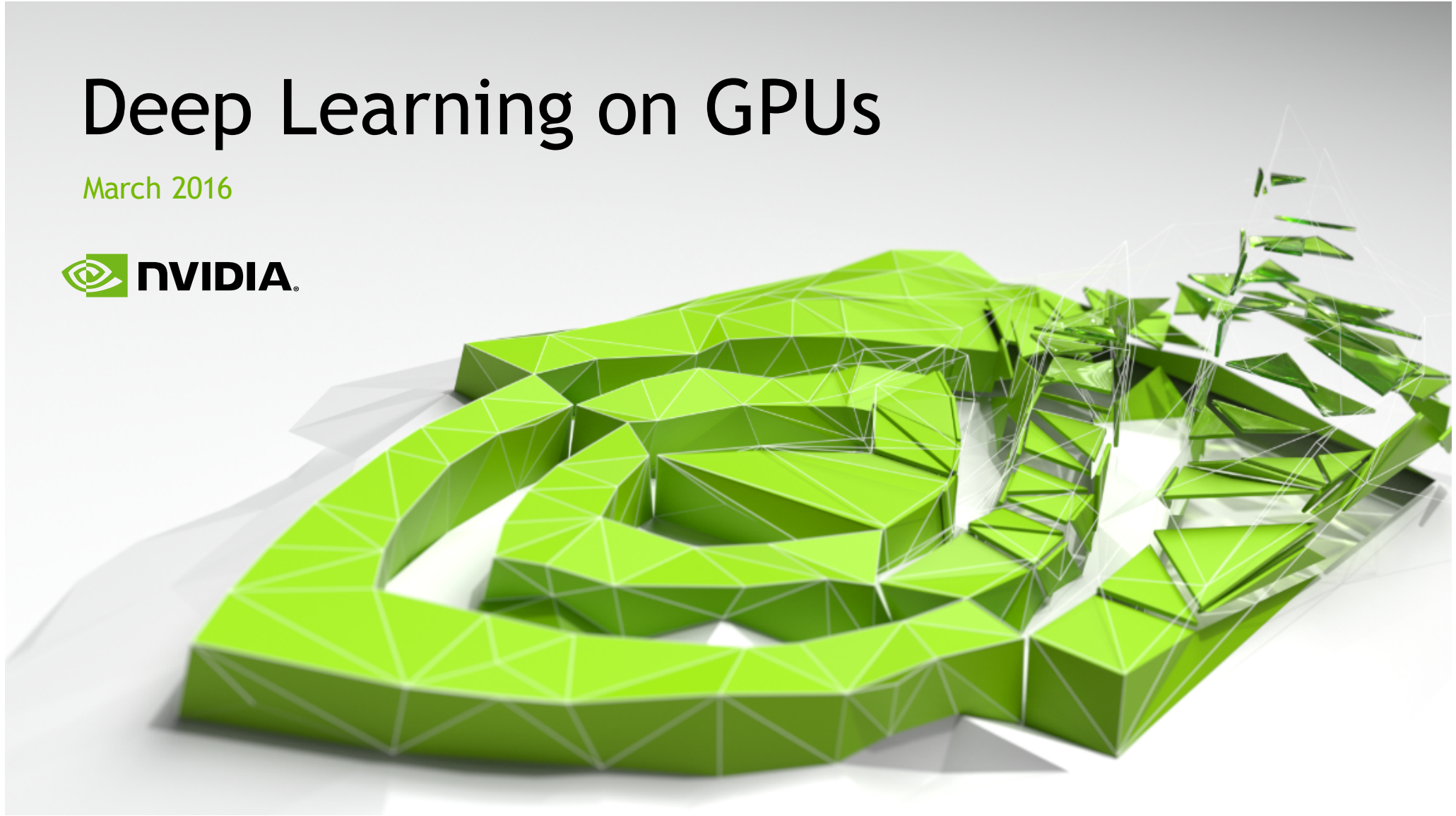


Deep Learning on GPUs

March 2016



AGENDA

What is Deep Learning?

GPUs and DL

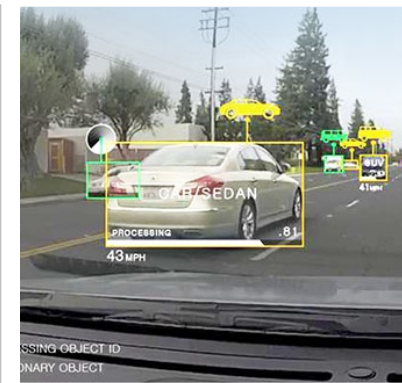
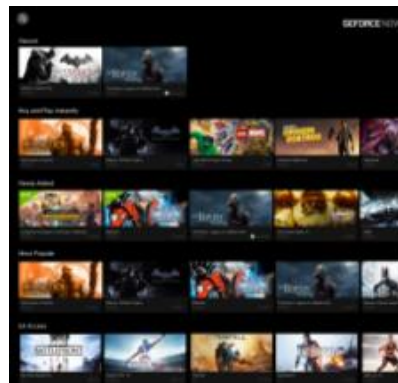
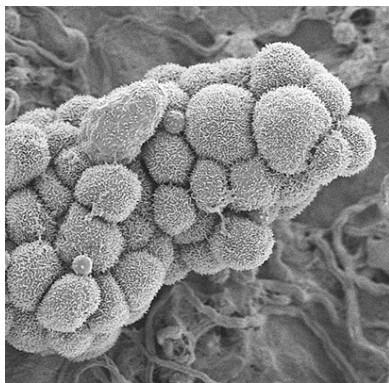
DL in practice

Scaling up DL

What is Deep Learning?



DEEP LEARNING EVERYWHERE



INTERNET & CLOUD

Image Classification
Speech Recognition
Language Translation
Language Processing
Sentiment Analysis
Recommendation

MEDICINE & BIOLOGY

Cancer Cell Detection
Diabetic Grading
Drug Discovery

MEDIA & ENTERTAINMENT

Video Captioning
Video Search
Real Time Translation

SECURITY & DEFENSE

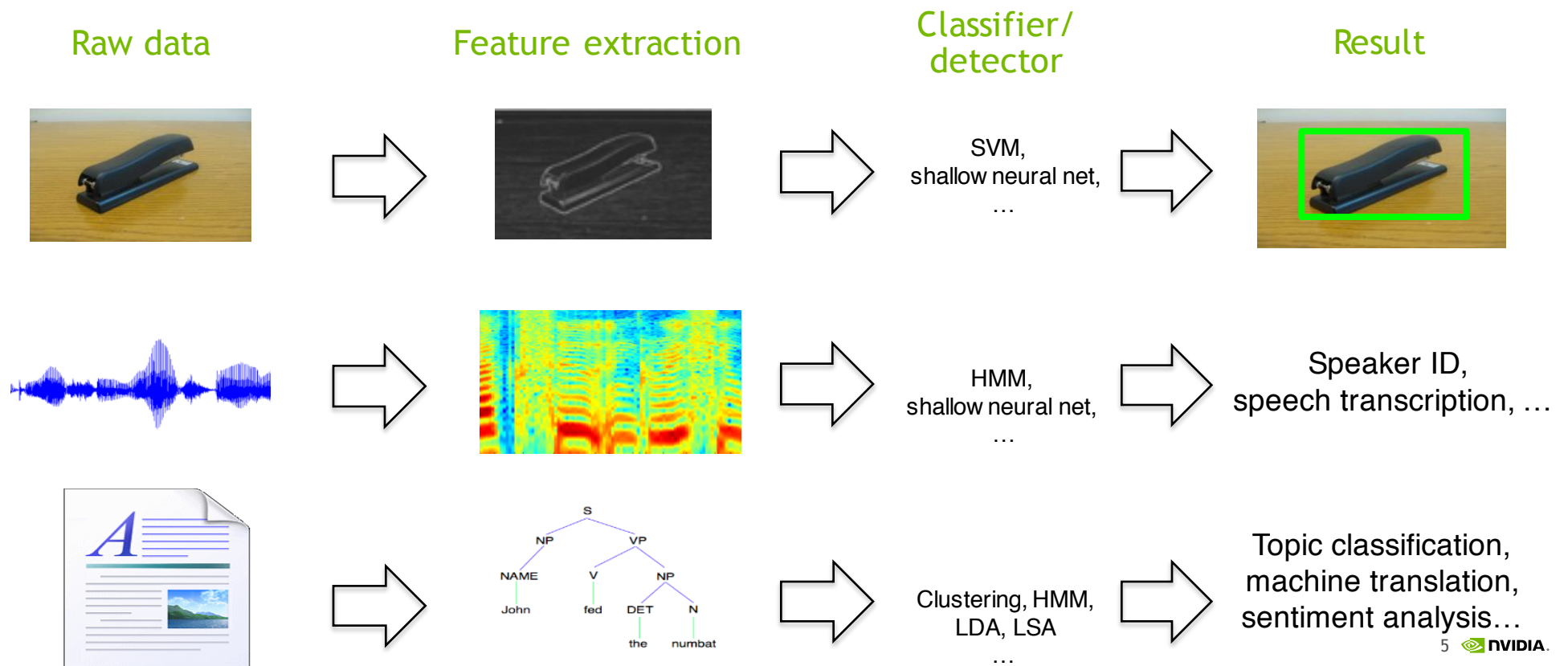
Face Detection
Video Surveillance
Satellite Imagery

AUTONOMOUS MACHINES

Pedestrian Detection
Lane Tracking
Recognize Traffic Sign

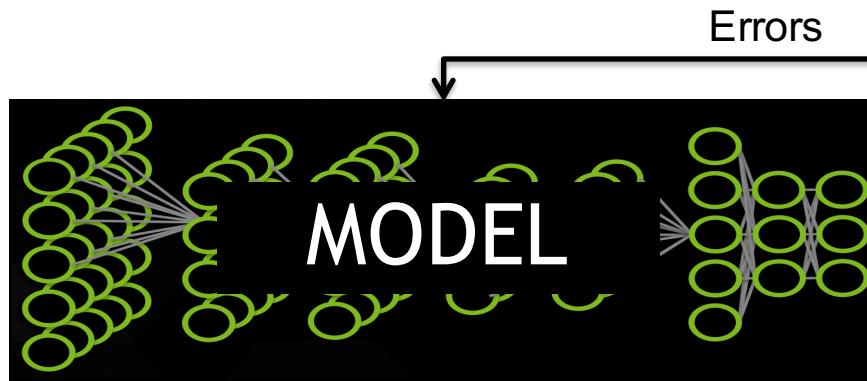
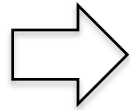
Traditional machine perception

Hand crafted feature extractors

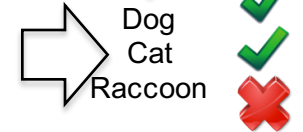


Deep learning approach

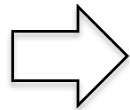
Train:



Errors

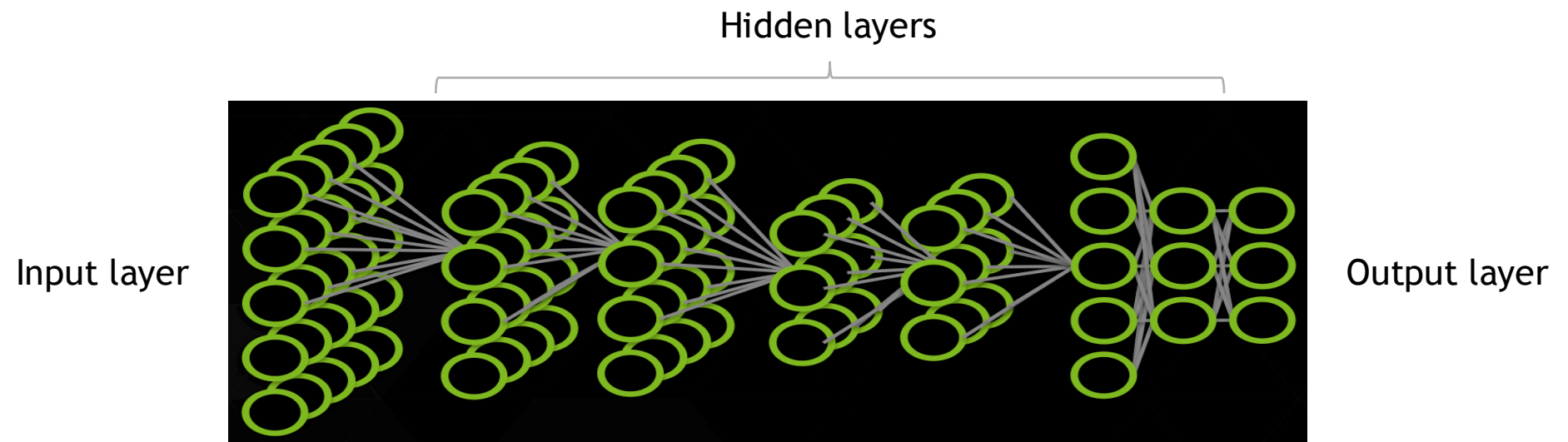


Deploy:



Artificial neural network

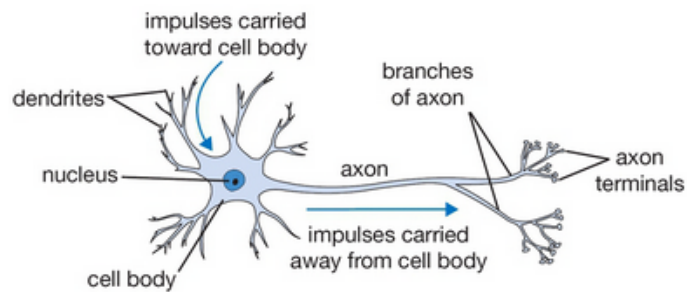
A collection of simple, trainable mathematical units that collectively learn complex functions



Given sufficient training data an artificial neural network can approximate very complex functions mapping raw data to output decisions

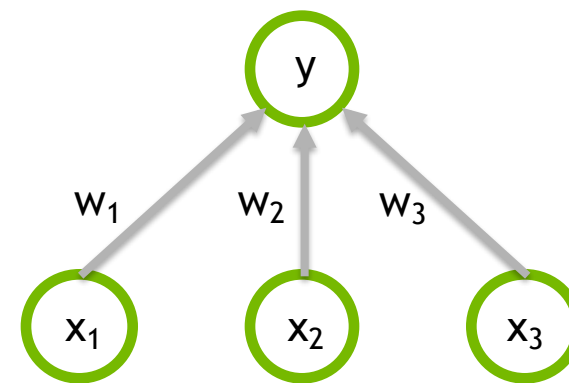
Artificial neurons

Biological neuron



From Stanford cs231n lecture notes

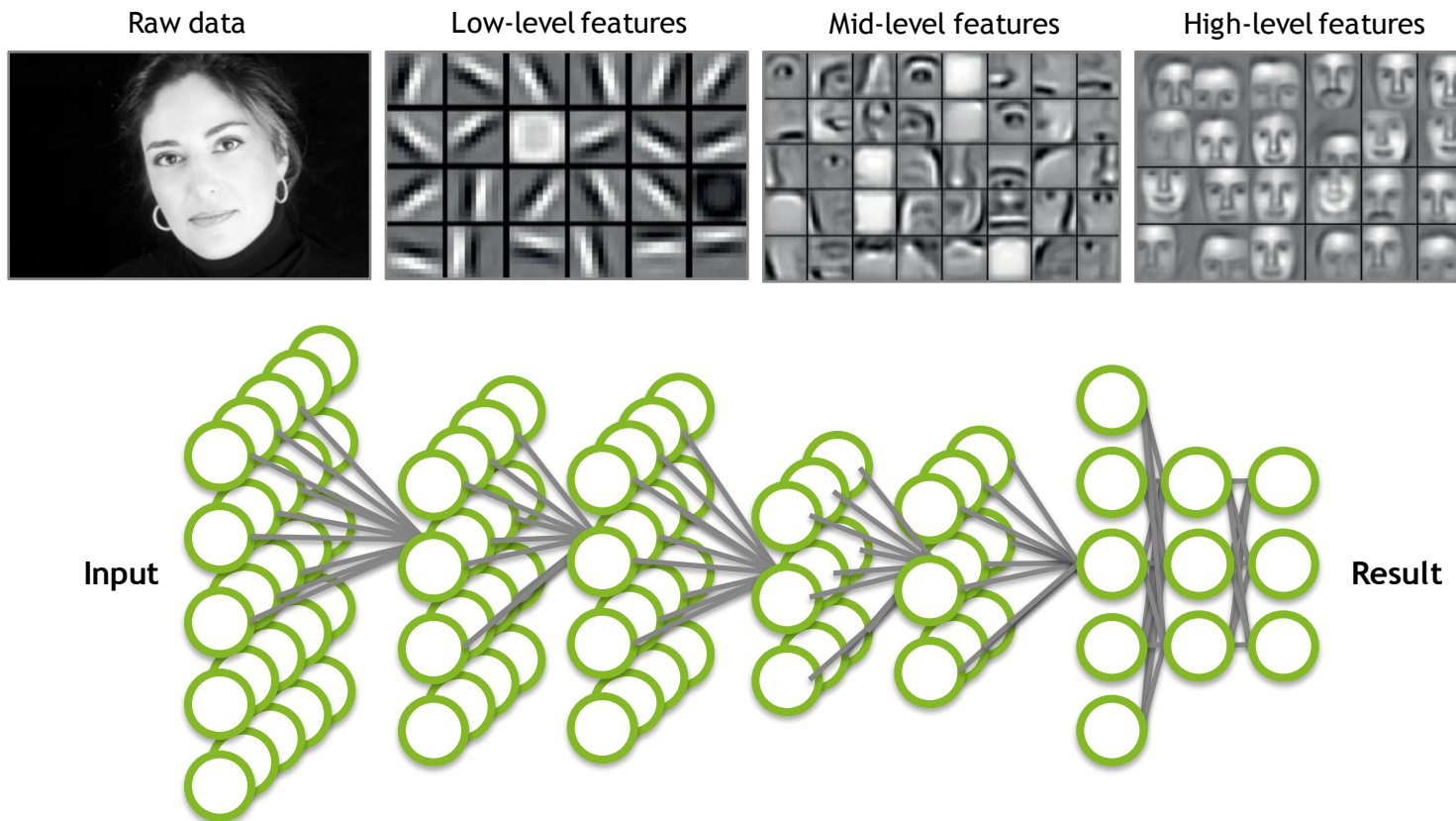
Artificial neuron



$$y = F(w_1x_1 + w_2x_2 + w_3x_3)$$

$$F(x) = \max(0, x)$$

Deep neural network (dnn)



Application components:

Task objective

e.g. Identify face

Training data

10-100M images

Network architecture

~ 10 layers

1B parameters

Learning algorithm

~30 Exaflops

~30 GPU days

Deep learning benefits

- **Robust**
 - No need to design the features ahead of time - features are automatically learned to be optimal for the task at hand
 - Robustness to natural variations in the data is automatically learned
- **Generalizable**
 - The same neural net approach can be used for many different applications and data types
- **Scalable**
 - Performance improves with more data, method is massively parallelizable

Baidu Deep Speech 2

End-to-end Deep Learning for English and Mandarin Speech Recognition



English and Mandarin speech recognition

Transition from English to Mandarin made simpler by end-to-end DL

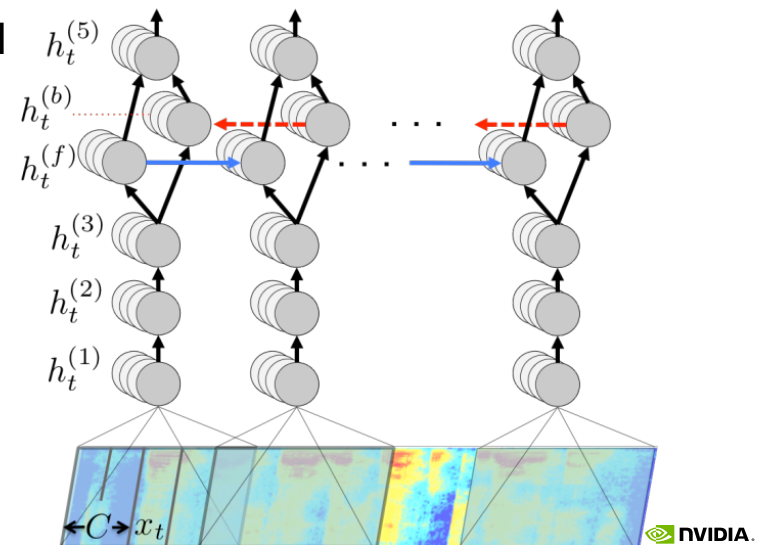
No feature engineering or Mandarin-specifics required

More accurate than humans

Error rate 3.7% vs. 4% for human tests

<http://svail.github.io/mandarin/>

<http://arxiv.org/abs/1512.02595>



AlphaGo

First Computer Program to Beat a Human Go Professional

Training DNNs: 3 weeks, 340 million training steps on 50 GPUs

Play: Asynchronous multi-threaded search

Simulations on CPUs, policy and value DNNs in parallel on GPUs

Single machine: 40 search threads, 48 CPUs, and 8 GPUs

Distributed version: 40 search threads, 1202 CPUs and 176 GPUs

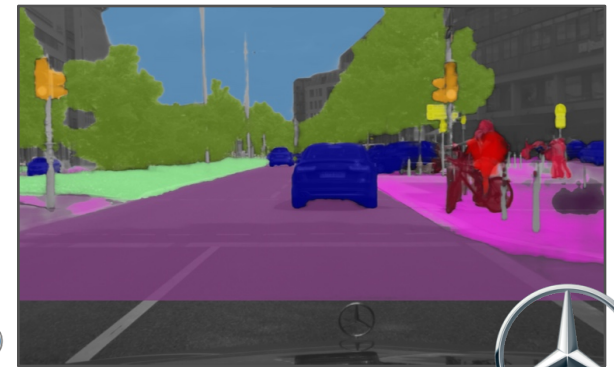
Outcome: Beat both European and World Go champions in best of 5 matches

<http://www.nature.com/nature/journal/v529/n7587/full/nature16961.html>

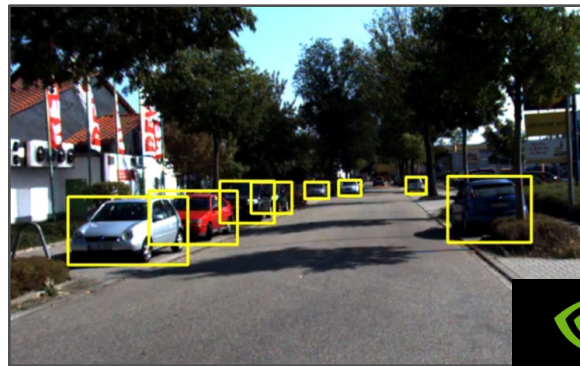
<http://deepmind.com/alpha-go.html>



Deep Learning for Autonomous vehicles



Audi



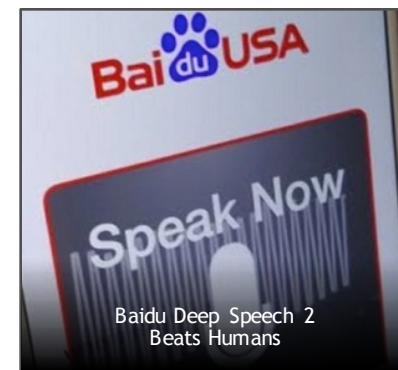
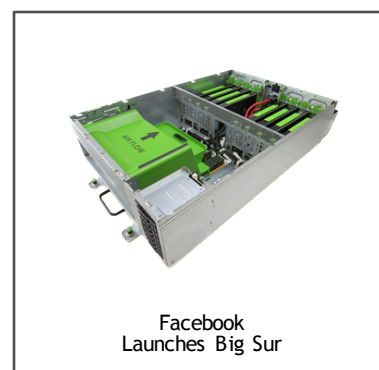
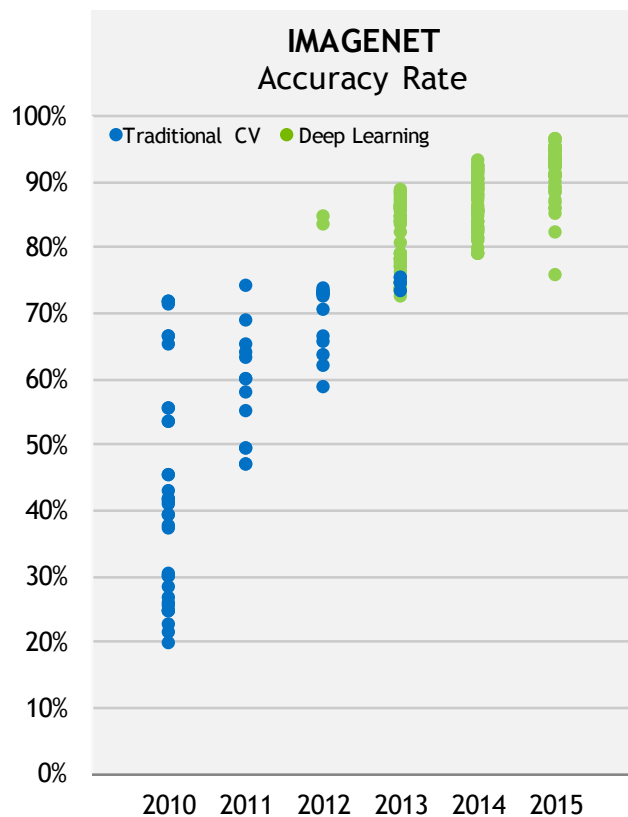
Audi

Deep Learning Synthesis

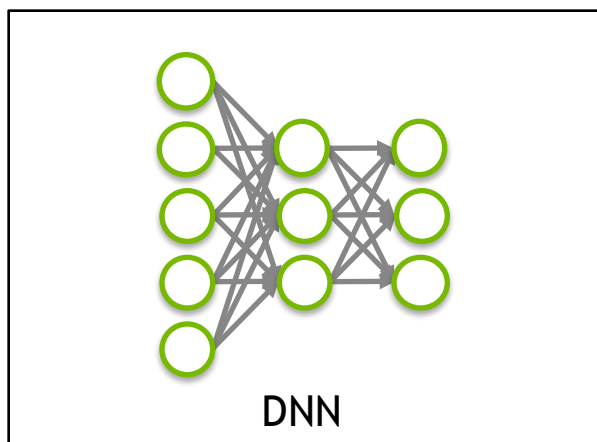


Texture synthesis and transfer using CNNs. Timo Aila et al., NVIDIA Research

THE AI RACE IS ON



The Big Bang in Machine Learning



“Google’s AI engine also reflects how the world of computer hardware is changing. (It) depends on machines equipped with GPUs... And it depends on these chips more than the larger tech universe realizes.”

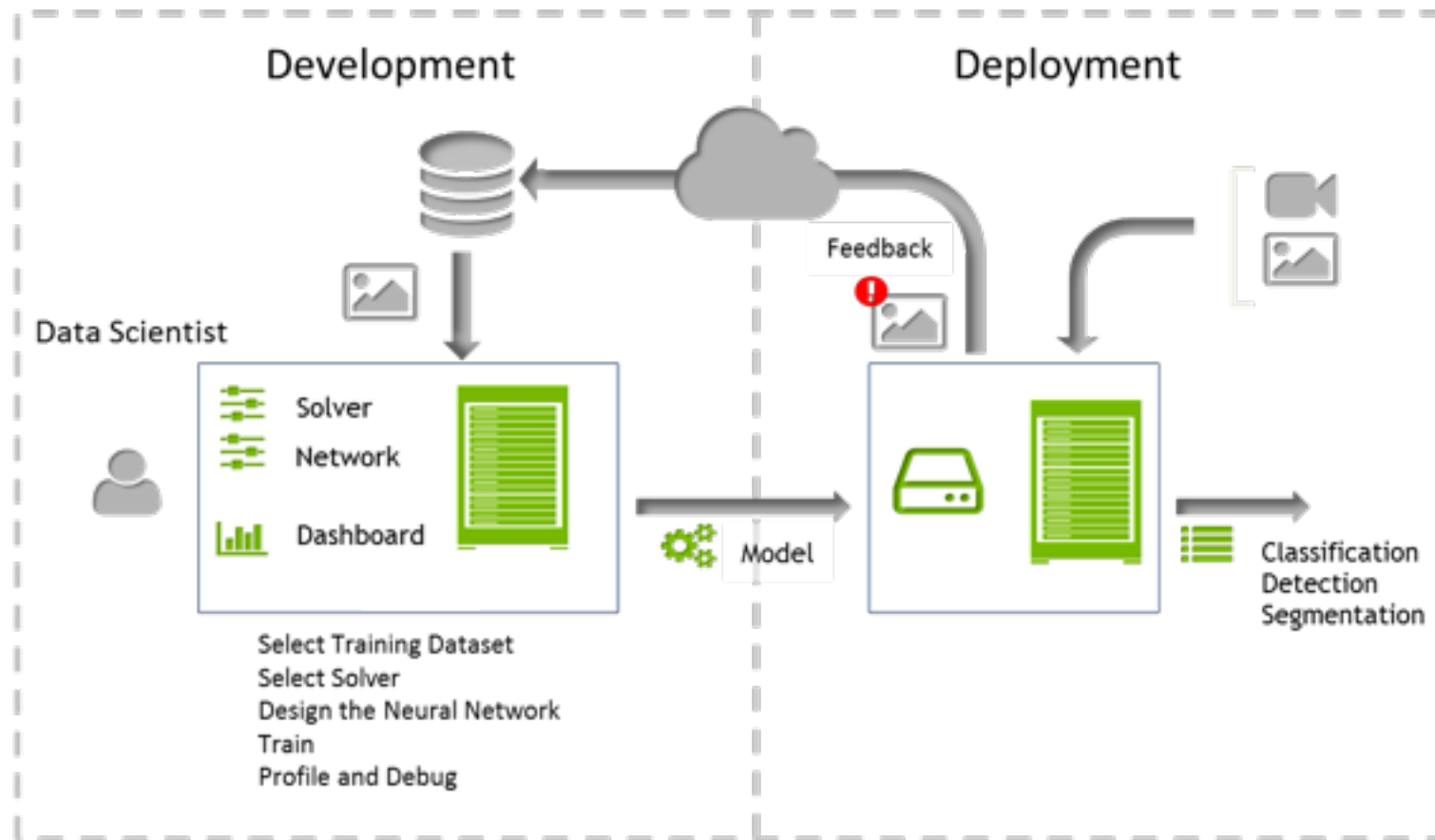
WIRED



GPUs and DL

USE MORE PROCESSORS TO GO FASTER

Deep learning development cycle

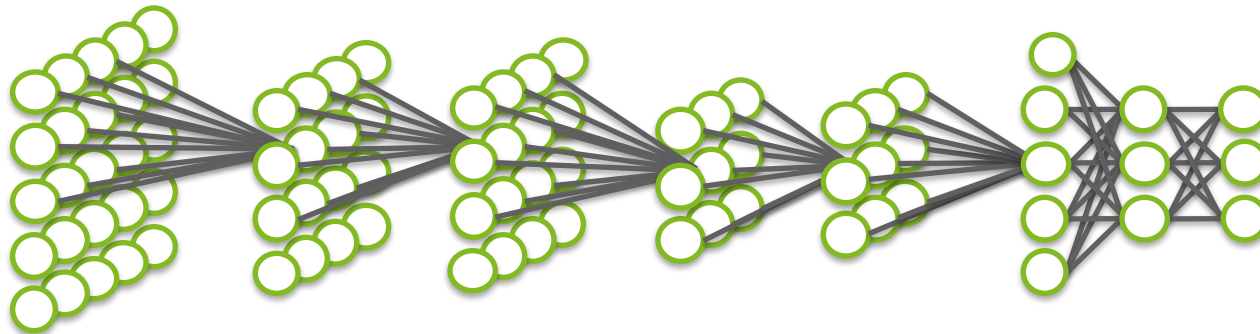


Three Kinds of Networks

DNN - all fully connected layers

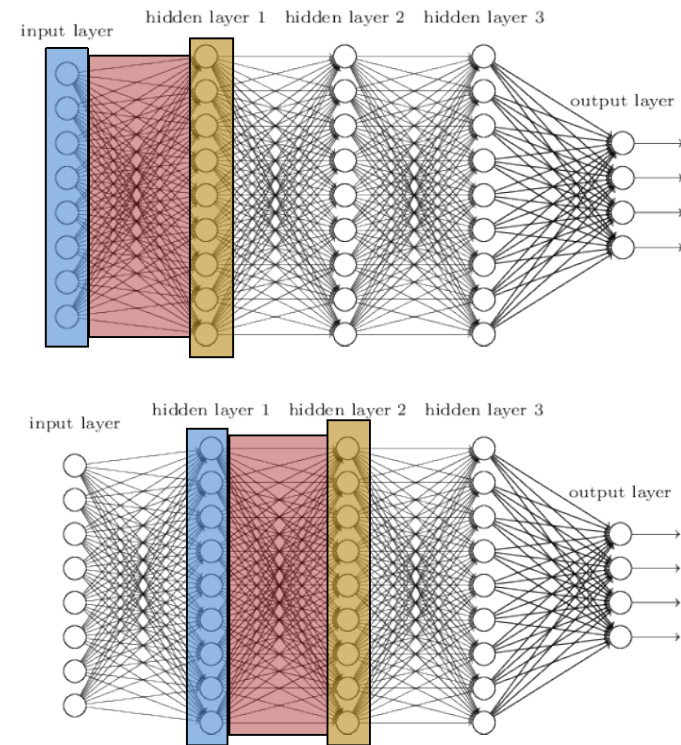
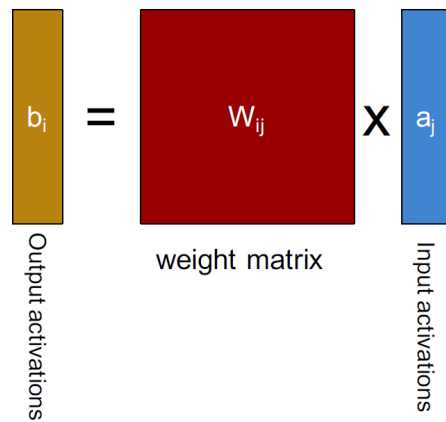
CNN - some convolutional layers

RNN - recurrent neural network, LSTM



DNN

Key operation is dense $M \times V$



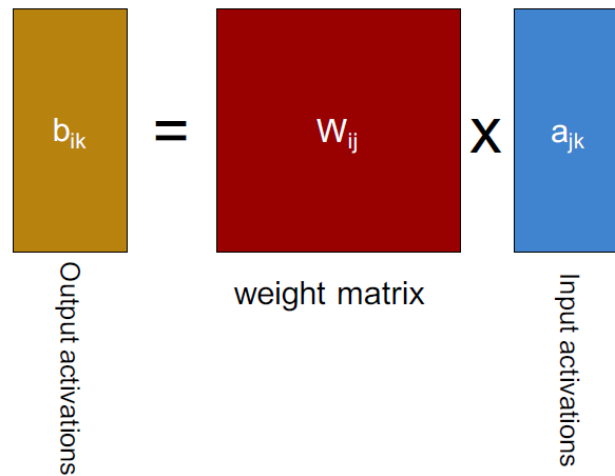
Repeat for each layer

Backpropagation uses dense matrix-matrix multiply starting from softmax scores

DNN

Batching for training and latency insensitive.

$M \times M$



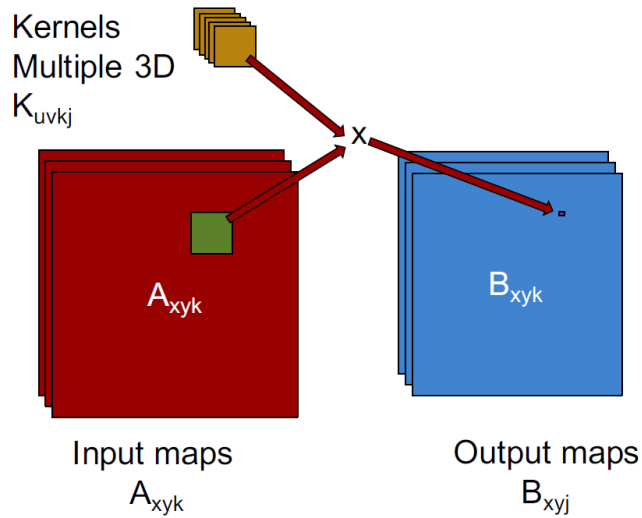
Batched operation is $M \times M$ - gives re-use of weights.

Without batching, would use each element of Weight matrix once.

Want 10-50 arithmetic operations per memory fetch for modern compute architectures.

CNN

Requires convolution and $M \times V$



Filters conserved through plane

Multiply limited - even without batching.

6D Loop

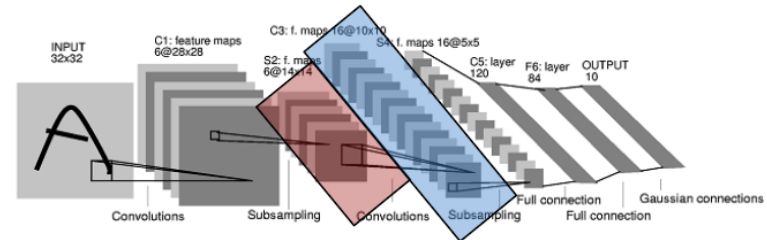
For each output map j

For each input map k

For each pixel x, y

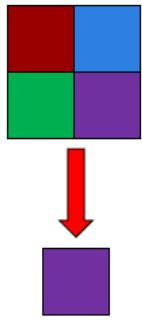
For each kernel element u, v

$$B_{xyj} += A_{(x-u)(y-v)k} \times K_{uvkj}$$



Other Operations

To finish building a DNN



Pooling



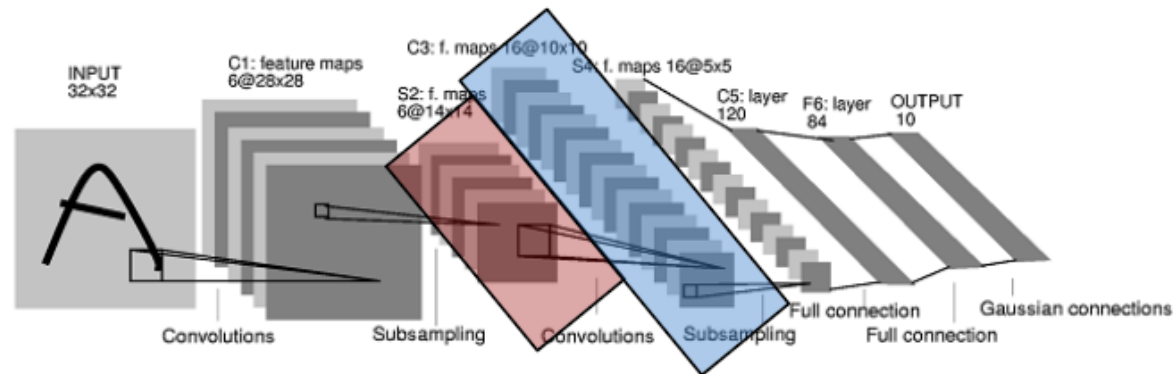
ReLU
(or other non-linear function)

$$w_{ij} += \alpha a_j g_i$$

Weight Update

These are not limiting factors with appropriate GPU use
Complex networks have hundreds of millions of weights.

Lots of Parallelism Available in a DNN

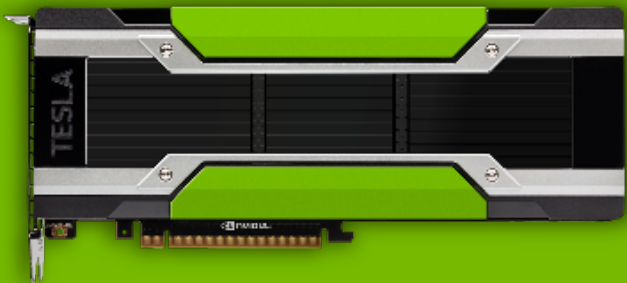


- Inputs
- Points of a feature map
- Filters
- Elements within a filter

- Multiplies within layer are independent
- Sums are reductions
- Only layers are dependent
- No data dependent operations
=> can be statically scheduled

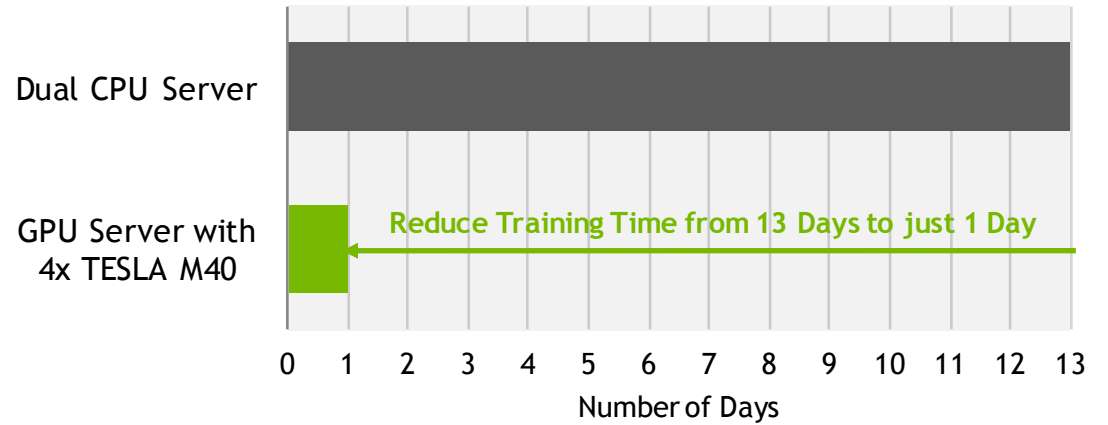
TESLA M40

World's Fastest Accelerator
for Deep Learning Training



28 Gflop/W

13x Faster Training Caffe

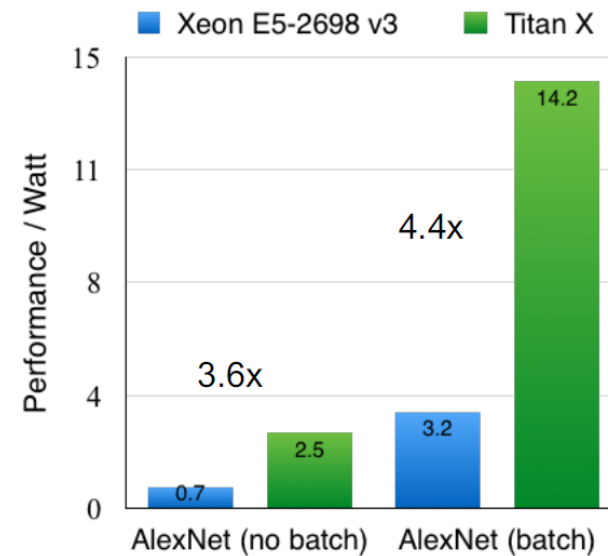
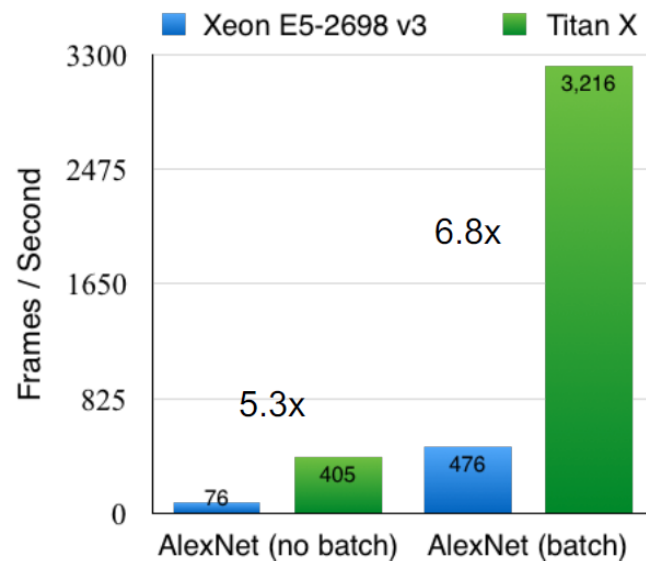


CUDA Cores	3072
Peak SP	7 TFLOPS
GDDR5 Memory	12 GB
Bandwidth	288 GB/s
Power	250W

*Note: Caffe benchmark with AlexNet,
CPU server uses 2x E5-2680v3 12 Core 2.5GHz CPU, 128GB System Memory, Ubuntu 14.04*

Comparing CPU and GPU - server class

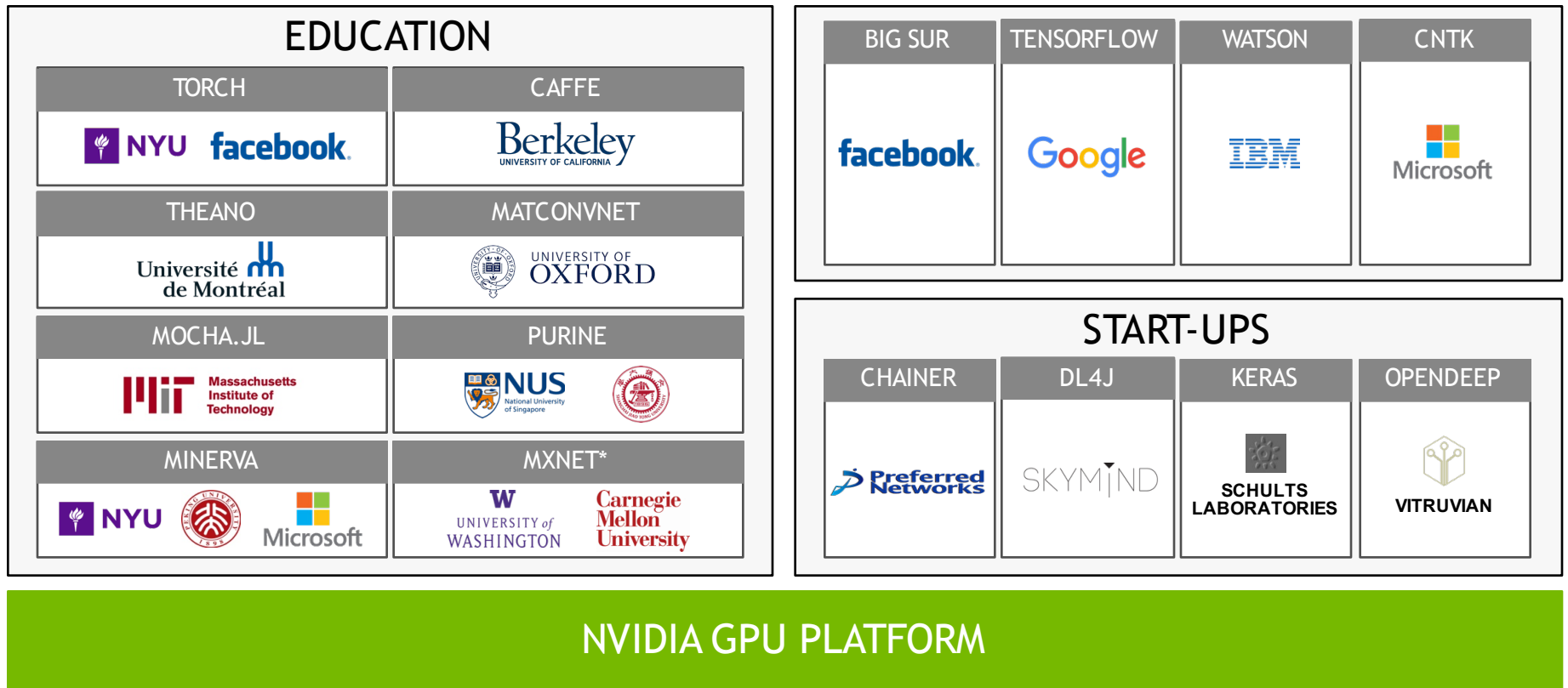
Xeon E5-2698 and Tesla M40



NVIDIA Whitepaper “GPU based deep learning inference: A performance and power analysis.”

DL in practice

The Engine of Modern AI



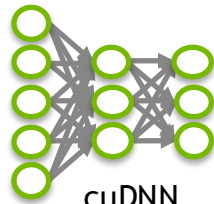
*U. Washington, CMU, Stanford, TuSimple, NYU, Microsoft, U. Alberta, MIT, NYU Shanghai

CUDA for Deep Learning Development

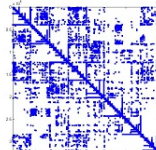
DEEP LEARNING SDK



DIGITS



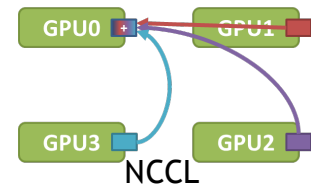
cuDNN



cuSPARSE



cuBLAS



NCCL

TITAN X

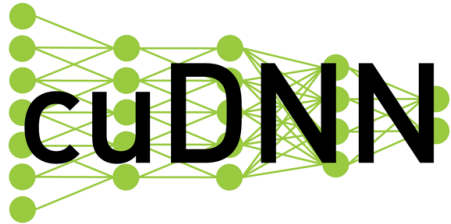


DEVBOX



GPU CLOUD





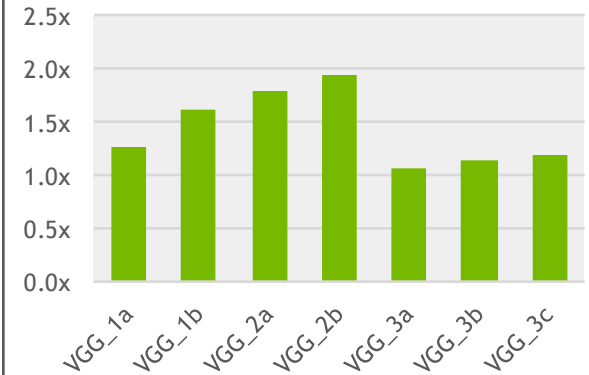
Deep Learning Primitives

Accelerating
Artificial Intelligence

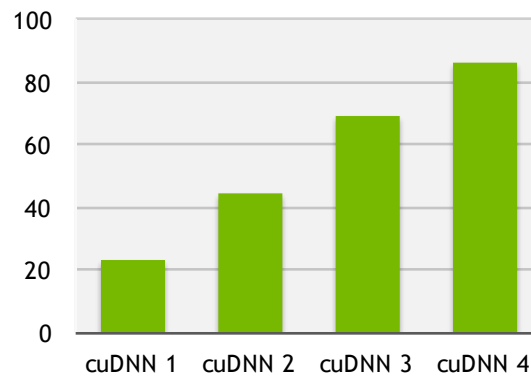
developer.nvidia.com/cudnn

- GPU-accelerated Deep Learning subroutines
- High performance neural network training
- Accelerates Major Deep Learning frameworks: Caffe, Theano, Torch, TensorFlow
- Up to 3.5x faster AlexNet training in Caffe than baseline GPU

Tiled FFT up to 2x faster than FFT



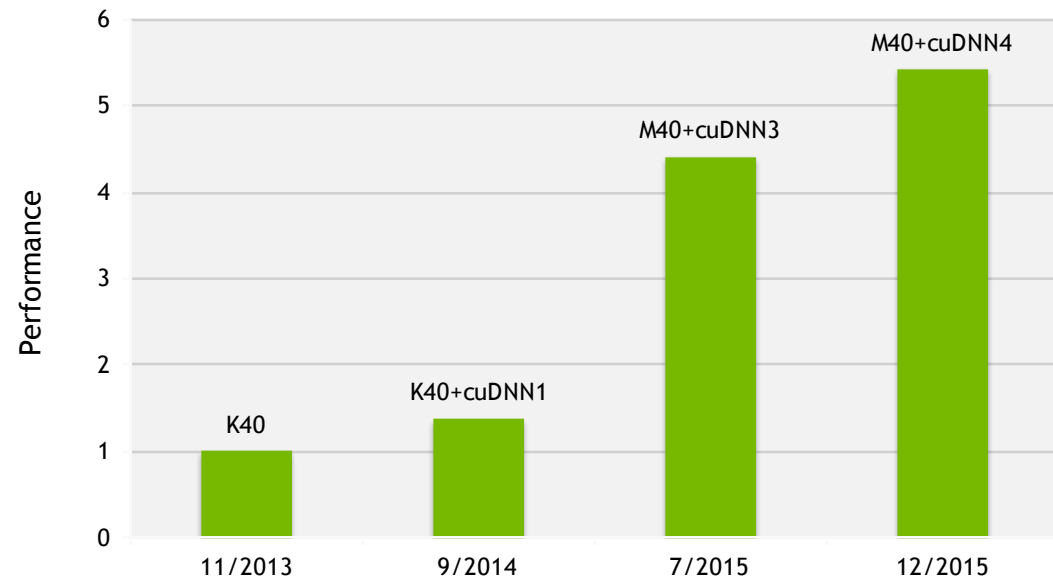
Millions of Images Trained Per Day



CUDA BOOSTS
DEEP LEARNING

5X IN 2 YEARS

Caffe Performance



*AlexNet training throughput based on 20 iterations,
CPU: 1x E5-2680v3 12 Core 2.5GHz. 128GB System Memory, Ubuntu 14.04*

NVIDIA DIGITS

Interactive Deep Learning GPU Training System

Process Data

Job Information

Job Directory
/home/michaelo/.digits
/jobs/20150311-171431-e0d8

Image Type
Color

Image Dimensions
256x256

Resize Mode
half_crop

Create DB (train)

Input file
train.txt

DB Entries
26759

Configure DNN

Select Dataset

PASCAL VOC
ILSVRC 2012
MNIST Dataset

Solver Options

Training epochs
30

Validation interval (in epochs)
1

Batch size
100

Base Learning Rate
0.01

Show advanced learning rate options

Custom Network

```
layer {  
  name: "conv1"  
  type: "Convolution"  
  bottom: "data"  
  top: "conv1"  
  param {  
    f_mult: 1  
    decay_mult: 1  
  }  
}
```

Monitor Progress

Solver
solver.prototxt
Network (train/val)
train_val.prototxt
Network (deploy)
deploy.prototxt

Dataset
voc_cropped@256x256
Done Wed Mar 11, 05:16:57 PM

Image Size
256x256

Image Type
COLOR

Create DB (train)
26759 images

Create DB (val)
8917 images

Loss and Accuracy Graph

Epoch	Loss (train)	Loss (val)	Accuracy (val)
0.0	3.5	3.5	0
2.5	1.8	1.8	60
5.0	1.5	1.5	65
7.5	1.4	1.4	68
10.0	1.3	1.3	70

Visualize Layers

Predictions

8
3
0
6
4

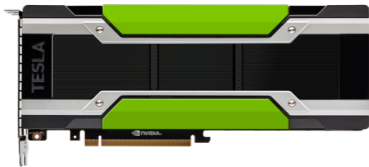
Layer Activations

conv1

pool1

developer.nvidia.com/digits

ONE ARCHITECTURE — END-TO-END AI



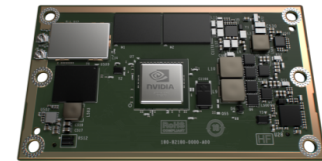
Tesla
for Cloud



Titan X
for PC



DRIVE PX
for Auto



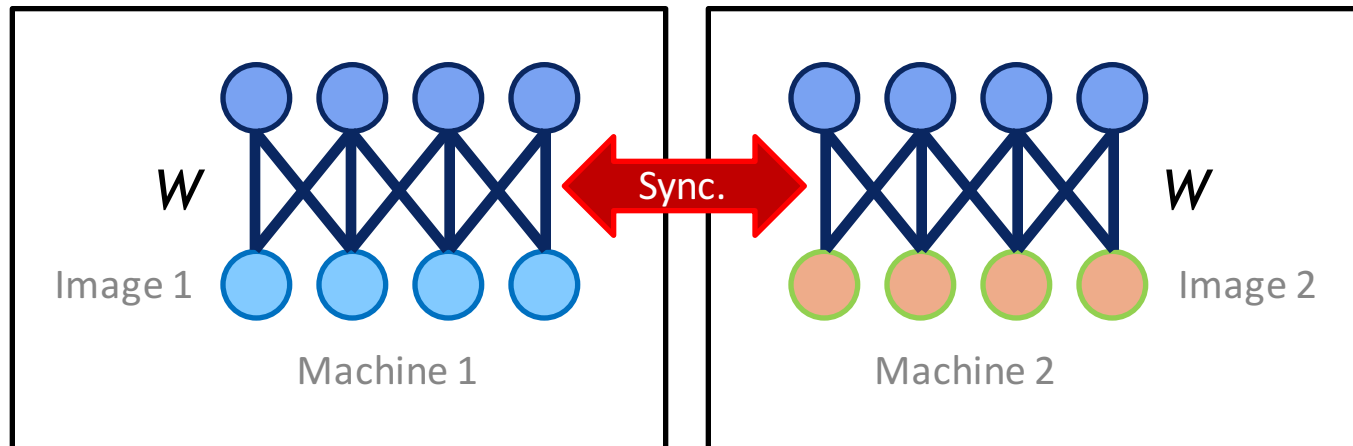
Jetson
for Embedded

Scaling DL



Scaling Neural Networks

Data Parallelism



Notes:

Need to sync model across machines.

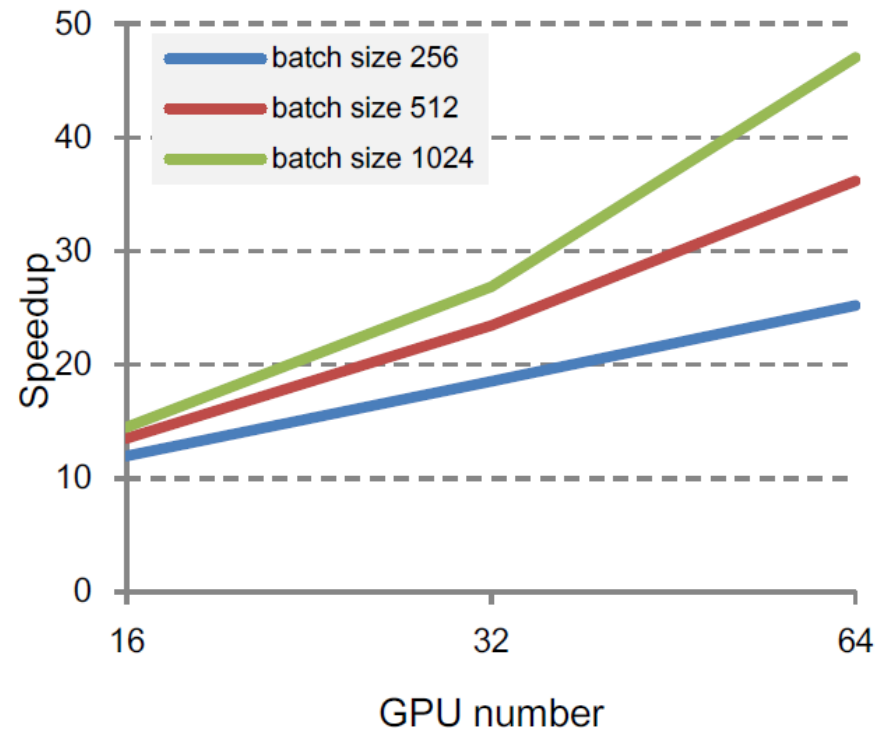
Largest models do not fit on one GPU.

Requires P-fold larger batch size.

Works across many nodes – parameter server approach – linear speedup.

Multiple GPUs

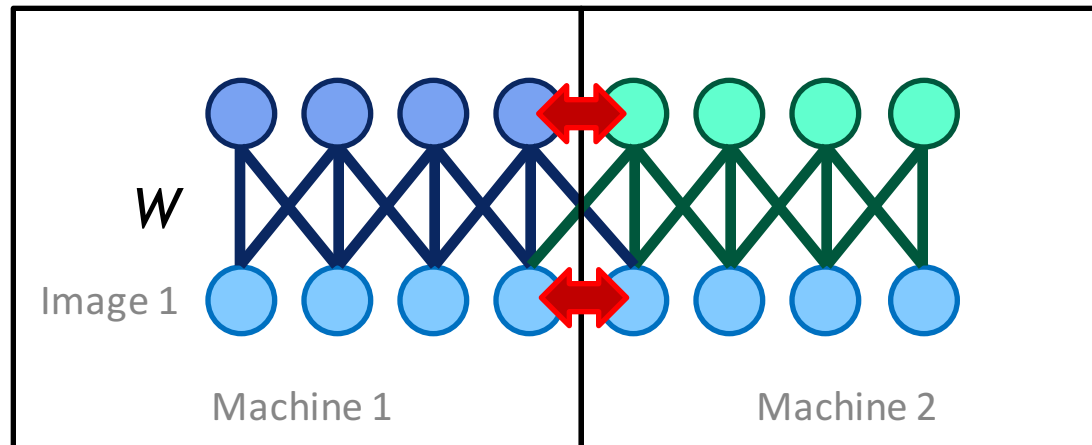
Near linear scaling - data parallel.



Ren Wu et al, Baidu, "Deep Image: Scaling up Image Recognition." arXiv 2015

Scaling Neural Networks

Model Parallelism



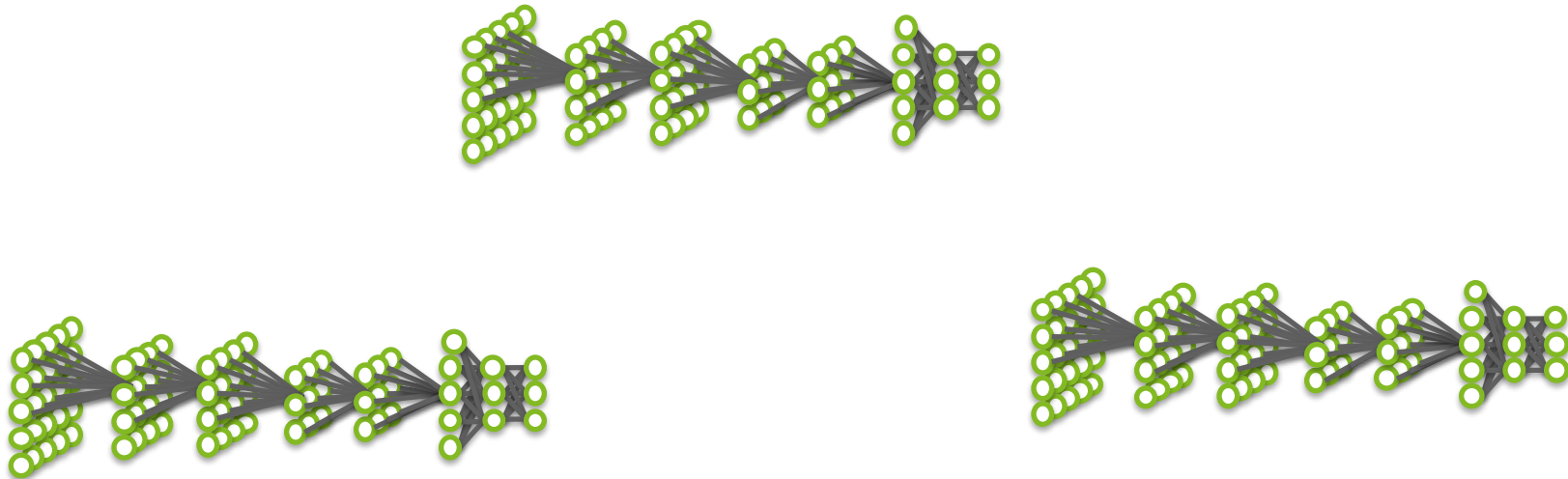
Notes:

- Allows for larger models than fit on one GPU.
- Requires much more frequent communication between GPUs.
- Most commonly used within a node – GPU P2P.
- Effective for the fully connected layers.

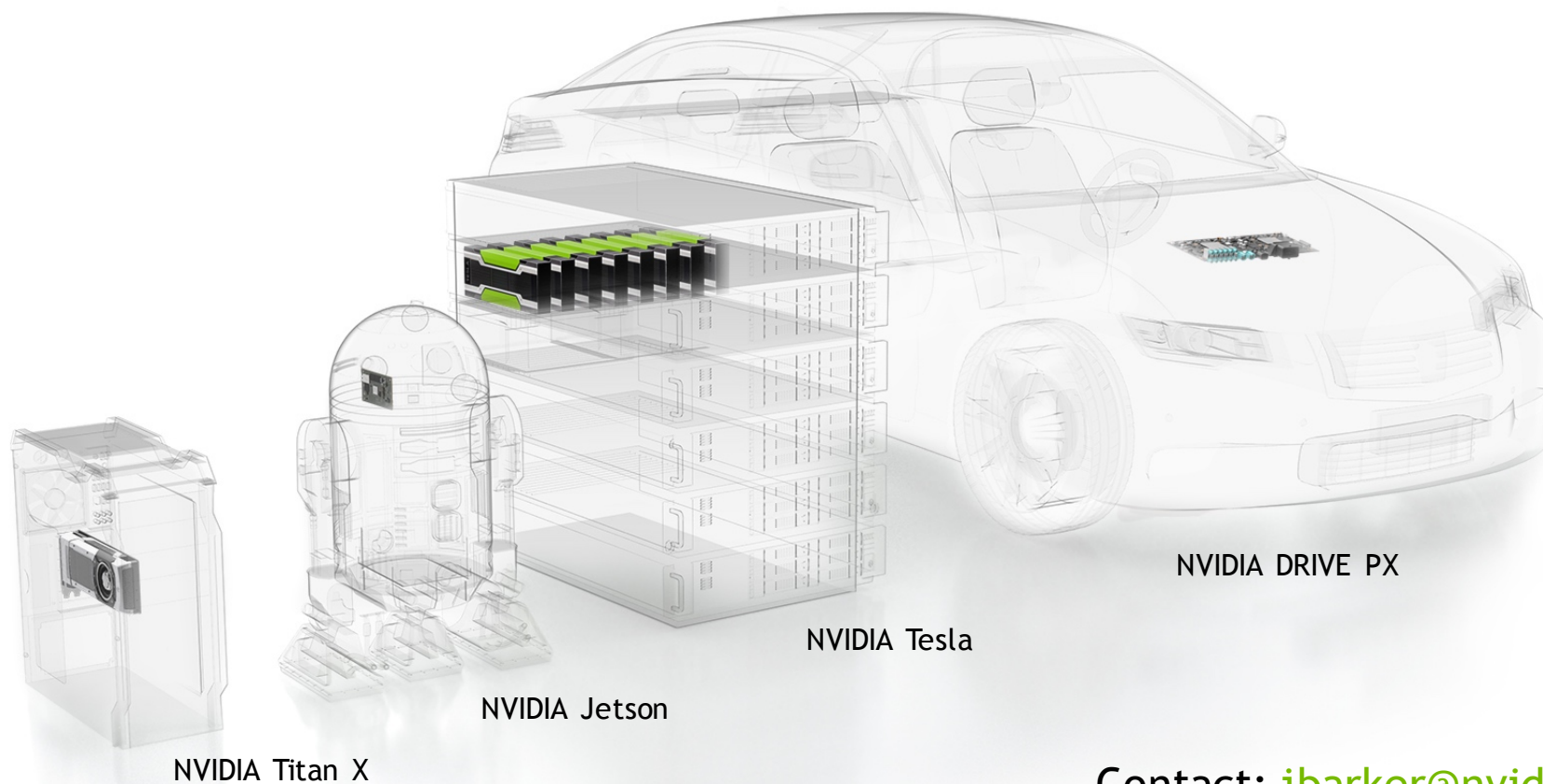
Scaling Neural Networks

Hyper Parameter Parallelism

Try many alternative neural networks in parallel – on different CPU / GPU / Machines.
Probably the most obvious and effective way!



Deep Learning Everywhere



NVIDIA Titan X

NVIDIA Jetson

NVIDIA Tesla

NVIDIA DRIVE PX

Contact: jbarker@nvidia.com