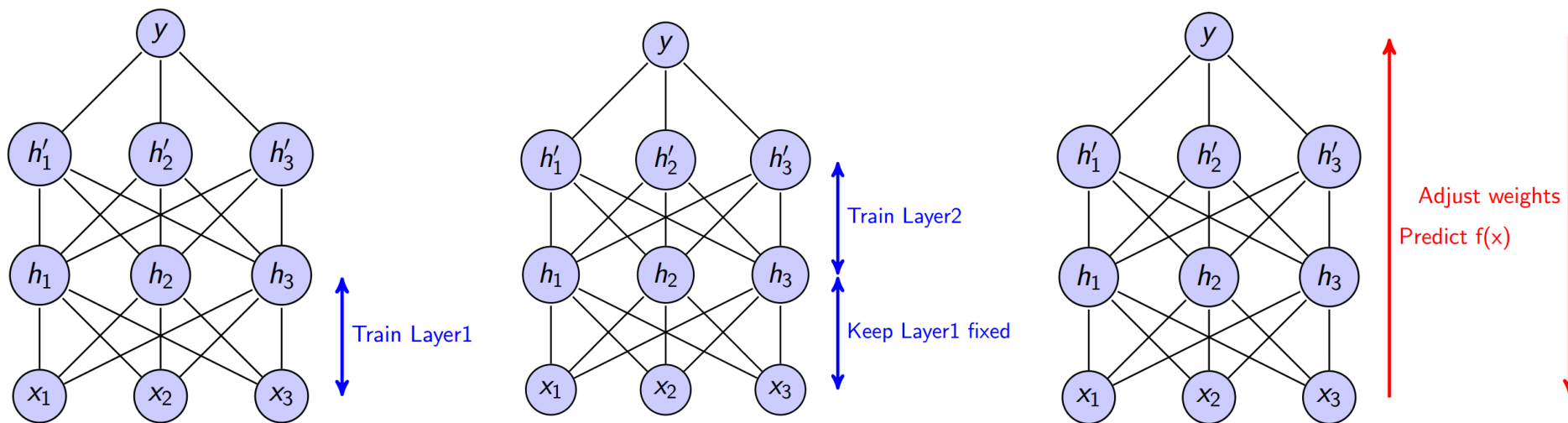# UNSUPERVISED LEARNING

# Topics

- Layer-wise (unsupervised) pre-training
  - Restricted Boltzmann Machines
  - Auto-encoders

# LAYER-WISE (UNSUPERVISED) PRE-TRAINING

# Breakthrough in 2006

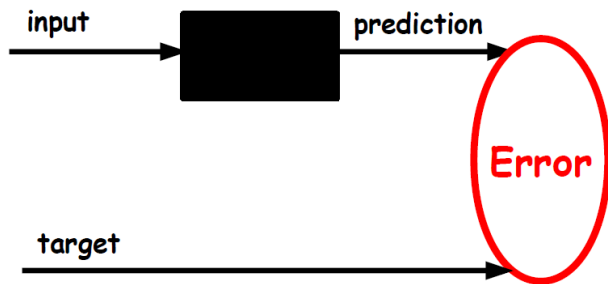- Layer-wise (unsupervised) pre-training

# Breakthrough in 2006

- Key idea:
  - Focus on modeling the input $P(X)$ better with each successive layer. Worry about optimizing the task $P(Y|X)$ later.
    - X=observation, Y=label
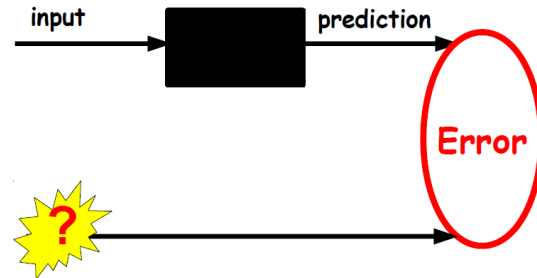  - Build generative mode $P(X)$ with a unsupervised manner from unlabeled data.

*"If you want to do computer vision, first learn computer graphics." -- Geoffrey Hinton*
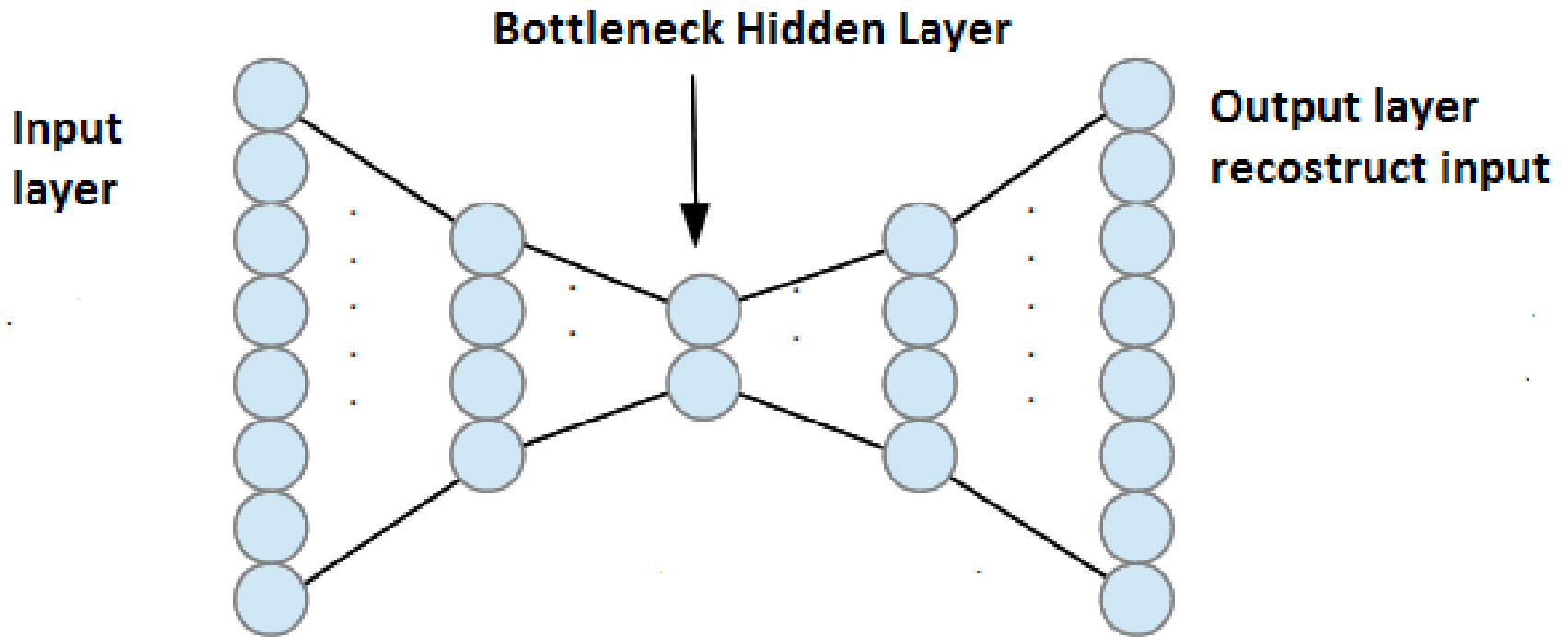
# Unsupervised learning

- Supervised learning

input → [ ] → prediction

target →

**Error**

- Unsupervised learning

input → [ ] → prediction

? →

**Error**

# Unsupervised learning

- How to constrain the model to represent training samples better than other data points?
  - [Restricted Boltzmann Machine]
    - Make models that defines the distribution of samples
  - [Auto-encoder with bottleneck]
    - reconstruct the input from the code & make code compact
  - [Sparse auto-encoder]
    - reconstruct the input from the code & make code sparse
  - [Denoising auto-encoder]
    - Add noise to the input or code

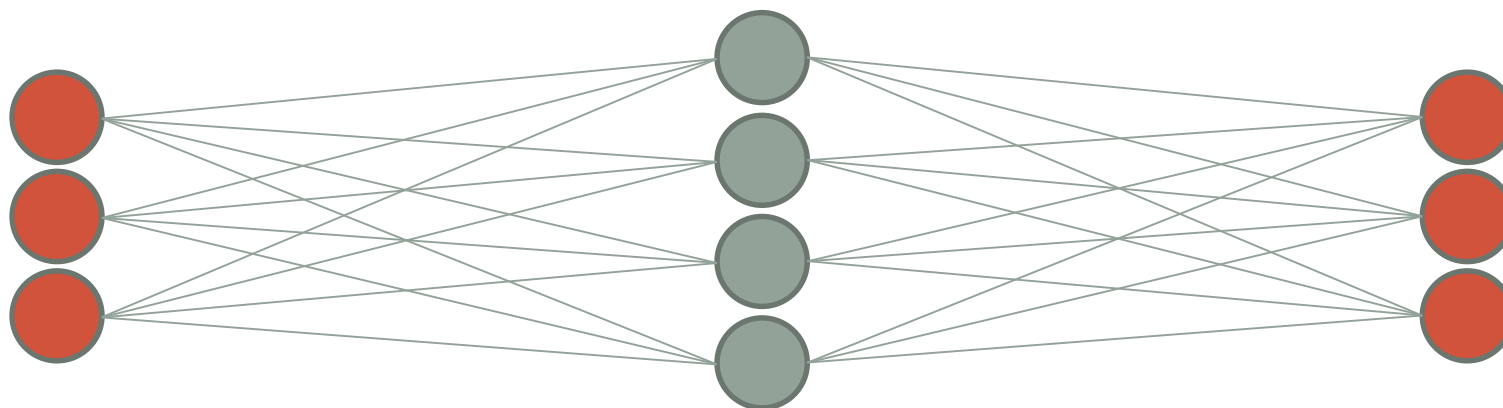# Auto-encoder with bottleneck

# Sparse auto-encoder

Input layer                    Hidden layer                    Output layer
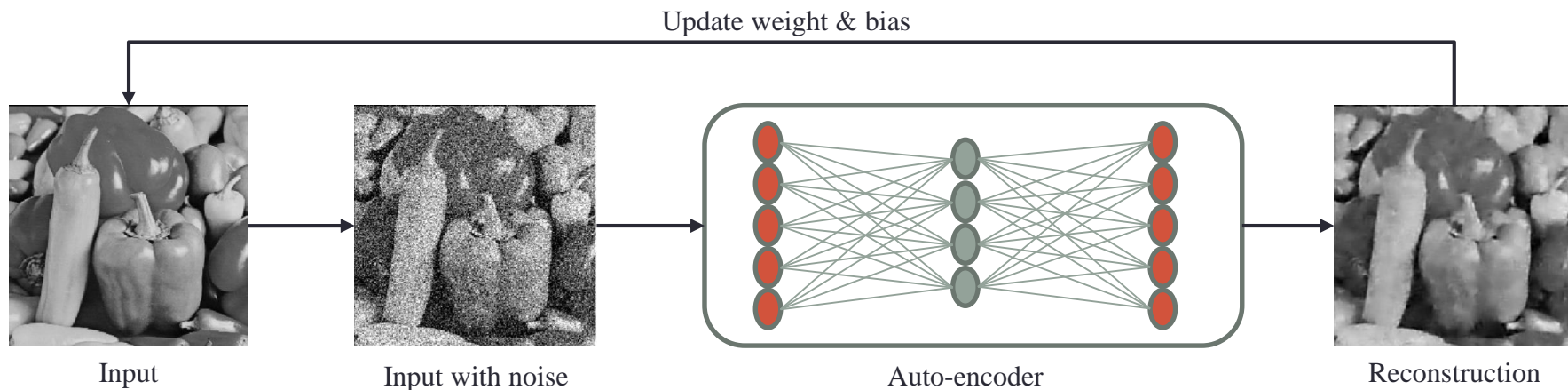


Input : MNIST Dataset
Input dimension : 784(28x28)
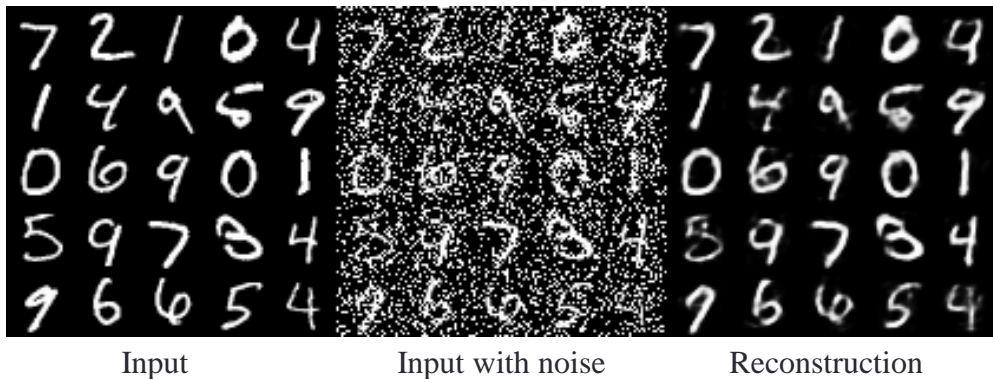
Reconstructions
Output dimension : 784(28x28)

# Denoising auto-encoder
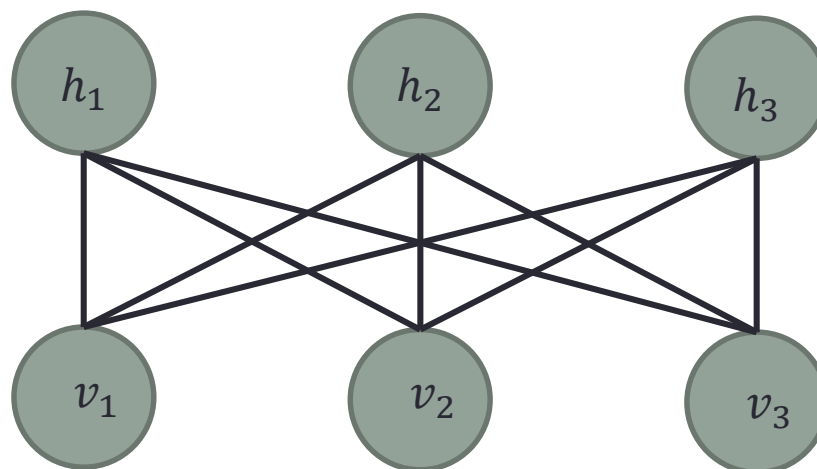
Training the denoising auto-encoder



Update weight & bias

Input     Input with noise     Auto-encoder     Reconstruction

Result



Input     Input with noise     Reconstruction

# RESTRICTED BOLTZMANN MACHINE

# Restricted Boltzmann Machine (RBM)

- RBM is a simple energy-based model: $p(v, h) = \frac{1}{Z}\exp(-E_\theta(v, h))$
  - With only $h - v$ interactions: $E_\theta(v, h) = -v^T W h - b^T v - d^T h$
  - here, we assume $h_j$ and $v_i$ are binary variables.
  - Normalizer $Z = \sum_{v,h}\exp(-E_\theta(v, h))$ is called a partition function

# Restricted Boltzmann Machine (RBM)

- The probability that the network assigns to a visible vector v:
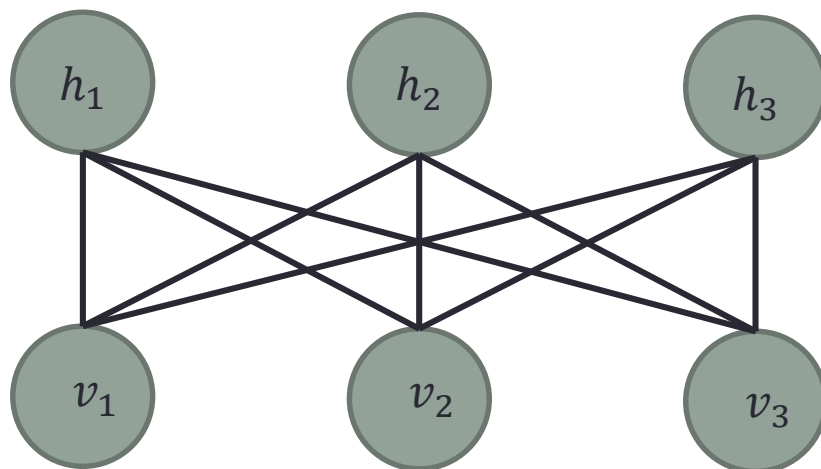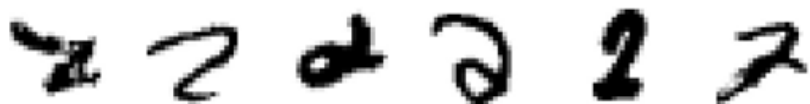  - $p(\mathrm{v}) = \frac{1}{Z}\sum_h \exp(-E_\theta(\mathrm{v}, \mathrm{h}))$
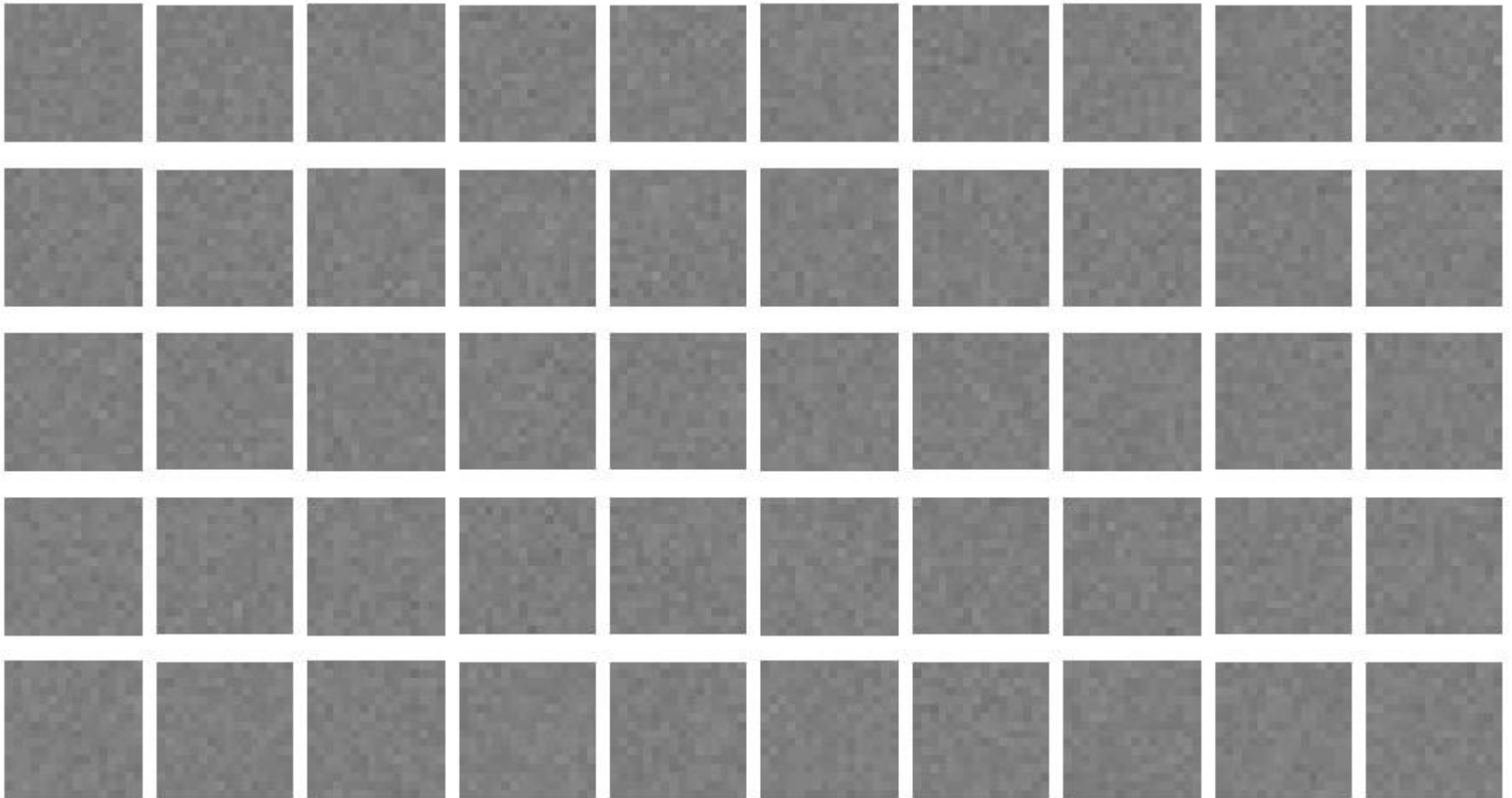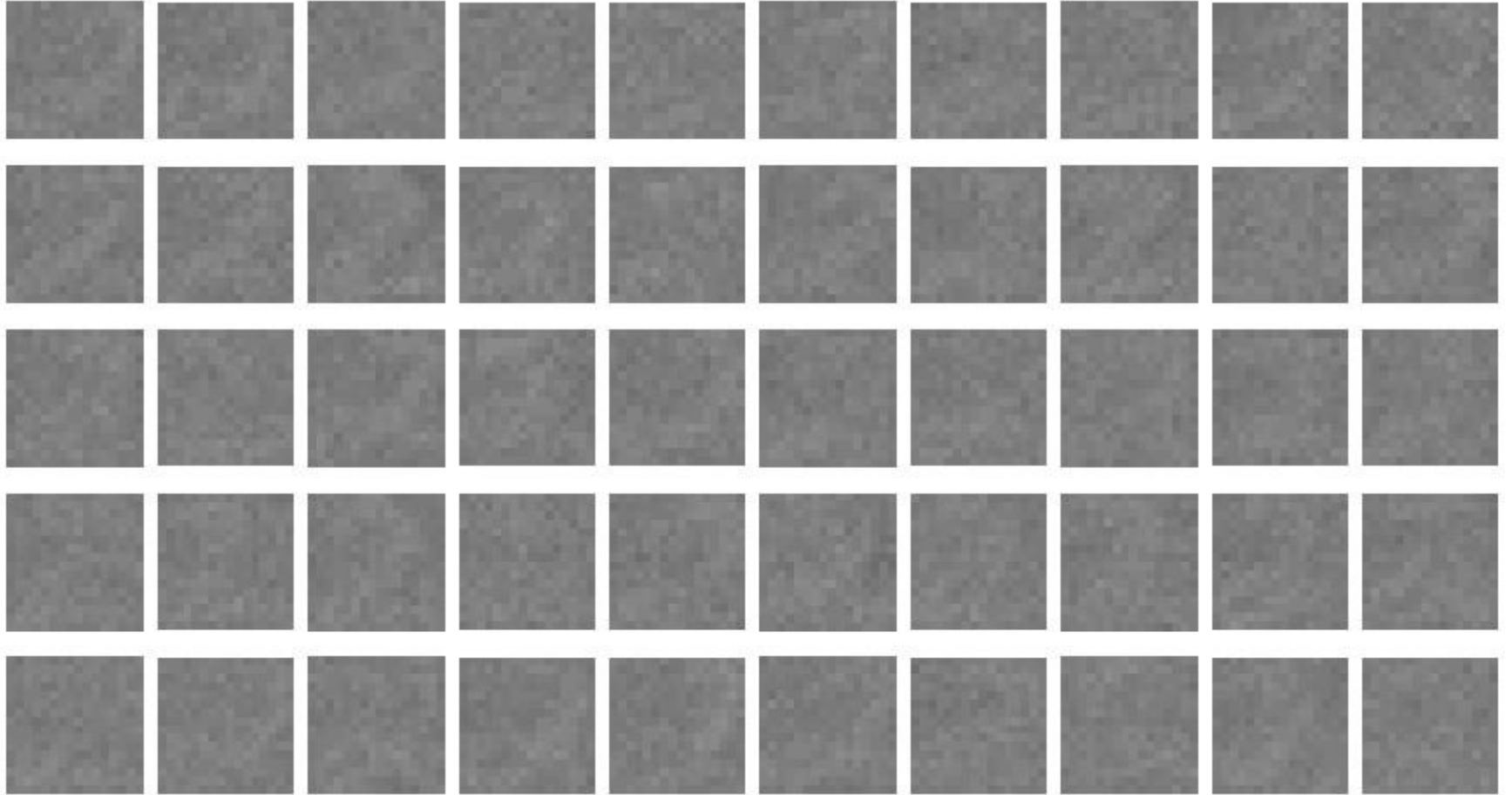
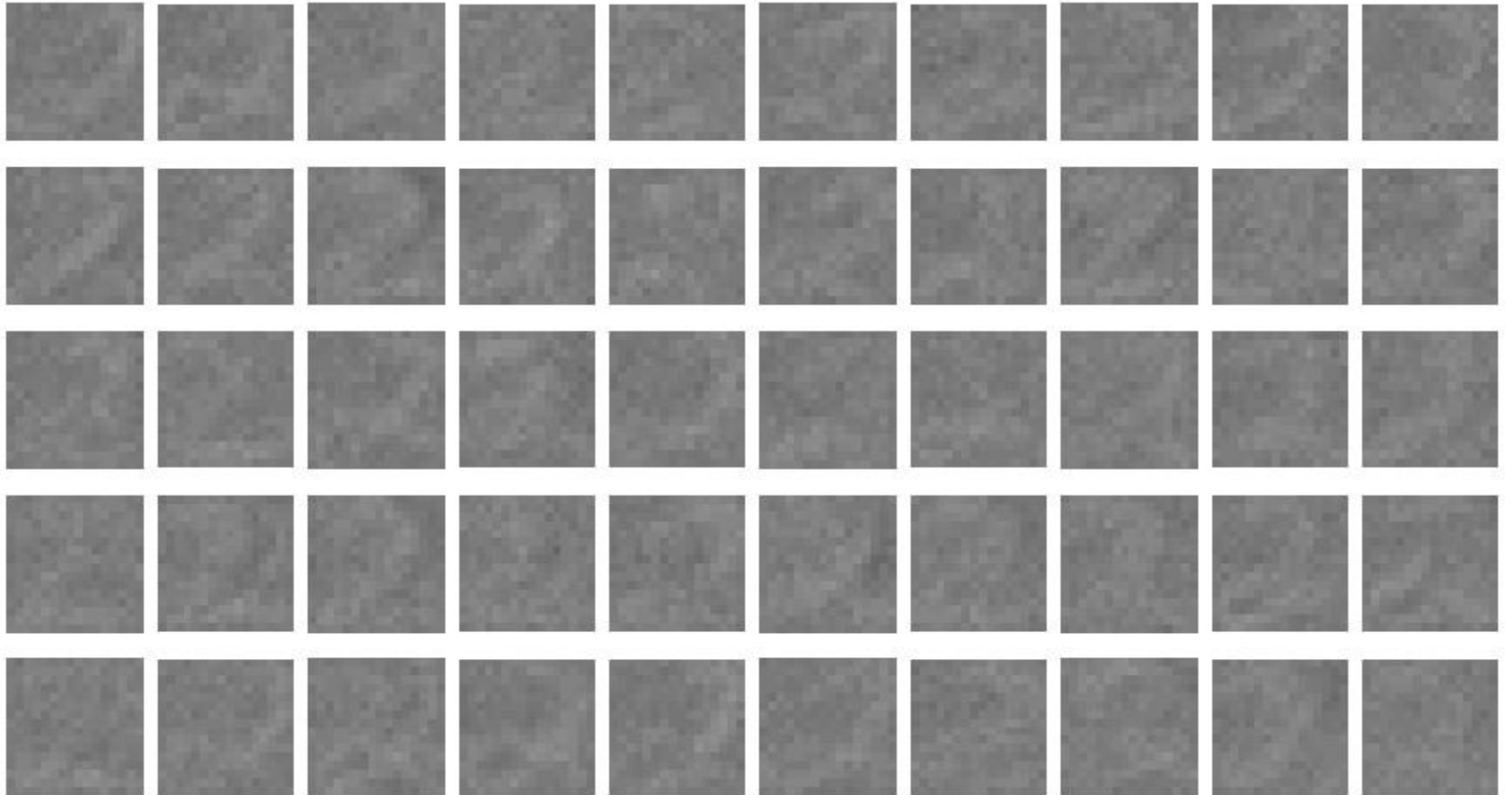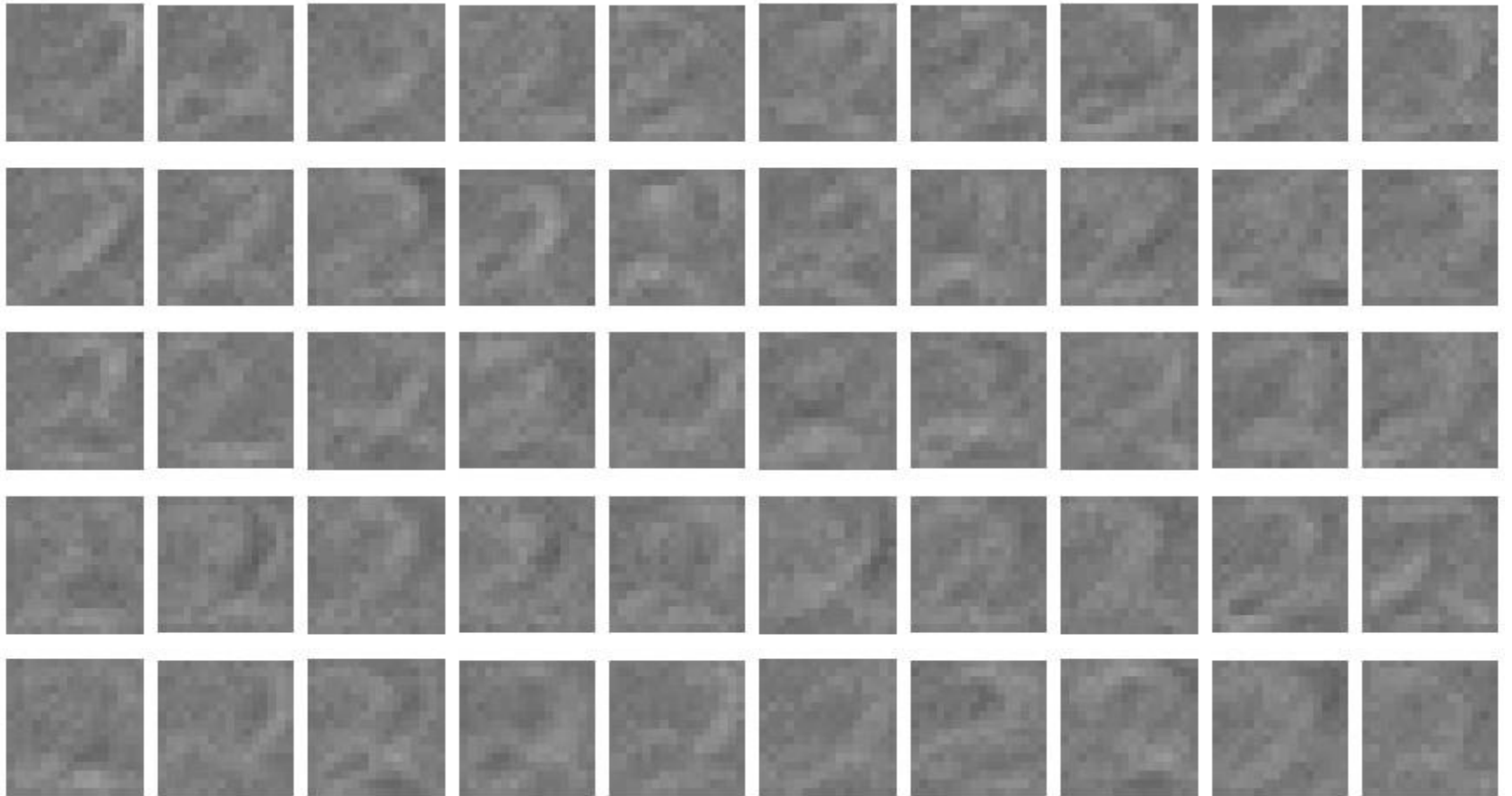h | 50 binary feature neurons |

v | 16 x 16 pixel image |
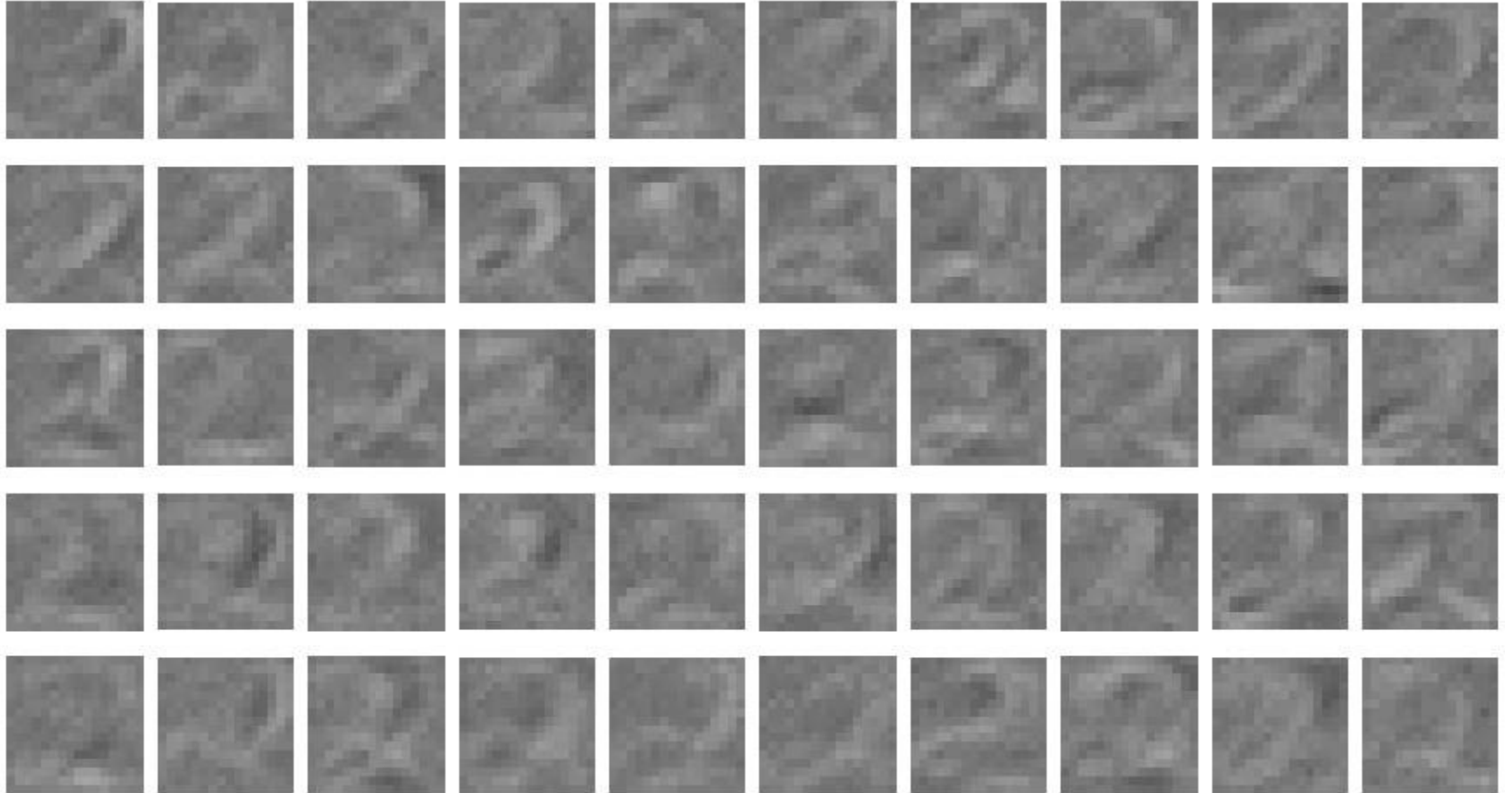
# The weights of the 50 feature detectors



We start with small random weights to break symmetry

# The final 50 x 256 weights



Each neuron grabs a different feature.

# Efficient learning procedure for RBMs

- The probability that the network assigns to a visible vector v:
  - $p(v) = \frac{1}{Z} \sum_h \exp(-E_\theta(v, h))$
  - $\hat{w} = \arg\max_w \sum_n \log p(v^n)$

- The derivative of the log probability w. r. t. a weight:
  - $\frac{1}{N} \sum_{n=1}^N \frac{\partial \log p(v^n)}{\partial w_{ij}} = <v_i h_j>_{data} - <v_i h_j>_{model}$

- Stochastic steepest ascent
  - $\Delta w_{ij} = \epsilon(<v_i h_j>_{data} - <v_i h_j>_{model})$

# Efficient learning procedure for RBMs

- Computing posteriors
  - $p(\mathrm{h}|\mathrm{v}) = \prod_{j=1} p(\mathrm{h}_j|\mathrm{v})$
  - $p(\mathrm{v}|\mathrm{h}) = \prod_{i=1} p(\mathrm{v}_i|h)$

- The absence of direct connections between hidden units in an RBM allows us to have
  - $p(h_j = 1|\mathrm{v}) = \sigma(d_j + \sum w_{ij} v_i)$

- The absence of direct connects between visible units in an RBM allows us to have
  - $p(v_i = 1|\mathrm{h}) = \sigma(b_i + \sum w_{ij} h_j)$

# Efficient learning procedure for RBMs



$$\Delta w_{ij} = \epsilon(<v_i h_j>_{data} - <v_i h_j>_{model})$$
$$= \epsilon(<v_i h_j>_{data} - <v_i h_j>^{\infty})$$

# Efficient learning procedure for RBMs



$$\Delta w_{ij} = \epsilon(< v_i h_j >_{data} - < v_i h_j >_{model})$$
$$= \epsilon(< v_i h_j >_{data} - < v_i h_j >^{\infty})$$
$$\simeq \epsilon(< v_i h_j >_{data} - < v_i h_j >^{1})$$

# Training RBMs with Contrastive Divergence

- The term $< v_i h_j >_{model}$ is expensive because it requires sampling $(v, h)$ from the model

- Gibbs sampling (sample $v$ then $h$ iteratively) works, but waiting for convergence at each gradient step is slow.

- Contrastive Divergence is faster: initialize with training point and wait only a few (usually 1) sampling steps
  - Let $v$ be a training point.
  - Sample $\widehat{h_j} \in \{0,1\}$ from $p\big(h_j = 1 \big| v\big) = \sigma(d_j + \sum w_{ij} v_i)$
  - Sample $\breve{v}_i \in \{0,1\}$ from $p\big(v_i = 1 \big| \widehat{h}\big) = \sigma(b_i + \sum w_{ij} \widehat{h}_j)$
  - Sample $\breve{h}_j \in \{0,1\}$ from $p\big(h_j = 1 \big| \breve{v}\big) = \sigma\big(b_i + \sum w_{ij} \breve{v}_i\big)$
  - $w_{ij} \leftarrow w_{ij} + \epsilon(v_i \widehat{h}_j - \breve{v}_i \breve{h}_j)$

# DEEP BELIEF NETWORK

Ruslan Salakhutdinov and Geoffrey Hinton,
"Deep Boltzmann Machine," 2009

# Deep Belief Nets (DBN) = Stacked RBM



- DBN defines a probabilistic generative model
  - $p(x) = \sum_{h,h',h''} p(x|h)p(h|h') \, p(h', h'')$

- Stacked RBMs can also be used to initialize a Deep Neural Network (DNN).

# Generating Data from a Deep Generative Model

- After training on 20k images, the generative model of can generate random  mages (dimension=8976) that are amazingly realistic!



**Deep Boltzmann Machine**

# DEEP AUTOENCODER

Hinton and Salakhutdinov, "Reducing the Dimensionality of Data with Neural Networks," Science, 2006

# Auto-encoder

- An auto-encoder is trained to encode the input in some representation so that the input can be reconstructed from that representation. Hence the target output is the input itself

$$x \longrightarrow \boxed{\text{encoder}} \xrightarrow{\ \ h\ \ } \boxed{\text{decoder}} \longrightarrow r \simeq x$$

# Deep auto-encoder



- Pre-training consists of learning **a stack of RBMs**, each having only one layer of feature detectors. The learned feature activations of one RBM are used as the "data" for training the next RBM in the stack. After the pre-training, the RBMs are "unrolled" to create a deep auto-encoder, which is then fine-tuned using **back-propagation** of error derivatives.

# Deep auto-encoder

- Can be used to reduce the dimensionality of the data
  - Better reconstruction than PCA

# Deep auto-encoder (digit data)

- PCA

- Deep auto-encoder
  - 784-1000-500-250-2

# Deep auto-encoder (documents)

- Visualization of the data (dimension reduction to 2D)
  - 2000-500-250-125-2

# Summary

- Layer-wise pre-training is the innovation that rekindled interest in deep architectures.

- Pre-training focuses on optimizing likelihood on the data, not the target label. First model $p(v)$ to do better $p(y|v)$.

- Why RBM? $p(h|x)$ is tractable, so it's easy to stack.

- RBM training can be expensive.
  - Solution: contrastive divergence

- DBN formed by stacking RBMs is a probabilistic generative mode

# RBM + SUPERVISED LEARNING

# UNSUPERVISED PRE-TRAINING

- We will use a greedy, layer-wise procedure
  - train one layer at a time, from first to last, with unsupervised criterion
  - fix the parameters of previous hidden layers
  - previous layers viewed as feature extraction

# FINE-TUNING

- Once all layers are pre-trained
  - add output layer
  - train the whole network using supervised learning
- Supervised learning is performed as in a regular feed-forward network
  - forward propagation, backpropagation and update
- We call this last phase fine-tuning
  - all parameters are "tuned" for the supervised task at hand
  - representation is adjusted to be more discriminative

# ONE LEARNING ALGORITHM HYPOTHESIS? GRANDMOTHER CELL HYPOTHESIS?

THERE'S A THEORY that human intelligence stems from a single algorithm.

The idea arises from experiments suggesting that the portion of your brain dedicated to processing sound from your ears could also handle sight for your eyes. This is possible only while your brain is in the earliest stages of development, but it implies that the brain is — at its core — a general-purpose machine that can be tuned to specific tasks.

About seven years ago, Stanford computer science professor Andrew Ng stumbled across this theory, and it changed the course of his career, reigniting a passion for artificial intelligence, or AI. "For the first time in my life," Ng says, "it made me feel like it might be possible to make some progress on a small part of the AI dream within our lifetime."

# THE MAN BEHIND THE GOOGLE BRAIN: ANDREW NG AND THE QUEST FOR THE NEW AI



Stanford professor Andrew Ng, the man at the center of the Deep Learning movement.
*Photo: Ariel Zambelich/Wired*



Auditory cortex learns to see

Andrew Ng's laptop explains Deep Learning. *Photo: Ariel Zambelich/Wired*

**SCIENCE** | Artificial Intelligence — cats — computer science — Google X

# Google's Artificial Brain Learns to Find Cat Videos

BY WIRED UK  06.26.12  |  11:15 AM  |  PERMALINK

Share  20    Tweet  1    8+1  497    Share    Pin it



*By Liat Clark, Wired UK*

When computer scientists at Google's mysterious X lab built a neural network of 16,000 computer processors with one billion connections and let it browse YouTube, it did what many web users might do — it began to look for cats.

# BUILDING HIGH-LEVEL FEATURES USING LARGE SCALE UNSUPERVISED LEARNING

Quoc V. Le et al, ICML, 2012

# Motivating Question

- Is it possible to learn high-level features (e.g. face detectors) using only unlabeled images?

- "Grandmother Cell" Hypothesis
  - Grandmother cell: A neuron that lights up when you see or hear your grandmother
    - Lots of interesting (controversial) discussions in the neuroscience literature

# Architecture (≃sparse deep auto-encoder)



Input to another layer above
(image with 8 channels)

Number of output
channels = 8

LCN Size = 5

Pooling Size = 5

Number
of maps = 8

RF size = 18

Number of input
channels = 3

H

W

One layer

Image Size = 200

$$\min_{W_d, W_e} \sum_m ||W_d W_e x^{(m)} - x^{(m)}|| \quad (1)$$

$$+ \sum_{m,k} \sqrt{\epsilon + P_k(W_e x^{(m)})^2} \quad (2)$$

(1): auto-encoder
(2): pooling

Repeated 3 times to form Deep Architecture
$x^{(m)}$ = image of 200x200 pixels x3 channels

# Training

- Using a deep network of 1 billion parameters
  - 10 million images (sampled from Youtube)
  - 1000 machines (16,000 cores) x 3 week.

- Model parallelism
  - Distributing the local weights $W_d, W_e$ in different machines
  - Asynchronous SGD

# Face neuron



Top stimuli from the test set

Optimal stimulus by numerical optimization
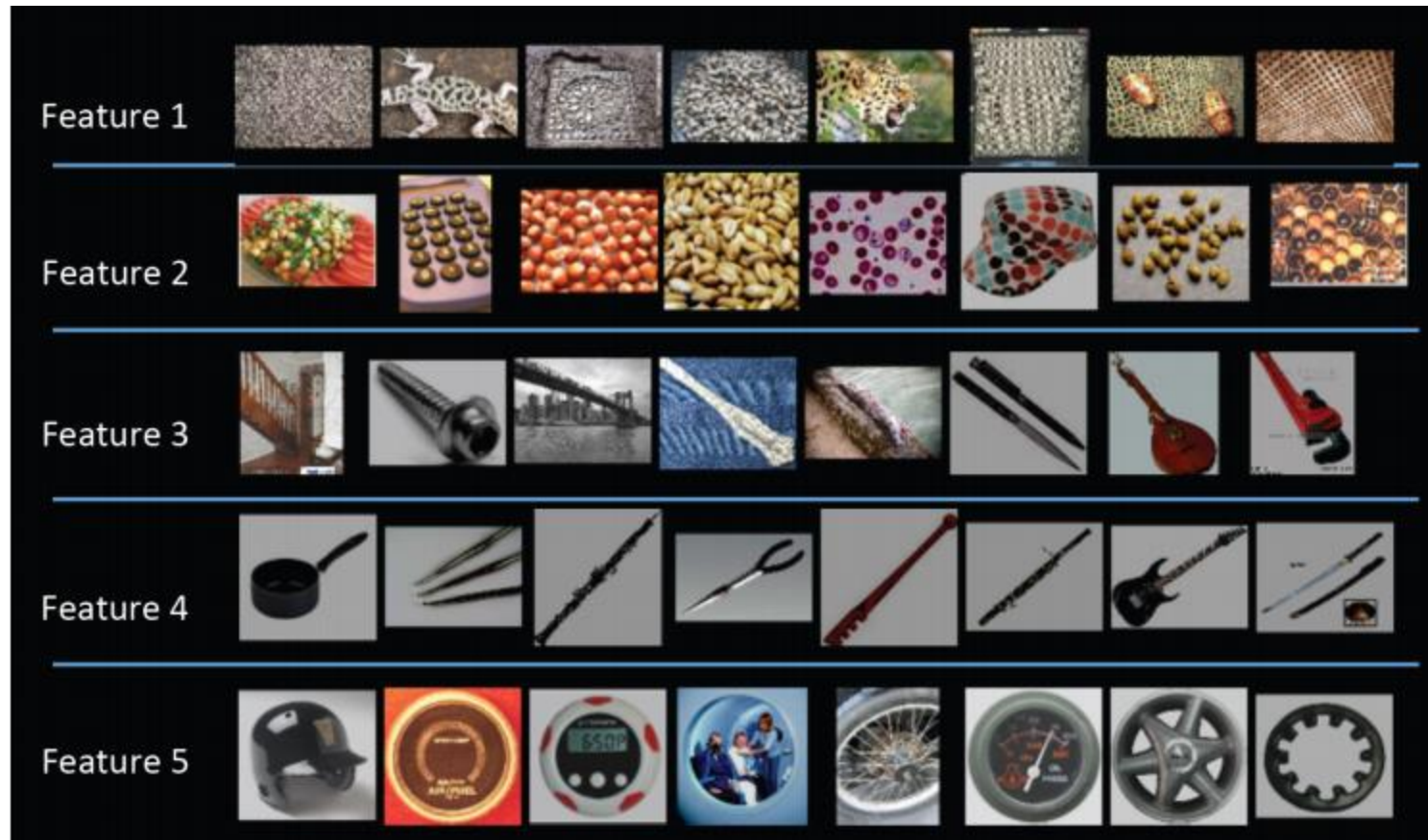
# Cat neuron



Top stimuli from the test set

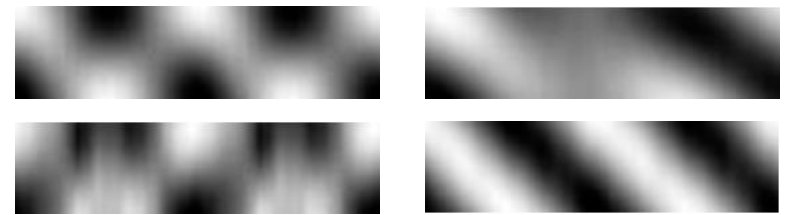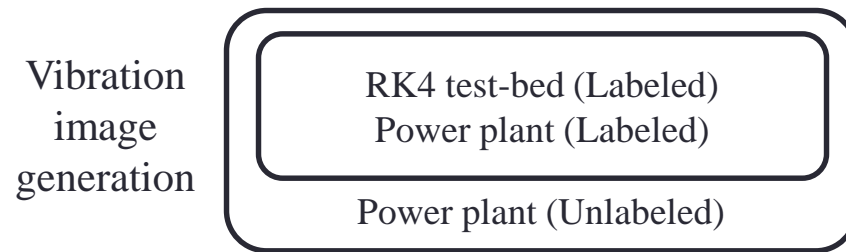Optimal stimulus
by numerical optimization

# More examples

# APPLICATION

| Methods | | | | | Tasks | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | ADAS | | | | | Smart factory |
| | | | | | Self Driving | | | Driver state | Vehicle Diagnosis | |
| | | | | | Localization | Perception | Planning/Control | | | |
| Methods | Traditional | Machine-Learning based method | Non-machine Learning | | GPS, SLAM | | Optimal control | | | |
| | Deep-Learning based | | Supervised | SVM MLP | | Pedestrian detection (HOG+SVM) | | | | |
| | | | | CNN | | Detection/Segmentation/Classification | End-to-end Learning | | | |
| | | | | RNN (LSTM) | | Dry/wet road classification | End-to-end Learning | Behavior Prediction/Driver identification | | * |
| | | | | DNN | | | | | * | * |
| | | | Reinforcement | | | | * | | | |
| | | | Unsupervised | | | | | | | * |

# Rotor System Diagnosis

**Vibration image generation**

RK4 test-bed (Labeled)
Power plant (Labeled)

Power plant (Unlabeled)

Vibration images using ODR

**High–level feature abstraction**

- Greedy layer-wise unsupervised pre-training :
Restricted Boltzmann Machine in deep belief network (DBN)

$v \quad h^1 \quad h^2 \quad h^{N-1} \quad h^N$

**Health reasoning**

- Classification :
multi-layer perceptron (MLP)
- Clustering : self-organizing map

Clustered_con1
Clustered_con2
Clustered_con3
Clustered_con4
Clustered_con5
Clustered_con6
Clustered_con7
Clustered_con8
…

Normal
Rubbing
Misalignment
Oil whirl
…

Known

U_Anomaly 1
U_Anomaly 2
U_Anomaly 3
…

Unknown