# Inference and learning for sensory data
-
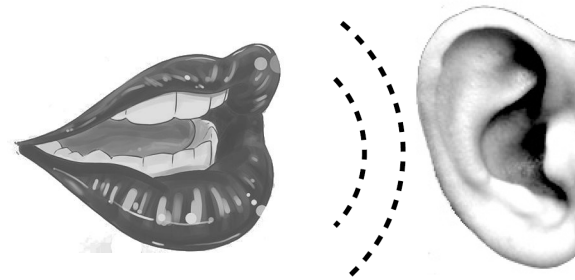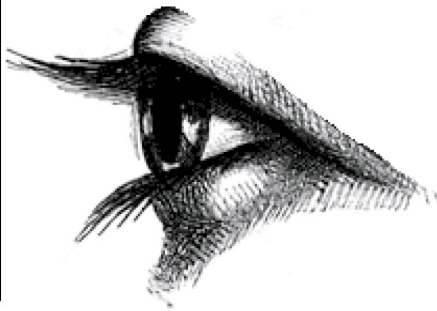# New priors, non-linear features, and the challenge of masking and occlusion

## Jörg Lücke

**Machine Learning**

Cluster of Excellence Hearing4all, Dept for Medical Physics and Acoustics,
Carl-von-Ossietzky University Oldenburg, Germany

Oldenburg, March 6, 2015

# Sensory Inference
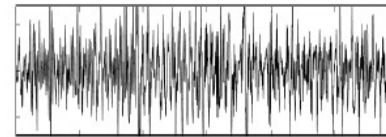


$\vec{y}^{(1)}$

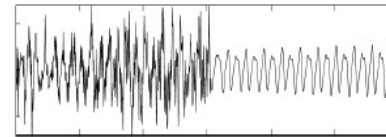$\vec{y}^{(2)}$

$\vec{y}^{(3)}$
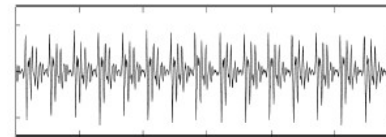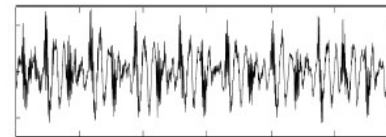
$\vec{y}^{(4)}$

$\vec{y}^{(5)}$

$\vec{y}^{(1)}$

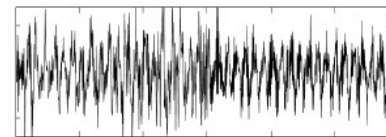$\vec{y}^{(2)}$

$\vec{y}^{(n)}$

Jörg Lücke

# Sensory Inference



**unconscious inference**

Helmholtz, *HB physiol. Optik*, 1867



$$\vec{y}^{(1)}$$

$$\vec{y}^{(2)}$$

$$\vdots$$

$$\vec{y}^{(n)}$$

$$\vdots$$

CARL VON OSSIETZKY *universität* OLDENBURG

Hearing 4all

# Sensory Inference: Example

# Sensory Inference: Example



Inference Task:

Which piano keys were pressed?

Build an artificial system that solves the task.

# Sensory Inference: Example

$$(\vec{y}^{(n)})^T$$

We re-express data point $\vec{y}^{(n)}$ :

$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

# Sensory Inference: Example



$$(\vec{y}^{(n)})^T$$

We re-express data point $\vec{y}^{(n)}$ :

$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

For piano data, we would choose:

 $\vec{W}_1^T$

 $\vec{W}_2^T$

 $\vec{W}_3^T$

 $\vec{W}_4^T$

etc.

# Sensory Inference: Example



$(\vec{y}^{(n)})^T$

We re-express data point $\vec{y}^{(n)}$:

$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

For piano data, we would choose:

 $\vec{W}_1^T$

 $\vec{W}_2^T$

 $\vec{W}_3^T$

 $\vec{W}_4^T$

etc.



$\vec{s}^{(n)}$

$p(\vec{s}\,|\,\vec{y}^{(n)})$

# Sensory Inference: Example



$(\vec{y}^{(n)})^T$

We re-express data point $\vec{y}^{(n)}$:

$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

For piano data, we would choose:

 $\vec{W}_1^T$

 $\vec{W}_2^T$

 $\vec{W}_3^T$

 $\vec{W}_4^T$

etc.

Real data causes/components:



Estimates of causes/components:



$\vec{s}^{(n)}$

$p(\vec{s} \mid \vec{y}^{(n)})$

# Sensory Inference: Example



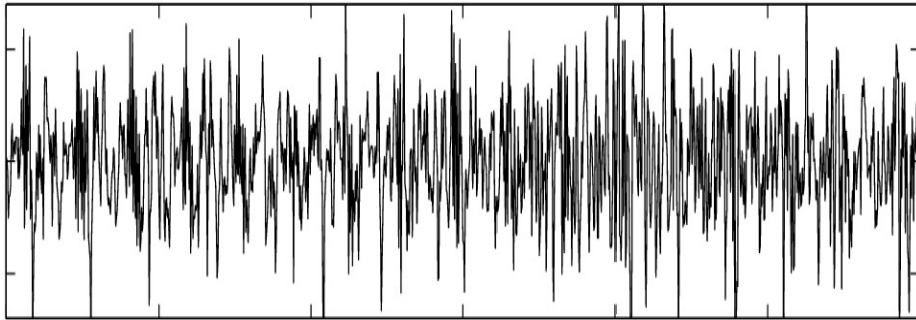$(\vec{y}^{(n)})^T$

We re-express data point $\vec{y}^{(n)}$ :

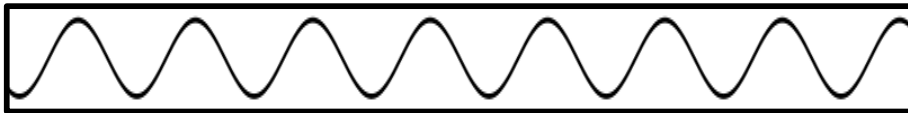$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

For piano data, we would choose:

$\vec{W}_1^T$

$\vec{W}_2^T$

$\vec{W}_3^T$

$\vec{W}_4^T$

etc.

Real data causes/components:

Estimates of causes/components:

$\vec{s}^{(n)}$

$p(\vec{s} \,|\, \vec{y}^{(n)})$

# Sensory Inference: Example



$(\vec{y}^{(n)})^T$
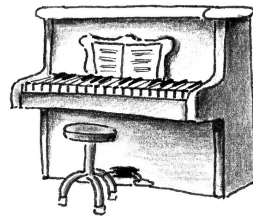
We re-express data point $\vec{y}^{(n)}$:

$$\vec{y}^{(n)} \approx \sum_{h=1}^{H} s_h^{(n)} \vec{W}_h$$

$$\boxed{\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}}$$

For piano data, we would choose:

$\vec{W}_1^T$

$\vec{W}_2^T$

$\vec{W}_3^T$

$\vec{W}_4^T$

etc.

Real data causes/components:



Estimates of causes/components:

$\vec{s}^{(n)}$

$p(\vec{s} \,|\, \vec{y}^{(n)})$

# Sensory Inference: Example



$$(\vec{y}^{(n)})^T$$

We re-express data point $\vec{y}^{(n)}$ :

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

For piano data, we would choose:



$$\vec{W}_1^T$$
$$\vec{W}_2^T$$
$$\vec{W}_3^T$$
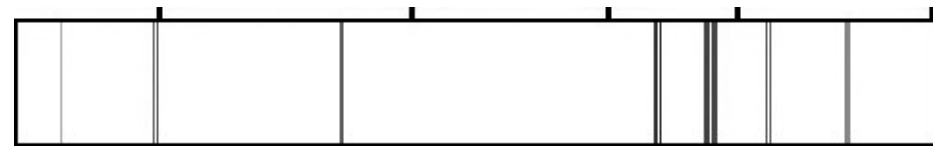$$\vec{W}_4^T$$
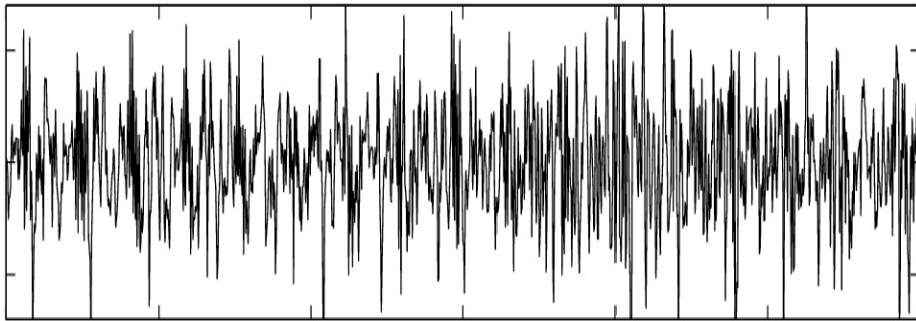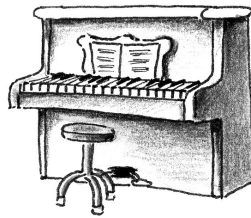
etc.

Real data causes/components:



Estimates of causes/components:



$$\vec{s}^{(n)}$$

$$p(\vec{s} \,|\, \vec{y}^{(n)})$$

# Sensory Inference: Example



$$(\vec{y}^{(n)})^T$$

We re-express data point $\vec{y}^{(n)}$ :

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

For piano data, we would choose:

 $\vec{W}_1^T$

 $\vec{W}_2^T$

 $\vec{W}_3^T$

 $\vec{W}_4^T$

etc.

Estimates of causes/components:



$$\vec{s}^{(n)}$$

$$p(\vec{s} \,|\, \vec{y}^{(n)})$$

# Sensory Inference: Example

We re-express data point $\vec{y}^{(n)}$ :

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

$(\vec{y}^{(n)})^T$

Probabilistic **generative model** (e.g., SC):

$$p(\vec{s} \mid \Theta) = \prod_h \frac{1}{\pi (1 + s_h^2)}$$

$$p(\vec{y} \mid \vec{s}, \Theta) = \mathcal{N}(\vec{y}; \sum_h s_h^{(n)} \vec{W}_h, \sigma^2 \mathbb{1})$$

For piano data, we would choose:

$\vec{W}_1^T$

$\vec{W}_2^T$

$\vec{W}_3^T$

$\vec{W}_4^T$

⋮

**dictionary**

Estimates of causes/components:

$\vec{s}^{(n)}$

$p(\vec{s} \mid \vec{y}^{(n)})$

# Sensory Inference: Example



We re-express data point $\vec{y}^{(n)}$:

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

$(\vec{y}^{(n)})^T$

Probabilistic **generative model** (e.g., SC):

$$p(\vec{s} \mid \Theta) = \prod_h \frac{1}{\pi\,(1 + s_h^2)}$$

$$p(\vec{y} \mid \vec{s}, \Theta) = \mathcal{N}(\vec{y}; \sum_h s_h^{(n)} \vec{W}_h, \sigma^2\, \mathbb{1})$$

For piano data, we would choose:

$\vec{W}_1^T$

$\vec{W}_2^T$

$\vec{W}_3^T$

$\vec{W}_4^T$

⋮

**dictionary**

Estimates of causes/components:

$\vec{s}^{(n)}$

$p(\vec{s} \mid \vec{y}^{(n)})$

# Sensory Inference: Example

We re-express data point $\vec{y}^{(n)}$:

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

$(\vec{y}^{(n)})^T$

Probabilistic **generative model** (e.g., SC):

$$p(\vec{s} \,|\, \Theta) = \prod_h \frac{1}{\pi \, (1 + s_h^2)}$$

$$p(\vec{y} \,|\, \vec{s}, \Theta) = \mathcal{N}(\vec{y}; \sum_h s_h^{(n)} \vec{W}_h, \sigma^2 \, \mathbb{1})$$

For a keyboard instrument:

$\vec{W}_1^T$

$\vec{W}_2^T$

Estimates of causes/components:

$\vec{W}_3^T$

$\vec{W}_4^T$

$\vec{s}^{(n)}$

**?**

**Solution: learn dictionary from data**

$p(\vec{s} \,|\, \vec{y}^{(n)})$

# Non-linear components



| generating system | combination of causes | model assumptions for $\vec{y}$ | dictionary $W$ |
|---|---|---|---|
| (piano) | $\approx\approx\approx$ $\oplus$ $\sim\!\!\sim$ | $\sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$ | $\approx\approx\approx$ |
| street scene | (house) (car) $\textcircled{?}$ (house+car) | $f(s_h^{(n)}, W) + \vec{\eta}$ | (house) (car) |

Jörg Lücke

# Non-linear components

Change model parameters $W$ until:

real data $\approx$ model data

Measure for this similarity:
**Data Likelihood**

| model assumptions for $\vec{y}$ | dictionary $W$ |
|---|---|
| $\sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$ |  |
| $f(s_h^{(n)}, W) + \vec{\eta}$ | |

CARL VON OSSIETZKY universität OLDENBURG

# Non-linear components

Change model parameters $W$ until:

real data $\approx$ model data

Measure for this similarity:
**Data Likelihood**



| model assumptions for $\vec{y}$ | dictionary $W$ |
|---|---|
| $\sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$ |  |
| $f(s_h^{(n)}, W) + \vec{\eta}$ |  |

Dayan & Zemel, 1996;
Lücke & Sahani, 2007, 2008;
Lücke, 2009; Lücke et al.
 etc.

E.g., **E**xpectation **M**aximization (EM) framework.
Dempster, 1977; Neal & Hinton, 1998.

Jörg Lücke

# Dictionary Examples



$$\vec{W}_1^T$$

$$\vec{W}_2^T$$

$$\vec{W}_3^T$$

$$\vec{W}_4^T$$

$$(\vec{y}^{(n)})^T \qquad \vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

CARL VON OSSIETZKY universität OLDENBURG

Hearing 4all

# Dictionary Examples



$$\vec{W}_1^T$$
$$\vec{W}_2^T$$
$$\vec{W}_3^T$$
$$\vec{W}_4^T$$

$$(\vec{y}^{(n)})^T \qquad \vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

$$\begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}^{\vec{s}}$$

Jörg Lücke

CARL VON OSSIETZKY universität OLDENBURG

Hearing 4all

# Dictionary Examples

$$\vec{y}^{(n)} = \sum_h s_h^{(n)} \vec{W}_h + \vec{\eta}$$

Olshausen & Field '96
… Sheikh et al., '14 …

Gaussian prior

Sparse prior

PCA / Factor Analysis
(vgl. principal axis transform)

Sparse coding / (ICA)



$$\vec{s} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

$$\vec{s} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

# Dictionary Examples



Spatial frequency varies

$$\vec{W}_1^T$$
$$\vec{W}_2^T$$
$$\vec{W}_3^T$$
$$\vec{W}_4^T$$

Image  DoG filtered  reconstruction  contributing components

SC  43 more

# Dictionary Examples



Olshausen & Field, *Nature* 1996



Lee, Battle, Raina, Ng, *NIPS* 2006;
Bornschein, Henniges, Lücke, *PLOS Comp Biology* 2013

Jörg Lücke

# Dictionary Examples

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$



SC ... 43 more

$$\vec{W}_1^T \ \vec{W}_2^T \ ...$$

# Dictionary Examples



$$\vec{W}_1^T \qquad \vec{W}_2^T \qquad \vec{W}_3^T$$

linearity assumption not realistic

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

SC  43 more

$$\vec{W}_1^T \ \vec{W}_2^T \ \dots$$

Hearing 4all

# Dictionary Examples



$$\vec{W}_1^T \qquad \vec{W}_2^T \qquad \vec{W}_3^T$$

$$\vec{y} = \max_h \{ s_h \vec{W}_h \} + \vec{\eta}$$

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

SC



43 more

$$\vec{W}_1^T \ \vec{W}_2^T \ \ldots$$

# Dictionary Examples



$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$

$$\vec{y} = \sum_{h} s_h \vec{W}_h + \vec{\eta}$$

**PCA** – standard approach, Pearson, 1901
**FA** – standard approach, e.g., Gorsuch, '83
**ICA** – has own conference, Comon, 1994
**SC** – Olshausen & Field, *Nature*, 1996
**NMF** – Lee & Seung, *Nature*, 1999
etc.

# Dictionary Examples



$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$

$$\vec{y} = \sum_{h} s_h \vec{W}_h + \vec{\eta}$$

crucial for
**compressed sensing**

SC
$$\vec{W}_1^T \, \vec{W}_2^T \, \ldots$$

43 more

Jörg Lücke

# Dictionary Examples



$$\vec{W}_1^T \qquad \vec{W}_2^T \qquad \vec{W}_3^T$$

$$\vec{y} = \max_h\{s_h\,\vec{W}_h\} + \vec{\eta}$$

$$\vec{y} = \sum_h s_h\vec{W}_h + \vec{\eta}$$

Dai et al., *NIPS* 2013
Bornschein et al., *PLOS CB* 2013
Shelton et al., *NIPS* 2012
Puertas, Bronschein, Lücke, *NIPS 2010*
Lücke, Sahani, *J Mach Learn Res* 2008
...
Roweis, *Eurospeech 2003*
Roweis, *NIPS 2002*
Varga & Moore, *ICASSP 1990*

SC  43 more

$$\vec{W}_1^T\ \vec{W}_2^T\ ...$$

Jörg Lücke

# Cochleagram Dictionaries



$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$

Roweis, *Eurospeech 2003*
Roweis, *NIPS 2002*
Varga, Moore, *ICASSP 1990*



$\vec{W}_1^T$ $\quad$ $\vec{W}_2^T$ $\quad$ ...

Jörg Lücke

# Computational Challenges

Jörg Lücke

# Example: Binary Prior

Binary Sparse Coding (BSC):

$$p(\vec{s}\,|\,\Theta) = \prod_h \pi^{s_h}\,(1-\pi\;)^{1-s_h}$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y};\textstyle\sum_h s_h\,\vec{W}_h, \sigma^2\,\mathbb{1})$$

Henniges et al., 2010

Maximal Causes Analysis (MCA):

$$p(\vec{s}\,|\,\Theta) = \prod_h \pi^{s_h}\,(1-\pi\;)^{1-s_h}$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y};\max_h\{s_h\,\vec{W}_h\}, \sigma^2\,\mathbb{1})$$

FIAS Frankfurt Institute for Advanced Studies

Jörg Lücke

# Example: Binary Prior

Binary Sparse Coding (BSC):

$$p(\vec{s} \mid \Theta) = \prod_h \pi^{s_h} (1 - \pi)^{1 - s_h}$$

$$p(\vec{y} \mid \vec{s}, \Theta) = \mathcal{N}(\vec{y}; \sum_h s_h \vec{W}_h, \sigma^2 \mathbb{1})$$

Henniges et al., 2010

## Optimization using Expectation Maximization (EM)

$$W^{\text{new}} = \left( \sum_n y^{(n)} \langle s \rangle_p^T \right) \left( \sum_n \langle s\, s^T \rangle_p \right)^{-1}$$

$$\sigma^{\text{new}} = \sqrt{\frac{1}{ND} \sum_n \left\langle \|y^{(n)} - W s\|^2 \right\rangle_p}$$

$$\pi^{\text{new}} = \frac{1}{ND} \sum_n \langle |s| \rangle_{p^n}$$

$$\text{with } |s| = \sum_{h=1}^{H} s_h$$

FIAS Frankfurt Institute for Advanced Studies

Jörg Lücke

# Example: Binary Prior

Binary Sparse Coding (BSC):

$$p(\vec{s} \mid \Theta) = \prod_h \pi^{s_h} (1 - \pi)^{1-s_h}$$

$$p(\vec{y} \mid \vec{s}, \Theta) = \mathcal{N}(\vec{y}; \sum_h s_h \vec{W}_h, \sigma^2 \mathbb{1})$$

Henniges et al., 2010

Expectation values scale exponentially.

$$\langle g(\boldsymbol{s}) \rangle_p = \frac{\sum_{\boldsymbol{s}} p(\boldsymbol{s}, \boldsymbol{y}^{(n)} \| \Theta^{\text{old}}) \, g(\boldsymbol{s})}{\sum_{\tilde{\boldsymbol{s}}} p(\tilde{\boldsymbol{s}}, \boldsymbol{y}^{(n)} \mid \Theta^{\text{old}})}$$

## Optimization using Expectation Maximization (EM)

$$W^{\text{new}} = \left( \sum_n \boldsymbol{y}^{(n)} \langle \boldsymbol{s} \rangle_p^T \right) \left( \sum_n \langle \boldsymbol{s}\,\boldsymbol{s}^T \rangle_p \right)^{-1}$$

$$\sigma^{\text{new}} = \sqrt{\frac{1}{ND} \sum_n \left\langle \left\| \boldsymbol{y}^{(n)} - W\,\boldsymbol{s} \right\|^2 \right\rangle_p}$$

$$\pi^{\text{new}} = \frac{1}{ND} \sum_n \langle |\boldsymbol{s}| \rangle_{p^n}$$

$$\text{with } |\boldsymbol{s}| = \sum_{h=1}^{H} s_h$$

# Example: Binary Prior

Binary Sparse Coding (BSC):

Expectation values scale exponentially.

$$p(\vec{s}\,|\,\Theta) = \prod_h \pi^{s_h}(1-\pi\,)^{1-s_h}$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y}; \textstyle\sum_h s_h \vec{W}_h, \sigma^2\,\mathbb{1})$$

$$\langle g(\boldsymbol{s})\rangle_p = \frac{\sum\limits_{\boldsymbol{s}} p(\boldsymbol{s}, \boldsymbol{y}^{(n)}\,\|\,\Theta^{\mathrm{old}})\, g(\boldsymbol{s})}{\sum\limits_{\tilde{\boldsymbol{s}}} p(\tilde{\boldsymbol{s}}, \boldsymbol{y}^{(n)}\,|\,\Theta^{\mathrm{old}})}$$

Henniges et al., 2010

$$\langle g(\boldsymbol{s})\rangle_p = \frac{\sum_a p(\vec{s}_a, \vec{y}^{(n)}\,|\,\Theta')\, g(\vec{s}) + \sum\limits_{\substack{a,b \\ a<b}} p(\vec{s}_{ab}, \vec{y}^{(n)}\,|\,\Theta')\, g(\vec{s}) + \ldots}{p(\vec{0}, \vec{y}^{(n)}\,|\,\Theta') + \sum_a p(\vec{s}_a, \vec{y}^{(n)}\,|\,\Theta') + \sum\limits_{\substack{a,b \\ a<b}} p(\vec{s}_{ab}, \vec{y}^{(n)}\,|\,\Theta') + \ldots}$$

**Idea: Truncate the sums**

where $\quad \vec{s}_a := (0,\ldots,0,1,0,\ldots,0)$ with only $s_a = 1$

$\vec{s}_{ab} := (0,\ldots,0,1,0,\ldots,0,1,0,\ldots,0)$ with only $s_a = 1, s_b = 1, a \neq b$,

and $\vec{s}_{abc}$ etc. are defined analogously.

# Example: Binary Prior

Binary Sparse Coding (BSC):

Expectation values scale exponentially.

$$
p(\vec{s}\,|\,\Theta) = \prod_h \pi^{s_h}\,(1-\pi\,)^{1-s_h}
$$

$$
p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y};\textstyle\sum_h s_h\,\vec{W}_h,\sigma^2\,\mathbb{1})
$$

$$
\left\langle g(\boldsymbol{s})\right\rangle_{\boldsymbol{q_n}} = \frac{\displaystyle\sum_{\boldsymbol{s}\in\mathcal{K}_n} p(\boldsymbol{s},\boldsymbol{y}^{(n)}\,\|\,\Theta^{\mathrm{old}})\,g(\boldsymbol{s})}{\displaystyle\sum_{\tilde{\boldsymbol{s}}\in\mathcal{K}_n} p(\tilde{\boldsymbol{s}},\boldsymbol{y}^{(n)}\,|\,\Theta^{\mathrm{old}})}
$$

$$
\langle g(\boldsymbol{s})\rangle_p = \frac{\sum_a p(\vec{s}_a,\vec{y}^{(n)}\,|\,\Theta')\,g(\vec{s}) + \sum_{\substack{a,b\\a<b}} p(\vec{s}_{ab},\vec{y}^{(n)}\,|\,\Theta')\,g(\vec{s}) + \dots}{p(\vec{0},\vec{y}^{(n)}\,|\,\Theta') + \sum_a p(\vec{s}_a,\vec{y}^{(n)}\,|\,\Theta') + \sum_{\substack{a,b\\a<b}} p(\vec{s}_{ab},\vec{y}^{(n)}\,|\,\Theta') + \dots}
$$

Such a truncation of sums is equivalent to a variational approximation:

$$
\tilde{q}^{(n)}(\vec{s};\Theta^{\mathrm{old}}) = \frac{p(\vec{s},\vec{y}^{(n)}\,|\,\Theta^{\mathrm{old}})}{\sum_{\vec{s}'\in\mathcal{K}_n} p(\vec{s}',\vec{y}^{(n)}\,|\,\Theta^{\mathrm{old}})}\,\delta(\vec{s}\in\mathcal{K}_n)
$$

variational distribution (not factored)

Jörg Lücke

# Example: Binary Prior

Binary Sparse Coding (BSC):

Expectation values scale exponentially.

$$p(\vec{s}\,|\,\Theta) = \prod_h \pi^{s_h}(1-\pi)^{1-s_h}$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y}; \textstyle\sum_h s_h \vec{W}_h, \sigma^2 \mathbb{1})$$

Henniges et al., 2010

$$\langle g(\boldsymbol{s})\rangle_p = \frac{\sum_{\boldsymbol{s}} p(\boldsymbol{s},\boldsymbol{y}^{(n)}\,\|\,\Theta^{\mathrm{old}})\,g(\boldsymbol{s})}{\sum_{\tilde{\boldsymbol{s}}} p(\tilde{\boldsymbol{s}},\boldsymbol{y}^{(n)}\,|\,\Theta^{\mathrm{old}})}$$

**Idea: Truncate the sums**

$$\langle g(\boldsymbol{s})\rangle_p = \frac{\sum_a p(\vec{s}_a,\vec{y}^{(n)}\,|\,\Theta')\,g(\vec{s}) + \sum_{\substack{a,b \\ a<b}} p(\vec{s}_{ab},\vec{y}^{(n)}\,|\,\Theta')\,g(\vec{s}) + \dots}{p(\vec{0},\vec{y}^{(n)}\,|\,\Theta') + \sum_a p(\vec{s}_a,\vec{y}^{(n)}\,|\,\Theta') + \sum_{\substack{a,b \\ a<b}} p(\vec{s}_{ab},\vec{y}^{(n)}\,|\,\Theta') + \dots}$$

**Expectation Truncation:**

$$q_n(\vec{s};\Theta) = \frac{1}{A}\,p(\vec{s}\,|\,\vec{y}^{(n)},\Theta)\,\delta(\vec{s}\in\mathcal{K}_n)$$

FIAS Frankfurt Institute for Advanced Studies

Jörg Lücke

# Relation to Other Approximations

exact: $\qquad q_n(\vec{s}; \Theta) = p(\vec{s} \mid \vec{y}^{(n)}, \Theta)$

MAP: $\qquad q_n(\vec{s}; \Theta) = \delta(\vec{s} - \vec{s}^{\max})$

Laplace: $\qquad q_n(\vec{s}; \Theta) = \mathcal{N}(\vec{s}; \vec{s}^{\max}, \Sigma)$

mean-field: $\qquad q_n(\vec{s}; \Theta) = \prod_h q_{h, \vec{\lambda}_n}^{(n)}(s_h; \Theta)$

truncated: $\qquad q_n(\vec{s}; \Theta) = \frac{1}{A} p(\vec{s} \mid \vec{y}^{(n)}, \Theta) \, \delta(\vec{s} \in \mathcal{K}_n)$

# Expectation Truncation

$$p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$$

latent space

$$\mathcal{K}_n$$

data point

$$\vec{y}^{(n)}$$

**ET**

# Expectation Truncation

$$p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$$

preselect subset $\mathcal{K}_n$

evaluate all states in $\mathcal{K}_n$



latent space

$\mathcal{K}_n$

data point $\vec{y}^{(n)}$

**ET**

# Expectation Truncation

Lücke, Eggert, *JMLR* 2010

preselect subset $\mathcal{K}_n$

evaluate all states in $\mathcal{K}_n$



$\mathcal{K}_n$

$\vec{y}^{(n)}$

**ET**

$p(\vec{s} \mid \vec{y}^{(n)}, \Theta)$

$\vec{y}^{(n)}$

**optimal case**

# Expectation Truncation

$p(\vec{s} \,|\, \vec{y}^{(n)}, \Theta)$

**deterministic**          **ET**                    **optimal case**

# Expectation Truncation

$$p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$$

**Exact:** $\quad q_n(\vec{s}; \Theta) \;=\; p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$

**MAP:** $\quad q_n(\vec{s}; \Theta) \;=\; \delta(\vec{s} - \vec{s}^{\,\mathrm{max}})$
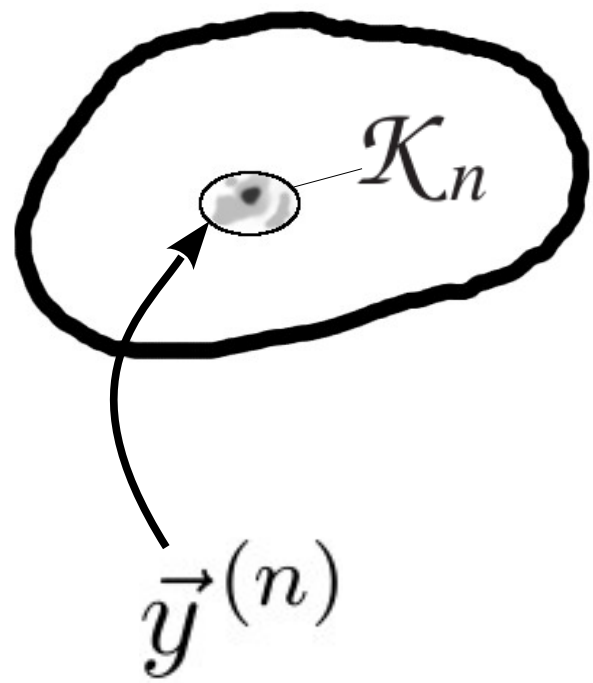
FIAS Frankfurt Institute for Advanced Studies

Jörg Lücke

# Expectation Truncation

Lücke, Eggert, *JMLR* 2010



$$p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$$

**Exact:** $\quad q_n(\vec{s};\Theta) \;=\; p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)$

**ET:** $\quad q_n(\vec{s};\Theta) \;=\; \frac{1}{A}\, p(\vec{s}\,|\,\vec{y}^{(n)}, \Theta)\, \delta(\vec{s}\in\mathcal{K}_n)$

**MAP:** $\quad q_n(\vec{s};\Theta) \;=\; \delta(\vec{s}-\vec{s}^{\,\mathrm{max}})$

# Expectation Truncation (ET)



preselect subset $\mathcal{K}_n$ $\quad\quad \vec{y}^{(n)} \quad\quad$ evaluate states in $\mathcal{K}_n$
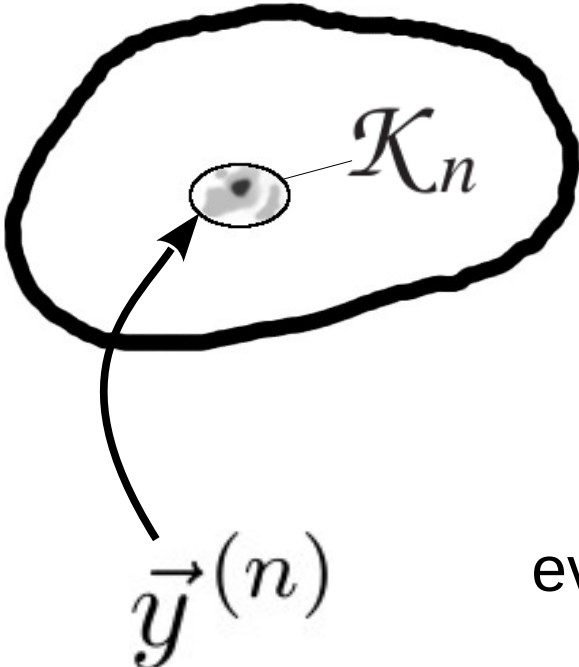
discriminative approaches $\longleftrightarrow$ generative approaches

# Expectation Truncation (ET)



Lücke, Eggert, *JMLR* 2010
Henniges et al., *LVA* 2010
Puertas et al., *NIPS* 2010
Shelton et al., *NIPS* 2011
Exarchakis et al., *LVA* 2012
Dai, Lücke, *CVPR* 2012a
Dai, Lücke, *CVPR* 2012b
Shelton et al., *NIPS* 2012
Bornschein et al., *PLOS CB* 2013
Sheikh et al., *JMLR* 2014
Henniges et al., *JMLR* 2014

preselect subset $\mathcal{K}_n$          evaluate states in $\mathcal{K}_n$

discriminative approaches ⟷ generative approaches

# Expectation Truncation (ET)

Lücke, Eggert, *JMLR* 2010

| variational approximation | correlations | multiple modes | example papers |
|---|---|---|---|
| max a-posteriori (MAP) | no | no | Olshausen, Field, 1996; A. Ng et al. |
| Gaussian | yes | no | Opper et al.; Seeger, 2008 |
| mean-field | no | yes | Titsias et al., 2011; Goodfellow … Bengio, 2012; |
|  |  |  |  |

FIAS Frankfurt Institute for Advanced Studies

Jörg Lücke

# Expectation Truncation (ET)

Lücke, Eggert, *JMLR* 2010

| variational approximation | correlations | multiple modes | example papers |
| --- | --- | --- | --- |
| max a-posteriori (MAP) | no | no | Olshausen, Field, 1996; A. Ng et al. |
| Gaussian | yes | no | Opper et al.; Seeger, 2008 |
| mean-field | no | yes | Titsias et al., 2011; Goodfellow … Bengio, 2012; |
| expectation truncation | yes | yes | Sheikh, Shelton, Lücke, 2012; Puertas … '10; Dai&Lücke, '12; |

**FIAS** Frankfurt Institute for Advanced Studies
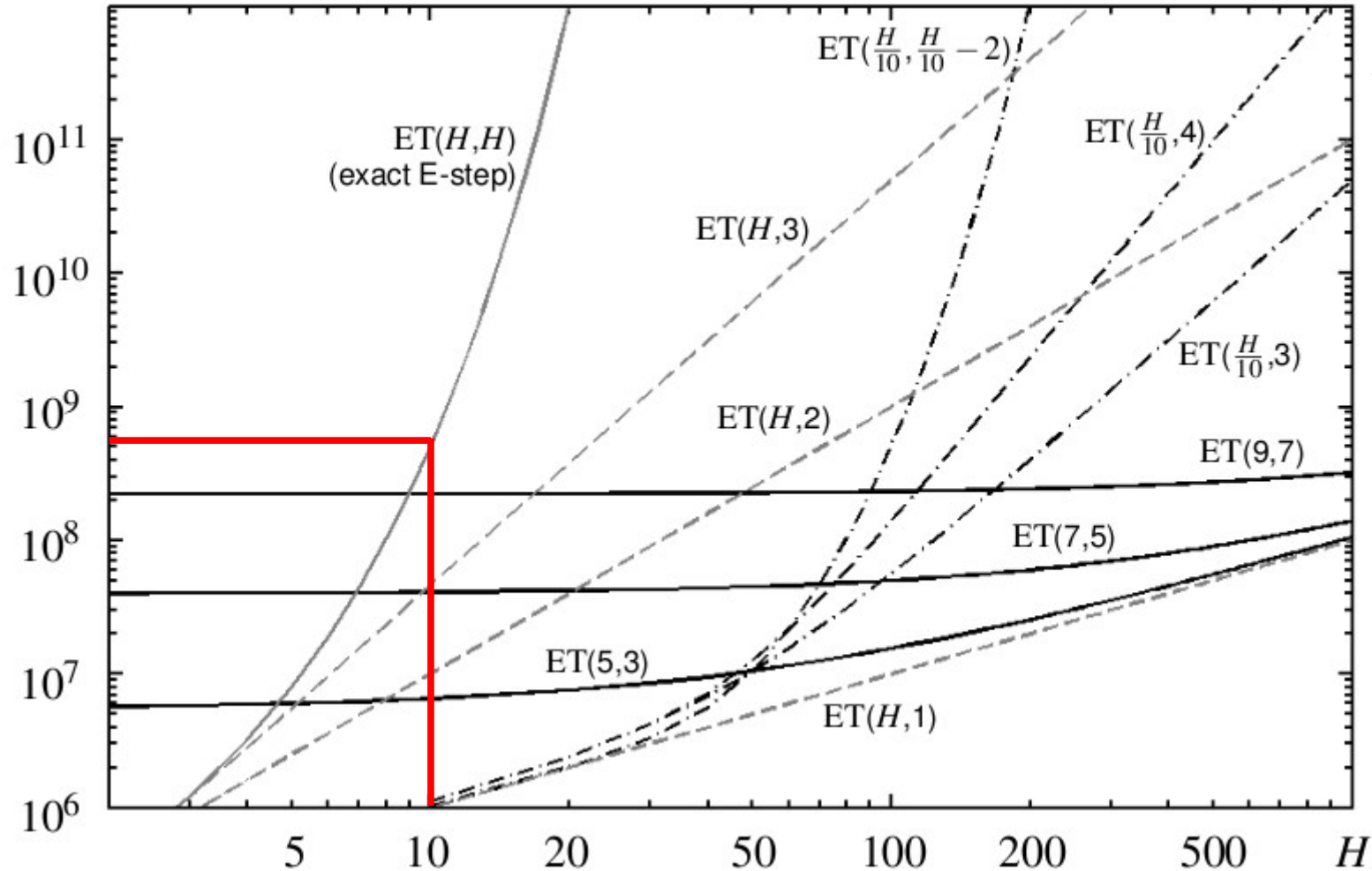
Jörg Lücke

# Expectation Truncation (ET)

Lücke, Eggert, *JMLR* 2010

| variational approximation | correlations | multiple modes | broad posterior distributions |
|---|---|---|---|
| max a-posteriori (MAP) | no | no | no |
| Gaussian | yes | no | yes |
| mean-field | no | yes | yes |
| expectation truncation | yes | yes | no |

# Expectation Truncation



E-step complexity

Lücke, Eggert, *JMLR* 2010;

exact EM

$$\mathcal{O}(e^{H})$$

FIAS Frankfurt Institute for Advanced Studies

# Expectation Truncation



E-step complexity

Lücke, Eggert, *JMLR* 2010;

ET parametrizes accuracy

**ET-EM**

$\mathcal{O}(H)$

**Expectation Truncation**

FIAS Frankfurt Institute for Advanced Studies

# Expectation Truncation



E-step complexity

Lücke, Eggert, *JMLR* 2010;

$ET(\frac{H}{10}, \frac{H}{10} - 2)$

$ET(H,H)$ (exact E-step)

$ET(\frac{H}{10}, 4)$

$ET(H,3)$

$ET(\frac{H}{10}, 3)$

$ET(H,2)$

$ET(9,7)$

$ET(7,5)$

$ET(5,3)$

$ET(H,1)$

**ET-EM**

$\mathcal{O}(H)$

## ET allows for optimizing prior parameters

Puertas et al., *NIPS* 2010;      Shelton et al., *NIPS* 2011
Henniges et al., *LVA/ICA* 2010;      Dai & Lücke, *CVPR* 2012 a & b
*Lücke, Eggert, JMLR* 2010;      Shelton *et al., NIPS* 2012

**FIAS** Frankfurt Institute for Advanced Studies

Jörg Lücke

E-step complexity



MCA generative model:

$$p(\vec{s}\,|\,\Theta) = \prod_h \pi_h^{s_h} (1 - \pi_h)^{1-s_h}$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y}; \max_h \{s_h^{(n)} \vec{W}_h\}, \sigma^2 \, \mathbb{1})$$

- multiple modes $\Big\}$ $\Big\{$ non-linear models
- correlations           advanced linear

**Applicable to large-scale since 2010.**

up to 4000 cores or GOLD (16 GPUs)

**- Problem: local likelihood optima**
→ simulated annealing

**- Problem: no closed-form M-steps**
→ general sol. for non-linear models
Dai, Lücke, *TPAMI* 2014
Lücke et al., *NIPS* 2009
Lücke, Sahani, *JMLR* 2008

**- Problem: E-step comput. intractable**
Dai et al.,, *NIPS* 2013
Shelton et al., *NIPS* 2012
Shelton et al., *NIPS* 2011
Puertas et al., *NIPS* 2010
Lücke, Eggert, *JMLR* 2010
Lücke et al., *NIPS* 2009
Lücke, Sahani, *JMLR* 2008

# Bars Test

$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$



The Bars Test, Földiák, 1990

Jörg Lücke

# Bars Test

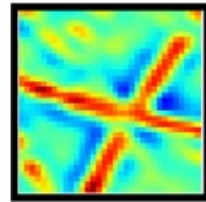$$\vec{y} = \max_h\{s_h \vec{W}_h\} + \vec{\eta}$$
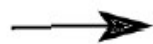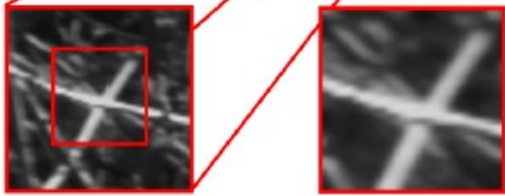
The Bars Test, Földiák, 1990
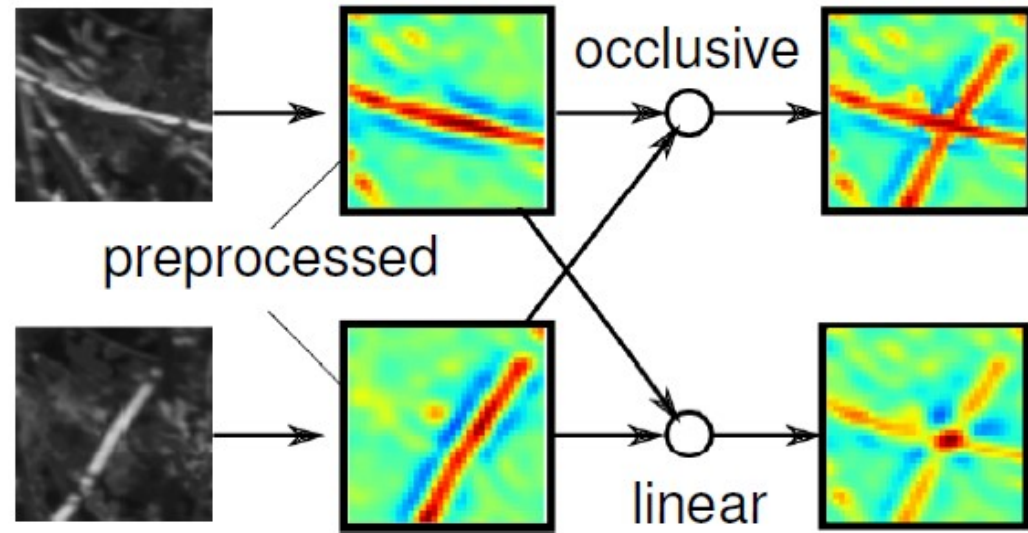
Lücke and Sahani, *JMLR* 2008
Lücke and Eggert, *JMLR* 2010

Jörg Lücke

# Encoding of Visual Scenes



$$\vec{y} = \max_h \{ s_h \vec{W}_h \} + \vec{\eta}$$

occlusive

preprocessed

linear

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

LGN

CARL VON OSSIETZKY universität OLDENBURG
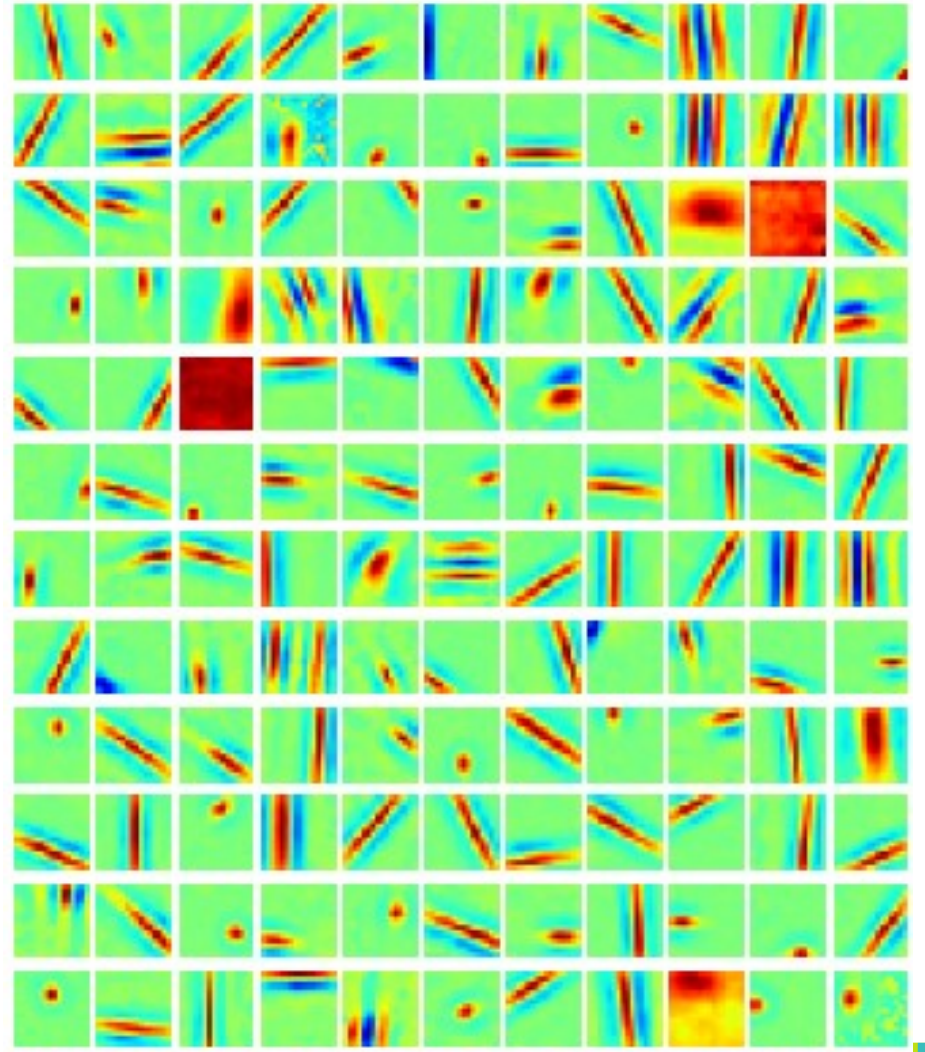
Hearing 4all

# Sparse Coding

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

# Non-linear Sparse Coding

$$\vec{y} = \max_h \{ s_h \vec{W}_h \} + \vec{\eta}$$

Puertas et al. 2010; Bornschein et. al., 2013



Response: Spike Times

Stimulus:

Olshausen & Field
*Nature* 1996

Nobel Prize in Physiology 1981
Hubel & Wiesel

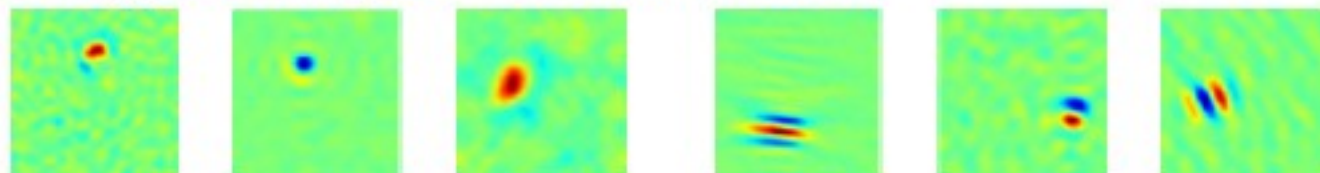## Sparse Coding

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

## Non-linear Sparse Coding

$$\vec{y} = \max_h \{s_h \vec{W}_h\} + \vec{\eta}$$

Puertas et al. 2010; Bornschein et. al., 2013
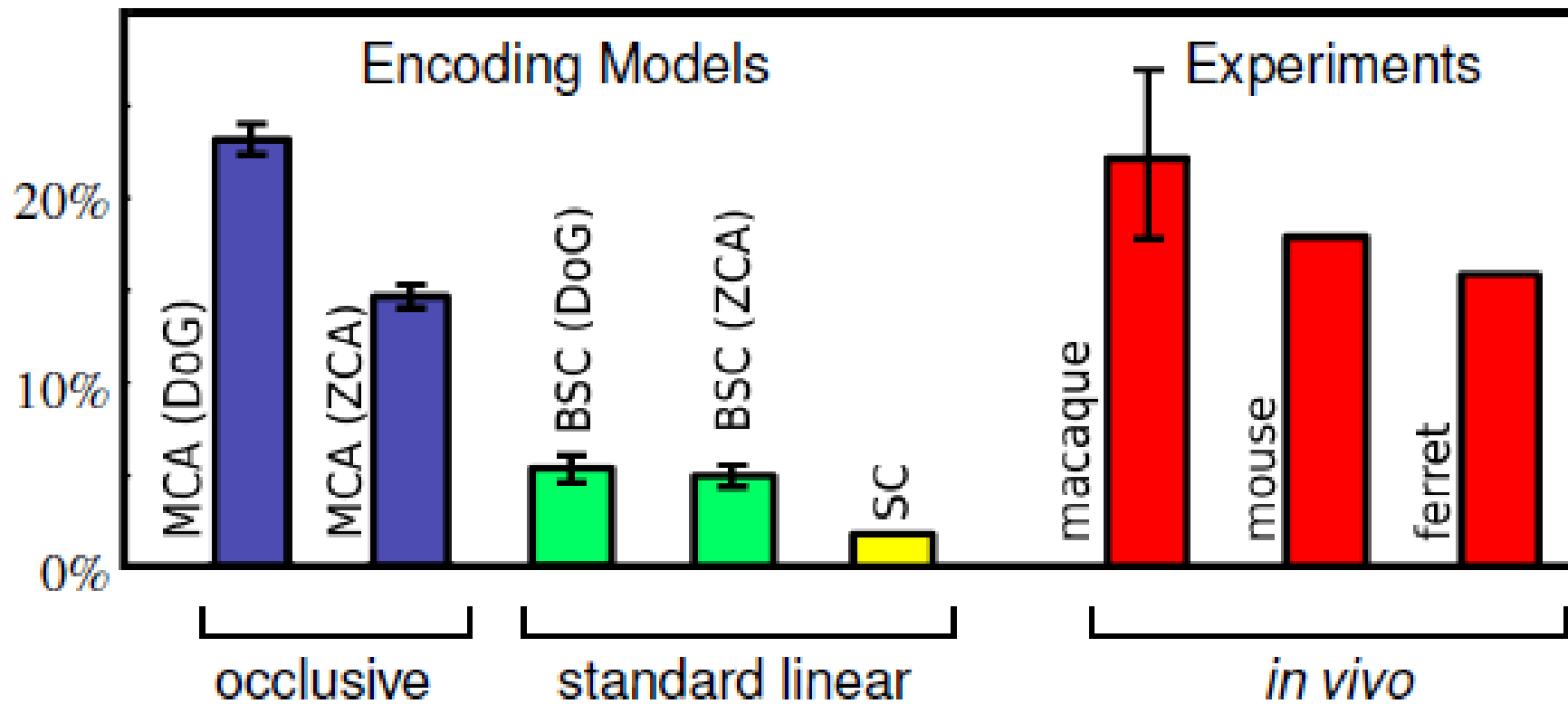
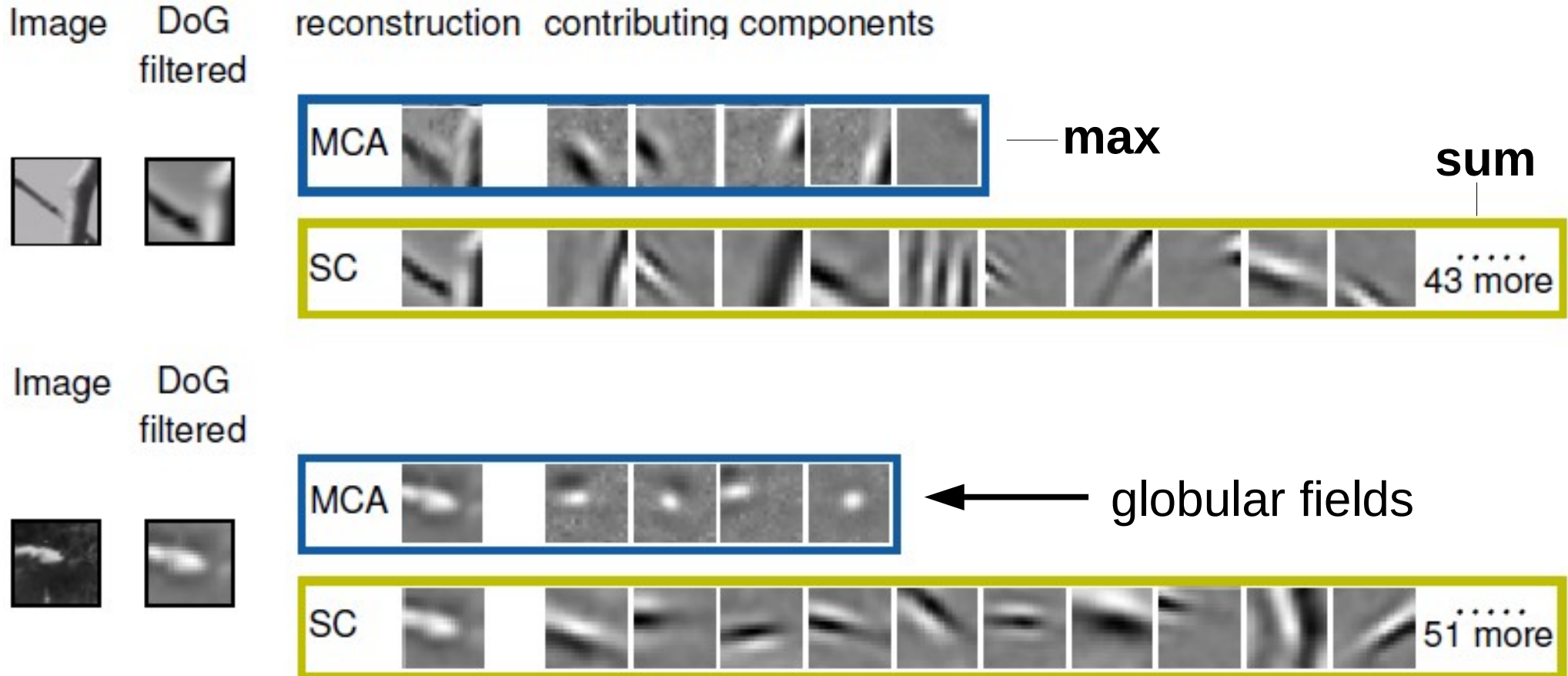Two types of simple cell RFs are measured

globular fields    standard Gabors    macaque

non-linear model (MCA)
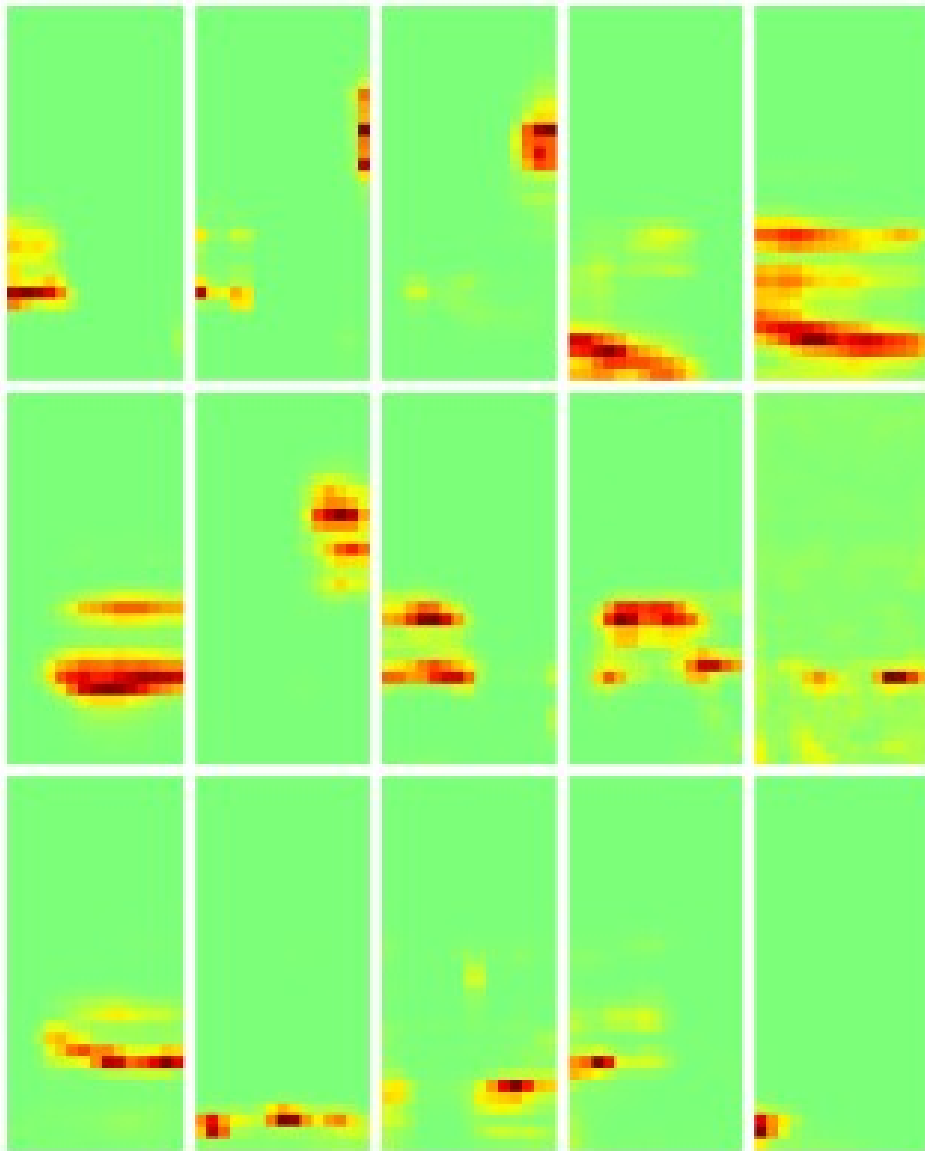
Predicted and Measured Percentages of Globular RFs

Encoding Models    Experiments

occlusive    standard linear    in vivo

Bornschein, Henniges, Lücke, *PLOS Comp Biology* 2013

# Encoding of Visual Scenes

("Viterbi" for sparse coding)
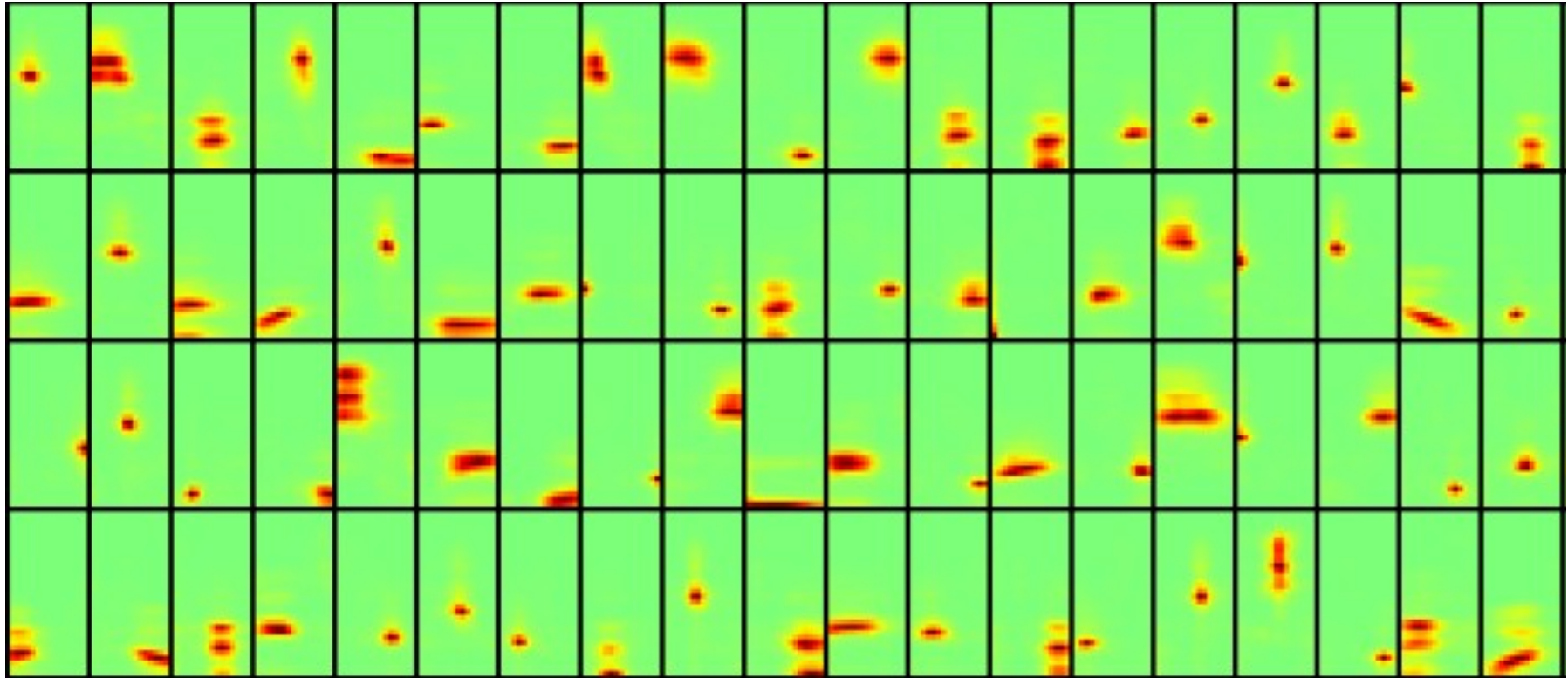
Jörg Lücke

# Encoding of Acoustic Scenes



Apply non-linear model:

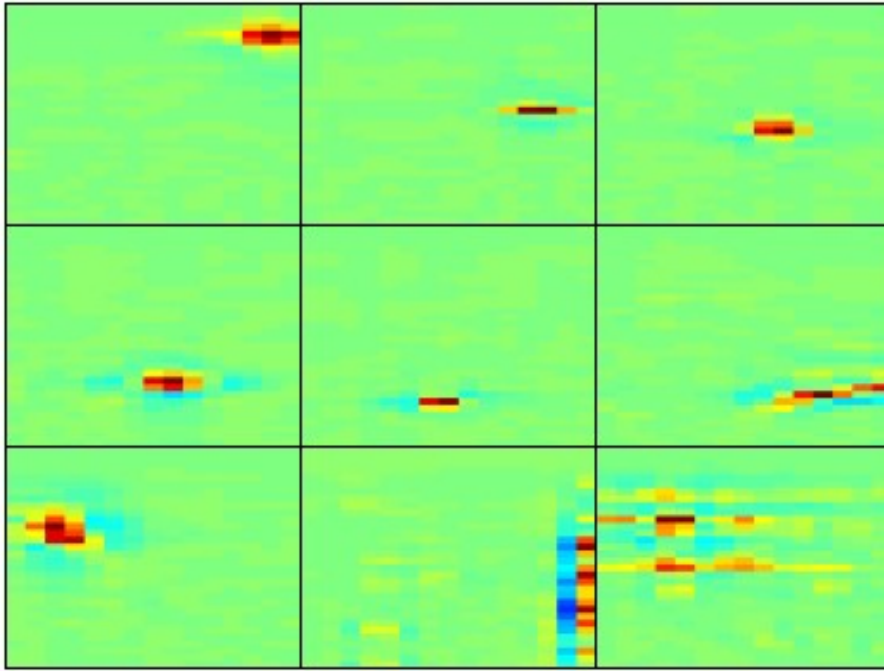$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$
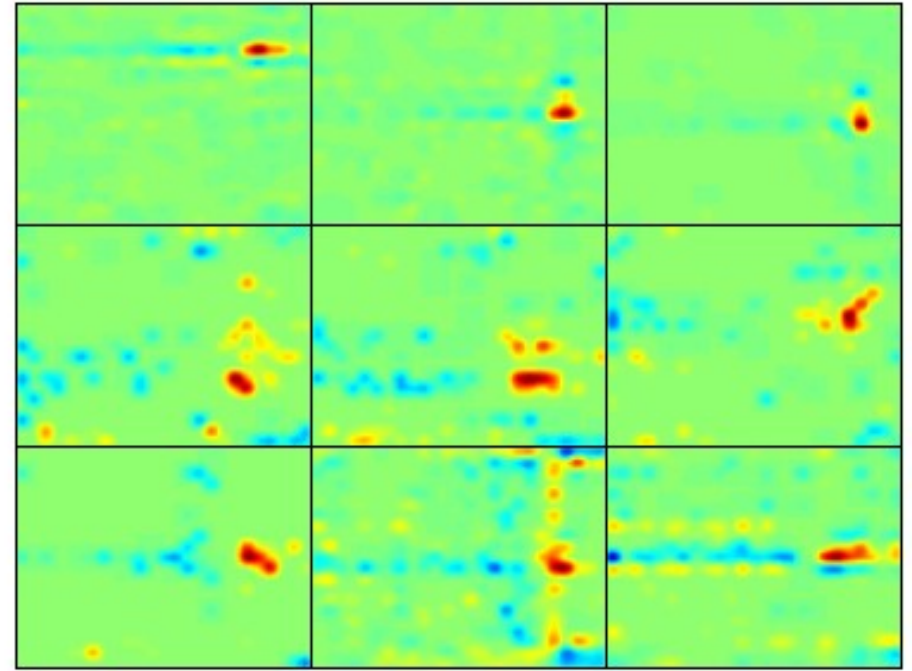


Jörg Lücke

# Encoding of Acoustic Scenes



generative fields (selection of H=1000 fields)

**Currently ongoing work** (Univ. Oldenburg / UC Berkeley / Univ. Cambridge / TU Berlin)

Jörg Lücke

# Encoding of Acoustic Scenes



STRFs learned by the model



Ferret A1 recordings

Both estimated using (regularized) reverse correlation.

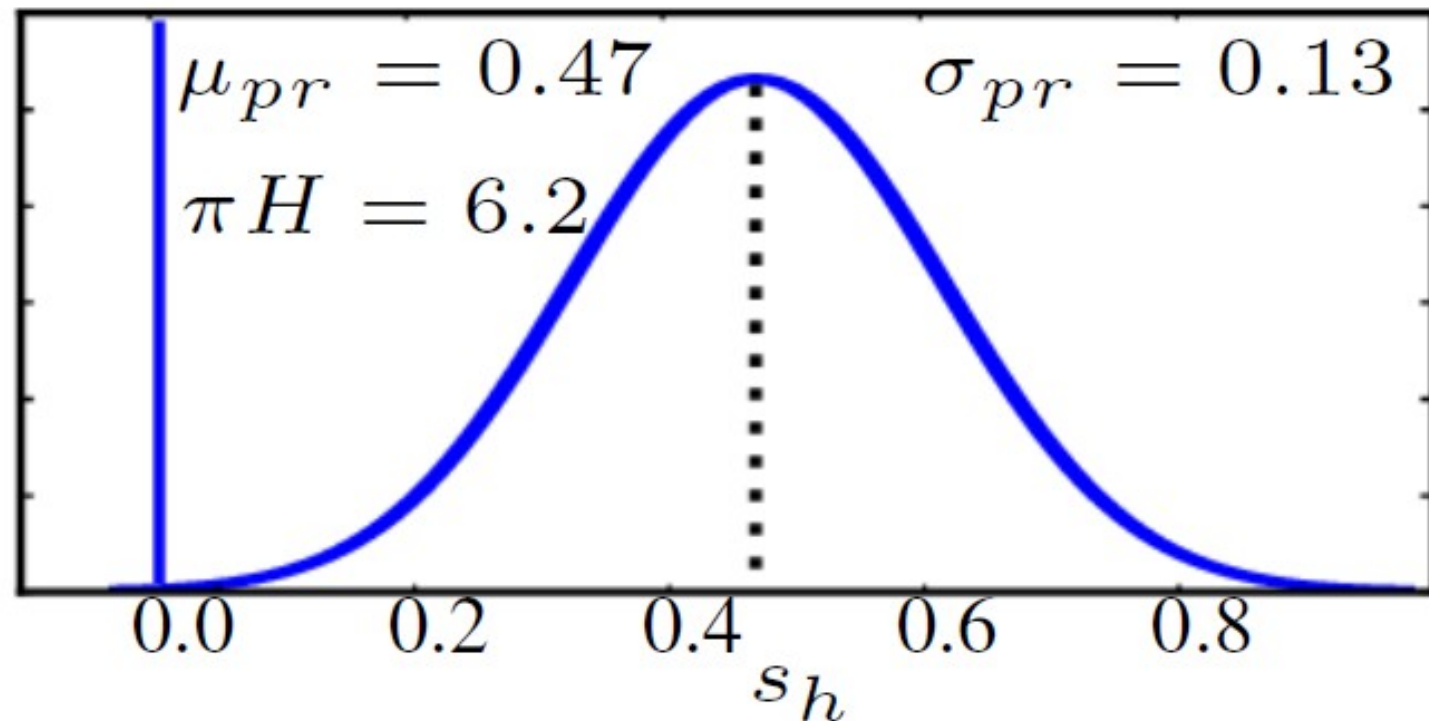**Currently ongoing work** (Univ. Oldenburg / UC Berkeley / Univ. Cambridge / TU Berlin)

Jörg Lücke

# Example: Spike-and-Slab Sparse Coding (GSC)

$$p(\vec{s}|\Theta) = \mathcal{B}(\vec{s}; \vec{\pi}) = \prod_{h=1}^{H} \pi_h^{s_h}(1-\pi_h)^{1-s_h}$$

$$p(\vec{z}|\Theta) = \mathcal{N}(\vec{z}; \vec{\mu}, \Psi)$$

$$p(\vec{y}\,|\,\vec{s}, \vec{z}, \Theta) = \mathcal{N}(\vec{y}; W(\vec{s}\odot\vec{z}), \Sigma)$$

Titsias, Lazaro-Gredilla, *NIPS* '11
Sheikh, Lücke, *LVA* '12
Goodfellow et al., *ICML* '12
Shelton et al., *NIPS* '12
… and more



$\mu_{pr} = 0.47 \qquad \sigma_{pr} = 0.13$

$\pi H = 6.2$

# Example: Gaussian Sparse Coding (GSC)

$$p(\vec{s}\,|\,\Theta) \;=\; \mathcal{B}(\vec{s};\vec{\pi}) = \prod_{h=1}^{H} \pi_h^{s_h}\,(1-\pi_h)^{1-s_h}$$

$$p(\vec{z}\,|\,\Theta) \;=\; \mathcal{N}(\vec{z};\vec{\mu},\Psi) \quad p(\vec{y}\,|\,\vec{s},\vec{z},\Theta) = \mathcal{N}(\vec{y};W(\vec{s}\odot\vec{z}),\Sigma)$$

| Noise | Noisy img | MTMKL$^{exp.}$ | K-SVD$^{mis.}$ | *K-SVD$^{match}$ | Beta pr. | GSC (H=64) | GSC (H=256) |
|---|---|---|---|---|---|---|---|
| | | | | | | | PSNR (dB) |
| $\sigma=15$ | 24.59 | **34.29** | 30.67 | 34.22 | 34.19 | 32.68 (H'=10,$\gamma$=8) | 33.78 (H'=18,$\gamma$=3) |
| $\sigma=25$ | 20.22 | 31.88 | 31.52 | 32.08 | 31.89 | 31.10 (H'=10,$\gamma$=8) | **32.01** (H'=18,$\gamma$=3) |
| $\sigma=50$ | 14.59 | 28.08 | 19.60 | 27.07 | 27.85 | 28.02 (H'=10,$\gamma$=8) | **28.35** (H'=10,$\gamma$=8) |

Lücke, Sheikh, *LVA* 2012;
Sheikh, Shelton, Lücke, *JMLR* 2014.



**GSC is state-of-the-art in denoising.**

… but denoising is just one tasks.

# Maximal Causes



cause 1

} cause

} cause

cause 2

**Maximal Causes Analysis (MCA) is state-of-the-art**
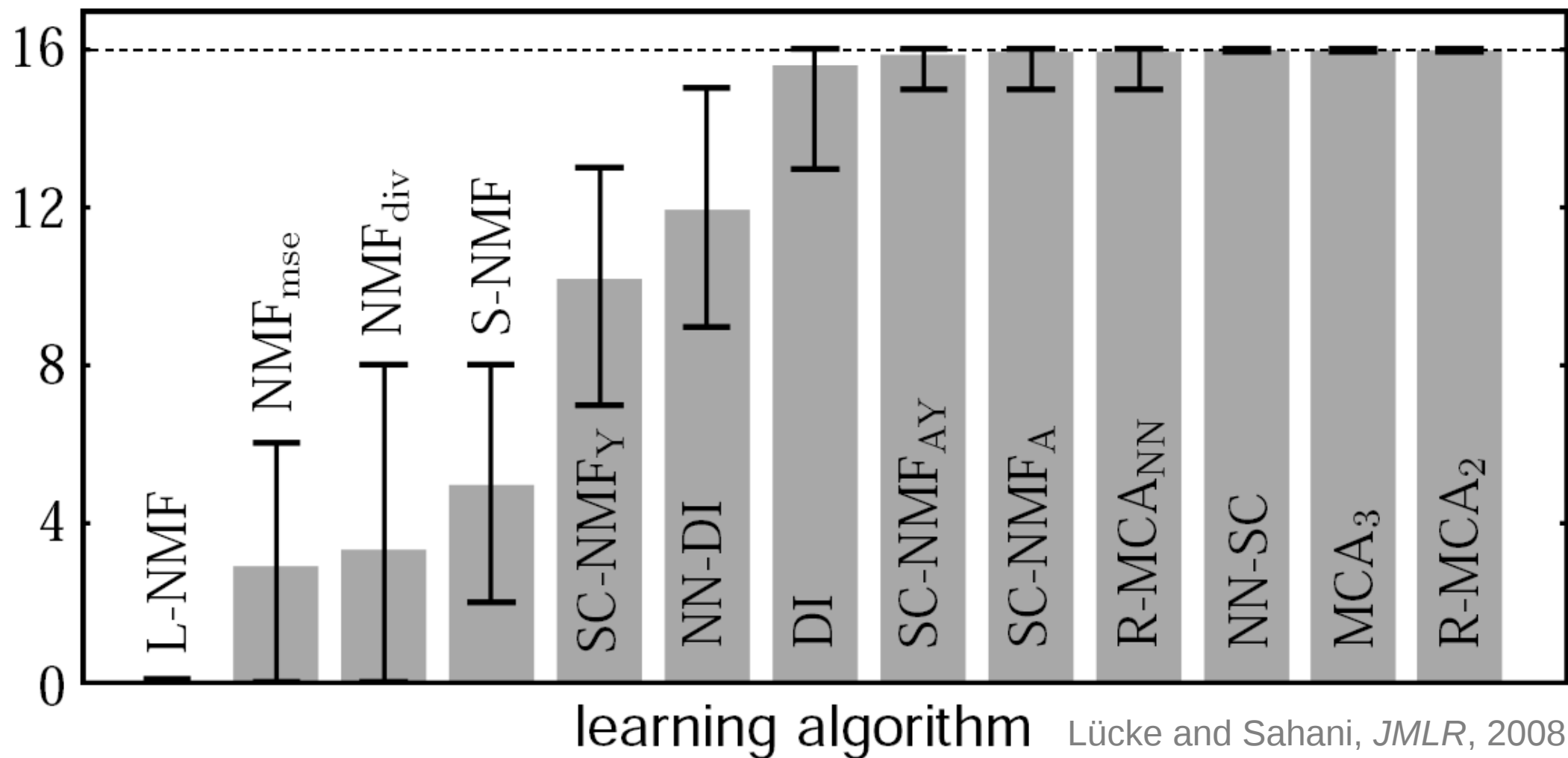
Source:
Hoyer, *J Mach Learning Res*, 2004
Spratling, *J Mach Learning Res*, 2006
Lücke/Sahani, *J Mach Learning Res*, 2008



bars

learning algorithm

# Selection of Statistical Models

**noisy-OR**

J. Bornschein

**occlusion**

M. Henniges

**exclusion**

Z. Dai

**mixtures**
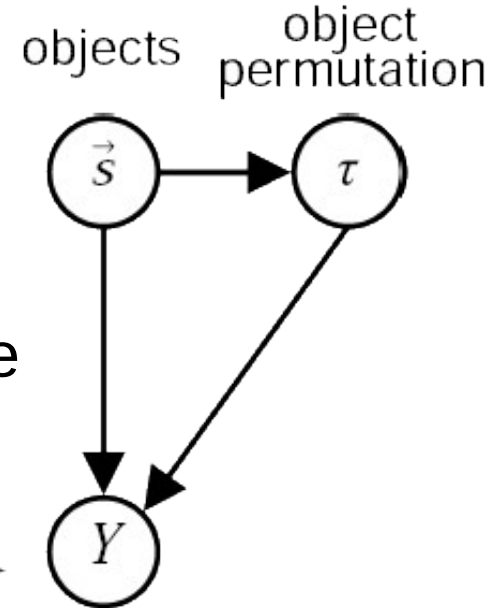
C. Keck,
S. Sheikh,
C. Savin

**linear SC**

A.-S. Sheikh,
M. Henniges
J. Shelton



objects     object permutation

Mask    Feature

$$\vec{\mathcal{T}}_d(S;\Theta) = W_{h_o d}\,\vec{T}_{h_o}$$
$$h_o = \mathrm{argmax}_h\{\tau(h)W_{hd}\}$$

Graphical Model

Lücke et al., *NIPS* 2009

Jörg Lücke

# Example App: Inpainting



Original image          80% lost pixels          reconstruction

spike-and-slab prior

Shelton et al., *NIPS 2012*

$$\vec{y} = \max_{h}\{s_h \vec{W}_h\} + \vec{\eta}$$

# Example: Structured Noise Removal



Jörg Lücke

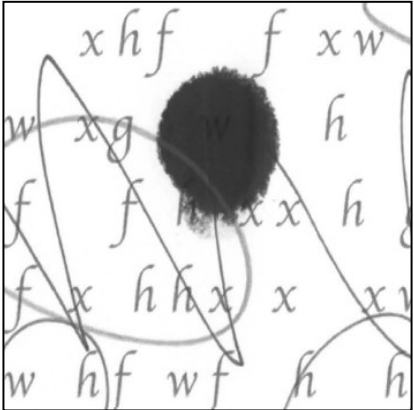# Example: Structured Noise Removal

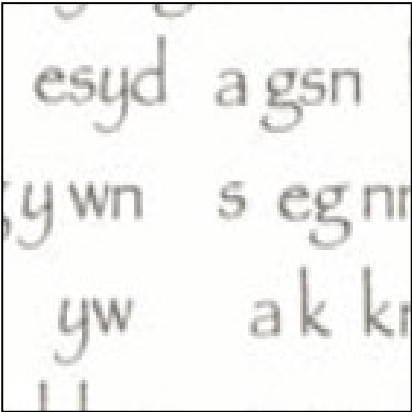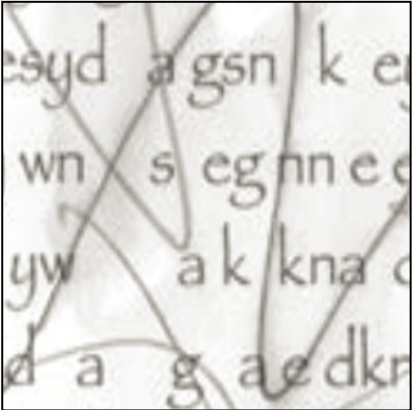# Example: Structured Noise Removal

| | 9chars | Klingon | Rotated, random placed | Occluded |
|---|---|---|---|---|

original

# Thanks to:

**Cluster of Excellence Hearing4all**
Deutsche Forschungsgemeinschaft

**Non-linear Probabilistic Models for Representational Recognition and Unsupervised Learning in Vision**
Deutsche Forschungsgemeinschaft

SPONSORED BY THE

Federal Ministry of Education and Research

**Bernstein Focus Neurotechnology Frankfurt**
Bundesministerium für Bildung und Forschung

**LOROC**
Honda Research Institute Europe

Jörg Lücke

# Thanks to:



Christian Keck

Jörg Bornschein

Marc Henniges

Dennis Forster

Abdul Saboor Sheikh

Jacquelyn Shelton

Jörg Lücke

Zhenwen Dai

Georgios Exarchakis

Project Researchers
Univ. Oldenburg, TU Berlin, Univ. Frankfurt

Jörg Lücke

# Thanks to:

**Richard Turner**
**Cristina Savin** (now IST Austria)
Cambridge University, UK

**Roland Memisevic**
**Jörg Bornschein**
University of Montreal, CA

**Maneesh Sahani**
Gatsby Unit, UCL, UK

**Bruno Olshausen**
**Nicol Harper**
Berkeley University, US

**Jörg Lücke**
University of Oldenburg

**Julian Eggert**
Honda-RI, Europe

**Pietro Berkes**
Enthought, Ltd., UK
Brandeis Univ., USA

**Marc-Thilo Figge**
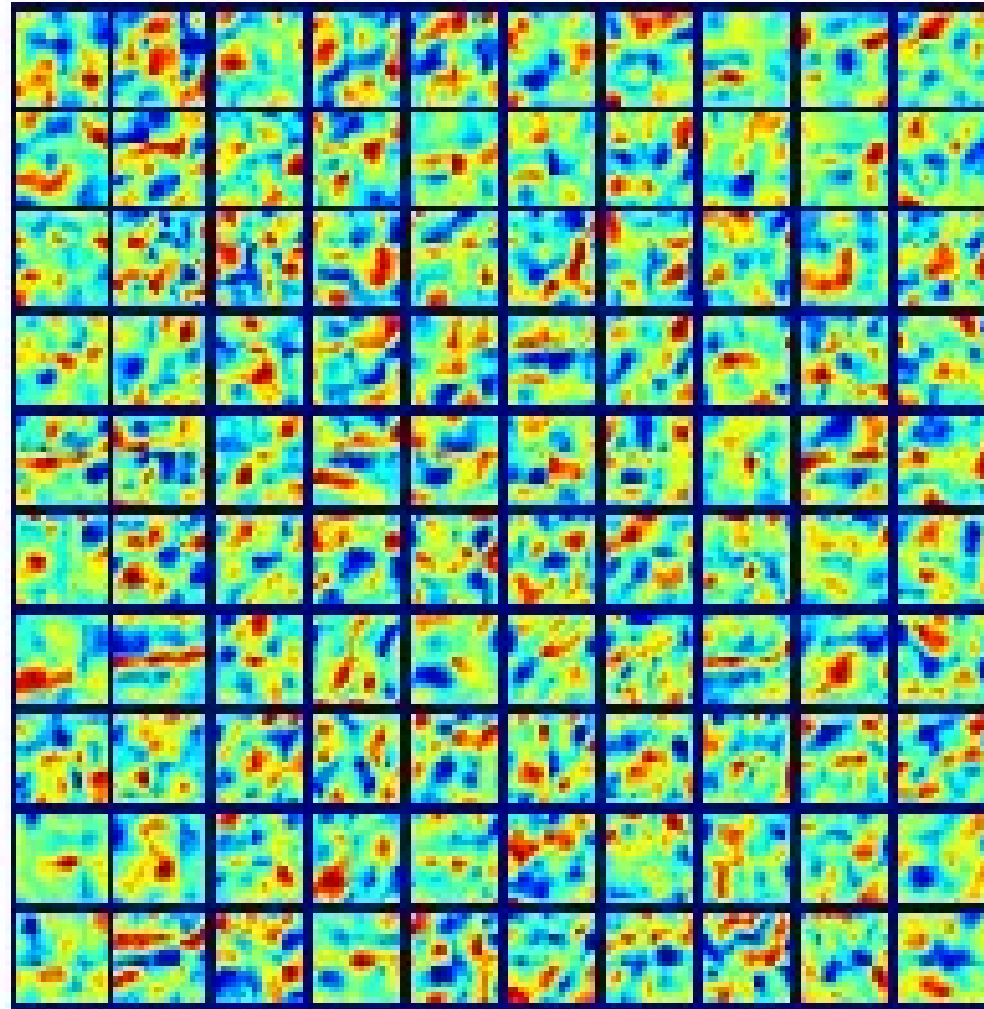Universität Jena
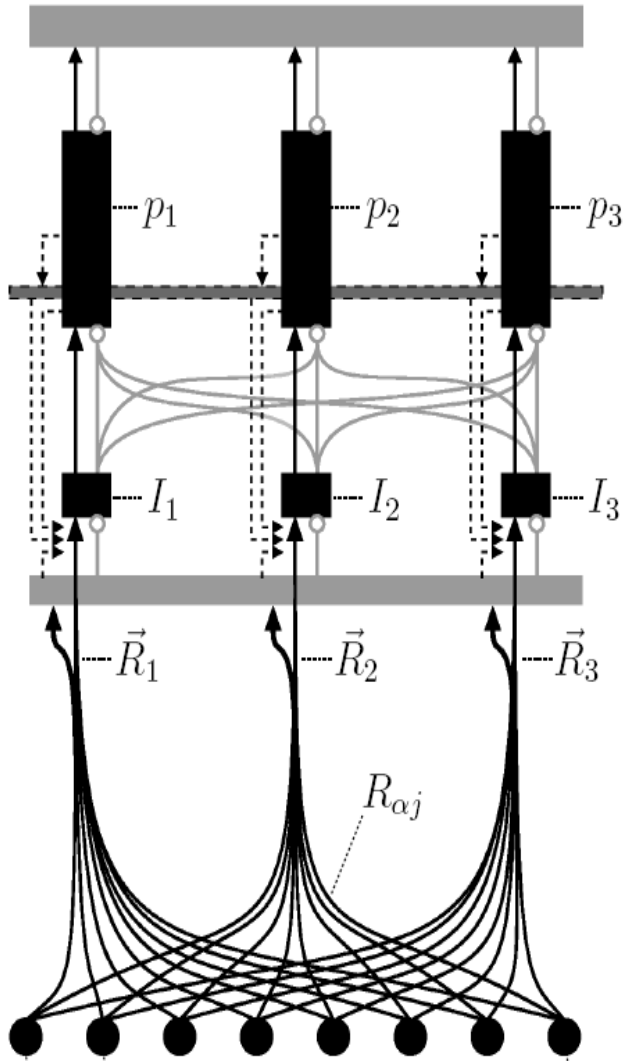
Jörg Lücke

# Thank you.

# Cortical Circuits



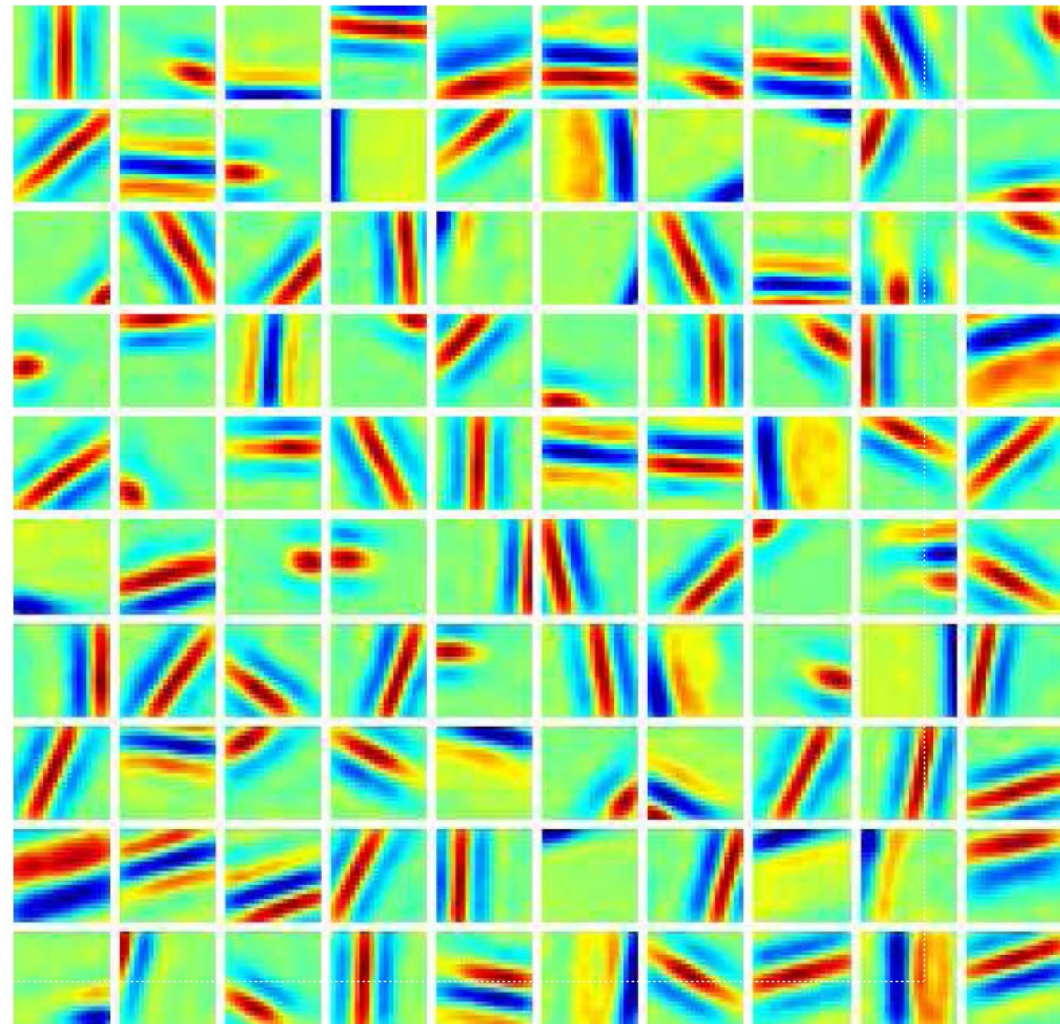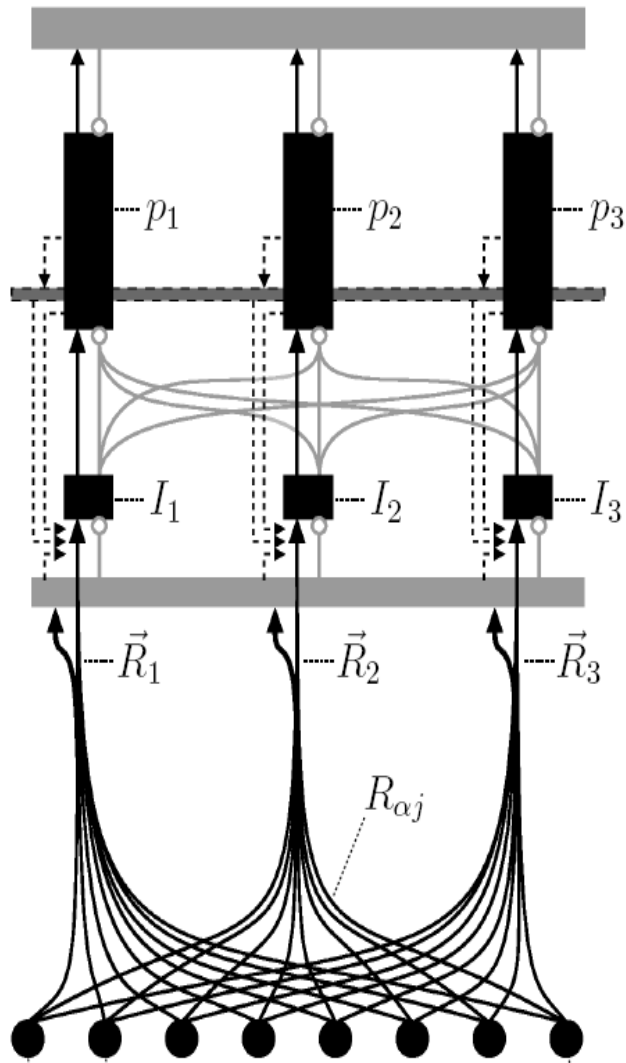Use Natural Image Patches as Input

# Cortical Circuits





Numerical simulations of stochastic and non-linear differential equations.

# Cortical Circuits
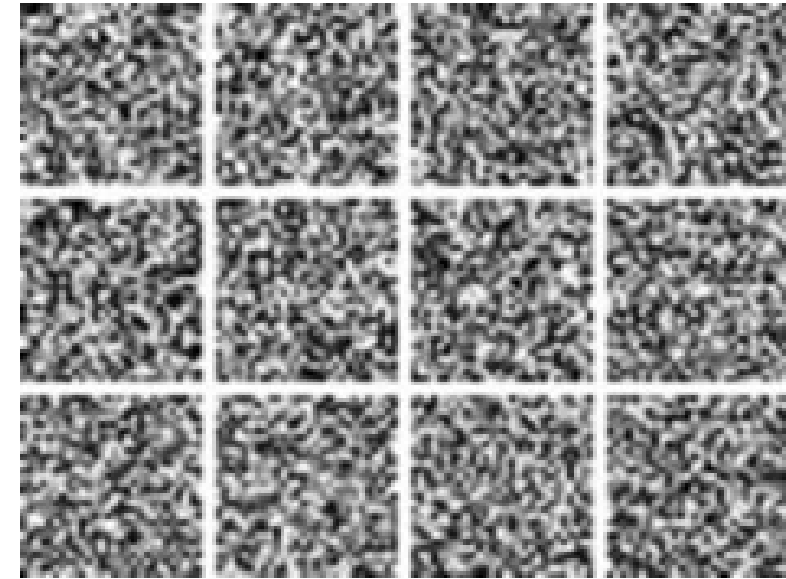


A  RFs for DoG images

Lücke, *Neural Computation*, 2009.

Jörg Lücke

# Application NMF

$$W \leftarrow W \odot \frac{\sum_n \vec{y}^{(n)} <\vec{s}^T>_{q_n}}{\sum_n W <\vec{s}\vec{s}^T>_{q_n}}$$

$$\left\langle g(\vec{s}, \Theta^{\text{old}}) \right\rangle_{q_n} \approx \frac{\sum_{\vec{s} \in \mathcal{K}_n} p(\vec{s}, \vec{y}^{(n)} \mid \Theta^{\text{old}}) \, g(\vec{s}, \Theta^{\text{old}})}{\sum_{\vec{s}' \in \mathcal{K}_n} p(\vec{s}', \vec{y}^{(n)} \mid \Theta^{\text{old}})}$$
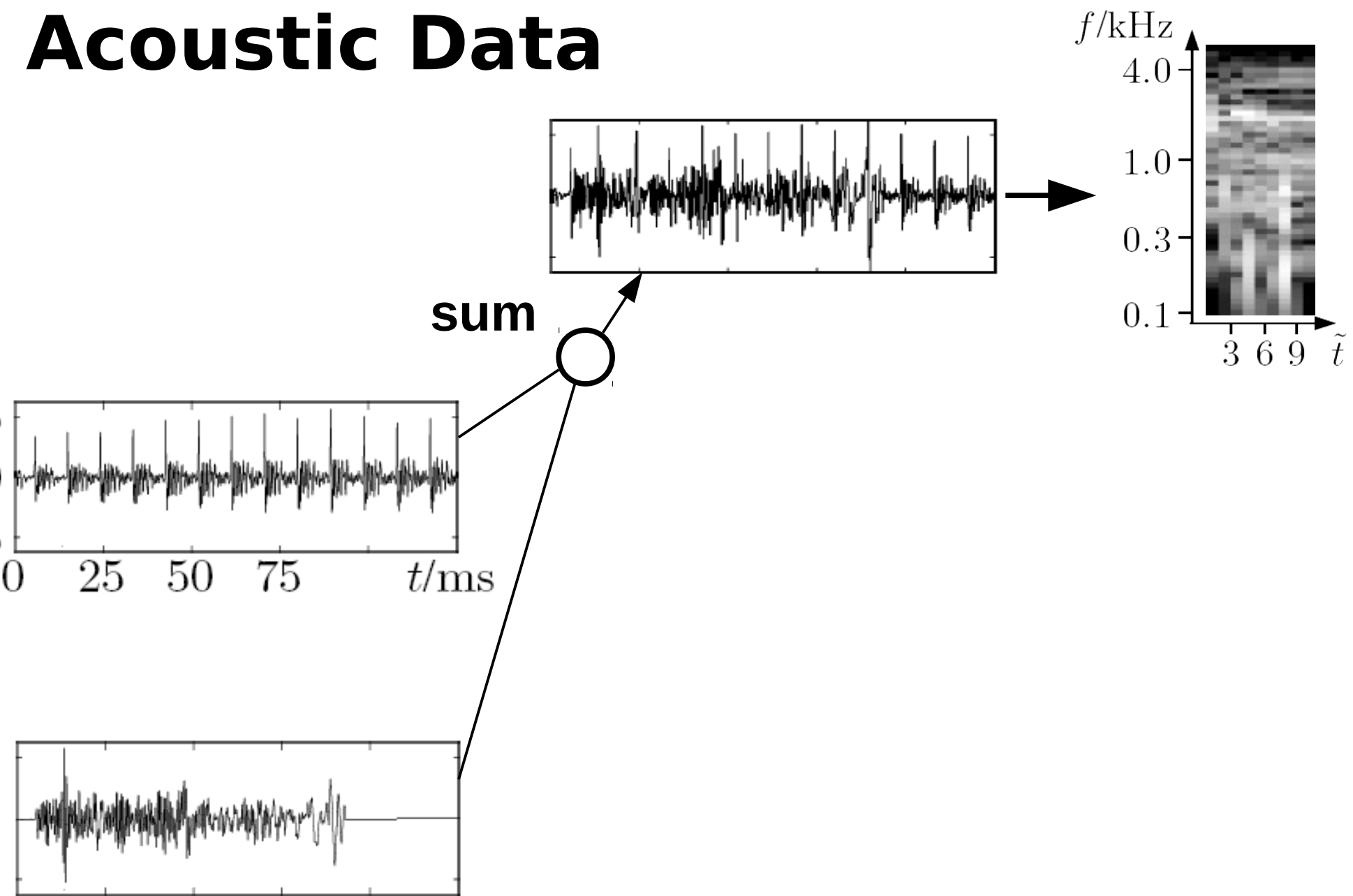


Application to MNIST data.

The model can be constraint by allowing only for positive $\vec{s}$ and positive $W$.
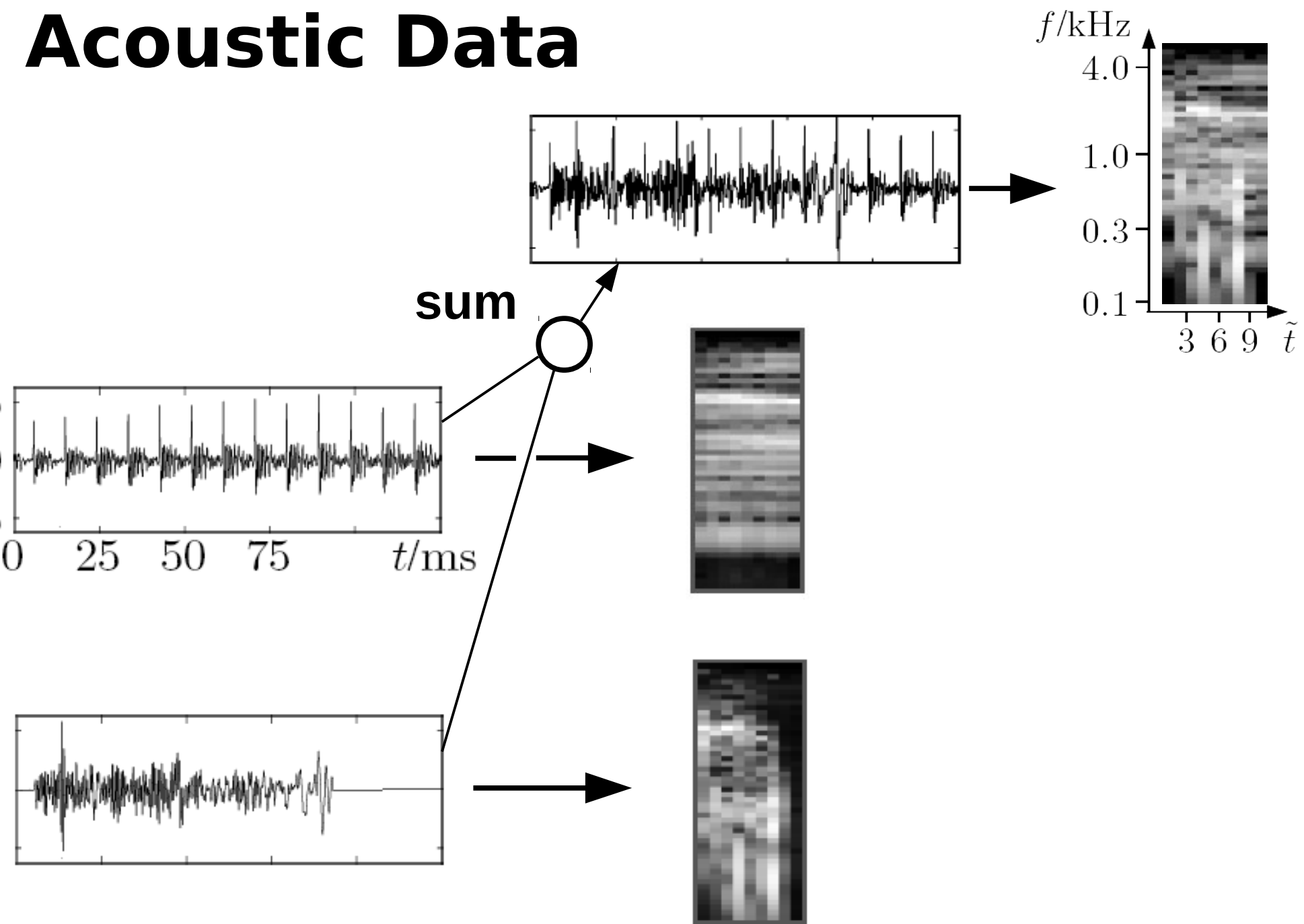
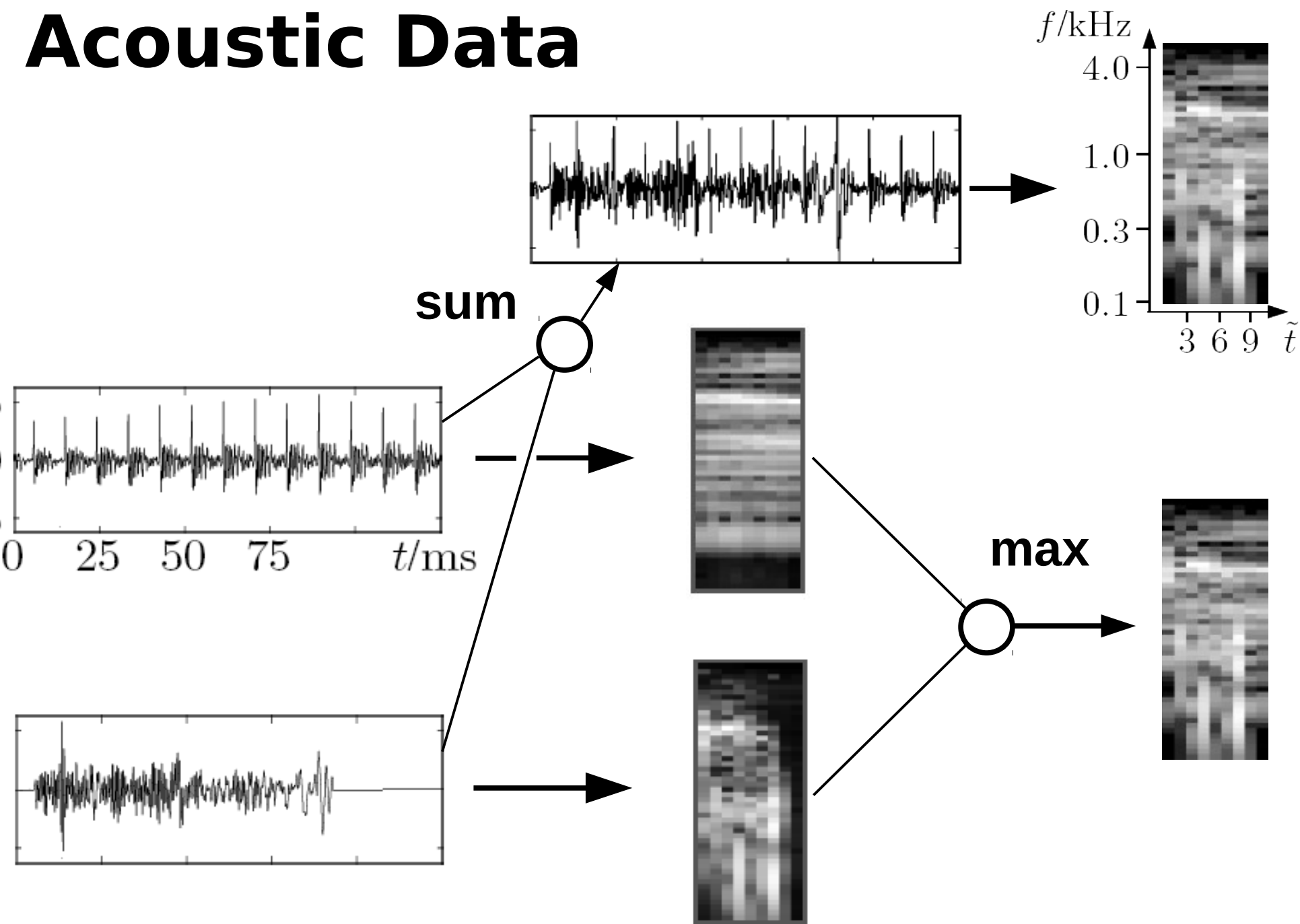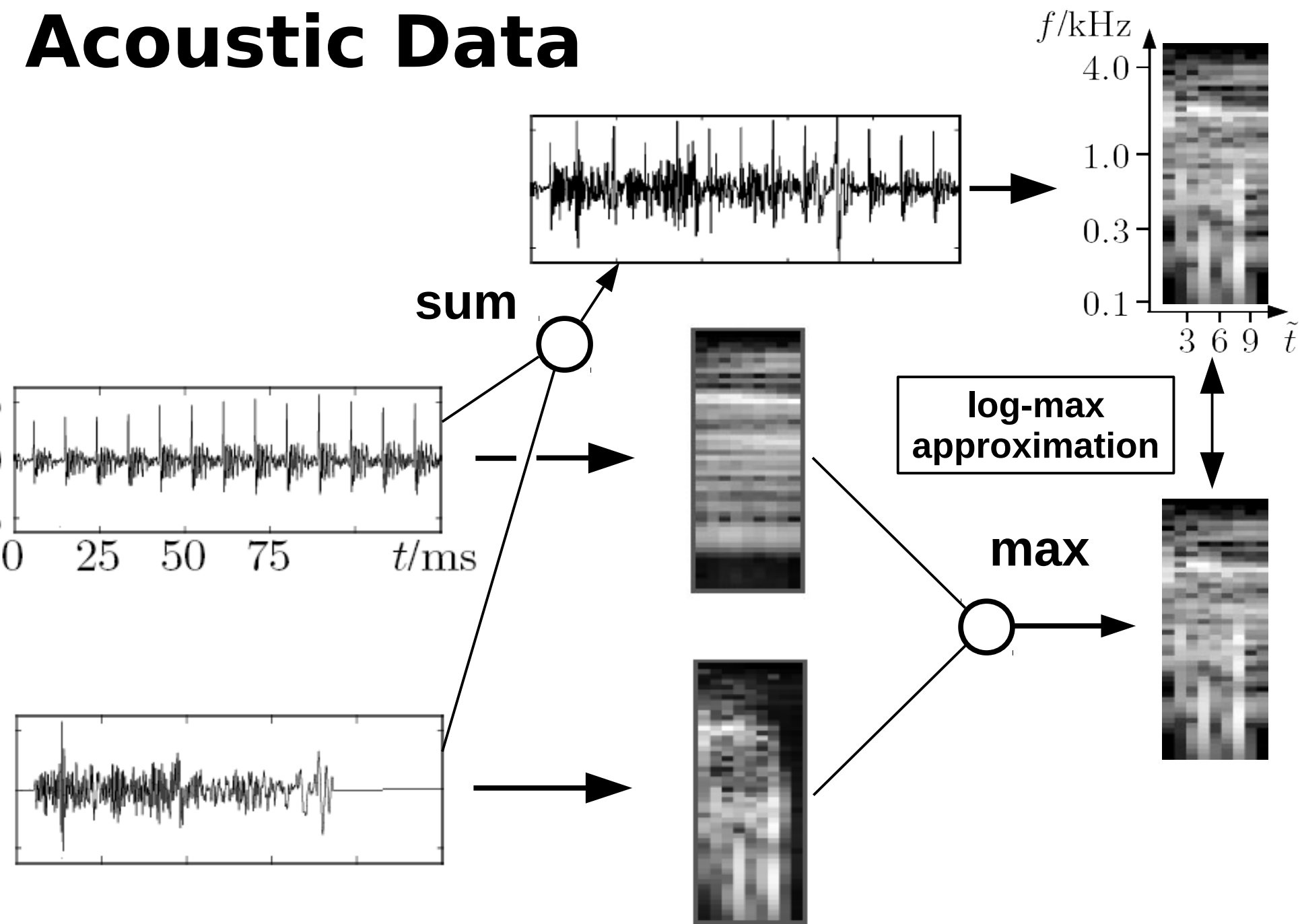We obtain a generative version of NMF (with binary latents).

Jörg Lücke

# Acoustic Data



**sum**

# Acoustic Data



**sum**

$f/\text{kHz}$

4.0

1.0

0.3

0.1

3  6  9  $\tilde{t}$

25  50  75  $t/\text{ms}$

CARL
VON
OSSIETZKY
*universität* OLDENBURG

Hearing
4all

# Acoustic Data



sum

max

$f/\mathrm{kHz}$

4.0

1.0

0.3

0.1

3 6 9 $\tilde{t}$

0   25   50   75   $t/\mathrm{ms}$

# Acoustic Data



sum

log-max
approximation
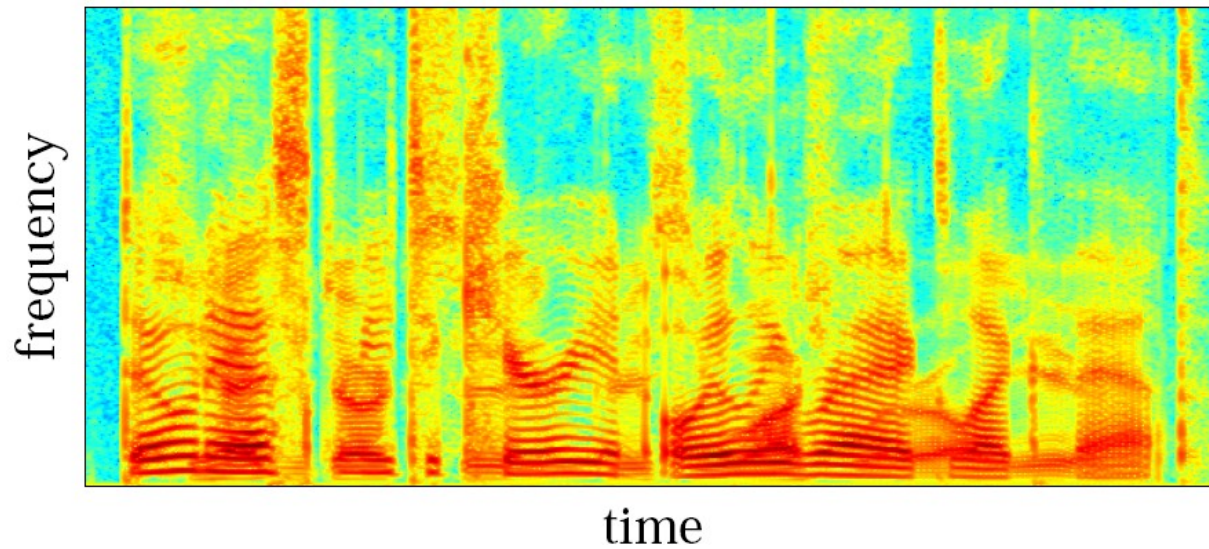
max

CARL
VON
OSSIETZKY
universität OLDENBURG

Hearing
4all

# Acoustic Data



Log-spectrogram
of a mixture of
two sound sources.



Max of the two individual
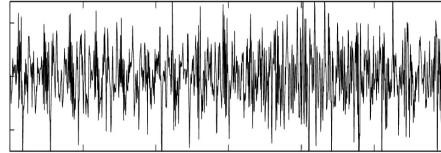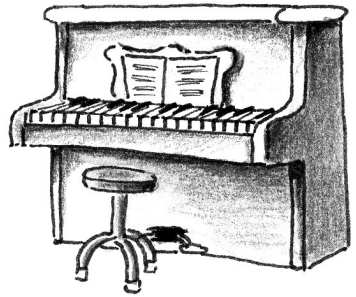log-spectrograms.

Source: Roweis, 2004

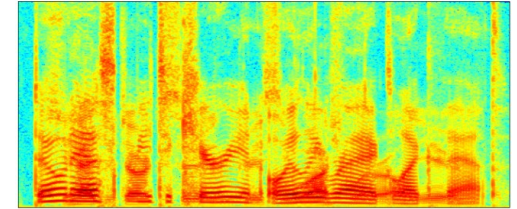Log-max approximation (Moore, 1983)

Jörg Lücke

# Acoustic Data

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

$$\vec{y} = f(s_{1:H} \, \vec{W}_{1:H}) + \vec{\eta}$$



$$\vec{y}$$

$$\vec{y}$$

Jörg Lücke
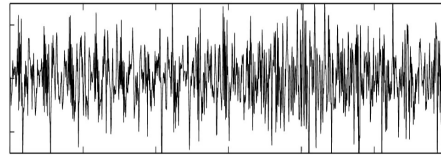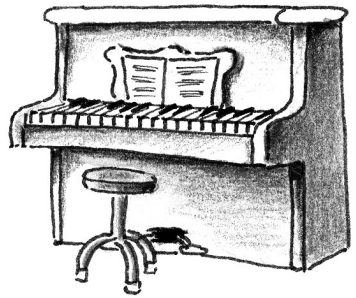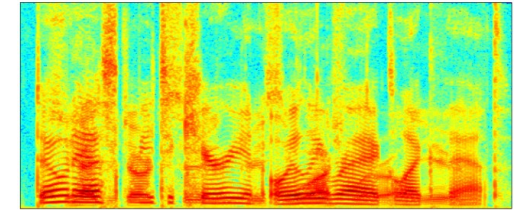
# Acoustic Data

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

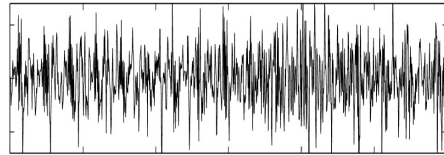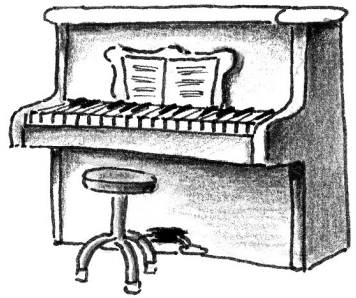$$\vec{y} = \max_h \{ s_h \vec{W}_h \} + \vec{\eta}$$

$\vec{y}$

$\vec{y}$

Bornschein et al., *PLOS CB* 2013
Shelton et al., *NIPS* 2012
Puertas, Bronschein, Lücke, *NIPS 2010*
Lücke, Sahani, *J Mach Learn Res 2008*
...
Roweis, *Eurospeech 2003*
Roweis, *NIPS 2002*
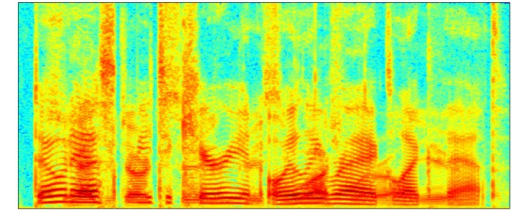Varga, Moore, *ICASSP 1990*

Jörg Lücke

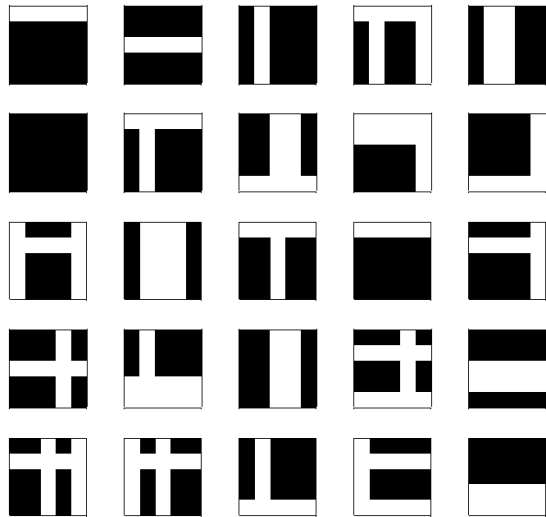# Acoustic Data

$$\vec{y} = \sum_h s_h \vec{W}_h + \vec{\eta}$$

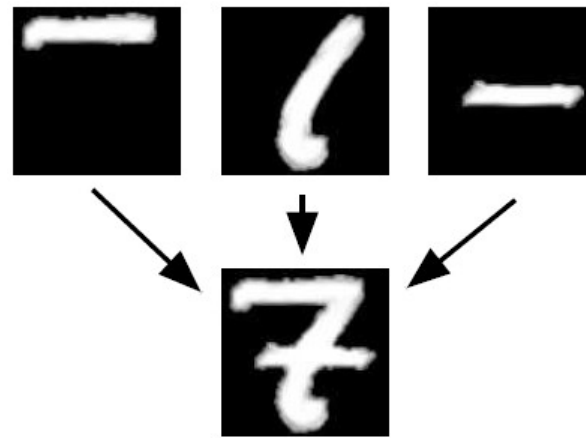$$\vec{y} = \max_h \{ s_h \vec{W}_h \} + \vec{\eta}$$



$\vec{y}$

$\vec{y}$



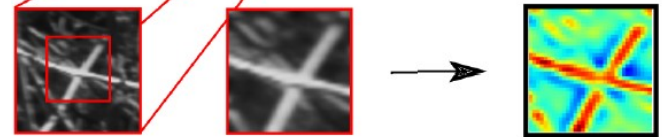The Bars Test.
Földiák, 1990

hand-written digits
(e.g., MNIST)

Natural image patches

**10s hidden**

**10s-100s hidden**

**100s-1000s hidden**

Jörg Lücke

CARL VON OSSIETZKY
universität OLDENBURG

Hearing 4all

# Other projects



Microscopy
Image Analysis

Deep Learning Architectures
for Pattern Recognition
Keck, Savin, Lücke
*PLOS Comp Bio*, 2012

weights

recurrent
stage

feed-forward



Jörg Lücke

# Other projects



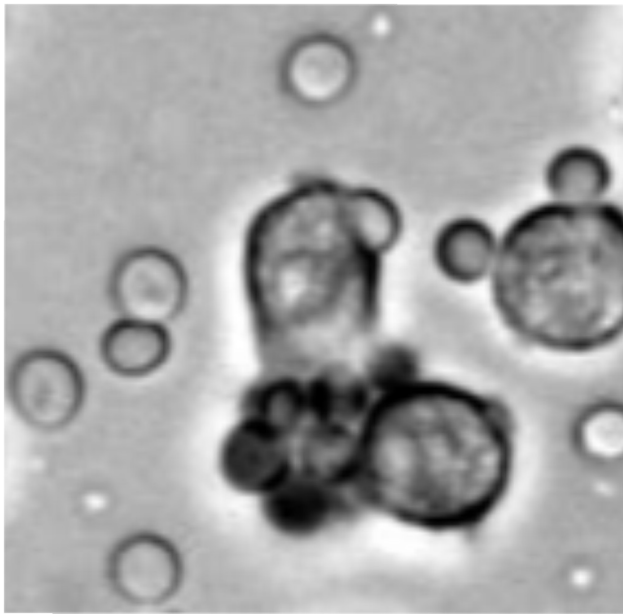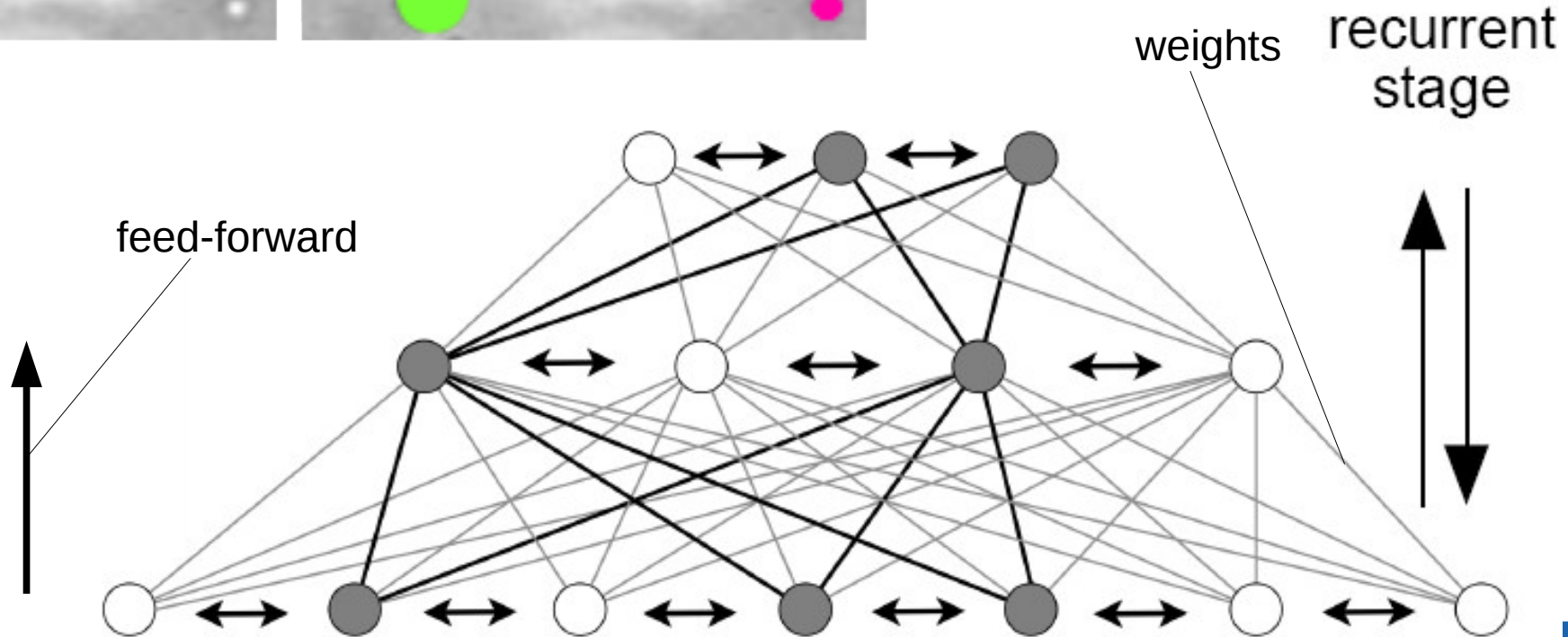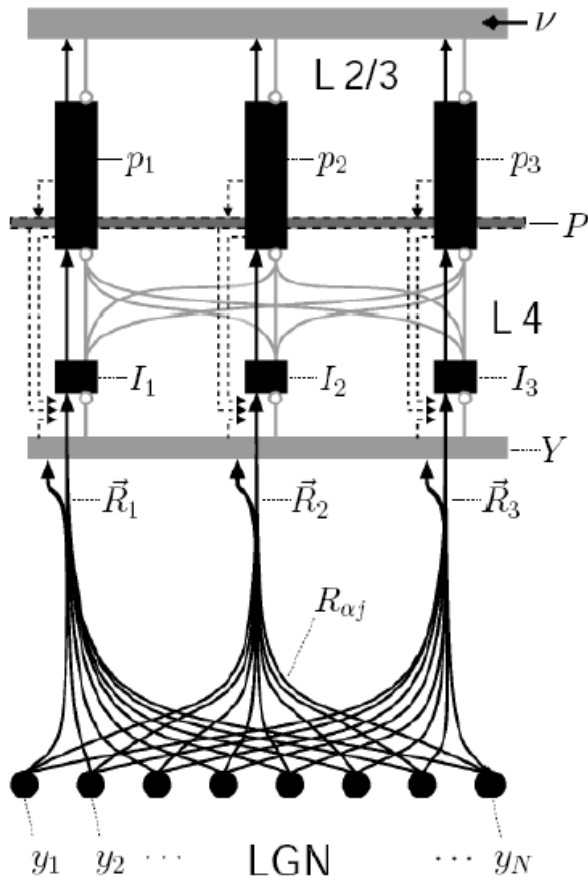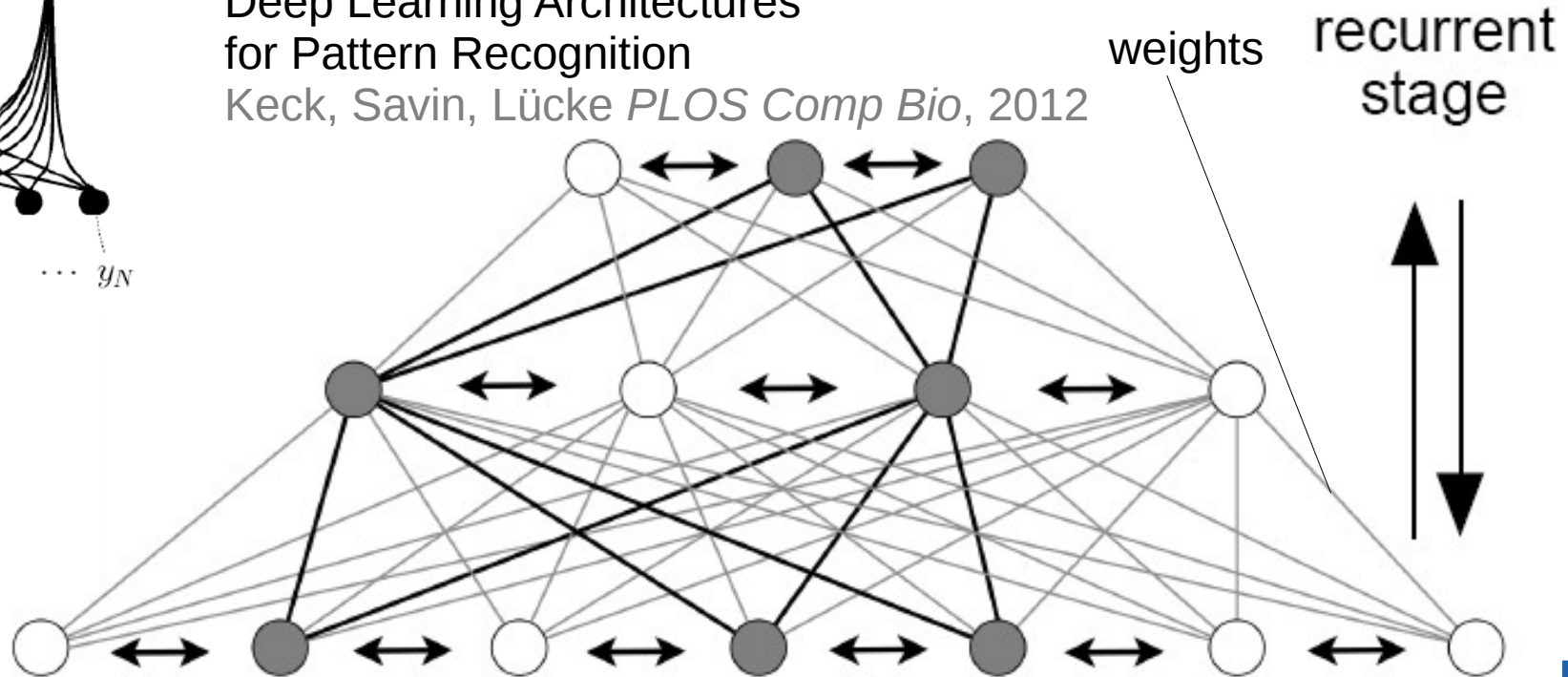Deep Learning Architectures
for Pattern Recognition
Lücke, *Neural Comp* 2009
Lücke, *ICANN* 2005-2007
Lücke, *Neural Networks* 2004
Lücke, Malsburg, *Neural Comp* 2004
...

Deep Learning Architectures
for Pattern Recognition
Keck, Savin, Lücke *PLOS Comp Bio*, 2012

weights

recurrent stage

Jörg Lücke

probability theory / applied mathematics / computer science

learning

generative models

approximate inference

signal processing; computer hearing
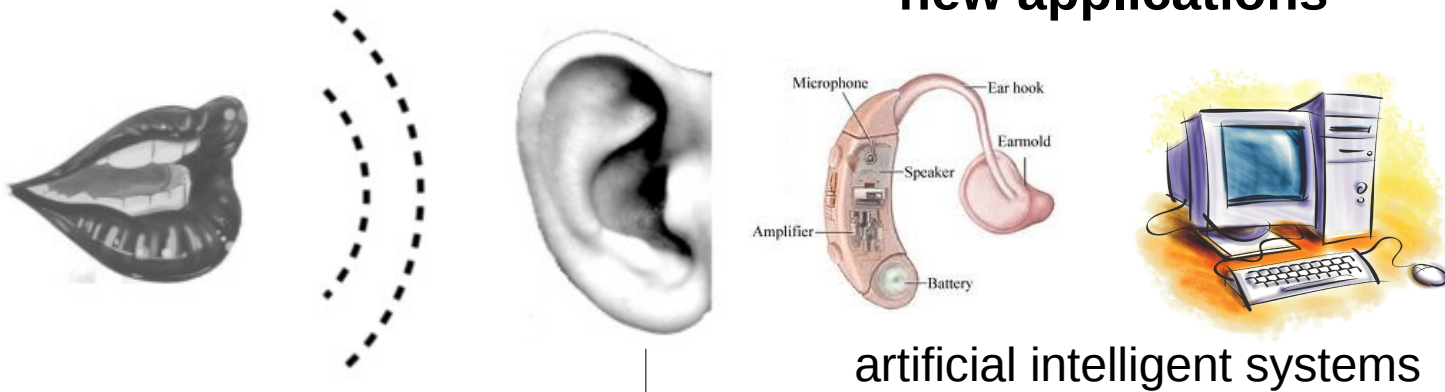
computational neuroscience

The World

The Learner

$$\vec{y} = \max_h \{s_h \, \vec{W}_h\} + \vec{\eta}$$

World State   Projection   Observed State   Inverse   Model State

Physics

Neuroscience

CARL VON OSSIETZKY universität OLDENBURG

Jörg Lücke

Hearing 4all

**new applications**

artificial intelligent systems

probability theory / applied mathematics / computer science

**learning**

generative models ⟷ approximate inference

signal processing; computer hearing

computational neuroscience

The World

The Learner

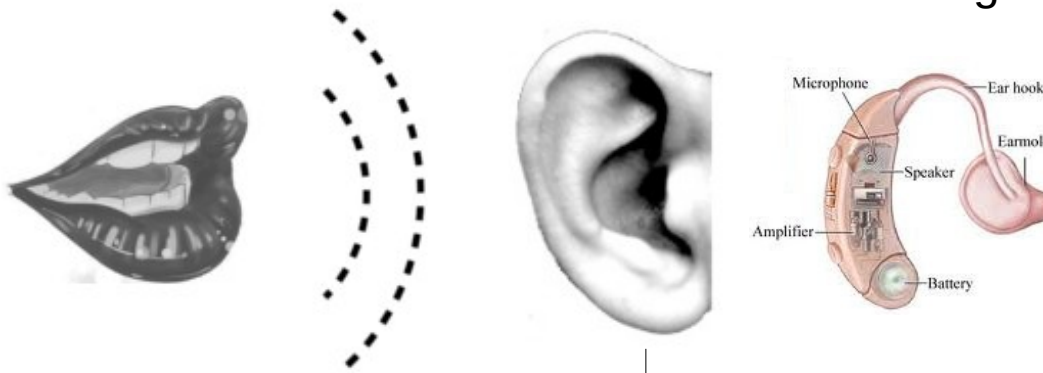$$\vec{y} = \max_h \{s_h \, \vec{W}_h\} + \vec{\eta}$$

algorithms for hearing instruments

personalizing hearing devices

HörTech

Hearing4all area B

probability theory / applied mathematics / computer science

learning

generative models

approximate inference

signal processing; computer hearing
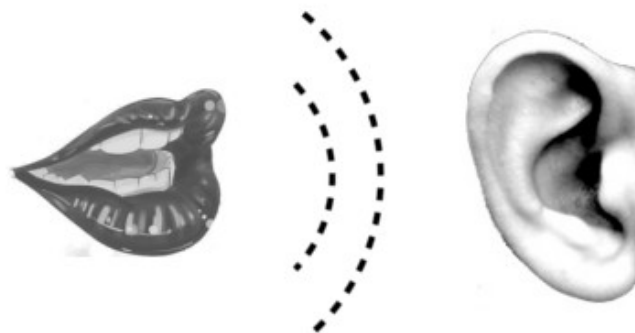
Individualisierte Hörakustik

Hearing4all area A

computational neuroscience

The World

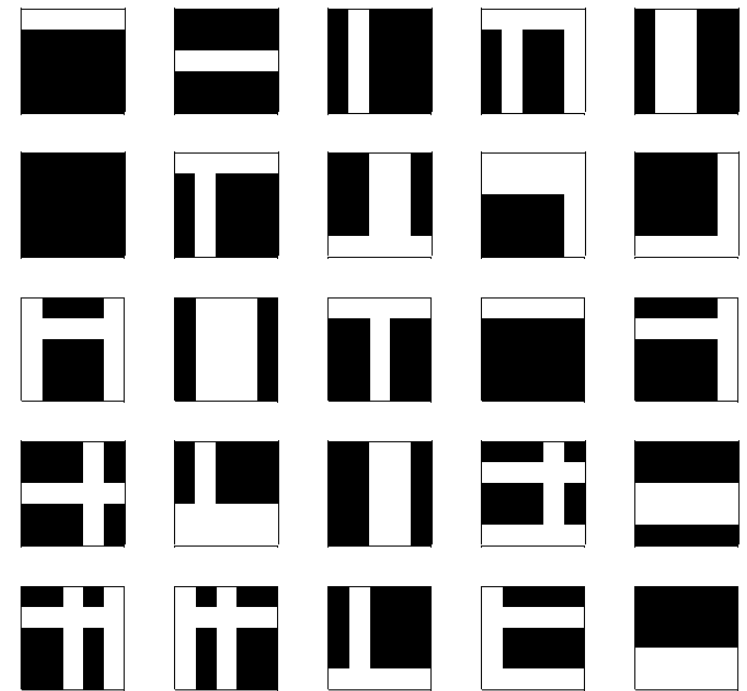The Learner

Zentrum für Neurosensorik

# Linear Causes

Linear generative models:

$$p(\vec{s}\,|\,\Theta) = \prod_h \cdots$$

$$p(\vec{y}\,|\,\vec{s},\Theta) = \mathcal{N}(\vec{y}; \textstyle\sum_h s_h^{(n)}\,\vec{W}_h, \sigma^2\,\mathbb{1})$$

Obtain basis functions:

The Bars Test, Földiák, 1990

**PCA**

**ICA 15**

Source:
Hochreiter &
Schmidhuber,
*Neural Comp*,
1999

Jörg Lücke