

Lecture 8:

Deep Neural Network Evaluation

**Visual Computing Systems
CMU 15-769, Fall 2016**

Training/evaluating deep neural networks

Technique leading to many high-profile AI advances in recent years

Speech recognition/natural language processing

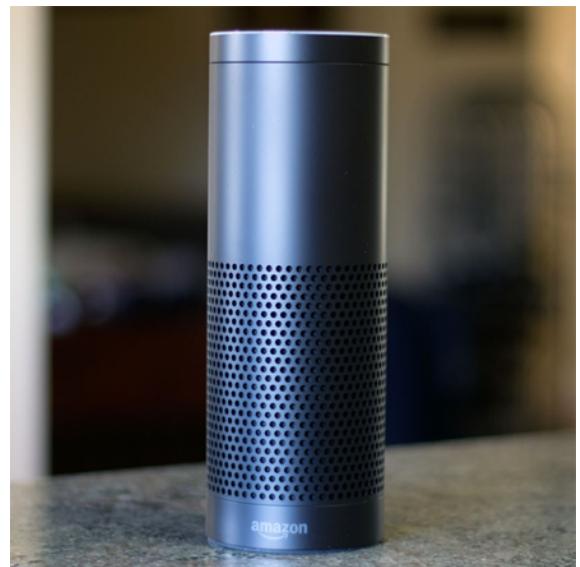
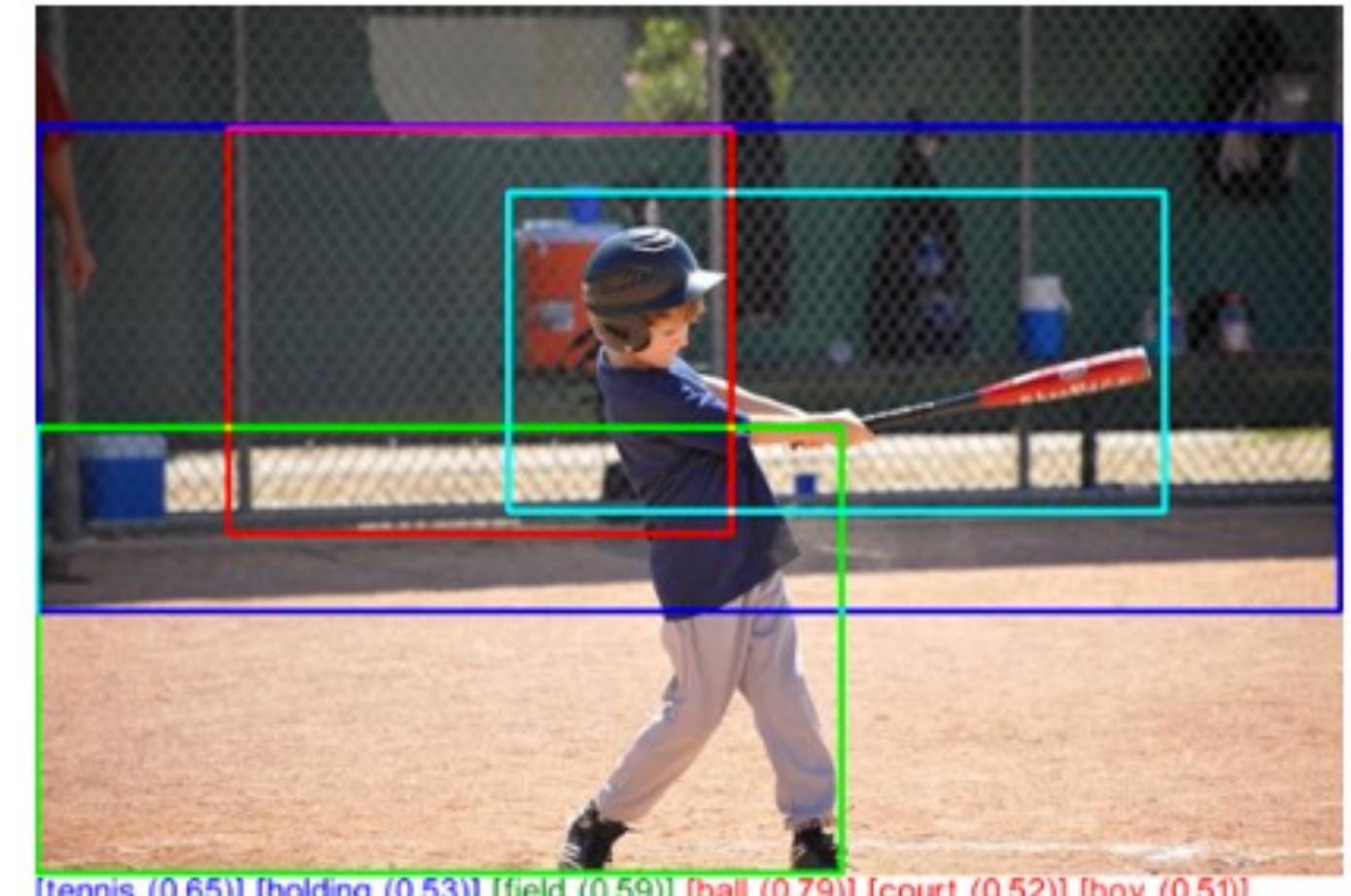


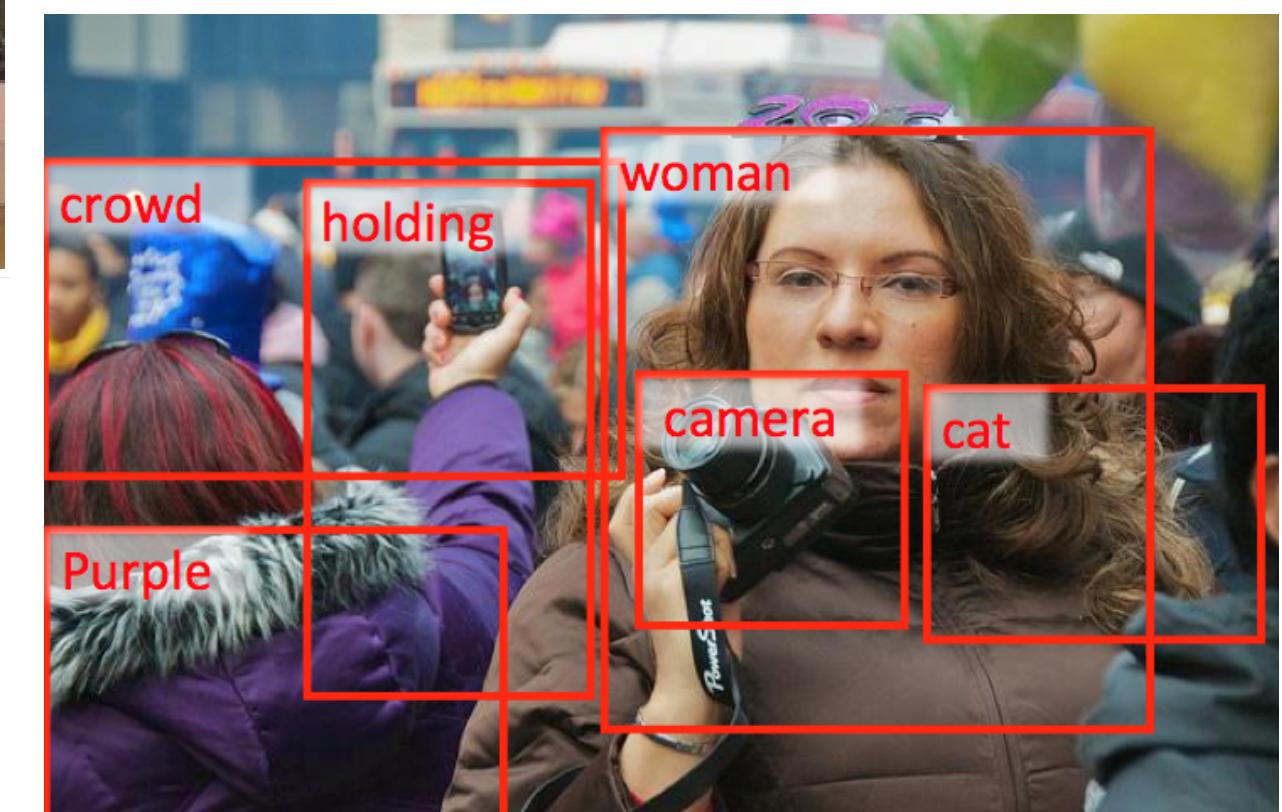
Image interpretation and understanding



Game AI



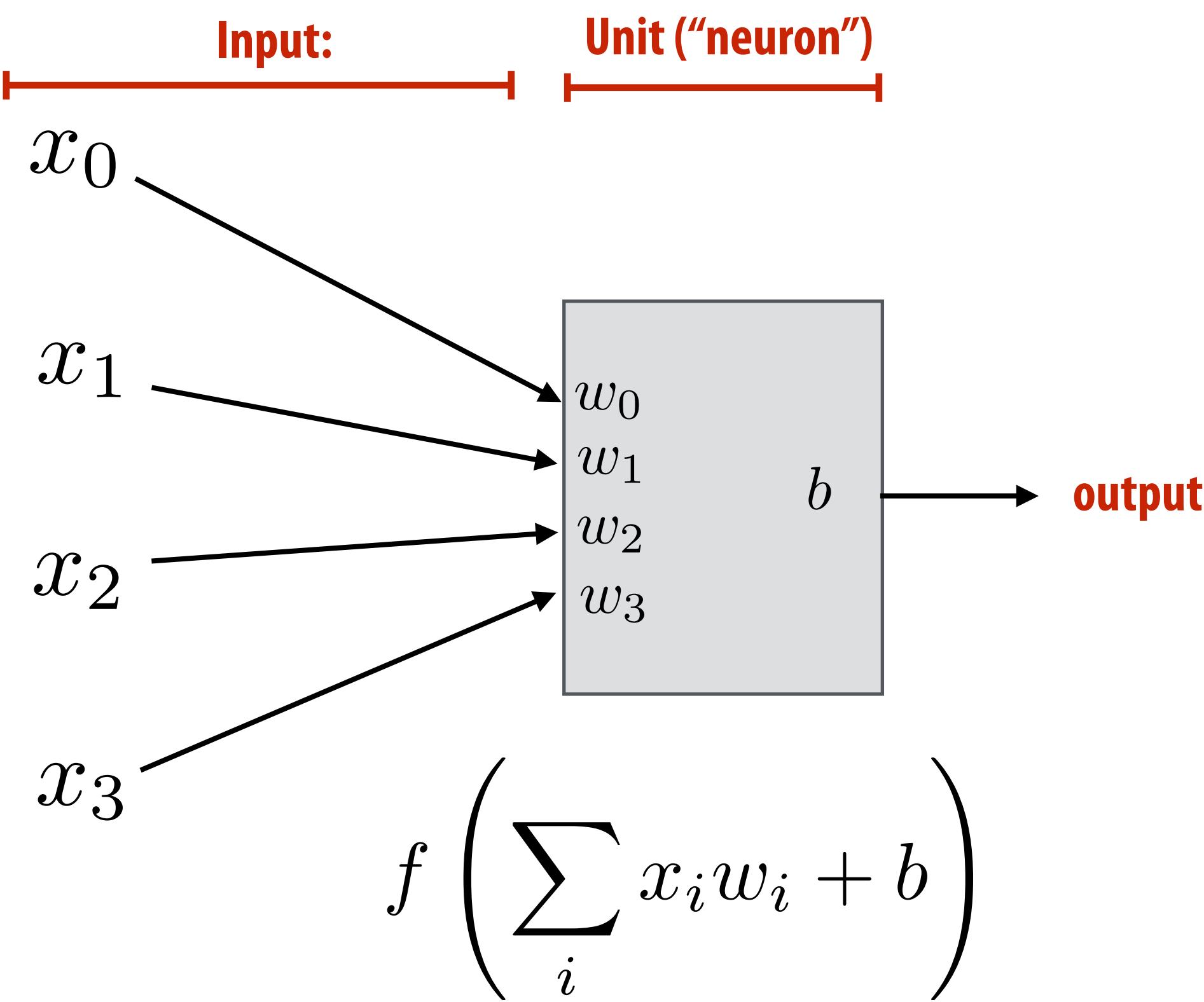
a baseball player swinging a bat at a ball
a boy is playing with a baseball bat



What is a deep neural network?

A basic unit:

Unit with n inputs described by $n+1$ parameters
(weights + bias)

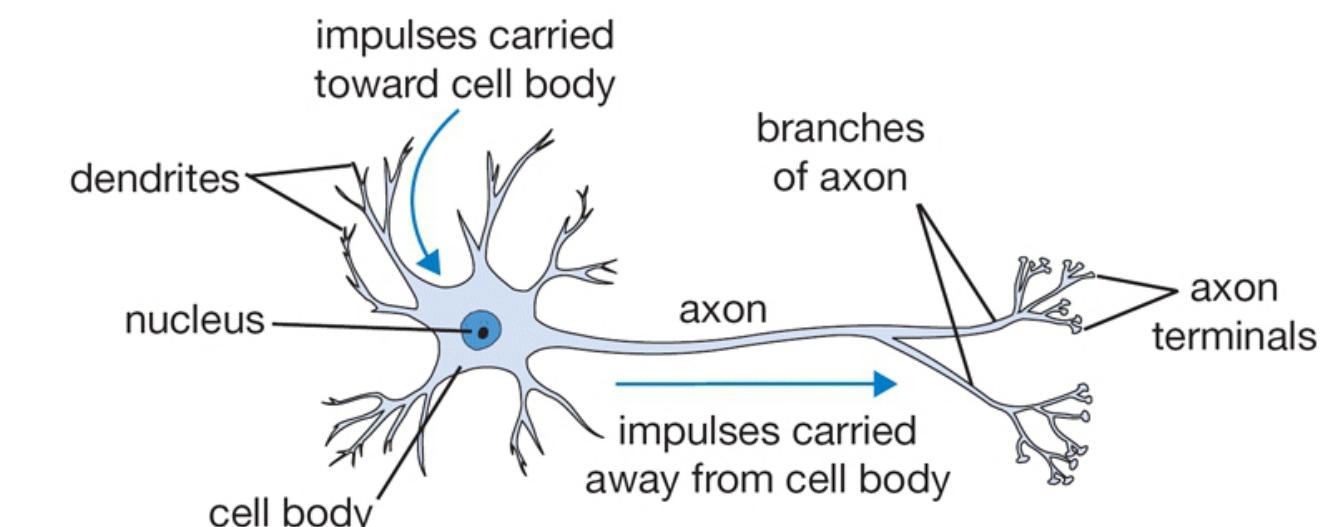


Example: rectified linear unit (ReLU)

$$f(x) = \max(0, x)$$

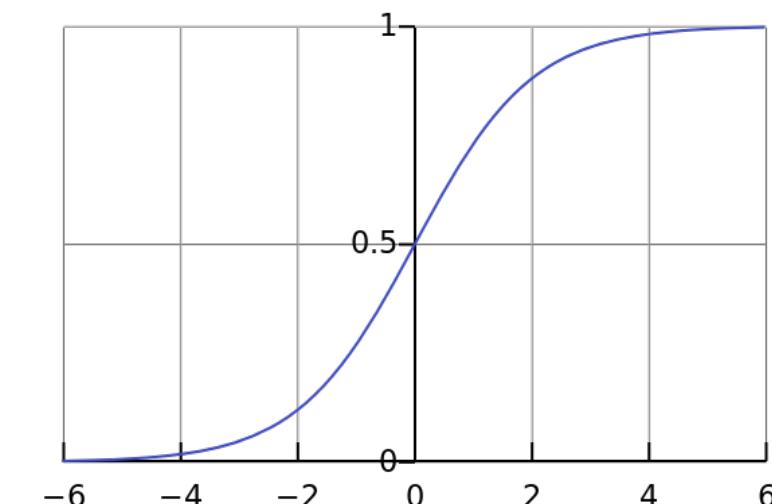
Basic computational interpretation:
It's just a circuit!

Biological inspiration:
unit output corresponds loosely to activation of neuron



Machine learning interpretation:
binary classifier: interpret output as the probability of one class

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



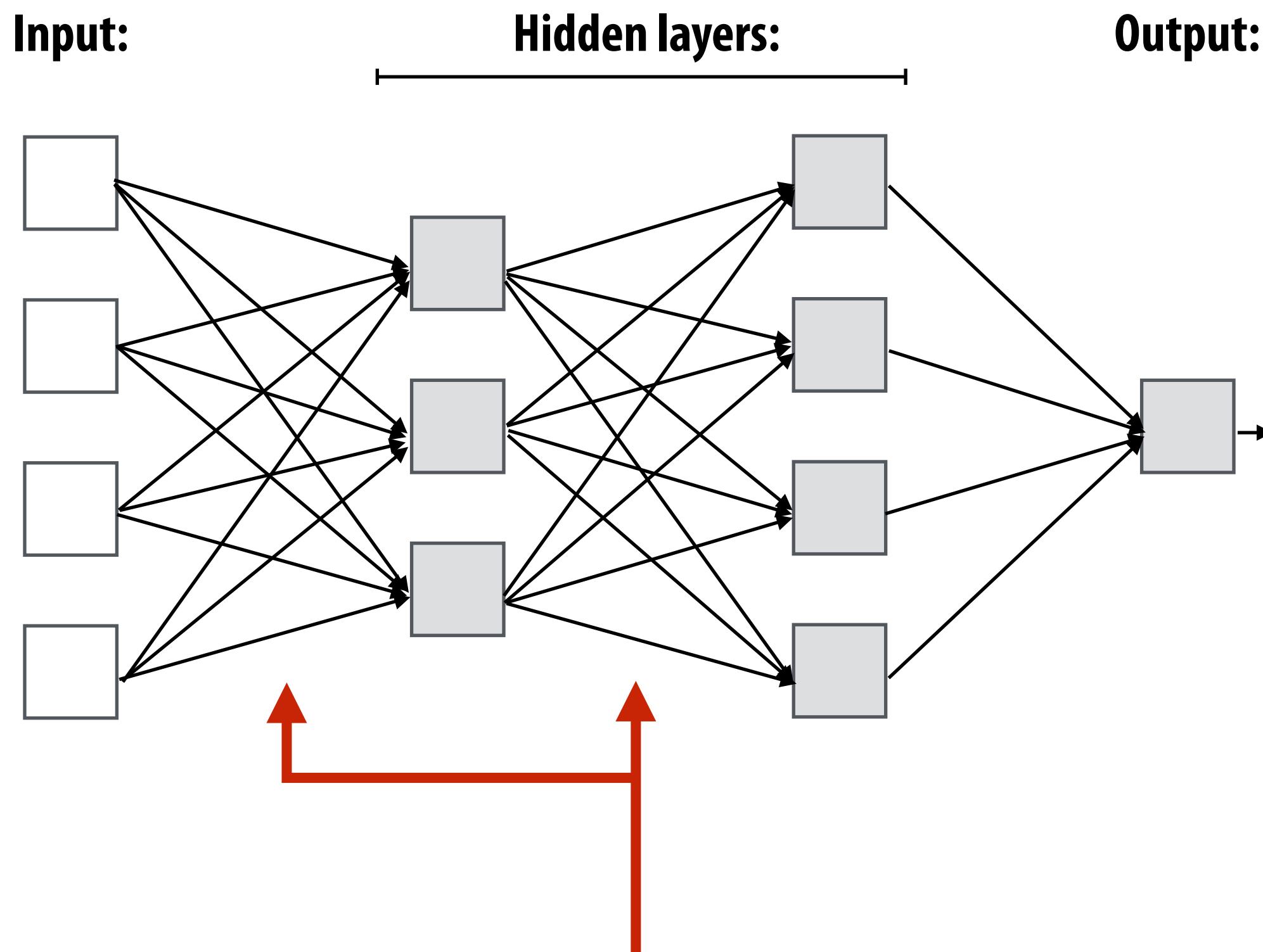
What is a deep neural network: topology

This network has: 4 inputs, 1 output, 7 hidden units

“Deep”= at least one hidden layer

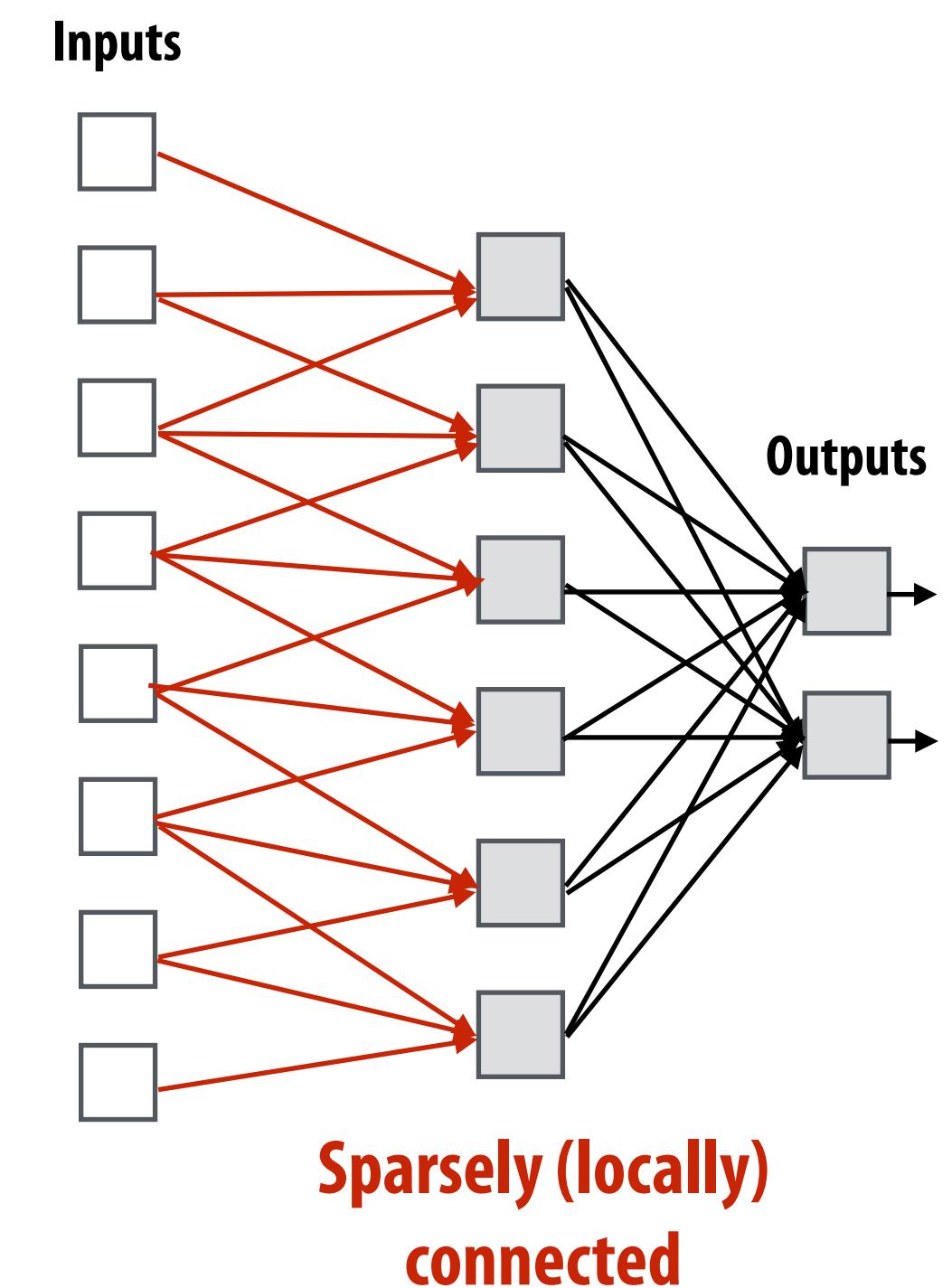
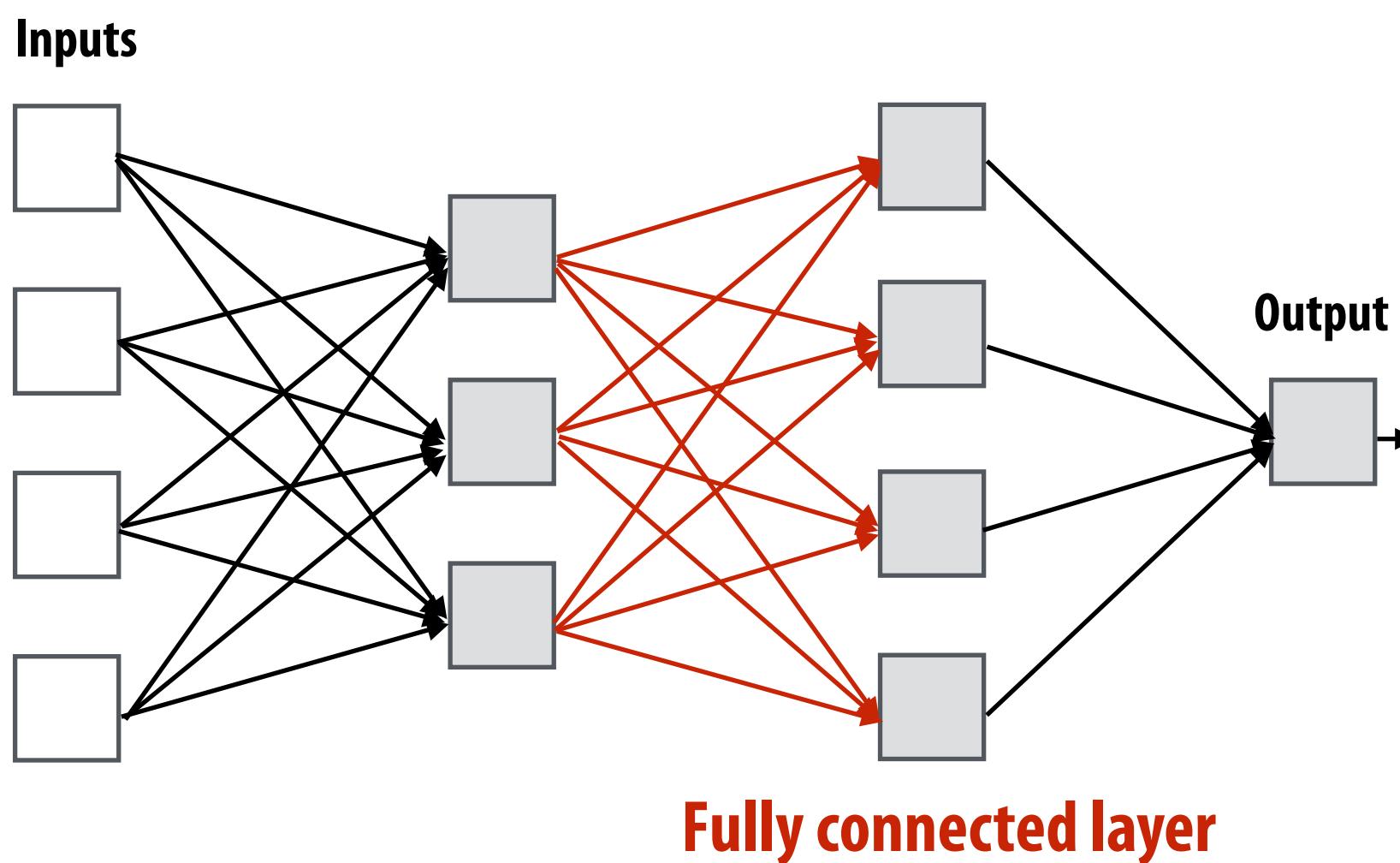
Hidden layer 1: 3 units x (4 weights + 1 bias) = 15 parameters

Hidden layer 2: 4 units x (3 weights + 1 bias) = 16 parameters



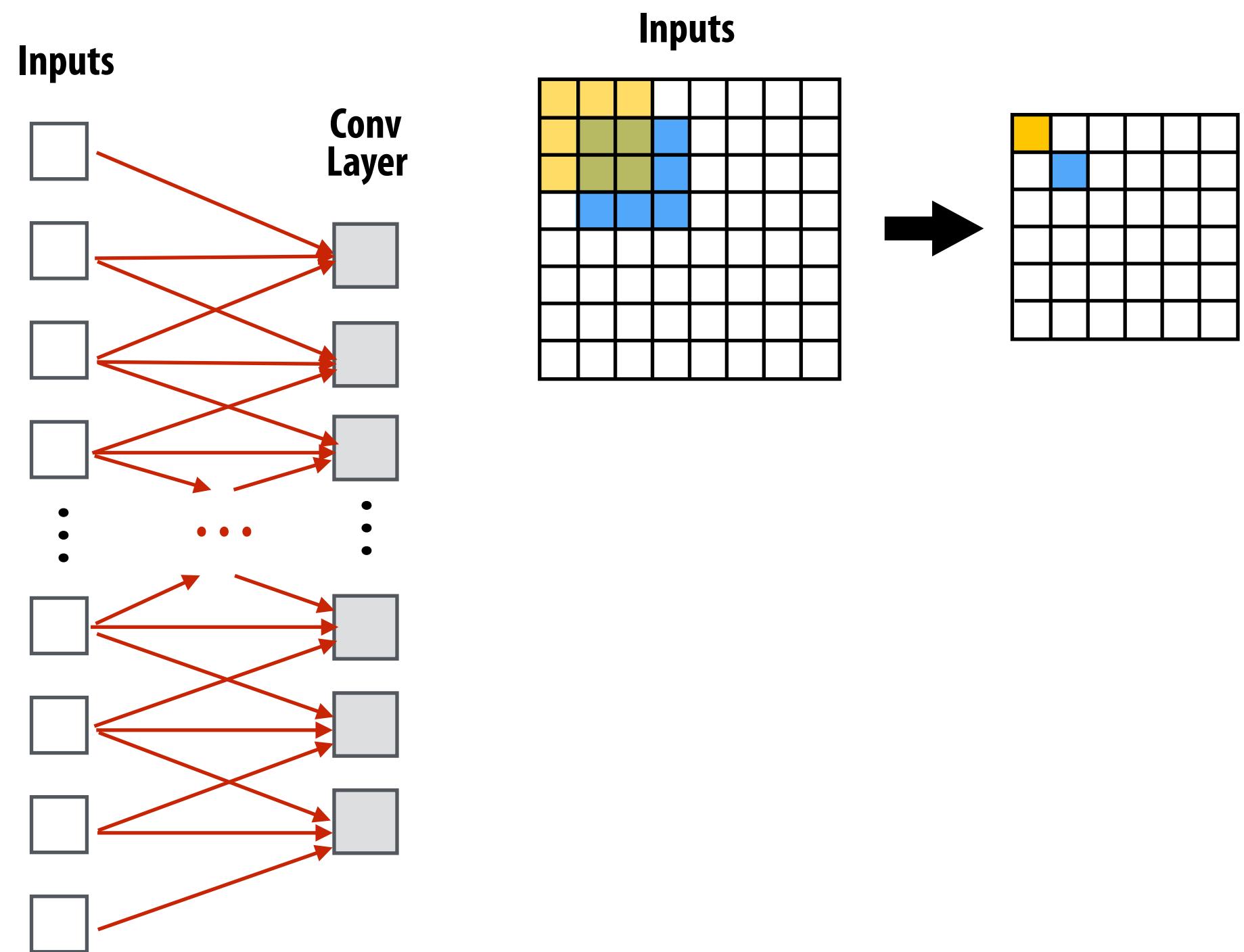
Note “fully-connected” topology in this example

What is a deep neural network: topology



Recall image convolution (3x3 conv)

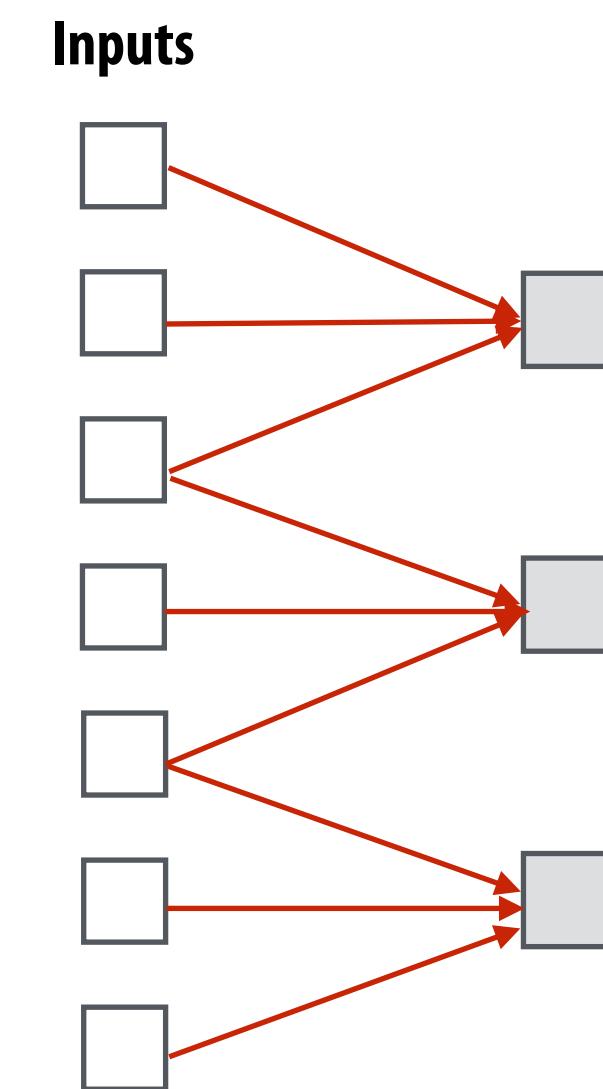
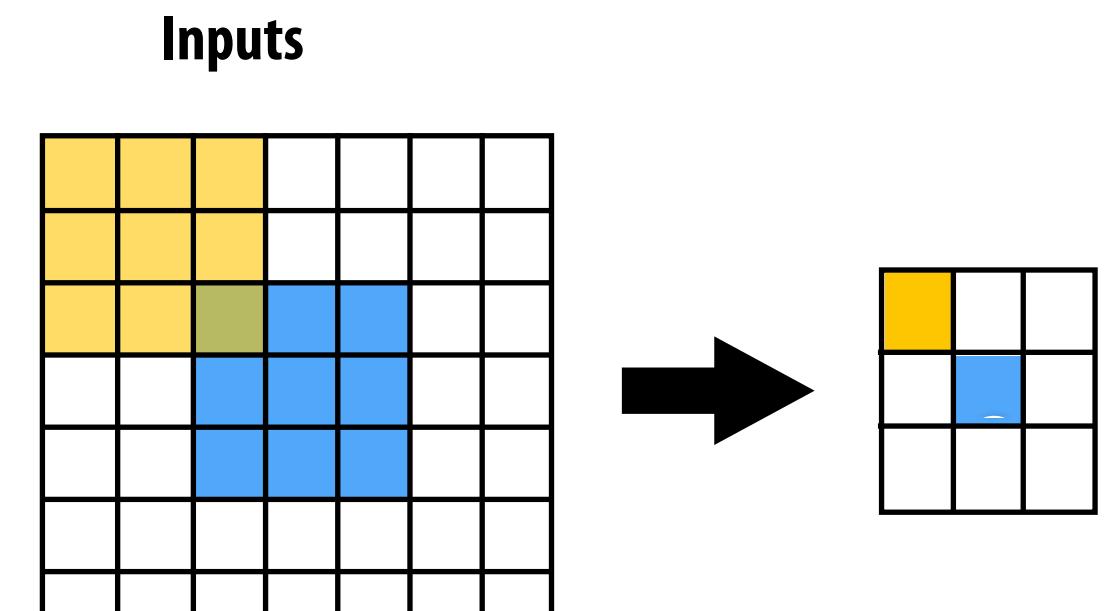
```
int WIDTH = 1024;  
int HEIGHT = 1024;  
float input[(WIDTH+2) * (HEIGHT+2)];  
float output[WIDTH * HEIGHT];  
  
float weights[] = {1.0/9, 1.0/9, 1.0/9,  
                   1.0/9, 1.0/9, 1.0/9,  
                   1.0/9, 1.0/9, 1.0/9};  
  
for (int j=0; j<HEIGHT; j++) {  
    for (int i=0; i<WIDTH; i++) {  
        float tmp = 0.f;  
        for (int jj=0; jj<3; jj++)  
            for (int ii=0; ii<3; ii++)  
                tmp += input[(j+jj)*(WIDTH+2) + (i+ii)] * weights[jj*3 + ii];  
        output[j*WIDTH + i] = tmp;  
    }  
}
```



Convolutional layer: locally connected AND all units in layer share the same parameters (same weights + same bias):
(note: network diagram only shows links for a 1D conv: a.k.a. one iteration of *i* in loop)

Strided 3x3 convolution

```
int WIDTH = 1024;  
int HEIGHT = 1024;  
int STRIDE = 2;  
float input[(WIDTH+2) * (HEIGHT+2)];  
float output[(WIDTH/STRIDE) * (HEIGHT/STRIDE)];  
  
float weights[] = {1.0/9, 1.0/9, 1.0/9,  
                   1.0/9, 1.0/9, 1.0/9,  
                   1.0/9, 1.0/9, 1.0/9};  
  
for (int j=0; j<HEIGHT; j+=STRIDE) {  
    for (int i=0; i<WIDTH; i+=STRIDE) {  
        float tmp = 0.f;  
        for (int jj=0; jj<3; jj++)  
            for (int ii=0; ii<3; ii++) {  
                tmp += input[(j+jj)*(WIDTH+2) + (i+ii)] * weights[jj*3 + ii];  
            }  
        output[(j/STRIDE)*WIDTH + (i/STRIDE)] = tmp;  
    }  
}
```

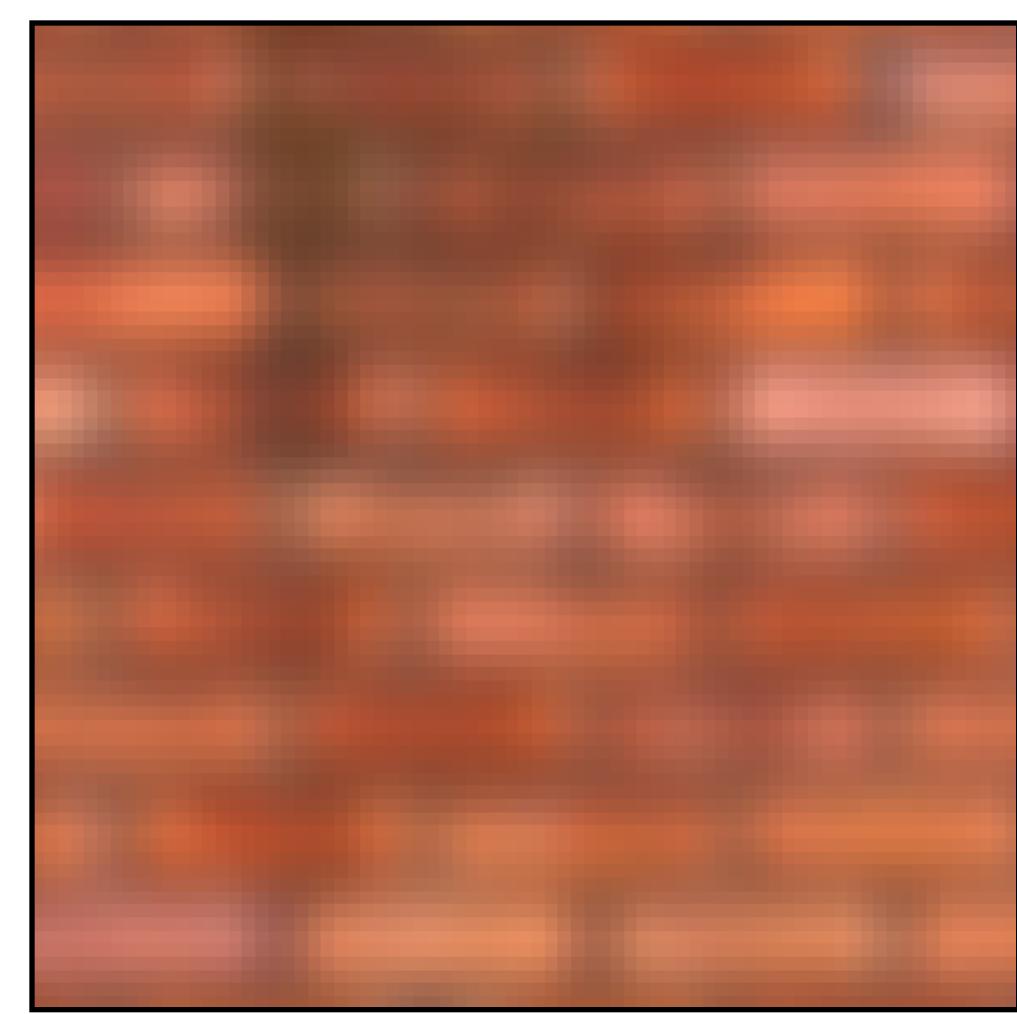
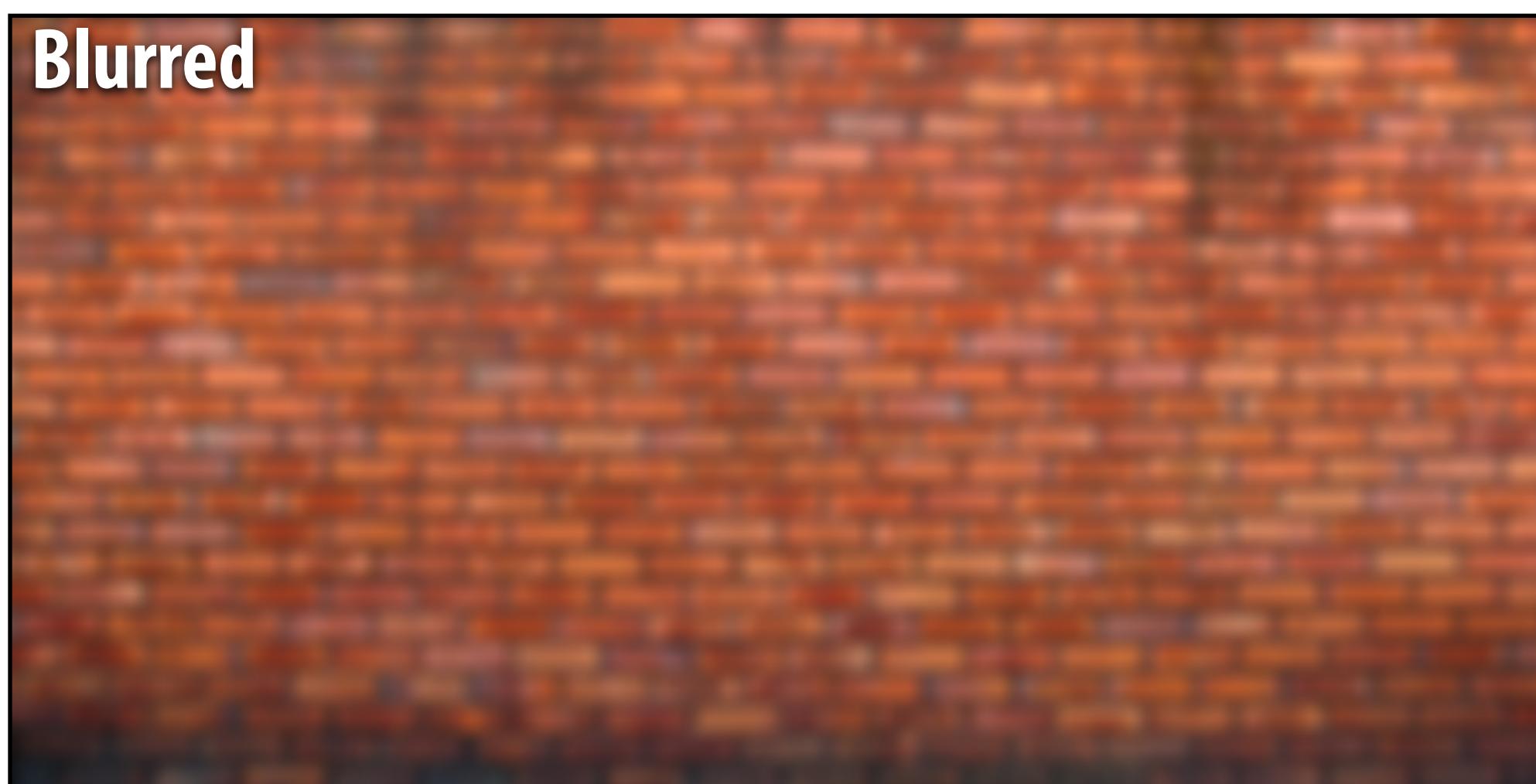
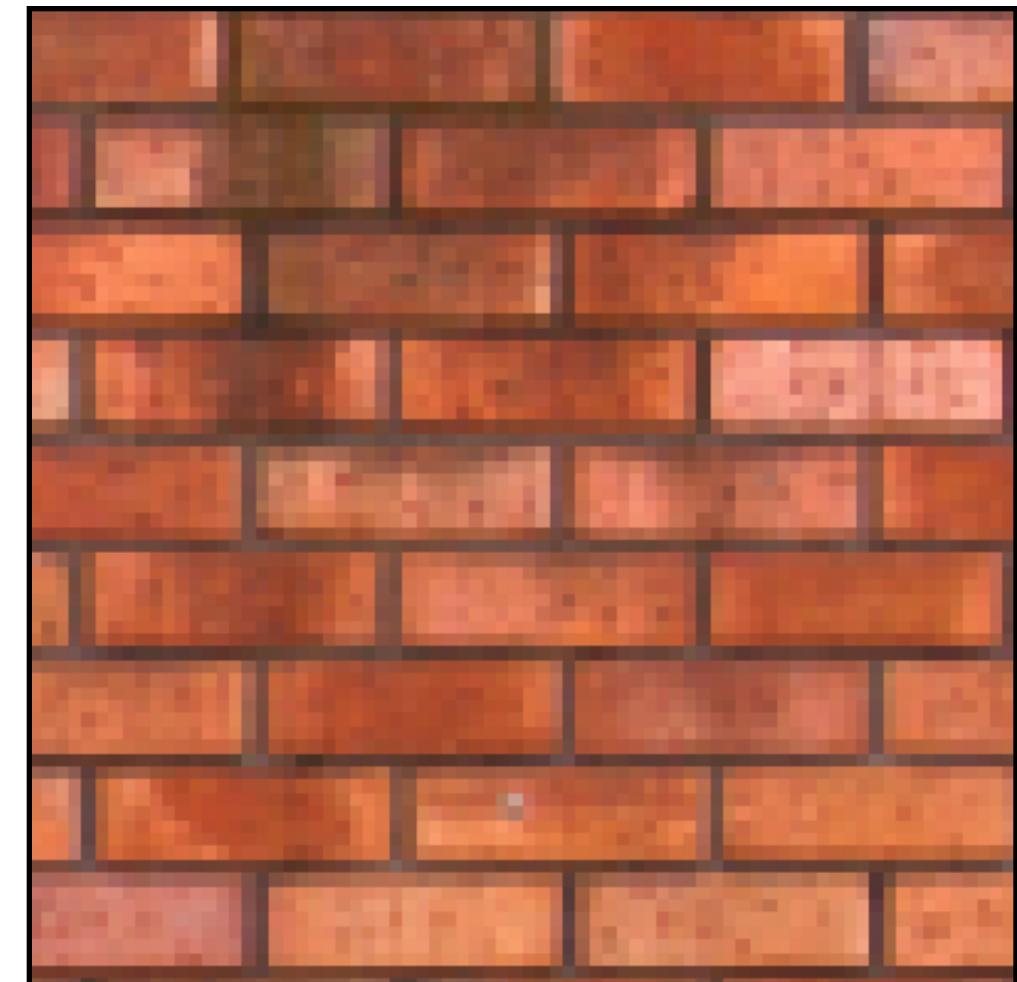
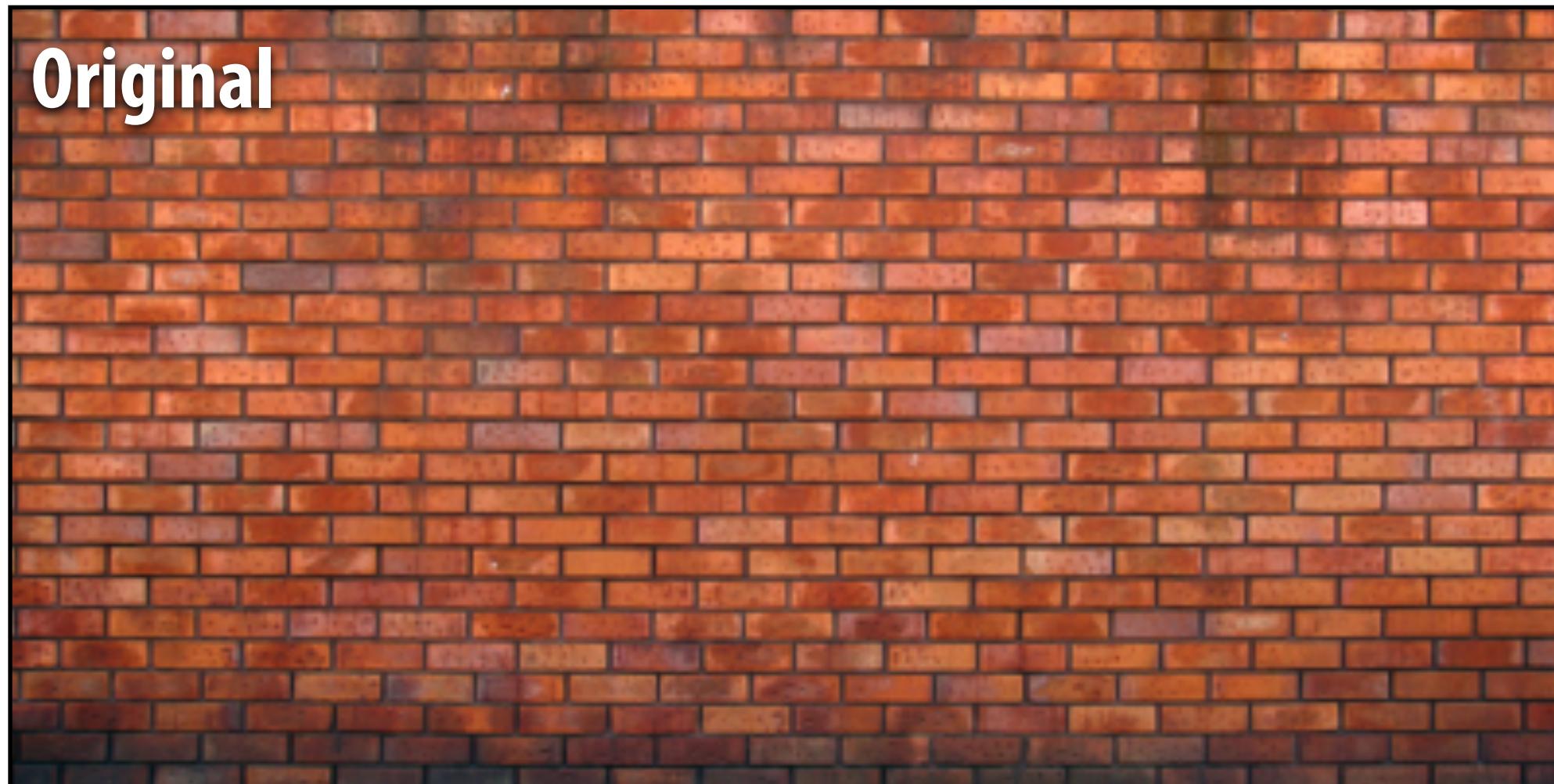


Convolutional layer with stride 2

What does convolution using these filter

$$\begin{bmatrix} .075 & .124 & .075 \\ .124 & .204 & .124 \\ .075 & .124 & .075 \end{bmatrix}$$

“Gaussian Blur”



What does convolution with these filters do?

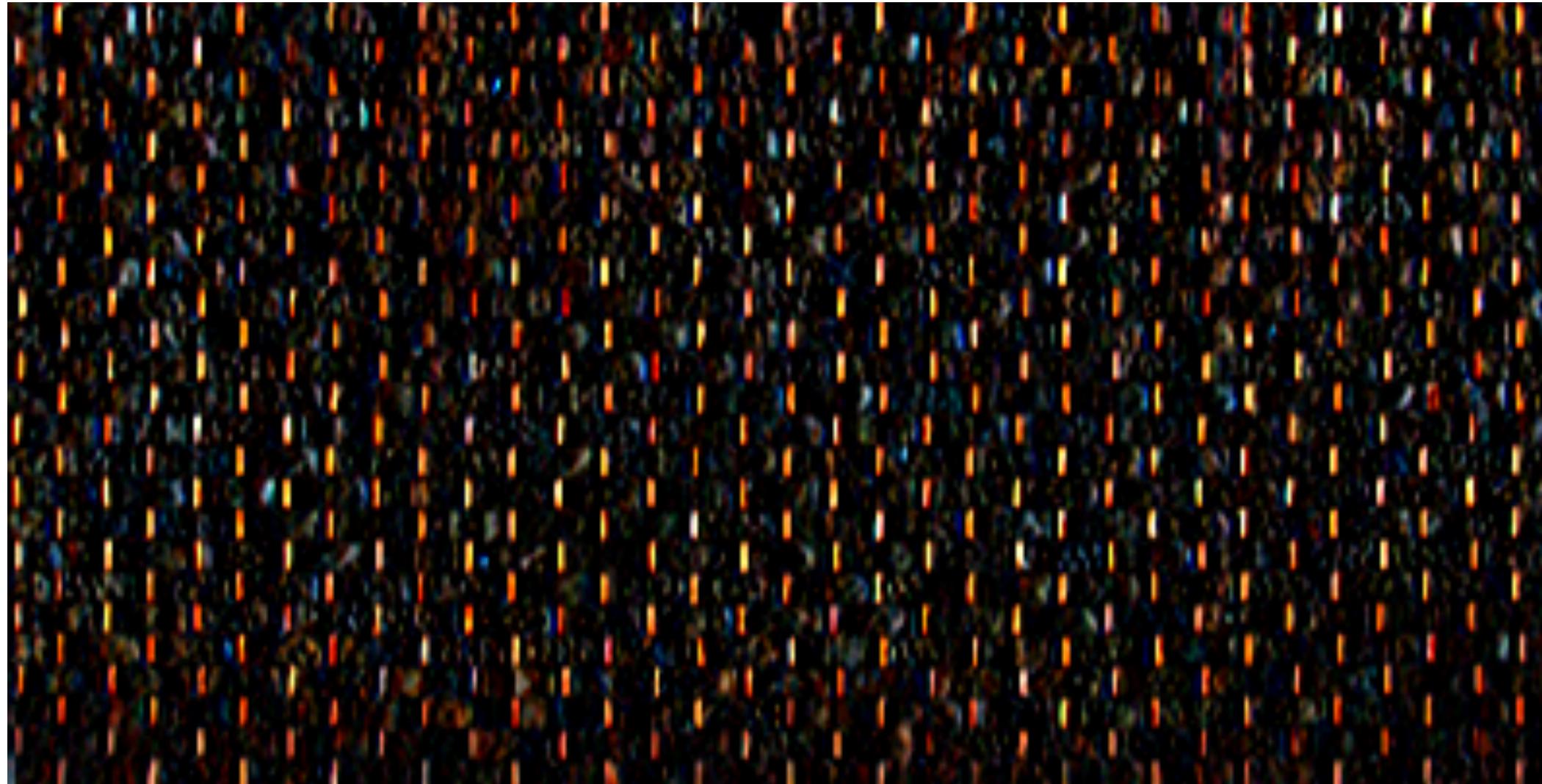
$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

**Extracts horizontal
gradients**

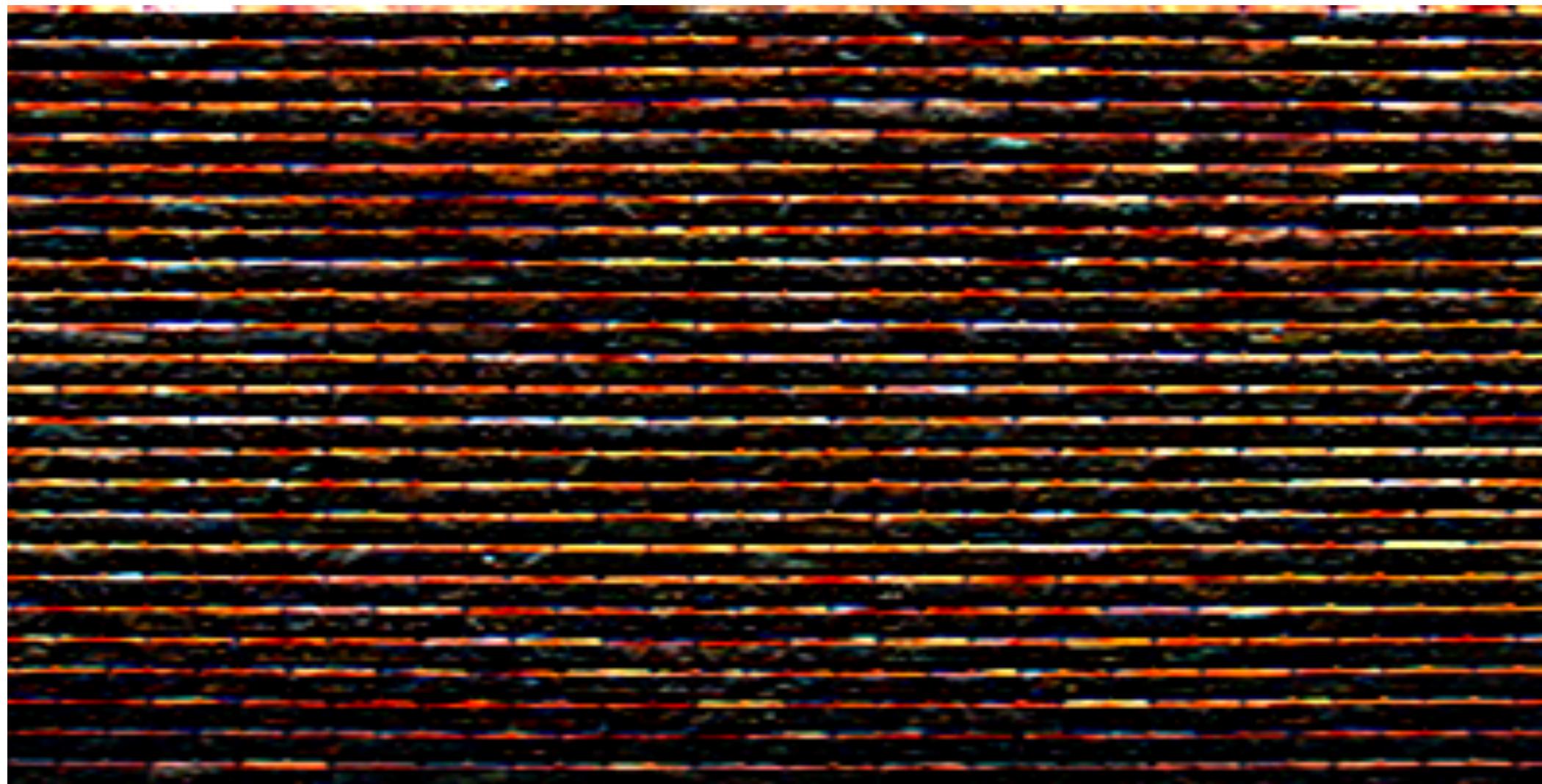
$$\begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$$

**Extracts vertical
gradients**

Gradient detection filters



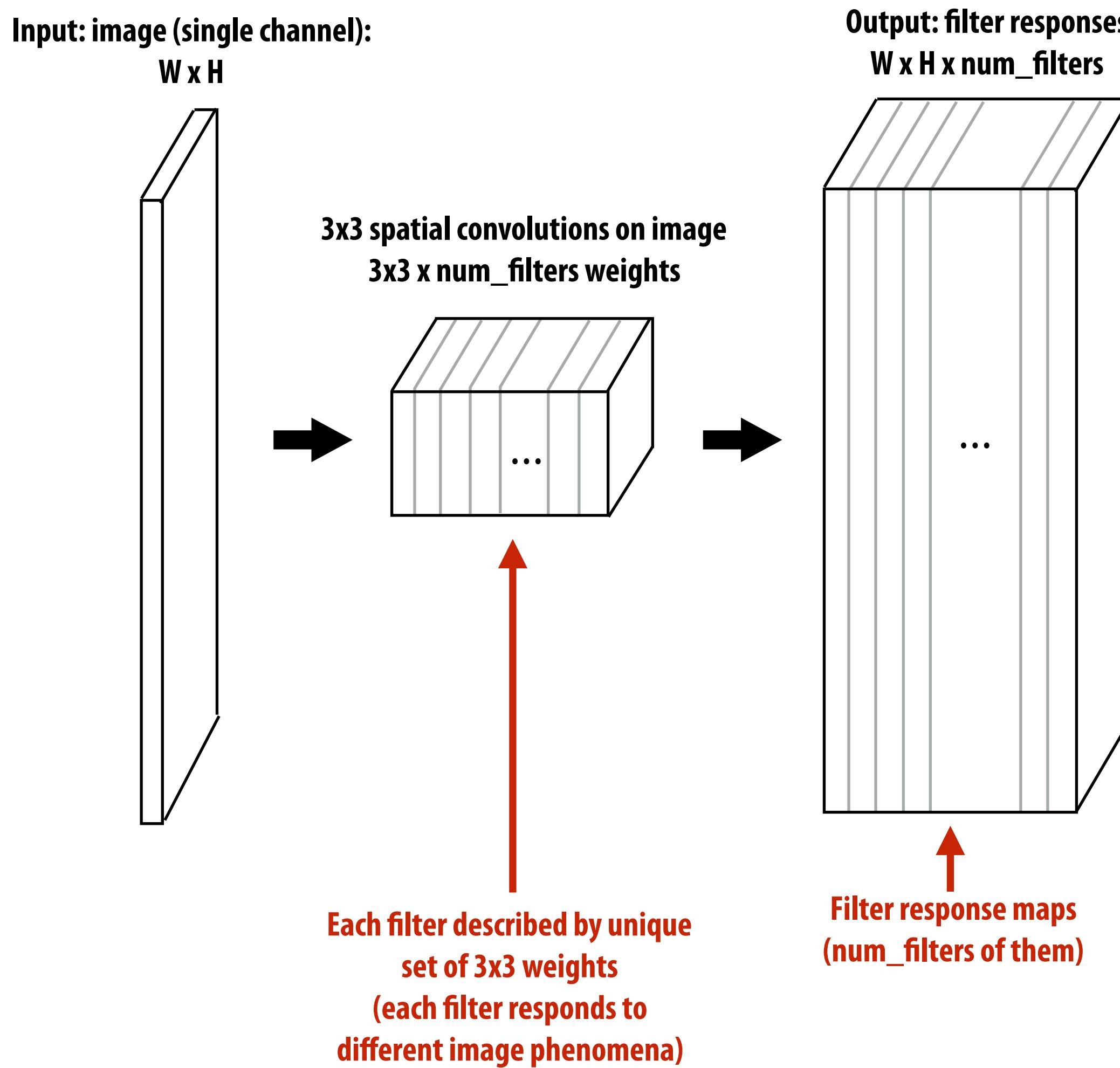
Horizontal gradients



Vertical gradients

Note: you can think of a filter as a “detector” of a pattern, and the magnitude of a pixel in the output image as the “response” of the filter to the region surrounding each pixel in the input image

Applying many filters to an image at once

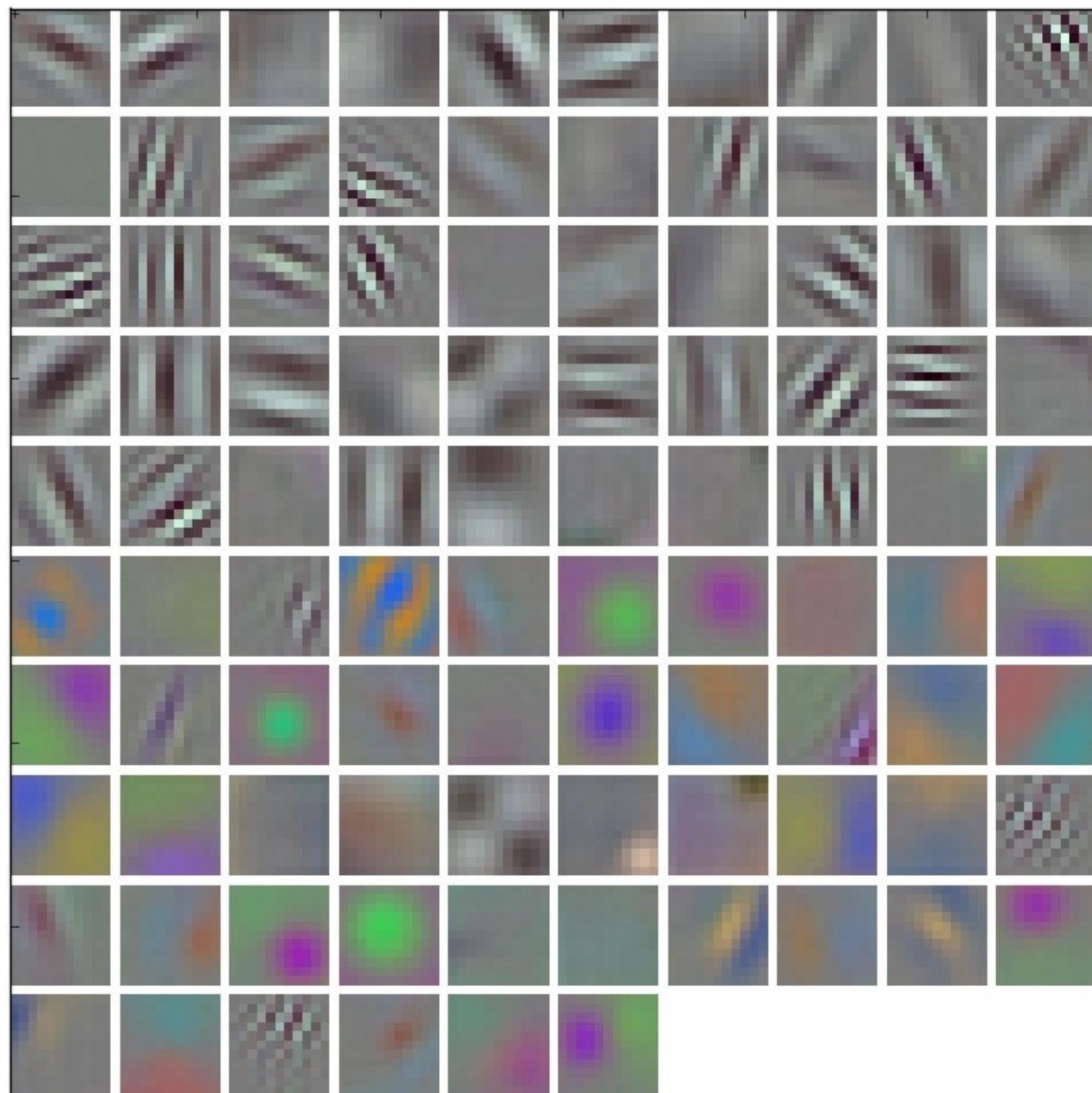


Applying many filters to an image at once

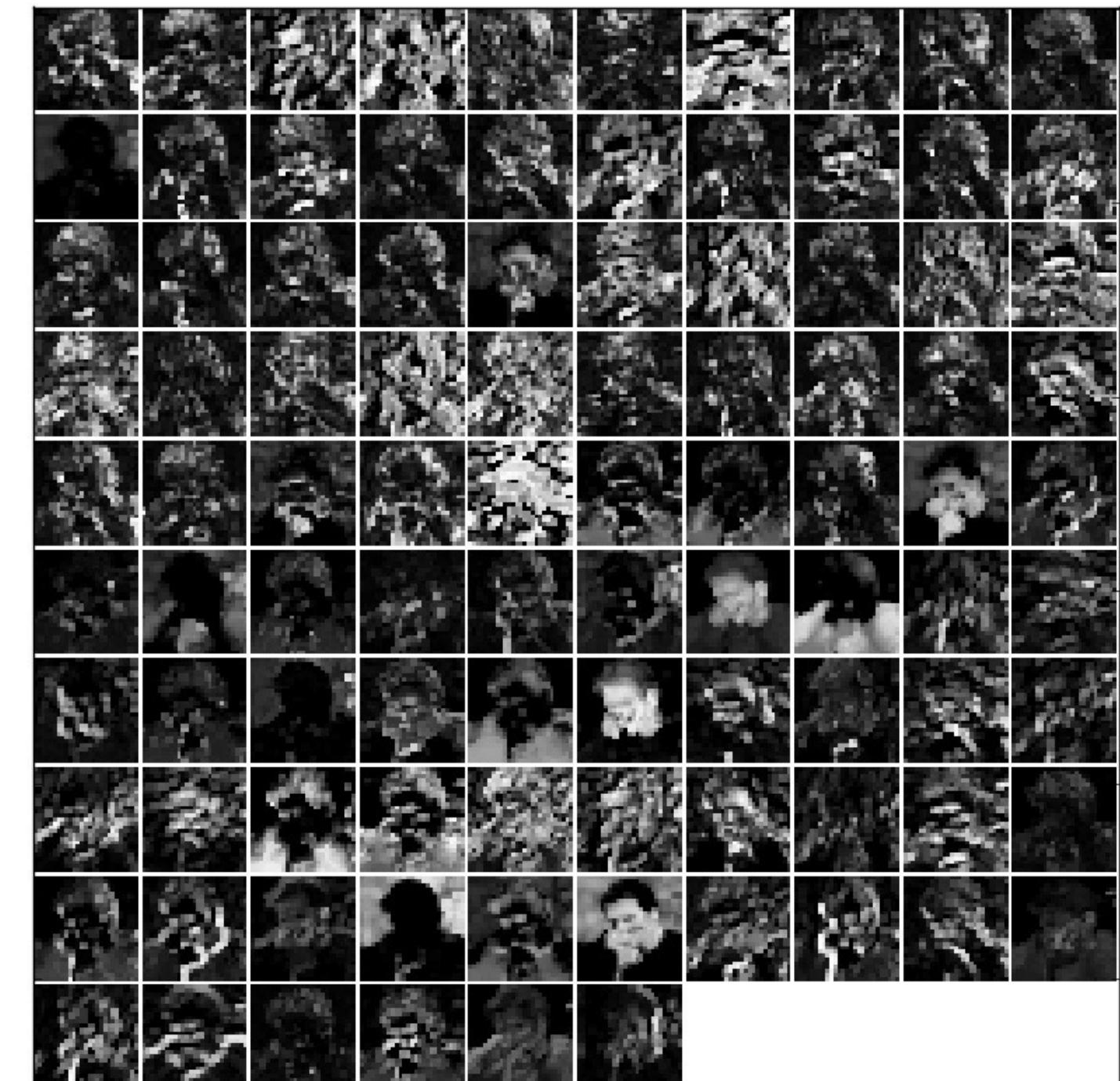
Input RGB image ($W \times H \times 3$)



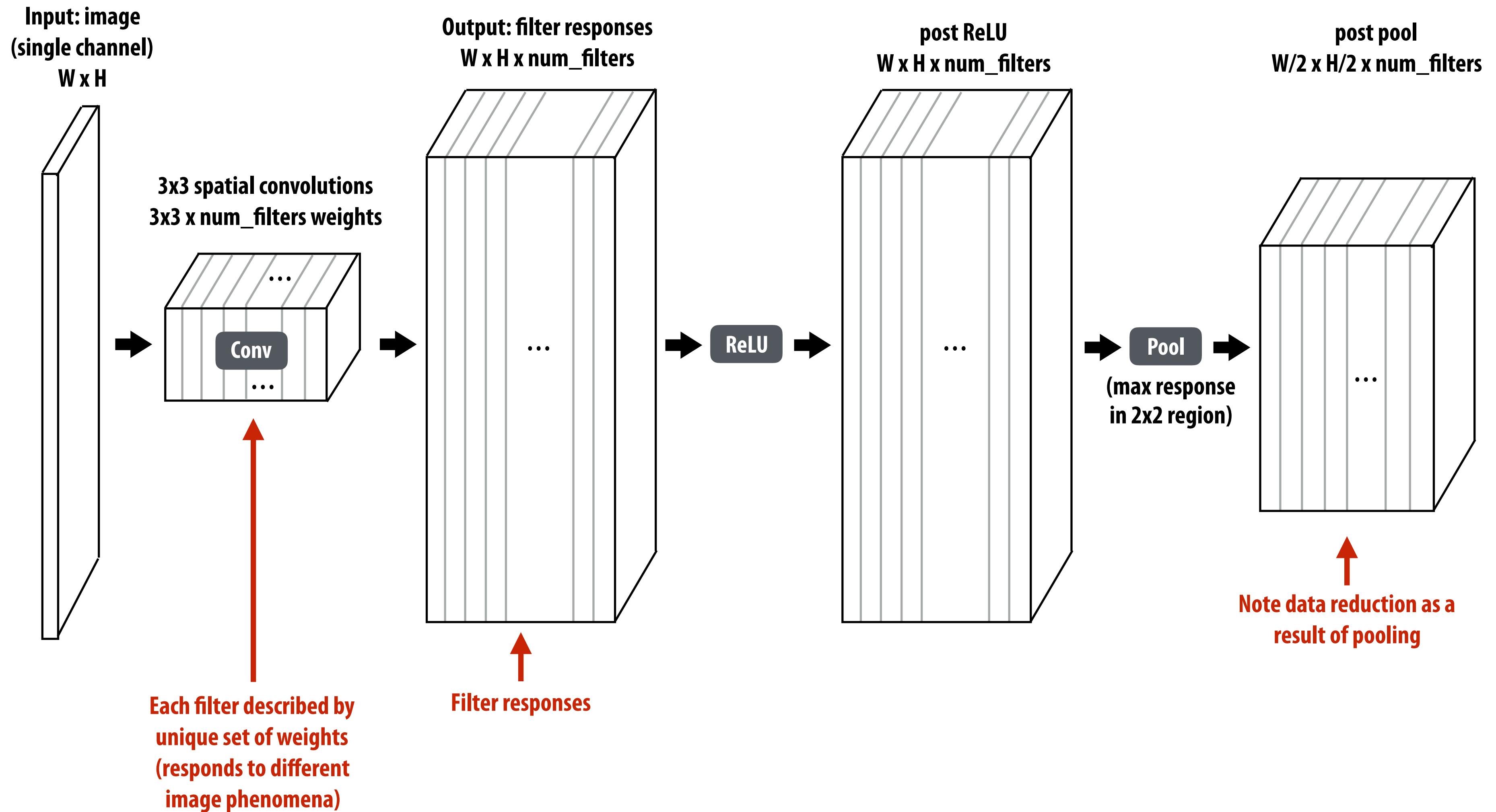
96 11x11x3 filters
(operate on RGB)



96 responses (normalized)



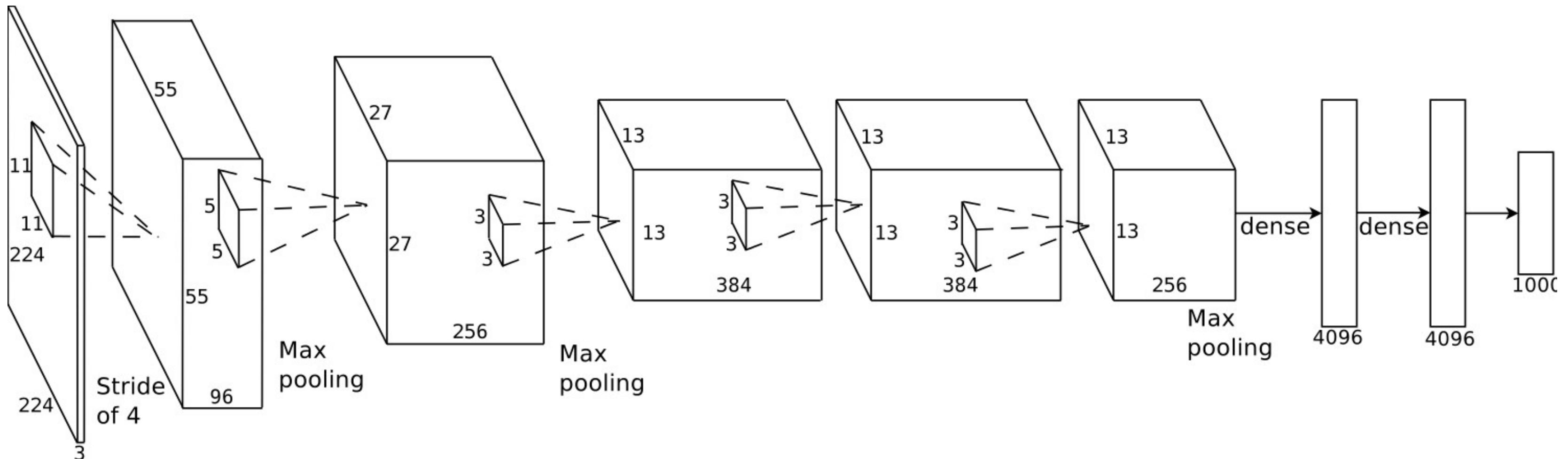
Adding additional layers



Example: “AlexNet” object detection network

Sequences of conv + ReLU + pool (optional) layers

Example: AlexNet [Krizhevsky12]: 5 convolutional layers + 3 fully connected layers



Another example: VGG-16 [Simonyan15]: 13 convolutional layers

input: 224 x 224 RGB

conv/ReLU: 3x3x3x64

conv/ReLU: 3x3x64x64

maxpool

conv/ReLU: 3x3x64x128

conv/ReLU: 3x3x128x128

maxpool

conv/ReLU: 3x3x128x256

conv/ReLU: 3x3x256x256

conv/ReLU: 3x3x256x256

maxpool

conv/ReLU: 3x3x256x512

conv/ReLU: 3x3x512x512

conv/ReLU: 3x3x512x512

maxpool

conv/ReLU: 3x3x512x512

conv/ReLU: 3x3x512x512

conv/ReLU: 3x3x512x512

maxpool

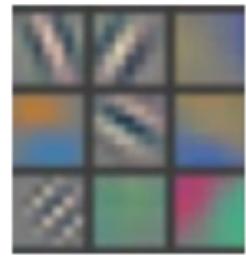
fully-connected 4096

fully-connected 4096

fully-connected 1000

soft-max

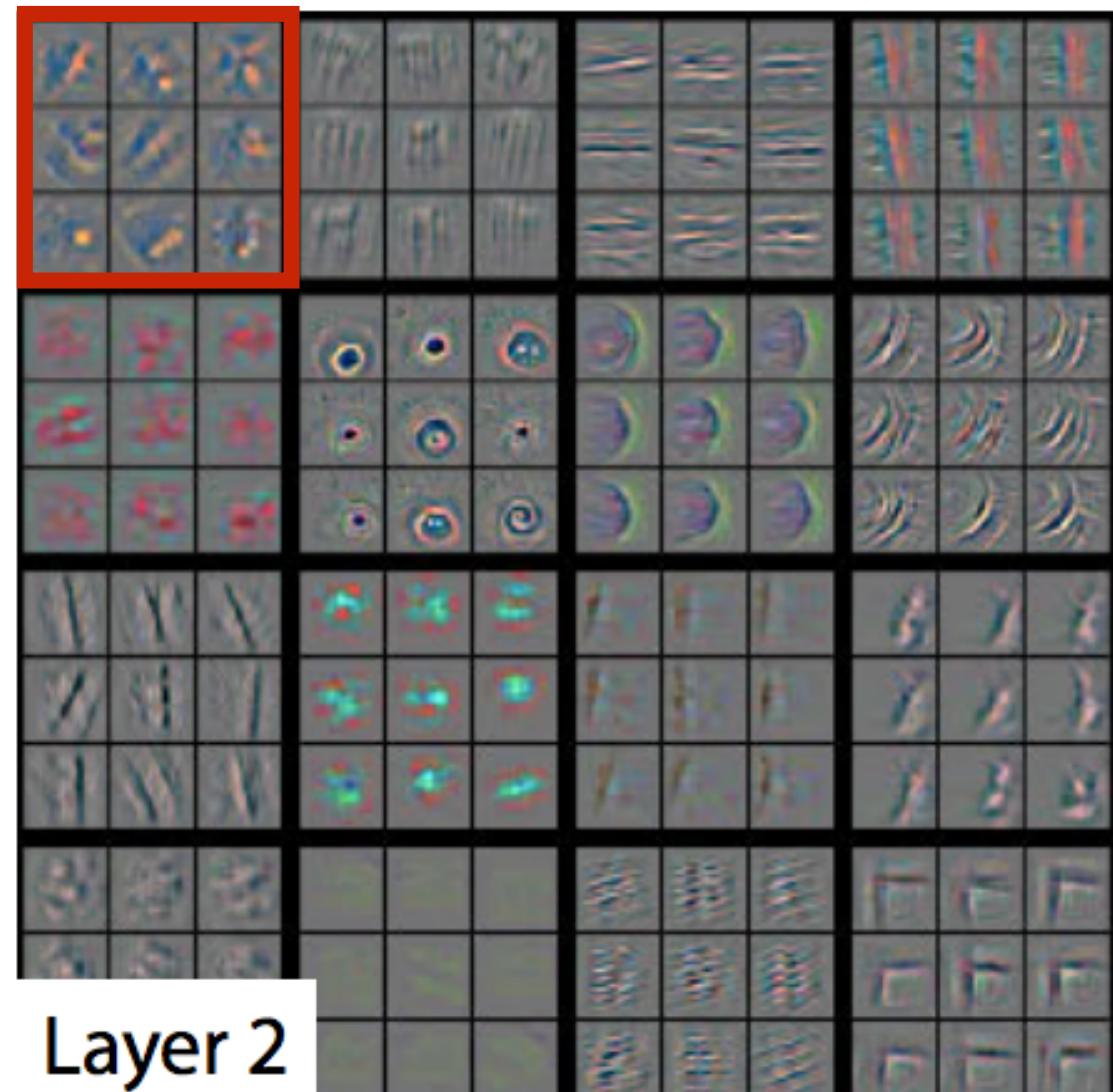
Why deep?



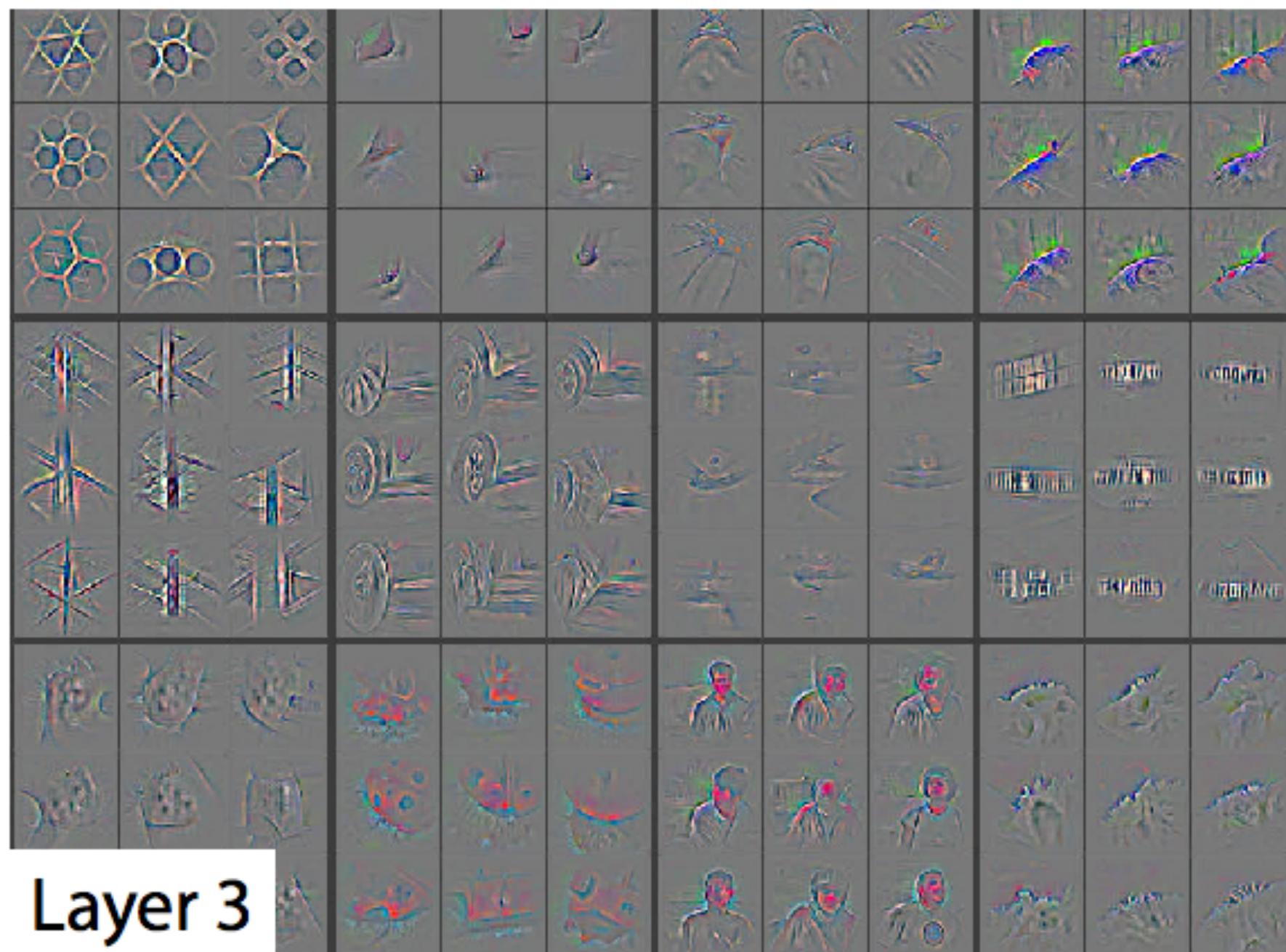
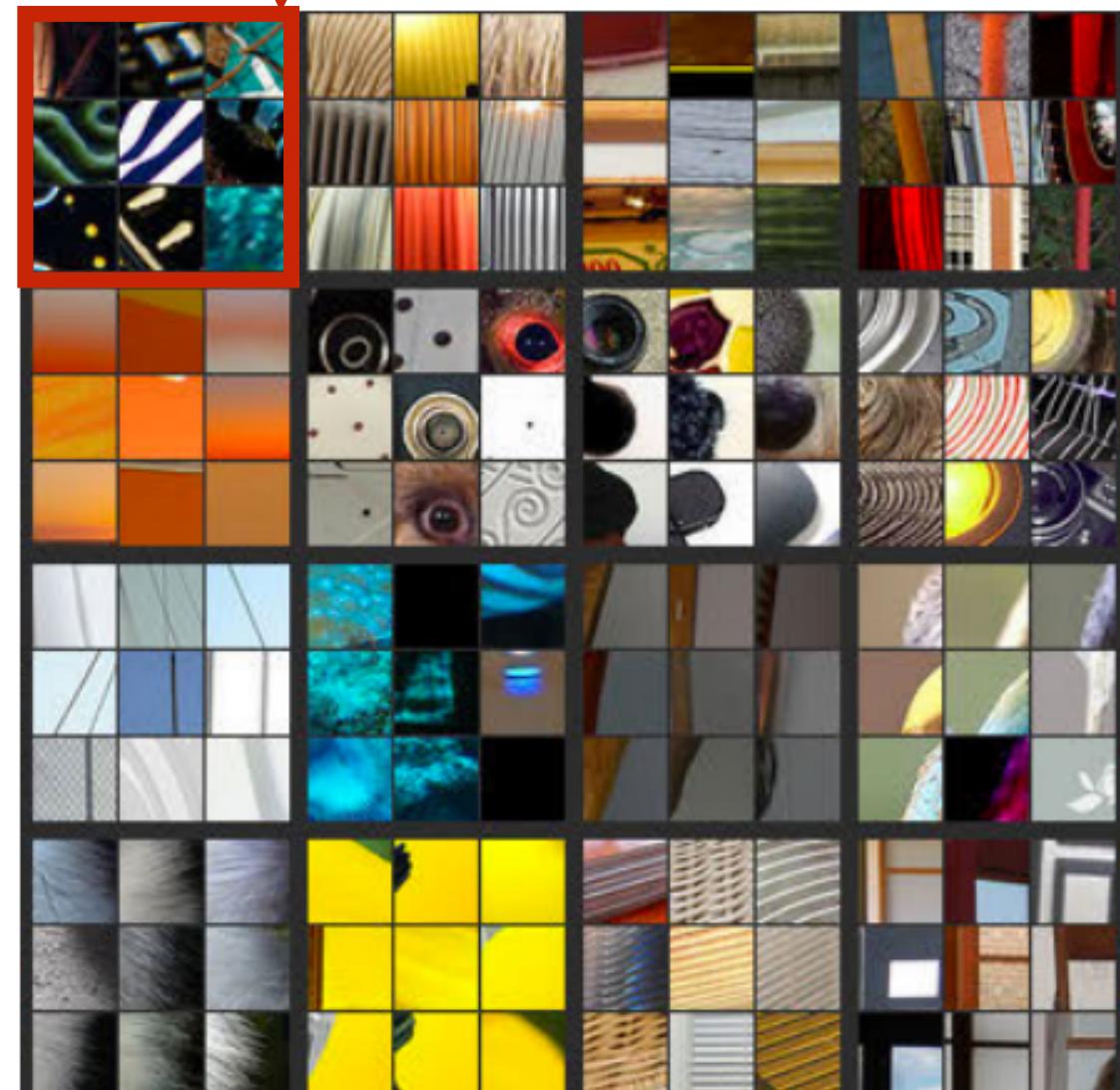
Layer 1



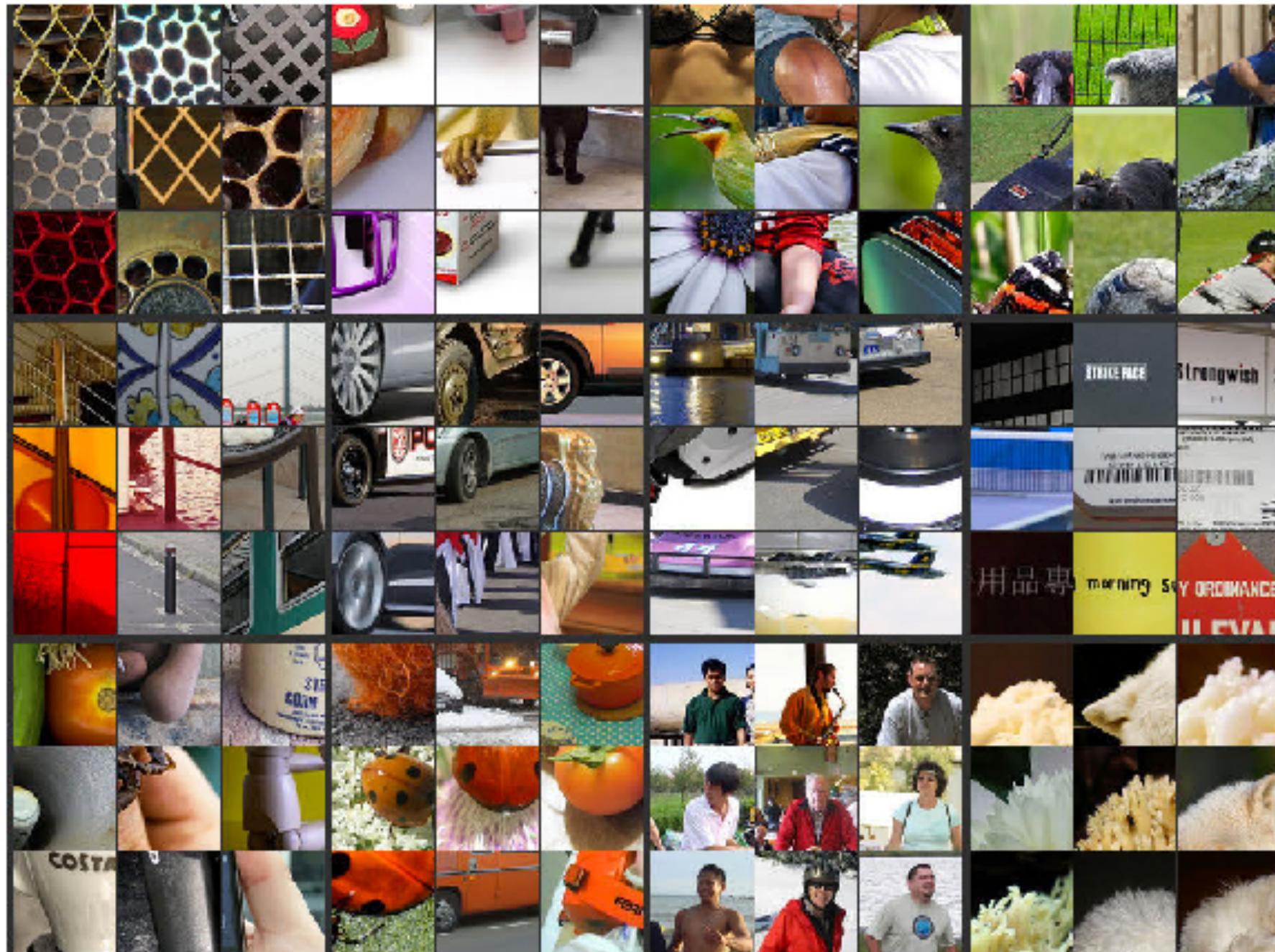
Left: what pixels trigger the response
Right: images that generate strongest response for filters at each layer



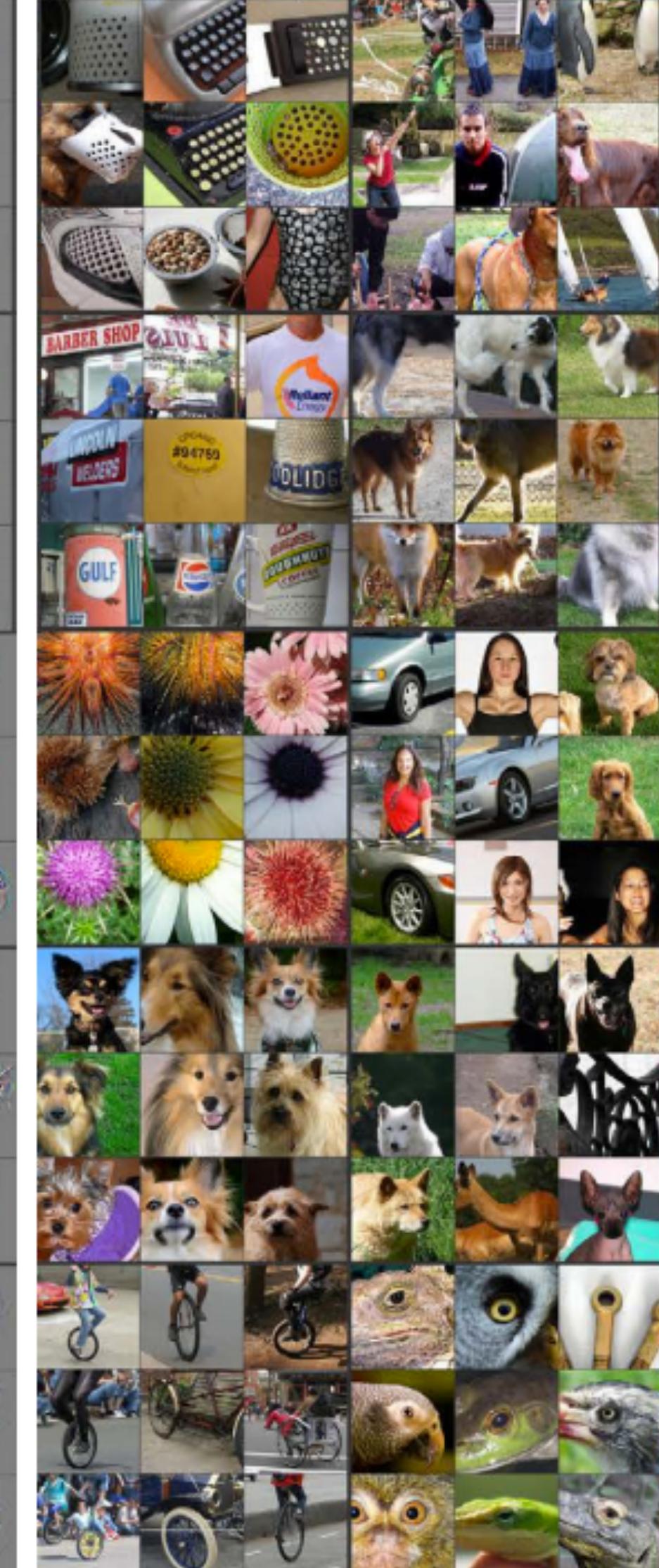
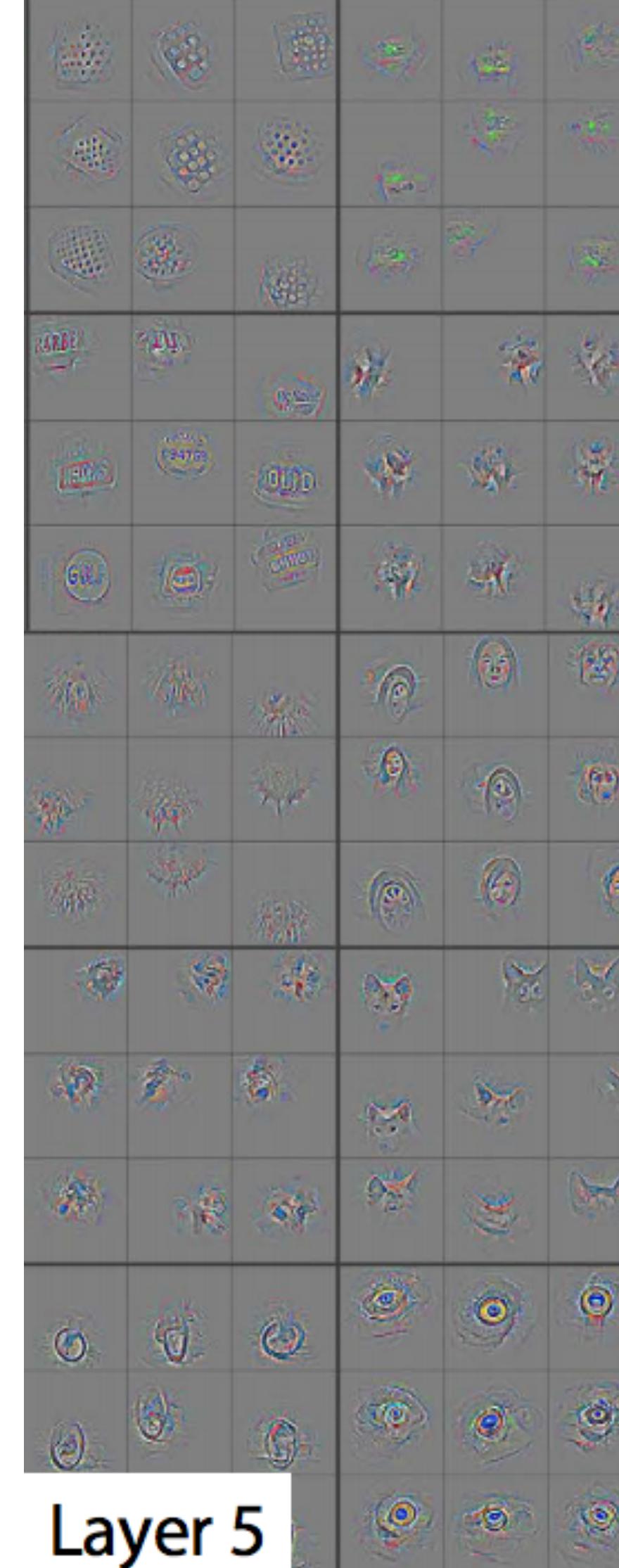
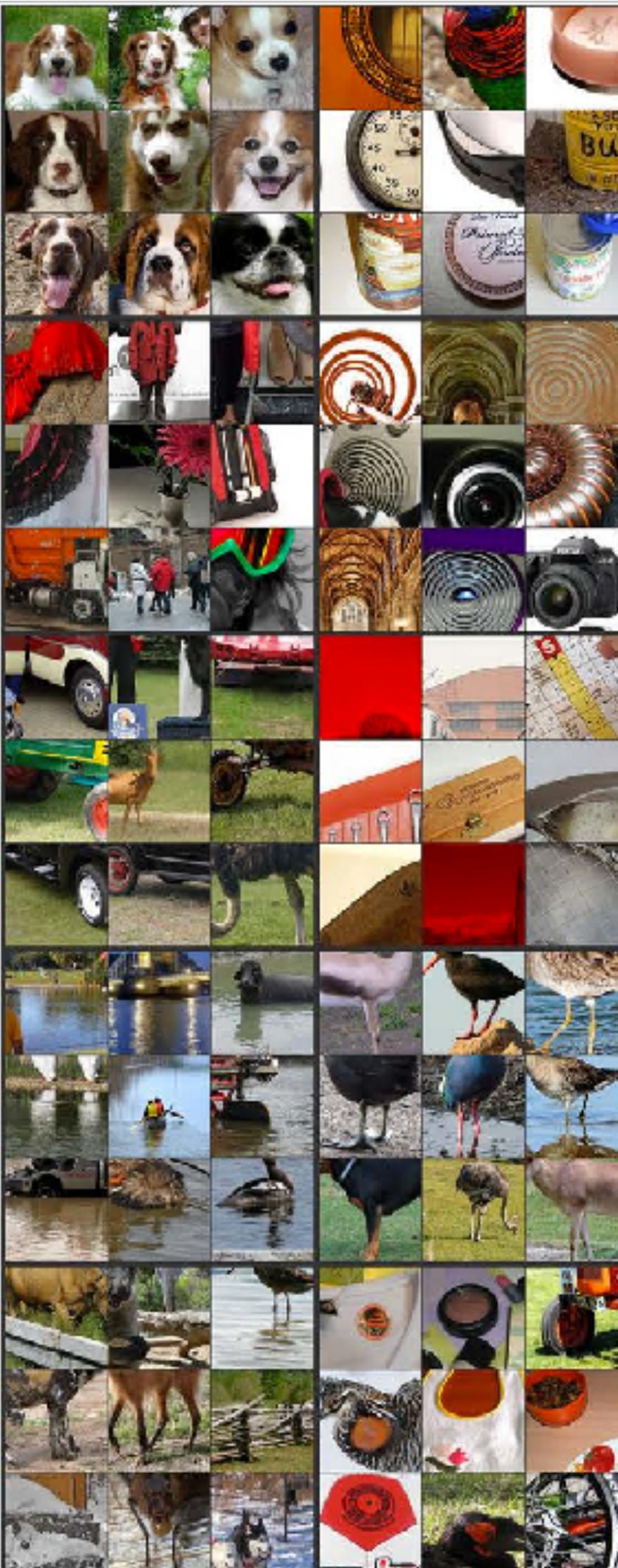
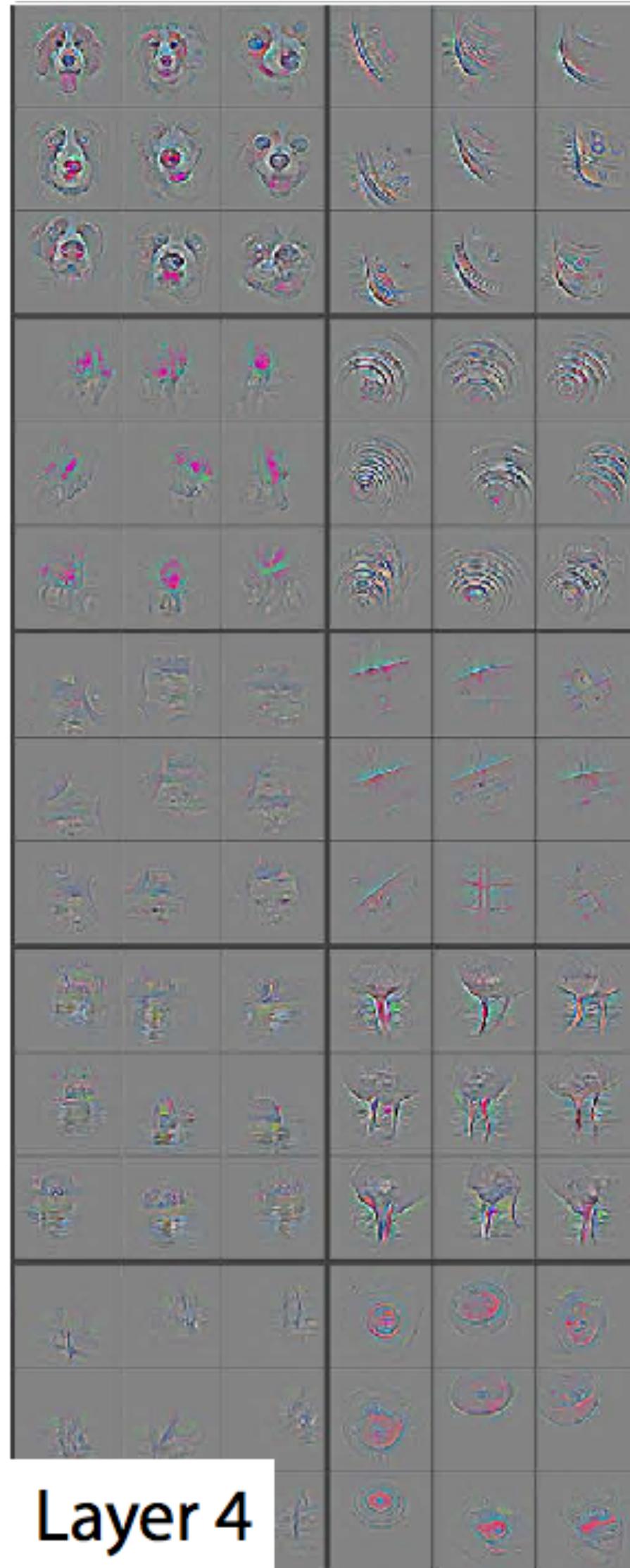
Layer 2



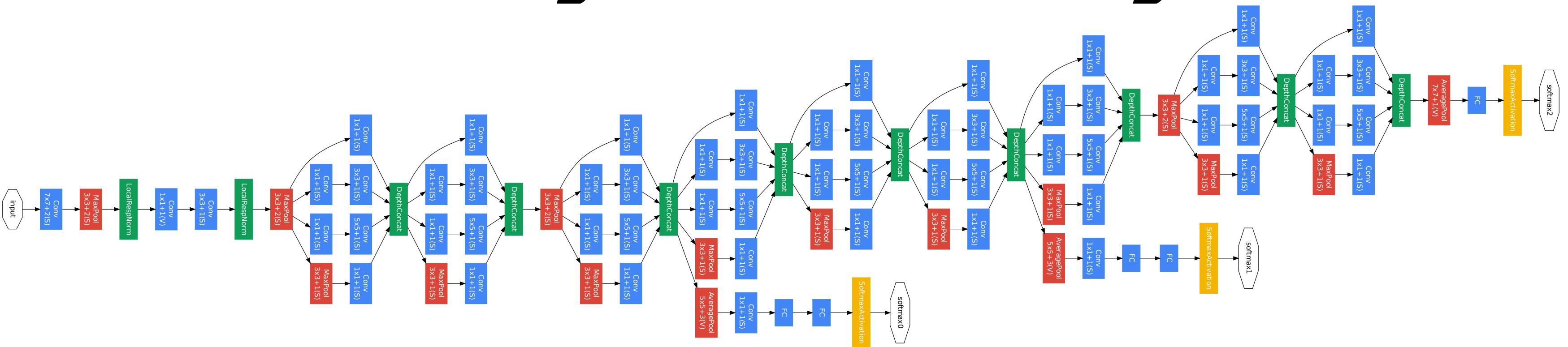
Layer 3



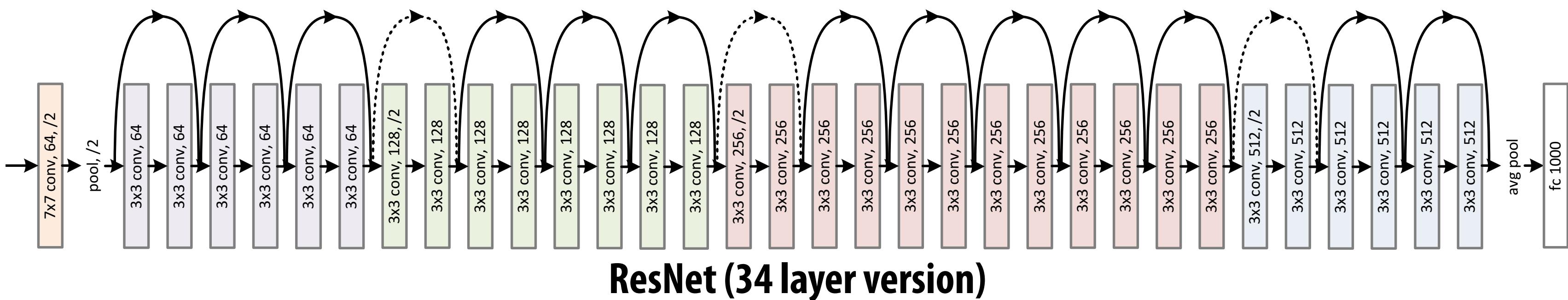
Why deep?



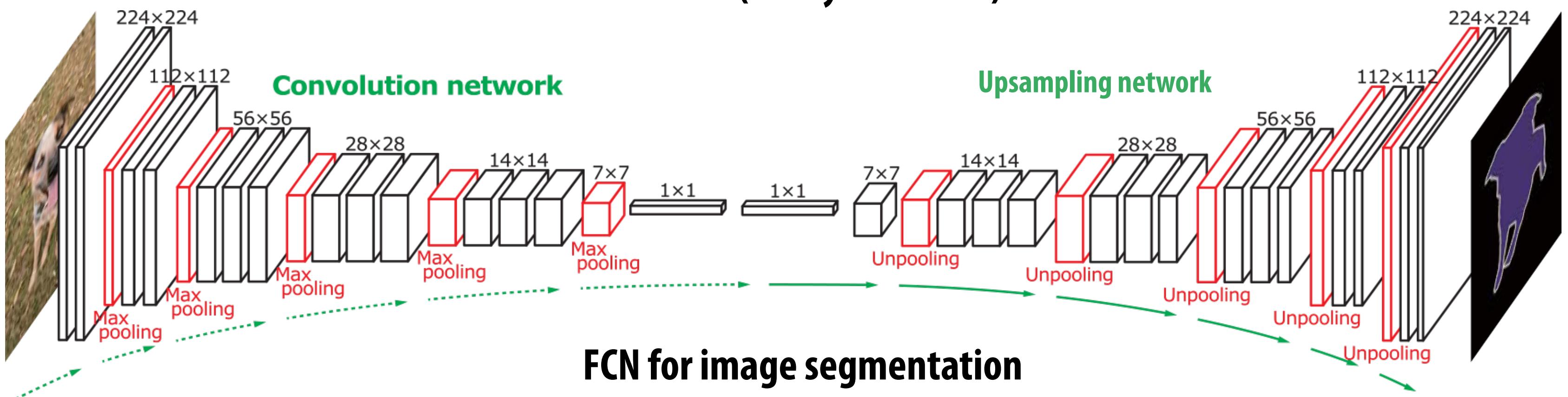
More recent image understanding networks



Inception (GoogleLeNet)



ResNet (34 layer version)



Deep networks learn useful representations

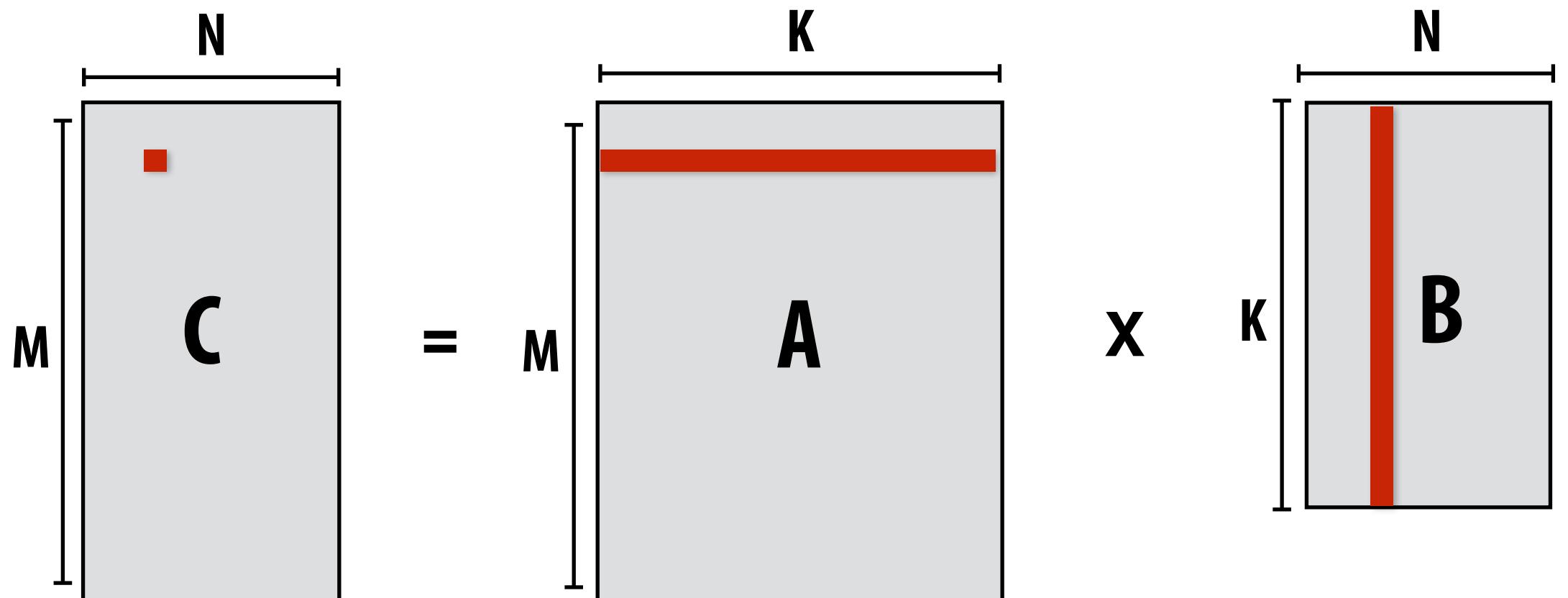
- **Simultaneous, multi-scale learning of useful concepts for task at hand**
 - Example on previous slides: subparts detectors emerged in network for object classification
- **“Fine tuning”**
 - Common practice of adjusting previous representations (weights) learned on task A, for new (but similar) task B
 - Reduces training time for task B
 - Useful when little training data for task B exists
 - Example: keep first N layers the same, tune weights in highest layers based on task B training data

Efficiently implementing convolution layers

Dense matrix multiplication

```
float A[M][K];
float B[K][N];
float C[M][N];

// compute C += A * B
#pragma omp parallel for
for (int j=0; j<M; j++)
    for (int i=0; i<N; i++)
        for (int k=0; k<K; k++)
            C[j][i] += A[j][k] * B[k][i];
```



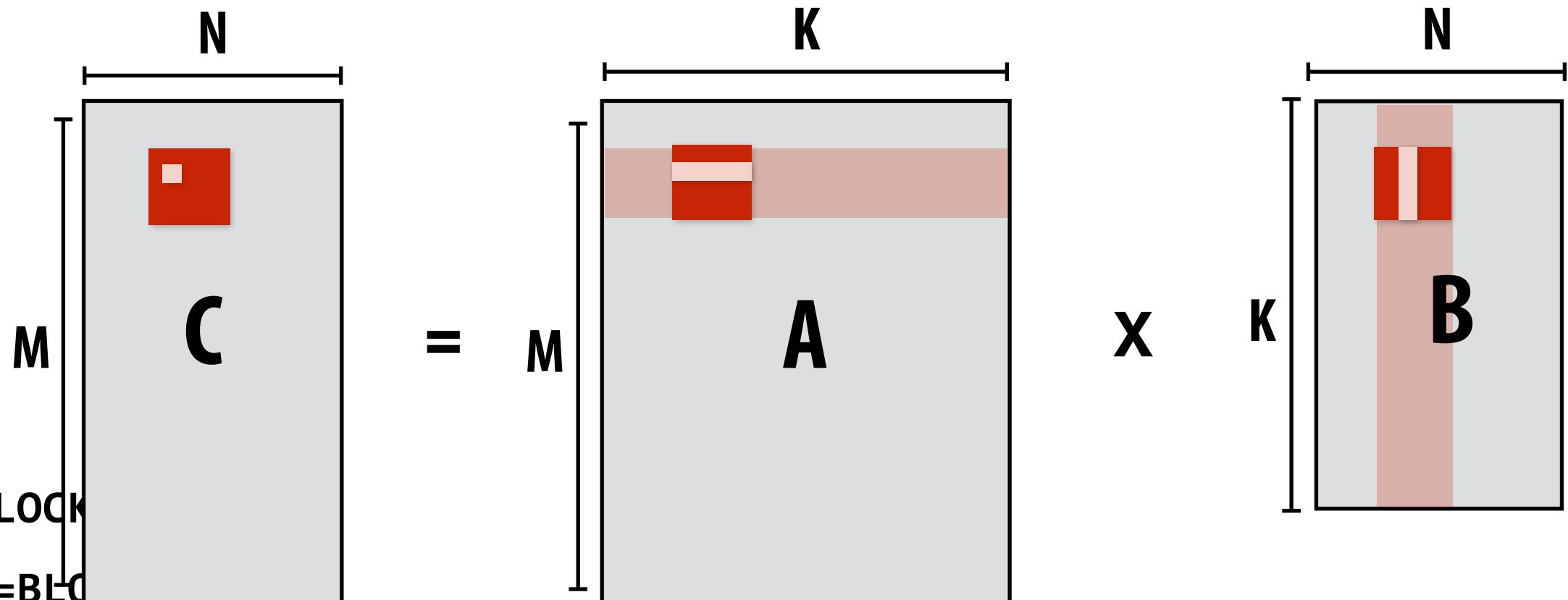
What is the problem with this implementation?

Low arithmetic intensity (does not exploit temporal locality in access to A and B)

Blocked dense matrix multiplication

```
float A[M][K];
float B[K][N];
float C[M][N];

// compute C += A * B
#pragma omp parallel for
for (int jblock=0; jblock<M; jblock+=BLOCKSIZE_M)
    for (int iblock=0; iblock<N; iblock+=BLOCKSIZE_N)
        for (int kblock=0; kblock<K; kblock+=BLOCKSIZE_K)
            for (int j=0; j<BLOCKSIZE_J; j++)
                for (int i=0; i<BLOCKSIZE_I; i++)
                    for (int k=0; k<BLOCKSIZE_K; k++)
                        C[jblock+j][iblock+i] += A[jblock+j][kblock+k] * B[kblock+k][iblock+i];
```



Idea: compute partial result for block of C while required blocks of A and B remain in cache
(Assumes BLOCKSIZE chosen to allow block of A, B, and C to remain resident)

Self check: do you want as big a BLOCKSIZE as possible? Why?

Hierarchical blocked matrix mult

Exploit multi-level of memory hierarchy

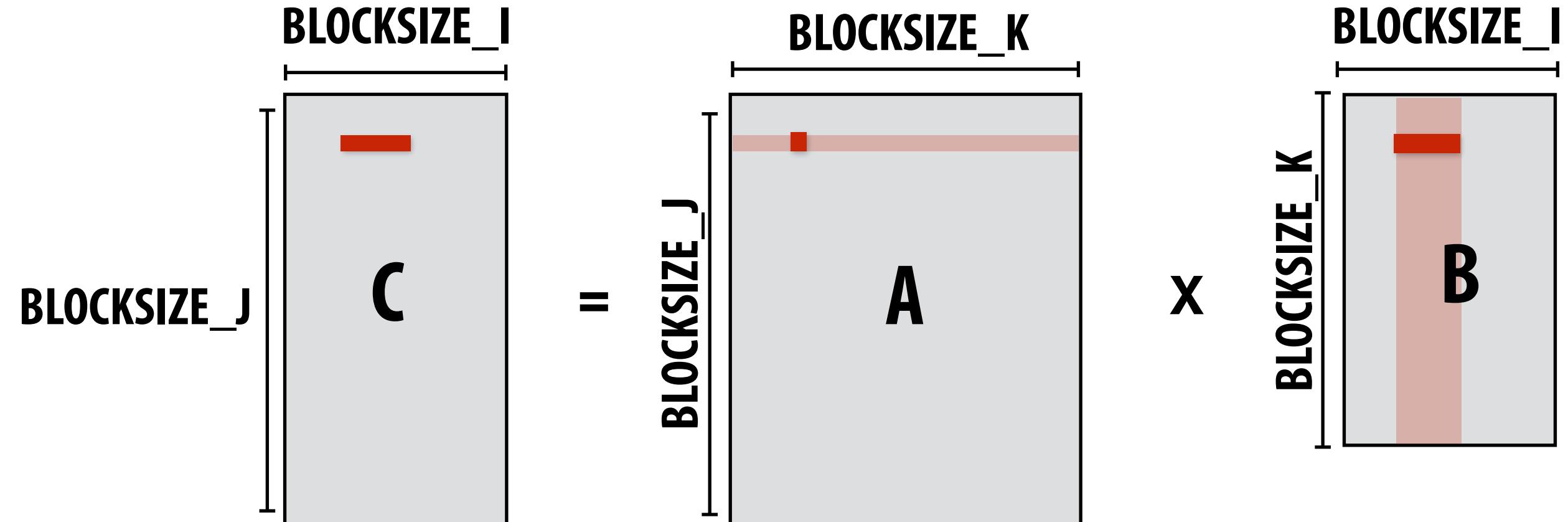
```
float A[M][K];
float B[K][N];
float C[M][N];

// compute C += A * B
#pragma omp parallel for
for (int jblock2=0; jblock2<M; jblock2+=L2_BLOCKSIZE_J)
    for (int iblock2=0; iblock2<N; iblock2+=L2_BLOCKSIZE_I)
        for (int kblock2=0; kblock2<K; kblock2+=L2_BLOCKSIZE_K)
            for (int jblock1=0; jblock1<L1_BLOCKSIZE_J; jblock1+=L1_BLOCKSIZE_J)
                for (int iblock1=0; iblock1<L1_BLOCKSIZE_I; iblock1+=L1_BLOCKSIZE_I)
                    for (int kblock1=0; kblock1<L1_BLOCKSIZE_K; kblock1+=L1_BLOCKSIZE_K)
                        for (int j=0; j<BLOCKSIZE_J; j++)
                            for (int i=0; i<BLOCKSIZE_I; i++)
                                for (int k=0; k<BLOCKSIZE_K; k++)
...
...
```

Not shown: final level of “blocking” for register locality...

Blocked dense matrix multiplication (1)

Consider SIMD parallelism
within a block



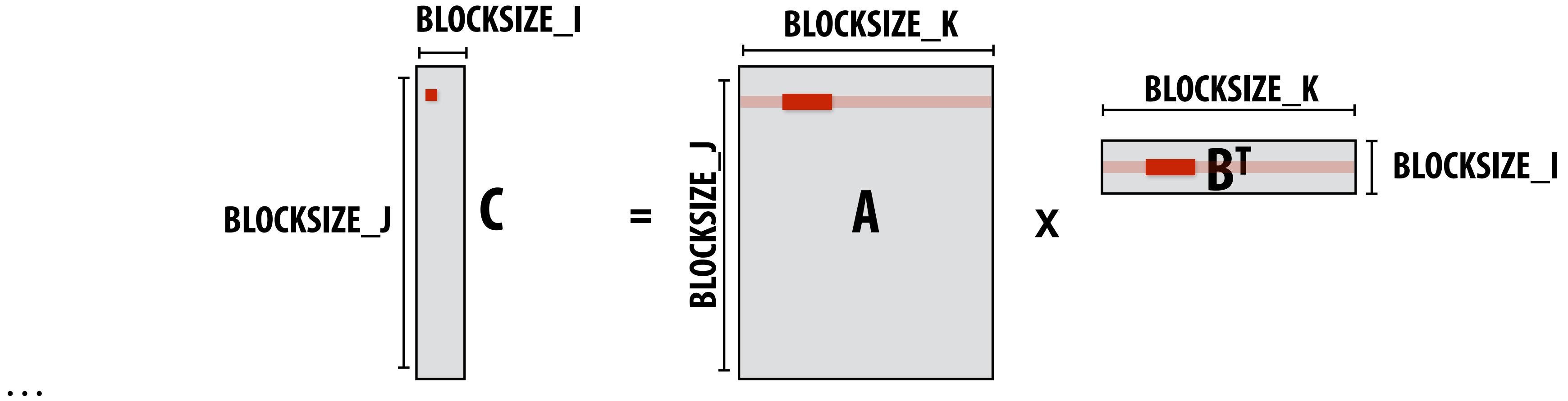
```
...
for (int j=0; j<BLOCKSIZE_J; j++) {
    for (int i=0; i<BLOCKSIZE_I; i+=SIMD_WIDTH) {
        SIMD_vec C_accum = vec_load(&C[jblock+j][iblock+i]);
        for (int k=0; k<BLOCKSIZE_K; k++) {
            // C = A*B + C
            SIMD_vec A_val = splat(&A[jblock+j][kblock+k]); // load a single element in vector register
            SIMD_muladd(A_val, vec_load(&B[kblock+k][iblock+i]), C_accum);
        }
        vec_store(&C[jblock+j][iblock+i], C_accum);
    }
}
```

Vectorize i loop

Good: also improves spatial locality in access to B

Bad: working set increased by SIMD_WIDTH, still walking over B in large steps

Blocked dense matrix multiplication (2)

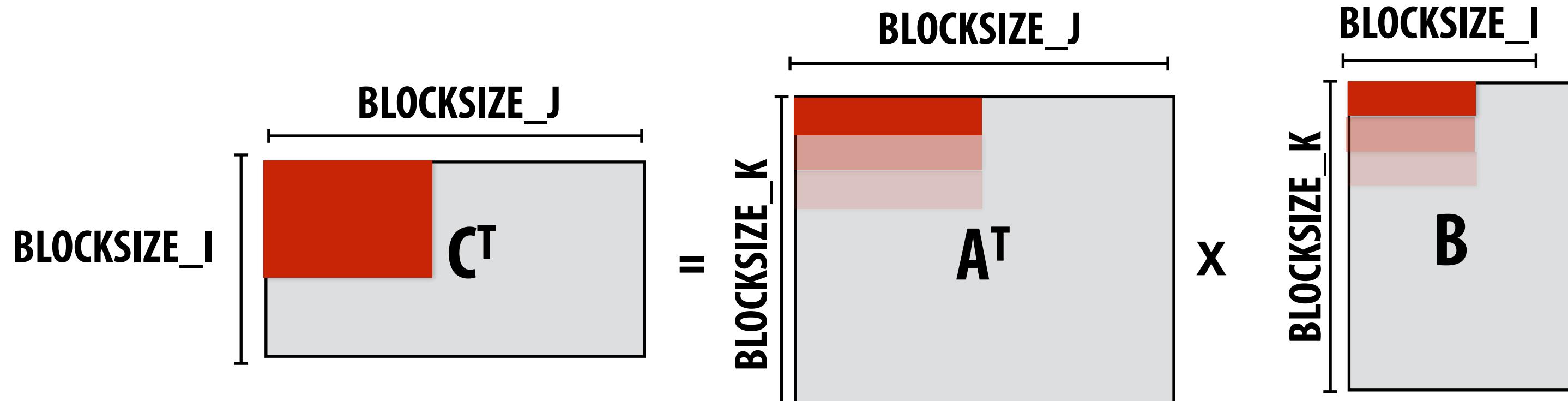


...

```
for (int j=0; j<BLOCKSIZE_J; j++)  
    for (int i=0; i<BLOCKSIZE_I; i++) {  
        float C_scalar = C[jblock+j][iblock+i];  
        for (int k=0; k<BLOCKSIZE_K; k+=SIMD_WIDTH) {  
            // C_scalar = dot(A,B) + C_scalar  
            C_scalar += simd_dot(vec_load(&A[jblock+j][kblock+k]), vec_load(&Btrans[iblock+i][[kblock+k]));  
        }  
        C[jblock+j][iblock+i] = C_scalar;  
    }
```

Assume i dimension is small. Previous vectorization scheme (1) would not work well.
Pre-transpose block of B (copy block of B to temp buffer in transposed form)
Vectorize innermost loop

Blocked dense matrix multiplication (3)



```
// assume blocks of A and C are pre-transposed as Atrans and Ctrans
for (int j=0; j<BLOCKSIZE_J; j+=SIMD_WIDTH) {
    for (int i=0; i<BLOCKSIZE_I; i+=SIMD_WIDTH) {

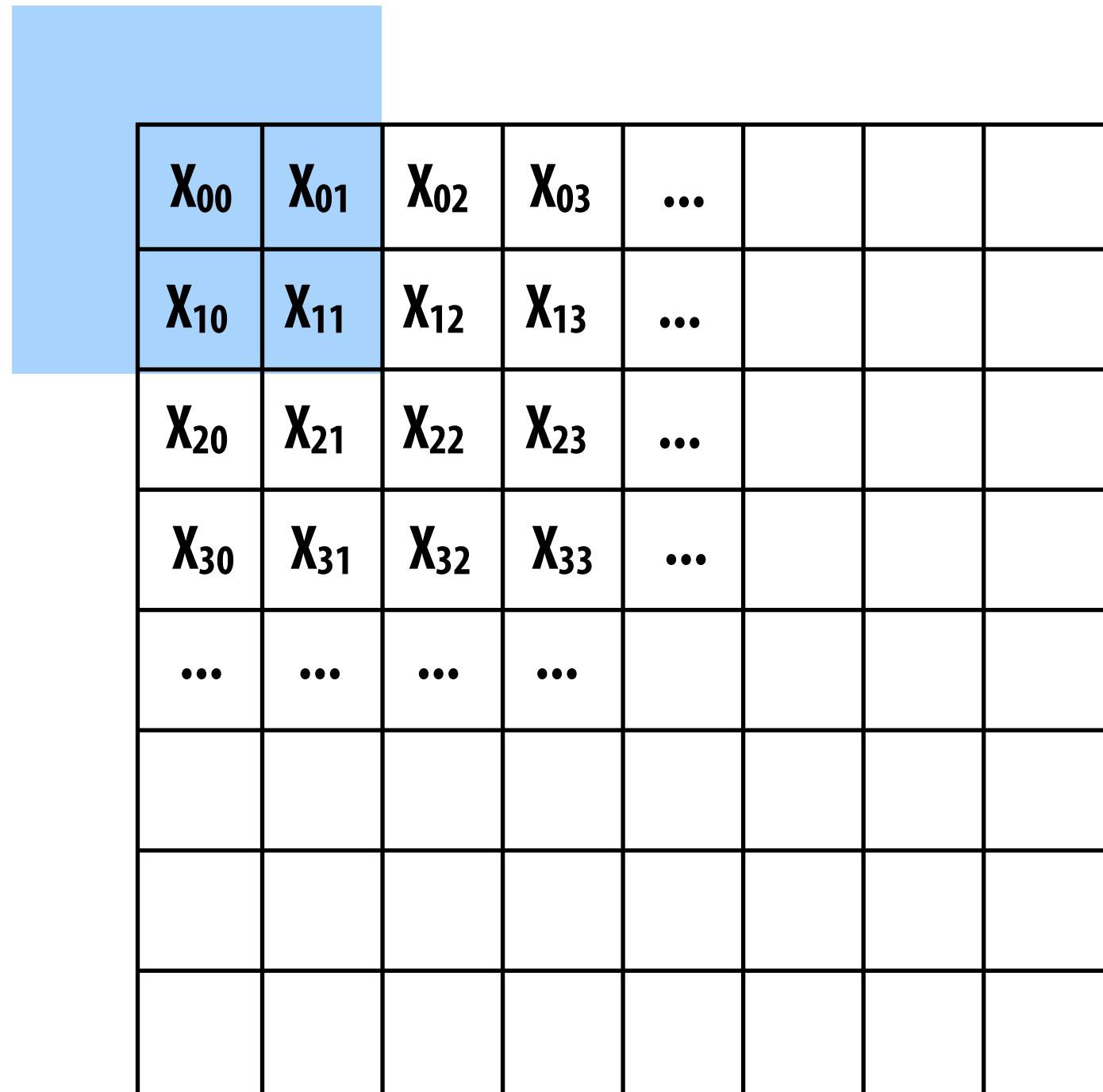
        SIMD_vec C_accum[SIMD_WIDTH];
        for (int k=0; k<SIMD_WIDTH; k++) // load C_accum
            C_accum[k] = vec_load(&Ctrans[iblock+i+k][jblock+j]);

        for (int k=0; k<BLOCKSIZE_K; k++) {
            SIMD_vec bvec = vec_load(&B[kblock+k][iblock+i]);
            for (int kk=0; kk<SIMD_WIDTH; kk++) // innermost loop items not dependent
                SIMD_muladd(vec_load(&Atrans[kblock+k][jblock+j], splat(bvec[kk])), C_accum[kk]);
        }

        for (int k=0; k<SIMD_WIDTH; k++)
            vec_store(&Ctrans[iblock+i+k][jblock+j], C_accum[k]);
    }
}
```

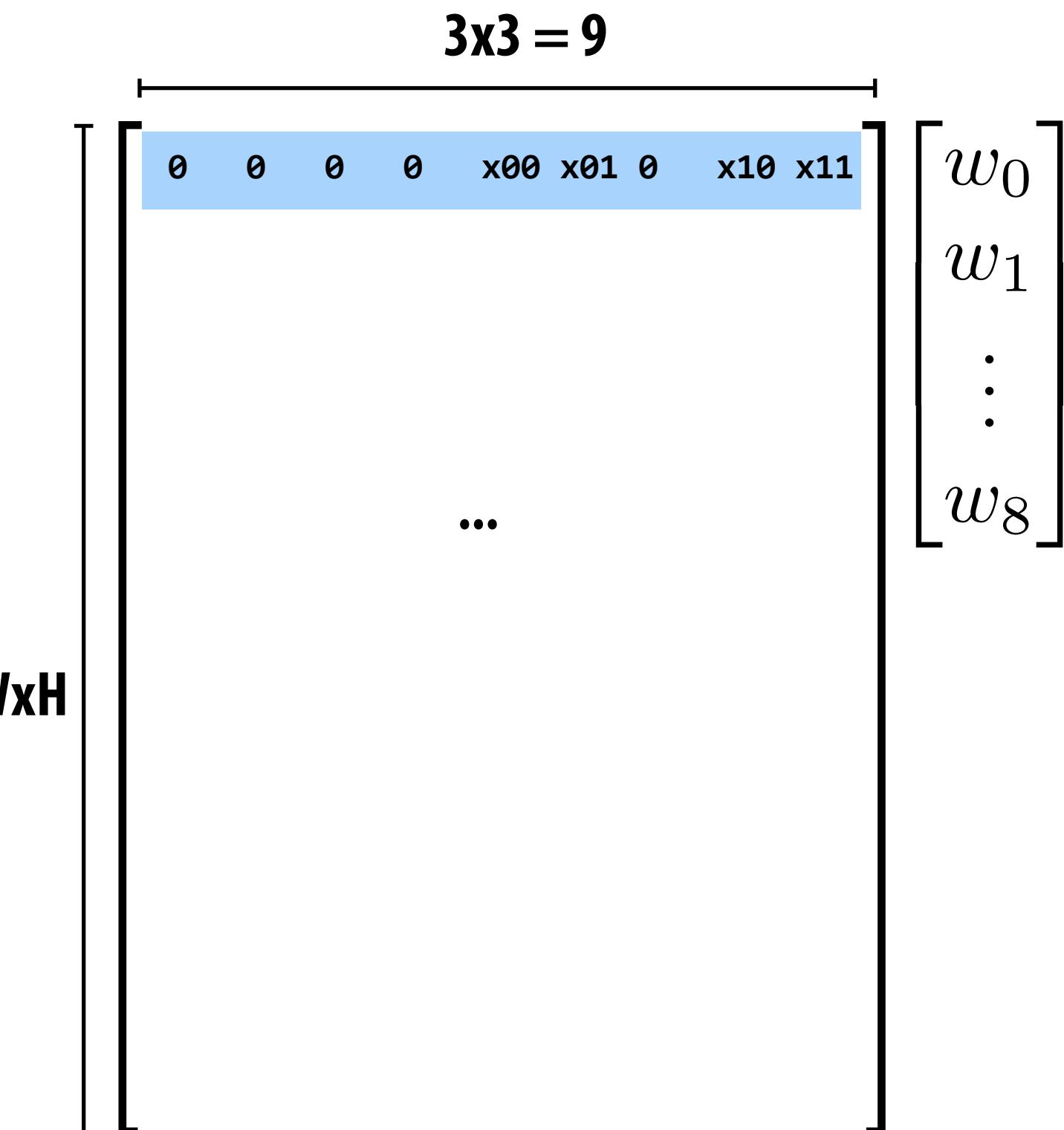
Convolution as matrix-vector product

Construct matrix from elements of input image



X_{00}	X_{01}	X_{02}	X_{03}	...				
X_{10}	X_{11}	X_{12}	X_{13}	...				
X_{20}	X_{21}	X_{22}	X_{23}	...				
X_{30}	X_{31}	X_{32}	X_{33}	...				
...					

$O(N)$ storage multiplier for filter with N elements
Must construct input data matrix



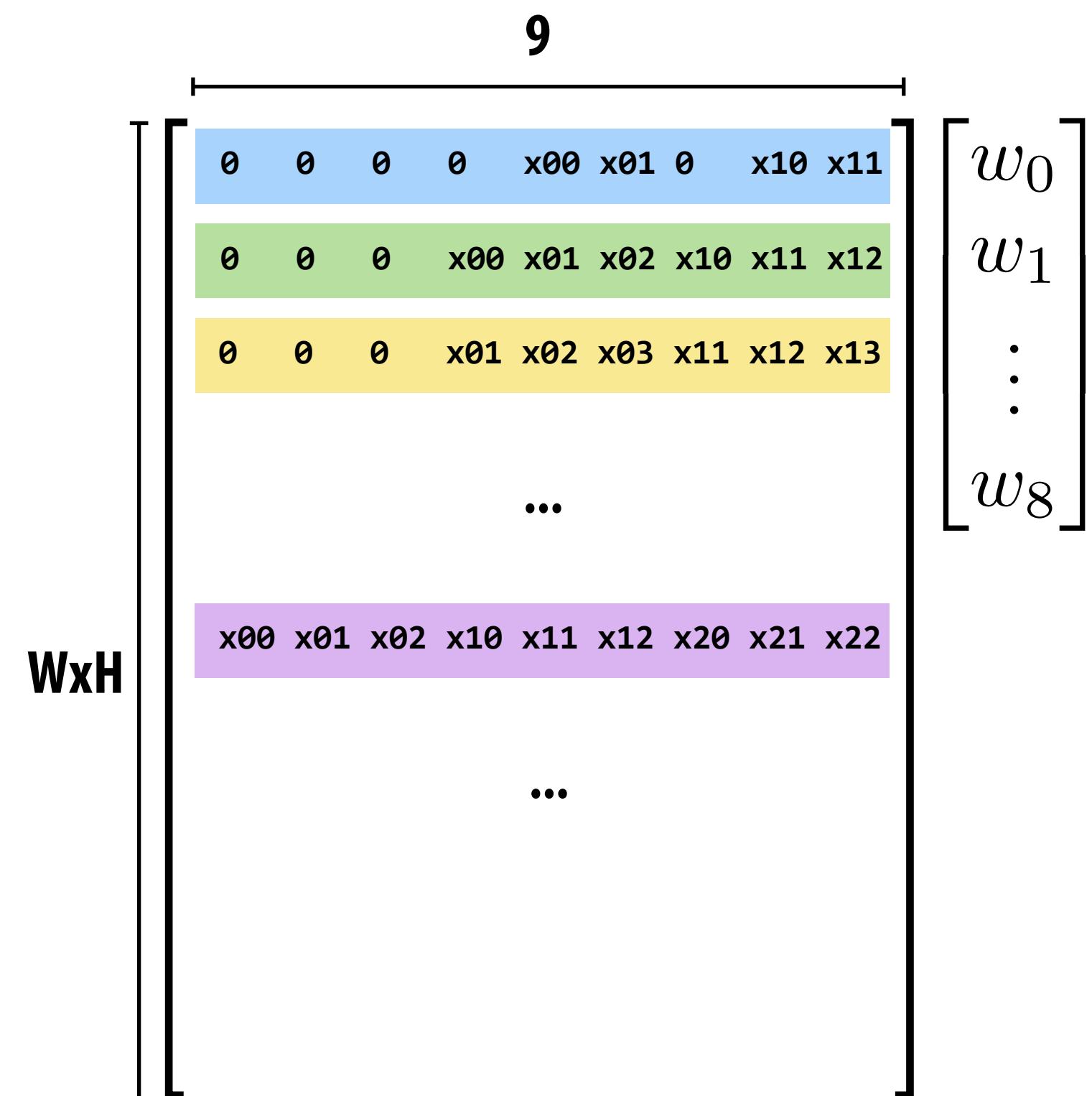
Note: 0-pad matrix

3x3 convolution as matrix-vector product

Construct matrix from elements of input image

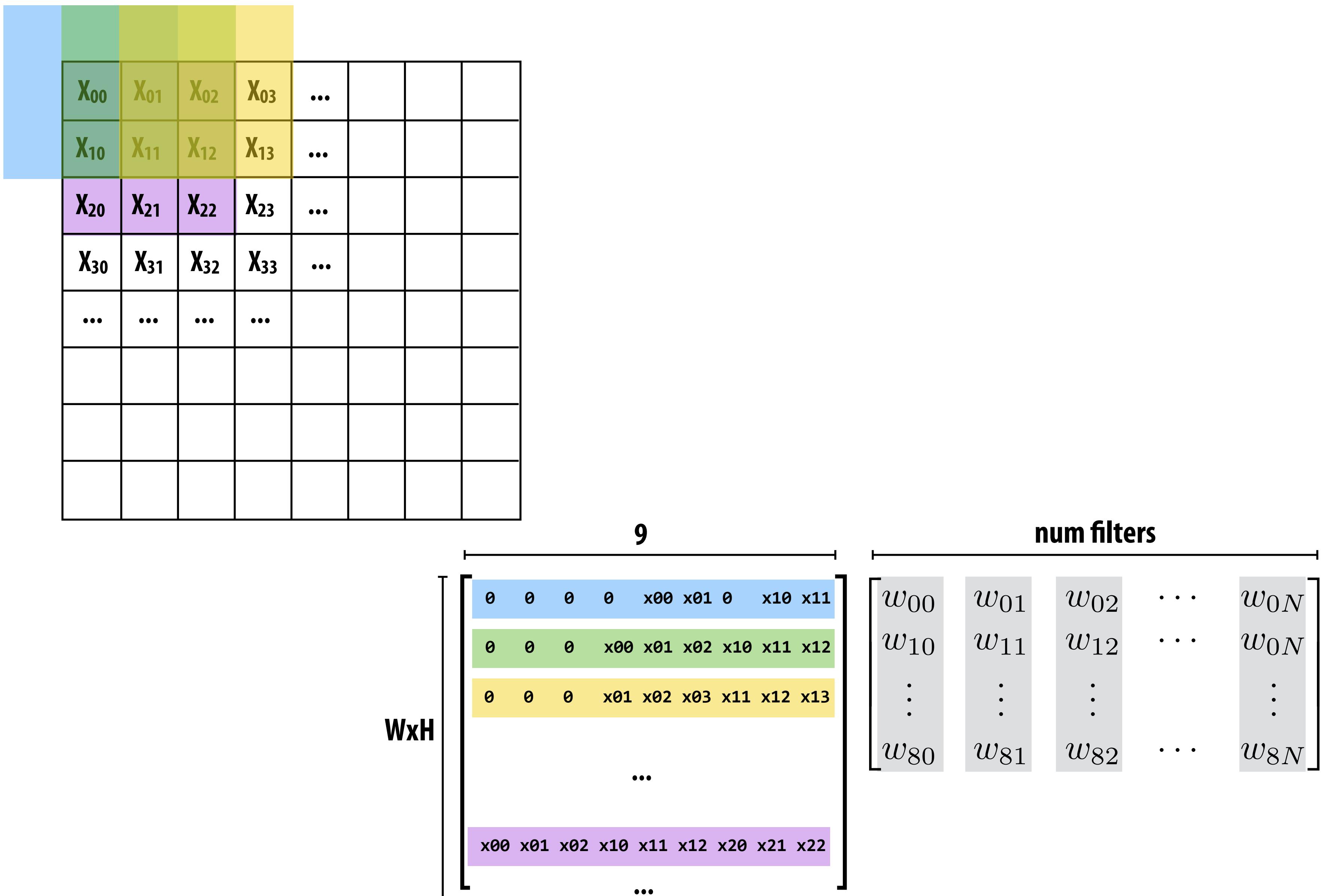
	green	yellow	yellow	yellow				
X_{00}	X_{01}	X_{02}	X_{03}	\dots				
X_{10}	X_{11}	X_{12}	X_{13}	\dots				
X_{20}	X_{21}	X_{22}	X_{23}	\dots				
X_{30}	X_{31}	X_{32}	X_{33}	\dots				
\dots	\dots	\dots	\dots					

$O(N)$ storage overhead for filter with N elements
Must construct input data matrix

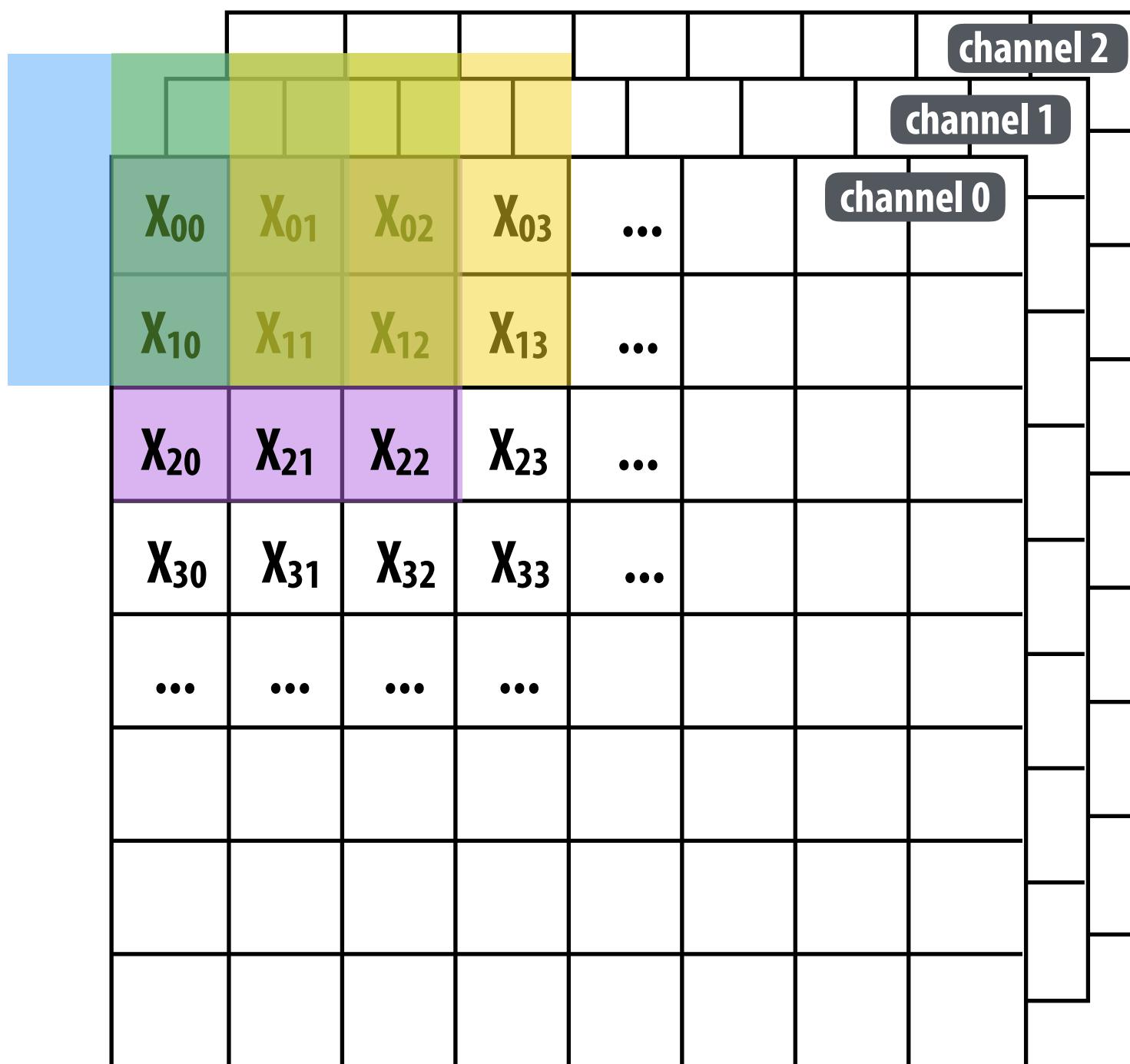


Note: 0-pad matrix

Multiple convolutions as matrix-matrix mult

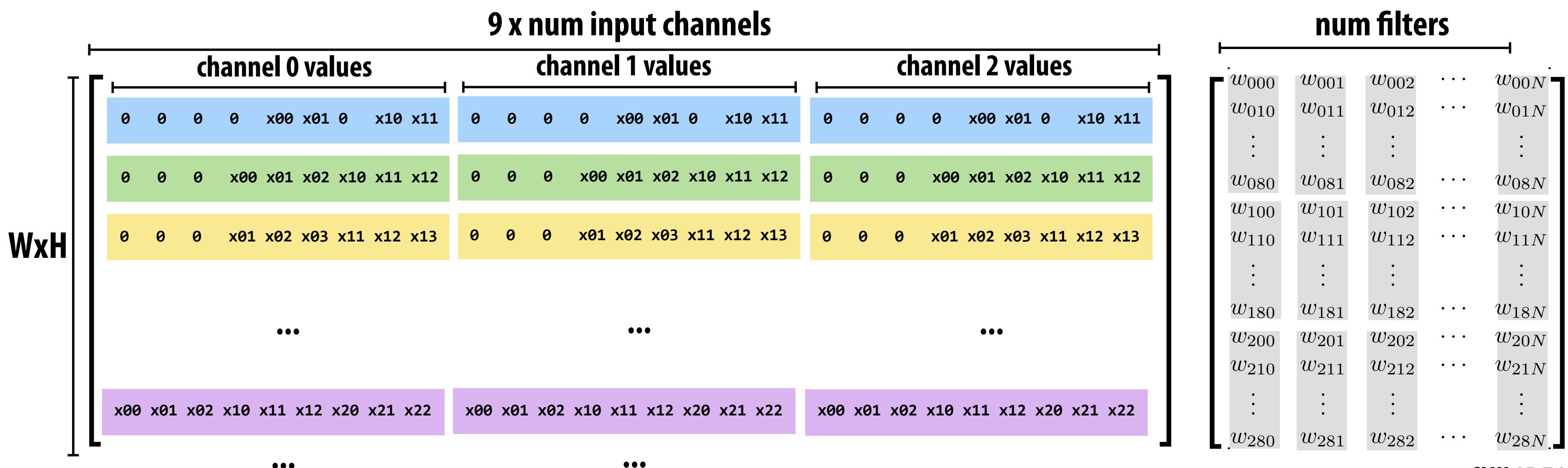


Multiple convolutions on multiple input channels



For each filter, sum responses over input channels

Equivalent to $(3 \times 3 \times \text{num_channels})$ convolution
on $(W \times H \times \text{num_channels})$ input data



VGG memory footprint

Calculations assume 32-bit values (image batch size = 1)

	weights mem:	output size (per image)	(mem)	
input: 224 x 224 RGB image	—	224x224x3	150K	
conv: (3x3x3) x 64	6.5 KB	224x224x64	12.3 MB	
conv: (3x3x64) x 64	144 KB	224x224x64	12.3 MB	
maxpool	—	112x112x64	3.1 MB	
conv: (3x3x64) x 128	228 KB	112x112x128	6.2 MB	
conv: (3x3x128) x 128	576 KB	112x112x128	6.2 MB	
maxpool	—	56x56x128	1.5 MB	
conv: (3x3x128) x 256	1.1 MB	56x56x256	3.1 MB	
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB	
conv: (3x3x256) x 256	2.3 MB	56x56x256	3.1 MB	
maxpool	—	28x28x256	766 KB	
conv: (3x3x256) x 512	4.5 MB	28x28x512	1.5 MB	
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB	
conv: (3x3x512) x 512	9 MB	28x28x512	1.5 MB	
maxpool	—	14x14x512	383 KB	
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB	
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB	
conv: (3x3x512) x 512	9 MB	14x14x512	383 KB	
maxpool	—	7x7x512	98 KB	
fully-connected 4096	392 MB	4096	16 KB	
fully-connected 4096	64 MB	4096	16 KB	
fully-connected 1000	15.6 MB	1000	4 KB	
soft-max		1000	4 KB	

inputs/outputs get multiplied by image batch size

output size
(per image)

multiply by next layer's conv window size to form input matrix to next conv layer!!! (for VGG's 3x3 convolutions, this is 9x data amplification)

Direct implementation of conv layer

```
float input[INPUT_HEIGHT][INPUT_WIDTH][INPUT_DEPTH];
float output[INPUT_HEIGHT][INPUT_WIDTH][LAYER_NUM_FILTERS];
float layer_weights[LAYER_CONVY][LAYER_CONVX][INPUT_DEPTH];

// assumes convolution stride is 1
for (int img=0; img<IMAGE_BATCH_SIZE; img++)
    for (int j=0; j<INPUT_HEIGHT; j++)
        for (int i=0; i<INPUT_WIDTH; i++)
            for (int f=0; f<LAYER_NUM_FILTERS; f++) {
                output[j][i][f] = 0.f;
                for (int kk=0; kk<INPUT_DEPTH; kk++) // sum over filter responses of input channels
                    for (int jj=0; jj<LAYER_CONVY; jj++) // spatial convolution
                        for (int ii=0; ii<LAYER_CONVX; ii++) // spatial convolution
                            output[j][i][f] += layer_weights[f][jj][ii][kk] * input[j+jj][i+ii][kk];
            }
}
```

Seven loops with significant input data reuse: reuse of filter weights (during convolution), and reuse of input values (across different filters)

Avoids O(N) footprint increase by avoiding materializing input matrix

In theory loads O(N) times less data (potentially higher arithmetic intensity... but matrix mult is typically compute-bound)

But must roll your own highly optimized implementation of complicated loop nest.

Conv layer in Halide

```
int in_w, in_h, in_ch = 4;           // input params: assume initialized
Func in_func;                      // assume input function is initialized

int num_f, f_w, f_h, pad, stride;   // parameters of the conv layer

Func forward = Func("conv");
Var x("x"), y("y"), z("z"), n("n"); // n is minibatch dimension

// This creates a padded input to avoid checking boundary
// conditions while computing the actual convolution
f_in_bound = BoundaryConditions::repeat_edge(in_func, 0, in_w, 0, in_h);

// Create image buffers for layer parameters
Image<float> W(f_w, f_h, in_ch, num_f)
Image<float> b(num_f);

RDom r(0, f_w, 0, f_h, 0, in_ch);
// Initialize to bias
forward(x, y, z, n) = b(z);
forward(x, y, z, n) += W(r.x, r.y, r.z, z) *
    f_in_bound(x*stride + r.x - pad,
                y*stride + r.y - pad,
                r.z, n);
```

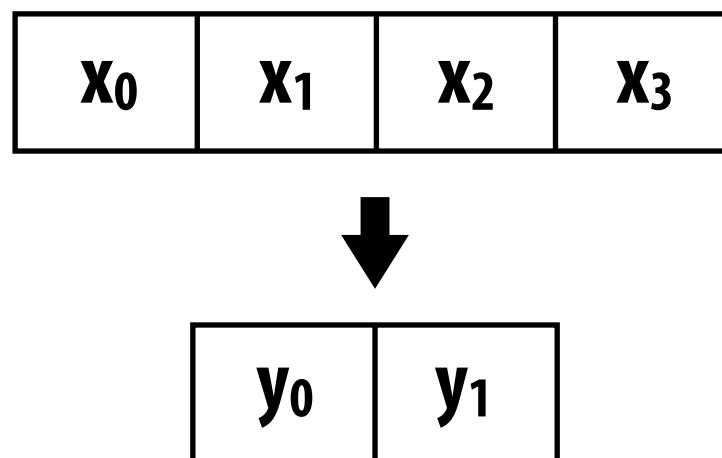
Consider scheduling this seven-dimensional loop nest.

Algorithmic improvements

- Direct convolution can be implemented efficiently in Fourier domain (convolution → element-wise multiplication)
 - Overhead: FFT to transform inputs into Fourier domain, inverse FFT to get responses back to spatial domain ($N \lg N$)
 - Inverse transform amortized over all input channels (due to summation over inputs)

- Direct convolution using work-efficient Winograd convolutions

1D example: consider producing two outputs of a 3-tap 1D convolution with weights: $w_0 w_1 w_2$



1D 3-tap total cost:
4 multiplies
8 additions
(4 to compute m's + 4 to produce final result)

Direct convolution: 6 multiples, 4 adds
In 2D can notably reduce multiplications
(3x3 filter: 2.25x fewer multiples for 2x2 block of output)

$$\begin{bmatrix} y_0 \\ y_1 \end{bmatrix} = \begin{bmatrix} x_0 & x_1 & x_2 \\ x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} m_1 + m_2 + m_3 \\ m_2 - m_3 - m_4 \end{bmatrix}$$

$$m_1 = (x_0 - x_1)w_0$$

$$m_2 = (x_1 + x_2) \frac{w_0 + w_1 + w_2}{2}$$

$$m_3 = (x_2 - x_1) \frac{w_0 - w_1 + w_2}{2}$$

$$m_4 = (x_1 - x_3)w_2$$

Filter dependent
(can be precomputed)

Deep neural networks on GPUs

- Today, many-performant DNN implementations target GPUs
 - High arithmetic intensity computations (computational characteristics similar to dense matrix-matrix multiplication)
 - Benefit from flop-rich architectures
 - Highly-optimized library of kernels exist for GPUs (NVIDIA's cuDNN)
 - Most CPU-based implementations use basic matrix-multiplication-based formulation (good implementations could run faster!)



Emerging architectures for deep learning?

- **NVIDIA Pascal (NVIDIA's latest GPU)**
 - Adds double-throughput 16-bit floating point ops
 - This feature is already common on mobile GPUs
- **Intel Xeon Phi (Knights Landing)**
 - Flop rich 68-core x86 processor for scientific computing and machine learning
- **FPGAs, ASICs**
 - Not new: FPGA solutions have been explored for years
 - Significant amount of ongoing industry and academic research
 - Google's TPU
 - We will discuss several recent papers in an upcoming class

Summary: efficiently evaluating deep nets

■ Computational structure

- **Convlayers: high arithmetic intensity, significant portion of cost of evaluating a network**
- **Similar data access patterns to dense-matrix multiplication (exploiting temporal reuse is key)**
- **But straight reduction to matrix-matrix multiplication is often sub-optimal**
- **Work-efficient techniques for convolutional layers (FFT-based, Winograd convolutions)**

■ Large numbers of parameters: significant interest in reducing size of networks for both training and evaluation

- **Will discuss pruning/quantization in upcoming classes**

■ Many ongoing studies of specialized hardware architectures for efficient evaluation

- **Future CPUs/GPUs, ASICs, FPGAs, ...**
- **Specialization will be important to achieving “always on” applications**