

# **Natural TTS Synthesis by Conditioning WaveNet on Mel Spectrogram Predictions**

Подготовил:  
Косса Николай Евгеньевич, БПМИ202

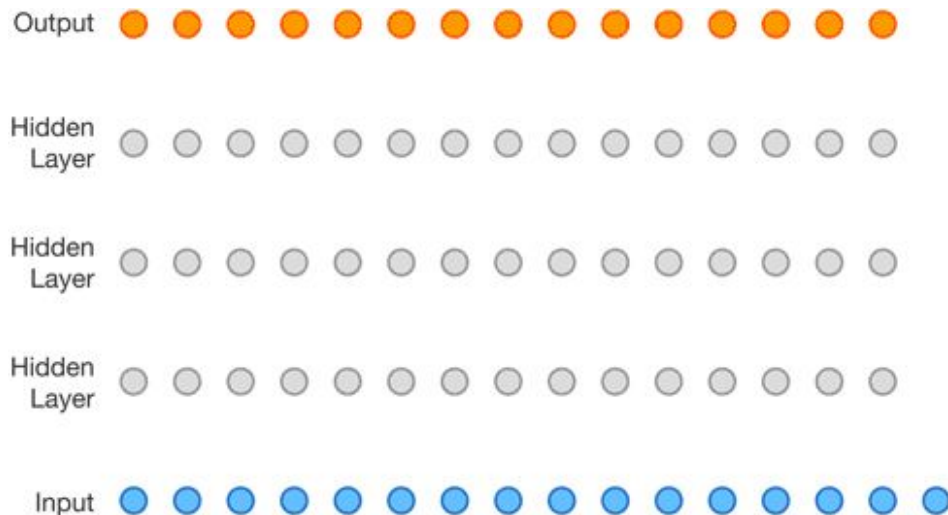
# План



1. WaveNet
2. Tacotron
3. Спектрограммы
4. Tacotron 2:
  - архитектура
  - modified WaveNet vocoder
  - эксперименты

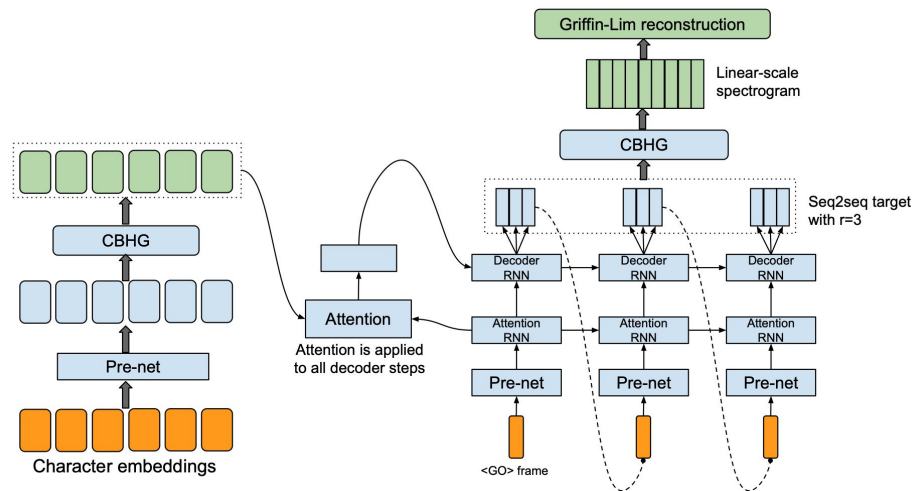
# WaveNet

- Входные данные: лингвистические признаки,  $f_0$ , длительность фонем.
- Средняя оценка MOS: 4.34.



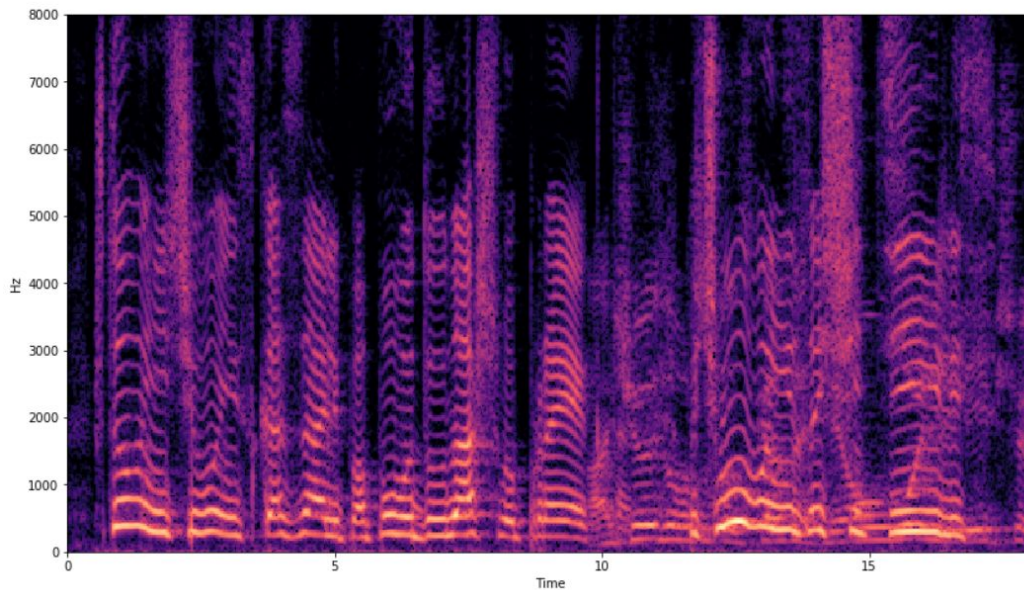
# Tacotron

- Входные данные: текст.
- Выход: mel-спектрограмма.
- Алгоритм Гриффина-Лима.
- Средняя оценка MOS: 4.0.



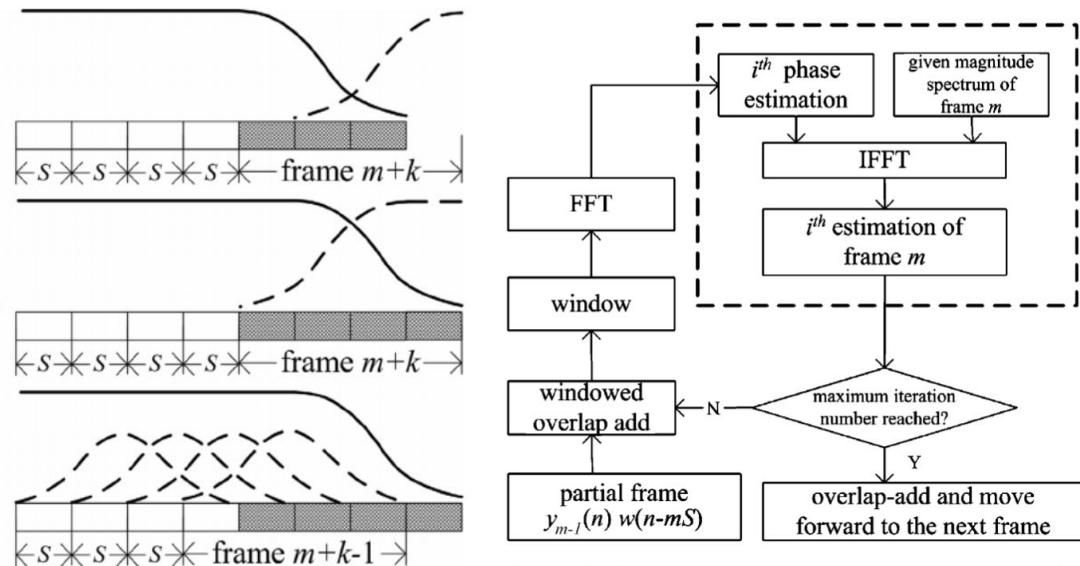
# Спектрограммы

- Спектрограмма не содержит информации о фазе.
- Фаза зависит от времени.



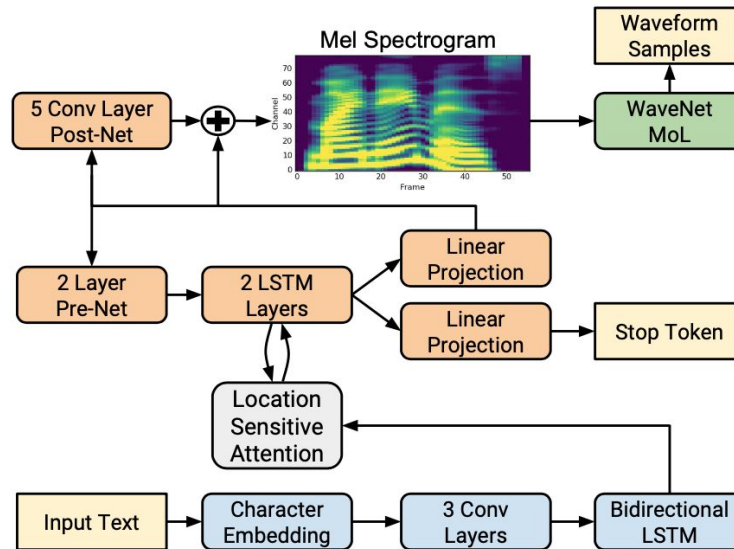
# Спектрограммы

- Алгоритм Гриффина-Лима (GLA) пользуется избыточностью STFT.



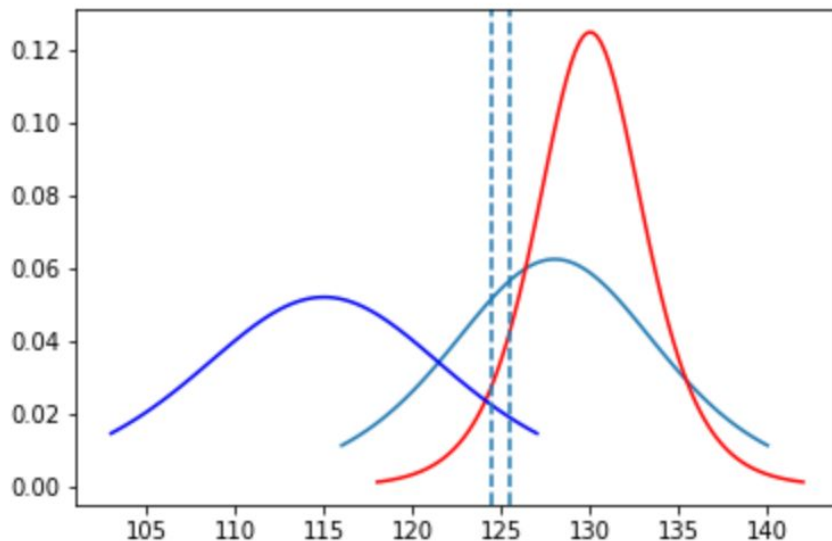
# Tacotron 2

- В отличие от первой версии, tacotron 2 не нуждается в GLA.
- Аудио предсказывает модифицированная WaveNet.



# Tacotron 2

- Вместо softmax используется 10-компонентная смесь логистических распределений (MoL).





# Tacotron 2

- Качество итоговой модели выше baseline.

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	$4.582 \pm 0.053$
Tacotron 2 (this paper)	<b><math>4.526 \pm 0.066</math></b>

# Tacotron 2

- Обучение на предсказанных спектрограммах лучше, чем на реальных.

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

# Tacotron 2

- WaveNet качественнее GLA.
- Mel-спектрограммы строго лучше линейных.

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>

# Материалы



- <https://arxiv.org/pdf/1712.05884.pdf>
- <https://arxiv.org/pdf/1609.03499v2.pdf>
- [an-explanation-of-discretized-logistic-mixture-likelihood](#)
- <https://wiki.aalto.fi>

# Вопросы?

