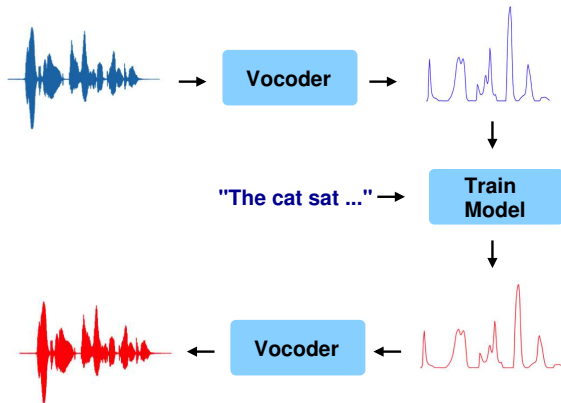


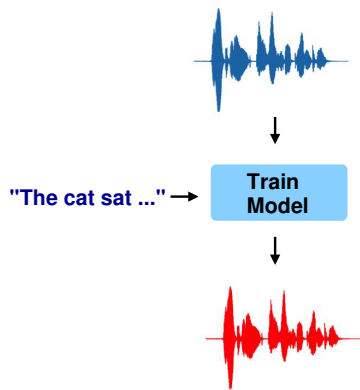
# MLMI10: Waveform Level Synthesis

Mark Gales

Lent 2019



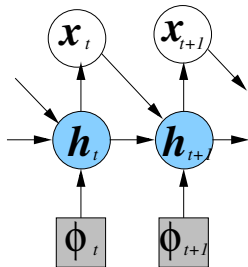
- Quality of vocoder fundamentally limits performance



- Challenges:
  - 16KHz/24KHz predictions
  - long-span dependencies
  - 16KHz/24KHz labels

# Long-Span Dependencies

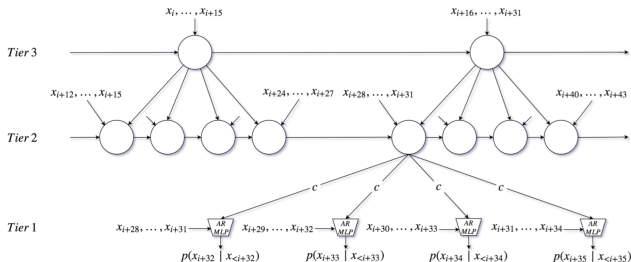
- Speech events act over about 250ms (roughly)
  - 5ms frame rate: this is 50 frames
  - 16Hz samples: this is 4000 samples



- Consider recurrent model:
  - history vector representation
  - rule-of-thumb  $< 100$  samples  
6ms history!
  - insufficient history memory

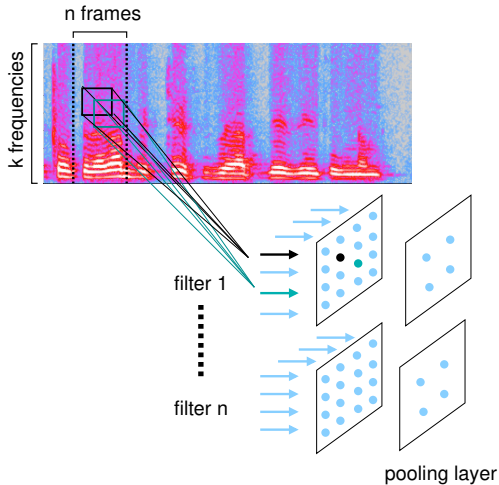
Alternative model (probably) required!

# Clockwork RNNs and SampleRNN [15, 10]



- Use a hierarchical RNN architecture - multiple tiers
  - each tier operates at a different timescale: in diagram  
Tier1: 16KHz Tier2: 4KHz Tier3: 1KHz
  - like an (analogue) clock: hours, minutes, seconds
- Allows long-term dependencies in RNN architectures

# Convolutional Neural Networks [12]



- Various parameters control the form of the CNN:
  - **number** (depth): how many filters to use
  - **receptive field** (filter size): height/width/depth ( $h \times w \times d$ )
  - **stride**: how far filter moves in the convolution
  - **dilation**: “gaps” between filter elements
  - **zero-padding**: do you pad the edges of the “image” with zeroes
- Filter output can be stacked to yield depth for next layer

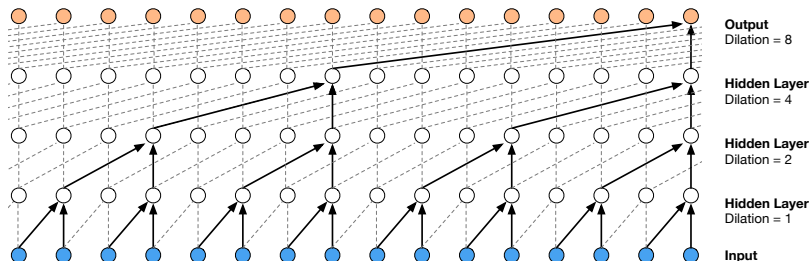
- Various parameters control the form of the CNN:
  - **number** (depth): how many filters to use
  - **receptive field** (filter size): height/width/depth ( $h \times w \times d$ )
  - **stride**: how far filter moves in the convolution
  - **dilation**: “gaps” between filter elements
  - **zero-padding**: do you pad the edges of the “image” with zeroes
- Filter output can be stacked to yield depth for next layer
- To illustrate the impact consider 1-dimensional case - default:
  - **zero-padding**, **stride=1**, **dilation=0**

**sequence:** 1, 2, 3, 4, 5, 6      **filter:** 1, 1, 1

**default:** 3, 6, 9, 12, 15, 11      **no padding:** 6, 9, 12, 15

**dilation=1:** 4, 6, 9, 12, 8, 10      **stride=2:** 3, 9, 15





- Use CNN acting over the one-dimensional waveform sequence
  - dilation used to control growth in number of parameters
  - different dilations at different levels
- A **gated** activation function is used (layer  $k$ , filter  $f$ )

$$\tanh\left(\mathbf{w}_f^{(k)} * \mathbf{x}_{1:T}\right) \odot \sigma\left(\mathbf{w}_g^{(k)} * \mathbf{x}_{1:T}\right)$$

- Simplest training criterion - minimise

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \log(p(x_t | \mathbf{x}_{1:t-1}; \theta)) \propto \sum_{t=1}^T \frac{(\hat{x}_t - x_t)^2}{\sigma^2}$$

- $\hat{x}_t = f(\mathbf{x}_{1:t-1}; \theta)$  is the prediction
  - Gaussian noise added to systematic system - is this good?
- Adopt “ASR-style” estimation - run as classification
  - quantise the values to predict  $x_t \rightarrow x_t^q$  ( $\mu$ -law)
  - minimise the cross-entropy to the real target  $\mathbf{x}_{1:T} \rightarrow \mathbf{x}_{1:T}^q$

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \sum_{k=1}^K \delta(\omega_k, x_t^q) \log(P(\omega_k | \mathbf{x}_{1:t-1}; \theta))$$

- Quantising the waveform allows classification training
  - no assumption about Gaussian cost functions (important)
- Training criterion interesting mix of quantised and continuous

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \sum_{k=1}^K \delta(\omega_k, x_t^q) \log(P(\omega_k | \mathbf{x}_{1:t-1}; \theta))$$

- predict quantised value  $x_t^q$
  - history is continuous  $\mathbf{x}_{1:t-1}$
- **Sample generation** follows simple process:
  - obtain discrete distribution  $P(x_t^q | \hat{\mathbf{x}}_{1:t-1}; \theta)$
  - sample from discrete distribution to yield  $\hat{x}_t^q$
  - convert from discrete to continuous  $\hat{x}_t^q \rightarrow \hat{x}_t$
  - append to history  $\hat{\mathbf{x}}_{1:t} = \{\hat{\mathbf{x}}_{1:t-1}, \hat{x}_t\}$ ,  $t = t + 1$ , repeat

# Conditional Synthesis

- The model described so far have no conditioning on the text
  - generates speech like sound, but not speech
- Additional information (beyond context labels) used:
  - **duration**: LSTM-RNN-based phone duration
  - **log-F0**: autoregressive CNN-based log-F0 prediction models
- Add conditioning vector  $\lambda$  yields information about
  - **text context**  $\lambda_{1:T}^{(q)}$ : up-sampled phonetic information from text
  - **speaker**  $\lambda^{(s)}$ : speaker representation

- Activation function then has the form (ignoring gating term)
  - for filter  $f$  of layer  $k$

$$\tanh \left( \mathbf{w}_f^{(k)} * \mathbf{x}_{1:T} + \mathbf{v}_f^{(k)} * \lambda_{1:T}^{(q)} + [\mathbf{u}_f^{(k)\top} \lambda^{(s)}]_{1:T} \right) \odot \sigma (...)$$

- where  $[\dots]_{1:T}$  indicates repeating the output from 1 to  $T$

# Parallel Wavenet

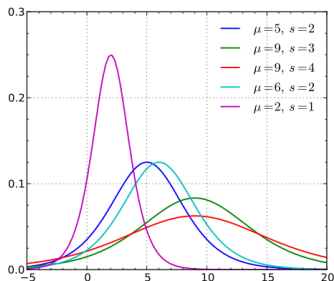
- WaveNet generates high quality samples
  - but orders of magnitude too much latency
  - limits sampling frequency to 16KHz
  - only able to support 8-bit resolution
- Rather than swapping neural vocoder
  - modified output distribution to support 16-bit resolution
  - modified synthesis to allow parallel generation of samples
  - different architectures for training and synthesis

# Output Distribution: Logistic Distribution

- Interested in a distribution  $\mathbb{L}(x; \mu, s)$  where

$$\int_{-\infty}^a \mathbb{L}(x; \mu, s) dx = \sigma\left(\frac{a - \mu}{s}\right); \quad \sigma(x) = \frac{1}{1 + \exp(-x)}$$

- $\sigma()$  is the **logistic function** - a valid CDF
- This is the **Logistic distribution**



- the form of the PDF is then

$$\mathbb{L}(x; \mu, s) = \frac{\exp(-(x - \mu)/s)}{s(1 + \exp(-(x - \mu)/s))^2}$$

- $\mu$  is the mean (**location**) of the distribution
- $s > 0$  is the **scale** parameter

- Modelling using a **softmax** for the distribution limiting
  - doesn't scale well with the number of bin values
- Using the CDF of Logistic distributed variables

$$\int_a^b \mathbb{L}(x; \mu, s) dx = \sigma\left(\frac{b - \mu}{s}\right) - \sigma\left(\frac{a - \mu}{s}\right)$$

- Extending to a **mixture model** yields and examining a bin value

$$P(x^q | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{s}) = \sum_{m=1}^M \pi_m \left( \sigma\left(\frac{x^q + 0.5 - \mu_m}{s_m}\right) - \sigma\left(\frac{x^q - 0.5 - \mu_m}{s_m}\right) \right)$$

- recognises that bins 127 and 129 are near 128
- allows expansion of targets with increasing cost



- Generate  $T$ , independent Logistic distributed, samples  $\mathbf{z}_{1:T}$ 
  - apply (causal) transformation from  $\mathbf{z}_{1:T} \rightarrow \mathbf{x}_{1:T}$

$$\log(p(\mathbf{x}_{1:T})) = \log(p(\mathbf{z}_{1:T})) - \log\left(\left|\frac{d\mathbf{x}_{1:T}}{d\mathbf{z}_{1:T}}\right|\right); \quad x_t = f(\mathbf{z}_{1:t})$$

- possible to show that (from the causal aspect of transform)

$$\log\left(\left|\frac{d\mathbf{x}_{1:T}}{d\mathbf{z}_{1:T}}\right|\right) = \sum_{t=1}^T \log\left(\frac{\partial f(\mathbf{z}_{1:t})}{\partial z_t}\right)$$

- The following transformation is applied

$$x_t = z_t f_s(\mathbf{z}_{1:t-1}; \boldsymbol{\theta}) + f_\mu(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})$$

- $f_s(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})$ : scale parameter prediction from  $\mathbf{z}_{1:t-1}$
- $f_\mu(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})$ : location parameter prediction from  $\mathbf{z}_{1:t-1}$

# Inverse Autoregressive Flows (IAFs) (cont)

- As  $z_t$  is Logistic ( $\mathbb{L}(0, 1)$ ) distributed this means that

$$p(x_t | \mathbf{z}_{1:t-1}; \boldsymbol{\theta}) = \mathbb{L}(x_t; f_\mu(\mathbf{z}_{1:t-1}; \boldsymbol{\theta}), f_s(\mathbf{z}_{1:t-1}; \boldsymbol{\theta}))$$

- $f_\mu(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})$  and  $f_s(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})$  from WaveNet style network
- Allows all samples  $\mathbf{x}_{1:T}$  to be generated in parallel
  - $x_t$  no explicit dependence on previous generated samples  $\mathbf{x}_{1:t-1}$
  - $x_t$  only depends on the noise samples  $\mathbf{z}_{1:t-1}$
  - can generate  $\mathbf{x}_{1:T}$  in parallel given  $\mathbf{z}_{1:T}$  - [Parallel WaveNet](#)
- Transforms random sequence  $\mathbf{z}_{1:T}$  to structured  $\mathbf{x}_{1:T}$ 
  - too challenging to achieve with a single “flow”
  - stack multiple flows together for parallel WaveNet

- Training parallel WaveNet directly is slow
  - training log-likelihood calculation is sequential (see IAF)
  - WaveNet training efficient (prediction  $x_t$  uses real data  $\mathbf{x}_{1:t-1}$ )
- Given a trained WaveNet  $p_w()$  train  $p(\mathbf{x}; \theta)$  to minimise

$$\mathcal{L}(\theta) = \mathcal{D}_{\text{kl}}(p() \| p_w()) = \sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}; \theta) \log \left( \frac{p(\mathbf{x}; \theta)}{p_w(\mathbf{x})} \right)$$

- $\mathcal{X}$  is the set of all waveforms
- contrast to teacher-student training

$$\mathcal{L}(\theta) = \sum_{\mathbf{x} \in \mathcal{X}} p_w(\mathbf{x}) \log (p(\mathbf{x}; \theta))$$

- Two terms in criterion (both can be computed efficiently)
  - student entropy
  - cross-entropy

- Student entropy for a sequence of length  $T$

$$\begin{aligned}\sum_{\mathbf{x}_{1:T} \in \mathcal{X}} p(\mathbf{x}_{1:T}; \boldsymbol{\theta}) \log(p(\mathbf{x}_{1:T}; \boldsymbol{\theta})) &= -\mathbb{E} \left\{ \sum_{t=1}^T \log(p(x_t | \mathbf{z}_{1:t-1}; \boldsymbol{\theta})) \right\}_{p(\mathbf{z}_{1:T})} \\ &= \mathbb{E} \left\{ \sum_{t=1}^T \log(f_s(\mathbf{z}_{1:t-1}; \boldsymbol{\theta})) \right\}_{p(\mathbf{z}_{1:T})} + 2T\end{aligned}$$

- uses equality for entropy of a logistic distribution
- can be computed without generating samples

# Cross-Entropy (to WaveNet model)

- The cross-entropy term has the following form

$$\sum_{\mathbf{x}_{1:T} \in \mathcal{X}} p(\mathbf{x}_{1:T}; \boldsymbol{\theta}) \log(p_w(\mathbf{x}_{1:T}))$$

- this requires drawing sample sequences from  $p(\mathbf{x}; \boldsymbol{\theta})$
- For a  $T$  length sequence this can be written as (see paper)

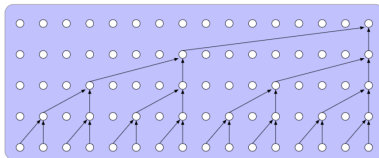
$$\sum_{t=1}^T \mathbb{E} \{ p(x_t | \mathbf{x}_{1:t-1}; \boldsymbol{\theta}) \log(p_w(x_t | \mathbf{x}_{1:t-1})) \}_{p(\mathbf{x}_{1:t-1}; \boldsymbol{\theta})}$$

- this is efficient (see paper!)

# Probability Density Distillation Overview

## WaveNet Teacher

Linguistic features ----->

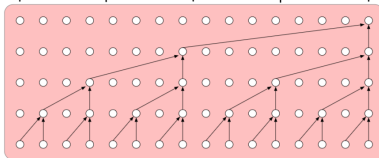


Teacher Output

$$P(x_i | x_{<i})$$

## WaveNet Student

Linguistic features ----->



Generated Samples

$$x_i = g(z_i | z_{<i})$$

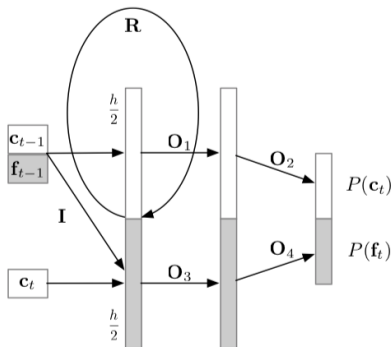
Student Output

$$P(x_i | z_{<i})$$

Input noise

$$z_i$$

# WaveRNN



- Modification to original WaveNet configuration
  - simplified (recurrent) architecture
  - split representation coarse ( $c_t$ ) and fine ( $f_t$ )
  - efficient, improved, resolution



# WaveRNN Diagram Notation

- Split the 16-bit output into:
  - coarse,  $c_t$  8 most significant bits
  - fine,  $f_t$  8 least significant bits (conditioned on coarse)
- Diagram (from paper) uses notation
  - $O_1$  and  $O_2$ : “coarse” weight matrices
  - $O_3$  and  $O_4$ : “fine” weight matrices
  - $I$  weight matrix (not identity!) from paper
- For the network layers
  - first layer recurrent (variant on GRU)
  - second layer ReLU
  - third layer soft-max

(Additional speed-ups described in paper: [sparse WaveRNN](#))

- Split the prediction 16-bit prediction into 2 parts

$$\mathbf{x}_t = \begin{bmatrix} \mathbf{c}_t \\ \mathbf{f}_t \end{bmatrix}$$

- coarse**  $\mathbf{c}_t$  (8 bits - 256 levels): top 8 bits of waveform

$$\mathbf{c}_t = \mathcal{F}_c(\mathbf{x}_{t-1}, \mathbf{h}_{t-1}) = \mathcal{F}_c(\mathbf{c}_{t-1}, \mathbf{f}_{t-1}, \mathbf{h}_{t-1})$$

- fine**  $\mathbf{f}_t$  (8 bits - 256 levels): bottom 8 bits of waveform

$$\mathbf{f}_t = \mathcal{F}_c(\mathbf{c}_{t-1}, \mathbf{f}_{t-1}, \mathbf{h}_{t-1}, \mathbf{c}_t)$$

- Equations can be split into distinct parts (ignoring biases)

$$\begin{aligned}P(\mathbf{c}_t) &= \text{softmax}(\mathbf{O}_2 \text{relu}(\mathbf{O}_1 \mathbf{z}_t^c)) \\ \mathbf{z}_t^c &= \mathbf{u}_t^c \odot \mathbf{h}_{t-1}^c + (\mathbf{1} - \mathbf{u}_t^c) \odot \mathbf{e}_t^c \\ \mathbf{e}_t^c &= \tanh(\mathbf{r}_t^c \odot (\mathbf{W}_e^c \mathbf{h}_{t-1}) + \mathbf{W}_x^c \mathbf{x}_{t-1})\end{aligned}$$

- where the gates have the form

$$\begin{aligned}\mathbf{r}_t^c &= \text{sigmoid}(\mathbf{W}_{rh}^c \mathbf{h}_{t-1} + \mathbf{W}_{rx}^c \mathbf{x}_{t-1}) \\ \mathbf{u}_t^c &= \text{sigmoid}(\mathbf{W}_{uh}^c \mathbf{h}_{t-1} + \mathbf{W}_{ux}^c \mathbf{x}_{t-1})\end{aligned}$$

- and the input and the history have the form

$$\mathbf{x}_{t-1} = \begin{bmatrix} \mathbf{c}_{t-1} \\ \mathbf{f}_{t-1} \end{bmatrix}; \quad \mathbf{h}_{t-1} = \begin{bmatrix} \mathbf{z}_{t-1}^c \\ \mathbf{z}_{t-1}^f \end{bmatrix}$$

- Initially the coarse elements are generated

$$\mathbf{c}_t \sim P(\mathbf{c}_t)$$

- Equations can be split into distinct parts

$$P(\mathbf{c}_t) = \text{softmax}(\mathbf{O}_4 \text{relu}(\mathbf{O}_3 \mathbf{z}_t^f))$$

$$\mathbf{z}_t^f = \mathbf{u}_t^f \odot \mathbf{h}_{t-1}^f + (\mathbf{1} - \mathbf{u}_t^f) \odot \mathbf{e}_t^f$$

$$\mathbf{e}_t^f = \tanh(\mathbf{r}_t^f \odot (\mathbf{W}_e^f \mathbf{h}_{t-1}) + \mathbf{W}_x^f \mathbf{x}_{t-1} + \mathbf{W}_{uc} \mathbf{c}_t)$$

- where the “fine” gates have the form

$$\mathbf{r}_t^f = \text{sigmoid}(\mathbf{W}_{rh}^f \mathbf{h}_{t-1} + \mathbf{W}_{rx}^f \mathbf{x}_{t-1} + \mathbf{W}_{rc} \mathbf{c}_t)$$

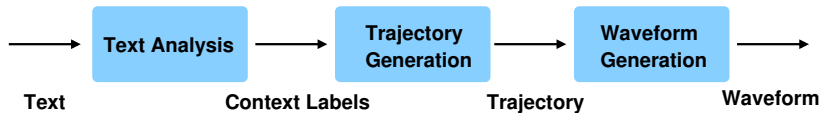
$$\mathbf{u}_t^f = \text{sigmoid}(\mathbf{W}_{uh}^f \mathbf{h}_{t-1} + \mathbf{W}_{ux}^f \mathbf{x}_{t-1} + \mathbf{W}_{uc} \mathbf{c}_t)$$

- input,  $\mathbf{x}_{t-1}$ , and history,  $\mathbf{h}_{t-1}$ , same as coarse

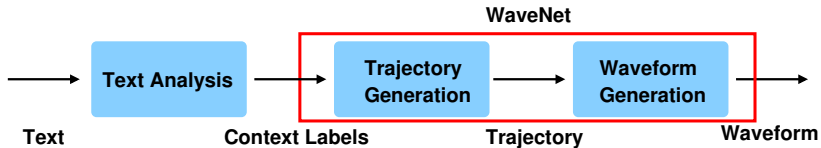
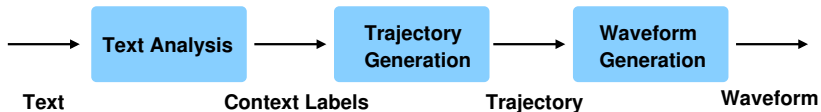
# Conditioning Vectors

- To generate speech require conditioning on text:
  - convert text to context features
  - generate durations for each of model states
  - append “position” features to context
- Yields context features at 200Hz (5ms frame rate)
  - up-sample to yield required (16KHz) context features  $\lambda_{1:T}^{(q)}$
- Requires linguistic information:
  - text processing and POS tags
  - phonetic lexicon (may introduce errors)
  - phone/state duration models

# “Traditional” Synthesis

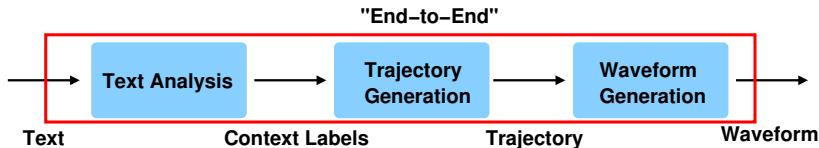
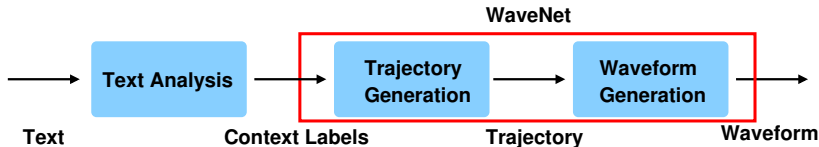
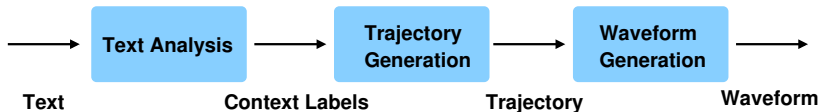


# “Traditional” Synthesis





# "Traditional" Synthesis



# “End-to-End” Synthesis

- Use deep-learning to derive  $\lambda_{1:T}^{(q)}$ 
  - map word sequence  $\omega_{1:L}$  to character sequence  $\mathbf{c}_{1:P}$
  - map character sequence  $\mathbf{c}_{1:P}$  to context features  $\lambda_{1:T}^{(q)}$
  - optimise mapping for synthesis performance
  - note increase in sample rate  $\approx 10\text{-}20\text{ Hz}$  to  $16\text{Hz}$
- To handle difference in required sample rate use **attention**
  - use **decode state**,  $\mathbf{s}_{\tau-1}$  as key for attention at time  $\tau$

$\mathbf{c}_{1:P} \rightarrow \mathbf{h}_{1:P}$  (character encoding)

$$\alpha_{\tau} = \mathcal{F}_{\text{att}}(\mathbf{s}_{\tau-1}, \alpha_{\tau-1}, \mathbf{h}_{1:P}); \quad \mathbf{g}_{\tau} = \sum_{i=1}^P \alpha_{\tau i} \mathbf{h}_i$$

$$\mathbf{y}_{\tau} = \text{generate}(\mathbf{s}_{\tau-1}, \mathbf{g}_{\tau}); \quad \mathbf{s}_{\tau} = \mathcal{F}_{\text{rec}}(\mathbf{s}_{\tau-1}, \mathbf{g}_{\tau}, \mathbf{y}_{\tau})$$

- form of  $\mathcal{F}_{\text{att}}(\cdot)$  very general

- **Content-based** attention - seen before
  - operates between the key and elements of the encoding

$$\alpha_{\tau} = \mathcal{F}_{\text{att}}(\mathbf{s}_{\tau-1}, \mathbf{h}_{1:P}); \quad \alpha_{\tau,i} = \frac{\exp(\mathbf{f}(\mathbf{s}_{\tau-1}, \mathbf{h}_i))}{\sum_{j=1}^P \exp(\mathbf{f}(\mathbf{s}_{\tau-1}, \mathbf{h}_j))}$$
$$\mathbf{f}(\mathbf{s}_{\tau-1}, \mathbf{h}_i) = \mathbf{w}^T \tanh(\mathbf{W}_s \mathbf{s}_{\tau-1} + \mathbf{W}_h \mathbf{h}_i)$$

- **Location-based** attention
  - operates between key and previous attention vector

$$\alpha_{\tau} = \mathcal{F}_{\text{att}}(\mathbf{s}_{\tau-1}, \alpha_{\tau-1});$$
$$\mathbf{f}(\mathbf{s}_{\tau-1}, \alpha_{\tau-1}) = \mathbf{w}^T \tanh(\mathbf{W}_s \mathbf{s}_{\tau-1} + \mathbf{W}_{\alpha} (\mathbf{F} \star \alpha_{\tau-1}))$$

- $\mathbf{F}$  is the matrix for progression (location) of attention
- **Hybrid attention** combines attributes of both

- For sequence data left-to-right emphasis useful
  - reduces chance of inappropriate attention for long sequences
- Use modified version of location attention mechanism
  - convolve with an  $M$ -component mixture of Gaussian functions

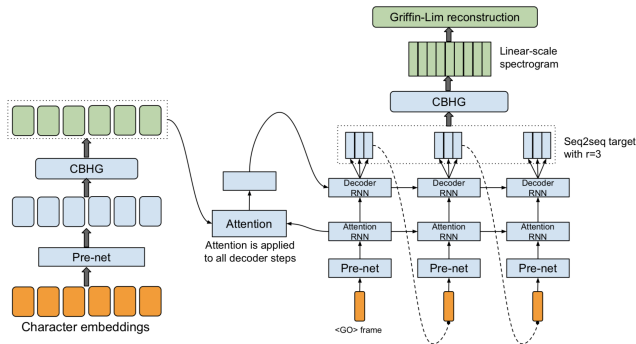
$$\alpha_{\tau i} = \sum_{k=1}^M a_t^{(k)} \exp\left(-b_t^{(k)}(\kappa_t^{(k)} - i)^2\right); \quad \mathbf{g}_{\tau} = \sum_{i=1}^P \alpha_{\tau i} \mathbf{h}_i$$

- where  $\hat{a}_t^{(k)}$ ,  $\hat{b}_t^{(k)}$  and  $\hat{\kappa}_t^{(k)}$  are predicted from the network

$$\begin{aligned} a_t^{(k)} &= \exp(\hat{a}_t^{(k)}); & b_t^{(k)} &= \exp(\hat{b}_t^{(k)}) \\ \kappa_t^{(k)} &= \kappa_{t-1}^{(k)} + \exp(\hat{\kappa}_t^{(k)}) \end{aligned}$$

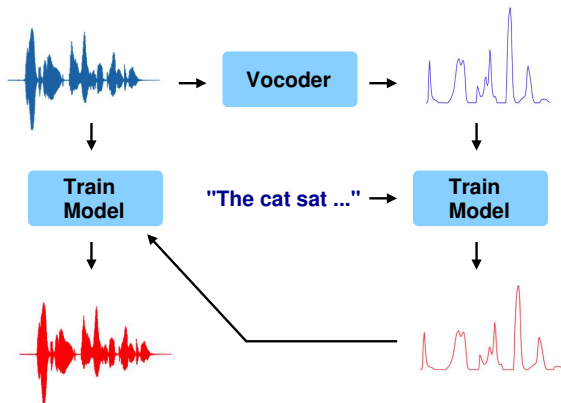
- has been used for end-to-end synthesis and recognition

# Tacotron [23] (Not Waveform-Level)

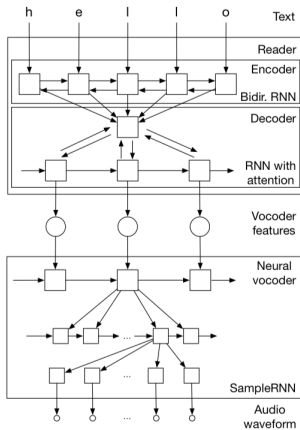


- End-to-end configurations - characters as input
  - content-based attention mechanism
  - CBHG layer added for post-processing
  - predicts linear spectral magnitude (not vocoder parameters)
  - Griffin Lim is used as the synthesiser

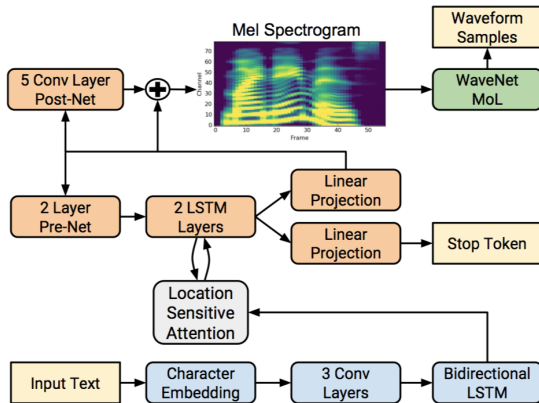
# Trajectory Conditional Models



- Condition neural vocoder on trajectory parameters
  - interpretable layer (the trajectories)
  - leverage adaptation approaches - interesting direction



- Early **end-to-end** architecture
- Uses location attention for context
  - predicts vocoder features
  - neural vocoder for synthesis
- SampleRNN neural vocoder used



- Integrates (modified) Tacotron with (modified) WaveNet



- Changes to the original Tacotron paper:
  - predict Mel-spectrogram not linear
  - simplified network configuration for encoder/decoder
  - location-based attention mechanism
- Changes to the original WaveNet paper:
  - condition on Mel-spectrogram (not linguistic features)
  - uses parallel WaveNet for speed

# Speaker Representations

# Variable Length Mapping

- Range of applications make use of speaker representations
  - speaker clustering
  - speaker recognition/verification
  - speaker adaptation
- Simplest form is multi-speaker (1-of-K coding)
- More generally require a fixed-length representation
  - variable length sequence  $\mathbf{x}_{1:T}^{(s)} \rightarrow$  speaker representation  $\boldsymbol{\lambda}^{(s)}$

$$\boldsymbol{\lambda}^{(s)} = \phi(\mathbf{x}_{1:T}^{(s)})$$

- in 4F10 already seen application using SVMs
  - makes use of [Fisher Kernel](#)
- Can make use of a UBM ( $\boldsymbol{\theta}_{\text{ubm}}$ )
  - large GMM used to represent all speakers
  - MAP adapt model to target speaker  $\boldsymbol{\theta}_{\text{ubm}}^{(s)}$

- From 4F10 lectures

$$\phi(\mathbf{x}_{1:T}) = \begin{bmatrix} \log(p(\mathbf{x}_{1:T}|\boldsymbol{\theta}_{\text{ubm}}^{(s)})) - \log(p(\mathbf{x}_{1:T}|\boldsymbol{\theta}_{\text{ubm}})) \\ \nabla_{\boldsymbol{\theta}} \log(p(\mathbf{x}_{1:T}|\boldsymbol{\theta}))|_{\boldsymbol{\theta}_{\text{ubm}}^{(s)}} \end{bmatrix}$$

- the first term is the standard GMM-based score
- the second term is Fisher score for the speaker model
- only derivatives wrt the mean parameters used
- If only the derivative part is used then

$$\phi(\mathbf{x}_{1:T}) = \begin{bmatrix} \sum_{t=1}^T P(1|\mathbf{x}_t, \boldsymbol{\theta}_{\text{ubm}}^{(s)}) \boldsymbol{\Sigma}_{\text{ubm}}^{(s1)-1} (\mathbf{x}_t - \boldsymbol{\mu}_{\text{ubm}}^{(s1)}) \\ \vdots \\ \sum_{t=1}^T P(M|\mathbf{x}_t, \boldsymbol{\theta}_{\text{ubm}}^{(s)}) \boldsymbol{\Sigma}_{\text{ubm}}^{(sM)-1} (\mathbf{x}_t - \boldsymbol{\mu}_{\text{ubm}}^{(sM)}) \end{bmatrix}$$

- Fisher Information Matrix is sometimes used as a metric

- Rather than taking derivative, it is possible to use parameters
  - consider the means of the speaker adapted UBM

$$\lambda^{(s)} = \phi(\mathbf{x}_{1:T}^{(s)}) = \begin{bmatrix} \mu_{\text{ubm}}^{(s1)} \\ \vdots \\ \mu_{\text{ubm}}^{(sm)} \\ \vdots \\ \mu_{\text{ubm}}^{(sM)} \end{bmatrix}$$

- Both this form and Fisher Kernel yield large spaces
  - if only means used  $M \times d$  elements
  - originally used for SVM-based systems (see 4F10)
- Can we make the speaker information more compact?

- The actual observed data is impacted by multiple factors
  - speaker (desired variability to model)
  - channel/session attributes (not desired)
- Decomposing the mean supervector yields

$$\lambda^{(s)} = \mu_{\text{si}} + \mathbf{V}\lambda_{\text{sp}}^{(s)} + \mathbf{U}\lambda_{\text{ch}}^{(s)} + \mathbf{D}\mathbf{z}$$

- $\mathbf{V}$  and  $\mathbf{U}$  and loading matrices
- $\mu_{\text{si}}$  is the speaker-independent mean
- $\lambda_{\text{sp}}^{(s)}$  point in speaker-space (prior  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- $\lambda_{\text{ch}}^{(s)}$  point in channel/session-space (prior  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- $\mathbf{D}$  the noise matrix,  $\mathbf{z}$  noise term (prior  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- Effectively a large Gaussian distribution: typical dimensions
  - $\lambda^{(s)}$ : 20000;  $\lambda_{\text{sp}}^{(s)}$ : 300;  $\lambda_{\text{ch}}^{(s)}$ : 100
  - iterative training process - see paper

- Identity Vector (iVector): simplify JFA merge speaker/channel

$$\lambda^{(s)} = \mu_{\text{si}} + \mathbf{T} \lambda_{\text{sp}}^{(s)}$$

- $\mathbf{T}$  is the **total variability** matrix
- $\lambda_{\text{sp}}^{(s)}$  point in **speaker-space** (prior  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ )
- This is similar to **Factor Analysis**: use EM
  - unobserved: speaker  $\lambda_{\text{sp}}$ , component at  $t$   $P(m|\lambda_{\text{sp}}, \mathbf{x}_t^{(s)}; \theta)$

$$\begin{aligned} Q(\theta, \hat{\theta}) = & \sum_{s=1}^S \int p(\lambda_{\text{sp}} | \theta, \mathbf{x}_{1:T}^{(s)}) \sum_{t=1}^T \sum_{m=1}^M P(m | \lambda_{\text{sp}}, \mathbf{x}_t^{(s)}; \theta) \\ & \log \left( \mathcal{N}(\mathbf{x}_t^{(s)}; \hat{\mu}_{\text{si}}^{(m)} + \hat{\mathbf{T}}^{(m)} \lambda_{\text{sp}}, \hat{\Sigma}^{(m)}) \right) d\lambda_{\text{sp}} \end{aligned}$$

- new model parameters  $\hat{\theta} = \left\{ \dots, \hat{\mathbf{T}}^{(m)}, \hat{\mu}_{\text{si}}^{(m)}, \hat{\Sigma}^{(m)}, \dots \right\}$
- for simplicity  $P(m|\lambda_{\text{sp}}, \mathbf{x}_t^{(s)}; \theta)$  often fixed for training

- At test-time iVector extracted using

$$\hat{\lambda}_{\text{sp}}^{(s)} = \arg \max_{\lambda_{\text{sp}}} \left\{ p(\lambda_{\text{sp}} | \mathbf{x}_{1:T}^{(s)}, \theta) \right\}$$

- again EM is used to find iVector
- Model related to [CAT](#) and [EigenVoices](#)
  - point estimate of  $\lambda_{\text{sp}}^{(s)}$  used, rather than distribution
  - treated as part of the parameter estimation stage

$$\begin{aligned} \mathcal{Q}(\theta, \hat{\theta}) = & \sum_{s=1}^S \sum_{m=1}^M \sum_{t=1}^T P(m | \lambda_{\text{sp}}^{(s)}, \mathbf{x}_t^{(s)}; \theta) \left[ \log(P(\hat{\lambda}_{\text{sp}}^{(s)})) \right. \\ & \left. + \log \left( \mathcal{N}(\mathbf{x}_t^{(s)}; \hat{\mu}_{\text{si}}^{(m)} + \hat{\mathbf{T}}^{(m)} \hat{\lambda}_{\text{sp}}^{(s)}, \hat{\Sigma}^{(m)}) \right) \right] \end{aligned}$$

- possible to factorise  $\lambda_{\text{sp}}^{(s)}$  (JFA) include [orthogonality constraint](#)



# iVectors for Speaker Recognition

- Extract iVectors for all enrolled speakers,  $\lambda_{sp}^{(1)}, \dots, \lambda_{sp}^{(S)}$ 
  - extract for test speaker  $\lambda_{sp}$
  - need to select “closest” enrolled speaker
- For speed look at distances between iVectors

$$\hat{s} = \arg \min_s \left\{ d(\lambda_{sp}, \lambda_{sp}^{(s)}) \right\}$$

- euclidean distance:

$$d(\lambda_{sp}, \lambda_{sp}^{(s)}) = \|\lambda_{sp} - \lambda_{sp}^{(s)}\|^2$$

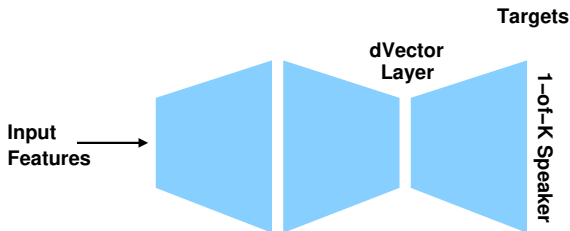
- (-) cosine distance:

$$d(\lambda_{sp}, \lambda_{sp}^{(s)}) = - \frac{\lambda_{sp}^T \lambda_{sp}^{(s)}}{\sqrt{\lambda_{sp}^T \lambda_{sp} \lambda_{sp}^{(s)T} \lambda_{sp}^{(s)}}}$$

popular choice (empirically good!)

# dVector Representation [21]

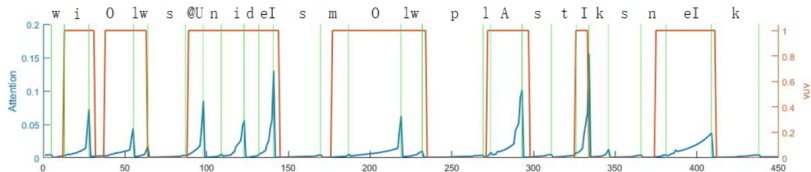
- Train vector to **discriminate** between speakers
  - related to bottleneck features for ASR



- Targets are a 1-of-K coding of speaker
  - wide window of features to yield good performance
- Simple approach used to handle temporal aspect of signal
  - $\lambda_t^{(s)}$  is the vector for frames centered at time  $t$

$$\lambda^{(s)} = \frac{1}{T} \sum_{t=1}^T \lambda_t^{(s)}$$

# Attention-based dVector Representation [22]

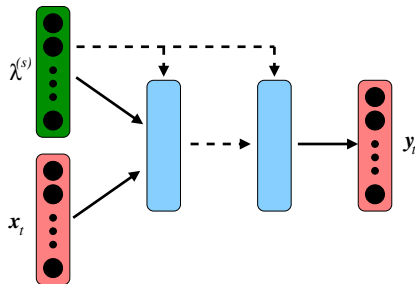


- Rather than using simple averaging of d-Vectors, use attention

$$\lambda^{(s)} = \sum_{t=1}^T \alpha_t \lambda_t^{(s)}$$

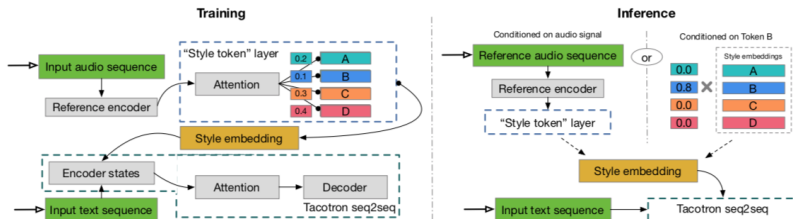
- self-attention mechanism used based on context features
- model trained in an integrated fashion (including attention)
- Attention focuses on voiced regions

# Speaker Representations for Adaptation



- Speaker representation can be used as auxiliary information
  - simple for of speaker adaptation
  - no initial hypothesis required
  - can be optionally be applied to other layers of network

# General Style Tokens [24]



- In addition to speaker, also interested in [style](#)
  - speakers prosody varies with context
  - narrative speech, emotions, audiobooks
- It is hard to consistently label "style"
  - unsupervised approaches increasingly popular
  - [general style tokens](#) for Tacotron one example

- [1] W. Campbell, D. Sturim, D. Reynolds, and A. Solomonoff, "SVM based speaker verification using a GMM supervector kernel and NAP variability compensation," in *Proc. ICASSP*, 2006.
- [2] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," *CoRR*, vol. abs/1506.07503, 2015. [Online]. Available: <http://arxiv.org/abs/1506.07503>
- [3] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Trans. Audio, Speech & Language Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- [4] M. J. F. Gales, "Cluster adaptive training of hidden Markov models," *IEEE Transactions Speech and Audio Processing*, vol. 8, pp. 417–428, 2000.
- [5] A. Graves, "Generating sequences with recurrent neural networks," *CoRR*, vol. abs/1308.0850, 2013. [Online]. Available: <http://arxiv.org/abs/1308.0850>
- [6] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. van den Oord, S. Dieleman, and K. Kavukcuoglu, "Efficient neural audio synthesis," *CoRR*, vol. abs/1802.08435, 2018. [Online]. Available: <http://arxiv.org/abs/1802.08435>
- [7] P. Karanasou, Y. Wang, M. Gales, and P. Woodland, "Adaptation of deep neural network acoustic models using factorised i-vectors," in *Proceedings of Interspeech'14*, 2014.
- [8] P. Kenny, G. Boulianne, P. Ouellet, and P. Dumouchel, "Joint factor analysis versus eigenchannels in speaker recognition," *IEEE Trans. Audio, Speech & Language Processing*, vol. 15, no. 4, pp. 1435–1447, 2007.
- [9] D. P. Kingma, T. Salimans, and M. Welling, "Improving variational inference with inverse autoregressive flow," *CoRR*, vol. abs/1606.04934, 2016. [Online]. Available: <http://arxiv.org/abs/1606.04934>
- [10] J. Koutník, K. Greff, F. J. Gomez, and J. Schmidhuber, "A clockwork RNN," *CoRR*, vol. abs/1402.3511, 2014. [Online]. Available: <http://arxiv.org/abs/1402.3511>
- [11] R. Kuhn, P. Nguyen, J.-C. Junqua, L. Goldwasser, N. Niedzielski, S. Fincke, K. Field, and M. Contolini, "Eigenvoices for speaker adaptation," in *Proceedings ICSLP*, 1998, pp. 1771–1774.
- [12] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.

- [13] Z.-H. Ling, S.-Y. Kang, H. Zen, A. Senior, M. Schuster, X.-J. Qian, H. M. Meng, and L. Deng, "Deep learning for acoustic modeling in parametric speech generation: A systematic review of existing techniques and future trends," *IEEE Signal Processing Magazine*, vol. 32, no. 3, pp. 35–52, 2015.
- [14] C. Longworth and M. Gales, "Derivative and parametric kernels for speaker verification," in *Proceedings InterSpeech*, September 2007.
- [15] S. Mehri, K. Kumar, I. Gulrajani, R. Kumar, S. Jain, J. Sotelo, A. C. Courville, and Y. Bengio, "SampleRNN: An unconditional end-to-end neural audio generation model," *CoRR*, vol. abs/1612.07837, 2016. [Online]. Available: <http://arxiv.org/abs/1612.07837>
- [16] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma, "PixelCNN++: Improving the pixelCNN with discretized logistic mixture likelihood and other modifications," *CoRR*, vol. abs/1701.05517, 2017. [Online]. Available: <http://arxiv.org/abs/1701.05517>
- [17] J. Shen, R. Pang, R. J. Weiss, M. Schuster, N. Jaitly, Z. Yang, Z. Chen, Y. Zhang, Y. Wang, R. J. Skerry-Ryan, R. A. Saurous, Y. Agiomyrgiannakis, and Y. Wu, "Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions," *CoRR*, vol. abs/1712.05884, 2017. [Online]. Available: <http://arxiv.org/abs/1712.05884>
- [18] J. Sotelo, S. Mehri, K. Kumar, J. F. Santos, K. Kastner, A. Courville, and Y. Bengio, "Char2wav: End-to-end speech synthesis," 2017.
- [19] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," *CoRR*, 2016. [Online]. Available: <https://arxiv.org/pdf/1609.03499.pdf>
- [20] A. van den Oord, Y. Li, I. Babuschkin, K. Simonyan, O. Vinyals, K. Kavukcuoglu, G. van den Driessche, E. Lockhart, L. C. Cobo, F. Stimberg, N. Casagrande, D. Grewe, S. Noury, S. Dieleman, E. Elsen, N. Kalchbrenner, H. Zen, A. Graves, H. King, T. Walters, D. Belov, and D. Hassabis, "Parallel wavenet: Fast high-fidelity speech synthesis," *CoRR*, vol. abs/1711.10433, 2017. [Online]. Available: <http://arxiv.org/abs/1711.10433>
- [21] E. Variani, X. Lei, E. McDermott, I. Lopez-Moreno, and J. Gonzalez-Dominguez, "Deep neural networks for small footprint text-dependent speaker verification," *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4052–4056, 2014.

- [22] M. Wan, G. Degottex, and M. J. F. Gales, "Integrated speaker-adaptive speech synthesis," in *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec 2017, pp. 705–711.
- [23] Y. Wang, R. J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, Q. V. Le, Y. Agiomyrgiannakis, R. Clark, and R. A. Saurous, "Tacotron: A fully end-to-end text-to-speech synthesis model," *CoRR*, vol. abs/1703.10135, 2017. [Online]. Available: <http://arxiv.org/abs/1703.10135>
- [24] Y. Wang, D. Stanton, Y. Zhang, R. J. Skerry-Ryan, E. Battenberg, J. Shor, Y. Xiao, F. Ren, Y. Jia, and R. A. Saurous, "Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis," *CoRR*, vol. abs/1803.09017, 2018. [Online]. Available: <http://arxiv.org/abs/1803.09017>
- [25] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *CoRR*, vol. abs/1511.07122, 2015. [Online]. Available: <http://arxiv.org/abs/1511.07122>