

Tacotron: Towards End-to-End Speech Synthesis

Speech Synthesis

- Задача синтеза – перевод текста в речь.
- Речь должна:
 - Передовать смысл сообщения.
 - Не допускать ошибок.
 - Звучать натурально.
- Критерий качества – Mean Opinion Score (MOS)

Подходы к синтезу

- Concatenative synthesis
- Parametric synthesis

ПОДХОДЫ К СИНТЕЗУ

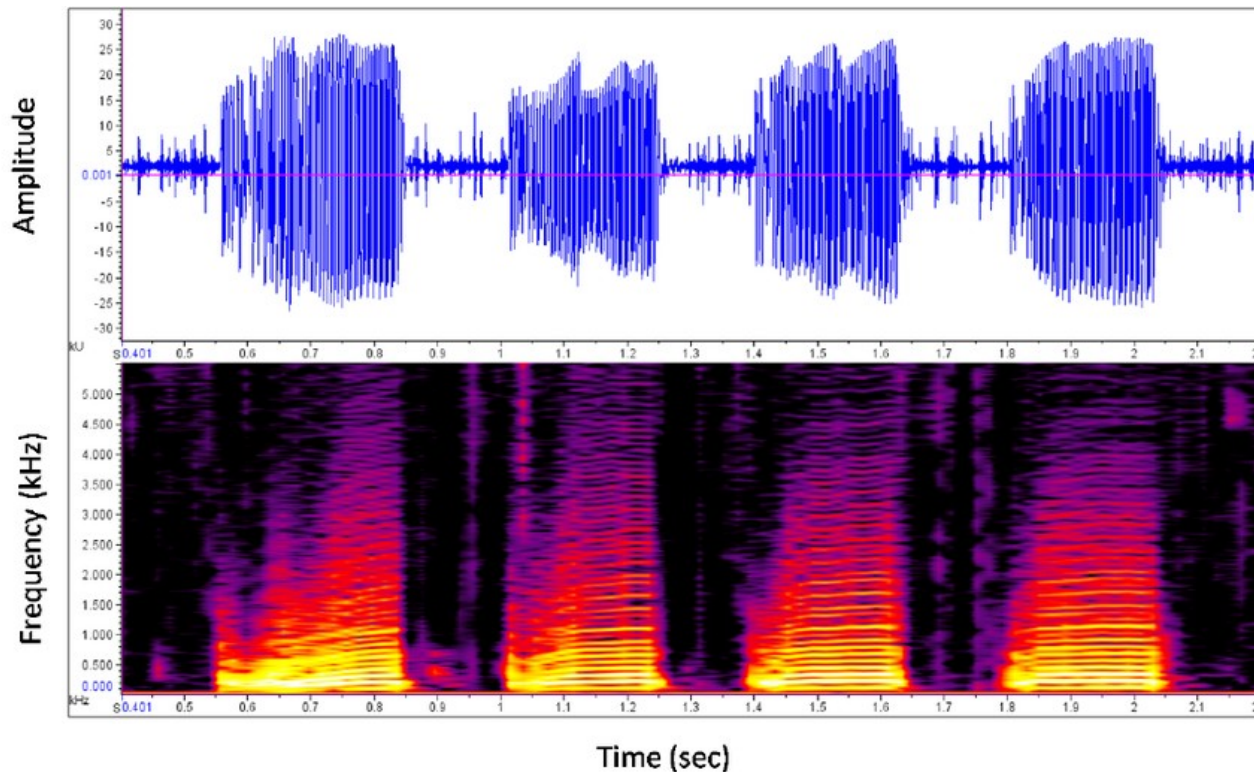
- Concatenative synthesis
- Parametric synthesis
 - Extract linguistic features
 - Duration model
 - Acoustic feature prediction model
 - Signal-processing-based vocoder

Данные

- На вход можно: текст, фонемы, спец. СИМВОЛЫ.
- Аудио.

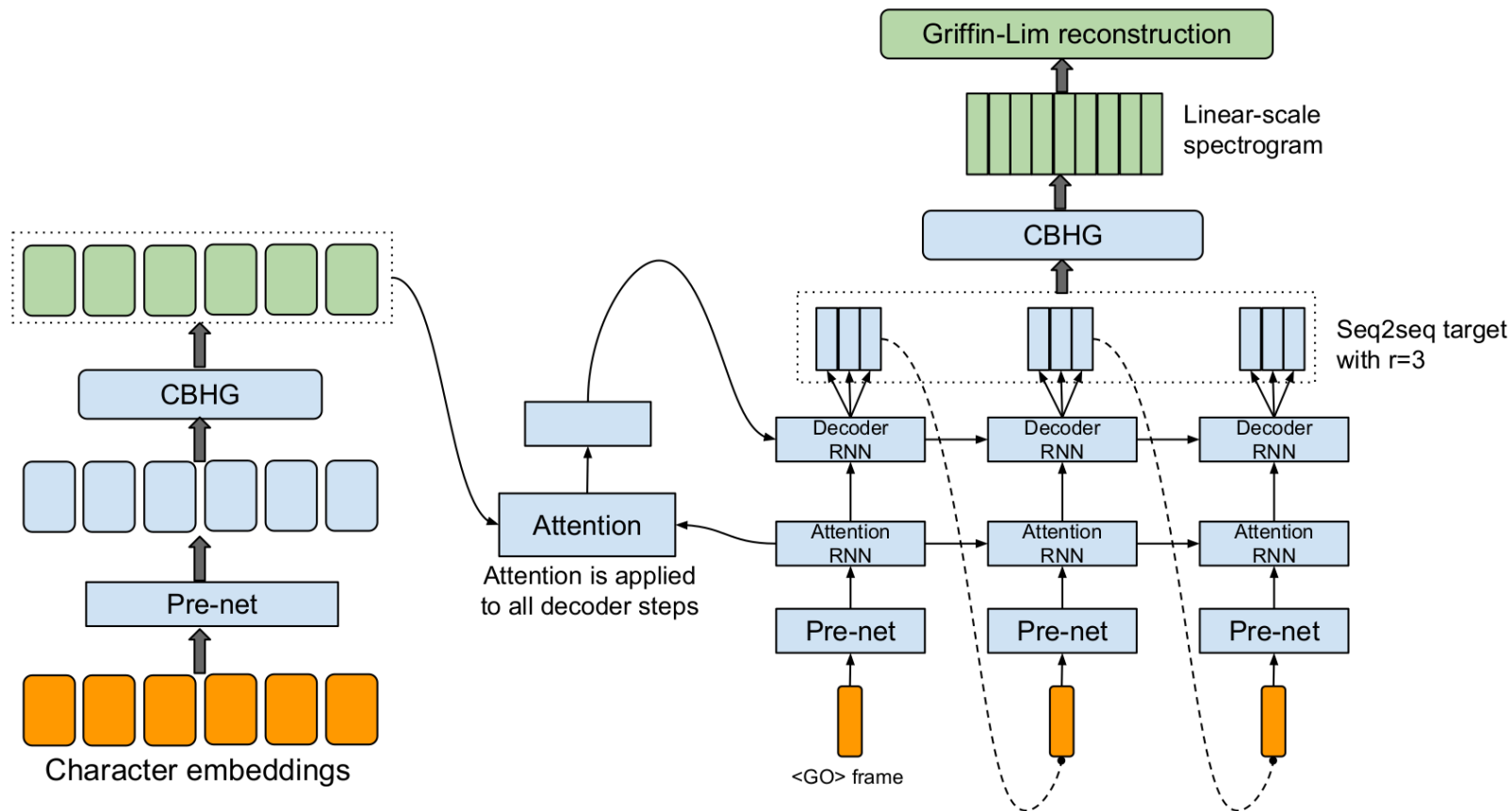
Spectrogram

- Спектрограмма – представление спектра частот звука.

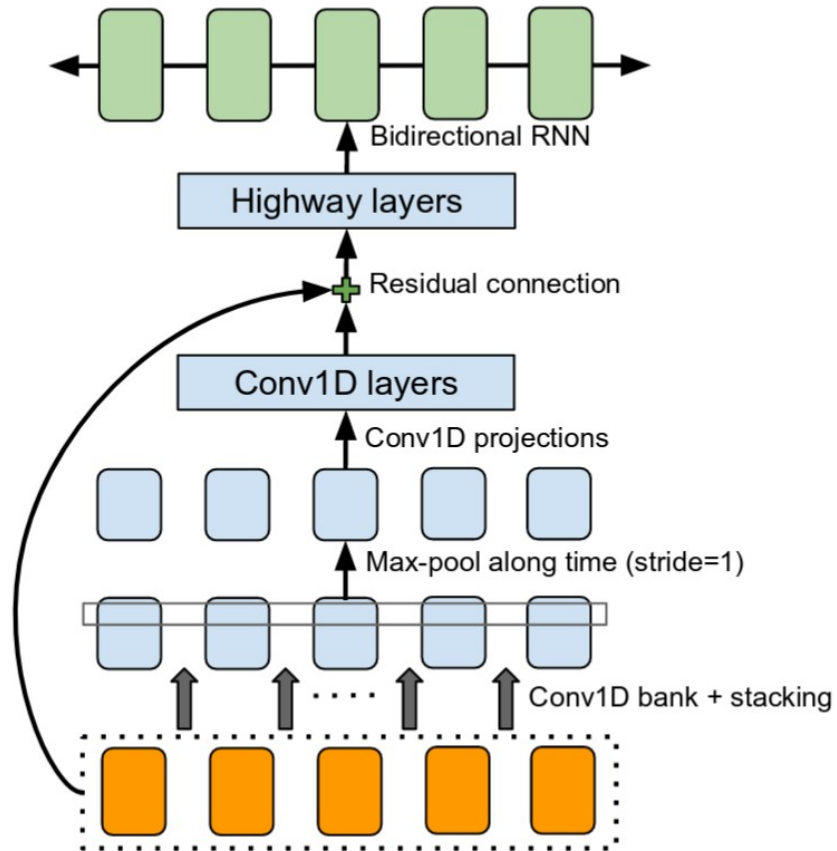


Модель Такотрона

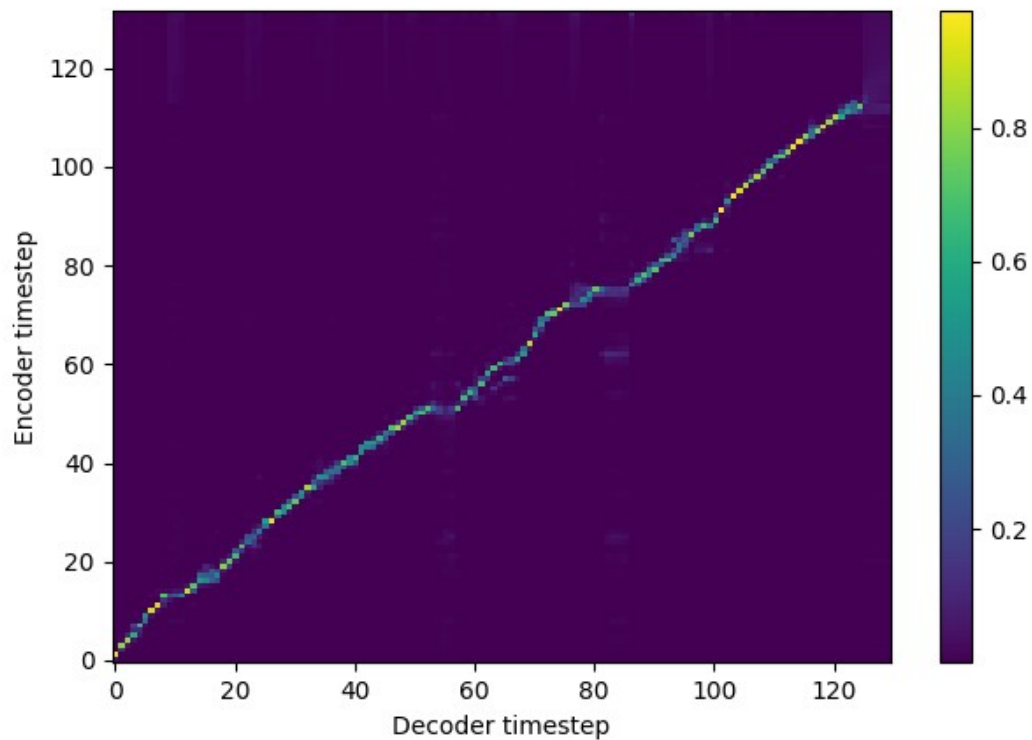
- Encoder
- Attention-based decoder
- Post-processing net



CBHG block



Attention



Tacotron, 2018-03-31 21:58, step=3400, loss=0.88621

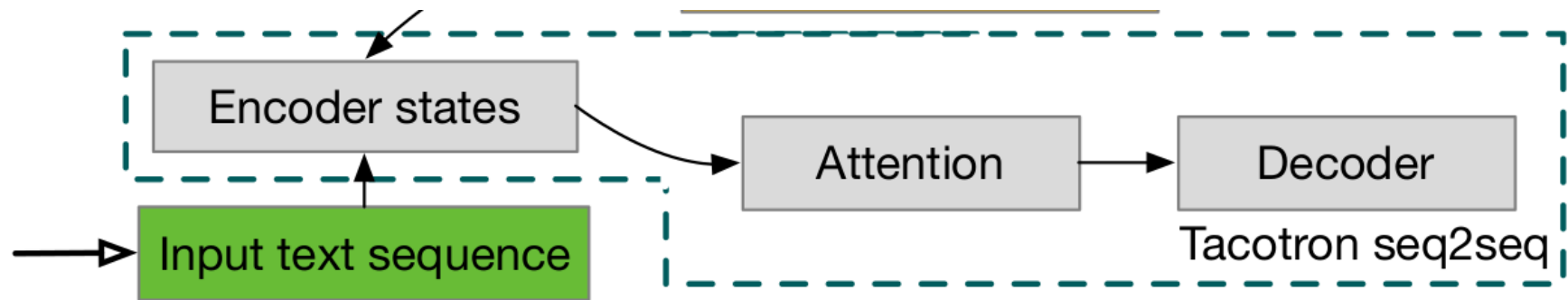
Вокодеры

- Griffin-Lim
- WaveNet
- Parallel WaveNet

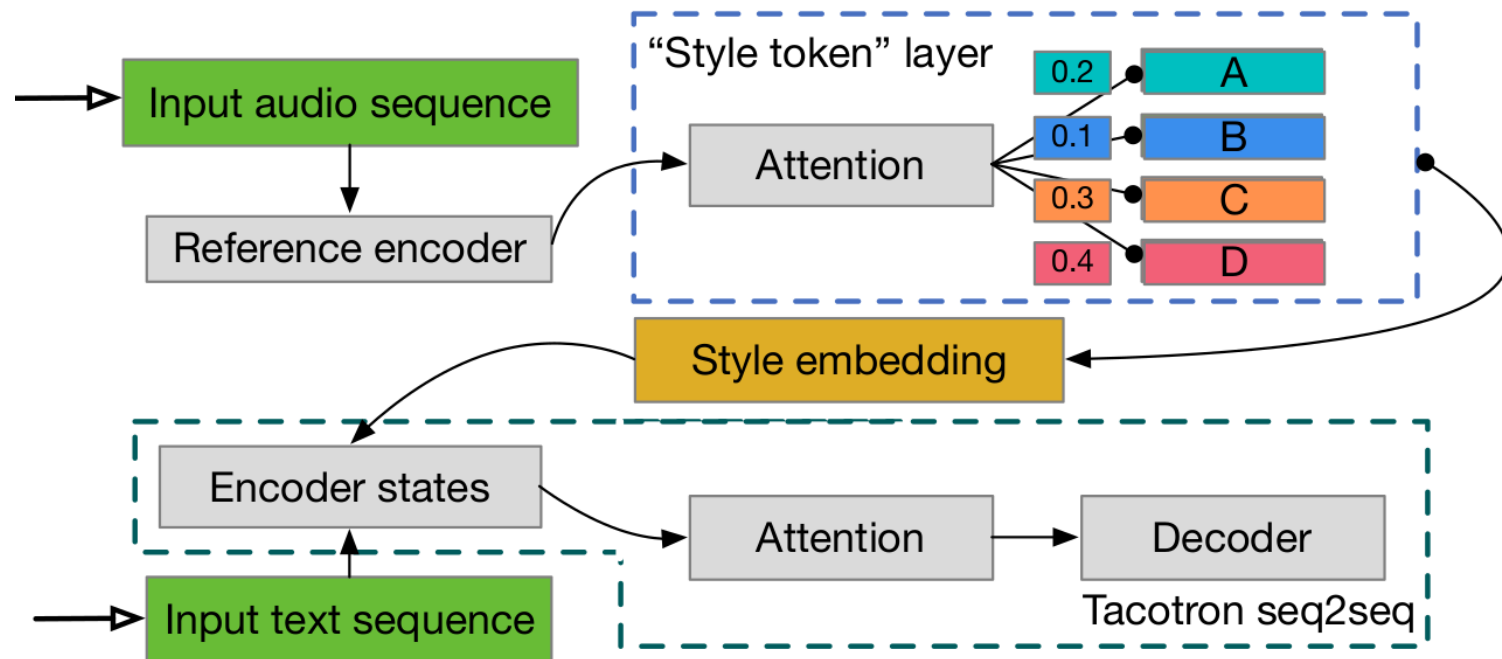
System	MOS
Parametric	3.492 ± 0.096
Tacotron (Griffin-Lim)	4.001 ± 0.087
Concatenative	4.166 ± 0.091
WaveNet (Linguistic)	4.341 ± 0.051
Ground truth	4.582 ± 0.053
Tacotron 2 (this paper)	4.526 ± 0.066

Prosody modeling

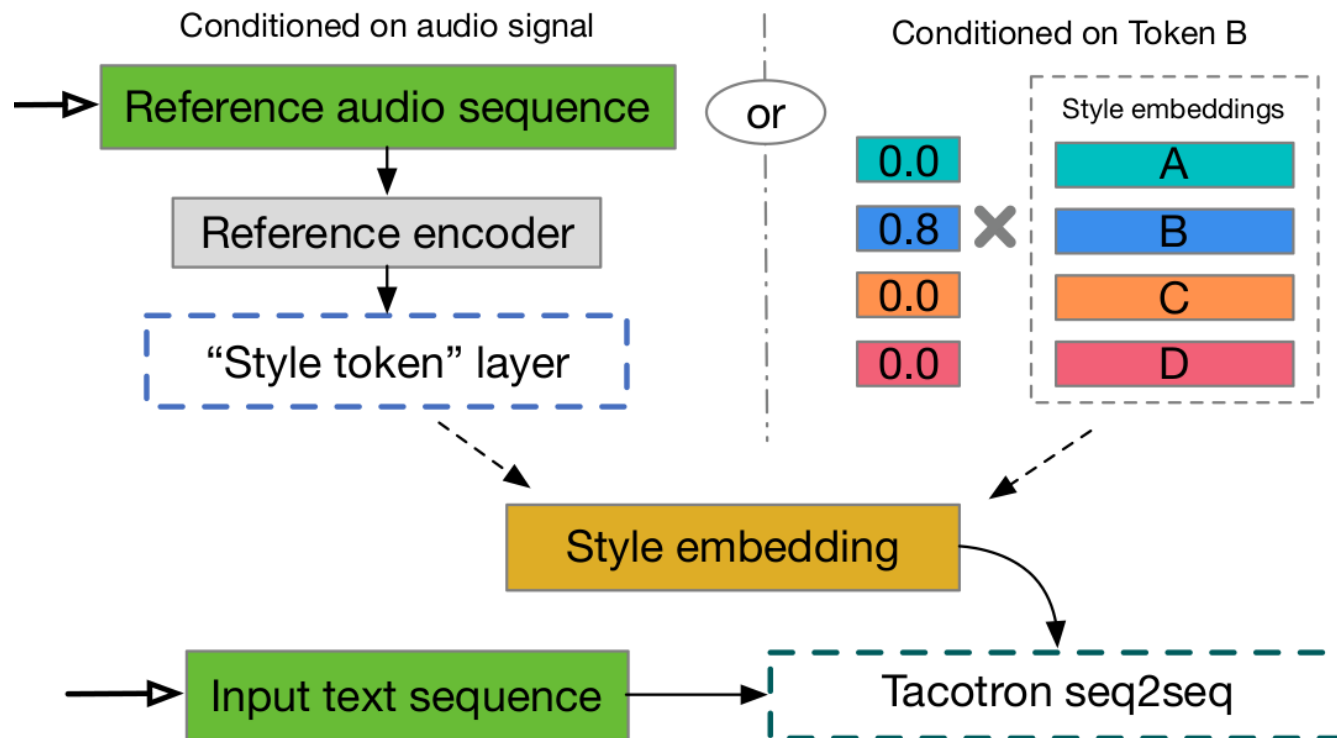
- В тексте нет информации о:
 - Интонациях
 - Особенности голоса
 - Эффектов окружения

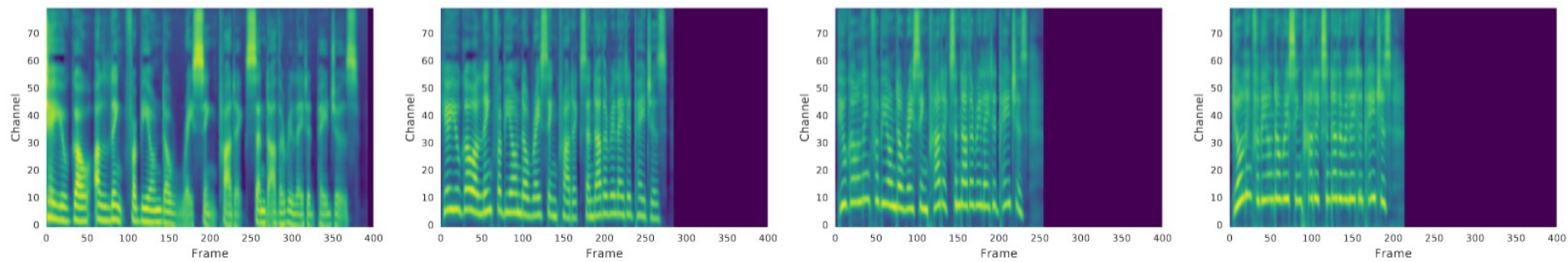


Training

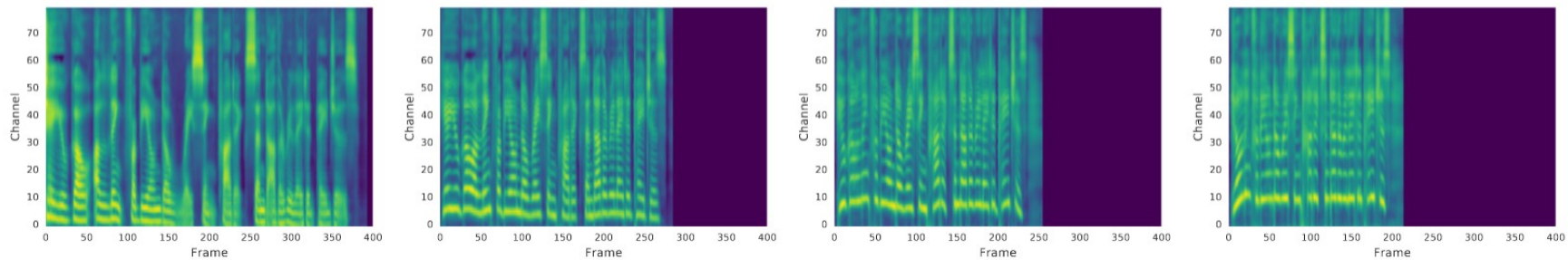


Inference

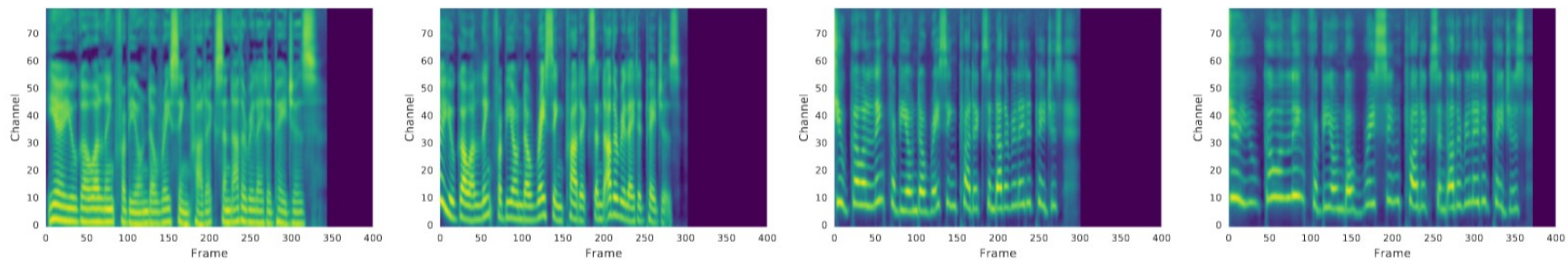




(a) Token A (speed)



(a) Token A (speed)



(b) Token B (animated)

NOISE %	BASELINE TACOTRON	GST
50%	2.819 ± 0.269	4.080 ± 0.075
75%	1.819 ± 0.227	3.993 ± 0.074
90%	1.609 ± 0.131	4.031 ± 0.082
95%	1.353 ± 0.090	3.997 ± 0.066

- <https://google.github.io/tacotron/> - статъи по такотрону.