images from *Geri's Game* (Pixar, 1997)

SIU2017

# Adversarial Training
## Attacks on Deep Networks and Generative Adversarial Networks

**Erkut Erdem** • **Aykut Erdem** • **Levent Karacan**
*Computer Vision Lab, Hacettepe University*

HACETTEPE UNIVERSITY COMPUTER VISION LAB

HACETTEPE UNIVERSITY

# Outline

- **Part 1:** Attacks on Deep Networks

- **Part 2:** Generative Adversarial Networks (GANs)

——— 10 Minutes Break

- **Part 3:** Image Editing with GANs

HACETTEPE UNIVERSITY COMPUTER VISION LAB

HACETTEPE UNIVERSITY

John Carpenter's The Thing (1982)

# Part 1 – Attacks on Deep Networks
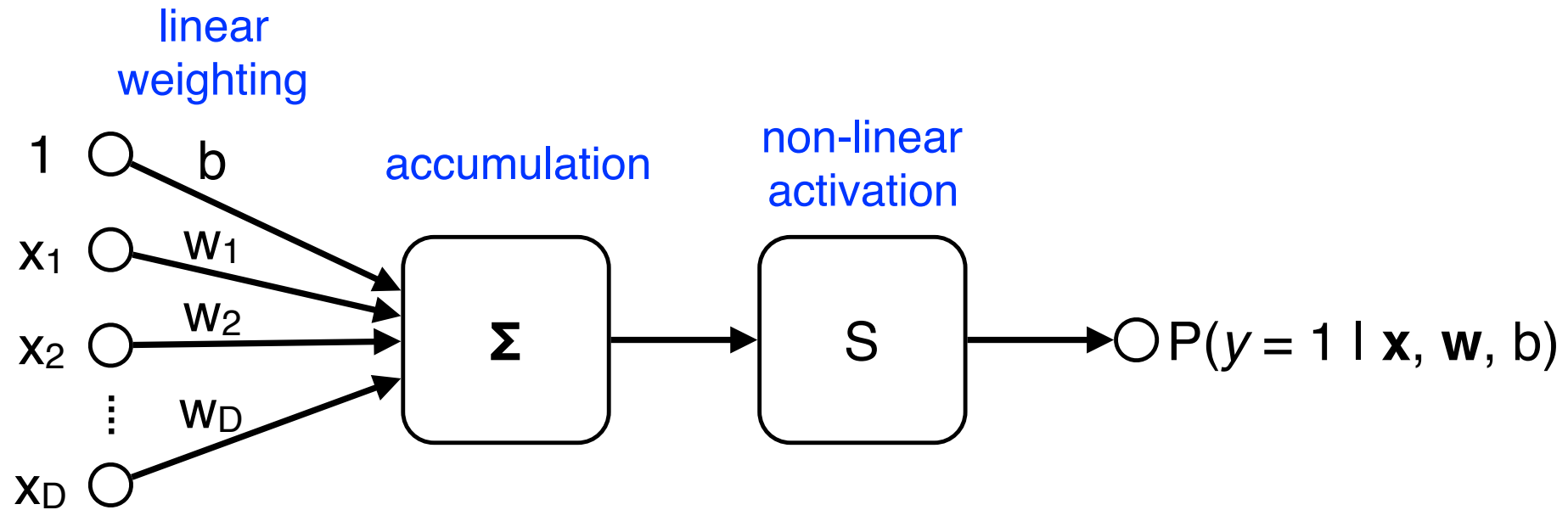
**Erkut Erdem**

*Computer Vision Lab, Hacettepe University*

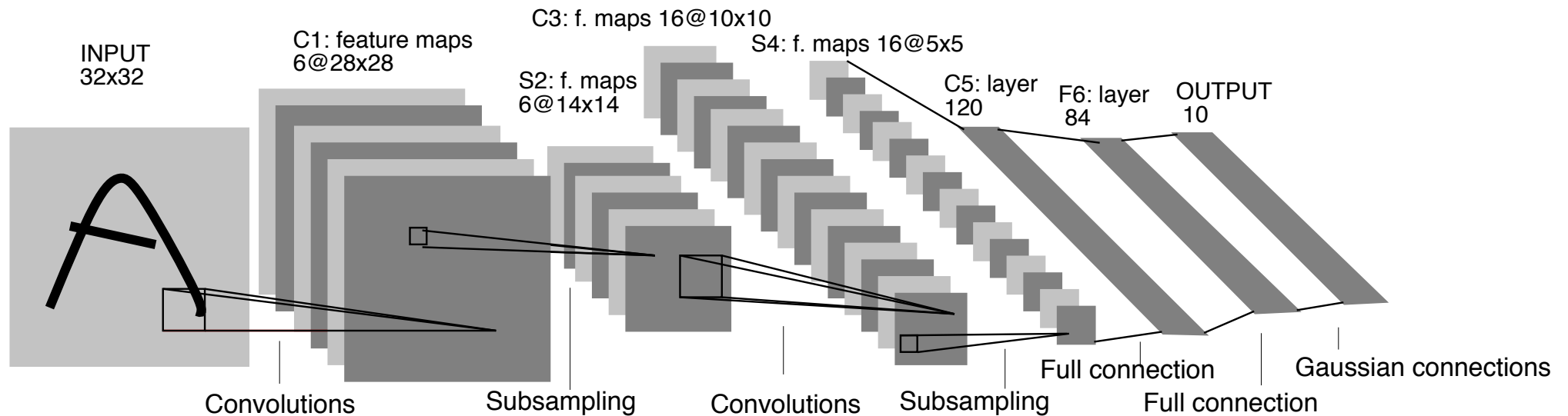# Deep Convolutional Networks in 10 mins

# 1st Era (1940's-1960's): Invention

- Connectionism (Hebb 1940's): complex behaviors arise from interconnected networks of simple units

- Artificial neurons (Hebb, McCulloch and Pitts 1940's-1950's)

- Perceptron (Rosenblatt 1950's): Single layer with learning rule

# 2nd Era (1980's-1990's): Multi-layered Networks

- Back-propagation (Rumelhart, Hinton and Williams 1986 +others): effective way to train multi-layered networks

- Convolutional networks (LeCun et al. 1989): architecture adapted for images (inspired by Hubel and Wiesel's simple/complex cells)

# The Deep Learning Era (2011-present)

- Big gains in performance on perceptual tasks:
  - Vision
  - Speech understanding
  - Natural language processing

- Three ingredients:
  1. Deep neural network models (supervised training)
  2. Big labeled datasets
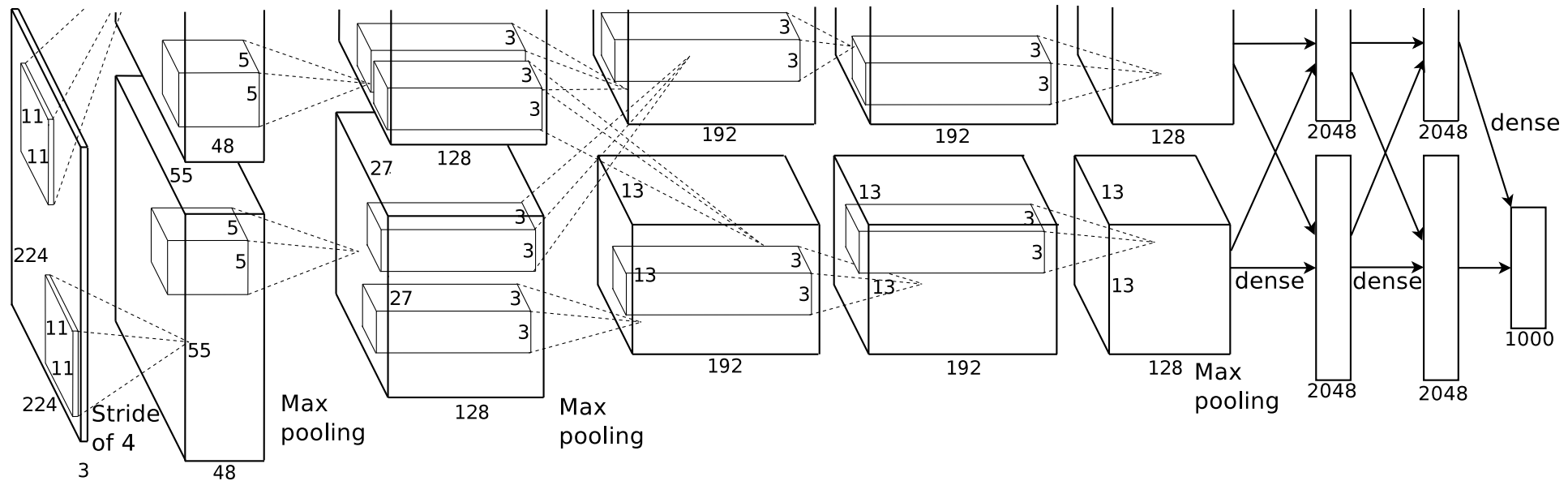  3. Fast GPU computation

# Powerful Hardware

- Deep neural nets highly amenable to implementation on Graphics Processing Units (GPUs)
  - Matrix multiplication
  - 2D convolution

- Latest generation nVidia GPUs (Pascal) deliver 10 Tflops
  - Faster than fastest computer in the world in 2000
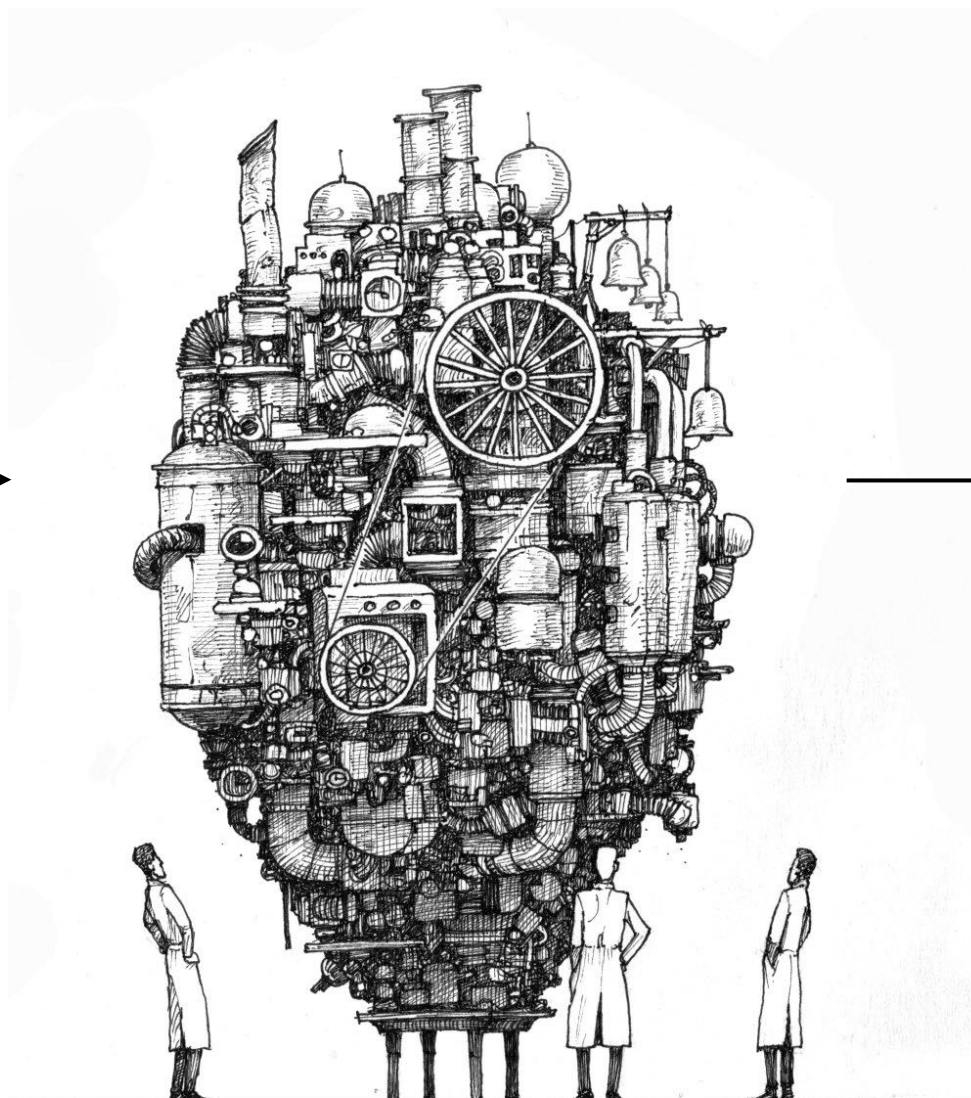  - 10 million times faster than 1980's Sun workstation

# AlexNet: The Model That Changed The History

- Krizhevsky, Sutskever and Hinton (2012)
  - 8 layer Convolutional network model [LeCun et al. 1989]
  - 7 hidden layers, 650,000 neurons, ~60,000,000 parameters
  - Trained on 1.2 million ImageNet images (with labels)
  - GPU implementation (50x speedup over CPU)
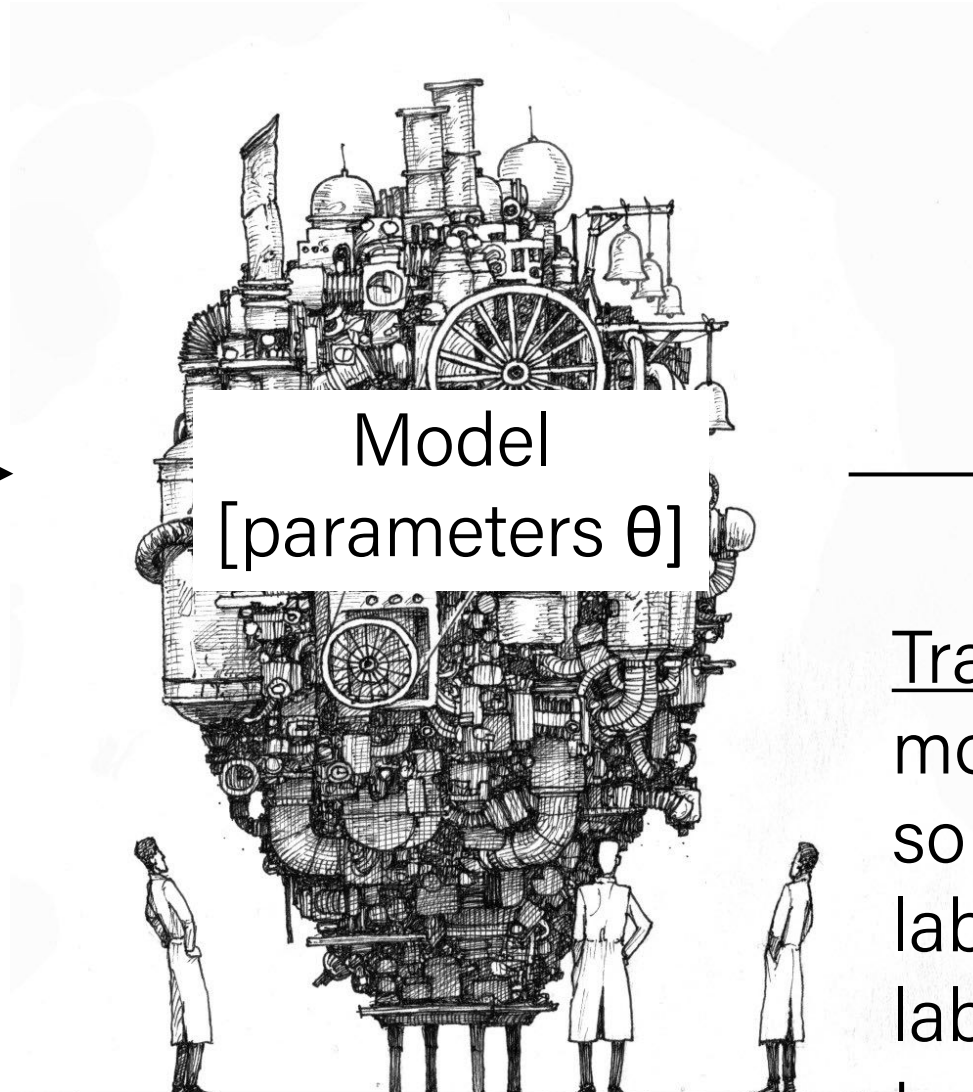  - Training time: 1 week on pair of GPUs

# Supervised Learning: Image Classification



"Cat"

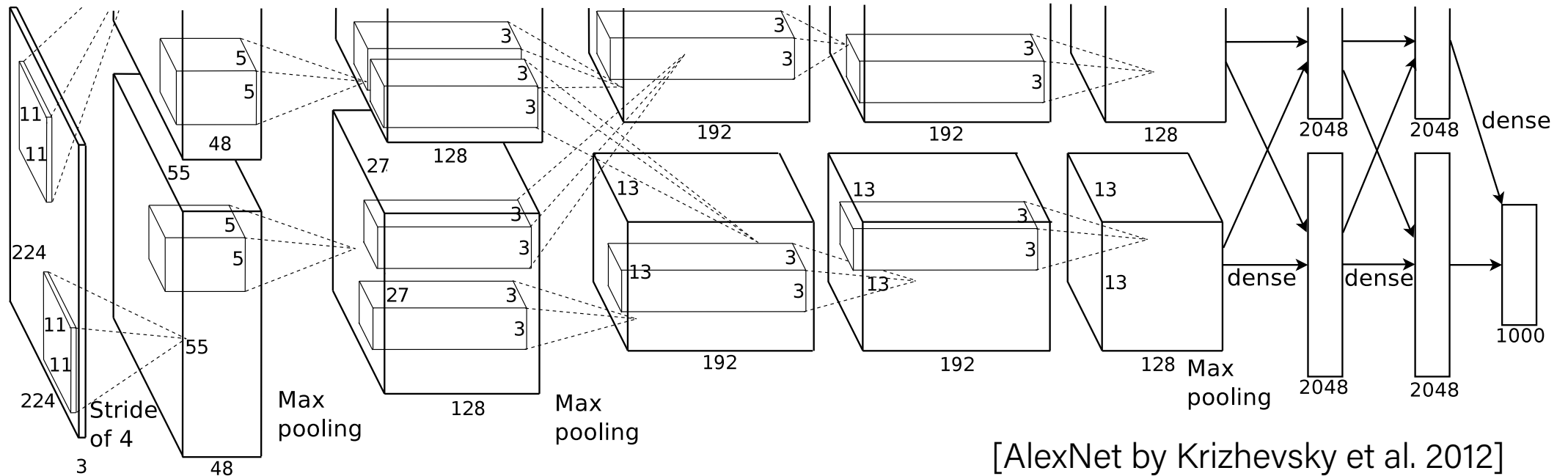Joshua Drewe

# Supervised Learning: Image Classification



Model
[parameters θ]

"Cat"

Training: Adjust model parameters θ so predicted labels match true labels across training set

Joshua Drewe

# Modern Convolutional Nets



[AlexNet by Krizhevsky et al. 2012]

Excellent **performance** in most image understanding tasks

Learn a sequence of **general-purpose representations**

Millions of parameters learned from data

The "**meaning**" of the representation is unclear
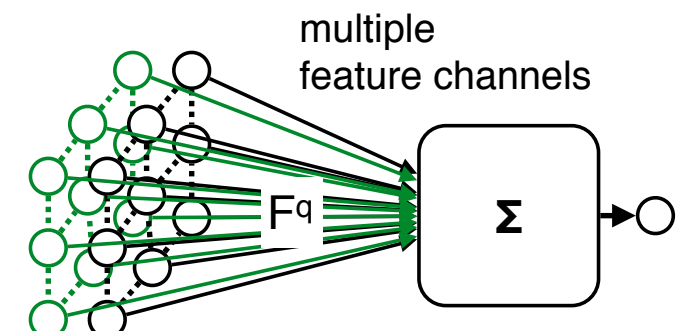
# Convolutions with Filters
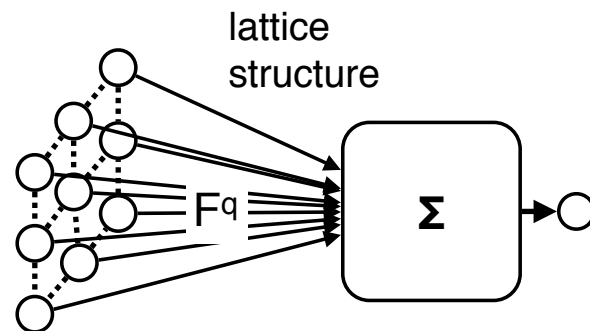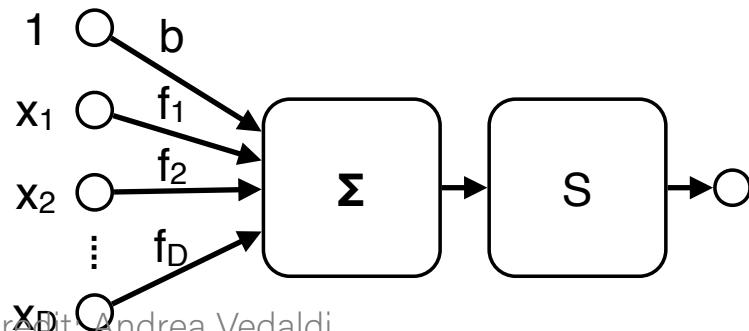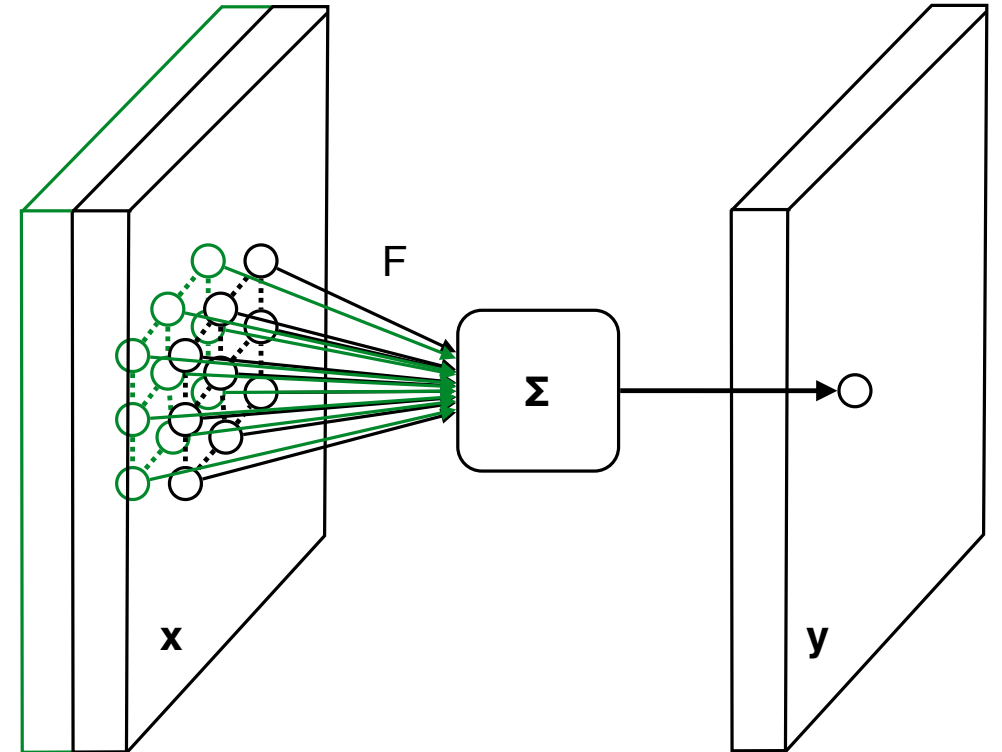
- Each filter acts on multiple input channel

  – **Convolution is local**
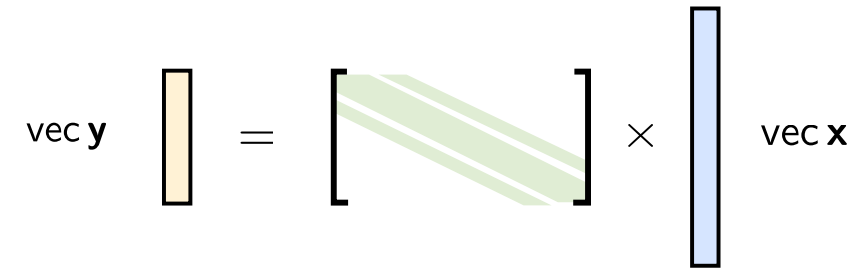  Filters look locally
  Parameter sharing

  – **Translation invariant**
  Filters act the same everywhere

# Convolution



- Convolution = Spatial filtering

$$(a \star b)[i, j] = \sum_{i', j'} a[i', j'] b[i - i', j - j']$$

**Convolution transpose**

**Transposed**

- Different filters (weights) reveal a different characteristics of the input.





vec **y** = × vec **x**

Transposed matrix

# Convolution



- Convolution = Spatial filtering

$$(a \star b)[i,j] = \sum_{i',j'} a[i',j']b[i-i',j-j']$$

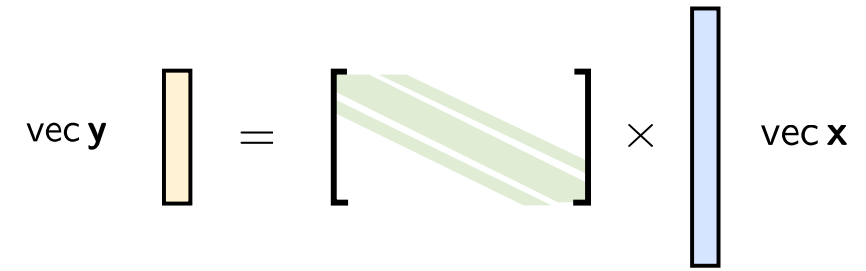- Different filters (weights) reveal a different characteristics of the input.



| 0 | -1 | 0 |
|---|----|---|
| -1 | 4 | -1 |
| 0 | -1 | 0 |

Transposed

vec **y** = × vec **x**

Transposed matrix

15

# Convolution

x    *    y

- Convolution = Spatial filtering    F

$$(a \star b)[i,j] = \sum_{i',j'} a[i',j']b[i-i',j-j']$$

vec **y**  =  [ ]  ×  [ ] vec **x**

Banded matrix equivalent to  F

- Different filters (weights) reveal a different characteristics of the input.

x    *ᵀ    y

F

| 1 | 0 | -1 |
|---|---|----|
| 2 | 0 | -2 |
| 1 | 0 | -1 |

*

vec **y**  =  ×  vec **x**

Transposed matrix
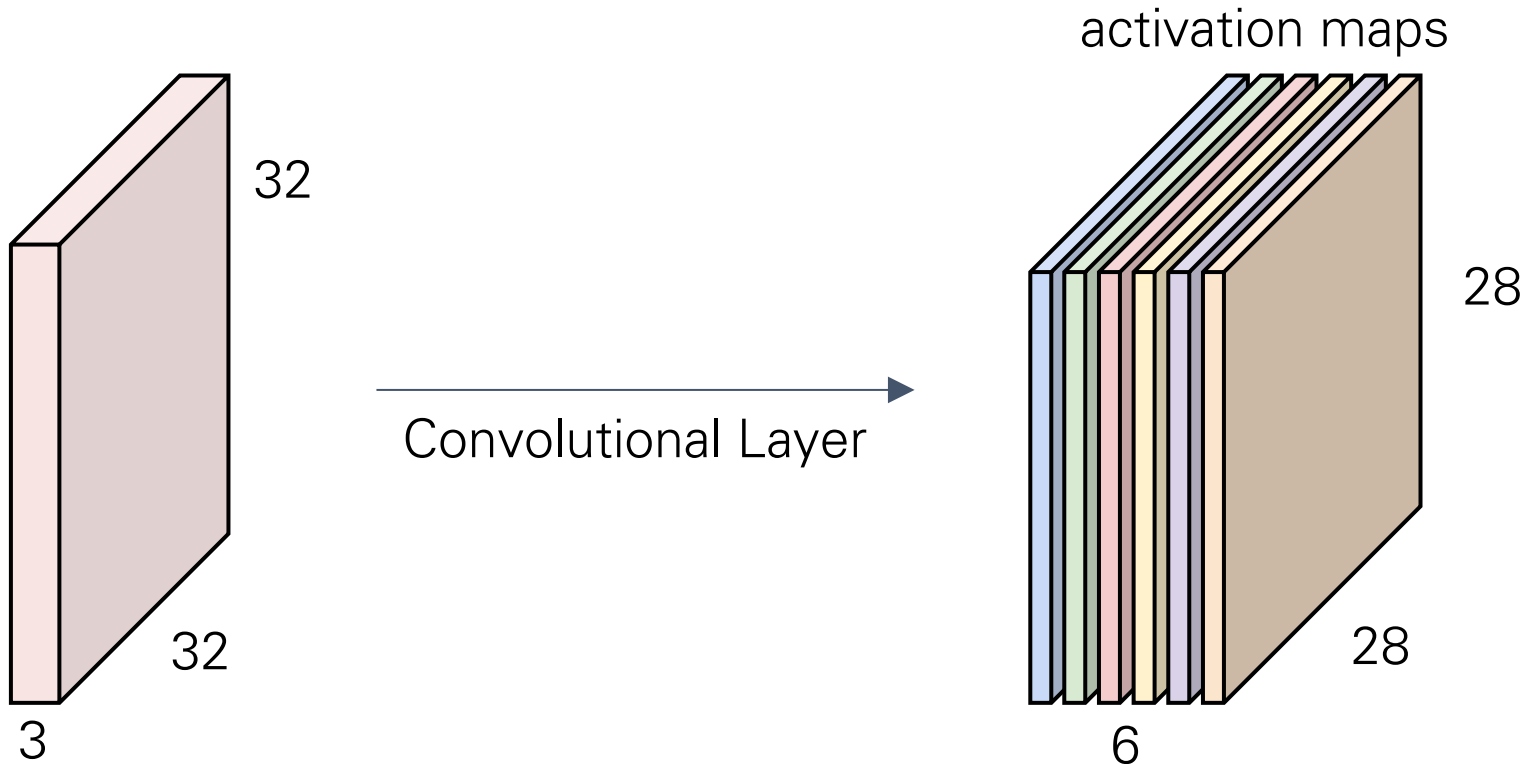
# Convolutional Layer
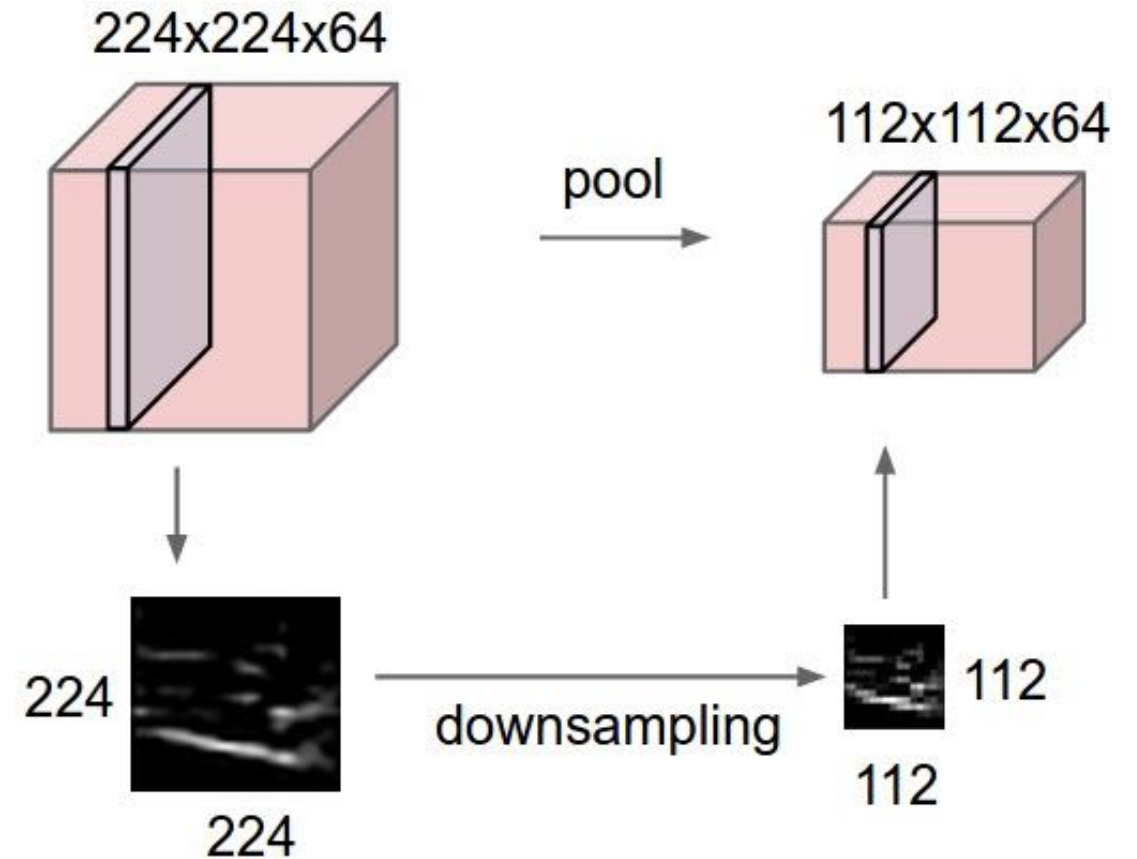
- Multiple filters produce multiple output channels
- For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

activation maps

32

32

3

Convolutional Layer

28

28

6

We stack these up to get an output of size 28x28x6.

17

# Pooling Layer

- makes the representations smaller and more manageable

- operates over each activation map independently:

- Max pooling, average pooling, etc.



224x224x64

pool →

112x112x64

224 × 224 → downsampling → 112 × 112

Single depth slice

x ↑

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

max pool with 2x2 filters and stride 2 →

| 6 | 8 |
|---|---|
| 3 | 4 |

→ y

18

# Fully Connected Layer

- contains neurons that connect to the entire input volume, as in ordinary Neural Networks

20

# Feature Learning

- Hierarchical layer structure allows to learn hierarchical filters (features).



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]
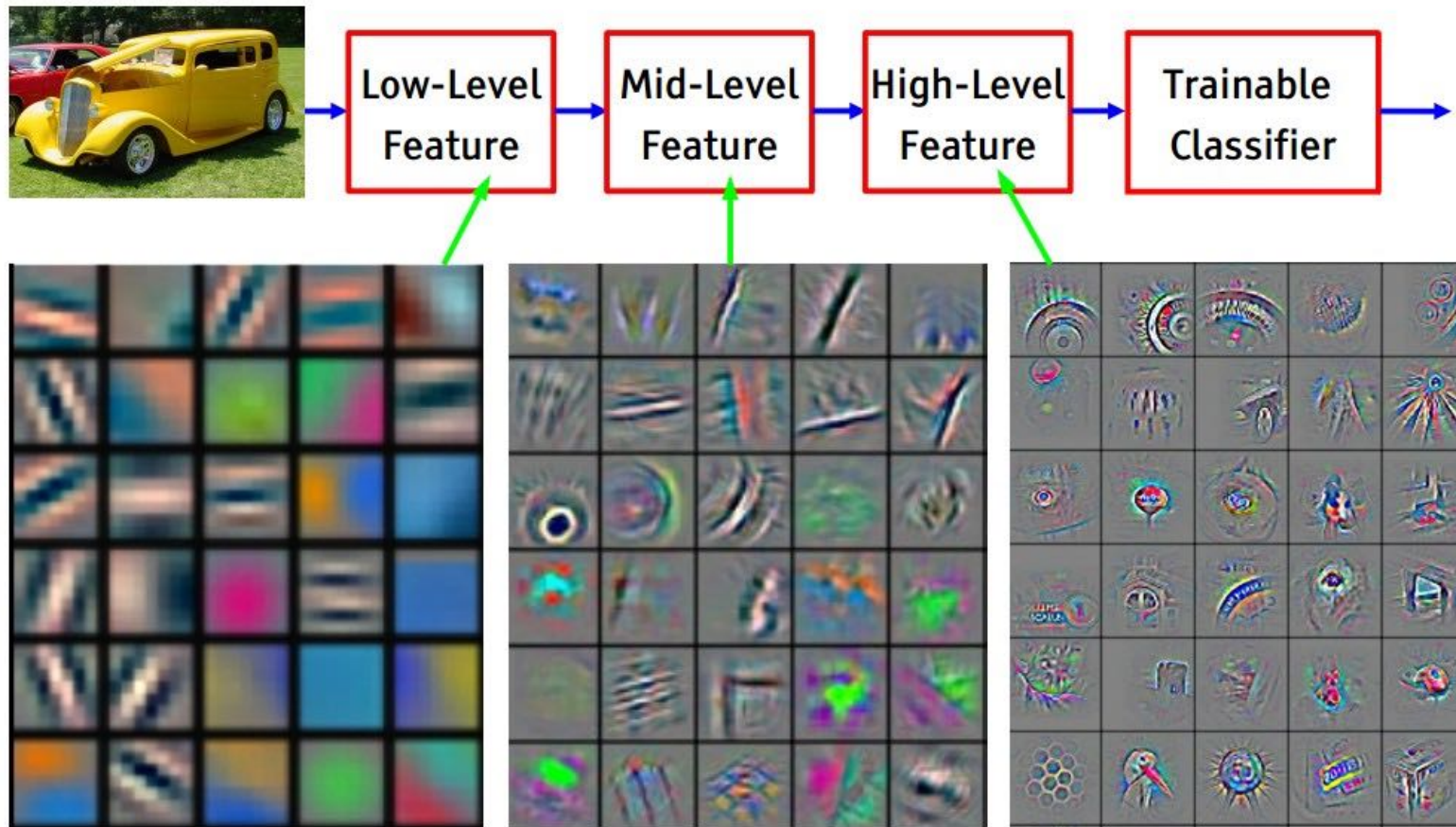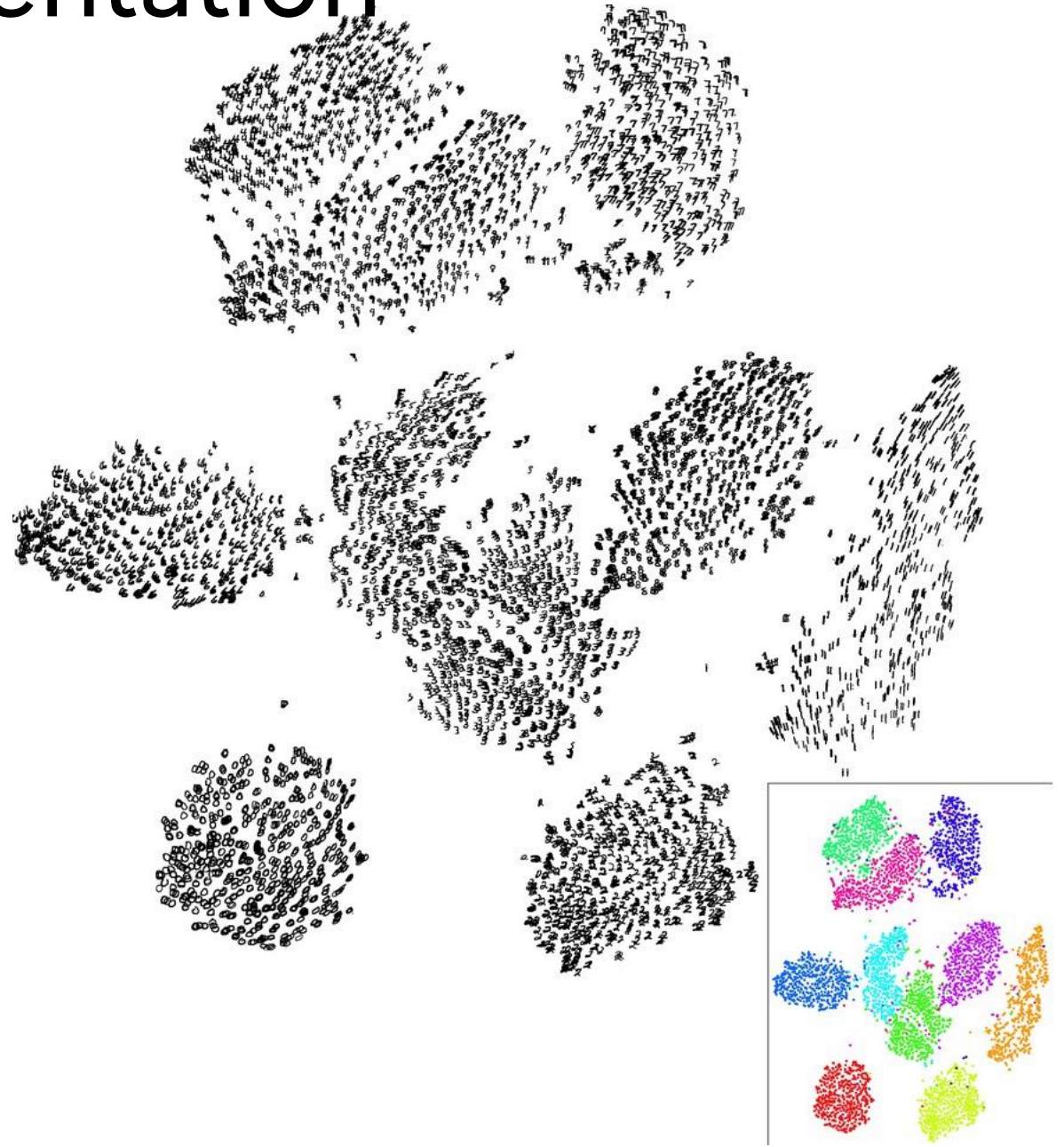
20

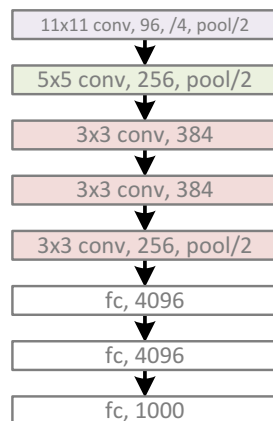# Visualizing The Representation

## t-SNE visualization
*(van der Maaten & Hinton)*

- Embed high-dimensional points so that locally, pairwise distances are conserved

- i.e. similar things end up in similar places. dissimilar things end up wherever

- **Right**: Example embedding of MNIST digits (0-9) in 2D

# Three Years of Progress

AlexNet, 8 layers

(ILSVRC 2012)

| 11x11 conv, 96, /4, pool/2 |
|---|
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

VGG, 19 layers

(ILSVRC 2014)

| 3x3 conv, 64 |
|---|
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

- Very deep
- Simply deep

GoogLeNet,
22 layers

(ILSVRC 2014)

- Branching
- Bottleneck
- Skip connection

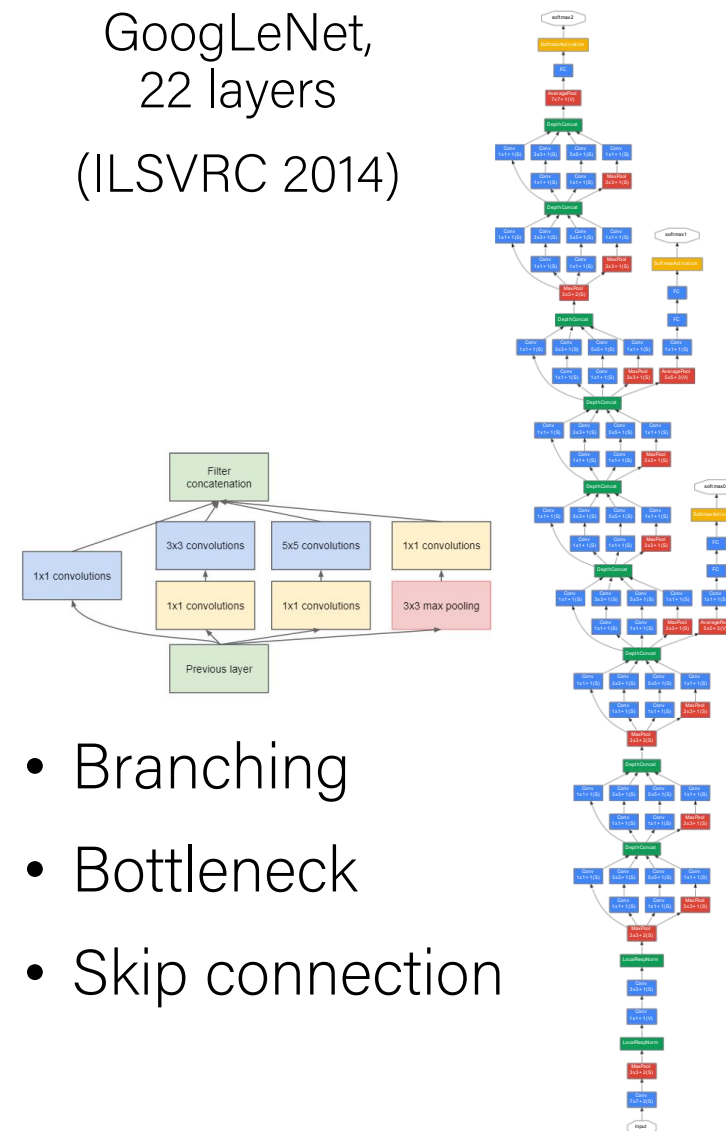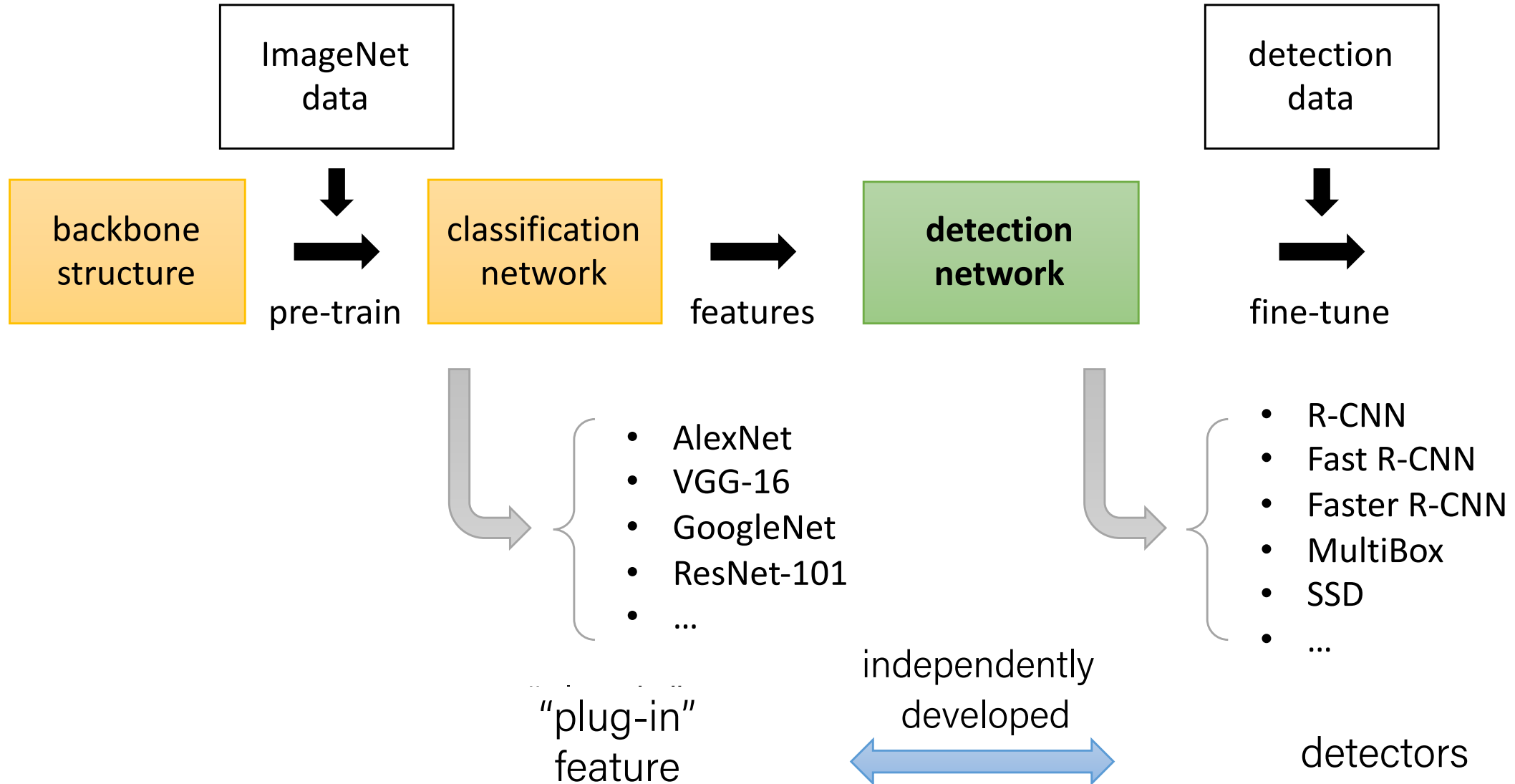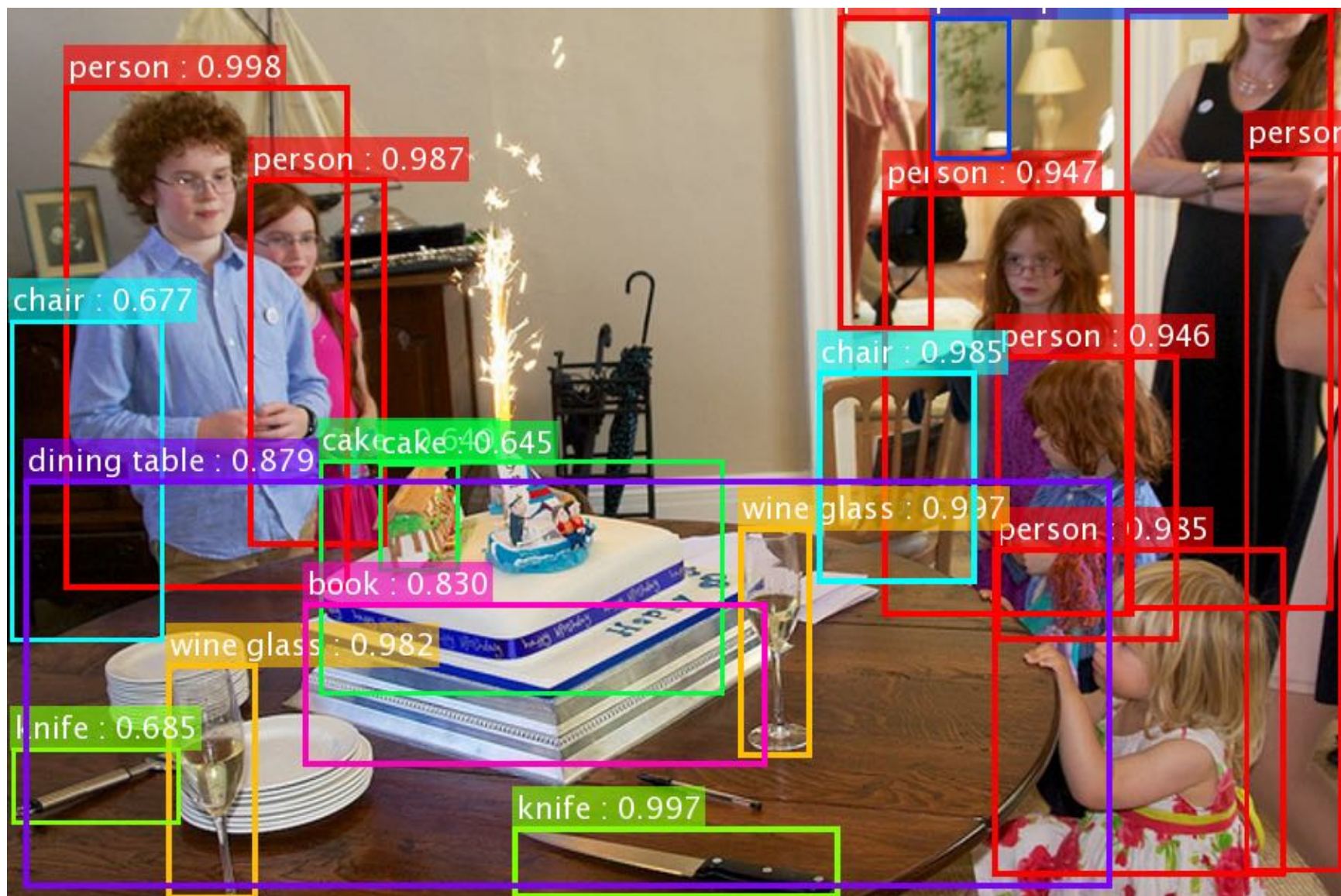# Training Deep Neural Networks

- The network is trained by stochastic gradient descent.
- Backpropagation is used similarly as in a fully connected network.
- Pass gradients through element-wise activation function.
- We also need to pass gradients through the convolution operation and the pooling operation.

# Object Detection Networks



| ImageNet data | | detection data |

backbone structure →(pre-train)→ classification network →(features)→ **detection network** →(fine-tune)→

classification network:
- AlexNet
- VGG-16
- GoogleNet
- ResNet-101
- …

"plug-in" feature

detection network:
- R-CNN
- Fast R-CNN
- Faster R-CNN
- MultiBox
- SSD
- …

detectors

independently developed

# ResNet's Object Detection Results on COCO



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. Deep Residual Learning for Image Recognition. CVPR 2016.

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. NIPS 2015.
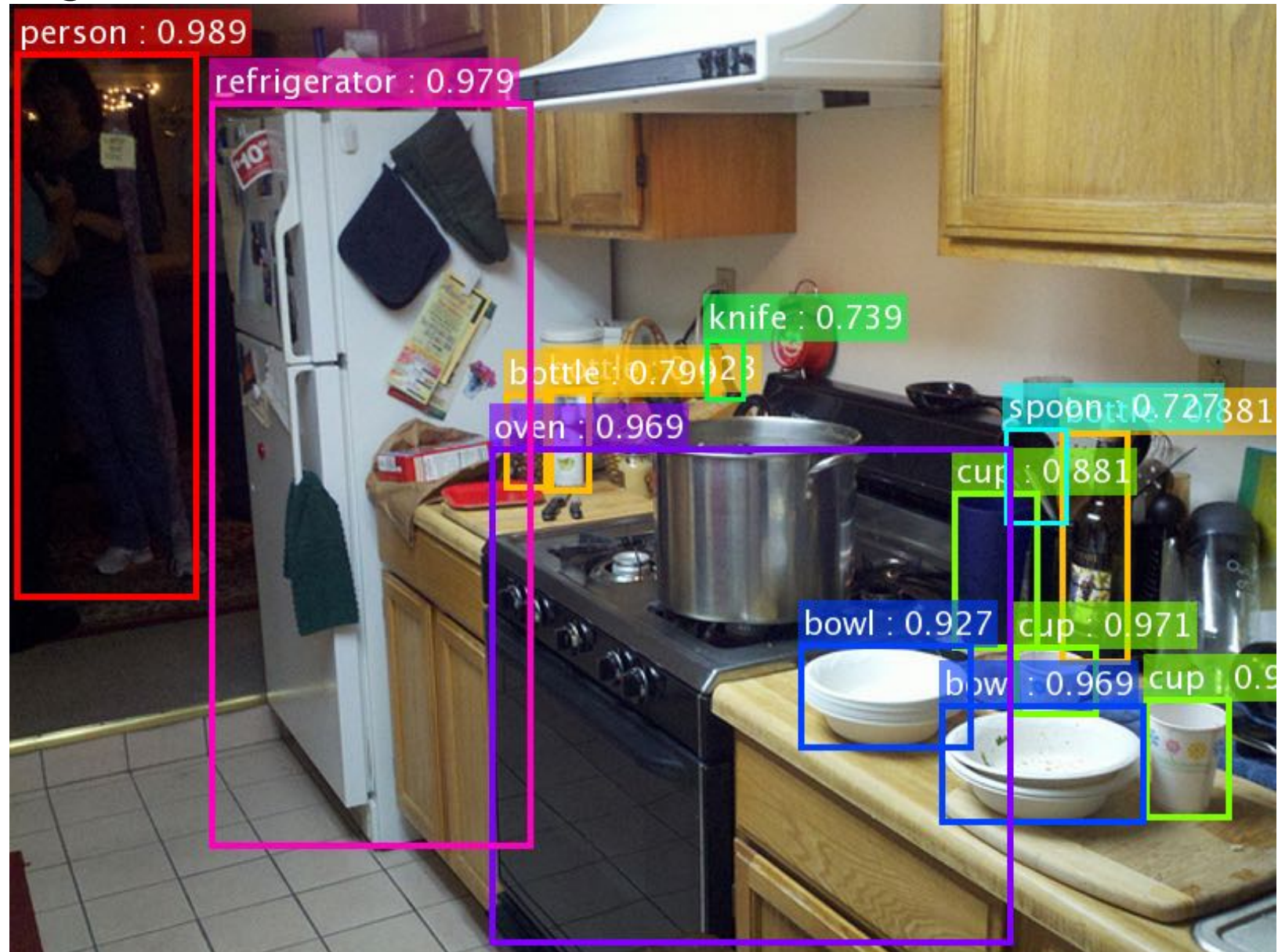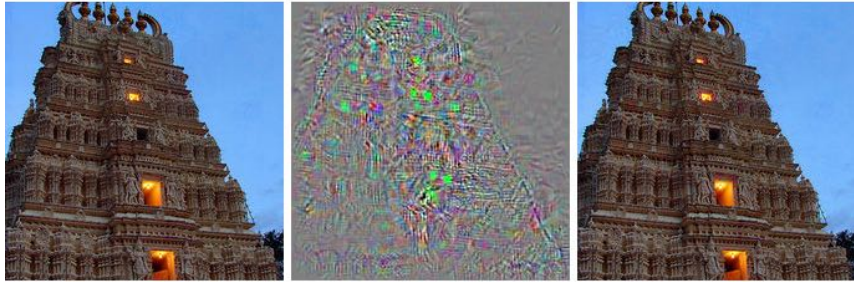
# ResNet's Object Detection Results on COCO



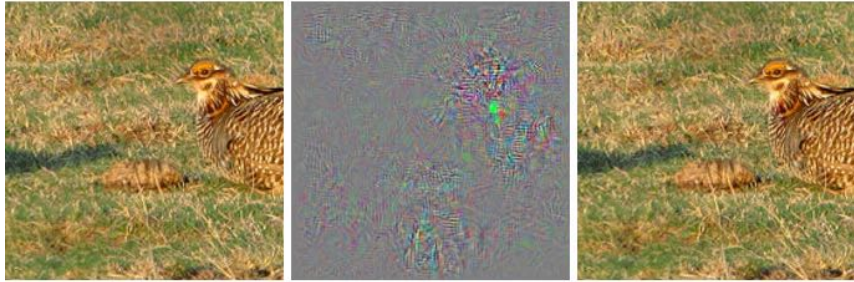Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. Deep Residual Learning for Image Recognition. CVPR 2016.

Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. NIPS 2015.

# Story isn't over yet!

# Story isn't over yet!
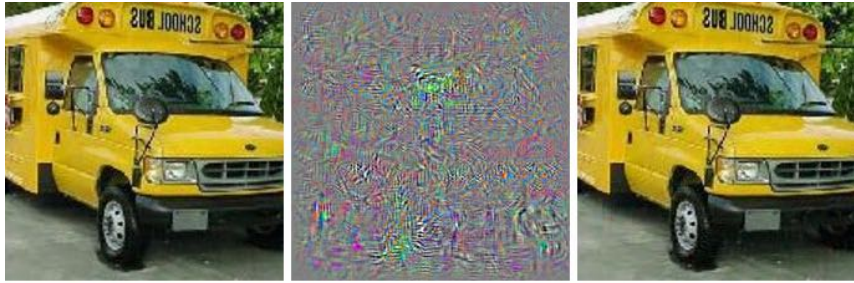
... we have reached
the point where ML works,
but let's see how it can be
easily fooled.

# Adversarial Examples

# Machine Learning System



Sample x

Machine
Learning
System
**f**

"Cat"

$f(x) = y_{true}$

Joshua Drewe

# Adversarial Examples



Adversarial example a
(indistinguishable from x)

Machine Learning System f

"Dog"

$f(a) \neq y_{true}$

Joshua Drewe

# Adversarial Examples in The Human Brain



These are concentric circles, not intertwined spirals.

(Pinna and Gregory, 2002)

# Adversarial Examples

- Adversarial examples pose potential security threats for practical machine learning systems.

- e.g., hypothetical attacks on autonomous vehicles

# Adversarial Examples

- Two types of adversaries (Papernot and Goodfellow 2016):

1. Poisoning training sets
   - interfere with the integrity of the training process
   - make modifications to existing training data, or insert additional data in the existing training set
   - increases the prediction error

2. Forcing models to make mistakes instantly with adversarial examples
   - perturb the inputs on which the model makes predictions (after training, during the inference phase)
   - generate "visually random" images that make a lot of sense to a machine learning system, but no sense at all to us

# Not just for neural nets

- Linear models

  - Logistic regression

  - Softmax regression

  - SVMs

- Decision trees

- Nearest neighbors

# Lets fool a binary linear classifier: (logistic regression)

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$



Since the probabilities of class 1 and 0 sum to one, the probability for class 0 is $P(y = 0 \mid x; w, b) = 1 - P(y = 1 \mid x; w, b)$. Hence, an example is classified as a positive example (y = 1) if $\sigma(w^T x + b) > 0.5$, or equivalently if the score $w^T x + b > 0$.

# Lets fool a binary linear classifier:

| x | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 | ← input example |
|---|---|----|---|----|---|---|---|----|---|---|------------------|
| w | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 | ← weights |

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **x** | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |
| **w** | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |

← input example

← weights

class 1 score = dot product:
= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3
=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474
i.e. the classifier is 95% certain that this is class 0 example.

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| x | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |

← input example

**W**

| -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |
|---|---|---|---|---|---|---|---|---|---|

← weights

adversarial **x**

| ? | ? | ? | ? | ? | ? | ? | ? | ? | ? |
|---|---|---|---|---|---|---|---|---|---|

class 1 score = dot product:
= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3
=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474
i.e. the classifier is 95% certain that this is class 0 example.

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **X** | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 | ← input example |
| **W** | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 | ← weights |
| adversarial **x** | 1.5 | -1.5 | 3.5 | -2.5 | 2.5 | 1.5 | 1.5 | -3.5 | 4.5 | 1.5 | |

class 1 score before:

-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3

=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474

-1.5+1.5+3.5+2.5+2.5-1.5+1.5-3.5-4.5+1.5 = 2

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

=> probability of class 1 is now 1/(1+e^(-(2))) = 0.88

**i.e. we improved the class 1 probability from 5% to 88%**

# Lets fool a binary linear classifier:

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |

**x** &larr; input example

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |

**w** &larr; weights

adversarial **x**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1.5 | -1.5 | 3.5 | -2.5 | 2.5 | 1.5 | 1.5 | -3.5 | 4.5 | 1.5 |

class 1 score before:
-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3
=> probability of class 1 is $1/(1+e^{-(-3)}) = 0.0474$

-1.5+1.5+3.5+2.5+2.5-1.5+1.5-3.5-4.5+1.5 = 2

=> probability of class 1 is now $1/(1+e^{-(2)}) = 0.88$
**i.e. we improved the class 1 probability from 5% to 88%**

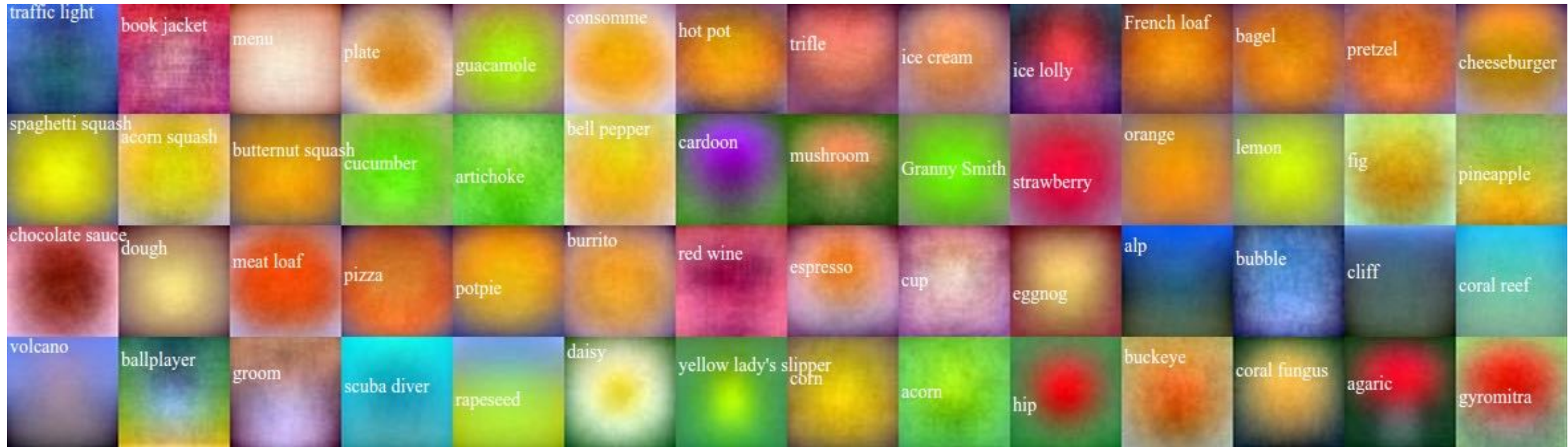This was only with 10 input dimensions. A 224x224 input image has 150,528.

(It's significantly easier with more numbers, need smaller nudge for each)

# Blog post: Breaking Linear Classifiers on ImageNet

Recall CIFAR-10 linear classifiers:
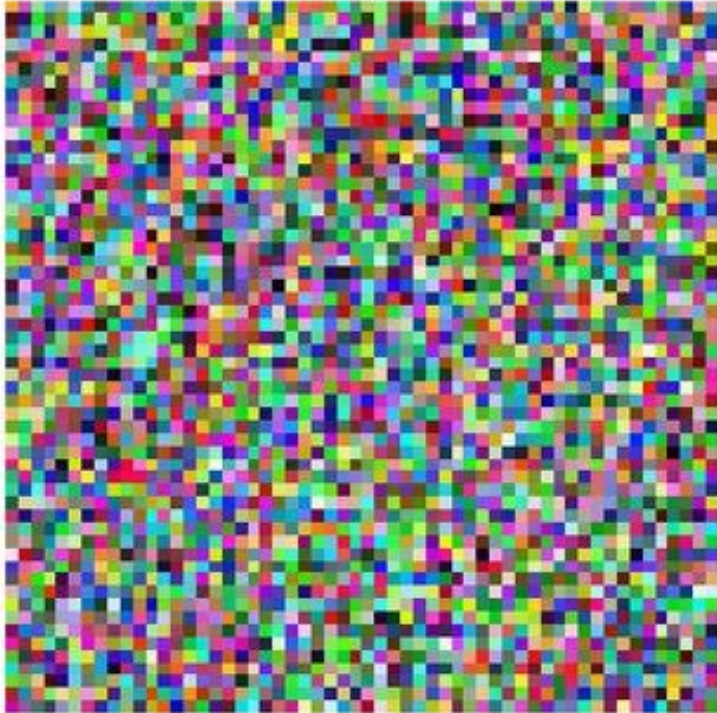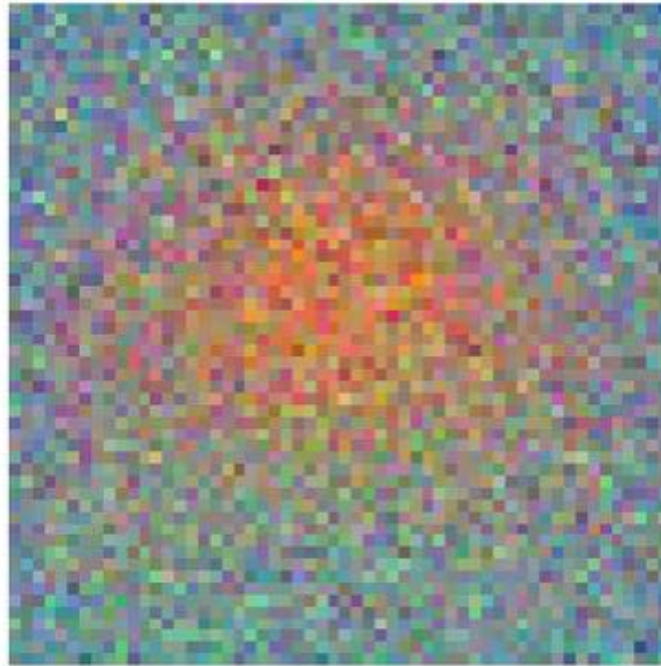


ImageNet classifiers:



http://karpathy.github.io/2015/03/30/breaking-convnets/

# Breaking Linear Classifiers on ImageNet



mix in a tiny bit of
Goldfish classifier weights

**100% Goldfish**

43

# Breaking Linear Classifiers on ImageNet

# Breaking Linear Classifiers on ImageNet

# Intriguing Properties of Neural Networks
(Szegedy et al., 2013)



correct  +distort  ostrich  correct  +distort  ostrich

Minimize $\|r\|_2$ subject to:
1. $f(x + r) = l$
2. $x + r \in [0, 1]^m$

f: classifier function
x: input image
r: distortion
l: target label

Minimize $c|r| + \text{loss}_f(x + r, l)$
subject to $x + r \in [0, 1]^m$

# Explaining and Harnessing Adversarial Examples
(Goodfellow et al., 2014)



"panda"
57.7% confidence

$+.007 \times$

"nematode"
8.2% confidence

$=$

"gibbon"
99.3 % confidence

# Explaining and Harnessing Adversarial Examples

(Goodfellow et al., 2014)



"panda"
57.7% confidence

$+ .007 \times$

"nematode"
8.2% confidence

$=$

"gibbon"
99.3 % confidence

$$X^{adv} = X + \epsilon \, \text{sign} \left( \nabla_X J(X, y_{true}) \right)$$

Score of label $y_{true}$, given input image X
*e.g. cross entropy loss*

# Explaining and Harnessing Adversarial Examples
(Goodfellow et al., 2014)

- Perturbation is computed to minimize a specific norm in the input domain while increasing the model's prediction error

$$J(\tilde{x}, \theta) \approx J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x).$$

Maximize

$$J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x)$$

subject to

$$||\tilde{x} - x||_\infty \leq \epsilon$$

$$\Rightarrow \tilde{x} = x + \epsilon \operatorname{sign}(\nabla_x J(x)).$$

The Fast Gradient Sign Method



- - - - - Task decision boundary
———— Model decision boundary

✖ Test point for class 1
✖ Adversarial example for class 1

✖ Training points for class 1
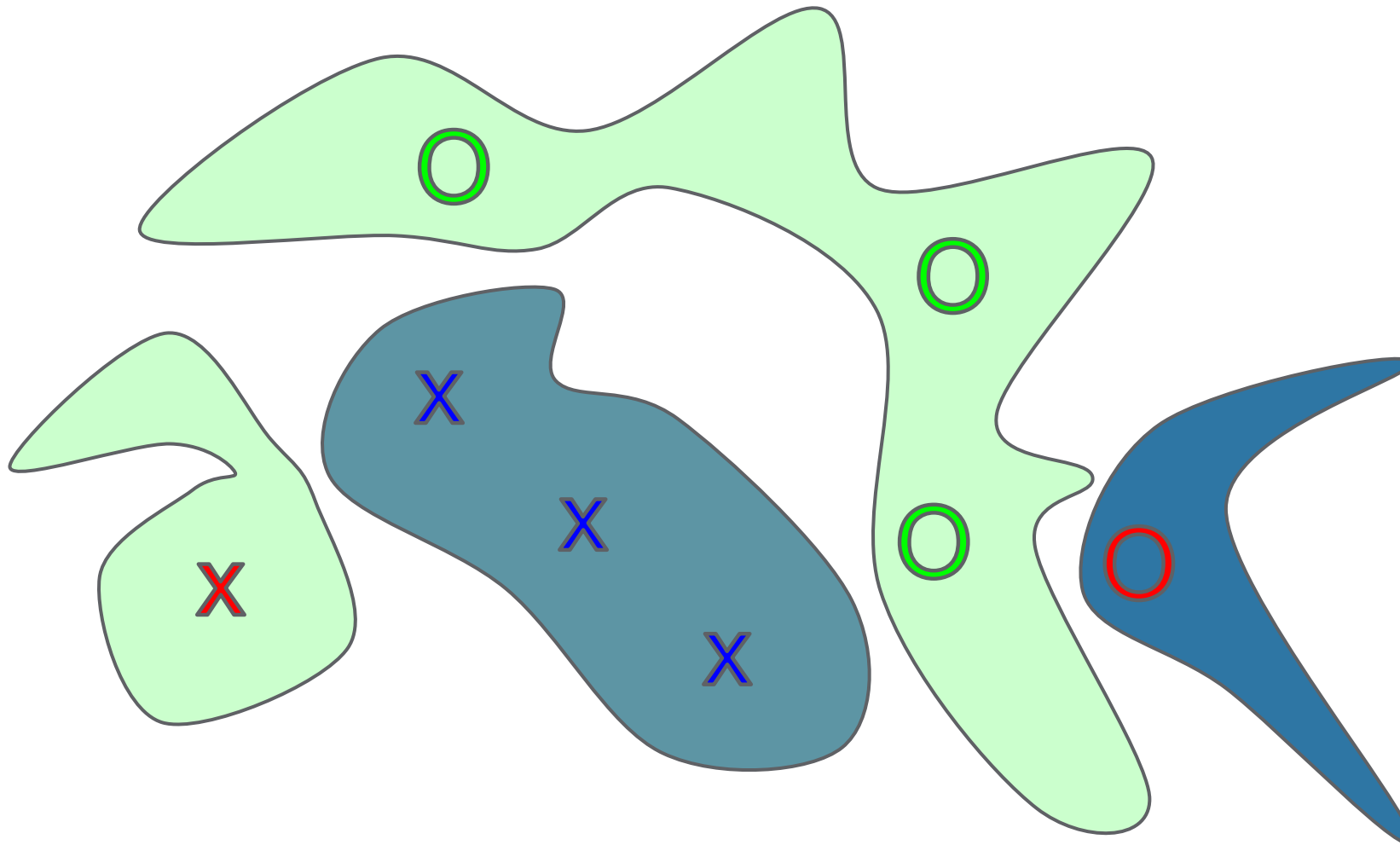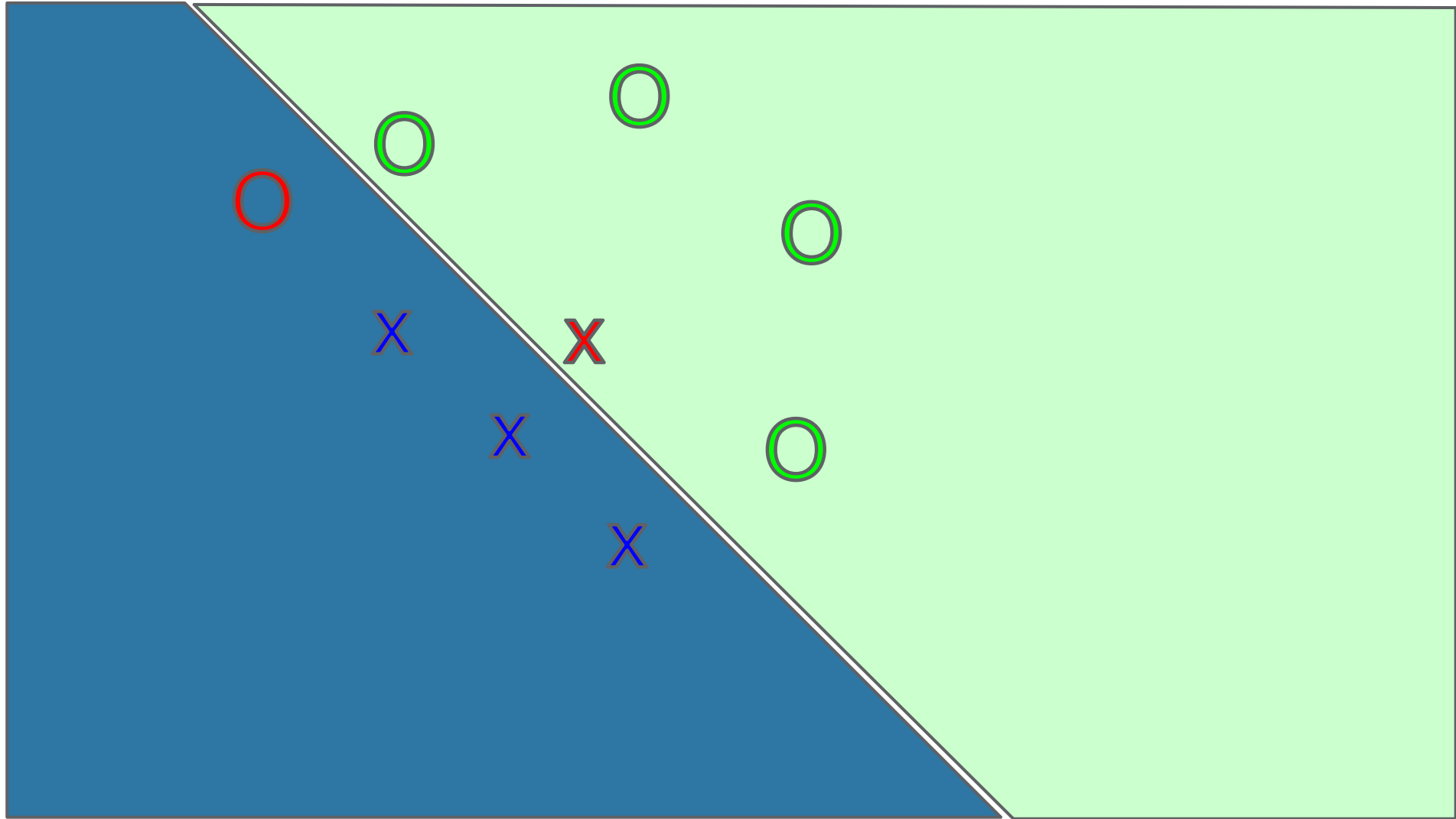● Training points for class 2

● Test point for class 2
● Adversarial example for class 2
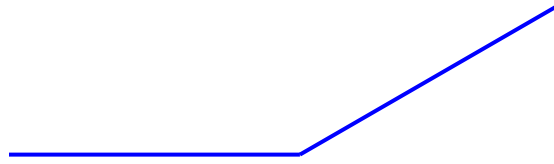
# Adversarial Examples from Overfitting

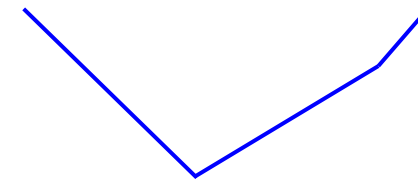# Adversarial Examples from Excessive Linearity
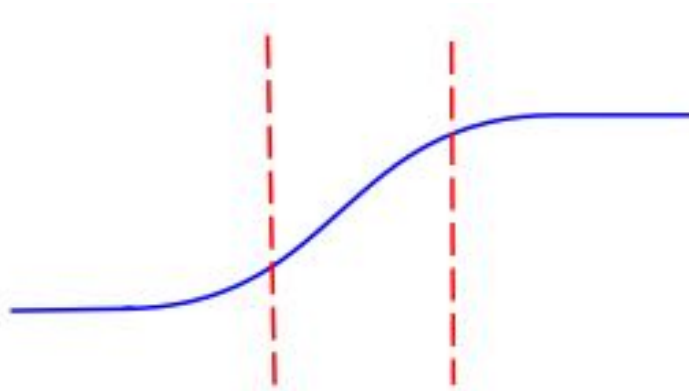
# Modern deep nets are very piecewise linear

Rectified linear unit

Maxout

Carefully tuned sigmoid

LSTM

| t=1 | t=2 | t=3 |
|-----|-----|-----|
| out | out | out |
| hid | hid | hid |
| inp | inp | inp |

# Gradient-based Adversarial Examples

- Fast Gradient Sign (Goodfellow et al., 2014)

$$\boldsymbol{X}^{adv} = \boldsymbol{X} + \epsilon \, \text{sign}\big(\nabla_X J(\boldsymbol{X}, y_{true})\big)$$

- Basic Iterative Method (Kurakin et al., 2017)

$$\boldsymbol{X}_0^{adv} = \boldsymbol{X}, \quad \boldsymbol{X}_{N+1}^{adv} = Clip_{X,\epsilon}\Big\{\boldsymbol{X}_N^{adv} + \alpha \, \text{sign}\big(\nabla_X J(\boldsymbol{X}_N^{adv}, y_{true})\big)\Big\}$$

$$Clip_{X,\epsilon}\{\boldsymbol{X}'\}(x,y,z) = \min\Big\{255, \boldsymbol{X}(x,y,z)+\epsilon, \max\big\{0, \boldsymbol{X}(x,y,z)-\epsilon, \boldsymbol{X}'(x,y,z)\big\}\Big\}$$

- Iterative Least-Likely Class Method (Kurakin et al., 2017)

$$y_{LL} = \arg\min_y\big\{p(y|\boldsymbol{X})\big\}$$

$$\boldsymbol{X}_0^{adv} = \boldsymbol{X}, \quad \boldsymbol{X}_{N+1}^{adv} = Clip_{X,\epsilon}\big\{\boldsymbol{X}_N^{adv} - \alpha \, \text{sign}\big(\nabla_X J(\boldsymbol{X}_N^{adv}, y_{LL})\big)\big\}$$

# Gradient-based Adversarial Examples

Original
sample

Fast
Gradient

Basic
Iterative
Method

Iterative
Least-Likely
Class Method

# Adversarial Examples in the Physical World
## (Kurakin, Goodfellow, Bengio, 2017)



(a) Printout

(b) Photo of printout

(c) Cropped image

Adversarial Examples In The Physical World
Kurakin A., Goodfellow I., Bengio S., 2016

# Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images
(Nguyen, Yosinski, Clune, 2014)



State-of-the-art DNNs can recognize real images with high confidence

2 But DNNs are also easily fooled: images can be produced that are unrecognizable to humans, but DNNs believe with 99.99% certainty are natural objects

# Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images
(Nguyen, Yosinski, Clune, 2014)

>99.6% confidences



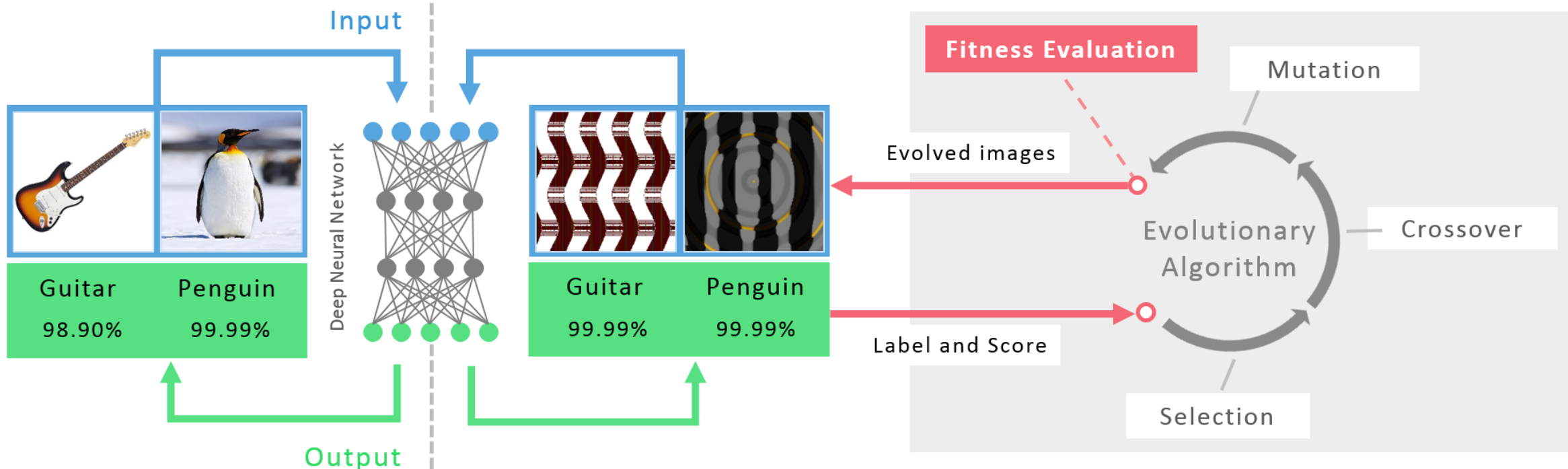| robin | cheetah | armadillo | lesser panda |
| centipede | peacock | jackfruit | bubble |

# Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images
(Nguyen, Yosinski, Clune, 2014)

>99.6% confidences

# Adversarial Learning – Failed Defenses

Generative
pretraining

Removing perturbation
with an autoencoder

Adding noise
at test time

Ensembles

Confidence-reducing
perturbation at test time

Error correcting
codes

Multiple glimpses

Weight decay

Double backprop

Adding noise
at train time

Various
non-linear units

Dropout

# Adversarial Learning – Defense Techniques

- Two defense techniques

1. Adversarial training (Szegedy et al., 2013)
    - a brute force solution where adversarial examples are generated and the model is explicitly trained not to be fooled by each of them.
    - improves the generalization of a machine learning model

2. Defensive distillation (Hinton et al., 2015; Papernot and McDaniel, 2016)
    - a strategy where the model is trained to output probabilities of different classes, rather than hard decisions about which class to output
    - smooths the model's decision surface in adversarial directions exploited by the adversary

# Adversarial Training

Labeled as bird

Still has same label (bird)



Decrease
probability
of bird class

# Adversarial Training

- Generate adversarial examples and use them while training
- Introduce an adversarial regularization term to the general loss function

$$\tilde{J}(\boldsymbol{\theta}, \boldsymbol{x}, y) = \alpha J(\boldsymbol{\theta}, \boldsymbol{x}, y) + (1 - \alpha) J(\boldsymbol{\theta}, \boldsymbol{x} + \epsilon \text{sign}\left(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y)\right)$$

training target          adversarial regularization

# Virtual Adversarial Training

Unlabeled; model
guesses it's probably
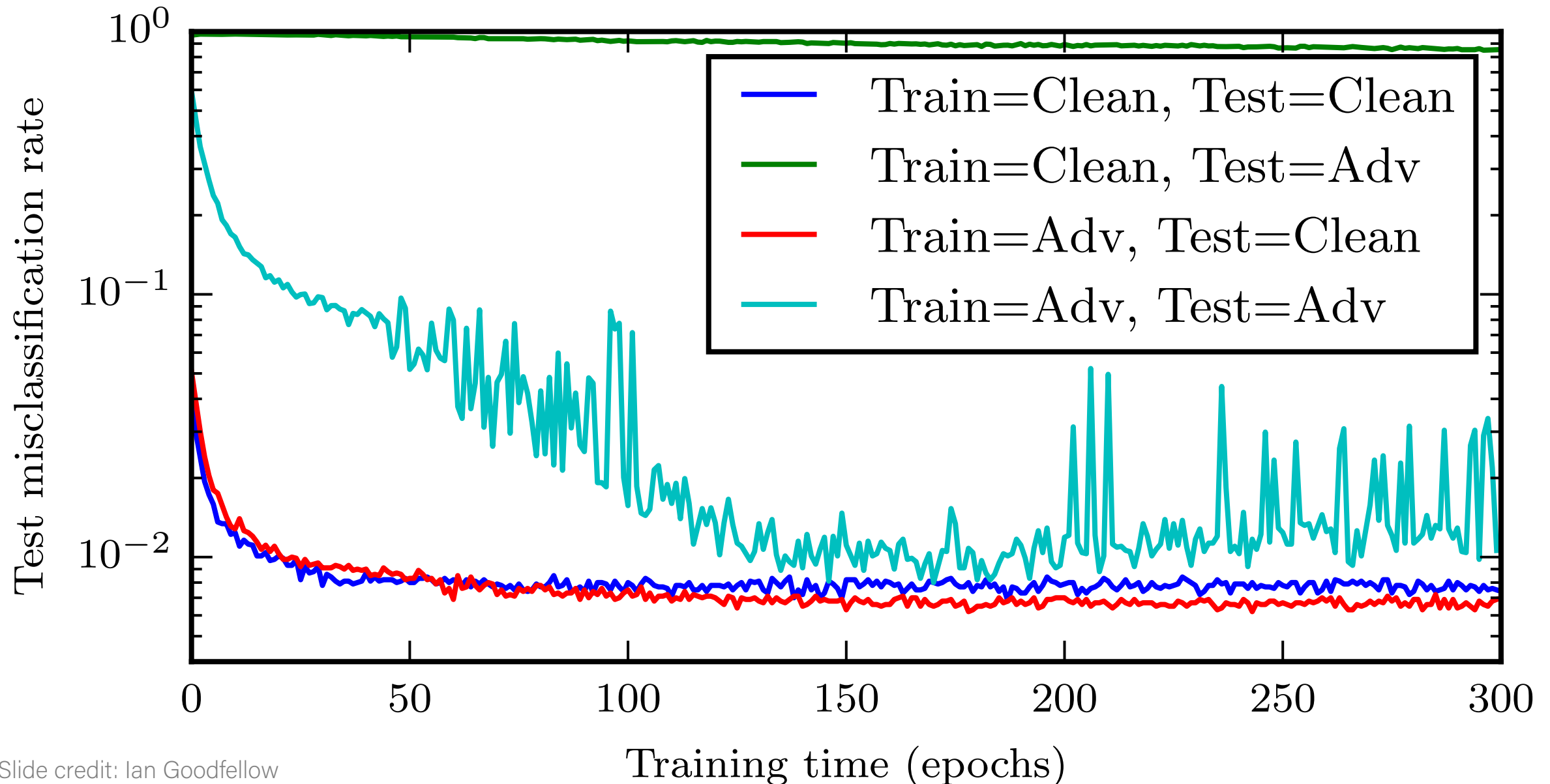a bird, maybe a plane

New guess should
match old guess
(probably bird, maybe plane)



Adversarial
perturbation
intended to
change the guess

# Training on Adversarial Examples



Slide credit: Ian Goodfellow

# Adversarial Training of Other Models

- Linear models: SVM / linear regression cannot learn a step function, so adversarial training is less useful, very similar to weight decay

- k-NN: adversarial training is prone to overfitting.

- Takeway: neural nets can actually become more secure than other models.

- Adversarially trained neural nets have the best empirical success rate on adversarial examples of any machine learning model.
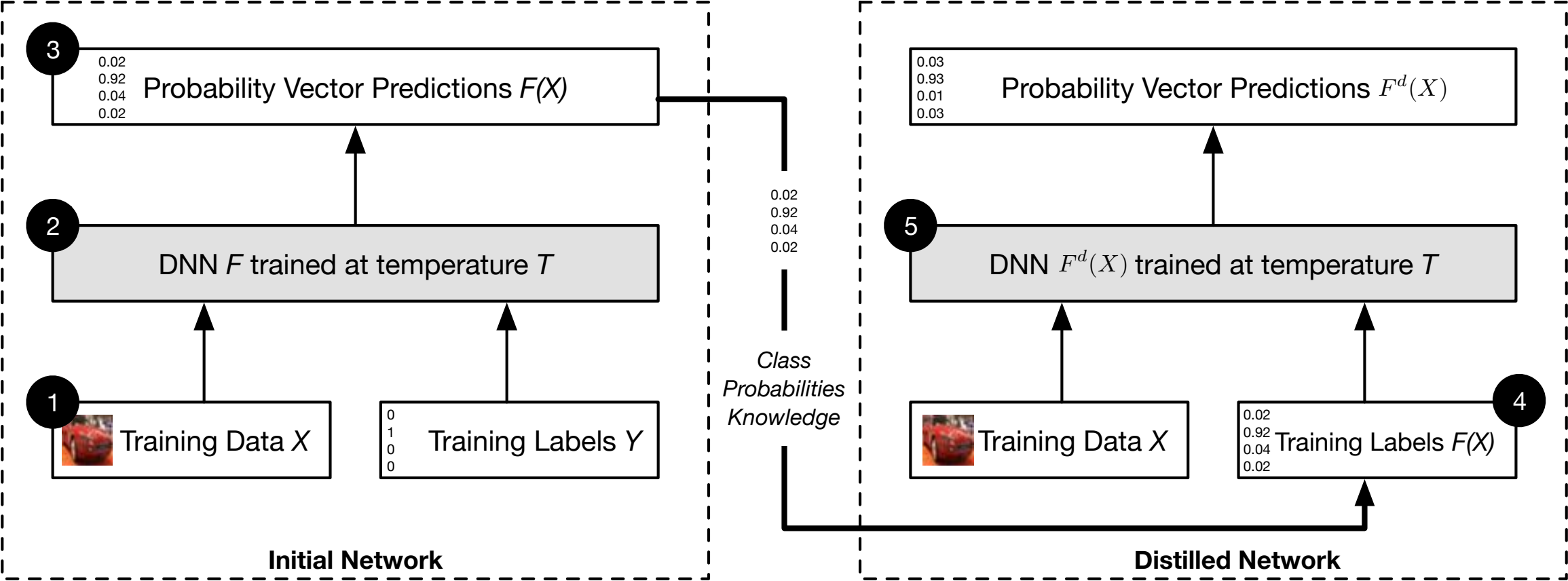
# Defensive Distillation

- Neural networks typically produce class probabilities by using a "softmax" output layer:

$$q_i = \frac{exp(z_i/T)}{\sum_j exp(z_j/T)}$$

T: a temperature that is normally set to 1

$q_i$: class probability

- Defensive distillation changes the training procedure essentially by re-configuring this "softmax" layer.

- It smooths the model's decision surface, eliminates overfitting, and thus increase robustness of the deep neural network model.

- Simplest form: Use the original model's predictions as the groundtruth labels to train the distilled model.

# Defensive Distillation



$$\frac{\partial C}{\partial z_i} = \frac{1}{T}\left(q_i - p_i\right) = \frac{1}{T}\left(\frac{e^{z_i/T}}{\sum_j e^{z_j/T}} - \frac{e^{v_i/T}}{\sum_j e^{v_j/T}}\right)$$
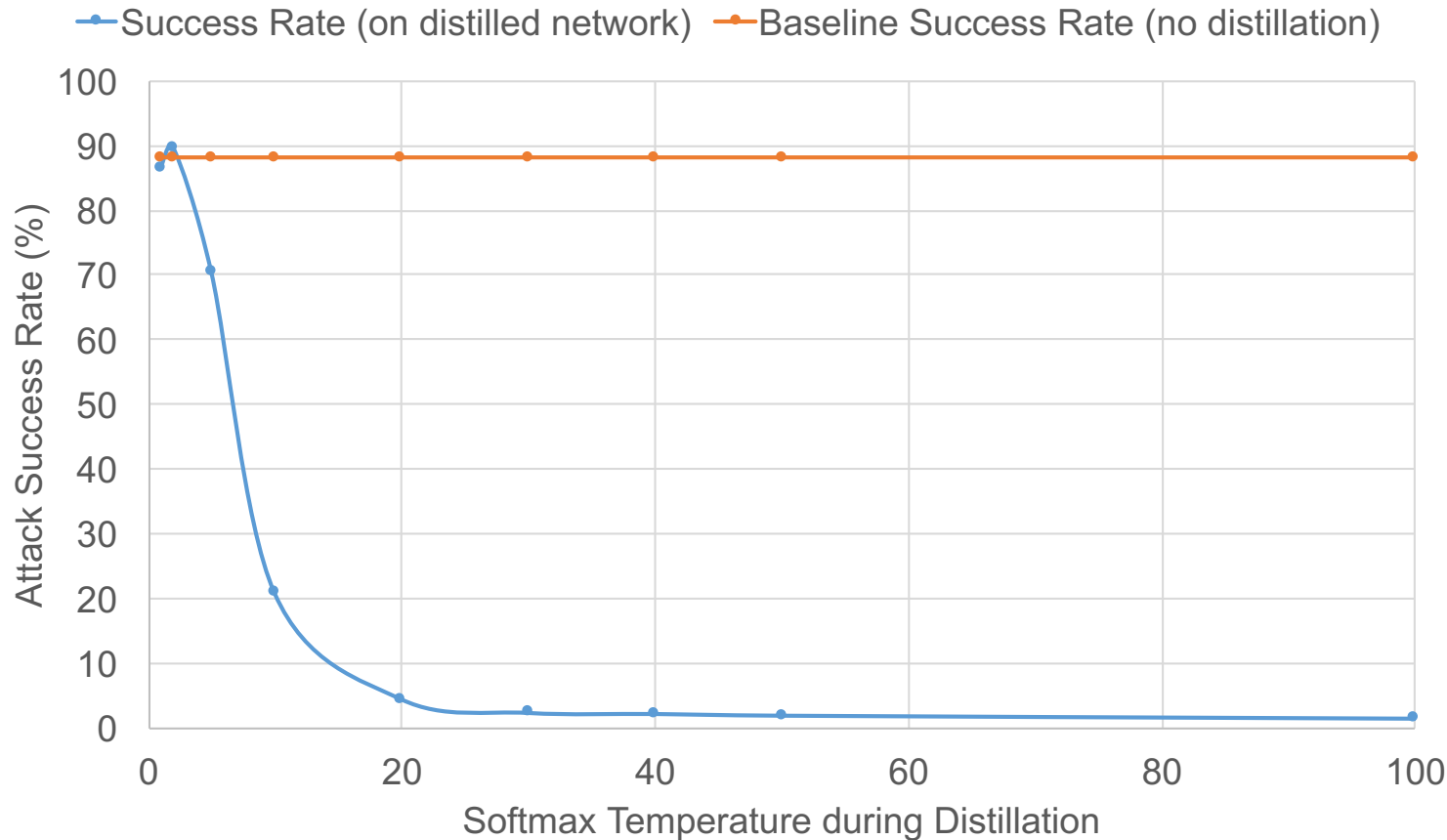
Papernot et al., 2016

# Defensive Distillation

- This strategy include the following steps (Papernot and McDaniel, 2016):

1.  Train a first instance of the neural network by using the training data (X, Y) where the labels Y indicate the correct class of samples X.

2.  Infer predictions of the training data and provide a new training dataset (X,f(X)) where the new class labels are the probability vectors quantifying the likeliness of X being in each class.

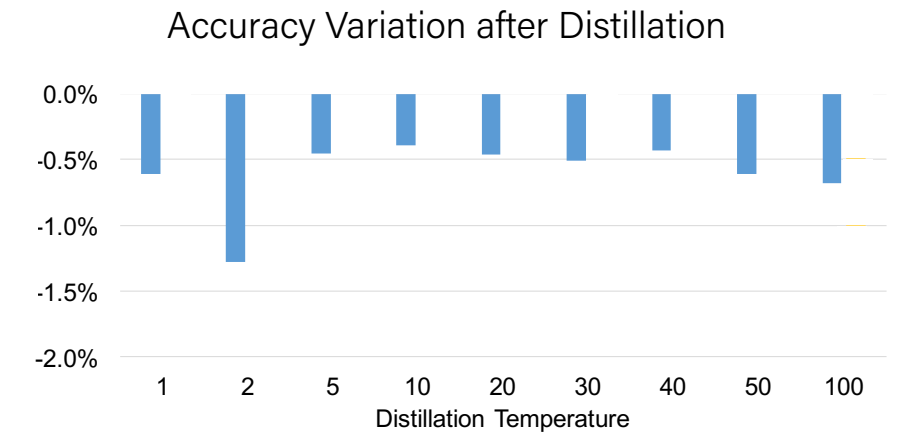3.  Train a *distilled* instance of the neural network f using this newly labeled dataset (X,f(X)).

While training the first and the distilled network, use the same high T values. At test time, deploy the distilled network by setting T back to 1.

# Defensive Distillation

- Distillation at high temperatures improves the smoothness of the network, and reduces its sensitivity to small input variations.
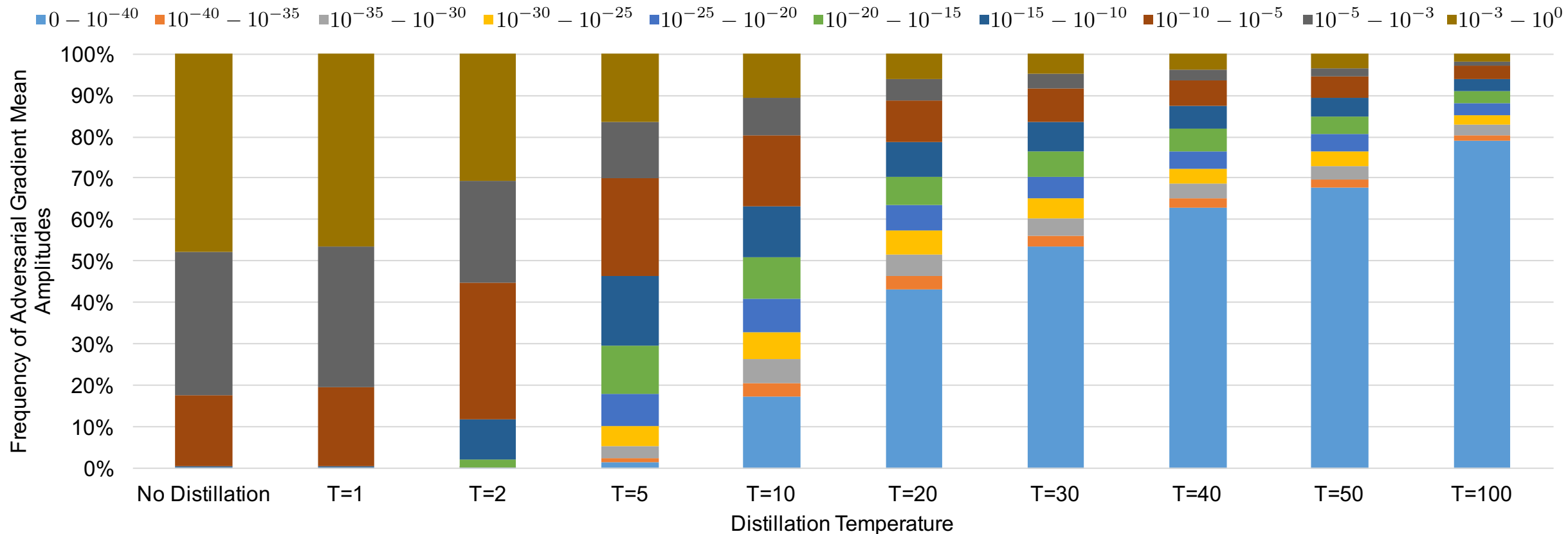


| | | |
|---|---|---|
| 2 | 82.23 | 87.67 |
| 5 | 24.67 | 67 |
| 10 | 6.78 | 47.56 |
| 20 | 1.34 | 18.23 |
| 30 | 1.44 | 13.23 |
| 40 | 0.45 | 9.34 |
| 50 | 1.45 | 6.23 |
| 100 | 0.45 | 5.11 |
| No distillation | 95.89 | 87.89 |

On the MNIST model - a 9 layer deep neural network with a 99.5% test accuracy (Papernot and McDaniel, 2016)

**Success Rate (on distilled network)** **Baseline Success Rate (no distillation)**

Attack Success Rate (%) vs Softmax Temperature during Distillation

Accuracy Variation after Distillation

Distillation Temperature

# Distillation and Sensitivity

- Distillation reduces gradients exploited by adversaries to craft perturbations.



10,000 samples from the CIFAR10 test set into bins according to the mean value of their adversarial gradient amplitude. (Papernot et al., 2016)

# Summary

- Big gains in performance on perceptual tasks by using deep neural network models.

- Machine learning has not yet reached true human-level performance.

- Adversarial examples show that many modern machine learning algorithms can be easily fooled.

- Many different ways of attacking deep neural network models.

- Very few ways of defending deep neural network models.

- Recent work (Papernot et al., 2017) considers more realistic threat models
  - The adversary have no knowledge of the machine learning architecture and model parameters.