

# The Future of Machine Learning Acceleration

Elliott Delaye  
Distinguished Engineer

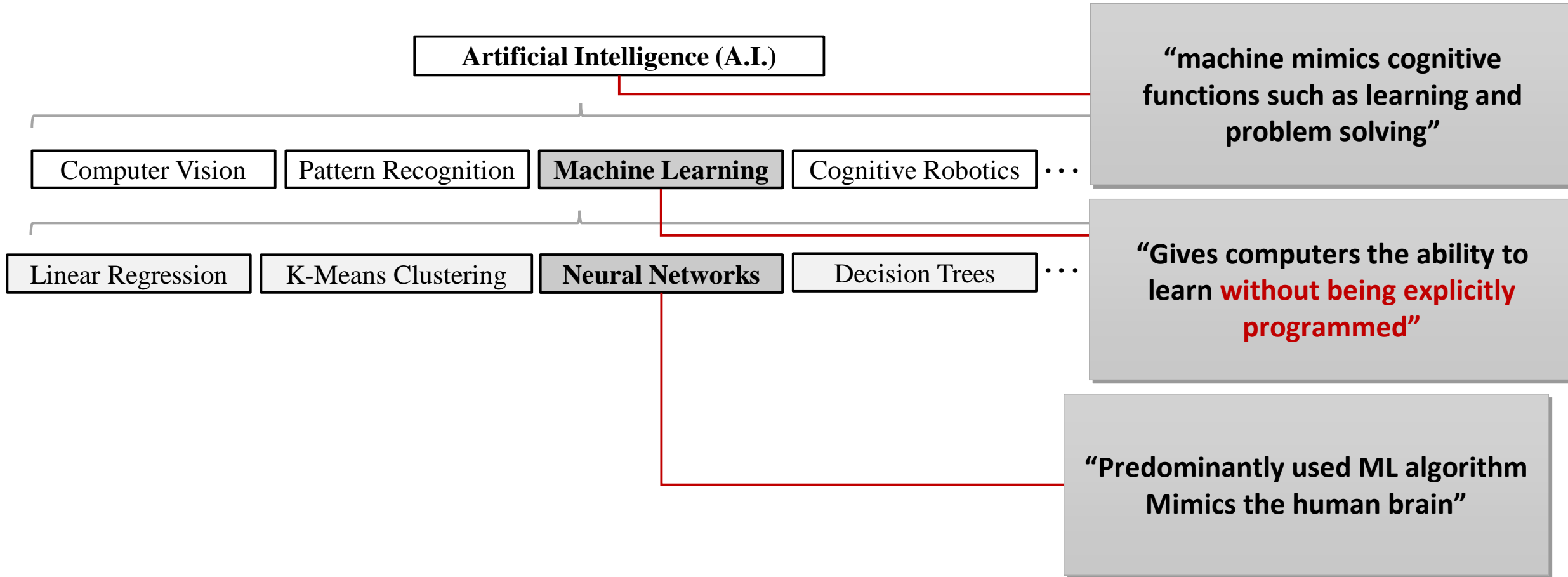


- > **Machine Learning - Neural Networks**
- > **Computation & Memory Requirements**
- > **Algorithmic Optimization Techniques**
- > **Hardware Architectures**

# Neural Networks



# A.I. – Machine Learning - Neural Networks



# Neural Networks (Deep Neural Networks, etc.)

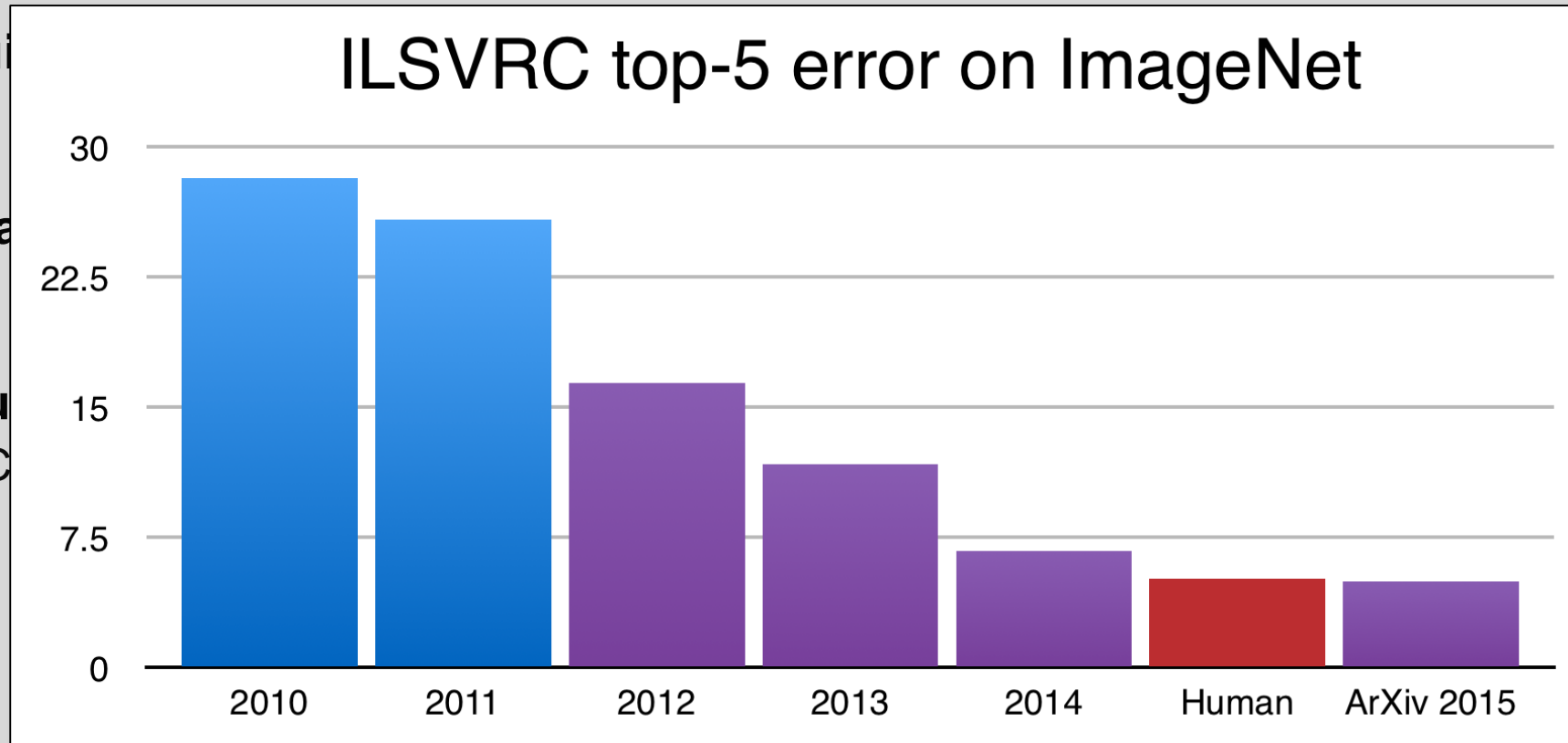
*Why are they so popular?*

> Require

> NNs a

> If you

>> C



place other

e simple rules can  
em

viously  
iters

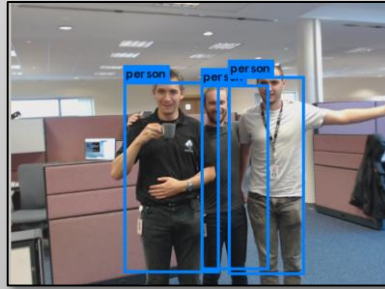
ely unsolved

networks old?

# Increasing Range of Applications



Image Classification



Object Detection



Semantic Segmentation

Computer Vision  
CNNs



Speaker  
Diarization



Speech  
Recognition

Speech Recognition  
RNNs, LSTMs



Translation

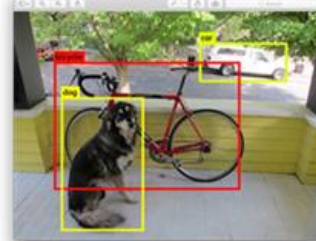


Sentiment Analysis

Natural Language Processing  
Sequence to sequence



Recommender



GamePlay

Many more emerging...

Others

# Popular Neural Networks

ResNet50, VGG,  
AlexNet, InceptionV3



Image Classification

Faster R-CNN,  
SSD, YoloV2



Object Detection

Mask-R-CNN,  
SSD

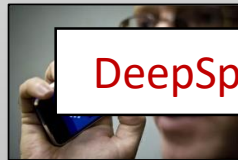


Semantic Segmentation

Computer Vision  
CNNs



Speaker  
Diarization

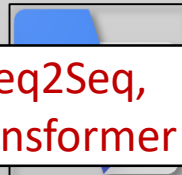


Speech  
Recognition

DeepSpeech2

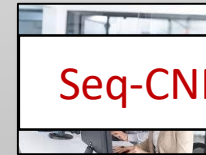
Speech Recognition  
RNNs, LSTMs

Seq2Seq,  
Transformer



Translation

Seq-CNN



Sentiment Analysis

Natural Language Processing  
Sequence to sequence

NCF



Recommender

MiniGo,  
DeepQ, A3C

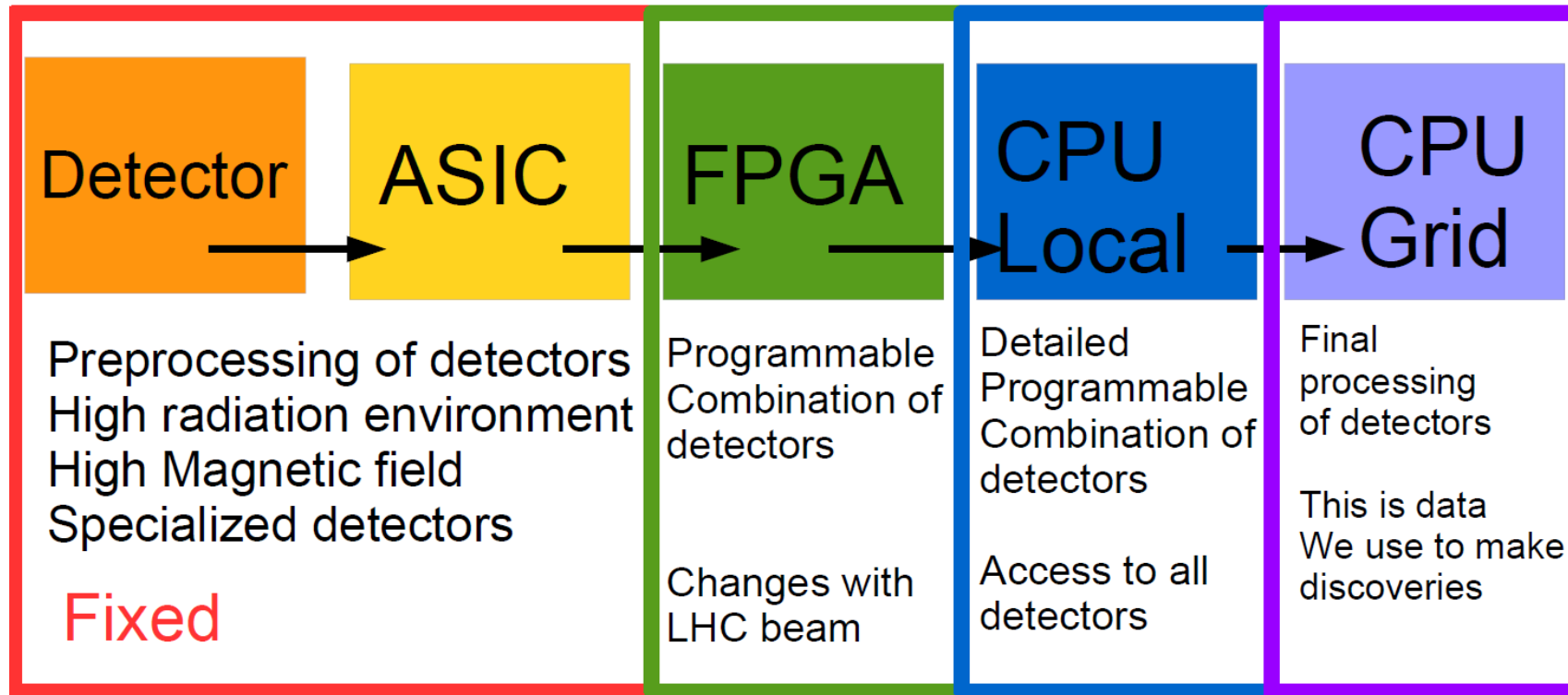
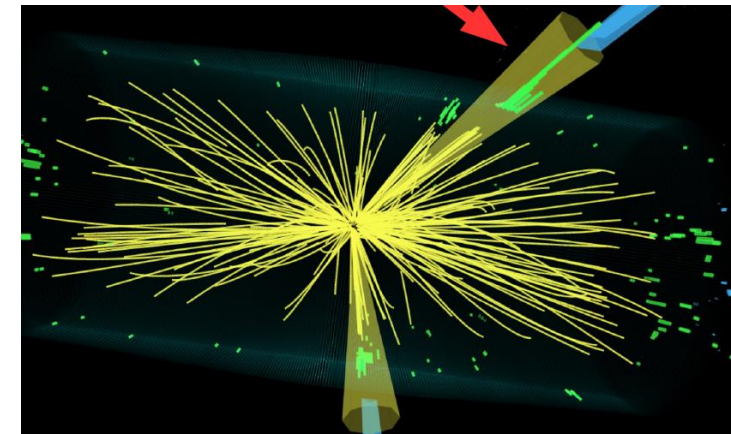


GamePlay

Others

# New uses – High Energy Physics

Collision rate is 40 MHz  
A new collision every 25ns



Latency : 10µs

100ms

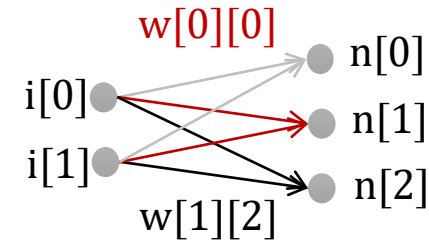
10s



# Convolutional Neural Networks (CNNs)

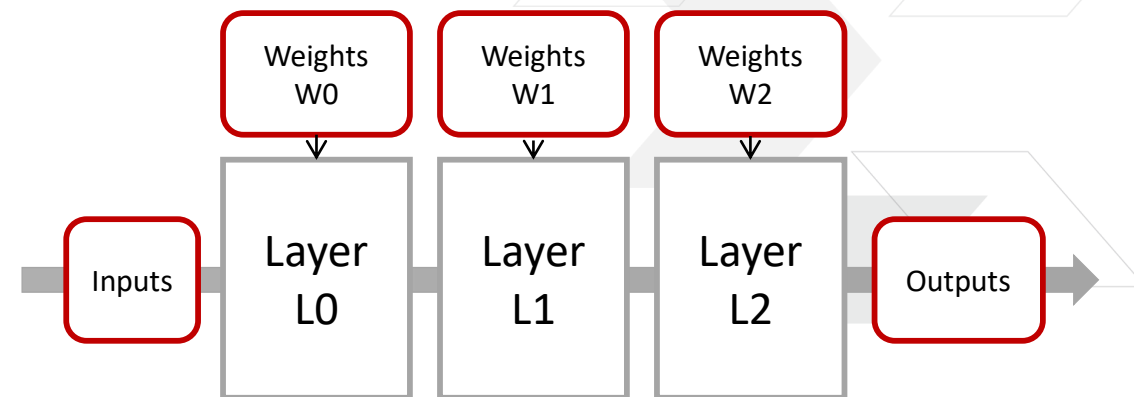
## *from a computational point of view*

- > CNNs are usually feed forward\* computational graphs constructed from one or more layers
  - >> Up to 1000s of layers
- > Each layer consists of “neurons”  $n[i]$  which are interconnected with synapses, associated with weights  $w[i][j]$
- > Each neuron computes:
  - >> Typically a dot-product (multiply-accumulate)
  - >> Followed by a non-linear “activation” function
  - >> Without non-linear function, L0, L1, L2 could be collapsed into a single layer



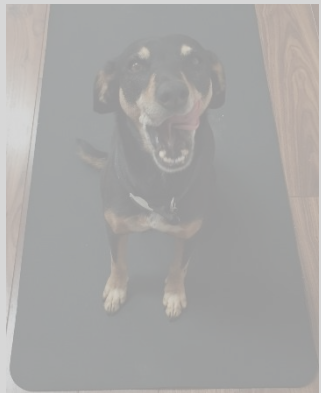
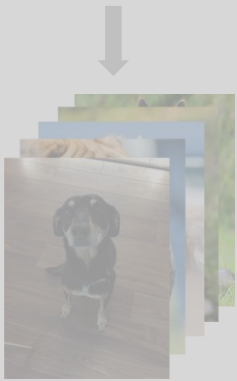
“Synapse” with weight  $w[j][i]$   
“Neuron”  $n[i]$

$$n[0] = \text{Activation}(w[0][0]*i[0] + w[1][0]*i[1])$$

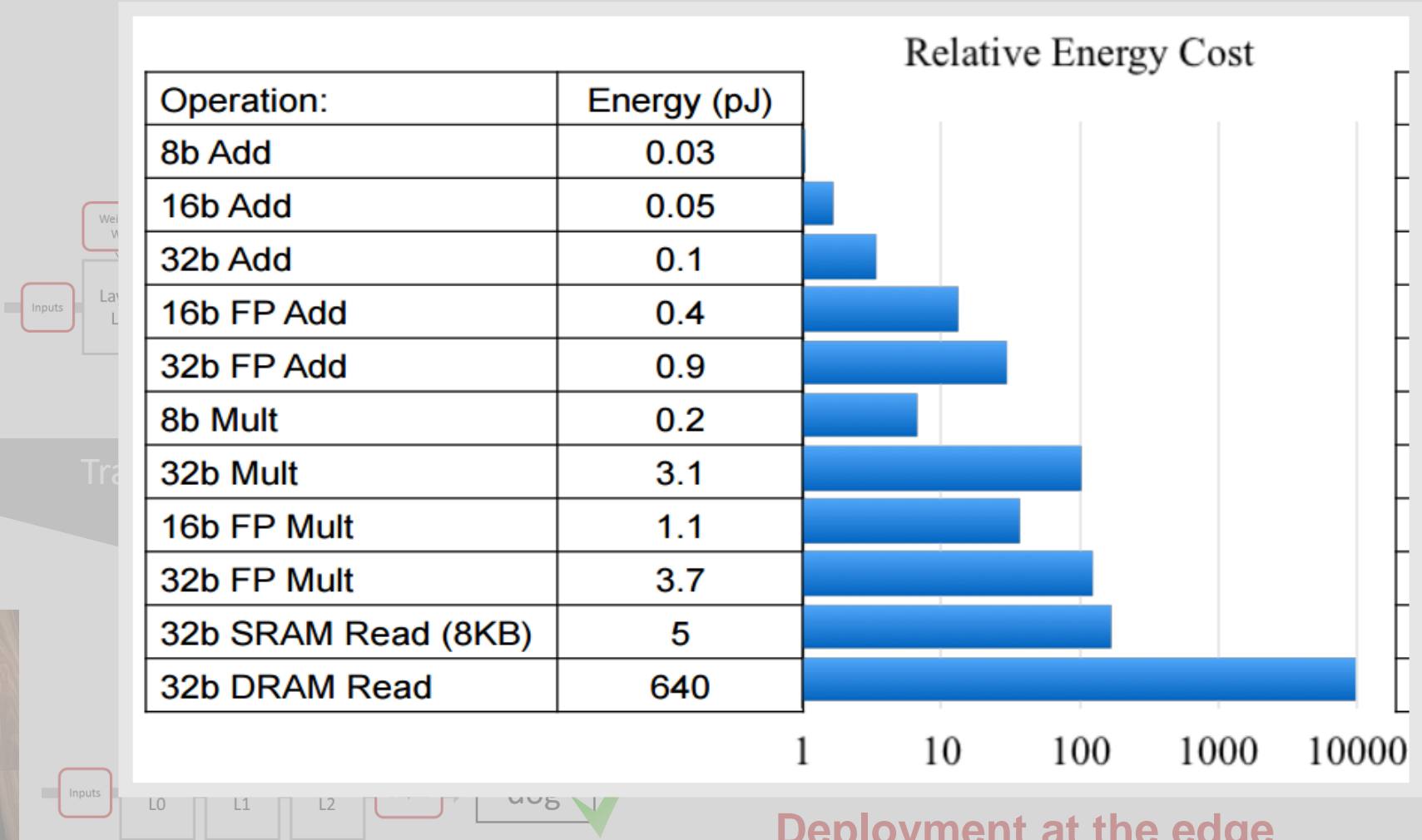


# From Training to Inference

Training dataset

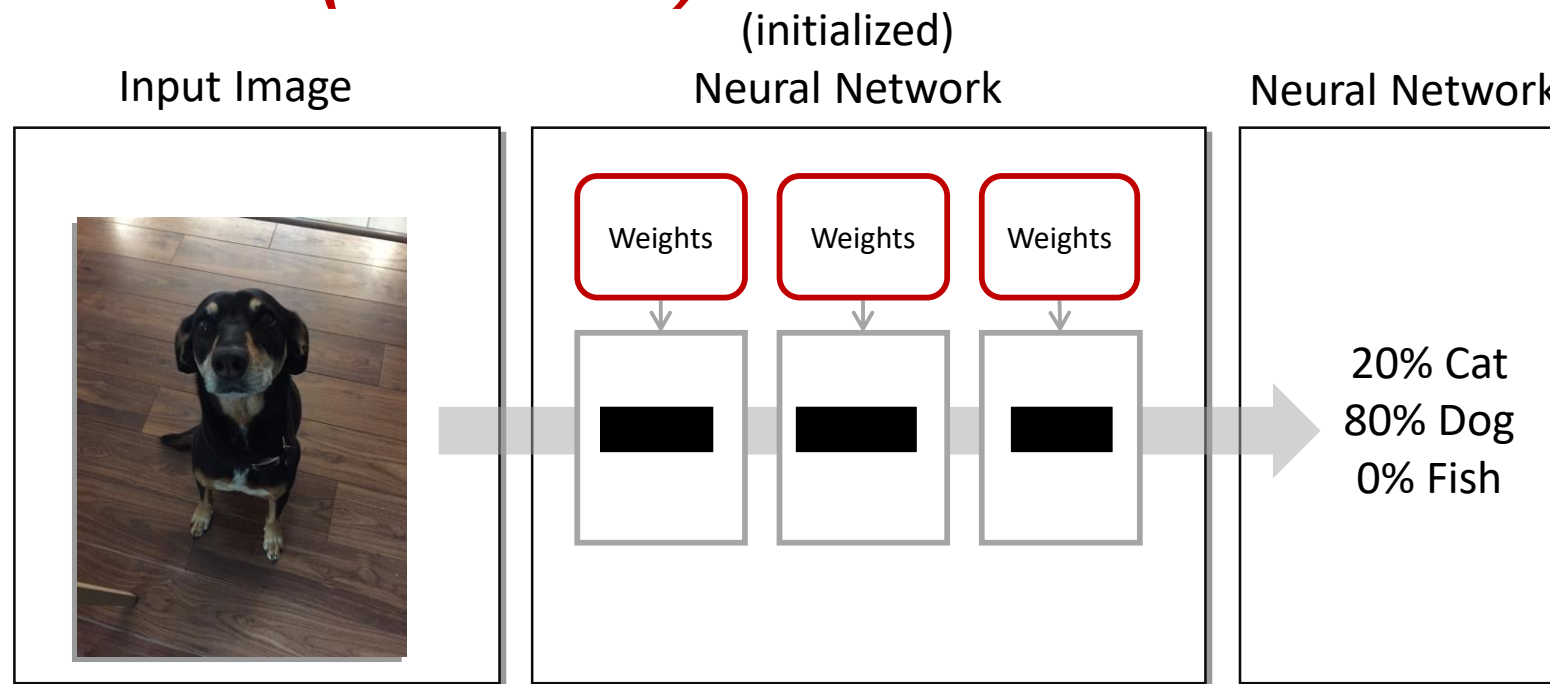


>> 10



# Example: ResNet50

## *Forward Pass (Inference)*



For ResNet50:

70 Layers

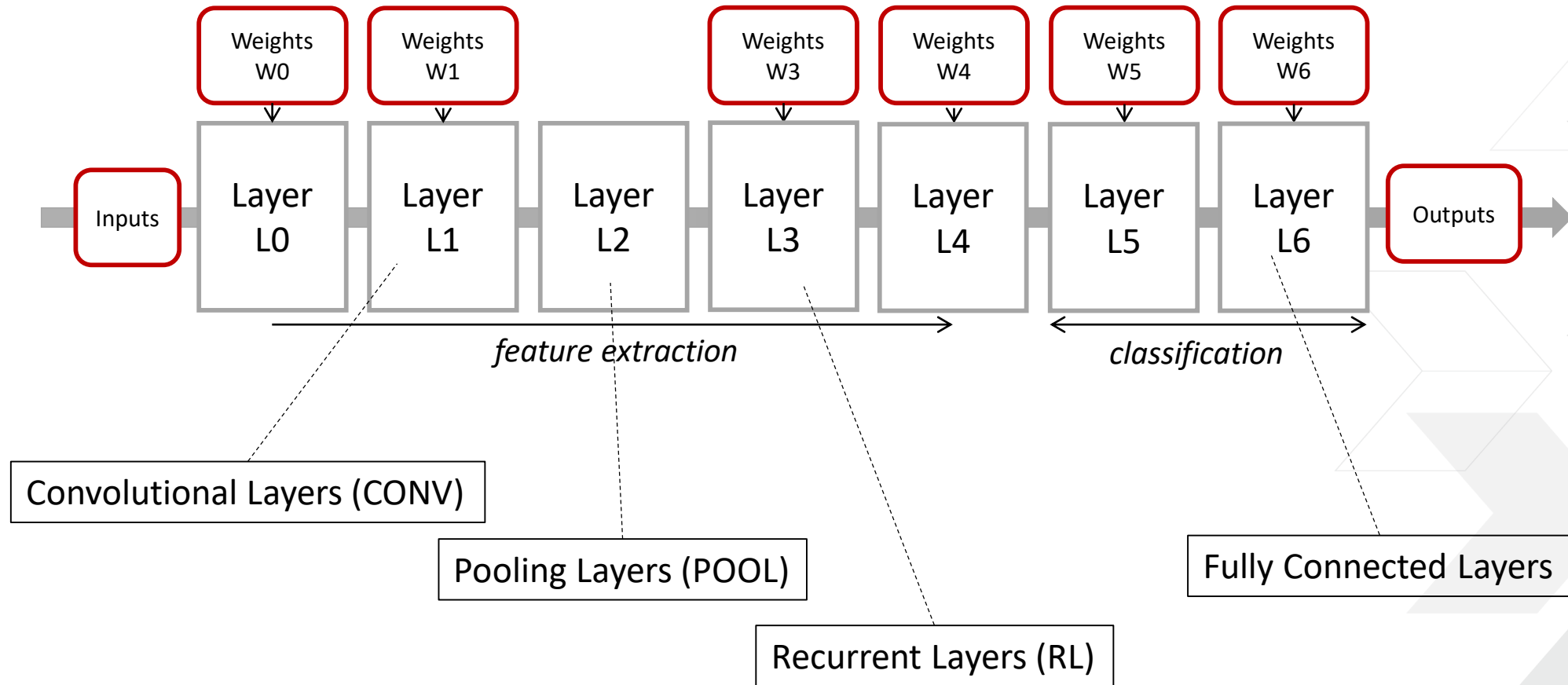
7.7 Billion operations

25.5 MBytes of weight storage\*

10.1 MBytes for activations\*

*\*Assuming int8*

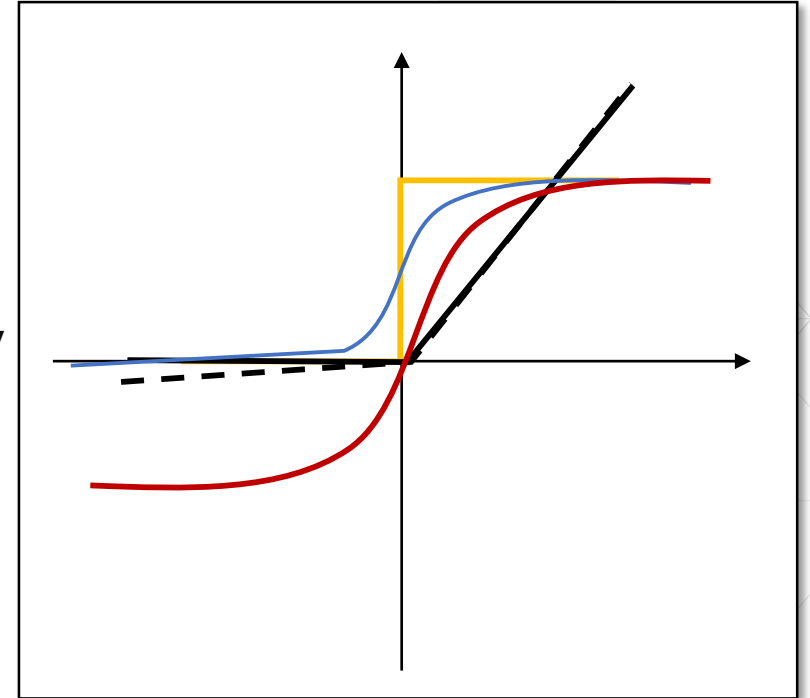
# NNs in More Detail



Activation & Batch Normalization (RELU, PRELU, BN, SCALE, etc.)

# Activation Functions

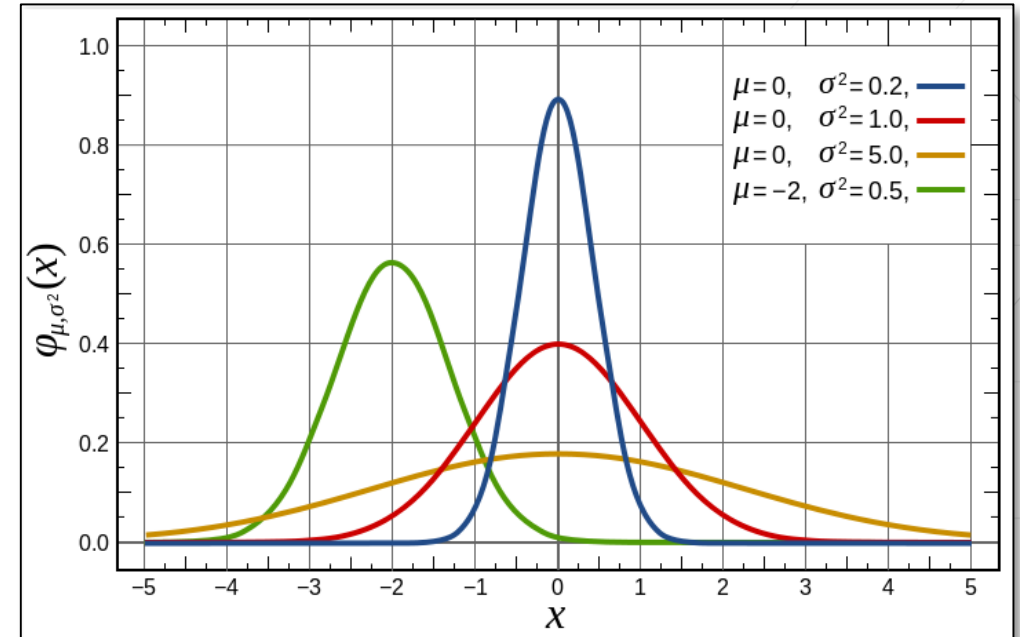
- > **Implements the concept of “Firing”**
  - >> Non-linear so we can approximate more complex functions
- > **Most popular for CNN: Rectified Linear Unit (ReLU)\***
  - >> Popular as it propagates gradients better than bounded and easy to compute,  $x = \max(0, x)$
- > **Other common ones include: **tanh**, leaky ReLU, **sigmoid**, **threshold functions** for quantized neural networks**
- > **Implementation:**
  - >> Support for special functions as well as some level of flexibility



*\*Nair, V. and Hinton, G.E., 2010. Rectified linear units improve restricted boltzmann machines. In Proceedings of the 27th international conference on machine learning (ICML-10) (pp. 807-814).*

# Batch Normalization

- > Normalizes the statistics of activation values across layers
- > Significantly reduces the training time of networks, can improve accuracy and makes it less sensitive to initialization
- > Compute:
  - >> Lightweight at inference
  - >> Heavy duty during training
    - Subtract mean, divide by standard deviation to achieve zero-centered distribution with unit variance



[https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution)

# Fully Connected Layers

(aka inner product or dense layers)

> Each input activation is connected to every output activation

>> Why it's termed "Fully Connected"

> Can be written as a matrix-vector product with an element-wise non-linearity applied afterwards.

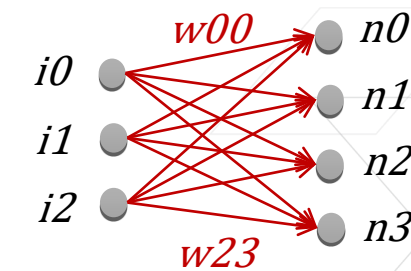
> Implementation Challenges

>> Connectivity

>> High weight memory requirement: #IN \* #OUT \* BITS

>> Low arithmetic intensity assuming weights off-chip

$$2 * \#IN * \#OUT / \#IN * \#OUT * BITS/8$$



$$\begin{bmatrix} i_0 & i_1 & i_2 \end{bmatrix} \times \begin{pmatrix} W_{00} & W_{01} & W_{02} & W_{03} \\ W_{10} & W_{11} & W_{12} & W_{13} \\ W_{20} & W_{21} & W_{22} & W_{23} \end{pmatrix} = \begin{bmatrix} n_0' & n_1' & n_2' & n_3' \end{bmatrix}$$

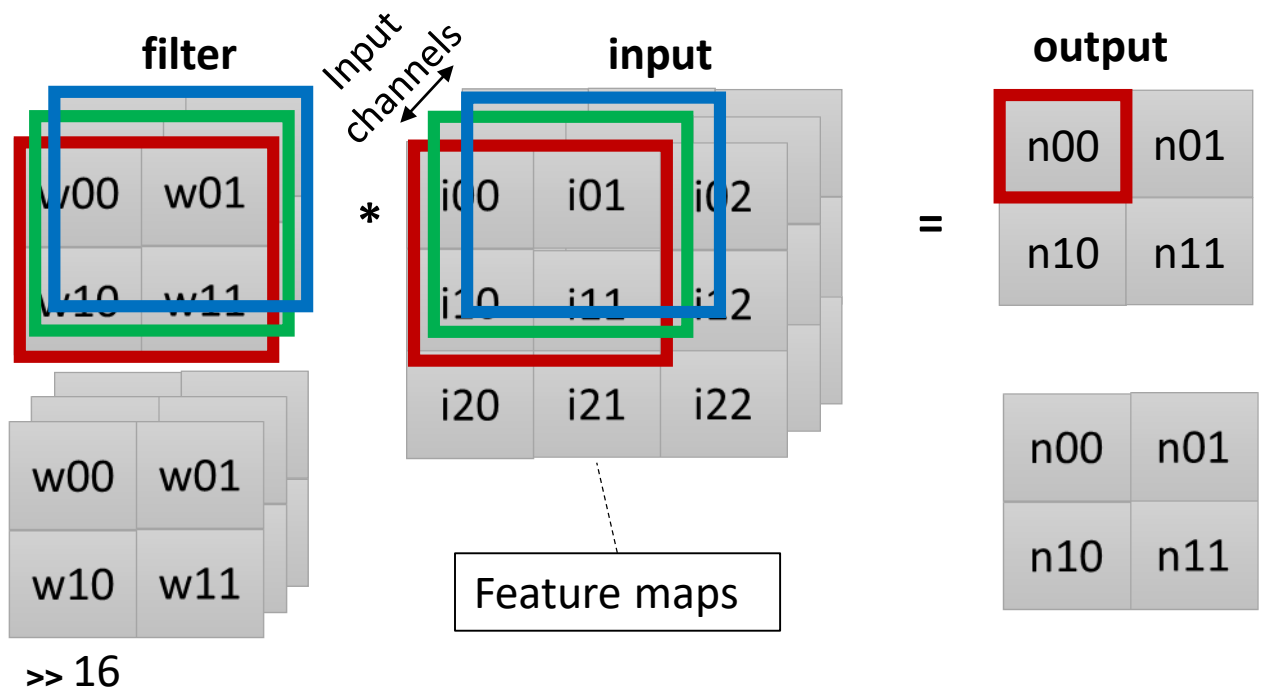
$$(n_0 \ n_1 \ n_2 \ n_3) = Act(n_0' \ n_1' \ n_2' \ n_3')$$

MODEL	CONV WEIGHTS (M)	FC WEIGHTS (M)
ResNet50	23.454912	2.048
AlexNet	2.332704	58.621952
VGG16	14.710464	123.633664

# Convolutional Layers

## Example 2D Convolution

- > Convolutions capture some kind of locality, spatial or temporal, that we know exists in the domain
- > Input of each “neuron” reduced – Be
  - >> Applying convolution to all images in the previous layer
- > Weights represent the filters used for convolutions



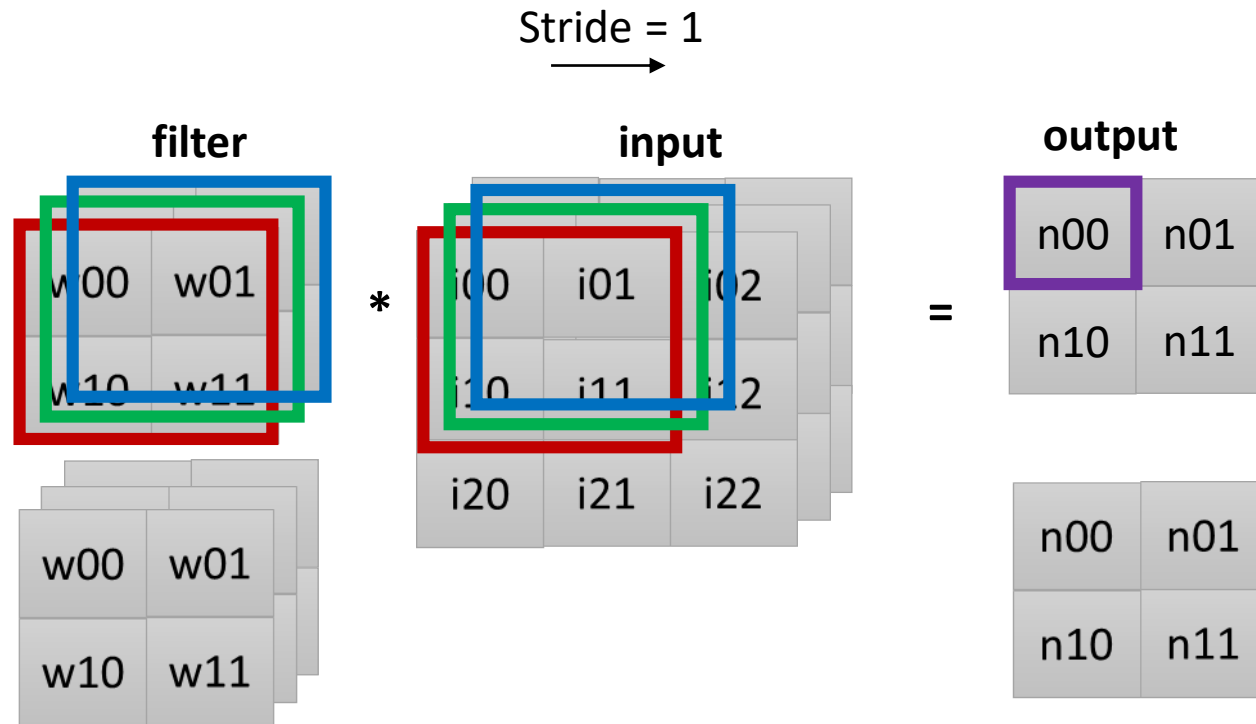
$$n_{00} = Act(w_{00} * i_{00} + w_{01} * i_{01} + w_{10} * i_{10} + w_{11} * i_{11} + w_{00} * i_{00} + w_{01} * i_{01} + w_{10} * i_{10} + w_{11} * i_{11} + w_{00} * i_{00} + w_{01} * i_{01} + w_{10} * i_{10} + w_{11} * i_{11})$$

Input channel 0



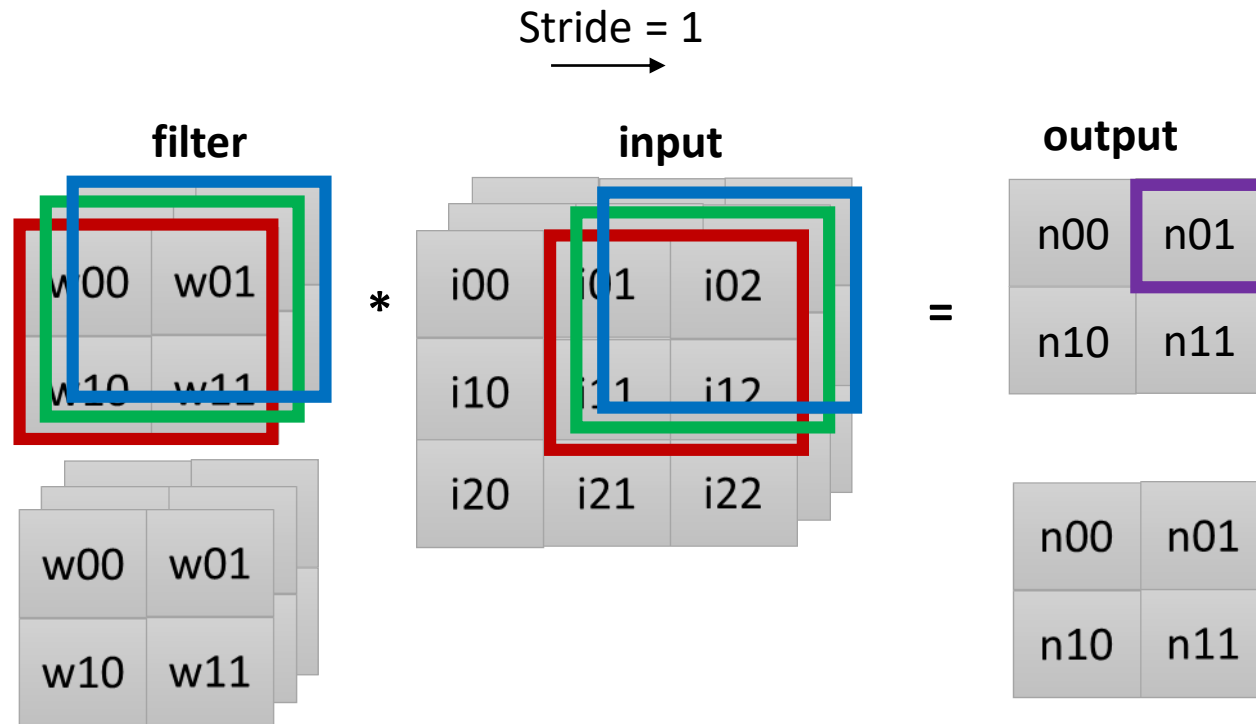
# 2D Convolutional Layers

- > Slide the window till one feature map is complete
  - >> With a given stride size



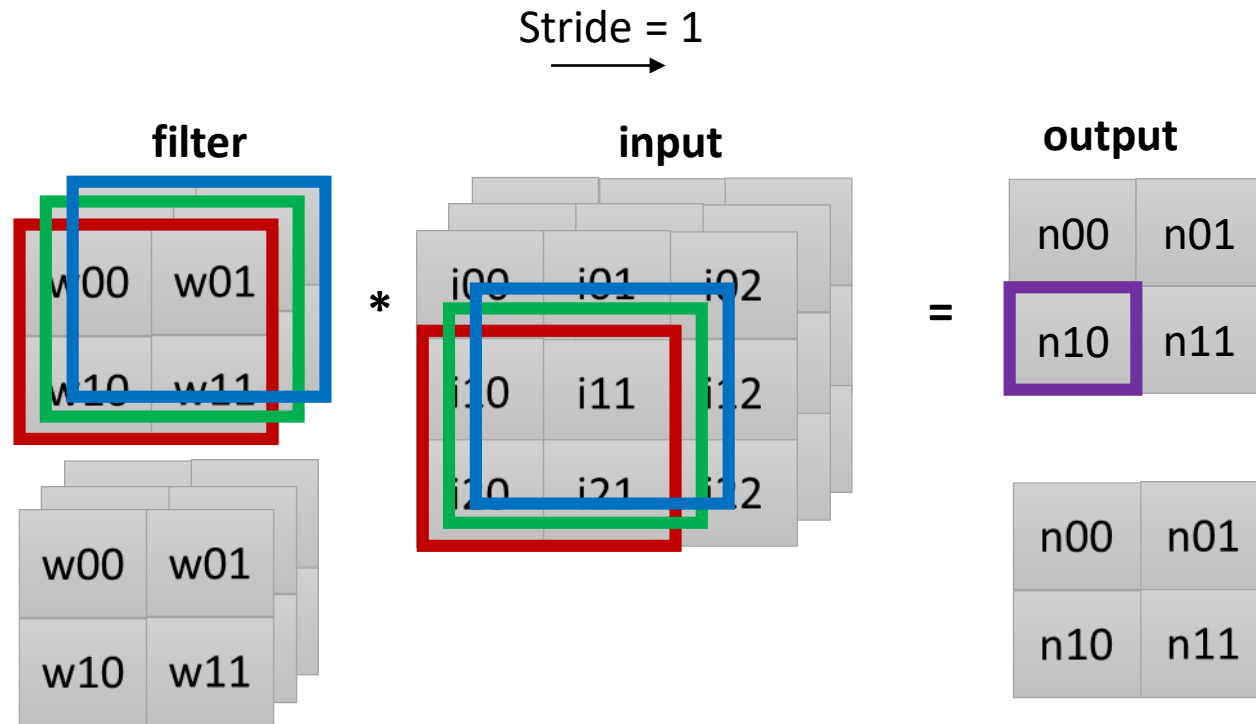
# 2D Convolutional Layers

- > Slide the window till one feature map is complete
  - >> With a given stride size



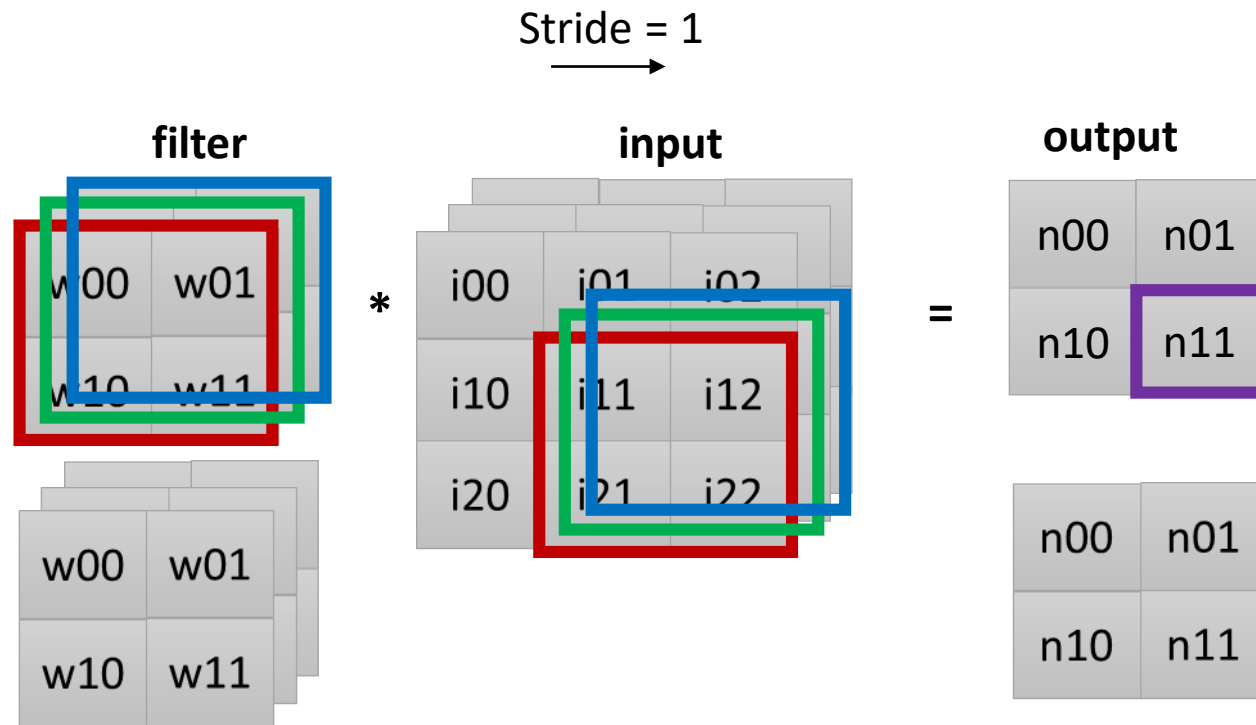
# 2D Convolutional Layers

- > Slide the window till one feature map is complete
  - >> With a given stride size



# 2D Convolutional Layers

- > Slide the window till one feature map is complete
  - >> With a given stride size

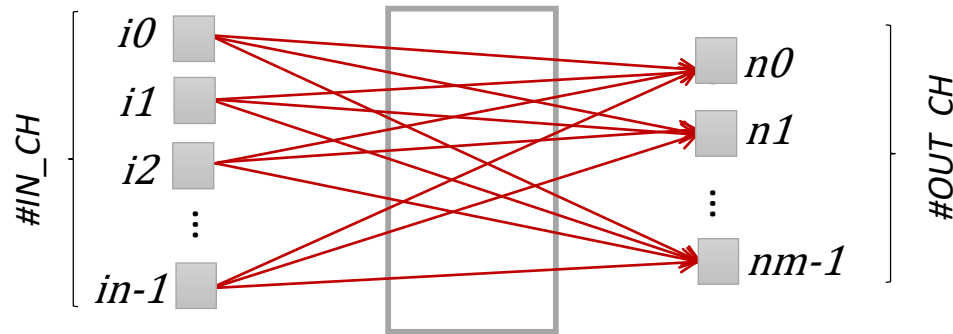


# Convolutions

## Challenges

### > Channel connectivity issue

- >> Every input channel information broadcasts to every output channel



100s to 1000 channels

### > Huge amounts of compute

- >> Dense convolutions account for the majority of the compute
- >> Image based so inputs could be HD video, 1920x1080

### > Lots of memory used

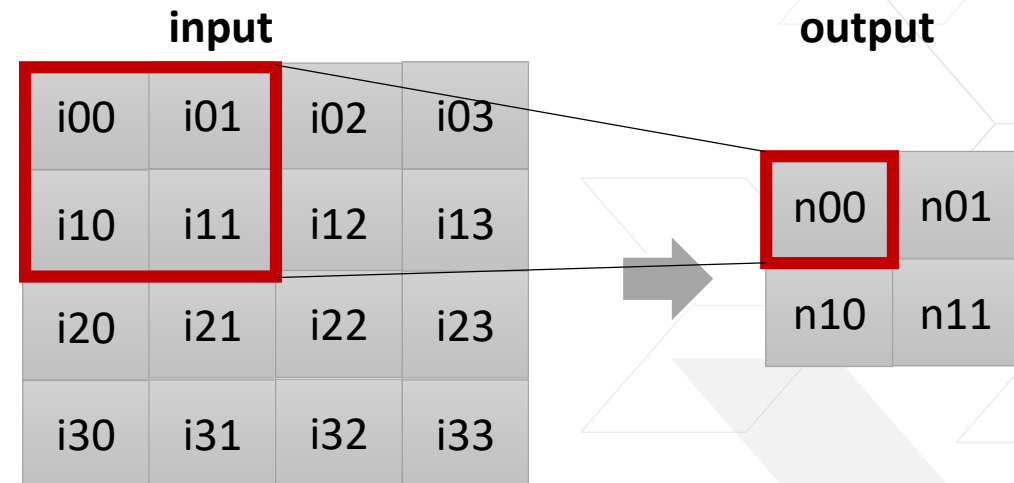
- >> 1920x1080x3 = 6.2MB
- >> YOLOv2 First Output, 1920x1080x32 = 66.3MB

MODEL	CONV [GOPS]	FC [GOPS]
ResNet50	7.712	0.004
AlexNet	1.332	0.044
VGG16	30.693	0.247

# Pooling Layer

- > Down-samplers of images
- > Reduces compute in subsequent layers
- > May use MAX or AVERAGE
- > Compute:
  - >> Low amount of compute
  - >> Potentially replaceable with larger strides in previous convolution

Max pool with 2x2 filters and stride of 2:



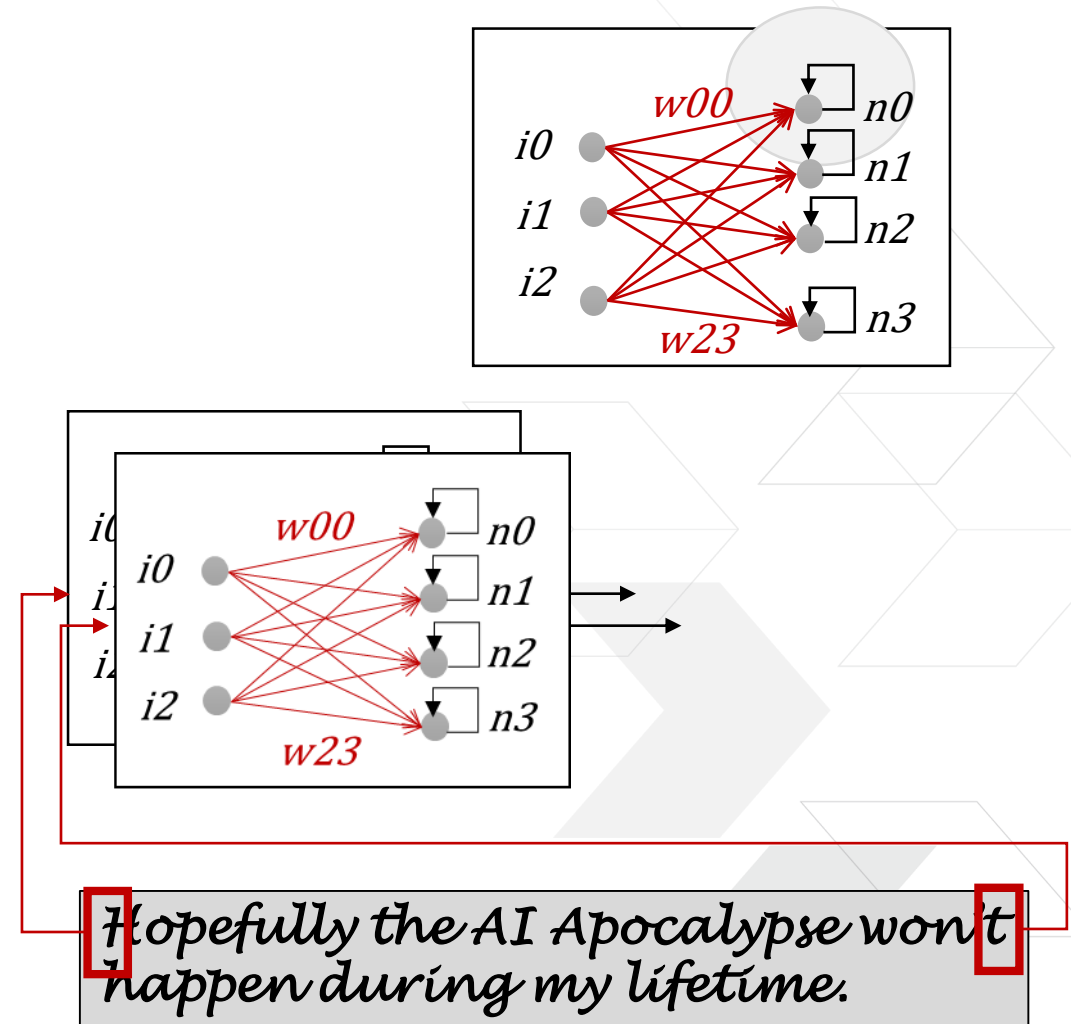
$$n_{00} = \text{MAX}(i_{00}, i_{01}, i_{10}, i_{11})$$

Or

$$n_{00} = \text{AVG}(i_{00}, i_{01}, i_{10}, i_{11})$$

# Recurrent Layer Types

- > **Contain state for processing sequences**
  - >> For example needed in speech or optical character recognition
  - >> “Apocal???”
- > **Uni-directional or bi-directional**
  - >> “I ????? You”
- > **More sophisticated types to address the vanishing gradients problem for learning more than 5-10 timesteps**
  - >> GRU (gated recurrent unit)
  - >> LSTM (long short term memory)



# Recurrent Layers

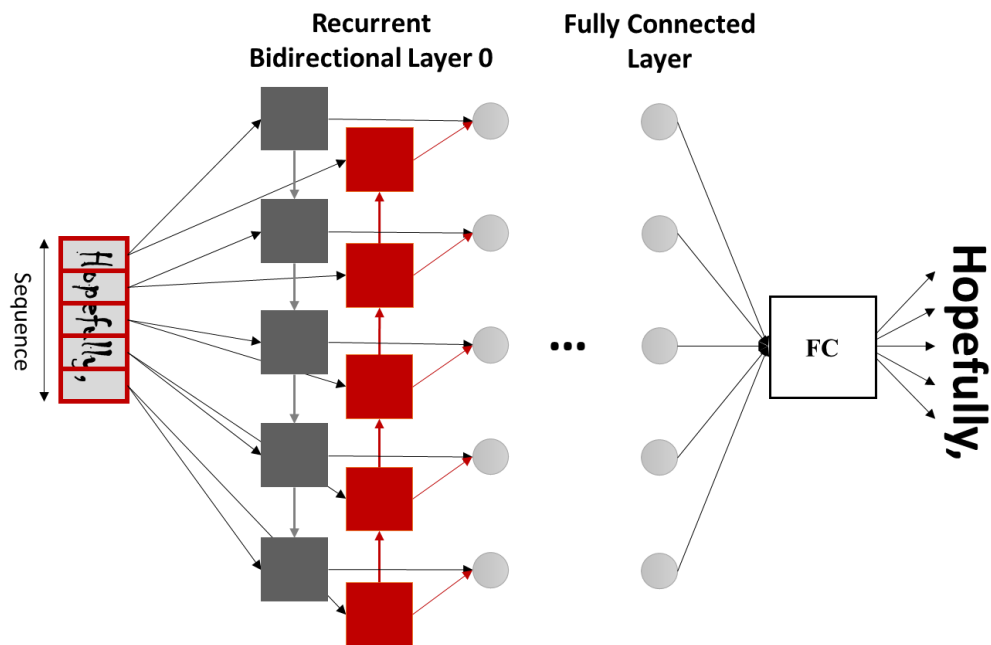
## *Challenges in Additional Data Dependencies*

### > Input sequence

- >> Unlike batch, additional data dependencies between inputs of the same sequence and state

### > Bi-directional NNs

- >> Full sequence needs to be completed before the next layer.
- >> Fewer opportunities for parallelism
- >> Operations generally stress memory bandwidth (compared to CONV)



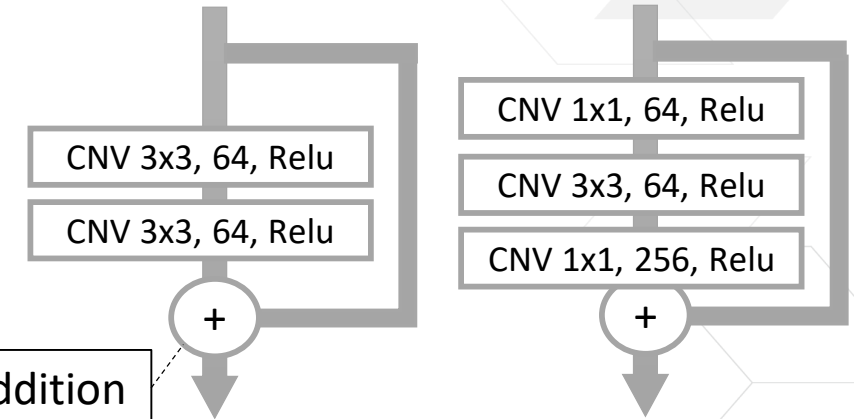


# Meta-Layers

## > Residual layers (ResNets \*)

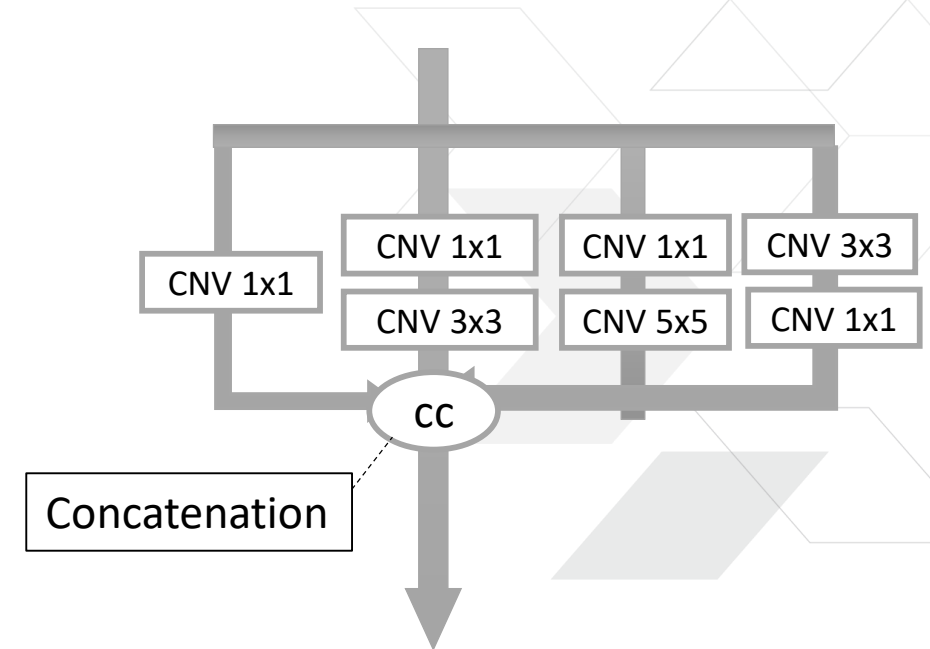
- >> Introduced to make larger networks more trainable
- >> Better gradient propagation through skip connections during training
- >> Plus 1x1 convolutions to reduce dimensionality and save compute

Elementwise addition



## > Inception layers (GoogleNet\*\*)

- >> Huge variation in spatial features => combining different size convolutions in one layer
- >> Plus 1x1 convolutions to reduce dimensionality and save compute
- >> Later on additional factorization to reduce compute
  - 3x3 = 1x3 and 3x1



## > Many more...

## > Implementation: support for non-linear topologies!

# Computation & Memory Requirements



# Compute and Memory Requirements

## Architecture Neutral, Per Layer

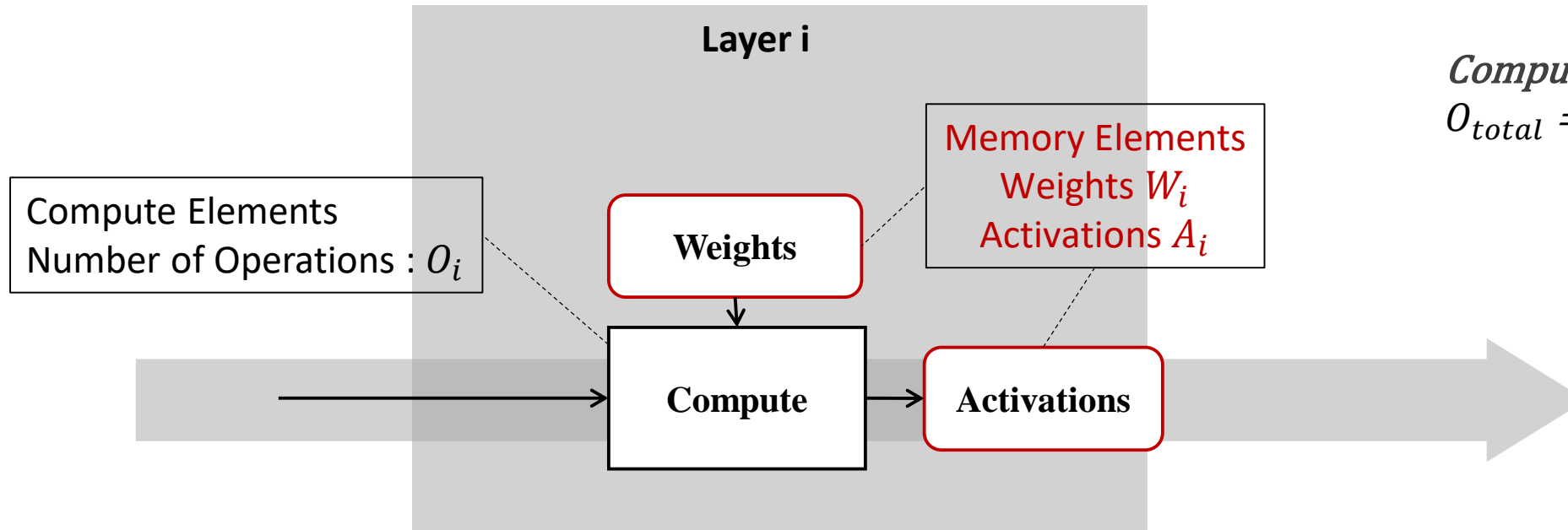
*Memory Requirements:*

$$A_{total} = \sum A_i$$

$$W_{total} = \sum W_i$$

*Compute Requirements:*

$$O_{total} = \sum O_i$$



*IN, IN\_CH:*

*number of inputs and input channels*

*OUT, OUT\_CH:*

*number of outputs and output channels*

*F\_DIM, FM\_DIM:*

*filter and feature map dimensions (assumed square)*

*BATCH:*

*batch size*

*BITS:*

*bit precision in data types*

*GATES:*

*number of gates in RNNs:*

*STATES:*

*worst case*

*SEQ:*

*sequence length*

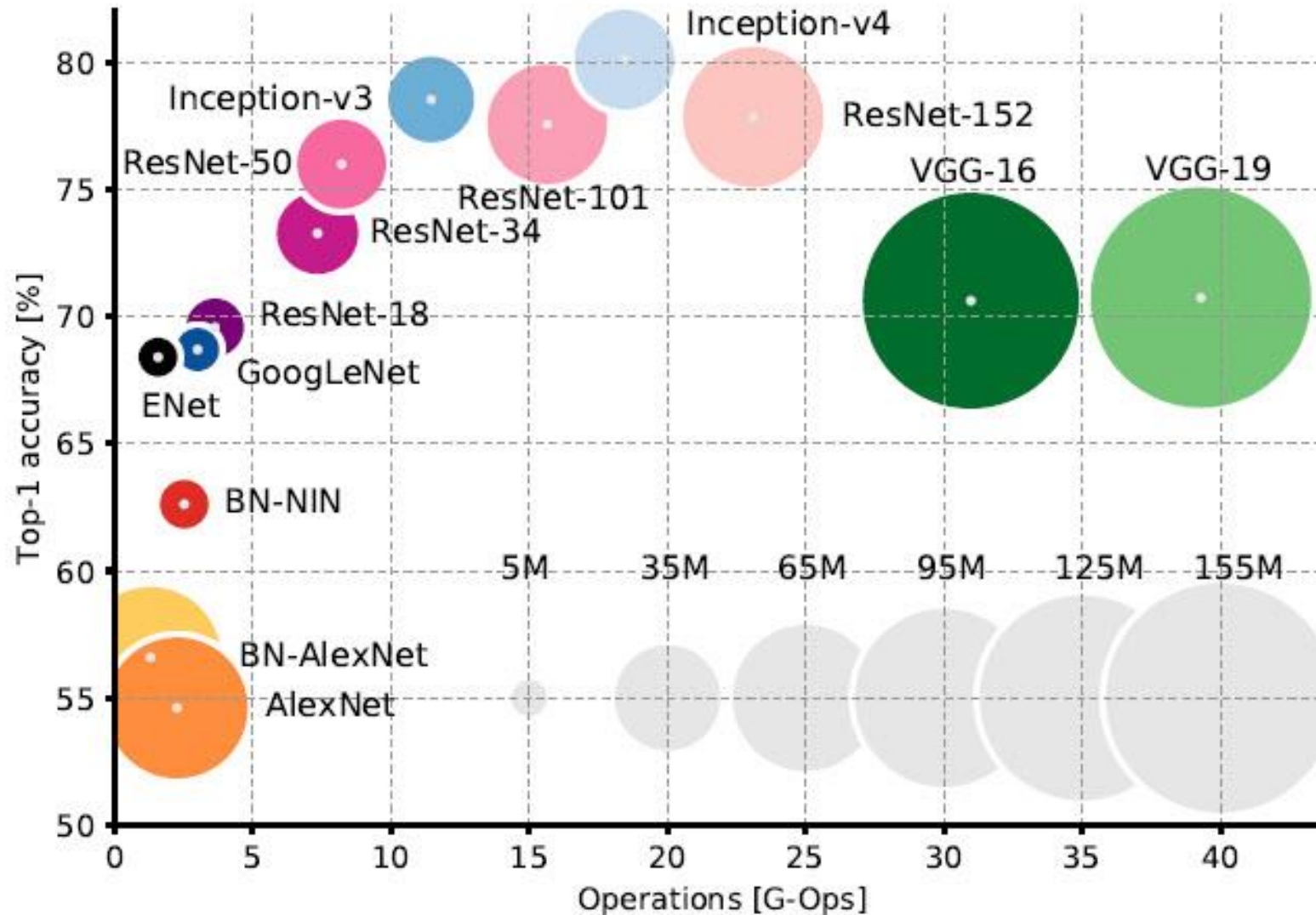
*HID:*

*hidden size (state + output from each state)*

*DIRS:*

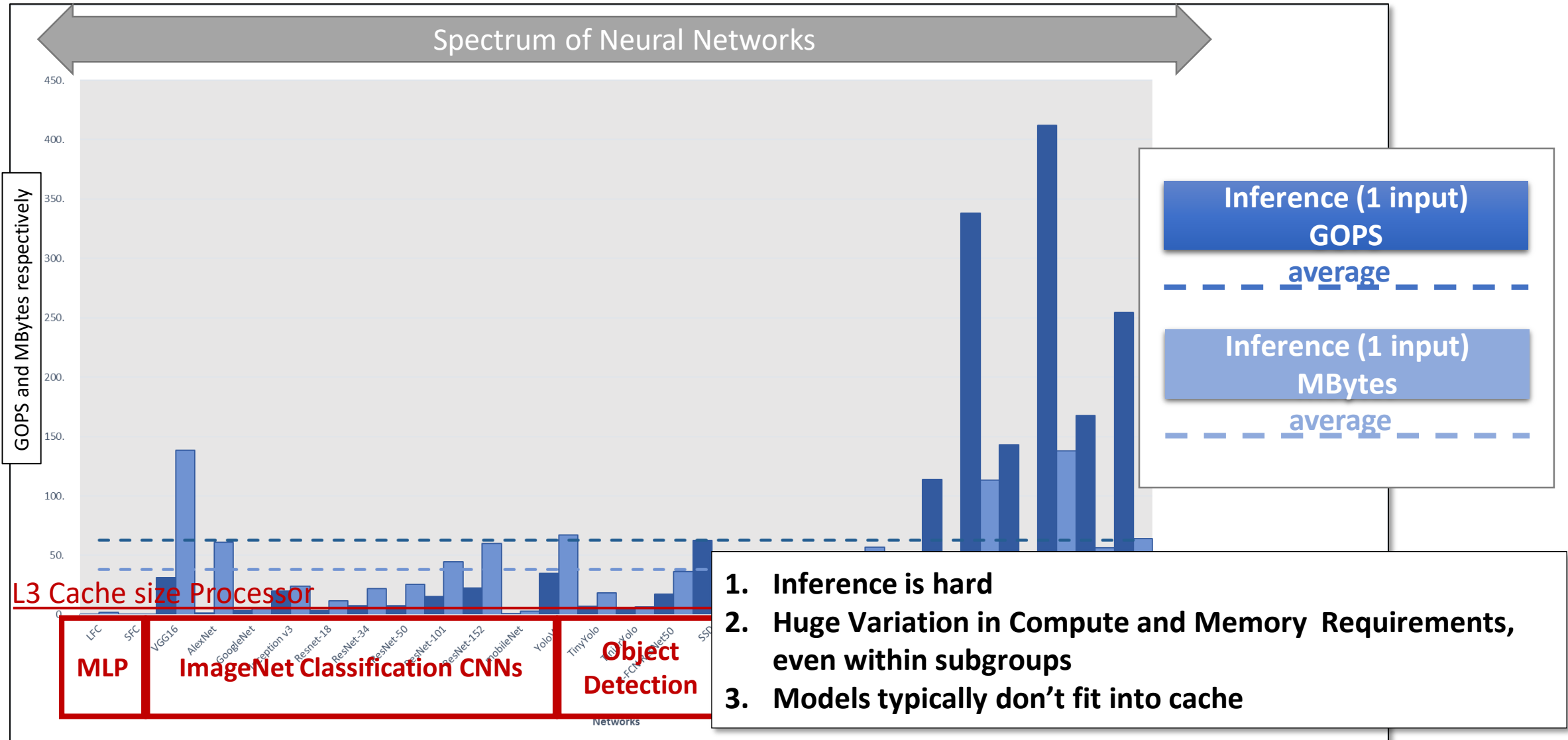
*1 for unidirectional and 2 for bidirectional RNN*

# Network Complexity by Operations & Weights



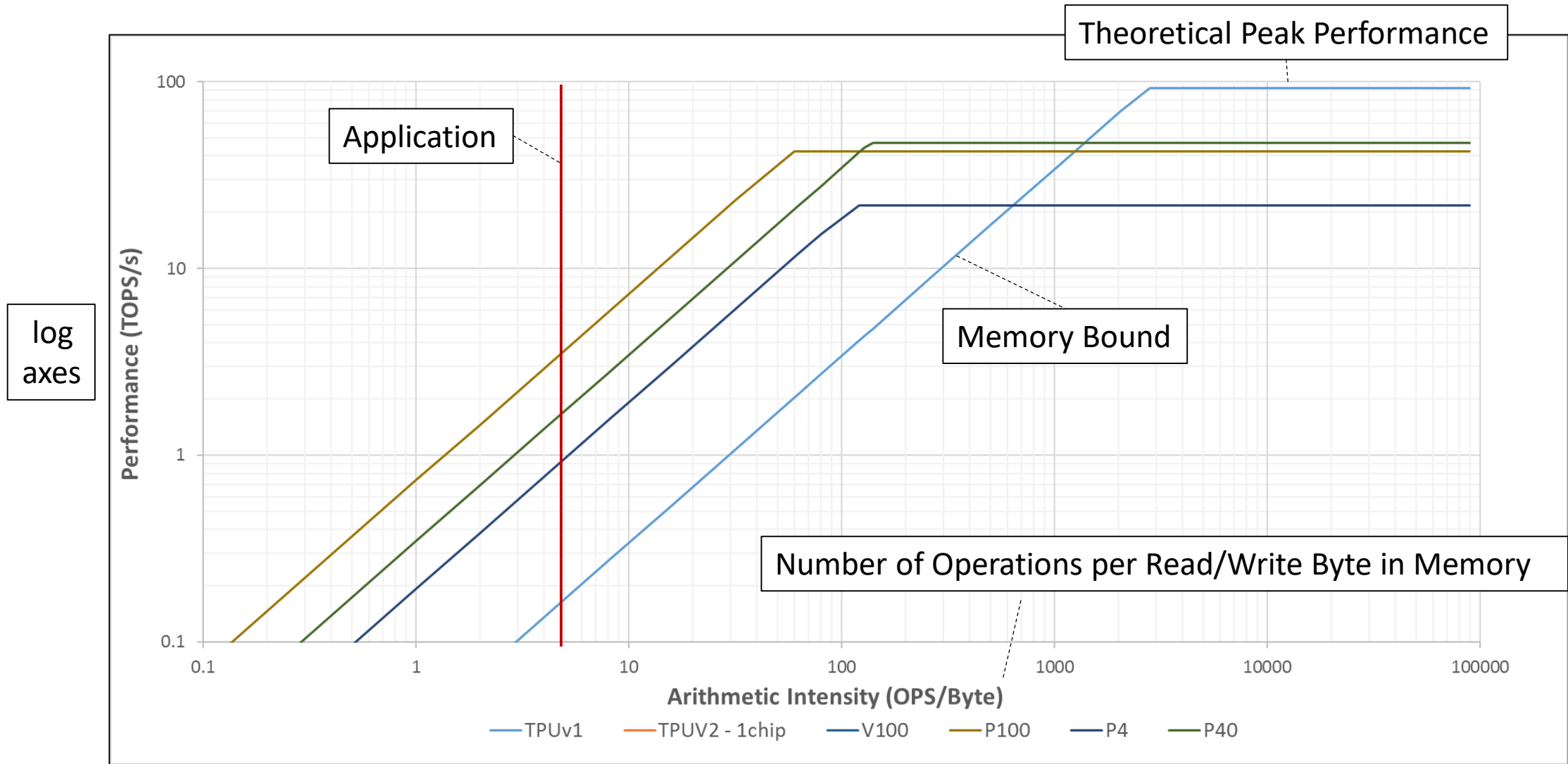
# Inference Compute and Memory Across a Spectrum of Neural Networks

\*architecture independent  
 \*\*1 image forward  
 \*\*\* batch = 1  
 \*\*\*\* int8



1. Inference is hard
2. Huge Variation in Compute and Memory Requirements, even within subgroups
3. Models typically don't fit into cache

# Rooflines\*

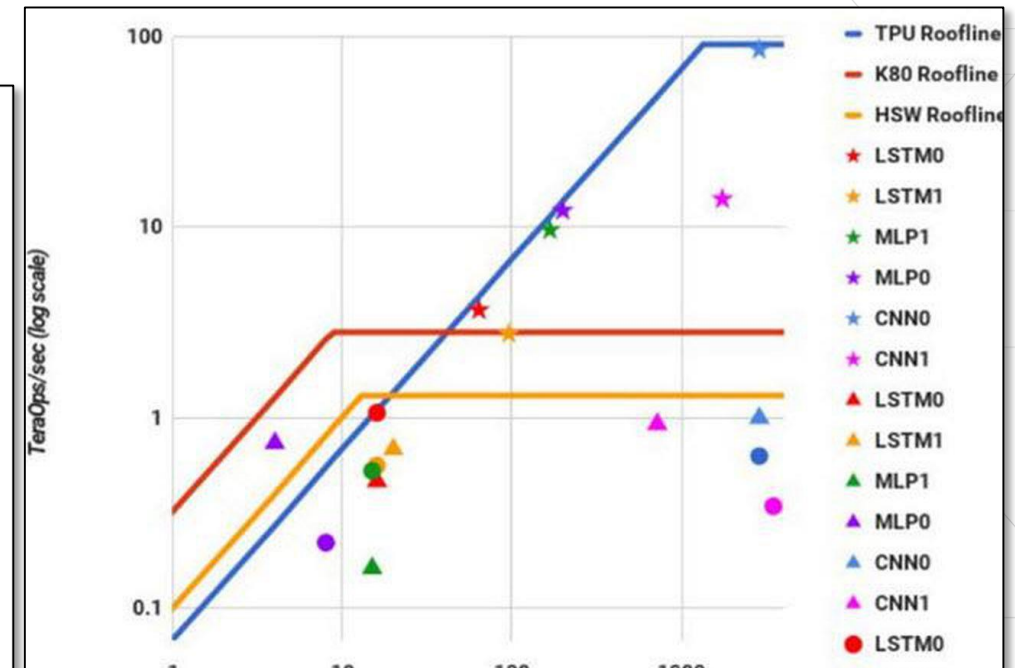
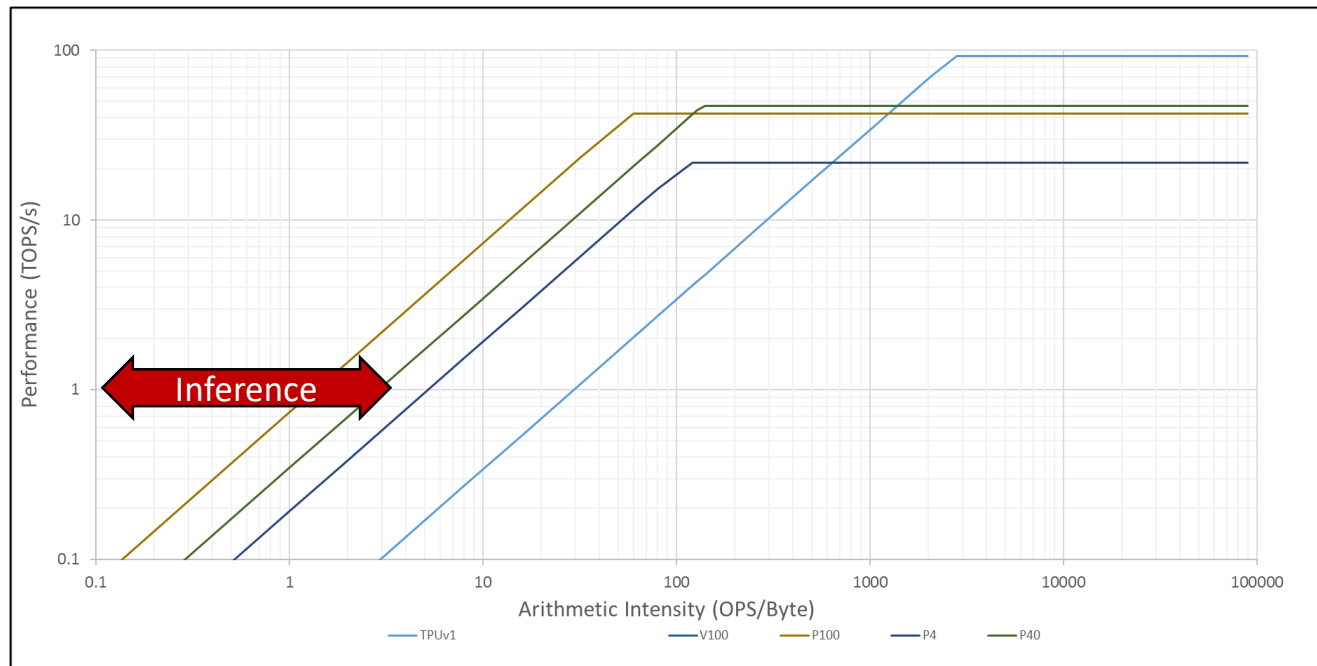


# Arithmetic Intensity

## Across a Spectrum of Neural Networks

\* batch = 1  
 \*\* with respect to weights assuming weights are off-chip

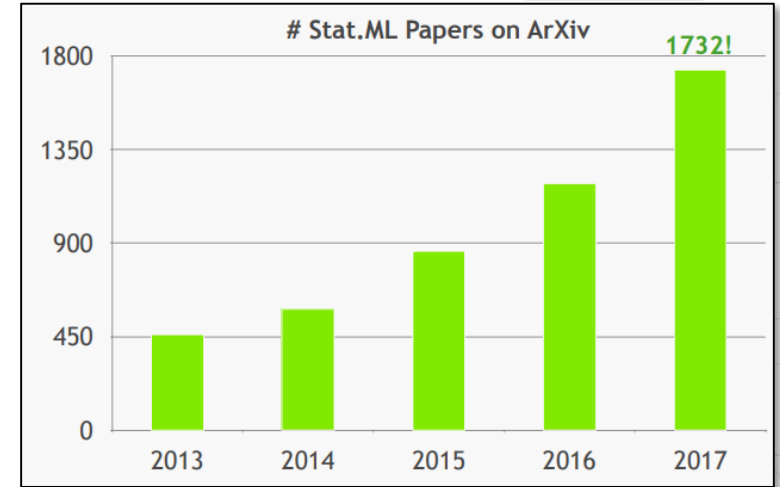
- > Memory requirement for weights, activations are beyond typically available on-chip memory
- > This yields low arithmetic intensity
  - >> For example for inference, assuming weights off-chip and naïve implementation, majority of networks is below 6OPS:Byte



Jouppi, N.P., Young, C., Patil, N., Patterson, D., Agrawal, G., Bajwa, R., Bates, S., Bhatia, S., Boden, N., Borchers, A. and Boyle, R., 2017, June. In-datascenter performance analysis of a tensor processing unit. ISCA'2017

# In Summary: CNNs are associated with...

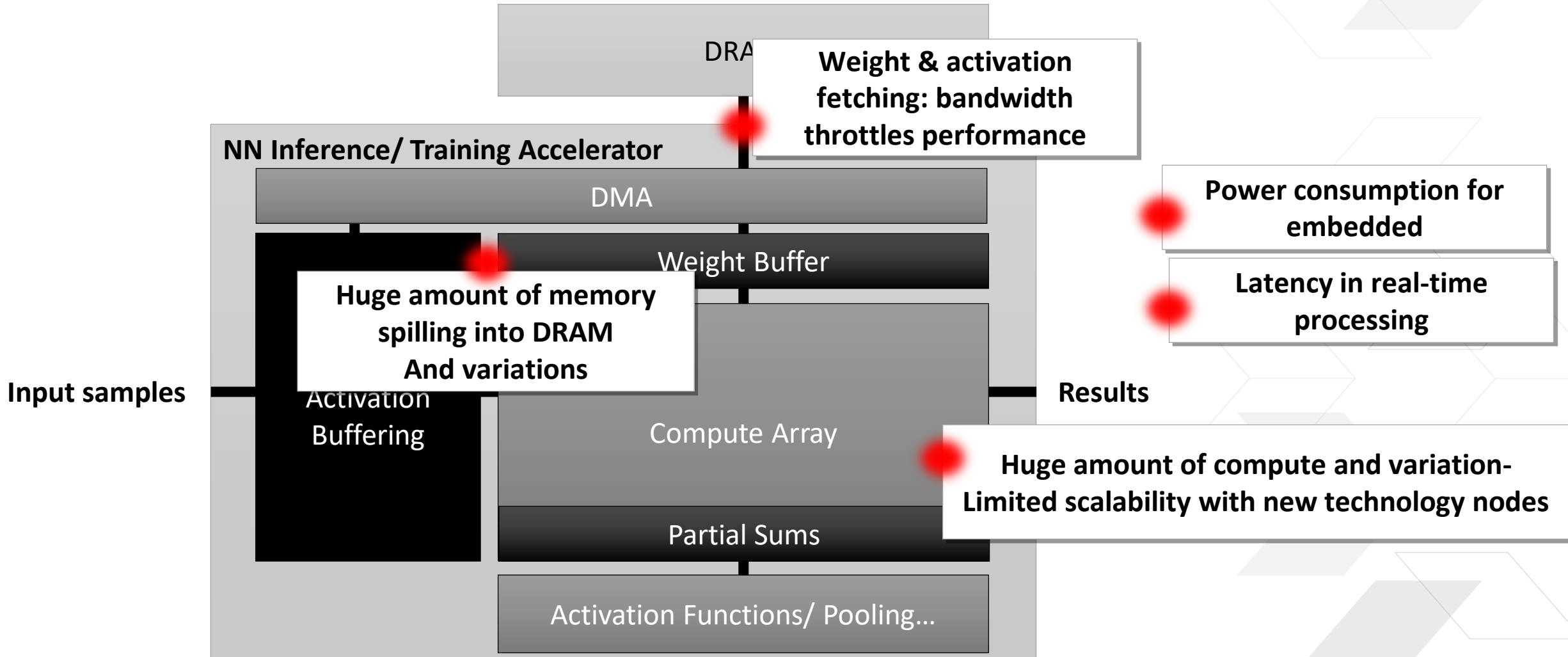
- > **Significant amounts of memory and computation**
- > **Huge variation between topologies and within them**
- > **Fast changing algorithms**
- > **Special functions, non-linear topologies**
- > **However, incredibly parallel!**
  - >> For convolutions: filter dimensions, feature map dimensions, input & output channels, batches, layers, and even precisions (discussed later)



Adopted from Ce Zhang, ETH Zurich, Systems Group Retreat



# Architectural Challenges/ Pain Points



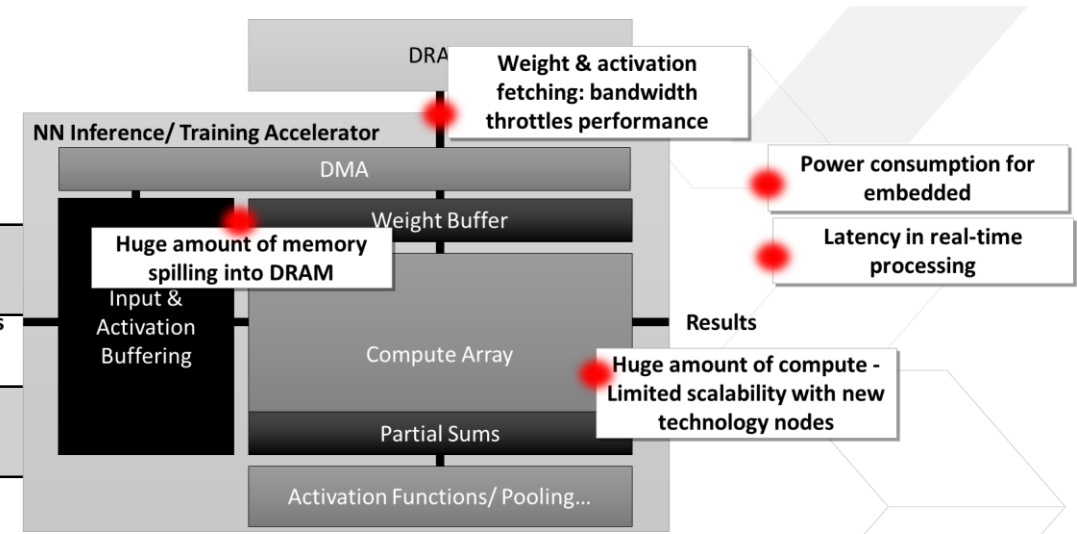
Requires algorithmic & architectural innovation

# Algorithmic Optimization Techniques



# Optimization Techniques

- Loop transformations to minimize memory access\*
- Pruning
- Compression
- Winograd, Strassen and FFT
- Novel layer types (squeeze, shuffle, shift)
- Numerical Representations & Reducing Precision



\*Chen, Y.H., Krishna, T., Emer, J.S. and Sze, V., 2017. Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE Journal of Solid-State Circuits*, 52(1), pp.127-13

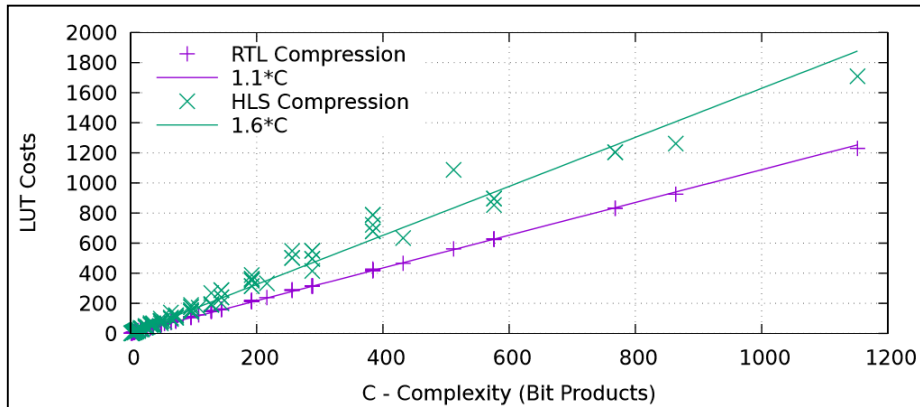
# Example: Reducing Bit-Precision

## > Linear reduction in memory footprint

- >> Reduces weight fetching memory bandwidth
- >> NN model may even stay on-chip

## > Reducing precision shrinks inherent arithmetic cost in both ASICs and FPGAs

- >> Instantiate **100x** more compute within the same fabric and thereby scale performance

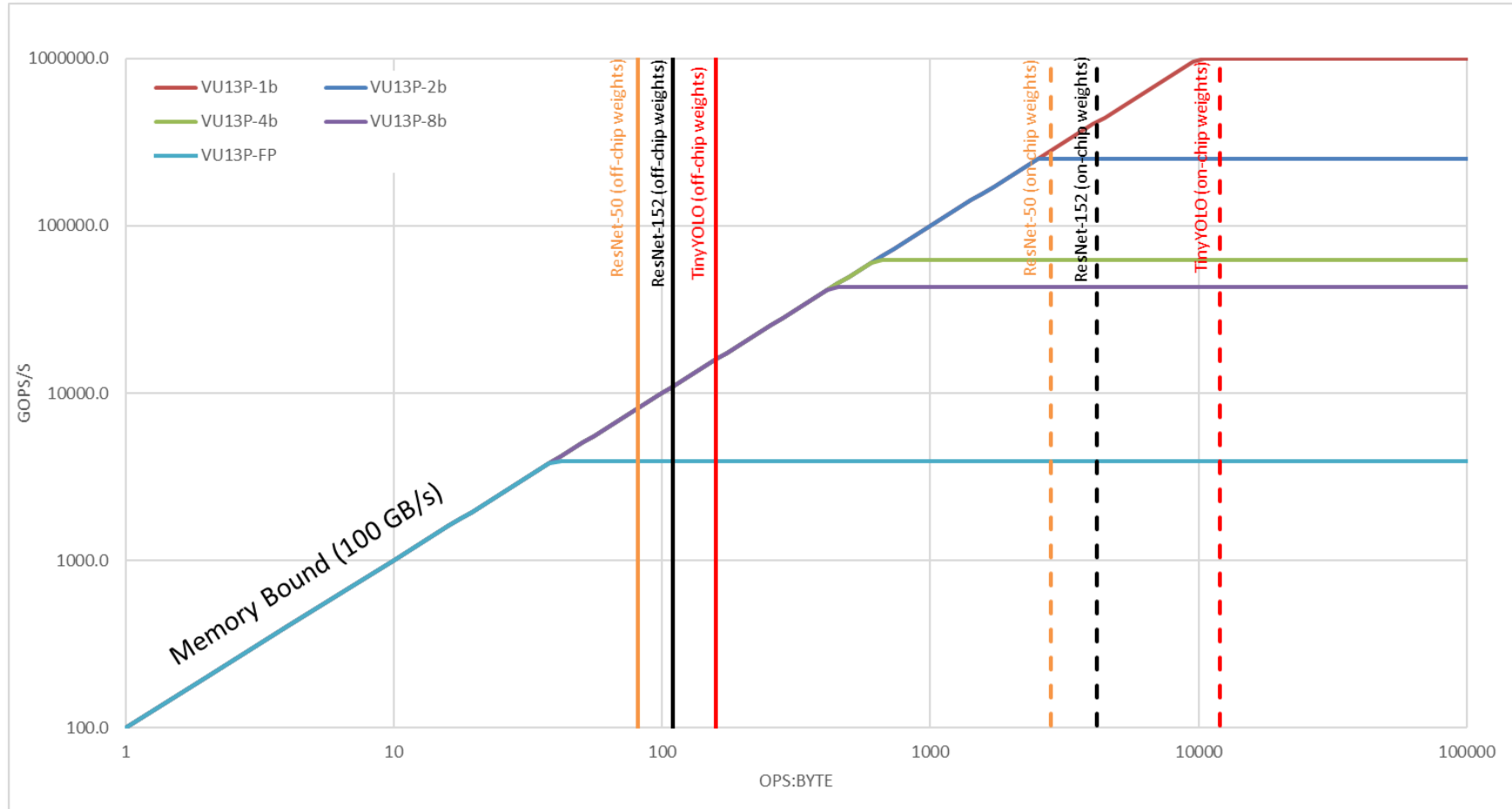


$C = \text{size of accumulator} * \text{size of weight} * \text{size of activation}$   
(to appear in ACM TRETSE on DL, FINN-R)

Precision	Modelsize [MB] (ResNet50)
1b	3.2
8b	25.5
32b	102.5

# Reducing Precision provides Performance Scalability

*Example: ResNet50, ResNet152 and TinyYolo*



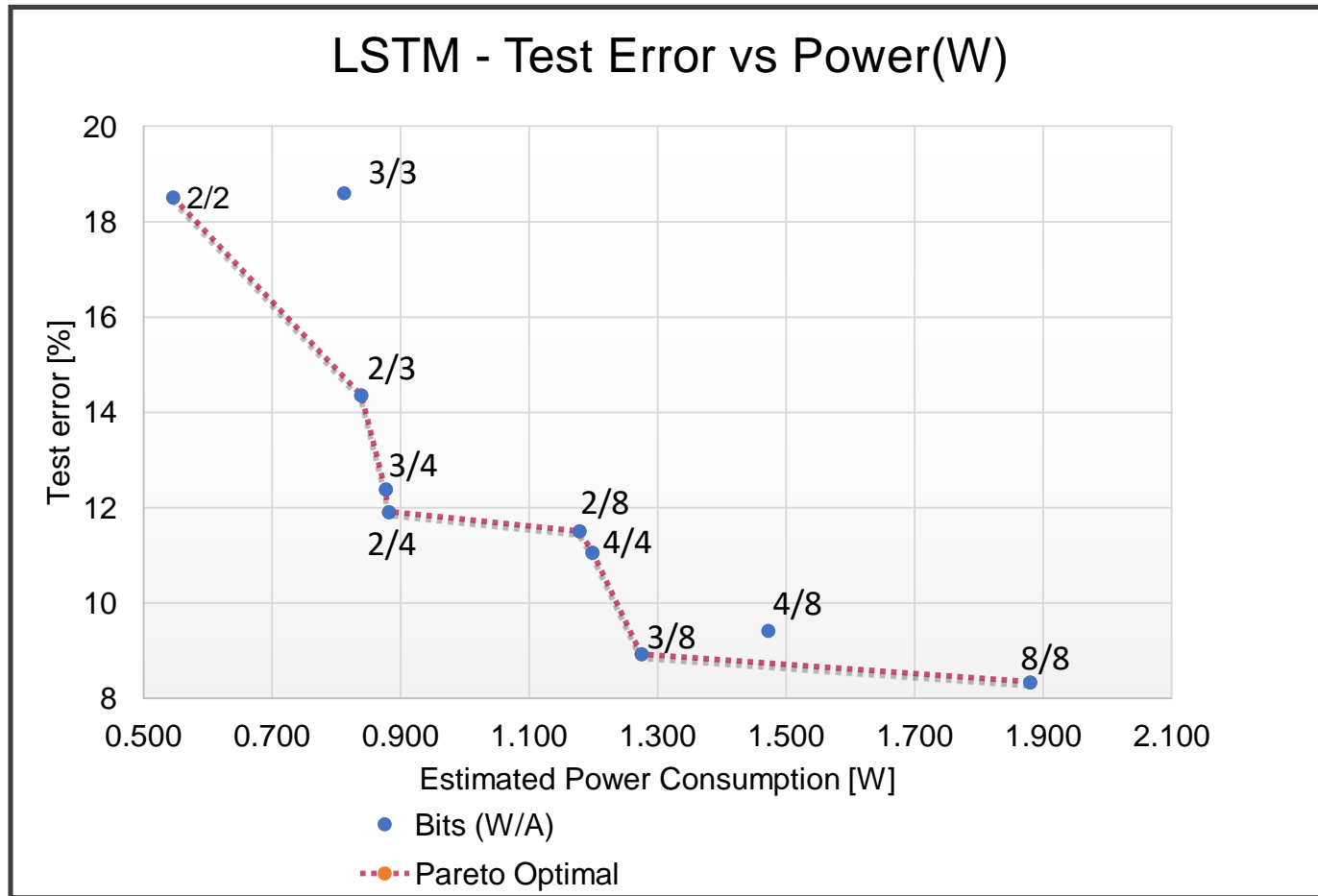
RP scales compute performance

*Theoretical Peak Performance for a VU13P with different Precision Operations*  
*Assumptions: Application can fill device to 90% (fully parallelizable) 710MHz*

RP reduces model size=> to stay on-chip

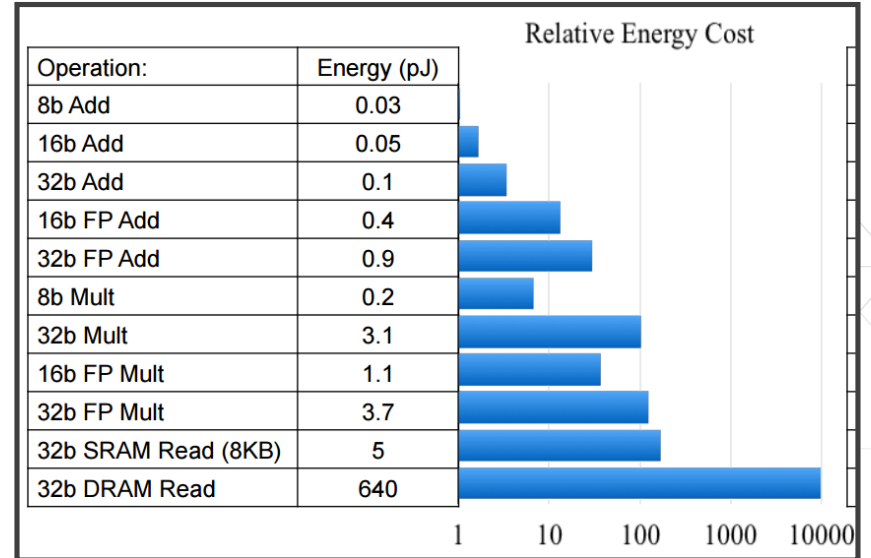
# Reducing Precision Inherently Saves Power

## FPGA:



Target Device ZU7EV • Ambient temperature: 25 °C • 12.5% of toggle rate • 0.5 of Static Probability • Power reported for PL accelerated block only

## ASIC:

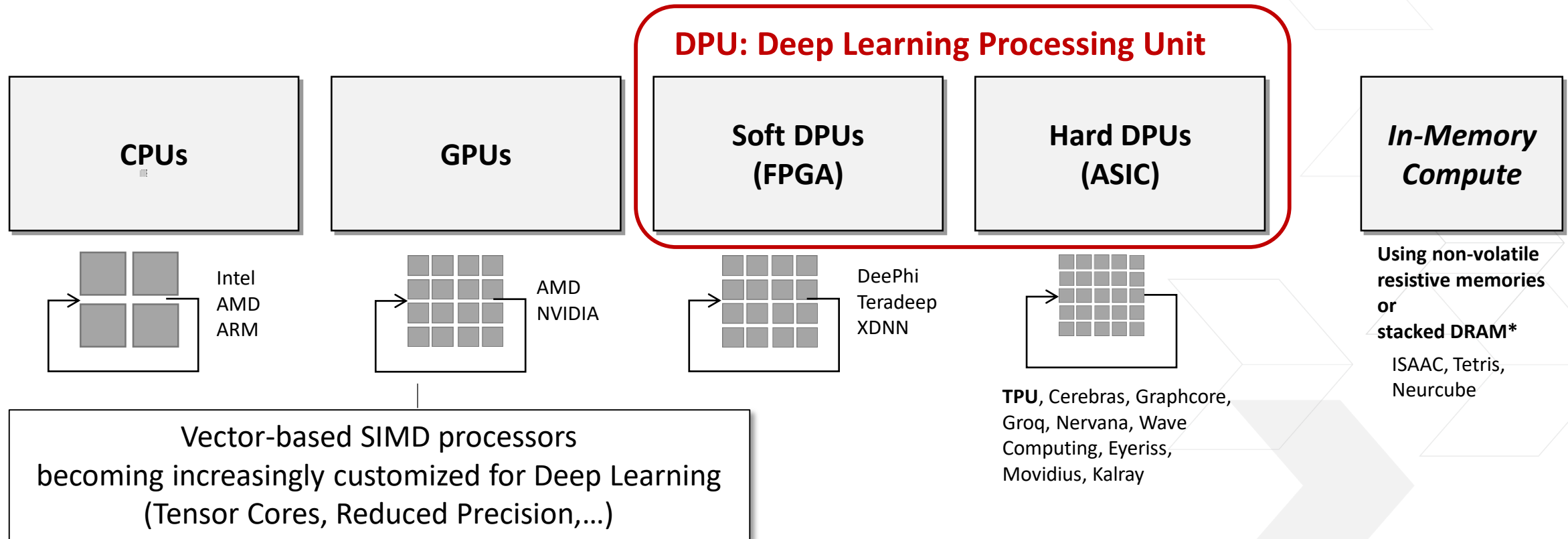


Source: Bill Dally (Stanford), Cadence Embedded Neural Network Summit, February 1, 2017

# Hardware Architectures and their Specialization Towards CNN Workloads

*Exciting Times in Computer  
Architecture Research!*

# Spectrum of New Architectures for Deep Learning



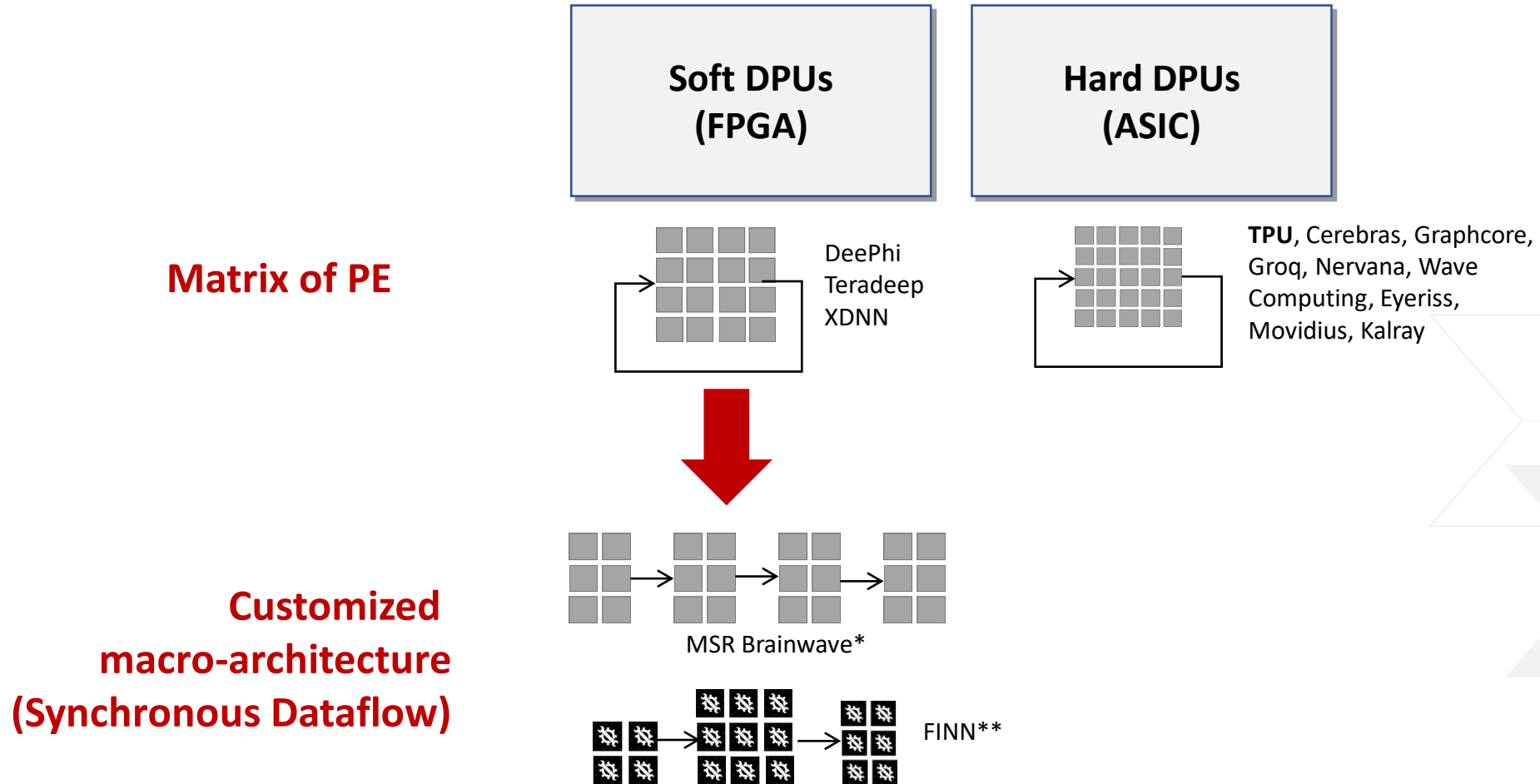
\*Shafiee, A., Nag, A., Muralimanohar, N., Balasubramonian, R., Strachan, J.P., Hu, M., Williams, R.S. and Srikumar, V., 2016. ISAAC: A convolutional neural network accelerator with in-situ analog arithmetic in crossbars. ACM SIGARCH

Chi, P., Li, S., Xu, C., Zhang, T., Zhao, J., Liu, Y., Wang, Y. and Xie, Y., 2016, June. Prime: A novel processing-in-memory architecture for neural network computation in rram-based main memory. In ACM SIGARCH

Chen, Y., Luo, T., Liu, S., Zhang, S., He, L., Wang, J., Li, L., Chen, T., Xu, Z., Sun, N. and Temam, O., 2014, December. Dadiannao: A machine-learning supercomputer. In Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture (pp. 609-622). IEEE Computer Society.



# Architectural Choices – Macro-Architecture

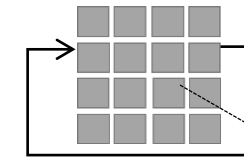
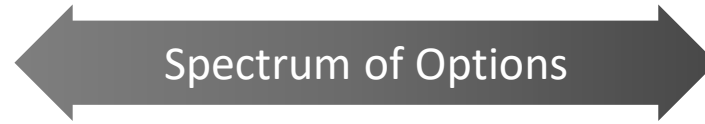
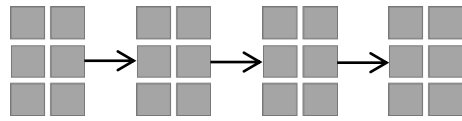


\*Chung, E., Fowers, J., Ovtcharov, K., Papamichael, M., Caulfield, A., Massengill, T., Liu, M., Lo, D., Alkalay, S., Haselman, M. and Abeydeera, M. Serving DNNs in Real Time at Datacenter Scale with Project Brainwave. *IEEE Micro*, 38(2)

<https://www.microsoft.com/en-us/research/uploads/prod/2018/06/ISCA18-Brainwave-CameraReady.pdf>

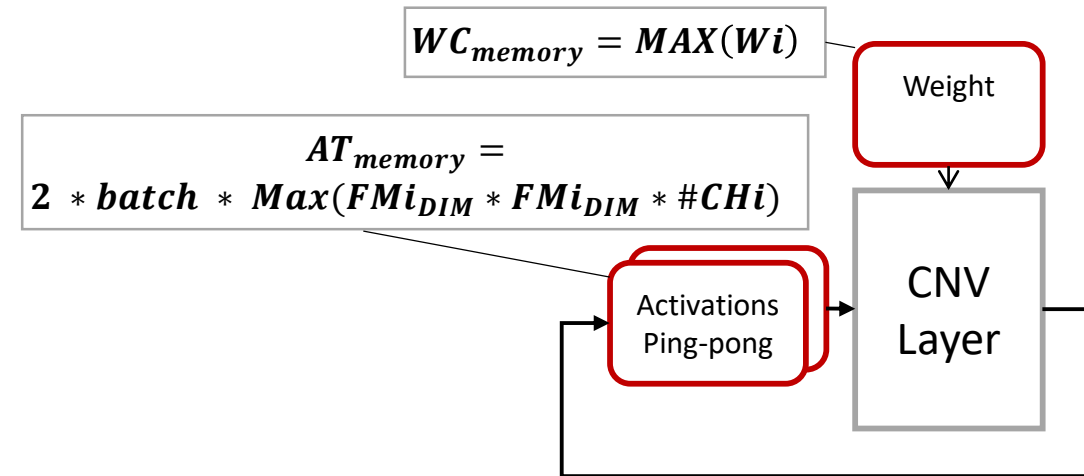
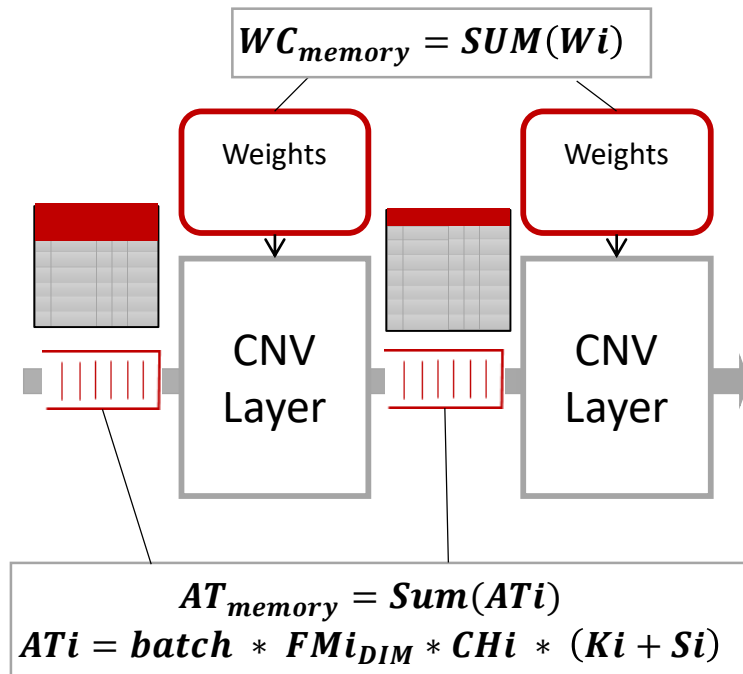
\*\*Umuroglu, Yaman, Umuroglu, Y., Fraser, N.J., Gambardella, G., Blott, M., Leong, P., Jahre, M. and Vissers, K. "FINN: A framework for fast, scalable binarized neural network inference." *ISFPGA'2017*

# Synchronous Dataflow (SDF) vs Matrix of Processing Elements (MPE)

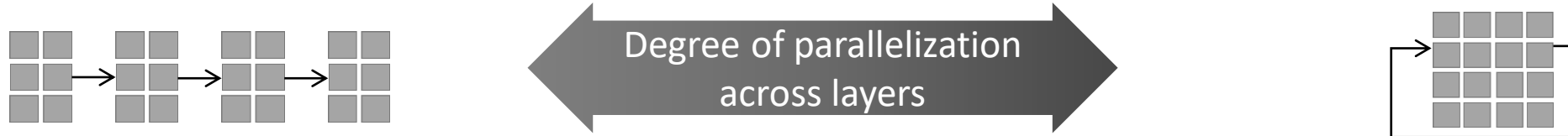


MAC, Vector Processor

>> End points are pure layer-by-layer compute and feed-forward dataflow architecture



# Synchronous Dataflow (SDF) vs Matrix of Processing Elements (MPE)



- Requires less activation buffering
- Higher compute and memory efficiency due to custom-tailored hardware design
- Less flexibility
- Less latency (reduced buffering)
- No control flow (static schedule)

- Requires less on-chip weight memory, but more activation buffers
- Efficiency of memory for weights and activations depends on how well balanced the topology is
- Flexible hardware, which can scale to arbitrary large networks
- Compute efficiency is a scheduling problem  
=> generating sophisticated scheduling algorithms

# Summary



# Summary

- > **Deep NNs are increasingly being adopted for new workloads and key to the current industrial revolution and perhaps the next**
- > **Associated with significant challenges to increase performance**
- > **Requires algorithmic and architectural innovation (co-designed)**
- > **Emerging: Huge spectrum of algorithms and increasingly diverse & heterogenous hardware architectures**
- > **Clear metrics for comparison needed**
  - >> Hardware performance always tying back to application performance (accuracy) to allow for algorithmic optimizations
  - >> Always trading off accuracy, performance (op/sec), latency (ms), power consumption (W)

# Exciting Times for our Community:

## Many New Architectures Evolving - Programmable and Hardened

