

Features for Audio-Visual Speech Recognition

Doctoral thesis by Iain Matthews,
University of East Anglia reviewed
by Tobias Andersen,
Lab. of Computational Eng., HUT

Audiovisual speech primer

- Why is visual speech interesting?
 - ◆ It improves the "SNR" of speech
 - ◆ It helps distinguish speakers in a multi-speaker environment
 - ◆ The best speech decoder – the human brain – always try to use it

Audiovisual speech primer



Review overview

- Visual speech recognition
 - ◆ low level
 - ◆ Area
 - ◆ Height and width
 - ◆ Area histogram
 - ◆ Scale histogram
 - ◆ high level
 - ◆ Active shape model (ASP)
 - ◆ Active appearance model (AAP)
- HMM
- AV integration

Low level speech recognition

- Sieves
 - ◆ Erosion/dilation
 - ◆ Opening/closing
 - ◆ Granules and scales

Erosion

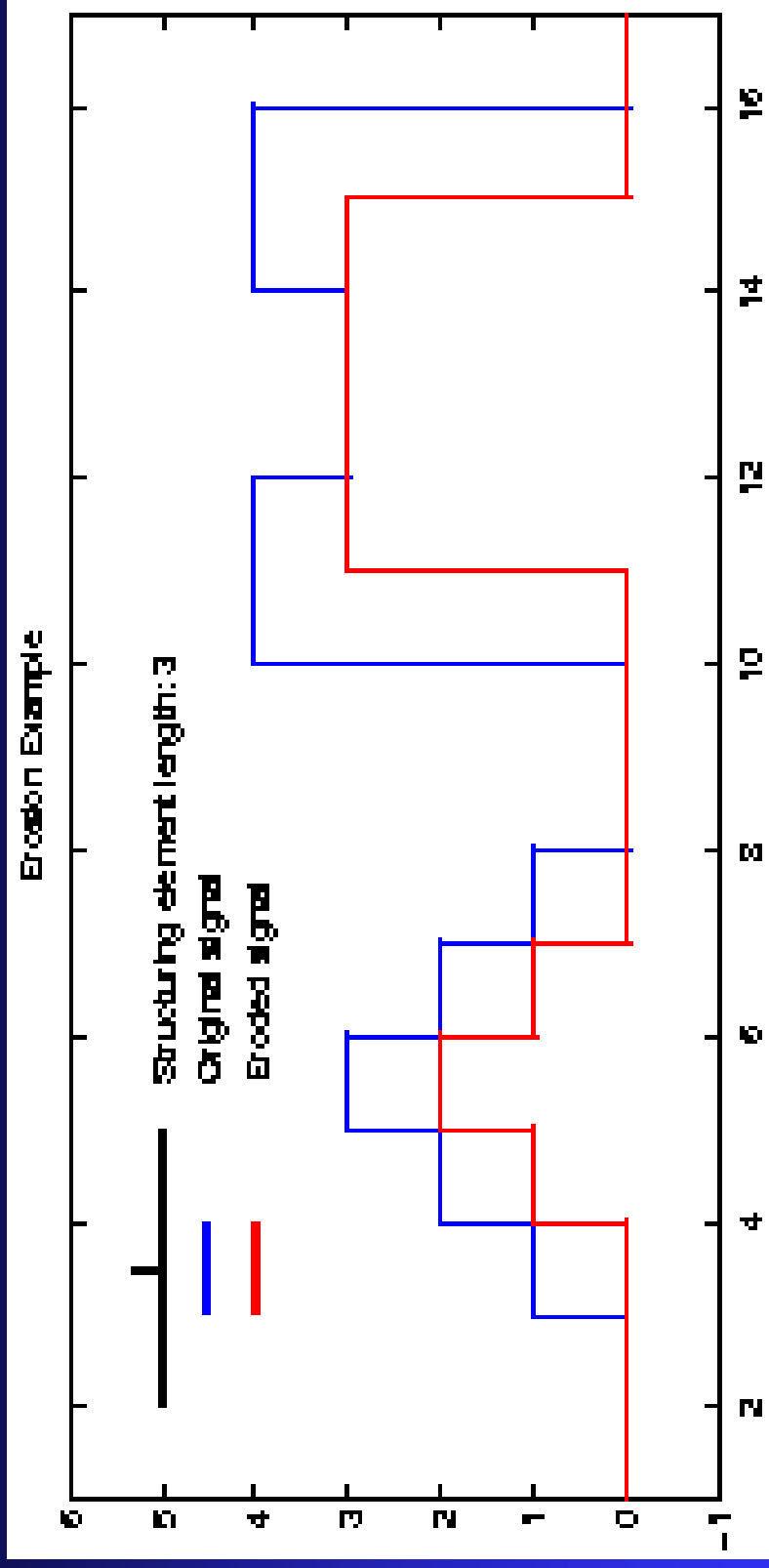


Figure 4.2: Example of greyscale erosion.

Dilation

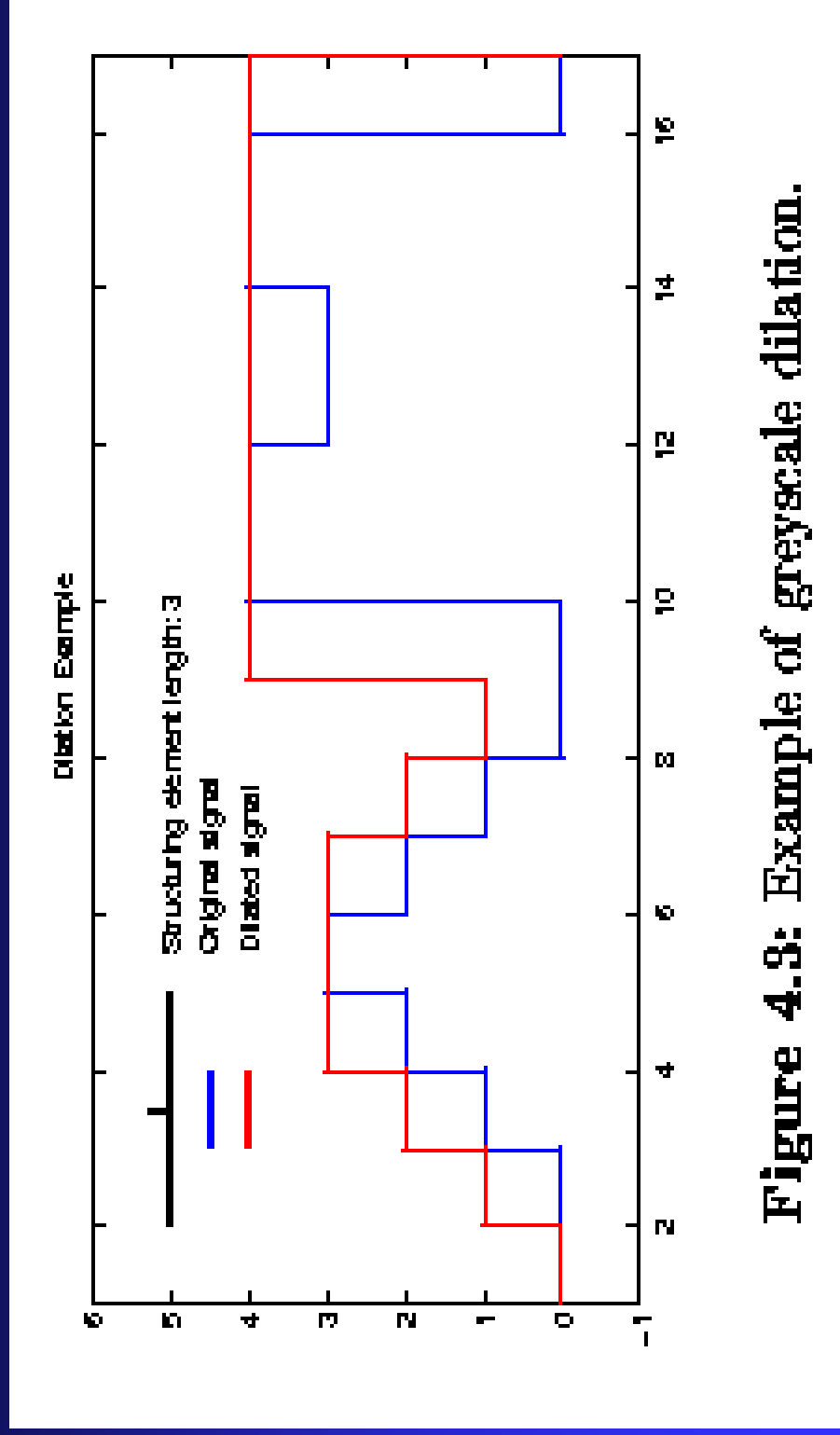


Figure 4.3: Example of grayscale dilation.

Opening and closing

- Opening
 - ◆ Erosion followed by dilation
 - ◆ Removes positive extrema
- Closing
 - ◆ Dilation followed by erosion
 - ◆ Removes negative extrema

Recursive median filter

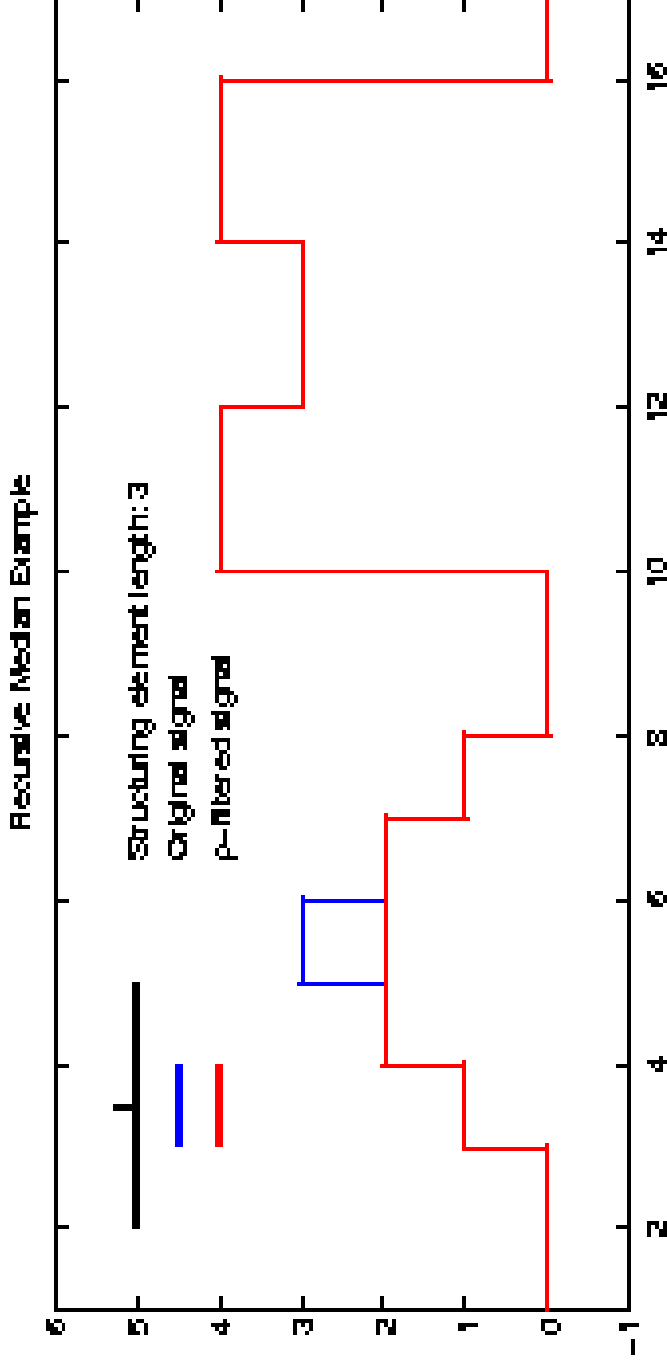
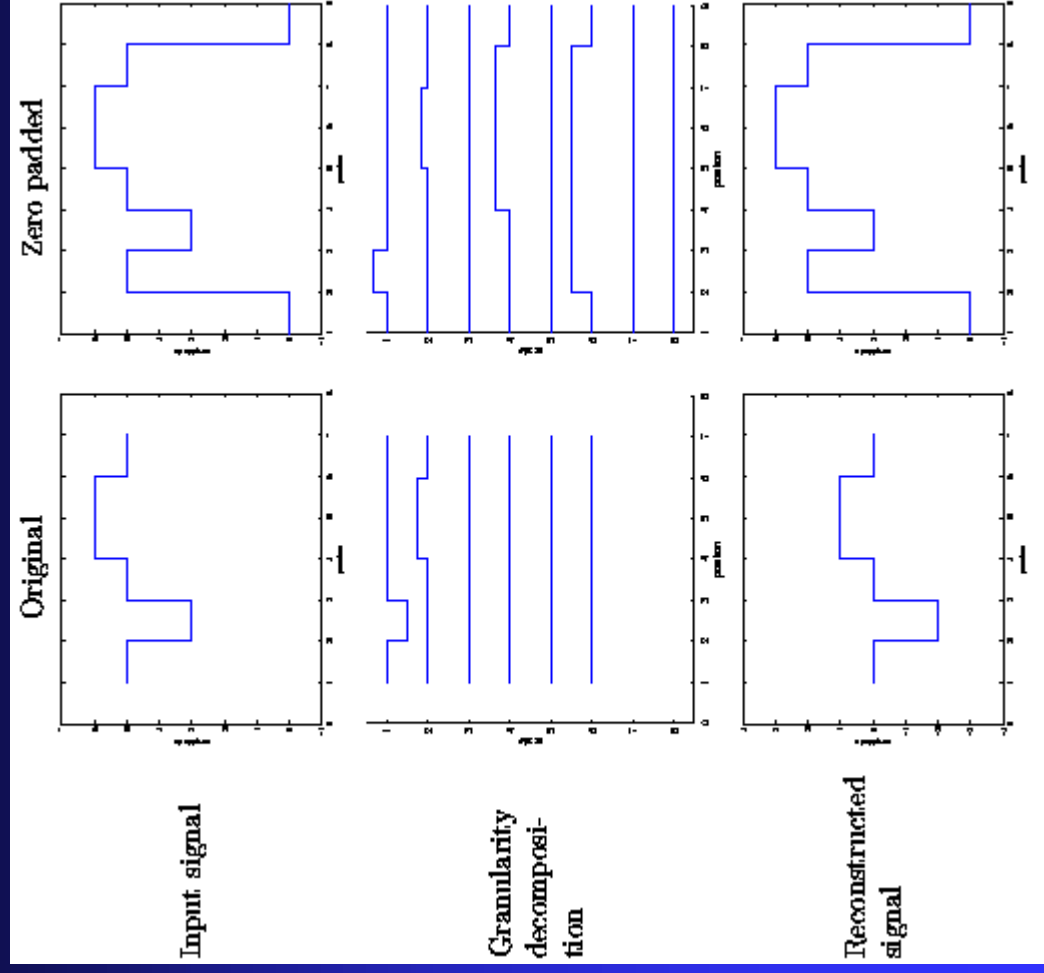


Figure 4.8: Example of a recursive median filter with three

Preservation of baseling



Sieve decomposition

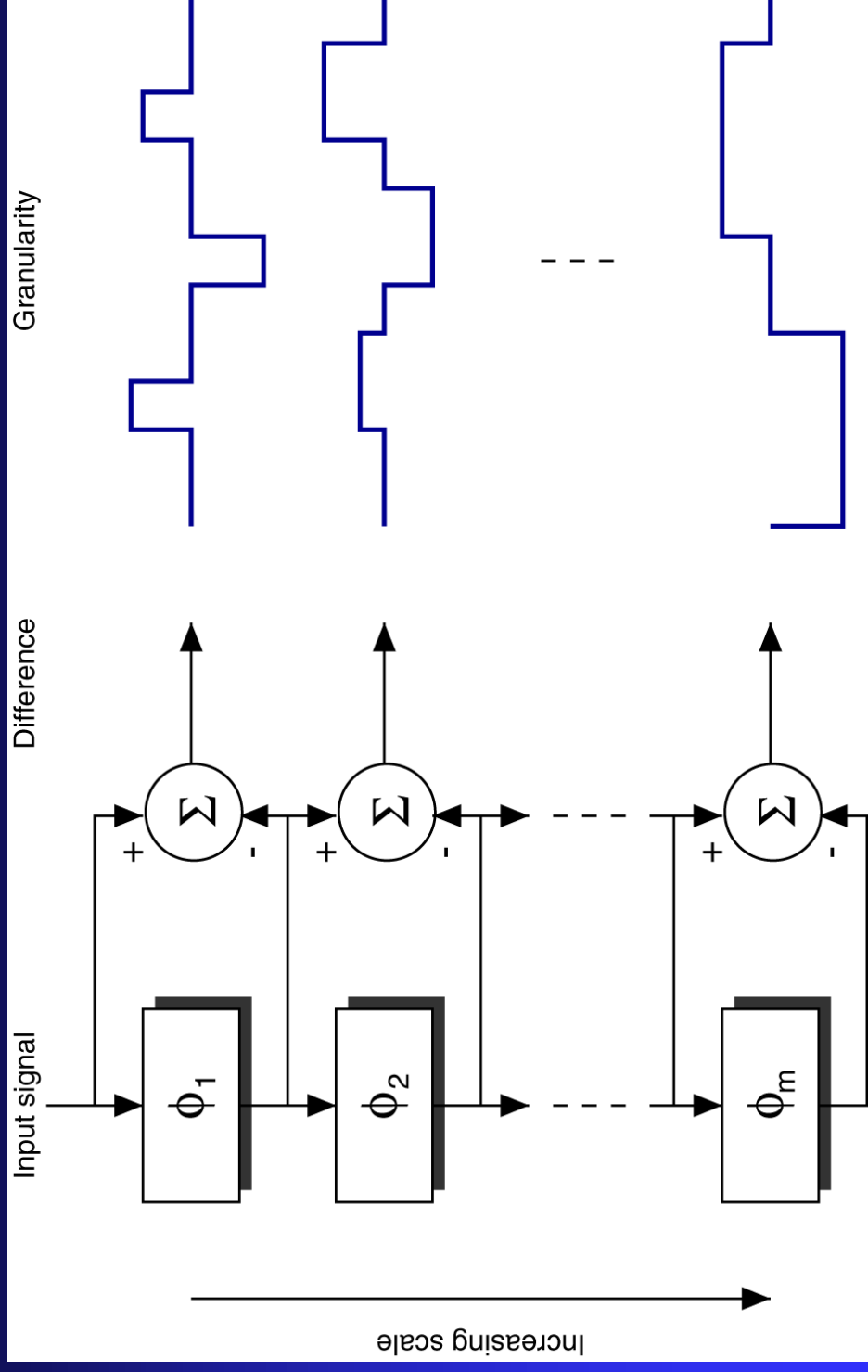


Figure 4.1: Sieve structure.

Sieve decomposition

- The size of the scaling element determines the scale
- At each scale, *granules* may or may not be extracted
- Scales are binned into *channels* to lower computational load

2D sieve decomposition

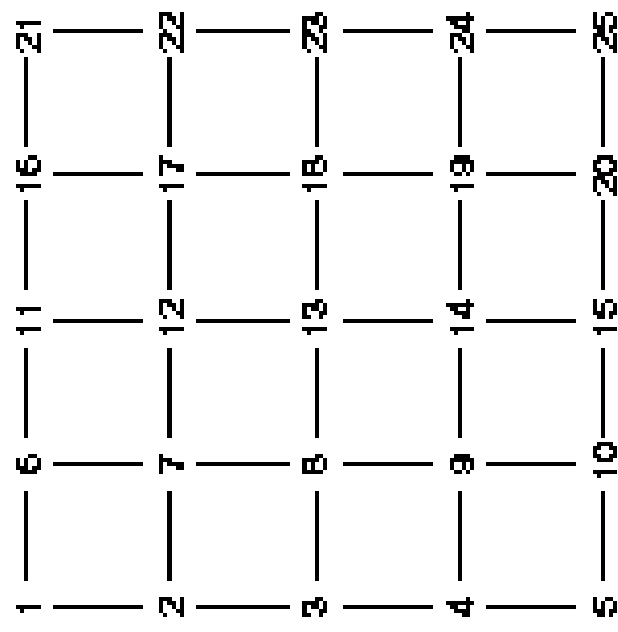
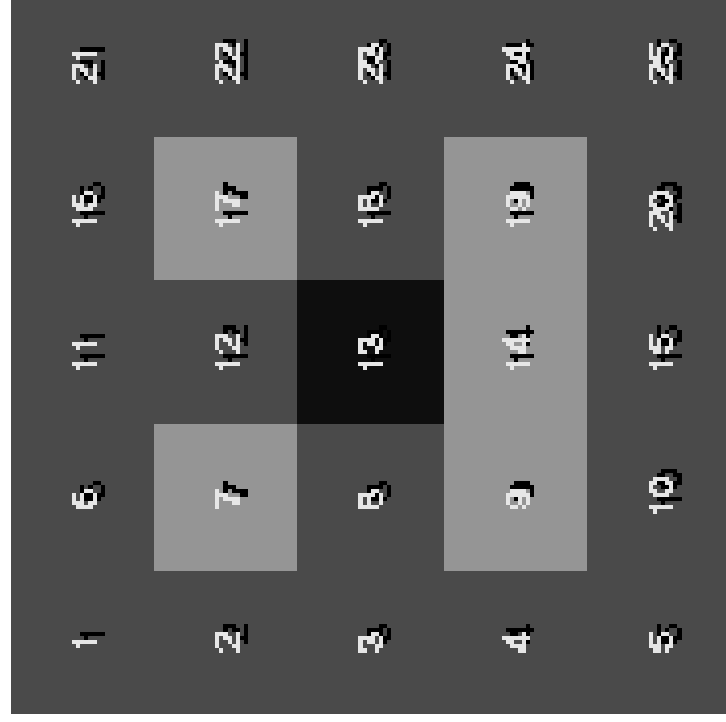
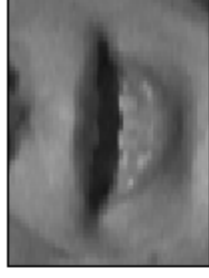


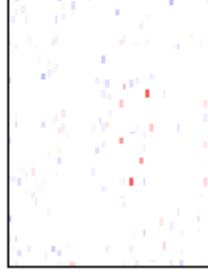
Figure 4.8: Image represented as a four-connected graph

Area tracking

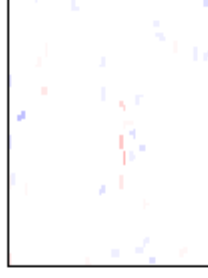
Example image



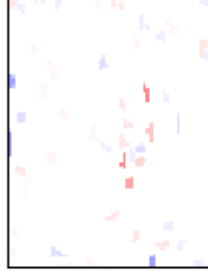
2D M-sieve: 1 to 2



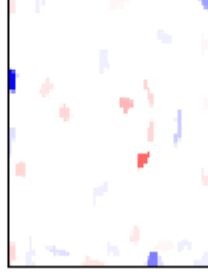
2D M-sieve: 3 to 4



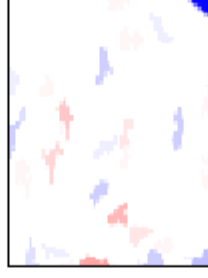
2D M-sieve: 5 to 8



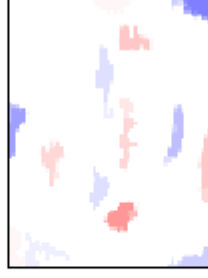
2D M-sieve: 9 to 16



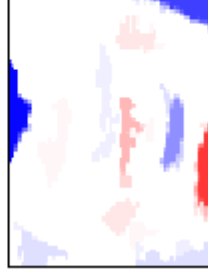
2D M-sieve: 17 to 32



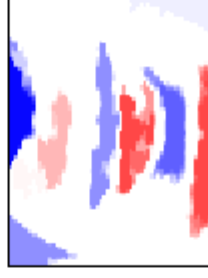
2D M-sieve: 33 to 64



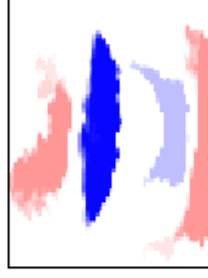
2D M-sieve: 65 to 128



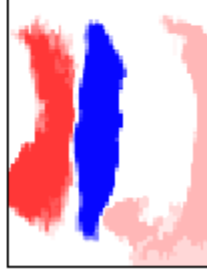
2D M-sieve: 129 to 256



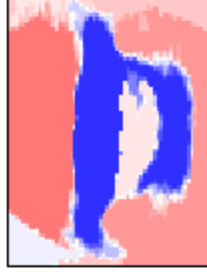
2D M-sieve: 257 to 512



2D M-sieve: 513 to 1024



2D M-sieve: 1025 to 2048



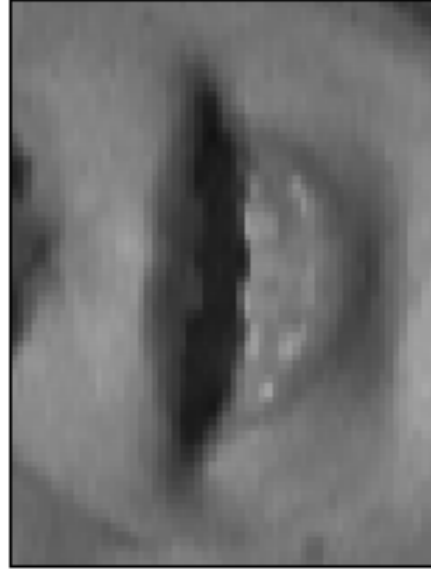
2D M-sieve: 2049 to 4096



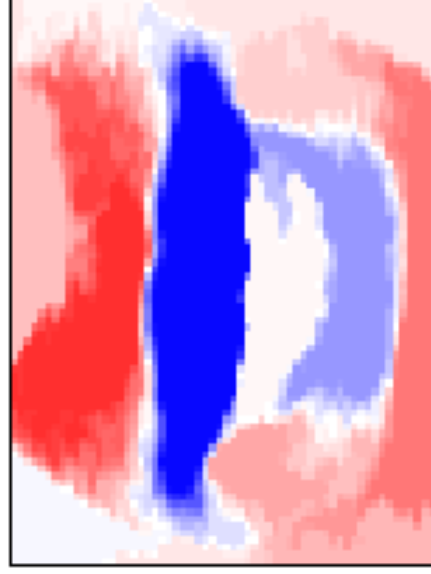
Area tracking

- Idea: Some scale/channel will find the area of the mouth opening
- Problem: what's the right scale?
- Solution: 300-2000

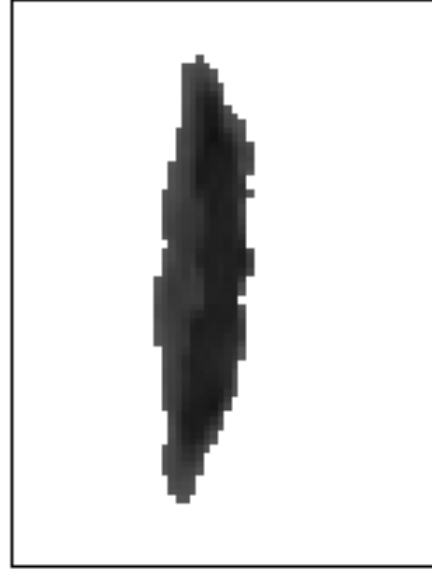
Area tracking



(a) image



(b) bandpass 300 to 2000



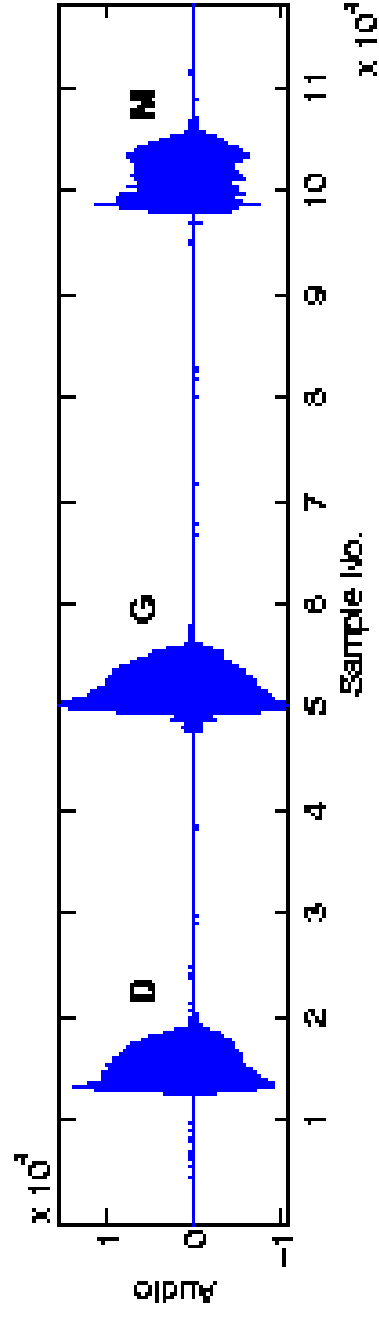
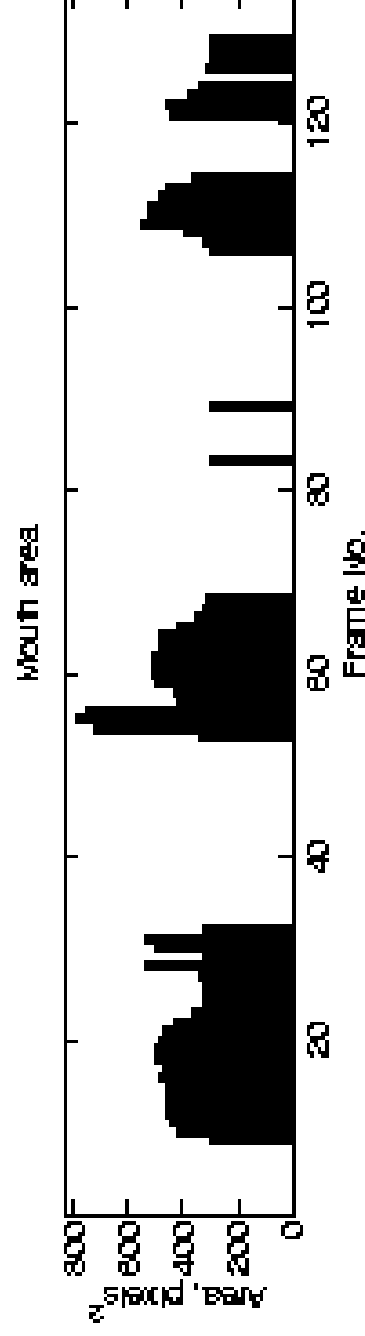
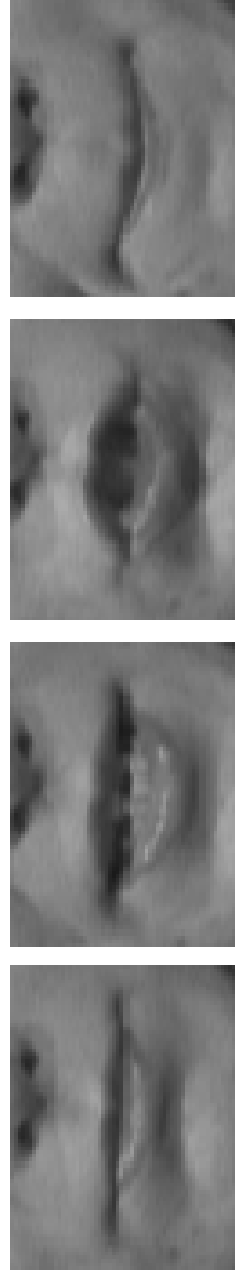
(d) masked image



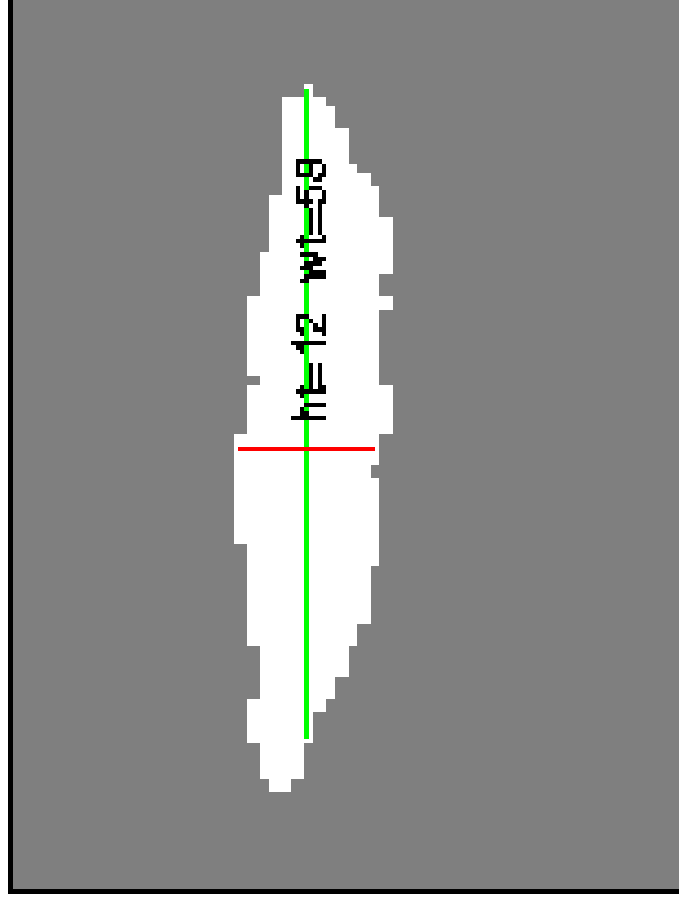
(c) ≤ -25 threshold mask

Area tracking

Example sequence 'D-G-M'



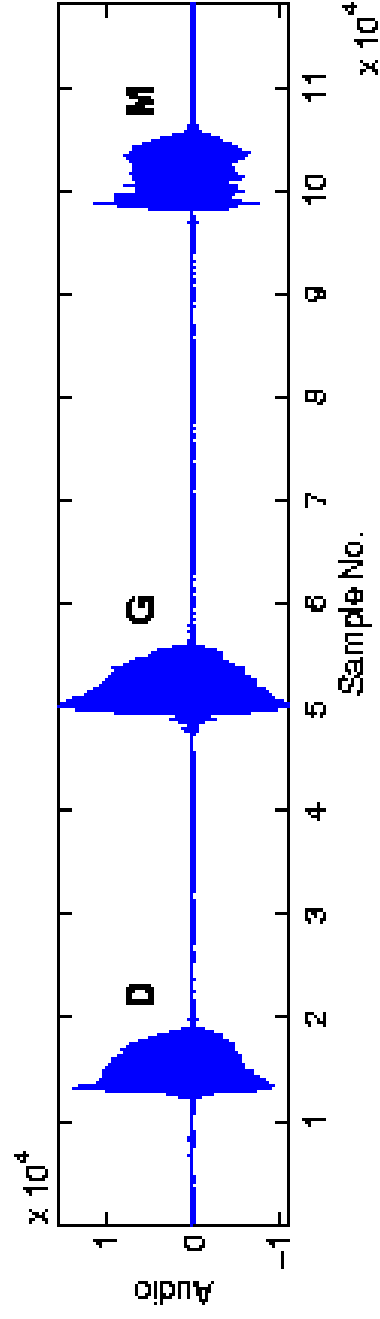
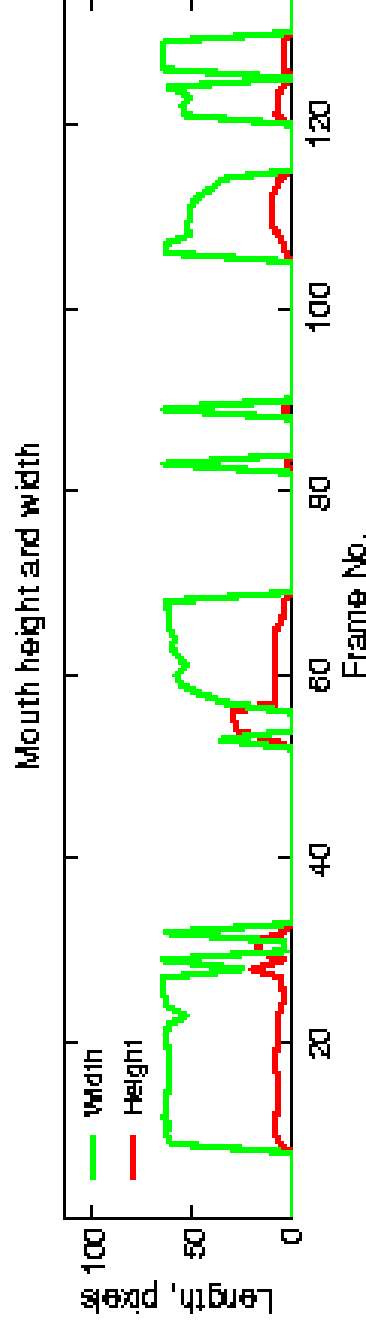
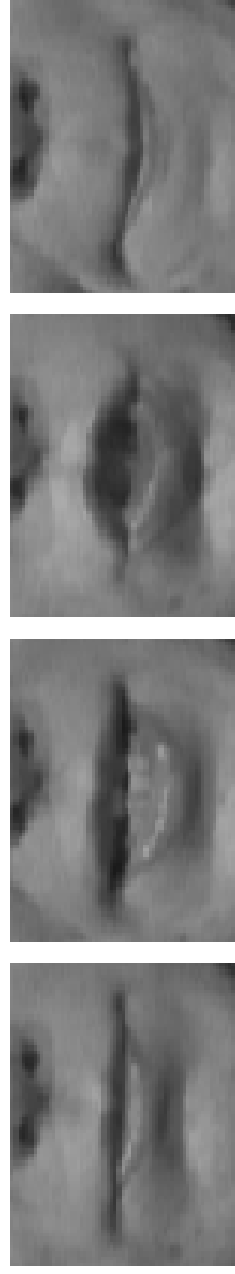
Height and width tracking



(c) measured height and width

Height and width tracking

Example sequence 'D-G-M'

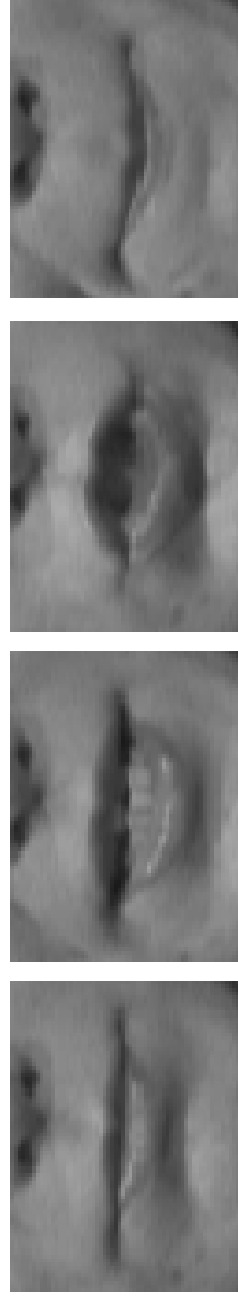


Area histogram

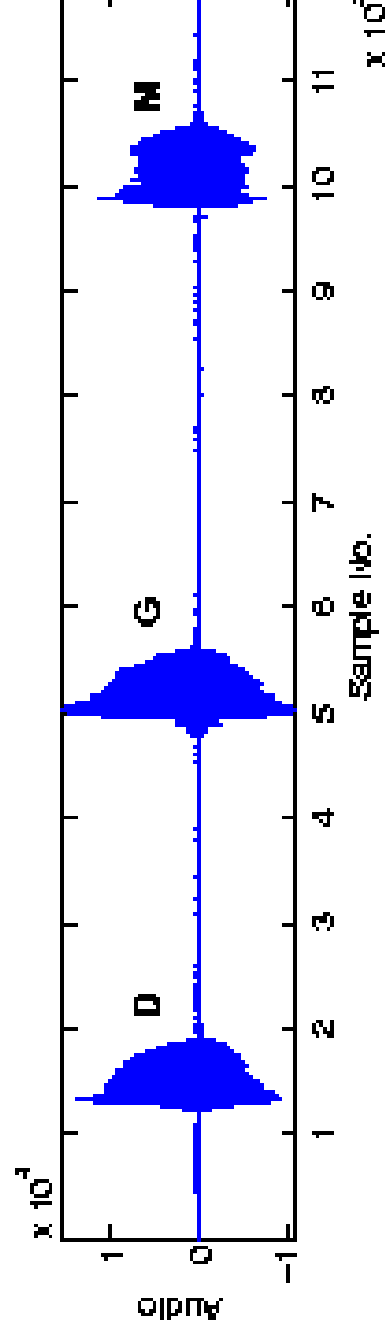
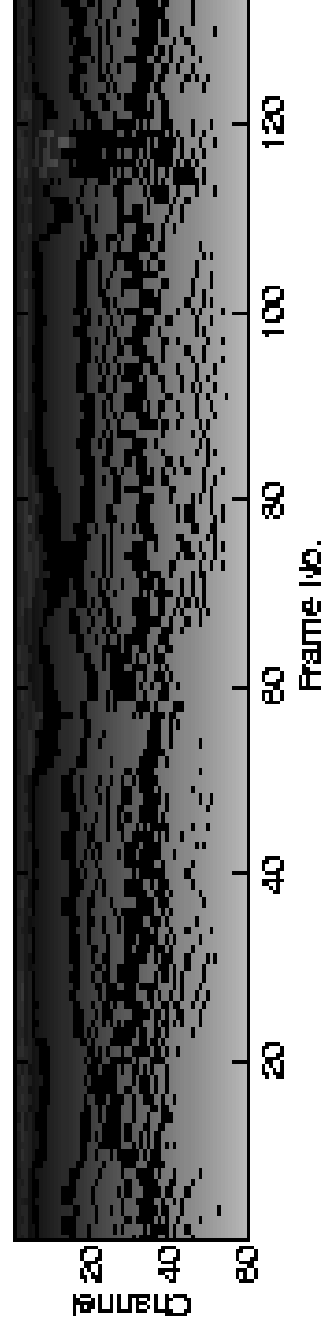
- To overcome the arbitrary choice of channel
- Choose all channels

Area histogram

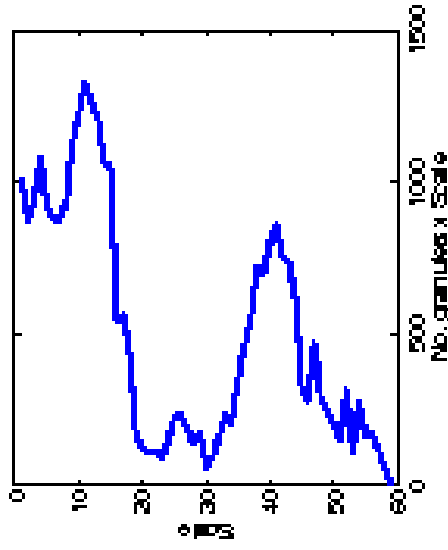
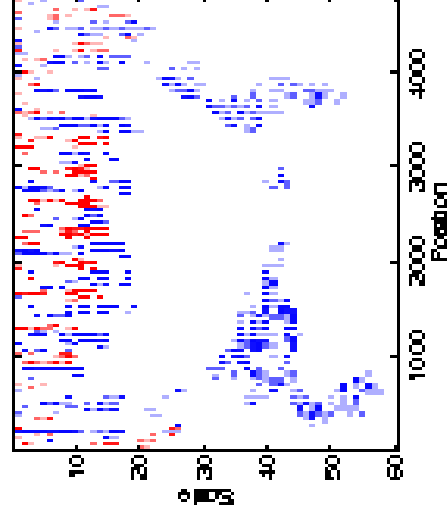
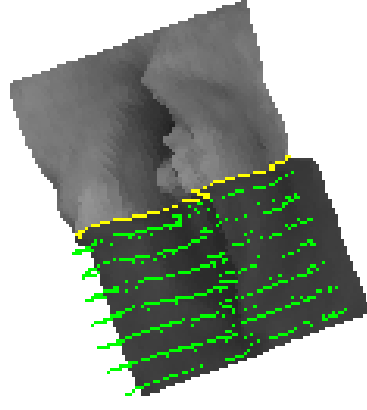
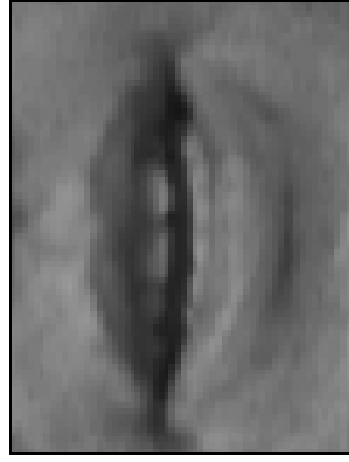
Example sequence 'D-G-M'



Area c-sleve histogram, 60 linear channels (linear scaled)



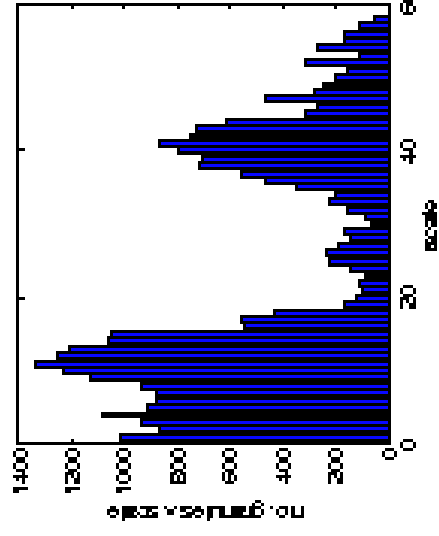
Scale histogram



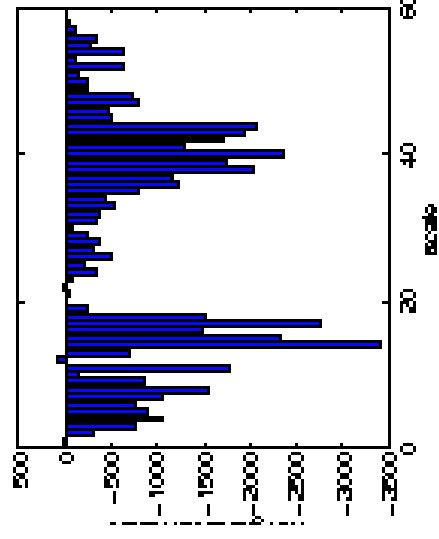
PCA

- To further reduce the number of channels
- Problem: Covariance or correlation?
- Problem: How many PCs?

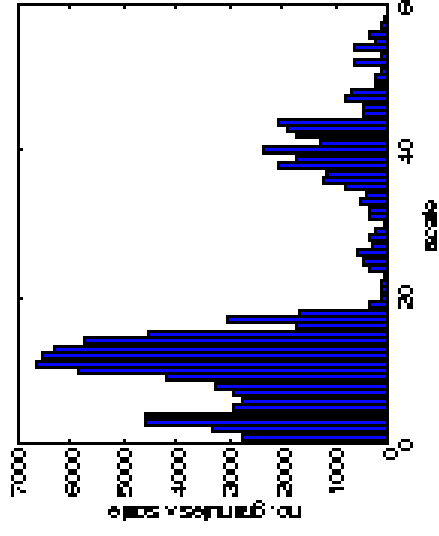
Scale histogram



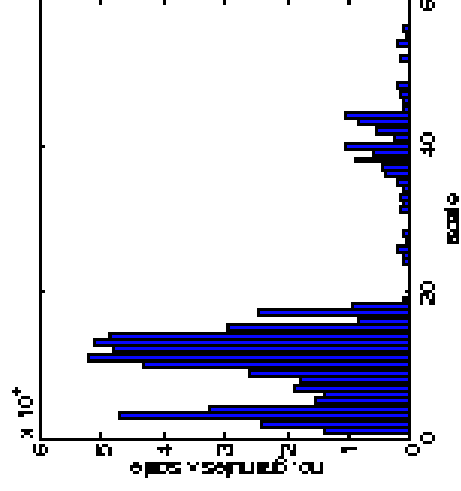
(c) nH , scale count



(d) α , amplitude sum



(e) $|a|$, magnitude sum



(f) β , amplitude² sum

The Full investigation

Attribute	Settings
Sieve type	median, m opening, o closing, c
DC Baseline	preserve ignore
Interpolate	N, 25Hz I, 50Hz
PCA components	10 20
PCA type	covariance correlation
HMM states	5 7 9
Gaussian modes	1 3

Key	Type	Details
sf	scale count	count of granules at each scale
a	amplitude sum	sum of amplitudes at each scale
$ a $	magnitude sum	sum of absolute amplitude at each scale
a^2	power sum	sum of squared amplitude at each scale

Table 7.8: Key to scale histogram types.

The Full investigation

- Yet another parameter: Scaling - Multiplication of count by scale
- 2304 combinations of 9 parameters
- Tested on two databases

Evaluation of low level methods

- Area tracking and height/width tracking gave less than 10% correct identification across talkers
- Area histogram was better, but still not good. Less than 20% correct identification across talkers
- All three methods not further discussed

Evaluation of low level methods

- Consistent results (two databases):
 - ◆ Covariance better than correlation for PCA
 - ◆ Opening performs consistently poorer than median or closing

Evaluation of low level methods

- Inconsistent results (two databases):
- Three Gaussian modes better than one for one training set – vice versa for another
- Preserving the DC is better for one set – poorer for another
- Linear scaling helps for one database – not for another. Conclusion: Better do it.
- 20 PCs better than 10 for one database – vice versa for another

Evaluation of low level methods

- Robustness of results:
 - ◆ Linear scaling affects correlation PCA indicating that the results are "statistically fragile"

High level methods

- Active shape model (ASP)
- Active appearance model (AAP)

Active shape model (ASP)

- Based on a Point Distribution Model
- Landmark and secondary points placed manually on training set
- Spline fitted to points
- Secondary points repositioned evenly along spline

Active shape model (ASP)

- Elimination of pose variation by minimization of

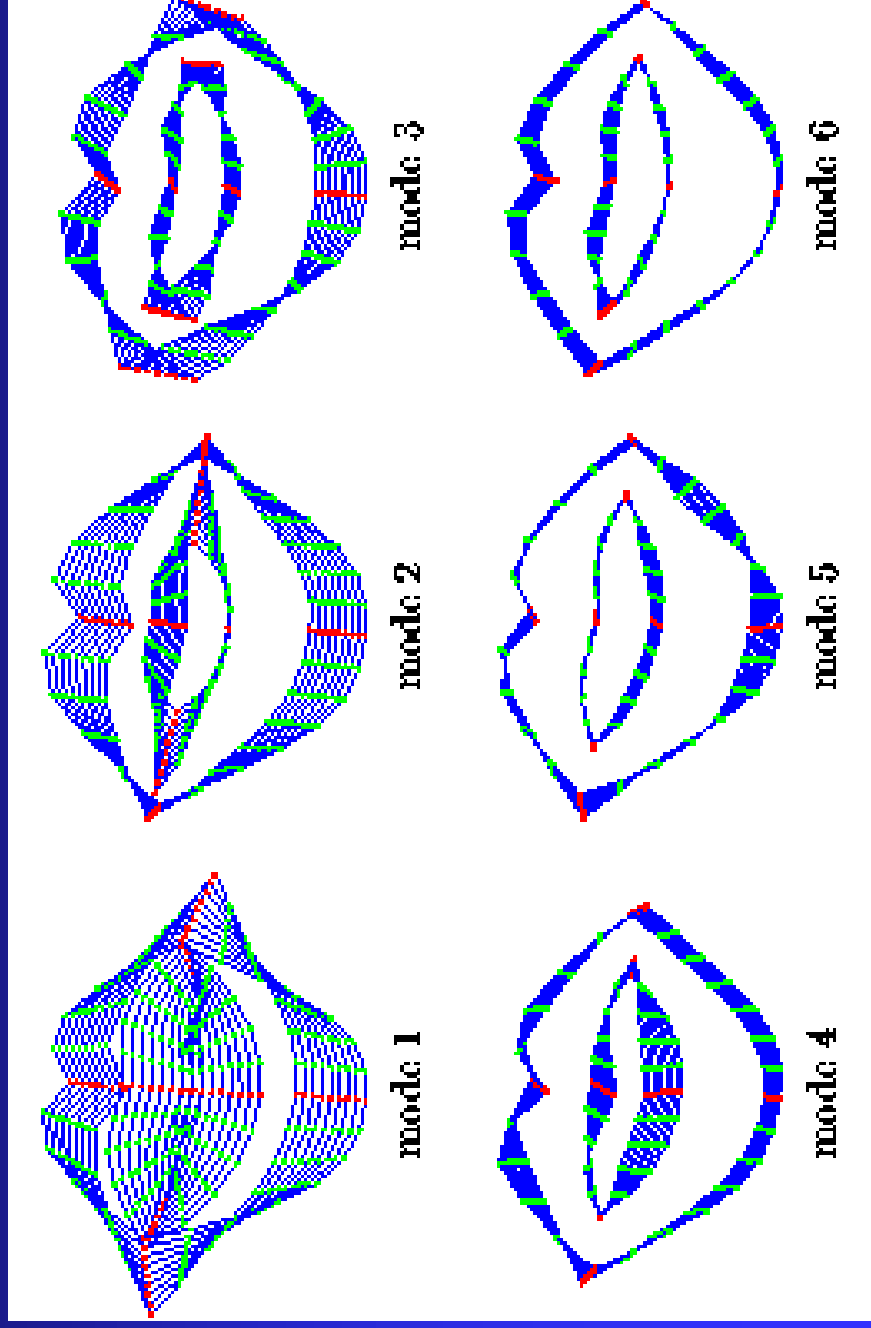
$$E = (\mathbf{x}_1 - M(a, \theta)[\mathbf{x}_2] - \mathbf{t})^T W (\mathbf{x}_1 - M(a, \theta)[\mathbf{x}_2] - \mathbf{t})$$

$$M(a, \theta) \begin{bmatrix} x_{jE} \\ y_{jE} \end{bmatrix} = \begin{pmatrix} (a \cos \theta) x_{jE} - (a \sin \theta) y_{jE} \\ (a \sin \theta) x_{jE} + (a \cos \theta) y_{jE} \end{pmatrix}$$
$$\mathbf{t} = (t_{x1}, t_{y1}, \dots, t_{xN}, t_{yN})$$

- Weights highest for stable points

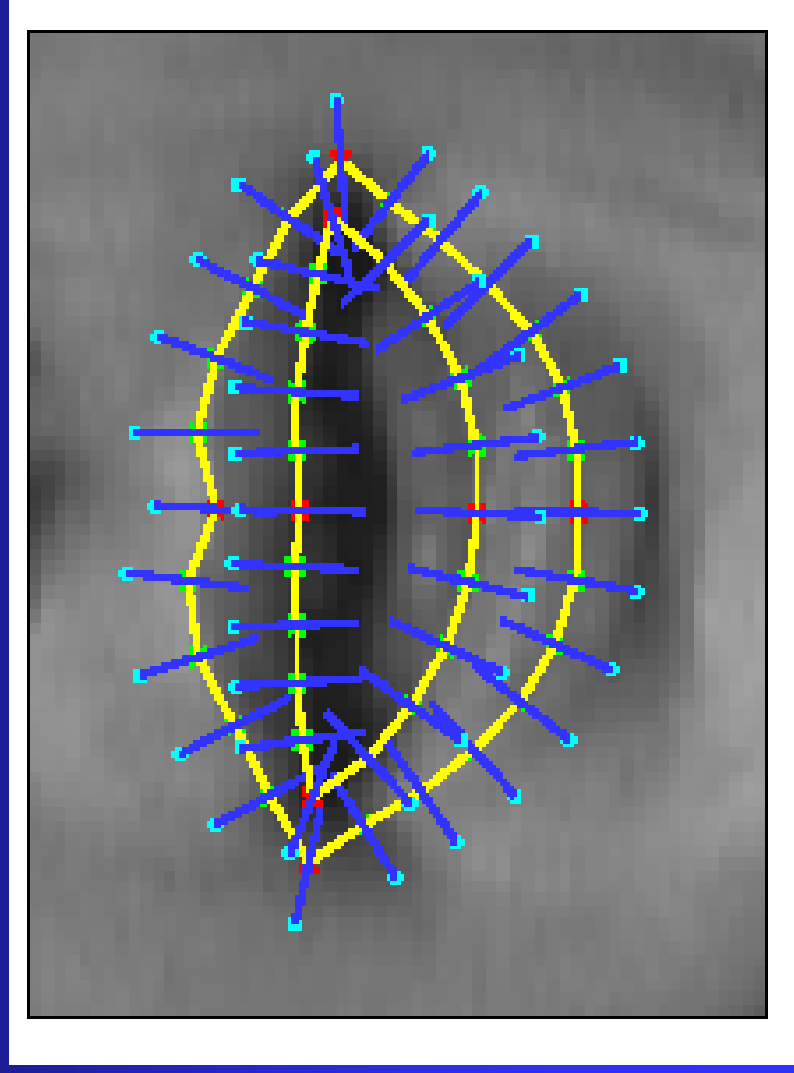
Active shape model (ASP)

- 95% PCA filtering

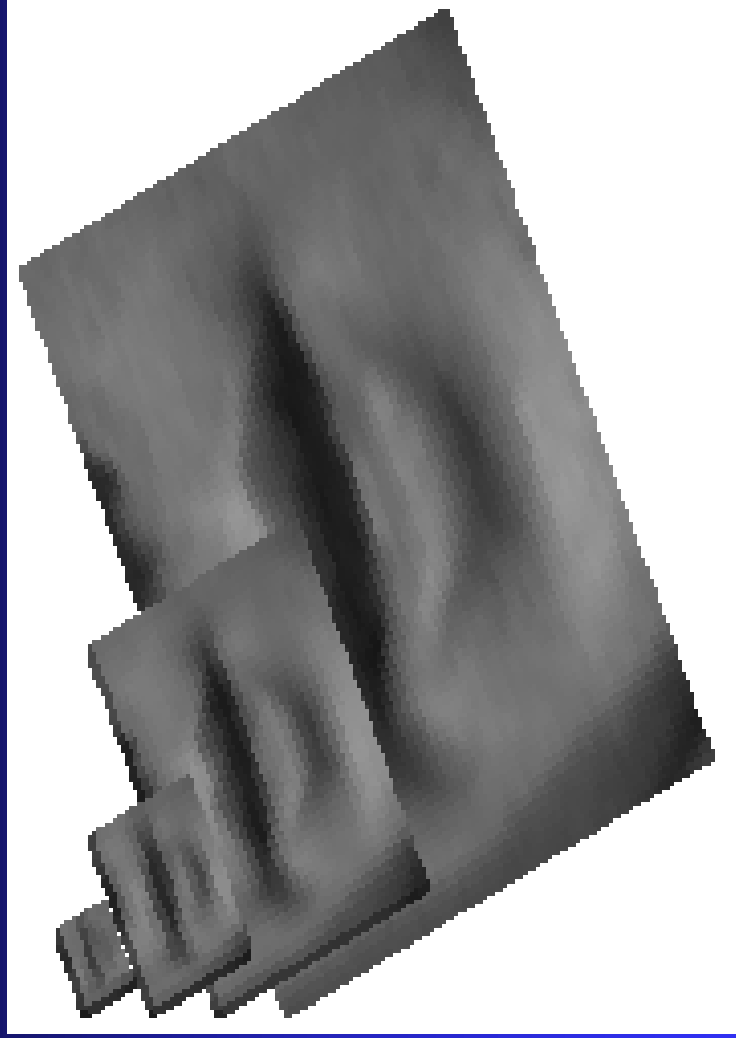


Active shape model (ASP)

- Grey Level Profile Distribution Model (GLPD)



Active shape model (ASP)



- Multi-resolution Fitting
- GLDM starts at coarse scale – moves to finer
- May yield better fit

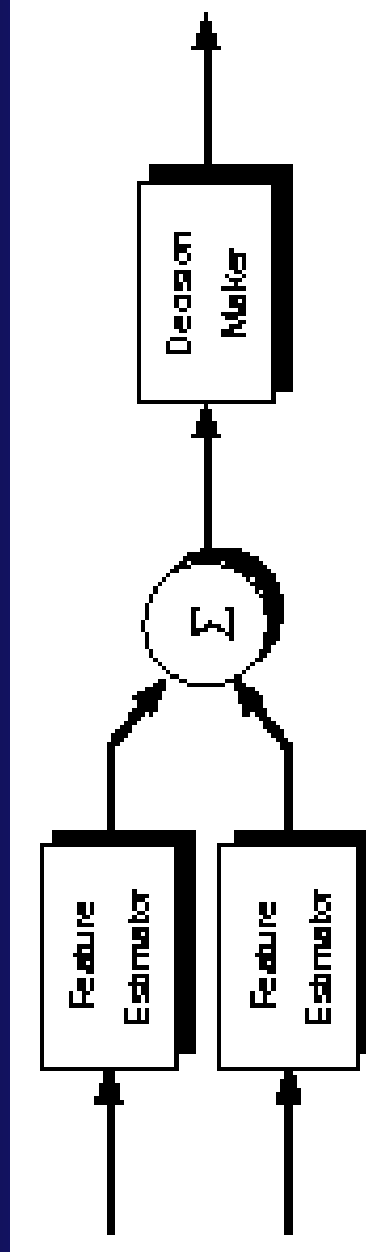
Active Appearance model (AMP)

- Combines the ASM with a greylevel appearance model (GAM)
- GAM:
 - ◆ Warp images to mean shape
 - ◆ Pixel values within mean shape forms the GAM

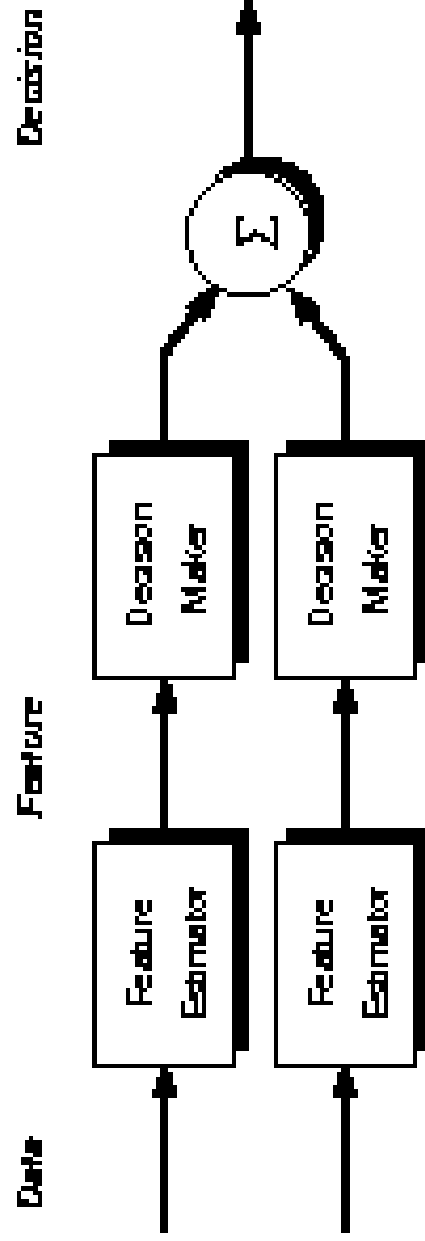
Evaluation of high level methods

- ASM poorer than MSA
- Multi-resolution better
- Longer Greylevel Profiles better
- Initialising optimisation at mean shape better than at previous shape – difficult to break out of local minima

Audiovisual integration



(b) Feature fusion (early integration).



(c) Decision fusion (late integration).

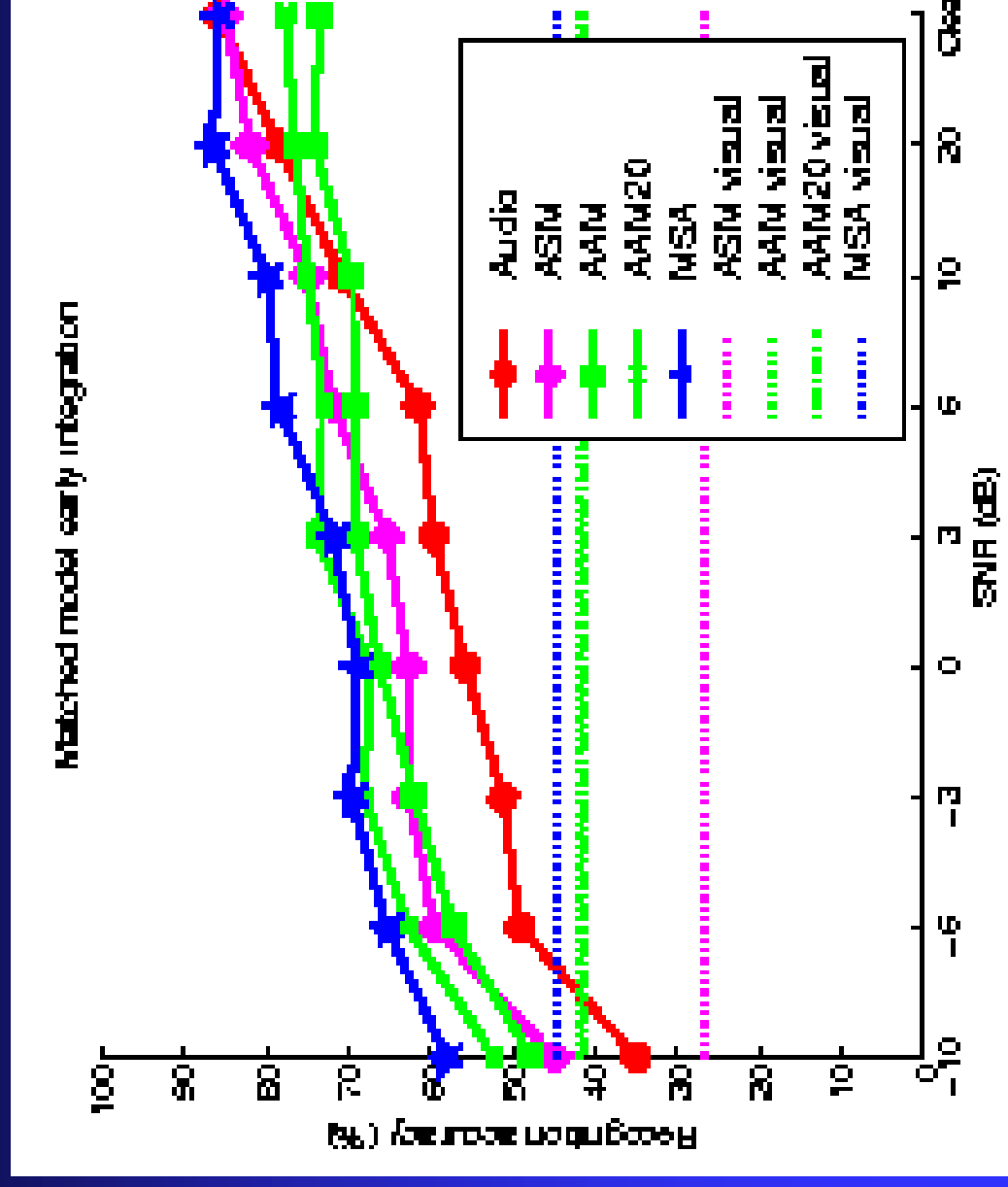
Early audiovisual integration

- SNR variation
- Matched condition
 - ◆ Trained at each SNR – tested at each SNR
 - ◆ Best possible performance
- Unmatched condition
 - ◆ Trained at highest SNR and tested at all SNRs

Early audiovisual integration

- Visual methods tested: MSA, ASM & AAM
- Measure of success: improvement in performance *compared to audio only*
- Conclusion
 - ◆ AAM best with degraded audio, but *lower* than audio only at highest SNR
 - ◆ MSA and ASM gains only (and very little) with more HMM Gaussian modes than possible with audio only (due to less training data).
- Overall: Poor performance

Early audiovisual integration

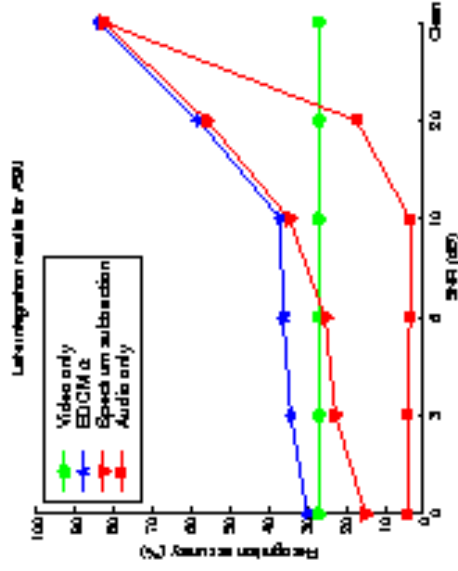


Late audiovisual integration

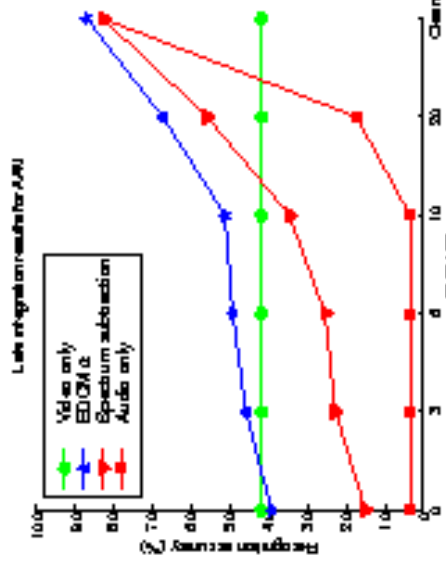
■ Integration rule

$$w_s^* = \underset{s: 1, 2, \dots, N}{\text{arg max}} \{ \alpha \mathcal{L}(w_s | A) + (1 - \alpha) \mathcal{L}(w_s | V) \}$$

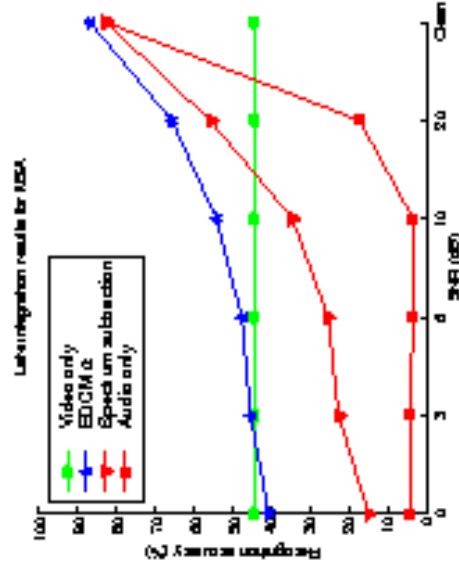
Late audiovisual integration



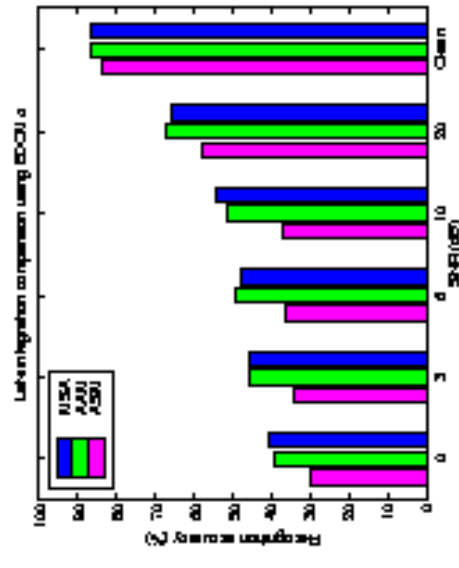
(a) ASM visual features



(b) AAM visual features



(c) MSA visual features



(d) EDCM integration comparison

Homework

- Derive Bayes' rule for two conditions
 - ◆ $P(\text{percept} | \text{Audio, Visual}) = ?$
- Compare Matthew's late integration rule with that of Fritsch et al.'s (Katagiri, ch. 7.4)
 - ◆ Is any (or both) statistically optimal (Bayesian)? Justify.
 - ◆ How does each respond to incongruent (wrongly dubbed) audiovisual speech?
- The McGurk effect was demonstrated during this talk. It is an illusion occurring from incongruent audiovisual speech. It may seem mysterious why the brain would allow such an illusion. Give a reason for why this might not be a mystery.