



Acoustic Modeling for Speech Synthesis

Heiga Zen

Dec. 14th, 2015@ASRU

Outline

Background

HMM-based acoustic modeling

- Training & synthesis

- Limitations

ANN-based acoustic modeling

- Feedforward NN

- RNN

Conclusion



Outline

Background

HMM-based acoustic modeling

- Training & synthesis
- Limitations

ANN-based acoustic modeling

- Feedforward NN
- RNN

Conclusion



Text-to-speech as sequence-to-sequence mapping

Automatic speech recognition (ASR)

Speech (real-valued time series) \rightarrow Text (discrete symbol sequence)



Text-to-speech as sequence-to-sequence mapping

Automatic speech recognition (ASR)

Speech (real-valued time series) \rightarrow Text (discrete symbol sequence)

Statistical machine translation (SMT)

Text (discrete symbol sequence) \rightarrow Text (discrete symbol sequence)



Text-to-speech as sequence-to-sequence mapping

Automatic speech recognition (ASR)

Speech (real-valued time series) \rightarrow Text (discrete symbol sequence)

Statistical machine translation (SMT)

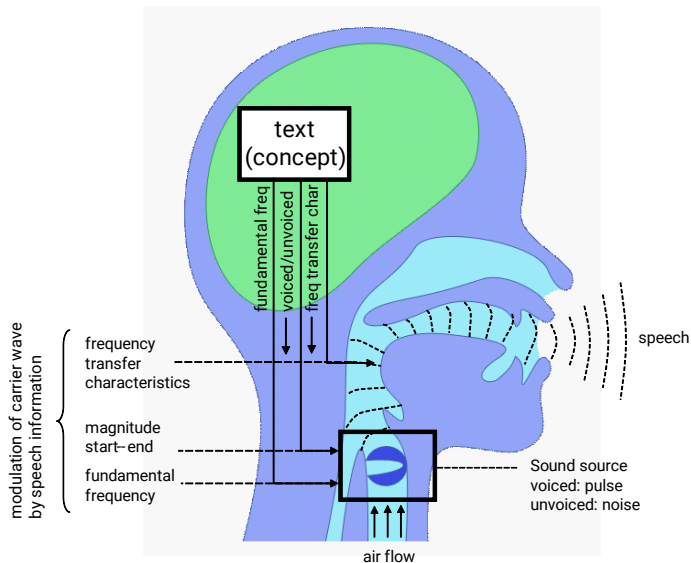
Text (discrete symbol sequence) \rightarrow Text (discrete symbol sequence)

Text-to-speech synthesis (TTS)

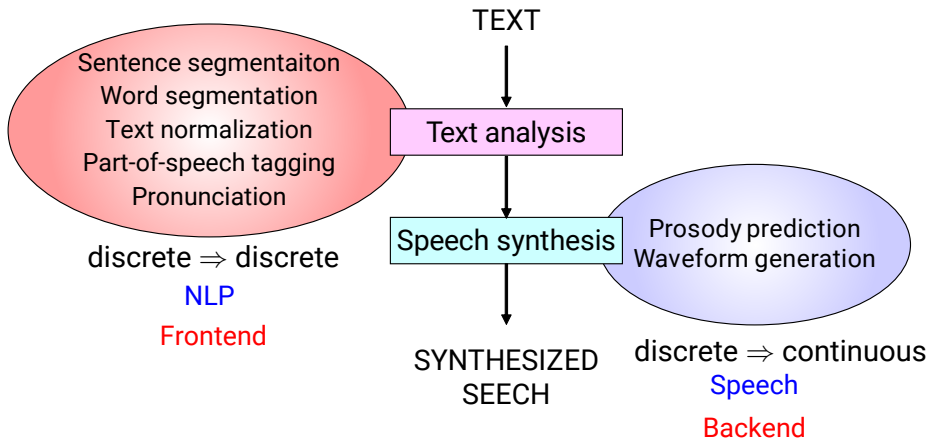
Text (discrete symbol sequence) \rightarrow Speech (real-valued time series)



Speech production process



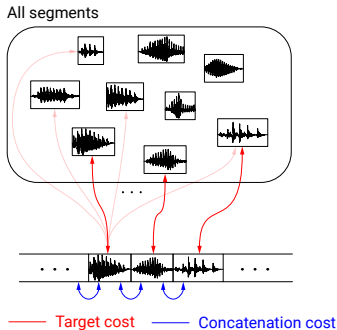
Typical flow of TTS system



This presentation mainly talks about backend



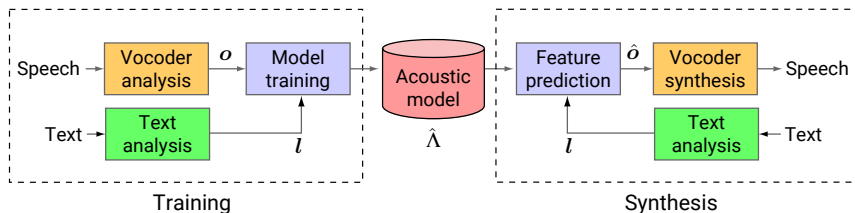
Concatenative speech synthesis



- Concatenate actual small speech segments from database
→ **Very high segmental naturalness**
- Single segment per unit (e.g., diphone) → diphone synthesis [1]
- Multiple segments per unit → unit selection synthesis [2]



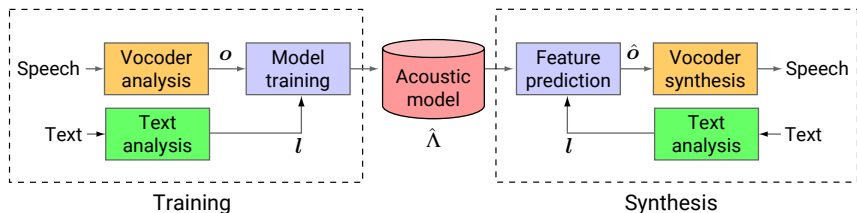
Statistical parametric speech synthesis (SPSS) [4]



- Parametric representation rather than waveform
- Model relationship between linguistic & acoustic features
- Predict acoustic features then reconstruct waveform



Statistical parametric speech synthesis (SPSS) [4]

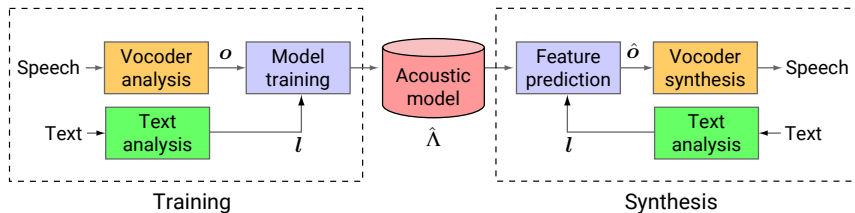


- Parametric representation rather than waveform
- Model relationship between linguistic & acoustic features
- Predict acoustic features then reconstruct waveform

SPSS can use any acoustic model, but HMM-based one is very popular
→ [HMM-based speech synthesis \[3\]](#)



Statistical parametric speech synthesis (SPSS) [4]



Pros

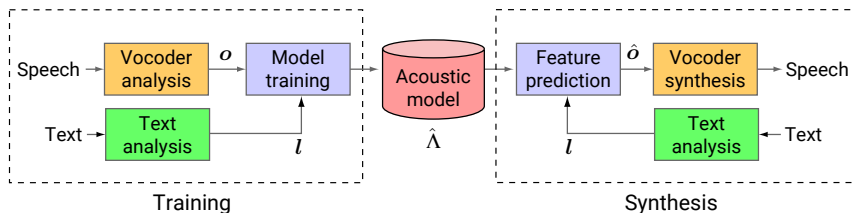
- Small footprint
- Flexibility to change voice characteristics
- Robust to data sparsity and noise/mistakes in data

Cons

- Segmental naturalness



Major factors for naturalness degradation



- **Vocoder analysis/synthesis**
 - *How to parameterize speech?*
- **Acoustic model**
 - *How to represent relationship between speech & text?*
- **Oversmoothing**
 - *How to generate speech from model?*



Outline

Background

HMM-based acoustic modeling

Training & synthesis

Limitations

ANN-based acoustic modeling

Feedforward NN

RNN

Conclusion



Formulation of SPSS

Training

- Extract linguistic features l & acoustic features o
- Train acoustic model Λ given (o, l)

$$\hat{\Lambda} = \arg \max_{\Lambda} p(o | l, \Lambda)$$



Formulation of SPSS

Training

- Extract linguistic features l & acoustic features o
- Train acoustic model Λ given (o, l)

$$\hat{\Lambda} = \arg \max_{\Lambda} p(o | l, \Lambda)$$

Synthesis

- Extract l from text to be synthesized
- Generate most probable o from $\hat{\Lambda}$ then reconstruct waveform

$$\hat{o} = \arg \max_o p(o | l, \hat{\Lambda})$$



Formulation of SPSS

Training

- Extract linguistic features l & acoustic features o
- Train acoustic model Λ given (o, l)

$$\hat{\Lambda} = \arg \max_{\Lambda} p(o | l, \Lambda)$$

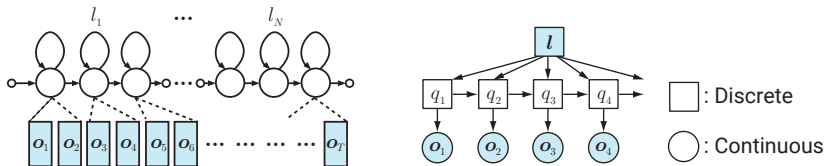
Synthesis

- Extract l from text to be synthesized
- Generate most probable o from $\hat{\Lambda}$ then reconstruct waveform


$$\hat{o} = \arg \max_o p(o | l, \hat{\Lambda})$$



Training – HMM-based acoustic modeling



$$\begin{aligned} p(\mathbf{o} | \mathbf{l}, \Lambda) &= \sum_{\forall \mathbf{q}} p(\mathbf{o} | \mathbf{q}, \Lambda) P(\mathbf{q} | \mathbf{l}, \Lambda) \quad \mathbf{q}: \text{hidden states} \\ &= \sum_{\forall \mathbf{q}} \prod_{t=1}^T p(\mathbf{o}_t | q_t, \Lambda) P(\mathbf{q} | \mathbf{l}, \Lambda) \quad q_t: \text{hidden state at } t \\ &= \sum_{\forall \mathbf{q}} \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) P(\mathbf{q} | \mathbf{l}, \Lambda) \end{aligned}$$

ML estimation of HMM parameters → Baum-Welch (EM) algorithm [5] 

Training – Linguistic features

Linguistic features: phonetic, grammatical, & prosodic features

- **Phoneme**

phoneme identity, position

- **Syllable**

length, accent, stress, tone, vowel, position

- **Word**

length, POS, grammar, prominence, emphasis, position, pitch accent

- **Phrase**

length, type, position, intonation

- **Sentence**

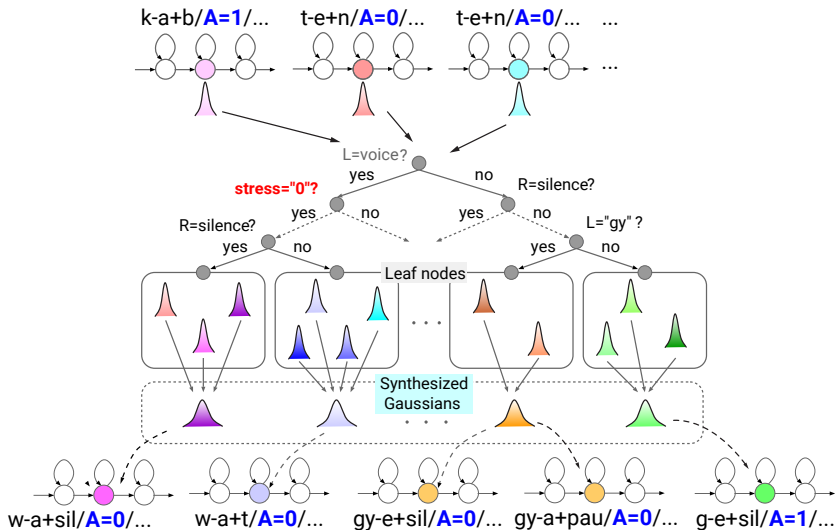
length, type, position

...

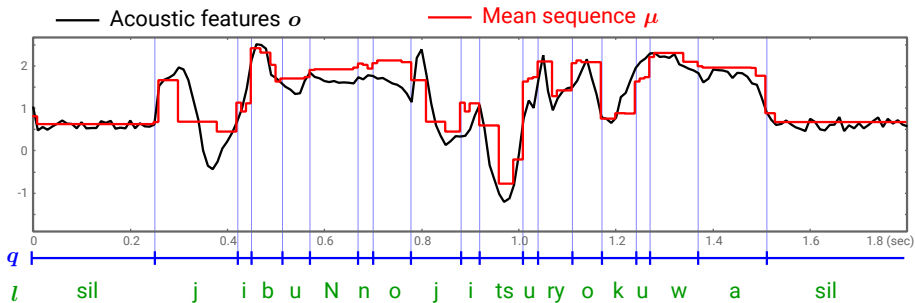
→ Impossible to have enough data to cover all combinations



Training – ML decision tree-based state clustering [6]



Training – Example



Formulation of SPSS

Training

- Extract linguistic features l & acoustic features o
- Train acoustic model Λ given (o, l)

$$\hat{\Lambda} = \arg \max_{\Lambda} p(o | l, \Lambda)$$

Synthesis

- Extract l from text to be synthesized
- Generate most probable o from $\hat{\Lambda}$ then reconstruct waveform

$$\hat{o} = \arg \max_o p(o | l, \hat{\Lambda})$$

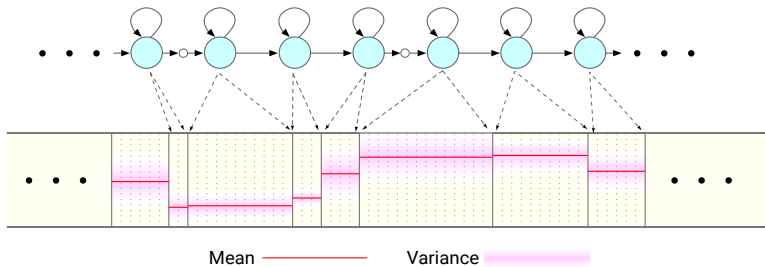


Synthesis – Predict most probable acoustic features

$$\begin{aligned}\hat{o} &= \arg \max_{\mathbf{o}} p(\mathbf{o} | \mathbf{l}, \hat{\Lambda}) \\ &= \arg \max_{\mathbf{o}} \sum_{\forall \mathbf{q}} p(\mathbf{o}, \mathbf{q} | \mathbf{l}, \hat{\Lambda}) \\ &\approx \arg \max_{\mathbf{o}} \max_{\mathbf{q}} p(\mathbf{o}, \mathbf{q} | \mathbf{l}, \hat{\Lambda}) \\ &= \arg \max_{\mathbf{o}} \max_{\mathbf{q}} p(\mathbf{o} | \mathbf{q}, \hat{\Lambda}) P(\mathbf{q} | \mathbf{l}, \hat{\Lambda}) \\ &\approx \arg \max_{\mathbf{o}} p(\mathbf{o} | \hat{\mathbf{q}}, \hat{\Lambda}) \quad s.t. \quad \hat{\mathbf{q}} = \arg \max_{\mathbf{q}} P(\mathbf{q} | \mathbf{l}, \hat{\Lambda}) \\ &= \arg \max_{\mathbf{o}} \mathcal{N}(\mathbf{o}; \boldsymbol{\mu}_{\hat{\mathbf{q}}}, \boldsymbol{\Sigma}_{\hat{\mathbf{q}}}) \\ &= \boldsymbol{\mu}_{\hat{\mathbf{q}}} \\ &= \left[\boldsymbol{\mu}_{\hat{q}_1}^\top, \dots, \boldsymbol{\mu}_{\hat{q}_T}^\top \right]^\top\end{aligned}$$



Synthesis – Most probable acoustic features given HMM



$\hat{o} \rightarrow$ step-wise \rightarrow discontinuity can be perceived



Synthesis – Using dynamic feature constraints [7]

$$o_t = \begin{bmatrix} c_t^\top & \Delta c_t^\top \end{bmatrix}^\top$$

$\left(\begin{array}{c} \text{blue box} \\ \text{red box} \end{array} \right)_{2M}$

 $\left(\text{blue box} \right)_M$

 $\left(\text{red box} \right)_M$

$$\Delta c_t = c_t - c_{t-1}$$

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



Synthesis – Speech parameter generation algorithm [7]

$$\hat{o} = \arg \max_{\mathbf{o}} p(\mathbf{o} | \hat{\mathbf{q}}, \hat{\Lambda}) \quad s.t. \quad \mathbf{o} = \mathbf{W}\mathbf{c}$$

$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\hat{\mathbf{q}}}, \boldsymbol{\Sigma}_{\hat{\mathbf{q}}})$$

$$= \arg \max_{\mathbf{c}} \log \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\hat{\mathbf{q}}}, \boldsymbol{\Sigma}_{\hat{\mathbf{q}}})$$



Synthesis – Speech parameter generation algorithm [7]

$$\hat{o} = \arg \max_{\mathbf{o}} p(\mathbf{o} | \hat{q}, \hat{\Lambda}) \quad s.t. \quad \mathbf{o} = \mathbf{W} \mathbf{c}$$

$$\hat{c} = \arg \max_{\mathbf{c}} \mathcal{N}(\mathbf{W} \mathbf{c}; \mu_{\hat{q}}, \Sigma_{\hat{q}})$$

$$= \arg \max_{\mathbf{c}} \log \mathcal{N}(\mathbf{W} \mathbf{c}; \mu_{\hat{q}}, \Sigma_{\hat{q}})$$

$$\frac{\partial}{\partial \mathbf{c}} \log \mathcal{N}(\mathbf{W} \mathbf{c}; \mu_{\hat{q}}, \Sigma_{\hat{q}}) \propto \mathbf{W}^{\top} \Sigma_{\hat{q}}^{-1} \mathbf{W} \mathbf{c} - \mathbf{W}^{\top} \Sigma_{\hat{q}}^{-1} \mu_{\hat{q}}$$

$$\mathbf{W}^{\top} \Sigma_{\hat{q}}^{-1} \mathbf{W} \mathbf{c} = \mathbf{W}^{\top} \Sigma_{\hat{q}}^{-1} \mu_{\hat{q}}$$

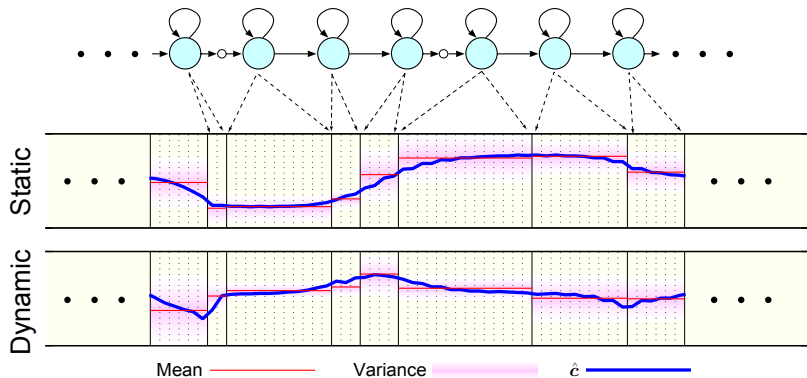
where

$$\mu_{\mathbf{q}} = \left[\mu_{q_1}^{\top}, \mu_{q_2}^{\top}, \dots, \mu_{q_T}^{\top} \right]^{\top}$$

$$\Sigma_{\mathbf{q}} = \text{diag} [\Sigma_{q_1}, \Sigma_{q_2}, \dots, \Sigma_{q_T}]$$

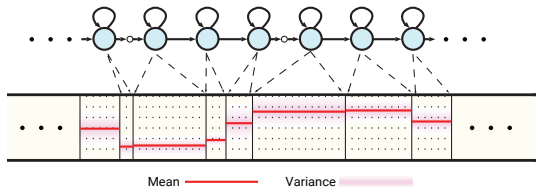
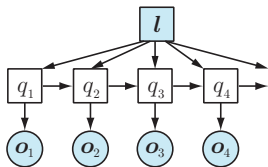


Synthesis – Most probable acoustic features under constraints between static & dynamic features



HMM-based acoustic model – Limitations (1)

Stepwise statistics

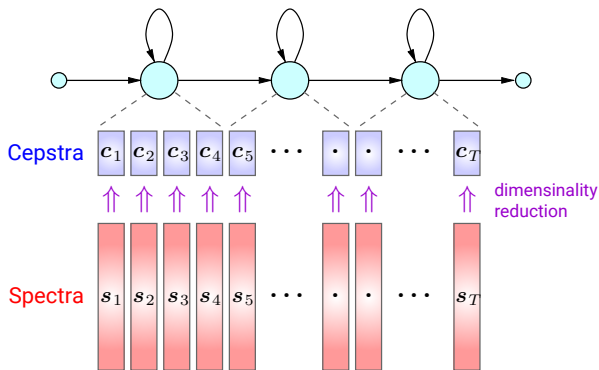


- Output probability only depends on the current state
- Within the same state, statistics are constant
→ **Step-wise statistics**
- Using dynamic feature constraints
→ **Ad hoc & introduces inconsistency betw. training & synthesis [8]**



HMM-based acoustic model – Limitations (2)

Difficulty to integrate feature extraction & modeling

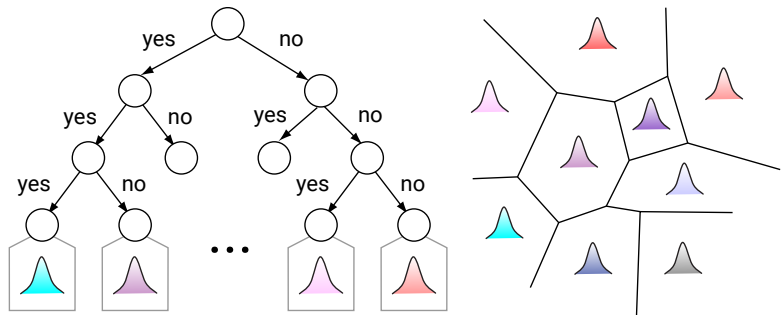


- Spectra or waveforms are high-dimensional & highly correlated
- Hard to be modeled by HMMs with Gaussian + diagonal covariance
→ Use low dimensional approximation (e.g., cepstra, LSPs)



HMM-based acoustic model – Limitations (3)

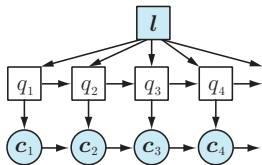
Data fragmentation



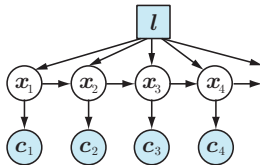
- Trees split input into clusters & put representative distributions
→ **Inefficient to represent dependency betw. ling. & acoust. feats.**
- Minor features are never used (e.g., word-level emphasis [9])
→ **Little or no effect**



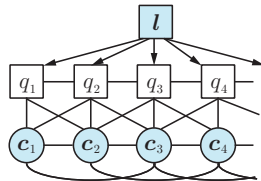
Alternatives – Stepwise statistics



ARHMM



LDM



Trajectory HMM

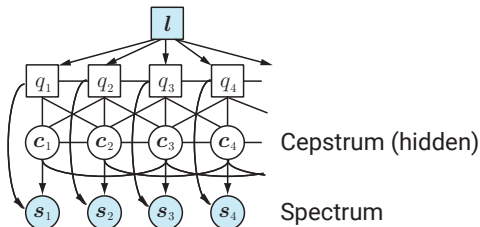
- Autoregressive HMMs (ARHMMs) [10]
- Linear dynamical models (LDMs) [11, 12]
- Trajectory HMMs [8]
- ...

Most of them use clustering → **Data fragmentation**

Often employ trees from HMM → **Sub-optimal**



Alternatives – Difficulty to integrate feature extraction



- Statistical vocoder [13]
- Minimum generation error with log spectral distortion [14]
- Waveform-level model [15]
- Mel-cepstral analysis-integrated HMM [16]

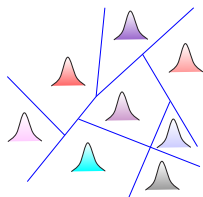
Use clustering to build tying structure → **Data fragmentation**

Often employ trees from HMM → **Sub-optimal**

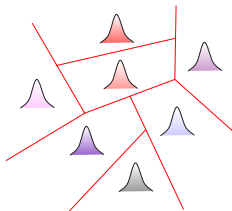


Alternatives – Data fragmentation

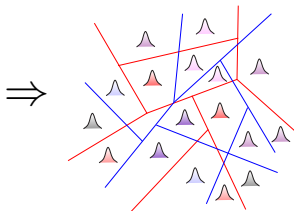
Tree1 (8 classes)



Tree2 (7 classes)



Combined (17 classes)



- Factorized decision tree [9, 17]
- Product of experts [18]

Each tree/expert still has data fragmentation → **Data fragmentation**
Fix other trees while building one tree [19, 20] → **Sub-optimal**



Outline

Background

HMM-based acoustic modeling

- Training & synthesis
- Limitations

ANN-based acoustic modeling

- Feedforward NN
- RNN

Conclusion



Linguistic → Acoustic mapping

- **Training**

Learn relationship between linguistic & acoustic features



Linguistic → Acoustic mapping

- **Training**

Learn relationship between linguistic & acoustic features

- **Synthesis**

Map linguistic features to acoustic ones



Linguistic → Acoustic mapping

- **Training**
Learn relationship between linguistic & acoustic features
- **Synthesis**
Map linguistic features to acoustic ones
- **Linguistic features used in SPSS**
 - Phoneme, syllable, word, phrase, utterance-level features
 - Around 50 different types
 - Sparse & correlated

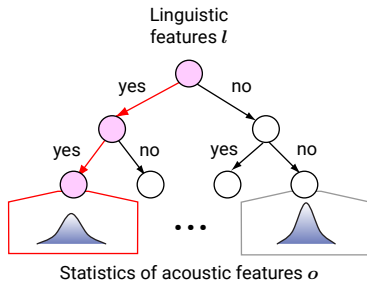
Effective modeling is essential



Decision tree-based acoustic model

HMM-based acoustic model & alternatives

→ Actually decision tree-based acoustic model



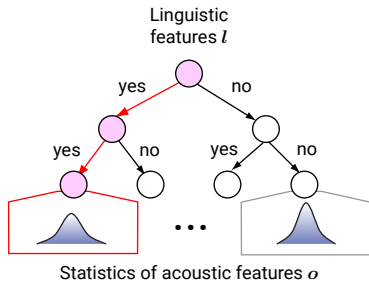
Regression tree: linguistic features → Stats. of acoustic features



Decision tree-based acoustic model

HMM-based acoustic model & alternatives

→ Actually decision tree-based acoustic model



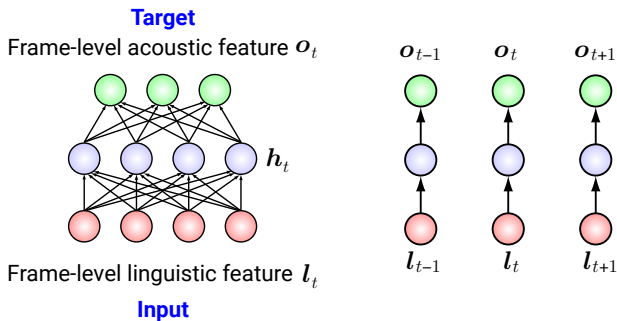
Regression tree: linguistic features → Stats. of acoustic features

Replace the tree with a general-purpose regression model

→ **Artificial neural network**



ANN-based acoustic model [21] – Overview



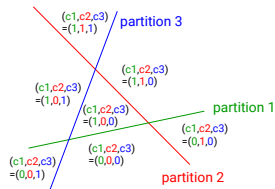
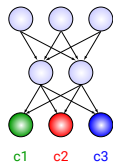
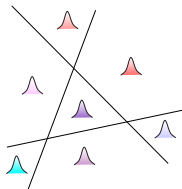
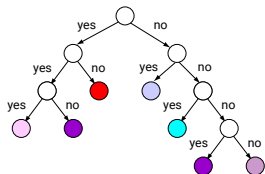
$$\mathbf{h}_t = f(\mathbf{W}_{hl}\mathbf{l}_t + \mathbf{b}_h) \quad \hat{\mathbf{o}}_t = \mathbf{W}_{oh}\mathbf{h}_t + \mathbf{b}_o$$
$$\hat{\Lambda} = \arg \min_{\Lambda} \sum_t \|\mathbf{o}_t - \hat{\mathbf{o}}_t\|_2 \quad \Lambda = \{\mathbf{W}_{hl}, \mathbf{W}_{oh}, \mathbf{b}_h, \mathbf{b}_o\}$$

$\hat{\mathbf{o}}_t \approx \mathbb{E}[\mathbf{o}_t | \mathbf{l}_t] \rightarrow$ Replace decision trees & Gaussian distributions



ANN-based acoustic model [21] – Motivation (1)

Distributed representation [22, 23]

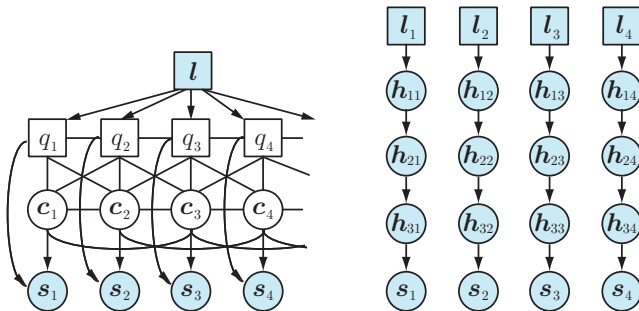


- **Fragmented:** n terminal nodes $\rightarrow n$ classes (linear)
- **Distributed:** n binary units $\rightarrow 2^n$ classes (exponential)
- **Minor features** (e.g., word-level emphasis) can affect synthesis



ANN-based acoustic model [21] – Motivation (2)

Integrate feature extraction [24, 25, 26]

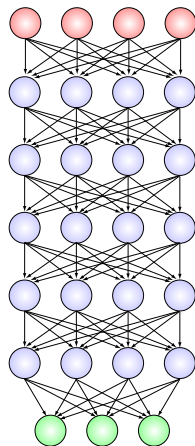
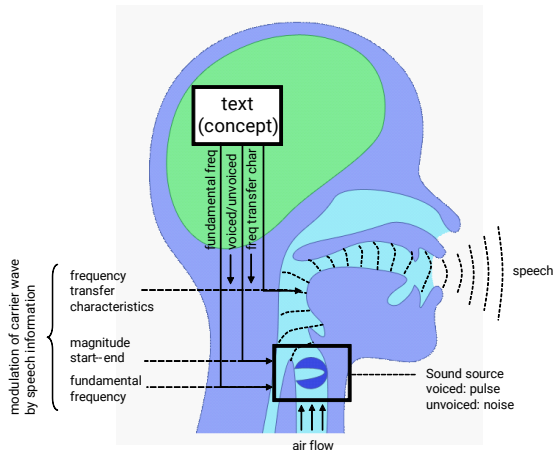


- Layered architecture with non-linear operations
- Can model high-dimensional/correlated linguistic/acoustic features
→ Feature extraction can be embedded in model itself



ANN-based acoustic model [21] – Motivation (3)

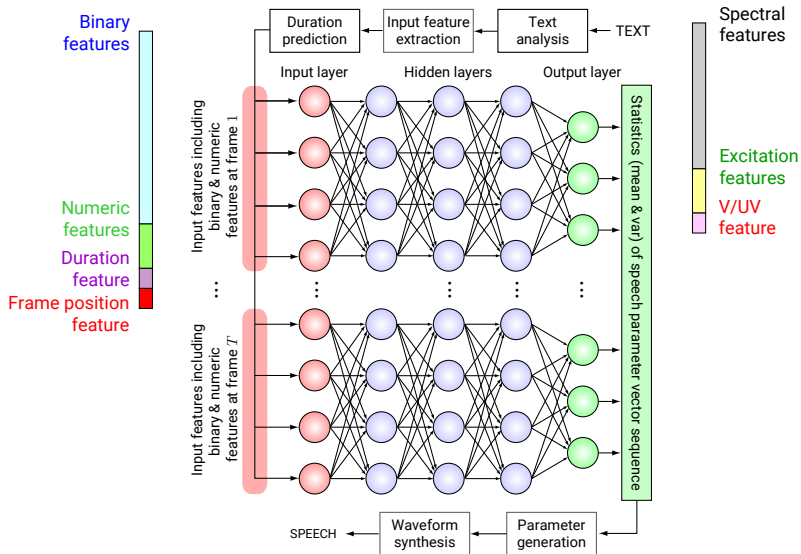
Implicitly mimic layered hierarchical structure in speech production



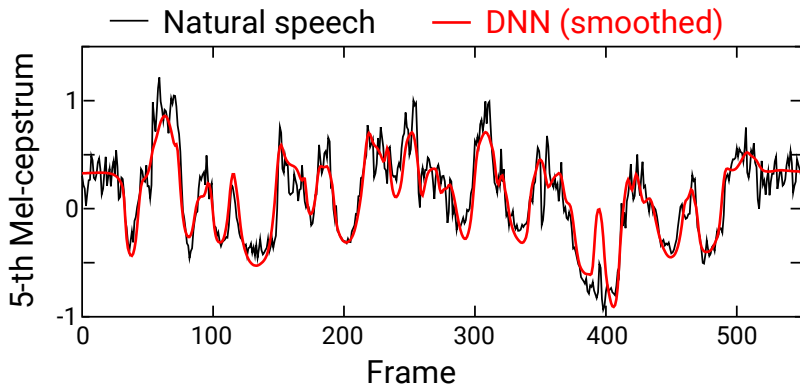
Concept → Linguistic → Articulator → Vocal tract → Waveform



DNN-based speech synthesis [21] – Implementation



DNN-based speech synthesis [21] – Example



DNN-based speech synthesis [21] – Subjective eval.

Compared HMM- & DNN-based TTS w/ similar # of parameters

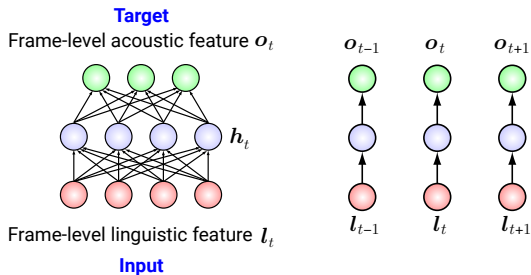
- US English, professional speaker, 30 hours of speech data
- Preference test
- 173 test sentences, 5 subjects per pair
- Up to 30 pairs per subject
- Crowd-sourced

Preference scores (higher one is better)

HMM	DNN	No pref.	#layers × #units
15.8%	38.5%	45.7%	4 × 256
16.1%	27.2%	56.7%	4 × 512
12.7%	36.6%	50.7%	4 × 1024



Feedforward NN-based acoustic model – Limitation

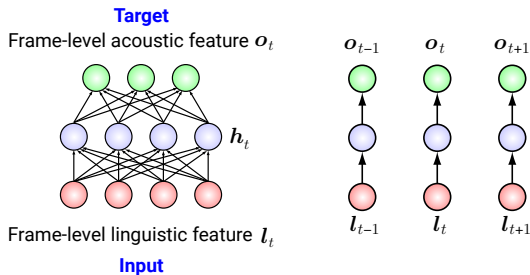


Each frame is mapped independently → **Smoothing is still essential**

Preference scores (higher one is better)		
DNN with dyn	DNN without dyn	No pref.
67.8%	12.0%	20.0%



Feedforward NN-based acoustic model – Limitation



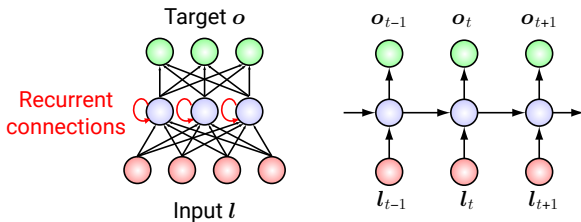
Each frame is mapped independently → **Smoothing is still essential**

Preference scores (higher one is better)		
DNN with dyn	DNN without dyn	No pref.
67.8%	12.0%	20.0%

Recurrent connections → **Recurrent NN (RNN) [27]**



RNN-based acoustic model [28, 29]

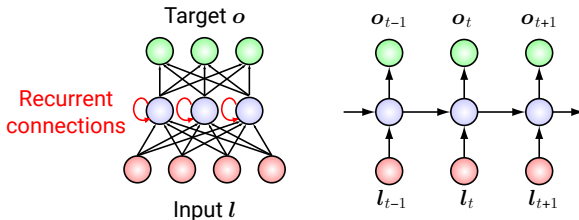


$$h_t = f(W_{hl}l_t + W_{hh}h_{t-1} + b_h) \quad \hat{o}_t = W_{oh}h_t + b_o$$
$$\hat{\Lambda} = \arg \min_{\Lambda} \sum_t \|o_t - \hat{o}_t\|_2 \quad \Lambda = \{W_{hl}, W_{hh}, W_{oh}, b_h, b_o\}$$

- DNN: $\hat{o}_t \approx \mathbb{E}[o_t | l_t]$
- RNN: $\hat{o}_t \approx \mathbb{E}[o_t | l_1, \dots, l_t]$



RNN-based acoustic model [28, 29]

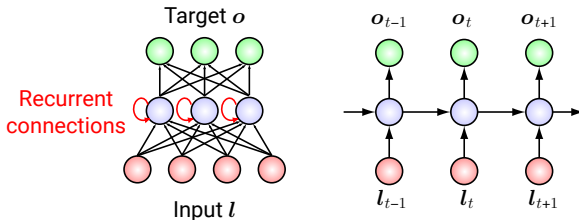


- Only able to use previous contexts

→ Bidirectional RNN [27]: $\hat{o}_t \approx \mathbb{E}[o_t | l_1, \dots, l_T]$



RNN-based acoustic model [28, 29]



- Only able to use previous contexts

→ Bidirectional RNN [27]: $\hat{o}_t \approx \mathbb{E}[o_t | l_1, \dots, l_T]$

- Trouble accessing long-range contexts

- Information in hidden layers loops quickly decays over time
- Prone to being overwritten by new information from inputs
- Long short-term memory (LSTM) [30]



LSTM-RNN-based acoustic model [29]

Subjective preference test (same US English data)

DNN: 3 layers, 1024 units

LSTM: 1 layer, 256 LSTM units

DNN with dyn	LSTM with dyn	No pref.
18.4%	34.9%	47.6%



LSTM-RNN-based acoustic model [29]

Subjective preference test (same US English data)

DNN: 3 layers, 1024 units

LSTM: 1 layer, 256 LSTM units

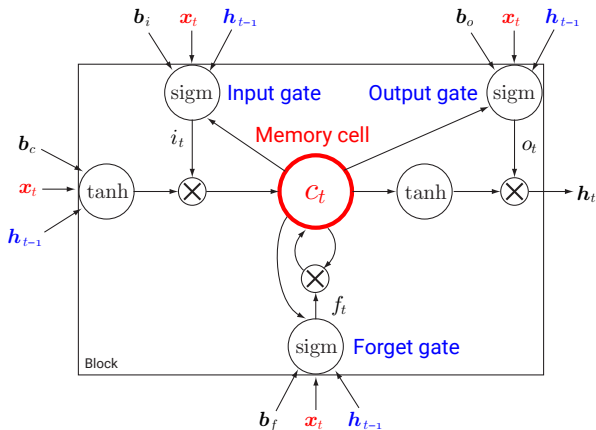
DNN with dyn	LSTM with dyn	No pref.
18.4%	34.9%	47.6%

LSTM with dyn	LSTM without dyn	No pref.
21.0%	12.2%	66.8%

→ Smoothing was still effective



Why?



Gate output: 0 -- 1

Input gate == 1
→ Write memory

Forget gate == 0
→ Reset memory

Output gate == 1
→ Read memory

- Gates in LSTM units: 0/1 switch controlling information flow
- Can produce rapid change in outputs
→ **Discontinuity**



How?

- Using loss function incorporating continuity



How?

- Using loss function incorporating continuity
- Integrate smoothing → Recurrent output layer [29]

$$h_t = \text{LSTM}(l_t) \quad \hat{o}_t = \mathbf{W}_{oh}h_t + \mathbf{W}_{oo}\hat{o}_{t-1} + \mathbf{b}_o$$



How?

- Using loss function incorporating continuity
- Integrate smoothing → Recurrent output layer [29]

$$h_t = \text{LSTM}(l_t) \quad \hat{o}_t = \mathbf{W}_{oh}h_t + \mathbf{W}_{oo}\hat{o}_{t-1} + \mathbf{b}_o$$

Works pretty well

LSTM with dyn (Feedforward)	LSTM without dyn (Recurrent)	No pref.
21.8%	21.0%	57.2%



How?

- Using loss function incorporating continuity
- Integrate smoothing → Recurrent output layer [29]

$$h_t = \text{LSTM}(l_t) \quad \hat{o}_t = W_{oh}h_t + W_{oo}\hat{o}_{t-1} + b_o$$

Works pretty well

LSTM with dyn (Feedforward)	LSTM without dyn (Recurrent)	No pref.
21.8%	21.0%	57.2%

Having two smoothing together doesn't work well → Oversmoothing?

LSTM with dyn (Recurrent)	LSTM without dyn (Recurrent)	No pref.
16.6%	29.2%	54.2%



Low-latency TTS by unidirectional LSTM-RNN [29]

HMM / DNN

- Smoothing by dyn. needs to solve set of T linear equations

$$\mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \mathbf{W} \mathbf{c} = \mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \boldsymbol{\mu}_{\hat{q}} \quad T: \text{Utterance length}$$



Low-latency TTS by unidirectional LSTM-RNN [29]

HMM / DNN

- Smoothing by dyn. needs to solve set of T linear equations

$$\mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \mathbf{W} \mathbf{c} = \mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \boldsymbol{\mu}_{\hat{q}} \quad T: \text{Utterance length}$$

- Order of operations to determine the first frame c_1 (latency)
 - Cholesky decomposition [7] $\rightarrow \mathcal{O}(T)$
 - Recursive approximation [31] $\rightarrow \mathcal{O}(L)$ L : lookahead, $10 \sim 30$



Low-latency TTS by unidirectional LSTM-RNN [29]

HMM / DNN

- Smoothing by dyn. needs to solve set of T linear equations

$$\mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \mathbf{W} \mathbf{c} = \mathbf{W}^\top \Sigma_{\hat{q}}^{-1} \boldsymbol{\mu}_{\hat{q}} \quad T: \text{Utterance length}$$

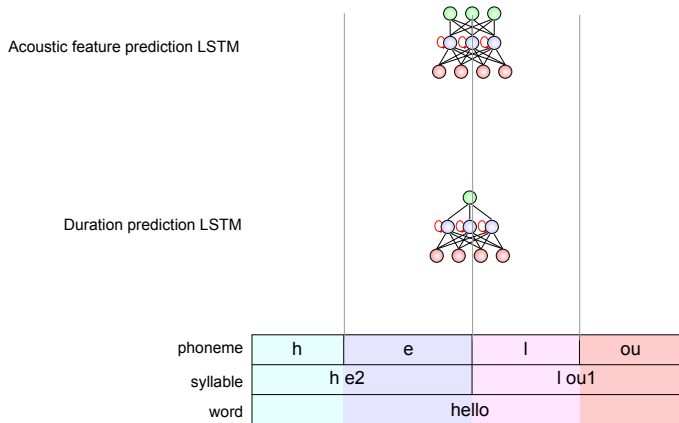
- Order of operations to determine the first frame c_1 (latency)
 - Cholesky decomposition [7] $\rightarrow \mathcal{O}(T)$
 - Recursive approximation [31] $\rightarrow \mathcal{O}(L)$ L : lookahead, $10 \sim 30$

Unidirectional LSTM with recurrent output layer [29]

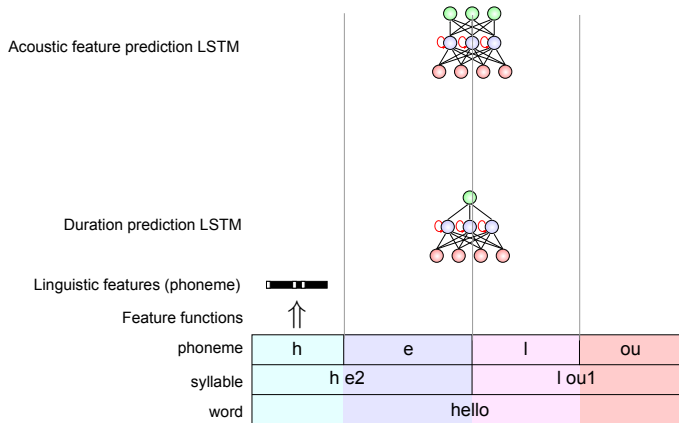
- No smoothing required, fully time-synchronous w/o lookahead
- Order of latency $\rightarrow \mathcal{O}(1)$



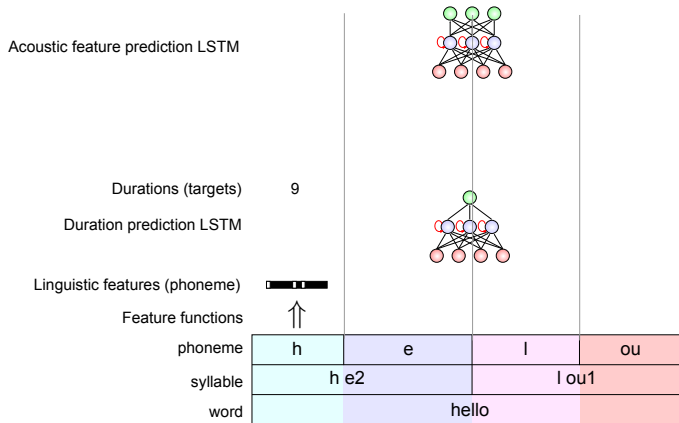
Low-latency TTS by LSTM-RNN [29] – Implementation



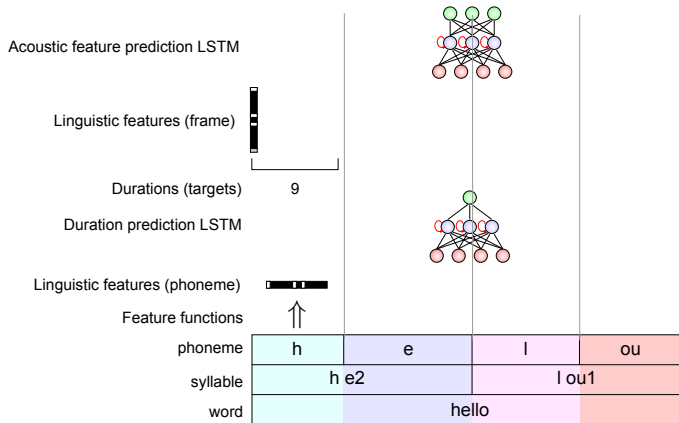
Low-latency TTS by LSTM-RNN [29] – Implementation



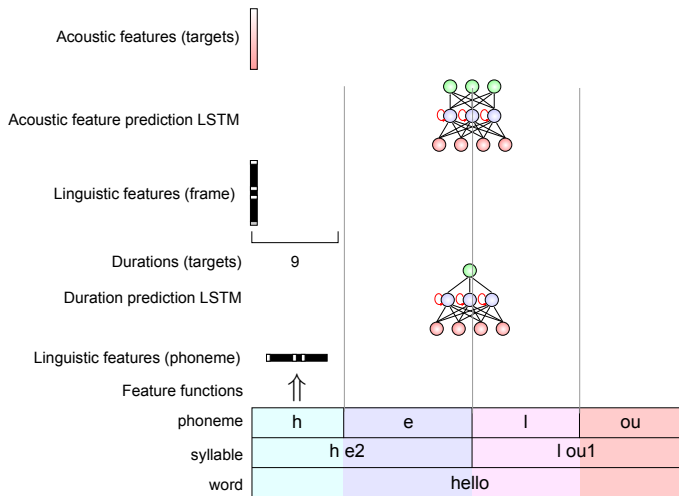
Low-latency TTS by LSTM-RNN [29] – Implementation



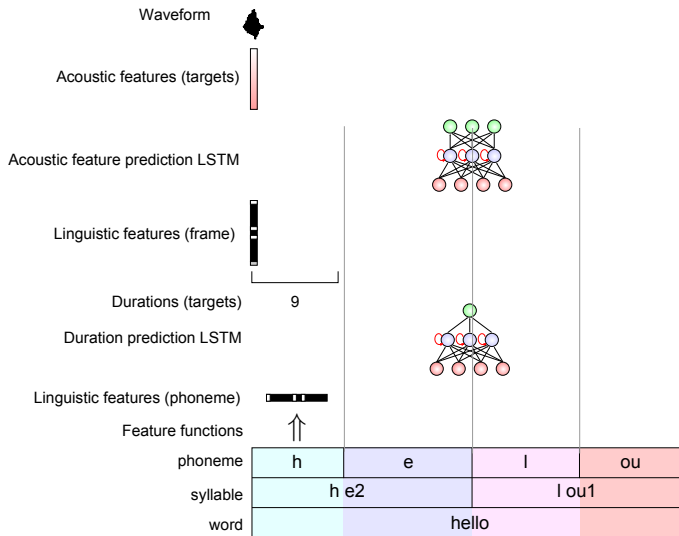
Low-latency TTS by LSTM-RNN [29] – Implementation



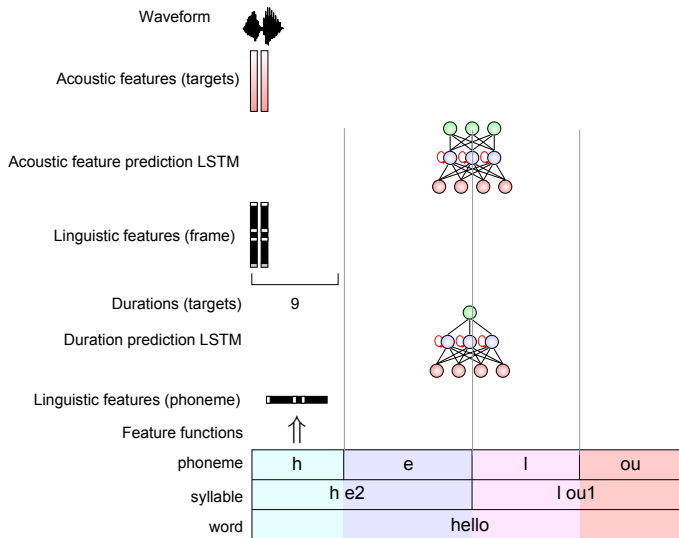
Low-latency TTS by LSTM-RNN [29] – Implementation



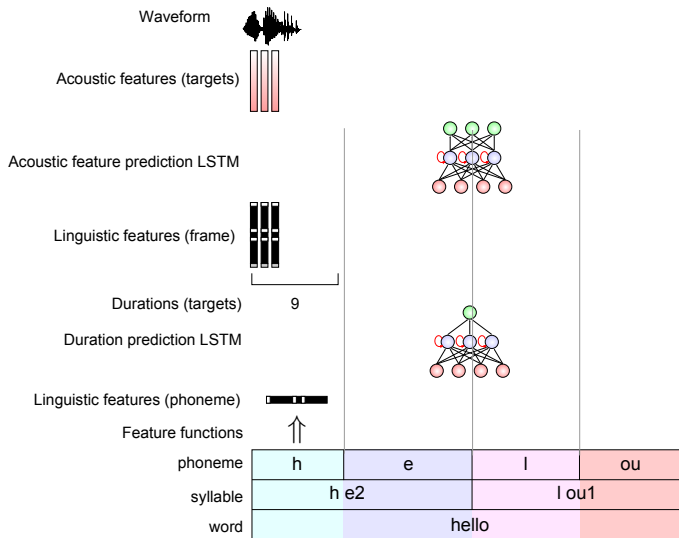
Low-latency TTS by LSTM-RNN [29] – Implementation



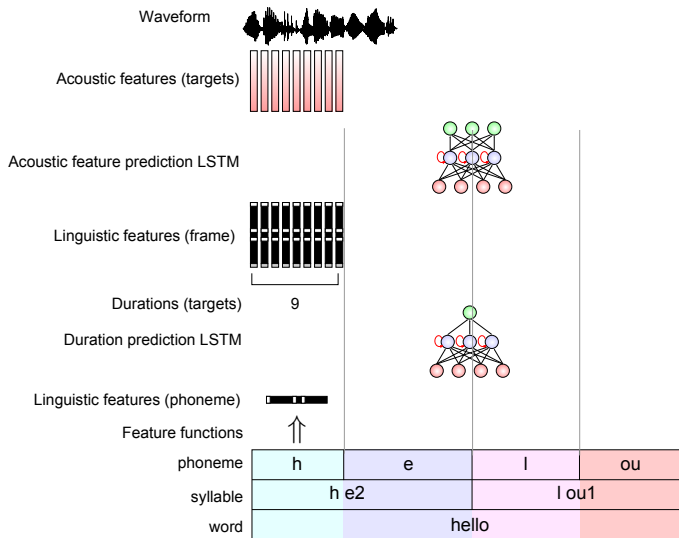
Low-latency TTS by LSTM-RNN [29] – Implementation



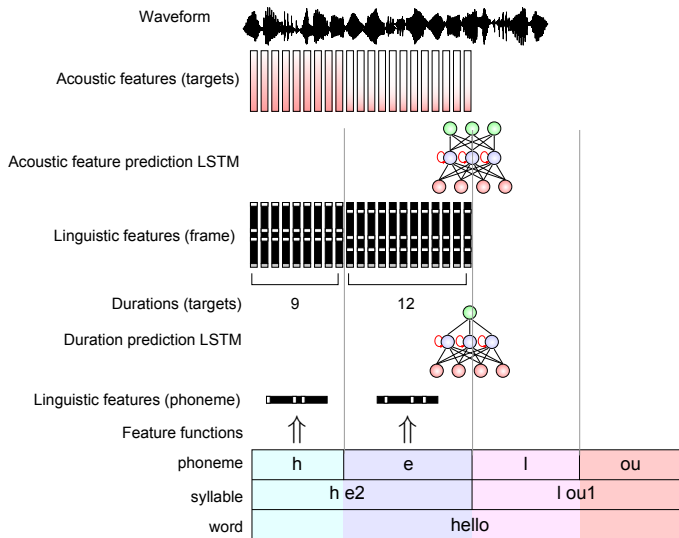
Low-latency TTS by LSTM-RNN [29] – Implementation



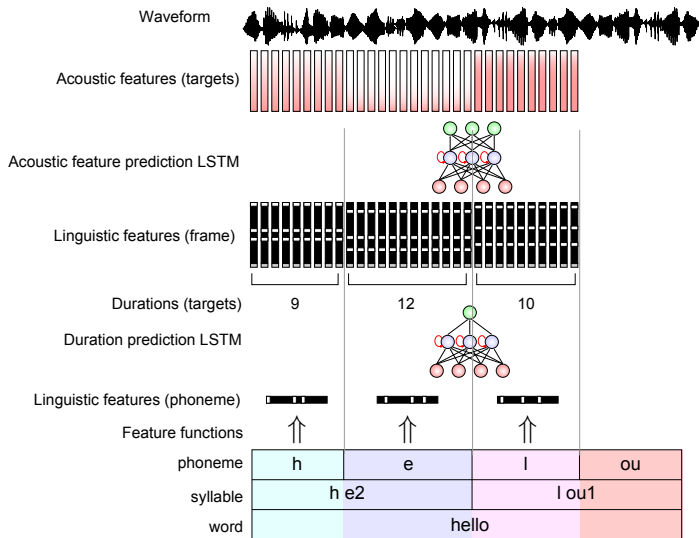
Low-latency TTS by LSTM-RNN [29] – Implementation



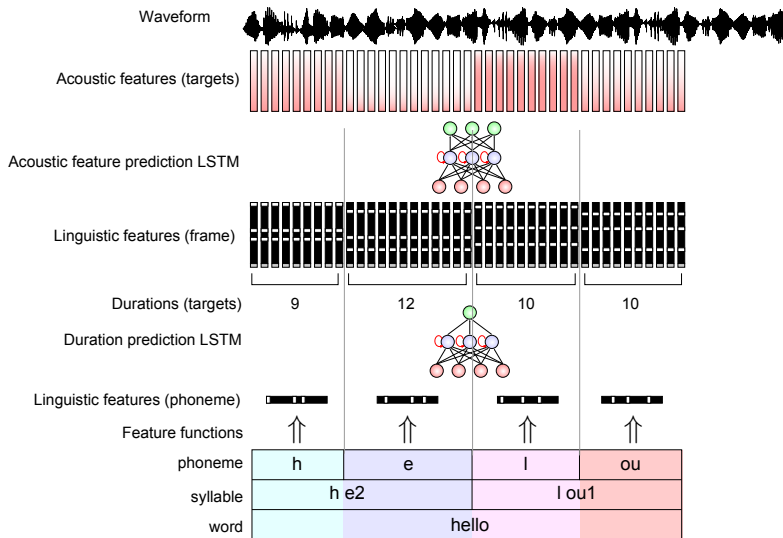
Low-latency TTS by LSTM-RNN [29] – Implementation



Low-latency TTS by LSTM-RNN [29] – Implementation



Low-latency TTS by LSTM-RNN [29] – Implementation



Some comments

Is this new? . . . no

- Feedforward NN-based speech synthesis [32]
- RNN-based speech synthesis [33]



Some comments

Is this new? . . . no

- Feedforward NN-based speech synthesis [32]
- RNN-based speech synthesis [33]

What's the difference?

- More layers, data, computational resources
- Better learning algorithm
- Modern SPSS techniques



Making LSTM-RNN-based TTS into production

Client-side (local) TTS for Android



Google Text-to-speech

Google Inc. Tools

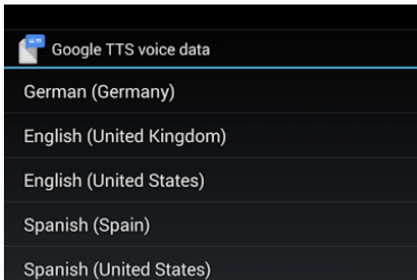
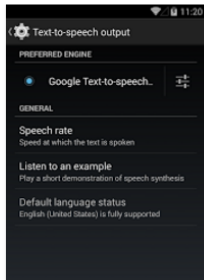
Top Developer

★★★★★ 546,569

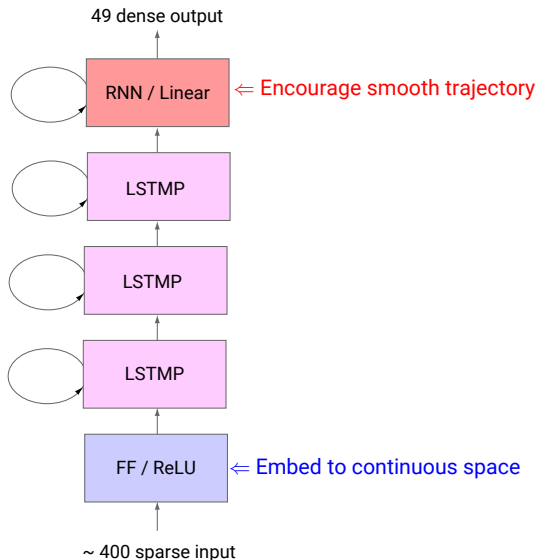
PEGI 3

This app is compatible with all of your devices.

Installed

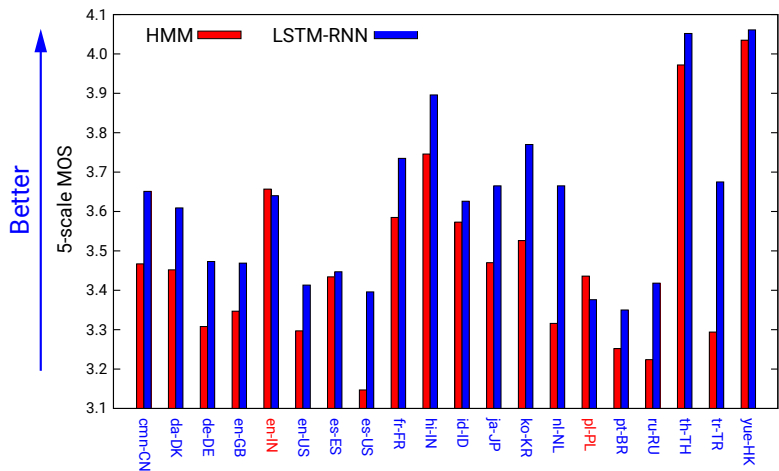


Network architecture



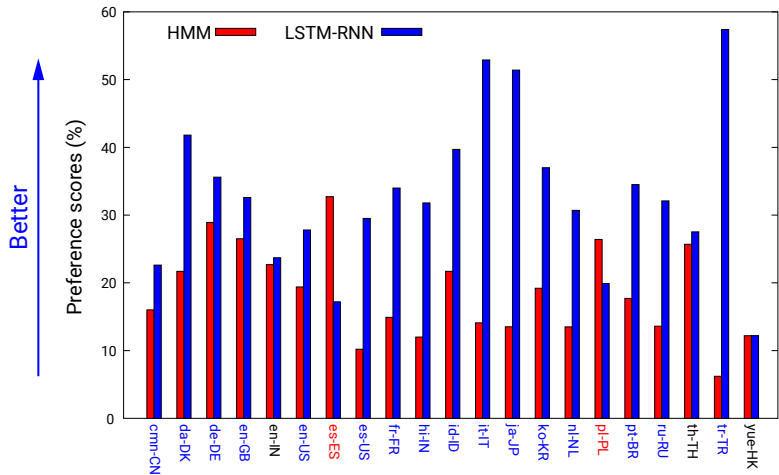
Results – HMM / LSTM-RNN

Subjective 5-scale Mean Opinion Score test (i18n)



Results – HMM / LSTM-RNN

Subjective preference test (i18n)



Results – HMM / LSTM-RNN

Latency & Battery/CPU usage

Latency (Nexus 7 2013)

Sentence	Average/Max latency (ms)	
	HMM	LSTM-RNN
very short (1 character)	26/30	37/72
short (~30 characters)	123/172	63/88
long (~80 characters)	311/418	118/190

CPU usage

HMM → LSTM-RNN: **+48%**

Battery usage (Daily usage by a blind Googler)

HMM: **2.8%** of 1475 mAH → LSTM-RNN: **4.8%** of 1919 mAH



Results – HMM / LSTM-RNN

Summary

- **Naturalness**
LSTM-RNN > HMM
- **Latency**
LSTM-RNN < HMM
- **CPU/Battery usage**
LSTM-RNN > HMM

LSTM-RNN-based TTS is in production at Google



Outline

Background

HMM-based acoustic modeling

- Training & synthesis
- Limitations

ANN-based acoustic modeling

- Feedforward NN
- RNN

Conclusion



Acoustic models for speech synthesis – Summary

- **HMM**
 - Discontinuity due to step-wise statistics
 - Difficult to integrate feature extraction
 - Fragmented representation



Acoustic models for speech synthesis – Summary

- **HMM**
 - Discontinuity due to step-wise statistics
 - Difficult to integrate feature extraction
 - Fragmented representation
- **Feedforward NN**
 - Easier to integrate feature extraction
 - Distributed representation
 - Discontinuity due to frame-by-frame independent mapping



Acoustic models for speech synthesis – Summary

- **HMM**
 - Discontinuity due to step-wise statistics
 - Difficult to integrate feature extraction
 - Fragmented representation
- **Feedforward NN**
 - Easier to integrate feature extraction
 - Distributed representation
 - Discontinuity due to frame-by-frame independent mapping
- **(LSTM) RNN**
 - Smooth → Low latency



Acoustic models for speech synthesis – Future topics

- **Visualization for debugging**
 - Concatenative → Easy to debug
 - HMM → Hard
 - ANN → Harder



Acoustic models for speech synthesis – Future topics

- **Visualization for debugging**

- Concatenative → Easy to debug
- HMM → Hard
- ANN → Harder

- **More flexible voice-based user interface**

- Concatenative → Record all possibilities
- HMM → Weak/rare signals (input) are often ignored
- ANN → Weak/rare signals can contribute



Acoustic models for speech synthesis – Future topics

- **Visualization for debugging**

- Concatenative → Easy to debug
- HMM → Hard
- ANN → Harder

- **More flexible voice-based user interface**

- Concatenative → Record all possibilities
- HMM → Weak/rare signals (input) are often ignored
- ANN → Weak/rare signals can contribute

- **Fully integrate feature extraction**

- Current: Linguistic features → Acoustic features
- Goal: Character sequence → Speech waveform



Thanks!



References I

- [1] E. Moulines and F. Charpentier.
Pitch synchronous waveform processing techniques for text-to-speech synthesis using diphones.
Speech Commn., 9:453–467, 1990.
- [2] A. Hunt and A. Black.
Unit selection in a concatenative speech synthesis system using a large speech database.
In *Proc. ICASSP*, pages 373–376, 1996.
- [3] T. Yoshimura, K. Tokuda, T. Masuko, T. Kobayashi, and T. Kitamura.
Simultaneous modeling of spectrum, pitch and duration in HMM-based speech synthesis.
In *Proc. Eurospeech*, pages 2347–2350, 1999.
- [4] H. Zen, K. Tokuda, and A. Black.
Statistical parametric speech synthesis.
Speech Commn., 51(11):1039–1064, 2009.
- [5] L. Rabiner.
A tutorial on hidden Markov models and selected applications in speech recognition.
In *Proc. IEEE*, volume 77, pages 257–285, 1989.
- [6] J. Odell.
The use of context in large vocabulary speech recognition.
PhD thesis, Cambridge University, 1995.
- [7] K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura.
Speech parameter generation algorithms for HMM-based speech synthesis.
In *Proc. ICASSP*, pages 1315–1318, 2000.
- [8] H. Zen, K. Tokuda, and T. Kitamura.
Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic features.
Comput. Speech Lang., 21(1):153–173, 2007.



References II

- [9] K. Yu, F. Mairesse, and S. Young.
Word-level emphasis modelling in HMM-based speech synthesis.
In *Proc. ICASSP*, pages 4238–4241, 2010.
- [10] M. Shannon, H. Zen, and W. Byrne.
Autoregressive models for statistical parametric speech synthesis.
IEEE Trans. Acoust. Speech Lang. Process., 21(3):587–597, 2013.
- [11] C. Quillen.
Kalman filter based speech synthesis.
In *Proc. ICASSP*, pages 4618–4621, 2010.
- [12] V. Tsirias, R. Maia, V. Diakouloukas, Y. Stylianou, and V. Digalakis.
Linear dynamical models in speech synthesis.
In *Proc. ICASSP*, pages 300–304, 2014.
- [13] T. Toda and K. Tokuda.
Statistical approach to vocal tract transfer function estimation based on factor analyzed trajectory hmm.
In *Proc. ICASSP*, pages 3925–3928, 2008.
- [14] Y.-J. Wu and K. Tokuda.
Minimum generation error training with direct log spectral distortion on LSPs for HMM-based speech synthesis.
In *Proc. Interspeech*, pages 577–580, 2008.
- [15] R. Maia, H. Zen, and M. Gales.
Statistical parametric speech synthesis with joint estimation of acoustic and excitation model parameters.
In *Proc. ISCA SSW7*, pages 88–93, 2010.
- [16] K. Nakamura, K. Hashimoto, Y. Nankaku, and K. Tokuda.
Integration of spectral feature extraction and modeling for HMM-based speech synthesis.
IEICE Trans. Inf. Syst., E97-D(6):1438–1448, 2014.



References III

- [17] K. Yu, H. Zen, F. Mairesse, and S. Young.
Context adaptive training with factorized decision trees for HMM-based statistical parametric speech synthesis.
Speech Commn., 53(6):914–923, 2011.
- [18] H. Zen, M. Gales, Y. Nankaku, and K. Tokuda.
Product of experts for statistical parametric speech synthesis.
IEEE Trans. Audio Speech Lang. Process., 20(3):794–805, 2012.
- [19] K. Saino.
A clustering technique for factor analysis-based eigenvoice models.
Master thesis, Nagoya Institute of Technology, 2008.
(in Japanese).
- [20] H. Zen, N. Braunschweiler, S. Buchholz, M. Gales, K. Knill, S. Krstulovic, and J. Latorre.
Statistical parametric speech synthesis based on speaker and language factorization.
IEEE Trans. Audio, Speech, Lang. Process., 20(6):1713–1724, 2012.
- [21] H. Zen, A. Senior, and M. Schuster.
Statistical parametric speech synthesis using deep neural networks.
In *Proc. ICASSP*, pages 7962–7966, 2013.
- [22] G. Hinton, J. McClelland, and D. Rumelhart.
Distributed representation.
In D. Rumelhart, J. McClelland, and the PDP Research Group, editors, *Parallel distributed processing: Explorations in the microstructure of cognition*. MIT Press, 1986.
- [23] Y. Bengio.
Deep learning: Theoretical motivations.
<http://www.iro.umontreal.ca/~bengioy/talks/dlss-3aug2015.pdf>, 2015.



References IV

- [24] C. Valentini-Botinhao, Z. Wu, and S. King.
Towards minimum perceptual error training for DNN-based speech synthesis.
In Proc. Interspeech, pages 869–873, 2015.
- [25] S. Takaki, S.-J. Kim, J. Yamagishi, and J.-J. Kim.
Multiple feed-forward deep neural networks for statistical parametric speech synthesis.
In Interspeech, pages 2242–2246, 2015.
- [26] K. Tokuda and H. Zen.
Directly modeling speech waveforms by neural networks for statistical parametric speech synthesis.
In Proc. ICASSP, pages 4215–4219, 2015.
- [27] M. Schuster and K. Paliwal.
Bidirectional recurrent neural networks.
IEEE Trans. Signal Process., 45(11):2673–2681, 1997.
- [28] Y. Fan, Y. Qian, and F. Soong.
TTS synthesis with bidirectional LSTM based recurrent neural networks.
In Proc. Interspeech, pages 1964–1968, 2014.
- [29] H. Zen and H. Sak.
Unidirectional long short-term memory recurrent neural network with recurrent output layer for low-latency speech synthesis.
In Proc. ICASSP, pages 4470–4474, 2015.
- [30] S. Hochreiter and J. Schmidhuber.
Long short-term memory.
Neural Comput., 9(8):1735–1780, 1997.



References V

- [31] K. Koishida, K. Tokuda, T. Masuko, and T. Kobayashi.
Vector quantization of speech spectral parameters using statistics of dynamic features.
In *Proc. ICSP*, pages 247–252, 1997.
- [32] O. Karaali, G. Corrigan, and I. Gerson.
Speech synthesis with neural networks.
In *Proc. World Congress on Neural Networks*, pages 45–50, 1996.
- [33] C. Tuerk and T. Robinson.
Speech synthesis using artificial neural networks trained on cepstral coefficients.
In *Proc. Eurospeech*, pages 1713–1716, 1993.

