**10-418 / 10-618 Machine Learning for Structured Data**

Machine Learning Department
School of Computer Science
Carnegie Mellon University

**ML**
MACHINE LEARNING
D E P A R T M E N T

# Convolutional Neural Networks

Matt Gormley
Lecture 16
Oct. 21, 2019

# Reminders

- **Homework 3: Structured SVM**
  - **Out: Tue, Oct. 18**
  - **Due: Mon, Nov. 4 at 11:59pm**
- **Midterm Exam Viewing**
- **Midsemester Grades**

aka. Max-Margin Markov Networks (M³Ns)

# STRUCTURED SVM

# Structured Perceptron

**Whiteboard**

– Warmup: Binary SVM

– Warmup: Binary SVM Hinge Loss

– Structured Large Margin

– Structured Hinge Loss

– Gradient of Structured Hinge Loss

– SGD for Structured SVM

– Loss Augmented MAP Inference

# Max vs "Soft-Max" Margin

- ## SVMs:

$$\min_{\mathbf{w}} k||\mathbf{w}||^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}} \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(y) \right) \right)$$

Hard (Penalized) Margin

- ## Maxent:

$$\min_{\mathbf{w}} \quad k||w||^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) \right) \right)$$

Soft Margin

- ## Very similar!  Both try to make the true score better than a function of the other scores.
  - The SVM tries to beat the augmented runner-up
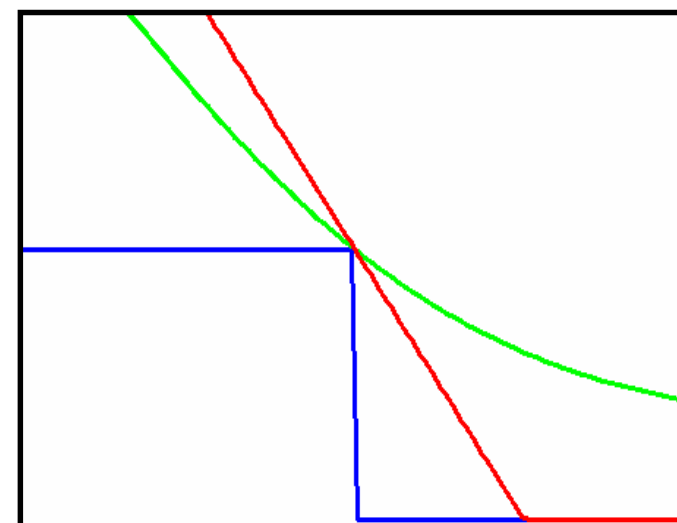  - The maxent classifier tries to beat the "soft-max"

# Hinge Loss

- Consider the per-instance SVM objective:

$$\min_{\mathbf{w}} k||\mathbf{w}||^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}} \left[ \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(y) \right] \right)$$
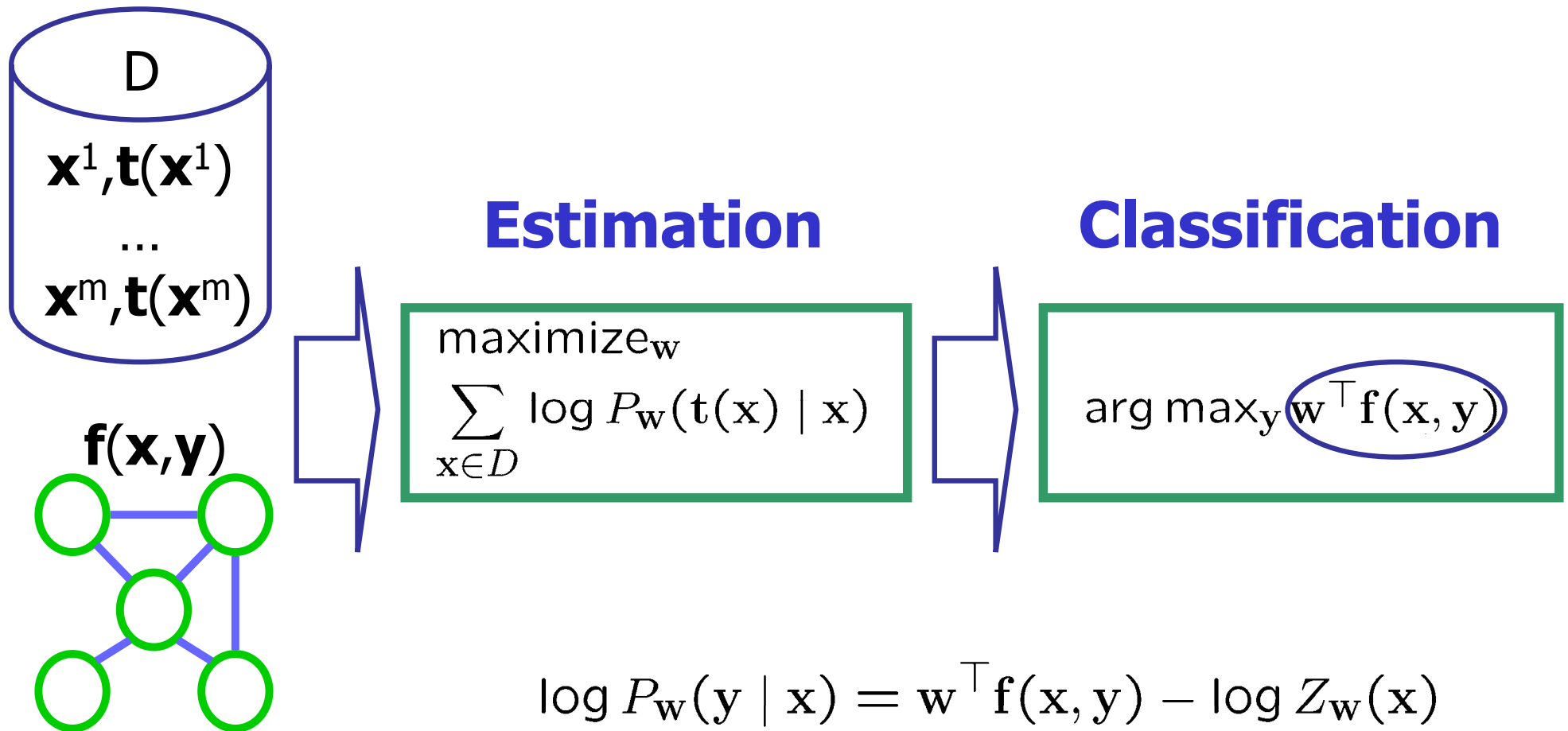
- This is called the "hinge loss"
  - Upper bounds zero-one loss
  - Unlike maxent / log loss, you stop gaining objective once the true label wins by enough
  - You can start from here and derive the SVM objective

$$\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y} \neq \mathbf{y}^i} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})$$

# Max (Conditional) Likelihood

$$D$$

$$\mathbf{x}^1, \mathbf{t}(\mathbf{x}^1)$$

$$\ldots$$

$$\mathbf{x}^m, \mathbf{t}(\mathbf{x}^m)$$

$$\mathbf{f}(\mathbf{x}, \mathbf{y})$$

**Estimation**

$$\text{maximize}_{\mathbf{w}}$$
$$\sum_{\mathbf{x} \in D} \log P_{\mathbf{w}}(\mathbf{t}(\mathbf{x}) \mid \mathbf{x})$$

**Classification**

$$\arg\max_{\mathbf{y}} \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y})$$

$$\log P_{\mathbf{w}}(\mathbf{y} \mid \mathbf{x}) = \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y}) - \log Z_{\mathbf{w}}(\mathbf{x})$$

Don't need to learn entire distribution!
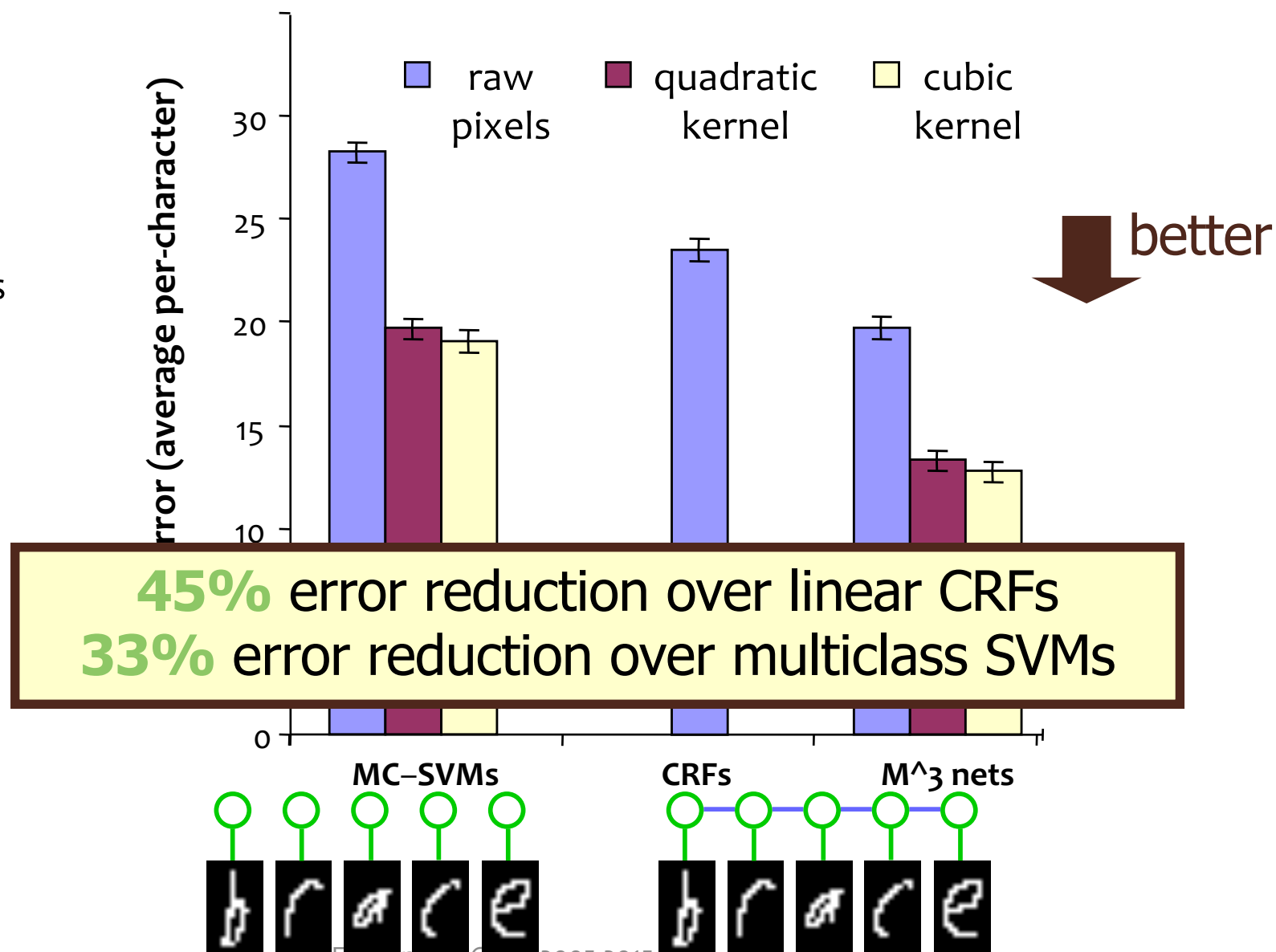
# Structured SVM

The original name for **Structured SVM:**

- **Max-Margin Markov Networks**
- **abbreviated as** M$^3$Ns

# Results: Handwriting Recognition

Length: ~8 chars
Letter: 16x8 pixels
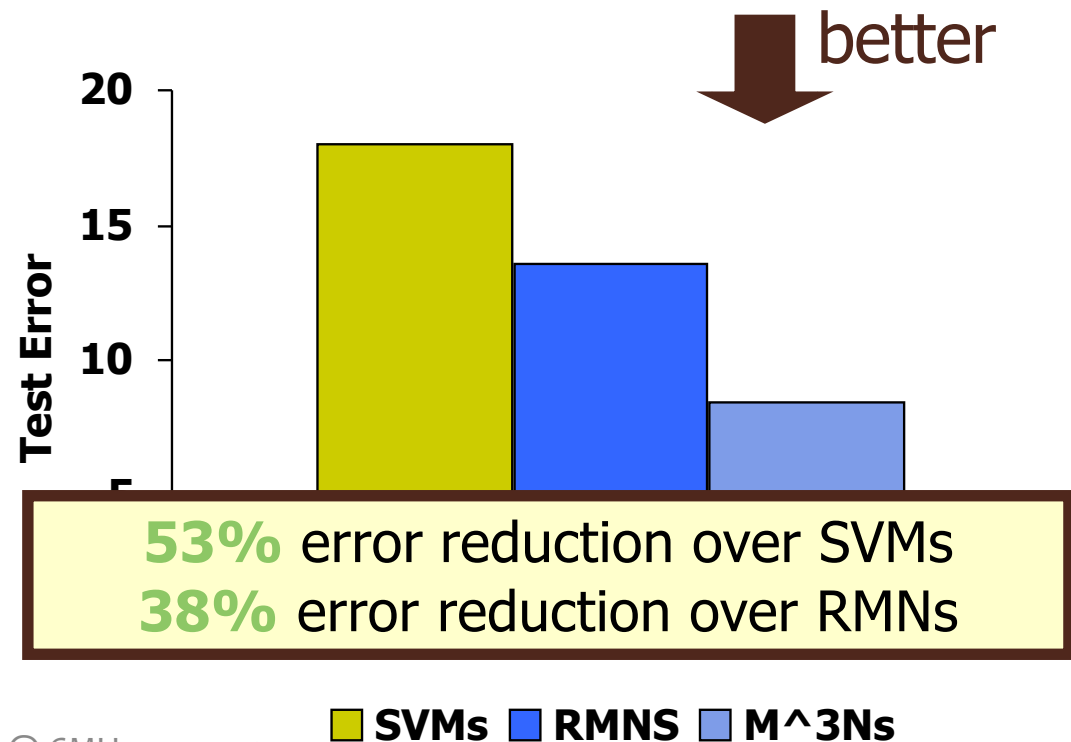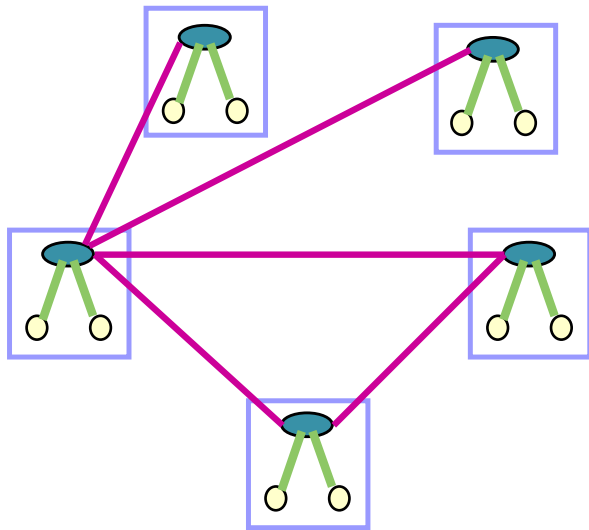10-fold Train/Test
5000/50000 letters
600/6000 words

Models:
  Multiclass-SVMs*
  CRFs
  M³ nets



better

**45%** error reduction over linear CRFs
**33%** error reduction over multiclass SVMs
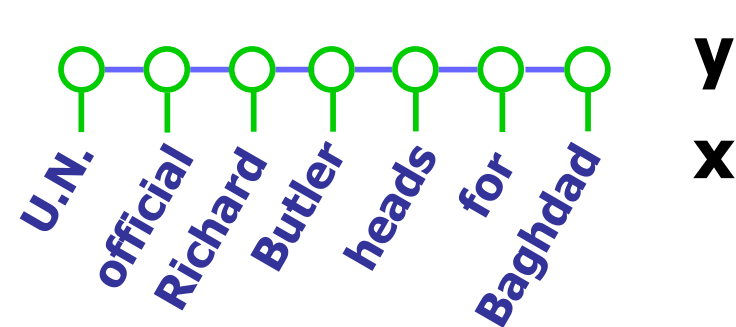
*Crammer & Singer 01

# Results: Hypertext Classification

- WebKB dataset
  - Four CS department websites: 1300 pages/3500 links
  - Classify each page: faculty, course, student, project, other
  - Train on three universities/test on fourth
- Inference: loopy belief propagation
- Learning: relaxed dual



*Taskar et al 02

**53%** error reduction over SVMs
**38%** error reduction over RMNs

better

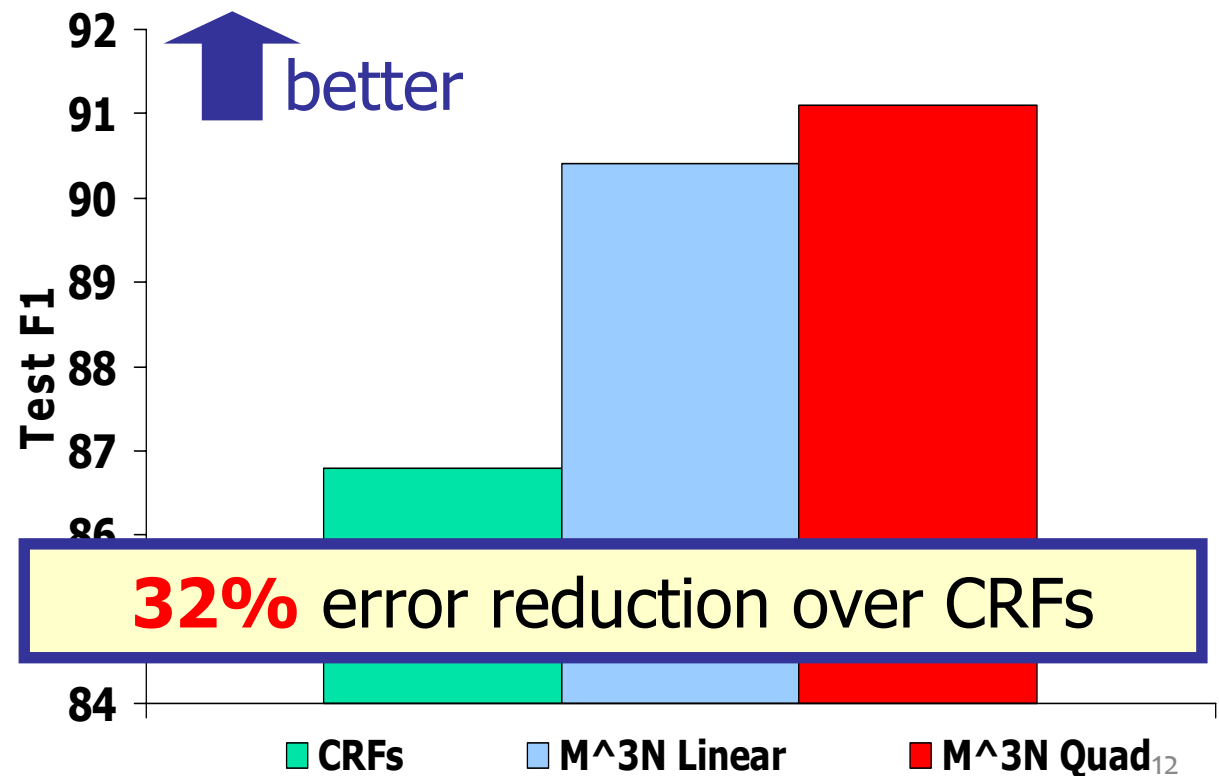SVMs  RMNS  M^3Ns

Test Error

# Named Entity Recognition

- Locate and classify named entities in sentences:
  - 4 categories: organization, person, location, misc.
  - e.g. "U.N. official Richard Butler heads for Baghdad".
- CoNLL 03 data set (200K words train, 50K words test)

$y_i$ = org/per/loc/misc/none

$\mathbf{f}(y_i, x) = [...,$
$I(y_i=org, x_i="U.N."),$
$I(y_i=per, x_i=capitalized),$
$I(y_i=loc, x_i=known\ city),$
$..., ]$

**32%** error reduction over CRFs

CRFs  M^3N Linear  M^3N Quad

# Associative Markov networks

$$P(\mathbf{y} \mid \mathbf{x}) \propto \prod_i \phi_i(y_i, \mathbf{x}_i) \prod_{ij} \phi_{ij}(y_i, y_j, \mathbf{x}_{ij}) = \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$

$\underbrace{\phantom{\prod_i \phi_i(y_i, \mathbf{x}_i)}}$ **Point features** spin-images, point height

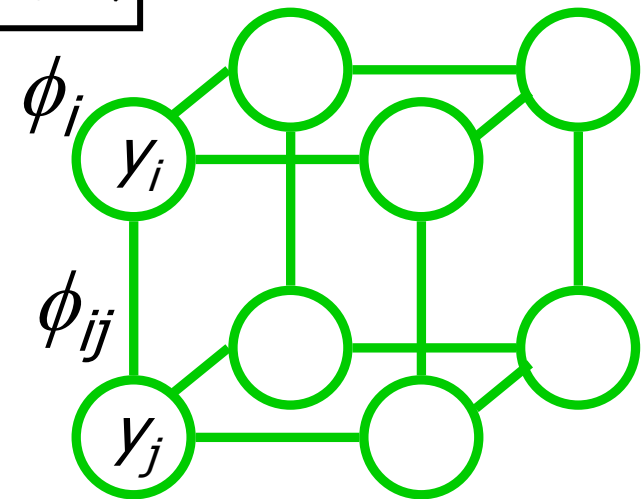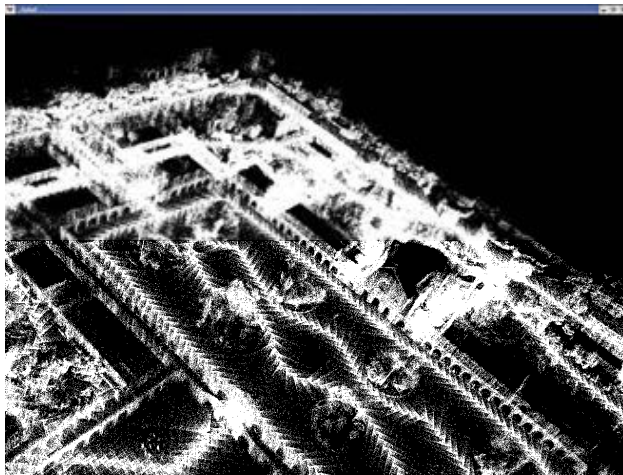$\underbrace{\phantom{\prod_{ij} \phi_{ij}}}$ **Edge features** length of edge, edge orientation

"associative" restriction

$$\phi_{ij}(y_i, y_j) = \begin{bmatrix} \phi_{ij}(1,1) & & 1 \\ & \ddots & \\ 1 & & \phi_{ij}(K,K) \end{bmatrix}$$

bonus

$$\phi_{ij}(k, k) \geq 1$$



$\phi_i$ $y_i$
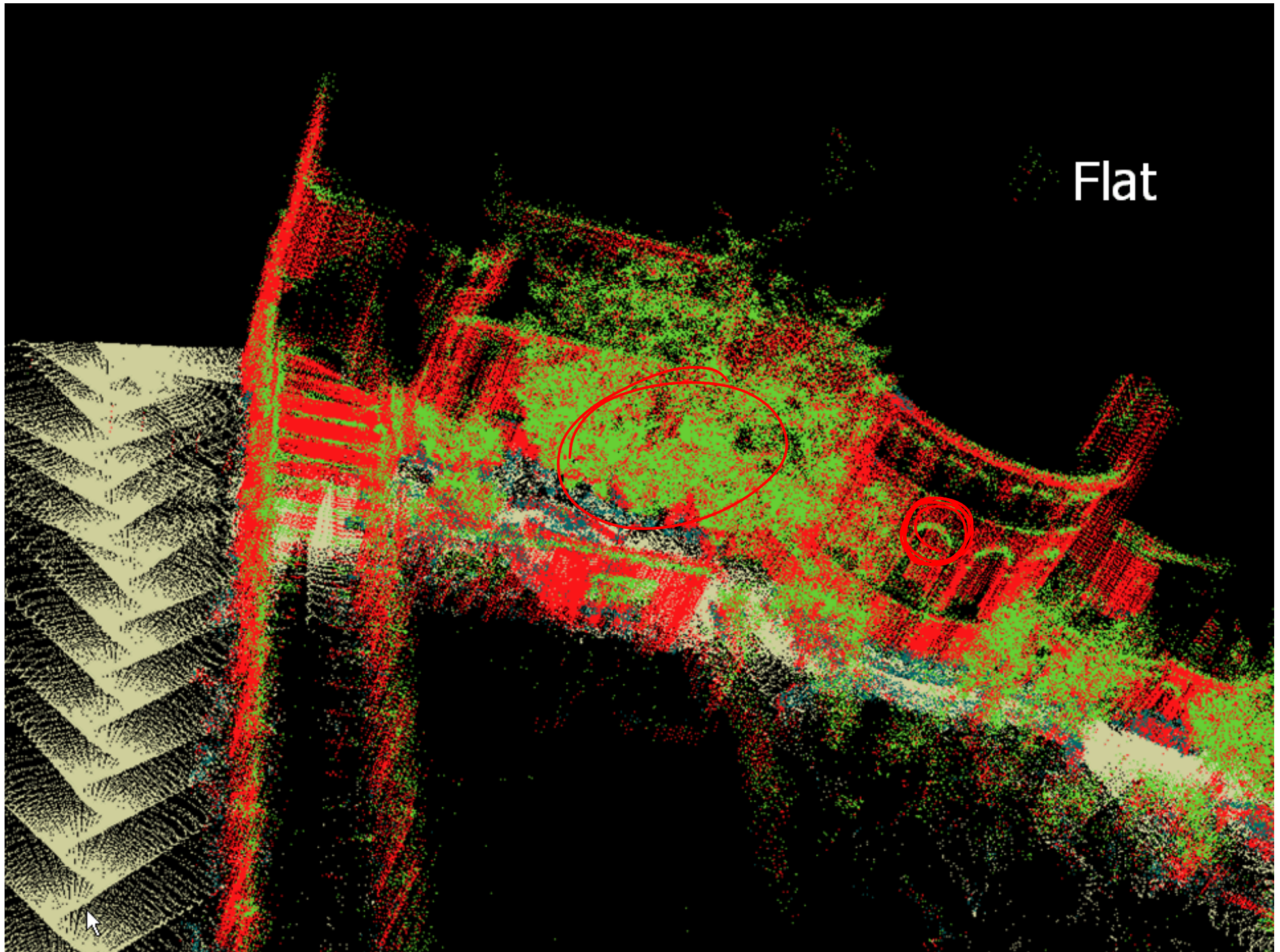
$\phi_{ij}$

$y_j$

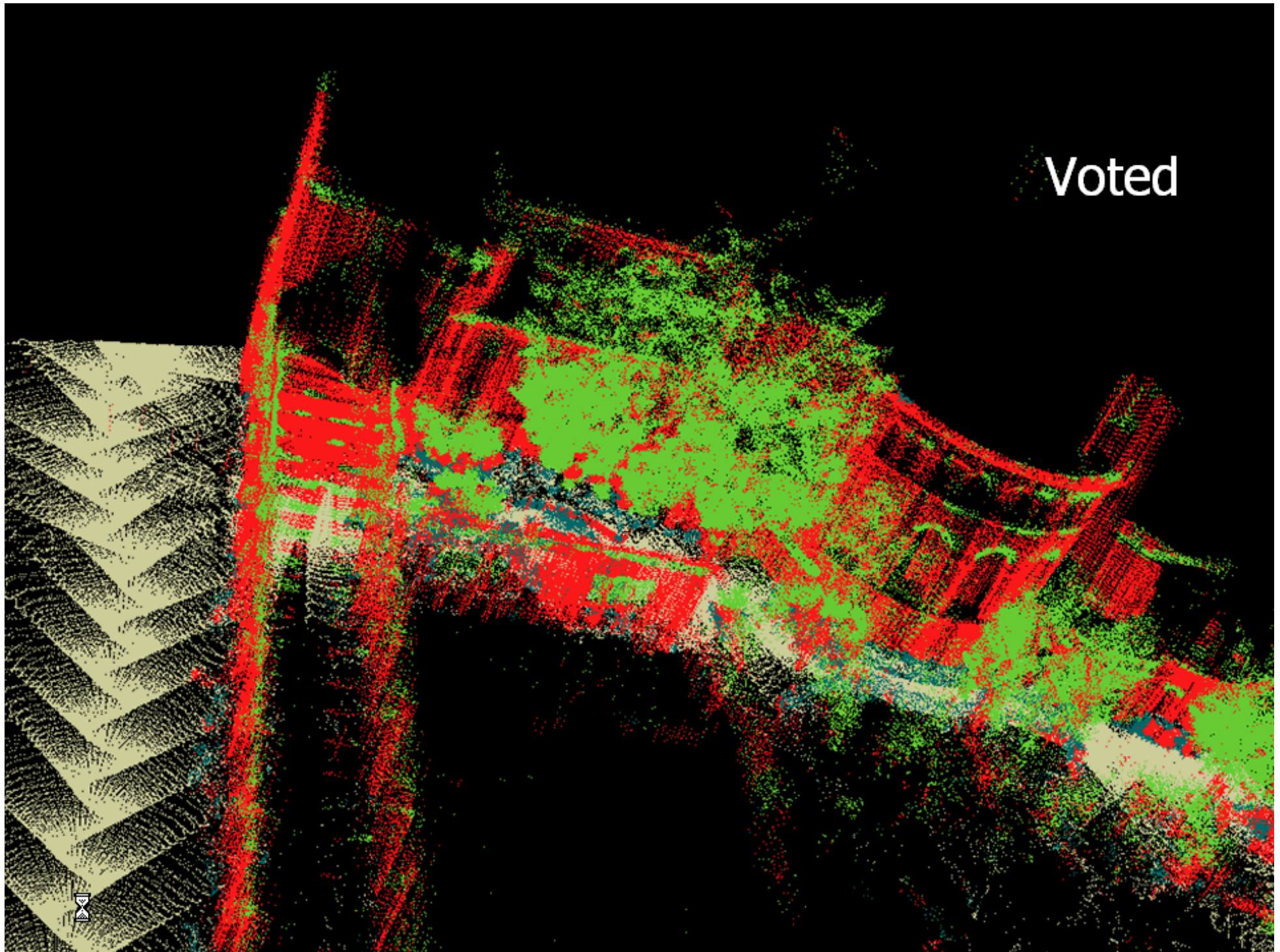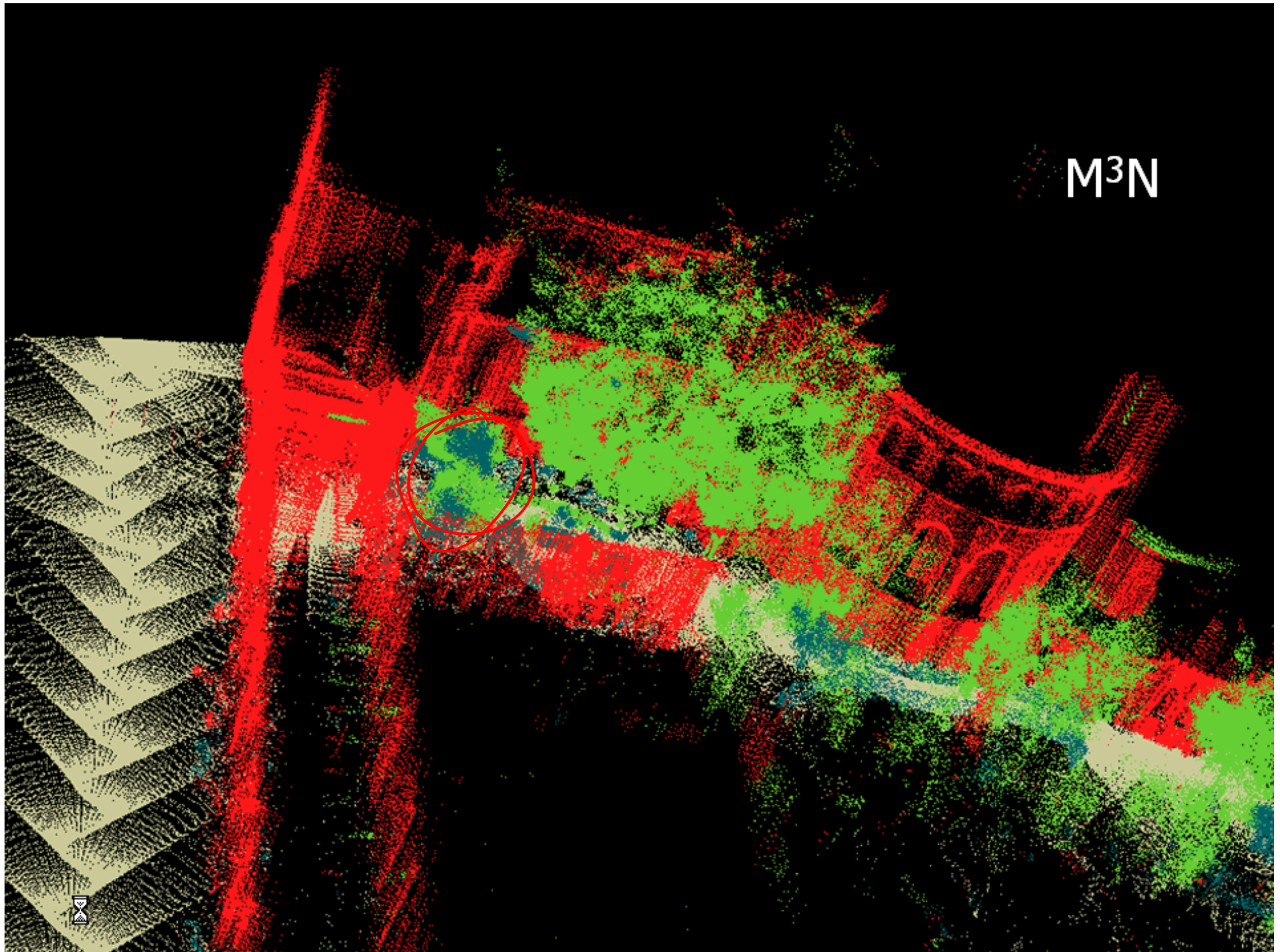# Max-margin AMNs results



Label: ground, building, tree, shrub

Training: 30 thousand points     Testing: 3 million points

Flat

Voted

M³N

# Segmentation results

Hand labeled 180K test points

| Model | Accuracy |
|-------|----------|
| SVM | 68% |
| V-SVM | 73% |
| M³N | 93% |

# CNNs Outline

- **Background: Computer Vision**
  - Image Classification
  - ILSVRC 2010 - 2016
  - Traditional Feature Extraction Methods
  - Convolution as Feature Extraction
- **Convolutional Neural Networks (CNNs)**
  - Learning Feature Abstractions
  - Common CNN Layers:
    - Convolutional Layer
    - Max-Pooling Layer
    - Fully-connected Layer (w/tensor input)
    - Softmax Layer
    - ReLU Layer
  - Background: Subgradient
  - Architecture: LeNet
  - Architecture: AlexNet
  - Architecture: ResNet
- **Training a CNN**
  - SGD for CNNs
  - Backpropagation for CNNs

# BACKGROUND: COMPUTER VISION

# Example: Image Classification

- ImageNet LSVRC-2011 contest:
  - **Dataset**: 1.2 million labeled images, 1000 classes
  - **Task**: Given a new image, label it with the correct class
  - **Multiclass** classification problem
- Examples from http://image-net.org/

# German iris, Iris kochi

iris of northern Italy having deep blue-purple flowers; similar to but smaller than Iris germanica

# Court, courtyard

An area wholly or partly surrounded by walls or buildings. "the house was built around an inner court"

# Feature Engineering for CV

Edge detection (Canny)



Corner Detection (Harris)

Scale Invariant Feature Transform (SIFT)

Figure from Lowe (1999) and Lowe (2004)

# Example: Image Classification

**CNN for Image Classification**
(Krizhevsky, Sutskever & Hinton, 2012)
15.3% error on ImageNet LSVRC-2012 contest

Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax

# CNNs for Image Recognition

Slide from Kaiming He

# CONVOLUTION

# What's a convolution?

- Basic idea:
  - Pick a 3x3 matrix F of weights
  - Slide this over an image and compute the "inner product" (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

- Key point:
  - Different convolutions extract different types of low-level "features" from an image
  - All that we need to vary to generate these different features is the weights of F

$\underline{Ex:}$ 1 input channel, 1 output channel

$\underline{Input}$

| $X_{11}$ | $X_{12}$ | $X_{13}$ |
|---|---|---|
| $X_{21}$ | $X_{22}$ | $X_{23}$ |
| $X_{31}$ | $X_{32}$ | $X_{33}$ |

$\longrightarrow$

$\underline{Conv}$

| $\alpha_{11}$ | $\alpha_{12}$ |
|---|---|
| $\alpha_{21}$ | $\alpha_{22}$ |

$\longrightarrow$

$\underline{Output}$

| $Y_{11}$ | $Y_{12}$ |
|---|---|
| $Y_{21}$ | $Y_{22}$ |

$$Y_{11} = \alpha_{11} X_{11} + \alpha_{12} X_{12} + \alpha_{21} X_{21} + \alpha_{22} X_{22} + \alpha_0$$

$$Y_{12} = \alpha_{11} X_{12} + \alpha_{12} X_{13} + \alpha_{21} X_{22} + \alpha_{22} X_{23} + \alpha_0$$

$$Y_{21} = \alpha_{11} X_{21} + \alpha_{12} X_{22} + \alpha_{21} X_{31} + \alpha_{22} X_{32} + \alpha_0$$

$$Y_{22} = \alpha_{11} X_{22} + \alpha_{12} X_{23} + \alpha_{21} X_{32} + \alpha_{22} X_{33} + \alpha_0$$

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

### Input Image

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

### Convolution

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

### Convolved Image

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

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

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

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| 2 | 0 | 2 | 1 | 0 |
| 2 | 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| 2 | 0 | 2 | 1 | 0 |
| 2 | 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

### Input Image

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

### Convolution

### Convolved Image

| | | | | |
|---|---|---|---|---|
| 3 | 2 | 2 | 3 | 1 |
| 2 | 0 | 2 | 1 | 0 |
| 2 | 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| 2 | 0 | 2 | 1 | 0 |
| 2 | 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

Convolution

Convolved Image

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.



Input Image

Convolution
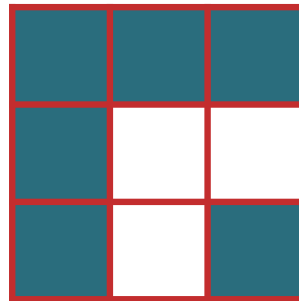
Convolved Image

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

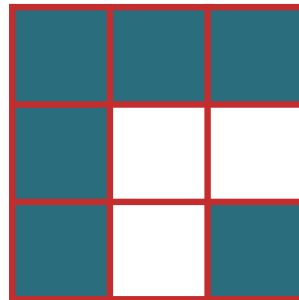| 3 | 2 | 2 | | |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

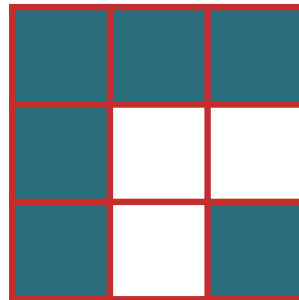| 3 | 2 | 2 | 3 |   |
|---|---|---|---|---|
|   |   |   |   |   |
|   |   |   |   |   |
|   |   |   |   |   |
|   |   |   |   |   |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| | | | | |
| | | | | |
| | | | | |
| | | | | |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| 2 |   |   |   |   |
|   |   |   |   |   |
|   |   |   |   |   |
|   |   |   |   |   |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

Convolution

Convolved Image

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Convolution

Convolved Image

| 3 | 2 | 2 | 3 | 1 |
|---|---|---|---|---|
| 2 | 0 | 2 | 1 | 0 |
| 2 | 2 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Identity Convolution

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

Convolved Image

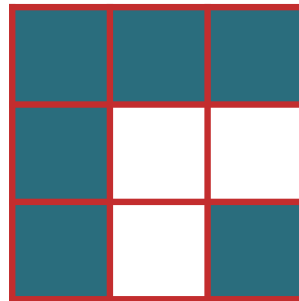| 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |

# Background: Image Processing

A **convolution matrix** is used in image processing for tasks such as edge detection, blurring, sharpening, etc.

Input Image

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

Blurring
Convolution

| .1 | .1 | .1 |
|----|----|----|
| .1 | .2 | .1 |
| .1 | .1 | .1 |

Convolved Image

| .4 | .5 | .5 | .5 | .4 |
|----|----|----|----|----|
| .4 | .2 | .3 | .6 | .3 |
| .5 | .4 | .4 | .2 | .1 |
| .5 | .6 | .2 | .1 | 0 |
| .4 | .3 | .1 | 0 | 0 |

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo

# What's a convolution?

- Basic idea:
  - Pick a 3x3 matrix F of weights
  - Slide this over an image and compute the "inner product" (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

- Key point:
  - Different convolutions extract different types of low-level "features" from an image
  - All that we need to vary to generate these different features is the weights of F

Ex: 1 input channel, 1 output channel

Input

| $X_{11}$ | $X_{12}$ | $X_{13}$ |
| $X_{21}$ | $X_{22}$ | $X_{23}$ |
| $X_{31}$ | $X_{32}$ | $X_{33}$ |

Conv

| $\alpha_{11}$ | $\alpha_{12}$ |
| $\alpha_{21}$ | $\alpha_{22}$ |

Output

| $Y_{11}$ | $Y_{12}$ |
| $Y_{21}$ | $Y_{22}$ |

$$Y_{11} = \alpha_{11} X_{11} + \alpha_{12} X_{12} + \alpha_{21} X_{21} + \alpha_{22} X_{22} + \alpha_0$$

$$Y_{12} = \alpha_{11} X_{12} + \alpha_{12} X_{13} + \alpha_{21} X_{22} + \alpha_{22} X_{23} + \alpha_0$$

$$Y_{21} = \alpha_{11} X_{21} + \alpha_{12} X_{22} + \alpha_{21} X_{31} + \alpha_{22} X_{32} + \alpha_0$$

$$Y_{22} = \alpha_{11} X_{22} + \alpha_{12} X_{23} + \alpha_{21} X_{32} + \alpha_{22} X_{33} + \alpha_0$$

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | | |
|---|---|---|
| | | |
| | | |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | 3 | |
|---|---|---|
| | | |
| | | |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| | | |
|---|---|---|
| 3 | 3 | 1 |
| | | |
| | | |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | 3 | 1 |
|---|---|---|
| 3 |   |   |
|   |   |   |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

### Input Image

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

### Convolution

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

### Convolved Image

| 3 | 3 | 1 |
|---|---|---|
| 3 | 1 |   |
|   |   |   |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | 3 | 1 |
|---|---|---|
| 3 | 1 | 0 |
|   |   |   |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | 3 | 1 |
|---|---|---|
| 3 | 1 | 0 |
| 1 |   |   |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image

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

Convolution

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

Convolved Image

| 3 | 3 | 1 |
|---|---|---|
| 3 | 1 | 0 |
| 1 | 0 |   |

# Downsampling

- Suppose we use a convolution with stride 2
- Only 9 patches visited in input, so only 9 pixels in output

Input Image



Convolution

Convolved Image

# CONVOLUTIONAL NEURAL NETS

# A Recipe for Machine Learning

**1. Given training data:**

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

**2. Choose each of these:**

- Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

- Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

**3. Define goal:**

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

**4. Train with SGD:**

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# A Recipe for Machine Learning

1.

2. Choose each of these:
   – Decision function

$$\hat{y} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

   – Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

- Convolutional Neural Networks (CNNs) provide another form of **decision function**
- Let's see what they look like…

# Convolutional Neural Network (CNN)

- Typical layers include:
  - Convolutional layer
  - Max-pooling layer
  - Fully-connected (Linear) layer
  - ReLU layer (or some other nonlinear activation function)
  - Softmax
- These can be arranged into arbitrarily deep topologies
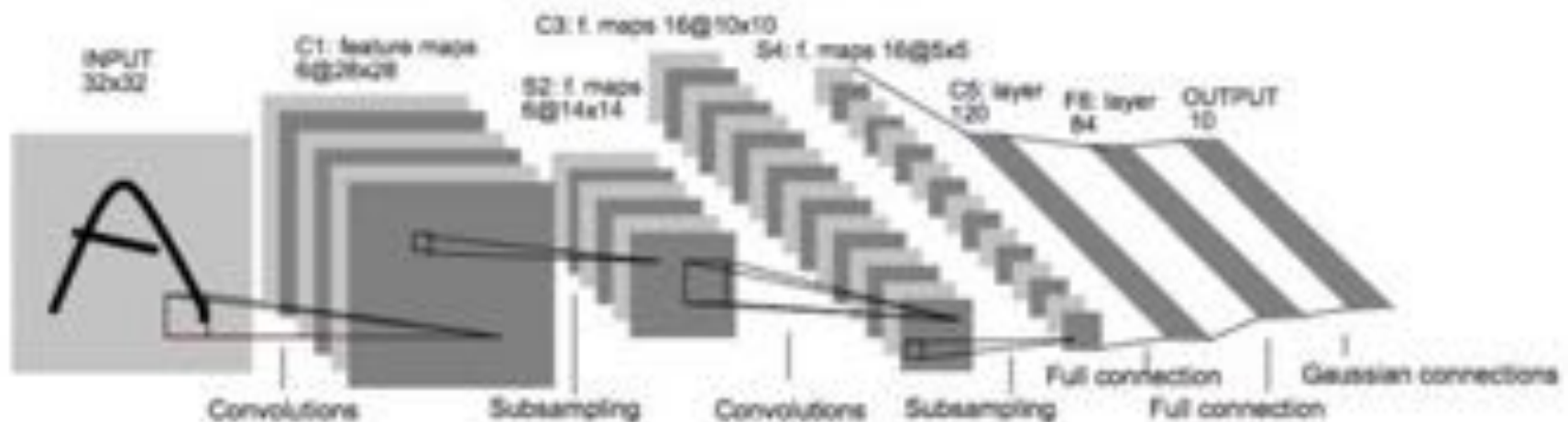
# Architecture #1: LeNet-5



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.

# Convolutional Layer

**CNN** key idea:
Treat convolution matrix as
parameters and learn them!

Input Image

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

Learned
Convolution

| $\theta_{11}$ | $\theta_{12}$ | $\theta_{13}$ |
|---|---|---|
| $\theta_{21}$ | $\theta_{22}$ | $\theta_{23}$ |
| $\theta_{31}$ | $\theta_{32}$ | $\theta_{33}$ |

Convolved Image

| .4 | .5 | .5 | .5 | .4 |
|---|---|---|---|---|
| .4 | .2 | .3 | .6 | .3 |
| .5 | .4 | .4 | .2 | .1 |
| .5 | .6 | .2 | .1 | 0 |
| .4 | .3 | .1 | 0 | 0 |

# Downsampling by Averaging

- Downsampling by averaging **used to be** a common approach
- This is a special case of convolution where the weights are fixed to a uniform distribution
- The example below uses a stride of 2

Input Image

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

Convolution

| 1/4 | 1/4 |
|-----|-----|
| 1/4 | 1/4 |

Convolved Image

| 3/4 | 3/4 | 1/4 |
|-----|-----|-----|
| 3/4 | 1/4 | 0 |
| 1/4 | 0 | 0 |

# Max-Pooling

- Max-pooling is another (common) form of downsampling
- Instead of averaging, we take the max value within the same range as the equivalently-sized convolution
- The example below uses a stride of 2

Input Image

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

Max-pooling

| $x_{i,j}$ | $x_{i,j+1}$ |
|---|---|
| $x_{i+1,j}$ | $x_{i+1,j+1}$ |

Max-Pooled Image

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

$$y_{ij} = \max(x_{ij},$$
$$x_{i,j+1},$$
$$x_{i+1,j},$$
$$x_{i+1,j+1})$$

# TRAINING CNNS

# A Recipe for Machine Learning

1. Given training data:
$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

2. Choose each of these:
   – Decision function
   $$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$
   – Loss function
   $$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

1. Given training data:

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

3. Define goal:

2. Choose each of th

– Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

opposite the gradient)

$$\boldsymbol{\theta}^{(t)} \quad {}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

- Q: Now that we have the CNN as a decision function, how do we compute the gradient?

- A: Backpropagation of course!

# SGD for CNNs

Ex: Architecture:     Given $\vec{x}, y^*$

$J = \ell(y, y^*)$

$y = \text{softmax}(z^{(5)})$

$z^{(5)} = \text{linear}(z^{(4)}, W)$

$z^{(4)} = \text{relu}(z^{(3)})$

$z^{(3)} = \text{conv}(z^{(2)}, \beta)$

$z^{(2)} = \text{max-pool}(z^{(1)})$

$z^{(1)} = \text{conv}(\vec{x}, \alpha)$

Parameters $\vec{\Theta} = [\alpha, \beta, W]$

## SGD:

① Init $\vec{\Theta}$

② While not converged:

    Sample $i \in \{1, \dots, N\}$

    Forward: $y = h_\Theta(\vec{x}^{(i)})$, $J_i(\Theta) = \ell(y, y^*)$

    Backward: $\nabla_{\vec{\Theta}} J_i(\Theta) = \dots$

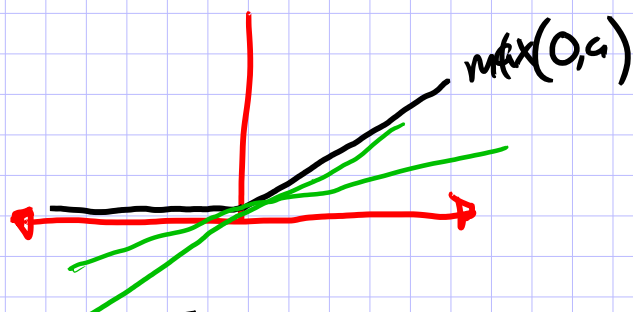    Update: $\vec{\Theta} \leftarrow \vec{\Theta} - \lambda \nabla_{\vec{\Theta}} J_i(\Theta)$

# LAYERS OF A CNN

# ReLU Layer

ReLU Layer    Input: $\vec{x} \in \mathbb{R}^k$    Output: $\vec{y} \in \mathbb{R}^k$

Forward:

$\vec{y} = \sigma(\vec{x})$   ← element-wise

$\sigma(a) = \max(0, a)$

max(0,a)

Backward:

$$\frac{dJ}{dx_i} = \frac{dJ}{dy_i} \frac{dy_i}{dx_i}$$

subderivative

where $\frac{dy_i}{dx_i} = \begin{cases} 1 & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$

# Softmax Layer

Softmax Layer

Input: $\vec{x} \in \mathbb{R}^K$   Output: $\vec{y} \in \mathbb{R}^K$

Forward:

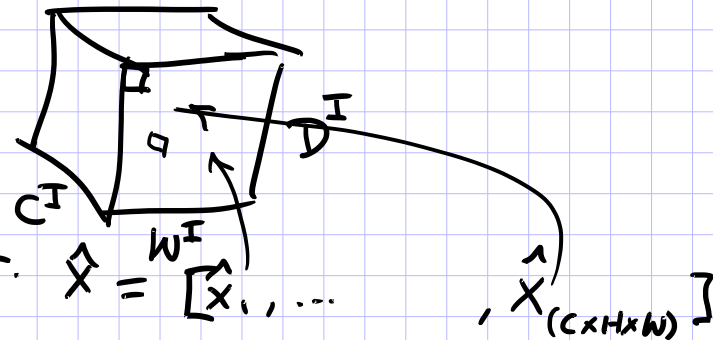$$y_i = \frac{\exp(x_i)}{\sum_{k=1}^{K} \exp(x_k)}$$

Backward:

$$\frac{dJ}{dx_j} = \sum_{i=1}^{K} \frac{dJ}{dy_i} \frac{dy_i}{dx_j}$$

where $\dfrac{dy_i}{dx_j} = \begin{cases} y_i(1-y_i) & \text{if } i=j \\ -y_i y_j & \text{otherwise} \end{cases}$

# Fully-Connected Layer

Fully Connected Layer (w/ tensor input)

- Suppose input is a 3D Tensor: $X = $



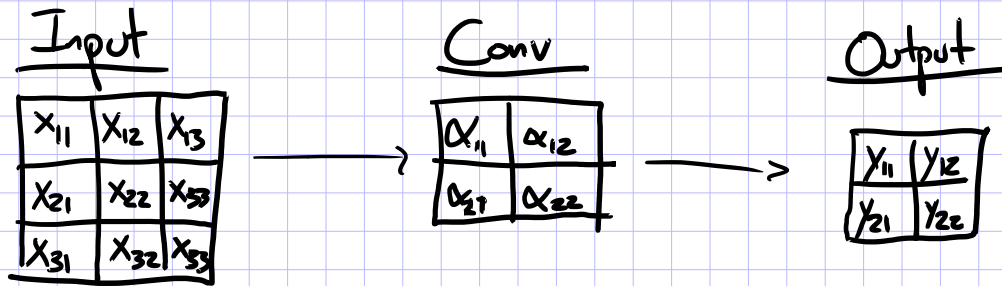- Stretch out into a long vector. $\hat{x} = [\hat{x}_1, \ldots , \hat{x}_{(C \times H \times W)}]$
- then standard linear layer:

$$y = \alpha^T \hat{x} + \alpha_0 \quad \text{where} \quad \alpha \in \mathbb{R}^{A \times B}$$
$$|\hat{x}| = A \quad , |y| = B$$

# Convolutional Layer

**Ex:** 1 input channel, 1 output channel

Input
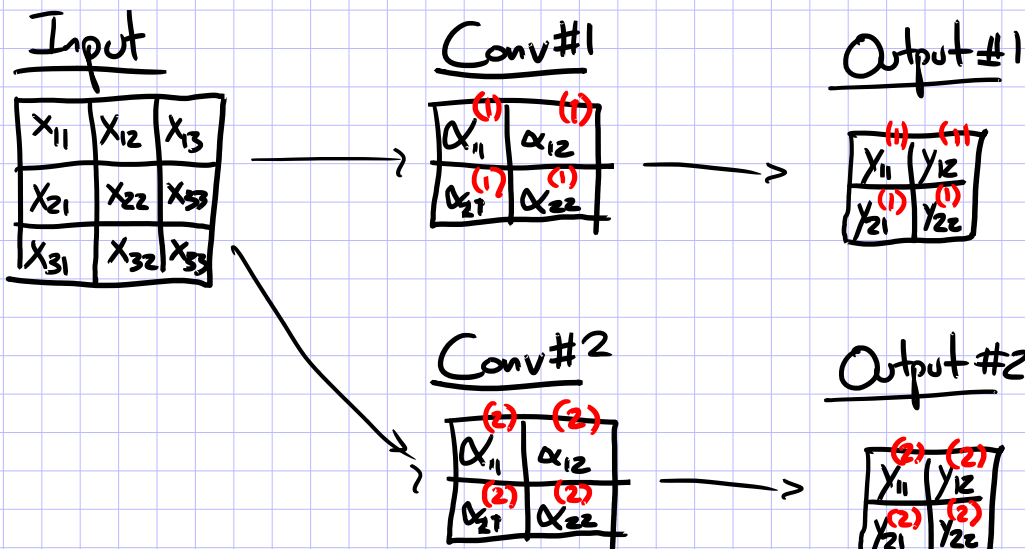
| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{53}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\rightarrow$

Conv

| $\alpha_{11}$ | $\alpha_{12}$ |
|---|---|
| $\alpha_{21}$ | $\alpha_{22}$ |

$\rightarrow$

Output

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$$y_{11} = \alpha_{11} x_{11} + \alpha_{12} x_{12} + \alpha_{21} x_{21} + \alpha_{22} x_{22} + \alpha_0$$

$$y_{12} = \alpha_{11} x_{12} + \alpha_{12} x_{13} + \alpha_{21} x_{22} + \alpha_{22} x_{23} + \alpha_0$$

$$y_{21} = \alpha_{11} x_{21} + \alpha_{12} x_{22} + \alpha_{21} x_{31} + \alpha_{22} x_{32} + \alpha_0$$

$$y_{22} = \alpha_{11} x_{22} + \alpha_{12} x_{23} + \alpha_{21} x_{32} + \alpha_{22} x_{33} + \alpha_0$$

**Ex:** 1 input channel, 2 output channels

Input

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{53}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

Conv#1

| $\alpha_{11}^{(1)}$ | $\alpha_{12}^{(1)}$ |
|---|---|
| $\alpha_{21}^{(1)}$ | $\alpha_{22}^{(1)}$ |

$\rightarrow$

Output#1

| $y_{11}^{(1)}$ | $y_{12}^{(1)}$ |
|---|---|
| $y_{21}^{(1)}$ | $y_{22}^{(1)}$ |

$$y_{11}^{(1)} = \alpha_{11}^{(1)} x_{11} + \alpha_{12}^{(1)} x_{12} + \alpha_{21}^{(1)} x_{21} + \alpha_{22}^{(1)} x_{22} + \alpha_0^{(1)}$$

$$y_{12}^{(1)} = \ldots$$

$$y_{21}^{(1)} = \ldots$$

$$y_{22}^{(1)} = \alpha_{11}^{(1)} x_{22} + \alpha_{12}^{(1)} x_{23} + \alpha_{21}^{(1)} x_{32} + \alpha_{22}^{(1)} x_{33} + \alpha_0^{(1)}$$

Conv#2

| $\alpha_{11}^{(2)}$ | $\alpha_{12}^{(2)}$ |
|---|---|
| $\alpha_{21}^{(2)}$ | $\alpha_{22}^{(2)}$ |

$\rightarrow$

Output #2

| $y_{11}^{(2)}$ | $y_{12}^{(2)}$ |
|---|---|
| $y_{21}^{(2)}$ | $y_{22}^{(2)}$ |

$$y_{11}^{(2)} = \alpha_{11}^{(2)} x_{11} + \alpha_{12}^{(2)} x_{12} + \alpha_{21}^{(2)} x_{21} + \alpha_{22}^{(2)} x_{22} + \alpha_0^{(2)}$$
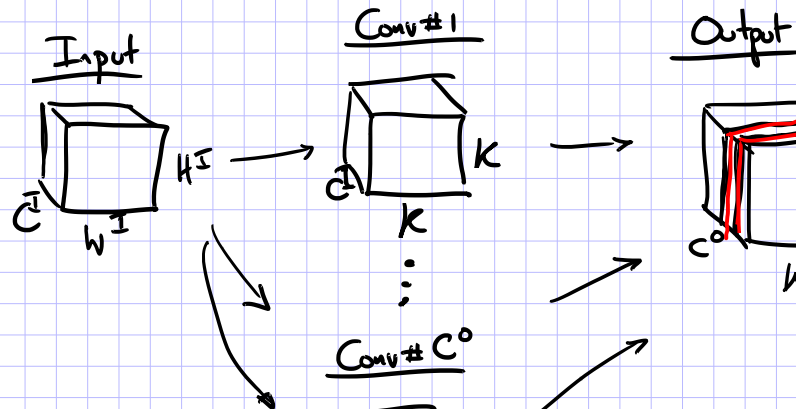
$$y_{12}^{(2)} = \ldots$$

$$y_{21}^{(2)} = \ldots$$

$$y_{22}^{(2)} = \alpha_{11}^{(2)} x_{22} + \alpha_{12}^{(2)} x_{23} + \alpha_{21}^{(2)} x_{32} + \alpha_{22}^{(2)} x_{33} + \alpha_0^{(2)}$$

# Convolutional Layer

Ex: $C^I$ input channels, $C^I$ output channels

Input



Conv#1

Conv# $C^O$

Output

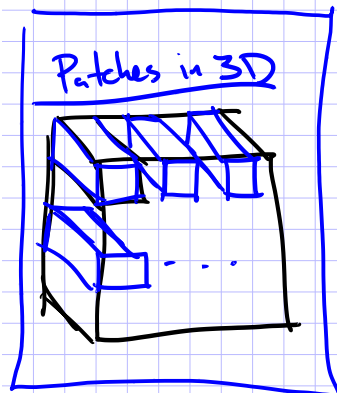*j-th slice is output from j-th convolution matrix*

$$H^O = \lfloor (H^I + 2p - K)/s + 1 \rfloor$$

$$W^O = \lfloor (W^I + 2p - K)/s + 1 \rfloor$$

where $p$ = # pixels of padding on input
$k$ = size of conv. matrix
$s$ = stride length

Patches in 3D

Forward:

$$y_{ij}^{(k)} = \alpha_0^{(k)} + \sum_{c=1}^{C^I} \sum_{q=1}^{K} \sum_{r=1}^{K} \alpha_{qr}^{(k)} x_{mn}^{(c)} \quad \text{where} \quad \begin{array}{l} m = s(i-1) + q \\ n = s(j-1) + r \end{array}$$

Backward:

$$\frac{dJ}{d\alpha_0^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_0^{(k)}}$$

$$\frac{dJ}{d\alpha_{qr}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_{qr}^{(k)}}$$

$$\frac{dJ}{dx_{mn}^{(c)}} = \sum_i \sum_j \sum_k \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(c)}}$$
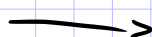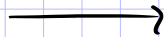
*just some calculus*

84

# Max-Pooling Layer



Ex: 1 input channel, 1 output channel, stride of 1

**Input**

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|----------|----------|----------|
| $x_{21}$ | $x_{22}$ | $x_{53}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

**Pool Size**

**Output**

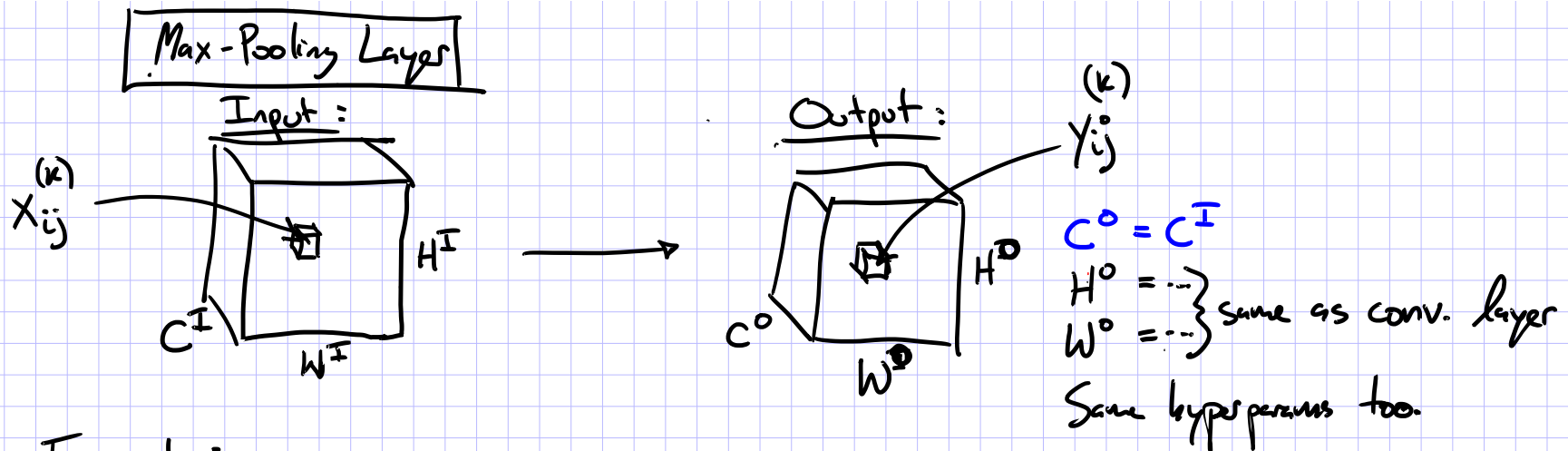| $y_{11}$ | $y_{12}$ |
|----------|----------|
| $y_{21}$ | $y_{22}$ |

$$y_{11} = \max(x_{11}, \ x_{12}, \ x_{21}, \ x_{22})$$

$$y_{12} = \max(x_{12}, \ x_{13}, \ x_{22}, \ x_{23})$$

$$y_{21} = \max(x_{21}, \ x_{22}, \ x_{31}, \ x_{32})$$

$$y_{22} = \max(x_{22}, \ x_{23}, \ x_{32}, \ x_{33})$$

# Max-Pooling Layer



Max-Pooling Layer

Input:

$X_{ij}^{(k)}$

$H^I$, $C^I$, $W^I$

Output:

$Y_{ij}^{(k)}$

$H^O$, $C^O$, $W^O$

$C^O = C^I$

$H^O = \cdots$
$W^O = \cdots$ } same as conv. layer

Same hyperparams too.

**Forward:**

$$Y_{ij}^{(k)} = \max_{\substack{q \in \{1,\dots,K\} \\ r \in \{1,\dots,K\}}} X_{mn}^{(k)} \quad \text{where} \quad \begin{array}{l} m = s(i-1)+q \\ n = s(j-1)+r \end{array}$$

**Backward:**

$$\frac{dJ}{dx_{mn}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \boxed{\frac{dy_{ij}^{(k)}}{dx_{mn}^{(k)}}}$$

**Subderivatives**

+ Max() is not differentiable, but subdifferentiable.
+ There are a __set__ of derivatives and we can just choose __one__ for SGD.

$y = \max(a,b)$

$$\Rightarrow \frac{dJ}{da} = \frac{dJ}{dy}\frac{dy}{da} \quad \text{where} \quad \frac{dy}{da} = \begin{cases} 1 & \text{if } a > b \\ 0 & \text{otherwise} \end{cases}$$

86

# Convolutional Neural Network (CNN)

- Typical layers include:
  - Convolutional layer
  - Max-pooling layer
  - Fully-connected (Linear) layer
  - ReLU layer (or some other nonlinear activation function)
  - Softmax
- These can be arranged into arbitrarily deep topologies
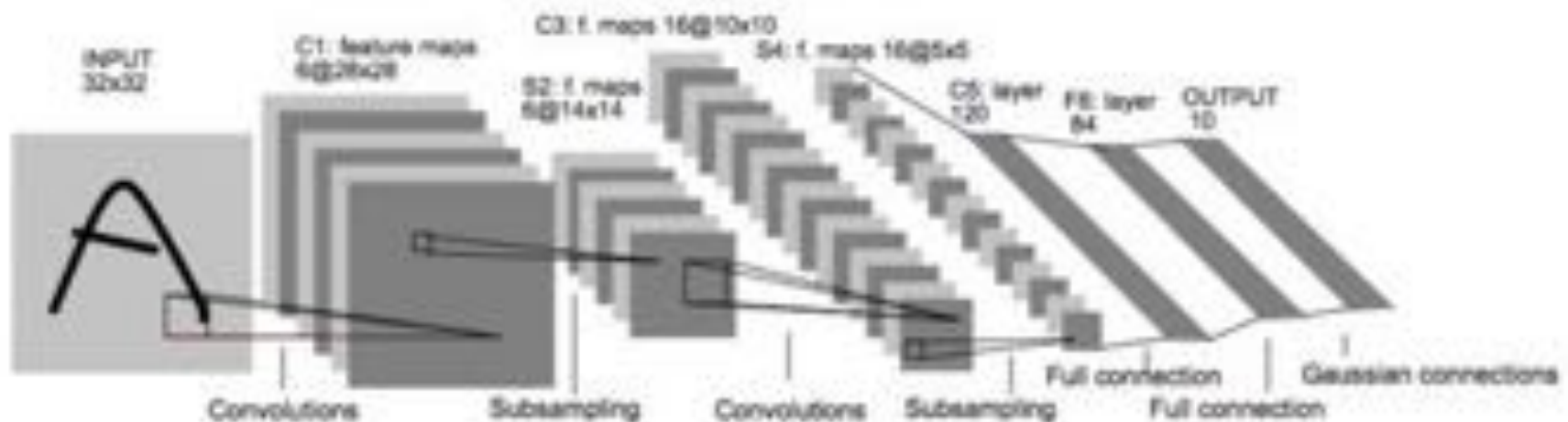
# Architecture #1: LeNet-5



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.

# Architecture #2: AlexNet

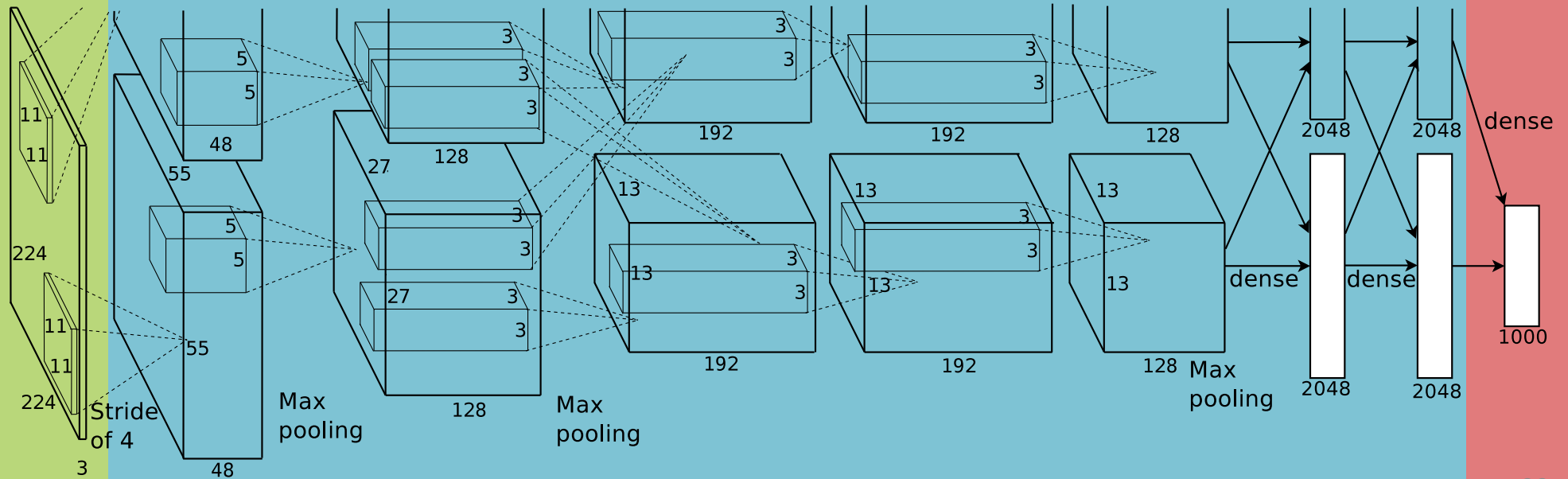**CNN for Image Classification**
(Krizhevsky, Sutskever & Hinton, 2012)
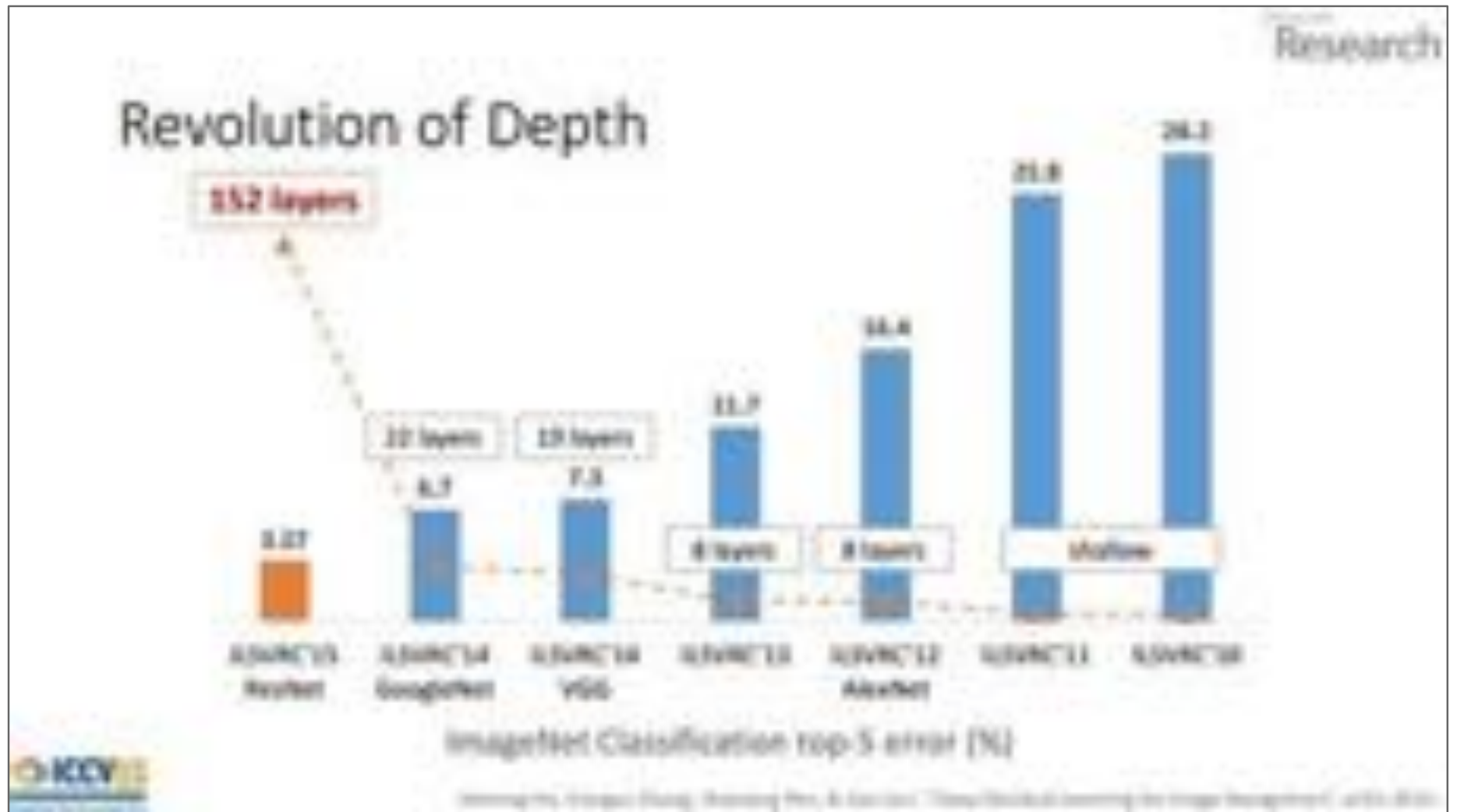15.3% error on ImageNet LSVRC-2012 contest

Input image (pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax

# CNNs for Image Recognition

Slide from Kaiming He

The key building block of ResNet

# RESIDUAL CONNECTIONS

# Slides in this section from...

Microsoft Research

# Deep Residual Learning

MSRA @ ILSVRC & COCO 2015 competitions

Kaiming He

with Xiangyu Zhang, Shaoqing Ren, Jifeng Dai, & Jian Sun

Microsoft Research Asia (MSRA)

ICCV15
International Conference on Computer Vision

# Revolution of Depth

**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

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Revolution of Depth

**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, 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 |

**GoogleNet, 22 layers**
**(ILSVRC 2014)**



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)

VGG, 19 layers
(ILSVRC 2014)

ResNet, 152 layers
(ILSVRC 2015)
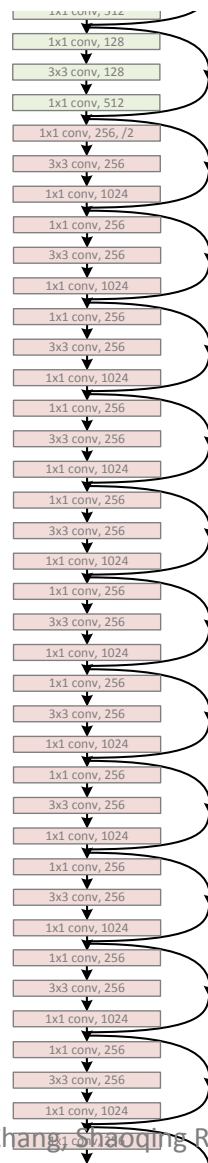
Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Revolution of Depth

ResNet, 152 layers



7x7 conv, 64, /2, pool/2
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x1 conv, 128, /2
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
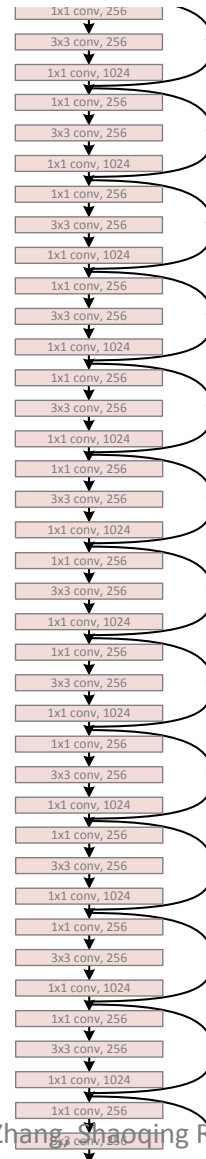1x1 conv, 512
1x1 conv, 256, /2
3x3 conv, 256

(there was an animation here)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
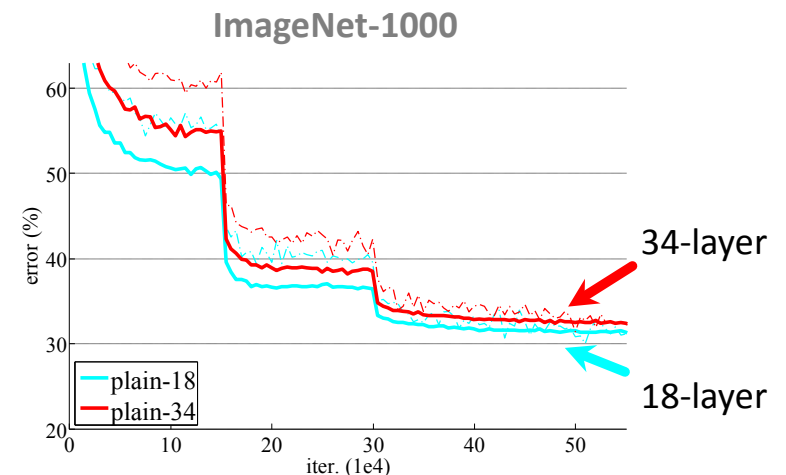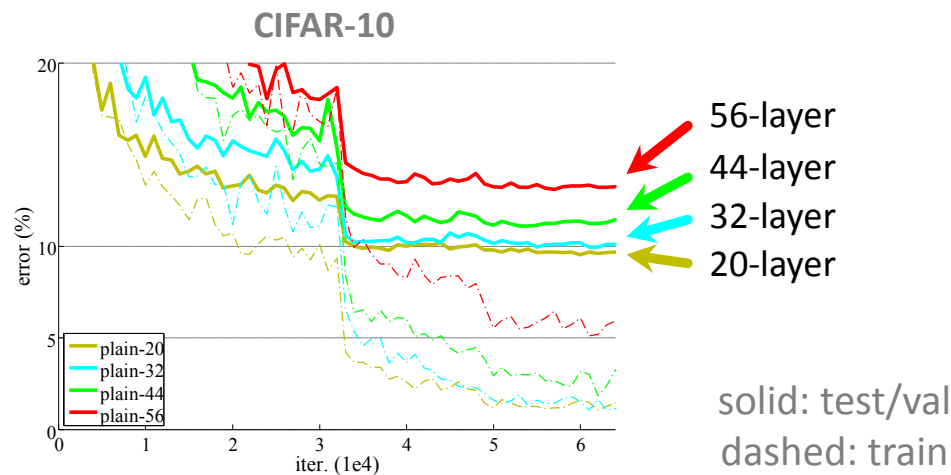
# Revolution of Depth

ResNet, 152 layers



(there was an animation here)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Revolution of Depth

ResNet, 152 layers



(there was an animation here)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
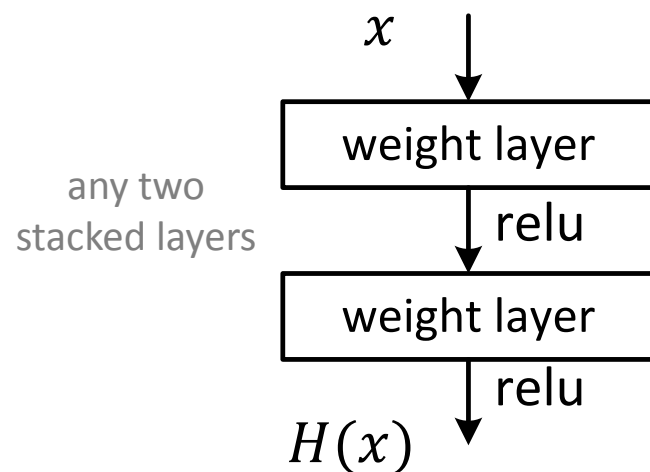
# Revolution of Depth

ResNet, 152 layers



(there was an animation here)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Simply stacking layers?



CIFAR-10

56-layer
44-layer
32-layer
20-layer

solid: test/val
dashed: train

ImageNet-1000

34-layer
18-layer

- "Overly deep" plain nets have **higher training error**
- A general phenomenon, observed in many datasets

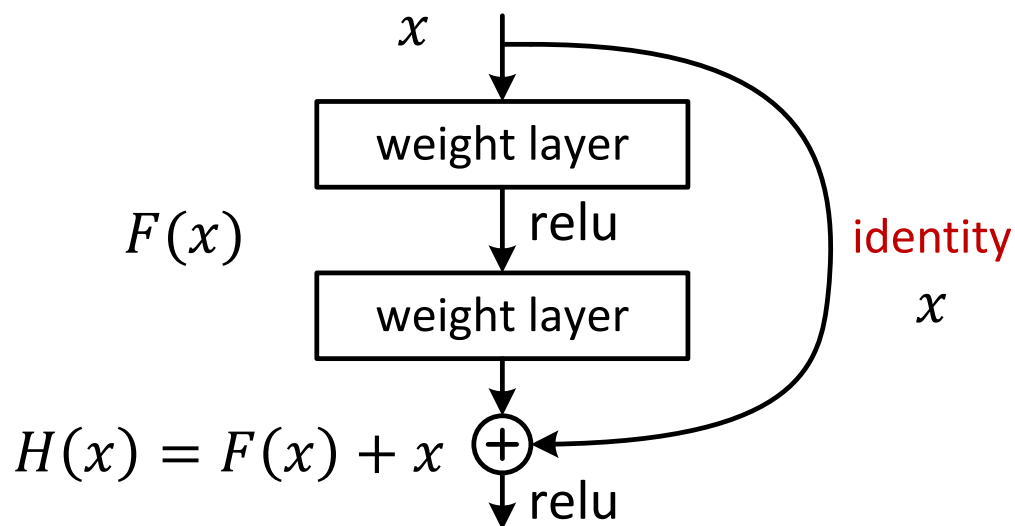Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Deep Residual Learning

- Plaint net

$x$

| weight layer |

relu

| weight layer |

relu

$H(x)$

any two
stacked layers

$H(x)$ is any desired mapping,

hope the 2 weight layers fit $H(x)$

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Deep Residual Learning

- **Residual** net



$x$

weight layer

$F(x)$    relu

weight layer

identity

$x$

$H(x) = F(x) + x$ ⊕

relu

$H(x)$ is any desired mapping,

~~hope the 2 weight layers fit $H(x)$~~

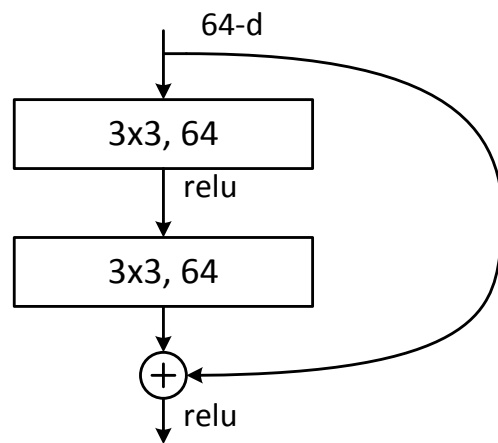hope the 2 weight layers fit $F(x)$

let $H(x) = F(x) + x$

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Deep Residual Learning

- $F(x)$ is a residual mapping w.r.t. identity



$$x$$

weight layer

relu

$F(x)$
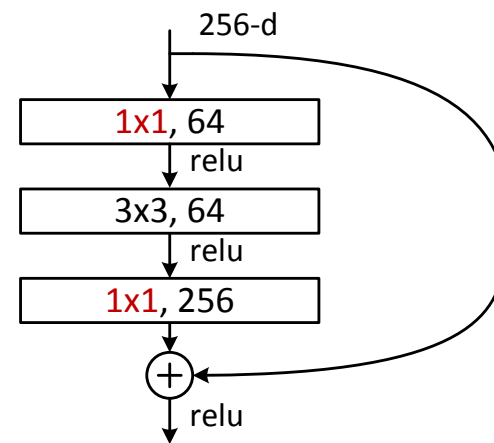
weight layer

identity

$x$

$$H(x) = F(x) + x \quad \oplus$$

relu

- If identity were optimal, easy to set weights as 0

- If optimal mapping is closer to identity, easier to find small fluctuations

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# ImageNet experiments

- A practical design of going deeper



all-3x3 ⟷ similar complexity ⟷ bottleneck
(for ResNet-50/101/152)

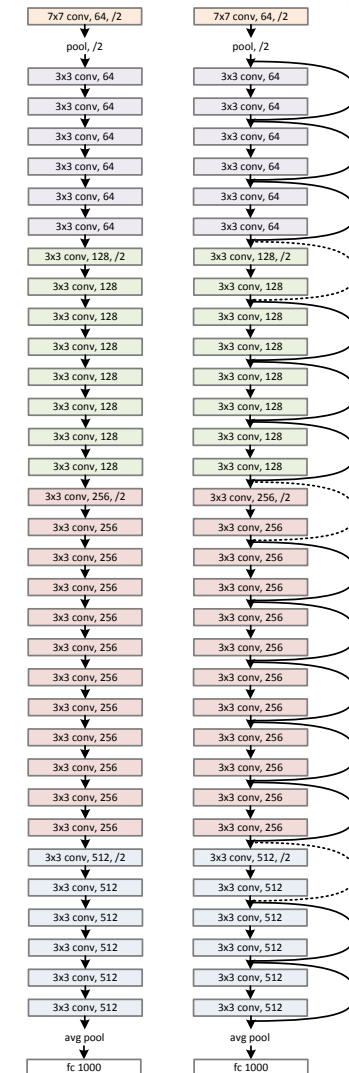Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Network "Design"

- Keep it simple

- Our basic design (VGG-style)
  - all 3x3 conv (almost)
  - spatial size /2 => # filters x2
  - Simple design; just deep!

- Other remarks:
  - no max pooling (almost)
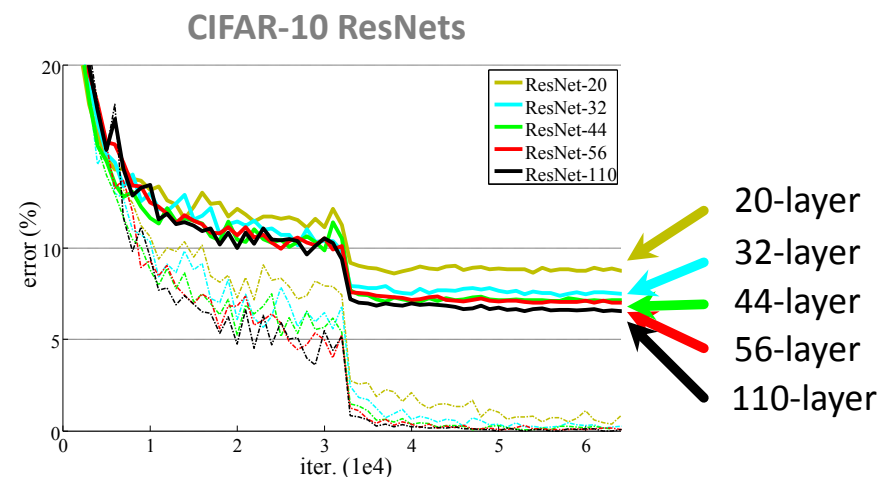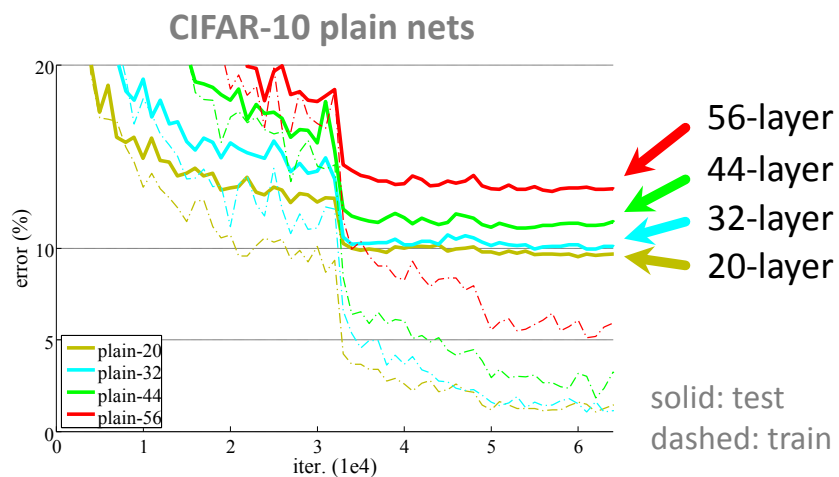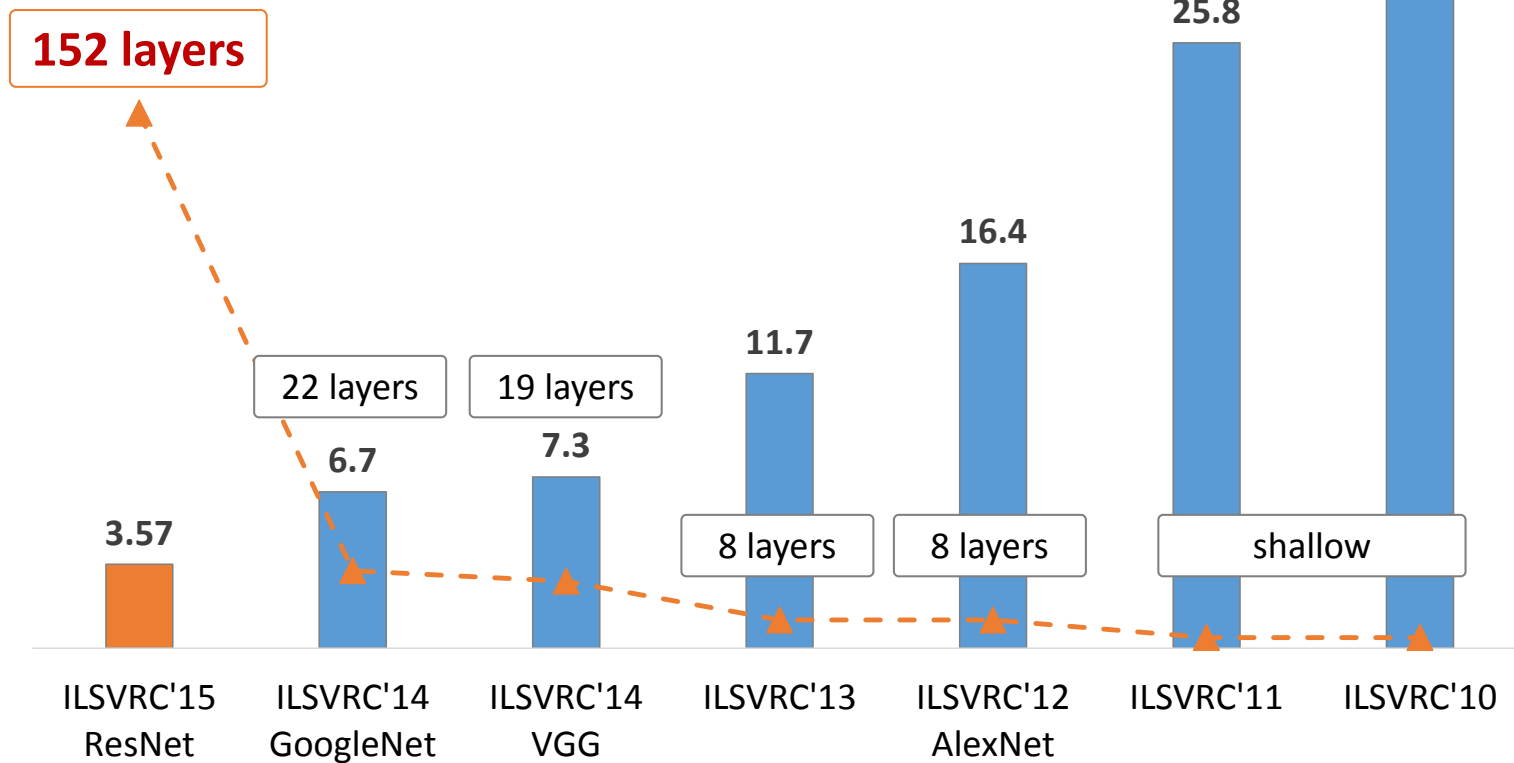  - no hidden fc
  - no dropout

plain net

ResNet

# CIFAR-10 experiments



- Deep ResNets can be trained without difficulties
- Deeper ResNets have **lower training error**, and also lower test error

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# ImageNet experiments
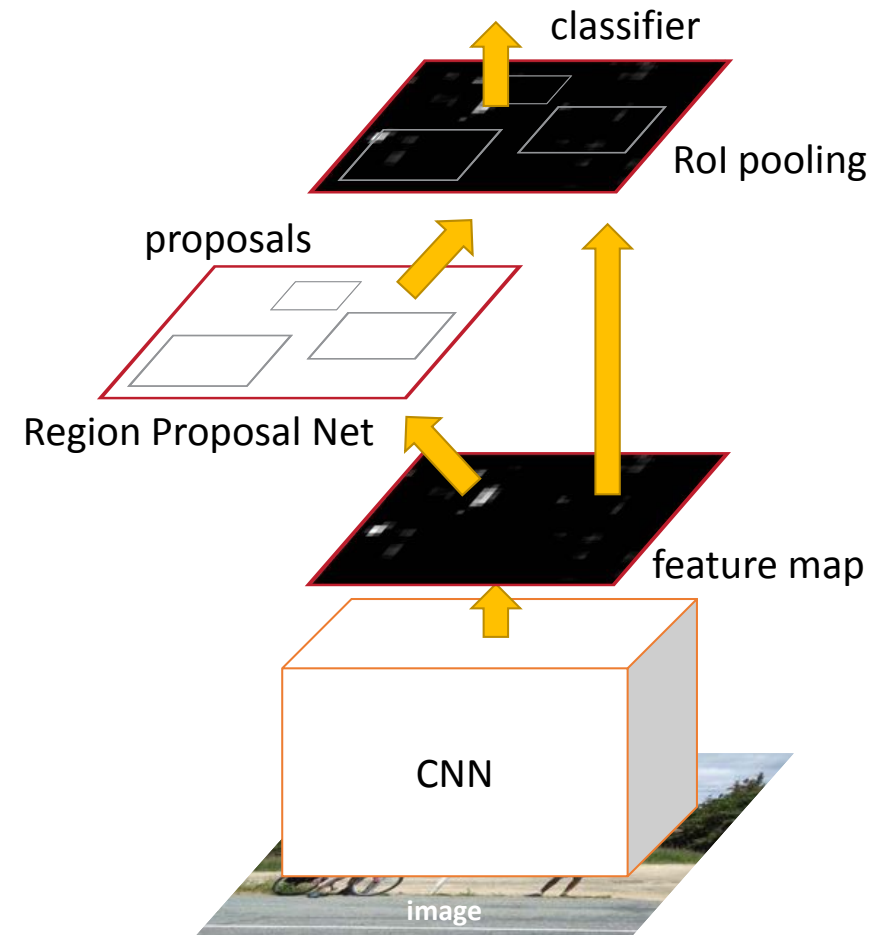


ImageNet Classification top-5 error (%)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.

# Object Detection (brief)

- Simply "Faster R-CNN + ResNet"

| Faster R-CNN baseline | mAP@.5 | mAP@.5:.95 |
|---|---|---|
| VGG-16 | 41.5 | 21.5 |
| ResNet-101 | **48.4** | **27.2** |

COCO detection results
(ResNet has 28% relative gain)

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.
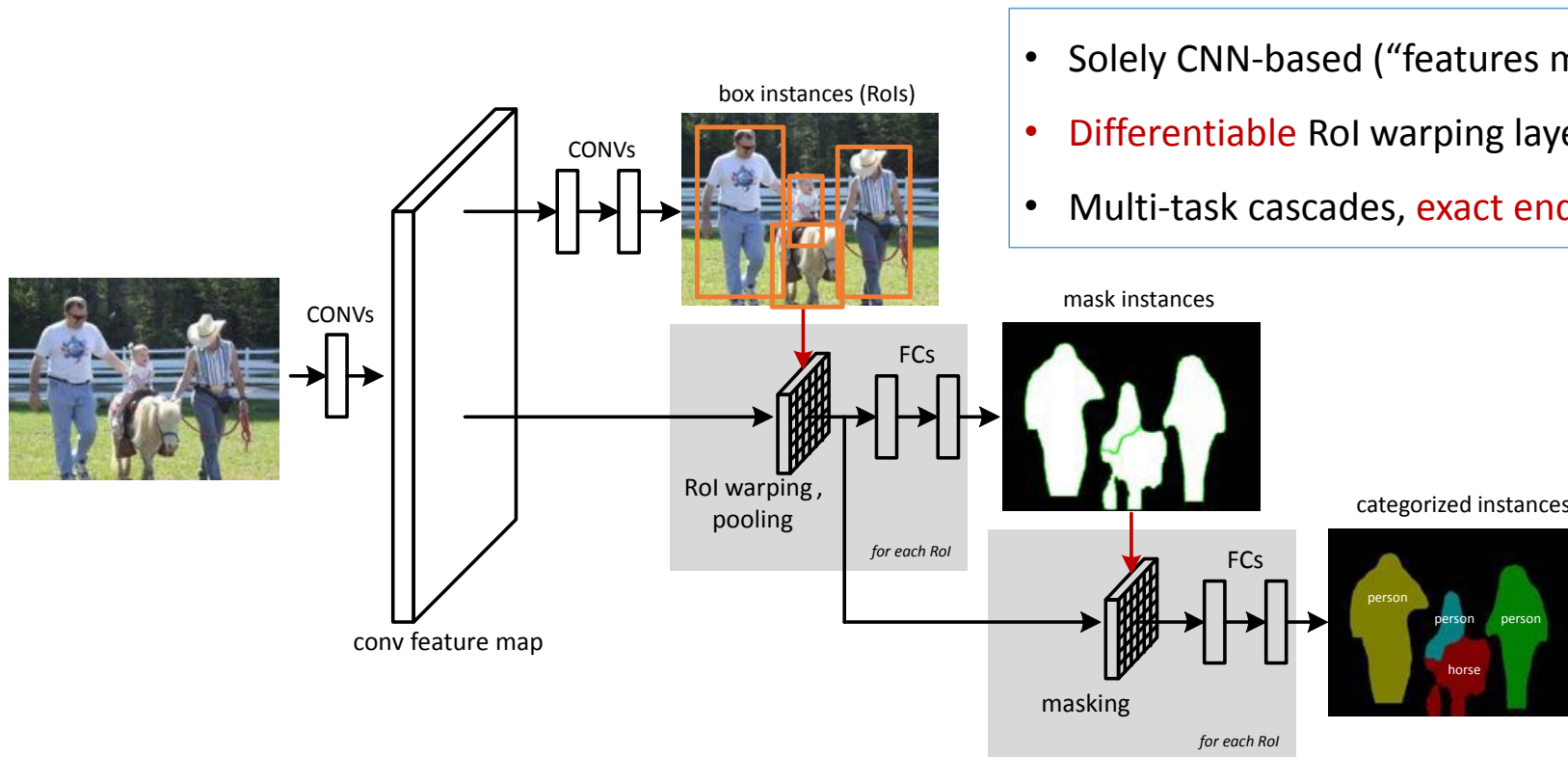
Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.
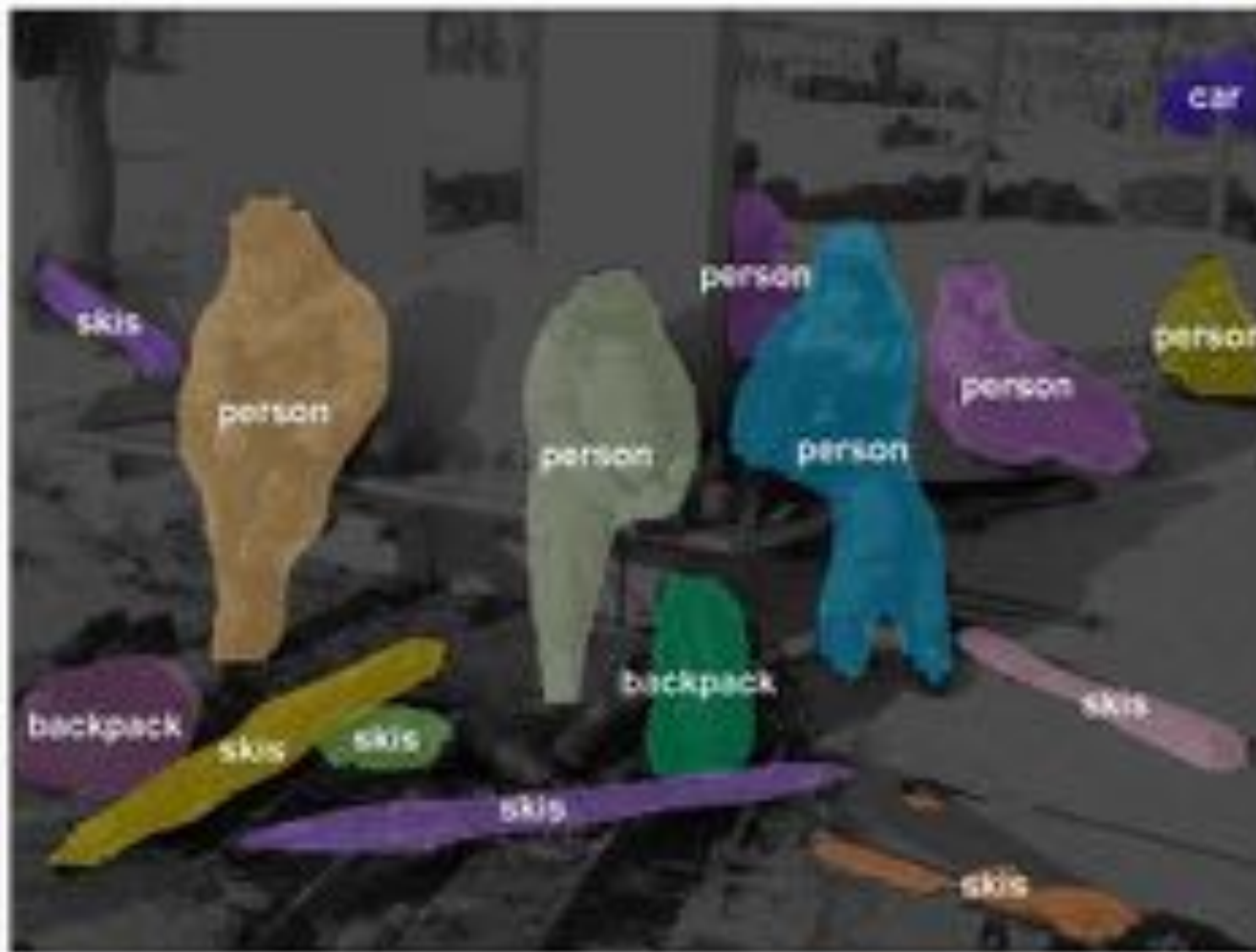
# Instance Segmentation (brief)



- Solely CNN-based ("features matter")
- Differentiable RoI warping layer (w.r.t box coord.)
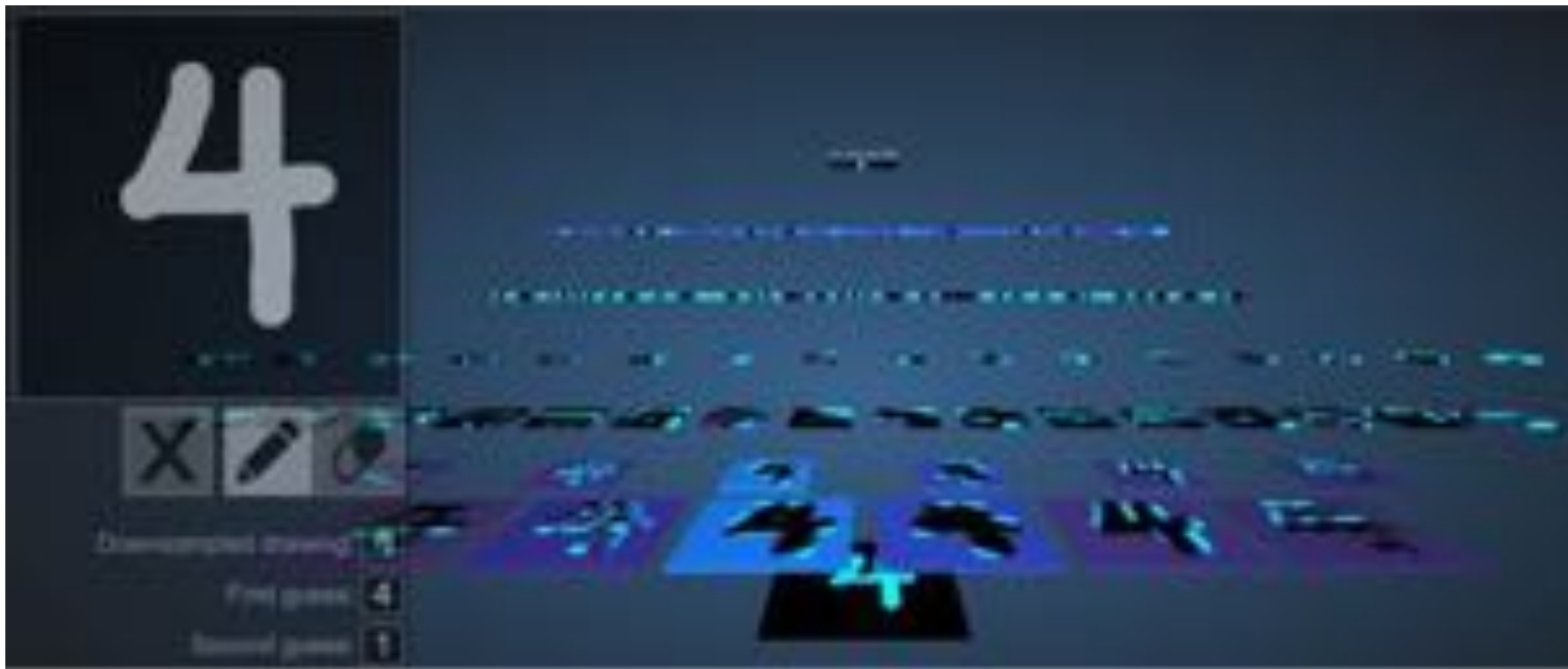- Multi-task cascades, exact end-to-end training

box instances (RoIs)

CONVs

CONVs

conv feature map

RoI warping , pooling

*for each RoI*

FCs

mask instances

FCs

masking

*for each RoI*

categorized instances

person · person · person · horse

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Jifeng Dai, Kaiming He, & Jian Sun. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". arXiv 2015.

input

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Jifeng Dai, Kaiming He, & Jian Sun. "Instance-aware Semantic Segmentation via Multi-task Network Cascades". arXiv 2015.

# CNN VISUALIZATIONS

# 3D Visualization of CNN

http://scs.ryerson.ca/~aharley/vis/conv/
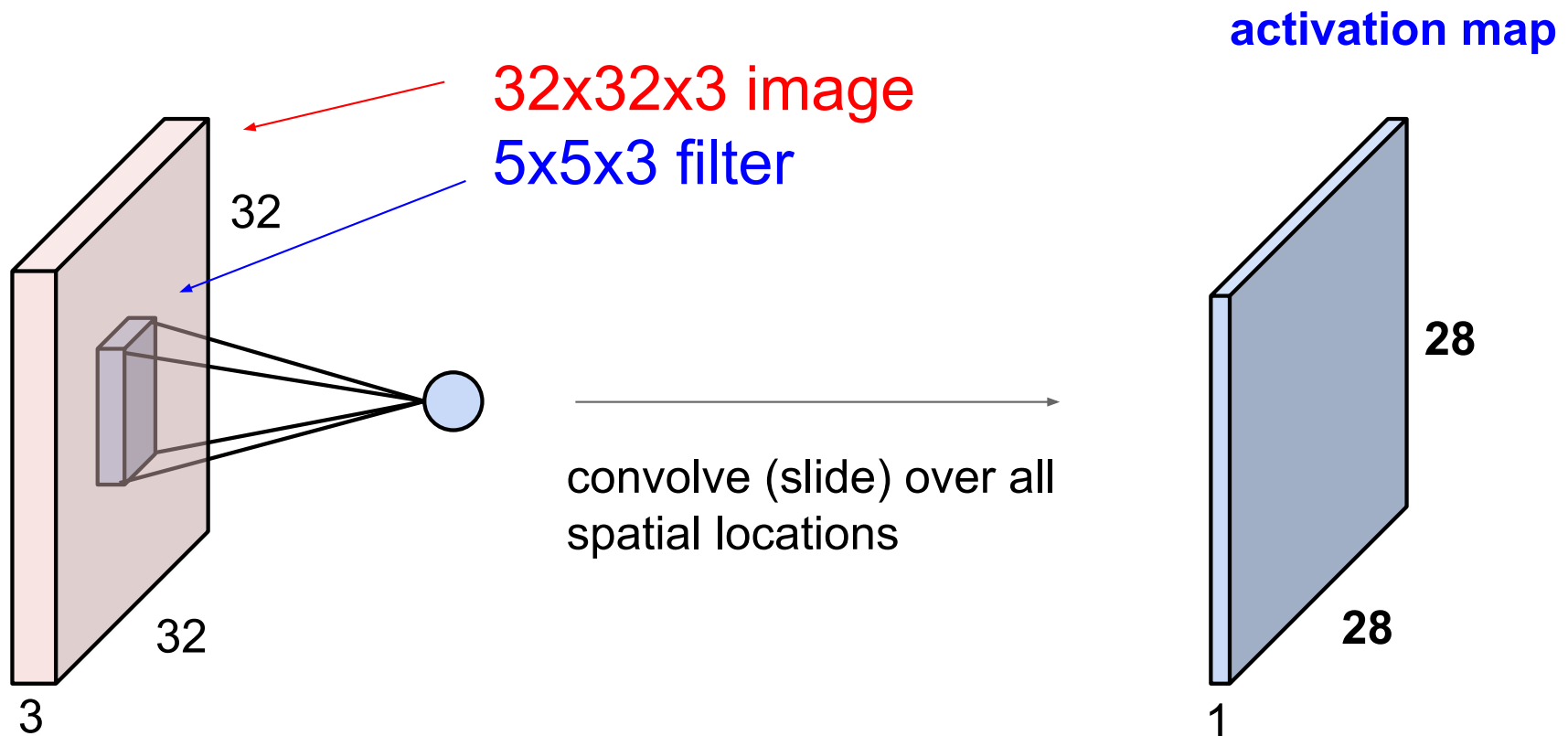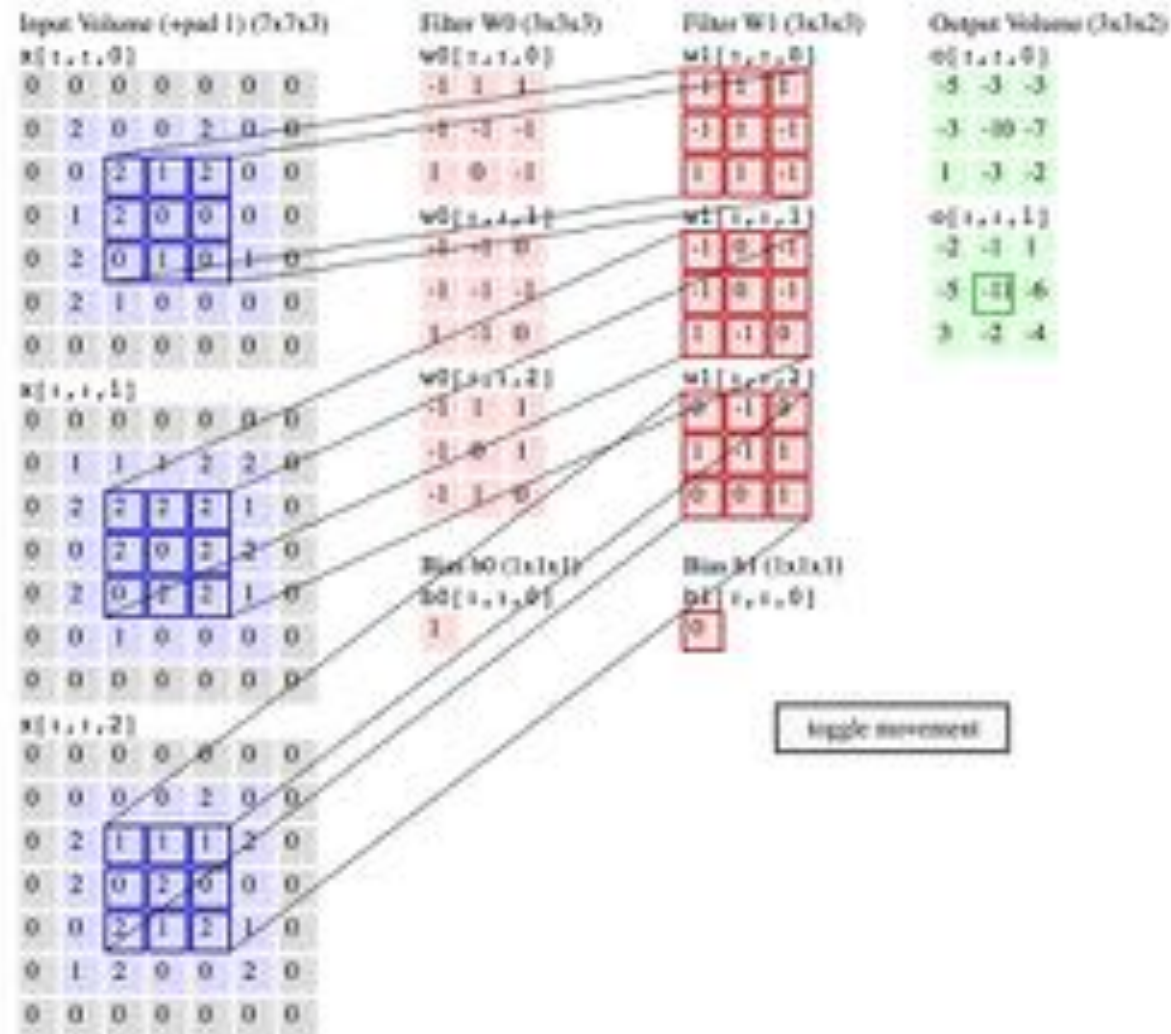
# Convolution of a Color Image

- Color images consist of 3 floats per pixel for RGB (red, green blue) color values
- Convolution must also be 3-dimensional



**activation map**

32x32x3 image
5x5x3 filter

convolve (slide) over all spatial locations

# Animation of 3D Convolution

http://cs231n.github.io/convolutional-networks/



Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)

# MNIST Digit Recognition with CNNs (in your browser)

https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html

Figure from Andrej Karpathy

# CNN Summary

**CNNs**

- Are used for all aspects of **computer vision**, and have won numerous pattern recognition competitions

- Able learn **interpretable features** at different levels of abstraction

- Typically, consist of **convolution** layers, **pooling** layers, **nonlinearities**, and **fully connected** layers