

COM4511 Speech Technology: Speech Synthesis

Anton Ragni

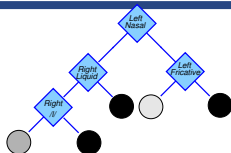
March 11, 2020



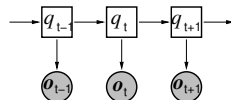
- ▶ One of key interfaces of speech technology
 - ▶ dialogue systems, assistants, synthetic voices and many more!
- ▶ Can be regarded as an **inverse problem** to speech recognition
 - ▶ generate speech given observed latent representation (word sequence)
- ▶ Previous lectures examined multiple generative models
 - ▶ **why do we need to talk about speech synthesis?**



(a) Standard ASR topology



(b) Phonetic decision trees



(c) Dynamic Bayesian Network

- Probability density function of observation sequences given word sequence

$$p(\mathbf{O}_{1:T} | \mathbf{w}_{1:L}) = \sum_{\mathbf{q}_{1:T} \in Q_{1:T}^{(\mathbf{w}_{1:L})}} p(\mathbf{O}_{1:T} | \mathbf{q}_{1:T}) P(\mathbf{q}_{1:T} | \mathbf{w}_{1:L})$$

- Simple form of conditional density if state output distributions are Gaussians

$$p(\mathbf{O}_{1:T} | \mathbf{q}_{1:T}) = \prod_{t=1}^T \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{q_t}, \boldsymbol{\Sigma}_{q_t}) = \mathcal{N}(\mathbf{O}_{1:T}; \boldsymbol{\mu}_{\mathbf{q}_{1:T}}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}})$$

- "synthesise speech" by maximising likelihood or sampling (how?)
 - what is the form of $\boldsymbol{\mu}_{\mathbf{q}_{1:T}}$ and $\boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}$?
- BUT ...

$$\begin{bmatrix} \vdots \\ \mathbf{o}_{t+1} \\ \mathbf{o}_t \\ \mathbf{o}_{t-1} \\ \vdots \end{bmatrix} = \begin{bmatrix} \dots & \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & -\mathbf{I} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & -\mathbf{I} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} & \dots \\ \dots & \mathbf{0} & \mathbf{0} & -\mathbf{I} & \mathbf{0} & \mathbf{I} & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \mathbf{o}_{t+1}^{(s)} \\ \mathbf{o}_t^{(s)} \\ \mathbf{o}_{t-1}^{(s)} \\ \vdots \end{bmatrix}$$

- Introduce **dynamic features** to relax conditional independence assumption

$$\mathbf{o}_t = \begin{bmatrix} \mathbf{o}_t^{(s)} \\ \Delta \mathbf{o}_t^{(s)} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{o}_{t+1}^{(s)} \\ \mathbf{o}_t^{(s)} \\ \mathbf{o}_{t-1}^{(s)} \end{bmatrix}$$

- possible to incorporate higher order "derivatives" (used by all non-NN HMMs)
- Yield linear transformation of the underlying speech parameterisation (e.g. MFCC)

$$\mathbf{O}_{1:T} = \mathbf{W} \mathbf{O}_{1:T}^{(s)}$$

- any statistical model defined over $\mathbf{O}_{1:T}$ rather than $\mathbf{O}_{1:T}^{(s)}$ is **inconsistent** (inc. HMM)
- derive sequence mean and covariance for static distribution $\mathcal{N}(\mathbf{O}_{1:T}^{(s)}; \boldsymbol{\mu}_{\mathbf{q}_{1:T}}^{(s)}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{(s)})$

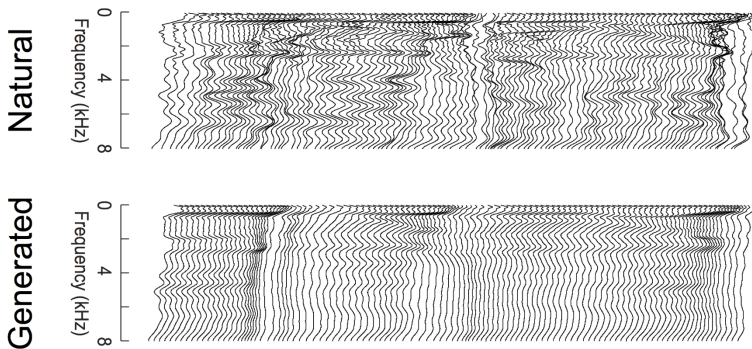
- Express distribution of static features

$$\begin{aligned}\frac{1}{Z_{\mathbf{q}_{1:T}}} \mathcal{N}(\mathbf{O}_{1:T}; \boldsymbol{\mu}_{\mathbf{q}_{1:T}}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}) &= \frac{1}{Z_{\mathbf{q}_{1:T}}} \mathcal{N}(\mathbf{W}\mathbf{O}_{1:T}^{(s)}; \boldsymbol{\mu}_{\mathbf{q}_{1:T}}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{W}\mathbf{O}_{1:T}^{(s)} - \boldsymbol{\mu}_{\mathbf{q}_{1:T}})^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} (\mathbf{W}\mathbf{O}_{1:T}^{(s)} - \boldsymbol{\mu}_{\mathbf{q}_{1:T}})\right) \\ &\propto \exp\left(-\frac{1}{2}\left(\mathbf{O}_{1:T}^{(s)T} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \mathbf{W} \mathbf{O}_{1:T}^{(s)} - \boldsymbol{\mu}_{\mathbf{q}_{1:T}}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \mathbf{W} \mathbf{O}_{1:T}^{(s)} - \mathbf{O}_{1:T}^{(s)T} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \boldsymbol{\mu}_{\mathbf{q}_{1:T}} + \boldsymbol{\mu}_{\mathbf{q}_{1:T}}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \boldsymbol{\mu}_{\mathbf{q}_{1:T}}\right)\right) \\ &= \mathcal{N}\left(\mathbf{O}_{1:T}^{(s)}; \left(\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \mathbf{W}\right)^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \boldsymbol{\mu}_{\mathbf{q}_{1:T}}, \left(\mathbf{W}^T \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{-1} \mathbf{W}\right)^{-1}\right) = \mathcal{N}(\mathbf{O}_{1:T}^{(s)}; \boldsymbol{\mu}_{\mathbf{q}_{1:T}}^{(s)}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{(s)}) = p(\mathbf{O}_{1:T}^{(s)} | \mathbf{q}_{1:T})\end{aligned}$$

- complex mean and covariance structure breaks conditional independence assumptions
- BUT** training preserves these assumptions (inconsistency)
- Maximise probability distribution of static features during training ([Trajectory HMM](#))

$$p(\mathbf{O}_{1:T}^{(s)} | \mathbf{w}_{1:L}) = \sum_{\mathbf{q}_{1:T} \in \mathbf{Q}_{1:T}^{(\mathbf{w}_{1:L})}} p(\mathbf{O}_{1:T}^{(s)} | \mathbf{q}_{1:T}) P(\mathbf{q}_{1:T} | \mathbf{w}_{1:L})$$

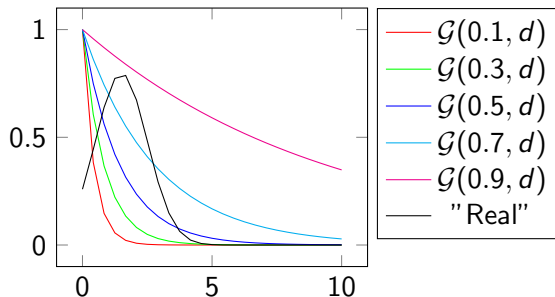
- training and inference expensive due to dependency on the whole sequence



- ▶ Synthesised speech sounds "muffled" (low amplitude broadened formants)
- ▶ Enforce expected variance on generated speech

$$\hat{\mathbf{O}}_{1:T}^{(s)} = \arg \max_{\mathbf{O}_{1:T}^{(s)}} \left\{ \mathcal{N} \left(\mathbf{O}_{1:T}^{(s)}, \boldsymbol{\mu}_{\mathbf{q}_{1:T}}^{(s)}, \boldsymbol{\Sigma}_{\mathbf{q}_{1:T}}^{(s)} \right) \mathcal{N} \left(\frac{1}{T} \sum_{t=1}^T (\mathbf{o}_t^{(s)} - \boldsymbol{\mu})^2; \boldsymbol{\mu}^{(gv)}, \boldsymbol{\Sigma}^{(gv)} \right)^\alpha \right\}$$

- ▶ global variance statistics, $\boldsymbol{\mu}^{(gv)}$ and $\boldsymbol{\Sigma}^{(gv)}$, estimated on training data
- ▶ what does it remind you of?



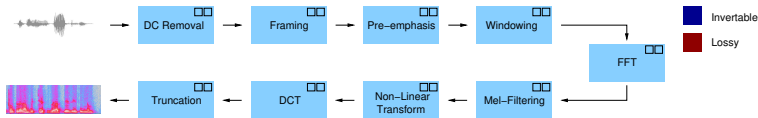
- Probability of latent variable sequence in left-to-right topology

$$P(\mathbf{q}_{1:T} | \mathbf{w}_{1:L}) = a_{1,2}^{d_{1,2}} a_{2,2}^{d_{2,2}} a_{2,3}^{d_{2,3}} \dots a_{K-1,K-1}^{d_{K-1,K-1}} a_{K-1,K}^{d_{K-1,K}}$$

- transition probability $a_{i,j}$ and count $d_{i,j}$
- Probability of staying d times in state j is proportionate to

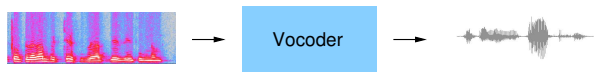
$$P(q_t \neq j, q_{t+1} = j, q_{t+2} = j, \dots, q_{t+d} = j, q_{t+d+1} \neq j) \propto a_{j,j}^d = \mathcal{G}(a_{j,j}, d)$$

- a geometric distribution — very poor model of phone durations
- Need a proper duration distribution — [Hidden Semi-Markov Model](#)



Mel Frequency Cepstral Coefficients (European Standards Organisation ETSI ES 201 108)

- ▶ Standard speech parameterisations are lossy
 - ▶ fundamental frequency F_0 (pitch) filtered out
 - ▶ phase removed due to perceived indifference to human ear
- ▶ Augment spectral features with needed information
 - ▶ care needed due to non-trivial nature of pitch and phase
 - ▶ what have you learnt about pitch and phase?
- ▶ Example: colour code MFCC blocks as invertible/non-invertible, lossy/lossless



- ▶ Standard speech parameterisations are not invertable
 - ▶ impossible to recover the original waveform
- ▶ Multiple options available:
 - ▶ iterative Griffin-Lim algorithm

$$\mathbf{C}_{1:T}^{(n+1)} = \left(\mathbf{G} \mathbf{G}^\dagger \left(\mathbf{A}_{1:T} \odot \exp(\mathbf{j} \angle \mathbf{C}_{1:T}^{(n)}) \right) \right), \quad \mathbf{s}_{1:fT} = \mathbf{G}^\dagger \mathbf{C}_{1:T}^{(N)}$$

- ▶ complex spectrum $\mathbf{C}_{1:T}$ (amplitude $\mathbf{A}_{1:T}$, phase $\angle \mathbf{C}_{1:T}$), \mathbf{G} STFT, \mathbf{G}^\dagger inverse STFT
 - ▶ low-quality phase reconstruction algorithm
- ▶ advanced signal processing techniques (MELP, STRAIGHT)
- ▶ neural networks provide powerful alternative (later)

Type	ASR	TTS	Description
Neighbours	✓	✓	previous/following phones
Position	x [†]	✓	within syllable, word, phrase
Stress/Accent	x [†]	✓	degree, "stress" distances
Linguistic Role	x	✓	POS
Suprasegmental	x	✓	length, tone

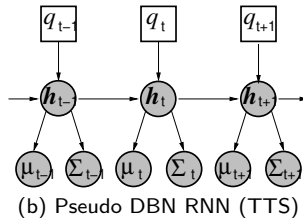
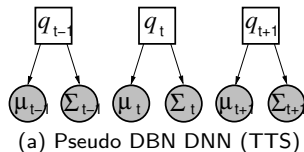
Contextual information usage in ASR and TTS

- ▶ Standard choices of modelling units in speech recognition
 - ▶ phones (/aa/, /ah/, ..., /zh/), graphemes (a, b, ..., y, z)
 - ▶ co-articulation effects handled by context-dependent units (name each type)
- ▶ BUT no expressive characteristic is explicitly modelled — poor synthesis quality
- ▶ Increase specificity of phonetic units

aa^aa-v+dh=ax @ 2_1 / A:1_0_1 / B:1-0-2@1-1&11-3#9-2\$2-1!1-2;8-2|aa / C:0+0+2 / D:content_1 / E:in+1@10+3&7+1#1+2 / F:det_1 / G:0_0 / H:13=12@1=1|L-L% / I:0=0 / J:13+12-1

- ▶ decision trees enable robust estimates BUT do fragment available data (options?)

[†] often used with "limited" resource (non-English) languages



- Predict Gaussian mixture model given current (DNN) or all past (RNN) latent variables

$$p(\mathbf{o}_t | q_t, \dots, q_1) = \sum_{m=1}^M c_{t,m} \mathcal{N}(\mathbf{o}_t; \boldsymbol{\mu}_{t,m}, \boldsymbol{\Sigma}_{t,m})$$

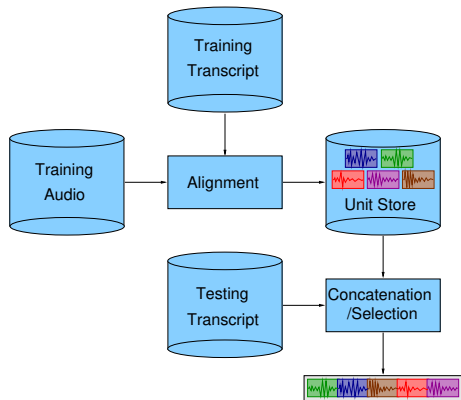
- mean and variance parameters

$$\boldsymbol{\mu}_{t,m} = \phi^{(m)}(\mathbf{A}^{(m)} \mathbf{h}_t + \mathbf{b}^{(m)}), \quad \text{vec}(\boldsymbol{\Sigma}_{t,m}) = \phi^{(\nu)}(\mathbf{A}^{(\nu)} \mathbf{h}_t + \mathbf{b}^{(\nu)})$$

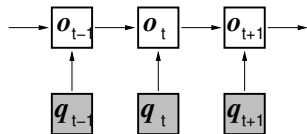
- constraints required to ensure positive semi-definite covariance matrices

- ▶ Role varies among applications
 - ▶ ASR ("law" of large numbers)
 - ▶ one case, "small" vocabulary, optional number and abbr. expansion, punct. removal
 - ▶ TTS ("law" of small numbers)
 - ▶ preserve as much information as possible
- ▶ Examples:

St.	→	Street
St.	→	Saint
Ms	→	Miss
MS	→	Marks and Spencer
MS	→	Microsoft
£1984	→	one thousand, nine hundred and eighty four pounds
1984	→	nineteen eighty four
1984	→	one thousand nine hundred eighty four
- ▶ Accurate text normalisation is context-dependent



- ▶ Align transcripts to audio to create a unit store
 - ▶ need to decide on the nature of units, alignment process and audio representation
- ▶ (If necessary) select and concatenate units to "synthesise" a waveform
 - ▶ need to decide on the nature of selection and concatenation



Discrete elements:

- ▶ latent "observations" \mathbf{o}_t (unit store)
- ▶ observed "latent" variables \mathbf{q}_t (lexical description)

- ▶ Total cost of audio representation for given word sequence

$$\pi(\mathbf{O}_{1:T}, \mathbf{w}_{1:L}) = \sum_{\mathbf{Q}_{1:T} \in \mathcal{Q}_{1:T}^{(\mathbf{w}_{1:L})}} \pi(\mathbf{O}_{1:T}, \mathbf{Q}_{1:T}, \mathbf{w}_{1:L}) = \sum_{\mathbf{Q}_{1:T} \in \mathcal{Q}_{1:T}^{(\mathbf{w}_{1:L})}} \prod_{t=1}^T \pi(\mathbf{q}_t, \mathbf{o}_t) \pi(\mathbf{o}_{t-1}, \mathbf{o}_t)$$

- ▶ target cost $\pi(\mathbf{q}_t, \mathbf{o}_t)$, concatenation cost $\pi(\mathbf{o}_{t-1}, \mathbf{o}_t)$
- ▶ what does this form remind you of?
- ▶ Use dynamic programming to infer optimal observation sequence

$$\hat{\mathbf{O}}_{1:T} = \arg \max_{\mathbf{O}_{1:T}} \left\{ \pi(\mathbf{O}_{1:T}, \mathbf{Q}_{1:T}, \mathbf{w}_{1:L}) \right\}$$

- ▶ use "smart" approaches to reduce search complexity

- ▶ Use advanced sequence models to map linguistic units into waveforms

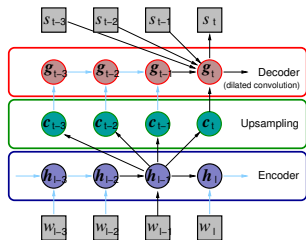
$$\mathbf{w}_{1:L} \longrightarrow \mathbf{s}_{1:fT}$$

- ▶ (discrete) waveform samples $\mathbf{s}_{1:fT}$ typically taken at 16-24 kHz
 - ▶ what is f for 16 kHz and 24 kHz?
- ▶ Alternatively, introduce intermediate latent representation

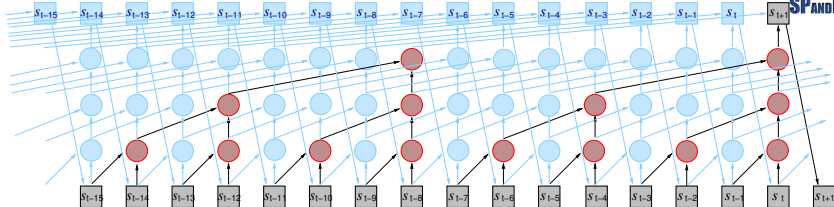
$$\mathbf{w}_{1:L} \longrightarrow \mathbf{Q}_{1:T}$$

$$\mathbf{Q}_{1:T} \longrightarrow \mathbf{s}_{1:fT}$$

- ▶ spectrogram $\mathbf{Q}_{1:T}$ provides interpretable representation
- ▶ In both cases need to decide if intermediate sub-word representation needed
 - ▶ "flat" or "rich" phonemes/graphemes



(a) Pseudo DBN

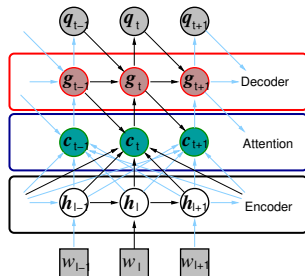


(b) Dilated Convolution

- Predict waveform conditioning on rich linguistic units

$$P(\mathbf{s}_{1:fT} | \mathbf{w}_{1:L}) \approx \prod_{t=1}^{fT} P(s_t | \mathbf{s}_{t-n+1:t-1}, \mathbf{w}_{1:L}) \approx \prod_{t=1}^{fT} P(s_t | \mathbf{g}_t)$$

- dilated convolution handles high speech rate (16 kHz — 24 kHz)
- upsampling layer handles low symbol rate (2 Hz words, 14 Hz phones)
- **BUT** autoregressive-like nature and high rate makes inference/training slow
 - additionally relies on carefully chosen linguistic units and fundamental frequency



Write down:

$$\mathbf{q}_t = \dots$$

$$\mathbf{g}_t = \dots$$

$$\mathbf{c}_t = \dots$$

$$\mathbf{h}_l = \dots$$

- Predict spectrogram (100 Hz) conditioning on phone/grapheme sequence (14 Hz)

$$p(\mathbf{Q}_{1:T} | \mathbf{w}_{1:L}) = \prod_{t=1}^T p(\mathbf{q}_t | \mathbf{Q}_{1:t-1}, \mathbf{w}_{1:L}) \approx \prod_{t=1}^T p(\mathbf{q}_t | \mathbf{g}_t)$$

- attention-based encoder-decoder architecture
- Predict waveform conditioning on spectrogram
 - use any suitable neural vocoder (e.g. WaveNet)

- ▶ Function composition (revisited)

$$\mathbf{y} = \mathbf{f}^{(K)} \odot \mathbf{f}^{(K-1)} \odot \dots \odot \mathbf{f}^{(1)}(\mathbf{x}) = \mathbf{f}(\mathbf{x})$$

- ▶ if $p_{\mathcal{X}}(\mathbf{x})$ is known, which constraints \mathbf{f} must obey to represent valid distribution $p_{\mathcal{Y}}(\mathbf{y})$
- ▶ Change of variables formula

$$p_{\mathcal{Y}}(\mathbf{y}) = p_{\mathcal{X}}(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right|$$

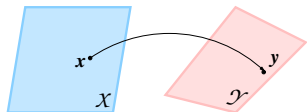
- ▶ if \mathbf{f} is invertible

$$p_{\mathcal{Y}}(\mathbf{y}) = p_{\mathcal{X}}(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{y}^{-1}}{\partial \mathbf{x}} \right) \right| = p_{\mathcal{X}}(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right) \right|^{-1}$$

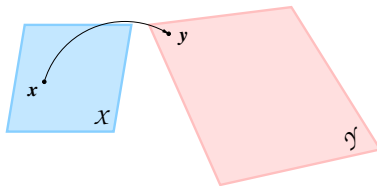
- ▶ Jacobian of composition

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}^{(K)}}{\partial \mathbf{f}^{(K-1)}} \cdot \frac{\partial \mathbf{f}^{(K-1)}}{\partial \mathbf{f}^{(K-2)}} \cdot \dots \cdot \frac{\partial \mathbf{f}^{(1)}}{\partial \mathbf{x}}$$

- ▶ Self-normalising compositions enable efficient sampling and training
 - ▶ provided Jacobian can be efficiently evaluated (conditions?)



(a) Rotation & Contraction



(b) Rotation & Expansion

- ▶ Linear transformation of Gaussian distributed random variable $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$$

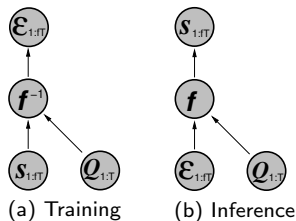
- ▶ inverse transformation (assuming \mathbf{A} invertable)

$$\mathbf{x} = \mathbf{f}^{-1}(\mathbf{y}) = \mathbf{A}^{-1}(\mathbf{y} - \mathbf{b})$$

- ▶ Probability density of transformed variable as a function original density

$$p_Y(\mathbf{y}) = p_X(\mathbf{x}) \left| \det \left(\frac{\partial \mathbf{A}\mathbf{x} + \mathbf{b}}{\partial \mathbf{x}} \right) \right|^{-1} = \frac{1}{|\mathbf{A}|} p_X(\mathbf{x}) = \frac{1}{|\mathbf{A}|} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$$

- ▶ compensate for contraction/expansion (discuss $\mathbf{A} = \text{diag}(2)$, $\mathbf{A} = \text{diag}(\frac{1}{2})$)



Symbols:

$\epsilon_{1:fT}$ – sample

$\mathbf{Q}_{1:T}$ – spectrogram

$\mathbf{s}_{1:fT}$ – waveform

f – invertible composition

Invertible compositions:

- ▶ affine coupling
- ▶ 1×1 convolution
- ▶ many many more!

- ▶ Synthesise speech by passing noise (guess!?) through self-normalising compositions
 - ▶ training

$$p_{\mathcal{E}}(\epsilon_{1:fT}) = p_{\mathcal{S}}(\mathbf{s}_{1:fT}) \left| \det \left(\frac{\partial \mathbf{f}^{-1}(\mathbf{s}_{1:fT}; \mathbf{Q}_{1:T})}{\partial \mathbf{s}_{1:fT}} \right) \right|^{-1}$$

- ▶ inference

$$p_{\mathcal{S}}(\mathbf{s}_{1:fT}) = p_{\mathcal{E}}(\epsilon_{1:fT}) \left| \det \left(\frac{\partial \mathbf{f}(\epsilon_{1:fT}; \mathbf{Q}_{1:T})}{\partial \epsilon_{1:fT}} \right) \right|^{-1}$$

- ▶ Multiple normalising flows examined: inverse auto-regressive flow, WaveGlow

- ▶ Simple example of invertable composition

$$\mathbf{f}(\epsilon_{1:fT}) = \begin{bmatrix} \epsilon_{1:\tau} \\ \epsilon_{\tau+1:fT} \odot \exp(\mathbf{f}^{(a)}(\epsilon_{1:\tau})) + \mathbf{f}^{(b)}(\epsilon_{1:\tau}) \end{bmatrix}, \quad \text{where} \quad \epsilon_{1:fT} = \begin{bmatrix} \epsilon_{1:\tau} \\ \epsilon_{\tau+1:fT} \end{bmatrix}$$

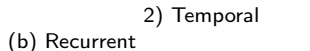
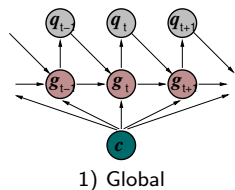
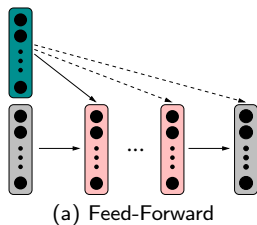
- ▶ scale $\mathbf{f}^{(a)}$ and bias $\mathbf{f}^{(b)}$ compositions need not be invertable
- ▶ Restricted form yields simple Jacobian

$$\frac{\partial \mathbf{f}(\epsilon_{1:fT})}{\partial \epsilon_{1:fT}} = \begin{bmatrix} \mathbf{1}_{1:\tau} & \mathbf{0} \\ \frac{\partial \mathbf{f}(\epsilon_{1:fT})_{\tau+1:fT}}{\partial \epsilon_{1:\tau}} & \text{diag}(\exp(\mathbf{f}^{(a)}(\epsilon_{1:\tau}))) \end{bmatrix}$$

- ▶ Inverse composition

$$\mathbf{f}^{-1}(\mathbf{s}_{1:fT}) = \begin{bmatrix} \mathbf{s}_{1:\tau} \\ (\mathbf{s}_{\tau+1:fT} - \mathbf{f}^{(b)}(\mathbf{s}_{1:\tau})) \odot \exp(-\mathbf{f}^{(a)}(\mathbf{s}_{1:\tau})) \end{bmatrix}$$

- ▶ **Example:** prove that affine coupling is invertable ($\mathbf{f}^{-1}(\mathbf{f}(\epsilon_{1:fT})) = \epsilon_{1:fT}$?)



- ▶ Advanced models so far conditioned on simple symbol sequences (or nothing!)
 - ▶ may need to include other information (fundamental frequency, duration, speaker)
- ▶ Function composition makes conditioning simple
 - ▶ feed-forward unit

$$y = f(x; c)$$

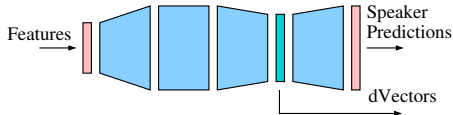
- ▶ recurrent unit

$$g_t = f(g_{t-1}, q_t; c_t)$$

- ▶ BUT need to know how to create conditioning vectors

$$\phi^{(s)}(\mathbf{O}_{1:T}) = \begin{bmatrix} \log\left(\frac{p(\mathbf{O}_{1:T}; \theta^{(s)})}{p(\mathbf{O}_{1:T}; \theta)}\right) \\ \nabla_{\theta} \log(p(\mathbf{O}_{1:T}; \theta))|_{\theta=\theta^{(s)}} \end{bmatrix}$$

(a) Fisher kernel



(b) dVector

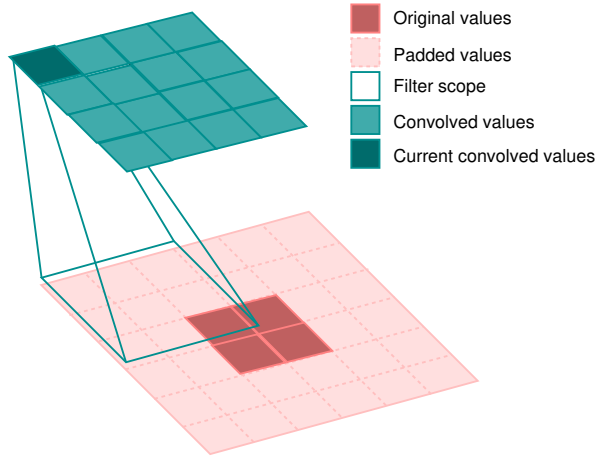
- ▶ Range of applications make use of speaker representations
 - ▶ speaker identification, verification, recognition, clustering, adaptation
- ▶ Typically continuous fixed dimensional representations

$$\mathbf{c}^{(s)} = \phi^{(s)}(\mathbf{O}_{1:T})$$

- ▶ Fisher kernel, joint factor analysis, ?Vectors
- ▶ Consider neural network based dVector representation

$$\mathbf{c} = \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t^{(L-1)}$$

- ▶ alternatively use non-uniform attention weights $\alpha_{1:T}$
 - ▶ what does dVector network topology reminds you of?



- Use transposed convolution ("de-convolution") to learn optimal up-sampling
 - options to choose filter, padding values, stride,

- ▶ General training criterion — Minimum Mean Squared Error

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \theta) &= -\frac{1}{R} \sum_{r=1}^R \sum_{t=1}^{T_r} \log(p(\hat{\mathbf{z}}_t^{(r)} | \hat{\mathbf{z}}_{1:t-1}^{(r)}, \hat{\mathbf{w}}_{1:L_r}^{(r)})) \\ &\propto -\frac{1}{R} \sum_{r=1}^R \sum_{t=1}^{T_r} (\hat{\mathbf{z}}_t^{(r)} - \mathbf{z}_t^{(r)})^\top \Sigma^{-1} (\hat{\mathbf{z}}_t^{(r)} - \mathbf{z}_t^{(r)})\end{aligned}$$

- ▶ reference $\hat{\mathbf{z}}_t^{(r)}$ and model $\mathbf{z}_t^{(r)}$ prediction, noise covariance Σ (issues?)
- ▶ Alternatively, use classification style objective with discrete waveform samples

$$\mathcal{L}(\mathcal{D}; \theta) = -\frac{1}{R} \sum_{r=1}^R \sum_{t=1}^T \log(p(\hat{\mathbf{s}}_t^{(r)} | \hat{\mathbf{s}}_{1:t-1}^{(r)}, \hat{\mathbf{w}}_{1:L_r}^{(r)}))$$

- ▶ discrete reference $\hat{\mathbf{s}}_t^{(r)}$ and quantised model $\tilde{\mathbf{s}}_t^{(r)}$ predictions (issues?)
- ▶ quantise into 256 discrete values (why?) following μ -law compression

$$\tilde{s}_t^{(r)} = \text{int} \left(127 \text{sign}(s_t^{(r)}) \frac{\log(1 + \mu |s_t^{(r)}|)}{\log(1 + \mu)} + 127 \right)$$

- ▶ All generative models examined so far assume some parametric distribution
 - ▶ no matter how flexible not "true" distributions
 - ▶ in many cases unnecessary assumption!
- ▶ Often only samples needed so could forfeit ability to evaluate probability density
 - ▶ need to be able to judge samples as "good" or "bad" to train sample generator
- ▶ Pit sample generator against sample discriminator in a minimax "game"

$$\mathcal{L}(\boldsymbol{\theta}^{(d)}, \boldsymbol{\theta}^{(g)}) = \mathcal{E}_{\mathbf{O} \sim p_{\mathbf{O}}} \left\{ \log(d(\mathbf{O}; \boldsymbol{\theta}^{(d)})) \right\} + \mathcal{E}_{\mathbf{N} \sim p_{\mathbf{N}}} \left\{ \log(1 - d(\underbrace{g(\mathbf{N}; \boldsymbol{\theta}^{(g)})}_{\tilde{\mathbf{O}}}; \boldsymbol{\theta}^{(d)})) \right\}$$

- ▶ true distribution $p_{\mathbf{O}}$ assumed to exist but not explicitly modelled (training data)
 - ▶ noise \mathbf{N} sampled from noise distribution $p_{\mathbf{N}}$ converted to speech $\tilde{\mathbf{O}}$ by generator g
 - ▶ discriminator d classifies true and noise samples, generator makes this task harder
- ▶ **BUT** training such adversarial networks is not trivial
 - ▶ initialisation, stability, etc.

- ▶ One possible form of GAN objective function for sequence data

$$\mathcal{L}_{\text{gan}}(\mathcal{D}; \theta^{(d)}, \theta^{(g)}) = \frac{1}{R} \sum_{r=1}^R \sum_{t=1}^{T_r} \log(\sigma(d(\mathbf{o}_t^{(r)}; \theta^{(d)}))) + \frac{1}{R'} \sum_{r'=1}^{R'} \sum_{t'=1}^{T_{r'}} \log(1 - \sigma(d(\tilde{\mathbf{o}}_{t'}^{(r')}; \theta^{(d)})))$$

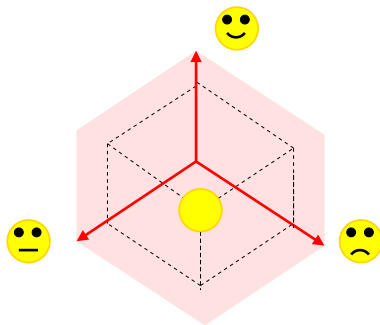
- ▶ samples can be drawn from parametric and nonparametric generators
 - ▶ average per-frame classification accuracy
- ▶ Incorporate GAN objective function in a soft fashion

$$\mathcal{L}(\mathcal{D}; \theta^{(d)}, \theta^{(g)}) = \mathcal{L}_{\text{mge}}(\mathcal{D}; \theta^{(g)}) + \alpha \mathcal{L}_{\text{gan}}(\mathcal{D}; \theta^{(d)}, \theta^{(g)})$$

- ▶ minimum generation error objective function

$$\mathcal{L}_{\text{mge}}(\mathcal{D}; \theta^{(g)}) = \frac{1}{R} \sum_{r=1}^R \|\mathbf{o}_{1:T_r}^{(r)} - \tilde{\mathbf{o}}_{1:T_r}^{(r)}\|_2^2$$

- ▶ Many other options available



- ▶ **Paralinguistic information** crucial for naturalistic speech generation
 - ▶ emotions, speaking style, attitude
- ▶ Controlled modification complicated
 - ▶ not clear how to define expressive "state"
 - ▶ neural network adaptation challenging

- ▶ Objective measures
 - ▶ unlike speech recognition lacks commonly accepted measure
 - ▶ distortion (and other signal processing) style measures only indicative
- ▶ Subjective measures
 - ▶ Mean Opinion Score (MOS): Likert scale (0-5)
 - ▶ Preference Tests: better, worse, no preference
 - ▶ ABX Tests: distance to supplied reference (closer, further)
 - ▶ Transcription Test: transcribe what was said (and compare!)
- ▶ Issues with subjective measures
 - ▶ need to formulate precise criteria
 - ▶ need to ensure the criteria have been followed
 - ▶ can you name any other issue?

- ▶ These two lectures explored speech synthesis
 - ▶ hidden Markov models (HMM)
 - ▶ concatenative (unit selection) speech synthesis
 - ▶ advanced neural network approaches
 - ▶ evaluation
- ▶ Focus on issues surrounding generating realistic speech using HMMs
 - ▶ inconsistency
 - ▶ duration modelling
 - ▶ smoothness
 - ▶ vocoder
 - ▶ modelling units and text normalisation
- ▶ Also discussed advanced forms of neural networks
 - ▶ WaveNet and Tacotron
 - ▶ Normalising Flows
 - ▶ Generative Adversarial Networks