DTU



Lars Kai Hansen lkai@dtu.dk

## Sensing the deep structure of signals





What other computer scientists think I do



from theano import

DaveStmonds

What mathematicians think I do

What I think I do

What I actually do

# the brain design

Three design principles

i) Division of labor
Centers for vision, hearing, smell etc
ii) Neural networks of simple computers
iii) Learning – adaptivity -plasticity

But how does this structure represent / index the world, how does is rank importance? etc etc



# Outline



What is deep structure?

Cognitive components and attention modeling Ecology of audio signals

Is structure is determined by the environment: Statistics/ physics / mechanisms?

Uniqueness of perception in the brain & Uniqueness in deep neural networks

What about higher order cognition, social cognition?

# Attention & human optimality

"... the withdrawal from some things in order to deal effectively with others" William James (1890)

"... <u>To behave adaptively in a complex world, an</u> <u>animal must select, from the wealth of information</u> <u>available to it, the information that is most relevant</u> <u>at any point in time</u>. This information is then evaluated in working memory, where it can be analyzed in detail, decisions about that information can be made, and plans for action can be elaborated. The mechanisms of attention are responsible for selecting the information that gains access to working memory."

Eric I. Knudsen (2007)

things...?

SPLINE 2016





W. James, *The Principles of Psychology, Vol. 1, Dover Publications,* 1880/1950.
E.I. Knudsen, "Fundamental Components of Attention," *Annual Review of Neuroscience, vol. 30, no. 1, pp. 57–78, 2007.*

## Deep structure needed to predict the future



Processing in the brain is based on extremely well-informed / optimized representations and mechanisms –

A key issue is selective attention...

Fundamental questionWhat can you attend to?or...What is an object / chunkof information?

# Cognitive component analysis ...what we can attend to

- The object / chunk is a key notion in cognitive psychology
  - ...number of objects in short time memory, objects "race to short term memory"
  - Miller, G.A. (1956), The Magical Number Seven, Plus or Minus Two: Some Limits on our Capacity for Processing Information. Psychological Review, 63, 81-97
  - Bundesen, C., Habekost, T. and Kyllingsbæk, S., 2005. A neural theory of visual attention: bridging cognition and neurophysiology. Psychological review, 112(2), p.291).
  - Miller: "...we are not very definite about what constitutes a chunk of information."
  - A pragmatic definition of an object could be: An object is a signal source with independent behavior in a given environment (...imagined?)
- Theoretical issues: The relation between supervised and un-supervised learning. Related to the discussion of the utility of unlabeled examples in supervised learning and fast/one sample learning...

Practical Issues: Can we predict which digital media components a user will pay attention to? -a key challenge for cognitive systems.





# Vector space representation

- Abstract representation can be used for all digital media
- A "cognitive event" is represented as a point in a high-dimensional "feature space" – document similarity ~ spatial proximity in a given metric
- Text: Term/keyword histogram, N-grams
- Image: Color histogram, texture measures
- Video: Object coordinates (tracking), active appearance models
- Sound: Spectral coefficients, mel cepstral coefficients, gamma tone filters

Contexts can be identified by their feature associations ( = Latent semantics )

S. Deerwester et al. *Indexing by latent semantic analysis*. Journal of the American Society for Information Science, 41(6), 391-407, (1990)

## The independent component hypothesis

Challenge: Presence of multiple agents/contexts
 Need to "blindly" separate source signals = learn contexts
 ICA, NMF, tensor factorization provides (almost) unique solutions to...



# Linear mixing generative model ICA - "Synthesis" simplistic model incorporating sparsity and independence



# Protocol for comparing supervised and unsupervised learning

- Use the "unsupervised-then-supervised" scheme to implement a classifier:
  - Train the unsupervised scheme, eg., ICA
  - Freeze the ICA representation (A matrix)
  - Train a simple (e.g. Naïve Bayes) classifier using the features obtained in unsupervised learning Use
- Compare with supervised classifier
  - Error rates of the two systems
  - Compare posterior probabilities

DTU

#### **Phoneme classification**

#### Nasal vs oral: "Esprit project ROARS" (Alinat et al., 1993)



Binary classification

Error rates: 0.23 (sup.), 0.22 (unsup.) Bitrates: 0.48 (sup.), 0.39 (unsup.)

# Cognitive components of speech

- Basic representation: Mel weigthed cepstral coefficients (MFCCs)
- Modeling at different time scales 20 msec – 1000 msec



- Phonemes
- Gender
- Speaker identity





Figure 3: The latent space is formed by the two first principal components of data consisting of four separate utterances representing the sounds 's', 'o', 'f', 'a'. The structure clearly shows the sparse component mixture, with 'rays' emanating from the origin (0,0). The ray embraced in a rectangle contains a mixture of 's' and 'f' features, a cognitive component associated with the vowel /e/ sound.

TRAINING DATA Mel weighted cepstral coeff. (MFCC) 10 12 14 700 400 300 600 TEST DATA 4 0 10 12 14 16 100 **S** Α CLIPPED CEPSTRALS: |z| > 1.7 °° 2000 € 0.2 0.1 Ο # 146# C -0.1 4 r \*\*\* -0.2 8 [a] PHONEME IN 'S' AND 'F' •S 💊 -0.3 P --6 8 -0.4 0 0.2 0.4 PC1 0.6 0.8 -0.2 0 1

SPLINE 2016

#### Error rate comparison

For the given time scales and thresholds, data locate around y = x, and the correlation coefficient  $\rho = 0.67$ , p < 1.38e - 09.





#### Sample-to-sample correlation

Three groups: vowels eh, ow;
fricatives s, z, f, v; and stops k, g, p, t.
25-d MFCCs; EBS to keep 99%

energy; PCA reduces dimension to 6.

- Two models had a similar pattern of making correct predictions and mistakes, and the percentage of matching between supervised and unsupervised learning was 91%.

DTU

## Longer time scales



Time integrated (1000ms) MFCC's: text independent speaker recognition....

Feng & Hansen (CIMCA, 2005)

Error rate correlations for super/unsupervised learning for different cognitive time scales and events

Challenged by degree of sparsity and time averaging

Gender, Identity, Heigth etc are the Audio Gist vars



**Fig. 4**. Figure shows test error rates of both supervised and unsupervised learning on four topics: phonemes, gender, height and identity. Solid lines indicate y = x in the coordinate systems. All data located along this line, meaning high correlation between supervised and unsupervised learning.



doi:10.1038/nature11020

SPLINE 2016

Selective cortical representation of attended speaker in multi-talker speech perception

Nima Mesgarani<sup>1</sup> & Edward F. Chang<sup>1</sup>

# Human brain mechanisms Attention in speech mixtures



Figure 1 | Acoustic and neural reconstructed spectrograms for speech from a single speaker or a mixture of speakers. a, b, Example acoustic waveform and auditory spectrograms of speaker one (male; a) and speaker two (female; **b**). **c**, Waveform and spectrogram of the mixture of the two shows highly overlapping energy distributions. d, Difference spectrogram highlights the mixture regions where speaker one (blue) or two (red) has more acoustic energy. e, f, Neuralpopulation-based stimulus reconstruction of speaker one (e) and speaker two (f) alone shows similar spectrotemporal features as the original spectrograms in **a** and **b**. **g**, **h**, The reconstructed spectrograms from the same mixture sound when attending to either speaker one (g) or two (**h**) highly resemble the single speaker reconstructions, shown in e and f, respectively.

i, Overlay of the a maximum energy spectrograms in a

P value  $10 \times 10^{-6}$  $10 \times 10^{-2}$ 

# Uniqueness of representations?

Modern society's deep specialization requires efficient shared representations

You know what I mean - right?

Does machine learning also develop shared representations and if so - why?







JP Dmochowski, P Sajda, J Dias, LC Parra, "Correlated components of ongoing EEG point to emotionally laden attention -a possible marker of engagement?" Frontiers of Human Neuroscience, 6:112, 2012. JP Dmochowski, MA. Bezdek, BP. Abelson, JS. Johnson, EH Schumacher, LC Parra, "Audience preferences are predicted by temporal reliability of neural processing", Nature Communications 5:4567, 2014. AT Poulsen, S, Kamronn, J Dmochowski, LC Parra, LK, Hansen: . "Measuring engagement in a classroom: Synchronised neural recordings during a video presentation". *arXiv preprint arXiv:1604.03019 (2016)*.

## What is the joint attention signal?

Driven by early visual response hich is modulated by attention...







#### Real-time feasible in (sub)-groups, correlate with computed saliency...

Hiliard et al. Sensory gain control (amplification) as a mechanism of selective attention Phil.Trans. R. Soc. Lond. B (1998) 353, 1257^1270

DTU





Important for engineering proxies for human information processing... Cf. efficient coding of "context-to-action" mapping

DTU

# Deep networks



#### Lars Kai Hansen, DTU Compute

## **Reducing the Dimensionality of Data with Neural Networks**

G. E. Hinton\* and R. R. Salakhutdinov





Fig. 1. Pretraining consists of learning a stack of restricted Boltzmann machines (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 pretraining, the RBMs are "unrolled" to create a deep autoencoder, which is then fine-tuned using backpropagation of error derivatives.

Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313, no. 5786 (2006): 504-507.

#### SPLINE 2016

How well determined are the representations by ecology... sensitivity analysis



(a) Error increase versus strength of damage (100 networks).

RJ Aaskov, LK Hansen "On the resilience of deep neural networks to weight damage." *In review* (2016)

#### On the resilience of deep neural networks to weight damage

Rasmus Jessen Aaskov Dept. of Applied Mathematics and Computer Science Technical University of Denmark Lars Kai Hansen Dept. of Applied Mathematics and Computer Science Technical University of Denmark





(a) Example of a well performing network (20.25% final test error rate). (b) Example of a network with poor performance (23.71% final test error rate).

Figure 3. Two examples of how our method for estimating loss behaves in different cases. In both figures, the average loss of a repeated damage experiment and our estimated expected loss is shown for each of the individual layers. The expected loss is illustrated as a dotted line with distinctive markers for each layer.

# CONVERGENT LEARNING: DO DIFFERENT NEURAL NETWORKS LEARN THE SAME REPRESENTATIONS?

Yixuan Li<sup>1</sup>\*, Jason Yosinski<sup>1</sup>\*, Jeff Clune<sup>2</sup>, Hod Lipson<sup>3</sup>, & John Hopcroft<sup>1</sup>



Figure 1: Correlation matrices for the conv1 layer, displayed as images with minimum value at black and maximum at white. (**a,b**) Within-net correlation matrices for Net1 and Net2, respectively. (**c**) Between-net correlation for Net1 vs. Net2. (**d**) Between-net correlation for Net1 vs. a version of Net2 that has been permuted to approximate Net1's feature order. The partially white diagonal of this final matrix shows the extent to which the alignment is successful; see Figure 3 for a plot of the Uses along this diagonal and further discussion.



L, Yixuan, J Yosinski, J Clune, H Lipson, J Hopcroft. "Convergent Learning: Do different neural networks learn the same representations?." *arXiv preprint arXiv:1511.07543* (2015)

DTU

# What about "higher order" cognition?



#### Lars Kai Hansen, DTU Compute

# Independent contexts in document databases

• x(j,t) is the occurence of the j'th word in the t'th document.

• s(k,t) quantifies how much the k'th context is expressed in t'th document.

• A(j,k) quantifies the typical importance of the j'th word in the k'th context

Data Stream Word histograms (Term / doc matrix) Data extraction Filter Normalize Modeling PCA ICA Classification Group topics Time flow Keywords Analysis chat join pm cnn board message allpolitics visit check america ... susan smith mother children kid life .... people census elian state clinton government good father ...

ICA in text Isbell and Viola (1999) Kolenda, Hansen, (2000)

DTU

## PCA vs ICA document scatterplots

Terms	Documenta								
	el	2	c3	c4	ch	ml	m2	m3	<b>m</b> 4
computer	1	1	0		0	0	0	0	0
EPS	10	0	1	1	9	0	8	8	
numen.	24 S	0 A	÷.	1	2	0	10	90 A	4
Interisce	A	- <u>N</u>	÷.		10	. V.	N A	9	N A
response		÷.	÷.	14	14 (A)	- 10 - A	2		
ayatem	0	- 20	÷.	2	- 12	0	0		0
Line	10 10	- 22	÷.	1		- N	10	8	4
mark	0	6	6		0	1	÷.	1	10
mincea	0	6	<u> </u>	0	0	o i	in in	- Ç.	
STATES ST	0	1			6	0	0	õ	1
trees	0	õ	0	0	ö	1	- S	1	ô
Law a		17					-		

#### Lars Kai Hansen, DTU Compute

## Linear mixture of independent agents in term-document scatterplots





Linear mixture of independent contexts observed in short time features (mel-ceptrum) in a music database.

SPLINE 2016

## Social networks: Linear mixtures of independent communities?



Genre patterns in expert's opinion on similar music artists

(AMG400, Courtesy D. Ellis)

"Movie actor network" - A collaborative small world network 128.000 movies 380.000 actors



## Independent contexts in multi-media

- Organizing webpages in categories
   Labels obtained from
- Labels obtained from Yahoo's directory
- Features: Text, color, and texture subsets of MPEG image features



Feature / document matrix

L.K. Hansen, J. Larsen and T. Kolenda "On Independent Component Analysis for Multimedia Signals". In L. Guan et al.: *Multimedia Image and Video Processing*, CRC Press, Ch. 7, pp. 175-199, 2000.



Performance of the system trained by associating unsupervised independent components with labels – generalization based on Yahoo cathegories

Modality	Classification Error
Color	23.0%
Texture	18.0%
Texture/Color	11.5%
Text	5.7%
Combined (texture/color/text)	2.8%



Fig. 3. Scatterplots of the text and image multimedia data, projected to a two-dimensional subspace found by PCA. Grey value of points corresponds to the three classes considered, see Fig. 4. The ray like structure strongly suggest an ICA interpretation, however, the relevance of this representation can only be determined by a subsequent inspection of the recovered source signals. As we will see in section 4.6, it turns out that there is an interesting alignment of the source signals and a manual labeling of the multimedia documents.



#### Lars Kai Hansen, DTU Compute





4.5

89.75

9



Combined errorrate: 2.8% Single best errorrate: 5.7%





#### CASTSEARCH - CONTEXT BASED SPEECH DOCUMENT RETRIEVAL

Lasse Lohilahti Mølgaard, Kasper Winther Jørgensen, and Lars Kai Hansen

Informatics and Mathematical Modelling Technical University of Denmark Richard Petersens Plads Building 321, DK-2800 Kongens Lyngby, Denmark



**Fig. 1**. The system setup. The audio stream is first processed using audio segmentation. Segments are then using an automatic speech recognition (ASR) system to produce text segments. The text is then processed using a vector representation of text and apply non-negative matrix factorization (NMF) to find a topic space.



Fig. 3. Figure 3(a) shows the manual segmentation of the news show into 7 classes. Figure 3(b) shows the distribution  $p(k|d^*)$  used to do the actual segmentation shown in figure 3(c). The NMF-segmentation is in general consistent with the manual segmentation. Though, the segment that is manually segmented as 'crime' is labeled 'other' by the NMF-segmentation

Mølgaard et al. 2007



## castsearch.imm.dtu.dk



A. Meng, P. Ahrendt, J. Larsen, L.K. Hansen: *Temporal Feature Integration for Music Genre Classification*. *IEEE* Transactions on Audio and Speech and Language Processing 15(5): 1654-1664 (2007)

T. Lehn-Schiøler, J. Arenas-García, K.B. Petersen and L.K. Hansen: A Genre Classification Plug-in for Data Collection. Proc. 7th Intl. Conf. on Music Information Retrieval, ISMIR 2006, pp. 320-321, Victoria, Canada, Oct. (2006).

L.K. Hansen, T. Lehn-Schiøler, K.B. Petersen, J. Arenas-Garcia, J. Larsen, and S.H. Jensen: *Learning and clean-up in a large music database. EUSIPCO 2007, European Conference on Signal Processing, Poznan (2007).* 

\_ 8 ×

🟠 🕶 🦉

#### muzeeker

Wikipedia based common sense Wikipedia used as a proxy for the music users mental model Implementation: Filter retrieval using

Wikipedia's article/ categories

muzeeker.com



S. Halling, M.K. Sigurdsson, J.E. Larsen, S. Knudsen, L.K. Hansen: MuZeeker: A domain SpecificWikipedia-based Search Engine. In Proc. First International Workshop on Mobile Multimedia Processing. Tampa, USA (2008).

🖉 MuZeeker Search#top - Windows Internet Explorer

J.E. Larsen, S. Halling, M. Sigurdsson and L.K. Hansen: MuZeeker - Adapting a music search engine for mobile phones. To appear in Springer Lecture Notes in Computer Science 'Mobile Multimedia Processing: Fundamentals, Methods, and Applications', Selected papers from First International Workshop on Mobile Multimedia Processing, Tampa, USA. (2010).

# So representations are optimal, what about attention?



#### Top down vs bottom up attention

## Bottom up

Attention determined by feature of the input

#### Audio

Cocktail party effect (Pollack+ Pickett, 57)

#### Visual

Classical spatial novelty saliency (Itti+Koch, 04)

## Top down

Attention determined by state of the observer

Audio

Cocktail party problem (Cherry, 64)

Visual

ambiguous pictures eye tracking

See e.g. J.M. Wolfe et al. "How fast can you change your mind? The speed of top-down guidance in visual search" Vision Research 44 (2004) 1411–1426



## Visual system hierarchy



DJ Felleman, DC. Van Essen. "Distributed hierarchical processing in the primate cerebral cortex." *Cerebral cortex* 1, no. 1 (1991): 1-47.







**Fig. 3.** Convergence and divergence in visual processing. Arrows represent major lines of information flow from subcortical P and M streams (bottom) to the selectivities represented among neurons at early stages of cortical analysis (middle) and from there to two general tasks of vision (top level). The hatched portion of the M cell curve represents their nonlinear component of processing. The processing streams associated with each property in the middle row are assigned on the basis of a high incidence of selectivity recorded physiologically (6, 7).



Itti, Dhavale & Pighin, Proc. SPIE, 2003

# ML model of "optimal" top down attention

**1**. The <u>task</u> is implemented as decision problem

**2.** Attention is represented as the <u>choice</u> over set of detailed features

## Standard probabilistic classifier

Model of posterior probability

Two sets of features

i) Features setting the context'the gist' (x)

(Friedman, 79; Torralba et al., 04)

ii) Potential features (z)considered by the attentionmechanism

A. Friedman: Framing pictures: the role of knowledge in automatized encoding and memory of gist. Journal of Experimental Psychology: General 1979;108:316–355.

2

5

Focur

# Mathematical model

We are interested in a partial observation **x** under a decision task: Choose among "C" actions

$$p(c|\mathbf{x}) = \int p(c, \mathbf{z}|\mathbf{x}) d\mathbf{z}$$
$$= \frac{\int p(c, \mathbf{x}, \mathbf{z}) d\mathbf{z}}{\sum_{c=1}^{C} \int p(c, \mathbf{x}, \mathbf{z}) d\mathbf{z}}$$

# versus getting additional information though z

$$p(c|\boldsymbol{x}, z_j) = \sum_{c=1}^{C} \int p(c, \boldsymbol{z} | \boldsymbol{x}) \prod_{i \neq j} dz_i$$
$$= \frac{\int p(c, \boldsymbol{x}, \boldsymbol{z}) \prod_{i \neq j} dz_i}{\sum_{c=1}^{C} \int p(c, \boldsymbol{x}, \boldsymbol{z}) \prod_{i \neq j} dz_i}$$

# Measure the information gain

First used by Lindley (1956) for experimental design..

#### ON A MEASURE OF THE INFORMATION PROVIDED BY AN EXPERIMENT<sup>1, 2</sup>

#### BY D. V. LINDLEY

University of Cambridge and University of Chicago

1. Summary. A measure is introduced of the information provided by an experiment. The measure is derived from the work of Shannon [10] and involves the knowledge prior to performing the experiment, expressed through a prior probability distribution over the parameter space. The measure is used to compare some pairs of experiments without reference to prior distributions; this method of comparison is contrasted with the methods discussed by Black-

or experimental design, where the object of experimentation is not to reach decisions but rather to gain knowledge about the world.

D. V. Lindley, "On a measure of the information provided by an experiment," Annals Mathematical Statistics, vol. 4, pp. 986–1005, 1956.

# Information theoretical model

$$\Delta S_j(\boldsymbol{x}, z_j) = -\sum_{c=1}^C \int \log p(c, \boldsymbol{z} | \boldsymbol{x}) p(c, \boldsymbol{z} | \boldsymbol{x}) d\boldsymbol{z} + \sum_{c=1}^C \log p(c | \boldsymbol{x}, z_j) p(c | \boldsymbol{x}, z_j) G_j(\boldsymbol{x}) \equiv \int \Delta S_j(\boldsymbol{x}, z_j) p(z_j | \boldsymbol{x}) dz_j = \sum_{c=1}^C \int \log p(c | \boldsymbol{x}, z_j) p(c, z_j | \boldsymbol{x}) dz_j$$

C

c=1

$$\int \log p(c, \boldsymbol{z} | \boldsymbol{x}) p(c, \boldsymbol{z} | \boldsymbol{x}) d\boldsymbol{z}.$$



# Gaussian-Discrete distribution



# Information gain by requesting the j'th feature

$$G_{j}(\boldsymbol{x}) = \sum_{c=1}^{C} \sum_{k=1}^{K} p(c|k)p(k|\boldsymbol{x}) \times \frac{\text{distribution pdf}}{\int \log \left[ p(c, \boldsymbol{x}, z_{j}) \right] p(z_{j}|\boldsymbol{x}, k) dz_{j}}$$
  
$$- \sum_{k=1}^{K} p(k|\boldsymbol{x}) \int \log \left[ p(\boldsymbol{x}, z_{j}) \right] p(z_{j}|\boldsymbol{x}, k) dz_{j}$$
  
$$+ \text{ const.}$$



#### Visual attention: Binary decision based on GIST / FOCUS

#### 27 image features in total

- 1-9: GIST....whole image NMF factors
- 10-27: Foci: 10x10 patch NMF factors

 $N_{train} = 2000$   $N_{test} = 600$ Attention = 28% Random = 41% (p<0.01)





Fig. 2. Visual domain example. The visual data  $(30 \times 30 \text{ pixels})$  is designed as a  $3 \times 3$  array of locations in which one of two binary sparse patterns can be located (two signal classes C = 2). The gist consists of 9 NMF features computed from a large training set of images, while the low-level features are 18 additional NMF features trained on local  $10 \times 10$  patches.

# The Role of Top-Down Attention in the Cocktail Party: Revisiting Cherry's Experiment after Sixty Years

Letizia Marchegiani\*<sup>†</sup>, Seliz G. Karadoğan\*, Tobias Andersen\*, Jan Larsen\* and Lars Kai Hansen\*<sup>‡</sup>

## Weak and Strong Top-Down Attention

$$p(y|k) \rightarrow p(y|k,\beta) = \frac{p(y|k)^{\beta}}{\sum_{c} p(y|k)^{\beta}}$$

$$P(j) = \frac{\exp(\gamma G_j)}{\sum_j \exp(\gamma G_{j'})}$$





**Figure 7.3:** The illustrations of temporal and spectral overlap definitions, the bins represent time-frequency regions of an IBM (frequency bins are not equally spaced, gammatone filtering is used). Only black regions represent overlapped parts on (c).



**Figure 7.6:** Correlation for different LC values, WinLength = 20ms and Num-Chan=32. Left to right: Undirected and Directed. Top to bottom: Temporal and Spectral.

DTU

## Summary



Evidence that cognitive components are the "chunks" of attention

Optimality: Representations are quite uniquely determined by statistical and physical properties of the environment

Attention a function of sensory representations and the mental state/goal/task of the beholder

Optimality: A simple information optimizing mechanism can use task information to determine what to do next and thereby improve decision making

Research supported by Innovation Fund Denmark, Dansh Research Councils, the Lundbeck Foundation, the Novo Nordisk Foundation

# Conclusions & outlook

Evidence that phonemes, gender, identity are independent components 'objects' in the (time stacked) MFCC representation
Evidence that human categorization is based on sparse independent components in social networks, text, digital media
Conjecture: Objects in digital media can be identified as independent components: The brain uses old tricks from perception to solve complex "modern" problems.



## **Acknowledgments**

- Danish Re
- EU Com
- NIH Hum





TU: Toolbox

