

# Superhuman Speech Analysis? Getting Broader, Deeper & Faster.



## Björn W. Schuller

Head GLAM, Imperial College London

Chair EHW, University of Augsburg

CEO audEERING



Superhuman?

# Superhuman? ASR.

- **Human: ASR**

Misses ~1-2 words in 20 → 5-10% “Word Error Rate” (WER)

→ 1 minute conversation ~16 words

- **Machine: ASR**

*Switchboard*: 2.4k (260 hrs), 543 speakers

**1995: 43%** (IBM), **2004: 15.2%** (IBM), **2016: 8%** (IBM), **6.3%** (Microsoft)

**2017: 5.5%** (IBM) **5.1%** (Microsoft/IBM)

**Human: 5.9%** WER (single) **5.1% WER** (multiple pro transcribers)

AM: CNN-BLSTM, LM: entire history of a dialog session

# Superhuman? Paralings.

- **Speech Analysis (CP): Objective Tasks**

Alcohol Intoxication

16 speakers from ALC, 47 listeners: **71.7%** UAR (human)

Interspeech 2011 Challenge full ALC: **72.2%** UAR (system fusion)

Agglomeration (Weninger et al. 2011) >80%

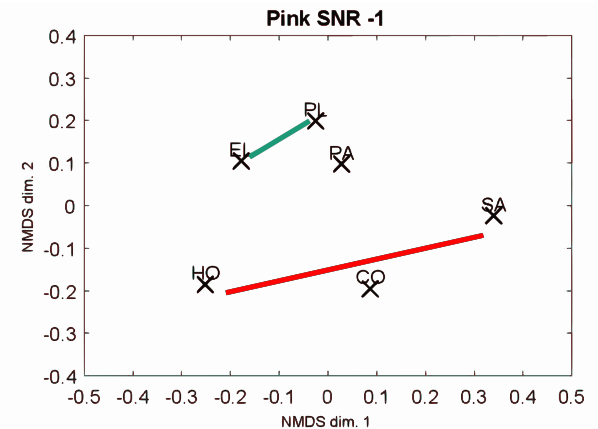
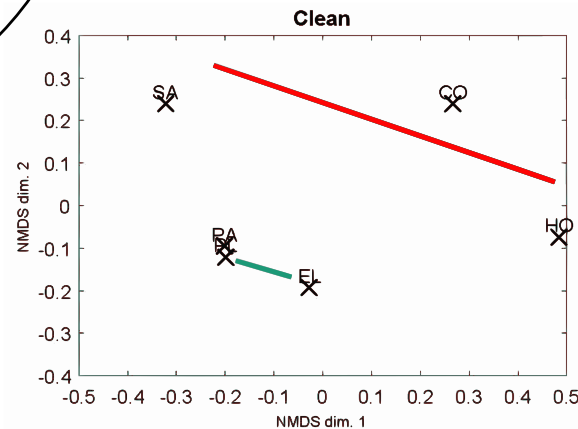
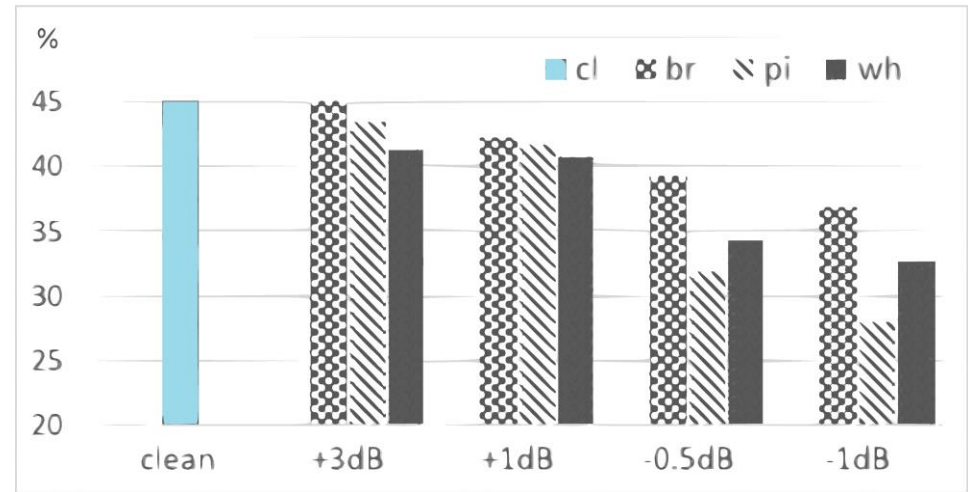
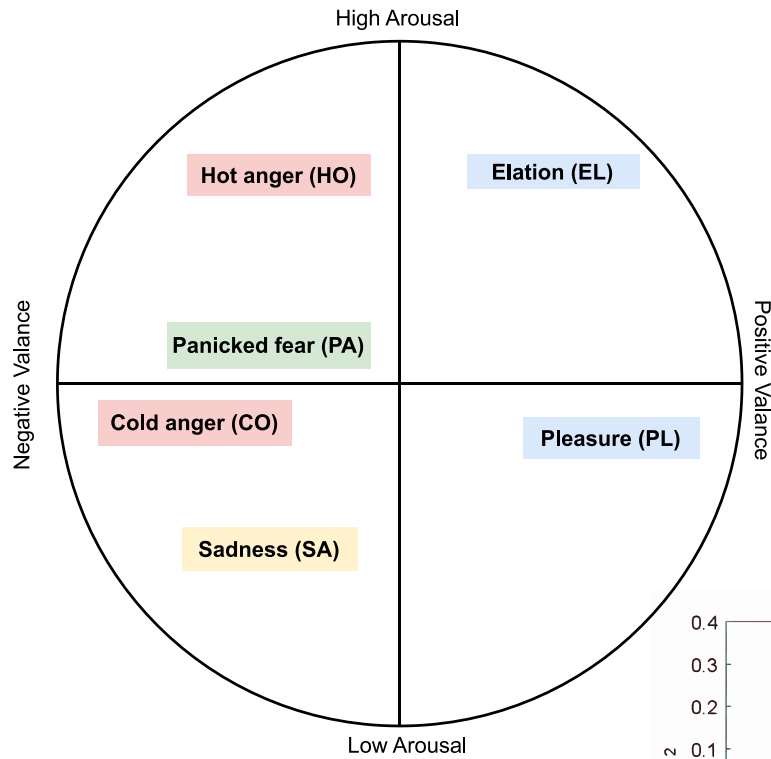
Heart Rate, Skin Conductance, Health State, ...

- **Speech Analysis (CP): Subjective Tasks**

Ground Truth?

Emotion, Personality, Likability, ...?

# Human Performance?



“The Perception of noisified non-sense speech in the noise”, *Interspeech*, 2017.

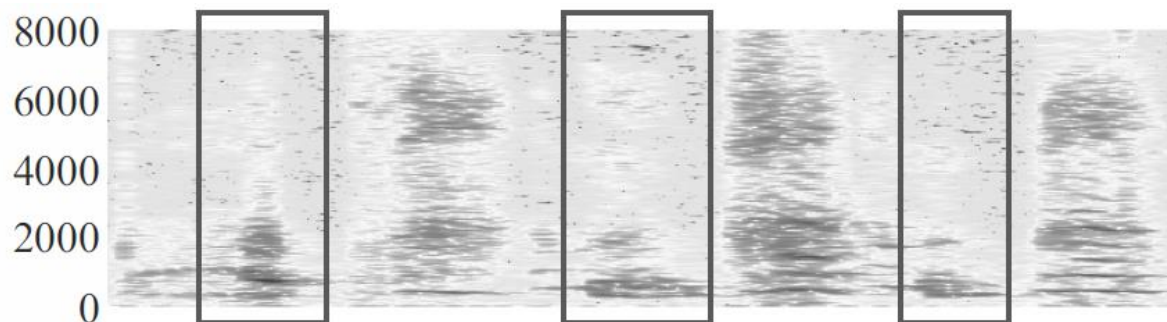
# Rett & ASC.

- **Rett & ASC Early Diagnosis**

16 hours of home videos

6-12 / 10 months

Vocal cues: e.g., inspiratory vocalisation



	%UA
Rett Syndrome	76.5
ASC	75.0



*"A Novel Way to Measure and Predict Development: A Heuristic Approach to Facilitate the Early Detection of Neurodevelopmental Disorders", Current Neurology and Neuroscience Reports, 2017.*

*"Earlier Identification of Children with Autism Spectrum Disorder: An Automatic Vocalisation-based Approach", Interspeech, 2017.*

Getting Broader.

## Speaker ID & Verification

Speech Recognition

Language Understanding

Deep Paralings

Sentiment Analysis

**Speech  
Analysis**

Gender Recognition

Broad Paralings

Emotion Recognition

Language ID

Health Classification

Speaker Diarisation

Personality Recognition



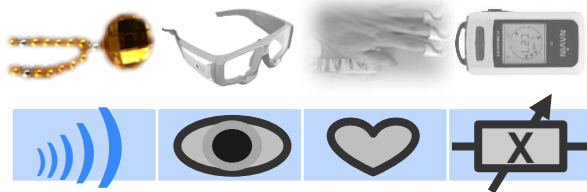
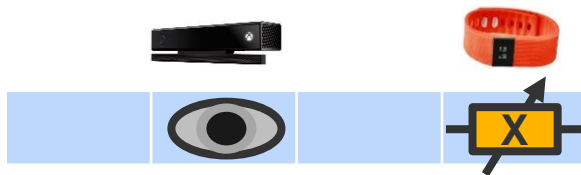
# Paralings.

## INTERSPEECH COMPARE

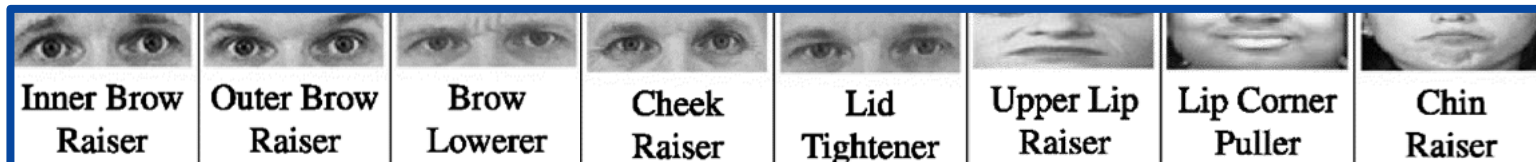
		# Classes	%UA/*AUC/+CC
	Addressee	2	70.6
	Cold	2	72.0
	Snoring	4	70.5
	Deception	2	72.1
	Sincerity	[0,1]	65.4+
	2018	11	82.2
	Affect: Atypical	[-1,1] ?	[0,1] 43.3+
	Affect: Self-Ass.	[-1,1] ?	[0,100] 54.0+
	Crying	3 ?	7 62.7
	Heart Beats	3 ?	3 61.6
	Age	Physical Load	2 71.9
	Gender	Social Signals	2x2 92.7*
	Interest	Conflict	2 85.9
	Emotion	Emotion	12 46.1
	Negativity	Autism	4 69.4

# Broad Paralings.

- Pseudo Multimodality**

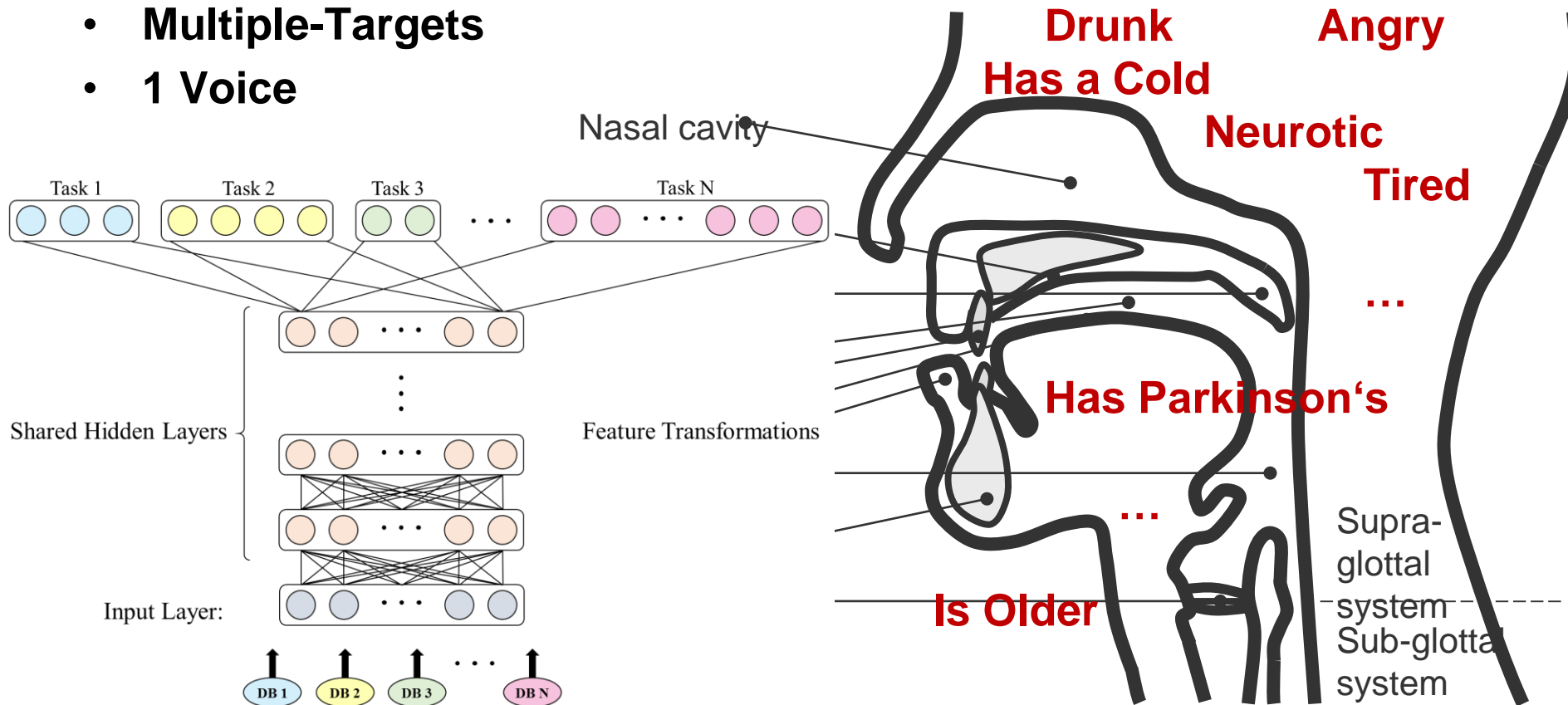


	*MAE
	+CC
	%UA
Heart Rate	8.4*
Skin Conductance	.908+
Facial Action Units	65.0
Eye-Contact	67.4



# Broad Paralings.

- **Multiple-Targets**
- **1 Voice**



# Broad Paralings.

- Cross-Task Self-Labeling**

%UA	Base	CTL
Extraversion	71.7	+1.8
Agreeableness	58.6	+4.5
Neuroticism	63.3	+3.0
Likability	57.2	+2.9

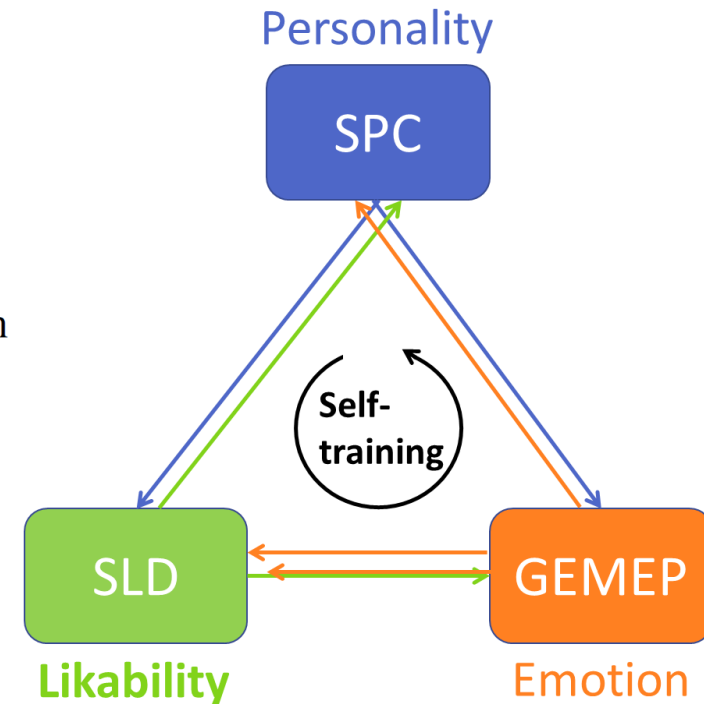
---

**Algorithm:** *Cross-Task Labelling*

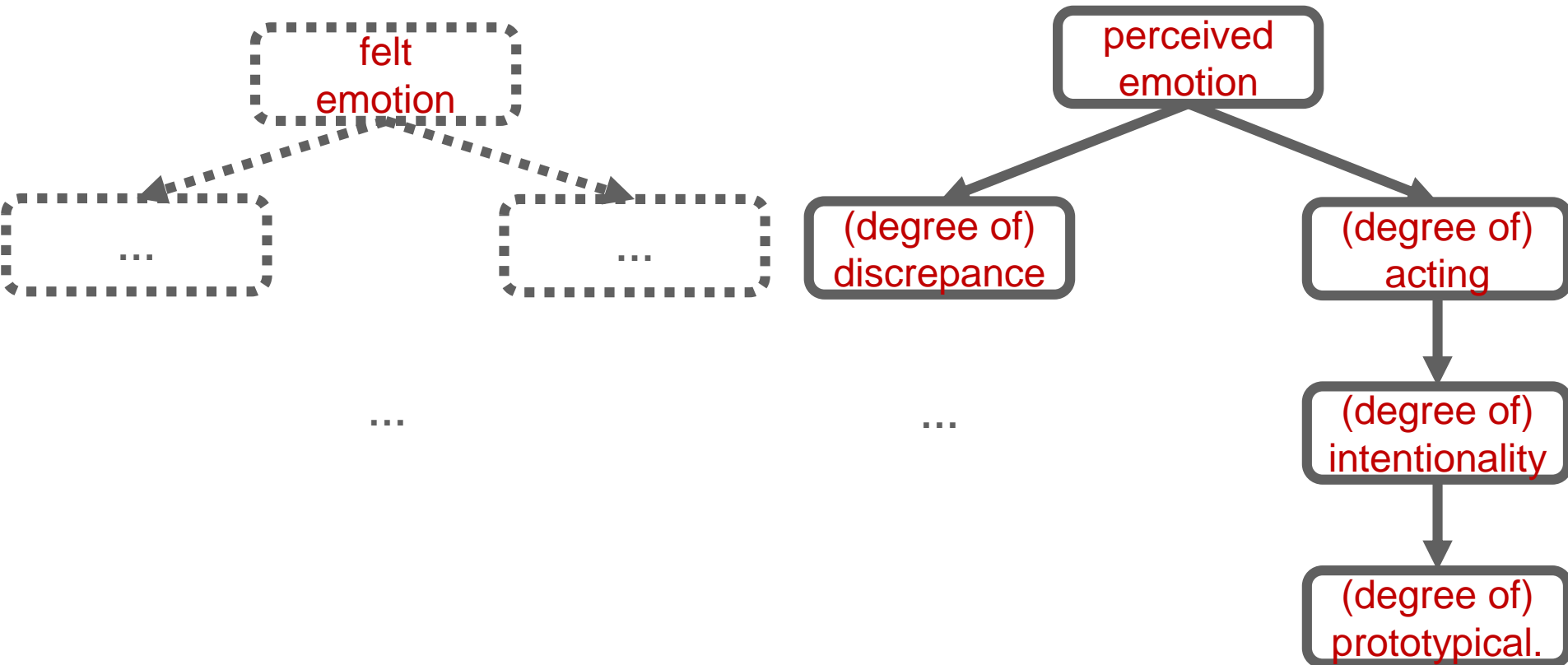
**Repeat for each task:**

**Repeat until  $\mathcal{U} \in \{\}$ :**

- (Optional) Upsample training set  $\mathcal{L}$  to even class distribution  $\mathcal{L}_D$
  - Use  $\mathcal{L}/\mathcal{L}_D$  to train classifier  $\mathcal{H}$ , then classify  $\mathcal{U}$
  - Select a subset  $\mathcal{N}_{st}$  that contains those instances predicted with the highest confidence values
  - Remove  $\mathcal{N}_{st}$  from the unlabelled set  $\mathcal{U}$ ,  $\mathcal{U} = \mathcal{U} \setminus \mathcal{N}_{st}$
  - Add  $\mathcal{N}_{st}$  to the labelled set  $\mathcal{L}$ ,  $\mathcal{L} = \mathcal{L} \cup \mathcal{N}_{st}$
- 



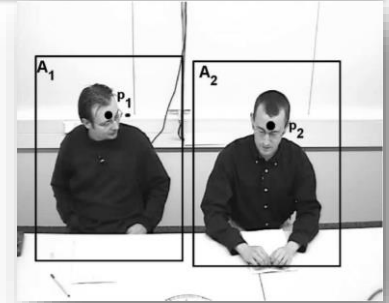
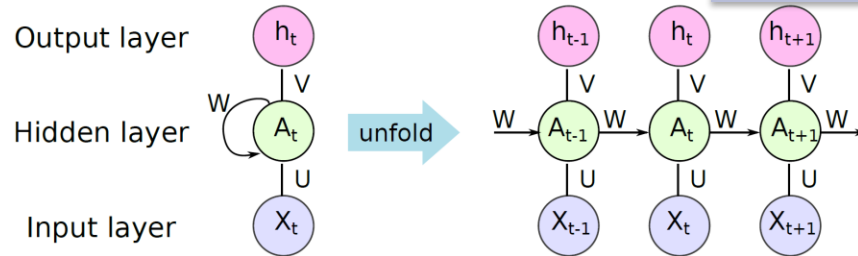
# Deep Paralings.



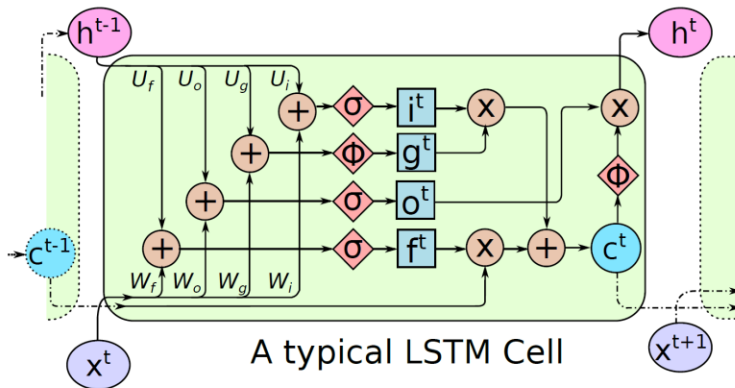
Getting Deeper.

# Deep Recurrent Nets.

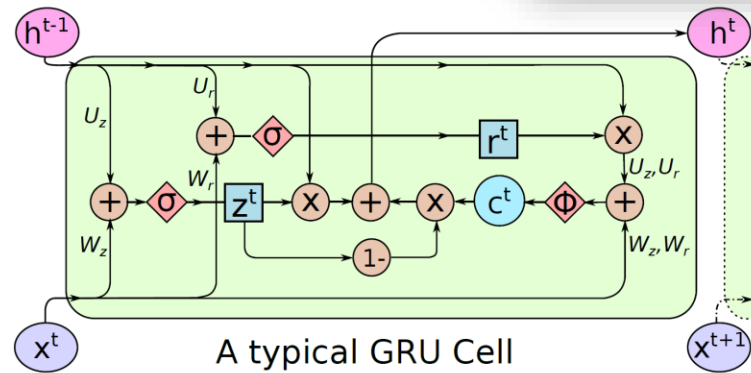
Arousal	CC
HMM	83.5
HMM+LSTM-RNN	87.2
(LSTM-RNN)	96.3



Recurrent neural network unfolded



A typical LSTM Cell



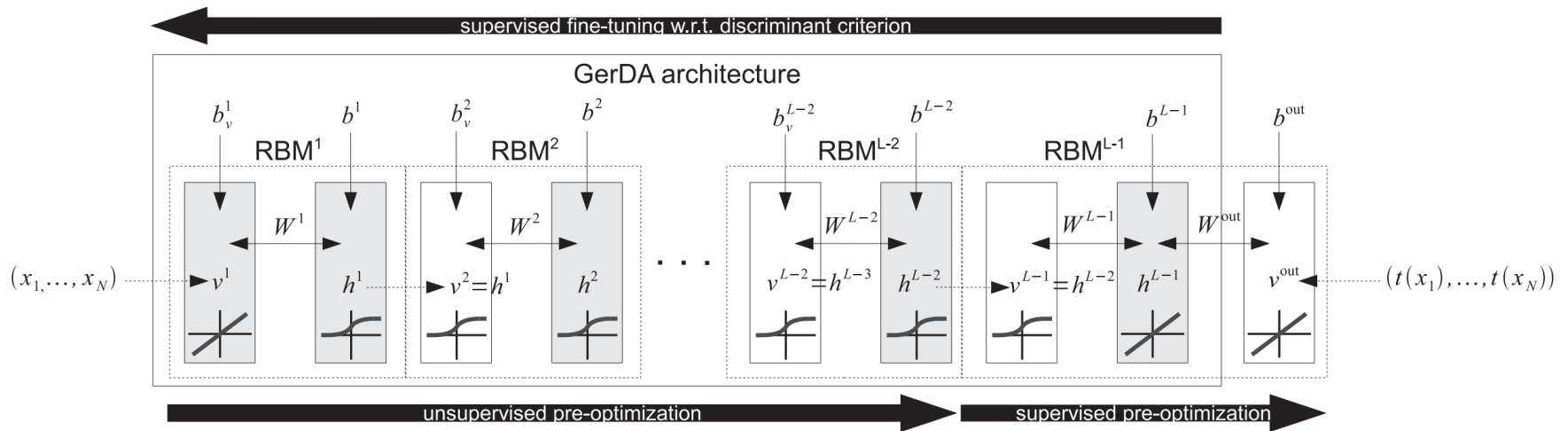
A typical GRU Cell

x Multiplication  
 + Addition  
 1- Subtract from 1  
 f Apply function  $f$   $\xrightarrow{A}$  Multiply by weight A

“A Combined LSTM-RNN-HMM Approach to Meeting Event Segmentation and Recognition”, ICASSP, 2006.

“Abandoning Emotion Classes – Towards Continuous Emotion Recognition with Modelling of Long-Range Dependencies”, Interspeech, 2008.

# Deep Recurrent Nets.

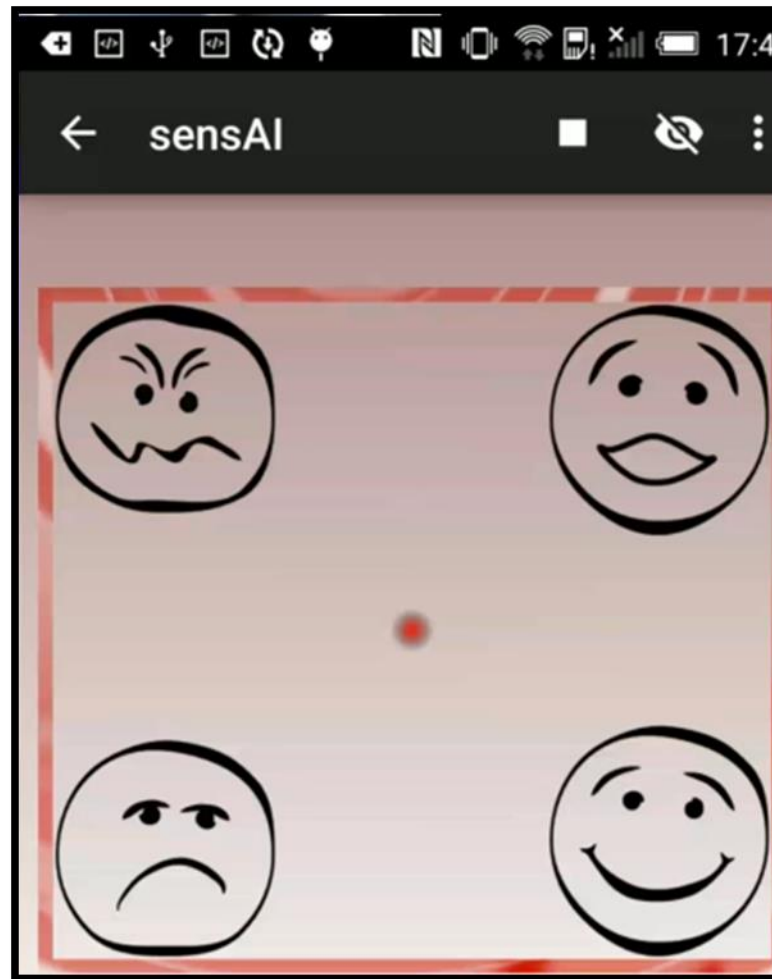


*"Robust discriminative keyword spotting for emotionally colored spontaneous speech using bidirectional LSTM networks", ICASSP, 2009.*

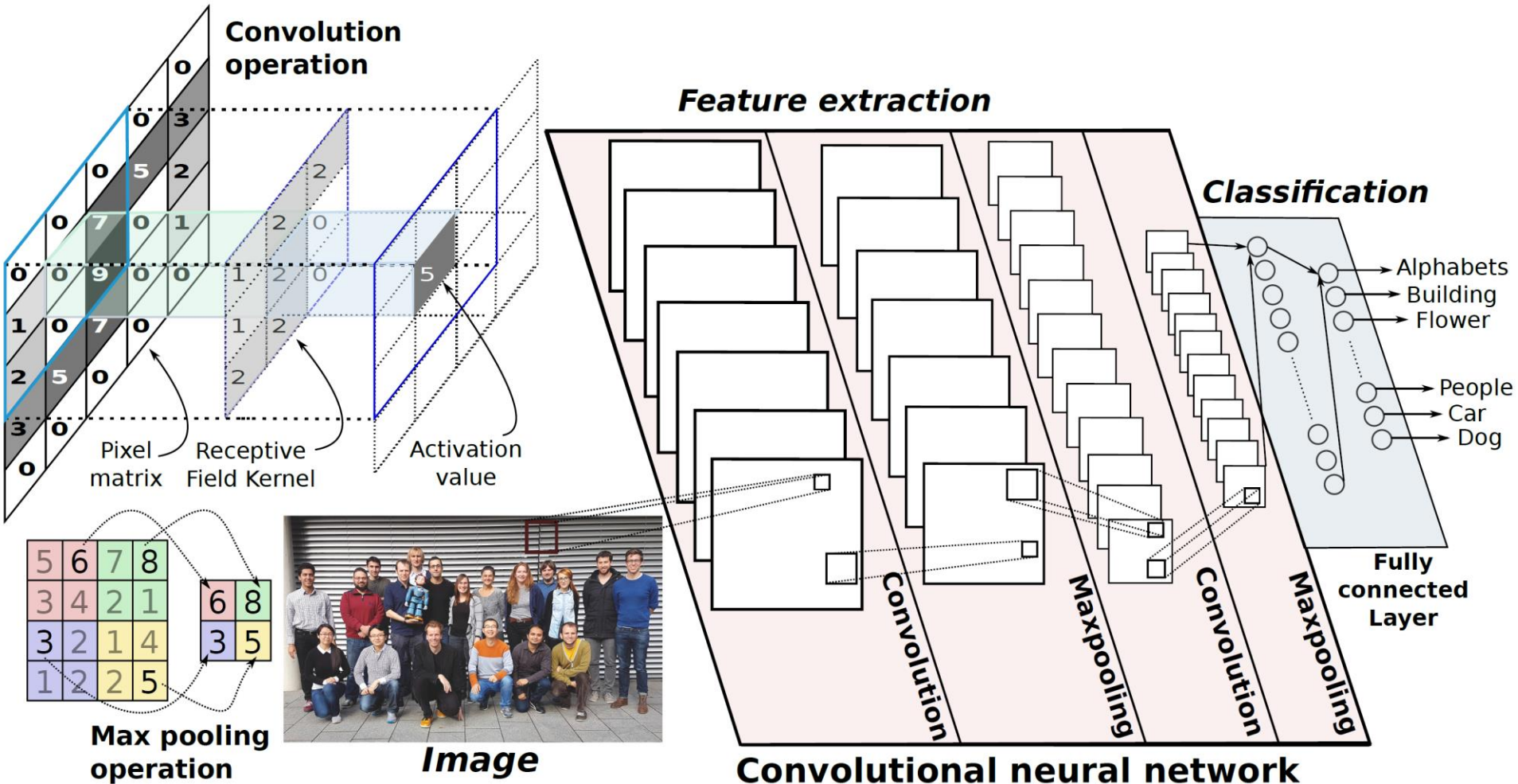
*"Deep neural networks for acoustic emotion recognition: raising the benchmarks", ICASSP, 2011.*



# Deep Recurrent Nets.



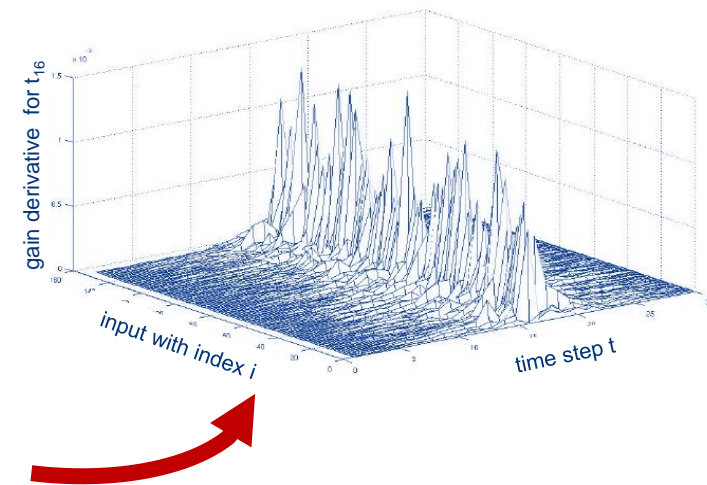
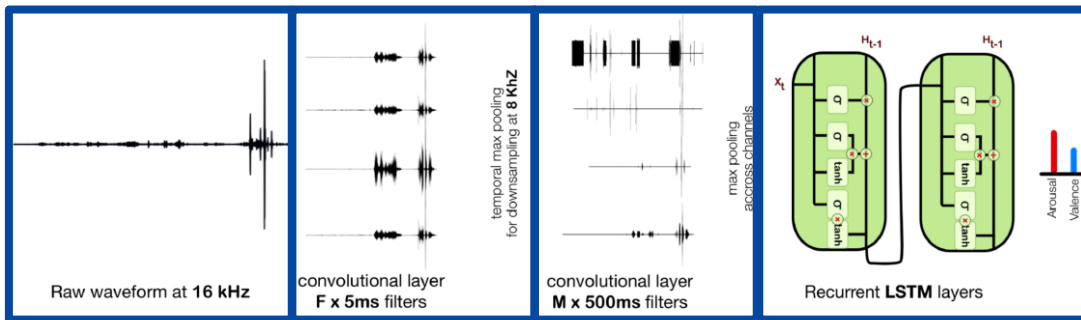
# Convolutional Neural Nets.



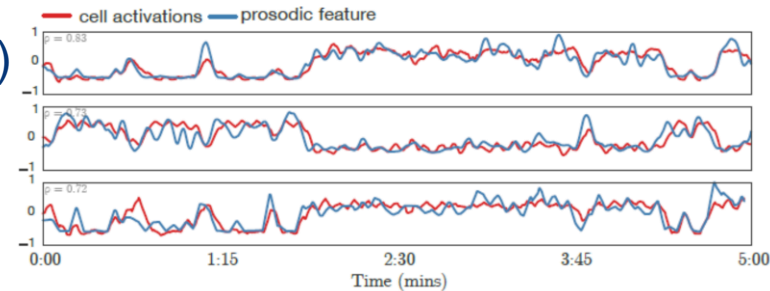
# End-to-End.

- CNN + LSTM RNN

CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261



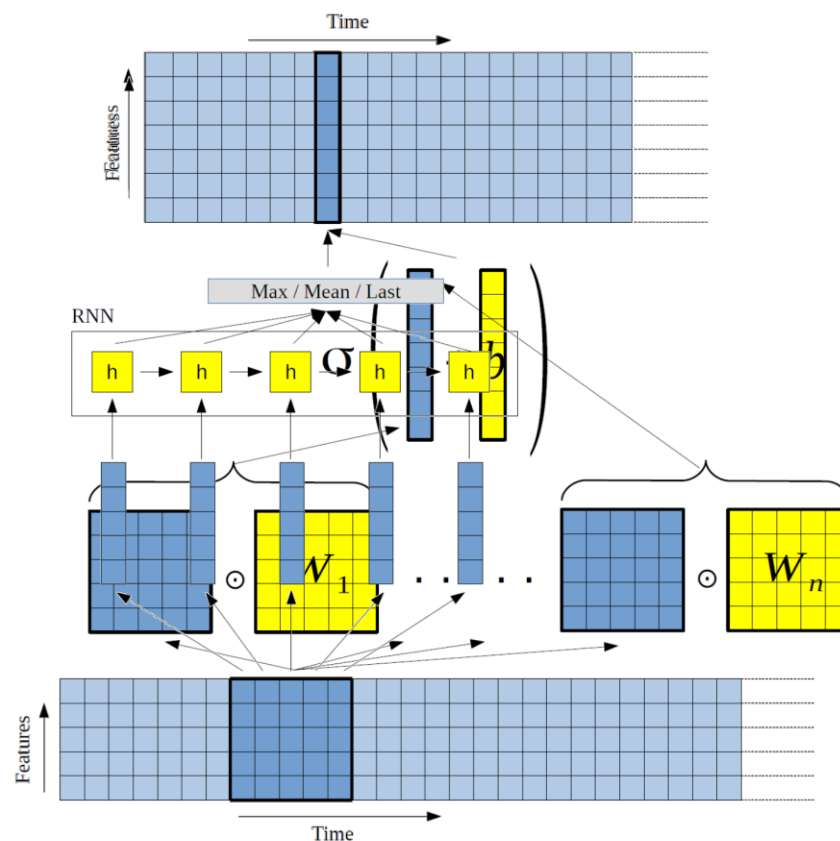
- energy range (.77)
- loudness (.73)
- F0 mean (.71)



# End-to-End.

- CNN + LSTM → CLSTM ?

CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261

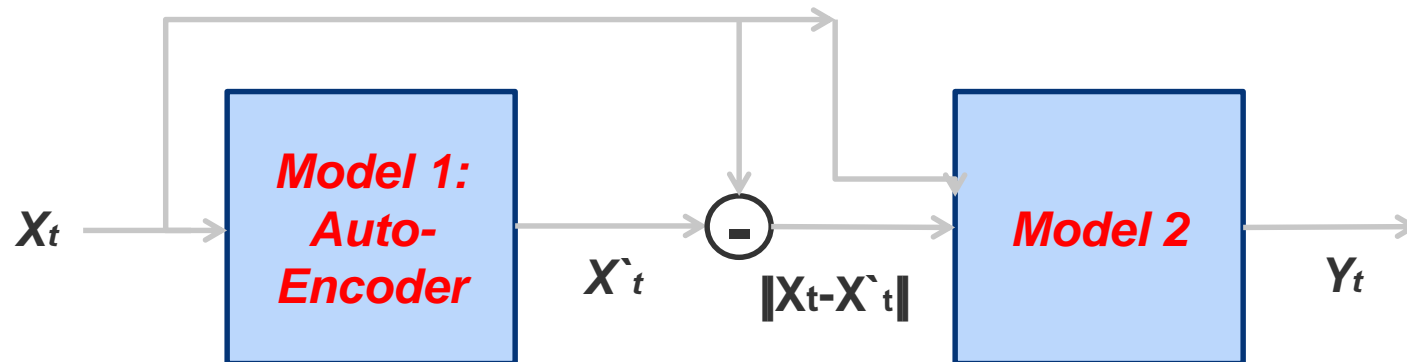


# Learning by Errors.

CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261
Reconstruction Error	.729	.360

- **Reconstruction Error**

RE of Auto-Encoder as additional input feature



Either: Low Level Descriptors (LLD) or Statistical functionals

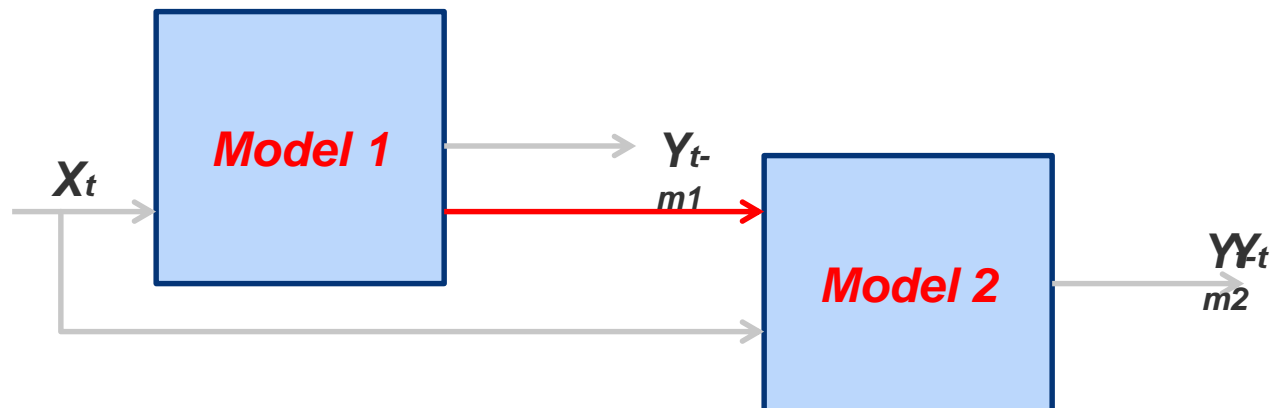
Deep BLSTM RNN

# Prediction-based.

- **Tandem Learning**

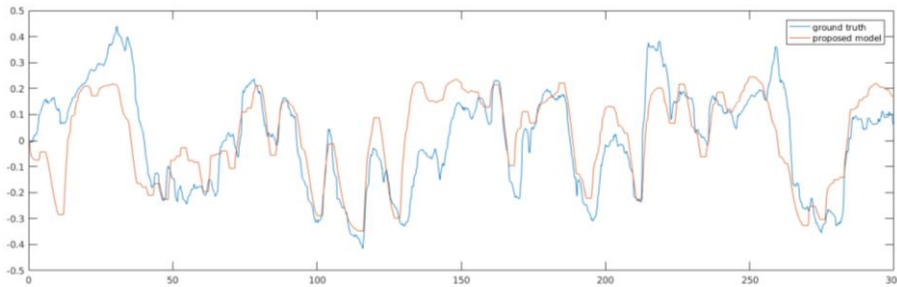
concatenate two models  
for combined strengths

CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261
Reconstruction Error	.729	.360
Prediction-based	.744	.377

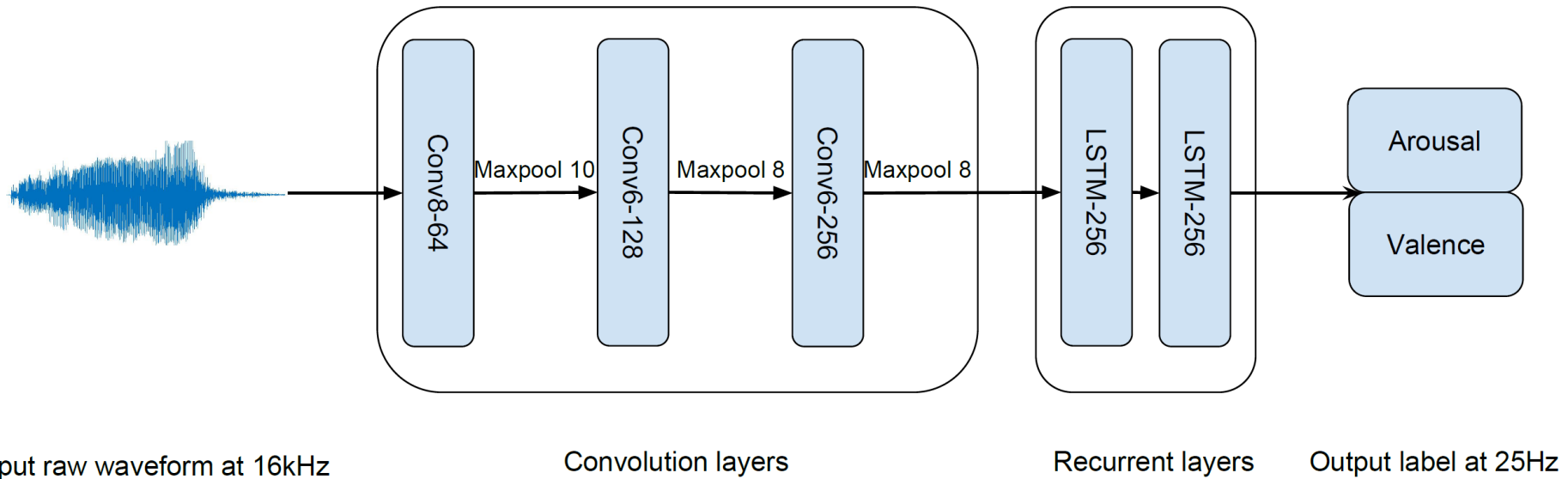


# End-to-End.

- CNN + LSTM RNN**



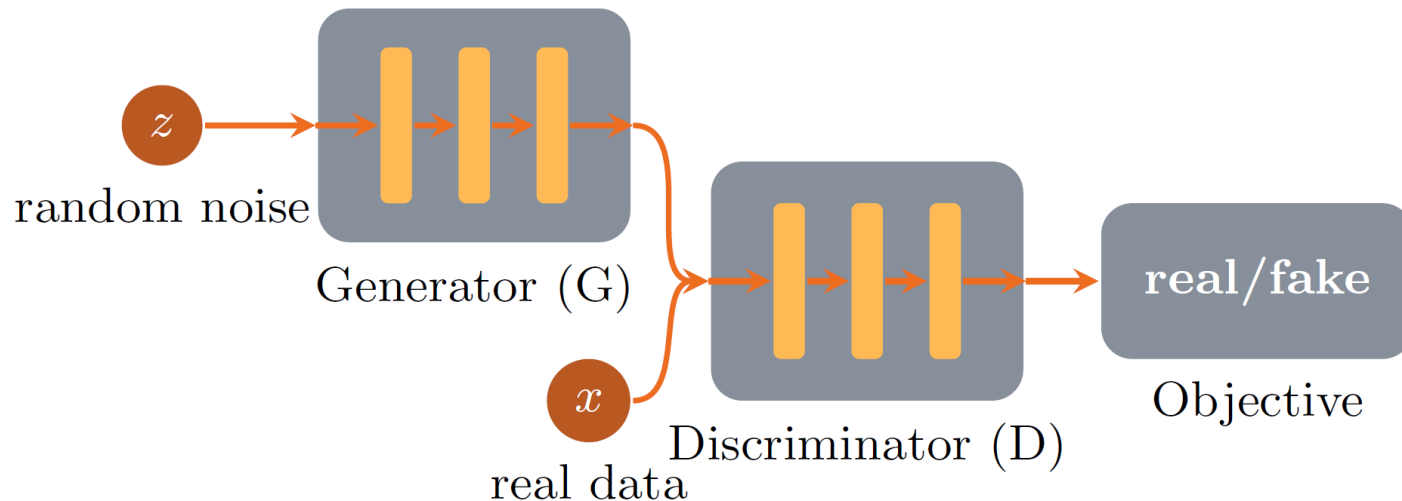
CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261
Reconstruction Error	.729	.360
Prediction-based	.744	.377
BoAW	.753	.430
e2e (submitted)	.787	.440



# Adversarial Nets.

- **Conditional Adversarial Nets**

CCC Recola	Arousal	Valence
ComParE+LSTM	.382	.187
e2e (2016)	.686	.261
Reconstruction Error	.729	.360
Prediction-based	.744	.377
BoAW	.753	.430
e2e (submitted)	.787	.440
CAN (submitted)	.737	.455





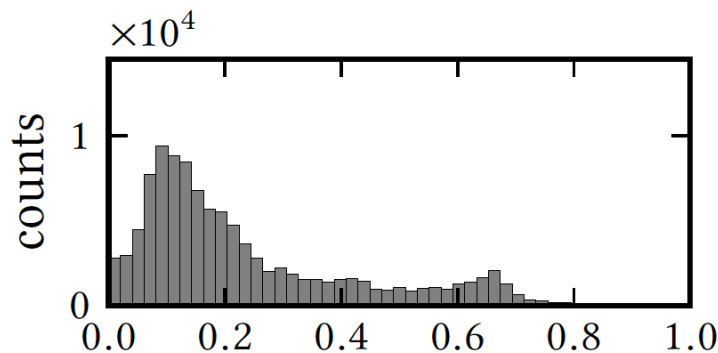
# Co-Learning Trust.

- **Multi-task Learning of Subjective / Uncertain Ground Truth**

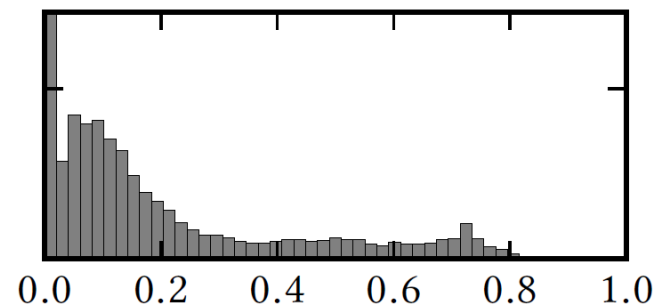
Example: Arousal / Valence (SEWA data of AVEC 2017)

Perception uncertainty (K ratings):

$$\sigma_n^{(i)} = \sqrt{\frac{1}{K-1} \sum_{k=1}^K (e_{n,k}^{(i)} - e_n^{\text{MLE},(i)})^2}$$



(a) arousal



(b) valence

# Co-Learning Trust.

CCC SEWA	Arousal	Valence
Single	.234	.267
Multiple (+conf)	.275	.292
Single (A/V)	.386	.478
Multiple (+conf, A/V)	.450	.515



(a)  $E^{(V)} = 0.08$ ,  $\sigma^{(V)} = 0.17$



(b)  $E^{(V)} = 0.08$ ,  $\sigma^{(V)} = 0.79$



(c)  $E^{(V)} = 0.79$ ,  $\sigma^{(V)} = 0.47$



(d)  $E^{(V)} = 0.69$ ,  $\sigma^{(V)} = 0.70$

# Audio = Images?

CNN+LSTM  
Functionals

%UA  
40.3  
58.8

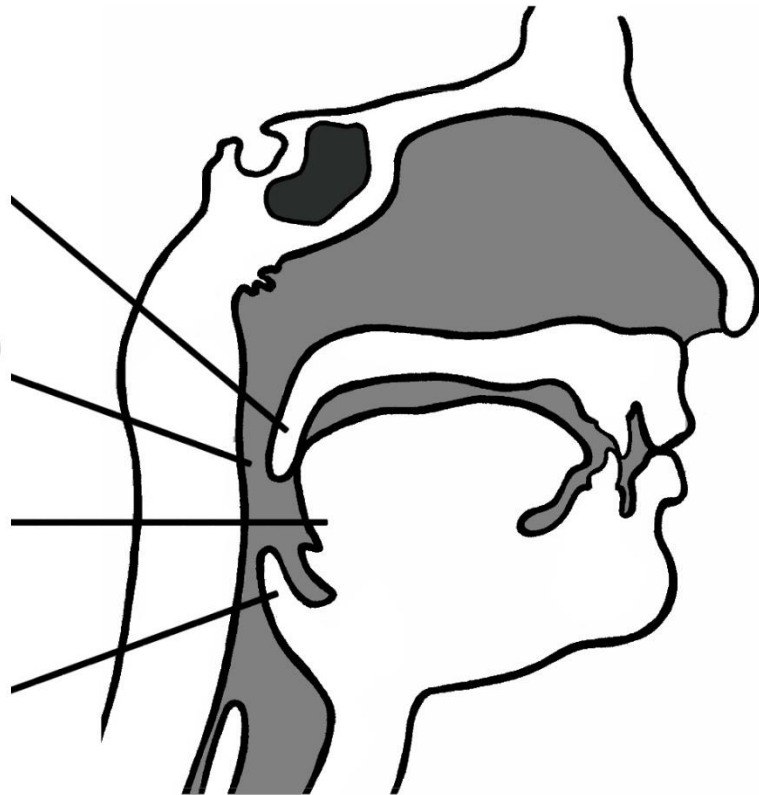
- **VOTE Snoring Classification**

Velum, soft palate **V**

Oropharyngeal **O**

Tongue base **T**

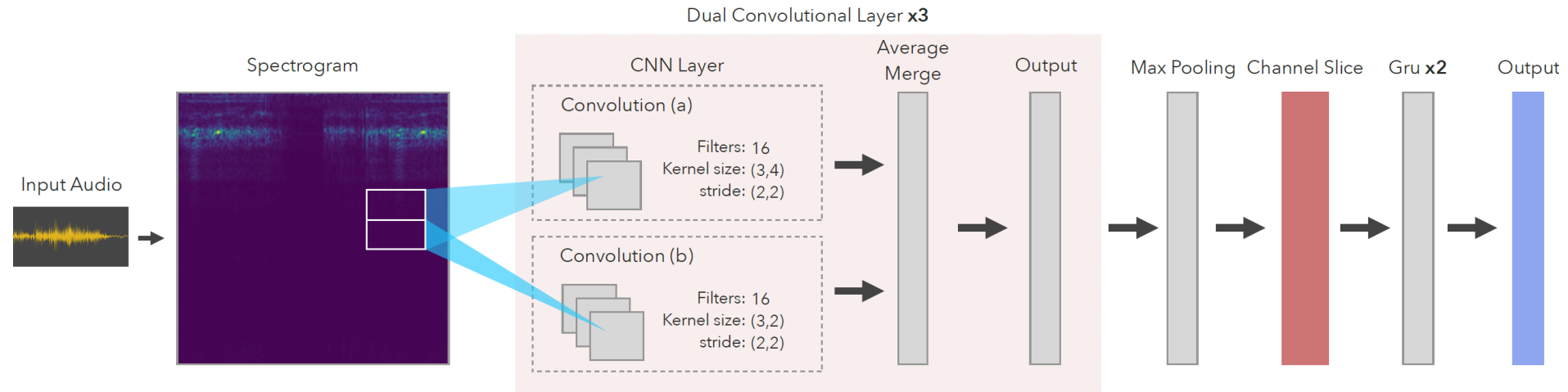
Epiglottis **E**



# Audio = Images?

- VOTE Snoring Classification**

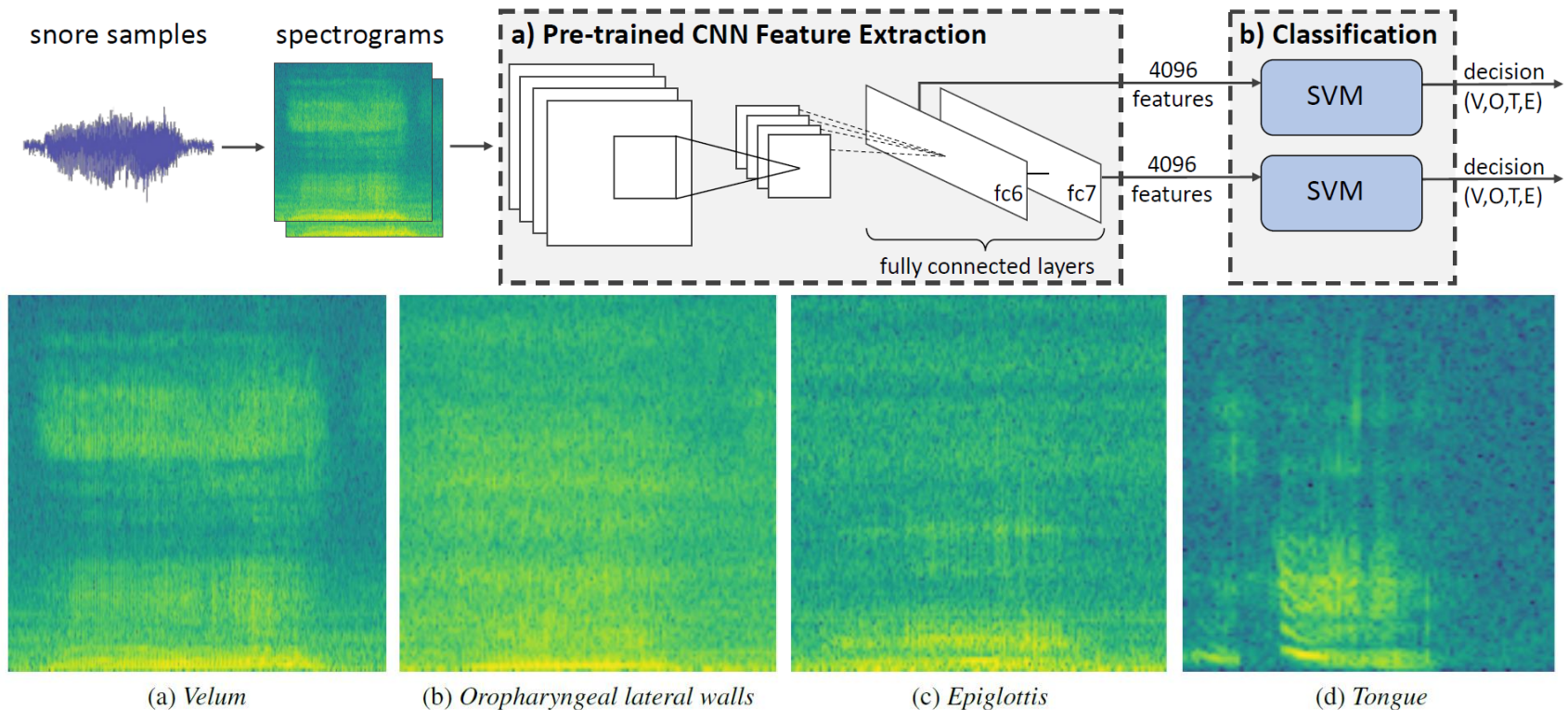
	%UA
CNN+LSTM	40.3
Functionals	58.8
CNN+GRU	63.8



# Audio = Images?

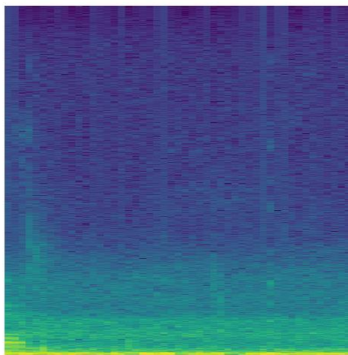
- VOTE Snoring Classification**

	%UA
CNN+LSTM	40.3
Functionals	58.8
CNN+GRU	63.8
Deep Spec	67.0

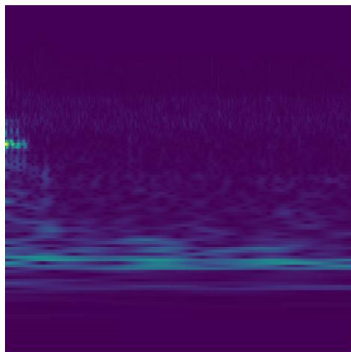


# Audio = Images?

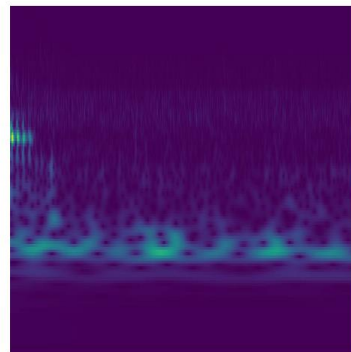
- **Wavelets vs STFT via VGG16**



(a) STFT



(b) bump wavelet



(c) morse wavelet

---

Input: 224×224 RGB image

---

2×conv size: 3; ch: 64  
Maxpooling

---

2×conv size: 3; ch: 128  
Maxpooling

---

3×conv size: 3; ch: 256  
Maxpooling

---

3×conv size: 3; ch: 512  
Maxpooling

---

3×conv size: 3; ch: 512  
Maxpooling

---

Fully connected layer *fc6* with 4096 neurons  
Fully connected layer *fc7* with 4096 neurons  
Fully connected layer with 1000 neurons

---

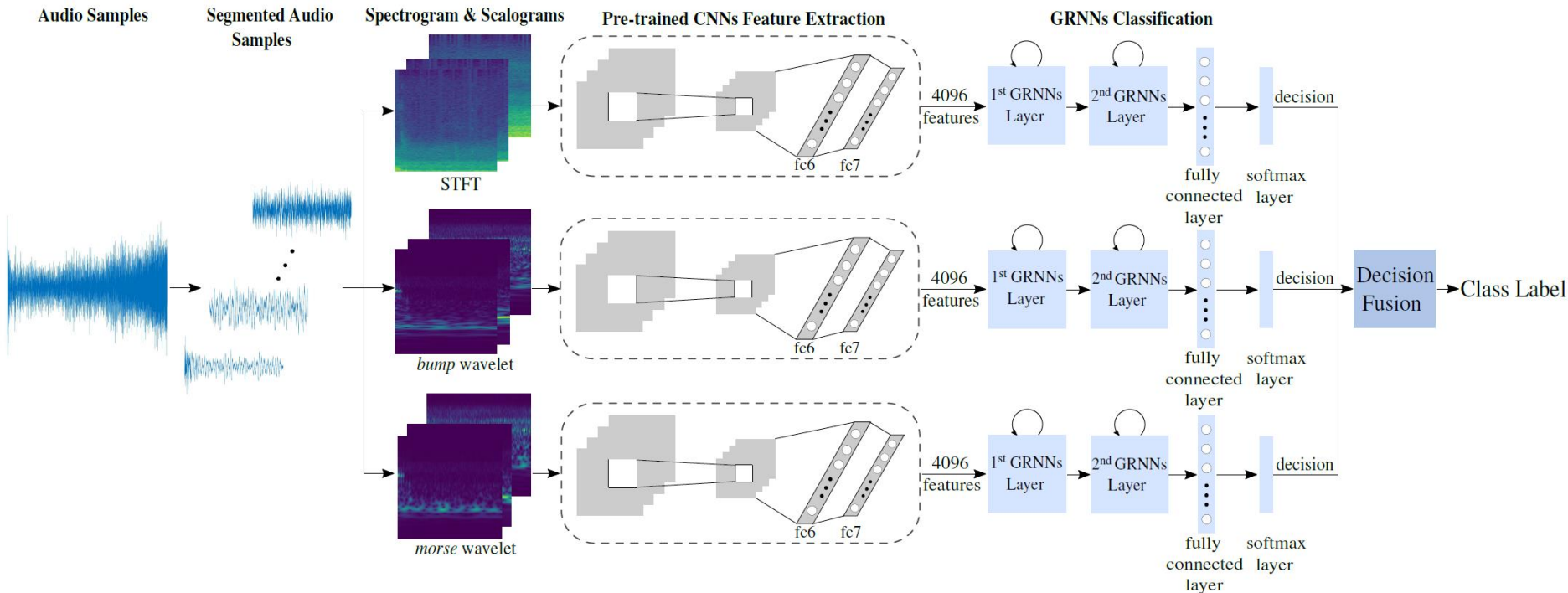
Output: softmax layer of probabilities for 1000 classes

---

# Audio = Images?

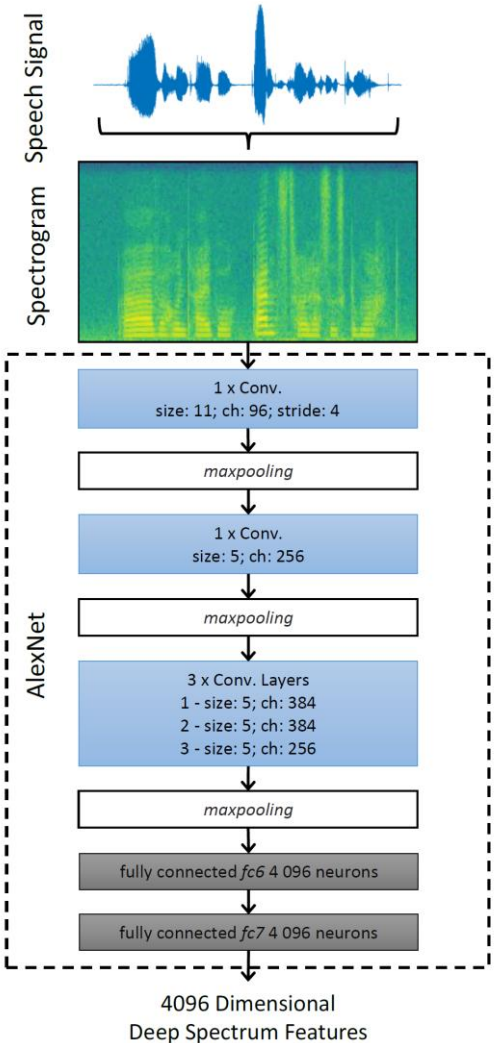
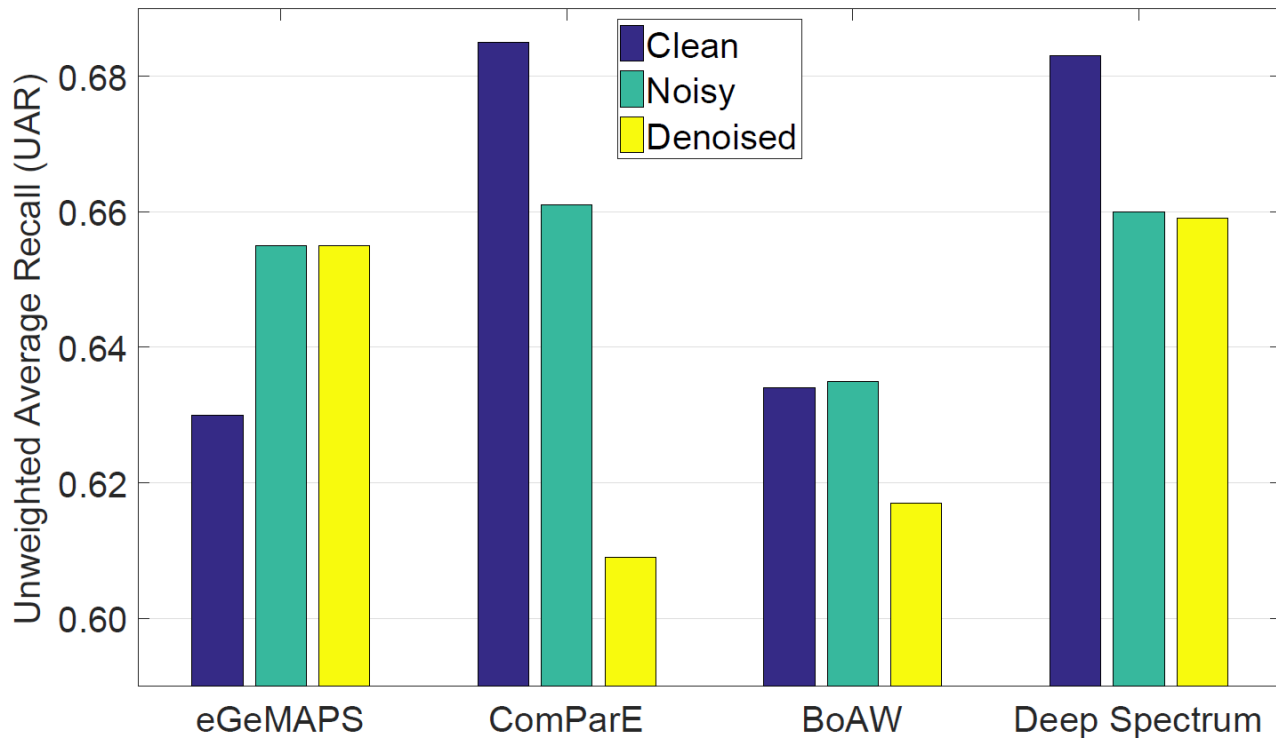
- Wavelets vs STFT via VGG16

DCASE 2017	%WA
STFT	76.5
STFT+bump	79.8
STFT+morse	76.9
All	80.9



# Speech = Images?

- Emotion with Image Nets**  
IS Emotion Challenge task – 2 classes

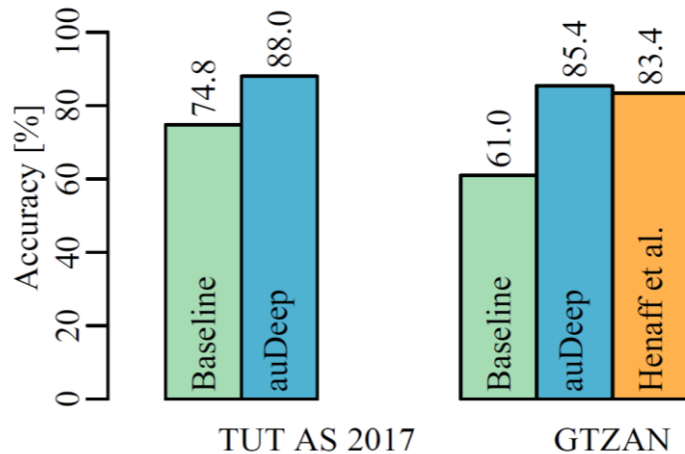
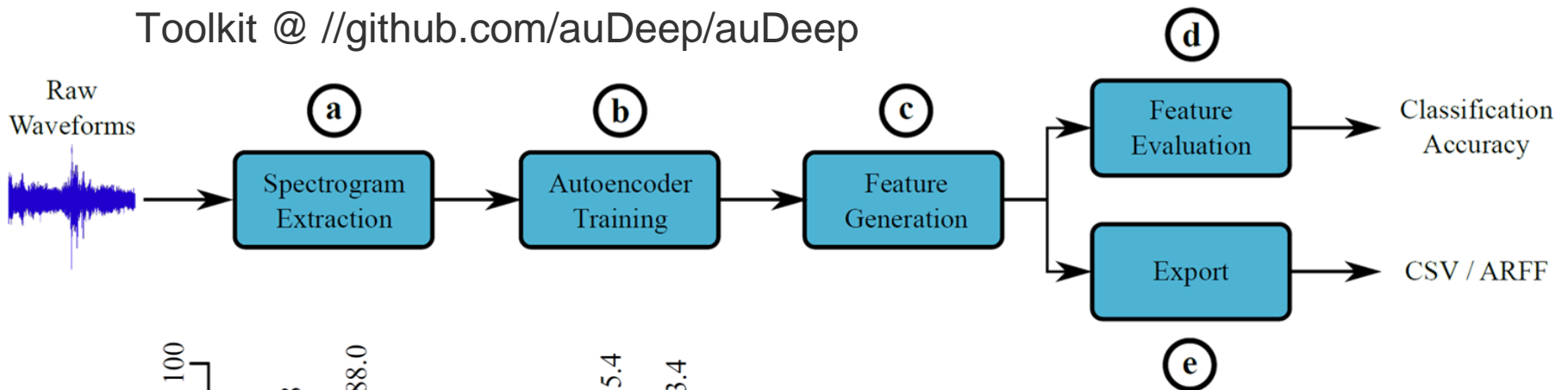




# Speech = Images?

- audDeep**

Toolkit @ [//github.com/auDeep/auDeep](https://github.com/auDeep/auDeep)



Getting Faster.

# Data?

- **2.0 Yet?**



0-1 years:

1 – 100 hrs

**ASA (~10 hrs)**



2-3 years:

~1000 hrs



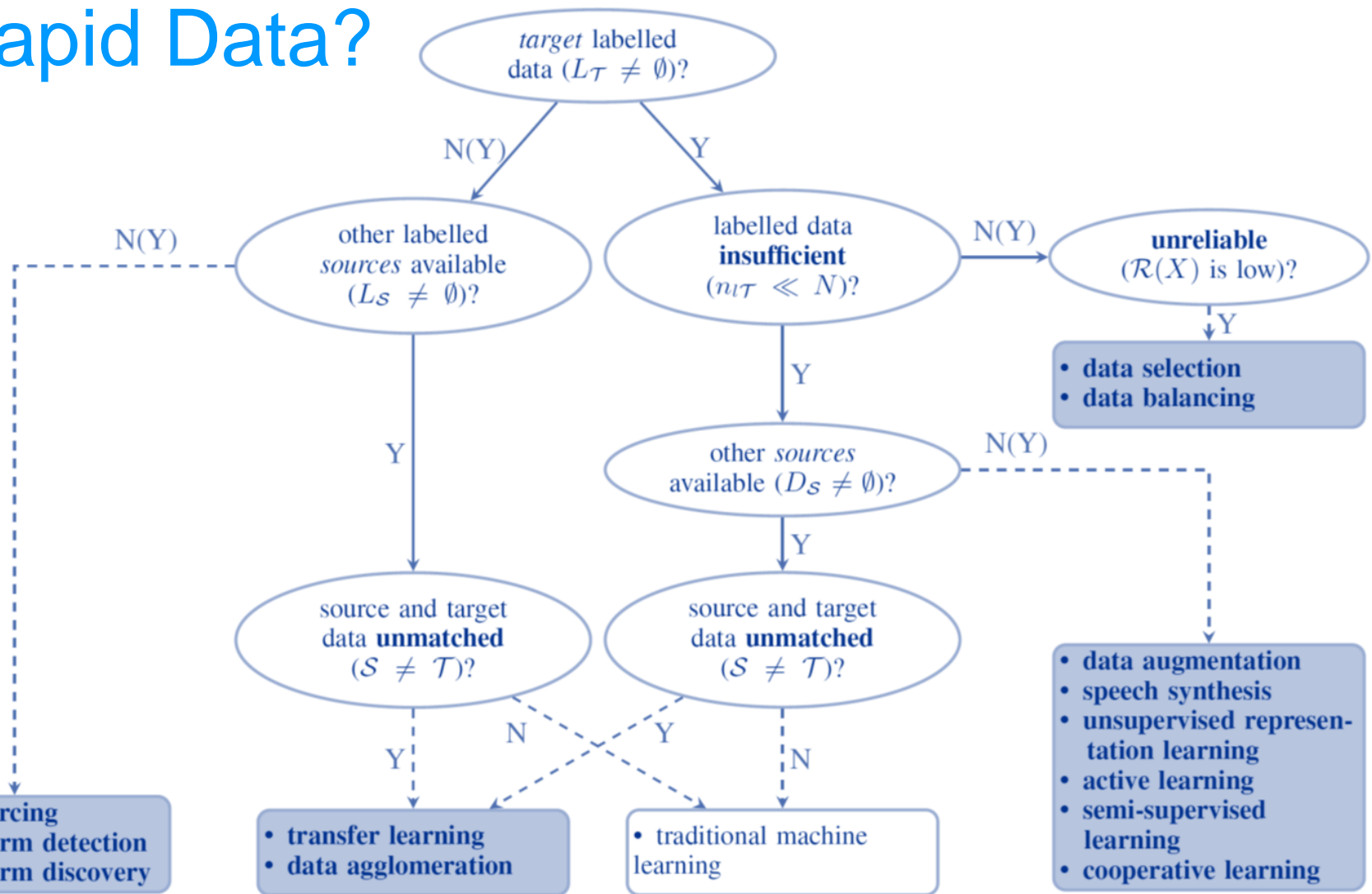
10-x years:

~10000 hrs

**ASR (2000+ hrs)**

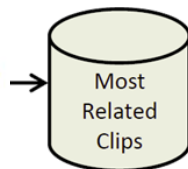
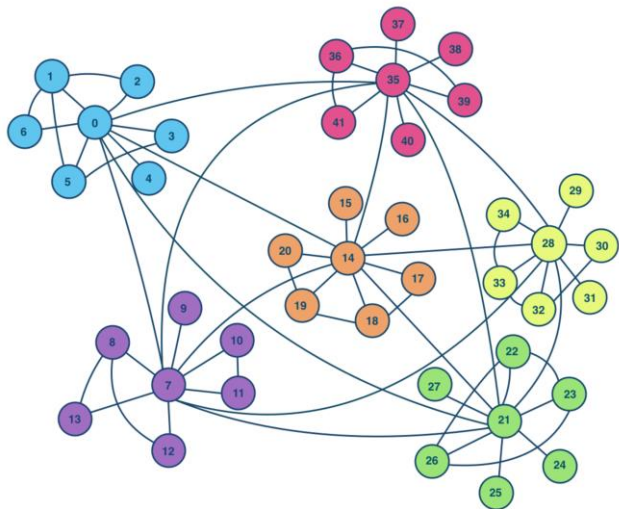
→ Recognise states/traits independent of person, content, language, cultural background, acoustic disturbances at human parity?

# Rapid Data?

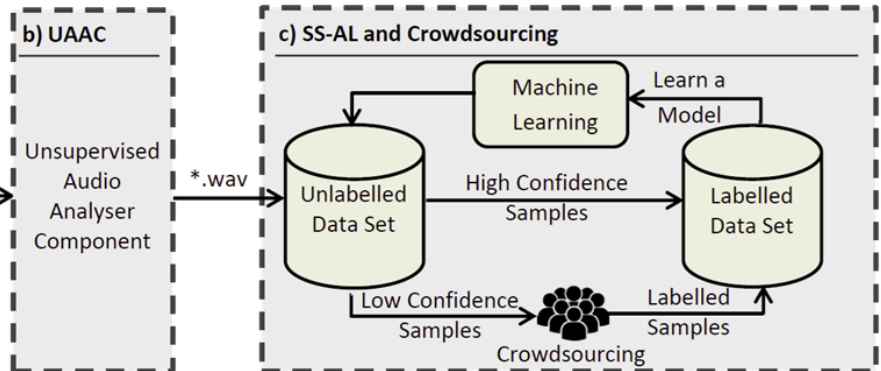


# Rapid Data.

- **YouTube?**  
300 h/min videos  
3k videos for new tasks  
only 3 h/task



%UA	oSMILE	oXBOW	CNN
Freezing	<u>70.2</u>	67.5	57.0
Intoxication	64.7	<u>72.6</u>	66.8
Screaming	89.2	<u>97.0</u>	89.2
Threatening	<u>73.8</u>	67.0	71.9
Coughing	95.4	<u>97.6</u>	95.4
Sneezing	79.2	79.8	<u>85.2</u>

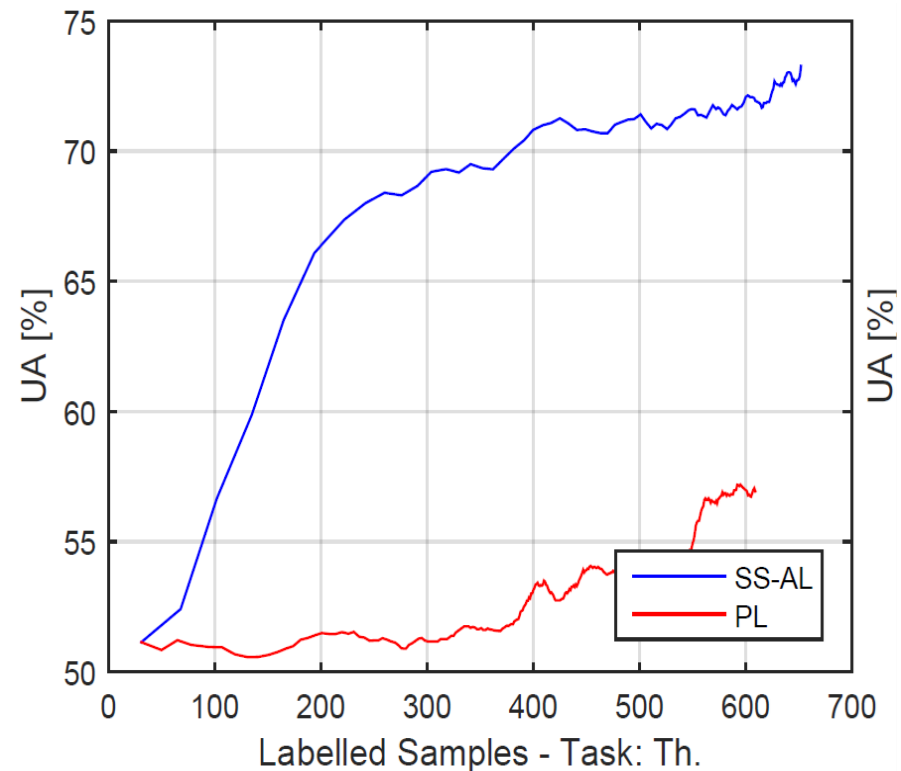
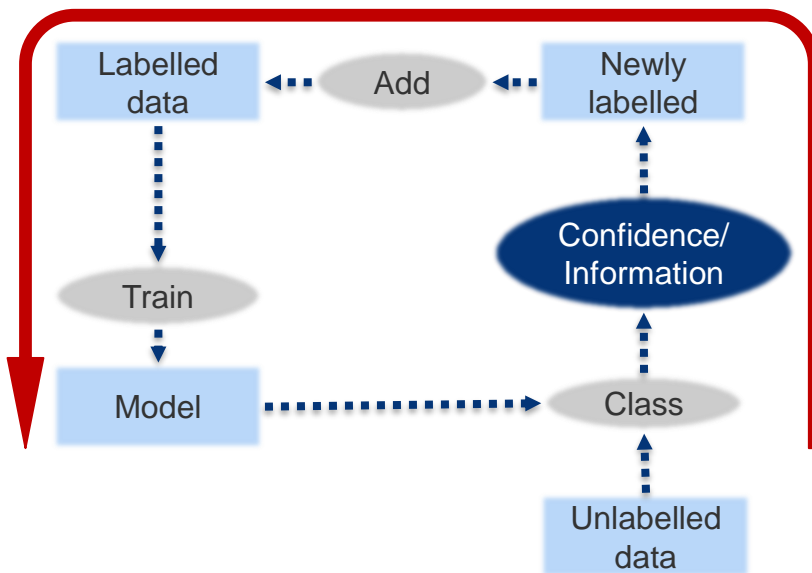


# Rapid Data.



- **Intelligent Labelling**

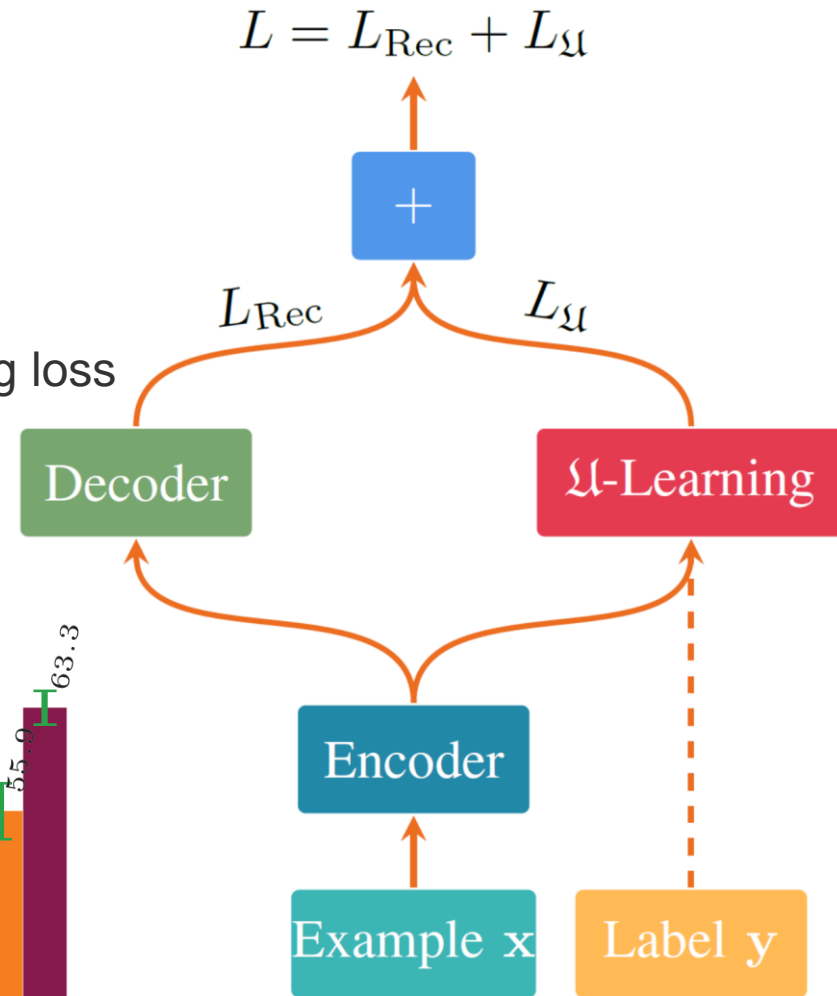
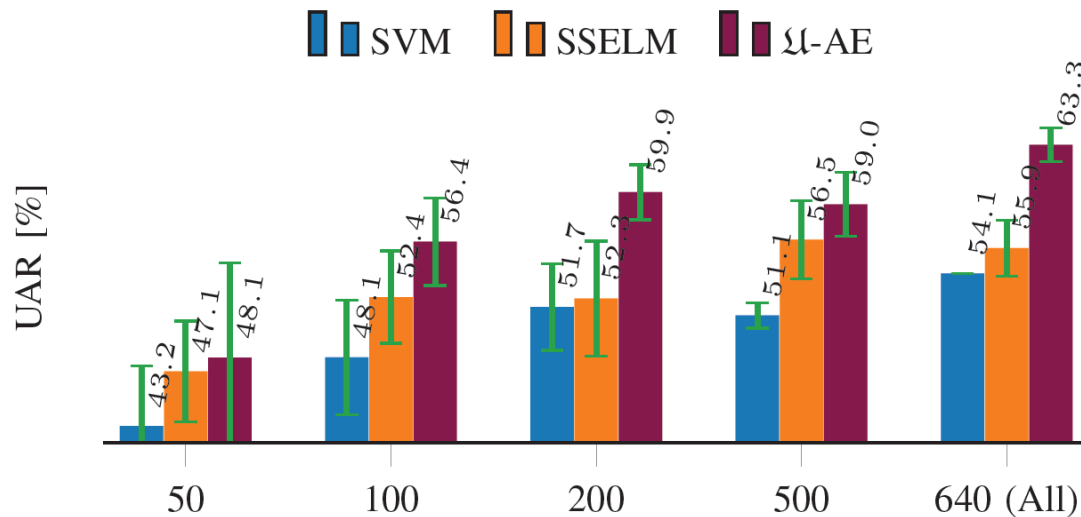
- 0) Transfer Learning
- 1) Dynamic Active Learning
- 2) Semi-Supervised Learning



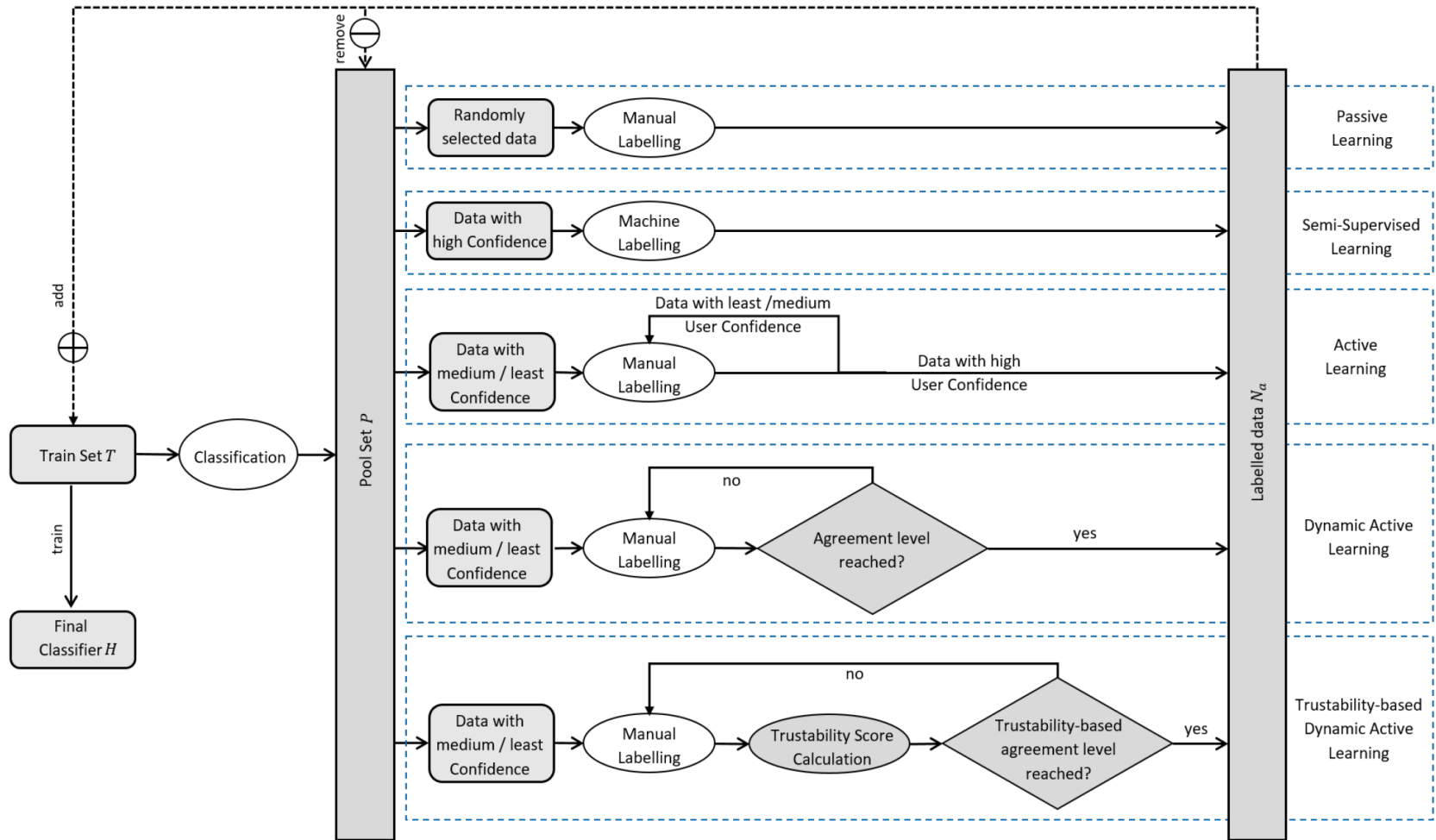
# Rapid Data: TL.

- TL: Universum Autoencoders**

Jointly minimise reconstruction error & universum (unlabelled dataset) learning loss  
Whispered → TRANSFER → normal  
GeWEC (4 class) + Unlabelled: ABC

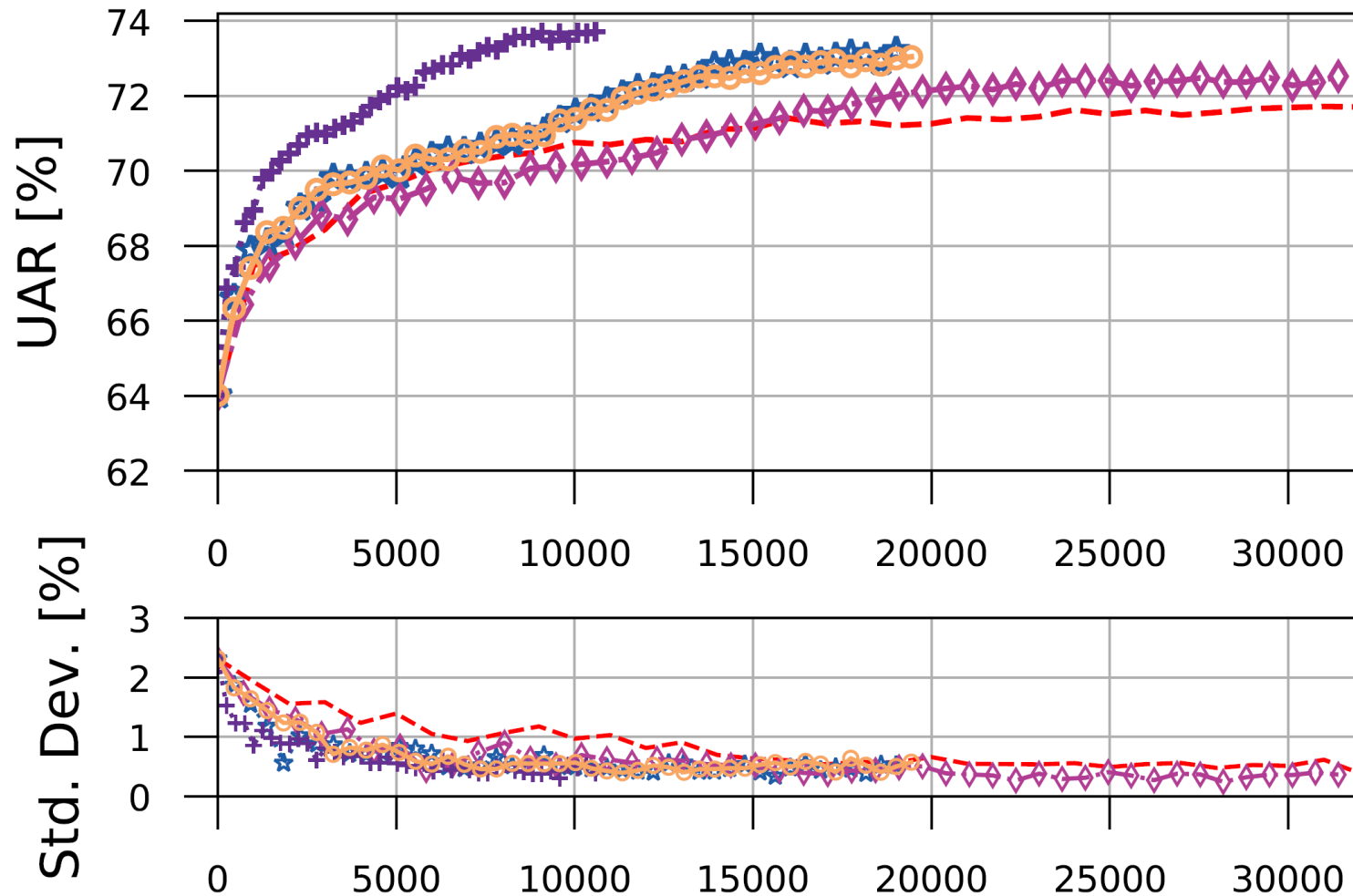


# Rapid Data: AL





# Rapid Data: AL.



# Rapid Data: AL+CS.

THEAR PLAY

Home Play Leaderb

Progress of database: Eating  
16%

The North Wind and the Sun were disputing which was the wrapped in a warm cloak.

[Report a problem](#)

Top players

Last 7 days Last 30 days All time

#	Username	Rank	Gamerscore
1	Maryna	Intermediate	★ 30828
2	max	Intermediate	★ 29848
3	isa	Intermediate	★ 22630
4	zixing	Novice	★ 10100
5	jing	Novice	★ 10092
6	Christoph	Novice	★ 9075
7	Hesy	Beginner	★ 2552
8	Simone	Beginner	★ 2035



list FAQ Contact Your Profile Logout

That answer was okay. Guess. Points earned: 100

Badge Name	Conditions
Early Bird	Answer 100 questions between 00:00 and 06:00
Night Owl	Answer 100 questions between 22:00 and 06:00
Expert	Reach a score of 5000 Points
Master	Reach a score of 20000 Points
Powerman	Collect 100 Bonus Items (in total)
Regular Customer	Have a constant log-in streak of 7 days
Way to go	Answer 100 questions in total
Autobiographer	Fill out own bibliography
Chatterbox (hidden)	Used the contact form 5 times

Personal Multiplier (?:) 3.1  
Answered questions: 1  
awarded at March 21, 2016, 10:01 a.m.

Dataset of the week

ASPA (nativeness)

This dataset is a collection of 30 second excerpts of various scientific talks. Here we would like to know how you would rate the speaker's proficiency of the English language.

[Play this dataset](#)

Alcoholic Samples  
Samples from drunk people

Current Multiplier 1x

Available Audiodata 2

Available Questions 3

Your Progress 9%

# Rapid Data: SSL.

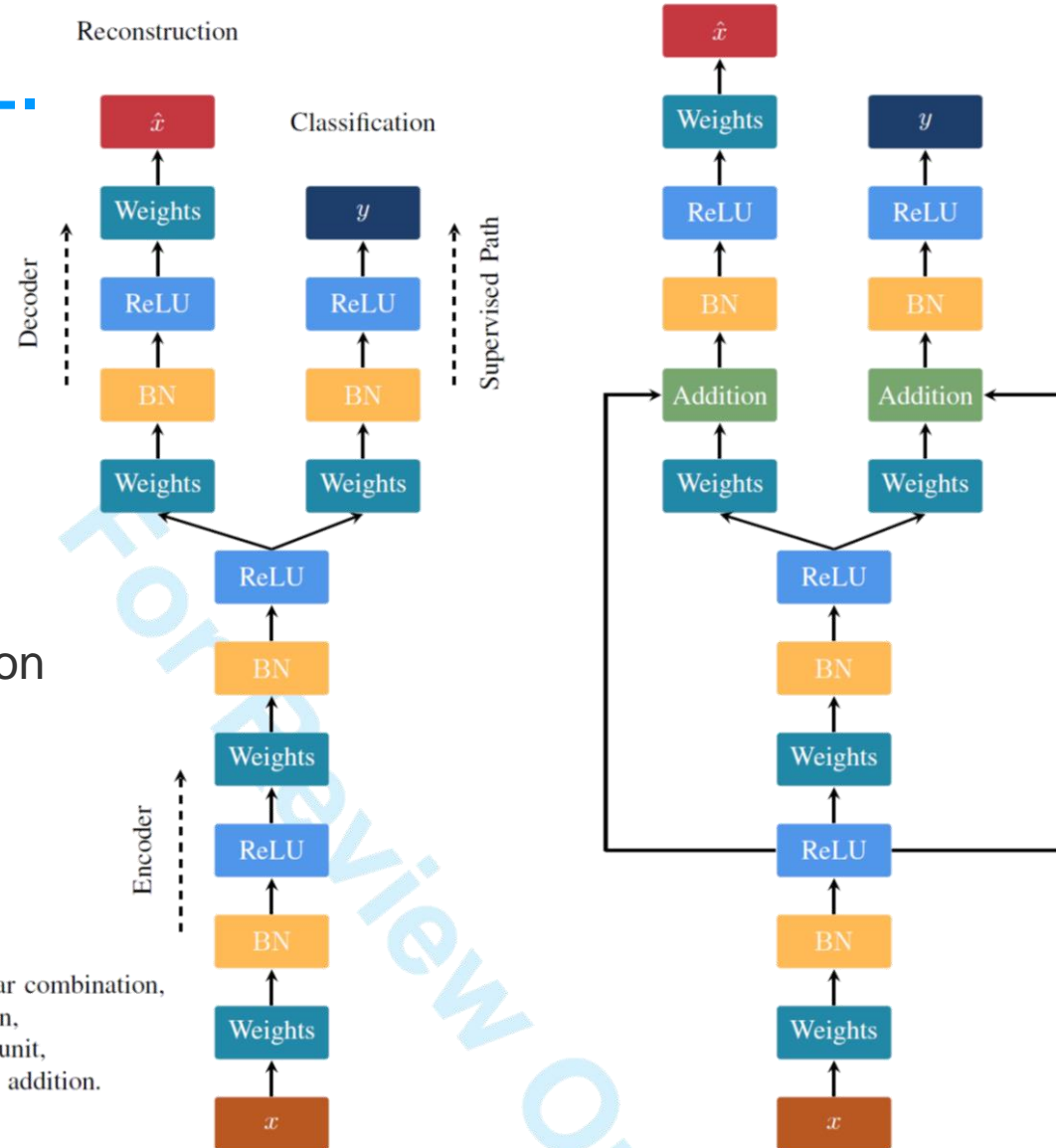
- AEs for SSL**

Supervised Learning:  
Keep only relevant info

Unsupervised AEs:  
Keep all info for reconstruction

w/o (left) or w/ (right)  
skip compensation

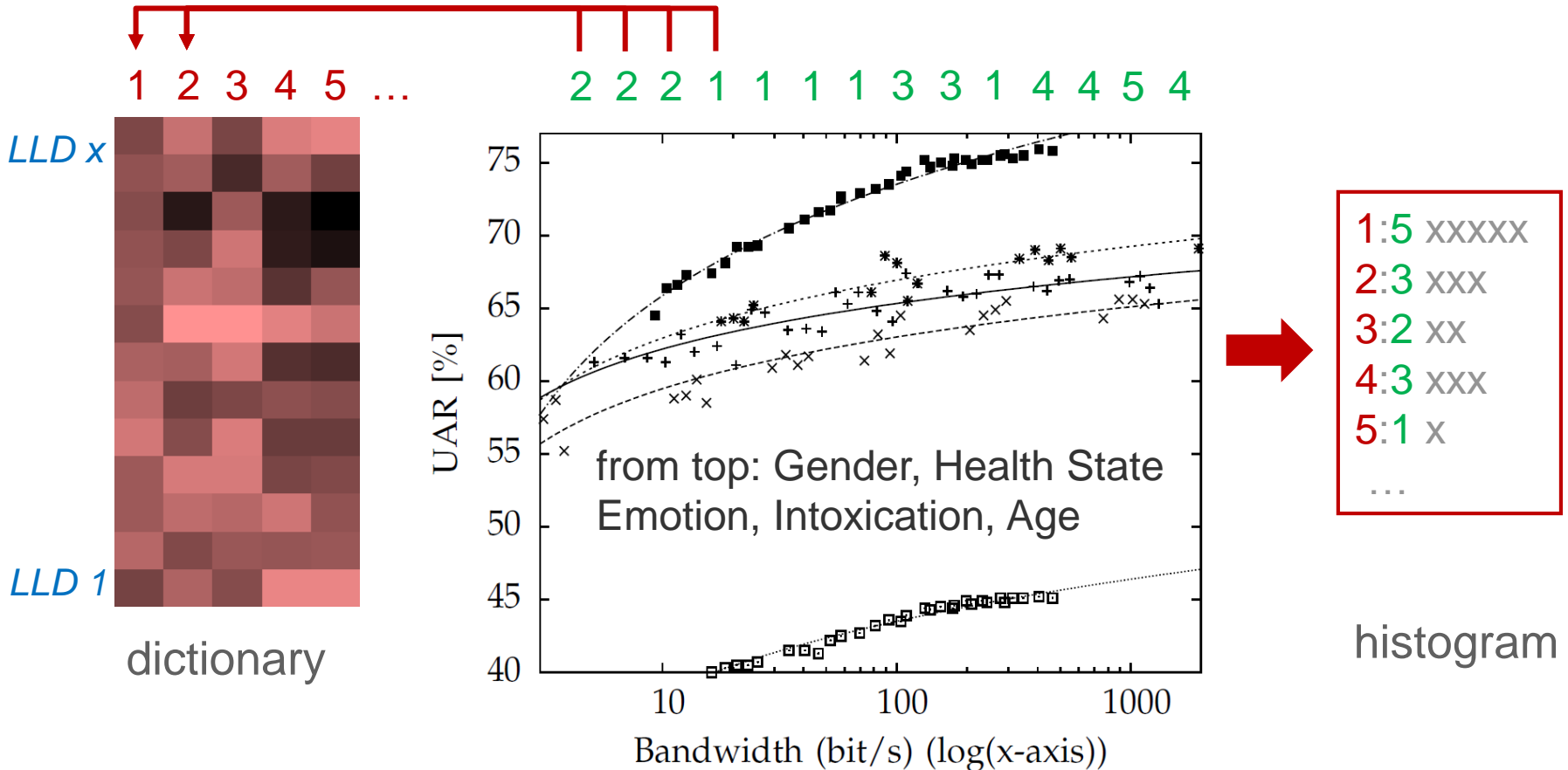
**Weights:** Weighted linear combination,  
**BN:** Batch normalisation,  
**ReLU:** Rectified linear unit,  
**Addition:** Element-wise addition.



# Fast Transmission.

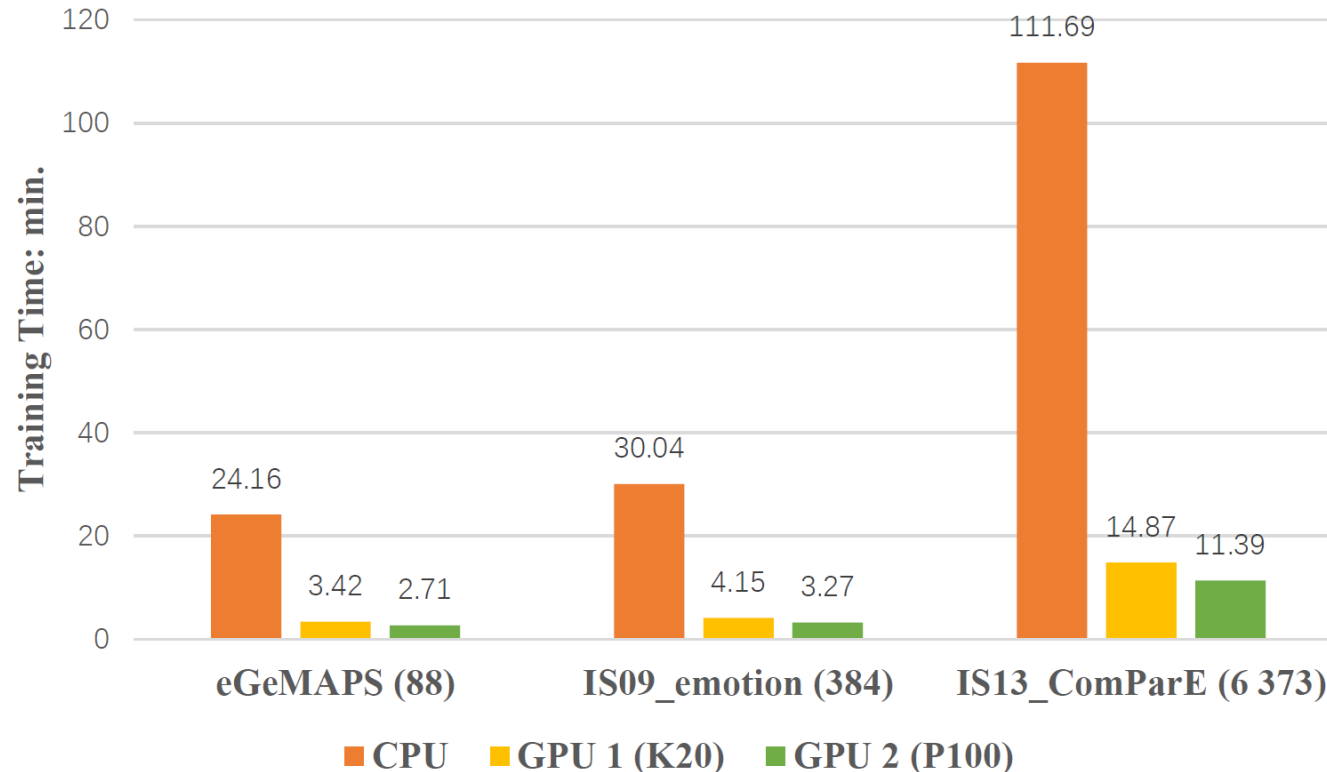
openXBOW -|)→

## vector quantisation



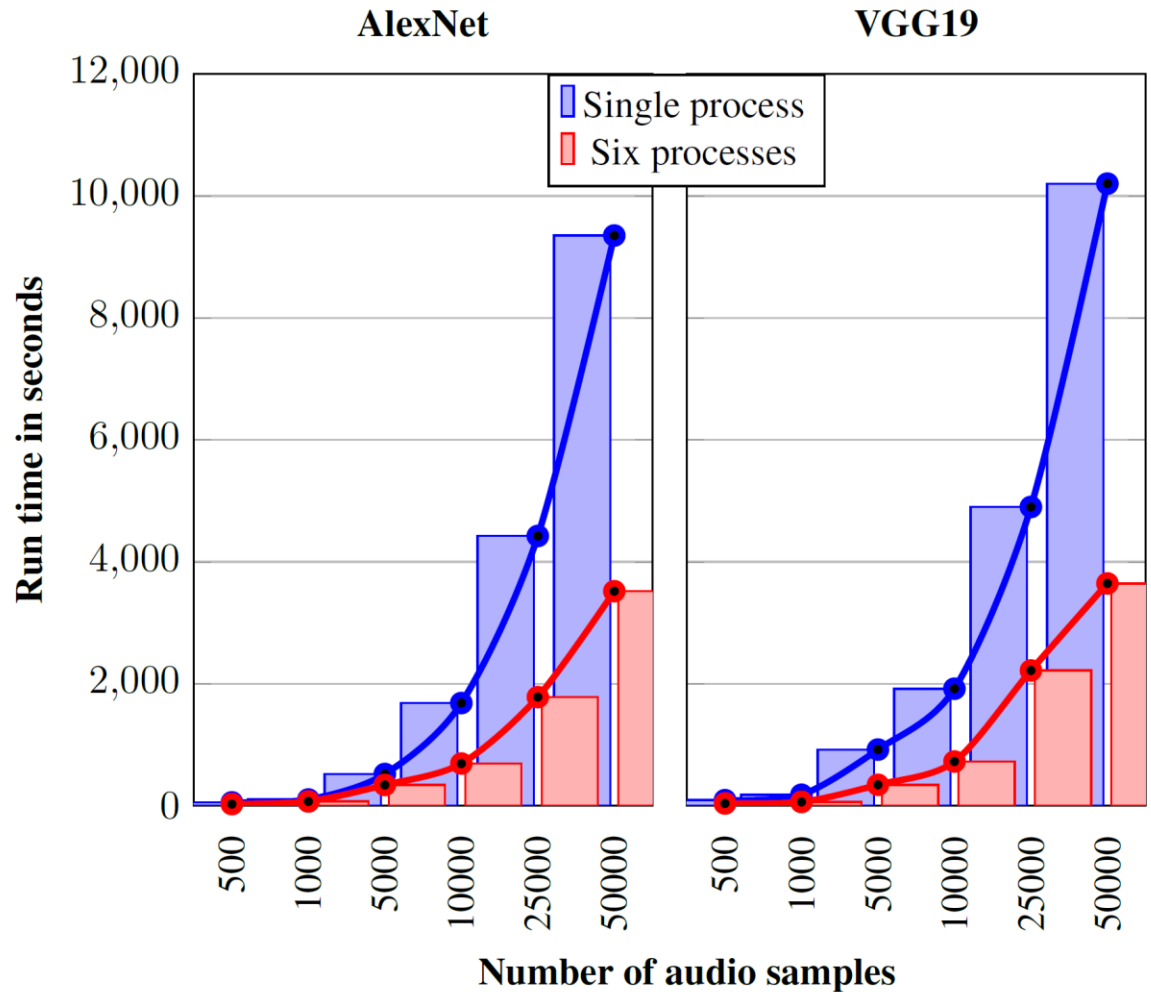
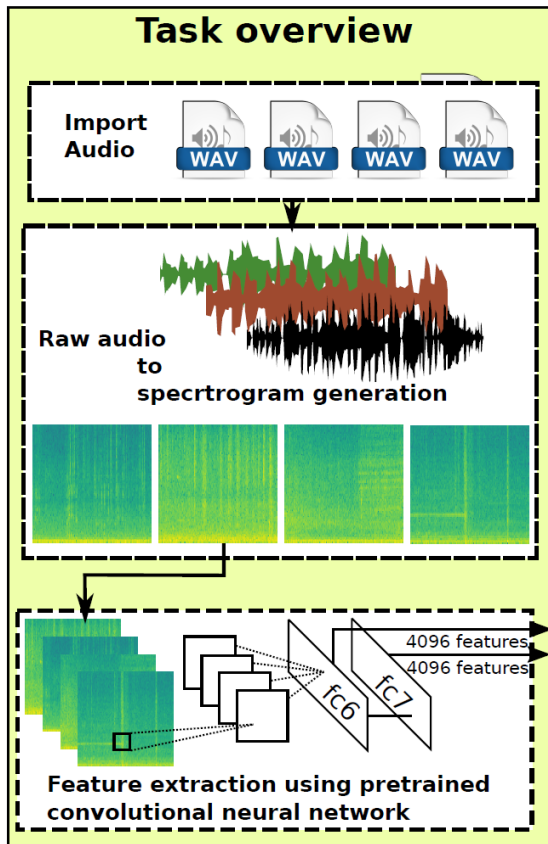
# Fast Processing.

- GPU feature extraction



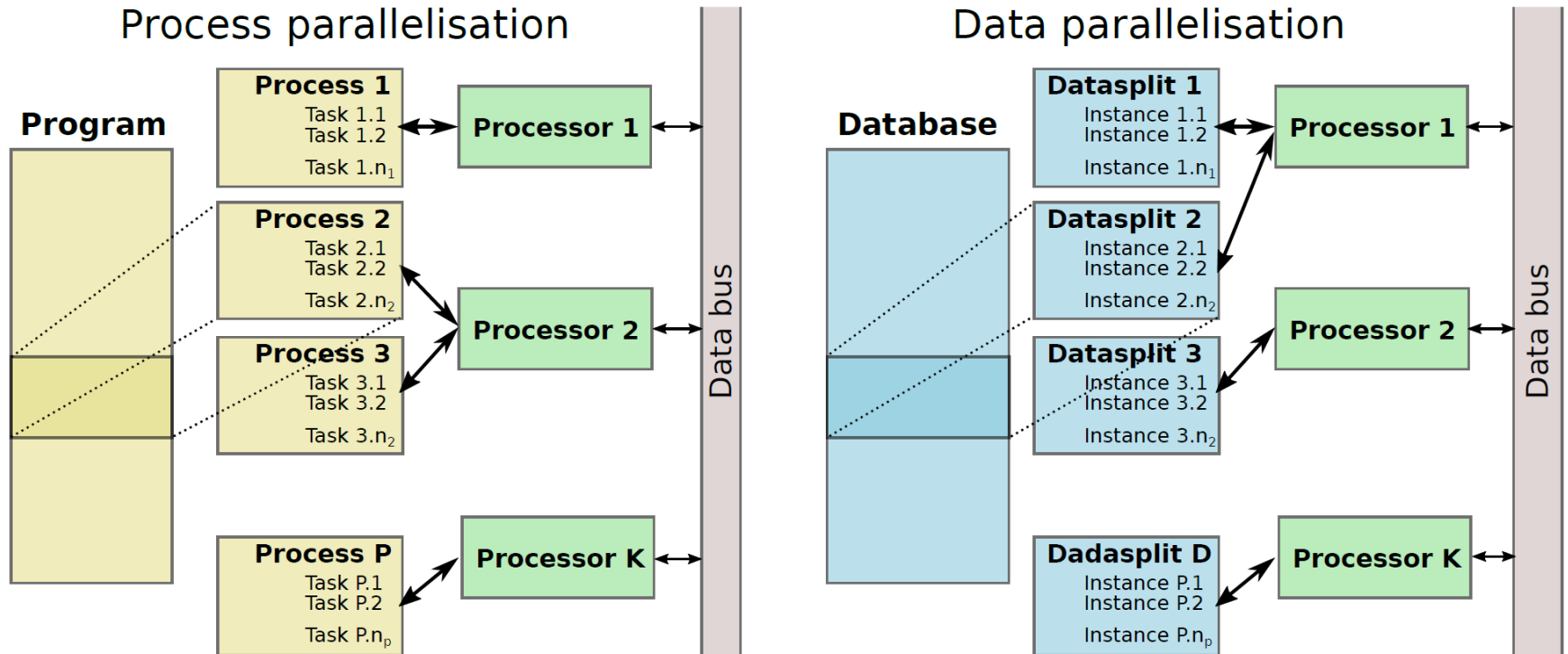
# Fast Processing.

- Parallel...



# Fast Processing.

- Parallelisation**



Application?



# Speaker Verification.

- **Robstness against Paralinguistics?**

- **Example: Alcohol Intoxication**

Negative influence

Improvement by multi-condition training

Larger effect for female speakers

UBM	Tgt	True	Imp.	EER
S	S	S	S	8.1
S	S	A	S	12.9
S	S	A	A	12.3
S	A	S	S	10.9
S	A	A	S	8.1
S	A	A	A	7.9

sober

alcoholised

Spontaneous



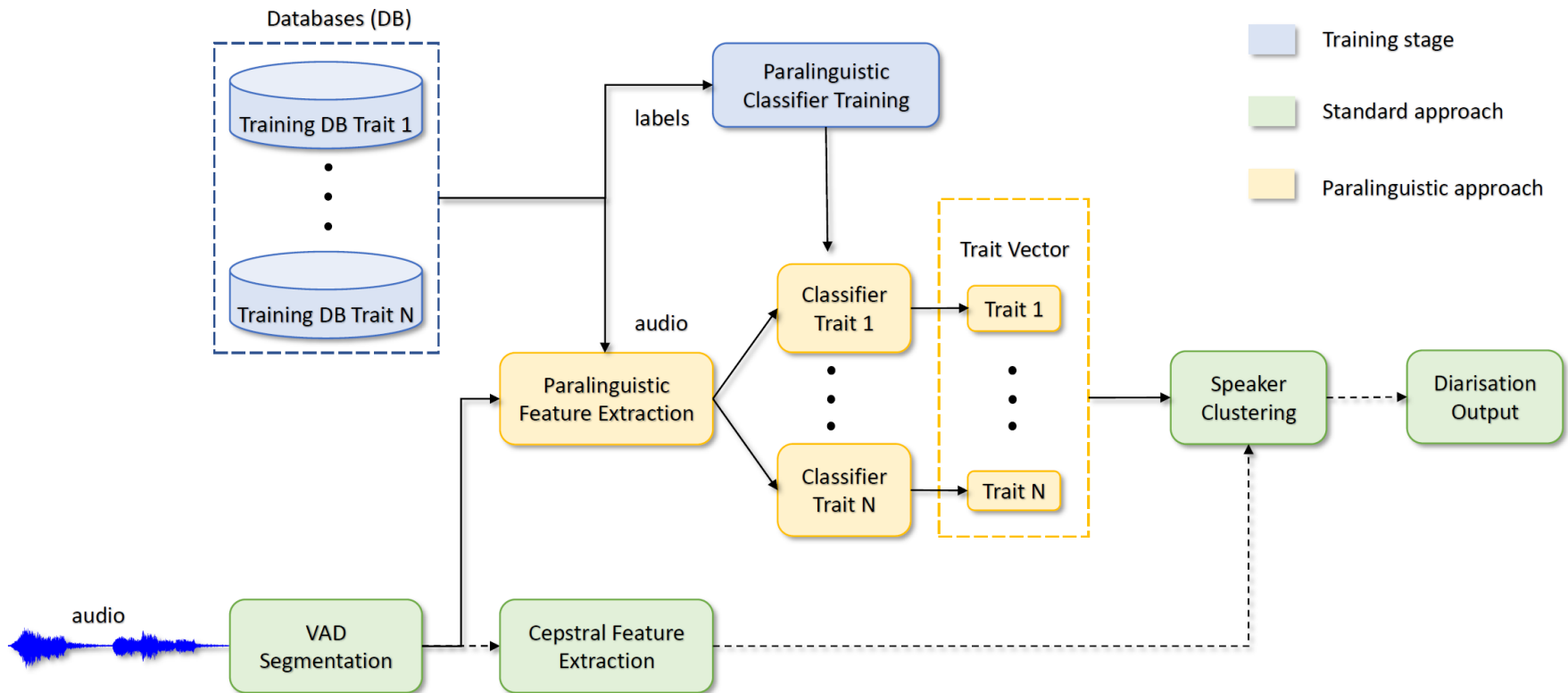
Command & control



# Diarisation.

- Paralings for Diarisation  
SEWA database

System	Miss	sperr
LIUM	6.3	39.0
sensAI	15.2	23.4
Paralings	6.3	38.0



# So?

- **Superhuman in several objective tasks**
- **More (independent) perception studies for subjective tasks!**
- **Increased realism and performance (up to 2x)**
- **Good progress by improved Deep Architectures**
- **Still even many low hanging fruits!**

# Vision.

- **Tighter Coupling w/ Synthesis**
- **Embedding in Dialogues**
- **Reinforcement Learning**
- **NPU optimized Solutions**



Thank You.

# Events.

## »»CROSSROADS OF **S**PEECH AND **L**ANGUAGE««

GENERAL CHAIRS: G. Kubin, Graz | Z. Kačič, Maribor  
TECHNICAL CHAIRS: T. Hain, Sheffield | B. Schuller, Passau/London



**S**  
INTERSPEECH  
2019

GRAZ – AUSTRIA  
SEPTEMBER 15<sup>th</sup> – 19<sup>th</sup> 2019



# Events.



aaac  
emotion-research.net



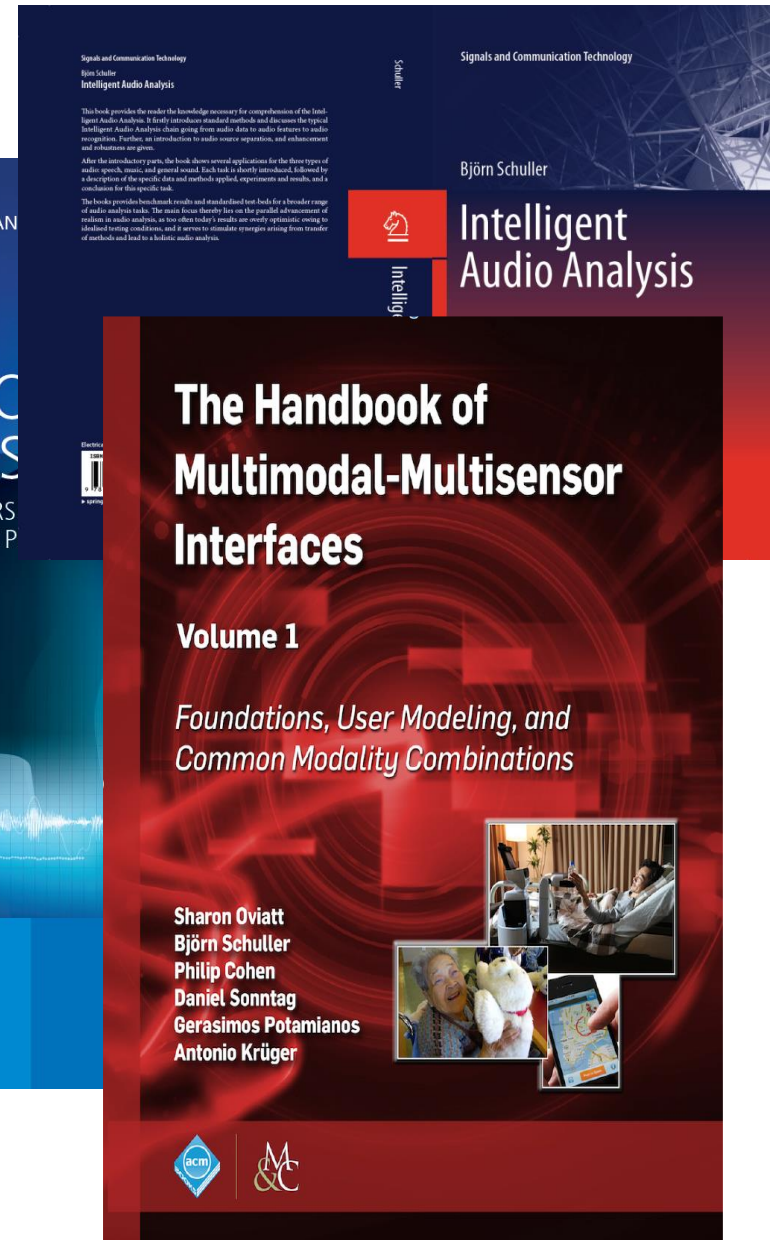
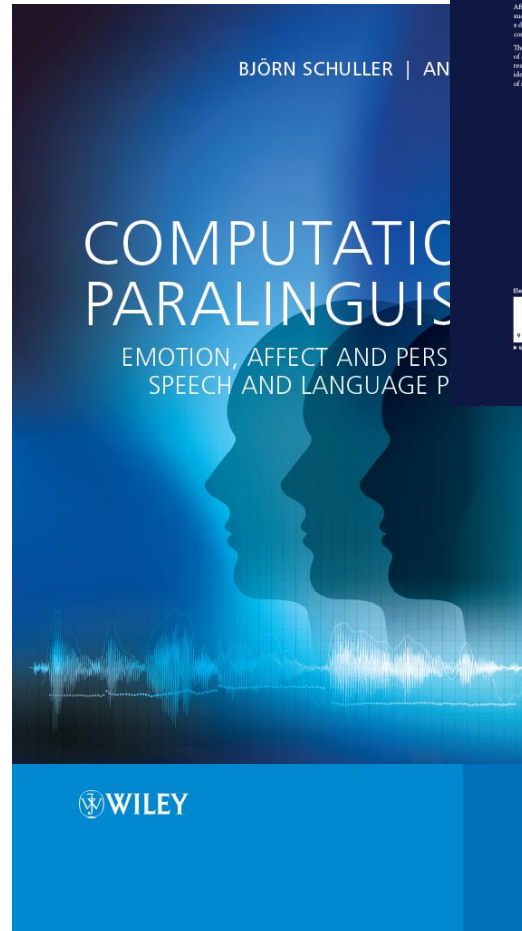
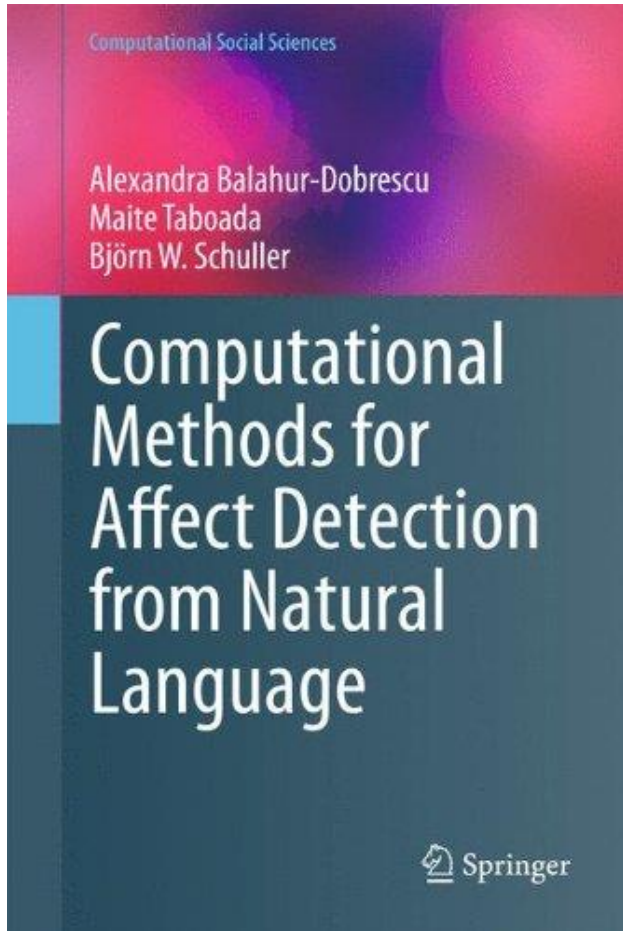
**CALL FOR PAPERS - ACII2019, Cambridge, UK**

**8th International Conference on Affective Computing and Intelligent Interaction**

**3-6 September, 2019**

[www.acii2019.org](http://www.acii2019.org)

# Books.



# Abstract & CV

Human performance is often appearing as a glass ceiling when it comes to automatic speech and speaker analysis. In some tasks, such as health monitoring, however, automatic analysis has successfully started to break this ceiling. The field has benefited from more than a decade of deep neural learning approaches such as recurrent LSTM nets and deep RBMs by now; however, recently, a further major boost could be witnessed. This includes the injection of convolutional layers for end-to-end learning, as well as active and autoencoder-based transfer learning and generative adversarial network topologies to better cope with the ever-present bottleneck of severe data scarcity in the field. At the same time, multi-task learning allowed to broaden up on tasks handled in parallel and include the often met uncertainty in the gold standard due to subjective labels such as emotion or perceived personality of speakers. This talk highlights the named and further latest trends such as increasingly deeper nets and the usage of deep image nets for speech analysis on the road to 'holistic' superhuman speech analysis 'seeing the whole picture' of the person behind a voice. At the same time, increasing efficiency is shown for an ever 'bigger' data and increasingly mobile application world that requires fast and resource-aware processing. The exploitation in ASR and SLU is featured throughout.

Björn W. Schuller heads Imperial College London's/UK Group on Language Audio & Music (GLAM), is a CEO of audEERING, and a Full Professor at University of Augsburg/Germany in CS. He further holds a Visiting Professorship at the Harbin Institute of Technology/China. He received his diploma, doctoral, and habilitation degrees from TUM in Munich/Germany in EE/IT. Previous positions of his include Visiting Professor, Associate, and Scientist at VGTU/Lithuania, University of Geneva/Switzerland, Joanneum Research/Austria, Marche Polytechnic University/Italy, and CNRS-LIMSI/France. His 650+ technical publications (15000+ citations, h-index 59) focus on machine intelligence for audio and signal analysis. He is the Editor in Chief of the IEEE Transactions on Affective Computing, a General Chair of ACII 2019, and a Technical Chair of Interspeech 2019 among various further roles.