

# Speech Processing

Mohamed CHETOUANI  
Professeur des Universités

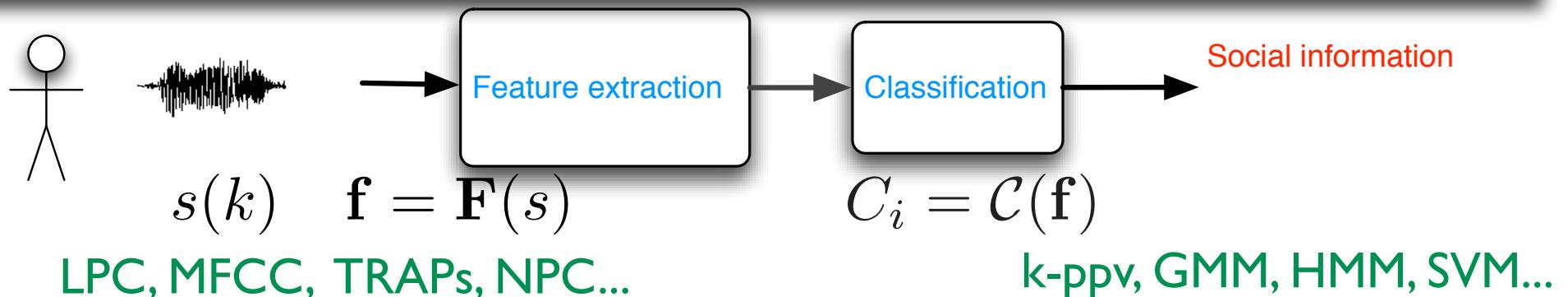
Institut des Systèmes Intelligents et de Robotique (ISIR)  
Sorbonne Université

[mohamed.chetouani@sorbonne-universite.fr](mailto:mohamed.chetouani@sorbonne-universite.fr)

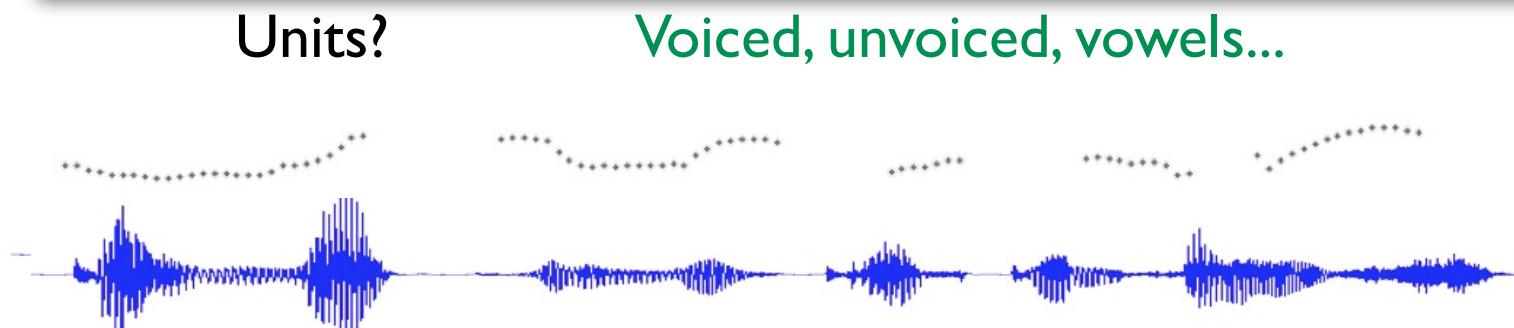


# Speech processing in Social Interactions

Problem statement: From speech samples to social information:  
Identity, Emotion, Intention, Pathology...



Idea: Characterization of speech signals by feature extraction,  
time-scale and classification methods



Units?  
Statistical nature of speech signals?

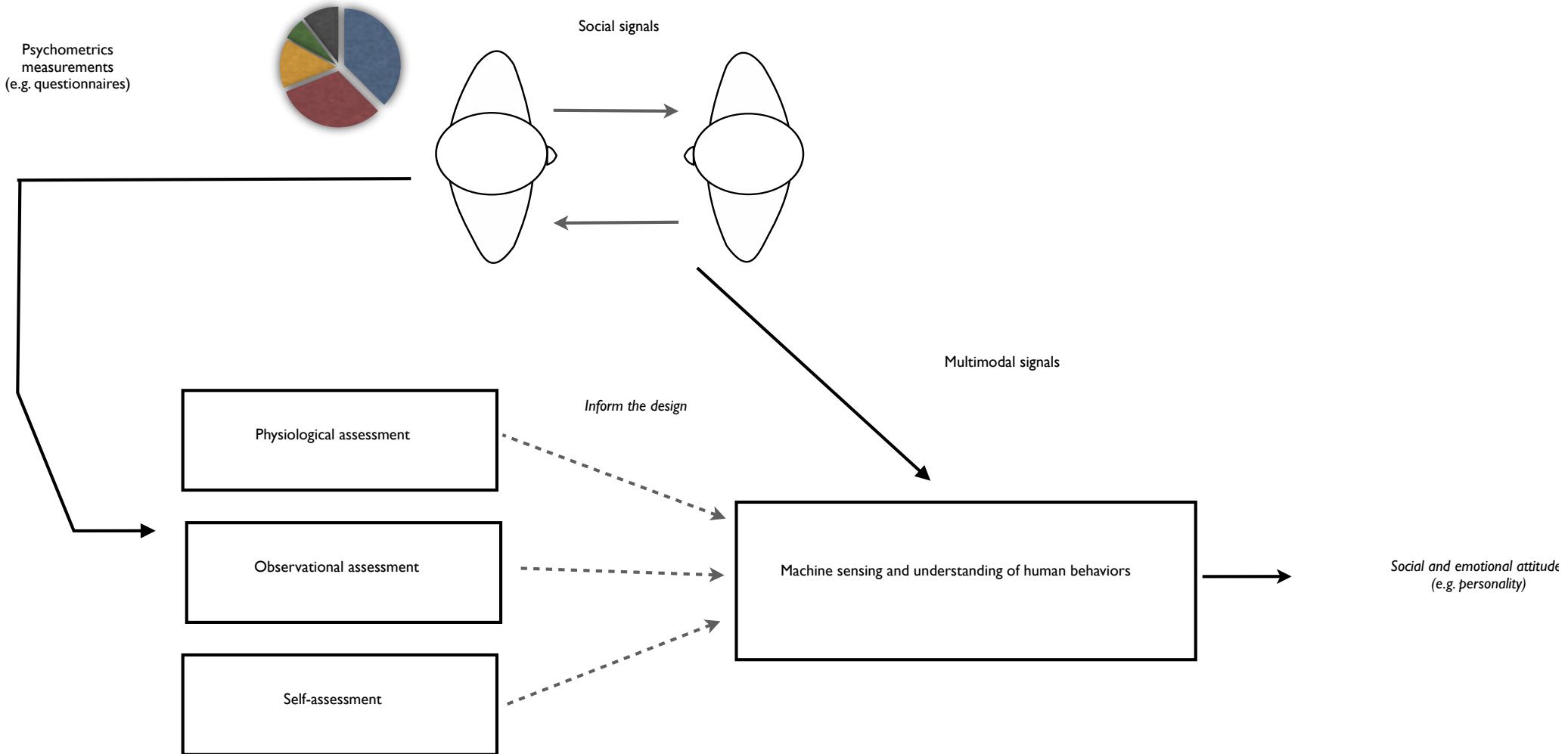
Voiced, unvoiced, vowels...

Gaussian vs Non-gaussian,  
Stationarity...

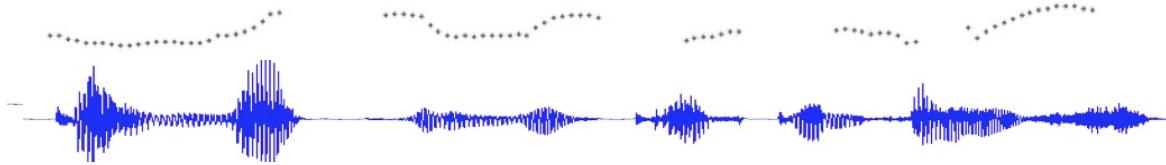
Subjectivity?

Manual vs  
Automatic  
annotations

# Affective Computing



# Speech Affective Computing Methodology



## (1) Extraction of Low-Level Descriptors (LLD):

- Pitch (fundamental frequency)
- Energy
- Spectrum
- ...

## (2) Application of functionals

- Statistics
- Duration
- Linear prediction
- ...

## (3) Classification

- Speech turn
- Utterance turn

# Speech Affective Computing

## Lessons learnt from the First Computational Paralinguistics Challenge

Table 6

ComParE acoustic feature set: 65 provided **low-level descriptors** (LLD).

<b>4 Energy Related LLD</b>	<b>Group</b>
Sum of Auditory Spectrum (Loudness)	Prosodic
Sum of RASTA-Style Filtered Auditory Spectrum	Prosodic
RMS Energy, Zero-Crossing Rate	Prosodic
<b>55 Spectral LLD</b>	<b>Group</b>
RASTA-Style Auditory Spectrum, Bands 1–26 (0–8 kHz)	Spectral
MFCC 1–14	Cepstral
Spectral Energy 250–650 Hz, 1 k–4 kHz	Spectral
Spectral Roll Off Point 0.25, 0.50, 0.75, 0.90	Spectral
Spectral Flux, Centroid, Entropy, Slope, Harmonicity	Spectral
Spectral Psychoacoustic Sharpness	Spectral
Spectral Variance, Skewness, Kurtosis	Spectral
<b>6 Voicing Related LLD</b>	<b>Group</b>
$F_0$ (SHS & Viterbi Smoothing)	Prosodic
Probability of Voicing	Sound Quality
Log. HNR, Jitter (Local, Delta), Shimmer (Local)	Sound Quality

# Speech Affective Computing

## Lessons learnt from the First Computational Paralinguistics Challenge

Table 7

ComParE acoustic feature set: **functionals** applied to LLDs as defined in [Table 6](#).

<b>Mean Values</b>
Arithmetic Mean $A^{\Delta}, B$ , Arithmetic Mean of Positive Values $A^{\Delta}, B$
Root-Quadratic Mean, Flatness
<b>Moments:</b> Standard Deviation, Skewness, Kurtosis
<b>Temporal Centroid</b> $A^{\Delta}, B$
<b>Percentiles</b>
Quartiles 1–3, Inter-Quartile Ranges 1–2, 2–3, 1–3
1%-tile, 99%-tile, Range 1–99%
<b>Extrema</b>
Relative Position of Maximum and Minimum, Full Range (Maximum–Minimum)
<b>Peaks and Valleys</b> <sup>A</sup>
Mean of Peak Amplitudes,
Difference of Mean of Peak Amplitudes to Arithmetic Mean
Mean of Peak Amplitudes Relative to Arithmetic Mean
Peak to Peak Distances: Mean and Standard Deviation
Peak Range Relative to Arithmetic Mean
Range of Peak Amplitude Values
Range of Valley Amplitude Values Relative to Arithmetic Mean
Valley-Peak (Rising) Slopes: Mean and Standard Deviation
Peak-Valley (Falling) Slopes: Mean and Standard Deviation
<b>Up-Level Times:</b> 25%, 50%, 75%, 90%
<b>Rise and Curvature Time</b>
Relative Time in which Signal is Rising
Relative Time in which Signal has Left Curvative
<b>Segment Lengths</b> <sup>A</sup>
Mean, Standard Deviation, Minimum, Maximum
<b>Regression</b> $A^{\Delta}, B$
Linear Regression: Slope, Offset, Quadratic Error
Quadratic Regression: Coefficients $a$ and $b$ , Offset $c$ , Quadratic Error
<b>Linear Prediction</b>
LP Analysis Gain (Amplitude Error), LP Coefficients 1–5

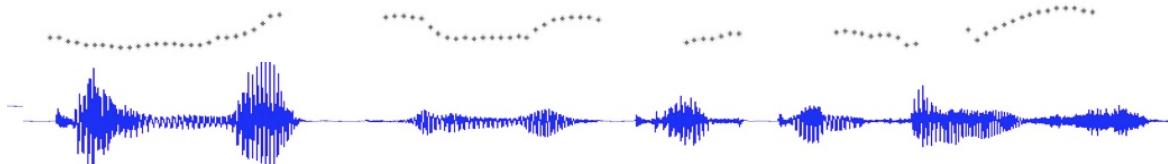
<sup>A</sup>Functionals applied only to energy related and spectral LLDs (group A)

<sup>B</sup>Functionals applied only to voicing related LLDs (group B)

$\Delta$ Functionals applied only to  $\Delta$ LLDs

$\Delta\Delta$  Functionals **not** applied to  $\Delta$ LLDs

# Speech Affective Computing Methodology



## **INTERSPEECH 2013 configuration:**

6373 features

Group A: 4 energy related LLDs + 55 spectral LLDs

Group B: 6 voicing related LLDs

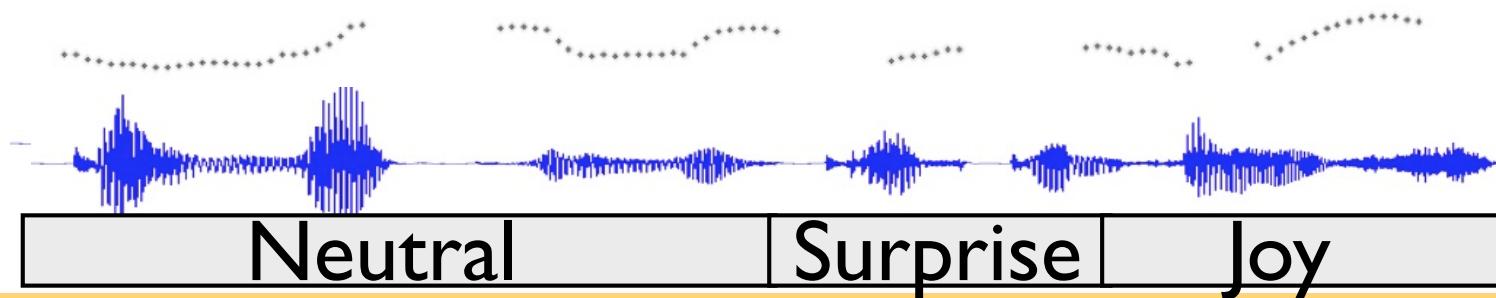
54 functionals applied to Group A + 46 applied to Delta LLDs -> 5900 features

39 functionals applied to Group B and Delta LLDs -> 468 features

5 temporal static descriptors for voiced segments: mean length, standard deviation of the segment length, minimum and maximum of voiced segments, ratio of non-zero F0

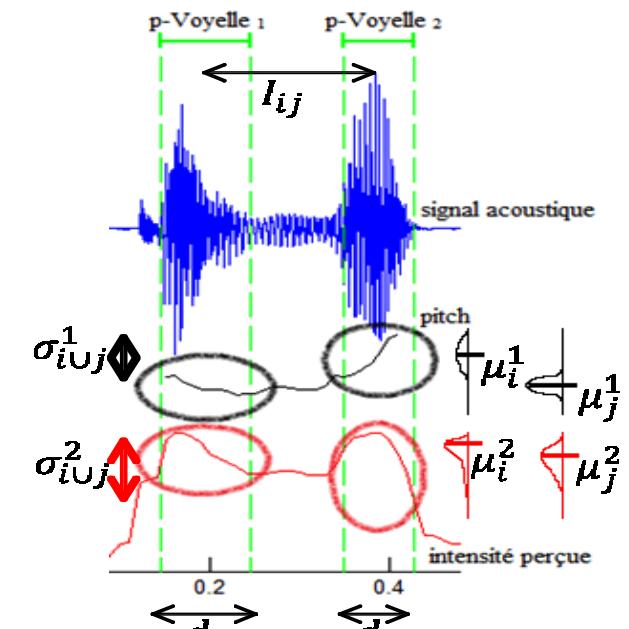
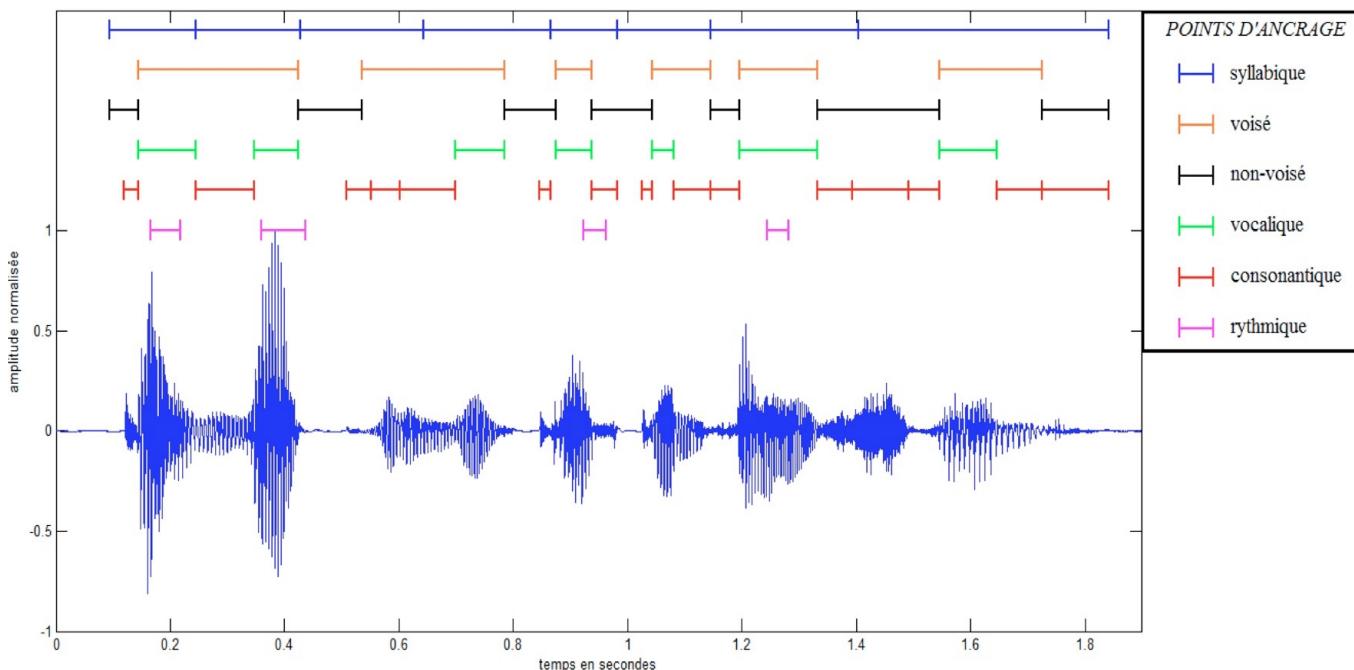
**Total: 5900 + 468 + 5 = 6373 features**

# Speech Affective Computing



Units for emotional speech processing:

- Acoustical, Prosodic and Rhythmic prominence
- Dynamics of units: Huang-Hilbert transform, Hotelling distance...

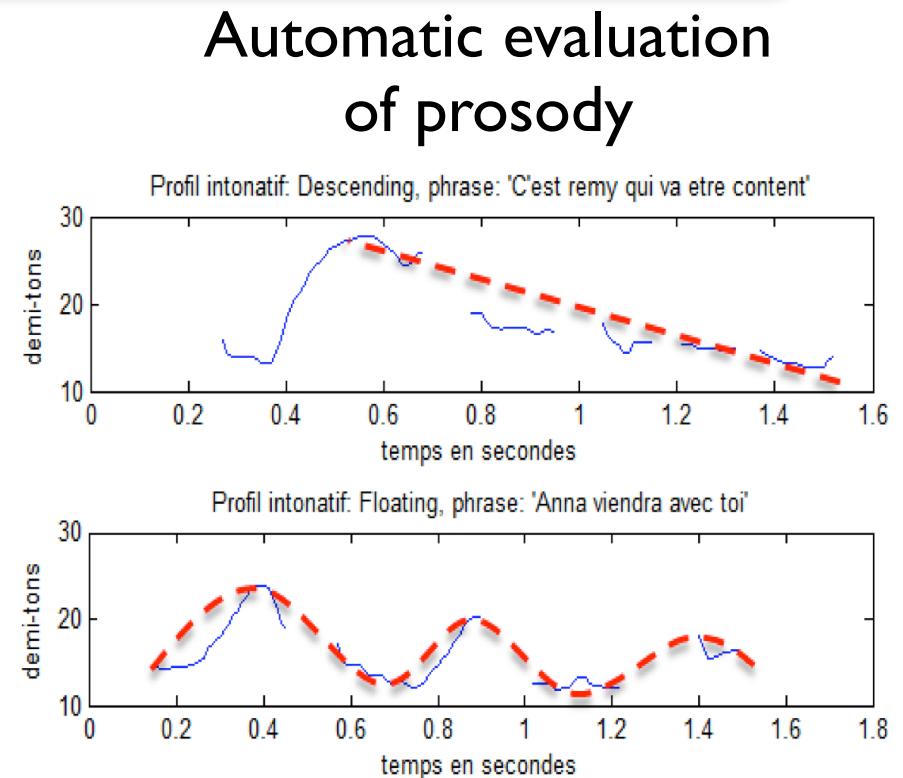
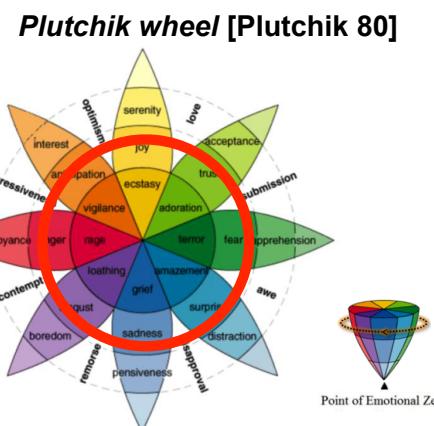
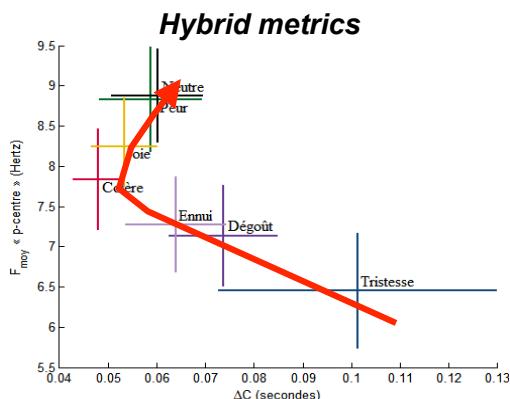
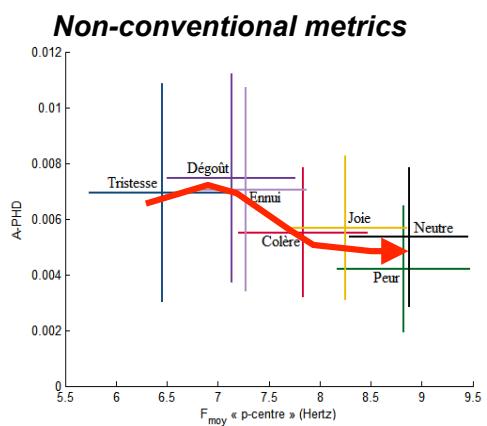
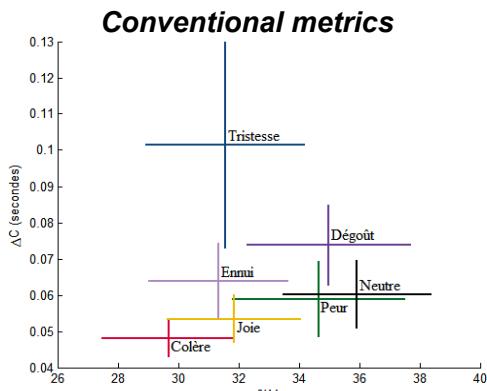


$$HD_{ij} = \frac{d_i d_j}{d_i + d_j} \left[ (\mu_i - \mu_j)^T \Sigma_{i \cup j}^{-1} (\mu_i - \mu_j) \right]$$

# Speech Affective Computing

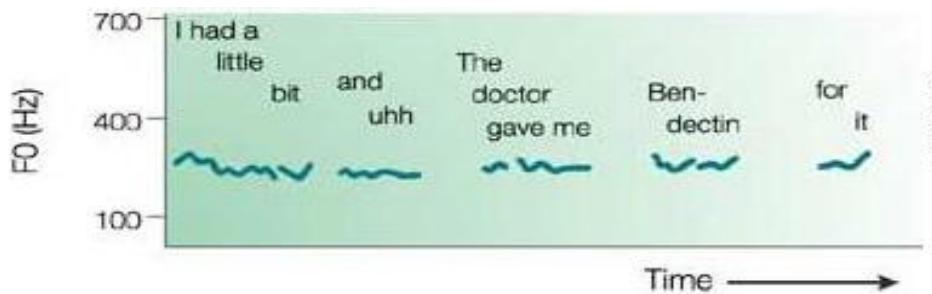
## Applications

- ▶ Acted emotion recognition
- ▶ Children assessment: Autism, Pervasive Developmental Disorders, Language impairment
- ▶ Spontaneous and atypical emotions

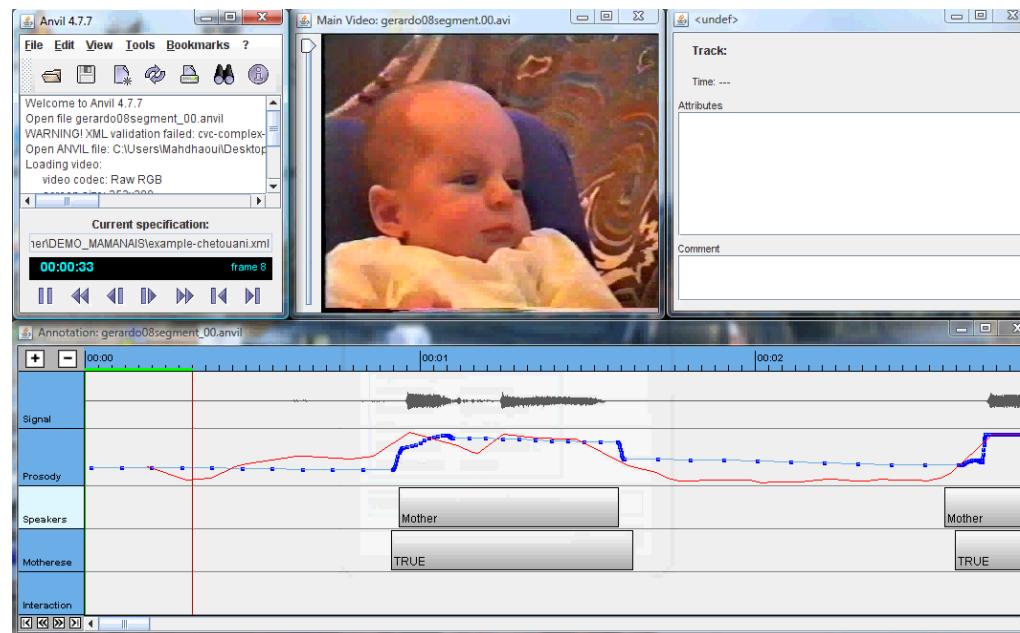
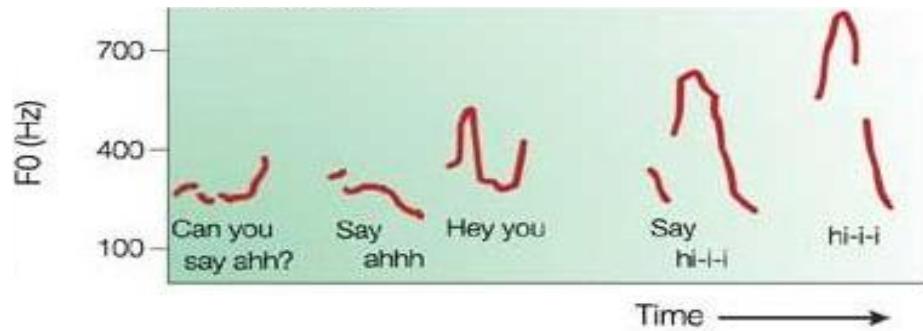


# Speech Affective Computing Methodology

Adult-directed speech



Infant-directed speech



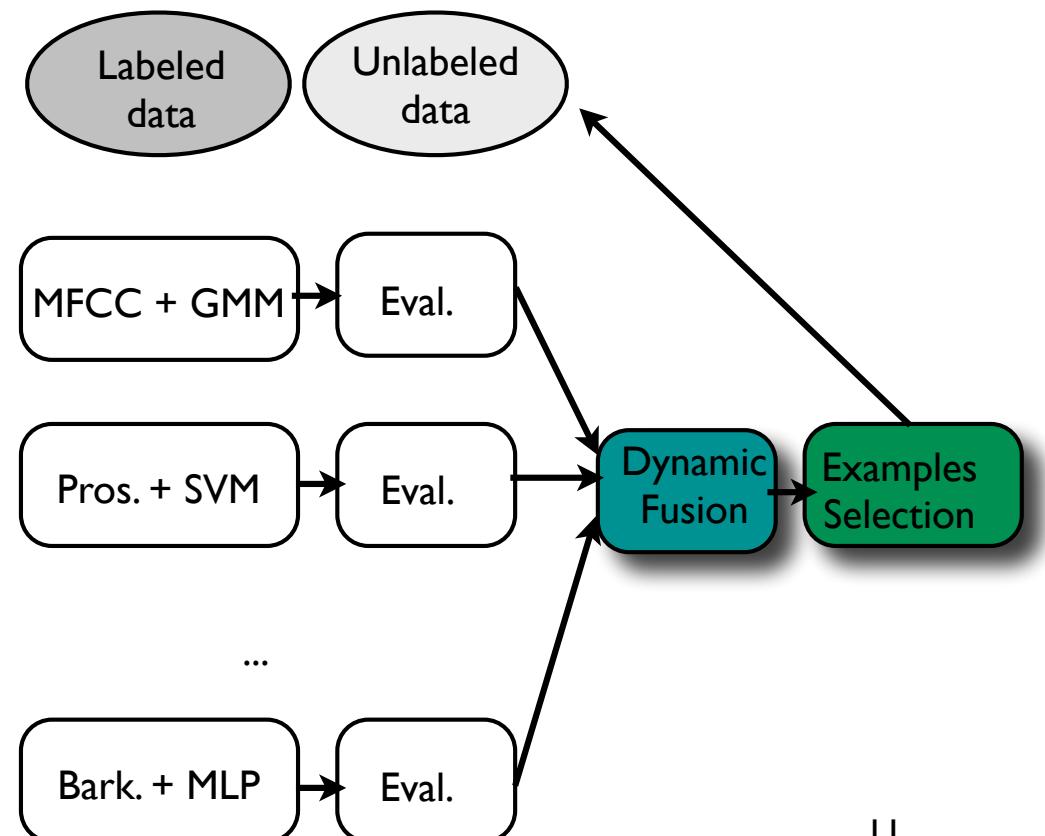
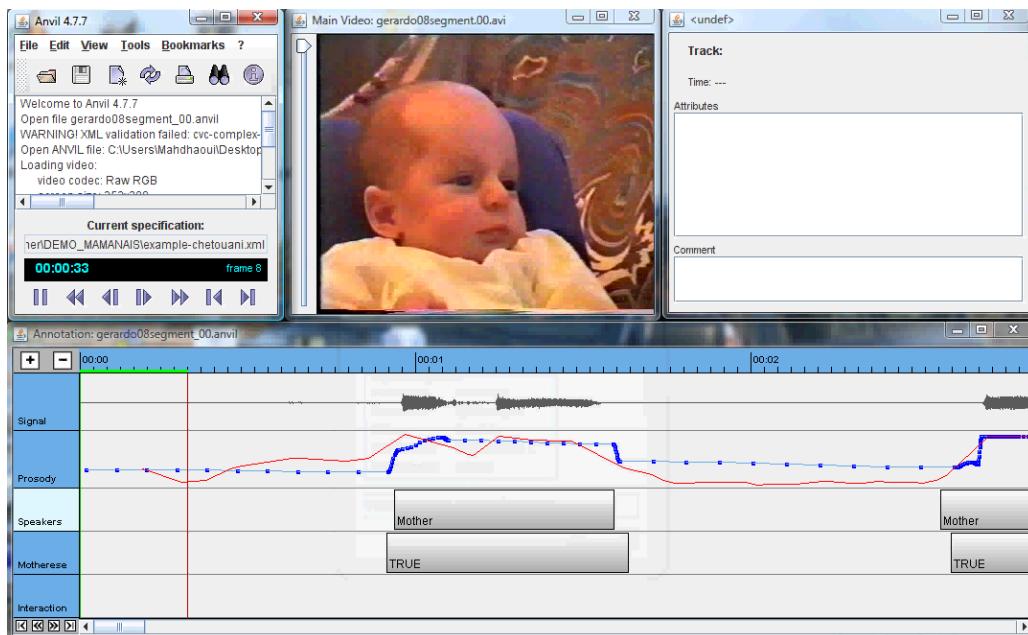
C. Saint-georges et al. : Motherese in Interaction: At the Cross-Road of Emotion and Cognition? (A Systematic Review) *Plos One*, 2013

# Speech Affective Computing

## Co-Learning approaches

### Annotation and Subjectivity of social signals

- ▶ Semi-supervised learning: combining labeled and unlabeled data
- ▶ Co-training and fusion: Multi-view characterization



# Speech Affective Computing

## Co-Learning approaches

Table 3: Self-training algorithm

<p><b>Given:</b></p> <ul style="list-style-type: none"><li>a set <math>L</math> of Labelled examples</li><li>a set <math>U</math> of Unlabelled examples</li><li>a number <math>n</math> of examples to be added to <math>L</math> in each iteration</li></ul> <p><b>Loop:</b></p> <ul style="list-style-type: none"><li>Use <math>L</math> to train the classifier <math>h</math></li><li>Allow <math>h</math> to label <math>U</math></li><li>Let <math>T</math> be the <math>n</math> examples in <math>U</math> on which <math>h</math> makes the most confident predictions</li><li>Add <math>T</math> to <math>L</math></li><li>Remove <math>T</math> from <math>U</math></li></ul> <p><b>End</b></p>
---

Table 4: Co-Training algorithm

<p><b>Given:</b></p> <ul style="list-style-type: none"><li>a set <math>L</math> of Labelled examples</li><li>a set <math>U</math> of Unlabelled examples</li></ul> <p><b>Loop:</b></p> <ul style="list-style-type: none"><li>Use <math>L</math> to train each classifier <math>h_1</math></li><li>Use <math>L</math> to train each classifier <math>h_2</math></li><li>Allow <math>h_1</math> to label <math>p_1</math> positive and <math>n_1</math> negative examples from <math>U</math></li><li>Allow <math>h_2</math> to label <math>p_2</math> positive and <math>n_2</math> negative examples from <math>U</math></li><li>Add these self-labelled examples to <math>L</math></li><li>Remove these self-labelled examples from <math>U</math></li></ul> <p><b>End</b></p>
---

# Speech Affective Computing

## Co-Learning approaches

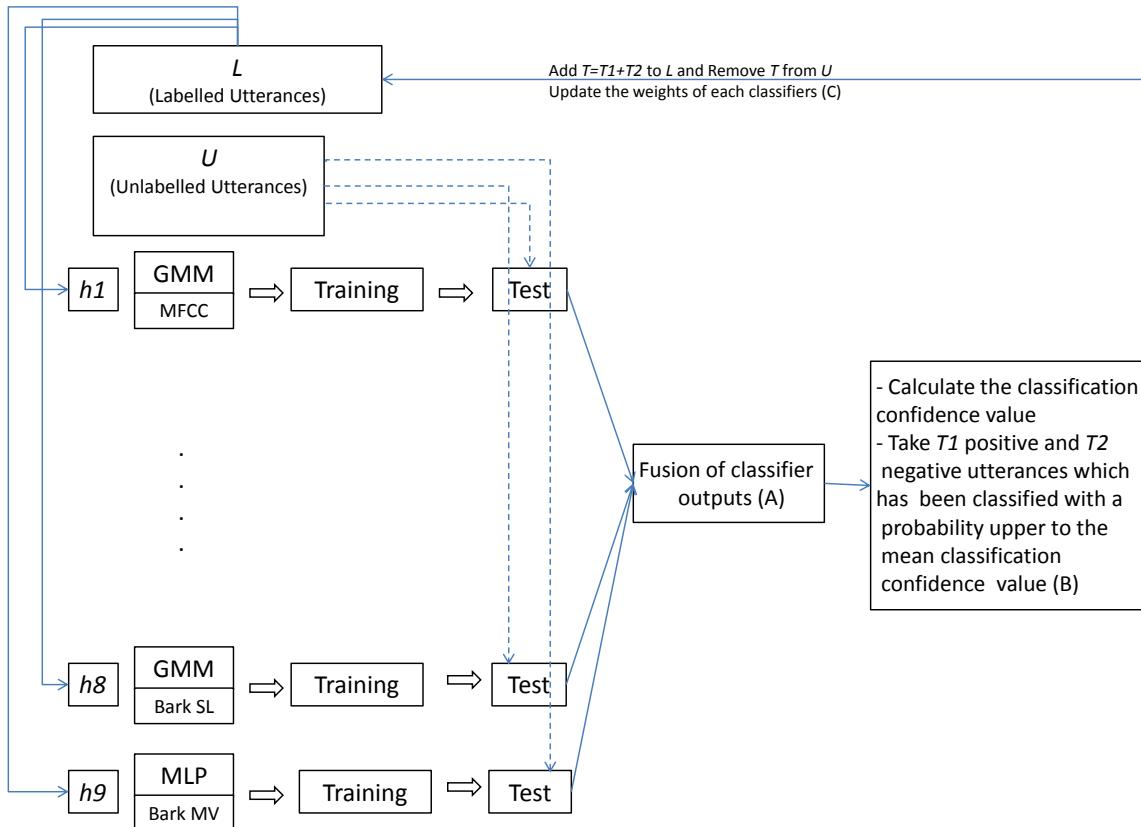


Figure 5: Structure of the proposed Co-training algorithm

A. Mahdhaoui and M. Chetouani : Supervised and semi-supervised infant-directed speech classification for parent-infant interaction analysis. *Speech Communication* 53(9) 1149-1161 (2011).

# Speech Affective Computing

## Co-Learning approaches

Table 5: The proposed Co-Training algorithm

**Given:**

a set  $L$  of  $m$  Labelled examples  $\{(l_1^1, \dots, l_v^1, y_1), \dots, (l_m^1, \dots, l_m^v, y_m)\}$  with labels  $y_i = \{1, 2\}$   
 a set  $U$  of  $n$  Unlabelled examples  $\{(x_1^1, \dots, x_1^v), \dots, (x_n^1, \dots, x_n^v)\}$   
 $v$  = number of view (classifier)

**Initialization:**

$\omega_k$  (weights of classifier) =  $1/v$  for all the view

**While U not empty**

**A. Classify all the example of the test database:**

Do for  $k = 1, 2, \dots, v$

1. Use  $L$  to train each classifier  $h_k$
2. Classify all examples of  $U$  by each  $h_k$
3. Calculate the probability of classification for each example  $x_i$  from  $U$ ,  
 $p(C_j|x_i) = \sum_{k=1}^v \omega_k \times h_k(C_j|x_i^k)$
4.  $Labels(x_i) = argmax(p(C_j|x_i))$

End for

**B. Update the training ( $L$ ) and test ( $U$ ) databases:**

$U_j = \{z_1, \dots, z_{n_j}\}$  the ensemble of example classified  $C_j$

Do for  $i = 1, 2, \dots, n_j$

$$p(C_j|z_i) = \frac{\sum_{k=1}^v \omega_k \times h_k(C_j|z_i^k)}{\sum_{k=1}^v \omega_k}$$

End for

$$margin_j = \frac{\sum_1^{n_j} p(C_j|z_i)}{n_j}$$

Take  $T_j$  from  $U_j$  the examples which has classified on  $C_j$  with a probability upper to  $margin_j$ .

$$T = \sum T_j$$

Add  $T$  to  $L$  and remove it from  $U$

**C. Update weights:**  $\omega_k = \frac{\sum_{i=1}^{size(T)} h_k(z_i^k)}{\sum_{k=1}^v \sum_{i=1}^{size(T)} h_k(z_i^k)}$

**End While**

A. Mahdhaoui and M. Chetouani : Supervised and semi-supervised infant-directed speech classification for parent-infant interaction analysis. *Speech Communication* 53(9) 1149-1161 (2011).

# Speech Affective Computing

## Co-Learning approaches

Table 8: Classification accuracy with different numbers of annotations

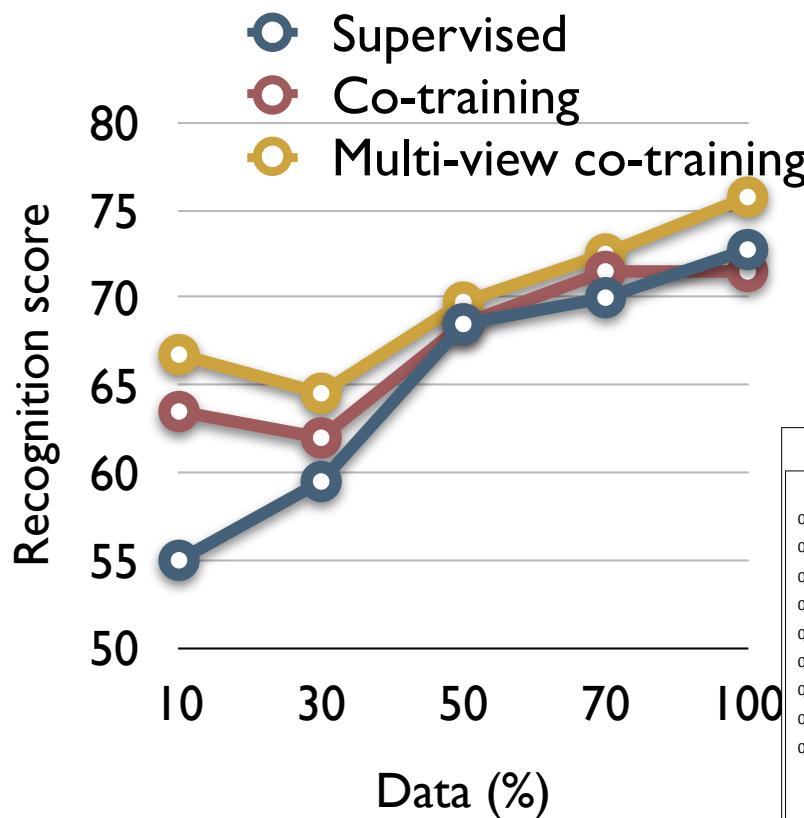
Number of annotations	10	20	30	40	50	60	70	80	90	100
Proposed Co-training method	66.75	65.25	63.5	67	69.75	72.25	72.5	71.75	74	75.75
<i>Co-training</i> standard (using h1 and h4)	63.5	62.5	62	64.5	68.5	69.75	71.25	69.5	71	71.5
<i>Co-training</i> standard (using all the classifiers h1-h9)	57	58.5	58.5	61	64	67	67.25	68	69	68.5
<i>Self-training</i> (using h1: MFCC-GMM)	52	50	50	54	55	62.5	61	65	69	70.25
<i>Self-training</i> (using h2: prosody-GMM)	54	52.5	53	52	53.5	58	59	62	64.5	67.75
Supervised method: MFCC-GMM (best configuration)	55	59.25	59.5	61.5	68.5	71	70	69.75	71.5	72.75

A. Mahdhaoui and M. Chetouani : Supervised and semi-supervised infant-directed speech classification for parent-infant interaction analysis. *Speech Communication* 53(9) 1149-1161 (2011).

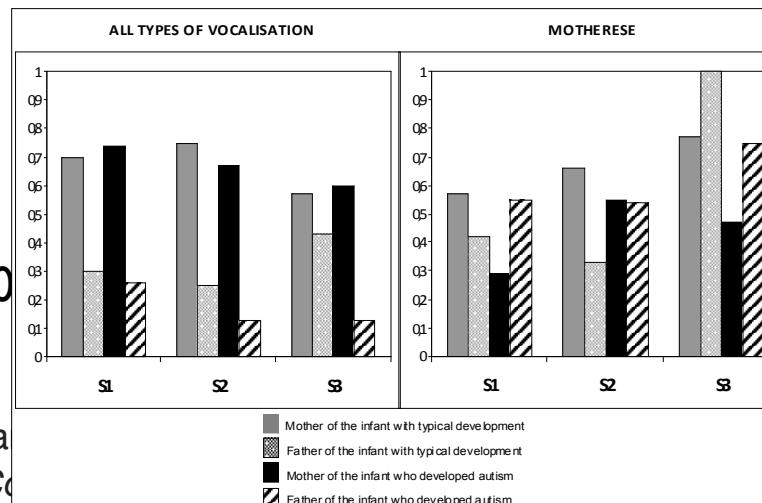
# Speech Affective Computing

## Co-Learning approaches

- ▶ Motherese detection in Family Home Movies
- ▶ Early signs of Autism Spectrum Disorders



Motherese  
detection



A. Mahdhaoui and M. Chetouani : Supervised approach for parent-infant interaction analysis. *Speech Communication*

R. Cassel et al.: Course of maternal prosodic incitation (motherese) during early development in autism: an exploratory study. *Interaction Studies*. Vol 14 Pages 480-496 2014.

# Lessons learnt from the First Computational Paralinguistics Challenge

B. Schuller et al. / Computer Speech & Language 53 (2019) 156–180

159

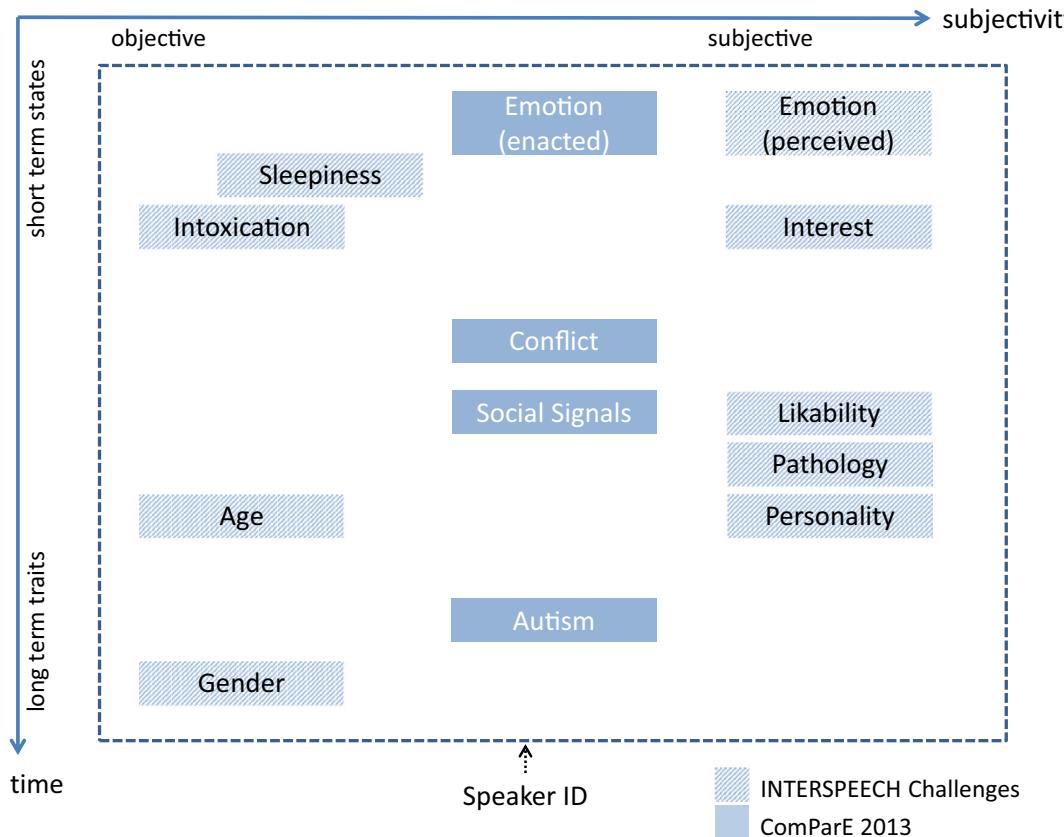


Fig. 1. Speaker characteristics investigated in the INTERSPEECH Challenges 2009–2012 and the First Computational Paralinguistics Challenge (ComParE) 2013.

Björn Schuller, Felix Weninger, Yue Zhang, Fabien Ringeval, Anton Batliner, Stefan Steidl, Florian Eyben, Erik Marchi, Alessandro Vinciarelli, Klaus Scherer, **Mohamed Chetouani**, Marcello Mortillaro, Affective and behavioural computing: Lessons learnt from the First Computational Paralinguistics Challenge, Computer Speech & Language (2019)

# Lessons learnt from the First Computational Paralinguistics Challenge

Table 4

Partitioning of the GEMEP database into train, dev(velopment), and test set for 12-way classification by emotion category, and binary classification by pos(itive)/neg(ative) arousal (A) and valence (V).

#	Train	Dev	Test	A	V	$\Sigma$
Admiration <sup>+</sup>	20	2	8	pos	pos	30
Amusement	40	20	30	pos	pos	90
Anxiety	40	20	30	neg	neg	90
Cold anger	42	12	36	neg	neg	90
Contempt <sup>+</sup>	20	6	4	neg	neg	30
Despair	40	20	30	pos	neg	90
Disgust <sup>+</sup>	20	2	8	—*	—*	30
Elation	40	12	38	pos	pos	90
Hot anger	40	20	30	pos	neg	90
Interest	40	20	30	neg	pos	90
Panic fear	40	12	38	pos	neg	90
Pleasure	40	20	30	neg	pos	90
Pride	40	12	38	pos	pos	90
Relief	40	12	38	neg	pos	90
Sadness	40	12	38	neg	neg	90
Shame <sup>+</sup>	20	2	8	pos	neg	30
Surprise <sup>+</sup>	20	6	4	—*	—*	30
Tenderness <sup>+</sup>	20	6	4	neg	pos	30
$\Sigma$	602	216	442			1260

<sup>+</sup> Mapped to ‘other’ and excluded from evaluation in 12-class task.

\* Mapped to ‘undefined’ and excluded from evaluation in binary tasks.

# Deep Learning for Human Affect Recognition: Insights and New Developments

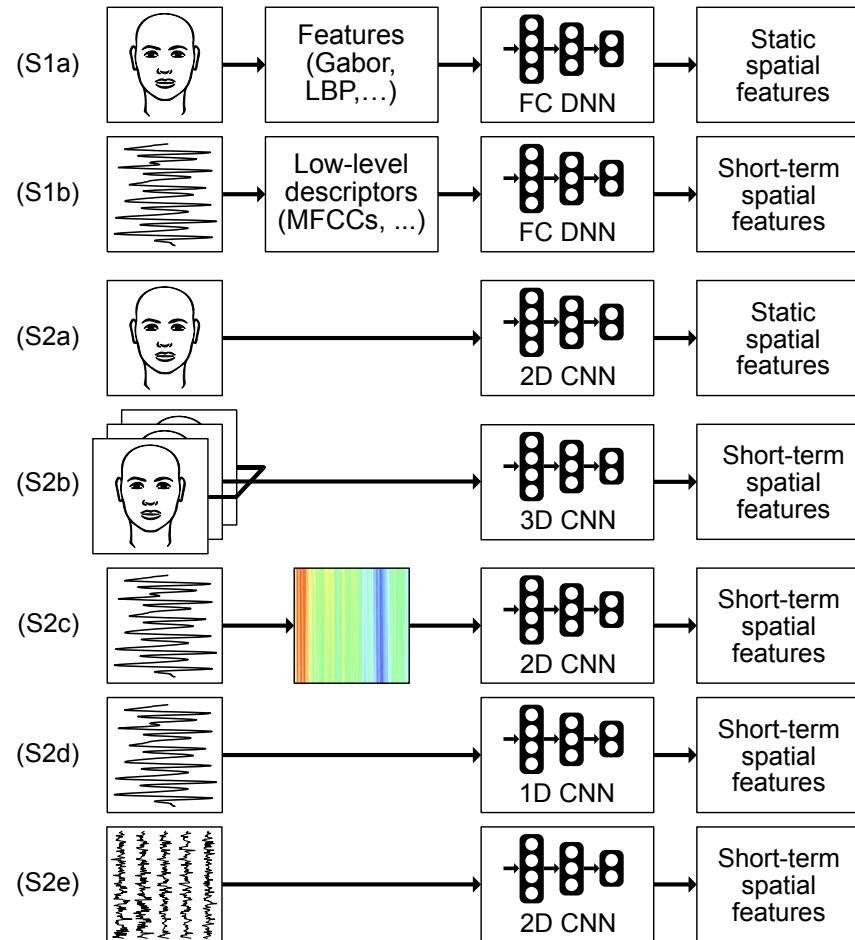


Fig. 2. Applications of deep learning for spatial feature learning with fully-connected DNNs (S1a–S1b) and CNNs (S2a–S2e).

# Deep Learning for Human Affect Recognition: Insights and New Developments

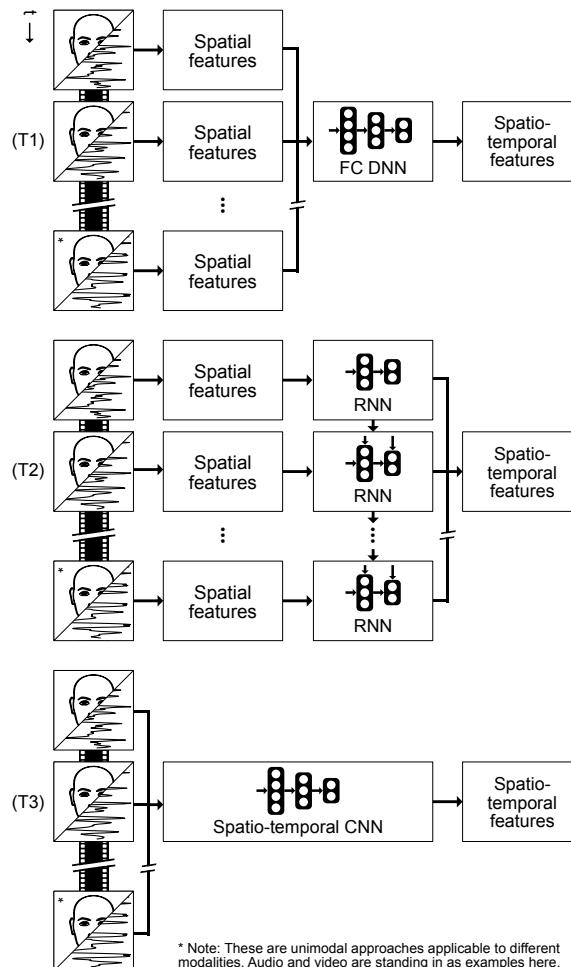


Fig. 5. Applications of deep learning for temporal feature learning with fully-connected DNNs (T1), RNNs (T2), and CNNs (T3).

# Deep Learning for Human Affect Recognition: Insights and New Developments

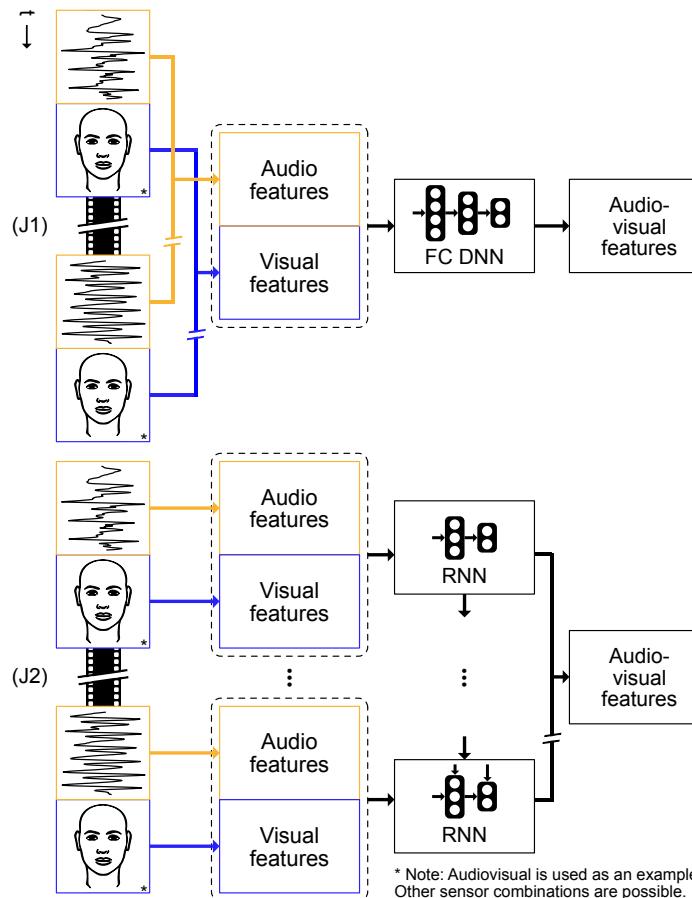
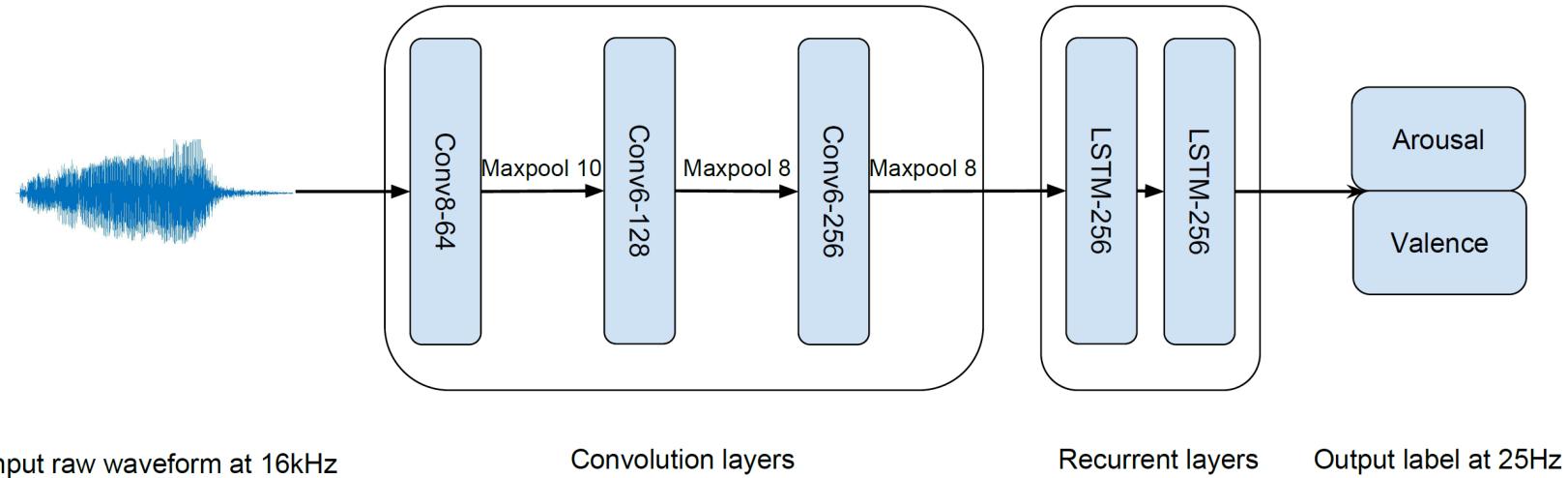


Fig. 6. Applications of deep learning for joint multimodal feature learning with fully-connected fusion DNNs (J1) and fusion RNNs (J2).

# End-to-end approach

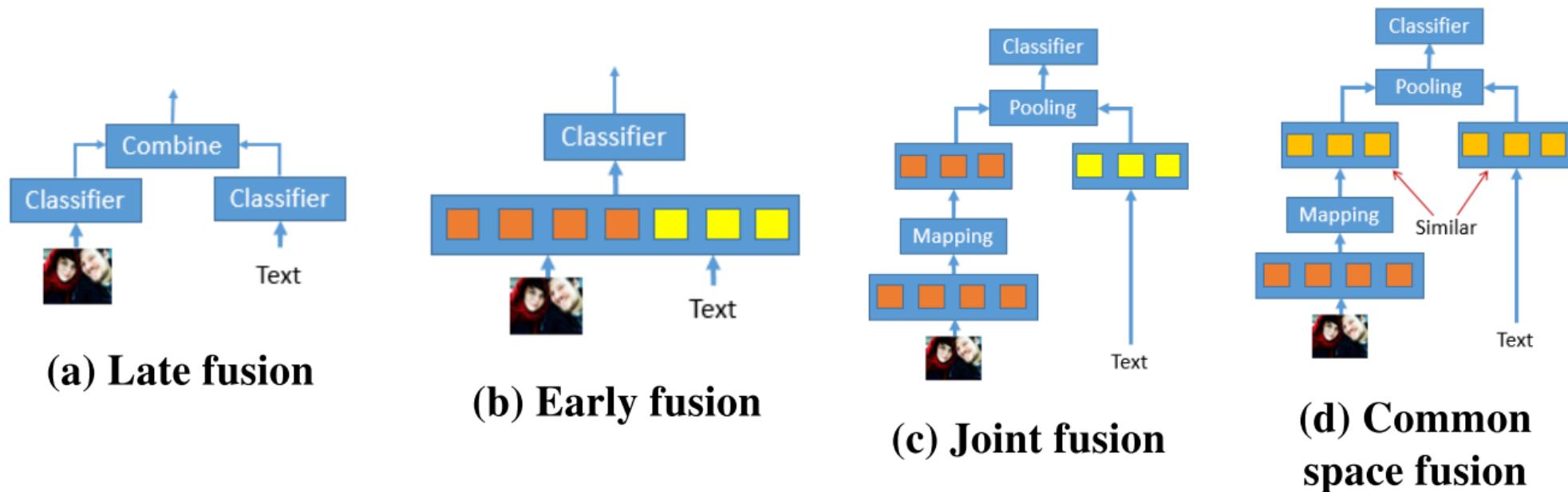


## Convolution operation:

$$(f \star h)(i, j) = \sum_{k=-T}^T \sum_{m=-T}^T f(k, m) \cdot h(i - k, j - m) \quad (1)$$

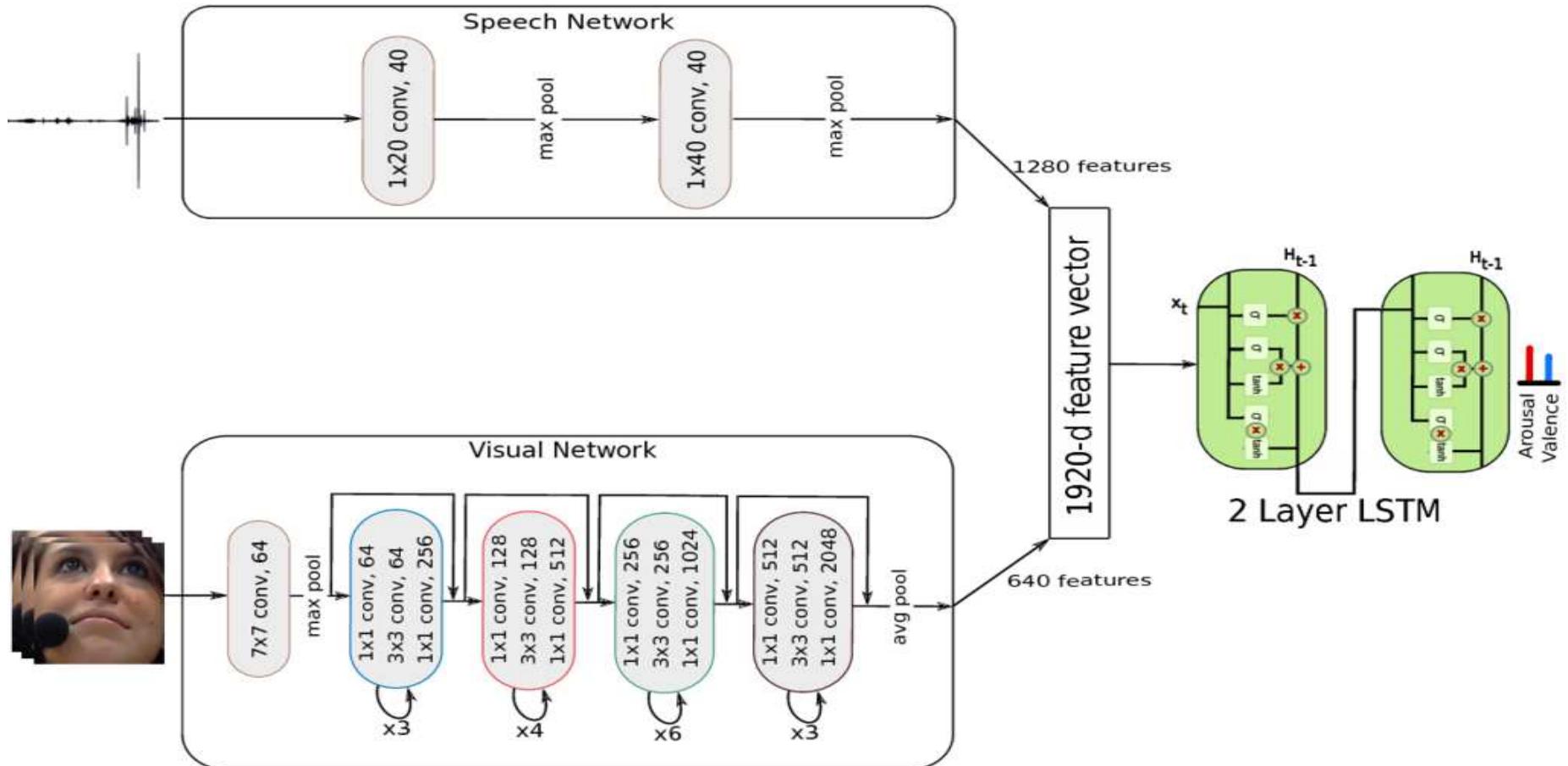
P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. W. Schuller and S. Zafeiriou, "End-to-End Multimodal Emotion Recognition Using Deep Neural Networks," in IEEE Journal of Selected Topics in Signal Processing, vol. 11, no. 8, pp. 1301-1309, Dec. 2017, doi: 10.1109/JSTSP.2017.2764438.

# End-to-End Multimodal Emotion Recognition



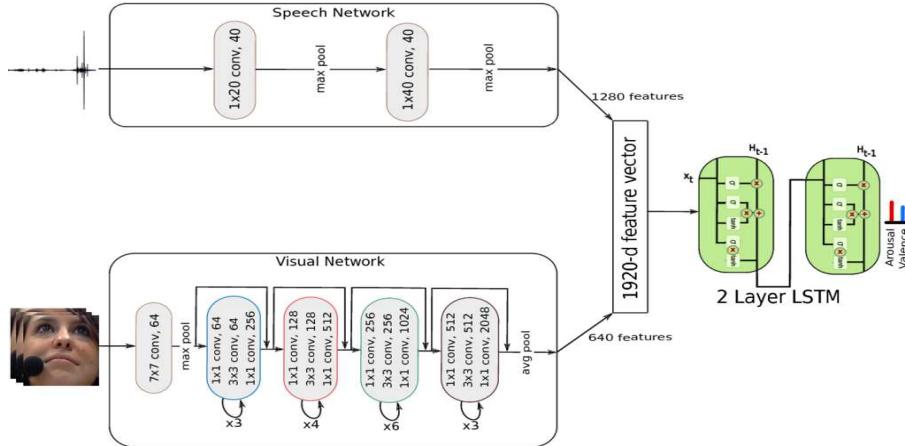
P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. W. Schuller and S. Zafeiriou, "End-to-End Multimodal Emotion Recognition Using Deep Neural Networks," in IEEE Journal of Selected Topics in Signal Processing, vol. 11, no. 8, pp. 1301-1309, Dec. 2017, doi: 10.1109/JSTSP.2017.2764438.

# End-to-End Multimodal Emotion Recognition



P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. W. Schuller and S. Zafeiriou, "End-to-End Multimodal Emotion Recognition Using Deep Neural Networks," in IEEE Journal of Selected Topics in Signal Processing, vol. 11, no. 8, pp. 1301-1309, Dec. 2017, doi: 10.1109/JSTSP.2017.2764438.

# End-to-End Multimodal Emotion Recognition



ResNet

## Speech Network:

- Raw waveform (6s)
- Temporal Convolution ( $F=20$ ): space time finite impulse filters with 5ms window in order to extract fine-scale spectral information
- Pooling across time: Half-wave rectifier (analogous to the cochlear transduction step in the human ear) and then downsampled to 8kHz (pool size = 2)
- Temporal convolution ( $F=40$ ): Extract more long-term characteristics of the speech signal
- Max pooling across channels
- Dropout

# End-to-End Multimodal Emotion Recognition

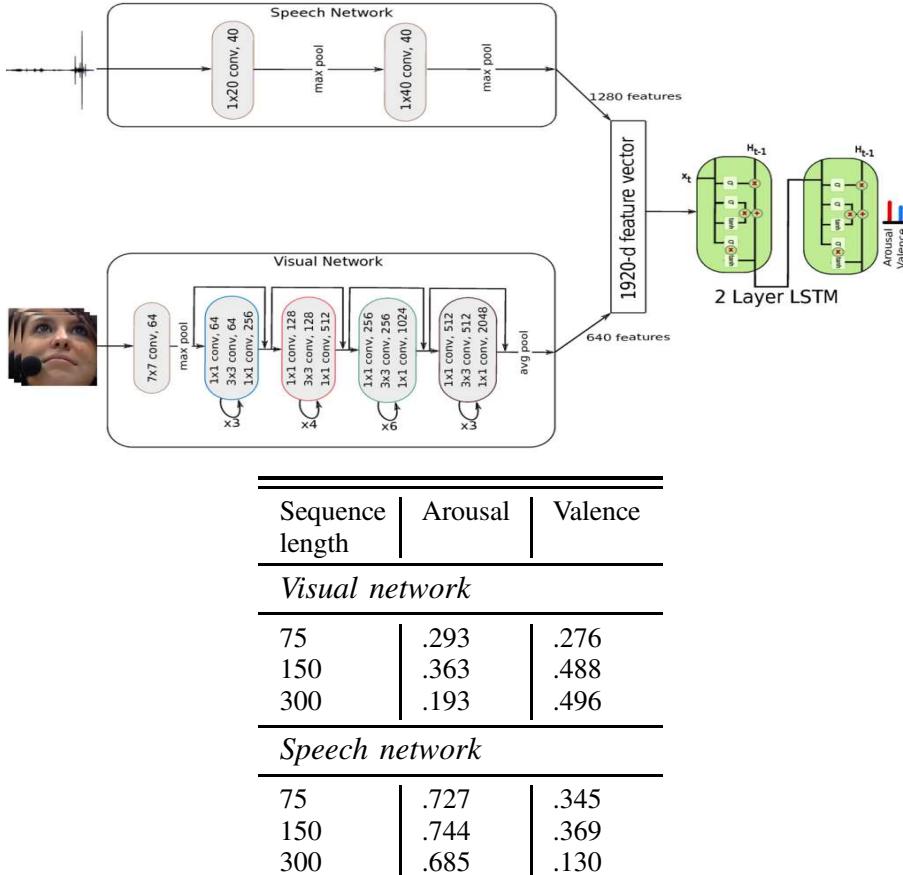


Table II: Results (in terms of  $\rho_c$ ) on arousal and valence after 60 epochs when varying sequence length for speech and visual networks.

**Objective function:**  
Instead of MSE, correlation

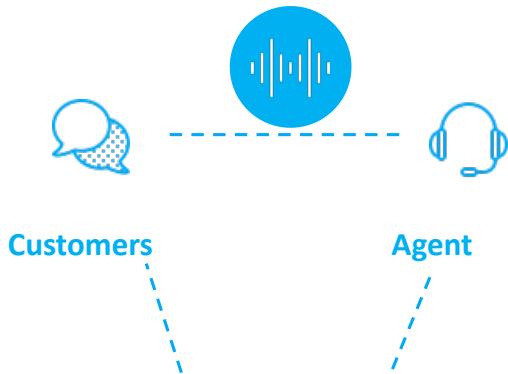
$$\begin{aligned} \mathcal{L}_c &= 1 - \rho_c = 1 - \frac{2\sigma_{xy}^2}{\sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2} \\ &= 1 - 2\sigma_{xy}^2 \psi^{-1} \end{aligned} \quad (3)$$

where  $\psi = \sigma_x^2 + \sigma_y^2 + (\mu_x - \mu_y)^2$  and  $\mu_x = \mathbb{E}(\mathbf{x})$ ,  $\mu_y = \mathbb{E}(\mathbf{y})$ ,  $\sigma_x^2 = \text{var}(\mathbf{x})$ ,  $\sigma_y^2 = \text{var}(\mathbf{y})$  and  $\sigma_{xy}^2 = \text{cov}(\mathbf{x}, \mathbf{y})$ . Thus, to minimise  $\mathcal{L}_c$  (or maximise  $\rho_c$ ), we backpropagate the gradient of the last layer weights with respect to  $\mathcal{L}_c$ ,

$$\frac{\partial \mathcal{L}_c}{\partial \mathbf{x}} \propto 2 \frac{\sigma_{xy}^2 (\mathbf{x} - \mu_y)}{\psi^2} + \frac{\mu_y - \mathbf{y}}{\psi}, \quad (4)$$

where all vector operations are done element-wise.

# Customer Satisfaction & Emotion



A typical 5-point Likert scale:

- Completely Dissatisfied
- Somewhat Dissatisfied
- Neutral
- Somewhat Satisfied
- Completely Satisfied

- **H1 Customers' emotions and CSAT response**
  - **H1a:** Customers expressing positive emotions respond more to CSAT questionnaires;
  - **H1b:** Customers expressing negative emotions respond less to CSAT questionnaires;
  - **H1c:** Customers expressing anger respond less to CSAT questionnaires;
- **H2 Customers' emotions and self-reported satisfaction**
  - **H2a:** Customers expressing positive emotions report higher satisfaction;
  - **H2b:** Customers expressing negative emotions report lower satisfaction;
  - **H2c:** Customers expressing anger report lower satisfaction;
- **H3 Customers' emotional profiles and CSAT response rate**
  - **H3a:** Customers manifesting upward positive valence dynamics (more positive emotions towards the end of the call) exhibit a higher CSAT response rate compared to flat or negative dynamics;
  - **H3b:** Customers manifesting downward negative valence dynamics (fewer negative emotions towards the end of the call) exhibit a higher CSAT response rate compared to flat or positive dynamics;
  - **H3c:** Customers manifesting downward anger dynamics (fewer anger events towards the end of the call) exhibit a higher CSAT response rate;
- **H4 Customers' emotional profiles and self-reported satisfaction**
  - **H4a:** Customers manifesting upward positive valence dynamics report higher satisfaction;
  - **H4b:** Customers manifesting downward negative valence dynamics report higher satisfaction;
  - **H4c:** Customers manifesting downward anger dynamics report higher satisfaction;

# Customer Satisfaction & Emotion

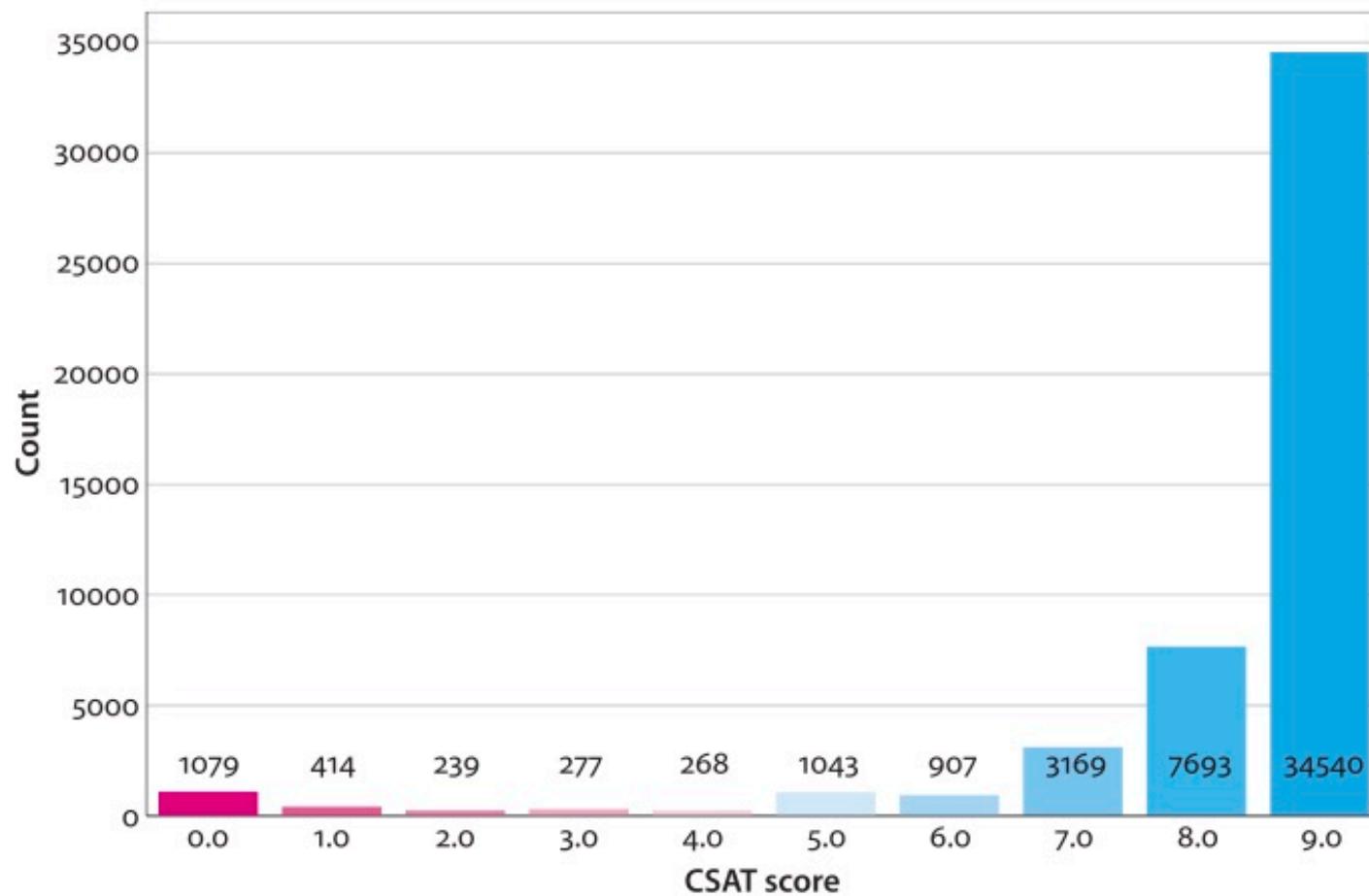
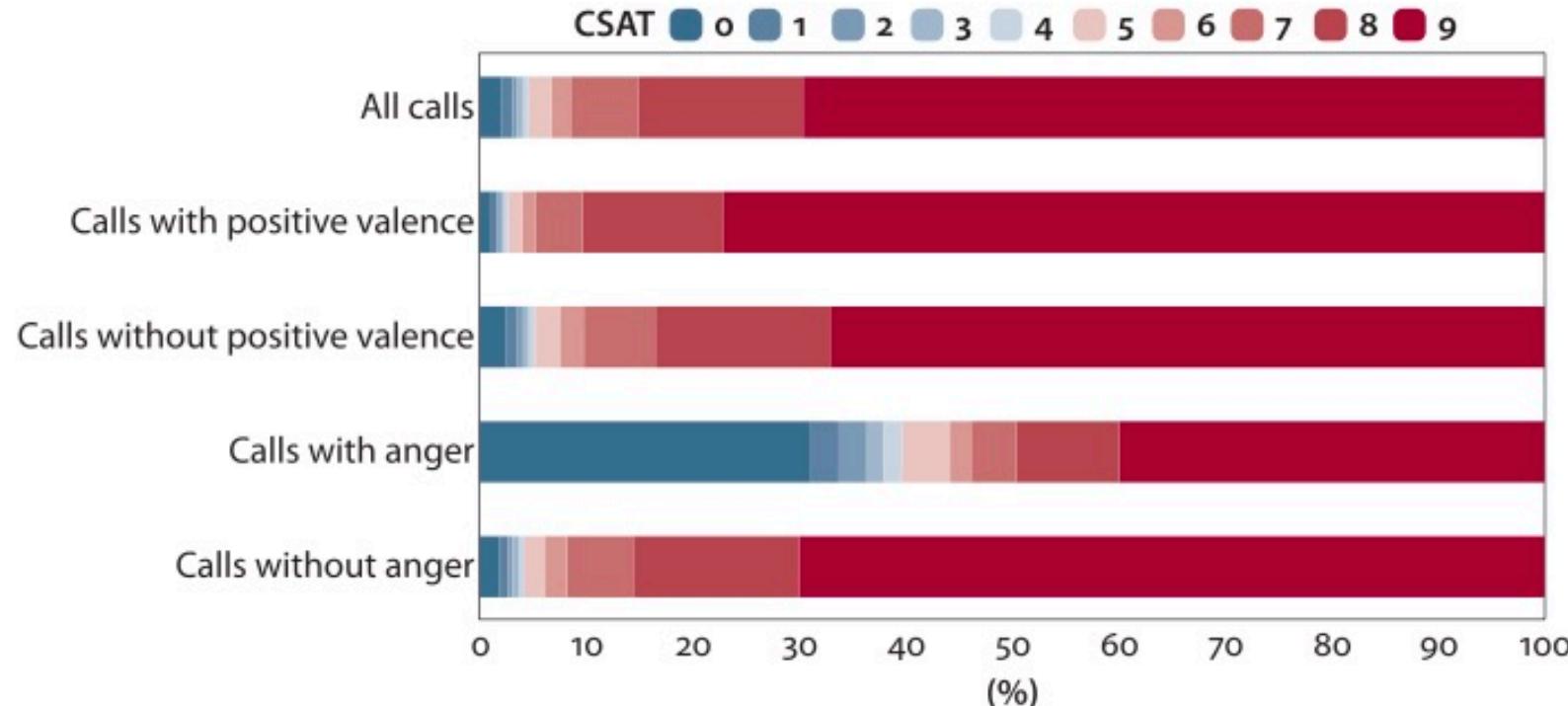


Figure 1. CSAT scoring distribution

Eric Bolo, Muhammad Samoul, Nicolas Seichepine, Mohamed Chetouani. Quietly Angry, Loudly Happy: Self-Reported Customer Satisfaction Vs. Automatically Detected Emotion In Contact Center Calls. *Interaction Studies*, 2023, 24 (1), pp.168-192. 10.1075/is.22038.bol. hal-04109350

# Customer Satisfaction & Emotion



**Figure 3.** CSAT scoring distribution with respect to emotional analysis. Number of calls are reported in Table 6

# Customer Satisfaction & Emotion

$$\text{CSAT response rate} = \frac{\text{Nbr of Calls with CSAT scoring}}{\text{Total Nbr of Calls}}$$

$$\Delta = \frac{N_{end} - N_{beg}}{N_{beg} + N_{mid} + N_{end}}$$

$$\text{aggH-CSAT} = \frac{\sum_{i=6}^9 CSATscore(i)}{\sum_{i=1}^9 CSATscore(i)}$$

Table 4: CSAT response rate and emotion in respect to the emotion analysis. \* indicates a significant difference with the baseline ( $p<0.05$ )

# of calls	With CSAT scoring	Without CSAT Scoring	CSAT response rate
All calls	49 629	111 001	0.29
Calls with positive valence	13 491	27 212	<b>0.33*</b>
Calls with negative valence	31 722	67 907	<b>0.32*</b>
Calls with anger	892	3 586	<b>0.20*</b>

# Customer Satisfaction & Emotion

Table 5: Proportion of detected emotion in calls with and without CSAT scoring.

Proportion in	With CSAT scoring ( $\pi_1$ )	Without CSAT Scoring ( $\pi_2$ )	$\pi_1$ vs. $\pi_2$
Calls with positive valence	0.27	0.24	$z = 11.36, p < 0.05$
Calls with negative valence	0.63	0.61	$z = 10.46, p < 0.05$
Calls with anger	0.01	0.03	$z = -16.12, p < 0.05$

Table 6: Number of calls per CSAT scoring in respect of affective computing analysis (see also figure 3) and the Aggregate High Customer Satisfaction score (aggH-CSAT) in respect to the emotion analysis. \* indicates a significant difference with the baseline ( $p < 0.05$ )

	0	1	2	3	4	5	6	7	8	9	aggH-CSAT
All calls	1079	414	239	277	268	1043	907	3169	7693	34540	0.93
Calls with positive valence	150	94	44	43	53	182	165	602	1814	10344	<b>0.96*</b>
Calls without positive valence	929	320	195	234	215	861	742	2567	5879	24196	<b>0.92*</b>
Calls with negative valence	883	299	180	184	193	747	556	2012	4819	21849	<b>0.92*</b>
Calls without negative valence	196	115	59	93	75	296	351	1157	2874	12691	<b>0.95*</b>
Calls with anger	184	17	17	11	10	39	26	51	115	422	<b>0.69*</b>
Calls without anger	895	397	222	266	258	1004	881	3118	7578	34118	<b>0.94*</b>

# Customer Satisfaction & Emotion

Table 7: CSAT response rate with respect to the customers' emotional profiles. \* indicates a significant difference with the baseline ( $p<0.05$ )

# of calls	With CSAT scoring	Without CSAT Scoring	CSAT response rate
All calls	49 629	111 001	0.29
Calls with upward positive valence dynamics	11 896	24 577	<b>0.32*</b>
Calls with downward positive valence dynamics	302	521	<b>0.36*</b>
Calls with upward negative valence dynamics	15 489	40 017	<b>0.27*</b>
Calls with downward negative valence dynamics	5 854	9 621	<b>0.37*</b>
Calls with upward anger dynamics	437	2 155	<b>0.16*</b>
Calls with downward anger dynamics	82	230	<b>0.38*</b>

# Summary

► Multiple features

► Deep learning architectures for speech processing

► Interplay between speech and machine learning approaches

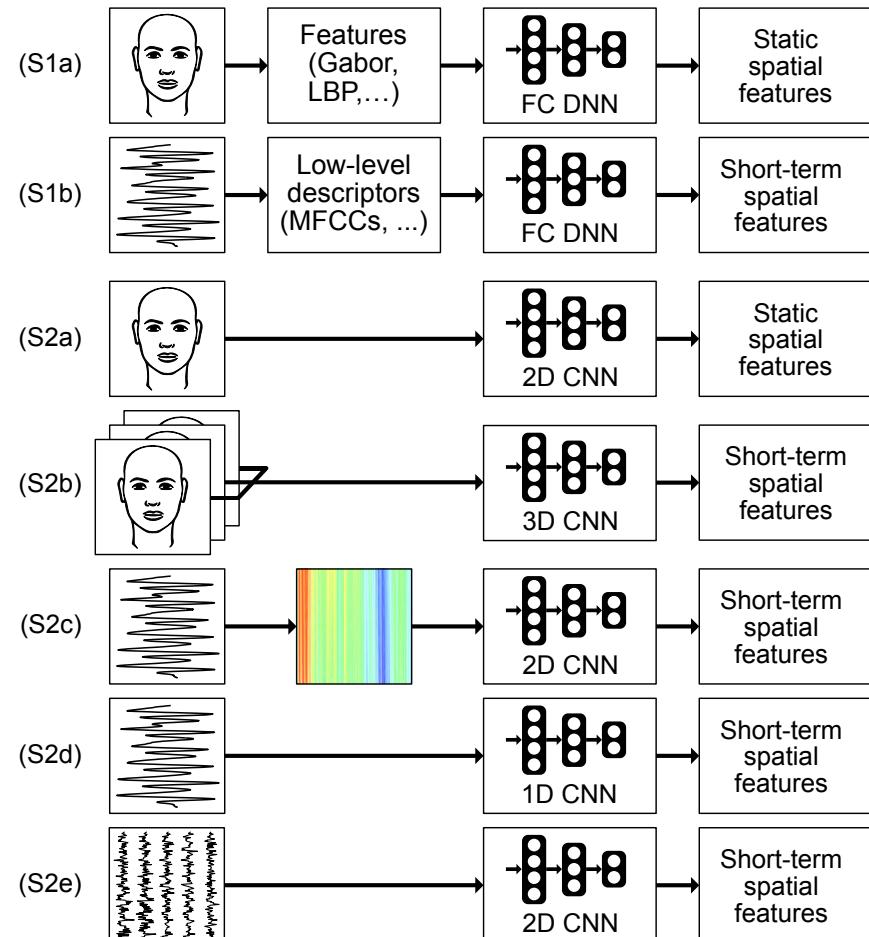


Fig. 2. Applications of deep learning for spatial feature learning with fully-connected DNNs (S1a–S1b) and CNNs (S2a–S2e).

# Next steps

## Practicals

### Intention recognition from speech signals

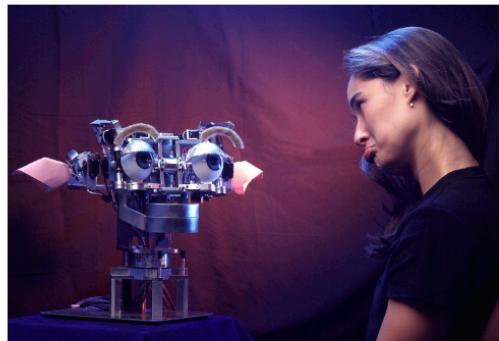


Fig. 2. Kismet is an expressive robotic creature designed for natural social interaction with people. See text.

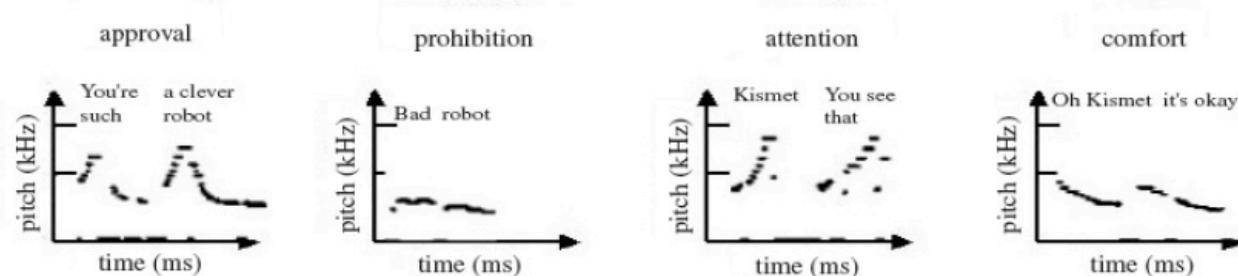


Fig. 2. Fernald's prototypical prosodic contours shown in robot directed speech for approval, attentional bid, prohibition, and soothing.

# Next steps

## Practicals

### Intention recognition from speech signals

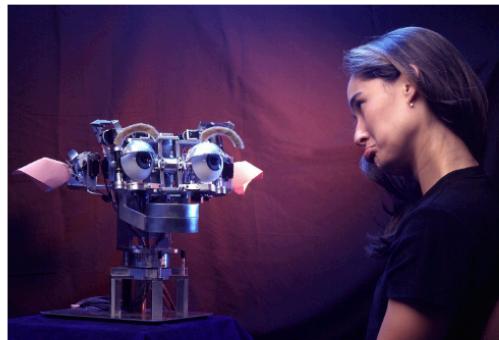


Fig. 2. Kismet is an expressive robotic creature designed for natural social interaction with people. See text.

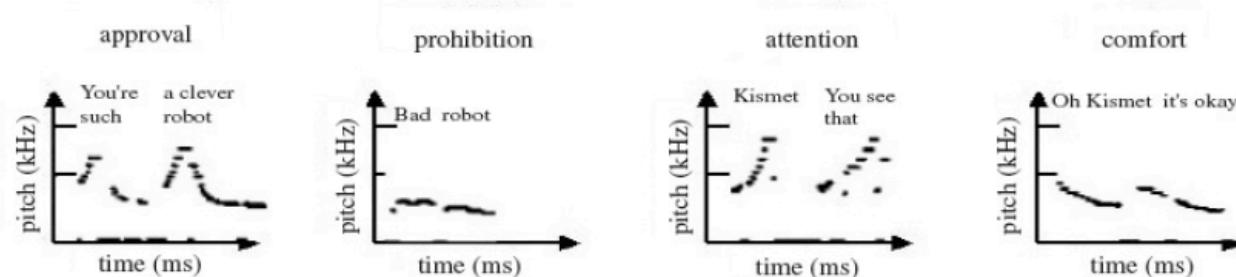


Fig. 2. Fernald's prototypical prosodic contours shown in robot directed speech for approval, attentional bid, prohibition, and soothing.

# Thank you for your attention



Questions?

©FP7 Michelangelo