

MLMI14: Speech Synthesis

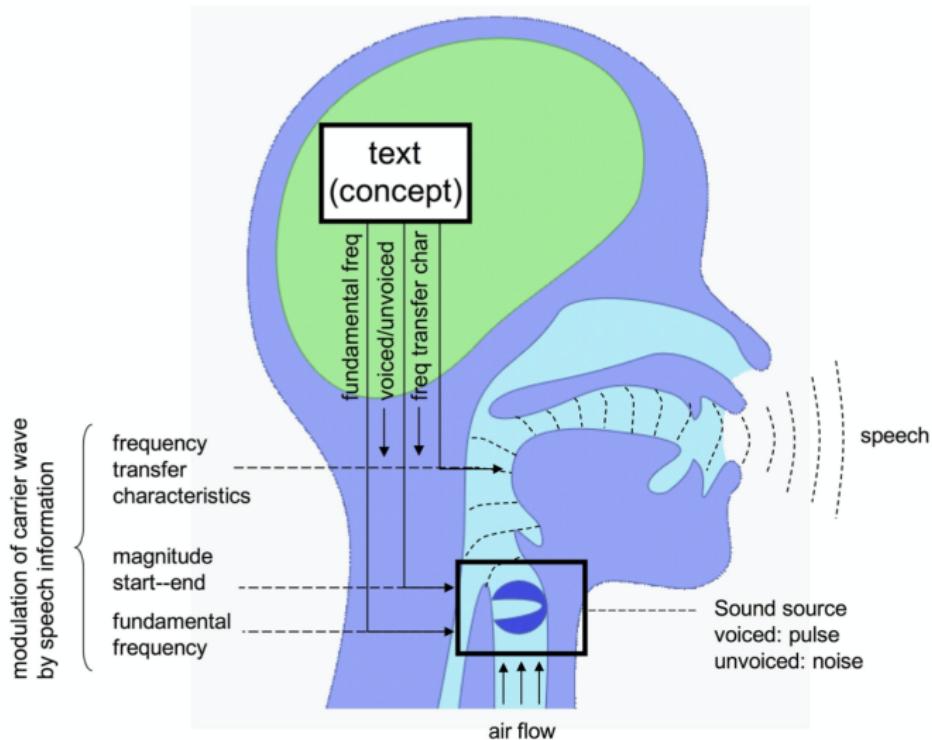
Mark Gales

Lent 2021

Sequence Mapping Process

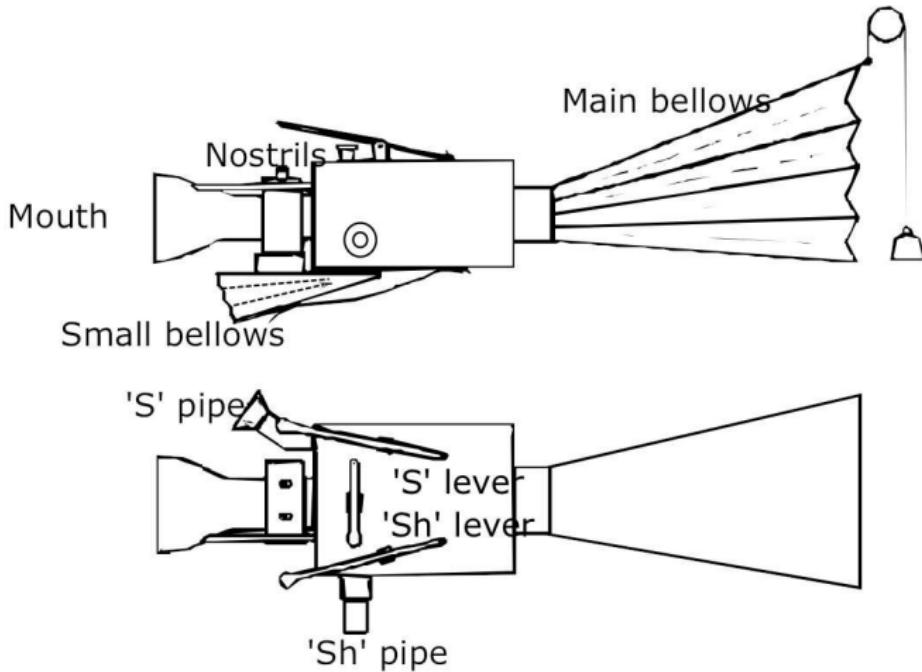
- For speech and language processing need to handle **variable length sequences**
 - the lengths of the source and target sequences not necessarily the same
 - number of possible sequences vast
 - often sequence-to-sequence mappings required
- Automatic Speech Recognition (ASR)
Speech (continuous) → **Text** (discrete)
- Machine Translation (MT)
Text (discrete) → **Text** (discrete)
- Text-to-Speech Synthesis (TTS)
Text (discrete) → **Speech** (continuous)

Speech Production

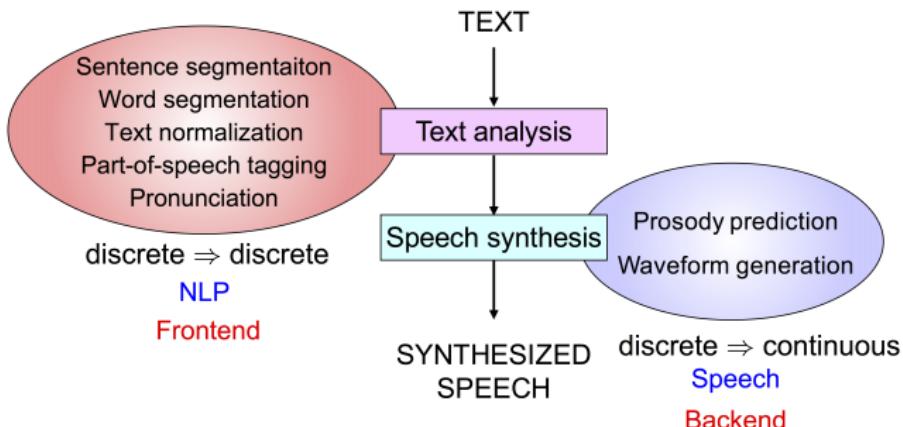


Speech Synthesis: von Kempelen (1791)

- Long history of development:



Speech Synthesis Processing Pipeline



- Typical processing pipeline: combines
 - natural language processing (NLP)
 - signal processing
 - machine learning (for statistical speech synthesis)

- Not as trivial as you might imagine ...
 - **abbreviation expansion**: standard abbreviations
 - **numbers**: dates, weights etc
 - **specialist vocabulary**: chemical symbols, texting ...
- Example mappings

St. → Street

St. → Saint

Mr. → Mister

£1984 → one thousand, nine hundred and eighty four pounds

1984 → nineteen eighty four

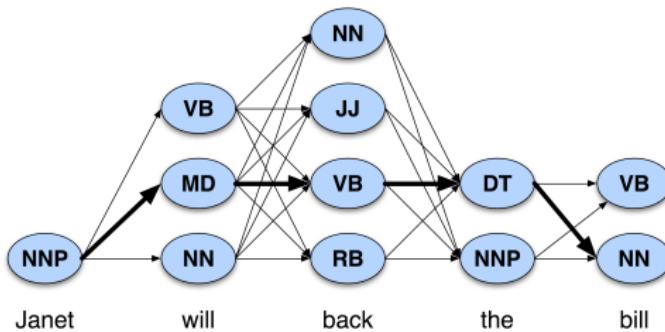
1984 → one thousand nine hundred and eighty four

AT&T → A. T. and T.

- Normalisation approach needs to be **context dependent**

Part of Speech (PoS) Tagger

- Given a word-sequence, $w_{1:L}$, extract the PoS tags $t_{1:L}$
Janet/NNP will/MD back/VB the/DT bill/NN
 - however individual words are ambiguous (from Jurafsky & Martin)



- pronunciation (and intonation) impacted by PoS

Penn TreeBank PoS Tags

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	‘ or “
POS	possessive ending	<i>'s</i>	”	right quote	’ or ”
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	[, (, {, <
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis],), }, >
RB	adverb	<i>quickly, never</i>	,	comma	,
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	! ?
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	: ; ... - -
RP	particle	<i>up, off</i>			

Pronunciation Lexicon

- Manual lexicon used for TTS, but typically more detailed
 - need to select from multiple pronunciations (use e.g. PoS)
 - only **one** pronunciation can be synthesised!
- Example lexicon entries (from Festival)

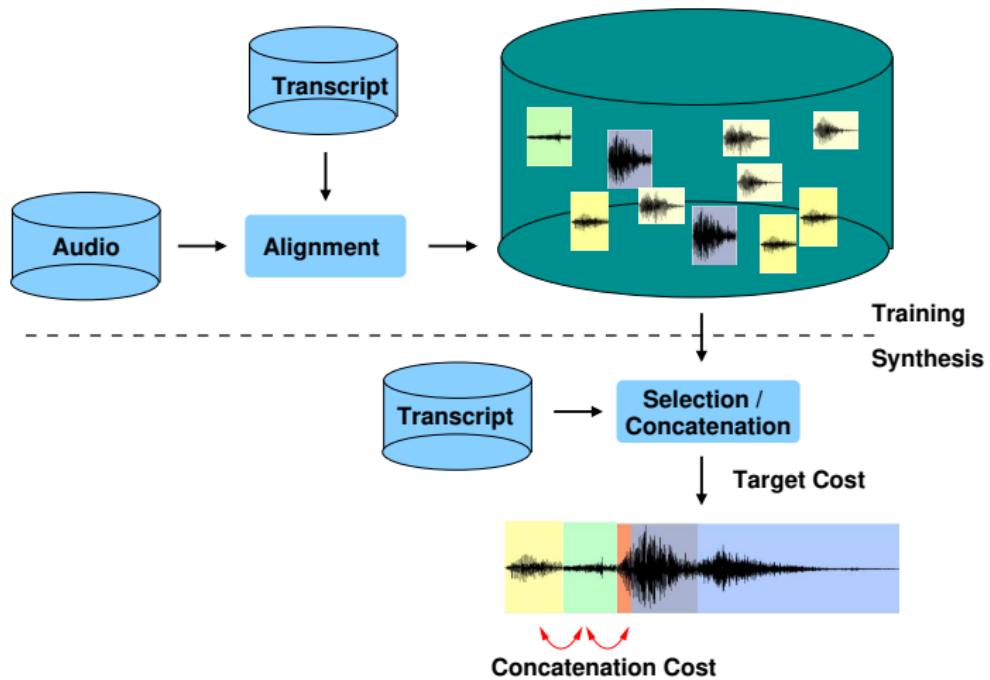
```
( "walkers" n ((( w oo ) 1) (( k @ z ) 0)) )  
( "present" v ((( p r e ) 0) (( z @ n t ) 1)) )  
( "lives" n ((( l ai v z ) 1)) )  
( "lives" v ((( l i v z ) 1)) )
```

- phone sequence (as in ASR)
- syllable boundary information
- stress-markers

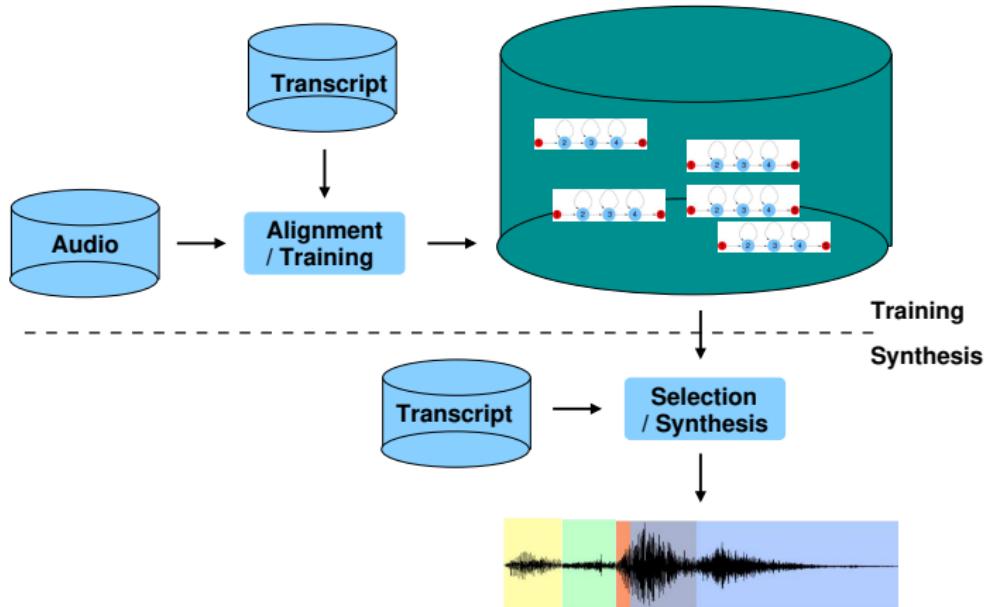
Text Analysis to Speech Synthesis

- Typical output of the text analysis stage is discrete:
 - set of phones
 - syllable positions
 - stress markers
 - word identity and part-of-speech
- These are passed to the synthesis system for:
 - aligning speech data with text
 - unit extraction/selection: **concatenative** synthesis
 - decision tree construction acoustic model training: **statistical** synthesis
- For concatenative systems additional information from text analysis
 - duration: target duration for each unit
 - F0/power: intonation pattern for target utterance

Concatenative Speech Synthesis



Statistical Speech Synthesis



Concatenative vs Statistical

- Concatenative Speech Synthesis:
 - high quality segmental accuracy
 - need large corpus of segments to cover target domain
- Statistical Speech Synthesis
 - (potentially very) small footprint
 - flexibility (able to control speaker characteristics)
 - issues with segmental naturalness
- Hybrid Approach: combining attributes of both schemes
 - use statistical approaches for unit selection
 - combine synthesis units with extracted units

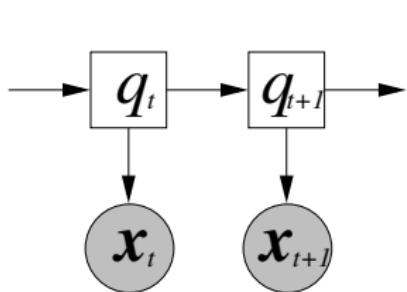
Concatenative vs Statistical

- Concatenative Speech Synthesis:
 - high quality segmental accuracy
 - need large corpus of segments to cover target domain
- Statistical Speech Synthesis
 - (potentially very) small footprint
 - flexibility (able to control speaker characteristics)
 - issues with segmental naturalness
- Hybrid Approach: combining attributes of both schemes
 - use statistical approaches for unit selection
 - combine synthesis units with extracted units
- Then shift to purely deep-learning style approaches: WaveNet
 - removed need for signal processing (and natural language processing)

HMM Trajectory Generation

Hidden Markov Models

- An important sequence model **hidden Markov model** (HMM)
 - an example of a **dynamic Bayesian network** (DBN)
 - consider a sequence of observations $\mathbf{x}_1, \dots, \mathbf{x}_T$



- discrete **latent variables**
 - q_t describes discrete state-space
 - conditional independence assumptions

$$P(q_t|q_0, \dots, q_{t-1}) = P(q_t|q_{t-1})$$

$$p(\mathbf{x}_t|\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, q_0, \dots, q_t) = p(\mathbf{x}_t|q_t)$$

- The likelihood of the data is

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{\mathbf{q} \in \mathbf{Q}_T} P(q_0) \prod_{t=1}^T P(q_t|q_{t-1}) p(\mathbf{x}_t|q_t)$$

$\mathbf{q} = \{q_0, \dots, q_T\}$ and \mathbf{Q}_T is all T -length state sequences

- Maximum Likelihood (ML) training discussed for ASR
 - Expectation-Maximisation (EM) often used

$$\boldsymbol{\mu}_m = \frac{\sum_{t=1}^T \gamma_m(t) \mathbf{x}_t}{\sum_{t=1}^T \gamma_m(t)} \quad \boldsymbol{\Sigma}_m = \frac{\sum_{t=1}^T \gamma_m(t) (\mathbf{x}_t - \boldsymbol{\mu}_m)(\mathbf{x}_t - \boldsymbol{\mu}_m)^T}{\sum_{t=1}^T \gamma_m(t)}$$

- $\gamma_m(t) = P(q_t = m | \mathbf{x}_1, \dots, \mathbf{x}_T, \boldsymbol{\lambda})$: $\boldsymbol{\lambda}$ “current” parameters
- Feature vector often comprises static/delta/delta-delta

$$\begin{bmatrix} \mathbf{x}_t \\ \Delta \mathbf{x}_t \\ \Delta^2 \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} - \mathbf{x}_{t-1} \\ \Delta \mathbf{x}_{t+1} - \Delta \mathbf{x}_{t-1} \end{bmatrix}$$

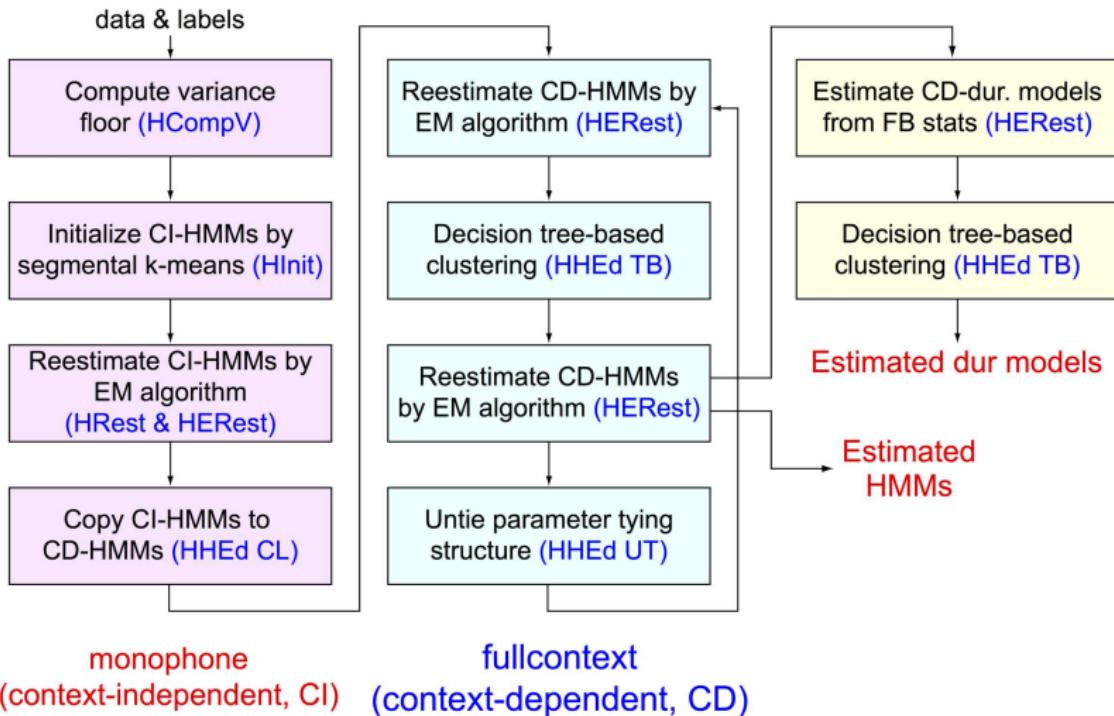
- normalisation terms for delta and deltas ignored

- Configurations for synthesis similar to ASR - major differences
 - (often) single Gaussian component per state
 - rich decision tree questions and deep trees
 - state-position roots of decision trees
 - multiple streams
 - duration model
 - multi-space distributions
- Baseline feature vector - 5ms frame-rate, 25ms window-size

$$\begin{bmatrix} \text{mel-cepstrum} - 40 - 60 \text{ dim} \\ \text{log-F0} - 1 \text{ dim} \end{bmatrix}$$

- delta and delta-delta parameters can be added as usual

Synthesis HMM Training Process



Multiple Stream Systems

- Multiple stream states are standard features in HTK
 - rarely used for speech recognition
 - used for statistical speech synthesis with HMMs
- Two multi-stream options - each has **stream weight** w_s
 - multiple component streams

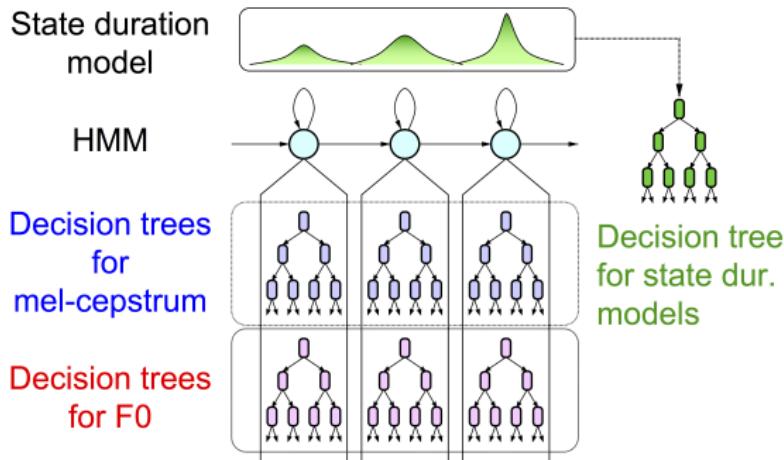
$$p\left(\left[\begin{array}{ccc} \mathbf{x}_t^{(1)\top} & \dots & \mathbf{x}_t^{(S)\top} \end{array}\right]^\top \middle| q_t\right) \propto \prod_{s=1}^S \left(\sum_{m=1}^{M_s} c_{q_t}^{(sm)} \mathcal{N}(\mathbf{x}_t^{(s)}; \boldsymbol{\mu}_{q_t}^{(sm)}, \boldsymbol{\Sigma}_{q_t}^{(sm)}) \right)^{w_s}$$

- multiple decision trees

$$p\left(\left[\begin{array}{ccc} \mathbf{x}_t^{(1)\top} & \dots & \mathbf{x}_t^{(S)\top} \end{array}\right]^\top \middle| q_t^{(1)}, \dots, q_t^{(S)}\right) \propto \prod_{s=1}^S \left(\mathcal{N}(\mathbf{x}_t^{(s)}; \boldsymbol{\mu}_{q_t^{(s)}}^{(s)}, \boldsymbol{\Sigma}_{q_t^{(s)}}^{(s)}) \right)^{w_s}$$

where $q_t^{(s)}$ indicates the context for stream s at time t

Multiple Decision Trees for Synthesis



- Root of decision tree is a state position (phone-state in ASR)
 - separate decision trees for spectral parameters/(log)F0
 - separate trees also used for delta and delta-delta (log)F0
 - separate decision tree for duration

- The standard duration (τ) model for HMMs is an exponential decay

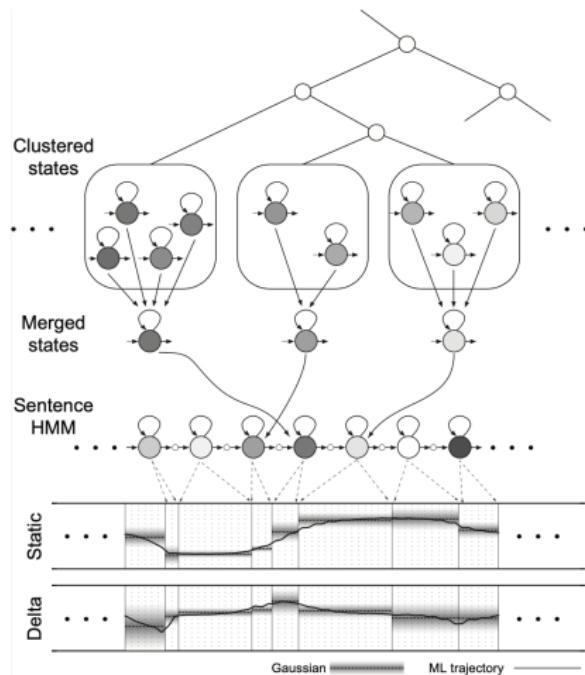
$$P(\tau|s_i) = a_{ii}^{\tau-1}(1 - a_{ii})$$

- not a realistic model for actual durations
- efficient as Viterbi can be directly applied
- important for synthesis to get duration correct
- Hidden Semi Markov Models: explicit duration model
 - for synthesis single Gaussian used for duration

$$p(\tau|s_i) = \mathcal{N}(\tau; \mu_i^{\text{dur}}, \Sigma_i^{\text{dur}})$$

- can be trained on final iteration, or integrated in training

Synthesis Trajectory Generation

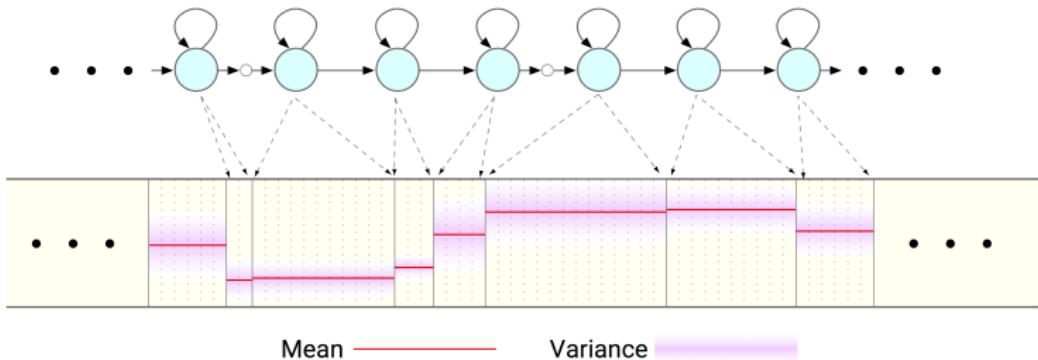


- Text analysis yield a sequence of phones and attributes
 - use decision tree to obtain sentence state sequence s_1, \dots, s_L
 - state sequence will be different for each stream (decision tree)
- First stage is to generate the duration for each of the L states
 - simplest: take the mean of the Gaussian for state s_i : μ_i^{dur}
 - total sentence length is

$$T = \sum_{i=1}^L \mu_i^{\text{dur}}$$

- Yields T length vector, q , indicating the state at each time
 - q can now used to generate the parameter trajectories

HMM Trajectory - Standard Model



- Given the sequence of T states \mathbf{q} associated with a word sequence
 - need to generate distribution over the trajectory
 - assume that observation vector \mathbf{x} is just the cepstra then

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{q}) = \prod_{t=1}^T p(\mathbf{x}_t | q_t) = \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t})$$

- Given the distribution $p(\mathbf{x}_{1:T}|\mathbf{q})$ need to generate a trajectory
 - standard solution select **maximum likelihood** estimate

$$\hat{\mathbf{x}}_{1:T} = \arg \max \{ p(\mathbf{x}_{1:T}|\mathbf{q}) \}$$

- Using the standard distribution yields

$$\hat{\mathbf{x}}_{1:T} = \arg \max \left\{ \prod_{t=1}^T \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) \right\} = [\boldsymbol{\mu}_{q_1}^\top \dots \boldsymbol{\mu}_{q_T}^\top]^\top$$

- most likely sequence simply the sequence of means

HMM Trajectory - Static/Delta

- Modify the “observations” to include delta parameters, $\Delta \mathbf{x}$,
 - possible to write for a single observation

$$\bar{\mathbf{x}}_t = \begin{bmatrix} \mathbf{x}_t \\ \Delta \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} 0 & +\mathbf{I} & 0 \\ -\mathbf{I} & 0 & +\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{t-1} \\ \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{bmatrix}$$

$$p\left(\begin{bmatrix} \mathbf{x}_t \\ \Delta \mathbf{x}_t \end{bmatrix} \middle| q_t\right) = \mathcal{N}\left(\bar{\mathbf{x}}_t; \bar{\boldsymbol{\mu}}_{q_t}, \bar{\boldsymbol{\Sigma}}_{q_t}\right)$$

- and for the complete sequence

$$\begin{bmatrix} \vdots \\ \mathbf{x}_t \\ \Delta \mathbf{x}_t \\ \mathbf{x}_{t+1} \\ \Delta \mathbf{x}_{t+1} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots & \dots \\ \dots & \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & -\mathbf{I} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & -\mathbf{I} & \mathbf{0} & \mathbf{I} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \mathbf{x}_{t-1} \\ \mathbf{x}_t \\ \mathbf{x}_{t+1} \\ \mathbf{x}_{t+2} \\ \vdots \end{bmatrix} = \mathbf{W} \mathbf{x}_{1:T}$$

- need the **distribution** of $\mathbf{x}_{1:T}$

HMM Trajectory (cont)

- Possible to write

$$p(\mathbf{Wx}_{1:T} | \mathbf{q}) = \mathcal{N}(\mathbf{Wx}_{1:T}; \bar{\mu}_\mathbf{q}, \bar{\Sigma}_\mathbf{q},)$$

where

$$\bar{\mu}_\mathbf{q} = \begin{bmatrix} \vdots \\ \bar{\mu}_{q_t} \\ \bar{\mu}_{q_{t+1}} \\ \vdots \end{bmatrix}, \bar{\Sigma}_\mathbf{q} = \begin{bmatrix} \dots & \vdots & \vdots & \vdots & \dots \\ \dots & \mathbf{0} & \bar{\Sigma}_{q_t} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \bar{\Sigma}_{q_{t+1}} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

- but we need a distribution for $\mathbf{x}_{1:T}$ not $\mathbf{Wx}_{1:T}$
- Note: an example of a **product of experts** for details see
 - H. Zen, M. J. F. Gales, Y. Nankaku, and K. Tokuda, "Product of experts for statistical parametric speech synthesis," IEEE Trans. Audio, Speech & Language Processing, 2012

HMM Trajectory (cont)

- Expanding out the form of distribution

$$\begin{aligned} p(\mathbf{Wx}_{1:T}|\boldsymbol{q}) &\propto \exp\left(-\frac{1}{2}(\mathbf{Wx}_{1:T} - \bar{\boldsymbol{\mu}}_{\boldsymbol{q}})^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} (\mathbf{Wx}_{1:T} - \bar{\boldsymbol{\mu}}_{\boldsymbol{q}})\right) \\ &\propto \exp\left(-\frac{1}{2}\left(\mathbf{x}_{1:T}^T \mathbf{W}^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} \mathbf{W} \mathbf{x}_{1:T} - 2\bar{\boldsymbol{\mu}}_{\boldsymbol{q}}^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} \mathbf{W} \mathbf{x}_{1:T}\right)\right) \end{aligned}$$

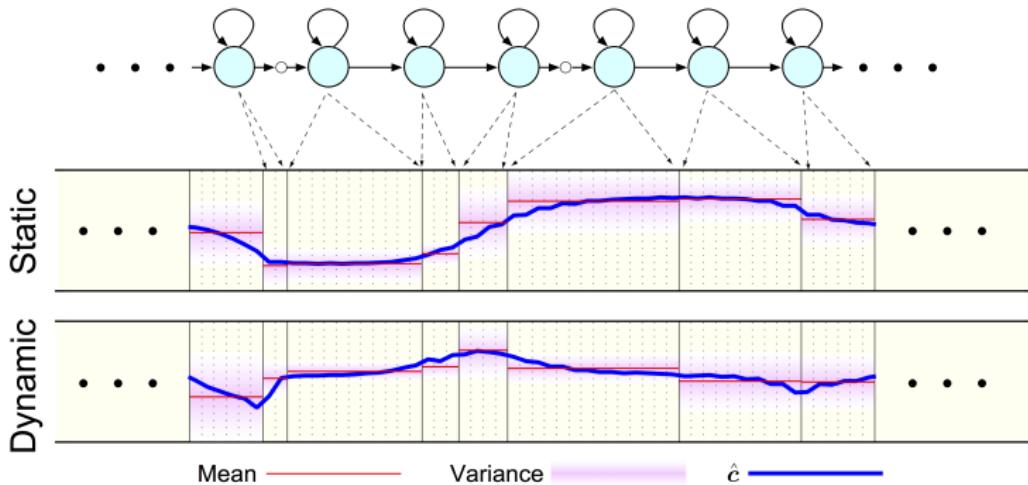
- Equating co-efficients yields

$$\begin{aligned} p(\mathbf{x}_{1:T}|\boldsymbol{q}) &= \mathcal{N}\left(\mathbf{x}_{1:T}, \left(\mathbf{W}^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} \mathbf{W}\right)^{-1} \mathbf{W}^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} \bar{\boldsymbol{\mu}}_{\boldsymbol{q}}, \left(\mathbf{W}^T \bar{\boldsymbol{\Sigma}}_{\boldsymbol{q}}^{-1} \mathbf{W}\right)^{-1}\right) \\ &= \mathcal{N}(\mathbf{x}_{1:T}; \boldsymbol{\mu}_{\boldsymbol{q}}, \boldsymbol{\Sigma}_{\boldsymbol{q}}) \end{aligned}$$

- Again estimate the ML trajectory

$$\hat{\mathbf{x}}_{1:T} = \boldsymbol{\mu}_{\boldsymbol{q}}$$

HMM Trajectory (cont)



- Final covariance matrix **not diagonal**
 - introduces correlation over time - a smooth trajectory
- Final distribution over the complete sequence is **Gaussian**
 - able to generate a distribution for **any** sequence length

Multiple Component Acoustic Models

- Previous estimation has assumed a single Gaussian per state
 - possible to generalise to M components per state

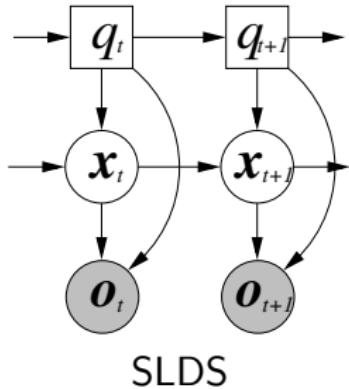
$$p(\bar{\mathbf{x}}_t | q_t) = \sum_{m=1}^M c_m \mathcal{N}(\bar{\mathbf{x}}_t, \bar{\boldsymbol{\mu}}_{q_t}^{(m)}, \bar{\boldsymbol{\Sigma}}_{q_t}^{(m)})$$

- Again consider **ML** estimation of trajectory

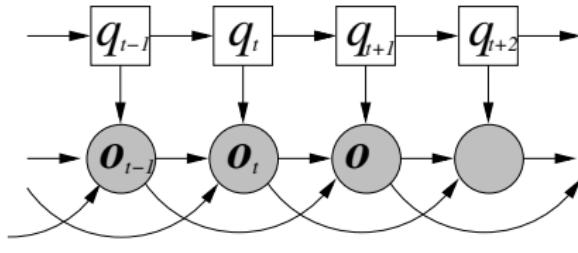
$$\hat{\mathbf{x}}_{1:T} = \arg \max \{ p(\mathbf{x}_{1:T} | \mathbf{q}, \mathbf{m}) \}$$

- maximisation is over $\mathbf{x}_{1:T}$ and \mathbf{m}
- \mathbf{m} is T dimensional vector of component at each time instance
- Computationally expensive (and limited performance gains)

Alternative Acoustic Models



SLDS



AR-HMM

- Acoustic models also investigated for synthesis:
 - **Switching Linear Dynamical Systems:**
 - additional continuous latent variable
 - **Auto-Regressive HMM:**
 - additional dependencies on observations
- Both yield a smoother trajectory (without dynamics)

SLDS Speech Trajectory Modelling

Frames from phrase:
SHOW THE GRIDLEY'S...

Legend

- True
- HMM
- SLDS

