

Deep learning en pratique

Jean Savinien

Outline

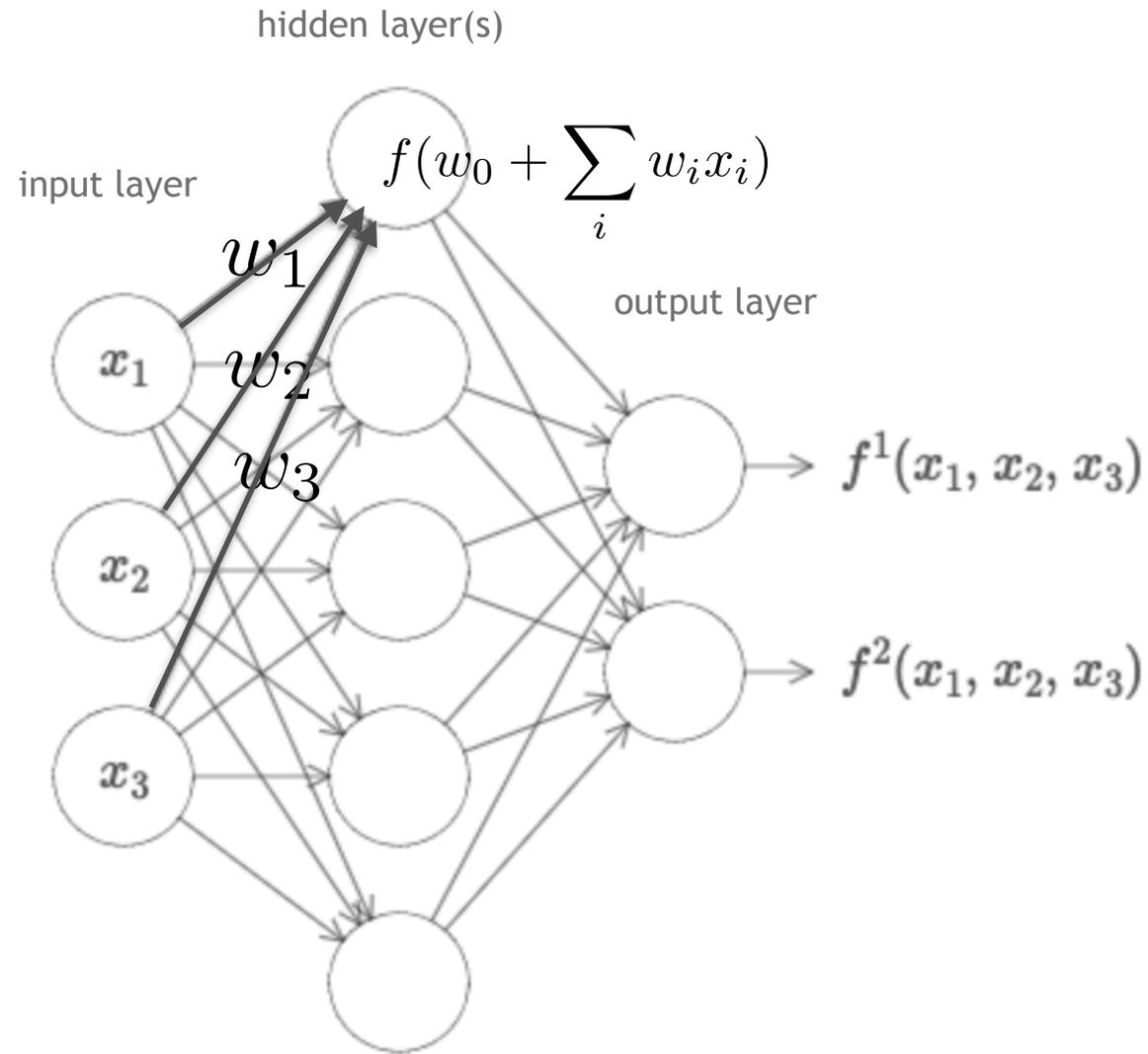
- **What are Neural Nets?**
- **Main Applications**
 - **Computer Vision**
 - **NLP**
 - **Tabular Data**
 - **Collaborative Filtering**
- **In Practice**
 - **Data Augmentation**
 - **Gradients Problems & Init**
 - **Learning Rates**
 - **Overfitting**
 - **GPUs & NN libraries**

“What is it ?”

“What is it for ?”

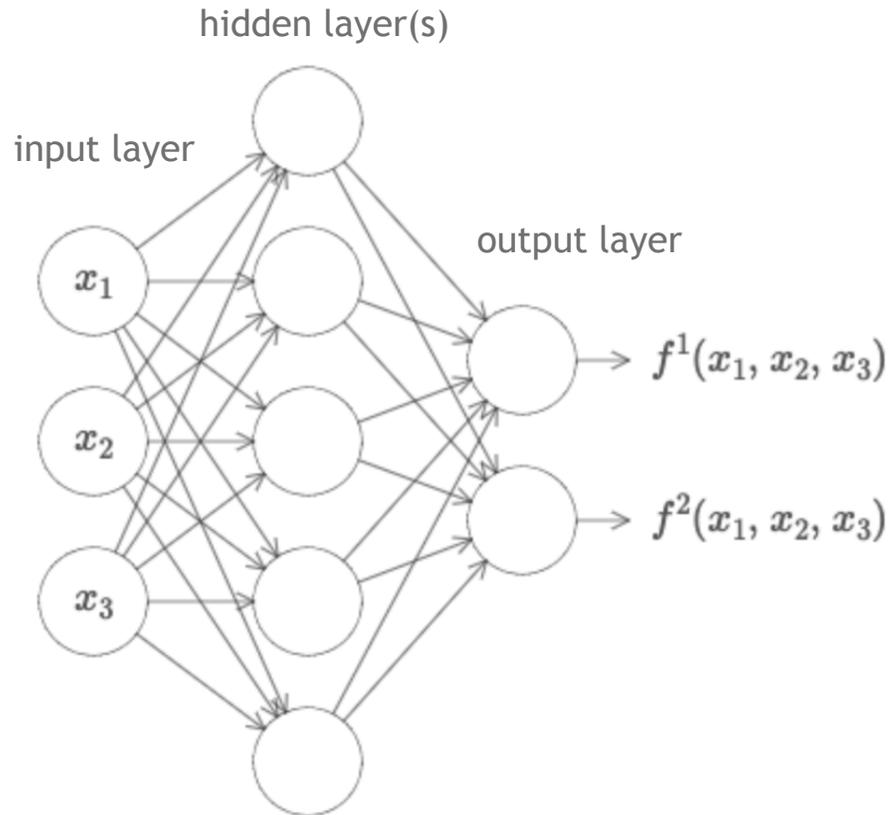
“How do you use them ?”

Neural Networks



NN is *deep* when
 $\#\{\text{hidden layers}\}$
is large

Universal Approximation Theorems



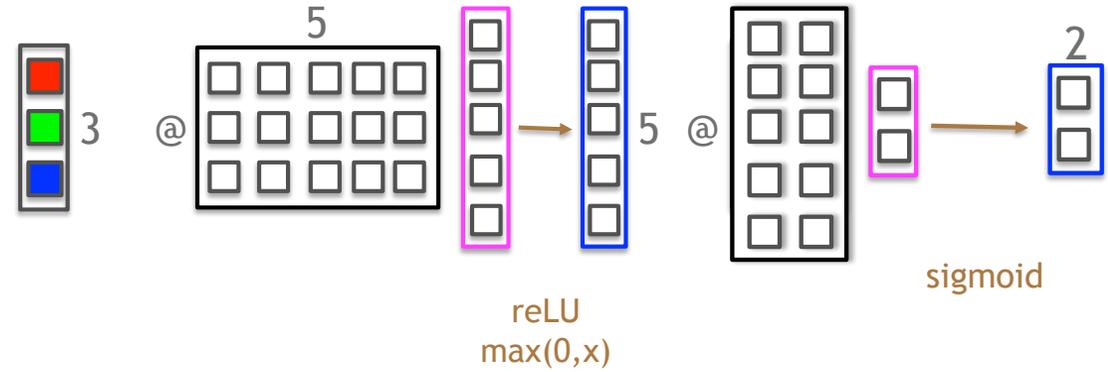
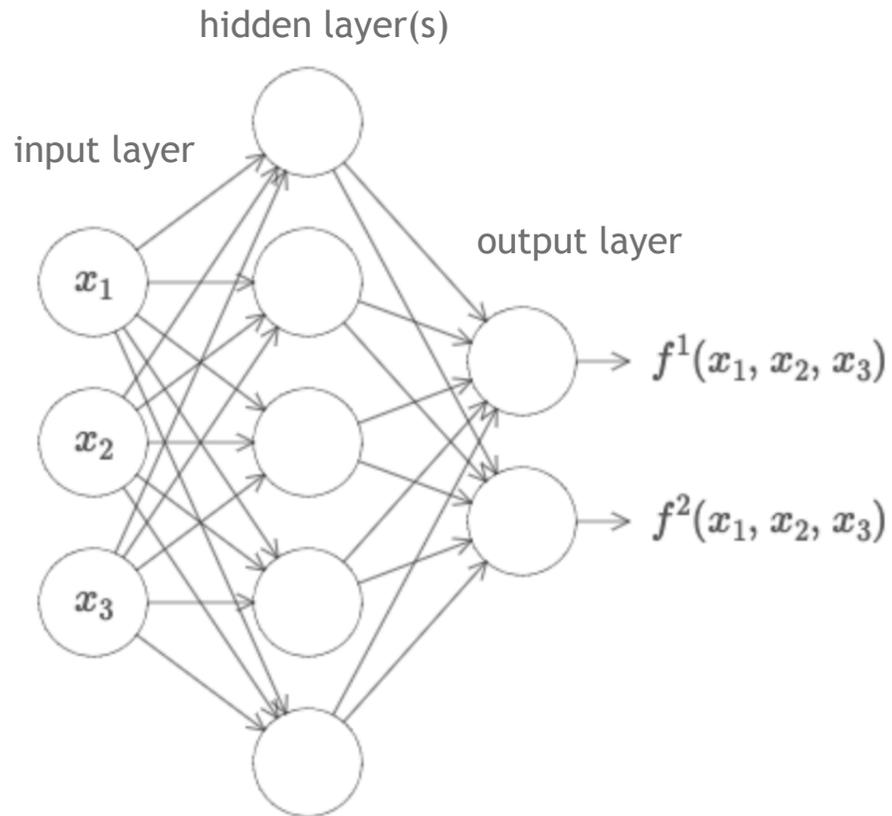
NN functions are dense in classes of functions of interest

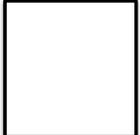
eg:

NN of depth 1 with continuous activations are dense in $C^0(\mathbb{R}^n, \mathbb{R}^m)$ (compact topology).

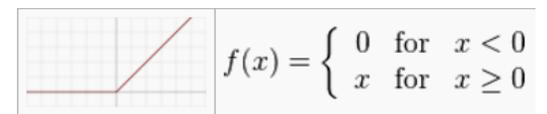
ref *G. Sibenko '89, K. Hornik '91, ... A. Kratsios 2020*

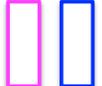
Neural Networks

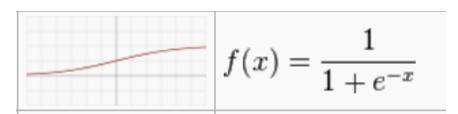


 parameters (weights, biases)

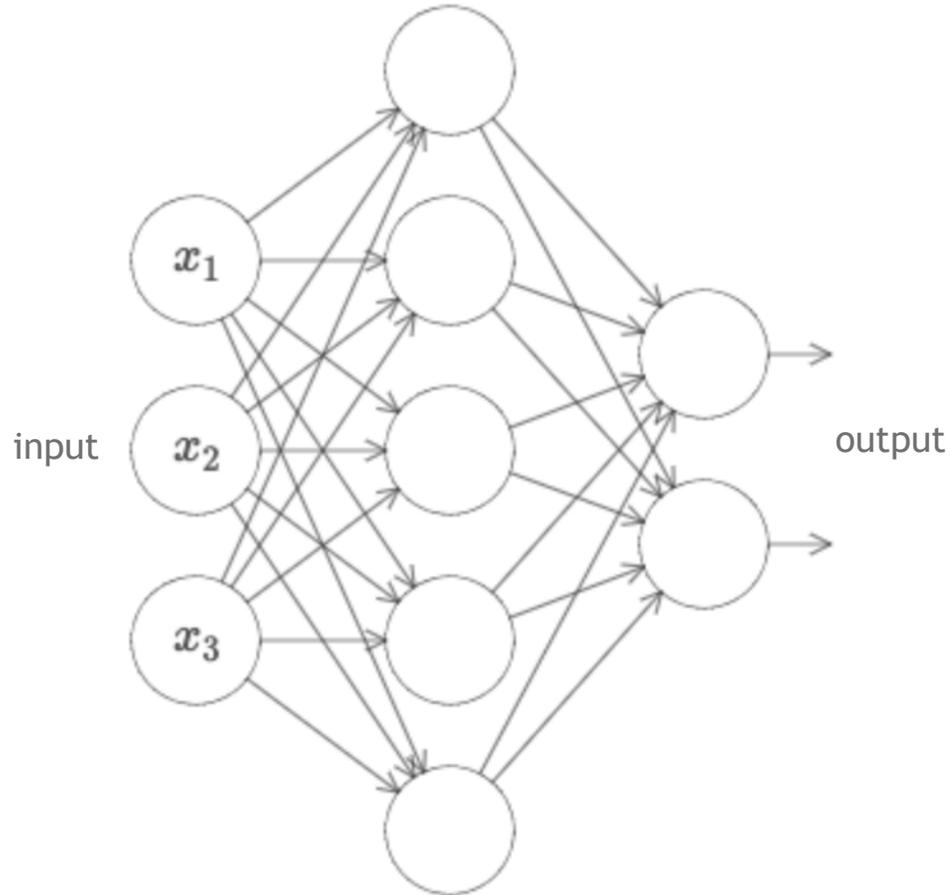
reLU, sigmoid activation functions



 activations



Training



$$\text{Loss } L(w_i\text{'s}) = \sum_{\text{inputs } x} l_x = \sum_{\text{batches } b} l_b$$

learning rate

$$\text{Gradient descend: } w_i \leftarrow w_i - \text{lr} \frac{\partial L}{\partial w_i}$$

performed “backwards” (*backprop*)

SGD: computations one batch at a time

$$w_i \leftarrow w_i - \text{lr} \frac{\partial l_b}{\partial w_i}$$

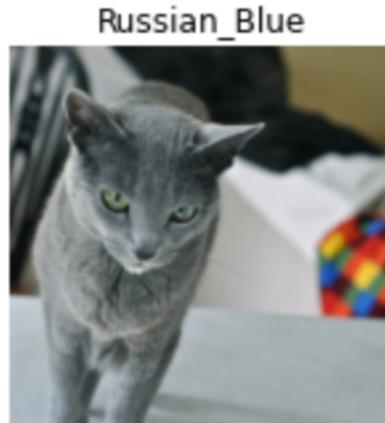
1 epoch : when all input/batches have been gone through

What are the NN Main Applications?

- **Computer Vision**
- **NLP**
- **Tabular Data**
- **Collaborative Filtering**

Main Applications - Computer Vision

single-label
classification



multi-label
classification

...
localization
object detection

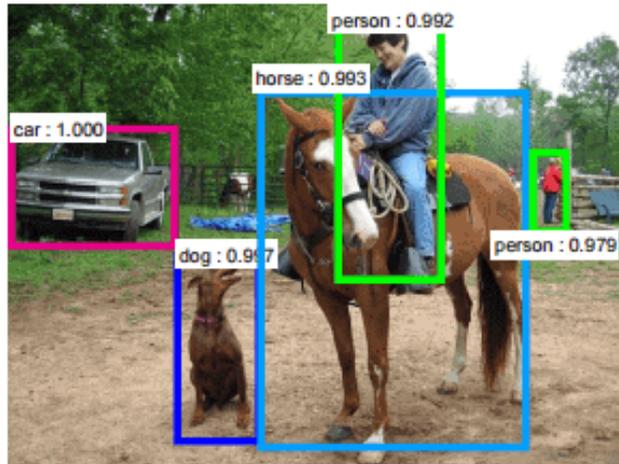
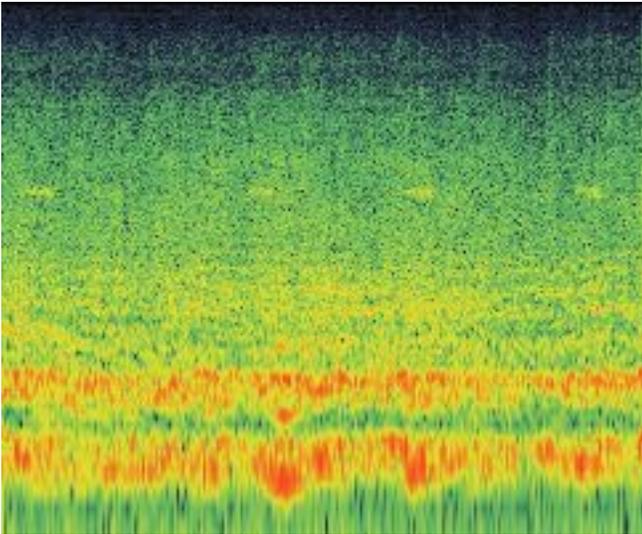


image segmentation

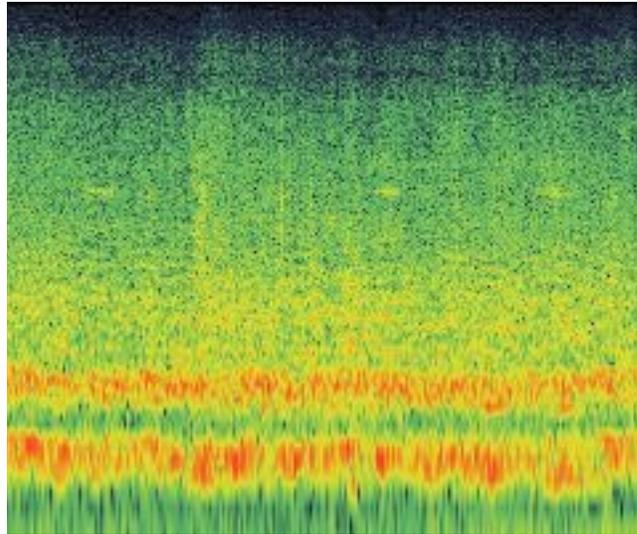


Conversion into image classification

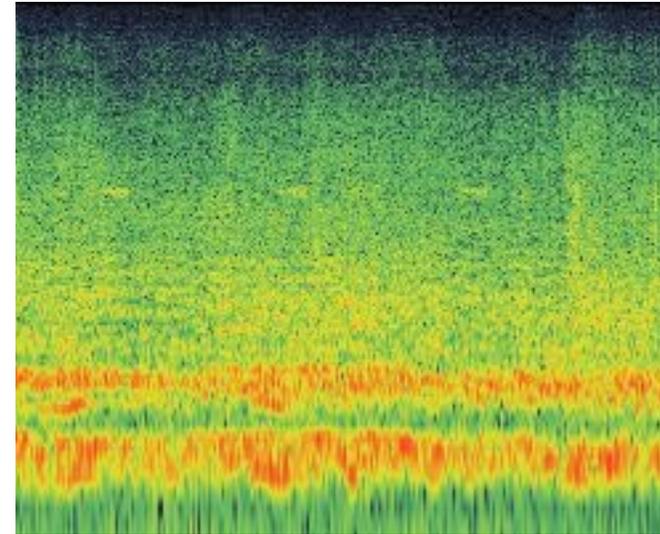
- **Detection of beehive's states...**



no queen



queen



swarming

emlyon student's project with the makers lab (Ammar, Hadjur, Radisson, Savinien IOT2019)

CNN - visualizing activations in CONVnets

Visualizing and Understanding Convolutional Networks

Matthew D. Zeiler and Rob Fergus

Dept. of Computer Science,
New York University, USA
{zeiler,fergus}@cs.nyu.edu

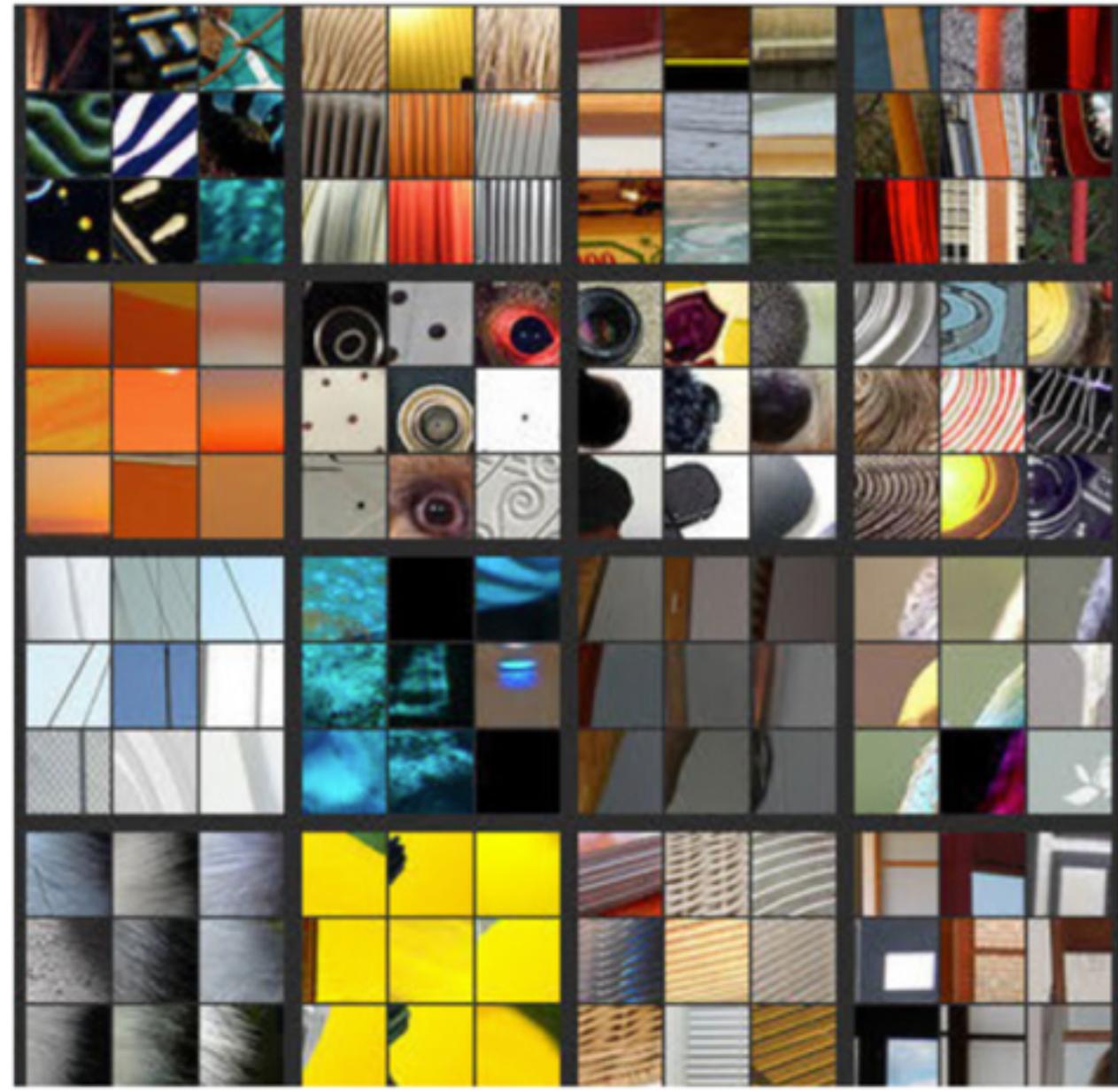
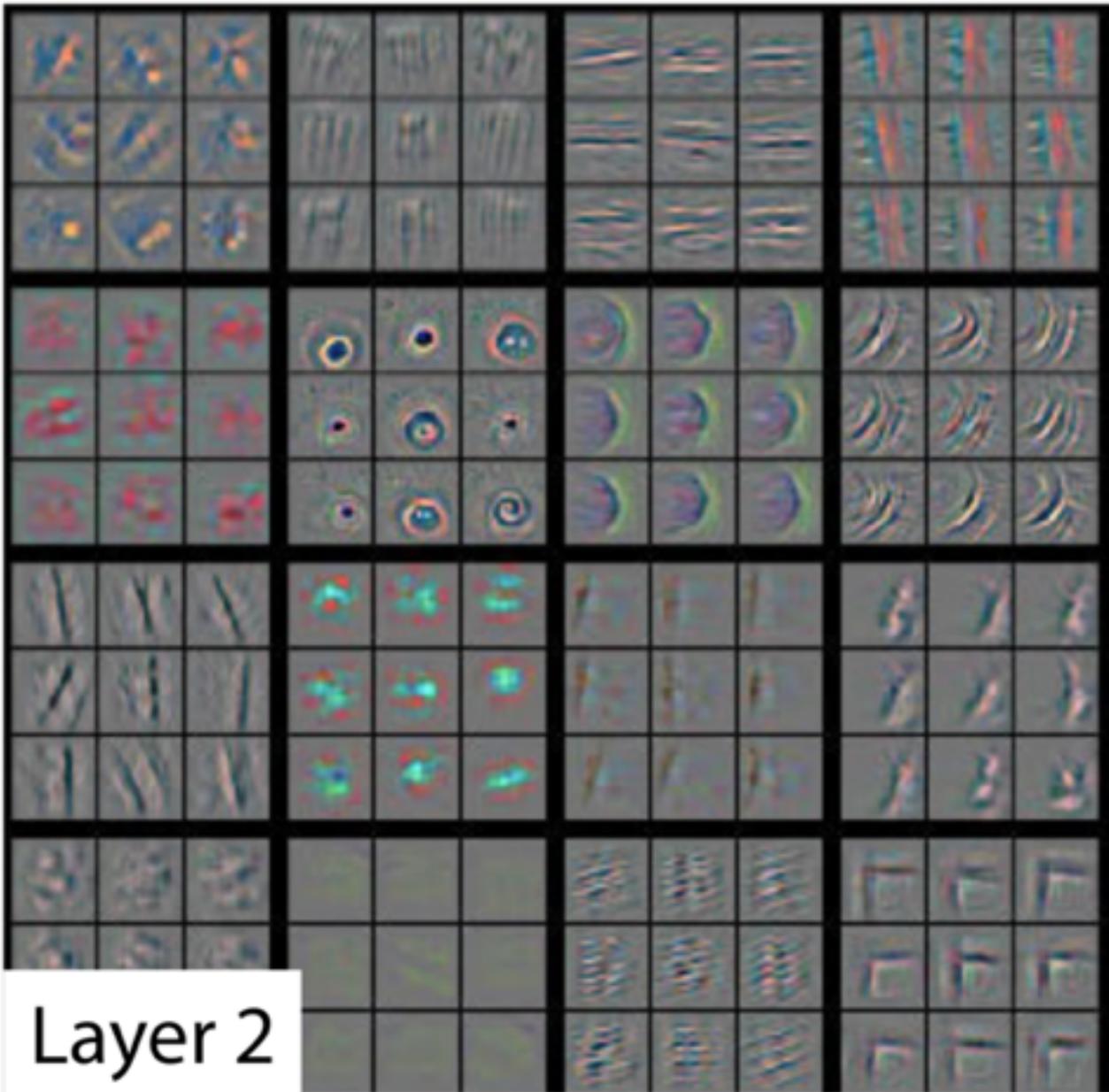
ECCV 2014



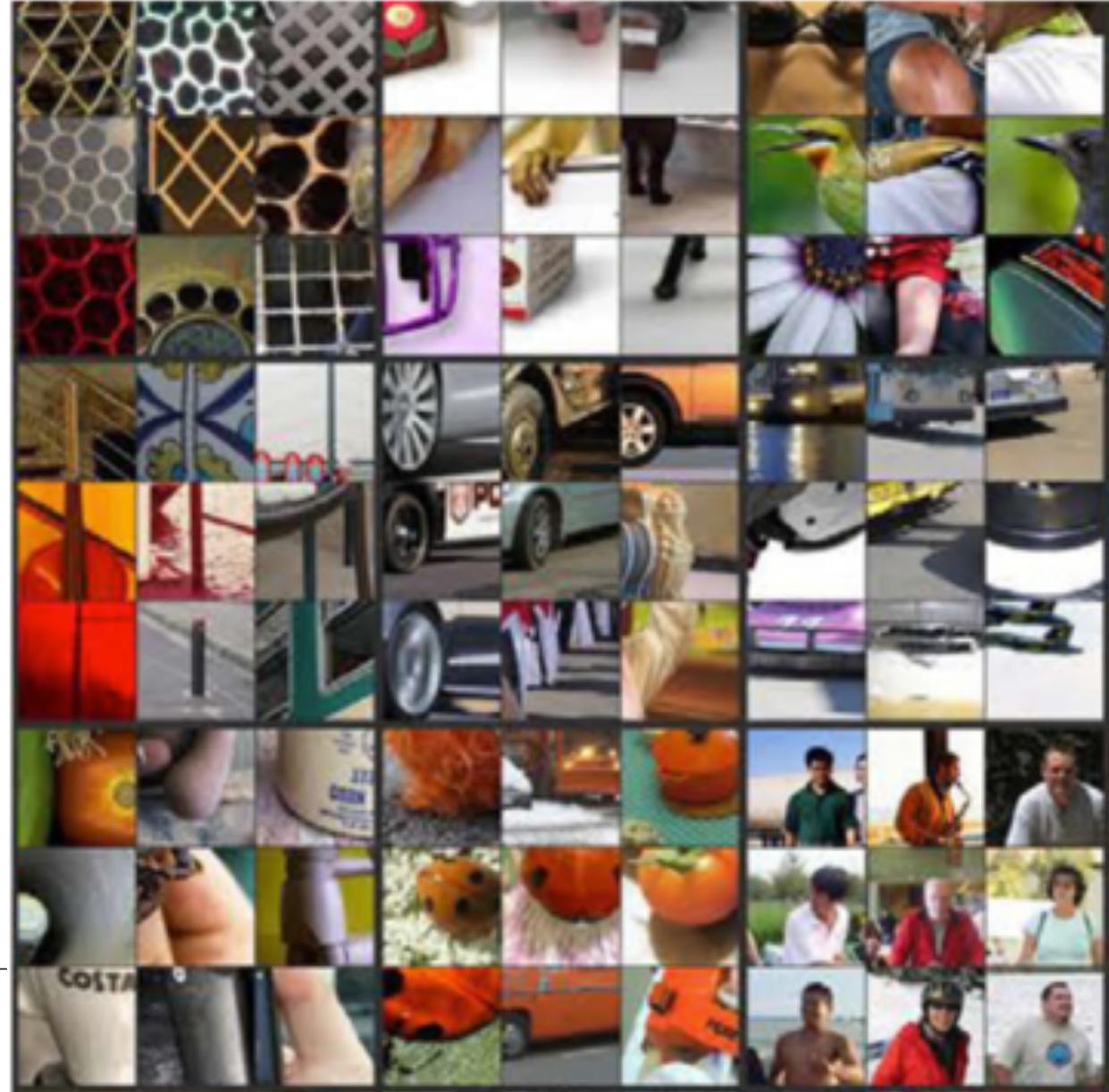
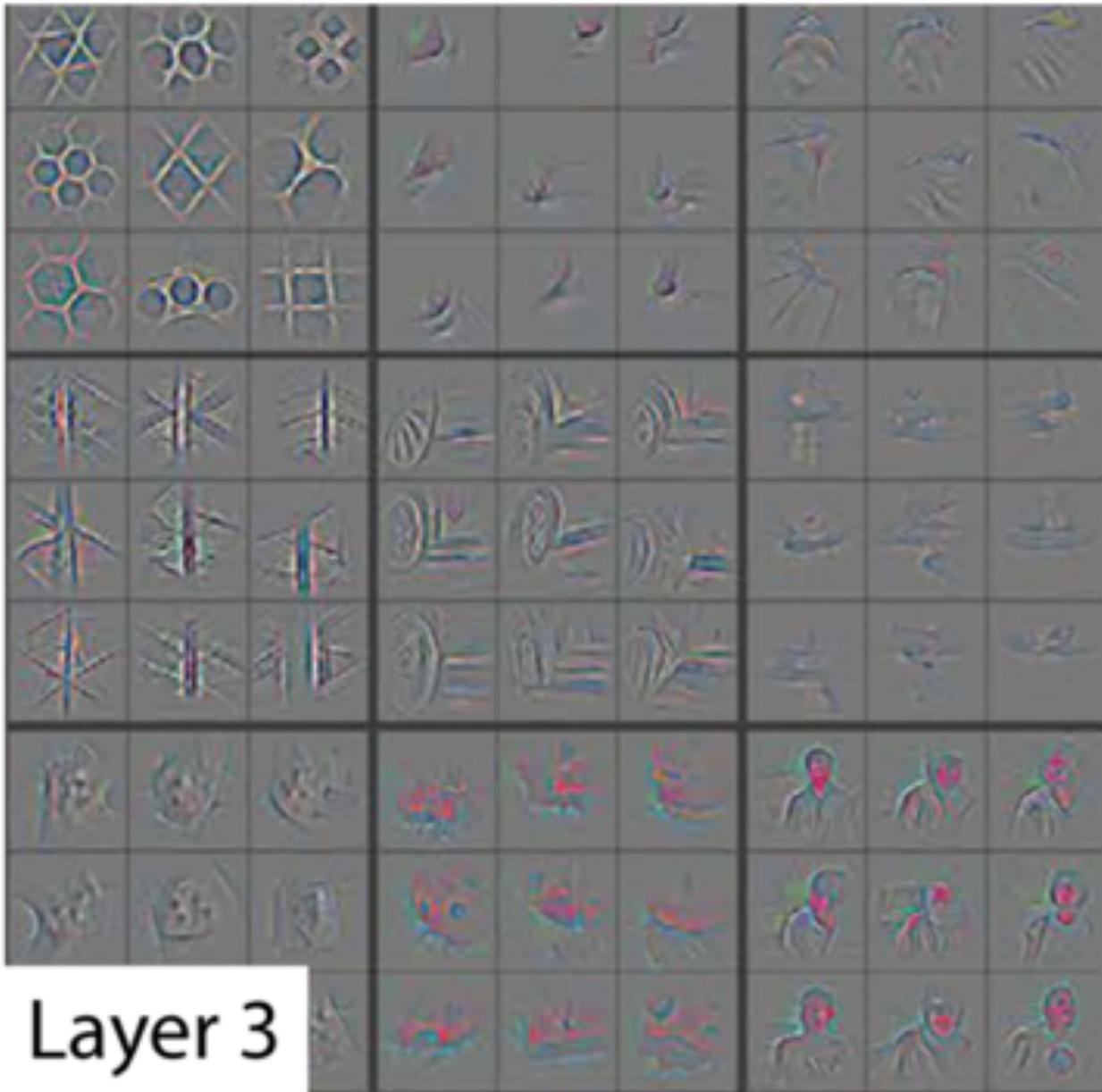
Layer 1



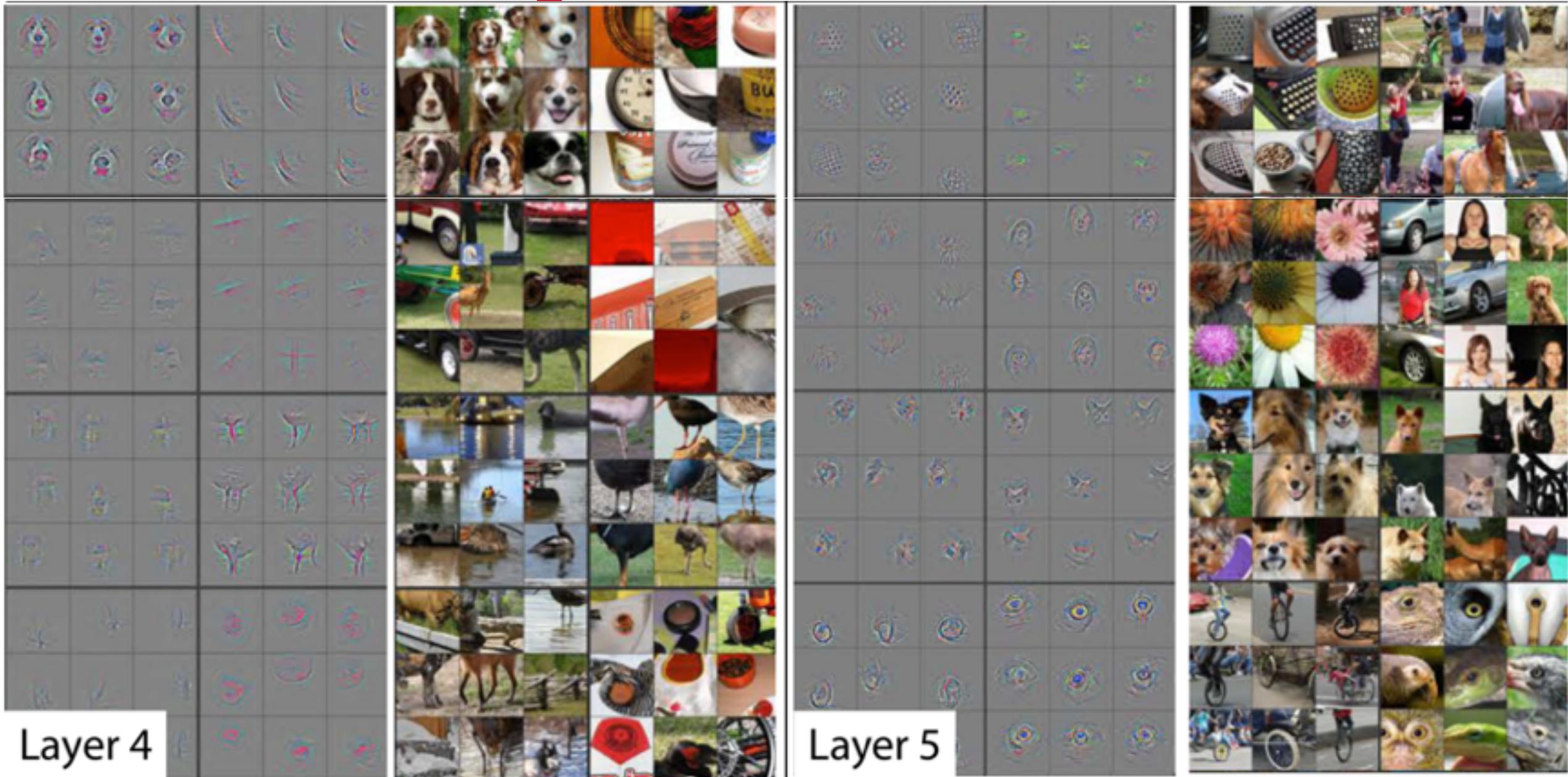
CNN - visualizing activations in CONVnets



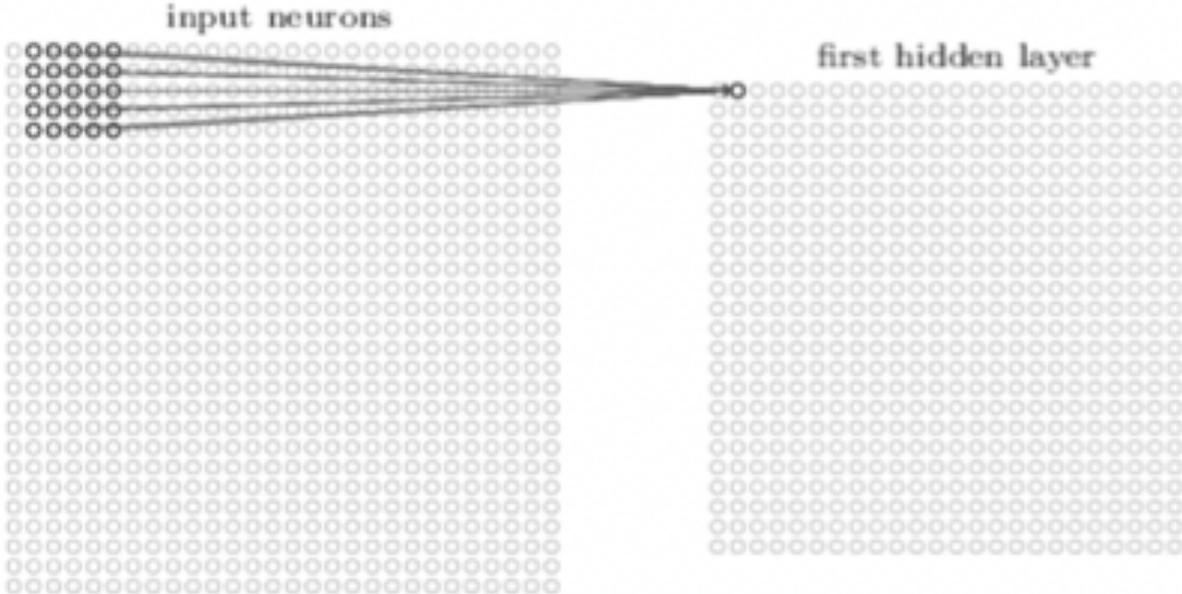
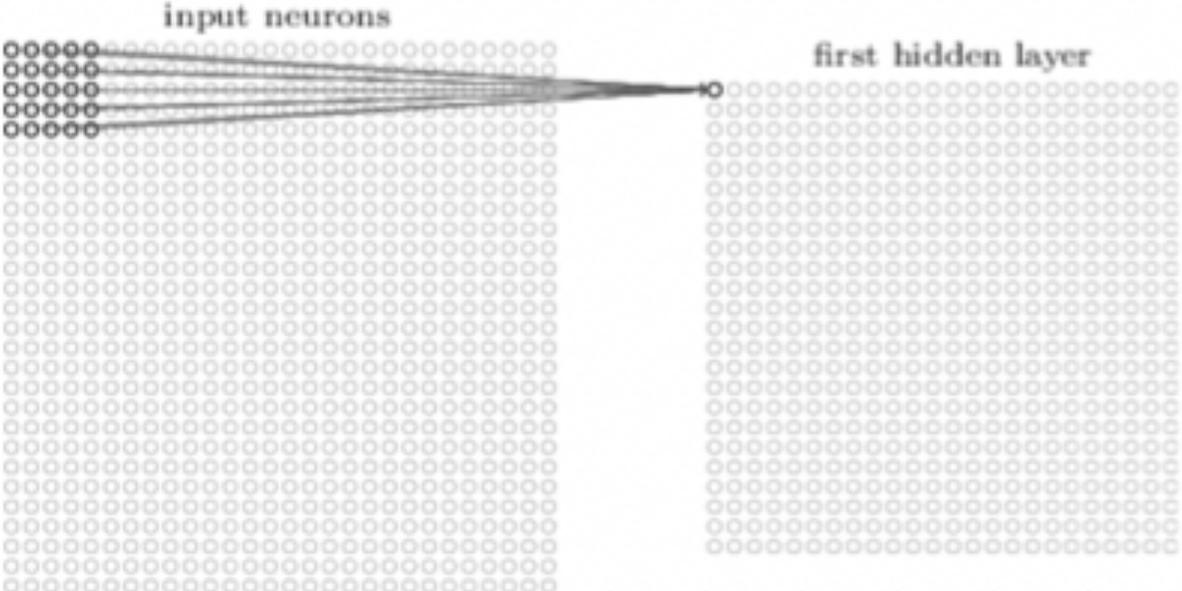
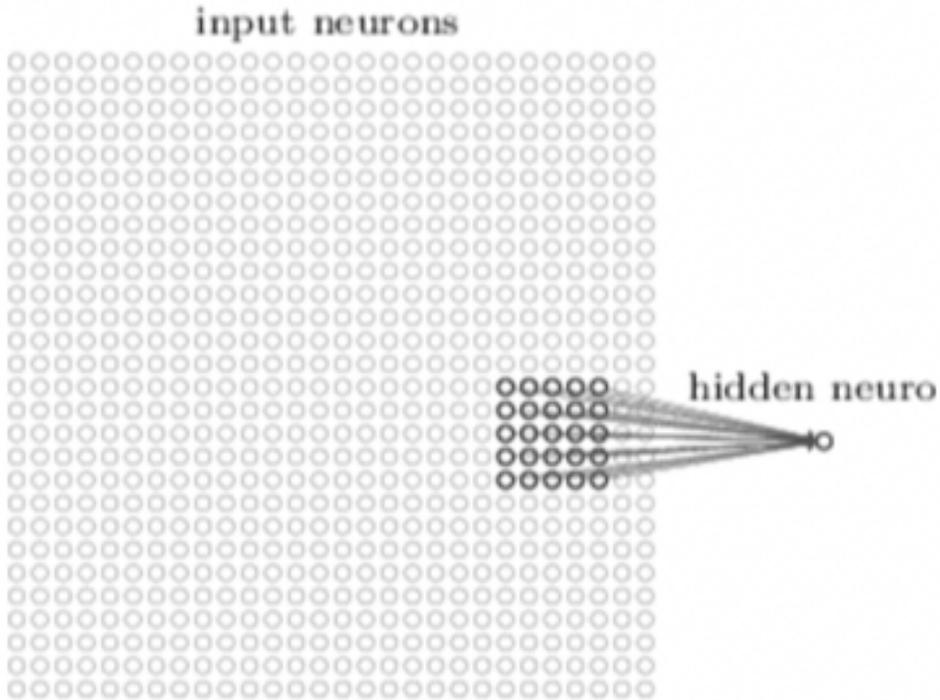
CNN - visualizing activations in CONVnets



CNN - visualizing activations in CONVnets

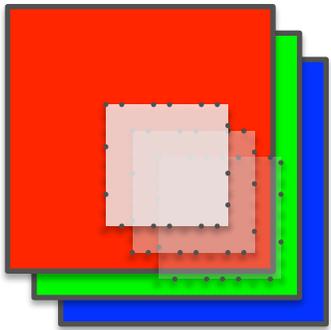


CNN - convolution



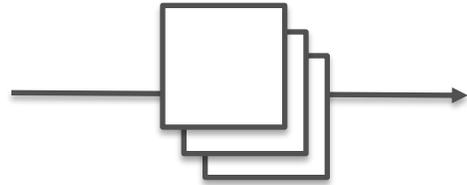
CNN - feature maps

input image



3x5x5
3 channels
RGB

kernel



3x3x3

channel output

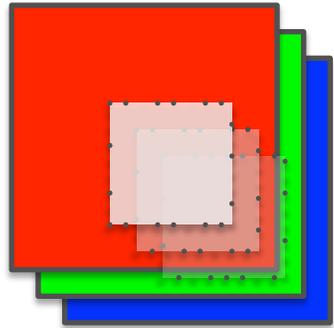


5x5

learning 1 elementary pattern,
e.g. edge, color gradient, etc.

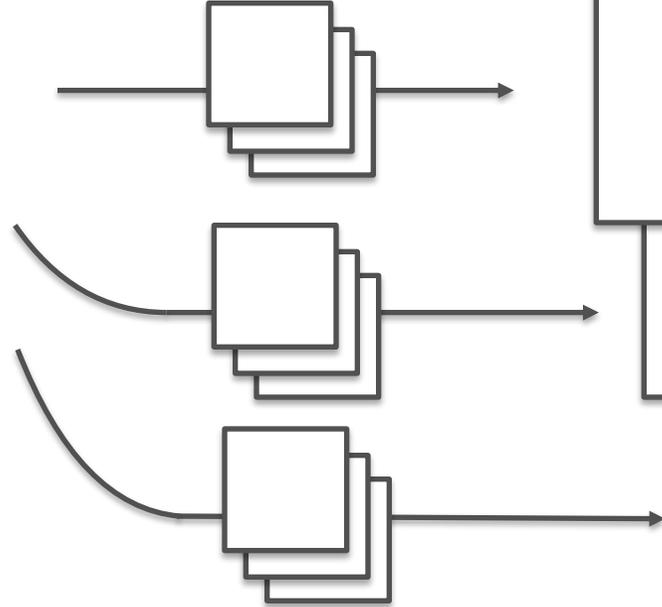
CNN - feature maps

input image

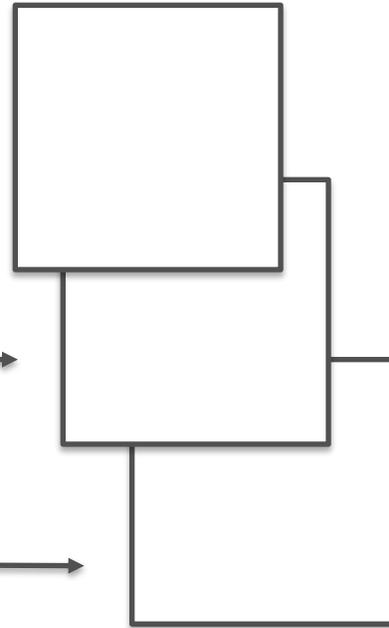


3x5x5
3 channels
RGB

kernels

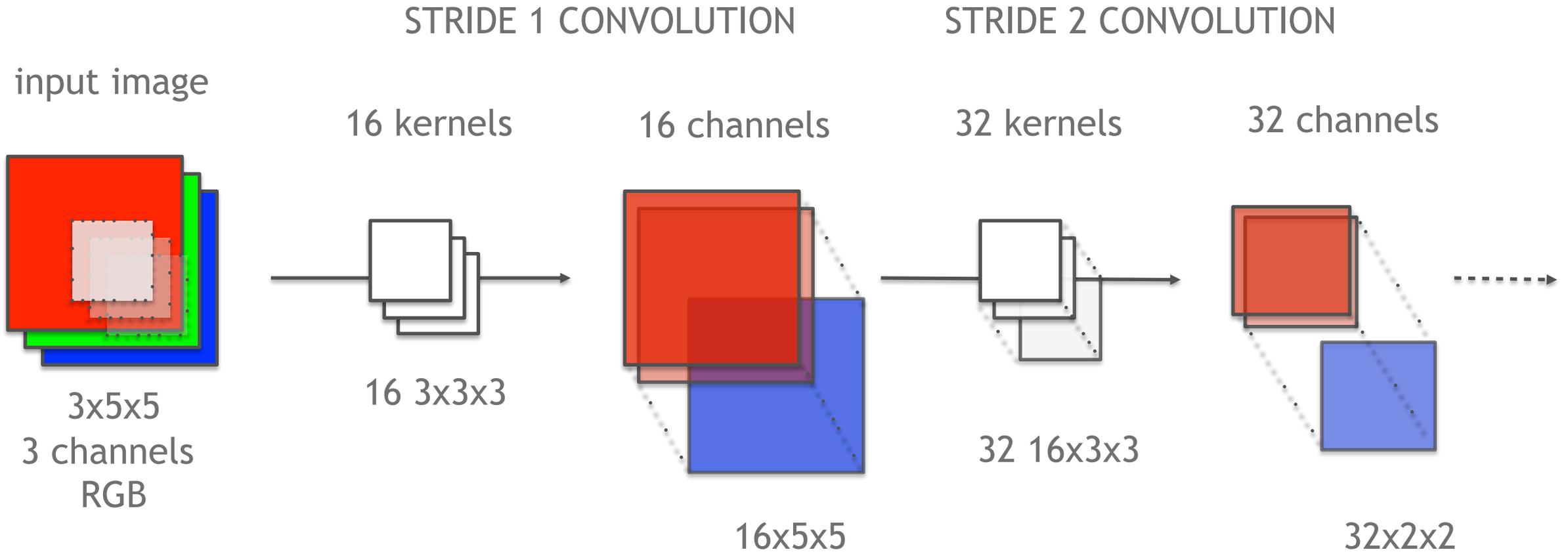


channel outputs



learning 3 elementary patterns

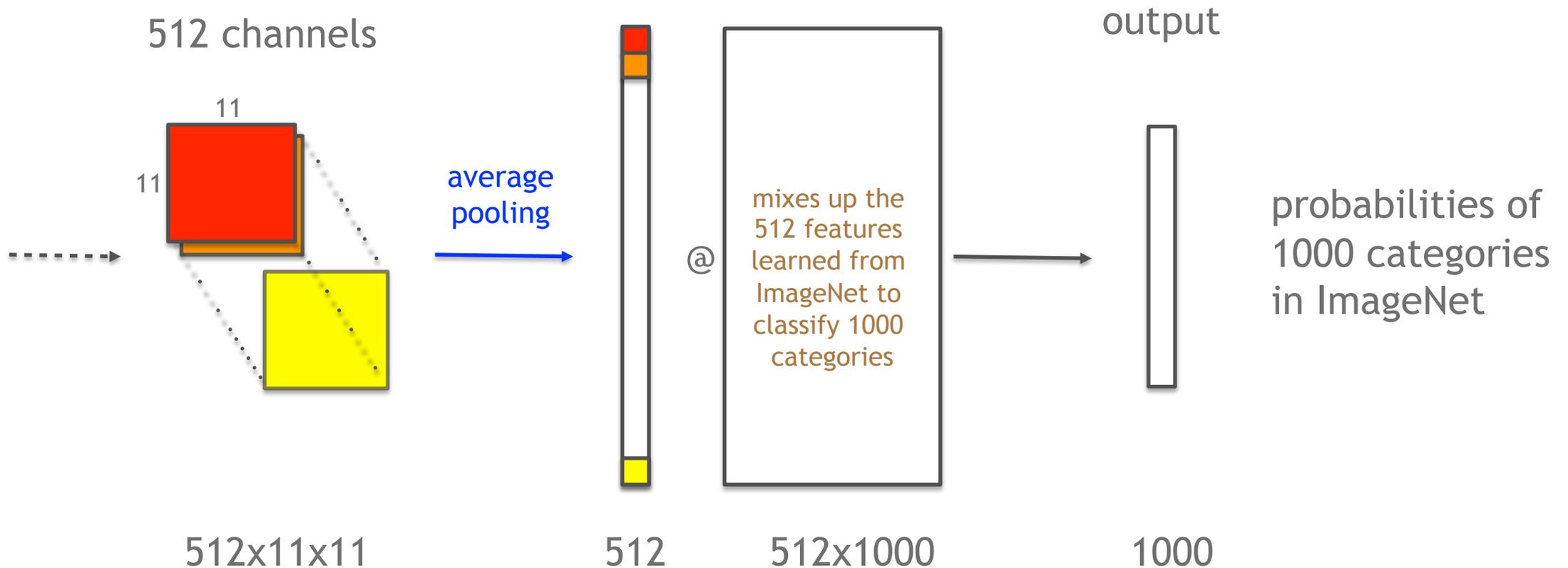
CNN - CONVnet



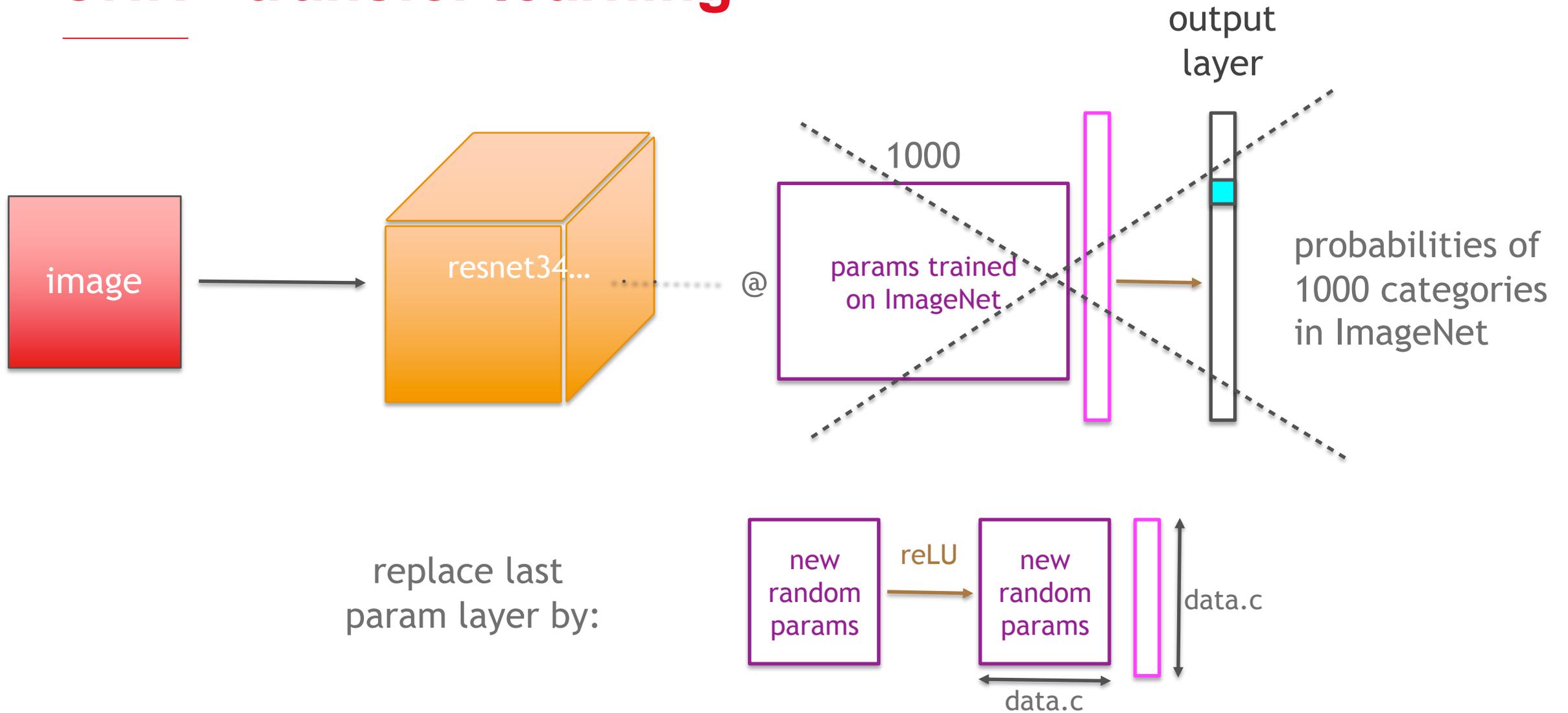
*halves picture sizes
& doubles channels*

CNN - CONVnet

alternate layers of stride 1 & stride 2 convs, to end up with 512 channels

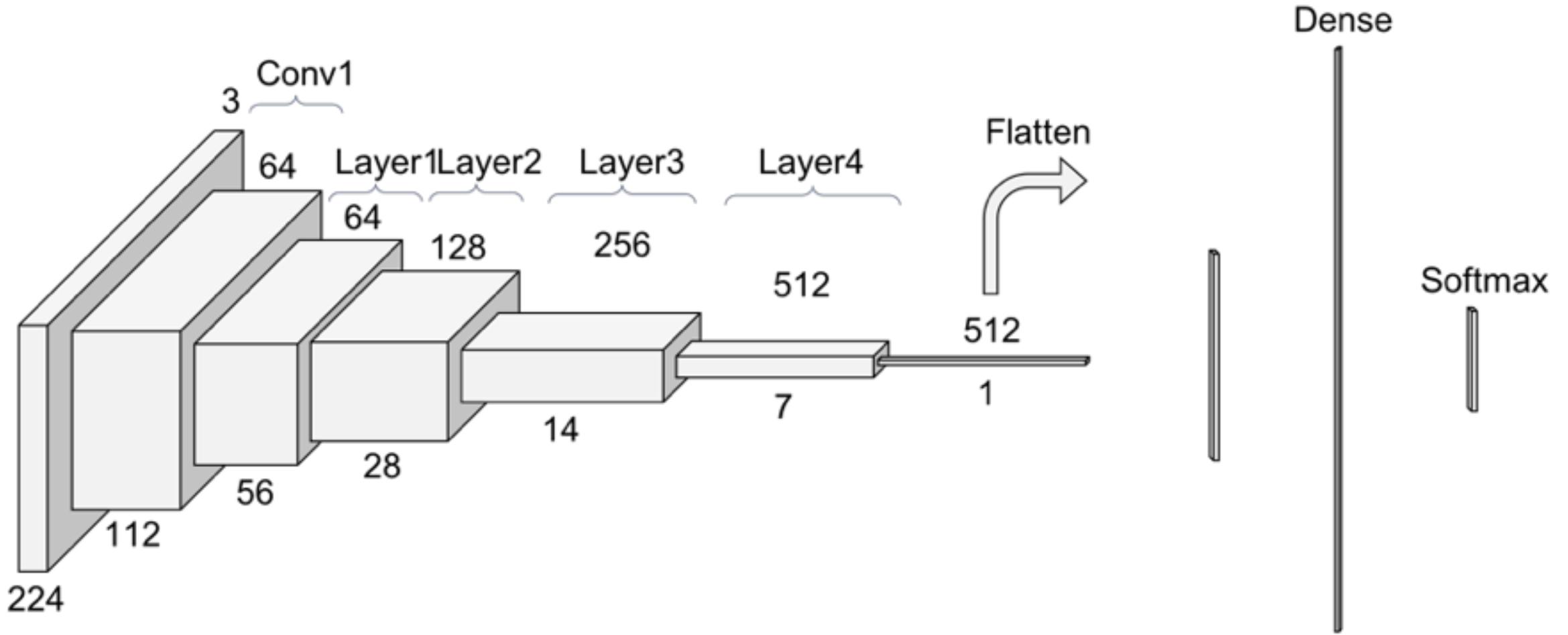


CNN - transfer learning

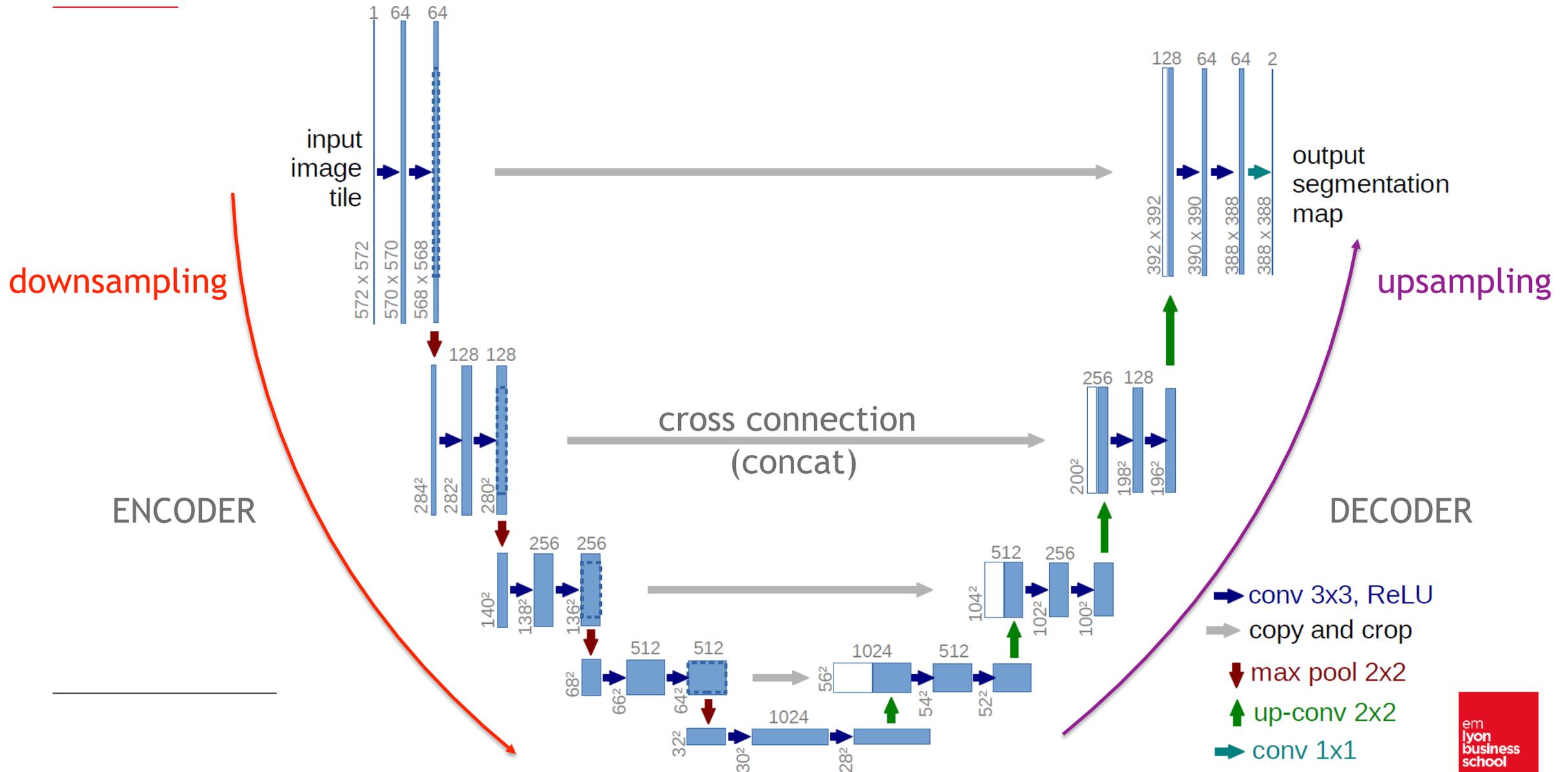


CNN - resnets

alternate layers of stride 1 & stride 2 convs, to end up with 512 channels



CNN - Unets



CNN - resnets

Deep Residual Learning for Image Recognition

IEEE CVPR 2016

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

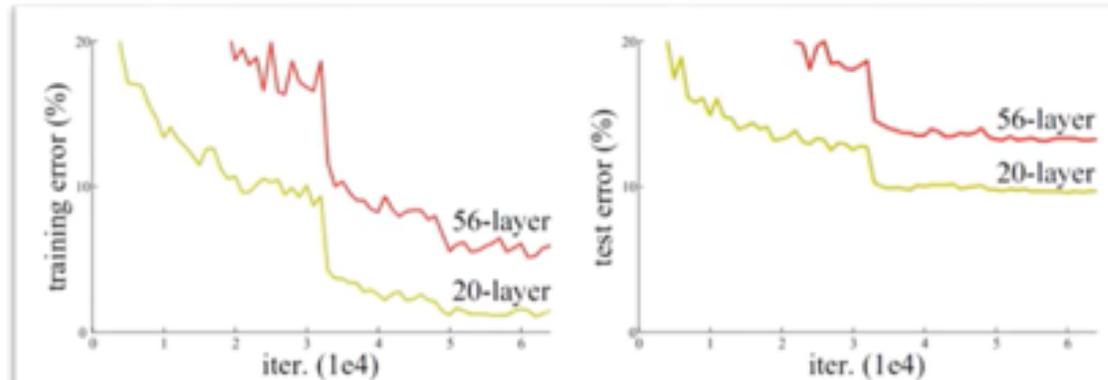


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

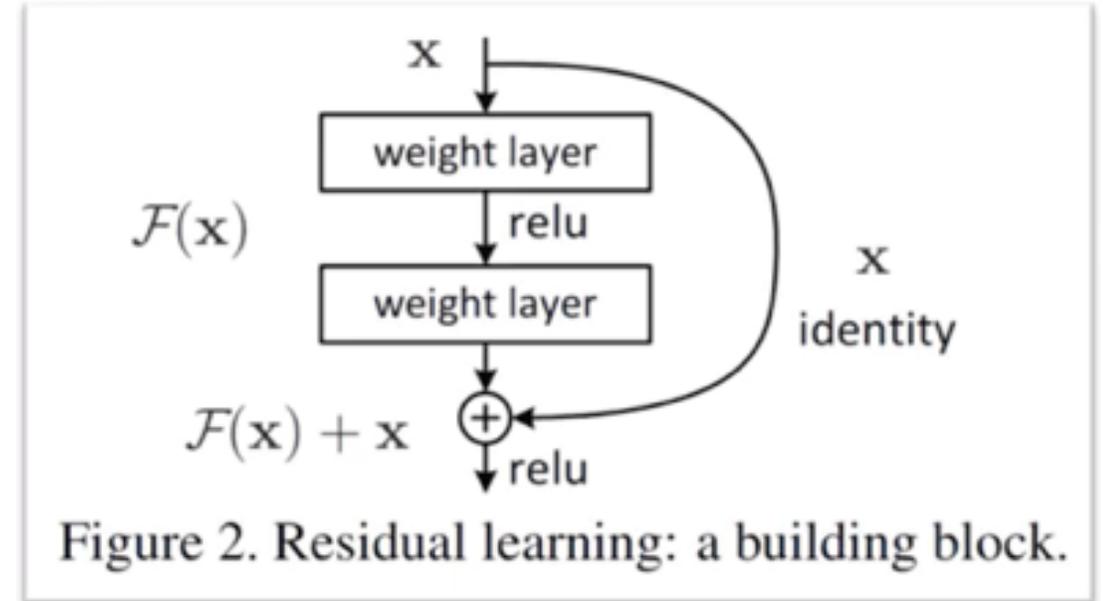
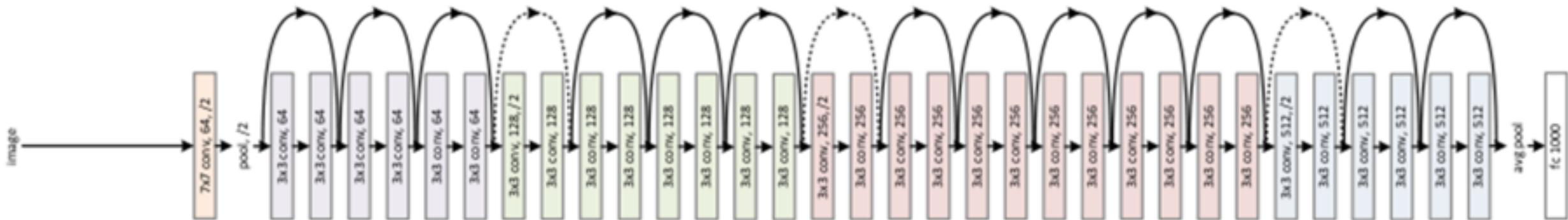


Figure 2. Residual learning: a building block.

$$\text{output} = x + \text{conv2}(\text{conv1}(x))$$

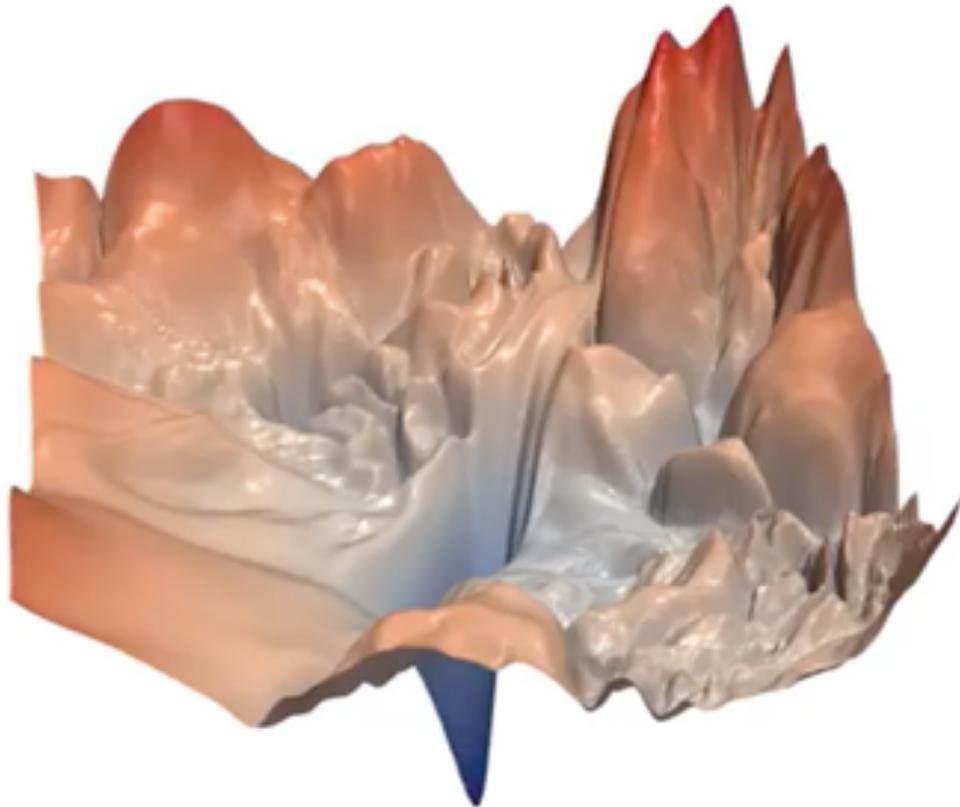


CNN - resnets

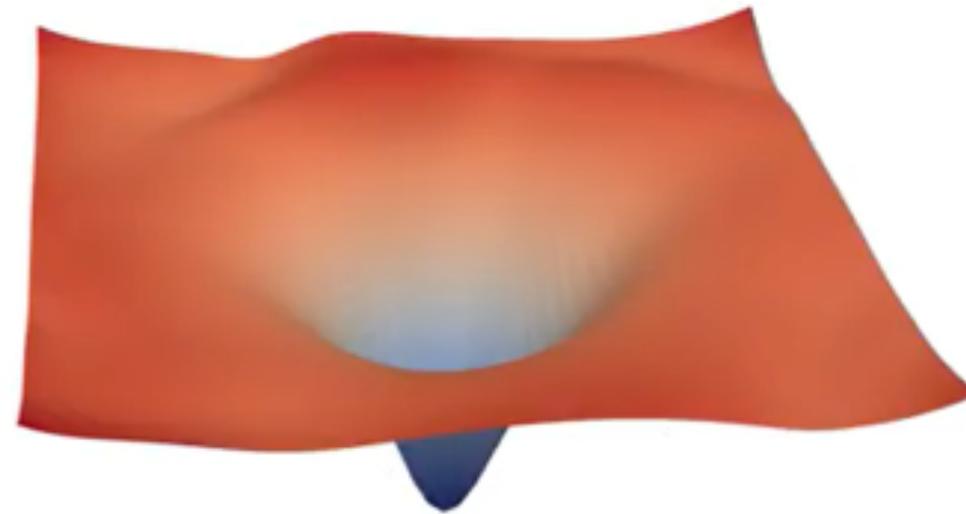
Visualizing the Loss Landscape of Neural Nets

NIPS'18

Hao Li¹, Zheng Xu¹, Gavin Taylor², Christoph Studer³, Tom Goldstein¹



(a) without skip connections



(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

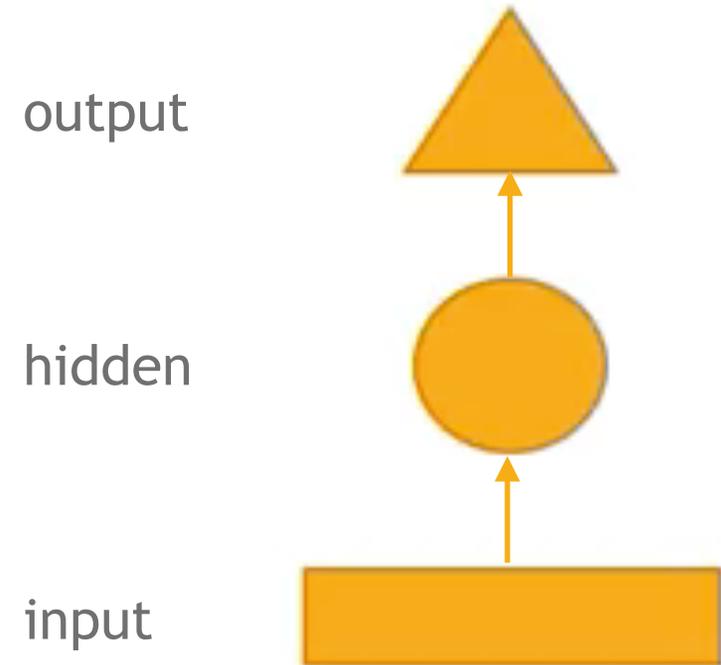
Main Applications - NLP

See Julien's talk :)

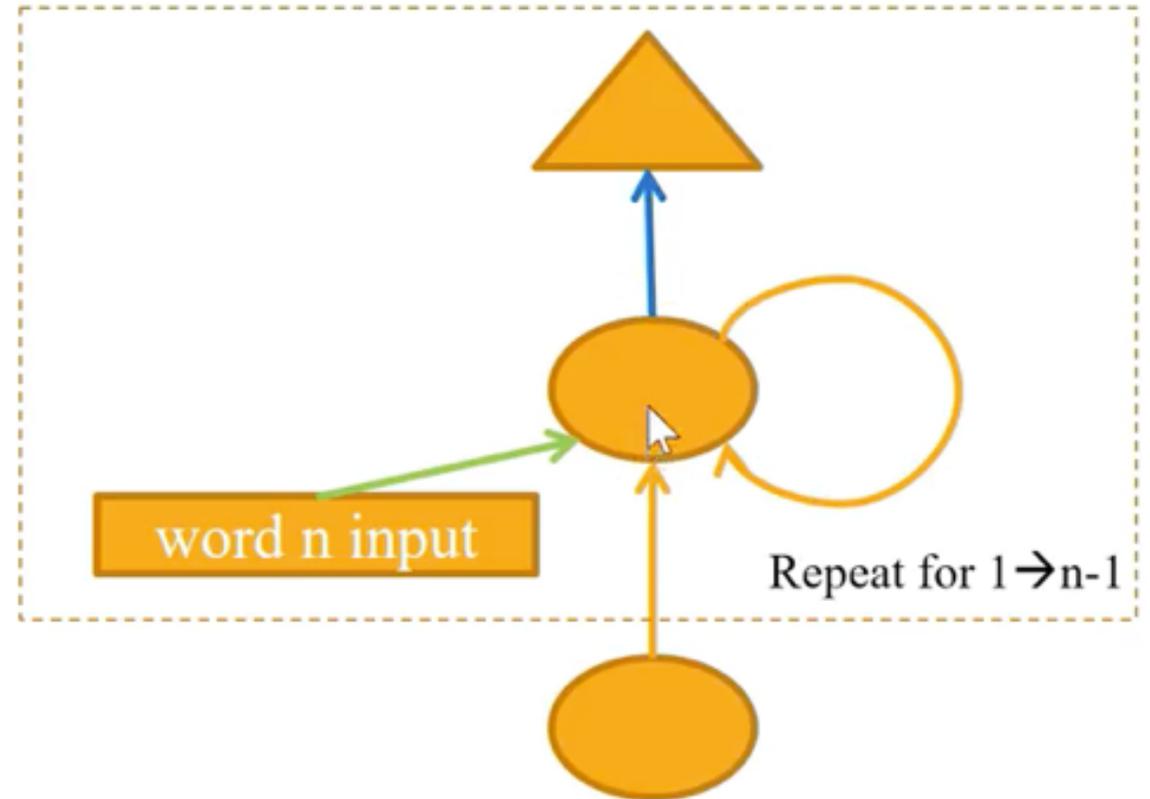
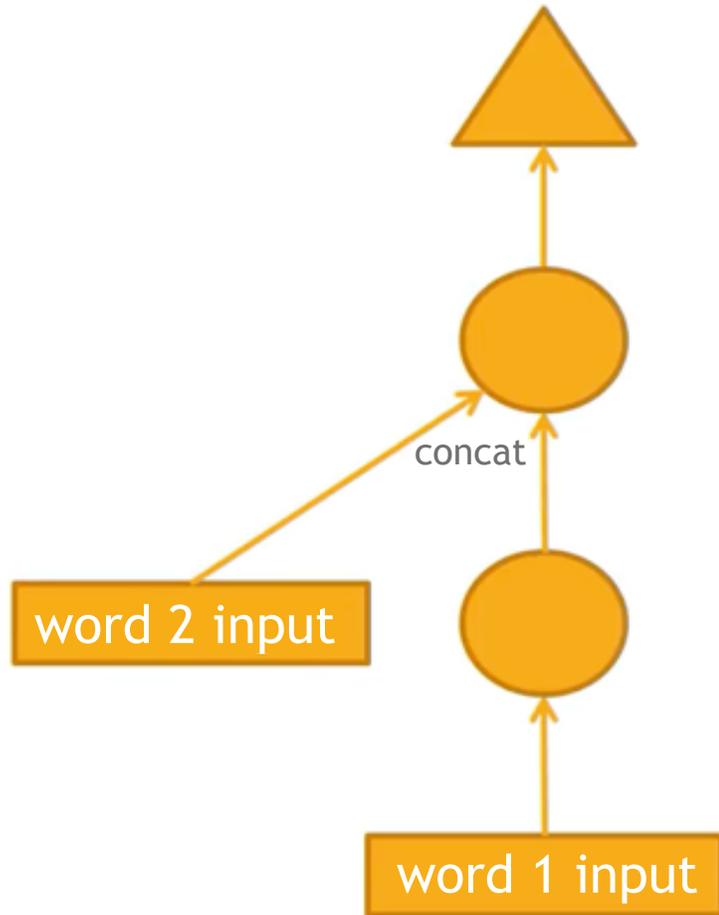
Recurrent Neural Nets (RNNs)

What is an RNN?

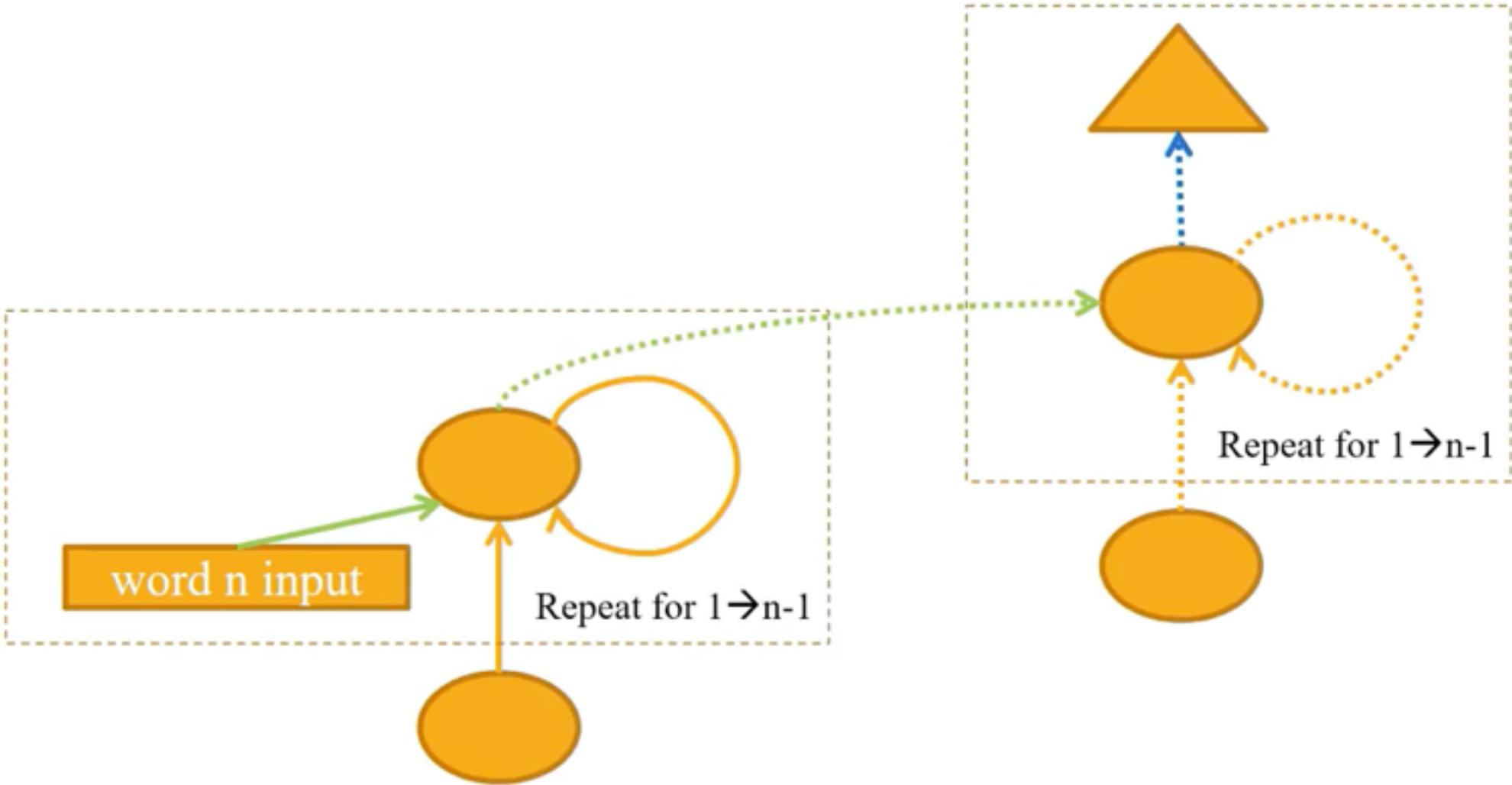
essentially a *fully connected NN*
trained in a specific way



RNNs

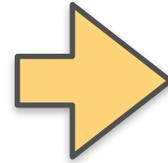


Stacked RNNs

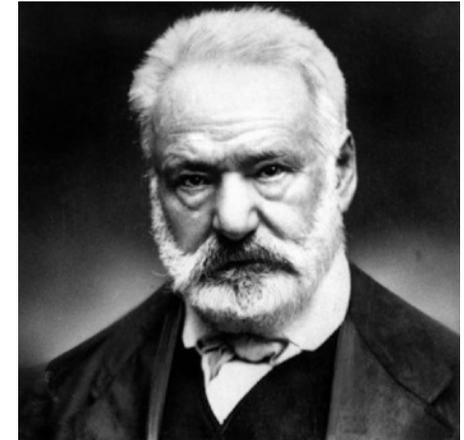


Transfer Learning for NLP

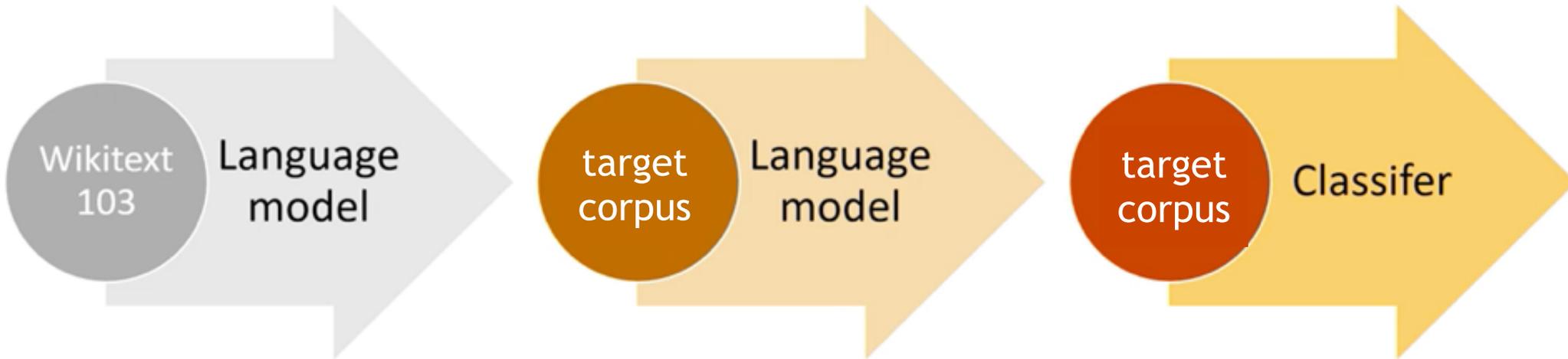
La terre a vu jadis errer des paladins;
Ils flamboyaient ainsi que des éclairs soudains,
Puis s'évanouissaient, laissant sur les visages
La crainte, et la lueur de leurs brusques passages ;



or



???



Main Applications - Tabular Data

Embeddings & NN for tabular data

Entity Embeddings of Categorical Variables

Cheng Guo* and Felix Berkhahn†
Neokami Inc.
(Dated: April 25, 2016)

ranked 2nd at Rosmann
kaggle competition

method	MAPE	MAPE (with EE)
KNN	0.290	0.116
random forest	0.158	0.108
gradient boosted trees	0.152	0.115
neural network	0.101	0.093

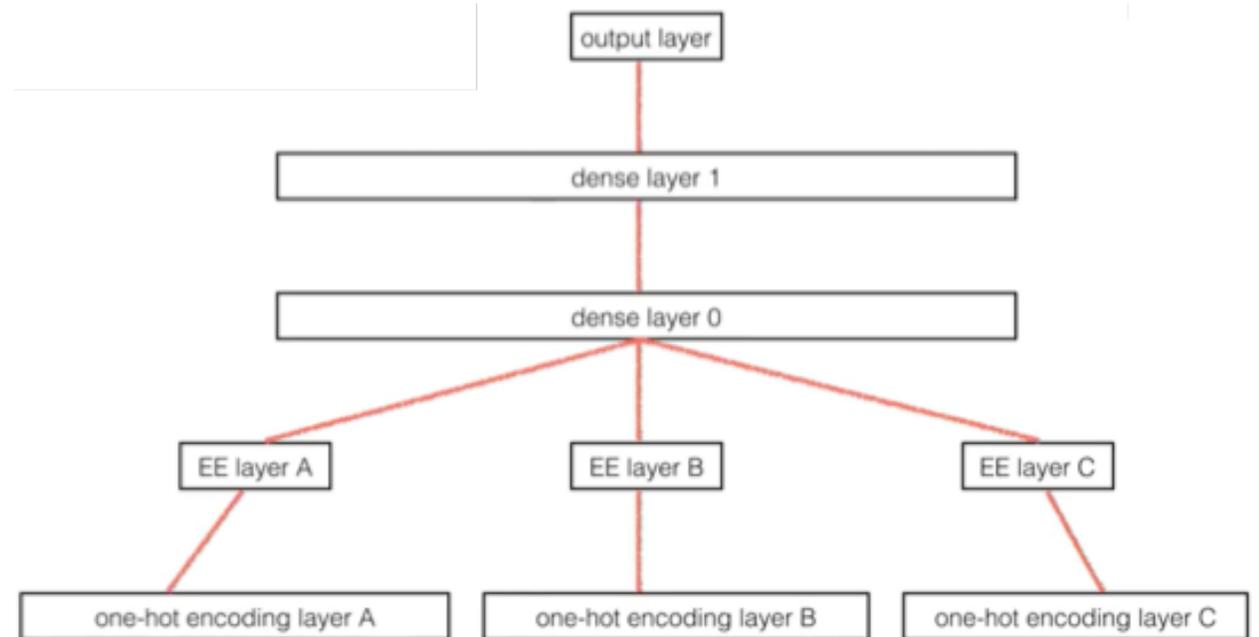
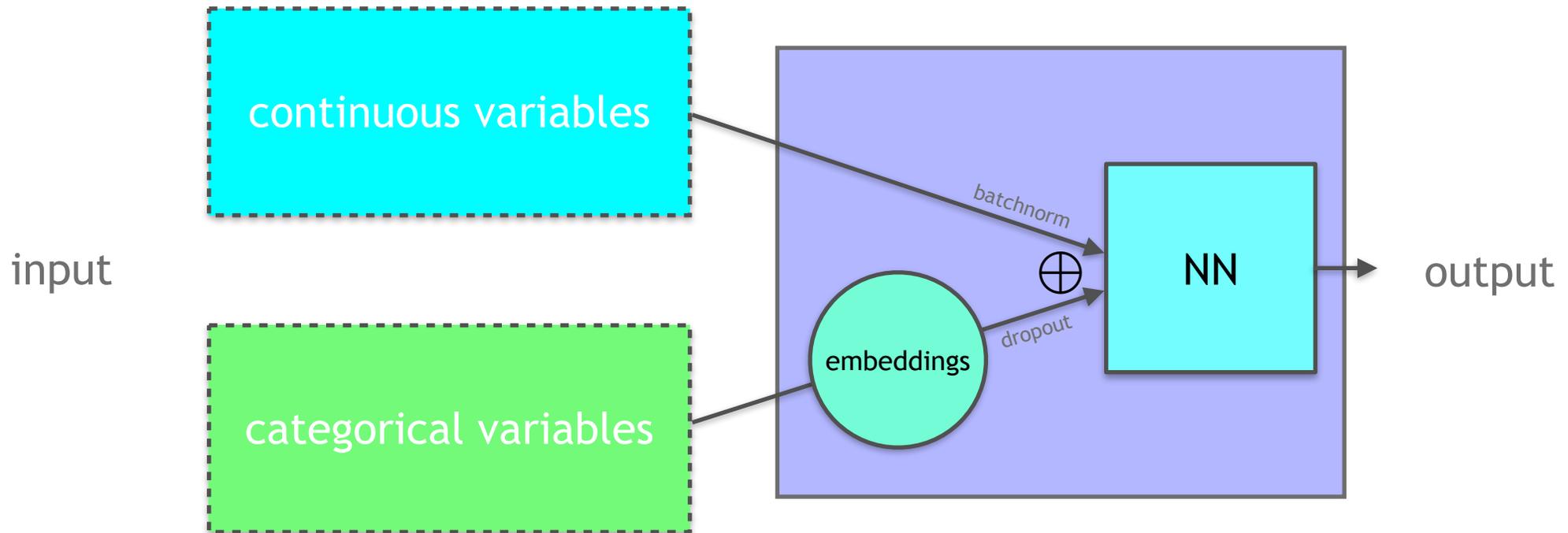
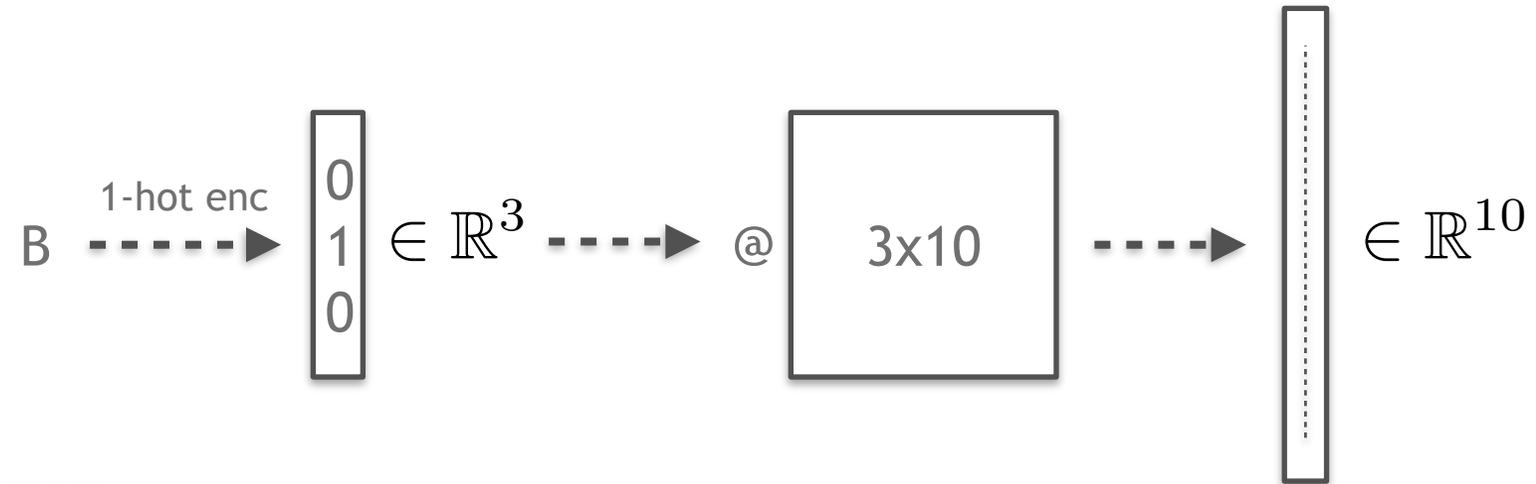


FIG. 1. Illustration that entity embedding layers are equivalent to extra layers on top of each one-hot encoded input.

Embeddings

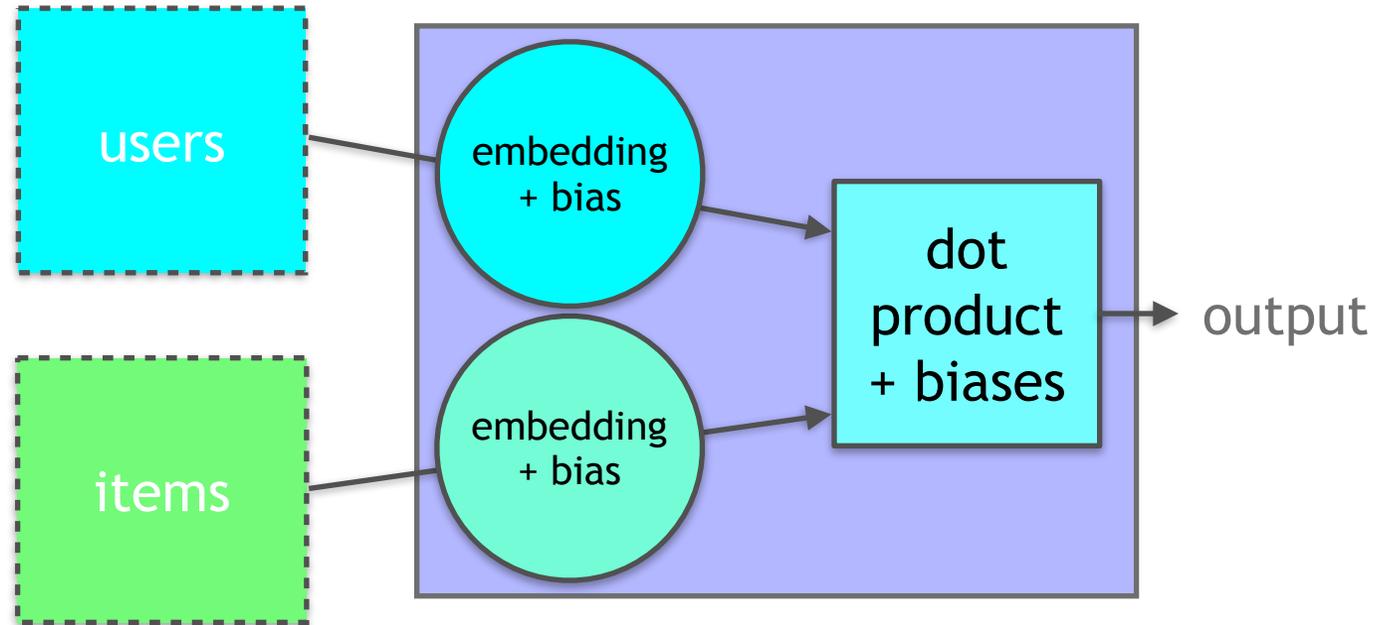
10-dim embedding of a card 3 categorical variable (A,B,C):



Main Applications - Collaborative Filtering

userId	movieId	review
1	302	★
2	5709	★★★
...

input

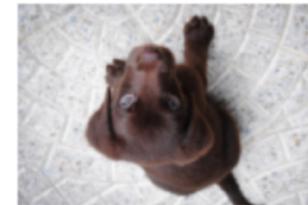
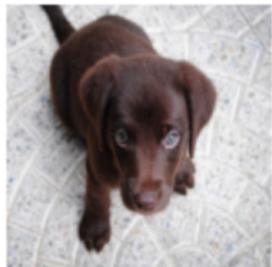
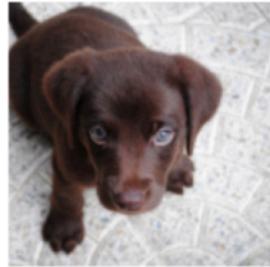
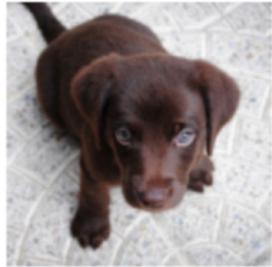
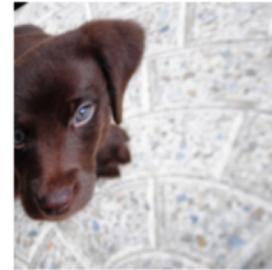
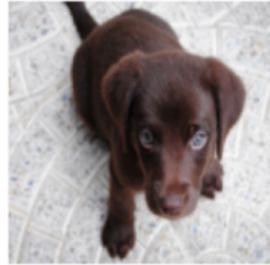


In practice

- **Data Augmentation**
- **Gradients Problems & Init**
- **Learning Rates**
- **Overfitting**
- **GPUs & NN libraries**

Data Augmentation

resize, crop, pad, squish, warp, zoom, flip, rotate,...
vary lightning, grayscale, contrast, blur, *resolution*,...



Gradients Problems & Initialization

problems

vanishing & exploding signals ?

vanishing & exploding gradients ?

culprits

NN depth

choice of activations

bad init

cures...

resnets
(skip connections)

LSTM

params init

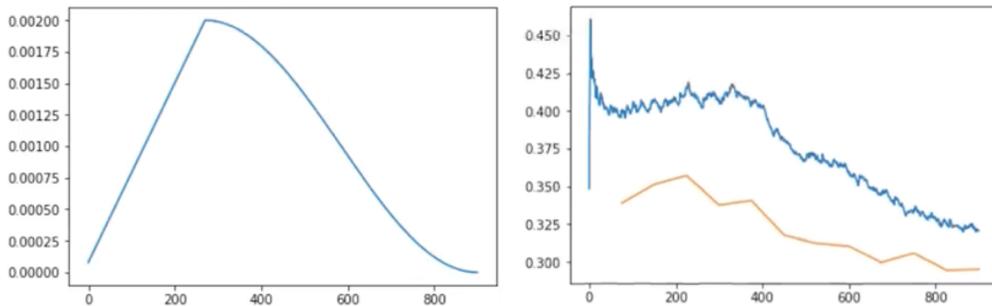
Learning rates

dynamic learning rates:

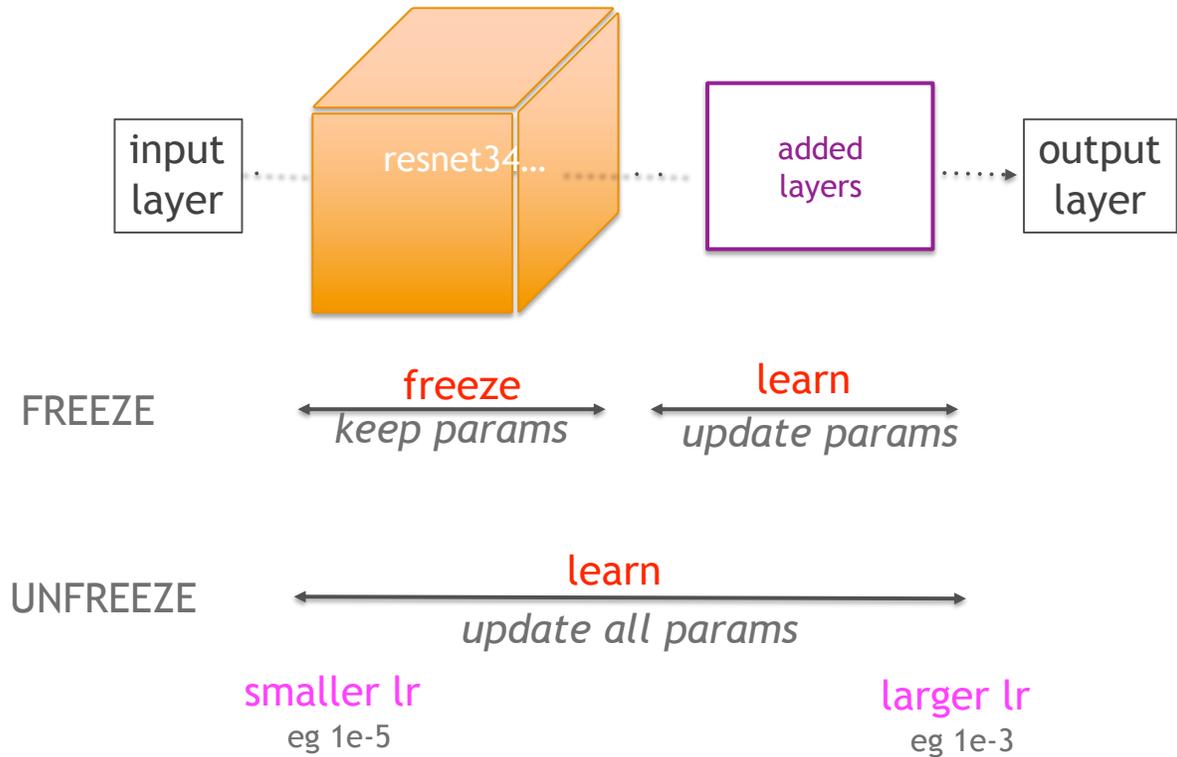
- momentum $M_t = (1 - \beta)\nabla_t + \beta M_{t-1}$
- RMSprop $R_t = (1 - \beta)\nabla_t^2 + \beta R_{t-1}$
- Adam:

$$w_i^{(t)} \leftarrow w_i^{(t-1)} - \text{lr} \frac{M_t}{\sqrt{R_t}}$$

learning rate annealing



discriminative learning rate



Overfitting

Deep NN have **HUGE** numbers of learnable parameters !

eg models widely used today:

- resnet50 model for image recognition: over **23 millions**
- “BERT base” model for NLP: over **110 millions**

Standard ways to reduce overfitting:

- weight decay (typically L2-reg)
- dropout

Dropout

“I went to my bank. The tellers kept changing and I asked one of them why. He said he didn’t know but they got moved around a lot. I figured it must be because it would require cooperation between employees to successfully defraud the bank. This made me realize that randomly removing a different subset of neurons on each example would prevent conspiracies and thus reduce overfitting”

Hinton: Reddit AMA

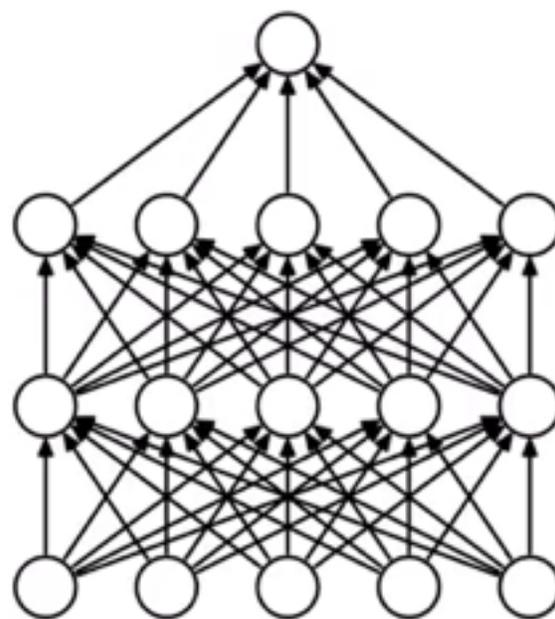
Dropout

Dropout: A Simple Way to Prevent Neural Networks from Overfitting

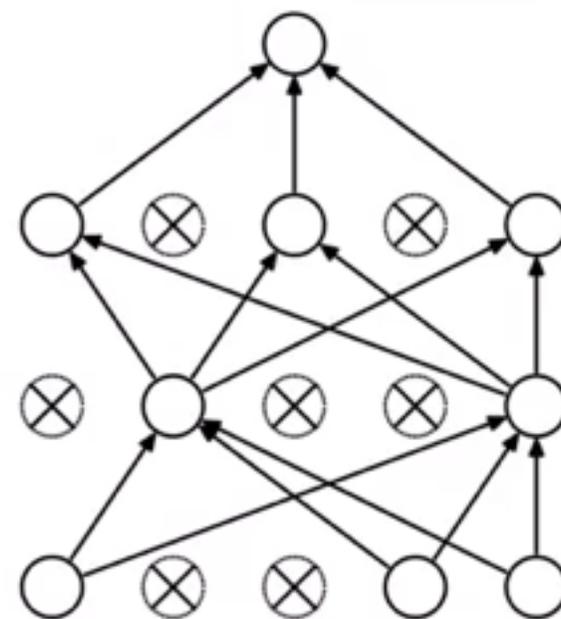
JMLR 2014

Nitish Srivastava
Geoffrey Hinton
Alex Krizhevsky
Ilya Sutskever
Ruslan Salakhutdinov

NITISH@CS.TORONTO.EDU
HINTON@CS.TORONTO.EDU
KRIZ@CS.TORONTO.EDU
ILYA@CS.TORONTO.EDU
RSALAKHU@CS.TORONTO.EDU



(a) Standard Neural Net



(b) After applying dropout.

Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

Batchnorm

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

ICML'2015

Sergey Ioffe

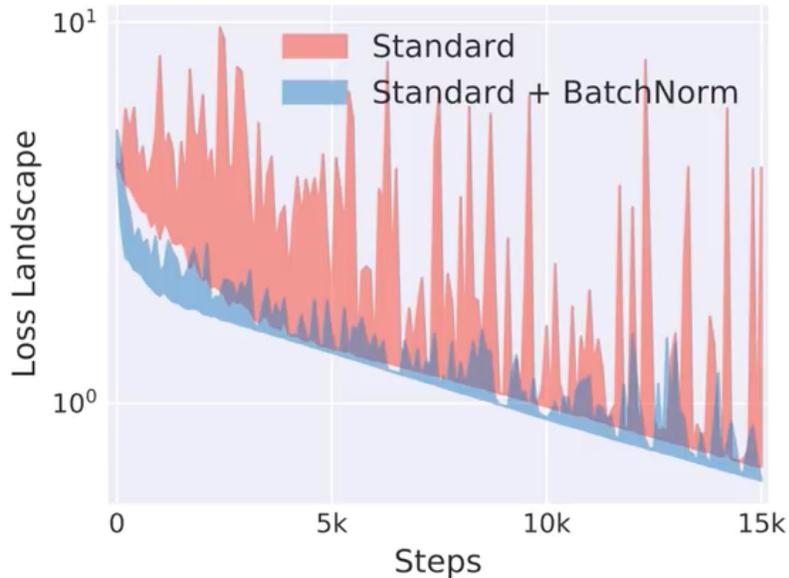
Google Inc., sioffe@google.com

Christian Szegedy

Google Inc., szegedy@google.com

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$


batchnorm
flattens out
loss surface
=> higher lr
=> faster train

Eg at current stage of learning, last activations predict outputs

$$\hat{y} = f(w_i \text{'s}; x)$$

in (-1,1)... but target is in range 1 to 5!

BN scales and shifts output
via 2 extra learnable parameters

$$\hat{y} = \gamma f(w_i \text{'s}; x) + \beta$$

S. Santukar *et al.* MIT

GPUs & NN libraries

GPUs

free

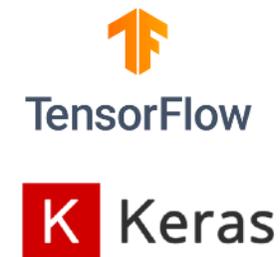


professional solutions



300\$ free credit
education packs
0.3-1.2€/hour

NN libraries



VS



early
makers

em
lyon
business
school