

Linear Gaussian models for speech recognition

Antti-Veikko Rosti

November 2002



Cambridge University Engineering Department

SVR Speech Seminar

Overview

- State-space models
- Linear Gaussian models
- Bayesian networks
- Factor analysed HMMs
- Segment models
- Conclusions

State-space models

Generative model of a state-space model:

$$\begin{aligned} \mathbf{o}_t &= g(\mathbf{x}_t, \mathbf{v}_t) \\ \mathbf{x}_{t+1} &= f(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{w}_t) \end{aligned}$$

- k -dimensional state vectors \mathbf{x}_t ;

- d -dimensional observation vectors \mathbf{o}_t ;

- state evolution noise \mathbf{w}_t ;

- observation noise \mathbf{v}_t .

Model for speech production

- articulator positions in state vectors;

- acoustic realisations in observation vectors;

- mappings $f(\cdot)$ and $g(\cdot)$ non-linear!

What about the mappings...

- state vector components are covariance model factors;
- observations are linear combinations of state vectors and noise.

Interpretation:

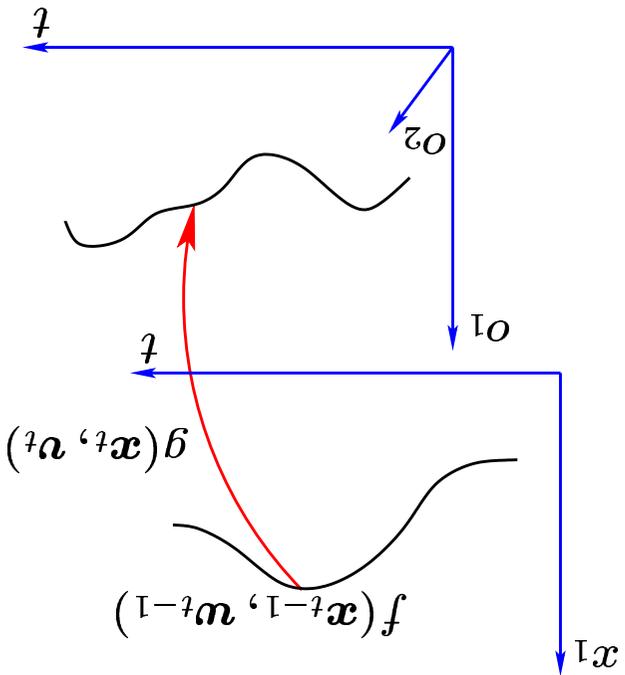
- not a good model for speech production!

$$x_{t+1} = f(x_t, w_t);$$

- usually Markov in state-space; i.e.,
- Gaussian distributed noise w_t and v_t ;
- linear mappings $f(\cdot)$ and $g(\cdot)$;

State-space models with

Linear Gaussian models



Observation process

Factor analysis (FA)

$$\begin{aligned}
 o_t &= C^t x_t + v_t \\
 v_t &\sim \mathcal{N}(m_{(o)}^t, \Sigma_{(o)}^t)
 \end{aligned}$$

- x_t are $\mathcal{N}(\mathbf{0}, I)$ distributed factors;

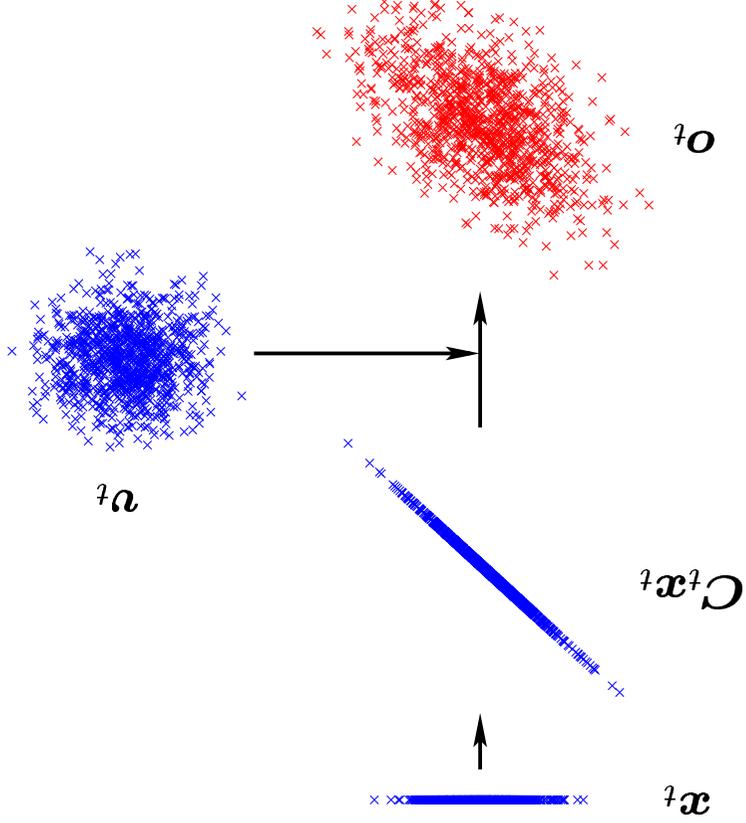
- C^t are matrices of factor loadings;

- v_t are errors with diagonal $\Sigma_{(o)}^t$;

- observation covariance $C^t C^{t'} + \Sigma_{(o)}^t$;

- fewer parameters than full covariance if $k > (p - 1)/2$.

Extensions: SFA, IFA, ...



Linear discriminant analysis

Generative model:

$$o_t = C_t \begin{bmatrix} x_t \\ v_t \end{bmatrix} \sim \mathcal{N}(o_t, \Sigma_{(o)})$$

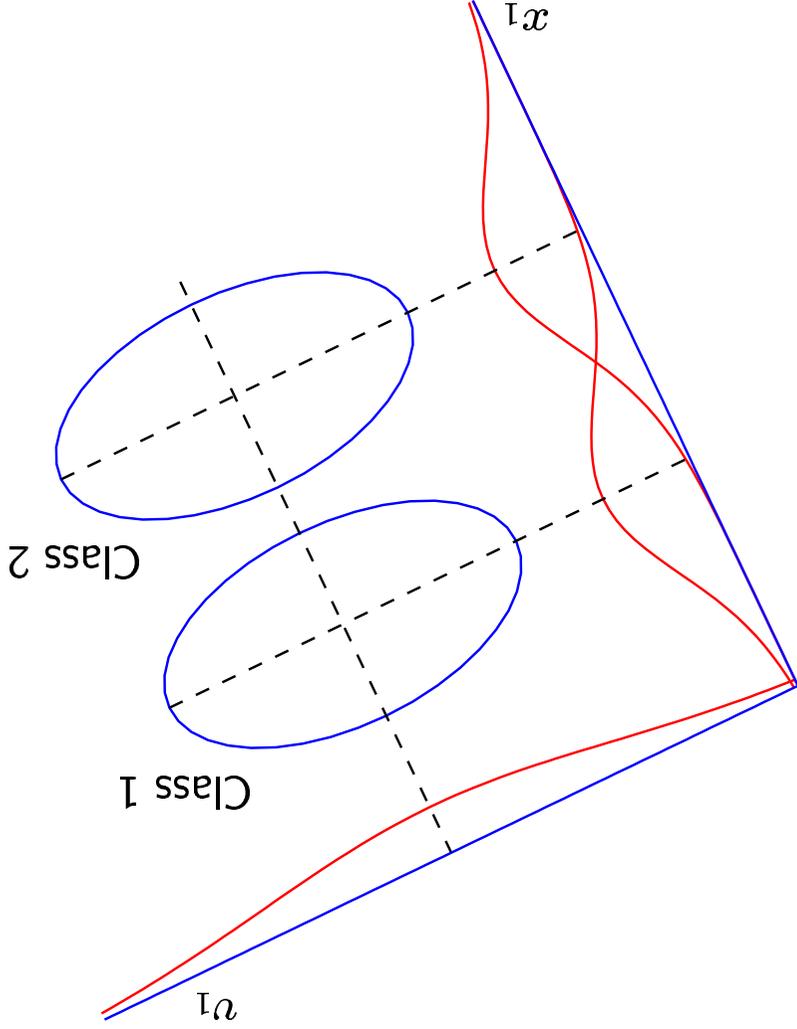
Projection scheme which

- minimises within class variance;
- maximises between class variance.

Example:

- two classes of 2-d observations;
- project down to optimal axis, x_1 .

Extensions: HLDA, MLDA, HDA, ...



State evolution process

Hidden Markov model

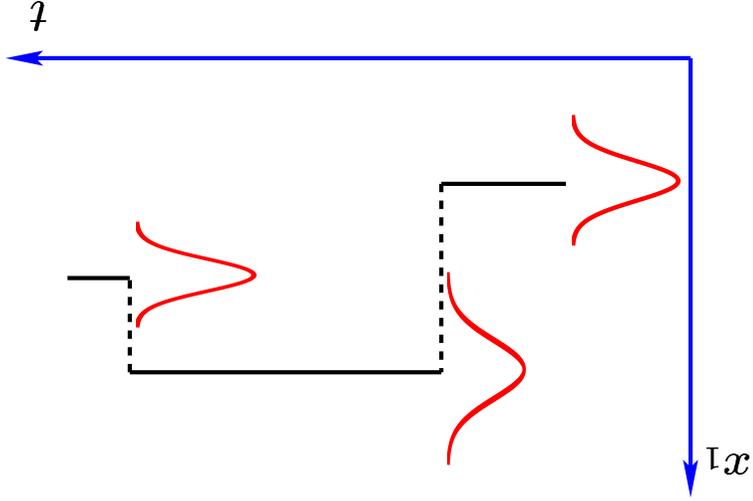
$$\mathcal{M}_{hmm} \sim \mathcal{M}_{hmm} \{ \mu_{ij}^j, \Sigma_{(x)}^j \}$$

- hidden discrete states $q_t = j, j \in [1, N_s];$

- diagonal covariances $\Sigma_{(x)}^j$.

Properties:

- piecewise constant trajectory;
- state conditional independence assumption.



Linear first-order Gauss-Markov

Generative model:

$$\begin{aligned}
 x_{t+1} &= A^t x_t + w_t \\
 w_t &\sim \mathcal{N}(w_t, \Sigma_{(x)}^t)
 \end{aligned}$$

- transition matrix A^t ;

- diagonal covariances $\Sigma_{(x)}$;

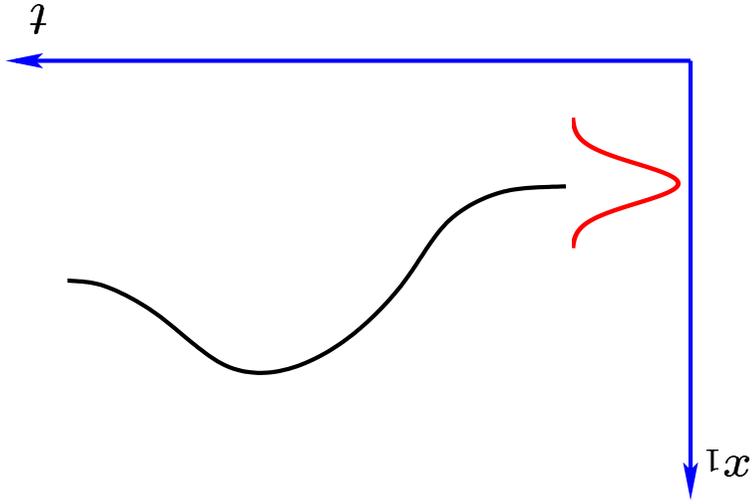
- initial distribution:

$$x_1 = \mathcal{N}(x_1, \Sigma_{(x)}^1)$$

Properties:

- continuous trajectory;

- models temporal correlation explicitly.



Bayesian networks

Graphical models to illustrate independence assumptions

- round node: continuous variable;

- square node: discrete variable;

- shaded node: observable;

- no arrow: conditional independence.

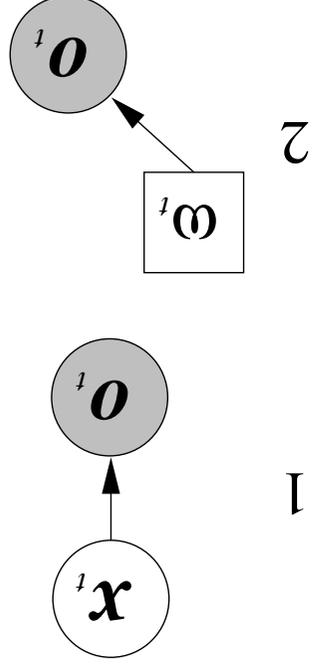
Examples:

1. Factor Analysis:

$$p(\mathbf{o}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{o}_t; \mathbf{C}_t \mathbf{x}_t + \boldsymbol{\mu}_t^{(o)}, \boldsymbol{\Sigma}_t^{(o)})$$

2. Mixture of Gaussians:

$$p(\mathbf{o}_t | \omega_t = n) = \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$$

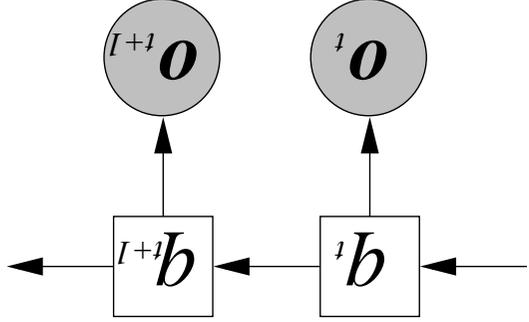


Dynamic Bayesian networks

Hidden Markov model

- hidden discrete Markov sequence: $P(q_{t+1}|q_1, \dots, q_t)$

- state conditional observation distribution: $p(o_t|q_t = j) = \mathcal{N}(o_t; \mu_j, \Sigma_j)$



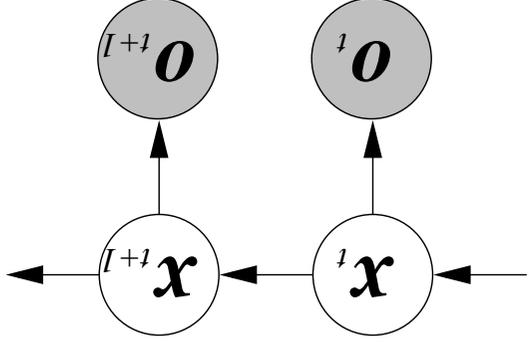
Linear dynamical system

- hidden Gauss-Markov sequence:

$$p(x_{t+1}|x_t) = \mathcal{N}(x_{t+1}; A x_t + u, \Sigma_x)$$

- factor analysis observation process:

$$p(o_t|x_t) = \mathcal{N}(o_t; C x_t + h, \Sigma_o)$$



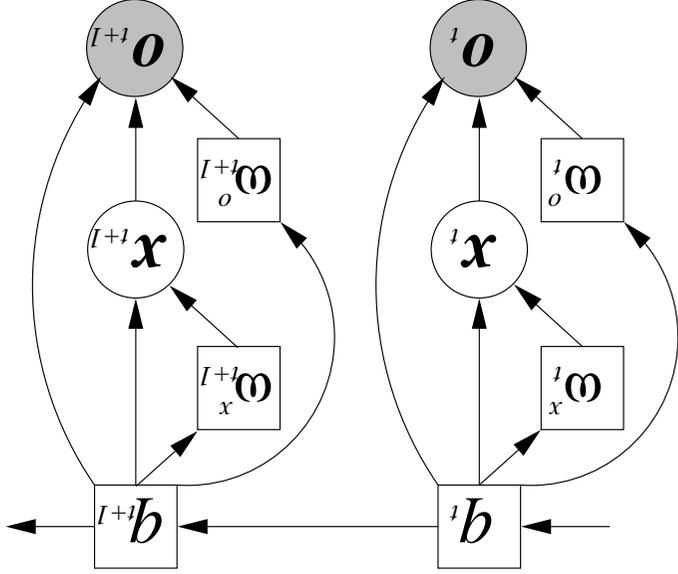
Factor analysed HMM

Generative model:

$$x_t = \mathcal{M}_{hmm}$$

$$o_t = C^t x_t + v_t$$

- mixture of Gaussians HMM: $\mathcal{M}_{hmm} = \{a_{ij}, c_{ij}^n, \mu_{(o)}^{jm}, \Sigma_{(o)}^{jm}\}$
- mixture of Gaussians noise: $v_t \sim \sum_m c_{(o)}^{jm} \mathcal{N}(o_{(o)}^{jm}, \Sigma_{(o)}^{jm})$



Effectively $M(o)M(x)$ component full covariance matrix system:

- mean vectors $\mu_{hmm}^j = C_{(x)}^{jn} \mu_{(o)}^{jn} + \mu_{(o)}^{jm}$

- covariance matrices $\Sigma_{hmm}^j = C_{(x)}^{jn} \Sigma_{(o)}^{jn} C_{(x)}^{jn} + \Sigma_{(o)}^{jm}$

Properties of FAHMM

Number of free parameters:

- order of $2(M^{(x)} - 1)k + pk + 2M^{(o)}d$ per state;
- arbitrary tying of model parameters.

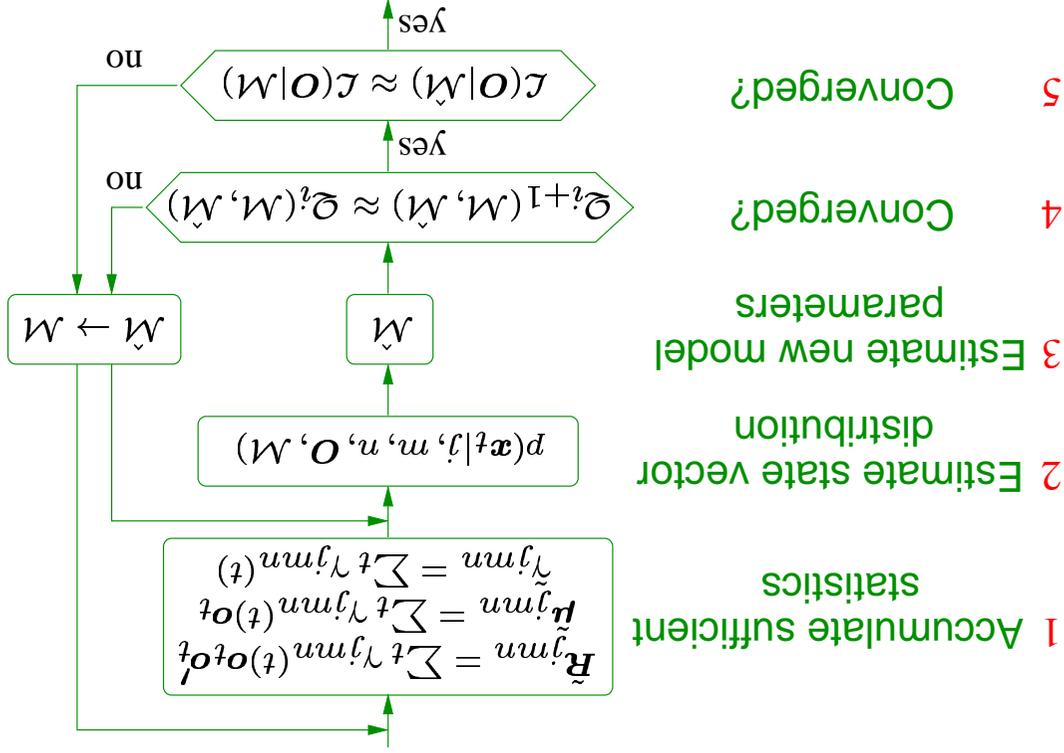
Related standard systems include:

- $M^{(x)} = 0$ \leftrightarrow hidden Markov model;
- $M^{(x)} = 1$ \leftrightarrow shared factor analysis;
- $M^{(o)} = 1$ \leftrightarrow dynamic version of independent factor analysis;
- $k = d$ and $a_t = 0$ \leftrightarrow semi-tied covariance matrix HMM (STC);
- $k > d$ and $a_t = 0$ \leftrightarrow extended MLLT (restricted).

Training FAHMM

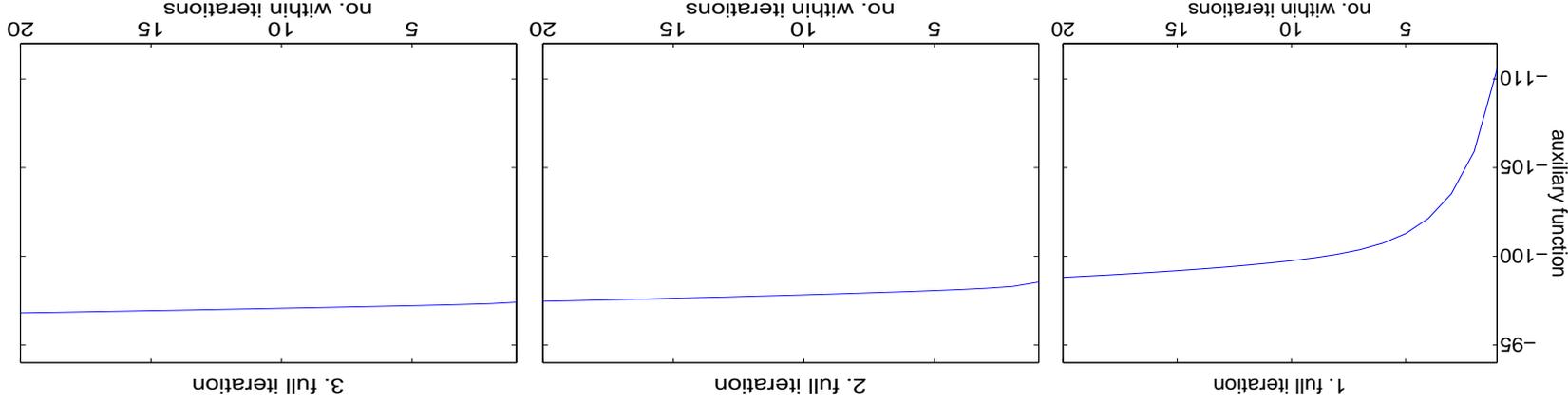
Expectation maximisation algorithm:

- auxiliary function $Q(M, \hat{M}) = E \{ \log p(O, \Omega^o, X, \Omega^x, \hat{Q} | \hat{M}) | O, M \}$;
- forward-backward algorithm gives $\gamma_{jmn}(t) = P(j, m, n | O, M)$



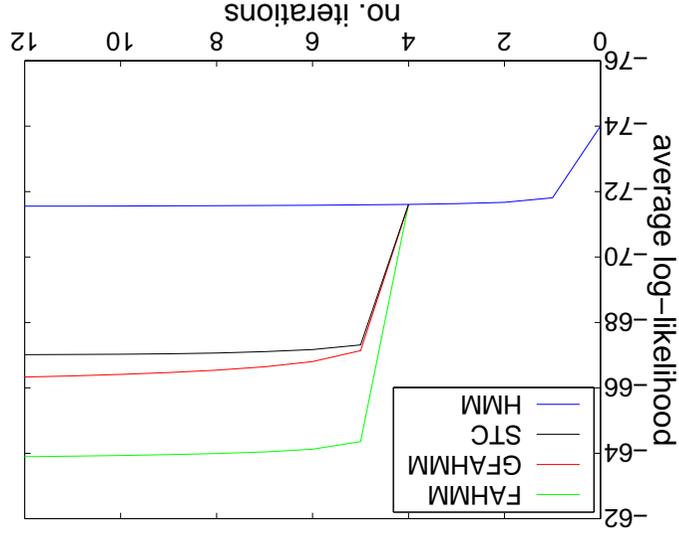
Training FAHMM continued

Multiple within iterations increase log-likelihood faster:



Log-likelihoods vs. full iterations:

- standard HMM;
- global full C STC;
- global C FAHMM ($k = 39$);
- untied FAHMM ($k = 13$).



Resource Management experiments

- 3990 sentence training data set (train+dev-aug);
- 39 dimensional feature vectors (MFCC+E+D+A);
- decision tree clustered cross-word triphone models;
- 1200 sentence evaluation data set (feb89+oct89+feb91+sep92);
- word-pair grammar;
- best HMM performance 3.99% (6 comps);
- global full transform STC;
- global observation matrix FAHMM with $k = 39$.

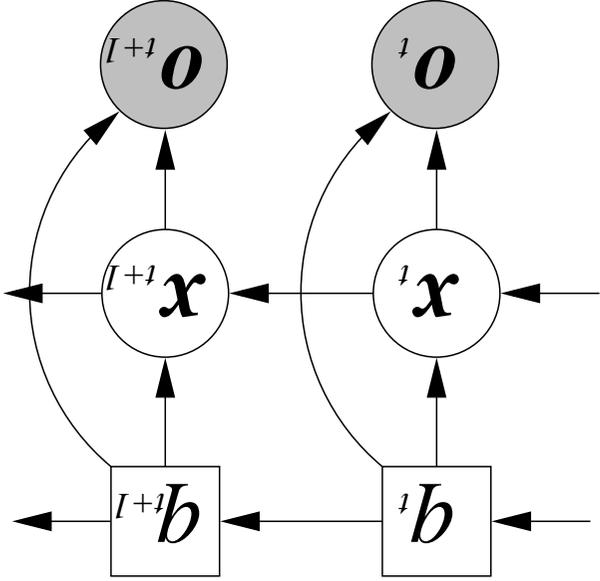
no. comps	1	5
HMM	7.79%	4.09%
STC	7.06%	3.83%
FAHMM	6.52%	3.68%

Switchboard Minitrain experiments

- 18 hour training data set;
- 39 dimensional feature vectors (MF-PLP+A+D);
- VTLN + mean and variance side based normalisation;
- decision tree clustered cross-word triphone models;
- 10 conversation sides of SWBD-2 and 10 of CHE evaluation set;
- untied FAHMM with $k = 13$;
- best baseline performance with 12 components;
- baseline performance exceeded with $M^{(o)} = 2$ and $M^{(x)} = 4$ (nfp 741);
- best performance with $M^{(o)} = 2$ and $M^{(x)} = 6$.

no. comps	WER	nfp
FAHMM	50.7%	793
HMM	51.0%	936

Switching linear dynamical system



$$\begin{aligned} x_{t+1} &= A^t x_t + w_t \\ o_t &= C^t x_t + v_t \end{aligned}$$

Generative model:

- HMM indexes standard LDS parameters: $\mu_{(i)}^j, \Sigma_{(i)}^j, A_j, \mu_{(o)}^j, \Sigma_{(o)}^j, C_j, \mu_j^j, \Sigma_j^j$;
- continuous state vector evolution.

Intractable inference!

Given the segmentation

- Kalman filtering for likelihood computation;
- Kalman smoothing for parameter estimation.

Stochastic segment model vs. SLDS

Stochastic segment model

- first segment model with standard LDS as the trajectory model;
- segments assumed independent given the discrete state; i.e., state vectors initialised on the segment boundaries;
- training alternates between segmentation and parameter estimation.

Switching LDS

- relaxes independent segment assumption by propagating state vector posteriors over the segment boundaries;
- exact inference results in exponential growth of paths (N^s).

Both model unimodal distributions!

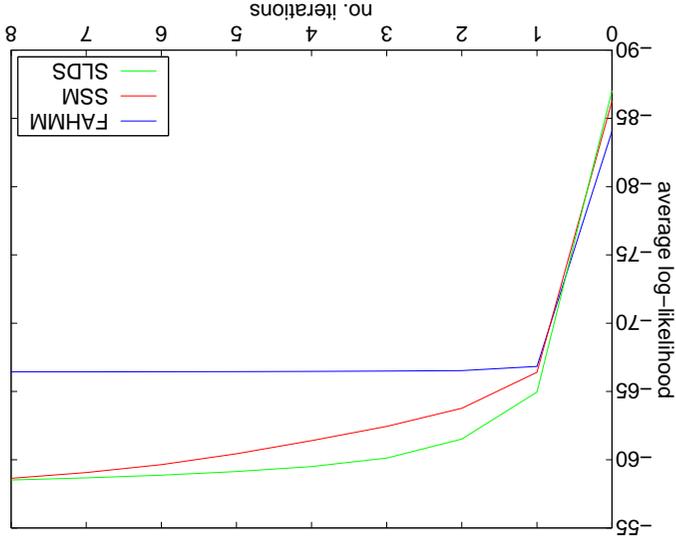
Resource Management experiments

Aim to compare discrete vs. continuous state evolution processes

- decision tree clustered cross-word triphone system at model level;
- single state per model in SSM and SLDS, $k = 13$;
- three state FAHMM, $k = 13$, all parameters except $\mu_j^{(x)}$ tied at model level;
- model aligned training data using alignment FAHMM.

Forced alignment training

- flat start tied three state **FAHMM**;
- flat start single state **SSM**;
- flat start single state **SLDS**.



Resource Management results

20-best re-scoring

- hypotheses and alignments using alignment FAHMM;
- hypotheses re-scored using FAHMM, SSM and SLDS.

Test data word error rates:

FAHMM	SSM	SLDS
11.00 %	10.69 %	13.45 %

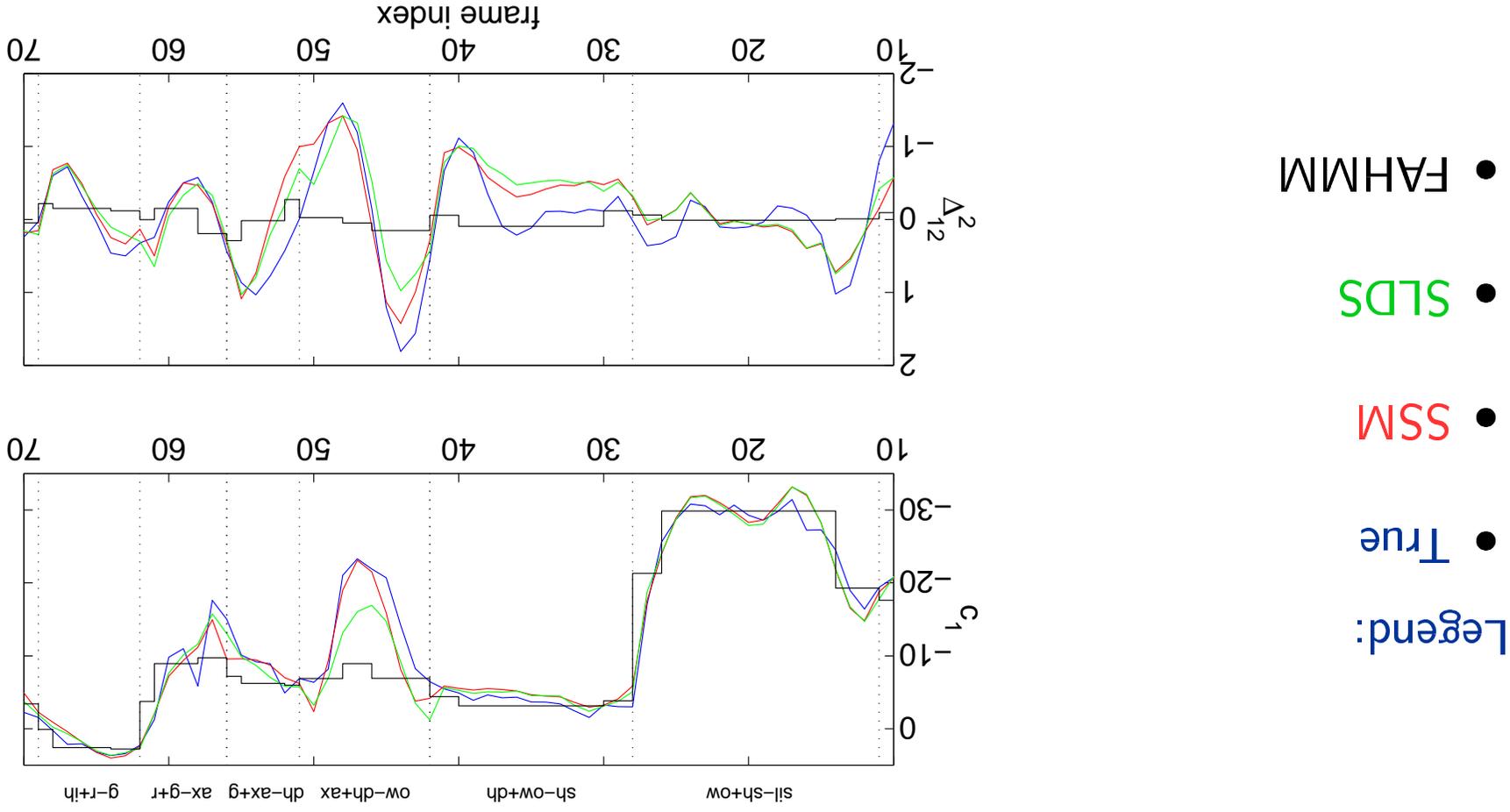
Subset of training data word error rates:

FAHMM	SSM	SLDS
4.30 %	3.96 %	4.30 %

SSM outperforms SLDS!

True and estimated trajectories

- estimates given by $\hat{o}_t = C^t \hat{x}_t + \mu_t^{(o)}$, where $\hat{x}_t = E\{x_t | O\}$.



Conclusions

- general framework of linear Gaussian models developed;
- developed and evaluated factor analysed HMMs;
- initial experiments with SSM and SLDS.

Future work

- develop decoding algorithms for SSM and SLDS;
- investigate multi-modal distributions for SSM and SLDS.

Any suggestions ...