

AI & ML | C++

Ashot Vardanian

Intro

Who am I?

Ashot Vardanian, 23

Independent Software Developer
AI Researcher

Astrophysics Department,
St. Petersburg State Polytechnic
University

Basic, Pascal, JS, C, Objective-C,
Fortran, C++

[linkedin.com/in/ashvardanian](https://www.linkedin.com/in/ashvardanian)
[fb.com/ashvardanian](https://www.facebook.com/ashvardanian)



Plan

1. Theory behind Neural Networks
2. Demo: building `neural::mlp` from scratch
3. State-of-the-art
 1. ANN models
 2. AI Hardware
 3. ML Libraries
4. Demo: `std::pow("AI", 3)`
5. Future of C++

Main theme of the talk:

**Neural Networks are easy to
research and describe, but**

**Hard to implement and
optimise for hardware!**

Related Material

The image consists of three vertically stacked video frames from different CPPCON events. The left frame shows a terminal window with orange text output:

```
>>> c++ deep-learning.cpp  
>>> ./a.out data.csv  
>>> 42
```

The middle frame is a video still from CPPCON 2017. It features a man standing behind a dark podium, speaking into a microphone. The podium has a nameplate that reads "PETER GOLDSBOROUGH". To the right of the speaker, there is a slide with the title "A Tour of Deep Learning With C++". The top of the slide has the CPPCON 2017 logo and the text "THE C++ CONFERENCE • BELLEVUE, WASHINGTON". The right frame is a video still from CPPCON 2018. It shows the same man, Peter Goldsborough, now wearing a dark t-shirt with a yellow logo, standing and speaking. The top of the slide has the CPPCON 2018 logo and the text "THE C++ CONFERENCE • BELLEVUE, WASHINGTON". The slide title is "Machine Learning in C++ with PyTorch".

Theory

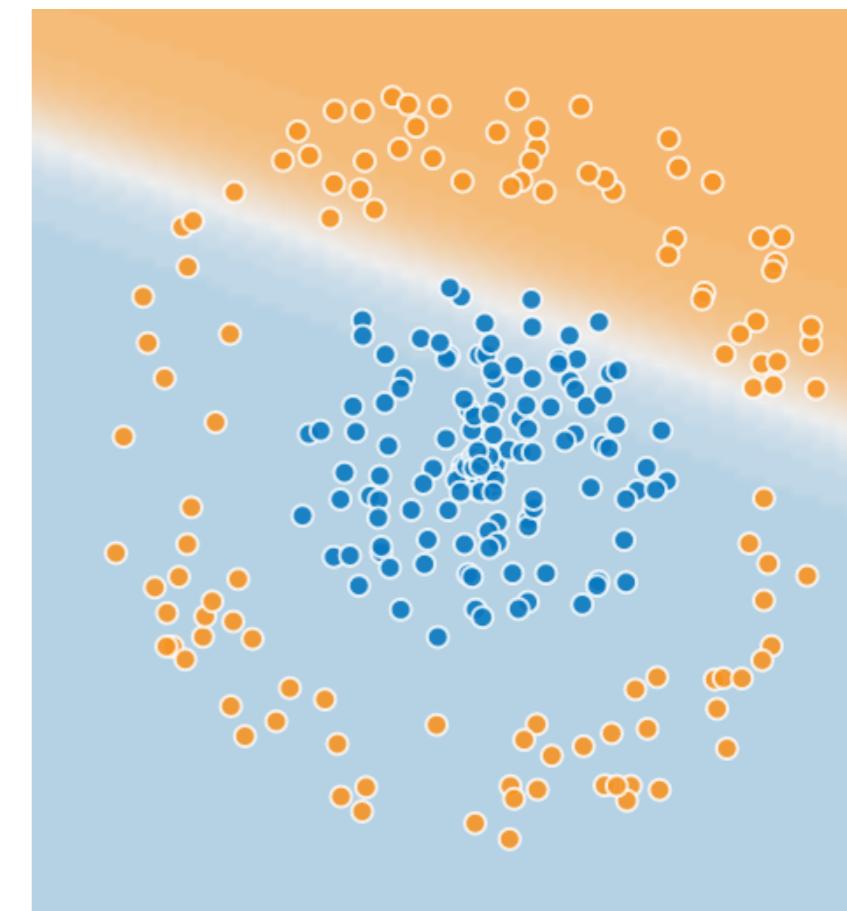
How can we split the points?



KD-Trees

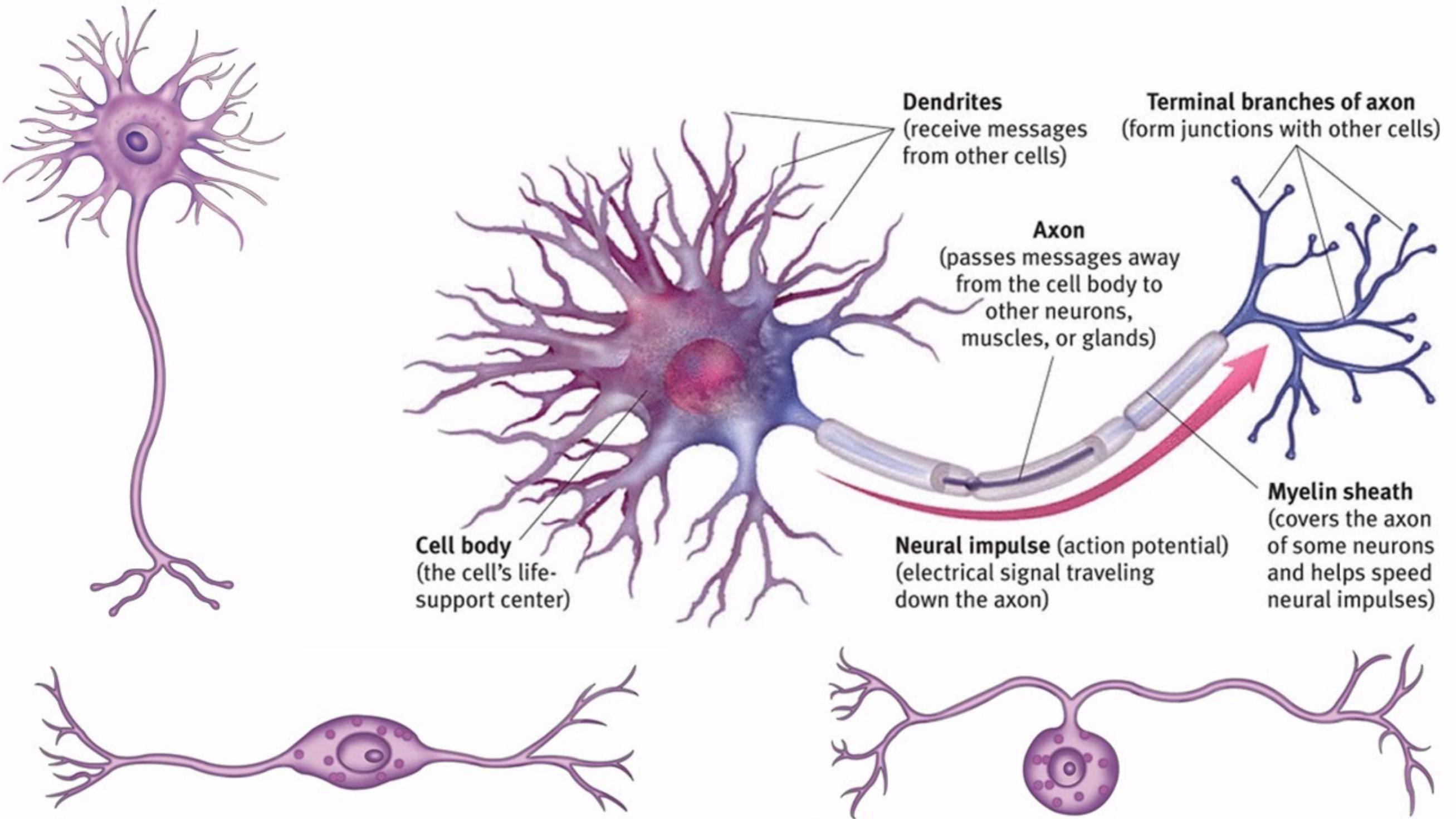


Linear Models

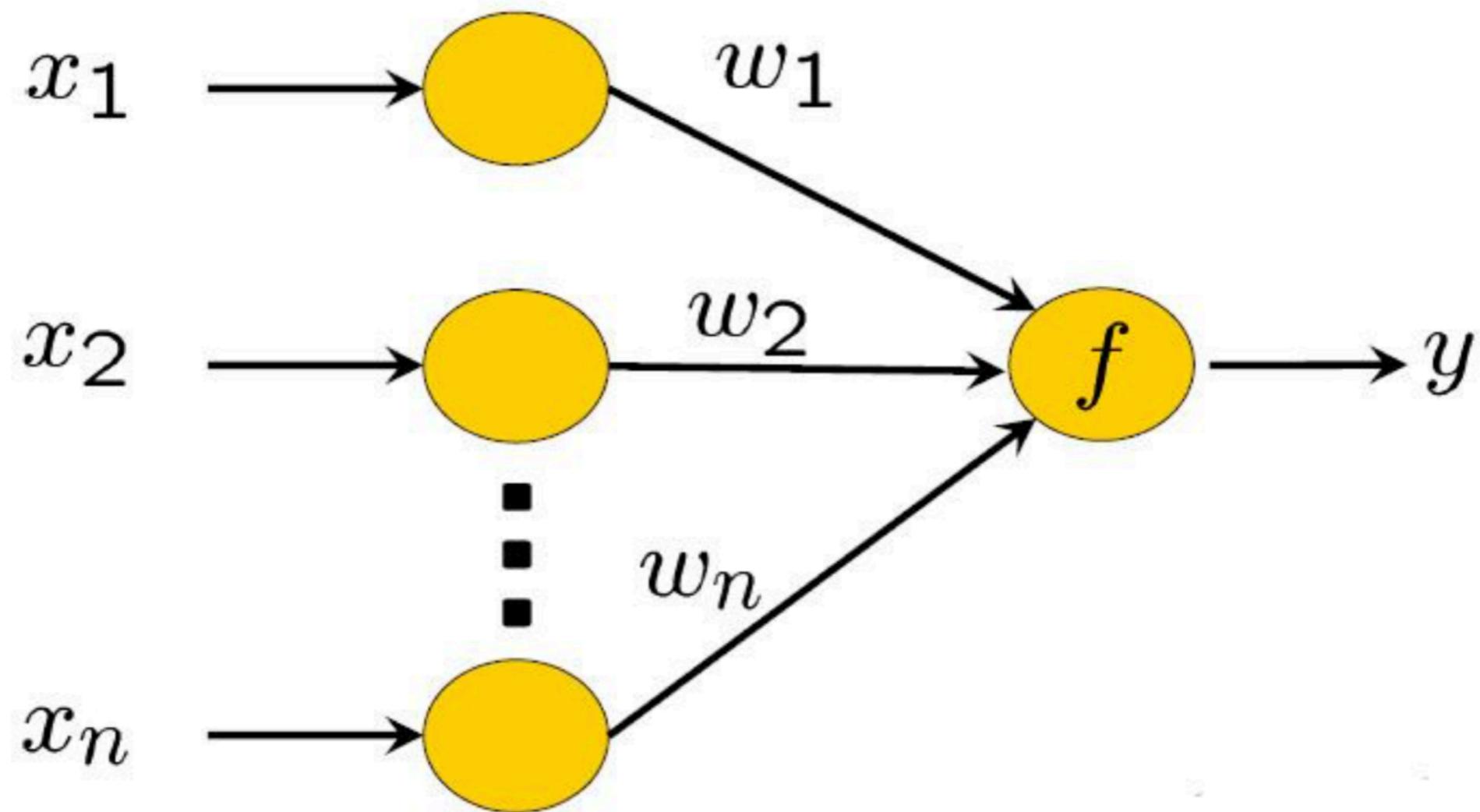


Non-Linear

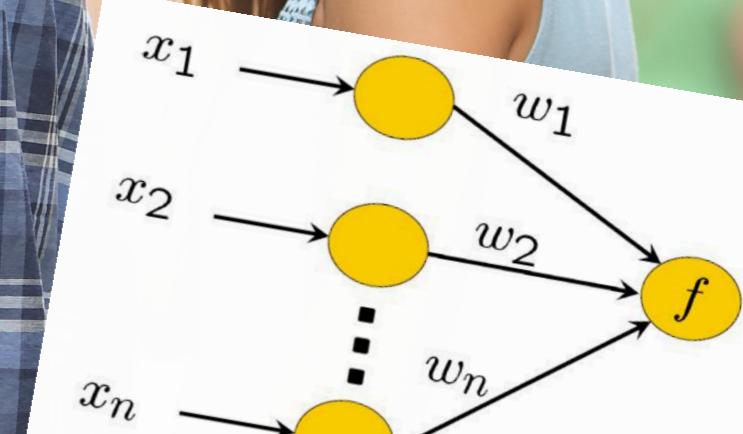
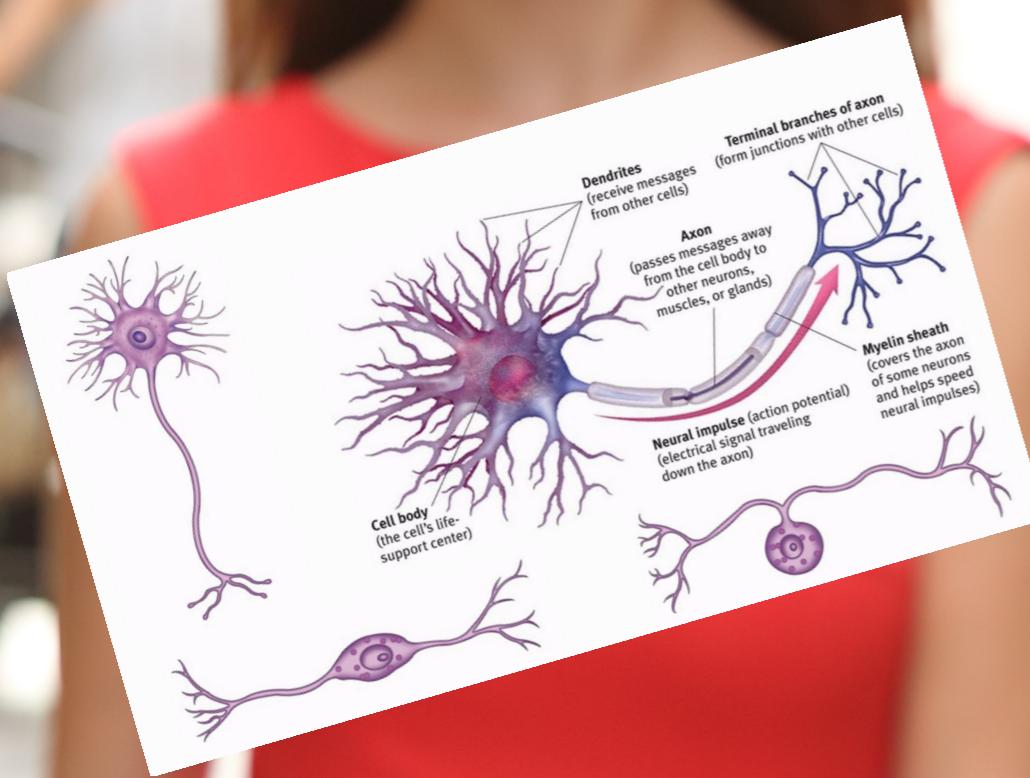
Inspiration



Incarnation



Expectation vs Reality

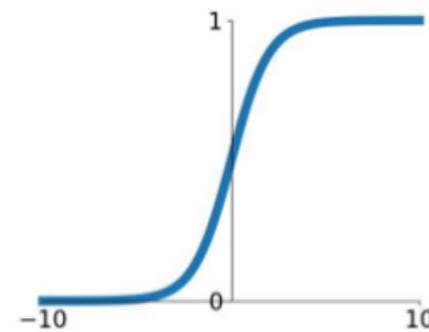


Activation Functions

The most common differentiable “non-linearities” in Neural Networks

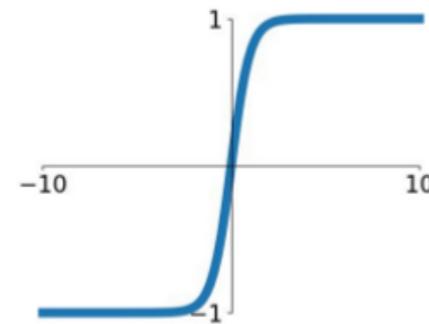
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



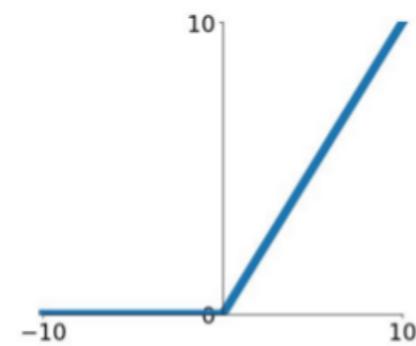
tanh

$$\tanh(x)$$



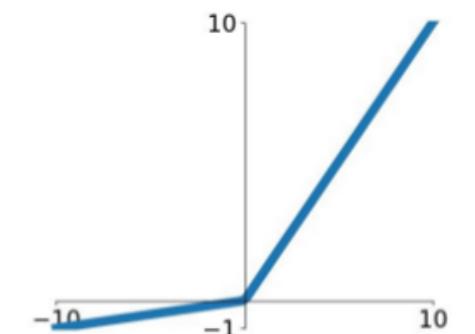
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

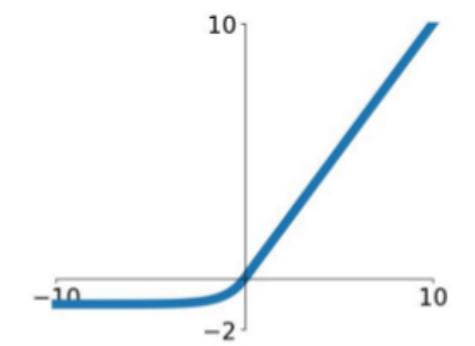


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Archean period

Single cell life

Progression

1943 - 1963

- First Neuron Model (1943), Pitts & McCulloch
- Perceptron (1958), Rosenblatt
- Perceptron convergence theorem (1962)
- **Back-propagation (1963), Bryson**

Degression

1963 - 1982

- Perceptron can't even learn the XOR function
- **Nobody does back-propagation, so no MLP training**

Protozoic period

Eucaryotes

Progression

1983 - 1993

- Hopfield Nets (1982)
- Boltzmann Machines (1985)
- Good understanding of “non-linearities”, but still poor understanding of network topology, layers complexity

Degression

1993-2006

- SVM (1995) are killing ANNs
- **Training deeper networks consistently yields poor results**

- Deep convolutional neural networks (1998), Yann LeCun

Phanerozoic period

Complex Multicellular Life

- **Hardware** improvements empower larger models
- Deep Belief Networks (2006), Hinton
- Deep Autoencoder based networks (2007), Bengio

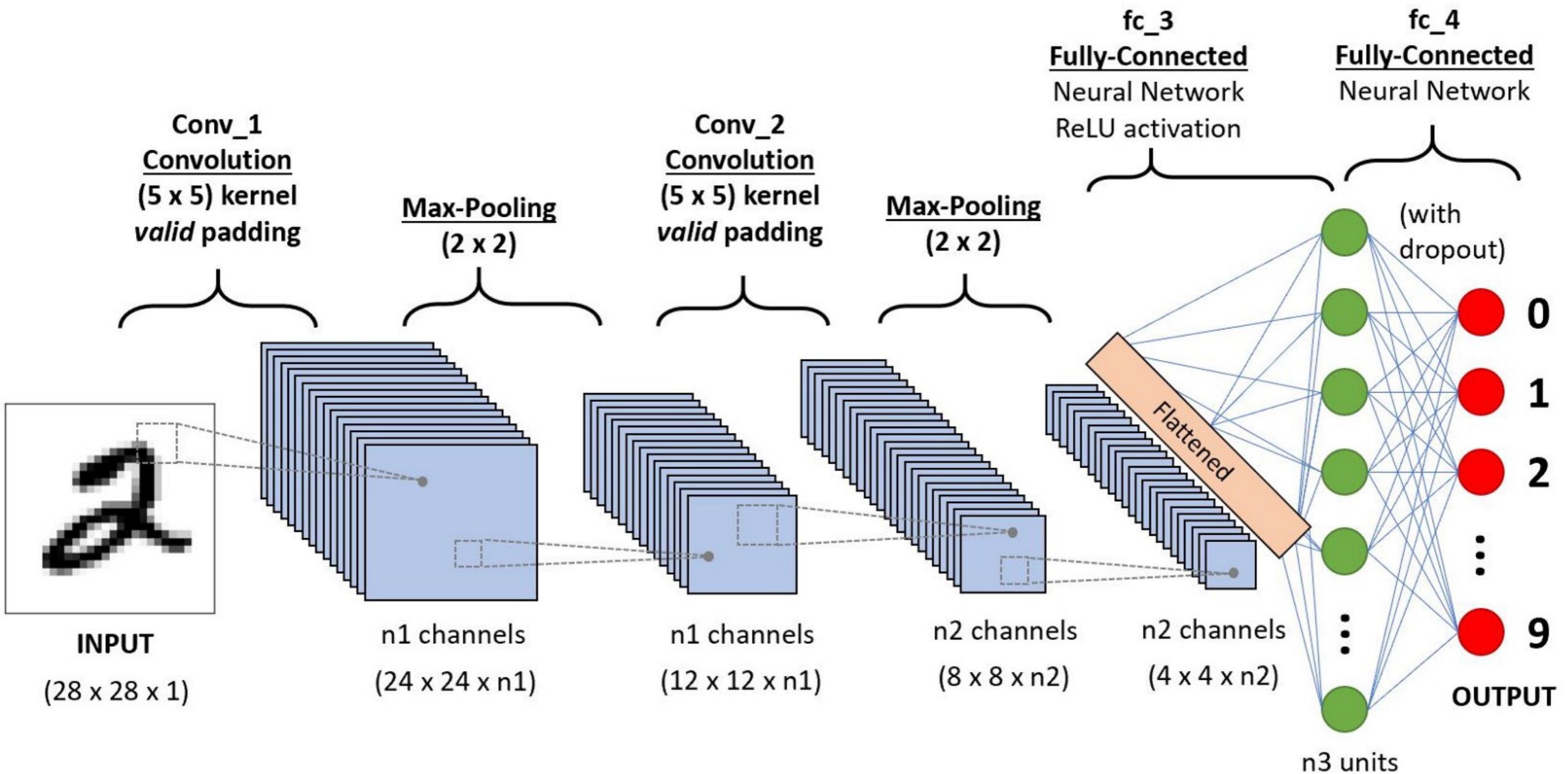


Simple Demo: Building MLPs with **STL**

State of the Art: Neural Networks

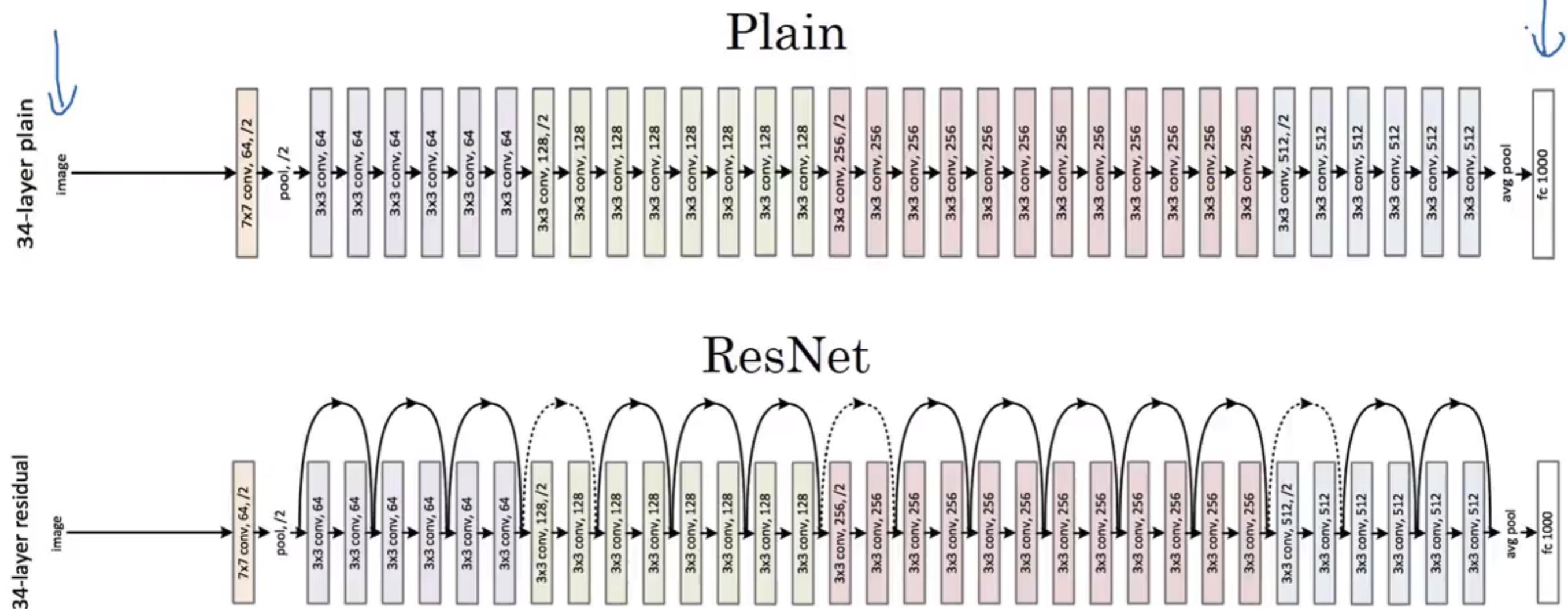
CNNs

A sequential model with convolutional layers



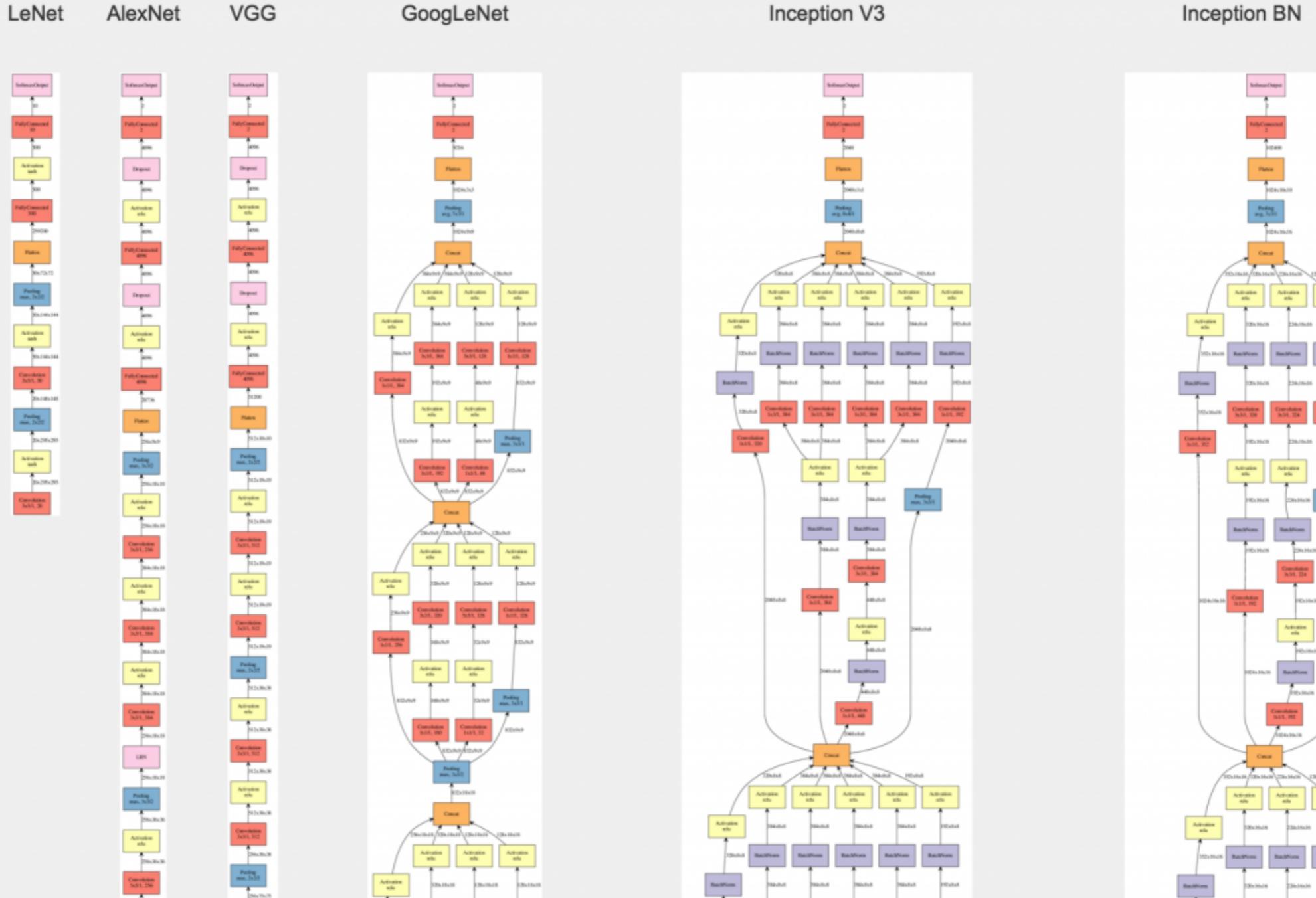
ResNet

A less sequential model with convolutional layers and skip-connections

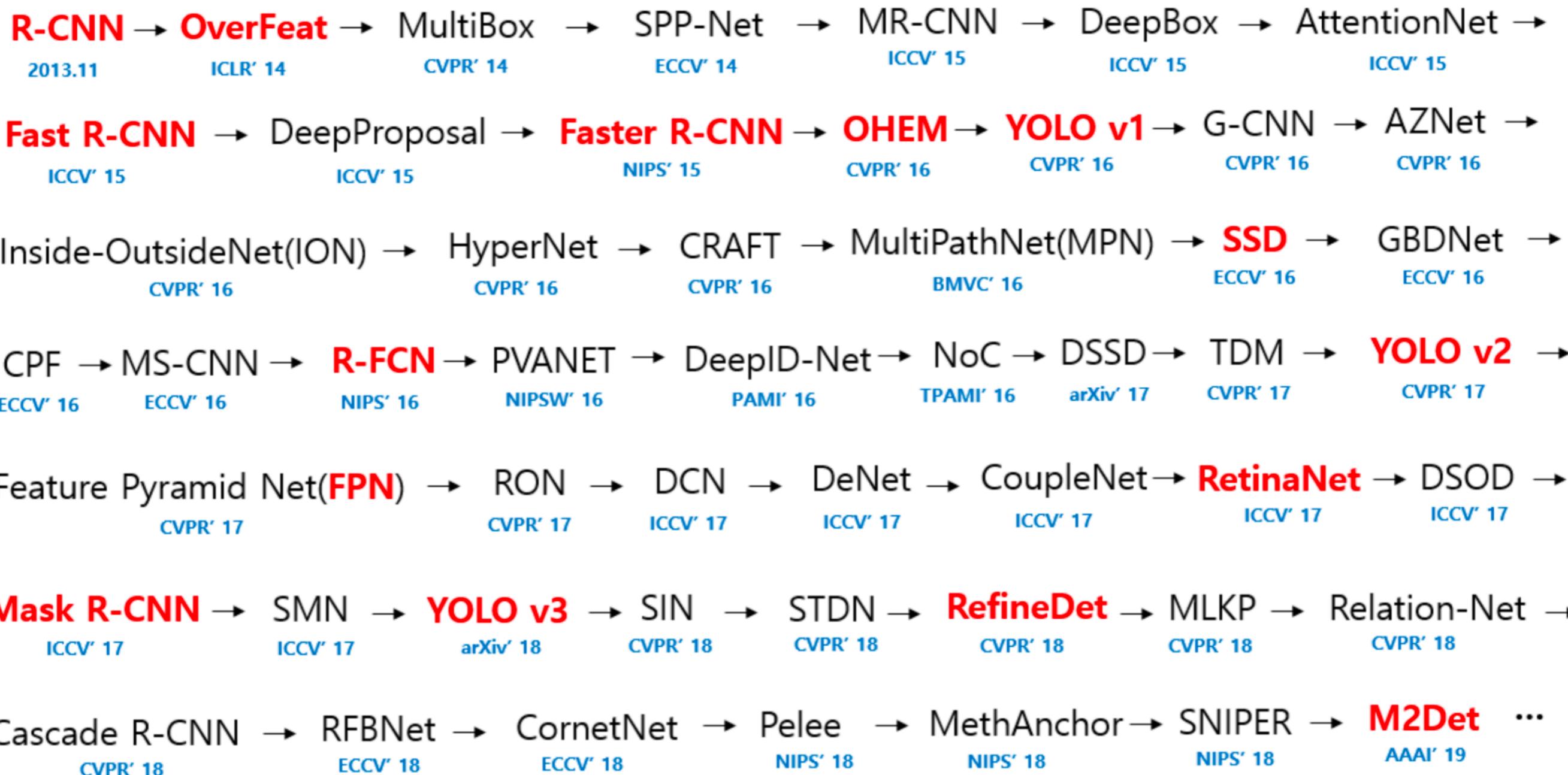


Complex Convolutional Networks

Directed acyclic computational graphs models used in production

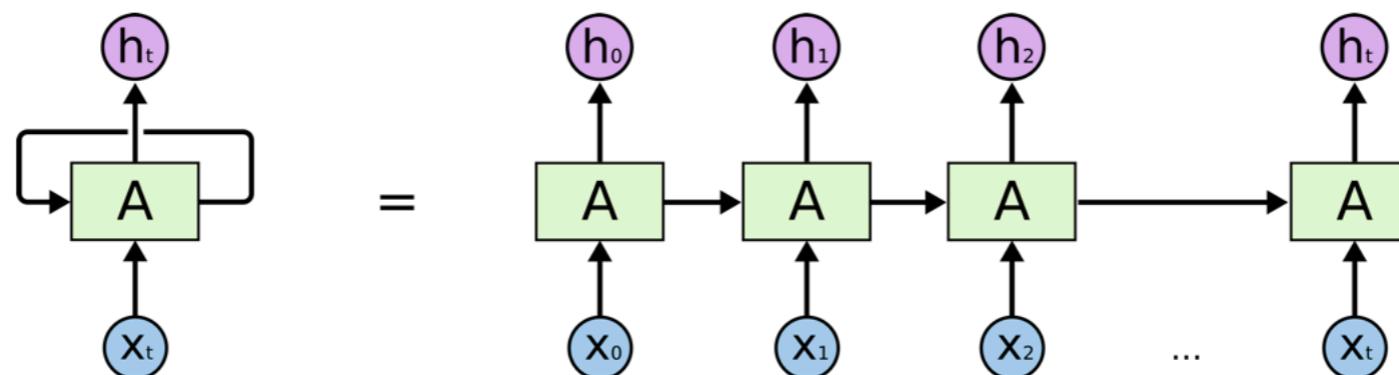


CNNs for Object Detection

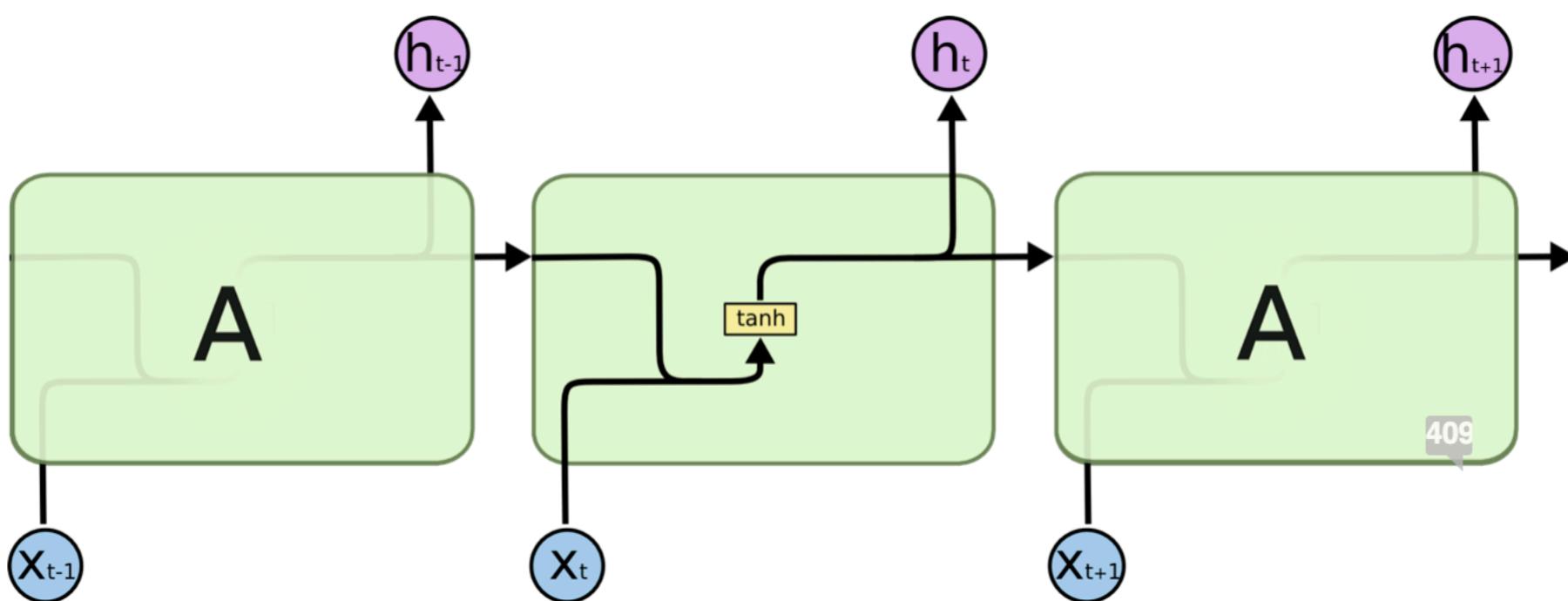


Recurrent Neural Networks

Back-propagation “thought-time” as a memory bottleneck and limitation for parallel execution



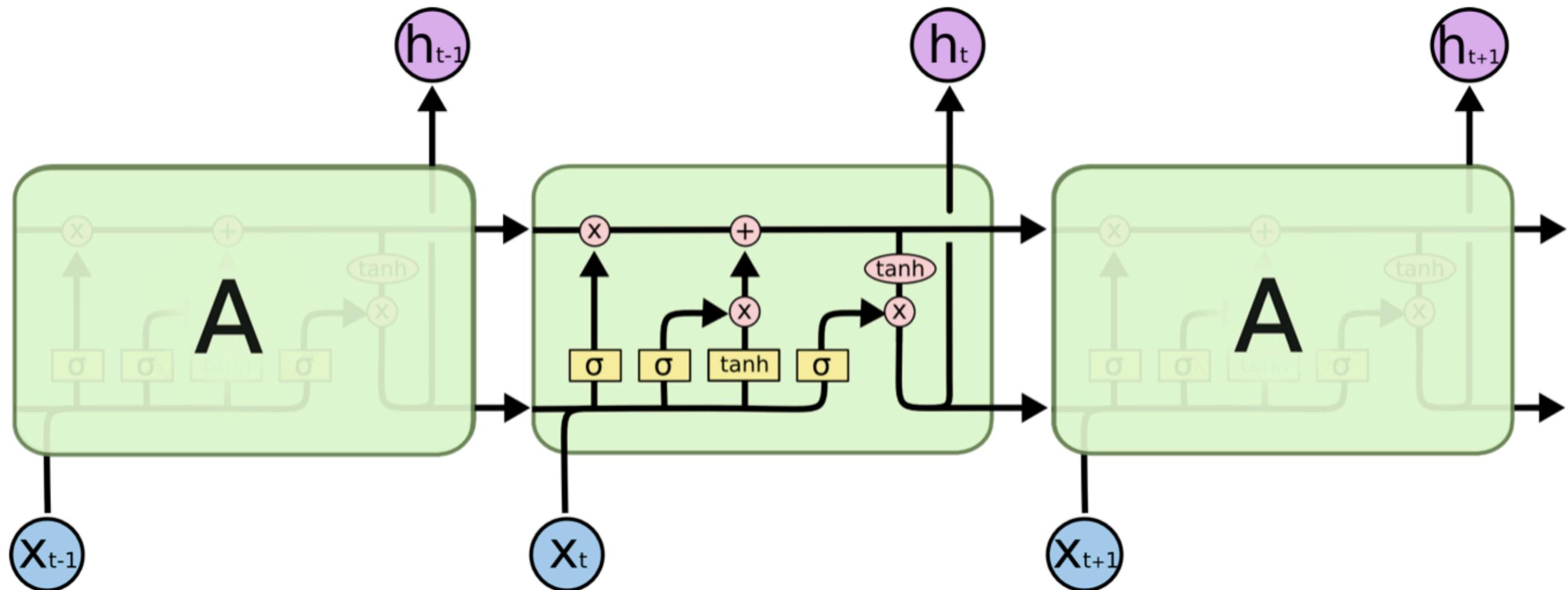
An unrolled recurrent neural network.



The repeating module in a standard RNN contains a single layer.

Long Short-Term Memory

Further increasing problem size
with “forget gate” & “modulation gate”



The repeating module in an LSTM contains four interacting layers.

Alpha Zero & Reinforcement

Anatomy of a world champion chess engine

Domain knowledge, extensions, heuristics in 2016 TCEC world champion Stockfish:

Board Representation:

Bitboards with Little-Endian Rank-File Mapping (LERF), Magic Bitboards, BMI2 - PEXT Bitboards, Piece-Lists,

Search:

Iterative Deepening, Aspiration Windows, Parallel Search using Threads, YBWC, Lazy SMP, Principal Variation Search.

Transposition Table:

Shared Hash Table, Depth-preferred Replacement Strategy, No PV-Node probing, Prefetch

Move Ordering:

Countermove Heuristic, Counter Moves History, History Heuristic, Internal Iterative Deepening, Killer Heuristic, MVV/LVA, SEE,

Selectivity:

Check Extensions if SEE ≥ 0 , Restricted Singular Extensions, Futility Pruning, Move Count Based Pruning, Null Move Pruning, Dynamic Depth Reduction based on depth and value, Static Null Move Pruning, Verification search at high depths, ProbCut, SEE Pruning, Late Move Reductions, Razoring, Quiescence Search

Evaluation:

Tapered Eval, Score Grain, Point Values Midgame: 198, 817, 836, 1270, 2521, Endgame: 258, 846, 857, 1278, 2558, Bishop Pair, Imbalance Tables, Material Hash Table, Piece-Square Tables, Trapped Pieces, Rooks on (Semi) Open Files, Outposts, Pawn Hash Table, Backward Pawn, Doubled Pawn, Isolated Pawn, Phalanx, Passed Pawn, Attacking King Zone, Pawn Shelter, Pawn Storm, Square Control, Evaluation Patterns,

Endgame Tablebases:

Syzygy TableBases

Same system can play any 2-player game

Chess



4 Hours

AlphaZero
surpasses StockFish

Shogi



2 Hours
AlphaZero
surpasses Elmo

8 Hours
AlphaZero
surpasses AlphaGo

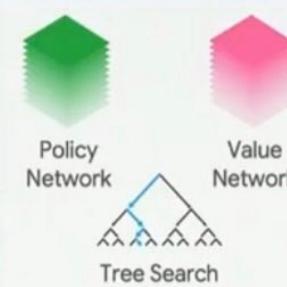
AlphaZero (AlphaGoZero)

~100K self-play games



Synthetic Training Data
Use new training data to train new network

~40M games
~3s per game
~5000 games at a time



New Policy Network
New Value Network

55% win rate

Amount of search per decision

Human Grandmaster
100's of moves

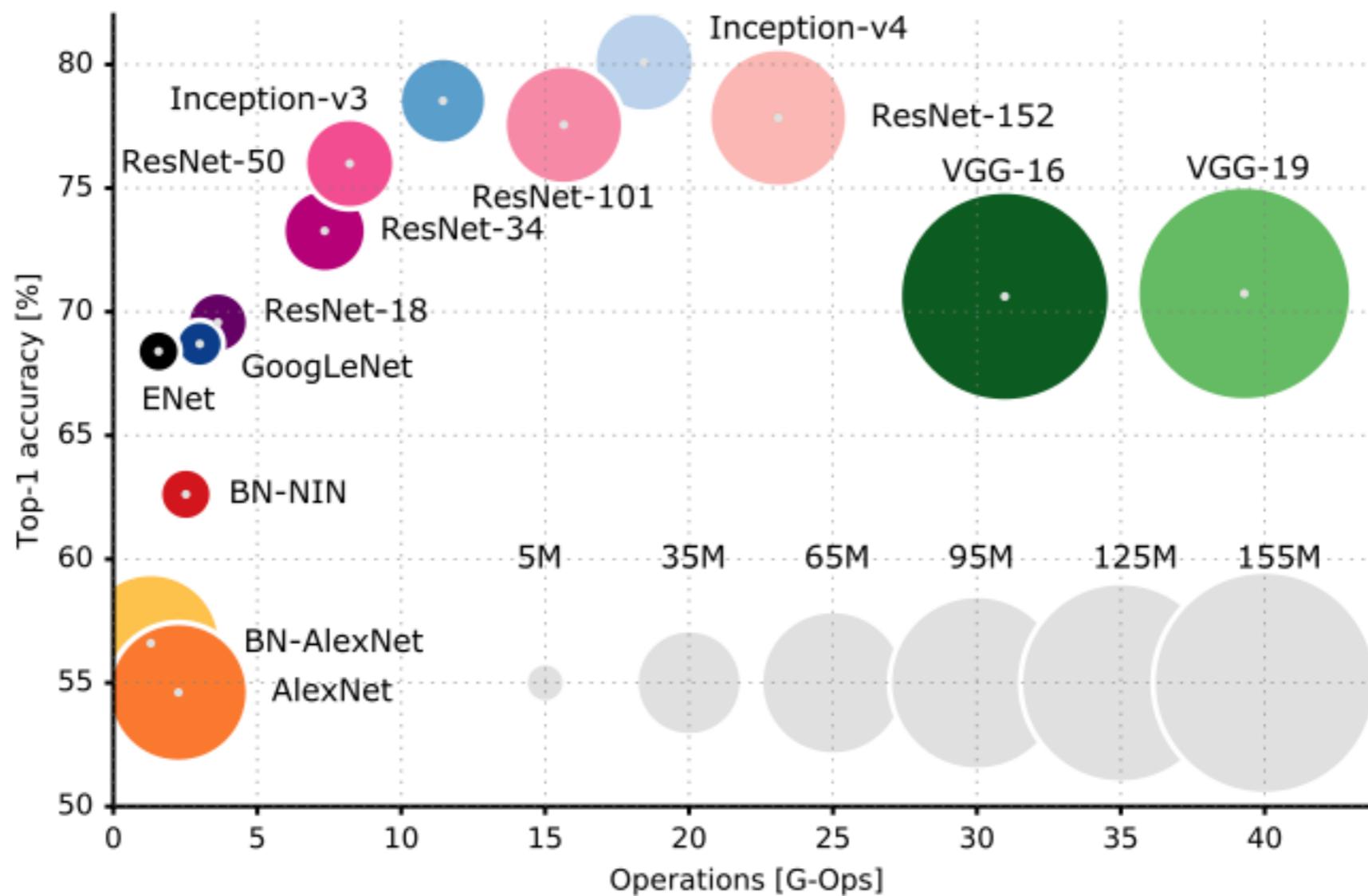
AlphaZero
10,000's of moves

State-of-the-art chess engine
10,000,000s of moves

Let's look into some
“ONNX” graphs

SotA Models Sizes

The deep models have enormous number of weights, which implies
long training time and potentially higher accuracy!



Model Training Price

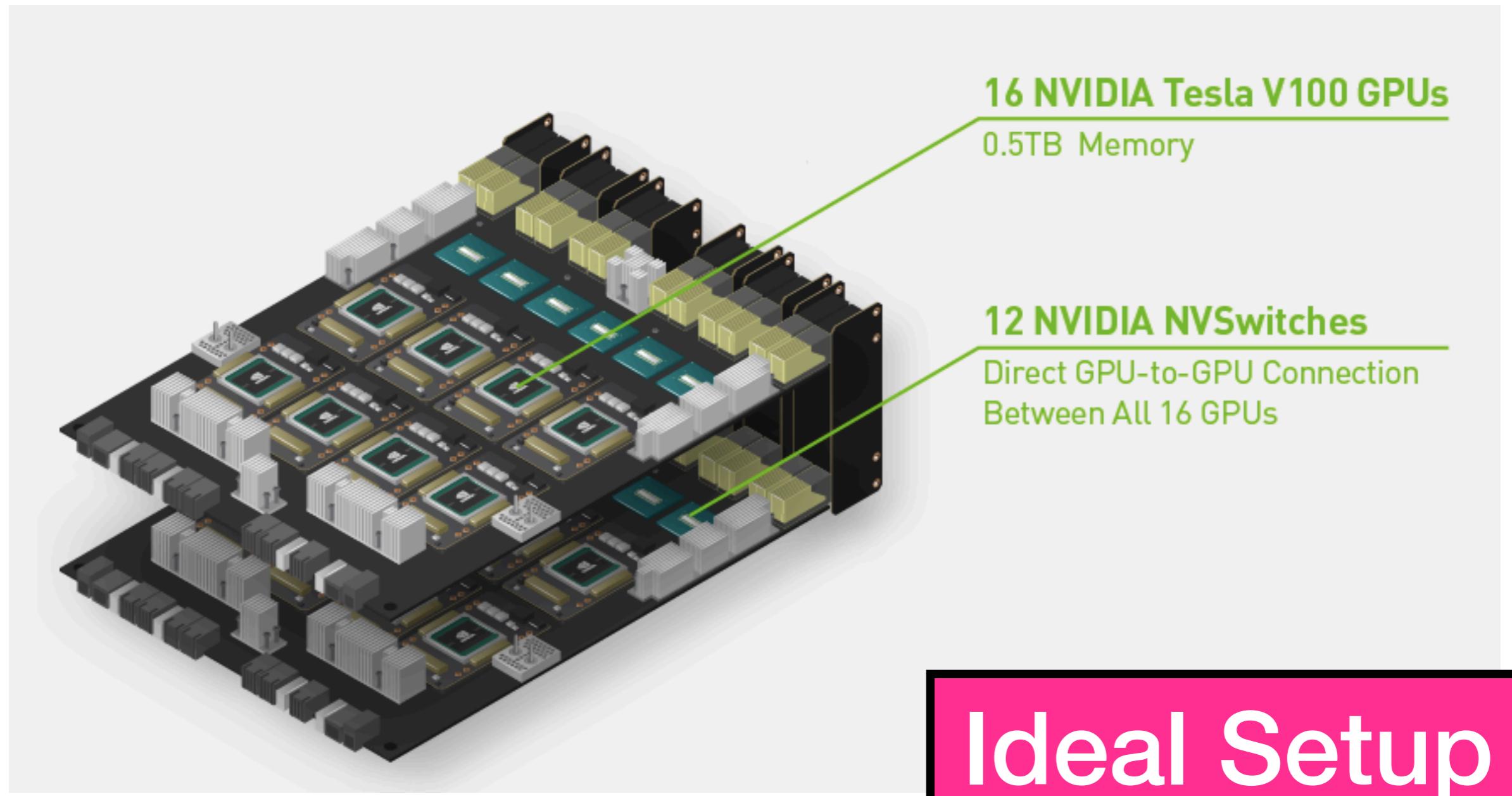
Inference of ANNs is relatively cheap, but training a single SotA model may cost over 100K USD in the cloud...

So the hardware must be optimized!

Product (CPU/GPU)	GPUs	On Demand \$/Hour	3 Year Reserve \$/Hour	3 Year Cost 100% Utilization
p2.8xlarge (64 vCores/K80)	8	\$7.20	\$3.40	\$89,352.00 to \$189,216.00
p3.16xlarge (64 vCores/V100)	8	\$24.48	\$9.87	\$259,383.60 to \$643,334.40
On-Premises (2x22 Xeon/P100)	8	~ \$2.59	-	~ \$68,000.00
On-Premises (Xeon/V100)	8	~ \$3.81	-	~ \$100,000.00

State of the Art: AI Hardware

Why are those instances so expensive?



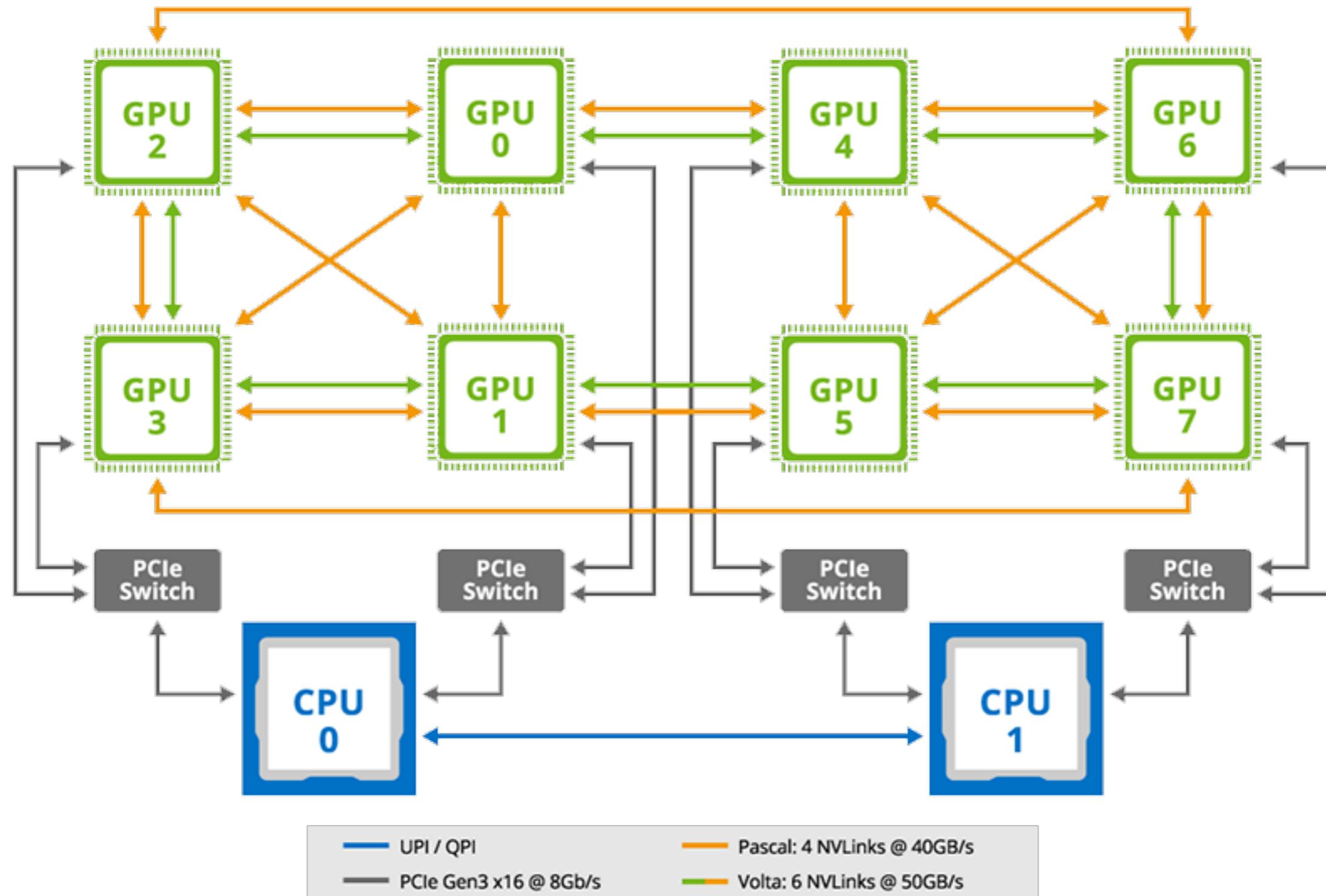
Good Setup





Bad Setup

Proximity of memory buffers defines the performance limitations.
Every framework must track the **NUMA** layout to fully use **RDMA**.



Performance Metrics

	Intel Xeon Platinum 8180M	Nvidia Titan V	Google TPU v3
Year	2017	2017	2018
Memory	< 1.5 TB DDR4	12 GB HBM2	4x16 GB HBM
Memory Bandwidth	< 150 GB/s	650 GB/s	2.4 TB/s
TDP	205 W	250 W	~ 250 W
Fabrication Process	14 nm	12 nm	12-20 nm
Floating point Performance	< 10 TOPs	15 (120) TOPs	180 TOPs

State of the Art: AI Libraries

Library Zoo



C++

Python

Manual Memory Control

Reference Counting for every variable

Compiled

Interpreted

Scope-limited Lifetime

Variables are accessible outside the loop

Statically-typed

Dynamically-typed

Floating point types are **float** and **double**

Floats are implementation specific, but normally have the capacity of **double**

Integer types are hardware defined

At least 16 bytes for zero integer, 24 bytes in most cases, 36 bytes or more for flexible

No proper modules

> 80 K easily pluggable modules

Complicated Optimised Code

Simple Slow Code

	C++	Python		
Tasks	Time	Memory	Time	Memory
Mandelbrot	2 s	26 Kb	260 s	52 Kb
N-Body	8 s	1 Kb	797 s	8 Kb
Spectral Norm	2 s	1 Kb	173	52 Kb
Fasta	1 s	2 Kb	63	683 Kb
Fannkuch Redux	10 s	1 Kb	472	50 Kb
Binary Trees	4 s	117 Kb	80	450 Kb
K-Nucleotide	4 s	156 Kb	73	188 Kb
Regex Redux	2 s	203 Kb	18	446 Kb
Rev Complement	2 s	996 Kb	18	1007 Kb
Pi Digits	2 s	4 Kb	3	10 Kb

Popular AI projects

	Tensorflow	PyTorch	MxNet	Keras	CNTK	Horovod
First Release	Nov, 2015	Jan, 2012	Apr, 2015	Mar, 2015	Aug, 2014	Aug, 2017
GitHub Stars	130 K	28 K	17 K	42 K	16 K	6 K
GitHub Watchers	8 K	1 K	1 K	2 K	1 K	280
Commits	57 K	18 K	9 K	5 K	16 K	300
Forks	75 K	7 K	6 K	16 K	4 K	1 K
Used by	41 K		400	27 K		58

Audience

	Tensorflow	PyTorch	MxNet	Keras	CNTK	Horovod
GitHub Repos	60 K	23 K	2 K	24 K	500	80
Arxiv Mentions	339	119	34	51	15	6
Google Trends	-3%	79%	2%	0%	-78%	42%
Interest Regions	China, Korea	China, Korea	China	Indonesia, China, Malaysia, Korea.	China, Indonesia, Korea	China, Korea

Codebases

	Tensorflow	PyTorch	MxNet	Keras	CNTK	Horovod
Last Commit ID	d38eaa19b3b656a 4d4fd4ae388004a8 1a1c30094	c6255a57e4c15b8 8ded24195a9983a 4aac401c22	f4598e743aa70e18 9db66e3d696f4cb 52441d098	9d33a024e3893ec 2a4a15601261f447 25c6715d1	8ed70502bf1019e0 ddeab8420d8babd d52bc7941	73d860f239632176 1e0f5ef6fe934130a fd69094
Total Files	10,778	4,715	2,856	252	1,727	161
Code Lines	2,251,532	710,449	406,488	46,618	331,206	17,557
Com. Lines	555,516	100,223	119,447	15,730	58,753	4,892
C++	53%	56%	35%	0%	59%	48%
Python	31%	26%	24%	93%	16%	44%
CUDA/ OpenCL	0%	13%	4%	0%	5%	0%

**Let's take a look at
some samples!**

Using Keras for MLP

```
model = Sequential()
model.add(Dense(512, activation='relu', input_shape=(784,)))
model.add(Dropout(0.2))
model.add(Dense(512, activation='relu'))
model.add(Dropout(0.2))
model.add(Dense(num_classes, activation='softmax'))

model.summary()
model.compile(loss='categorical_crossentropy',
              optimizer=RMSprop(),
              metrics=['accuracy'])
history = model.fit(x_train, y_train,
                      batch_size=batch_size,
                      epochs=epochs,
                      verbose=1,
                      validation_data=(x_test, y_test))
score = model.evaluate(x_test, y_test, verbose=0)
print('Test loss:', score[0])
print('Test accuracy:', score[1])
```

Using Keras for CNN

```
model = Sequential()
model.add(Conv2D(32, activation='relu',
                 kernel_size=(3,3),
                 input_shape=input_shape))
model.add(Conv2D(64, (3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(num_classes, activation='softmax')))

model.summary()
model.compile(loss=losses.categorical_crossentropy,
              optimizer=optimizers.Adadelta(),
              metrics=['accuracy'])
```

Using TF Python API for CNN

Using TF C++ for CNN

```
DataSet data_set("/path/normalized_car_features.csv");

Tensor x_data(DataTypeToEnum<float>::v(),
             TensorShape { data_set.x().size()/3, 3 });
std::copy_n(data_set.x().begin(), data_set.x().size(),
            x_data.flat<float>().data());

Tensor y_data(DataTypeToEnum<float>::v(),
             TensorShape { data_set.y().size(), 1 });
std::copy_n(data_set.y().begin(), data_set.y().size(),
            y_data.flat<float>().data());
```

Using TF C++ for MLP (1/3)

```
// This contains the graph.  
Scope scope = Scope::NewRootScope();  
  
// Every operation injects itself into the graph.  
auto x = Placeholder(scope, DT_FLOAT);  
auto y = Placeholder(scope, DT_FLOAT);  
  
// Make weights matrices.  
auto w1 = Variable(scope, {3, 3}, DT_FLOAT);  
auto assign_w1 = Assign(scope, w1, RandomNormal(scope, {3, 3}, DT_FLOAT));  
  
auto w2 = Variable(scope, {3, 2}, DT_FLOAT);  
auto assign_w2 = Assign(scope, w2, RandomNormal(scope, {3, 2}, DT_FLOAT));  
  
auto w3 = Variable(scope, {2, 1}, DT_FLOAT);  
auto assign_w3 = Assign(scope, w3, RandomNormal(scope, {2, 1}, DT_FLOAT));  
  
// Make bias vectors.  
auto b1 = Variable(scope, {1, 3}, DT_FLOAT);  
auto assign_b1 = Assign(scope, b1, RandomNormal(scope, {1, 3}, DT_FLOAT));  
  
auto b2 = Variable(scope, {1, 2}, DT_FLOAT);  
auto assign_b2 = Assign(scope, b2, RandomNormal(scope, {1, 2}, DT_FLOAT));  
  
auto b3 = Variable(scope, {1, 1}, DT_FLOAT);  
auto assign_b3 = Assign(scope, b3, RandomNormal(scope, {1, 1}, DT_FLOAT));
```

Using TF C++ for MLP (2/3)

```
auto layer_1 = Tanh(scope, Add(scope, MatMul(scope, x, w1), b1));
auto layer_2 = Tanh(scope, Add(scope, MatMul(scope, layer_1, w2), b2));
auto layer_3 = Tanh(scope, Add(scope, MatMul(scope, layer_2, w3), b3));

auto regularization = AddN(scope, {
    L2Loss(scope, w1),
    L2Loss(scope, w2),
    L2Loss(scope, w3)
});

auto loss = Add(scope,
    ReduceMean(scope, Square(scope, Sub(scope, layer_3, y)), {0, 1}),
    Mul(scope, Cast(scope, 0.01, DT_FLOAT), regularization));
```

Using TF C++ for MLP (3/3)

```
// Backward pass:  
std::vector<Output> grad_outputs;  
AddSymbolicGradients(scope, { loss }, { w1, w2, w3, b1, b2, b3 }, &grad_outputs);  
  
// Update the weights and bias using gradient descent.  
auto apply_w1 = ApplyGradientDescent(scope, w1, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[0]});  
auto apply_w2 = ApplyGradientDescent(scope, w2, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[1]});  
auto apply_w3 = ApplyGradientDescent(scope, w3, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[2]});  
auto apply_b1 = ApplyGradientDescent(scope, b1, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[3]});  
auto apply_b2 = ApplyGradientDescent(scope, b2, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[4]});  
auto apply_b3 = ApplyGradientDescent(scope, b3, Cast(scope, 0.01, DT_FLOAT), {grad_outputs[5]});  
  
ClientSession session(scope);  
std::vector<Tensor> outputs;  
  
// Init the weights and biases by running the assigns nodes.  
session.Run({assign_w1, assign_w2, assign_w3, assign_b1, assign_b2, assign_b3}, nullptr);  
  
// Training steps.  
for (int i = 0; i < 5000; i++) {  
    session.Run({{x, x_data}, {y, y_data}}, {loss}, &outputs);  
    session.Run({{x, x_data}, {y, y_data}}, {apply_w1, apply_w2, apply_w3, apply_b1, apply_b2,  
    apply_b3, layer_3}, nullptr);  
}
```

Using PyTorch C++

```
/// Automatically differentiable, CPU/GPU-enabled arrays.  
torch::Tensor;  
  
/// Composable modules for ANNs.  
torch::nn;  
  
/// Optimization algorithms: SGD, Adam or RMSprop...  
torch::optim;  
  
/// Datasets, pipelines and multi-threaded, async data loader.  
torch::data;  
  
/// Serialization API for model checkpoints.  
torch::serialize;  
  
/// Glue to bind your C++ models into Python.  
torch::python;  
  
/// Pure C++ access to the TorchScript JIT compiler.  
torch::jit;
```

The front-ends of all major libraries look the same, but is that all?
The complexity arises in the back-end part!

Chaotic universe of deep learning



Vendor 1

Vendor 2

Vendor n

CPU

Many-core

DSP

FPGA

GPU

ASIC/ASIP

GM's world famous secret network

Caffe

SSD

libDNN-clBLAS

clBLAS

VGG

CLBlast

cuBLAS (BVLC)

YOLO

cuDNN (BVLC)

GoogleNet

Theano

cuDNN (NVIDIA)

AlexNet

TensorFlow

Fast RCNN

Torch

libDNN-viennaCL

cuBLAS fp16 (NVIDIA)

ResNet

cuDNN fp16 (NVIDIA)

viennaCL

RCNN

cuBLAS (NVIDIA)

libDNN-cuBLAS

CNTK

libDNN-CLBlast

TensorRT fp16

Mask RCNN

OpenBLAS



Deep Learning Stack

AI frameworks span across all those layers and often incorporate **multiple libraries at every layer** for maximum deployment flexibility!

Deep Learning

Keras, TF, cuDNN, MKL-DNN

Computational Graphs*

Eigen, TF, VexCL, ArrayFire

Math

Eigen, MKL, VexCL, cuBLAS, ArrayFire, Boost.Compute

Parallelism

Intel TBB, OpenML, OpenACC

Language & Extensions*

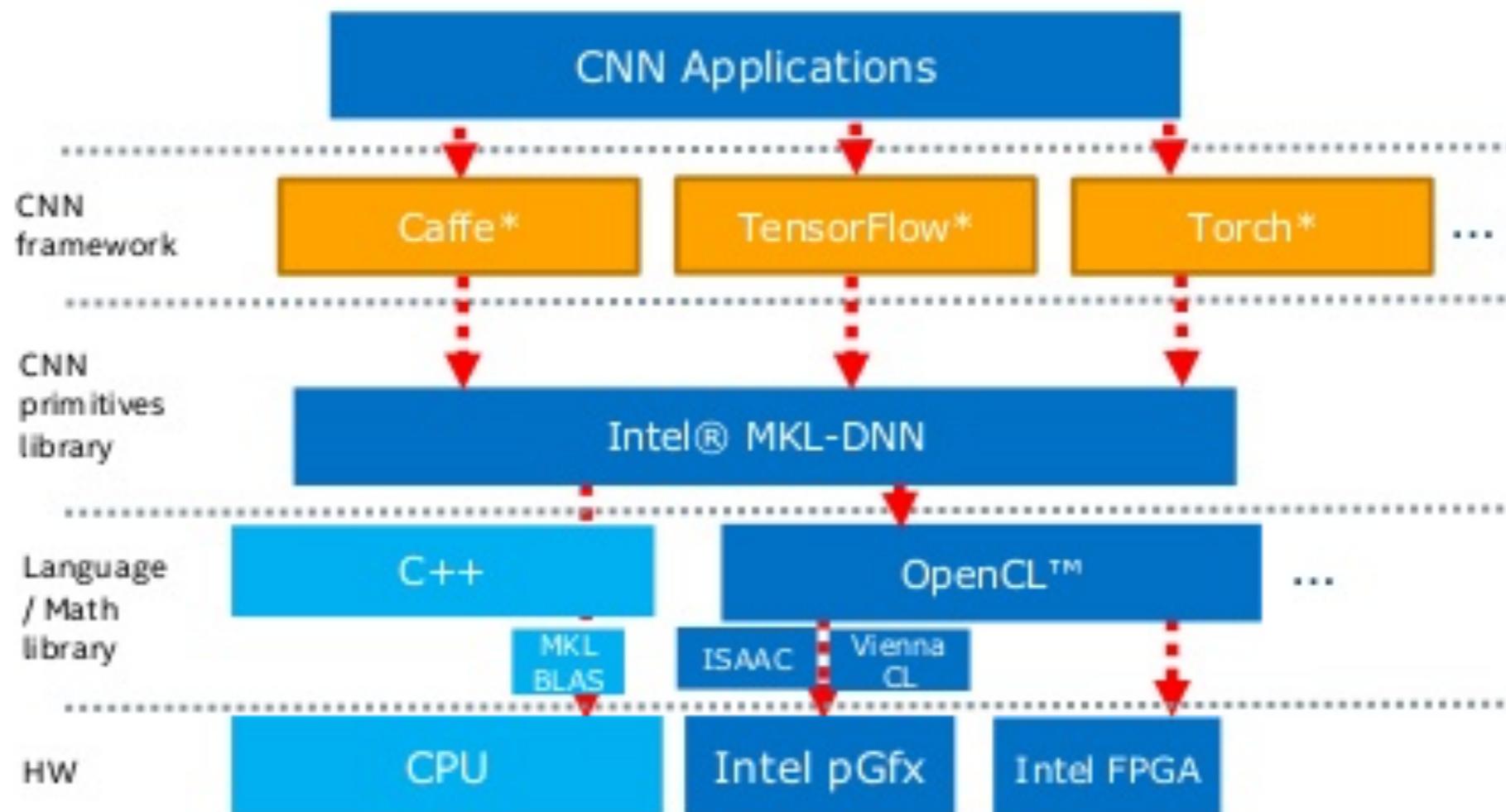
CUDA, OpenCL, OpenML, OpenACC

Compilers*

LLVM, TVM, GCC

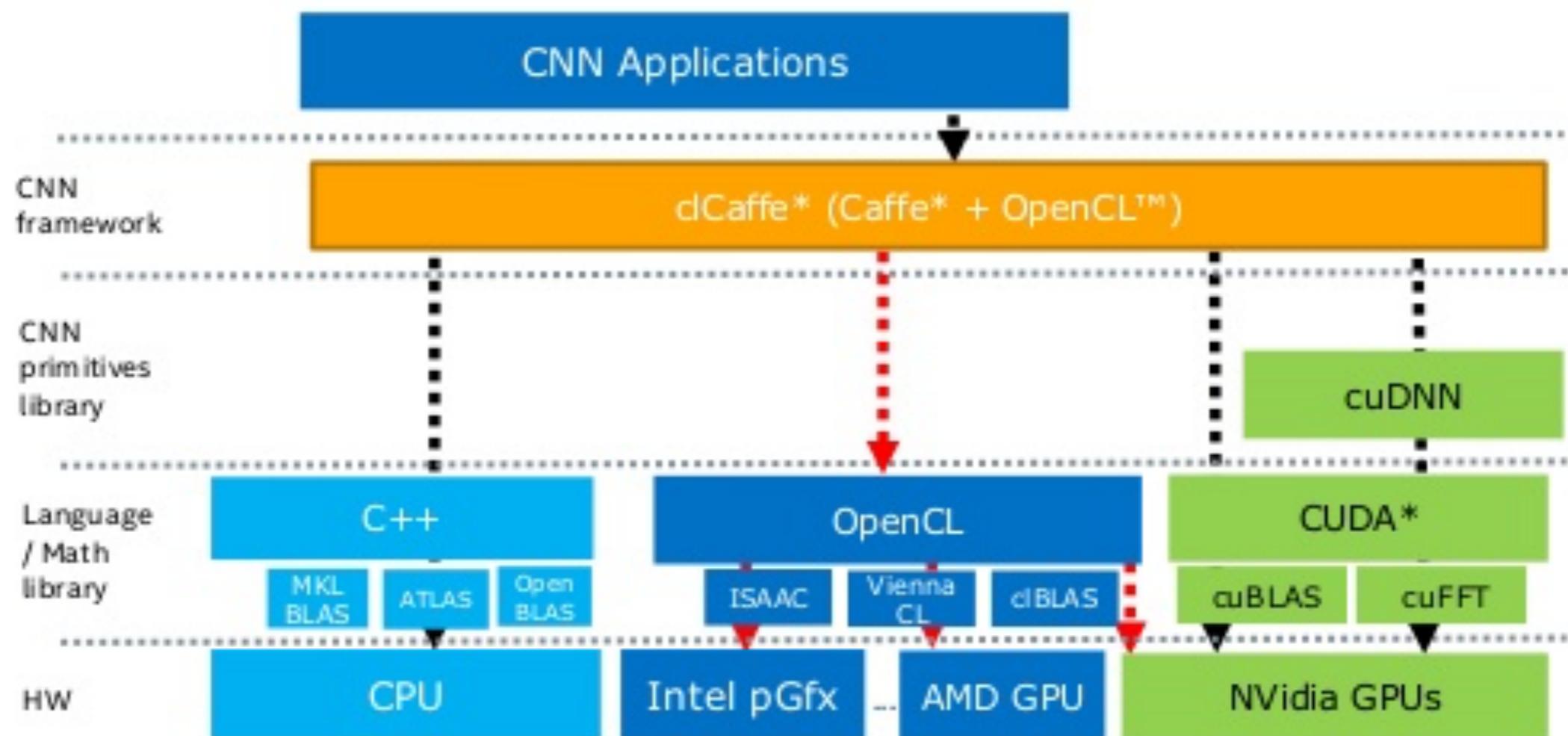
Intel Software Stack

Intel provides multiple libraries optimised for their CPUs.
They have different names, but their functionality is often similar.
MKL, MKL-DNN, BLAS, TBB.



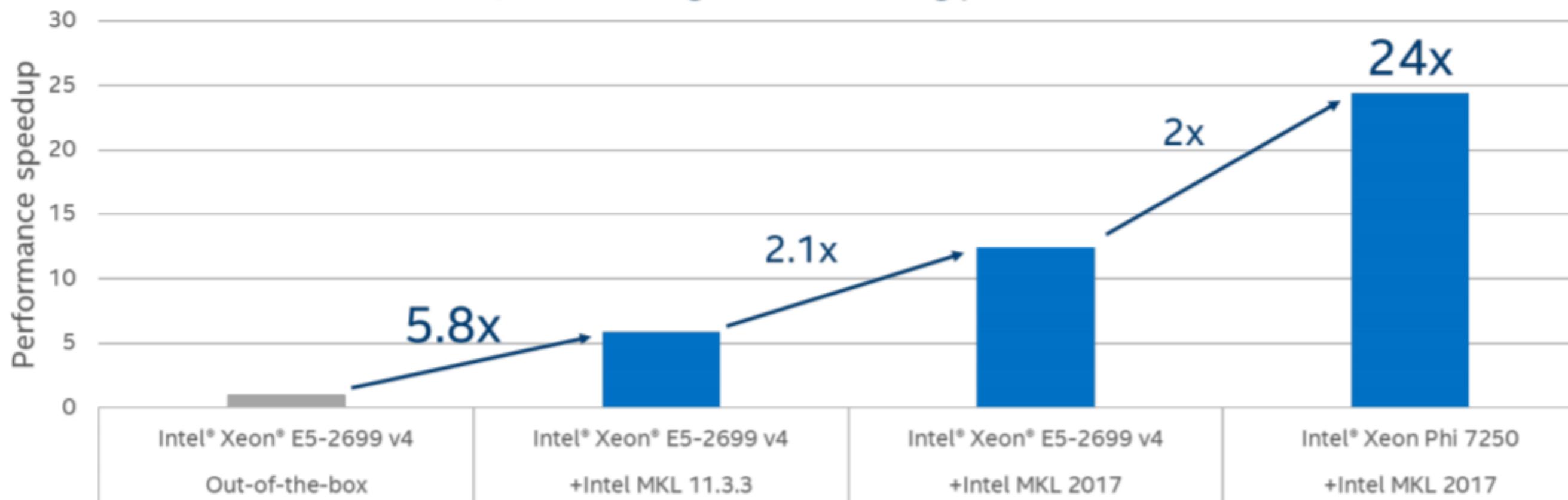
Intel GPU Stack

Intel supported OpenCL as primary language for their integrated GPUs.
ViennaCL, cIBLAS, Vulkan.



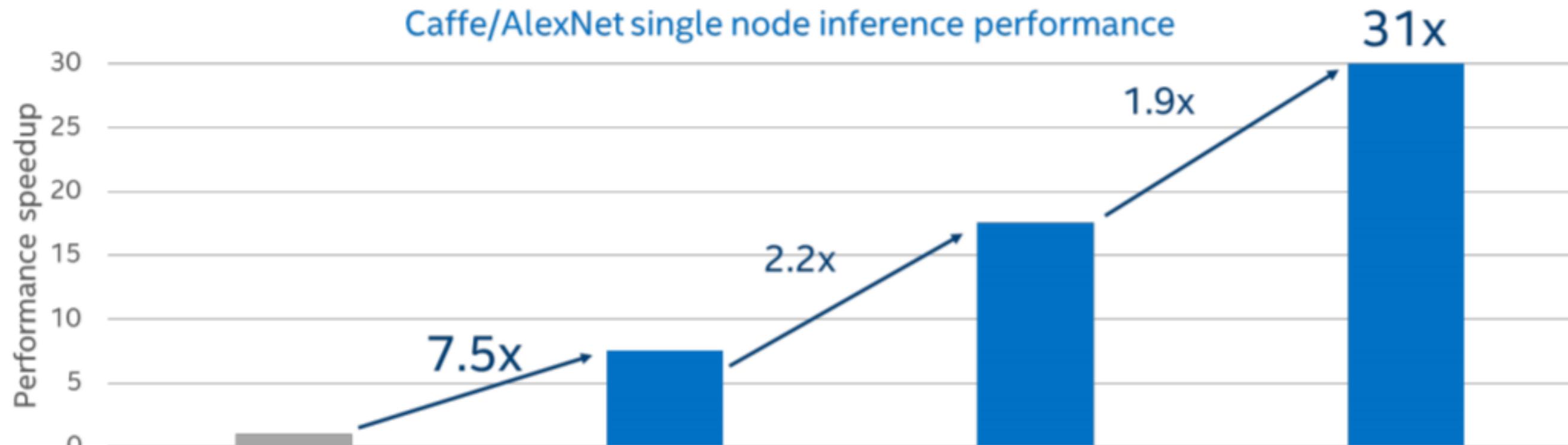
Improved Deep Neural Network training performance using Intel® Math Kernel Library (Intel® MKL)

Caffe/AlexNet single node training performance



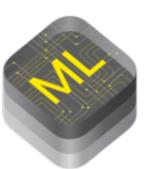
Improved Deep Neural Network inference performance using Intel® Math Kernel Library (Intel® MKL)

Caffe/AlexNet single node inference performance



Nvidia Software Stack

1. Math:
 1. **cuBLAS**: Linear Algebra
 2. cuSPARSE: Sparse MAtrix Operations
 3. **cuDNN**: Deep Learning Primitives
2. Device-Specific:
 1. **Automatic Mixed Precision**: Tensor Cores Support on Volta
 2. OpticalFlowSDK: Optical Flow for Video Inference on Turing
3. Networking:
 1. **NCCL**: Multi-GPU Communication
4. Pre/Post-processing:
 1. **TensorRT**: Deep Learning Inference Engine
 2. DALI: Input Data Processing
 3. DeepStreamSDK: Deep Learning for Video Analytics
 4. AI-Assisted Annotation SDK: AI enabled Annotation for Medical Imaging
 5. Transfer-Learning Toolkit



Front End

Graph Optimizer

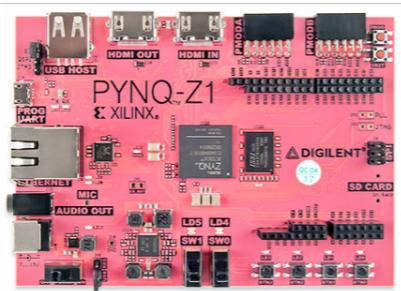
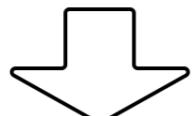
Tensor Optimizer

VTA JIT Runtime

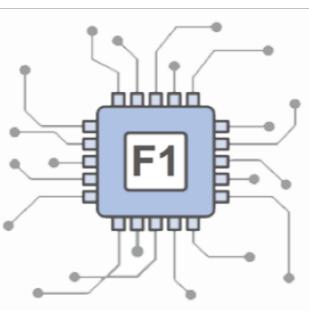
VTA ISA

VTA Micro-Architecture

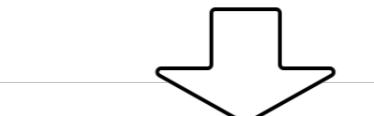
Edge FPGA



Cloud FPGA (in progress)



Simulator



Compiler and
Code Generation

Hardware
Design

Hardware
Deployment

Complex Demo: PyTorch, Clang, Git

Future of C++

Proposal P1019.

Executors with STL Algo

```
auto sum = std::reduce(std::execution::par, data);  
  
std::execution::static_thread_pool pool(4);  
sum = std::reduce(par.on(pool.executor()), data);  
  
my_gpu_executor my_executor;  
sum = std::reduce(par.on(my_executor), data);
```

Proposal P0009.

mdspan: K-Dimensional Range

```
vector<float> own(X * Y);
mdspan<float, dimensions<dyn, dyn>> ref(data.data(), X, Y);

expect(ref.rank() == 2);
expect(ref.size() == X * Y);
expect(ref.extent(0) == X);
expect(ref.extent(1) == Y);
expect(ref.is_regular() == true);
expect(ref.stride(0) == 1);
expect(ref.stride(1) == 1);
expect(ref.span() == X * Y);

for (int i = 0; i < ref.extent(0); i++) {
    for (int j = 0; j < ref.extent(1); j++) {
        ref(i, j) = rand();
    }
}
```

Proposal P1436.

Affinity in Executors

```
// Current platform-dependant versions:  
// Solaris: pbind()  
// Linux: sched_setaffinity()  
// Windows: SetThreadAffinityMask()  
  
enum class bulk_execution_affinity {  
    none, balanced, scatter, compact  
};  
  
executor exec;  
auto exec_subs = execution::require(exec, execution::bulk,  
                                     bulk_execution_affinity.scatter);  
exec_subs.bulk_execute([](std::size_t i, shared s) {  
    func(i);  
}, 8, shared_factory);
```

Proposal P0798.

Monadic std::optional

```
std::optional<image> get_einstein_pic(image const & img) {
    return crop_to_face(img)
        .and_then(add_bow_tie)
        .and_then(make_eyes_sparkle)
        .transform(add_glasses)
        .or_else([] {
            return [=] { std::cout << "Failed!"; };
        });
}
```

More proposals for numerics

- P1368: Multiplication and division of fixed-point numbers
- P0943: Support C atomics in C++
- P1438: A Rational Number Library for C++

Contacts



Ashot Vardanian

[linkedin.com/in/ashvardanian](https://www.linkedin.com/in/ashvardanian)
[fb.com/ashvardanian](https://www.facebook.com/ashvardanian)