

AUDIOVISUAL DISCRIMINATION BETWEEN LAUGHTER AND SPEECH

Stavros Petridis, Maja Pantic
sp104@doc.ic.ac.uk, m.pantic@imperial.ac.uk

Department of Computing

Berlin 27/2/2009

Motivation

- Powerful affective and social signal
- Most frequently annotated non-verbal behaviour in recorded natural speech
 - 8.6% of the time a person vocalizes is spent on laughing
 - 0.8% is spent on laughing while talking

Previous work

- Audio-only Laughter Detection
 - Audio features only (MFCC, prosody)
 - Hidden Markov Models (HMMs) / SVM /GMM
- Audiovisual laughter Detection
 - Audio features (MFCC)
 - Visual features (lip lengths, lip angles)
 - Improved performance
 - Tested on 3 sequences

Objective

- To use standard tools to discriminate laughter from speech using an audiovisual approach

Audiovisual Recognition

- Audiovisual speech recognition
 - Audio + Visual features (e.g. lip contours)
 - Hidden Markov Models (HMMs) + Variants of HMMs
 - Feature level
 - Enhanced performance
- Audiovisual affect recognition
 - Recognition of six basic + non basic emotions
 - HMMs, GMMs, RNN, Bayesian Networks
 - Decision level, model-based fusion

Dataset

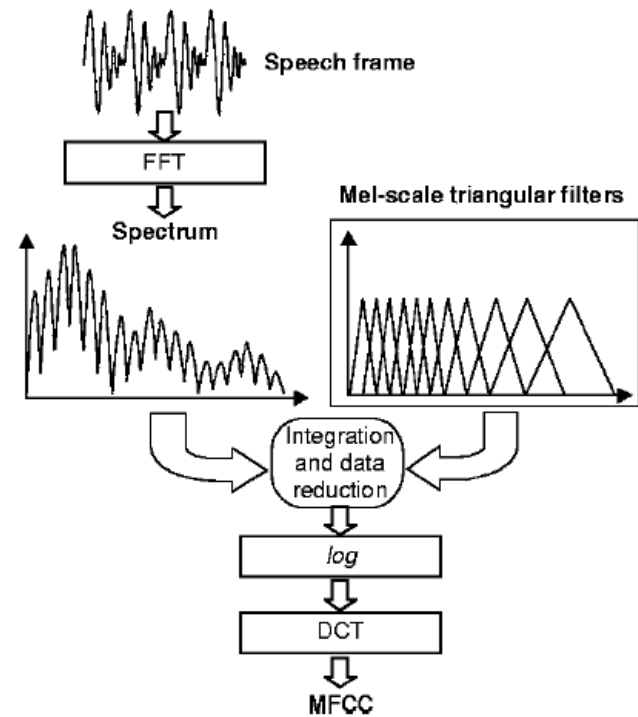
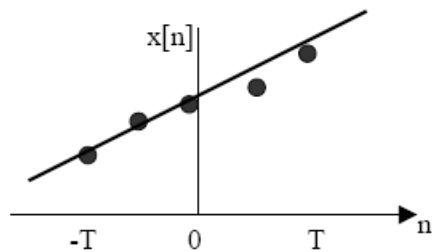
- AMI Meeting Corpus
- 10 subjects – spontaneous expressions
- 67 voiced laughters (100.88 sec)
- 48 unvoiced laughters (60.76 sec)
- 93 speech segments (177.86 sec)
- 720 x 576 pixels, 25 FPS





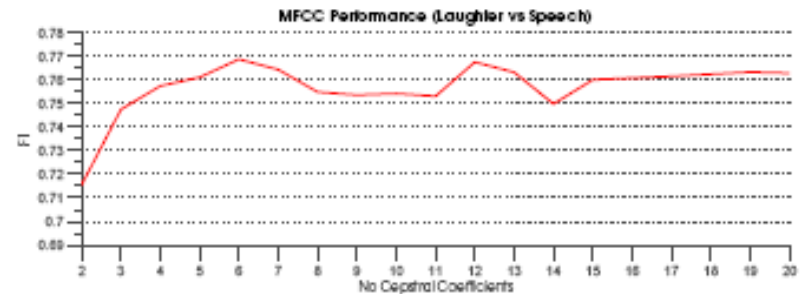
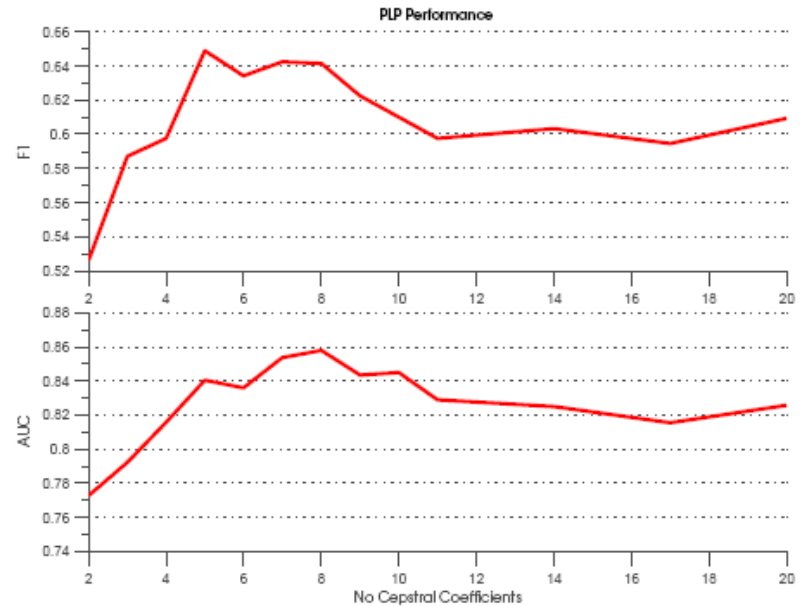
Audio Features

- Spectral Features
 - MFCC
 - PLP
- Non Spectral Features (Prosodic)
 - Pitch
 - Energy
- Delta Features



AUDIO FEATURES

- Spectral Features
 - 6 MFCC coeff. + 6 delta
 - Same results with PLP
 - Frame rate: 50 FPS,
50% overlap
- Prosodic Features
 - Pitch
 - Energy



VISUAL FEATURES - TRACKING

- 20 facial points
- Eyes: 4 points each
- Eyebrows: 2 points each
- Nose: 3 points
- Mouth: 4 points
- Chin: 1 point



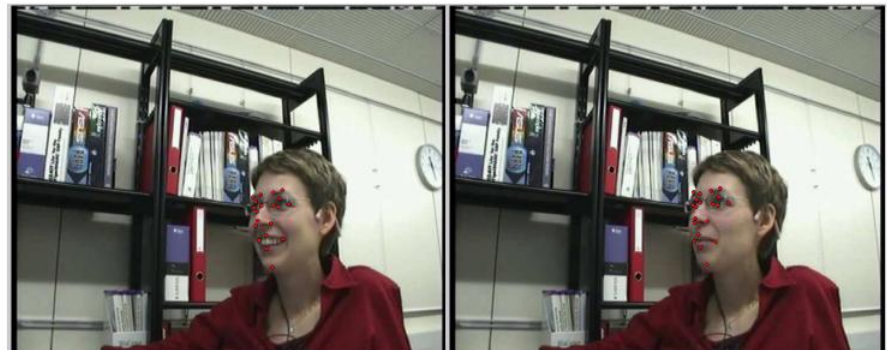
(a) Frame 2

(b) Frame 42



(c) Frame 82

(d) Frame 122

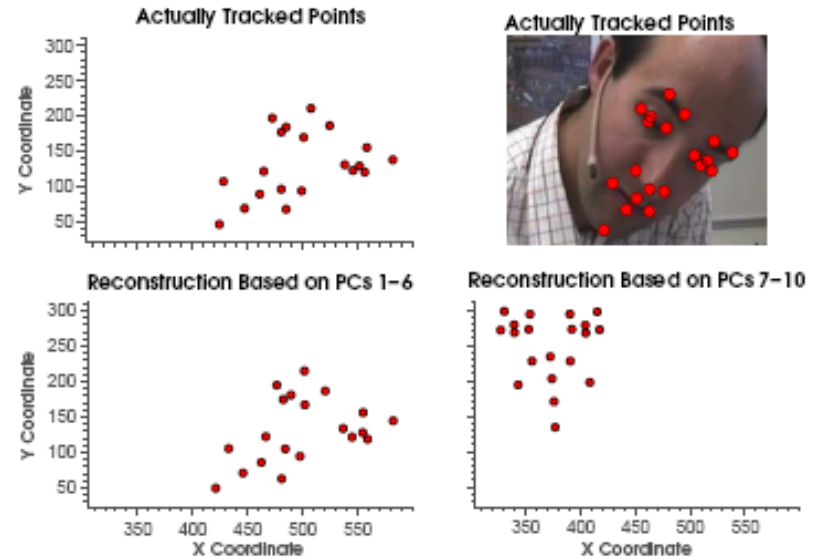


(e) Frame 162

(f) Frame 176

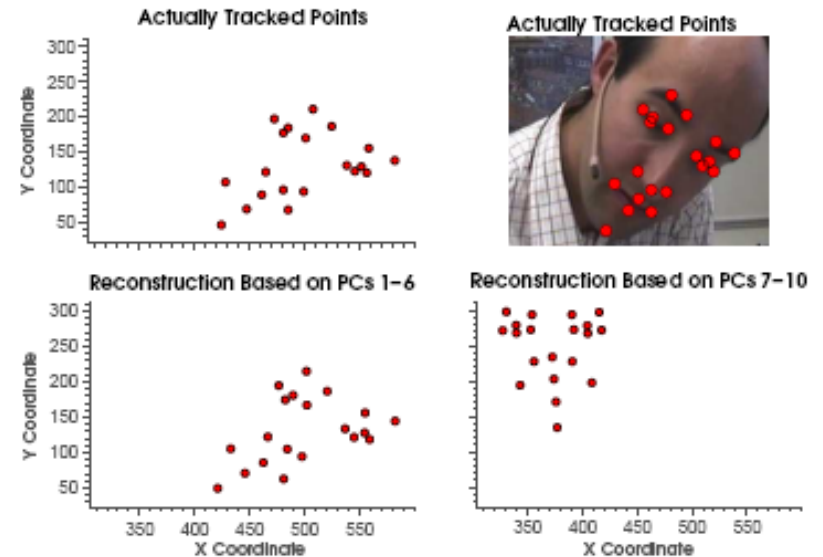
VISUAL FEATURES – DECOUPLING OF HEAD AND FACE

- Concatenate the (x,y) coord.
→ 40-dim vector
- k frames → k x 40 data matrix
- Find the PCs of the data matrix (PCA)
- Point Distribution Model with no alignment



VISUAL FEATURES – DECOUPLING OF HEAD AND FACE

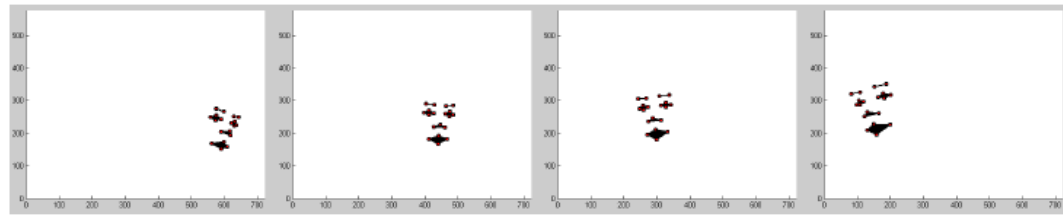
- Greatest variance of the data lies on the first PCs
- Greatest variance = head (rigid) movements
- First N(=5) PCs: reflect head movements (large variation)
- N+1...M(=6-10) PCs: reflect facial expressions (small variation)



$$b = (x - \bar{x})P$$

$$x \approx \hat{x} = \bar{x} + bP^T$$

- P: eigenvector matrix (40xN)
- b: shape parameters, (1xN)

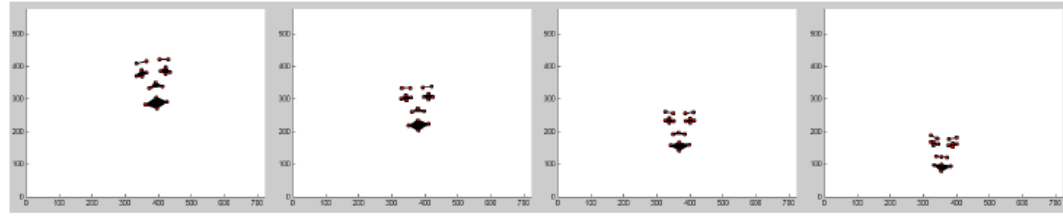


(a)

(b)

(c)

(d)

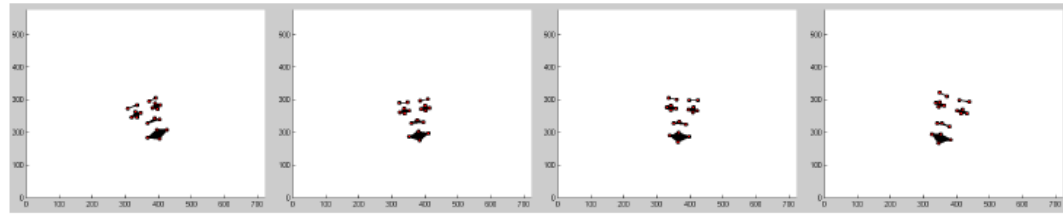


(e)

(f)

(g)

(h)

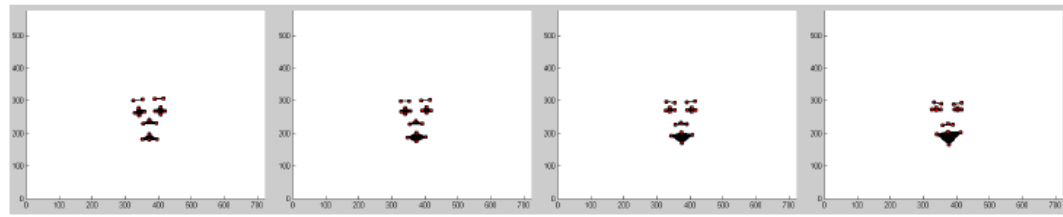


(i)

(j)

(k)

(l)

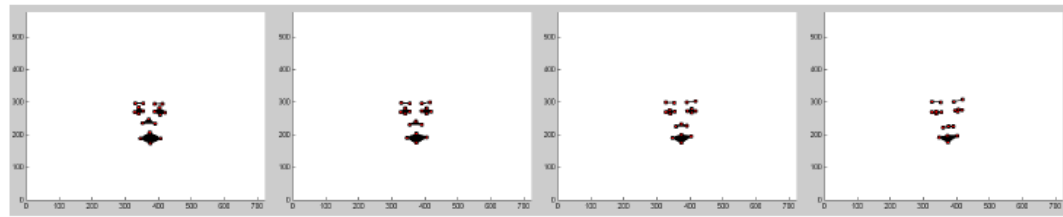


(m)

(n)

(o)

(p)

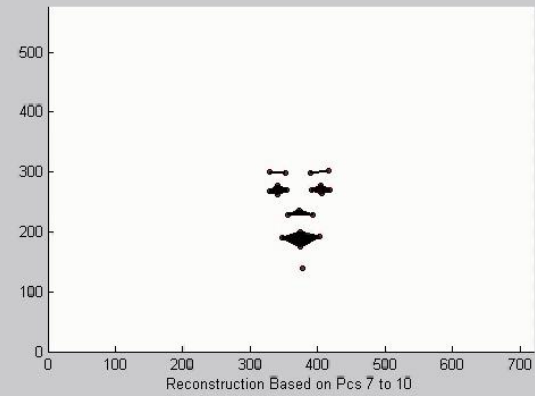
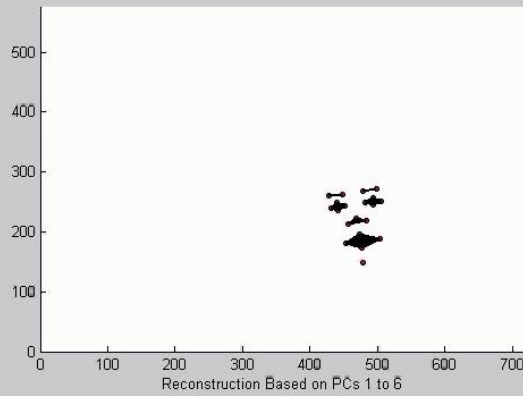
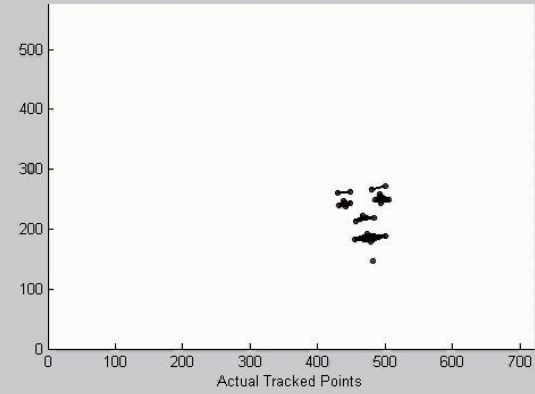
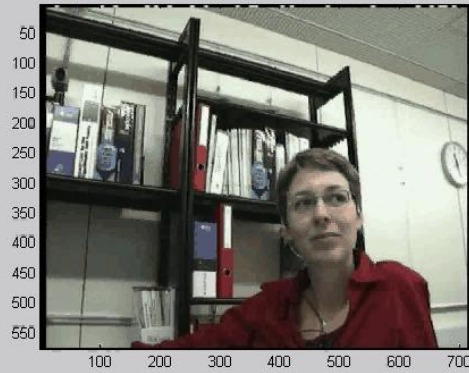


(q)

(r)

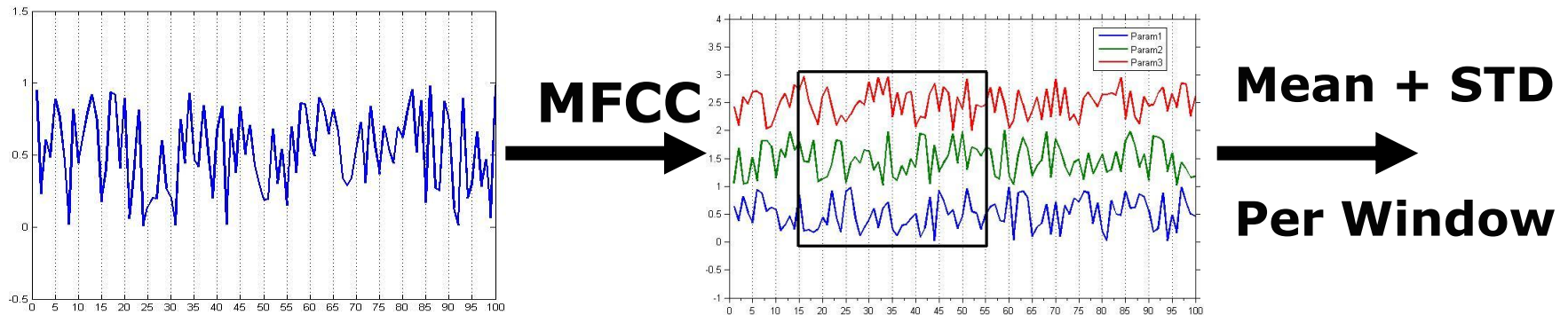
(s)

(t)



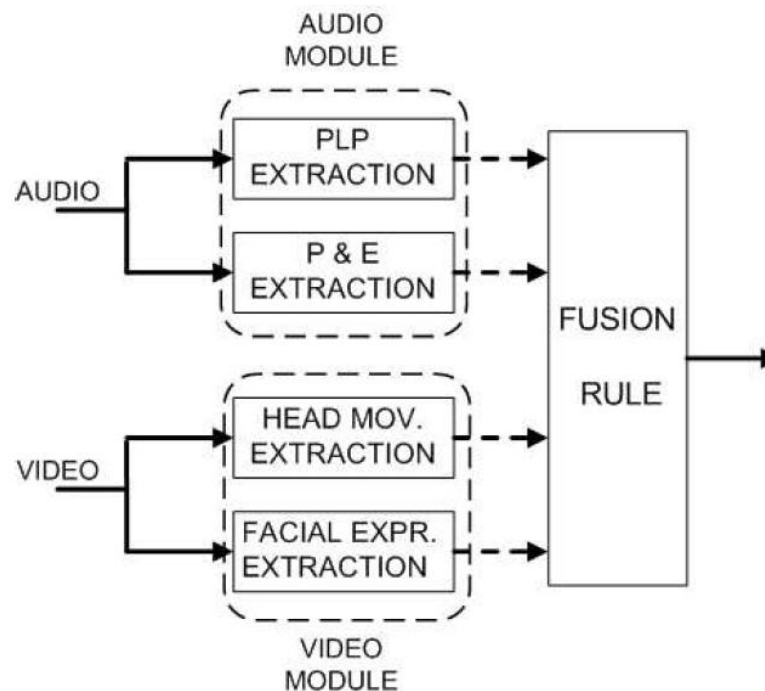
FEATURE EXTRACTION

- Visual features
 - b1 – b5, b6 – b10 per frame
- Audio features
 - Mean + STD of MFCC, Pitch, Energy over 320ms window
 - Pitch Unvoiced Ratio



SYSTEM OVERVIEW

- Leave-one-subject-out cross validation
- 2 / 3 classes: Laughter vs Speech
- Decision / Feature Level Fusion

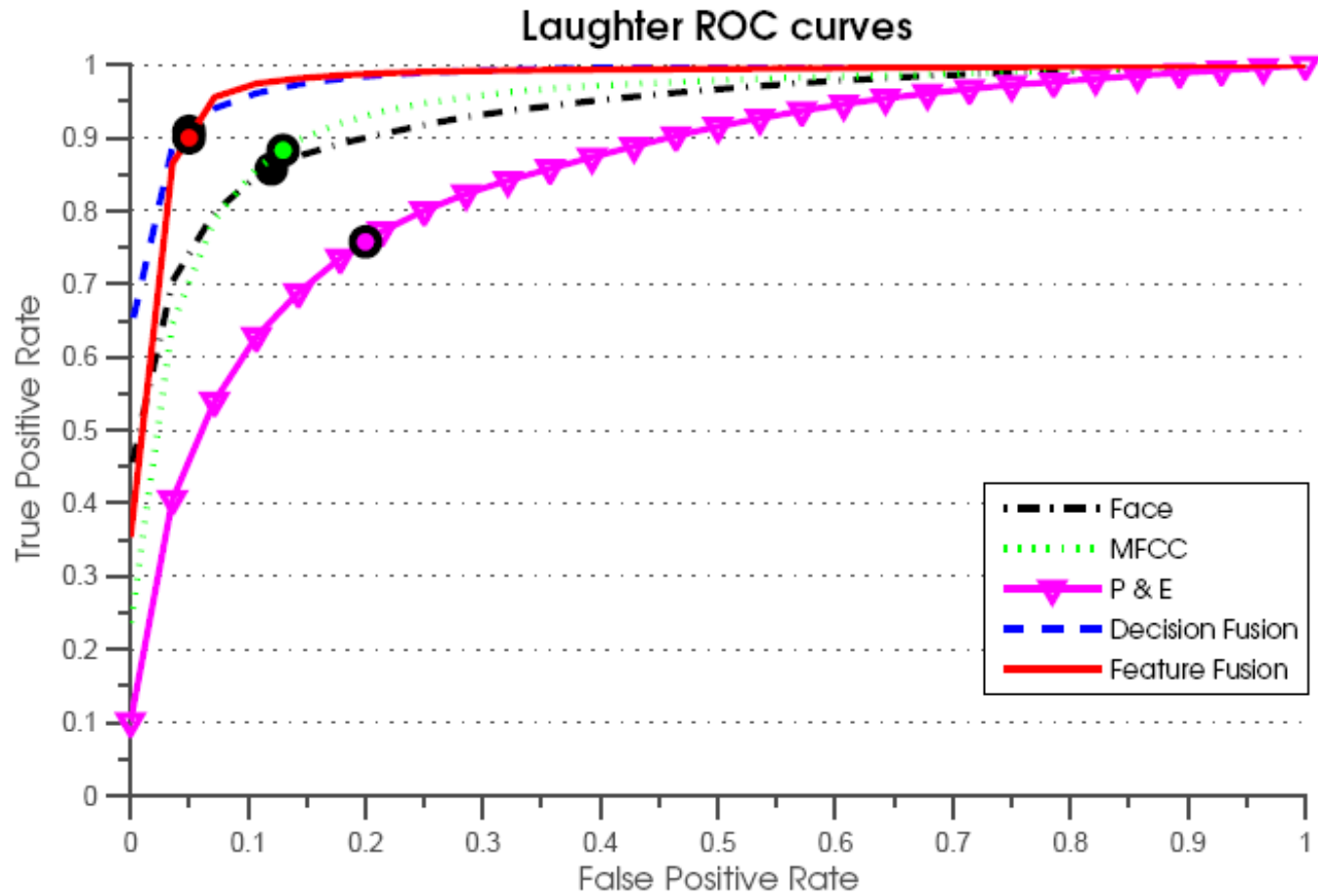


RESULTS – LAUGHTER vs SPEECH

<i>Cues</i>	<i>Features</i>	<i>F1</i>	<i>CR</i>
<i>FACE</i>	<i>b6 – b10</i>	<i>84.28</i>	<i>85.06</i>
<i>HEAD</i>	<i>b1 – b5</i>	<i>50.71</i>	<i>51.02</i>
<i>PROSODY</i>	<i>Pitch + energy</i>	<i>72.69</i>	<i>74.43</i>
<i>CEPSTRAL FEATURES</i>	<i>MFCC</i>	<i>86.16</i>	<i>86.69</i>

<i>Fusion</i>	<i>F1</i>	<i>CR</i>
<i>FACE + MFCC (DF)</i>	<i>93.17</i>	<i>93.37</i>
<i>FACE + MFCC + P&E (FF)</i>	<i>94.37</i>	<i>94.59</i>

RESULTS – AUDIOVISUAL



RESULTS – Voiced Laughter vs Unvoiced Laughter vs Speech

<i>Cues</i>	<i>Features</i>	<i>F1 Unvoiced</i>	<i>F1 Voiced.</i>	<i>F1 Speech</i>	<i>CR</i>
<i>FACE</i>	<i>b6 – b10</i>	<i>37.21</i>	<i>57.48</i>	<i>84.39</i>	<i>67.03</i>
<i>HEAD</i>	<i>b1 – b5</i>	<i>19.32</i>	<i>34.56</i>	<i>49.72</i>	<i>37.27</i>
<i>PROSODY</i>	<i>Pitch + energy</i>	<i>57.24</i>	<i>59.65</i>	<i>74.39</i>	<i>66.20</i>
<i>CEPSTRAL FEATURES</i>	<i>MFCC</i>	<i>57.97</i>	<i>69.14</i>	<i>85.68</i>	<i>74.83</i>

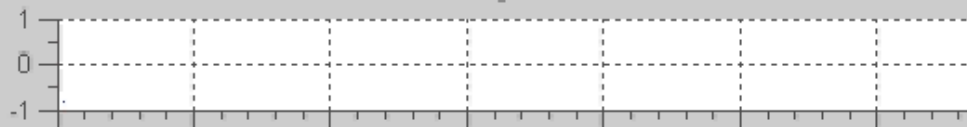
<i>Fusion</i>	<i>F1 Unvoiced</i>	<i>F1 Voiced</i>	<i>F1 Speech</i>	<i>CR</i>
<i>FACE + MFCC + PE (DF)</i>	<i>66.14</i>	<i>74.62</i>	<i>91.89</i>	<i>81.00</i>
<i>FACE + MFCC + P&E + HEAD(FF)</i>	<i>66.80</i>	<i>78.86</i>	<i>94.14</i>	<i>83.51</i>

FEATURE SELECTION

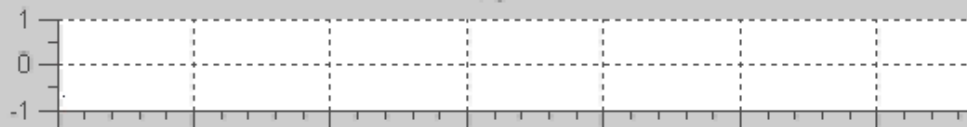
Speech	Voiced	Unvoiced
B7	B7	Mean Energy
Mean MFCC2	Std Energy	Mean MFCC6
Std MFCC2	Mean MFCC1	Mean MFCC1
Std MFCC4	B10	B7
B10	B8	Mean MFCC2
Std MFCC3	B4	B6
Mean Pitch	Mean MFCC3	Std Pitch



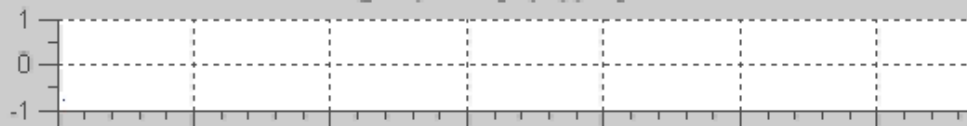
Audio



Video



Decision Level Fusion



Frames

Conclusions

- Audiovisual fusion leads to improved performance
- Best cues: cepstral features, facial expressions
- High accuracy for speech vs laughter
- Discriminating 2 types of laughter is challenging

THANK YOU! 😊