

TOSHIBA

Leading Innovation >>>

Globally Normalized Model for Statistical Speech Synthesis



Heiga ZEN

Toshiba Research Europe Ltd.
Cambridge Research Lab.

Speech Synthesis Seminar Series @ CUED, Cambridge, UK
January 26th, 2011

Outline

Trajectory HMM

- Speech parameter generation
- Derivation of trajectory HMM
- Relationship between parameter generation & trajectory HMM
- Trajectory HMM as globally normalized model
- Parameter estimation

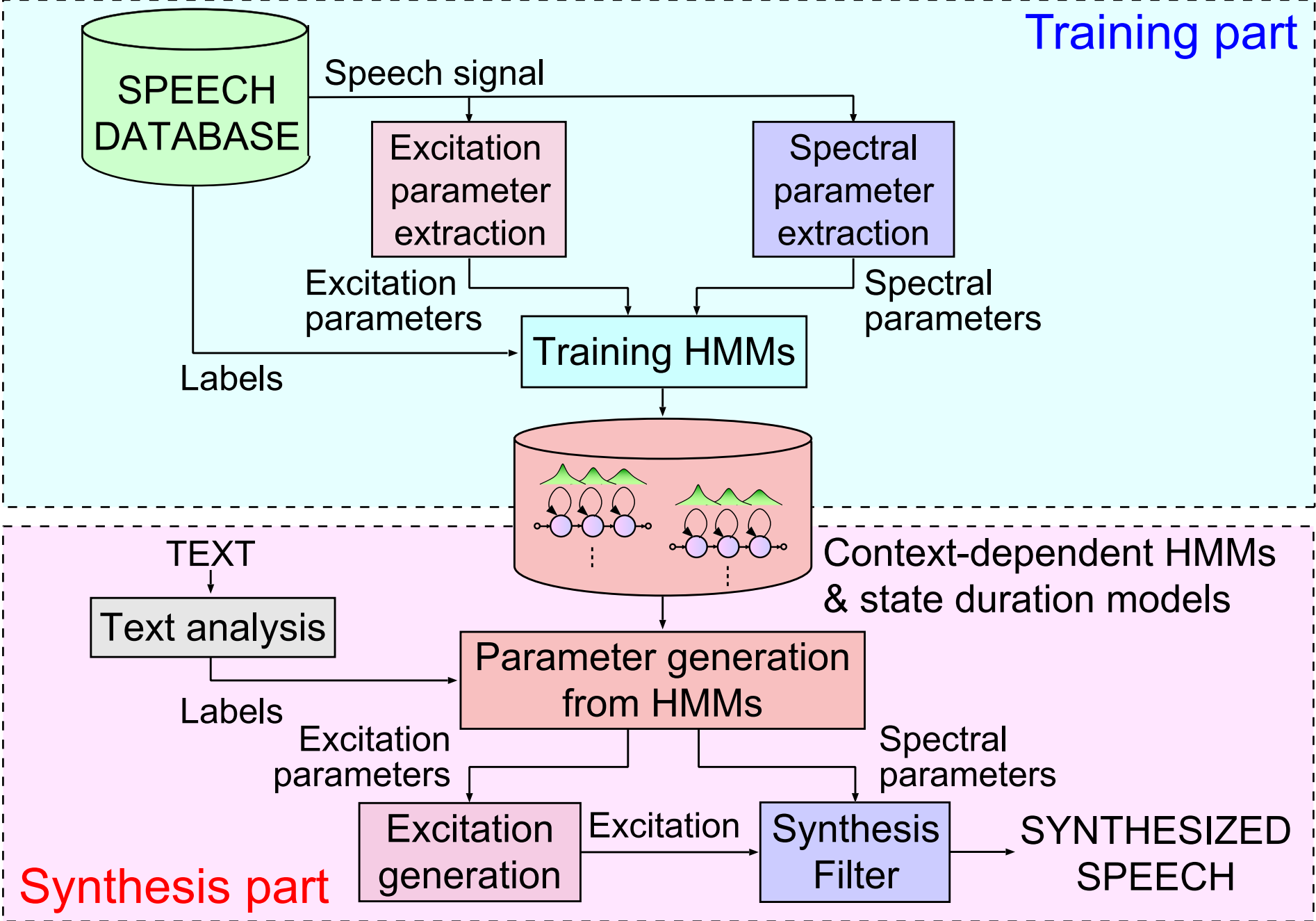
Min generation error (MGE) training & trajectory HMM

- Relationship
- Properties of MGE

(If time remains) Product of Experts (PoE)

- Combination of multiple AMs as PoE
- PoE & trajectory HMM

HMM-based speech synthesis system



Speech parameter generation algorithm

Determine a speech parameter vector sequence that maximizes its output probability given label l & HMM λ

$$\begin{aligned}\hat{\mathbf{o}} &= \arg \max_{\mathbf{o}} p(\mathbf{o} \mid l, \hat{\lambda}) \\ &= \arg \max_{\mathbf{o}} \sum_{\forall \mathbf{q}} p(\mathbf{o} \mid \mathbf{q}, \hat{\lambda}) p(\mathbf{q} \mid l, \hat{\lambda}) \\ &= \arg \max_{\mathbf{o}, \mathbf{q}} p(\mathbf{o} \mid \mathbf{q}, \hat{\lambda}) p(\mathbf{q} \mid l, \hat{\lambda}) \\ \hat{\mathbf{q}} &= \arg \max_{\mathbf{q}} p(\mathbf{q} \mid l, \hat{\lambda}) \\ \hat{\mathbf{o}} &= \arg \max_{\mathbf{o}} p(\mathbf{o} \mid \hat{\mathbf{q}}, \hat{\lambda})\end{aligned}$$

Output prob of \mathbf{o} given \mathbf{l} & HMM λ

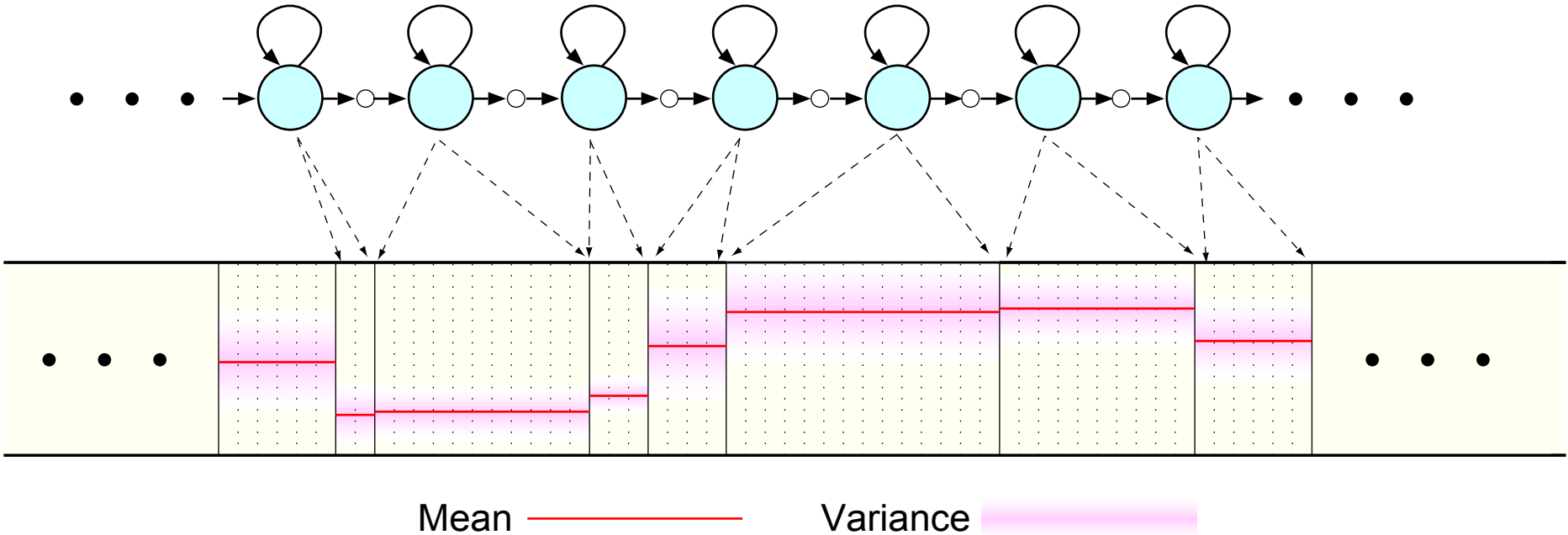
$$p(\mathbf{o} | \mathbf{l}, \lambda) = \sum_{\forall \mathbf{q}} \underbrace{p(\mathbf{o} | \mathbf{q}, \lambda)}_{\text{state-output}} \underbrace{P(\mathbf{q} | \mathbf{l}, \lambda)}_{\text{state-transition}}$$

$$p(\mathbf{o} | \mathbf{q}, \lambda) = \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) \quad \Leftarrow \text{single Gaussian}$$

$$= \mathcal{N} \left(\begin{bmatrix} \mathbf{o}_1 \\ \mathbf{o}_2 \\ \vdots \\ \mathbf{o}_T \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_{q_1} \\ \boldsymbol{\mu}_{q_2} \\ \vdots \\ \boldsymbol{\mu}_{q_T} \end{bmatrix}, \underbrace{\begin{bmatrix} \boldsymbol{\Sigma}_{q_1} & & & \mathbf{0} \\ & \boldsymbol{\Sigma}_{q_2} & & \\ & & \ddots & \\ \mathbf{0} & & & \boldsymbol{\Sigma}_{q_T} \end{bmatrix}}_{\text{diagonal}} \right)$$
$$= \mathcal{N}(\mathbf{o}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$

Generated trajectory

$$\begin{aligned}\hat{o} &= \arg \max_o p(o \mid \hat{q}, \hat{\lambda}) \\ &= \arg \max_o \mathcal{N}(o ; \mu_{\hat{q}}, \Sigma_{\hat{q}}) \\ &= \mu_{\hat{q}} \quad \Leftarrow \text{mean vector sequence}\end{aligned}$$



Relationship between o and c

$$o_t = \begin{bmatrix} c_t \\ \Delta c_t \end{bmatrix} \begin{matrix} \leftarrow \text{static} \\ \leftarrow \text{dynamic} \end{matrix} \quad \Delta c_t = c_t - c_{t-1}$$

$$\begin{matrix} o \\ o_{t-1} \\ o_t \\ o_{t+1} \end{matrix} \begin{bmatrix} \vdots \\ c_{t-1} \\ \Delta c_{t-1} \\ c_t \\ \Delta c_t \\ c_{t+1} \\ \Delta c_{t+1} \\ \vdots \end{bmatrix} = \begin{matrix} W \\ \dots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \dots \end{matrix} \begin{bmatrix} \dots & \vdots & \vdots & \vdots & \vdots & \dots \\ \dots & 0 & I & 0 & 0 & \dots \\ \dots & -I & I & 0 & 0 & \dots \\ \dots & 0 & 0 & I & 0 & \dots \\ \dots & 0 & -I & I & 0 & \dots \\ \dots & 0 & 0 & 0 & I & \dots \\ \dots & 0 & 0 & -I & I & \dots \\ \dots & \vdots & \vdots & \vdots & \vdots & \dots \end{bmatrix} \begin{matrix} c \\ \vdots \\ c_{t-2} \\ c_{t-1} \\ c_t \\ c_{t+1} \\ \vdots \end{matrix}$$

Speech parameter generation algorithm

$$\begin{aligned}\hat{o} &= \arg \max_o p(o \mid \hat{q}, \hat{\lambda}) \Big|_{o=Wc} \\ &= \arg \max_o \mathcal{N}(o; \mu_{\hat{q}}, \Sigma_{\hat{q}}) \Big|_{o=Wc}\end{aligned}$$

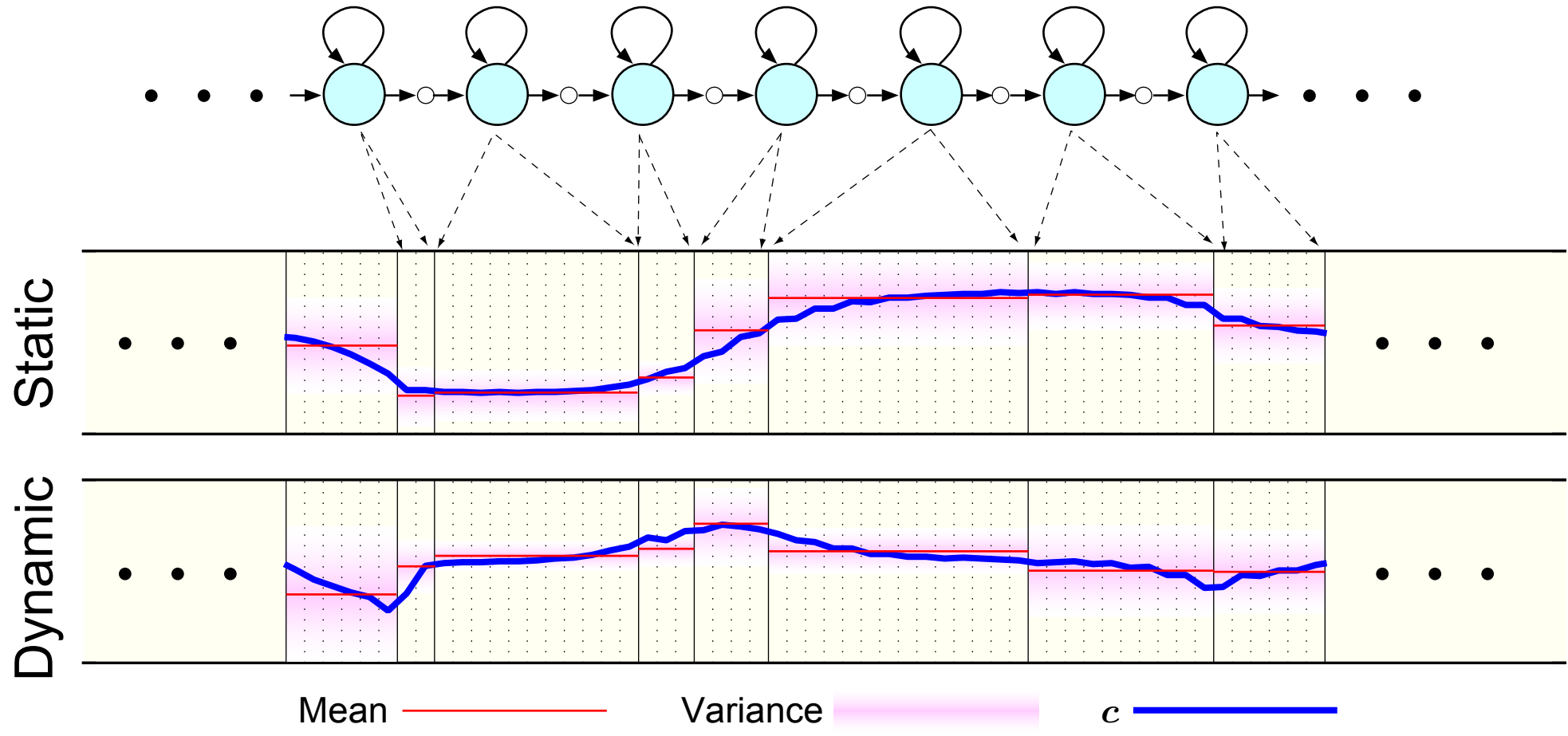


$$\hat{c} = \arg \max_c \mathcal{N}(Wc; \mu_{\hat{q}}, \Sigma_{\hat{q}})$$



$$W^T \Sigma_q^{-1} W \hat{c} = W^T \Sigma_q^{-1} \mu_q$$

Generated trajectory



Inconsistency between training & synthesis

Training & synthesis parts are inconsistent

- Training part

- * Baum-Welch training
- * Labels are often given manually
- * Model training model w/o dynamic feature constraints

- Synthesis part

- * Viterbi (single-path) approximation
- * Labels are often given automatically (by text analysis)
- * Parameter generation w/ dynamic feature constraints

How about introducing dyn feature constraints to training?

Output prob of \mathbf{o} given \mathbf{l} & HMM λ

$$p(\mathbf{o} | \mathbf{l}, \lambda) = \sum_{\forall \mathbf{q}} \underbrace{p(\mathbf{o} | \mathbf{q}, \lambda)}_{\text{state-output}} \underbrace{P(\mathbf{q} | \mathbf{l}, \lambda)}_{\text{state-transition}}$$

$$p(\mathbf{o} | \mathbf{q}, \lambda) = \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) \quad \Leftarrow \text{single Gaussian}$$

$$= \mathcal{N} \left(\begin{bmatrix} \mathbf{o}_1 \\ \mathbf{o}_2 \\ \vdots \\ \mathbf{o}_T \end{bmatrix}; \begin{bmatrix} \boldsymbol{\mu}_{q_1} \\ \boldsymbol{\mu}_{q_2} \\ \vdots \\ \boldsymbol{\mu}_{q_T} \end{bmatrix}, \underbrace{\begin{bmatrix} \boldsymbol{\Sigma}_{q_1} & & & \mathbf{0} \\ & \boldsymbol{\Sigma}_{q_2} & & \\ & & \ddots & \\ \mathbf{0} & & & \boldsymbol{\Sigma}_{q_T} \end{bmatrix}}_{\text{diagonal}} \right)$$
$$= \mathcal{N}(\mathbf{o}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$

Inconsistency in HMM w/ dynamic features

Under $o = Wc$

$$p(o | q, \lambda) = \mathcal{N}(Wc; \mu_q, \Sigma_q) \Rightarrow \text{incorrect!}$$

Why?

$$\int_c \mathcal{N}(Wc; \mu_q, \Sigma_q) dc \neq 1 \Rightarrow \text{integral over } c \text{ must be 1 to be a valid PDF}$$

Why does this happen?

Static features \Rightarrow random variables

Dynamic features \Rightarrow Not random variables!!

Normalization

Normalized to achieve valid PDF

$$\begin{aligned} Z_q &= \int_{\mathbf{c}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q) d\mathbf{c} \\ &= \frac{\sqrt{(2\pi)^{MT} |\mathbf{P}_q|}}{\sqrt{(2\pi)^{2MT} |\boldsymbol{\Sigma}_q|}} \exp \left\{ -\frac{1}{2} \left(\boldsymbol{\mu}_q^\top \boldsymbol{\Sigma}_q^{-1} \boldsymbol{\mu}_q - \mathbf{r}_q^\top \mathbf{P}_q \mathbf{r}_q \right) \right\} \end{aligned}$$

$\mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q) \Rightarrow$ **invalid PDF!**

$\frac{1}{Z_q} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q) \Rightarrow$ **valid PDF!!**

Definition of trajectory HMM

Use normalized Gaussian \Rightarrow trajectory HMM is defined

$$p(\mathbf{c} | \mathbf{l}, \lambda) = \sum_{\forall \mathbf{q}} \underbrace{p(\mathbf{c} | \mathbf{q}, \lambda)}_{\text{state-output}} \underbrace{P(\mathbf{q} | \mathbf{l}, \lambda)}_{\text{state-transition}}$$

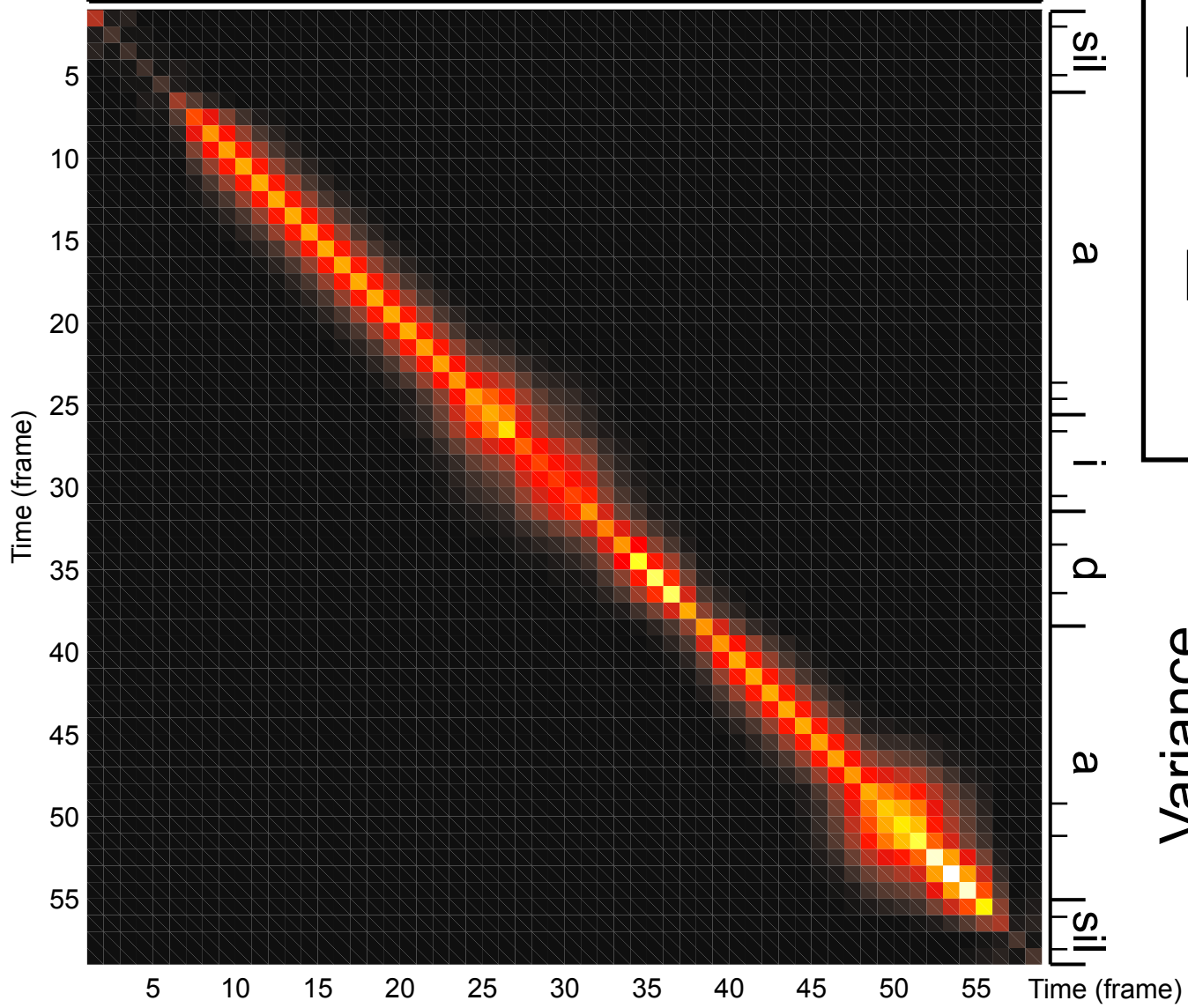
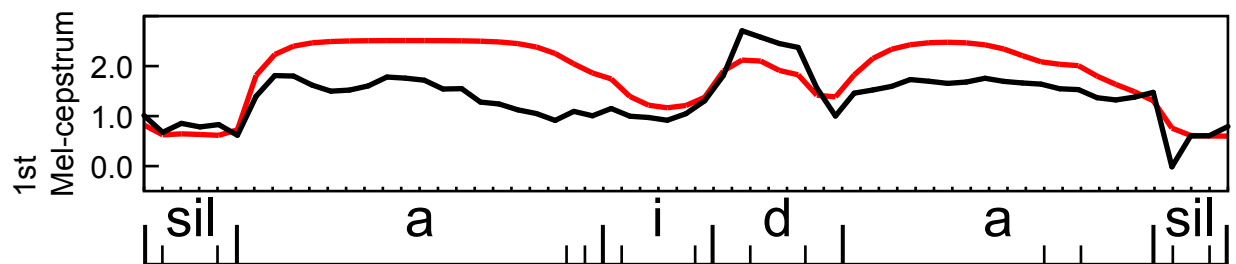
$$p(\mathbf{c} | \mathbf{q}, \lambda) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}}) \Leftarrow \text{normalized Gaussian}$$

$$= \mathcal{N}(\mathbf{c}; \underbrace{\bar{\mathbf{c}}_{\mathbf{q}}}_{\substack{\text{mean} \\ \text{vec}}}, \underbrace{P_{\mathbf{q}}}_{\substack{\text{cov} \\ \text{mat}}}) \Leftarrow \text{Gaussian over } \mathbf{c}$$

$$\mathbf{R}_{\mathbf{q}} \bar{\mathbf{c}}_{\mathbf{q}} = \mathbf{r}_{\mathbf{q}}$$

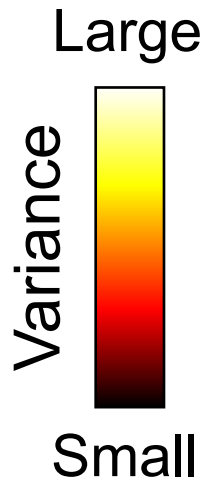
$$\mathbf{R}_{\mathbf{q}} = \mathbf{W}^{\top} \boldsymbol{\Sigma}_{\mathbf{q}}^{-1} \mathbf{W} = P_{\mathbf{q}}^{-1}$$

$$\mathbf{r}_{\mathbf{q}} = \mathbf{W}^{\top} \boldsymbol{\Sigma}_{\mathbf{q}}^{-1} \boldsymbol{\mu}_{\mathbf{q}}$$



Mean & variance
 \Rightarrow varies in a state

Frame correlation
 \Rightarrow captured by P_q



Covariance matrix P_q

Trajectory HMM & speech parameter generation

Mean vector of trajectory HMM

$$\mathbf{W}^\top \Sigma_q^{-1} \mathbf{W} \bar{\mathbf{c}}_q = \mathbf{W}^\top \Sigma_q^{-1} \boldsymbol{\mu}_q$$

Trajectory by speech parameter generation algorithm

$$\mathbf{W}^\top \Sigma_q^{-1} \mathbf{W} \hat{\mathbf{c}} = \mathbf{W}^\top \Sigma_q^{-1} \boldsymbol{\mu}_q$$

⇒ they are identical

Trajectory HMM as globally normalized model

HMM \Rightarrow locally (frame-level) normalized model

$$\begin{aligned} p(\mathbf{o} \mid \mathbf{q}, \lambda) &= \prod_{t=1}^T p(\mathbf{o}_t \mid q_t, \lambda) \\ &= \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t ; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) \end{aligned}$$

Trajectory HMM \Rightarrow globally (utt-level) normalized model

$$\begin{aligned} p(\mathbf{c} \mid \mathbf{q}, \lambda) &= \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c} ; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}}) \\ &= \frac{1}{Z_{\mathbf{q}}} \prod_{t=1}^T \mathcal{N}\left(\left[\mathbf{c}_t^{\top}, \Delta \mathbf{c}_t^{\top}\right]^{\top} ; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}\right) \end{aligned}$$

Estimating trajectory HMM parameters

ML estimation of trajectory HMM

$$\hat{\lambda} = \arg \max_{\lambda} p(\mathbf{c} \mid \mathbf{l}, \lambda)$$

Locally normalized model

Parameter estimation for each state can be done separately

Globally normalized model

Parameter estimation of all states have to be done jointly

$$\boldsymbol{\mu} = [\boldsymbol{\mu}_1^{\top}, \boldsymbol{\mu}_2^{\top}, \dots, \boldsymbol{\mu}_N^{\top}]^{\top} \quad : \text{all mean vectors}$$

$$\boldsymbol{\phi} = [\boldsymbol{\Sigma}_1^{-1}, \boldsymbol{\Sigma}_2^{-1}, \dots, \boldsymbol{\Sigma}_N^{-1}] \quad : \text{all precision matrices}$$

Parameter update formulae

$$\sum_{\forall \mathbf{q}} p(\mathbf{q} | \mathbf{c}, \lambda') \mathbf{S}_q^\top \Sigma_q^{-1} \mathbf{W} \mathbf{P}_q \mathbf{W}^\top \Sigma_q^{-1} \mathbf{S}_q \boldsymbol{\mu}$$
$$= \sum_{\forall \mathbf{q}} p(\mathbf{q} | \mathbf{c}, \lambda') \mathbf{S}_q^\top \Sigma_q^{-1} \mathbf{W} \mathbf{c}$$

mean vectors \Rightarrow closed form

$$\frac{\partial Q(\lambda, \lambda')}{\partial \phi} = \sum_{\forall \mathbf{q}} p(\mathbf{q} | \mathbf{c}, \lambda') \left\{ \frac{1}{2} \mathbf{S}_q^\top \text{diag}^{-1} \left(\mathbf{W} \mathbf{P}_q \mathbf{W}^\top - \mathbf{W} \mathbf{c} \mathbf{c}^\top \mathbf{W}^\top \right. \right.$$
$$\left. \left. + \mathbf{W} \bar{\mathbf{c}}_q \bar{\mathbf{c}}_q^\top \mathbf{W}^\top + \boldsymbol{\mu}_q \mathbf{c}^\top \mathbf{W}^\top + \mathbf{W} \mathbf{c} \boldsymbol{\mu}_q^\top - \boldsymbol{\mu}_q \bar{\mathbf{c}}_q^\top \mathbf{W}^\top - \mathbf{W} \bar{\mathbf{c}}_q \boldsymbol{\mu}_q^\top \right) \right\}$$

covariance matrices \Rightarrow numerical optimization

Drawback of trajectory HMM training

Exact EM is intractable

- Computing posterior prob of q is intractable
- Single-path (Viterbi) or Monte Carlo approximation

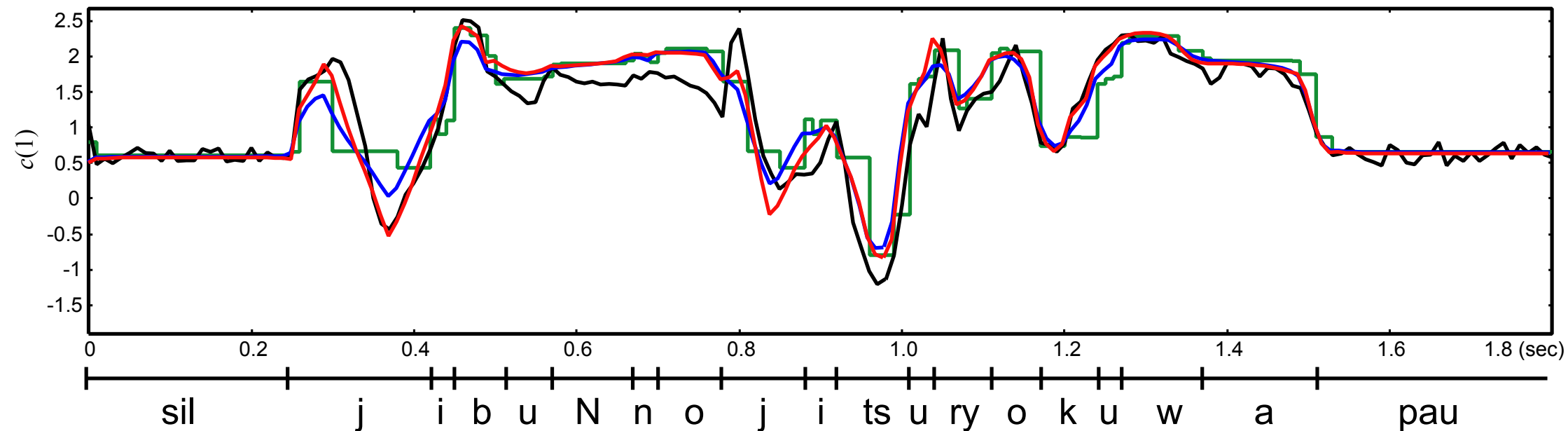
Exact tree-based clustering is also intractable

- Splitting one nodes affects the other nodes
- Trees built for HMMs are often used

Computationally & memory intensive

- High dimensional matrix operations
- Numerical optimization

Effect of parameter reestimation



— Training data

— Mean sequence of the HMM

— Mean sequence of the trajectory HMM (w/o update)

— Mean sequence of the trajectory HMM (with update)

Outline

Trajectory HMM

- Speech parameter generation
- Derivation of trajectory HMM
- Relationship between parameter generation & trajectory HMM
- Trajectory HMM as globally normalized model
- Parameter estimation

Min generation error (MGE) training & trajectory HMM

- Relationship
- Properties of MGE

(If time remains) Product of Experts (PoE)

- Combination of multiple AMs as PoE
- PoE & trajectory HMM

MGE training & trajectory HMM (1)

ML training & MGE training w/ Euclidean dist [Wu;'06]

$$\begin{aligned}\hat{\lambda}_{\text{ML}} &= \arg \max_{\lambda} p(\mathbf{c} \mid \mathbf{q}, \lambda) \\ &= \arg \max_{\lambda} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_q, \mathbf{P}_q)\end{aligned}$$

$$\begin{aligned}\hat{\lambda}_{\text{MGE}} &= \arg \min_{\lambda} \mathcal{E}(\mathbf{c}; \mathbf{q}, \lambda) \\ &= \arg \min_{\lambda} \underbrace{\|\mathbf{c} - \bar{\mathbf{c}}_q\|_2}_{\text{Euclidean distance}} \Leftarrow \text{MMSE estimation} \\ &= \arg \max_{\lambda} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_q, \mathbf{I}) \Leftarrow \text{Identity covariance matrix}\end{aligned}$$

MGE training & trajectory HMM (3)

Performance of ML & MGE w/ Euc is similar, why?

⇒ Due to speech parameter generation algorithm

$$\hat{\mathbf{c}}_{\text{ML}} = \arg \max_{\mathbf{c}} p(\mathbf{c} \mid \hat{\mathbf{q}}, \hat{\lambda}_{\text{ML}})$$

$$= \arg \max_{\mathbf{c}} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\hat{\mathbf{q}}}, \mathbf{P}_{\hat{\mathbf{q}}})$$

$$= \bar{\mathbf{c}}_{\hat{\mathbf{q}}}$$

$$\hat{\mathbf{c}}_{\text{MGE}} = \arg \max_{\mathbf{c}} p(\mathbf{c} \mid \hat{\mathbf{q}}, \hat{\lambda}_{\text{MGE}})$$

$$= \arg \max_{\mathbf{c}} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\hat{\mathbf{q}}}, \mathbf{I})$$

$$= \bar{\mathbf{c}}_{\hat{\mathbf{q}}}$$

MGE training & trajectory HMM (4)

Random sampling from ML & MGE w/ Euc distance

ML

$$\tilde{\mathbf{c}}_{\text{ML}} \sim \mathcal{N}(\bar{\mathbf{c}}_{\hat{q}}, \mathbf{P}_{\hat{q}})$$

⇒ Temporal correlations will be kept

MGE

$$\tilde{\mathbf{c}}_{\text{MGE}} \sim \mathcal{N}(\bar{\mathbf{c}}_{\hat{q}}, \mathbf{I})$$

⇒ Temporal correlations will be discarded

MGE training & trajectory HMM (5)

Which is better, ML or MGE?

- w/ parameter generation, MGE is more reasonable

* MGE $\Rightarrow \mu$ & Σ to represent mean trajectory

* ML $\Rightarrow \mu$ for mean trajectory, Σ for mean trj & temporal cov

\Rightarrow MGE can focus on modeling mean trajectory

- w/ random sampling, ML is more reasonable

* MGE *ignores* temporal correlations

* ML *models* temporal correlations

Summary

Trajectory HMM

- Derived from HMM w/ dynamic feature constraints
- Can be viewed as a globally normalized model
- All states need to be estimated jointly
- Generated params = mean vector of trajectory HMM

MGE training

- MGE w/ Euclid distance = MMSE estimation of trajectory HMM
- w/ speech parameter generation algorithm (ML parm gen)
 - ⇒ ML & MGE work similarly
- w/ random sampling
 - ⇒ MGE won't work well

Outline

Trajectory HMM

- Speech parameter generation
- Derivation of trajectory HMM
- Relationship between parameter generation & trajectory HMM
- Trajectory HMM as globally normalized model
- Parameter estimation

Min generation error (MGE) training & trajectory HMM

- Relationship
- Properties of MGE

(If time remains) Product of Experts (PoE)

- Combination of multiple AMs as PoE
- PoE & trajectory HMM

Combination of multiple acoustic models

Combine multiple AMs to reduce over-smoothing

- * Training; estimate multiple-level AMs *individually*

$$\hat{\lambda}_i = \arg \max_{\lambda_i} p(f_i(\mathbf{c}) | \lambda_i) \quad i = 1, \dots, M$$

- * Synthesis; generate \mathbf{c} that *jointly* maximize output probs from AMs

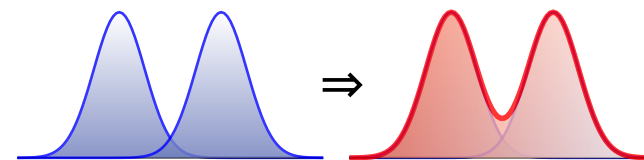
$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} \sum_{i=1}^M \alpha_i \log p(f_i(\mathbf{c}) | \hat{\lambda}_i)$$

- * Feature function, $f_i(\mathbf{c})$, extracts acoustic feats for i -th AM from \mathbf{c}
 - e.g., dynamic feats, DCT, average, summation, global variance
- * Parameters of AMs, λ_i , are trained *independently*
 - Use weights to control balance among AMs
- * Weights, α_i , are determined by *held-out data* (or tuned manually)

Mixture model vs Product model

Mixture of experts

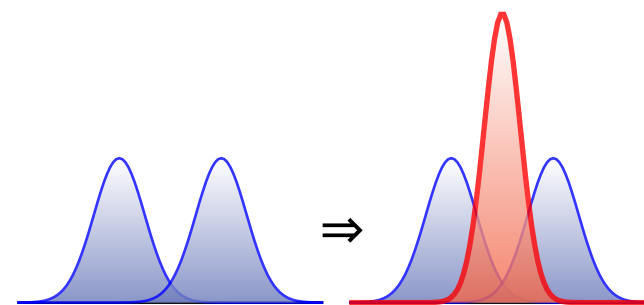
$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \sum_{i=1}^M \alpha_i p(f_i(\mathbf{c}) \mid \lambda_i)$$



- * Data is generated from **union** of experts
- * Robust for modeling data with many variations
- * GMM \rightarrow Mixture of Gaussians

Product of experts [Hinton;'02]

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$



- * Data is generated from **intersection** of experts
- * Efficient for modeling data with many constraints
- * PoG \rightarrow Product of Gaussians

Combination of multiple AMs as PoE

Combination of multiple AMs can be viewed as PoE

$$\begin{aligned}\hat{\mathbf{c}} &= \arg \max_{\mathbf{c}} p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \arg \max_{\mathbf{c}} \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} \\ &= \arg \max_{\mathbf{c}} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} = \arg \max_{\mathbf{c}} \sum_{i=1}^M \alpha_i \log p(f_i(\mathbf{c}) \mid \lambda_i)\end{aligned}$$

* Generating \mathbf{c} from combination of multiple AMs

→ Equivalent to generating \mathbf{c} from PoE consisting of AMs

* Regarding combination of multiple AMs as PoE

→ **Jointly** estimate multiple AMs

$$\{\hat{\lambda}_1, \dots, \hat{\lambda}_M\} = \arg \max_{\lambda_1, \dots, \lambda_M} \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$

Product of Gaussians

Product of Gaussians (PoG)

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \mathcal{N}(f_i(\mathbf{c}); \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

- * Special case of PoE; All experts are Gaussian
- * If all feature functions are linear
 - PoG also becomes Gaussian
 - Normalization constant

$$Z = \int \prod_{i=1}^M \mathcal{N}(f_i(\mathbf{c}); \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{c}$$

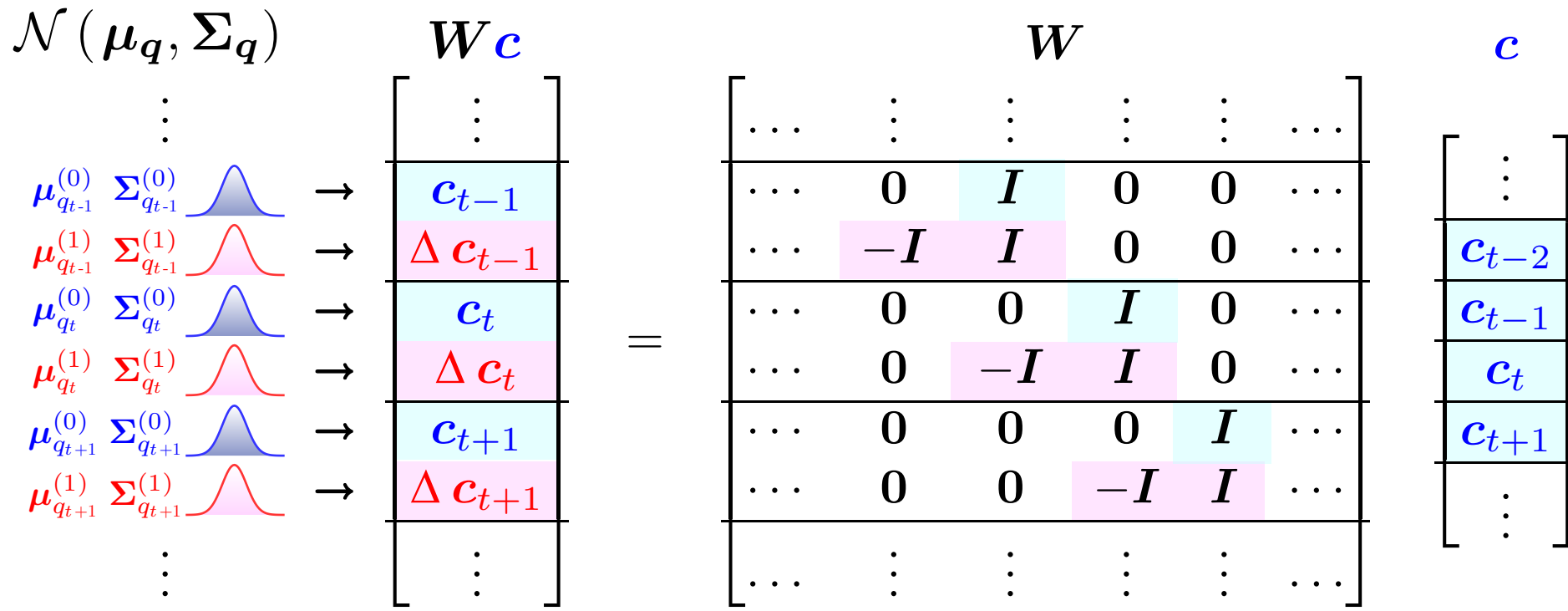
can be computed in **closed form**

Trajectory HMM as product of Gaussians

Trajectory HMM can be viewed as PoG [Williams;'05, Zen;'07]

$$p(\mathbf{c} | \lambda) = \sum_{\forall \mathbf{q}} p(\mathbf{c} | \mathbf{q}, \lambda) p(\mathbf{q} | \lambda)$$

$$p(\mathbf{c} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$



Trajectory HMM as product of Gaussians

Trajectory HMM can be viewed as PoG [Williams;'05, Zen;'07]

$$p(\mathbf{c} | \lambda) = \sum_{\forall \mathbf{q}} p(\mathbf{c} | \mathbf{q}, \lambda) p(\mathbf{q} | \lambda)$$

$$p(\mathbf{c} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$

$$= \frac{1}{Z_{\mathbf{q}}} \prod_{t=1}^T \prod_{d=0}^2 \mathcal{N}\left(f_t^{(d)}(\mathbf{c}); \boldsymbol{\mu}_{q_t}^{(d)}, \boldsymbol{\Sigma}_{q_t}^{(d)}\right) \quad f_t^{(d)}(\mathbf{c}) : d\text{-th dyn feat at frame } t$$

$$Z_{\mathbf{q}} = \int \prod_{t=1}^T \prod_{d=0}^2 \mathcal{N}\left(f_t^{(d)}(\mathbf{c}); \boldsymbol{\mu}_{q_t}^{(d)}, \boldsymbol{\Sigma}_{q_t}^{(d)}\right) d\mathbf{c}$$

- * Experts are Gaussians, feature functions are dynamic features
- * Gaussian experts are multiplied over time

General PoE (non-linear feat or non-Gaussian)

General form of PoE

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$

- * Feature functions can be non-linear, experts can be non-Gaussian
- * Normalization term has no closed form
- * Training is complex, usually normalization term is approximated

Example: global variance (GV) [Toda;'07]

$$p(\mathbf{c} \mid \mathbf{q}, \lambda, \lambda_{GV}) = \frac{1}{Z_q} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_q, \mathbf{P}_q)^\alpha \mathcal{N}(f_v(\mathbf{c}); \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)$$

$$f_v(\mathbf{c}) = \frac{1}{T} \sum_{t=1}^T \text{diag} \left[(\mathbf{c}_t - \bar{\mathbf{c}}) (\mathbf{c}_t - \bar{\mathbf{c}})^\top \right] : \text{intra-utt variance, } \mathbf{quadratic}$$

Contrastive divergence learning

Contrastive divergence learning [Hinton;'02]

- * Training algorithm for general PoE
- * Combination of sampling & gradient methods

1. Draw J samples from PoE

$$\mathbf{c}^{(j)} \sim p(\mathbf{c} | \boldsymbol{\lambda}) \quad j = 1, \dots, J \quad \boldsymbol{\lambda} = \{\lambda_1, \dots, \lambda_M\}: \text{PoE model params}$$

2. Compute approximated derivative of log likelihood w.r.t. $\boldsymbol{\lambda}$

$$\frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \approx \underbrace{\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_0}}_{\text{expectation over data}} - \underbrace{\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_J}}_{\text{expectation over samples}}$$

3. Update model params using gradient method

$$\boldsymbol{\lambda}' = \boldsymbol{\lambda} - \eta \cdot \left(\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_0} - \left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_j} \right) \quad \eta : \text{learning rate}$$

$$\boldsymbol{\lambda} = \boldsymbol{\lambda}'$$

4. Iterate 1-3 until converge

Experiment - Multiple-Level Dur Models as PoE

Experimental conditions

- * Training data; 2,469 utterances
- * Development data; 137 utterances
 - Used to optimize weights in conventional method
 - Weights were optimized to minimize RMSE by grid search
 - Baseline & proposed method did not use development data
- * Almost the same training setup as Nitech-HTS 2005 [Zen;'06]
- * Test data; 137 utterances
- * State, phone, & syllable-level models were clustered individually
 - # of leaf nodes
 - * state; 607, phoneme; 1,364, syllable; 281

Experimental Results

Duration prediction results (RMSE in frame (rel imp))

Model	Phoneme	Syllable	Pause
Baseline (st)	5.08 (ref)	8.98 (ref)	35.0 (ref)
uPoE (st*ph)	4.62 (9.1%)	8.13 (9.5%)	31.8 (9.1%)
uPoE (st*ph*syl)	4.62 (9.1%)	8.11 (9.7%)	31.8 (9.1%)
PoE (st*ph)	4.60 (9.4%)	8.04 (10.5%)	31.9 (8.9%)
PoE (st*ph*syl)	4.57 (10.0%)	8.02 (10.7%)	31.9 (8.9%)

st; state only, st*ph; state & phoneme, st*ph*syl; state, phoneme, & syllable

uPoE; individually trained multiple-level duration models with optimized weights

PoE; jointly estimated multiple-level duration models

Experiment - Global Variance as PoE

Experimental conditions

- * Training data; 2,469 utterances
 - Training data was split into mini-batch (250 utterances)
 - 10 MCMC sampling at each contrastive divergence learning
 - * Hybrid Monte Carlo with 20 leap-frog steps
 - * Leap-frog size was adjusted adaptively
 - Learning rate was annealed at every 2,000 iterations
 - Momentum method was used to accelerate learning
 - Context-dependent logarithmic GV w/o silence was used
- * Test sentences; 70 sentences
 - Paired comparison test, # of subjects 7 (native English speaker)
 - 30 sentences per subject

Experimental Results

Paired comparison test result

Baseline	PoE	No preference
17.1	32.4	50.5

Baseline; conventional (not jointly estimated) GV

PoE; proposed (jointly estimated) GV

Difference was statistically significant at $p < 0.05$ level

Summary

Statistical parametric synthesis based on PoE

- **Combination of multiple-level AMs is formulated as PoE**
- **Jointly estimate multiple-level AMs as PoE**
 - * Linear feature function with Gaussian experts
 - Can be estimated in the same way as trajectory HMM
 - * Non-linear feature function and/or non-Gaussian experts
 - Contrastive divergence learning
- **Experiments**
 - * Jointly estimating multiple AMs as PoE improved performance

References

- [Hinton;'02] G. Hinton, "Training product of experts by minimizing contrastive divergence," *Neural Computation*, vol. 14, no. 8, pp. 1771-1800, 2002.
- [Williams;'05] C. Williams, "How to pretend that correlated variables are independent by using difference observations," *Neural Computation*, vol. 17, no. 1, pp. 1--6, 2005.
- [Zen;'07] H. Zen, et al., "Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic feature vector sequences," *Computer Speech & Language*, vol. 21, no. 1, pp.153-173, 2007.
- [Latorre;'08] J. Latorre & M. Akamine, "Multilevel parametric-base F0 model for speech synthesis," *Proc. Interspeech*, pp. 2274--2277, 2008.
- [Qian;'09] Y. Qian, et al., "Improved prosody generation by maximizing joint likelihood of state and longer units," *Proc. ICASSP*, pp. 1173--1176, 2009.
- [Wang;'08] C.-C. Wang, et al., "Multi-layer F0 modeling for HMM-based speech synthesis," *Proc. ISCSLP*, pp. 129--132, 2008.
- [Ling;'06] Z.-H. Ling, et al., "USTC system for Blizzard Challenge 2006 an improved HMM-based speech synthesis method," *Proc. Blizzard Challenge Workshop*, 2006.
- [Gao;'08] B.-H. Gao, et al., "Duration refinement by jointly optimizing state and longer unit likelihood," *Proc. Interspeech*, 2266--2269, 2008.
- [Toda;'07] T. Toda & K. Tokuda, "A speech parameter generation algorithm considering global variance for HMM-based speech synthesis," *IEICE Trans. Inf. Syst.*, vol. E90-D, no. 5, pp. 816--824, 2007.