

LISTEN. ATTENTION. SPELL

William Chan, Navdeep Jaitly, Quoc V. Le, Oriol Vinyals
Google, 2015

Подготовила Гаврилова А. А.

ПЛАН

- ***Введение в ASR***

- Препроцессинг звуковой дорожки
- Акустическая модель
- Языковая модель
- Пунктуационная модель
- Метрики

- ***Listen, Attend and Spell***

- Архитектура
- Функция потерь
- Beam-Search decoding
- Результаты

АУДИОФОРМАТЫ

- **Аналоговый сигнал**

непрерывное множество
значений



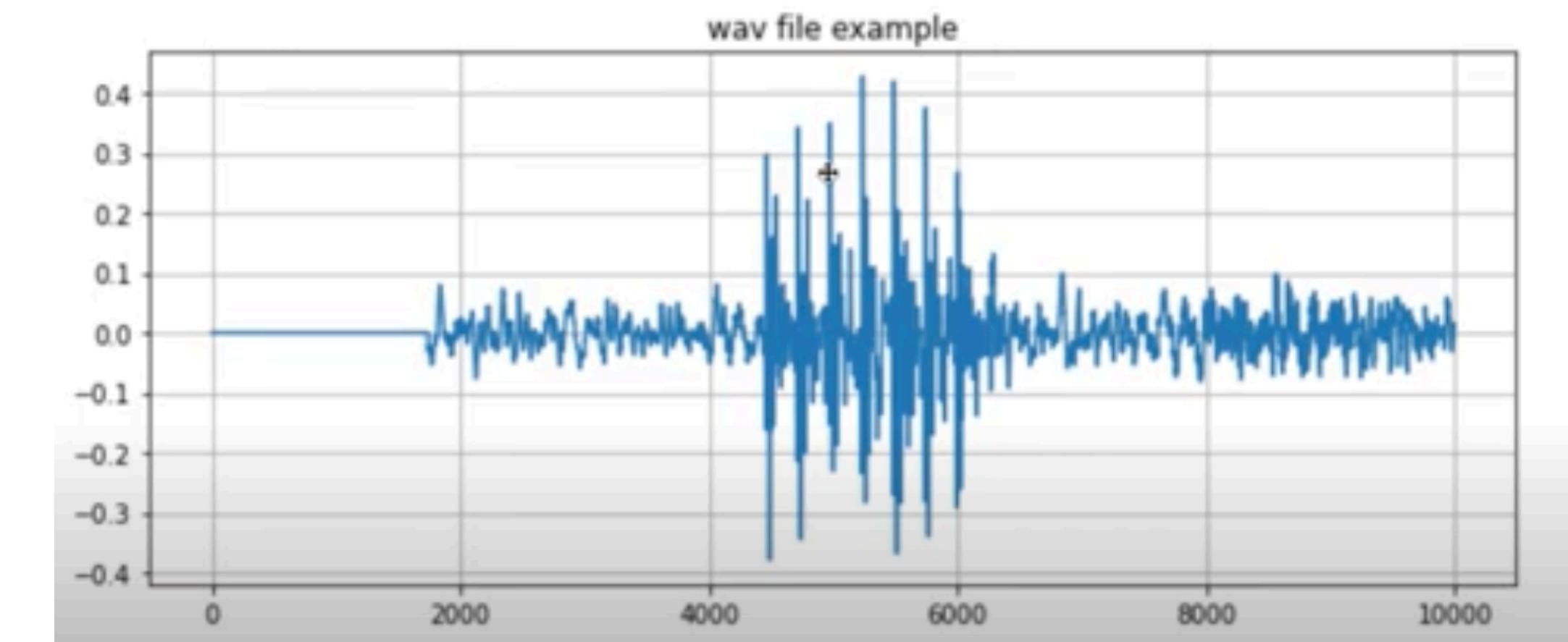
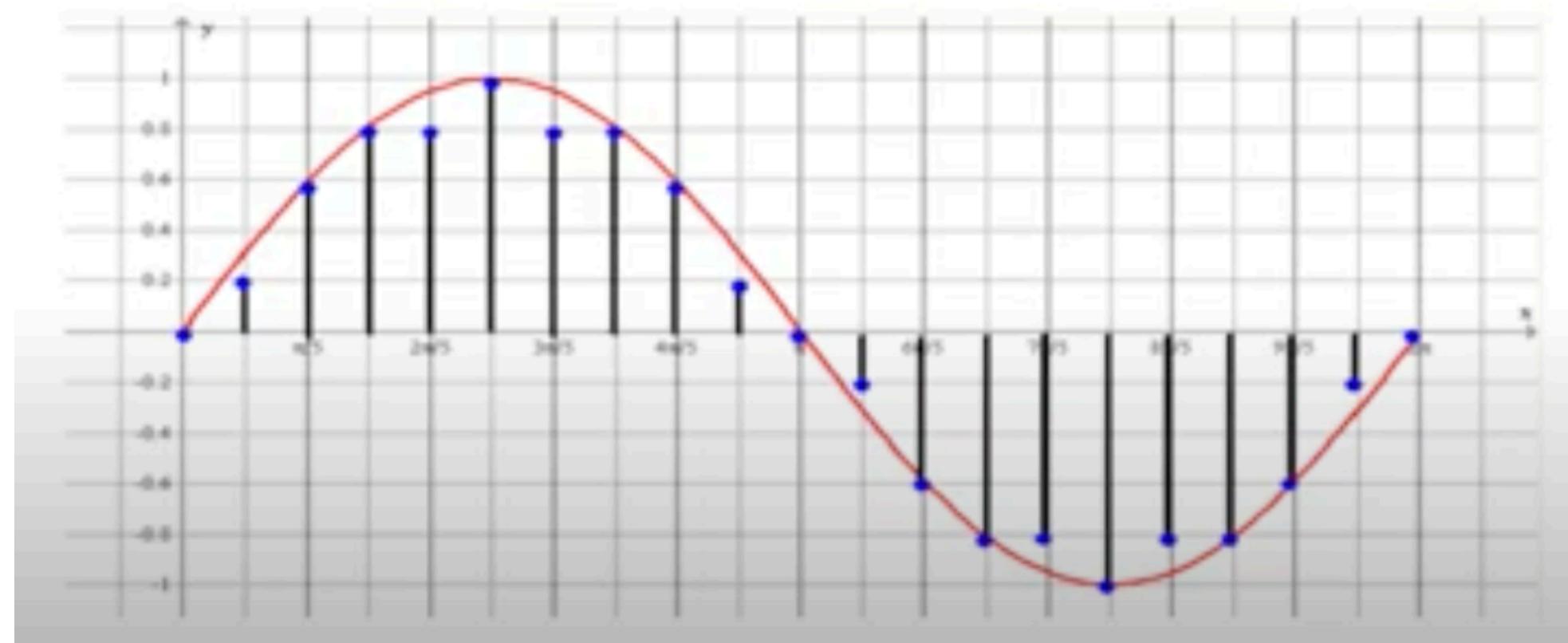
- **Цифровой сигнал**

последовательность
дискретных значений



КАК ХРАНИТСЯ АУДИО?

- sample_rate - число отсчетов в секунду (16000, 22050)
- Librosa - библиотека для работы со звуком



КАК ПЕРЕВЕСТИ АУДИО В ТЕНЗОР?

- Будем работать как с картинками (с ними умеем работать)

$$G(f) = \int_{-\infty}^{\infty} g(t) e^{-2\pi i f t} dt.$$

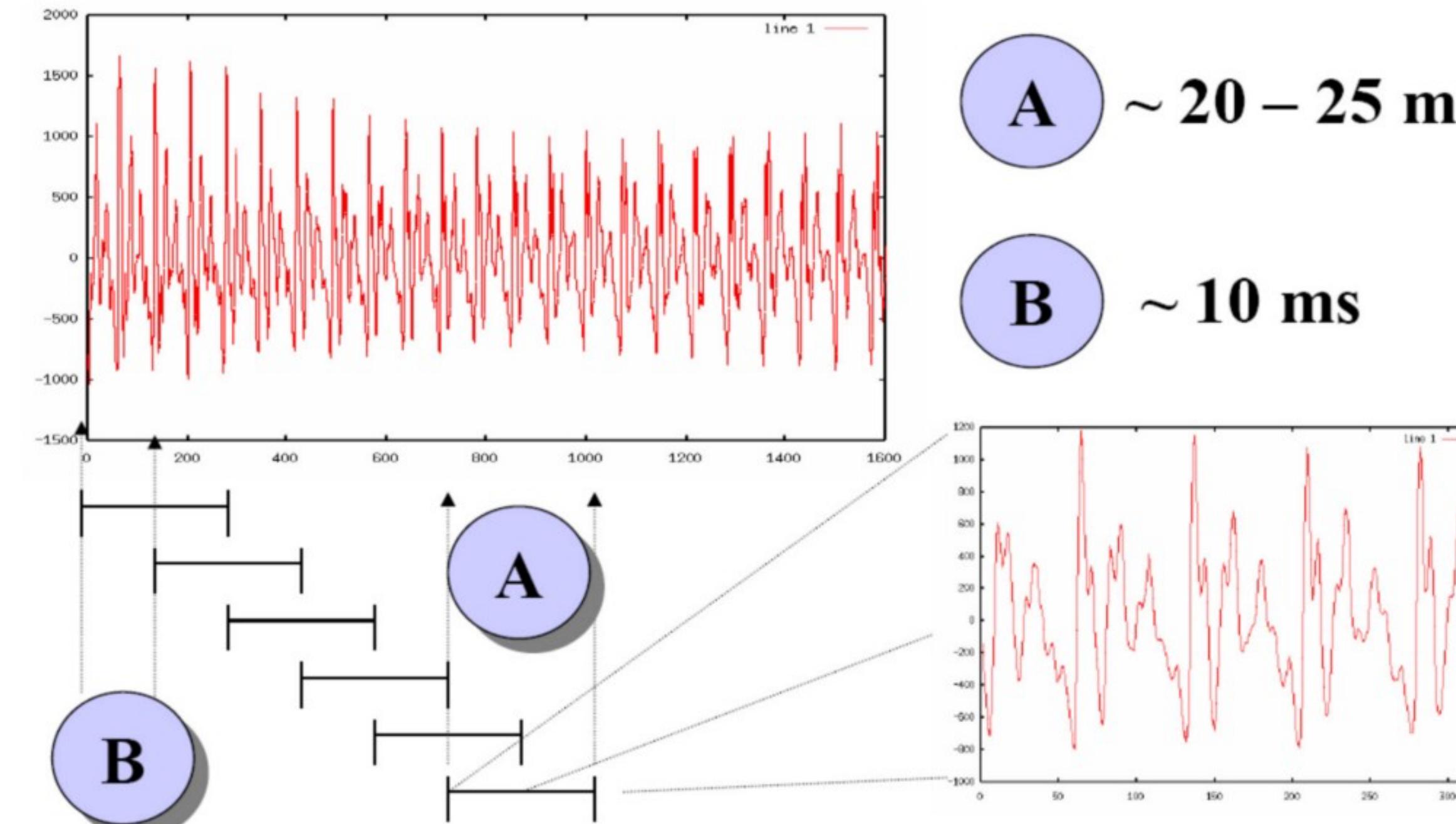
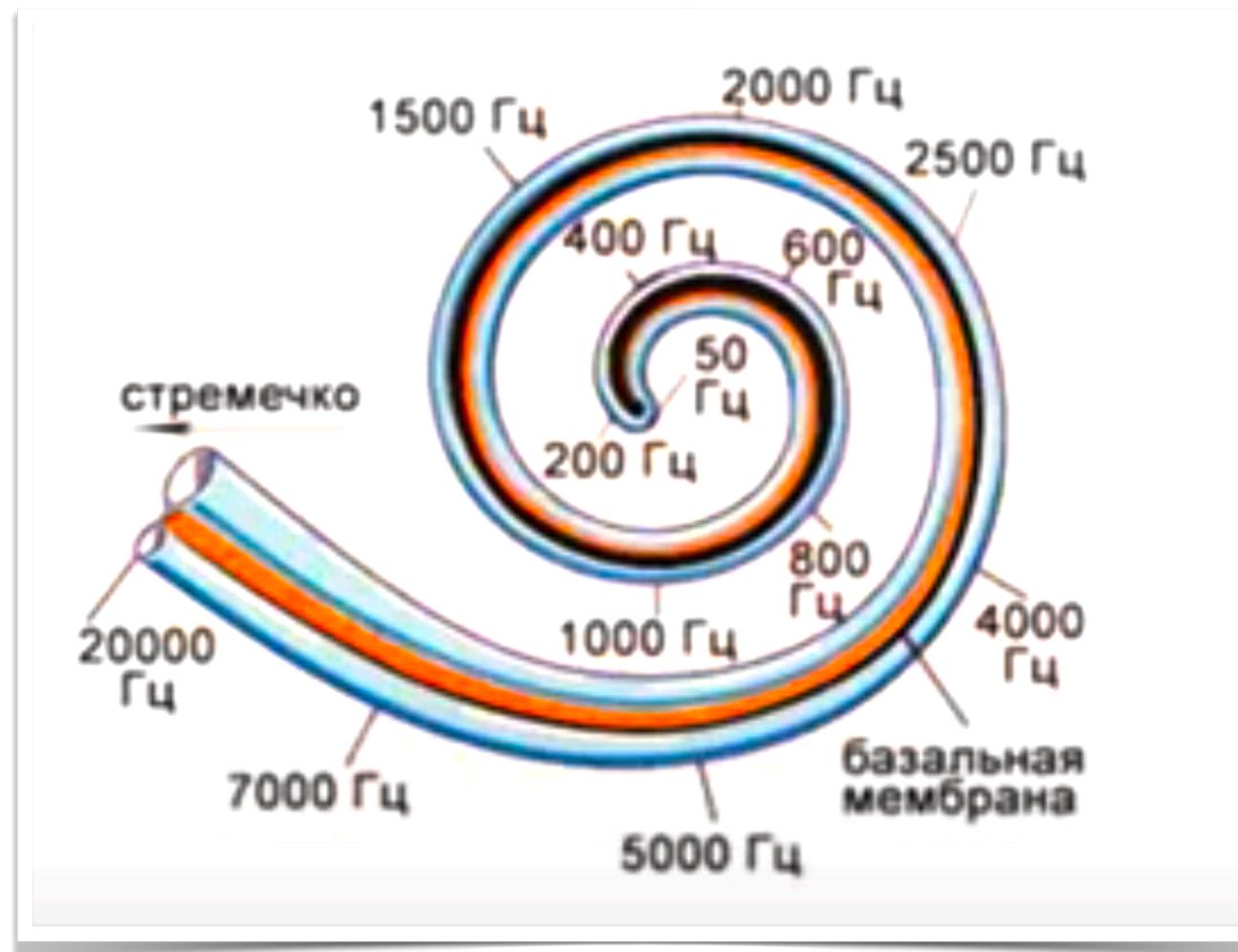
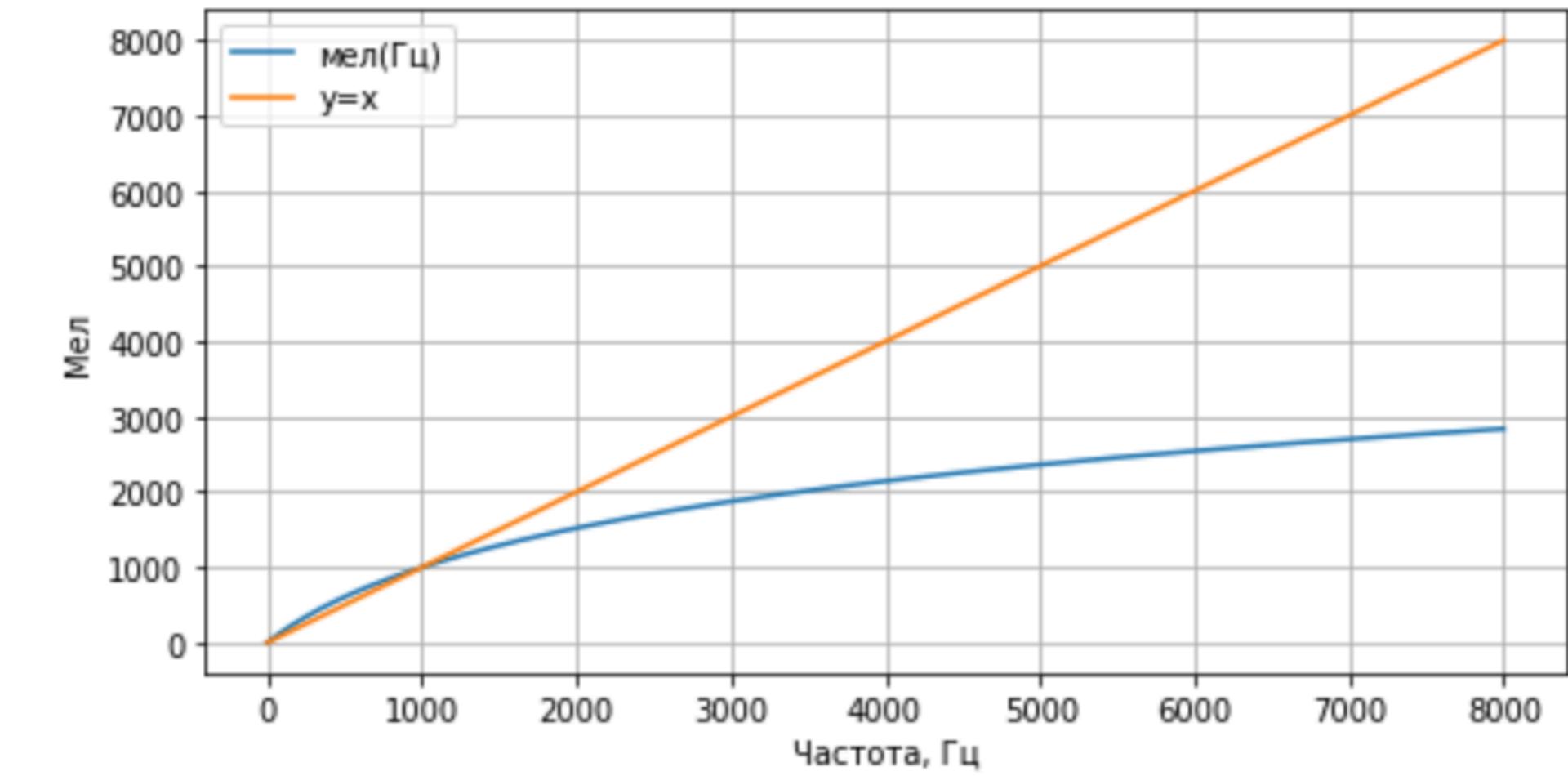
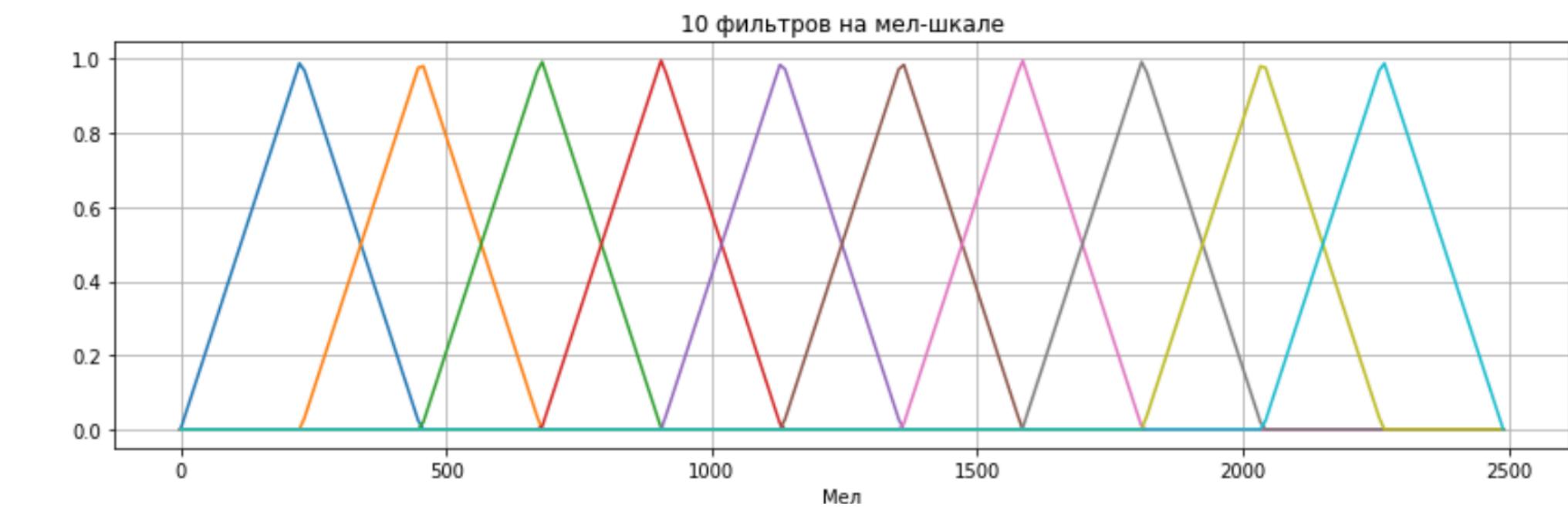
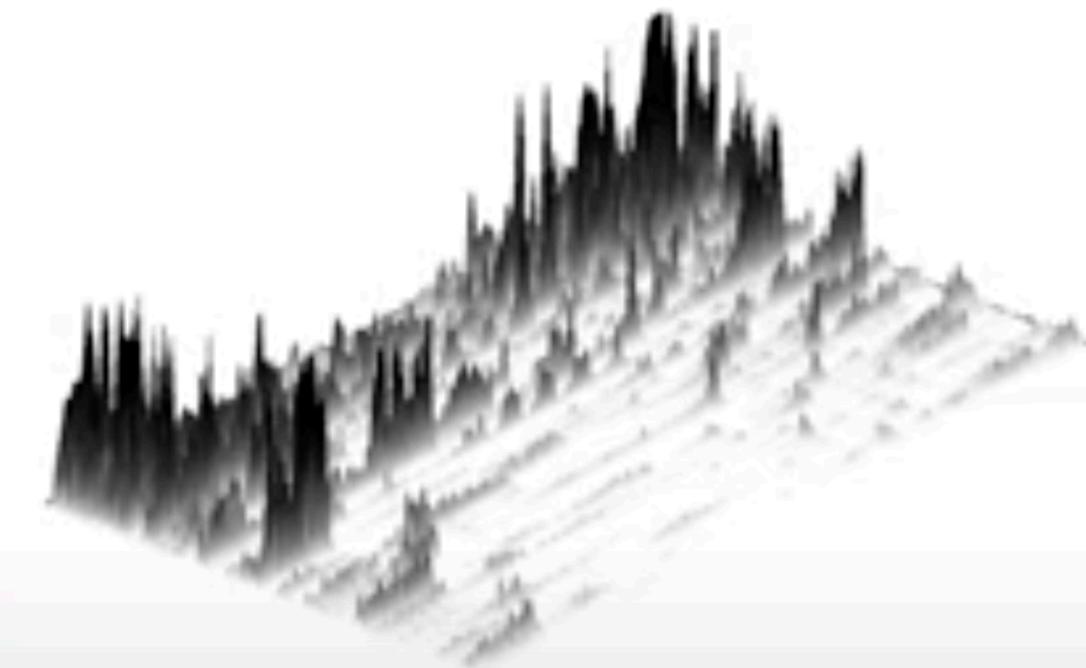
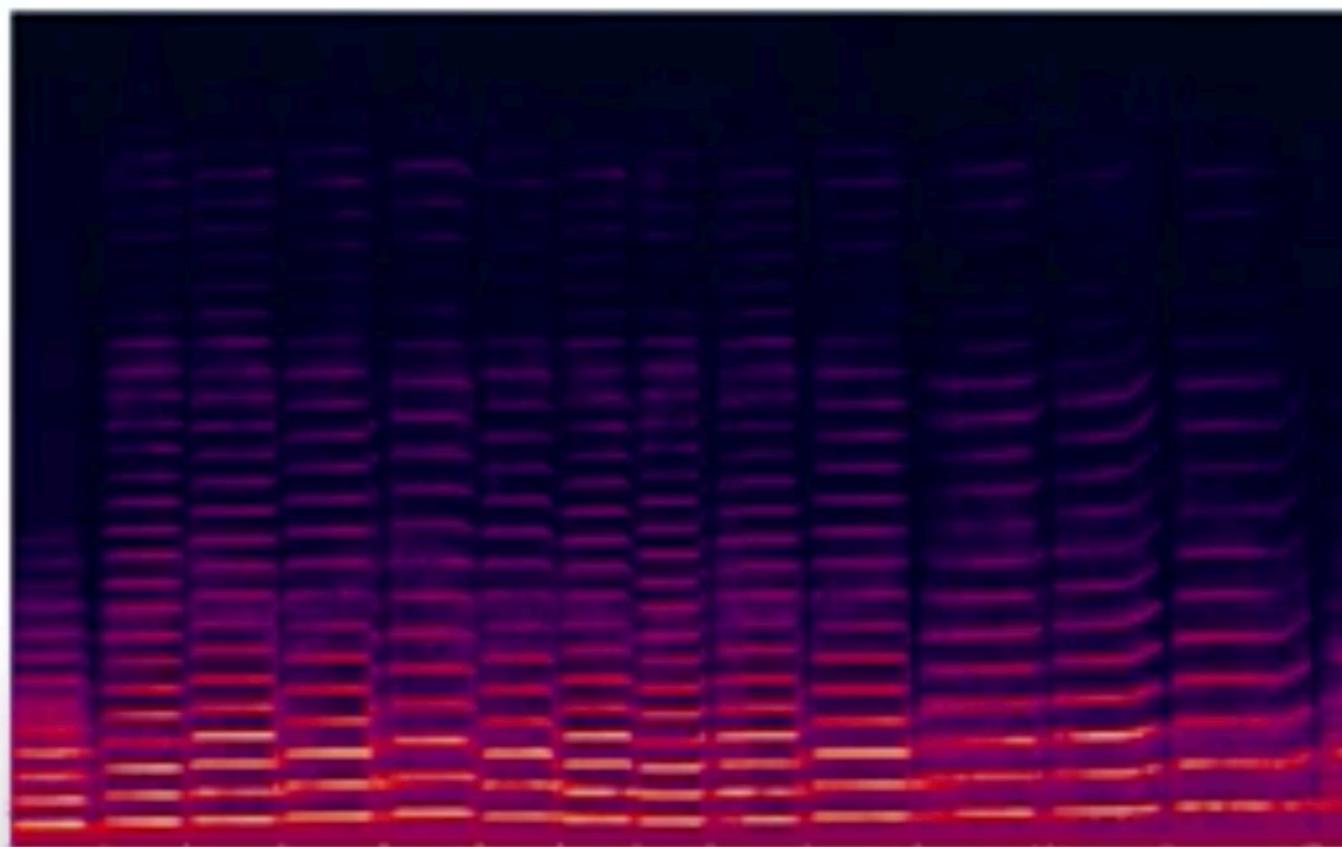
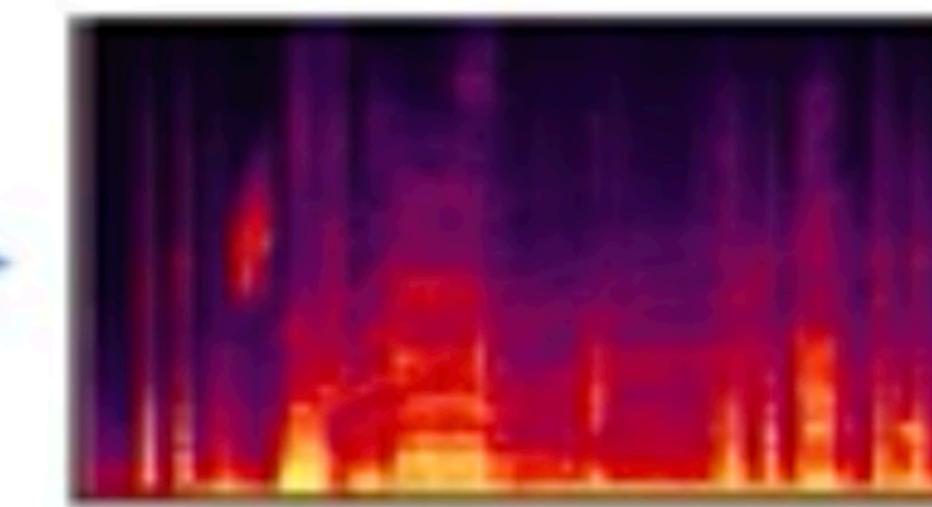


Image from Bryan Pellom

МЕЛ-СПЕКТРОГРАММА

$$mel = 1127.01048 \ln\left(1 + \frac{freq}{700}\right)$$





приве янадя

акустическая
модель



привет я наядя

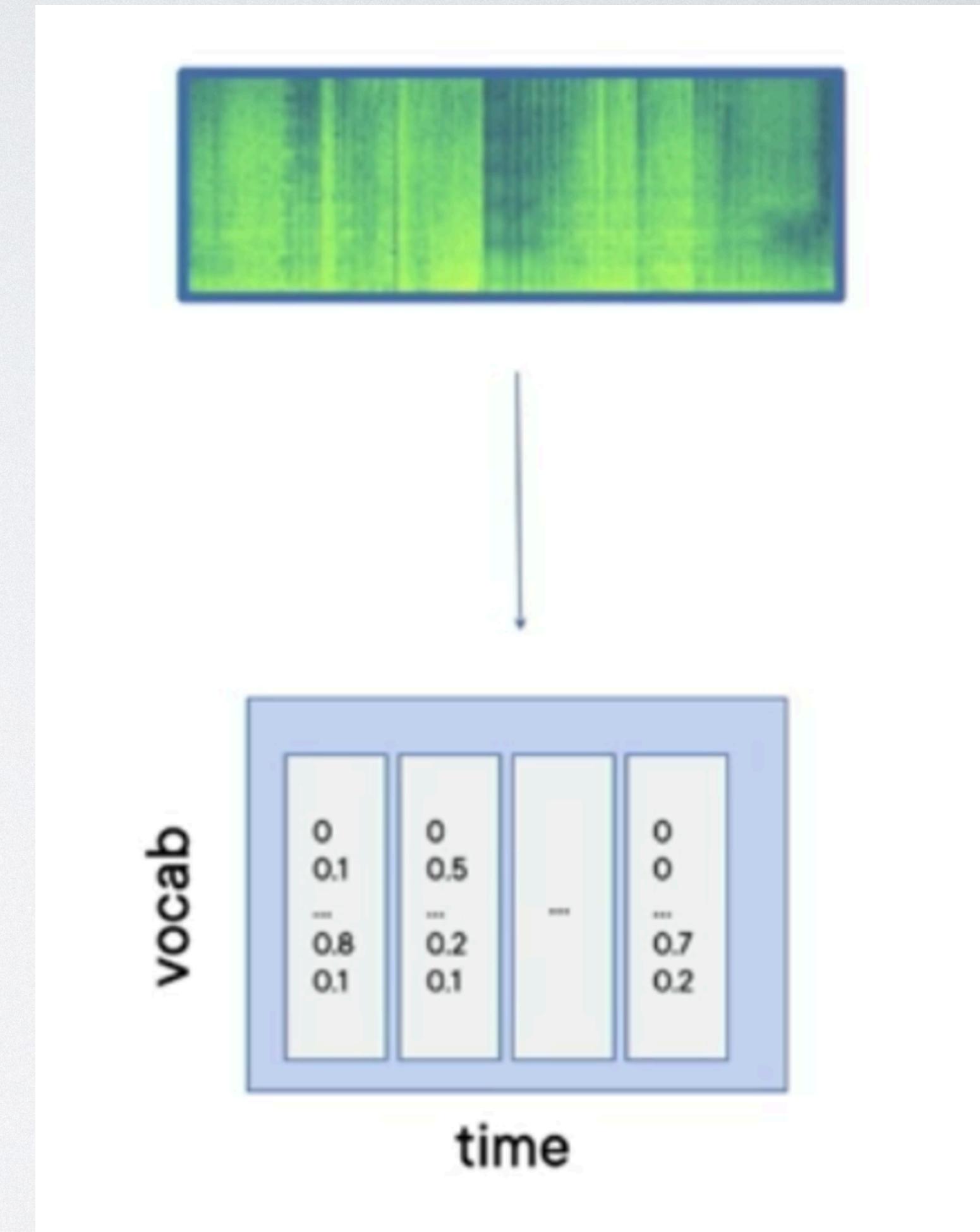
языковая
модель



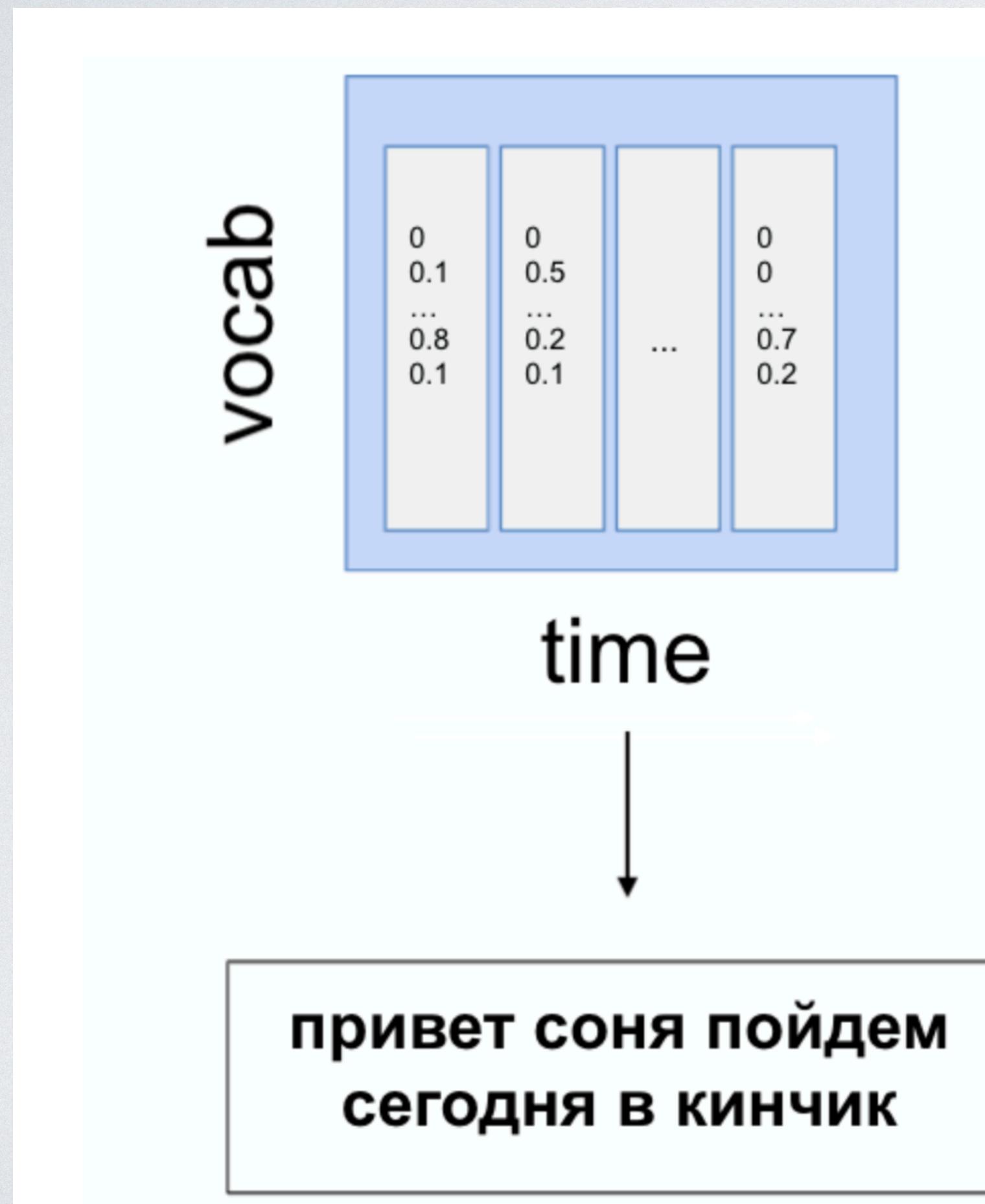
profit

Акустическая модель

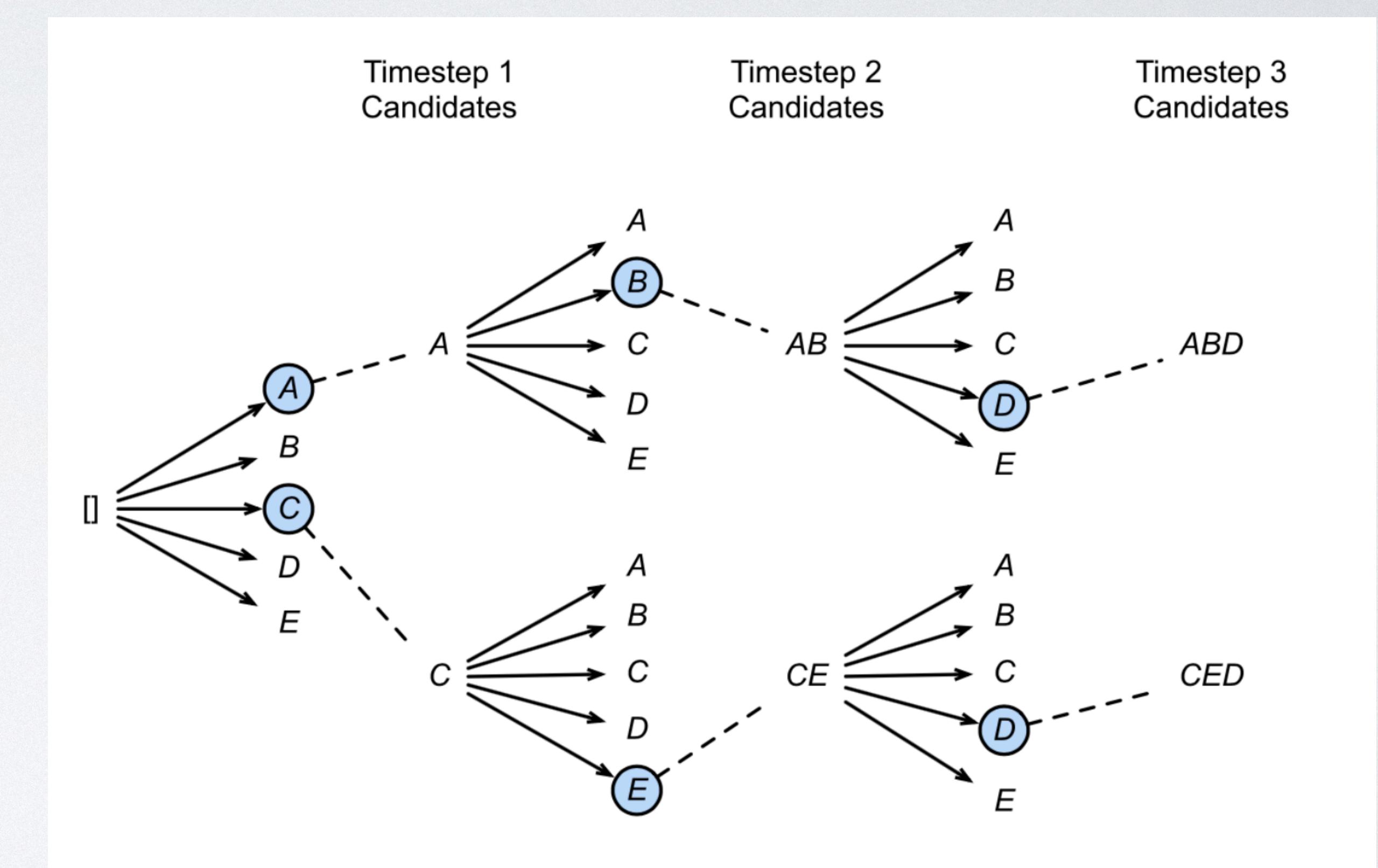
- В качестве простой акустической модели в задаче распознавания речи использовались **марковские модели**, затем нейросети
- На выходе матрица распределения по времени вероятностей каждой фонемы



языковая модель



Beam-search



ПУНКТУАЦИОННАЯ МОДЕЛЬ

привет соня пойдем
сегодня в кинчик



Привет, Соня. Пойдем
сегодня в кинчик?

МЕТРИКИ

WER

True: quick **brown** fox jumped over **a** lazy dog
Pred: quick **brow** an fox jumped over | lazy dog

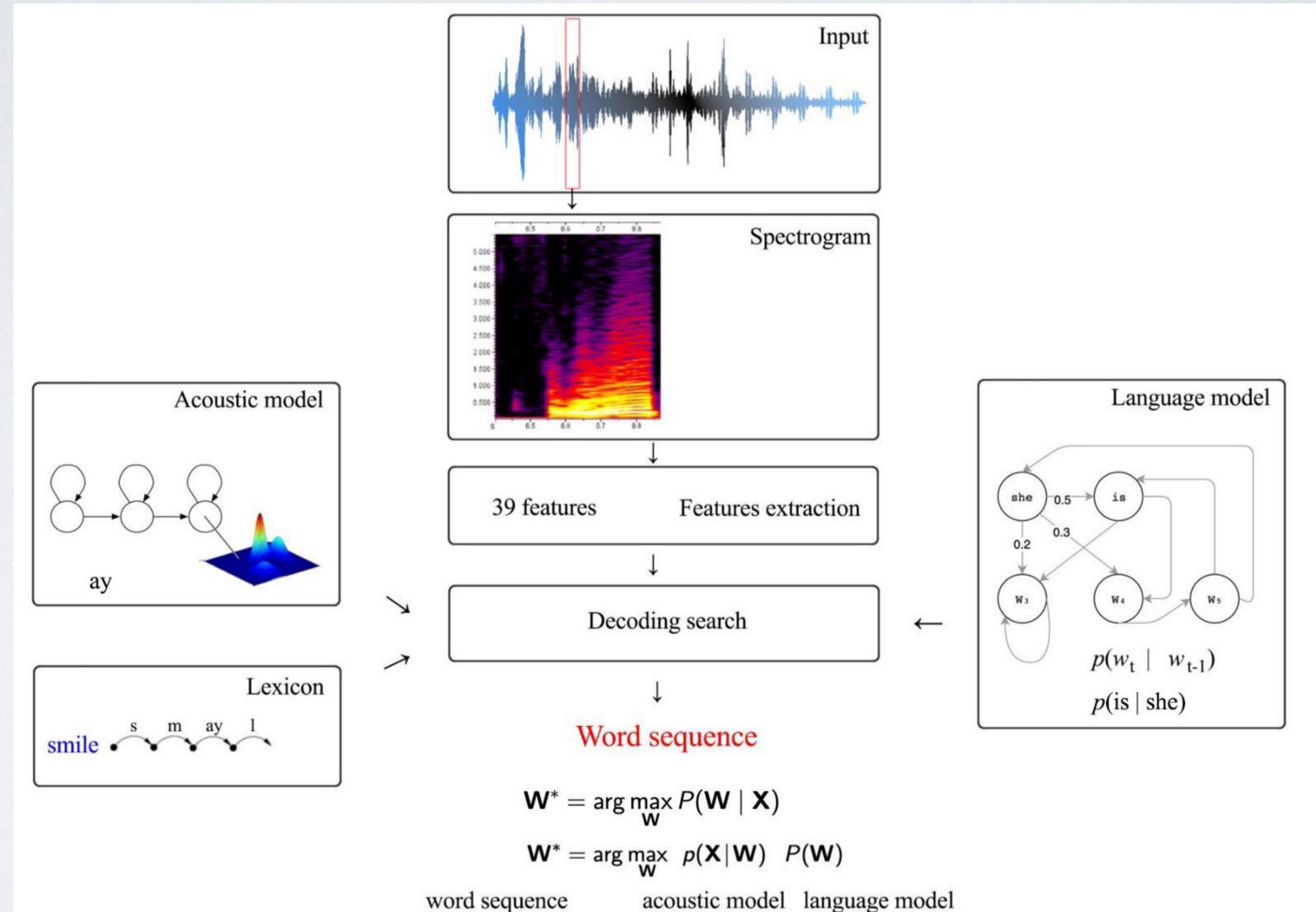
CER

То же, но с символами

$$WER = \frac{S + D + I}{N} = \frac{S + D + I}{S + D + C}$$

- S – число замен
 - D – число удалений
 - I – число вставок
 - C – число корректных слов
- N – всего слов

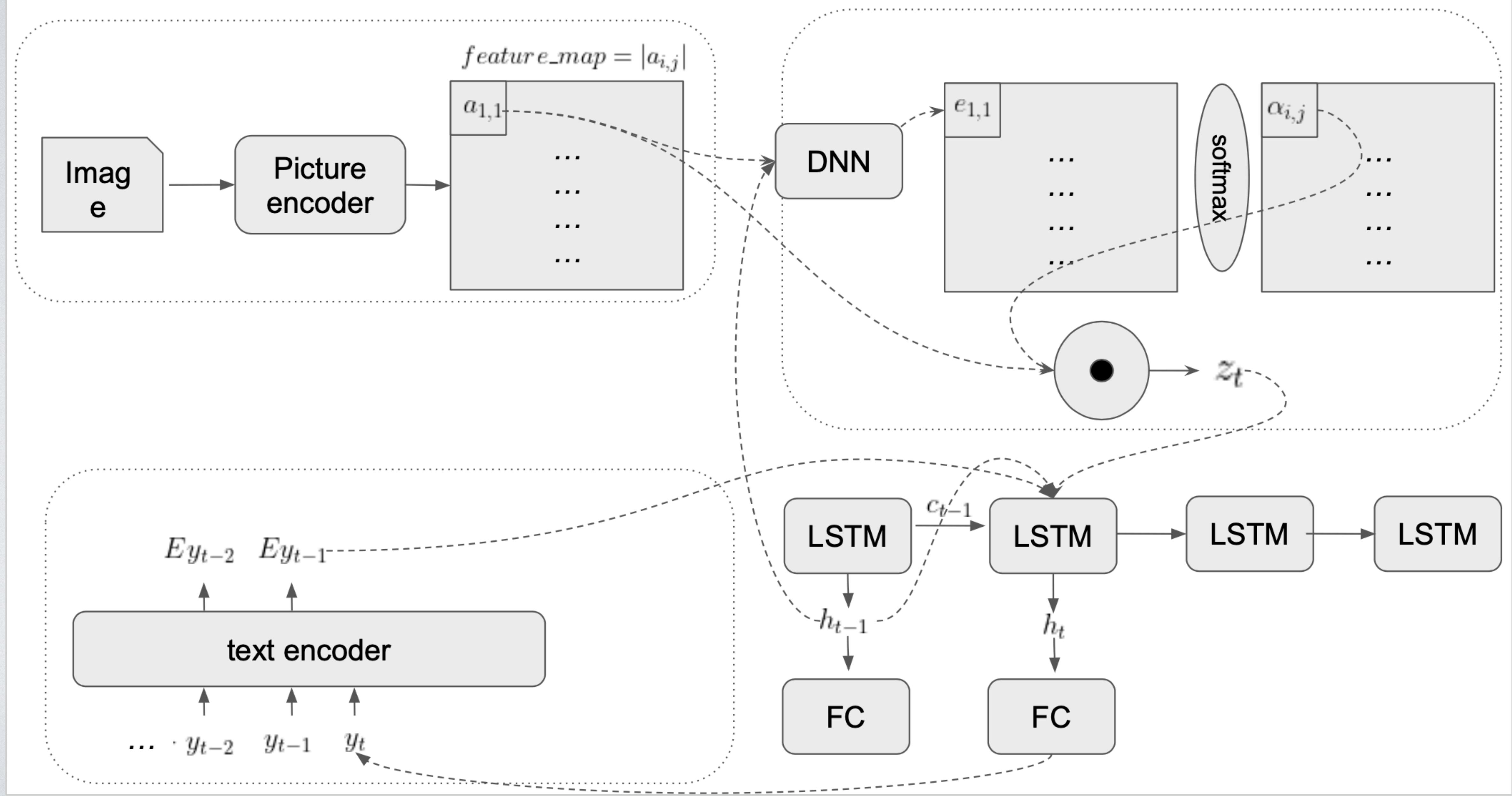
ДО НЕЙРОСЕТЕЙ

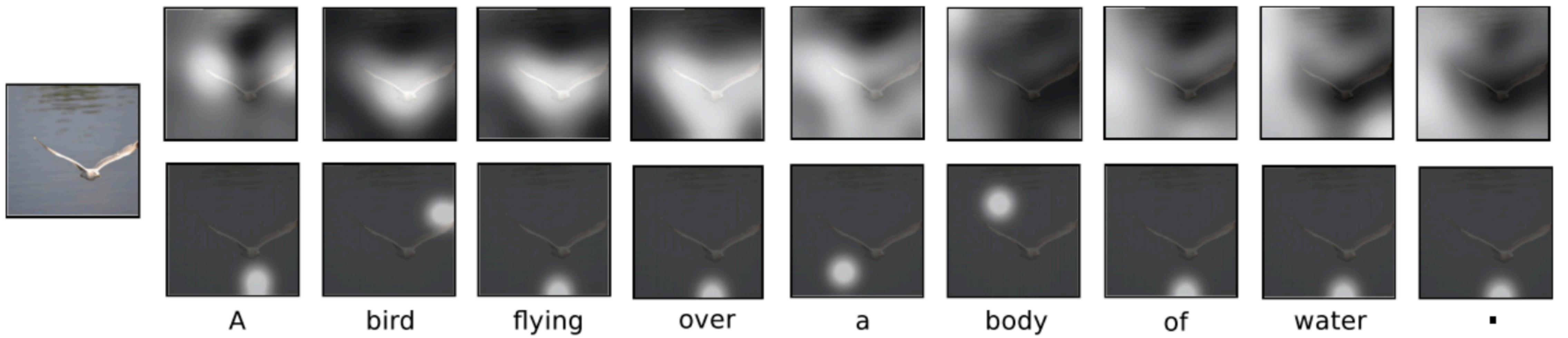


ВСПОМНИМ ATTENTION : IMAGES

A young boy is playing basketball. 	Two dogs play in the grass. 	A dog swims in the water. 
A group of people walking down a street. 	A group of women dressed in formal attire. 	Two children play in the water. 
A skier is skiing down a snowy hill. 	A little girl in a pink shirt is swinging. 	A dog jumps over a hurdle. 

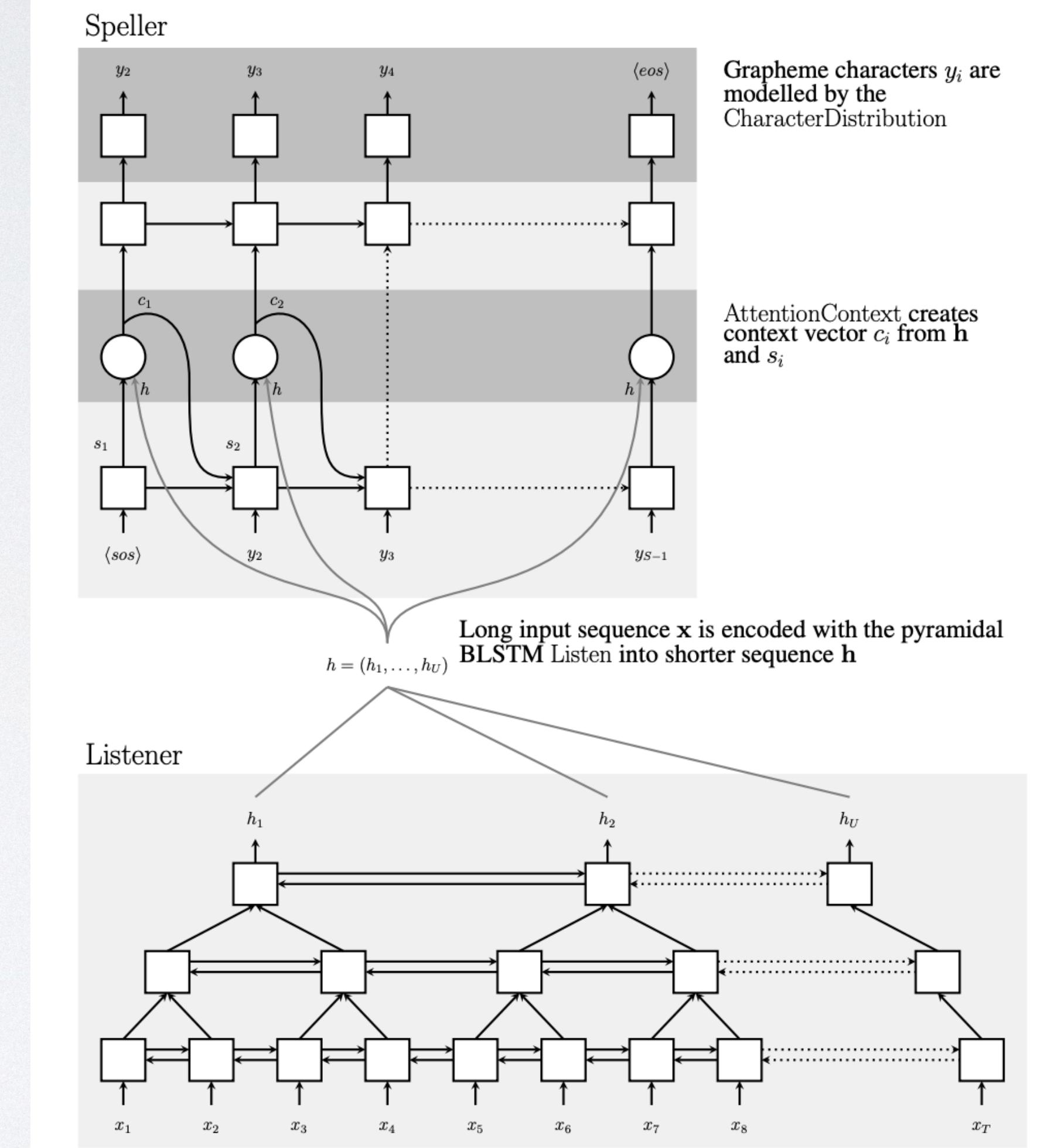
Show attend and tell [11] for image captioning task





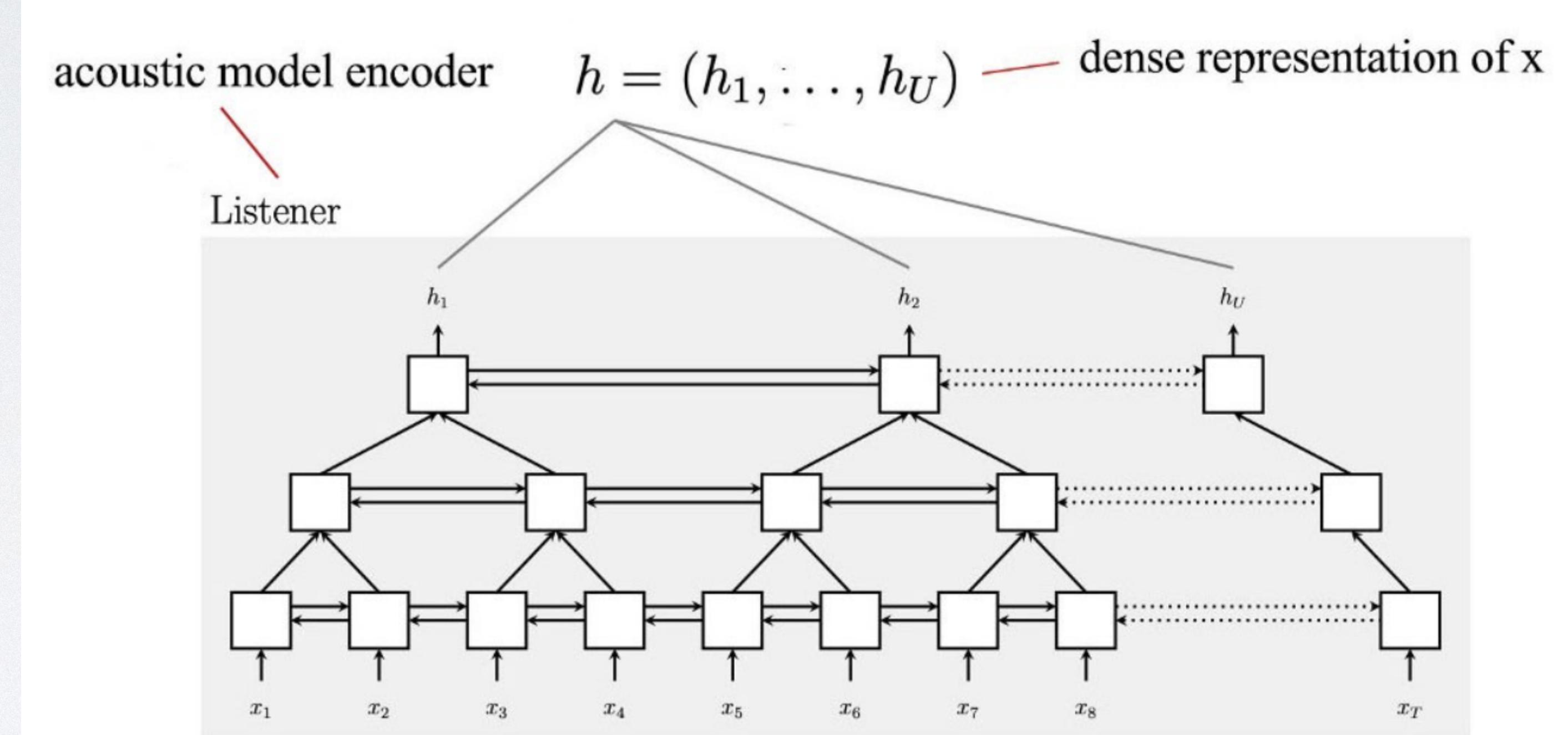
LISTEN. ATTEND. SPELL

- Цель - избавиться от независимости выходов
- Listener - енкодер
- Speller - декодер (получаем транскрипцию)
- Attend - сообщает декодеру, какая часть дорожки релевантна в данный шаг



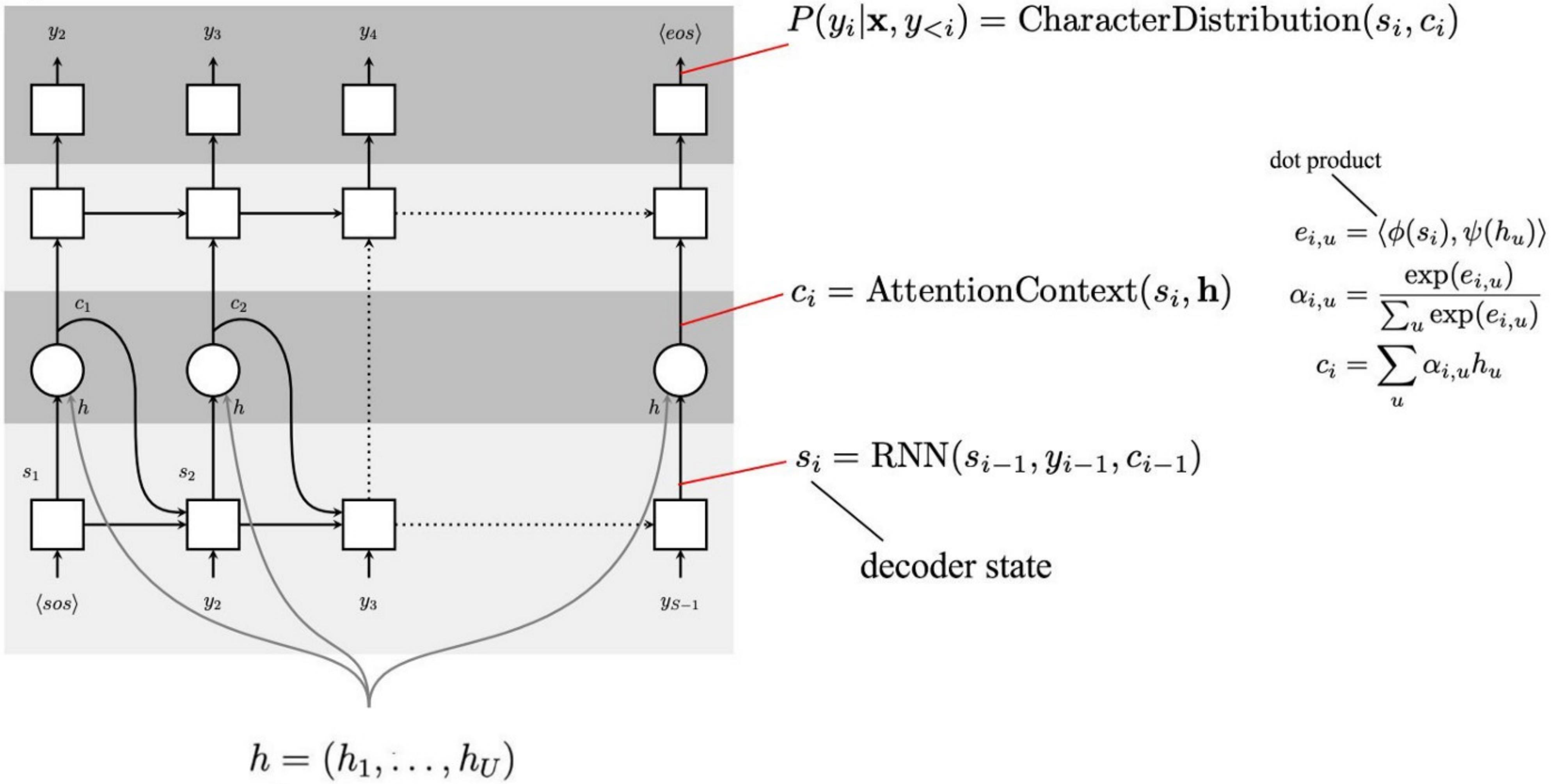
LISTENER

- Пирамидальная bi-LSTM

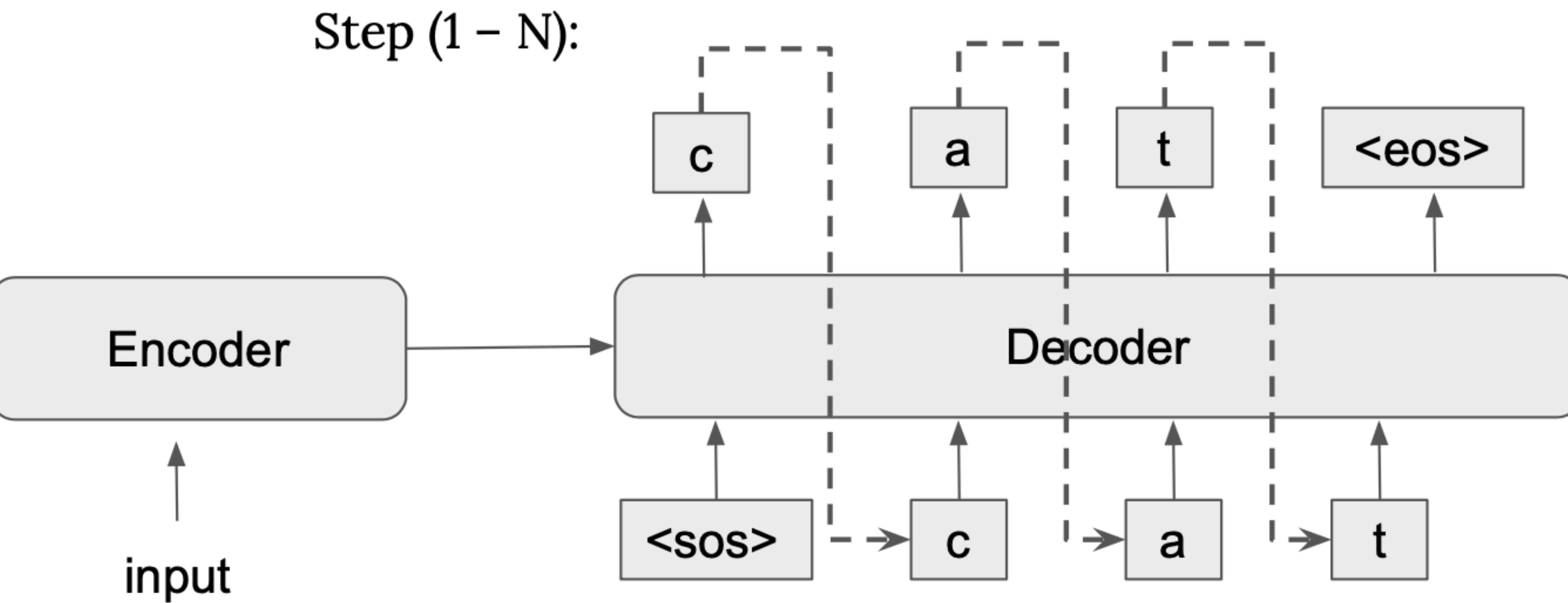


SPELLER

Speller



INFERENCE



- $\text{total time} = I \times \text{run_encoder_time} + T \times \text{run_decoder_time}$

CTC VS LAS

	CTC beam search	LAS beam search
Computation al cost	<ul style="list-style-type: none">• $1 * \text{run_encoder}()$• $T * \text{beam_size} * \text{expand_beam}()$	<ul style="list-style-type: none">• $1 * \text{run_encoder}()$• $T * \text{beam_size} * \text{run_decoder}()$
Path merging		
Overall experience	<ul style="list-style-type: none">• painful to code• joyful to run	<ul style="list-style-type: none">• painful to code• painful to run

РЕЗУЛЬТАТЫ

Method	Year	WER (test-clean)	WER (test-other)
Human	~ 200 000 b.c.	5.83	12.69
Deep Speech 2 [5]	2015	5.15	12.73
LAS (original paper [8])	2015	-----	-----
LAS [9*]	2019	3.2	9.8

ВОПРОСЫ

- Что подается на вход акустической модели?
- Чем отличается beam-search в СТС и LAS?
- В чем основная цель Listen, Attention, Spell подхода?

РЕСУРСЫ

- <https://www.youtube.com/watch?v=3MjlkWxXigM> lectures
- <https://arxiv.org/pdf/1508.01211.pdf> Listen.Attend.Spell
- <https://habr.com/ru/company/vk/blog/579412/> Как работает
распознавание речи ВК
- <https://www.youtube.com/watch?v=fodf4Pttve4> Введение в
обработку звука