

Speech Technologies

From basics to trends

HUAWEI TECHNOLOGIES CO. LTD (Russia)

Mikhail Kudinov

www.huawei.com

Contents

- 
- 1. Introduction**
 - 2. Speech representation and perception**
 - 3. Speech synthesis**
 - 4. Speech recognition**

Introduction

SPEECH TECHNOLOGIES

Speech Technologies in a Nutshell

- **Speech technology** is a type of computing technology that enables an electronic device to recognize, analyze and understand spoken word or audio containing human speech
- Subfields:
 - ⇒ Speech synthesis
 - ⇒ Speech recognition
 - ⇒ Speaker recognition/verification
 - ⇒ Speech coding
 - ⇒ Voice conversion

Speaker recognition

- **Speaker identification:** identify speaker from a known set by voice in the record
- **Speaker verification:** identify if the voice in the record belongs to the target speaker
 - ⇒ Voice biometry/access control for Barclays
- **Speaker diarization:** segment a long speech record into parts according to which person is speaking now

Speech synthesis / voice conversion

- **Speech synthesis** in narrow sense is a generation of a speech signal based on text input
- **Voice conversion** is a transformation of a speech signal preserving linguistic information
- **Accessibility** for visual impaired users
 - ⇒ Microsoft Windows accessibility screen reader
- **Text readers**
 - ⇒ Xinhua news presenter
- **Accessibility** for speech-impaired users
 - ⇒ Deepmind's initiative for speech-impaired users
- **Voice assistants:**
 - ⇒ Alexa, Siri, Cortana, Bixby...

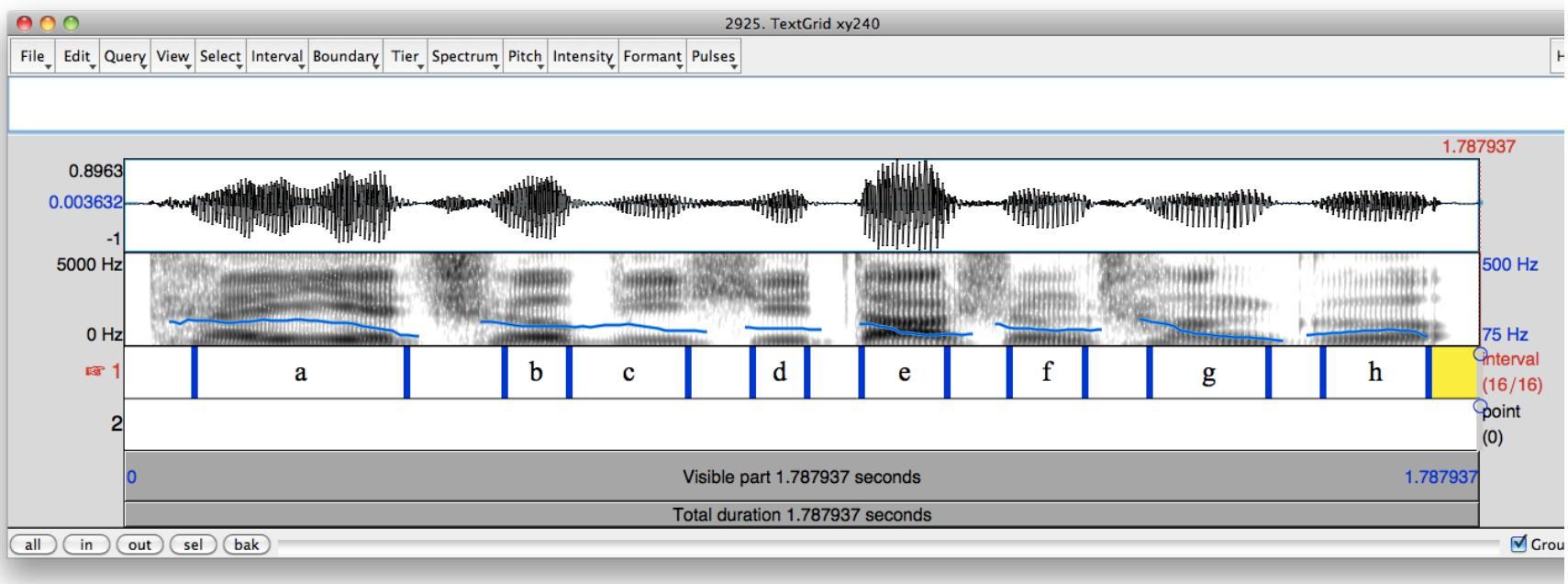
Speech recognition

- **Speech recognition** is a translation of a spoken utterance into text
- **Key-word spotting:**
 - ⇒ Wakeup phrases: Okay, Google!
- **Small vocabulary SR:**
 - ⇒ Voice dial
- **Large-vocabulary SR:**
 - ⇒ Speaker dependent: call center solution for Barclays
 - ⇒ Speaker independent: voice commands, voice search, voice assistants, dictation task

Speech representation and perception

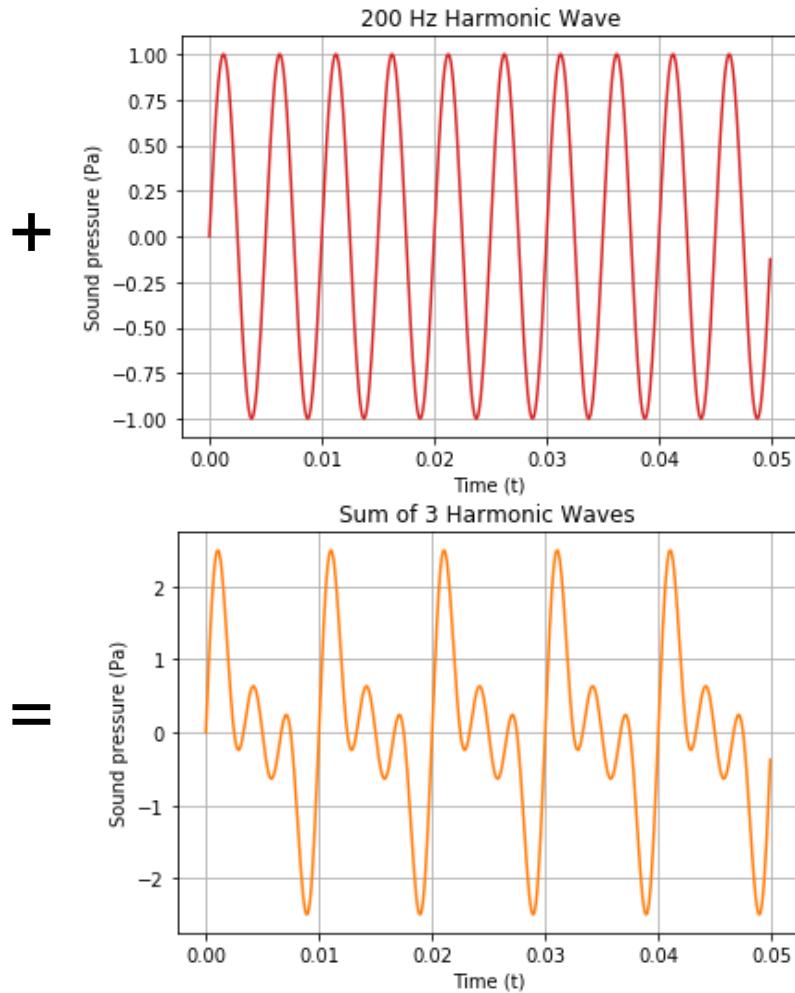
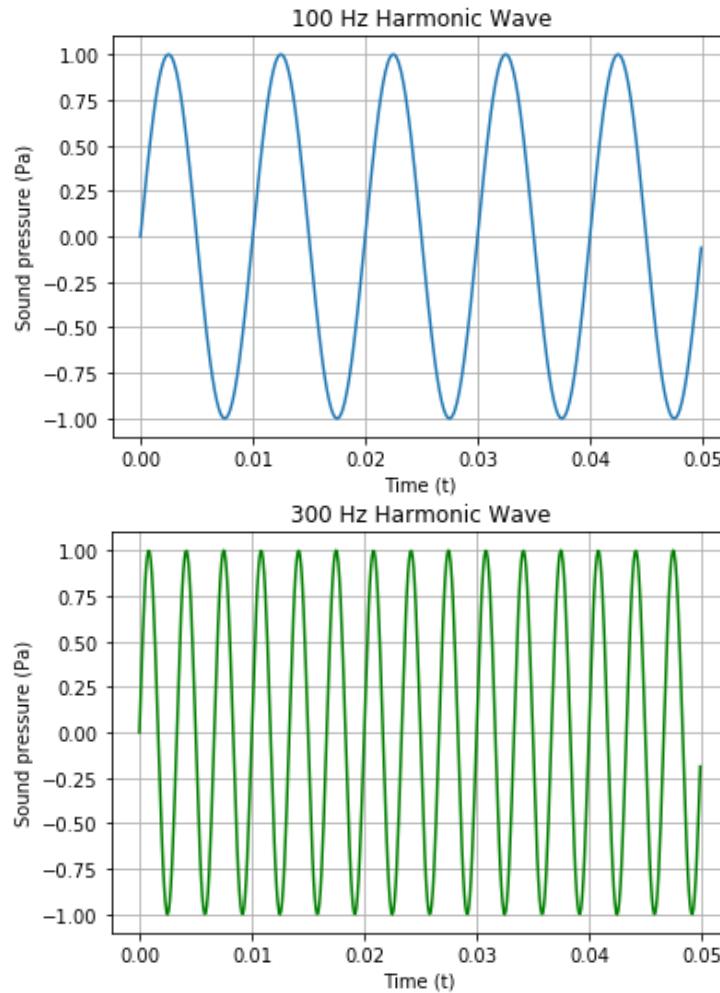
WHAT IS SPEECH?

Waveform – Spectrogram – Text

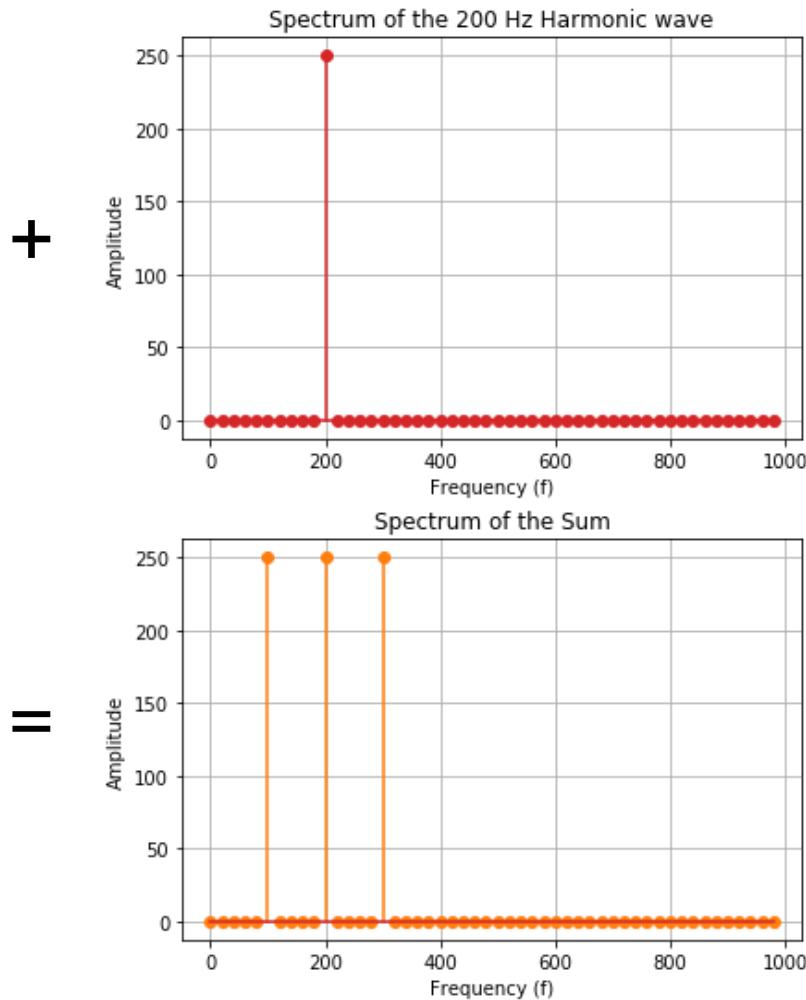
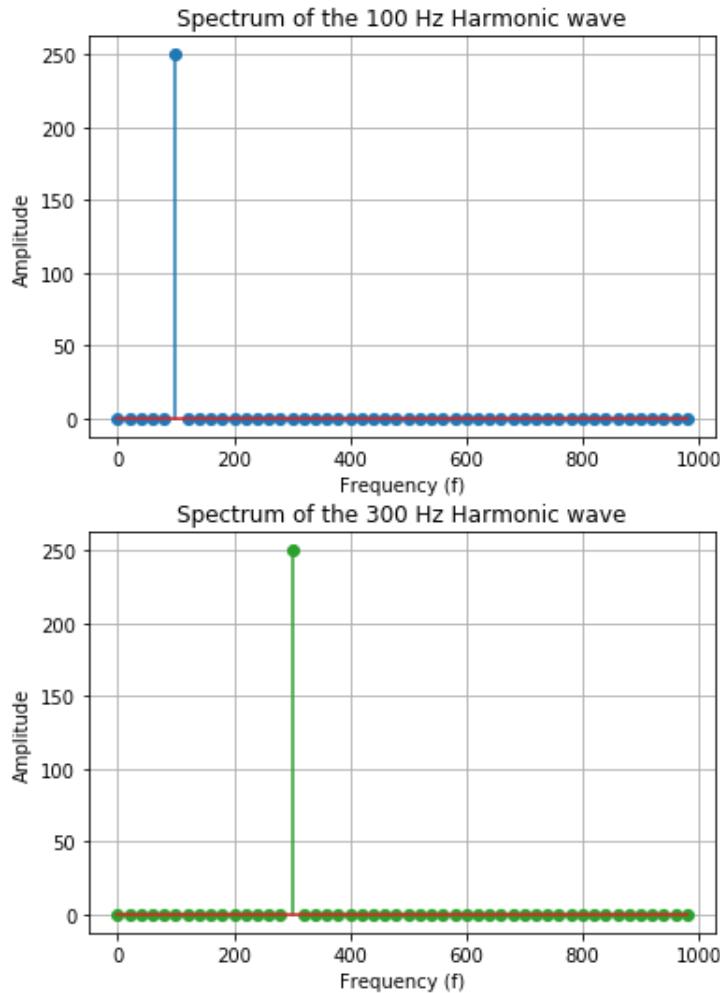


- Waveform: sound intensity as a function of time
- Spectrogram: short-time spectrum as a ‘function’ of a time step
- Text: linguistic meaning corresponding to a speech fragment

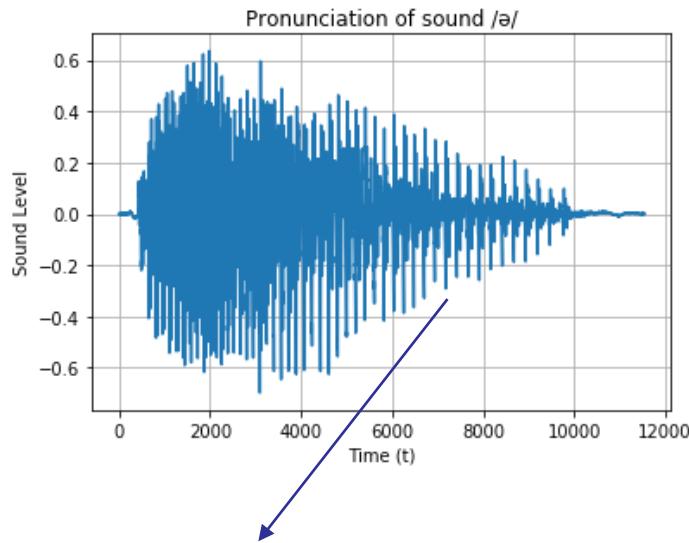
From sound wave to spectrogram: harmonic oscillations



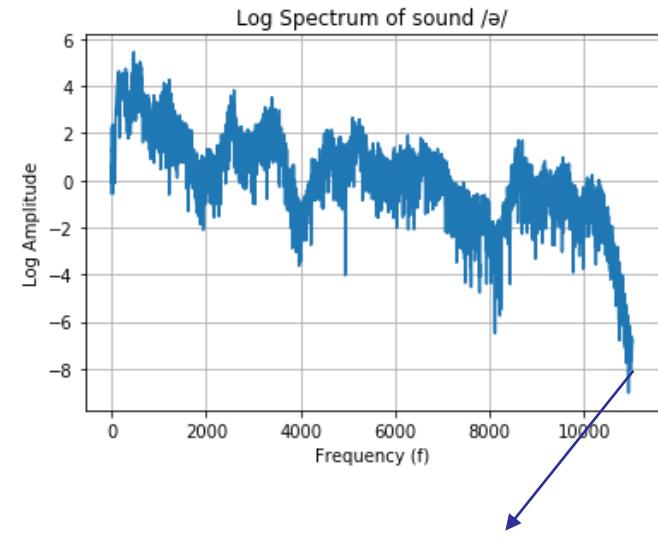
From sound wave to spectrogram: sum of spectra



From sound wave to spectrogram: phoneme /ə/



Discrete Fourier Transform (DFT)*



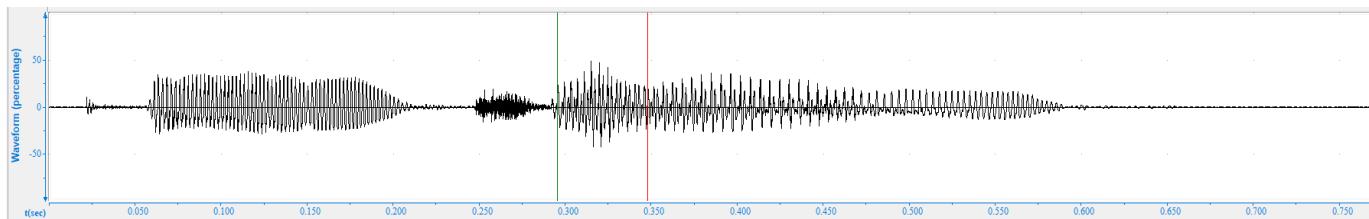
Digital signal is time-discretized so that each of the signal is represented by N samples. Number of samples per second is called *sampling rate*

11 kHz is the highest frequency (*Nyquist frequency*) which can be present in the signal at 22 kHz sampling rate

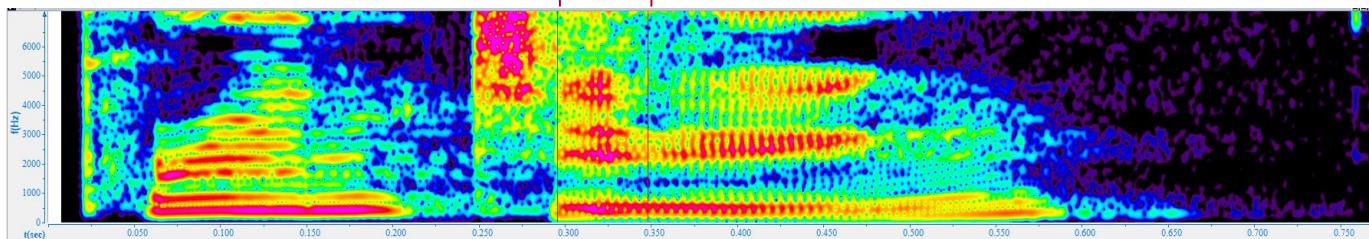
*actually this is half of $\log |FFT(x)|$

From sound wave to spectrogram: Short-time Fourier Transform

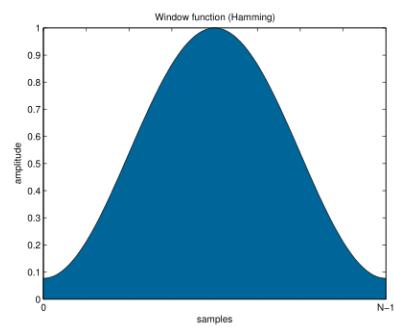
Spectrum (Fast Fourier Transform)



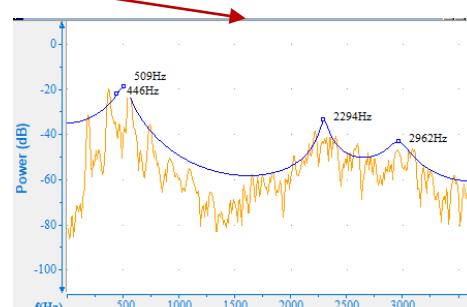
Waveform



Spectrogram



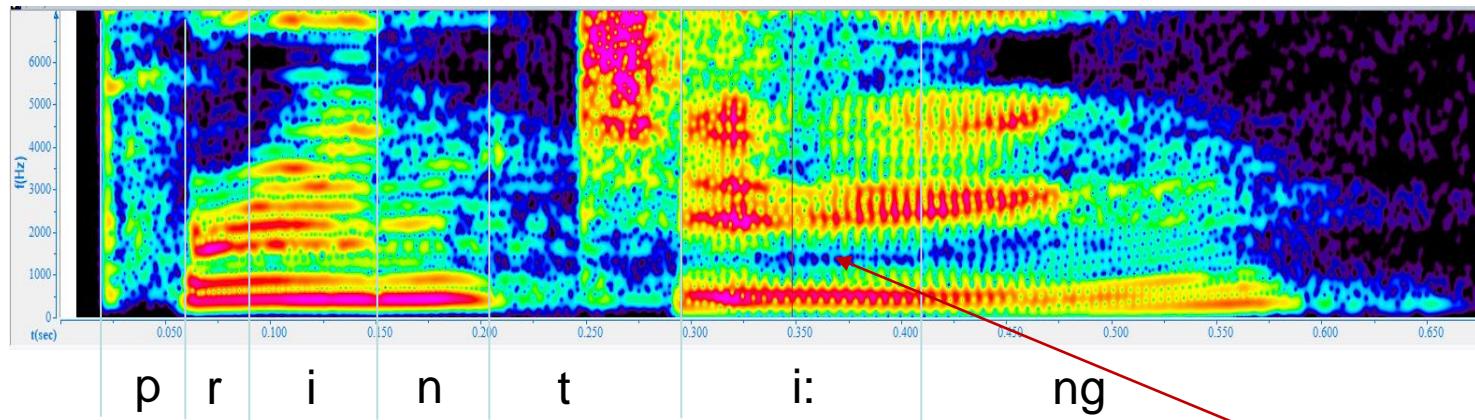
Windowing
function is used
to generate
'smooth' STFT



Fast Fourier Transform
(10-30 msec segment)

Spectrum of a
segment

Closer look at the spectrogram



p r i n t

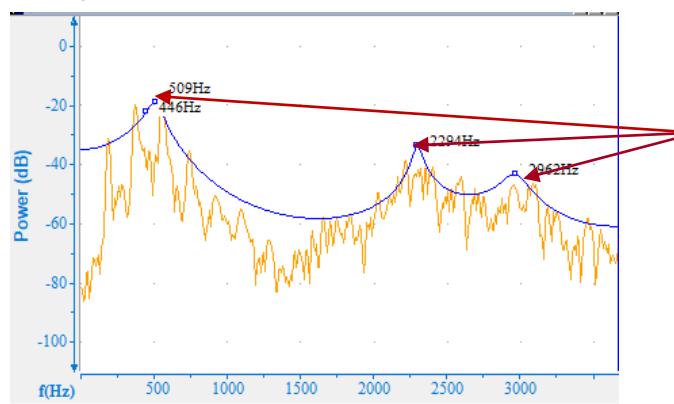
Stop sounds /p/, /t/ have a short pause and a noise segment

English phoneme /i:/

i:

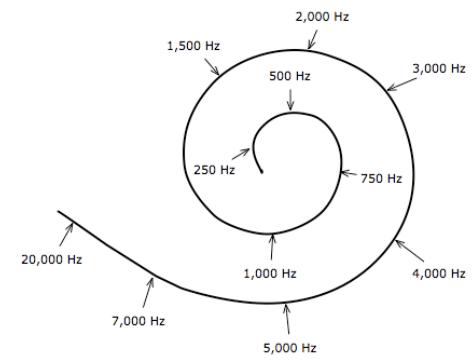
ng

There is no clear border between sounds
But the neighboring sounds affect spectra of each other
Voiced sounds have high peaks in a low frequency range



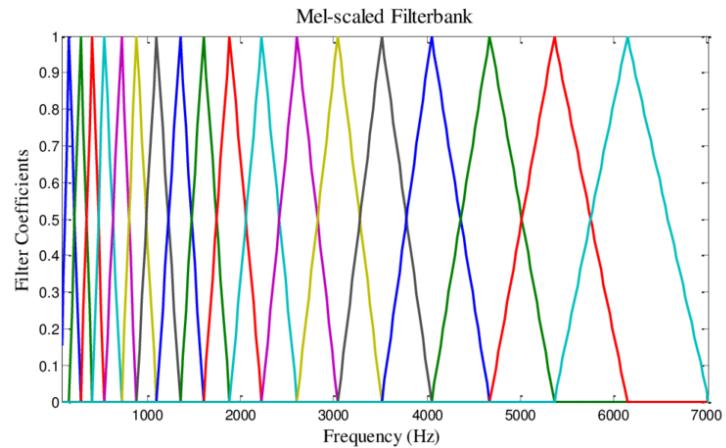
Formant frequencies

Cochlea carries out Fourier Transform in a human inner ear

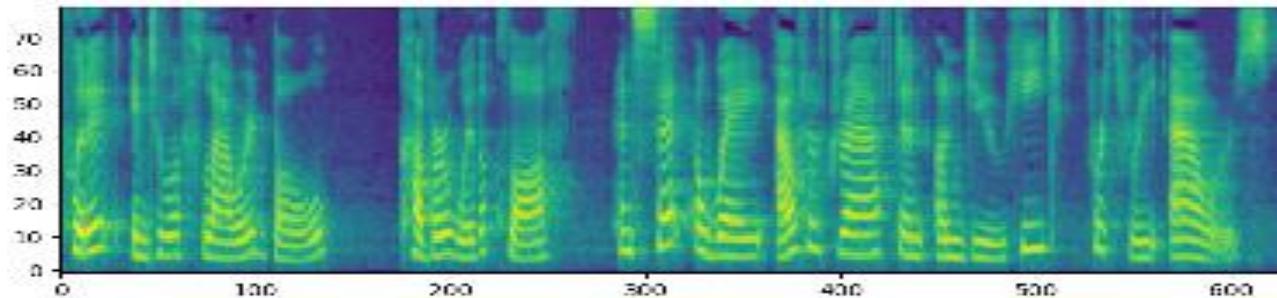


Perceptually motivated features: Mel-scale

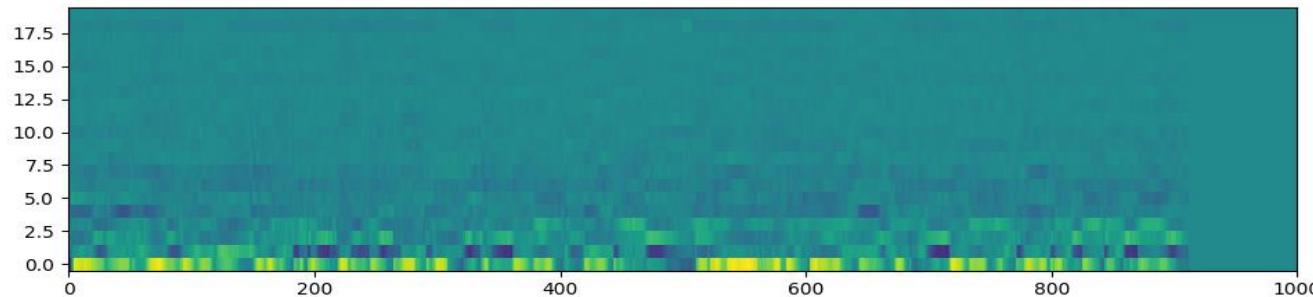
- FFT coefficients are redundant
- Sound perception by the human auditory system is highly non-linear
- Idea #1: compress spectrum into a small number of energies in critical subbands
 - ⇒ Mel-scale spectral coefficients (MFSC)
 - ⇒ Bark-scale spectral coefficients
- Idea #2: decorrelate $\log(\text{MFSC})$ or $\log(\text{BFSC})$ in PCA-like manner (use Discrete Cosine Transform)
 - ⇒ Mel-frequency cepstral coefficients: MFCC
 - ⇒ Bark-frequency cepstral coefficients: BFCC



Acoustic representations: MFSC vs. BFCC

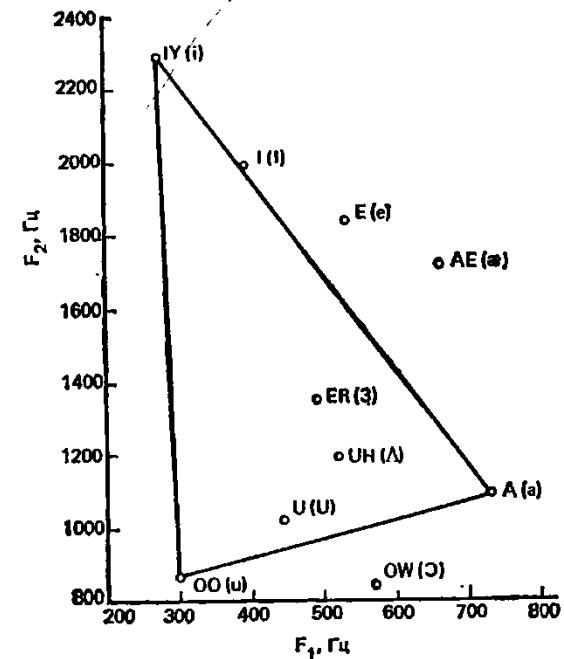
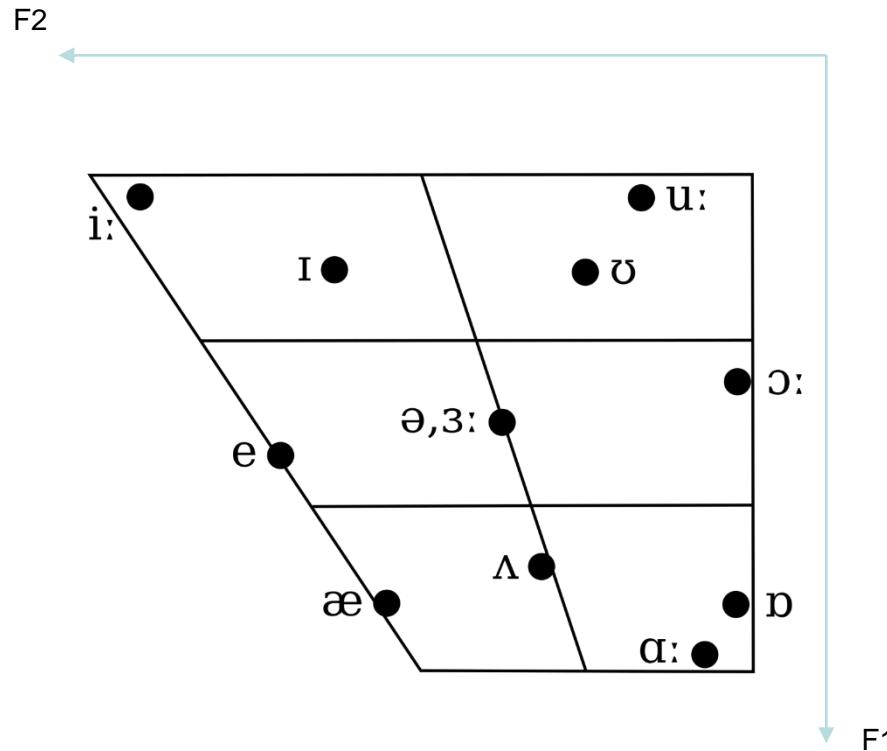


80-band MFSC



18-band Bark-scale cepstral coefficients

From spectrogram to text: phonetics (vowels)



From spectrogram to text: phonetics (consonants)

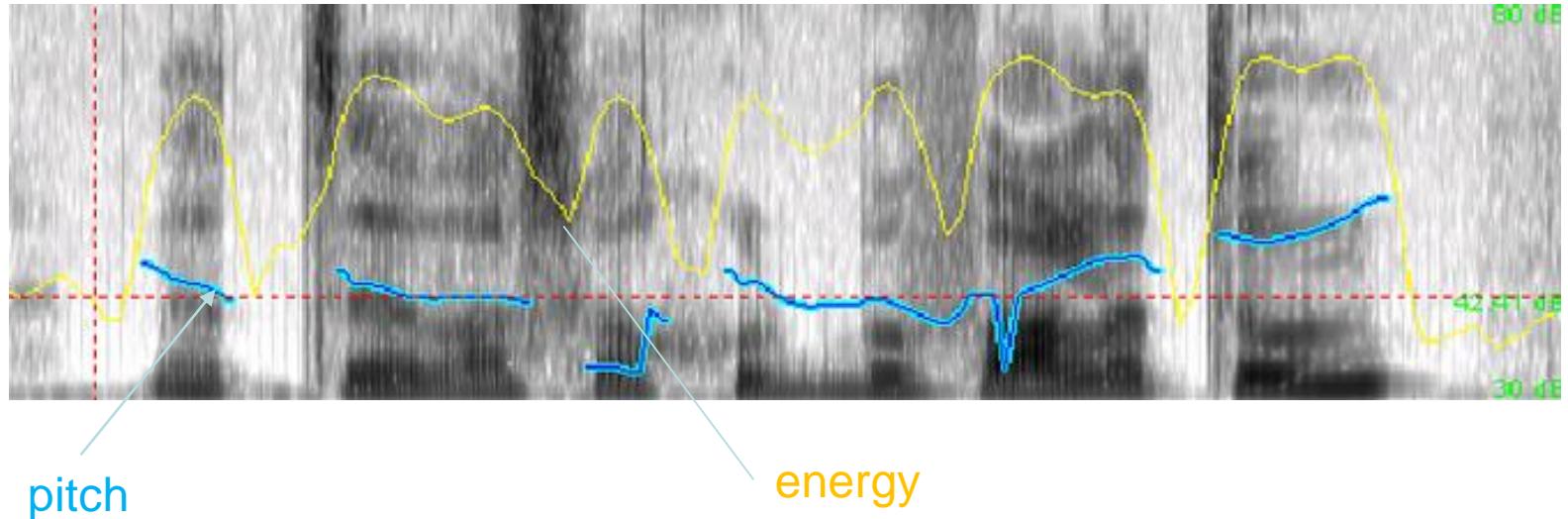
the international phonetic alphabet (2005)

consonants (pulmonic)	LABIAL		CORONAL				DORSAL				RADICAL		LARYNGEAL
	Bilabial	Labio-dental	Dental	Alveolar	Palato-alveolar	Retroflex	Alveolo-palatal	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	n̪		n		ɳ		j̪	ɳ	ɳ			
Plosive	p b			t d		t̪ d̪		c ɟ	k g	q ɣ		ʔ ʡ	ʔ ʡ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ڇ	ç ڇ	x ɣ	χ ڻ	h ڻ	h ڻ	h ڻ
Approximant		v		ɿ		ɬ		j	ɺ				
Tap, flap		v̚		r̚		ɬ̚							
Trill	B			r							R		R̚
Lateral fricative			ɬ ɭ		ɬ̚		ɻ		ɺ				
Lateral approximant			l̚		l̚		ɻ̚		ɺ̚				
Lateral flap			ɶ		ɶ̚								

Where symbols appear in pairs, the one to the right represents a modally voiced consonant, except for murmured h.

Shaded areas denote articulations judged to be impossible. Light grey letters are unofficial extensions of the IPA.

Beyond text: speech prosody



- Speech is not just a sequence of phonemes. The *supersegmental* features of the utterance are extremely important
 - Pitch (fundamental frequency)
 - Loudness (energy)
 - Rhythm (distribution of lengths of the spoken sounds)

Speech synthesis

TEXT-TO-SPEECH

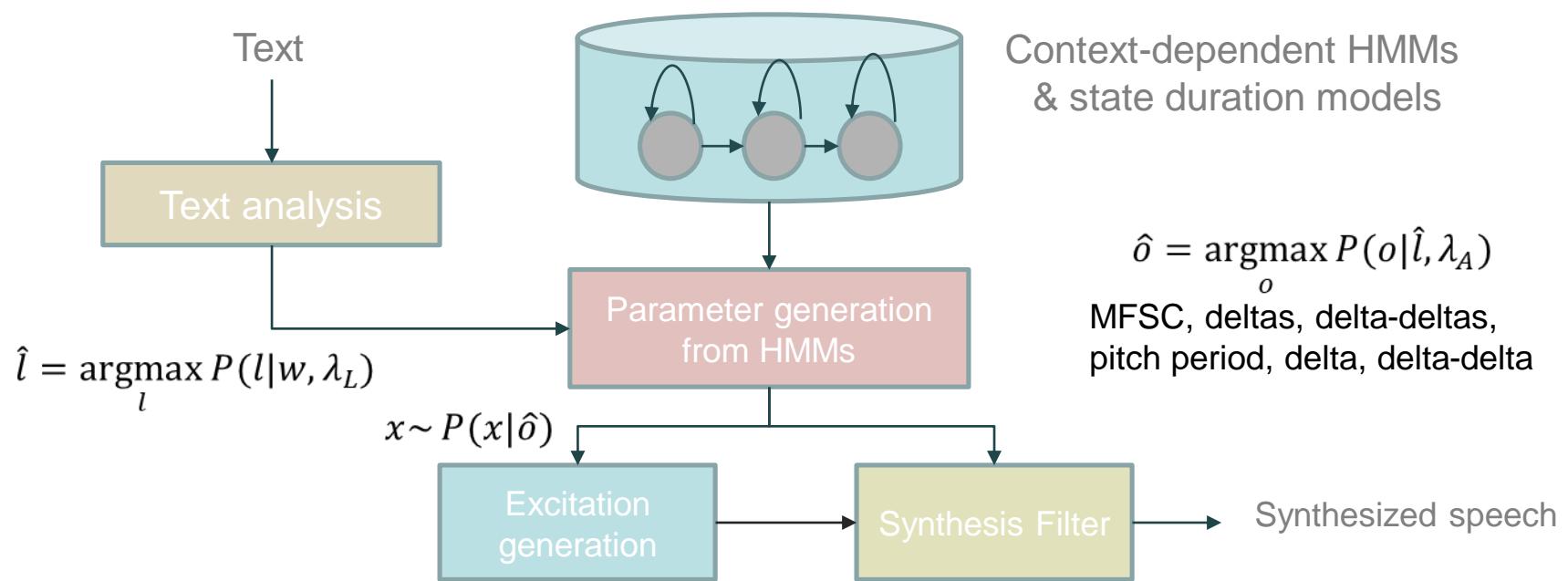
Speech synthesis approaches

- Rule-based, formant synthesis
 - ⇒ Phonetic units are generated according to hand-crafted rules
- Corpus-based, concatenative synthesis
 - ⇒ Concatenate speech units from a database
 - Diphone synthesis
 - Unit selection synthesis
- Corpus-based statistical synthesis
 - ⇒ Feature-generation+vocoder
 - HMM
 - DNN
 - ⇒ E2E systems

Statistical speech synthesis

- Given a training database consisting of pairs of speech waveforms X and texts W estimate a statistical model parameters λ for generating speech waveform x for text $w \notin W$: $x \sim P(x|w, \lambda)$
- Usually $P(x|w, \lambda)$ is decomposed into submodules
 - $\Rightarrow P(x|w, \lambda) = p(x|f) p(f|I, \lambda_A) p(I|w, \lambda_L)$
 - $\Rightarrow f$: parameteric representation of speech waveform x
 - $\Rightarrow I$: linguistic feature
 - $\Rightarrow \lambda = \{\lambda_A, \lambda_L\}$: generative model parameter
 - λ_A : acoustic model parameter
 - λ_L : linguistic model parameter

Statistical speech synthesis before Tacotron era

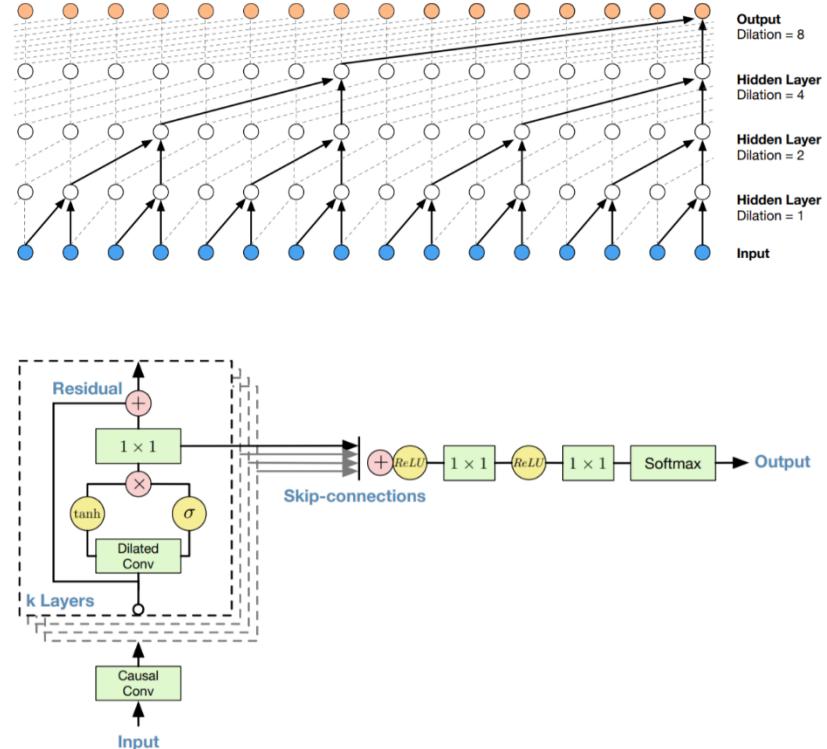


Neural vocoder: WaveNet [Oord17]

- Autoregressive model based
- Main building element: residual block with dilated convolutions
- 30 residual blocks
- Predicts mu-law companded outputs

$$\Rightarrow \frac{\ln(1+\mu|x|)}{1+\ln(1+\mu)}$$

- 16kHz sound
- Best RTF is ~3.0 (Nvidia P100) [Kal18]

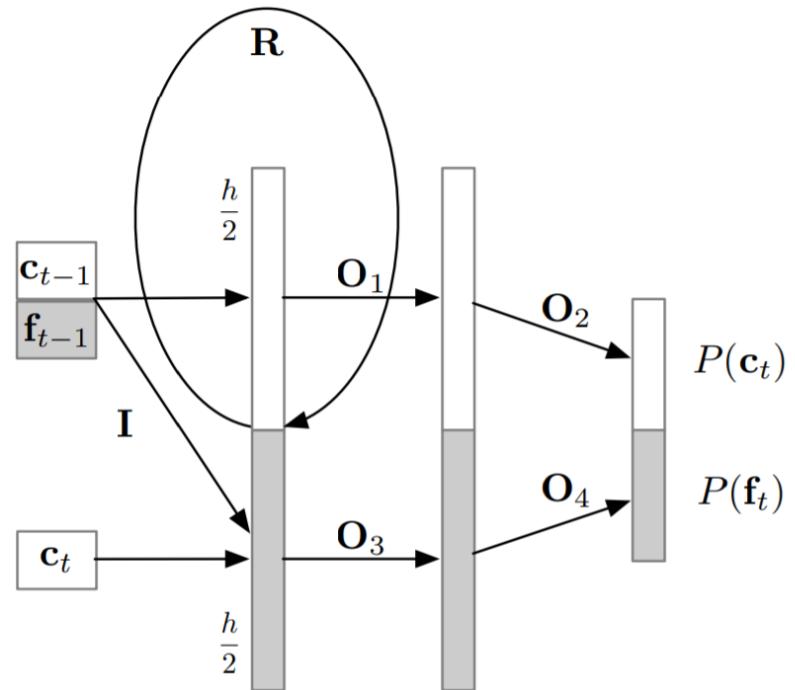


Parallel WaveNet, ClariNet (and WaveGlow)

- Idea #1: use trained autoregressive model as a reference distribution (teacher) to train a non-autoregressive vocoder (based on IAF) via KL-divergence minimization
- Parallel WaveNet [Oord&Li17]
 - ⇒ MoL
 - ⇒ Power loss, Style loss, Contrastive loss
- ClariNet [Ping19]
 - ⇒ Normal distribution
 - ⇒ STFT loss (almost like power loss)
 - ⇒ Tricks with KL divergence regularization
- WaveGlow [Prein18] uses Glow network and is trained from scratch

Recurrent Vocoder: WaveRNN

- Recurrent neural network conditioned on linguistic features (so, like original WaveNet it's not actually a vocoder) [Kal18]
- Quality comparable to WaveNet but RTF=0.25 on GPU at 24kHz
- 16-bit prediction via 2-step prediction: coarse (first 8-bit) and fine (second 8-bit) parts
- Block sparsification



LPCNet: speech coding meets deep learning

- We can estimate a_k by minimizing energy of the residual

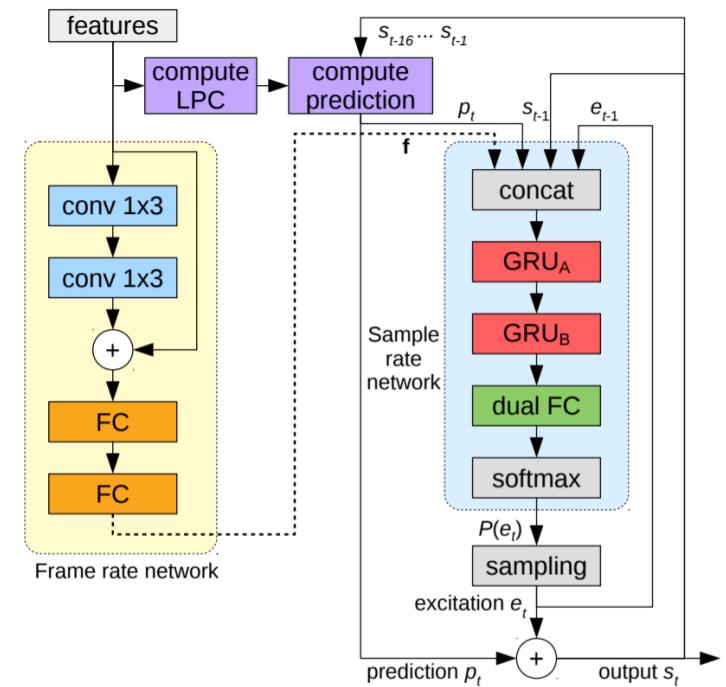
$$E = \sum_n e^2(n) = \sum_n (s(n) - \tilde{s}(n))^2 = \sum_n (s(n) - \sum_{k=1}^p a_k s(n-k))^2$$

⇒ which is equivalent to MSE minimization in the analysis window

- LPCs can be found in both time and frequency domain i.e. from spectrum
- If we know $e(n)$ and a_k for the speech frame we can restore $s(n)$
- Small module of $e(n)$ is a resource of compression in speech codecs

LPCNet

- Compute LPCs from Bark-scale cepstral coefficients
- Use frame-rate network to compute feature vector f
- Feed f to a sample rate recurrent network predicting excitation signal $e(n)$
- Output signal is computed as if it was restored from true LPCs and excitation signal
- Operates on 8kHz, 8-bit mu-law
- RTF: ~0.2 on CPU



Val[19]

HiFi-GAN [Kong20]

- Issue:
 - ⇒ GAN's improve the sampling efficiency/memory usage, but the quality is still not acceptable
- Problem statement:
 - ⇒ Propose GAN, that will generate high-fidelity waveforms comparable with AR and Flow-based models
- Solution:
 - ⇒ Modelling periodic patterns of audio for enhancing quality

HiFi-GAN: Architecture (generator)

- Generator:

- fully convolutional
- upsamples mel to match the length of audio
- ConvTranspose + Multi-Receptive Field Fusion module (MRF)
- MRF: resblock with diverse receptive field patterns

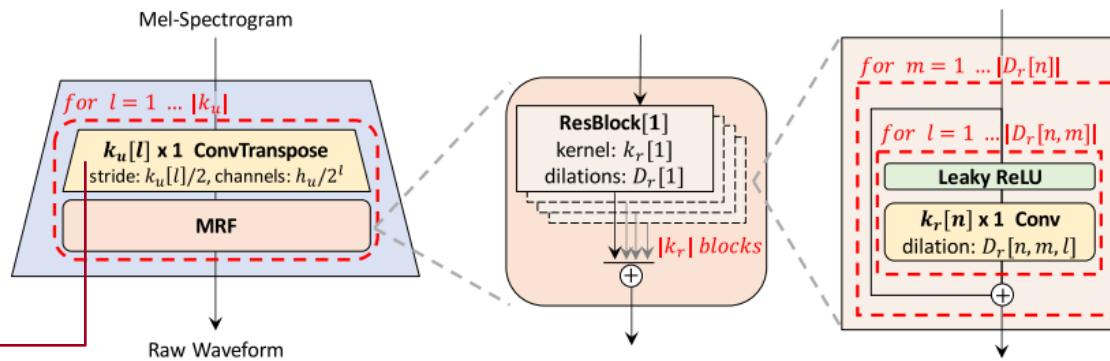


Figure 1: The generator upsamples mel-spectrograms up to $|k_u|$ times to match the temporal resolution of raw waveforms. A MRF module adds features from $|k_r|$ residual blocks of different kernel sizes and dilation rates. Lastly, the n -th residual block with kernel size $k_r[n]$ and dilation rates $D_r[n]$ in a MRF module is depicted.

$$L_{out} = (L_{in}-1) \times \text{stride} - 2 \times \text{padding} + \text{dilation} \times (\text{kernel_size} - 1) + \text{output_padding} + 1$$

HiFi-GAN: Architecture (discriminator)

- Multi-Period Discriminator:

- mixture of sub-discriminators
- each one works with given period p
- used periods: [2, 3, 5, 7, 11]

- Multi-Scale Discriminator:

- drawn from MelGAN paper
- mixture of sub-discriminators
- each one operates on different scales
- scales: [raw, $\times 2$, $\times 4$]

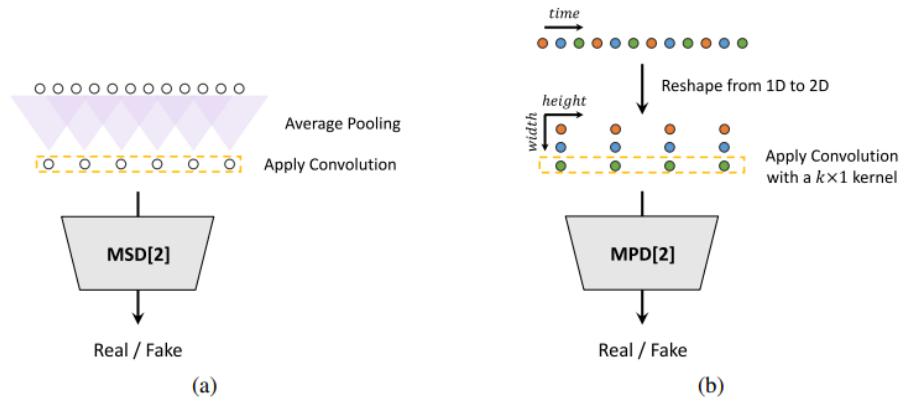


Figure 2: (a) The second sub-discriminator of MSD. (b) The second sub-discriminator of MPD with period 3.

HiFi-GAN: Results

Table 1: Comparison of the MOS and the synthesis speed. Speed of n kHz means that the model can generate $n \times 1000$ raw audio samples per second. The numbers in () mean the speed compared to real-time.

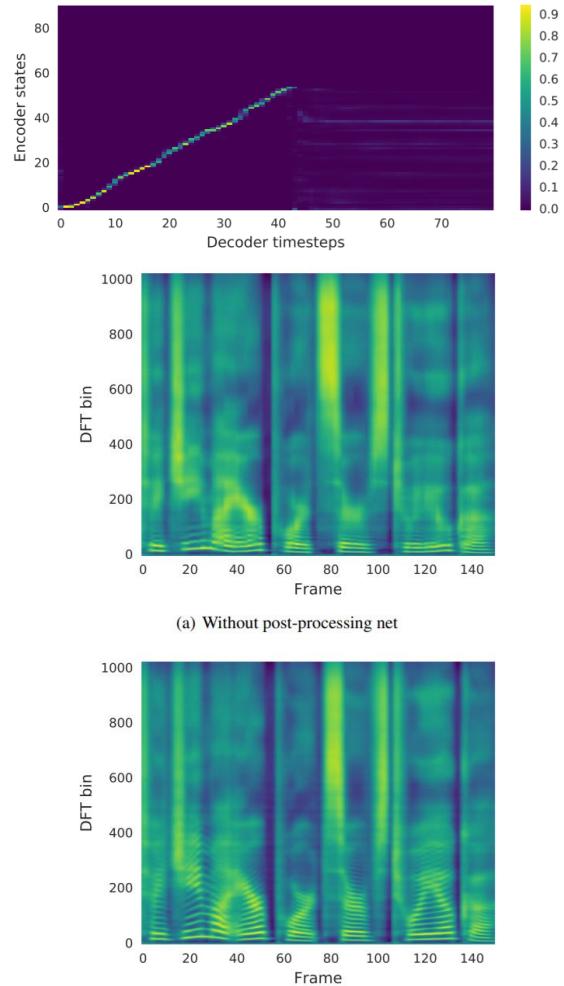
Model	MOS (CI)	Speed on CPU (kHz)	Speed on GPU (kHz)	# Param (M)
Ground Truth	4.45 (± 0.06)	—	—	—
WaveNet (MoL)	4.02 (± 0.08)	—	0.07 ($\times 0.003$)	24.73
WaveGlow	3.81 (± 0.08)	4.72 ($\times 0.21$)	501 ($\times 22.75$)	87.73
MelGAN	3.79 (± 0.09)	145.52 ($\times 6.59$)	14,238 ($\times 645.73$)	4.26
HiFi-GAN V1	4.36 (± 0.07)	31.74 ($\times 1.43$)	3,701 ($\times 167.86$)	13.92
HiFi-GAN V2	4.23 (± 0.07)	214.97 ($\times 9.74$)	16,863 ($\times 764.80$)	0.92
HiFi-GAN V3	4.05 (± 0.08)	296.38 ($\times 13.44$)	26,169 ($\times 1,186.80$)	1.46

Table 2: Ablation study results. Comparison of the effect of each component on the synthesis quality.

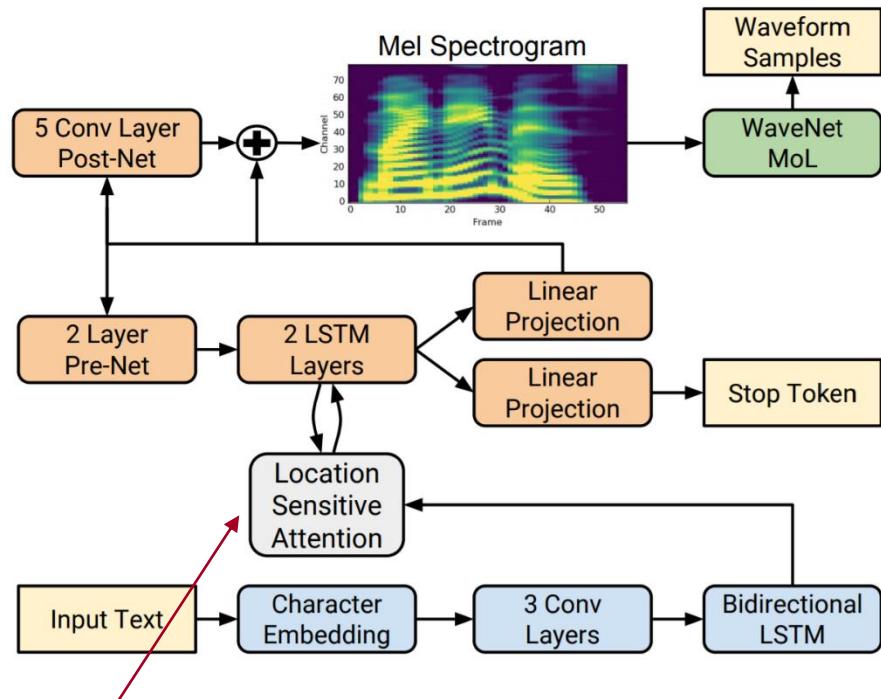
Model	MOS (CI)
Ground Truth	4.57 (± 0.04)
Baseline (HiFi-GAN V3)	4.10 (± 0.05)
w/o MPD	2.28 (± 0.09)
w/o MSD	3.74 (± 0.05)
w/o MRF	3.92 (± 0.05)
w/o Mel-Spectrogram Loss	3.25 (± 0.05)
MPD $p=[2,4,8,16,32]$	3.90 (± 0.05)
MelGAN	2.88 (± 0.08)
MelGAN with MPD	3.35 (± 0.07)

Neural feature-generation: Tacotron/Tacotron2

- Seq2Seq model with attention
 - ⇒ Content-based in Tacotron [Wang17]
 - ⇒ Location-sensitive attention in Tacotron2 [Shen18]
- Predicts 80-band MFSC
 - ⇒ Griffin-Lim vocoder in Tacotron
 - ⇒ WaveNet in Tacotron2
- Predicting N frames per inference step
 - ⇒ N = 2-3 in Tacotron
 - ⇒ N = 1 in Tacotron2 (as reported in [Shen18])
- Makes use of so-called PostNet increasing frequency resolution
- Tacotron2: separate stop-token prediction

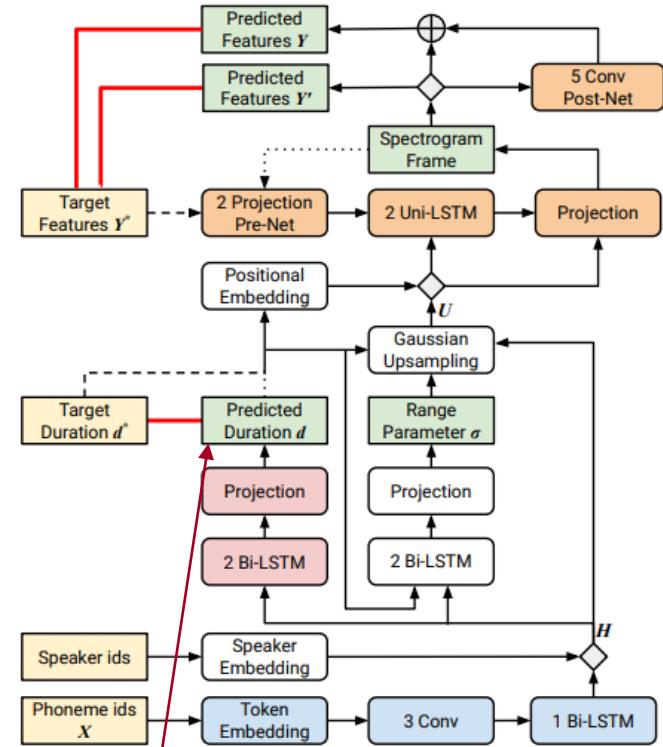


Tacotron2 and Non-attentive Tacotron



Attention mechanism

[Shen18]



Duration prediction

[Shen21]

Glow-TTS [Kim20]

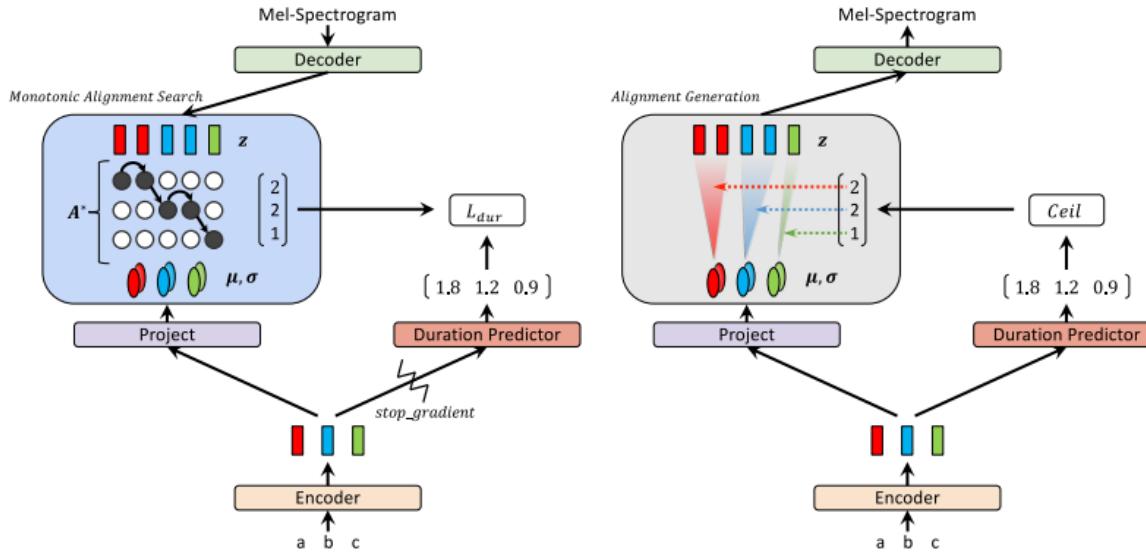
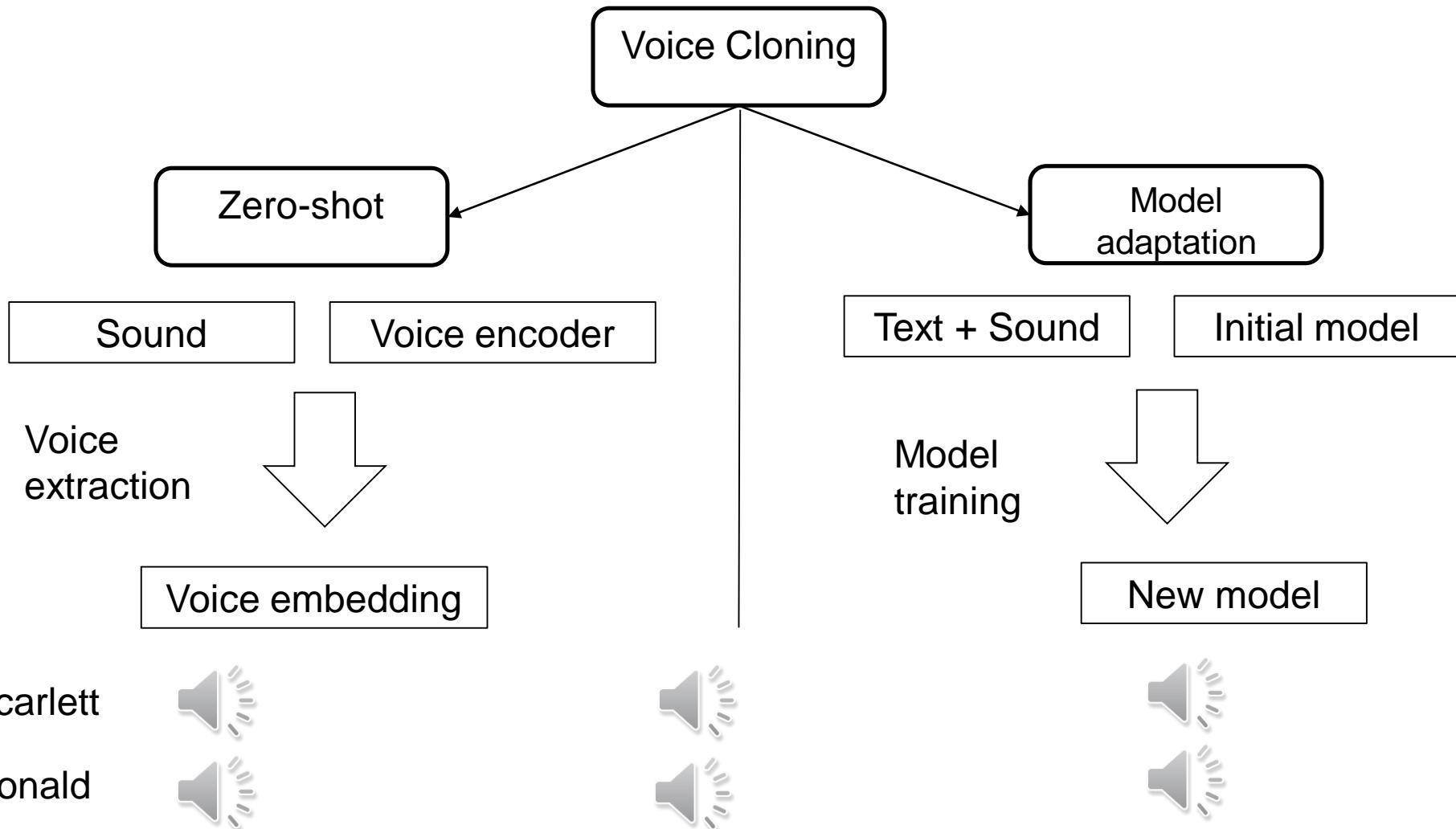


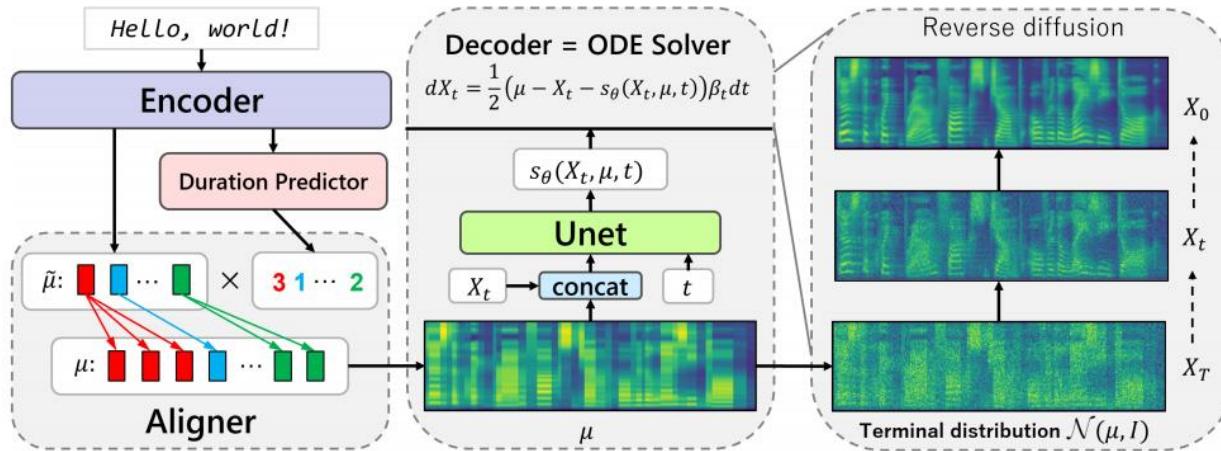
Figure 1: Training and inference procedures of Glow-TTS.

- Text encoder transforms text phoneme embeddings into a sequence of parameters of a Normal distribution $N(\mu, \sigma)$
- Then duration predictor upsamples the sequence to match the length of the target spectrogram
- Decoder is a Normalizing Flow (Glow) transforming samples from $N(\mu, \sigma)$ to spectra
- Duration predictor is trained iteratively to match the optimal alignment path between decoder and encoder outputs

Voice cloning



GRAD-TTS [Popov21]



- Encoder predicts “noisy spectra”, i.e. means of Normal distribution $N_0(\mu, I)$
- Decoder is trained to mimic a trajectory of a backward diffusion process restoring sample X_0 from X_T drawn from $N_0(\mu, I)$
- Forward diffusion trajectory ($X_0 \dots X_T$) can be calculated in closed form for any given X_0 so it can be used for “generating” target values for the backward diffusion
- Duration predictor is trained by means of Monotonic Alignment Search (MAS) from Glow-TTS paper

Accepted for ICML 2021

Speech recognition

SPEECH-TO-TEXT

Speech recognition as sequence labelling

- **Input sequence:** sequence of features (MFCC, MFSC or even raw wave samples)
- **Output sequence:** sequence of phonemes characters
- **Acoustic model:** feature sequence -> phoneme posteriors
- **Language model:** phoneme posteriors -> word sequence posteriors
- **Difference from the classical NLP sequence-labelling: lengths of the input and output sequences are different**
- Previous state-of-art: GMM-HMM models
- GMM acoustic models were replaced with DNN-based acoustic models: TDNN (actually CNN), LSTM, Bi-LSTM, Sequence-to-sequence, Transformer etc.
- Latest architectures can be trained to work without language models

ASR in 2010s

- DNN-HMM models
 - MFCC are replaced with MFSC
 - GMM-HMM system is trained for generating train set for DNN-based system
 - GMM is replaced with DNN as the acoustic model
 - DNN is usually 1d-convolutional (TDNN) or recurrent (GRU, LSTM or Bi-LSTM)
 - NN-based LM is used for hypotheses rescoring
- DNN-CTC models
 - Sequence-to-sequence neural models are trained from scratch with Connectionist Temporal Classification Loss (CTC)
- End2End models
 - No feature extractor is used. Input is raw waveform. No separate language model

Connectionist Temporal Classification (CTC) Intuition

- All operations in Forward-Backward procedure are differentiable
- So we can train a sequence model in a semi-supervised manner
End-to-End without HMMs
- Criterion: $\underset{\theta}{\operatorname{argmax}} P(\mathbf{l}|X, \theta)$ where, X is a sequence of features, \mathbf{l} is the corresponding text
- Let y_t^k is the k-th softmax output at time step k of the RNN
- Let π be a sequence of RNN outputs
- Define $\mathcal{B}(\pi_{1..t})$ as a result of a mapping from the sequence of RNN outputs onto the subsequence of the orthographic form \mathbf{l}
- We assume that all y_t^k are conditionally independent

Connectionist Temporal Classification (CTC)

Forward-backward recursion

Actual CTC aligns feature sequence and ‘extended’ sequence of characters with special blank symbol allowing to map the input sequence directly onto orthographic form: e.g. /apple/ -> “apple”. So we also consider extended orthographic form \mathbf{l}'

Forward procedure

$$\alpha_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T : \\ \mathcal{B}(\pi_{1:t}) = \mathbf{l}_{1:s}}} \prod_{t'=1}^t y_{\pi_{t'}}^{t'}$$

Initialization

$$\begin{aligned}\alpha_1(1) &= y_b^1 \\ \alpha_1(2) &= y_{\mathbf{l}_1}^1 \\ \alpha_1(s) &= 0, \quad \forall s > 2\end{aligned}$$

Recursion

$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{\mathbf{l}_s'}^t & \text{if } \mathbf{l}_s' = b \text{ or } \mathbf{l}_{s-2}' = \mathbf{l}_s' \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{\mathbf{l}_s'}^t & \text{otherwise} \end{cases}$$

where

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

Stop

$$p(\mathbf{l}|\mathbf{x}) = \alpha_T(|\mathbf{l}'|) + \alpha_T(|\mathbf{l}'|-1)$$

Backward procedure

$$\beta_t(s) \stackrel{\text{def}}{=} \sum_{\substack{\pi \in N^T : \\ \mathcal{B}(\pi_{t:T}) = \mathbf{l}_{s:|\mathbf{l}|}}} \prod_{t'=t}^T y_{\pi_{t'}}^{t'}$$

Initialization

$$\begin{aligned}\beta_T(|\mathbf{l}'|) &= y_b^T \\ \beta_T(|\mathbf{l}'|-1) &= y_{\mathbf{l}_{|\mathbf{l}|}}^T \\ \beta_T(s) &= 0, \quad \forall s < |\mathbf{l}'|-1\end{aligned}$$

Recursion

$$\beta_t(s) = \begin{cases} \bar{\beta}_t(s)y_{\mathbf{l}_s'}^t & \text{if } \mathbf{l}_s' = b \text{ or } \mathbf{l}_{s+2}' = \mathbf{l}_s' \\ (\bar{\beta}_t(s) + \beta_{t+1}(s+2))y_{\mathbf{l}_s'}^t & \text{otherwise} \end{cases}$$

where

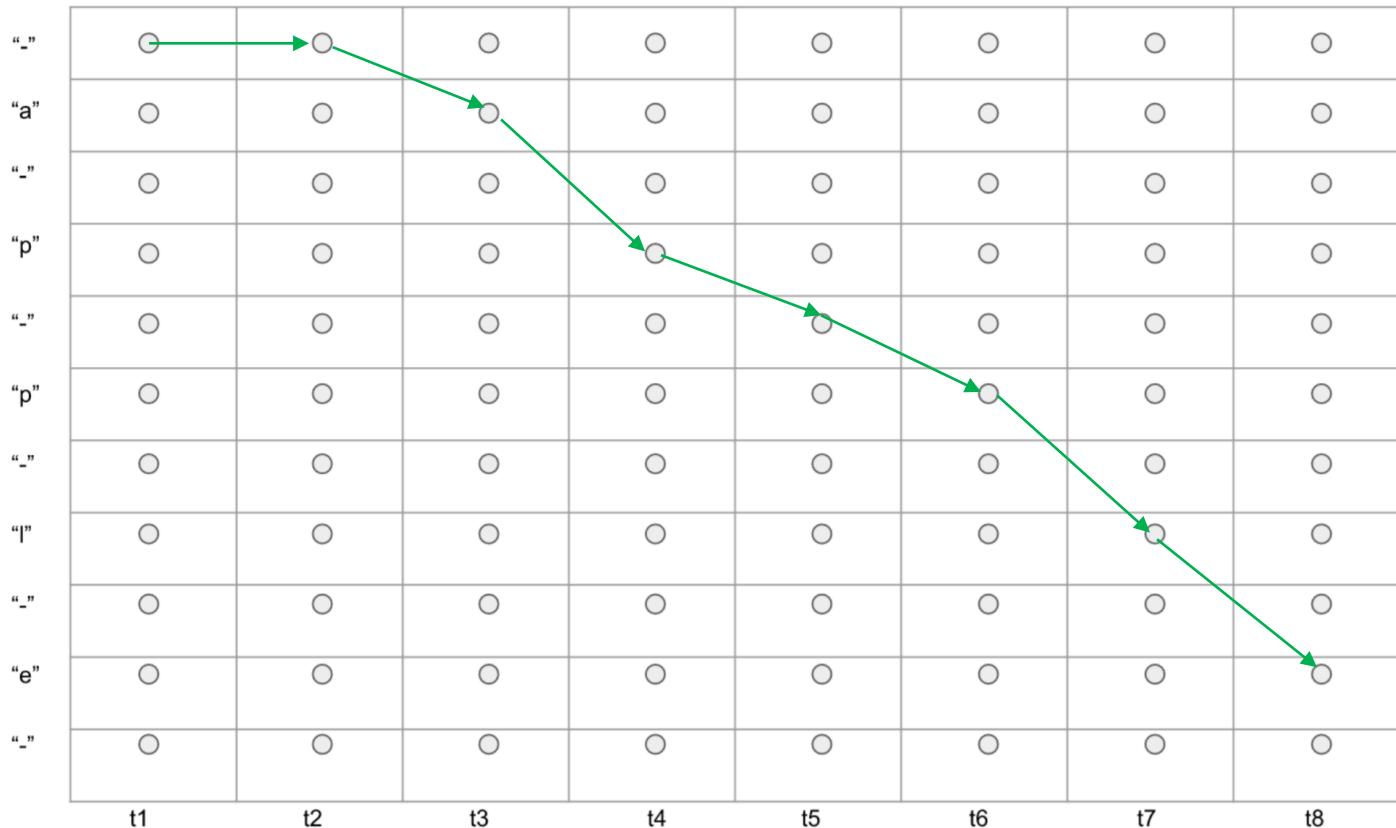
$$\bar{\beta}_t(s) \stackrel{\text{def}}{=} \beta_{t+1}(s) + \beta_{t+1}(s+1).$$

CTC example

“_”	○	○	○	○	○	○	○	○
“a”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“p”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“p”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“l”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“e”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○

CTC example

$$\mathcal{B}(''-ap-ple') = 'apple'$$



CTC example

Initialization

“_”	○	○	○	○	○	○	○	○
“a”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“p”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“p”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“l”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○
“e”	○	○	○	○	○	○	○	○
“_”	○	○	○	○	○	○	○	○

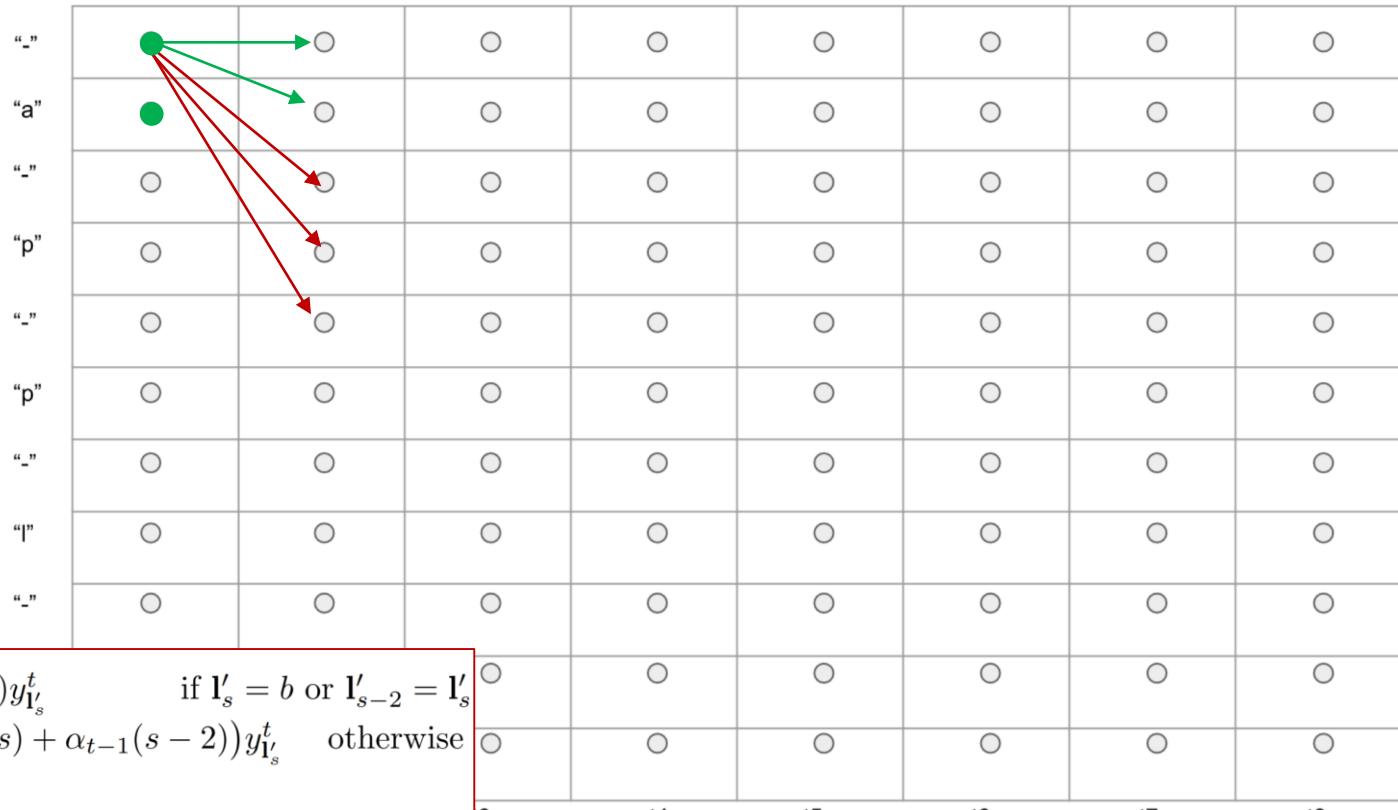
$$\alpha_1(1) = y_b^1$$

$$\alpha_1(2) = y_{I_1}^1$$

$$\alpha_1(s) = 0, \forall s > 2$$

CTC example

Recursion



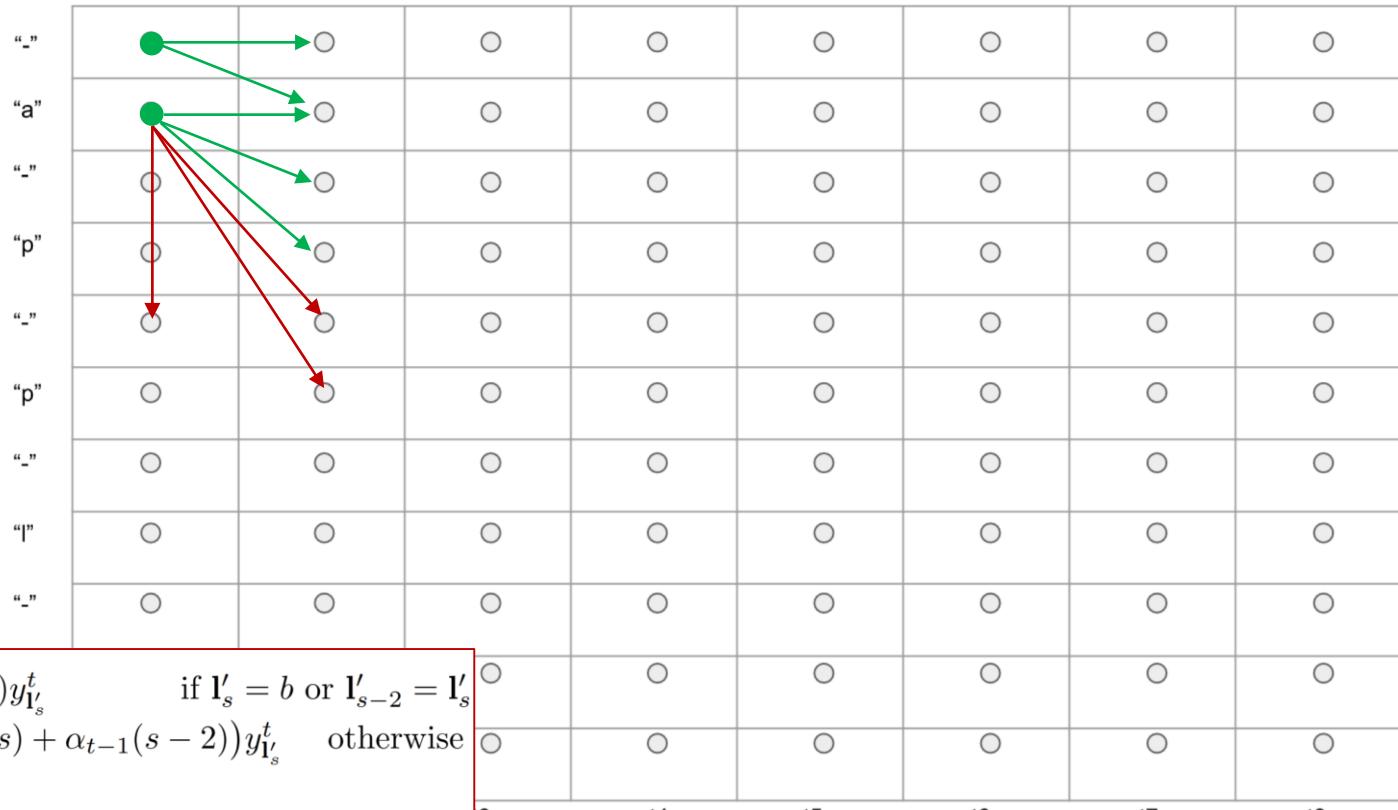
$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{l'_s}^t & \text{otherwise} \end{cases}$$

where

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

CTC example

Recursion



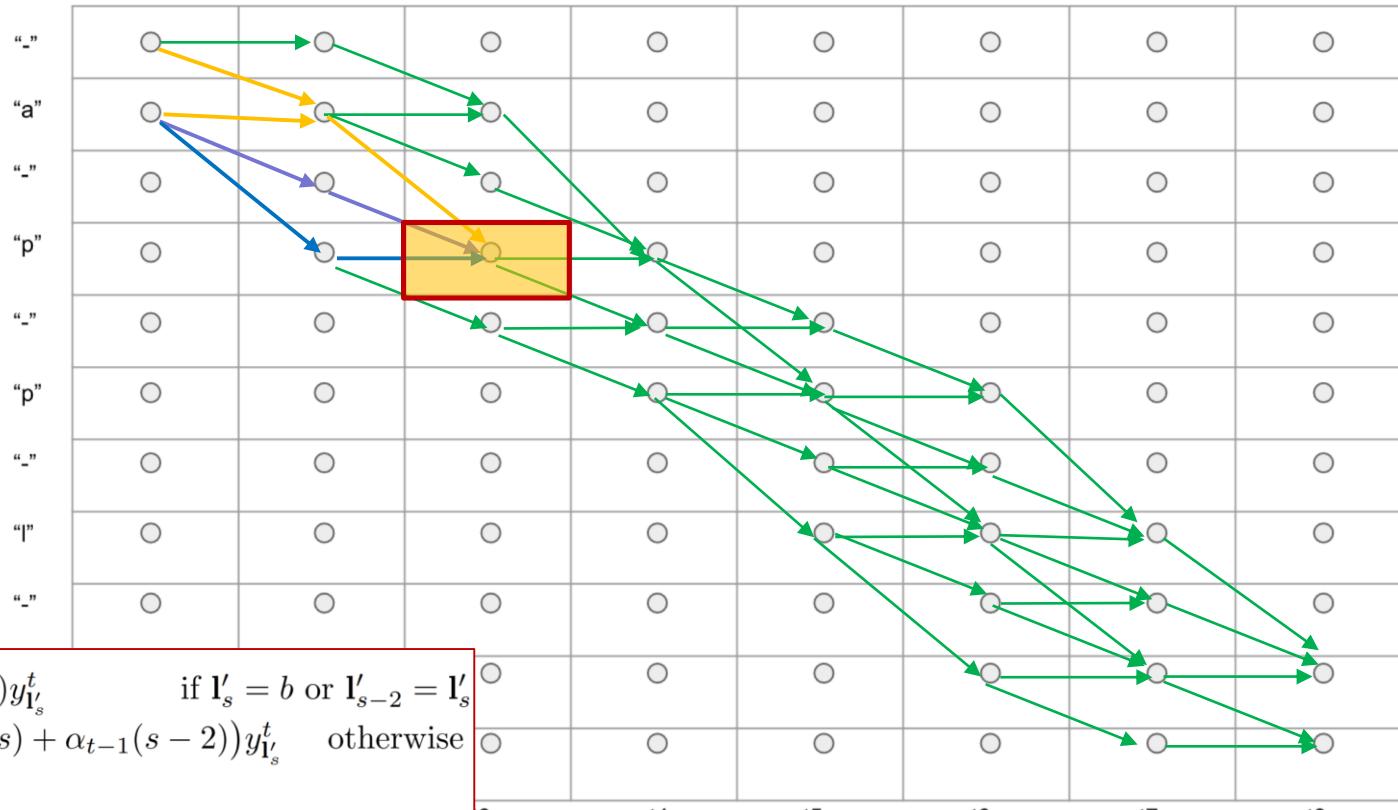
$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{l'_s}^t & \text{if } l'_s = b \text{ or } l'_{s-2} = l'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{l'_s}^t & \text{otherwise} \end{cases}$$

where

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

CTC example

$$\alpha_3(4) = p(" _ap") + p(" aap") + p(" a_p") + p(" app")$$



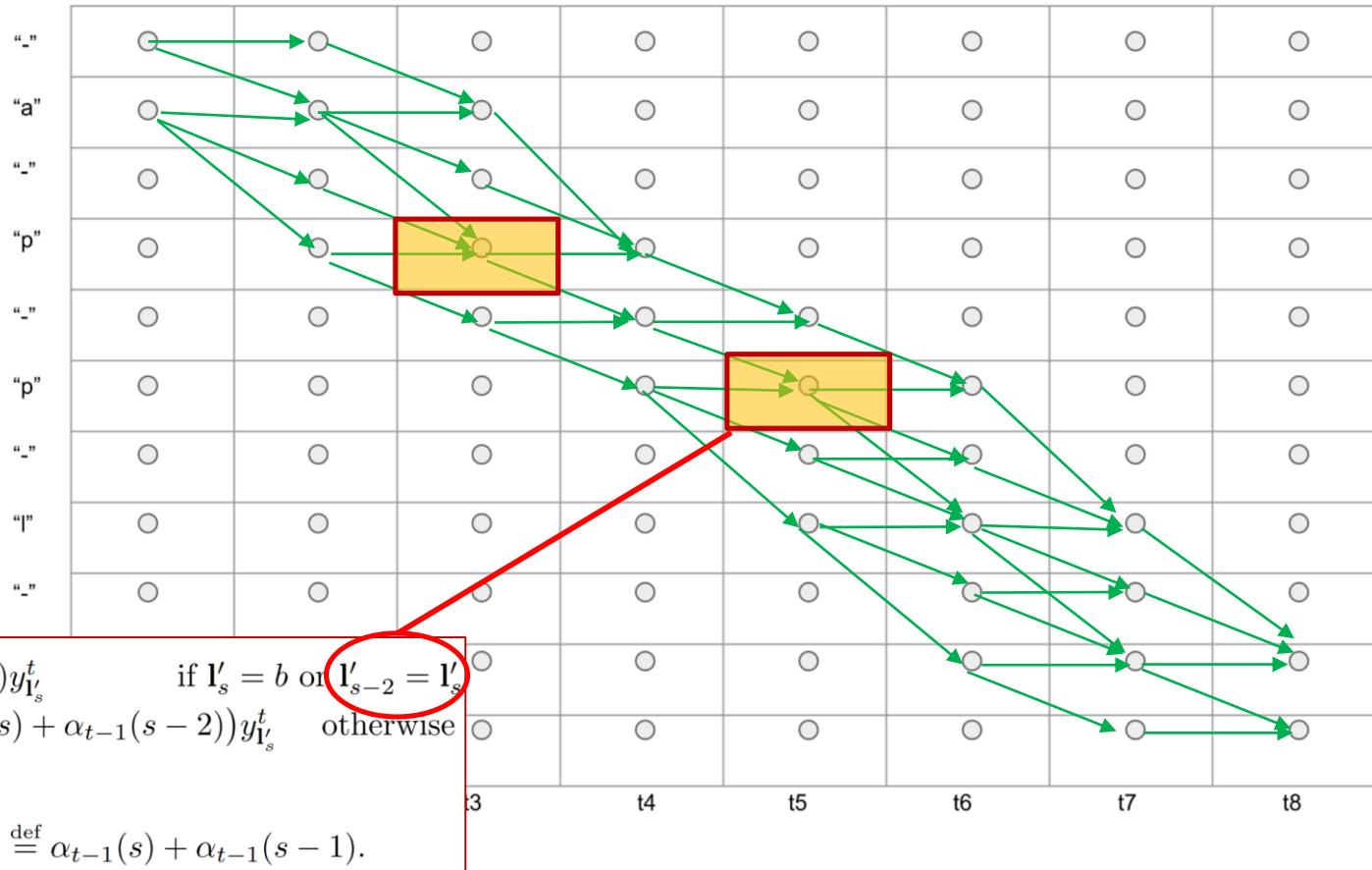
$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{\mathbf{l}'_s}^t & \text{if } \mathbf{l}'_s = b \text{ or } \mathbf{l}'_{s-2} = \mathbf{l}'_s \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{\mathbf{l}'_s}^t & \text{otherwise} \end{cases}$$

where

$$\bar{\alpha}_t(s) \stackrel{\text{def}}{=} \alpha_{t-1}(s) + \alpha_{t-1}(s-1).$$

CTC example

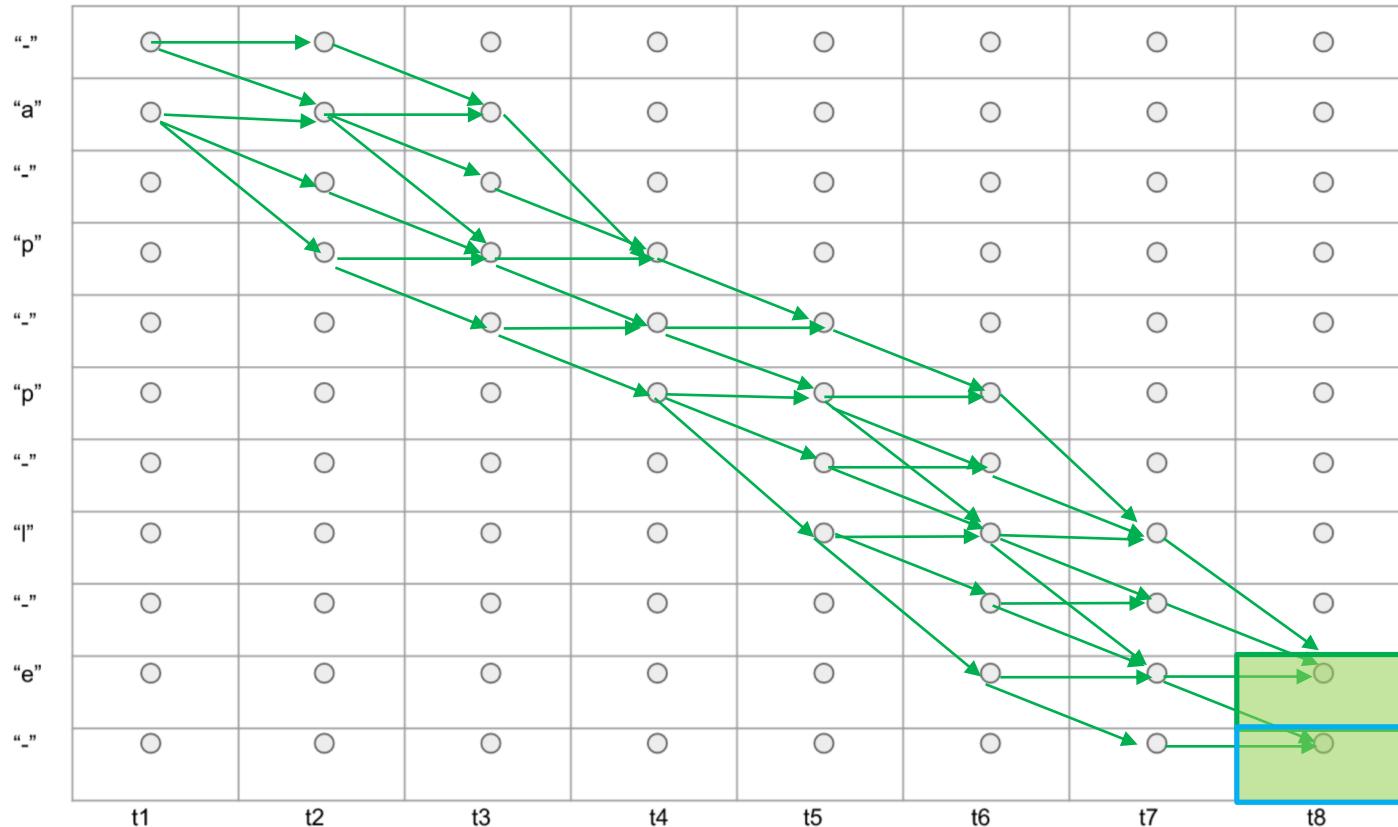
Handling double characters



CTC example

Stop

$$L = -\ln(\alpha_8(10) + \alpha_8(11)) = -\ln(l = 'apple' | X)$$



CTC example

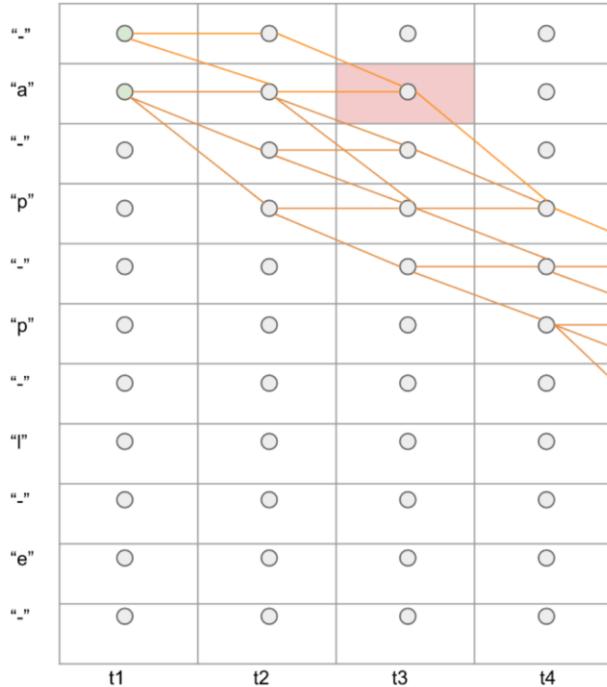


$$\alpha_3(2) = p("aa") + p("-aa") + p("aaa") = \\ = y_-^1 \cdot y_-^2 \cdot y_a^3 + y_-^1 \cdot y_a^2 \cdot y_a^3 + y_a^1 \cdot y_a^2 \cdot y_a^3$$

$$\beta_3(2) = p("apple") = y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

$$\Rightarrow \alpha_3(2) \cdot \beta_3(2) = y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 + \\ + y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 \\ + y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8$$

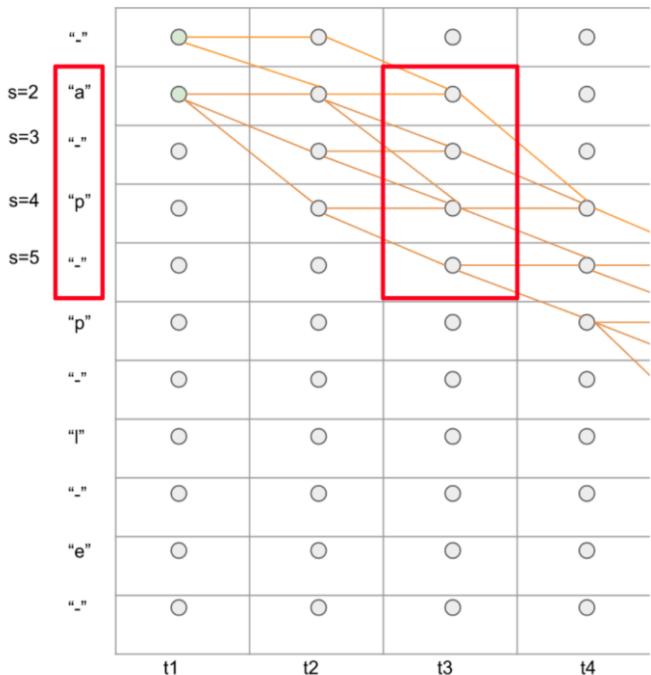
CTC example



$$\begin{aligned}
 & \alpha_3(2) \cdot \beta_3(2) = y_-^1 \cdot y_-^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 + \\
 & + y_-^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 \\
 & + y_a^1 \cdot y_a^2 \cdot y_a^3 \cdot y_a^3 \cdot y_p^4 \cdot y_-^5 \cdot y_p^6 \cdot y_l^7 \cdot y_e^8 = \\
 & = (p("--aap-ple") + p("-aap-ple") + p("aaap-ple")) * y_a^3
 \end{aligned}$$

$$\frac{\alpha_3(2) \cdot \beta_3(2)}{y_a^3} = (p("--aap-ple") + p("-aap-ple") + p("aaap-ple"))$$

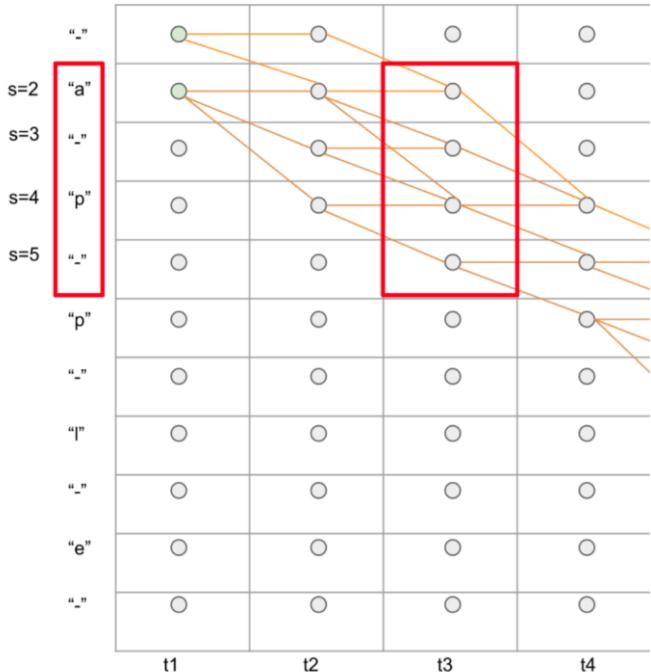
CTC example



$$p_3("apple") = \sum_{s=2}^5 \frac{\alpha_3(s)\beta_3(s)}{y_{l_s}^3}$$

$$p_t("apple") = \sum_{s=1}^{|l|} \frac{\alpha_t(s)\beta_t(s)}{y_{l_s}^t}$$

CTC example: gradient



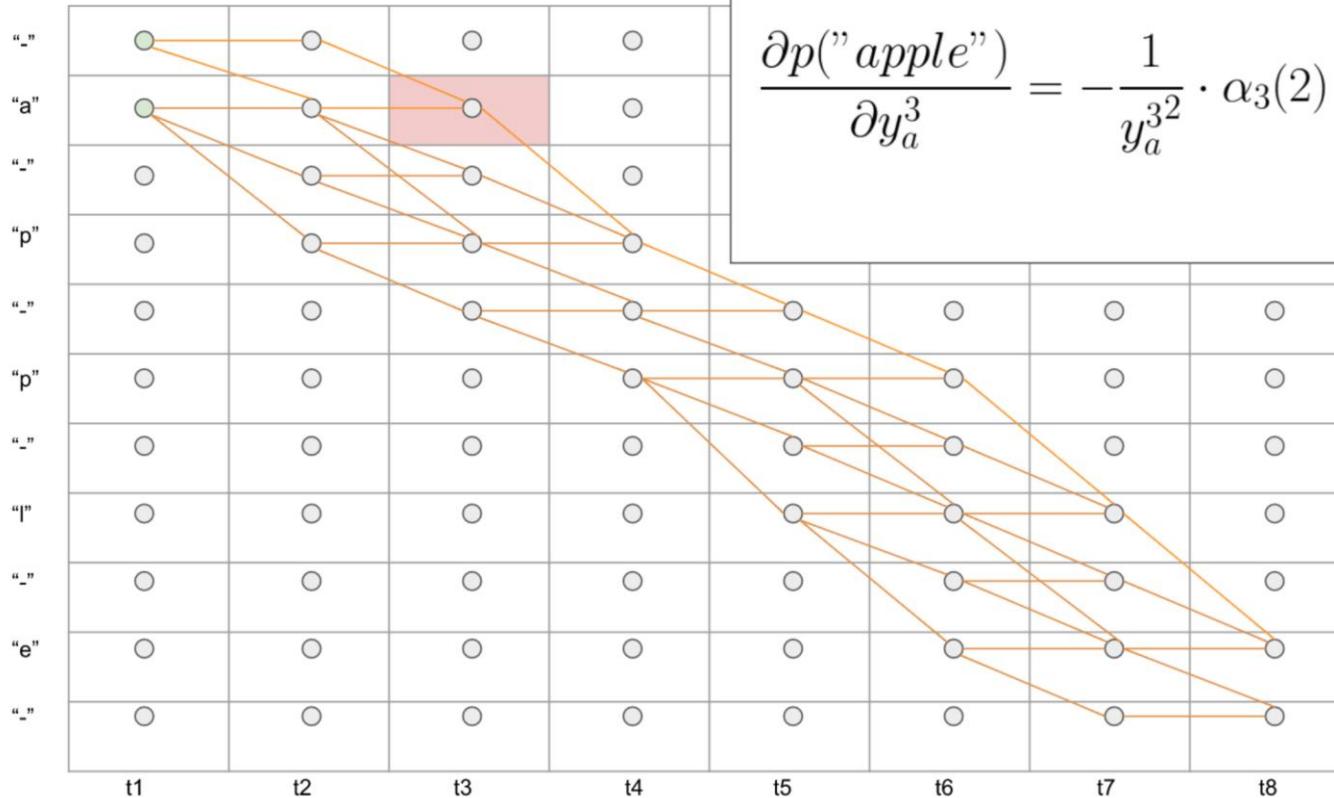
$$p("apple") = \sum_{s=1}^{|seq|} \frac{\alpha_t(s) \cdot \beta_t(s)}{y_{seq(s)}^t}$$

$$\frac{\partial p("apple")}{\partial y_k^t} = -\frac{1}{y_k^{t^2}} \cdot \sum_{s:seq(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

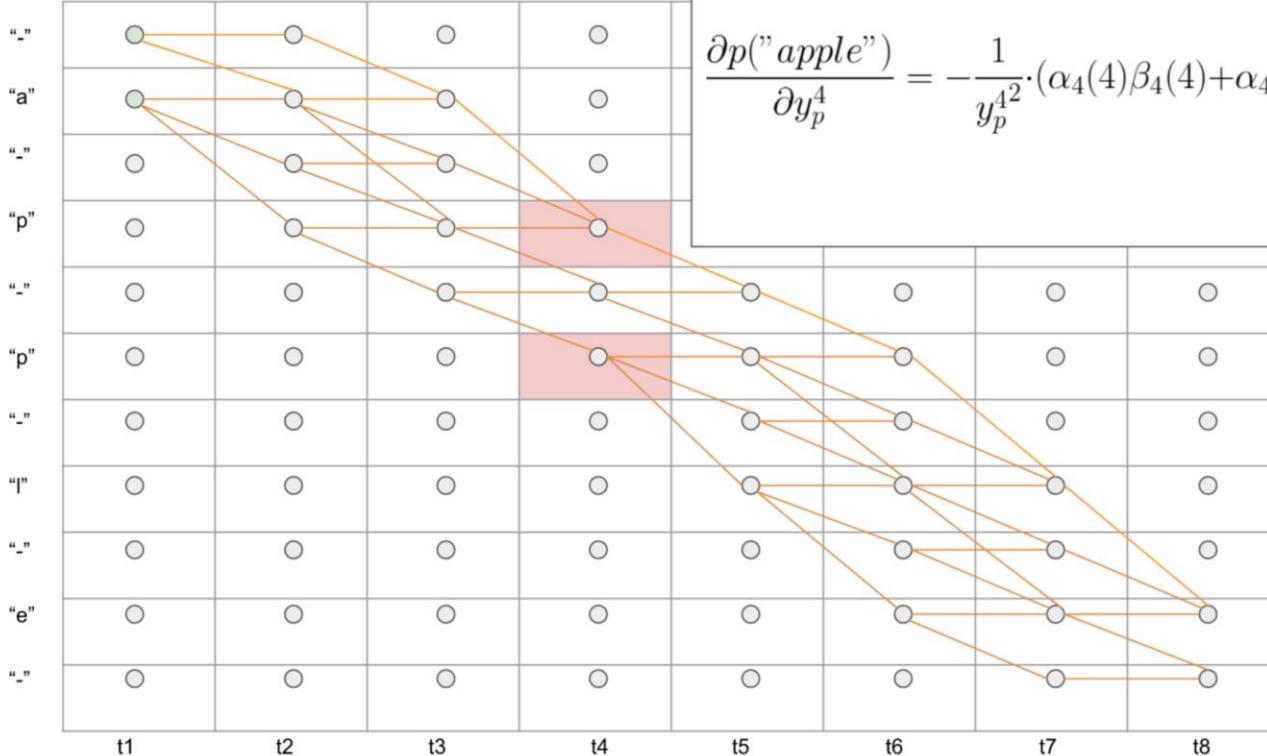
$$\frac{\partial (-\ln(p("apple")))}{\partial y_k^t} = -\frac{1}{p("apple")} \frac{\partial p("apple")}{\partial y_k^t}$$

$$\frac{\partial p("apple")}{\partial y_k^t} = -\frac{1}{y_k^{t^2}} \cdot \sum_{s:seq(s)=k} \alpha_t(s) \cdot \beta_t(s)$$

CTC example: gradient



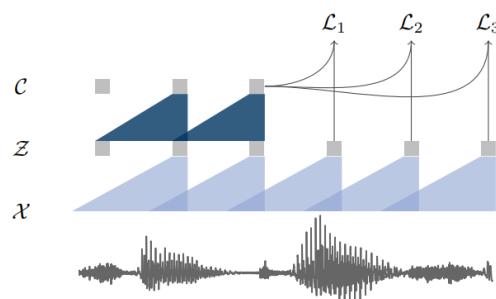
CTC example: gradient



Wav2Vec: evolution

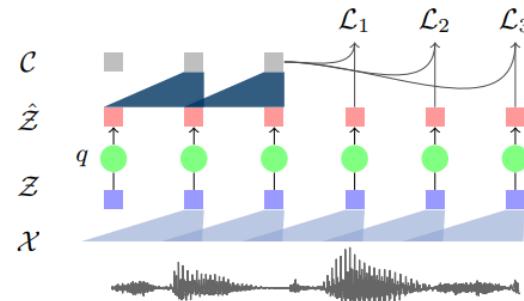
Common part:

- X is raw audio signal
- encoder network $f: X \rightarrow Z$
- context network $g: Z \rightarrow C$
- Objective: contrastive loss
- ~1000h base model pretraining



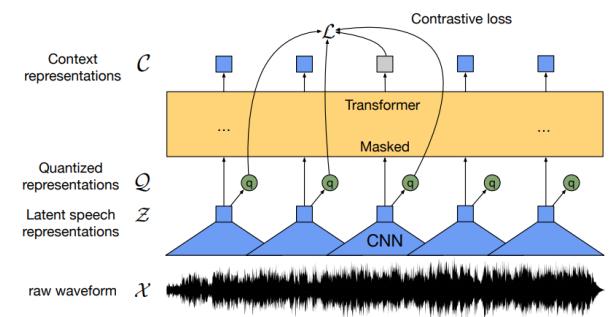
vq-wav2vec:

- quantizer $q: Z \rightarrow \hat{Z}$
- discretized pipeline with BERT



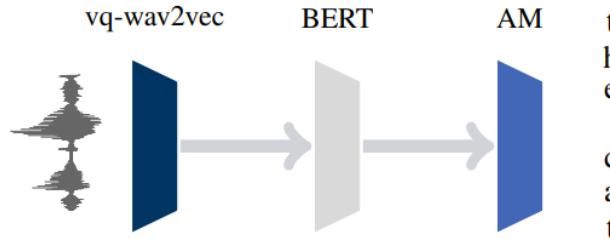
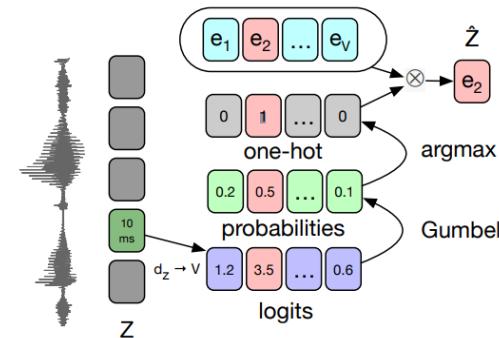
Wav2Vec 2.0:

- transformer $g: Z \rightarrow C$
- layer normalization
- masking
- + diversity loss
- Large model : 60k h pretraining



vq-wav2vec

- Codebook contains V representations of size d $\hat{z} = e_i, \quad e \in \mathbb{R}^{V \times d}$
- Gumbel-Softmax - differentiable way to select discrete codebook entity
- Dividing feature vector into several groups G



Pretraining BERT with discrete inputs of audio data q by masking spans

Wav2Vec Results

WER

LibriSpeech

	dev		test	
	clean	other	clean	other
vq-wav2vec Gumbel + Transformer Big	5.6	15.5	6.2	18.2
wav2vec2 LV-60 + Transformer	1.6	3.0	1.8	3.3

PER

CommonVoice

Model	D	#pt	#ft	es	fr	it	ky	nl	ru	sv	tr	tt	zh	Avg
Number of pretraining hours per language				168h 1h	353h 1h	90h 1h	17h 1h	29h 1h	55h 1h	3h 1h	11h 1h	17h 1h	50h 1h	793h 10h
XLSR-53		D ₅₃	53	1	2.9	5.0	5.7	6.1	5.8	8.1	12.2	7.1	5.1	18.3 7.6

WER CommonVoice

Model	#pt	#ft	en	de	nl	fr	es	it	pt	pl	
Number of training hours			44.7k	2k	1.6k	1.1k	918	247	161	104	
XLSR-53		53	1-100h	13.2	7.4	10.9	9.8	7.9	12.0	15.7	18.9

LibriSpeech test clean:

1	Conformer + Wav2vec 2.0 + SpecAugment-based Noisy Student Training with Libri-Light	1.4
2	Conv + Transformer + wav2vec2.0 + pseudo labeling	1.5
3	ContextNet + SpecAugment-based Noisy Student Training with Libri-Light	1.7
4	Multistream CNN with Self-Attentive SRU	1.75
5	wav2vec 2.0 with Libri-Light	1.8

Wav2Vec Results

Table 1: WER on the Librispeech dev and test sets for the Libri-light low-resource labeled data setups of 10 min, 1 hour and 10 hours. As unlabeled data we use the audio of Librispeech (LS-960) or the larger LibriVox (LV-60k). ST (s2s scratch) trains a sequence to sequence model with a word-piece vocabulary on the pseudo-labeled data from random initialization while as ST (ctc ft) fine-tunes wav2vec 2.0 with the pseudo-labels using CTC and a letter-based vocabulary. All results are with language models at inference time.

Model	Unlbld data	dev		test	
		clean	other	clean	other
10 min labeled					
Discr. BERT [27]	LS-960	15.7	24.1	16.3	25.2
wav2vec 2.0 [24]	LS-960	6.6	10.6	6.8	10.8
+ ST (s2s scratch)	LS-960	4.1	7.0	5.0	8.1
+ ST (ctc ft)	LS-960	3.6	6.6	4.0	7.2
wav2vec 2.0 [24]	LV-60k	5.0	8.4	5.2	8.6
+ ST (s2s scratch)	LV-60k	2.6	4.7	3.1	5.4
+ ST (ctc ft)	LV-60k	2.8	4.6	3.0	5.2
1h labeled					
Discr. BERT [27]	LS-960	8.5	16.4	9.0	17.6
wav2vec 2.0 [24]	LS-960	3.8	7.1	3.9	7.6
+ ST (s2s scratch)	LS-960	2.9	5.6	3.4	6.6
+ ST (ctc ft)	LS-960	2.8	5.5	3.1	6.3
10h labeled					
Discr. BERT [27]	LS-960	5.3	13.2	5.9	14.1
IPL [14]	LS-960	23.5	25.5	24.4	26.0
wav2vec 2.0 [24]	LS-960	2.9	5.7	3.2	6.1
+ ST (s2s scratch)	LS-960	2.5	5.1	3.5	5.9
+ ST (ctc ft)	LS-960	2.6	5.2	2.9	5.7

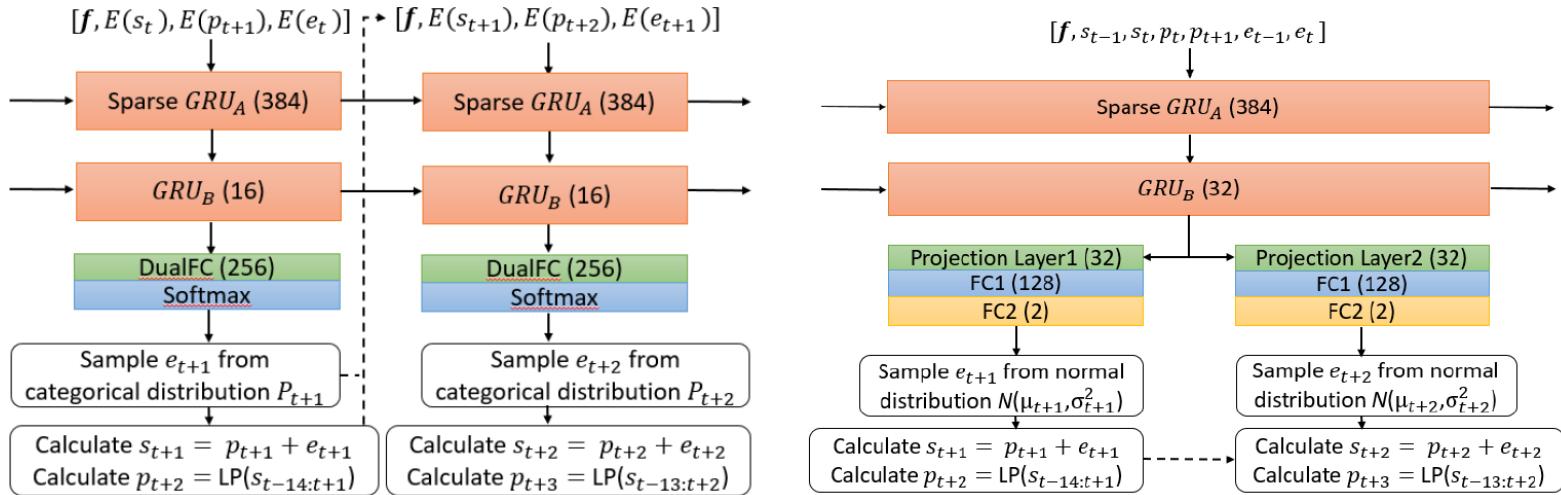
References

- [Oord17] A. van den Oord et al. WAVENET: A GENERATIVE MODEL FOR RAW AUDIO
- [Kal18] N. Kalchbrenner et al. EFFICIENT NEURAL AUDIO SYNTHESIS
- [Wang17] Y. Wang et al. TACOTRON: TOWARDS END-TO-END SPEECH SYNTHESIS
- [Shen18] J. Shen et al. NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS
- [Shen21] J. Shen et al. Non-Attentive Tacotron: Robust and Controllable Neural TTS Synthesis Including Unsupervised Duration Modeling
- [Oord&Li17] A. van den Oord et al. PARALLEL WAVENET: FAST HIGH-FIDELITY SPEECH SYNTHESIS
- [Ping18] W.Ping et al. CLARINET: PARALLEL WAVE GENERATION IN END-TO-END TEXT-TO-SPEECH
- [Pren19] R.Prenger et al. WAVEGLOW: A FLOW-BASED GENERATIVE NETWORK FOR SPEECH SYNTHESIS
- [Val19] J.-M. Valin, J. Skoglund LPCNET: IMPROVING NEURAL SYNTHESIS THROUGH LINEAR PREDICTION
- [Kim20] Glow-TTS: A Generative Flow for Text-to-Speech via Monotonic Alignment Search
- [Kong20] HiFi-GAN: Generative Adversarial Networks for Efficient and High Fidelity Speech Synthesis
- [Popov21] Grad-TTS: A Diffusion Probabilistic Model for Text-to-Speech

Thank You

www.huawei.com

Multisample Gaussian LPCNet



- Idea #1: predict $n > 1$ samples per step
 - ⇒ Speedup up to 50%
- Idea #2: predict parameters of Normal distribution
 - ⇒ 16-bit sound
 - ⇒ ~50% reduction in number of parameters

Connectionist Temporal Classification (CTC)

Forward-backward recursion (simple case)

- $P\{D(X) = Y\} =$

$$\begin{aligned}& \sum_k P\{D(X_1^t) = S_1^k, D(x_t) = s_k\} \cdot P\{D(X_{t+1}^T) = S_k^L \cup D(X_{t+1}^T) = S_{k+1}^L | D(x_t) = s_k\} = \\&= \sum_k P(D(X_1^t) = S_1^k, D(x_t) = s_k) \frac{[P(D(X_{t+1}^T) = S_k^L, D(x_t) = s_k) + P(D(X_{t+1}^T) = S_{k+1}^L, D(x_t) = s_k)]}{p(D(x_t) = s_k)} \\&= \sum_k \alpha_t(s_k) \beta_t(s_k) / y_t^k\end{aligned}$$

Forward procedure (no blanks)

Init: $\alpha_1(s_k) = y_1^k$

Recursion: $\alpha_t(s_k) = [\alpha_{t-1}(s_{k-1}) + \alpha_{t-1}(s_k)]y_t^k$

Stop: $\alpha_T(s_L) = P(D(X) = Y)$

Backward procedure (no blanks)

Init: $\beta_T(s_k) = y_T^k$

Recursion: $\beta_t(s_k) = [\beta_{t+1}(s_{k+1}) + \beta_{t+1}(s_k)]y_t^k$

Stop: $\beta_1(s_1) = P(D(X) = Y)$