

First of all what is Deep Learning?

- Composition of non-linear transformation of the data.
- Goal: **Learn useful representations, aka features, directly from data.**
- Many varieties, can be unsupervised or supervised.
- Today is about ConvNets, which is a **supervised** deep learning method.

Recap: Supervised Learning

$\{(\mathbf{x}^i, y^i), i=1 \dots P\}$ training dataset

\mathbf{x}^i i-th input training example

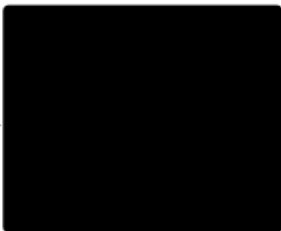
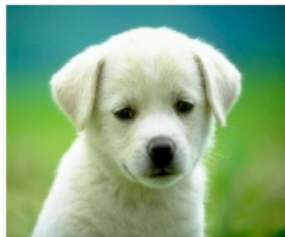
y^i i-th target label

P number of training examples



Supervised Learning: Examples

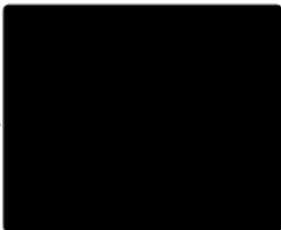
Classification



“dog”

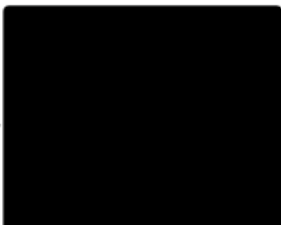
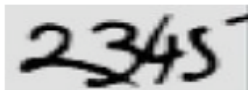
classification

Denoising



regression

OCR



“2 3 4 5”

structured prediction

Supervised Deep Learning

Classification

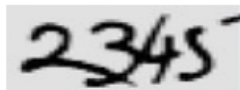


“dog”

Denoising



OCR

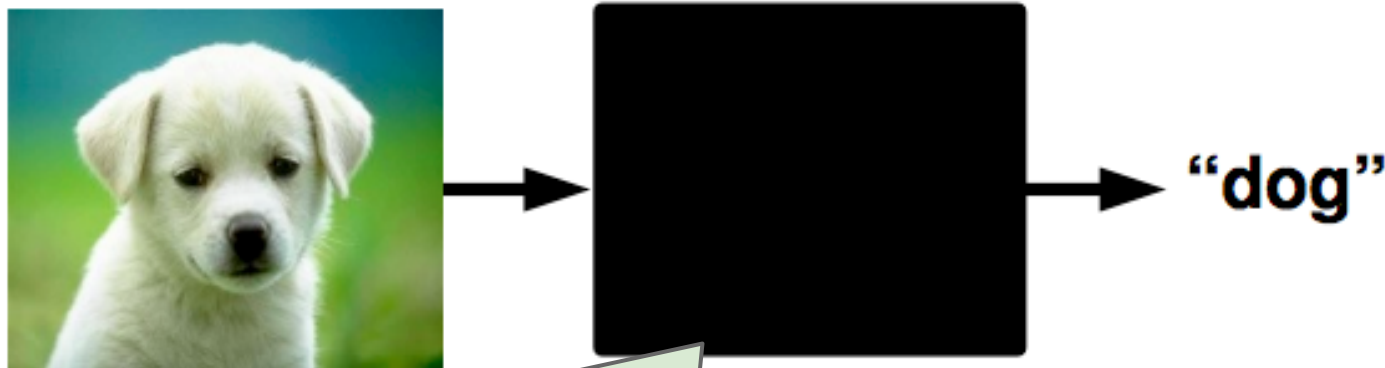


“2 3 4 5”

So deep learning is about learning
feature representation in a
compositional manner.

But wait,
why learn features?

The Black Box in a Traditional Recognition Approach



Preprocessing

Feature
Extraction
(HOG, SIFT, etc)

Post-processing
(Feature selection,
MKL etc)

Classifier
(SVM,
boosting, etc)

The Black Box in a Traditional Recognition Approach



“dog”

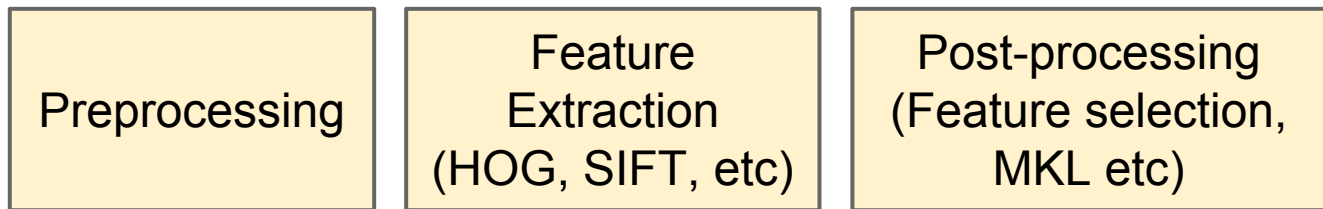
Hand Engineered

Preprocessing

Extract
(HOG, SIFT, etc)

Post-processing
(feature selection,
MKL etc)

Classifier
(SVM,
boosting, etc)



- Most critical for accuracy
 - Most time-consuming in development
 - What is the best feature???
 - What is next?? Keep on crafting better features?
- ⇒ Let's learn feature representation directly from data.

Learn features and classifier *together*

⇒ Learn an end-to-end recognition system.

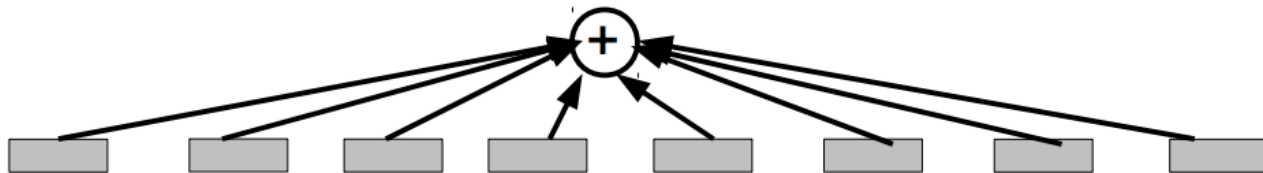
A non-linear map that takes raw pixels directly to labels.

Q: How can we build such a highly non-linear system?

A: By combining simple building blocks we can make more and more complex systems.

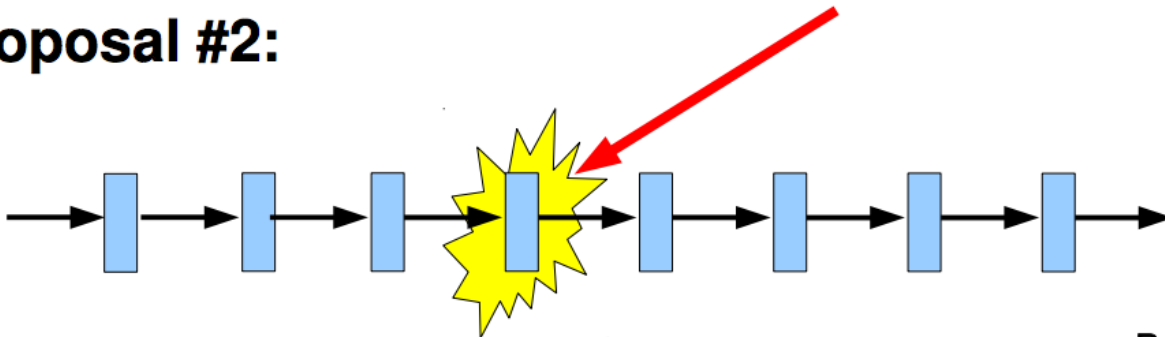
Building a complicated function

Proposal #1:



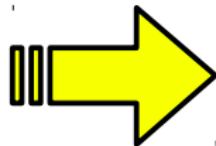
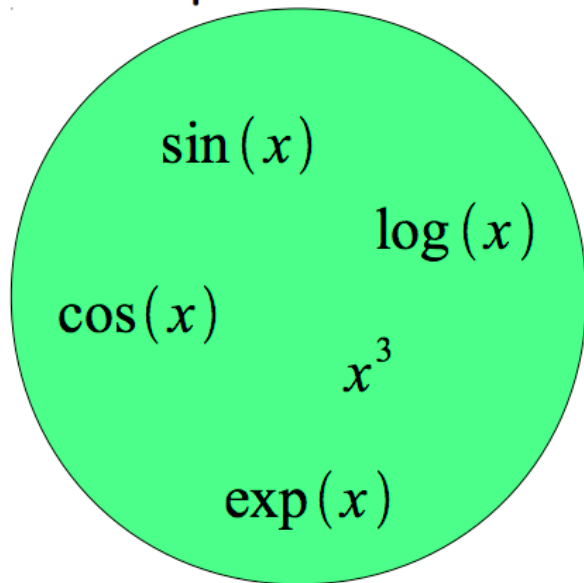
Each box is a simple nonlinear function

Proposal #2:

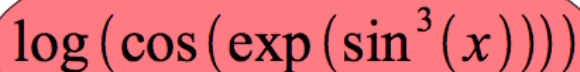


Building a complicated function

Simple Functions

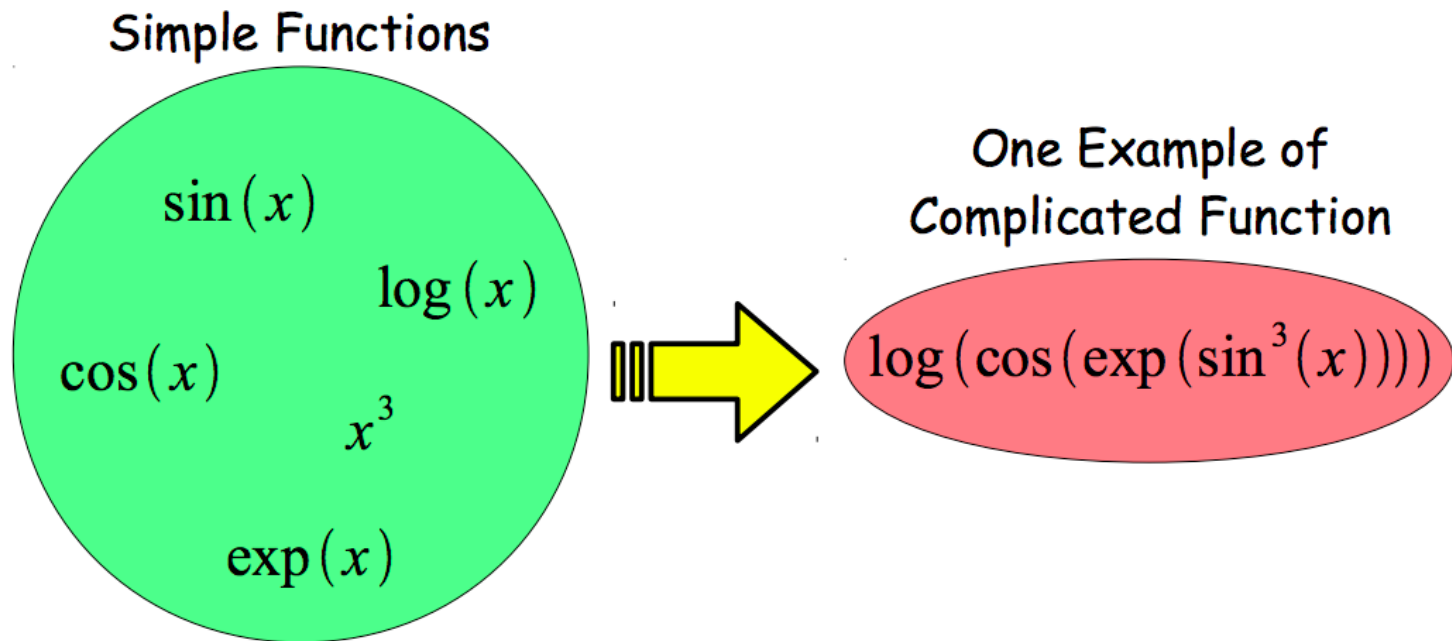


One Example of
Complicated Function



A red oval containing the complicated function expression: $\log(\cos(\exp(\sin^3(x))))$.

Building a complicated function

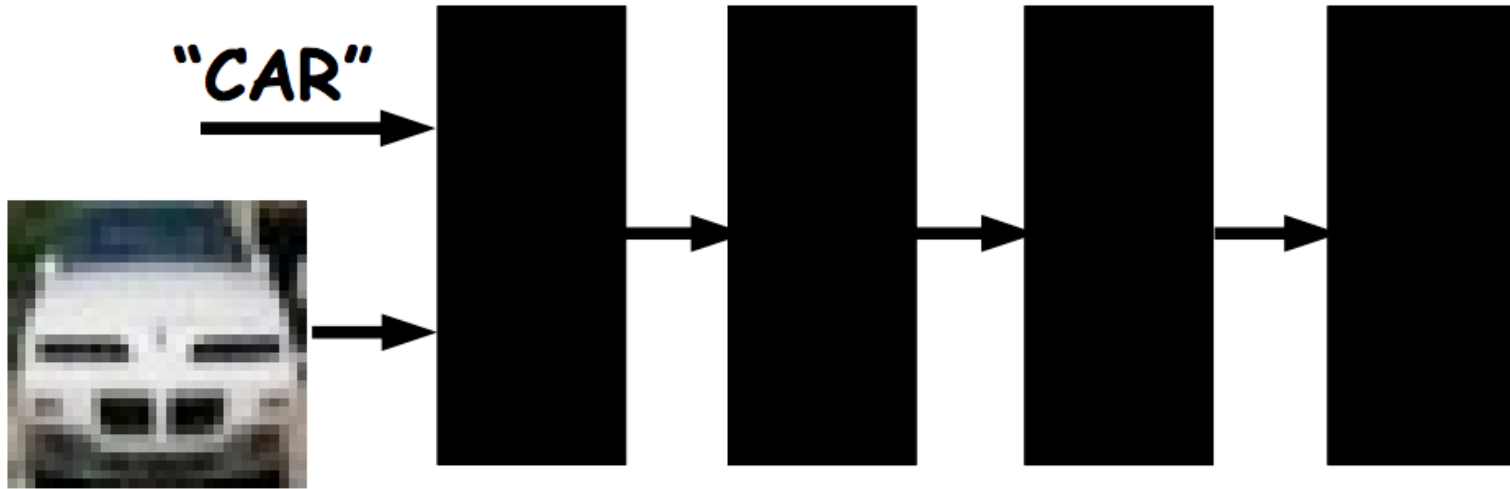


- Composition is at the core of deep learning methods
- Each “simple function” will have parameters subject to learning

Intuition behind Deep Neural Nets

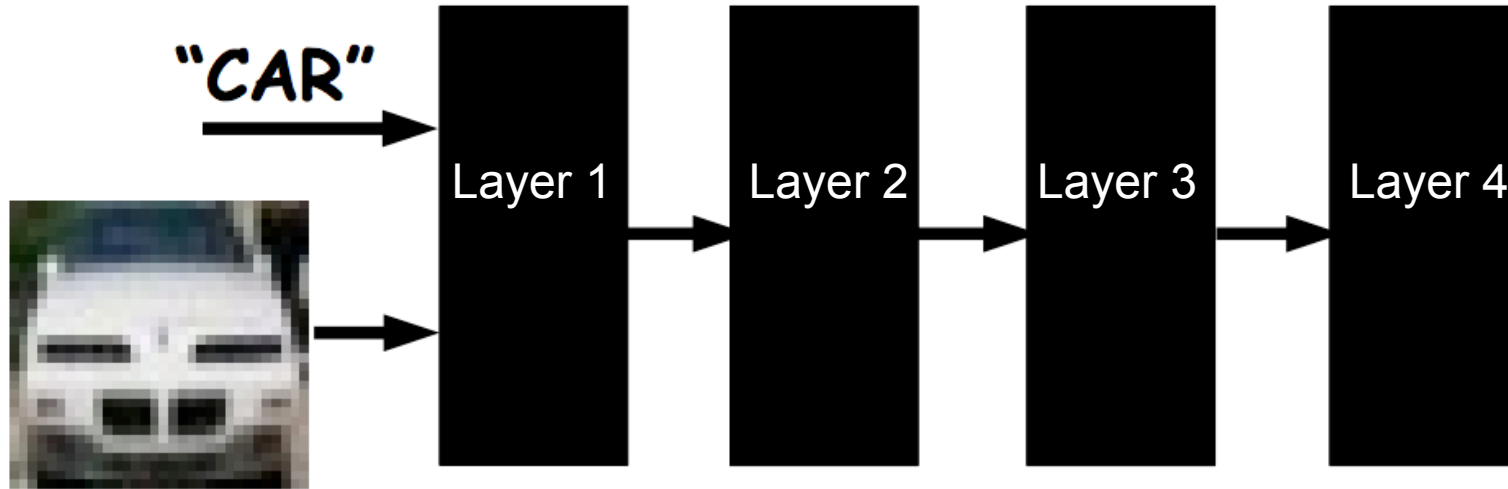


Intuition behind Deep Neural Nets



NOTE: Each black box can have trainable parameters.
Their composition makes a highly non-linear system.

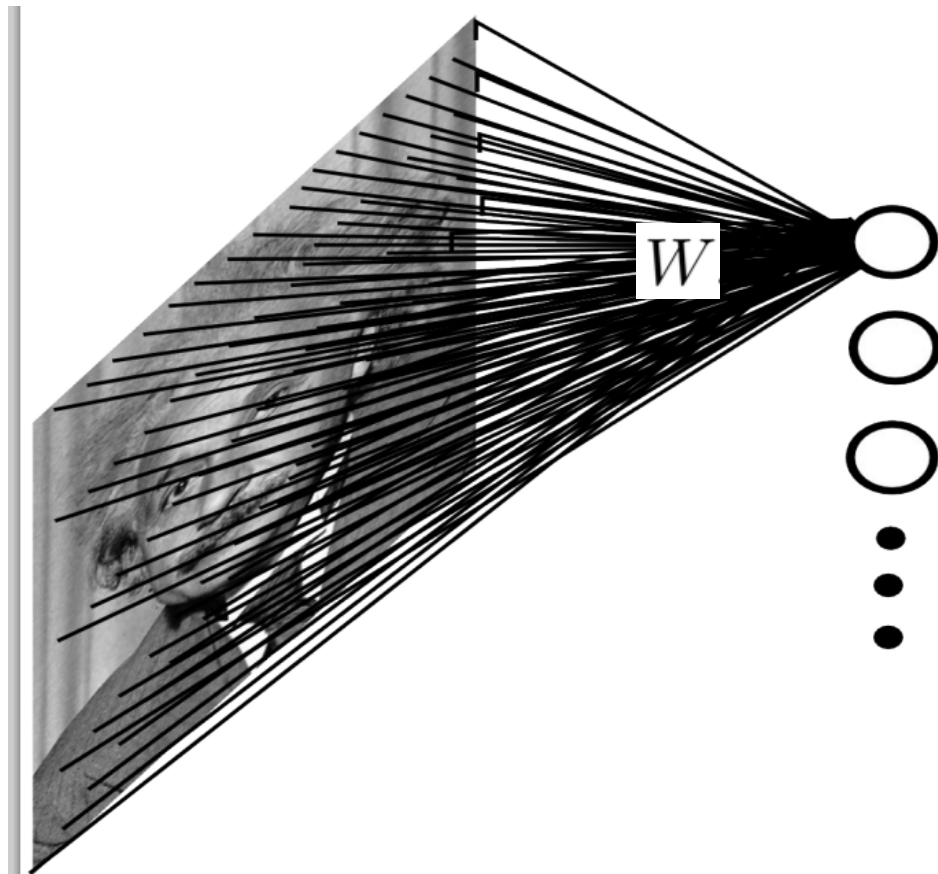
Intuition behind Deep Neural Nets



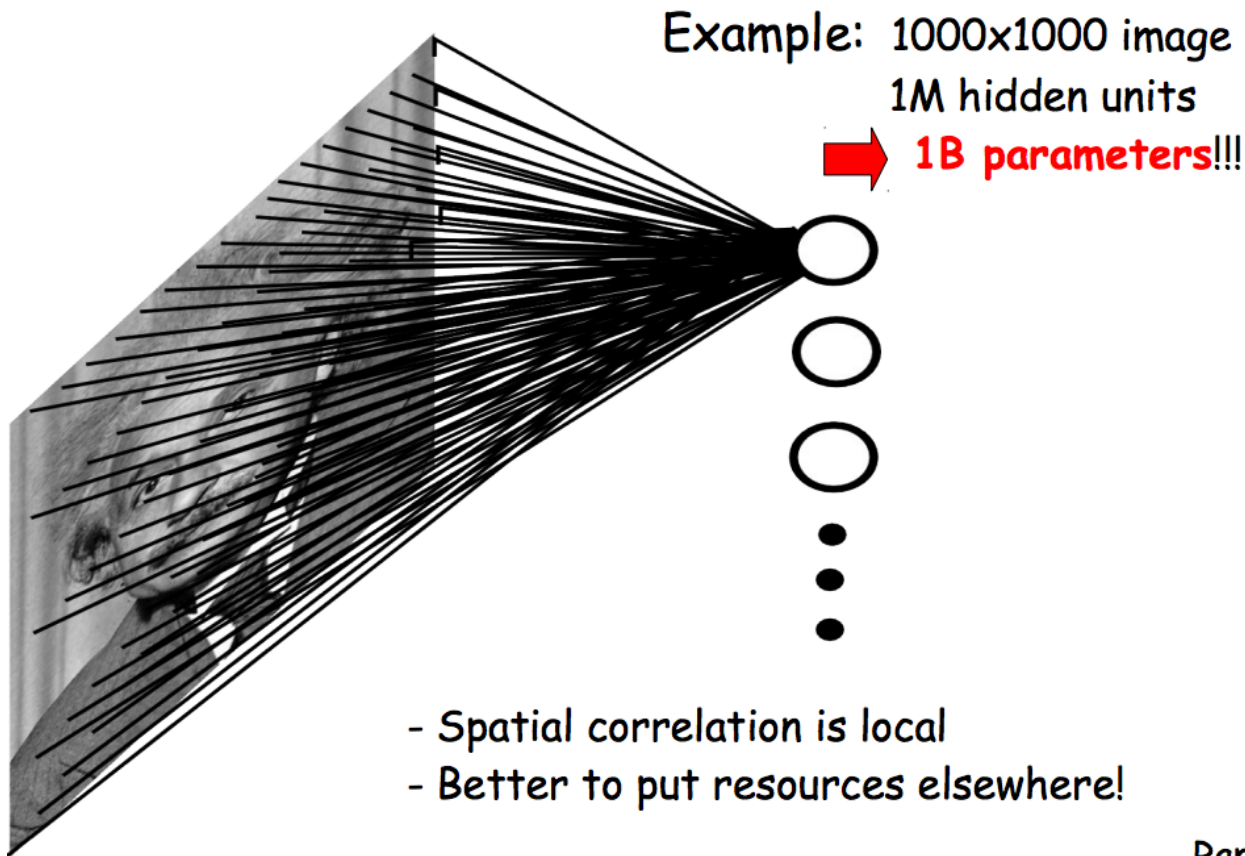
NOTE: Each black box can have trainable parameters.
Their composition makes a highly non-linear system.

The final layer outputs a probability distribution of categories.

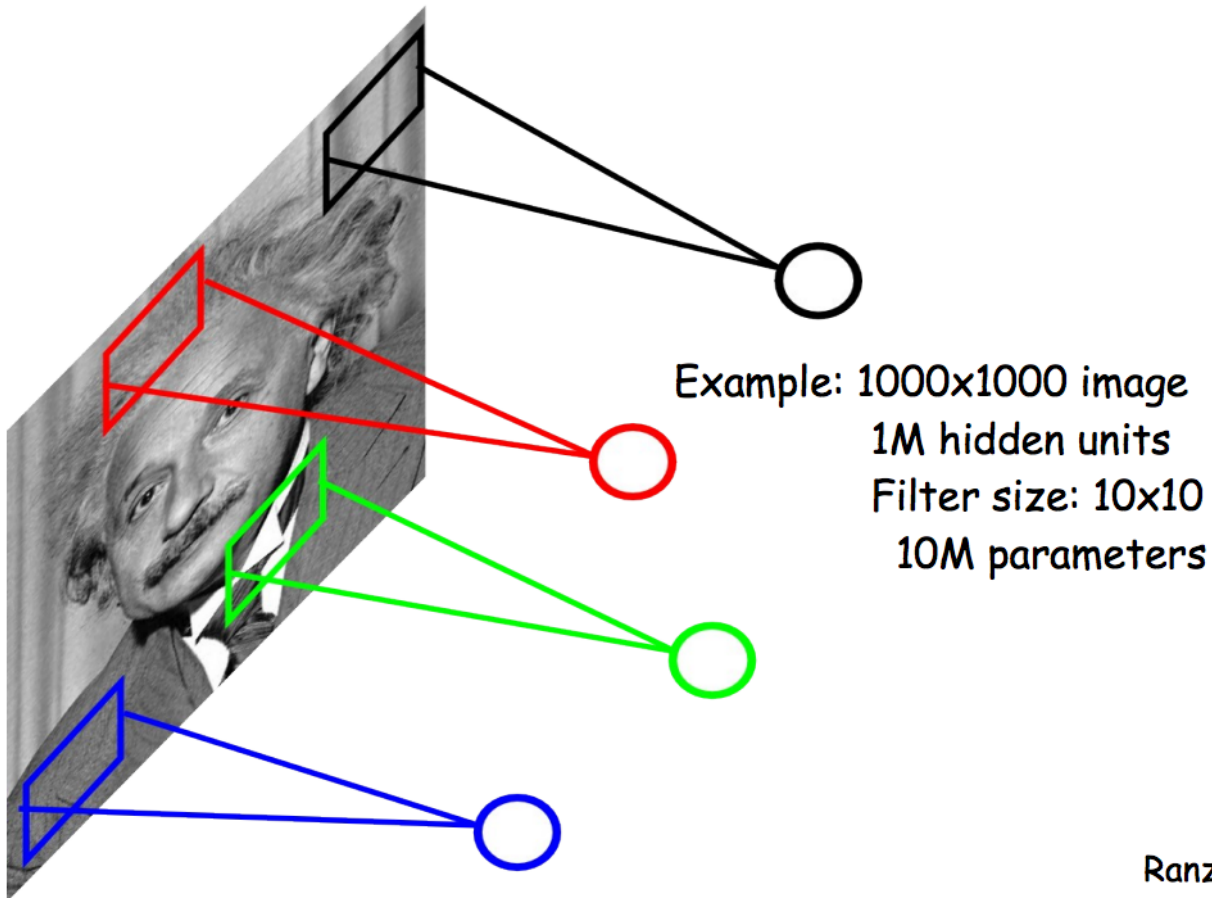
When the input data is an image..



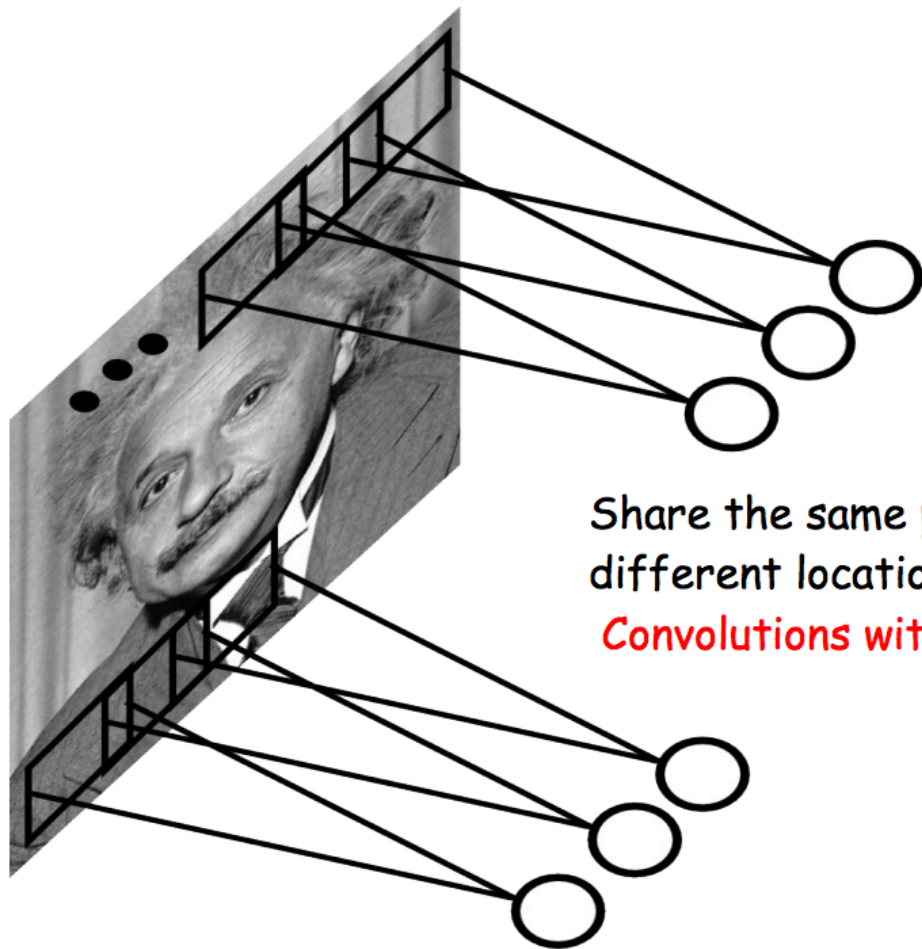
When the input data is an image..



Reduce connection to local regions



Reuse the same kernel everywhere

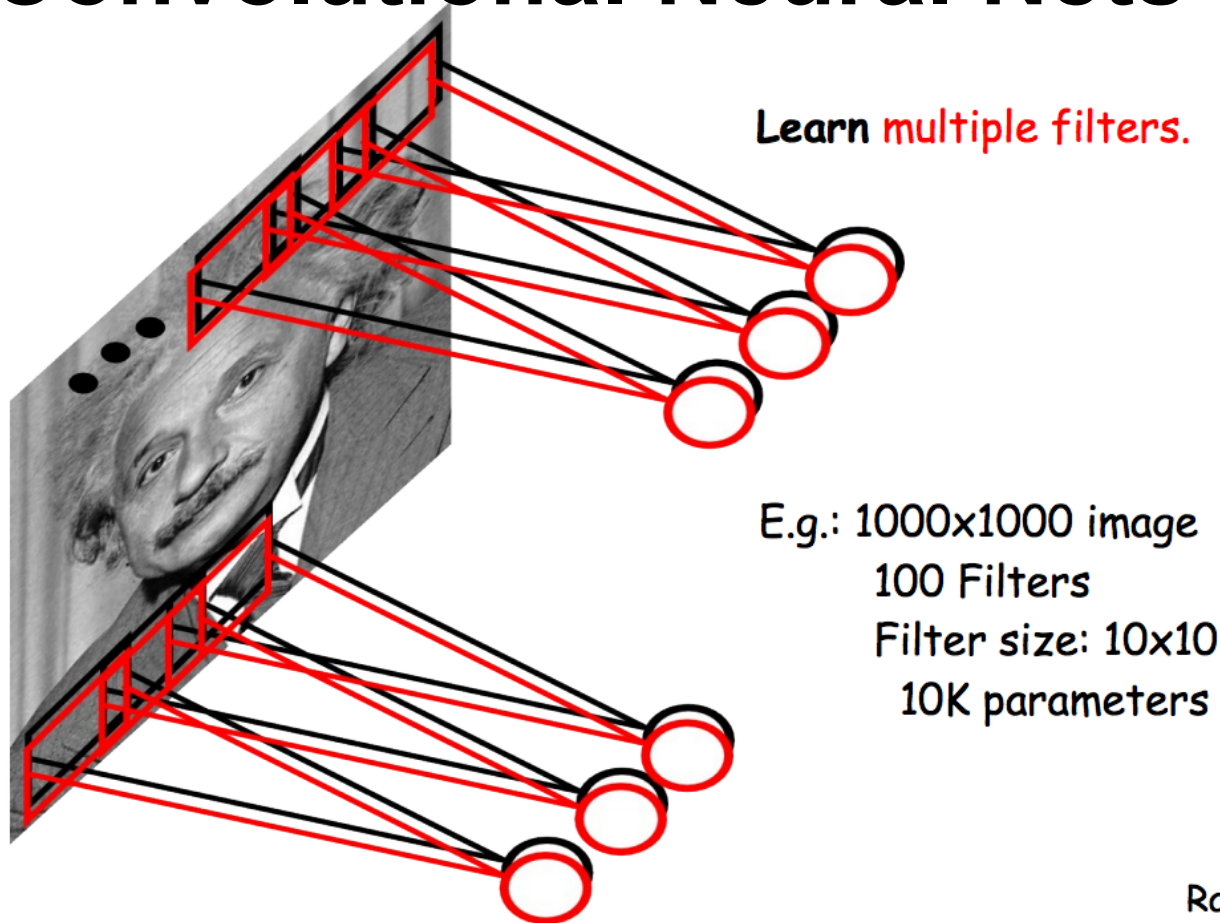


Because interesting features (edges) can happen at anywhere in the image.

Share the same parameters across different locations:

Convolutions with learned kernels

Convolutional Neural Nets

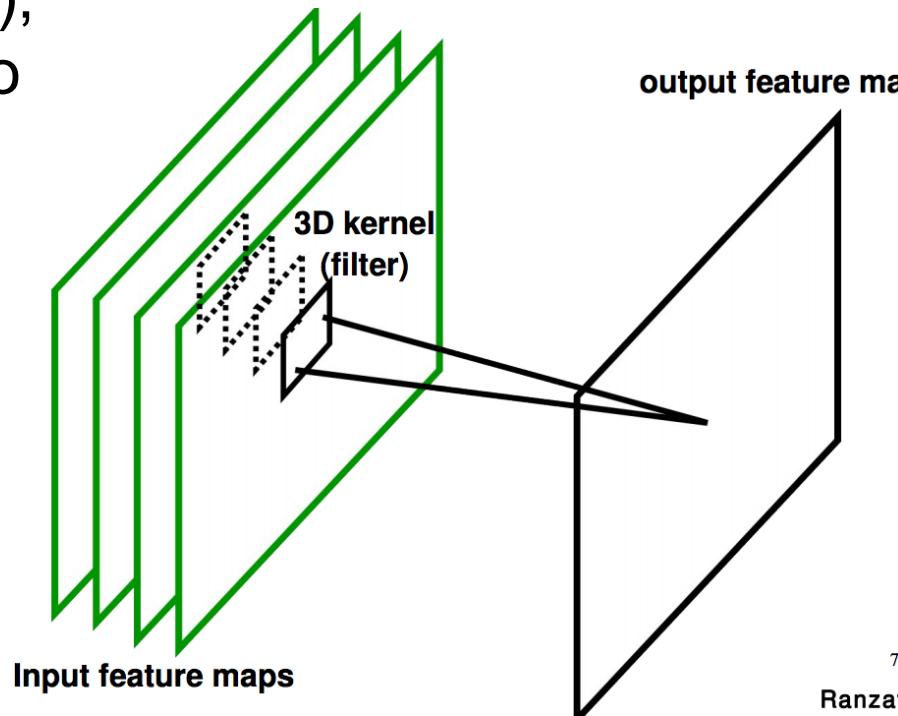


Detail

If the input has 3 channels (R,G,B), 3 separate k by k filter is applied to each channel.

Output of convolving 1 feature is called a *feature map*.

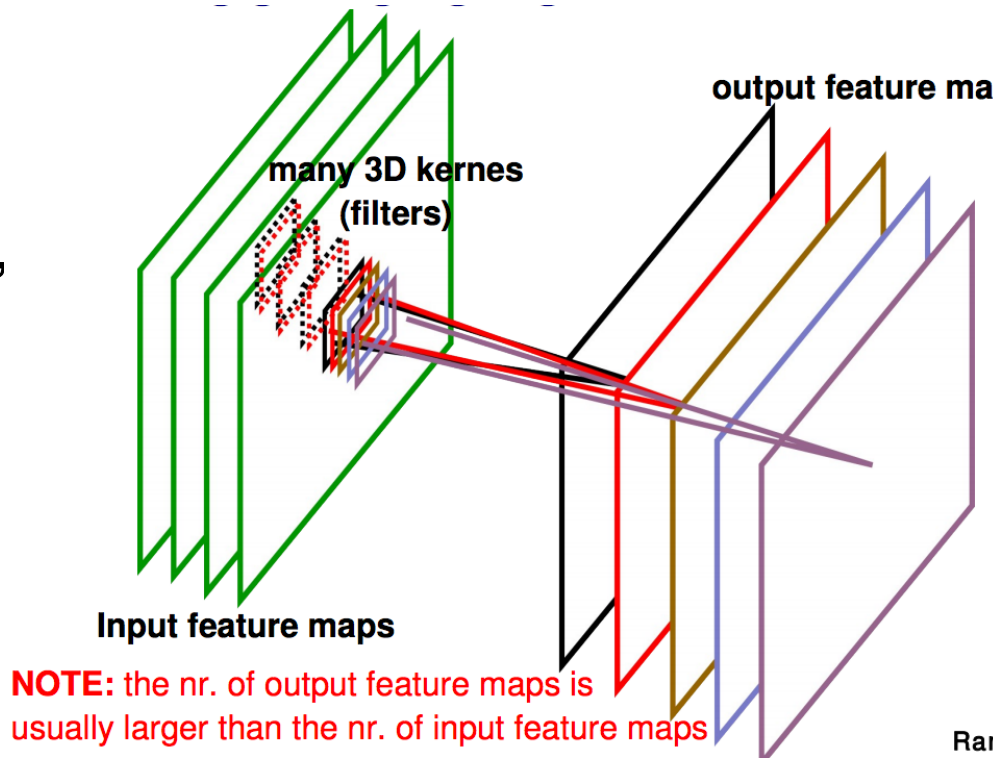
This is just sliding window, ex. the output of one part filter of DPM is a feature map



Using multiple filters

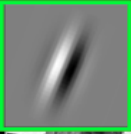
Each filter detects features in the output of previous layer.

So to capture different features, learn multiple filters.



Example of filtering

- Convolutional
 - Translation equivariance
 - Tied filter weights
(same at each position → few parameters)

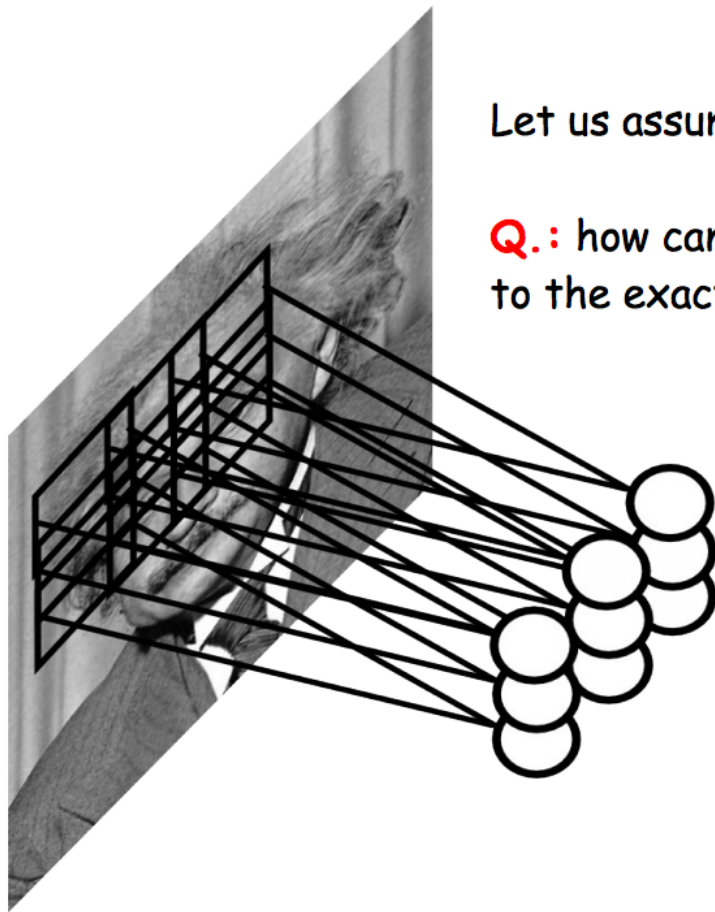


Input



Feature Map

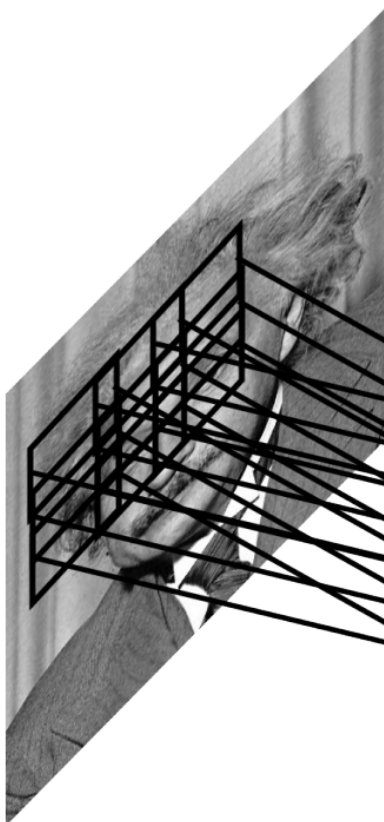
Building Translation Invariance



Let us assume filter is an “eye” detector.

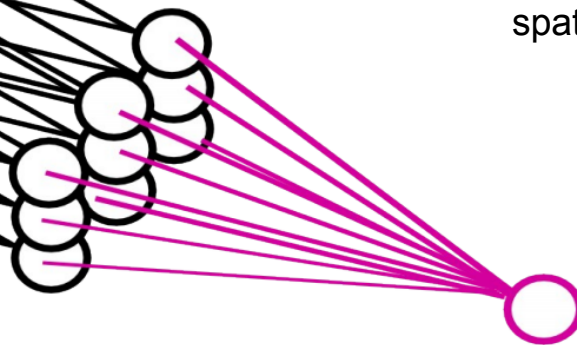
Q.: how can we make the detection robust to the exact location of the eye?

Building Translation Invariance via Spatial Pooling



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

Pooling also subsamples the image, allowing the next layer to look at larger spatial regions.



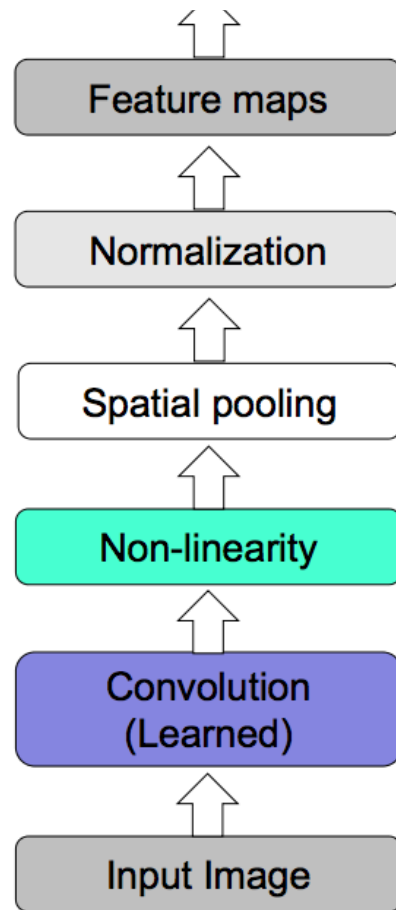
Summary of a typical convolutional layer

Doing all of this consists one layer.

- Pooling and normalization is optional.

Stack them up and train just like multi-layer neural nets.

Final layer is usually fully connected neural net with output size == number of classes



Revisiting the composition idea

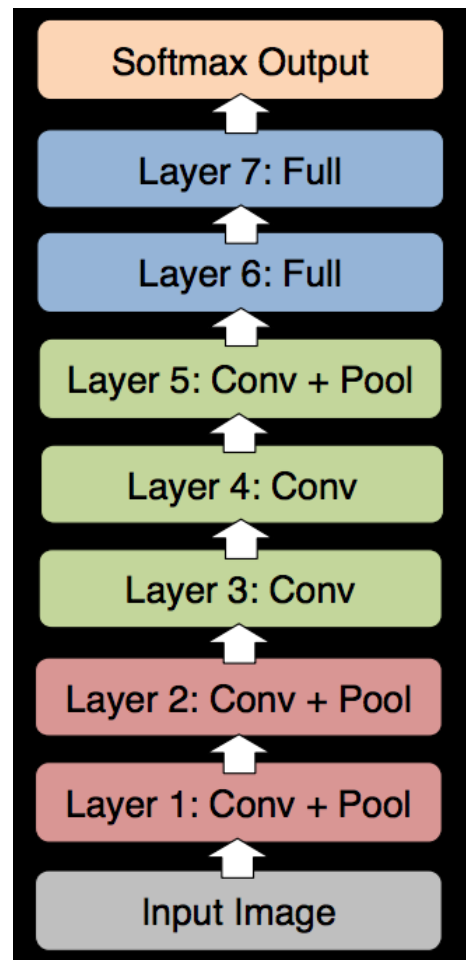
Every layer learns a feature detector by combining the output of the layer before.

⇒ More and more abstract features are learned as we stack layers.

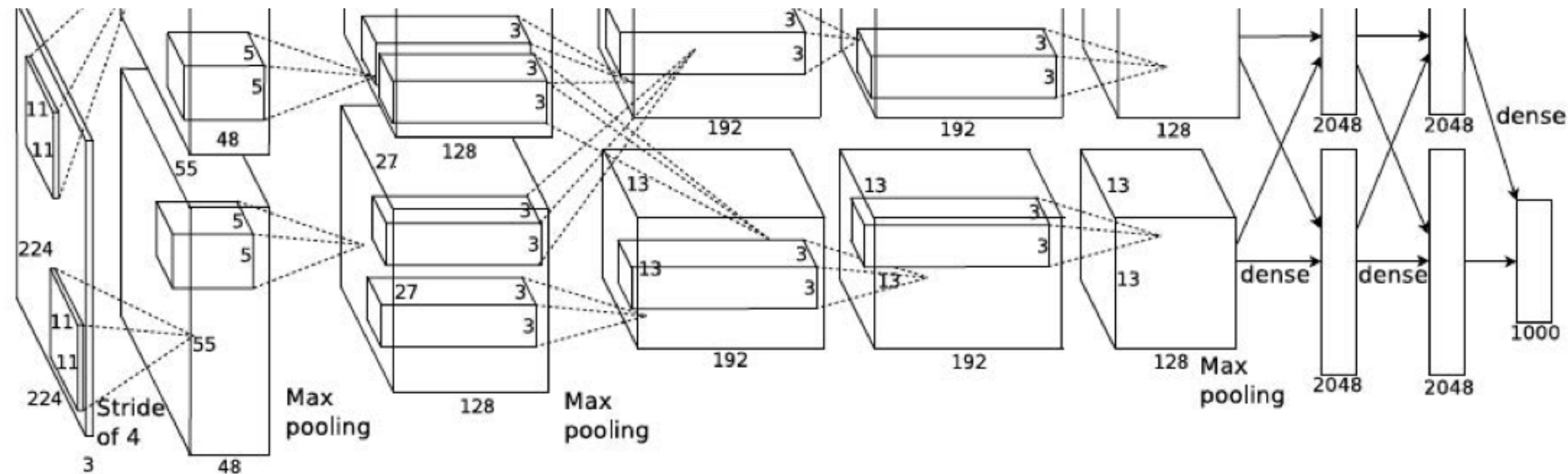
Keep this in mind and let's look at what kind of things ConvNets learn.

Architecture of Alex Krizhevsky et al.

- 8 layers total.
- Trained on Imagenet Dataset (1000 categories, 1.2M training images, 150k test images)
- 18.2% top-5 error
 - Winner of the ILSVRC-2012 challenge.



Architecture of Alex Krizhevsky et al.

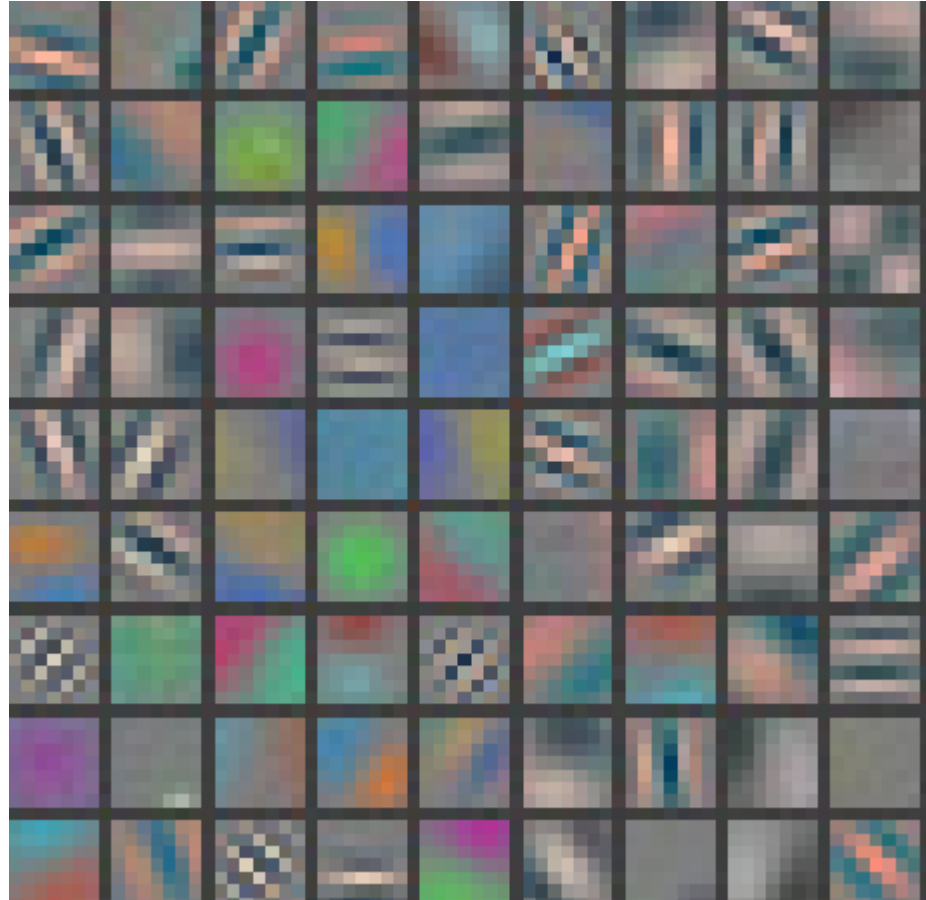


First layer filters

Showing 81 filters of
 $11 \times 11 \times 3$.

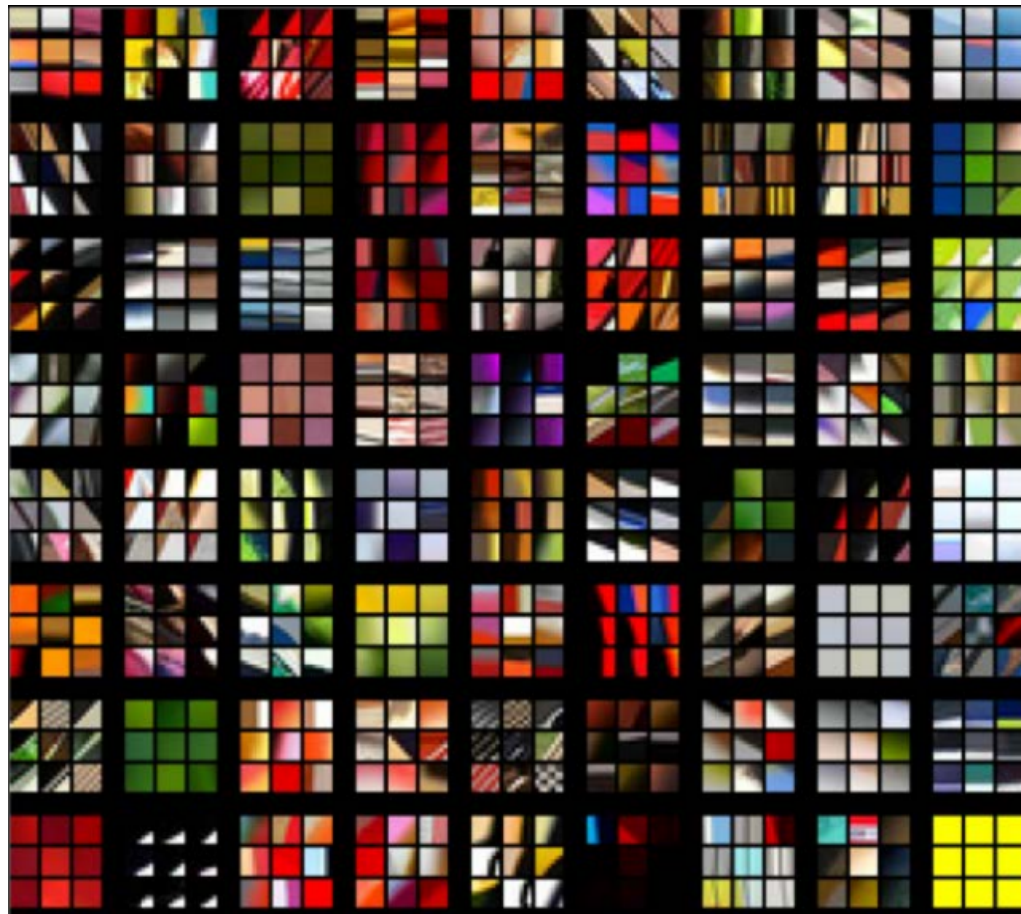
Capture low-level
features like oriented
edges, blobs.

Note these oriented edges are
analogous to what SIFT uses to
compute the gradients.



Top 9 patches that activate each filter in layer 1

Each 3x3 block shows the top 9 patches for one filter.



Layer 2: Top-9 Patches



Layer 2: Top-9 Patches

Note how the previous low-level features are combined to detect a little more abstract features like textures.





Layer 4: Top-9 Patches



Layer 5: Top-9 Patches

ConvNets as generic feature extractor

- A well-trained ConvNets is an excellent feature extractor.
- Chop the network at desired layer and use the output as a feature representation to train a SVM on some other vision dataset.

	Cal-101 (30/class)	Cal-256 (60/class)
SVM (1)	44.8 \pm 0.7	24.6 \pm 0.4
SVM (2)	66.2 \pm 0.5	39.6 \pm 0.3
SVM (3)	72.3 \pm 0.4	46.0 \pm 0.3
SVM (4)	76.6 \pm 0.4	51.3 \pm 0.1
SVM (5)	86.2 \pm 0.8	65.6 \pm 0.3
SVM (7)	85.5 \pm 0.4	71.7 \pm 0.2
Softmax (5)	82.9 \pm 0.4	65.7 \pm 0.5
Softmax (7)	85.4 \pm 0.4	72.6 \pm 0.1

- Improve further by taking a pre-trained ConvNet and re-training it on a different dataset. Called *fine-tuning*