

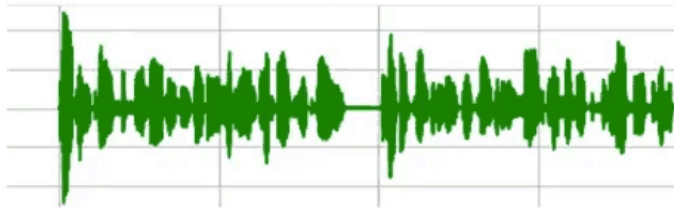
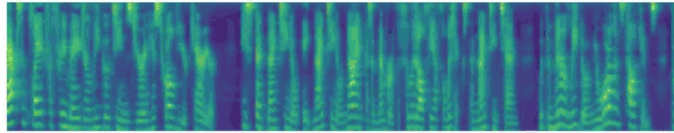
# Speech Synthesis

Milan Straka

 May 20, 2025

# Mel Spectrograms

"Hello, my name is Joseph."



Text Input



Mel Spectrogram



Wave Form

[https://miro.medium.com/v2/1\\*QxNK0\\_jgds7tKIDaHGRz7A.png](https://miro.medium.com/v2/1*QxNK0_jgds7tKIDaHGRz7A.png)

# Speech Generation Pipeline

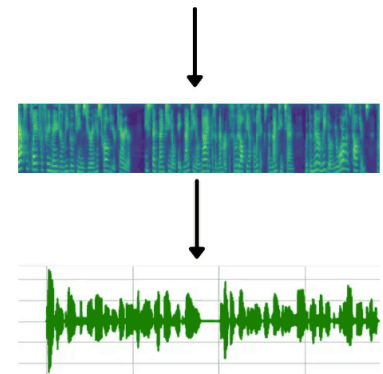
Traditionally, the input text was represented by **phonemes**.

VOWELS	monophthongs				diphthongs		<b>Phonemic Chart</b> voiced unvoiced		
	i: sheep	ɪ ship	ʊ good	u: shoot	ɪə here	eɪ wait			
	e bed	ə teacher	ɜ: bird	ɔ: door	ʊə tourist	ɔɪ boy			əʊ show
	æ cat	ʌ up	ɑ: far	ɒ on	eə hair	aɪ my			aʊ cow
CONSONANTS	p pea	b boat	t tea	d dog	tʃ cheese	dʒ June	k car	g go	
	f fly	v video	θ think	ð this	s see	z zoo	ʃ shall	ʒ television	
	m man	n now	ŋ sing	h hat	l love	r red	w wet	j yes	

**Phonemic Chart**  
voiced  
unvoiced

<https://www.englishclub.com/images/pronunciation/Phonemic-Chart.jpg>

"Hello, my name is Joseph."

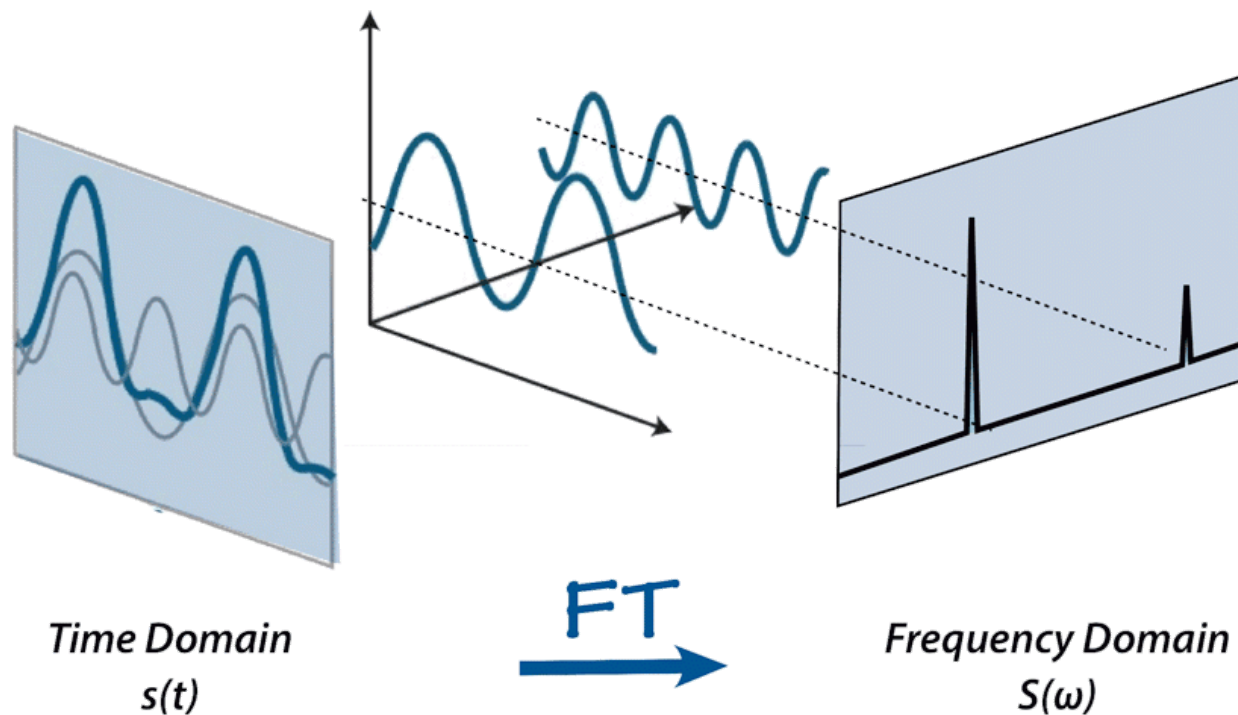


[https://miro.medium.com/v2/1\\*QxNK0\\_jgds7tKIDaHGRz7A.png](https://miro.medium.com/v2/1*QxNK0_jgds7tKIDaHGRz7A.png)

Text Input  
↓  
Mel Spectrogram  
↓  
Wave Form

However, nowadays the input is most commonly represented using **characters** (some systems still allow specifying some input words using phonemes to allow more control over the output).





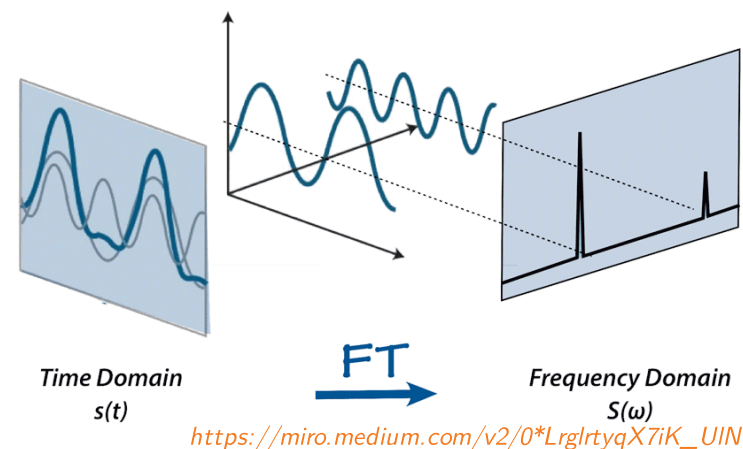
[https://miro.medium.com/v2/0\\*LrglrtyqX7iK\\_UIN](https://miro.medium.com/v2/0*LrglrtyqX7iK_UIN)

A sound can be represented in the **time domain**, as a sequence of amplitudes, i.e., how the signal changes over time. —  $\rightarrow$  *tabhe faryaqe samplash v milwfon*

Alternatively, it can be represented in a **frequency domain**, where a signal is represented by a *spectrum*: the magnitudes and phases of sinusoids with a collection of frequencies.

To convert a signal between time and frequency domains, Fourier transform can be used. In case of discrete-time signal, **discrete Fourier transform (DFT)** converts

- $N$  complex-valued equally-spaced samples in the time domain to
- $N$  complex-valued coefficients (magnitude and phase) of sinusoids with frequencies 1 Hz, 2 Hz, 3 Hz, ...,  $N$  Hz.



**Fast Fourier transform (FFT)** is an algorithm that can perform discrete Fourier transform in in  $\mathcal{O}(N \log N)$  time for  $N = 2^k$  a power of two.

*> možu pro realné čísla počet jen polku spektra*

When the input signal consists of only real values, the upper half of the frequency spectrum is determined by the lower half ( $x_{N-i} = \overline{x_i}$ , where  $\overline{x}$  denotes complex conjugate); therefore, a discrete signal of  $N$  samples can be represented by sinusoids with frequencies of 1, 2, 3, ...,  $\frac{N}{2}$ .

# Spectrogram

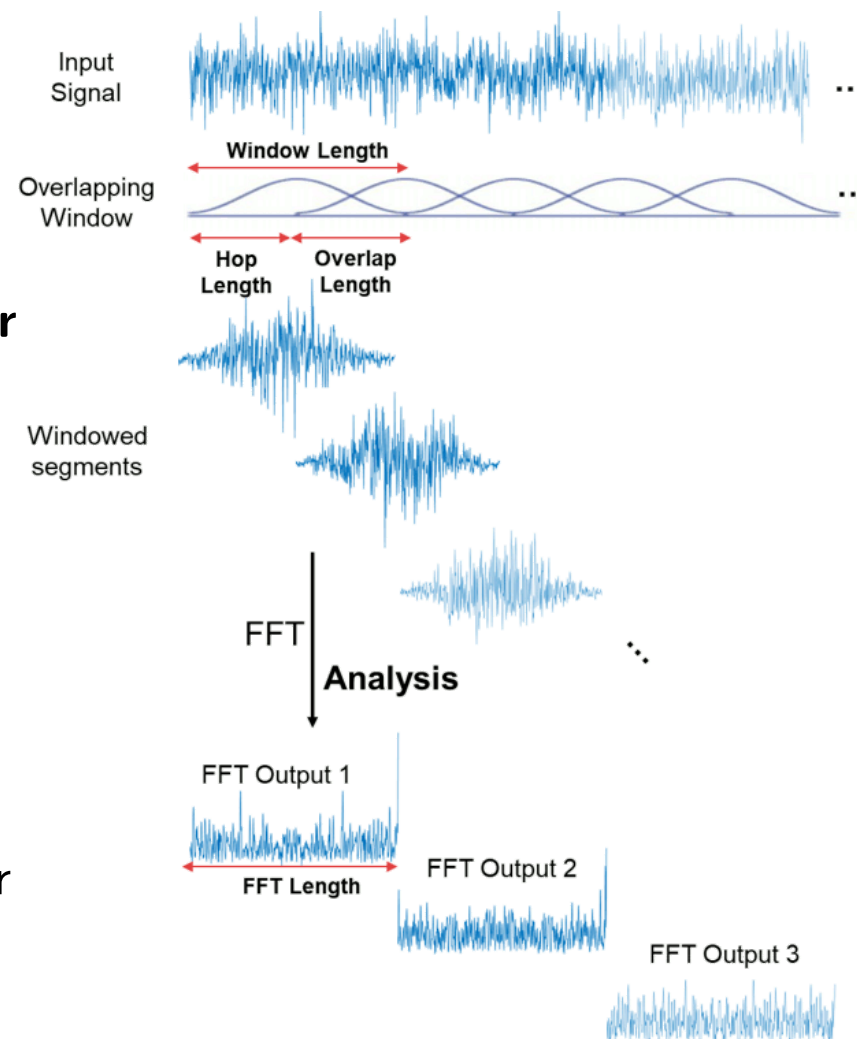
We can represent an arbitrary-length signal by a **spectrogram**, a series of fixed-length spectral densities (i.e., for every sinusoid, we compute some real-valued density, usually *power density*  $x^2$ ).

The spectrogram is computed using **short-time Fourier transform (STFT)**, where

- DFT is performed on windows of fixed size,
- neighboring windows are shifted by a fixed hop length.

For every window, the signal is composed of sinusoids with frequencies up to half of the window length.

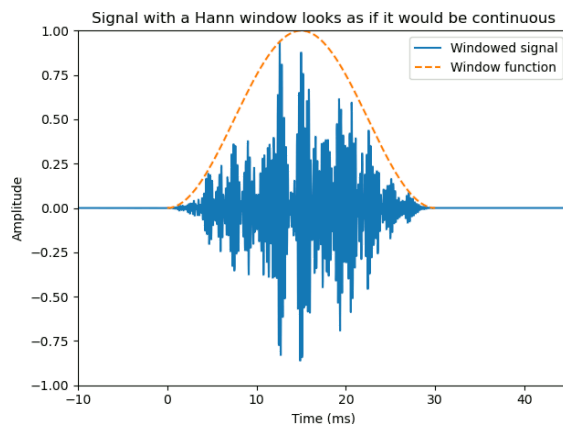
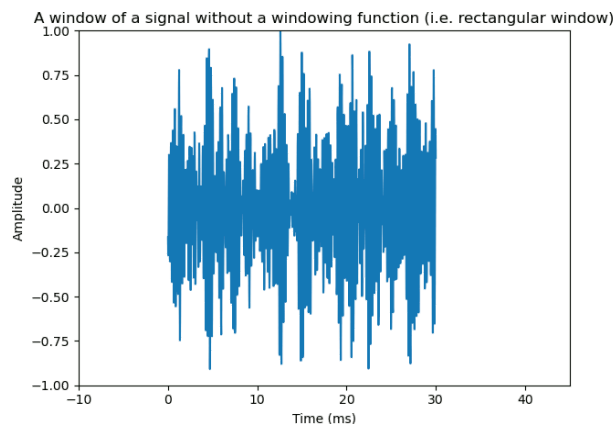
In speech synthesis, assuming a sampling frequency of 22.05kHz or 24kHz, commonly used values are 1024 for window length (denoted commonly as  $n\_fft$ ; ~45ms) and hop length of 256 (circa 11ms).



[https://www.mathworks.com/help/dsp/ref/stft\\_output.png](https://www.mathworks.com/help/dsp/ref/stft_output.png)

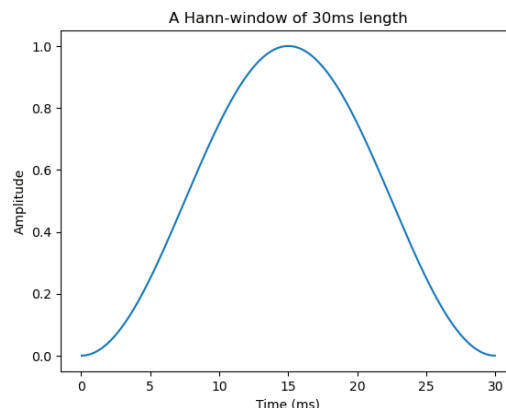
# Spectrogram

To avoid discontinuities at the borders of the windows, *windowing* is used.

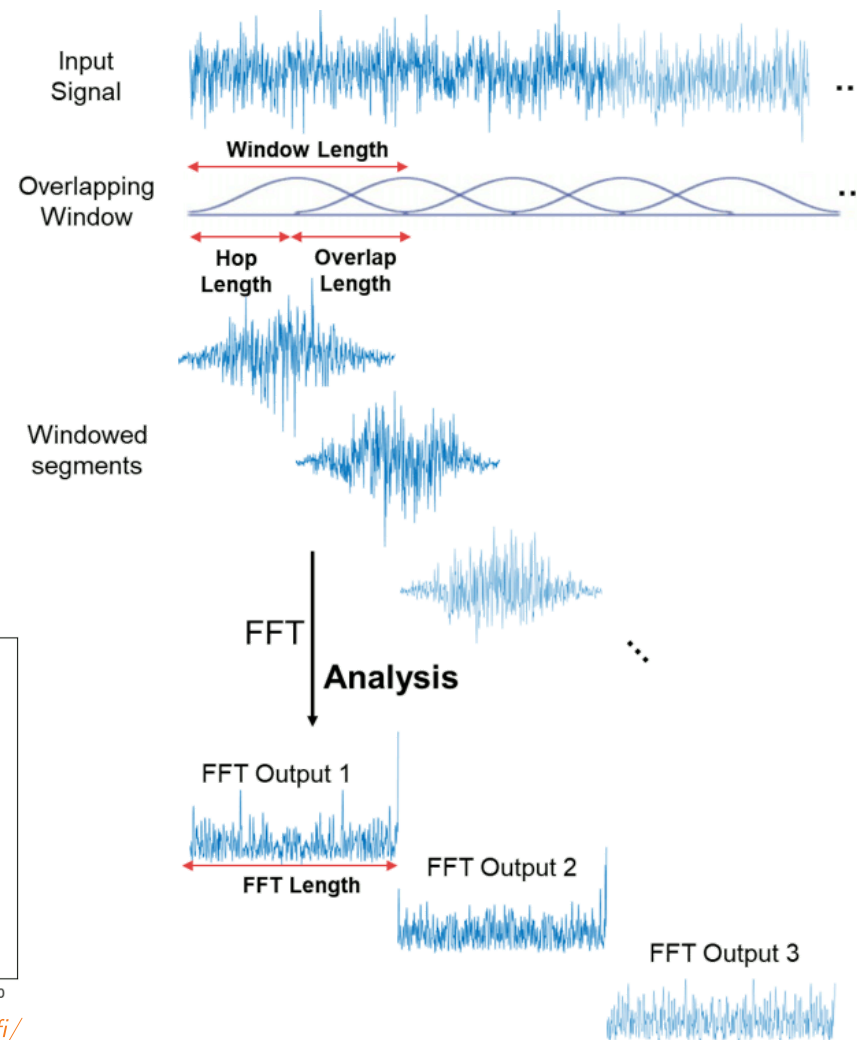


<https://speechprocessingbook.aalto.fi/Representations/Windowing.html>

Many possible window function exist; we will use the **Hann** window, which is as suitably scaled cosine function.



<https://speechprocessingbook.aalto.fi/Representations/Windowing.html>



[https://www.mathworks.com/help/dsp/ref/stft\\_output.png](https://www.mathworks.com/help/dsp/ref/stft_output.png)

# Mel Spectrogram

Linearly increasing frequencies are not perceived to be equal in distance from one another by humans.

← empiricky naměřeno

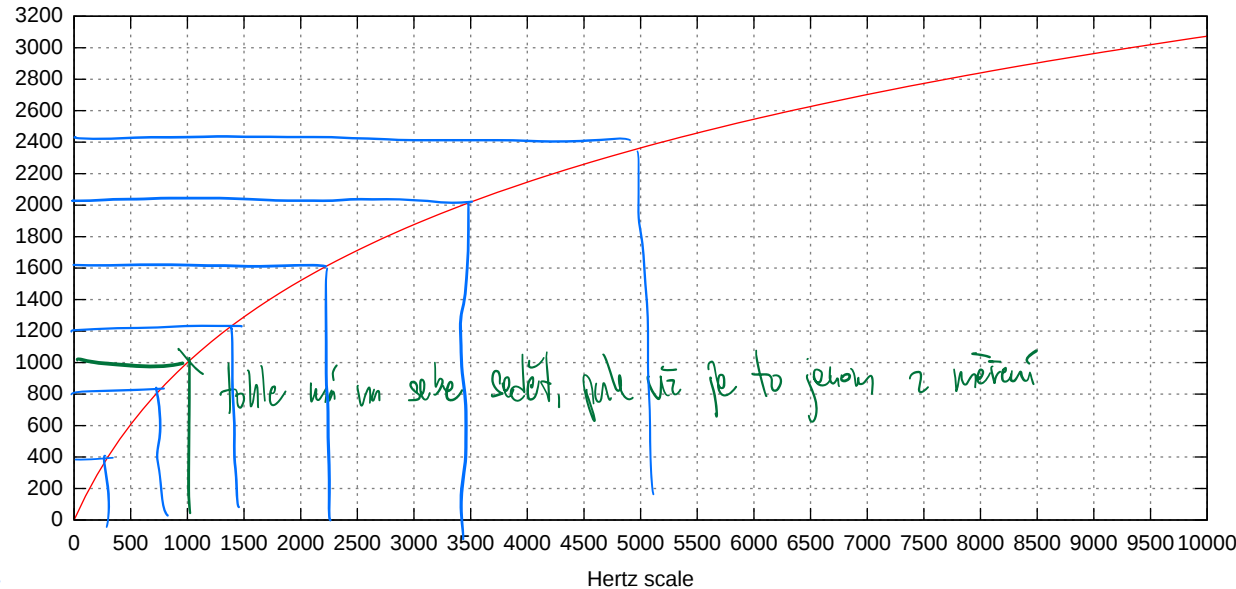
Therefore, the **mel scale** (from the word melody) was proposed, so that the pitches were judged by listeners to be in equal distance.

There are in fact several mel scales, the most commonly used is

$$m = 2595 \log_{10} \left( 1 + \frac{f}{700} \right).$$

→  
evidentní je hlavní focus ve  
spodních frekvencích, kde je potřeba  
i pro hodnoty m malý skok daleko  
větší odstoupení

→  
střed  
je to  
furt  
stejný  
odskok



→  
tohle už m sebe sedět, pak už je to jenom z měření

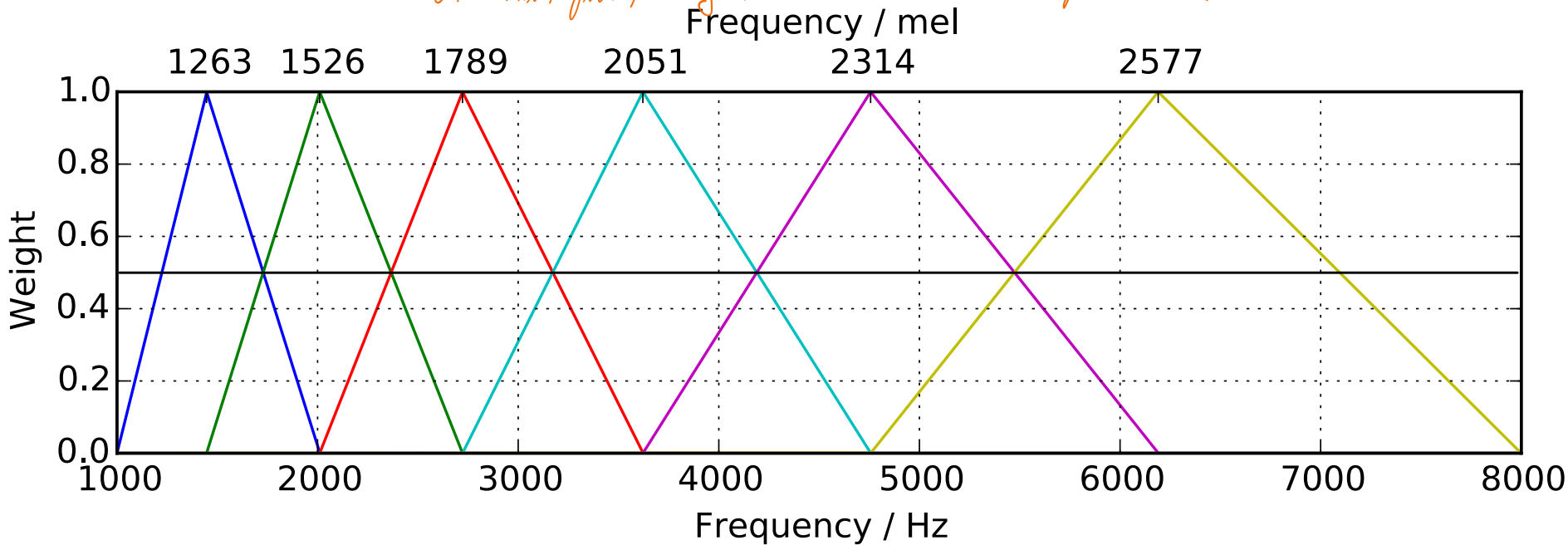
→  
postupně se  
mí prodlužuje  
odskok

[https://upload.wikimedia.org/wikipedia/commons/a/aa/Mel-Hz\\_plot.svg](https://upload.wikimedia.org/wikipedia/commons/a/aa/Mel-Hz_plot.svg)

# Mel Spectrogram

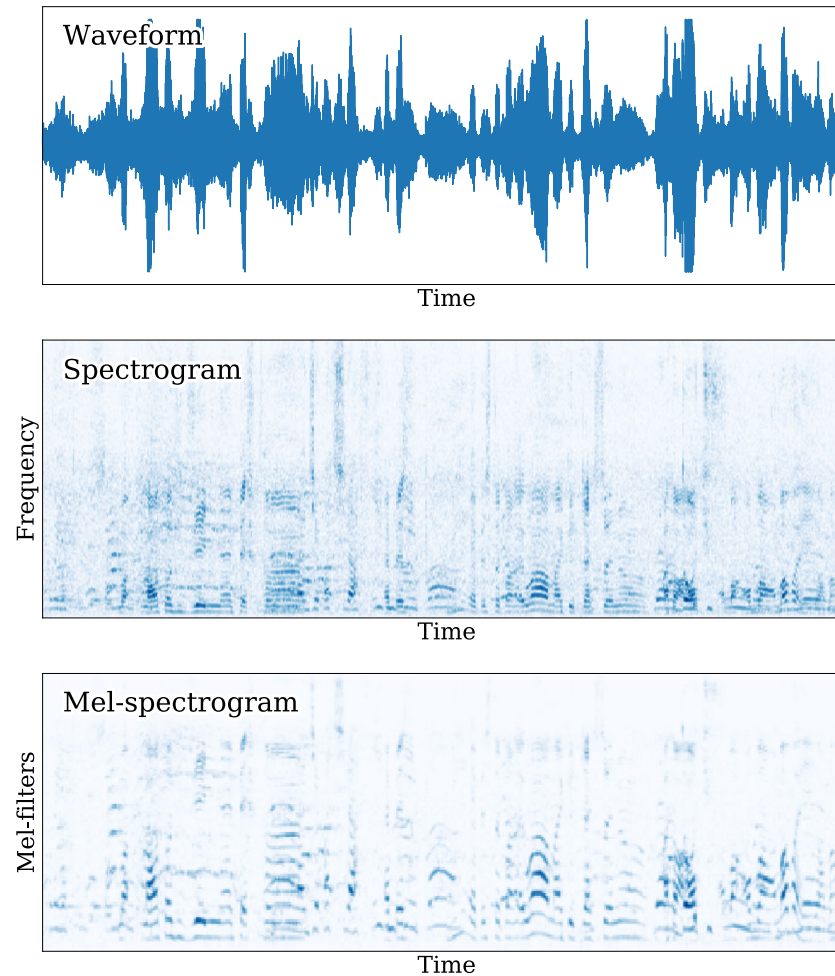
To convert spectrogram to mel spectrogram, a collection of triangular filters equally-distanced in the mel scale is used.

*Tohle mám proto, aby mohl vrát i mezi více pruhovými frekvencemi*



[https://sigggue.github.io/pyfilterbank//melbank-1\\_00.pdf](https://sigggue.github.io/pyfilterbank//melbank-1_00.pdf)

Finally, we usually take the logarithm of the power spectrum, given that humans perceive loudness on logarithmic scale (i.e., dBs are also a logarithmic scale).

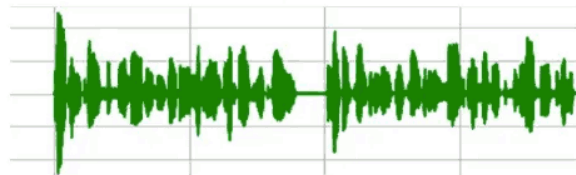
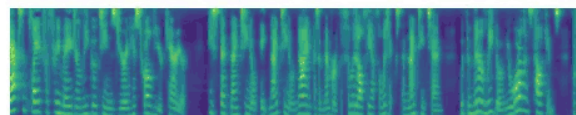


**Fig. 1** Waveform, spectrogram, and mel-spectrogram of a 10-s speech segment obtained from Google AudioSet. The mel-spectrogram, based on the auditory-based mel-frequency scale, provides better resolution for lower frequencies than the spectrogram

Figure 1 of "Exploring convolutional, recurrent, and hybrid deep neural networks for speech and music detection in a large audio dataset", <https://doi.org/10.1186/s13636-019-0152-1>



"Hello, my name is Joseph."



Text Input



Mel Spectrogram



Wave Form

[https://miro.medium.com/v2/1\\*QxNK0\\_jgds7tKIDaHGRz7A.png](https://miro.medium.com/v2/1*QxNK0_jgds7tKIDaHGRz7A.png)

In the context of speech synthesis, a **vocoder** (combination of words voice and encoder) is the second part of the pipeline, converting a (usually mel) spectrogram to a waveform.



# WaveNet

WaveNet, proposed in 2016, was one of the first neural architectures that produced high quality human speech.

Not many details were published in the original paper. From today's perspective, WaveNet is a vocoder capable of converting mel spectrograms to waveforms.

In the following, we start the description without the mel spectrogram conditioning, but we add it later.

Our goal is to model speech, using a convolutional auto-regressive model

$$P(\mathbf{x}) = \prod_t P(x_t | x_{t-1}, \dots, x_1).$$

*representace zmluvneho signalu v jednom casovom slozku*  
*decoder only*

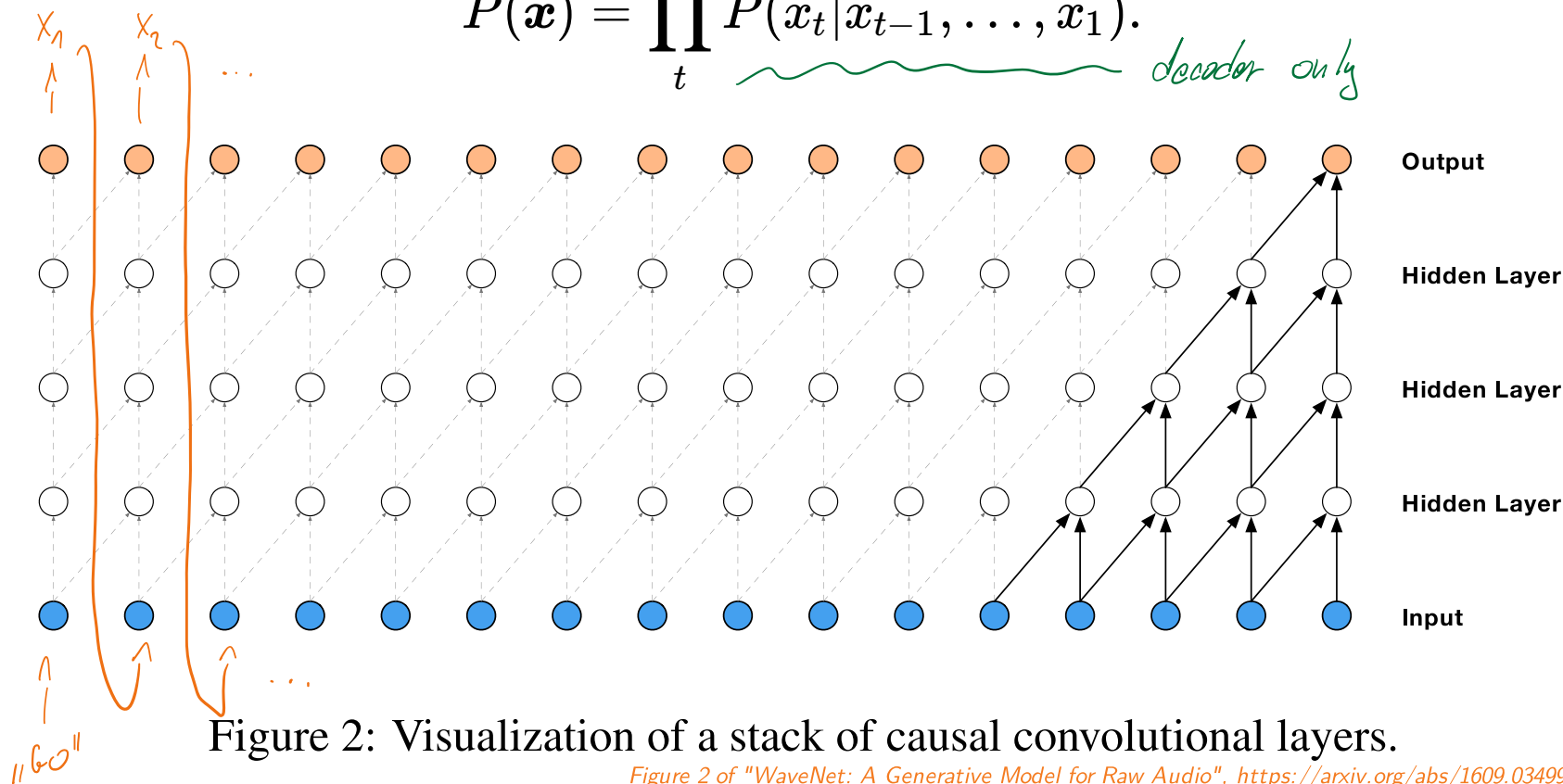


Figure 2 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

However, to achieve larger receptive field, we utilize **dilated** (or **atrous**) convolutions:

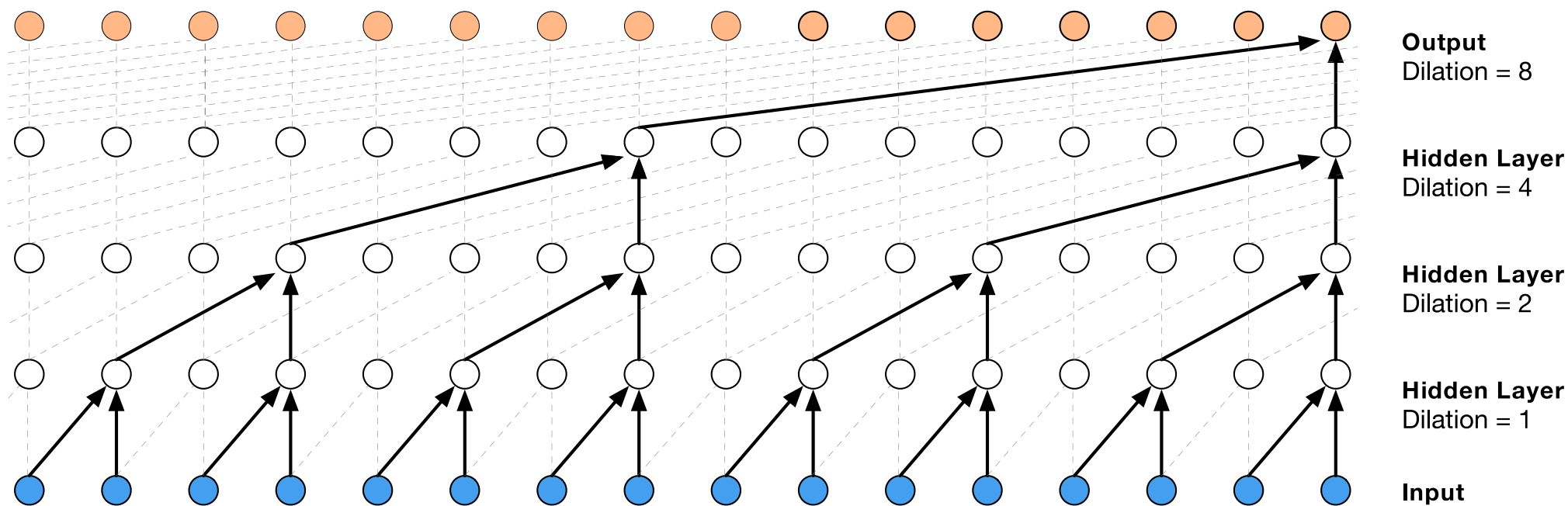
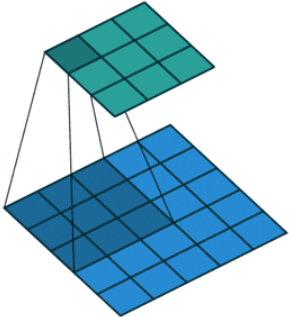
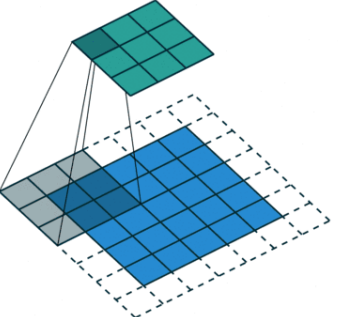
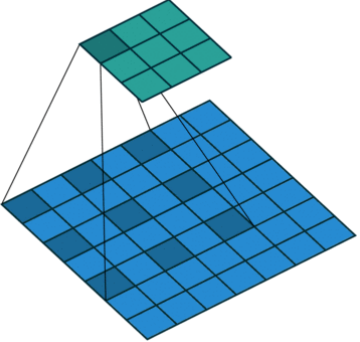
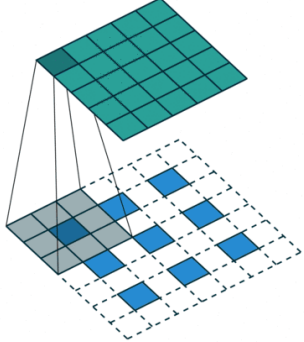


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Figure 3 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

# Dilated Versus Regular Versus Strided Convolutions

Regular Convolution	Strided Convolution
 <p data-bbox="340 583 869 611"><a href="https://github.com/vdumoulin/conv_arithmetic">https://github.com/vdumoulin/conv_arithmetic</a></p>	 <p data-bbox="1366 583 1895 611"><a href="https://github.com/vdumoulin/conv_arithmetic">https://github.com/vdumoulin/conv_arithmetic</a></p>
Dilated Convolution	Transposed Strided Convolution
 <p data-bbox="340 1115 869 1142"><a href="https://github.com/vdumoulin/conv_arithmetic">https://github.com/vdumoulin/conv_arithmetic</a></p>	 <p data-bbox="1366 1115 1895 1142"><a href="https://github.com/vdumoulin/conv_arithmetic">https://github.com/vdumoulin/conv_arithmetic</a></p>

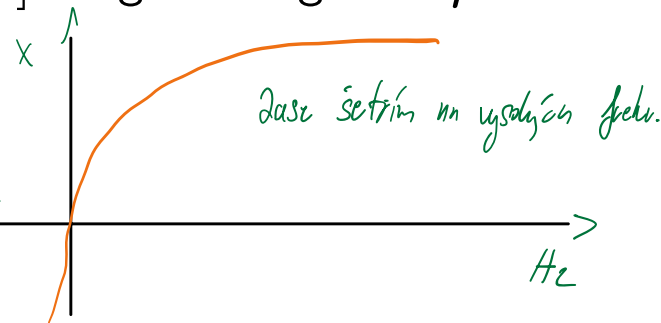
## Output Distribution *~ 70 years by x, later die gemacht?*

WaveNet generates audio with 16kHz frequency and 16-bit samples.

However, classification into 65 536 classes would not be efficient. Instead, WaveNet adopts the  $\mu$ -law transformation, which passes the input samples in  $[-1, 1]$  range through the  $\mu$ -law encoding

*→ positiv u negativ  
pro negativ pichos*

$$\text{sign}(x) \frac{\log(1 + 255|x|)}{\log(1 + 255)}, \sim$$



and the resulting  $[-1, 1]$  range is linearly quantized into 256 buckets.

The model therefore predicts each samples using classification into 256 classes, and then uses the inverse of the above transformation on the model predictions.



# WaveNet – Architecture

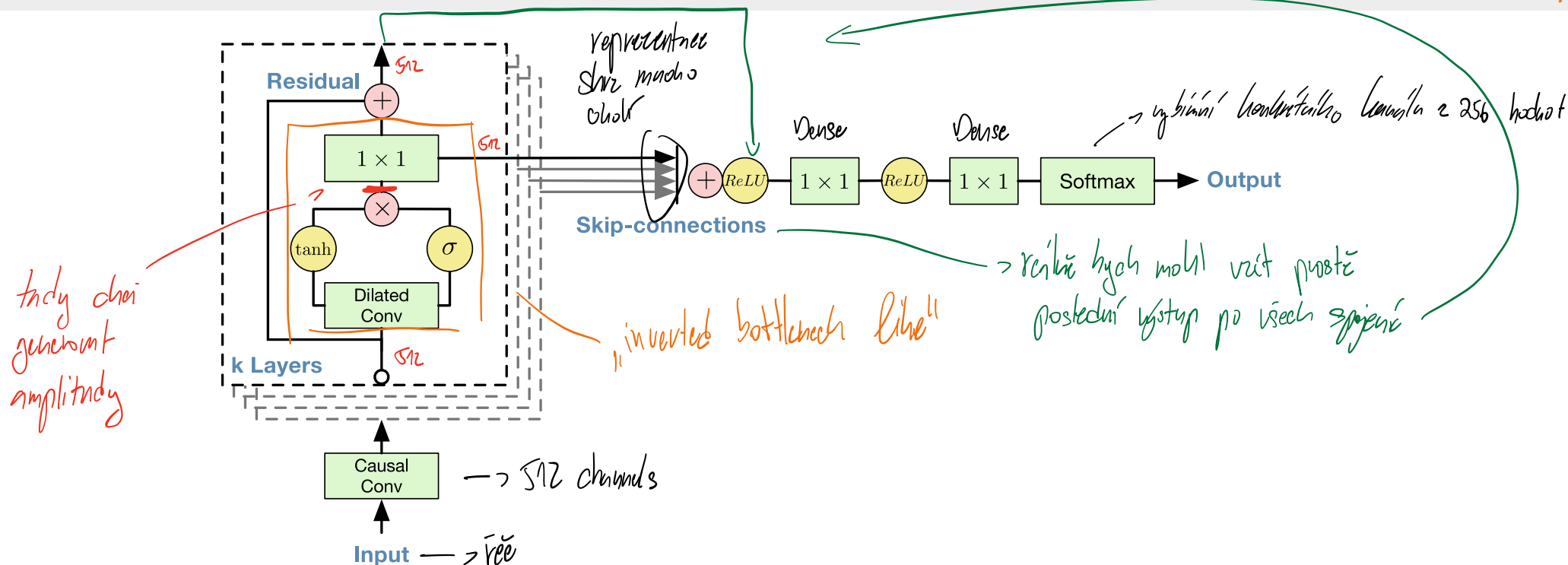


Figure 4: Overview of the residual block and the entire architecture.

Figure 4 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

The outputs of the dilated convolutions are passed through the *gated activation unit*:

$$z = \boxed{\tanh(W_f * x)} \odot \boxed{\sigma(W_g * x)}.$$

*jak by to mělo vypadat* (how it should look) and *jak více generovat?* (how to generate more?) are handwritten notes above the equation.

## Global Conditioning *→ vstupní pro CNN si rozšířme o naše conditions*

Global conditioning is performed by a single latent representation  $\mathbf{h}$ , changing the gated activation function to

$$\mathbf{z} = \tanh(\mathbf{W}_f * \mathbf{x} + \underline{\mathbf{V}_f \mathbf{h}}) \odot \sigma(\mathbf{W}_g * \mathbf{x} + \underline{\mathbf{V}_g \mathbf{h}}).$$

## Local Conditioning

For local conditioning, we are given a time series  $\mathbf{h}$ , possibly with a lower sampling frequency (for example the mel spectrogram). We first use transposed convolutions  $\mathbf{y} = f(\mathbf{h})$  to match resolution and then compute analogously to global conditioning

$$\mathbf{z} = \tanh(\mathbf{W}_f * \mathbf{x} + \mathbf{V}_f * \mathbf{y}) \odot \sigma(\mathbf{W}_g * \mathbf{x} + \mathbf{V}_g * \mathbf{y}).$$

The authors mention that using repetition instead of transposed convolution worked slightly worse.



The original paper did not mention hyperparameters, but later it was revealed that:

- 30 layers were used
  - grouped into 3 dilation stacks with 10 layers each
  - in a dilation stack, dilation rate increases by a factor of 2, starting with rate 1 and reaching maximum dilation of 512
- kernel size of a dilated convolution is 2 (and increased to 3 in Parallel WaveNet)
- residual connection has dimension 512
- gating layer uses 256+256 hidden units
- the  $1 \times 1$  convolutions in the output step produce 256 filters
- trained for 1 000 000 steps using Adam with a fixed learning rate of  $2e-4$

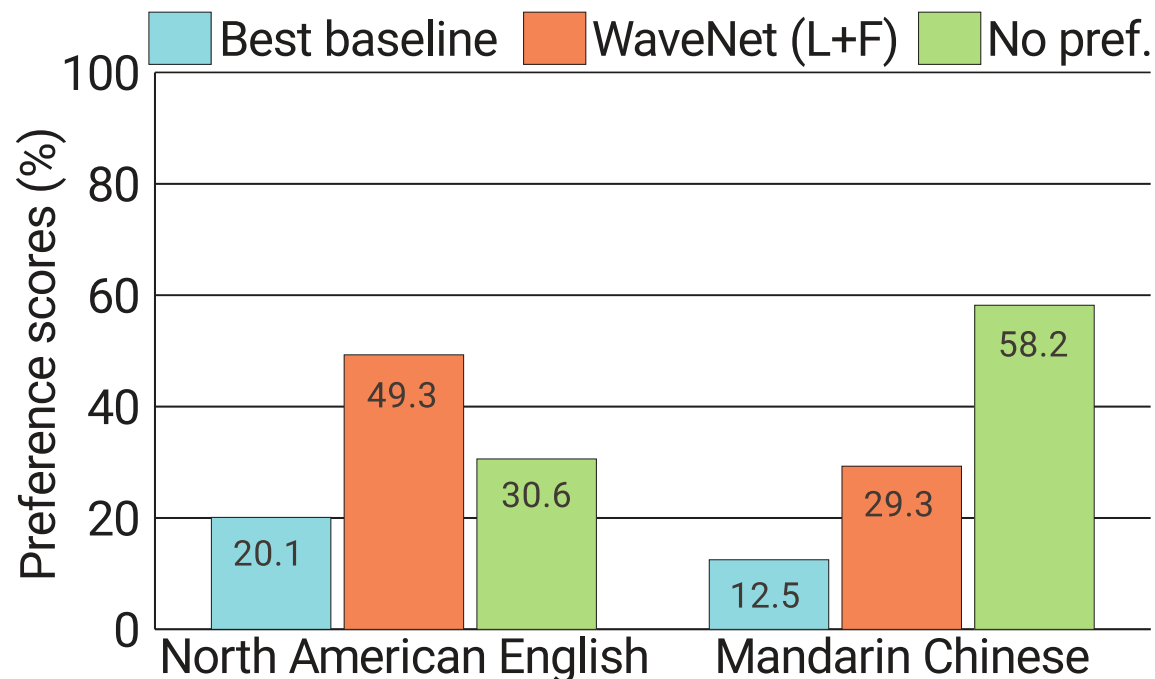


Figure 5: Subjective preference scores (%) of speech samples between (top) two baselines, (middle) two WaveNets, and (bottom) the best baseline and WaveNet. Note that LSTM and Concat correspond to LSTM-RNN-based statistical parametric and HMM-driven unit selection concatenative baseline synthesizers, and WaveNet (L) and WaveNet (L+F) correspond to the WaveNet conditioned on linguistic features only and that conditioned on both linguistic features and  $\log F_0$  values.

Figure 5 of "WaveNet: A Generative Model for Raw Audio", <https://arxiv.org/abs/1609.03499>

# Gated Activations in Transformers

# Gated Activations in Transformers

Similar gated activations seem to work the best in Transformers, in the FFN module.

Activation Name	Formula	$\text{FFN}(x; \mathbf{W}_1, \mathbf{W}_2)$
ReLU	$\max(0, x)$	$\max(0, x\mathbf{W}_1)\mathbf{W}_2$
GELU	$x\Phi(x)$	$\text{GELU}(x\mathbf{W}_1)\mathbf{W}_2$
Swish	$x\sigma(x)$	$\text{Swish}(x\mathbf{W}_1)\mathbf{W}_2$

There are several variants of the new gated activations:

Activation Name	Formula	$\text{FFN}(x; \mathbf{W}, \mathbf{V}, \mathbf{W}_2)$
GLU (Gated Linear Unit)	$\sigma(x\mathbf{W} + \mathbf{b}) \odot (x\mathbf{V} + \mathbf{c})$	$(\sigma(x\mathbf{W}) \odot x\mathbf{V})\mathbf{W}_2$
ReGLU	$\max(0, x\mathbf{W} + \mathbf{b}) \odot (x\mathbf{V} + \mathbf{c})$	$(\max(0, x\mathbf{W}) \odot x\mathbf{V})\mathbf{W}_2$
GEGLU	$\text{GELU}(x\mathbf{W} + \mathbf{b}) \odot (x\mathbf{V} + \mathbf{c})$	$(\text{GELU}(x\mathbf{W}) \odot x\mathbf{V})\mathbf{W}_2$
SwiGLU	$\text{Swish}(x\mathbf{W} + \mathbf{b}) \odot (x\mathbf{V} + \mathbf{c})$	$(\text{Swish}(x\mathbf{W}) \odot x\mathbf{V})\mathbf{W}_2$

# Gated Activations in Transformers

	Score Average	CoLA MCC	SST-2 Acc	MRPC F1	MRPC Acc	STSB PCC	STSB SCC	QQP F1	QQP Acc	MNLIm Acc	MNLImm Acc	QNLI Acc	RTE Acc	EM	F1
FFN <sub>ReLU</sub>	83.80	51.32	94.04	<b>93.08</b>	<b>90.20</b>	89.64	89.42	89.01	91.75	85.83	86.42	92.81	80.14	83.18	90.87
FFN <sub>GELU</sub>	83.86	53.48	94.04	92.81	<b>90.20</b>	89.69	89.49	88.63	91.62	85.89	86.13	92.39	80.51	83.09	90.79
FFN <sub>Swish</sub>	83.60	49.79	93.69	92.31	89.46	89.20	88.98	88.84	91.67	85.22	85.02	92.33	81.23	83.25	90.76
FFN <sub>GLU</sub>	84.20	49.16	94.27	92.39	89.46	89.46	89.35	88.79	91.62	86.36	86.18	92.92	<b>84.12</b>	82.88	90.69
FFN <sub>GEGLU</sub>	84.12	53.65	93.92	92.68	89.71	90.26	90.13	89.11	91.85	86.15	86.17	92.81	79.42	83.55	91.12
FFN <sub>Bilinear</sub>	83.79	51.02	<b>94.38</b>	92.28	89.46	90.06	89.84	88.95	91.69	<b>86.90</b>	<b>87.08</b>	92.92	81.95	<b>83.82</b>	91.06
FFN <sub>SwiGLU</sub>	84.36	51.59	93.92	92.23	88.97	<b>90.32</b>	<b>90.13</b>	<b>89.14</b>	<b>91.87</b>	86.45	86.47	<b>92.93</b>	83.39	83.42	91.03
FFN <sub>ReGLU</sub>	<b>84.67</b>	<b>56.16</b>	<b>94.38</b>	92.06	89.22	89.97	89.85	88.86	91.72	86.20	86.40	92.68	81.59	83.53	<b>91.18</b>

Table 2 of "GLU Variants Improve Transformer", <https://arxiv.org/abs/2002.05202> Table 4 of "GLU Variants Improve Transformer", <https://arxiv.org/abs/2002.05202>

Model	Params	Ops	Step/s	Early loss	Final loss	SGLUE	XSum	WebQ	WMT EnDe
Vanilla Transformer	223M	11.1T	3.50	2.182 ± 0.005	1.838	71.66	17.78	23.02	26.62
GeLU	223M	11.1T	3.58	2.179 ± 0.003	1.838	<b>75.79</b>	<b>17.86</b>	<b>25.13</b>	26.47
Swish	223M	11.1T	3.62	2.186 ± 0.003	1.847	<b>73.77</b>	17.74	<b>24.34</b>	<b>26.75</b>
ELU	223M	11.1T	3.56	2.270 ± 0.007	1.932	67.83	16.73	23.02	26.08
GLU	223M	11.1T	3.59	2.174 ± 0.003	<b>1.814</b>	<b>74.20</b>	<b>17.42</b>	24.34	<b>27.12</b>
GeGLU	223M	11.1T	3.55	2.130 ± 0.006	<b>1.792</b>	<b>75.96</b>	<b>18.27</b>	<b>24.87</b>	<b>26.87</b>
ReGLU	223M	11.1T	3.57	2.145 ± 0.004	<b>1.803</b>	<b>76.17</b>	<b>18.36</b>	<b>24.87</b>	<b>27.02</b>
SeLU	223M	11.1T	3.55	2.315 ± 0.004	1.948	68.76	16.76	22.75	25.99
SwiGLU	223M	11.1T	3.53	2.127 ± 0.003	<b>1.789</b>	<b>76.00</b>	<b>18.20</b>	<b>24.34</b>	<b>27.02</b>
LiGLU	223M	11.1T	3.59	2.149 ± 0.005	<b>1.798</b>	<b>75.34</b>	<b>17.97</b>	<b>24.34</b>	26.53
Sigmoid	223M	11.1T	3.63	2.291 ± 0.019	1.867	<b>74.31</b>	17.51	23.02	26.30
Softplus	223M	11.1T	3.47	2.207 ± 0.011	1.850	<b>72.45</b>	17.65	<b>24.34</b>	<b>26.89</b>

Table 1 of "Do Transformer Modifications Transfer Across Implementations and Applications?", <https://arxiv.org/abs/2102.11972>

Prostě uveďme, platí, že gated units jsou obecně lepší než jen obyčejná velká nelinearita.

# Parallel WaveNet

# Parallel WaveNet

Parallel WaveNet is an improvement of the original WaveNet by the same authors.

First, the output distribution was changed from 256  $\mu$ -law values to a Mixture of Logistic (suggested in another paper – PixelCNN++, but reused in other architectures since):

*Logistic je nie dostatočne blízko  $\mathcal{N}$*   $x \sim \sum_i \pi_i \text{Logistic}(\mu_i, s_i)$  *tyhle tři hodnoty jede jednotka musí generovat*

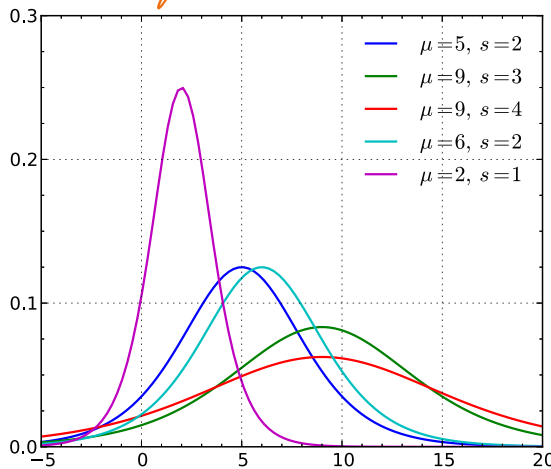
*↳ dělali bychom s  $\mathcal{N}$  (norm. rozdětím), ale potřebujeme vyjádřit PDF pomocí triviálních operací...*

The logistic distribution is a distribution with a  $\sigma$  as cumulative density function (where the mean and scale is parametrized by  $\mu$  and  $s$ ). Therefore, we can write

$$P(x|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{s}) = \sum_i \pi_i \left[ \sigma\left(\frac{x + 0.5 - \mu_i}{s_i}\right) - \sigma\left(\frac{x - 0.5 - \mu_i}{s_i}\right) \right],$$

where we replace  $-0.5$  and  $0.5$  in the edge cases by  $-\infty$  and  $\infty$ .

In Parallel WaveNet teacher, 10 mixture components are used.



<https://commons.wikimedia.org/wiki/File:Logisticpdffunction.svg>

# Parallel WaveNet - díky rekurentnímu zpracování trvalo jen. 1. sec. asi 100s.

Auto-regressive (sequential) inference is extremely slow in WaveNet.

Instead, we model  $P(x_t)$  as  $P(x_t | \mathbf{z}_{<t}) = \text{Logistic}(x_t; \mu^1(\mathbf{z}_{<t}), s^1(\mathbf{z}_{<t}))$  for a *random*  $\mathbf{z}$  drawn from a logistic distribution  $\text{Logistic}(\mathbf{0}, \mathbf{1})$ . Therefore, using the reparametrization trick,

$$x_t^1 = \mu^1(\mathbf{z}_{<t}) + \boxed{z_t} \cdot s^1(\mathbf{z}_{<t}).$$

určují ty parametry pouze

z všech dopředných dostupných  $x_{t+1}$

Usually, one iteration of the algorithm does not produce good enough results – consequently, 4 iterations were used by the authors. In further iterations,

je to sice o něco lepší,  
ale musíme to ještě paralelizovat

$$x_t^i = \mu^i(x_{<t}^{i-1}) + x_t^{i-1} \cdot s^i(x_{<t}^{i-1}).$$

z každé iterace bere výstup z předchozí iterace

After  $N$  iterations,  $P(\mathbf{x}_t^N | \mathbf{z}_{<t})$  is a logistic distribution with location  $\mu^{\text{tot}}$  and scale  $s^{\text{tot}}$ :

$$\mu_t^{\text{tot}} = \sum_{i=1}^N \mu^i(x_{<t}^{i-1}) \cdot \left( \prod_{j>i}^N s^j(x_{<t}^{j-1}) \right) \text{ and } s_t^{\text{tot}} = \prod_{i=1}^N s^i(x_{<t}^{i-1}),$$

where we have denoted  $\mathbf{z}$  as  $\mathbf{x}^0$  for convenience.



The consequences of changing the model from  $x_t = f(x_{t-1}, \dots, x_1)$  to

*a to je úplná katastrofa*

$$x_t^1 = \mu^1(\mathbf{z}_{<t}) + z_t \cdot s^1(\mathbf{z}_{<t})$$

$$x_t^i = \mu^i(\mathbf{x}_{<t}^{i-1}) + x_t^{i-1} \cdot s^i(\mathbf{x}_{<t}^{i-1})$$

$$\begin{array}{ccc} x_1 & x_2 & x_3 \\ | & | & | \\ z_0 & z_1 & z_2 \end{array}$$

are:

*to můžeme vzít z toho dříveho vstupu.*

*PAUL MÁM ALE AUTOREGRESSIVNÍ TRÉNDOVÁNÍ!*

*log p(x) → kde vzít to z toho z?*

- During inference, the prediction can be computed in parallel, because  $x_t^i$  depends only on  $\mathbf{x}_{<t}^{i-1}$ , not on  $\mathbf{x}_{<t}^i$ .
- However, we cannot perform training in parallel. If we try maximizing the log-likelihood of an input sequence  $\mathbf{x}^1$ , we need to find out which  $\mathbf{z}$  sequence generates it.
  - The  $z_1$  can be computed using  $x_1^1$ .
  - However,  $z_2$  depends not only on  $x_1^1$  and  $x_2^1$ , but also on  $z_1$ ; generally,  $z_t$  depends on  $\mathbf{x}^1$  and also on all  $\mathbf{z}_{<t}$ , and can be computed only sequentially.

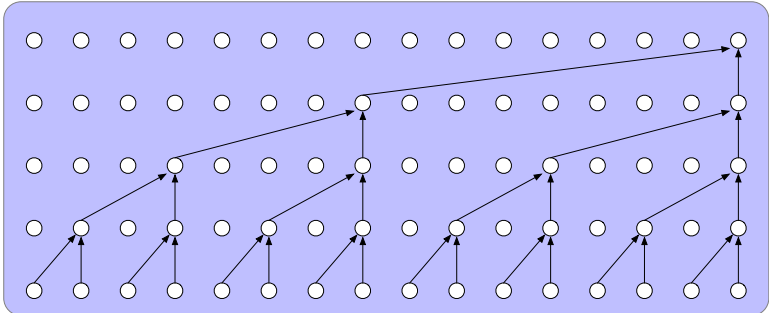
Therefore, WaveNet can perform parallel training and sequential inference, while the proposed model can perform parallel inference but sequential training.

# Probability Density Distillation

The authors propose to train the network by a **probability density distillation** using a teacher WaveNet (producing a mixture of logistic with 10 components) with KL-divergence as a loss.

## WaveNet Teacher

Linguistic features  $\dashrightarrow$



Teacher Output

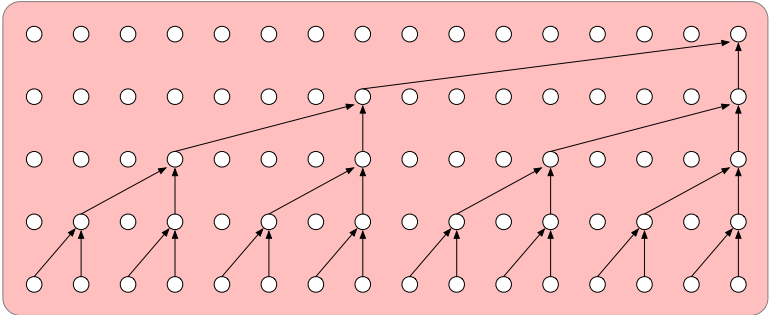
$$P(x_i | x_{<i})$$

Generated Samples

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

## WaveNet Student

Linguistic features  $\dashrightarrow$



Student Output

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

Input noise

$$z_i$$

Figure 2 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

# Probability Density Distillation

Therefore, instead of computing  $\mathbf{z}$  from some gold  $\mathbf{x}_g$ , we

- sample a random  $\mathbf{z}$ ;
- generate the output  $\mathbf{x}$ ;
- use the teacher WaveNet model to estimate the log-likelihood of  $\mathbf{x}$ ;
- update the student to match the log-likelihood of the teacher.

Denoting the teacher distribution as  $P_T$  and the student distribution as  $P_S$ , the loss is

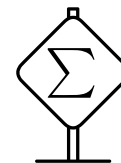
$$D_{\text{KL}}(P_S || P_T) = H(P_S, P_T) - \underbrace{H(P_S)}_{\text{tobľe se n o ale nem  e zohraovat, proto e student m  t h  j mno t .}}$$

Therefore, we do not only minimize cross-entropy, but we also try to keep the entropy of the student as high as possible – it is indeed crucial not to match just the mode of the teacher.

- Consider a teacher generating white noise, where every sample comes from  $\mathcal{N}(0, 1)$  – in this case, the cross-entropy loss of a constant **0**, complete silence, would be maximal.

In a sense, probability density distillation is similar to GANs. However, the teacher is kept fixed, and the student does not attempt to fool it but to match its distribution instead.

Because the entropy of a logistic distribution  $\text{Logistic}(\mu, s)$  is  $\log s + 2$ , the entropy term  $H(P_S)$  can be rewritten as follows:

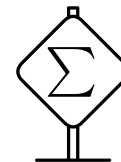


$$\begin{aligned} H(P_S) &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T -\log p_S(x_t | \mathbf{z}_{<t}) \right] \\ &= \mathbb{E}_{z \sim \text{Logistic}(0,1)} \left[ \sum_{t=1}^T \log s(\mathbf{z}_{<t}, \boldsymbol{\theta}) \right] + 2T. \end{aligned}$$

Therefore, this term can be computed without having to generate  $\mathbf{x}$ .

# Probability Density Distillation Details

However, the cross-entropy term  $H(P_S, P_T)$  requires sampling from  $P_S$  to estimate:



$$\begin{aligned}
 H(P_S, P_T) &= \int_{\mathbf{x}} -P_S(\mathbf{x}) \log P_T(\mathbf{x}) \\
 &= \sum_{t=1}^T \int_{\mathbf{x}} -P_S(\mathbf{x}) \log P_T(x_t | \mathbf{x}_{<t}) \\
 &= \sum_{t=1}^T \int_{\mathbf{x}} -\underbrace{P_S(\mathbf{x}_{<t})}_{\text{blue}} \underbrace{P_S(x_t | \mathbf{x}_{<t})}_{\text{green}} \underbrace{P_S(\mathbf{x}_{>t} | \mathbf{x}_{\leq t})}_{\text{red}} \log P_T(x_t | \mathbf{x}_{<t}) \\
 &= \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} \left[ \underbrace{\int_{x_t} -P_S(x_t | \mathbf{x}_{<t}) \log P_T(x_t | \mathbf{x}_{<t})}_{\text{green}} \underbrace{\int_{\mathbf{x}_{>t}} P_S(\mathbf{x}_{>t} | \mathbf{x}_{\leq t})}_{\text{red}} \right] \\
 &= \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} \left[ \overbrace{H\left(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t})\right)}^{\text{take to obtain just the green part of the sum}} \right]
 \end{aligned}$$

*takže to odvodiť jednoducho časťou vlním.*

$$H(P_S, P_T) = \sum_{t=1}^T \mathbb{E}_{P_S(\mathbf{x}_{<t})} H\left(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t})\right)$$

We can therefore estimate  $H(P_S, P_T)$  by:

- drawing a single sample  $\mathbf{x}$  from the student  $P_S$  [a *Logistic*( $\boldsymbol{\mu}^{\text{tot}}, \mathbf{s}^{\text{tot}}$ )],
- compute all  $P_T(x_t | \mathbf{x}_{<t})$  from the teacher in parallel [*mixture of logistic distributions*],
- and finally evaluate  $H(P_S(x_t | \mathbf{x}_{<t}), P_T(x_t | \mathbf{x}_{<t}))$  by sampling multiple different  $x_t$  from the  $P_S(x_t | \mathbf{x}_{<t})$ .

The authors state that this unbiased estimator has a much lower variance than naively evaluating a single sequence sample under the teacher using the original formulation.

# Parallel WaveNet

The Parallel WaveNet model consists of 4 iterations with 10, 10, 10, 30 layers, respectively. The dimension of the residuals and the gating units is 64 (compared to 512 in WaveNet).

The Parallel WaveNet generates over 500k samples per second, compared to ~170 samples per second of a regular WaveNet – more than a 1000 times speedup.

Method	Subjective 5-scale MOS
<b>16kHz, 8-bit <math>\mu</math>-law, 25h data:</b>	
LSTM-RNN parametric [27]	$3.67 \pm 0.098$
HMM-driven concatenative [27]	$3.86 \pm 0.137$
WaveNet [27]	$4.21 \pm 0.081$
<b>24kHz, 16-bit linear PCM, 65h data:</b>	
HMM-driven concatenative	$4.19 \pm 0.097$
Autoregressive WaveNet	$4.41 \pm 0.069$
Distilled WaveNet	$4.41 \pm 0.078$

Table 1 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

For comparison, using a single iteration with 30 layers achieve MOS of 4.21.

The Parallel WaveNet can be trained to generate speech of multiple speakers (using the global conditioning). Because such a model needs larger capacity, it used 30 layers in every iteration (instead of 10, 10, 10, 30).

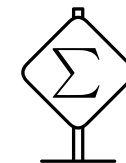
	Parametric	Concatenative	Distilled WaveNet
<b>English speaker 1</b> (female - 65h data)	3.88	4.19	4.41
<b>English speaker 2</b> (male - 21h data)	3.96	4.09	4.34
<b>English speaker 3</b> (male - 10h data)	3.77	3.65	4.47
<b>English speaker 4</b> (female - 9h data)	3.42	3.40	3.97
<b>Japanese speaker</b> (female - 28h data)	4.07	3.47	4.23

Table 2: Comparison of MOS scores on English and Japanese with multi-speaker distilled WaveNets. Note that some speakers sounded less appealing to people and always get lower MOS, however distilled parallel WaveNet always achieved significantly better results.

Table 2 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>



To generate high-quality audio, the probability density distillation is not entirely sufficient. The authors therefore introduce additional losses:

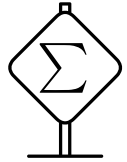


- **power loss:** ensures the power in different frequency bands is on average similar between the generated speech and human speech. For a conditioned training data  $(\mathbf{x}, \mathbf{c})$  and WaveNet student  $g$ , the loss is

$$\| \text{STFT}(g(\mathbf{z}, \mathbf{c})) - \text{STFT}(\mathbf{x}) \|^2.$$

- **perceptual loss:** apart from the power in frequency bands, we can use a pre-trained classifier to extract features from generated and human speech and add a loss measuring their difference. The authors propose the loss as squared Frobenius norm of differences between Gram matrices (uncentered covariance matrices) of features of a WaveNet-like classifier predicting phones from raw audio.
- **contrastive loss:** to make the model respect the conditioning instead of generating outputs with high likelihood independent on the conditioning, the authors propose a contrastive distillation loss ( $\gamma = 0.3$  is used in the paper):

$$D_{\text{KL}}(P_S(\mathbf{c}_1) || P_T(\mathbf{c}_1)) - \gamma D_{\text{KL}}(P_S(\mathbf{c}_1) || P_T(\mathbf{c}_2)).$$



Method	Preference Scores versus baseline concatenative system Win - Lose - Neutral
<b>Losses used</b>	
KL + Power	60% - 15% - 25%
KL + Power + Perceptual	66% - 10% - 24%
KL + Power + Perceptual + Contrastive (= default)	65% - 9% - 26%

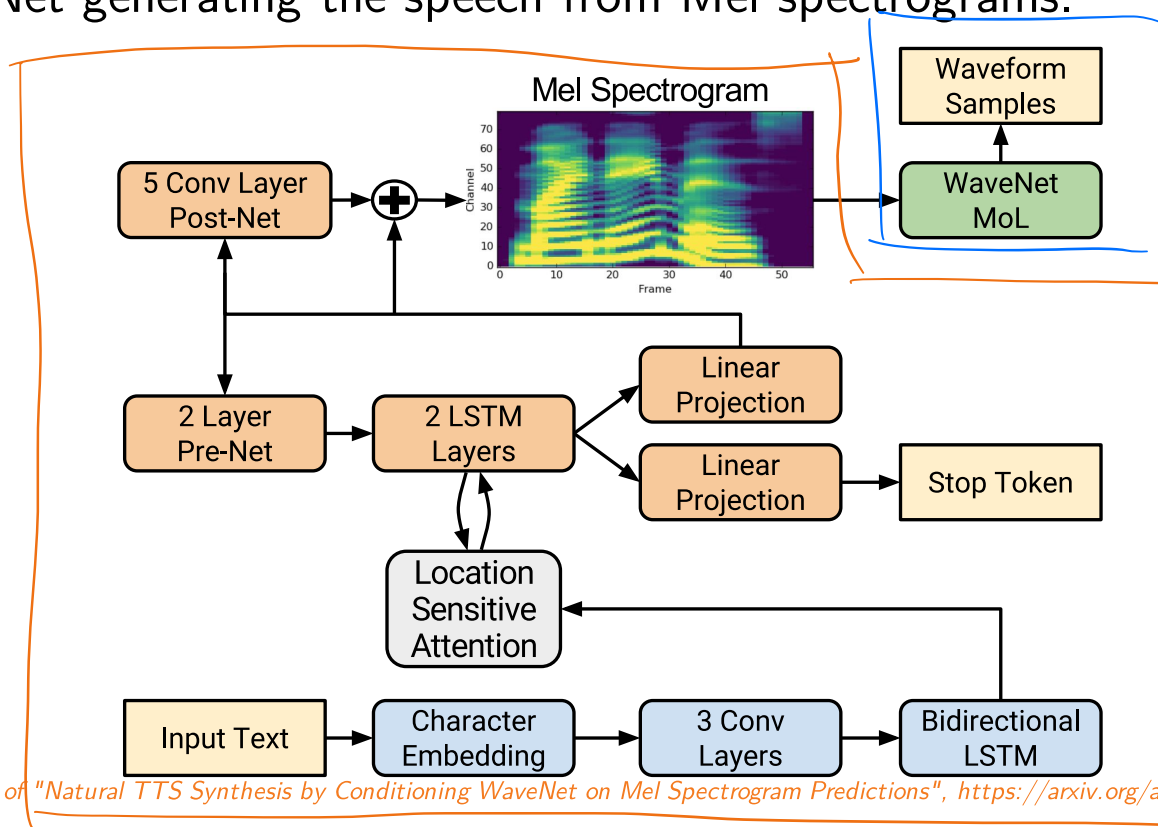
Table 3: Performance with respect to different combinations of loss terms. We report preference comparison scores since their mean opinion scores tend to be very close and inconclusive.

Table 3 of "Parallel WaveNet: Fast High-Fidelity Speech Synthesis", <https://arxiv.org/abs/1711.10433>

# Tacotron 2

Tacotron 2 model presents end-to-end speech synthesis directly from text. It consists of two components trained separately:

- a seq2seq model processing input characters and generating mel spectrograms;
- a Parallel WaveNet generating the speech from Mel spectrograms.



- oddělené systémy  
- vygenerují spektrum  
- pak převede na zvuk

Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

The Mel spectrograms used in Tacotron 2 are fairly standard:

- The authors propose a frame size of 50ms, 12.5ms frame hop, and a Hann window.
- STFT magnitudes are transformed into 80-channel Mel scale spanning 175Hz to 7.6kHz, followed by a log dynamic range compression (clipping input values to at least 0.01).

## Architecture

The characters are represented using 512-dimensional embeddings, processed by 3 Conv+BN+ReLU with kernel size 5 and 512 channels, following by a single bi-directional LSTM with 512 units.

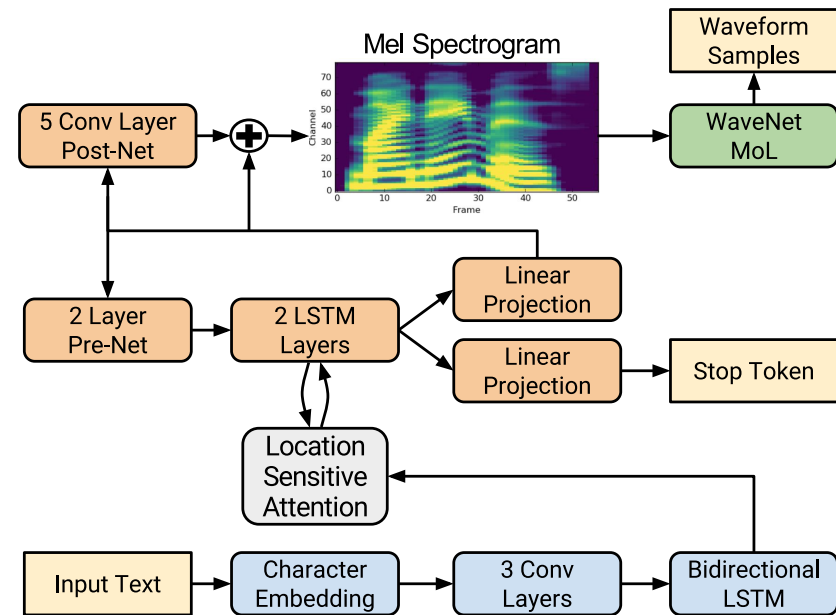


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

To make sequential processing of input characters easier, Tacotron 2 utilizes *location-sensitive attention*, which is an extension of the additive attention. While the additive (Bahdanau) attention computes

$$\alpha_i = \text{Attend}(s_{i-1}, h),$$

$$\alpha_{ij} = \text{softmax}(\mathbf{v}^\top \tanh(\mathbf{V}h_j + \mathbf{W}s_{i-1} + \mathbf{b})),$$

the location-sensitive attention also inputs the previous cumulative attention weights  $\hat{\alpha}_i$  into the current attention computation:

$$\alpha_i = \text{Attend}(s_{i-1}, h, \hat{\alpha}_{i-1}).$$

In detail, the previous attention weights are processed by a 1-D convolution with kernel  $\mathbf{F}$ :

$$\alpha_{ij} = \text{softmax}(\mathbf{v}^\top \tanh(\mathbf{V}h_j + \mathbf{W}s_{i-1} + (\mathbf{F} \overset{\text{conv.}}{*} \hat{\alpha}_{i-1})_j + \mathbf{b})),$$

using 32-dimensional 1D convolution with kernel size 31. *→ takhle dávám, aby attention věděla, na co attendovala v minulém kroku*

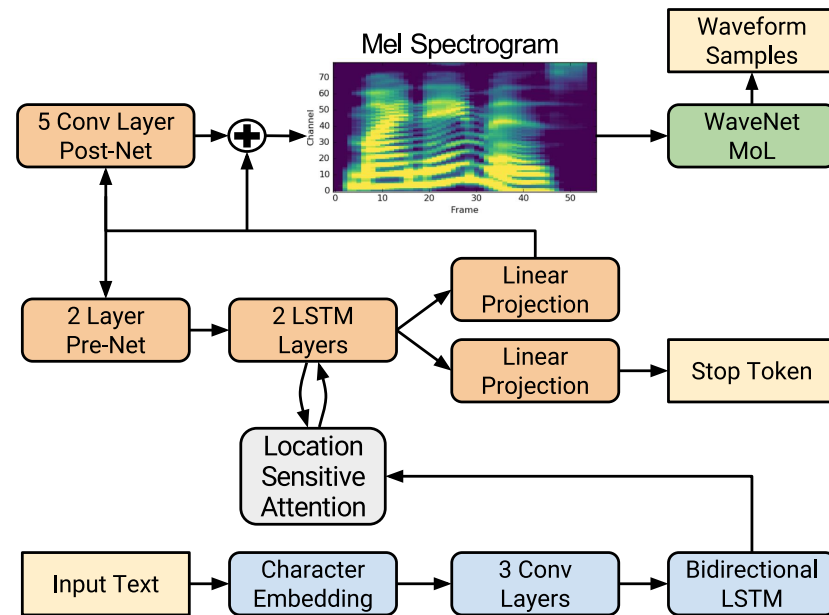


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

Usually, the current state in the attention  $\mathbf{s}_{i-1}$  is represented using the decoder cell state. However, in this model, it is beneficial for the attention to have access to the generated spectrogram.

Therefore, a separate *attention RNN* is used, which has to goal to represent the “mel spectrogram generated so far”: its inputs are the decoder input (the result of the pre-net) concatenated with the attention context vector  $\alpha_i^T \mathbf{h}$ , and the output of this attention RNN is used as  $\mathbf{s}_{i-1}$  in the attention computation.

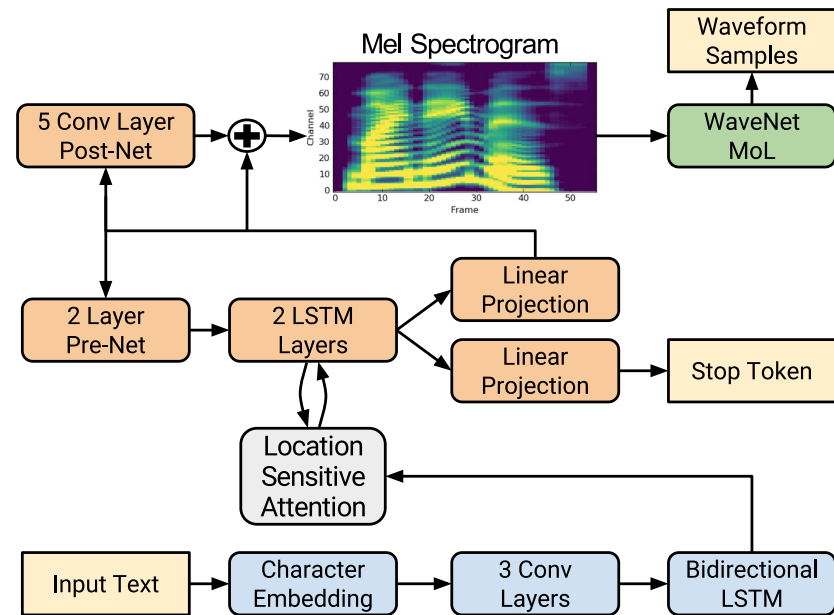


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

The decoder predicts the spectrogram one frame at a time. The predictions from the previous step are first passed through a *pre-net* composed of 2 fully-connected ReLU layers with 256 units and concatenated with attention context vector.

The decoder consists of 2 1024-dimensional LSTM cells and its output is linearly projected to the predicted frame.

The stop-token is predicted from a concatenation of decoder output and attention context vector, followed by a sigmoid activation.

The predicted spectrogram frame is post-processed by a *post-net* composed of 5 convolutional layers followed by batch normalization and tanh activation (on all but the last layers), with 512 channels and kernel size 5.

The target objective is the MSE error of the spectrogram and the decoder output both before and after the post-net.

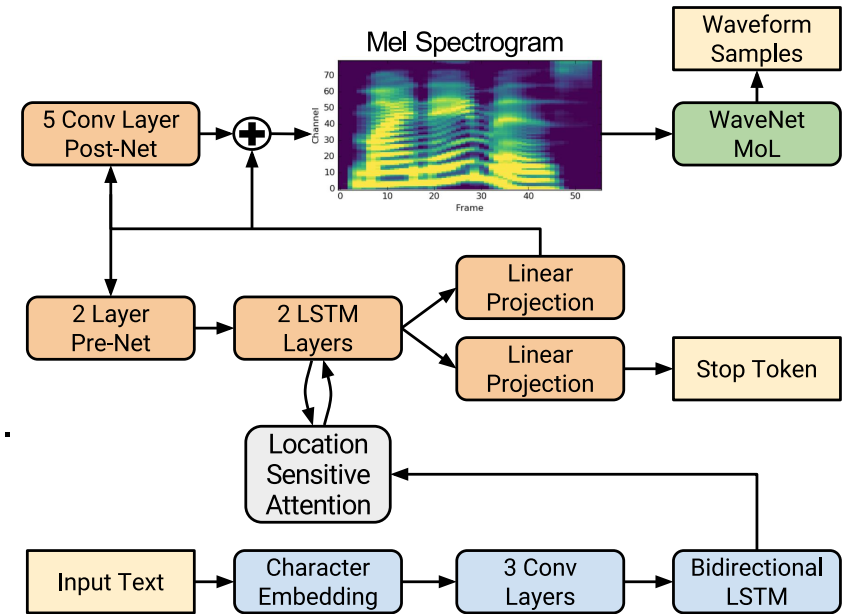


Figure 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>



System	MOS
Parametric	$3.492 \pm 0.096$
Tacotron (Griffin-Lim)	$4.001 \pm 0.087$
Concatenative	$4.166 \pm 0.091$
WaveNet (Linguistic)	$4.341 \pm 0.051$
Ground truth	$4.582 \pm 0.053$
<b>Tacotron 2 (this paper)</b>	<b><math>4.526 \pm 0.066</math></b>

Table 1 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

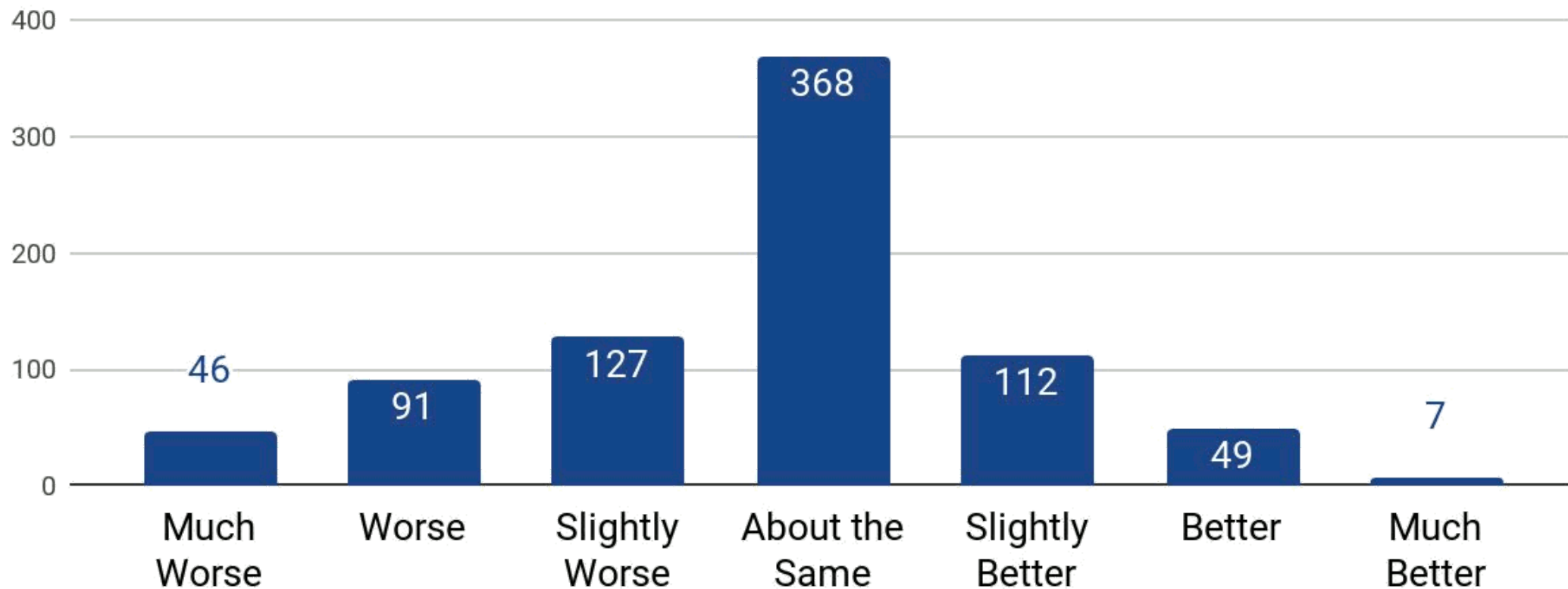


Figure 2 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

You can listen to samples at <https://google.github.io/tacotron/publications/tacotron2/>

Training	Synthesis	
	Predicted	Ground truth
Predicted	$4.526 \pm 0.066$	$4.449 \pm 0.060$
Ground truth	$4.362 \pm 0.066$	$4.522 \pm 0.055$

**Table 2.** Comparison of evaluated MOS for our system when WaveNet trained on predicted/ground truth mel spectrograms are made to synthesize from predicted/ground truth mel spectrograms.

Table 2 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

System	MOS
Tacotron 2 (Linear + G-L)	$3.944 \pm 0.091$
Tacotron 2 (Linear + WaveNet)	$4.510 \pm 0.054$
Tacotron 2 (Mel + WaveNet)	<b><math>4.526 \pm 0.066</math></b>

**Table 3.** Comparison of evaluated MOS for Griffin-Lim vs. WaveNet as a vocoder, and using 1,025-dimensional linear spectrograms vs. 80-dimensional mel spectrograms as conditioning inputs to WaveNet.

Table 3 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

Total layers	Num cycles	Dilation cycle size	Receptive field (samples / ms)	MOS
30	3	10	6,139 / 255.8	$4.526 \pm 0.066$
24	4	6	505 / 21.0	$4.547 \pm 0.056$
12	2	6	253 / 10.5	$4.481 \pm 0.059$
30	30	1	61 / 2.5	$3.930 \pm 0.076$

**Table 4.** WaveNet with various layer and receptive field sizes.

Table 4 of "Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions", <https://arxiv.org/abs/1712.05884>

# FastSpeech

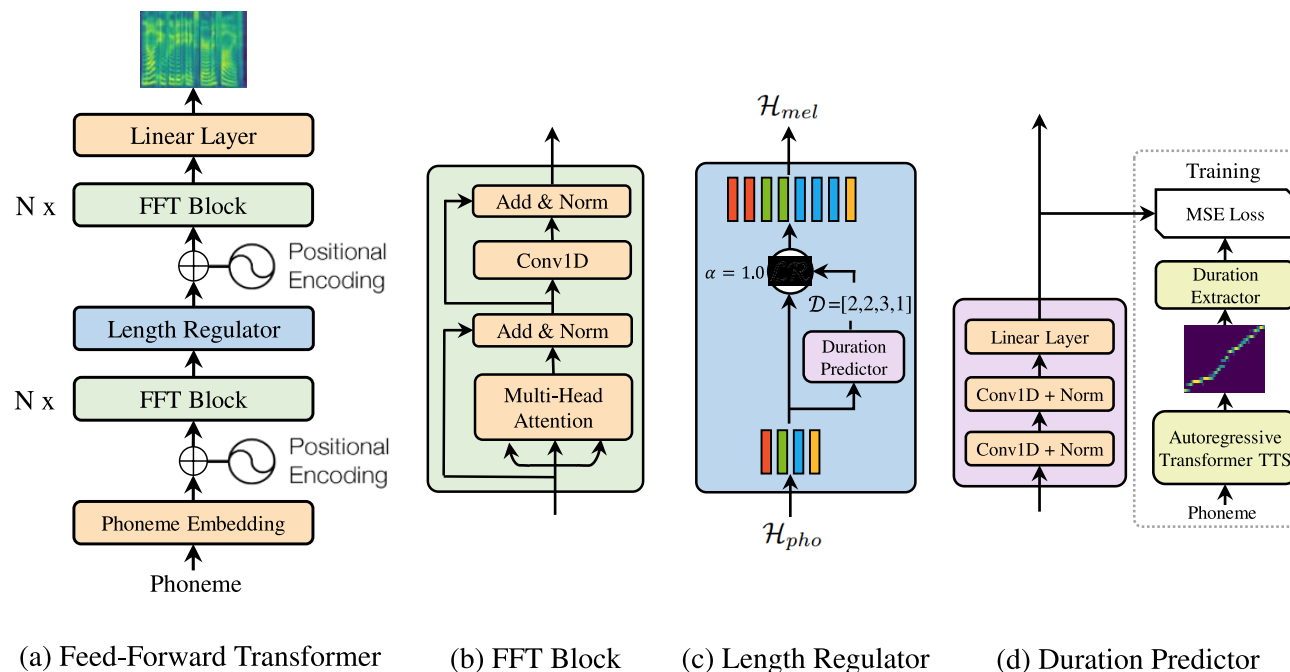


Figure 1: The overall architecture for FastSpeech. (a). The feed-forward Transformer. (b). The feed-forward Transformer block. (c). The length regulator. (d). The duration predictor. MSE loss denotes the loss between predicted and extracted duration, which only exists in the training process.

Figure 1 of "FastSpeech: Fast, Robust and Controllable Text to Speech", <https://arxiv.org/abs/1905.09263>

FastSpeech performs non-autoregressive decoding of mel spectrogram from the input text.

A hard character–spectrogram frame alignment is computed using some existing system, and the representation of input character is then repeated according to the alignment.

Method	MOS
<i>GT</i>	$4.41 \pm 0.08$
<i>GT (Mel + WaveGlow)</i>	$4.00 \pm 0.09$
<i>Tacotron 2 [22] (Mel + WaveGlow)</i>	$3.86 \pm 0.09$
<i>Merlin [28] (WORLD)</i>	$2.40 \pm 0.13$
<i>Transformer TTS [14] (Mel + WaveGlow)</i>	$3.88 \pm 0.09$
<i>FastSpeech (Mel + WaveGlow)</i>	$3.84 \pm 0.08$

Table 1: The MOS with 95% confidence intervals.

Table 1 of "FastSpeech: Fast, Robust and Controllable Text to Speech", <https://arxiv.org/abs/1905.09263>

Method	Repeats	Skips	Error Sentences	Error Rate
<i>Tacotron 2</i>	4	11	12	24%
<i>Transformer TTS</i>	7	15	17	34%
<i>FastSpeech</i>	0	0	0	0%

Table 3: The comparison of robustness between FastSpeech and other systems on the 50 particularly hard sentences. Each kind of word error is counted at most once per sentence.

Table 3 of "FastSpeech: Fast, Robust and Controllable Text to Speech", <https://arxiv.org/abs/1905.09263>

Later improved in FastPitch and FastSpeech 2 systems.

# One TTS Alignment To Rule Them All

# One TTS Alignment To Rule Them All

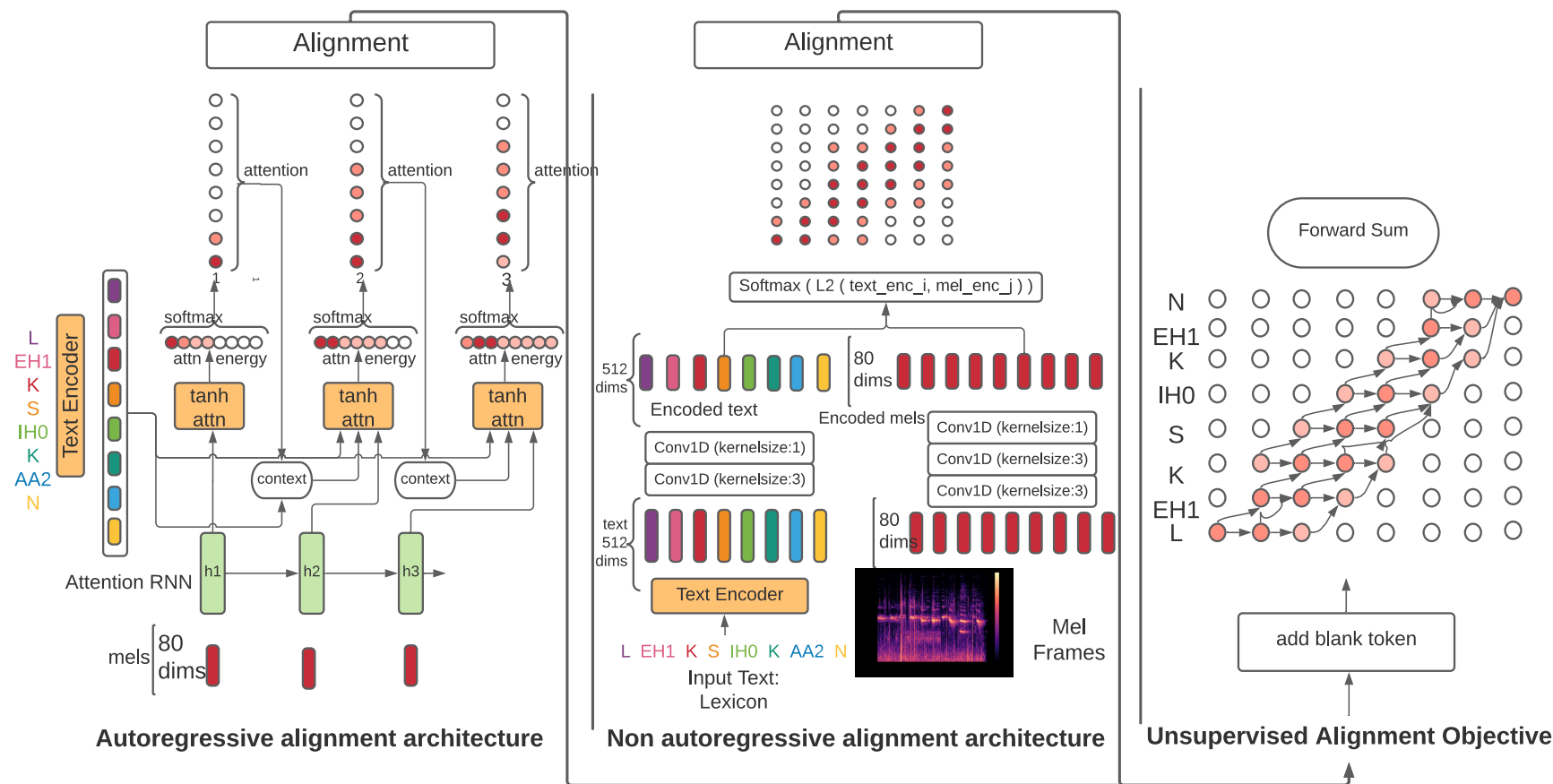


Figure 1: Overview of our Alignment Learning Framework: autoregressive models use a sequential attention mechanism to generate alignments between text and mels. Non-autoregressive models encode text and mels using simple 1D convolutions and use pairwise  $L_2$  distance to compute the alignments. The alignments represent the distribution  $P(s_t|x_t)$  and the alignment objective (Equation 1).

Figure 1 of "One TTS Alignment To Rule Them All", <https://arxiv.org/abs/2108.10447>



# One TTS Alignment To Rule Them All

Let  $s_i$  denotes the  $i$ -th input text character, and  $\mathbf{x}_t$  be the  $t$ -th mel spectrogram frame.

Assuming we have access to some alignment

$$P(s_i | \mathbf{x}_t; \theta).$$

We can train the alignment to maximize the probability of all monotonic alignments between the mel spectrogram and the input text:

$$\mathcal{L}(\theta; \mathbf{X}, \mathbf{s}) \stackrel{\text{def}}{=} \sum_{\substack{\mathbf{c} \text{ any valid} \\ \text{monotonic alignment}}} \prod_{t=1}^T P(c_t | \mathbf{x}_t; \theta).$$

Such a loss can be computed by exploiting the CTC loss (given that efficient implementations for it already exist).

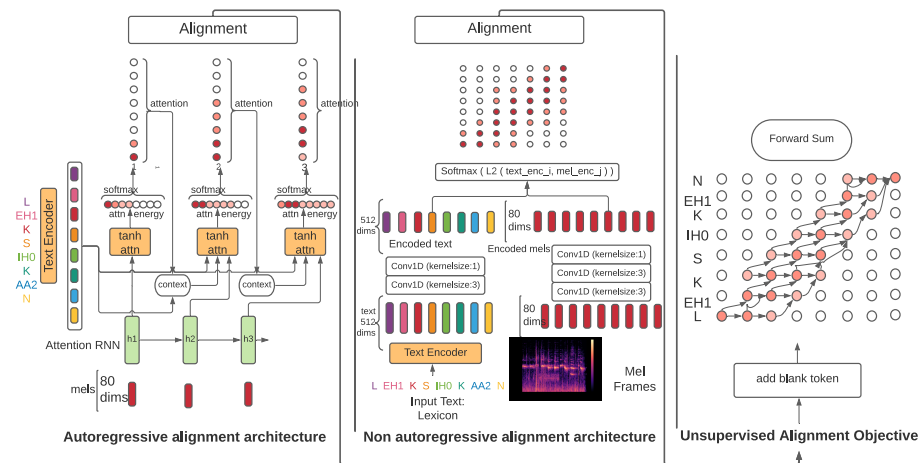


Figure 1: Overview of our Alignment Learning Framework: autoregressive models use a sequential attention mechanism to generate alignments between text and mels. Non-autoregressive models encode text and mels using simple 1D convolutions and use pairwise  $L_2$  distance to compute the alignments. The alignments represent the distribution  $P(s_t | \mathbf{x}_t)$  and the alignment objective (Equation 1).

Figure 1 of "One TTS Alignment To Rule Them All", <https://arxiv.org/abs/2108.10447>

# One TTS Alignment To Rule Them All

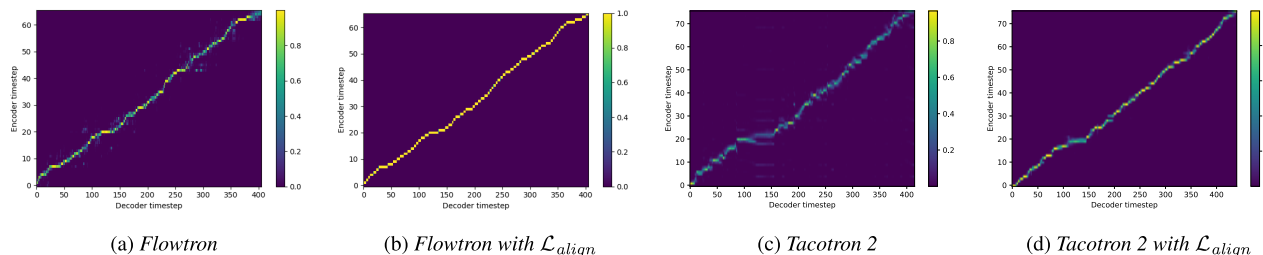


Figure 3: Converged soft alignments for Flowtron, Tacotron2. Alignment framework provides sharper and more connected alignments.

fidet and continuous alignments, and by extension, continuous speech without repeating or missing words.

Table 1: Pairwise preference scores judged by human raters, shown with 95% confidence intervals. Scores above 0.5 indicate models trained with  $\mathcal{L}_{align}$  were preferred by majority of raters.

Figure 3 of "One TTS Alignment To Rule Them All", <https://arxiv.org/abs/2108.10447>

Table 1: Pairwise preference scores judged by human raters, shown with 95% confidence intervals. Scores above 0.5 indicate models trained with  $\mathcal{L}_{align}$  were preferred by majority of raters.

Model	Alignment Framework vs Baseline
Tacotron 2	0.556 $\pm$ 0.068
Flowtron ( $\sigma = .5$ )	0.635 $\pm$ 0.065
RAD-TTS ( $\sigma = .5$ )	0.639 $\pm$ 0.066
FastPitch	0.565 $\pm$ 0.068
FastSpeech2	0.521 $\pm$ 0.067

Table 1 of "One TTS Alignment To Rule Them All", <https://arxiv.org/abs/2108.10447>

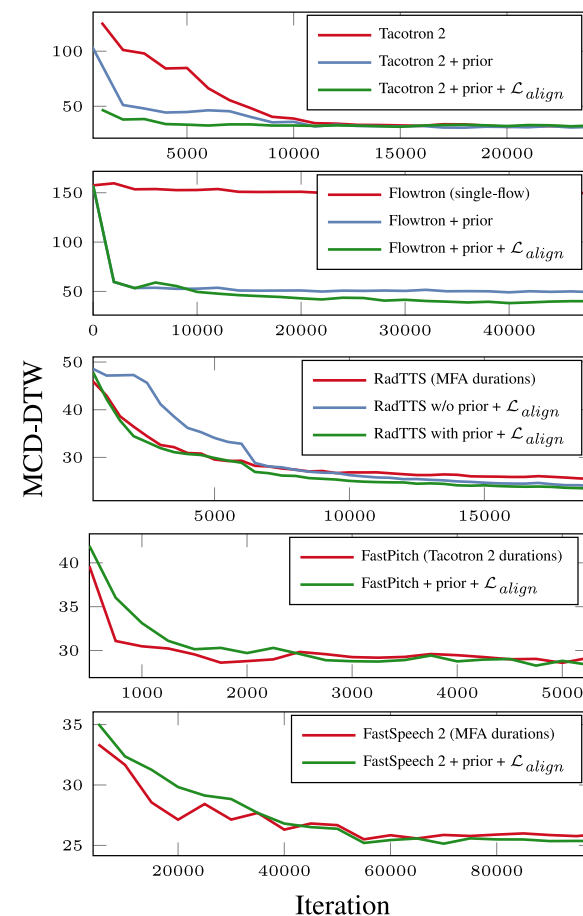
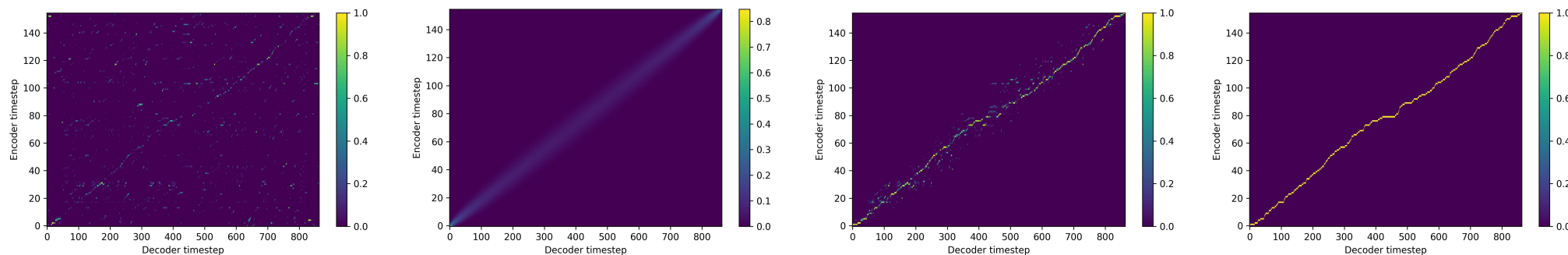


Figure 2: Convergence rate improvements in TTS models with the alignment learning framework

Figure 2 of "One TTS Alignment To Rule Them All", <https://arxiv.org/abs/2108.10447>

For faster alignment convergence (which makes the full model training faster), during training we might use a static 2D prior that is wider near the center and narrower near the corners.

The complete alignment is then the normalized product of the computed alignment and the alignment prior.



(a) Soft alignment  $\mathcal{A}_{soft}$

(b) Beta-binomial Prior

(c)  $\mathcal{A}_{soft}$  w/ Prior

(d) Hard alignment  $\mathcal{A}_{hard}$

Figure 2. Visualizations of the alignment attention matrices  $\mathcal{A}$ . The vertical axis represents text tokens from bottom to top. The horizontal axis represents mel-frames from left to right. Fig. 2a shows the baseline soft attention map. Fig. 2c combines Fig. 2a with Fig. 2b to penalize alignments straying too far from the diagonal. Fig. 2d is the most likely monotonic alignment extracted from Fig. 2c with Viterbi.

Figure 2 of "RAD-TTS: Parallel Flow-Based TTS with Robust Alignment Learning and Diverse Synthesis", <https://openreview.net/forum?id=0NQwnnwAORi>