

1



An Introduction to Convolutional Neural Networks

Yuan YAO
HKUST

Summary

- ▶ We had covered so far
 - ▶ Linear models (linear and logistic regression) – always a good start, simple yet powerful
 - ▶ Model Assessment and Selection – basics for all methods
 - ▶ Trees, Random Forests, and Boosting – good for high dim mixed-type heterogeneous features
 - ▶ Support Vector Machines – good for small amount of data but high dim geometric features
- ▶ Next, neural networks for unstructured data (image, language etc.):
 - ▶ **Convolutional Neural Networks** – image data
 - ▶ Recurrent Neural Networks, LSTM – sequence data
 - ▶ Transformer, BERT – machine translation etc.
 - ▶ Generative models and GANs – new unsupervised learning for image, etc.
 - ▶ Reinforcement Learning – Markov decision process, playing games, etc.

Kaggle survey: Top ML Methods

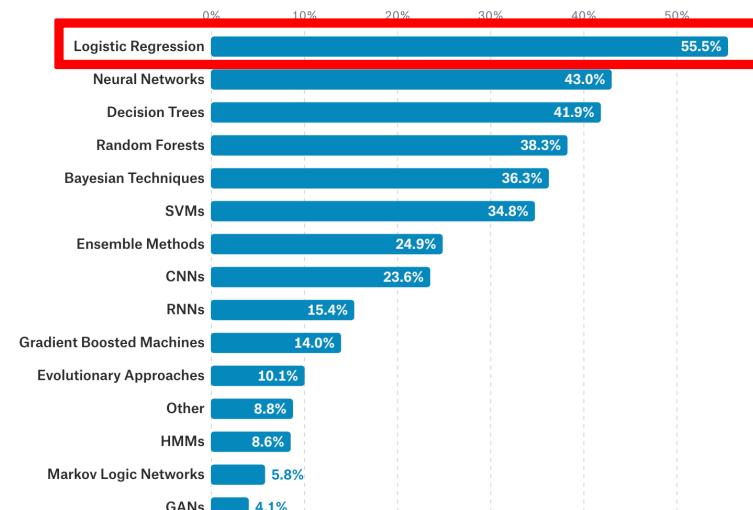
<https://www.kaggle.com/surveys/2017>

Academic

What data science methods are used at work?

Logistic regression is the most commonly reported data science method used at work for all industries except [Military and Security](#) where Neural Networks are used slightly more frequently.

Company Size ▾ Academic ▾ Job Title ▾



1,201 responses

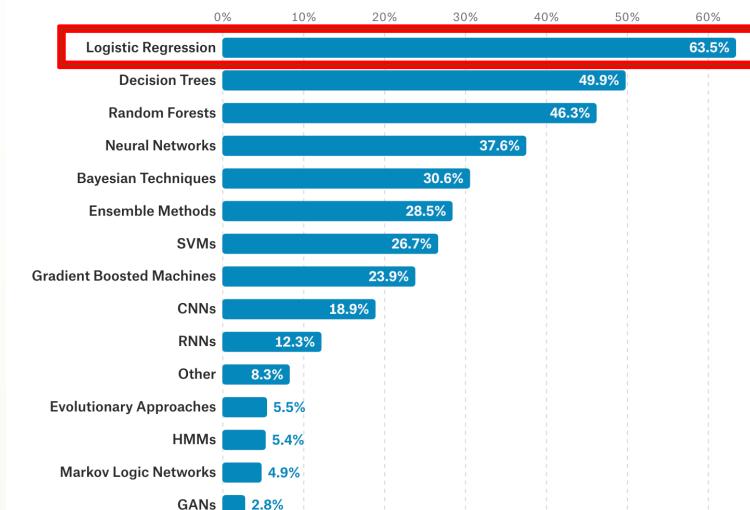
View code in Kaggle Kernels

Industry

What data science methods are used at work?

Logistic regression is the most commonly reported data science method used at work for all industries except [Military and Security](#) where Neural Networks are used slightly more frequently.

Company Size ▾ Industry ▾ Job Title ▾



7,301 responses

View code in Kaggle Kernels

What type of data is used at work?

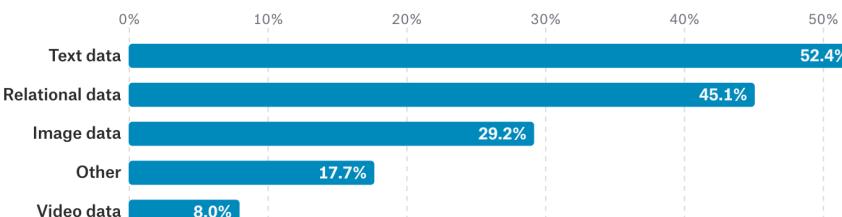
<https://www.kaggle.com/surveys/2017>

Academic

What type of data is used at work?

Relational data is the most commonly reported type of data used at work for all industries except for **Academia** and the **Military and Security** industry where text data's used more.

Company Size ▾ Academic ▾ Job Title ▾



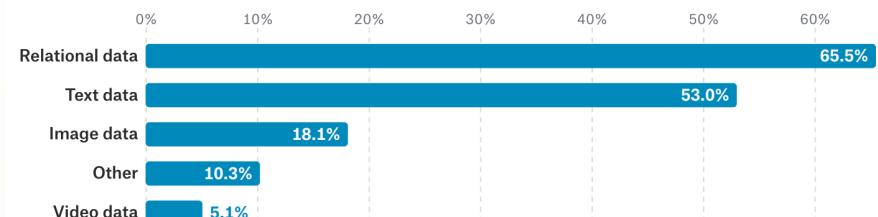
1,277 responses

Industry

What type of data is used at work?

Relational data is the most commonly reported type of data used at work for all industries except for **Academia** and the **Military and Security** industry where text data's used more.

Company Size ▾ Industry ▾ Job Title ▾



8,024 responses

Some reference books on Deep Learning

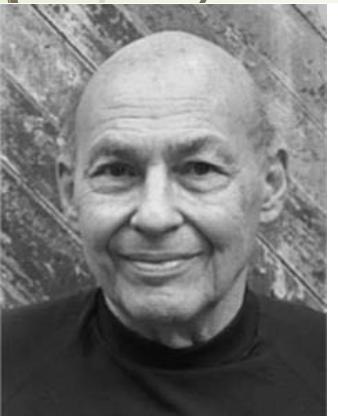
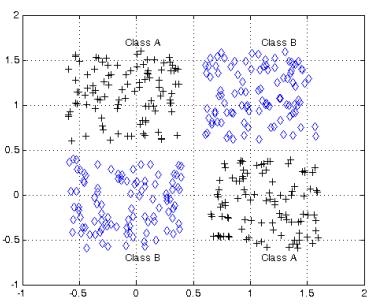
- ▶ Deep Learning with Python, Manning Publications 2017
 - ▶ by François Chollet
 - ▶ https://www.manning.com/books/deep-learning-with-python?a_aid=keras&a_bid=76564dff
- ▶ Deep Learning, MIT Press 2016
 - ▶ By Ian Goodfellow, Yoshua Bengio, and Aaron Courville,
 - ▶ <http://www.deeplearningbook.org/>
- ▶ Many other public resources

Locality or Sparsity of Computation

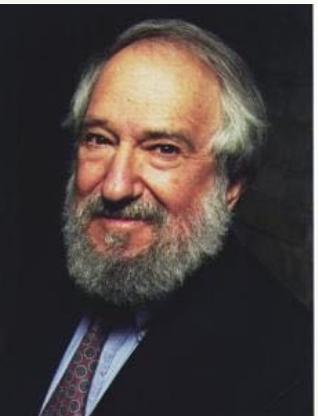
Minsky and Papert, 1969

Perceptron can't do **XOR** classification

Perceptron needs infinite global
information to compute **connectivity**



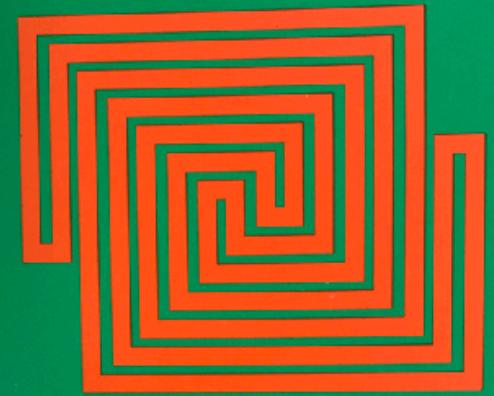
Marvin Minsky



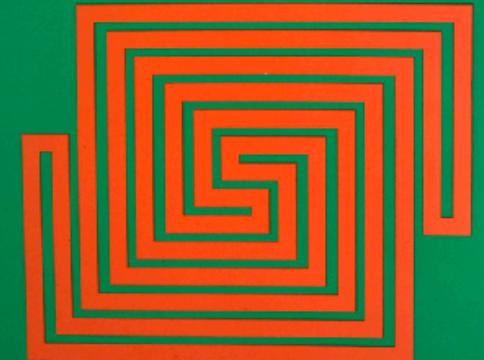
Seymour Papert

Locality or **Sparsity** is important:
Locality in time?
Locality in space?

Expanded Edition



Perceptrons



Marvin L. Minsky
Seymour A. Papert

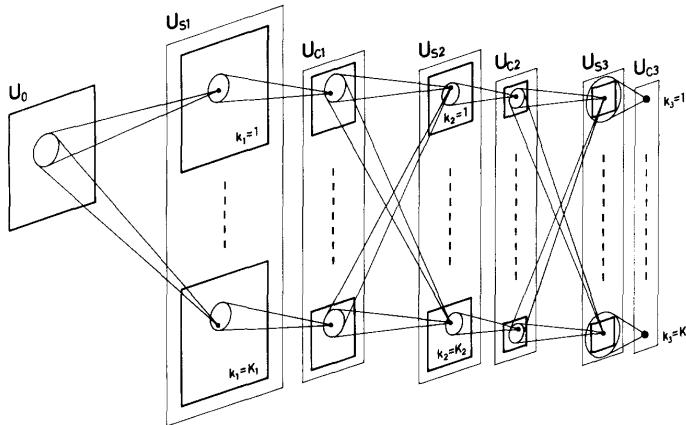
Convolutional Neural Networks: shift invariances and locality for images

Biol. Cybernetics 36, 193–202 (1980)

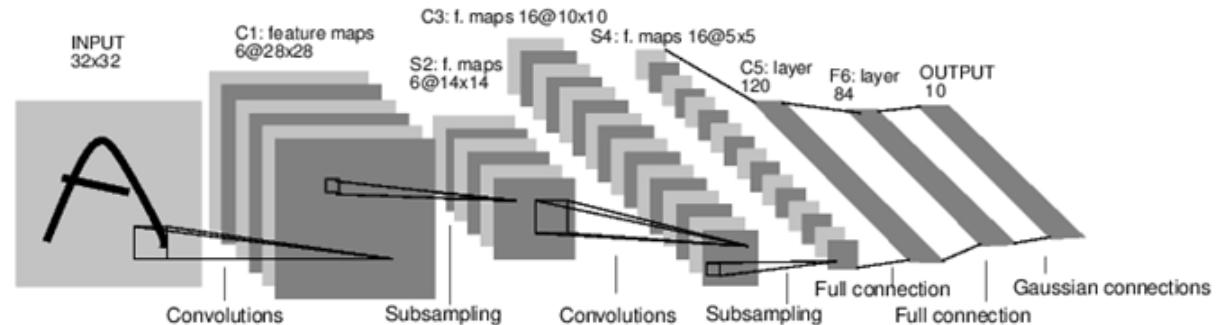
Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position

Kunihiko Fukushima

NHK Broadcasting Science Research Laboratories, Kinuta, Setagaya, Tokyo, Japan



- Can be traced to *Neocognitron* of Kunihiko Fukushima (1979)
- Yann LeCun combined convolutional neural networks with back propagation (1989)
- Imposes **shift invariance** and **locality** on the weights
- Forward pass remains similar
- Backpropagation slightly changes – need to sum over the gradients from all spatial positions



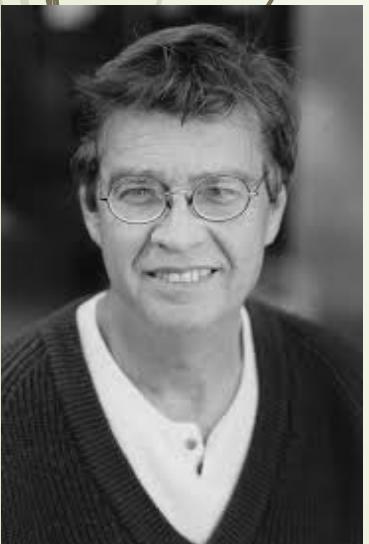
Multilayer Perceptrons (MLP) and Back-Propagation (BP) Algorithms

Rumelhart, Hinton, Williams (1986)

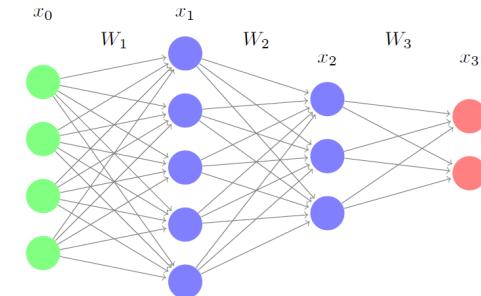
Learning representations by back-propagating errors, Nature, 323(9): 533-536

BP algorithms as **stochastic gradient descent** algorithms (**Robbins–Monro 1950; Kiefer–Wolfowitz 1951**) with Chain rules of Gradient maps

MLP classifies **XOR**, but the global hurdle on topology (connectivity) computation still exists



NATURE VOL. 323 9 OCTOBER 1986 LETTERS TO NATURE 533



Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton† & Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California, San Diego, La Jolla, California 92093, USA
† Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neuron-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output code to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-error-correcting rule.

There have been many attempts to design self-organizing neural networks. The aim is to find a powerful synaptic modification rule that will allow an arbitrarily connected neural network to develop an internal structure that is appropriate for a particular task domain. The task is specified by giving the desired state vector of the output units for each state vector of the input units. If the input units are directly connected to the output units it is relatively easy to find learning rules that iteratively adjust the relative strengths of the connections so as to progressively reduce the difference between the actual and desired output vectors*. Learning becomes more interesting but

more difficult when we introduce hidden units whose actual or desired states are not specified by the task. (In perceptrons, there are 'feature analysers' between the input and output that are not true hidden units because their input connections are fixed by hand, so their states are completely determined by the input vector: they do not learn representations.) The learning procedure must decide under what circumstances the hidden units should be active in order to help achieve the desired input-output behaviour. This amounts to deciding what these units should represent. We demonstrate that a general purpose and relatively simple procedure is powerful enough to construct appropriate learned representations.

The simplest form of the learning procedure is for layered networks which have a layer of input units at the bottom, any number of intermediate layers, and a layer of output units at the top. Connections within a layer or from higher to lower layers are forbidden, but connections can skip intermediate layers. An input vector is presented to the network by setting the states of the input units. Then the states of the units in each layer are determined by applying equations (1) and (2) to the connections coming from lower layers. All units within a layer have their states set in parallel, but different layers have their states set sequentially, starting at the bottom and working upwards until the states of the output units are determined.

The total input, x_j , to unit j is a linear function of the outputs, y_i , of the units that are connected to j and of the weights, w_{ij} , on these connections

$$x_j = \sum_i y_i w_{ij} \quad (1)$$

Units can be given biases by introducing an extra input to each unit which always has a value of 1. The weight on this extra input is called the bias and is equivalent to a threshold of the opposite sign. It can be treated just like the other weights.

A unit has a real-valued output, y_j , which is a non-linear function of its total input

$$y_j = \frac{1}{1 + e^{-x_j}} \quad (2)$$

* To whom correspondence should be addressed

BP Algorithm: Forward Pass

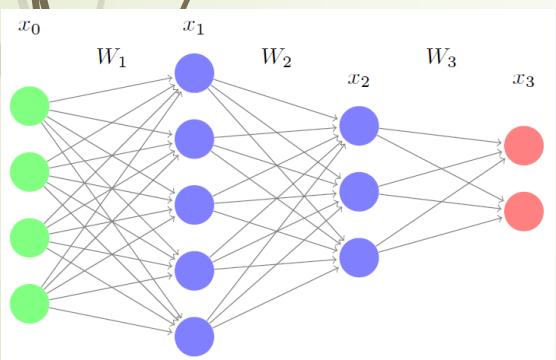
- Cascade of repeated [linear operation followed by coordinatewise nonlinearity]'s
- Nonlinearities: sigmoid, hyperbolic tangent, (recently) ReLU.

Algorithm 1 Forward pass

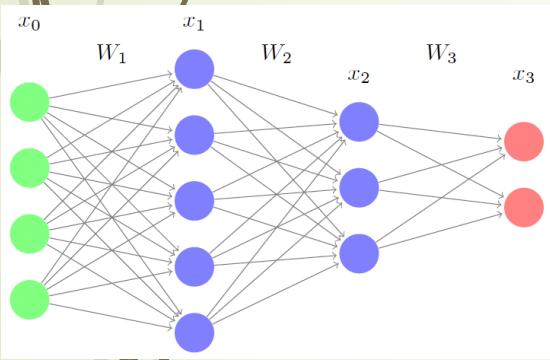
Input: x_0

Output: x_L

```
1: for  $\ell = 1$  to  $L$  do
2:    $x_\ell = f_\ell(W_\ell x_{\ell-1} + b_\ell)$ 
3: end for
```



BP algorithm = Gradient Descent Method



- Training examples $\{x_0^i\}_{i=1}^n$ and labels $\{y^i\}_{i=1}^n$
- Output of the network $\{x_L^i\}_{i=1}^m$
- Objective

$$J(\{W_l\}, \{b_l\}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|y^i - x_L^i\|_2^2 \quad (1)$$

Other losses include cross-entropy, logistic loss, exponential loss, etc.

- Gradient descent

$$W_l = W_l - \eta \frac{\partial J}{\partial W_l}$$

$$b_l = b_l - \eta \frac{\partial J}{\partial b_l}$$

In practice: use Stochastic Gradient Descent (SGD)

Derivation of BP: Lagrangian Multiplier

LeCun et al. 1988

Given n training examples $(I_i, y_i) \equiv (\text{input}, \text{target})$ and L layers

- Constrained optimization

$$\min_{W,x} \quad \sum_{i=1}^n \|x_i(L) - y_i\|_2$$

$$\text{subject to} \quad x_i(\ell) = f_\ell \left[W_\ell x_i(\ell-1) \right], \\ i = 1, \dots, n, \quad \ell = 1, \dots, L, \quad x_i(0) = I_i$$

- Lagrangian formulation (Unconstrained)

$$\min_{W,x,B} \mathcal{L}(W, x, B)$$

$$\mathcal{L}(W, x, B) = \sum_{i=1}^n \left\{ \|x_i(L) - y_i\|_2^2 + \sum_{\ell=1}^L B_i(\ell)^T \left(x_i(\ell) - f_\ell \left[W_\ell x_i(\ell-1) \right] \right) \right\}$$

back-propagation – derivation

- $\frac{\partial \mathcal{L}}{\partial B}$

Forward pass

$$x_i(\ell) = f_\ell \left[\underbrace{W_\ell x_i(\ell-1)}_{A_i(\ell)} \right] \quad \ell = 1, \dots, L, \quad i = 1, \dots, n$$

- $\frac{\partial \mathcal{L}}{\partial x}, z_\ell = [\nabla f_\ell] B(\ell)$

Backward (adjoint) pass

$$z(L) = 2\nabla f_L \left[A_i(L) \right] (y_i - x_i(L))$$

$$z_i(\ell) = \nabla f_\ell \left[A_i(\ell) \right] W_{\ell+1}^T z_i(\ell+1) \quad \ell = 0, \dots, L-1$$

- $W \leftarrow W + \lambda \frac{\partial \mathcal{L}}{\partial W}$

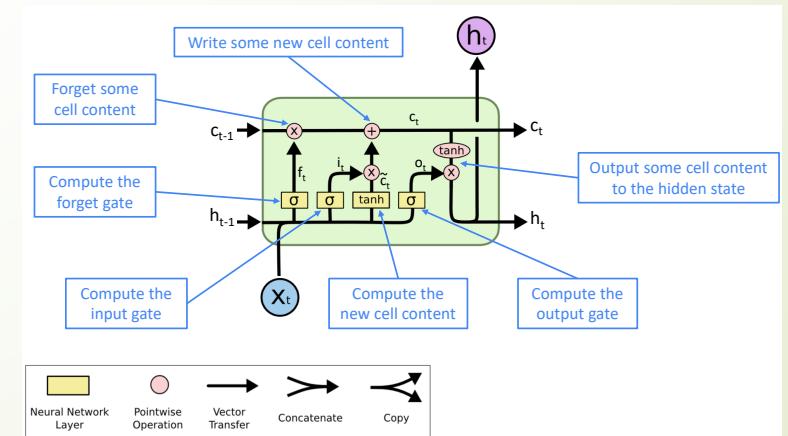
Weight update

$$W_\ell \leftarrow W_\ell + \lambda \sum_{i=1}^n z_i(\ell) x_i^T(\ell-1)$$

Long-Short-Term-Memory (LSTM, 1997)

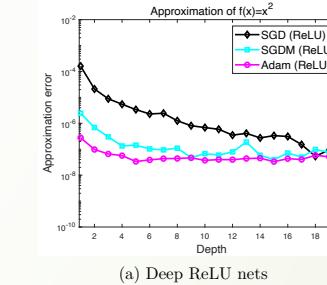
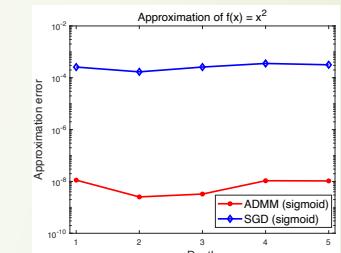
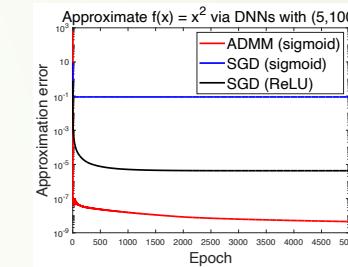
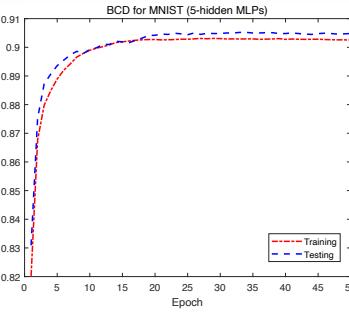
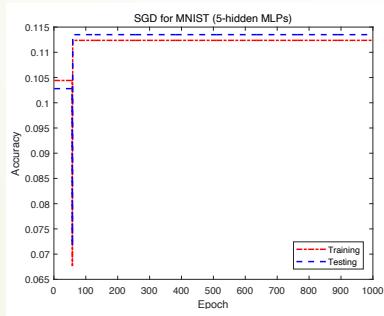
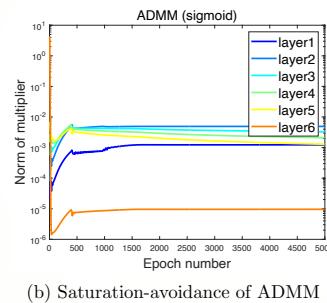
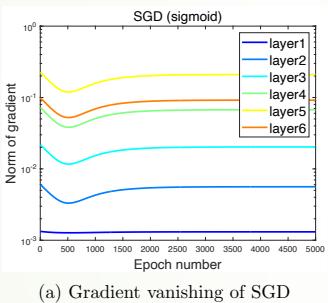


- ▶ Sepp Hochreiter; Jürgen Schmidhuber (1997). "Long short-term memory". *Neural Computation*. 9 (8): 1735–1780. (<https://www.bioinf.jku.at/publications/older/2604.pdf>)
- ▶ BP can not train deep networks due to gradient vanishing problem etc.
- ▶ Introduction of **short path** to train deep networks without vanishing gradient problem.
- ▶ This idea will come back to Convolutional Networks as **ResNet** in 2015.



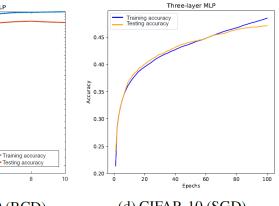
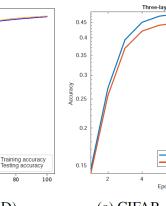
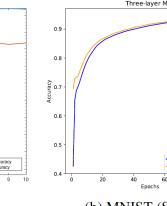
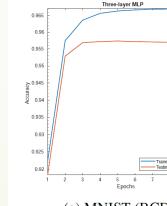
SGD vs. ADMM/BCD

- ▶ Stochastic Gradient Descent (SGD) suffers from the well-known *gradient vanishing issue* in deep learning
- ▶ ADMM/BCD may alleviate gradient vanishing



(b) Deep sigmoid nets

- High epoch efficiency of BCD at early stage



Notes on Algorithms

- ▶ Gradient descent (back propagation) can be derived via Lagrangian multiplier method [LeCun 1988, <http://yann.lecun.com/exdb/publis/pdf/lecun-88.pdf>]
- ▶ **ADMM** (Alternating Direction Method of Multipliers) is alternative primal-dual method via **Augmented Lagrangian** multipliers [Zeng-Lin-Y.-Zhou, JMLR 2021]
- ▶ **BCD** (Block-Coordinate-Descent) **drops** the dual update in Augmented Lagrangian multipliers [Zeng-Lau-Lin-Y., ICML 2019]
- ▶ Global convergence to KKT points from arbitrary initialization can be established with the aid of **Kurdyka-Łojasiewicz** framework.

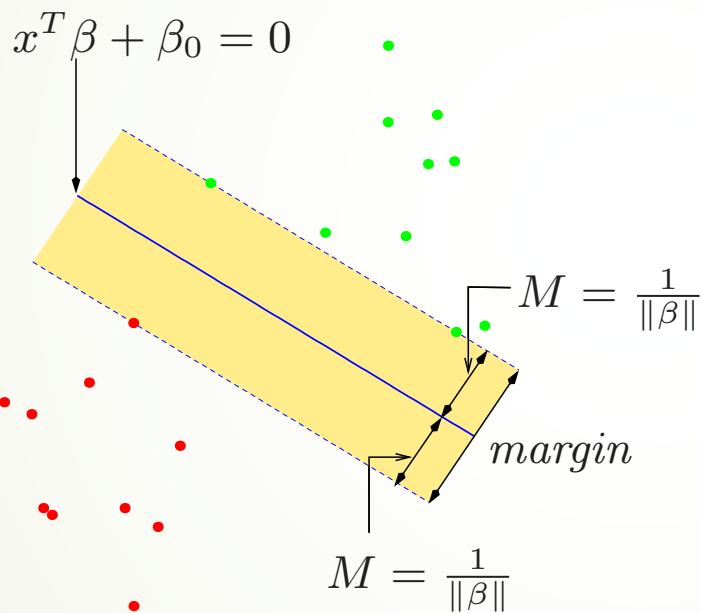
$$\begin{aligned} & \underset{\mathcal{W}, \mathcal{V}}{\text{minimize}} \quad \frac{1}{2} \|V_N - Y\|_F^2 + \frac{\lambda}{2} \sum_{i=1}^N \|W_i\|_F^2 \\ & \text{subject to} \quad V_i = \sigma(W_i V_{i-1}), \quad i = 1, \dots, N-1, \quad V_N = W_N V_{N-1}, \end{aligned}$$

Augmented Lagrangian function:

Lagrangian multiplier Λ_i

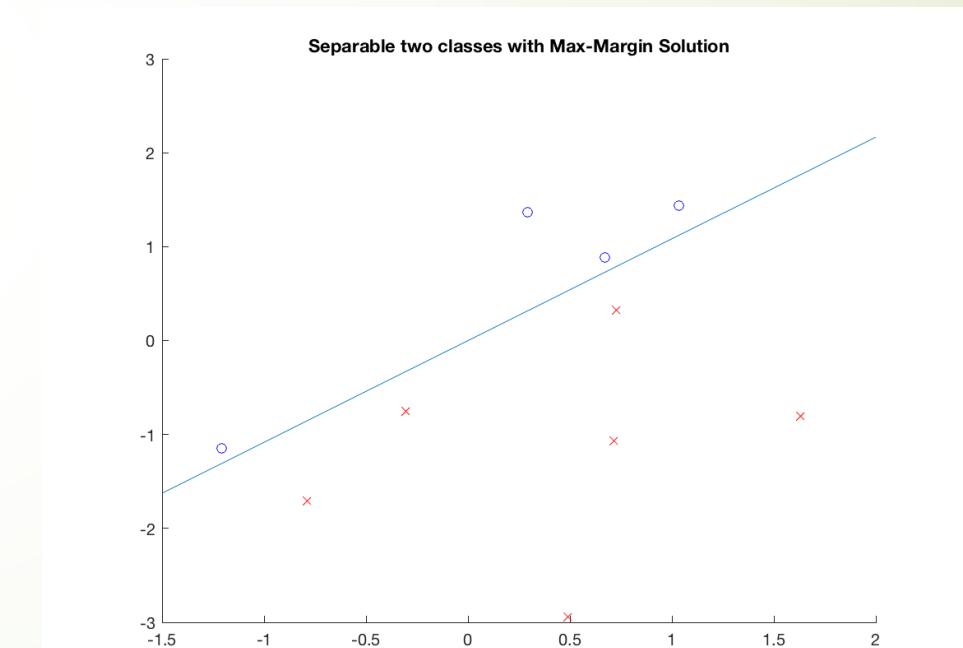
$$\begin{aligned} \mathcal{L}(\mathcal{W}, \mathcal{V}, \{\Lambda_i\}_{i=1}^N) := & \frac{1}{2} \|V_N - Y\|_F^2 + \frac{\lambda}{2} \sum_{i=1}^N \|W_i\|_F^2 \\ & + \sum_{i=1}^{N-1} \left(\frac{\beta_i}{2} \|\sigma(W_i V_{i-1}) - V_i\|_F^2 + \langle \Lambda_i, \sigma(W_i V_{i-1}) - V_i \rangle \right) \\ & + \frac{\beta_N}{2} \|W_N V_{N-1} - V_N\|_F^2 + \langle \Lambda_N, W_N V_{N-1} - V_N \rangle, \end{aligned}$$

Support Vector Machine (Max-Margin Classifier)



Vladimir Vapnik, 1994

$$\begin{aligned} & \text{minimize}_{\beta_0, \beta_1, \dots, \beta_p} \|\beta\|^2 := \sum_j \beta_j^2 \\ & \text{subject to } y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq 1 \text{ for all } i \end{aligned}$$



Convex optimization + Reproducing Kernel Hilbert Spaces (Grace Wahba etc.)

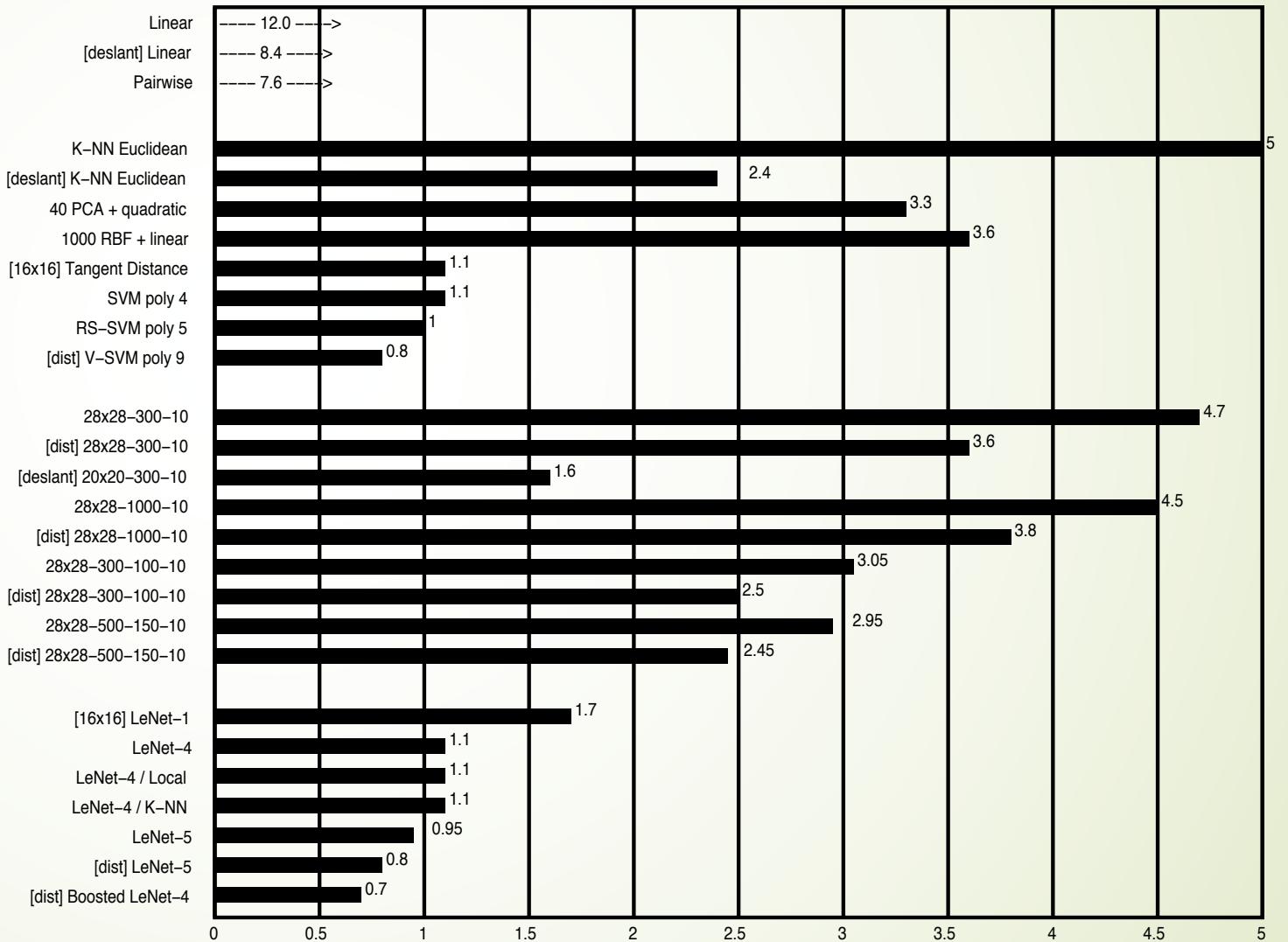
MNIST Challenge Test Error: SVM vs. CNN

LeCun et al. 1998



Simple SVM performs as well as Multilayer Convolutional Neural Networks which need careful tuning (LeNets)

Second dark era for NN:
2000s



LeNet

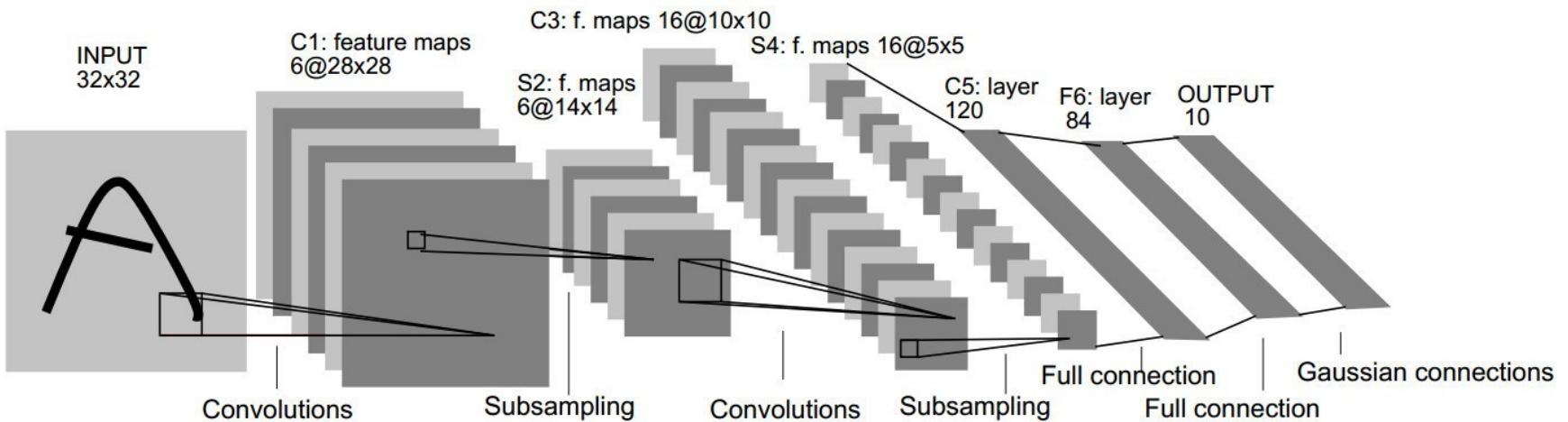


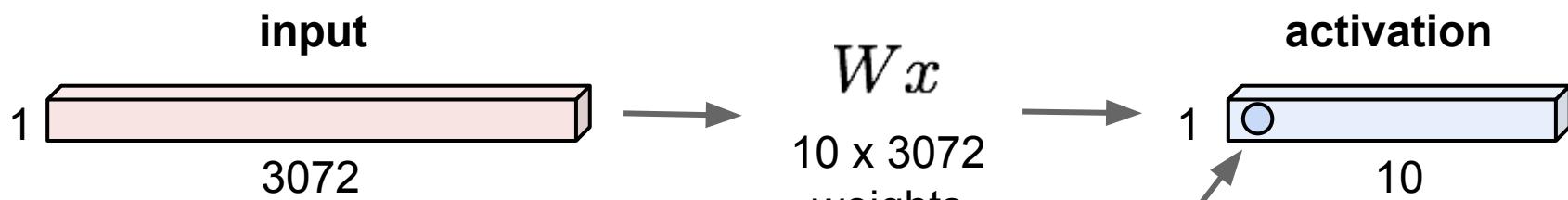
Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

<http://blog.csdn.net/Chenyukuai6625>

- **Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner.** Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, november 1998.

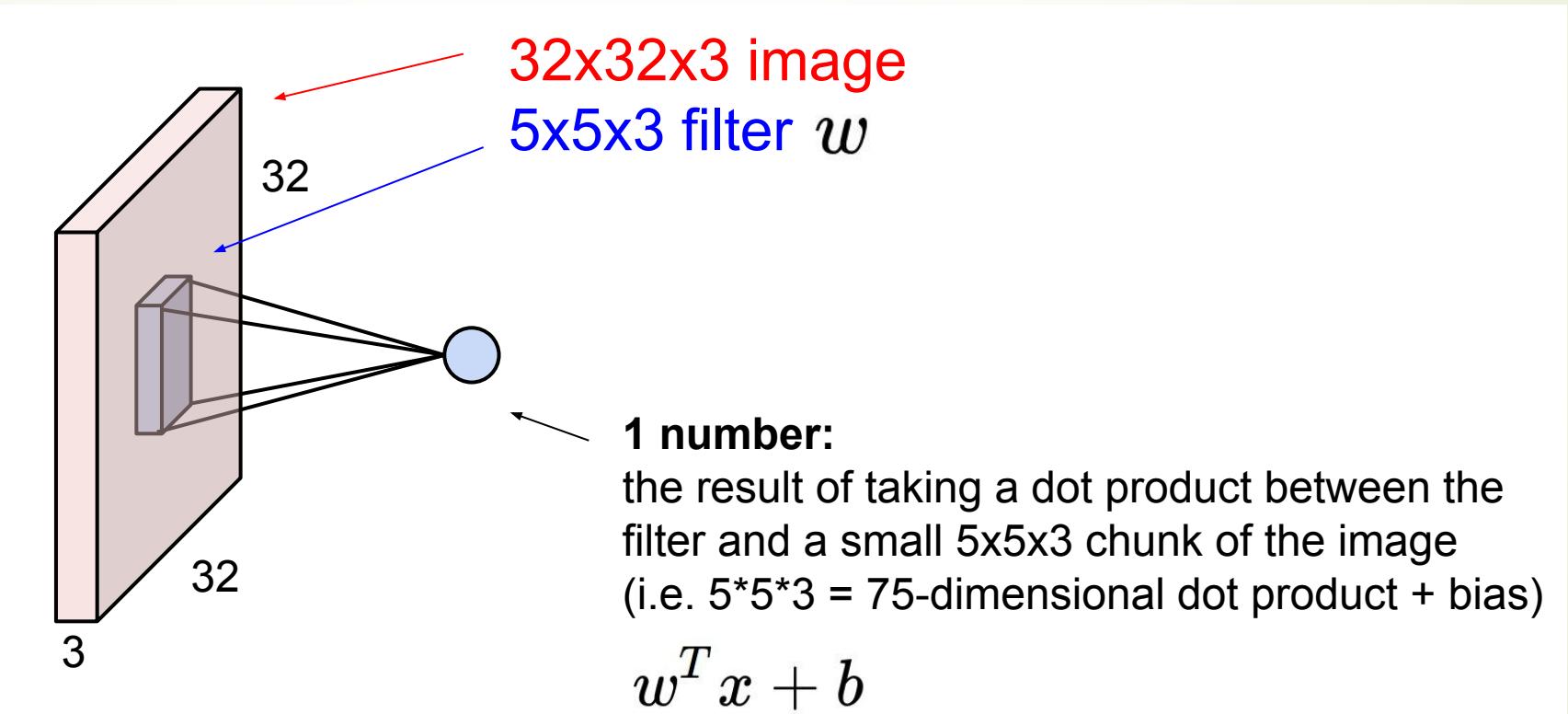
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

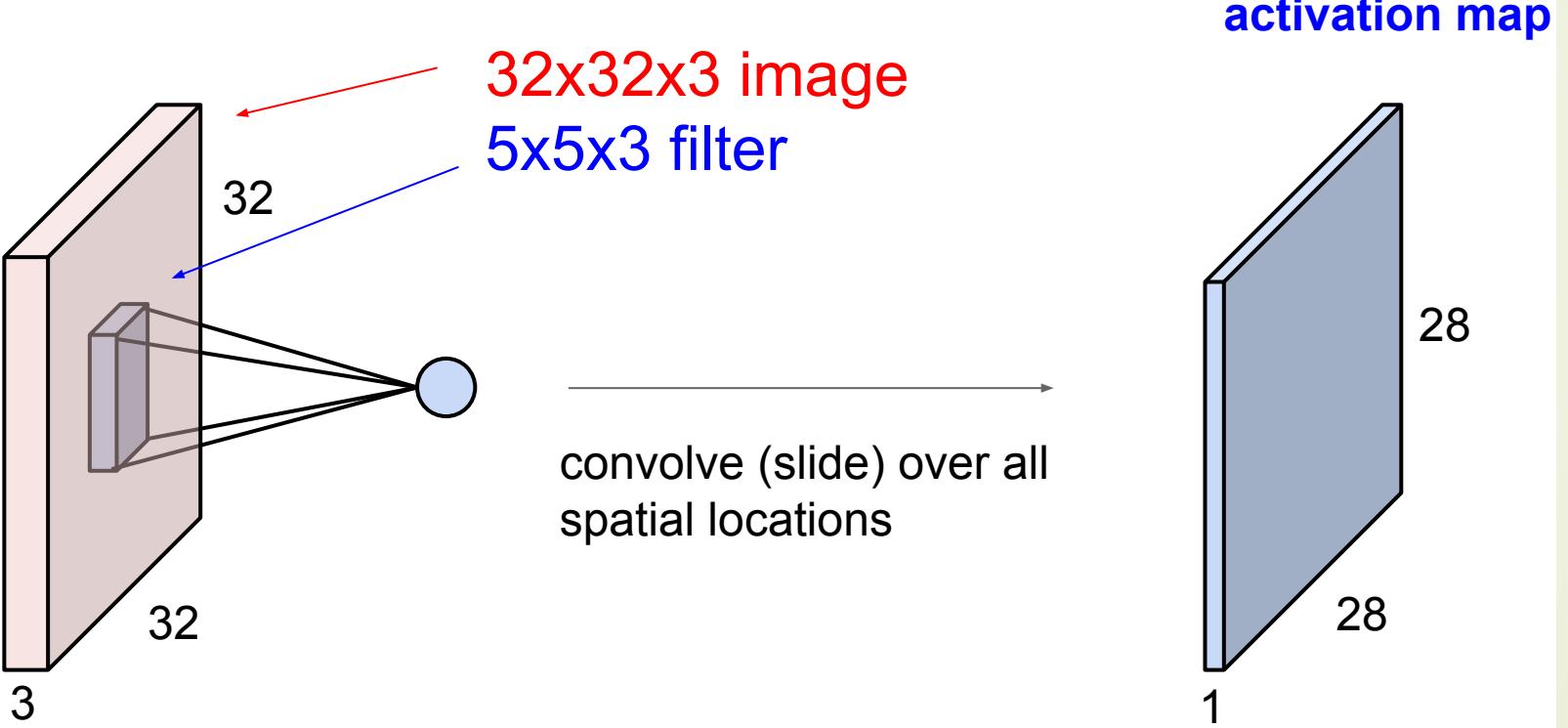


1 number:
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

Convolution



Convolution Layer: a first (blue) filter

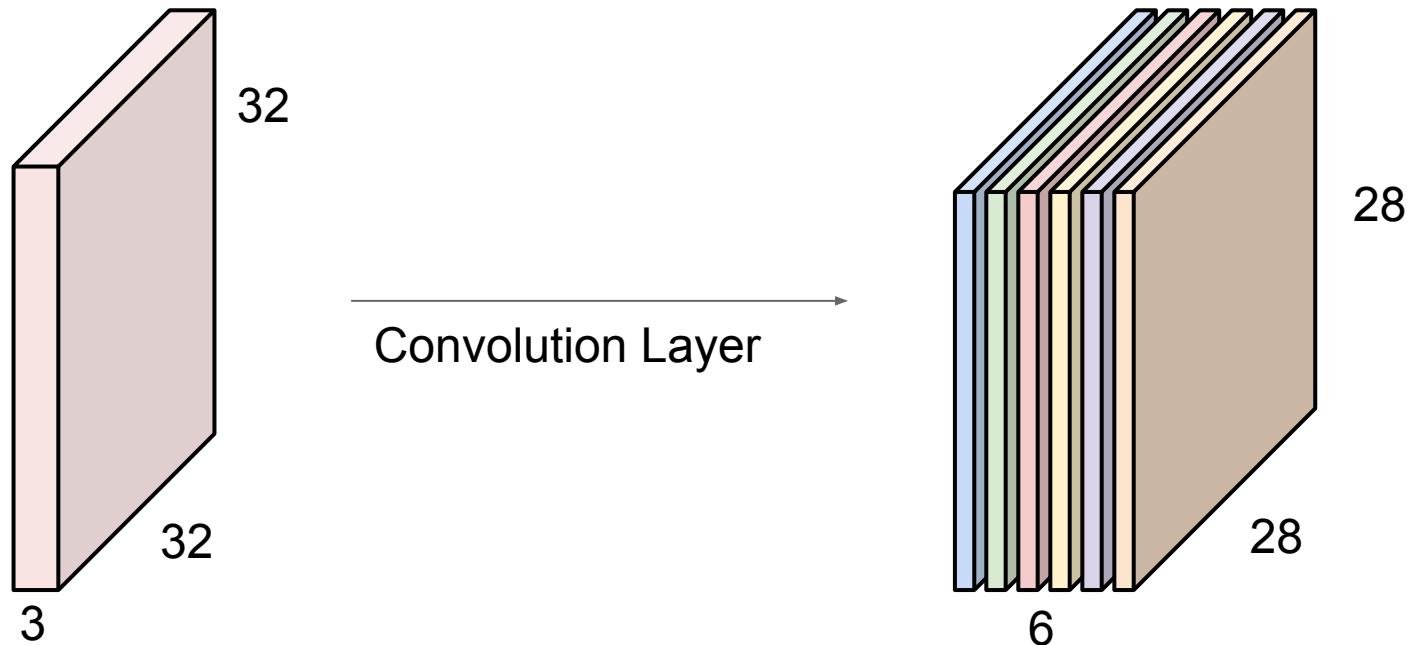


Convolution Layer: a second (green) filter



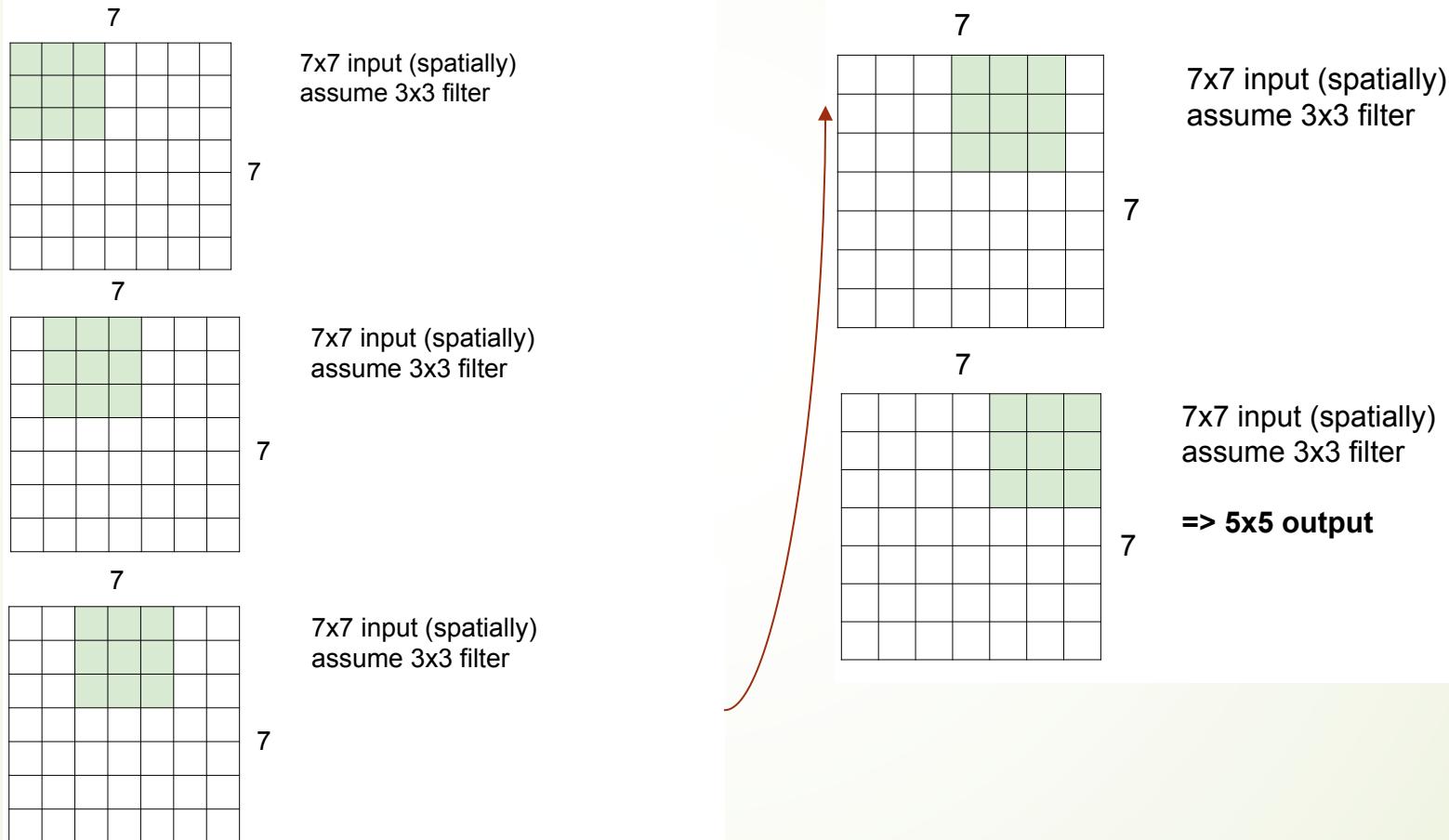
Convolution Layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

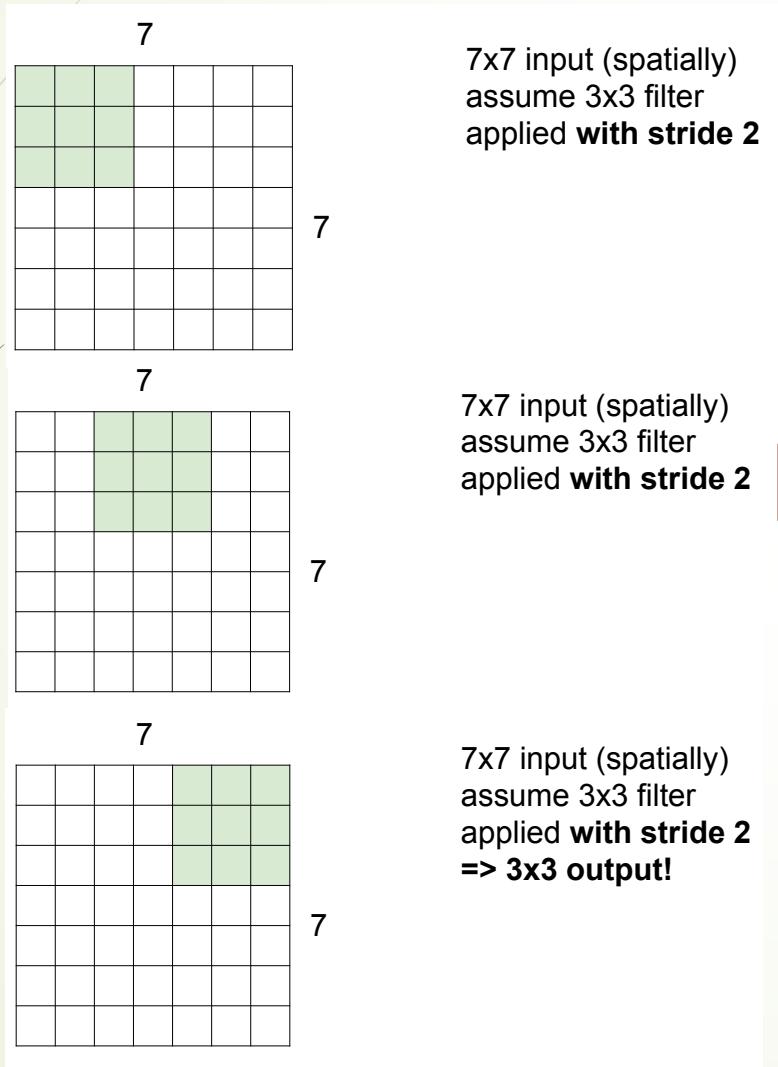


We stack these up to get a “new image” of size $28 \times 28 \times 6$!

A Closer Look at Convolution: stride=1



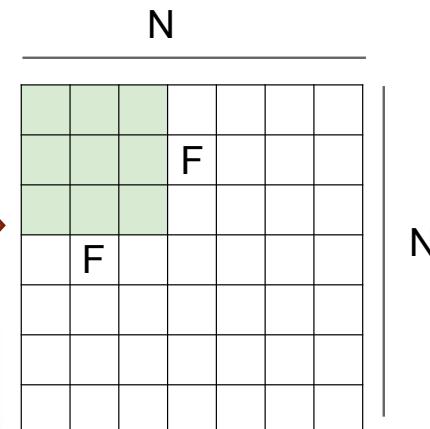
A Closer Look at Convolution: stride=2



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!



Output size:
(N - F) / stride + 1

e.g. N = 7, F = 3:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 \rightarrow$

A Closer Look at Convolution: Padding

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

7x7 output!

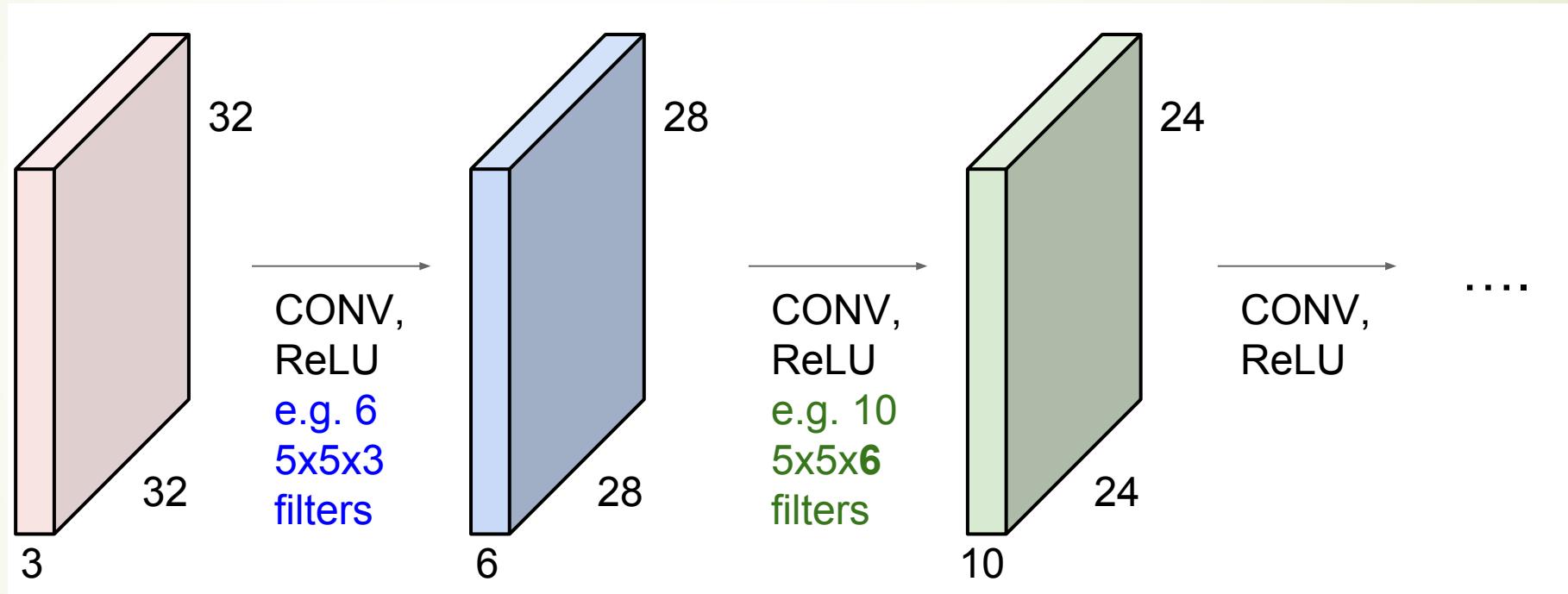
in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$. (will preserve size spatially)

e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

ConvNet:



Stride = 1
Padding = 0

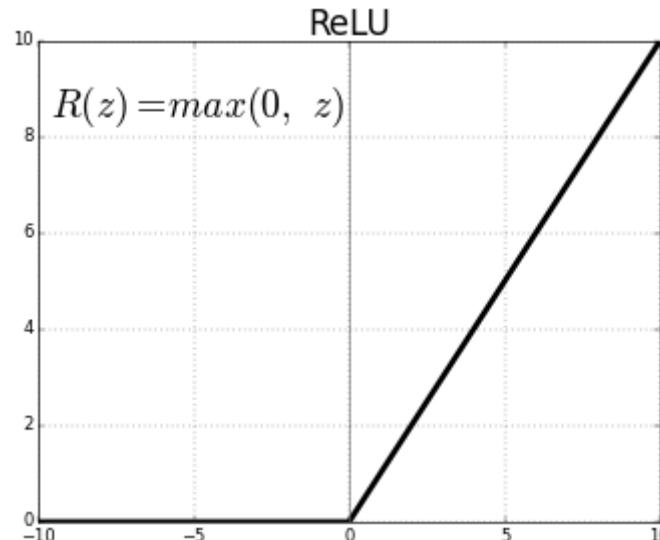
Formula: `NewImageSize = floor((ImageSize - Filter + 2*Padding)/Stride + 1)`

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

ReLU

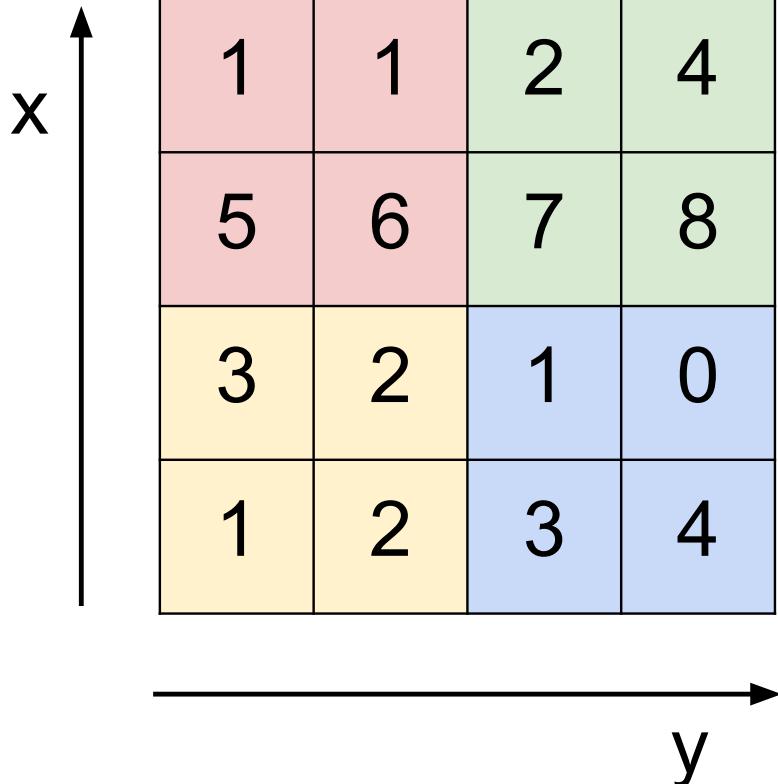
- Non-saturating function and therefore faster convergence when compared to other nonlinearities
- Problem of dying neurons



Source: https://ml4a.github.io/ml4a/neural_networks/

Max Pooling

Single depth slice



max pool with 2x2 filters
and stride 2

The result of the max pooling operation is a 2x2 output tensor:

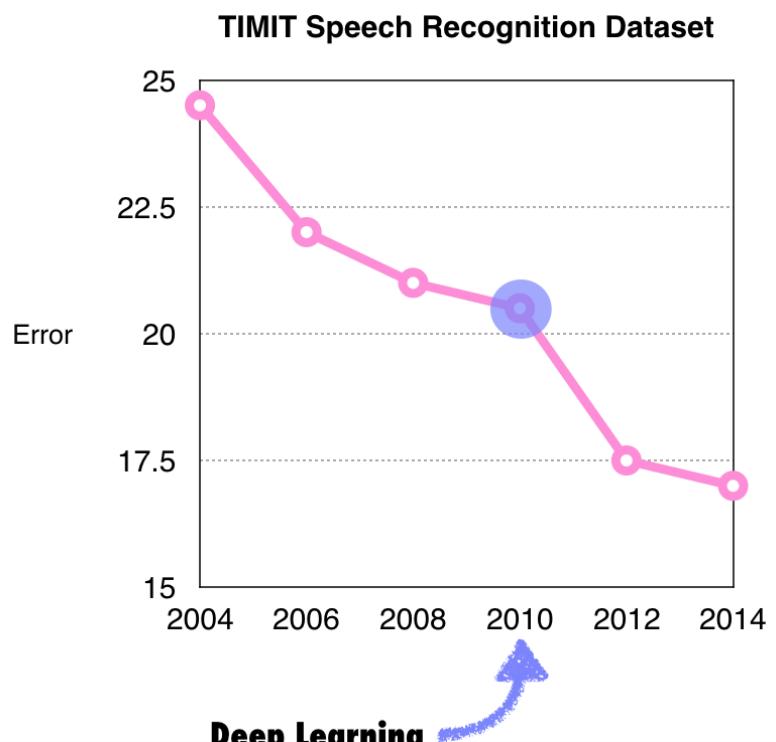
6	8
3	4

2000-2010: The Era of SVM, Boosting, ... as nights of Neural Networks



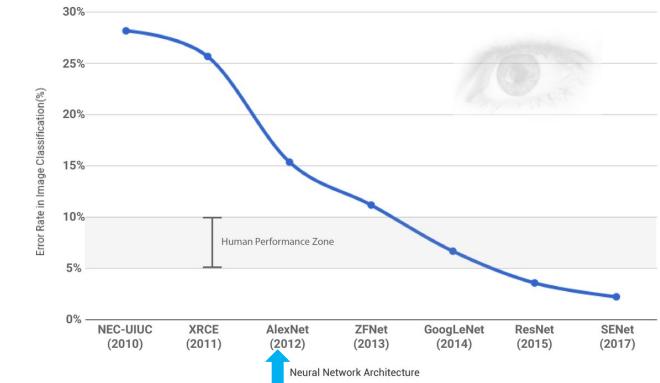
Around the year of 2012...

Speech Recognition: TIMIT



Computer Vision: ImageNet

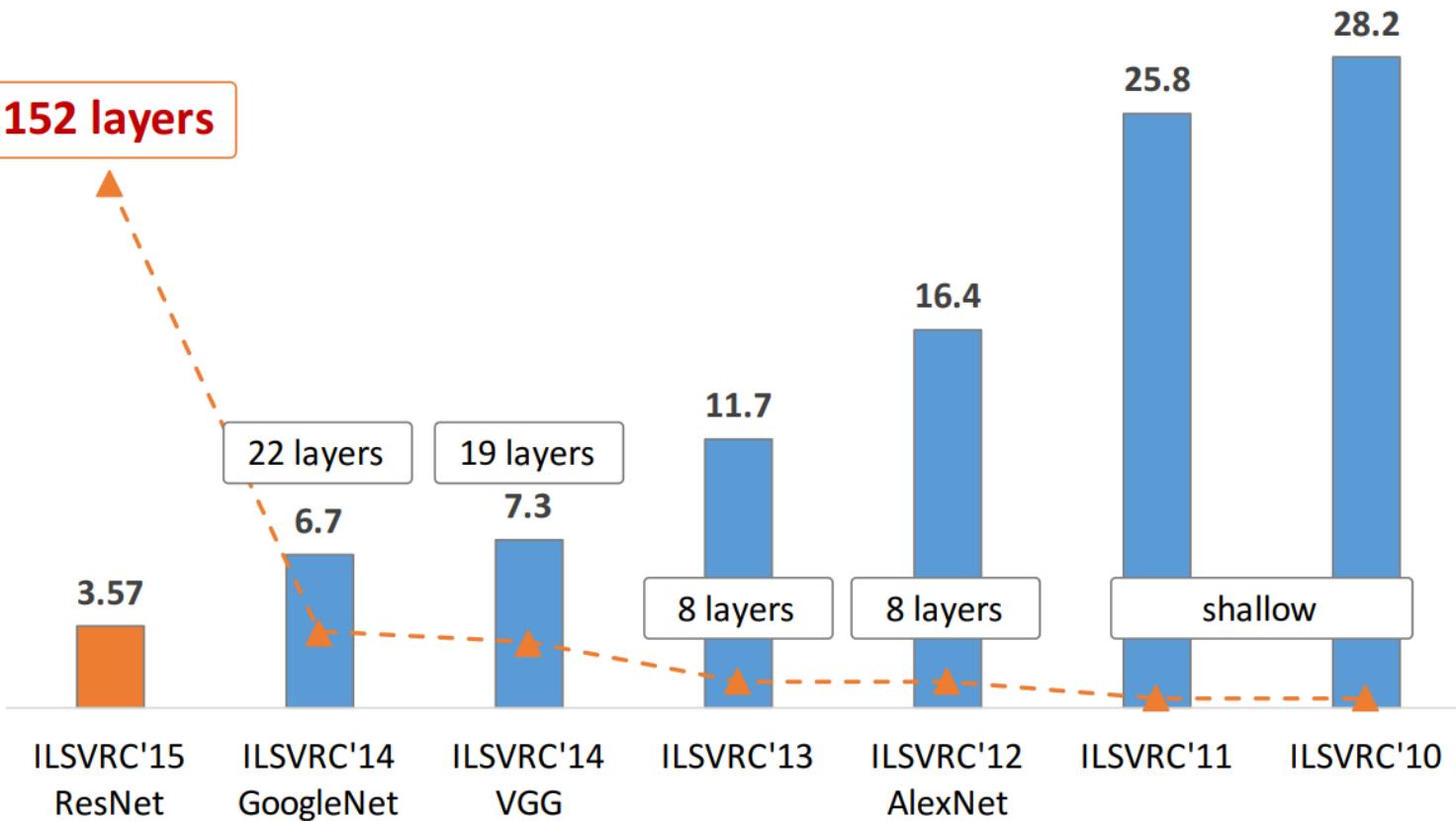
- ImageNet (subset):
 - 1.2 million training images
 - 100,000 test images
 - 1000 classes
- ImageNet large-scale visual recognition Challenge



source: <https://www.linkedin.com/pulse/must-read-path-breaking-papers-image-classification-muktabh-mayank>

Deep Learning

Depth as function of year

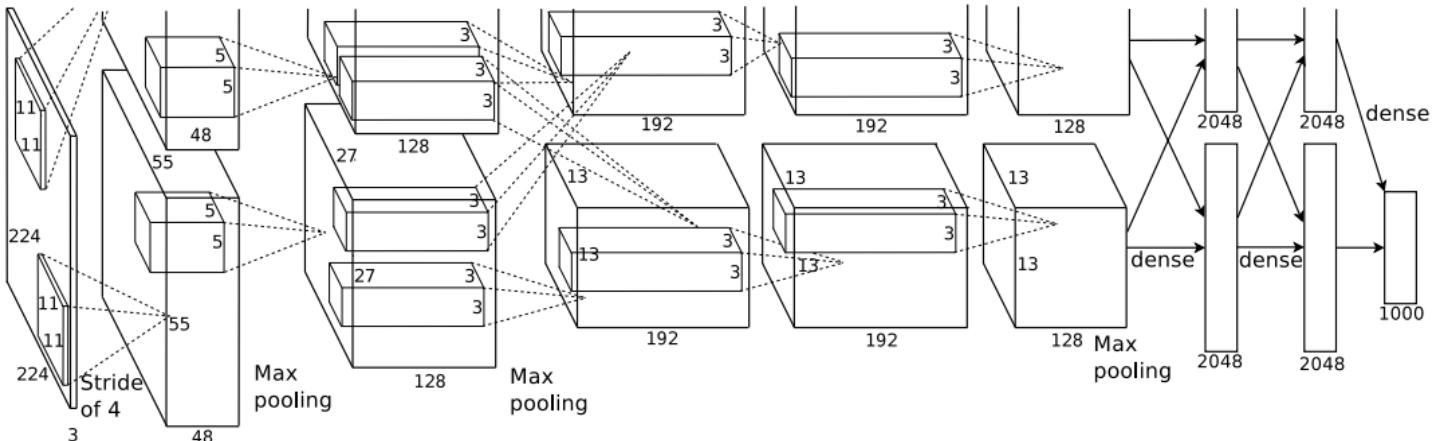


[He et al., 2016]

AlexNet (2012): Architecture



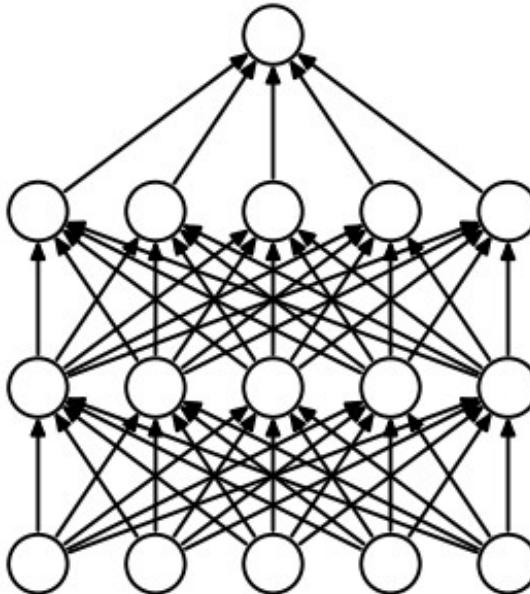
- 8 layers: first 5 convolutional, rest fully connected
- ReLU nonlinearity
- Local response normalization
- Max-pooling
- Dropout



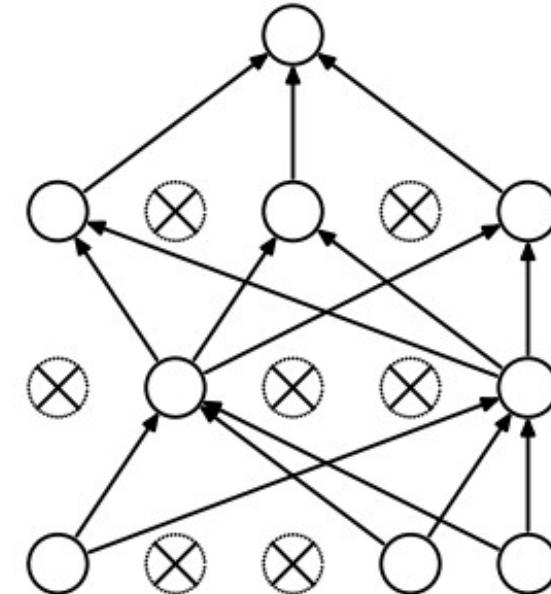
Source: [Krizhevsky et al., 2012]

<https://github.com/computerhistory/AlexNet-Source-Code>

AlexNet (2012): Dropout



(a) Standard Neural Net



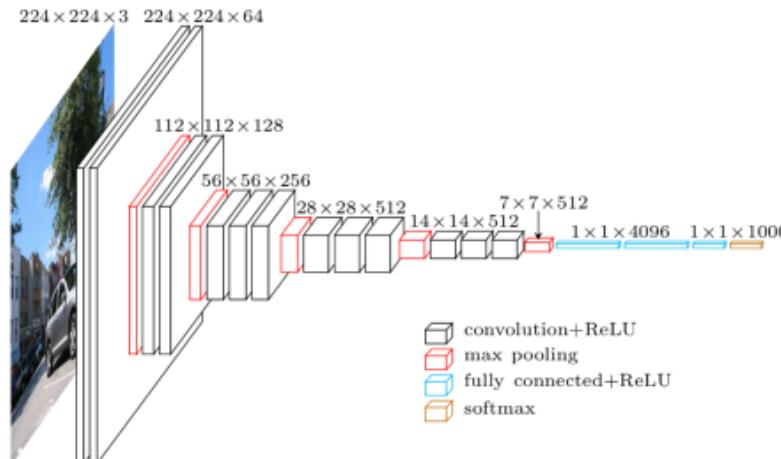
(b) After applying dropout.

Source: [Srivastava et al., 2014]

- Zero every neuron with probability $1 - p$
- At test time, multiply every neuron by p

VGG (2014) [Simonyan-Zisserman'14]

- Deeper than AlexNet: 11-19 layers versus 8
- No local response normalization
- Number of filters multiplied by two every few layers
- Spatial extent of filters 3×3 in all layers
- Instead of 7×7 filters, use three layers of 3×3 filters
 - Gain intermediate nonlinearity
 - Impose a regularization on the 7×7 filters

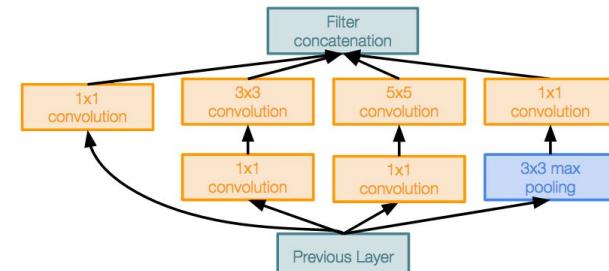


Stanford University

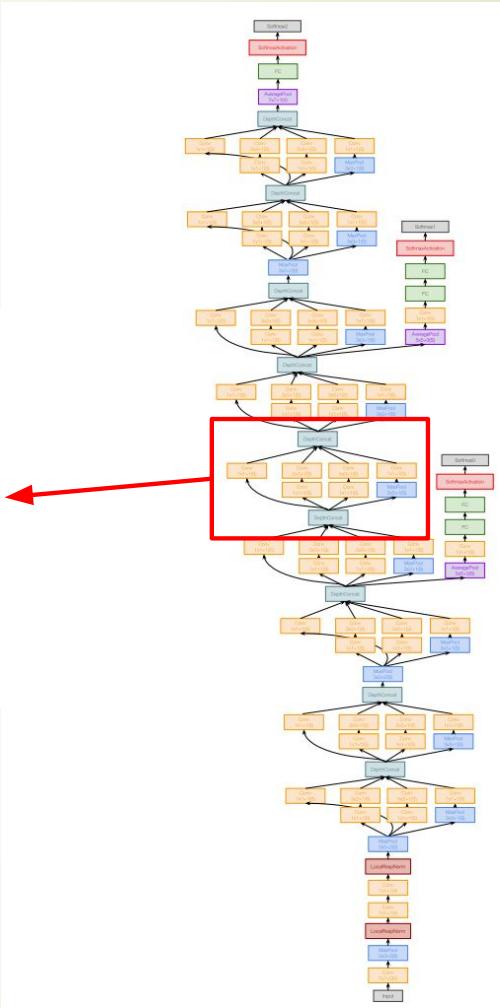
Source: <https://blog.heuritech.com/2016/02/29/>

GoogLeNet [Szegedy et al., 2014]

- ▶ 22 layers
- ▶ Efficient “Inception” module
- ▶ No FC layers
- ▶ Only 5 million parameters!
- ▶ 12x less than AlexNet
- ▶ ILSVRC’14 classification winner
(6.7% top 5 error)

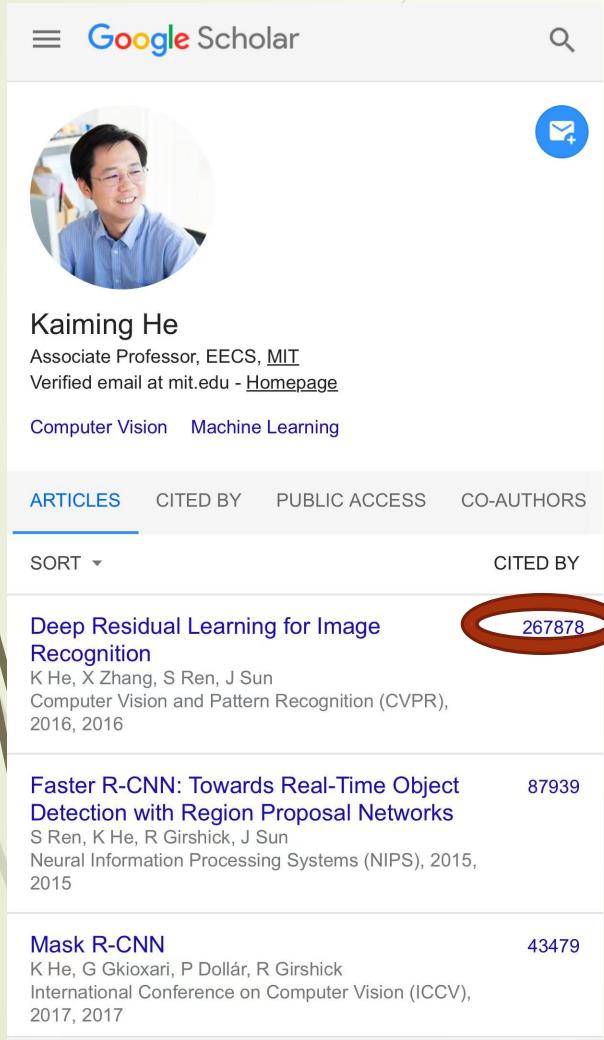


Inception module



ResNet (2015) [HGRS-15]

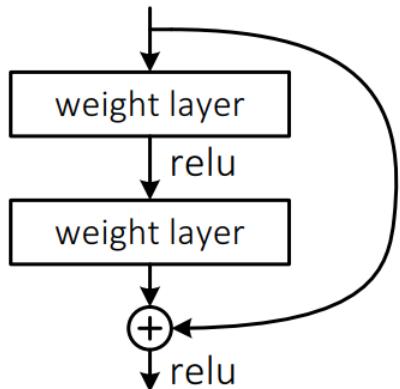
ILSVRC'15 classification winner
(3.57% top 5 error)



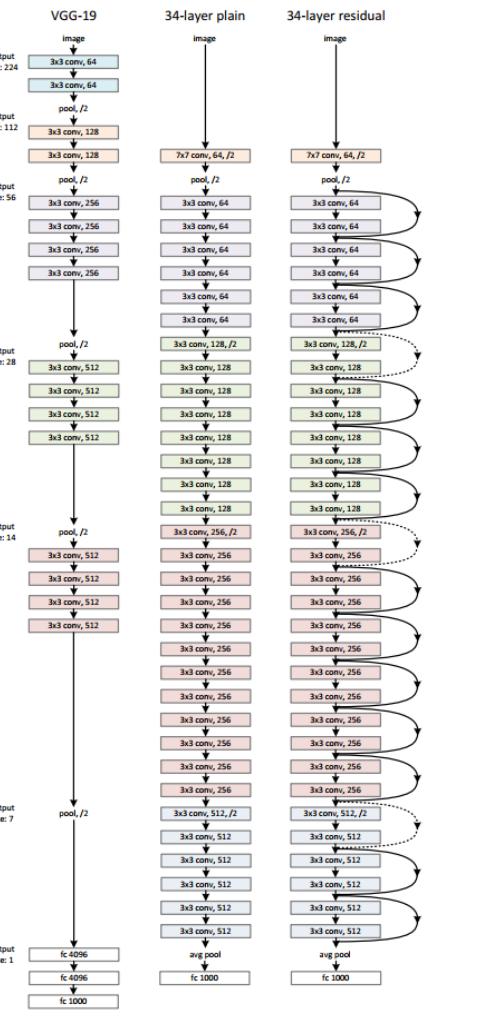
Kaiming He's Google Scholar profile shows his publications, citations, and research interests. His most cited paper, "Deep Residual Learning for Image Recognition," has 267,878 citations, which is circled in red.

ARTICLES	CITED BY	PUBLIC ACCESS	CO-AUTHORS
Deep Residual Learning for Image Recognition	267878		K He, X Zhang, S Ren, J Sun
Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks	87939		S Ren, K He, R Girshick, J Sun
Mask R-CNN	43479		K He, G Gkioxari, P Dollár, R Girshick

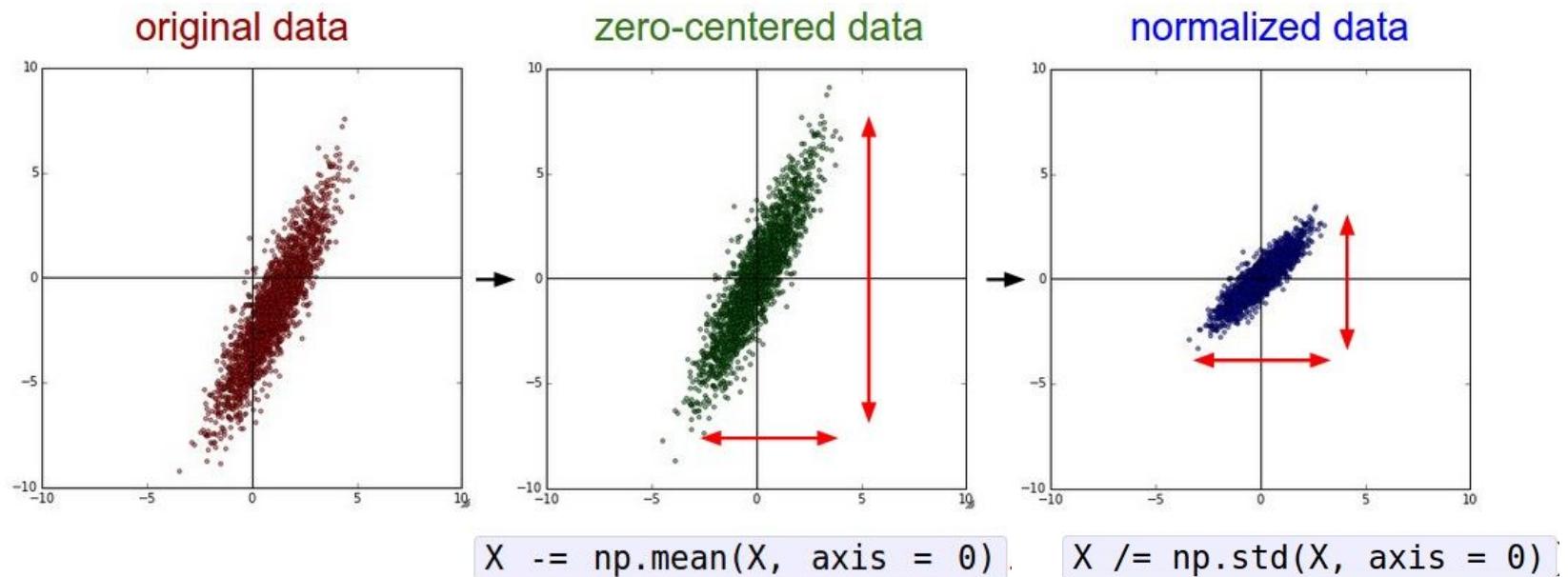
- Solves problem by adding skip connections
- Very deep: 152 layers
- No dropout
- Stride
- Batch normalization



Source: Deep Residual Learning for Image Recognition



Batch Normalization



(Assume X [NxD] is data matrix,
each example in a row)

Batch Normalization

Algorithm 2 Batch normalization [Ioffe and Szegedy, 2015]

Input: Values of x over minibatch $x_1 \dots x_B$, where x is a certain channel in a certain feature vector

Output: Normalized, scaled and shifted values $y_1 \dots y_B$

-
- 1: $\mu = \frac{1}{B} \sum_{b=1}^B x_b$
 - 2: $\sigma^2 = \frac{1}{B} \sum_{b=1}^B (x_b - \mu)^2$
 - 3: $\hat{x}_b = \frac{x_b - \mu}{\sqrt{\sigma^2 + \epsilon}}$
 - 4: $y_b = \gamma \hat{x}_b + \beta$
-

- Accelerates training and makes initialization less sensitive
- Zero mean and unit variance feature vectors

BatchNorm at Test

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

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

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

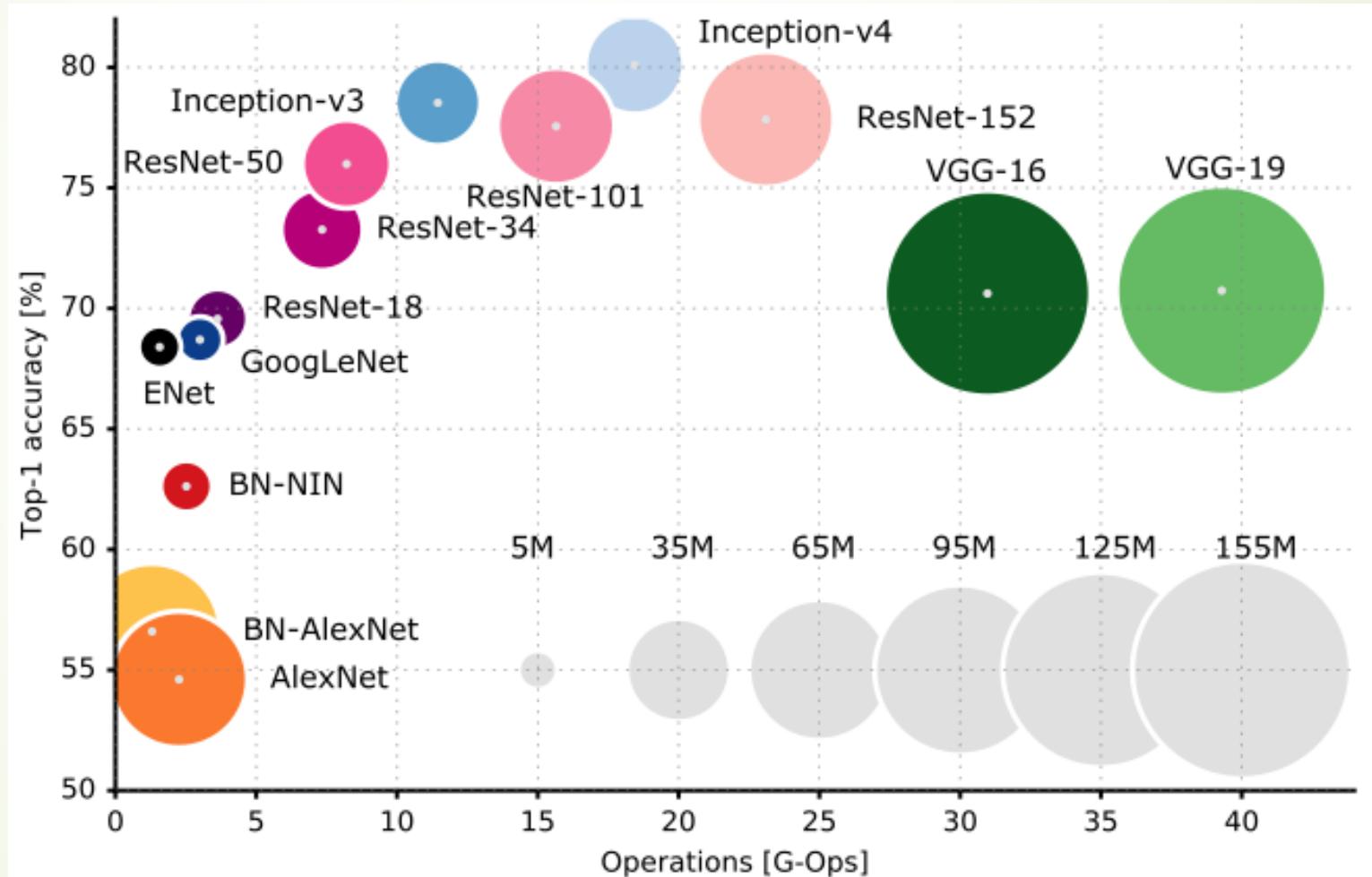
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Note: at test time BatchNorm layer functions differently:

The mean/std are not computed based on the batch. Instead, a single fixed empirical mean of activations during training is used.

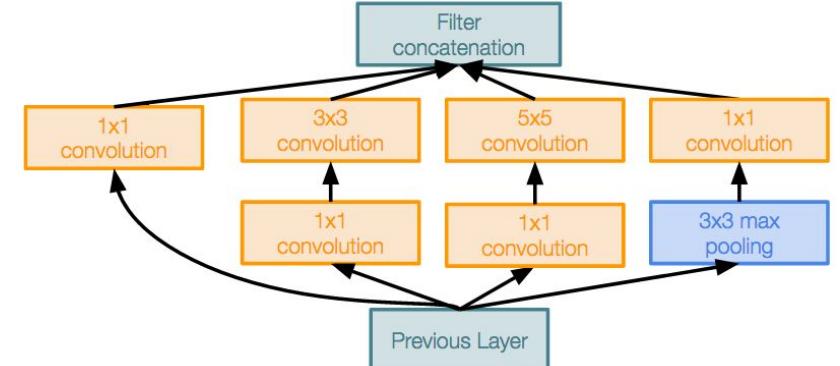
(e.g. can be estimated during training with running averages)

Complexity vs. Accuracy of Different Networks

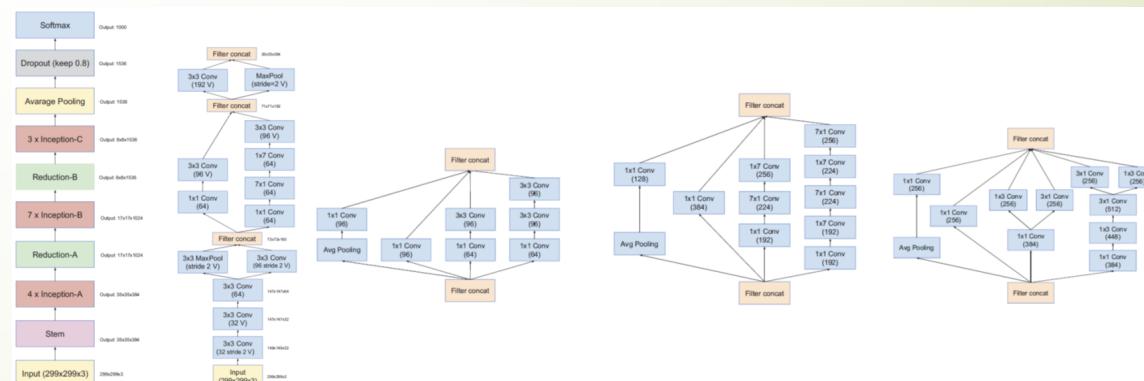


Inception-v4 = ResNet + Inception

- ▶ “Inception” module:
 - ▶ Introduced by Szegedy et al., 2014 in **GoogLeNet**
 - ▶ ILSVRC’14 classification winner (6.7% top 5 error)
- ▶ Apply parallel filter operations on the input from previous layer:
 - ▶ Dimensionality reduction (1×1 conv)
 - ▶ Multiple receptive field sizes for convolution (1×1 , 3×3 , 5×5)
 - ▶ Pooling operation (3×3)
- ▶ Concatenate all filter outputs together depth-wise



Inception module



Deep Learning Softwares

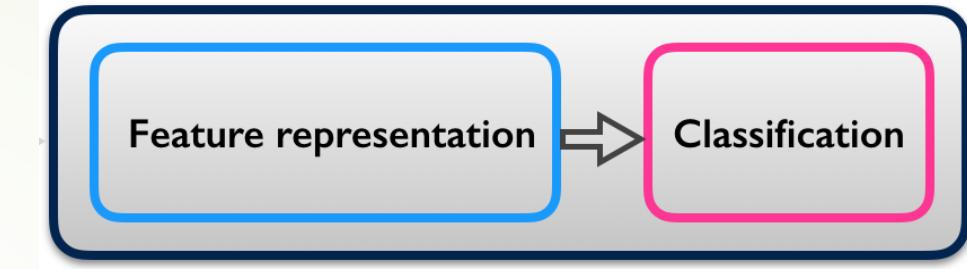
- ▶ **Pytorch** (developed by Yann LeCun and Facebook):
 - ▶ <http://pytorch.org/tutorials/>
- ▶ Tensorflow (developed by Google based on Caffe)
 - ▶ <https://www.tensorflow.org/tutorials/>
- ▶ Theano (developed by Yoshua Bengio)
 - ▶ <http://deeplearning.net/software/theano/tutorial/>
- ▶ **Keras (based on Tensorflow or Pytorch)**
 - ▶ https://www.manning.com/books/deep-learning-with-python?a_aid=keras&a_bid=76564dff

Show some examples by jupyter notebooks...



Transfer Learning: Feature Extraction and Fine Tuning

Transfer Learning?

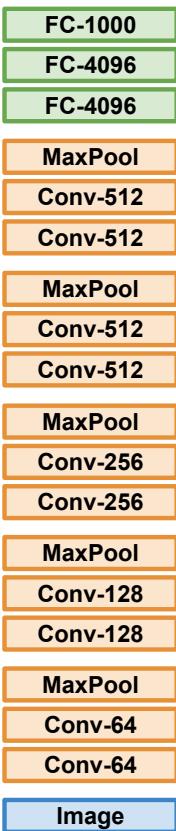


- Filters learned in first layers of a network are transferable from one task to another
- When solving another problem, no need to retrain the lower layers, just fine tune upper ones
- Is this simply due to the large amount of images in ImageNet?
- Does solving many classification problems simultaneously result in features that are more easily transferable?
- Does this imply filters can be learned in unsupervised manner?
- Can we characterize filters mathematically?

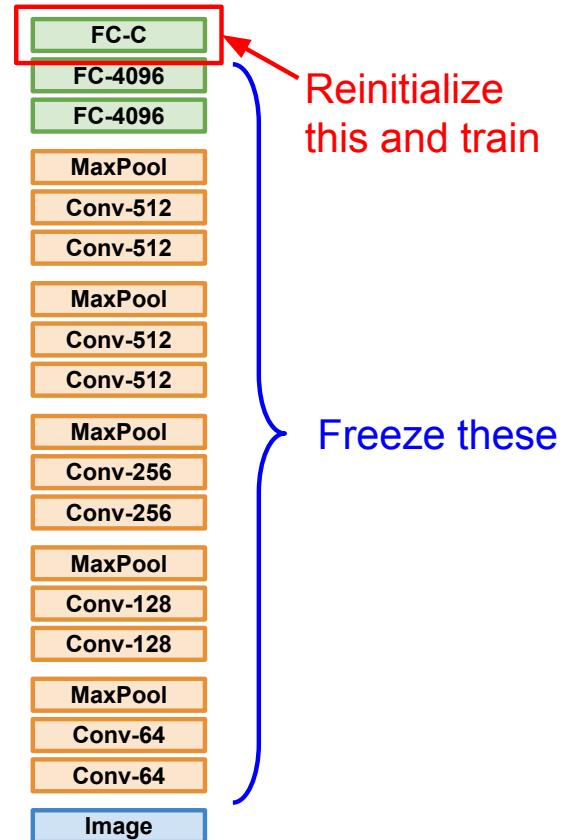
Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Transfer Learning with CNNs

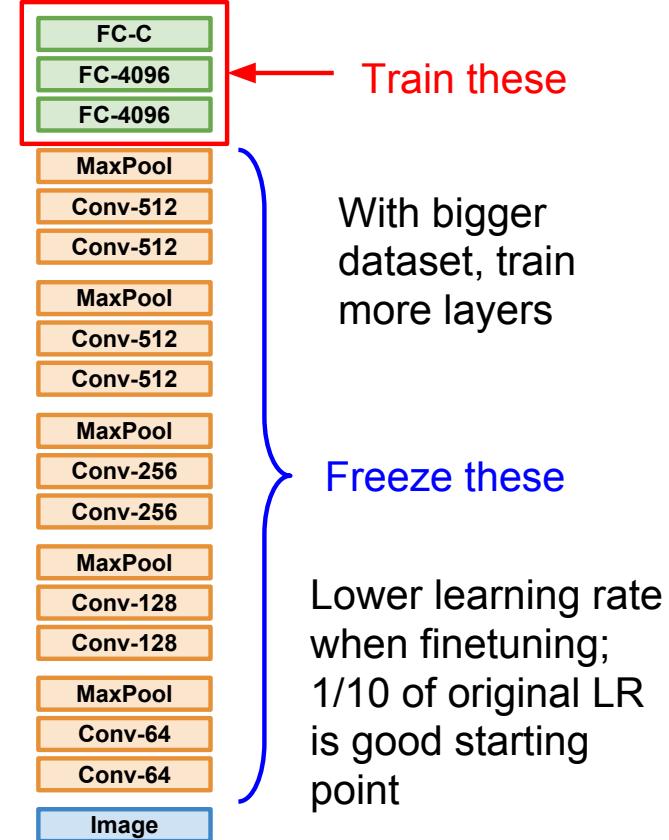
1. Train on Imagenet



2. Small Dataset (C classes)



3. Bigger dataset





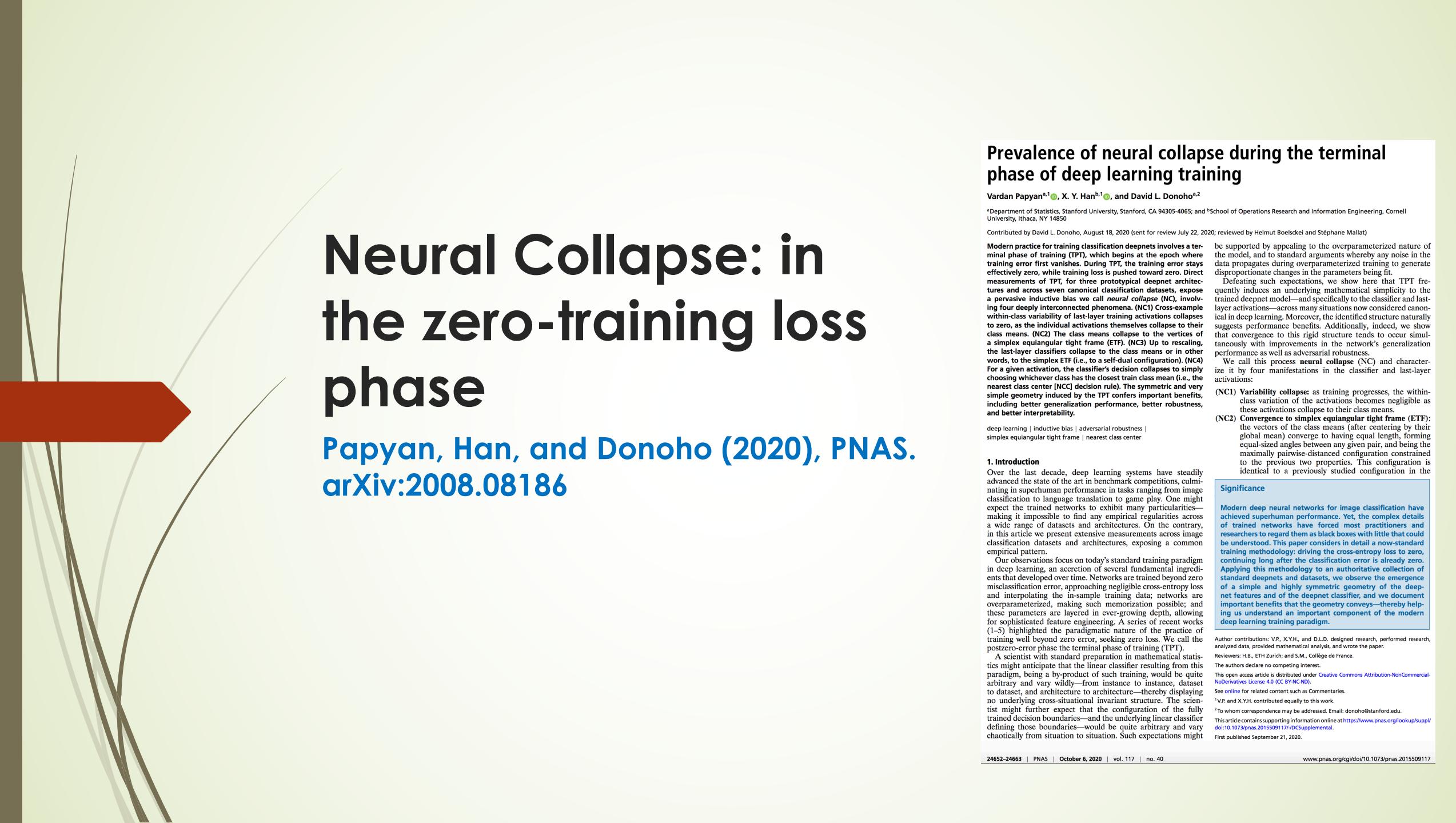
More specific

More generic

	very similar dataset	very different dataset
very little data	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
quite a lot of data	Finetune a few layers	Finetune a larger number of layers

Summary

- ▶ Feature Extraction vs. Fine-Tuning:
 - ▶ Feature extraction usually refers to freeze the bottom (early layers) and retrain the top (last) layer
 - ▶ Fine-Tuning usually refers to retrain the last few layers or the whole network initialized from pretrained parameters
 - ▶ They are both called transfer learning
- ▶ Jupyter notebook examples with pytorch:
 - ▶ https://github.com/aifin-hkust/aifin-hkust.github.io/blob/master/2020/notebook/finetuning_resnet.ipynb



Neural Collapse: in the zero-training loss phase

Papyan, Han, and Donoho (2020), PNAS.
arXiv:2008.08186

Prevalence of neural collapse during the terminal phase of deep learning training

Vardan Papyan^{a,1}, X. Y. Han^{b,1}, and David L. Donoho^{a,2}

^aDepartment of Statistics, Stanford University, Stanford, CA 94305-4065; and ^bSchool of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850

Contributed by David L. Donoho, August 18, 2020 (sent for review July 22, 2020; reviewed by Helmut Boelskei and Stéphane Mallat)

Modern practice for training classification deepnets involves a terminal phase of training (TPT), which begins at the epoch where training error first vanishes. During TPT, the training error stays effectively zero, while training loss is pushed toward zero. Direct measurements of TPT, for three prototypical deepnet architectures and across seven canonical classification datasets, expose a pervasive inductive bias we call *neural collapse* (NC), involving four deeply interconnected phenomena. (NC1) Cross-example within-class variability of last-layer training activations collapses to zero, as the individual activations themselves collapse to their class means. (NC2) The class means collapse to the vertices of a simplex equiangular tight frame (ETF). (NC3) Up to rescaling, the last-layer classifier collapses to the class mean or in other words, to the simplex ETF (i.e., to a self-dual configuration). (NC4) For a given activation, the classifier's decision collapses to simply choosing whichever class has the closest train class mean (i.e., the nearest class center [NCC] decision rule). The symmetric and very simple geometry induced by the TPT confers important benefits, including better generalization performance, better robustness, and better interpretability.

deep learning | inductive bias | adversarial robustness | simplex equiangular tight frame | nearest class center

1. Introduction

Over the last decade, deep learning systems have steadily advanced the state of the art in benchmark competitions, culminating in superhuman performance in tasks ranging from image classification to language translation to game play. One might expect the trained networks to exhibit many particularities—making it impossible to find any empirical regularities across a wide range of datasets and architectures. On the contrary, in this article we present extensive measurements across image classification datasets and architectures, exposing a common empirical pattern.

Our observations focus on today's standard training paradigm in deep learning, an accretion of several fundamental ingredients that developed over time. Networks are trained beyond zero misclassification error, approaching negligible cross-entropy loss and interpolating the in-sample training data; networks are overparameterized, making such memorization possible; and these parameters are layered in ever-growing depth, allowing for sophisticated feature engineering. A series of recent works (1–5) highlighted the paradigmatic nature of the practice of training well beyond zero error, seeking zero loss. We call the postzero-error phase the terminal phase of training (TPT).

A scientist with standard preparation in mathematical statistics might anticipate that the linear classifier resulting from this paradigm, being a by-product of such training, would be quite arbitrary and vary wildly—from instance to instance, dataset to dataset, and architecture to architecture—thereby displaying no underlying cross-situational invariant structure. The scientist might further expect that the configuration of the fully trained decision boundaries—and the underlying linear classifier defining those boundaries—would be quite arbitrary and vary chaotically from situation to situation. Such expectations might

be supported by appealing to the overparameterized nature of the model, and to standard arguments whereby any noise in the data propagates through overparameterized training to generate disproportionate changes in the parameters being fit.

Defeating such expectations, we show here that TPT frequently induces an underlying mathematical simplicity to the trained deepnet model—and specifically to the classifier and last-layer activations—across many situations now considered canonical in deep learning. Moreover, the identified structure naturally suggests performance benefits. Additionally, indeed, we show that convergence to this rigid structure tends to occur simultaneously with improvements in the network's generalization performance as well as adversarial robustness.

We call this process *neural collapse* (NC) and characterize it by four manifestations in the classifier and last-layer activations:

(NC1) **Variability collapse:** as training progresses, the within-class variation of the activations becomes negligible as these activations collapse to their class means.

(NC2) **Convergence to simplex equiangular tight frame (ETF):** the vectors of the class means (after centering by their global mean) converge to having equal length, forming equal-sized angles between any given pair, and being the maximally pairwise-distanced configuration constrained to the previous two properties. This configuration is identical to a previously studied configuration in the

Significance

Modern deep neural networks for image classification have achieved superhuman performance. Yet, the complex details of trained networks have forced most practitioners and researchers to regard them as black boxes with little that could be understood. This paper considers in detail a now-standard training methodology: driving the cross-entropy loss to zero, continuing long after the classification error is already zero. Applying this methodology to an authoritative collection of standard deepnets and datasets, we observe the emergence of a simple and highly symmetric geometry of the deepnet features and of the deepnet classifier, and we document important benefits that the geometry conveys—thereby helping us understand an important component of the modern deep learning training paradigm.

Author contributions: V.P., X.Y.H., and D.L.D. designed research, performed research, analyzed data, provided mathematical analysis, and wrote the paper.

Reviewers: H.B., ETH Zurich; and S.M., Collège de France.

The authors declare no competing interest.
This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

See online for related content such as Commentaries.

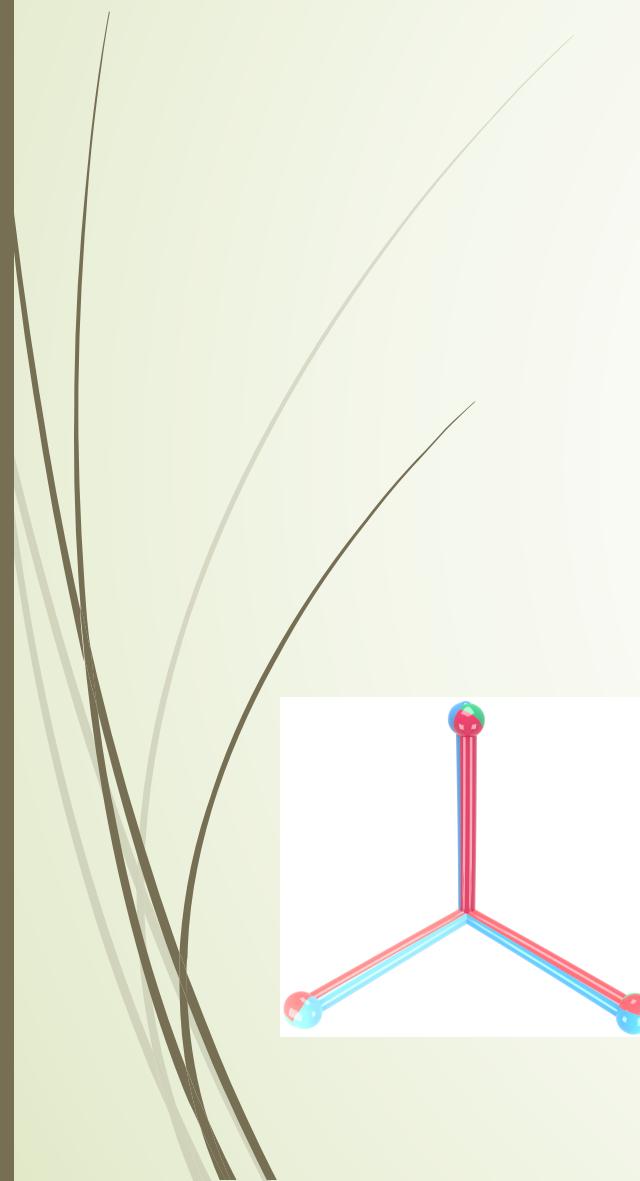
¹V.P. and X.Y.H. contributed equally to this work.

²To whom correspondence may be addressed. Email: donoho@stanford.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2015509117/DCSupplemental>.
First published September 21, 2020.

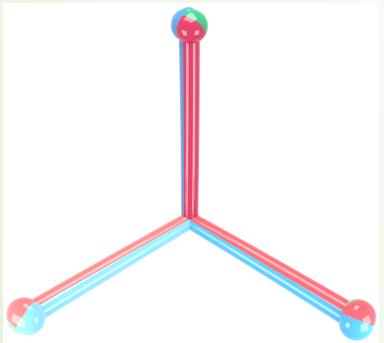
Neural Collapse phenomena, in post-zero-training-error phase

- ▶ **(NC1) Variability collapse:** As training progresses, the within-class variation of the activations becomes negligible as these activations collapse to their class-means.
 - ▶ **(NC2) Convergence to Simplex ETF:** The vectors of the class-means (after centering by their global-mean) converge to having equal length, forming equal-sized angles between any given pair, and being the maximally pairwise-distanced configuration constrained to the previous two properties. This configuration is identical to a previously studied configuration in the mathematical sciences known as Simplex **Equiangular Tight Frame (ETF)**.
-
- ▶ *Papyan, Han, and Donoho (2020), PNAS. arXiv:2008.08186*
 - ▶ Visualization: <https://purl.stanford.edu/br193mh4244>



Definition 1 (Simplex ETF). A *standard* Simplex ETF is a collection of points in \mathbb{R}^C specified by the columns of

$$\mathbf{M}^* = \sqrt{\frac{C}{C-1}} \left(\mathbf{I} - \frac{1}{C} \mathbf{1}\mathbf{1}^\top \right), \quad [1]$$



where $\mathbf{I} \in \mathbb{R}^{C \times C}$ is the identity matrix, and $\mathbf{1}_C \in \mathbb{R}^C$ is the ones vector. In this paper, we allow other poses, as well as rescaling, so the *general* Simplex ETF consists of the points specified by the columns of $\mathbf{M} = \alpha \mathbf{U} \mathbf{M}^* \in \mathbb{R}^{p \times C}$, where $\alpha \in \mathbb{R}_+$ is a scale factor, and $\mathbf{U} \in \mathbb{R}^{p \times C}$ ($p \geq C$) is a partial orthogonal matrix ($\mathbf{U}^\top \mathbf{U} = \mathbf{I}$).

Be careful!

- For imbalanced (long tail) classifications, **minority** classes may **collapse** and be absorbed by the majority classes
- Fang, He, Long, Su, PNAS 2021, 118(43):e2103091118

Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training

Cong Fang^{a,1}, Hangfeng He^a, Qi Long^b, and Weijie J. Su^{c,2}

^aDepartment of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104; ^bDepartment of Biostatistics, Epidemiology, and Informatics, University of Pennsylvania, Philadelphia, PA 19104; and ^cDepartment of Statistics and Data Science, University of Pennsylvania, Philadelphia, PA 19104

Edited by David L. Donoho, Stanford University, Stanford, CA, and approved August 30, 2021 (received for review February 15, 2021)

In this paper, we introduce the *Layer-Pealed Model*, a nonconvex, yet analytically tractable, optimization program, in a quest to better understand deep neural networks that are trained for a sufficiently long time. As the name suggests, this model is derived by isolating the topmost layer from the remainder of the neural network, followed by imposing certain constraints separately on the two parts of the network. We demonstrate that the Layer-Pealed Model, albeit simple, inherits many characteristics of well-trained neural networks, thereby offering an effective tool for explaining and predicting common empirical patterns of deep-learning training. First, when working on class-balanced datasets, we prove that any solution to this model forms a simplex equiangular tight frame, which, in part, explains the recently discovered phenomenon of neural collapse [V. Papyan, X. Y. Han, D. L. Donoho, *Proc. Natl. Acad. Sci. U.S.A.* 117, 24652–24663 (2020)]. More importantly, when moving to the imbalanced case, our analysis of the Layer-Pealed Model reveals a hitherto-unknown phenomenon that we term *Minority Collapse*, which fundamentally limits the performance of deep-learning models on the minority classes. In addition, we use the Layer-Pealed Model to gain insights into how to mitigate Minority Collapse. Interestingly, this phenomenon is first predicted by the Layer-Pealed Model before being confirmed by our computational experiments.

deep learning | neural collapse | class imbalance

Introduction

In the past decade, deep learning has achieved remarkable performance across a range of scientific and engineering domains (1–3). Interestingly, these impressive accomplishments were mostly achieved by heuristics and tricks, though often plausible, without much principled guidance from a theoretical perspective. On the flip side, however, this reality suggests the great potential a theory could have for advancing the development of deep-learning methodologies in the coming decade.

Unfortunately, it is not easy to develop a theoretical foundation for deep learning. Perhaps the most difficult hurdle lies in the nonconvexity of the optimization problem for training neural networks, which, loosely speaking, stems from the interaction between different layers of neural networks. To be more precise, consider a neural network for K -class classification (in logits), which in its simplest form reads¹

$$\begin{aligned} f(\mathbf{x}; \mathbf{W}_{\text{full}}) = & \mathbf{b}_L + \mathbf{W}_L \sigma(\mathbf{b}_{L-1} \\ & + \mathbf{W}_{L-1} \sigma(\dots \sigma(\mathbf{b}_1 + \mathbf{W}_1 \mathbf{x}) \dots)). \end{aligned}$$

Here, $\mathbf{W}_{\text{full}} := \{\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L\}$ denotes the weights of the L layers, $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_L\}$ denotes the biases, and $\sigma(\cdot)$ is a nonlinear activation function such as the rectified linear unit (ReLU). Owing to the complex and nonlinear interaction between the L layers, when applying stochastic gradient descent to the optimization problem

*The softmax step is implicitly included in the loss function, and we omit other operations such as max-pooling for simplicity.

$$\min_{\mathbf{W}_{\text{full}}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(f(\mathbf{x}_{k,i}; \mathbf{W}_{\text{full}}), \mathbf{y}_k) + \frac{\lambda}{2} \|\mathbf{W}_{\text{full}}\|^2, \quad [1]$$

with a loss function \mathcal{L} for training the neural network, it becomes very difficult to pinpoint how a given layer influences the output f (above, $\{\mathbf{x}_{k,i}\}_{i=1}^{n_k}$ denotes the training examples in the k -th class, with label \mathbf{y}_k , $N = n_1 + \dots + n_K$ is the total number of training examples, $\lambda > 0$ is the weight decay parameter, and $\|\cdot\|$ throughout the paper is the ℓ_2 norm). Worse, this difficulty in analyzing deep-learning models is compounded by an ever-growing number of layers.

Therefore, any attempt to develop a tractable and comprehensive theory for demystifying deep learning would presumably first need to simplify the interaction between a large number of layers. Following this intuition, in this paper, we introduce the following optimization program as a *surrogate* model for Eq. 1 with the goal of unveiling quantitative patterns of deep neural networks:

$$\begin{aligned} \min_{\mathbf{W}_L, \mathbf{H}} \quad & \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}(\mathbf{W}_L \mathbf{h}_{k,i}, \mathbf{y}_k) \\ \text{s.t.} \quad & \frac{1}{K} \sum_{k=1}^K \|\mathbf{w}_k\|^2 \leq E_W, \quad \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} \|\mathbf{h}_{k,i}\|^2 \leq E_H, \quad [2] \end{aligned}$$

Significance

The remarkable development of deep learning over the past decade relies heavily on sophisticated heuristics and tricks. To better exploit its potential in the coming decade, perhaps a rigorous framework for reasoning about deep learning is needed, which, however, is not easy to build due to the intricate details of neural networks. For near-term purposes, a practical alternative is to develop a mathematically tractable surrogate model, yet maintaining many characteristics of neural networks. This paper proposes a model of this kind that we term the *Layer-Pealed Model*. The effectiveness of this model is evidenced by, among others, its ability to reproduce a known empirical pattern and to predict a hitherto-unknown phenomenon when training deep-learning models on imbalanced datasets.

Author contributions: C.F., H.H., Q.L., and W.J.S. designed research; C.F., H.H., and W.J.S. performed research; C.F., H.H., and W.J.S. contributed new reagents/analytic tools; C.F., H.H., and W.J.S. analyzed data; and C.F., H.H., Q.L., and W.J.S. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.
This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

¹Present address: Department of Key Laboratory of Machine Perception, Peking University, Beijing 100871, China.

²To whom correspondence may be addressed. Email: suw@wharton.upenn.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2103091118/DCSupplemental>.

Published October 20, 2021.

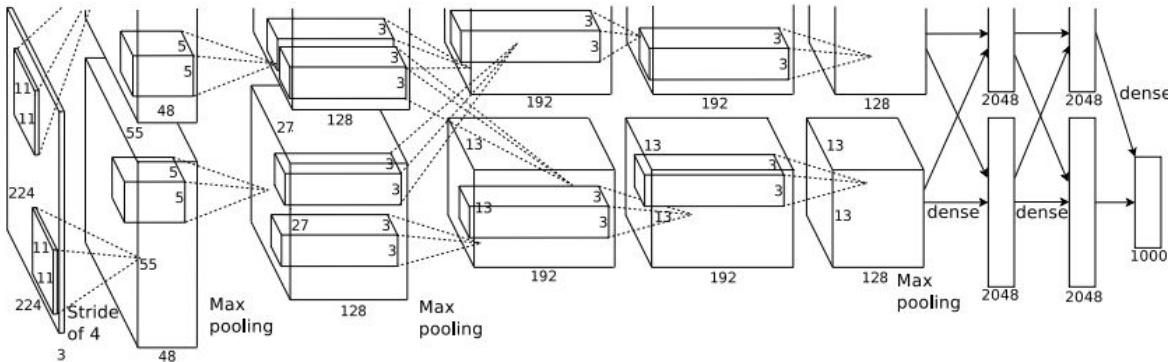


Visualizing Convolutional Networks

Understanding intermediate neurons?



Input Image:
3 x 224 x 224



What are the intermediate features looking for?

Class Scores:
1000 numbers

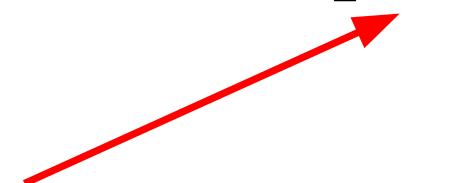
Visualizing CNN Features: Gradient Ascent

- **Gradient ascent:** Generate a synthetic image that maximally activates a neuron

$$I^* = \arg \max_I f(I) + R(I)$$

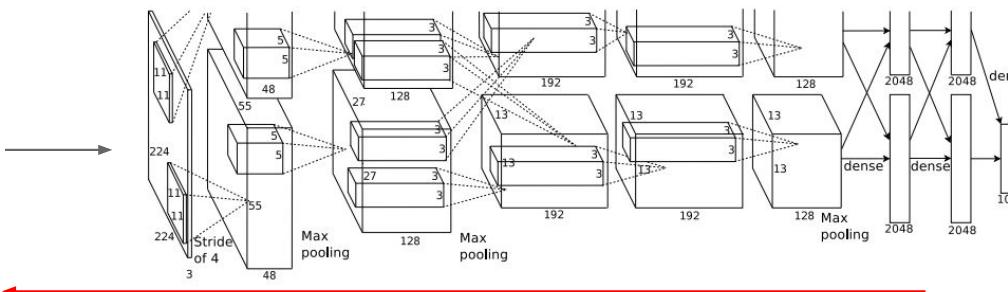
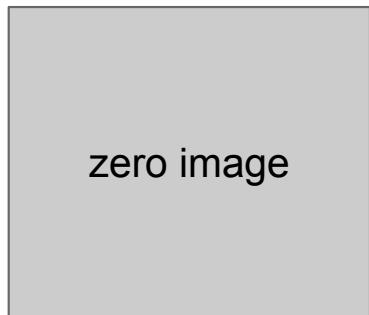
Neuron value

Natural image regularizer



Example: Class Visualization of CNN via Gradient Ascent

1. Initialize image to zeros



$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

score for class c (before Softmax)

Repeat:

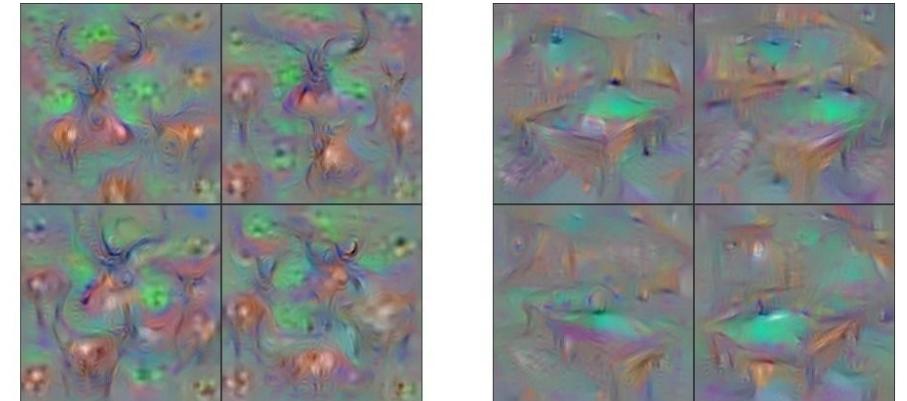
2. Forward image to compute current scores
3. Backprop to get gradient of neuron value with respect to image pixels
4. Make a small update to the image

Visualizing CNN Features: Gradient Ascent

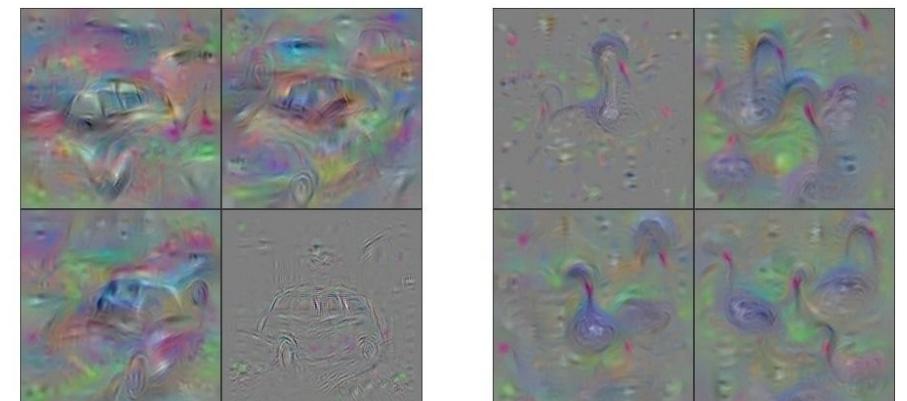
$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

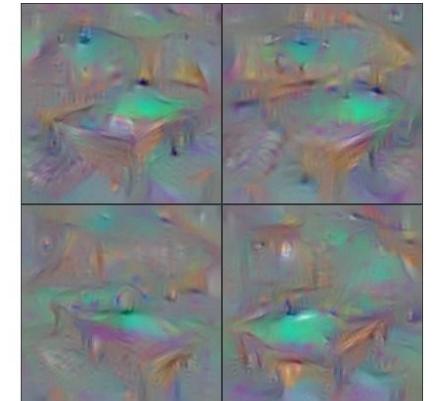
- (1) Gaussian blur image
- (2) Clip pixels with small values to 0
- (3) Clip pixels with small gradients to 0



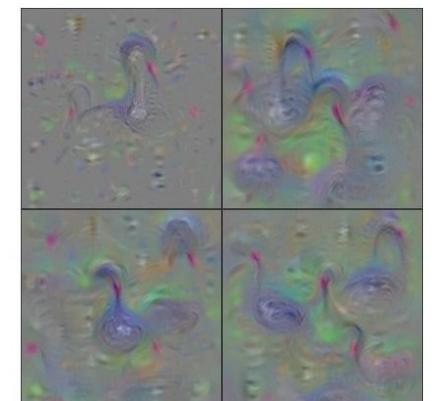
Hartebeest



Station Wagon



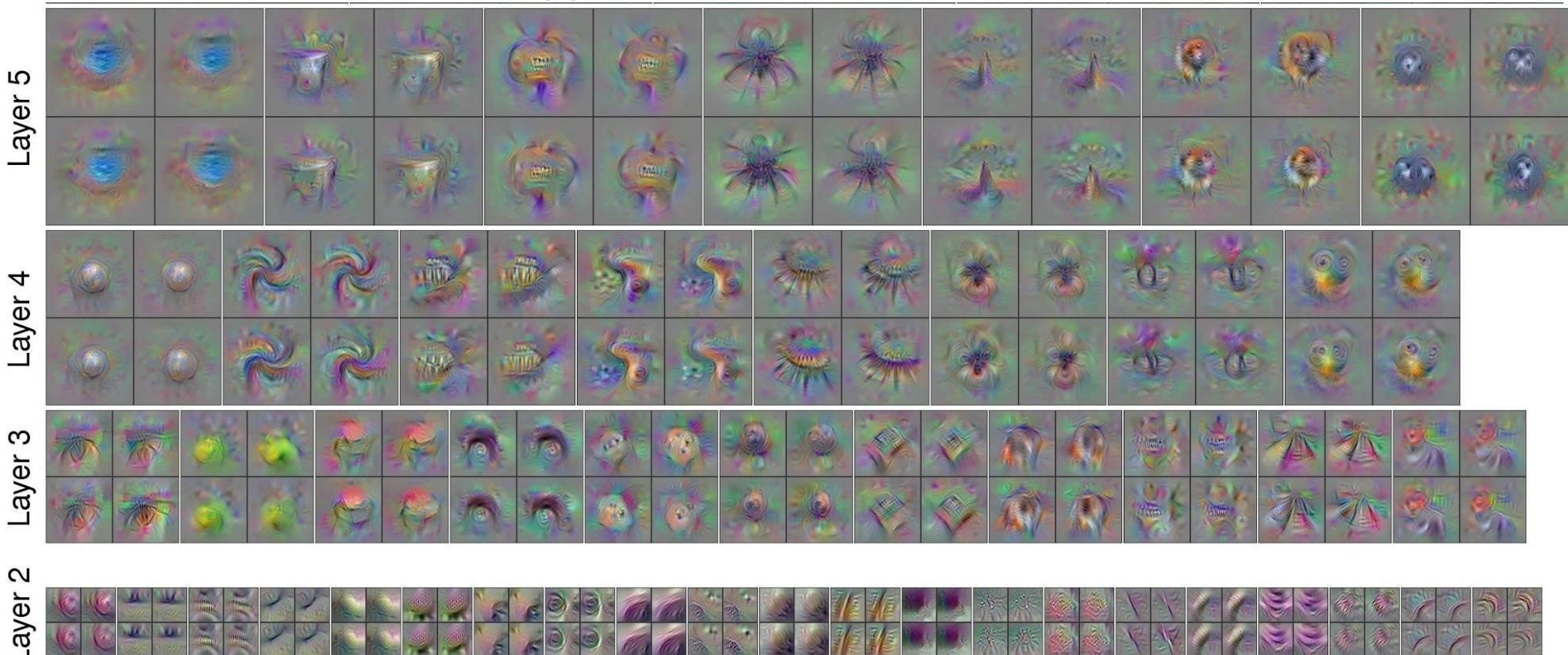
Billiard Table



Black Swan

Visualizing CNN Features: Gradient Ascent

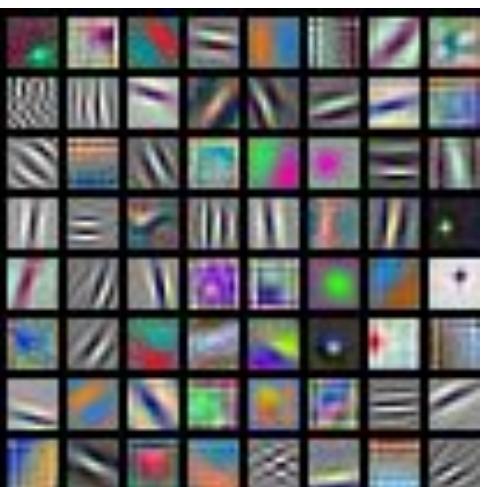
Use the same approach to visualize intermediate features



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014. Reproduced with permission.

It's easy to visualize early layers

First Layer: Visualize Filters



AlexNet:
 $64 \times 3 \times 11 \times 11$



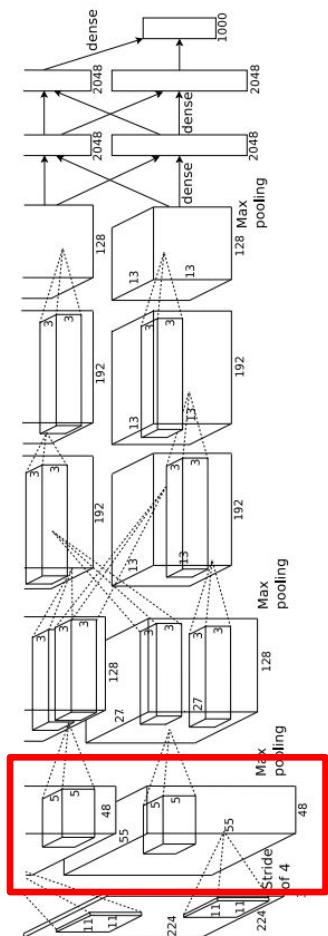
ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet-101:
 $64 \times 3 \times 7 \times 7$



DenseNet-121:
 $64 \times 3 \times 7 \times 7$



Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

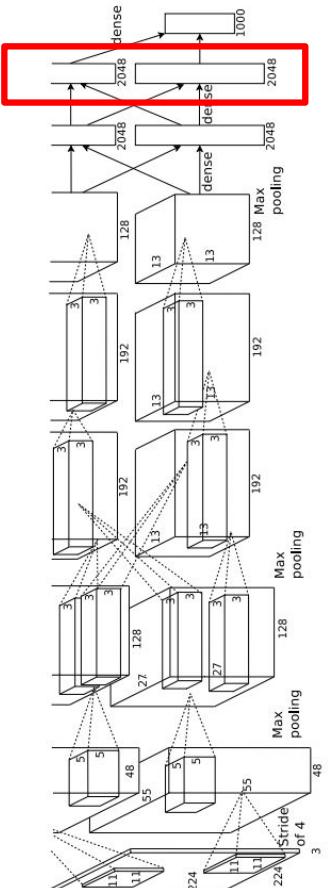
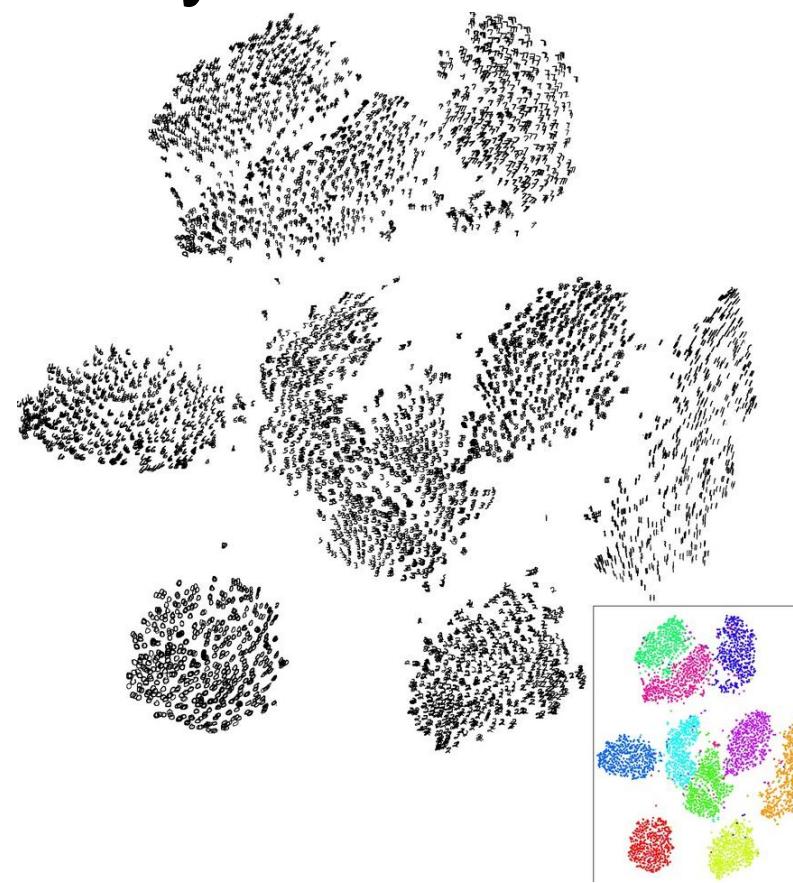
Last layers are hard to visualize

Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

Simple algorithm: Principle Component Analysis (PCA)

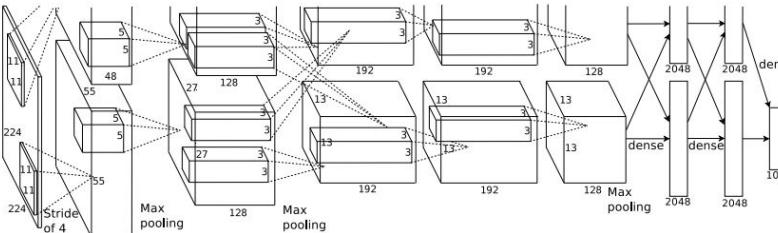
More complex: t-SNE



Van der Maaten and Hinton, “Visualizing Data using t-SNE”, JMLR 2008
Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.

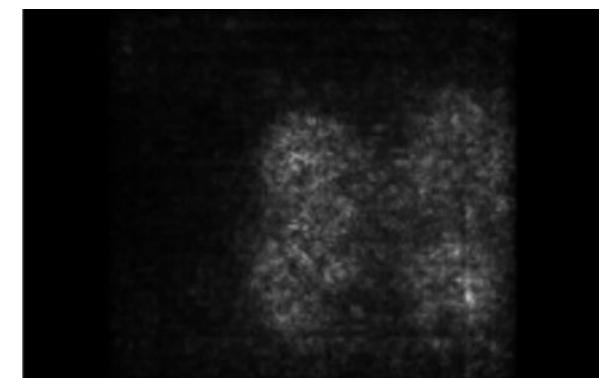
Saliency Maps

How to tell which pixels matter for classification?



Dog

Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels

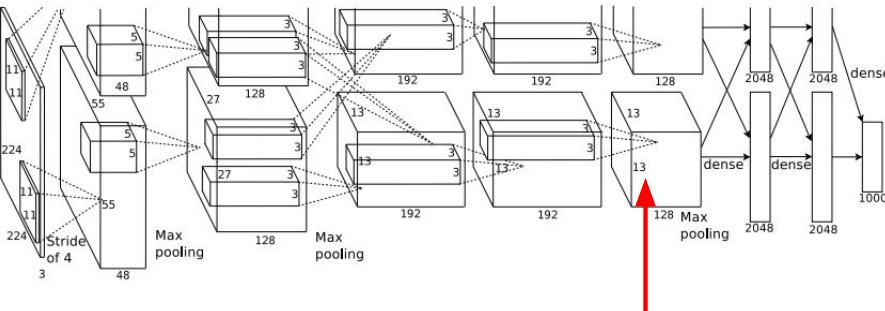


Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

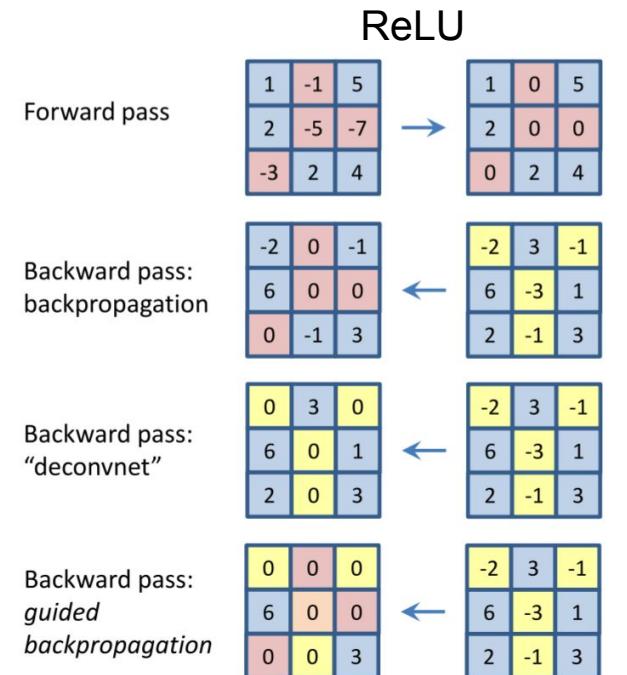
Guided BP

Intermediate features via (guided) backprop



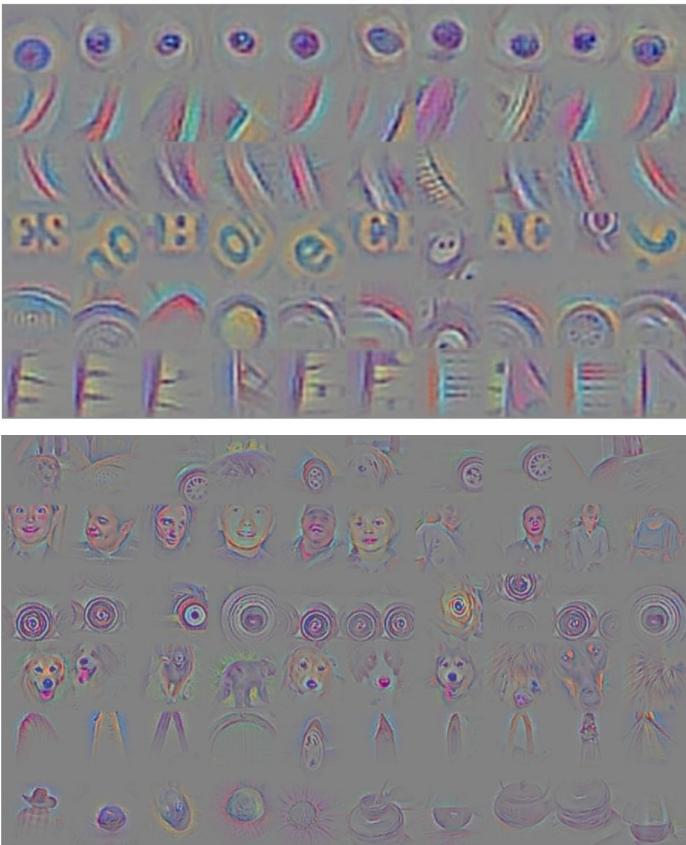
Pick a single intermediate neuron, e.g. one value in $128 \times 13 \times 13$ conv5 feature map

Compute gradient of neuron value with respect to image pixels



Images come out nicer if you only backprop positive gradients through each ReLU (guided backprop)

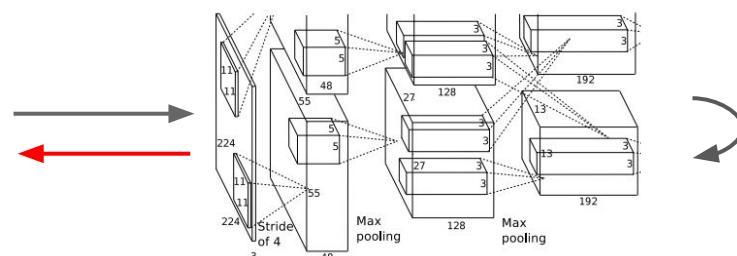
Intermediate features via Guided BP



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

DeepDream: amplifying features

Rather than synthesizing an image to maximize a specific neuron, instead try to **amplify** the neuron activations at some layer in the network

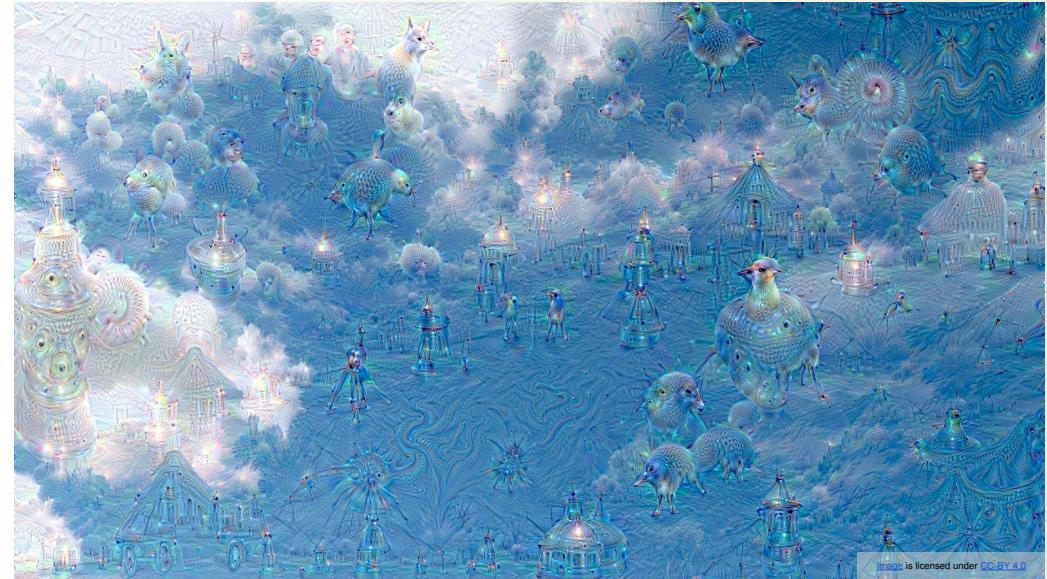


Choose an image and a layer in a CNN; repeat:

1. Forward: compute activations at chosen layer
2. Set gradient of chosen layer *equal to its activation*
3. Backward: Compute gradient on image
4. Update image

Equivalent to:
 $I^* = \arg \max_I \sum_i f_i(I)^2$

Example: DeepDream of Sky



"Admiral Dog!"



"The Pig-Snail"

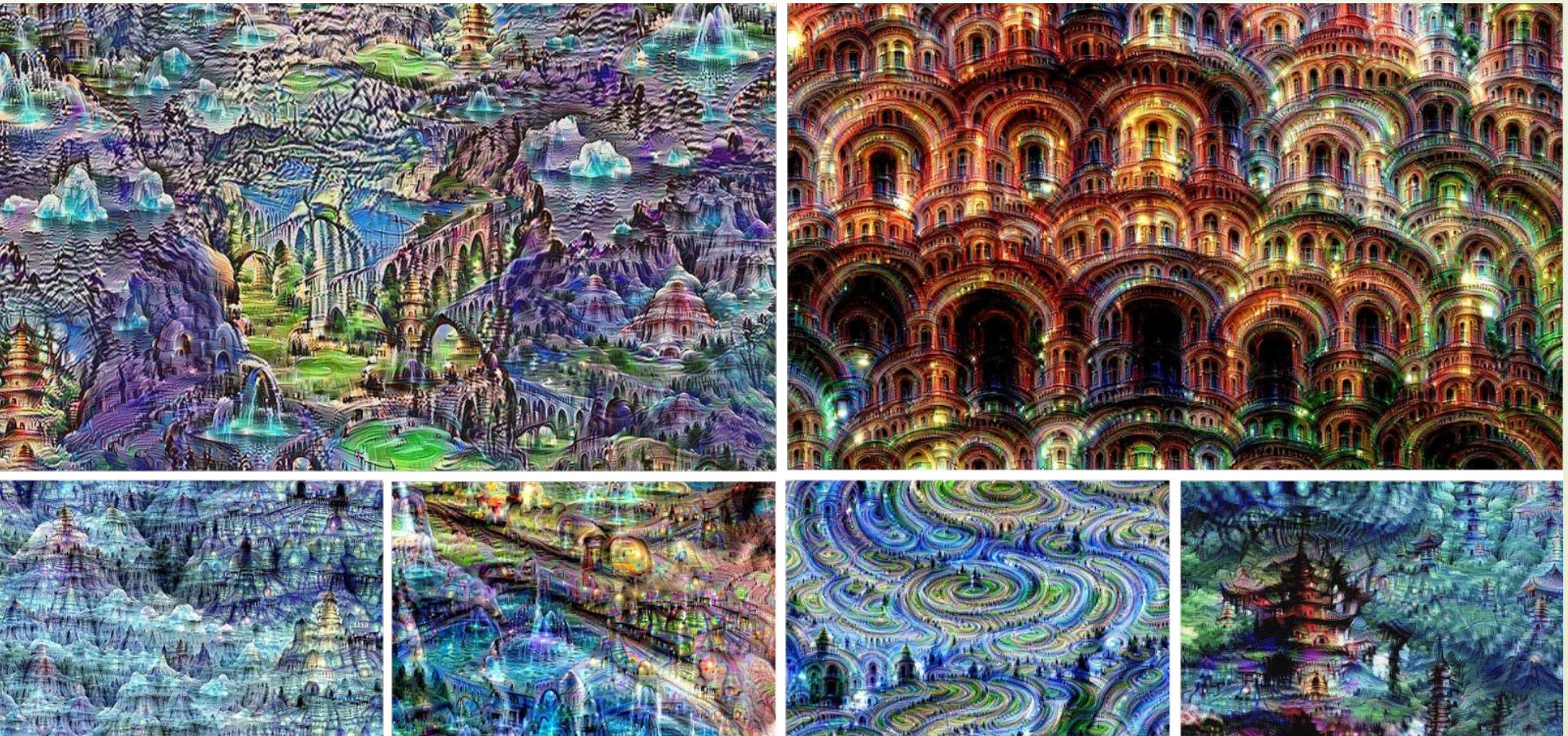


"The Camel-Bird"



"The Dog-Fish"

More Examples



[Image](#) is licensed under CC-BY 4.0

Python Notebooks

- ▶ An interesting Pytorch Implementation of these visualizatoion methods
 - ▶ <https://github.com/utkuozbulak/pytorch-cnn-visualizations>
- ▶ Some examples demo:
 - ▶ <https://github.com/aifin-hkust/aifin-hkust.github.io/blob/master/2020/notebook/vgg16-visualization.ipynb>
 - ▶ <https://github.com/aifin-hkust/aifin-hkust.github.io/blob/master/2020/notebook/vgg16-heatmap.ipynb>



Neural Style

Example: The Noname Lake in PKU

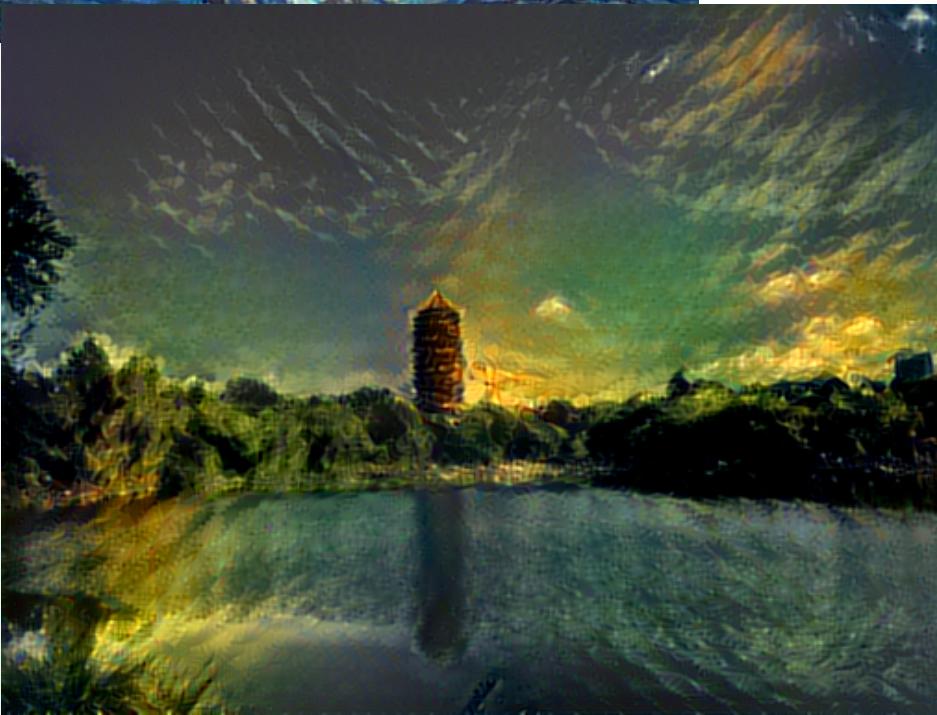
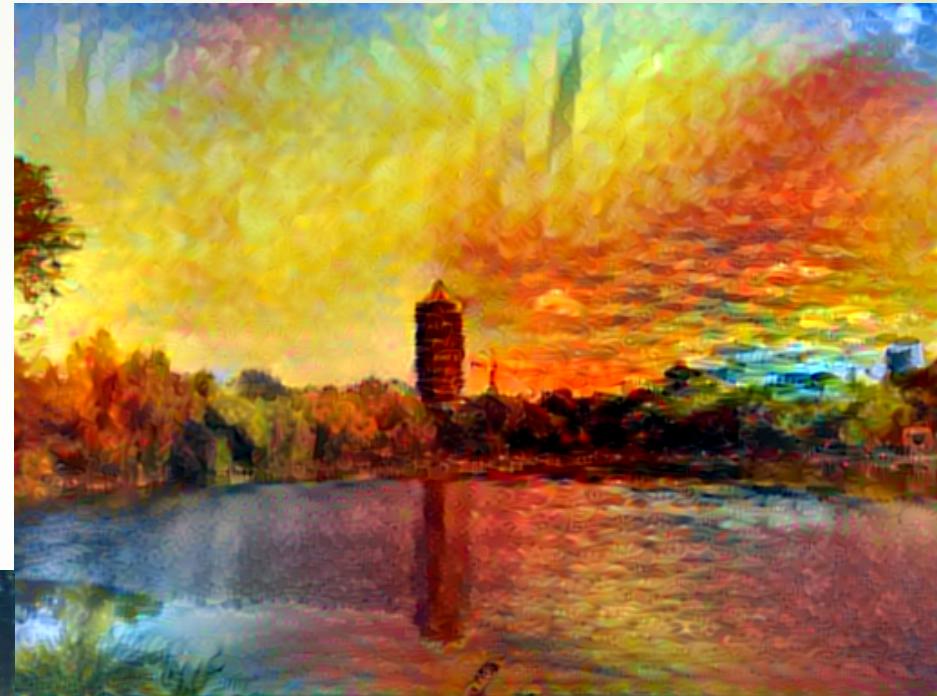
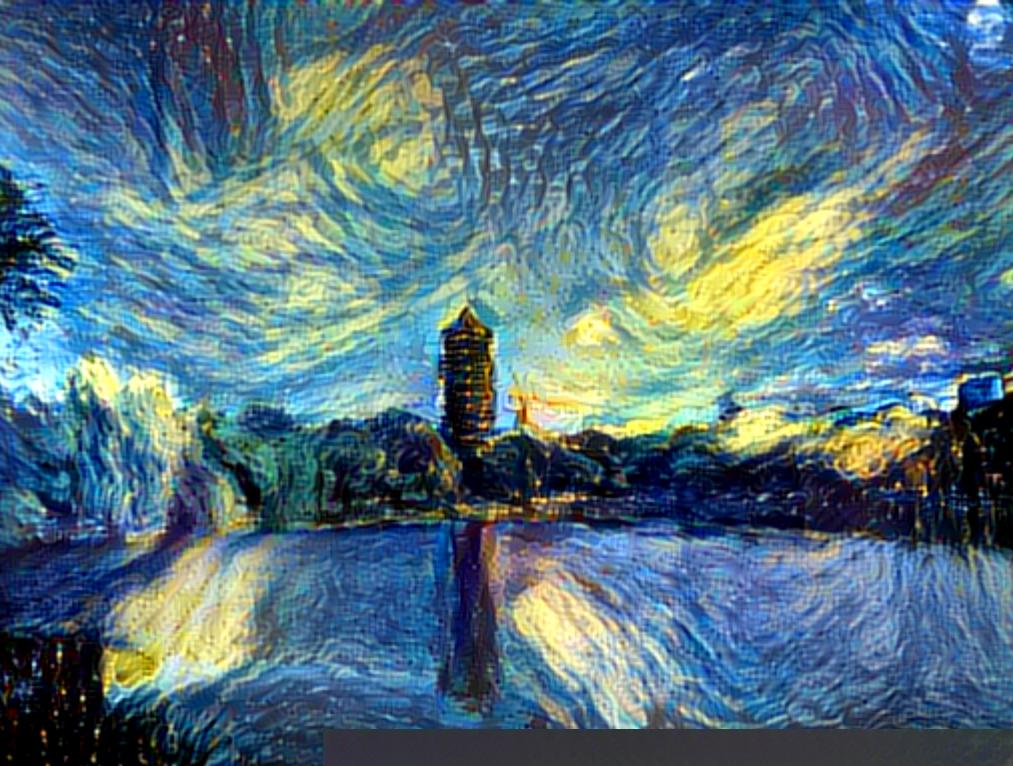




Left: Vincent Van Gogh, Starry Night
Right: Claude Monet, Twilight Venice
Bottom: William Turner, Ship Wreck



开始吧



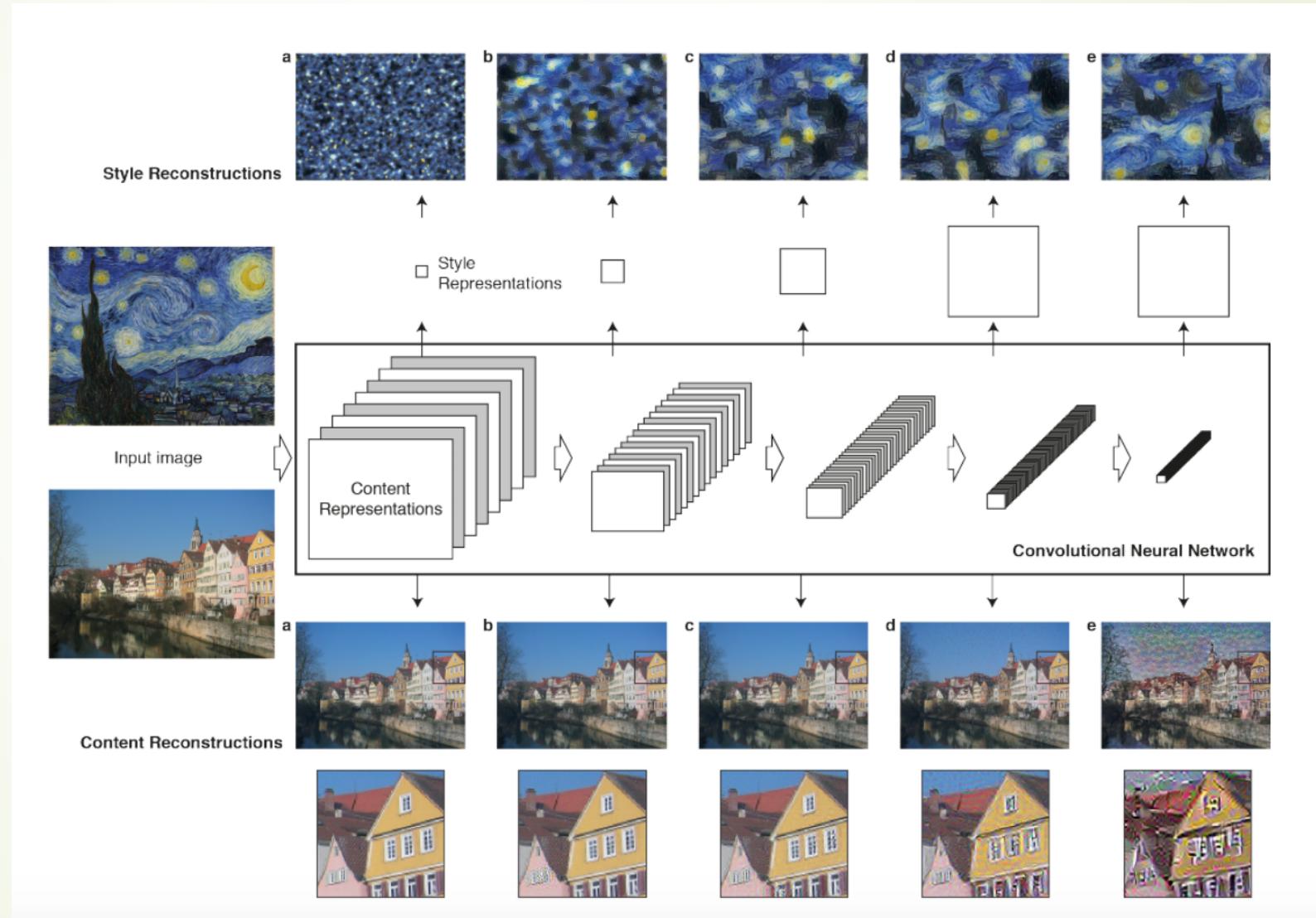
Application of Deep Learning:
Content-Style synthetic
pictures
By “neural-style”



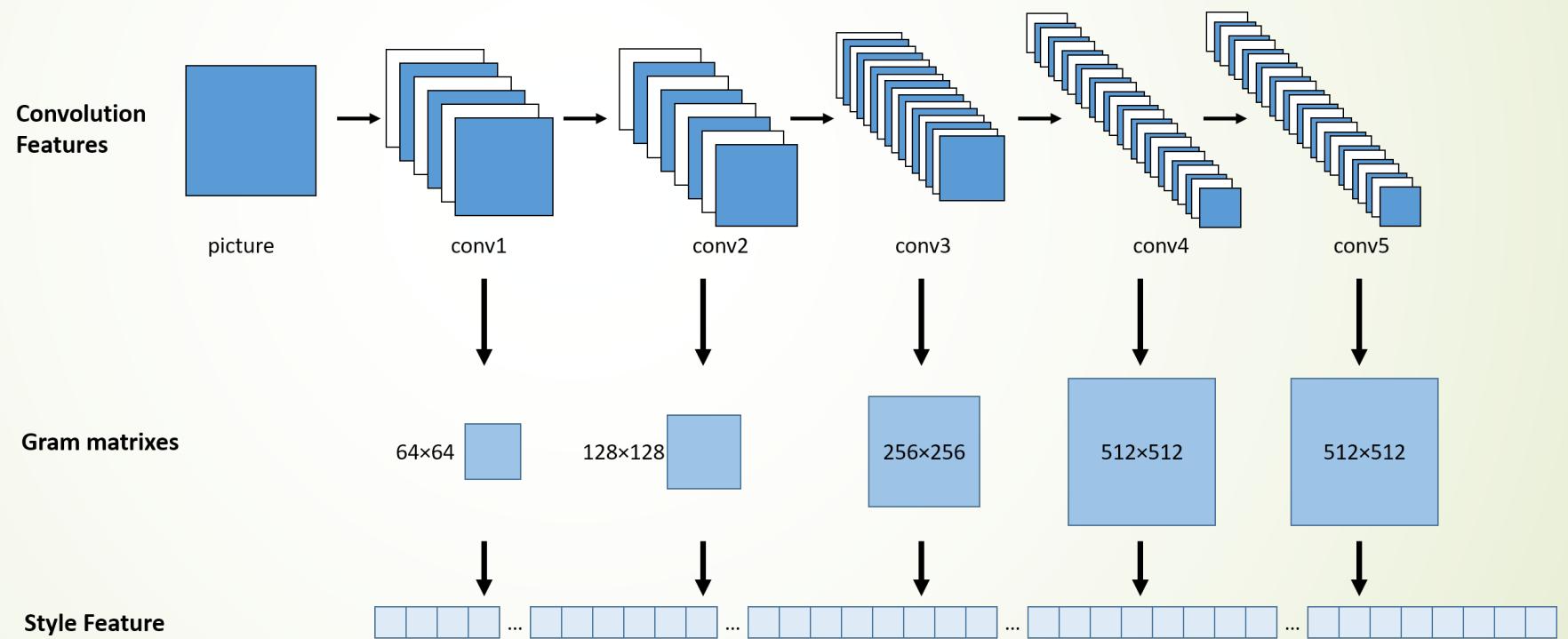
Neural Style

- ▶ J C Johnson's Website: <https://github.com/jcjohnson/neural-style>
- ▶ A torch implementation of the paper
 - ▶ *A Neural Algorithm of Artistic Style*,
 - ▶ by Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge.
 - ▶ <http://arxiv.org/abs/1508.06576>

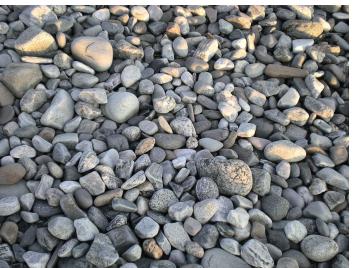
Style-Content Feature Extraction



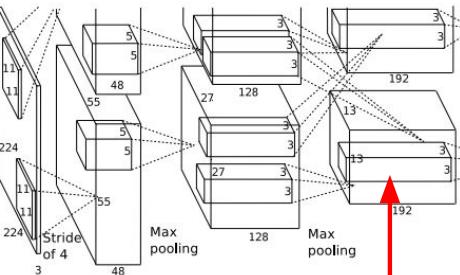
Style Features as Second Order Statistics



Gram Matrix as Style Features



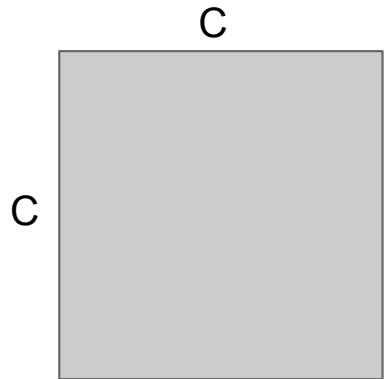
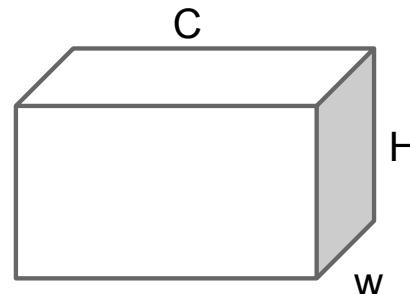
This image is in the public domain.



Each layer of CNN gives $C \times H \times W$ tensor of features; $H \times W$ grid of C -dimensional vectors

Outer product of two C -dimensional vectors gives $C \times C$ matrix measuring co-occurrence

Average over all HW pairs of vectors, giving **Gram matrix** of shape $C \times C$



Efficient to compute; reshape features from

$C \times H \times W$ to $=C \times HW$

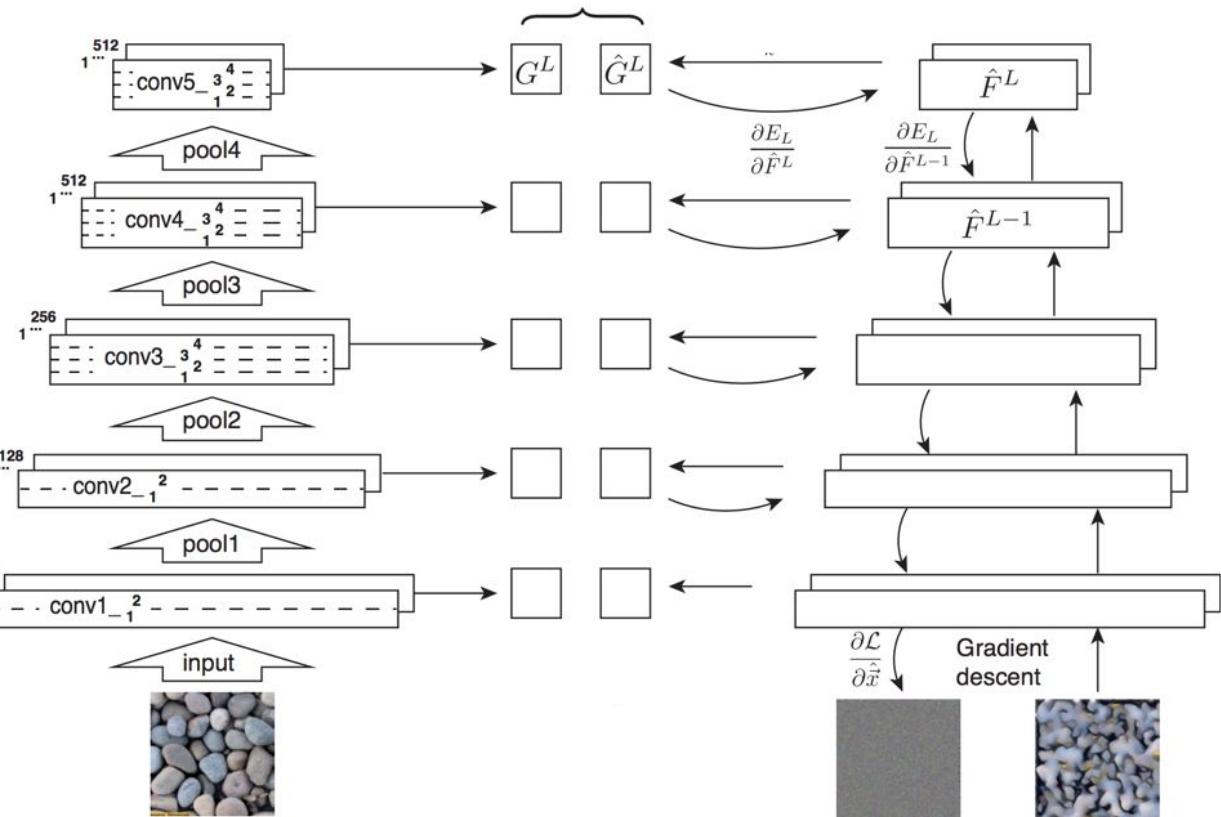
then compute $G = FF^T$

Neural Texture Synthesis

1. Pretrain a CNN on ImageNet (VGG-19)
2. Run input texture forward through CNN, record activations on every layer; layer i gives feature map of shape $C_i \times H_i \times W_i$
3. At each layer compute the *Gram matrix* giving outer product of features:
$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l \text{ (shape } C_i \times C_i\text{)}$$
4. Initialize generated image from random noise
5. Pass generated image through CNN, compute Gram matrix on each layer
6. Compute loss: weighted sum of L2 distance between Gram matrices
7. Backprop to get gradient on image
8. Make gradient step on image
9. GOTO 5

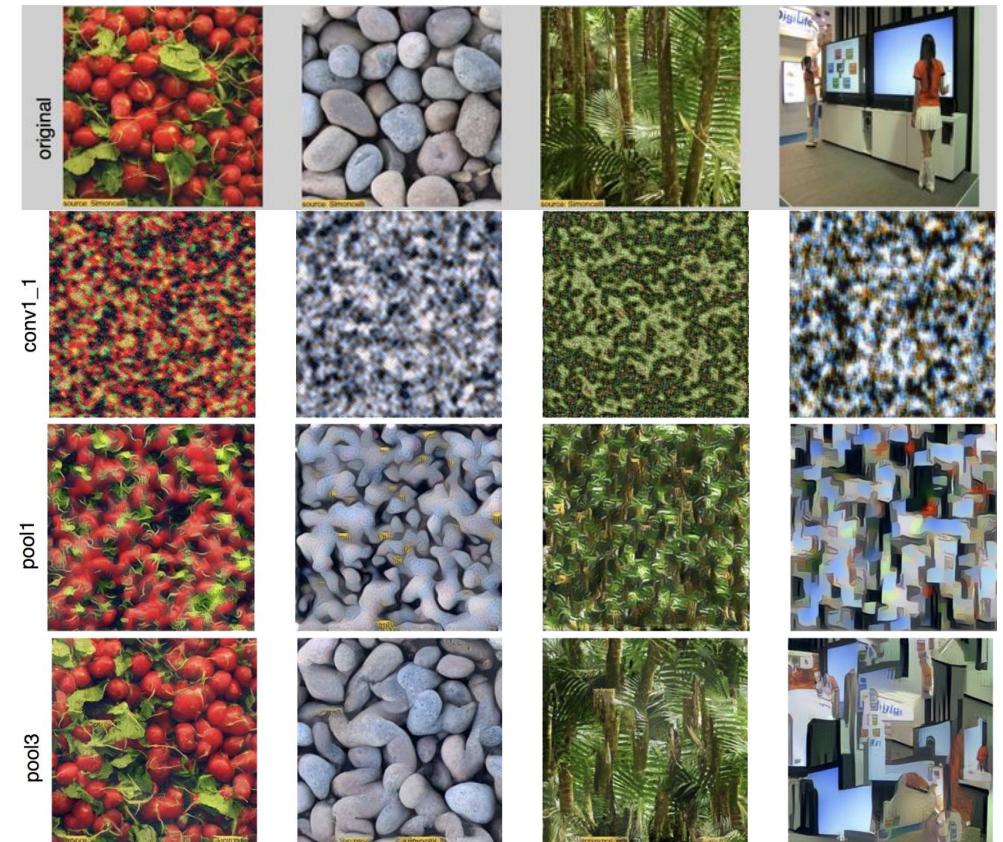
Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
 Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - \hat{G}_{ij}^l)^2 \quad \mathcal{L}(\vec{x}, \hat{\vec{x}}) = \sum_{l=0}^L w_l E_l$$



Neural Texture Synthesis

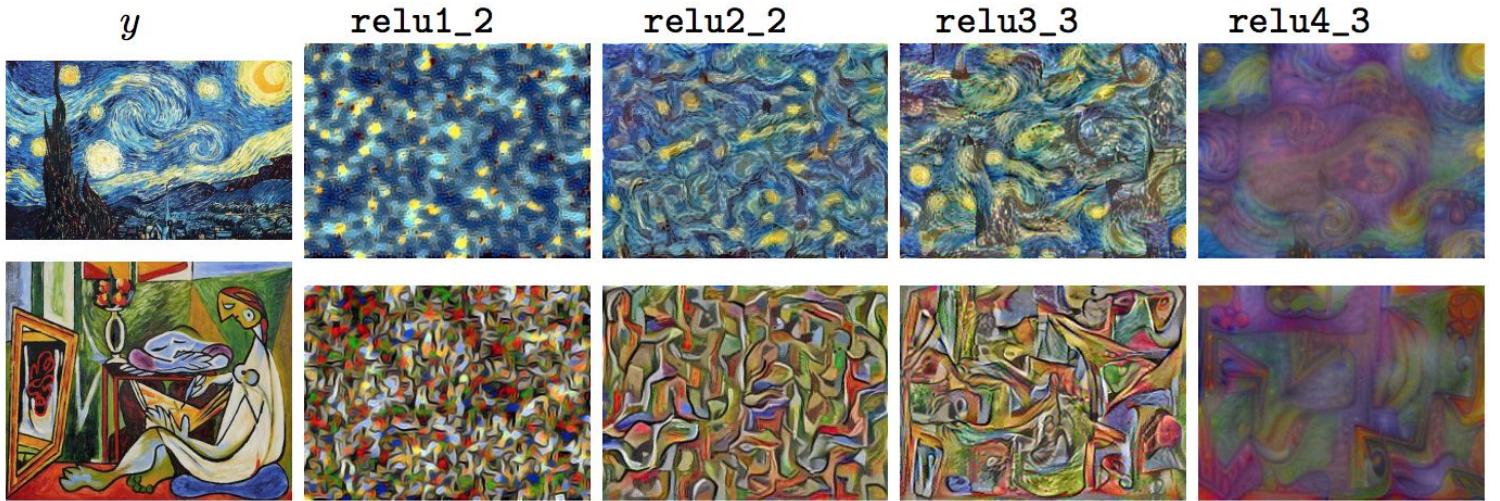
Reconstructing texture from higher layers recovers larger features from the input texture



Gatys, Ecker, and Bethge, "Texture Synthesis Using Convolutional Neural Networks", NIPS 2015
Figure copyright Leon Gatys, Alexander S. Ecker, and Matthias Bethge, 2015. Reproduced with permission.

Neural Texture Synthesis: Gram Reconstruction

Texture synthesis
(Gram
reconstruction)



Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector
- “looks natural” (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

Features of new image

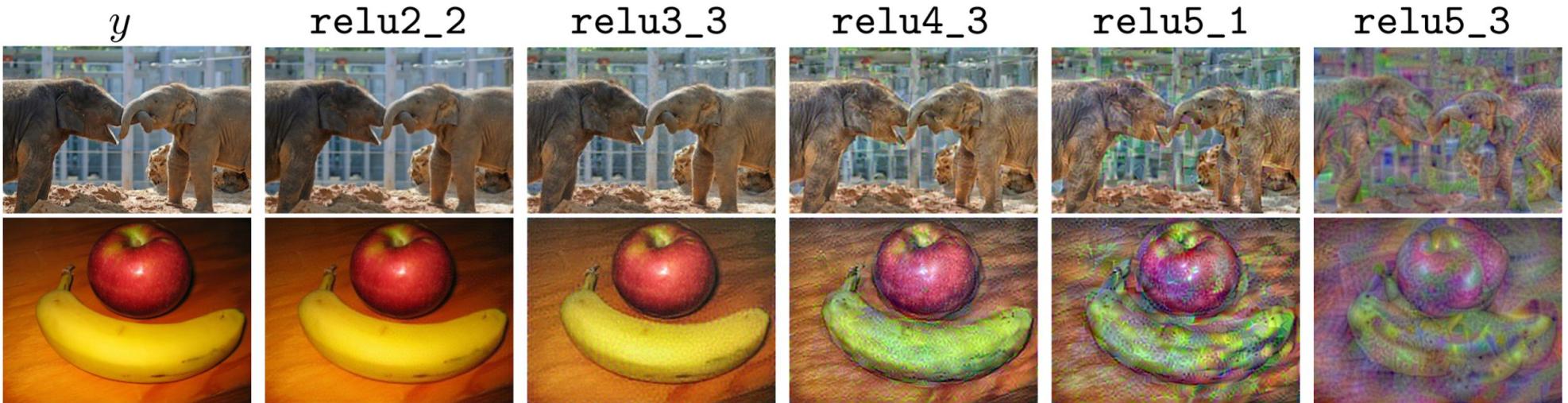
$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$

Total Variation regularizer
(encourages spatial smoothness)

Feature Inversion

Reconstructing from different layers of VGG-16



Mahendran and Vedaldi, "Understanding Deep Image Representations by Inverting Them", CVPR 2015
Figure from Johnson, Alahi, and Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution", ECCV 2016. Copyright Springer, 2016.
Reproduced for educational purposes.



Combined Loss for both Content (1st order statistics) and Style (2nd order statistics: Gram)

$$\mathcal{L}_{content}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2 .$$

$$\mathcal{L}_{style}(\vec{a}, \vec{x}) = \sum_{l=0}^L w_l E_l$$

where

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2 \quad G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Neural Style Transfer

Content Image



[This image](#) is licensed under CC-BY 3.0

+

Style Image



[Starry Night](#) by Van Gogh is in the public domain

=

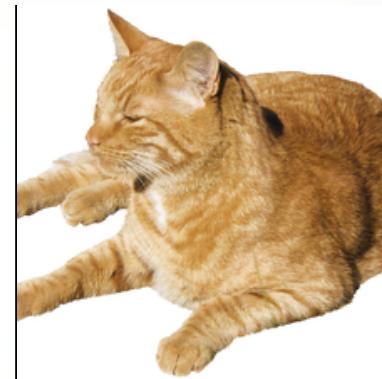


[This image](#) copyright Justin Johnson, 2015. Reproduced with permission.

CNN learns **texture** features, not shapes!



(a) Texture image
81.4% **Indian elephant**
10.3% indri
8.2% black swan



(b) Content image
71.1% **tabby cat**
17.3% grey fox
3.3% Siamese cat

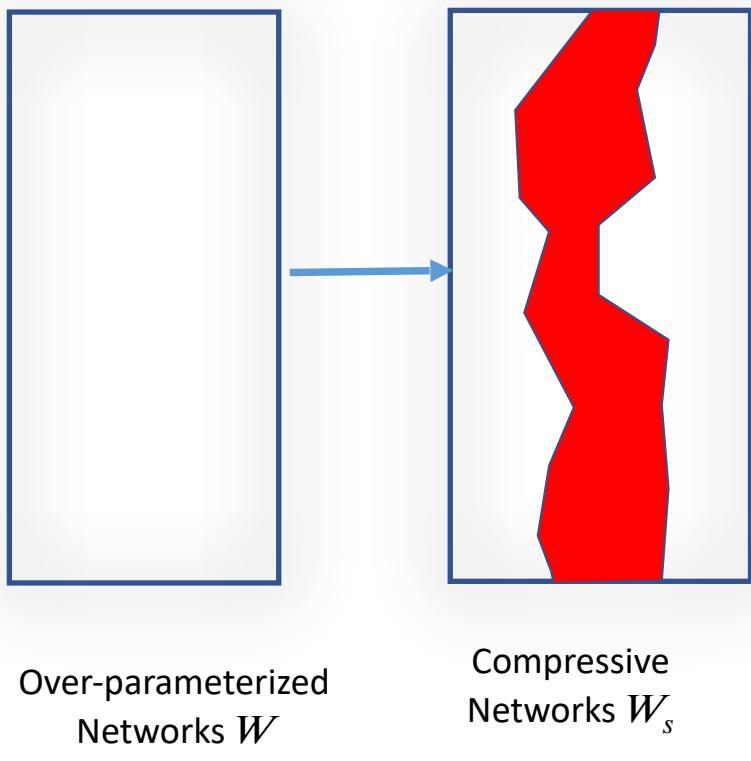


(c) Texture-shape cue conflict
63.9% **Indian elephant**
26.4% indri
9.6% black swan

Geirhos et al. ICLR 2019

<https://videoken.com/embed/W2HvLBMhCJQ?tocitem=46>

Lottery Ticket Hypothesis for Efficient Subnets in Deep Learning

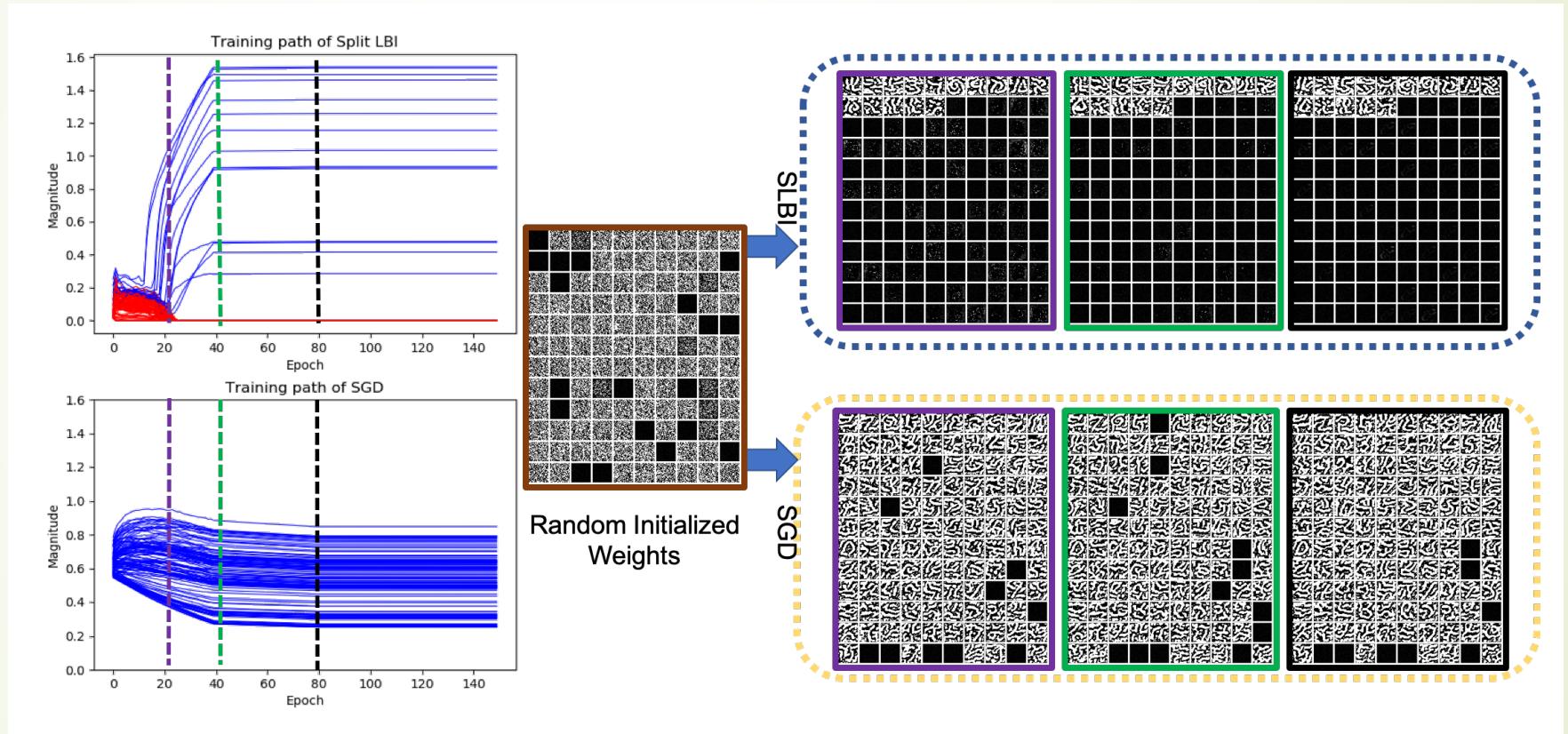


Lottery Ticket Hypothesis

- *Dense, randomly-initialized, feed-forward networks contain subnetworks (winning tickets) that – when trained in isolation – reach test accuracy comparable to the original network in a similar number of iterations. (Frankle & Carbin, 2019)*

Rewinding the network from the initialization, and find “winning ticket” subnet

Split LBI finds efficient sparse architecture



Yanwei Fu et al. TPAMI 45(2):1749-1765, 2023.

Yanwei Fu et al. DensiLBI, ICML 2020.

Texture bias in ImageNet training

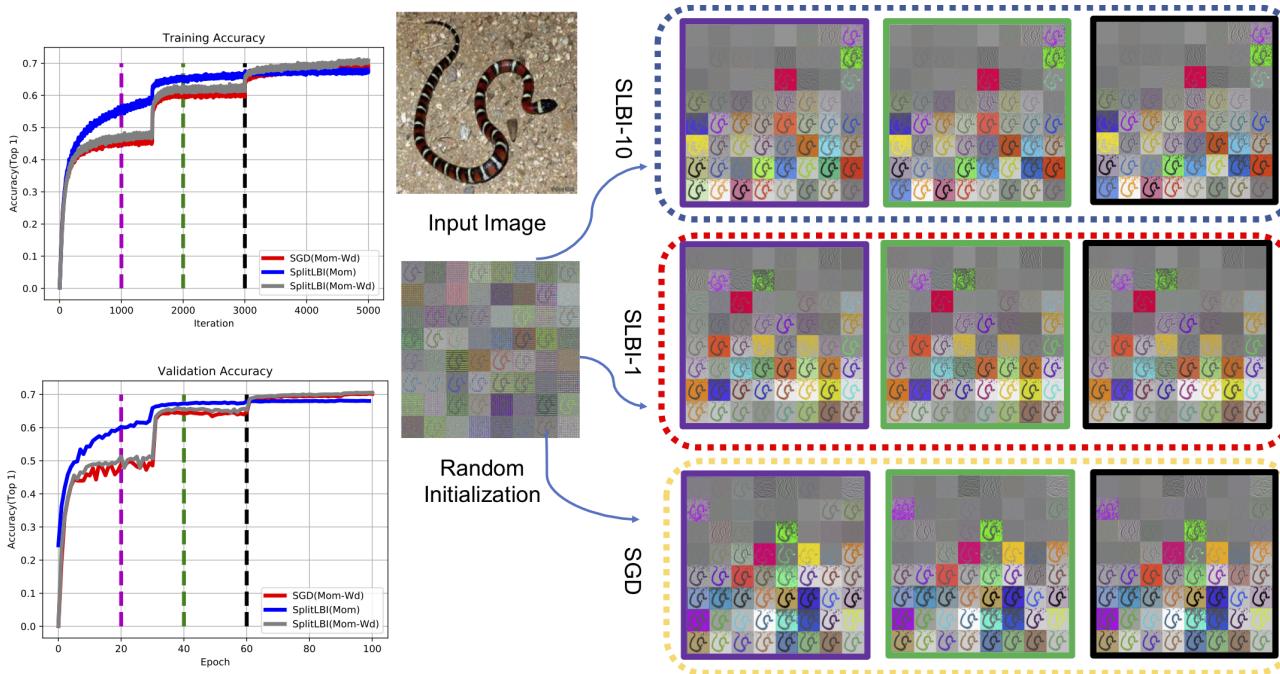
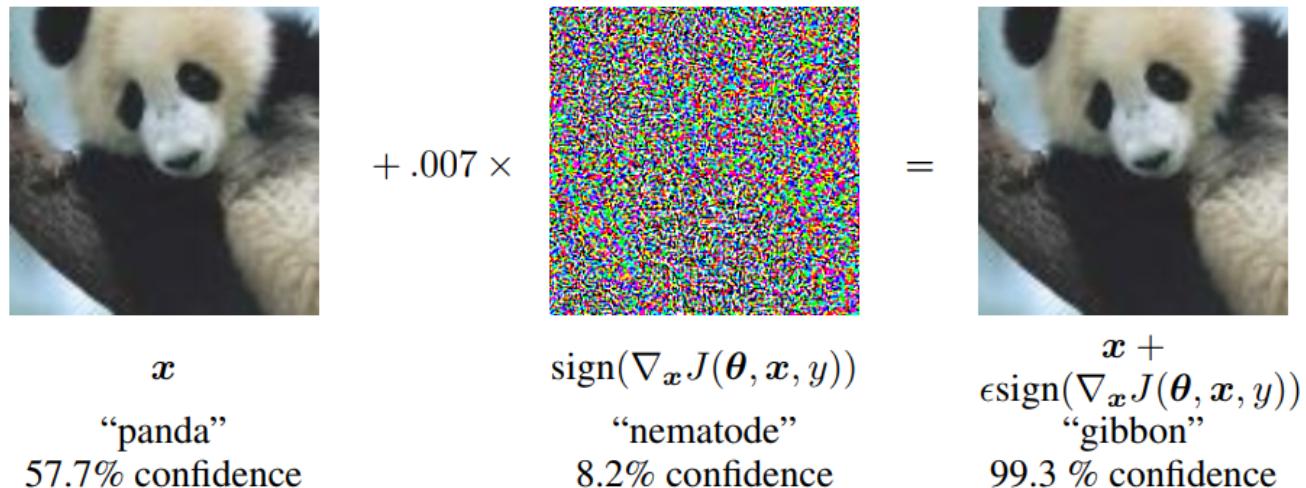


Figure: Visualization of the first convolutional layer filters of ResNet-18 trained on ImageNet-2012, where **texture** features are more important than **colour/shapes**. Given the input image and initial weights visualized in the middle, filter response gradients at 20 (purple), 40 (green), and 60 (black) epochs are visualized. SGD with Momentum (Mom) and Weight Decay (WD), is compared with SLBI.



Adversarial Examples and Robustness

Deep Learning may be fragile: adversarial examples



[Goodfellow et al., 2014]

- Small but malicious perturbations can result in severe misclassification
- Malicious examples generalize across different architectures
- What is source of instability?
- Can we robustify network?

Adversarial Examples: Fooling Images

- ▶ Start from an arbitrary image
- ▶ Pick an arbitrary class
- ▶ Modify the image to maximize the class
- ▶ Repeat until network is fooled

Fooling Images/Adversarial Examples

African elephant



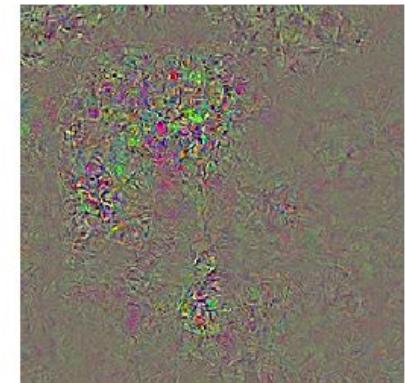
koala



Difference



10x Difference



schooner



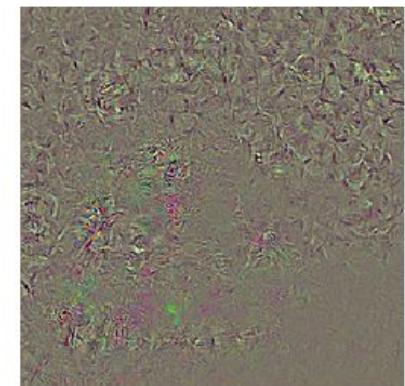
iPod



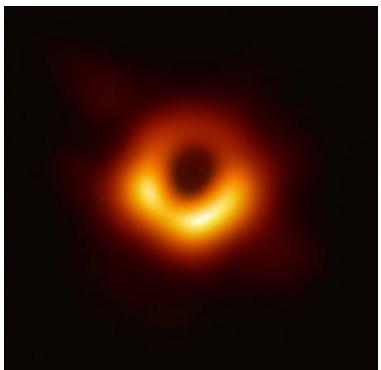
Difference



10x Difference



Convolutional Networks lack Robustness

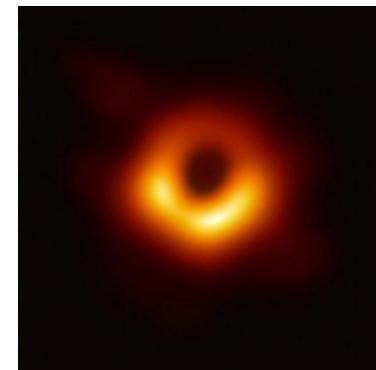


“black hole”
87.7% confidence

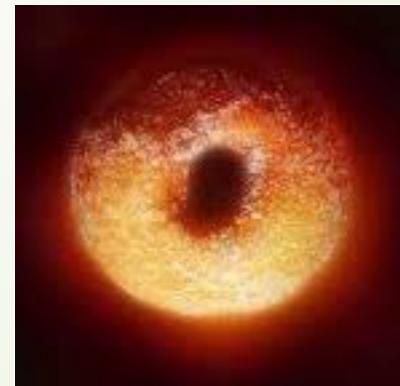
+ .007 ×



=

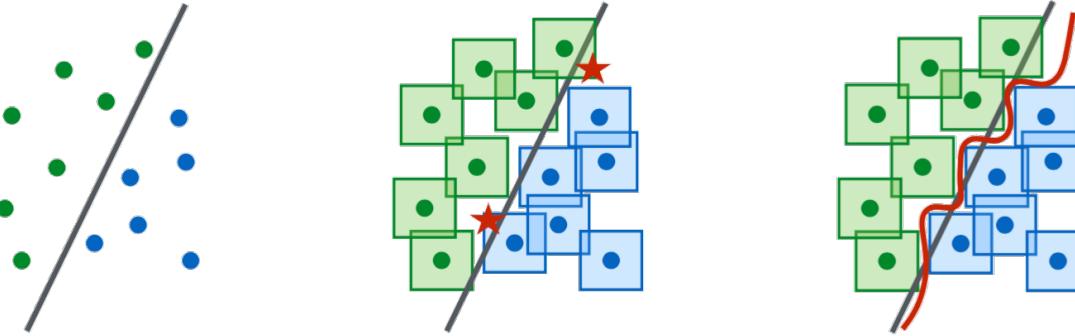


“donut”
99.3% confidence



Courtesy of Dr. Hongyang ZHANG.

Adversarial Robust Training



- Traditional training:

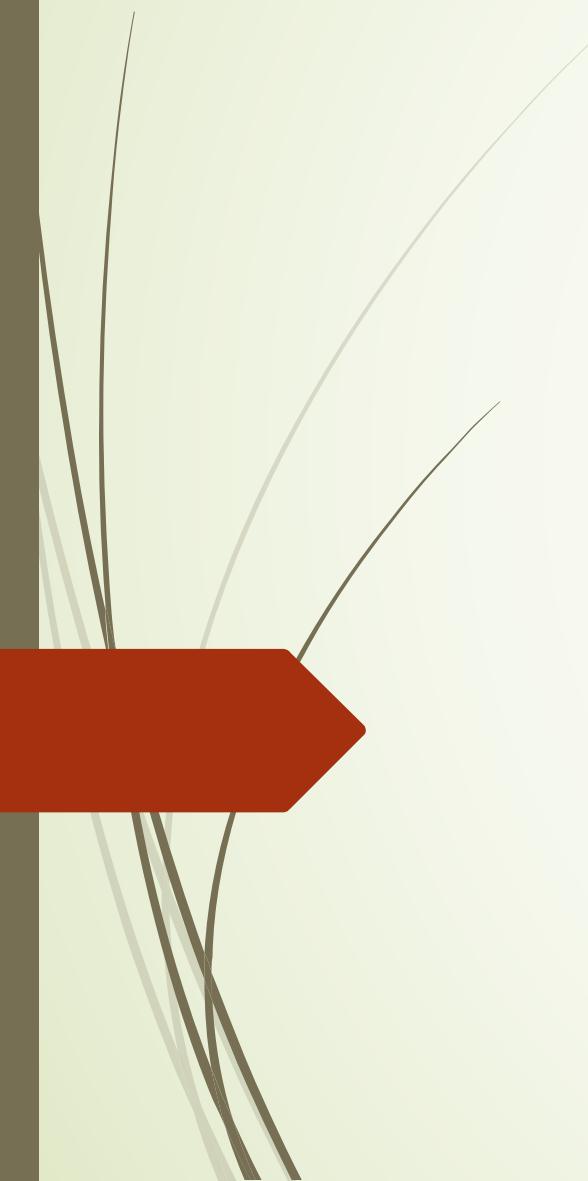
$$\min_{\theta} J_n(\theta, \mathbf{z} = (x_i, y_i)_{i=1}^n)$$

- e.g. square or cross-entropy loss as negative log-likelihood of logit models

- Robust optimization (Madry et al. ICLR'2018):

$$\min_{\theta} \max_{\|\epsilon_i\| \leq \delta} J_n(\theta, \mathbf{z} = (x_i + \epsilon_i, y_i)_{i=1}^n)$$

- robust to any distributions, yet computationally hard



Introduction to “(Re-)Imag(in)ing Price Trends”

By

Jingwen Jiang
University of Chicago

Bryan Kelly
Yale University, AQR Capital Management, and NBER

Dacheng Xiu
University of Chicago Booth School of Business

Brief Intro

In the empirical designs, they first embed 1D time-series data into 2D images depicting price and volumes.

Then they feed each training sample into CNN to estimate the probability of a positive subsequent return over short(5-day), medium(20-day), and long (60day) horizons.

Afterward, they use CNN-based out-of-sample predictions as signals in several asset pricing analyses.

Finally, they attempt to interpret the predictive patterns identified by the CNN.



Replication task

Mainly focus on:

- ▶ Data Preparation
- ▶ Model Design
- ▶ Workflow Design
- ▶ Performance Evaluation
- ▶ Interpretation

Data

- ▶ The sample runs from 1993-2019 shows daily opening, high, low prices. The original paper constructs datasets consisting of three scales of horizons (5-day, 20-day, 60-day). Here we just collect the **20-day** version. The total size of data is 8.6G.
- ▶ We already transferred the OHLC charts into images following the same procedures. Current images have the same resolution (64 * 60) and added with moving average lines(MA) and volume bars(VB).



Data

- Images labels take value **1** for positive returns ('up') and **0** for non-positive returns ('down'). In addition, we use **2** to mark the NaN value.

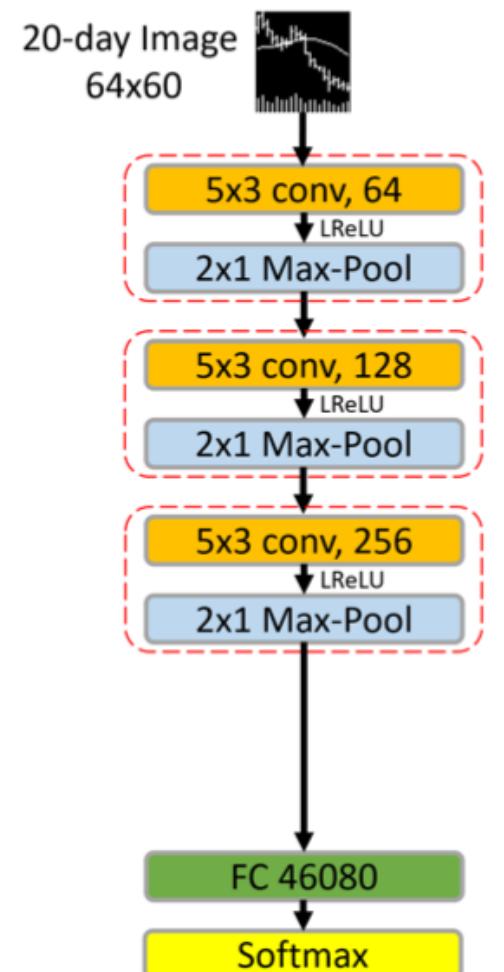
	Date	StockID	EWMA_vol	Retx	Retx_5d	Retx_20d	Retx_60d	Retx_week	Retx_month	Retx_quarter	...	Ret_week	Ret_month	Ret_quarter
0	2017-01-31	10001	0.000450	0.000000	4.370390e-07	-0.000002	-0.011860	NaN	-0.000002	NaN	...	NaN	-0.000002	NaN
1	2017-02-28	10001	0.000180	-0.003937	3.951997e-03	-0.003162	0.003953	NaN	0.003953	NaN	...	NaN	0.009953	NaN
2	2017-03-31	10001	0.000064	0.007936	-7.874612e-03	-0.015749	0.015748	-0.007875	-0.015749	0.017717	...	-0.007875	-0.015749	0.023704
3	2017-04-28	10001	0.000030	0.000000	9.999881e-03	0.016001	0.032002	0.010000	0.016001	NaN	...	0.010000	0.016001	NaN
4	2017-05-31	10001	0.000015	0.000000	4.370390e-07	0.015748	NaN	NaN	0.017717	NaN	...	NaN	0.023703	NaN

Data: Label Format

Retx_20d: < 0
Retx_20d_label: 0

Date	2017-01-31 00:00:00
StockID	10001
EWMA_vol	0.00045
Retx	0.0
Retx_5d	0.0
Retx_20d	-0.00002
Retx_60d	-0.01186
Retx_week	NaN
Retx_month	-0.00002
Retx_quarter	NaN
Retx_tstat	0.0
Retx_5d_tstat	0.00097
Retx_20d_tstat	-0.003558
Retx_60d_tstat	-26.333479
MarketCap	133078.0
Retx_label	0
Retx_5d_label	1
Retx_20d_label	0
Retx_60d_label	0
window_size	20
next_month_ret_0delay	-0.00002
next_month_ret_1delay	-0.00002
Ret	0.0
log_ret	0.0
cum_log_ret	2.075332
Ret_week	NaN
Ret_month	-0.00002
Ret_quarter	NaN
Ret_5d	0.0
Ret_20d	-0.00002
Ret_60d	-0.005954
Ret_65d	0.001998
Ret_180d	NaN
Ret_250d	NaN
Ret_260d	NaN

CNN Architecture Design



- ▶ A core building block consists of three operations:
 - ▶ convolution
 - ▶ activation
 - ▶ pooling
- ▶ In the paper, for 20-day images, they build a baseline CNN architecture with 3 conv blocks and connected with a fully connected layer as a classifier head.
- ▶ See the original paper for details (including the selection of the size of the convolution kernel, the selection of the convolution method, the design of the pooling layer and the selection of the activation function, etc.)

Data Split: Training, Validation, Testing

- ▶ Consider dividing the entire sample into **training**, **validation** and **testing** samples.
- ▶ In the original paper, they use **the first seven-year sample** (1993-1999) to **train** and **validate** the model, in which **70%** of the sample are randomly selected for training and the remaining **30%** for validation. The **remaining twenty years** of data comprise the out-of-sample test dataset.

Loss and evaluation

- ▶ You can simply treat the prediction analysis as a classification problem. Use Cross Entropy Loss

$$L_{CE}(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

- ▶ To measure the classification accuracy, a true positive (TP) or true negative (TN) occurs when a predicted “up” probability of greater than 50% coincides with a positive realized return and a probability less than 50% coincides with a negative return. False positives and negatives (FP and FN) are the complementary outcomes.

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

- ▶ For more evaluation metrics or methods, like **Sharpe Ratio**, please refer to the original paper.

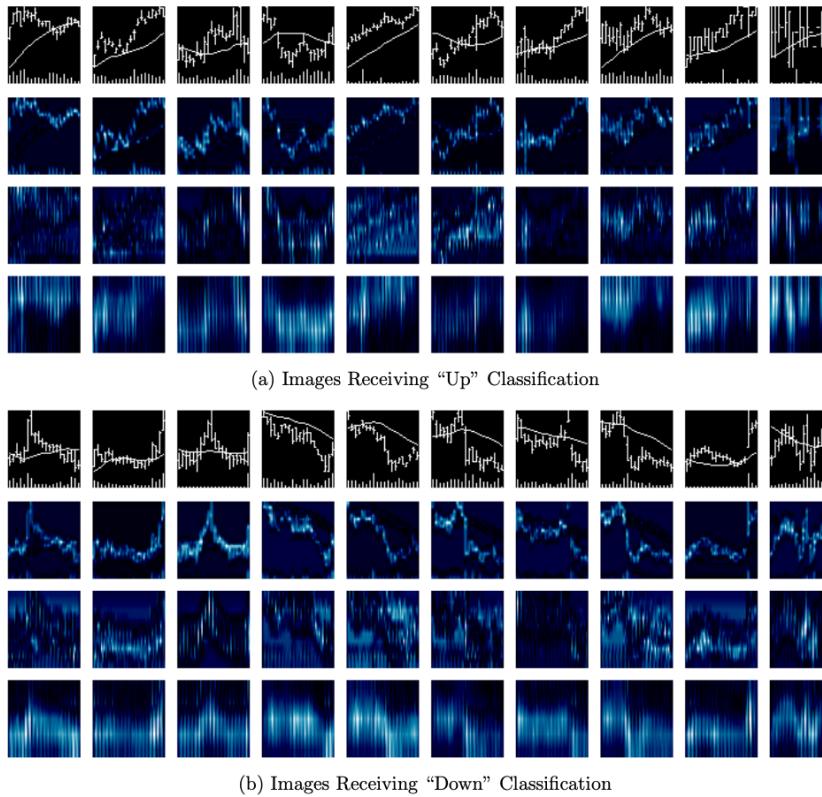
Training Process

- ▶ The author adopts several ways to combat over-fitting issue and aid efficient computation.
- ▶ They applied the Xavier initialization for weights in each layer, which guarantees faster convergence by scaling the initial weights.
- ▶ Other techniques like applying dropout, using batch normalization and early stopping also assists better performance.
- ❖ We recommend referring to the training details mentioned in the paper 3.3 when training the baseline model.

Exploring Interpretability

- ▶ Using a visualization method (Grad-CAM) to understand how different image examples activate the regions of the CNN to trigger 'up' or 'down' return predictions.

Figure 12: I20R20 Grad-CAM for 20 Images from 2019



Note: Brighter regions of the heatmap correspond to regions with the higher activation. For each panel, the first row is the original images followed by the grad-CAM for each layer in the CNN.

Extensions

► Ablation studies and test robustness,

- For example, you can perform the sensitivity analysis of the CNN prediction model to alternate choices in model architecture (e.g., varying the number of filters in each layer or varying the number of layers, like the paper shows in Table 18)

Table 18: Sensitivity to Model Structure and Estimation, I20R20

		Loss		Acc.		Correlation		Sharpe Ratio	
		V	T	V	T	Spearman	Pearson	EW	VW
Baseline		0.688	0.692	0.540	0.525	0.052	0.032	2.18	0.56
Filters (64)	32	0.688	0.691	0.541	0.526	0.053	0.032	1.91	0.43
	128	0.692	0.692	0.535	0.527	0.050	0.030	2.01	0.45
Layers (3)	2	0.688	0.692	0.540	0.525	0.047	0.028	1.61	0.19
	4	0.689	0.692	0.538	0.524	0.051	0.031	2.00	0.38
Dropout (0.50)	0.00	0.698	0.695	0.532	0.521	0.044	0.027	2.32	0.54
	0.25	0.691	0.693	0.536	0.525	0.050	0.030	2.11	0.56
	0.75	0.691	0.692	0.530	0.525	0.041	0.024	1.30	0.05
BN (yes)	no	0.685	0.691	0.549	0.528	0.059	0.037	2.44	0.58
Xavier (yes)	no	0.688	0.692	0.540	0.526	0.053	0.033	2.10	0.39
Activation (LReLU)	ReLU	0.689	0.693	0.538	0.520	0.052	0.031	1.71	0.35
Max Pool Size (2×1)	2×2	0.687	0.692	0.545	0.526	0.056	0.034	2.09	0.41
Filter Size (5×3)	3×3	0.689	0.693	0.538	0.522	0.050	0.028	1.39	0.28
	7×3	0.688	0.691	0.541	0.527	0.055	0.033	2.01	0.41
Dilation/Stride (2,1)/(3,1)	(2,1)/(1,1)	0.689	0.692	0.537	0.526	0.052	0.033	2.04	0.53
	(1,1)/(3,1)	0.689	0.692	0.538	0.526	0.049	0.030	1.66	0.17
	(1,1)/(1,1)	0.687	0.692	0.545	0.527	0.055	0.035	2.09	0.52



► What's more

- ❖ We encourage you not limited to simple binary classification tasks, since the label files we provided consist of more meaningful attributes, containing both categorical and numerical values.
 - For example, you can use the same 20-day horizon images to train your model to predict the return trend of different subsequent y-days even the detailed return values. (y can be 5, 20 even larger).
- ❖ Stock prices might be of high noise, so you may consider other financial markets, e.g. cryptocurrency market
 - For example, Kaggle Contest on G-Research Crypto Forecasting:
<https://www.kaggle.com/c/g-research-crypto-forecasting>

Thank you!

