SOUND AI

# Professor, PhD Jan Larsen

Section for Cognitive Systems
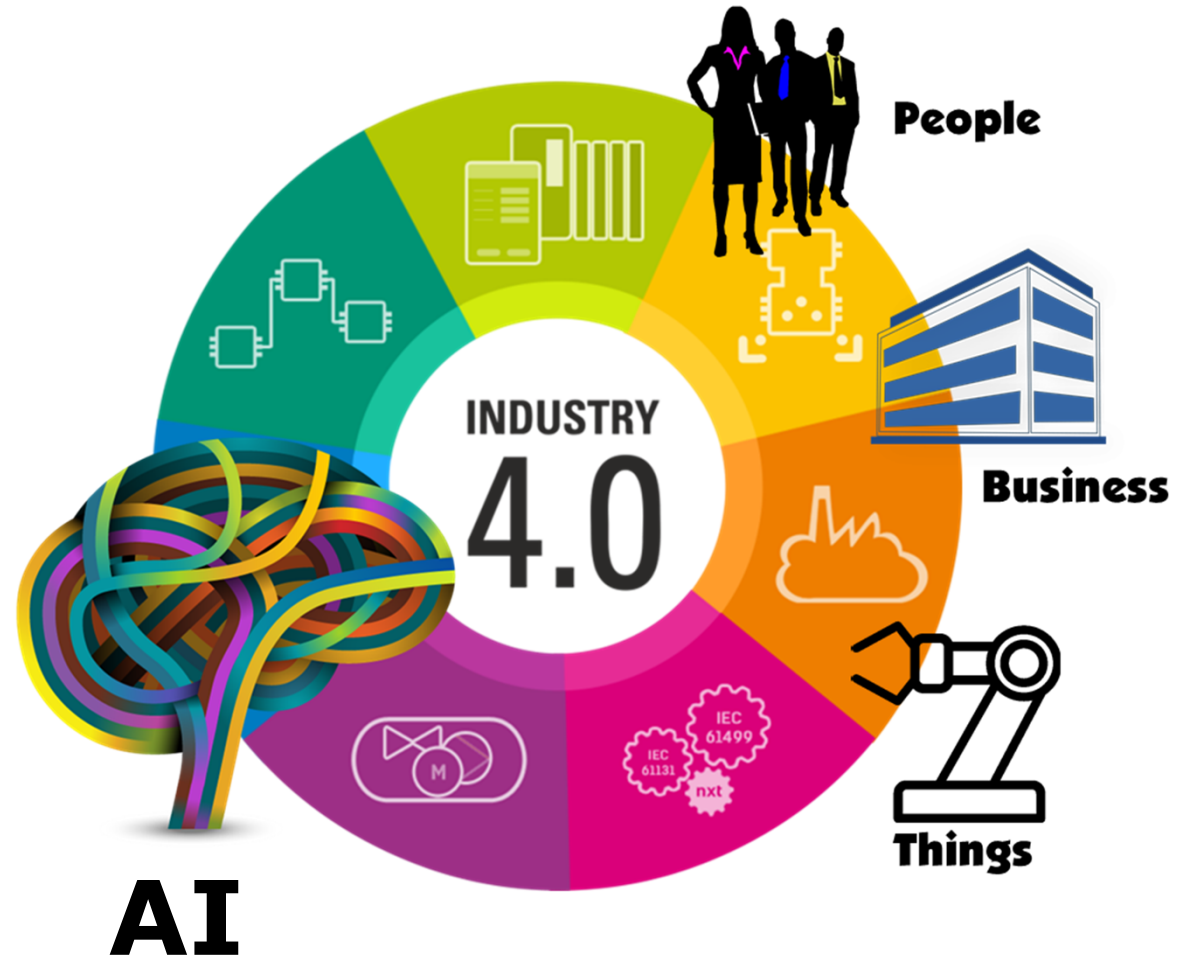DTU Compute, Technical University of Denmark

Participation in 17 international and national collaborative research projects.
Mentoring of 2 Senior Researchers and 9 Postdocs, 34 PhD, and 90 MSc students.
>60% of projects in collaboration with private companies and stakeholders.
Danish Sound Innovation Network national network.
12 commissioned RDI projects.

# My dream related to sound...

To create better quality of life by providing augmented and immersive sound experiences for people in society 4.0 using AI technology

A copy of the physical world through digitization makes it possible for cyber-physical systems to communicate and cooperate with each other and with humans in real time and perform decentralized decision-making



INDUSTRY 4.0

People

Business

Things

AI

https://en.wikipedia.org/wiki/Industry_4.0
B. Marr: Forbes, June 20, 2016, http://www.forbes.com/sites/bernardmarr/2016/06/20/what-everyone-must-know-about-industry-4-0/#4c979f804e3b
http://www.enterrasolutions.com/2015/10/industry-4-0-facing-the-coming-revolution.html

# Industry 4.0 = Civilization 4.0

It is a cognitive revolution that could be even more disruptive than earlier as it concerns not only the industry but the whole way we live our lives.

# AI
## Artificial Intelligence

# IA

## Intelligence Augmentation

# research focus

# CoSound

Machine learning based processing of audio data and related information, such as context, users' states, interaction, intention, and goals with the purpose of providing innovative services related to societal challenges in

**Transforming big audio data into semantically interoperable data assets and knowledge:** Enrichment and navigation in large sound archives such as broadcast

**Experience economy and edutainment:** New music services based on mood, optimization of sound systems

**Healthcare:** Music interventions to improve quality of life in relation to disorders such as e.g. stress, pain, and ADHD.
User-driven optimization of hearing aids.

# research focus

## MakeSense

Processing of sensor signals and related IoT data streams with the purpose of fostering innovative systems addressing societal challenges in

**Food:** Grain analysis

**Security**: Explosives and drug detection

**Health**: Blood and water analysis, intelligent drug delivery and sensing, e-health, personalized medicine

**Energy**: Wind mill maintenance

**Environment**: Exhaust gas analysis, large diesel engine predictive monitoring

**Resource efficiency**: Waste sorting

**Digital economy**: Event recommendation

# SOUND IS AFFECTIVE

# What are the mechanism? – the BRECVEM model

- **Brain stem reflexes** linked to acoustical properties, e.g. loudness

- **Evaluative conditioning** – association between music and emotion when they occur together

- **Emotional contagion** – emotion expressed in music, sad is e.g. linked low-pitches, slow, and quiet

- **Rhythmic entrainment** – movement synchronization with rhythm

- **Visual images** – creation of visual images

- **Episodic memories** – e.g. strong emotion when you hear a melody linked to an episode

- **Cognitive appraisal** -  mental analysis of music an creation of analytic or aesthetic pleasure (hit-songs)

- **Musical expectancy** -  balance between surprise and expectation

Ref: Juslin, P. N. and Västfäll, D. *Emotional response to music: The need to consider underlying mechanism. Behavioral and Brain Sciences, vol. 31, pp. 559–621, 2008*.
Line Gebauer & Peter Vuust, *Music interventions in Health Care, 2014.*

# AI IS EFFECTIVE

# What is machine learning?

Learning structures and patterns form from historical data to reliably predict outcome for new data.

Computers only do what they are programmed to do. ML infers new relations and patterns, which were not programmed. They learn and adapt to changing environment.

1. Computer systems that automatically improve through experience, or learns from data.
2. Inferential process that operate from representations that encode probabilistic dependencies among data variables capturing the likelihoods of relevant states in the world.
3. Development of fundamental statistical computational-information-theoretic laws that govern learning systems - including computers, humans, and other entities.

M. I. Jordan and T. M. Mitchell. *Machine learning: Trends, perspectives, and prospects*. Science, July 2015.
Samuel J. Gershman, Eric J. Horvitz, Joshua B. Tenenbaum. *Computational rationality: A converging paradigm for intelligence in brains, minds, and machines.* Science, July 2015.

# Brief history of AI

Late 40's Allan Touring: theory of computation

1948 Claude Shannon: A Mathematical Theory of Communication

1948 Norbert Wiener: Cybernetics - *Control and Communication in the Animal and the Machine*
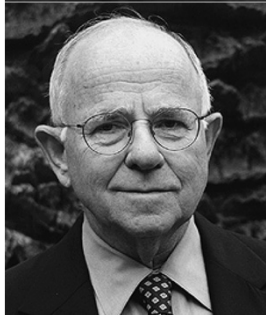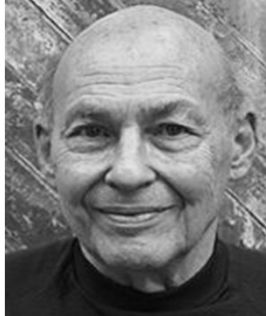
1950 The Touring test

1951 Marvin Minsky's analog neural networks (1st generation)

1956 Dartmouth conference: Artificial intelligence with aim of human like intelligence

1960 Bernard Widrow's ADALINE - adatpive linear systems

1956-1974 Many small scale "toy" projects in robotics, control and game solving

1974 Failure of success and Minsky's criticism of perceptron, lack of computational power, combinatorial explosion, Moravec's paradox: simple tasks are not easy to solve

1980's Expert systems useful in restricted domains

1980's Knowledge based systems – integration of diverse information sources

1980's The 2nd generation neural network revolution starts

Late 1980's Robotics and the role of embodiment to achieve intelligence

1990's AI and cybernetics research under new names such as machine learning, computational intelligence, evolutionary computing, neural networks, Bayesian networks, complex systems, game theory, deep neural networks (3rd generation) cognitive systems
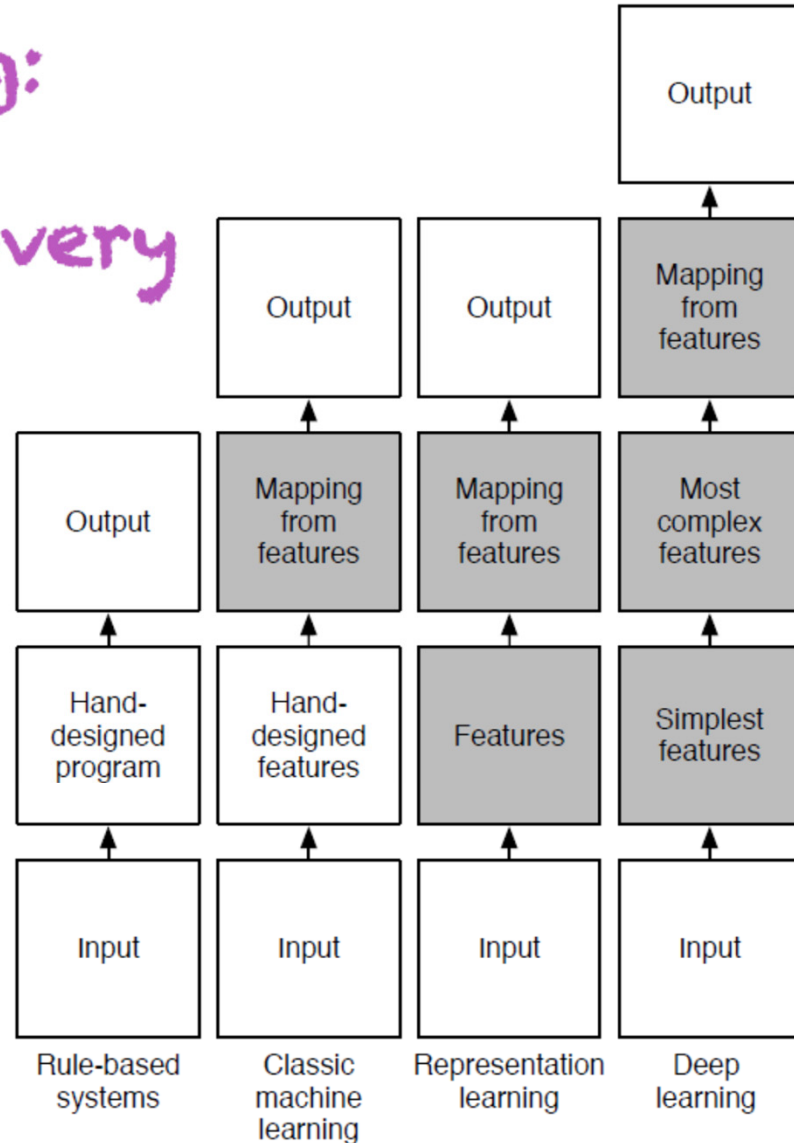
2010's deep neural networks (4rd generation) and cognitive systems, large scale data and computational frameworks, ML is commoditized
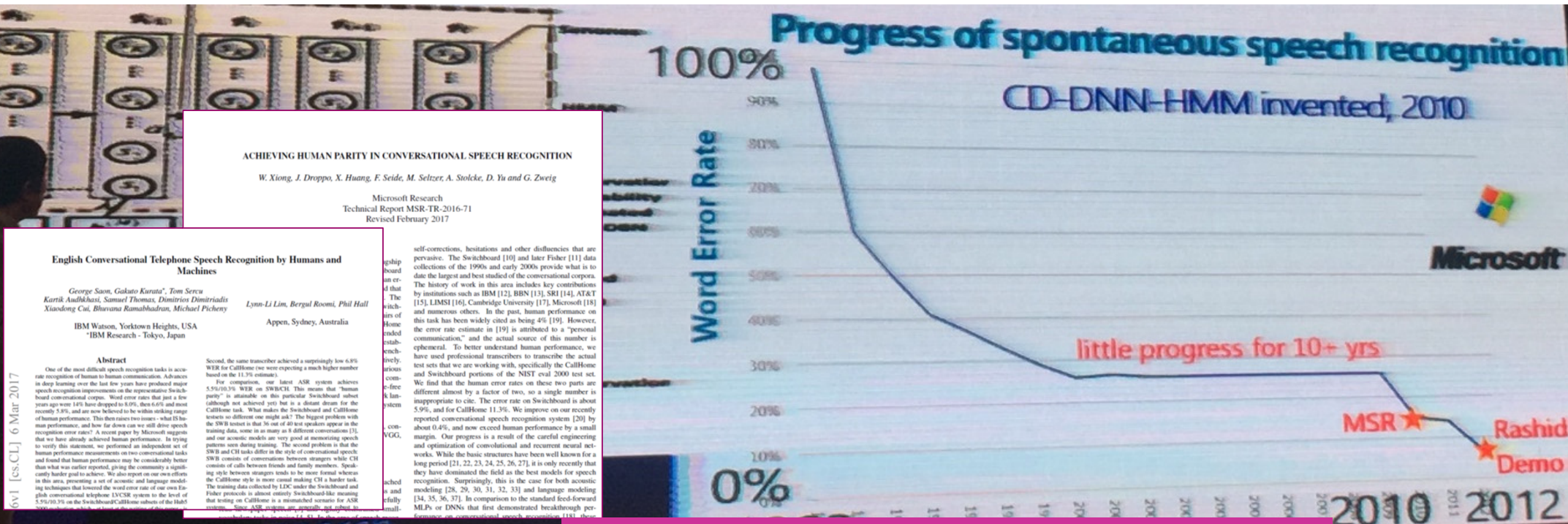
# Deep Learning: Automating Feature Discovery

Geoff Hinton, Yoshua Bengio, Yann LeCun, Deep Learning Tutorial, NIPS 2015, Montreal.

**Deep learning is a disruptive technology**



| Rule-based systems | Classic machine learning | Representation learning | Deep learning |
|---|---|---|---|
| | | | Output |
| | Output | Output | Mapping from features |
| Output | Mapping from features | Mapping from features | Most complex features |
| Hand-designed program | Hand-designed features | Features | Simplest features |
| Input | Input | Input | Input |

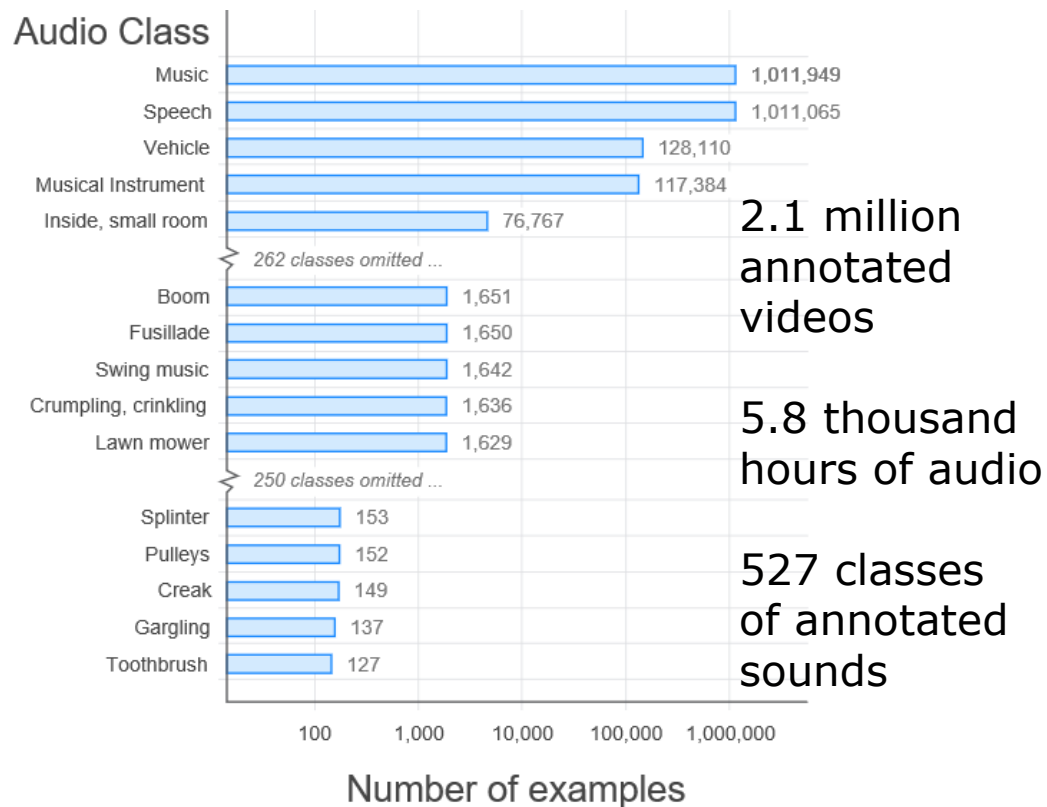# Machine learning is very successful for speech recognition and chat bots



Human parity is achieved Feb/March 2017

Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury. *Deep Neural Networks for Acoustic Modeling in Speech Recognition.* IEEE Signal Processing Magazine, 82, Nov. 2012.
George Saon, Gakuto Kurata, Tom Sercu, Kartik Audhkhasi, Samuel Thomas, Dimitrios Dimitriadis, Xiaodong Cui, Bhuvana Ramabhadran, Michael Picheny, Lynn-Li Lim, Bergul Roomi, Phil Hall. *English Conversational Telephone Speech Recognition by Humans and Machines, https://arxiv.org/abs/1703.02136, March 2017*
W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, G. Zweig. *Achieving Human Parity in Conversational Speech Recognition, https://arxiv.org/abs/1610.05256, October 2016.*

# Machine learning is very successful for audio classification



2.1 million annotated videos

5.8 thousand hours of audio

527 classes of annotated sounds

**Table 2:** Comparison of performance of several DNN architectures trained on 70M videos, each tagged with labels from a set of 3K. The last row contains results for a model that was trained much longer than the others, with a reduction in learning rate after 13 million steps.

| Architectures | Steps | Time | AUC | d-prime | mAP |
|---|---|---|---|---|---|
| Fully Connected | 5M | 35h | 0.851 | 1.471 | 0.058 |
| AlexNet | 5M | 82h | 0.894 | 1.764 | 0.115 |
| VGG | 5M | 184h | 0.911 | 1.909 | 0.161 |
| Inception V3 | 5M | 137h | **0.918** | **1.969** | 0.181 |
| ResNet-50 | 5M | 119h | 0.916 | 1.952 | **0.182** |
| ResNet-50 | 17M | 356h | **0.926** | **2.041** | **0.212** |

Mean average precision mAP is low because of low class prior $<10^{-4}$.

AUC is the area under curve of true positive rate vs. false positive rate.

Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, Marvin Ritter. *Audio Set: An ontology and human-labeled dataset for audio events*, IEEE ICASSP 2017, New Orleans, LA, March 2017.

Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron Weiss, Kevin Wilson. *CNN Architectures for Large-Scale Audio Classification*, ICASSP 2017, New Orleans, LA, March 2017.

# Machine learning is very successful for speech generation

**WaveNet** is a deep generative model of raw audio waveforms

WaveNets are able to generate speech which mimics any human voice and which sounds more natural than the best existing Text-to-Speech systems, reducing the gap with human performance by over 50%.
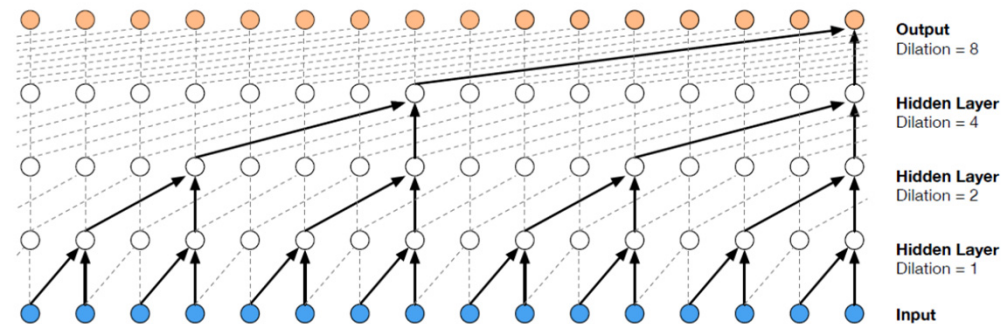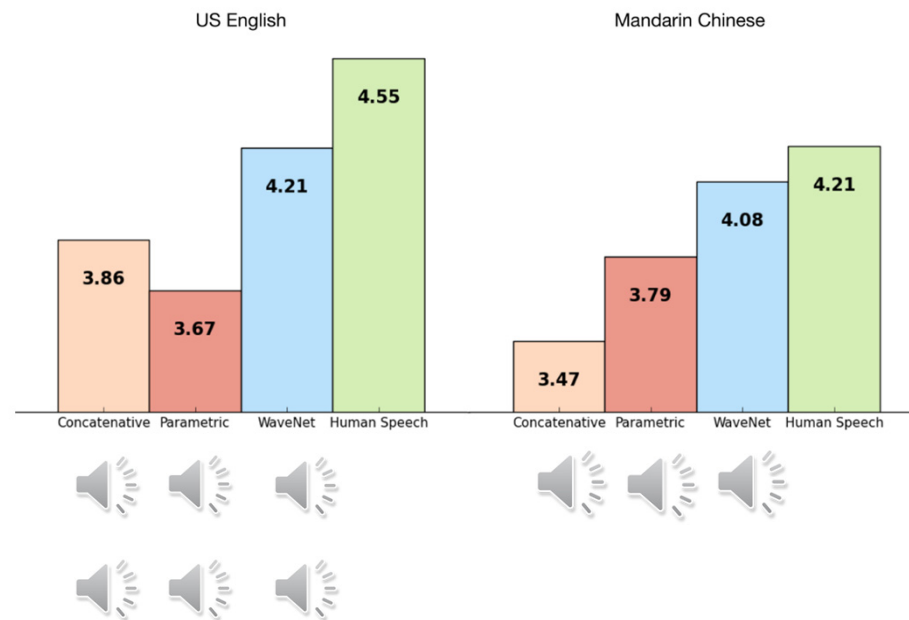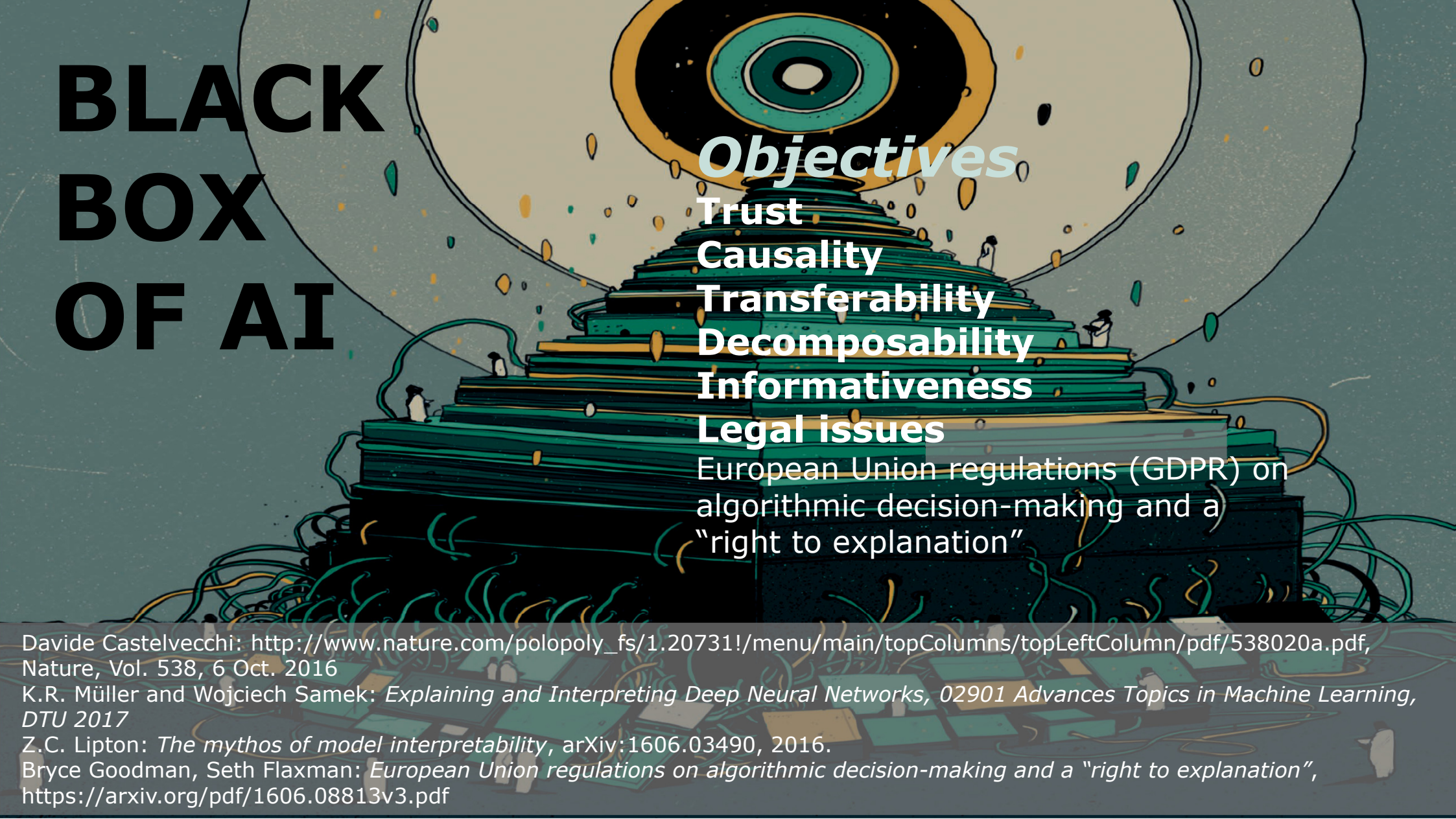


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu. *WAWENET: A Generative Model for Raw Audio,* https://arxiv.org/pdf/1609.03499.pdf, Sept 2016, https://deepmind.com/blog/wavenet-generative-model-raw-audio/

# BLACK BOX OF AI

## *Objectives*

**Trust**
**Causality**
**Transferability**
**Decomposability**
**Informativeness**
**Legal issues**
European Union regulations (GDPR) on algorithmic decision-making and a "right to explanation"

Davide Castelvecchi: http://www.nature.com/polopoly_fs/1.20731!/menu/main/topColumns/topLeftColumn/pdf/538020a.pdf, Nature, Vol. 538, 6 Oct. 2016
K.R. Müller and Wojciech Samek: *Explaining and Interpreting Deep Neural Networks, 02901 Advances Topics in Machine Learning, DTU 2017*
Z.C. Lipton: *The mythos of model interpretability*, arXiv:1606.03490, 2016.
Bryce Goodman, Seth Flaxman: *European Union regulations on algorithmic decision-making and a "right to explanation"*, https://arxiv.org/pdf/1606.08813v3.pdf

# Adversarial learning



(a) Original Input

| | blues | classical | country | disco | hiphop | jazz | metal | pop | reggae | rock | Pr |
|---|---|---|---|---|---|---|---|---|---|---|---|
| blues | 92.0 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 | 0.0 | 0.0 | 4.0 | 8.0 | 85.2 |
| classical | 0.0 | 84.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 100.0 |
| country | 0.0 | 4.0 | 92.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 | 92.0 |
| disco | 8.0 | 4.0 | 4.0 | 80.0 | 0.0 | 0.0 | 0.0 | 4.0 | 12.0 | 16.0 | 62.5 |
| hiphop | 0.0 | 0.0 | 0.0 | 0.0 | 76.0 | 0.0 | 0.0 | 0.0 | 8.0 | 0.0 | 90.5 |
| jazz | 0.0 | 8.0 | 0.0 | 0.0 | 0.0 | 92.0 | 0.0 | 0.0 | 4.0 | 0.0 | 88.5 |
| metal | 0.0 | 0.0 | 0.0 | 0.0 | 20.0 | 0.0 | 92.0 | 0.0 | 4.0 | 0.0 | 79.3 |
| pop | 0.0 | 0.0 | 0.0 | 12.0 | 0.0 | 4.0 | 0.0 | 92.0 | 0.0 | 16.0 | 74.2 |
| reggae | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 | 64.0 | 8.0 | 84.2 |
| rock | 0.0 | 0.0 | 4.0 | 8.0 | 4.0 | 0.0 | 8.0 | 0.0 | 4.0 | 48.0 | 63.2 |
| F | 88.5 | 91.3 | 92.0 | 70.2 | 82.6 | 90.2 | 85.2 | 82.1 | 72.7 | 54.5 | 81.2 |

(b) Input intercepted by adversary A1

| | blues | classical | country | disco | hiphop | jazz | metal | pop | reggae | rock | Pr |
|---|---|---|---|---|---|---|---|---|---|---|---|
| blues | 16.0 | 4.0 | 0.0 | 4.0 | 12.0 | 4.0 | 12.0 | 0.0 | 12.0 | 20.0 | 19.0 |
| classical | 16.0 | 8.0 | 12.0 | 16.0 | 4.0 | 12.0 | 8.0 | 4.0 | 16.0 | 4.0 | 8.0 |
| country | 8.0 | 8.0 | 4.0 | 12.0 | 4.0 | 0.0 | 12.0 | 4.0 | 4.0 | 8.0 | 6.2 |
| disco | 8.0 | 16.0 | 8.0 | 8.0 | 4.0 | 8.0 | 12.0 | 12.0 | 16.0 | 12.0 | 7.7 |
| hiphop | 0.0 | 4.0 | 16.0 | 8.0 | 20.0 | 40.0 | 16.0 | 20.0 | 0.0 | 4.0 | 15.6 |
| jazz | 20.0 | 16.0 | 12.0 | 12.0 | 0.0 | 8.0 | 0.0 | 20.0 | 8.0 | 8.0 | 7.7 |
| metal | 4.0 | 12.0 | 20.0 | 12.0 | 12.0 | 12.0 | 4.0 | 8.0 | 4.0 | 12.0 | 4.0 |
| pop | 4.0 | 4.0 | 8.0 | 8.0 | 20.0 | 8.0 | 28.0 | 20.0 | 8.0 | 8.0 | 17.2 |
| reggae | 0.0 | 16.0 | 8.0 | 4.0 | 12.0 | 4.0 | 4.0 | 8.0 | 8.0 | 20.0 | 9.5 |
| rock | 24.0 | 12.0 | 12.0 | 16.0 | 12.0 | 4.0 | 4.0 | 4.0 | 24.0 | 4.0 | 3.4 |
| F | 17.4 | 8.0 | 4.9 | 7.8 | 17.5 | 7.8 | 4.0 | 18.5 | 8.7 | 3.7 | 10.0 |

Corey Kereliuk, Bob L. Sturm, Jan Larsen: Deep Learning and Music Adversaries, IEEE Transactions on Multimedia, Nov. 2015

Corey Kereliuk, Bob L. Sturm, Jan Larsen: Deep Learning, Audio Adversaries, and Music Content Analysis, 2015 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, Oct. 2015

Corey Kereliuk, Bob L. Sturm, Jan Larsen: ?El Caballo Viejo? Latin Genre Recognition with Deep Learning and Spectral Periodicity, Fifth Biennial International Conference on Mathematics and Computation in Music (MCM2015), 2015.
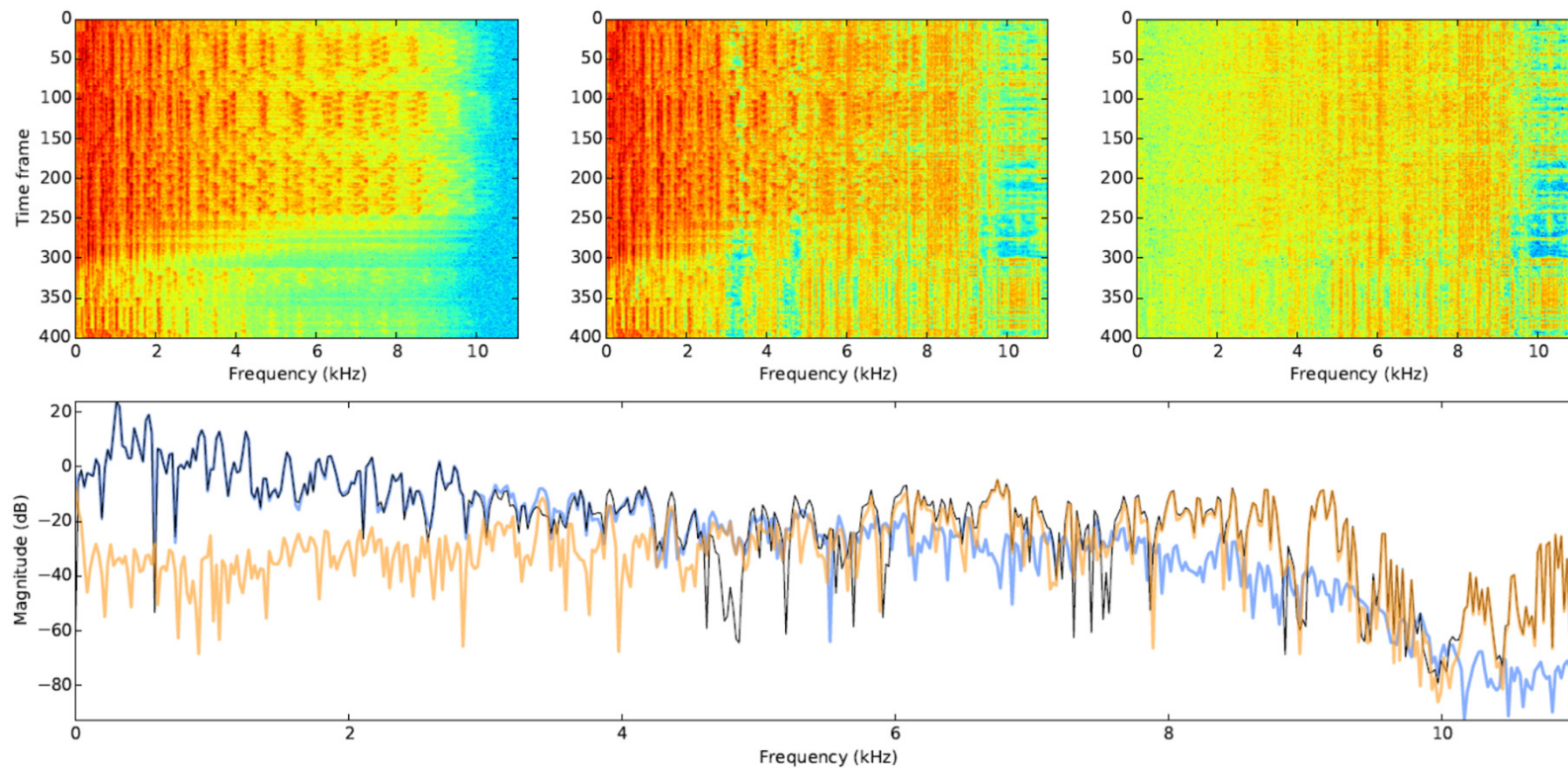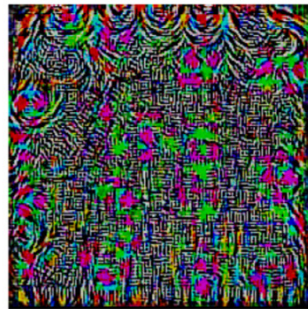
# Adversarial learning



Fig. 5. Top left: spectrogram excerpt from *GTZAN* Classical "21" (Mozart, Symphony No. 39 Finale) that the DNN-based system in Fig. 2(b) classifies as *Classical*. Top middle: spectrogram of adversarial example classified as *Reggae*. Top right: spectrogram of the difference of the two. Bottom: magnitude spectrum of one frame (1024 samples) of the original (light blue), adversarial example (black), and difference (orange). Note that all excerpts in *GTZAN* have a sampling rate of 22050 Hz. The SNR = 21.1dB.

Corey Kereliuk, Bob L. Sturm, Jan Larsen: Deep Learning and Music Adversaries, IEEE Transactions on Multimedia, Nov. 2015

Corey Kereliuk, Bob L. Sturm, Jan Larsen: Deep Learning, Audio Adversaries, and Music Content Analysis, 2015 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, Oct. 2015

Corey Kereliuk, Bob L. Sturm, Jan Larsen: ?El Caballo Viejo? Latin Genre Recognition with Deep Learning and Spectral Periodicity, Fifth Biennial International Conference on Mathematics and Computation in Music (MCM2015), 2015.
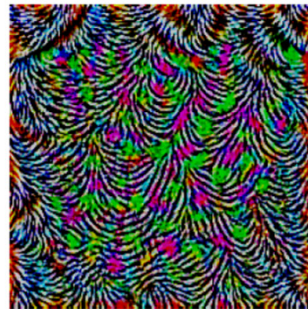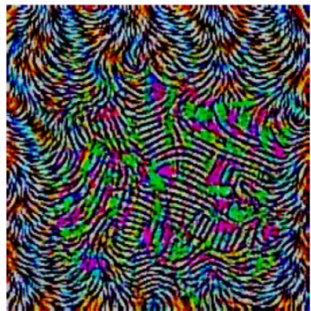
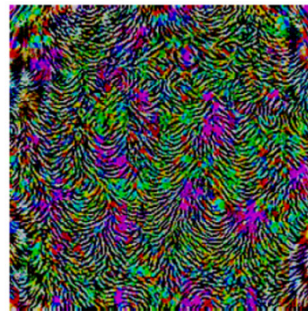# Universal Adversarial Learning



(a) CaffeNet  (b) VGG-F  (c) VGG-16
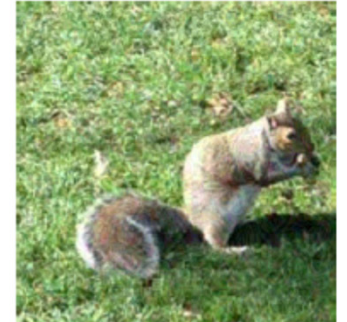(d) VGG-19  (e) GoogLeNet  (f) ResNet-152

wool    Indian elephant    Indian elephant

common newt    carousel    grey fox

Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, Pascal Frossard: Universal adversarial perturbations, arXiv:1610.08401. 2017

# What defines simple and complex problems - and how do we solve them them?

passive

active and autonoumous

exploration and summarization

prediction

continuous learning
reflection
pro-activeness
engagement
experimentation
creativity

**Unreasonable effectiveness of**

Mathematics E. Wigner, 1960

Data Halevy, Norvig, Pereira, 2009

RNNs Karpathy, 2015

**Experimentation and interaction through users-in-the-loop**

# INTERACTIVE MACHINE LEARNING IN SOUND

# Expressed emotions in music

- Jens Madsen, Jan Larsen. The Confidence Effect in Elicitation of Expressed Emotion in Music. To be submitted.

- Jens Madsen, Jan Larsen. Designing a Cognitive Music System. To be submitted.

- Jens Madsen, Bjø_____ure Representations

- Jens Madsen, Bjørn S_____ Music using Probabilistic Features Represen_____, submitted IEEE T-ASPL, 2015.

- Jens Madsen, Bjørn Sand Jensen, Jan Larsen and Jens Brehm Nielsen. *Towards Predicting Expressed Emotion in Music from Pairwise Comparisons*, 9th Sound and Music Computing Conference, 2012.

- Jens Madsen, Jens Brehm Nielsen, Bjørn Sand Jensen and Jan Larsen. *Modeling Expressed Emotions in Music using Pairwise Comparisons*, 9th International Symposium on Computer Music Modeling and Retrieval (CMMR) 2012.

- Madsen, J., Jensen, B.S., Larsen, J., Predictive modeling of expressed emotions in music using pairwise comparisons. M. Aramaki et al. (Eds.): CMMR 2012, LNCS 7900, pp. 253–277, 2013. Springer-Verlag Berlin Heidelberg 2013.

- Is it possible to model the users representation of expressed and induced emotion?
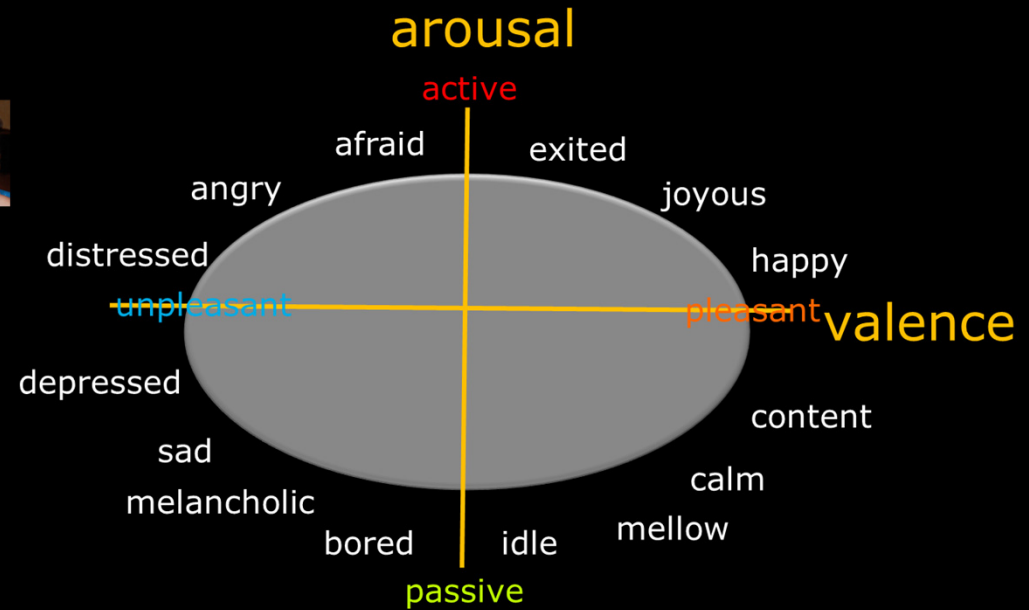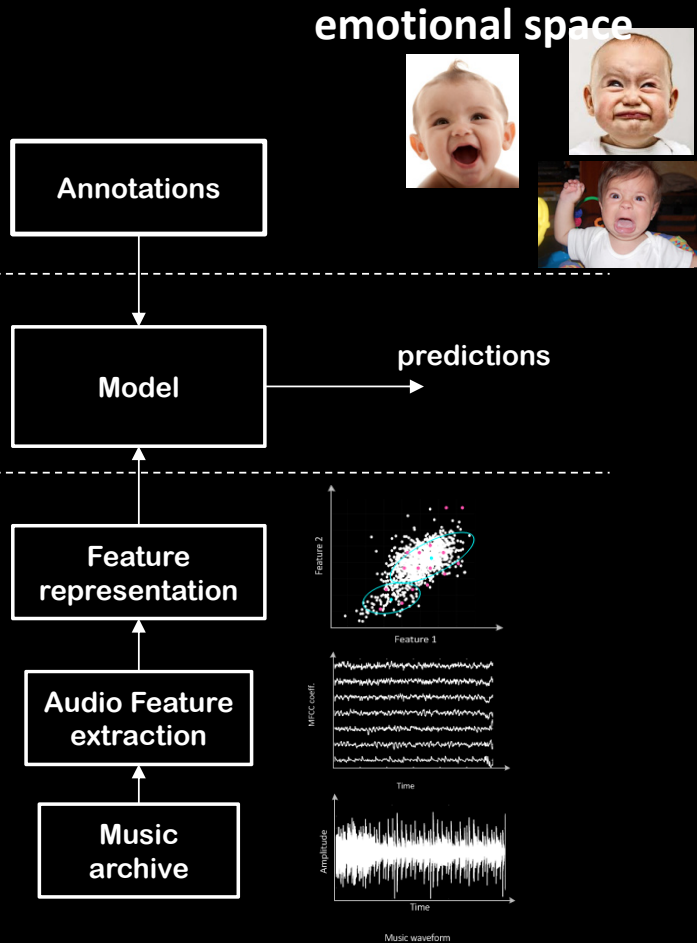
- Which scaling method should we use?

- Which role does mood play?

# Music Emotion Modeling

emotional space



## User modeling/ experimental paradigm

Annotations

## Machine learning

Model → predictions

## Audio signal processing/ Machine learning

Feature representation

Audio Feature extraction

Music archive



arousal

active

afraid | exited

angry | joyous

distressed | happy

unpleasant —————— pleasant valence

depressed

content

sad | calm

melancholic | mellow

bored | idle

passive

J. A. Russel: "A Circumplex Model of Affect," *Journal of Personality and Social Psychology*, 39(6):1161, 1980

J. A. Russel, M. Lewicka, and T. Niit, "A Cross-Cultural Study of a Circumplex Model of Affect," *Journal of Personality and Social Psychology*, vol. 57, pp. 848-856, 1989

# Learning curve modeling arousal shows nonlinear modelling is best

# How many pairwise comparisons do we need to model emotions?



**Using active learning**
15% for valence
9% for arousal

Madsen, J., Jensen, B.S., Larsen, J., Predictive modeling of expressed emotions in music using pairwise comparisons. M. Aramaki et al. (Eds.): CMMR 2012, LNCS 7900, pp. 253–277, 2013. Springer-Verlag Berlin Heidelberg 2013

# The power of human data

**Why - Humans as a measurement device**

**How - Humans in the loop**

**Who - Humans in the loop**

## Why - Humans as a measurement device

- With the purpose of individualization and dynamical response.
- With the purpose of group studies and population models.
- For eliciting perceptual, affective, and cognitive aspects.
- For acquiring other aspects e.g. behavioral and physical.
- For quality measurement and control.
- For obtaining shared cognitive and cultural information and contexts that helps disambiguation of meaning.

## How - Humans in the loop

- **Direct** measurement of physiological, cognitive and behavior states from physical devices.

- **Indirect** measurements from self-reports, experiments using direct, indirect and related scaling methods of objective or subjective information.

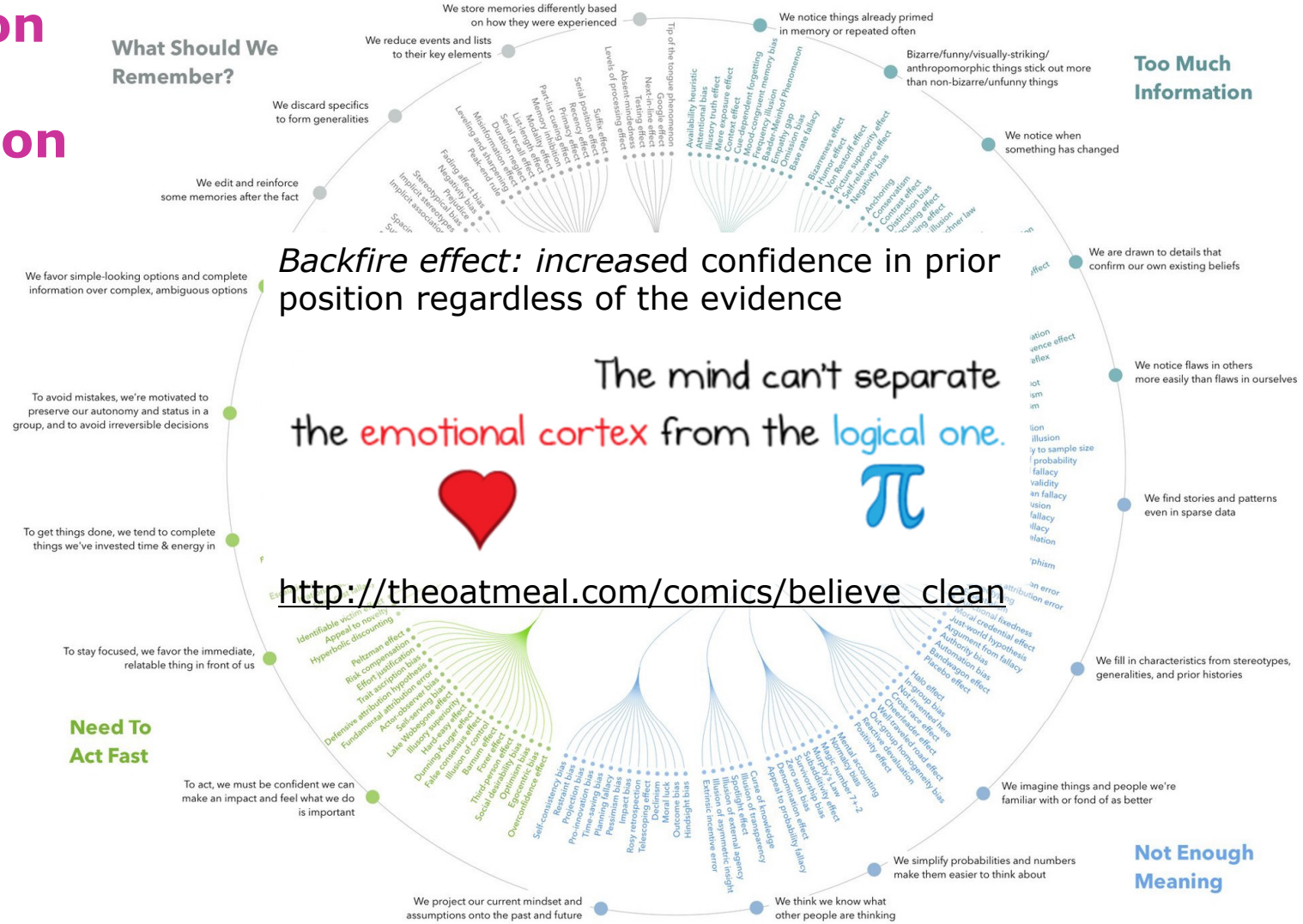Whether data are Experimental or Observational plays an important role!

## Who - Humans in the loop

- End-user
- Experimenter
- Developer
- Expert user
- Collaborative, transfer learning for crowds of humans

**Challenge: Robust adaptive learning and optimization from interaction with inconsistent, biased and often inattentive users**

- Modeling and/or knowledge of many aspects of the state of person(s) and the environment
- Modeling and representing uncertainty
  - concerning the "objective" (incl. needs, intentions, level of engagement)
  - concerning the interaction/answers/measurements from the subjects
- Support for varying complexity of a multi-aspect objective function
- Adaptive/online elicitation and learning of the objective function

# Human interaction with information

COGNITIVE BIAS CODEX, 2016

**What Should We Remember?**

**Too Much Information**

**Need To Act Fast**

**Not Enough Meaning**

*Backfire effect: increase*d confidence in prior position regardless of the evidence

The mind can't separate the emotional cortex from the logical one.

http://theoatmeal.com/comics/believe_clean

ALGORITHMIC LAYOUT + DESIGN BY JM3 · JOHN MANOOGIAN III // CONCEPT + METICULOUS CATEGORIZATION BY BUSTER BENSON // DEEP RESEARCH BY WIKIPEDIANS FAR + WIDE

# Interactive Learning / Sequential Experimental Design

**Generalization objective**
Eliciting and learning the entire model / objective function.
Expected change in relative entropy is derived from the posterior and predictive distribution.

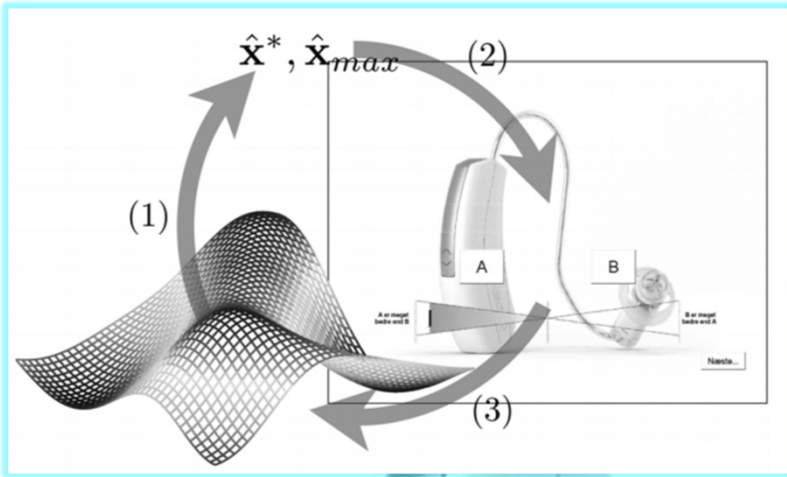**Optimization objective**
Learning and identifying optimum
The Expected Improvement of a new candidate sample (green points) is derived from the predictive distribution.



Which of the four green parameters settings/products/interface, x, should the user assess (rate/annotate/see/hear), or where do we need tp evaluate objective performance measurements

# General framework

State of users' mind

Users' profile

Intention/task/objective

Context

observation y

Subjective users' assessments or objective performance measurements

Interface

Probabilistic model

Sequential design

object(s)

features rep. object(s)

Systems/objects represented by features

proposed object(s), feature(s), user(s)

- Highly personalization needs.
- Dynamic environment and use with different needs.
- Latent, convoluted object functions which are difficult to express though verbal and motor actions.
- Users with disabilities – and often elderly people - provide inconsistent and noisy interactions.

**Optimization of hearing aids using Bayesian optimization**

Jens Brehm Nielsen, Jakob Nielsen: Efficient Individualization of Hearing and Processers Sound, ICASSP2013.
Jens Brehm Nielsen, Jakob Nielsen, Jan Larsen: Perception based Personalization of Hearing Aids using Gaussian Process and Active Learning, IEEE Trans. ASLP, vol. 23, no. 1, pp. 162 – 173, Jan 2015.
Maciej Korzepa, Michael Kai Petersen, Benjamin Johansen, Jan Larsen, Jakob Eg Larsen: Learning soundscapes from OPN sound navigator, poster 2017.
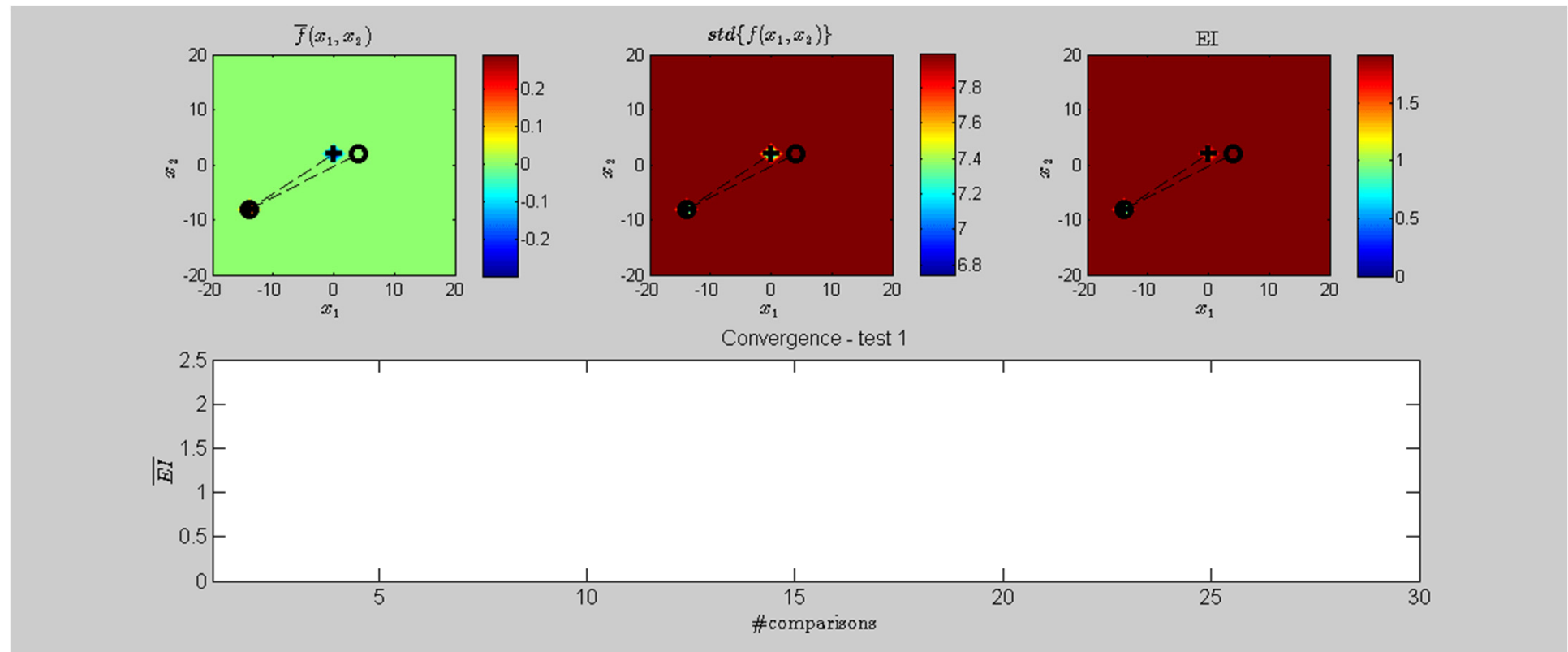
# Pairwise (2AFC) personalization of HA

# Hearing Aids Personalization

A real interactive optimization sequence in 30 iterations

# VOXVIP - smart crowdsourcing of the DR radio archive



voxvip.cosound.dk

Can smart crowdsourcing efficiently enrich radio archives with high quality metadata using machine learning and gamification?

Are model-based, active learning mechanisms suitable for smart crowdsourcing, and is optimal performance as regards time-use achieved?
Are age, sex, address relevant for recognition of specific voices?
Gamification: How does levels, difficulty and point assignment influence the quality and quantity of annotations?

# What is meta information?

Infinite number of aspects provides information about the individual clip/object or similarity between such objects

**Objective information**
- Who is speaking
- What is the topic discussed?
- Which objects are present in the clip?

**Subjective information**
- Which emotions are expressed in the clip?
- What is the sound quality?
- Which clip is preferred?
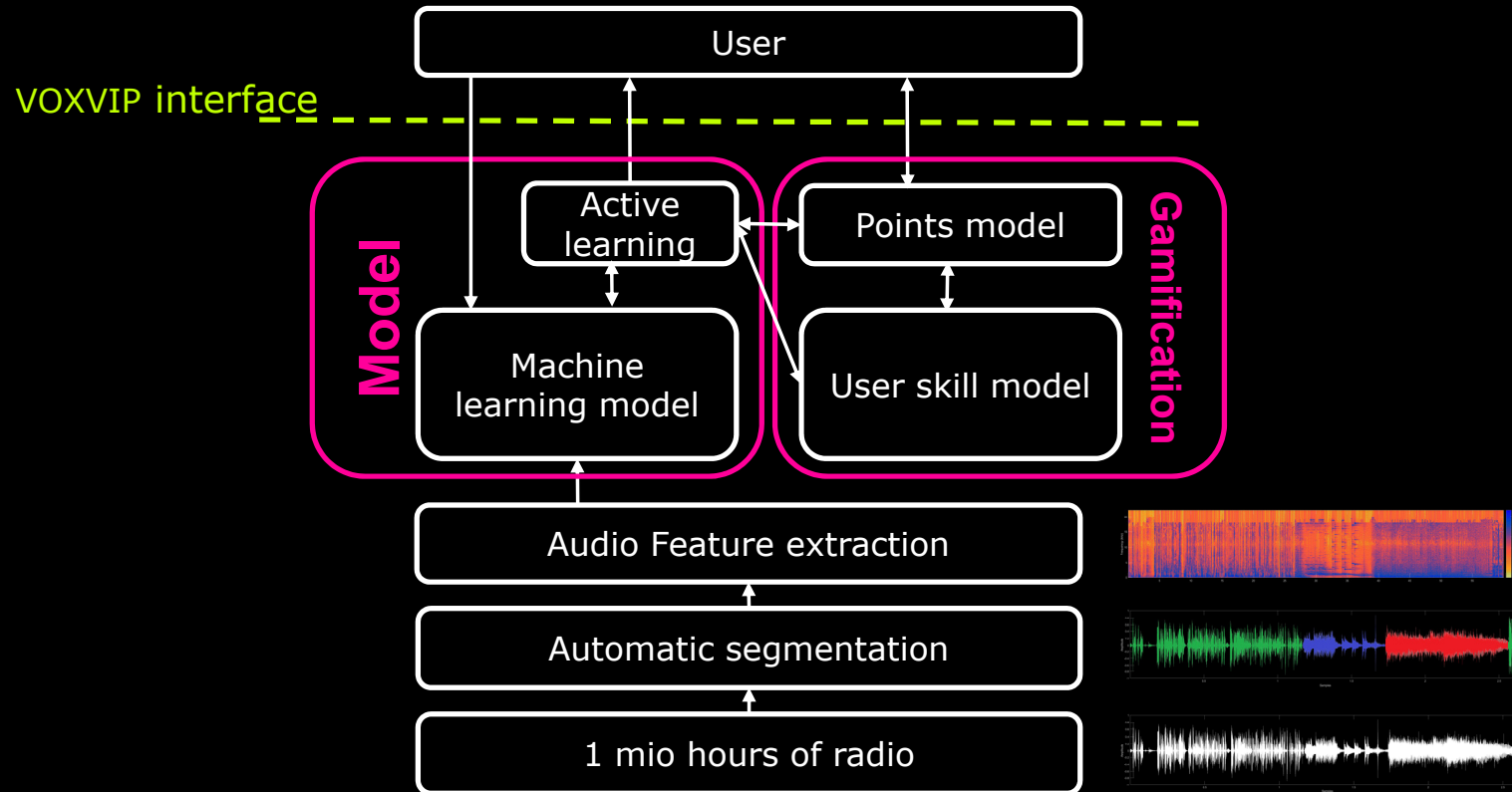
# How can meta information be created?

Lack of specific annotations requires prior knowledge

Manual annotation is limited or impossible due to the size of the archive, human resources, or annotators qualifications.

Semi-automatic machine learning can be used to predict information in the enture archive based on limited number of annotations.

Smart crowdsourcing exploits machine learning to predict information in the entire archive based on 'crowd annotators' annotations. The individual clip is selected based on uncertain information about the label, the annotators' qualifications and engagement based on active learning mechanisms.

# VOXVIP model

ALIVE INSIDE

by Michael Rossato Bennett, 2014
www.youtube.com/watch?v=5FWn4JB2YLU

HENRY

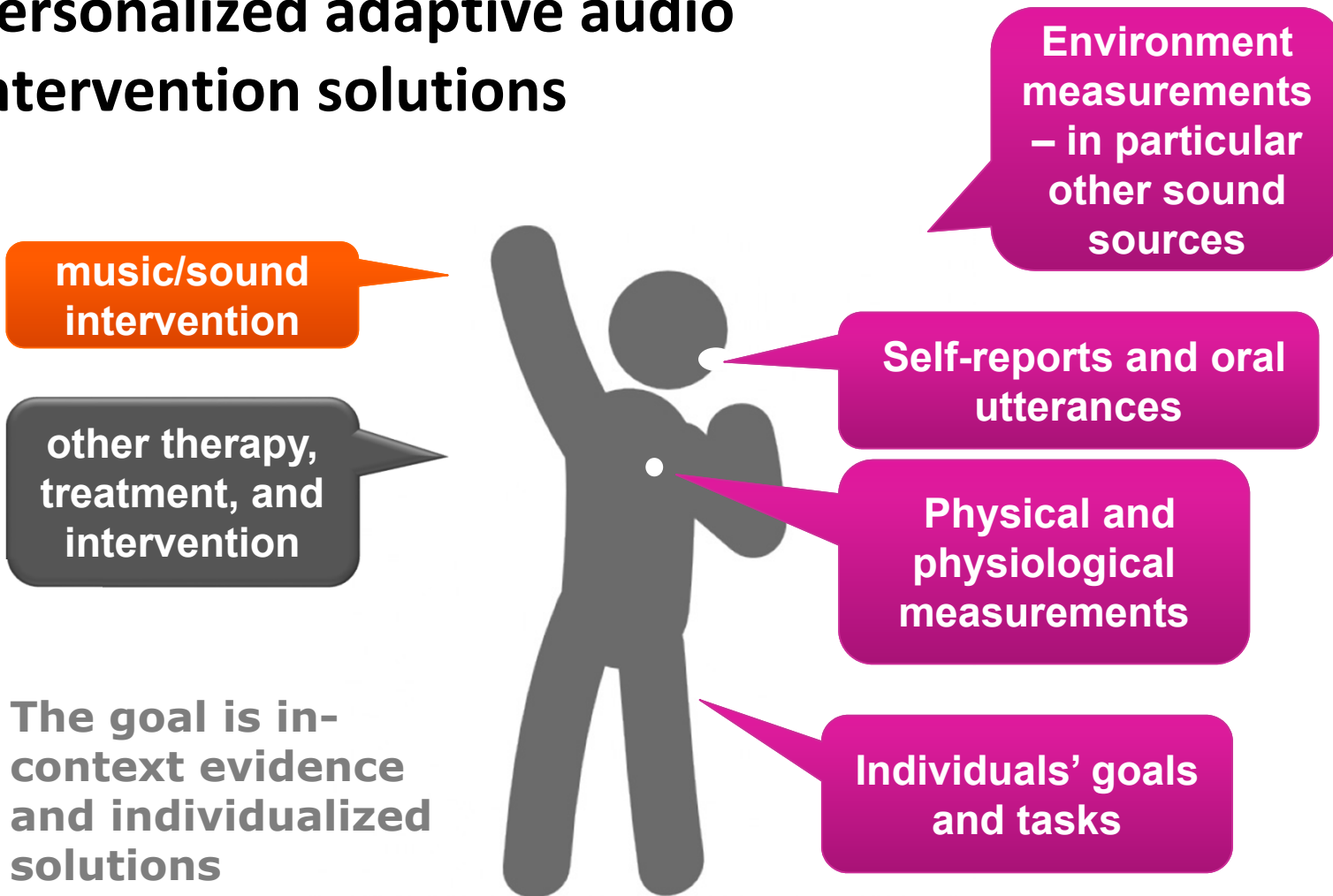# MUSIC AND SOUND INTERVENTION FOR IMPROVING SLEEP IN DEMENTIA PATIENTS

- Anecdotal reports
- Preserved ability to engage in musical activities
- Reduce social isolation
- Improve cognitive symptoms
- Reduce aggression
- Effects might not be specific to music

People highly absorbed in music (AIMS) listening to unfamiliar, but preferred music has higher recovery from a stress situation

S.L. Carstensen, J. Madsen, J. Larsen. *The Influence of Familiarity and Absorption on the Effectiveness of Music in Stress Reduction, in submission 2018.*
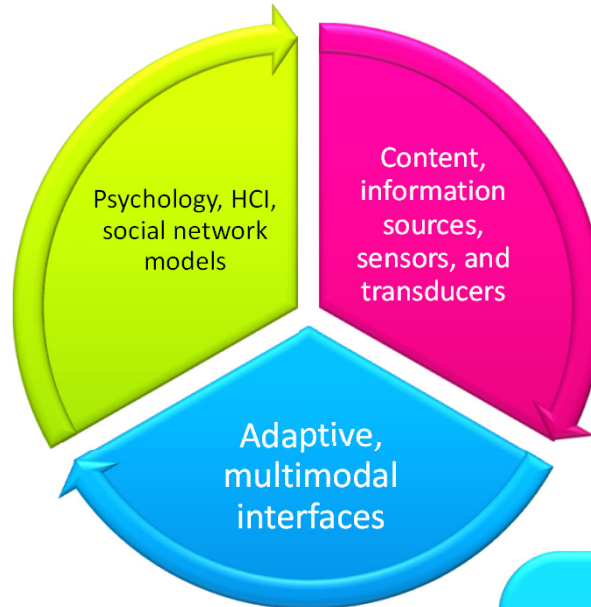
# FUTURE

# Cognizant audio systems
## *fully informed and aware systems*

**Context:**
who, where, what

**Users in the loop:**
direct and indirect

**Interactive dialog with the user enables long term/continuous behavior tracking, personalization, elicitation of perceptual and affective preferences, as well as adaptation**

**Listen in on audio and other sensor streams to segment, identify and understand**

Psychology, HCI, social network models

Content, information sources, sensors, and transducers

Adaptive, multimodal interfaces

**Flexible integration with other media modalities**

**Mixed modality experience: Use other modalities to enhance, substitute or provide complementary information**

TRANSFORMATION

Digital
Interactive - humans in the loop
Human centric
Cognitive
AI drives IA
Multimodal - IoT and hearables

# THE WAYS AHEAD

- Need for possibility to include co-creation and production.

- Need for more data across domains and situations.

- Need for systems and platforms that enables experimentation and direct user interaction.

- Need for better AI and machine learning methodology that can provides robust, interpretable, interactive learning from few examples.