

# Convolutional Neural Network for Computer Vision Feature Learning

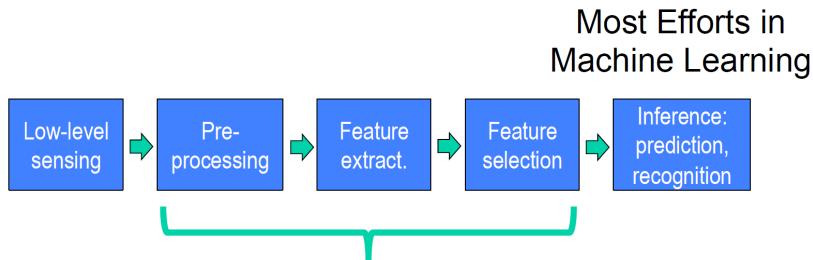
Dequan Wang

Fudan University

*dqwang12@fudan.edu.cn*

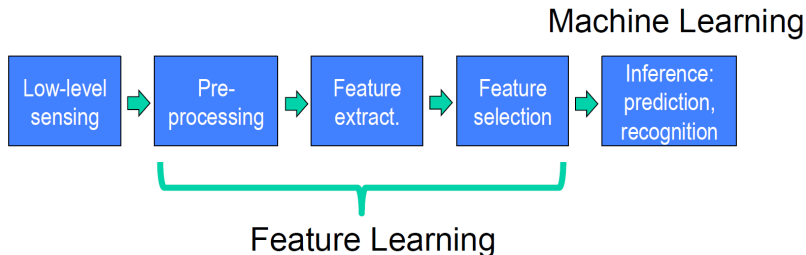
September 15, 2014

# Pipeline of Machine Visual Perception



- Most critical for accuracy
- Account for most of the computation for testing
- Most time-consuming in development cycle
- Often hand-craft in practice

# Learning with Structures

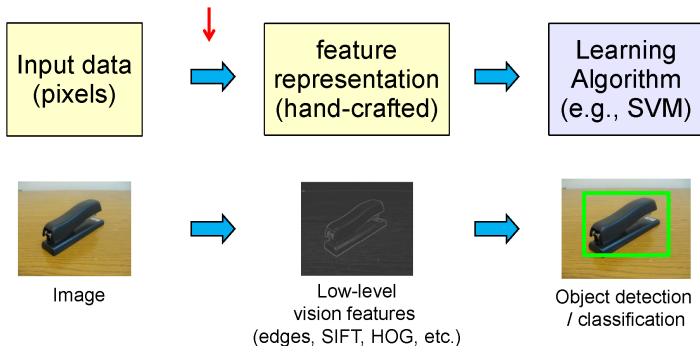


- ▶ Kernel Learning
- ▶ Transfer Learning
- ▶ Semi-Supervised Learning
- ▶ Manifold Learning
- ▶ Sparse Learning
- ▶ Structured Input-Output Prediction

# Traditional Recognition System

- ▶ Features are key to recent progress in recognition
- ▶ Multitude of hand-designed features currently in use

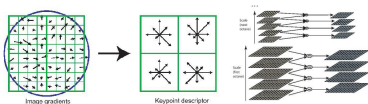
Features are not learned



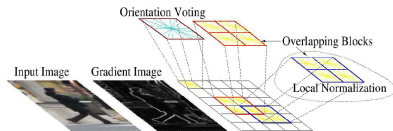


# Traditional Visual Feature

- Where next? Better classifiers or better feature?



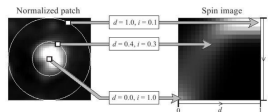
SIFT



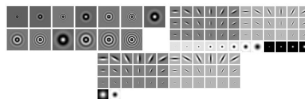
HoG

and many others:

SURF, MSER, LBP, Color-SIFT, Color histogram, GLOH, .....



Spin image

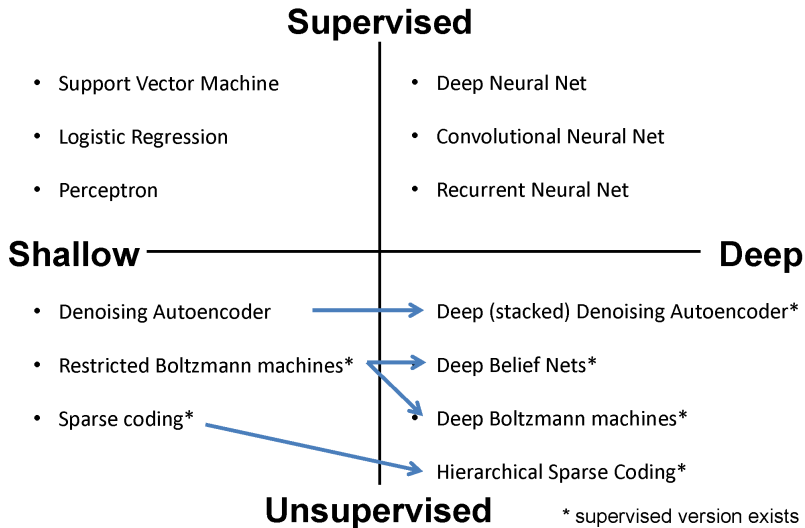


Textons

# Approaches to Learning Features

- Supervised Learning
  - End-to-end learning of deep architectures (e.g., deep neural networks) with back-propagation
  - Works well when the amounts of labels is large
  - Structure of the model is important (e.g. convolutional structure)
- Unsupervised Learning
  - Learn statistical structure or dependencies of the data from unlabeled data
  - Layer-wise training
  - Useful when the amount of labels is not large

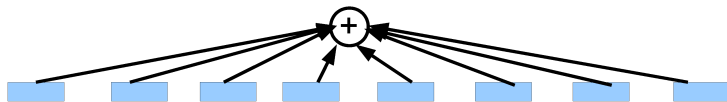
# Representation Learning



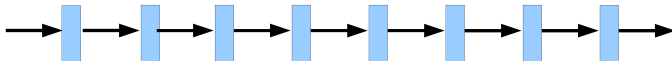
# Learning Non-Linear Features

Given a dictionary of simple non-linear functions:  $g_1, \dots, g_n$

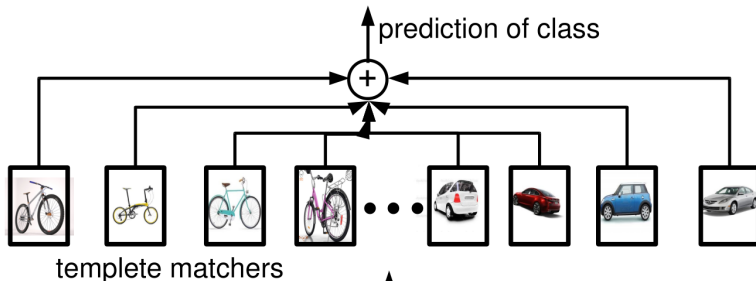
**Proposal #1: linear combination**  $f(x) \approx \sum_j g_j$



**Proposal #2: composition**  $f(x) \approx g_1(g_2(\dots g_n(x)\dots))$



# Linear Combination

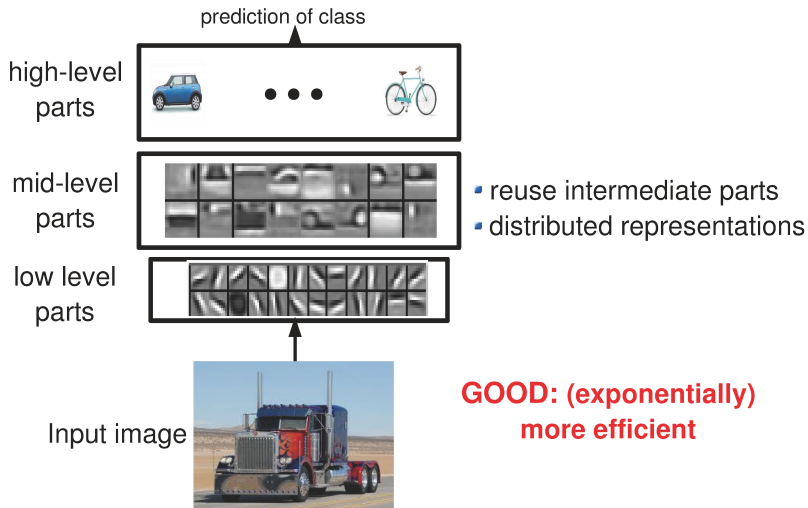


**BAD: it may require  
an exponential nr. of  
templates!!!**



Input image

# Composition



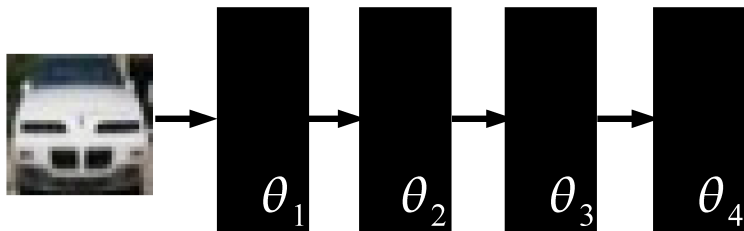
# Ranzato's Definition

**Deep Learning** is a method which makes predictions by using a sequence of non-linear processing stages. The resulting intermediate representations can be interpreted as feature hierarchies and the whole system is jointly learned from data.

- ▶ Some deep learning methods are probabilistic, others are loss-based, some are supervised, others are unsupervised ...

# Deep Learning in Practice

Optimization is easy, need to know a few tricks of the trade.



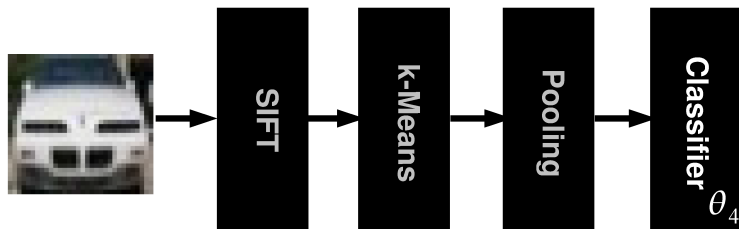
**Q:** What's the feature extractor? And what's the classifier?

**A:** No distinction, end-to-end learning!



# A Potential Problem with Deep Learning

Optimization is difficult: non-convex, non-linear system



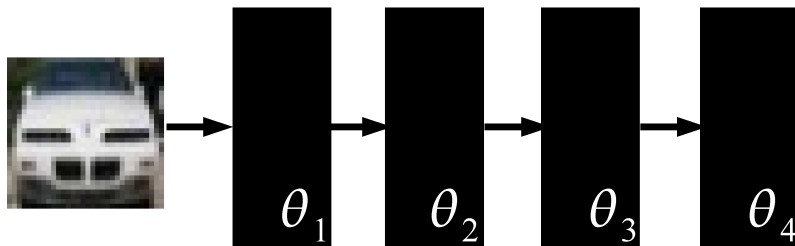
**Solution #1:** freeze first N-1 layer (engineer the features)

It makes it **shallow**!

- How to design features of features?
- How to design features for new imagery?

# A Potential Problem with Deep Learning

Optimization is difficult: non-convex, non-linear system



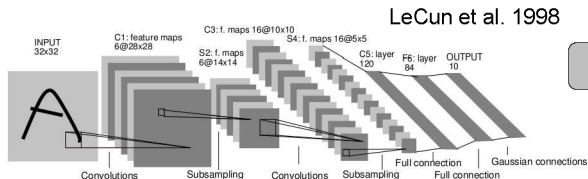
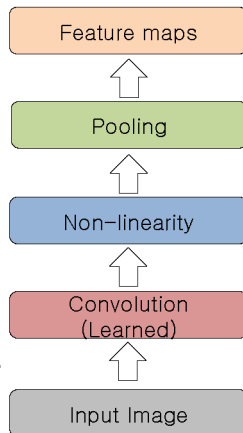
**Solution #2:** live with it!

It will converge to a local minimum.  
It is much more powerful!!

*Given lots of data, engineer less and learn more!!  
Just need to know a few tricks of the trade...*

# Convolutional Neural Network

- Feed-forward:
  - Convolve input
  - Non-linearity (rectified linear)
  - Pooling (local max)
- Supervised
- Train convolutional filters by back-propagating classification error



# Generalization Error Decomposition

- ▶ Given observations  $(X_i, Y_i), i = 1, \dots, n$
- ▶ Learn a predictive function  $f(X)$
- ▶ Generalization error  $E[L(f(X), Y)]$

$$E[L(f(X), Y)] = A + E + O$$

- ▶ Approximation error  $\rightarrow$  model class
- ▶ Estimation error  $\rightarrow$  data size
- ▶ Optimization error  $\rightarrow$  algorithm

# Generalization Error Decomposition

$$E[L(f(X), Y)] = A + E + O$$

- ▶ Approximation error → use complex model
- ▶ Estimation error → collect big data
- ▶ Optimization error → design an OK algorithm

## Trend of Machine Learning

- ▶ Approximation error: simple model → complex model
- ▶ Estimation error: small data → big data
- ▶ Optimization error: fine algorithm → scalable algorithm

# Why Today

- Non-convex & non-linear
- Intense computation
- Sensitive to initialization
- Overfitting
- Vanishing gradient



- Big data
- GPU
- Large scale parallel computation
- Layer-wise pre-training
- RELU, drop-out, better normalization, etc.

# Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton\* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

**D**imensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which

finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer “encoder” network

# Benchmark for ImageNet



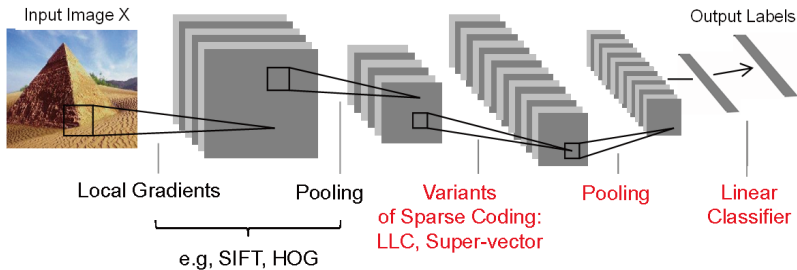
**72%, 2010**

**74%, 2011**

**85%, 2012**

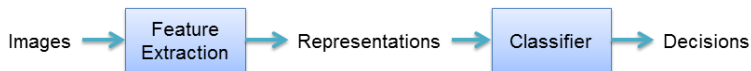


# Traditional System on ImageNet



# Drawback of Traditional System

- All these methods are not scalable
  - Handcrafted descriptors, time consuming, limited prior knowledge
  - Representations by unsupervised learning, not tunable
  - Cannot fully utilize large-scale training data



- We need to enable **end-to-end** learning.

# Deep Learning for Vision@CVPR 2012 Tutorial

## Buiding blocks

- ▶ RBMs, Autoencoder, Sparse Coding

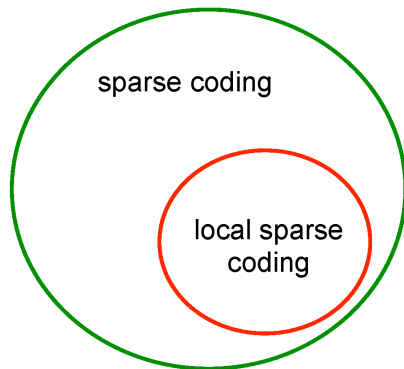
## Go deeper: Layerwise feature learning

- ▶ Layer-by-layer unsupervised training
- ▶ Layer-by-layer supervised training

## Fine tuning via Backpropogation

- ▶ If data are big enough, direct fine tuning is enough
- ▶ Sparsity on hidden layers are often useful

# Sparsity vs. Locality



- **Intuition:** similar data should get similar activated features
- **Local sparse coding:**
  - data in the same neighborhood tend to have shared activated features;
  - data in different neighborhoods tend to have different features activated.

# Challenge to Deep Learners

## Key Issues

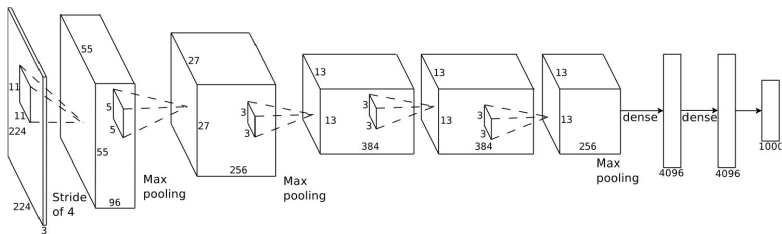
- ▶ What if no hand-craft features at all?
- ▶ What if use much deeper neural networks?

## Answer from Geoff Hinton

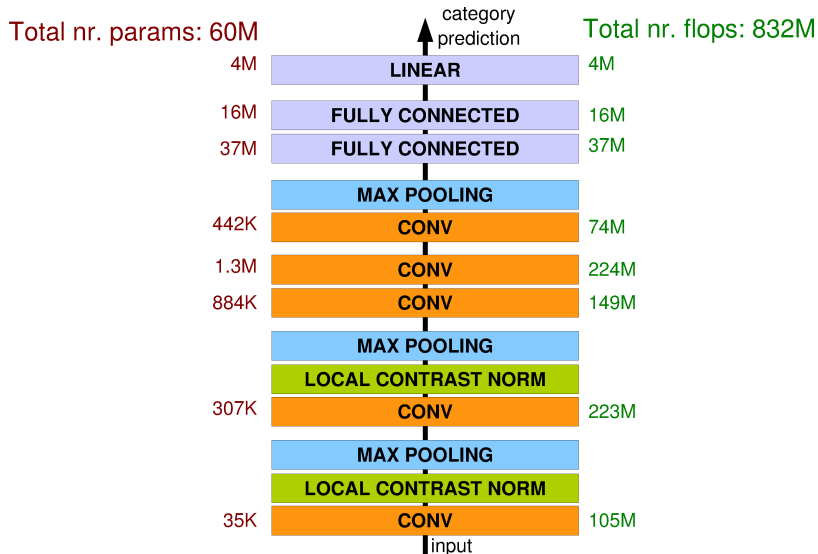
- ▶ Our chief critic, Jitendra Malik, has said that this competition is a good test of whether deep neural networks really do work well for object recognition.

# The Architecture

- 7 hidden layers not counting max pooling.
- Early layers are conv., last two layers globally connected.
- Uses rectified linear units in every layer.
- Uses competitive normalization to suppress hidden activities.



# The Architecture

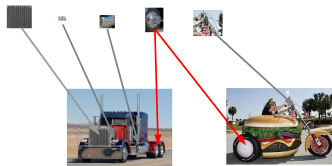


# Why Deep

- ▶ When input has hierarchical structure, the use of a hierarchical architecture is potentially more efficient because intermediate computations can be re-used.
- ▶ Architectures are efficient also because they use distributed representations which are shared across classes.

[1 1 0 0 0 1 0 **1** 0 0 0 0 1 1 0 1 ...] motorbike

[0 0 1 0 0 0 0 **1** 0 0 1 1 0 0 1 0 ...] truck



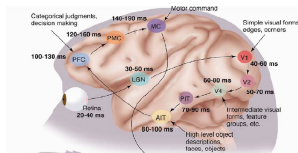


# Why Hierarchy

## Theoretical

- ▶ One limitation was based on the well-known depth-breadth tradeoff in circuits design Hastad [1987].
- ▶ This suggests that many functions can be much more efficiently represented with deeper architectures, often with a modest number of levels (e.g., logarithmic in the number of inputs).

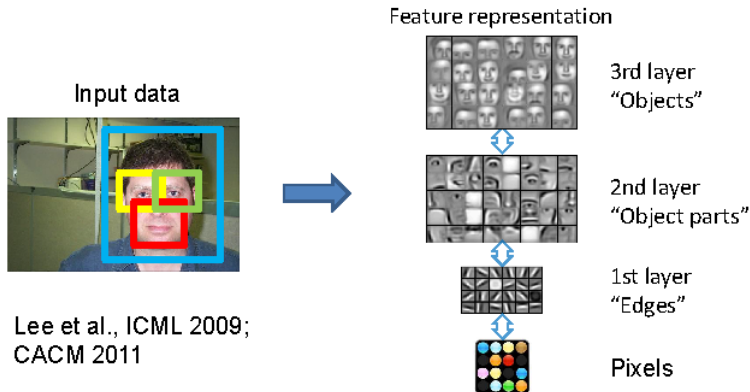
## Biological



- ▶ Visual cortex is hierarchical (Hubel-Wiesel Model)

# Learning Feature Hierarchy

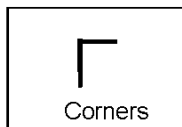
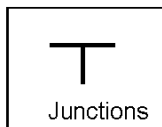
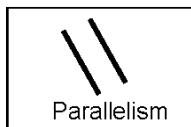
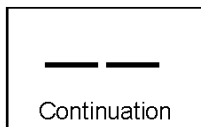
1. Learn **useful higher-level features** from images



2. **Fill in representation gap** in recognition

# Mid-Level Representation

- Mid-level cues



“Tokens” from Vision by D.Marr:

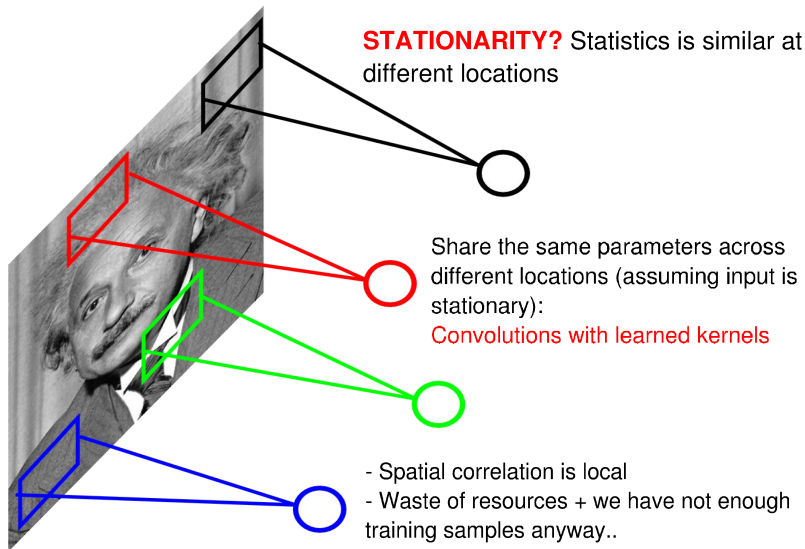


- Object parts:

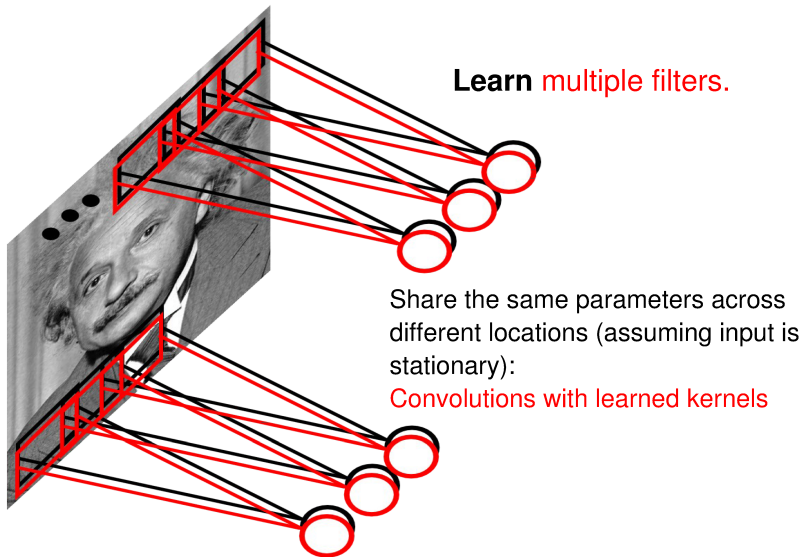


- Difficult to hand-engineer → What about learning them?

# Locally Connected Layer



# Convolutional Layer



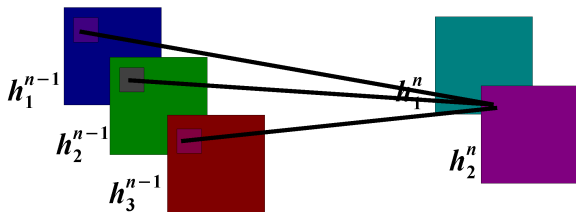
# Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output  
feature map

input feature  
map

kernel



# Convolutional Layer

A standard neural net applied to images:

- ▶ scales quadratically with the size of the input
- ▶ does not leverage stationarity

Solution:

- ▶ connect each hidden unit to a small patch of the input
- ▶ share the weight across space

# Convolutional Layer

## What is the size of the output?

It is proportional to the number of filters and depends on the stride. If kernels have size  $K \times K$ , input has size  $D \times D$ , stride is 1, and there are  $M$  input feature maps and  $N$  output feature maps then:

- ▶ the input has size  $M \times D \times D$
- ▶ the output has size  $N \times (D - K + 1) \times (D - K + 1)$
- ▶ the kernels have  $M \times N \times K \times K$  coefficients

## How many feature maps?

- ▶ Usually, there are more output feature maps than input feature maps. Convolutional layers can increase the number of hidden units by big factors

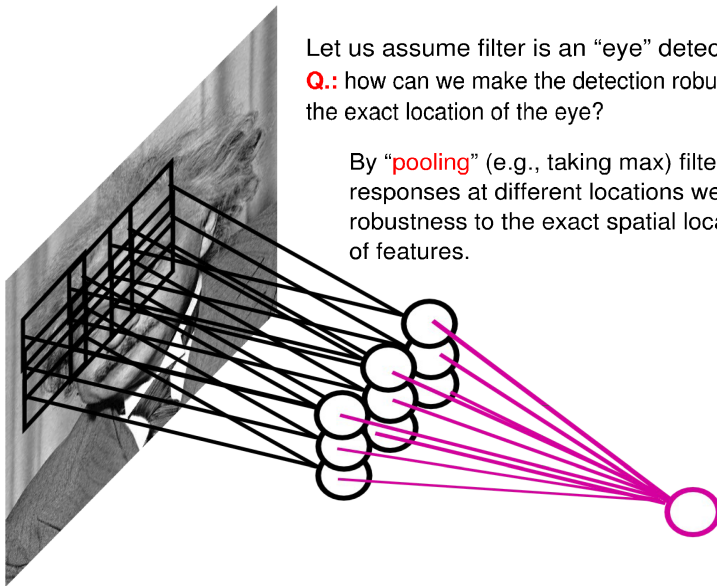


# Pooling Layer

Let us assume filter is an “eye” detector

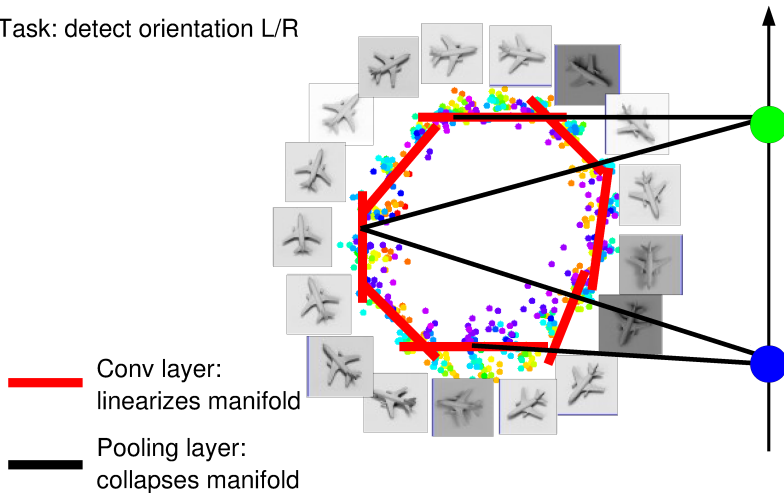
**Q.:** how can we make the detection robust to the exact location of the eye?

By “pooling” (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.

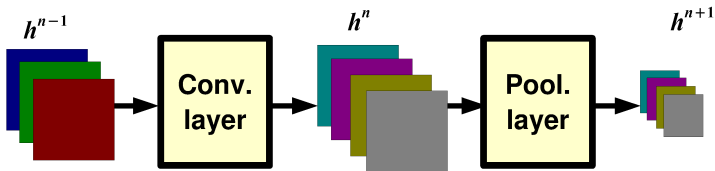


# Pooling Layer Interpretation

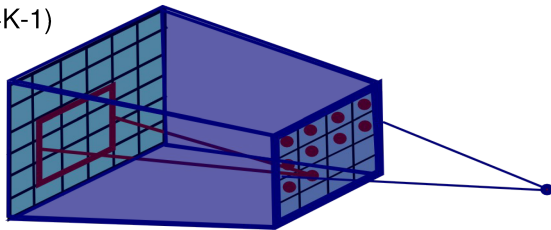
Task: detect orientation L/R



# Pooling Layer Receptive Field Size



If convolutional filters have size  $K \times K$  and stride 1, and pooling layer has pools of size  $P \times P$ , then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:  $(P+K-1) \times (P+K-1)$



# Pooling Layer Examples

Max-pooling:

$$h_j^n(x, y) = \max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_j^n(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})$$

L2-pooling:

$$h_j^n(x, y) = \sqrt{\sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_j^{n-1}(\bar{x}, \bar{y})^2}$$

L2-pooling over features:

$$h_j^n(x, y) = \sqrt{\sum_{k \in N(j)} h_k^{n-1}(x, y)^2}$$

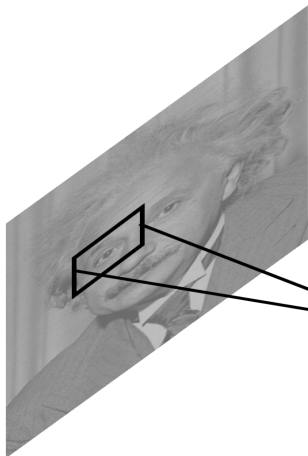
# Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$

Performed also across features  
and in the higher layers..

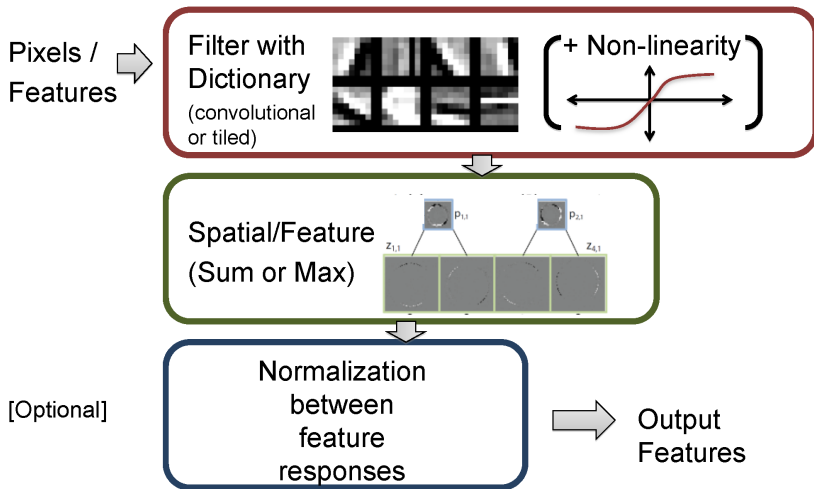
Effects:

- improves invariance
- improves optimization
- increases sparsity



**Note:** computational cost is  
negligible w.r.t. conv. layer.

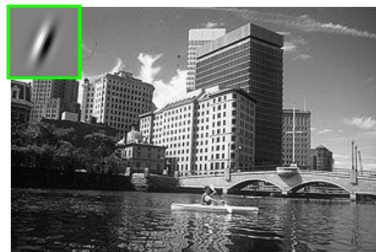
# Components of Each Layer



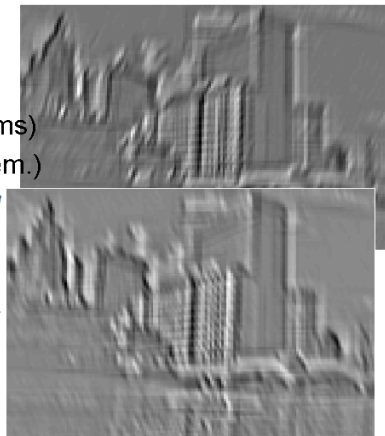
# Filtering

- **Convolutional**

- Dependencies are local
- Translation equivariance
- Tied filter weights (few params)
- Stride 1,2,... (faster, less mem.)



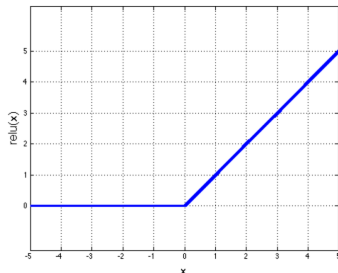
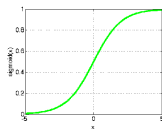
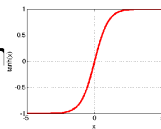
Input



Feature Map

# Non-Linearity

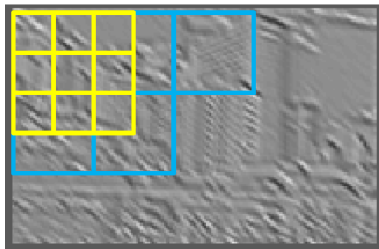
- Non-linearity
    - Per-element (independent)
    - Tanh
    - Sigmoid:  $1/(1+\exp(-x))$
    - Rectified linear
      - Simplifies backprop
      - Makes learning faster
      - Avoids saturation issues
- Preferred option



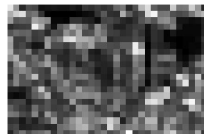


# Pooling

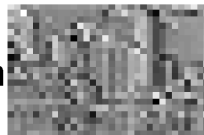
- Spatial Pooling
  - Non-overlapping / overlapping regions
  - Sum or max
  - Boureau et al. ICML'10 for theoretical analysis



Max



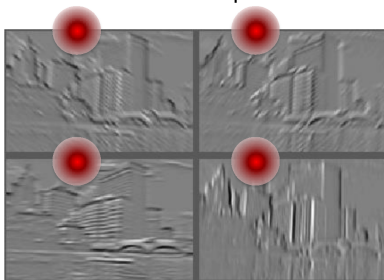
Sum



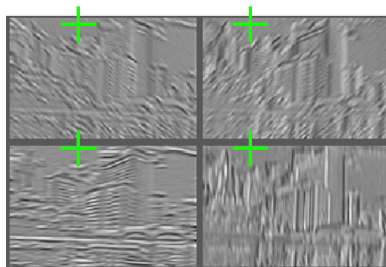
# Normalization

- Contrast normalization (across feature maps)
  - Local mean = 0, local std. = 1, “Local”  $\rightarrow$  7x7 Gaussian
  - Equalizes the features maps

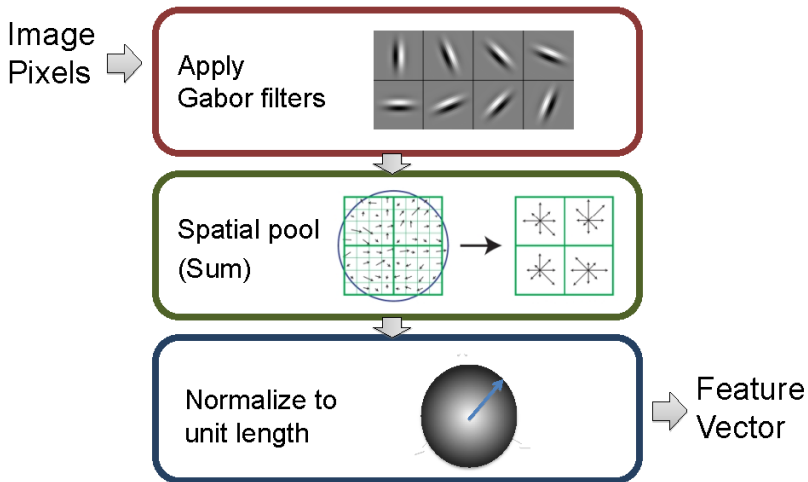
Feature Maps



Feature Maps  
After Contrast Normalization

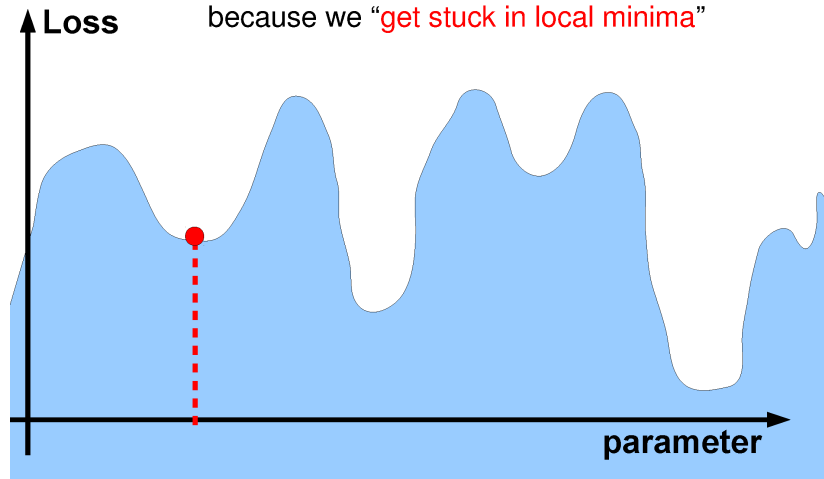


## Compare with SIFT Descriptor



# ConvNets: till 2012

Common wisdom: training does not work  
because we “get stuck in local minima”

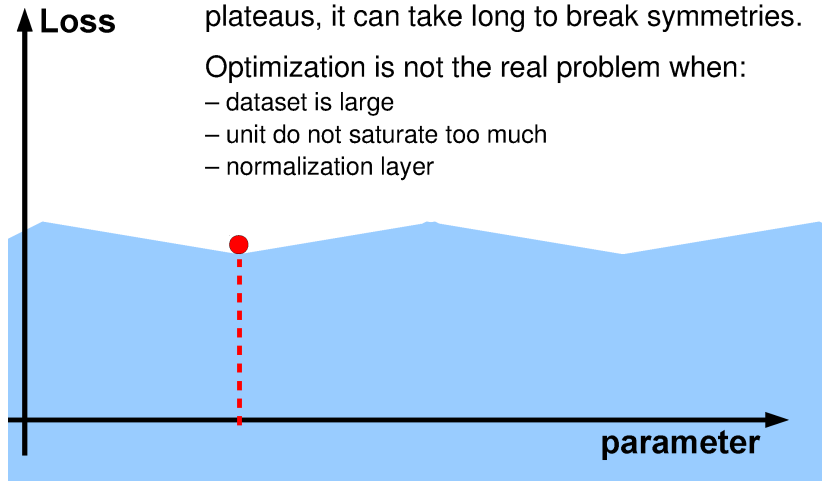


# ConvNets: today

Local minima are all similar, there are long plateaus, it can take long to break symmetries.

Optimization is not the real problem when:

- dataset is large
- unit do not saturate too much
- normalization layer

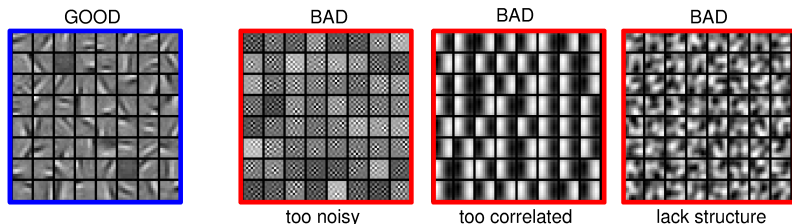


# Choosing the Architecture

- ▶ Task dependent
- ▶ Cross-validation
- ▶ [*Convolution*  $\rightarrow$  *Pooling*]\* + Fully Connected
- ▶ The more data, the more layers and kernels
  1. Look at the number of parameters at each layer
  2. Look at the number of flops at each layer
- ▶ Computational resources

## Good to Know

- ▶ Check gradients numerically by finite differences
- ▶ Visualize features  
feature maps need to be uncorrelated and have high variance
- ▶ Visualize parameters
- ▶ Measure error on both training and validation set
- ▶ Test on a small subset of the data and check the error  $\rightarrow 0$



**Good training:** learned filters exhibit structure and are uncorrelated.

# What if it does Not Work

- ▶ Training diverges:
  - ▶ Learning rate may be too large → decrease learning rate
  - ▶ Back Propagation is buggy → a numerical gradient checking
- ▶ Parameters collapse → Check loss function:
  - ▶ Is it appropriate for the task you want to solve
  - ▶ Does it have degenerate solutions? Check "pull-up" term
- ▶ Network is underperforming
  - ▶ Compute flops and nr. params. → if too small, make net larger
  - ▶ Visualize hidden units/params → fix optimization
- ▶ Network is too slow
  - ▶ Compute flops and nr. params.
  - ▶ GPU, distributed framework, make net smaller



# Summary

- ▶ Deep Learning = Learning Hierarchical Representations
- ▶ Supervised Learning: most successful set up today
- ▶ Optimization
  - ▶ Don't we get stuck in local minima? No, they are all the same
  - ▶ In large scale applications, local minima are even less of an issue
- ▶ Scaling
  - ▶ GPUs
  - ▶ Distributed framework
  - ▶ Better optimization techniques
- ▶ Generalization on small datasets (curse of dimensionality)
  - ▶ Data augmentation
  - ▶ Weight decay
  - ▶ Dropout
  - ▶ Unsupervised Learning
  - ▶ Multi-task Learning
  - ▶ Transfer Learning

Thank You

Q&A