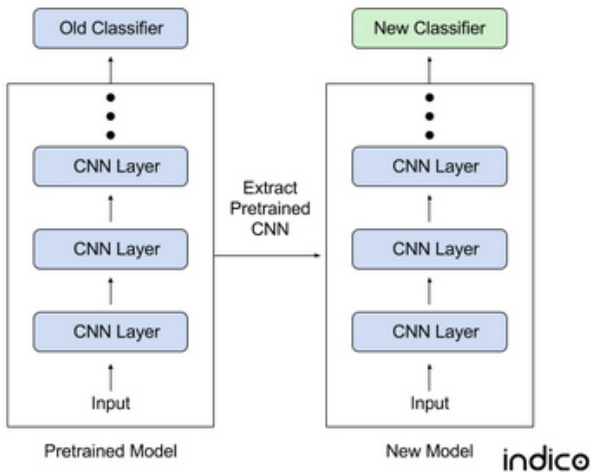# Deep Learning
## 9. Practical aspects of deep learning. Transfer learning. One shot learning. Hardware. Frameworks.

Viacheslav Dudar

Taras Shevchenko National University of Kyiv

2018

We assume that:

- Features extracted at early layers are general for all visual data
- Convolutional network forms a hierarchy of features: features in higher levels become more and more abstract
- To recognize different types of objects we need to change only the last layer of the network

# Tranfer learning use

- Small dataset, similar to original $->$ retrain last layer
- Large dataset, similar $->$ use original weights to initialize, finetune over the full net
- Small dataset, different $->$ train linear layer but earlier in the network
- Large dataset, different $->$ train from scratch (but weights initialization is still beneficial)
- If the input size is different, convolutional weights still can be used from another network
- Use smaller learning rates for finetuning

Small dataset: 1000 cats and 1000 dogs.

- Training from scratch: about 75-80% accuracy
- Finetuting fully connected layers of VGG: 90% accuracy
- Finetuting last convolutional and fully connected layers of VGG: 95%

Before NNs:

- Too small model $->$ underfitting
- Too big model $->$ overfitting
- Need to find optimal middle sized model and apply good regularization

- Divide into disjoint training, validation and test sets.
- Train on training set, check performance on validation set, modify model, repeat
- Use test set to estimate final performace

# How to modify model?

- High train error $->$ underfitting $->$ increase model size (increase layer size, add more layers)
- Low train error, high validation error $->$ overfitting $->$ increase regularization (dropout rate), use more data, use heavier data augmentation
- Low train error, low validation error $->$ done

Motivation:

- Deep learning works for large datasets. What to do with small datasets?
- Can we classify objects if we have only one example for each class?
- Application: face recognition

# Siamese network

|  |  |  |  |  |  |
|---|---|---|---|---|---|
|  |  | same | "cow" (speaker #1) | "cow" (speaker #2) | same |
|  |  | different | "cow" (speaker #1) | "cat" (speaker #2) | different |
|  |  | same | "can" (speaker #1) | "can" (speaker #2) | same |
|  |  | different | "can" (speaker #1) | "cab" (speaker #2) | different |

**Verification tasks (training)**

"cot" (speaker #4)   "cob" (speaker #4)   "cog" (speaker #4)

?     ?     ?

"cob" (speaker #3)

**One-shot tasks (test)**

# Siamese network



| Input layer | Hidden layer | Distance layer | Output layer |

$x_{1,1}$  $w_{1,1}^{(1)}$  $h_{1,1}$

$w_{3,1}^{(1)}$

$w_{1,N_1}^{(1)}$

$x_{1,N_1}$  $w_{3,N_1}^{(1)}$  $h_{1,N_2}$

$d_1$  $w_{3,1}^{(2)}$

$p$

$w_{3,N_2}^{(2)}$

$x_{2,1}$  $w_{1,1}^{(1)}$  $h_{2,1}$

$w_{3,1}^{(1)}$

$w_{1,N_1}^{(1)}$

$d_{N_2}$

$x_{2,N_1}$  $w_{3,N_1}^{(1)}$  $h_{2,N_2}$

# Siamese network

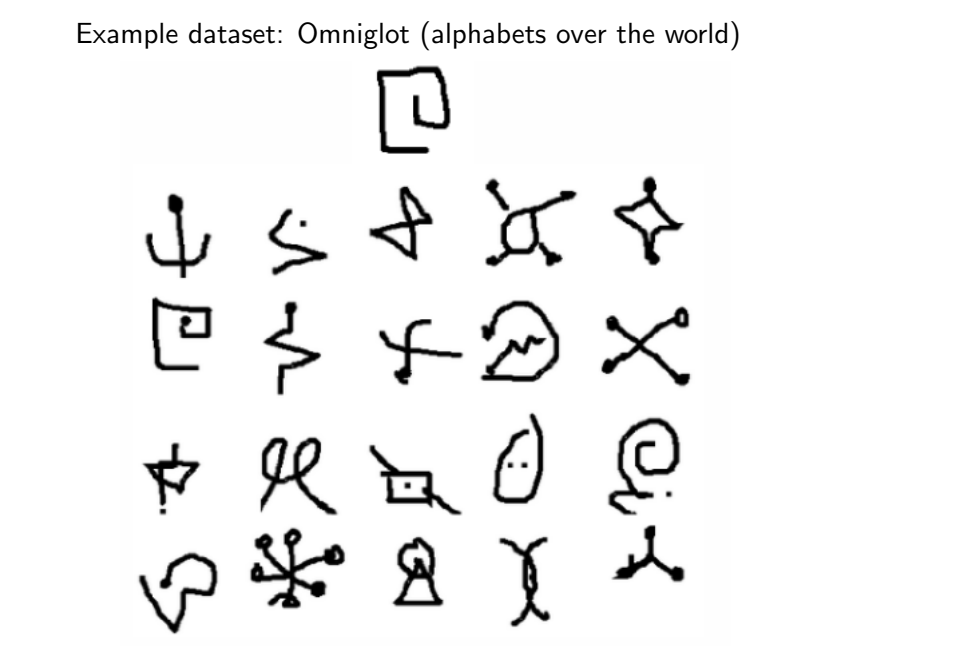- Generate feature vector for input image with convolutional network
- Maximize distance between feature vectors for different classes and minimize for the same classes
- If there is only one example for each class, minimize distance betwwen images and their affine transformations, subimages, and so on

# Example

Example dataset: Omniglot (alphabets over the world)

# Siamese network

| Method | Test |
| --- | --- |
| **Humans** | 95.5 |
| **Hierarchical Bayesian Program Learning** | 95.2 |
| **Affine model** | 81.8 |
| **Hierarchical Deep** | 65.2 |
| **Deep Boltzmann Machine** | 62.0 |
| **Simple Stroke** | 35.2 |
| **1-Nearest Neighbor** | 21.7 |
| **Siamese Neural Net** | 58.3 |
| **Convolutional Siamese Net** | 92.0 |

- Building blocks for designing and training neural nets
- High level programming interface
- Choose framework based on: ease of programming, running speed, support of GPUs
- Main function of most frameworks: automatic differentiation of the model

Many more...

# Some of deep learning frameworks

Feel free to explore wide variety of deep learning frameworks!

Keras (python): simple to use. Contains most popular layers, optimizers, pretrained networks. Good for exploiting already existing architectures.

PyTorch (python): (Python, former Torch, Lua language): able to create dynamical computational graphs. Good for trying own layers.

Theano: (Keras with Theano backend is faster on single GPU compared with Tensorflow backend)

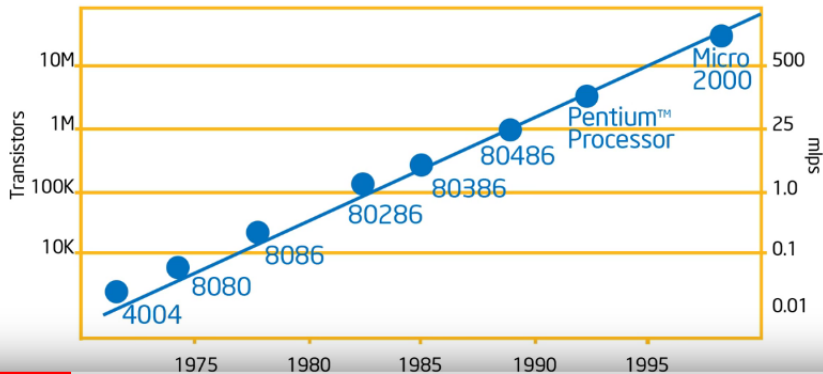Caffe: (C++, one of the fastest convnets implementations)

MXNET (Apache): supported by AWS, Azure

CNTK (Microsoft Cognitive Toolkit): optimized for Azure, integrated with Visual Studio.

Why parallel computing?



Moore's Law

The Good Old Days

??%/year

(SPECint)
Uniproccessor
Performance

Performance (vs. VAX-11/780)

Pentium 4, 3.0 GHz,
20 stage, 3 CISC
issue (6 uop issue)

Pentium 4, 3.6 GHz,
31 stage, 6 CISC
issue, 3 CISC issue

52%/year

Sparc V7 RISC
5-stage
Sun 4/260
16.7 MHz

PowerPC 604,
100 MHz
7 stage, 4 issue

25%/year

Vax "Nautilus",
CISC, Vax 8700

Vax "Star.,
CISC
Vax-11/780

From Hennessy and Patterson, *Computer Architecture: A Quantitative Approach*, 4th edition, Sept. 15, 2006

# Power wall



Partial Solution: Simple Low Power Cores

power = perf ^ 1.74

Mobile CPUs With Shallow Pipelines Use Less Power

Parallel architecture solves power problem at expense of need to parallelize computation.



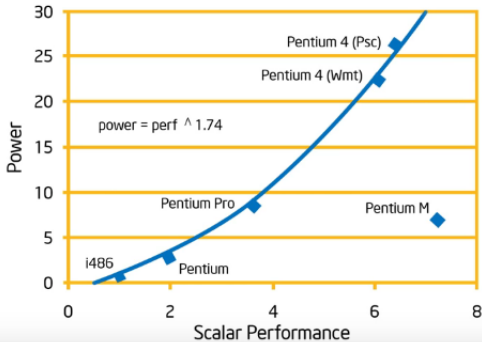Capacitance = 2.2C
Voltage = 0.6V
Frequency = 0.5f
Power = $0.396CV^2F$

CPU: efficient for few very complex instructions
GPU: efficient for many simple instructions



*The Difference between a CPU and GPU*

CPU          GPU

GPUs were used for rendering images on screen (performing in parallel many simple operations, like matrix multiplications) - this is exactly what we need for training neural net!

# GPGPU

GPGPU: General-purpose computing on graphics processing units
OpenCL (Open Computing Language): open source GPGPU
framework
CUDA: proprietary (NVIDIA) GPGPU framework (works faster)
Both are C-like.

# GPU-CPU comparison

# GPU – CPU servers comparison



Theoretical Peak Floating Point Operations per Clock Cycle, Double Precision

Legend:
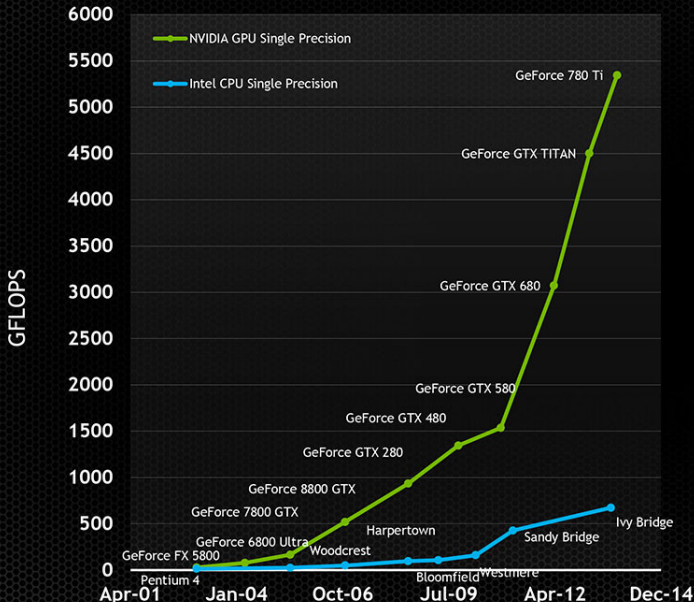- INTEL Xeon CPUs
- NVIDIA Tesla GPUs
- AMD Radeon GPUs
- INTEL Xeon Phis

NVIDIA Tesla GPUs: Tesla C1060, Tesla C1060, Tesla C2050, Tesla M2090, Tesla K20, Tesla K20X, Tesla K40, Tesla K40, Tesla P100

AMD Radeon GPUs: HD 3870, HD 4870, HD 5870, HD 6970, HD 6970, HD 7970 GHz Ed., HD 8970, FirePro W9100, FirePro S9150

INTEL Xeon CPUs: X5482, X5492, W5590, X5680, X5690, E5-2690, E5-2697 v2, E5-2699 v3, E5-2699 v3, E5-2699 v4

INTEL Xeon Phis: Xeon Phi 7120 (KNC), Xeon Phi 7290 (KNL)

Axes: FLOPs per Clock Cycle (vertical, $10^1$ to $10^4$) vs End of Year (horizontal, 2008, 2010, 2012, 2014, 2016)

# Modern GPU



## 2017: TESLA VOLTA V100

21B transistors
815 mm²

80 SM
5120 CUDA Cores
640 Tensor Cores

16 GB HBM2
900 GB/s HBM2
300 GB/s NVLink

*full GV100 chip contains 84 SMs

# Modern GPU

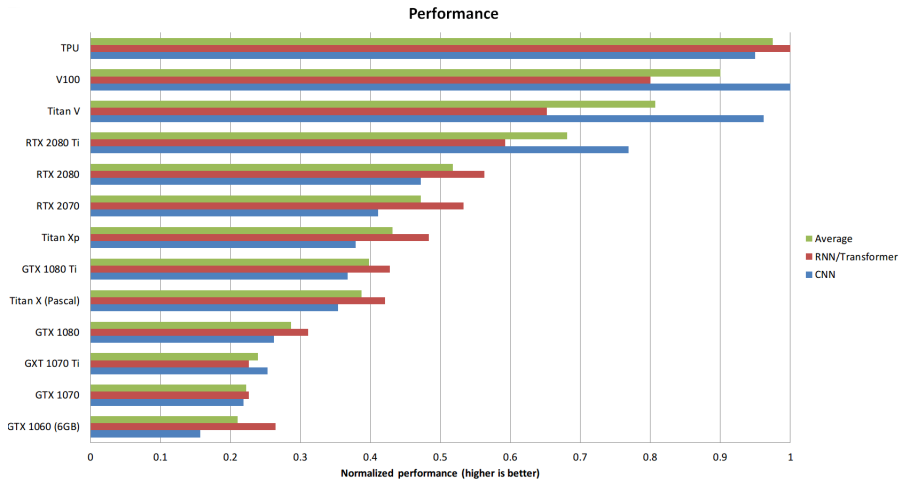|  | Tesla V100 for NVLink | Tesla V100 for PCIe |
|---|---|---|
| **PERFORMANCE**<br>with NVIDIA GPU Boost™ | DOUBLE-PRECISION<br>7.8 teraFLOPS<br><br>SINGLE-PRECISION<br>15.7 teraFLOPS<br><br>DEEP LEARNING<br>125 teraFLOPS | DOUBLE-PRECISION<br>7 teraFLOPS<br><br>SINGLE-PRECISION<br>14 teraFLOPS<br><br>DEEP LEARNING<br>112 teraFLOPS |
| **INTERCONNECT BANDWIDTH**<br>Bi-Directional | NVLINK<br>300 GB/s | PCIE<br>32 GB/s |
| **MEMORY**<br>CoWoS Stacked HBM2 | CAPACITY<br>32/16 GB HBM2 | |

# GPU comparison



**Performance**

Normalized performance (higher is better)

Legend: Average, RNN/Transformer, CNN

NVIDIA DGX-1 WITH TESLA V100
ESSENTIAL INSTRUMENT OF AI RESEARCH

960 Tensor TFLOPS | 8x Tesla V100 | NVLink Hybrid Cube
From 8 days on TITAN X to 8 hours
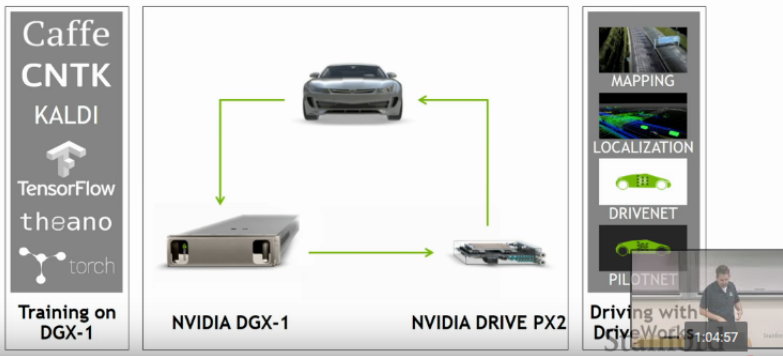400 servers in a box

NVIDIA DRIVE
END TO END SELF-DRIVING CAR PLATFORM

Parallelism

DEVIEW
2015

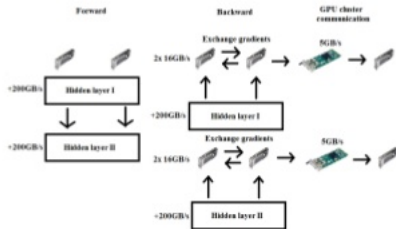## Data Parallelization

Forward

Backward

GPU cluster communication

Exchange gradients

2x 16GB/s

5GB/s

+200GB/s Hidden layer I

+200GB/s Hidden layer I

Exchange gradients

+200GB/s Hidden layer II

2x 16GB/s

5GB/s

+200GB/s Hidden layer II
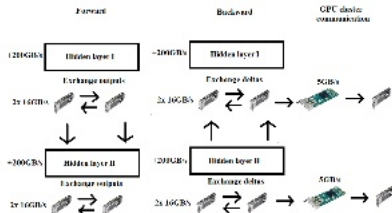
The good : Easy to implement
The bad  : Cost of sync increases with the number of GPU

## Model Parallelization

Forward

Backward

GPU cluster communication

+200GB/s Hidden layer I

+200GB/s Hidden layer I

Exchange outputs

Exchange deltas

2x 16GB/s

2x 16GB/s

5GB/s

+200GB/s Hidden layer II

+200GB/s Hidden layer II

Exchange outputs

Exchange deltas

2x 16GB/s

2x 16GB/s

5GB/s

The good : Larger network can be trained
The bad  : Sync is necessary in all layers