# Spiking Convolutional Deep Belief Networks for Unsupervised High Level Feature Extraction and Pattern Reconstruction

"Most of human and animal learning is unsupervised learning.
If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake.
We know how to make the icing and the cherry, but we don't know how to make the cake."

Reinforcement learning

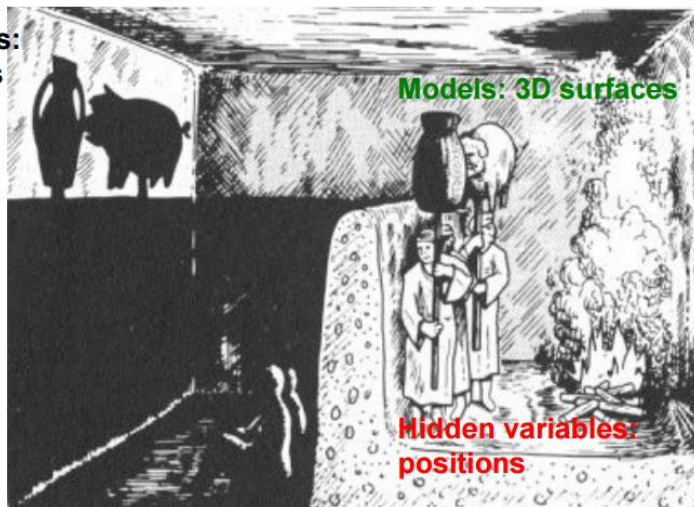Supervised learning

Unsupervised learning



Y Lecun, NIPS 2016

# "Learn the data"

"We believe that this can be achieved by learning a disentangled posterior distribution of the generative factors of the observed sensory input by leveraging the wealth of unsupervised data."[1]



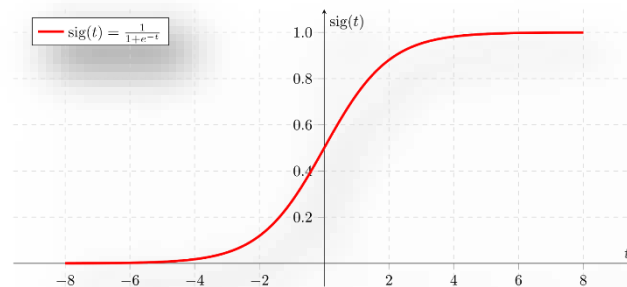[1] I. Higgins. Early Visual Concept Learning with Unsupervised Deep Learning. 2016
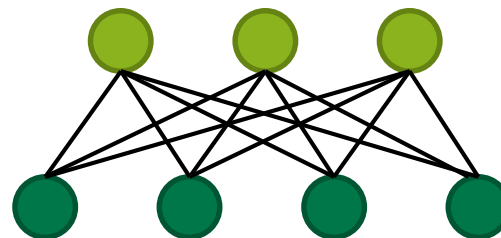
# Overview

- Restricted Boltzmann Machines (RBMs)
- Convolutional Deep Belief Networks

- Neural Sampling
- Event-Driven Contrastive Divergence
- Conversion

- Experiments & Results
- Conclusion

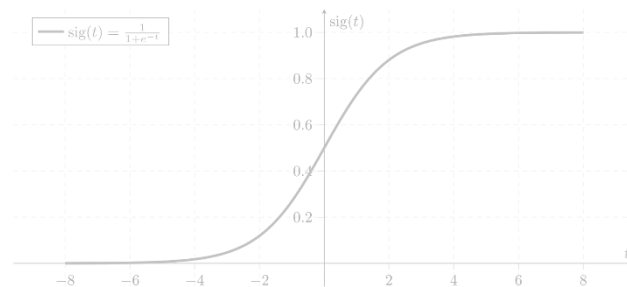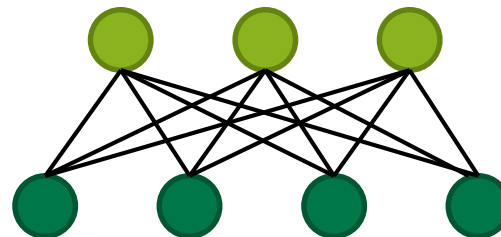# Restricted Boltzmann Machines

Neural Network:

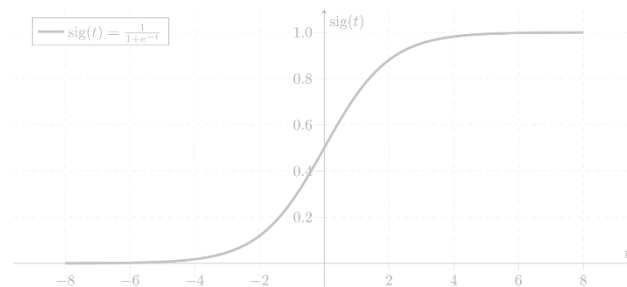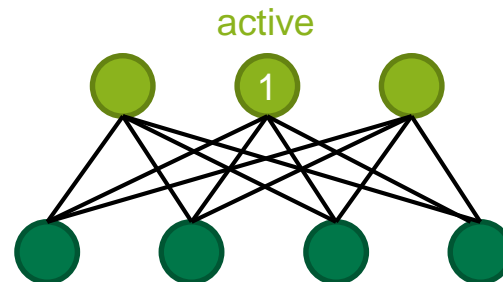- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p("active") = \sigma("input")$
- Energy-based Model

# Restricted Boltzmann Machines

Neural Network:


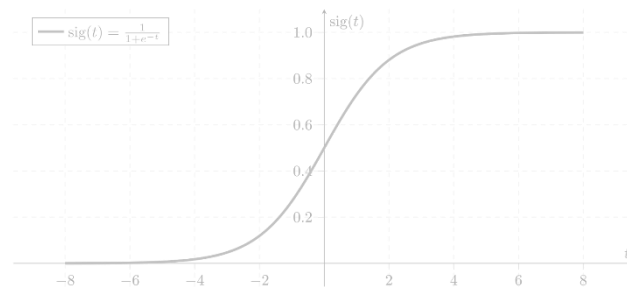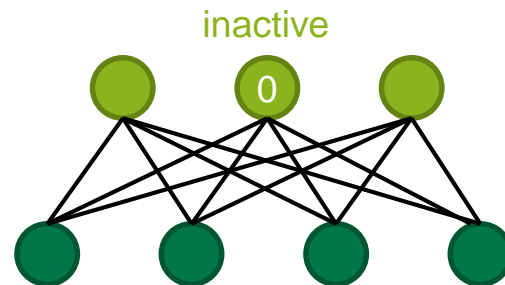
- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p("\,active\,") = \sigma("\,input\,")$
- Energy-based Model

# Restricted Boltzmann Machines

Neural Network:

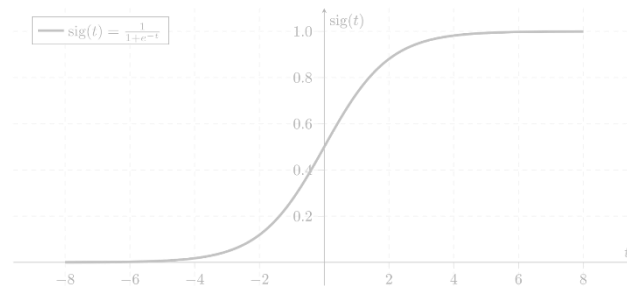- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p("\,active\,") = \sigma("\,input\,")$
- Energy-based Model

# Restricted Boltzmann Machines

Neural Network:

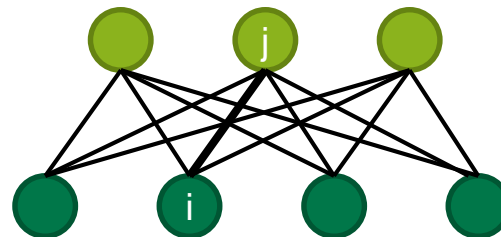- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p("\,active\,") = \sigma("\,input\,")$
- Energy-based Model



inactive



$\mathrm{sig}(t) = \frac{1}{1+e^{-t}}$

© FZI Forschungszentrum Informatik

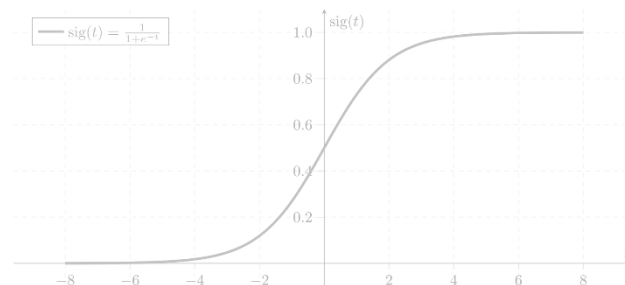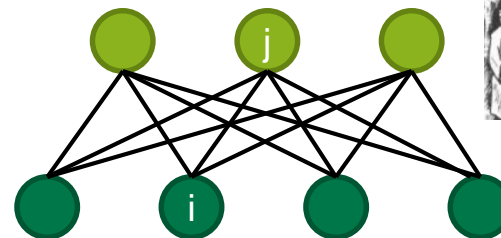# Restricted Boltzmann Machines

Neural Network:

- Binary Units
- **Symmetric Connections**
- Two bipartite Layers
- Stochastic Activations $p("\,active\,") = \sigma("\,input\,")$
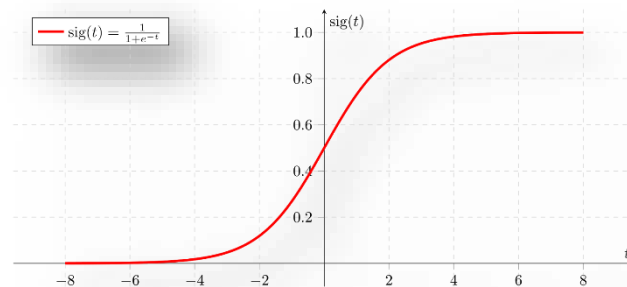- Energy-based Model





$\mathrm{sig}(t) = \frac{1}{1+e^{-t}}$

# Restricted Boltzmann Machines

Neural Network:

- Binary Units
- Symmetric Connections
- **Two bipartite Layers**
- Stochastic Activations $p("active") = \sigma("input")$
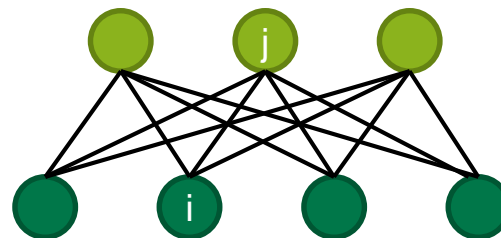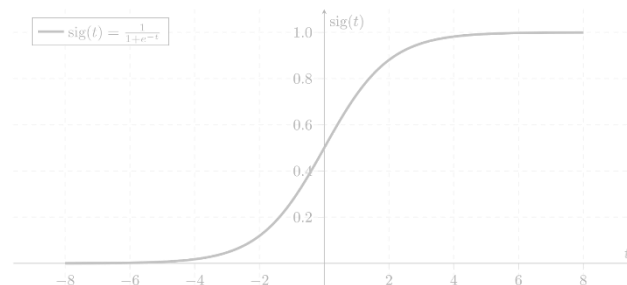- Energy-based Model

hidden

visible

# Restricted Boltzmann Machines

Neural Network:

- Binary Units
- Symmetric Connections
- Two bipartite Layers
- **Stochastic Activations** $p(" \, active \, ") = \sigma(" \, input \, ")$
- Energy-based Model



hidden

visible



$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

# Restricted Boltzmann Machines

Neural Network:

hidden

visible

- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p("active") = \sigma("input")$
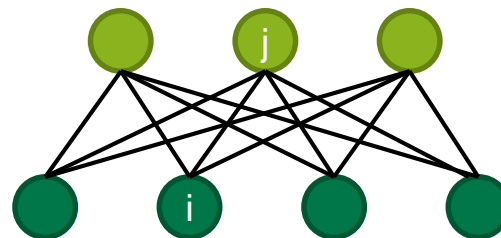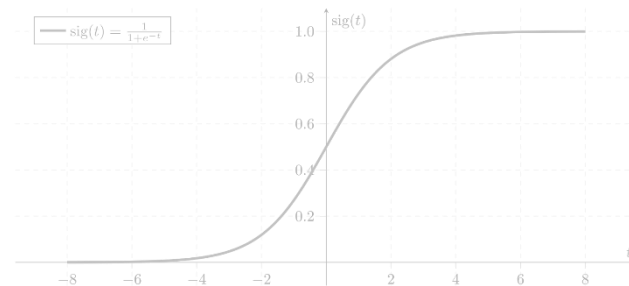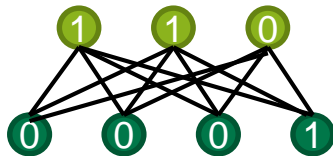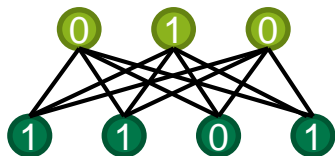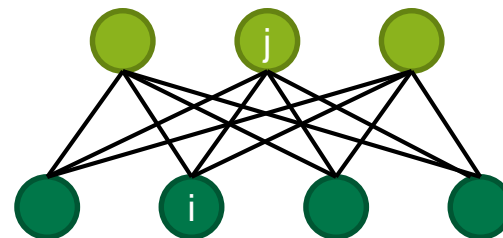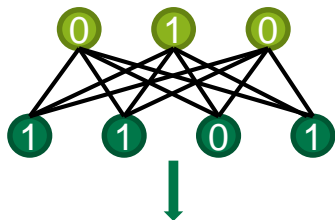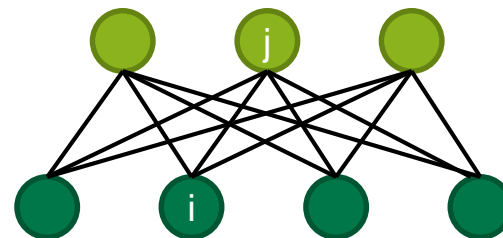- **Energy-based Model**

# Restricted Boltzmann Machines

Neural Network:

- Binary Units
- Symmetric Connections
- Two bipartite Layers
- Stochastic Activations $p(\text{"}active\text{"}) = \sigma(\text{"}input\text{"})$
- Energy-based Model



hidden

visible



$\text{sig}(t) = \frac{1}{1+e^{-t}}$

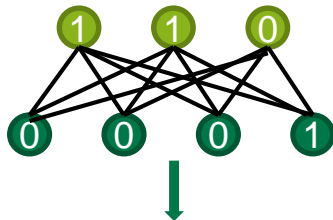# Restricted Boltzmann Machines

Neural Network:



hidden

visible

- Binary Units
- Symmetric Connections
- Two bipartite Layers
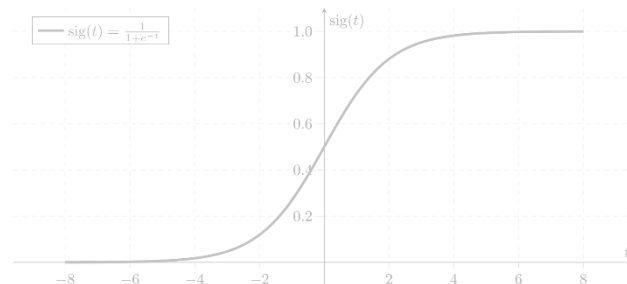- Stochastic Activations $p("\,active\,") = \sigma("\,input\,")$
- Energy-based Model



Energy = 5



Energy = 2

$$\text{sig}(t) = \frac{1}{1+e^{-t}}$$

# Training RBMs

→ Maximize probability of training samples ⇔ Minimize the Energy :

$$\frac{\partial E(v)}{\partial w} = \ ... = \left(v\,h^T\right)_{model} - \left(v\,h^T\right)_{data}$$

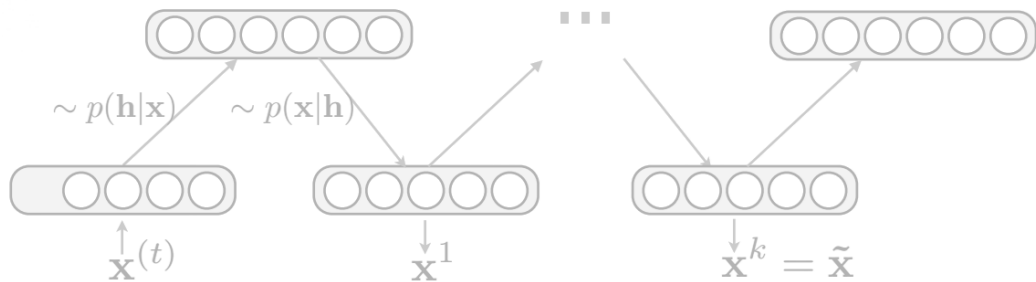$$\rightarrow w^{new} = w^{old} + \mu\left(\,(v\,h^T)_{data} - (v\,h^T)_{model}\right)$$

# Training RBMs

→ Maximize probability of training samples ⇔ Minimize the Energy :

$$\frac{\partial E(v)}{\partial w} = \ldots = \left(v\, h^T\right)_{\mathbf{model}} - \left(v\, h^T\right)_{\mathrm{data}}$$

$$\to w^{\mathrm{new}} = w^{\mathrm{old}} + \mu \left( \left(v\, h^T\right)_{data} - \left(v\, h^T\right)_{model} \right)$$

*

$\sim p(\mathbf{h}|\mathbf{x})$   $\sim p(\mathbf{x}|\mathbf{h})$

$\mathbf{x}^{(t)}$   $\mathbf{x}^1$   $\mathbf{x}^k = \tilde{\mathbf{x}}$

© FZI Forschungszentrum Informatik

# Training RBMs

→ Maximize probability of training samples ⇔ Minimize the Energy :

$$\frac{\partial E(v)}{\partial w} = \ldots = \left(v\,h^T\right)_{\text{model}} - \left(v\,h^T\right)_{\textbf{data}}$$

$$\rightarrow w^{\text{new}} = w^{\text{old}} + \mu\,(\,(v\,h^T)_{data} - (v\,h^T)_{model})$$
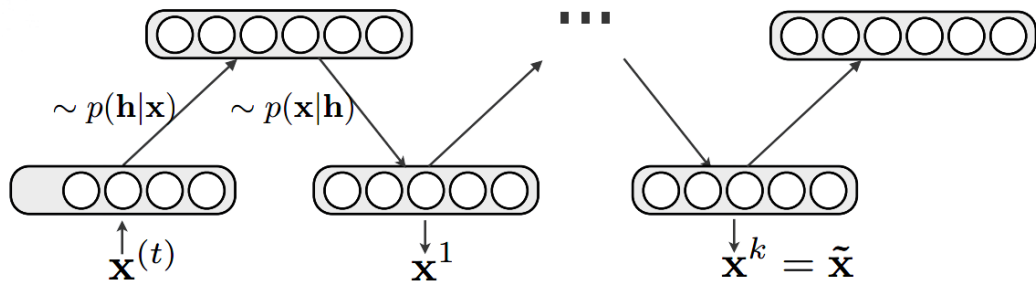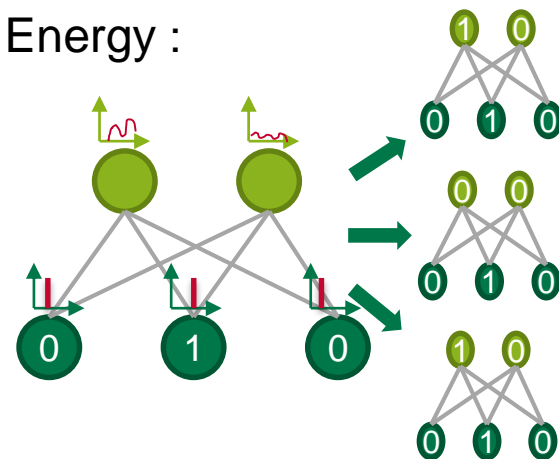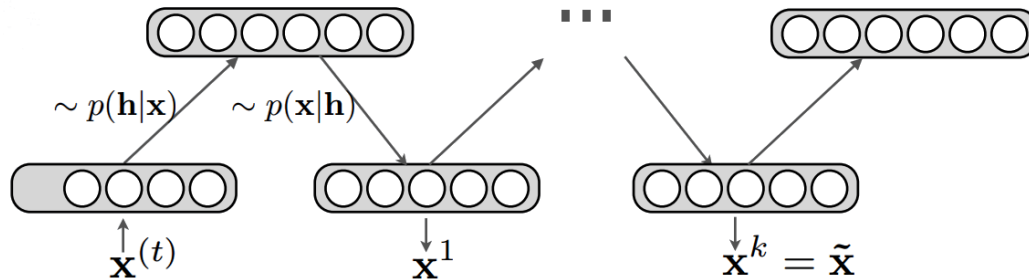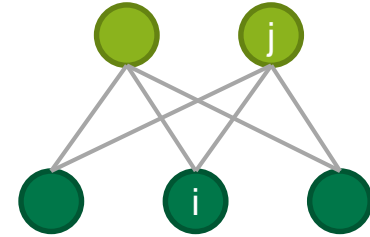
# Training RBMs

→ Maximize probability of training samples ⇔ Minimize the Energy :

$$\frac{\partial E(v)}{\partial w} = \ldots = \left( v\, h^T \right)_{\text{model}} - \left( v\, h^T \right)_{\textbf{data}}$$

$$\rightarrow w^{\text{new}} = w^{\text{old}} + \mu \left( (v\, h^T)_{data} - (v\, h^T)_{model} \right)$$





$\sim p(\mathbf{h}|\mathbf{x})$    $\sim p(\mathbf{x}|\mathbf{h})$

$\mathbf{x}^{(t)}$          $\mathbf{x}^1$          $\mathbf{x}^k = \tilde{\mathbf{x}}$

# Training RBMs

→ Maximize probability of training samples ⇔ Minimize the Energy :

$$\frac{\partial E(v)}{\partial w} = \ldots = \left(v\, h^{\mathrm{T}}\right)_{\mathrm{model}} - \left(v\, h^{\mathrm{T}}\right)_{\mathrm{data}}$$

$$\to w^{\mathrm{new}} = w^{\mathrm{old}} + \mu\left(\left(v\, h^{T}\right)_{data} - \left(v\, h^{T}\right)_{model}\right)$$

# Training RBMs

Contrastive Divergence: $\mathrm{w}^{\mathrm{new}} = \mathrm{w}^{\mathrm{old}} + \mu \left( \left( v\, h^T \right)_{data} - \left( v\, h^T \right)_{model} \right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu \left( (v\,h^T)_{data} - (v\,h^T)_{model} \right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu\left((v\,h^T)_{data} - (v\,h^T)_{model}\right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu \left( (v\,h^T)_{data} - (v\,h^T)_{model} \right)$

© FZI Forschungszentrum Informatik

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu\left((v\,h^T)_{data} - (v\,h^T)_{model}\right)$

# Training RBMs

Contrastive Divergence: $\mathrm{w}^{\mathrm{new}} = \mathrm{w}^{\mathrm{old}} + \mu \left( \left( v\, h^T \right)_{data} - \left( v\, h^T \right)_{model} \right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu\left(\left(v\,h^T\right)_{data} - \left(v\,h^T\right)_{model}\right)$

© FZI Forschungszentrum Informatik

# Training RBMs

Contrastive Divergence: $\text{w}^{\text{new}} = \text{w}^{\text{old}} + \mu \left( (v\,h^T)_{data} - (v\,h^T)_{model} \right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu\left(\left(v\,h^T\right)_{data} - \left(v\,h^T\right)_{model}\right)$

# Training RBMs

Contrastive Divergence: $\mathrm{w}^{\mathrm{new}} = \mathrm{w}^{\mathrm{old}} + \mu\left((v\,h^T)_{data} - (v\,h^T)_{model}\right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu\,(\,(v\,h^T)_{data} - (v\,h^T)_{model}\,)$

# Training RBMs

Contrastive Divergence: $\mathrm{w}^{\mathrm{new}} = \mathrm{w}^{\mathrm{old}} + \mu \left( (v\,h^T)_{data} - (v\,h^T)_{model} \right)$

# Training RBMs

Contrastive Divergence: $w^{new} = w^{old} + \mu \left( (v\,h^T)_{data} - (v\,h^T)_{model} \right)$
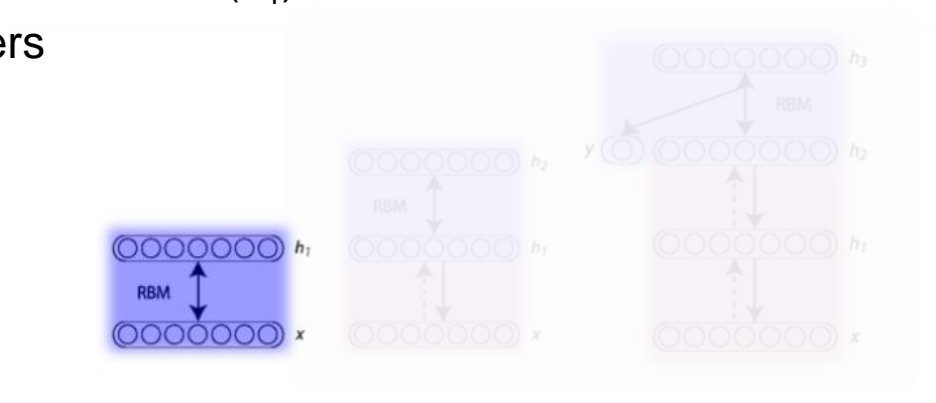
# "Diving in Deeper": Deep Belief Networks

- Adding Hidden Layers increases Representational Power
- DBNs can approximate any distribution over binary vectors



Low-level features: Edges

Input: Pixels

Image

# "Diving in Deeper": Deep Belief Networks

- Adding Hidden Layers increases Representational Power
- DBNs can approximate any distribution over binary vectors



Higher-level features:
Combination of edges

Low-level features:
Edges

Input: Pixels

Image

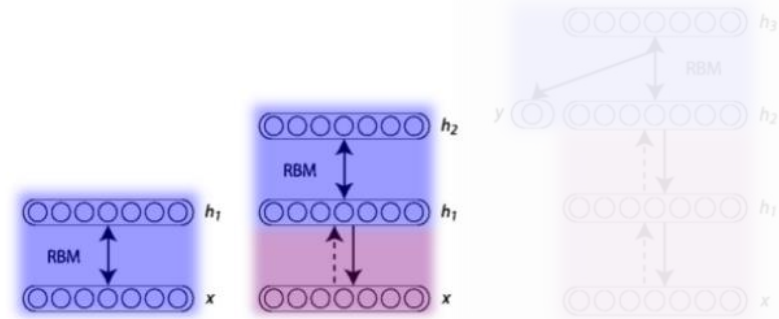# "Diving in Deeper": Deep Belief Networks

Built up by stacking up RBMs:

- Fit parameters $W_1$ of the 1st layer RBM to data (x)
- Freeze $W_1$ and use samples $h_1$ as data for the next layer
- Fit parameters $W_2$ of the 2nd layer RBM to data ($h_1$)
- Proceed recursively for the next layers

# "Diving in Deeper": Deep Belief Networks
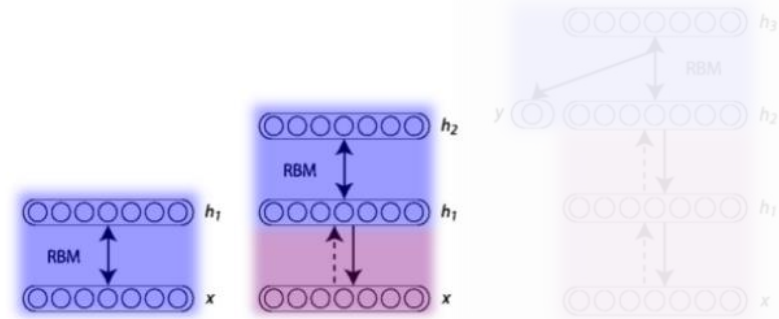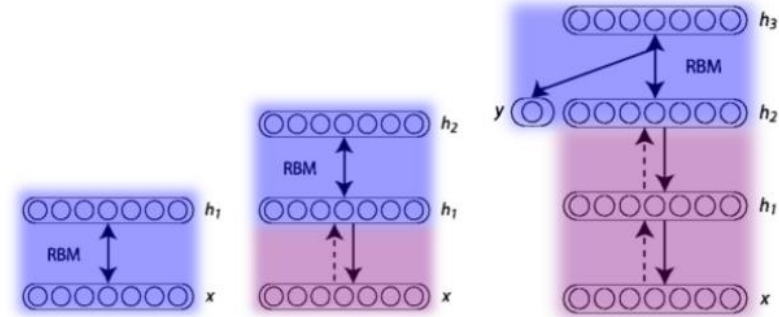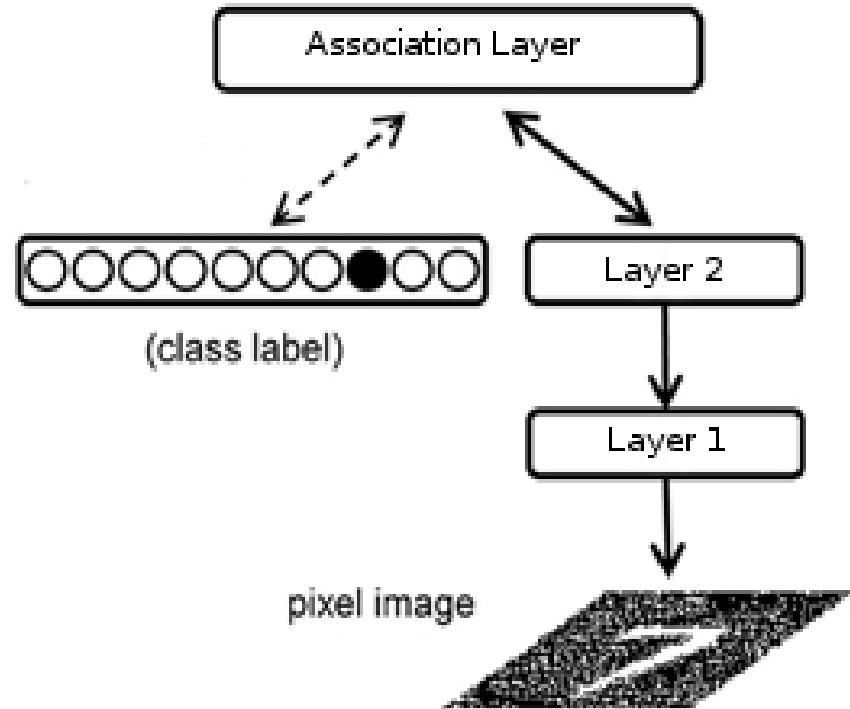
Built up by stacking up RBMs:

- Fit parameters $W_1$ of the 1st layer RBM to data (x)
- Freeze $W_1$ and use samples $h_1$ as data for the next layer
- Fit parameters $W_2$ of the 2nd layer RBM to data ($h_1$)
- Proceed recursively for the next layers

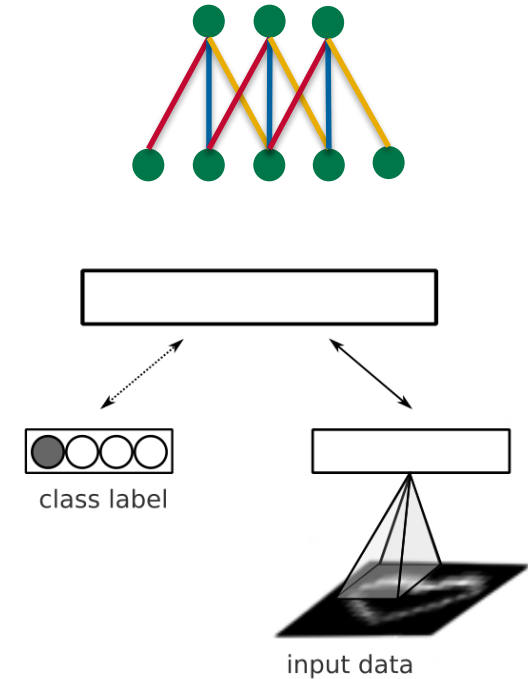# "Diving in Deeper": Deep Belief Networks

Built up by stacking up RBMs:

- Fit parameters $W_1$ of the 1st layer RBM to data (x)
- **Freeze $W_1$ and use samples $h_1$ as data for the next layer**
- Fit parameters $W_2$ of the 2nd layer RBM to data ($h_1$)
- Proceed recursively for the next layers

# "Diving in Deeper": Deep Belief Networks

Built up by stacking up RBMs:

- Fit parameters $W_1$ of the 1st layer RBM to data (x)
- Freeze $W_1$ and use samples $h_1$ as data for the next layer
- **Fit parameters $W_2$ of the 2nd layer RBM to data ($h_1$)**
- Proceed recursively for the next layers

# "Diving in Deeper": Deep Belief Networks

Built up by stacking up RBMs:

- Fit parameters $W_1$ of the 1st layer RBM to data (x)
- Freeze $W_1$ and use samples $h_1$ as data for the next layer
- Fit parameters $W_2$ of the 2nd layer RBM to data ($h_1$)
- Proceed recursively for the next layers

# Classification with DBNs

- Top Layer : Joint density for labels and images

- Bottom Layers : Feature Extraction

© FZI Forschungszentrum Informatik

# Convolutional DBNs

Using a convolutional architecture:

- Partially connected
- Shared weights

class label

input data

# Sampling in Spiking Neural Networks

- State of Neuron defined by its Firing

- State of Network defined by Neurons

# Sampling in Spiking Neural Networks

- State of Neuron defined by its Firing
  - Firing probability: $p(\mathrm{x}{=}1) \propto \sigma(Wz)$

- State of Network defined by Neurons

# Event-driven Contrastive Divergence

Contrastive Divergence:



hidden

visible

- For Spiking Neural Networks
- With Spike-time Dependent (Synaptic) Plasticity (STDP)

Data burn-in



Data distribution
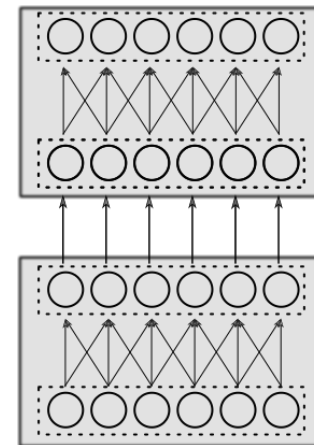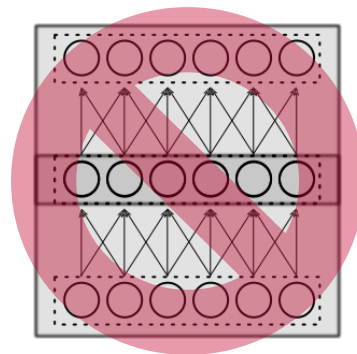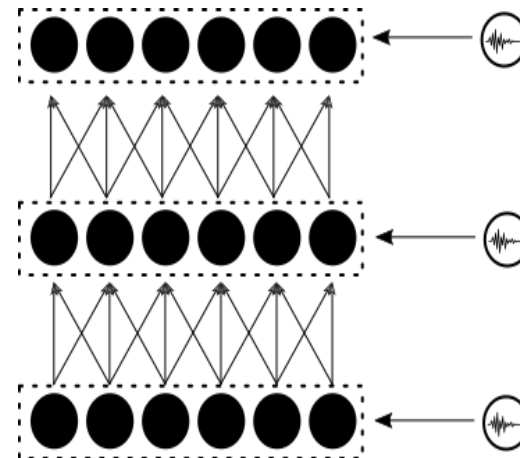


Model burn-in



Model distribution

# Spiking Convolutional DBN

- Layerwise Training with eCD
- Weight sharing with weight synchronization
- Inhibitory lateral Connections
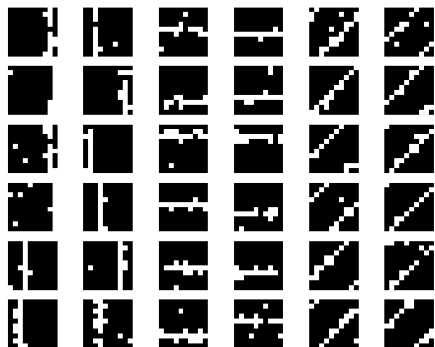- Forward Connections between RBMs

# Spiking Convolutional DBN

- **Layerwise Training with eCD**
- Weight sharing with weight synchronization
- Inhibitory lateral Connections
- Forward Connections between RBMs

# Spiking Convolutional DBN

- Layerwise Training with eCD
- **Weight sharing with weight synchronization**
- Inhibitory lateral Connections
- Forward Connections between RBMs



Training step

Weight synchronization

# Spiking Convolutional DBN

- Layerwise Training with eCD
- Weight sharing with weight synchronization
- **Inhibitory lateral Connections**
- Forward Connections between RBMs

# Spiking Convolutional DBN

- Layerwise Training with eCD
- Weight sharing with weight synchronization
- Inhibitory lateral Connections
- **Forward Connections between RBMs**

# Conversion

- Train a DBN
- Replace Binary Neurons with Spiking Neurons
- Use Synaptic Connections
- Scale Synaptic Weights
- Add external Poisson-Noise

# Experiments

## Datasets



| | | |
|:---:|:---:|:---:|
| **Stripes** | **Poker**<br>**(Event-based)** | **Ball-Can-Pen**<br>**(Event-based)** |
| 10 x 10 Pixel | 16 x 16 Pixel | 16 x 16 Pixel |

# eCD - Results

Accuracy:
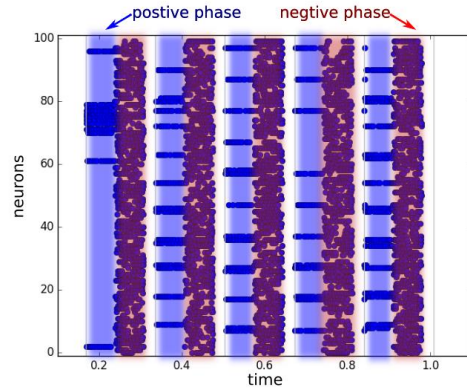
- Stripes: 1.00
- Poker: 0.94
- Ball-Can-Pen: 0.90



class label

input data

# eCD - Weights

© FZI Forschungszentrum Informatik

# eCD - Spikes

visible

hidden

Training

© FZI Forschungszentrum Informatik

# eCD - Reconstruction

Reconstruction



Bottom half missing

Right half missing

© FZI Forschungszentrum Informatik

# eCD - Reconstruction

Reconstruction



Bottom half missing

Right half missing

# Comparison

Artificial DBN and Spiking DBN trained with 100 Samples

Weights



Converted DBN          Event-based DBN

# Discussion

→ Train a spiking convolutional DBN

+ Unsupervised Learning

+ Spiking Neural Network

+ Event-based

+ Biological Plausibility

- Computational Resources

# Thanks !

*"We have truly autonomous cars when you tell it to drive to the office, and it decides to drive to the beach."*

*"Geoff Hinton doesn't disagree with you, he contrastively diverges."*

*"Geoff Hinton discovered how the brain really works. Once a year for the last 25 years."*

# References

Lecun, Yann; Predictive Learning; NIPS 2016 , https://drive.google.com/file/d/0BxKBnD5y2M8NREZod0tVdW5FLTQ/view

Kokkinos, Iasonas;  Introduction to Deep Learning;  http://cvn.ecp.fr/personnel/iasonas/course/DL5.pdf

Larochelle, Hugo; Neural networks: RBM - CD;  http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Lamblin, Pascal; Deep belief networks;  http://www.iro.umontreal.ca/~lisa/twiki/bin/view.cgi/Public/DeepBeliefNetworks

M. Zorzi, A. Testolin, and I. Stoianov. Modeling language and cognition with deep unsupervised learning: a tutorial overview.

L. Buesing, J. Bill, B. Nessler, and W. Maass. Neural dynamics as sampling: A model for stochastic computation in recurrent networks of spiking neurons.

M. A. Petrovici. Form Versus Function: Theory and Models for Neuronal Substrates. 2016.

K. A. Buchanan and J. R. Mellor. The activity requirements for spike timing-dependent plasticity in the hippocampus.

I. Higgins. Early Visual Concept Learning with Unsupervised Deep Learning. 2016

* Actually Bernoulli

# Appendix

$$C_m \frac{\partial u}{\partial t} = g_l (E_l - u(t)) + I^{syn} + I^{ext}$$

| | |
|---|---|
| Resting potential | -65 mV |
| Membrane capacity | 1.0 nF |
| Membrane time constant | 20.0 ms |
| Refractory period | 10.0 ms |
| Offset current | 1.0 nA |
| Reset potential | -53.0 mV |
| Spike threshold | -52.0 mV |
| Inhibitory reversal potential | 90.0 mV |
| Excitatory reversal potential | -0.0 mV |

# Appendix

$$STDP(v_i(t), h_j(t)) = v_i(t)A_{h_j}(t) + h_j(t)A_{v_i}(t),$$

$$A_{h_j}(t) = A \int_{-\infty}^{t} W(t-s)h_j(s)ds,$$

$$A_{v_i}(t) = A \int_{-\infty}^{t} W(t-s)v_i(s)ds.$$

$$W(x) = exp(\tfrac{x}{\tau}).$$

$$A_v = A_v \exp(\frac{-\Delta t}{\tau}) + a_\delta,$$

$$A_h = A_h \exp(\frac{-\Delta t}{\tau}),$$

$$\delta w = \mu g(t)A_v,$$

# Appendix

# Appendix

© FZI Forschungszentrum Informatik

# Appendix

| Simulated time | Spiking CNN | | CUBA LIF DBN | |
| --- | --- | --- | --- | --- |
| | Classification Accuracy | Runtime | Classification Accuracy | Runtime |
| 50 ms | 0.69 | 7.8 s | 0.81 | 8.1 s |
| 100 ms | 0.77 | 9.2 s | 0.89 | 10.5 s |
| 200 ms | 0.76 | 13.1 s | 0.89 | 14.6 s |
| 300 ms | 0.75 | 15.1 s | 0.91 | 18.5 s |
| 500 ms | 0.83 | 24.2 s | 0.93 | 30.6 s |

# Appendix

# Appendix

| | |
|---|---|
| $t_{burn-in}$ | 14 ms |
| $t_{learn}$ | 56 ms |
| $t_{flush}$ | 28 ms |
| Learn-rate | 1.0 |
| Weight-decay | 0.001 |
| Weight synchronization after $n$ samples | 1 |

# Appendix

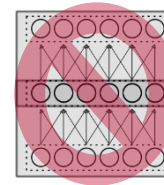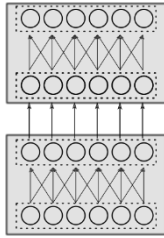| | Stripes | Poker | Ball-Pen-Can |
|---|---|---|---|
| Input | 100 | 256 | 256 |
| #1. Layer Params | 20 x 7 x 7 = 980 | 10 x 14 x 14 = 1960 | 20 x 14 x 14 = 3920 |
| 1. Layer | 20 x 4 x 4 = 320 | 10 x 3 x 3 = 90 | 20 x 3 x 3 = 180 |
| #2. Layer Params | (320 + 3) x 20 = 6460 | (90 + 4) x 10 = 940 | (180 + 4) x 10 = 1840 |
| 2. Layer | 20 | 10 | 10 |
| Labels | 3 | 4 | 4 |

# Appendix

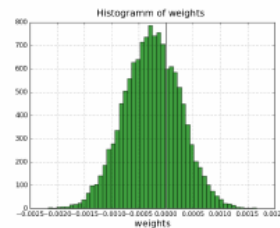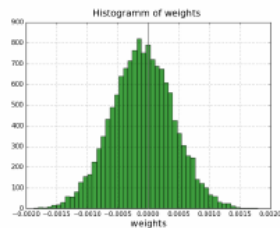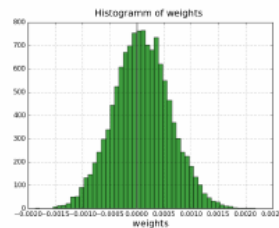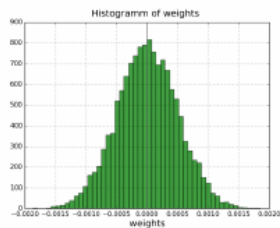|            | Stripes | Poker | Ball-Pen-Can |
|------------|---------|-------|--------------|
| Neurons    | 443     | 360   | 450          |
| Synapses   | 22140   | 18580 | 37120        |
| Parameters | 7440    | 2900  | 5760         |

© FZI Forschungszentrum Informatik

# Appendix

# Appendix

© FZI Forschungszentrum Informatik

# Appendix