

Современные методы управления просодией в задаче синтезе речи

Мылзенова Дарима, Just AI



План доклада

1. Что такое просодия и почему это важно для TTS
2. Контроль просодии через дополнительные признаки (с разметкой)
3. Нейросетевые архитектуры для контроля просодии

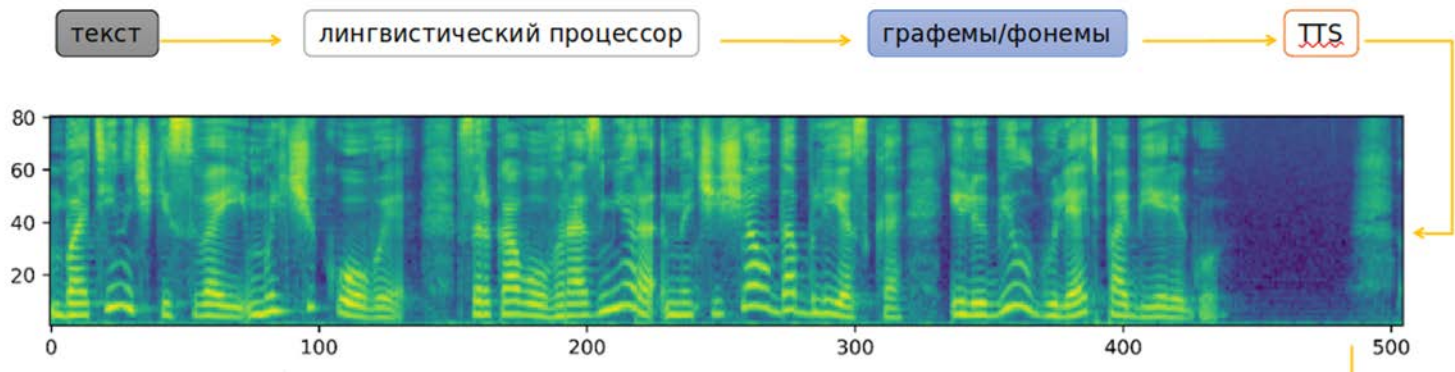


Задача синтеза речи





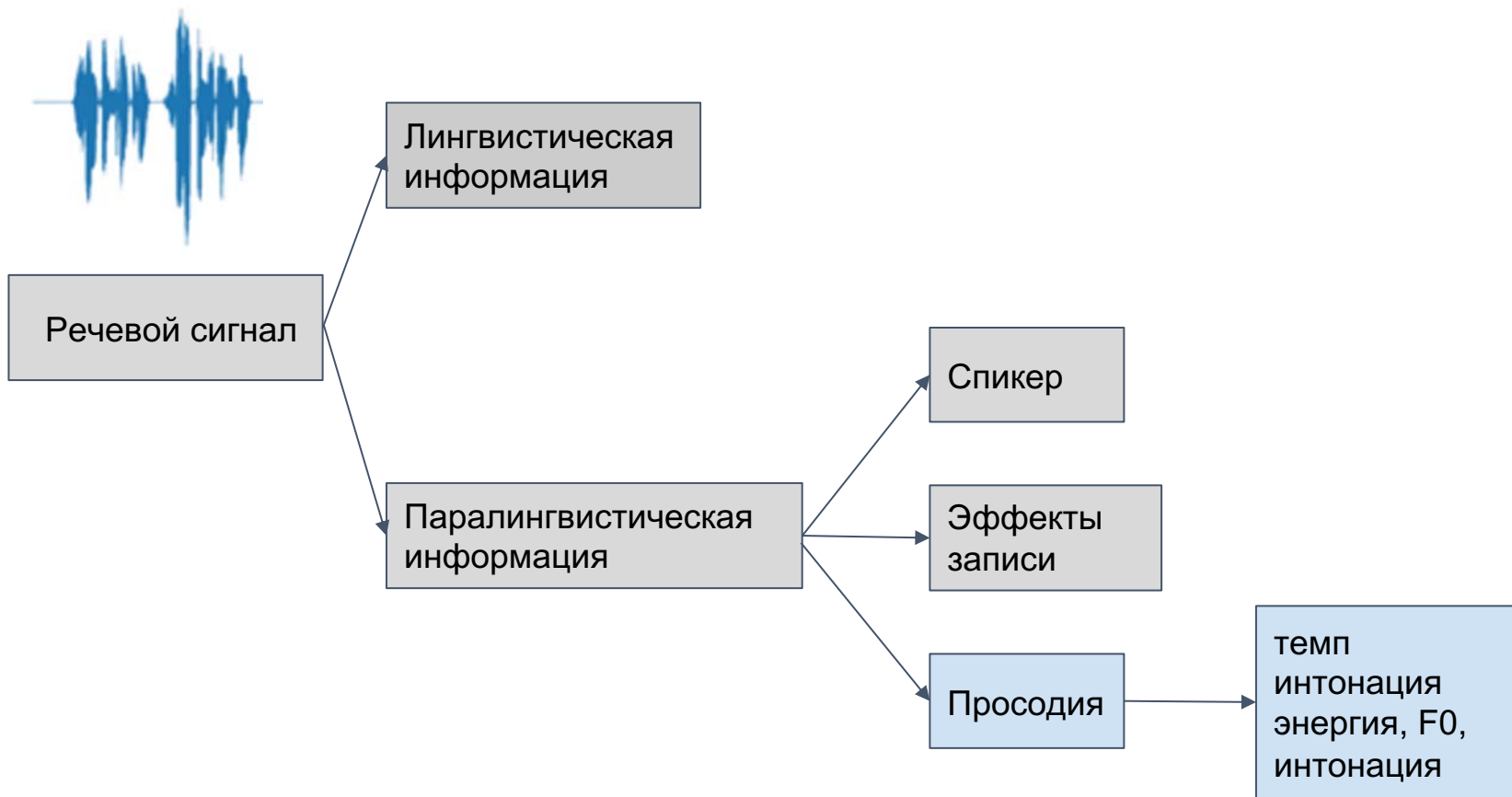
Задача синтеза речи



Проблема: задача one-to-many!



Задача синтеза речи



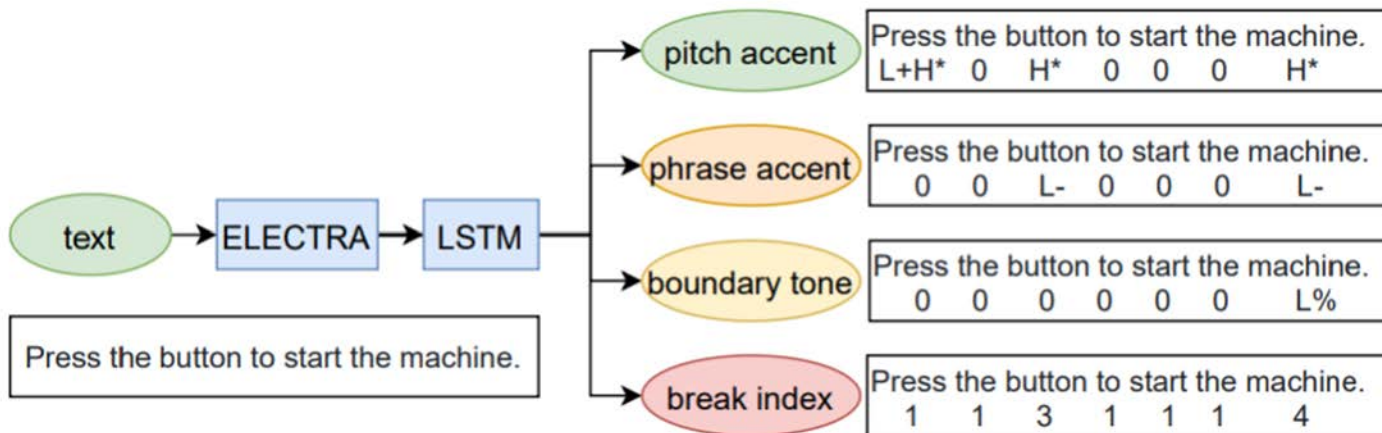


Как повлиять на просодию?

1. Разметить просодические фичи и подавать на вход модели
2. Без разметки - построить латентное представление просодии.

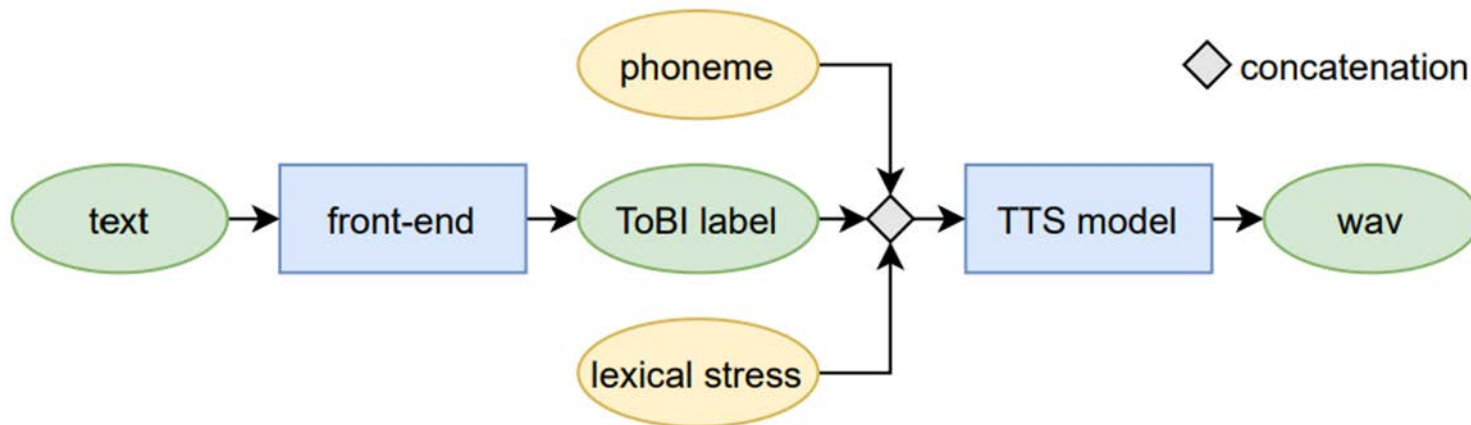


Supervised control: Tones and Breaks indices





Supervised control: Tones and Breaks indices





Supervised control: ИНТОНАЦИОННЫЕ КОНТУРЫ

ИК-1: Студенты вернулись.

ИК-2: Когда студенты вернулись?

ИК-3: Студенты вернулись?

ИК-4: (В среду я не могу прийти.) – А не в среду?

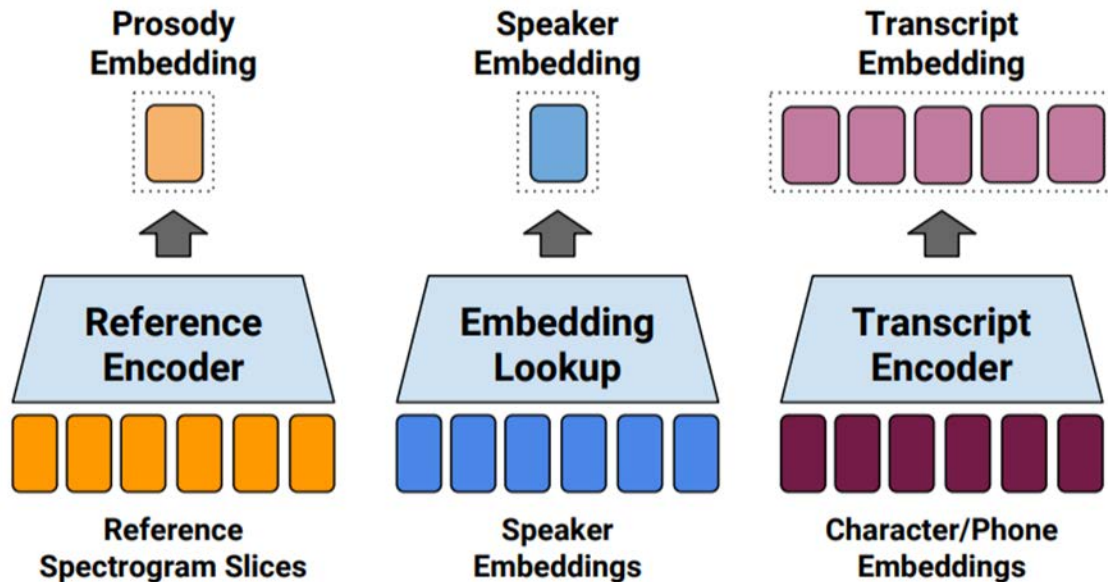
ИК-5: Сколько хорошего сделал!

ИК-6: А что мы сегодня узнали! (торжествующе)

ИК-7: А что он умеет! (разочарованно).



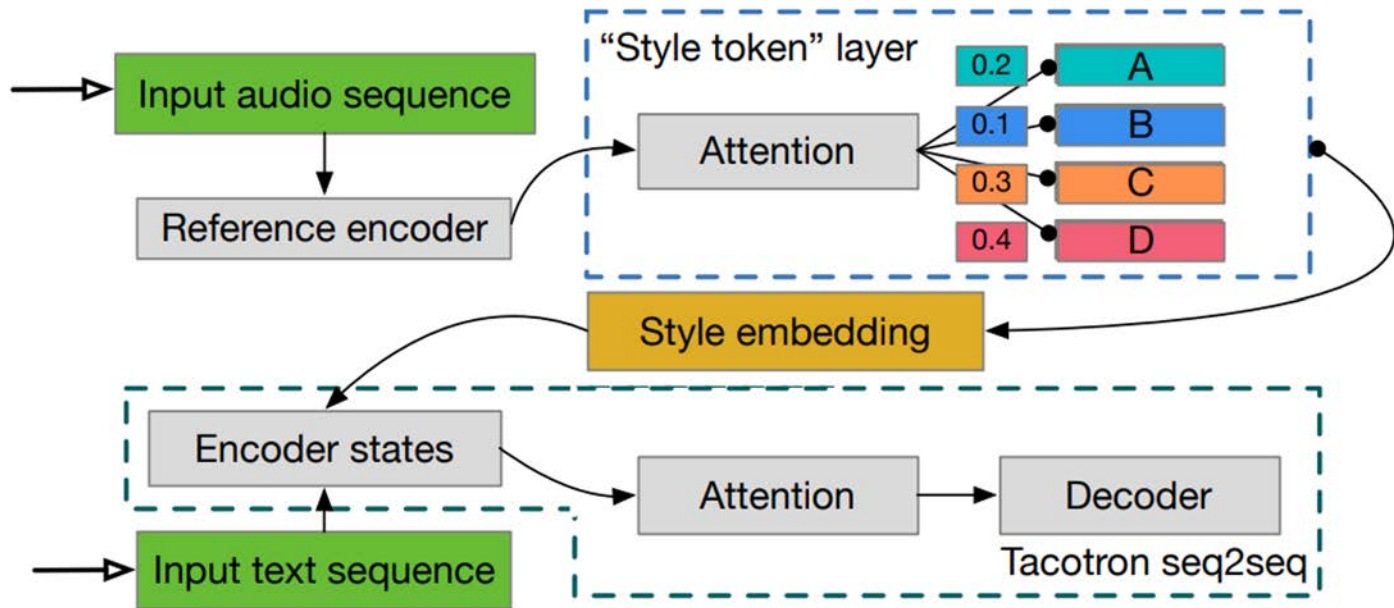
Unsupervised control: Tacotron 2



<https://arxiv.org/pdf/1803.09047.pdf>

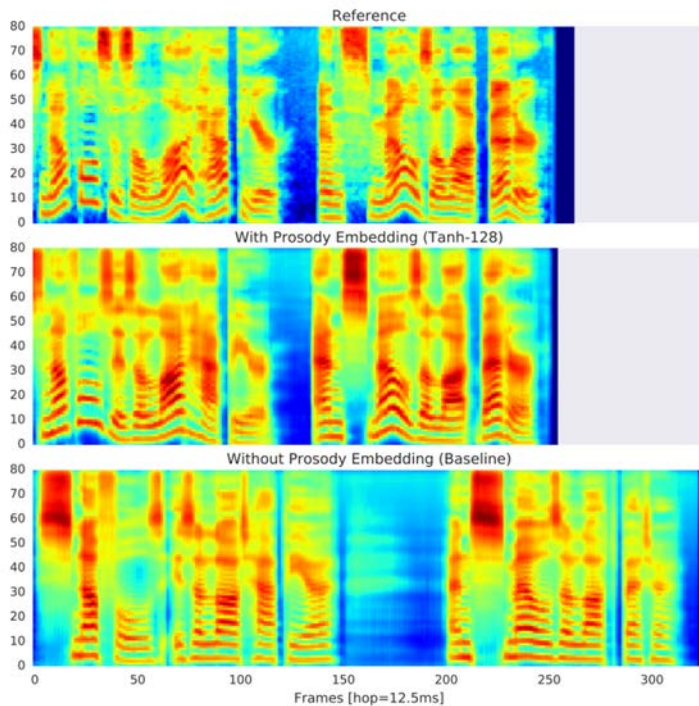


Unsupervised control with GST





Unsupervised control: Tacotron 2



референс

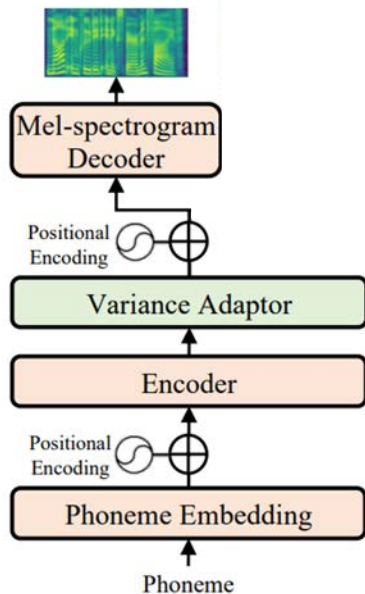
синтез с референсом

синтез безлайн модели

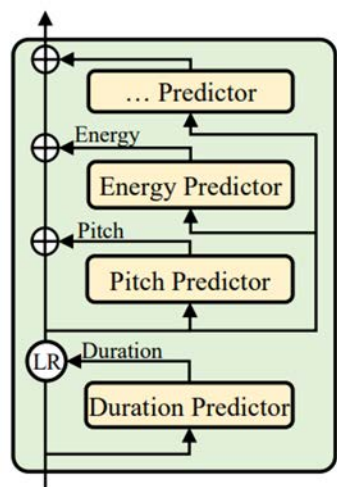
<https://arxiv.org/pdf/1803.09047.pdf>



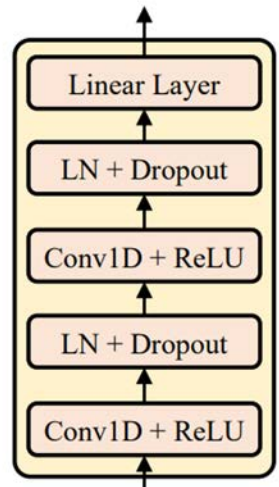
Unsupervised control: FastSpeech 2



(a) FastSpeech 2



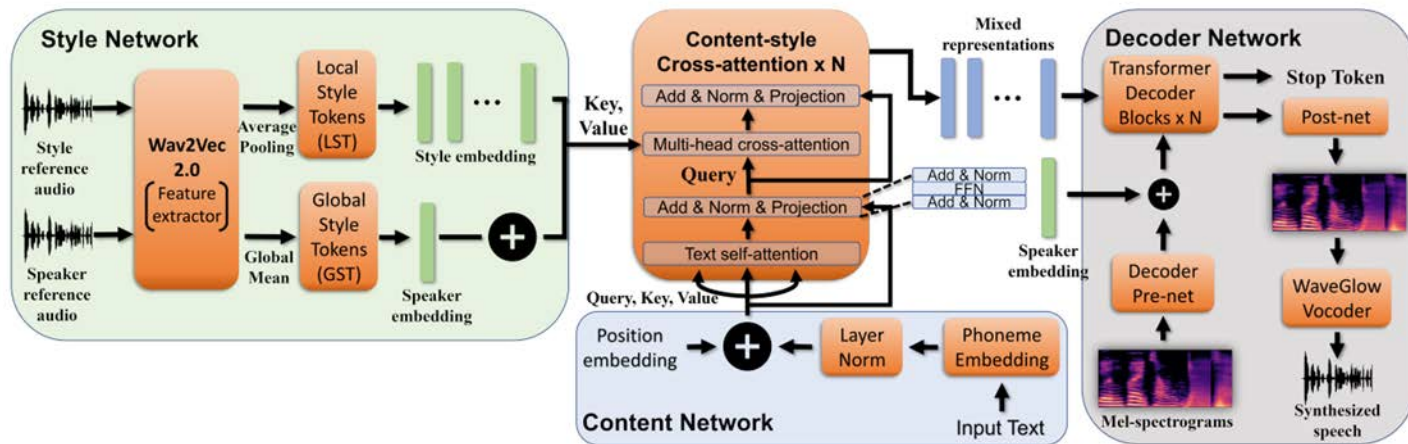
(b) Variance adaptor



(c)
Duration/pitch/energy
predictor



Fine-grained control with LST





Flowtron

- Normalizing Flows
 $p(x) \rightarrow p(z)$

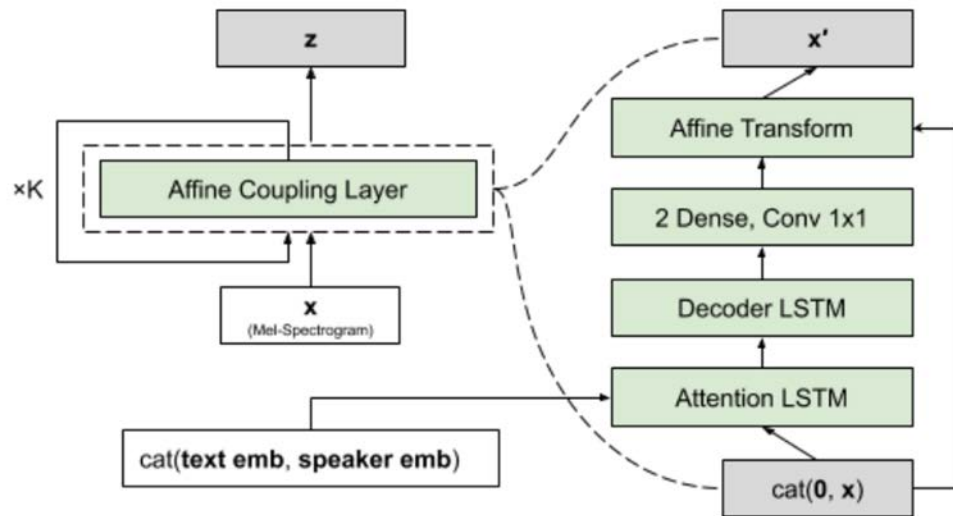
$$z = f_0 \circ f_1 \circ \dots \circ f_k(x)$$

$$z \sim N(z; 0, I)$$

x - фрейм спектра,
 f_i - обратимые

- Во время инференса сэмплируем случайный z , обращаем f_i и получаем x
 $x = f_{-1}$

$$x = f_k^{-1} \circ f_{k-1}^{-1} \circ \dots \circ f_0^{-1}(z)$$



<https://arxiv.org/pdf/2005.05957.pdf>



Flowtron

Сэмплируя с большей дисперсией можно добиться более экспрессивной и живой речи.

Обуславливая параметры распределения можно сделать качественный перенос просодии.

<https://developer.nvidia.com/blog/training-your-own-voice-font-using-flowtron/>



Summary

- Управление просодией с разметкой
- Предсказание латентного вектора просодии и стиля через Reference Encoder и GST
- Предсказание пичча, энергии и длительности с FastSpeech2
- Более детальное предсказание стиля - LST
- Экспрессивный синтез с FlowTron