

Trends and Developments in Deep Learning Research

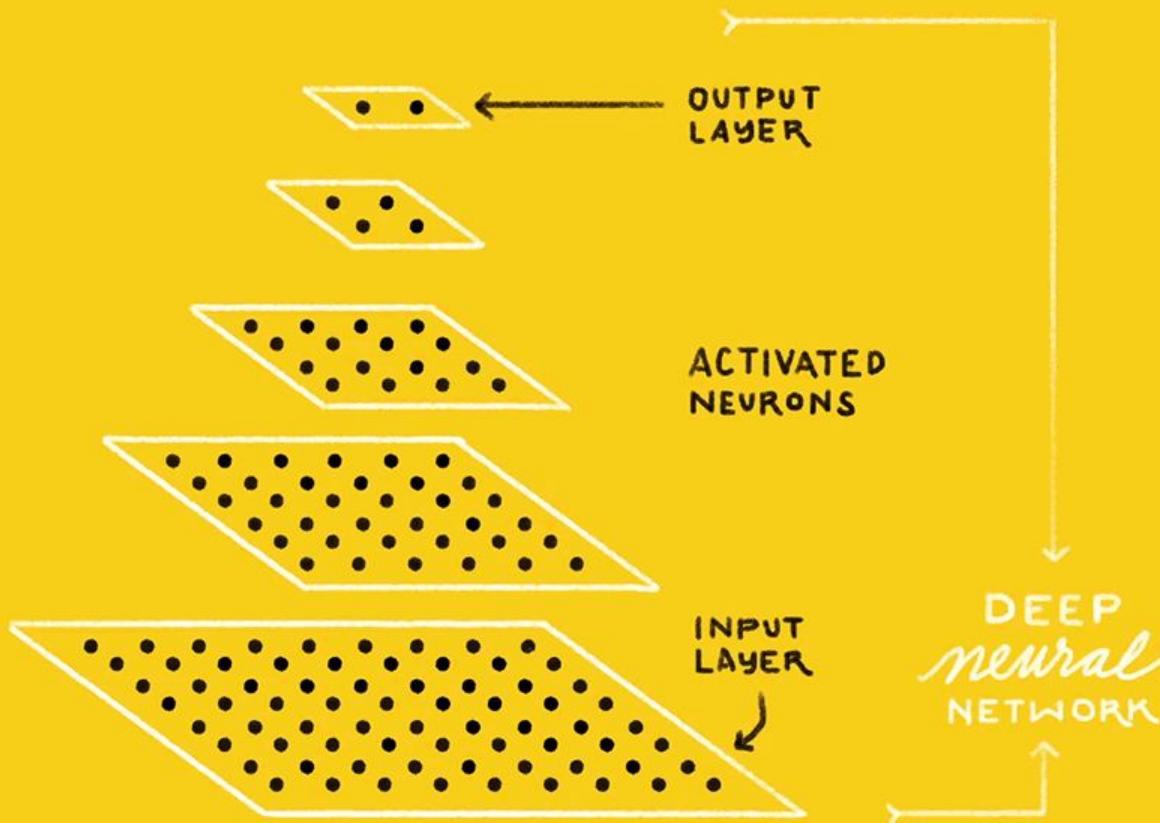
Jeff Dean
Google Brain team
g.co/brain

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

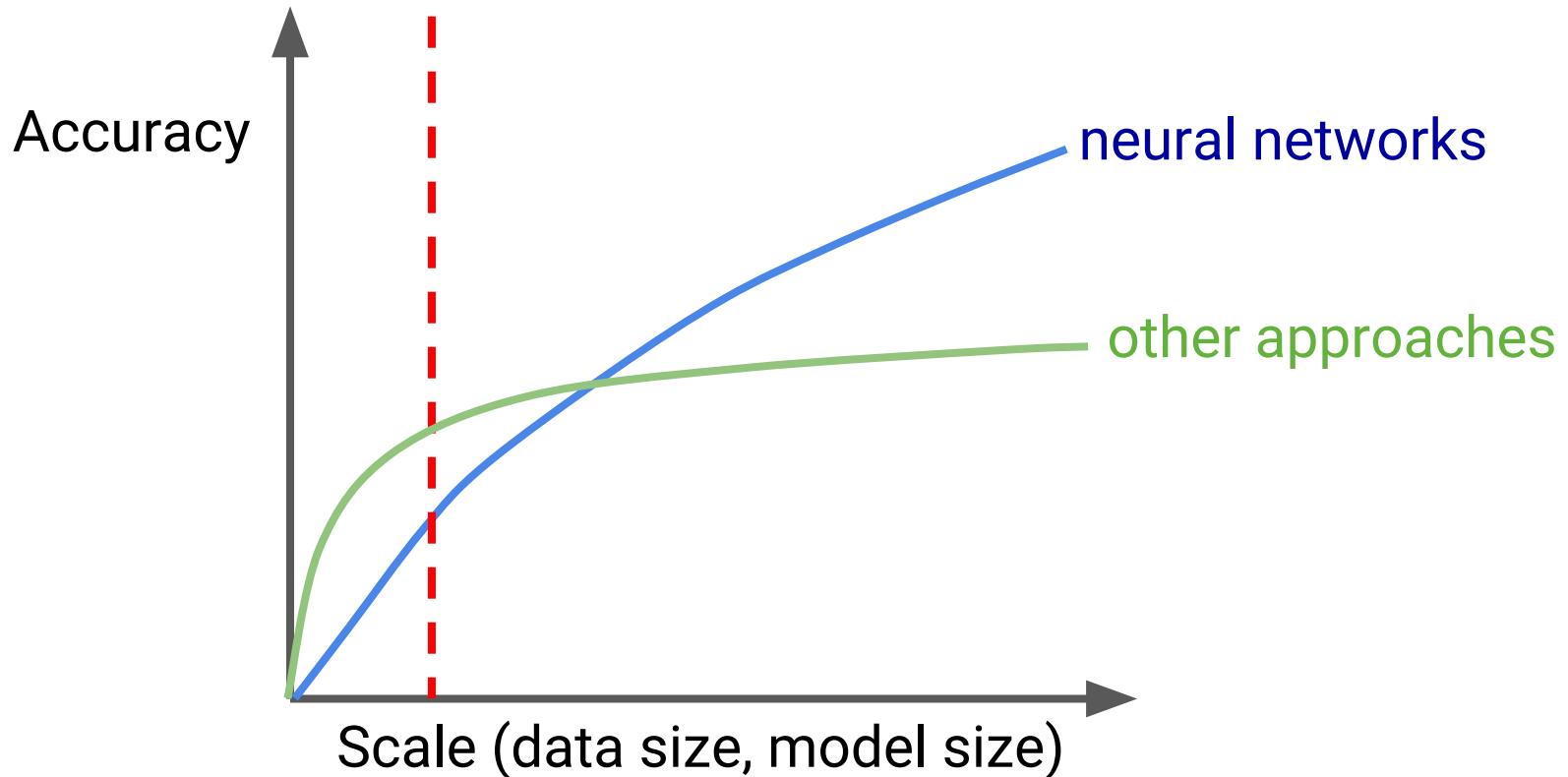
IS THIS A
CAT or DOG?



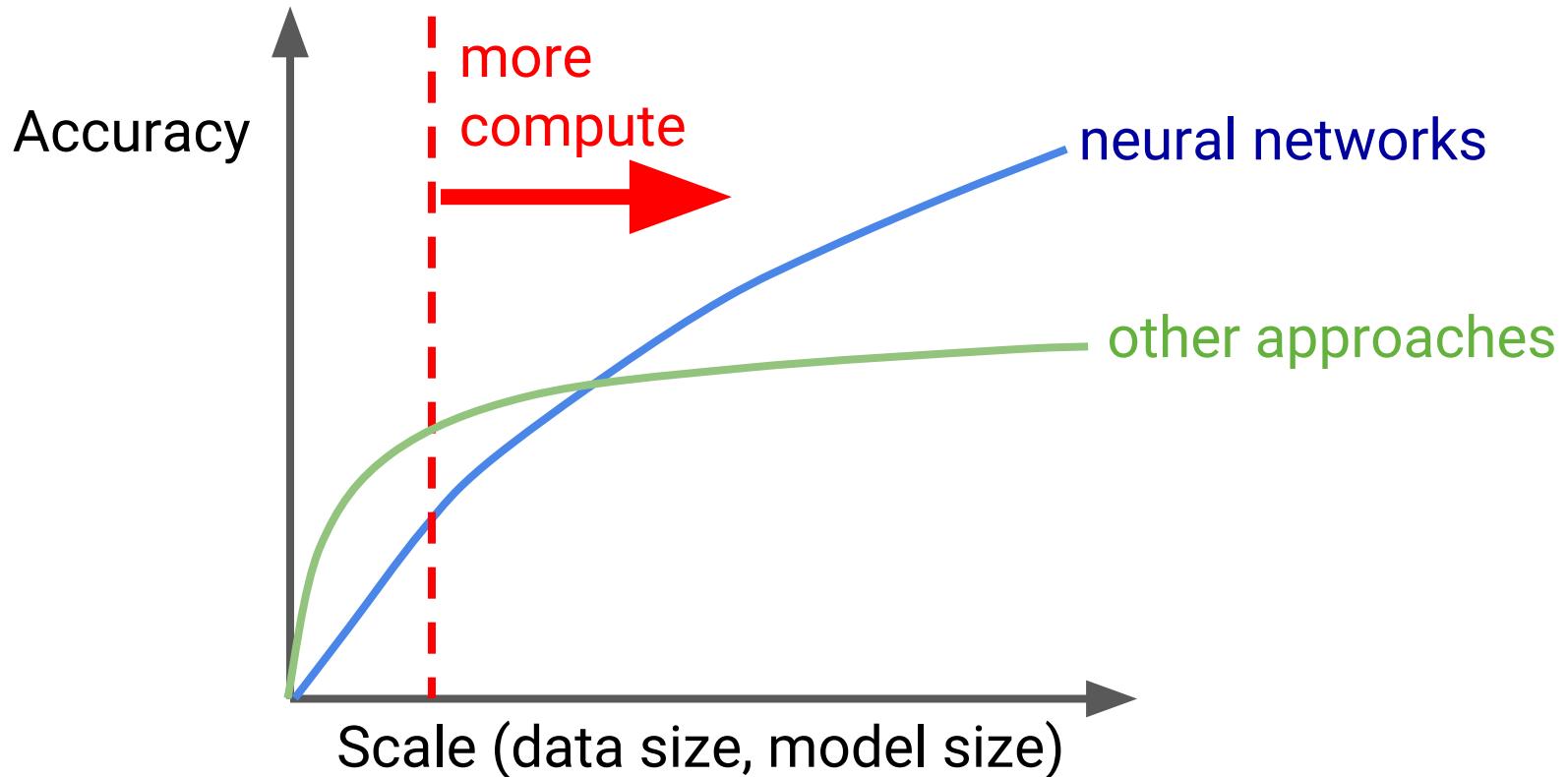
CAT DOG

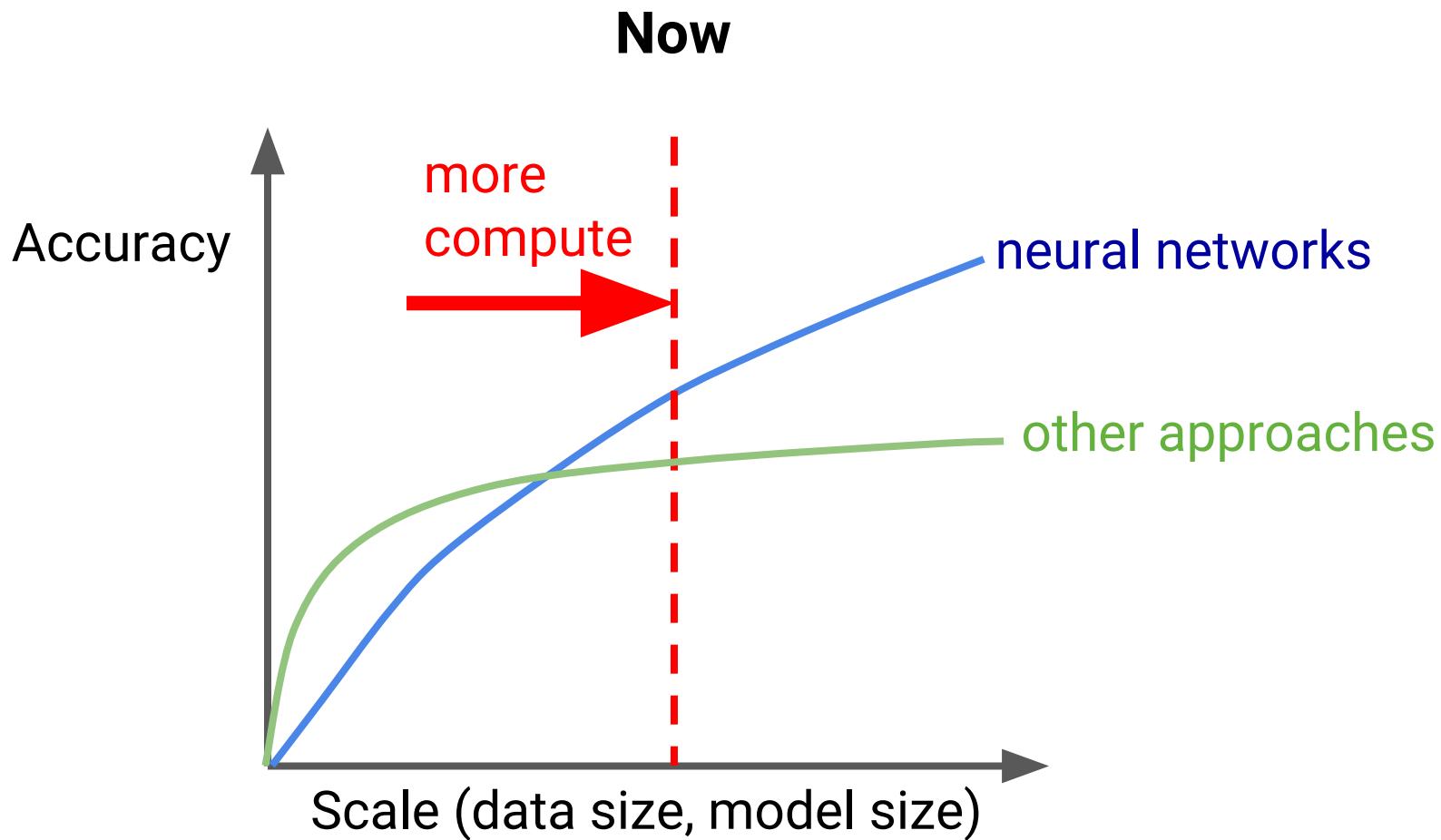


1980s and 1990s



1980s and 1990s



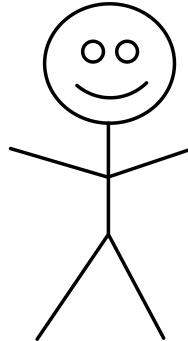


2011



26% errors

humans



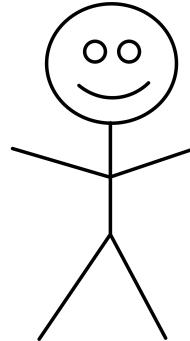
5% errors

2011



26% errors

humans



5% errors

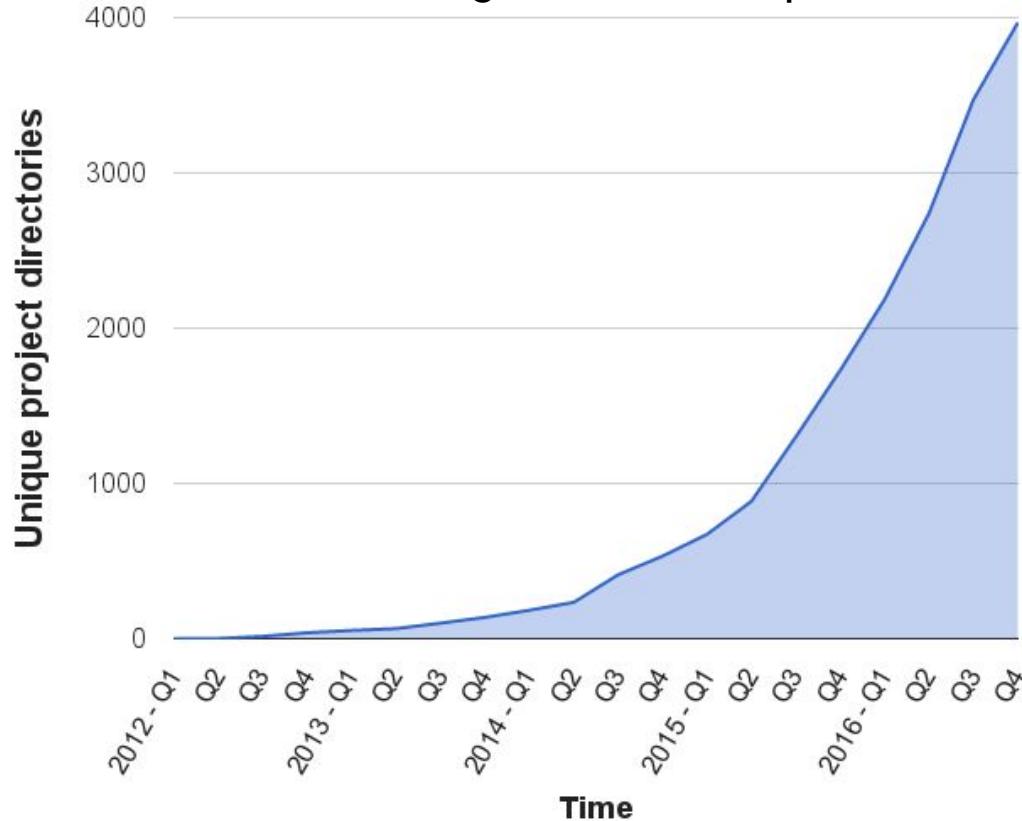
2016



3% errors

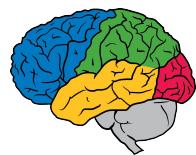
Growing Use of Deep Learning at Google

of directories containing model description files



Across many products/areas:

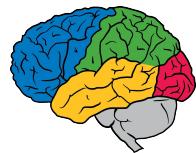
Android
Apps
drug discovery
Gmail
Image understanding
Maps
Natural language understanding
Photos
Robotics research
Speech
Translation
YouTube
... many others ...



Google Brain Team

- **Research:** 27 papers in ICML, NIPS, and ICLR in 2016, plus others in venues like ACL, ICASSP, CVPR, ICER, and OSDI; Brain Residency Program, hosted ~50 interns in 2016, ...
- **Product impact:** Many dozens of high impact collaborations in Google products/efforts like Search, Ads, Photos, Translate, GMail, Maps, Cloud ML, speech recognition, self-driving cars, robotics, ...
- **Open-source tools for ML:** TensorFlow, Magenta, visualization tools, ...

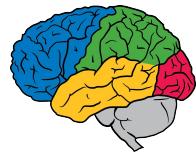
g.co/brain



Need to build the right tools

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

Guides

TUTORIALS HOW TO

Tutorials

Basic Neural Networks

MNIST For ML Beginners

Deep MNIST for Experts

TensorFlow Mechanics 101

Easy ML with tf.contrib.learn

tf.contrib.learn Quickstart

Large-scale Linear Models with TensorFlow

TensorFlow Linear Model Tutorial

TensorFlow Wide & Deep Learning Tutorial

Logging and Monitoring Basics with tf.contrib.learn

Building Input Functions with tf.contrib.learn

Creating Estimators in tf.contrib.learn

TensorFlow Serving

TensorFlow Serving

Image Processing

Convolutional Neural Networks

Image Recognition

Language and Sequence Processing

Vector Representations of Words

Recurrent Neural Networks

Sequence-to-Sequence Models

SyntaxNet

Non-ML Applications

Mandelbrot Set

Partial Differential Equations

TensorFlow Versions

Tutorials

Basic Neural Networks

The first few Tensorflow tutorials guide you through training and testing a simple neural network to classify handwritten digits from the MNIST database of digit images.

MNIST For ML Beginners

If you're new to machine learning, we recommend starting here. You'll learn about a classic problem, handwritten digit classification (MNIST), and get a gentle introduction to multiclass classification.

[View Tutorial](#)

Deep MNIST for Experts

If you're already familiar with other deep learning software packages, and are already familiar with MNIST, this tutorial will give you a very brief primer on TensorFlow.

[View Tutorial](#)

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 MNIST as the example.

[View Tutorial](#)

Easy ML with tf.contrib.learn

tf.contrib.learn Quickstart

A quick introduction to tf.contrib.learn, a high-level API for TensorFlow. Build, train, and evaluate a neural network with just a few lines of code.

[View Tutorial](#)

Contents

Basic Neural Networks

MNIST For ML Beginners

Deep MNIST for Experts

TensorFlow Mechanics 101

Easy ML with tf.contrib.learn

tf.contrib.learn Quickstart

Overview of Linear Models with tf.contrib.learn

Linear Model Tutorial

Wide and Deep Learning Tutorial

Logging and Monitoring Basics with tf.contrib.learn

Building Input Functions with tf.contrib.learn

Creating Estimators in tf.contrib.learn

TensorFlow Serving

TensorFlow Serving

Image Processing

Convolutional Neural Networks

Image Recognition

Deep Dream Visual Hallucinations

Language and Sequence Processing

Vector Representations of Words

Recurrent Neural Networks

Sequence-to-Sequence Models

SyntaxNet: Neural Models of Syntax

Non-ML Applications

Why Did We Build TensorFlow?

Wanted system that was **flexible**, **scalable**, and **production-ready**

DistBelief, our first system, was good on two of these, but lacked **flexibility**

Most existing open-source packages were also good on 2 of 3 but not all 3

TensorFlow Goals

Establish **common platform** for expressing machine learning ideas and systems

Make this platform the **best in the world** for both research and production use

Open source it so that it becomes a **platform for everyone**, not just Google

Facts and Figures

Launched on Nov. 9, 2015

Initial launch was reasonably fully-featured:

auto differentiation, queues, control flow, fairly comprehensive set of ops, ...

tutorials made system accessible

out-of-the-box support for CPUs, GPUs, multiple devices, multiple platforms

Some Stats

500+ contributors, most of them outside Google

12,000+ commits since Nov, 2015 (~30 per day)

1M+ binary downloads

#15 most popular repository on GitHub by stars (just passed Linux!)

Used in ML classes at quite a few universities now:

Toronto, Berkeley, Stanford, ...

Many companies/organizations using TensorFlow:

Google, DeepMind, OpenAI, Twitter, Snapchat, Airbus, Uber, ...



tensorflow

Search

We've found 5,172 repository results

Sort: Most stars ▾

[tensorflow/tensorflow](#)

Computation using data flow graphs for scalable machine learning

Python ★ 42,335 19,643 Updated 8 hours ago

[fchollet/keras](#)

Deep Learning library for Python. Convnets, recurrent neural networks, and more. Runs on Theano or *TensorFlow*.

Python ★ 10,924 3,638 Updated 16 hours ago

[tensorflow/models](#)

Models built with *TensorFlow*

Python ★ 10,501 2,792 Updated 7 hours ago

[aymericdamien/TensorFlow-Examples](#)

TensorFlow Tutorial and Examples for beginners

Jupyter Notebook ★ 8,051 1,986 Updated 13 days ago

autoencoder	merged changes from #25	10 months ago
compression	Update README with results for comparison.	2 months ago
differential_privacy	Update deep_cnn.py	7 hours ago
im2txt	Update GraphKeys.VARIABLES to GraphKeys.GLOBAL_VARIABLES.	a month ago
inception	Update losses.py	7 hours ago
lm_1b	Fix README	4 months ago
namigner	add the namigner model (#147)	8 months ago
neural_gpu	Add to neural_gpu documentation.	7 months ago
neural_programmer	edits to README	a month ago
next_frame_prediction	Add cross conv model for next frame prediction.	17 days ago
resnet	DOC: Typo in resnet documentation	13 days ago
slim	Updating README.md	2 months ago
street	Update vgsl_model.py	7 hours ago
swivel	Add sys.stdout.flush()	3 months ago
syntaxnet	Fix POS tagging score of Ling et al.(2005)	3 months ago
textsum	Update data.py	a month ago
transformer	Update cluttered_mnist.py	7 hours ago
tutorials	Update cifar10.py	7 hours ago
video_prediction	video prediction model code	3 months ago

December 2016 (0.6)



Python 3.3+



Faster on GPUs

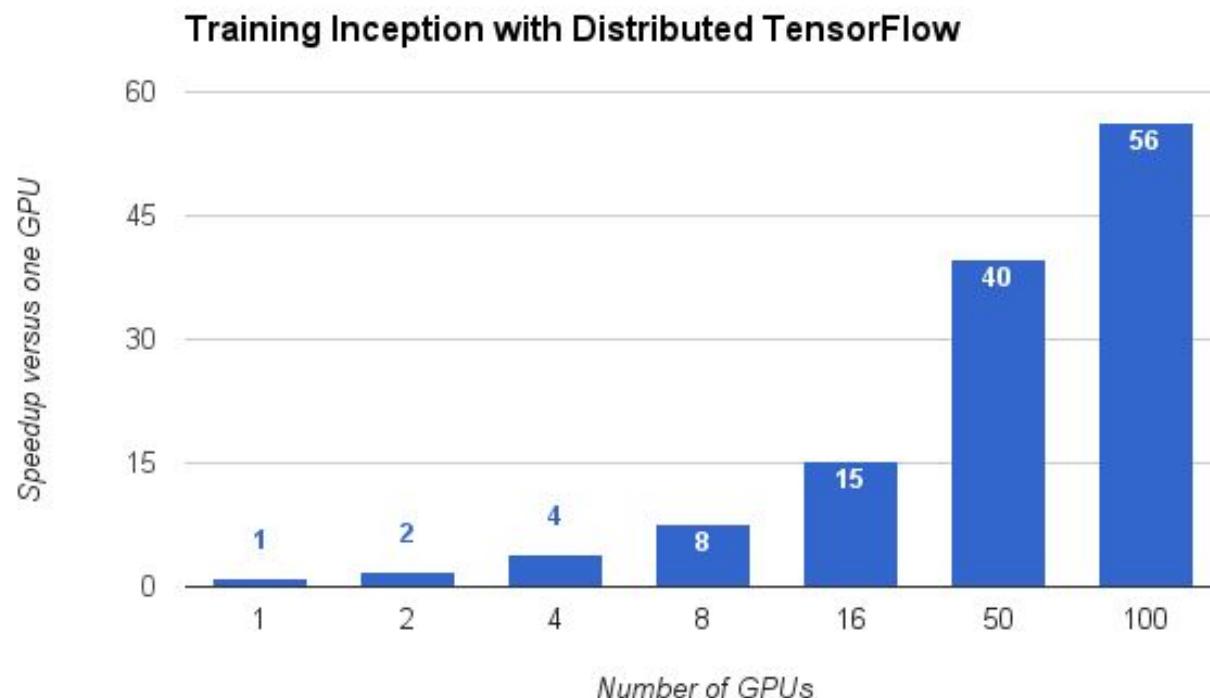
February 2016 (0.7)



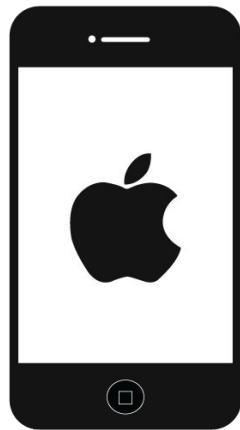
CuDNN v4

Dynamic
Loading of
Kernels

April 2016 (0.8)



June 2016 (0.9)

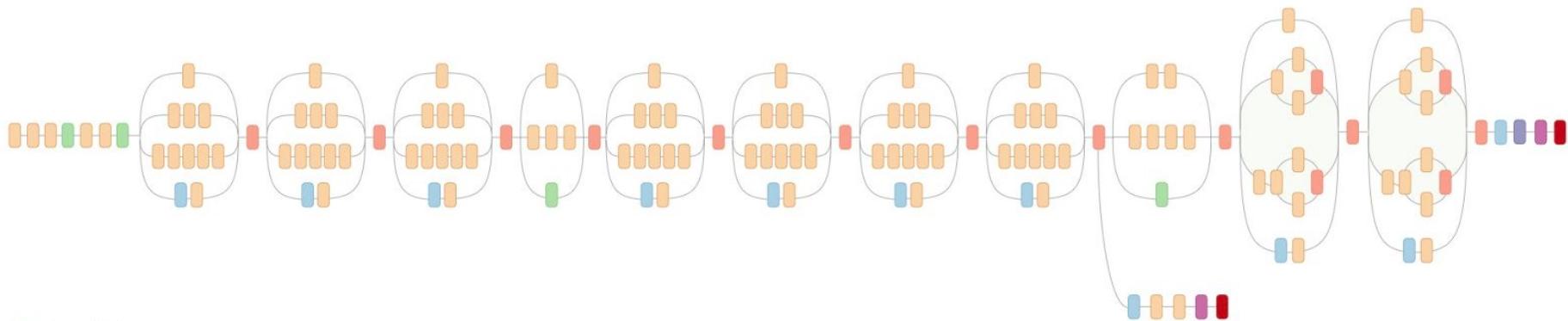


iOS



GPUs on Mac

August 2016 (0.10)



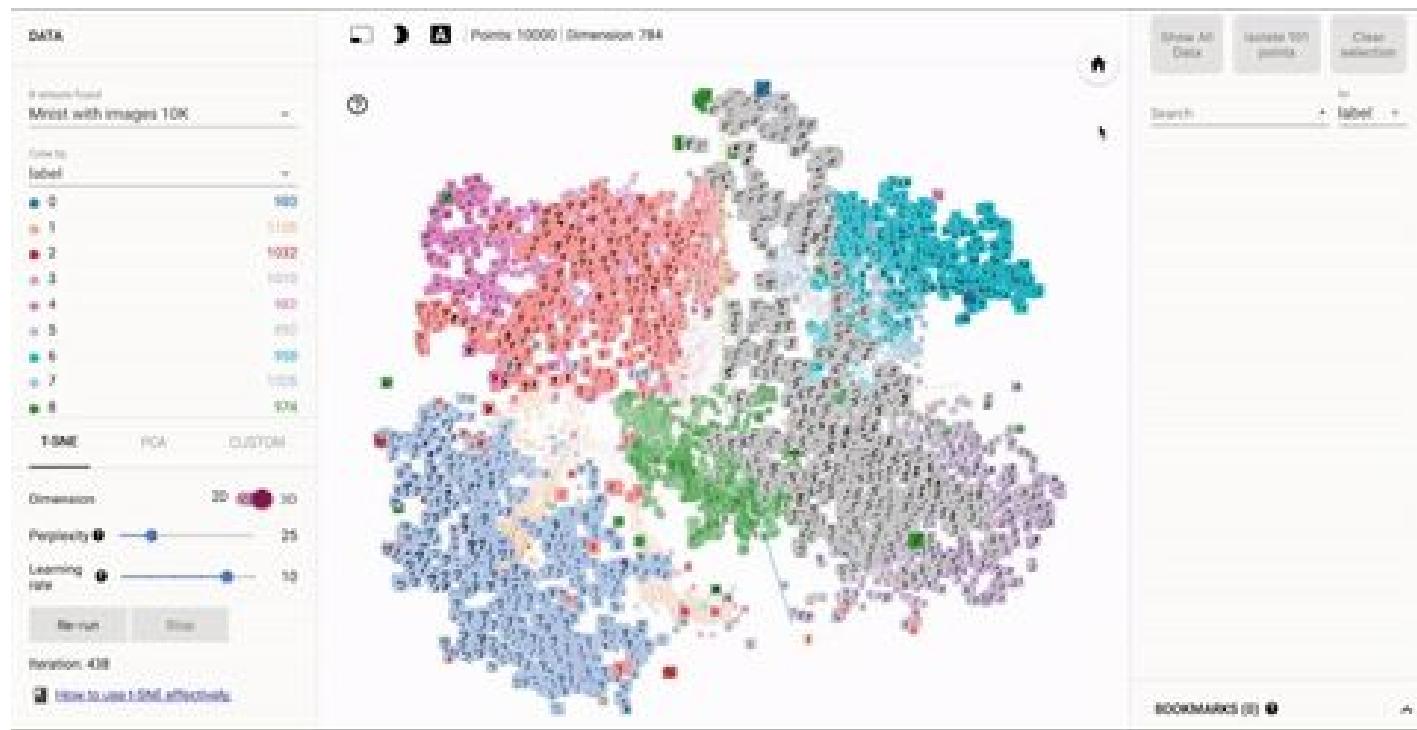
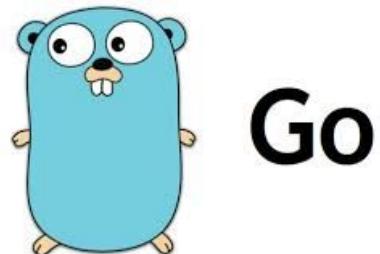
Slim

October 2016 (0.11)



CuDNN v5

November 2016 (0.12)



TensorFlow v0.11.0 RCO



gunan released this 24 days ago · 34 commits to r0.11 since this release

Major Features and Improvements

- cuDNN 5 support.
- HDFS Support.
- Adds Fused LSTM support via cuDNN 5 in `tensorflow/contrib/cudnn_rnn`.
- Improved support for NumPy style basic slicing including non-1 strides, ellipses, newaxis, and negative indices. For example complicated expressions like `foo[1, 2:4, tf.newaxis, ..., :-3:-1, :]` are now supported. In addition we have preliminary (non-broadcasting) support for sliced assignment to variables. In particular one can write `var[1:3].assign([1,11,111])`.
- Introducing `core/util/tensor_bundle` module: a module to efficiently serialize/deserialize tensors to disk. Will be used in TF's new checkpoint format.
- Added `tf.svd` for computing the singular value decomposition (SVD) of dense matrices or batches of matrices (CPU only).
- Added gradients for eigenvalues and eigenvectors computed using `self_adjoint_eig` or `self_adjoint_eigvals`.
- Eliminated `batch_*` methods for most linear algebra and FFT ops and promoted the non-batch version of the ops to handle batches of matrices.
- Tracing/timeline support for distributed runtime (no GPU profiler yet).
- C API gives access to inferred shapes with `TF_GraphGetTensorNumDims` and `TF_GraphGetTensorShape`.
- Shape functions for core ops have moved to C++ via `REGISTER_OP(...).SetShapeFn(...)`. Python shape inference `RegisterShape` calls use the C++ shape functions with `common_shapes.call_cpp_shape_fn`. A future release will remove `RegisterShape` from python.

Languages



C++



Go



Version 1.0 upcoming

Stability

Backwards Compatibility

Usability

Documentation

Libraries

Models

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*

<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

Appeared in OSDI 2016

<https://arxiv.org/abs/1605.08695>

Strong External Adoption



Adoption of Deep Learning Tools on GitHub



1M+ binary installs since November, 2015

Experiment Turnaround Time and Research Productivity

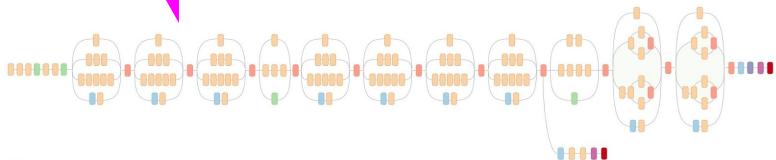
- **Minutes, Hours:**
 - **Interactive research! Instant gratification!**
- **1-4 days**
 - Tolerable
 - Interactivity replaced by running many experiments in parallel
- **1-4 weeks**
 - High value experiments only
 - Progress stalls
- **>1 month**
 - Don't even try



Just-In-Time Compilation

via TensorFlow's XLA, "Accelerated Linear Algebra" compiler

TF graphs go in,



Optimized & specialized assembly comes out.

```
0x00000000    movq    (%rdx), %rax  
0x00000003    vmovaps (%rax), %xmm0  
0x00000007    vmulps %xmm0, %xmm0, %xmm0  
0x0000000b    vmovaps %xmm0, (%rdi)  
...  
...
```

Let's explain that!

Demo: Inspect JIT code in TensorFlow iPython shell

XLA:CPU

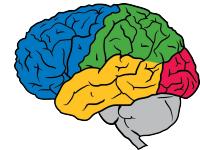
XLA:GPU

```
tensor( 0.0000e+00)
In [1]: %cpaste
Pasting code; enter '--' alone on the line to stop or use Ctrl-D.
:with tf.Session() as sess:
:    x = tf.placeholder(tf.float32, [4])
:    with tf.device("device:XLA_CPU:0"):
:        y = x * x
:    result = sess.run(y, {x: [1.5, 0.5, -0.5, -1.5]})

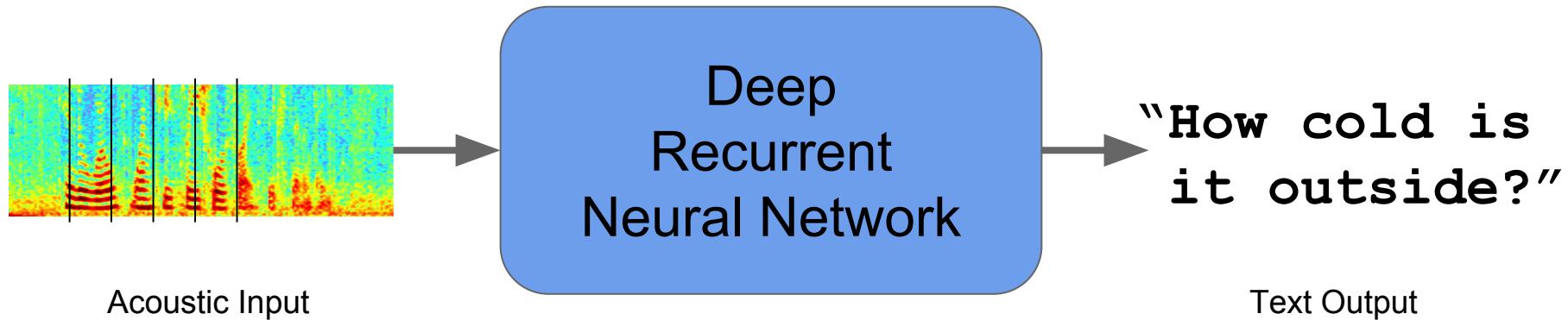

```

What are some ways that
deep learning is having
a significant impact at Google?

All of these examples implemented using TensorFlow
or our predecessor system



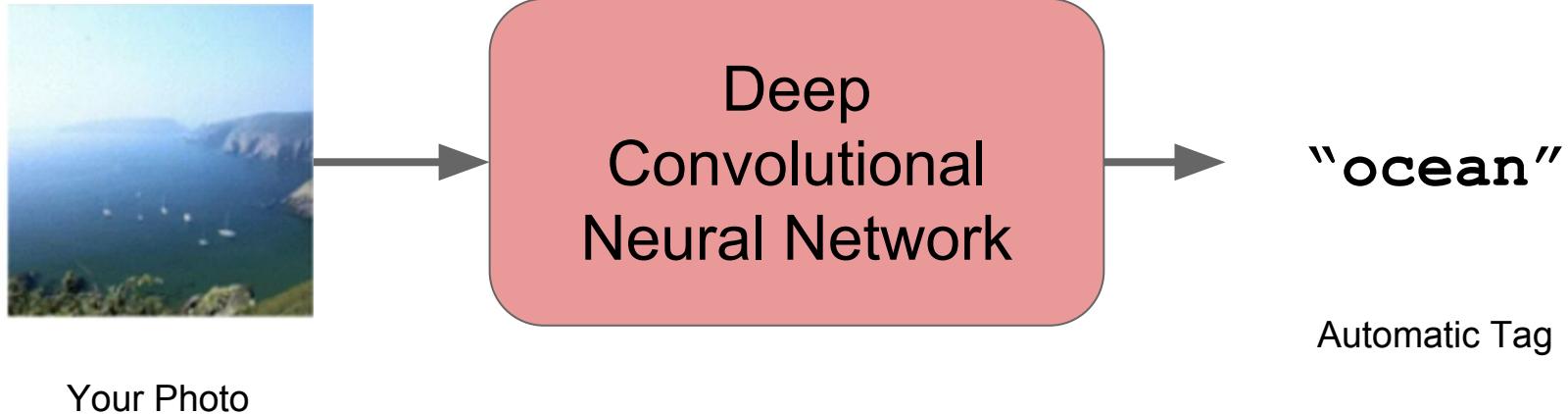
Speech Recognition



Reduced word errors by more than 30%

Google Research Blog - August 2012, August 2015

Google Photos Search



Search personal photos without tags.

Google Research Blog - June 2013

Google Photos Search

Things



Google my photos of siamese cats

Web Images Shopping Videos More ▾

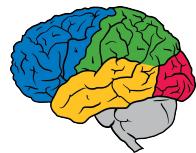
Your photos
Only you can see these results

The search results page shows a grid of 12 images of Siamese cats. The images include various poses and interactions with people, such as a cat sitting on a person's lap, a cat being held, and a cat sleeping. The images are arranged in three rows of four.

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 亞洲國際旅行社

Tel: 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 • VEHICLE INSPECTIONS •
- NEW CAR/ROADSIDE SERVICES • BRAKES • CLUTCHES •
- TYRES • SUSPENSION • TYRES • WHEEL ALIGNMENTS •
- AIR CONDITIONING • COOLANT/HYDRAULIC • OIL TREATMENT •
- FULL MAINTENANCE • BATTERIES • AUTO ELECTRICAL •

• Factory Trained Technicians



50

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

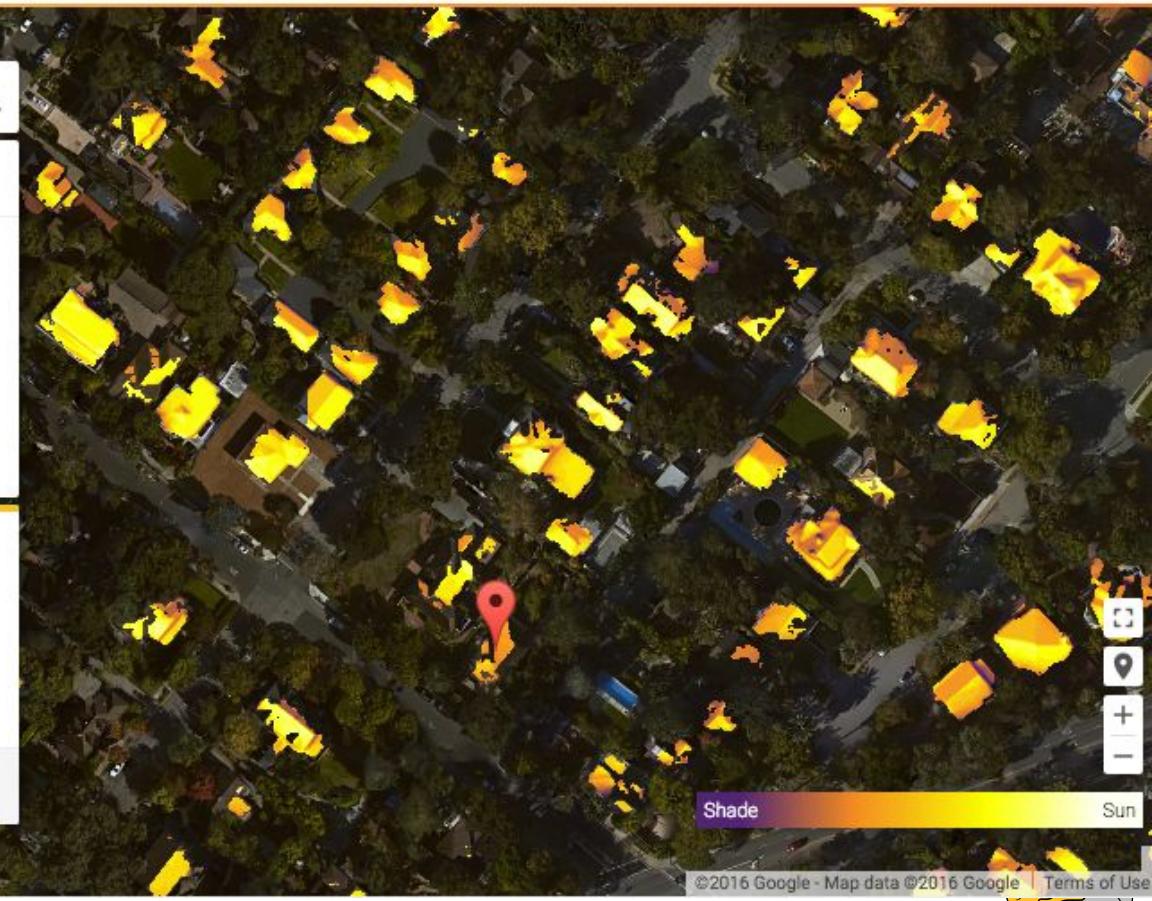
Based on 3D modeling of your roof and nearby trees

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.



MEDICAL IMAGING

Using similar model for detecting diabetic
retinopathy in retinal images

December 13, 2016

Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs

Varun Gulshan, PhD¹; Lily Peng, MD, PhD¹; Marc Coram, PhD¹; et al

» Author Affiliations

JAMA. 2016;316(22):2402-2410. doi:10.1001/jama.2016.17216

Performance **on par or slightly better** than the median of 8 U.S. board-certified ophthalmologists (F-score of 0.95 vs. 0.91).

<http://research.googleblog.com/2016/11/deep-learning-for-detection-of-diabetic.html>

Computers can now see

Large implications for machine learning for robotics

Combining Vision with Robotics

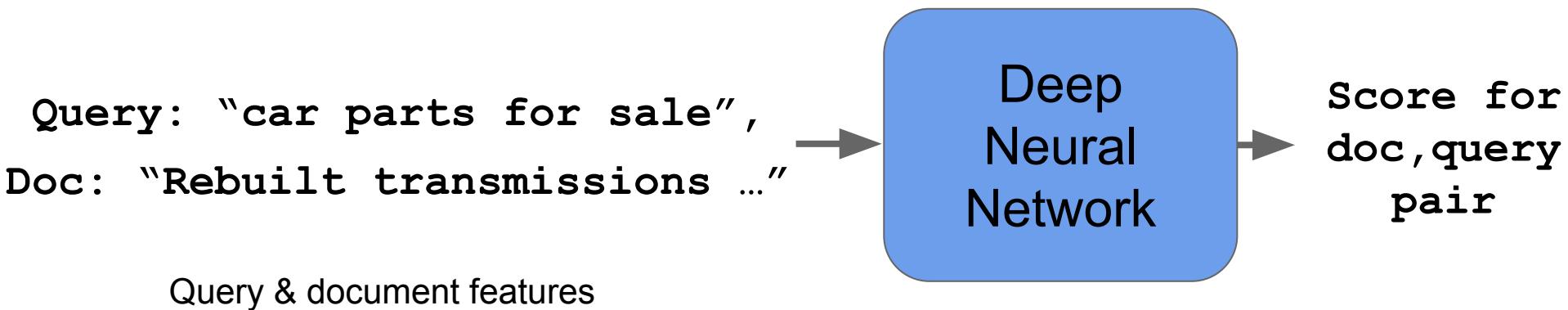
“Deep Learning for Robots: Learning from Large-Scale Interaction”, Google Research Blog, March, 2016

“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



Better language understanding

RankBrain in Google Search Ranking



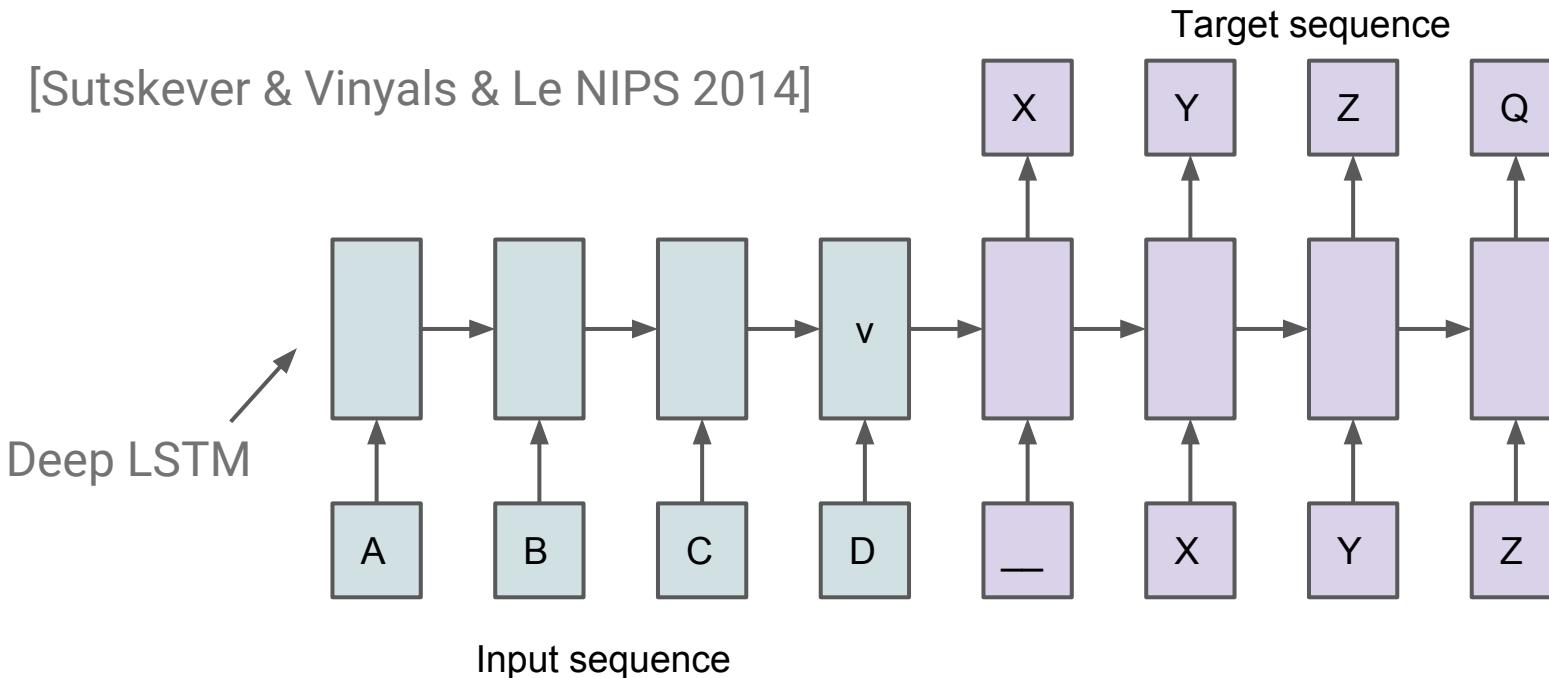
Launched in 2015

Third most important search ranking signal (of 100s)

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

Sequence-to-Sequence Model

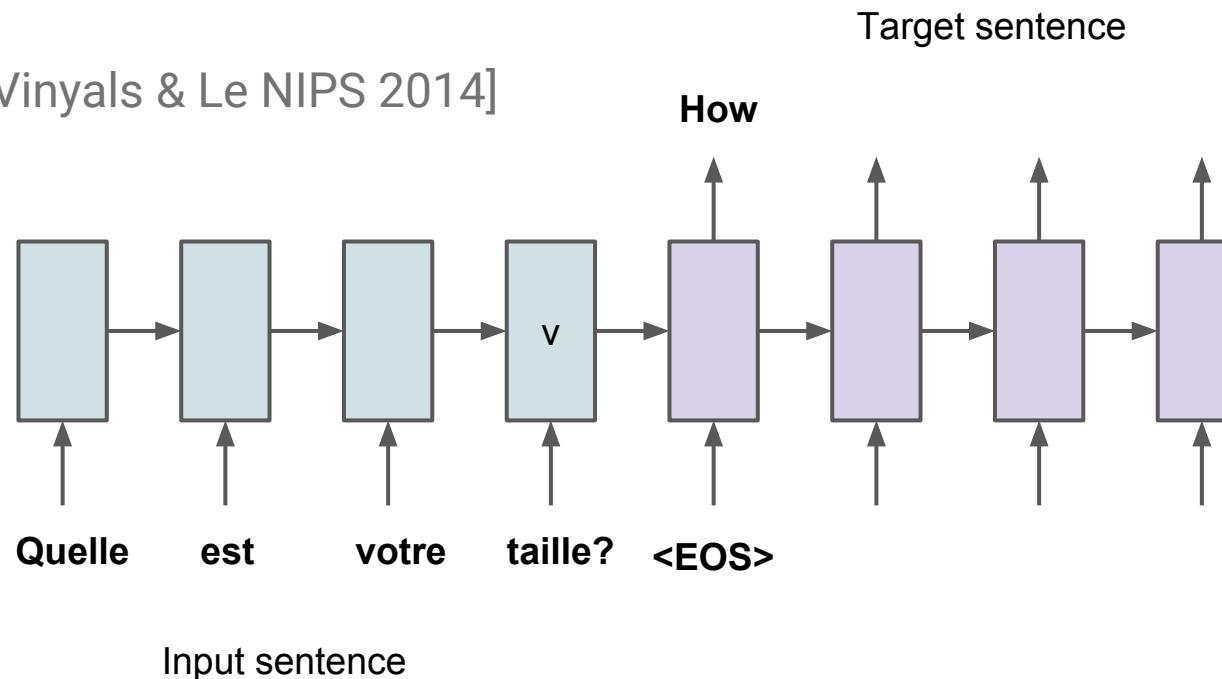
[Sutskever & Vinyals & Le NIPS 2014]



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

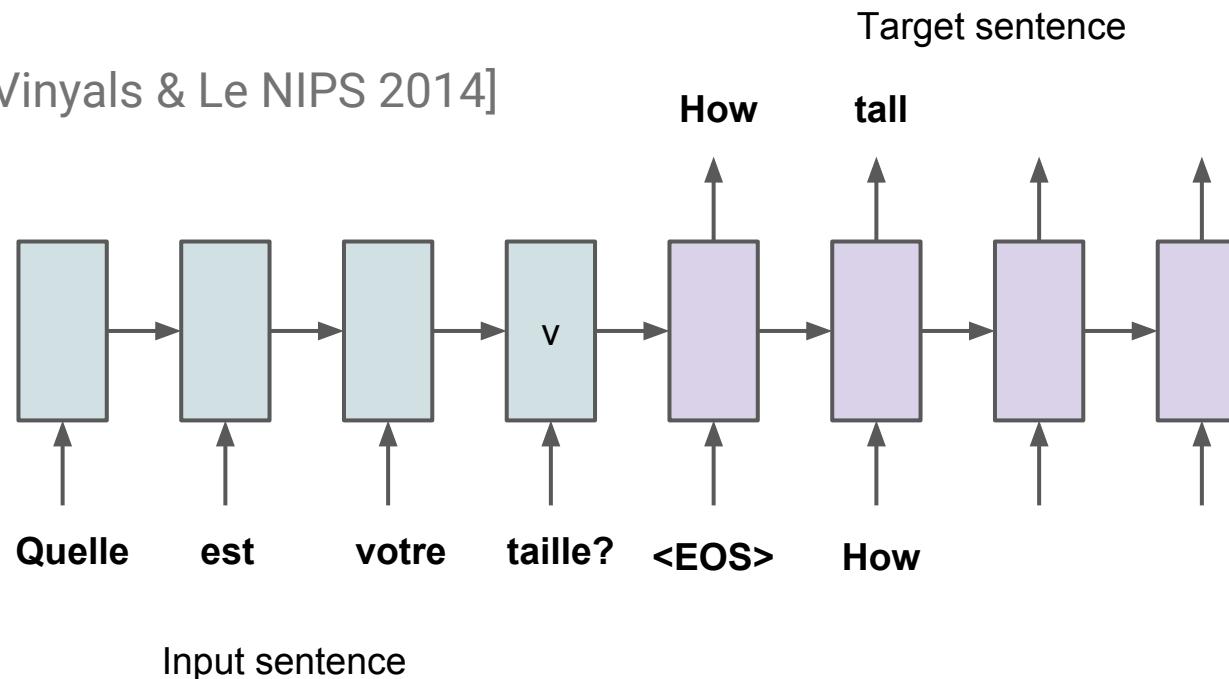
Sequence-to-Sequence Model: Machine Translation

[Sutskever & Vinyals & Le NIPS 2014]



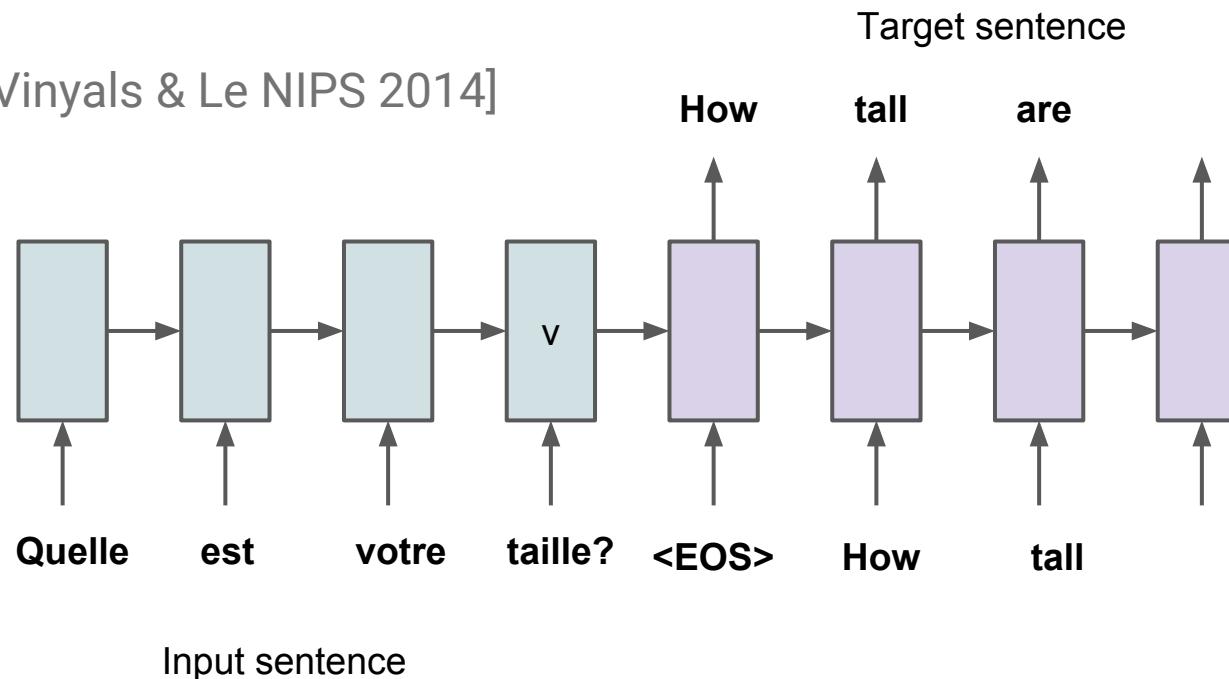
Sequence-to-Sequence Model: Machine Translation

[Sutskever & Vinyals & Le NIPS 2014]



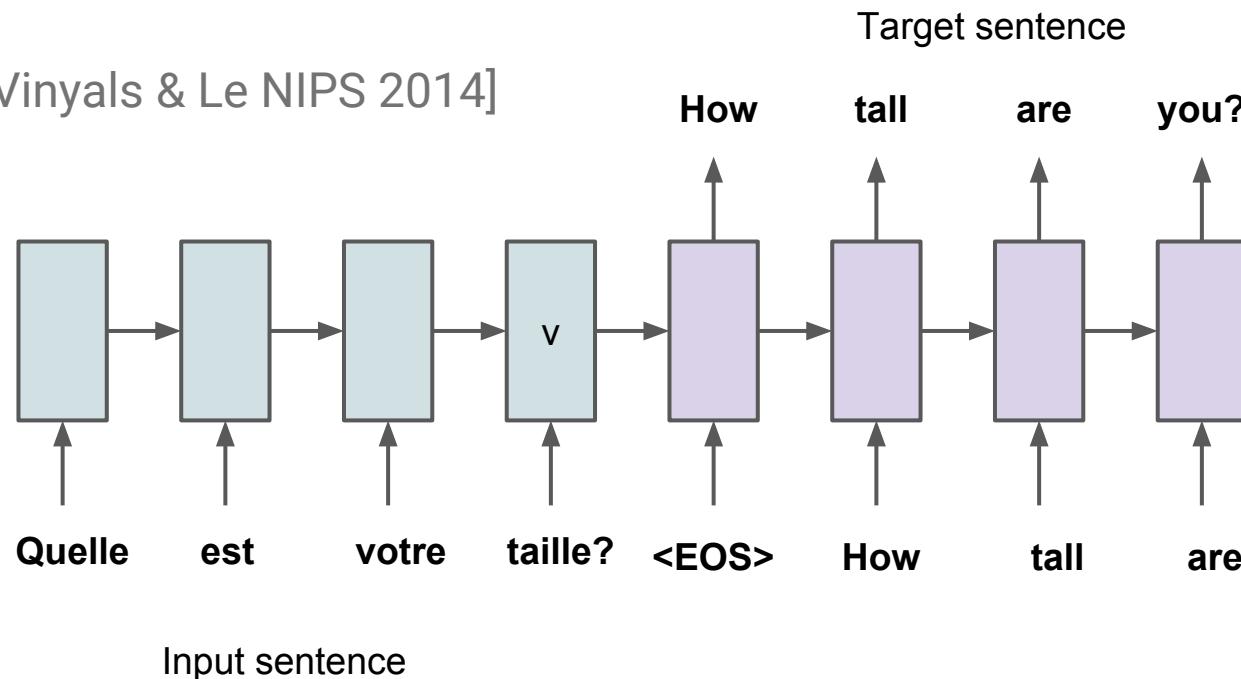
Sequence-to-Sequence Model: Machine Translation

[Sutskever & Vinyals & Le NIPS 2014]



Sequence-to-Sequence Model: Machine Translation

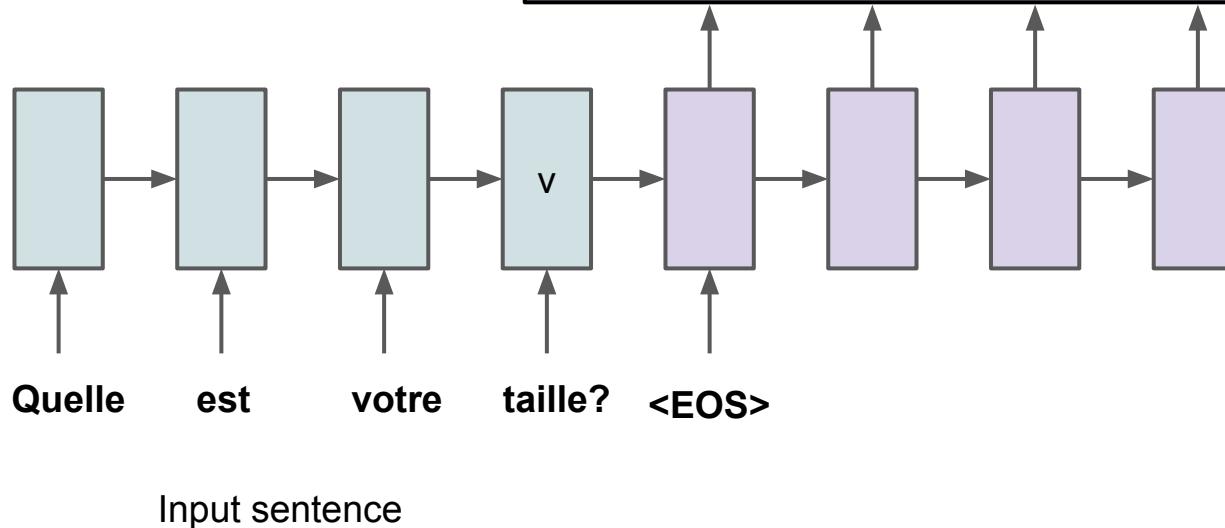
[Sutskever & Vinyals & Le NIPS 2014]



Sequence-to-Sequence Model: Machine Translation

[Sutskever & Vinyals & Le NIPS 2014]

**At inference time:
Beam search to choose most probable
over possible output sequences**



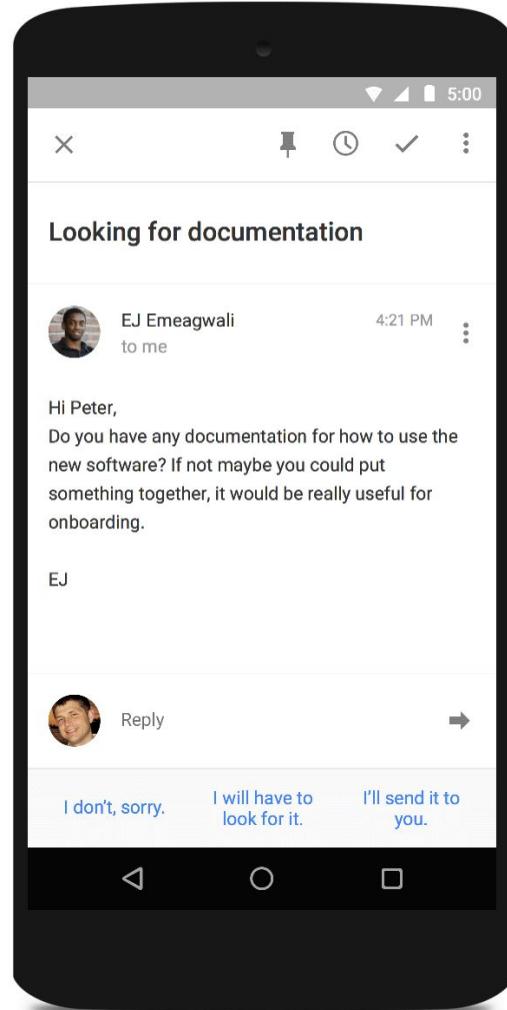


Smart Reply

April 1, 2009: April Fool's Day joke

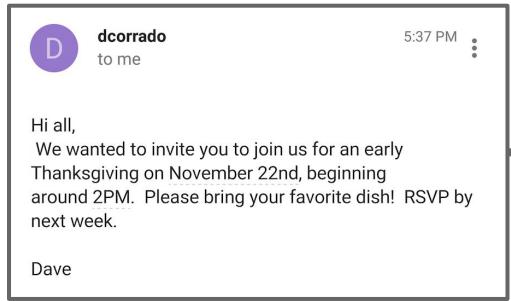
Nov 5, 2015: Launched Real Product

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



Incoming Email

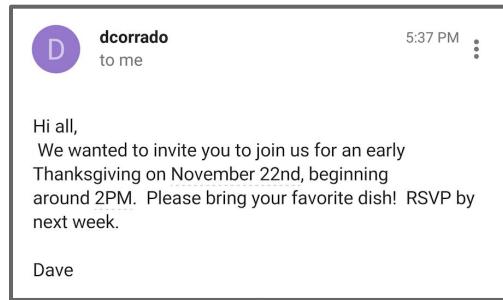
Smart Reply



Activate
Smart Reply?
yes/no

Smart Reply

Incoming Email



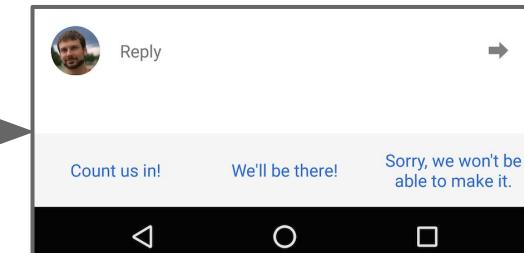
Small
Feed-Forward
Neural Network

Activate
Smart Reply?
yes/no



Deep Recurrent
Neural Network

Generated Replies



Combining vision and language

Image Captioning

[Vinyals et al., CVPR 2015]

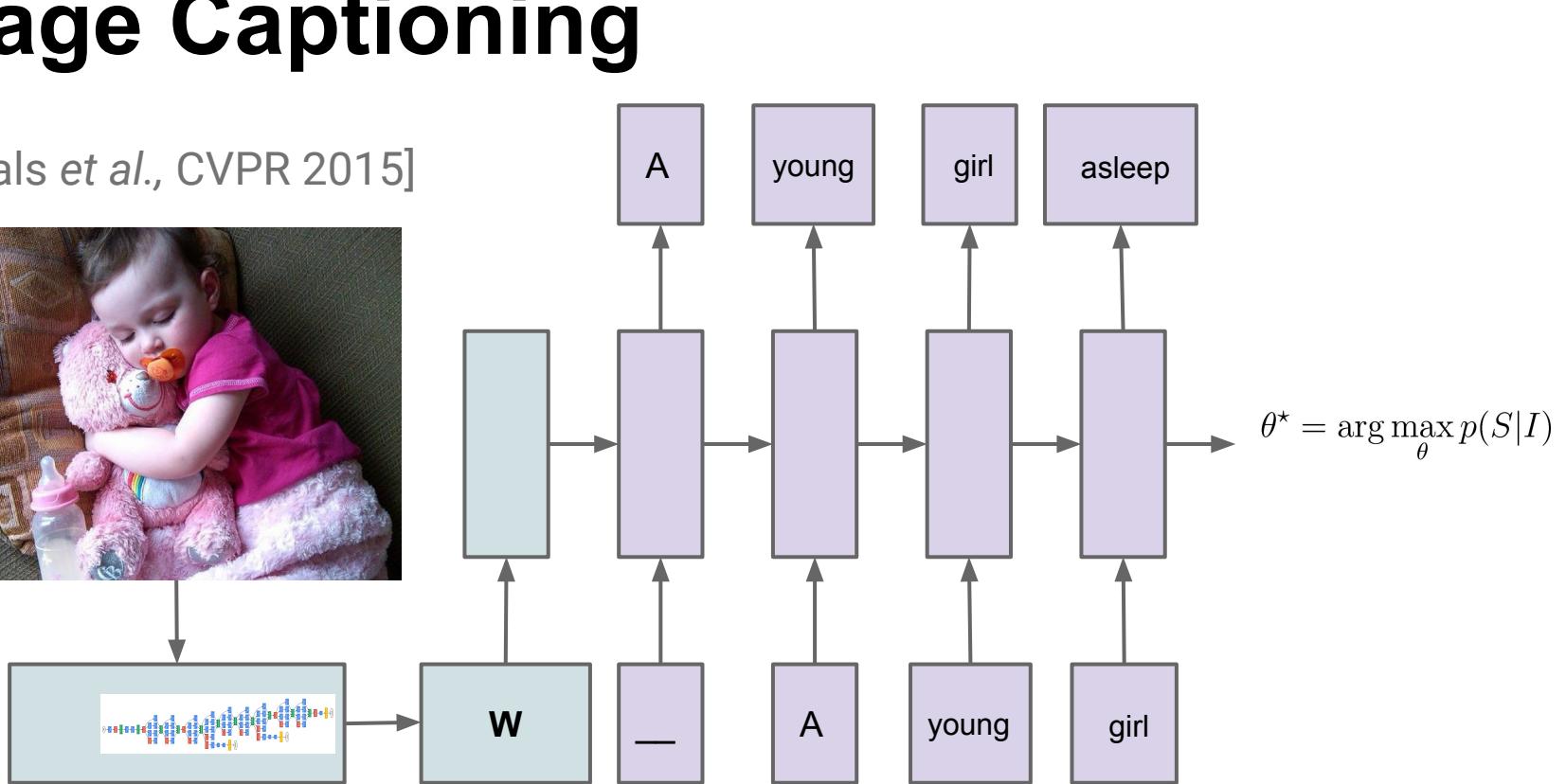


Image Captions Research



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.

Translation as a sign of
better language understanding

Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi
`yonghui,schuster,zhifengc,qvl,mnorouzi@google.com`

Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey,
Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser,
Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens,
George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa,
Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, Jeffrey Dean

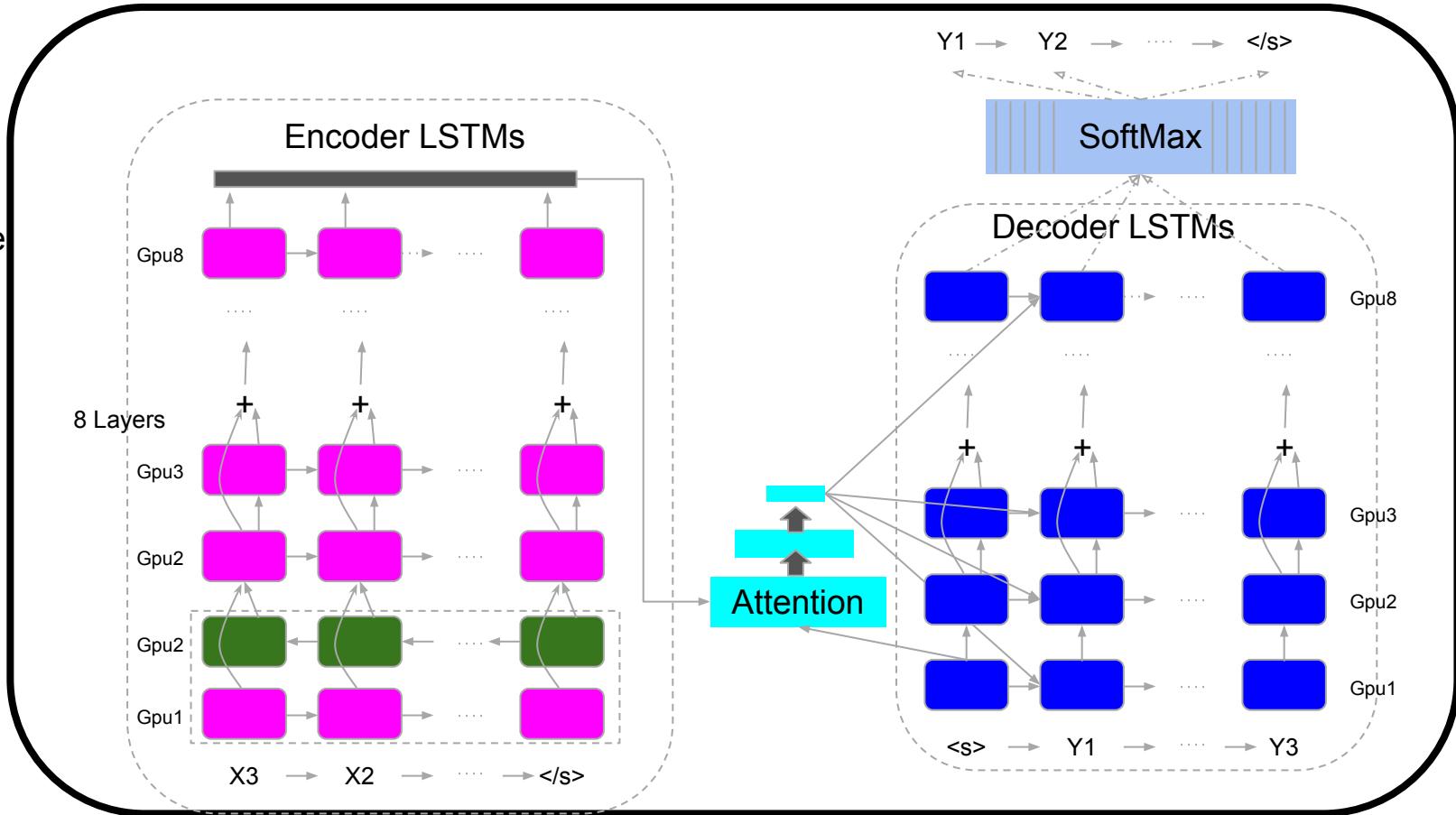
<https://arxiv.org/abs/1609.08144>

Great quality improvements

...but challenging scalability issues

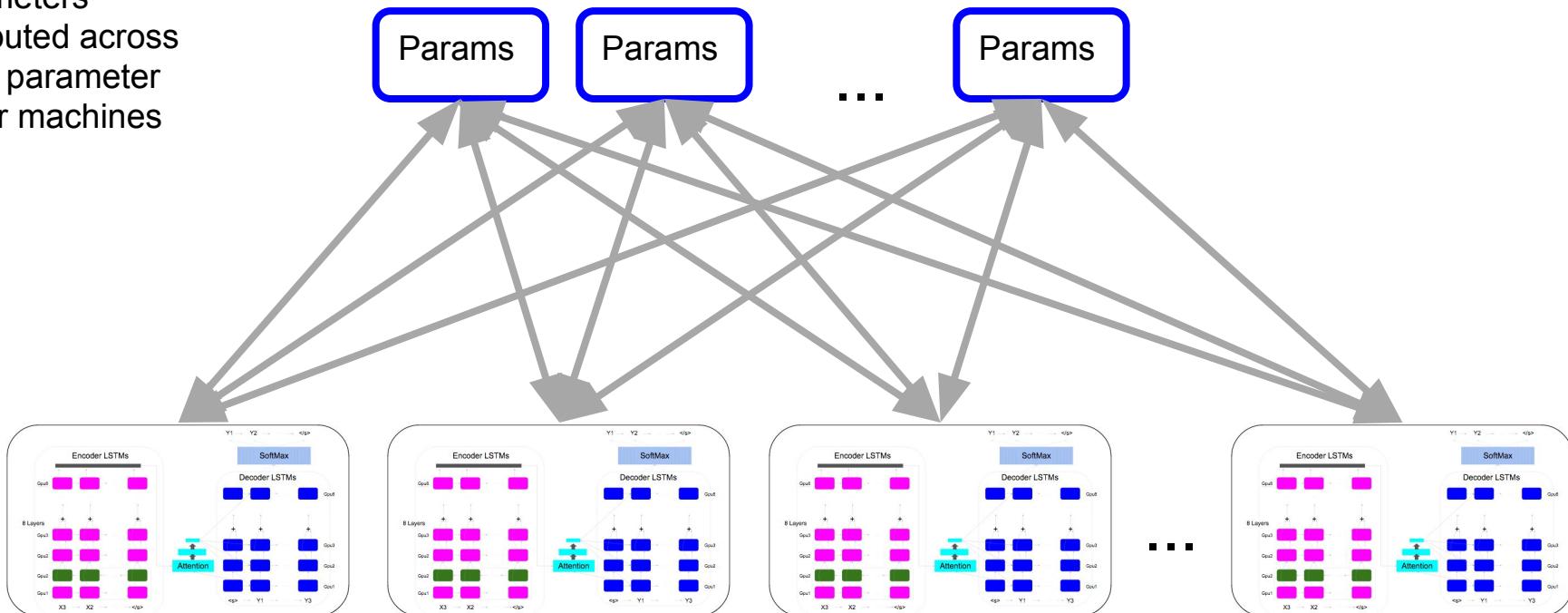
Google Neural Machine Translation Model

One model replica:
one machine
w/ 8 GPUs

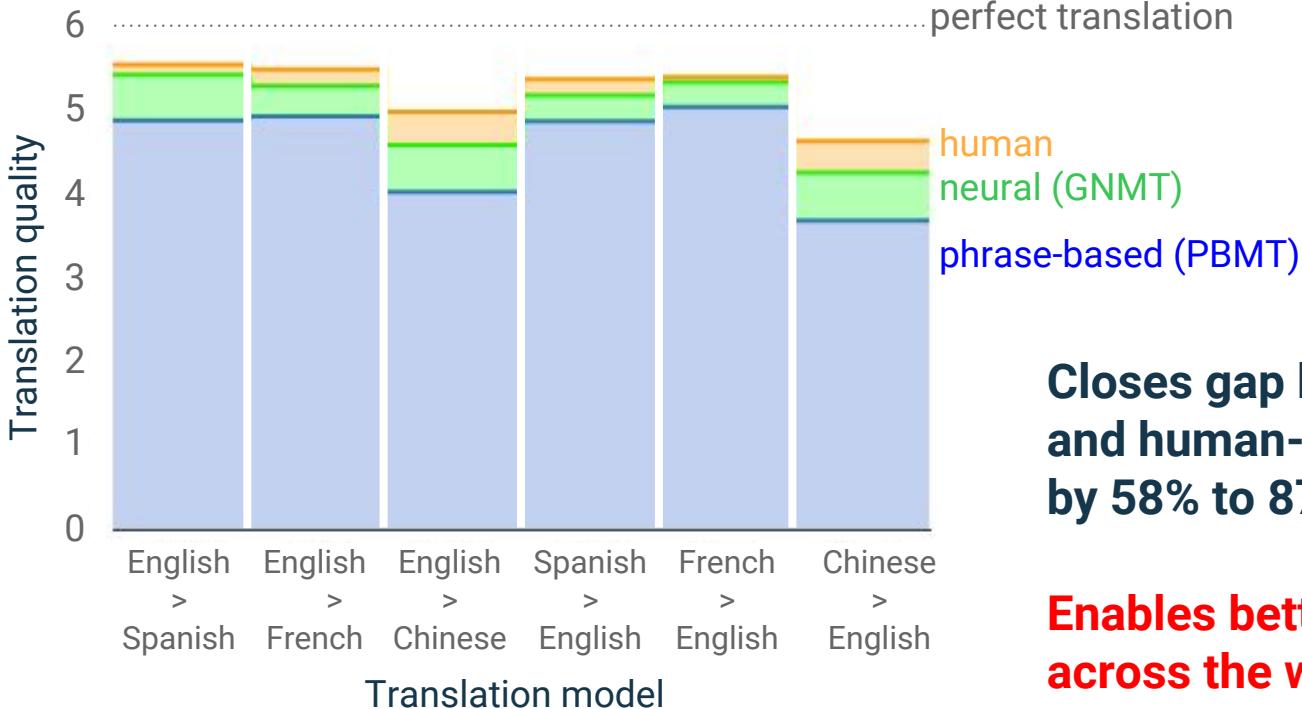


Model + Data Parallelism

Parameters
distributed across
many parameter
server machines



Neural Machine Translation



**Closes gap between old system
and human-quality translation
by 58% to 87%**

**Enables better communication
across the world**

More widespread use of:
Transfer and multi-task learning, zero-shot learning

Currently:

Most models are trained from scratch for a single task

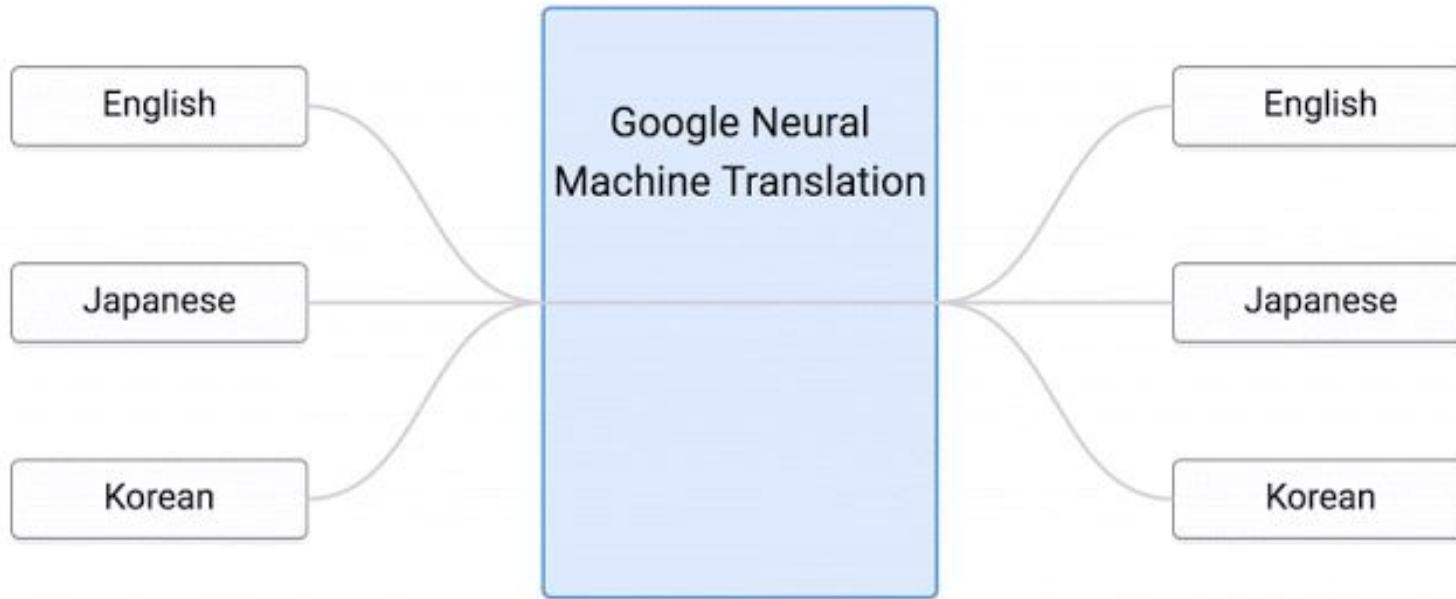
This is quite inefficient:

Data inefficient: needs lots of data for each task

Computation inefficient: starting from scratch is a lot of work

Human ML expert inefficient: substantial effort required for
each task

Training

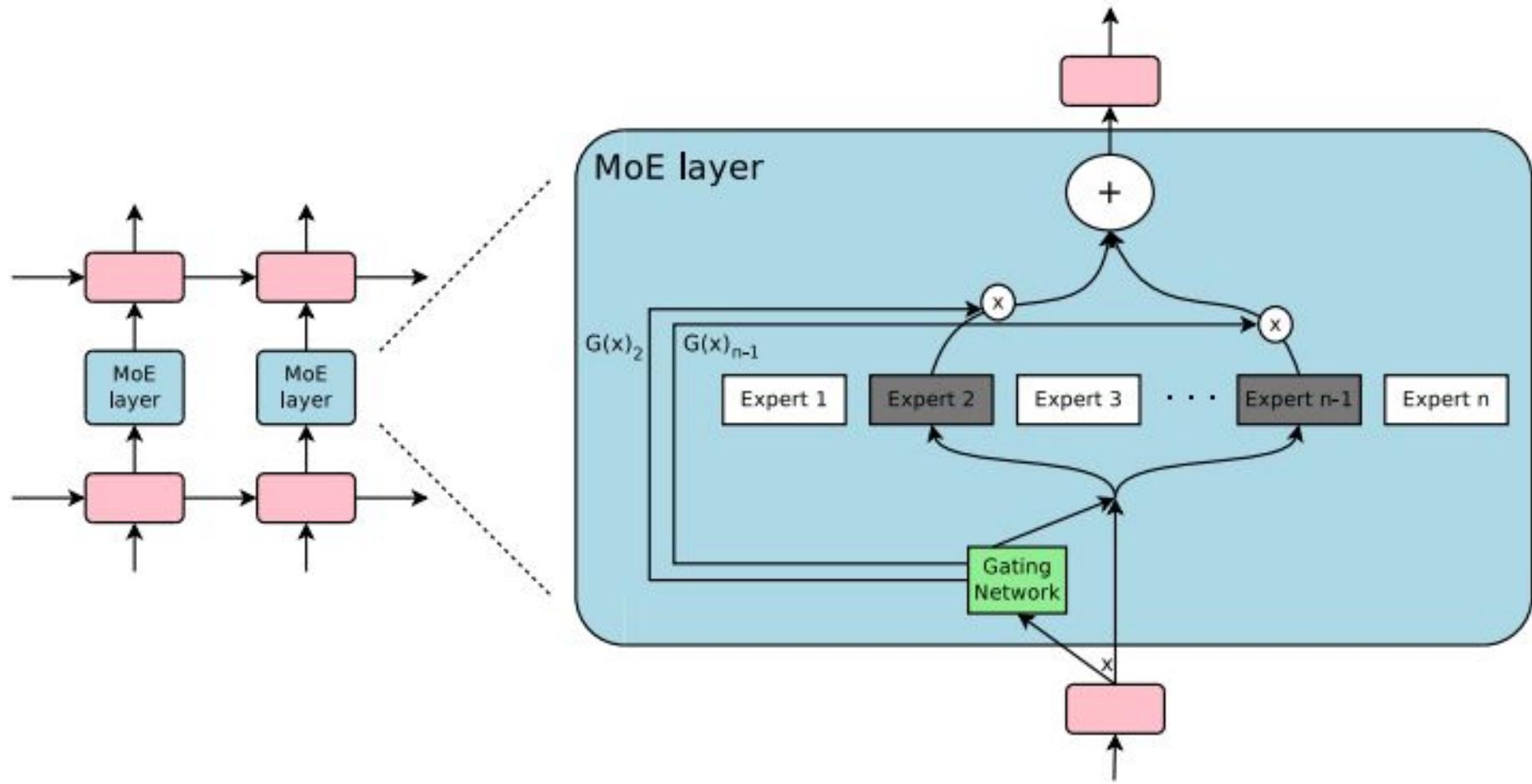


Google's Multilingual Neural Machine Translation System: Enabling Zero-Shot Translation,
Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat,
Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean
<https://arxiv.org/abs/1611.04558>

<https://research.googleblog.com/2016/11/zero-shot-translation-with-googles.html>

Bigger models, but sparsely activated

Per-Example Routing



Per-Example Routing

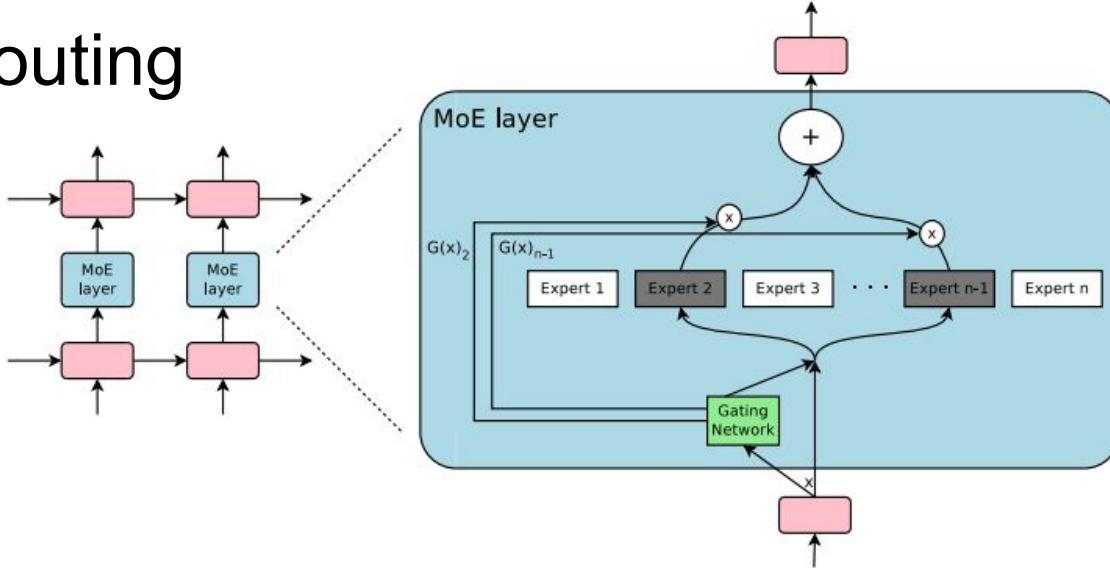


Table 7: Perplexity and BLEU comparison of our method against previous state-of-art methods on the Google Production En \rightarrow Fr dataset.

Model	Eval Perplexity	Eval BLEU	Test Perplexity	Test BLEU	Computation per Word	Total #Parameters	Training Time
MoE with 2048 Experts	2.60	37.27	2.69	36.57	100.8M	8.690B	1 day/64 k40s
GNMT (Wu et al., 2016)	2.78	35.80	2.87	35.56	214.2M	246.9M	6 days/96 k80s

Outrageously Large Neural Networks: The Sparsely-gated Mixture-of-Experts Layer,
Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le & Jeff Dean
Submitted to ICLR 2017, <https://openreview.net/pdf?id=B1ckMDqlg>

Automated machine learning (“learning to learn”)

Current:

Solution = ML expertise + data + computation

Can we turn this into:

Solution = data + 100X computation

???

NEURAL ARCHITECTURE SEARCH WITH REINFORCEMENT LEARNING

Barret Zoph,* Quoc V. Le

Google Brain

{barrettzoph, qvl}@google.com

Idea: model-generating model trained via RL

- (1) Generate ten models
- (2) Train them for a few hours
- (3) Use loss of the generated models as reinforcement learning signal

CIFAR-10 Image Recognition Task

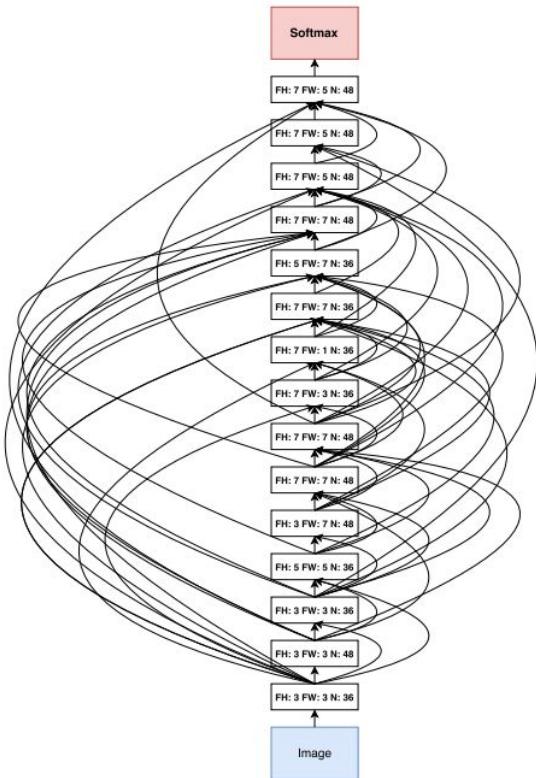


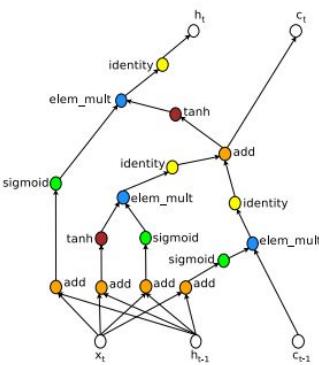
Figure 7: Convolutional architecture discovered by our method, when the search space does not have strides or pooling layers. FH is filter height, FW is filter width and N is number of filters.

Model	Depth	Parameters	Error rate (%)
Network in Network (Lin et al., 2013)	-	-	8.81
All-CNN (Springenberg et al., 2014)	-	-	7.25
Deeply Supervised Net (Lee et al., 2015)	-	-	7.97
Highway Network (Srivastava et al., 2015)	-	-	7.72
Scalable Bayesian Optimization (Snoek et al., 2015)	-	-	6.37
FractalNet (Larsson et al., 2016) with Dropout/Drop-path	21 21	38.6M 38.6M	5.22 4.60
ResNet (He et al., 2016a)	110	1.7M	6.61
ResNet (reported by Huang et al. (2016b))	110	1.7M	6.41
ResNet with Stochastic Depth (Huang et al., 2016b)	110 1202	1.7M 10.2M	5.23 4.91
Wide ResNet (Zagoruyko & Komodakis, 2016)	16 28	11.0M 36.5M	4.81 4.17
ResNet (pre-activation) (He et al., 2016b)	164 1001	1.7M 10.2M	5.46 4.62
DenseNet ($L = 40, k = 12$) Huang et al. (2016a)	40	1.0M	5.24
DenseNet($L = 100, k = 12$) Huang et al. (2016a)	100	7.0M	4.10
DenseNet ($L = 100, k = 24$) Huang et al. (2016a)	100	27.2M	3.74
Neural Architecture Search v1 no stride or pooling	15	4.2M	5.50
Neural Architecture Search v2 predicting strides	20	2.5M	6.01
Neural Architecture Search v3 max pooling	39	7.1M	4.47
Neural Architecture Search v3 max pooling + more filters	39	32.0M	3.84

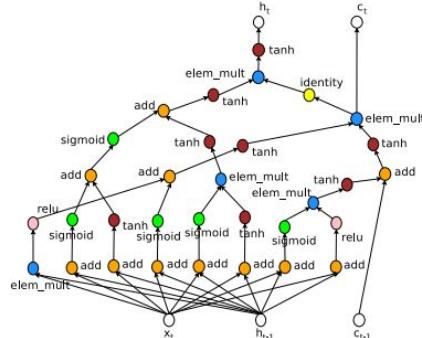
Table 1: Performance of Neural Architecture Search and other state-of-the-art models on CIFAR-10.

Penn Tree Bank Language Modeling Task

“Normal” LSTM cell



Cell discovered by
architecture search



Model	Parameters	Test Perplexity
Mikolov & Zweig (2012) - KN-5	2M [‡]	141.2
Mikolov & Zweig (2012) - KN5 + cache	2M [‡]	125.7
Mikolov & Zweig (2012) - RNN	6M [‡]	124.7
Mikolov & Zweig (2012) - RNN-LDA	7M [‡]	113.7
Mikolov & Zweig (2012) - RNN-LDA + KN-5 + cache	9M [‡]	92.0
Pascanu et al. (2013) - Deep RNN	6M	107.5
Cheng et al. (2014) - Sum-Prod Net	5M [‡]	100.0
Zaremba et al. (2014) - LSTM (medium)	20M	82.7
Zaremba et al. (2014) - LSTM (large)	66M	78.4
Gal (2015) - Variational LSTM (medium, untied)	20M	79.7
Gal (2015) - Variational LSTM (medium, untied, MC)	20M	78.6
Gal (2015) - Variational LSTM (large, untied)	66M	75.2
Gal (2015) - Variational LSTM (large, untied, MC)	66M	73.4
Kim et al. (2015) - CharCNN	19M	78.9
Press & Wolf (2016) - Variational LSTM, shared embeddings	24M	73.2
Merity et al. (2016) - Zoneout + Variational LSTM (medium)	20M	80.6
Merity et al. (2016) - Pointer Sentinel-LSTM (medium)	21M	70.9
Zilly et al. (2016) - Variational RHN, shared embeddings	24M	66.0
Neural Architecture Search with base 8	32M	67.9
Neural Architecture Search with base 8 and shared embeddings	25M	64.0
Neural Architecture Search with base 8 and shared embeddings	54M	62.4

Table 2: Single model perplexity on the test set of the Penn Treebank language modeling task. Parameter numbers with [‡] are estimates with reference to Merity et al. (2016).

More computational power needed

Deep learning is transforming how we design computers

Special computation properties

reduced
precision
ok

$$\begin{array}{r} \text{about 1.2} \\ \times \text{ about 0.6} \\ \hline \text{about 0.7} \end{array}$$

NOT

$$\begin{array}{r} 1.21042 \\ \times 0.61127 \\ \hline 0.73989543 \end{array}$$

Special computation properties

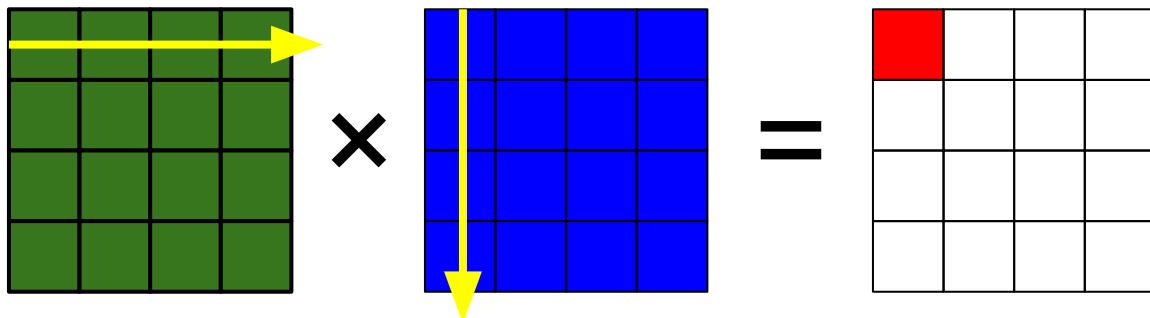
reduced
precision
ok

$$\begin{array}{r} \text{about 1.2} \\ \times \text{ about 0.6} \\ \hline \text{about 0.7} \end{array}$$

NOT

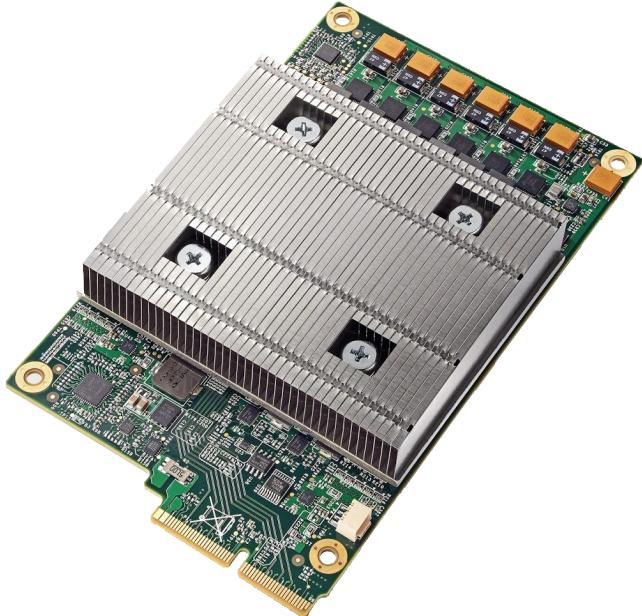
$$\begin{array}{r} 1.21042 \\ \times 0.61127 \\ \hline 0.73989343 \end{array}$$

handful of
specific
operations



Tensor Processing Unit

Custom Google-designed chip for neural net computations



In production use for >20 months: used on every search query, for neural machine translation, for AlphaGo match, ...



Example queries of the future

Which of these eye images shows symptoms of diabetic retinopathy?

Please fetch me a cup of tea from the kitchen

Describe this video in Spanish

Find me documents related to reinforcement learning for robotics and summarize them in German

Conclusions

Deep neural networks are making significant strides in speech, vision, language, search, robotics, healthcare, ...

If you're not considering how to use deep neural nets to solve some problems, **you almost certainly should be**



More info about our work

Main Research Areas

[Machine Learning Algorithms and Techniques](#)

[Healthcare](#)

[Computer Systems for Machine Learning](#)

[Robotics](#)

[Natural Language Understanding](#)

[Music and Art Generation](#)

[Perception](#)

[MORE PAPERS](#)

[BLOG POSTS](#)

Join the Team

Full Time Roles

We're looking for talented research scientists and software engineers enthusiastic about deep learning to join us.

[VIEW JOBS](#)

Brain Residency

This 12-month program is designed to jumpstart your career in deep learning, working with our scientists and engineers from the Google Brain Team.

[VIEW RESIDENCY](#)

Visiting Faculty

Visiting Faculty work closely with our scientists and engineers, and have the opportunity to explore projects at industrial scale with state-of-the-art technology.

[VIEW VISITING FACULTY](#)

Interns

Our interns work on projects utilizing the latest techniques in deep learning. In your application, indicate your research interests in the 'Cover letter/other notes' section, so it can be routed to the appropriate recruiter.

[VIEW INTERNSHIPS](#)

Thanks!