# Deep Learning for
# Object Category Recognition
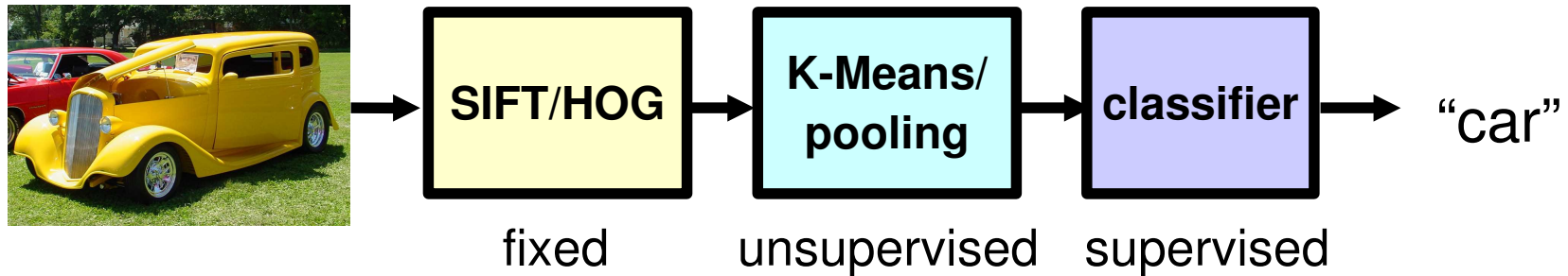
Marc'Aurelio Ranzato
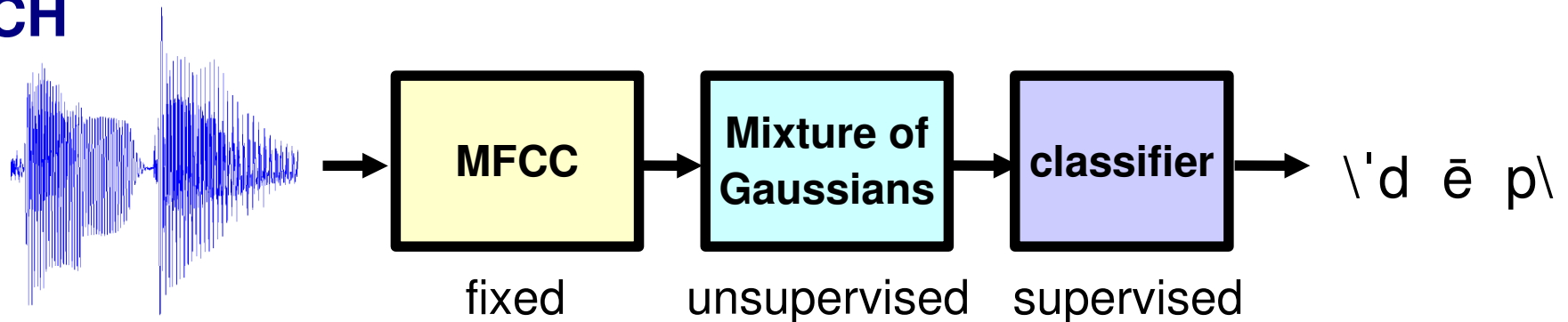
 Facebook, AI Group

# Traditional Pattern Recognition
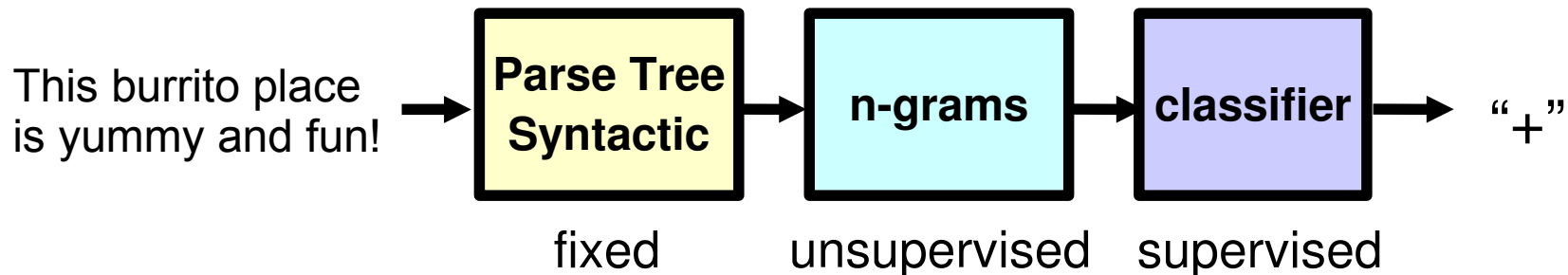
**VISION**



| SIFT/HOG | → | K-Means/ pooling | → | classifier | → "car" |
| fixed | | unsupervised | | supervised | |

**SPEECH**



| MFCC | → | Mixture of Gaussians | → | classifier | → \'d ē p\ |
| fixed | | unsupervised | | supervised | |

**NLP**

This burrito place is yummy and fun!

| Parse Tree Syntactic | → | n-grams | → | classifier | → "+" |
| fixed | | unsupervised | | supervised | |

Ranzato

# Hierarchical Compositionality (DEEP)

**VISION**

pixels ➡ edge ➡ texton ➡ motif ➡ part ➡ object

**SPEECH**

sample ➡ spectral band ➡ formant ➡ motif ➡ phone ➡ word

**NLP**

character ➡ word ➡ NP/VP/.. ➡ clause ➡ sentence ➡ story

3

**Ranzato**

# Deep Learning
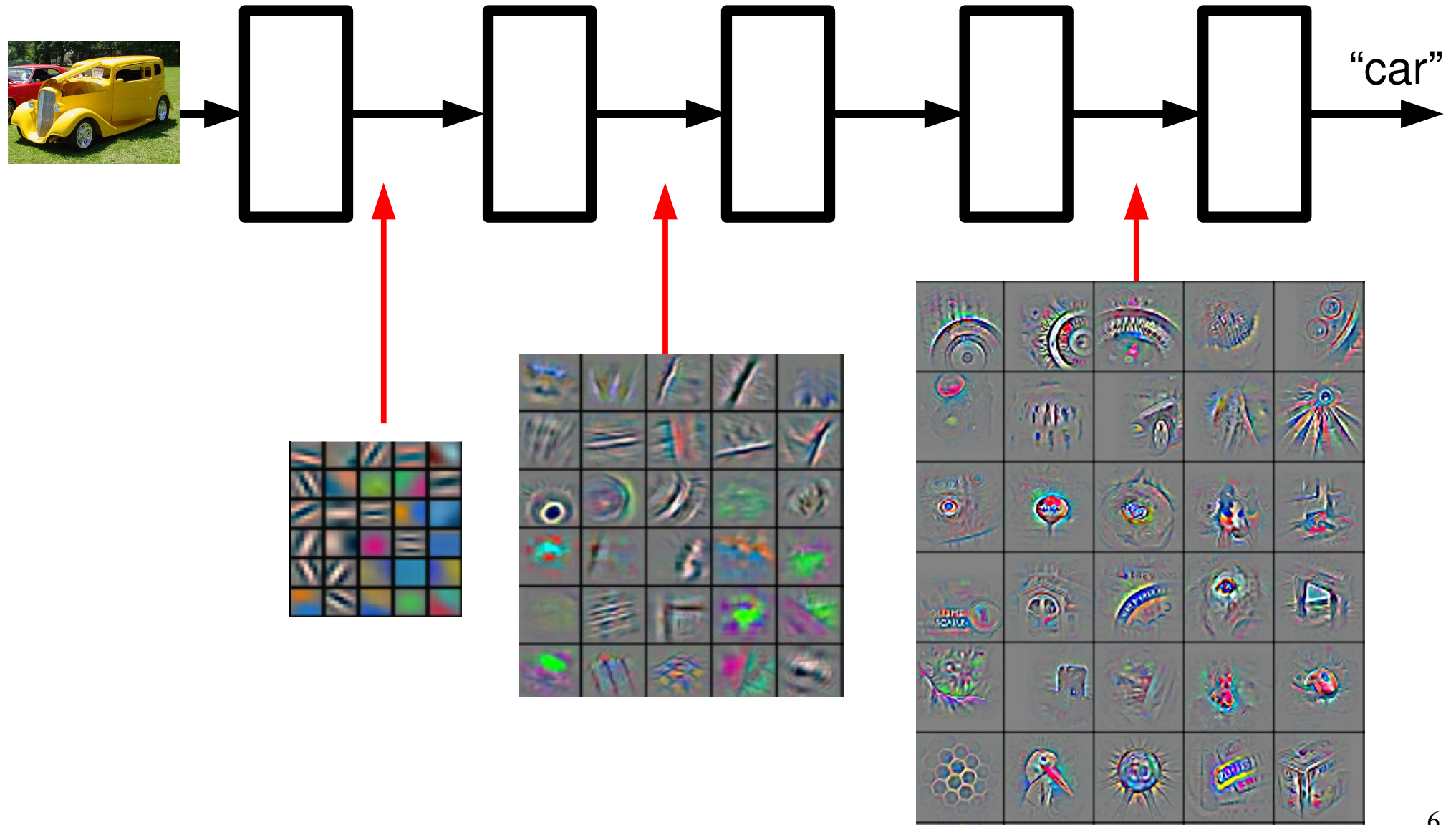


"car"

## What is Deep Learning

- Cascade of non-linear transformations
- End to end learning
- General framework (any hierarchical model is deep)

Ranzato

# Deep Learning  VS  Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels ➤ edge ➤ texton ➤ motif ➤ part ➤ object

**Ranzato**

# Deep Learning



"car"

Zeiler et al. "Visualizing and Understanding ConvNets" Arxiv 2013

Ranzato

6

# Deep Learning  VS  Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels → edge → texton → motif → part → object

- It is more efficient.

E.g.: Checking N-bit parity requires N-1 gates laid out on a tree of depth log(N-1). The same would require O(exp(N)) with a two layer architecture.

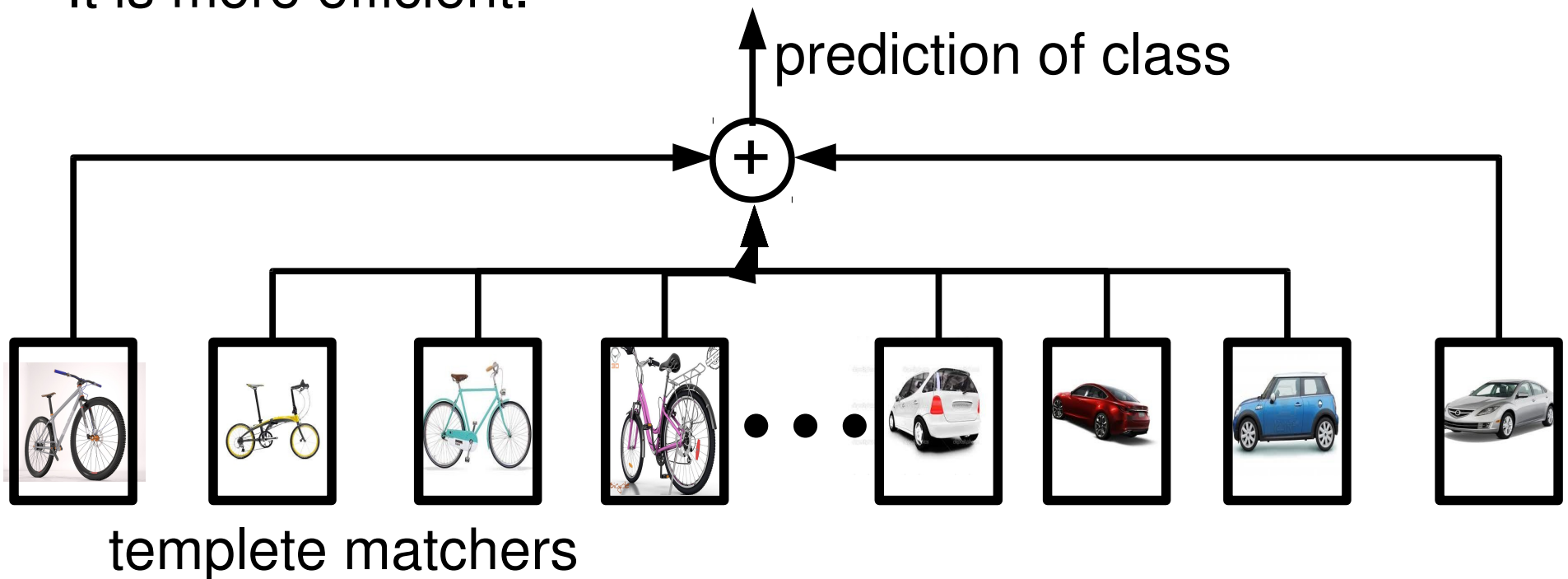$$p = \sum_i \alpha_i f_i(x) \qquad \text{VS} \qquad p = \alpha_n f_n(\alpha_{n-1} f_{n-1}(... \alpha_1 f_1(x)...))$$

Shallow learner is often inefficient: it requires exponential number of templates (basis functions).

Ranzato

# Deep Learning  VS  Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels ➡ edge ➡ texton ➡ motif ➡ part ➡ object

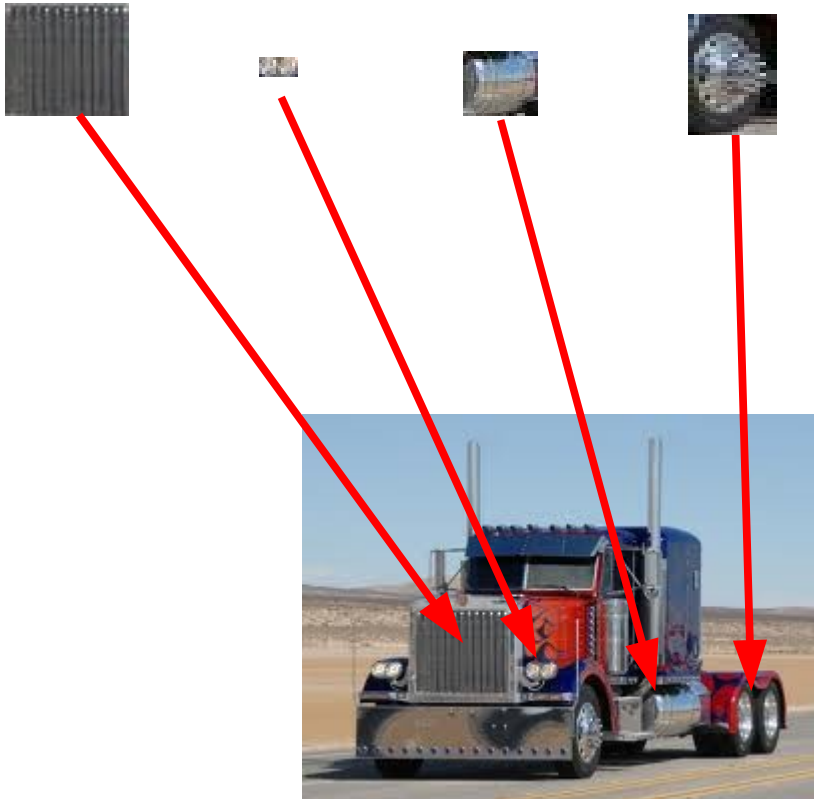- It is more efficient.



prediction of class

templete matchers

Shallow learner is inefficient.

Ranzato

# Composition: distributed representations

[0  0  **1**  0  0  0  0  **1**  0  0  **1**  **1**  0  0  **1**  0 … ]  truck feature
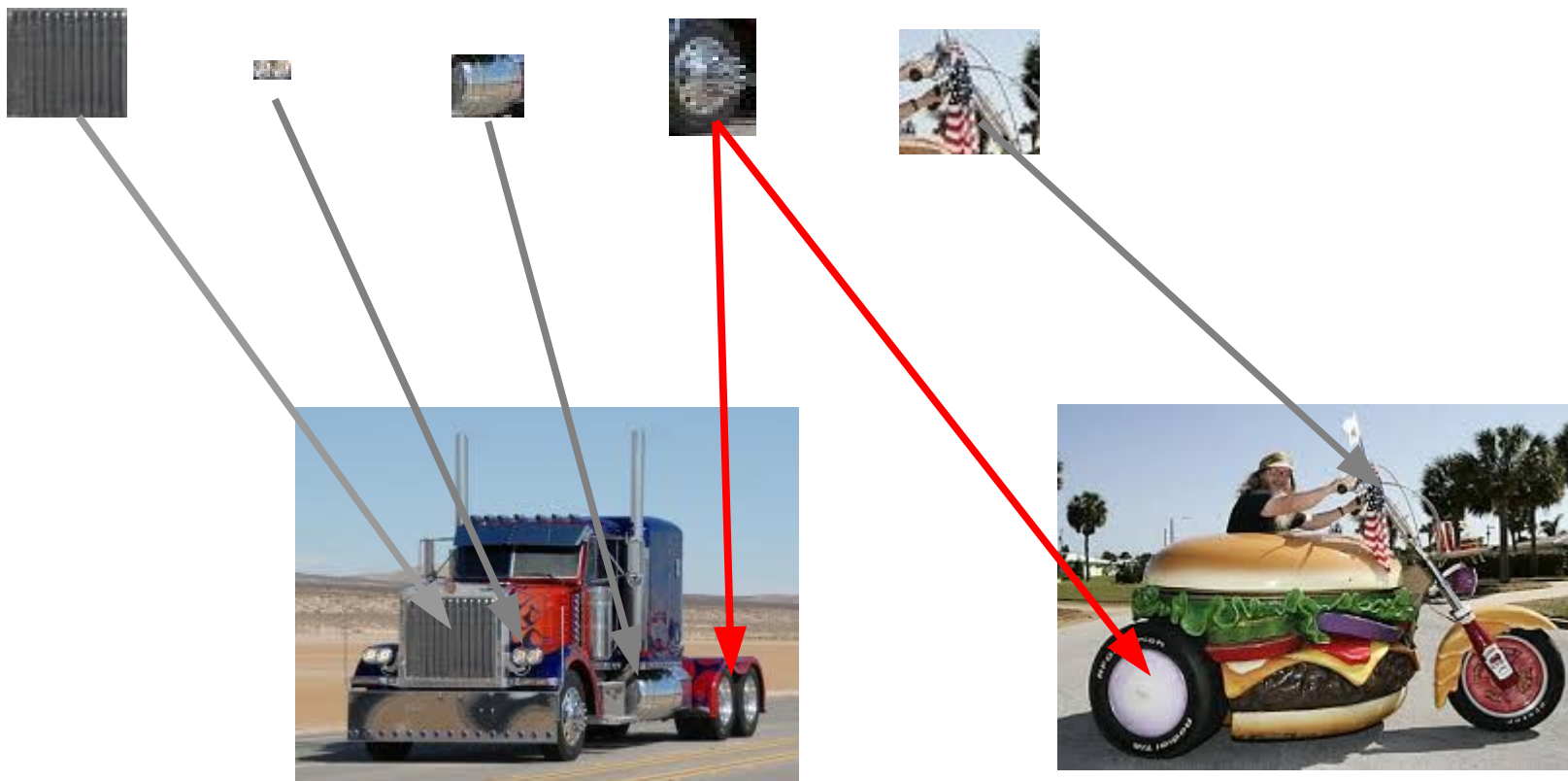


Exponentially more efficient than a 1-of-N representation (a la k-means)

**Ranzato**

# Composition: sharing
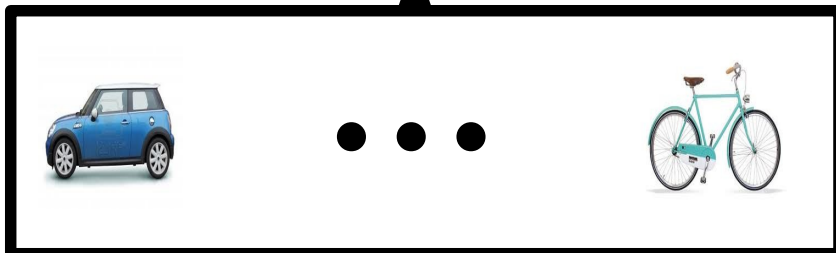
[1  1  0  0  0  1  0  **1**  0  0  0  0  0  **1**  **1**  0  **1**... ]  motorbike

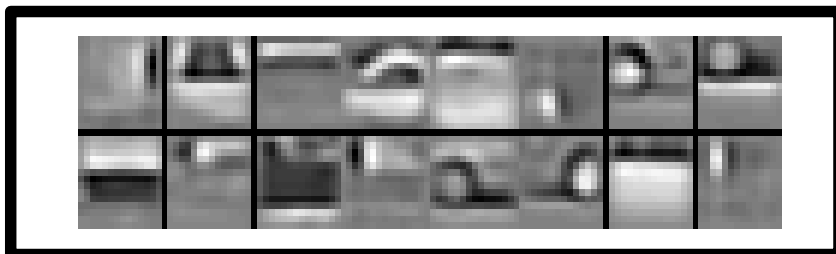[0  0  **1**  0  0  0  0  **1**  0  0  **1**  **1**  0  0  **1**  0 ... ]  truck
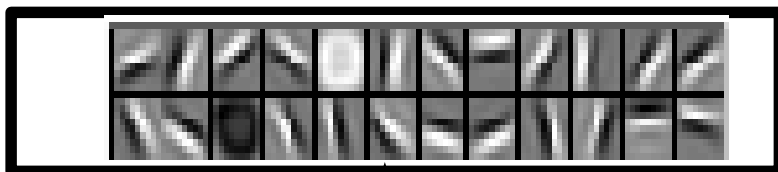
# Composition

prediction of class



high-level parts

mid-level parts

low level parts

Input image

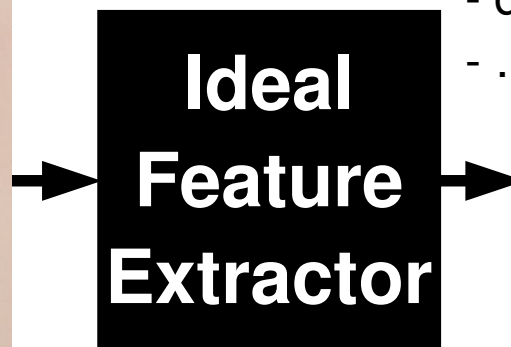- distributed representations
- feature sharing
- compositionality

**GOOD: (exponentially) more efficient**

Lee et al. "Convolutional DBN's ..." ICML 2009

Ranzato

# Deep Learning

# =

# Representation Learning

Ranzato

# Ideal Features



- window, right
- chair, left
- monitor, top of shelf
- carpet, bottom
- drums, corner
- ...

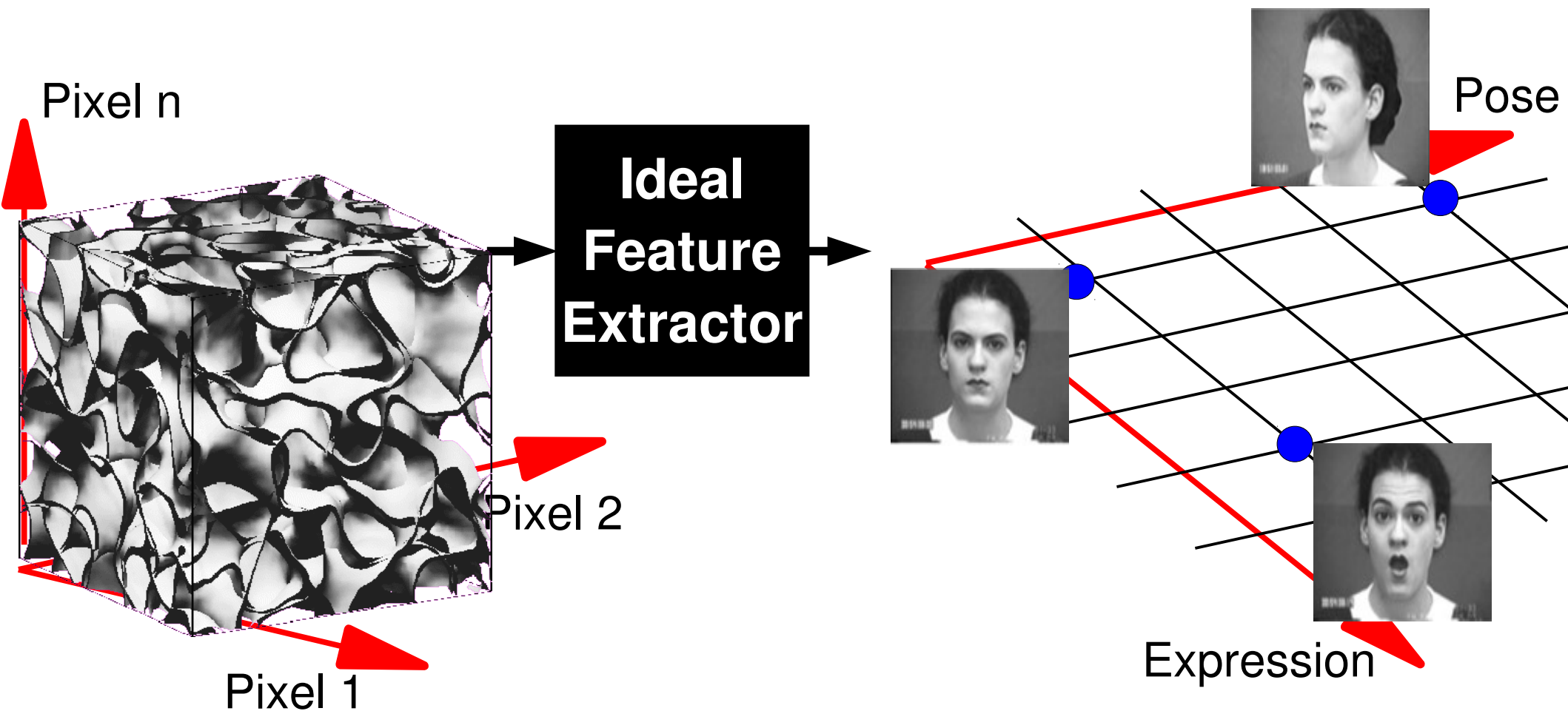**Ideal Feature Extractor**

- pillows on couch

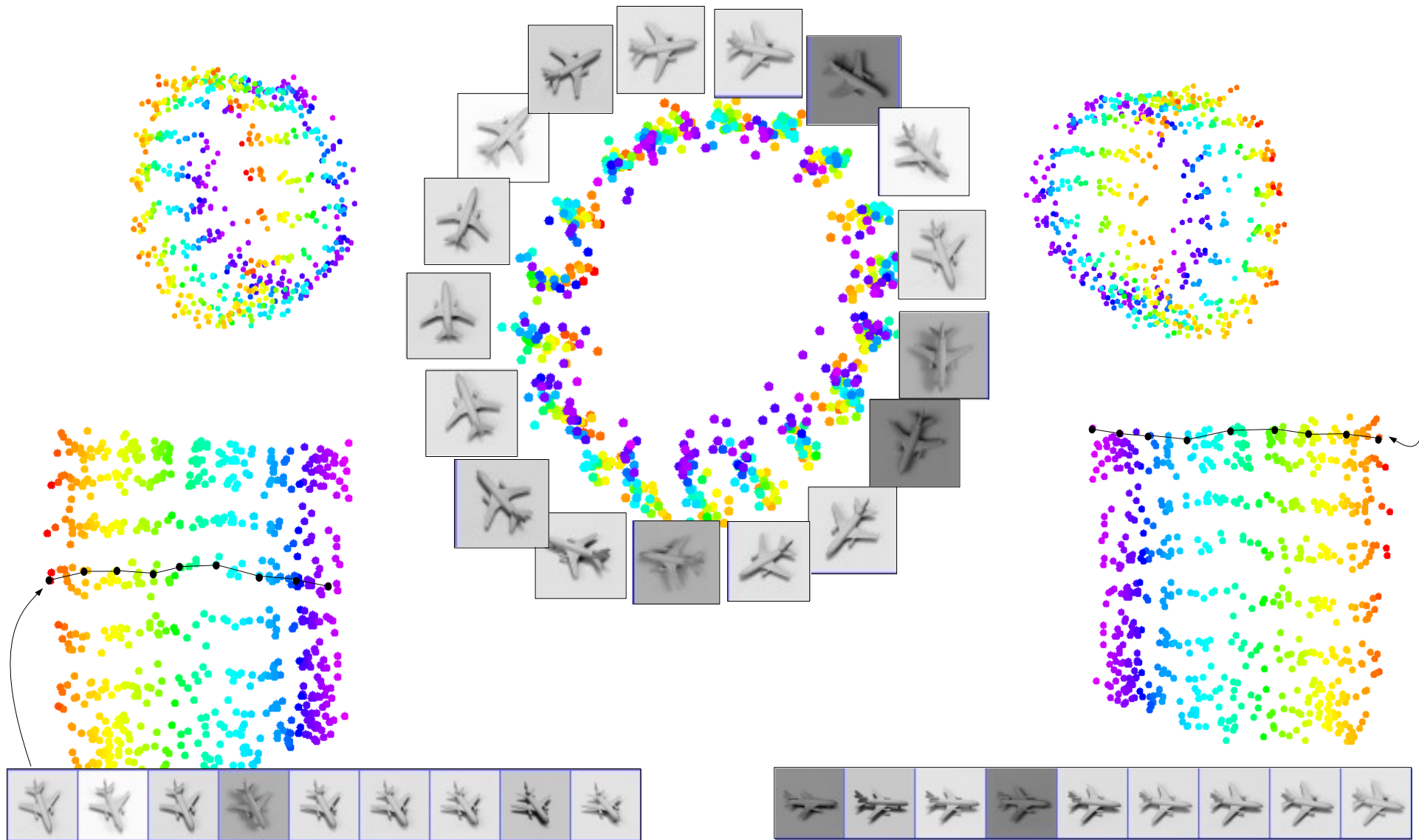**Q.:** What objects are in the image? Where is the lamp?
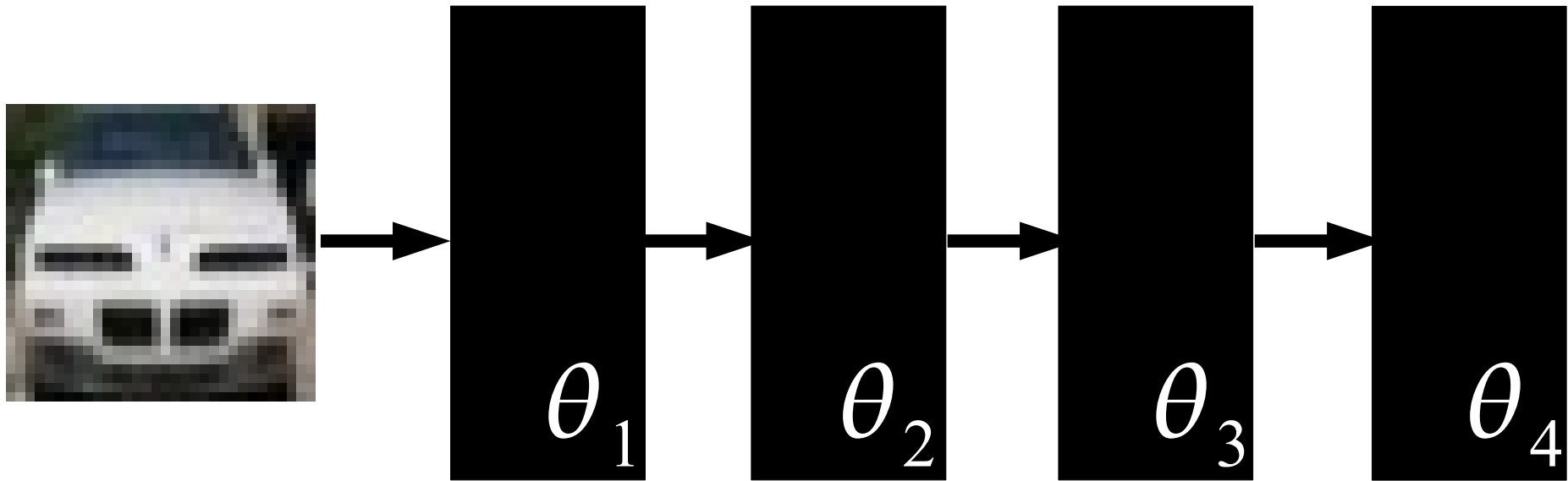What is on the couch? ...

# The Manifold of Natural Images

# Ideal Feature Extraction

E.g.: face images live in about 60-D manifold (x,y,z, pitch, yaw, roll, 53 muscles).

Pixel n

Pixel 2

Pixel 1

**Ideal Feature Extractor**

Pose

Expression

**Ranzato**

Hadsell et al. "Dimensionality reduction by learning an invariant mapping" CVPR 2006

# Deep Learning



$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4$

*Given lots of data, engineer less and learn more!!*
*Let the data find the structure (intrinsic dimensions).*

**Ranzato**

# Deep Learning in Practice

It works very well in practice:

Ranzato

# KEY IDEAS OF DEEP LEARNING

- Hierarchical non-linear system

  - Distributed representations
  - Sharing

- End-to-end learning

  - Joint optimization of features and classifier
  - Good features are learned as a side product of the learning process

**Ranzato**

# THE SPACE OF
# MACHINE LEARNING METHODS

Ranzato

Recurrent
Neural Net

Boosting

Convolutional
Neural Net

Neural Net

Perceptron

SVM

Deep (sparse/denoising)
Autoencoder

Autoencoder
Neural Net

Sparse Coding

ΣΠ

GMM

Deep Belief Net

Restricted BM

BayesNP

Disclaimer: showing only a
subset of the known methods

**SHALLOW**

**DEEP**

Recurrent
Neural Net

Convolutional
Neural Net

Neural Net

Deep (sparse/denoising)
Autoencoder

ΣΠ

Deep Belief Net

BayesNP

Boosting

Perceptron

SVM

Autoencoder
Neural Net

Sparse Coding

GMM

Restricted BM

22

**SHALLOW**

Recurrent
Neural Net

Boosting

Convolutional
Neural Net

Perceptron

Neural Net

SVM

**SUPERVISED**

**UNSUPERVISED**

Deep (sparse/denoising)
Autoencoder

Neural Net

Sparse Coding

ΣΠ

GMM

Deep Belief Net

Restricted BM

**DEEP**

BayesNP

23

**SHALLOW**

**SUPERVISED**

Recurrent Neural Net

Convolutional Neural Net

Neural Net

Boosting

Perceptron

SVM

**UNSUPERVISED**

Deep (sparse/denoising) Autoencoder

**PROBABILISTIC**

Neural Net

Sparse Coding

$\Sigma\Pi$

GMM

Deep Belief Net

Restricted BM

**DEEP**

BayesNP

24

# Deep Learning is B I G

- Main types of deep architectures

**feed-forward**
- Neural nets
- Conv Nets

input

**Feed-back**
- Hierar. Sparse Coding
- Deconv Nets

input

**Bi-directional**
- Stacked Auto-encoders
- DBM

input

**Recurrent**
- Recurrent Neural nets
- Recursive Nets
- LISTA

input

Ranzato

# Deep Learning is B I G

- Main types of deep architectures



feed-forward
- Neural nets
- Conv Nets

Feed-back
- Hierar. Sparse Coding
- Deconv Nets

Bi-directional
- Stacked Auto-encoders
- DBM

Recurrent
- Recurrent Neural nets
- Recursive Nets
- LISTA

input

26

Ranzato

# Deep Learning is B I G

- Main types of learning protocols

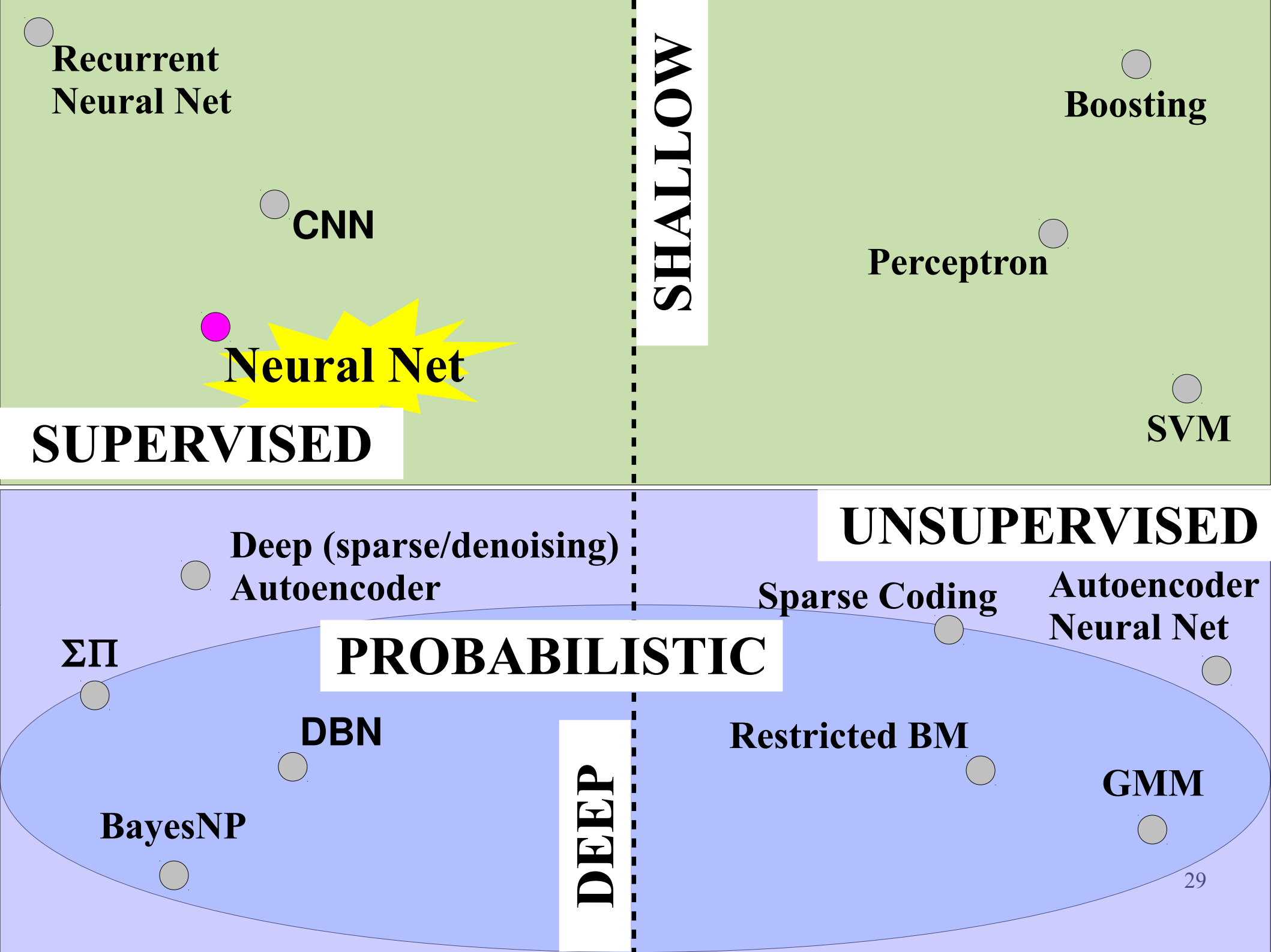  - Purely supervised
    - Backprop + SGD
    - Good when there is lots of labeled data.

  - Layer-wise unsupervised + superv. linear classifier
    - Train each layer in sequence using regularized auto-encoders or RBMs
    - Hold fix the feature extractor, train linear classifier on features
    - Good when labeled data is scarce but there is lots of unlabeled data.

  - Layer-wise unsupervised + supervised backprop
    - Train each layer in sequence
    - Backprop through the whole system
    - Good when learning problem is very difficult.

Ranzato

# Deep Learning is  B I G

- Main types of learning protocols

  - Purely supervised
    - Backprop + SGD
    - Good when there is lots of labeled data.

  - Layer-wise unsupervised + superv. linear classifier
    - Train each layer in sequence using regularized auto-encoders or RBMs
    - Hold fix the feature extractor, train linear classifier on features
    - Good when labeled data is scarce but there is lots of unlabeled data.

  - Layer-wise unsupervised + supervised backprop
    - Train each layer in sequence
    - Backprop through the whole system
    - Good when learning problem is very difficult.

28

Ranzato

**Recurrent Neural Net**

**Boosting**

**CNN**

**Perceptron**

**Neural Net**

**SVM**

**SUPERVISED**

**UNSUPERVISED**

**Deep (sparse/denoising) Autoencoder**

**Sparse Coding**

**Autoencoder Neural Net**

$\Sigma\Pi$

**PROBABILISTIC**

**DBN**

**Restricted BM**

**DEEP**

**GMM**

**BayesNP**

29

# Neural Nets

$$x \rightarrow \boxed{max\,(0, W_1\,x)} \xrightarrow{h_1} \boxed{max\,(0, W_2\,h_1)} \xrightarrow{h_2} \boxed{W_3\,h_2} \xrightarrow{h_3}$$

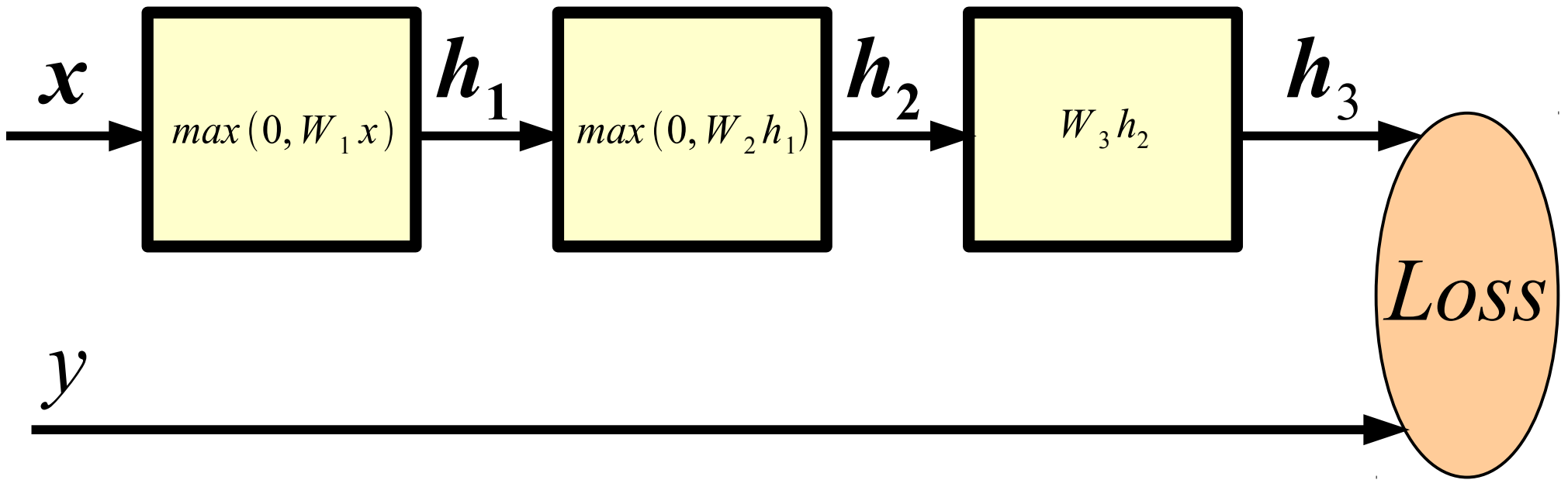**NOTE:** In practice, any (a.e. differentiable) non-linear transformation can be used.

# Computing Loss (example)



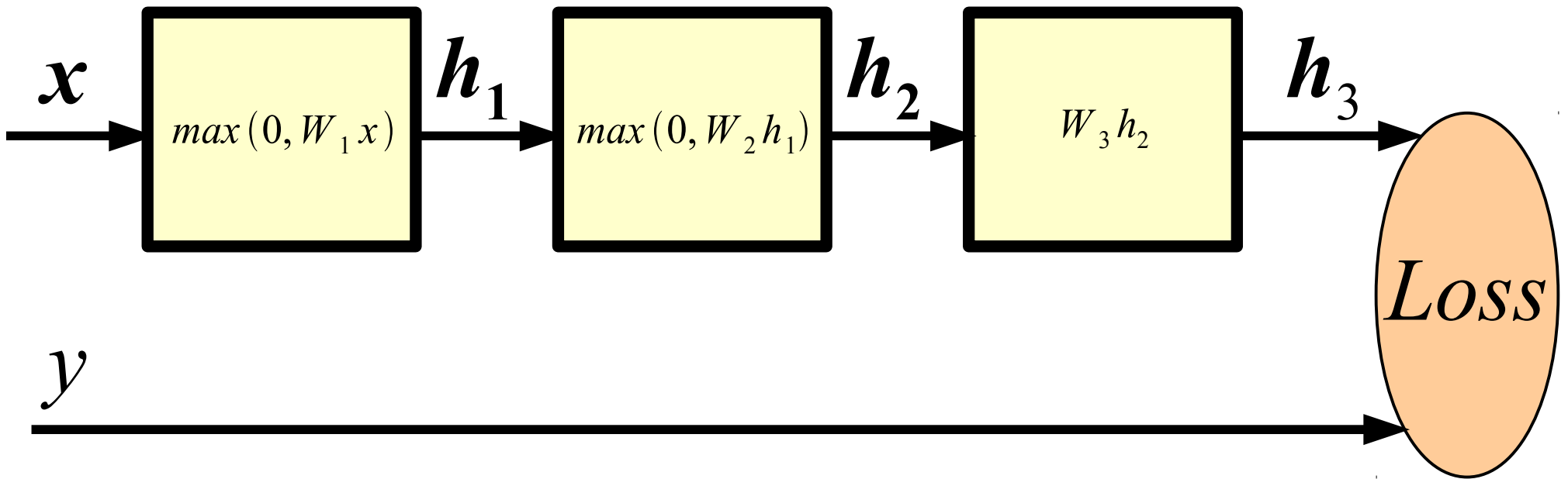$$p(c_k = 1 | x) = \frac{e^{h_{3_k}}}{\sum_j e^{h_{3_j}}}$$   Softmax: probability of class k given input.

$$L(x, y; \theta) = -\sum_j y_j \log p(c_j | x)$$   (Per-sample) Loss: negative log-likelihood.

$$\theta^* = arg\ min_\theta \sum_p L(x^p, y^p; \theta)$$   Learning: min loss (add some regularization).

# Loss



**Q.:** how to tune the parameters to decrease the loss?

If loss is (a.e.) differentiable we can compute gradients.

We can use chain-rule, a.k.a. **back-propagation**, to compute the gradients w.r.t. parameters at the lower layers.

Rumelhart et al. "Learning internal representations by back-propagating.." Nature 1986

# Computing derivative w.r.t. input softmax

$$p(c_k = 1 | x) = \frac{e^{h_{3_k}}}{\sum_j e^{h_{3_j}}}$$
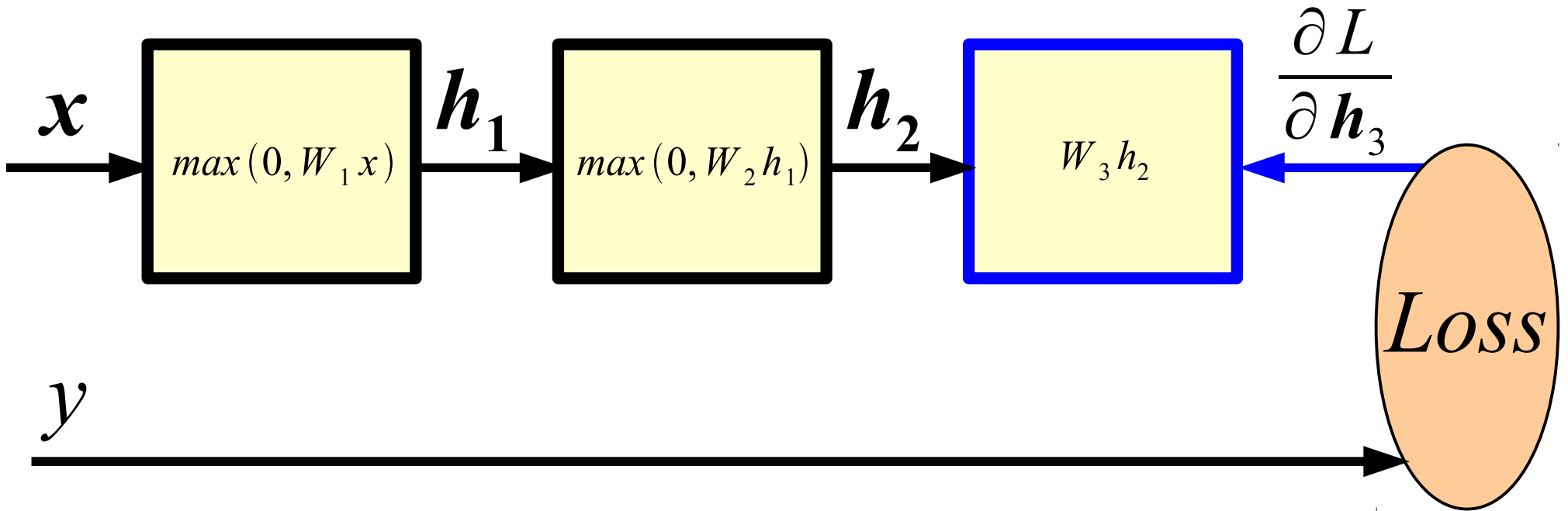
$$L(x, y; \theta) = -\sum_j y_j \log p(c_j | x)$$

By substituting the fist formula in the second, and taking the derivative w.r.t. $h_3$ we get:

$$\frac{\partial L}{\partial h_3} = p(c | x) - y$$
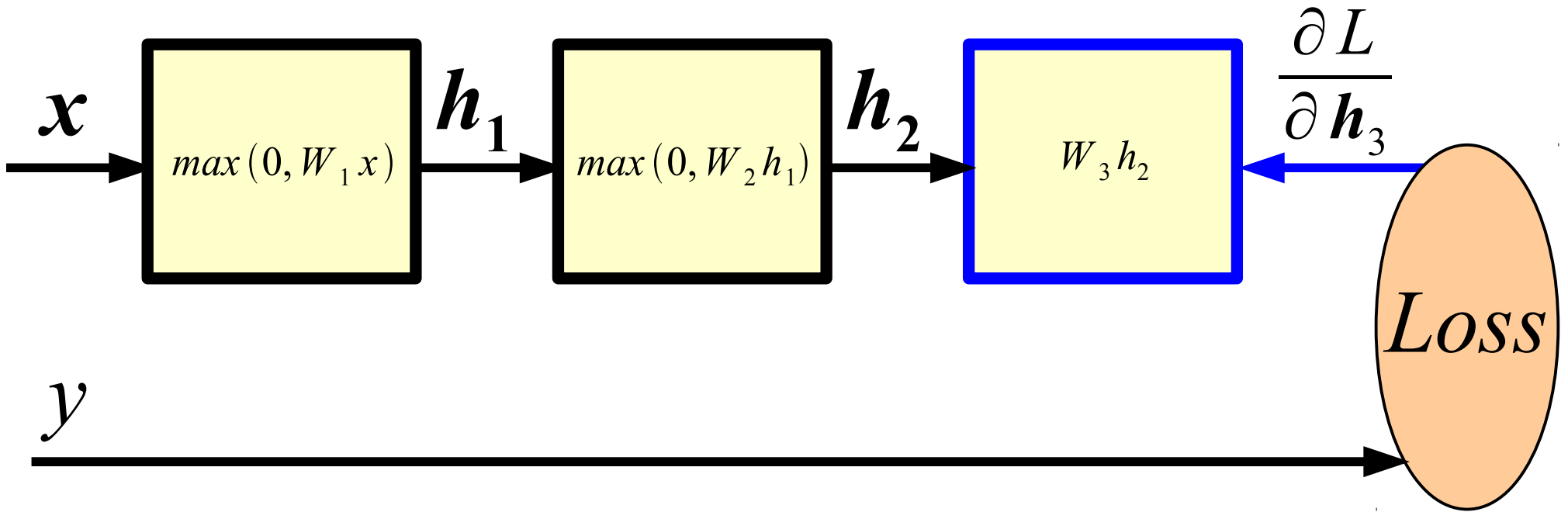
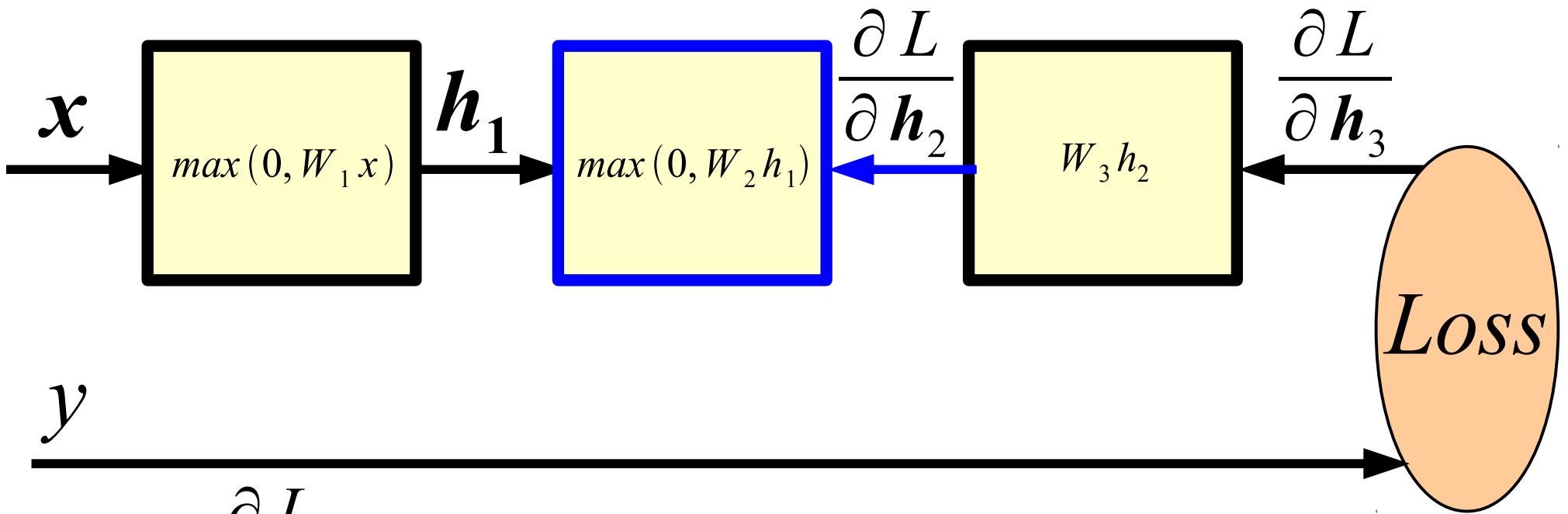**HOMEWORK: prove this equality!**

**Ranzato**

# Backward Propagation



Given $\partial L / \partial \boldsymbol{h}_3$ and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial \boldsymbol{h}_3} \frac{\partial \boldsymbol{h}_3}{\partial W_3} \qquad \frac{\partial L}{\partial \boldsymbol{h}_2} = \frac{\partial L}{\partial \boldsymbol{h}_3} \frac{\partial \boldsymbol{h}_3}{\partial \boldsymbol{h}_2}$$

# Backward Propagation



Given $\partial L / \partial \boldsymbol{h}_3$ and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial \boldsymbol{h}_3} \frac{\partial \boldsymbol{h}_3}{\partial W_3} \qquad \frac{\partial L}{\partial \boldsymbol{h}_2} = \frac{\partial L}{\partial \boldsymbol{h}_3} \frac{\partial \boldsymbol{h}_3}{\partial \boldsymbol{h}_2}$$

$$\frac{\partial L}{\partial W_3} = (p(c|\boldsymbol{x}) - y) \boldsymbol{h}_2^T \qquad \frac{\partial L}{\partial \boldsymbol{h}_2} = W_3^T (p(c|\boldsymbol{x}) - y)$$
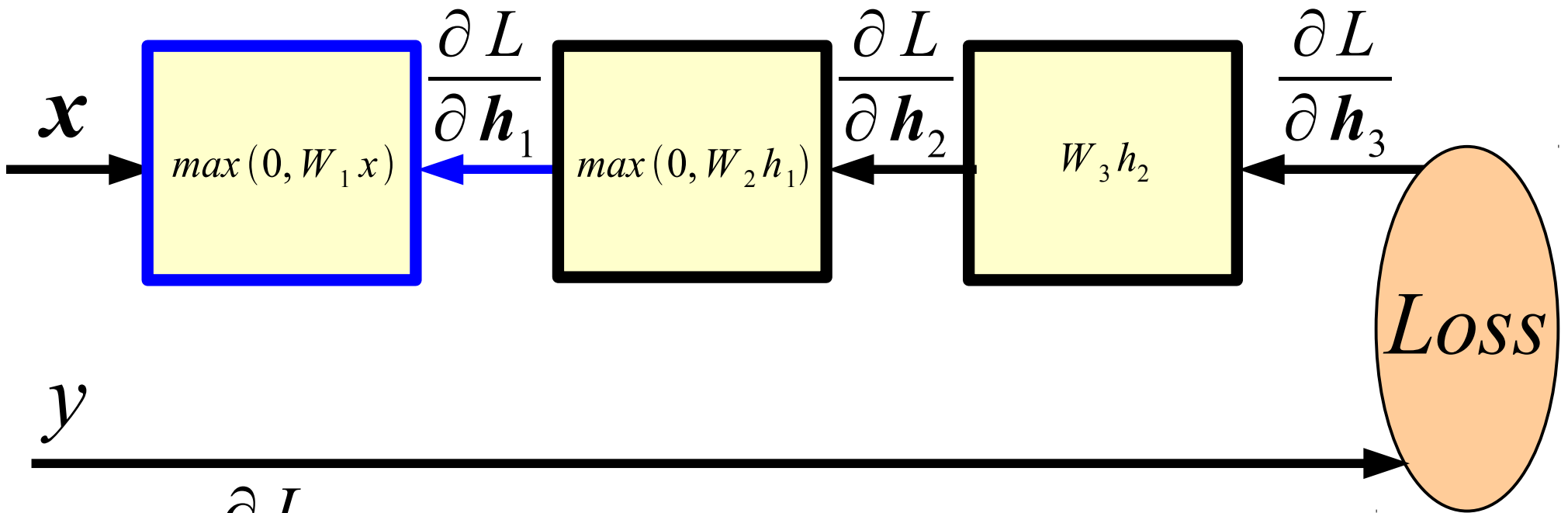
# Backward Propagation



Given $\dfrac{\partial L}{\partial \boldsymbol{h}_2}$ we can compute now:

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial \boldsymbol{h}_2} \frac{\partial \boldsymbol{h}_2}{\partial W_2} \qquad\qquad \frac{\partial L}{\partial \boldsymbol{h}_1} = \frac{\partial L}{\partial \boldsymbol{h}_2} \frac{\partial \boldsymbol{h}_2}{\partial \boldsymbol{h}_1}$$

**HOMEWORK: compute derivatives.**

**Ranzato**

# Backward Propagation



Given $\dfrac{\partial L}{\partial \boldsymbol{h}_1}$ we can compute now:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial \boldsymbol{h}_1} \frac{\partial \boldsymbol{h}_1}{\partial W_1}$$

Ranzato

# Optimization

**Stochastic Gradient Descent** (on mini-batches):

$$\theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta}, \, \eta \in R$$

**Stochastic Gradient Descent with Momentum:**

$$\theta \leftarrow \theta - \eta \Delta$$

$$\Delta \leftarrow 0.9 \Delta + \frac{\partial L}{\partial \theta}$$

**Note: there are many other variants...**

**Ranzato**

# Toy Code (Matlab): Neural Net Trainer

```matlab
% F-PROP
for i = 1 : nr_layers - 1
  [h{i} jac{i}] = nonlinearity(W{i} * h{i-1} + b{i});
end
h{nr_layers-1} = W{nr_layers-1} * h{nr_layers-2} + b{nr_layers-1};
prediction = softmax(h{l-1});


% CROSS ENTROPY LOSS
loss = - sum(sum(log(prediction) .* target)) / batch_size;


% B-PROP
dh{l-1} = prediction - target;
for i = nr_layers - 1 : -1 : 1
  Wgrad{i} = dh{i} * h{i-1}';
  bgrad{i} = sum(dh{i}, 2);
  dh{i-1} = (W{i}' * dh{i}) .* jac{i-1};
end


% UPDATE
for i = 1 : nr_layers - 1
  W{i} = W{i} - (lr / batch_size) * Wgrad{i};
  b{i} = b{i} - (lr / batch_size) * bgrad{i};
end
```
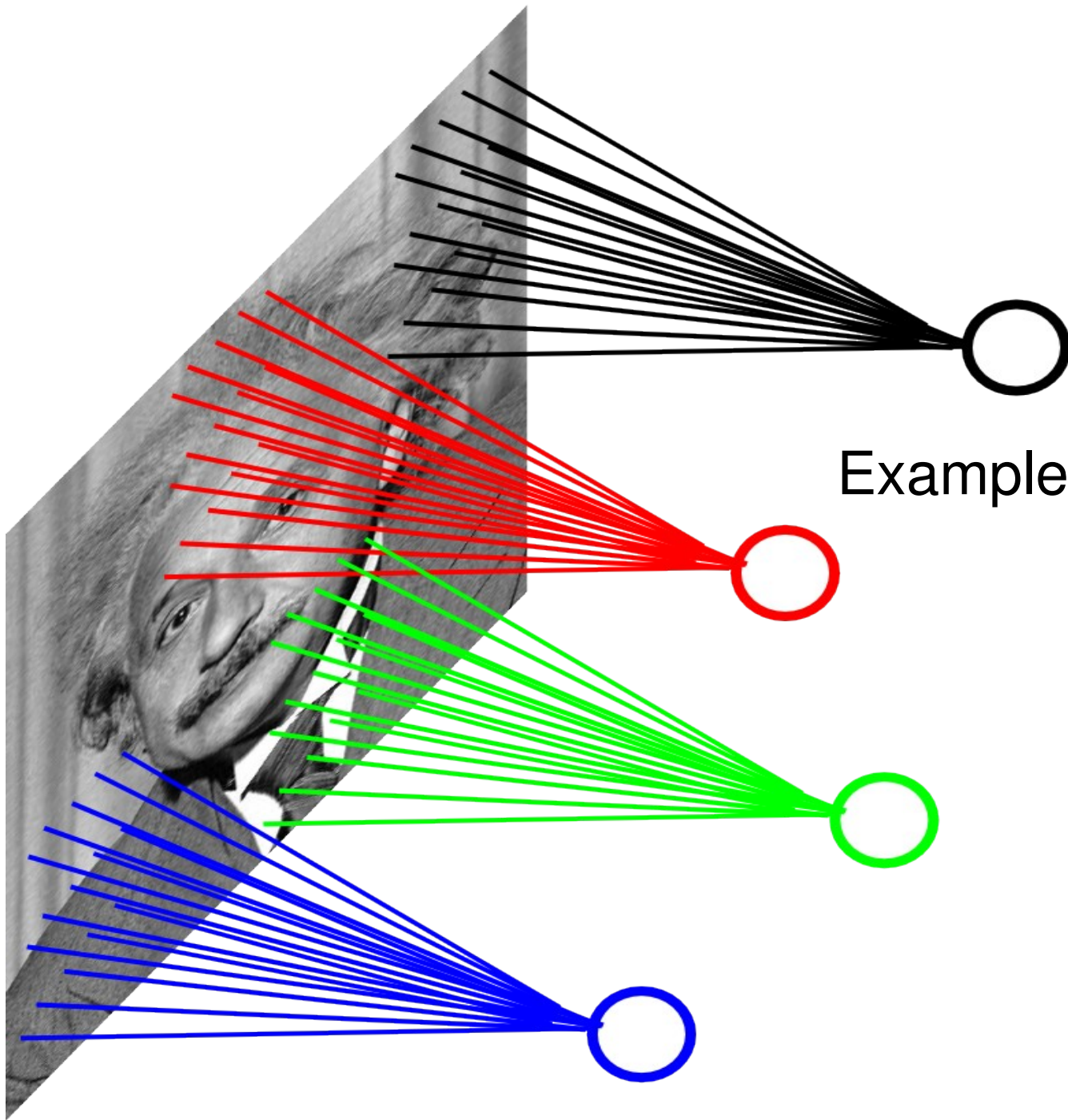
Ranzato

# FULLY CONNECTED NEURAL NET

Example:  200x200 image
40K hidden units
➡ **~2B parameters**!!!



- Spatial correlation is local
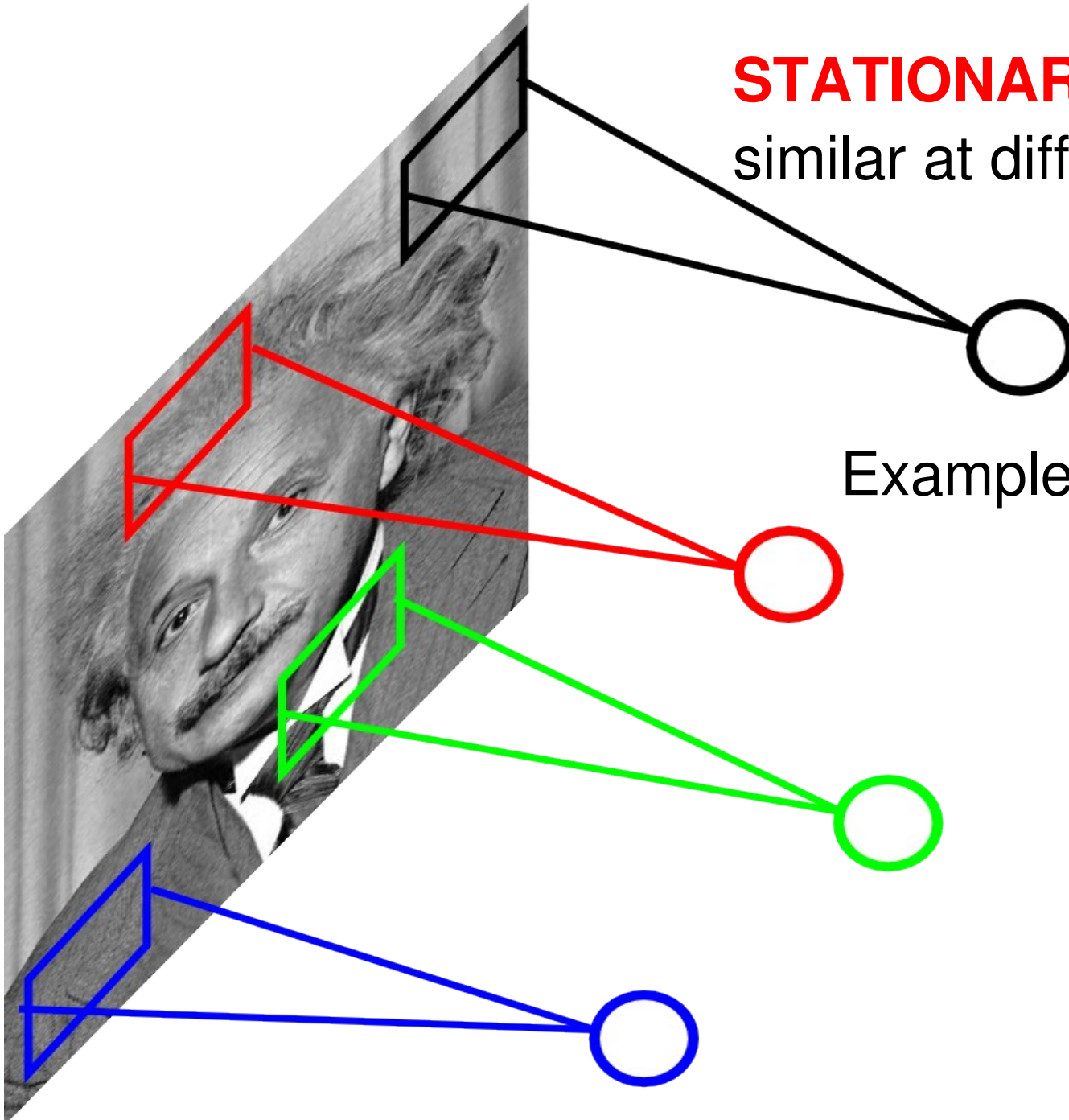- Better to put resources elsewhere!

**Ranzato**

# LOCALLY CONNECTED NEURAL NET

Example: 200x200 image
40K hidden units
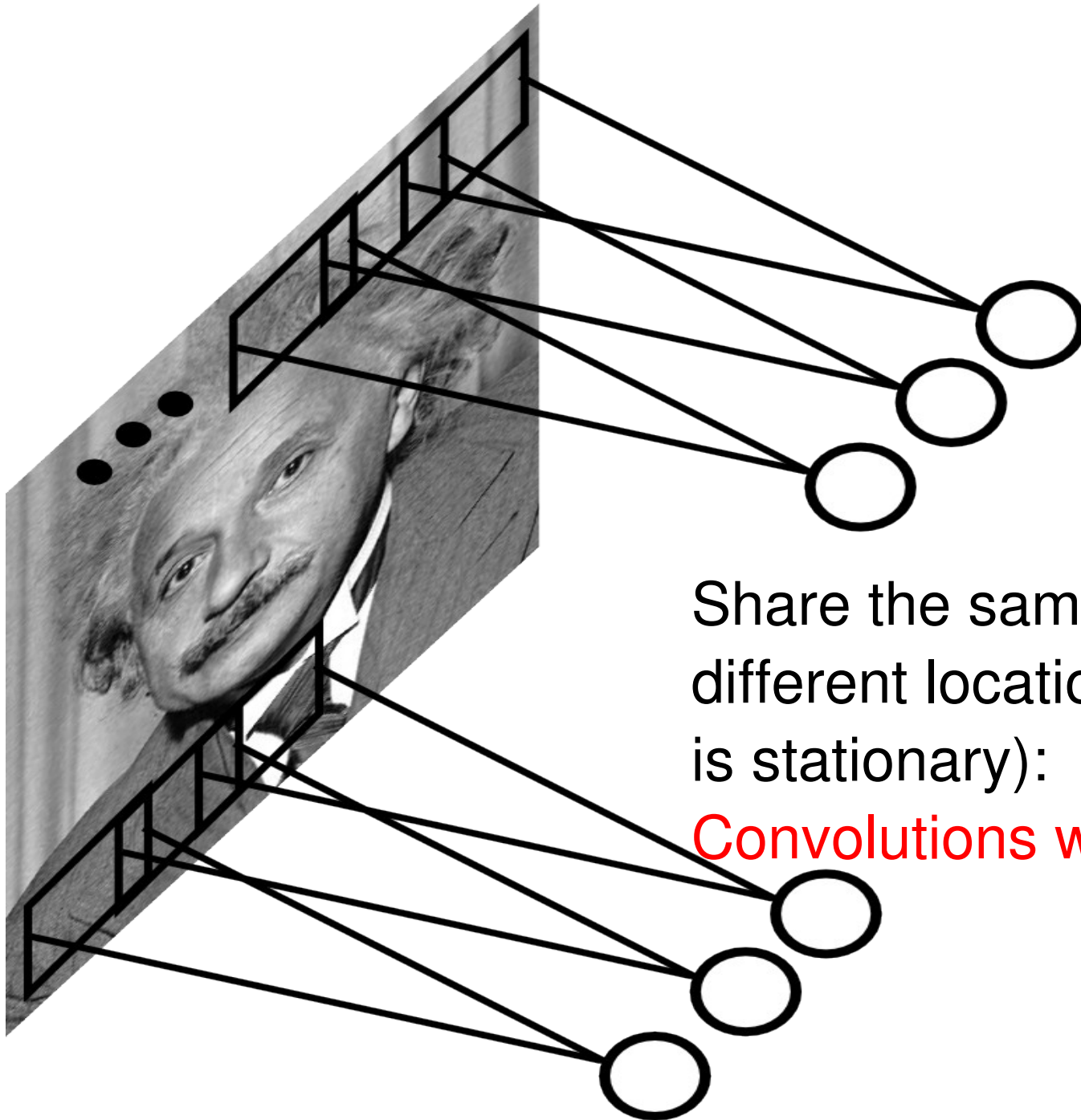Filter size: 10x10
4M parameters

**Ranzato**

# LOCALLY CONNECTED NEURAL NET



**STATIONARITY?** Statistics are similar at different locations

Example: 200x200 image
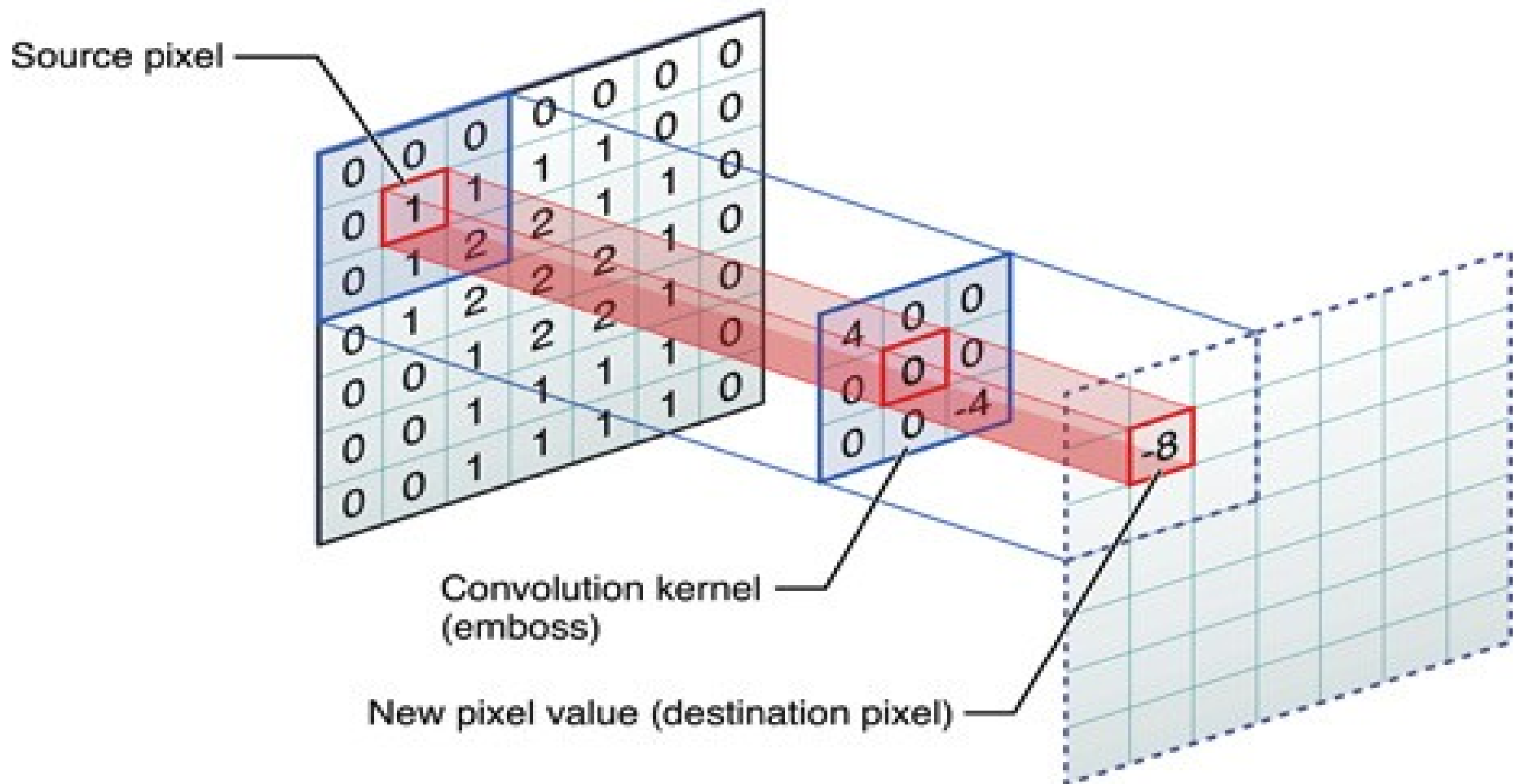40K hidden units
Filter size: 10x10
4M parameters

**Ranzato**

# CONVOLUTIONAL NET



Share the same parameters across different locations (assuming input is stationary):
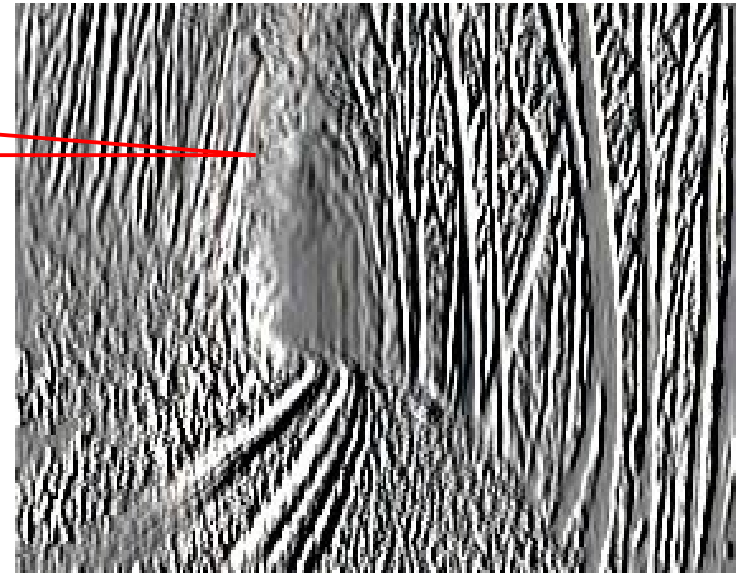
Convolutions with learned kernels

Ranzato

# Convolutional Layer



Source pixel

Convolution kernel
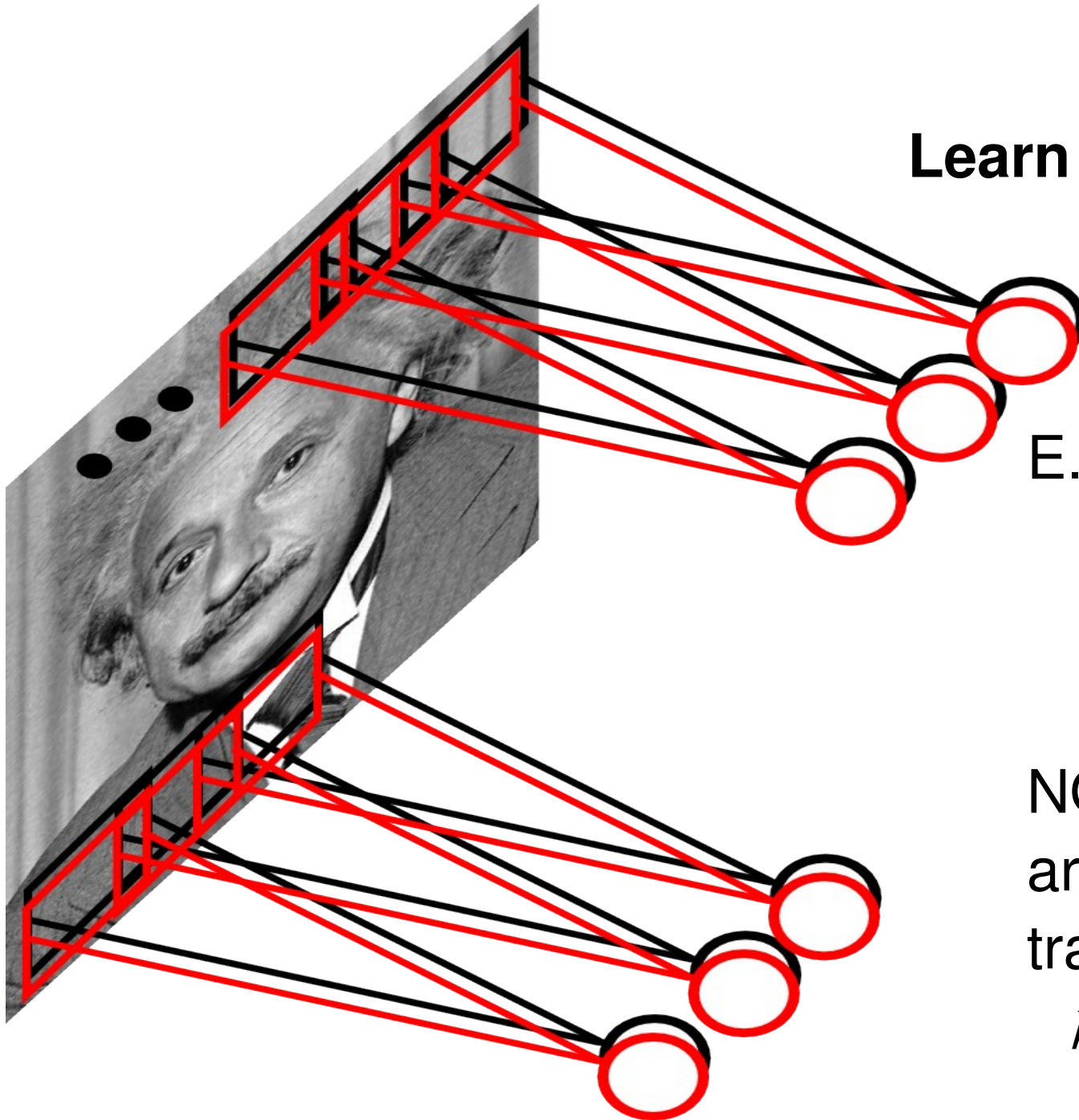(emboss)

New pixel value (destination pixel)

Ranzato

# Convolutional Layer



$$* \begin{matrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{matrix} =$$

# CONVOLUTIONAL NET



**Learn** multiple filters.

E.g.: 200x200 image
     100 Filters
     Filter size: 10x10
     10K parameters

NOTE: filter responses are non-linearly transformed:

$$h = max(0, x * w)$$

47

**Ranzato**

# KEY IDEAS

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

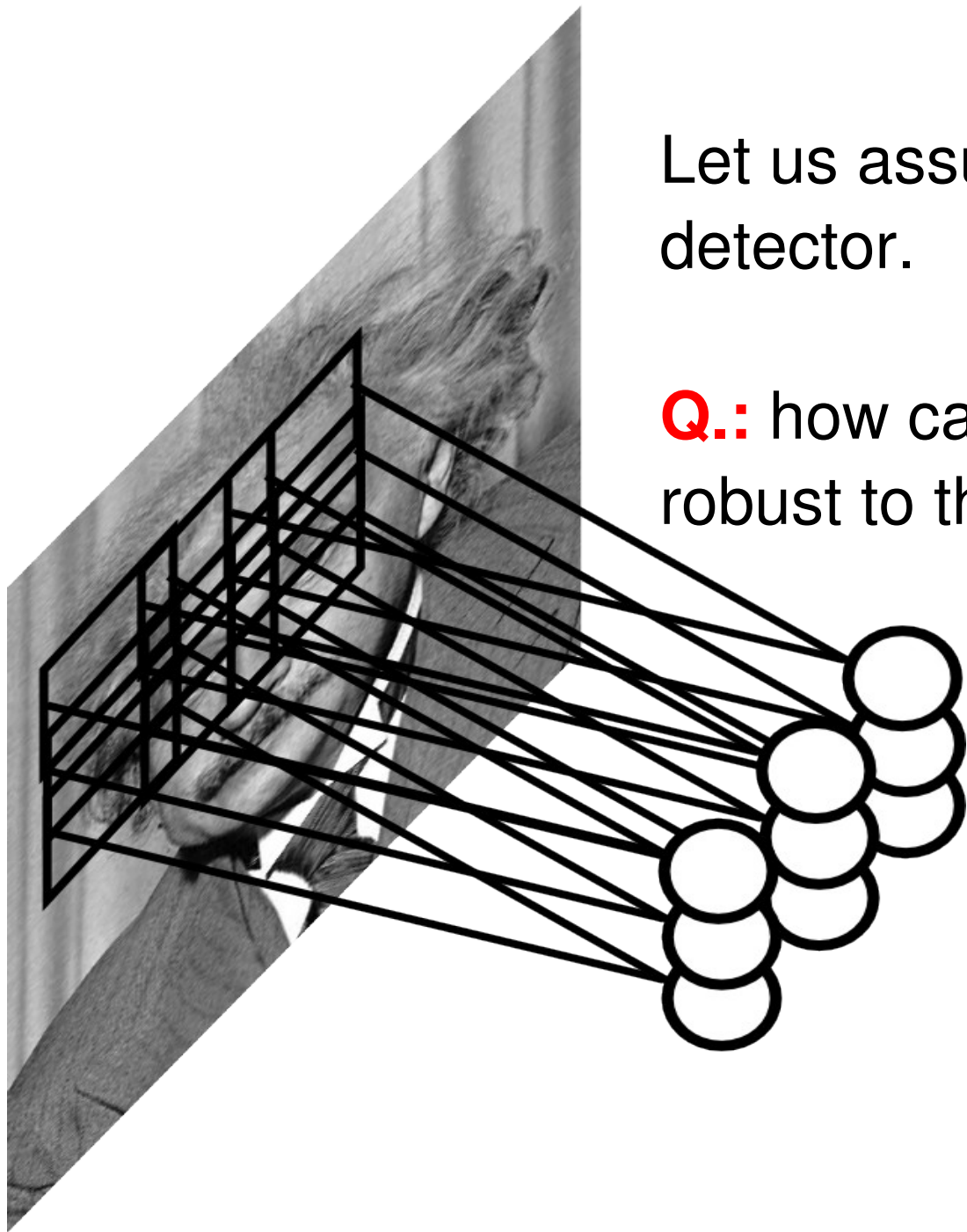- share the weight across hidden units

This is called: **convolutional layer.**
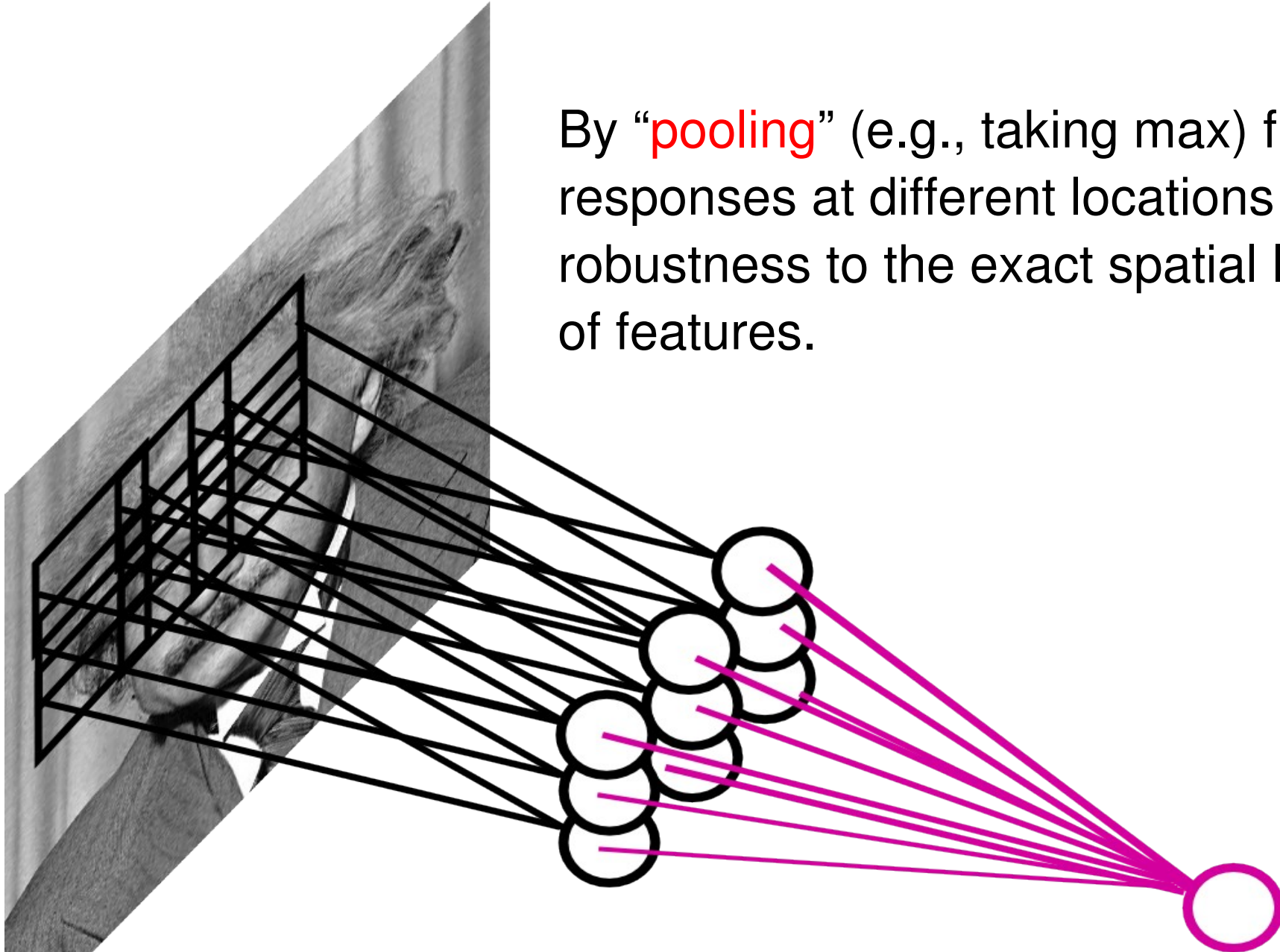A network with convolutional layers is called **convolutional network.**

LeCun et al. "Gradient-based learning applied to document recognition" IEEE 1998

# POOLING



Let us assume filter is an "eye" detector.

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

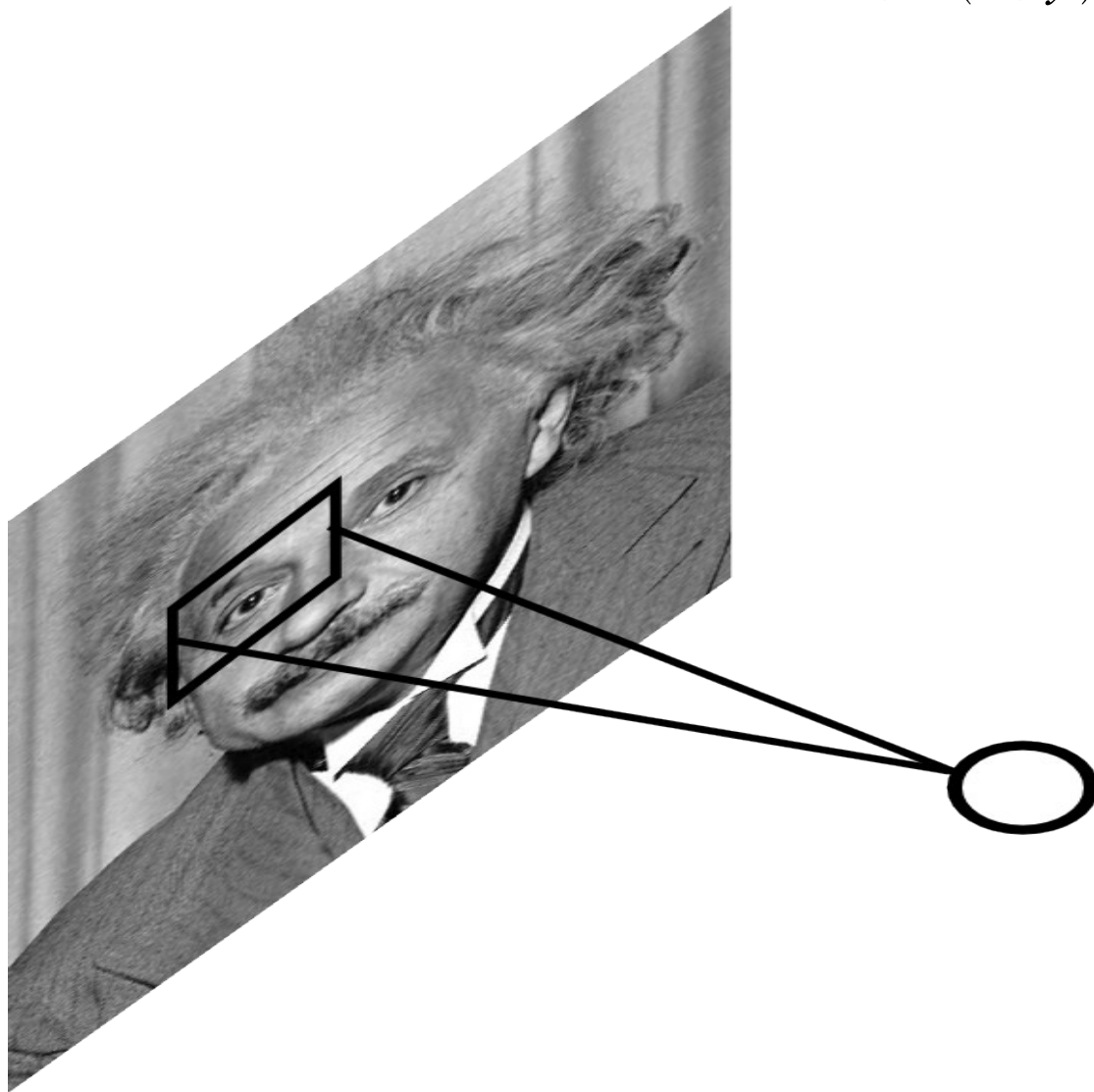# POOLING

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

**Ranzato**

# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$

**Ranzato**

# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$



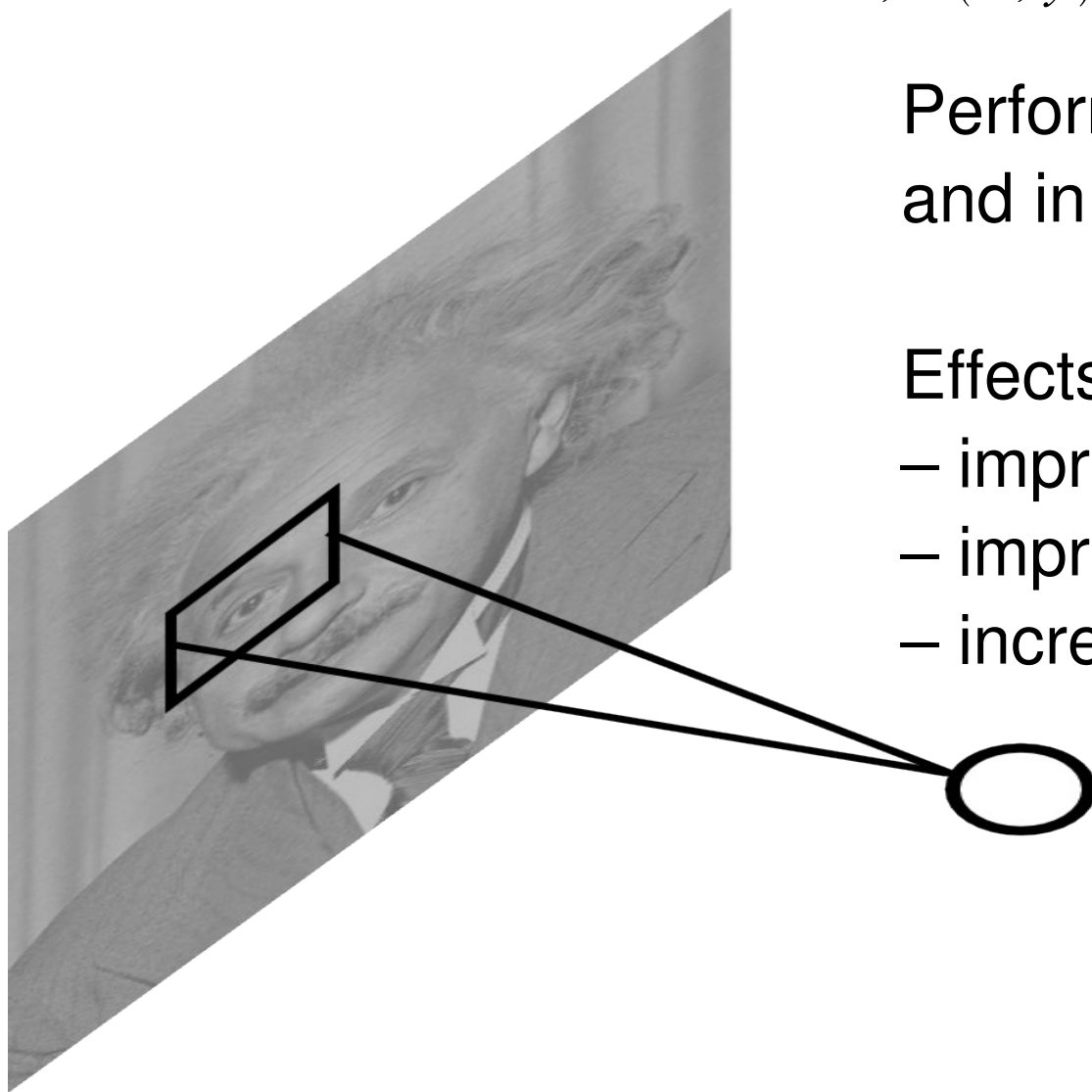We want the same response.

**Ranzato**

# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$



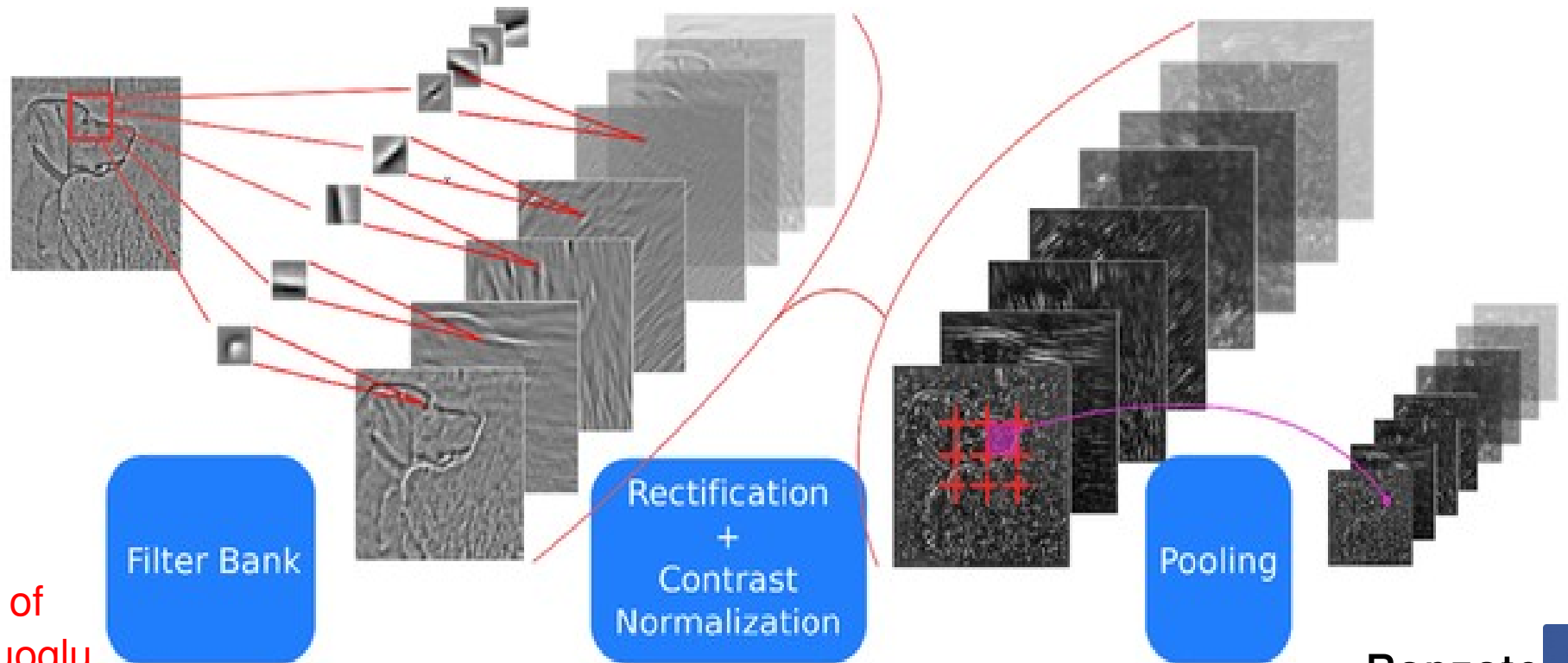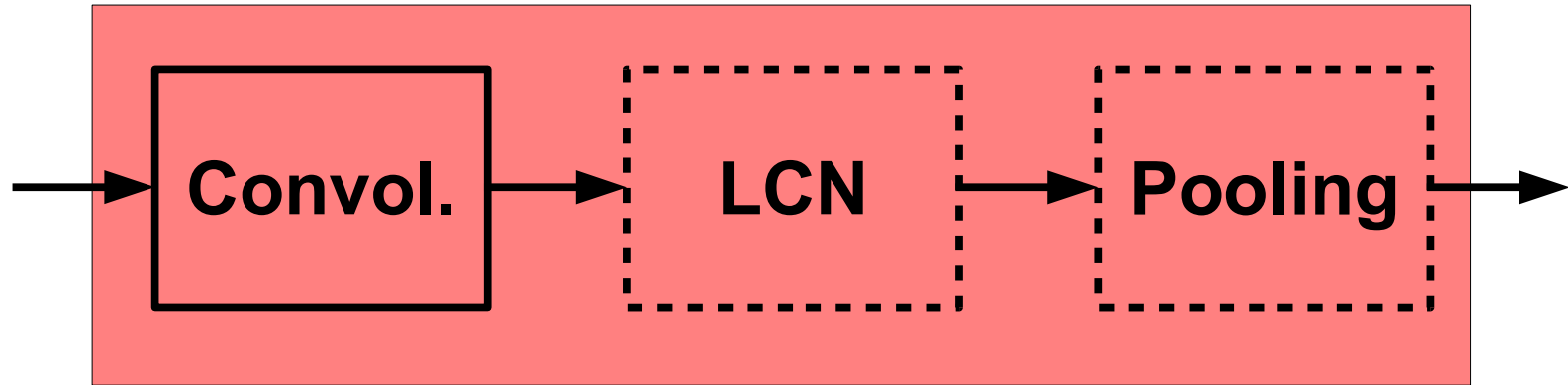Performed also across features and in the higher layers.

Effects:
– improves invariance
– improves optimization
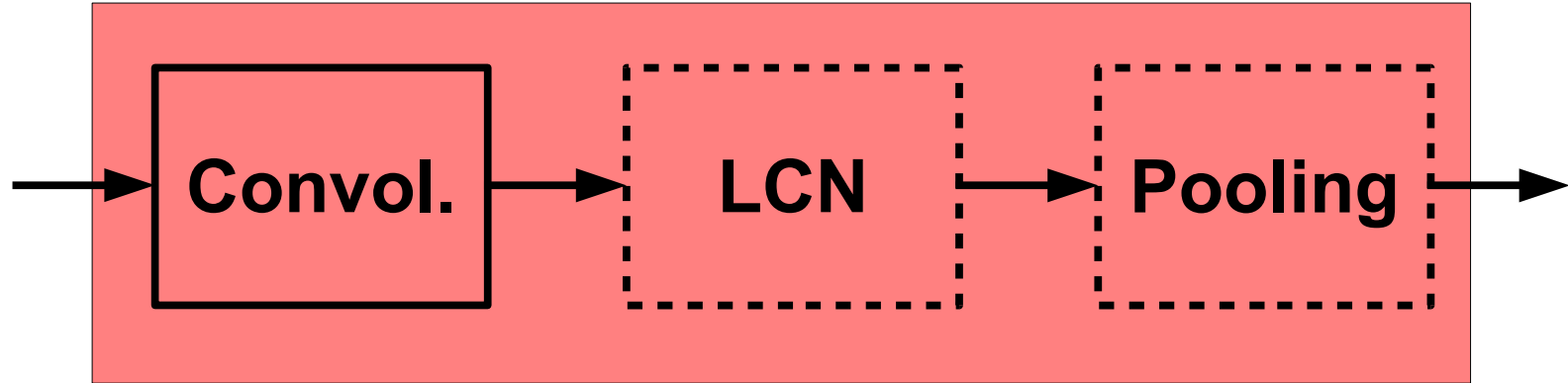– increases sparsity

**Ranzato**

# CONV NETS: TYPICAL ARCHITECTURE

**One stage (zoom)**

Ranzato

# CONV NETS: TYPICAL ARCHITECTURE

**One stage (zoom)**



Conceptually similar to: SIFT, HoG, etc.

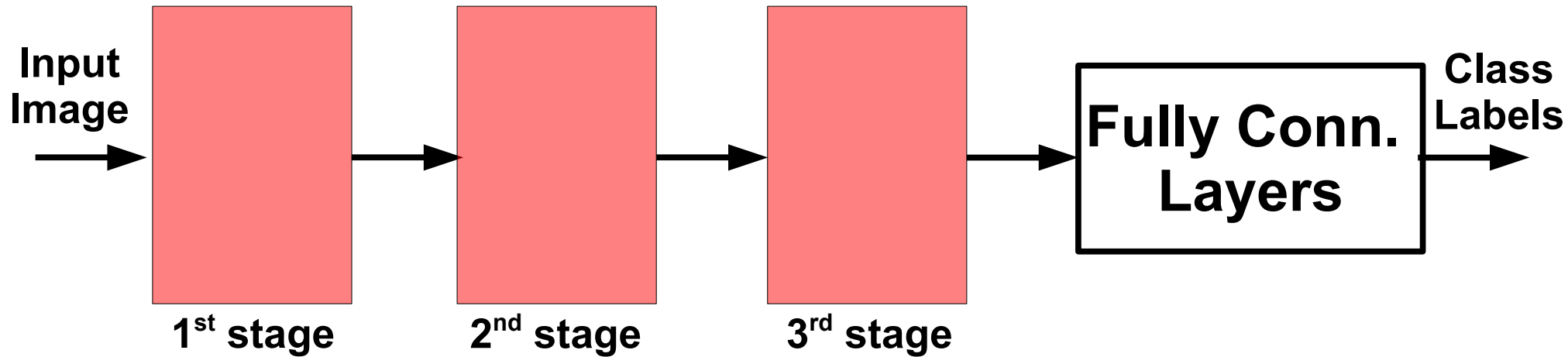**Ranzato**

# CONV NETS: TYPICAL ARCHITECTURE

**One stage (zoom)**



**Whole system**

Ranzato

# CONV NETS: TYPICAL ARCHITECTURE

**Whole system**



Input Image → 1st stage → 2nd stage → 3rd stage → Fully Conn. Layers → Class Labels

Conceptually similar to:

SIFT $\rightarrow$ K-Means $\rightarrow$ Pyramid Pooling $\rightarrow$ SVM

Lazebnik et al. "...Spatial Pyramid Matching..." CVPR 2006

SIFT $\rightarrow$ Fisher Vect. $\rightarrow$ Pooling $\rightarrow$ SVM

Sanchez et al. "Image classifcation with F.V.: Theory and practice" IJCV 2012

Ranzato

# CONV NETS: TRAINING

All layers are differentiable (a.e.).
We can use standard back-propagation.

**Algorithm:**
   **Given a small mini-batch**
   **- F-PROP**
   **- B-PROP**
   **- PARAMETER UPDATE**

**Ranzato**

# CONV NETS: EXAMPLES

- **OCR / House number & Traffic sign classification**

Ciresan et al. "MCDNN for image classification" CVPR 2012
Wan et al. "Regularization of neural networks using dropconnect" ICML 2013

# CONV NETS: EXAMPLES

- **Texture classification**

Sifre et al. "Rotation, scaling and deformation invariant scattering..." CVPR 2013

# CONV NETS: EXAMPLES

- **Pedestrian detection**

Sermanet et al. "Pedestrian detection with unsupervised multi-stage.." CVPR 2013
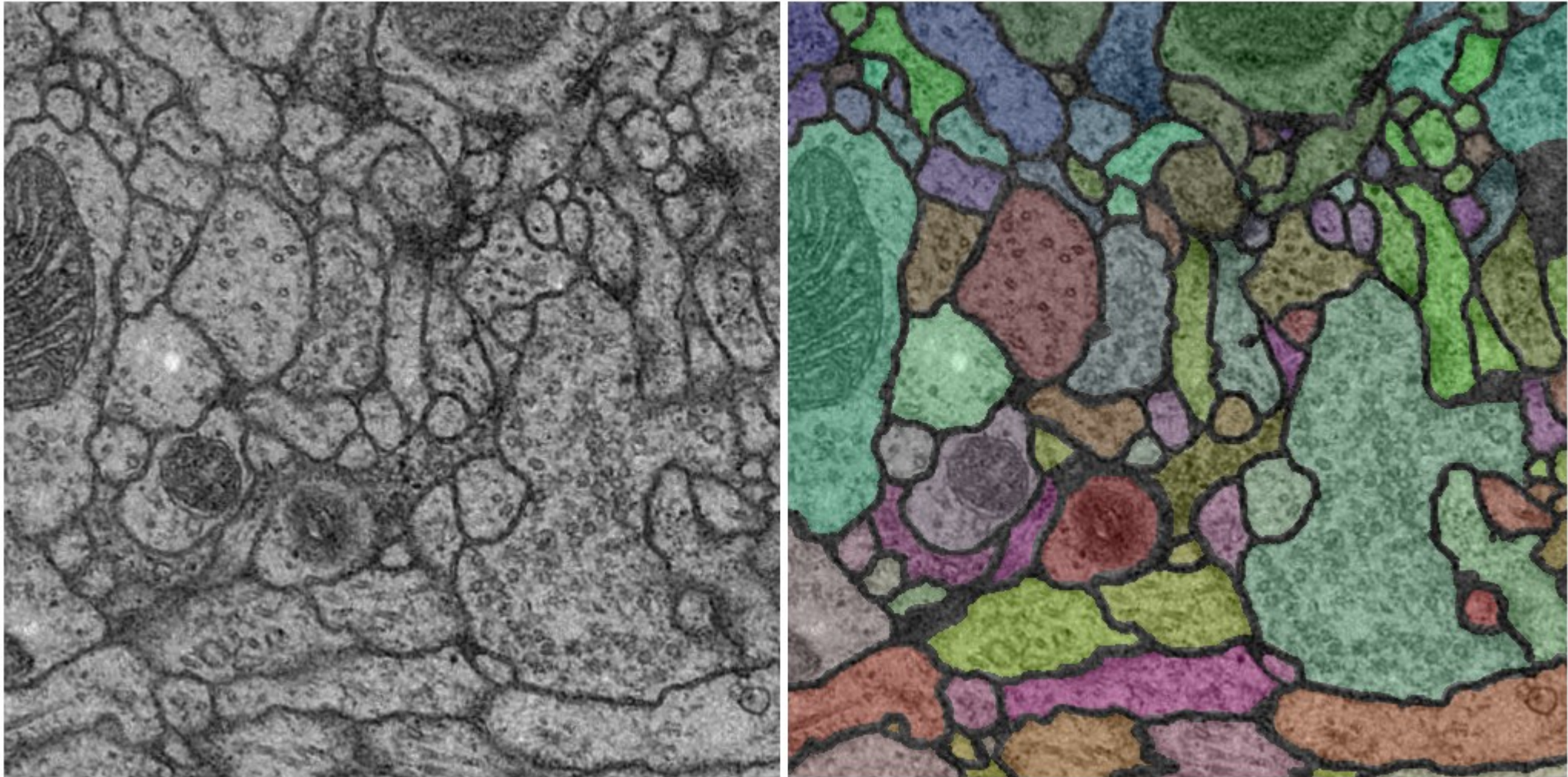
# CONV NETS: EXAMPLES

- **Scene Parsing**



Farabet et al. "Learning hierarchical features for scene labeling" PAMI 2013

Pinheiro et al. "Recurrent CNN for scene parsing" arxiv 2013

**Ranzato**

# CONV NETS: EXAMPLES

- **Segmentation 3D volumetric images**



Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012
Turaga et al. "Maximin learning of image segmentation" NIPS 2009

63

Ranzato

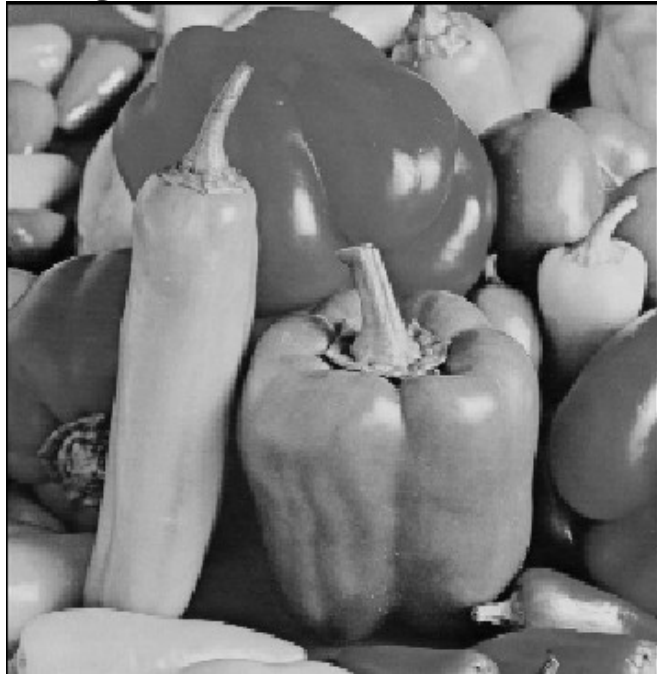# CONV NETS: EXAMPLES
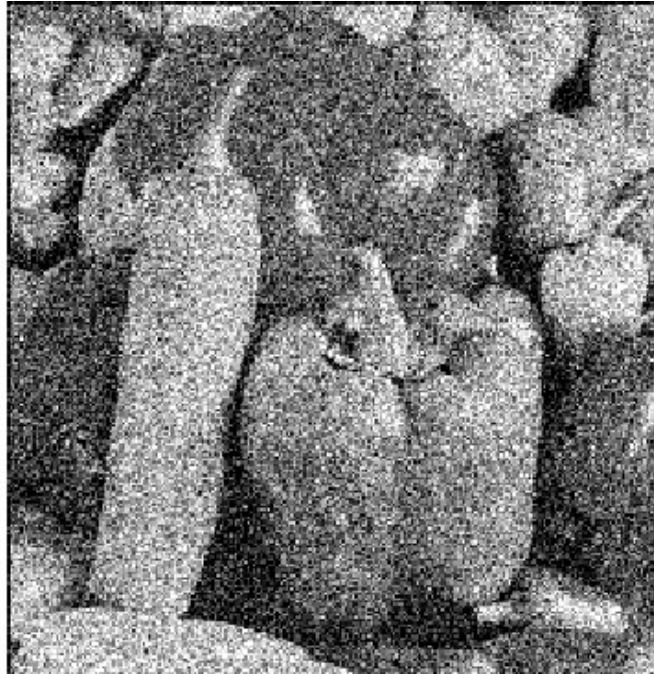
- Action recognition from videos



Taylor et al. "Convolutional learning of spatio-temporal features" ECCV 2010

# CONV NETS: EXAMPLES

- **Robotics**

Sermanet et al. "Mapping and planning ...with long range perception" IROS 2008

- **Denoising**

original    noised    denoised

Burger et al. "Can plain NNs compete with BM3D?" CVPR 2012

**Ranzato**

# CONV NETS: EXAMPLES

- **Dimensionality reduction / learning embeddings**

Hadsell et al. "Dimensionality reduction by learning an invariant mapping" CVPR 2006

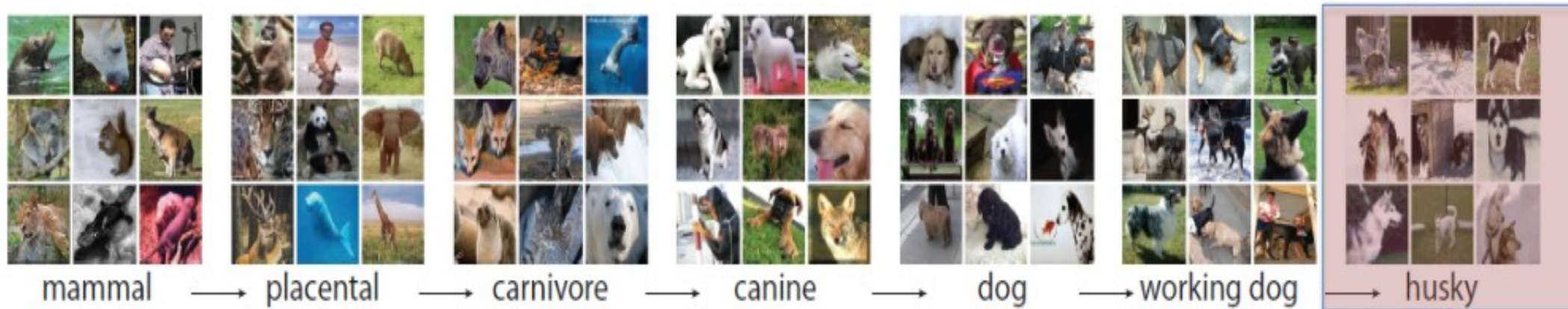# CONV NETS: EXAMPLES

- **Object detection**

Sermanet et al. "OverFeat: Integrated recognition, localization, ..." arxiv 2013
Girshick et al. "Rich feature hierarchies for accurate object detection..." arxiv 2013
Szegedy et al. "DNN for object detection" NIPS 2013

68
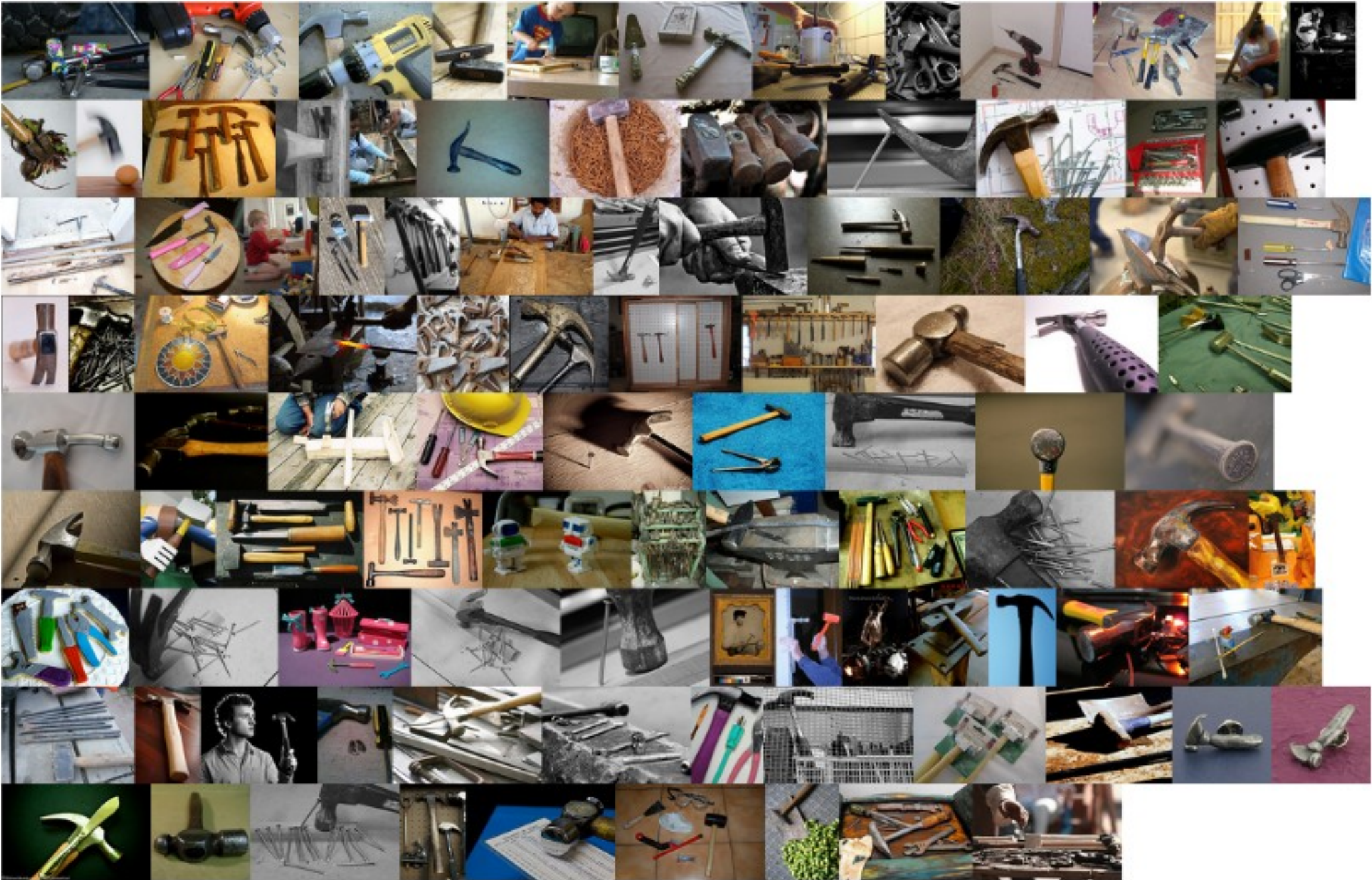
Ranzato

# Dataset: ImageNet 2012



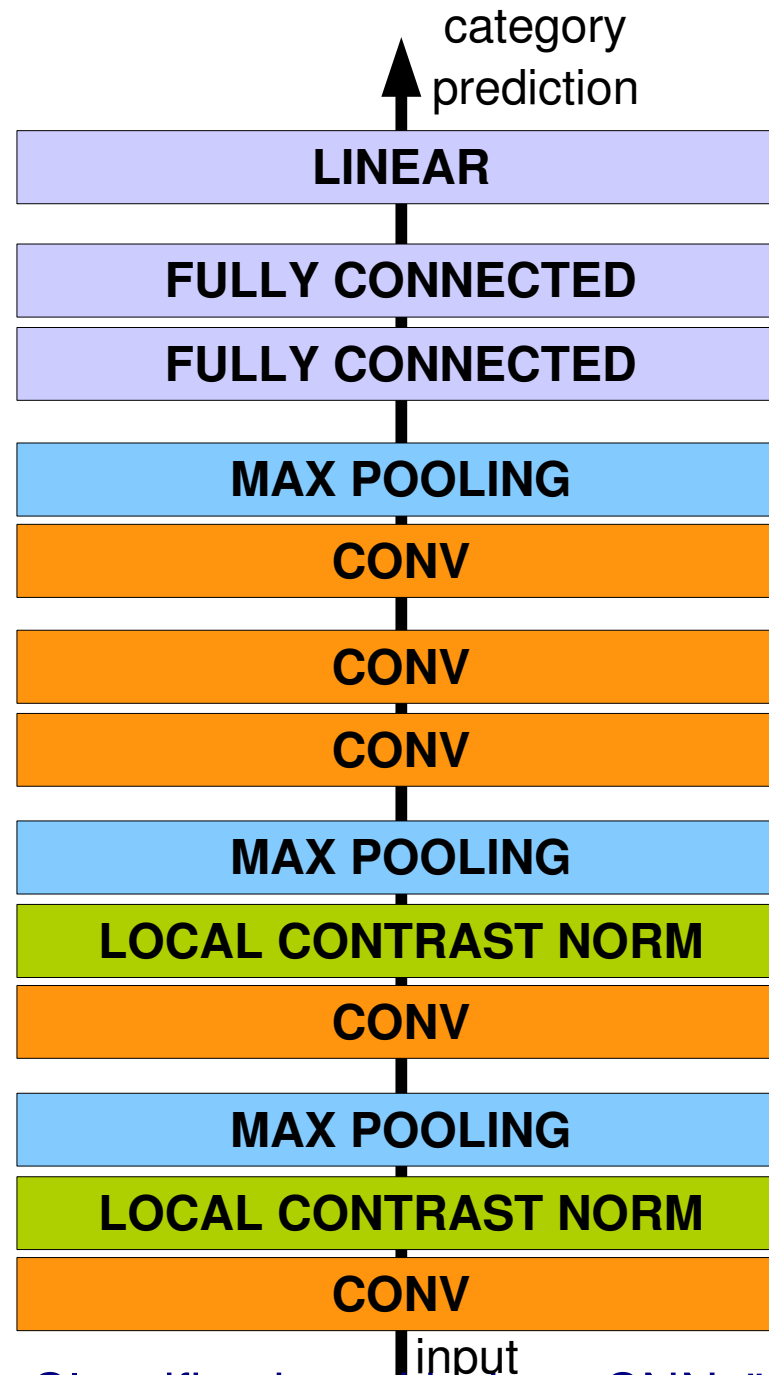mammal → placental → carnivore → canine → dog → working dog → husky

- S: (n) Eskimo dog, **husky** (breed of heavy-coated Arctic sled dog)
  - ○ *direct hypernym* / *inherited hypernym* / *sister term*
    - S: (n) working dog (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
      - S: (n) dog, domestic dog, Canis familiaris (a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) *"the dog barked all night"*
        - S: (n) canine, canid (any of various fissiped mammals with nonretractile claws and typically long muzzles)
          - S: (n) carnivore (a terrestrial or aquatic flesh-eating mammal) *"terrestrial carnivores have four or five clawed digits on each limb"*
            - S: (n) placental, placental mammal, eutherian, eutherian mammal (mammals having a placenta; all mammals except monotremes and marsupials)
              - S: (n) mammal, mammalian (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
                - S: (n) vertebrate, craniate (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
                  - S: (n) chordate (any animal of the phylum Chordata having a notochord or spinal column)
                    - S: (n) animal, animate being, beast, brute, creature, fauna (a living organism characterized by voluntary movement)
                      - S: (n) organism, being (a living thing that has (or can develop) the ability to act or function independently)
                        - S: (n) living thing, animate thing (a living (or once living) entity)
                          - S: (n) whole, unit (an assemblage of parts that is regarded as a single entity) *"how big is that part compared to the whole?"; "the team is a unit"*
                            - S: (n) object, physical object (a tangible and visible entity; an entity that can cast a shadow) *"it was full of rackets, balls and other objects"*
                              - S: (n) physical entity (an entity that has physical existence)
                                - S: (n) entity (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))
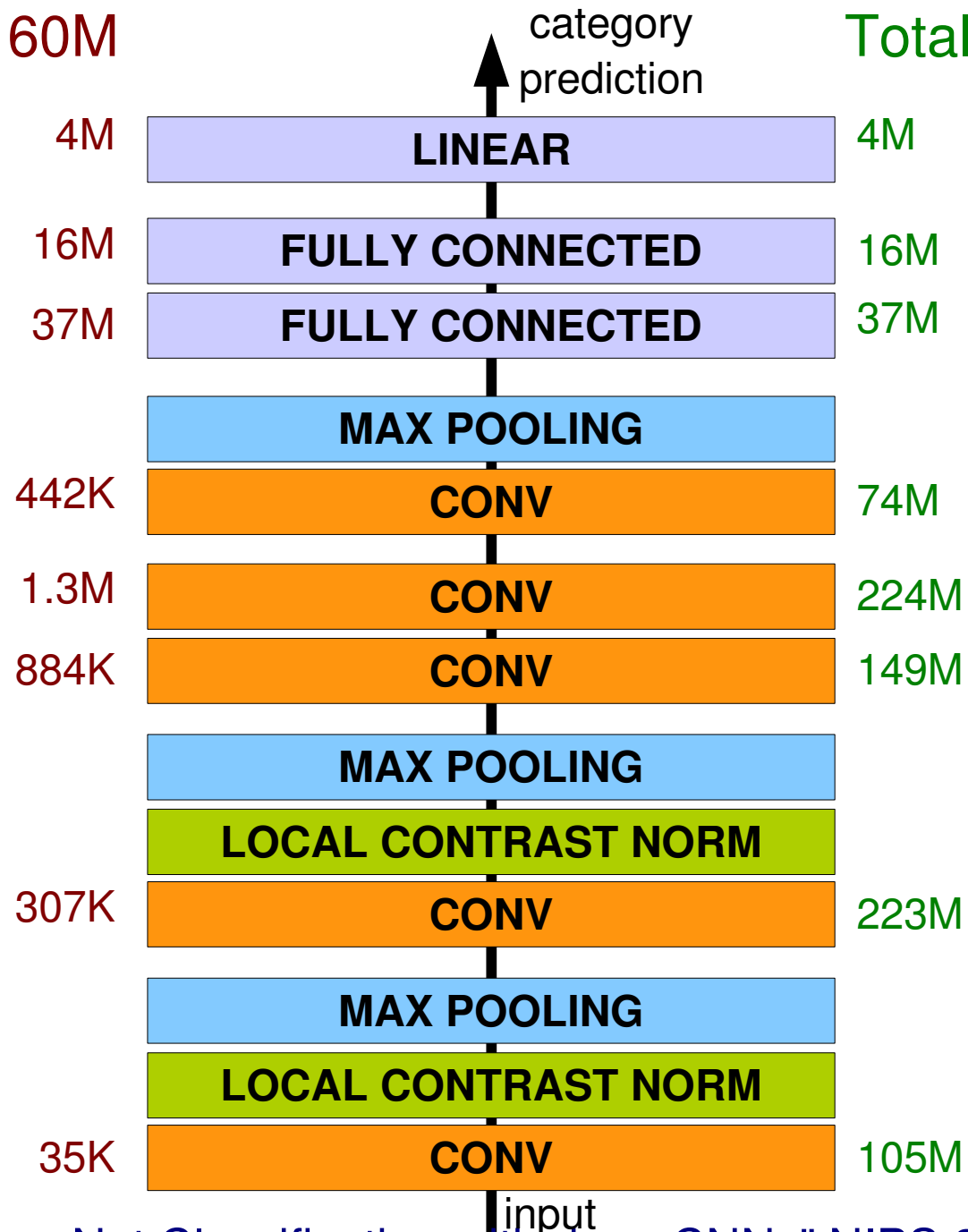
Deng et al. "Imagenet: a large scale hierarchical image database" CVPR 2009

# ImageNet

Examples of hammer:

# Architecture for Classification



category prediction

LINEAR

FULLY CONNECTED

FULLY CONNECTED

MAX POOLING

CONV

CONV

CONV

MAX POOLING

LOCAL CONTRAST NORM

CONV

MAX POOLING

LOCAL CONTRAST NORM

CONV

input

Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

71

Ranzato

# Architecture for Classification

Total nr. params: 60M

category prediction

Total nr. flops: 832M

| | | |
|---|---|---|
| 4M | **LINEAR** | 4M |
| 16M | **FULLY CONNECTED** | 16M |
| 37M | **FULLY CONNECTED** | 37M |
| | **MAX POOLING** | |
| 442K | **CONV** | 74M |
| 1.3M | **CONV** | 224M |
| 884K | **CONV** | 149M |
| | **MAX POOLING** | |
| | **LOCAL CONTRAST NORM** | |
| 307K | **CONV** | 223M |
| | **MAX POOLING** | |
| | **LOCAL CONTRAST NORM** | |
| 35K | **CONV** | 105M |

input

72

Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

**Ranzato**

# Optimization

**SGD with momentum**:

- Learning rate = 0.01

- Momentum = 0.9


**Improving generalization by**:

- Weight sharing (convolution)

- Input distortions

- Dropout = 0.5

- Weight decay = 0.0005

**Ranzato**

# Results: ILSVRC 2012

**Ranzato**

# Results



First layer learned filters (processing raw pixel values).

Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

**Ranzato**

**mite**

| | |
|---|---|
| | mite |
| | black widow |
| | cockroach |
| | tick |
| | starfish |

**container ship**

| | |
|---|---|
| | container ship |
| | lifeboat |
| | amphibian |
| | fireboat |
| | drilling platform |

**motor scooter**

| | |
|---|---|
| | motor scooter |
| | go-kart |
| | moped |
| | bumper car |
| | golfcart |

**leopard**

| | |
|---|---|
| | leopard |
| | jaguar |
| | cheetah |
| | snow leopard |
| | Egyptian cat |

**grille**

| | |
|---|---|
| | convertible |
| | grille |
| | pickup |
| | beach wagon |
| | fire engine |

**mushroom**

| | |
|---|---|
| | agaric |
| | mushroom |
| | jelly fungus |
| | gill fungus |
| | dead-man's-fingers |

**cherry**

| | |
|---|---|
| | dalmatian |
| | grape |
| | elderberry |
| | ffordshire bullterrier |
| | currant |

**Madagascar cat**

| | |
|---|---|
| | squirrel monkey |
| | spider monkey |
| | titi |
| | indri |
| | howler monkey |

# Demo of classifier by Matt Zeiler & Rob Fergus:

## *http://horatio.cs.nyu.edu/*

# Demo of classifier by Yangqing Jia & Trevor Darrell:

## *http://decafberkeleyvision.org/*

Ranzato

Figure 3: How well are we doing? (Left) Classification performance has seen steady improvement in the last few years, both in the number of categories on which algorithms are tested and in classification error rates. (Right) Performance of the best 2006 [Lazebnik et al., 2006] and the best 2007 algorithm [Varma, 2007] are compared here (classification error rates vs number of training examples). One may notice the significant year-on-year progress (see also left panel). Extrapolation enthusiasts may calculate that $10^8$ training examples would be sufficient to achieve 1% error rates with current algorithms. Furthermore, if the pace of year-on-year progress is constant on this log scale chart, 1% error rates with 30 training examples will be achieved in 8-10 years.

Excerpt from Perona Visual Recognition 2007

Donahue, Jia et al. DeCAF arXiv 1310.1531 2013

Ranzato

# CHOOSING THE ARCHITECTURE

- Task dependent

- Cross-validation

- [Convolution $\rightarrow$ LCN $\rightarrow$ pooling]* + fully connected layer

- The more data: the more layers and the more kernels
  - Look at the number of parameters at each layer
  - Look at the number of flops at each layer

- Computational cost

- Be creative :)

**Ranzato**

# HOW TO OPTIMIZE

- SGD (with momentum) usually works very well

- Pick learning rate by running on a subset of the data
  Bottou "Stochastic Gradient Tricks" Neural Networks 2012
  - Start with large learning rate and divide by 2 until loss does not diverge
  - Decay learning rate by a factor of ~1000 or more by the end of training

- Use _/ non-linearity

- Initialize parameters so that each feature across layers has similar variance. Avoid units in saturation.

Ranzato

# HOW TO IMPROVE GENERALIZATION

- Weight sharing (greatly reduce the number of parameters)

- Data augmentation (e.g., jittering, noise injection, etc.)

- Dropout
  Hinton et al. "Improving Nns by preventing co-adaptation of feature detectors" arxiv 2012

- Weight decay (L2, L1)

- Sparsity in the hidden units

- Multi-task (unsupervised learning)

**Ranzato**

# ConvNets: today

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



input/output invariant to permutations

breaking ties between parameters

Saturating units

Loss

parameter

# Neural Net Optimization is...



Like walking on a ridge between valleys

# ConvNets: today

**Loss**

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

**parameter**

# ConvNets: today

Today's belief is that the challenge is about:

– generalization

　　How many training samples to fit 1B parameters?

　　How many parameters/samples to model spaces with 1M dim.?

– scalability

# ConvNets: Why so successful today?

**capacity**

As time goes by, we get more data and more flops/s. The capacity of ML models should grow accordingly.

TIME

1K    1M        1B

**data**

1M

100M

TIME

10T

**flops/s**

Ranzato

# ConvNets: Why so successful today?



capacity

1B

flops/s

data

TIME

CNN were in many ways premature, we did not have enough data and flops/s to train them.

They would overfit and be too slow to tran (apparent local minima).

**NOTE:** methods have to be easily scalable!

Ranzato

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences

- Visualize features (feature maps need to be uncorrelated) and have high variance.



**samples** (vertical axis label)

**hidden unit** (horizontal axis label)

**Good training:** hidden units are sparse across samples and across features.

**Ranzato**

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences

- Visualize features (feature maps need to be uncorrelated) and have high variance.



**samples** (vertical axis)

**hidden unit** (horizontal axis)

**Bad training:** many hidden units ignore the input and/or exhibit strong correlations.

Ranzato

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences

- Visualize features (feature maps need to be uncorrelated) and have high variance.

- Visualize parameters

GOOD

BAD

BAD

BAD

too noisy

too correlated

lack structure

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

Ranzato

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences

- Visualize features (feature maps need to be uncorrelated) and have high variance.

- Visualize parameters
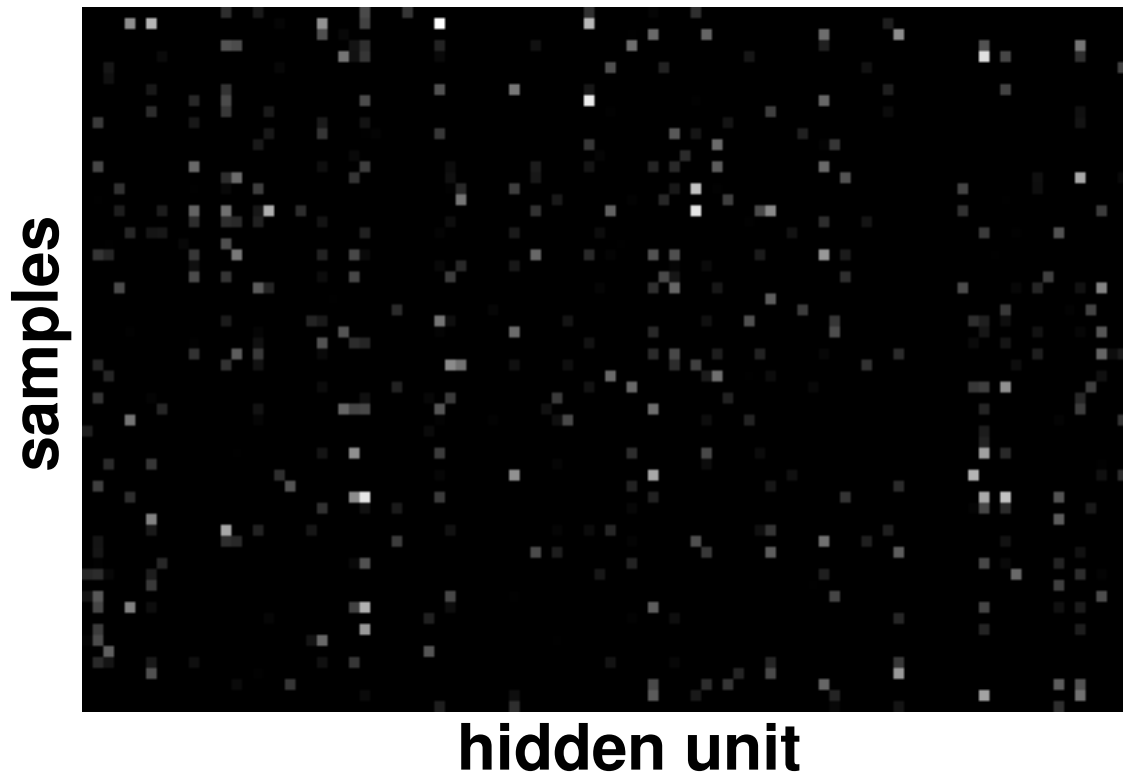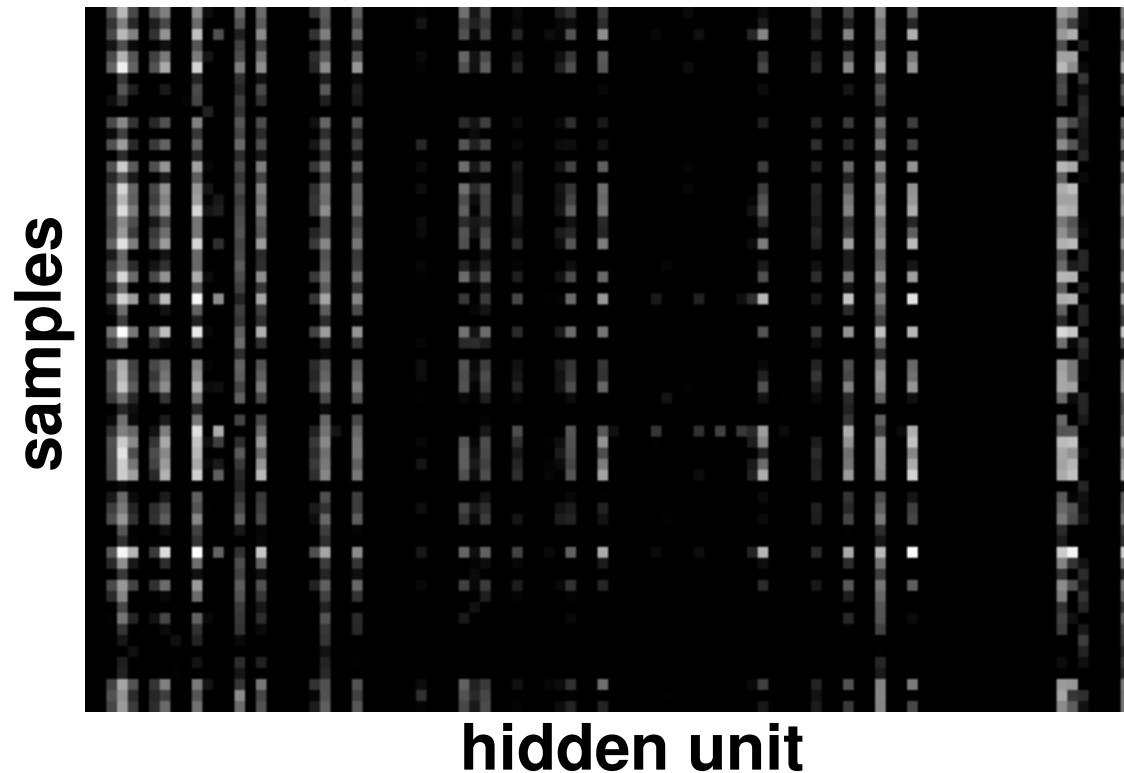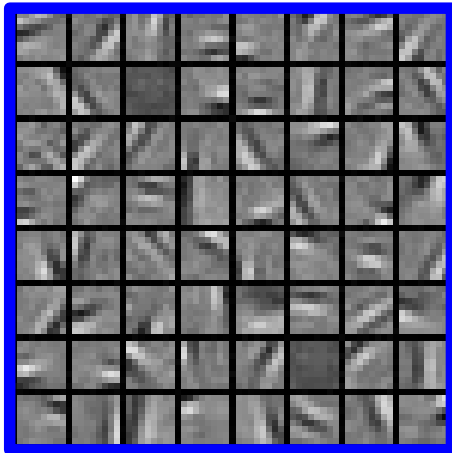
- Measure error on both training and validation set.

- Test on a small subset of the data and check the error $\rightarrow 0$.

**Ranzato**

# WHAT IF IT DOES NOT WORK?

- Training diverges:
  - Learning rate may be too large $\rightarrow$ decrease learning rate
  - BPROP is buggy $\rightarrow$ numerical gradient checking

- Parameters collapse / loss is minimized but accuracy is low
  - 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. $\rightarrow$ if too small, make net larger
  - Visualize hidden units/params $\rightarrow$ fix optmization

- Network is too slow
  - Compute flops and nr. params. $\rightarrow$ GPU,distrib. framework, make net smaller

**Ranzato**

Recurrent
Neural Net

Boosting

CNN

Perceptron

Neural Net

SVM

**SUPERVISED**

**UNSUPERVISED**

Deep (sparse/denoising)
Autoencoder

**Sparse Coding**

Autoencoder
Neural Net

ΣΠ

**PROBABILISTIC**

DBN

Restricted BM

**DEEP**

GMM

BayesNP

96

# Sparse Coding

$$E(\boldsymbol{x},\boldsymbol{h};W)=\frac{1}{2}\|\boldsymbol{x}-W\boldsymbol{h}\|_2^2+\lambda\|\boldsymbol{h}\|_1$$

$$\tilde{E}(\boldsymbol{x};W)=min_{\boldsymbol{h}}E(\boldsymbol{x},\boldsymbol{h};W)$$

$$L=\tilde{E}(\boldsymbol{x};W)$$

Olshausen & Field, Nature 1996                          Ranzato et al. NIPS 2006

# Inference in Sparse Coding

$$E(\boldsymbol{x}, \boldsymbol{h}) = \frac{1}{2}\|\boldsymbol{x} - W_2\boldsymbol{h}\|_2^2 + \lambda\|\boldsymbol{h}\|_1$$



Kavukcuoglu et al. "Predictive Sparse Decomposition" ArXiv 2008

Ranzato

# KEY IDEAS

- Inference can require expensive optimization

- We may approximate exact inference well by using a non-linear function (learn optimal approximation to perform fast inference)

- The original model and the fast predictor can be trained jointly

Kavukcuoglu et al. "Predictive Sparse Decomposition" ArXiv 2008
Kavukcuoglu et al. "Learning convolutonal feature hierarchies.." NIPS 2010
Gregor et al. "Structured sparse coding via lateral inhibition" NIPS 2011
Szlam et al. "Fast approximations to structured sparse coding..." ECCV 2012
Rolfe et al. "Discriminative Recurrent Sparse Autoencoders" ICLR 2013

99

Ranzato

Recurrent
Neural Net

Boosting

CNN

Perceptron

Neural Net

SVM

UNSUPERVISED

Deep (sparse/denoising)
Autoencoder

Sparse Coding

Autoencoder
Neural Net

ΣΠ

PROBABILISTIC

DBN

Restricted BM

GMM

BayesNP

DEEP
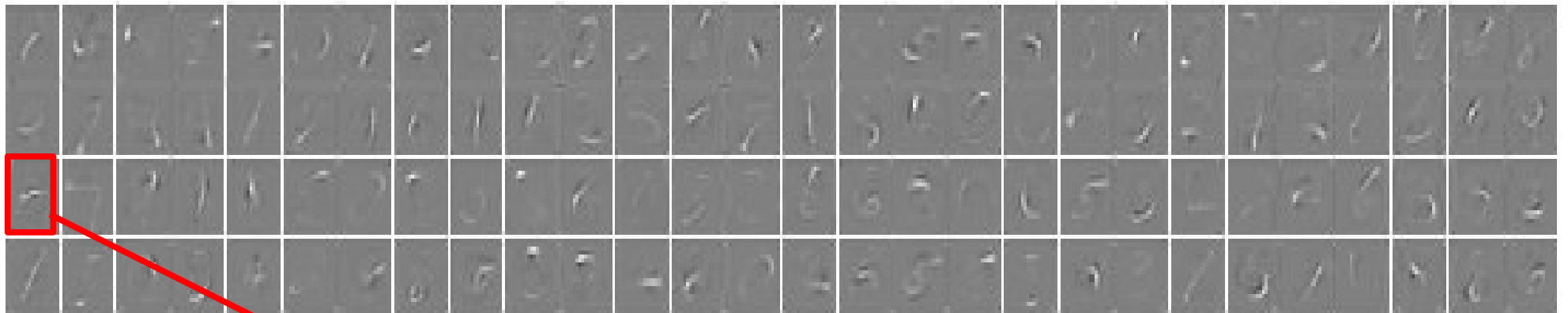
100

# Sampling After Training on Face Images



unconstrained samples

| Original | Input | 1st layer | 2nd layer | 3rd layer | 4th layer | 10 times |

conditional (on the left part of the face) samples

Ranzato

# Expression Recognition Under Occlusion

Ranzato et al. PAMI 2013

Ranzato

1   3   5   7   10   20   30   40   50

Tang et al. Robust BM for decognition and denoising CVPR 2012

**Ranzato**

# Pros

- Feature extraction is fast
- Unprecedented generation quality
- Advances models of natural images
- Trains without labeled data

# Cons

- Training is inefficient
  - Slow
  - Tricky
- Sampling scales badly with dimensionality
- What's the use case of generative models?

# Conclusion

- If generation is not required, other feature learning methods are more efficient (e.g., sparse auto-encoders).
- What's the use case of generative models?
- Given enough labeled data, unsup. learning methods have not produced more useful features.

**Ranzato**

# RNNs

*recurrent neural network handwriting generation demo*

Type a message into the text box, and the network will try to write it out longhand ([this paper](#) explains how it works). Be patient, it can take a while!

**Text** --- up to 100 characters, lower case letters work best

---

**Style** --- either let the network choose a writing style at random or prime it with a real sequence to make it mimic that writer's style.

○ *Take the broth away when they are*

○ *He dismissed the idea*

○ *prison welfare Officer complement*

○ *She looked closely as she*

○ *at Hunterscombe in being adapted for*

◉ random style

**http://www.cs.toronto.edu/~graves/handwriting.html**

# Structured Prediction



LeCun et al. "Gradient-based learning applied to document recognition" IEEE 1998

# Multi-Modal Learning



P A N D A

Frome et al. "DeVISE: A deep visual semantic embedding model" NIPS 2013

Socher et al. Zero-shot learning though cross modal transfer" NIPS 2013

**Ranzato**

# Multi-Modal Learning



(b) Bimodal Deep Autoencoder

Ngiam et al. "Multimodal deep learningl" ICML 2011

Srivastava et al. "Multi-modal learning with DBM" ICML 2012

Ranzato

# SUMMARY

- Deep Learning = Learning Hierarchical representations. Leverage compositionality to gain efficiency.

- Unsupervised learning: active research topic.

- Supervised learning: today it is the most successful set up.

- 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 (Google)
  - Better optimization techniques

- Generalization on small datasets (curse of dimensionality):
  - data augmentation
  - weight decay
  - dropout

**Ranzato**

# SOFTWARE

**Torch7: learning library that supports neural net training**

http://www.torch.ch

http://code.cogbits.com/wiki/doku.php  (tutorial with demos by C. Farabet)

**Python-based learning library  (U. Montreal)**

- http://deeplearning.net/software/theano/  (does automatic differentiation)

**Efficient CUDA kernels for ConvNets  (Krizhevsky)**

– code.google.com/p/cuda-convnet

**Caffe (Yangqing Jia)**

– http://caffe.berkeleyvision.org

Ranzato

# REFERENCES

## Convolutional Nets

– LeCun, Bottou, Bengio and Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998

- Krizhevsky, Sutskever, Hinton "ImageNet Classification with deep convolutional neural networks" NIPS 2012

– Jarrett, Kavukcuoglu, Ranzato, LeCun: What is the Best Multi-Stage Architecture for Object Recognition?, Proc. International Conference on Computer Vision (ICCV'09), IEEE, 2009

- Kavukcuoglu, Sermanet, Boureau, Gregor, Mathieu, LeCun: Learning Convolutional Feature Hierachies for Visual Recognition, Advances in Neural Information Processing Systems (NIPS 2010), 23, 2010

– see  yann.lecun.com/exdb/publis  for references on many different kinds of convnets.

– see http://www.cmap.polytechnique.fr/scattering/ for scattering networks (similar to convnets but with less learning and stronger mathematical foundations)

– see http://www.idsia.ch/~juergen/ for other references to ConvNets and LSTMs.

Ranzato

# REFERENCES

## Applications of Convolutional Nets

– Farabet, Couprie, Najman, LeCun. Scene Parsing with Multiscale Feature Learning, Purity Trees, and Optimal Covers", ICML 2012

– Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala and Yann LeCun: Pedestrian Detection with Unsupervised Multi-Stage Feature Learning, CVPR 2013

- D. Ciresan, A. Giusti, L. Gambardella, J. Schmidhuber. Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images. NIPS 2012

- Raia Hadsell, Pierre Sermanet, Marco Scoffier, Ayse Erkan, Koray Kavackuoglu, Urs Muller and Yann LeCun. Learning Long-Range Vision for Autonomous Off-Road Driving, Journal of Field Robotics, 26(2):120-144, 2009

– Burger, Schuler, Harmeling. Image Denoisng: Can Plain Neural Networks Compete with BM3D?, CVPR 2012

– Hadsell, Chopra, LeCun. Dimensionality reduction by learning an invariant mapping, CVPR 2006

– Bergstra et al. Making a science of model search: hyperparameter optimization in hundred of dimensions for vision architectures, ICML 2013

Ranzato

# REFERENCES

## Deep Learning in general

– deep learning tutorial slides at ICML 2013

– Yoshua Bengio, Learning Deep Architectures for AI, Foundations and Trends in Machine Learning, 2(1), pp.1-127, 2009.

– LeCun, Chopra, Hadsell, Ranzato, Huang: A Tutorial on Energy-Based Learning, in Bakir, G. and Hofman, T. and Schölkopf, B. and Smola, A. and Taskar, B. (Eds), Predicting Structured Data, MIT Press, 2006

Ranzato

# THANK YOU

*Ackknowledgements to Yann LeCun. Many slides from ICML 2013 and CVPR 2013 tutorial on deep learning.*

**Ranzato**