



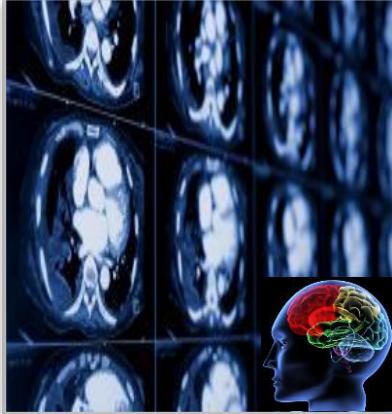
FRAMEWORKS FOR IMAGE CLASSIFICATION

Jonny Hancox
Deep Learning Solution Architect
October 2018

AGENDA

- What is a Deep Learning Framework?
- Overview of the field
- Some code examples
- Conclusions

DEEP LEARNING DRIVES INNOVATION

INTERNET SERVICES	HEALTHCARE	MEDIA & ENT.	SECURITY & DEFENSE	SELF DRIVING CARS
				
<ul style="list-style-type: none">• Image/Video classification• Speech recognition• Natural language processing	<ul style="list-style-type: none">• Cancer cell detection• Diabetic grading• Drug discovery	<ul style="list-style-type: none">• Video captioning• Content based search• Real time translation	<ul style="list-style-type: none">• Face recognition• Video surveillance• Cyber security	<ul style="list-style-type: none">• Pedestrian detection• Lane tracking• Recognize traffic signs

WHAT IS A DL FRAMEWORK?

A software framework is:

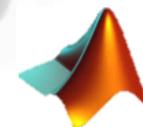
- a universal, reusable software platform used to develop applications, products and solutions
- an abstraction in which software providing generic functionality can be selectively changed by additional user-written code, thus providing application-specific software

A DL framework is therefore:

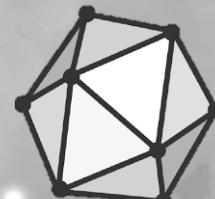
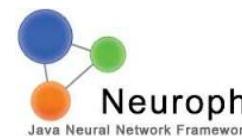
- A set of tools and libraries to enable a developer or data scientist to train and deploy DL solutions

DEEP LEARNING FRAMEWORKS

fast.ai



Caffe

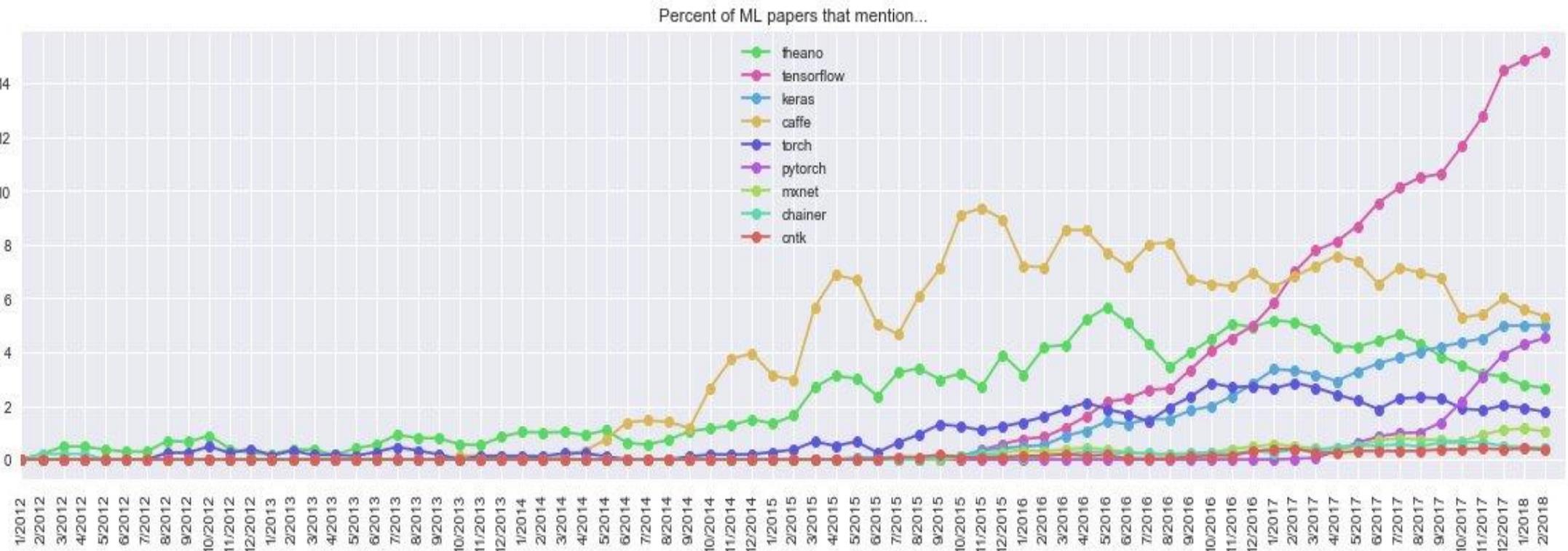


DEEP LEARNING FRAMEWORKS

Some important attributes to consider:

- Licensing Model
- Open or Closed Source
- Platform support
- Language Written In
- API languages
- CUDA support
- Multi-GPU (Scalability)
- Fine Control vs Ease of Use/Productivity
- Community Support

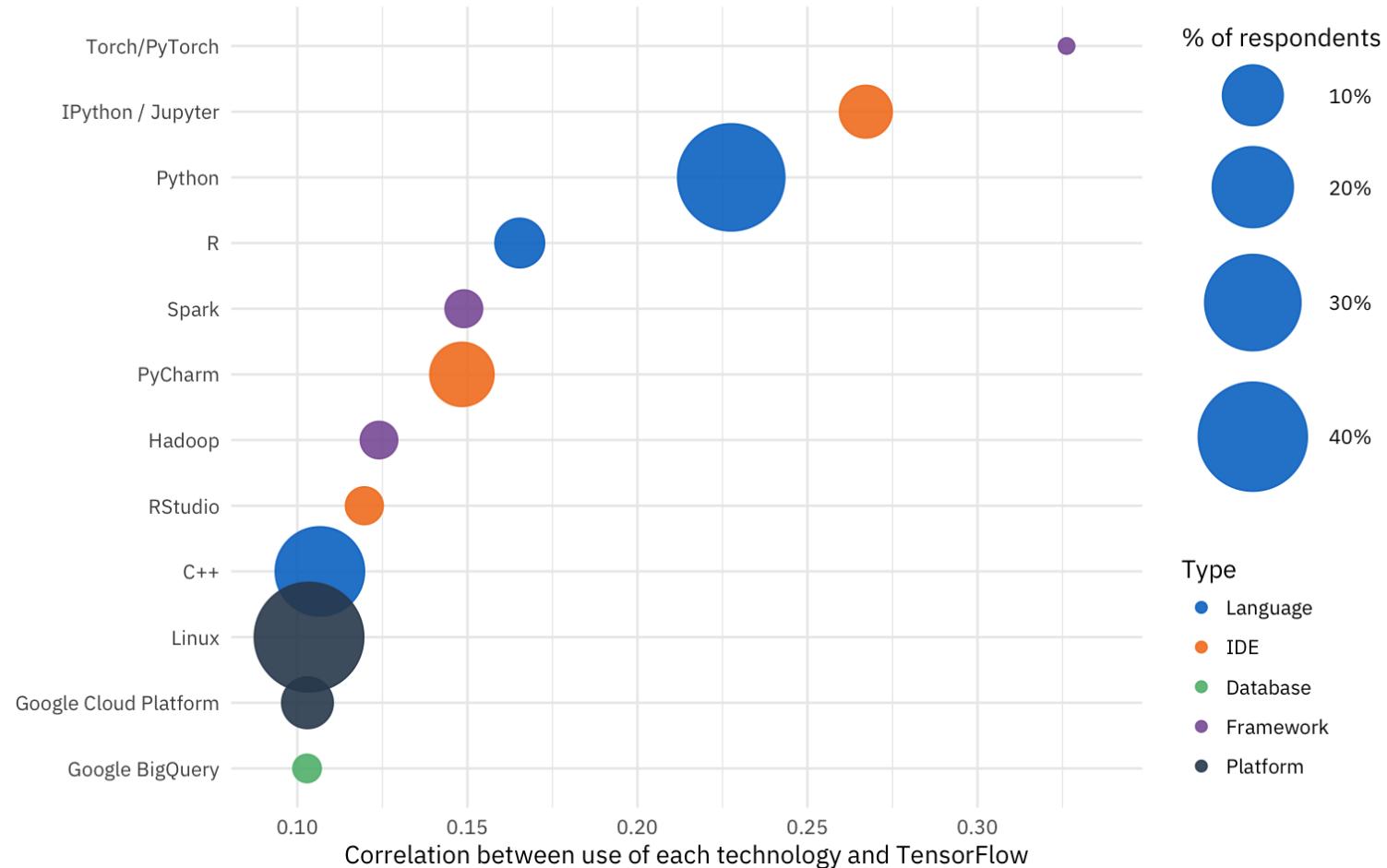
DEEP LEARNING FRAMEWORKS





Technologies used by developers who use TensorFlow

On the 2018 Stack Overflow survey



Source: [Stack Overflow 2018 Developer Survey](#)

THE DEEP LEARNING SOFTWARE STACK

GPU PROGRAMMING LANGUAGES

Numerical analytics ➤

MATLAB, Mathematica, LabVIEW

Fortran ➤

CUDA Fortran, OpenACC

C, C++ ➤

CUDA C++, OpenACC

Python ➤

CUDA Python, PyCUDA

C# ➤

Altimesh Hybridizer, Alea GPU

ACCELERATED DEEP LEARNING TRAINING STACK

IMAGENET

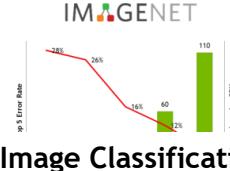


Image Classification



Object Detection

COMPUTER VISION



Voice Recognition



Language Translation



Recommendation Engines



Sentiment Analysis

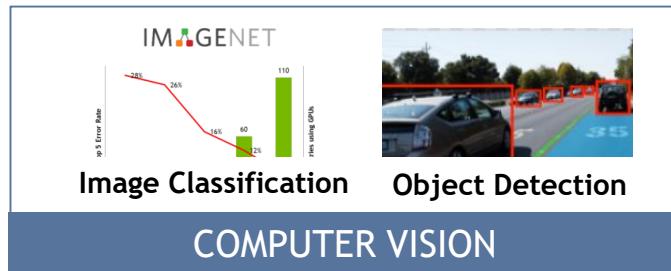
SPEECH AND AUDIO

NATURAL LANGUAGE PROCESSING



DEEP LEARNING FRAMEWORKS

ACCELERATED DEEP LEARNING TRAINING STACK



DL SW Libraries: Tensor/Graph Execution Engines (AKA Frameworks)

ACCELERATED DEEP LEARNING TRAINING STACK

IMGENET



Image Classification



Object Detection

COMPUTER VISION



Voice Recognition



Language Translation

SPEECH AND AUDIO



Recommendation
Engines



Sentiment Analysis

NATURAL LANGUAGE PROCESSING

Network description, Workflow, Hyper-parameter Sweep,
Experiment, Data and Job Management

DL SW Libraries: Tensor/Graph Execution Engines (AKA Frameworks)

Architecture Specific Optimization Layer

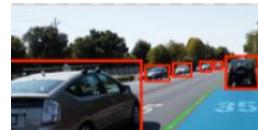


ACCELERATED DEEP LEARNING TRAINING STACK

IMAGENET



Image Classification



Object Detection

COMPUTER VISION



Voice Recognition



Language Translation

SPEECH AND AUDIO



Recommendation
Engines

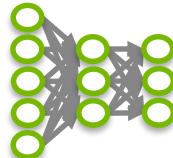


Sentiment Analysis

NATURAL LANGUAGE PROCESSING

Network description, Workflow, Hyper-parameter Sweep,
Experiment, Data and Job Management

DL SW Libraries: Tensor/Graph Execution Engines (AKA Frameworks)



cuDNN

DEEP LEARNING

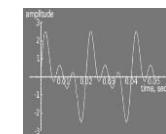
cuBLAS



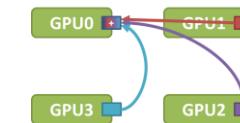
cuSPARSE



cuFFT



MATH LIBRARIES



MULTI-GPU



ACCELERATED DEEP LEARNING TRAINING STACK

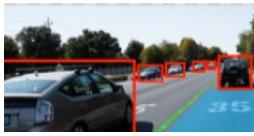
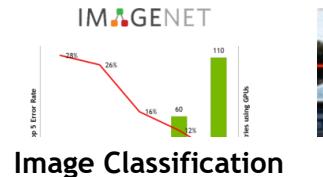


Image Classification

Object Detection

COMPUTER VISION



Voice Recognition



Language Translation

SPEECH AND AUDIO



Recommendation Engines



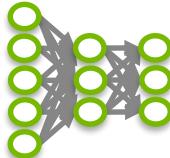
Sentiment Analysis

NATURAL LANGUAGE PROCESSING

Network description, Workflow, Hyper-parameter Sweep,
Experiment, Data and Job Management



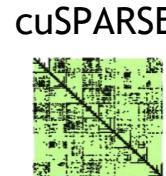
DEEP LEARNING FRAMEWORKS



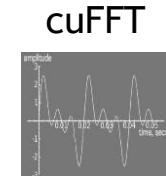
cuDNN



cuBLAS



cuSPARSE



cuFFT

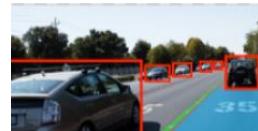
DEEP LEARNING

MATH LIBRARIES

MULTI-GPU



ACCELERATED DEEP LEARNING TRAINING STACK



Object Detection

COMPUTER VISION



Voice Recognition



Language Translation

SPEECH AND AUDIO



Recommendation Engines



Sentiment Analysis

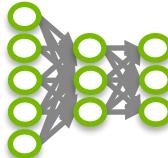
NATURAL LANGUAGE PROCESSING

Productivity Layer/Rapid experimentation: DIGITS, NVIDIA GPU Cloud

UI / JOB MANAGEMENT / DATASET VERSIONING/ VISUALIZATION



DEEP LEARNING FRAMEWORKS

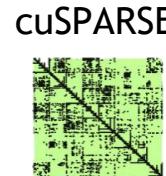


cuDNN

DEEP LEARNING



cuBLAS

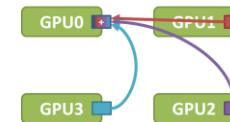


cuSPARSE



cuFFT

MATH LIBRARIES



MULTI-GPU



DEEP LEARNING PERFORMANCE CONSIDERATIONS

- Preprocess data outside of main training loop
 - Training will be run for many epochs, and repeated often during model development
- Use high level support for data loading and multi-GPU parallelization if available
 - Many pitfalls that may hamper performance due to data loading or CPU bottlenecks
 - **But:** make sure the interface is actually meant for high performance / production.
Example: TensorFlow `feed_dict` - bad performance but used in many examples
- Prefer versions of blocks / layers that is implemented on top of cuDNN if possible
- Prefer high level blocks if possible
 - Process all time steps of RNN using cuDNN based block instead of explicit loop

DEEP LEARNING PERFORMANCE CONSIDERATIONS

- All deep learning frameworks can achieve good (about equal) multi GPU performance and scaling
- **But:** Engineering effort differ significantly between frameworks and approaches
 - **Trade-off** between flexibility and low engineering effort
 - Monolithic frameworks offer less opportunity to “shoot oneself in the foot”
 - Many frameworks offer both high level interfaces and low level primitives for building any kind of multi GPU solution

NVIDIA cuDNN 7

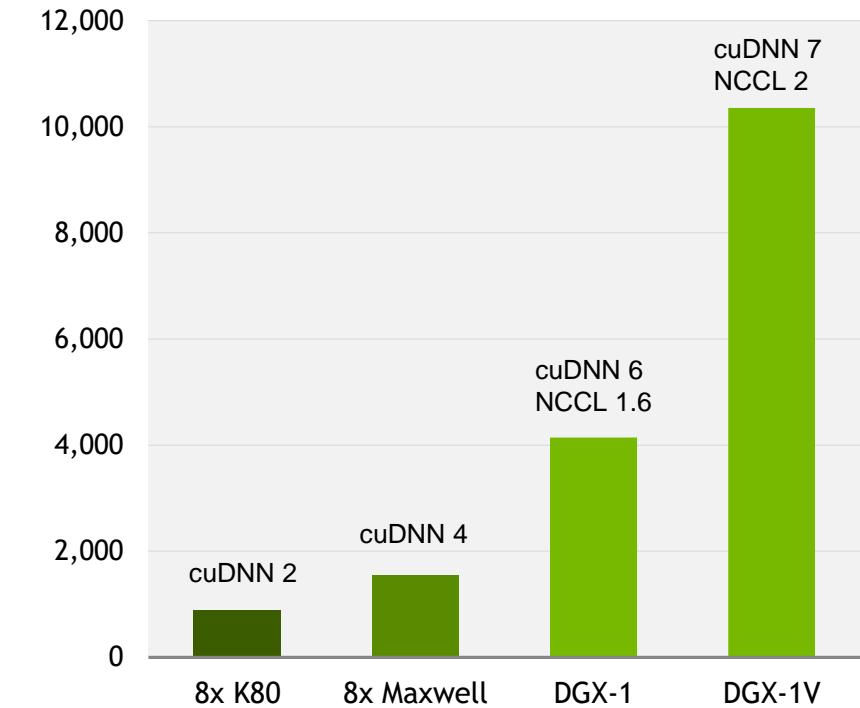
Deep Learning Primitives

High performance building blocks for deep learning frameworks

Drop-in acceleration for widely used deep learning frameworks such as Caffe2, Microsoft Cognitive Toolkit, PyTorch, Tensorflow, Theano and others

Accelerates industry vetted deep learning algorithms, such as convolutions, LSTM RNNs, fully connected, and pooling layers

Fast deep learning training performance tuned for NVIDIA GPUs



“ NVIDIA has improved the speed of cuDNN with each release while extending the interface to more operations and devices at the same time.”

— Evan Shelhamer, Lead Caffe Developer, UC Berkeley

NVIDIA Collective Communications Library (NCCL) 2

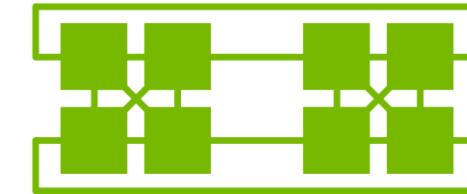
Multi-GPU and multi-node collective communication primitives

High-performance multi-GPU and multi-node collective communication primitives optimized for NVIDIA GPUs

Fast routines for multi-GPU multi-node acceleration that maximizes inter-GPU bandwidth utilization

Easy to integrate and MPI compatible. Uses automatic topology detection to scale HPC and deep learning applications over PCIe and NVLink

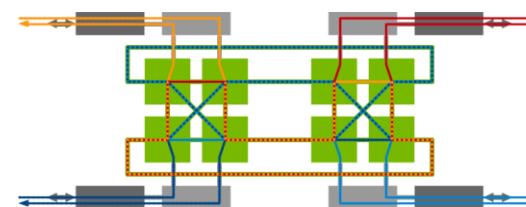
Accelerates leading deep learning frameworks such as Caffe2, Microsoft Cognitive Toolkit, MXNet, PyTorch and more



Multi-GPU:
NVLink
PCIe



Multi-Node:
InfiniBand verbs
IP Sockets

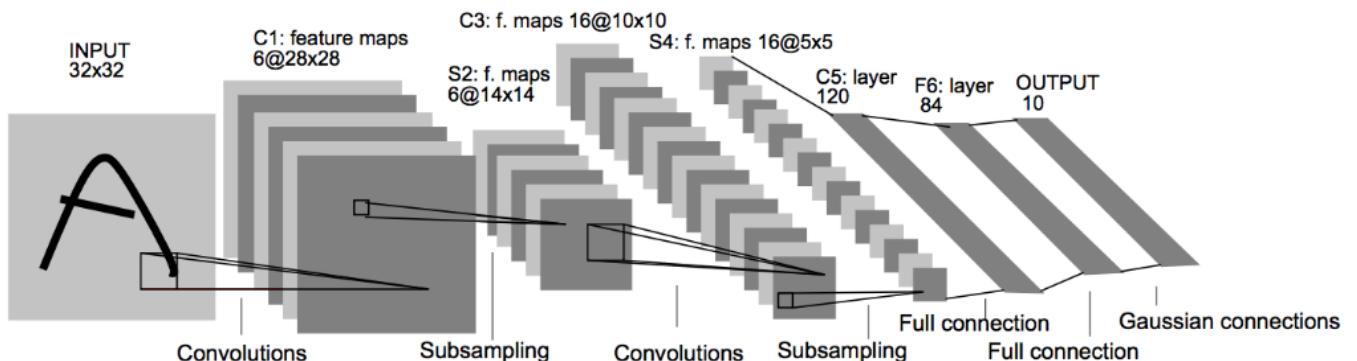
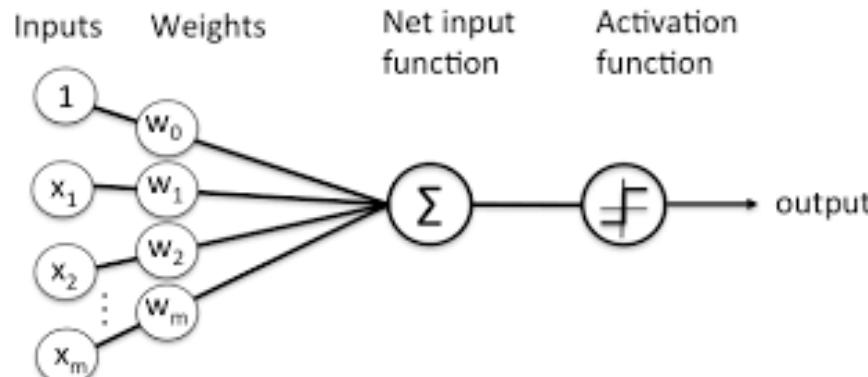


Automatic
Topology
Detection

DL PLATFORMS

Many different platforms and styles of interface but one unifying set of concepts:

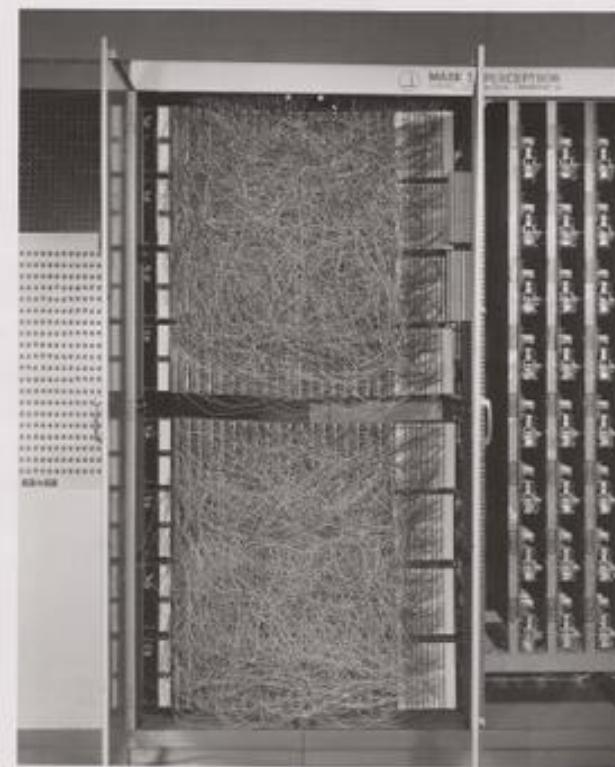
- Convolutions
- Pooling
- Activations
- Fully connected layers
- Back-propagation
- Etc.



CNN called LeNet by Yann LeCun (1998)

THE PERCEPTRON

The fundamental unit of artificial neural networks



NVIDIA DIGITS

Interactive Deep Learning GPU Training System

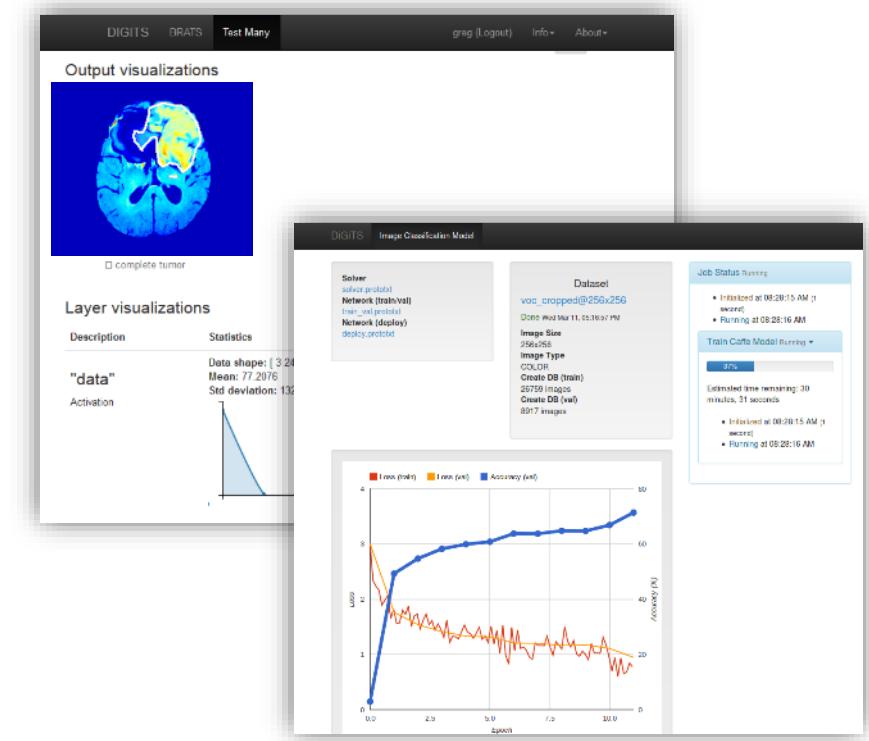
Interactive deep learning training application for engineers and data scientists

Simplify deep neural network training with an interactive interface to train and validate, and visualize results

Built-in workflows for image classification, object detection and image segmentation

Improve model accuracy with pre-trained models from the DIGITS Model Store

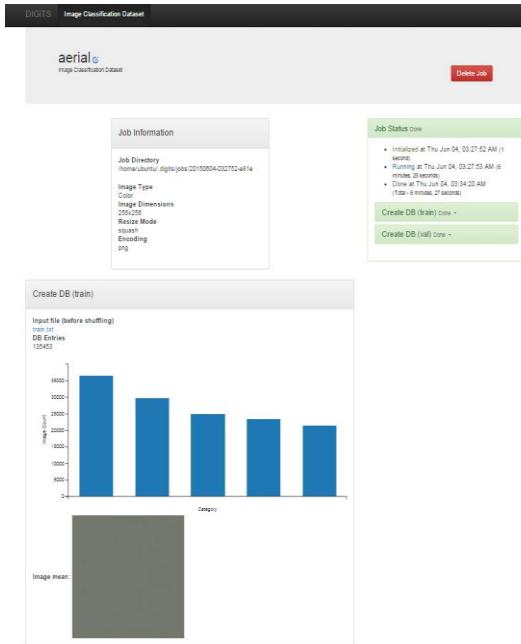
Faster time to solution with multi-GPU acceleration



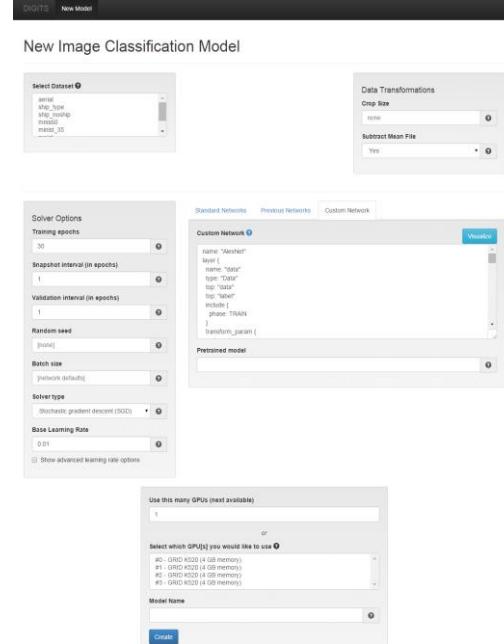
NVIDIA DIGITS

Interactive Deep Learning GPU Training System

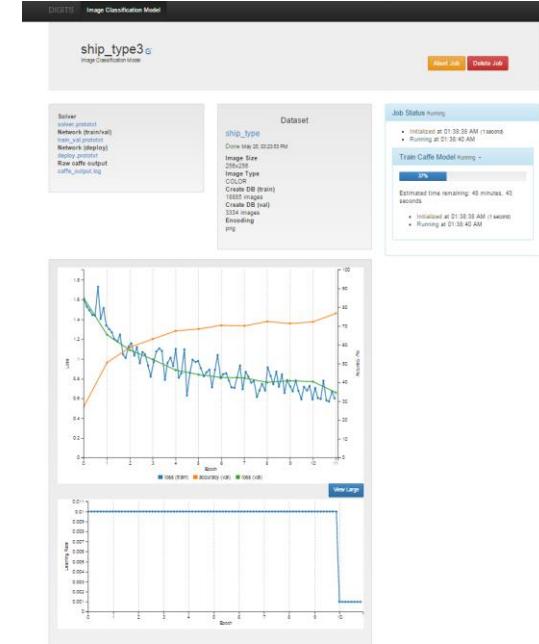
Process Data



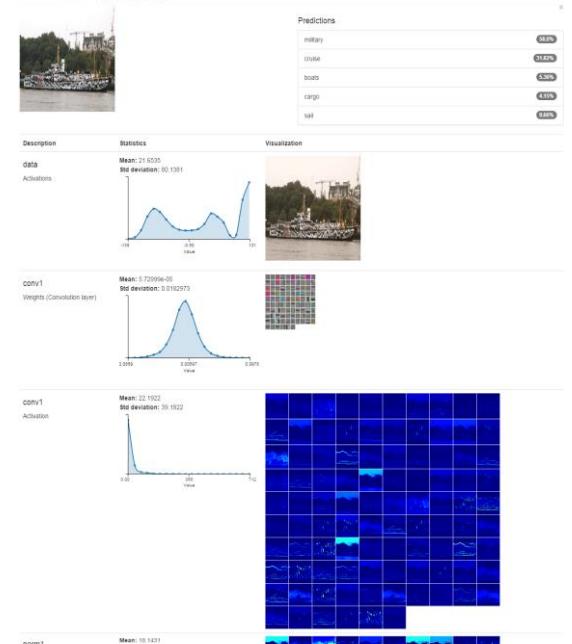
Configure DNN



Monitor Progress



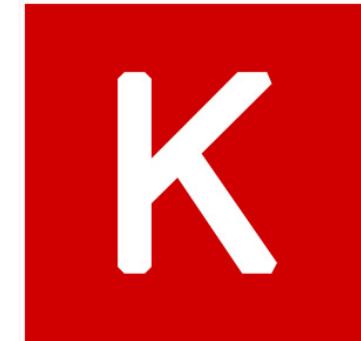
Visualization



DL FRAMEWORKS - KERAS

Enables fast experimentation with DNNs

- Supports multiple backends, including Tflow, CNTK, Theano
- Contains many commonly used layers types
- Features many tools to simplify working with text and images
- Keras can create deployments for iOS Andriod and the JVM
- Supports multi-GPU training



```
model = Sequential()  
model.add(Conv2D(64, (3, 3), activation='relu'))
```

DL FRAMEWORKS - TENSORFLOW

Powerful, popular and well-maintained

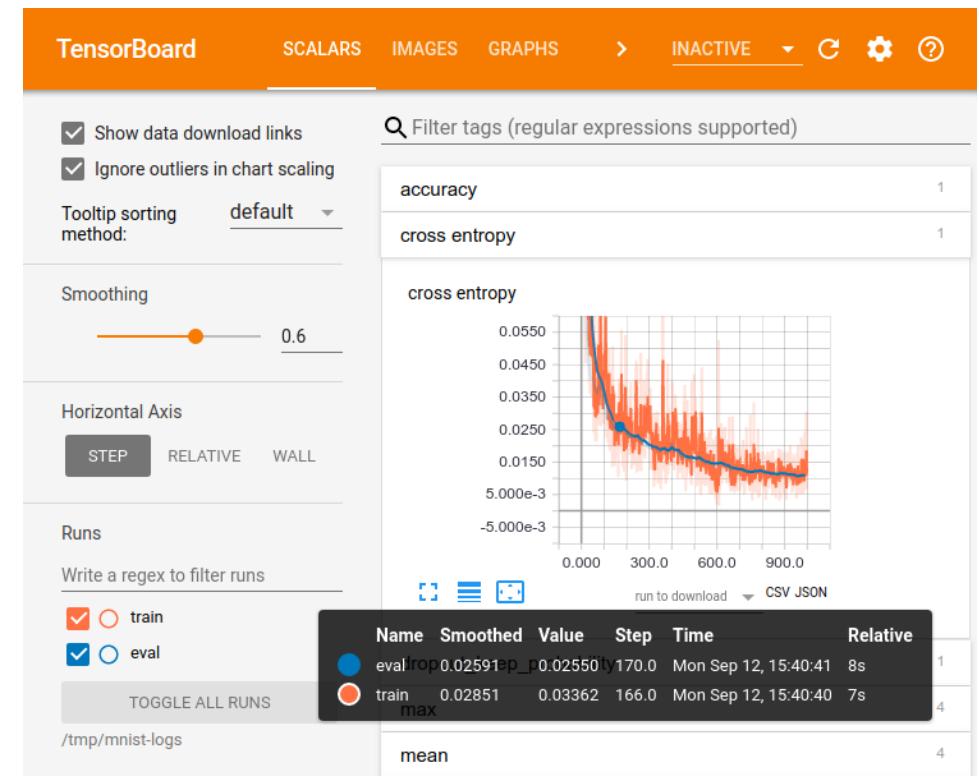
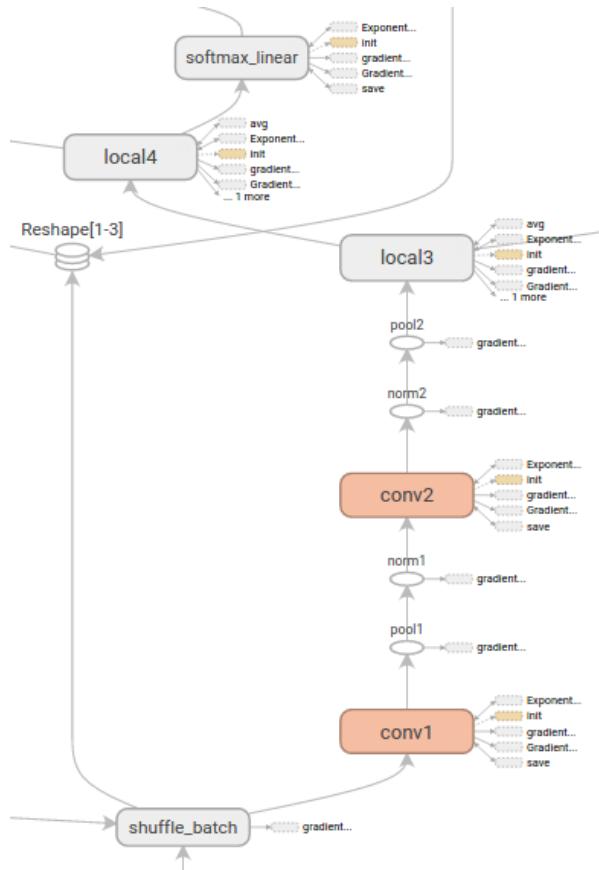


- Supports huge variety of DNN features
- Includes Tensorboard GUI for layer and output visualisation
- Many popular third-party integrations
- Python & C++ APIs
- Custom layers can be created
- Uses a declarative syntax for efficient graph execution

```
conv1 = tf.nn.conv2d(_X, _weights['wc1']
                    , strides=[1, 1, 1, 1], padding='SAME')
conv1 = tf.nn.bias_add(conv1 ,_biases['bc1'])
```

DL FRAMEWORKS - TENSORFLOW

Tensorboard - useful GUI



DL FRAMEWORKS - PYTORCH



Increasingly popular and easy to use

Primarily developed by Facebook

Good GPU Acceleration

Ideal for Python users (Pythonic implementation)

Dynamic computation graphs

Tensor usage similar to Numpy

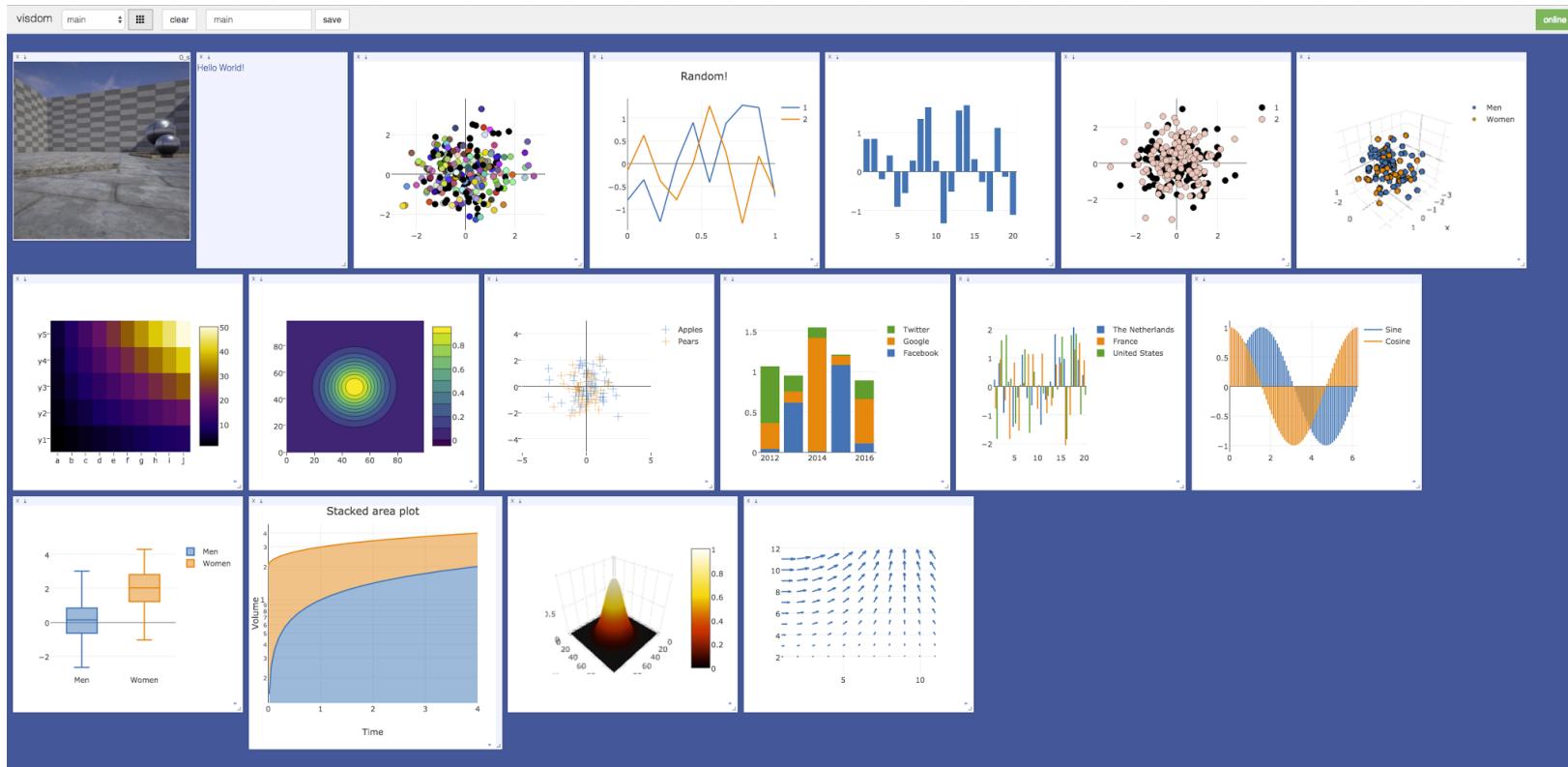
More suited to training than deployment

```
model = torch.nn.Sequential(  
    Torch.nn.add(Conv2D(64, (3, 3), activation='relu'))  
)
```

DL FRAMEWORKS - PYTORCH

PYTORCH

Visdom is Pytorch's equivalent to
TensorBoard



DL FRAMEWORKS - CAFFE2

Great performance and flexibility

- Primarily developed by Facebook
- Good GPU Acceleration
- Python & C++ APIs
- Support multi-GPU and multi-node execution
- Being merged with PyTorch



```
conv1 = brew.conv(model, data, 'conv1', 1, 20, 5)
pool1 = brew.max_pool(model, conv1, 'pool1', kernel=2, stride=2)
conv2 = brew.conv(model, pool1, 'conv2', 20, 50, 5)
```

DL FRAMEWORKS - MS COGNITIVE TOOLKIT

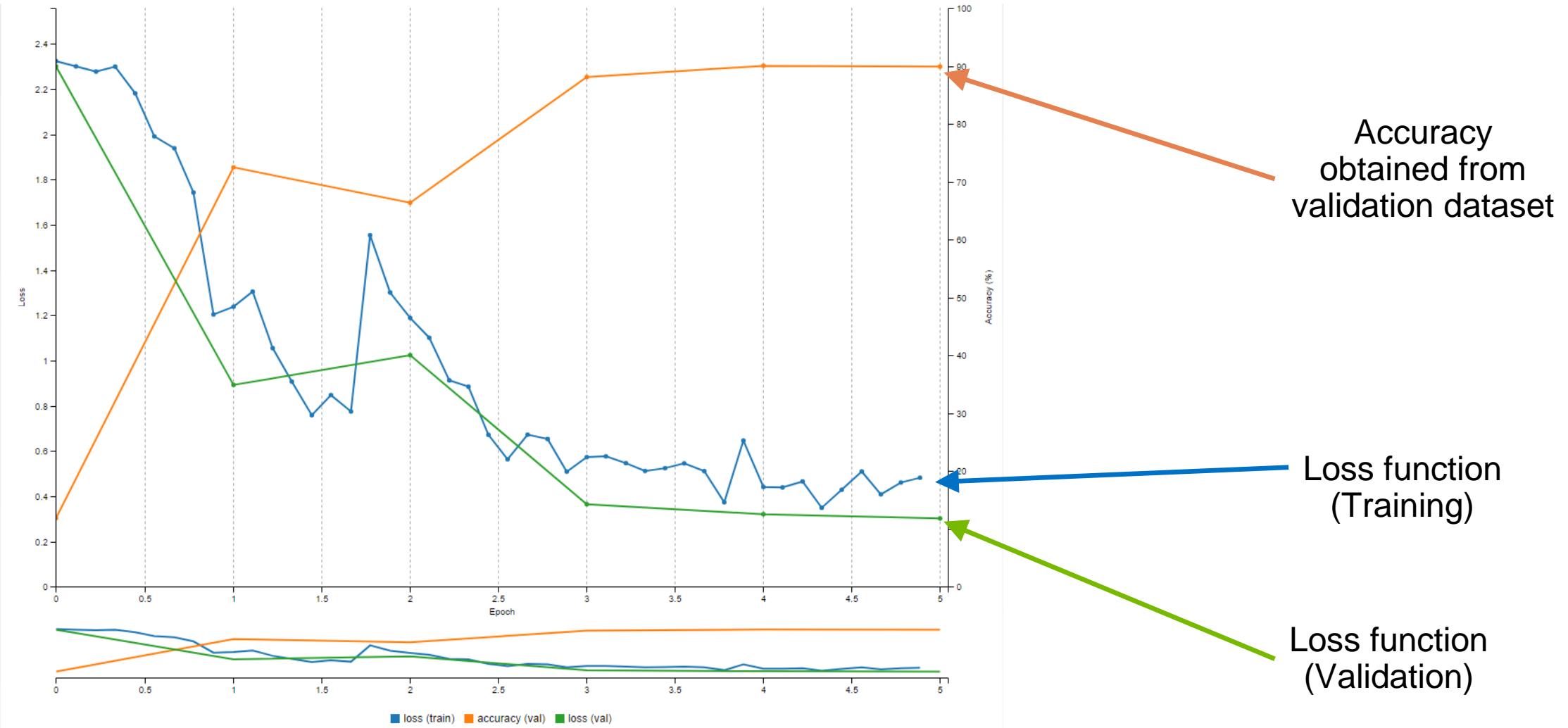
Great performance and flexibility

- Created by MS Speech researchers (formerly CNTK)
- Open-sourced since Jan 2016
- Multiple language support (C#, R, Python, C++)
- Supports multi-GPU, multi-node execution and Spark
- Composes computational graphs from NNs
- Lower precision support
- Fast execution!



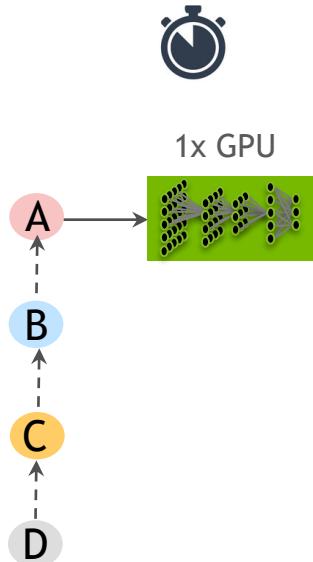
```
import cntk as C
input_dim = 784
num_output_classes = 10
feature = C.input_variable(input_dim)
label = C.input_variable(num_output_classes)
```

EVALUATE THE MODEL

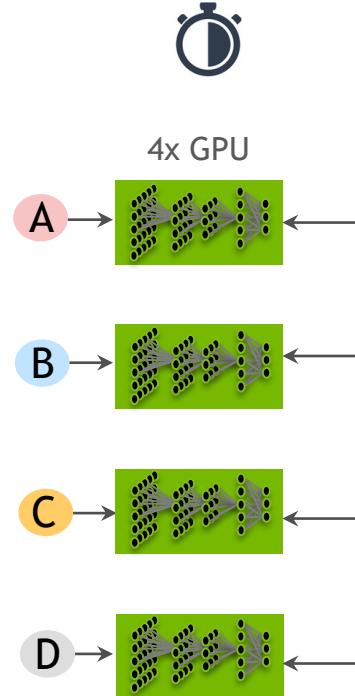


DISTRIBUTED DEEP LEARNING

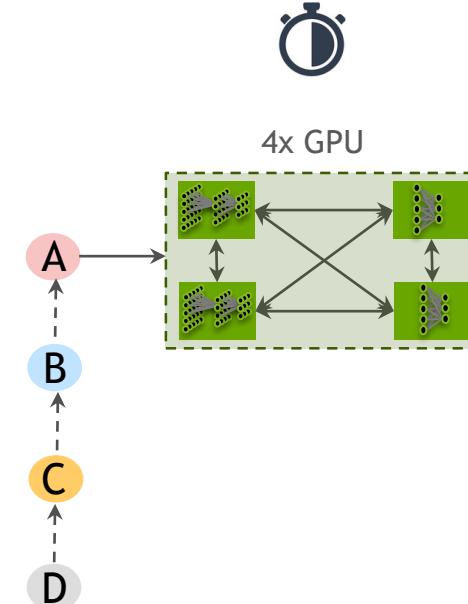
SINGLE GPU



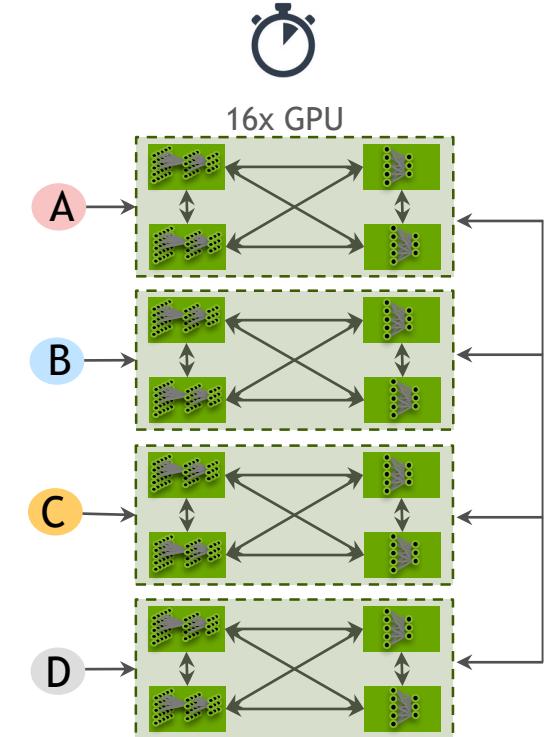
DATA PARALLEL



MODEL PARALLEL



DATA AND MODEL PARALLEL



Data & Model Parallel training yields increasingly faster time-to-solution

CONCLUSIONS

- There are many frameworks out there
- Most support the same operations
- There is no ‘best’ framework
- The right choice for you is dependent on several factors
- All the main frameworks support GPUs



QUESTIONS?

Thank you!

jhancox@nvidia.com



USING DEEP NEURAL NETWORKS FOR OBJECT DETECTION

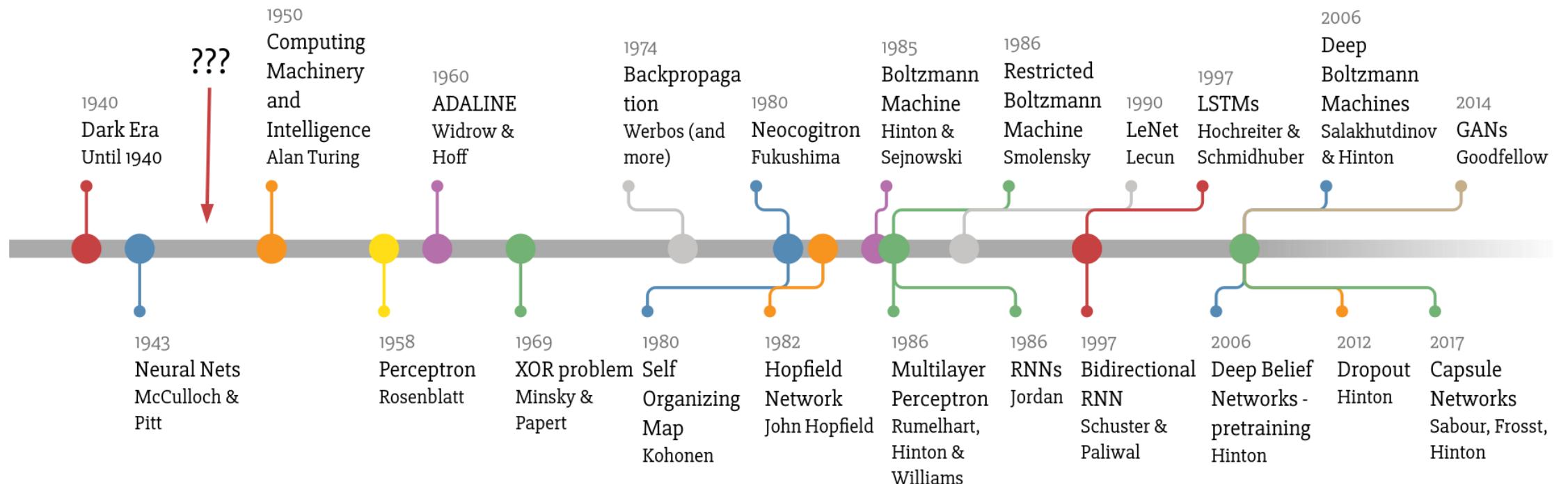
Jonny Hancox
Deep Learning Solution Architect
October 2018

HOW DL CAN BE APPLIED

INPUTS	BUSINESS QUESTION	AI/DL TASK	EXAMPLE OUTPUTS HEALTHCARE	EXAMPLE OUTPUTS RETAIL	EXAMPLE OUTPUTS FINANCE
 Text Data	Is “it” <u>present</u> or not?	Detection	Cancer Detection	Targeted ads	Cybersecurity
 Images	What <u>type</u> of thing is “it”?		Image Classification	Basket Analysis	Credit Scoring
 Video	To what <u>extent</u> is “it” present?		Tumor Size/Shape Analysis	Build 360° Customer View	Credit Risk Analysis
 Audio	What is the likely outcome?		Survivability Prediction	Sentiment & behavior recognition	Fraud Detection
	What will satisfy the objective?		Therapy Recommendation	Recommendation Engine	Algorithmic Trading
	What is the speaker saying?	Natural Language Processing	Expert diagnosis	Virtual personal assistants	Robo Advisors

NEURAL NETS

Some (recent) history



RECOGNIZING DIGITS

The canonical image classification problem

We have 10 digits - how do we recognise them and distinguish between them?

0123456789



0123456789



0123456789



RECOGNIZING DIGITS

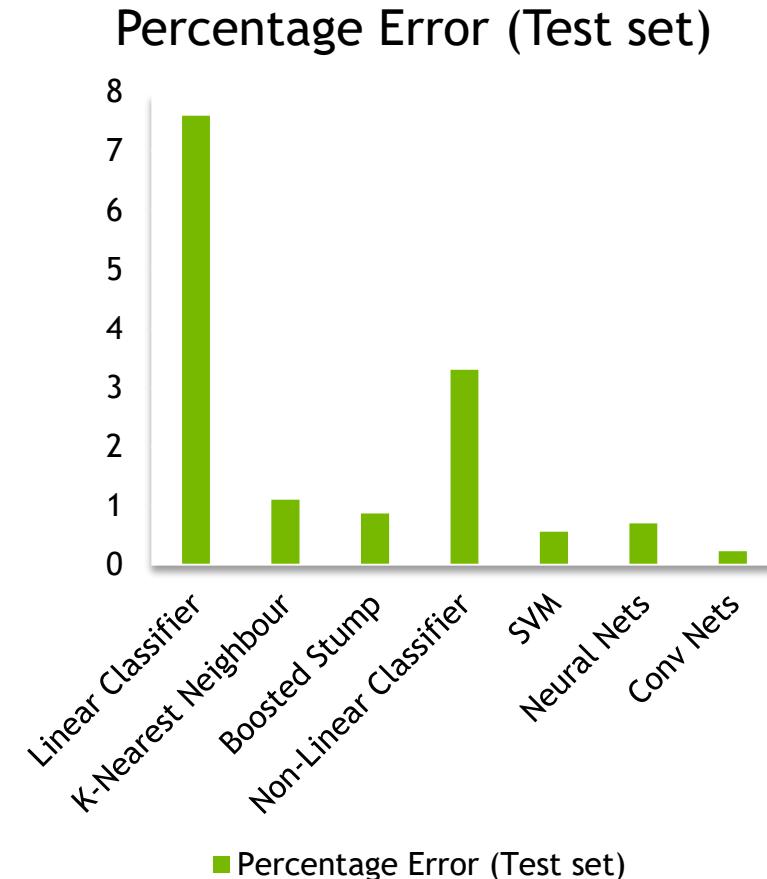
The MNIST Database of handwritten digits

[Yann LeCun](#) Courant Institute, NYU

[Corinna Cortes](#) Google Labs, New York

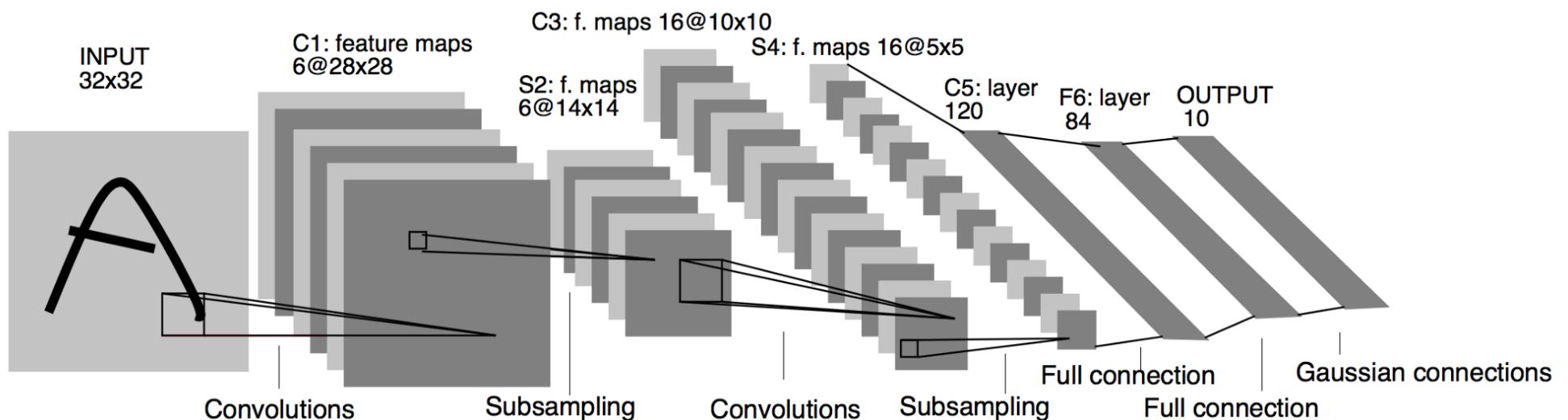
[Christopher J.C. Burges](#) Microsoft Research, Redmond

- A training set of 60,000 examples
- test set of 10,000 examples
- The digits have been size-normalised and centred



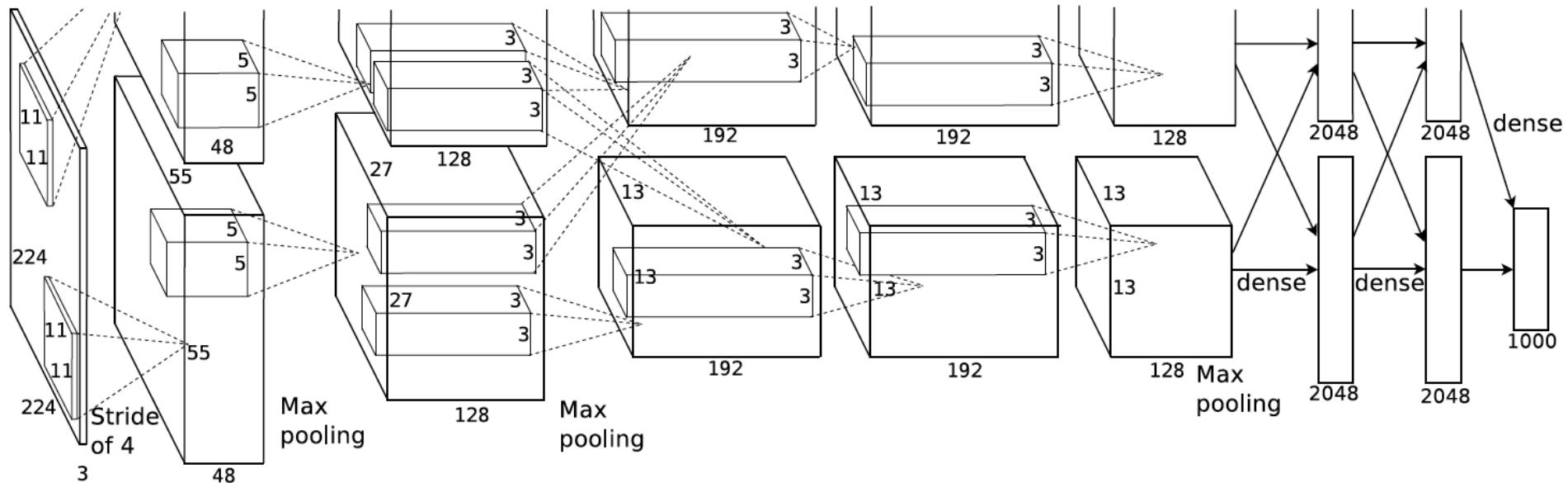
RECOGNIZING DIGITS

The LeNet-5 architecture as per Yann Lecun's original paper (1998)



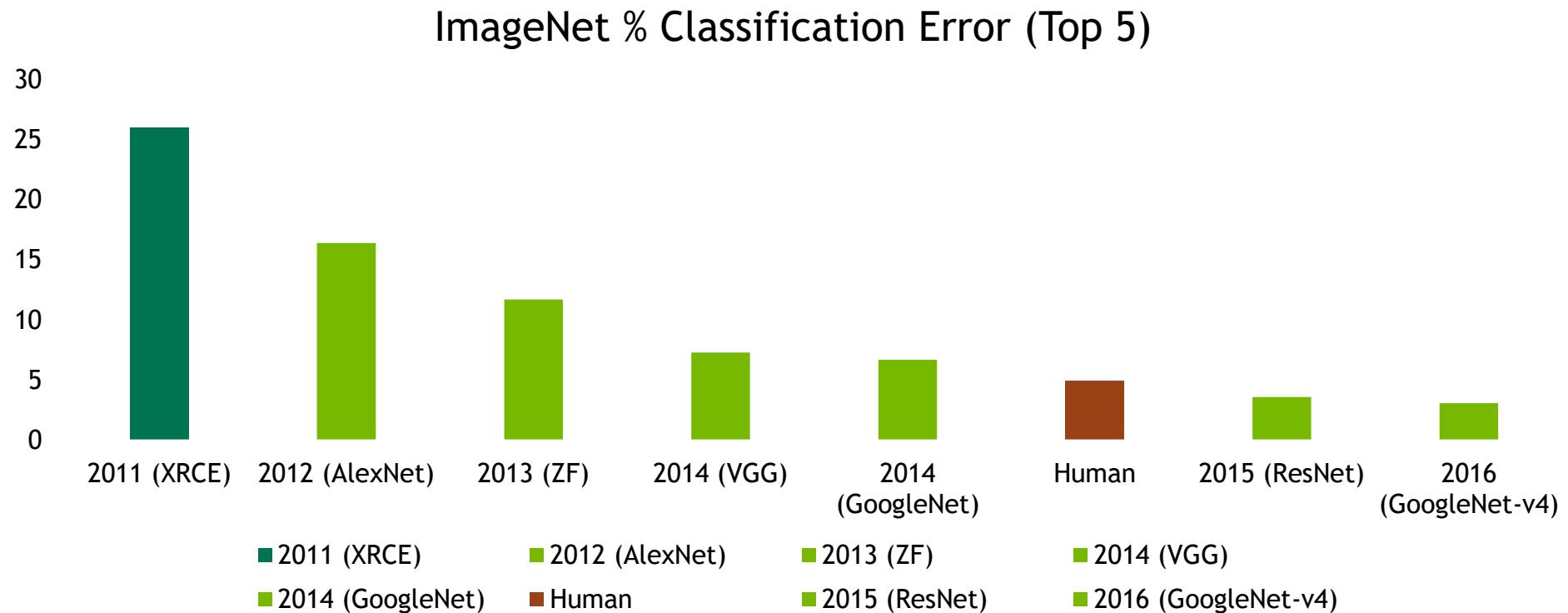
CLASSIFYING INTERNET IMAGES

The AlexNet architecture that won ImageNet (2012)



CLASSIFYING INTERNET IMAGES

Alex Krishevsky's 'AlexNet' spawned a new generation of DNNs



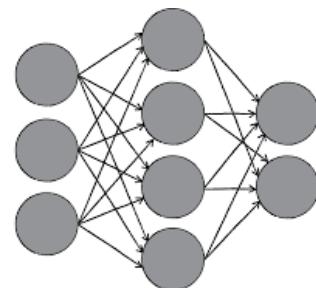
TRAINING A DNN

The three fundamental ingredients are the same for images or any other data

DATA



NETWORK

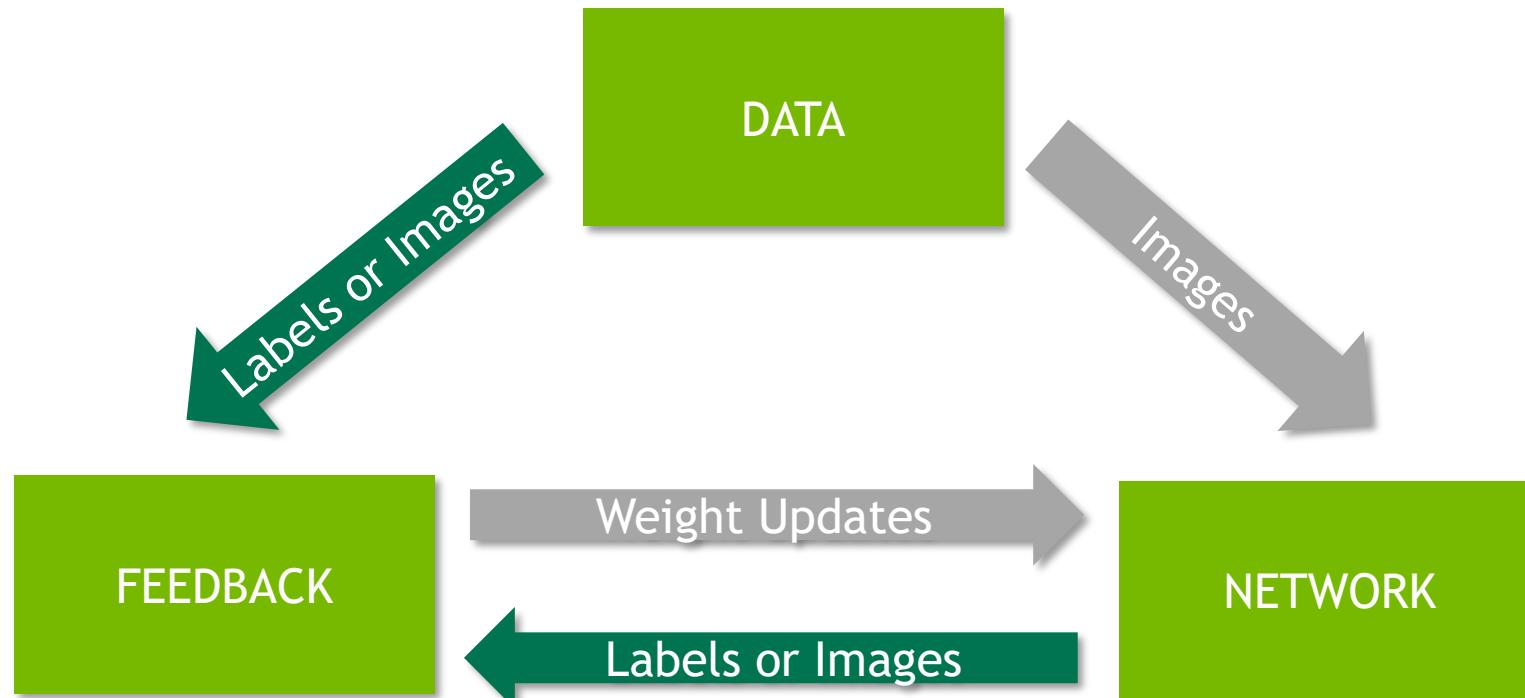


FEEDBACK



TRAINING A DNN

Classification, Segmentation & Detection



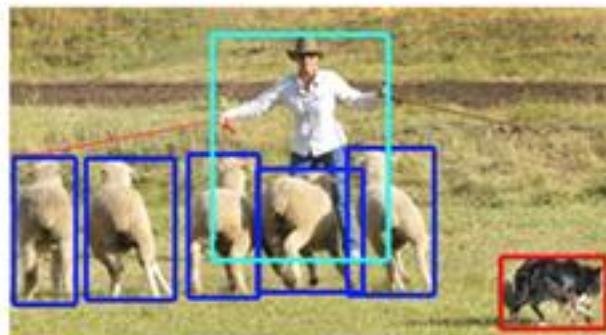
Compare Network
Output with
Expected Output

IMAGE DNN OUTPUTS

Classification, Detection & Segmentation require different levels of sophistication



(a) classification



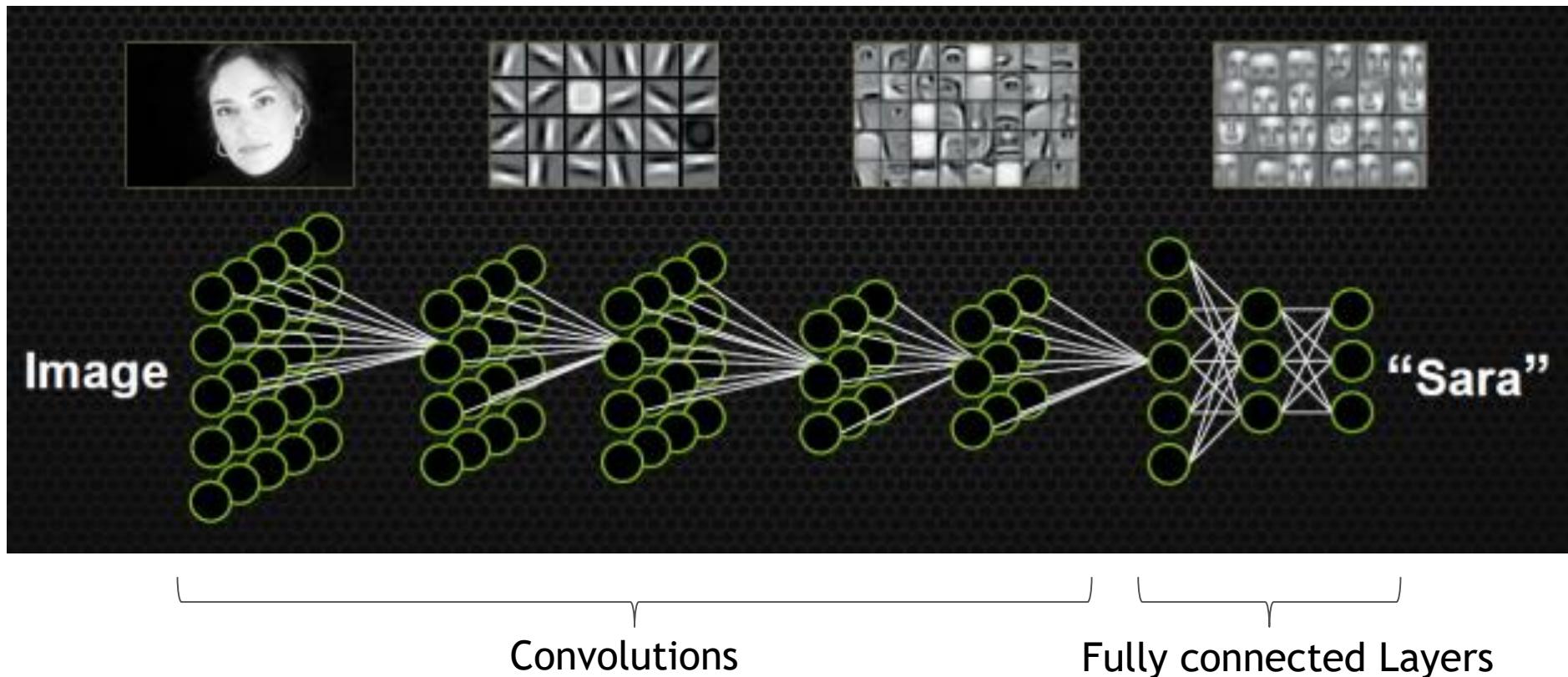
(b) detection



(c) segmentation

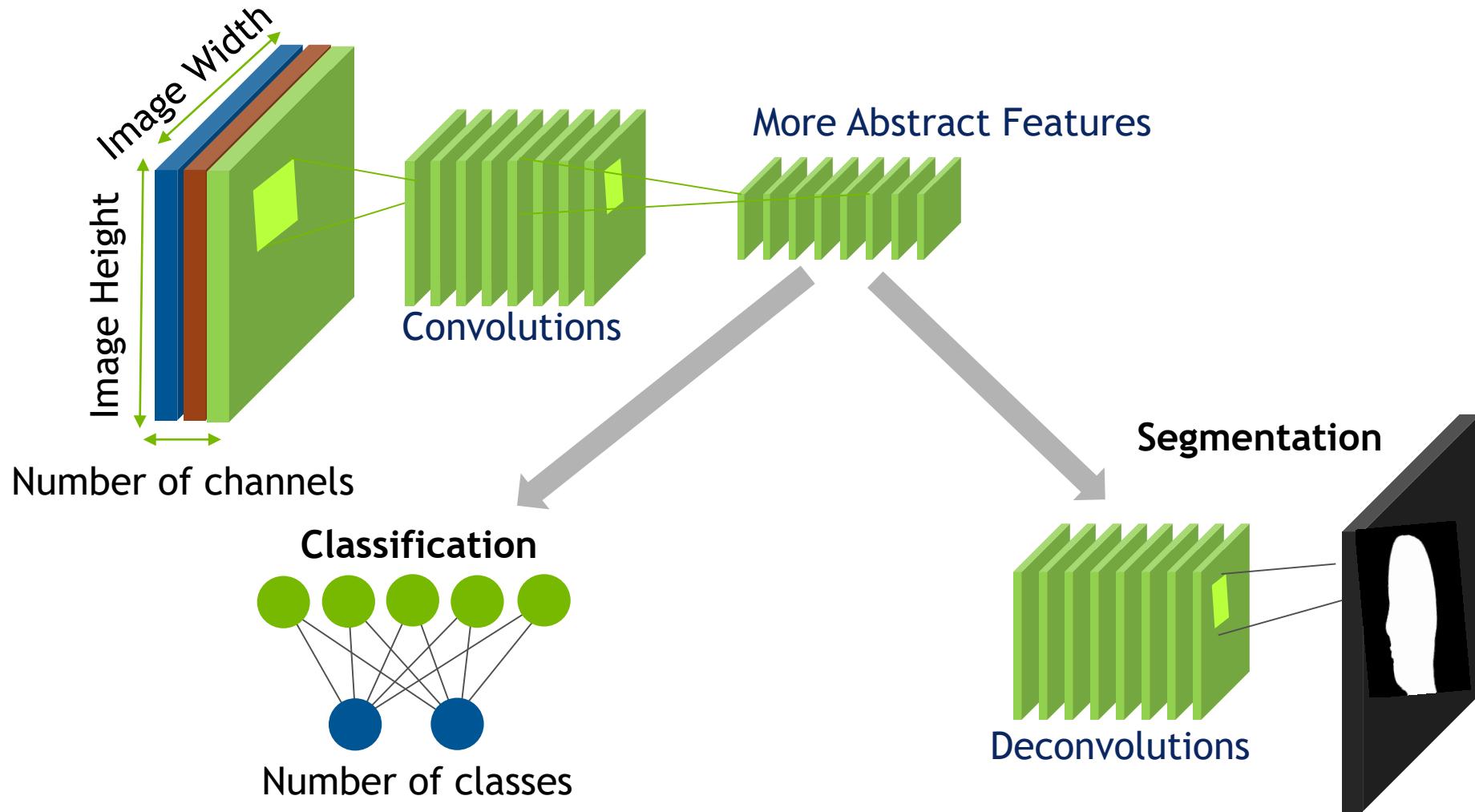
NETWORKS FOR IMAGES

Convolutions are the essential component, the rest depends on the task



NETWORKS FOR IMAGES

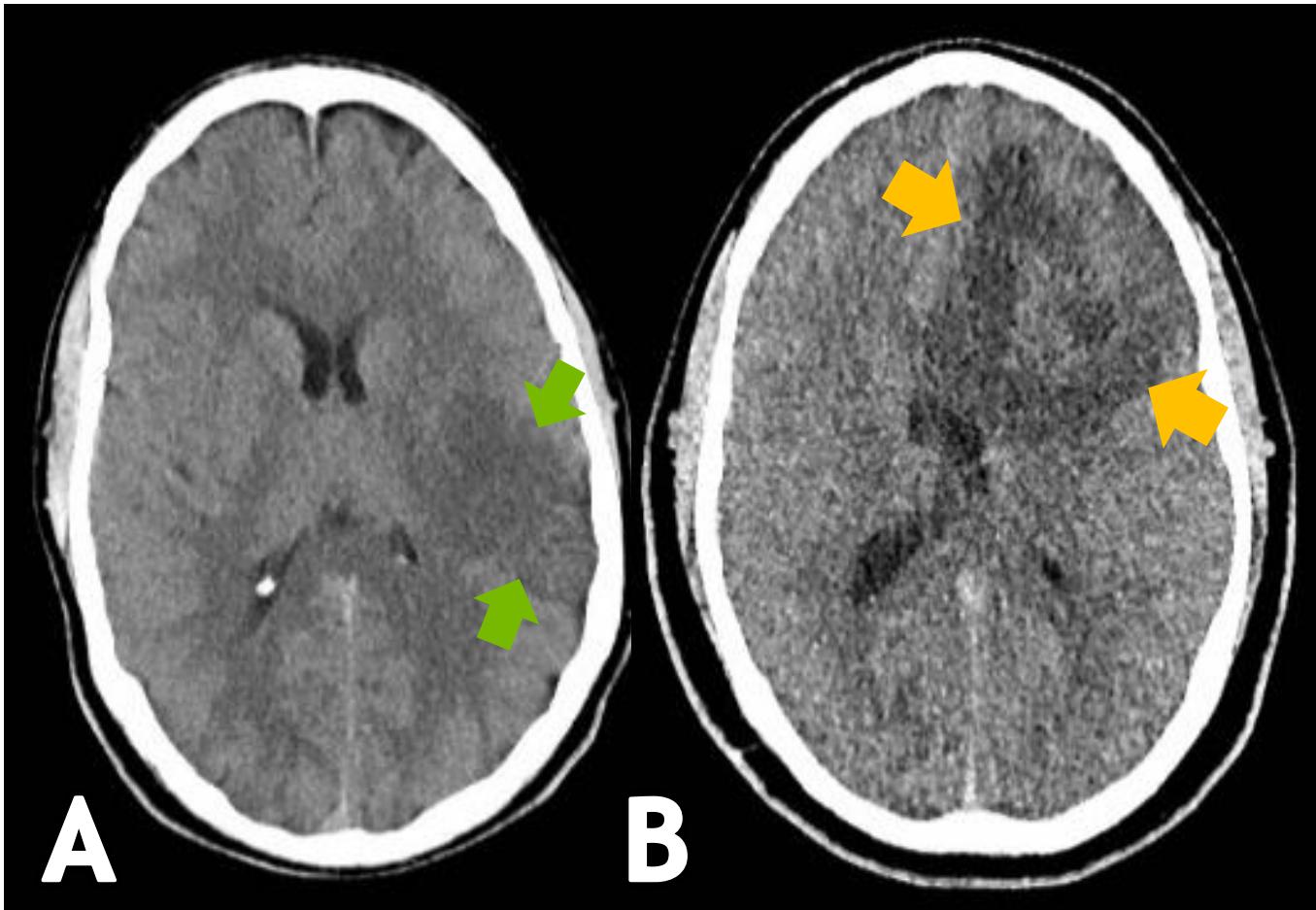
Convolutions are the essential component, the rest depends on the task



AI SPEEDS TIME TO CRITICAL CARE

The phrase “Time is Brain” means every minute counts after a stroke. A typical patient loses almost 2 million neurons per minute in which a stroke is untreated. Immediate treatment minimizes brain damage.

To help Radiologists diagnose the most urgent cases and enhance critical care, the OSU Department of Radiology used GPU-accelerated deep learning to develop an Automated Critical Test-Findings Identification and Online Notification System (ACTIONS). With GPUs, ACTIONS was trained in minutes vs. days. It identifies in seconds the most urgent cases of stroke, hydrocephalus, hemorrhage, and large tumors with an accuracy rate of 81% (stroke) and 91% (hydrocephalus, hemorrhage, large tumors), speeding time to critical care.



Examples of head CT examinations containing critical findings.

- A) A patient with a recent stroke involving the left cerebral hemisphere (green arrows).
- B) A patient with a large left frontal tumor compressing adjacent structures (orange arrows).



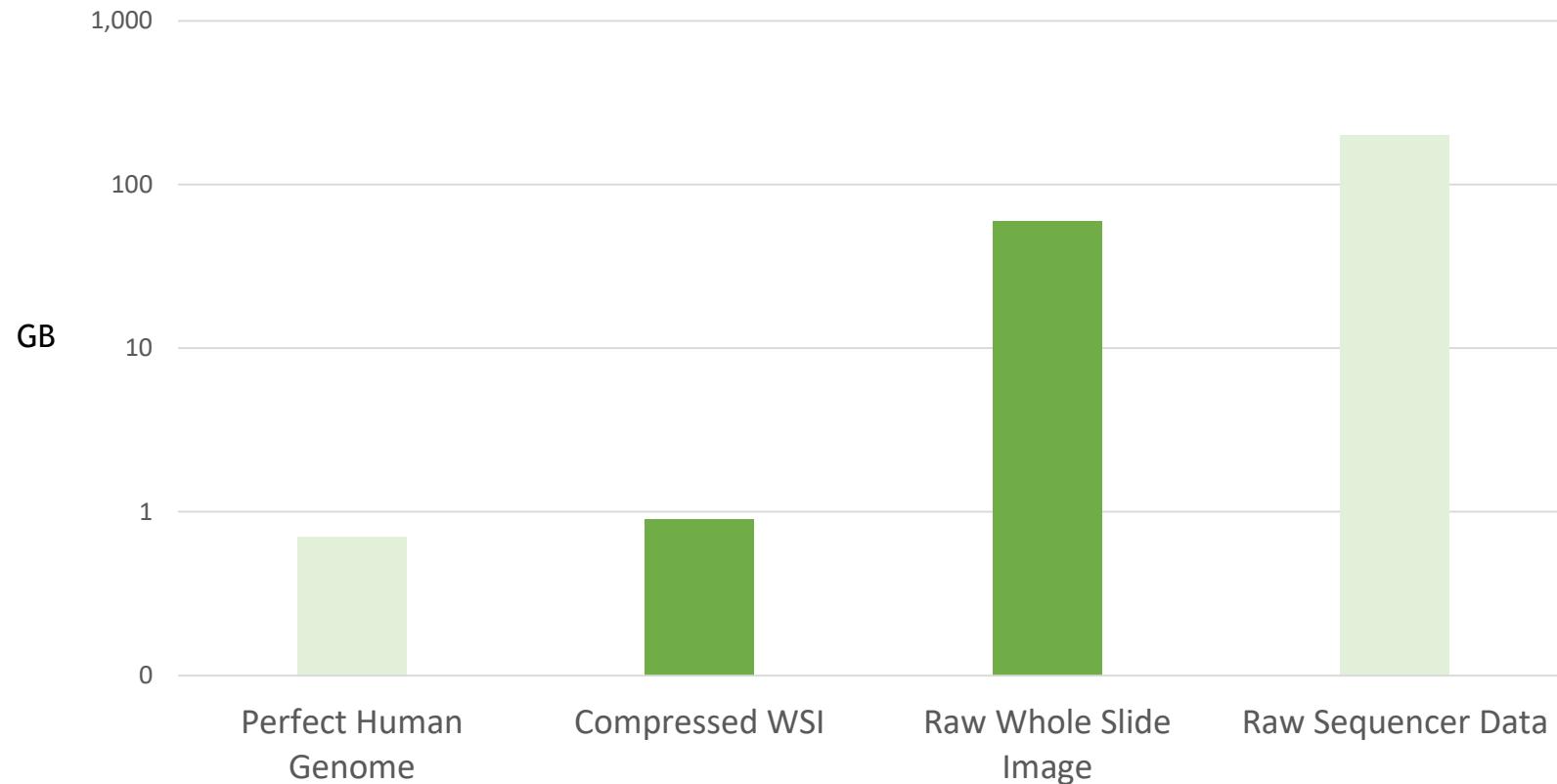
CHALLENGES IN DIGITAL PATHOLOGY

Although other modalities have their own data challenges, Whole Slide Images present a number of problems

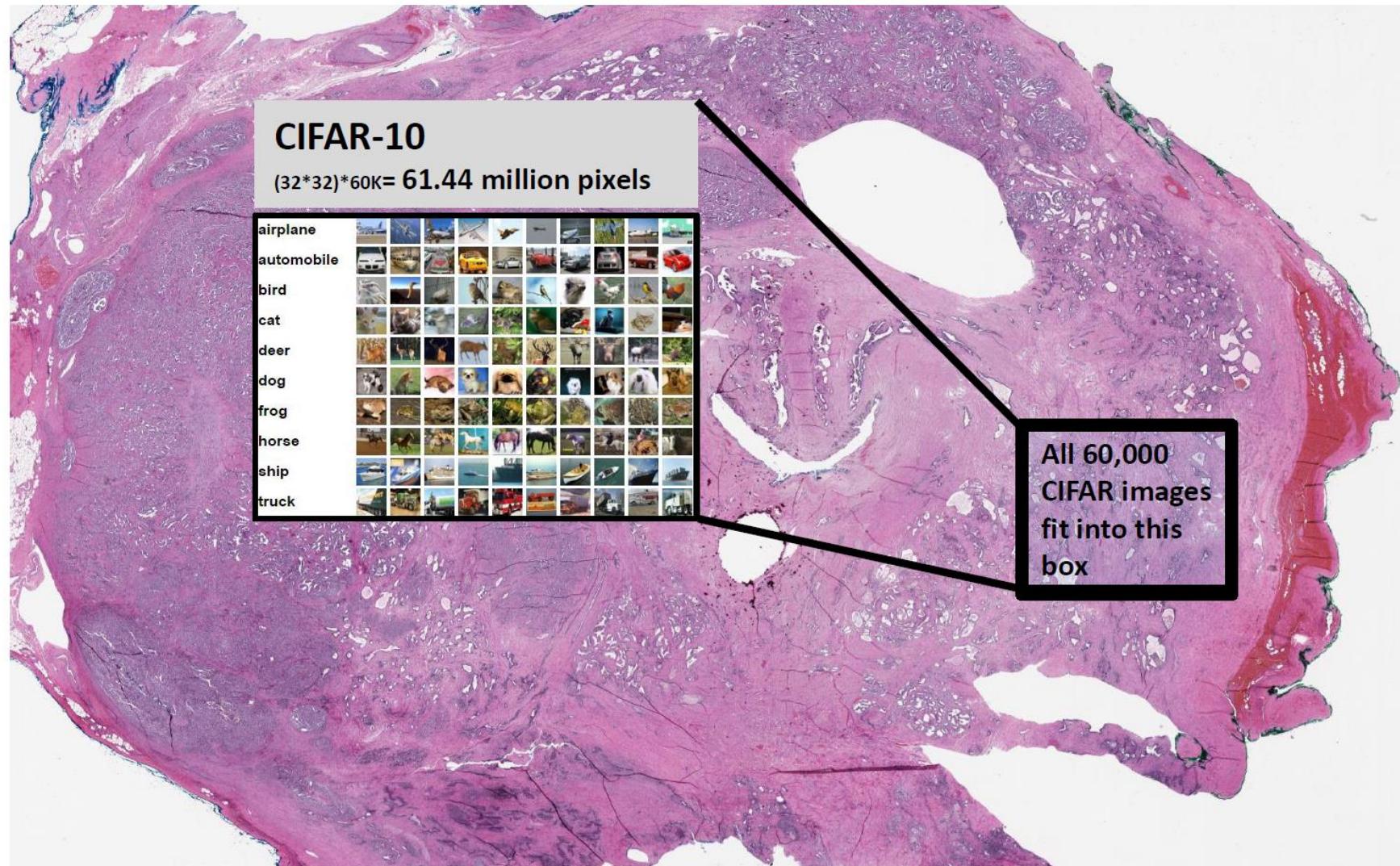
- VERY LARGE IMAGES
- Abnormalities may be limited to tiny islands in an ocean of ‘normal’ tissue
- Labels often only apply to whole image
- Limited public datasets and even fewer labelled ones
- A lot of background, which is not needed
- Only experts can identify abnormalities
- Subjective criteria used to classify - lots of inter (and intra) pathologist variability
- Lack of enthusiasm to digitise from many more ‘senior’ pathologists

DATA SIZES

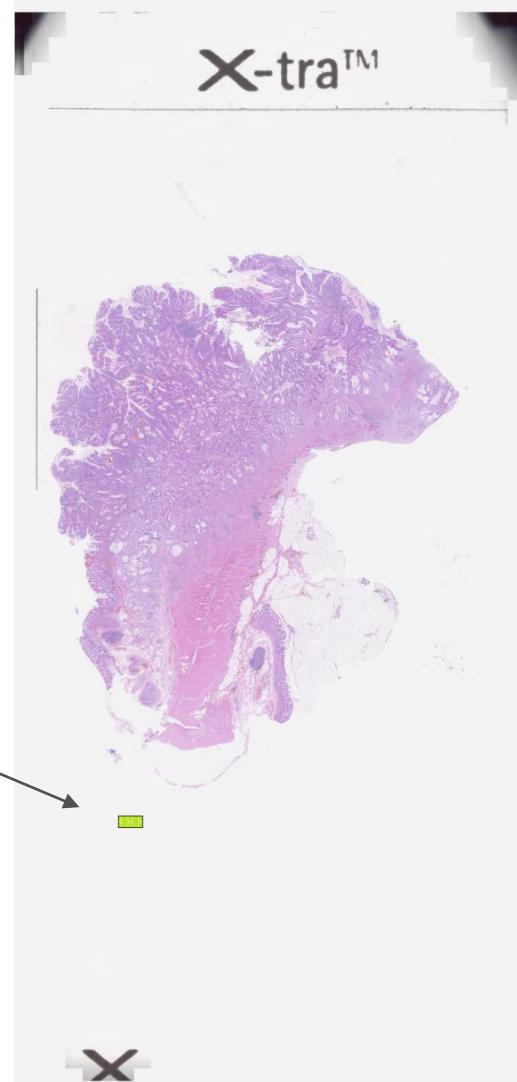
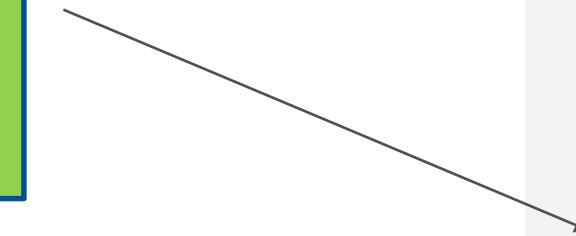
Genomics vs. Digital Pathology



1 Whole Slide
= 100,000 x 60,000
= 6 billion pixels

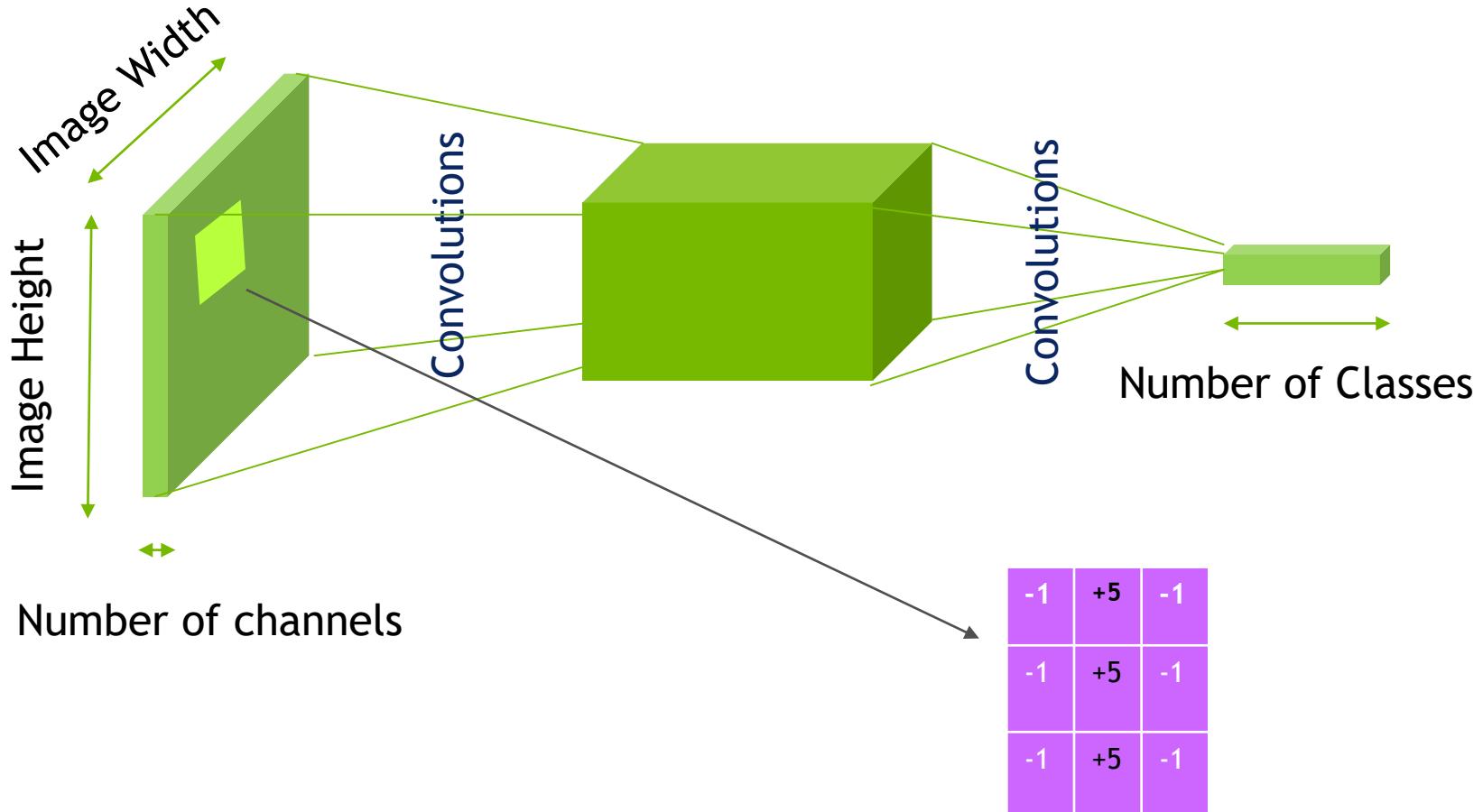


SIZE MATTERS



CONVOLUTIONAL NEURAL NETS

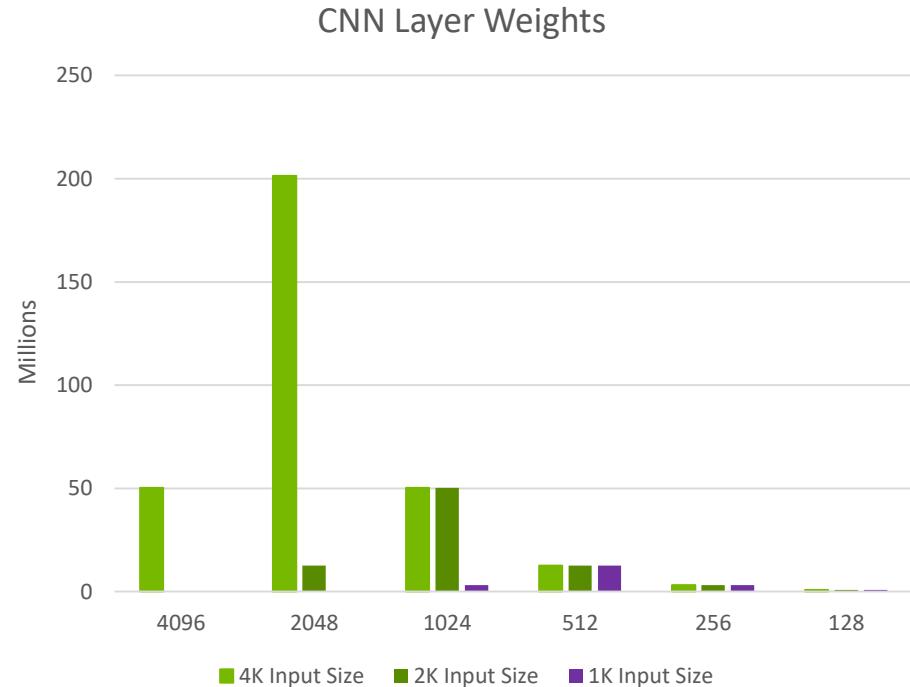
Why large images are problematic



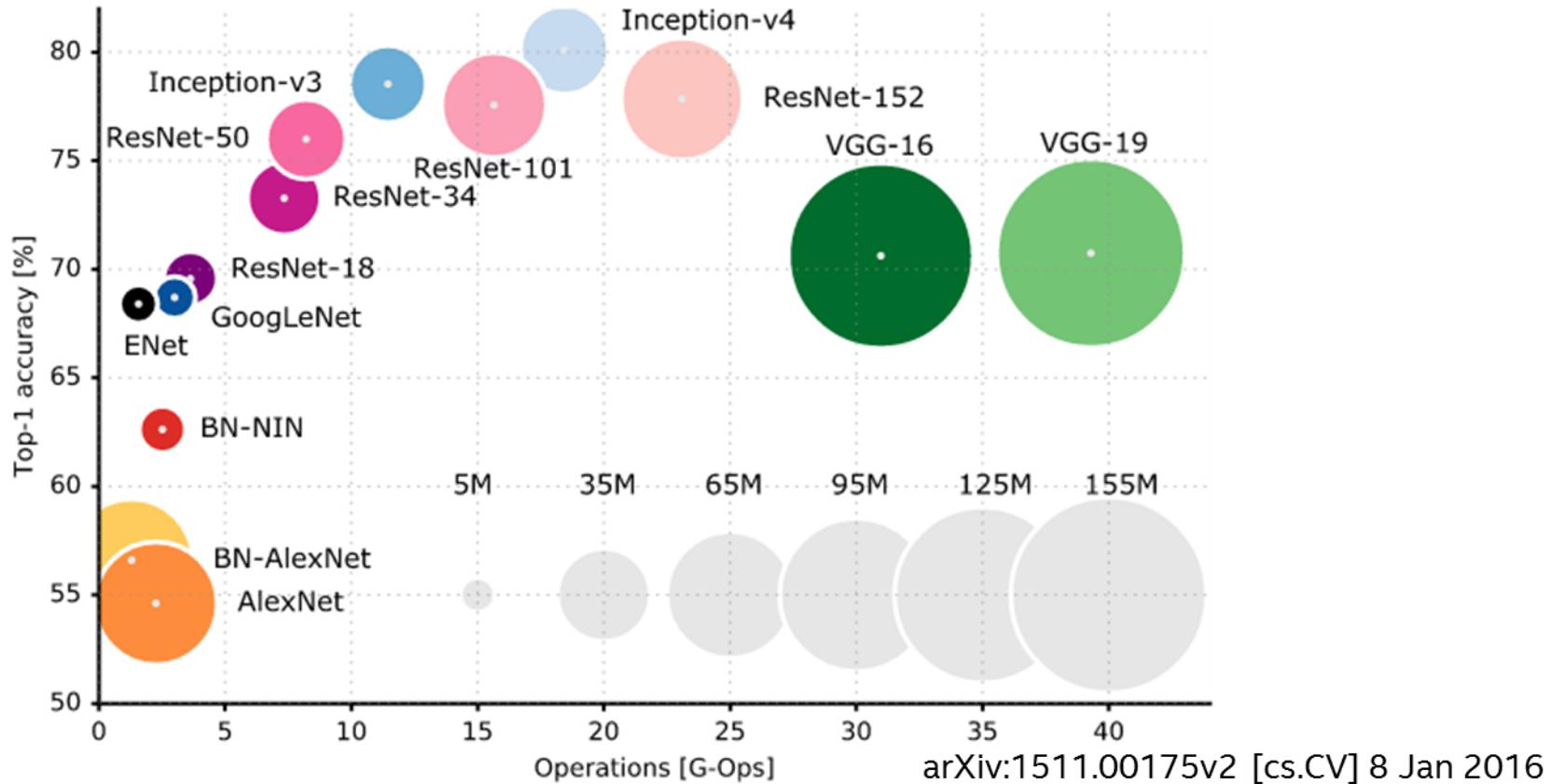
CONVOLUTIONAL NEURAL NETS

Adding convolution layers to ingest larger images increases memory exponentially

Input Size	Channels	Weights
4096	3	50M
2048	48	200M
1024	48	50M
512	48	12M
256	48	3M
128	48	790K

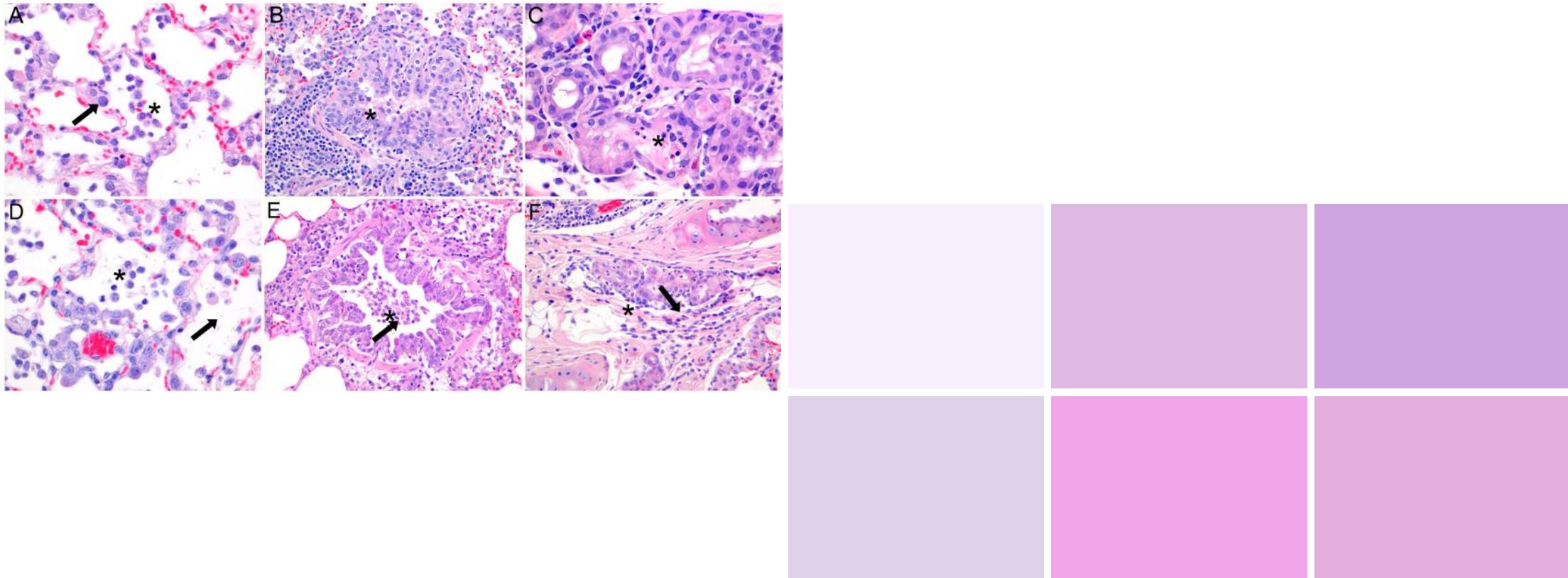


MODEL COMPLEXITY



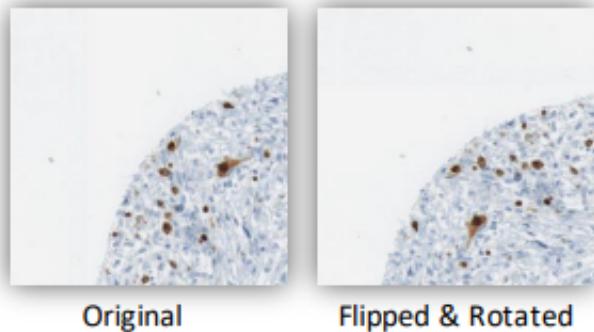
DOWN-SAMPLING

Down-sampling produces a significant loss of micro-features @ 40k to 512 pixels



EXPERT VARIABILITY

50 nuclei were repeated flipped and rotated to test the intra pathologist variability.

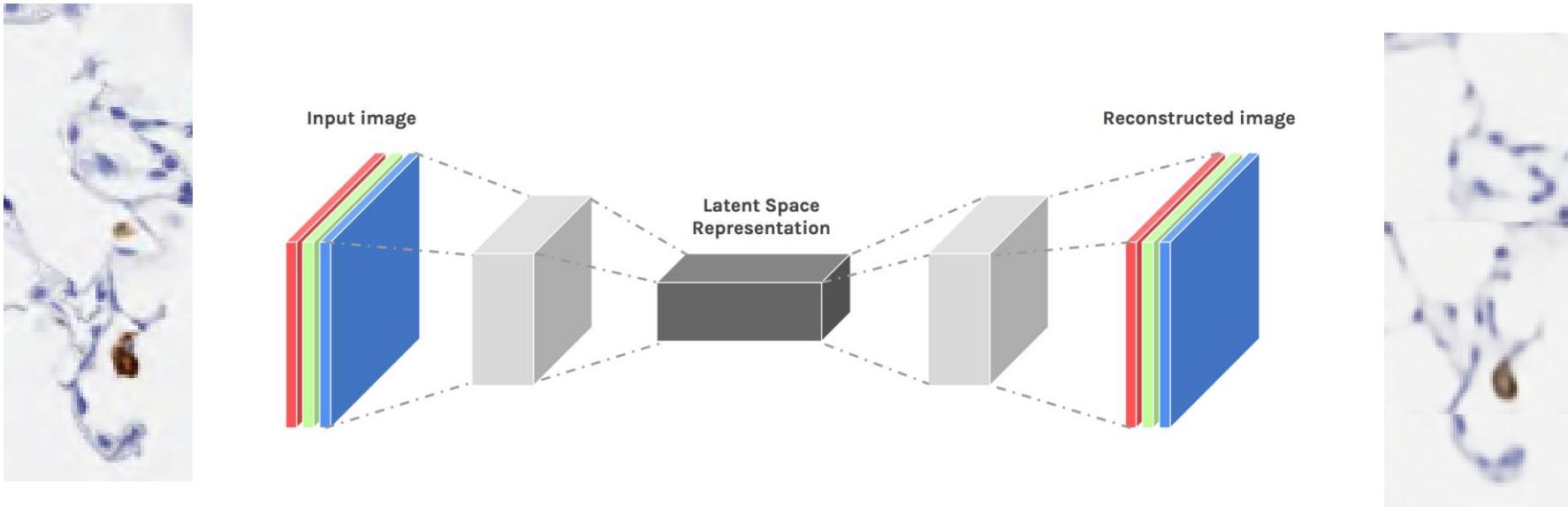


53/250
mismatches



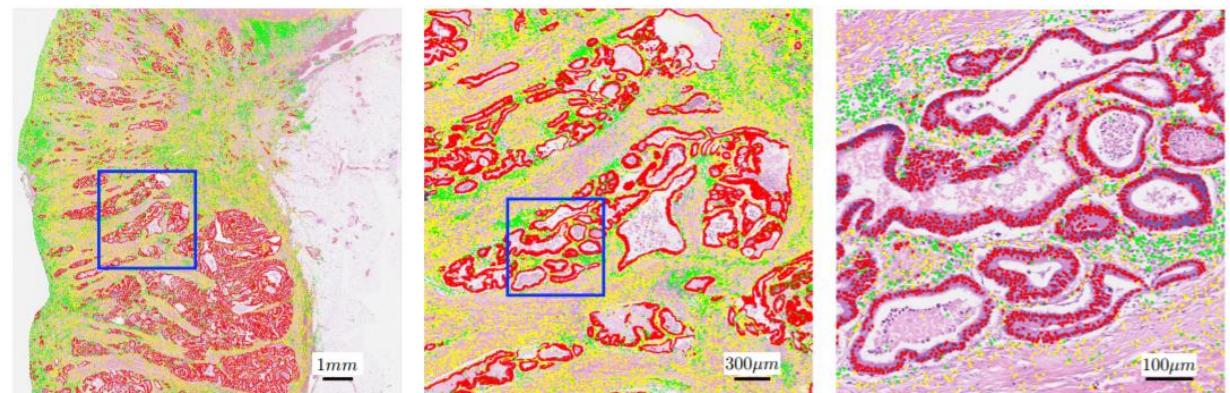
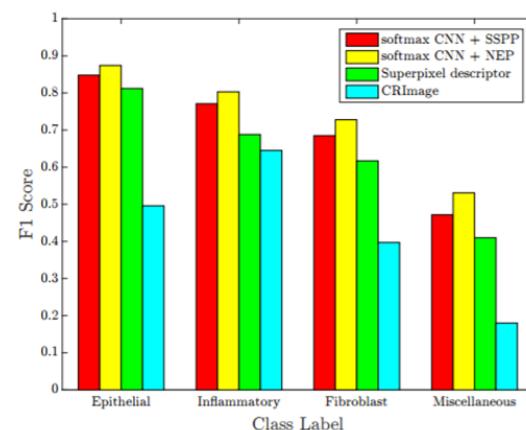
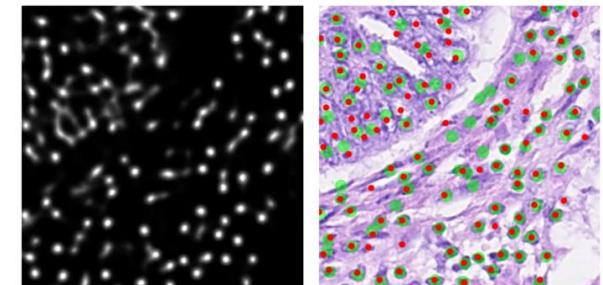
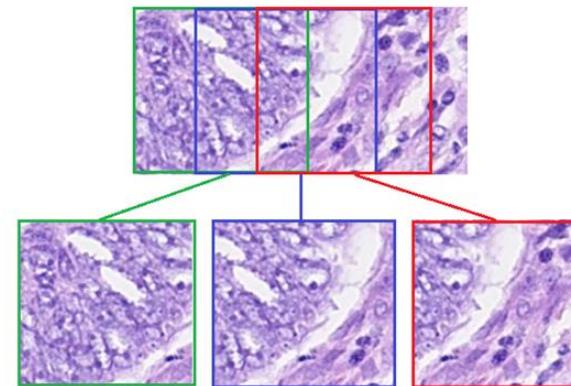
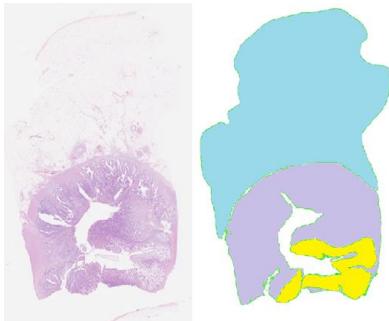
Baseline: Intra-Pathologists classification uncertainty of ~**20%**

SEMI-SUPERVISED DEEP LEARNING TO DETECT ABNORMALITIES





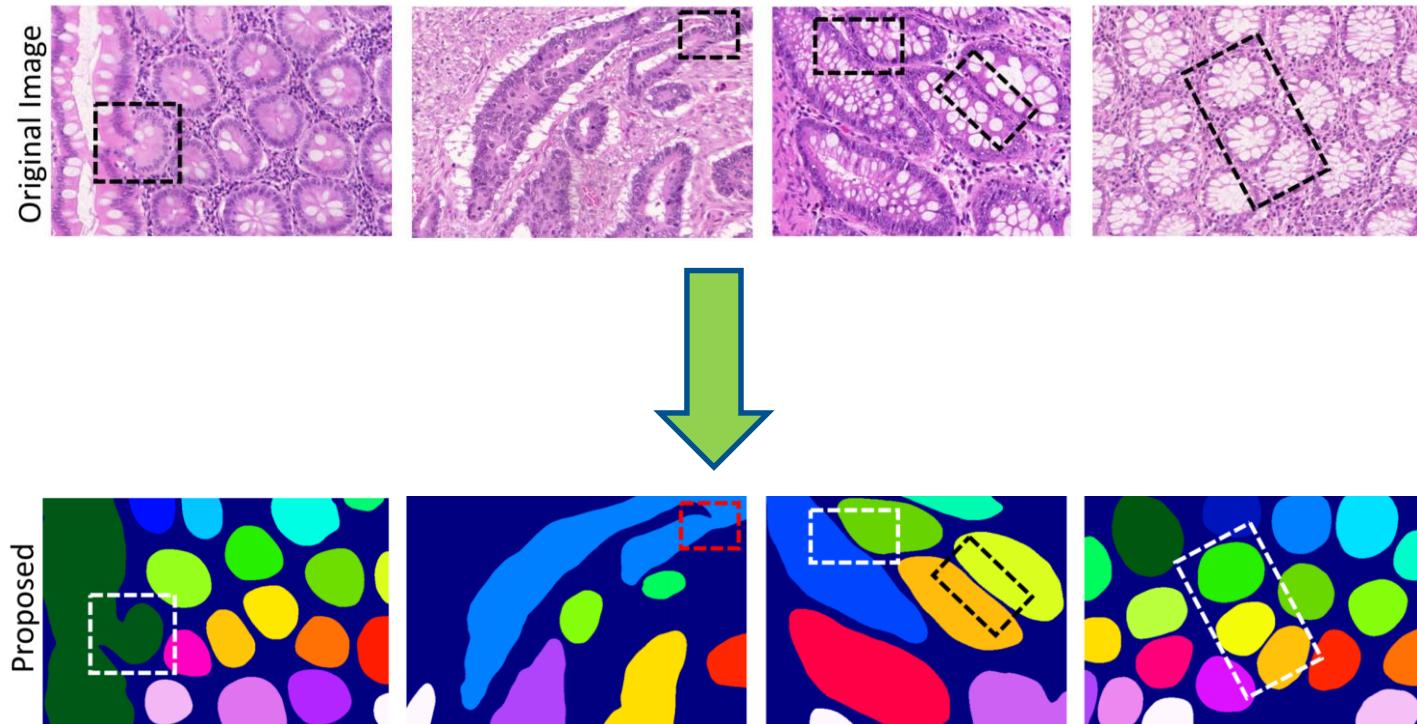
TISSUE IMAGE ANALYTICS LAB



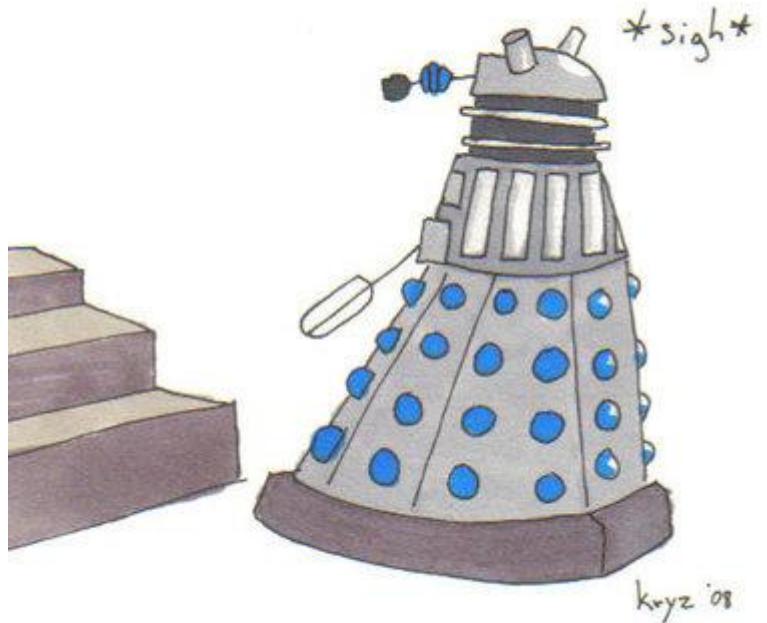


TISSUE IMAGE

ANALYTICS LAB



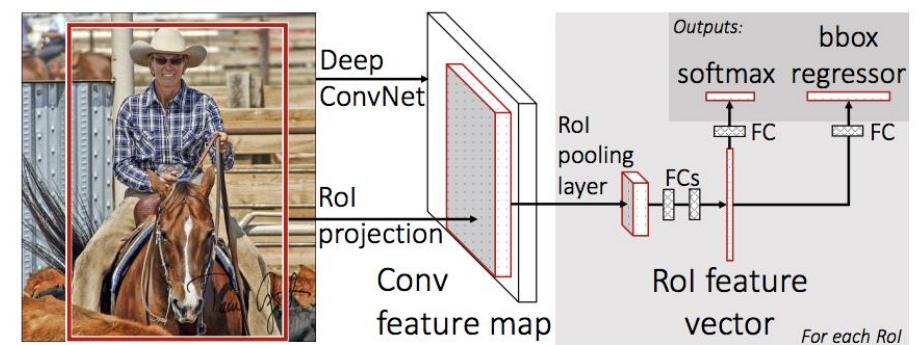
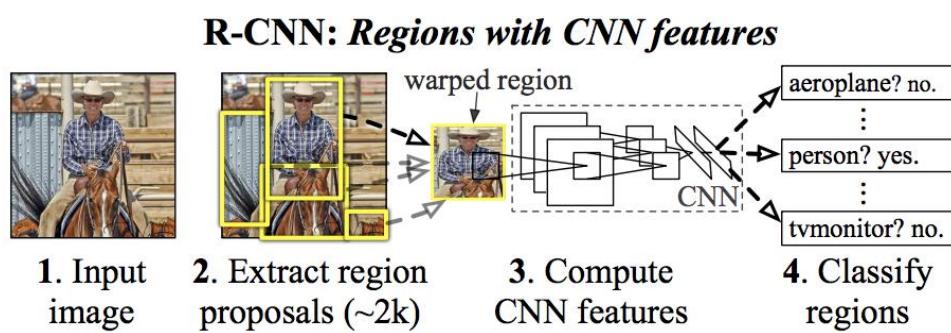
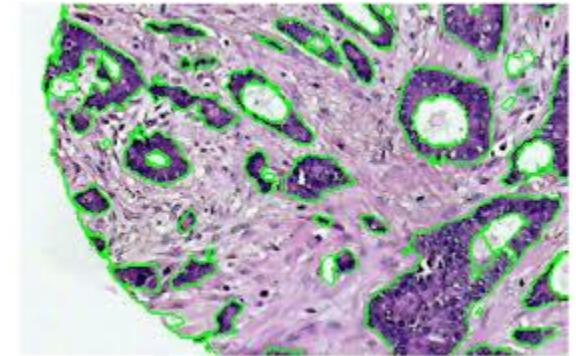
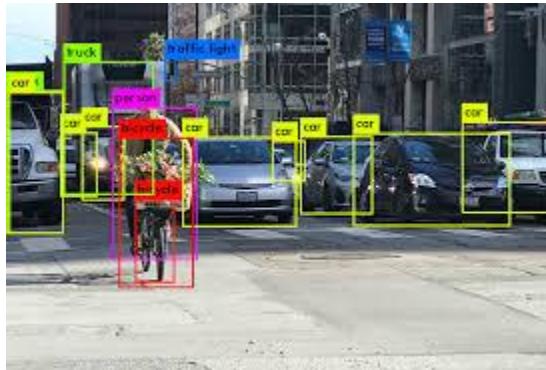
CHALLENGES



<https://www.deviantart.com/kryz-flavored/art/dalek-vs-stairs-99091404>

MULTIPLE INSTANCES

What if there are lots of objects on the same image?

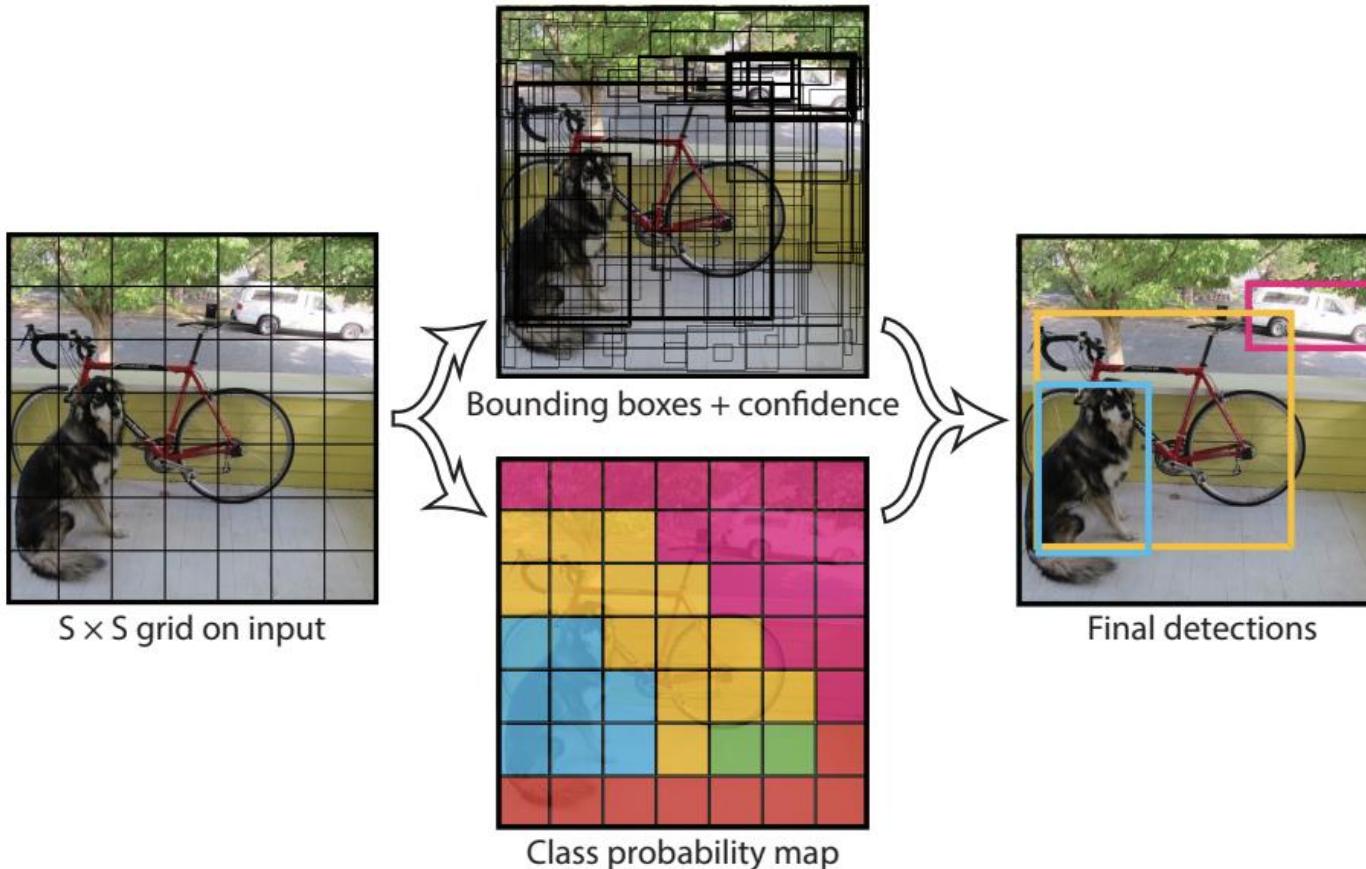


<https://arxiv.org/pdf/1311.2524.pdf>

<https://arxiv.org/pdf/1504.08083.pdf>

MULTIPLE INSTANCES

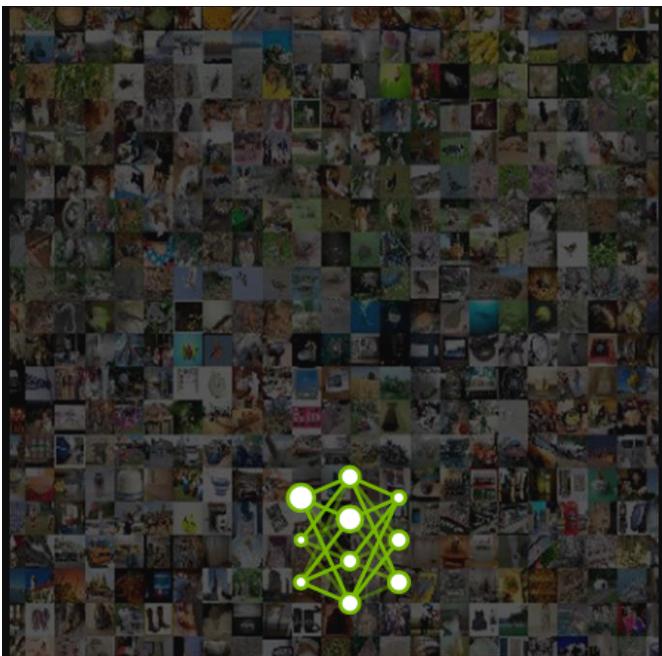
One recent and promising technique is called ‘YOLO’ (You only look once)



NEURAL NETWORK COMPLEXITY IS EXPLODING

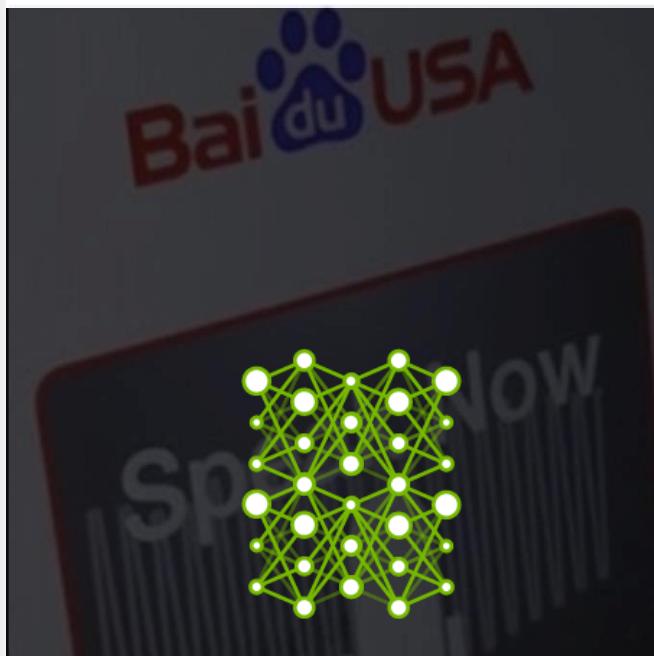
To Tackle Increasingly Complex Challenges

7 ExaFLOPS
60 Million Parameters



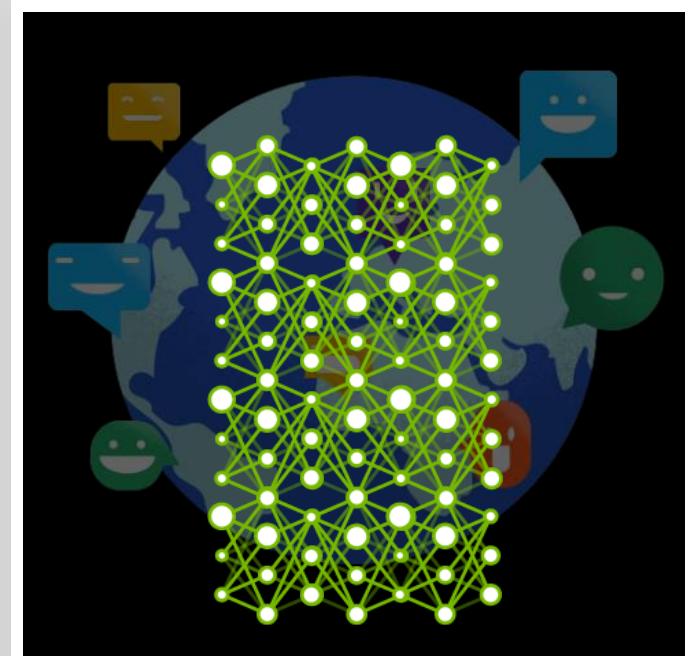
2015 - Microsoft ResNet
Superhuman Image Recognition

20 ExaFLOPS
300 Million Parameters



2016 - Baidu Deep Speech 2
Superhuman Voice Recognition

100 ExaFLOPS
8700 Million Parameters

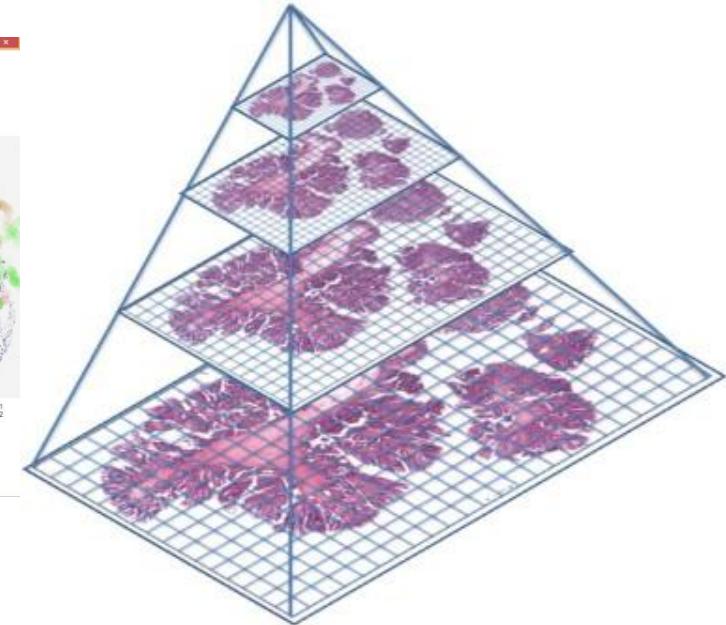


2017 - Google Neural Machine Translation
Near Human Language Translation

THE NEED FOR CONTEXT AND FOCUS...

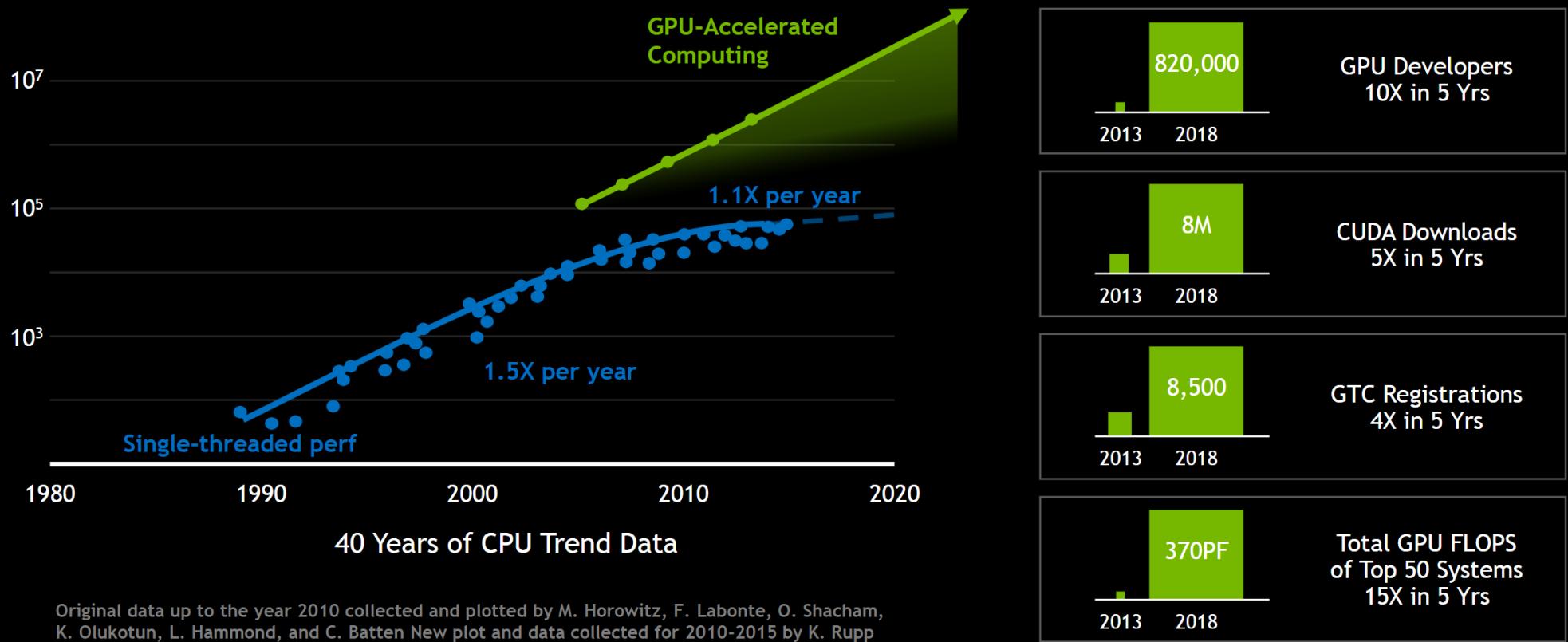


<https://blog.tcea.org/grants/>



<https://www.sciencedirect.com/science/article/pii/S1046202314002370>

RISE OF GPU COMPUTING



FUSING AI AND HIGH PERFORMANCE COMPUTING

5,120 CUDA cores

640 NEW Tensor cores

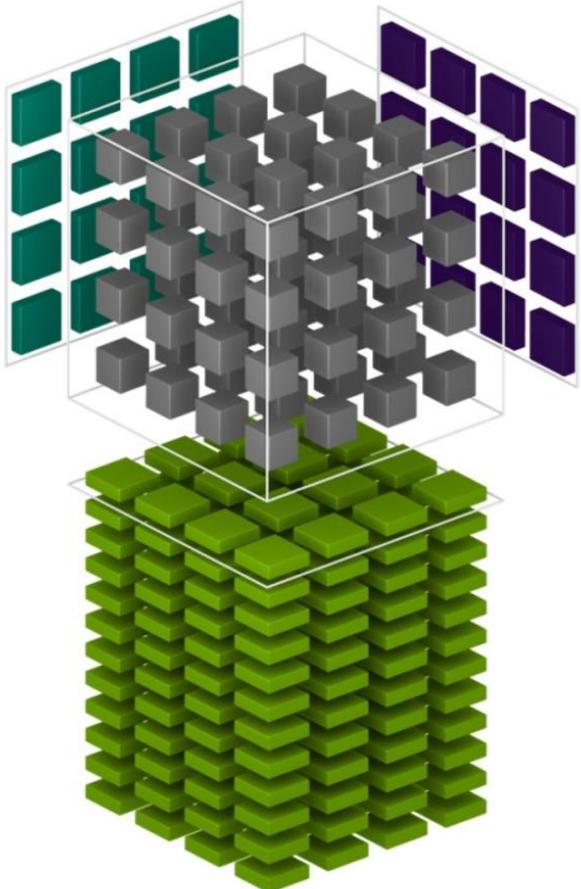
7.8 FP64 TFLOPS | 15.7 FP32 TFLOPS | 125 Tensor TFLOPS

20MB SM RF | 16MB Cache

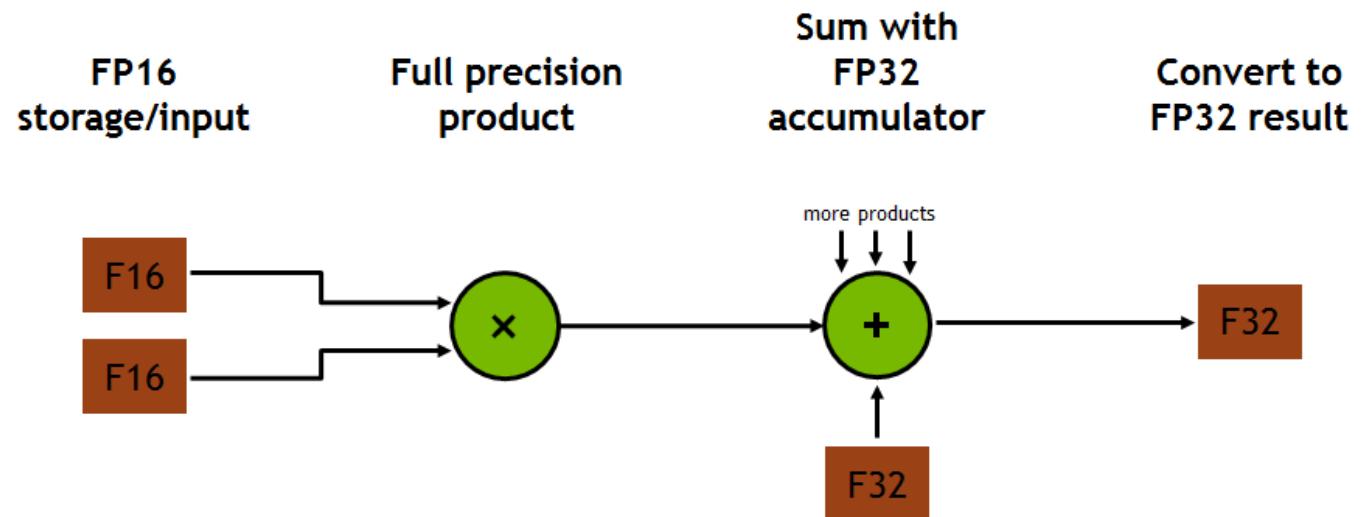
16GB/ 32GB HBM2 @ 900GB/s | 300GB/s NVLink



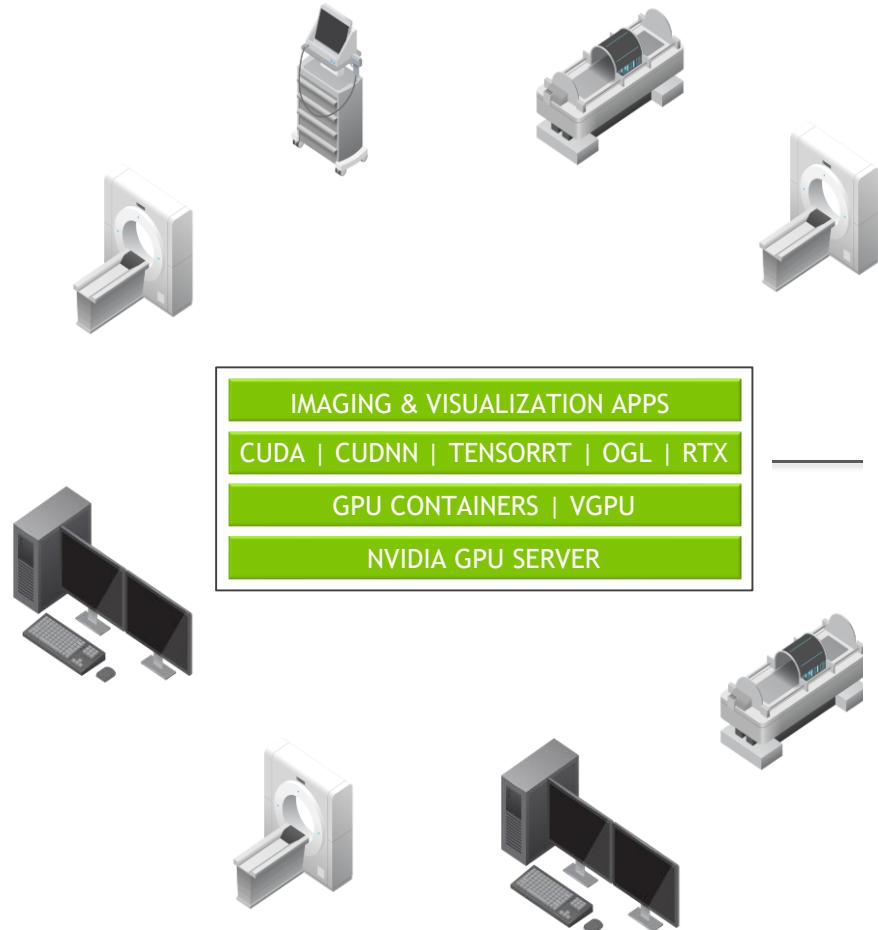
TENSOR CORES



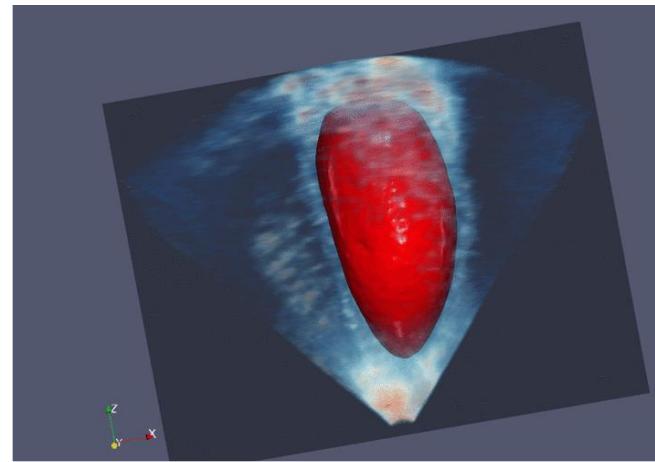
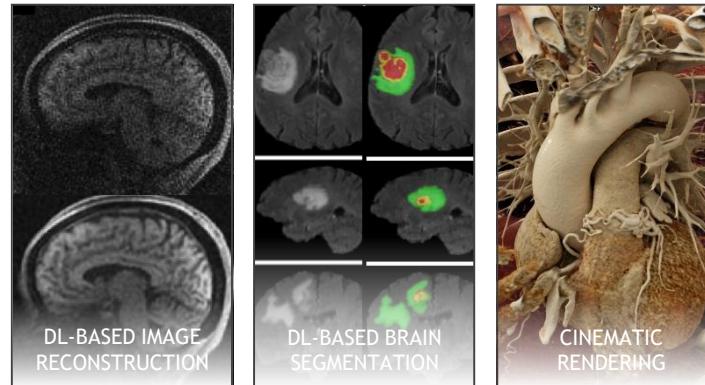
$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix}_{\text{FP16 or FP32}} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix}_{\text{FP16}} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}_{\text{FP16 or FP32}}$$

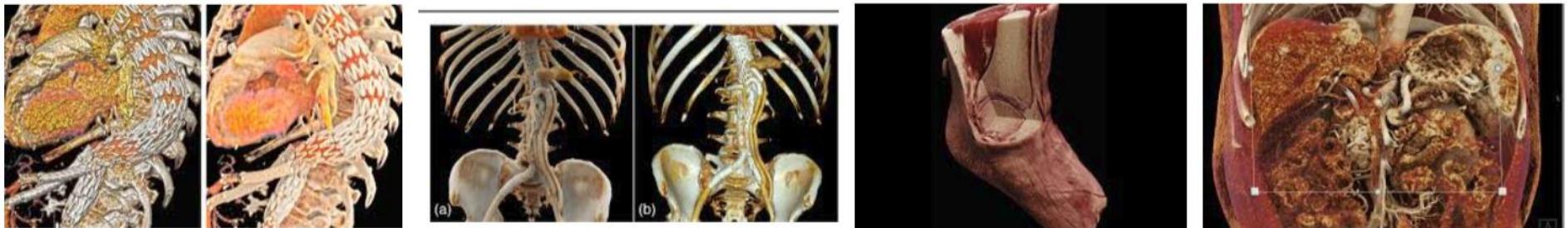


THE FUTURE



IMAGING & VISUALIZATION APPS
CUDA | CUDNN | TENSORRT | OGL | RTX
GPU CONTAINERS | VGPU
NVIDIA GPU SERVER





KEY TAKEAWAYS

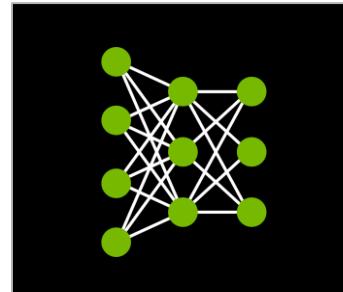
- Detection is really just a subtype of image classification
- CNNs are best for image-related work, but come in many flavours
- Imaging is not easy but is ripe for DL-based automation
- For large images, you need lots of GPU - with tensor cores !

NVIDIA DEEP LEARNING INSTITUTE

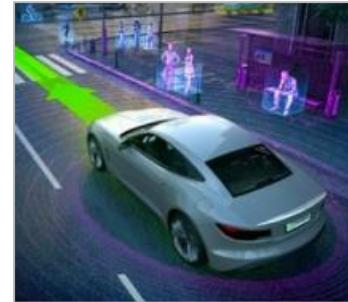
Hands-on, self-paced and instructor-led training in deep learning and accelerated computing for developers

Request onsite instructor-led workshops at your organization:
www.nvidia.com/requestdli

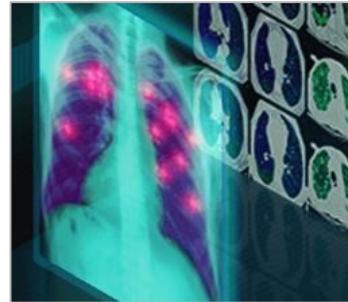
Take self-paced courses and electives online, view upcoming workshops, and learn about the University Ambassador Program: www.nvidia.com/dli



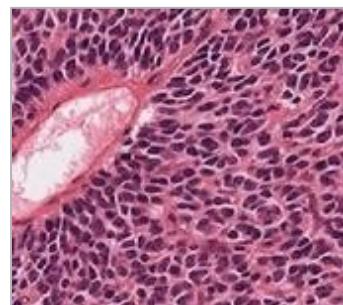
Deep Learning Fundamentals



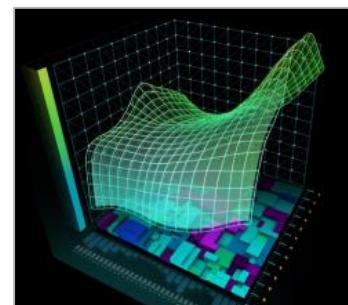
Autonomous Vehicles



Medical Image Analysis



Genomics



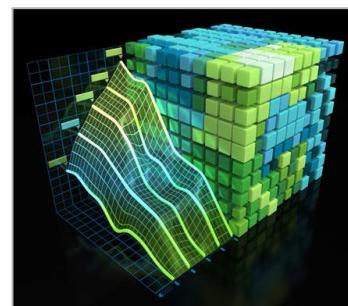
Finance



Digital Content Creation



Game Development



Accel. Computing Fundamentals

More industry-specific training coming soon...



QUESTIONS?

Thank you!

jhancox@nvidia.com