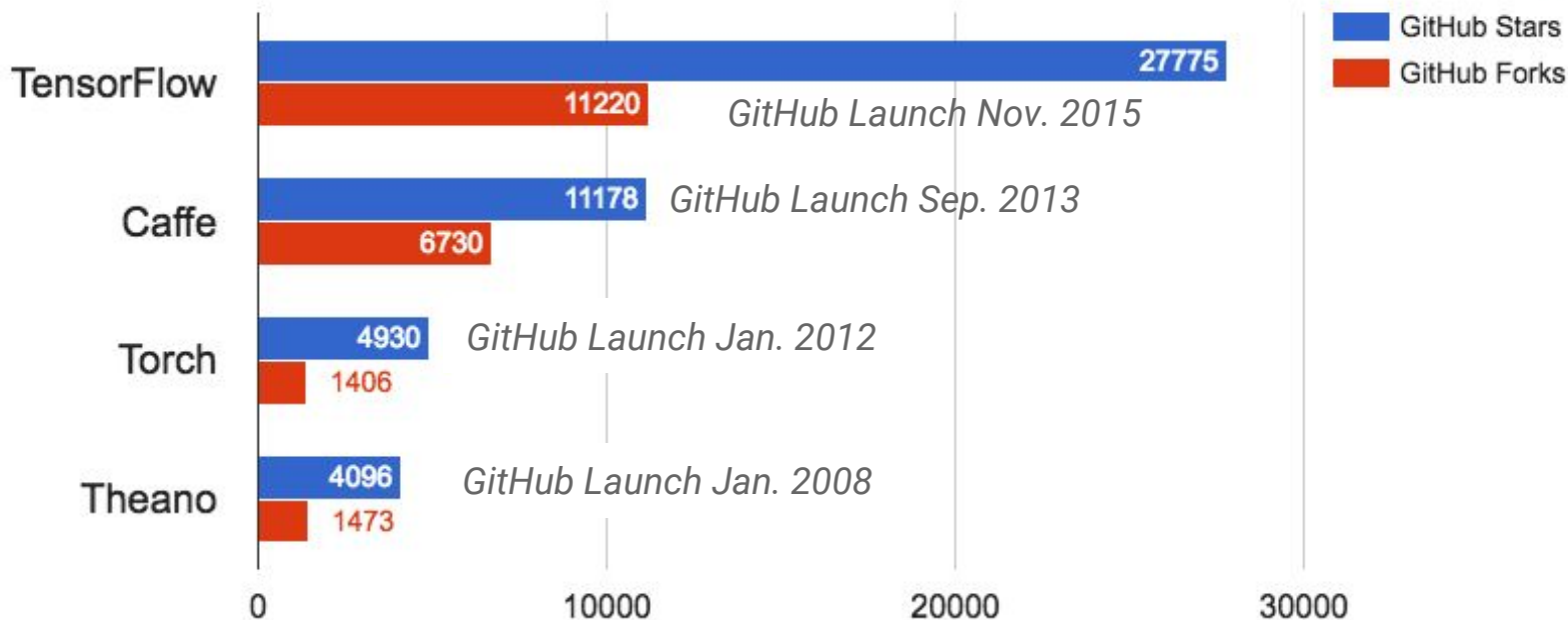


50,000+ binary installs in 72 hours, 500,000+ since November, 2015

Strong External Adoption



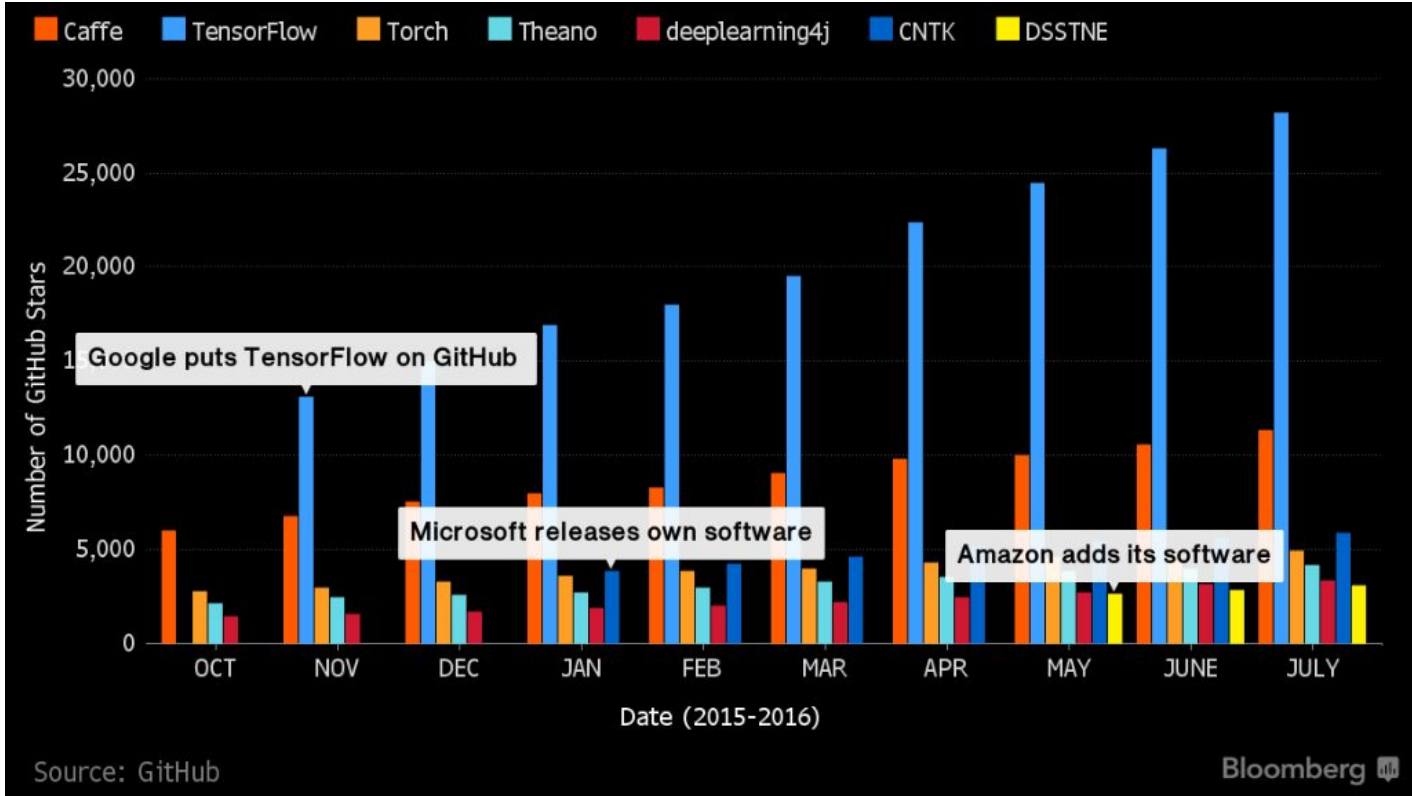
Adoption of Deep Learning Tools on GitHub



50,000+ binary installs in 72 hours, 500,000+ since November, 2015

Most forked new repo on GitHub in 2015 (despite only being available in Nov, '15)

Bloomberg Writes About Open Source Deep Learning Packages?



Source: GitHub

Bloomberg

Source: Bloomberg. www.bloomberg.com/news/articles/2016-07-21/google-sprints-ahead-in-ai-building-blocks-leaving-rivals-wary

Version:

MNIST For ML Beginners

- The MNIST Data
- Softmax Regressions
- Implementing the Regression
- Training
- Evaluating Our Model

Deep MNIST for Experts

- Setup
 - Load MNIST Data
 - Start TensorFlow InteractiveSession
- Build a Softmax Regression Model
 - Placeholders
 - Variables
 - Predicted Class and Cost Function
- Train the Model
 - Evaluate the Model
- Build a Multilayer Convolutional Network
 - Weight Initialization
 - Convolution and Pooling
 - First Convolutional Layer
 - Second Convolutional Layer
 - Densely Connected Layer
 - Readout Layer
 - Train and Evaluate the Model

TensorFlow Mechanics 101

- Tutorial Files
- Prepare the Data

TensorFlow Mechanics 101

This is a technical tutorial, where we walk you through the details of using TensorFlow infrastructure to train models at scale. We use again MNIST as the example.

[View Tutorial](#)

Convolutional Neural Networks

An introduction to convolutional neural networks using the CIFAR-10 data set. Convolutional neural nets are particularly tailored to images, since they exploit translation invariance to yield more compact and effective representations of visual content.

[View Tutorial](#)

Vector Representations of Words

This tutorial motivates why it is useful to learn to represent words as vectors (called word embeddings). It introduces the word2vec model as an efficient method for learning embeddings. It also covers the high-level details behind noise-contrastive training methods (the biggest recent advance in training embeddings).

[View Tutorial](#)

Recurrent Neural Networks

An introduction to RNNs, wherein we train an LSTM network to predict the next word in an English sentence. (A task sometimes called language modeling.)

[View Tutorial](#)

Sequence-to-Sequence Models

A follow on to the RNN tutorial, where we assemble a sequence-to-sequence model for machine translation. You will learn to build your own English-to-French translator, entirely machine learned, end-to-end.

[View Tutorial](#)

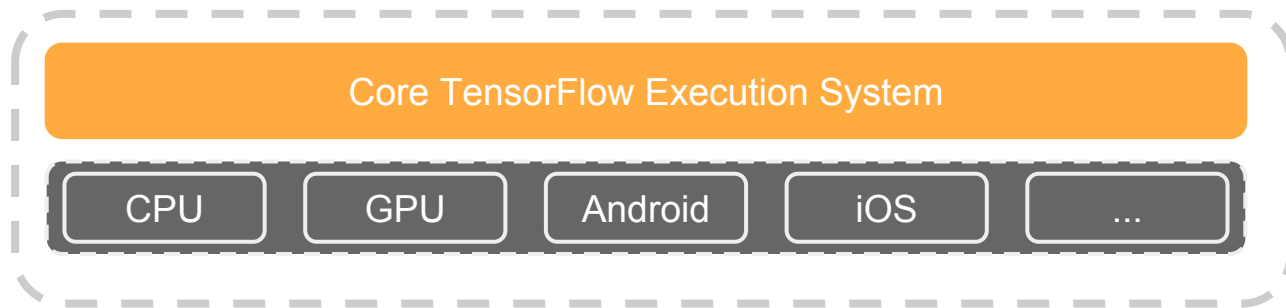


Motivations

- DistBelief (our 1st system) was the first scalable deep learning system, but not as flexible as we wanted for research purposes
- Better understanding of problem space allowed us to make some dramatic simplifications

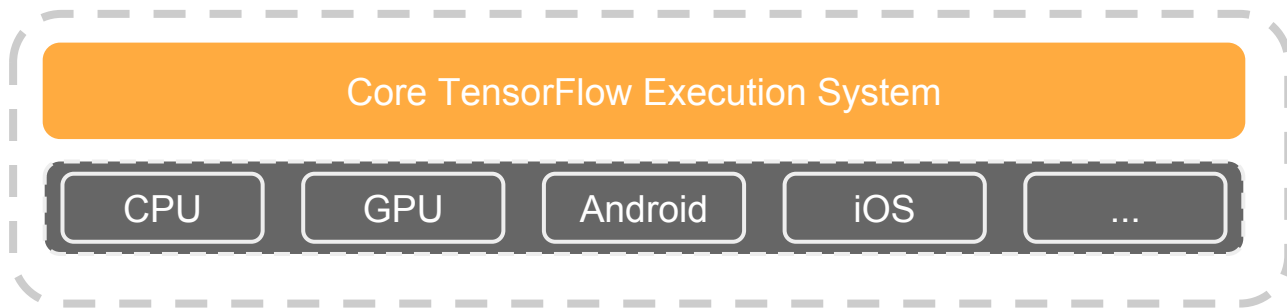
TensorFlow: Expressing High-Level ML Computations

- Core in C++
 - Very low overhead



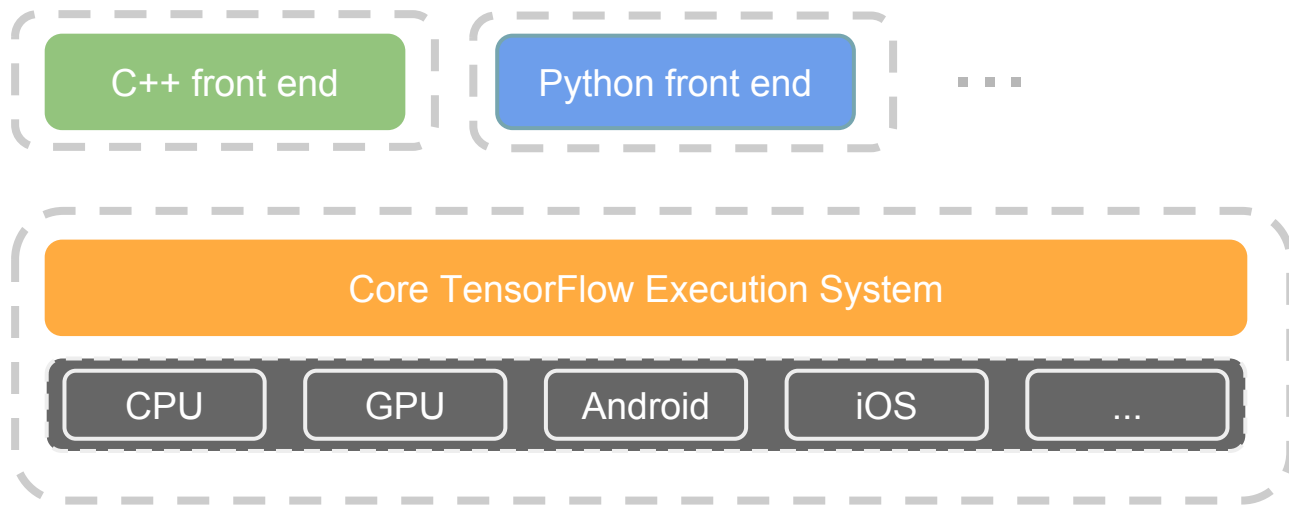
TensorFlow: Expressing High-Level ML Computations

- Core in C++
 - Very low overhead
- Different front ends for specifying/driving the computation
 - Python and C++ today, easy to add more

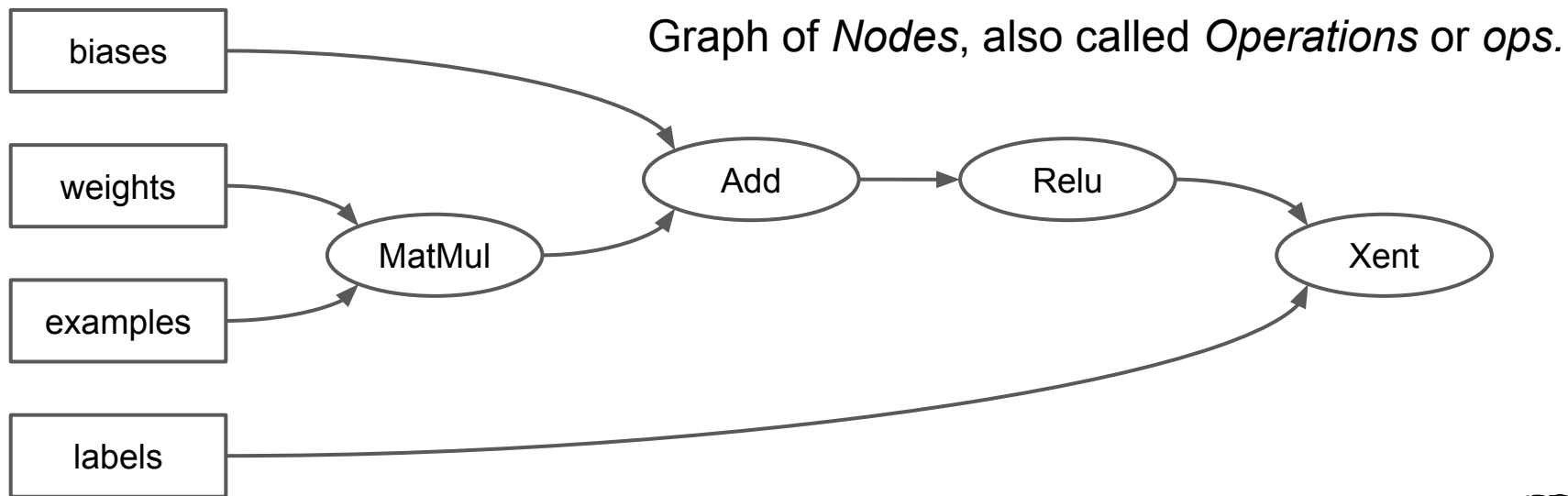


TensorFlow: Expressing High-Level ML Computations

- Core in C++
 - Very low overhead
- Different front ends for specifying/driving the computation
 - Python and C++ today, easy to add more

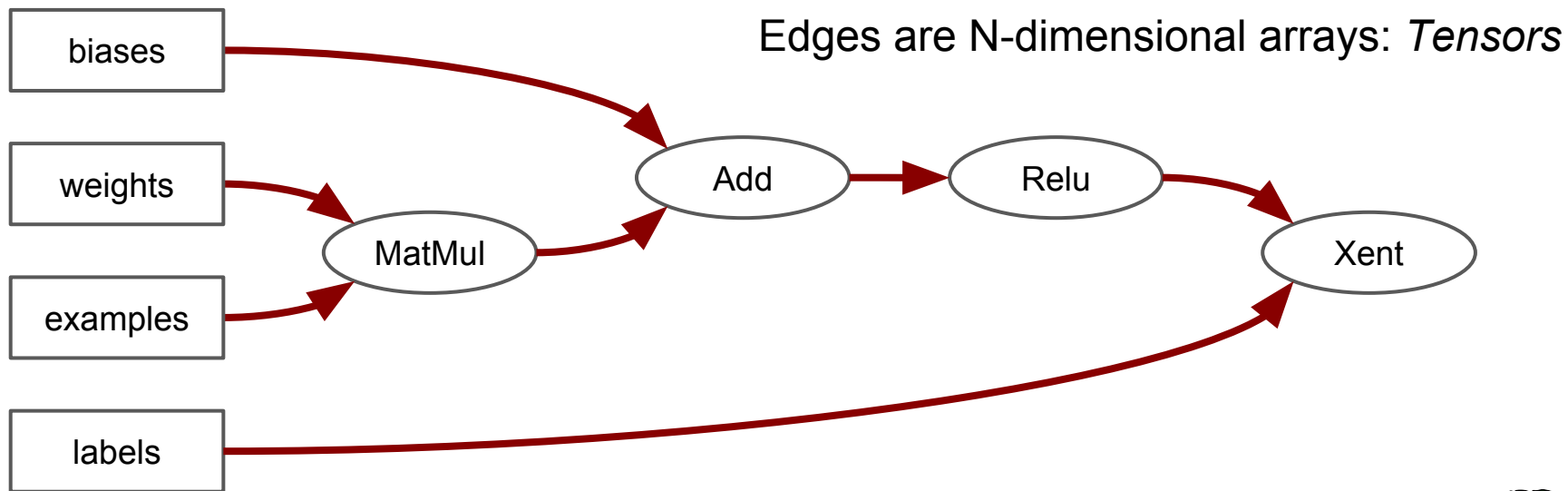


Computation is a dataflow graph



Computation is a dataflow graph

with tensors



Example TensorFlow fragment

- Build a graph computing a neural net inference.

```
import tensorflow as tf
from tensorflow.examples.tutorials.mnist import input_data

mnist = input_data.read_data_sets('MNIST_data', one_hot=True)
x = tf.placeholder("float", shape=[None, 784])
W = tf.Variable(tf.zeros([784,10]))
b = tf.Variable(tf.zeros([10]))
y = tf.nn.softmax(tf.matmul(x, W) + b)
```

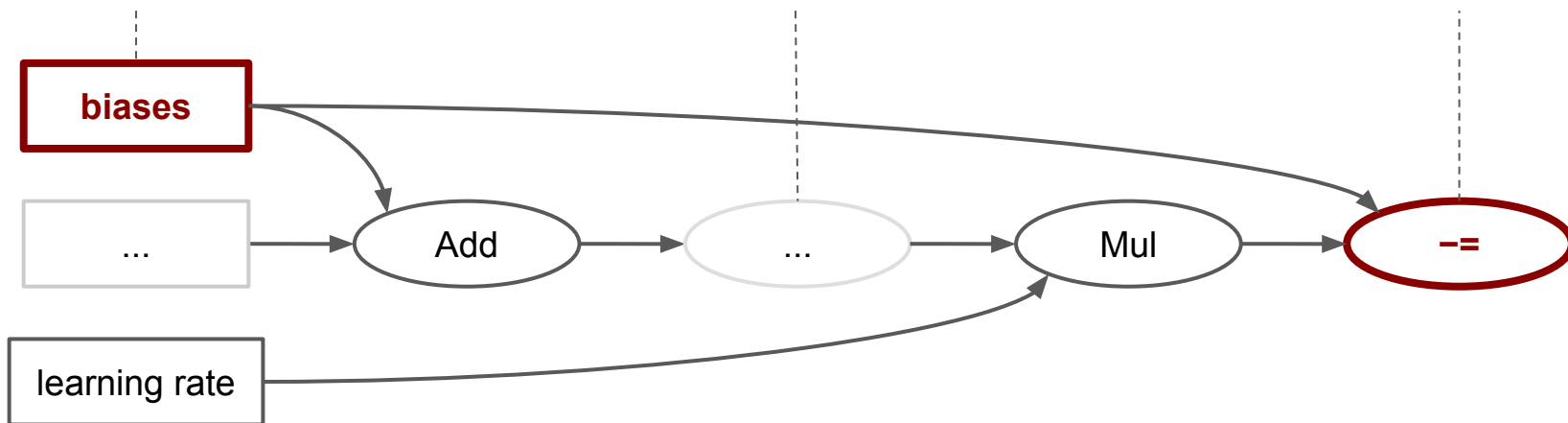
Computation is a dataflow graph

with state

'Biases' is a variable

Some ops compute gradients

--= updates biases



Symbolic Differentiation

- Automatically add ops to calculate symbolic gradients of variables w.r.t. loss function.
- Apply these gradients with an optimization algorithm

```
y_ = tf.placeholder(tf.float32, [None, 10])  
cross_entropy = -tf.reduce_sum(y_ * tf.log(y))  
opt = tf.train.GradientDescentOptimizer(0.01)  
train_op = opt.minimize(cross_entropy)
```

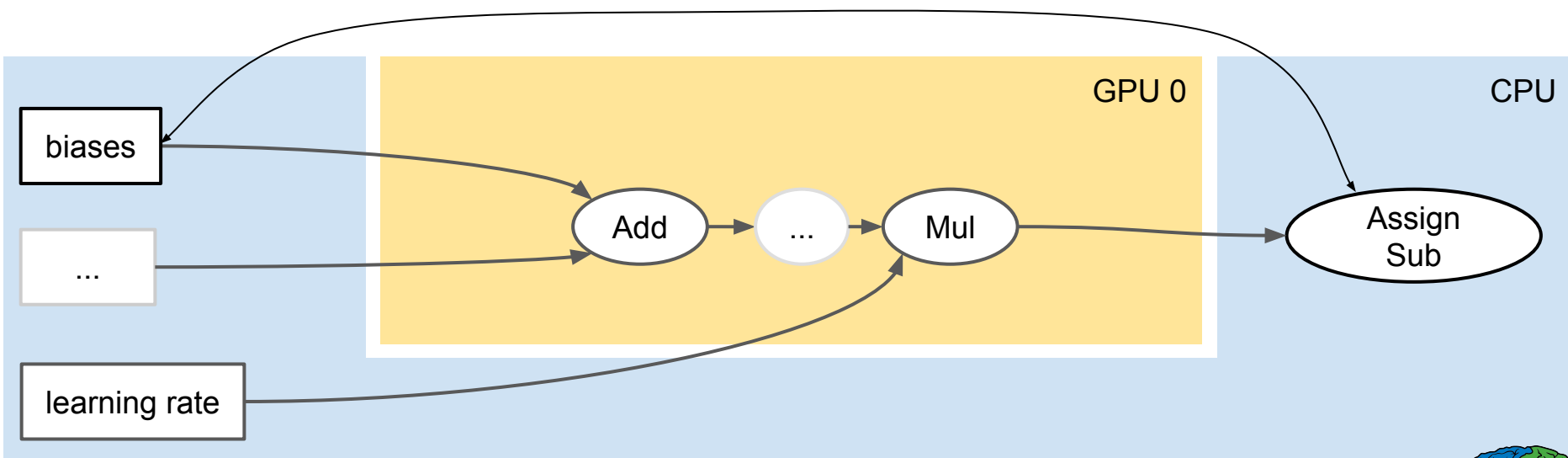

Define graph and then execute it repeatedly

- Launch the graph and run the training ops in a loop

```
init = tf.initialize_all_variables()
sess = tf.Session()
sess.run(init)
for i in range(1000):
    batch_xs, batch_ys = mnist.train.next_batch(100)
    sess.run(train_step, feed_dict={x: batch_xs, y_: batch_ys})
```

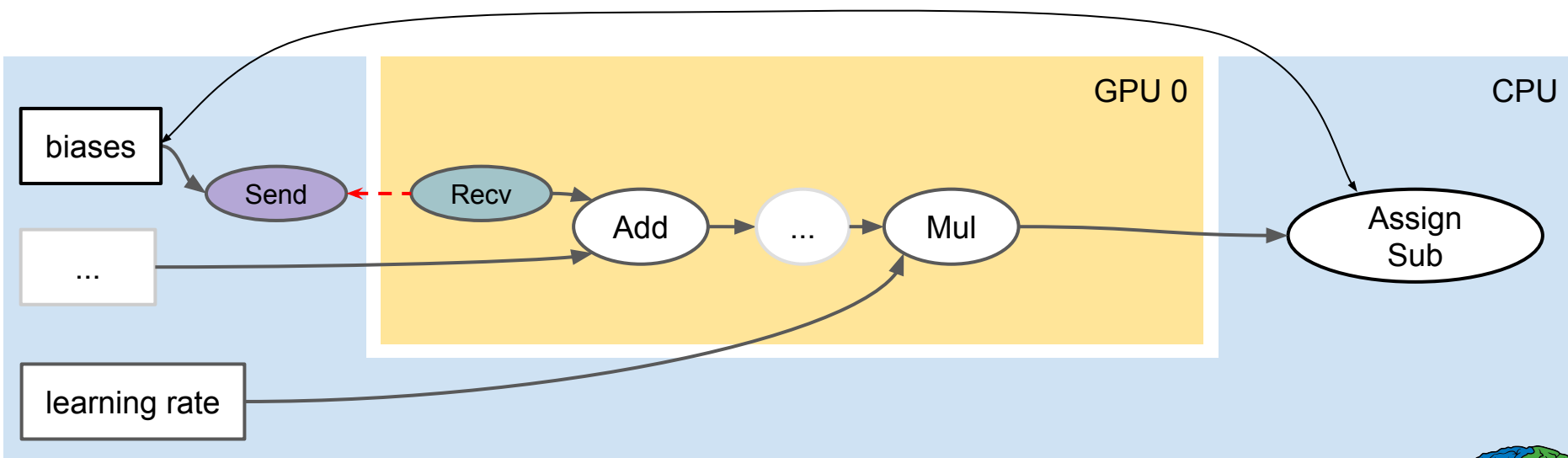
Computation is a dataflow graph

distributed



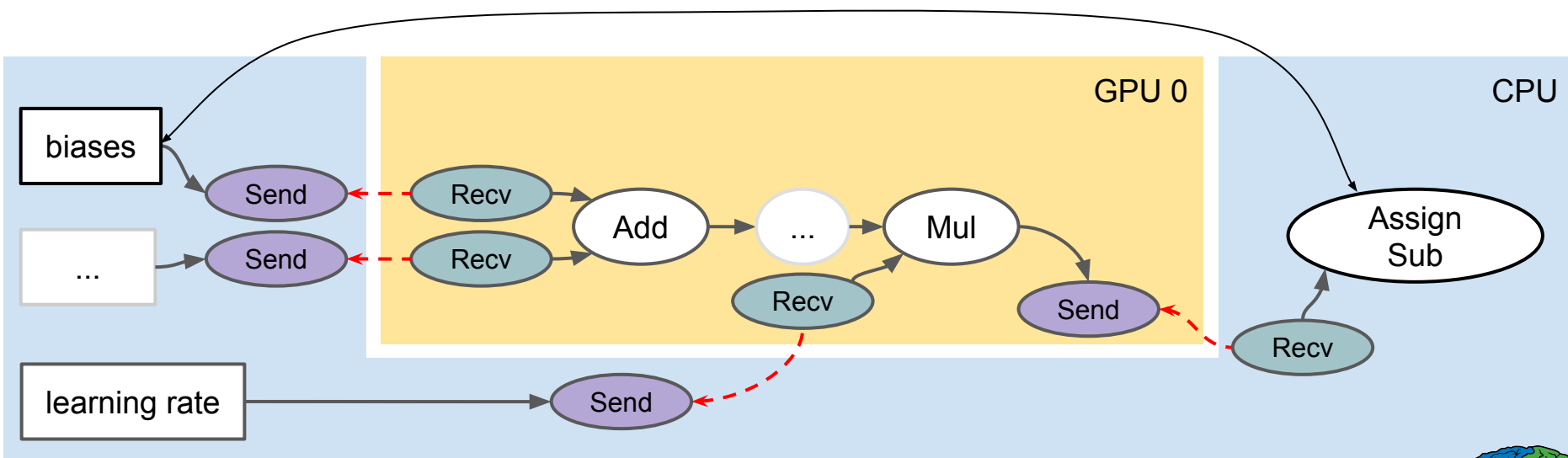
Assign *Devices* to Ops

- TensorFlow inserts *Send/Recv* Ops to transport tensors across devices
- *Recv* ops pull data from *Send* ops



Assign *Devices* to Ops

- TensorFlow inserts *Send/Recv* Ops to transport tensors across devices
- *Recv* ops pull data from *Send* ops



Send and Receive Implementations

- Different implementations depending on source/dest devices
- e.g. GPUs on same machine: **local GPU** → **GPU copy**
- e.g. CPUs on different machines: **cross-machine RPC**
- e.g. GPUs on different machines: **RDMA**



November 2015

🔗 **Release 0.5.0**

Initial release of TensorFlow.



December 2015

Release 0.6.0

Major Features and Improvements

- Python 3.3+ support via changes to python codebase and ability to specify python version via ./configure.
- Some improvements to GPU performance and memory usage: [convnet benchmarks](#) roughly equivalent with native cudnn v2 performance. Improvements mostly due to moving to 32-bit indices, faster shuffling kernels. More improvements to come in later releases.

Bug Fixes

- Lots of fixes to documentation and tutorials, many contributed by the public.
- 271 closed issues on github issues.



February 2016

Release 0.7.0

Major Features and Improvements

- Allow using any installed Cuda ≥ 7.0 and cuDNN $\geq R2$, and add support for cuDNN R4
- Added a `contrib/` directory for unsupported or experimental features, including higher level `layers` module
- Added an easy way to add and dynamically load user-defined ops
- Built out a good suite of tests, things should break less!
- Added `MetaGraphDef` which makes it easier to save graphs with metadata
- Added assignments for "Deep Learning with TensorFlow" udacity course

Bug Fixes and Other Changes

- Added a versioning framework for `GraphDef`s to ensure compatibility
- Enforced Python 3 compatibility
- Internal changes now show up as sensibly separated commits
- Open-sourced the doc generator



April 2016

Release 0.8.0

Major Features and Improvements

- Added a distributed runtime using GRPC
- Move skflow to `contrib/learn`
- Better linear optimizer in `contrib/linear_optimizer`
- Random forest implementation in `contrib/tensor_forest`
- CTC loss and decoders in `contrib/ctc`
- Basic support for `half` data type
- Better support for loading user ops (see examples in `contrib/`)
- Allow use of (non-blocking) Eigen threadpool with `TENSORFLOW_USE_EIGEN_THREADPOOL` define
- Add an extension mechanism for adding network file system support
- TensorBoard displays metadata stats (running time, memory usage and device used) and tensor shapes

Big Fixes and Other Changes

- Utility for inspecting checkpoints
- Basic tracing and timeline support
- Allow building against cuDNN 5 (not incl. RNN/LSTM support)



June 2016

Release 0.9.0

Major Features and Improvements

- Python 3.5 support and binaries
- Added iOS support
- Added support for processing on GPUs on MacOS
- Added makefile for better cross-platform build support (C API only)
- fp16 support and improved complex128 support for many ops
- Higher level functionality in contrib.{layers,losses,metrics,learn}
- More features to Tensorboard
- Improved support for string embedding and sparse features
- The RNN api is finally "official" (see, e.g., `tf.nn.dynamic_rnn`, `tf.nn.rnn`, and the classes in `tf.nn.rnn_cell`).
- TensorBoard now has an Audio Dashboard, with associated audio summaries.

Big Fixes and Other Changes

- Turned on CuDNN Autotune.
- Added support for using third-party Python optimization algorithms (contrib.opt).
- Google Cloud Storage filesystem support.
- HDF5 support



Activity

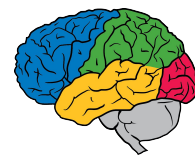
6,354 commits

12 branches

10 releases

335 contributors

Contributions to master, excluding merge commits



DeepMind moves to TensorFlow

Friday, April 29, 2016

Posted by Koray Kavukcuoglu, Research Scientist, Google DeepMind

At [DeepMind](#), we conduct state-of-the-art [research](#) on a wide range of algorithms, from deep learning and reinforcement learning to systems neuroscience, towards the goal of building [Artificial General Intelligence](#). A key factor in facilitating rapid progress is the software environment used for research. For nearly four years, the open source [Torch7](#) machine learning library has served as our primary research platform, combining excellent flexibility with very fast runtime execution, enabling rapid prototyping. Our team has been proud to contribute to the open source project in capacities ranging from occasional bug fixes to being core maintainers of several crucial components.

With Google's recent open source release of [TensorFlow](#), we initiated a project to test its suitability for our research environment. Over the last six months, we have re-implemented more than a dozen different projects in TensorFlow to develop a deeper understanding of its potential use cases and the tradeoffs for research. Today we are excited to announce that DeepMind will start using TensorFlow for all our future research. We believe that TensorFlow will enable us to execute our ambitious research goals at much larger scale and an even faster pace, providing us with a unique opportunity to further accelerate our research programme.

As one of the core contributors of Torch7, I have had the pleasure of working closely with an excellent community of developers and researchers, and it has been amazing to see all the great



Pre-trained Inception-v3 model released

<http://googleresearch.blogspot.com/2015/12/how-to-classify-images-with-tensorflow.html>

Dear TensorFlow community,

Today we are releasing our best image classifier trained on ImageNet data. As described in our recent Arxiv preprint at <http://arxiv.org/abs/1512.00567>, an ensemble of four of these models achieves 3.5% top-5 error on the validation set of the ImageNet whole image ILSVRC2012 classification task (compared with our ensemble from last year that won the 2014 ImageNet classification challenge with a 6.66% top-5 error rate).

In this release, we are supplying code and data files containing the trained model parameters for running the image classifier on:

- Both desktop and mobile environments
- Employing either a C++ or Python API.

In addition, we are providing a tutorial that describes how to use the image recognition system for a variety of use-cases.

http://www.tensorflow.org/tutorials/image_recognition/index.html



MNIST For ML Beginners

- The MNIST Data
- Softmax Regressions
- Implementing the Regression
- Training
- Evaluating Our Model

Deep MNIST for Experts

- Setup
 - Load MNIST Data
 - Start TensorFlow InteractiveSession
- Build a Softmax Regression Model
 - Placeholders
 - Variables
 - Predicted Class and Cost Function
- Train the Model
 - Evaluate the Model
- Build a Multilayer Convolutional Network
 - Weight Initialization
 - Convolution and Pooling
 - First Convolutional Layer
 - Second Convolutional Layer
 - Densely Connected Layer
 - Readout Layer
 - Train and Evaluate the Model

TensorFlow Mechanics 101

```
bazel build tensorflow/examples/label_image/...
```

That should create a binary executable that you can then run like this:

```
bazel-bin/tensorflow/examples/label_image/label_image
```

This uses the default example image that ships with the framework, and should output something similar to this:

```
I tensorflow/examples/label_image/main.cc:200] military uniform (866): 0.647296
I tensorflow/examples/label_image/main.cc:200] suit (794): 0.0477196
I tensorflow/examples/label_image/main.cc:200] academic gown (896): 0.0232411
I tensorflow/examples/label_image/main.cc:200] bow tie (817): 0.0157356
I tensorflow/examples/label_image/main.cc:200] bolo tie (940): 0.0145024
```

In this case, we're using the default image of [Admiral Grace Hopper](#), and you can see the network correctly identifies she's wearing a military uniform, with a high score of 0.6.



Experiment Turnaround Time and Research Productivity

- **Minutes, Hours:**
 - **Interactive research! Instant gratification!**
- **1-4 days**
 - Tolerable
 - Interactivity replaced by running many experiments in parallel
- **1-4 weeks**
 - High value experiments only
 - Progress stalls
- **>1 month**
 - Don't even try





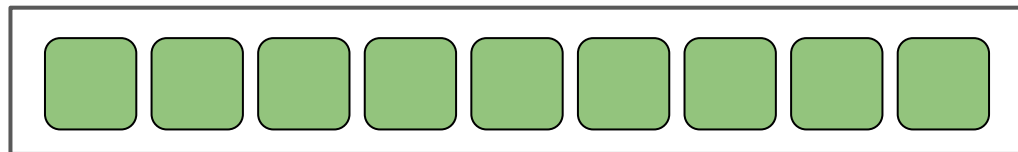
Data Parallelism

- Use multiple model replicas to process different examples at the same time
 - All collaborate to update model state (parameters) in shared parameter server(s)
- Speedups depend highly on kind of model
 - Dense models: 10-40X speedup from 50 replicas
 - Sparse models:
 - support many more replicas
 - often can use as many as 1000 replicas

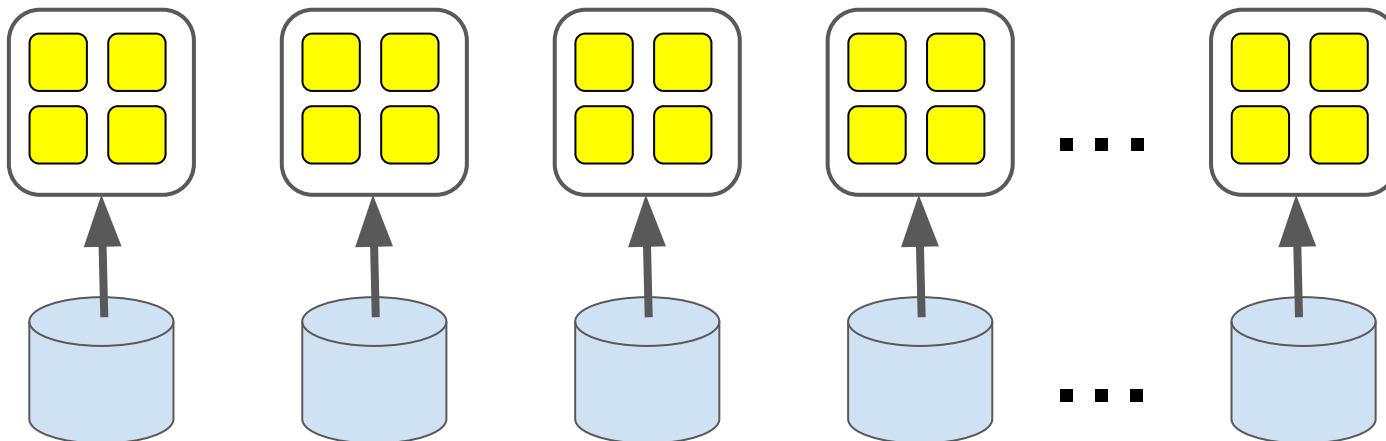


Data Parallelism

Parameter Servers



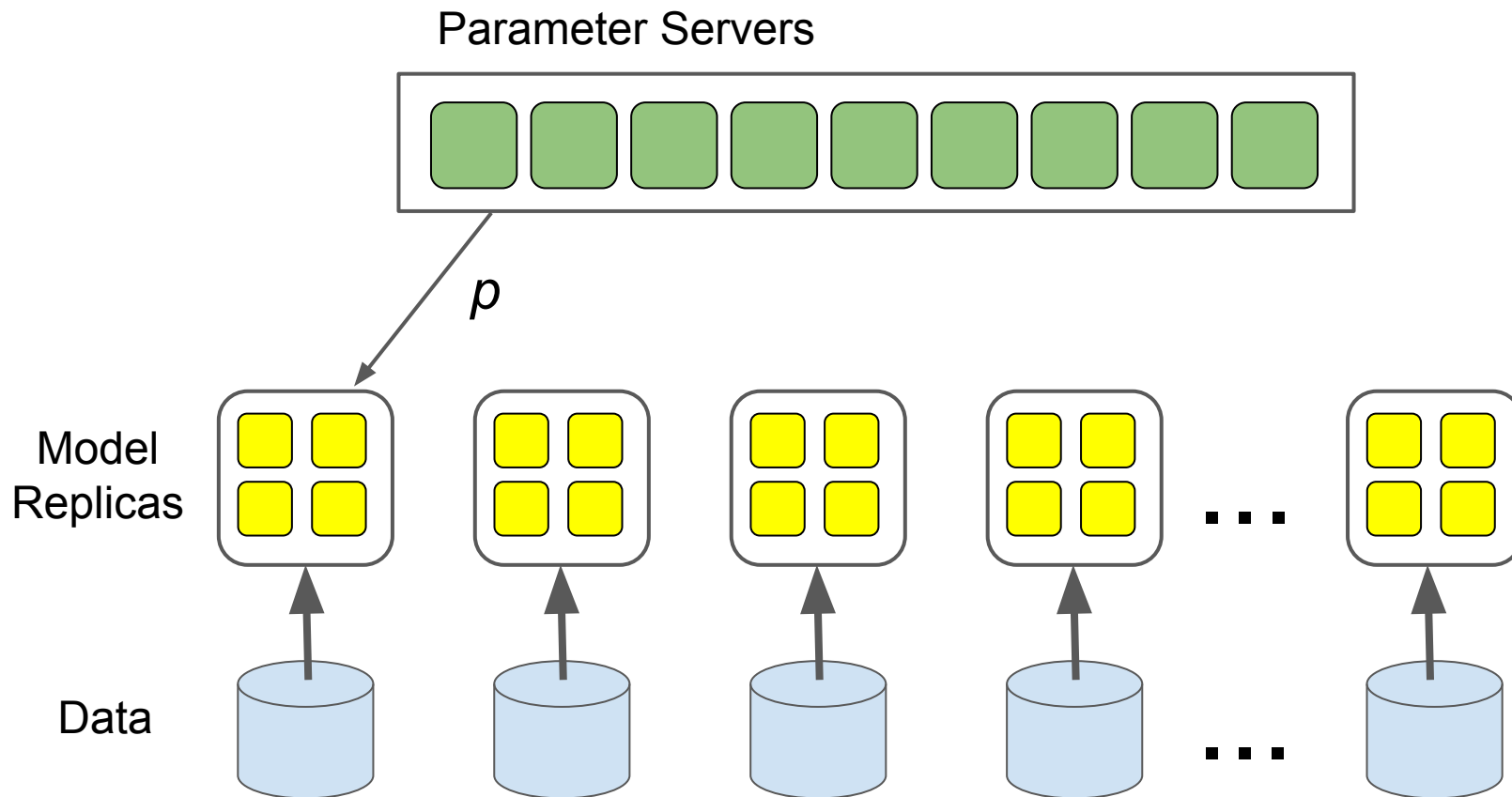
Model
Replicas



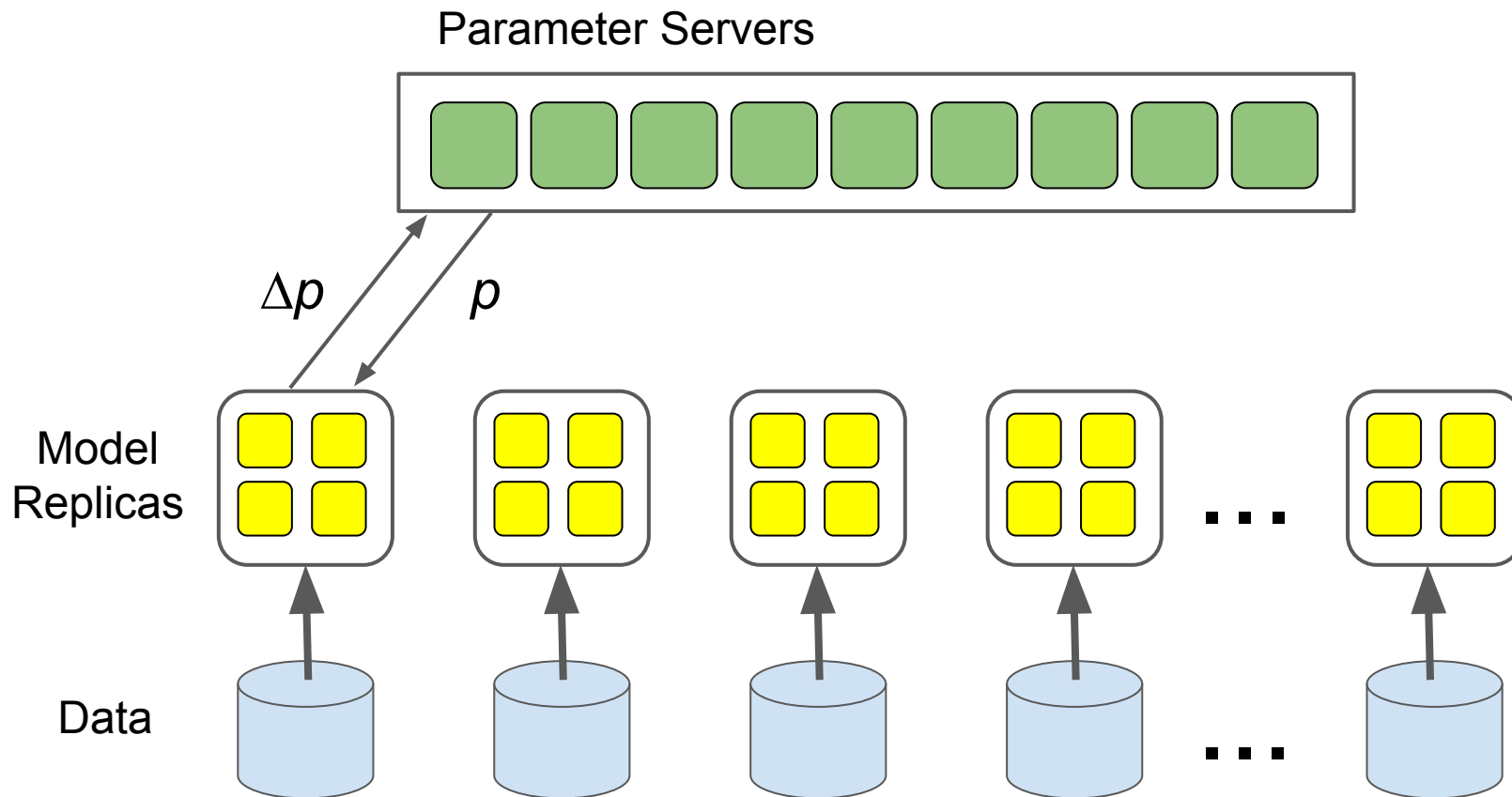
Data



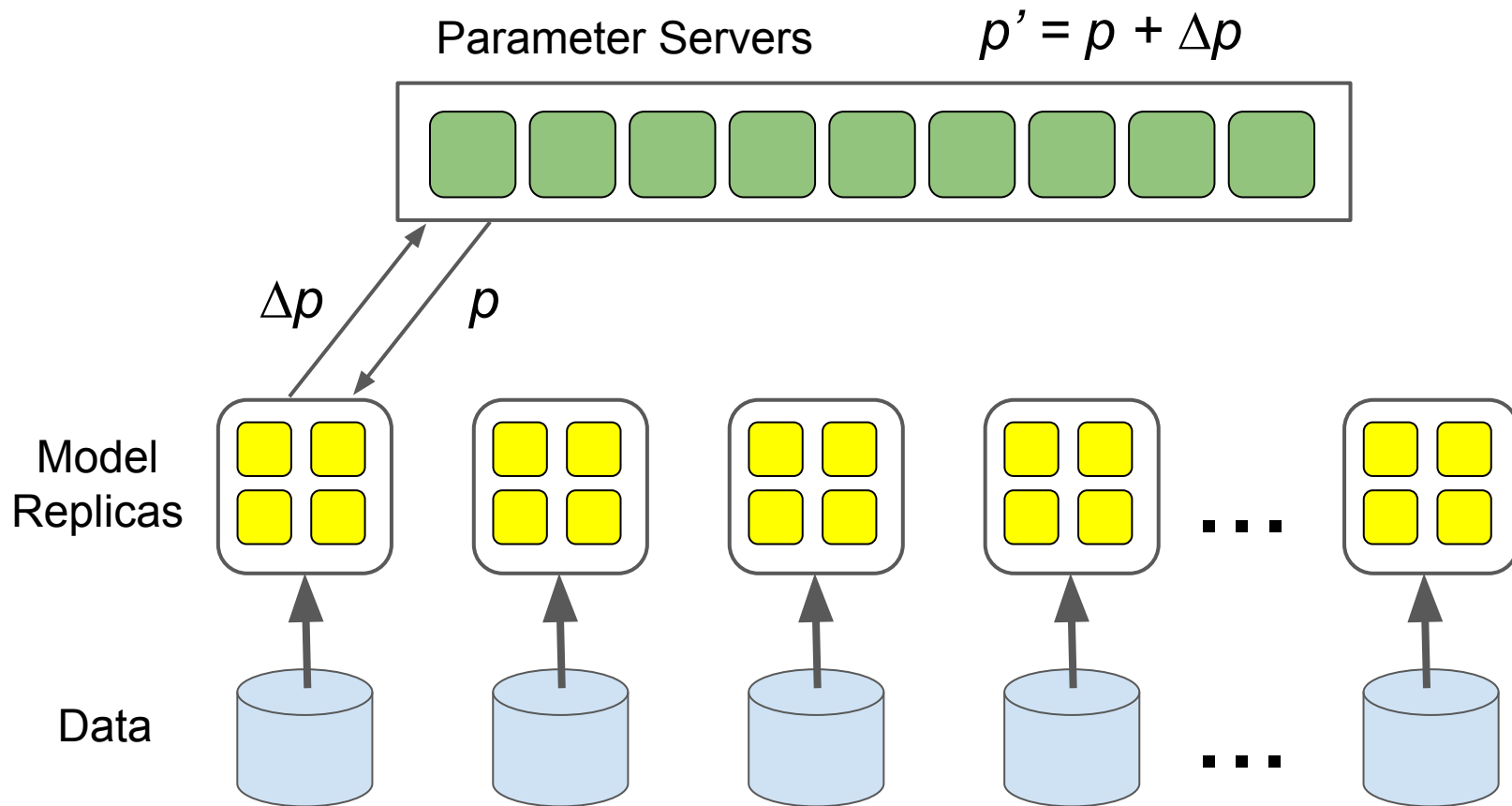
Data Parallelism



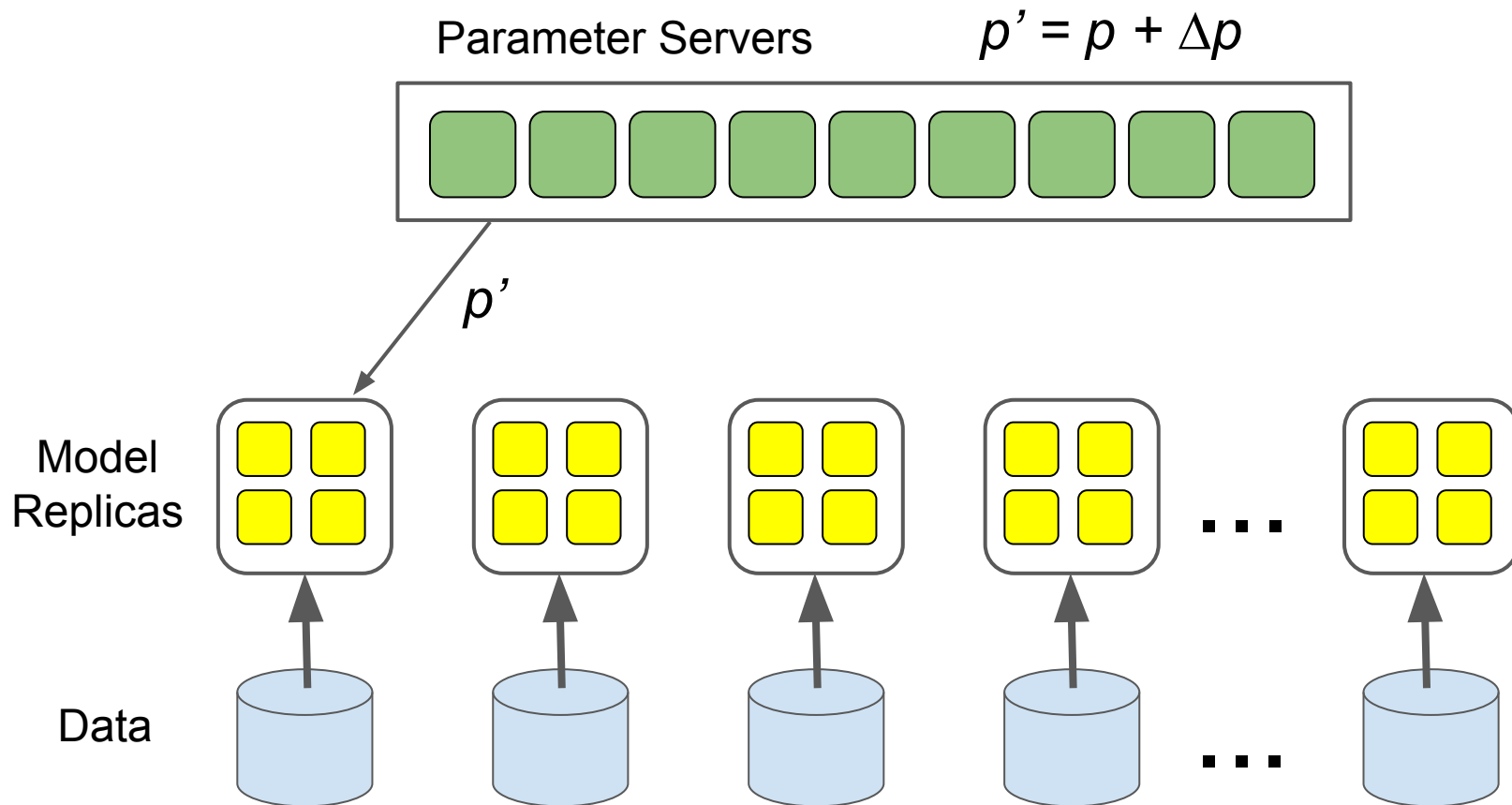
Data Parallelism



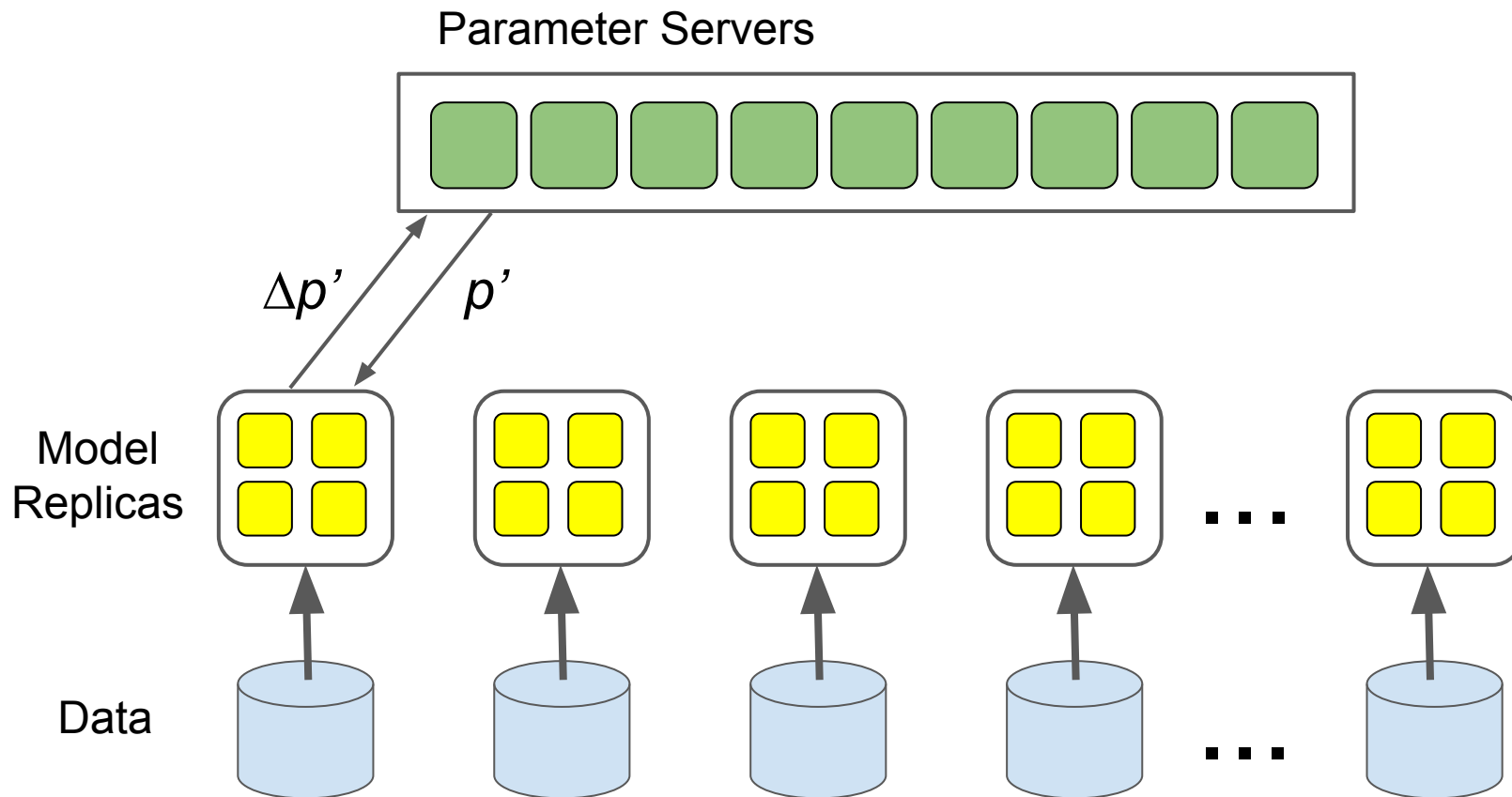
Data Parallelism



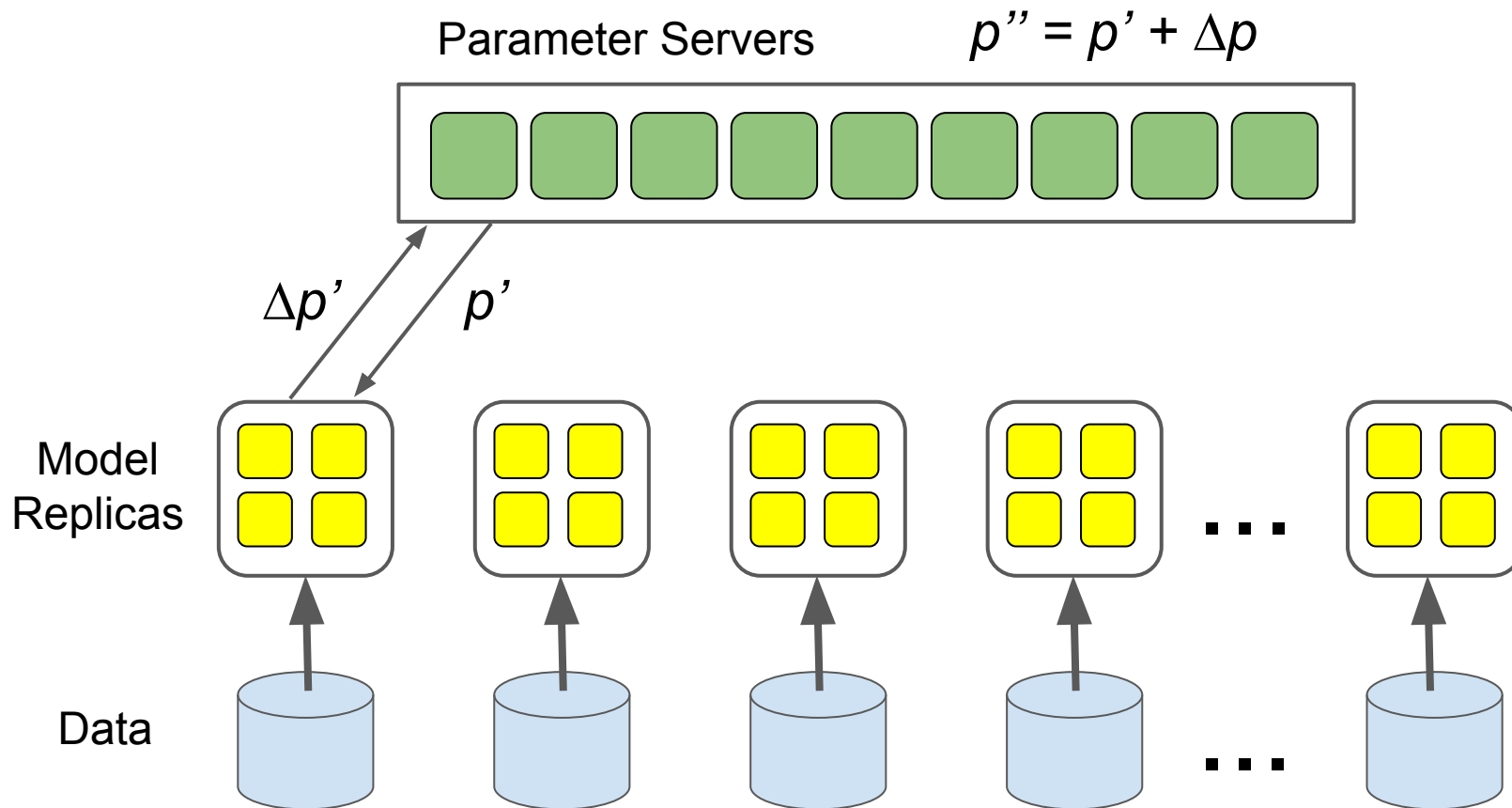
Data Parallelism



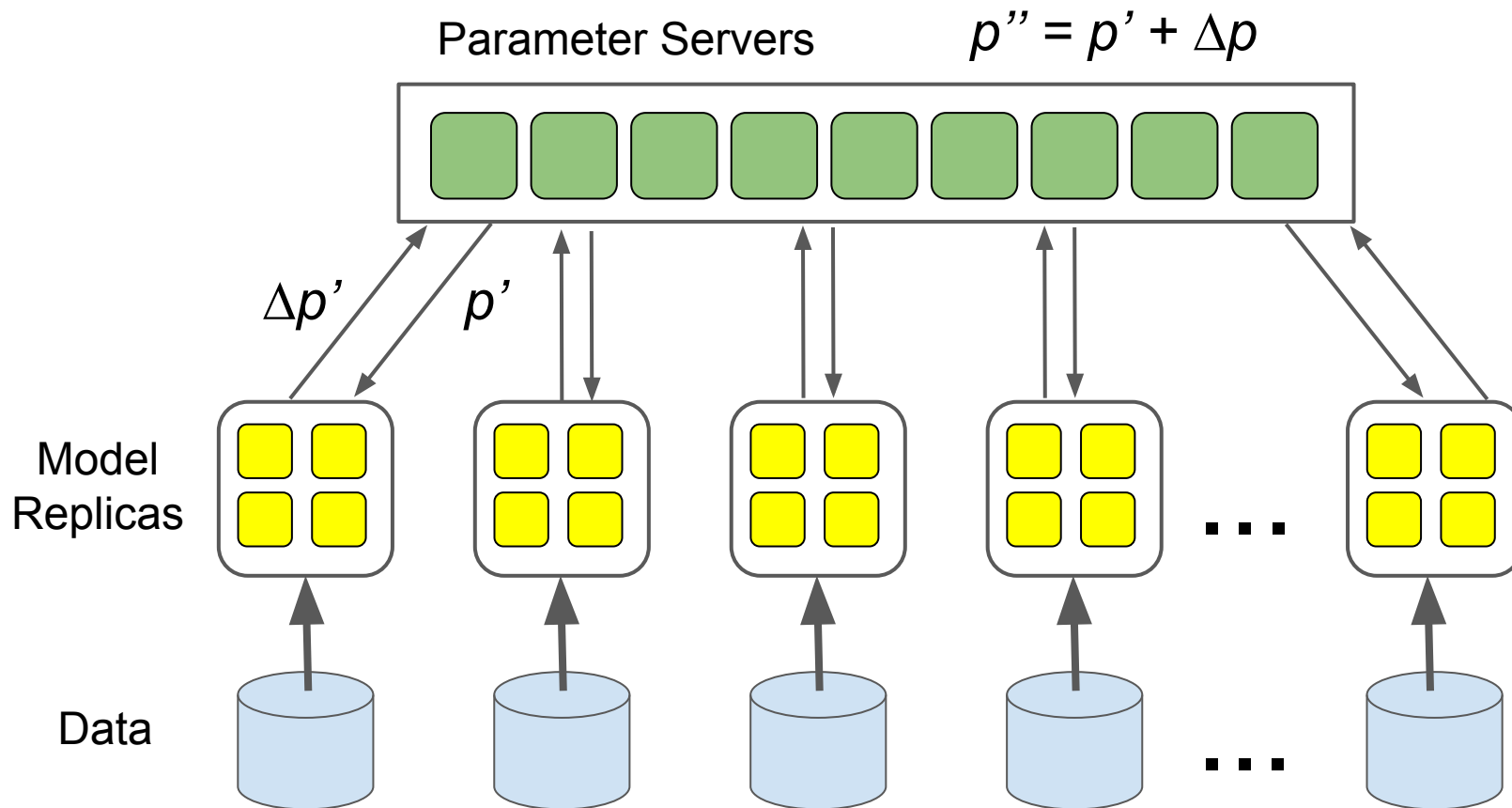
Data Parallelism



Data Parallelism



Data Parallelism



DistBelief: Separate Parameter Servers

Parameter update rules not the same programming model as the rest of the system

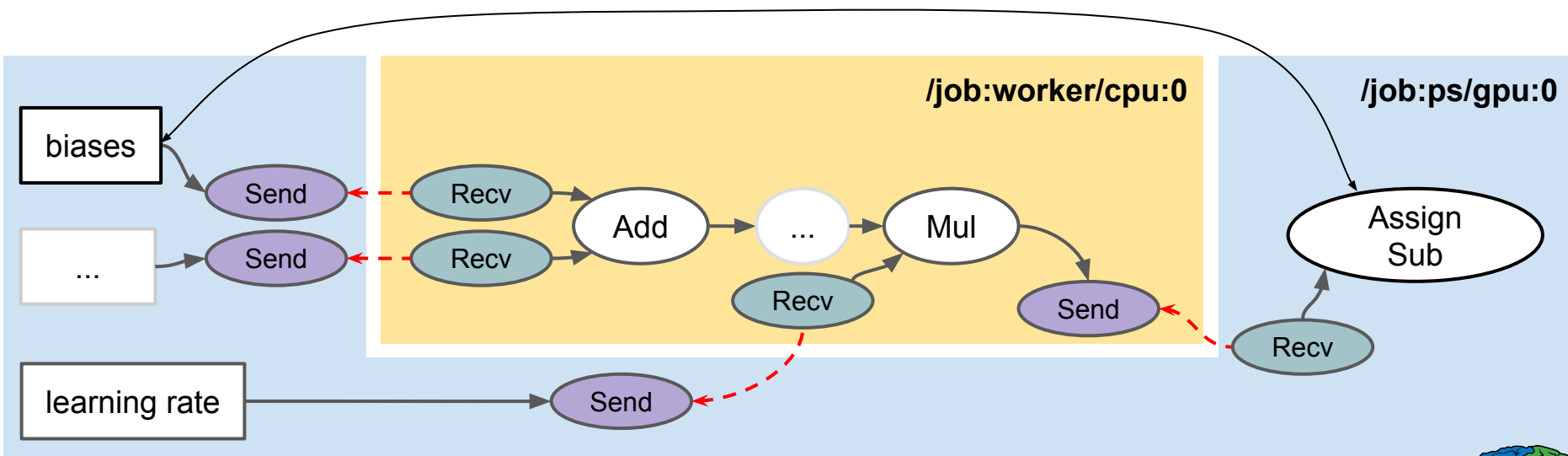
Separate code for parameter servers vs. rest of system

Lacked uniformity & was more complicated



Cross process communication is the same!

- Communication across machines over the network abstracted identically to cross device communication.



No specialized parameter server subsystem!



Data Parallelism Choices

Can do this **synchronously**:

- **N replicas** equivalent to an **N times larger batch size**
- Pro: No gradient staleness
- Con: Less fault tolerant (requires some recovery if any single machine fails)

Can do this **asynchronously**:

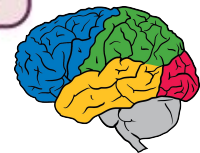
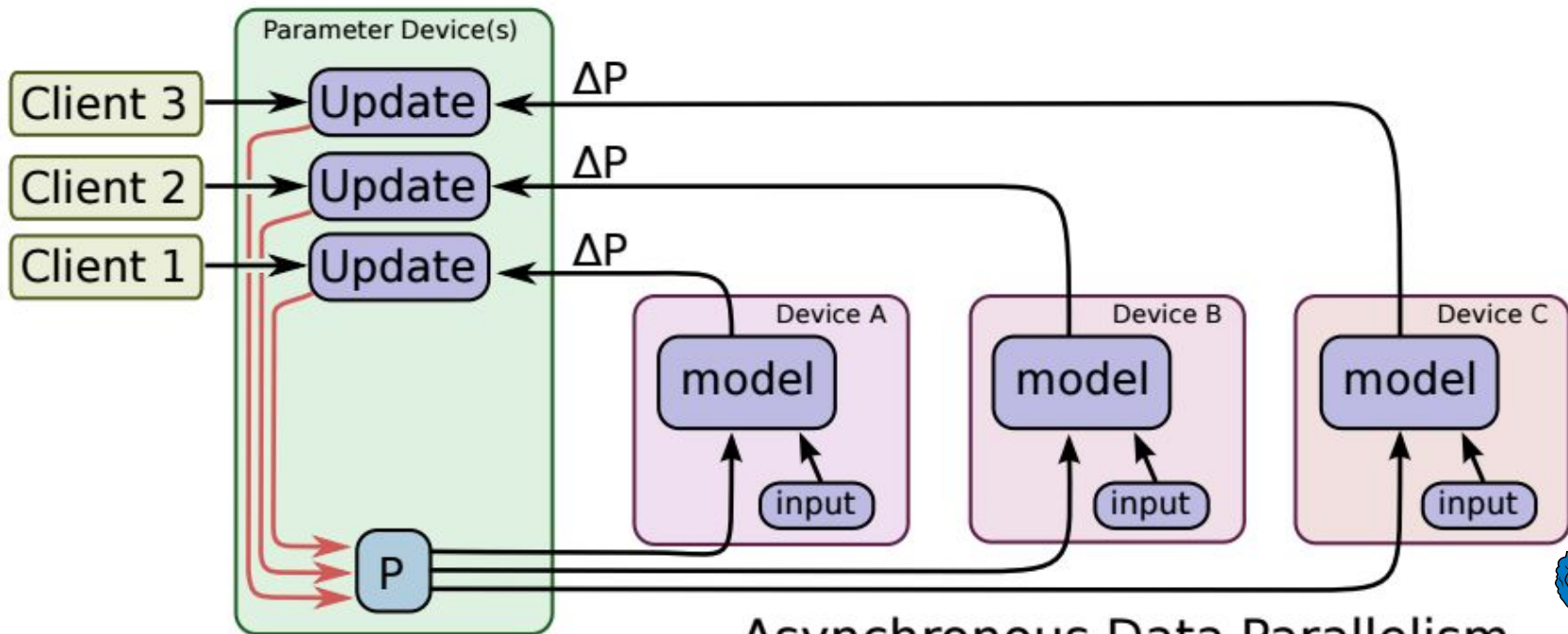
- Pro: Relatively fault tolerant (failure in model replica doesn't block other replicas)
- Con: Gradient staleness means each gradient less effective

(Or **hybrid**: M asynchronous groups of N synchronous replicas)

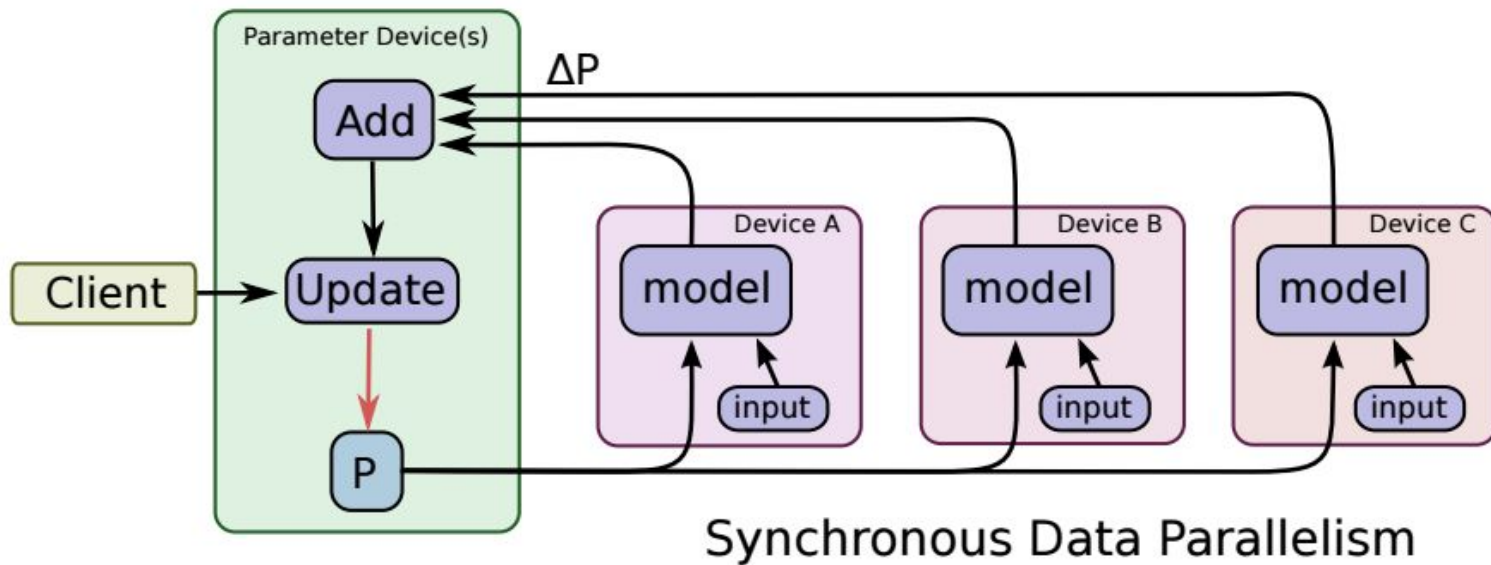


Asynchronous Training

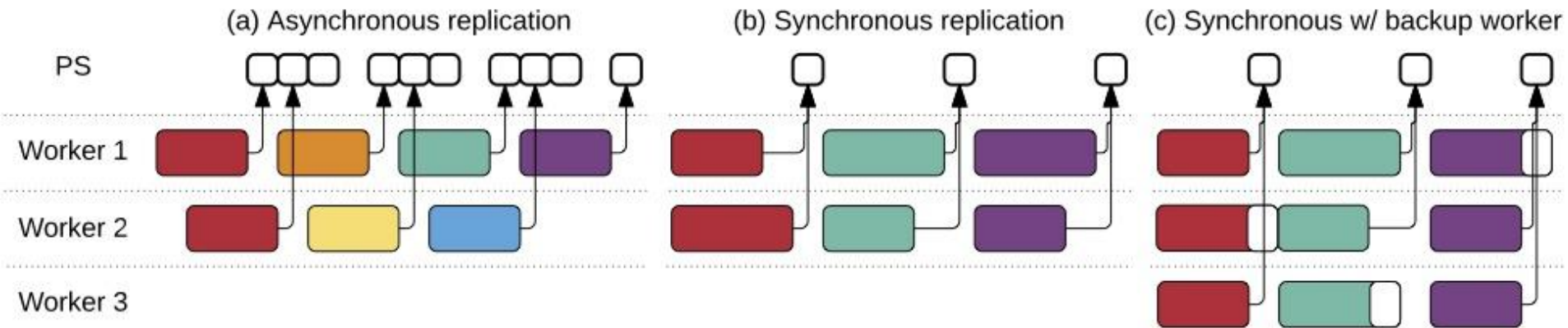
- Unlike DistBelief, no separate parameter server system:
 - Parameters are now just stateful nodes in the graph



Synchronous Variant



Synchronous vs. Asynchronous



Graph structure and low-level graph primitives (queues) allow us to play with synchronous vs. asynchronous update algorithms.

Data Parallelism Considerations

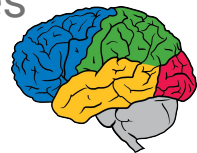
Want model computation time to be large relative to time to send/receive parameters over network

Models with fewer parameters, that reuse each parameter multiple times in the computation

- Mini-batches of size B reuse parameters B times

Certain model structures **reuse each parameter** many times within each example:

- **Convolutional models** tend to reuse hundreds or thousands of times per example (for different spatial positions)
- **Recurrent models** (LSTMs, RNNs) tend to reuse tens to hundreds of times (for unrolling through T time steps during training)



Success of Data Parallelism

- Data parallelism is **really important** for many of Google's problems (very large datasets, large models):
 - RankBrain uses 500 replicas
 - ImageNet Inception training uses 50 GPUs, ~40X speedup
 - SmartReply uses 16 replicas, each with multiple GPUs
 - State-of-the-art on LM "One Billion Word" Benchmark model uses both data and model parallelism on 32 GPUs



Image Model Training Time

Precision @ 1

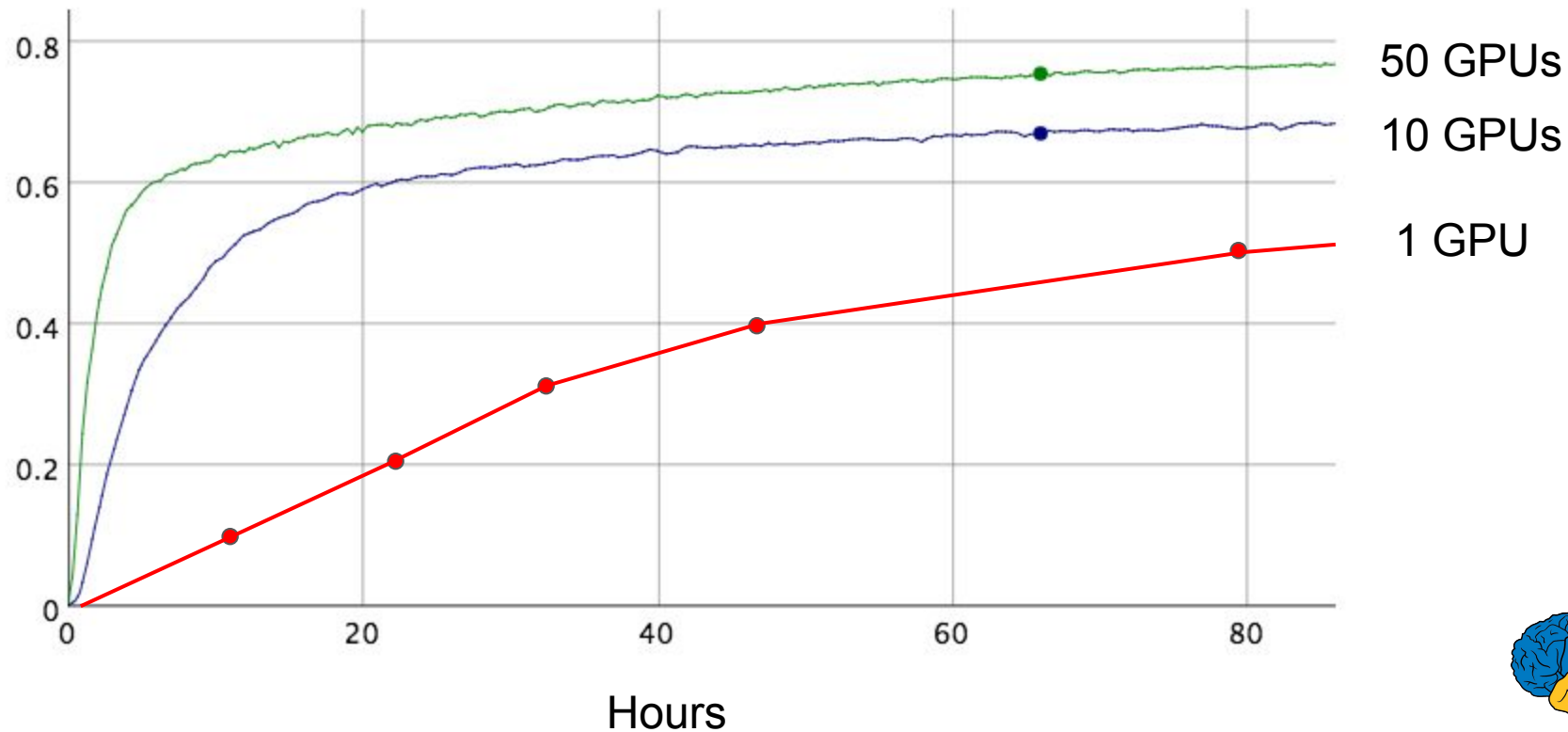
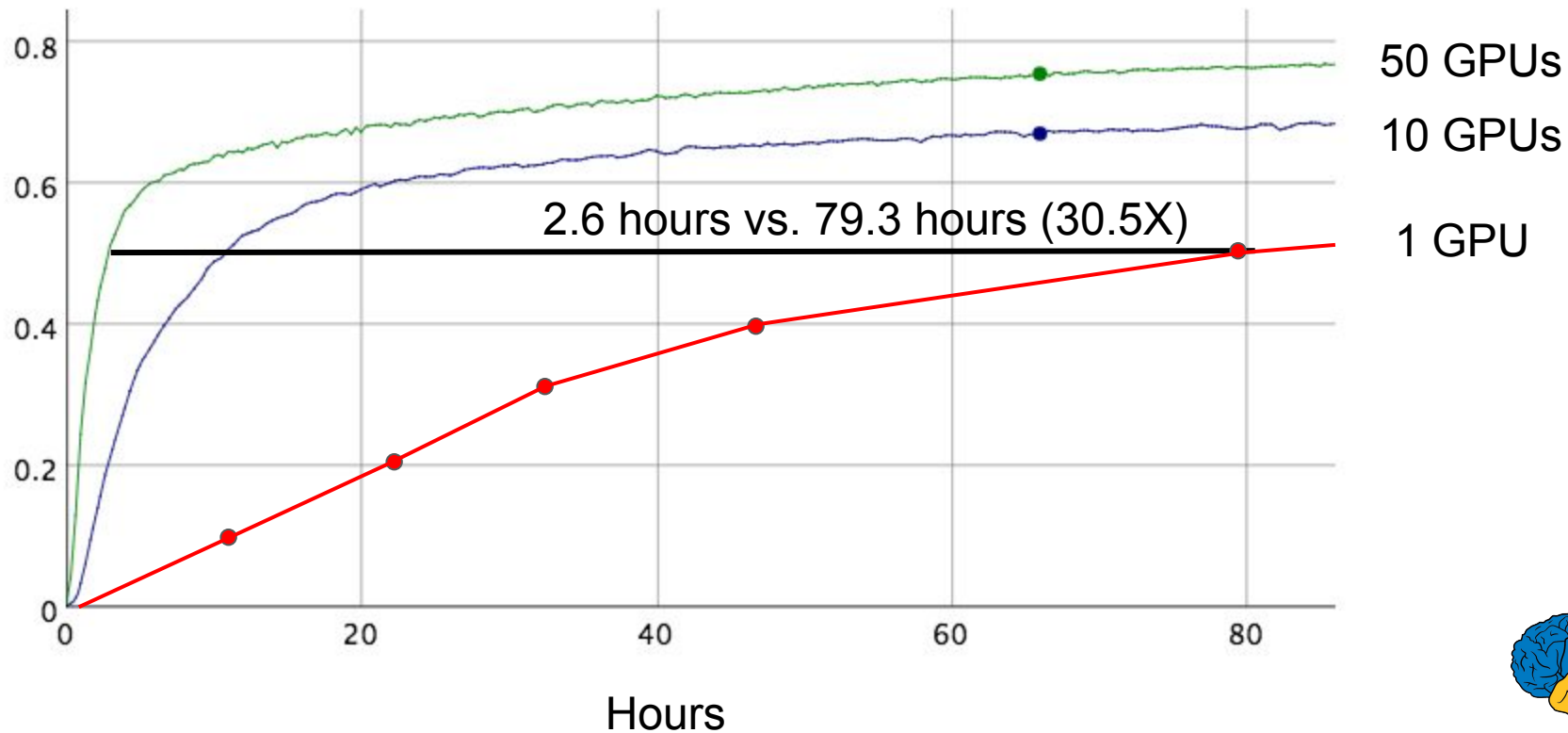
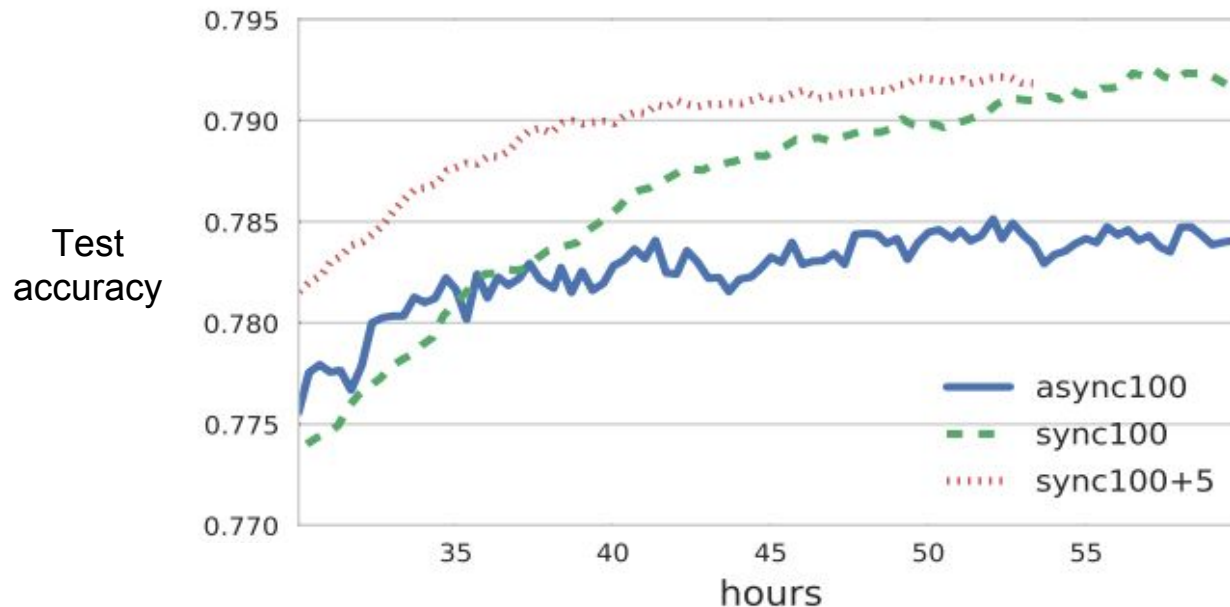


Image Model Training Time

Precision @ 1



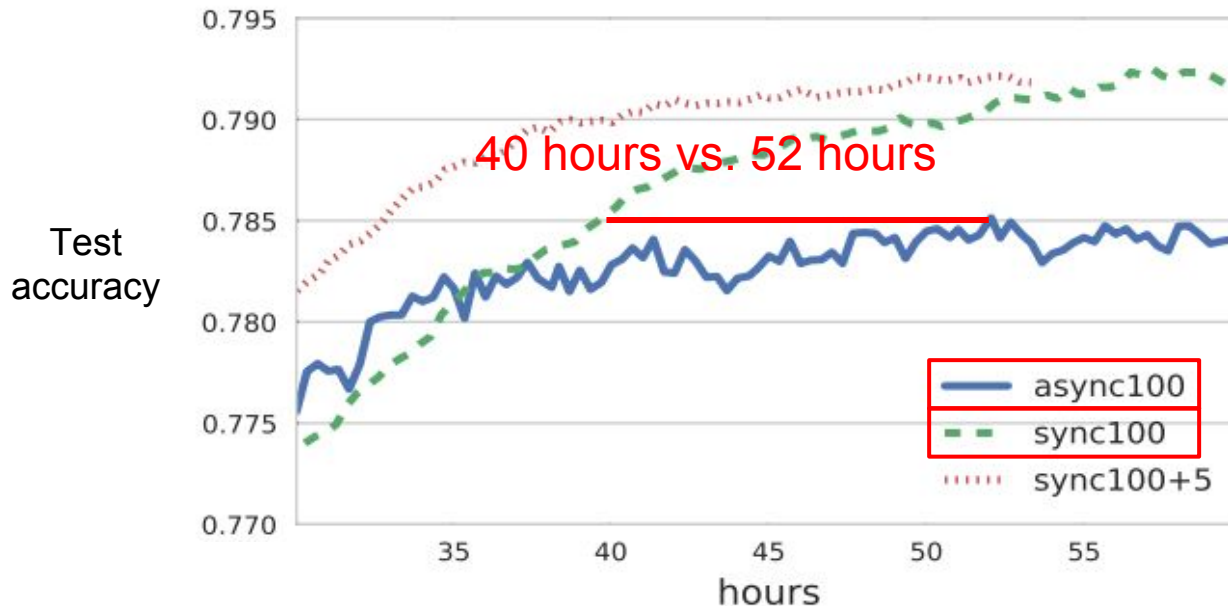
Synchronous converges faster (time to accuracy)



Synchronous updates (with backup workers) trains to higher accuracy faster
Better scaling to more workers (less loss of accuracy)

Revisiting Distributed Synchronous SGD, Jianmin Chen, Rajat Monga, Samy Bengio, Raal Jozefowicz, ICLR Workshop 2016, arxiv.org/abs/1604.00981

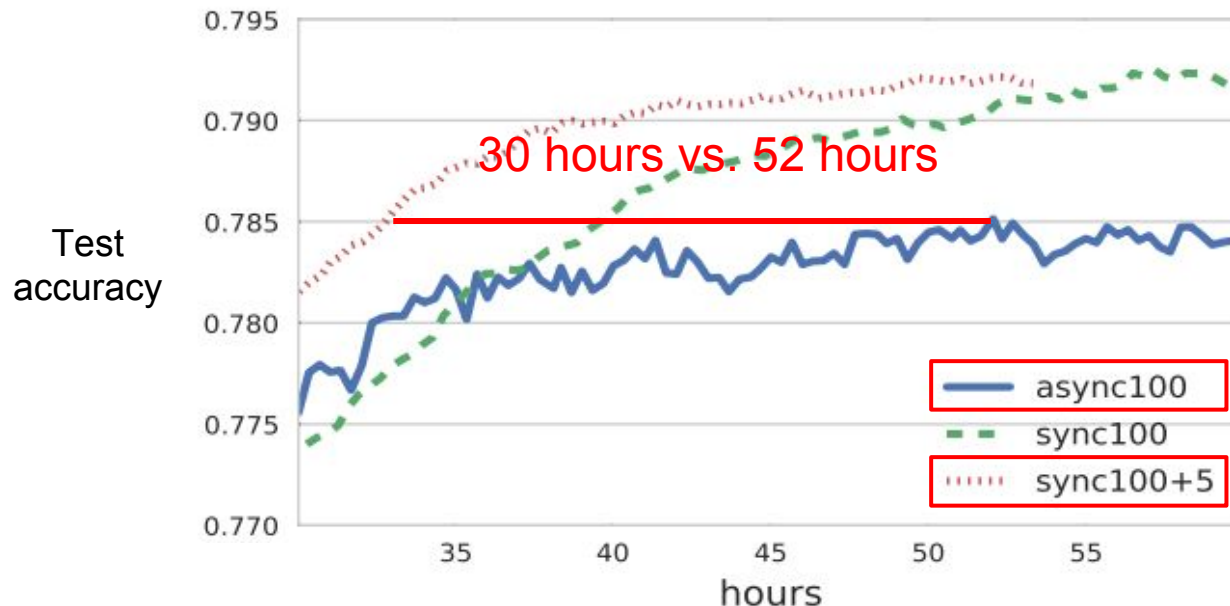
Synchronous converges faster (time to accuracy)



Synchronous updates (with backup workers) trains to higher accuracy faster
Better scaling to more workers (less loss of accuracy)

Revisiting Distributed Synchronous SGD, Jianmin Chen, Rajat Monga, Samy Bengio, Raal Jozefowicz, ICLR Workshop 2016, arxiv.org/abs/1604.00981

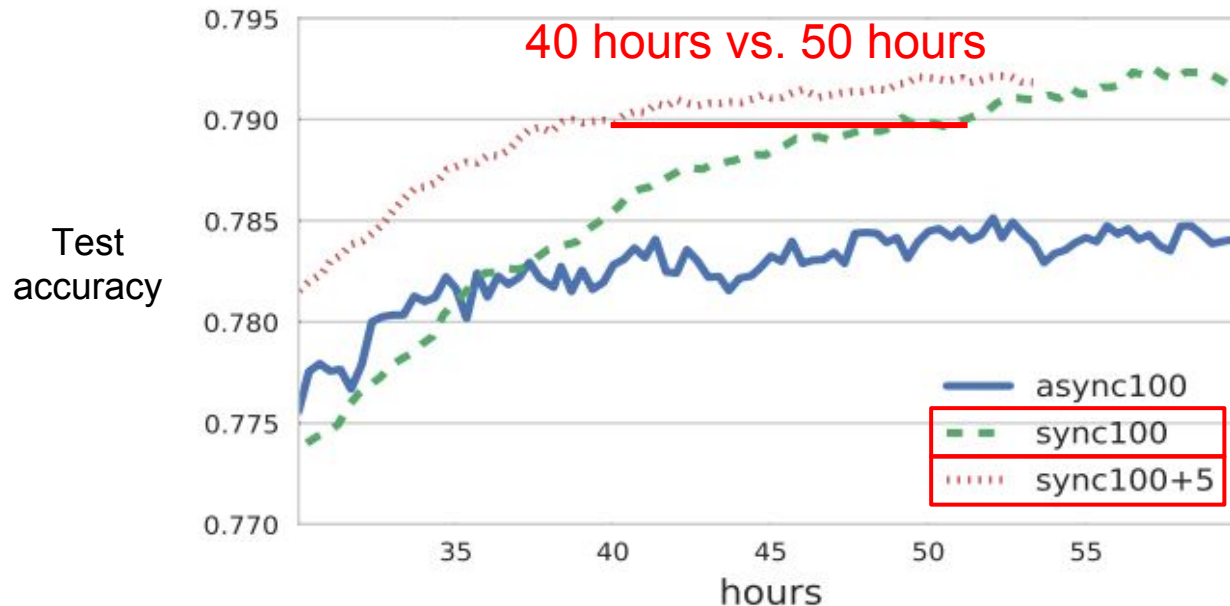
Synchronous converges faster (time to accuracy)



Synchronous updates (with backup workers) trains to higher accuracy faster
Better scaling to more workers (less loss of accuracy)

Revisiting Distributed Synchronous SGD, Jianmin Chen, Rajat Monga, Samy Bengio, Raal Jozefowicz, ICLR Workshop 2016, arxiv.org/abs/1604.00981

Synchronous converges faster (time to accuracy)



Synchronous updates (with backup workers) trains to higher accuracy faster
Better scaling to more workers (less loss of accuracy)

Revisiting Distributed Synchronous SGD, Jianmin Chen, Rajat Monga, Samy Bengio, Raal Jozefowicz, ICLR Workshop 2016, arxiv.org/abs/1604.00981

General Computations

Although we originally built TensorFlow for our uses around deep neural networks, it's actually quite flexible

Wide variety of machine learning and other kinds of numeric computations easily expressible in the computation graph model



Runs on Variety of Platforms

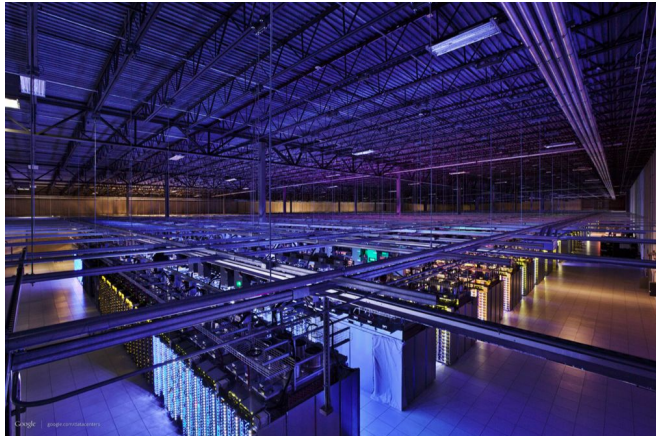
phones



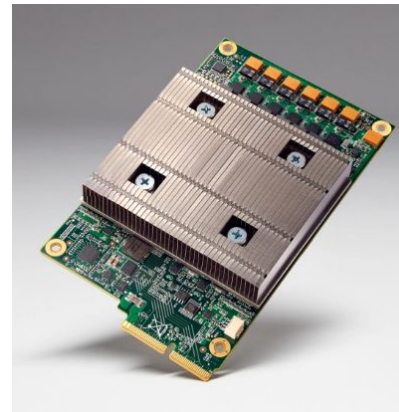
single machines (CPU and/or GPUs) ...



distributed systems of 100s
of machines and/or GPU cards



custom ML hardware



Trend: Much More Heterogeneous hardware

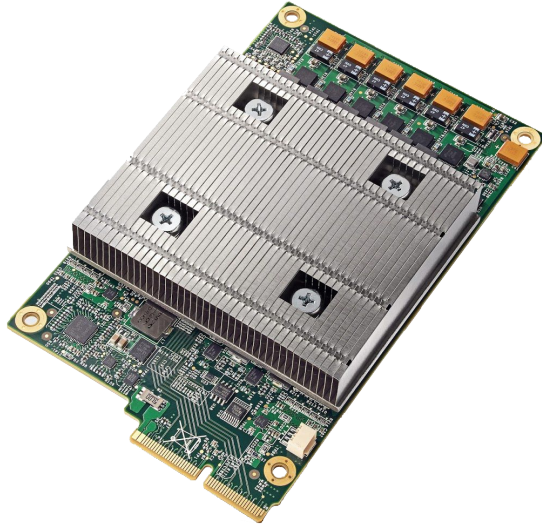
General purpose CPU performance scaling has slowed significantly

Specialization of hardware for certain workloads will be more important



Tensor Processing Unit

Custom machine learning ASIC



In production use for >16 months: used on every search query, used for AlphaGo match, many other uses, ...

See Google Cloud Platform blog: [Google supercharges machine learning tasks with TPU custom chip](#), by Norm Jouppi, May, 2016

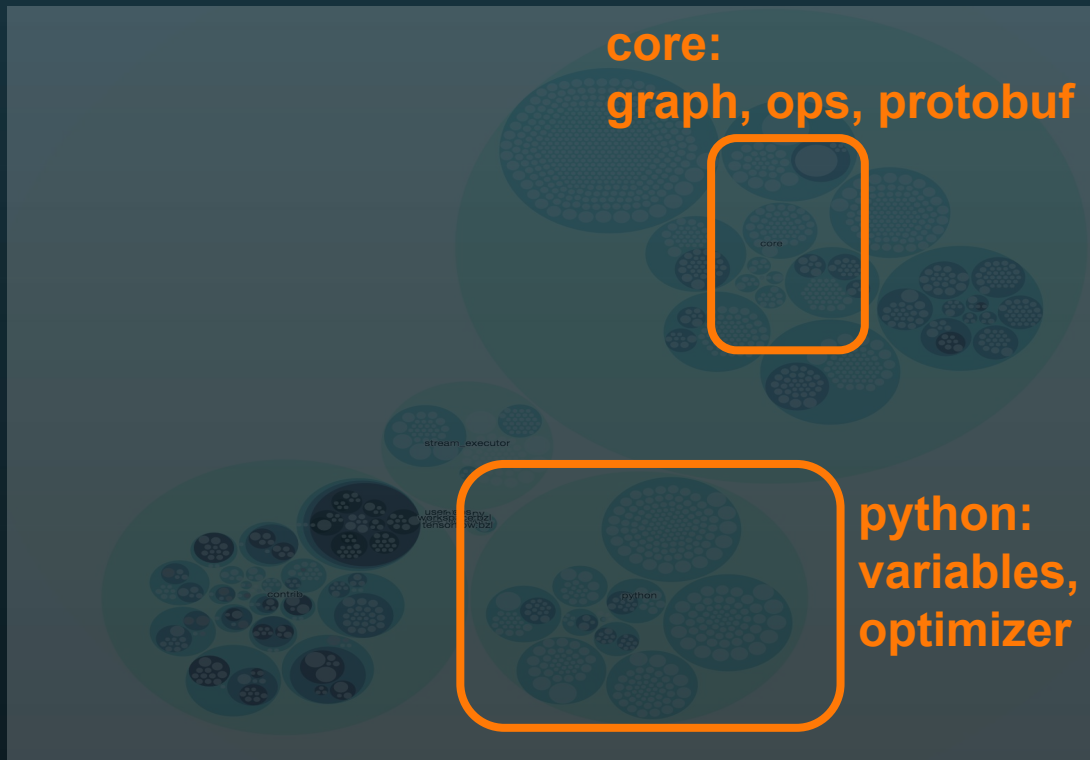
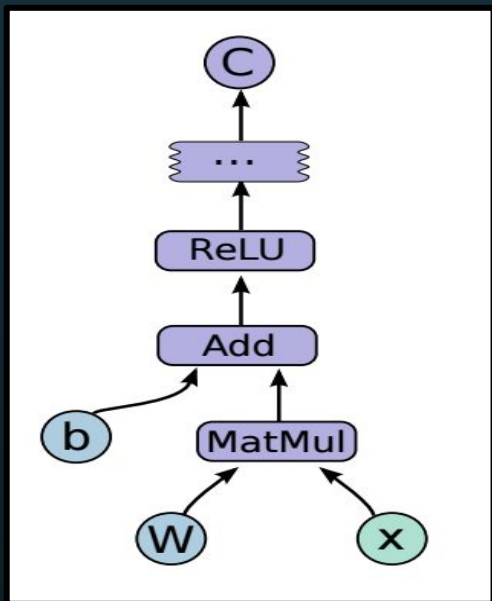
Extensible

- Core system defines a number of standard ***operations*** and ***kernels*** (device-specific implementations of operations)
- Easy to define new operators and/or kernels



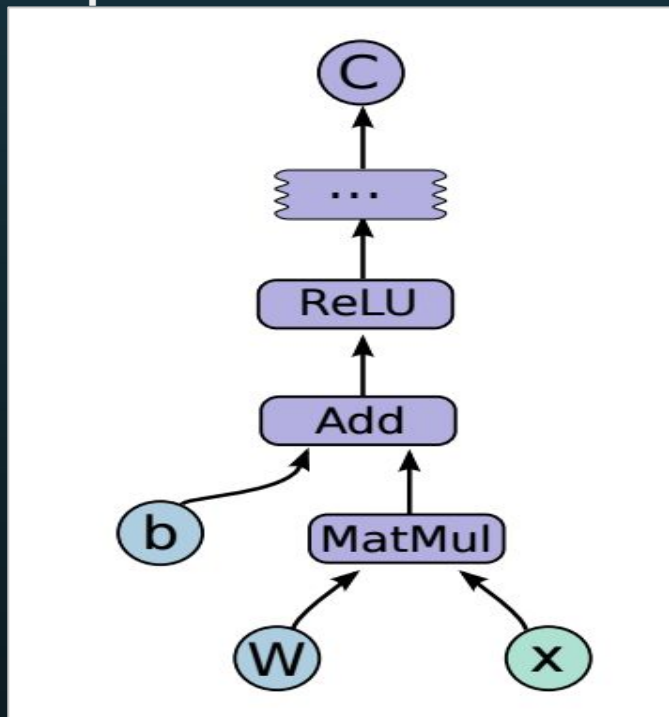
A tour through the TensorFlow codebase

1. Expressing graphs



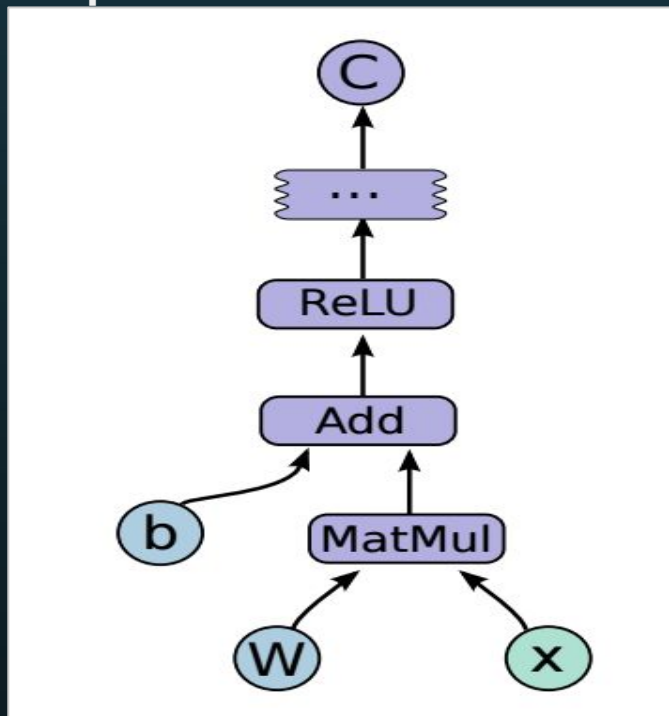
Expressing: Graphs and Ops

Graph

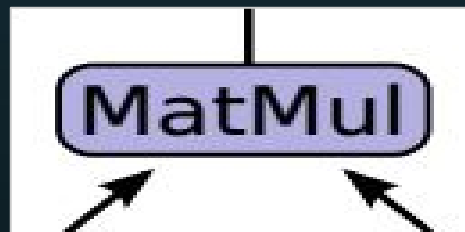
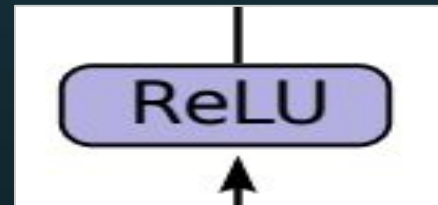


Expressing: Graphs and Ops

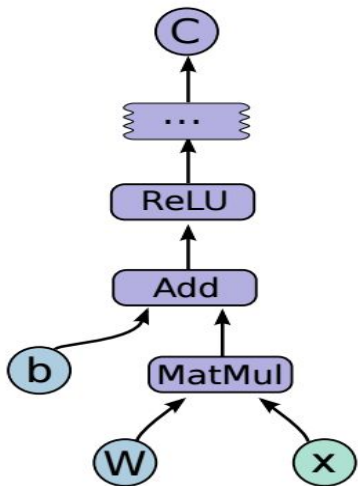
Graph



Ops

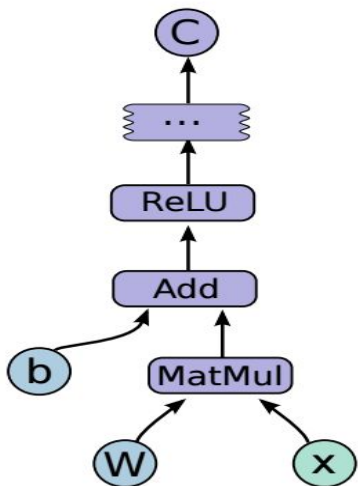


Expressing: Graphs and Ops



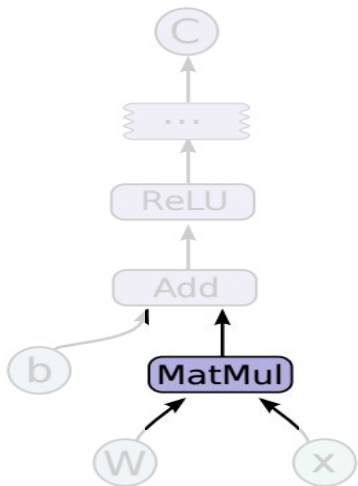
```
1 import tensorflow as tf
2
3 b = tf.Variable(tf.zeros([100]))
4 W = tf.Variable(tf.random_uniform([784,100],-1,1))
5 x = tf.placeholder(tf.float32, name="x")
6 relu = tf.nn.relu(tf.matmul(W, x) + b)
7 cost = # ...
8
9 s = tf.Session()
10 for step in xrange(0, 10):
11     input = # ...read in 100-D input array ...
12     result = s.run(cost, feed_dict={x: input})
13     print step, result
```


Expressing: Ops



```
1 import tensorflow as tf
2
3 b = tf.Variable(tf.zeros([100]))
4 W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
5 x = tf.placeholder(tf.float32, name="x")
6 relu = tf.nn.relu(tf.matmul(W, x) + b)
7 cost = # ...
8
9 s = tf.Session()
10 for step in xrange(0, 10):
11     input = # ...read in 100-D input array ...
12     result = s.run(cost, feed_dict={x: input})
13     print step, result
```

Expressing: Ops



```
1 import tensorflow as tf
2
3 b = tf.Variable(tf.zeros([100]))
4 W = tf.Variable(tf.random_uniform([784,100],-1,1))
5 x = tf.placeholder(tf.float32, name="x")
6 relu = tf.nn.relu(tf.matmul(W, x))
7 cost = # ...
8
9 s = tf.Session()
10 for step in xrange(0, 10):
11     input = # ...read in 100-D input array ...
12     result = s.run(cost, feed_dict={x: input})
13     print step, result
```

Expressing: Ops

```
tf.matmul(W, x)
```

in [math_ops.py#L1137](#)

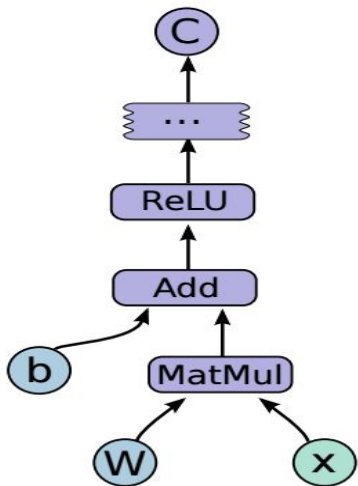
```
return gen_math_ops._mat_mul(a, b,  
                               transpose_a=transpose_a,  
                               transpose_b=transpose_b,  
                               name=name)
```

calls C++ wrappers generated by [cc/BUILD#L27](#)

OpDef interface defined in [math_ops.cc#L607](#)

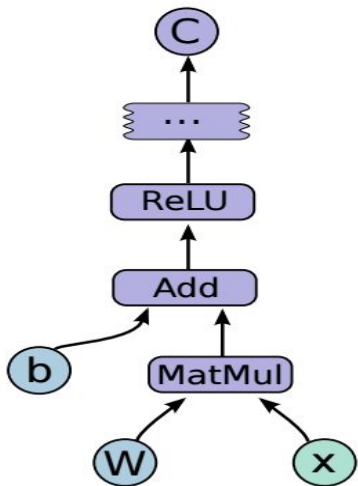
```
REGISTER_OP("MatMul")  
  .Input("a: T")  
  .Input("b: T")  
  .Output("product: T")  
  .Attr("transpose_a: bool = false")  
  .Attr("transpose_b: bool = false")  
  .Attr("T: {float, double, int32, complex64}")
```

Expressing: Graph



```
1 import tensorflow as tf
2
3 b = tf.Variable(tf.zeros([100]))
4 W = tf.Variable(tf.random_uniform([784,100],-1,1))
5 x = tf.placeholder(tf.float32, name="x")
6 relu = tf.nn.relu(tf.matmul(W, x))
7 cost = # ...
8
9 s = tf.Session()
10 for step in xrange(0, 10):
11     input = # ...read in 100-D input array ...
12     result = s.run(cost, feed_dict={x: input})
13     print step, result
```

Expressing: Graph



```
1 import tensorflow as tf
2
3 b = tf.Variable(tf.zeros([100]))
4 W = tf.Variable(tf.random_uniform([784,100],-1,1))
5 x = tf.placeholder(tf.float32, name="x")
6 relu = tf.nn.relu(tf.matmul(W, x) + b)
7 cost = # ...
8
9 s = tf.Session()
10 for step in xrange(0, 10):
11     input = # ...read in 100-D input array ...
12     result = s.run(cost, feed_dict={x: input})
13     print step, result
```

Expressing: Graph

Graph is built implicitly

[session.py#L896](#)

```
tf.matmul(W, x)
print(tf.get_default_graph().as_graph_def())
```

Expressing: Graph

Graph is built implicitly

[session.py#L896](#)

```
tf.matmul(W, x)
print(tf.get_default_graph().as_graph_def())
```

Variables add implicit ops

[variables.py#L146](#)

```
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
print(tf.get_default_graph().as_graph_def())
```

Expressing: Graph

Graph is built implicitly

[session.py#L896](#)

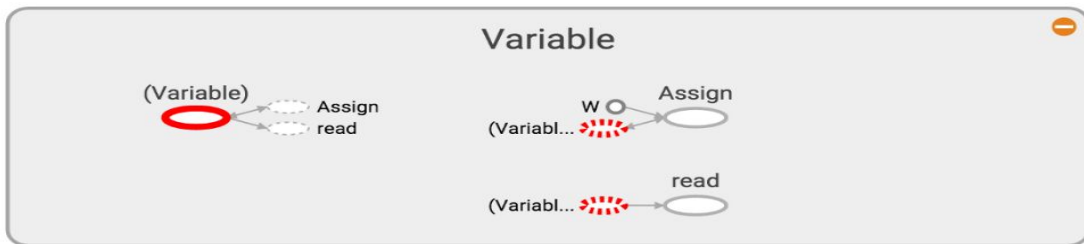
Variables add implicit ops

[variables.py#L146](#)

```
tf.matmul(W, x)
print(tf.get_default_graph().as_graph_def())
```

```
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
print(tf.get_default_graph().as_graph_def())
```

In TensorBoard:



Variable/(Variable)
Operation: Variable

Attributes (4)

- container {"s":""}
- dtype {"type":"DT_FLOAT"}
- shape {"shape":{"dim":[{"size":784}, {"size":10}]}}
- shared_name{"s":""}

Inputs (0)

Outputs (2)

- Variable/Assign
- Variable/read

Add to main graph

Expressing: Optimizers

Optimizer fns extend the graph

[optimizer.py:minimize#L155](#)

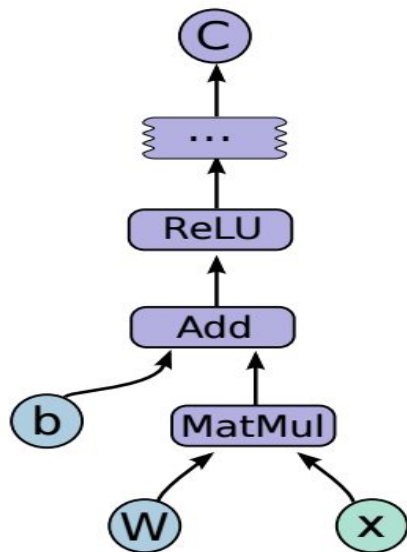
```
optimizer = tf.train.GradientDescentOptimizer(0.01)
train_step = optimizer.minimize(cross_entropy)
```

Expressing: Optimizers

Optimizer fns extend the graph
[optimizer.py:minimize#L155](#)

```
optimizer = tf.train.GradientDescentOptimizer(0.01)
train_step = optimizer.minimize(cross_entropy)
```

Trainable variables collected
[variables.py#L258](#)



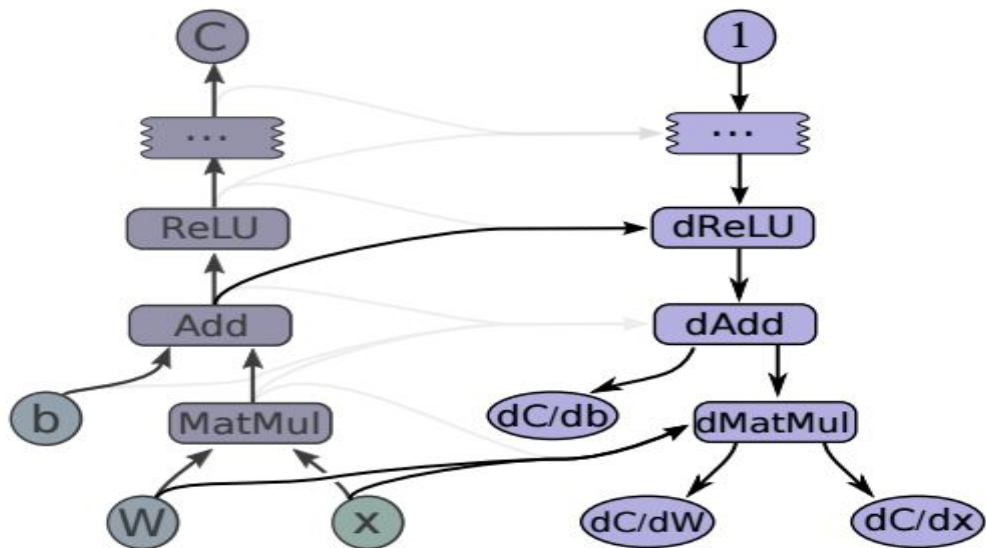
Expressing: Optimizers

Optimizer fns extend the graph
[optimizer.py:minimize#L155](#)

```
optimizer = tf.train.GradientDescentOptimizer(0.01)
train_step = optimizer.minimize(cross_entropy)
```

Trainable variables collected
[variables.py#L258](#)

Graph is extended with gradients
[gradients.py#L307](#)

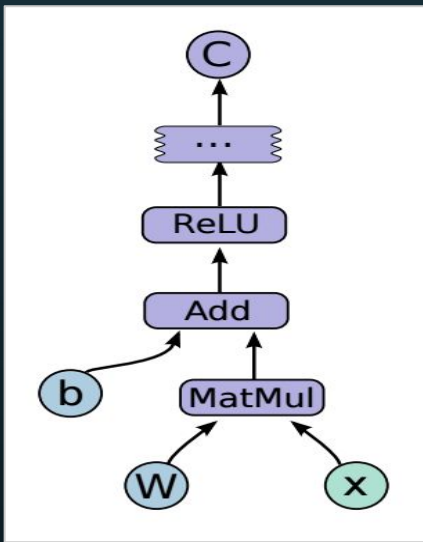


Expressing: Graph

Serialized as GraphDef

[graph.proto](#)

```
print(tf.get_default_graph().as_graph_def())
```

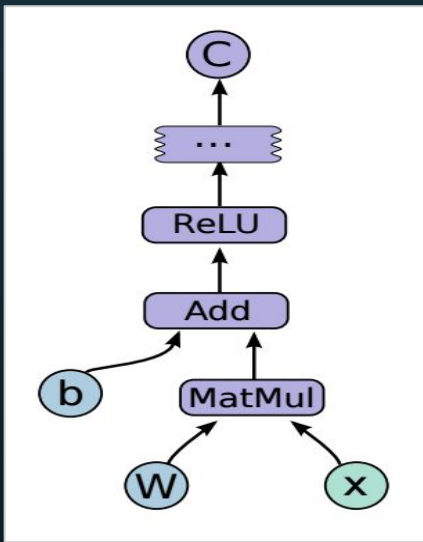


Expressing: Graph

Serialized as GraphDef

[graph.proto](#)

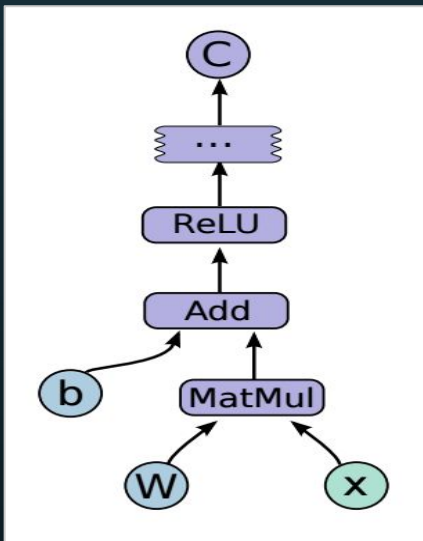
```
print(tf.get_default_graph().as_graph_def())
```



Expressing: Graph

Serialized as GraphDef

[graph.proto](#)



```
print(tf.get_default_graph().as_graph_def())
```

```
node {
  name: "MatMul"
  op: "MatMul"
  input: "W/read"
  input: "x"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
  attr {
    key: "transpose_a"
    value {
      b: false
    }
  }
  attr {
    key: "transpose_b"
    value {
      b: false
    }
  }
}
```

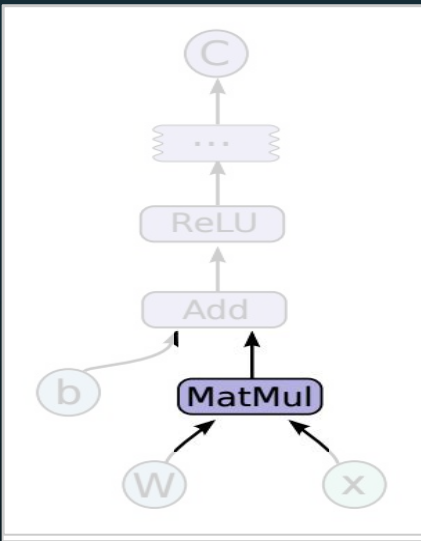
```
node {
  name: "add"
  op: "Add"
  input: "MatMul"
  input: "b/read"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
}
```

```
node {
  name: "Relu"
  op: "Relu"
  input: "add"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
}
```

Expressing: Graph

Serialized as GraphDef

[graph.proto](#)



```
print(tf.get_default_graph().as_graph_def())
```

```
node {
  name: "MatMul"
  op: "MatMul"
  input: "W/read"
  input: "x"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
  attr {
    key: "transpose_a"
    value {
      b: false
    }
  }
  attr {
    key: "transpose_b"
    value {
      b: false
    }
  }
}
```

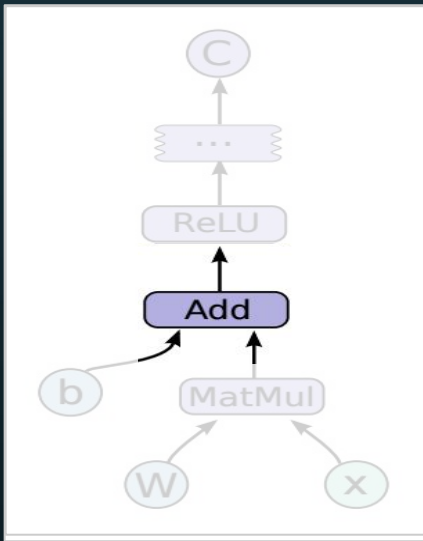
```
node {
  name: "add"
  op: "Add"
  input: "MatMul"
  input: "b/read"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
}
```

```
node {
  name: "Relu"
  op: "Relu"
  input: "add"
  attr {
    key: "T"
    value {
      type: DT_FLOAT
    }
  }
}
```

Expressing: Graph

Serialized as GraphDef

[graph.proto](#)



```
print(tf.get_default_graph().as_graph_def())
```

```
node {  
  name: "MatMul"  
  op: "MatMul"  
  input: "W/read"  
  input: "x"  
  attr {  
    key: "T"  
    value {  
      type: DT_FLOAT  
    }  
  }  
  attr {  
    key: "transpose_a"  
    value {  
      b: false  
    }  
  }  
  attr {  
    key: "transpose_b"  
    value {  
      b: false  
    }  
  }  
}
```

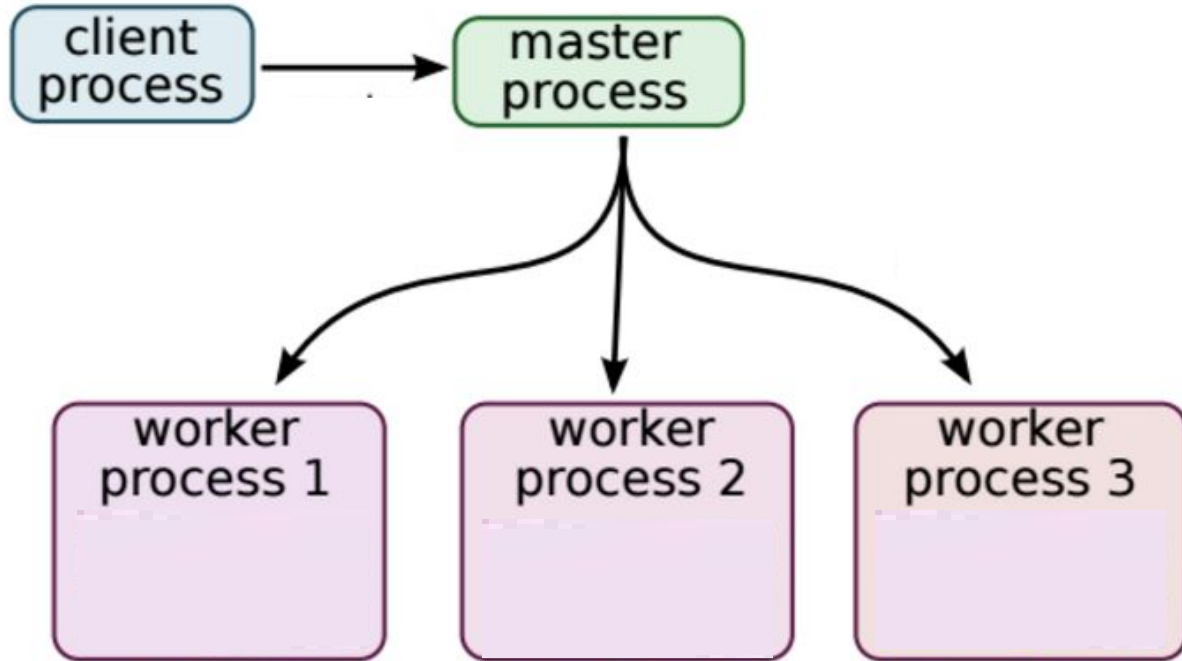
```
node {  
  name: "add"  
  op: "Add"  
  input: "MatMul"  
  input: "b/read"  
  attr {  
    key: "T"  
    value {  
      type: DT_FLOAT  
    }  
  }  
}
```

```
node {  
  name: "Relu"  
  op: "Relu"  
  input: "add"  
  attr {  
    key: "T"  
    value {  
      type: DT_FLOAT  
    }  
  }  
}
```

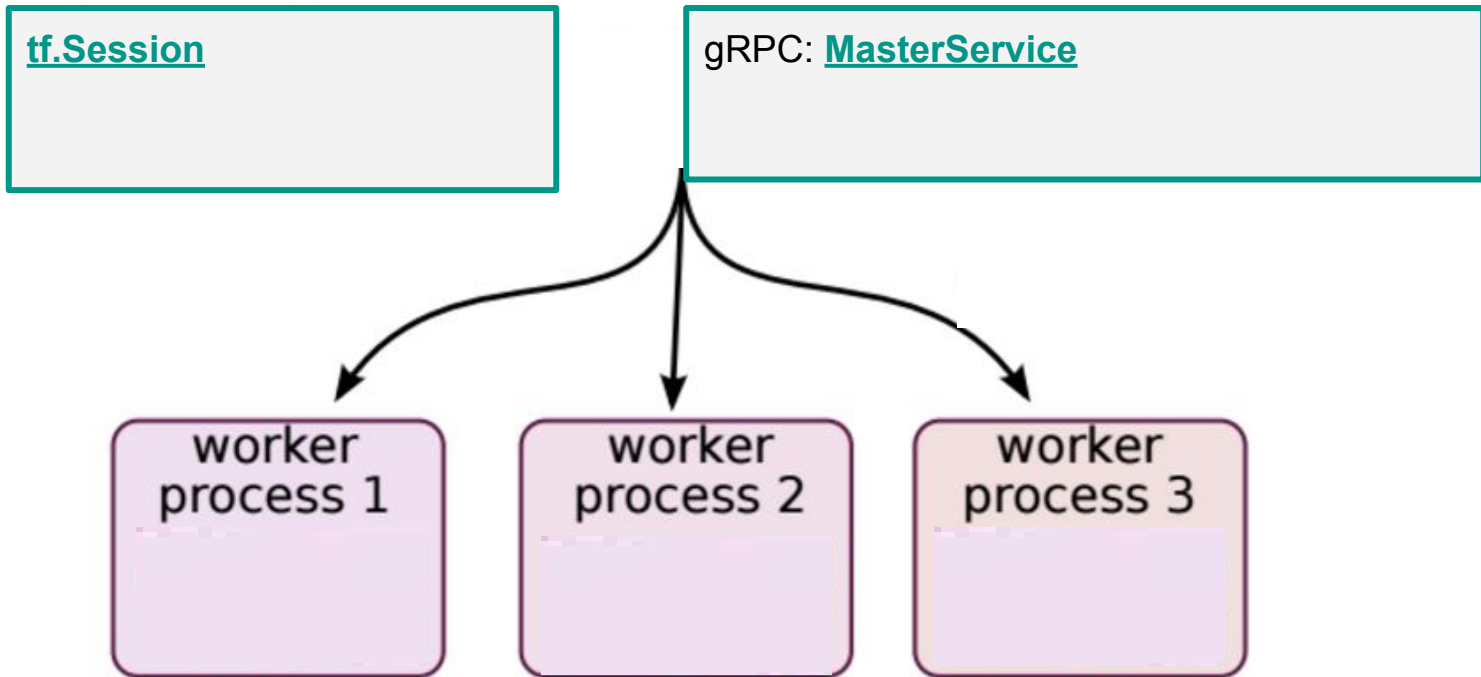

Distributing

- Sessions in distributed runtime
- Pruning
- Placing and Partitioning

Distributing: Creating a session



Distributing: Creating a session



Distributing: Creating a session

tf.Session

gRPC: MasterService

```
import tensorflow as tf

b = tf.Variable(tf.zeros([100]))
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
x = tf.placeholder(tf.float32, name="x")
relu = tf.nn.relu(tf.matmul(W, x) + b)
cost = # ...

s = tf.Session()
```

process 1

process 2

worker
process 3



Distributing: Creating a session

tf.Session

gRPC: Session

gRPC: MasterService

```
import tensorflow as tf

b = tf.Variable(tf.zeros([100]))
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
x = tf.placeholder(tf.float32, name="x")
relu = tf.nn.relu(tf.matmul(W, x) + b)
cost = # ...

s = tf.Session()
```

process 1

process 2

worker
process 3

Distributing: Creating a session

tf.Session

gRPC: Session

gRPC: MasterService
CreateSession(GraphDef)

```
import tensorflow as tf

b = tf.Variable(tf.zeros([100]))
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
x = tf.placeholder(tf.float32, name="x")
relu = tf.nn.relu(tf.matmul(W, x) + b)
cost = # ...

s = tf.Session()
```

process 1

process 2

worker
process 3



Distributing: Creating a session

`tf.Session`

gRPC: `Session`

gRPC: `MasterService`
`CreateSession(GraphDef)`

```
import tensorflow as tf

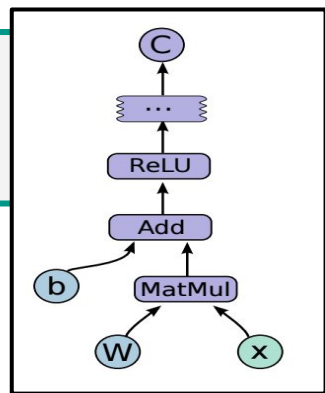
b = tf.Variable(tf.zeros([100]))
W = tf.Variable(tf.random_uniform([784, 100], -1, 1))
x = tf.placeholder(tf.float32, name="x")
relu = tf.nn.relu(tf.matmul(W, x) + b)
cost = # ...

s = tf.Session()
```

process 1

process 2

worker
process 3



Distributing: Running a session

`tf.Session`

gRPC: `Session`

gRPC: `MasterService`

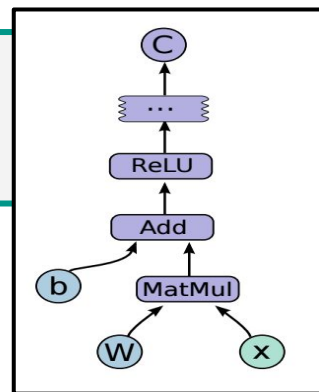
`CreateSession(GraphDef)`

```
result = s.run(cost, feed_dict={x: input})
```

worker
process 1

worker
process 2

worker
process 3



Distributing: Running a session

`tf.Session`

gRPC: `Session`

gRPC: `MasterService`

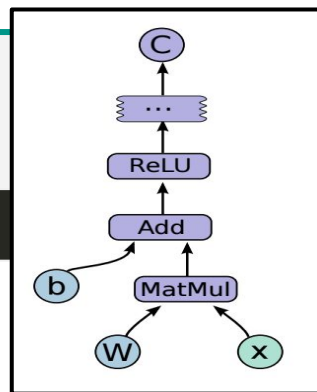
`CreateSession(GraphDef)`
`RunStep(feed, fetches)`

```
result = s.run(cost, feed_dict={x: input})
```

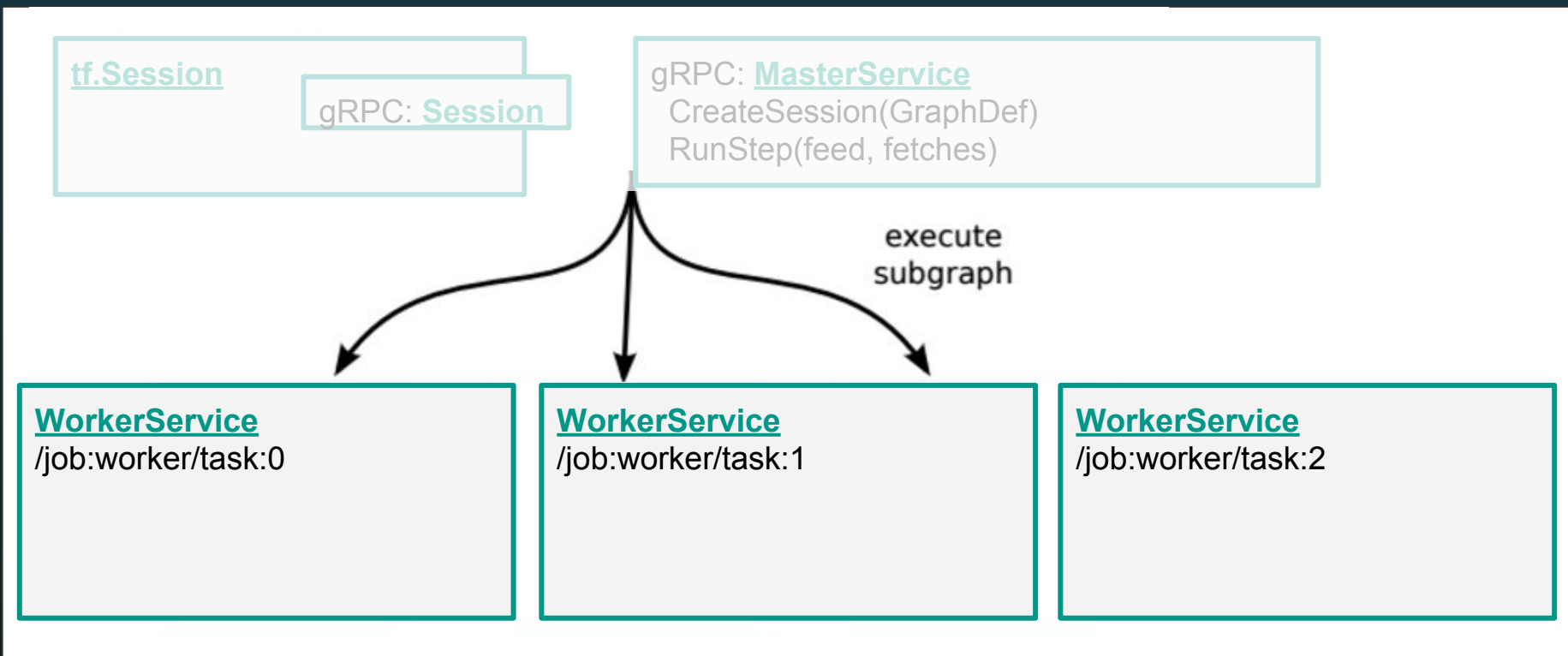
worker
process 1

worker
process 2

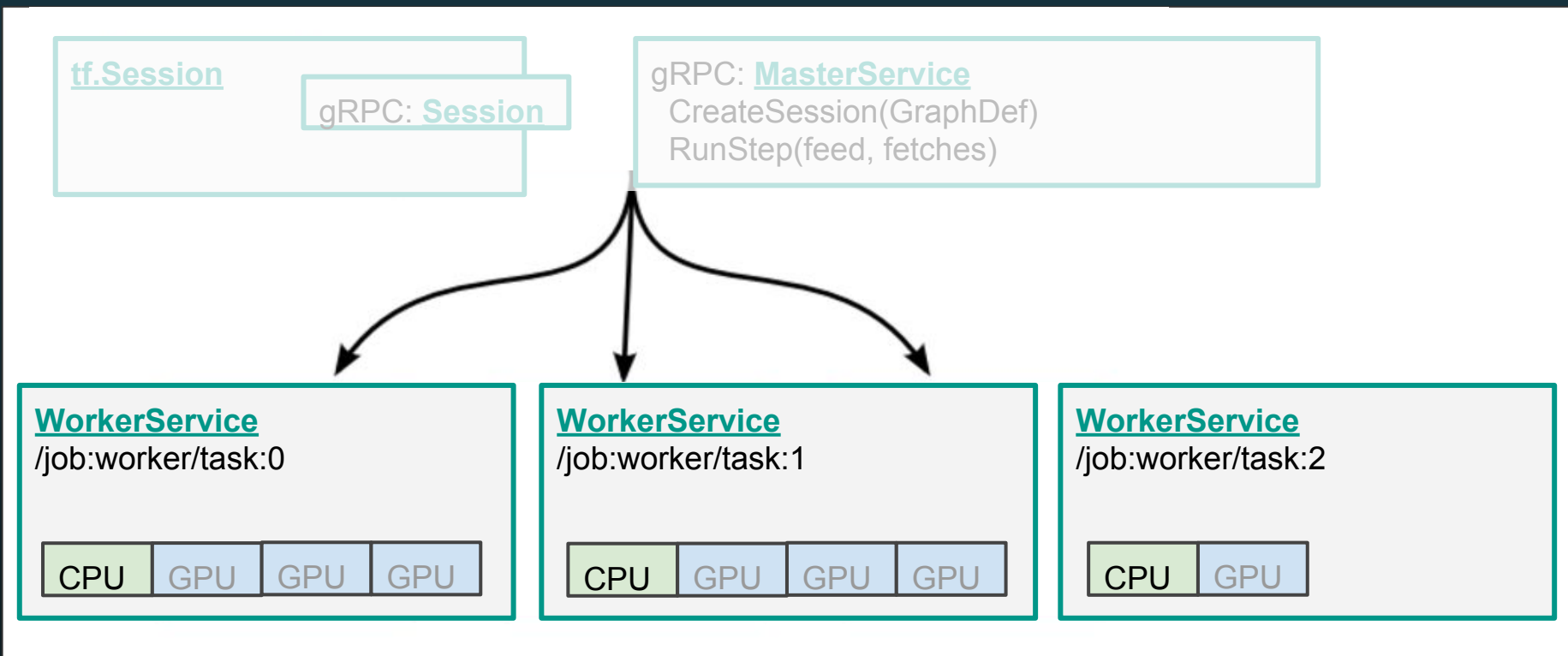
worker
process 3



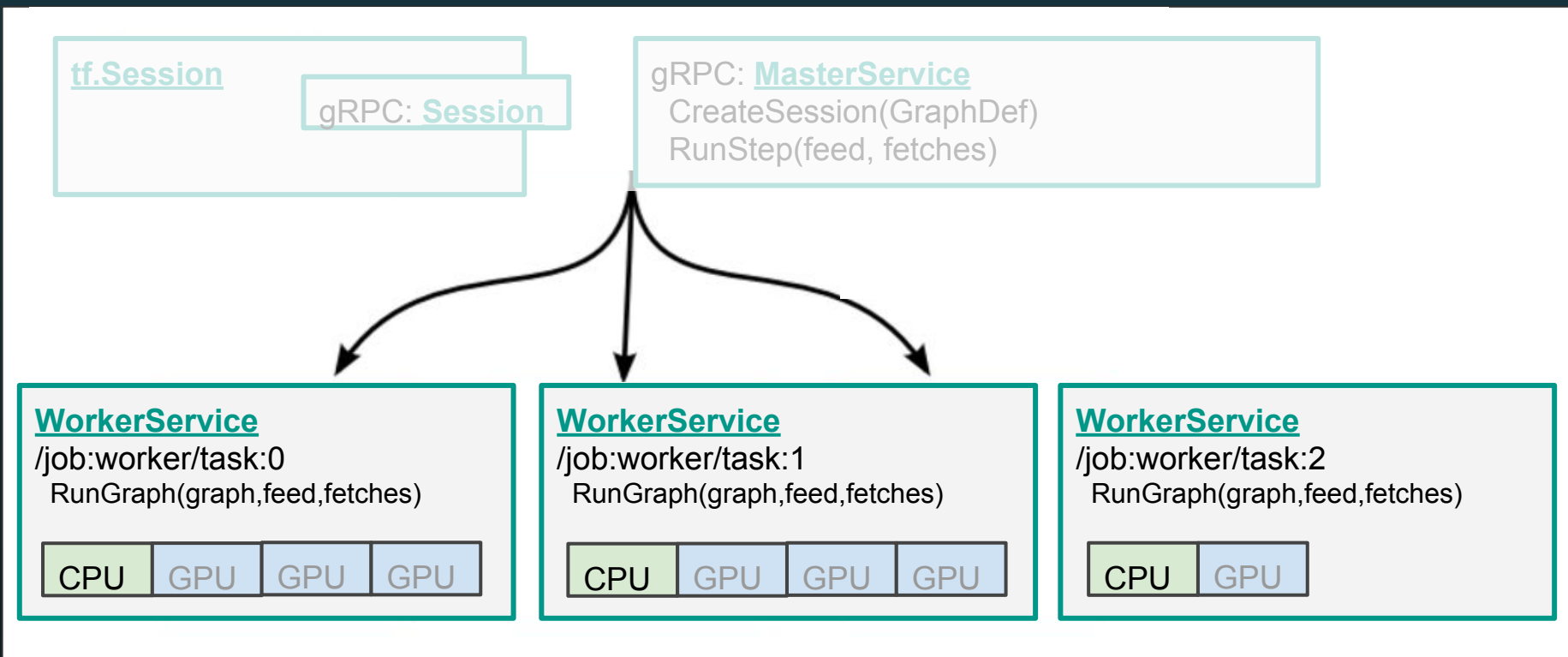
Distributing: Running a session



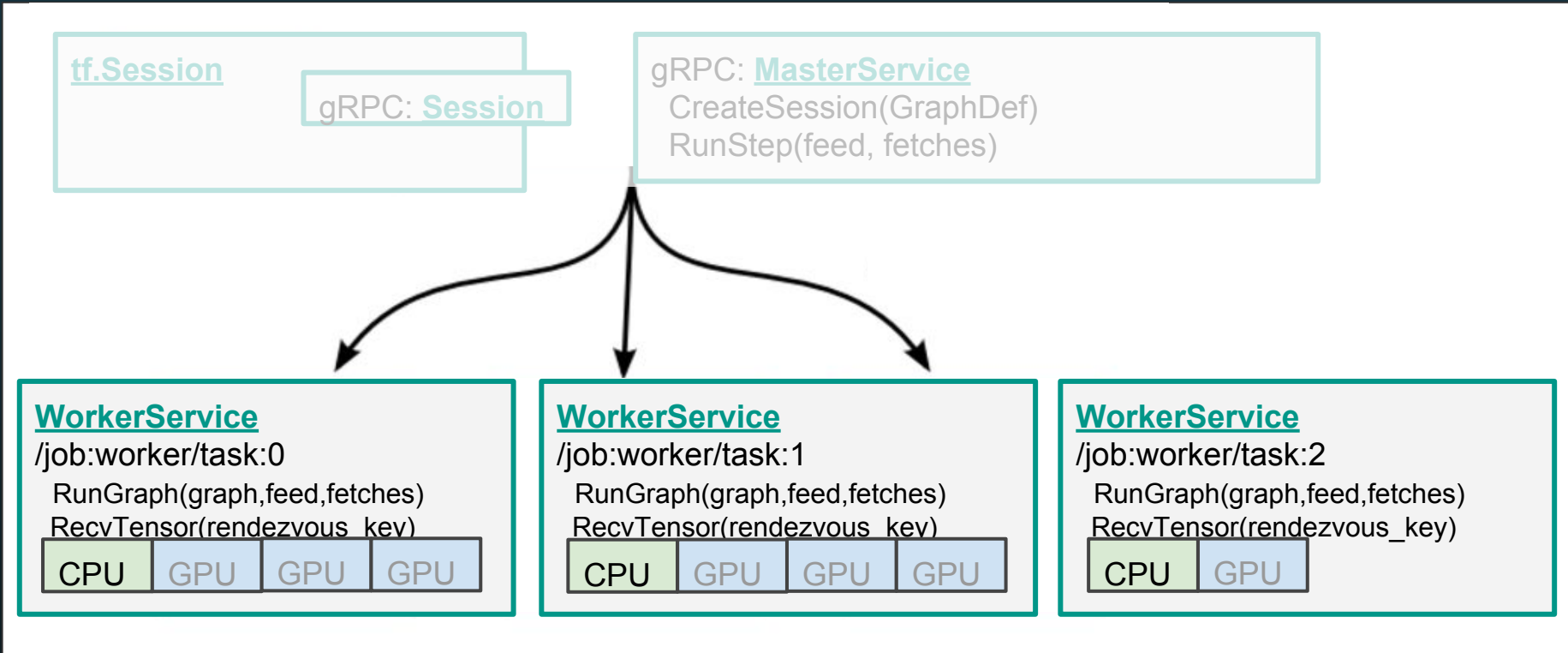
Distributing: Running a session



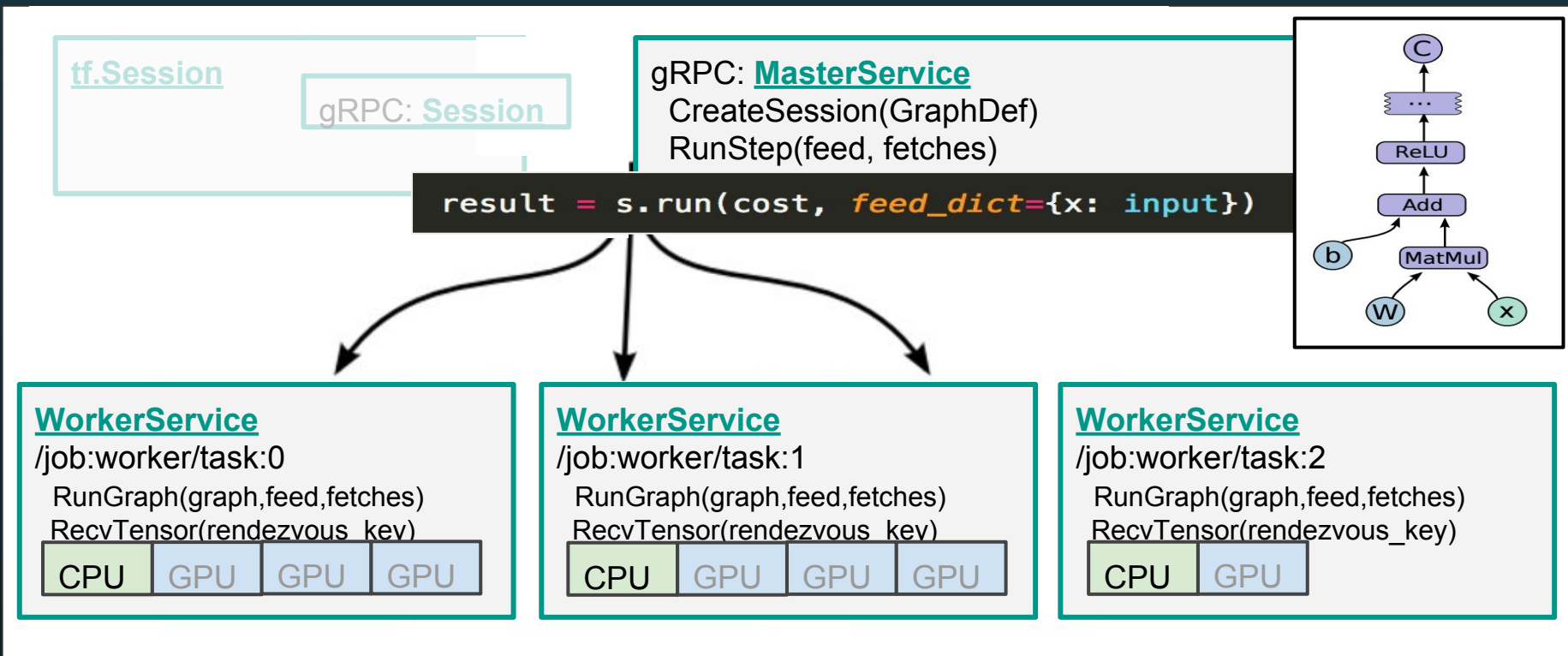
Distributing: Running a session



Distributing: Running a session



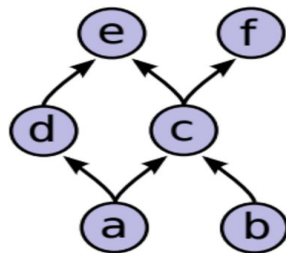
Distributing: Running a session



Distributing: Pruning

gRPC call to `Session::Run`
in [master_session.cc#L835](#)

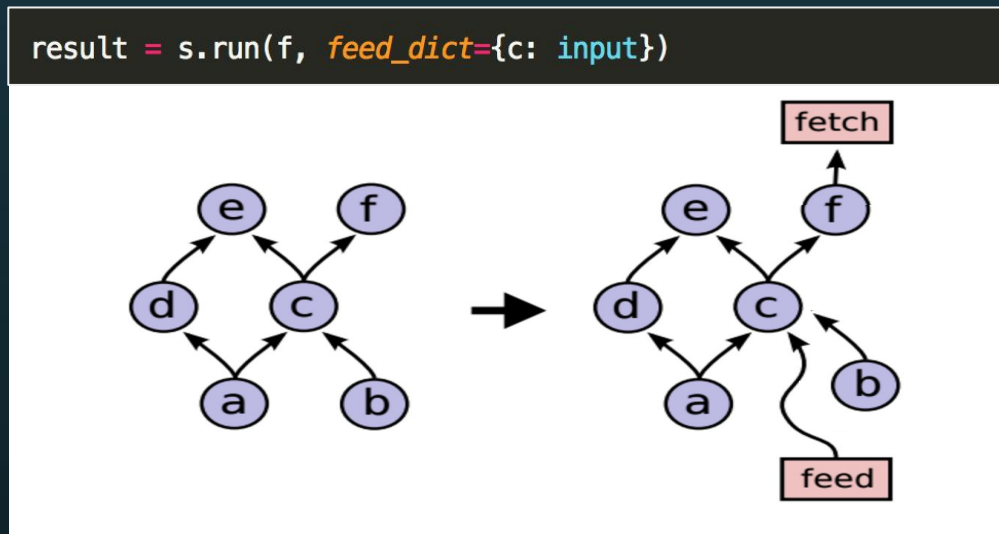
```
result = s.run(f, feed_dict={c: input})
```



Distributing: Pruning

gRPC call to `Session::Run`
in [master_session.cc#L835](#)

Rewrite with feed and fetch
[RewriteGraphForExecution](#)
in [graph/subgraph.cc#L225](#)

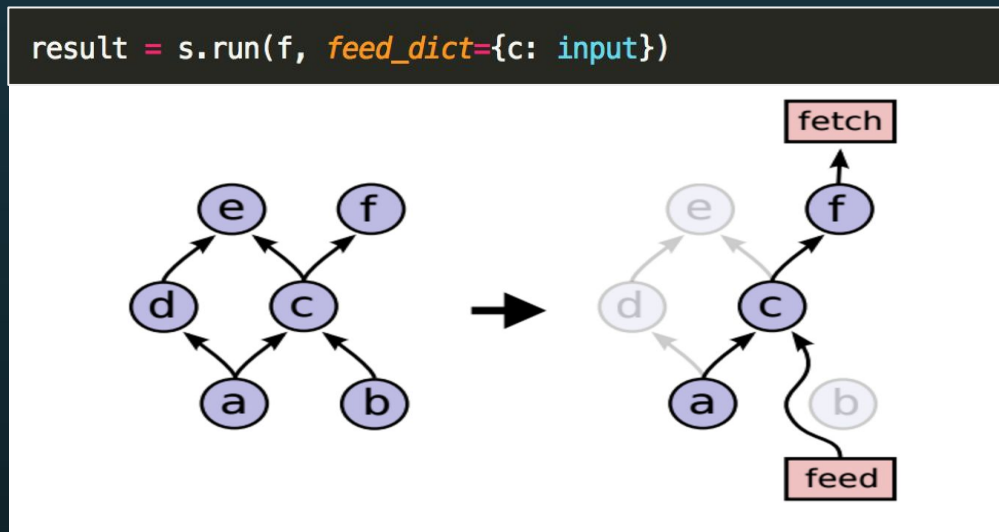


Distributing: Pruning

gRPC call to `Session::Run`
in [master_session.cc#L835](#)

Rewrite with feed and fetch
[RewriteGraphForExecution](#)
in [graph/subgraph.cc#L225](#)

Prune subgraph
[PruneForReverseReachability](#)
in [graph/algorithm.cc#L122](#)
tests in [subgraph_test.cc#142](#)



Distributing: Placing

Constraints from model

[DeviceSpec in device.py#L24](#)

```
with tf.device("/job:ps/task:0"):
    weights_1 = tf.Variable(...)
    biases_1 = tf.Variable(...)

with tf.device("/job:ps/task:1"):
    weights_2 = tf.Variable(...)
    biases_2 = tf.Variable(...)

with tf.device("/job:worker/task:7"):
    input, labels = ...
    layer_1 = tf.nn.relu(tf.matmul(input, weights_1))
    logits = tf.nn.relu(tf.matmul(layer_1, weights_2))
    # ...
```

Distributing: Placing

Constraints from model

[DeviceSpec in device.py#L24](#)

By device or colocation

[NodeDef in graph.proto](#)

```
with tf.device("/job:ps/task:0"):
    weights_1 = tf.Variable(...)
    biases_1 = tf.Variable(...)

with tf.device("/job:ps/task:1"):
    weights_2 = tf.Variable(...)
    biases_2 = tf.Variable(...)

with tf.device("/job:worker/task:7"):
    input, labels = ...
    layer_1 = tf.nn.relu(tf.matmul(input, weights_
    logits = tf.nn.relu(tf.matmul(layer_1, weights
    # ...
```

```
graph { node { device: "" }}
```

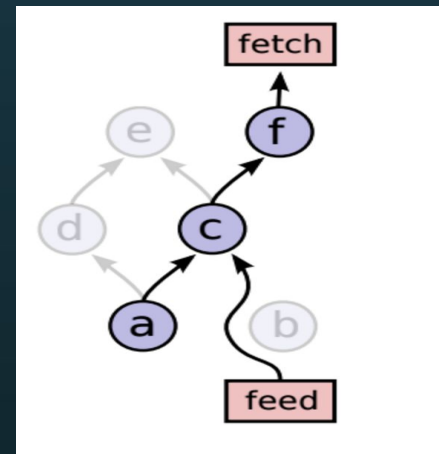
Distributing: Placing

Placing based on constraints

[SimplePlacer::Run](#)

in [simple_placer.cc#L558](#)

described in [simple_placer.h#L31](#)



WorkerService

/job:worker/task:0



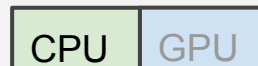
WorkerService

/job:worker/task:1



WorkerService

/job:worker/task:2



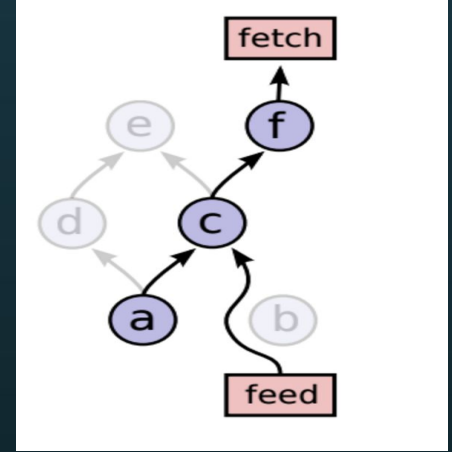
Distributing: Placing

Placing based on constraints

[SimplePlacer::Run](#)

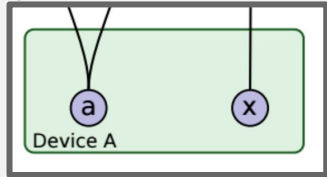
in [simple_placer.cc#L558](#)

described in [simple_placer.h#L31](#)



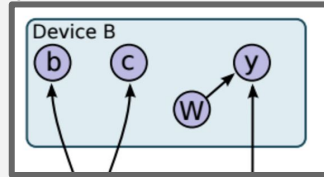
WorkerService

/job:worker/task:0



WorkerService

/job:worker/task:1



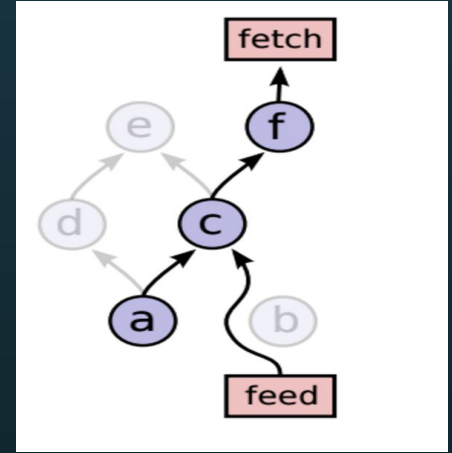
Distributing: Placing

Placing based on constraints

[SimplePlacer::Run](#)

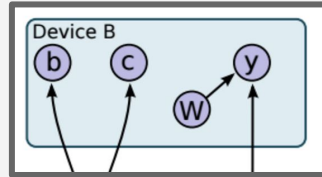
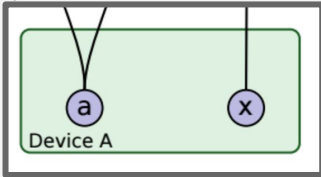
in [simple_placer.cc#L558](#)

described in [simple_placer.h#L31](#)



WorkerService

/job:worker/task:0



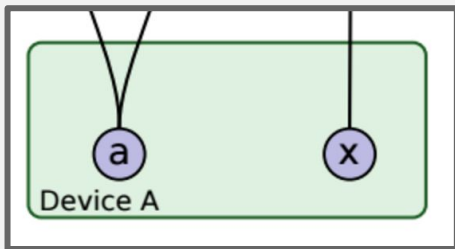
Distributing: Partitioning

Partition into subgraphs

in graph_partition.cc#L883

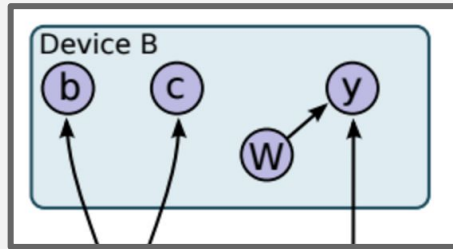
WorkerService

/job:worker/task:0



WorkerService

/job:worker/task:0



Distributing: Partitioning

Partition into subgraphs

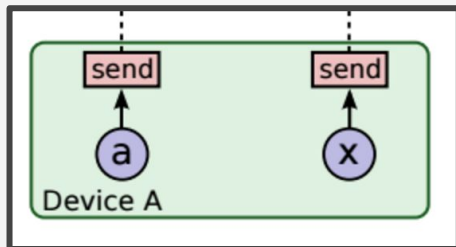
in [graph_partition.cc#L883](#)

Rewrite with **Send** and **Recv**

in [sendrecv_ops.cc#L56](#) and [#L97](#)

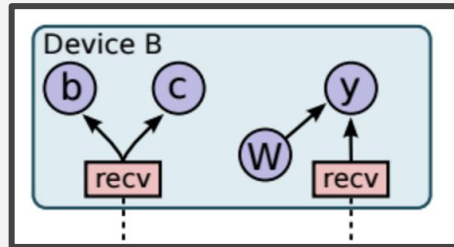
WorkerService

/job:worker/task:0



WorkerService

/job:worker/task:0



Distributing: Partitioning

Partition into subgraphs

in [graph_partition.cc#L883](#)

Rewrite with **Send** and **Recv**

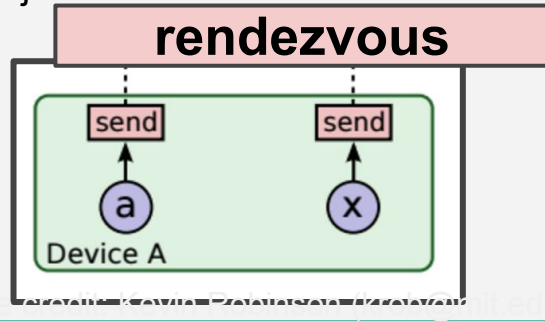
in [sendrecv_ops.cc#L56](#) and [#L97](#)

Rendezvous handles coordination

in [base_rendezvous_mgr.cc#L236](#)

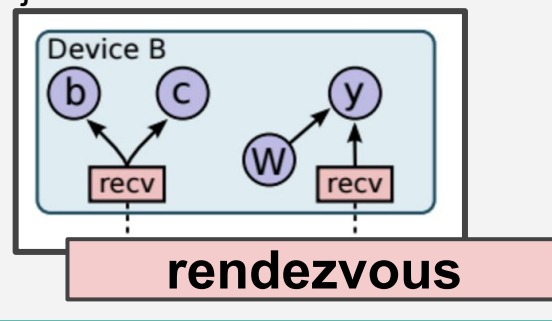
WorkerService

/job:worker/task:0



WorkerService

/job:worker/task:0



Distributing: Partitioning

Partition into subgraphs

in [graph_partition.cc#L883](#)

Rewrite with **Send** and **Recv**

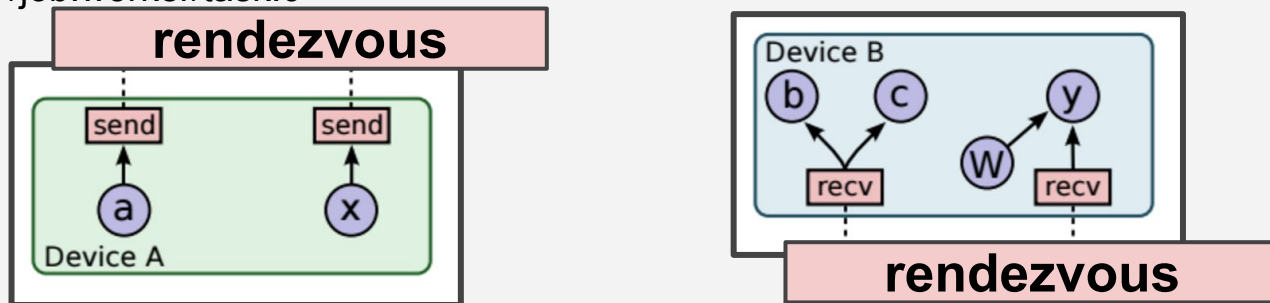
in [sendrecv_ops.cc#L56](#) and [#L97](#)

Rendezvous handles coordination

in [base rendezvous_mgr.cc#L236](#)

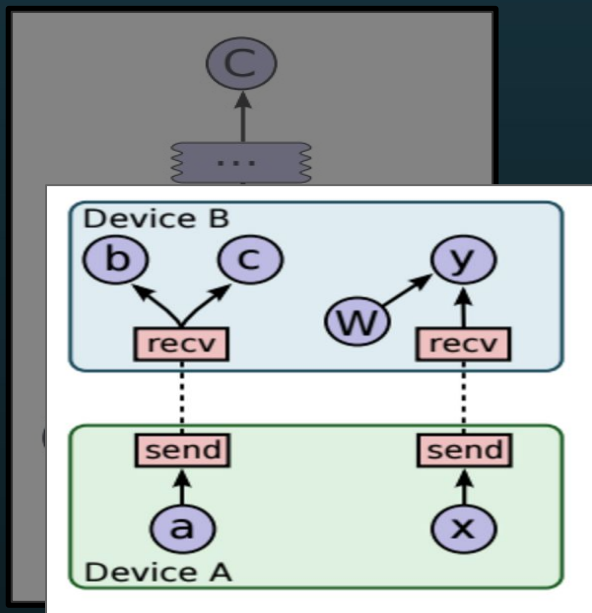
WorkerService

/job:worker/task:0

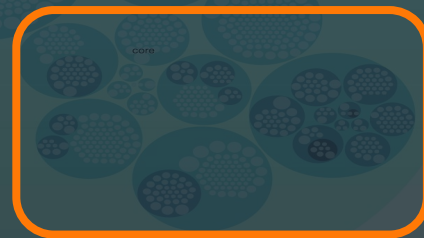


A tour through the TensorFlow codebase

2. Distributing the graph



core:
distributed_runtime
common_runtime



Executing: Executor

Parallelism on each worker

WorkerService

/job:worker/task:0

RunGraph(graph,feed,fetches)

RecvTensor(rendezvous_key)

CPU

GPU

GPU

GPU

Executing: Executor

Parallelism on each worker

```
WorkerService  
/job:worker/task:0  
RunGraph(graph,feed,fetches)  
RecvTensor(rendezvous_key)  
CPU GPU GPU GPU
```

Executing: Executor

Parallelism on each worker

WorkerService

/job:worker/task:0

RunGraph(graph,feed,fetches)

RecvTensor(rendezvous_key)

CPU

GPU

GPU

GPU

GraphMgr::ExecuteAsync

in [graph_mgr.cc#L283](#)

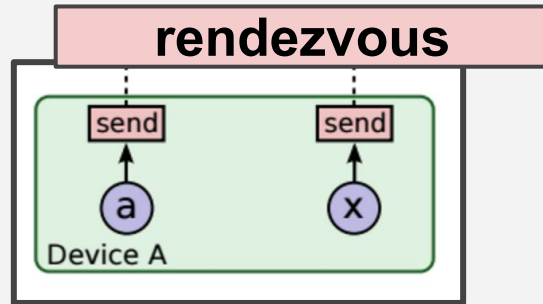
ExecutorState::RunAsync

in [executor.cc#L867](#)

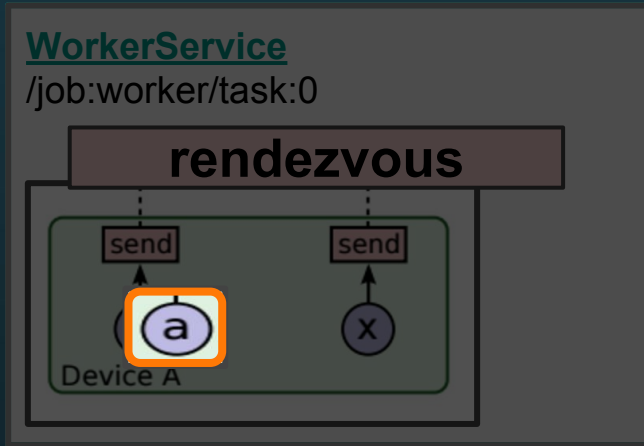
Executing: OpKernels

WorkerService

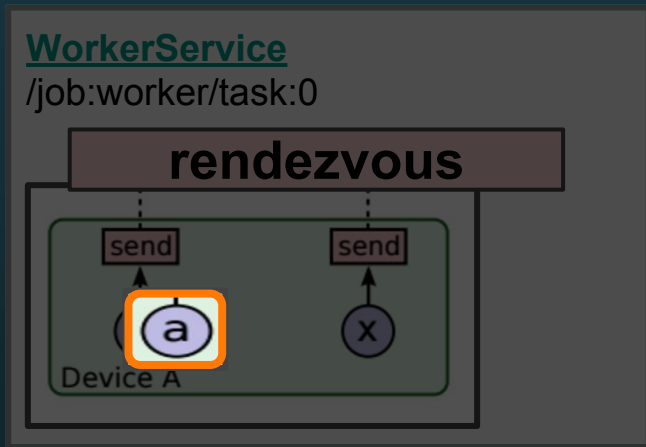
/job:worker/task:0



Executing: OpKernels

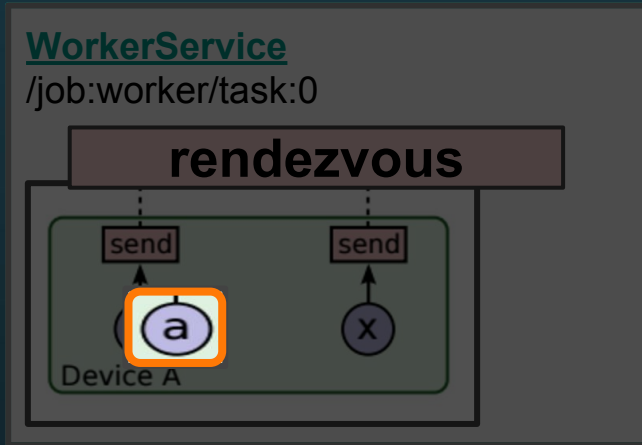


Executing: OpKernels



```
REGISTER_OP("MatMul")
  .Input("a: T")
  .Input("b: T")
  .Output("product: T")
  .Attr("transpose_a: bool = false")
  .Attr("transpose_b: bool = false")
  .Attr("T: {float, double, int32, complex64}")
  .Doc(R"doc(
Multiply the matrix "a" by the matrix "b".
```

Executing: OpKernels



Conv2D OpDef in [nn_ops.cc#L221](#)

```
REGISTER_OP("Conv2D")
  .Input("input: T")
  .Input("filter: T")
  .Output("output: T")
  .Attr("T: {float, double}")
  .Attr("strides: list(int)")
  .Attr("use_cudnn_on_gpu: bool = true")
  .Attr(GetPaddingAttrString())
```

Executing: OpKernels

Conditional build for OpKernels

```
#if GOOGLE_CUDA

// Registration of the GPU implementations.
REGISTER_KERNEL_BUILDER(
    Name("Conv2D").Device(DEVICE_GPU).TypeConstraint<float>("T"),
    Conv2DOp<GPUDevice, float>);

#endif // GOOGLE_CUDA
```

Executing: OpKernels

Conditional build for OpKernels

```
#if GOOGLE_CUDA

// Registration of the GPU implementations.
REGISTER_KERNEL_BUILDER(
    Name("Conv2D").Device(DEVICE_GPU).TypeConstraint<float>("T"),
    Conv2DOp<GPUDevice, float>);

#endif // GOOGLE_CUDA
```

CPU in [conv_ops.cc#L91](#)

GPU in [conv_ops.cc#L263](#)

Executing: OpKernels

OpKernels are **specialized** by device
adapted from [matmul_op.cc#L116](#)

```
template <typename Device, typename T, bool USE_CUBLAS>
class MatMulOp : public OpKernel {
public:
    explicit MatMulOp(OpKernelConstruction* ctx) : OpKernel(ctx) {
        OP_REQUIRES_OK(ctx, ctx->GetAttr("transpose_a", &transpose_a_));
        OP_REQUIRES_OK(ctx, ctx->GetAttr("transpose_b", &transpose_b_));
    }

    void Compute(OpKernelContext* ctx) override {
        const Tensor& a = ctx->input(0);
        const Tensor& b = ctx->input(1);

        //...

        LaunchMatMul<Device, T, USE_CUBLAS>::launch(ctx, this, a, b, dim_pair, out);
    }

private:
```

Executing: OpKernels

OpKernels are **specialized** by device
adapted from [matmul_op.cc#L116](#)

```
template <typename Device, typename T, bool USE_CUBLAS>
class MatMulOp : public OpKernel {
public:
    explicit MatMulOp(OpKernelConstruction* ctx) : OpKernel(ctx) {
        OP_REQUIRES_OK(ctx, ctx->GetAttr("transpose_a", &transpose_a_));
        OP_REQUIRES_OK(ctx, ctx->GetAttr("transpose_b", &transpose_b_));
    }

    void Compute(OpKernelContext* ctx) override {
        const Tensor& a = ctx->input(0);
        const Tensor& b = ctx->input(1);

        //...

        LaunchMatMul<Device, T, USE_CUBLAS>::launch(ctx, this, a, b, dim_pair, out);
    }

private:
```

Executing: OpKernels

OpKernels call into **Stream** functions

adapted from [matmul_op.cc#L71](#)

```
struct LaunchMatMul<GPUDevice, T, true /* USE_CUBLAS */> {
    static void launch(..., const Tensor& a, const Tensor& b, ..., Tensor* out) {
        const uint64 m = a.dim_size(1 - dim_pair[0].first);
        const uint64 k = a.dim_size(dim_pair[0].first);
        const uint64 n = b.dim_size(1 - dim_pair[0].second);
        // ...

        // Get a Stream for this GPUDevice
        auto* stream = ctx->op_device_context<GPUDeviceContext>()->stream();
        // ...

        // Launch the BLAS gemm kernel on the GPU stream
        bool blas_launch_status = stream->ThenBlasGemm(blas_transpose_b, blas_transpose_a,
                                                         n, m, k, 1.0f, b_ptr,
                                                         transpose_b ? k : n, a_ptr,
                                                         transpose_a ? m : k, 0.0f, &c_ptr,
                                                         n).ok();

        // ... return
    }
}
```

Executing: OpKernels

OpKernels call into **Stream** functions
adapted from [matmul_op.cc#L71](#)

```
struct LaunchMatMul<GPUDevice, T, true /* USE_CUBLAS */> {  
    static void launch(..., const Tensor& a, const Tensor& b, ..., Tensor* out) {  
        const uint64 m = a.dim_size(1 - dim_pair[0].first);  
        const uint64 k = a.dim_size(dim_pair[0].first);  
        const uint64 n = b.dim_size(1 - dim_pair[0].second);  
        // ...  
  
        // Get a Stream for this GPUDevice  
        auto* stream = ctx->op_device_context<GPUDeviceContext>()->stream();  
        // ...  
  
        // Launch the BLAS gemm kernel on the GPU stream  
        bool blas_launch_status = stream->ThenBlasGemm(blas_transpose_b, blas_transpose_a,  
                                                         n, m, k, 1.0f, b_ptr,  
                                                         transpose_b ? k : n, a_ptr,  
                                                         transpose_a ? m : k, 0.0f, &c_ptr,  
                                                         n).ok();  
  
        // ... return  
    }  
}
```


Executing: Stream functions

OpKernels call into **Stream** functions

in [conv_ops.cc#L292](#)

```
bool blas_launch_status =
  stream
    ->ThenBlasGemm(no_transpose, no_transpose, n, m, k, 1.0f, b_
                  n, a_ptr, k, 0.0f, &c_ptr, n)
    .ok();
```

Executing: Stream functions

OpKernels call into **Stream** functions

in [conv_ops.cc#L292](#)

```
bool blas_launch_status =
    stream
        ->ThenBlasGemm(no_transpose, no_transpose, n, m, k, 1.0f, b_
                       n, a_ptr, k, 0.0f, &c_ptr, n)
        .ok();
```

in [conv_ops.cc#L417](#)

```
CudnnScratchAllocator scratch_allocator(ConvolveScratchSize, ctx);
bool cudnn_launch_status =
    stream
        ->ThenConvolveWithScratch(input_desc, input_ptr, filter_desc,
                                   filter_ptr, conv_desc, output_desc,
                                   &output_ptr, &scratch_allocator)
        .ok();
```

Executing: Stream functions

Platforms provide GPU-specific implementations

cuBLAS

BlasSupport in [stream_executor/blas.h#L88](#)

DoBlasInternal in [cuda_blas.cc#L429](#)

Executing: Stream functions

Platforms provide GPU-specific implementations

cuBLAS

BlasSupport in [stream_executor/blas.h#L88](#)

DoBlasInternal in [cuda_blas.cc#L429](#)

cuDNN

DnnSupport in [stream_executor/dnn.h#L544](#)

DoConvolve in [cuda_dnn.cc#L629](#)

```
status = dynload::cudnnConvolutionForward(  
    parent_, ToHandle(dnn_handle_),  
    /*alpha=*/&alpha, /*srcDesc=*/input_4d.handle(),  
    /*srcData=*/input_data.opaque(), /*filterDesc=*/filter.handle(),  
    /*filterData=*/filter_data.opaque(), /*convDesc=*/conv.handle(),  
    /*algo=*/algo, /*workSpace=*/scratch.opaque(),  
    /*workSpaceSizeInBytes=*/scratch.size(), /*beta=*/&beta,  
    /*destDesc=*/output_4d.handle(), /*destData=*/output_data->opaque());
```

Session Interface

- `Extend`: add nodes to computation graph
- `Run`: execute an arbitrary subgraph
 - optionally feeding in Tensor inputs and retrieving Tensor output

Typically, setup a graph with one or a few `Extend` calls and then `Run` it thousands or millions of times



Single device performance important, but

....

biggest performance improvements come
from large-scale distributed systems with
model and data parallelism



Experiment Turnaround Time and Research Productivity

- **Minutes, Hours:**
 - **Interactive research! Instant gratification!**
- **1-4 days**
 - Tolerable
 - Interactivity replaced by running many experiments in parallel
- **1-4 weeks**
 - High value experiments only
 - Progress stalls
- **>1 month**
 - Don't even try



Transition

- How do you do this at scale?
- How does TensorFlow make distributed training easy?



Model Parallelism

- Best way to decrease training time: **decrease the step time**
- Many models have lots of inherent parallelism
- Problem is distributing work so communication doesn't kill you
 - local connectivity (as found in CNNs)
 - towers with little or no connectivity between towers (e.g. AlexNet)
 - specialized parts of model active only for some examples



Exploiting Model Parallelism

On a single core: Instruction parallelism (SIMD). Pretty much free.

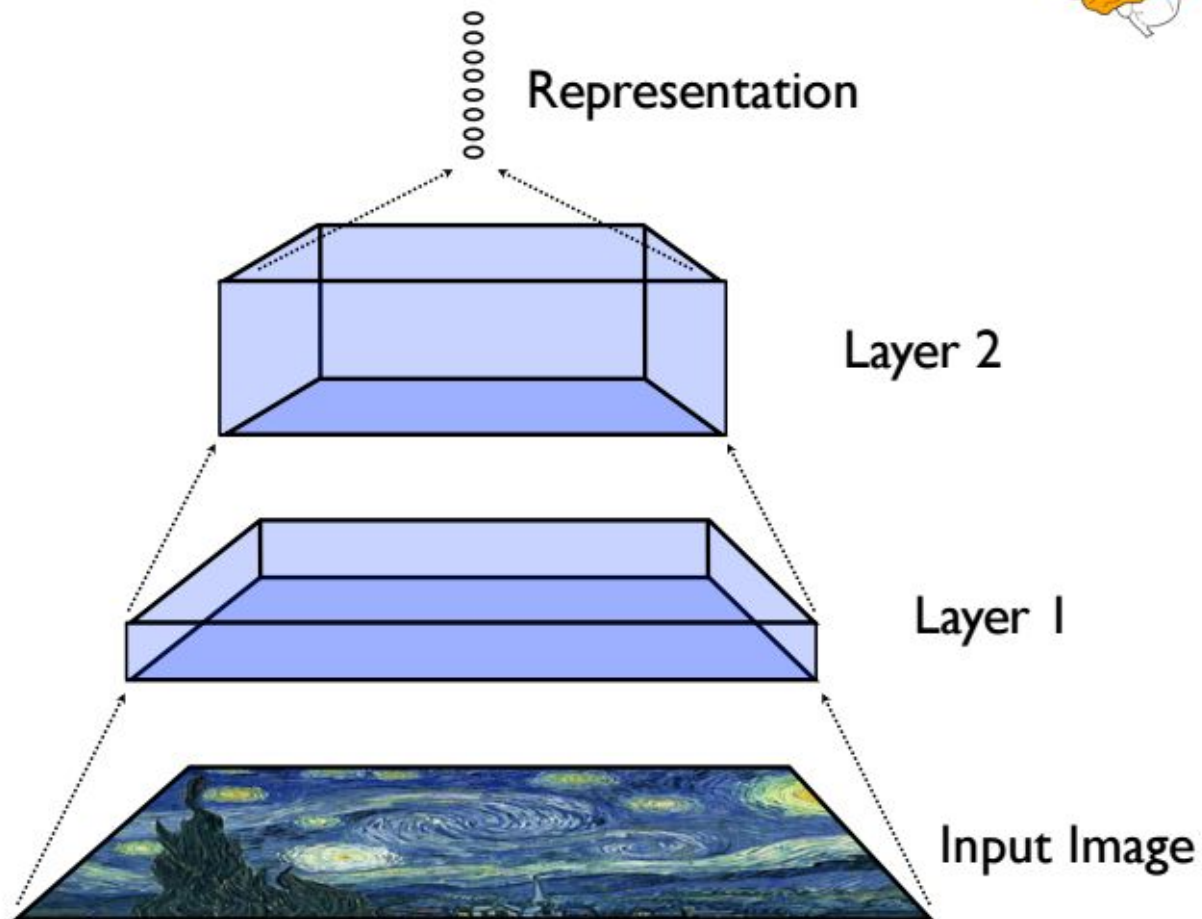
Across cores: thread parallelism. Almost free, unless across sockets, in which case inter-socket bandwidth matters (QPI on Intel).

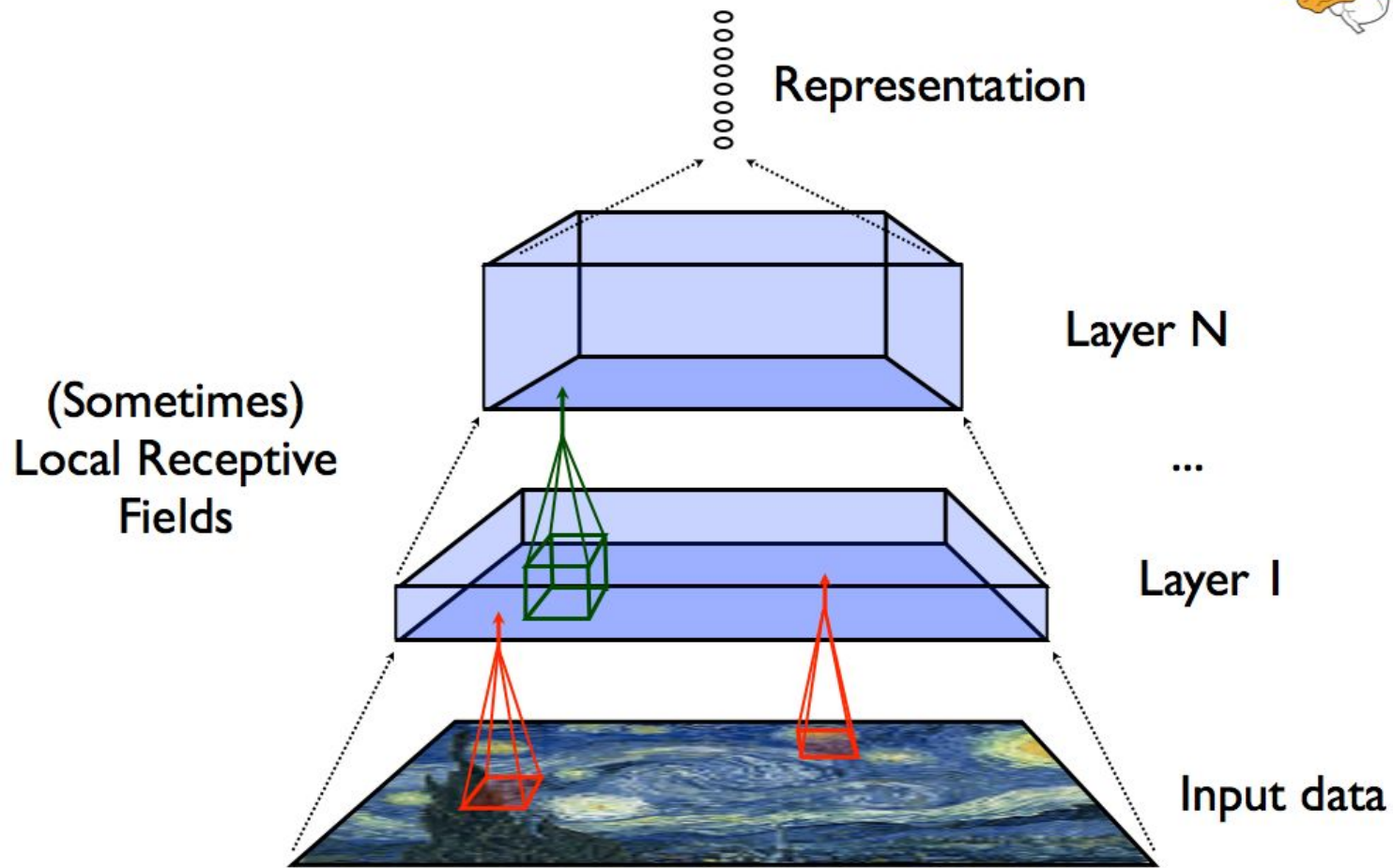
Across devices: for GPUs, often limited by PCIe bandwidth.

Across machines: limited by network bandwidth / latency

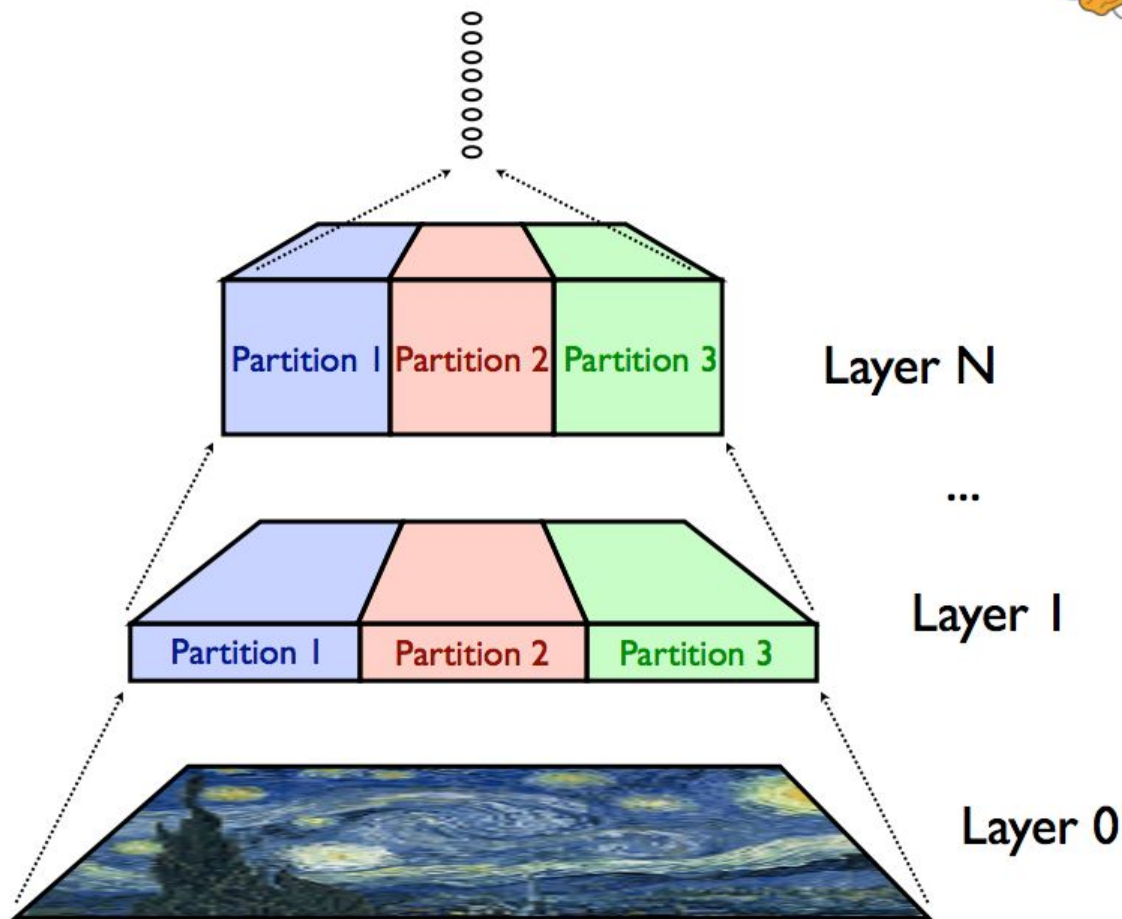


Model Parallelism

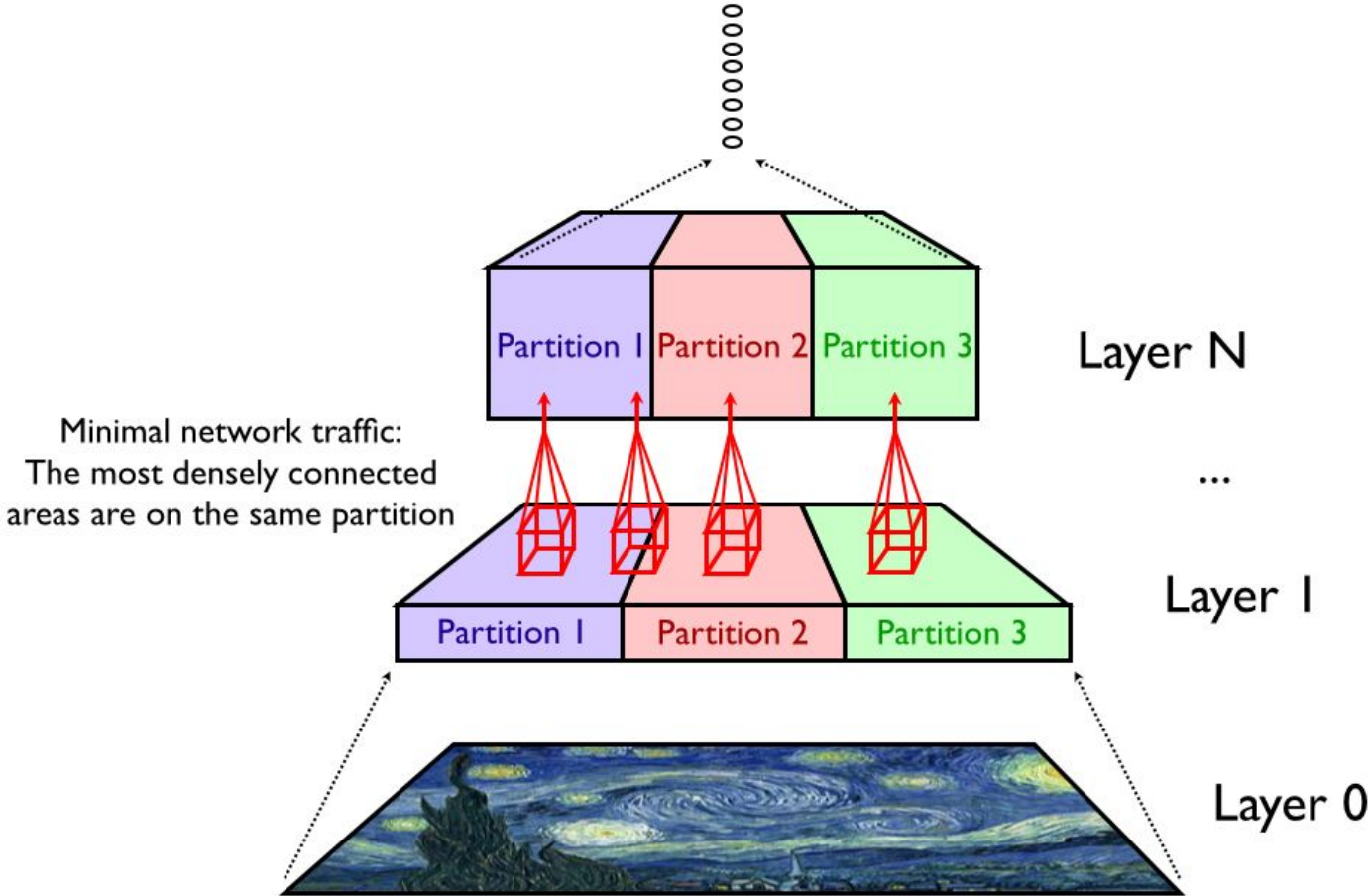




Model Parallelism: Partition model across machines



Model Parallelism: Partition model across machines



Using TensorFlow for Parallelism

Easy to express both model parallelism as well as data parallelism

- Very minimal changes to single device model code



Devices and Graph Placement

- Given a graph and set of devices, TensorFlow implementation must decide which device executes each node



Full and Partial Device Constraints (Hints)

Devices are named hierarchically:

```
/job:localhost/device:cpu:0
```

```
/job:worker/task:17/device:gpu:3
```

```
/job:parameters/task:4/device:cpu:0
```

Client can specify full or partial constraints for nodes in graph:

```
“Place this node on /job:localhost/device:gpu:2”
```

```
“Place this node on /device:gpu:*”
```



Placement Algorithm

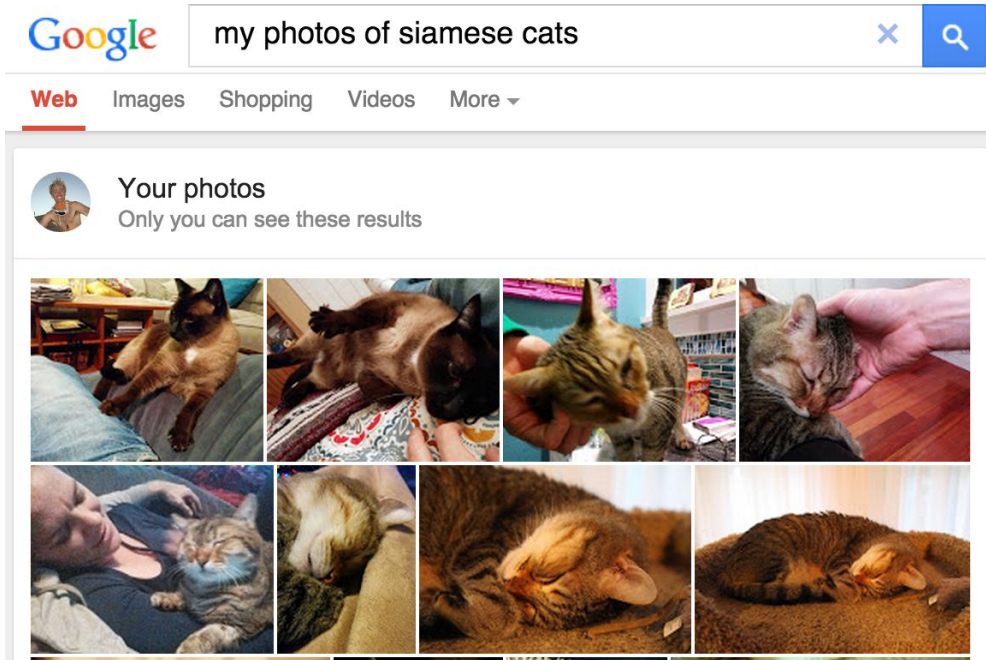
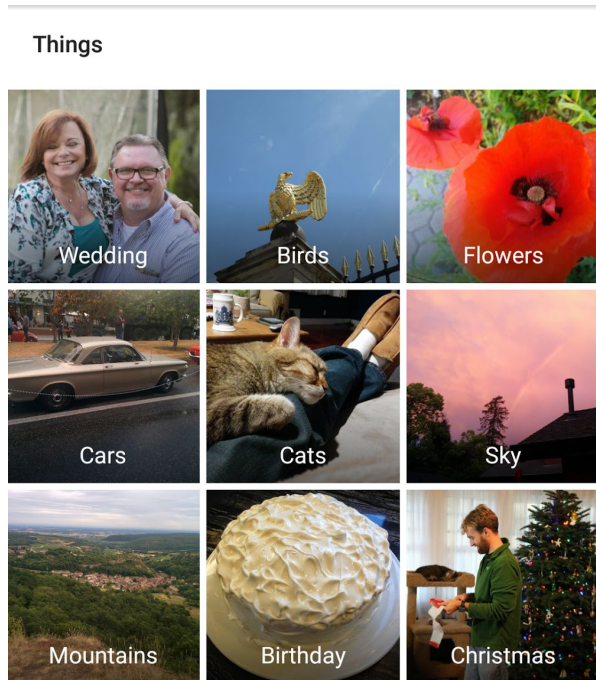
Given hints, plus a cost model (node execution time estimates and Tensor size estimates), make placement decisions

- Current relatively simple greedy algorithm
- Active area of work

Show CIFAR10 placement TensorBoard.



Google Photos Search



Reuse same model for completely different problems

**Same basic model structure
trained on different data,
useful in completely different contexts**

Example: given image → predict interesting pixels



ASIAWIDE TRAVEL 環宇國際旅遊

(61) (02) 9745 3355 1st Floor, 240 BURWOOD RD



Maria's Bakery Inn 超羣餅屋

Maria's Bakery Inn 超羣餅屋



CIANO MOTOR ENGINEERS

MECHANICAL REPAIRS TO ALL MAKES AND MODELS

Specialising In **BMW, MINI & TOYOTA**

8 REGATTA ROAD FIVE DOCK 9745 3173

88

- LATEST DIAGNOSTIC EQUIPMENT • REGO INSPECTIONS •
- NEW CAR/DOORER SERVICING • BRAKES • CLUTCHES •
- STEERING • SUSPENSION • TYRES • WHEEL ALIGNMENTS •
- RADIATORS • COILS • 2-1 COUPLERS • 5-1 TUNING •
- FUEL INJECTION SERVICING • BATTERIES • AUTO ELECTRICAL •

Factory Trained Technicians



58



1234 Bryant St, Palo Alto, CA 94301, USA



Analysis complete. Your roof has:



1,658 hours of usable sunlight per year

Based on day-to-day analysis of weather patterns



708 sq feet available for solar panels

Based on 3D modeling of your roof and nearby trees

If your electric bill is at least \$175/month, leasing solar panels could reduce it.

[FINE-TUNE ESTIMATE](#)

[SEE SOLAR PROVIDERS](#)

Wrong roof? Drag the marker to the right one.

Shade

Sun

A fundus photograph of a human retina, showing the optic disc on the right and a network of retinal vessels. Several areas are highlighted with semi-transparent segmentation masks in shades of purple and blue, indicating regions of interest or pathology. The background is a warm, orange-brown color.

MEDICAL IMAGING

Very good results using similar model for detecting diabetic retinopathy in retinal images

“Seeing” Go

Google’s AI just cracked the game that supposedly no computer could beat

By Mike Murphy | January 27, 2016



The screenshot shows the BBC News website interface. At the top, there are navigation tabs for News, Sport, Weather, iPlayer, TV, Radio, and More. A search bar is located on the right. Below the navigation is a red banner with the word 'NEWS' and a 'Find local news' button. The main content area is under the 'Technology' category. The article title is 'Google achieves AI 'breakthrough' at Go'. The sub-headline reads: 'An artificial intelligence program developed by Google beats Europe's top player at the ancient Chinese game of Go, about a decade earlier than expected.' The article is dated '27 January 2016' and is categorized under 'Technology'. There are two links: 'How did they do it?' and 'What is the game Go?'. A small image shows a hand placing a black Go stone on a board.

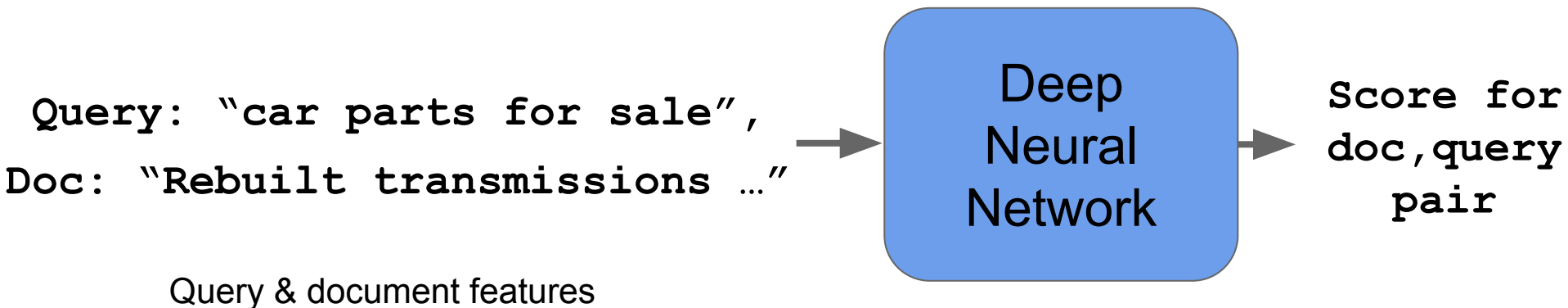


(Kiyoshi Ota)

... slowly started to encroach on activities we previously thought only a brilliantly sophisticated human brain could handle. In 1997, a supercomputer beat Grand Master Garry Kasparov at chess, and in 2011 IBM's Watson beat former human winners at the quiz game *Jeopardy*. But the ancient board game Go has long been one of the major goals of artificial intelligence research. It's understood to be one of the most difficult games for computers to handle due to the sheer number of possible moves a player can make at any given point. Until now, that is.



RankBrain in Google Search Ranking



Launched in 2015

Third most important search ranking signal (of 100s)

Bloomberg, Oct 2015: "Google Turning Its Lucrative Web Search Over to AI Machines"

Example: LSTM [Hochreiter et al, 1997][Gers et al, 1999]



$$i_t = W_{ix}x_t + W_{ih}h_{t-1} + b_i$$

$$j_t = W_{jx}x_t + W_{jh}h_{t-1} + b_j$$

$$f_t = W_{fx}x_t + W_{fh}h_{t-1} + b_f$$

$$o_t = W_{ox}x_t + W_{oh}h_{t-1} + b_o$$

$$c_t = \sigma(f_t) \odot c_{t-1} + \sigma(i_t) \odot \tanh(j_t)$$

$$h_t = \sigma(o_t) \odot \tanh(c_t)$$

```
def __call__(self, inputs, state, scope=None):
    """Long short-term memory cell (LSTM)."""
    with vs.variable_scope(scope or type(self).__name__): # "BasicLSTMCell"
        # Parameters of gates are concatenated into one multiply for efficiency.
        c, h = array_ops.split(1, 2, state)
        concat = linear([inputs, h], 4 * self._num_units, True)

        # i = input_gate, j = new_input, f = forget_gate, o = output_gate
        i, j, f, o = array_ops.split(1, 4, concat)

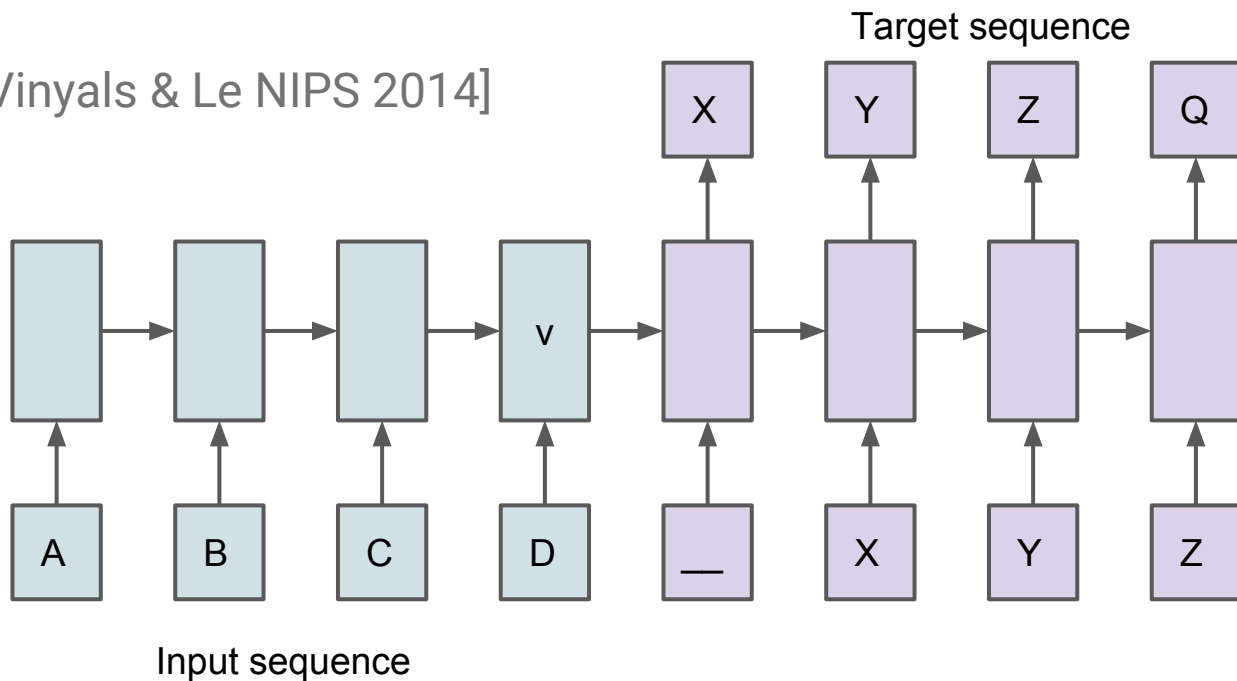
        new_c = c * sigmoid(f + self._forget_bias) + sigmoid(i) * tanh(j)
        new_h = tanh(new_c) * sigmoid(o)

    return new_h, array_ops.concat(1, [new_c, new_h])
```

Enables
long term
dependencies
to flow

Sequence-to-Sequence Model

[Sutskever & Vinyals & Le NIPS 2014]



$$P(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

Sequence-to-Sequence

- Active area of research
- Many groups actively pursuing RNN/LSTM
 - Montreal
 - Stanford
 - U of Toronto
 - Berkeley
 - Google
 - ...
- Further Improvements
 - Attention
 - NTM / Memory Nets
 - ...

Sequence-to-Sequence

- **Translation**: [Kalchbrenner *et al.*, EMNLP 2013][Cho *et al.*, EMLP 2014][Sutskever & Vinyals & Le, NIPS 2014][Luong *et al.*, ACL 2015][Bahdanau *et al.*, ICLR 2015]
- **Image captions**: [Mao *et al.*, ICLR 2015][Vinyals *et al.*, CVPR 2015][Donahue *et al.*, CVPR 2015][Xu *et al.*, ICML 2015]
- **Speech**: [Chorowsky *et al.*, NIPS DL 2014][Chan *et al.*, arxiv 2015]
- **Language Understanding**: [Vinyals & Kaiser *et al.*, NIPS 2015][Kiros *et al.*, NIPS 2015]
- **Dialogue**: [Shang *et al.*, ACL 2015][Sordoni *et al.*, NAACL 2015][Vinyals & Le, ICML DL 2015]
- **Video Generation**: [Srivastava *et al.*, ICML 2015]
- **Algorithms**: [Zaremba & Sutskever, arxiv 2014][Vinyals & Fortunato & Jaitly, NIPS 2015][Kaiser & Sutskever, arxiv 2015][Zaremba *et al.*, arxiv 2015]

How to do Image Captions?

$P(\text{English} | \text{French})$

Image Captioning

[Vinyals *et al.*, CVPR 2015]

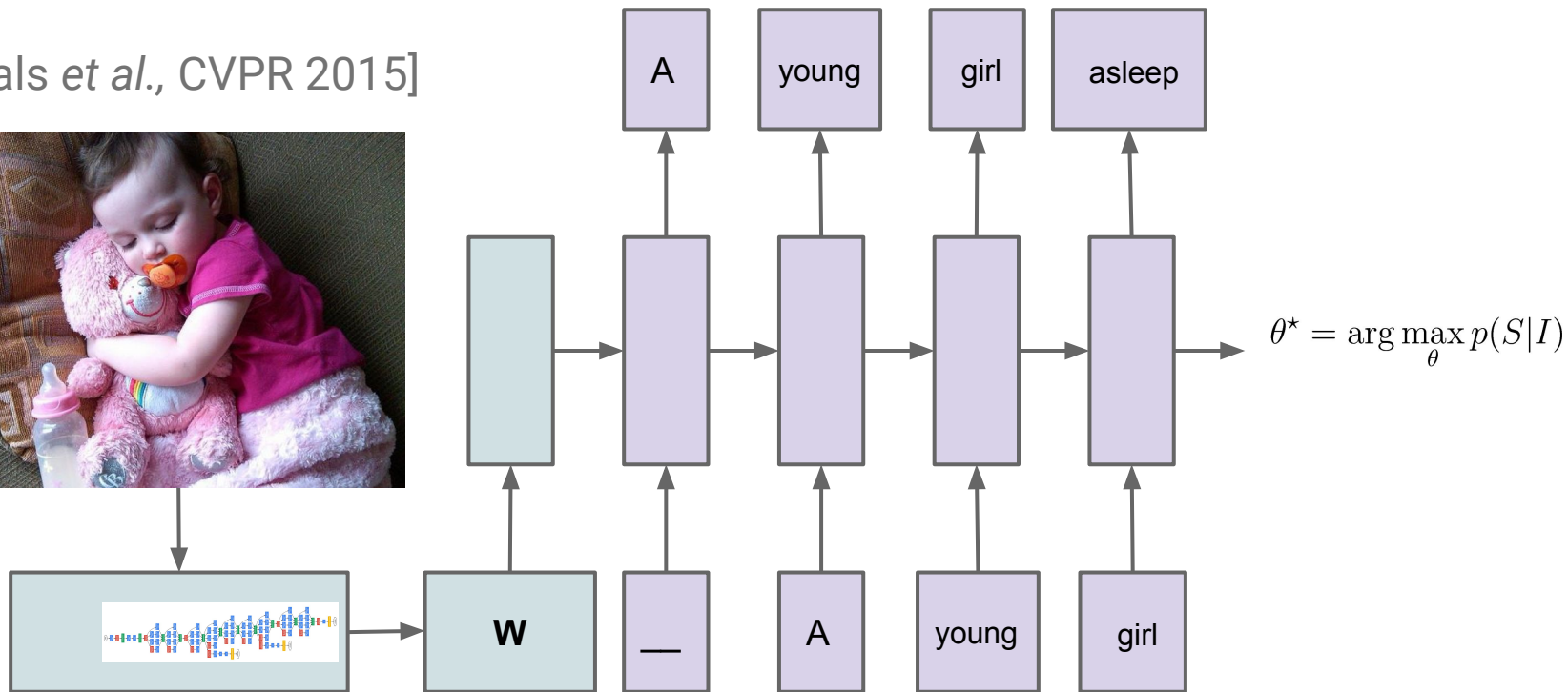
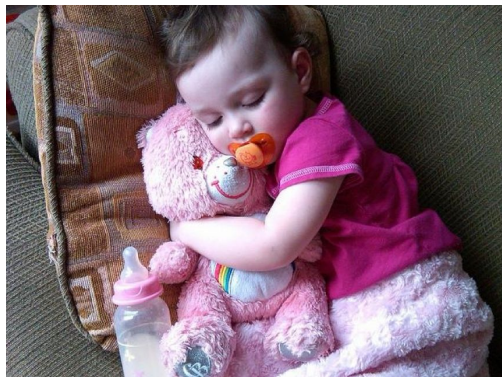
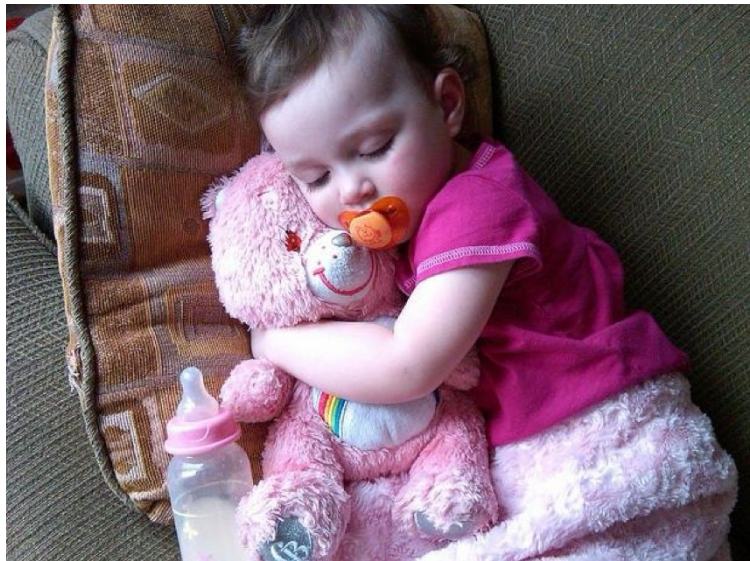


Image Captions Research



Human: A young girl asleep on the sofa cuddling a stuffed bear.

Model: A close up of a child holding a stuffed animal.

Model: A baby is asleep next to a teddy bear.



Human: A man outside cooking with a sub in his hand.

BestModel: A man is holding a sandwich in his hand.

InitialModel: A man cutting a cake with a knife.



Human: Someone is using a small grill to melt his sandwich.

BestModel: A person is cooking some food on a grill.

InitialModel: A pizza sitting on top of a white plate.



Human: A woman holding up a yellow banana to her face.

BestModel: A woman holding a banana up to her face.

InitialModel: A close up of a person eating a hot dog.



Human: A blue , yellow and red train travels across the tracks near a depot.

BestModel: A blue and yellow train traveling down train tracks.

InitialModel: A train that is sitting on the tracks.



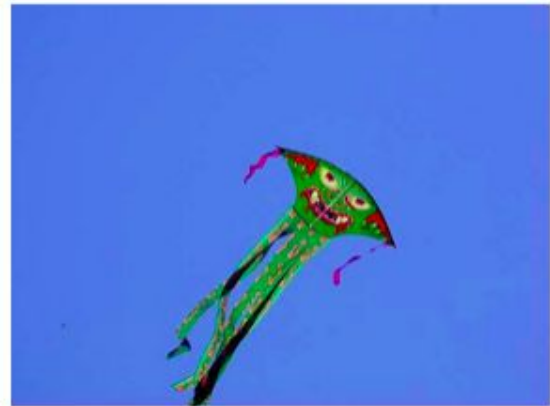
A man holding a tennis racquet on a tennis court.



Two pizzas sitting on top of a stove top oven



A group of young people playing a game of Frisbee

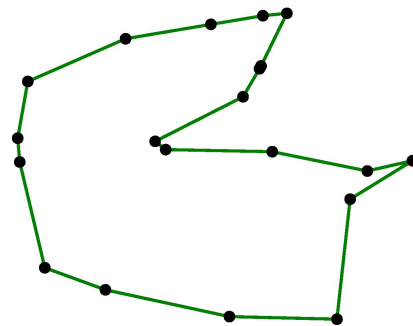
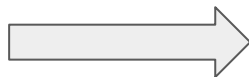
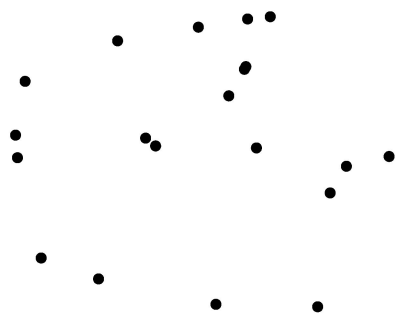


A man flying through the air while riding a snowboard

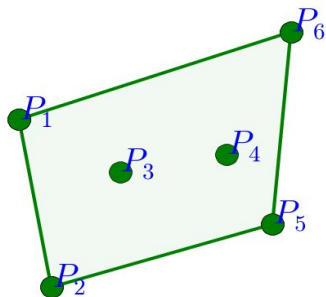


Pointer Networks

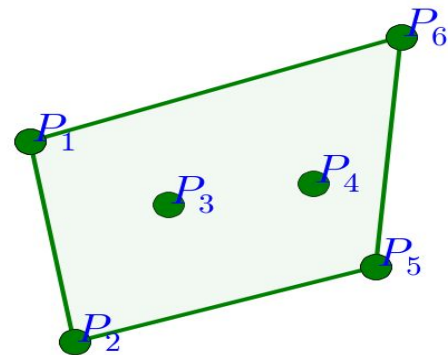
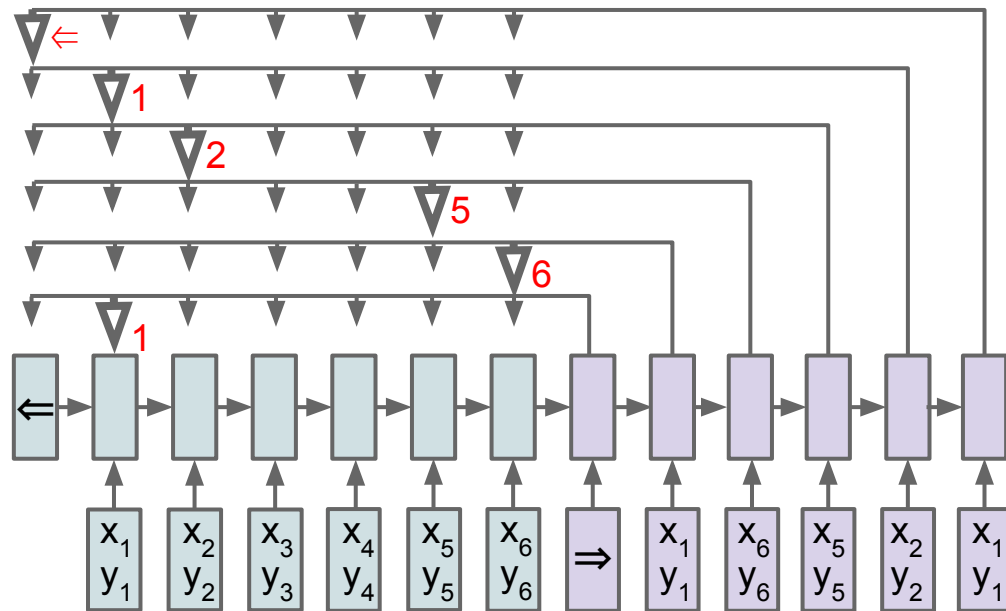
- Goal: Mappings where outputs are (sub)sets of inputs
- Travelling Salesman Problem



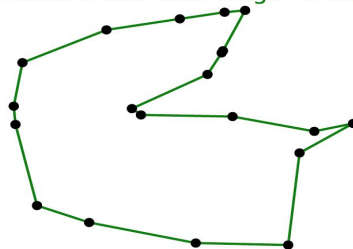
- Convex Hulls



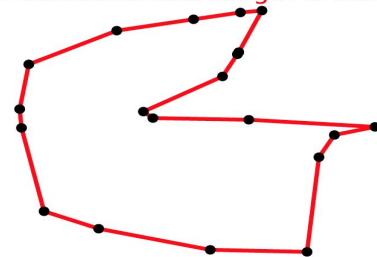
Pointer Networks



Ground Truth: tour length is 3.518



Predictions: tour length is 3.523



Neural Conversational Models

- Take movie subtitles (~900M words) or IT HelpDesk chats
- Predict the next dialog from history

i got to go .

no .

i get too emotional when i drink .

have another beer . i 've got to get up early .

no , you don 't . sit down .

i get too emotional when i drink .

will you have another beer ?

i 've got to go !

why ?

i got to get up early in the morning .

you 're drunk .

and emotional !

you got to go .

[Vinyals & Le ICML DL Workshop 2015]

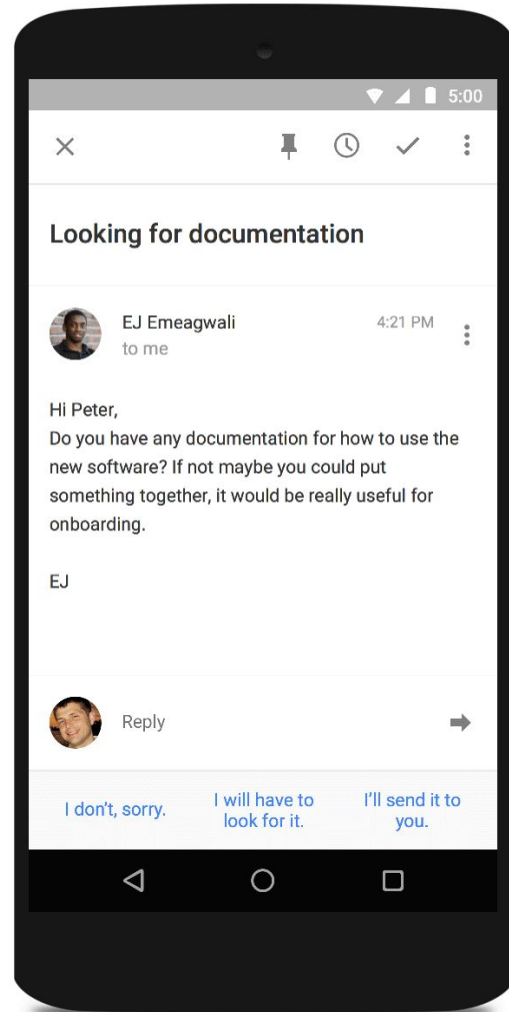


Smart Reply

April 1, 2009: April Fool's Day joke

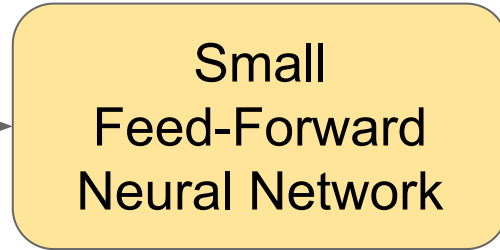
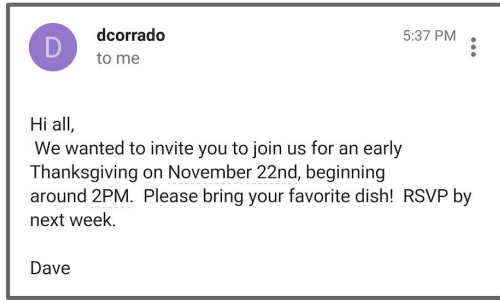
Nov 5, 2015: Launched Real Product

Feb 1, 2016: >10% of mobile Inbox replies



Smart Reply

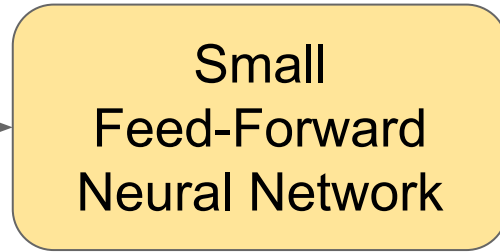
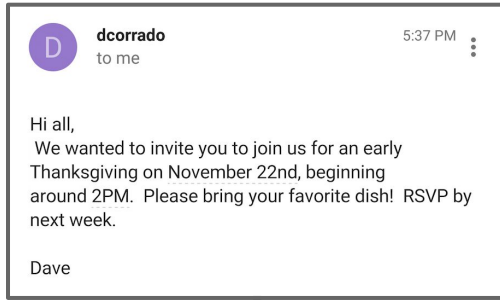
Incoming Email



Activate
Smart Reply?
yes/no

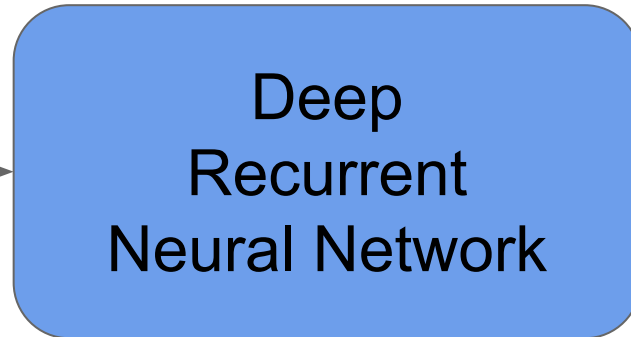
Smart Reply

Incoming Email

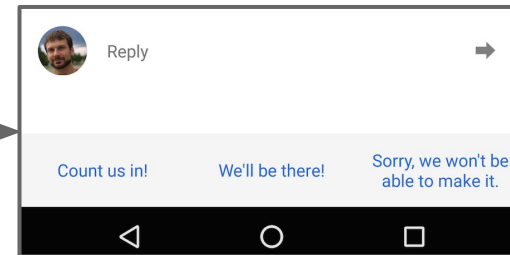


Activate Smart Reply?

yes/no



Generated Replies



Example: LSTM

```
for i in range(20):
```

```
    m, c = LSTMCell(x[i], mprev, cprev)
```

```
    mprev = m
```

```
    cprev = c
```



Example: Deep LSTM

for i in range(20):

for d in range(4): # d is depth

input = x[i] if d is 0 else m[d-1]

m[d], c[d] = LSTMCell(input, mprev[d], cprev[d])

mprev[d] = m[d]

cprev[d] = c[d]



Example: Deep LSTM

```
for i in range(20):  
    for d in range(4): # d is depth  
        input = x[i] if d is 0 else m[d-1]  
        m[d], c[d] = LSTMCell(input, mprev[d], cprev[d])  
        mprev[d] = m[d]  
        cprev[d] = c[d]
```



Example: Deep LSTM

```
for i in range(20):
```

```
    for d in range(4): # d is depth
```

```
        with tf.device("/gpu:%d" % d):
```

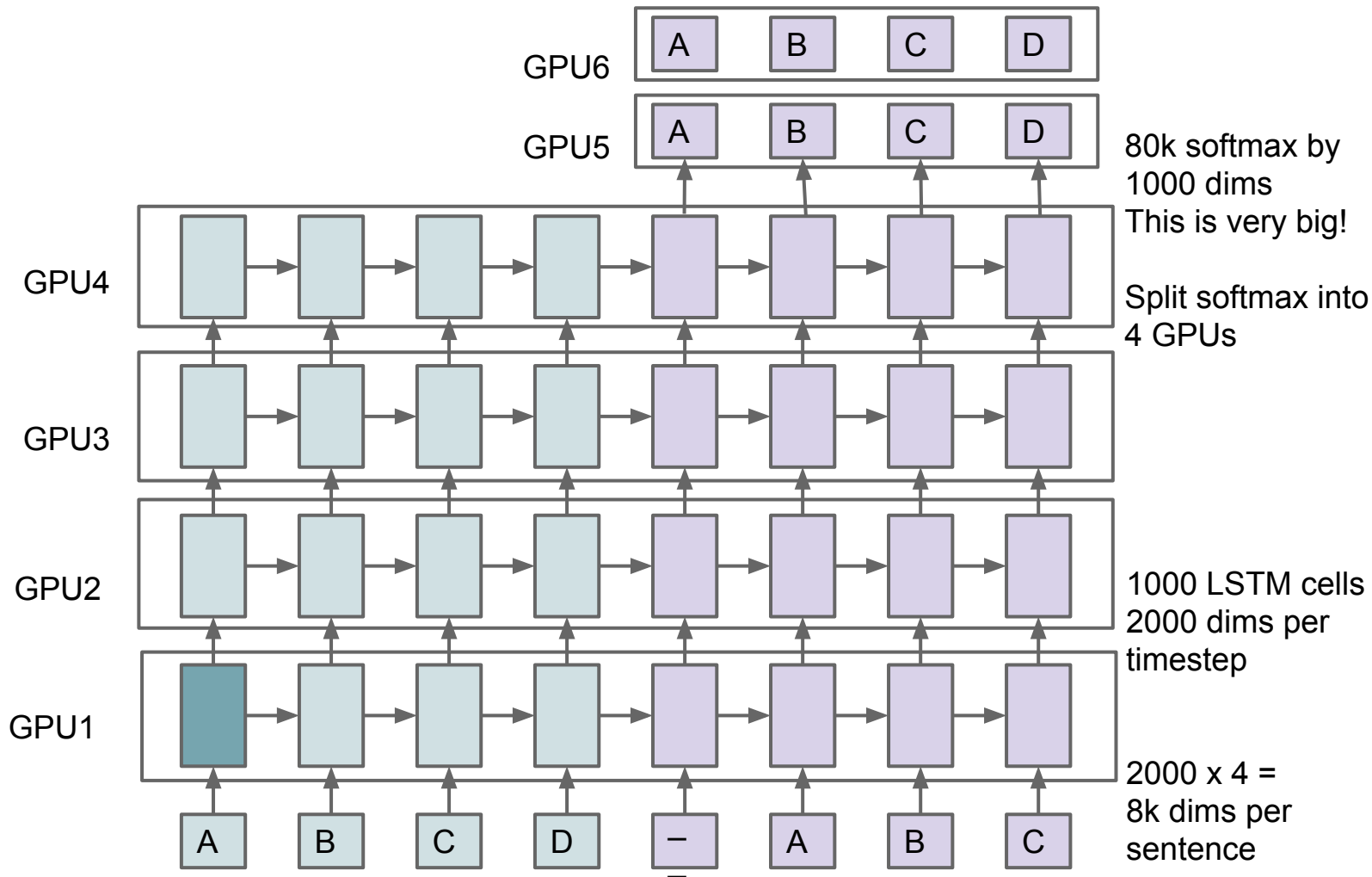
```
            input = x[i] if d is 0 else m[d-1]
```

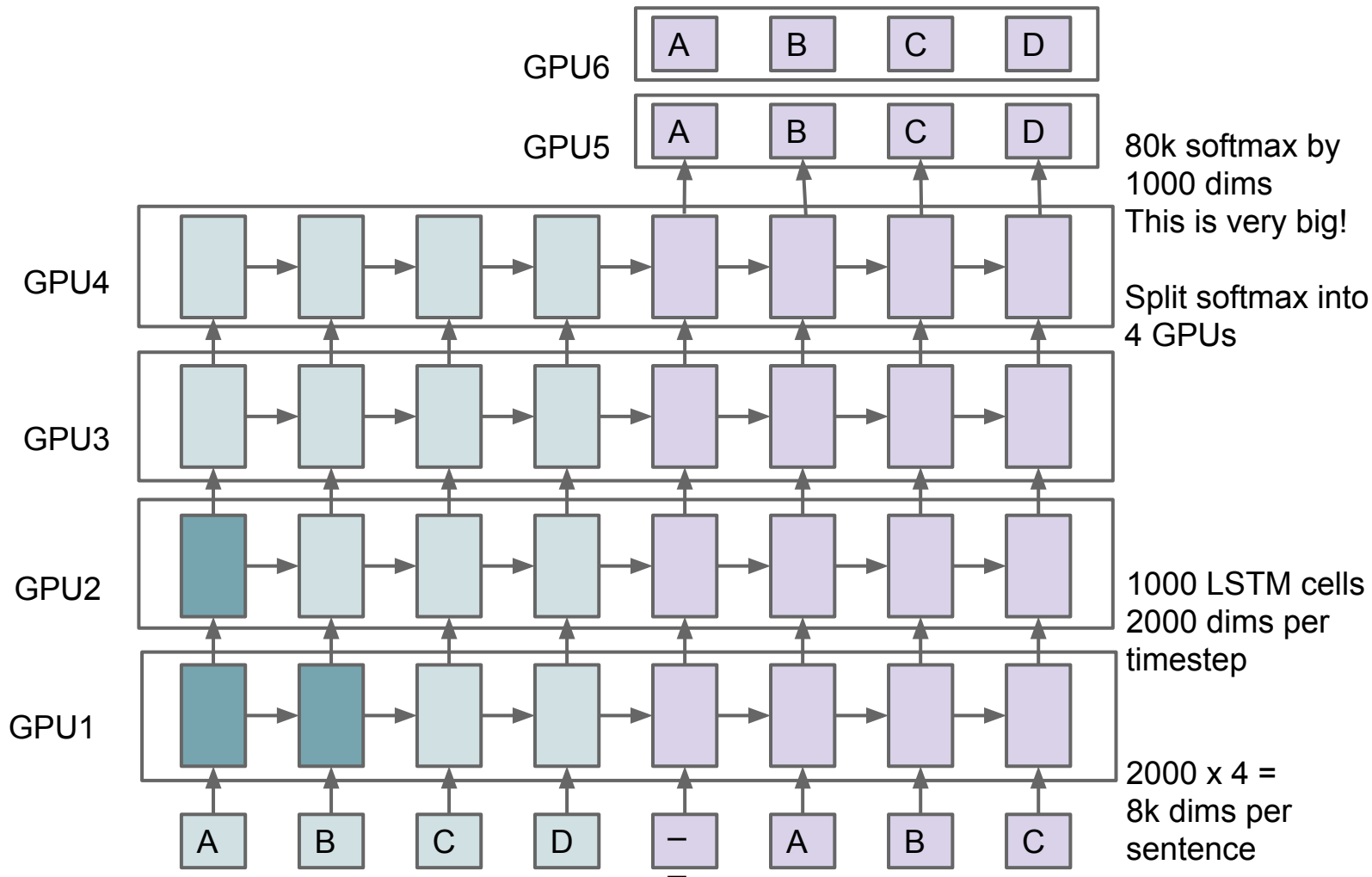
```
            m[d], c[d] = LSTMCell(input, mprev[d], cprev[d])
```

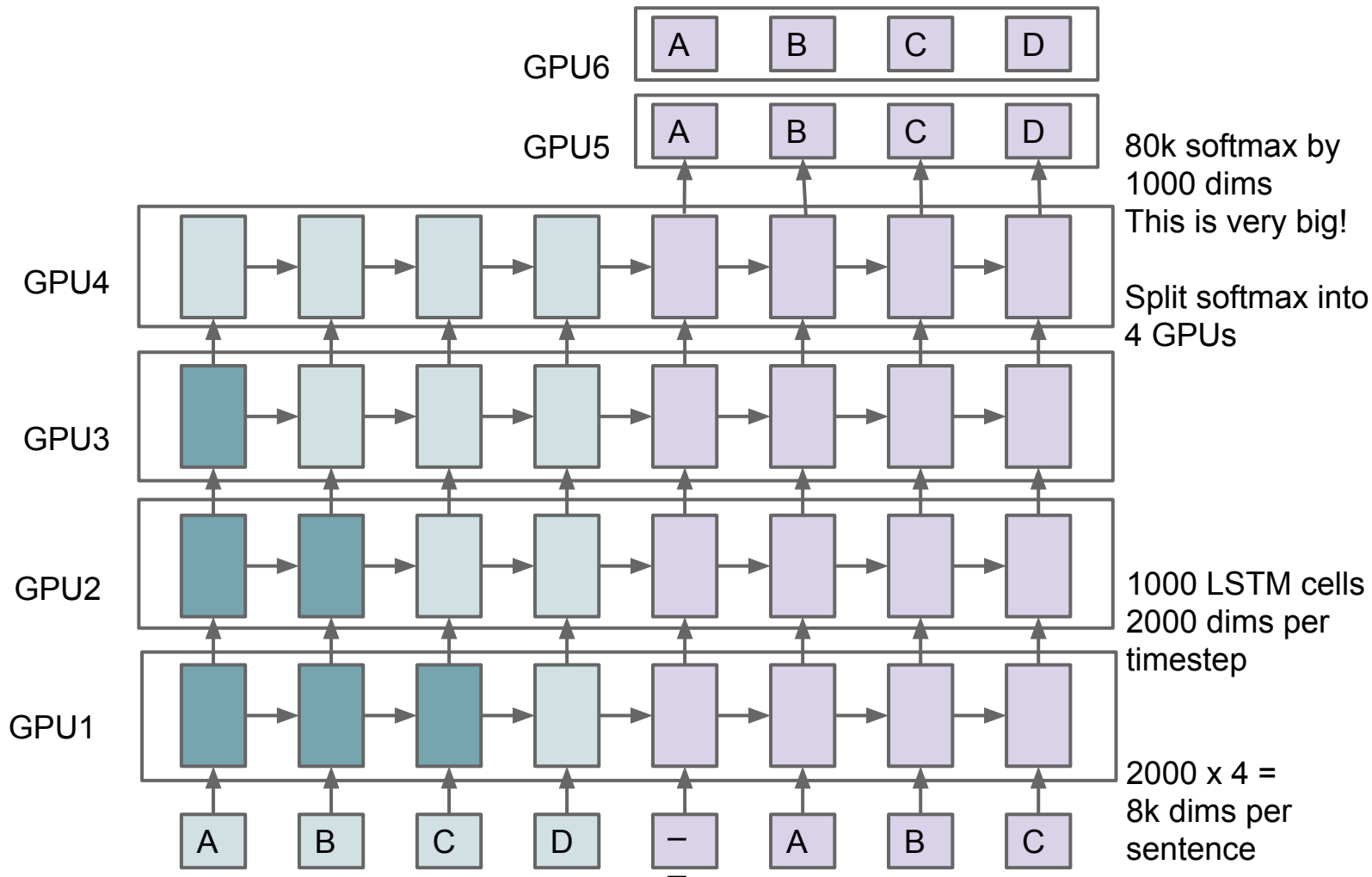
```
            mprev[d] = m[d]
```

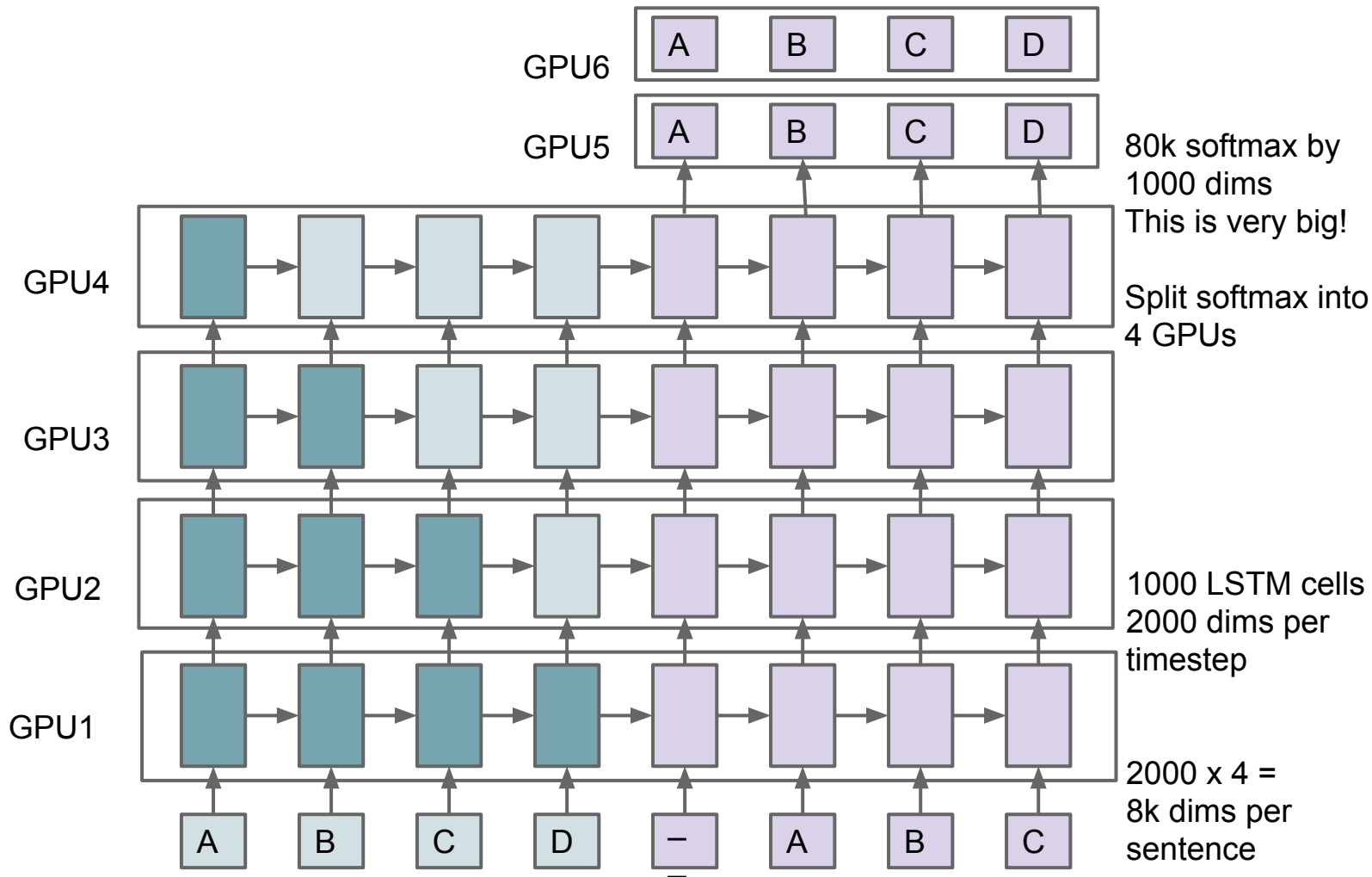
```
            cprev[d] = c[d]
```

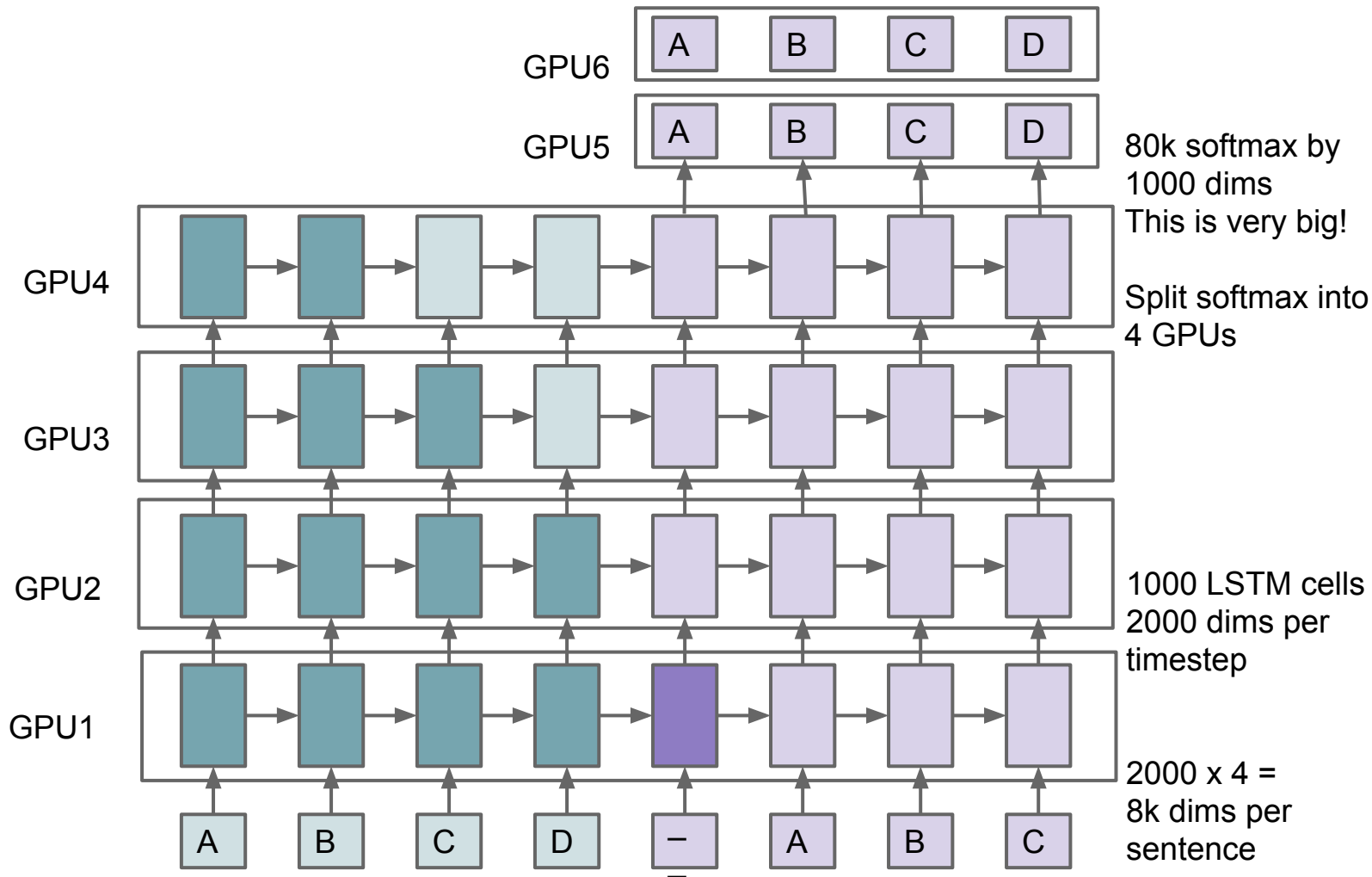


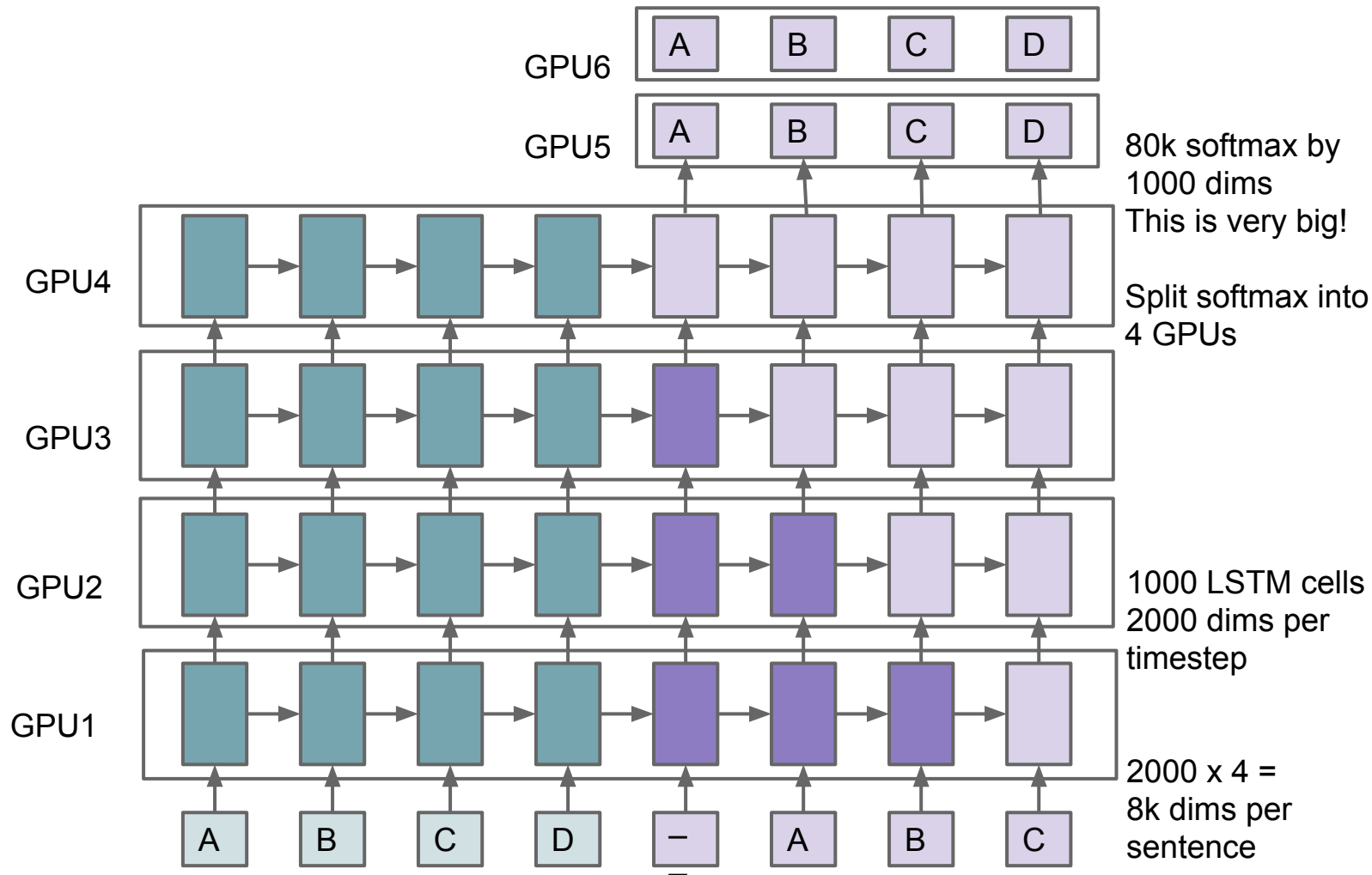


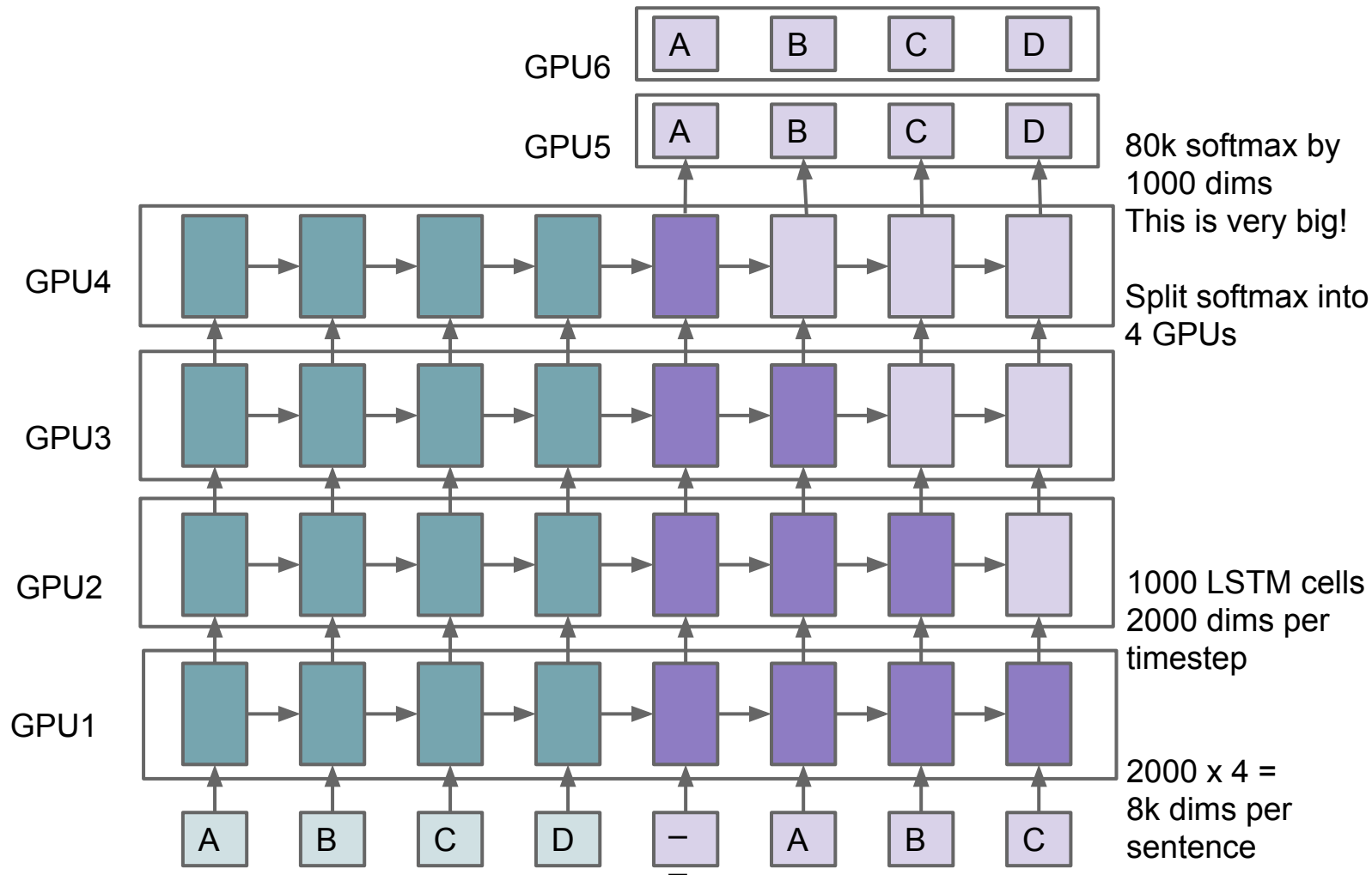


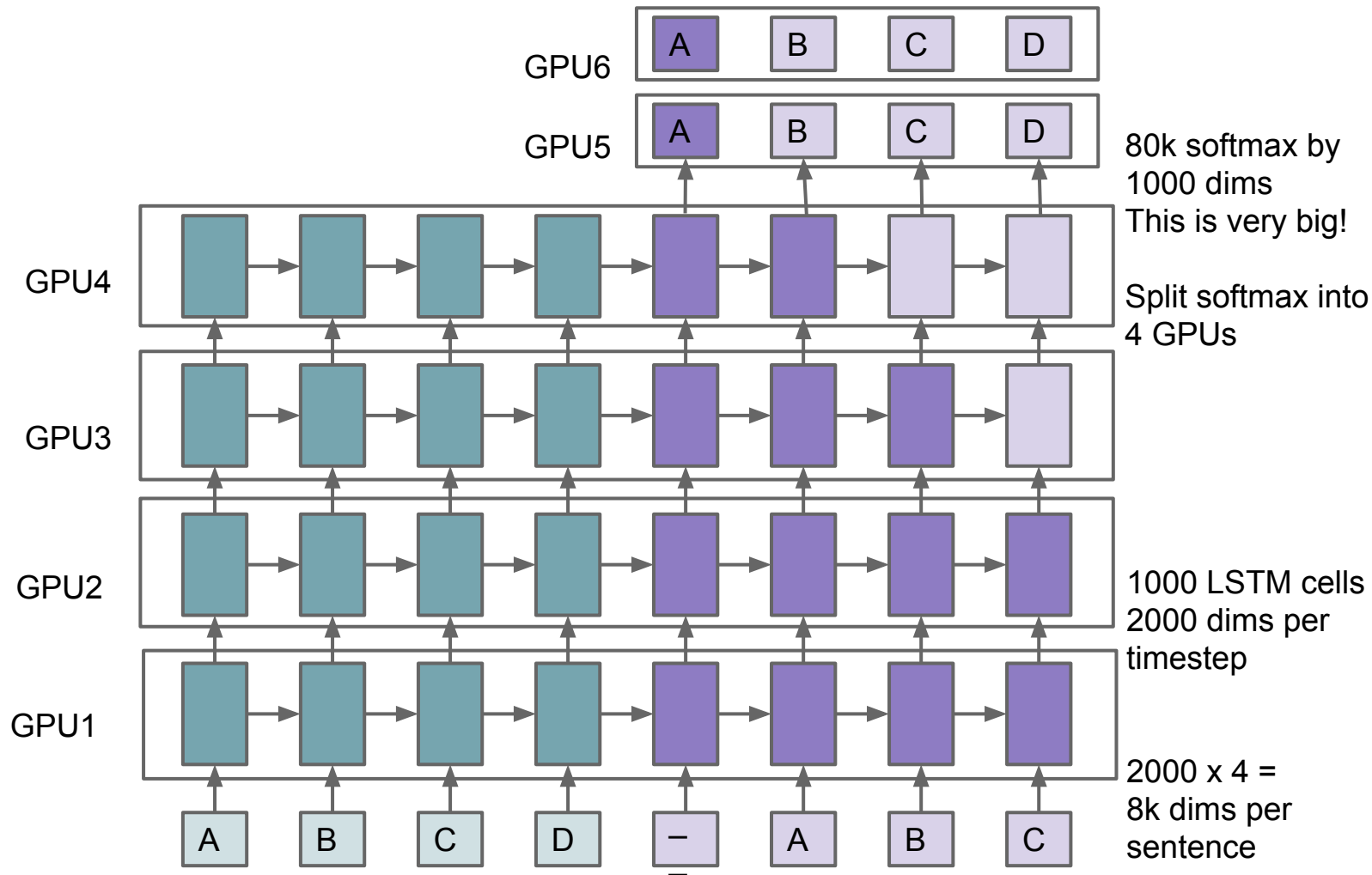


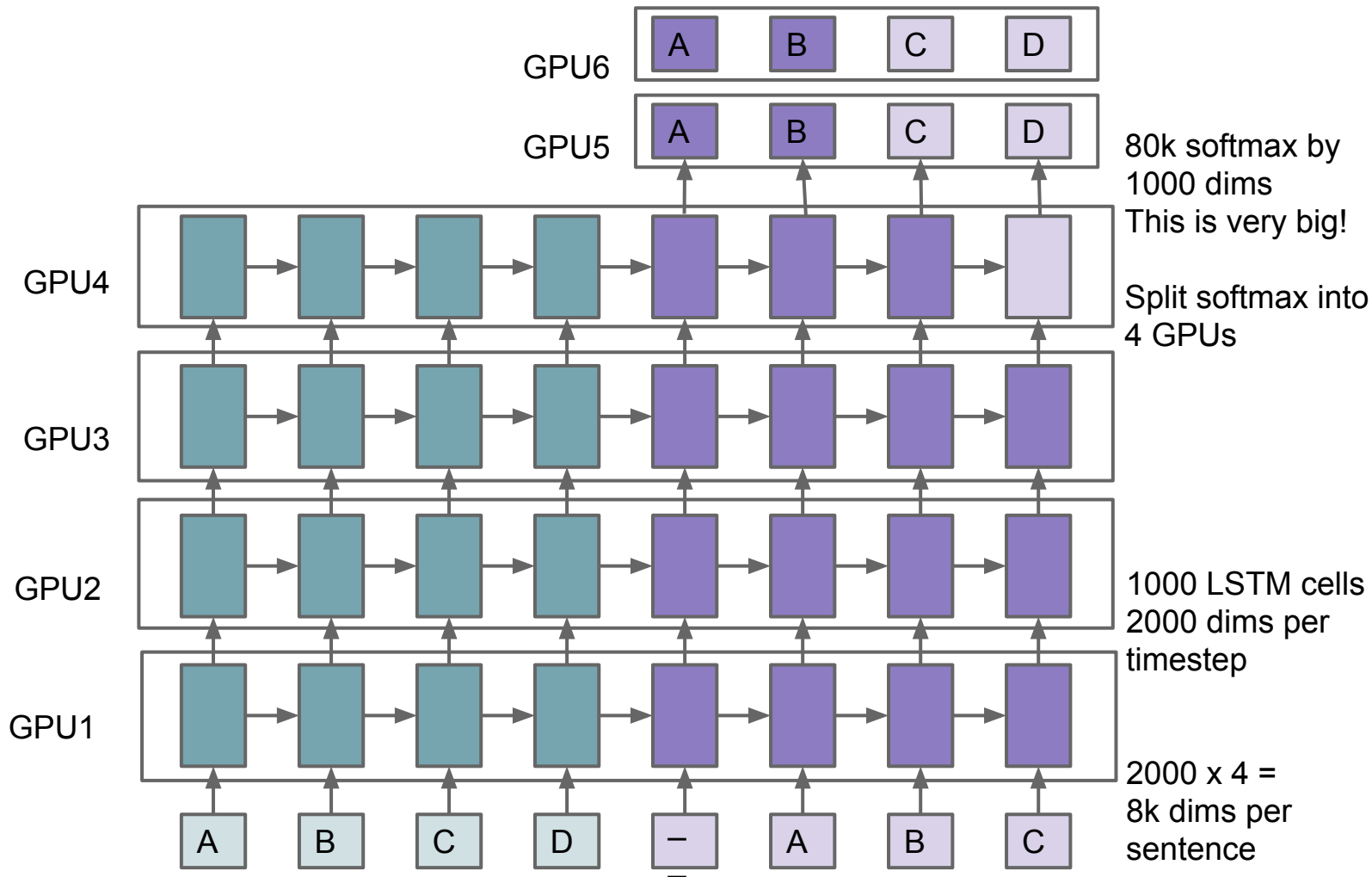


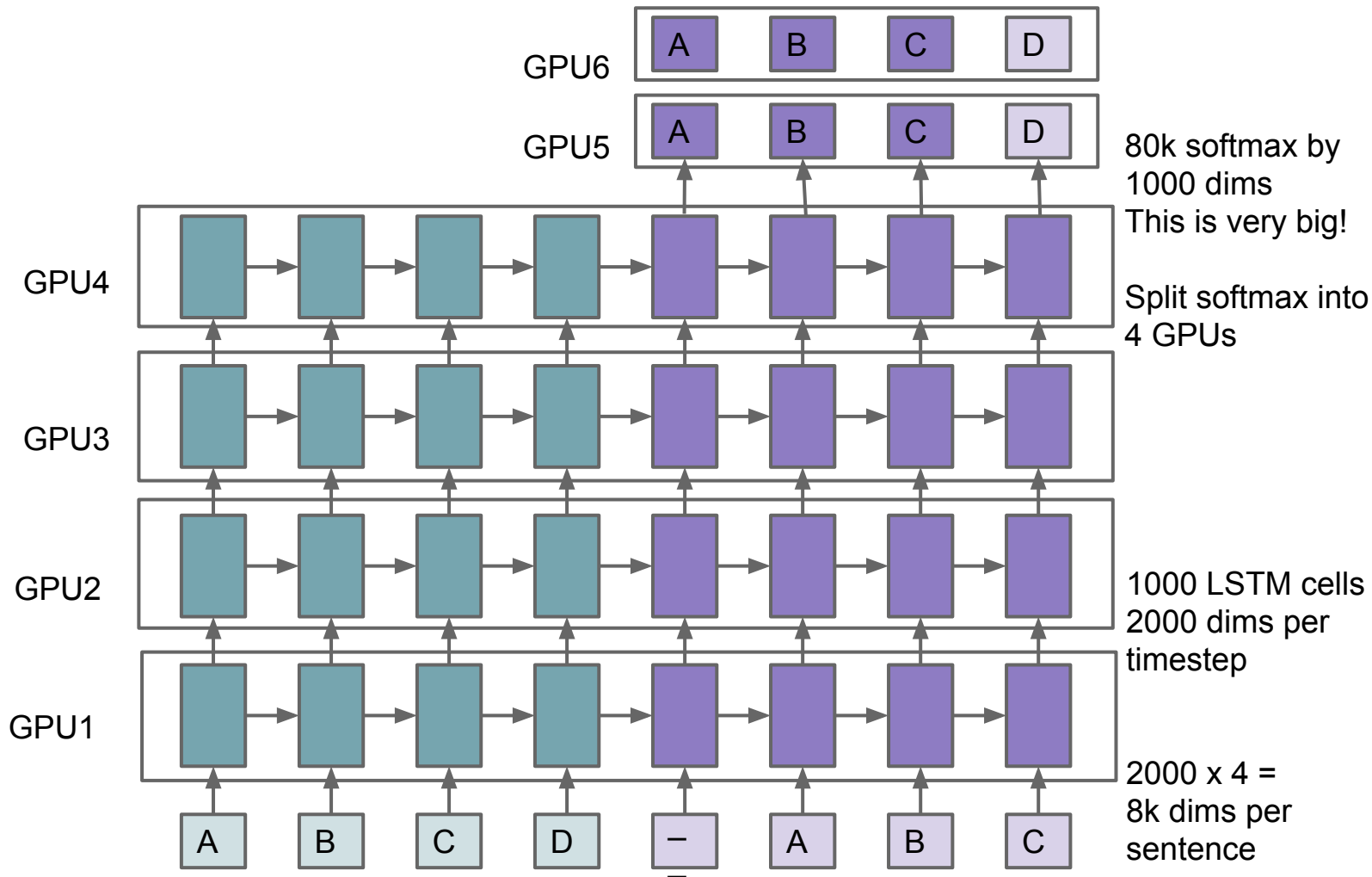


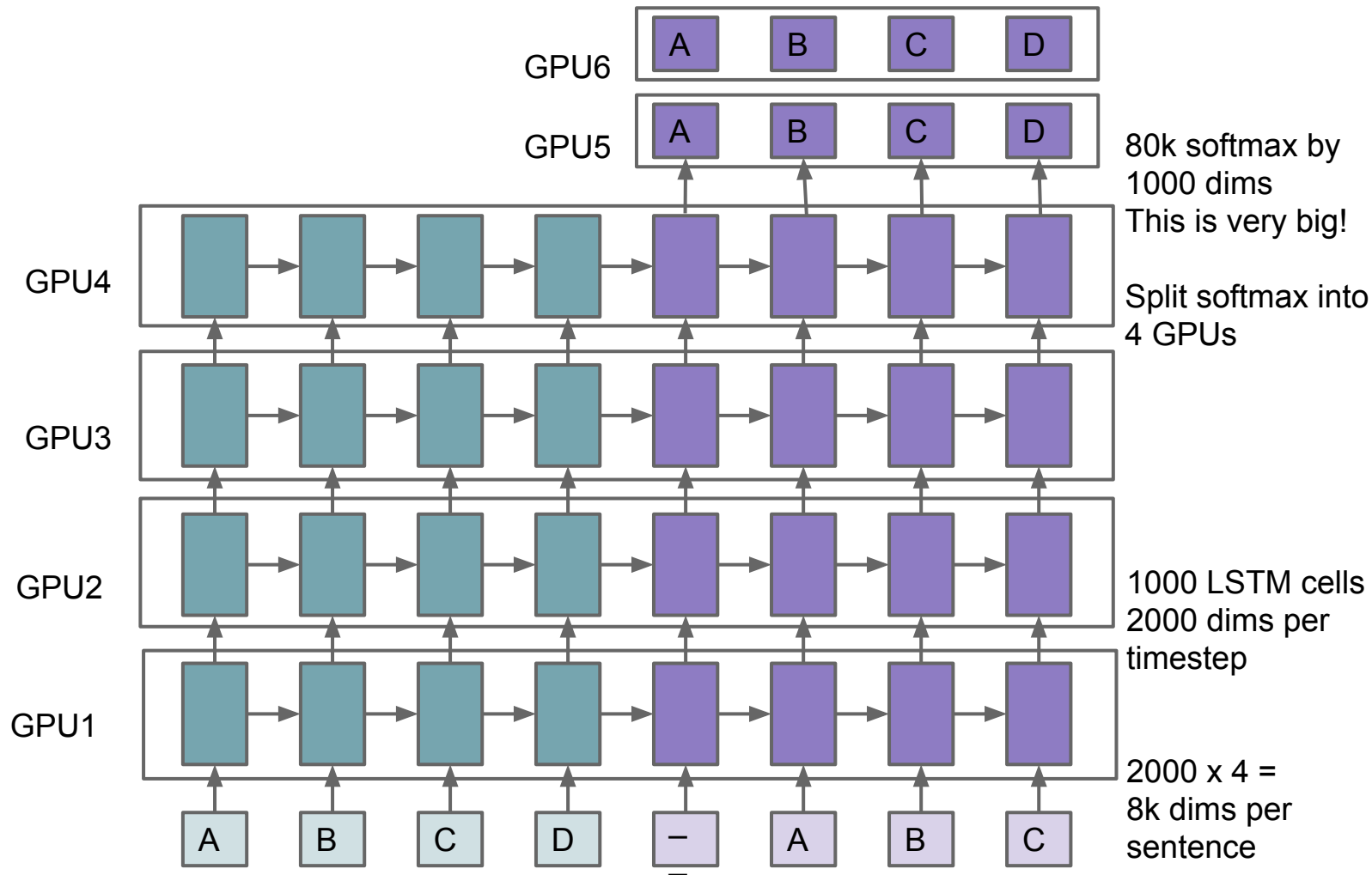










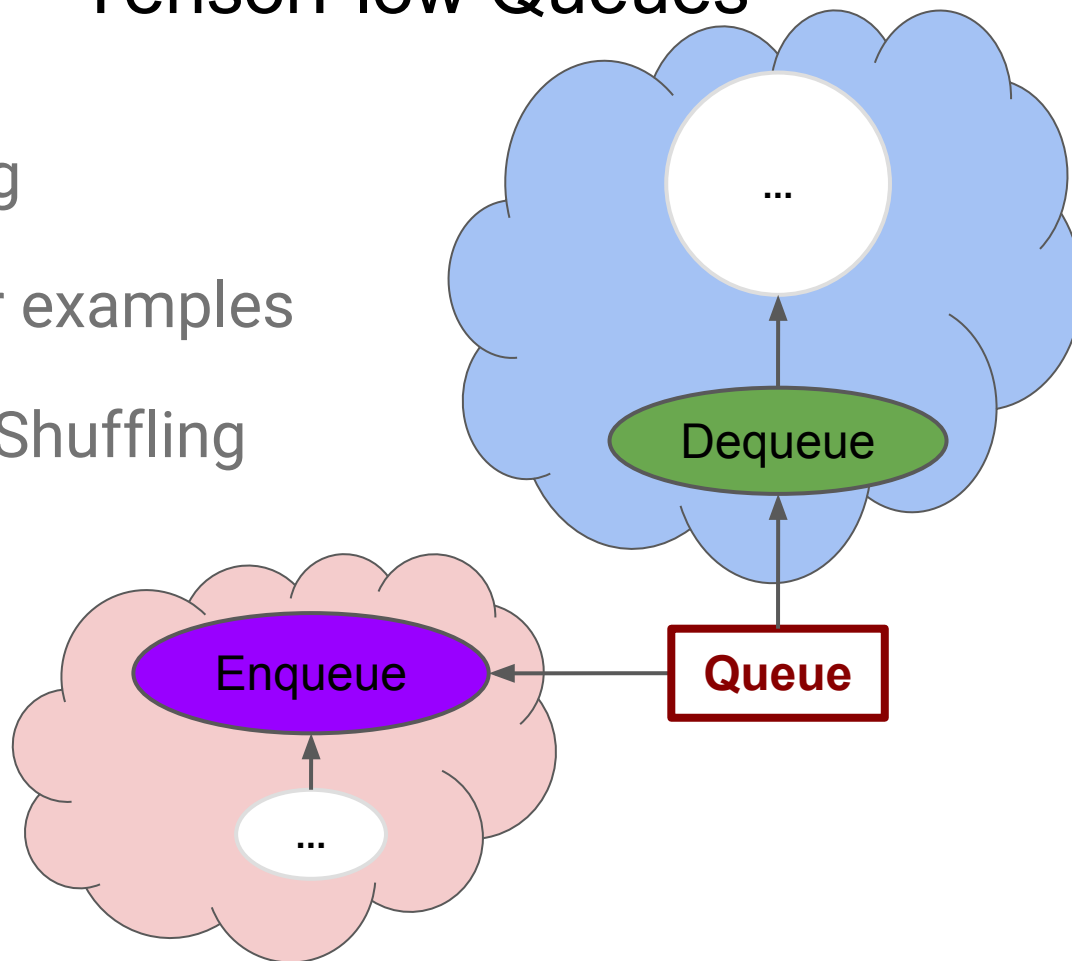


TensorFlow Queues

Input prefetching

Grouping similar examples

Randomization/Shuffling



Example: Deep LSTMs

- Wrinkles
 - Bucket sentences by length using a queue per length
 - Dequeue when a full batch of same length has accumulated
 - N different graphs for different lengths
 - Alternative: while loop



Expressing Data Parallelism

```
# We use the ReplicaDeviceSetter() device function to automatically
# assign Variables to the 'ps' jobs.
with tf.device("/cpu:0"):
    # Create the Mnist model.
    model = MnistModel(batch_size=16, hidden_units=200)

    # Get an initialized, and possibly recovered session.
    sess = tf.Session()

    # Train the model.
    for local_step in xrange(FLAGS.max_steps):
        _, loss, step = sess.run([model.train_op, model.loss, model.global_step])
        if local_step % 1000 == 0:
            print "step %d: %g" % (step, loss)
```



Expressing Data Parallelism

```
# We use the ReplicaDeviceSetter() device function to automatically
# assign Variables to the 'ps' jobs.
with tf.device(tf.ReplicaDeviceSetter(parameter_devices=10)):
    # Create the Mnist model.
    model = MnistModel(batch_size=16, hidden_units=200)

    # Create a Supervisor. It will take care of initialization, summaries,
    # checkpoints, and recovery. When multiple replicas of this program are running,
    # the first one, identified by --task=0 is the 'chief' supervisor (e.g., initialization, saving)
    supervisor = tf.Supervisor(is_chief=(FLAGS.task == 0), saver=model.saver)

    # Get an initialized, and possibly recovered session.
    sess = supervisor.PrepareSession(FLAGS.master_job)

    # Train the model.
    for local_step in xrange(int32_max):
        _, loss, step = sess.run([model.train_op, model.loss, model.global_step])
        if step >= FLAGS.max_steps:
            break
        if local_step % 1000 == 0:
            print "step %d: %g" % (step, loss)
```



Combining Vision with Robotics

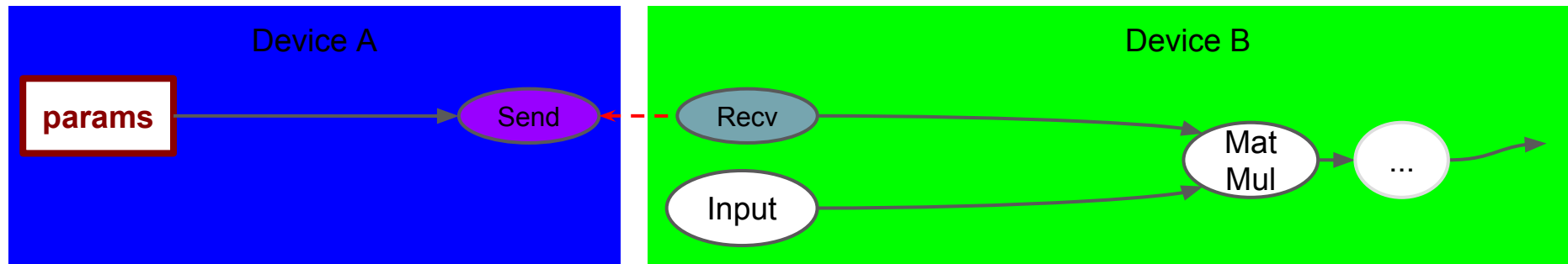
“Deep Learning for Robots: Learning from Large-Scale Interaction”, Google Research Blog, March, 2016

“Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection”, Sergey Levine, Peter Pastor, Alex Krizhevsky, & Deirdre Quillen, Arxiv, arxiv.org/abs/1603.02199



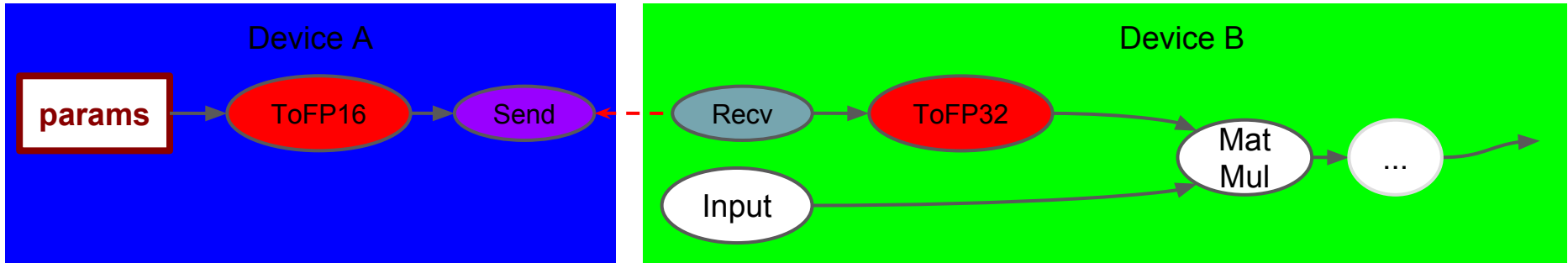
Network Optimizations

- Neural net training very tolerant of reduced precision
- e.g. drop precision to 16 bits across network



Network Optimizations

- Neural net training very tolerant of reduced precision
- e.g. drop precision to 16 bits across network



Quantization for Inference

- Need even less precision for inference
- 8-bit fixed point works well, but many ways of quantizing
- Critical for things like mobile devices
 - w/quantization, high-end smart phone can run Inception model at >6 frames per second (fps)



How Can You Get Started with Machine Learning?

Three ways, with varying complexity:

- (1) Use a Cloud-based API (Vision, Speech, etc.)
- (2) Use an existing model architecture, and retrain it or fine tune on your dataset
- (3) Develop your own machine learning models for new problems

More
flexible,
but more
effort
required



Use Cloud-based APIs



GOOGLE TRANSLATE API

Dynamically translate between thousands of available language pairs

cloud.google.com/translate



CLOUD SPEECH API ^{ALPHA}

Speech to text conversion powered by machine learning

cloud.google.com/speech



CLOUD VISION API

Derive insight from images with our powerful Cloud Vision API

cloud.google.com/vision

CLOUD TEXT API ^{ALPHA}

Use Cloud Text API for sentiment analysis and entity recognition in a piece of text.

cloud.google.com/text

Use Cloud-based APIs



GOOGLE TRANSLATE API

Dynamically translate between thousands of available language pairs

cloud.google.com/translate



CLOUD SPEECH API ^{ALPHA}

Speech to text conversion powered by machine learning

cloud.google.com/speech



CLOUD VISION API

Derive insight from images with our powerful Cloud Vision API

cloud.google.com/vision

CLOUD TEXT API ^{ALPHA}

Use Cloud Text API for sentiment analysis and entity recognition in a piece of text.

cloud.google.com/text

Google Cloud Vision API

<https://cloud.google.com/vision/>



"running", "score": 0.99803412,
"marathon", "score": 0.99482006



"joyLikelihood": "VERY_LIKELY"

"description": "ABIERTO\n",
"local": "es"

Google Cloud ML

Scaled service for training and inference w/TensorFlow



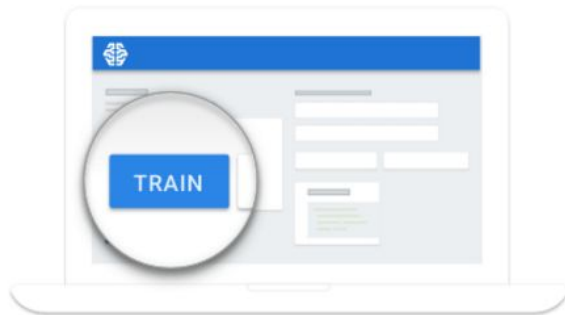
CLOUD MACHINE LEARNING ^{ALPHA}

Machine Learning on any data, of any size

[SIGN UP FOR THE LIMITED PREVIEW](#)

Managed scalable machine learning platform

Google Cloud Machine Learning is a managed platform that enables you to easily build machine learning models, that work on any type of data, of any size. Create your model with the powerful [TensorFlow](#) framework, that powers many Google products from [Google Photos](#), to [Google Cloud Speech](#). Build models of any size with our



A Few TensorFlow Community Examples

(From more than 2000 results for `'tensorflow'` on GitHub)

- DQN: github.com/nivwusquorum/tensorflow-deepq
- NeuralArt: github.com/woodrush/neural-art-tf
- Char RNN: github.com/sherjilozair/char-rnn-tensorflow
- Keras ported to TensorFlow: github.com/fchollet/keras
- Show and Tell: github.com/jazzsaxmafia/show_and_tell.tensorflow
- Mandarin translation: github.com/jikexueyuanwiki/tensorflow-zh

...



A Few TensorFlow Community Examples

(From more than ~~2000~~ 2100 results for 'tensorflow' on GitHub)

- DQN: github.com/nivwusquorum/tensorflow-deepq
- NeuralArt: github.com/woodrush/neural-art-tf
- Char RNN: github.com/sherjilozair/char-rnn-tensorflow
- Keras ported to TensorFlow: github.com/fchollet/keras
- Show and Tell: github.com/jazzsaxmafia/show_and_tell.tensorflow
- Mandarin translation: github.com/jikexueyuanwiki/tensorflow-zh

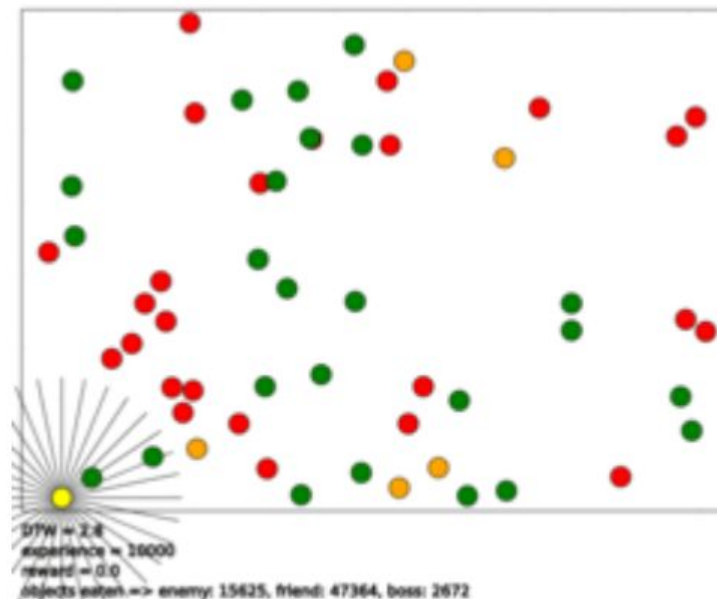
...



Reinforcement Learning using Tensor Flow

Quick start

Check out Karpathy game in `notebooks` folder.



The image above depicts a strategy learned by the DeepQ controller. Available actions are accelerating top, bottom, left or right. The reward signal is +1 for the green fellas, -1 for red and -5 for orange.



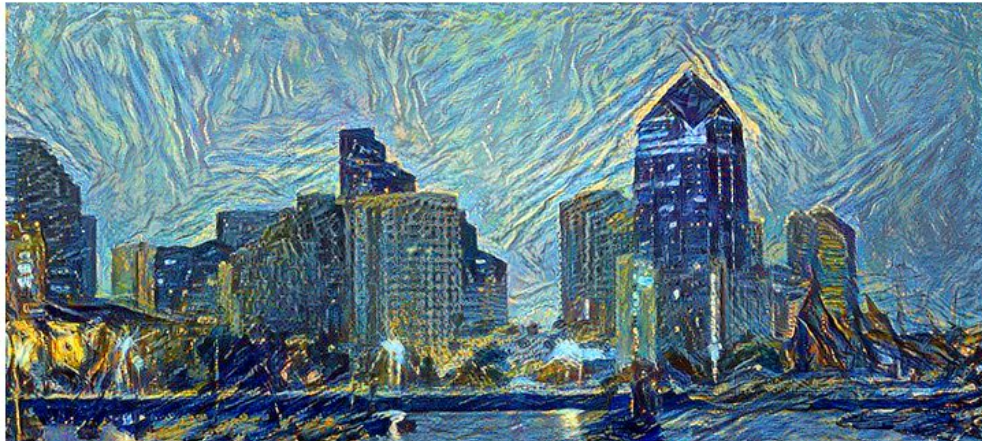
"Neural Art" in TensorFlow

An implementation of "A neural algorithm of Artistic style" in TensorFlow, for

- Introductory, hackable demos for TensorFlow, and
- Demonstrating the use of importing various Caffe cnn models (VGG and illustration2vec) in TF.

In this work, I put effort in putting the code simple as possible, for being a good introductory code to TF. For this reason, I also implemented very basic uses of TensorBoard (the visualizer). I also aimed on demonstrating the use of importing various Caffe models from *.caffemodel files into TensorFlow, especially models that seemed not to be imported by anybody yet in TF (as far as I know). Based on <https://github.com/ethereon/caffe-tensorflow>, I modified the importer so that it can import illustration2vec (<http://illustration2vec.net/>), which is another CNN available as a Caffe model. Using different CNNs yields different results, which reflects the characteristics of the model.

In the Neural Art problem setting, the weights of the CNN are fixed, and the input image into the CNN is the only "trainable" variable, making the code easy to understand (the optimized/trained image is the output image). I hope this example serves as a good introduction to TensorFlow as well as for entertainment purposes.



char-rnn-tensorflow

Multi-layer Recurrent Neural Networks (LSTM, RNN) for character-level language models in Python using Tensorflow.

Inspired from Andrej Karpathy's [char-rnn](#).

Requirements

- [Tensorflow](#)

Basic Usage

To train with default parameters on the tinyshakespeare corpus, run `python train.py`.

To sample from a checkpointed model, `python sample.py`.



Keras: Deep Learning library for Theano and TensorFlow

You have just found Keras.

Keras is a minimalist, highly modular neural networks library, written in Python and capable of running either on top of either [TensorFlow](#) or [Theano](#). It was developed with a focus on enabling fast experimentation. Being able to go from idea to result with the least possible delay is key to doing good research.

Use Keras if you need a deep learning library that:

- allows for easy and fast prototyping (through total modularity, minimalism, and extensibility).
- supports both convolutional networks and recurrent networks, as well as combinations of the two.
- supports arbitrary connectivity schemes (including multi-input and multi-output training).
- runs seamlessly on CPU and GPU.

Read the documentation at [Keras.io](https://keras.io).

Keras is compatible with: - **Python 2.7-3.5** with the Theano backend - **Python 2.7** with the TensorFlow backend



Neural Caption Generator

- Implementation of "Show and Tell" <http://arxiv.org/abs/1411.4555>
 - Borrowed some code and ideas from Andrej Karpathy's NeuralTalk.
- You need flickr30k data (images and annotations)

Code

- `make_flickr_dataset.py` : Extracting feats of flickr30k images, and save them in `'./data/feats.npy'`
- `model_tensorflow.py` : TensorFlow Version
- `model_theano.py` : Theano Version

Usage

- Flickr30k Dataset Download
- Extract VGG Features of Flickr30k images (`make_flickr_dataset.py`)
- Train: run `train()` in `model_tensorflow.py` or `model_theano.py`
- Test: run `test()` in `model_tensorflow.py` or `model_theano.py`.
 - parameters: VGG FC7 feature of test image, trained model path



TensorFlow is an Open Source Software
Library for Machine Intelligence

GET STARTED

你正在翻译的项目可能会比 **Android** 系统更加深远地影响着世界！

缘起

2015年11月9日，Google 官方在其博客上称，Google Research 宣布推出第二代机器学习系统 TensorFlow，针对先前的 DistBelief 的短板有了各方面的加强，更重要的是，它是开源的，任何人都可以用。

机器学习作为人工智能的一种类型，可以让软件根据大量的数据来对未来的情况进行阐述或预判。如今，领先的科技巨头无不在机器学习下予以极大投入。Facebook、苹果、微软，甚至国内的百度。Google 自然也在其中。「TensorFlow」是 Google



Concluding Remarks

- Model and Data Parallelism enable great ML work:
 - Neural Machine Translation: ~6x speedup on 8 GPUs
 - Inception / Imagenet: ~40x speedup on 50 GPUs
 - RankBrain: ~300X speedup on 500 machines
- TensorFlow open-source community vibrant and growing
- TensorFlow makes it easy to express ML computations



What Does the Future Hold?

Deep learning usage will continue to grow and accelerate:

- Across more and more fields and problems:
 - robotics, self-driving vehicles, ...
 - health care
 - video understanding
 - dialogue systems
 - personal assistance
 - ...





Google Brain Residency Program

One year immersion program in deep learning research

- First class started six weeks ago, planning for next year's class is underway

Learn to conduct deep learning research w/experts in our team

- Fixed one-year employment with salary, benefits, ...
- Goal after one year is to have conducted several research projects
- Interesting problems, TensorFlow, and access to computational resources

g.co/brainresidency



Google Brain Residency Program

Who should apply?

- people with BSc, MSc or PhD, ideally in CS, mathematics or statistics
- completed coursework in calculus, linear algebra, and probability, or equiv.
- programming experience
- motivated, hard working, and have a strong interest in deep learning

g.co/brainresidency



Google Brain Residency Program

Current class for June 2016 to May 2017

- $\frac{1}{3}$ B.S., $\frac{1}{3}$ M.S., $\frac{1}{3}$ Ph.D. or postdoc
- $\frac{1}{2}$ coming straight from school, $\frac{1}{2}$ with some post-school working experience
- Mix of backgrounds: computer scientists, math/stats, EE, physics, comp bio, ...

Applications for class for June 2017 to May 2018 will open in Fall 2016

g.co/brainresidency

Further Reading

- Dean, *et al.*, *Large Scale Distributed Deep Networks*, NIPS 2012, research.google.com/archive/large_deep_networks_nips2012.html.
- Mikolov, Chen, Corrado & Dean. *Efficient Estimation of Word Representations in Vector Space*, NIPS 2013, arxiv.org/abs/1301.3781.
- Sutskever, Vinyals, & Le, *Sequence to Sequence Learning with Neural Networks*, NIPS, 2014, arxiv.org/abs/1409.3215.
- Vinyals, Toshev, Bengio, & Erhan. *Show and Tell: A Neural Image Caption Generator*. CVPR 2015. arxiv.org/abs/1411.4555
- TensorFlow white paper, tensorflow.org/whitepaper2015.pdf (clickable links in bibliography)

g.co/brain (We're hiring! Also check out Brain Residency program at g.co/brainresidency)
www.tensorflow.org
research.google.com/people/jeff
research.google.com/pubs/BrainTeam.html

Questions?

