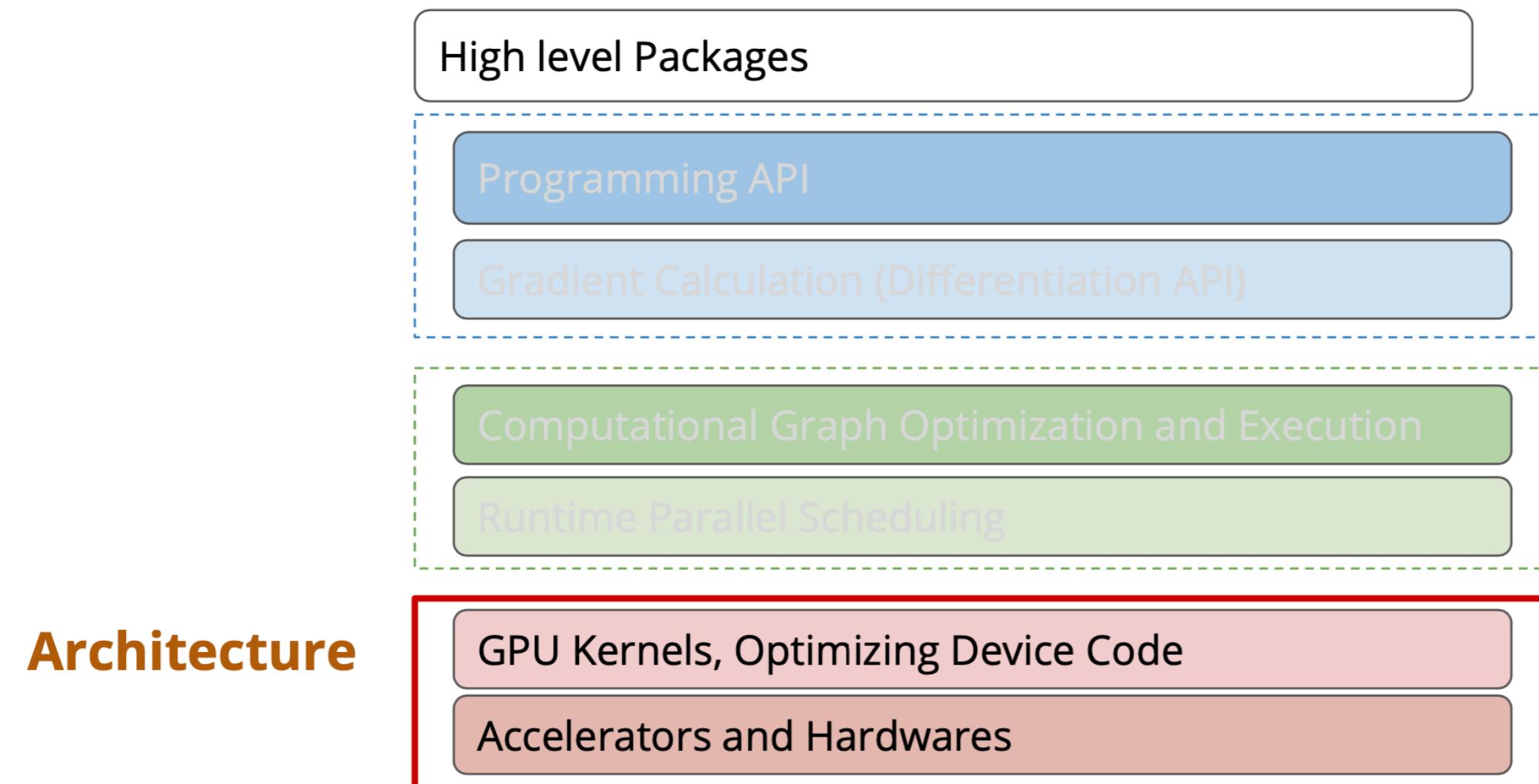


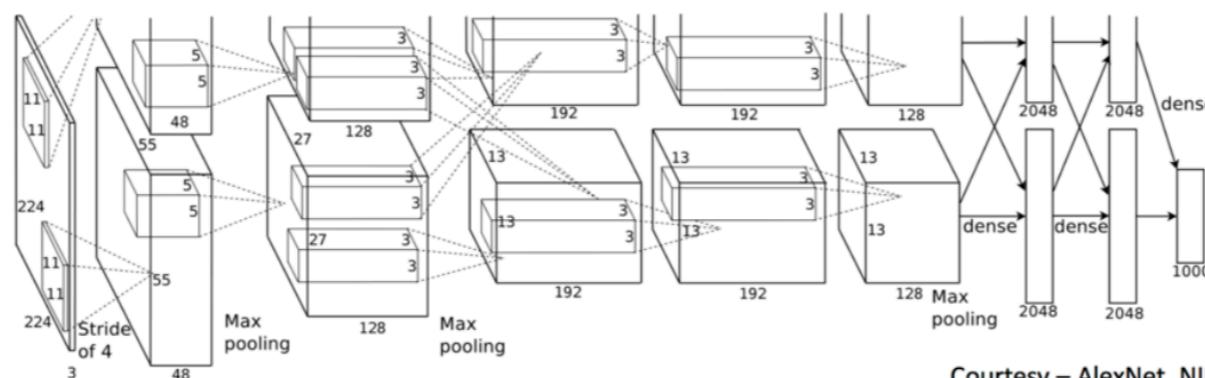
High-Performance Hardware for Deep Learning

Pooyan Jamshidi
UofSC

Typical Deep Learning System Stack



Hardware and Data enable DNNs



Courtesy – AlexNet, NIPS
2012

The need for Speed

- Larger data sets and models lead to better accuracy but also increase computation time. Therefore progress in deep neural networks is limited by how fast the networks can be computed.
- Likewise the application of convnets to low latency inference problems, such as pedestrian detection in self driving car video imagery, is limited by how fast a small set of images, possibly a single image, can be classified.

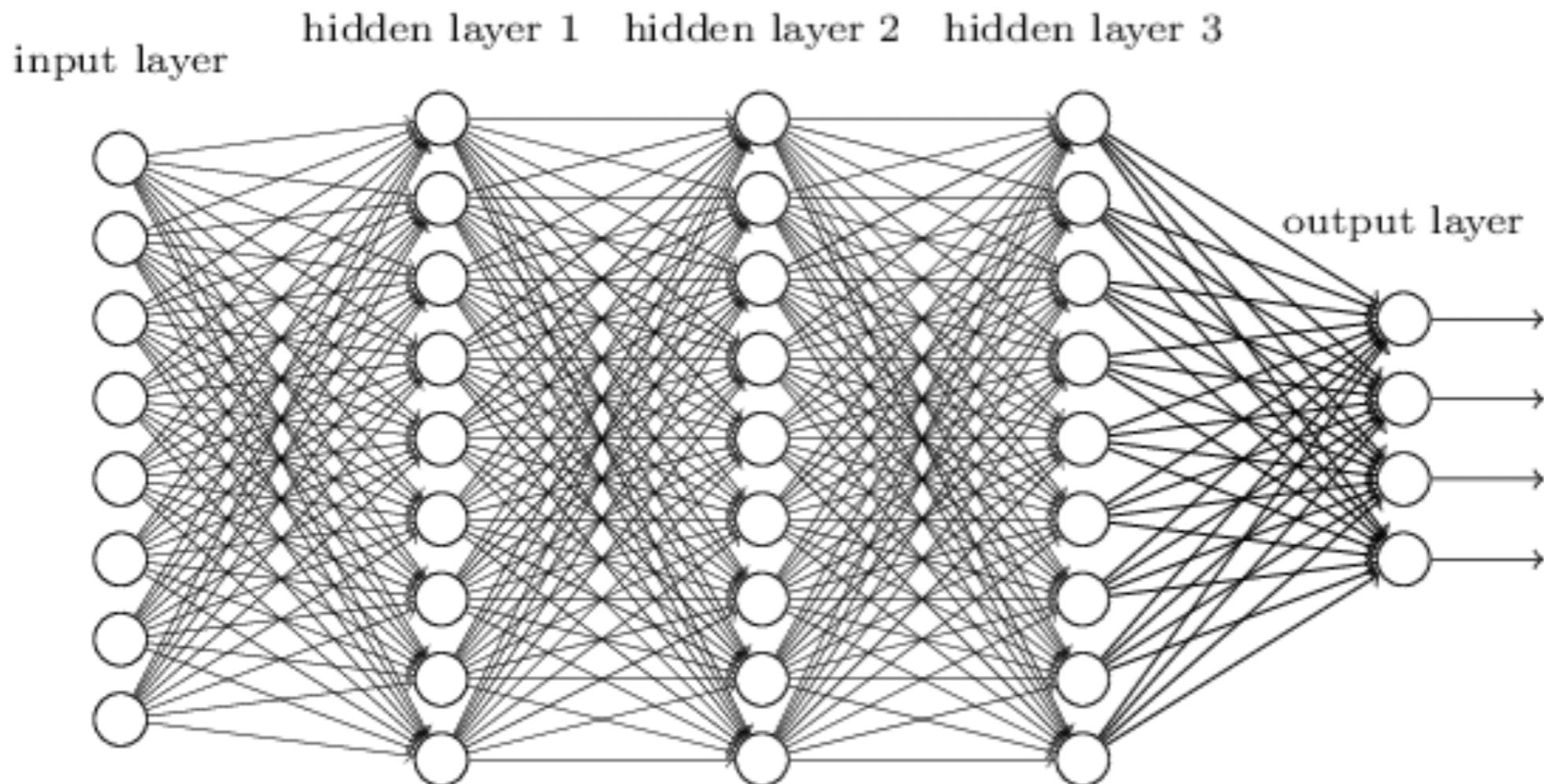
More data → Bigger Models → More Need for Compute
But Moore's law is no longer providing more compute...

Problem

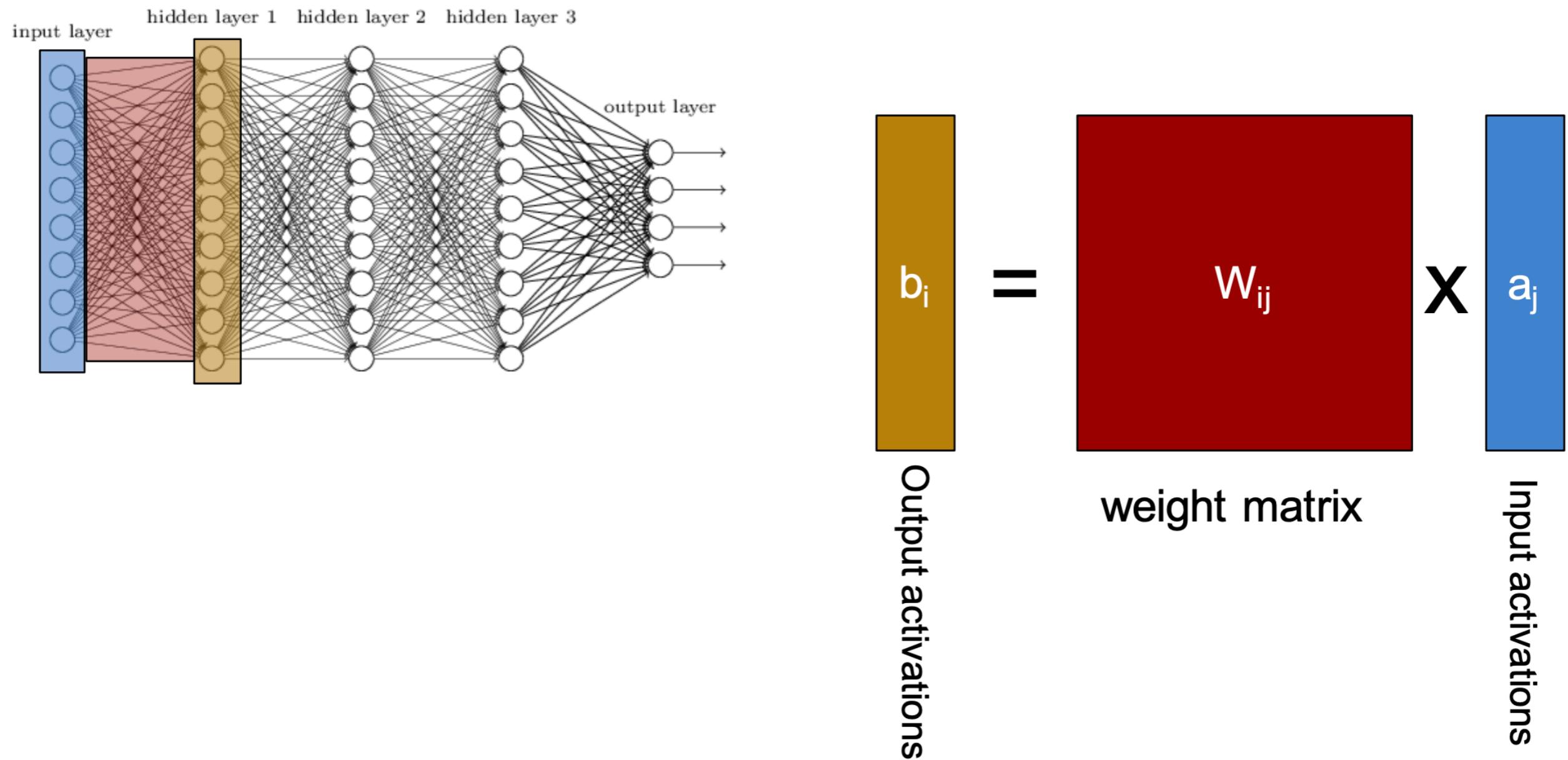
Acceleration

- Run a network faster (Performance, inf/s)
- Run a network more efficiently
 - Energy (inf/J)
 - Cost (inf/s\$)
- Inference
 - Just running the network forward
- Training
 - Running the network forward
 - Back-propagation of gradient
 - Update of parameters

What Network? DNNs, CNNs, and RNNs

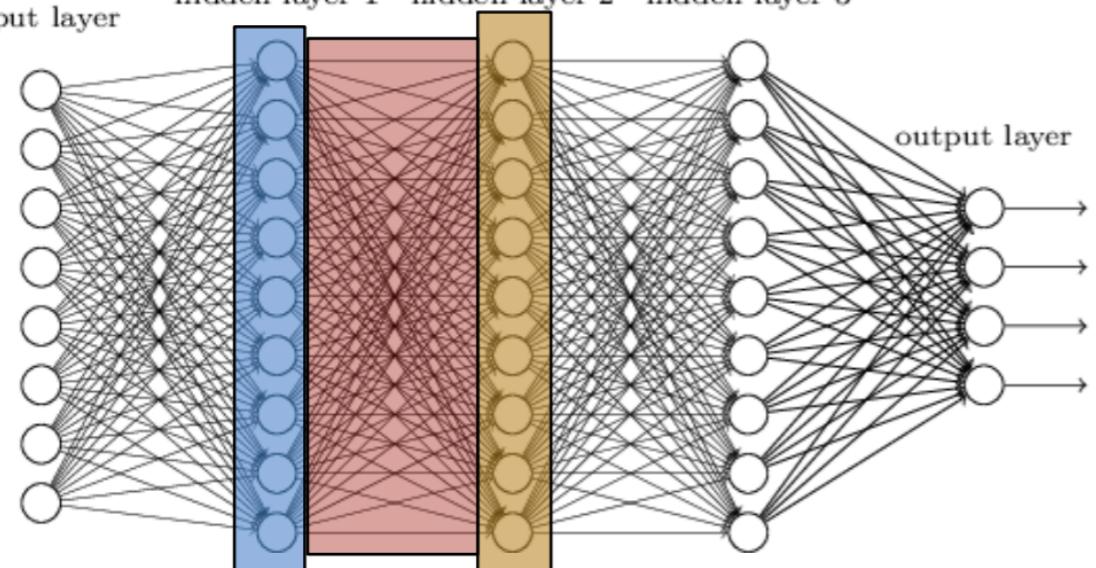


DNN, key operation is dense $M \times V$



DNN, key operation is dense $M \times V$

input layer hidden layer 1 hidden layer 2 hidden layer 3

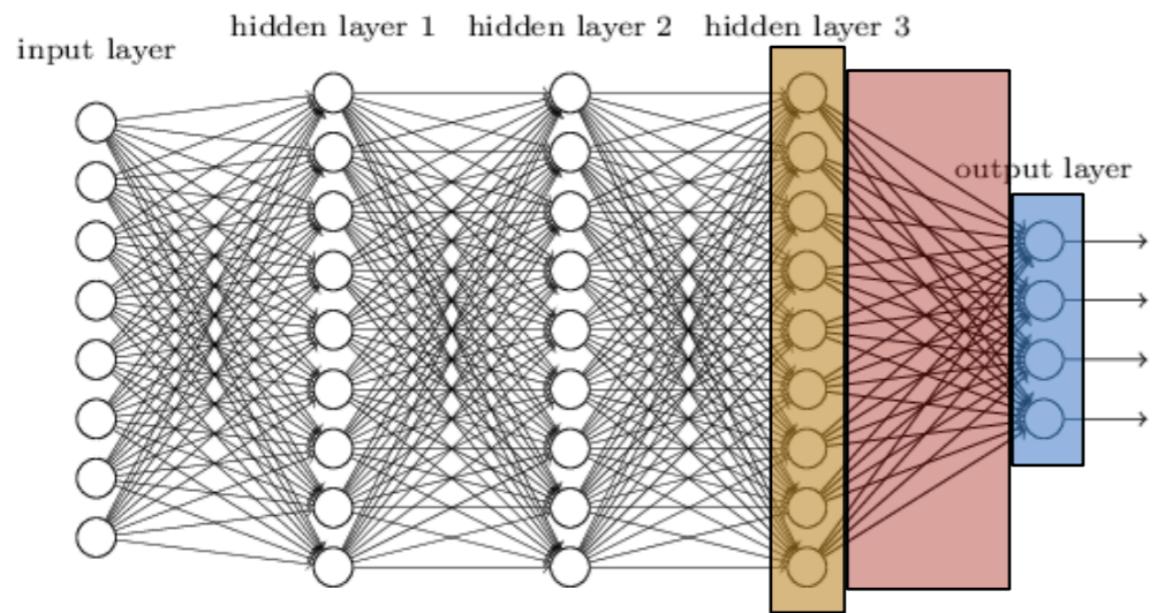


Repeat for each layer

$$b_i = W_{ij} \times a_j$$

Output activations weight matrix Input activations

DNN, key operation is dense $M \times V$

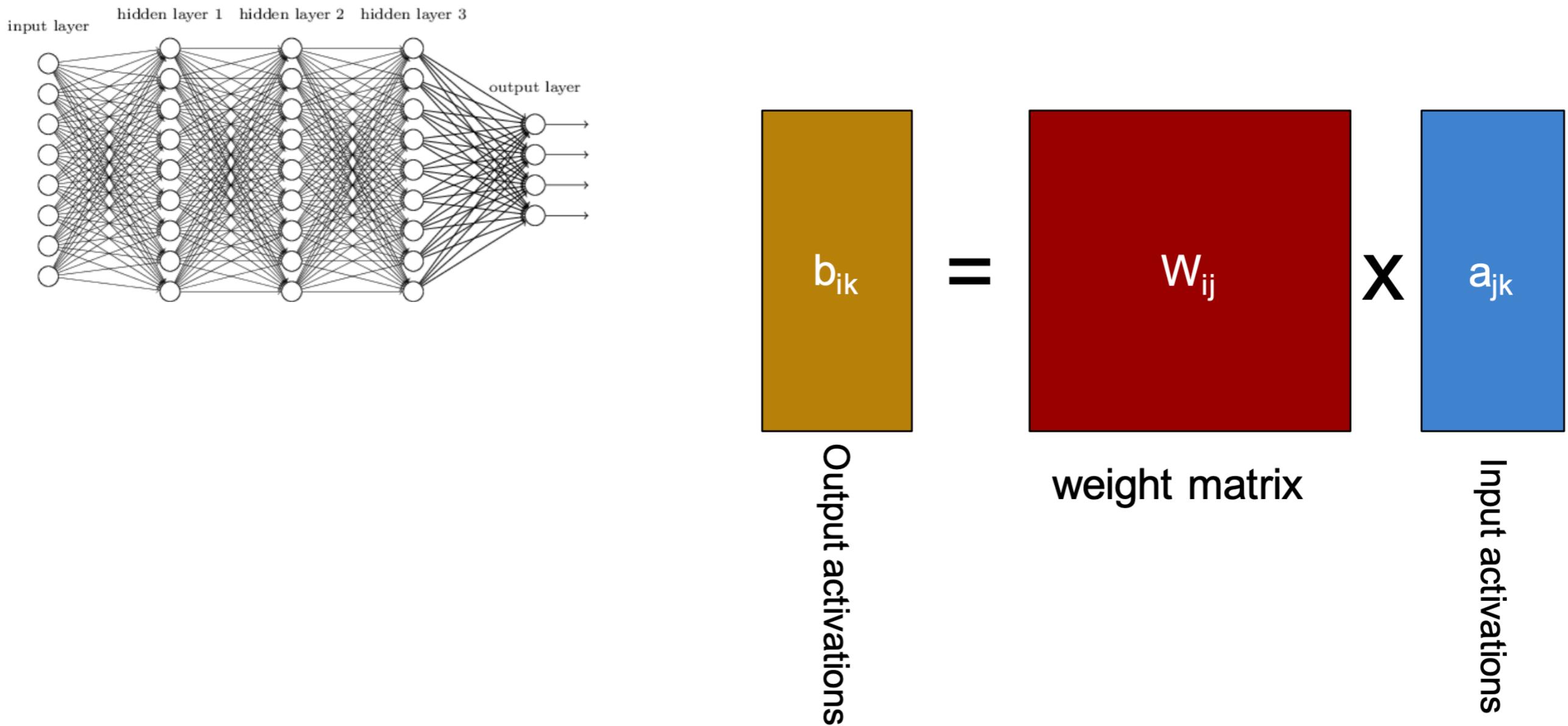


Backpropagation just does
this backward

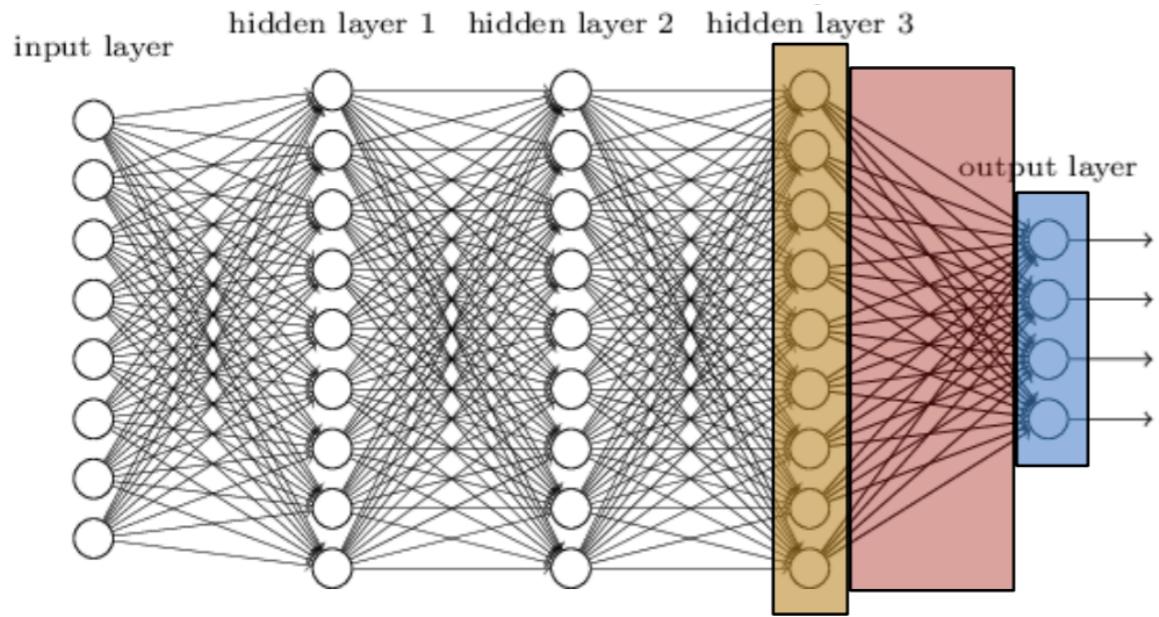
$$b_i = W_{ij} \times a_j$$

Input gradient weight matrix Output gradient

Training, and Latency Insensitive Networks can be Batched – operation is $M \times M$ – gives re-use of weights



For real time you can't batch
And there is sparsity in both weights and activations
key operations is $\text{spM} \times \text{spV}$

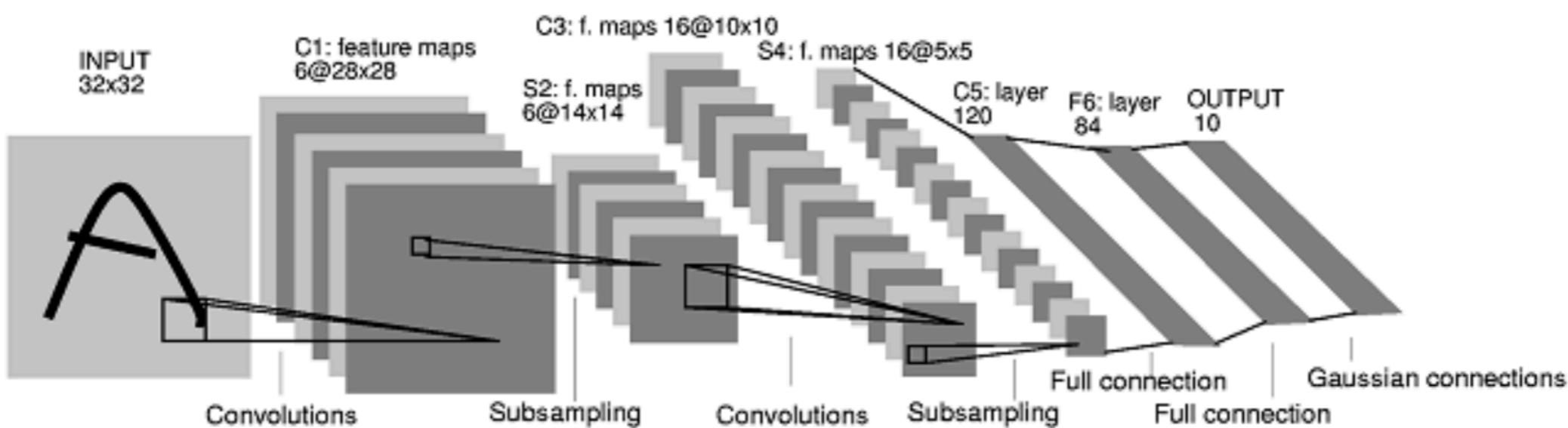


Backpropagation just does
this backward

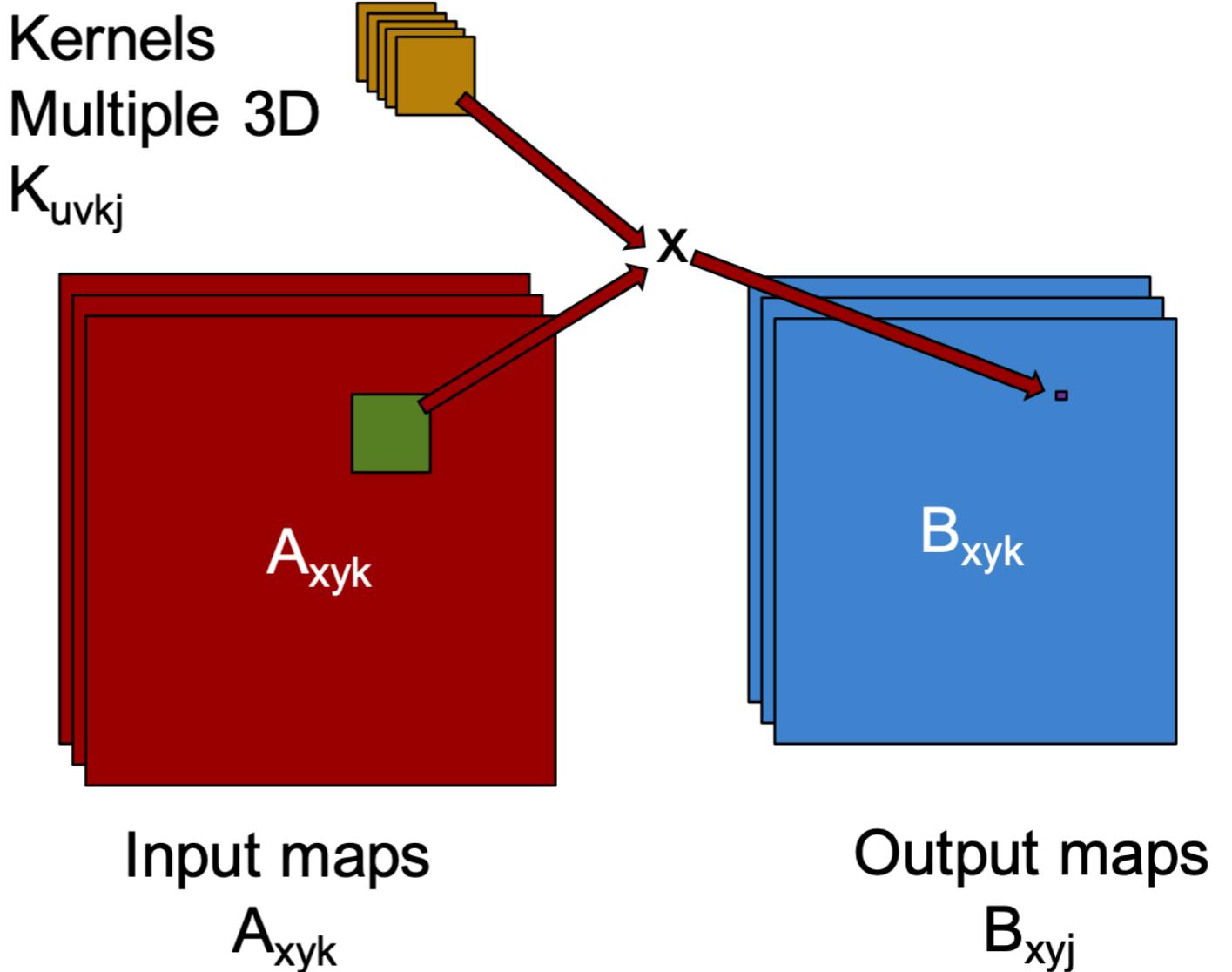
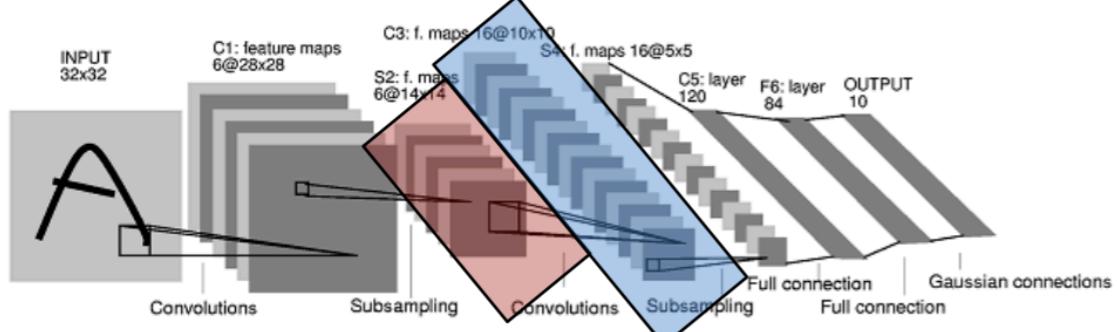
$$b_i = W_{ij} \times a_j$$

Input gradient weight matrix Output gradient

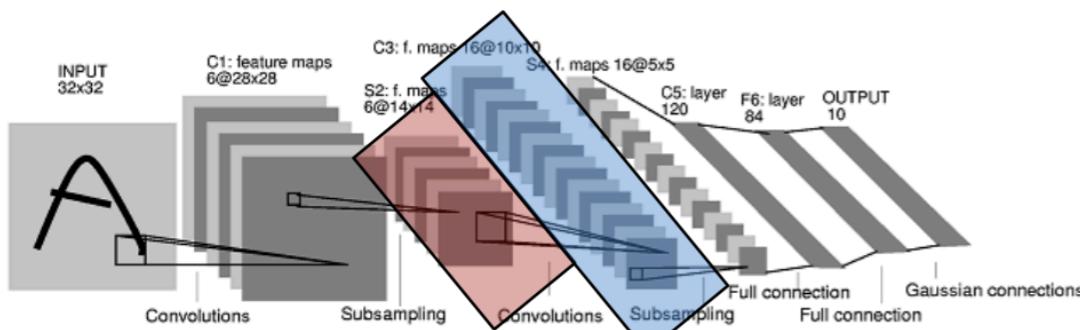
CNNs – For Image Inputs, Convolutional stages act as trained feature detectors



CNNs require Convolution in addition to $M \times V$

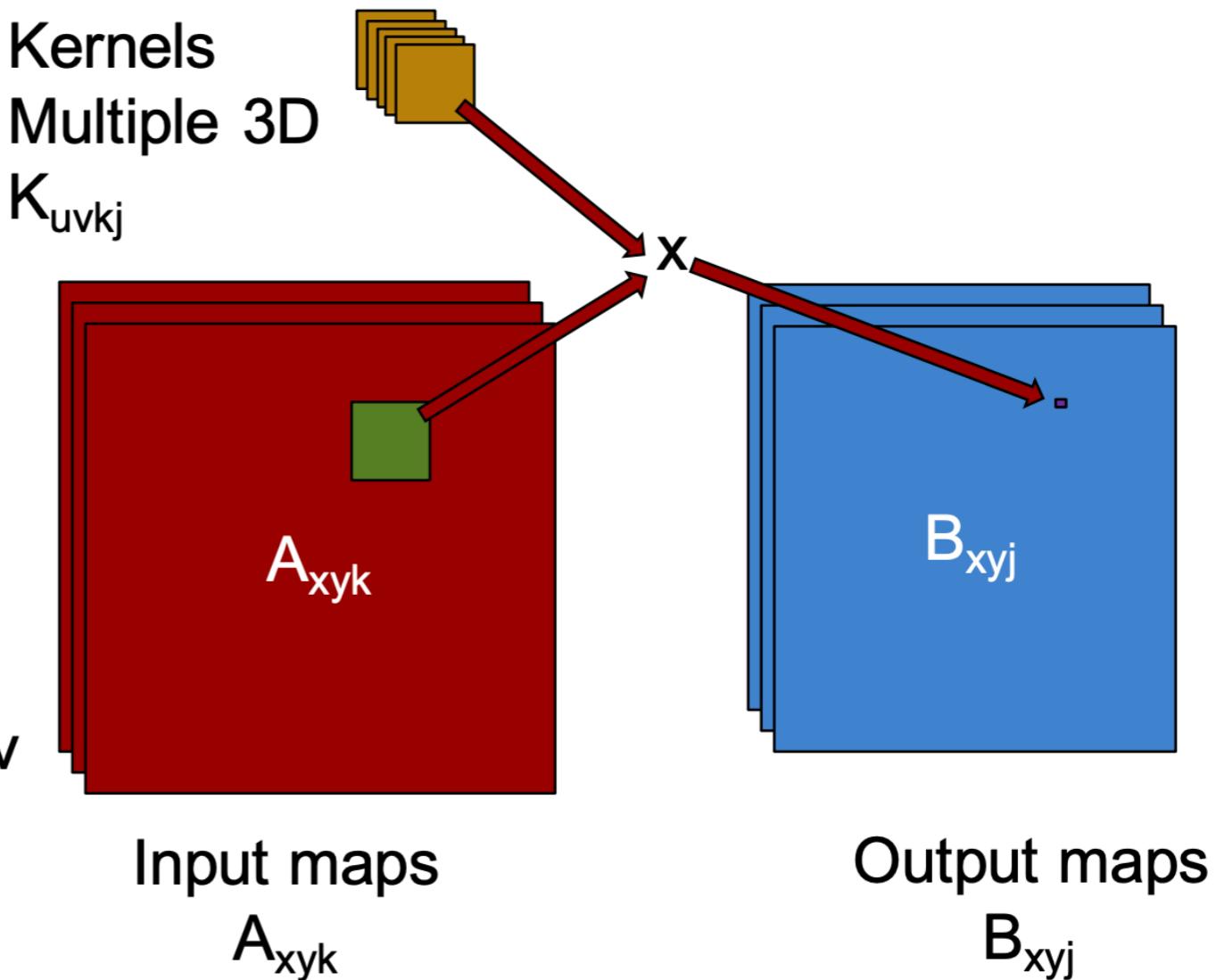


CNNs require Convolution in addition to $M \times V$

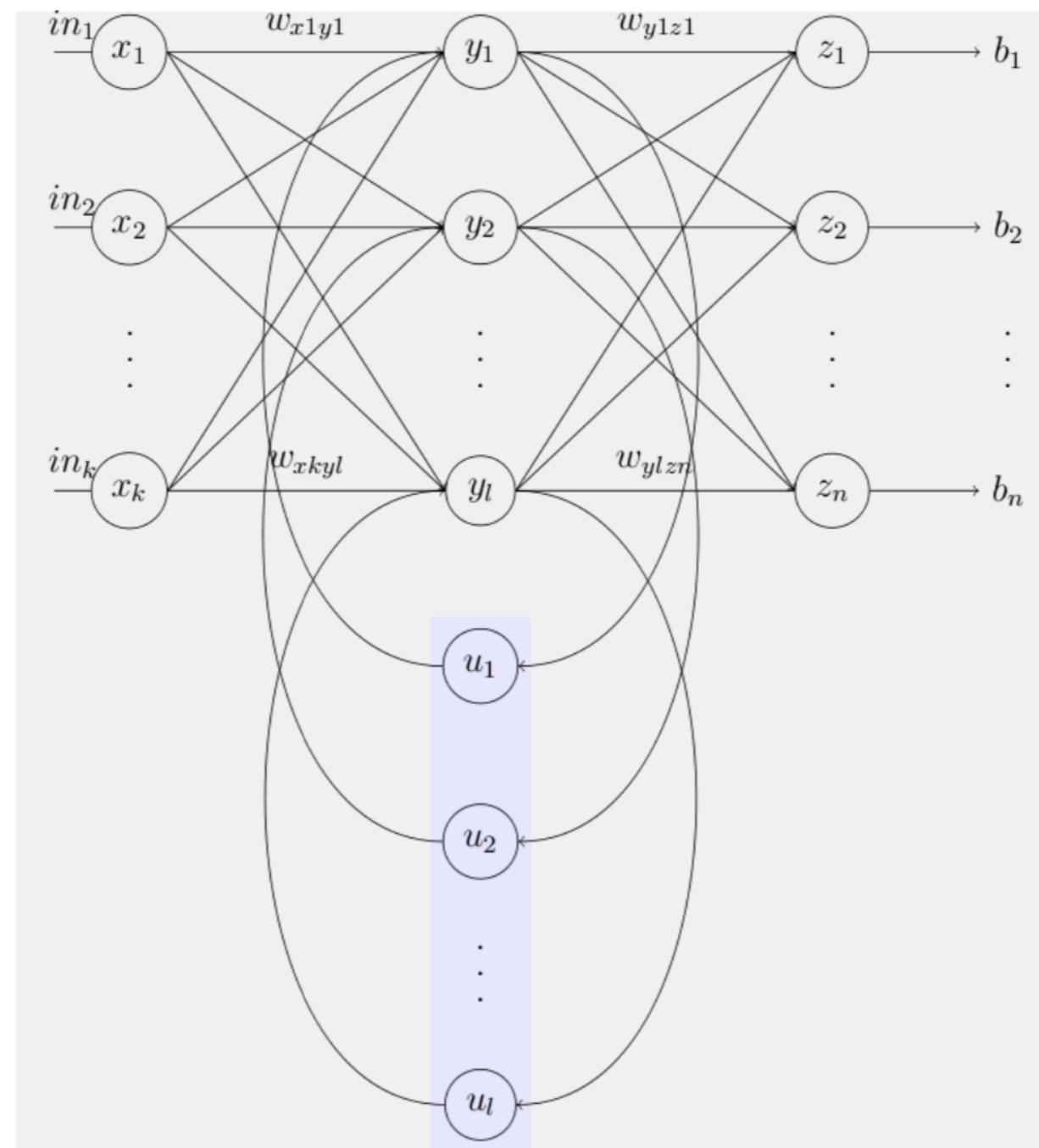


6D Loop
For each output map j
For each input map k
For each pixel x,y
For each kernel element u,v

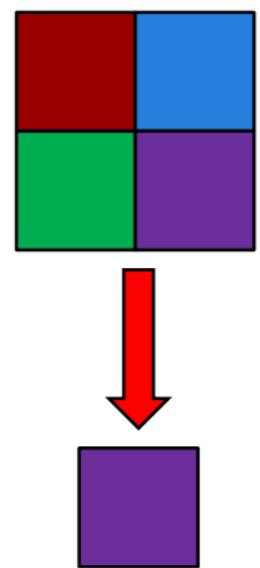
$$B_{xyj} += A_{(x-u)(y-v)k} \times K_{uvkj}$$



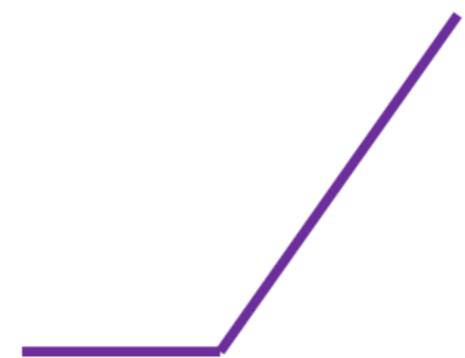
RNNs



Some Other Operations



Pooling

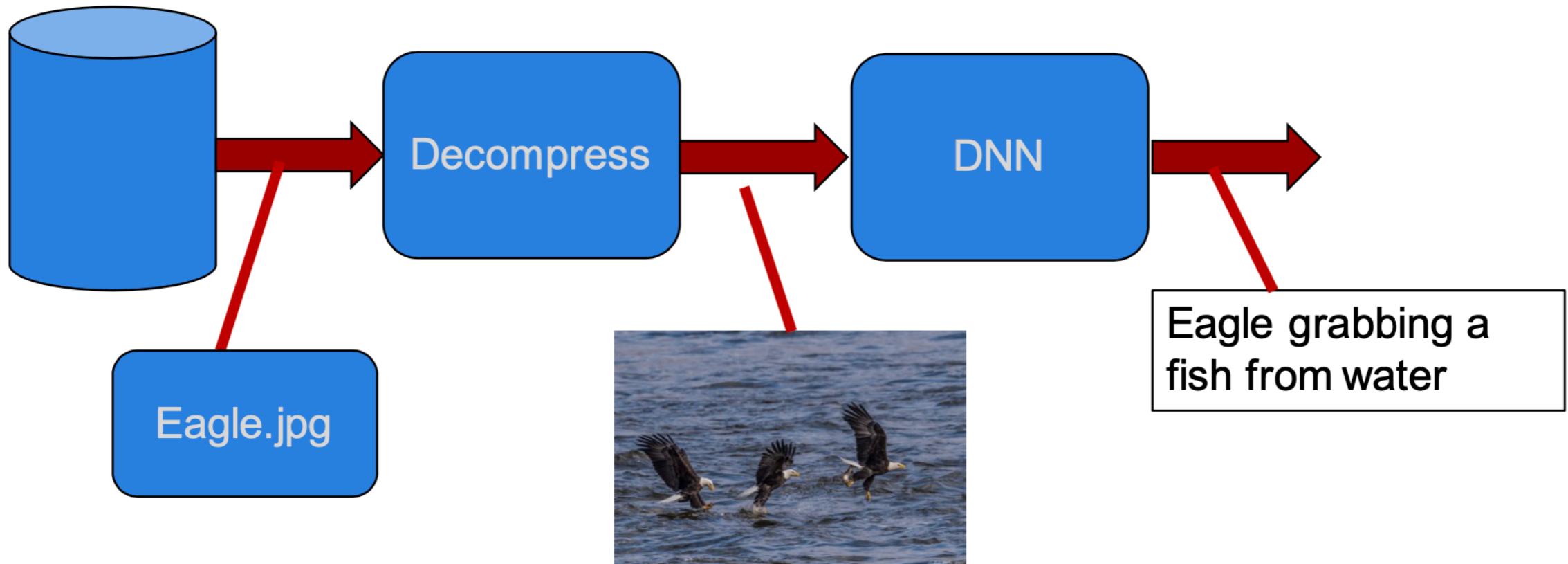


ReLU
(or other non-linear function)

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

Weight Update

Infrastructure



Summary of the Problem

- Run DNNs, CNNs, and RNNs
 - For training and inference
 - Can batch if not latency sensitive
- Optimize
 - Speed inf/s
 - Efficiency inf/J, inf/s\$
- Key operations are
 - MxV
 - MxM if batched
 - May be sparse (spM x spV)
 - Convolution
- Also
 - Pooling, non-linear operator (ReLU), weight update

Baseline

Baseline Performance Xeon E5-2698 – Single Core



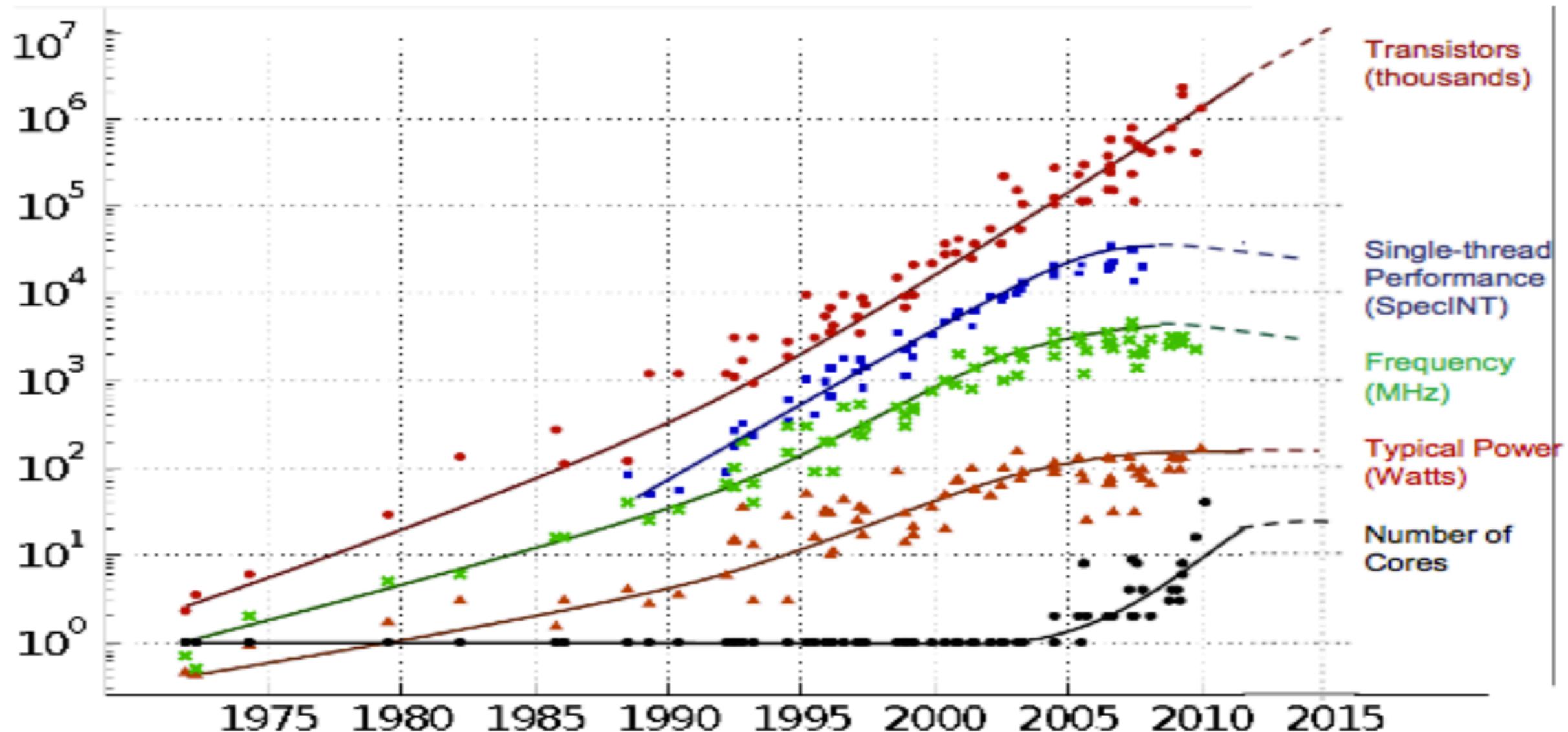
AlexNet – inference, batched

30 f/s

3.2 f/J

Most ops on AVX (SIMD) units

Moore's law made CPUs 300x faster than in 1990 But its over...

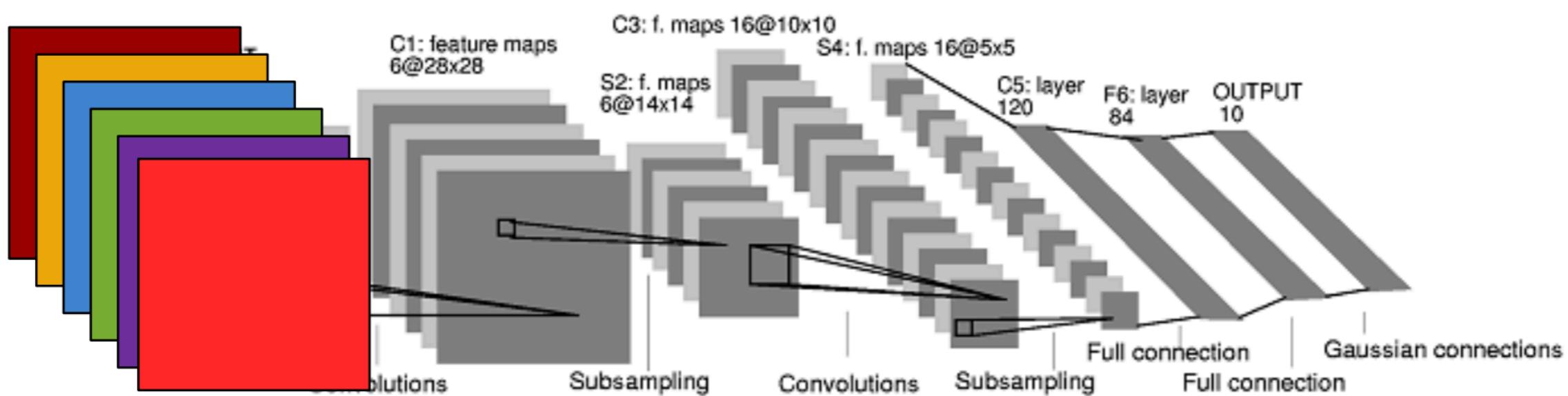


Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

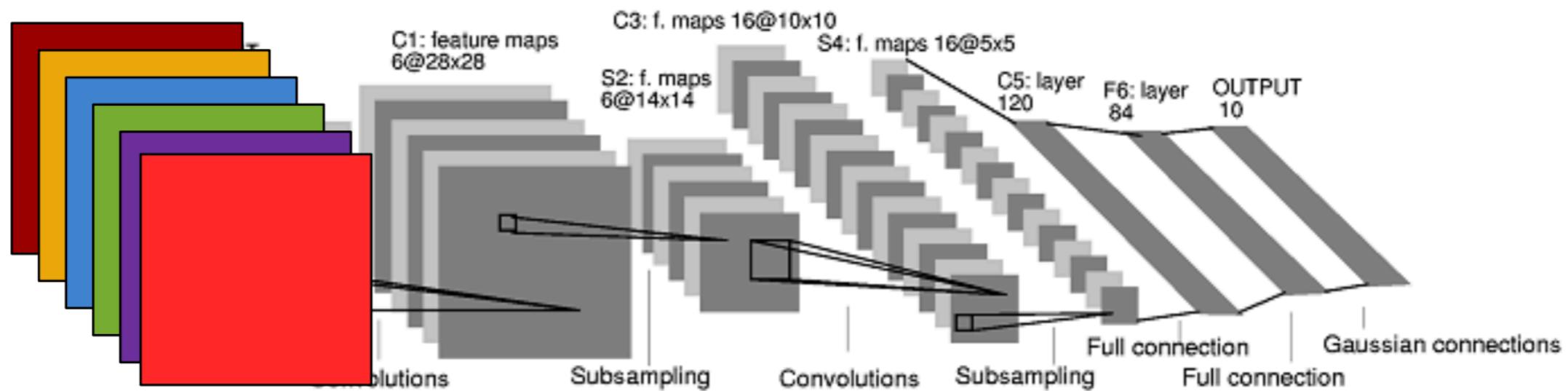
Parallelization

To go faster, use more
processors

Lots of parallelism in a DNN

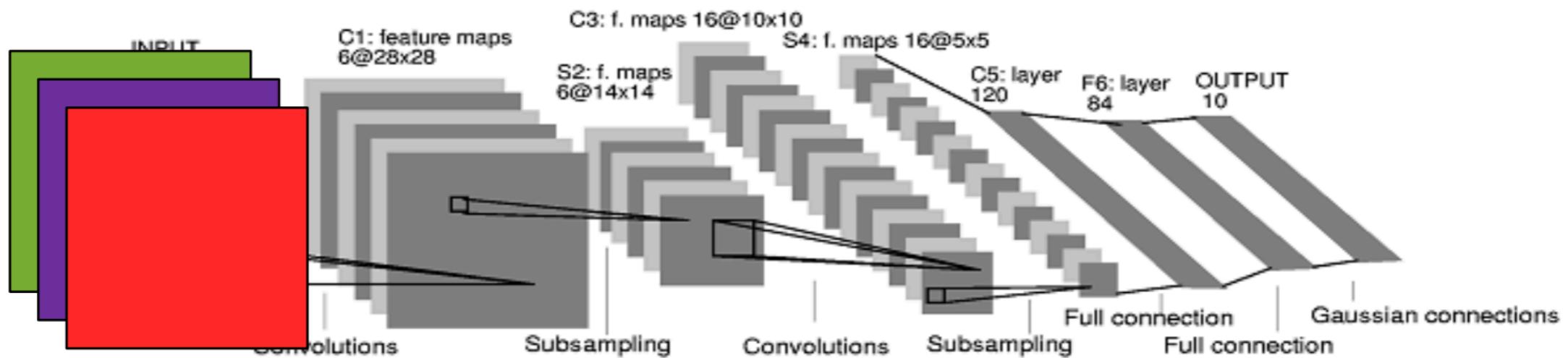
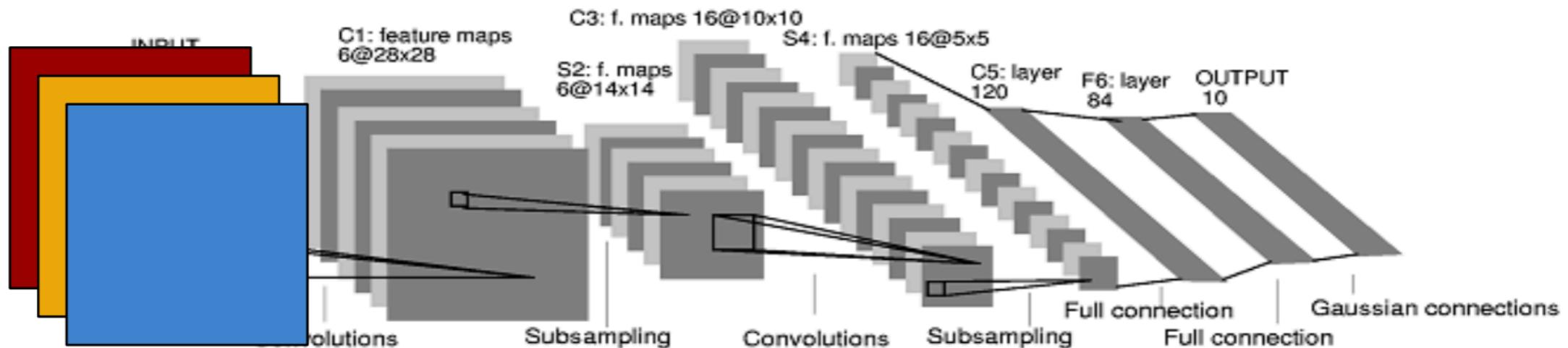


Lots of parallelism in a DNN

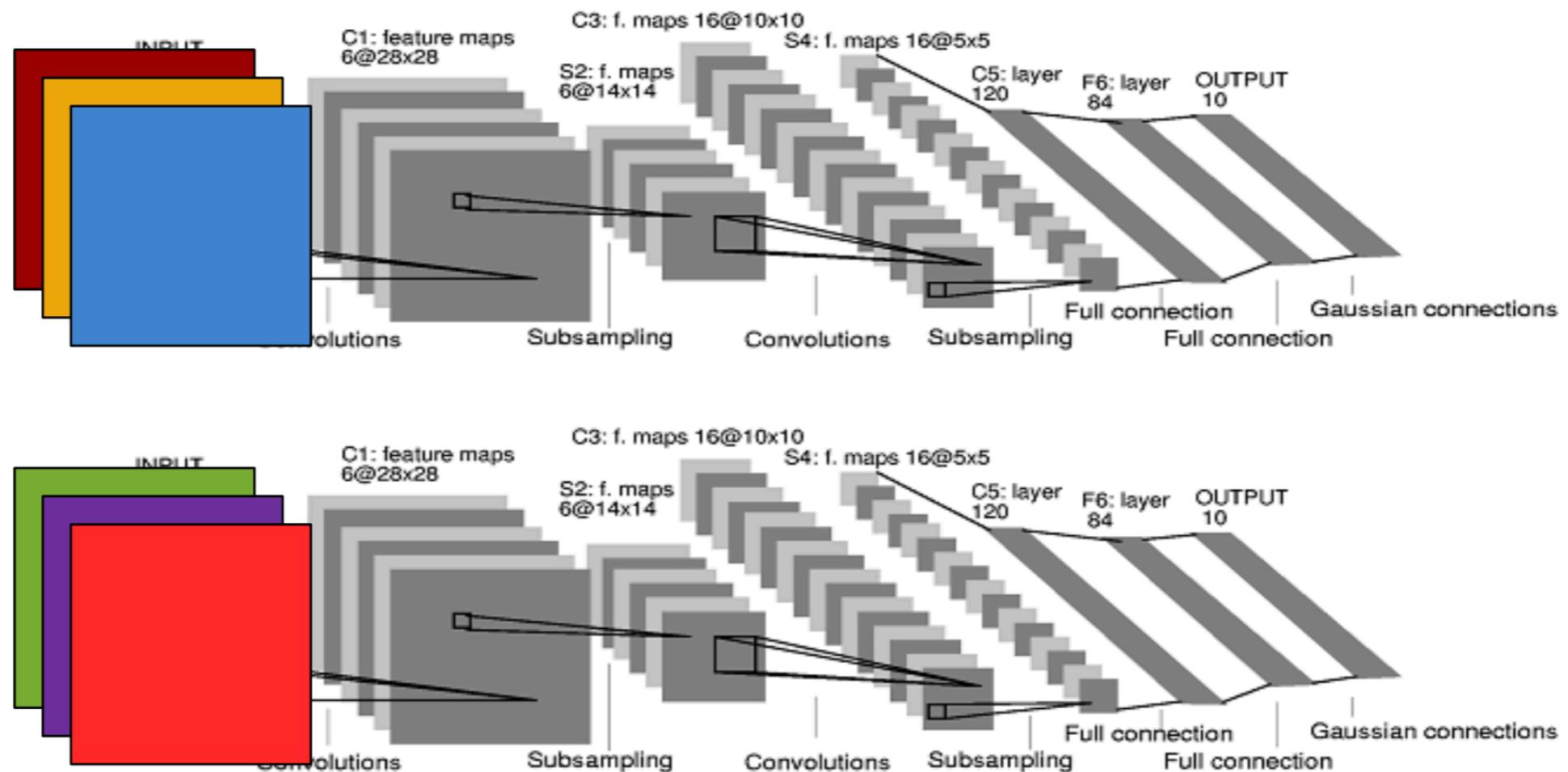


- Inputs
- Points of a feature map
- Filters
- Elements within a filter
- Multiplies within layer are independent
- Sums are reductions
- Only layers are dependent
- No data dependent operations
=> can be statically scheduled

Data Parallel – Run multiple inputs in parallel

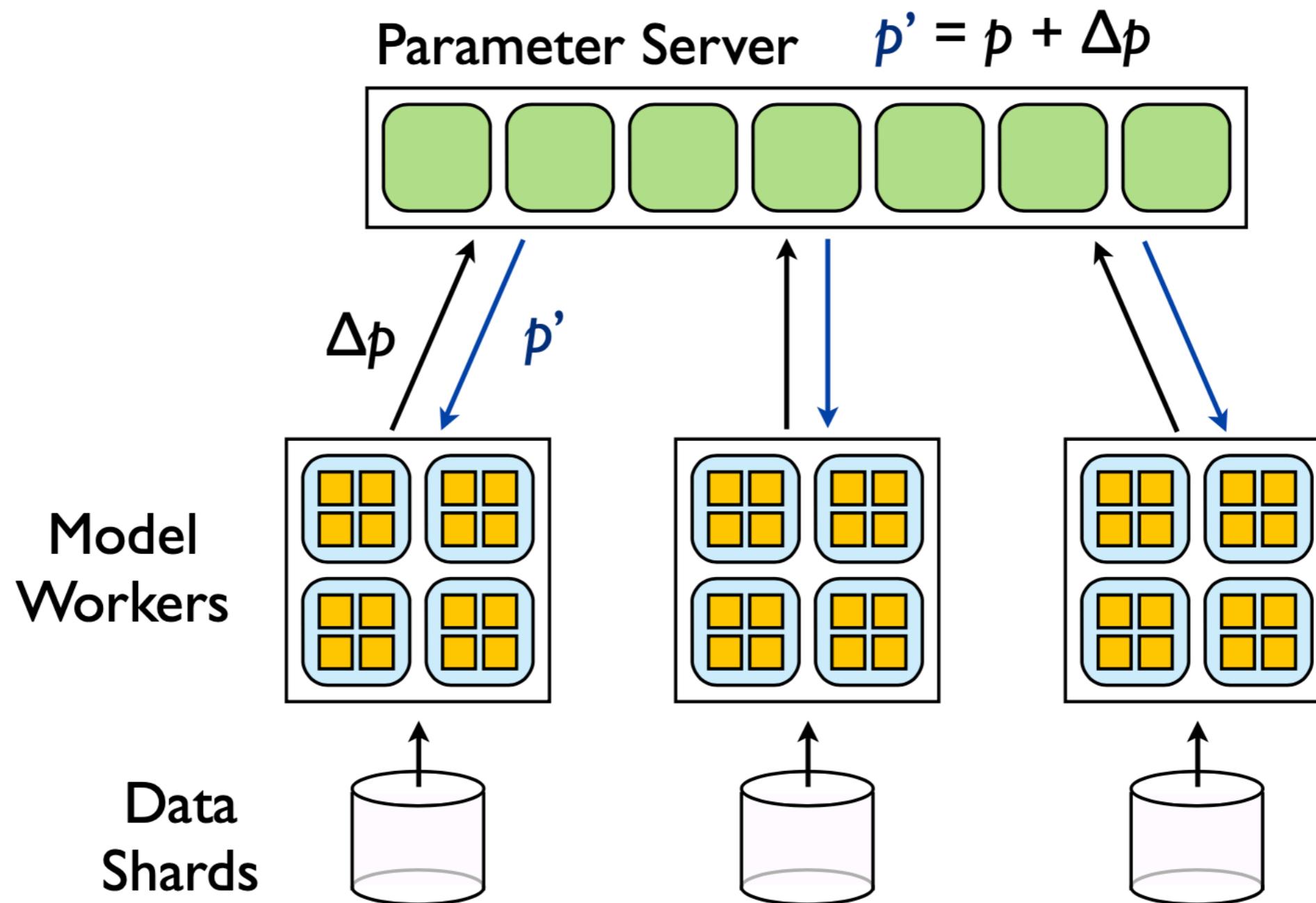


Data Parallel – Run multiple inputs in parallel



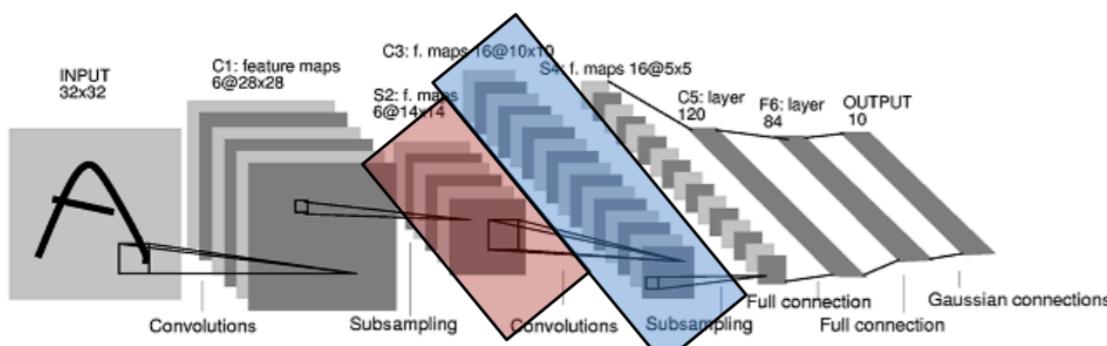
- Doesn't affect latency for one input
- Requires P-fold larger batch size
- For training requires coordinated weight update

Parameter Update



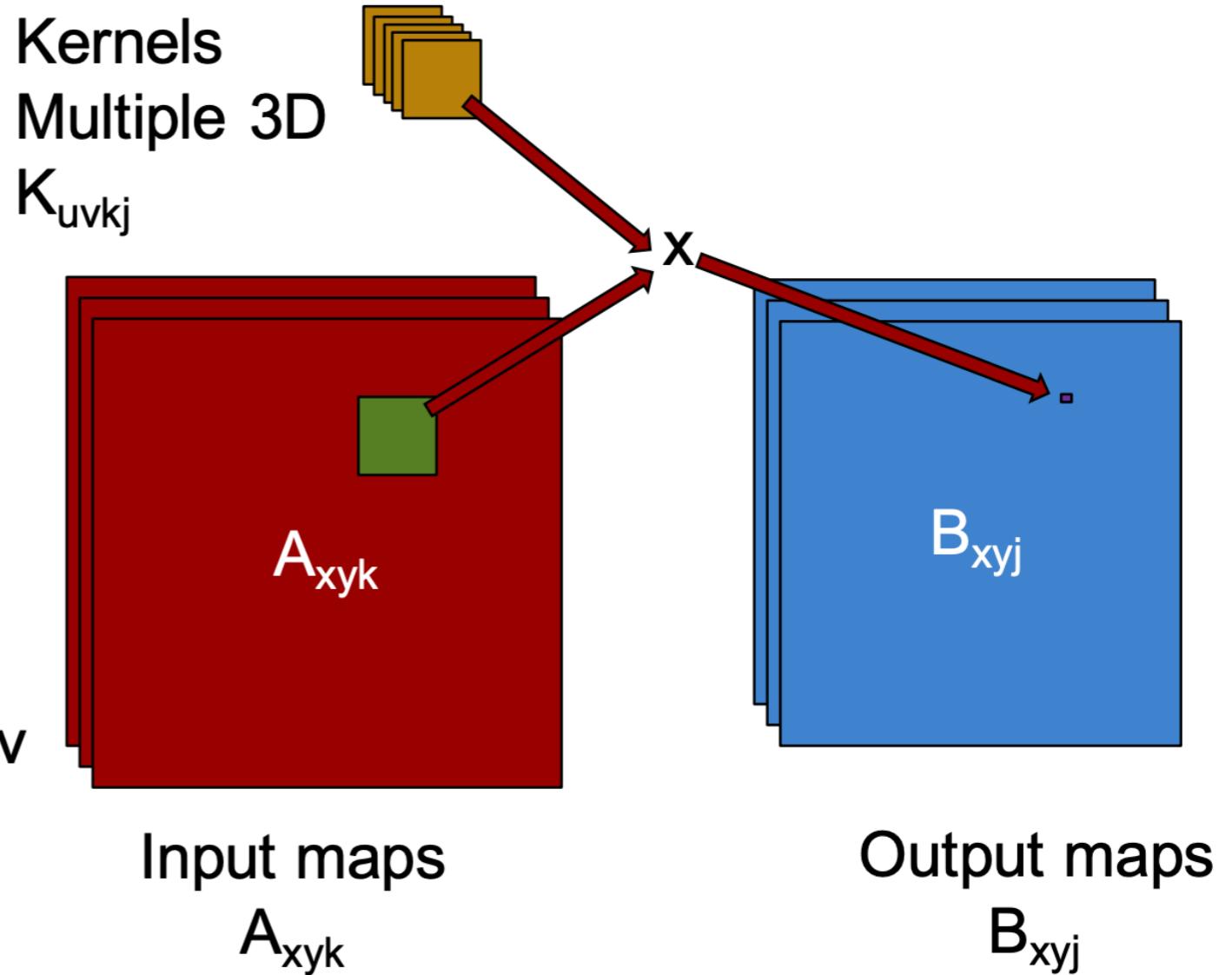
Model Parallel
Split up the Model – i.e.
the network

Model-Parallel Convolution



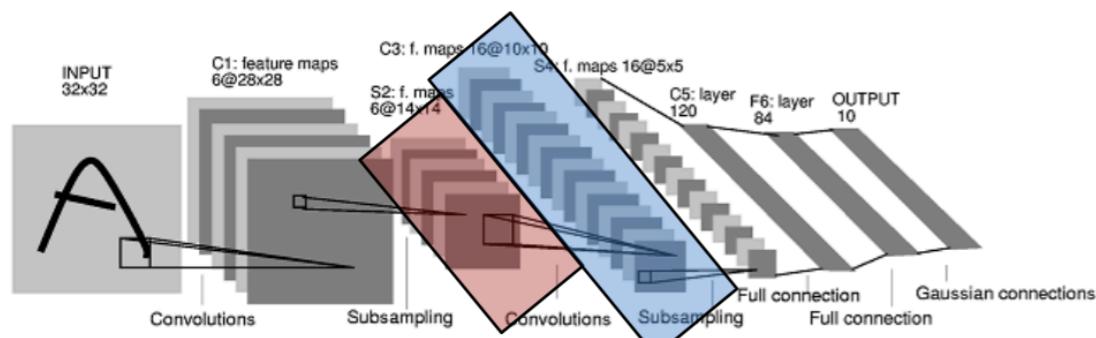
6D Loop
For each output map j
For each input map k
For each pixel x,y
For each kernel element u,v

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

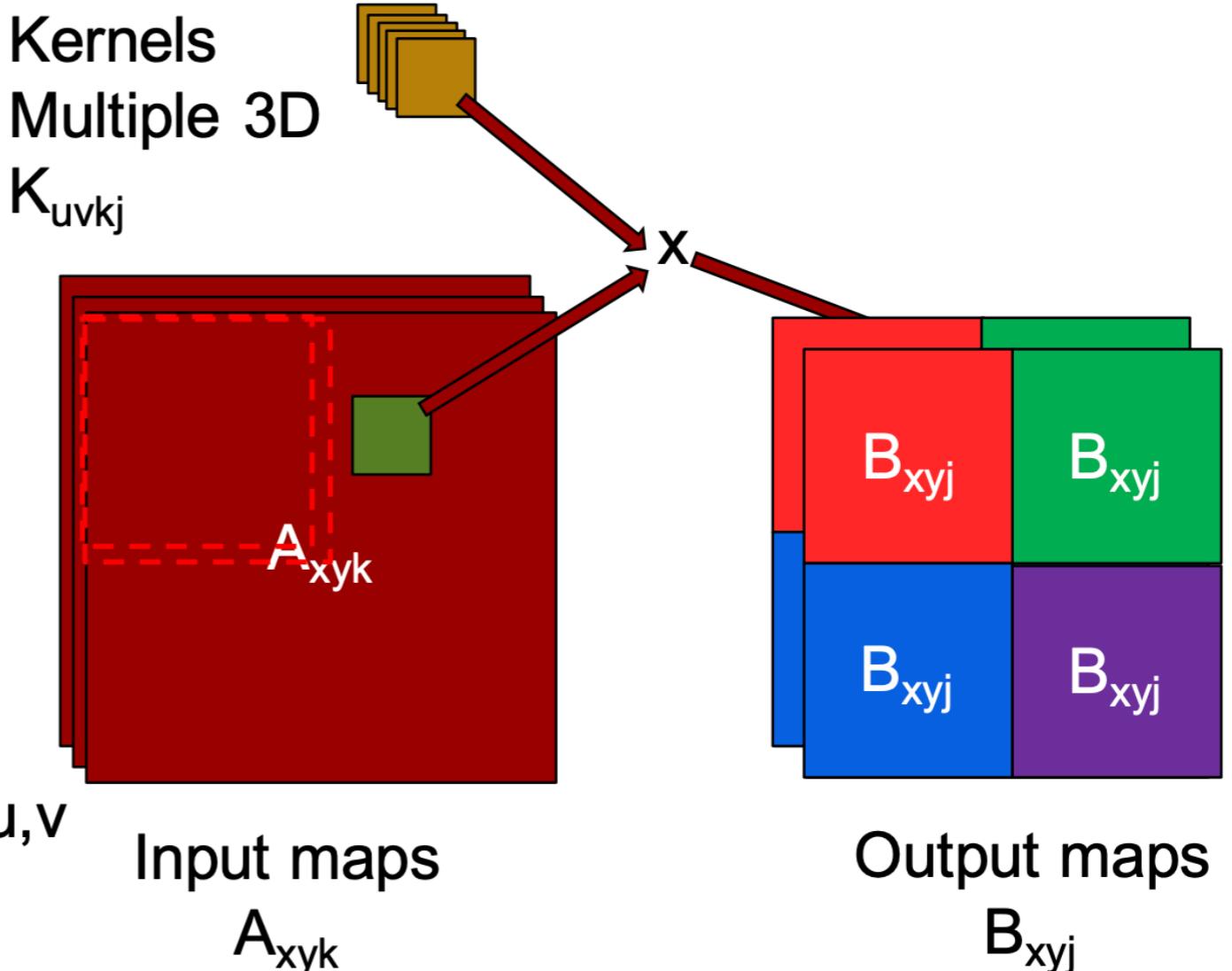


Model-Parallel Convolution

- by output region (x,y)

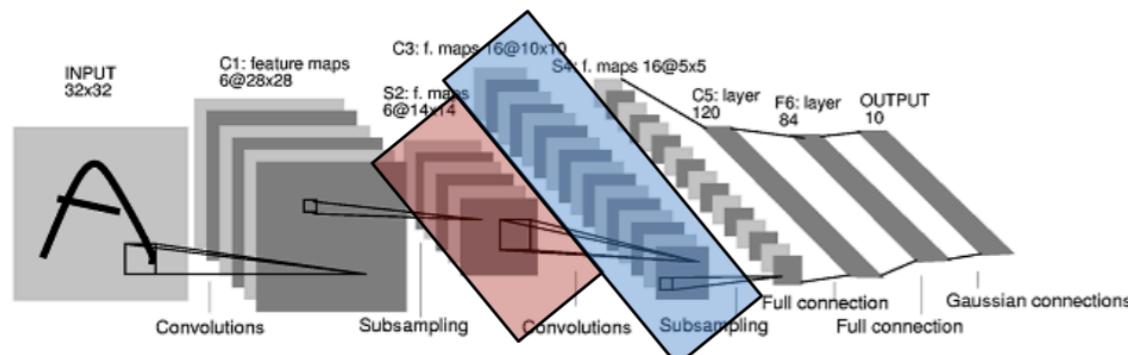


6D Loop
For all region XY
 For each output map j
 For each input map k
 For each pixel x,y in XY
 For each kernel element u,v
 $B_{xyj} += A_{(x-u)(y-v)k} \times K_{uvkj}$

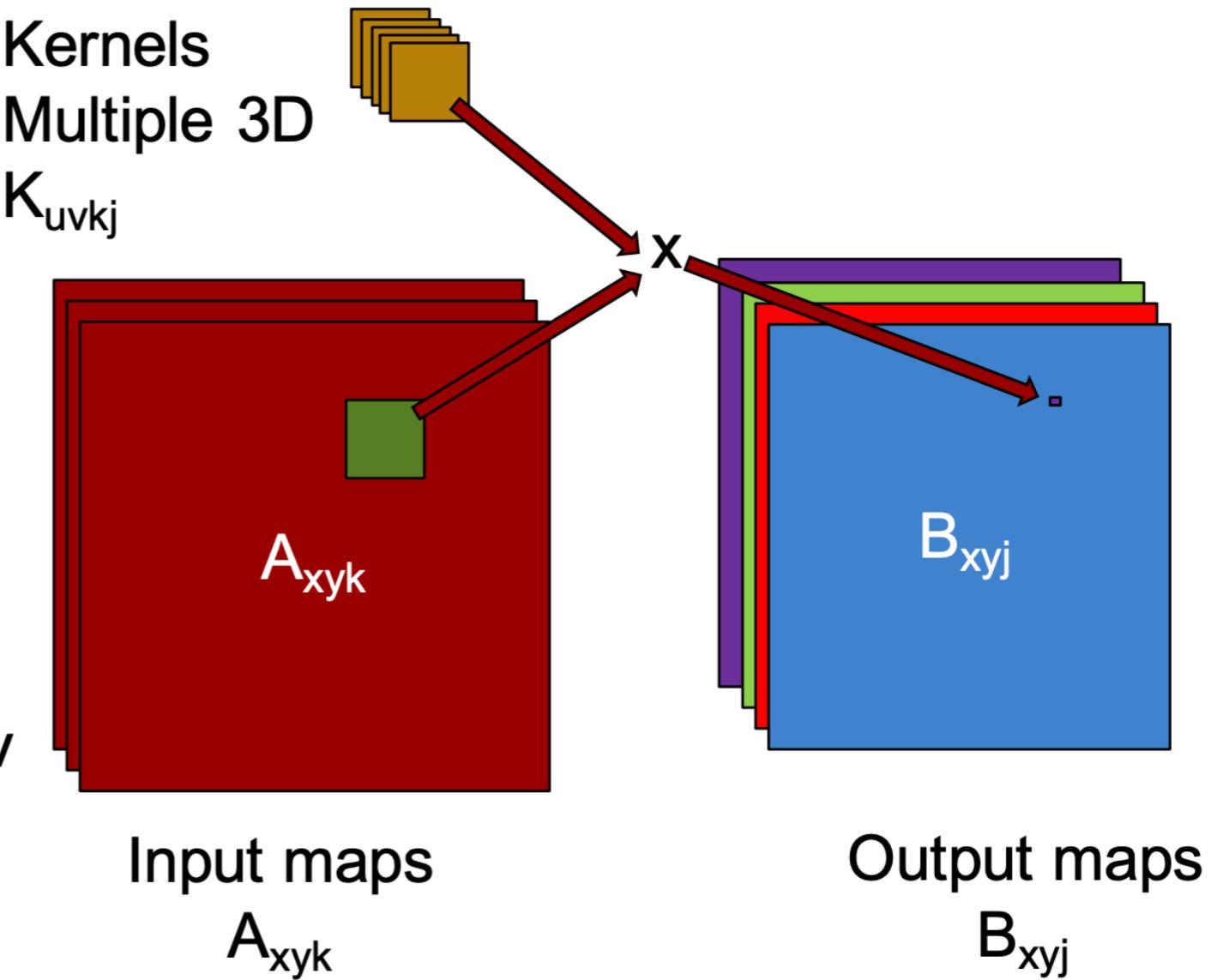


Model-Parallel Convolution

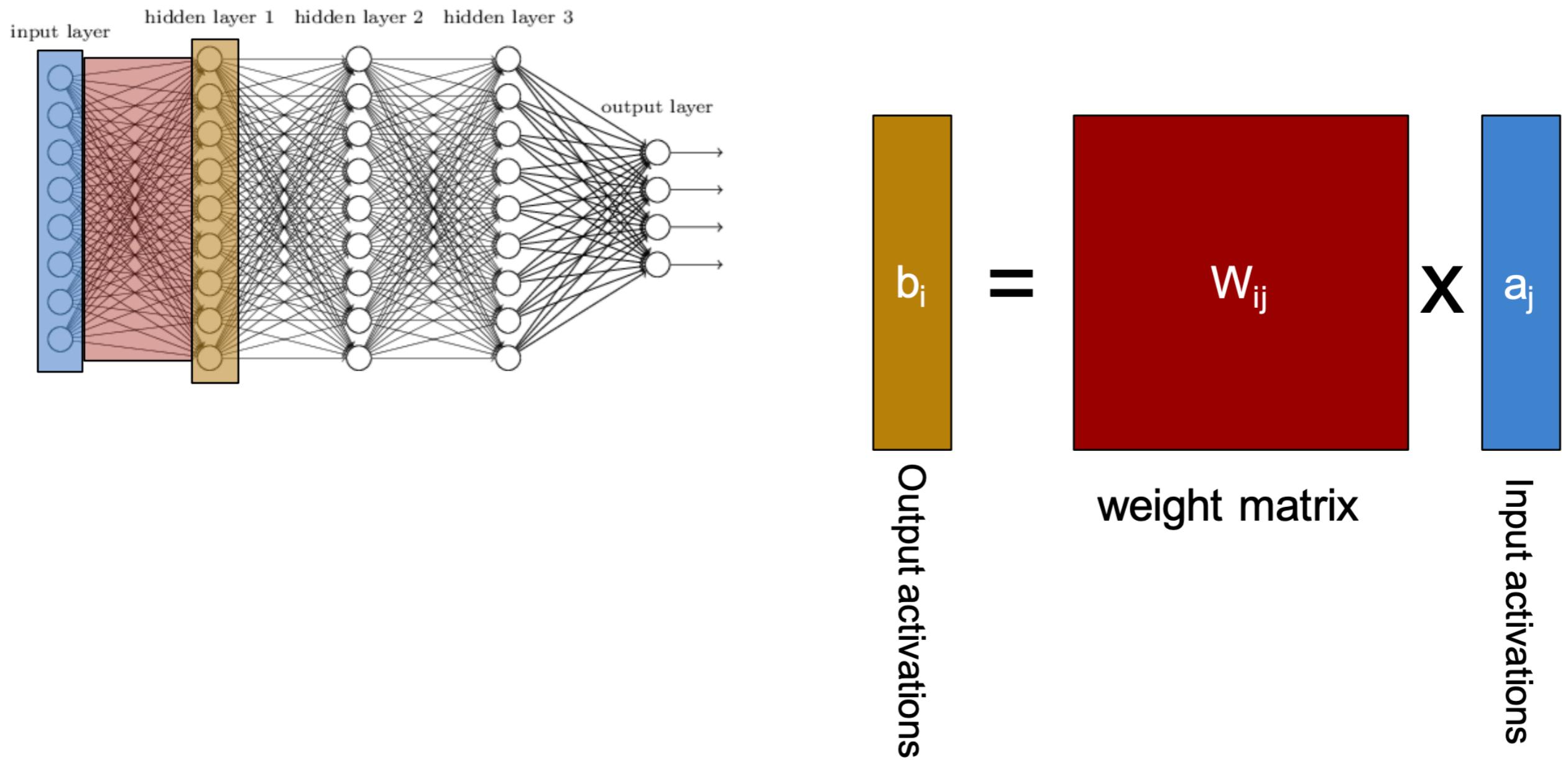
- By output map j (filter)



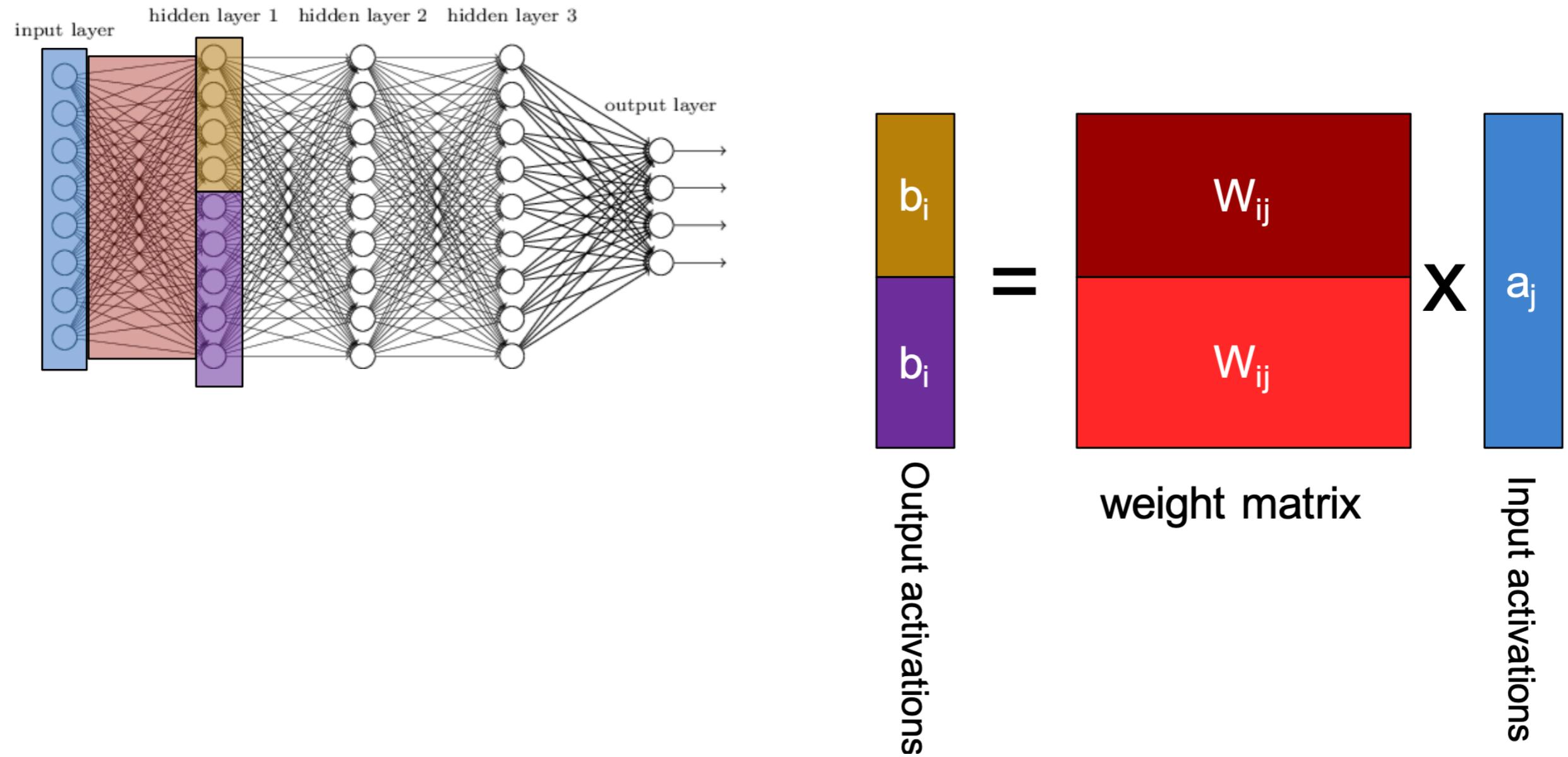
6D Loop
Forall output map j
 For each input map k
 For each pixel x,y
 For each kernel element u,v
 $B_{xyj} += A_{(x-u)(y-v)k} \times K_{uvkj}$



Model Parallel Fully-Connected Layer ($M \times V$)



Model Parallel Fully-Connected Layer ($M \times V$)

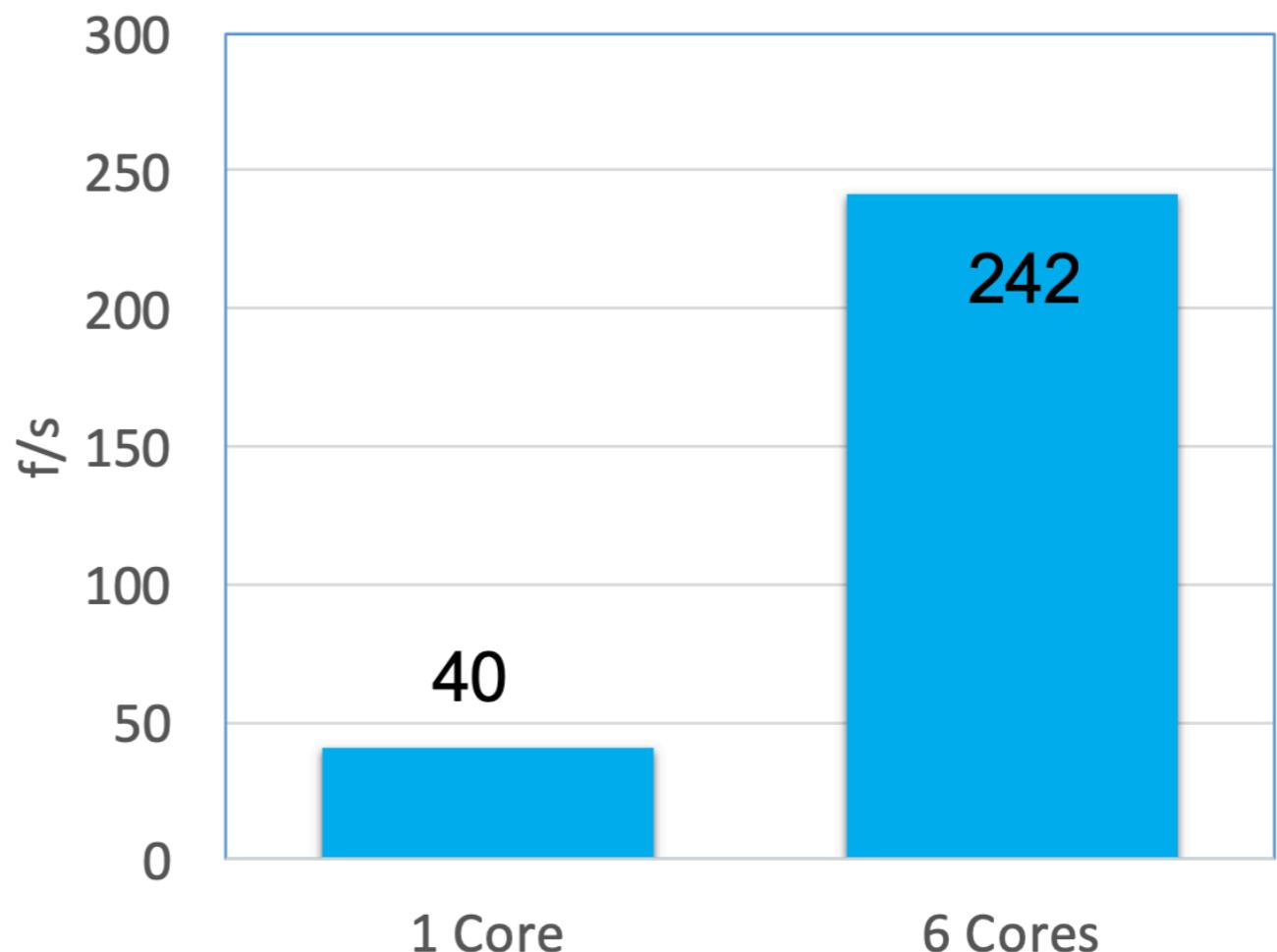


Hyper-Parameter Parallel

Try many alternative
networks in parallel

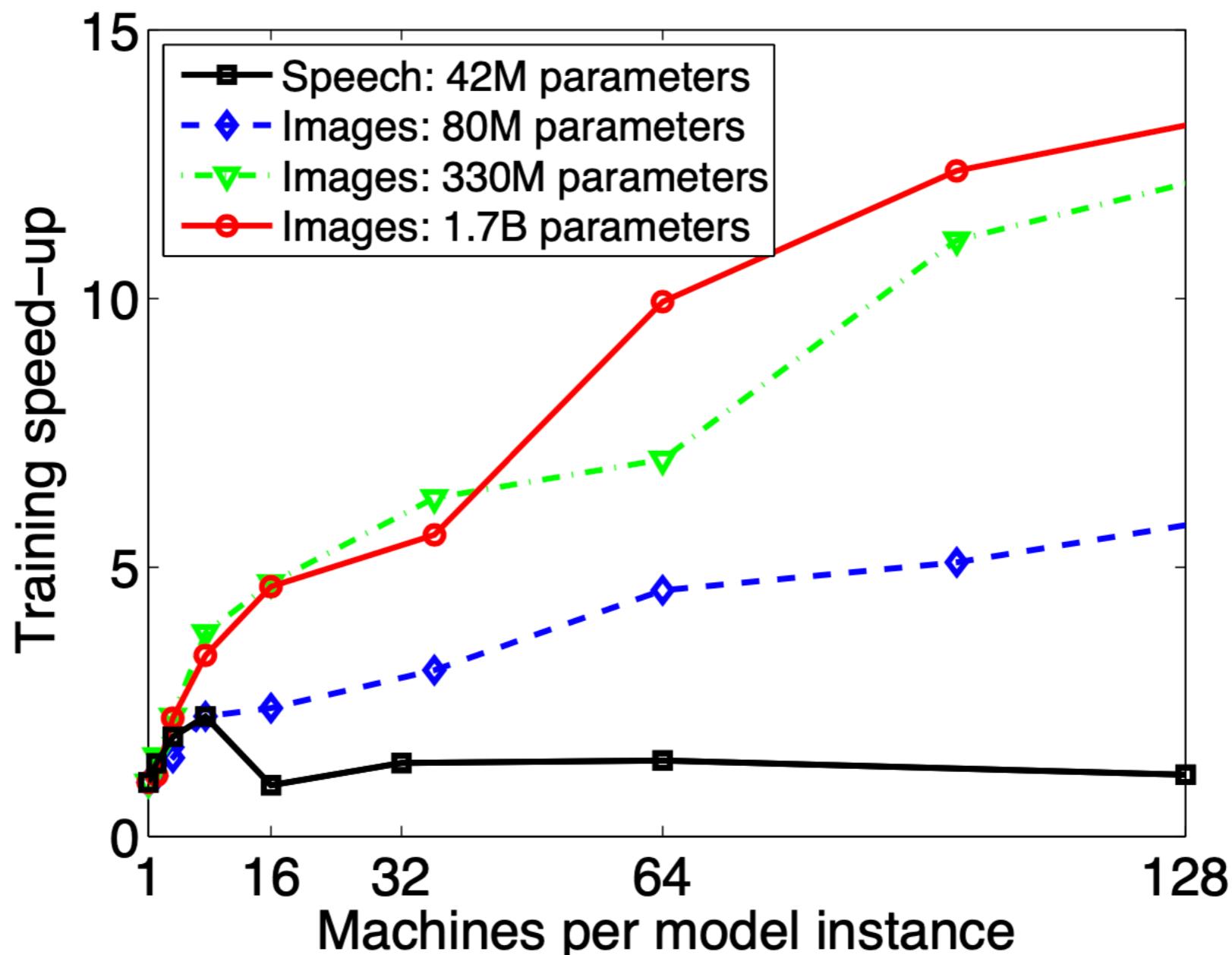
CPU Parallelism – Core i7 –

1 core vs 6 cores



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

Data and Model Parallel Performance



Summary of Parallelism

- Lots of parallelism in DNNs
 - 16M independent multiplies in one FC layer
 - Limited by overhead to exploit a fraction of this
- Data parallel
 - Run multiple training examples in parallel
 - Limited by batch size
- Modelparallel
 - Split model over multiple processors
 - By layer
 - Conv layers by map region
 - Fully connected layers by output activation
- Easy to get 16-64 GPUs training one model in parallel

GPUs

- To go fast, use multiple processors

- To go fast, use multiple processors
- To be efficient and fast, use GPUs

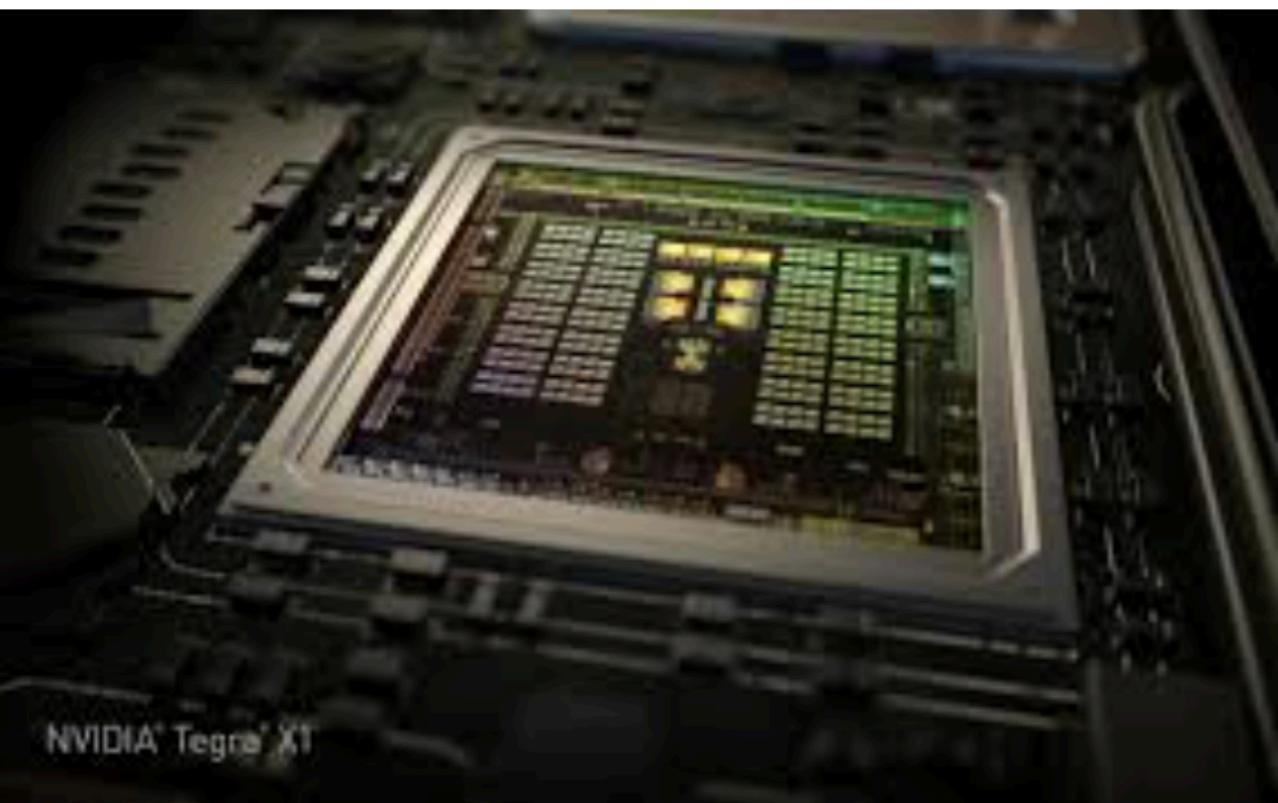
- To go fast, use multiple processors
- To be efficient and fast, use GPUs
- To be efficient and go really fast, use multiple GPUs

Titan X



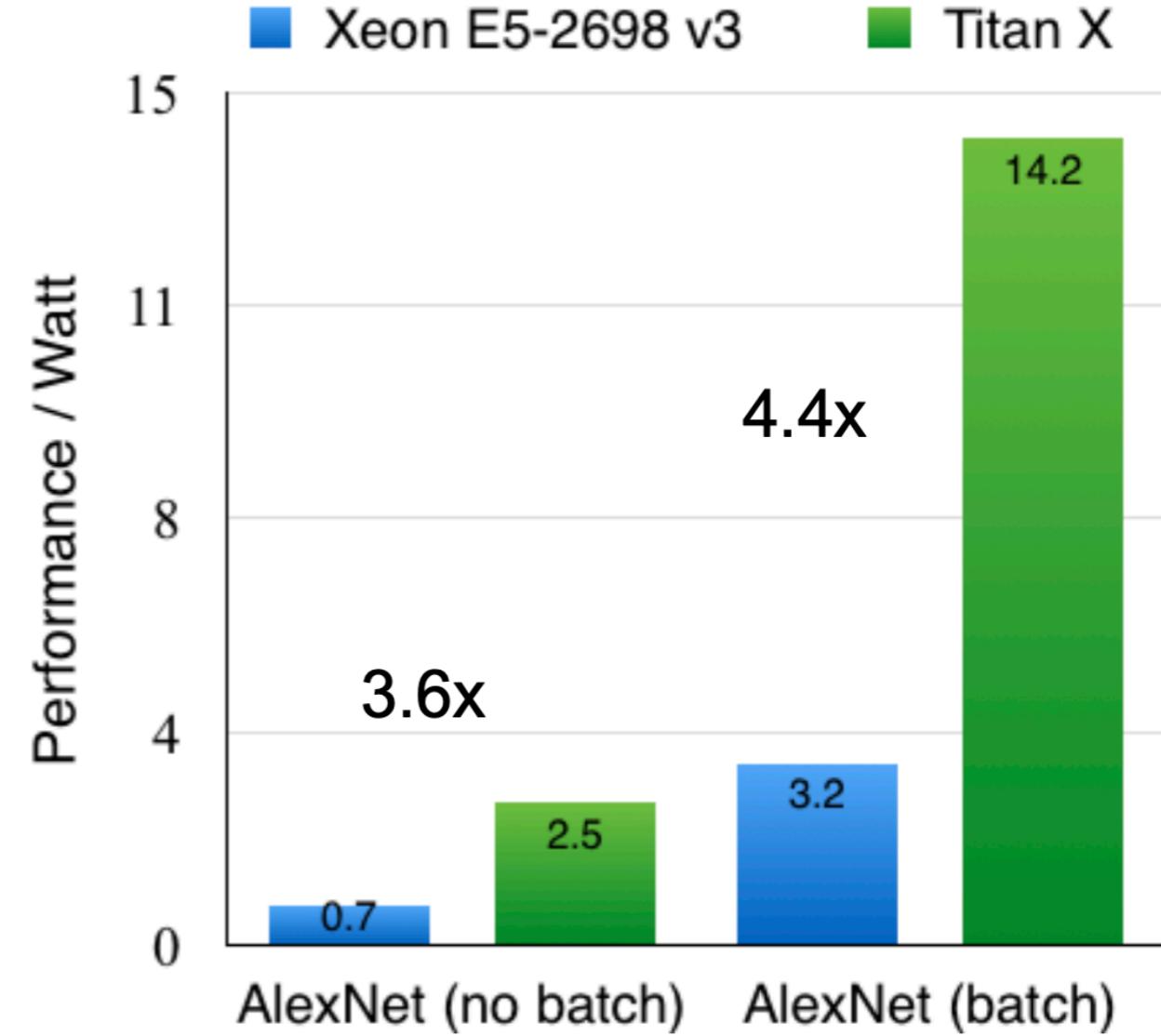
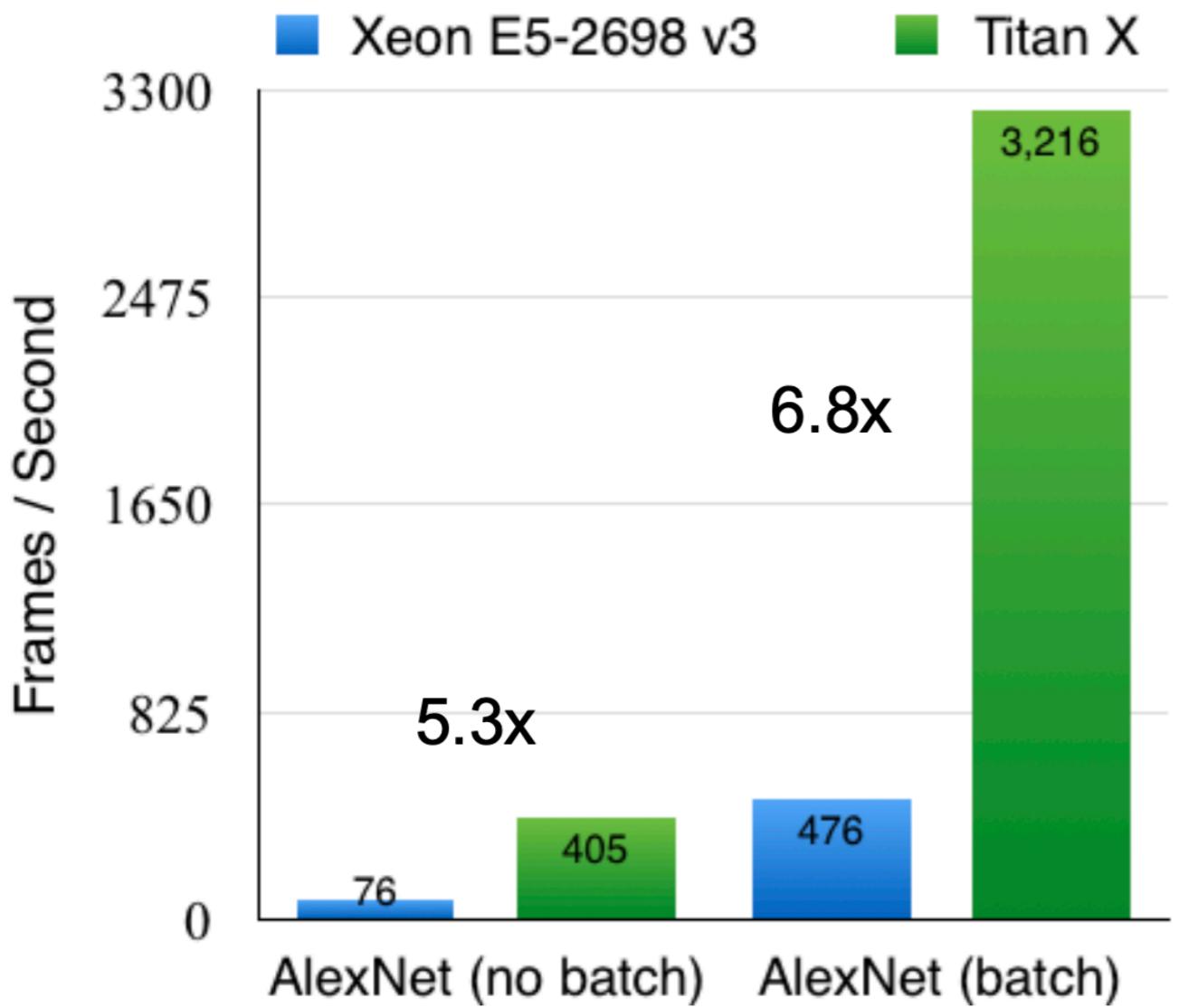
- 3072 CUDA cores @ 1 GHz
- 6 Teraflops FP32
- 12GB of GDDR5 @ 336 GB/sec
- 250W TDP
- 24GFLOPS/W
- 28nm process

Tegra X1



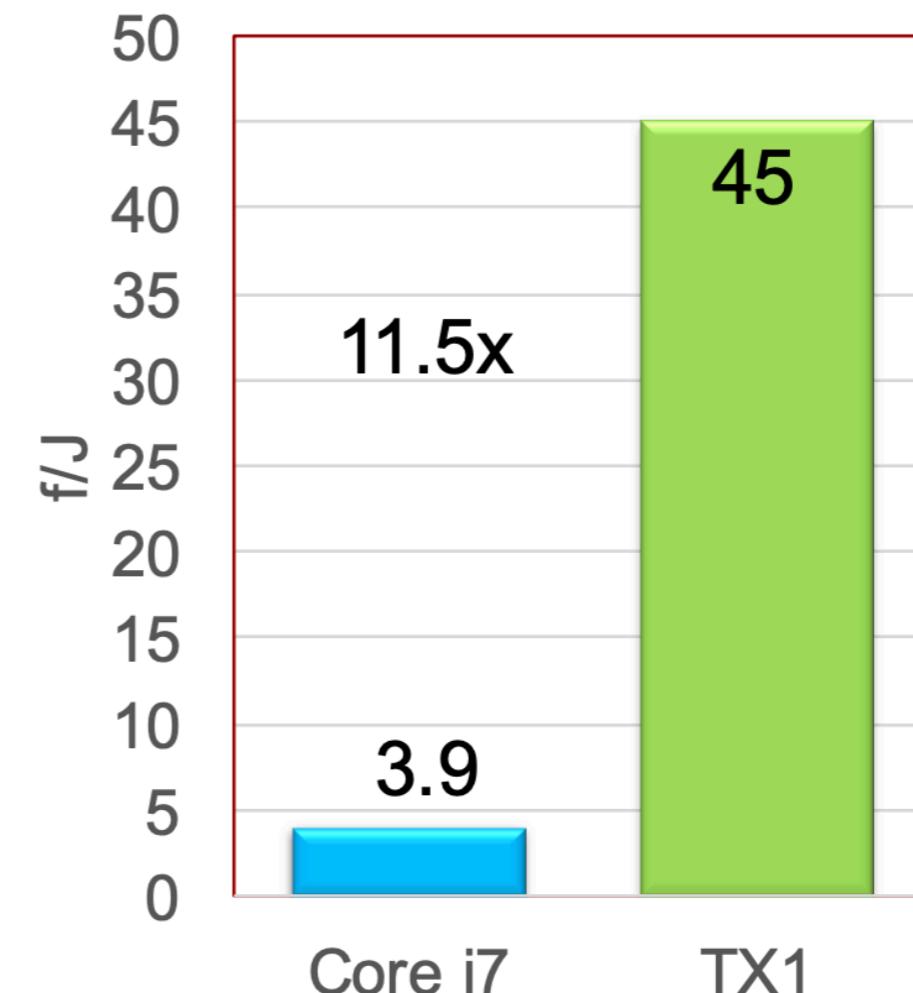
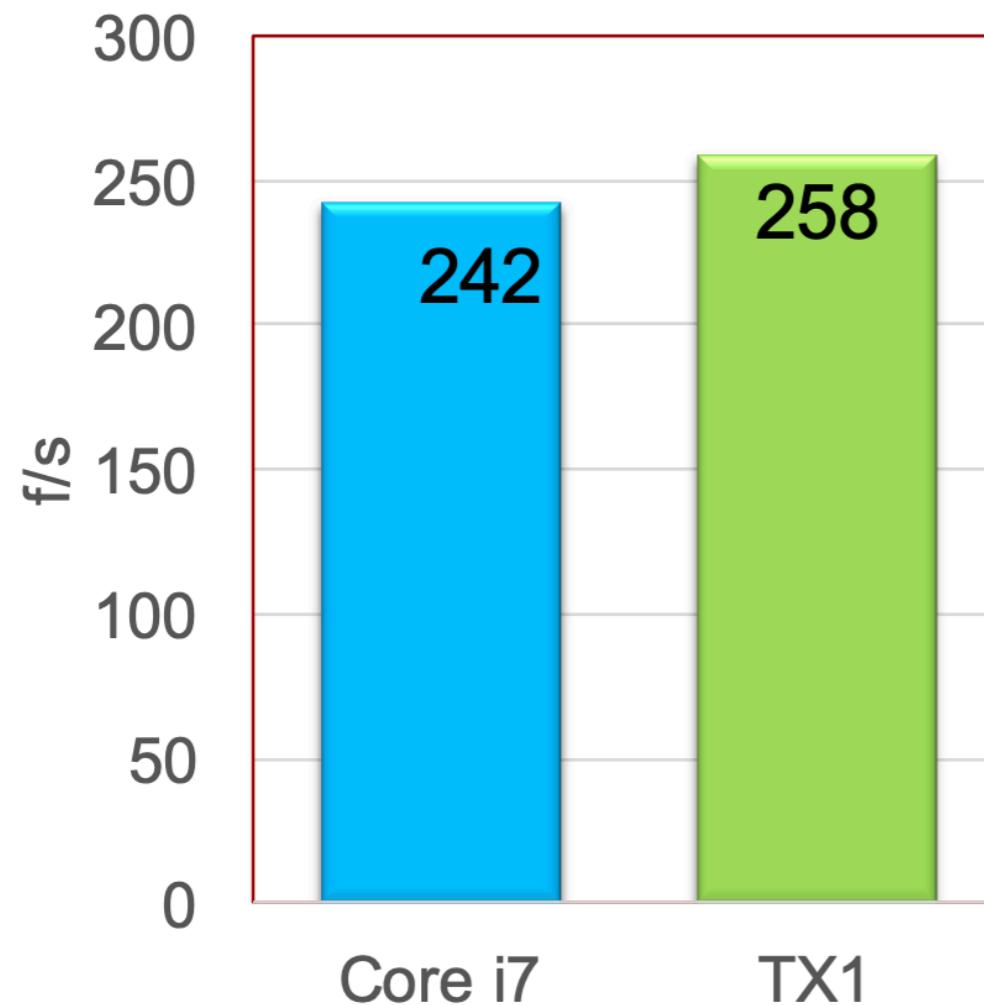
- 256 CUDA cores @ ~1 GHz
- 1 Teraflop FP16
- 4GB of LPDDR4 @ 25.6 GB/s
- 15 W TDP (1W idle, <10W typical)
- 100GFLOPS/W (FP16)
- 20nm process

Xeon E5-2698 CPU v.s. TitanX GPU



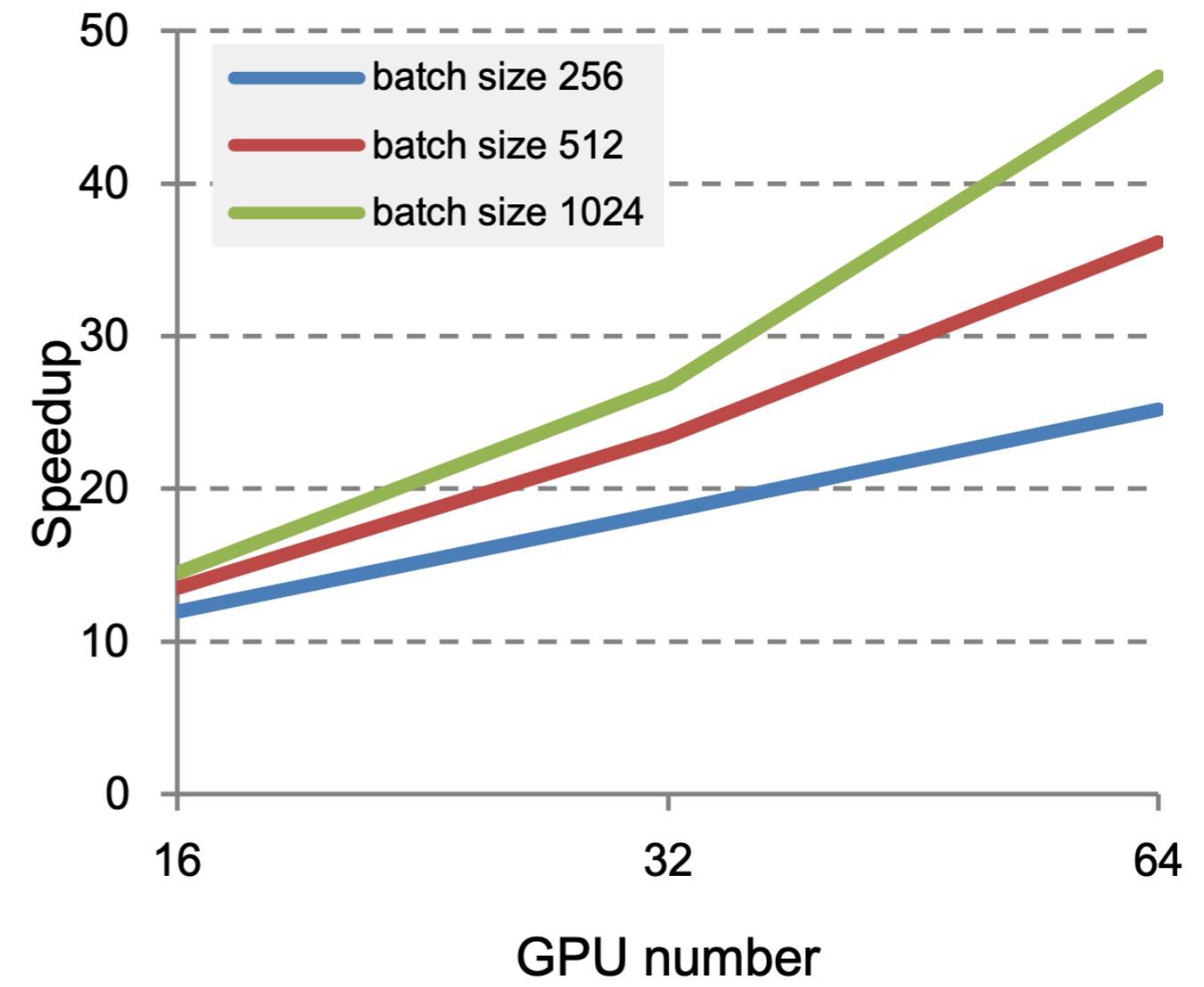
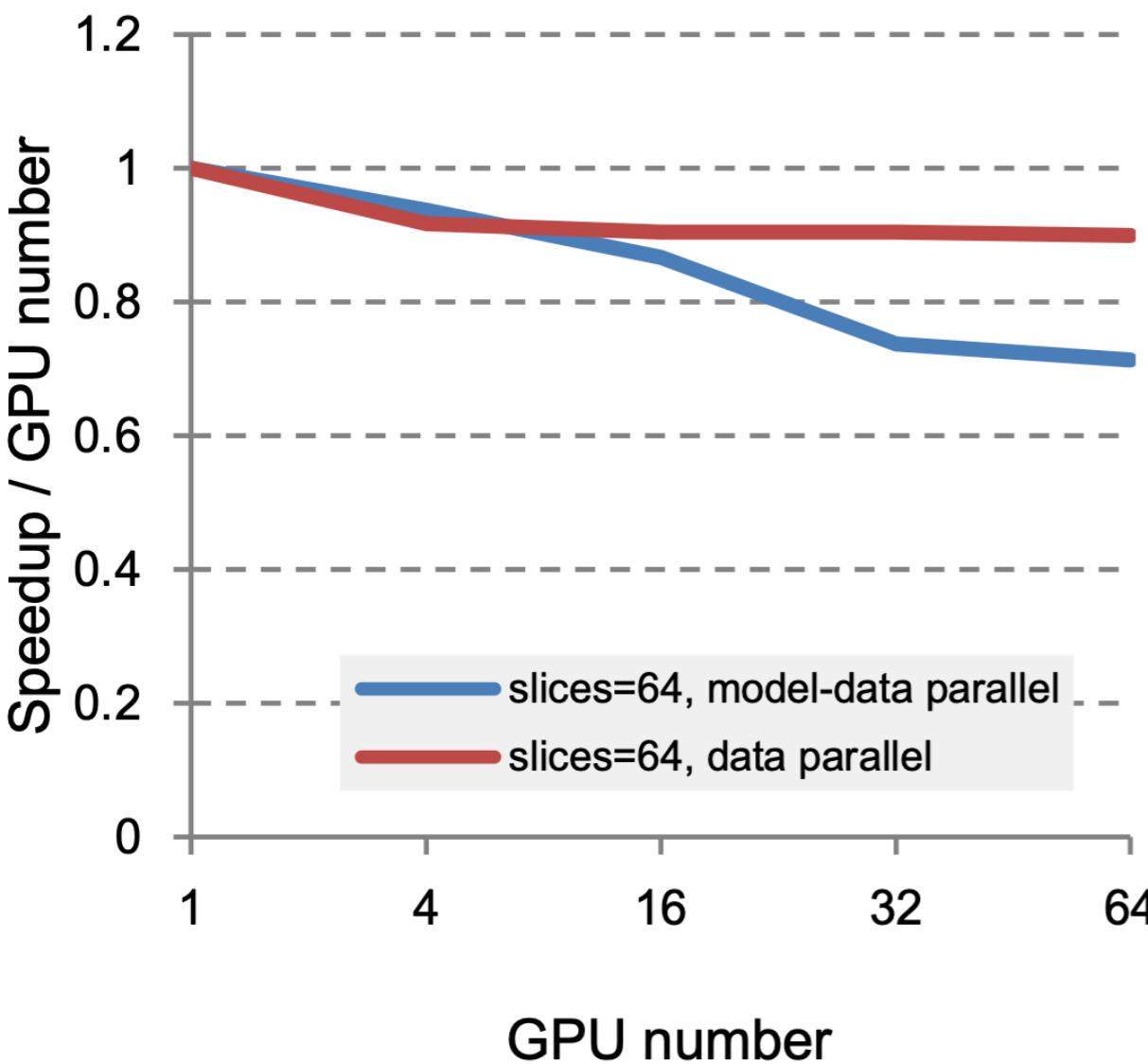
NVIDIA, "Whitepaper: GPU-based deep learning inference: A performance and power analysis."

Tegra X1 vs Core i7

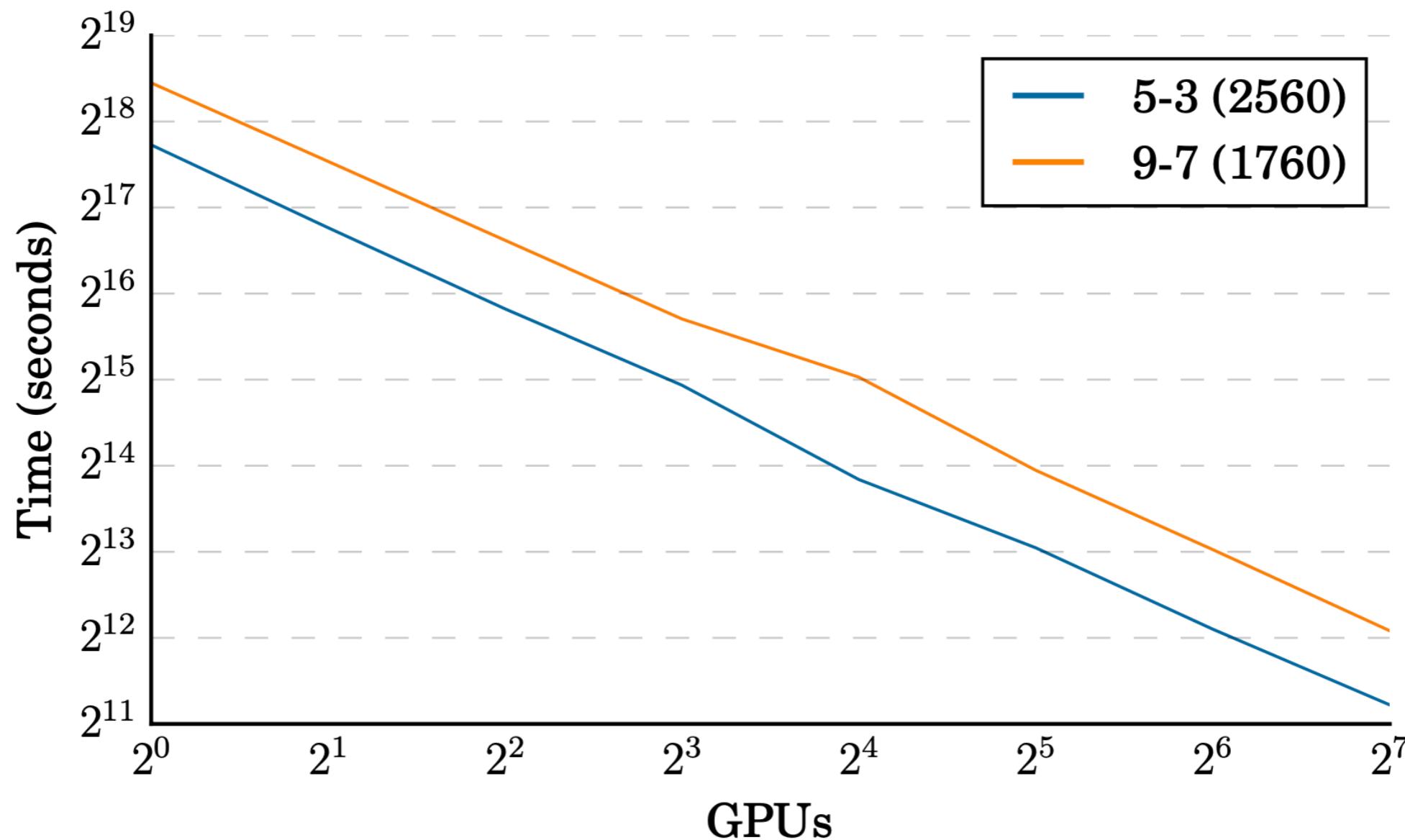


NVIDIA, "Whitepaper: GPU-based deep learning inference: A performance and power analysis."

Parallel GPUs



Parallel GPUs on Deep Speech 2



Summary of GPUs

- Titan X ~6x faster, 4x more efficient than Xeon E5
 - TX1 11.5x more efficient than Core i7
 - On inference
 - Larger gains on training
-
- Data parallelism scales easily to 16GPUs
 - With some effort, linear speedup to 128GPUs