# Ecological Language Acquisition via Incremental Model-Based Clustering

## Giampiero Salvi

KTH CSC TMH giampi@kth.se

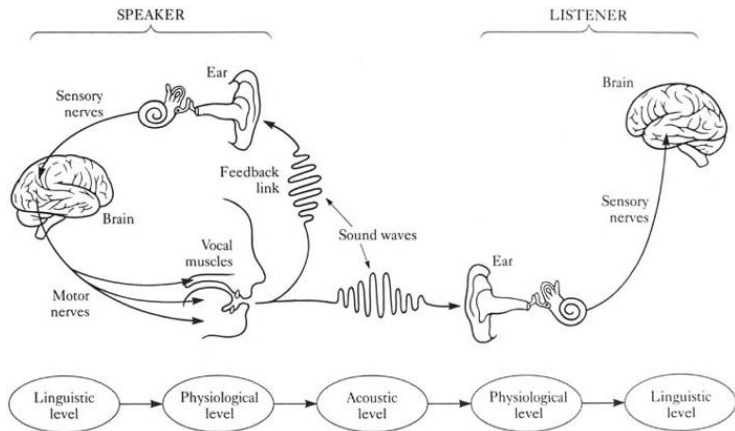### Nov. 2005

# The Speech Chain



Denes and Pinson (1993)
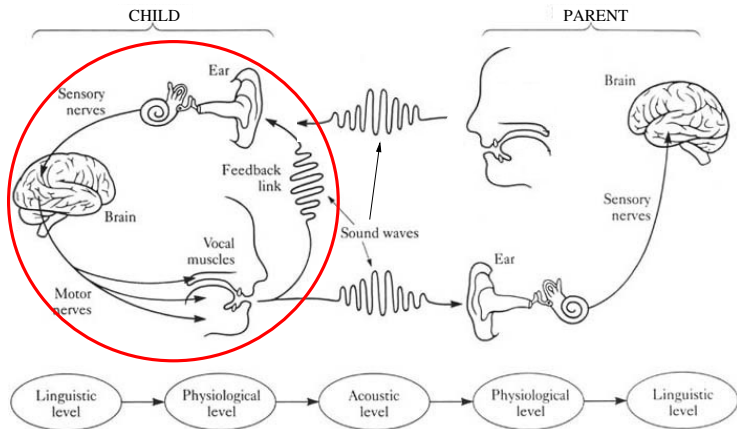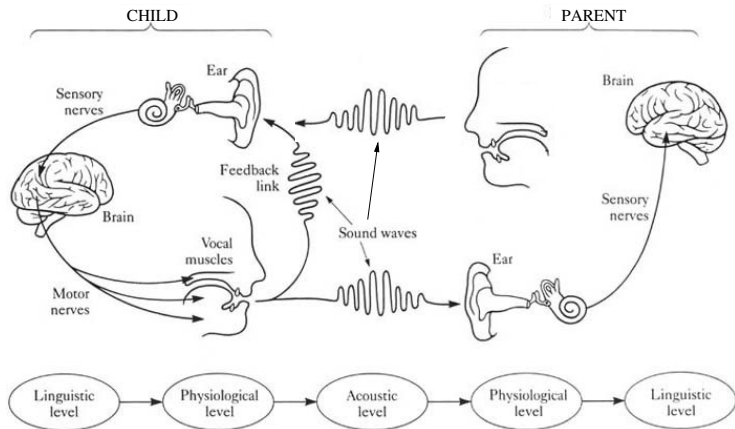
# The Speech Chain



Denes and Pinson (1993)

# The Speech Chain



Denes and Pinson (1993)

# The Speech Chain



Denes and Pinson (1993)

- Background: ecological theory of language acquisition (Lacerda et al., 2004)
  - the infant is naïve: no innate linguistic knowledge

- Background: ecological theory of language acquisition (Lacerda et al., 2004)
  - the infant is naïve: no innate linguistic knowledge
- Aim (long term): mathematical modelling of the learning process
  - acoustic features classification
  - time integration into meaningful sequences
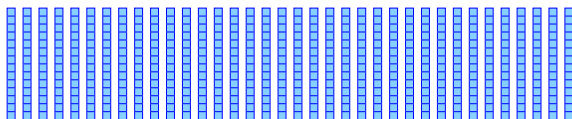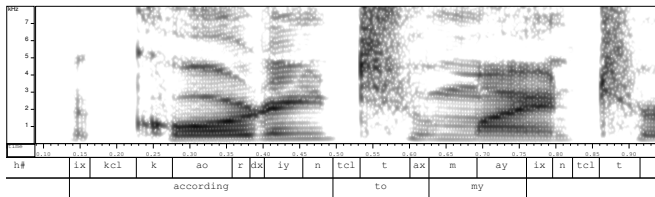  - integration of acoustic/visual information

- ▶ Background: ecological theory of language acquisition (Lacerda et al., 2004)
  - ▶ the infant is naïve: no innate linguistic knowledge
- ▶ Aim (long term): mathematical modelling of the learning process
  - ▶ acoustic features classification
  - ▶ time integration into meaningful sequences
  - ▶ integration of acoustic/visual information
- ▶ Aim Interspeech 2005 (Salvi, 2005): acoustic features classification
  - ▶ unsupervised
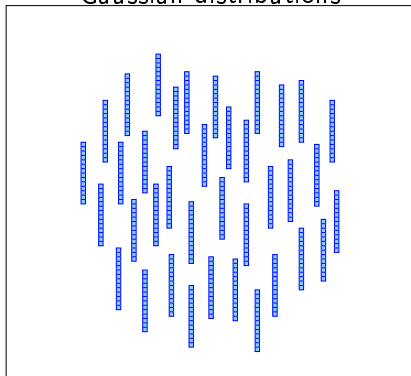  - ▶ incremental

# Acoustic features

Equally spaced windows of speech

# Assumption

Acoustic feature vectors independently drawn from mixture of Gaussian distributions

# Method

- Model-Based Clustering (Fraley and Raftery, 1998)
  - data modelled as mixture of probability distributions
  - each distribution represents a cluster
  - each data point belongs to each cluster with a certain probability
  - model parameters estimated via Expectation Maximisation
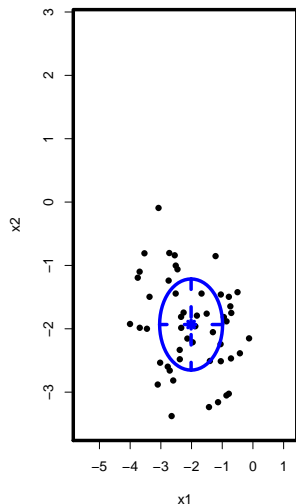  - different models compared via Bayes information criterion (BIC)

# Method

- Model-Based Clustering (Fraley and Raftery, 1998)
  - data modelled as mixture of probability distributions
  - each distribution represents a cluster
  - each data point belongs to each cluster with a certain probability
  - model parameters estimated via Expectation Maximisation
  - different models compared via Bayes information criterion (BIC)
- Incremental Model-Based Clustering (Fraley et al., 2003)
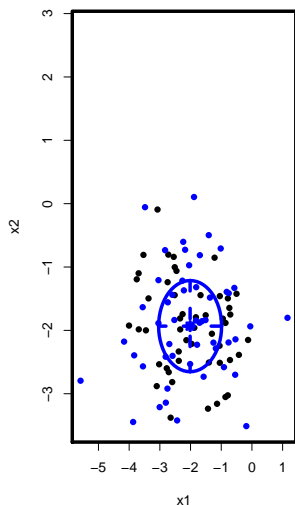  - introduced for large datasets

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
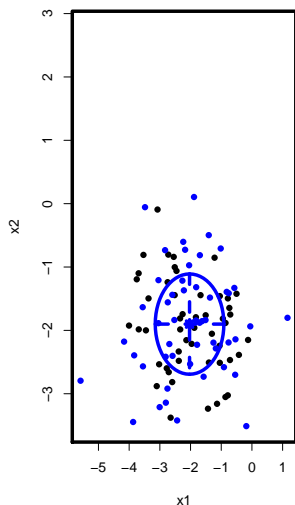6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
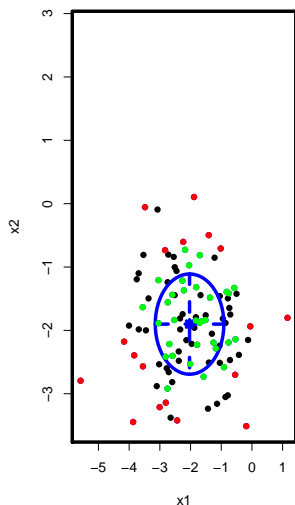6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
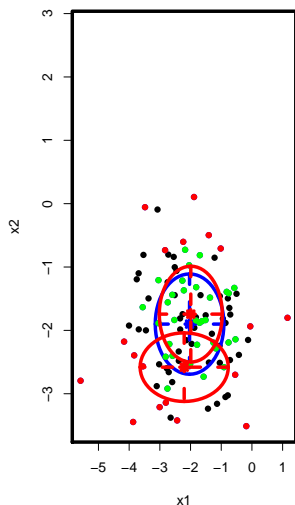6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
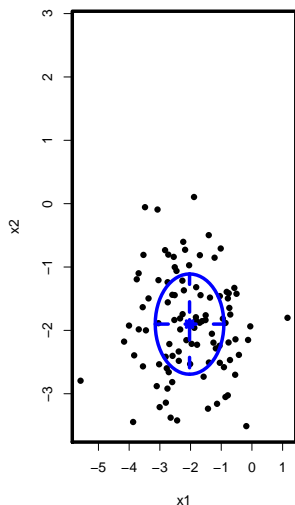6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into <span style="color:green">well</span> and <span style="color:red">poorly</span> modelled points
5. <span style="color:red">try a more complex model, if better BIC set as best and go back to 4</span>
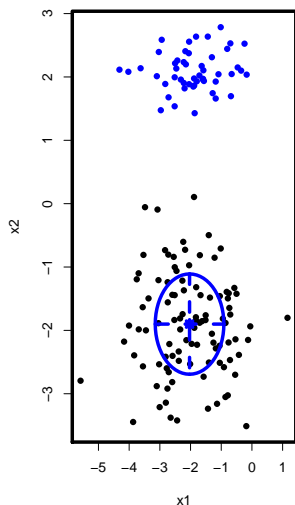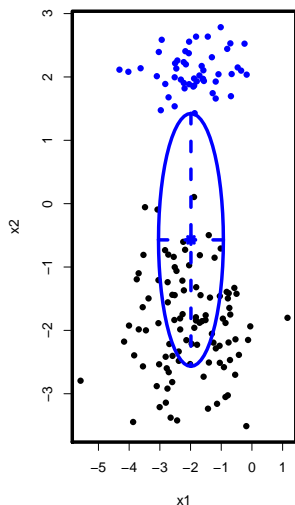6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
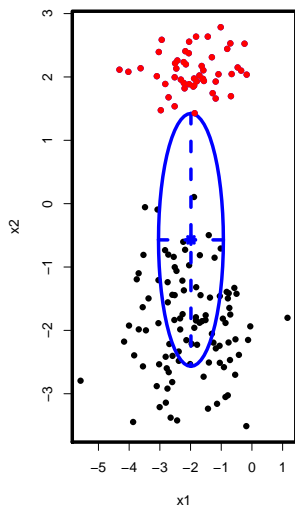6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
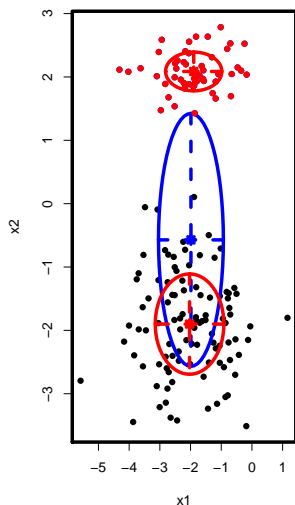6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
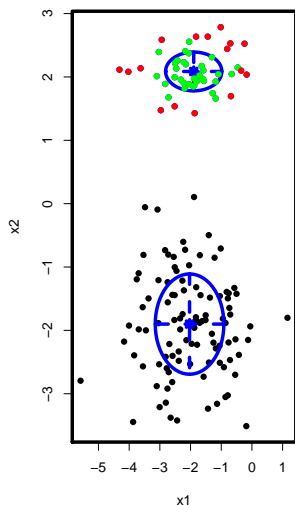6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
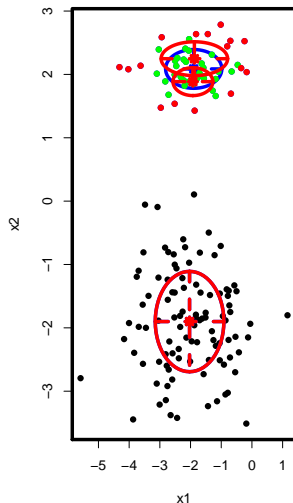6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
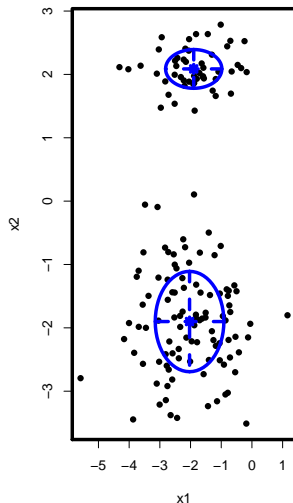6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
6. set the current best model and go back to 2

# Algorithm

1. start with a MCLUST model
2. get new data
3. adjust old model to new data
4. divide new data into well and poorly modelled points
5. try a more complex model, if better BIC set as best and go back to 4
6. set the current best model and go back to 2

# Experimental settings

- Data (ex1, ex2, ex3, ex4, ex5)
  - 12 minutes from the MILLE corpus
  - child directed speech (1 mother talking to her child)
  - Mel frequency cepstral coeffs computed every 10ms + differences of first and second order

# Experimental settings

- Data (ex1, ex2, ex3, ex4, ex5)
  - 12 minutes from the MILLE corpus
  - child directed speech (1 mother talking to her child)
  - Mel frequency cepstral coeffs computed every 10ms + differences of first and second order
- experimental factors
  - dimensionality of the data: from 3 to 39 dimensions
  - frame length: from 200msec to 3sec

# Evaluation

- problem: there is no reference (at the moment)

# Evaluation

- problem: there is no reference (at the moment)
- relative evaluation:

# Evaluation

- problem: there is no reference (at the moment)
- relative evaluation:
- time evolution of number of clusters
  - dependency with number of feature coefficients
  - dependency with frame length

# Evaluation

- problem: there is no reference (at the moment)
- relative evaluation:
- time evolution of number of clusters
  - dependency with number of feature coefficients
  - dependency with frame length
- agreement of classification in different conditions
  - variation of information (Meilă, 2002)

$$\mathbf{V}I(C, C') = H(C|C') + H(C'|C)$$

# Results



**effect of dimensionality (frame len = 50)**

Legend:
- nvars = 39
- nvars = 24
- nvars = 12
- nvars = 6
- nvars = 3

y-axis: # clusters
x-axis: time in s (1 s = 100 samples)

# Example

# Mismatch Child/Parent Voice
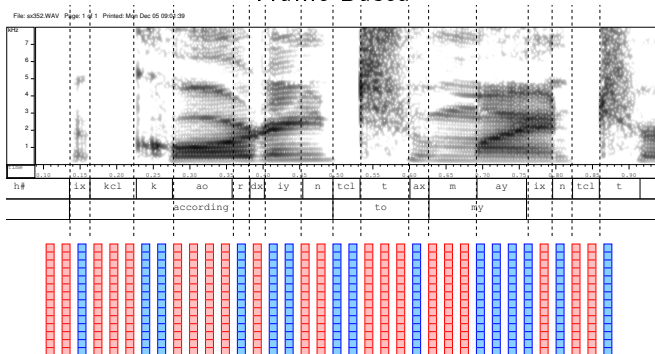
# Mismatch Child/Parent Voice



- ASR with children

# Mismatch Child/Parent Voice



- ASR with children
- Normalisation
  - VTLN: Vocal Tract Length Normalisation
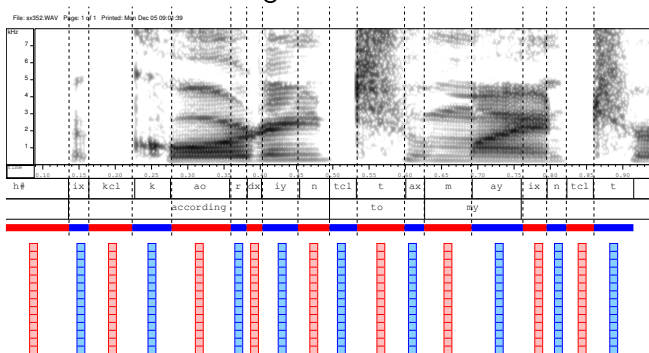  - Adaptation: hard in this context

# Mismatch Child/Parent Voice



- ASR with children
- Normalisation
    - VTLN: Vocal Tract Length Normalisation
    - Adaptation: hard in this context
- Relative Features
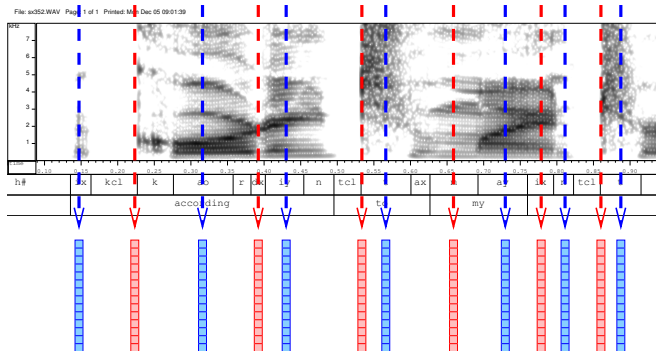
# Acoustic Features
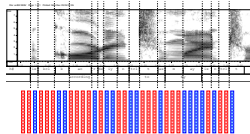
## Frame Based

# Acoustic Features

## Segment Based
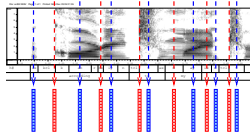
# Acoustic Features

## Landmark Based

# Consequences



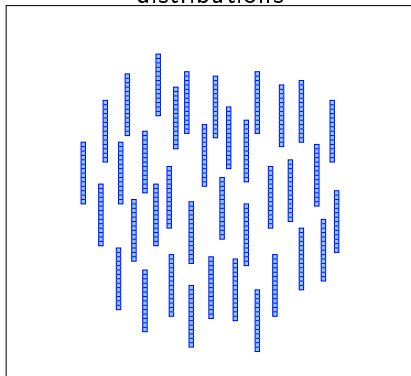$\Rightarrow$ **Sequence recognition (HMMs)**



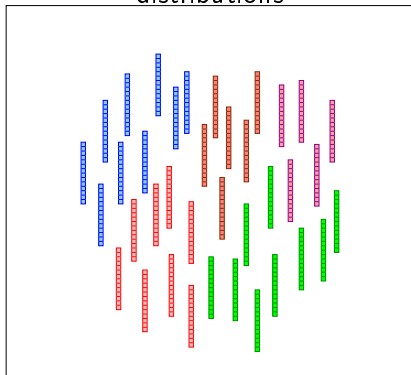$\Rightarrow$ **simpler relation acoustic categories/ linguistic units**
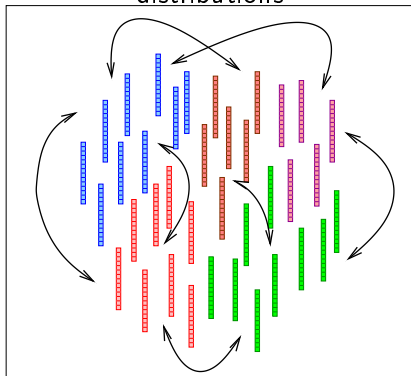
# Clustering Time Sequences

Acoustic vectors independently drawn from mixture of gaussian distributions

# Clustering Time Sequences

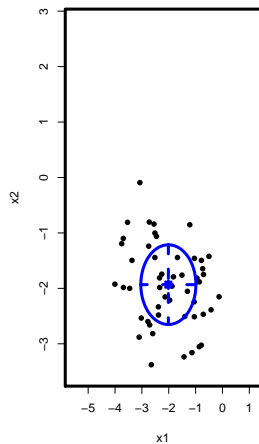Acoustic vectors independently drawn from mixture of gaussian distributions
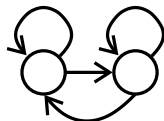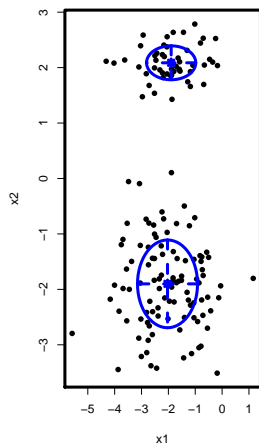
# Clustering Time Sequences

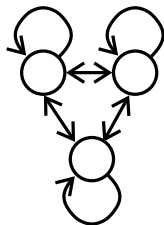Acoustic vectors independently drawn from mixture of gaussian distributions

# Modeling time evolution with Markov chains

# Modeling time evolution with Markov chains

# Modeling time evolution with Markov chains

# Modeling time evolution with Markov chains
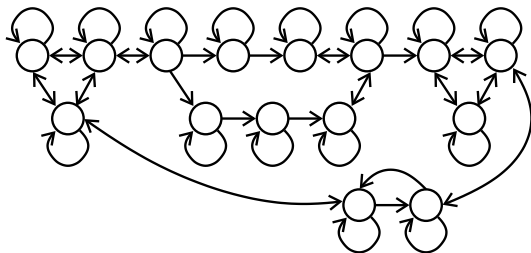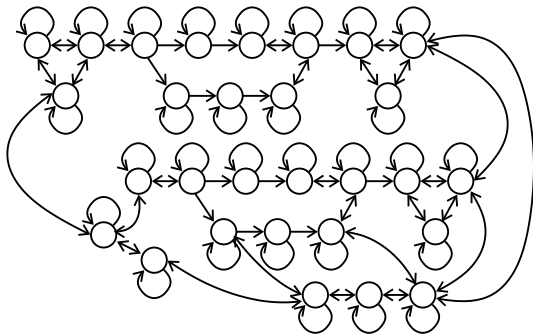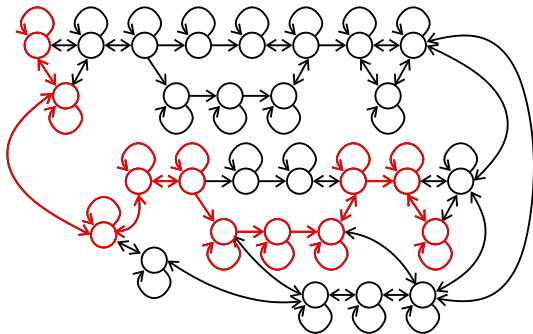
# Modeling time evolution with Markov chains
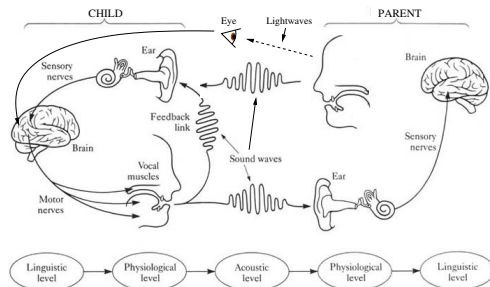
# Modeling time evolution with Markov chains

# The Visual Channel



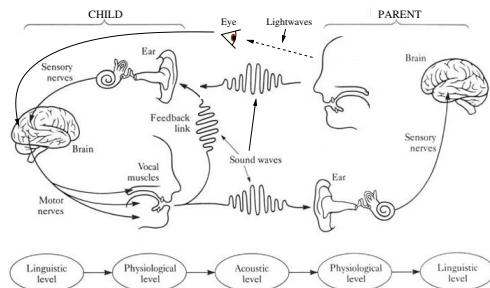- No one-to-one relation acoustic/visual info

# The Visual Channel



- No one-to-one relation acoustic/visual info
- Reinforcement Learning
  - perform match at higher levels (pseudo-words or -phrases)

# The Visual Channel

Perform visual/acoustic match on the Markov chain

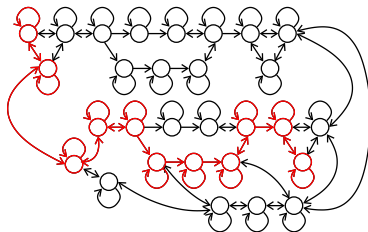Acoustic Event

Visual Event

# The Visual Channel

Perform visual/acoustic match on the Markov chain
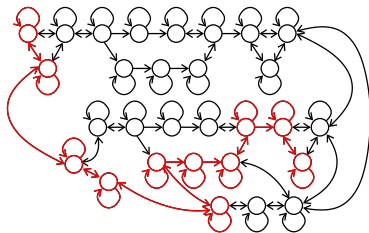
Acoustic Event

Visual Event

# The Visual Channel

Perform visual/acoustic match on the Markov chain

Acoustic Event

Visual Event

# The Visual Channel

Perform visual/acoustic match on the Markov chain

Acoustic Event

Visual Event

# The Final Question

- ► Are the acoustic blocks (categories) in a language learned out of their statistical occurrence or out of their contrastive use?

# The Final Question

- Are the acoustic blocks (categories) in a language learned out of their statistical occurrence or out of their contrastive use?

- in the first case: model based clustering and growing Markov chains are separate processes.

# The Final Question

- Are the acoustic blocks (categories) in a language learned out of their statistical occurrence or out of their contrastive use?

- in the first case: model based clustering and growing Markov chains are separate processes.

- in the second case: need to integrate everything

# Bibliography

http://www.speech.kth.se/~giampi

Denes, P. B. and Pinson, E. N. (1993). *The Speech Chain: Physics and Biology of Spoken Language*. W. H. Freeman.

Fraley, C., Raftery, A., and Wehrensy, R. (2003). Incremental model-based clustering for large datasets with small clusters. Technical Report 439, Department of Statistics, University of Washington.

Fraley, C. and Raftery, A. E. (1998). How many clusters? which clustering method? answers via model-based cluster analysis. *Computer Journal*, 41(8).

Lacerda, F., Klintfors, E., Gustavsson, L., Lagerkvist, L., Marklund, E., and Sundberg, U. (2004). Ecological theory of language acquisition. In *EPIROB*, pages 147–148.

Meilă, M. (2002). Comparing clusterings. Technical Report 418, Department of Statistics, University of Washington.

Salvi, G. (2005). Ecological language acquisition via incremental model-based clustering. In *Interspeech*, pages 1181–1184.