

CONVOLUTIONAL NETWORKS

Hint: Don't do this yourself!
Always steal from other people!

THIS IS YOUR MACHINE LEARNING SYSTEM?

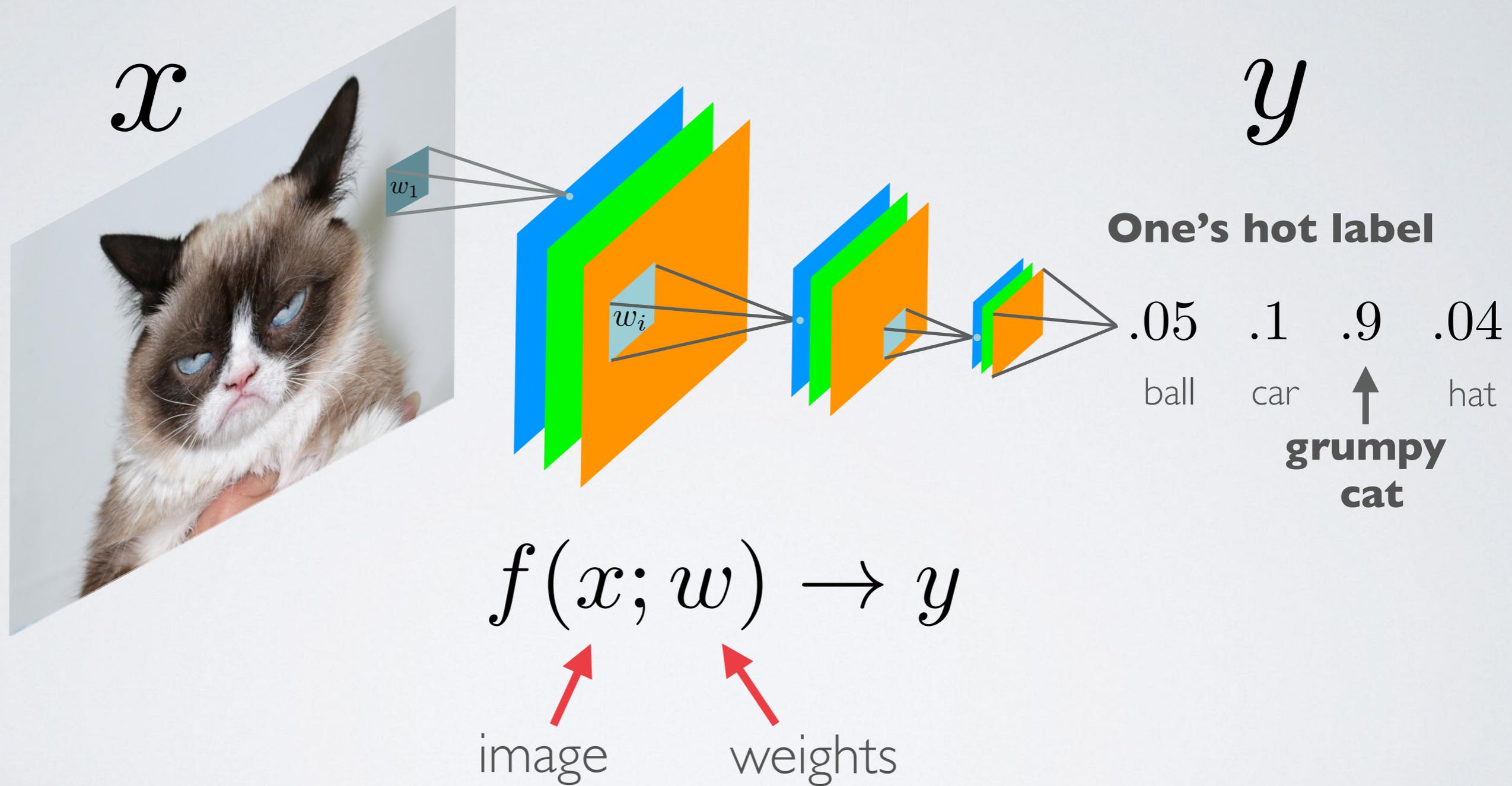
YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

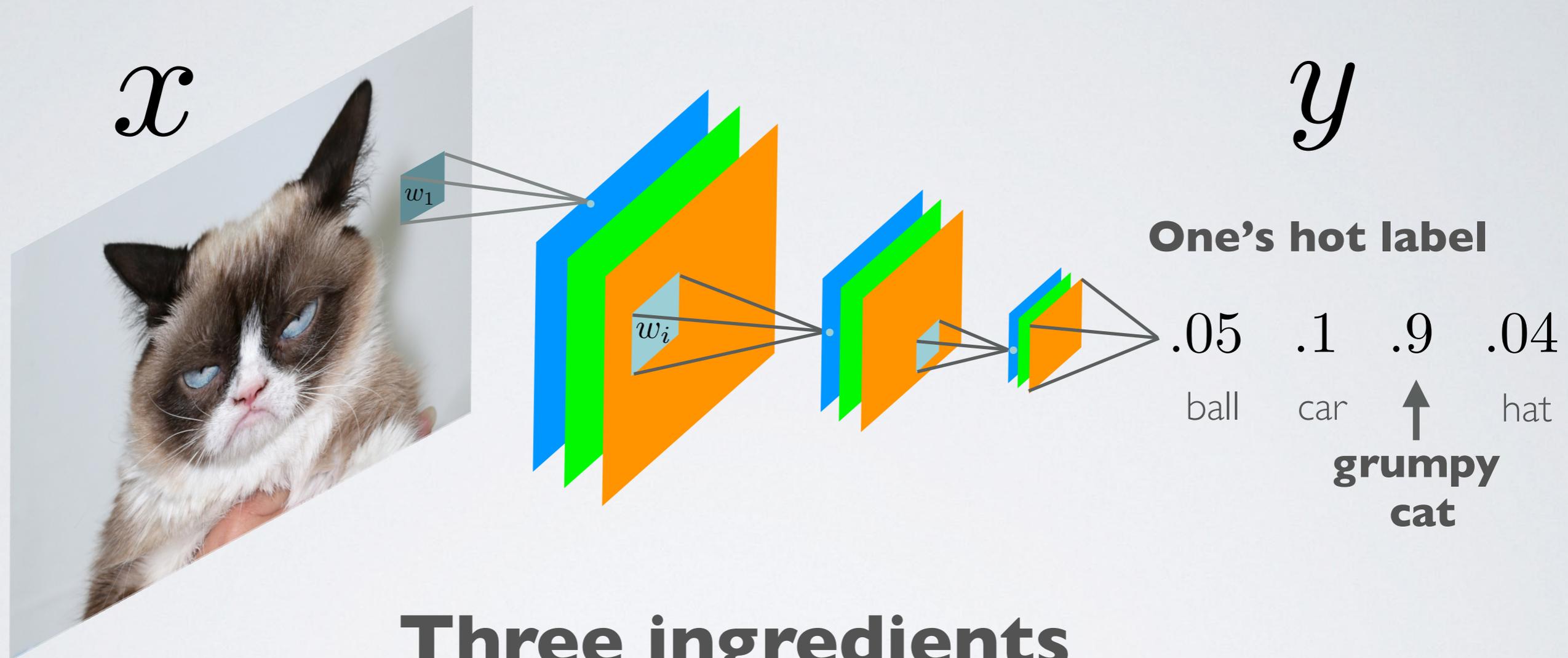
JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.



CONVOLUTIONAL NET



CONVOLUTIONAL NET



Three ingredients

- convolutions
- non-linearities
- pooling

CONVOLUTION

**Input
image**

\mathcal{X}

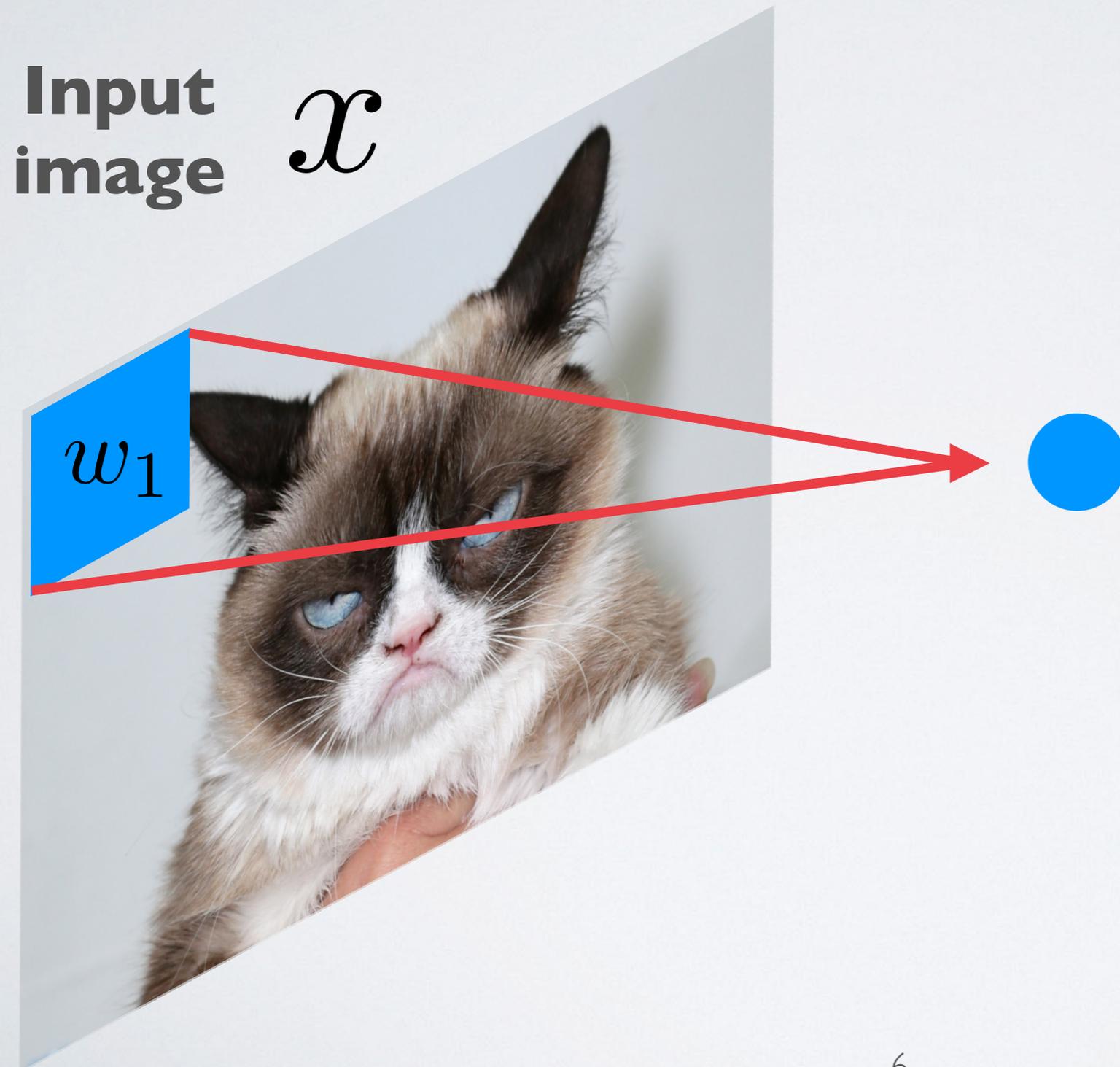


Convolutional filter

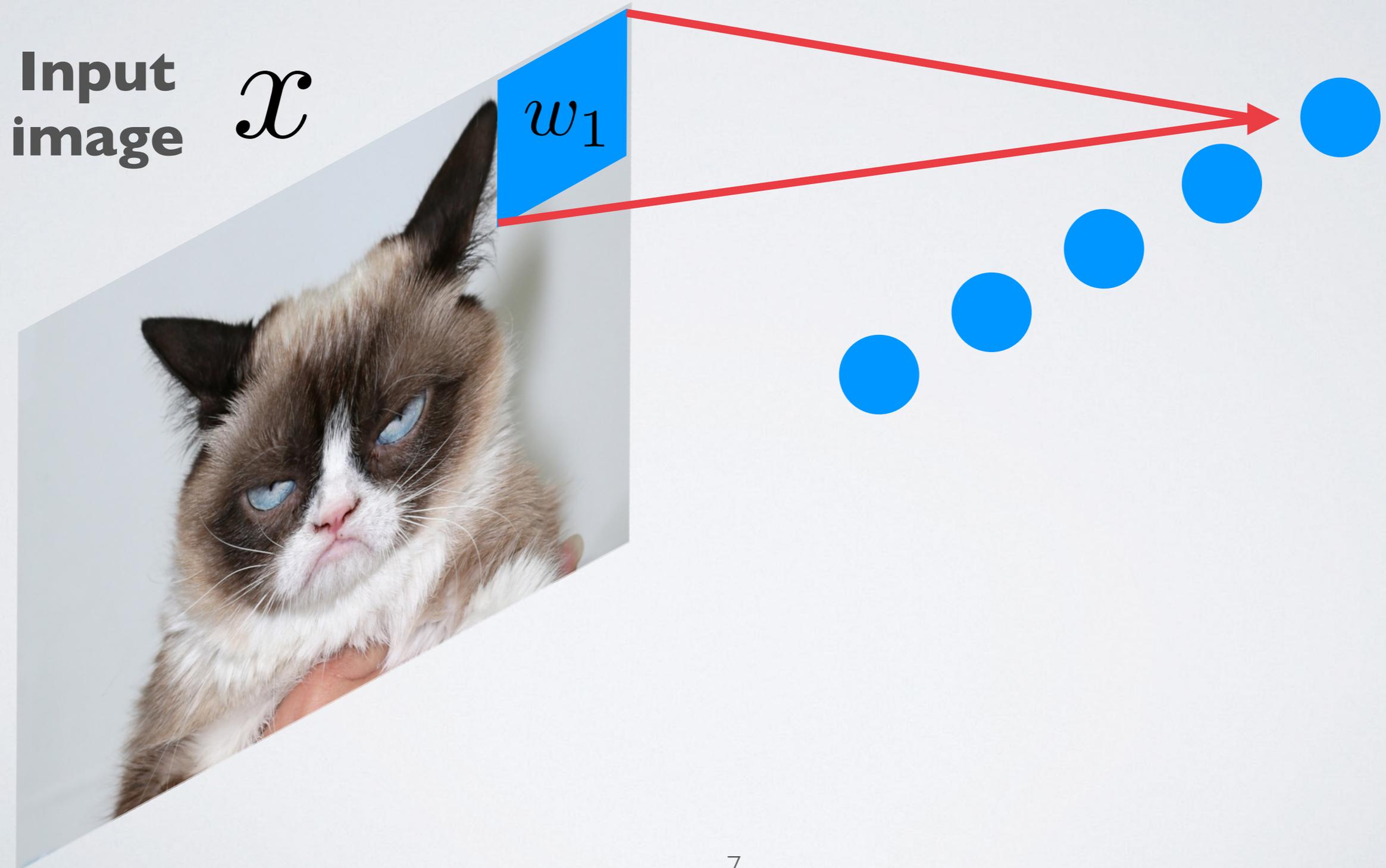
w_1

5	3	8
2	1	7
1	9	4

CONVOLUTION



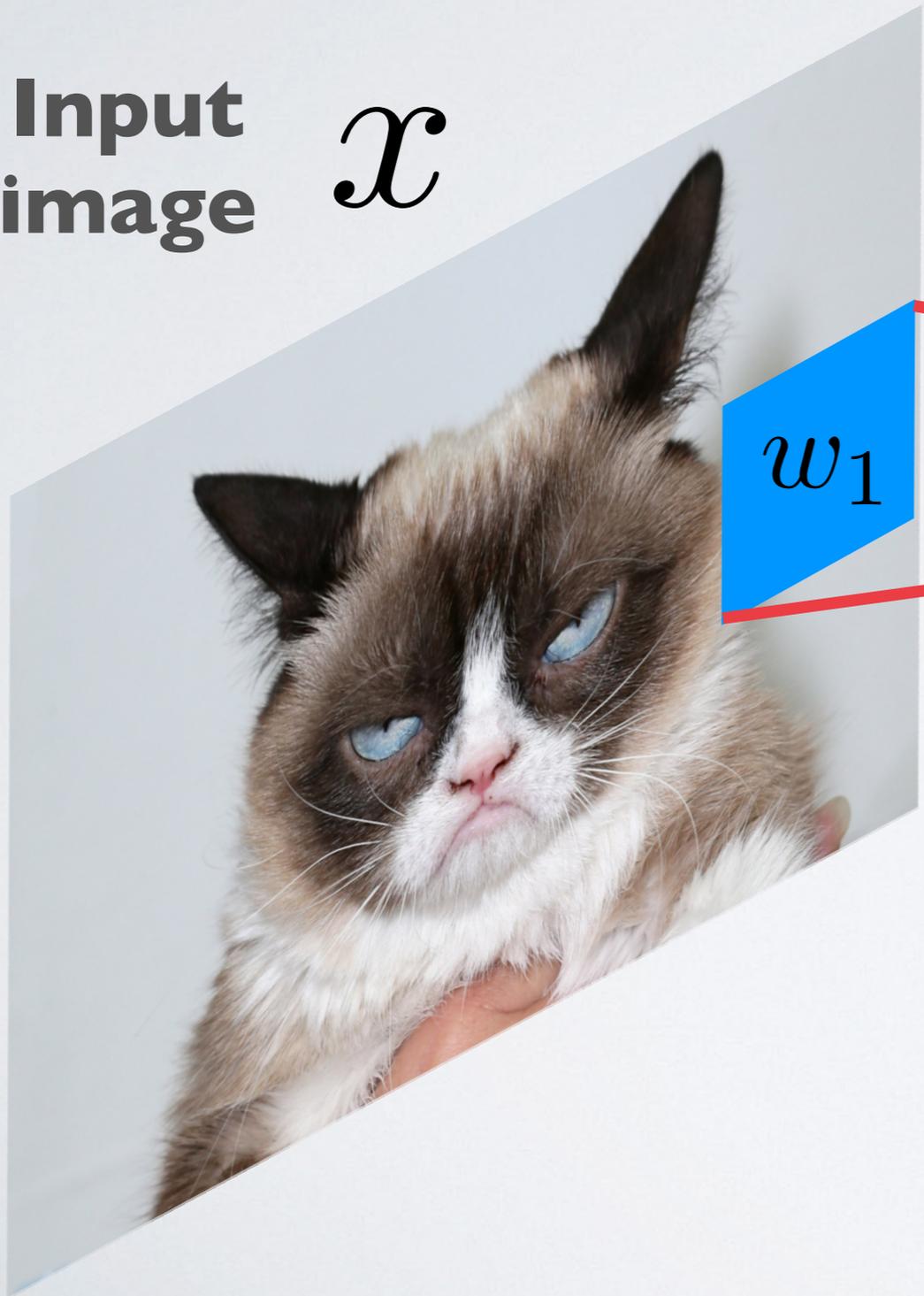
CONVOLUTION



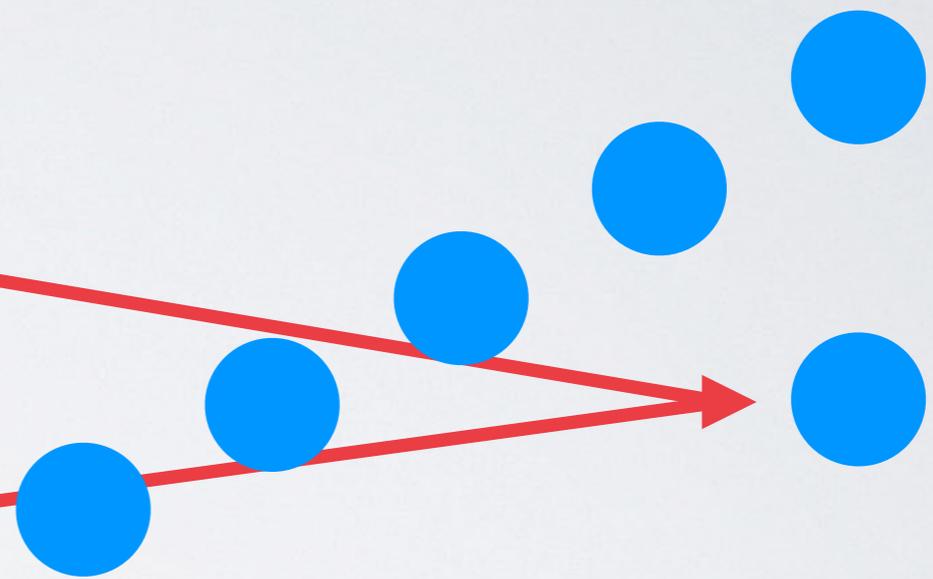
CONVOLUTION

Input
image

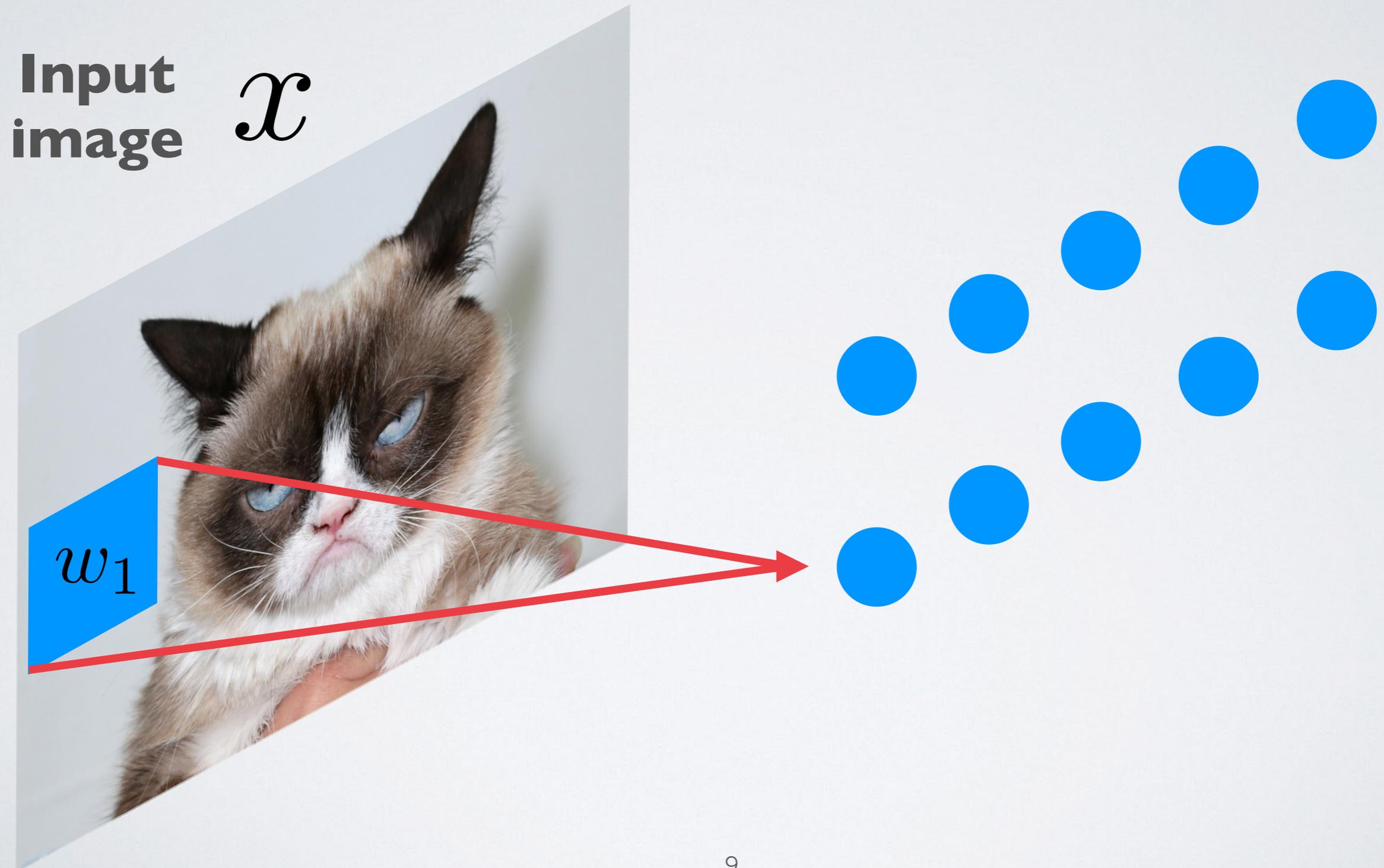
\mathcal{X}



w_1



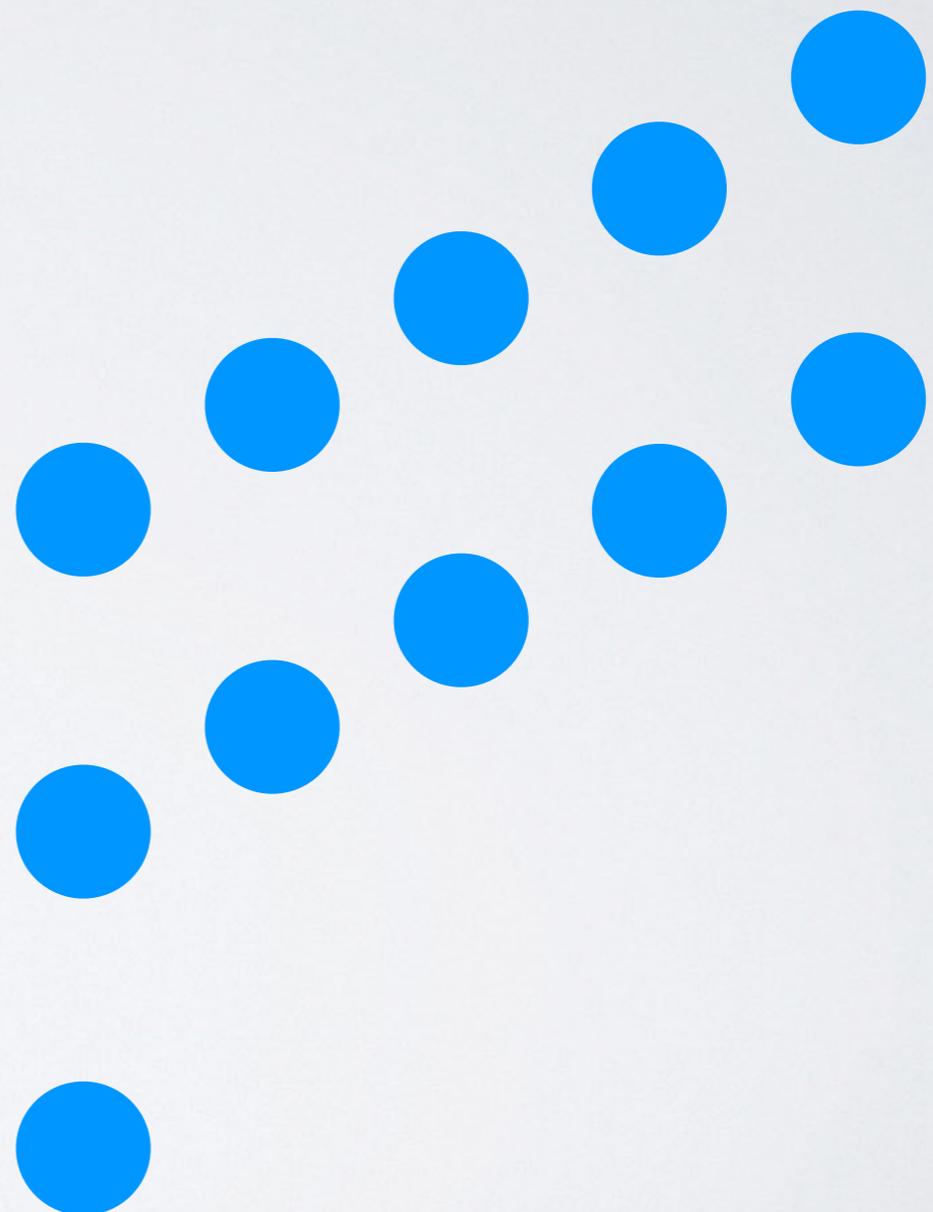
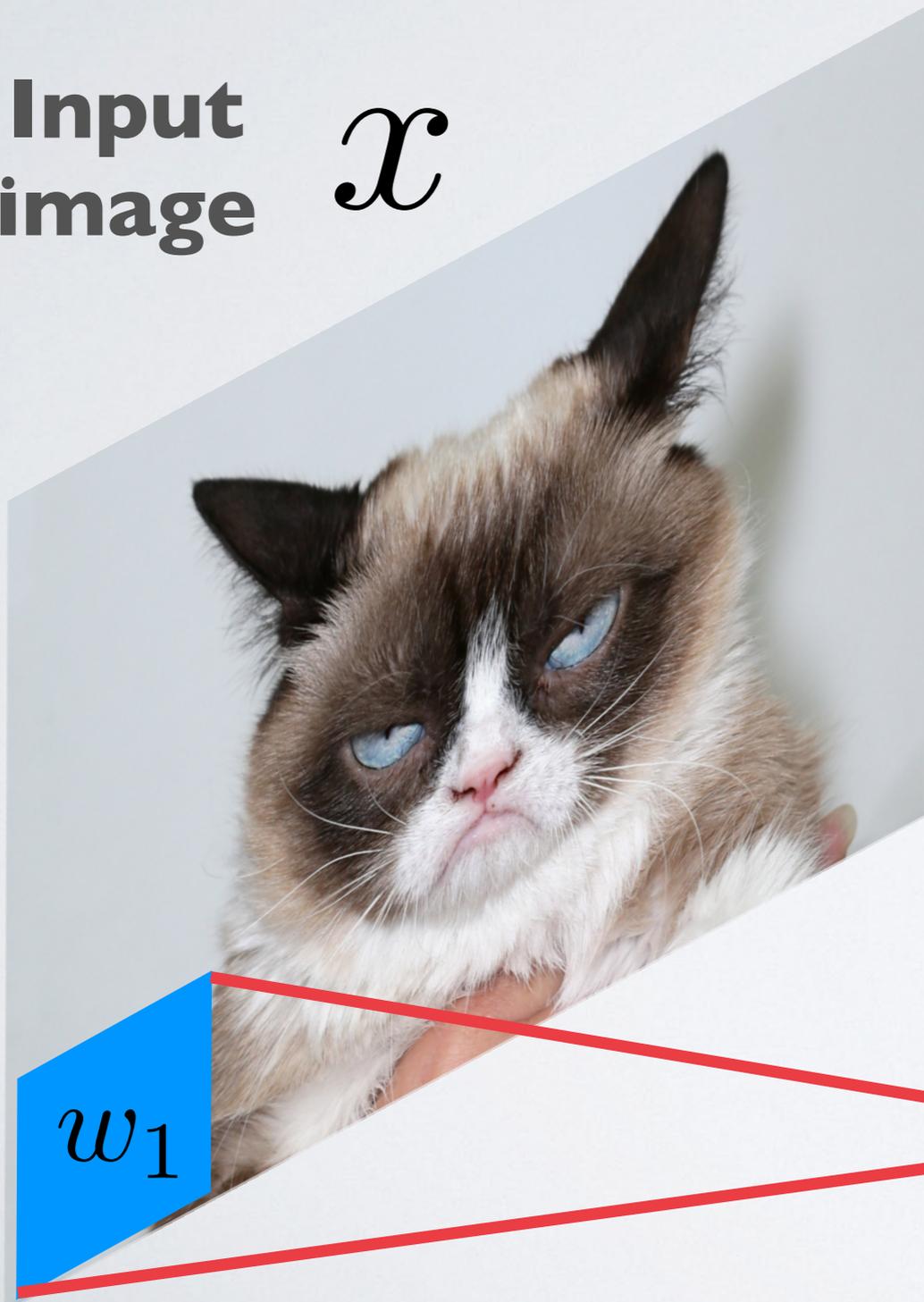
CONVOLUTION



CONVOLUTION

Input
image

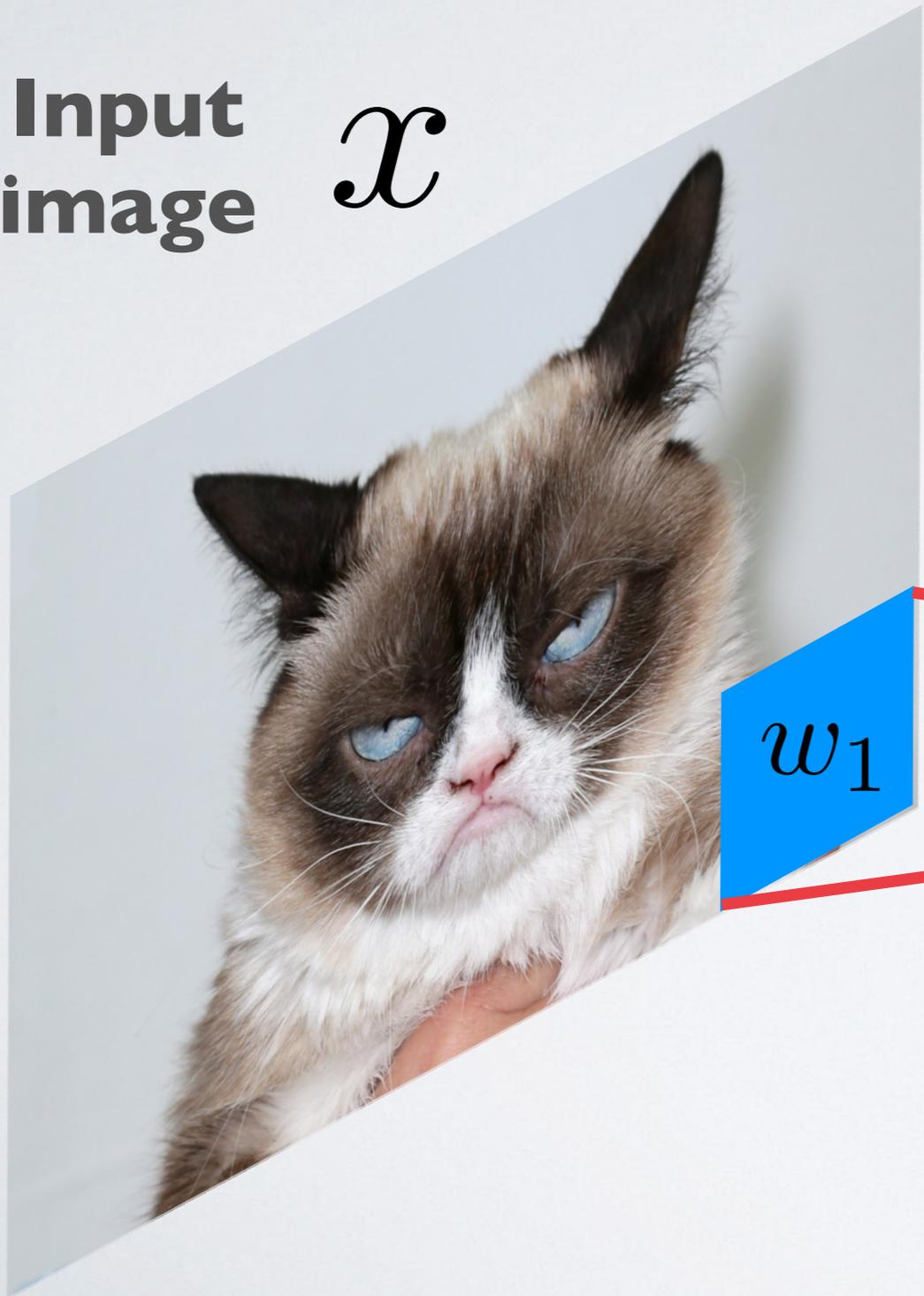
\mathcal{X}



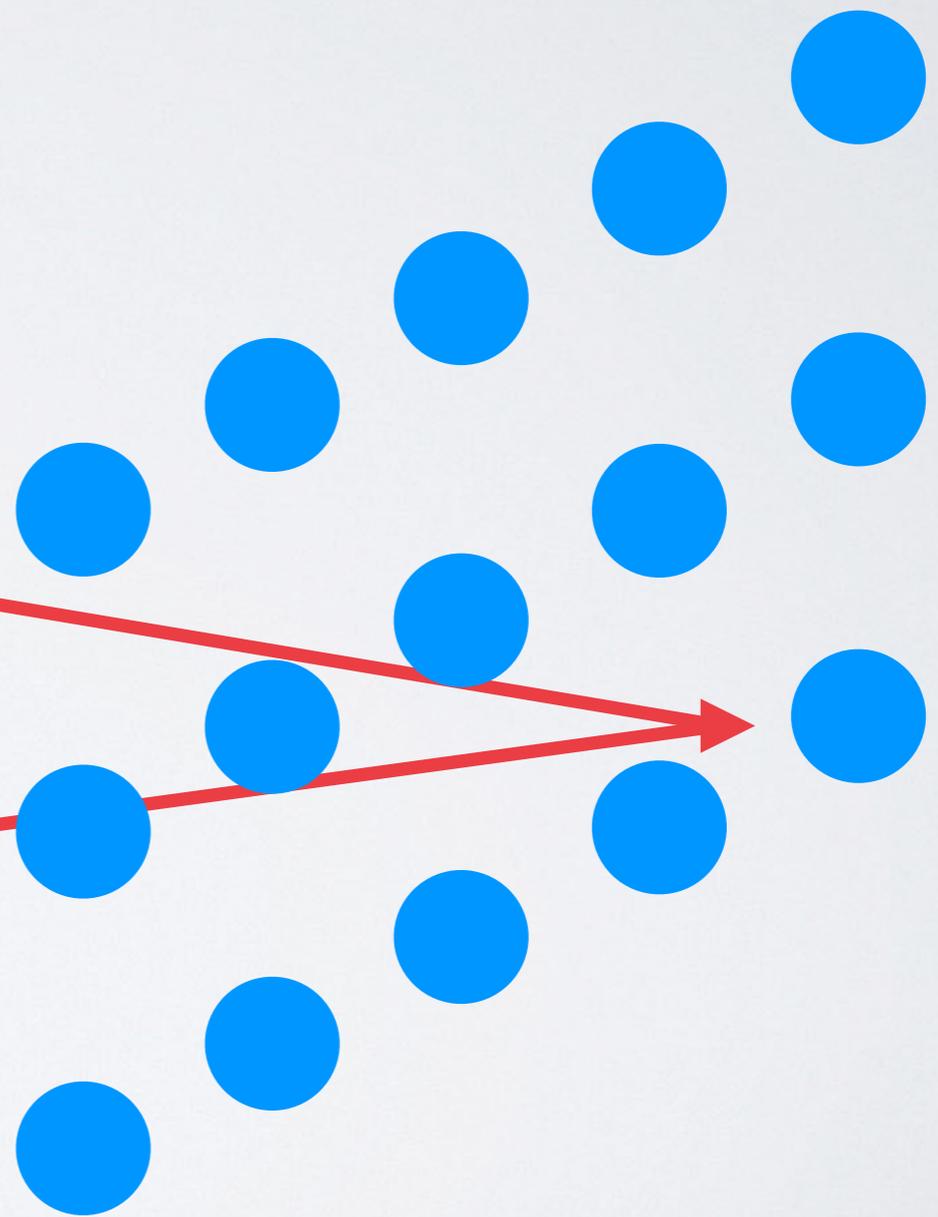
CONVOLUTION

Input
image

\mathcal{X}



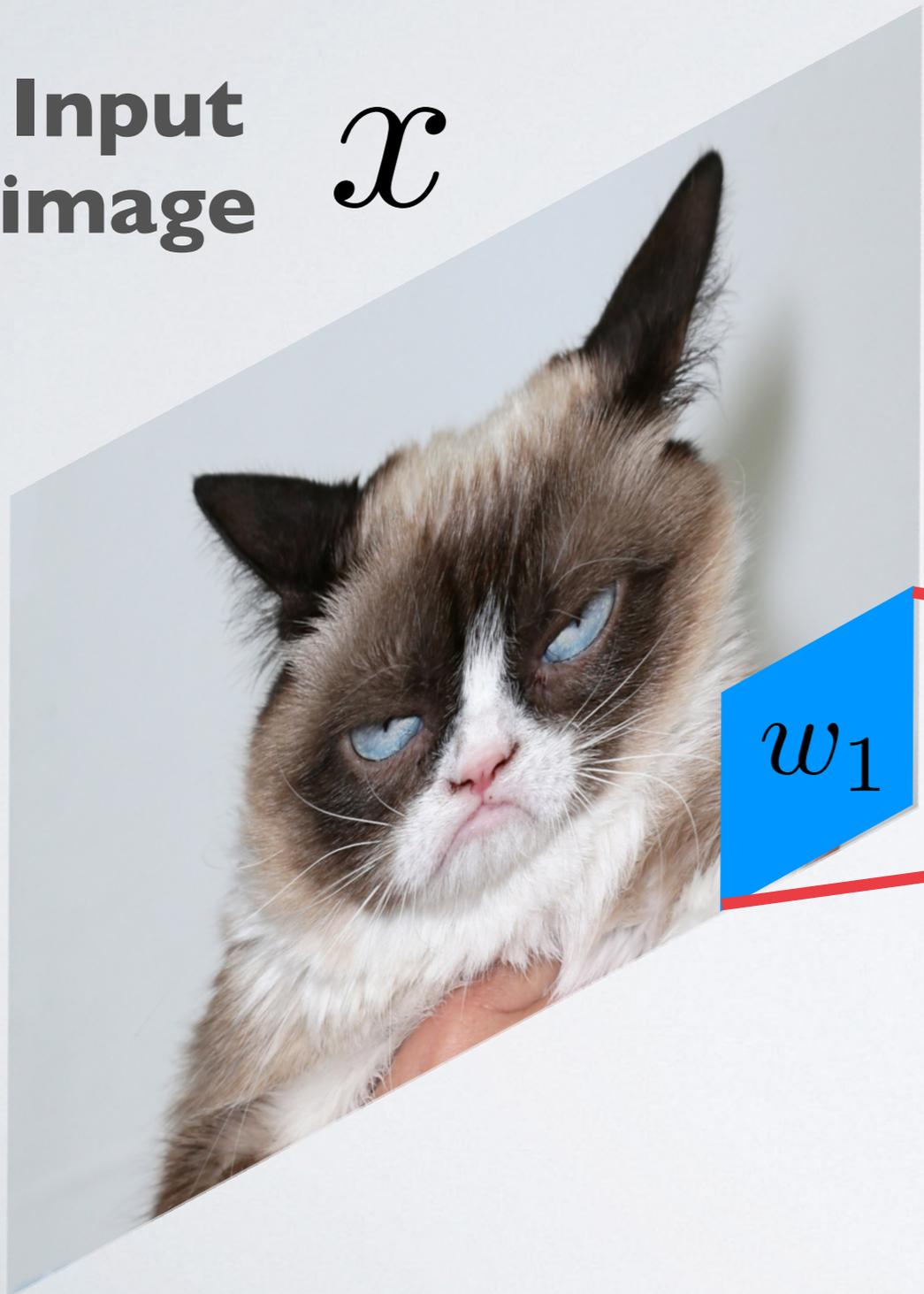
w_1



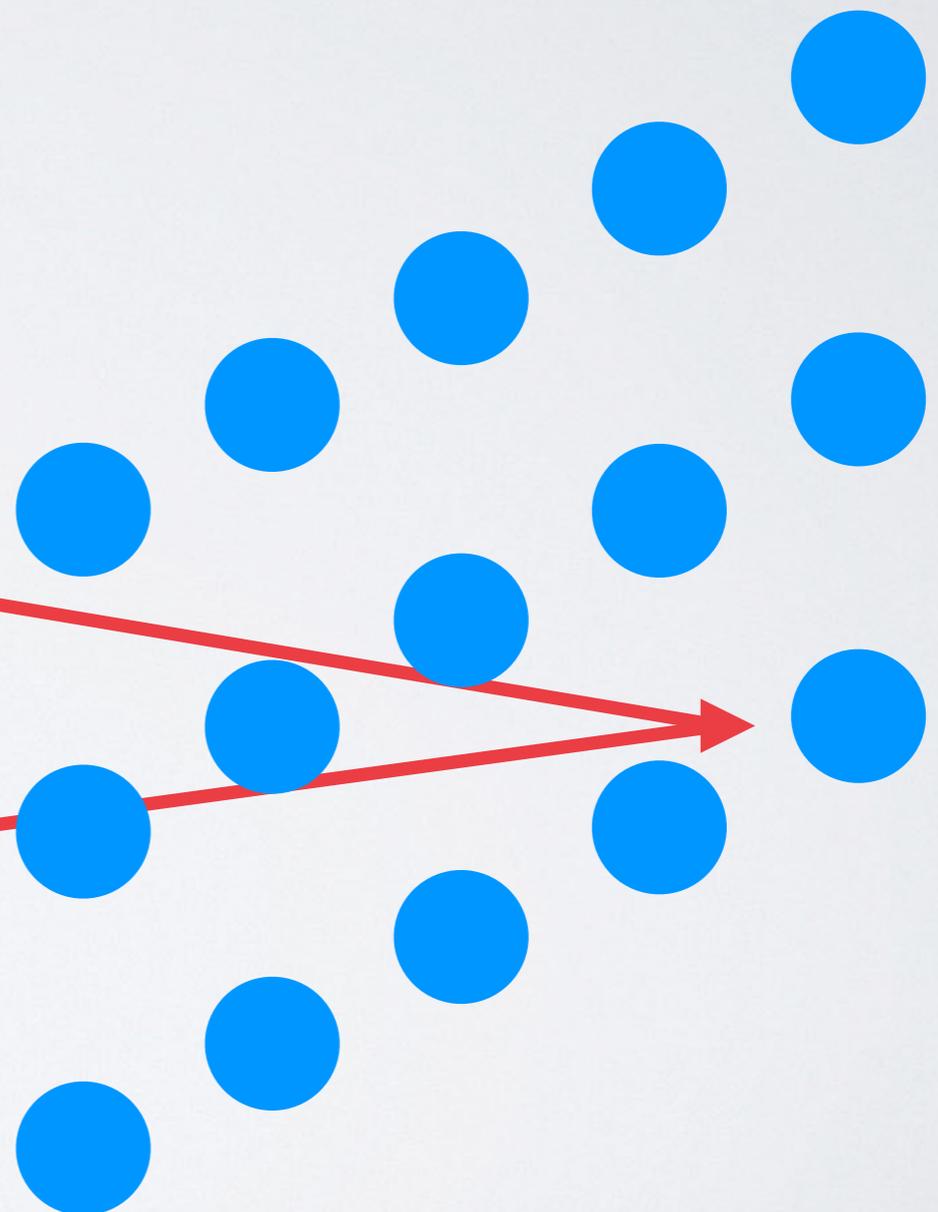
CONVOLUTION

Input
image

\mathcal{X}

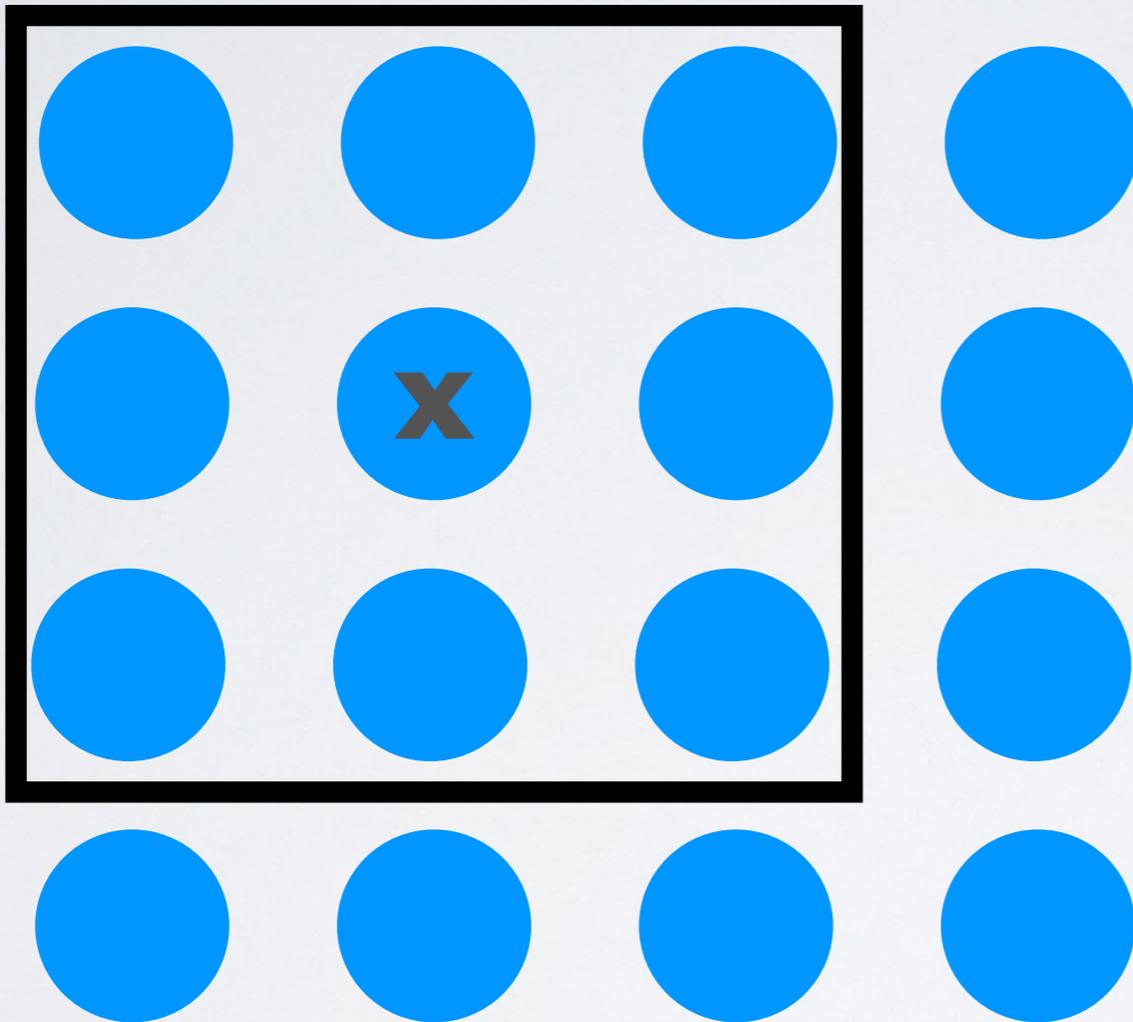


“Feature map”



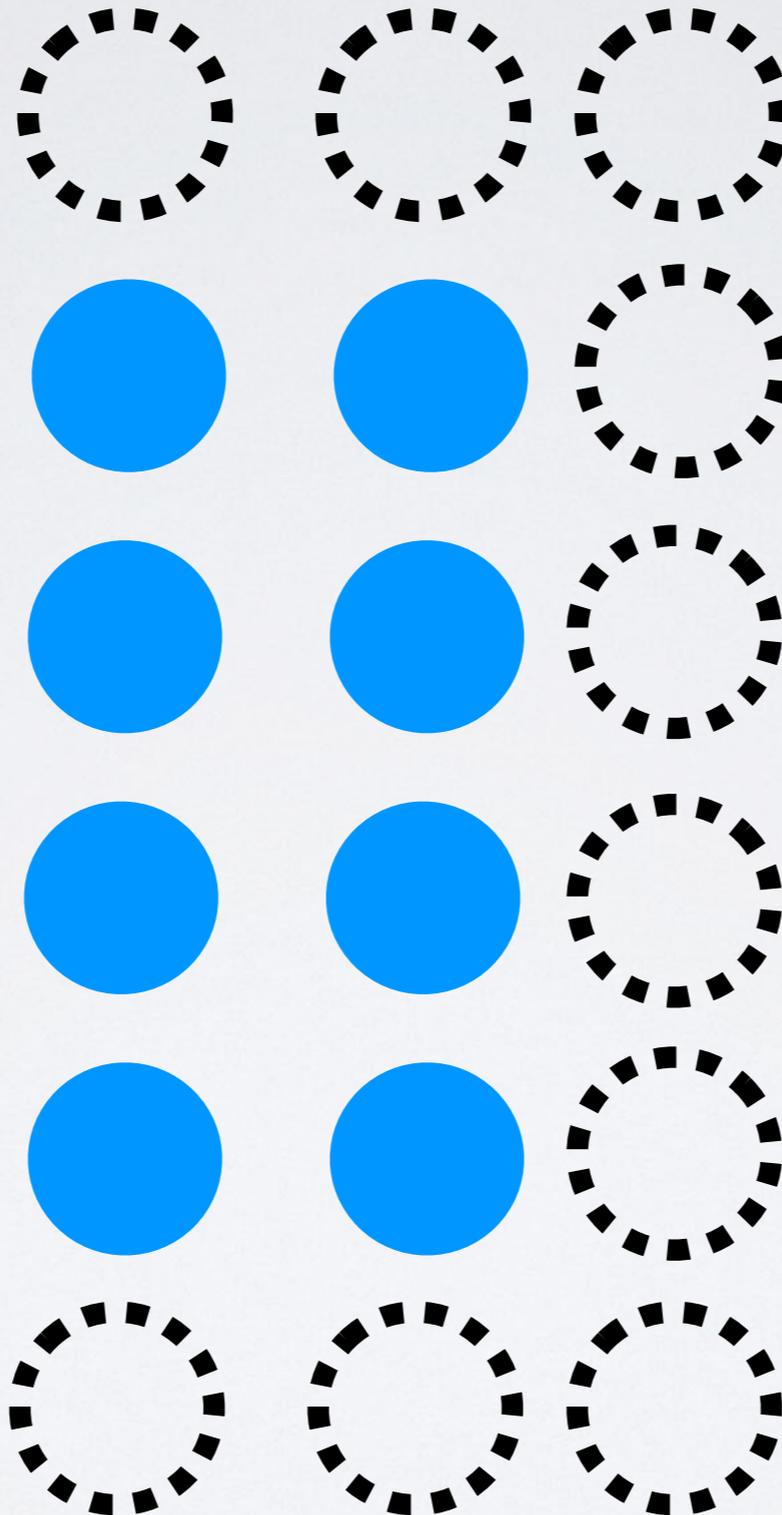
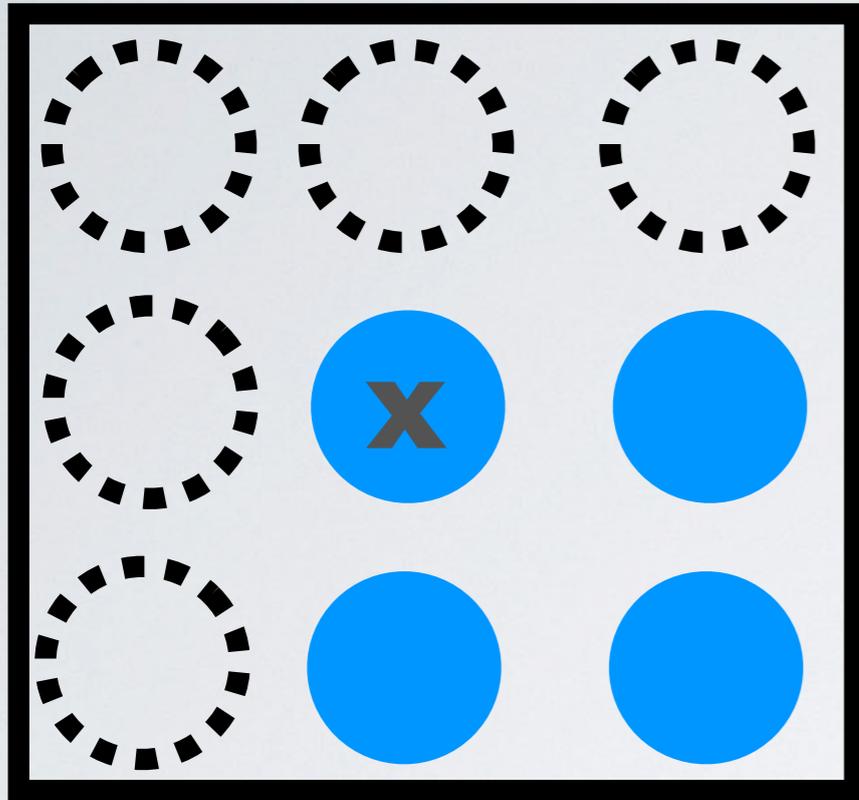
VARIANTS OF CONV

PADDING



**3x3 conv
produces
feature maps
that are 2 units
smaller**

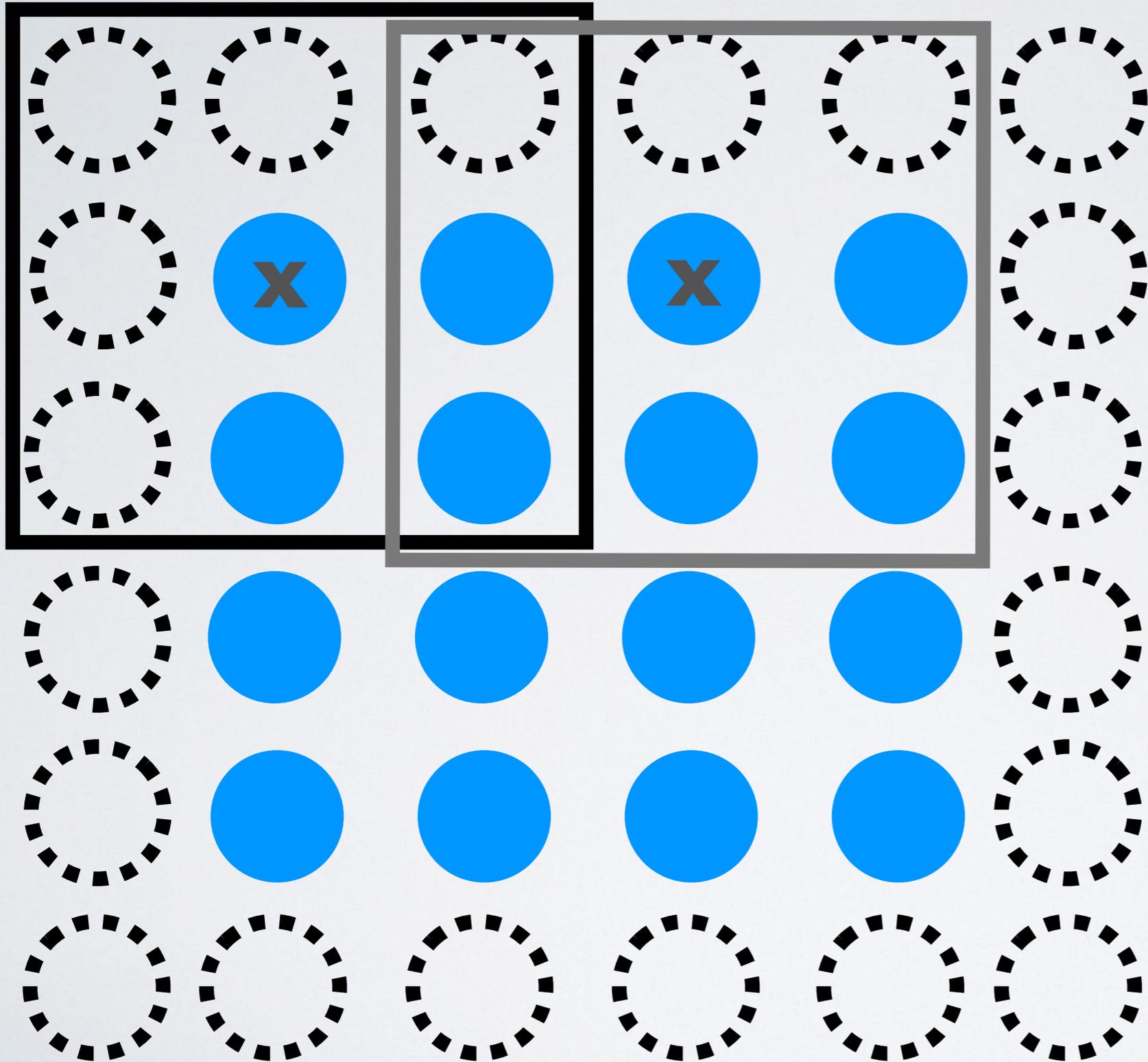
PADDING



padding=1

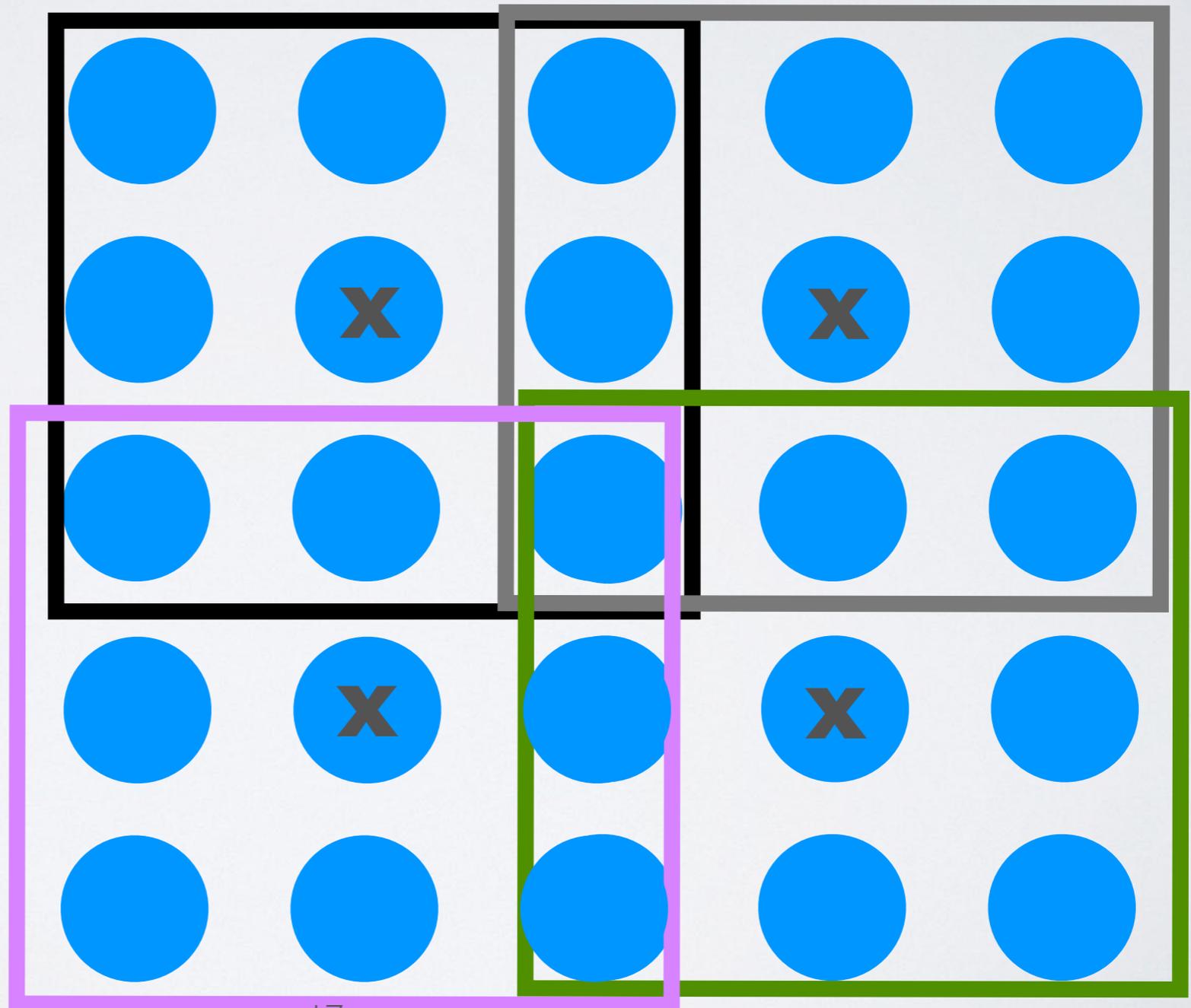
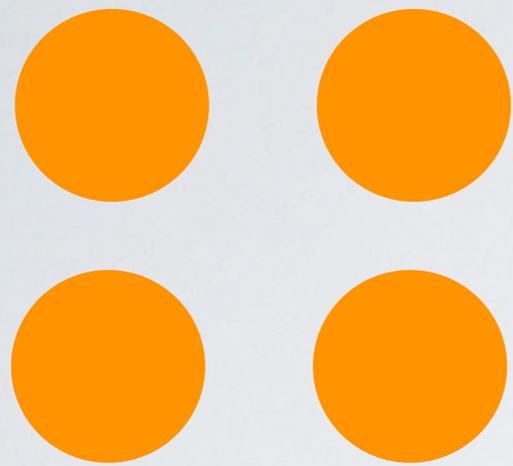
Output is same size
as input

STRIDE



**stride=2 cuts
resolution by
half**

TRANSPOSE / DECONVOLUTION

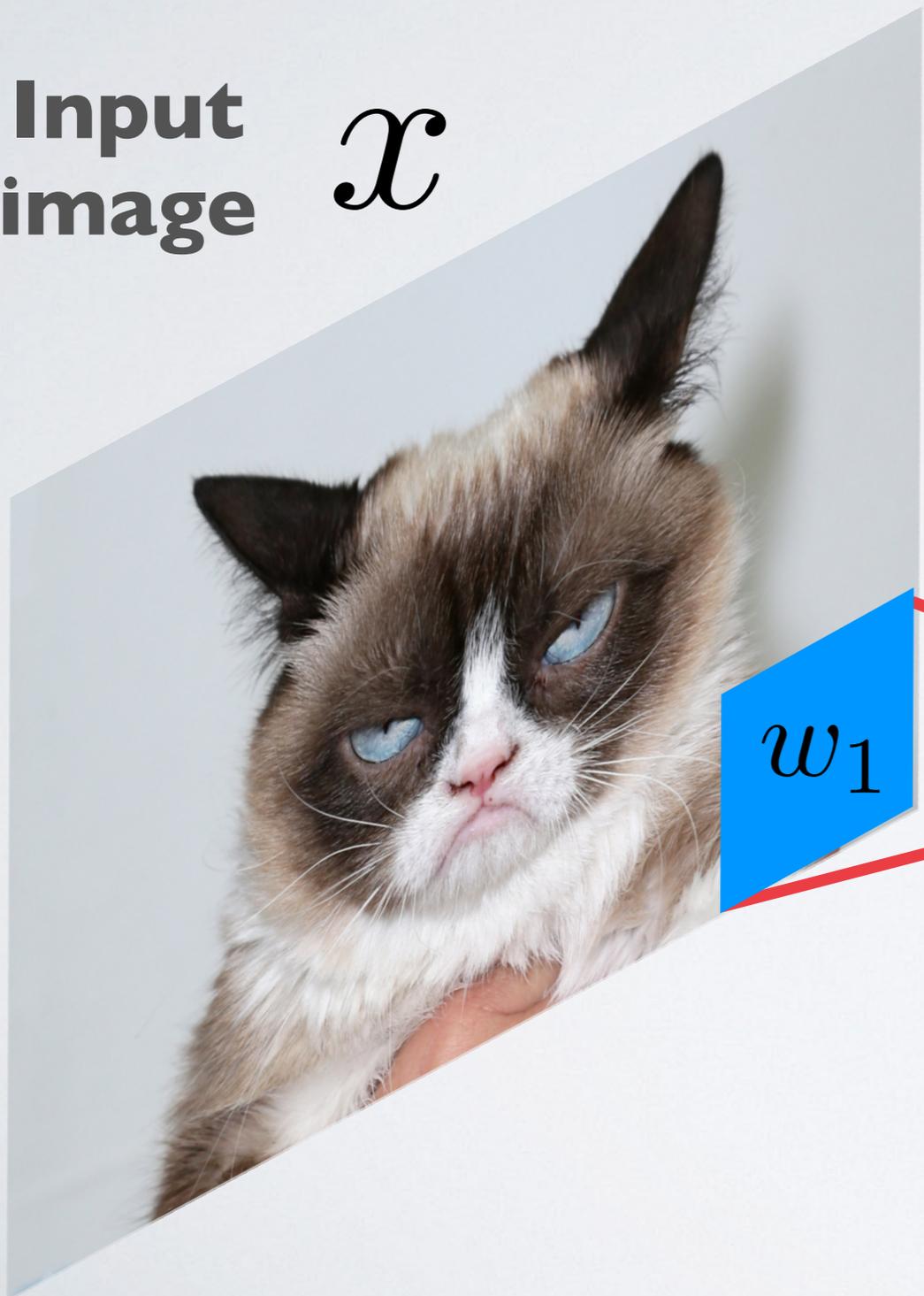


RELU

RELU NONLINEARITY

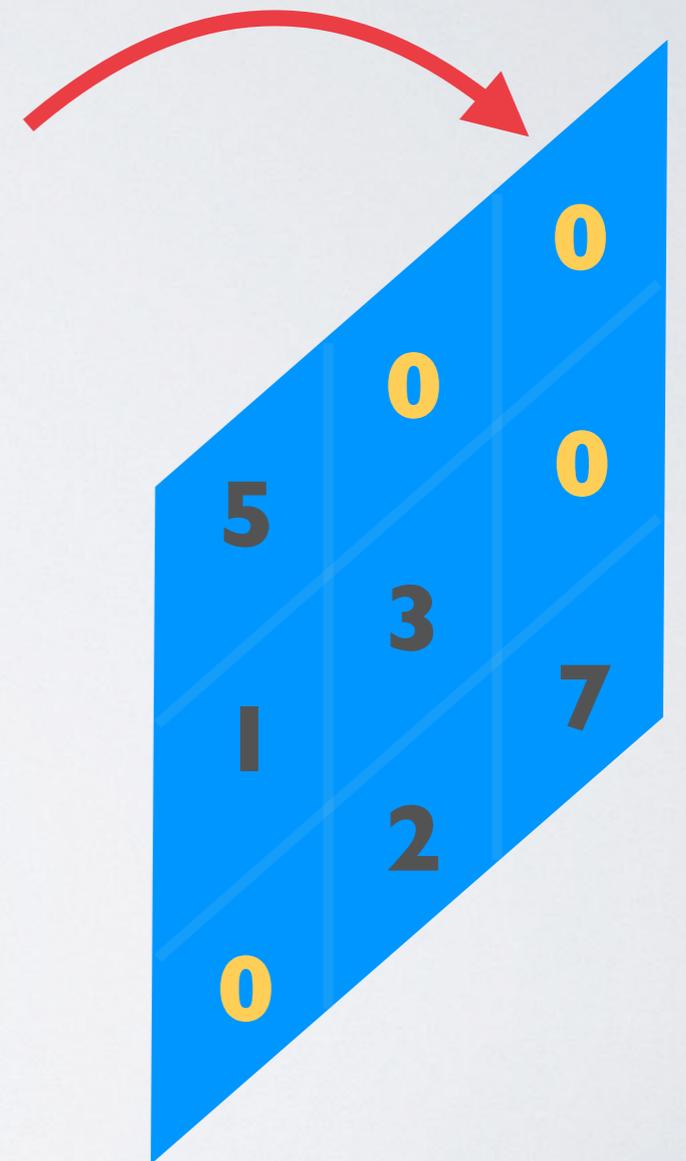
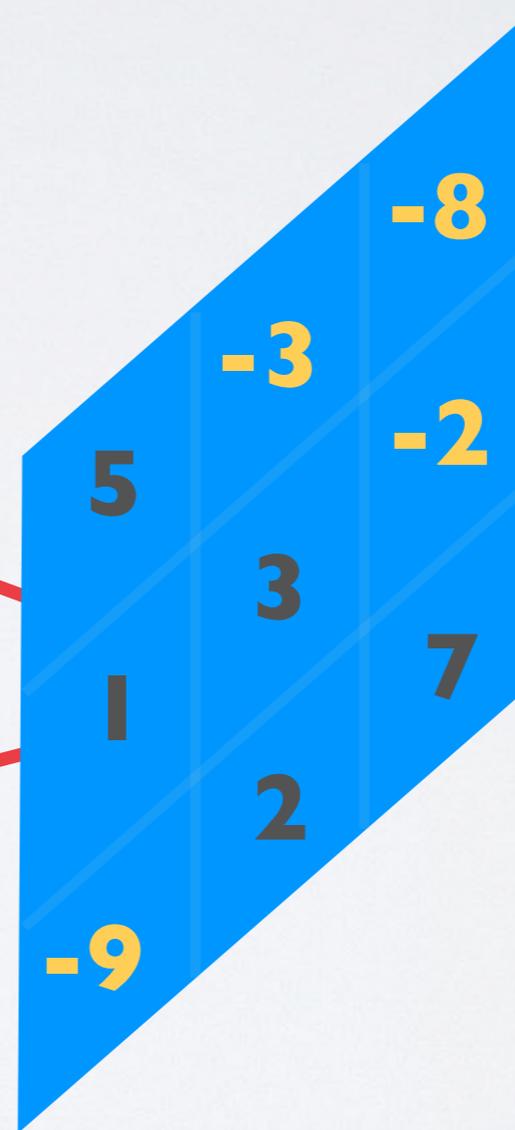
Input
image

\mathcal{X}



w_1

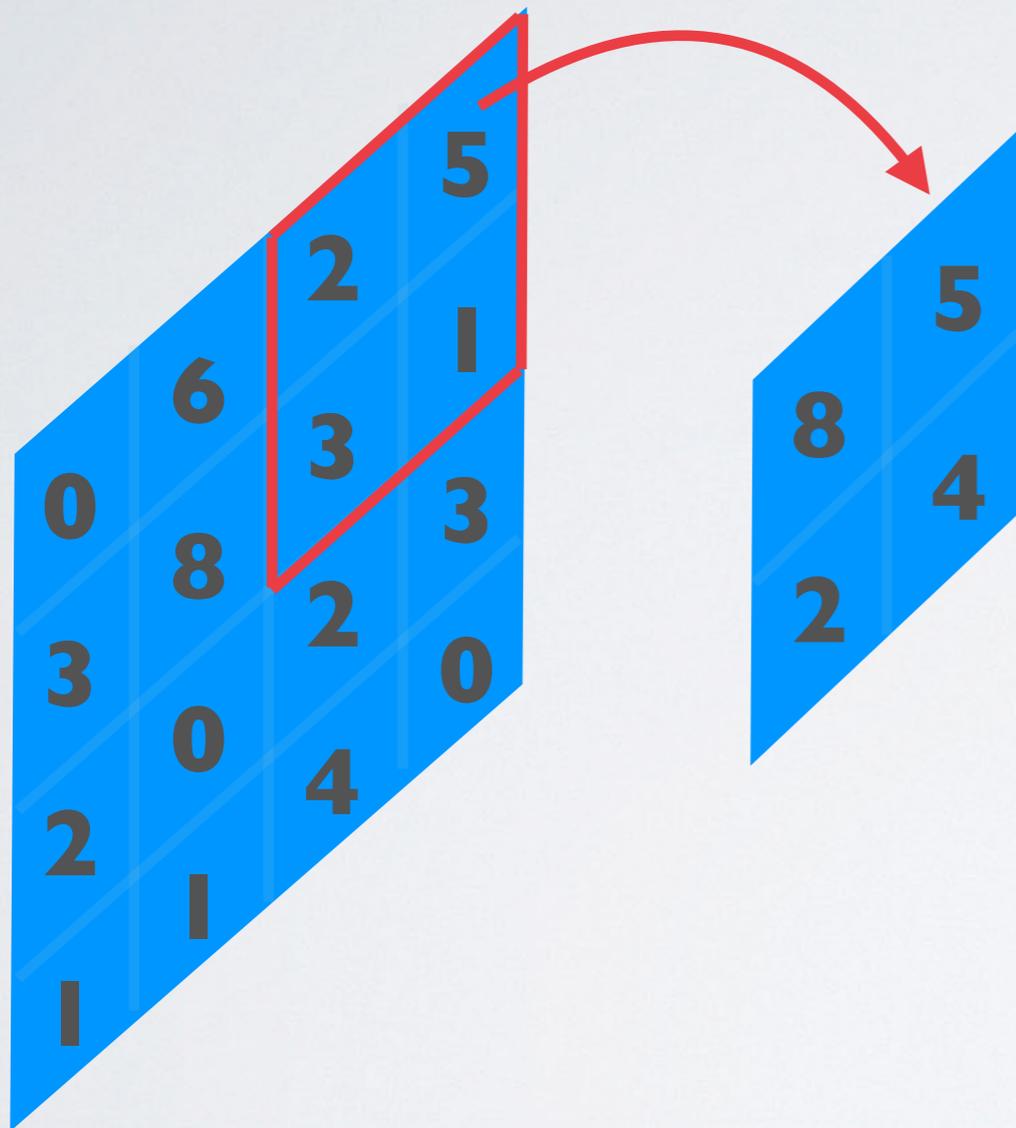
“Feature map” **remove negatives**



POOLING

POOLING

moving to lower resolutions



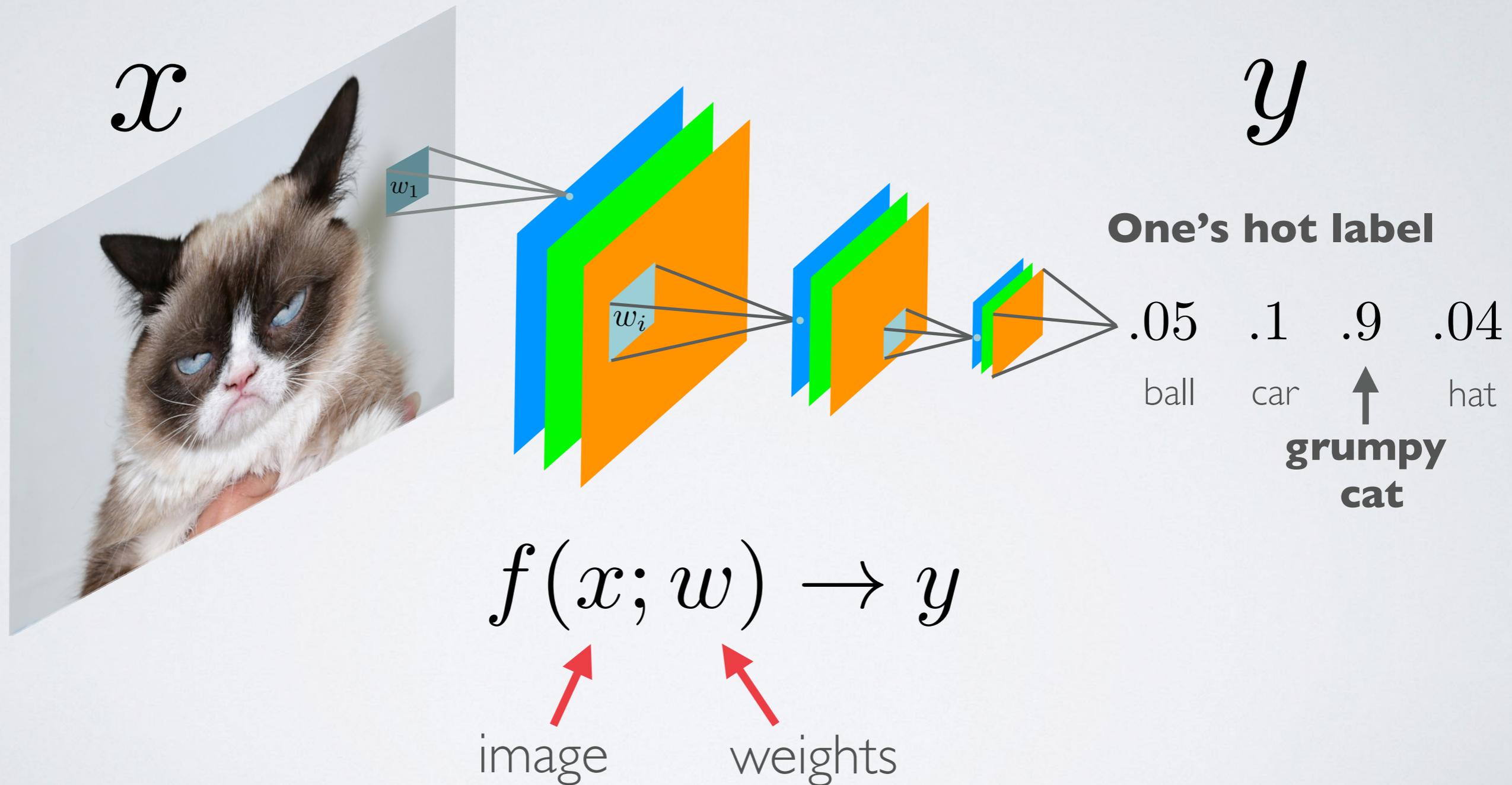
Max pooling

Average pooling

Strided convolution
(aka downsampling)

THAT'S ALL YOU NEED!

(sort of)



MNIST

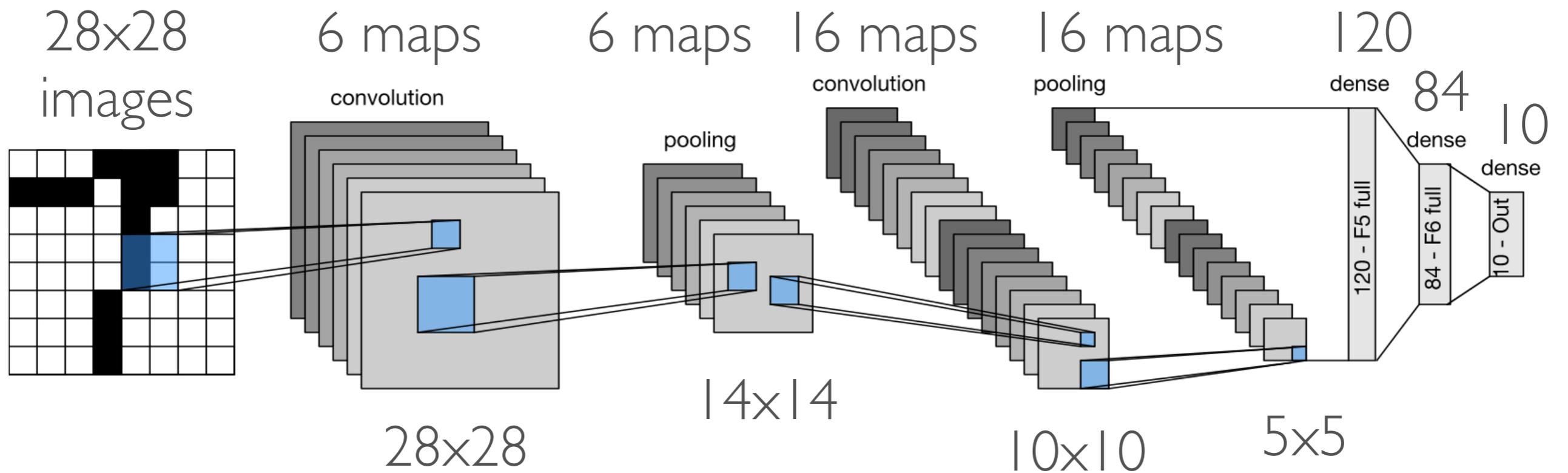
Modified National Institute of Standards and Technology



LENET

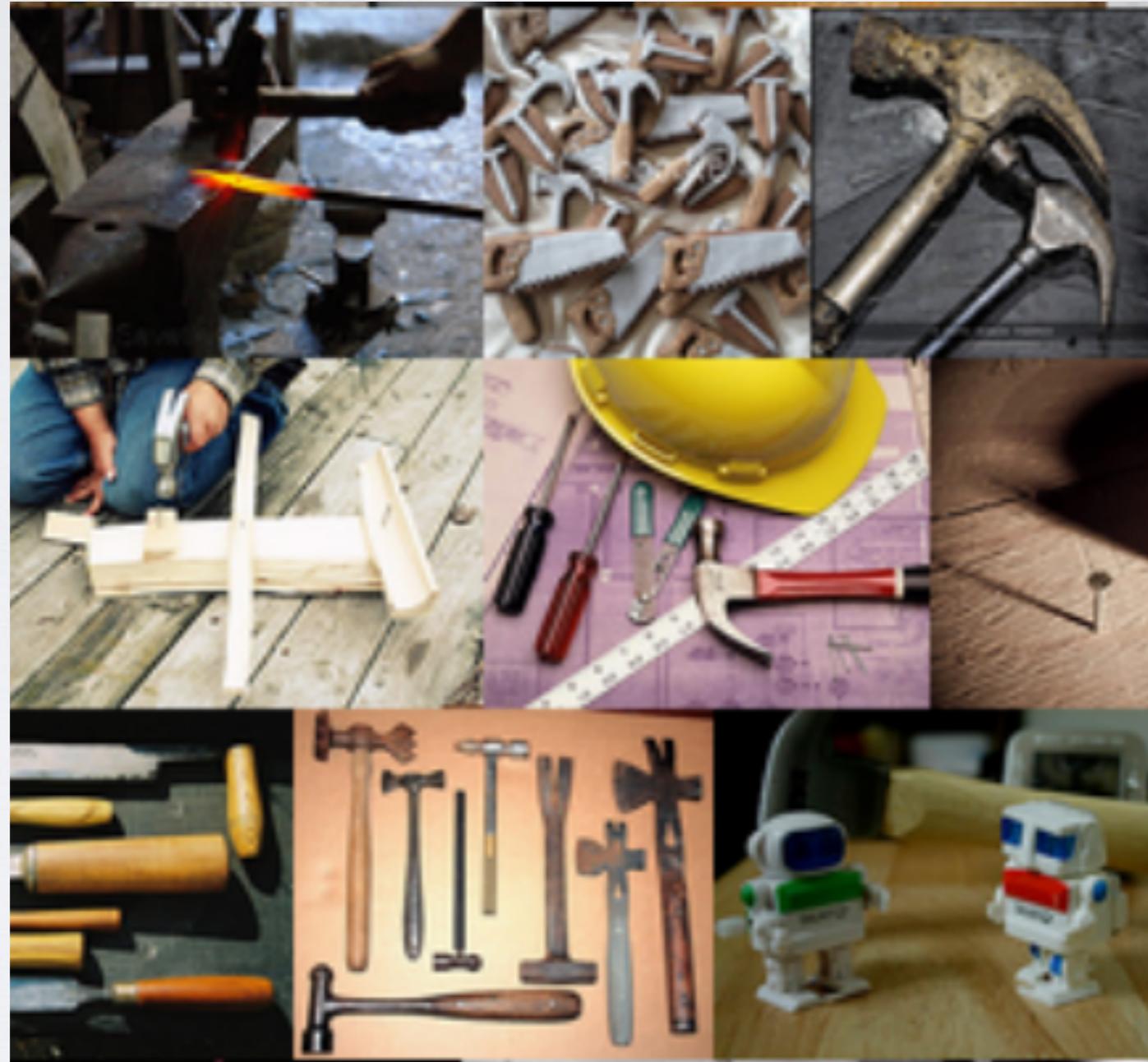
LeCun et al, 1998

The first successful convnet



IMAGENET DATASET

2010

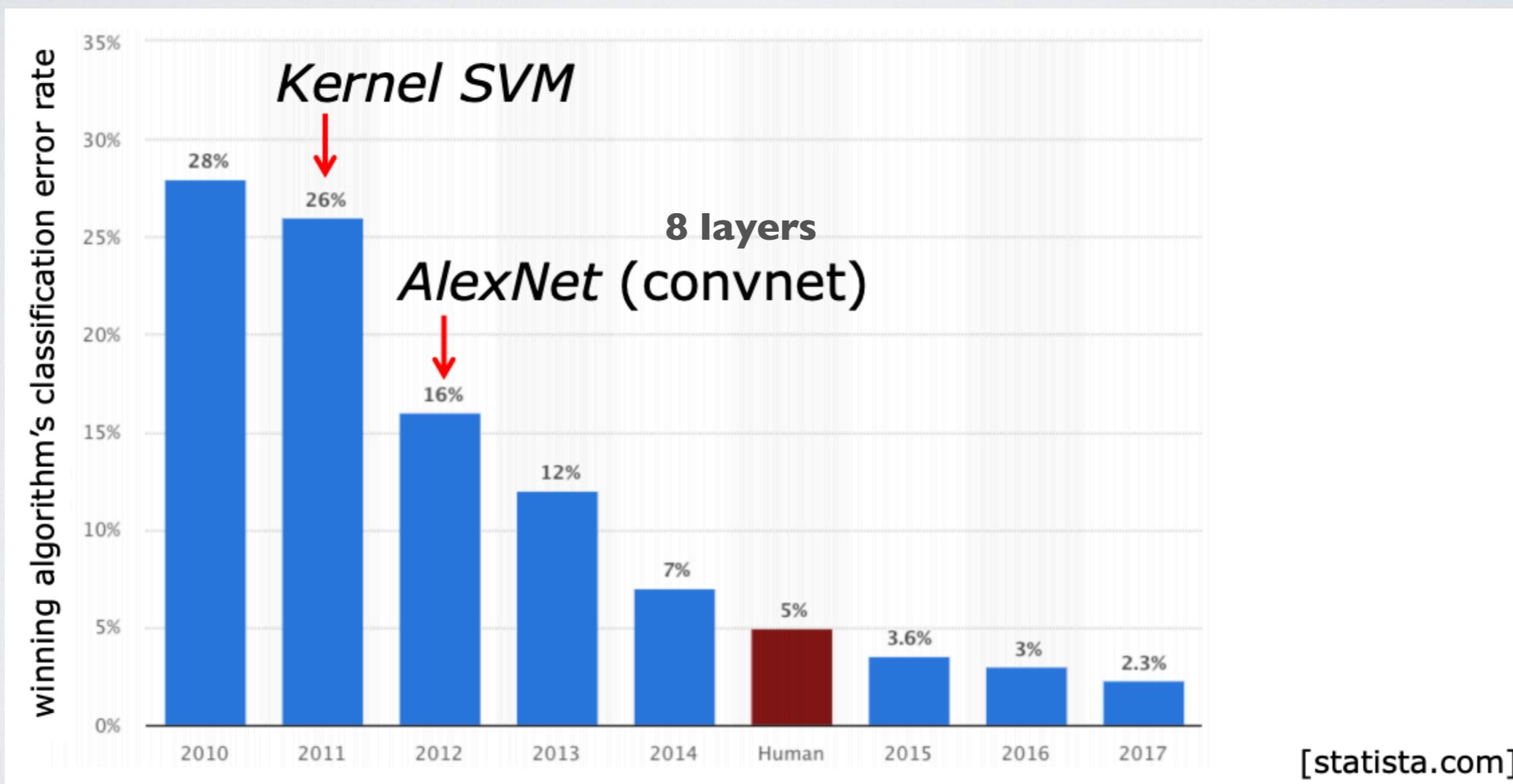


1 K classes

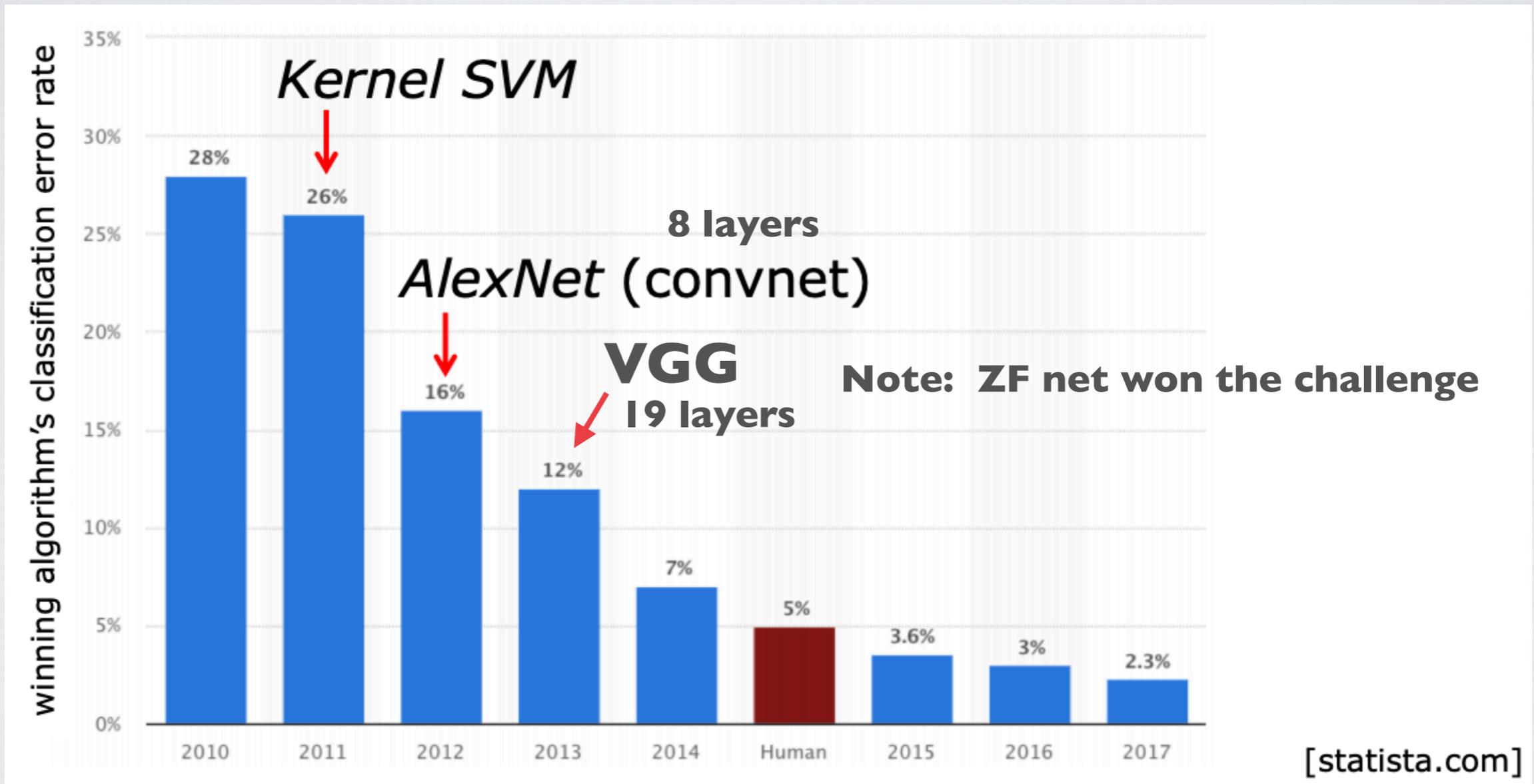
1 M Images

Usually 224x224 resolution

IMAGENET HISTORY

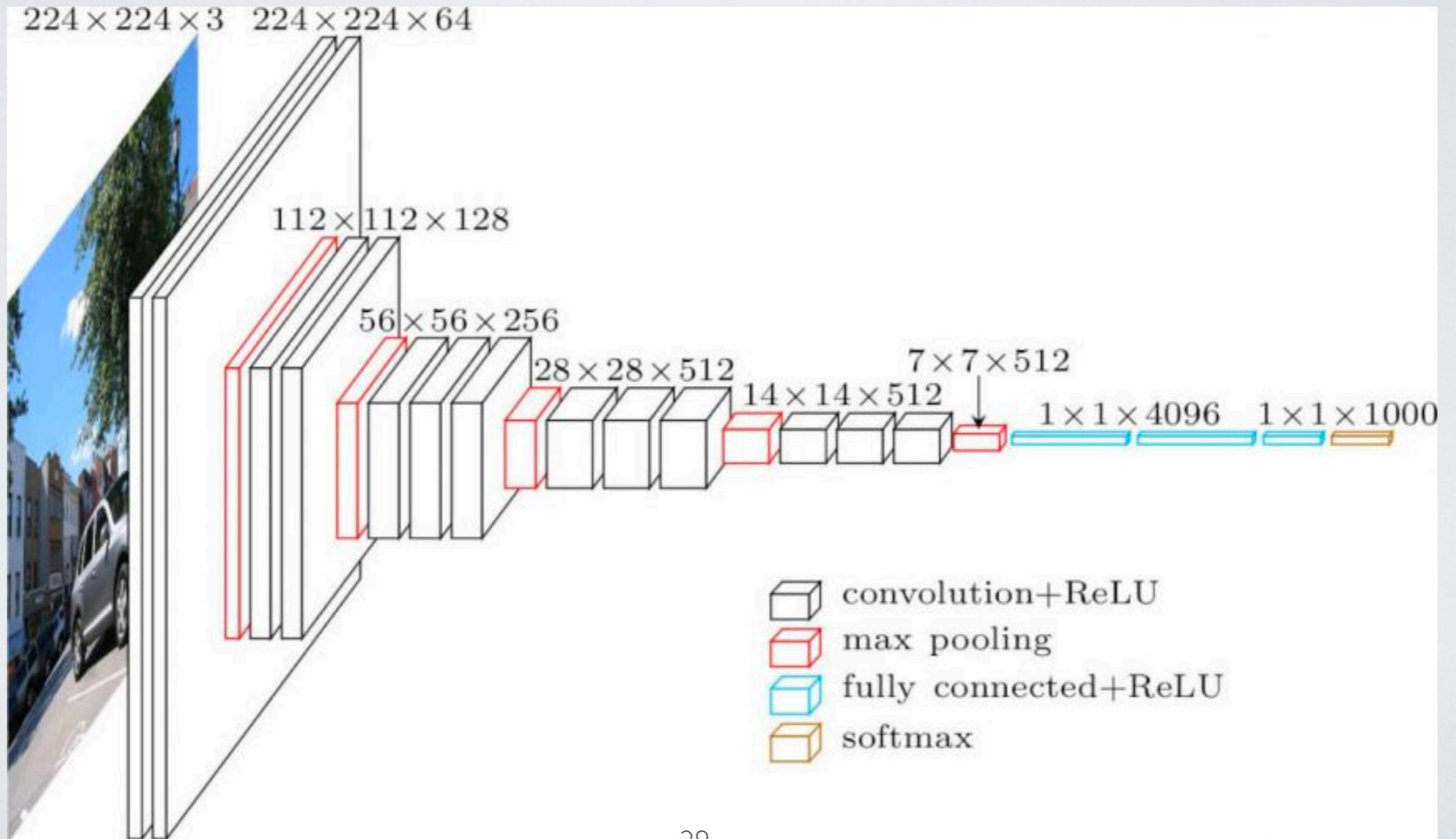


IMAGENET HISTORY

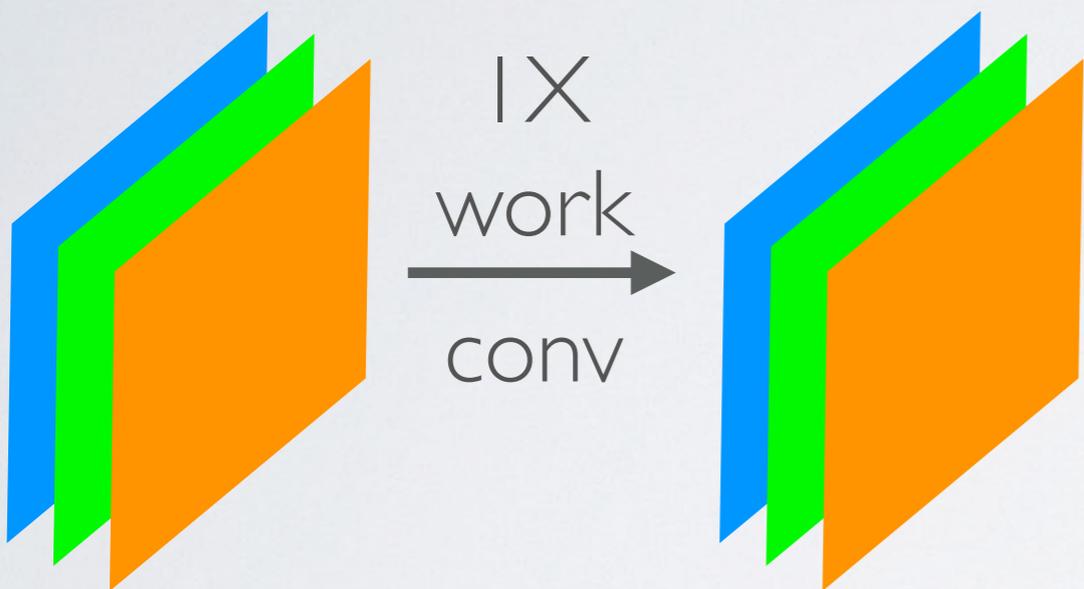


VGG

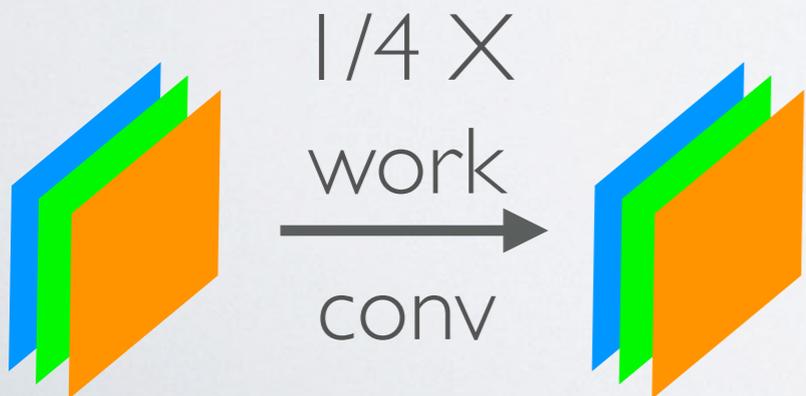
Visual graphics group, Oxford



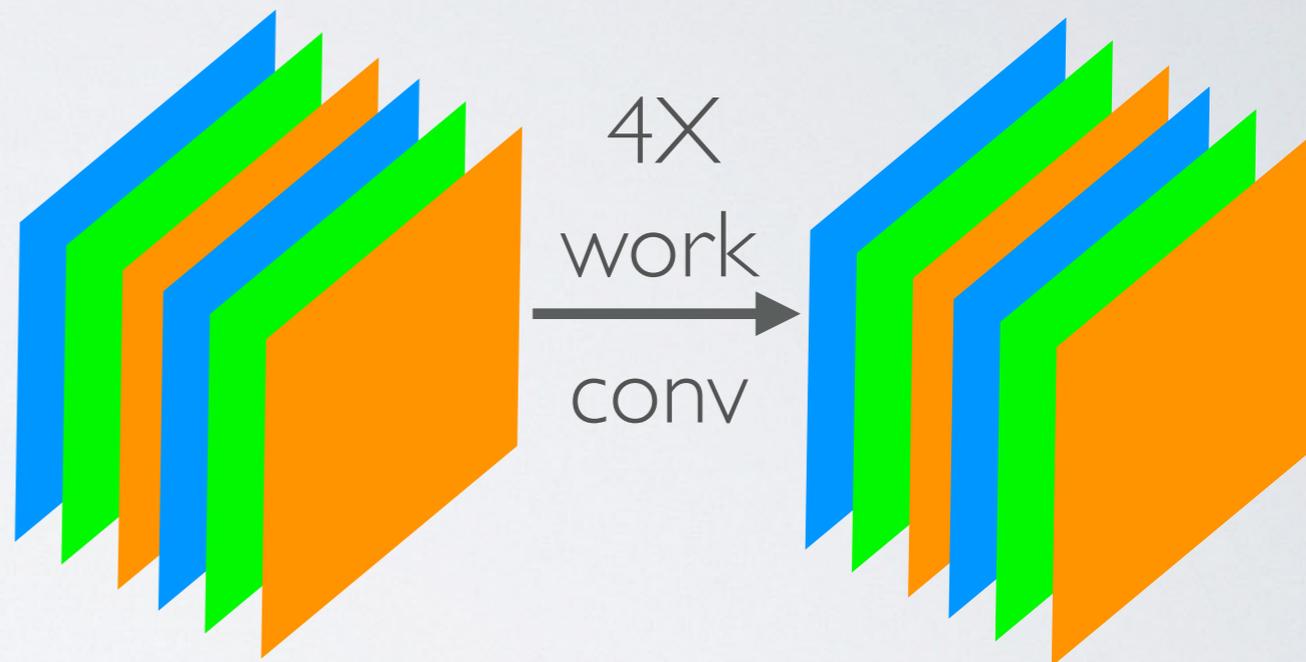
2D CONV NETS



Half the resolution

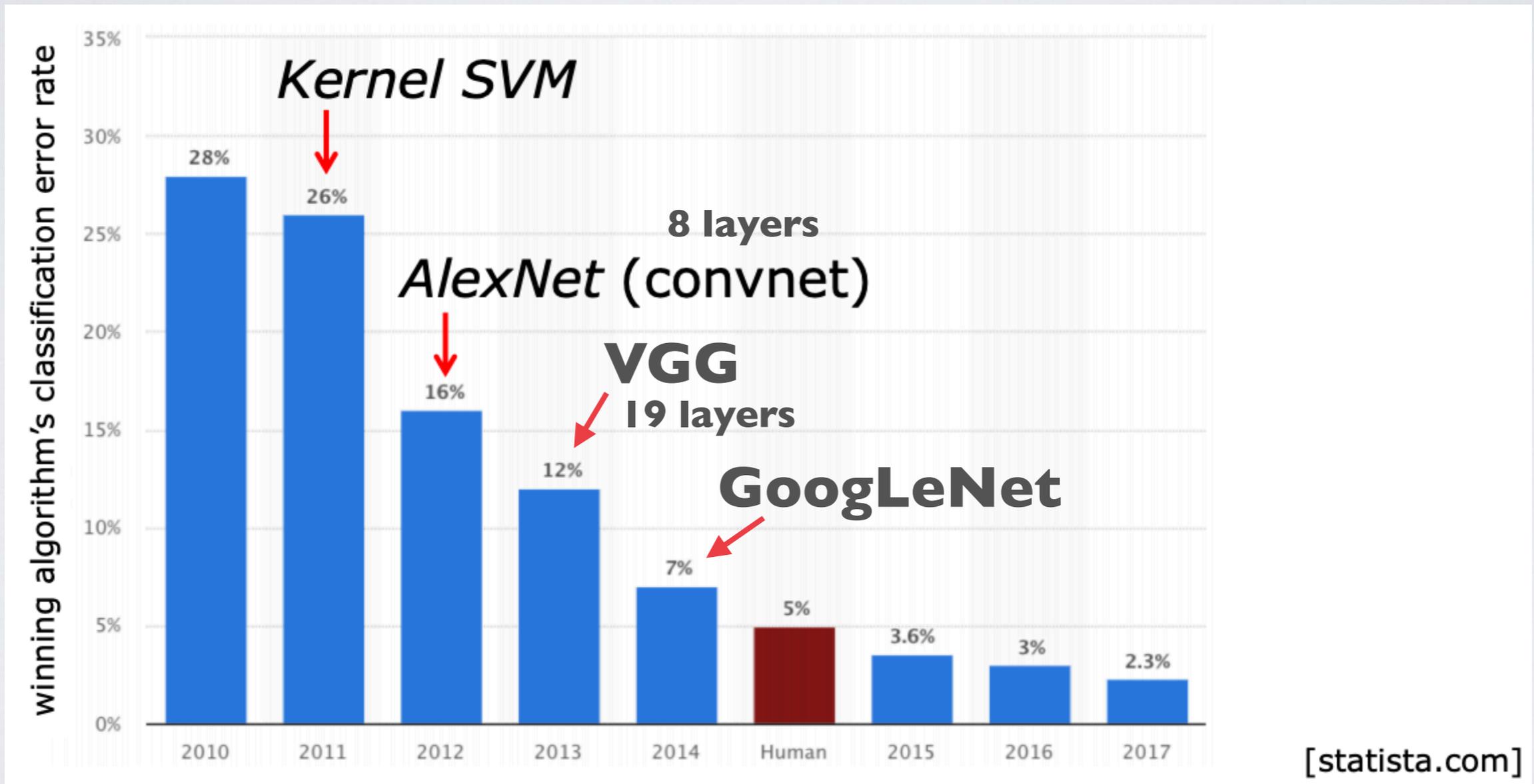


Twice as many maps

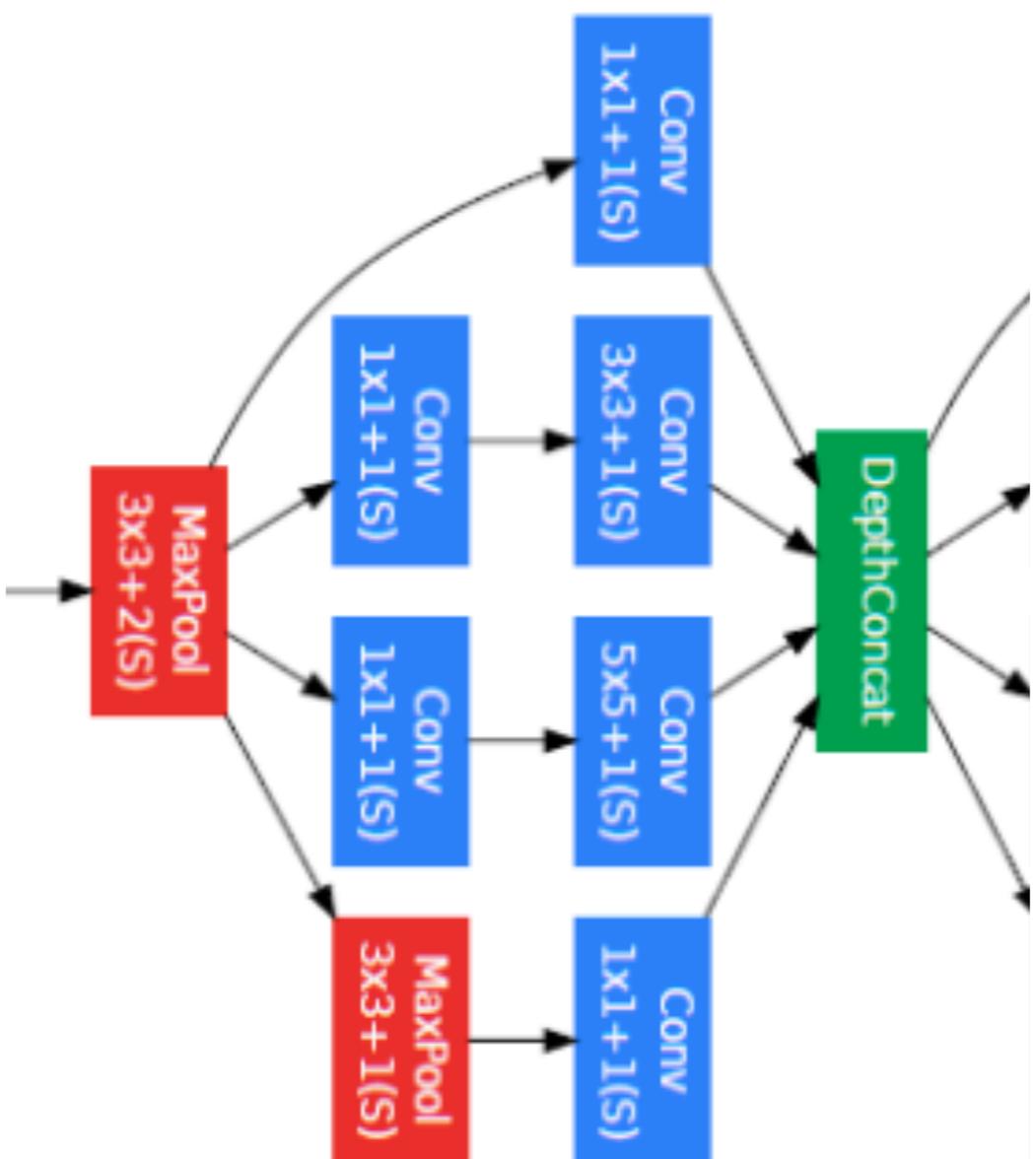


$(2X \text{ feature maps}) + (1/2 \text{ resolution})$
=
Same number of operations

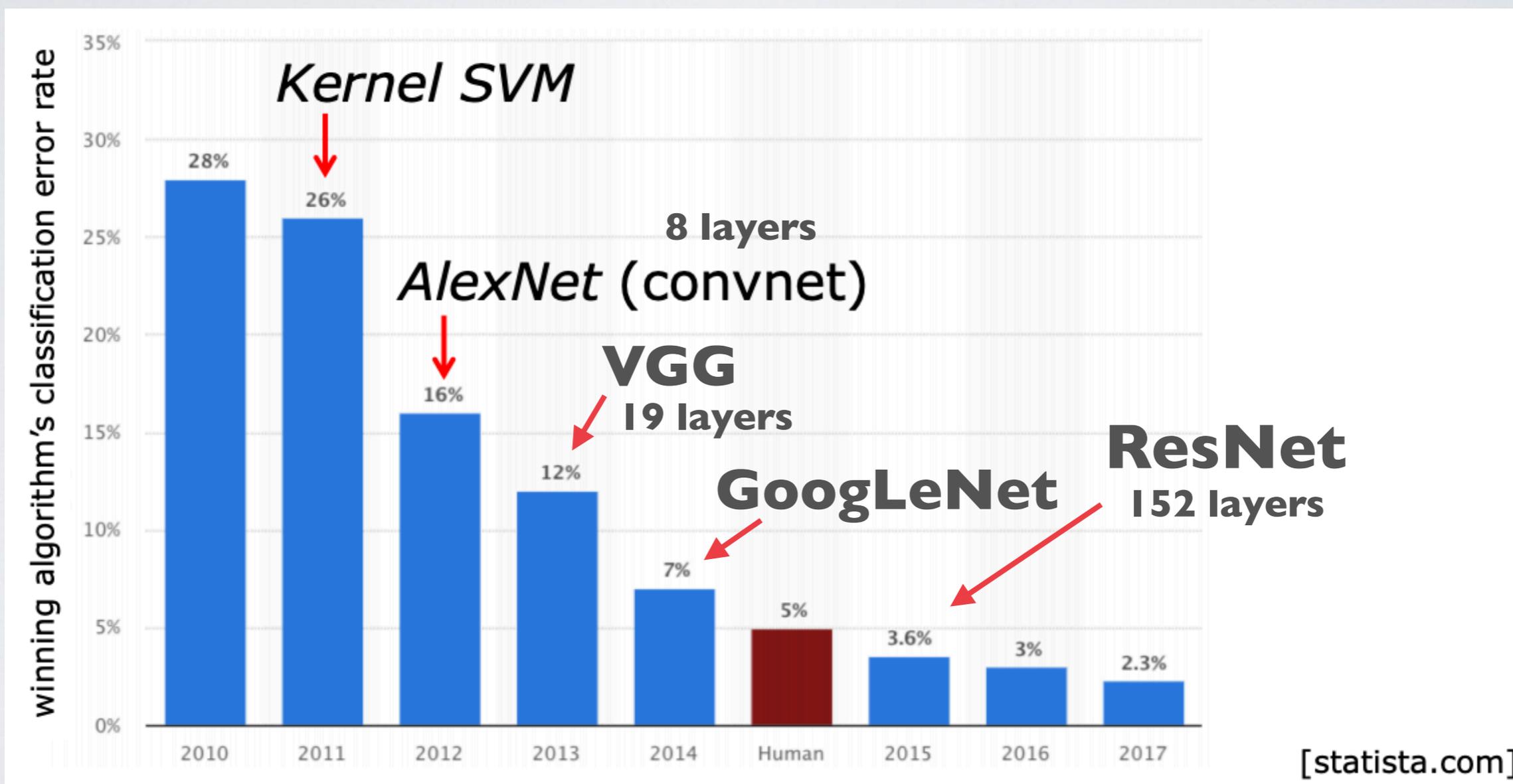
IMAGENET HISTORY



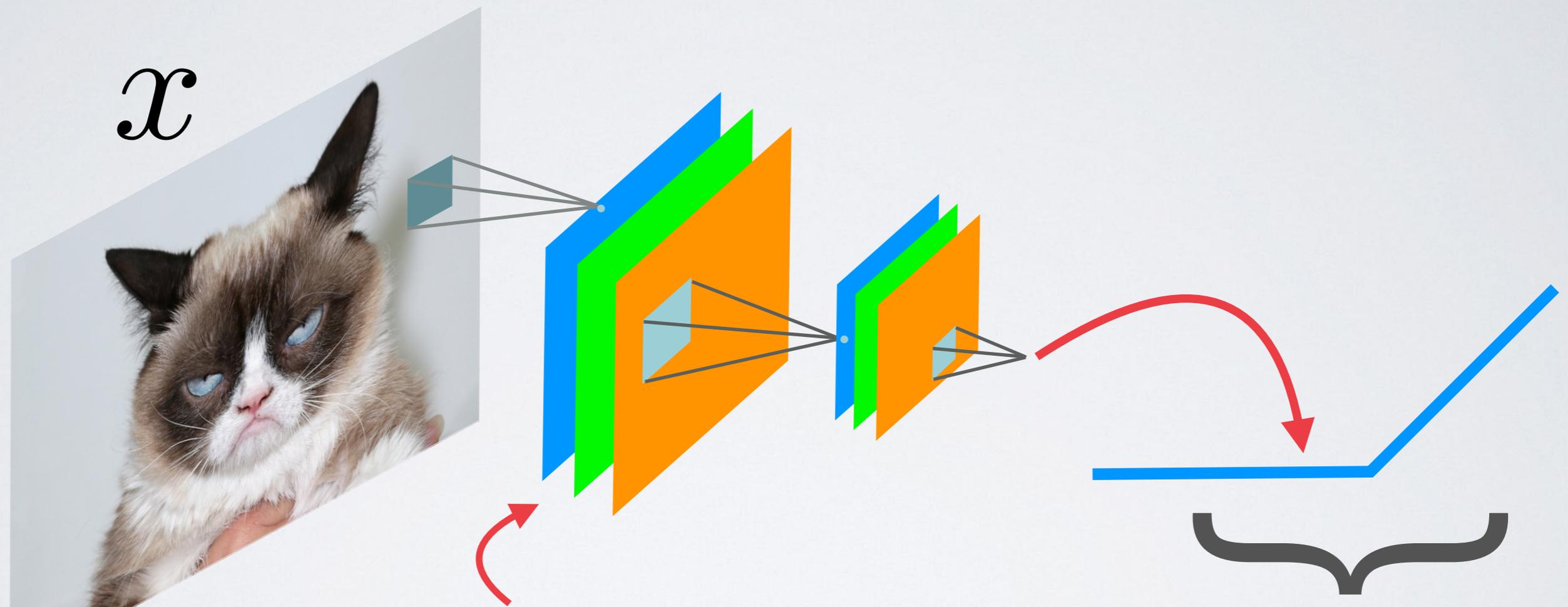
GOOGLENET



IMAGENET HISTORY

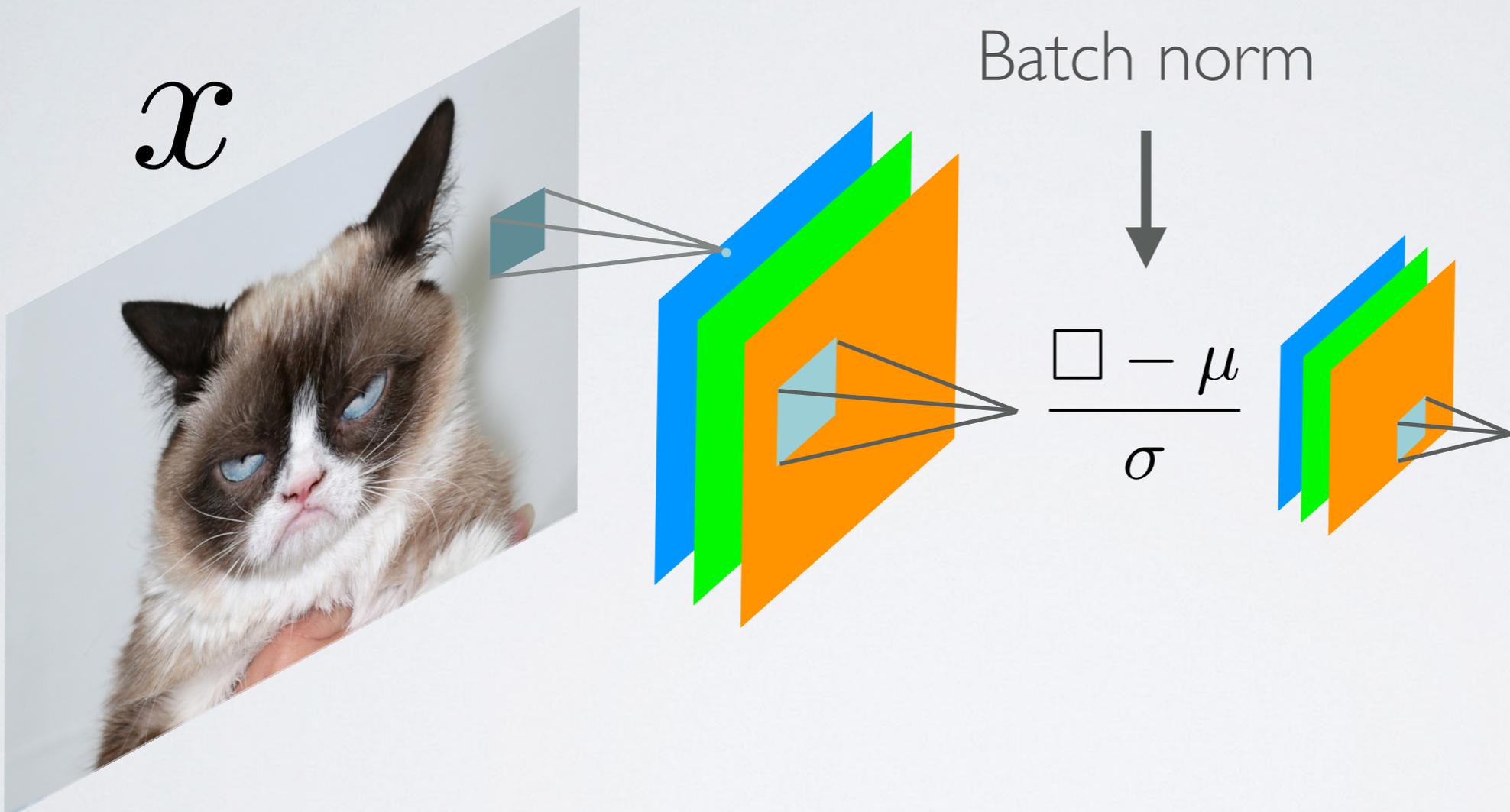


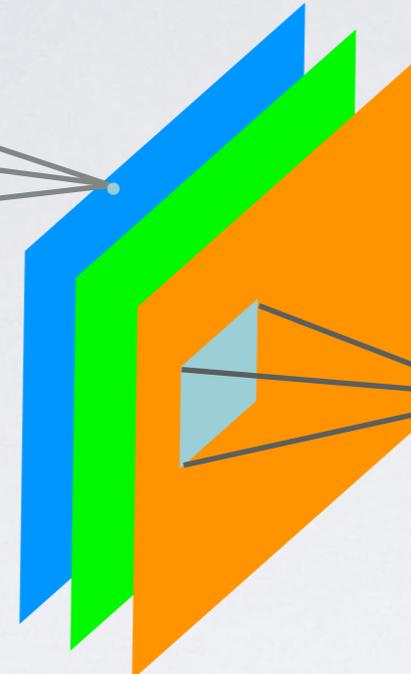
BATCH NORMALIZATION



Changes to shallow layers change
mean of deeper activations
“Internal covariant shift”

BATCH NORMALIZATION

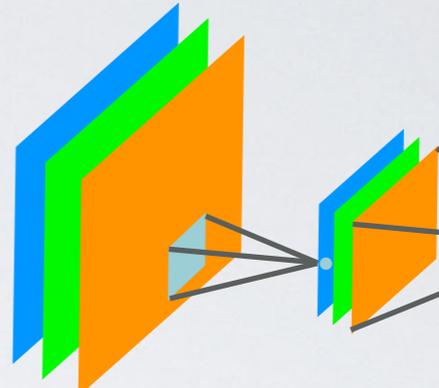




Batch norm



$$\frac{\square - \mu}{\sigma}$$



One's hot label

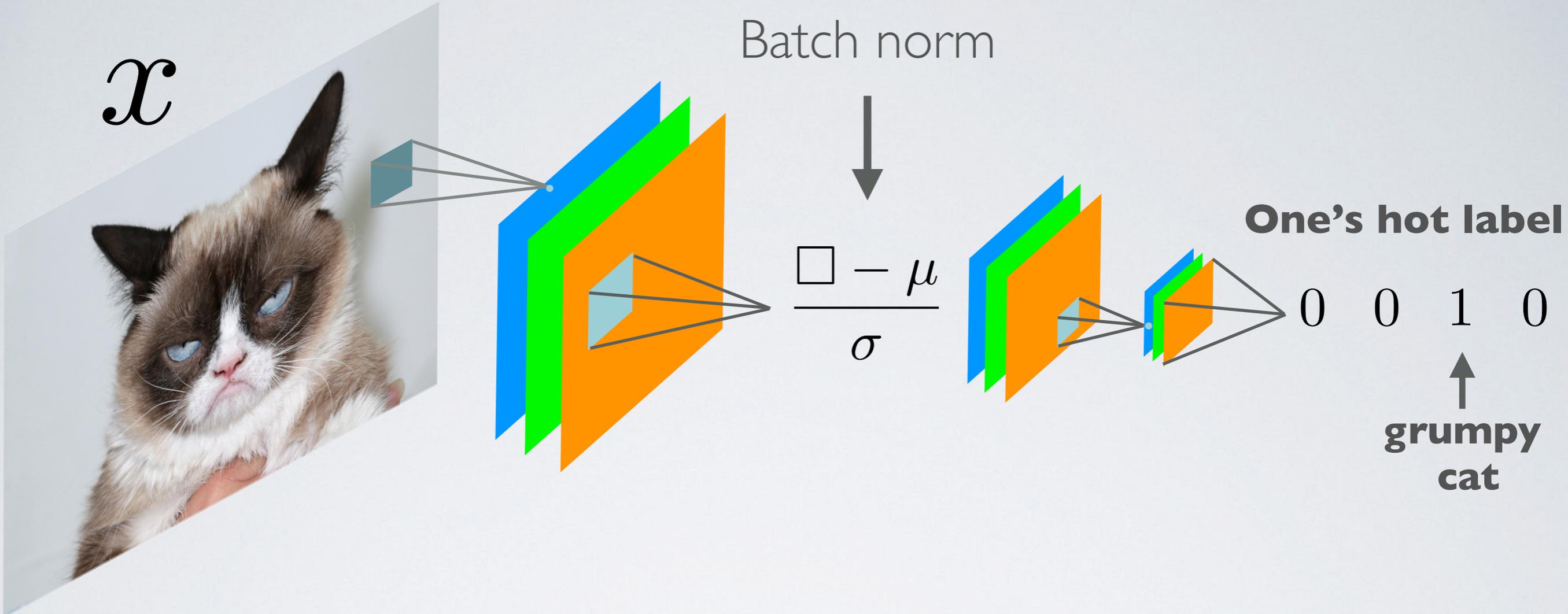
0 0 1 0

↑
grumpy
cat

Mean intensity of feature map over training batch

$$\gamma \frac{z - \mu}{\sigma} + \eta$$

STD over training batch



Trainable parameter

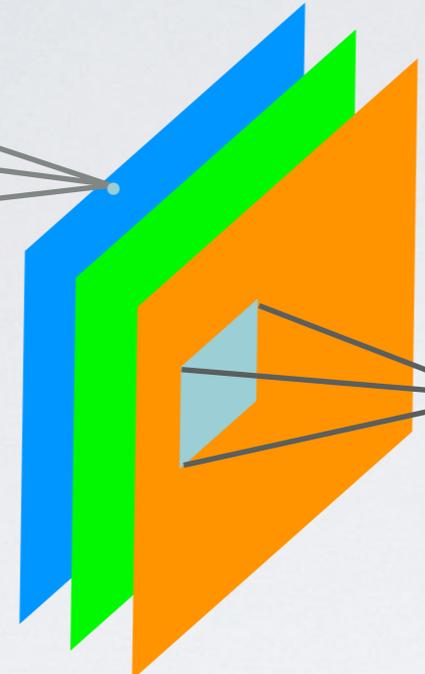
$$\gamma \frac{z - \mu}{\sigma} + \eta$$

Trainable parameter



X

Batch norm



$$\frac{\square - \mu}{\sigma}$$



One's hot label

0 0 1 0

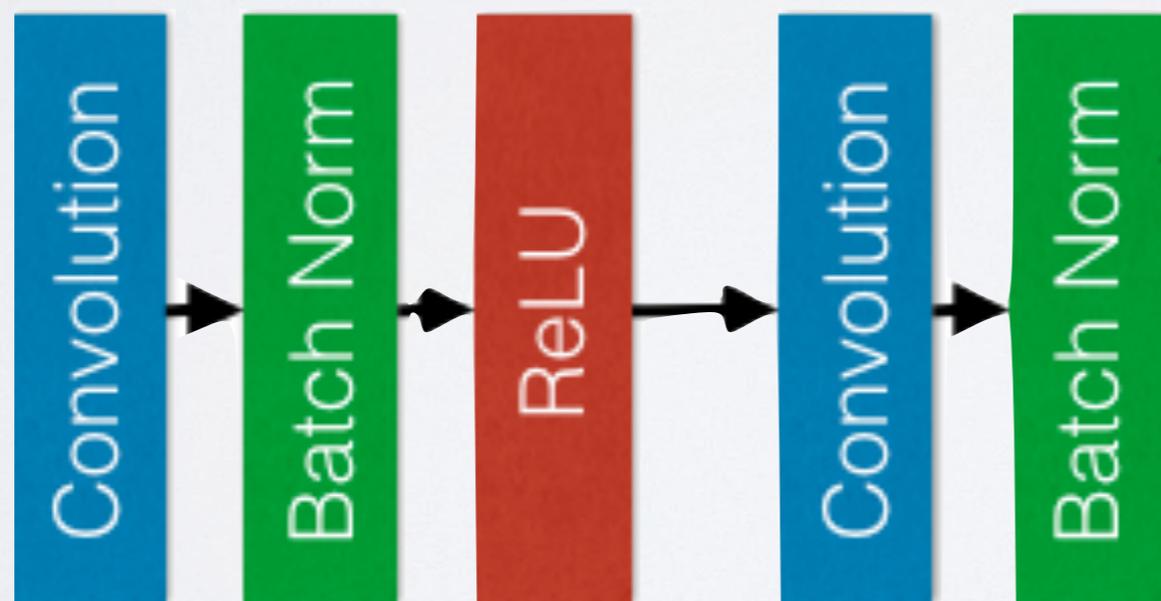
↑
grumpy
cat

Test time

Replace with "running mean"
and "running variance"

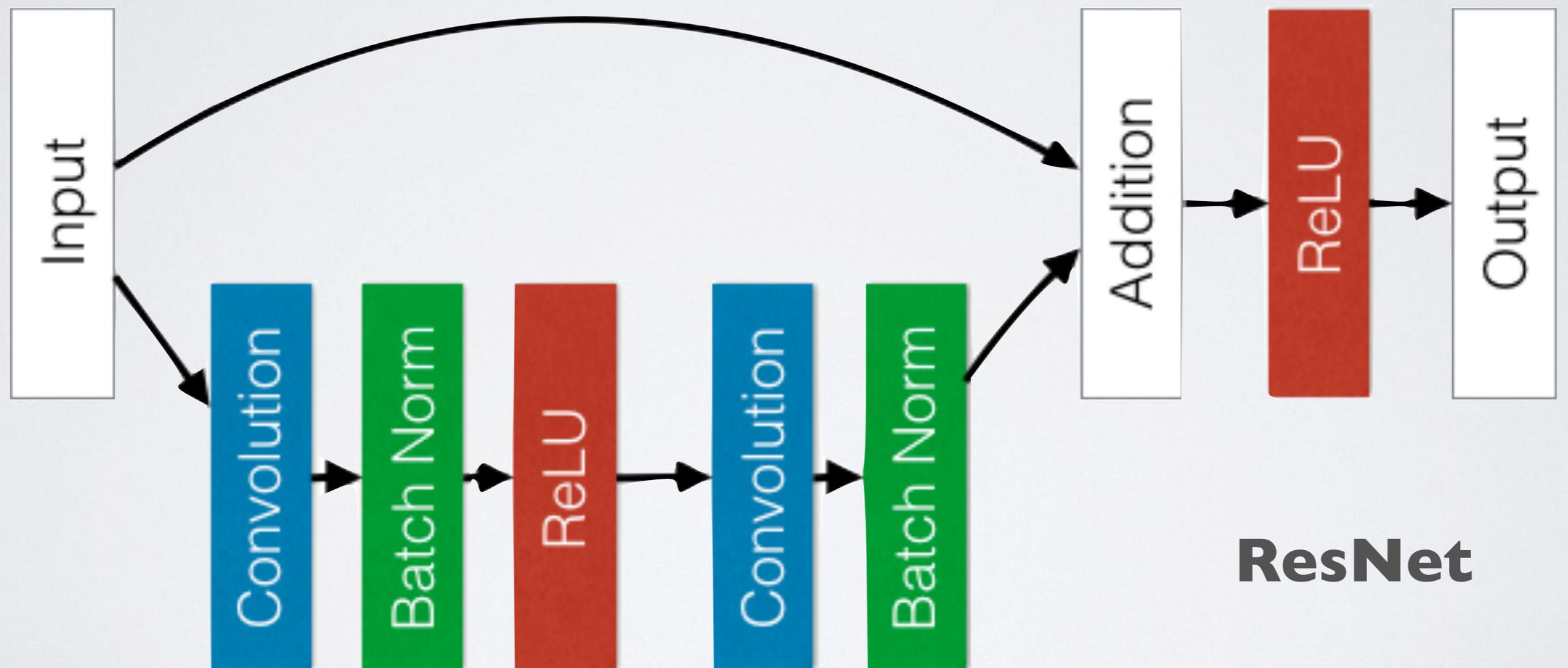
$$\gamma \frac{z - \mu}{\sigma} + \eta$$

RESIDUAL CONNECTIONS



**“VGG”-like
Net**

RESIDUAL CONNECTIONS



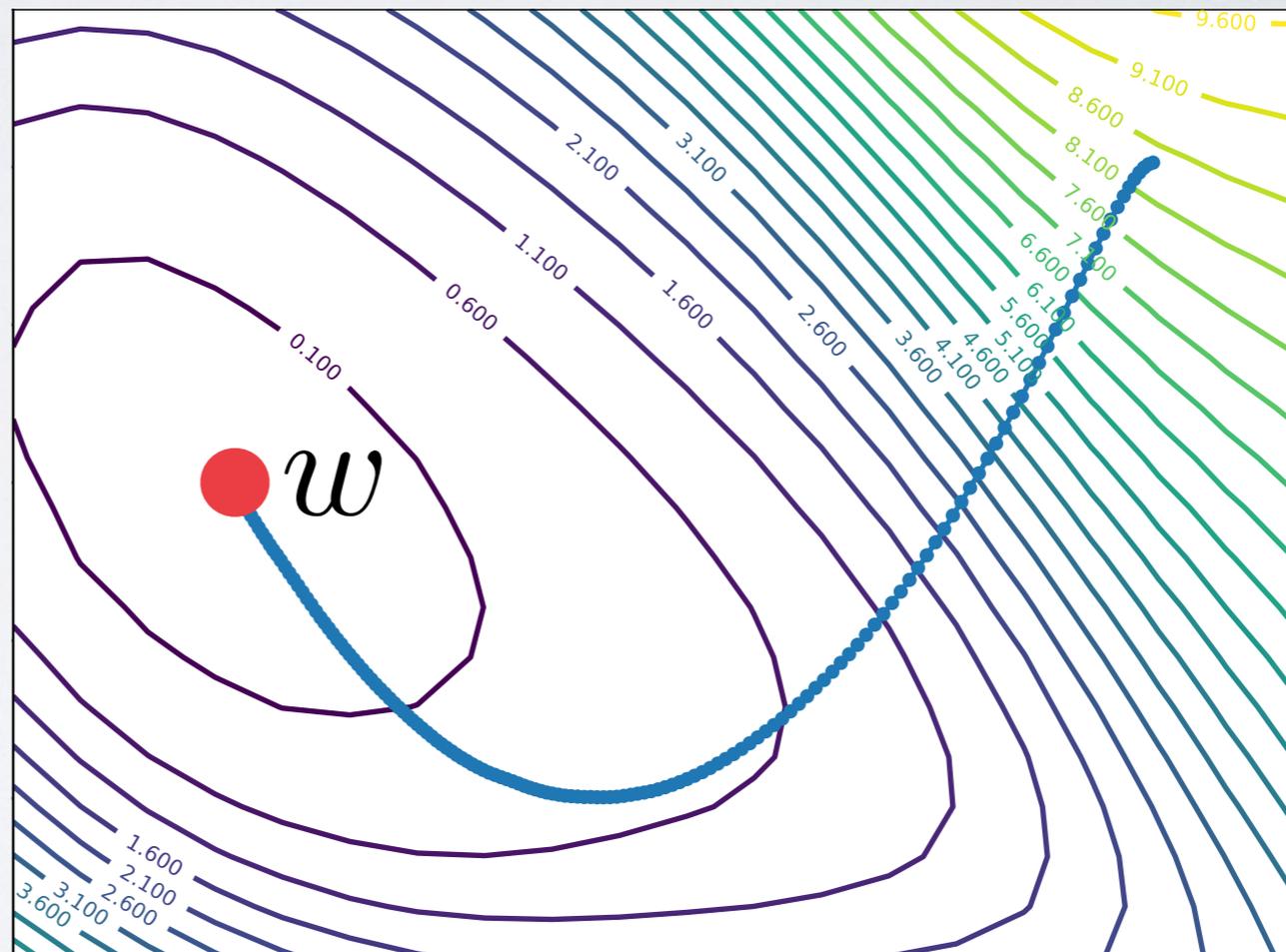
RESNET-34



WHY ARE SKIP
CONNECTIONS GOOD?

VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION

Step 1:
Find
minimizer

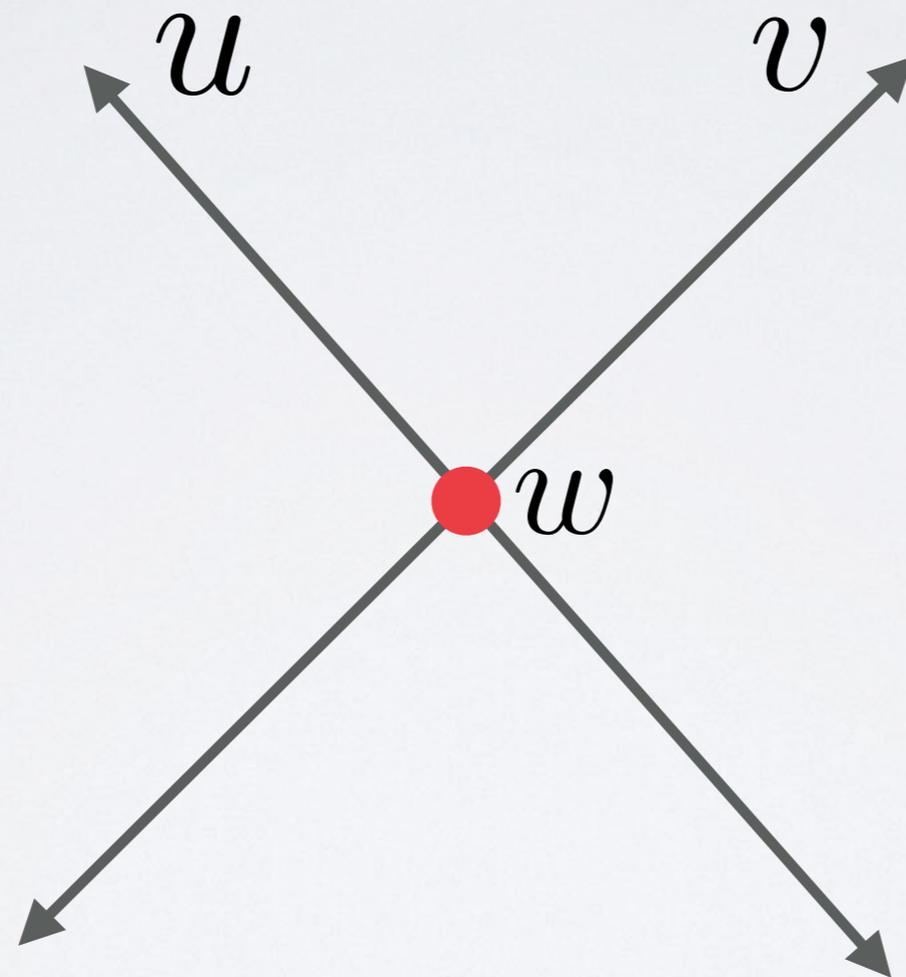


30 million dimensions

VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION

**Step 2:
Random
directions**

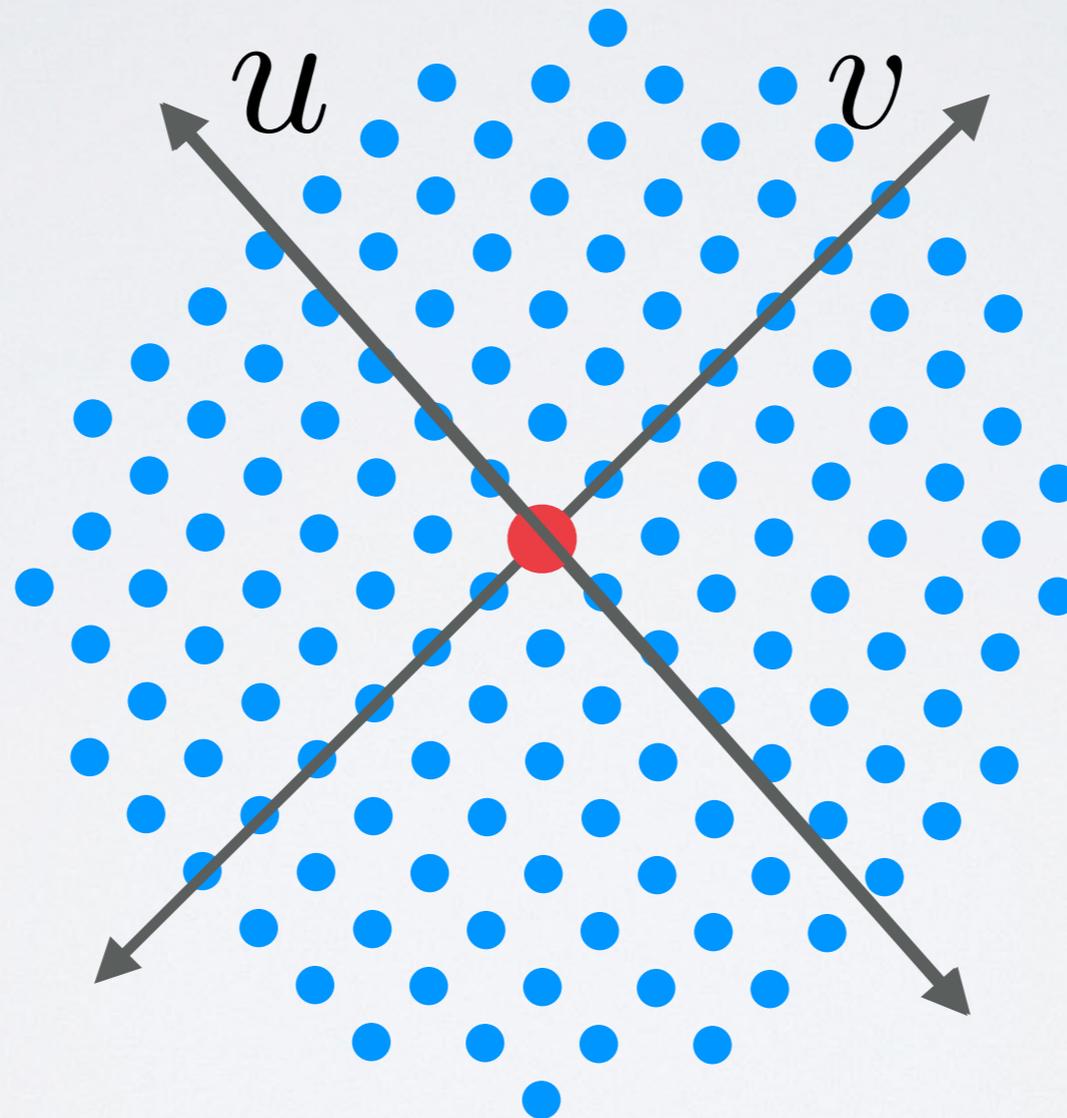
u, v



VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION

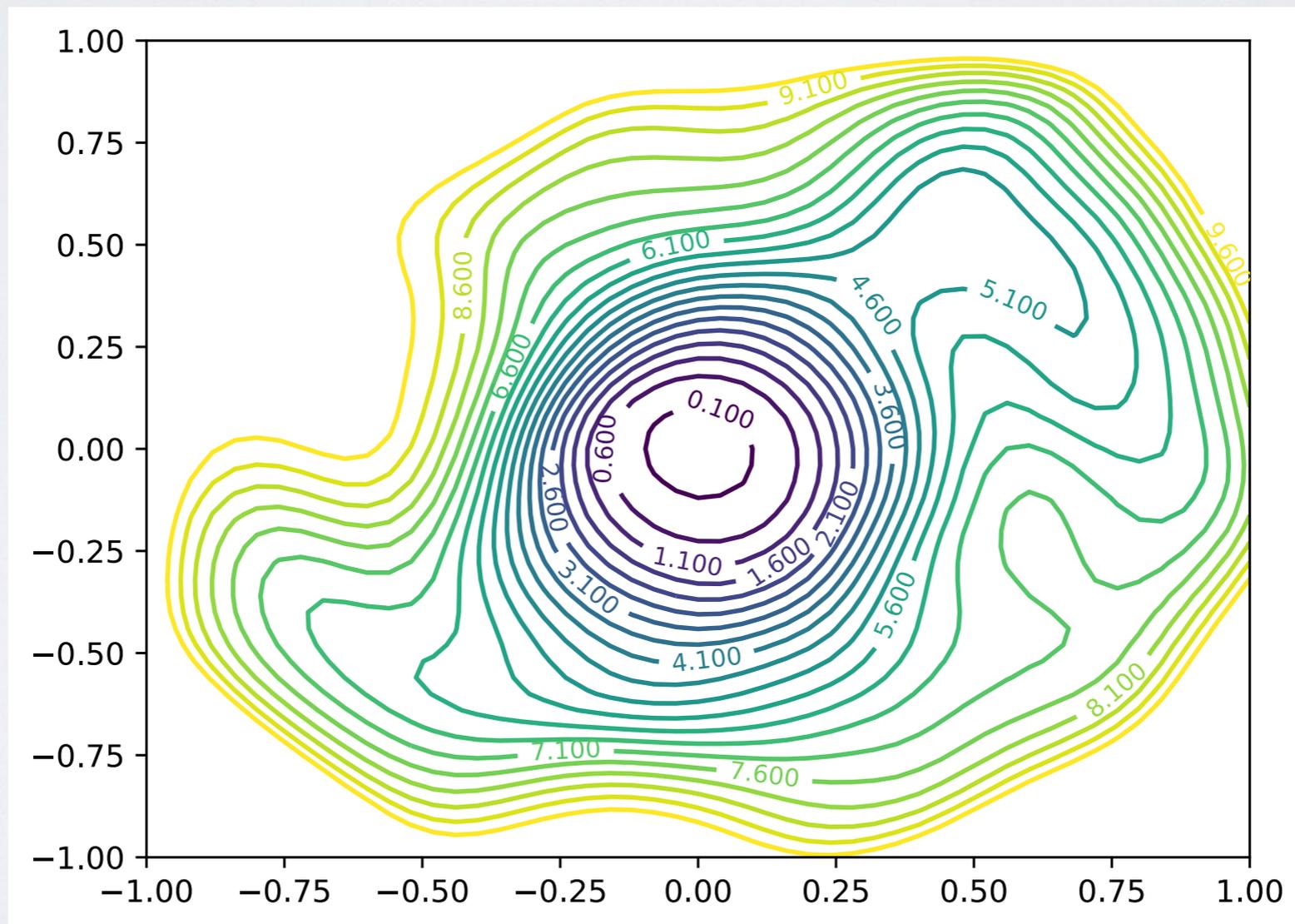
Step 3

Plot



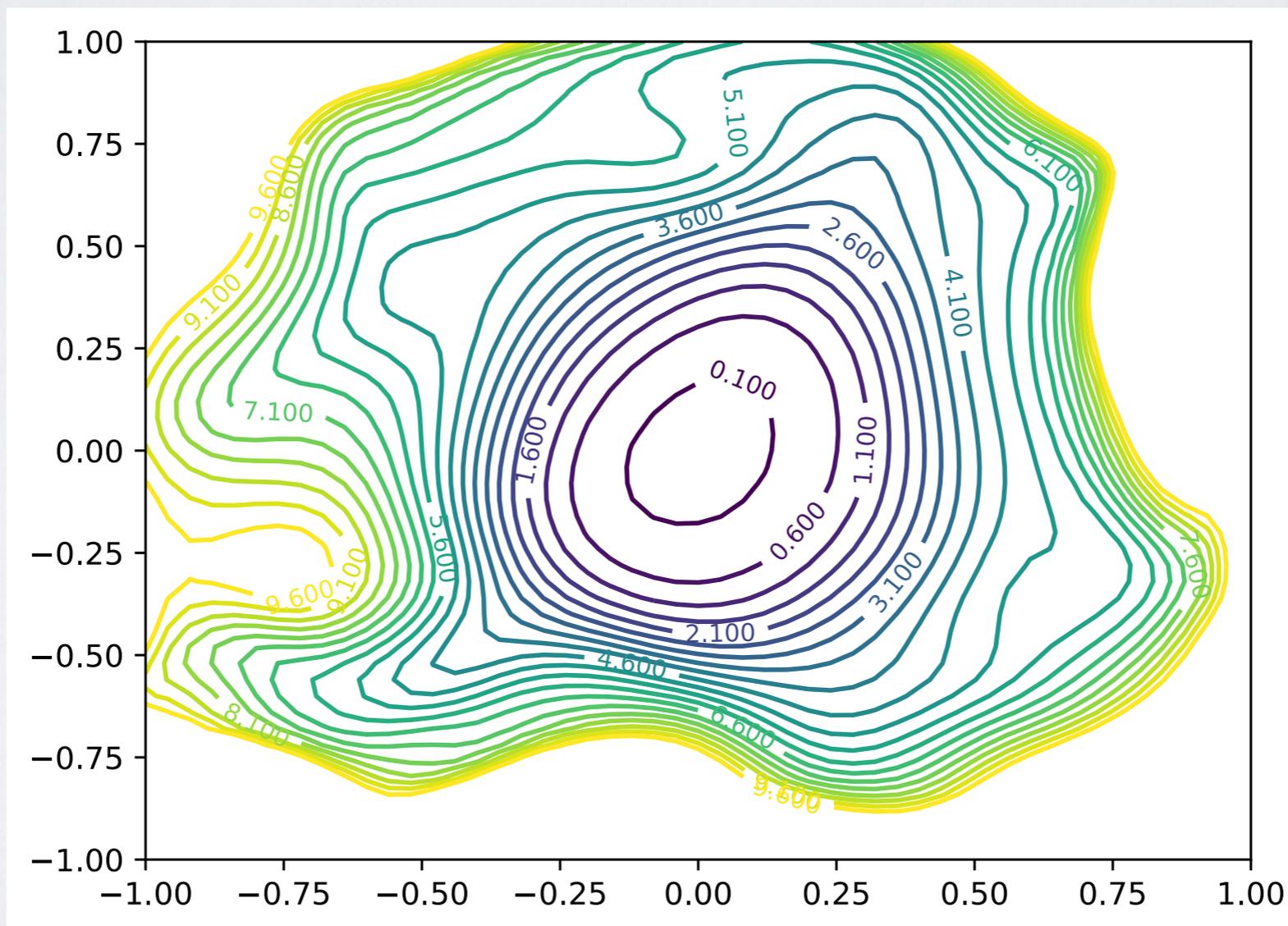
CHAOTIC TRANSITIONS

VGG-9



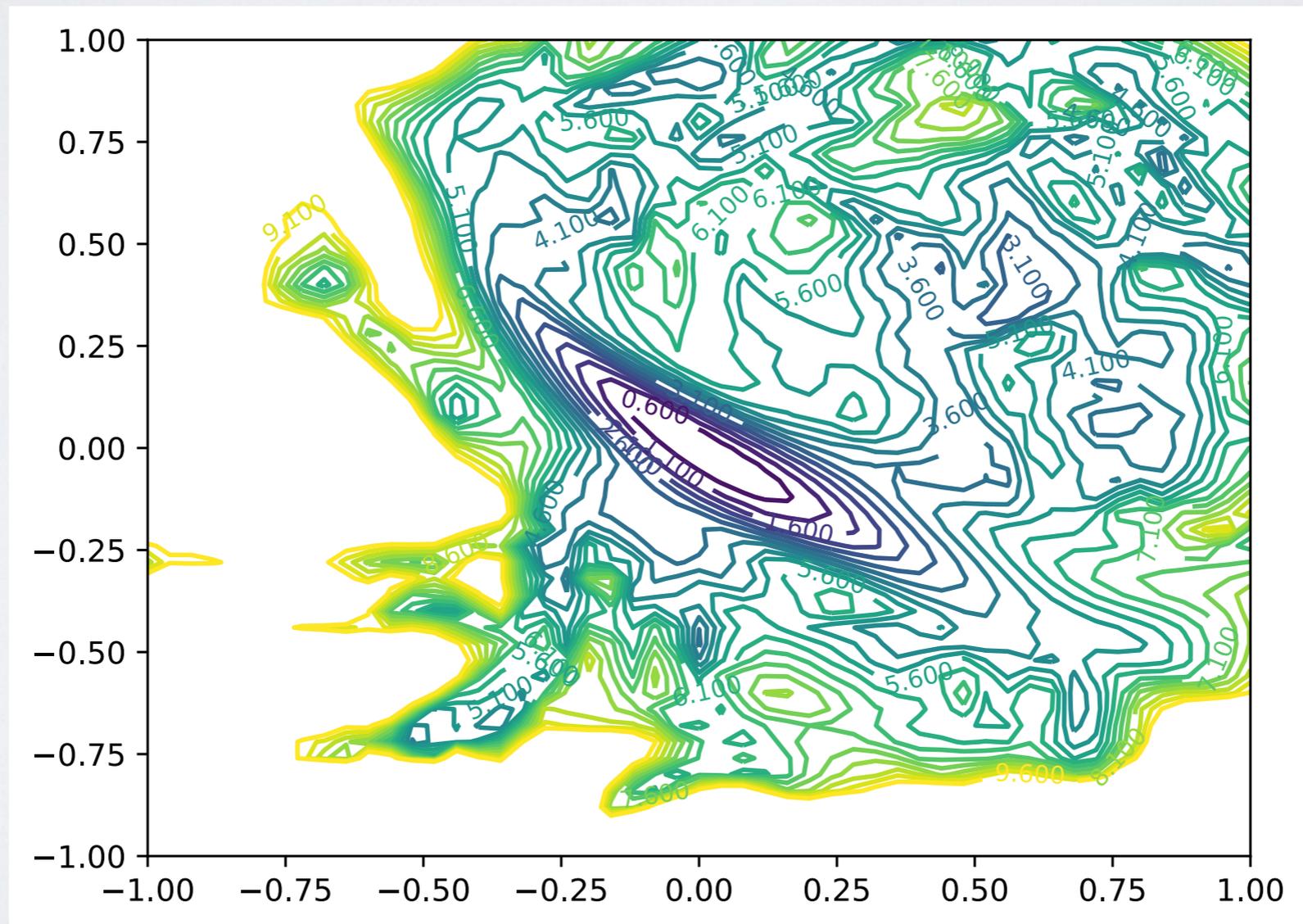
CHAOTIC TRANSITIONS

VGG-20



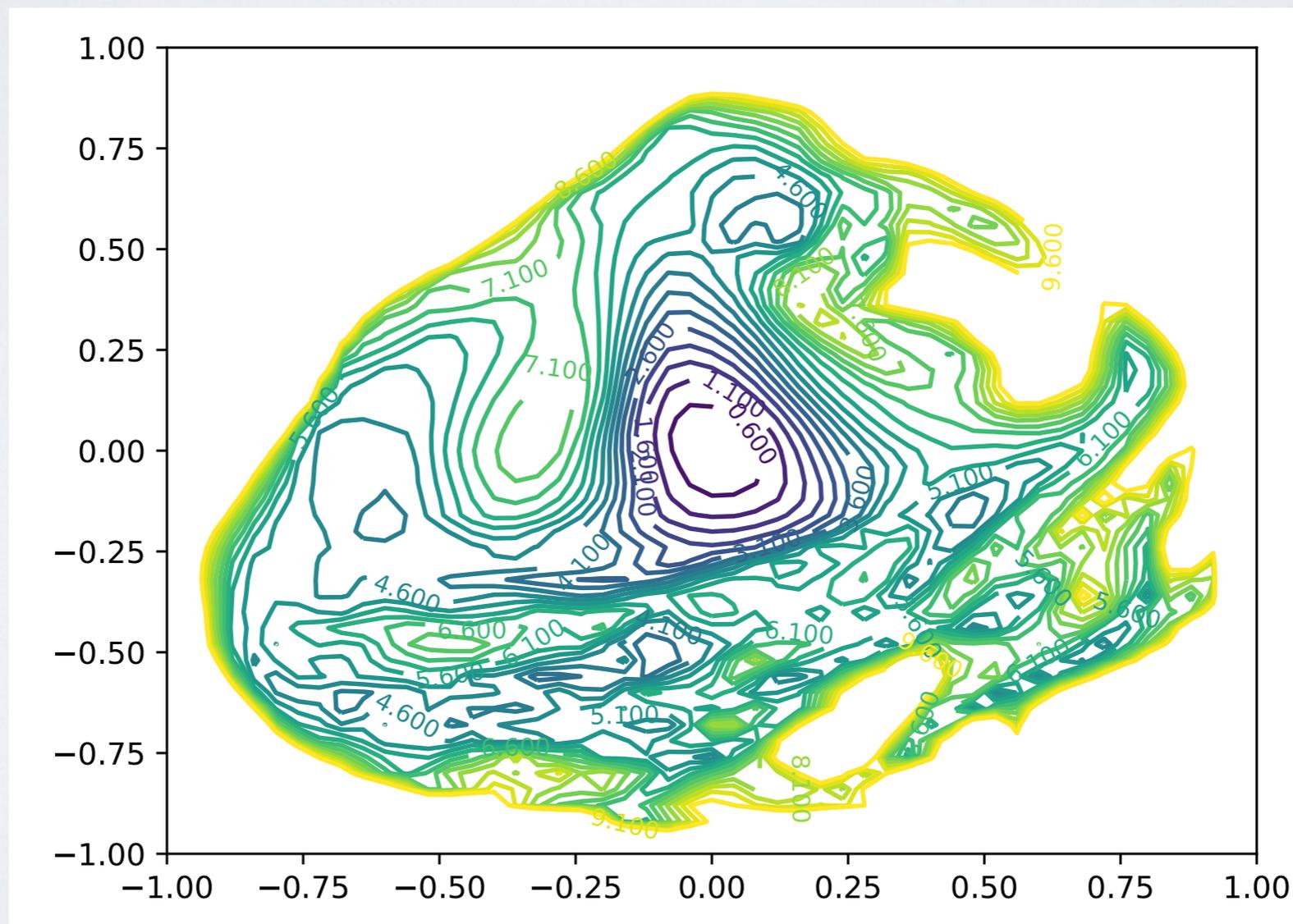
CHAOTIC TRANSITIONS

VGG-56



CHAOTIC TRANSITIONS

VGG-110

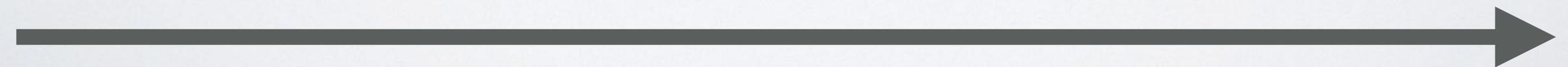
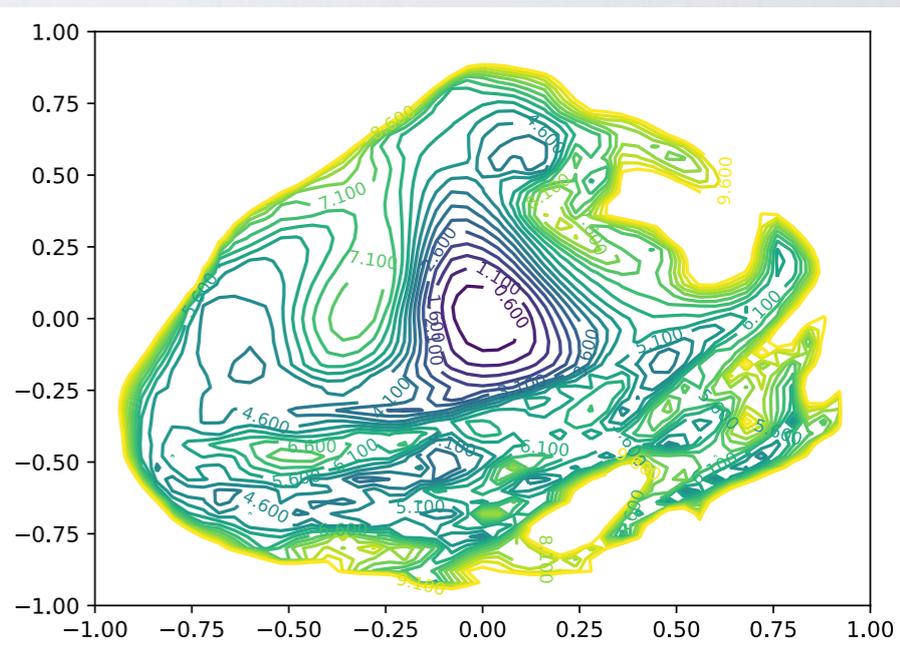
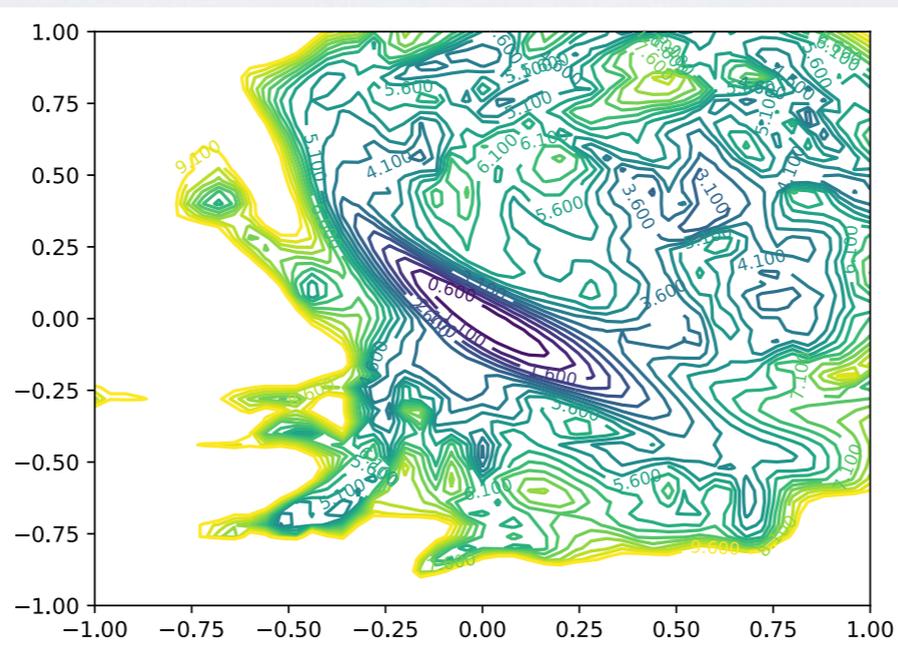
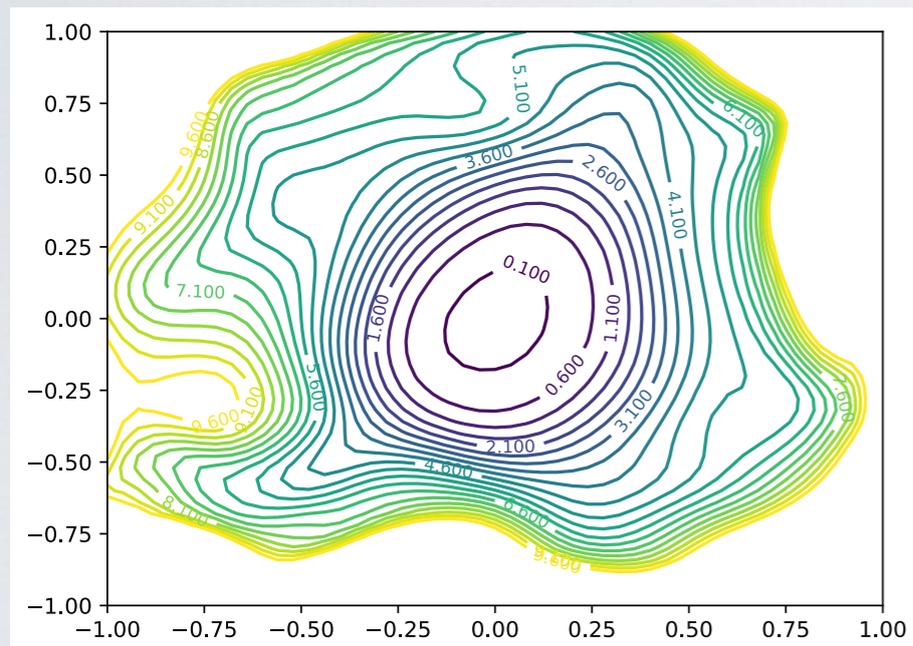


CHAOTIC TRANSITIONS

VGG-20

VGG-56

VGG-110



Convexity

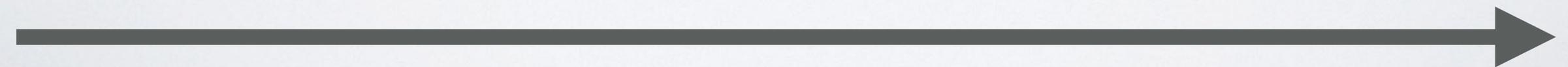
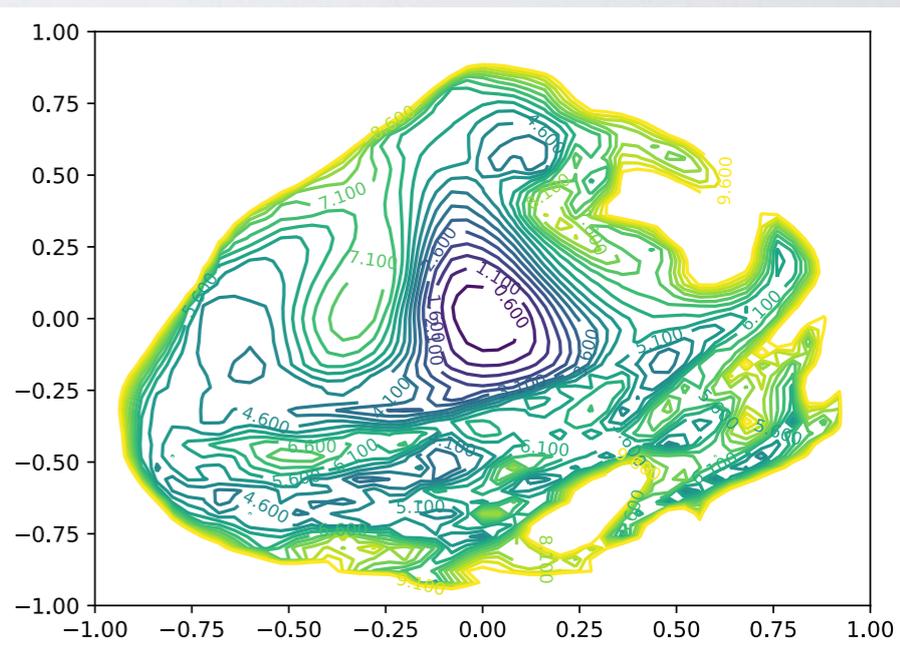
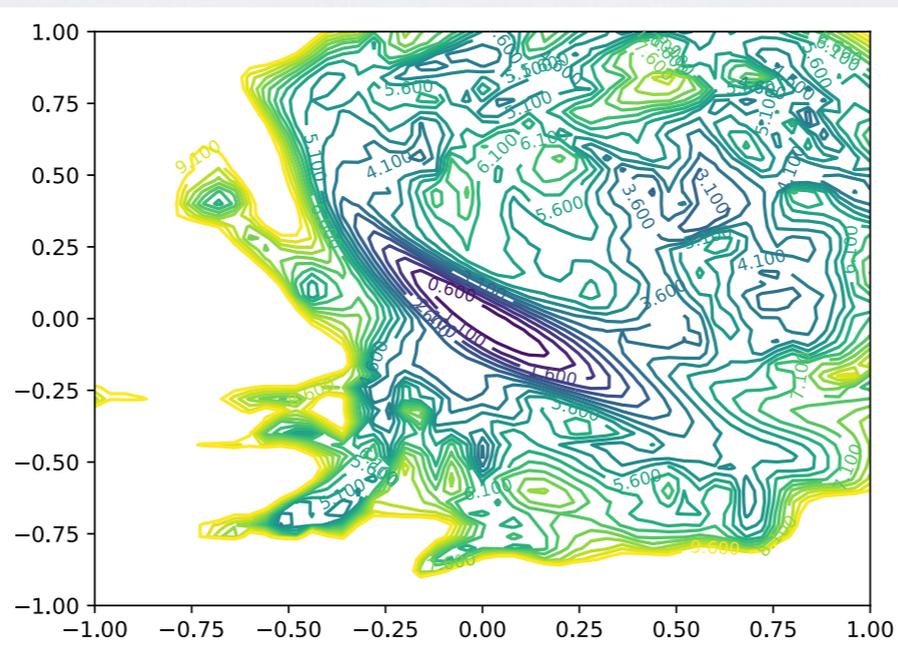
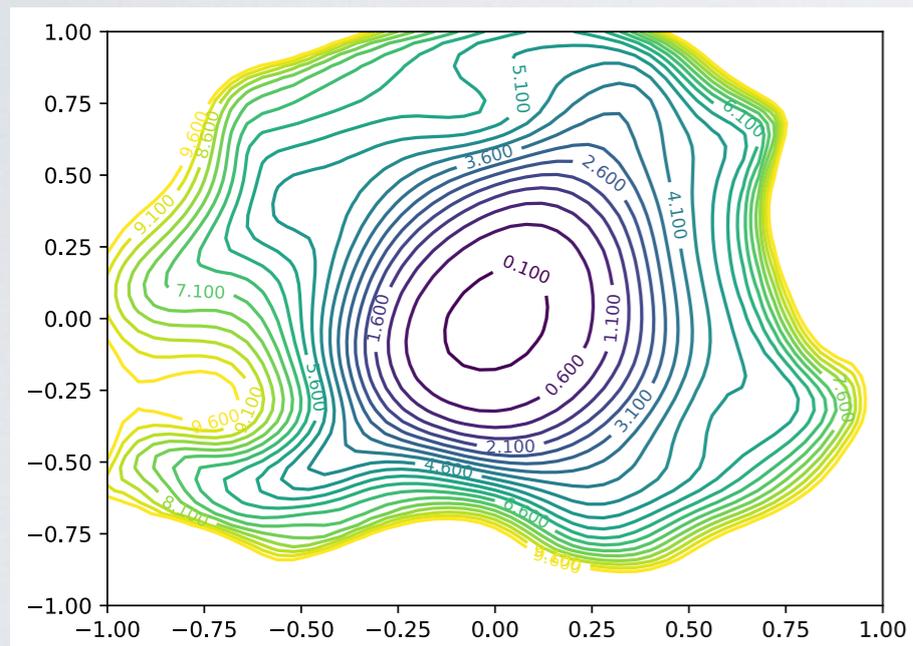
Chaos

CHAOTIC TRANSITIONS

VGG-20

VGG-56

VGG-110



Convexity

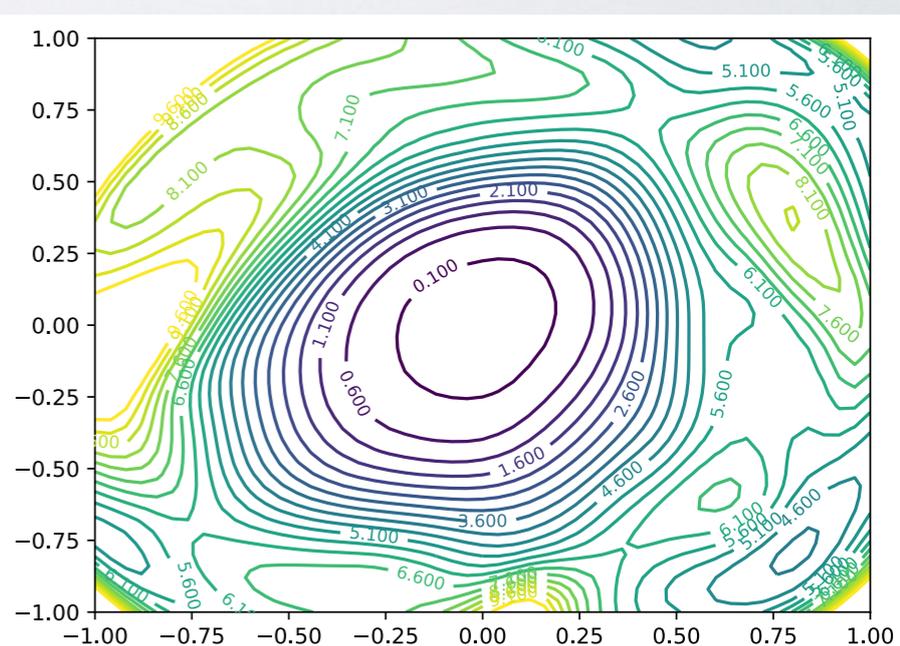
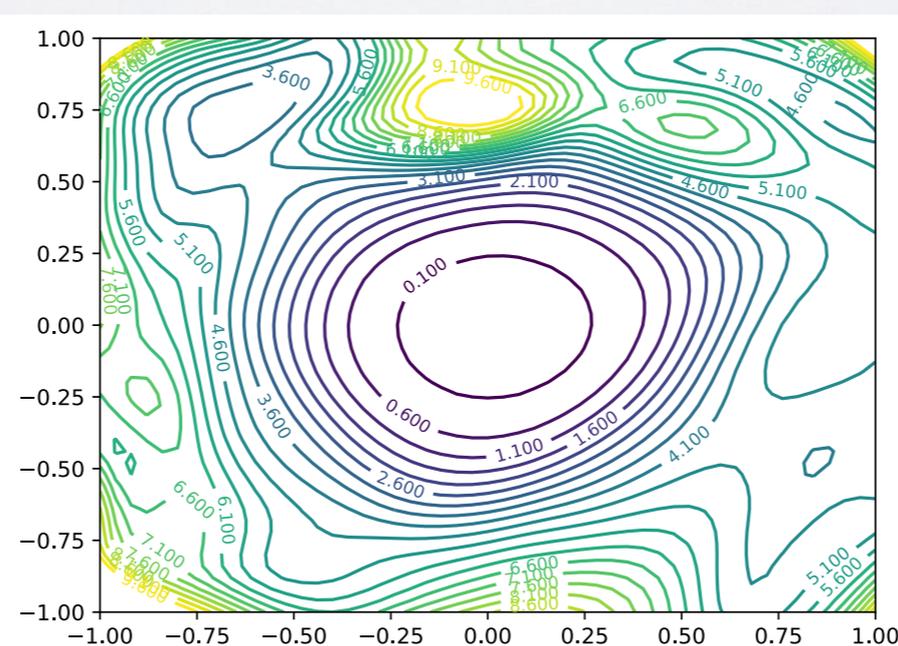
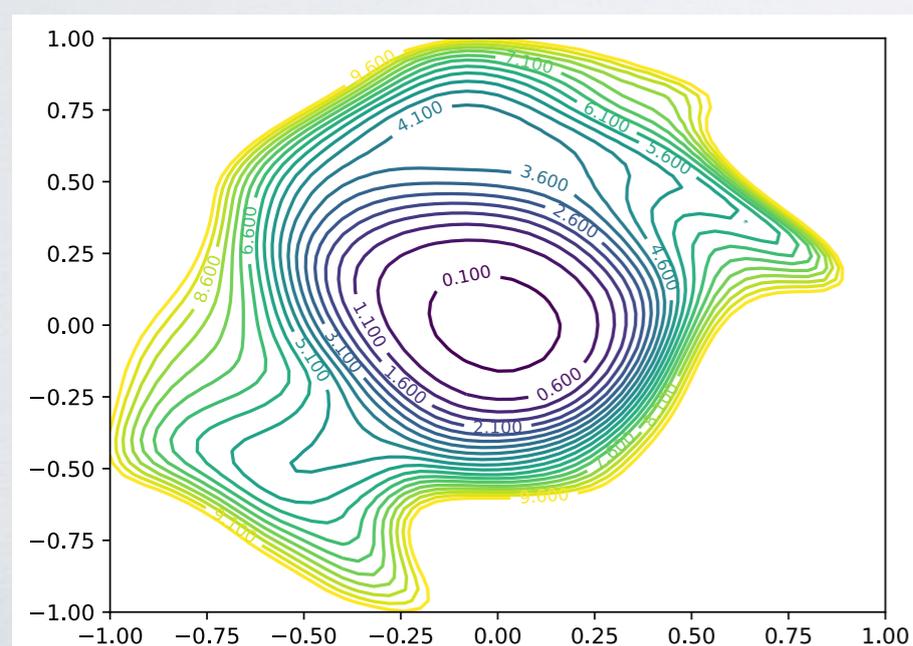
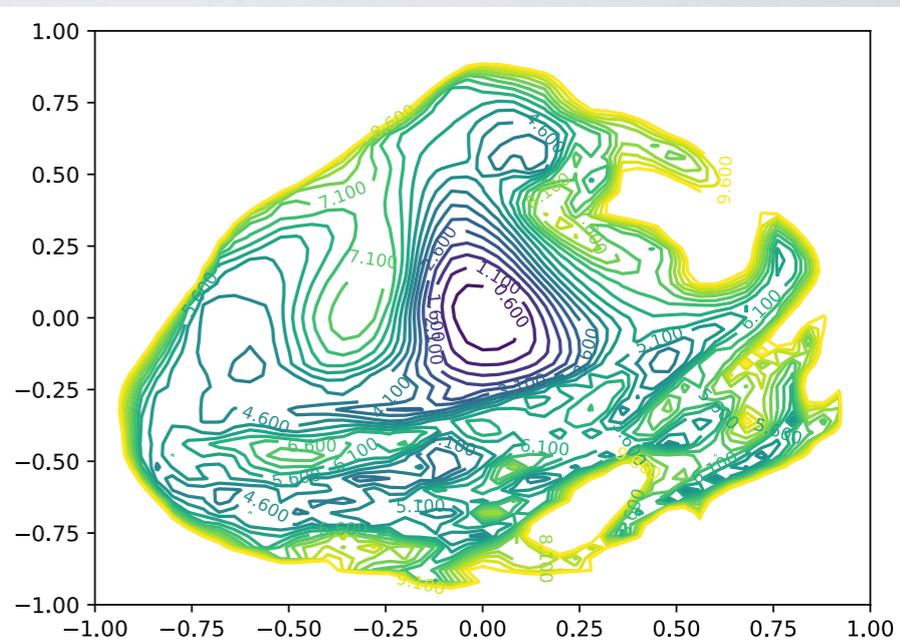
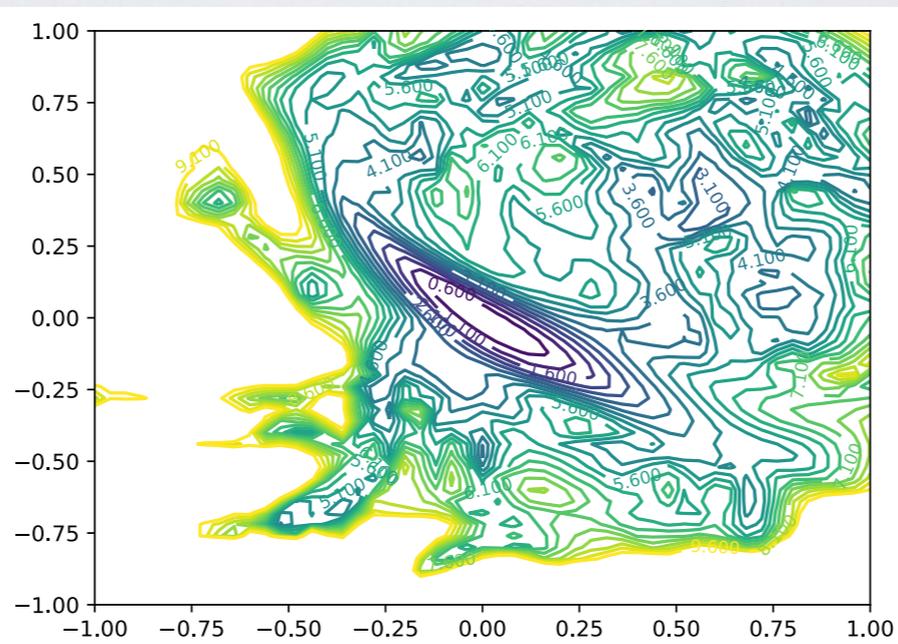
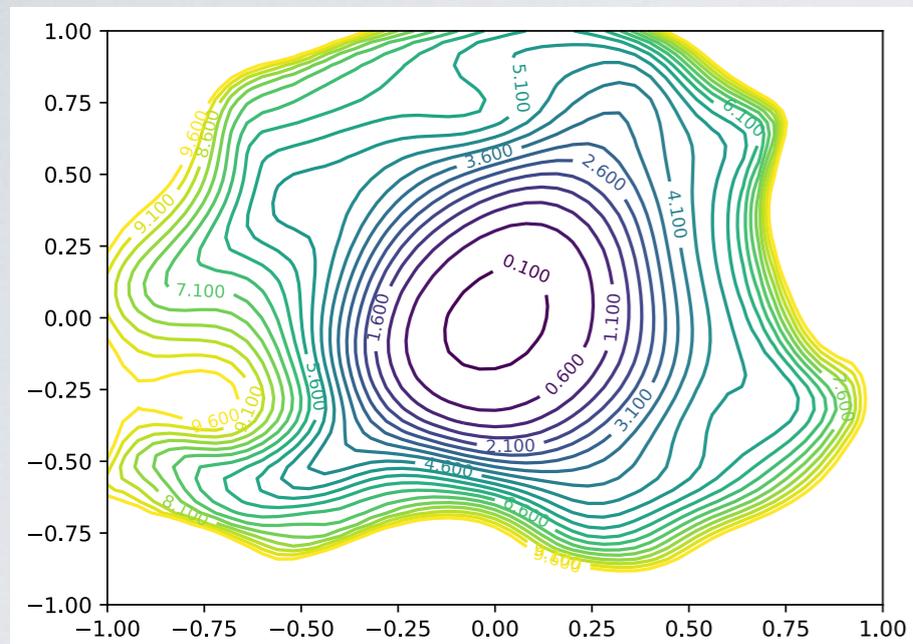
Chaos

CHAOTIC TRANSITIONS

VGG-20

VGG-56

VGG-110

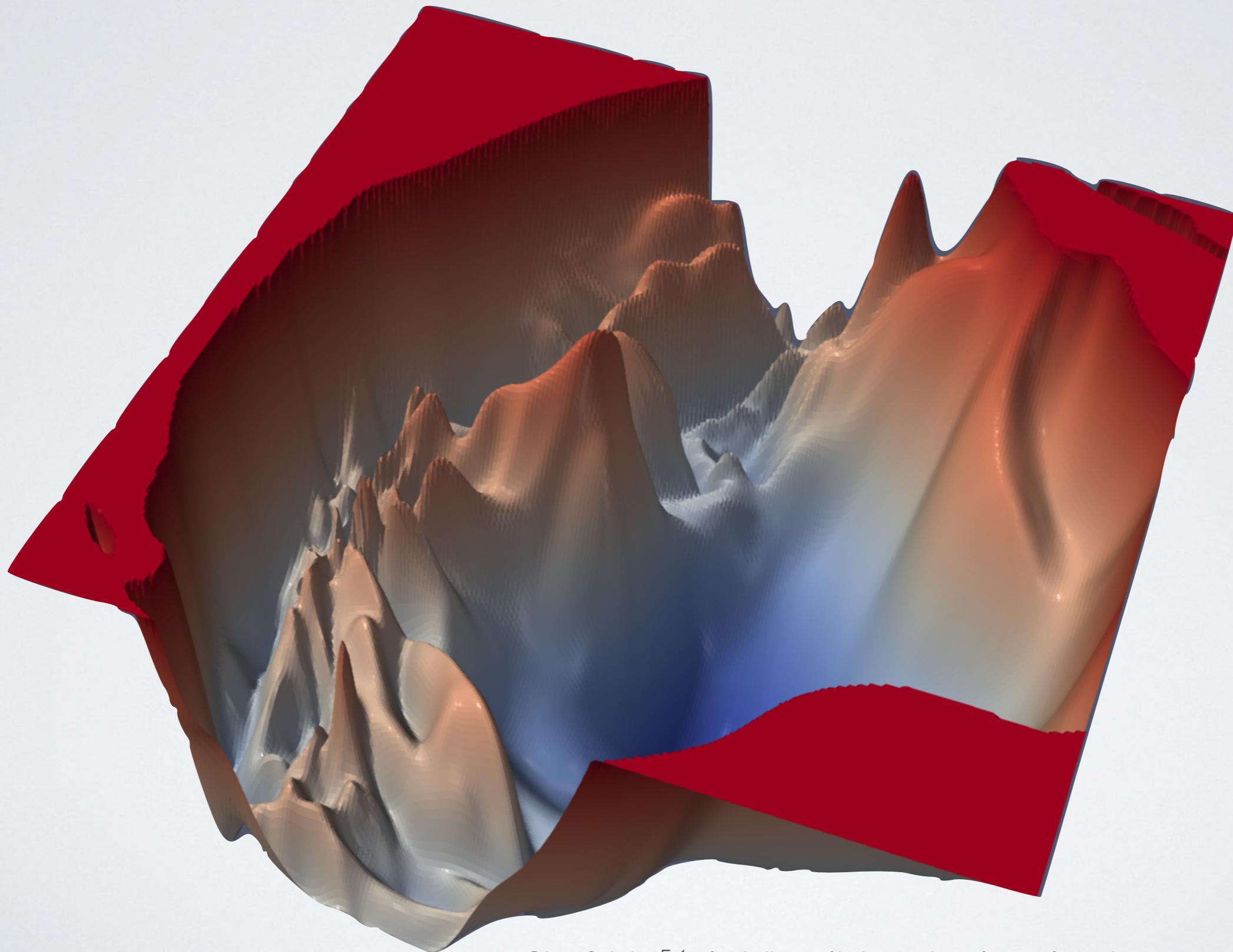


ResNet-20

ResNet-56

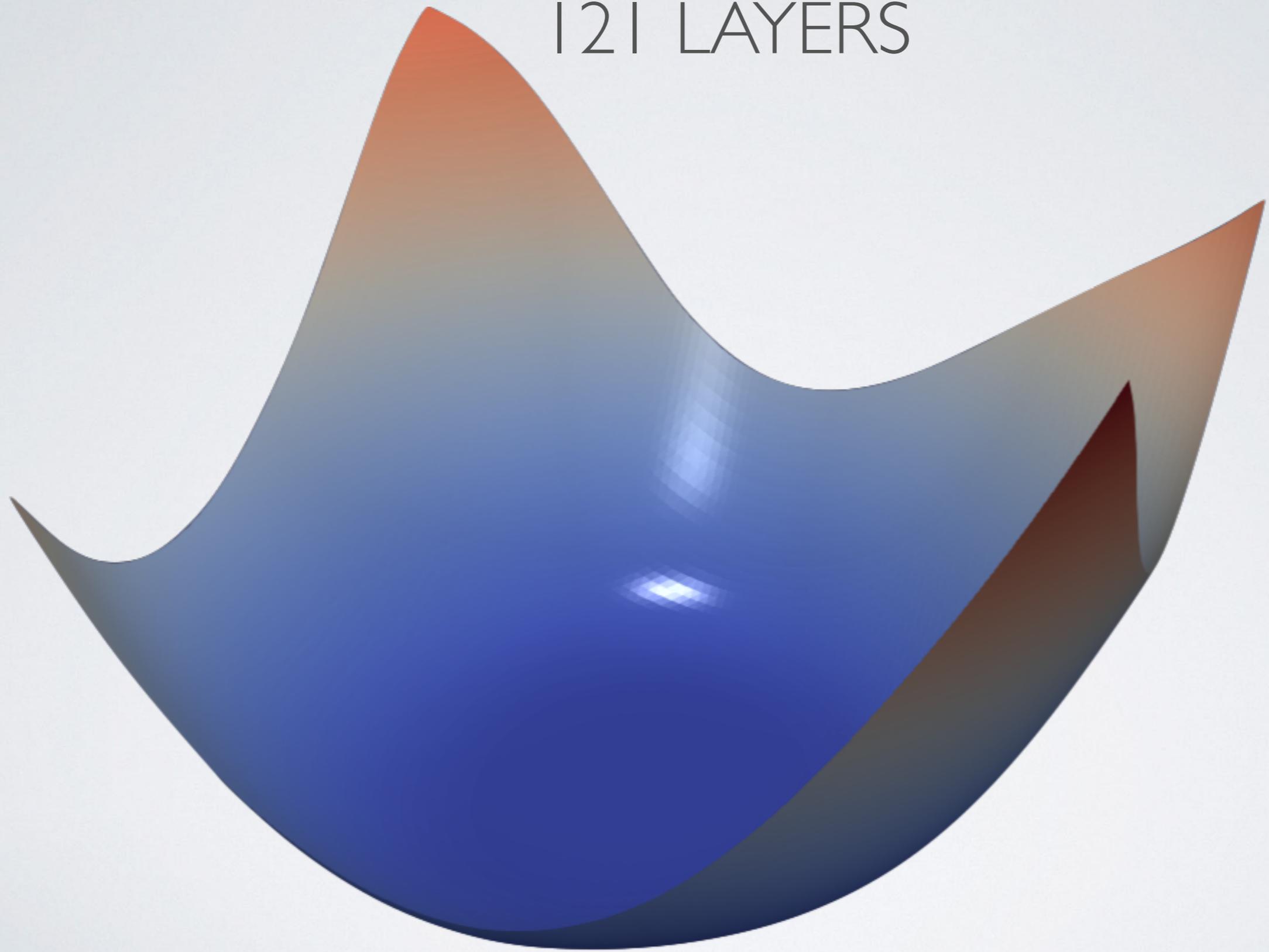
ResNet-110

VGG-110



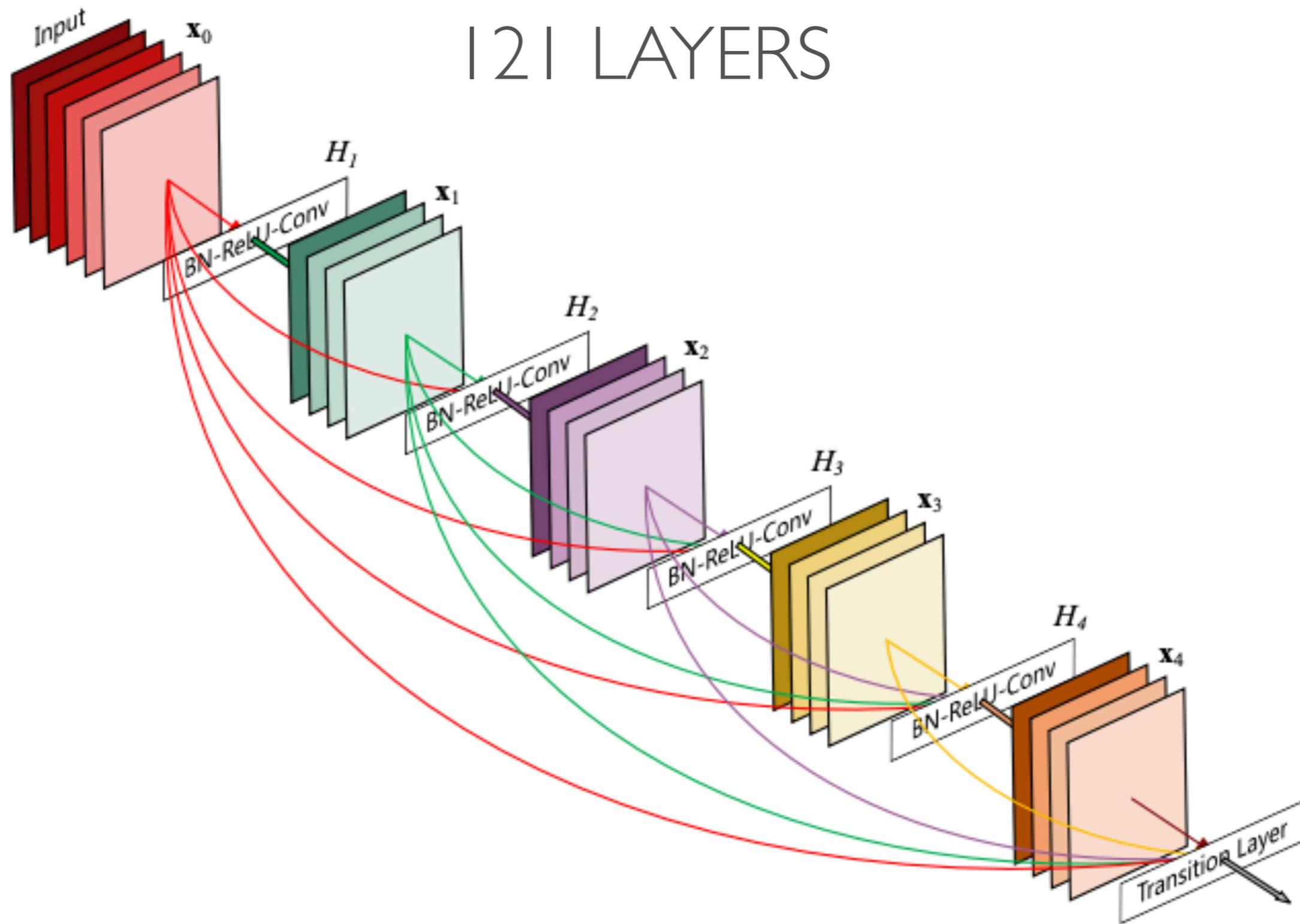
DENSENET

121 LAYERS



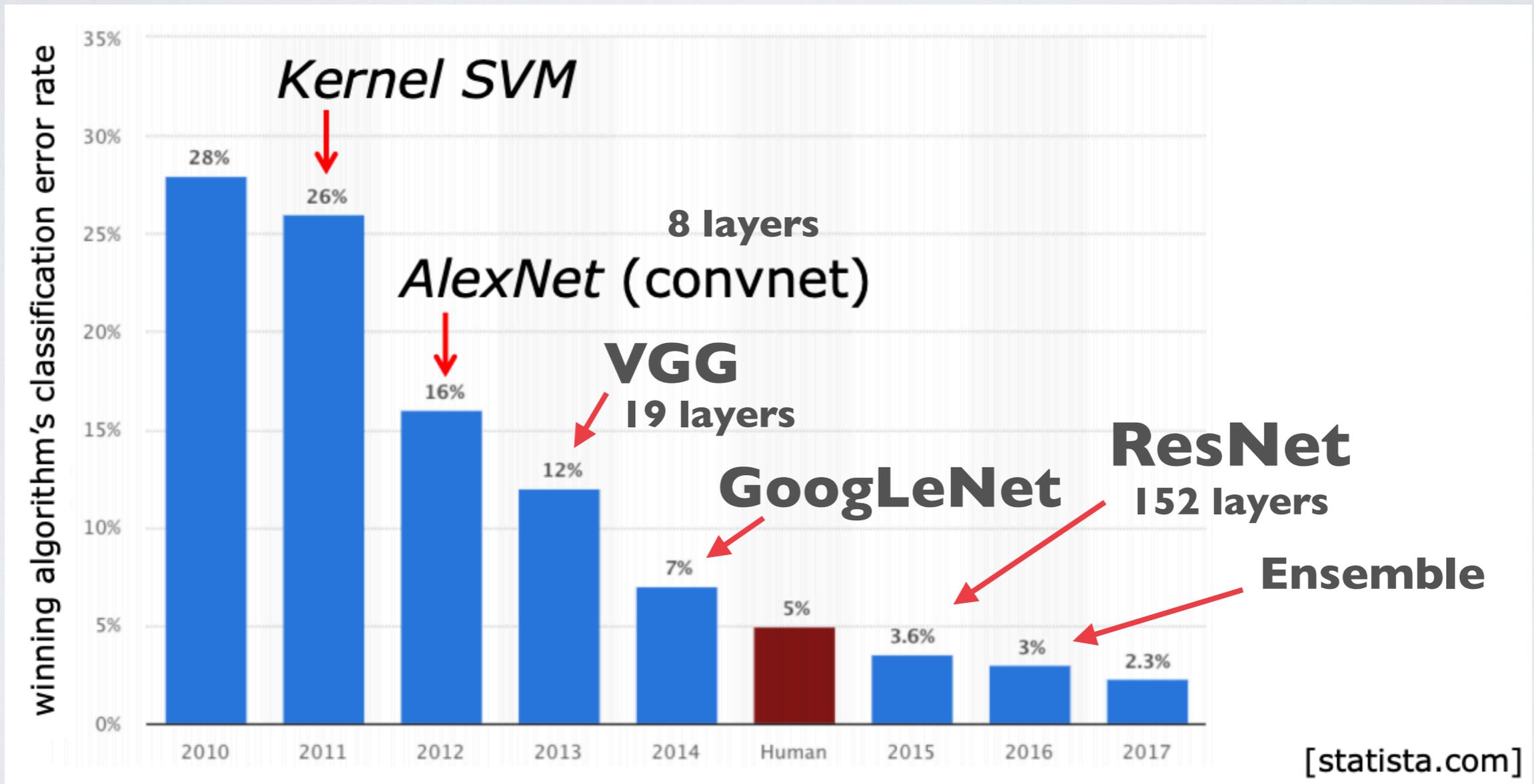
DENSENET

121 LAYERS

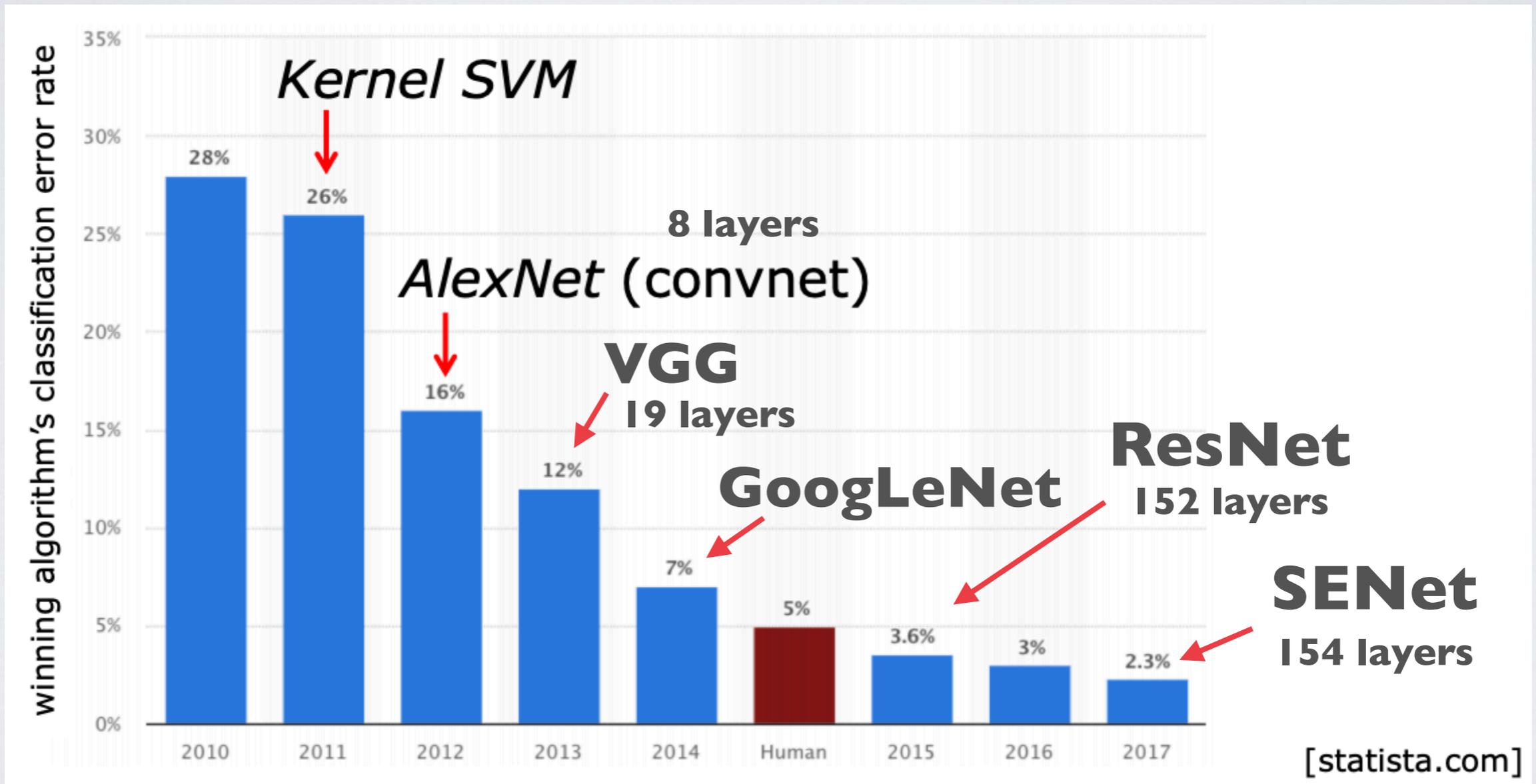


Huang et al, "Densely Connected," 2016

IMAGENET HISTORY



IMAGENET HISTORY



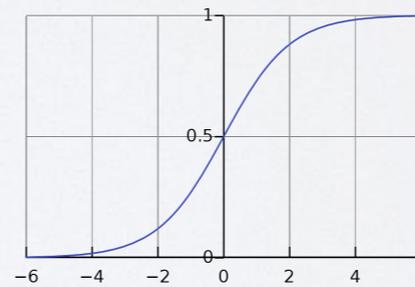
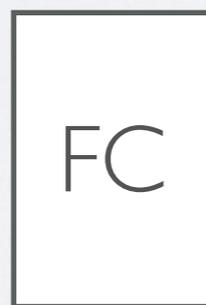
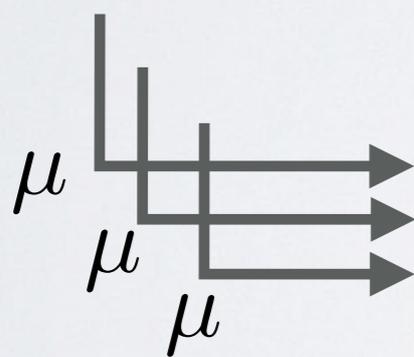
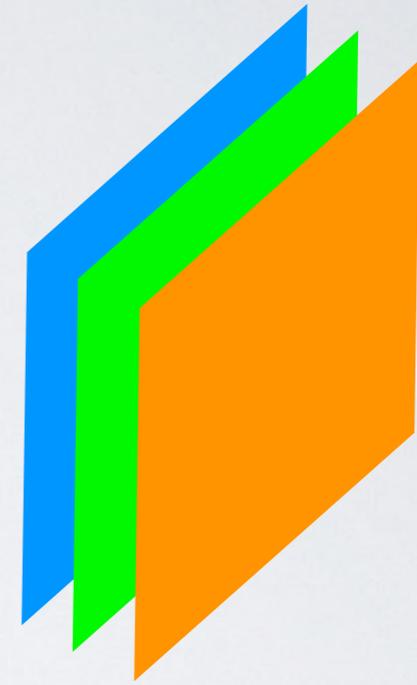
SENET

squeeze and excite network

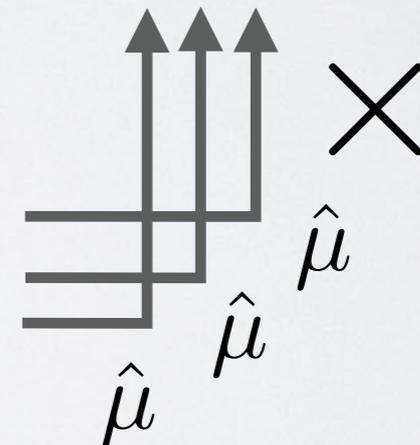
Hu et al 2017



SENet block



sigmoid



"squeeze"

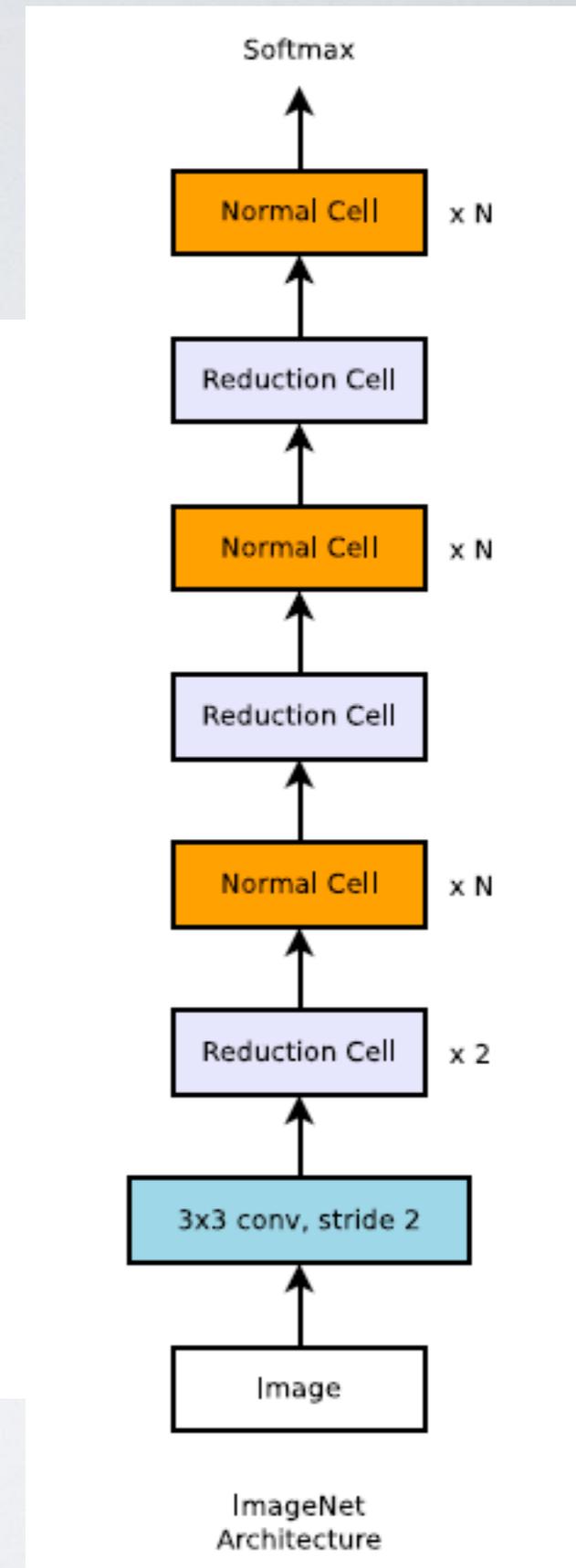
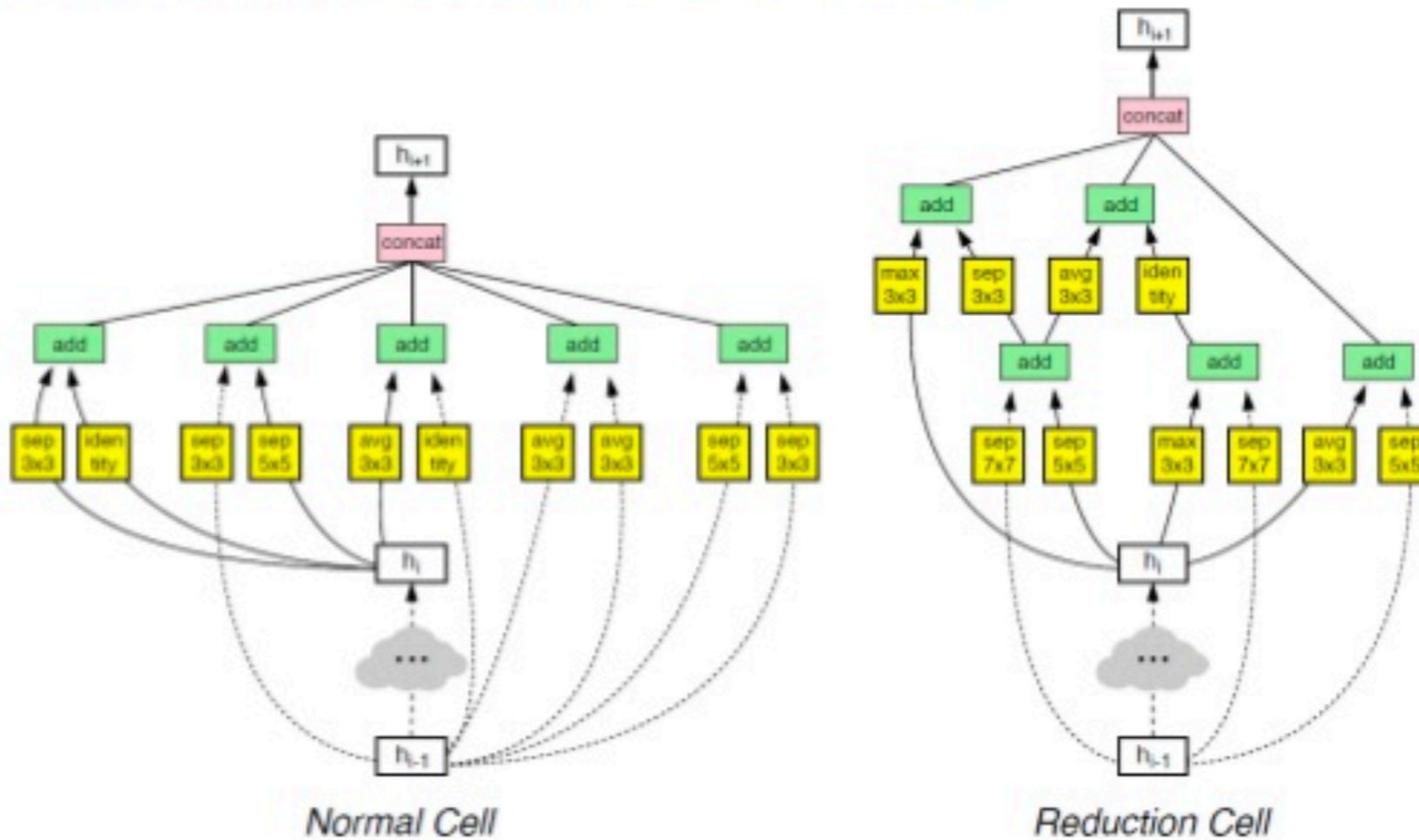
global average pooling

"excite"

multiply

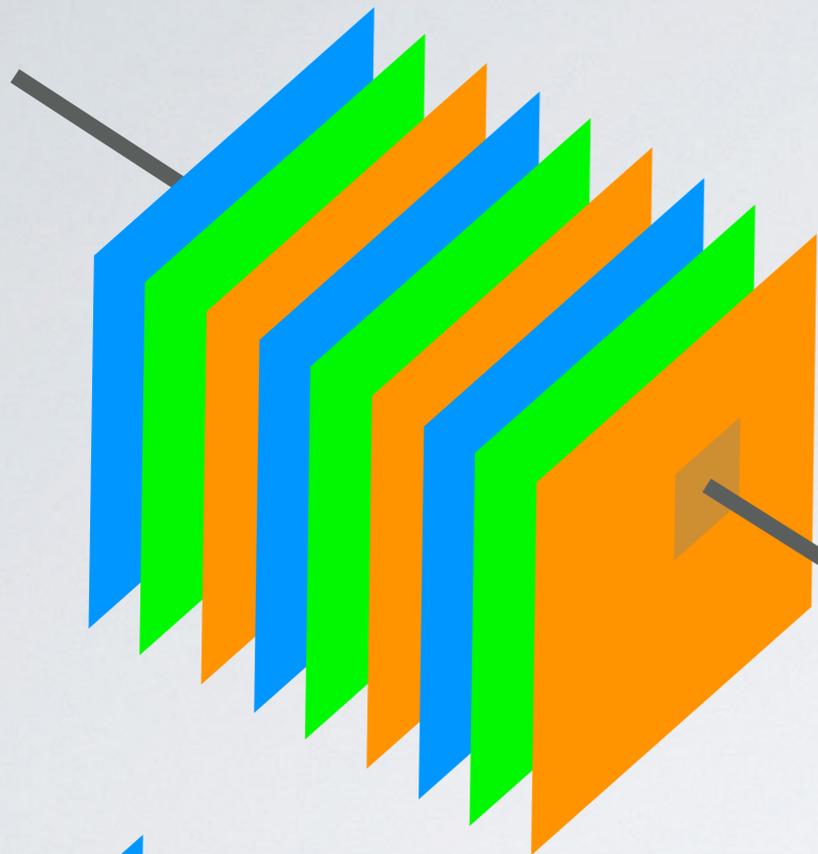
AUTO ML

Best Architecture (NASNet-A)

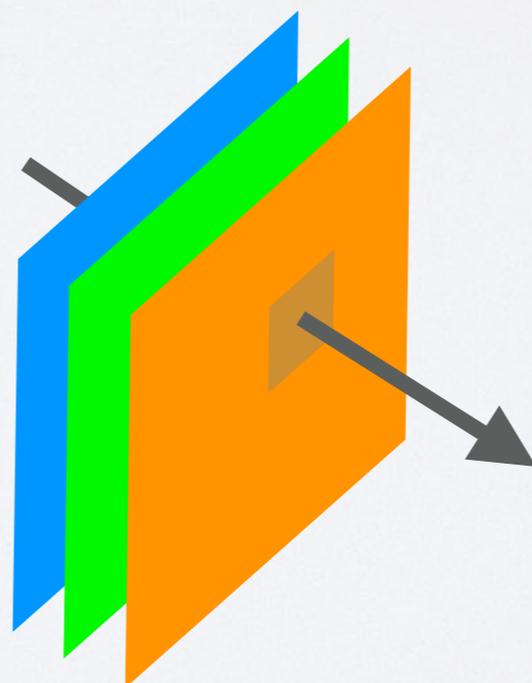
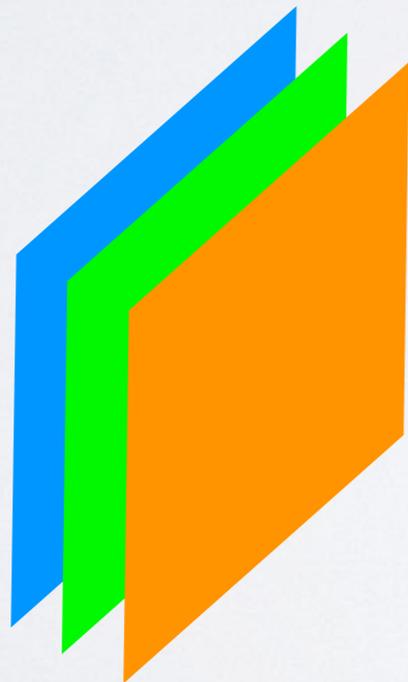
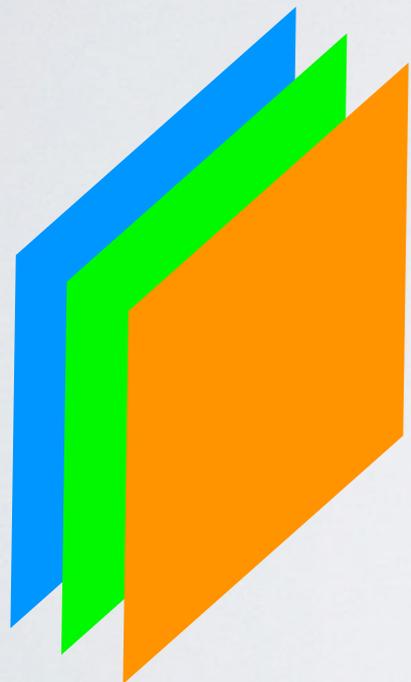


CONTROLLING COMPUTATIONAL COST

GROUPED CONVOLUTION



Feature cost
 $k^2 \times C$



Feature cost
 $k^2 \times (C/g)$

SEPARABLE CONV

MobileNet
2017

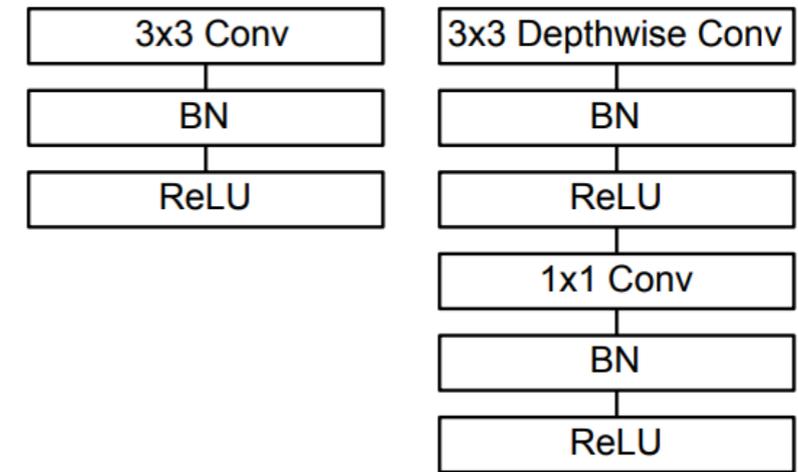
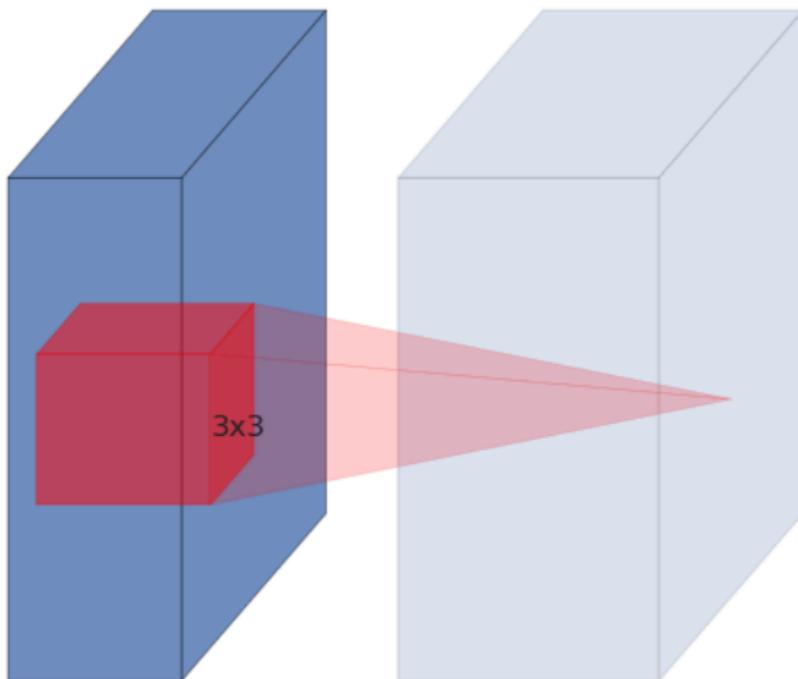
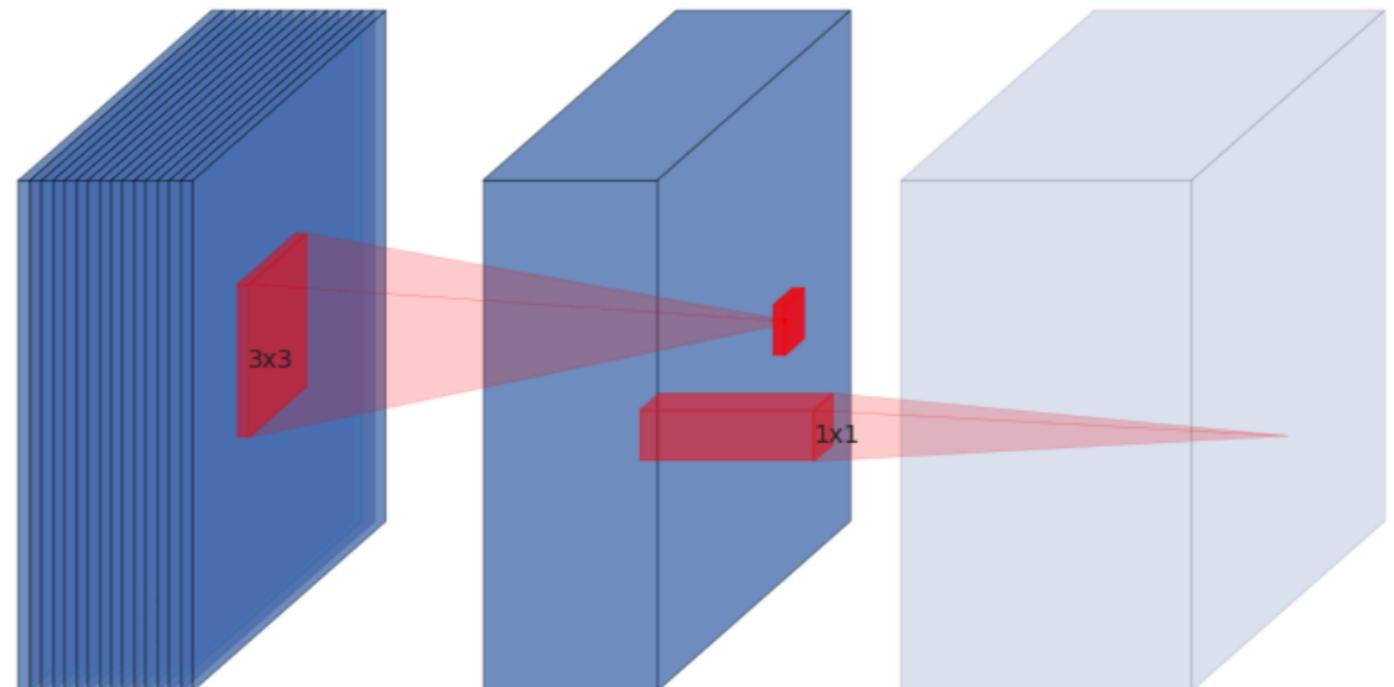


Figure 3. Left: Standard convolutional layer with batchnorm and ReLU. Right: Depthwise Separable convolutions with Depthwise and Pointwise layers followed by batchnorm and ReLU.

Regular Convolution



Separable Convolution Block



BOTTLENECKS

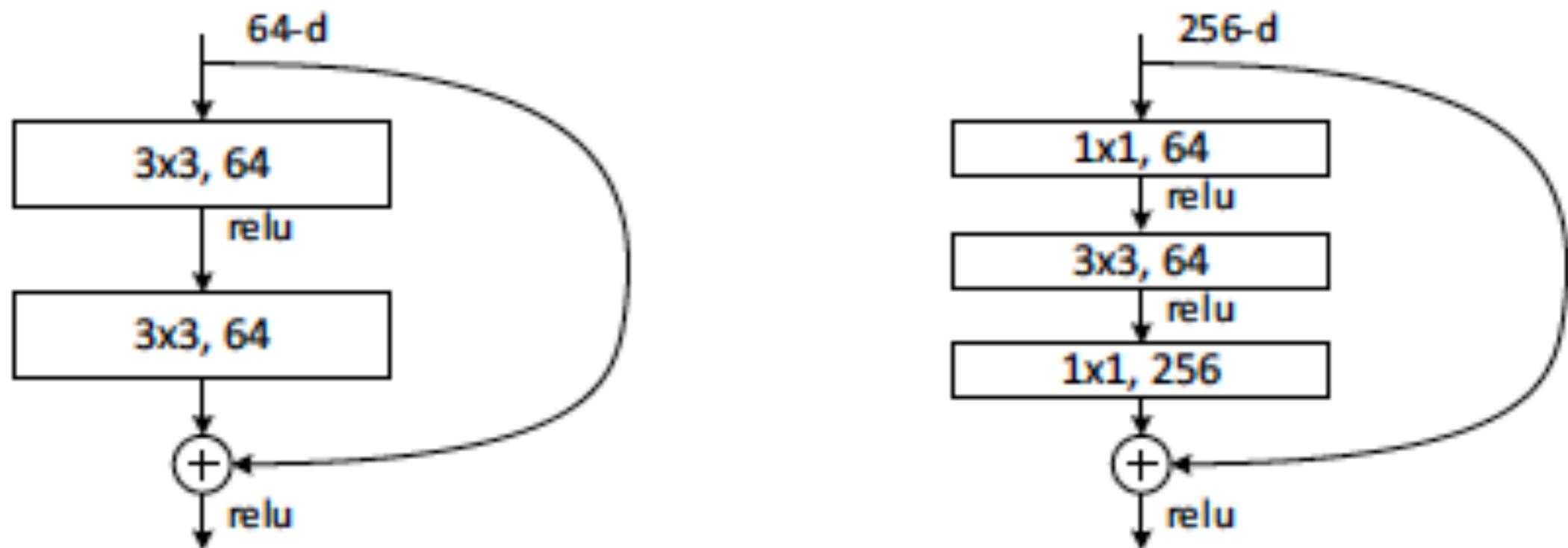


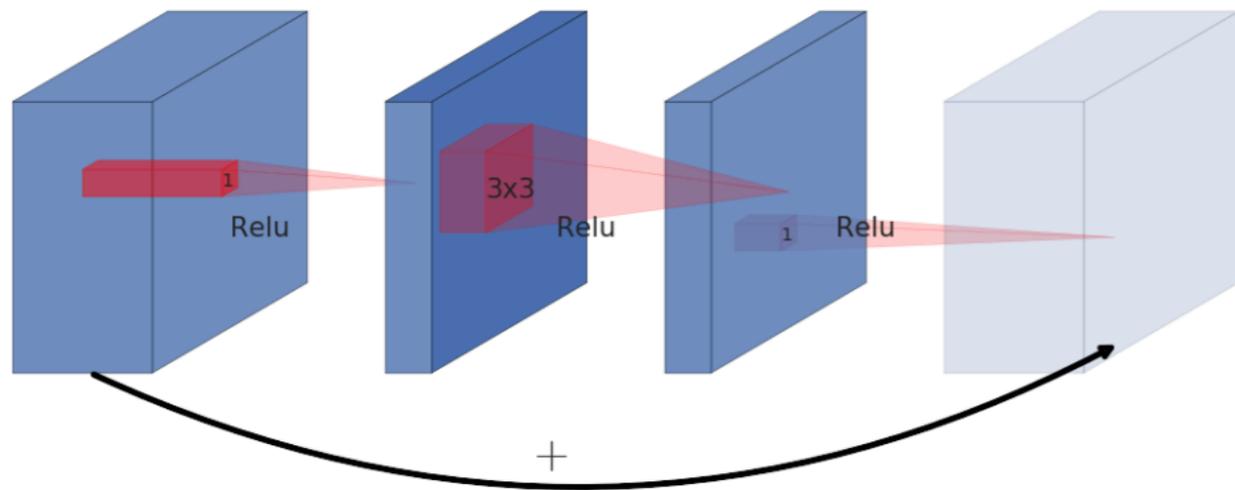
Figure 5. A deeper residual function \mathcal{F} for ImageNet. Left: a building block (on 56×56 feature maps) as in Fig. 3 for ResNet-34. Right: a “bottleneck” building block for ResNet-50/101/152.

INVERTED BOTTLENECK

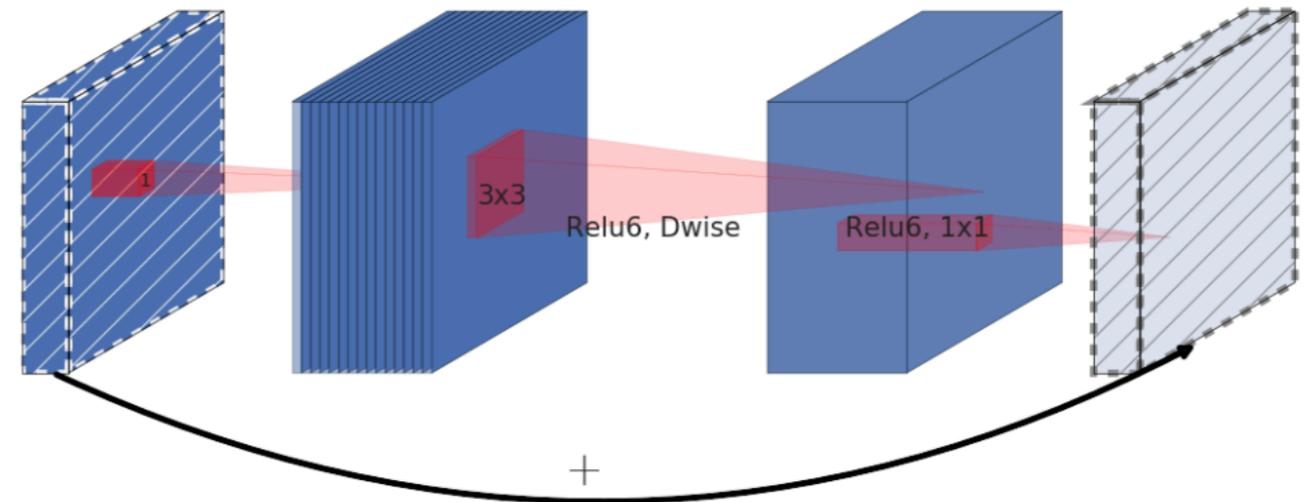
MobileNetv2

2019

(a) Residual block



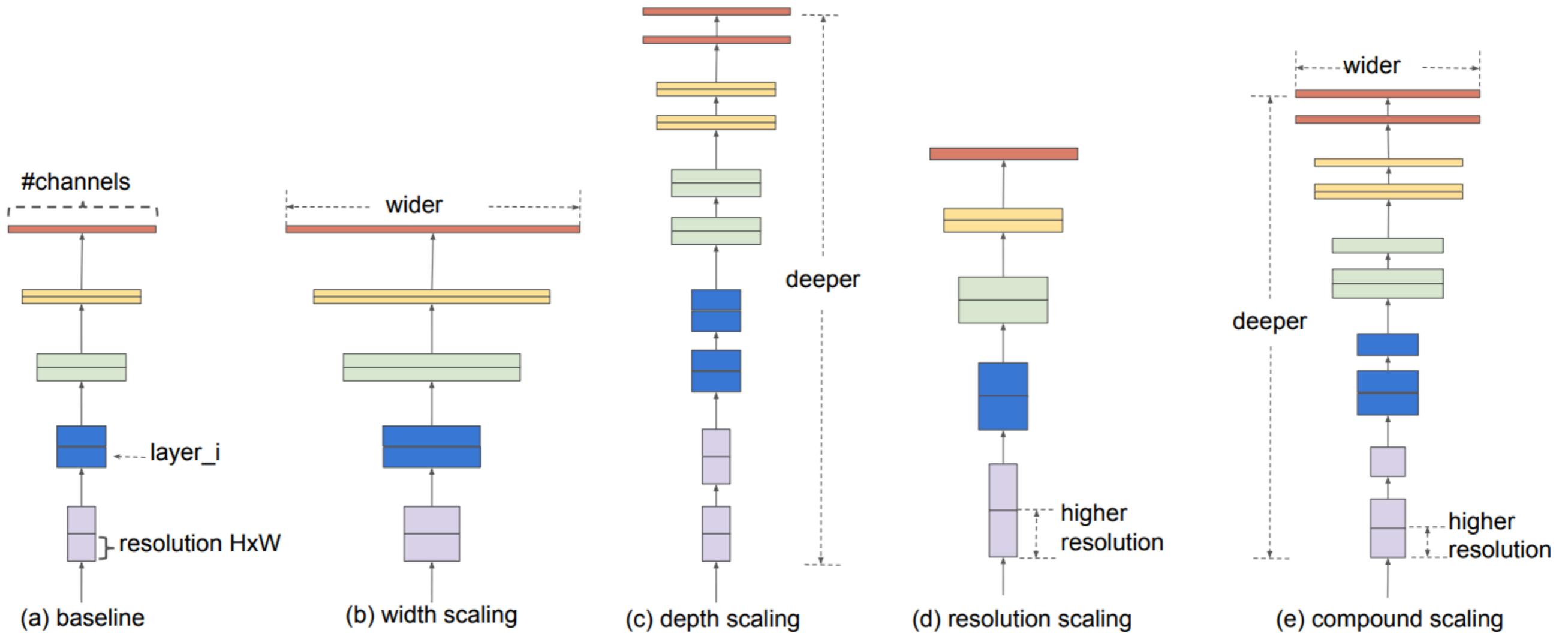
(b) Inverted residual block



Hashed boxed have no non-linearity

EfficientNet

2020



EfficientNet

2020

Optimize scaling parameters to get best performance for fixed cost

Do architecture search once on small models, then scale it up

$$\text{depth: } d = \alpha^\phi$$

$$\text{width: } w = \beta^\phi$$

$$\text{resolution: } r = \gamma^\phi$$

$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

EfficientNet-B0

Based on MobileNet

Stage i	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224×224	32	1
2	MBConv1, k3x3	112×112	16	1
3	MBConv6, k3x3	112×112	24	2
4	MBConv6, k5x5	56×56	40	2
5	MBConv6, k3x3	28×28	80	3
6	MBConv6, k5x5	14×14	112	3
7	MBConv6, k5x5	14×14	192	4
8	MBConv6, k3x3	7×7	320	1
9	Conv1x1 & Pooling & FC	7×7	1280	1

EfficientNet

2020

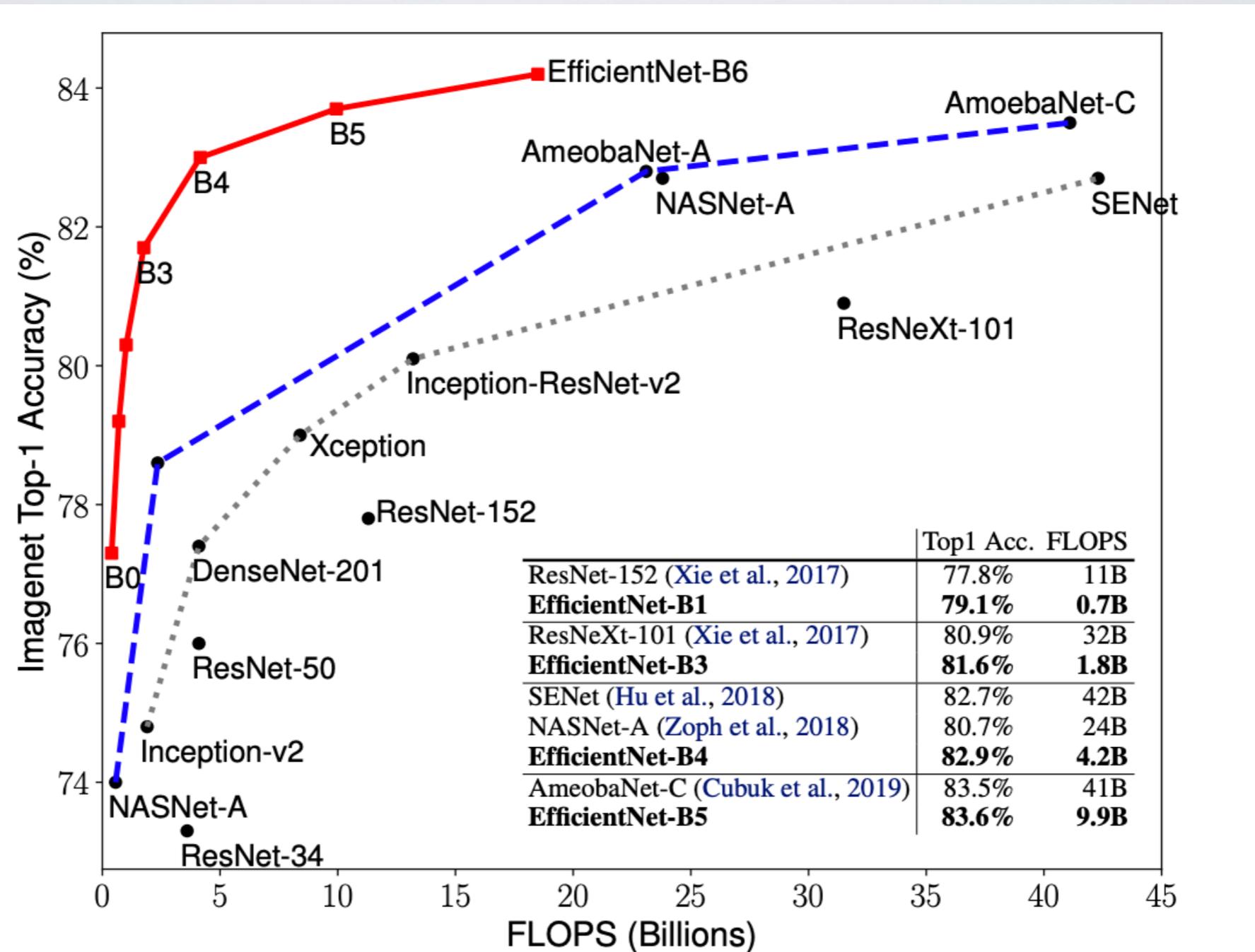


Figure 5. FLOPS vs. ImageNet Accuracy – Similar to Figure 1 except it compares FLOPS rather than model size.

Revisiting ResNets

2021

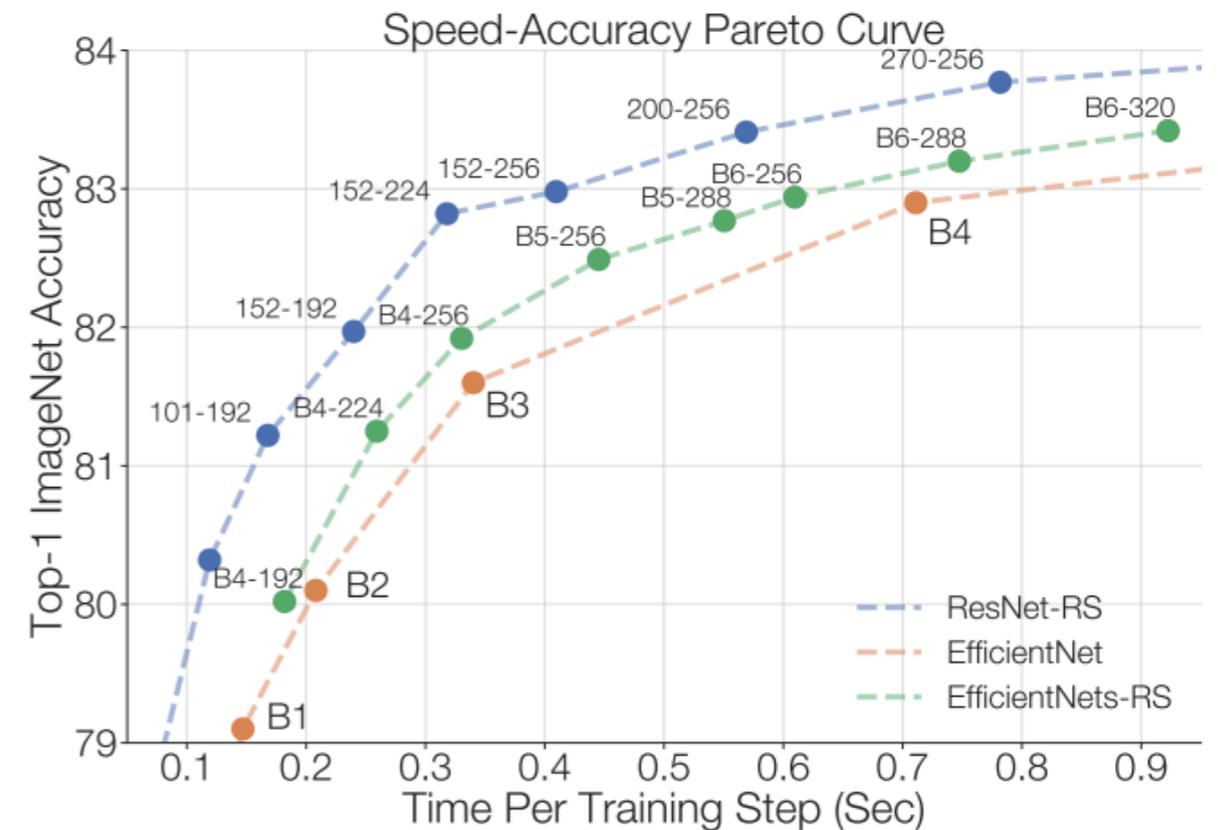
Train (modified) ResNets with modern tricks

The resolution scaling of EfficientNet is inefficient

ResNets are faster than EffNets despite having more FLOPs

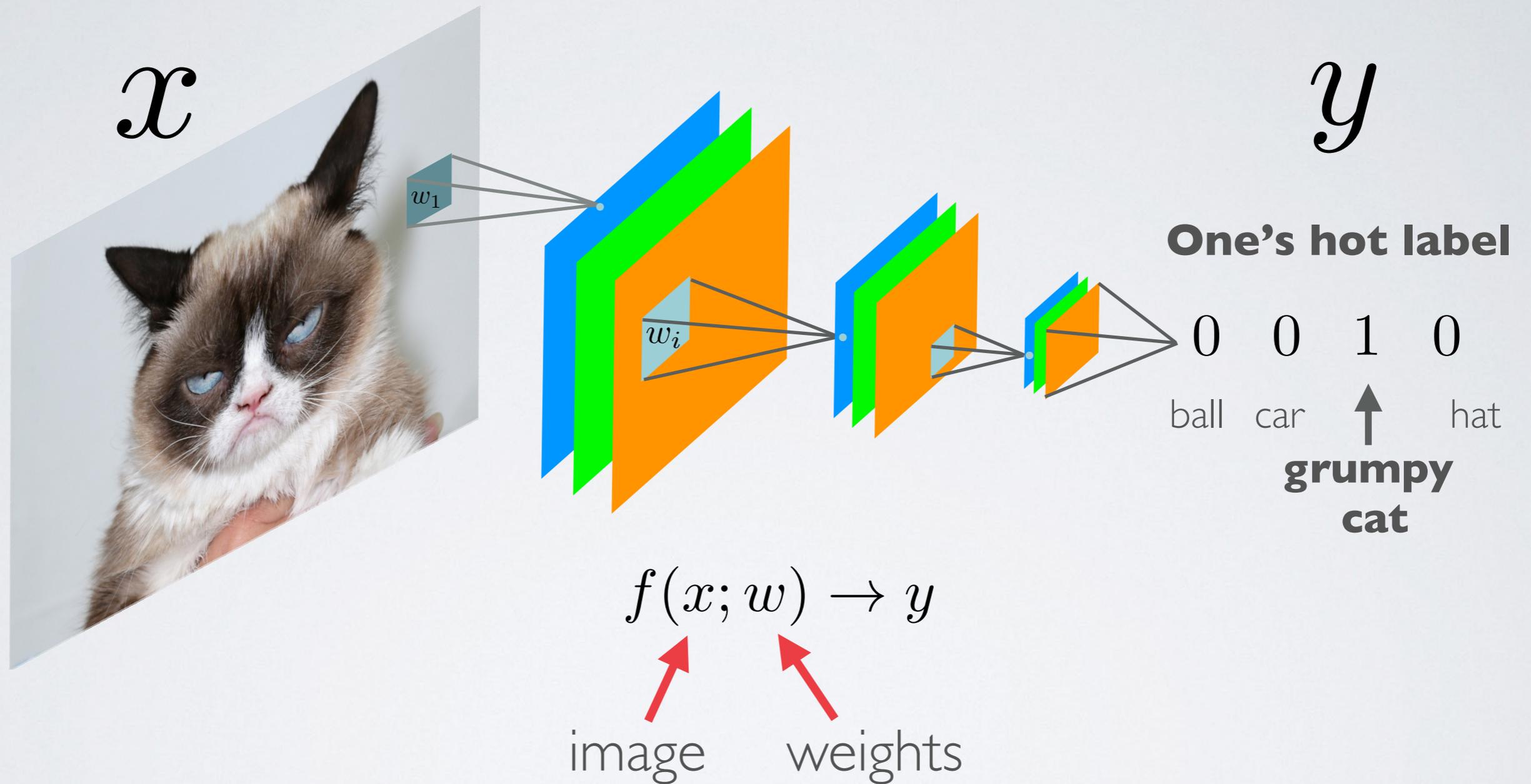
ResNet-RS trains faster than EffNet

Model	RS-350	ENet-B6	RS-420	ENet-B7
Resolution	256	528	320	600
Top-1 Acc.	84.0	84.0	84.4	84.7
Params (M)	164	43 (3.8x)	192	66 (2.9x)
FLOPs (B)	69	38 (1.8x)	128	74 (1.7x)
TPU-v3				
Latency (s)	1.1	3.0 (2.7x)	2.1	6.0 (2.9x)
Memory (GB)	7.3	16.6 (2.3x)	15.5	28.3 (1.8x)
V100				
Latency (s)	4.7	15.7 (3.3x)	10.2	29.9 (2.8x)

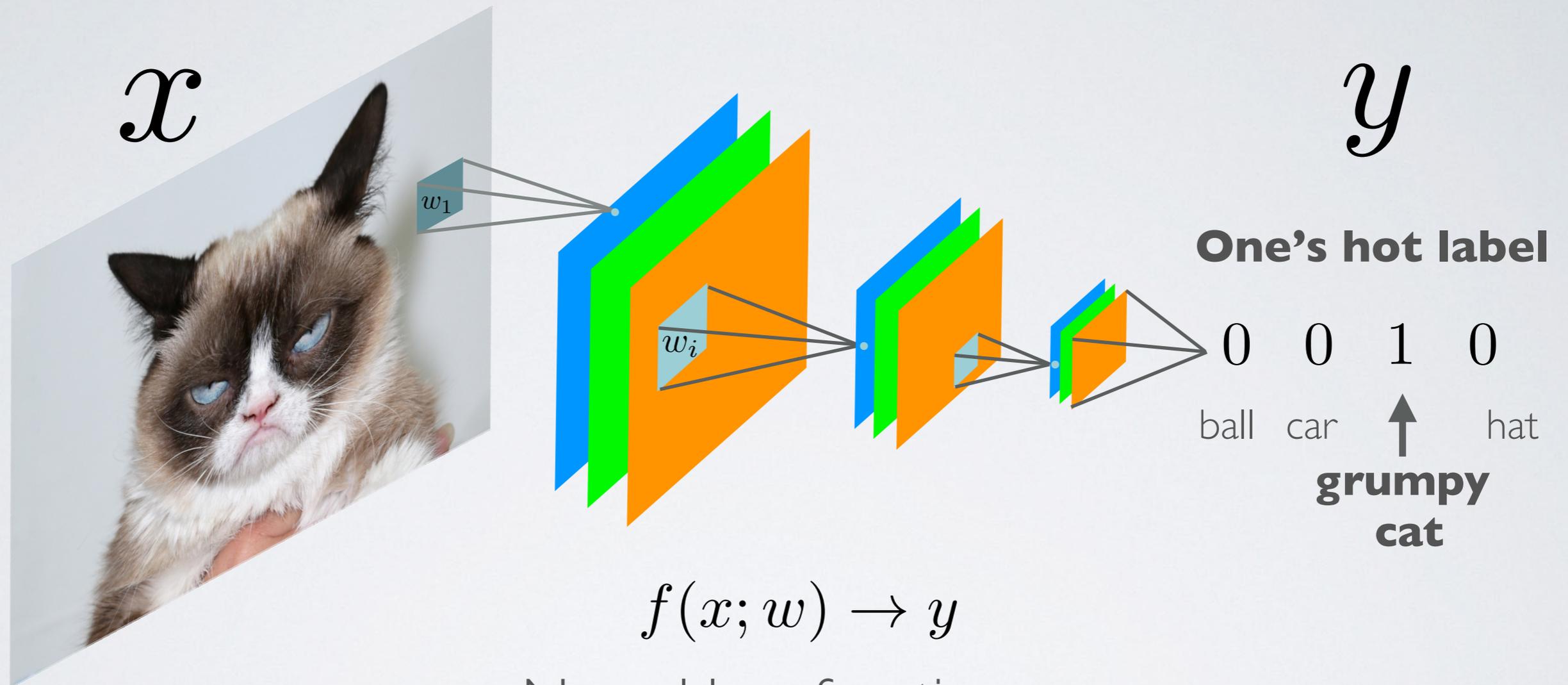


CONV NET SECRET SAUCE

FLEXIBILITY



FLEXIBILITY



$$f(x; w) \rightarrow y$$

Neural loss function

$$L(w) = \min_w \sum_i \|f(x_i; w) - y_i\|^2$$

IMPLICIT BIAS

Neural net loss landscapes have bad minima

But we don't find them

Neural nets “like” to learn natural image behavior

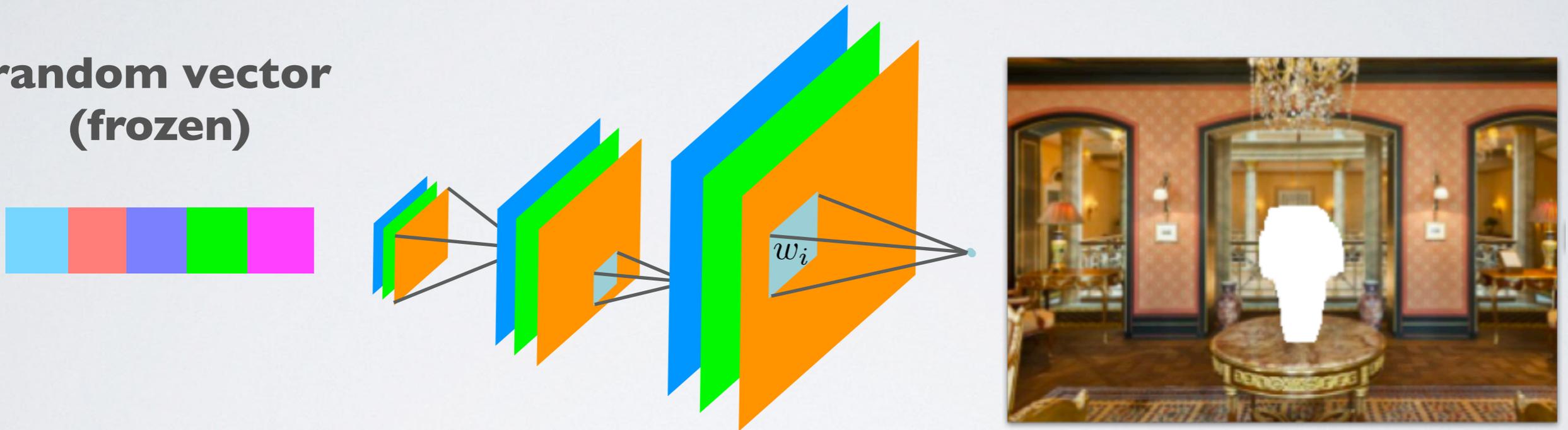
IMPLICIT BIAS



Corrupted

IMPLICIT BIAS

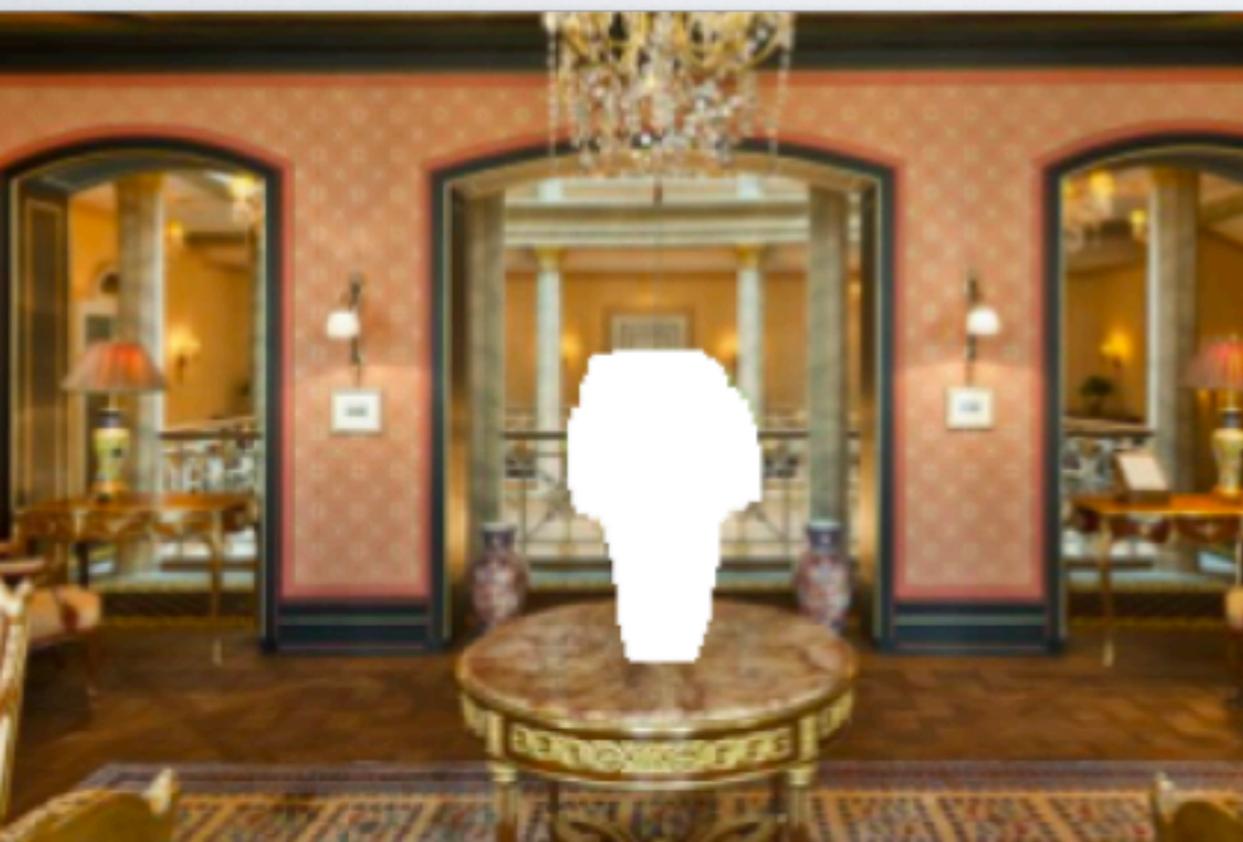
random vector
(frozen)



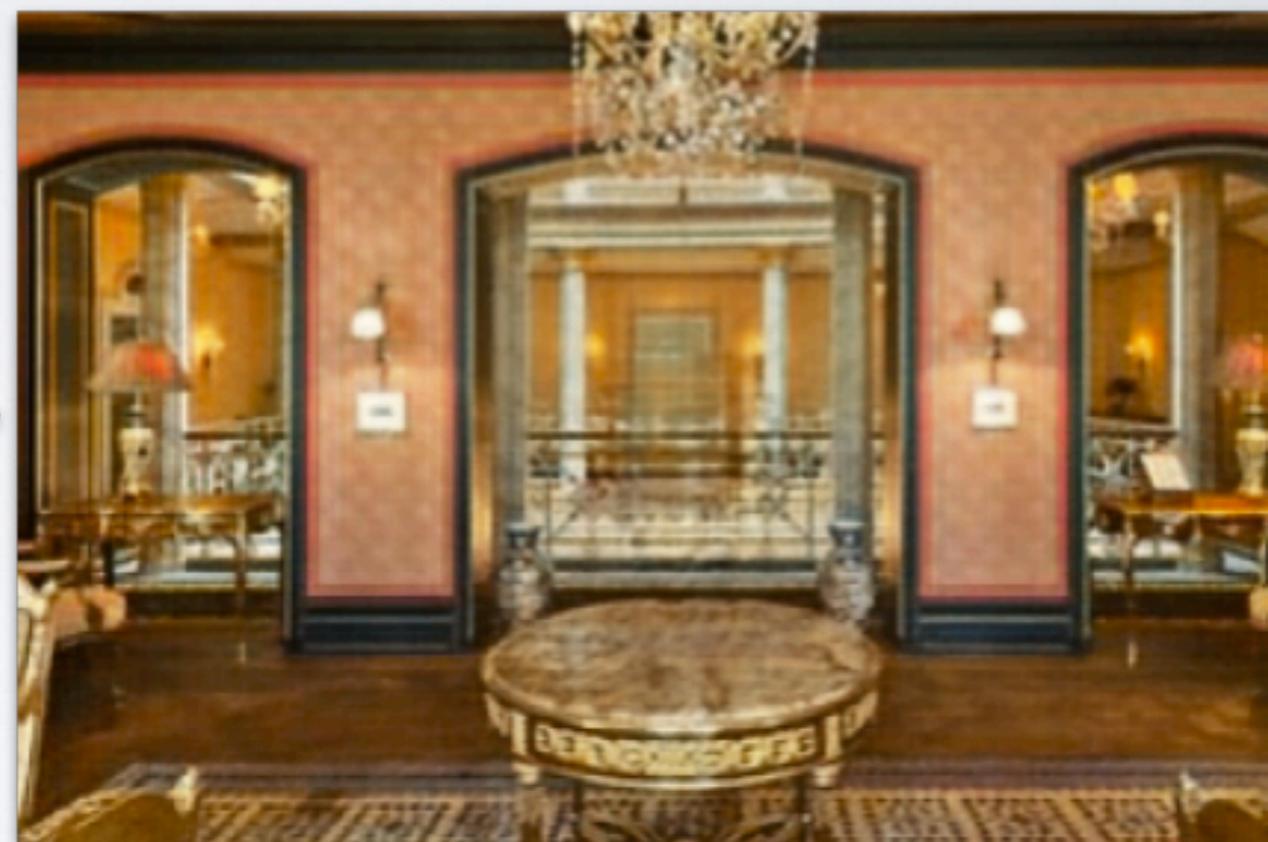
$$\min_{\text{network params}} \ell(\text{network output, known pixels})$$

IMPLICIT BIAS

Inpainting



Corrupted



Deep image prior

IMPLICIT BIAS

Inpainting



Corrupted



Deep image prior

TRANSFORMERS?

Not in this lecture