

Синтез речи

Лекция №1

Гриша Стерлинг, SberDevices

1. Введение
2. Метрики качества
3. Обзор задачи
4. Вокодеры (ч. 1)

Синтез речи

Задача:

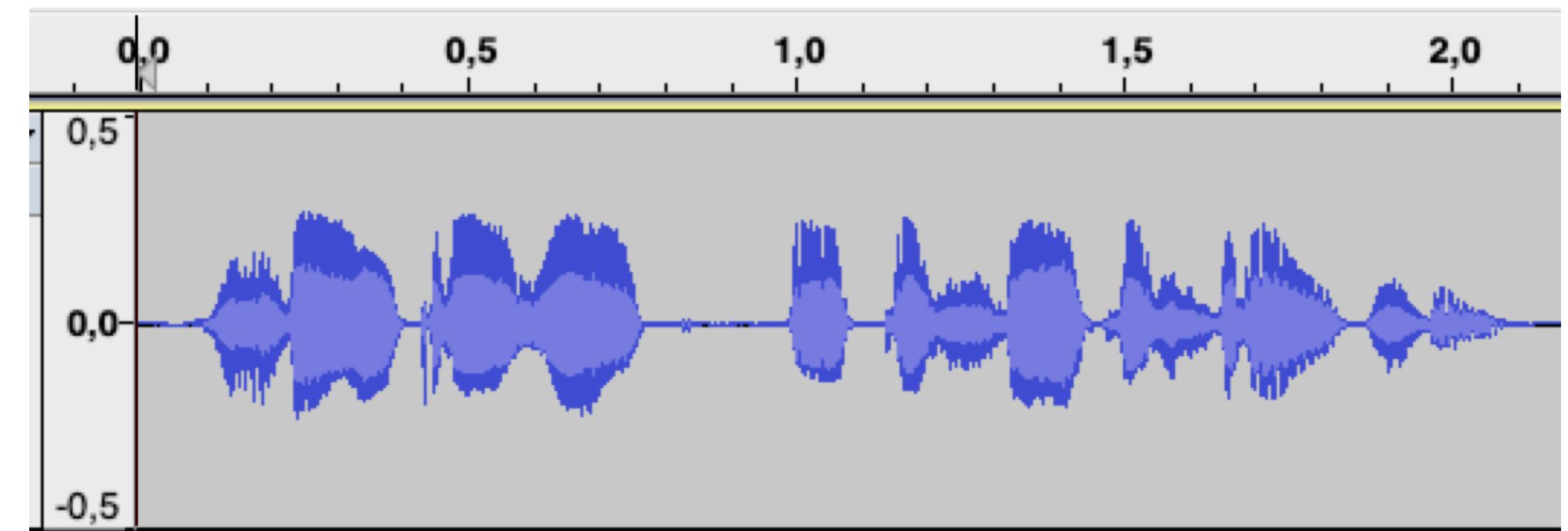
- озвучить заданный текст голосом

Задачи со звездочкой:

- style transfer
- несуществующим голосом
- controllable speech synthesis
- на другом языке
- эмоции
- шепот
- субвокализации, смех

«Всем привет, это синтез речи»

->

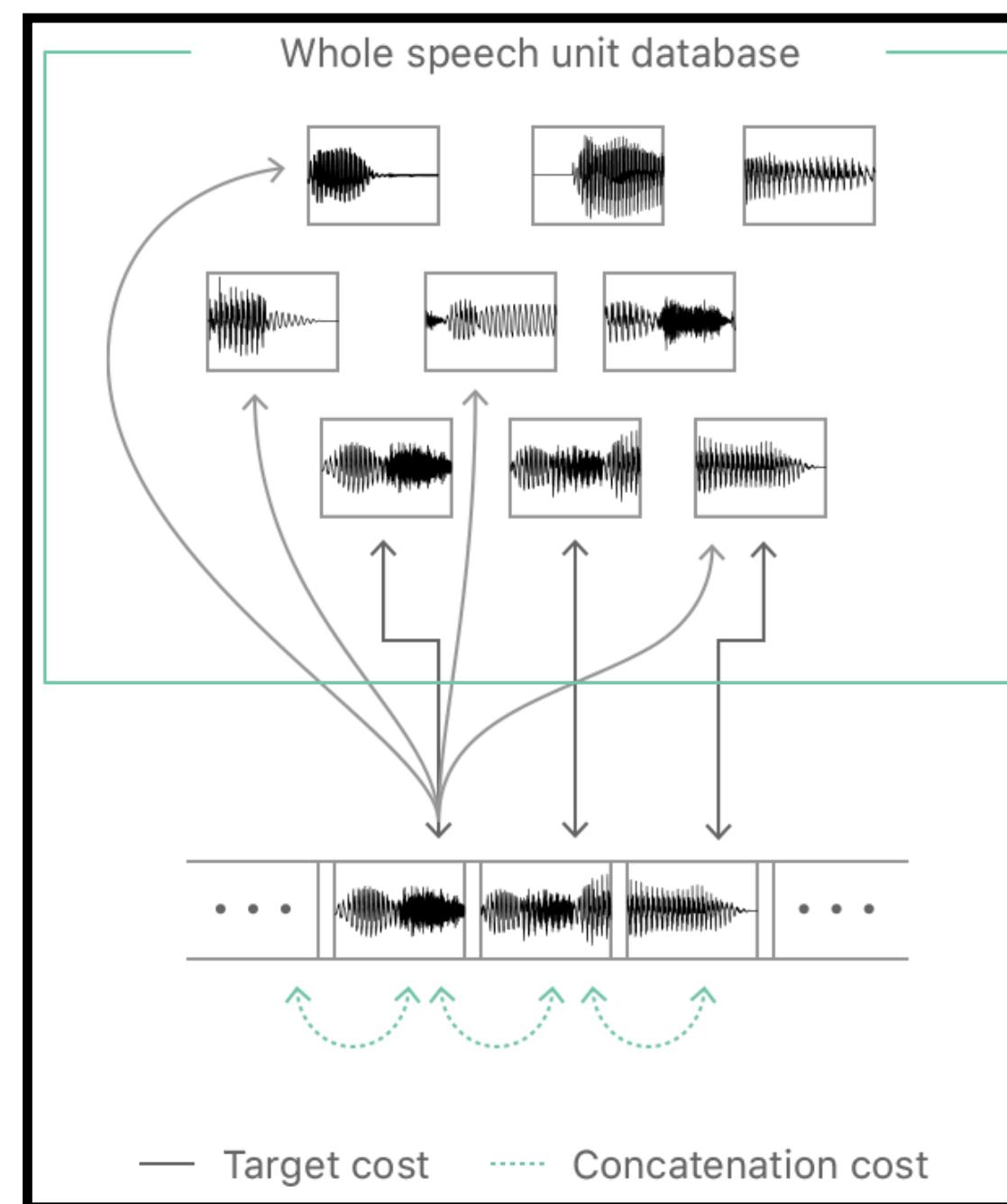


Синтез речи

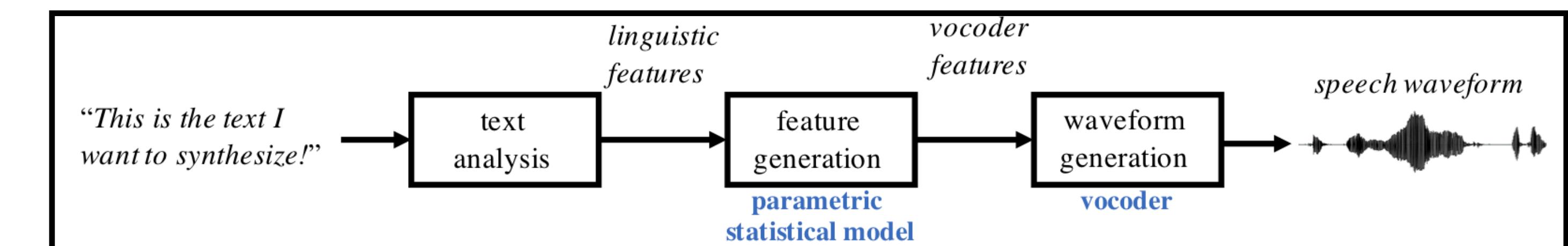


Concatenative
(unit selection)

Parametric



- 2 стадии:
- акустическая модель
 - вокодер (**voice coder**)



Метрики качества

Crowdsourced:

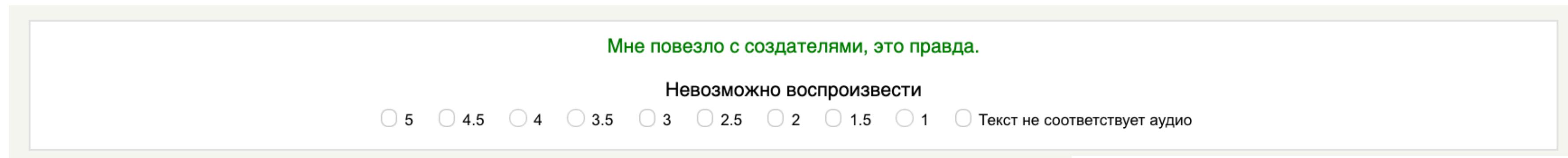
- MOS (mean opinion score)
- MUSHRA (Multiple Stimuli with Hidden Reference and Anchor)
- PSER (Pronunciation sentence error rate)
- SBS

Automatic:

- SNR (sound to noise ratio)
- Frechet distance
- nn-based mos
- MCD (mel-cepstral distortion)

Метрики качества

MOS:



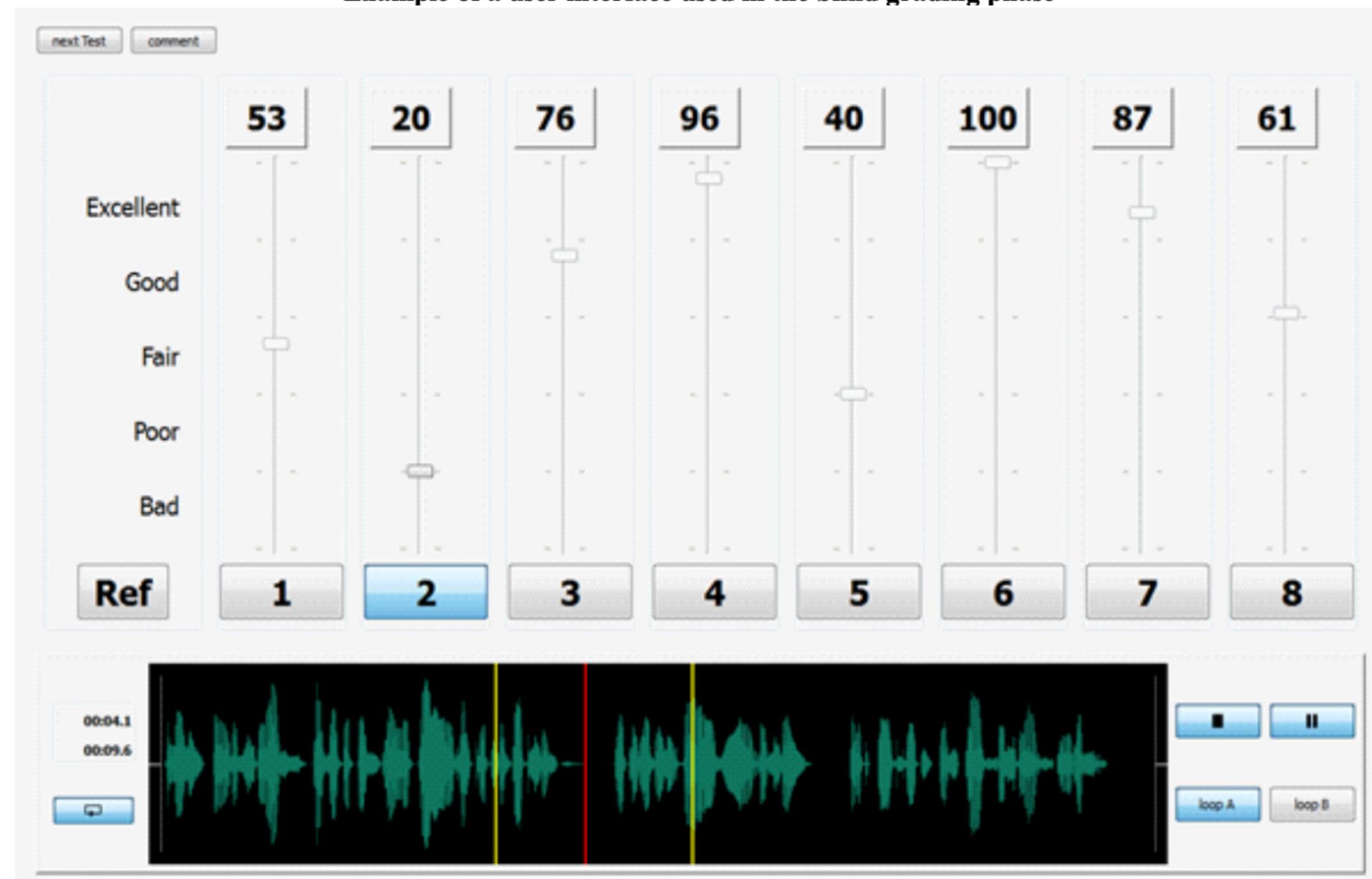
$$MOS = \frac{\sum_{n=1}^N R_n}{N}$$

Where R are the individual ratings for a given stimulus by N subjects.

FIGURE 4

Example of a user interface used in the blind grading phase

MUSHRA:



- + reference
- + hidden reference
- + few samples
- + 1-2 anchors

Метрики качества

SBS:

Я не смог подтвердить перевод. давайте попробуем позже

Вариант А

▶ 0:04 / 0:04

Вариант В

▶ 0:04 / 0:04

1 ○ Точно А 2 ○ Не могу выбрать 3 ○ Точно В

PSER:

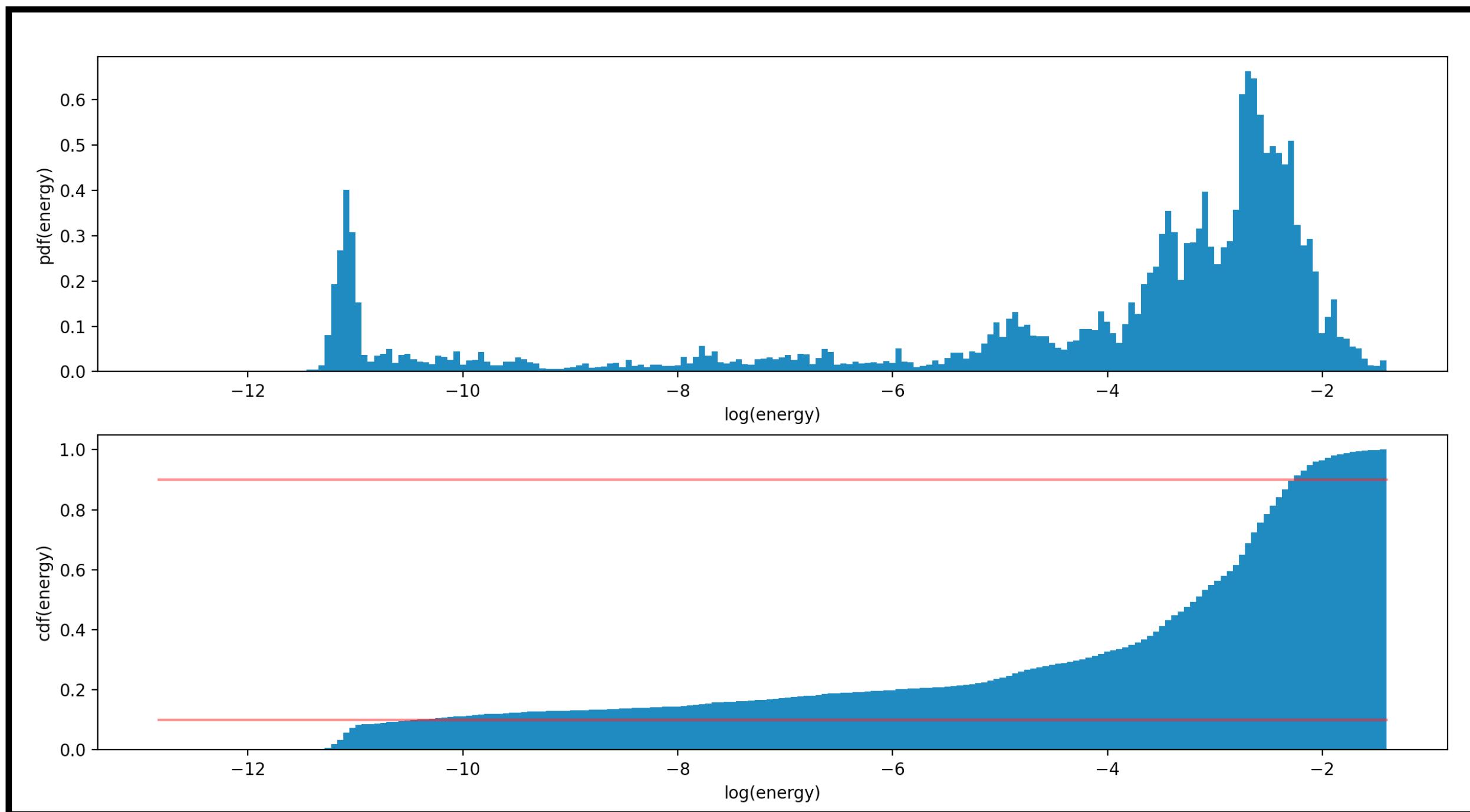
благодарю! мои разработчики будут рады!

▶ 0:00 / 0:03

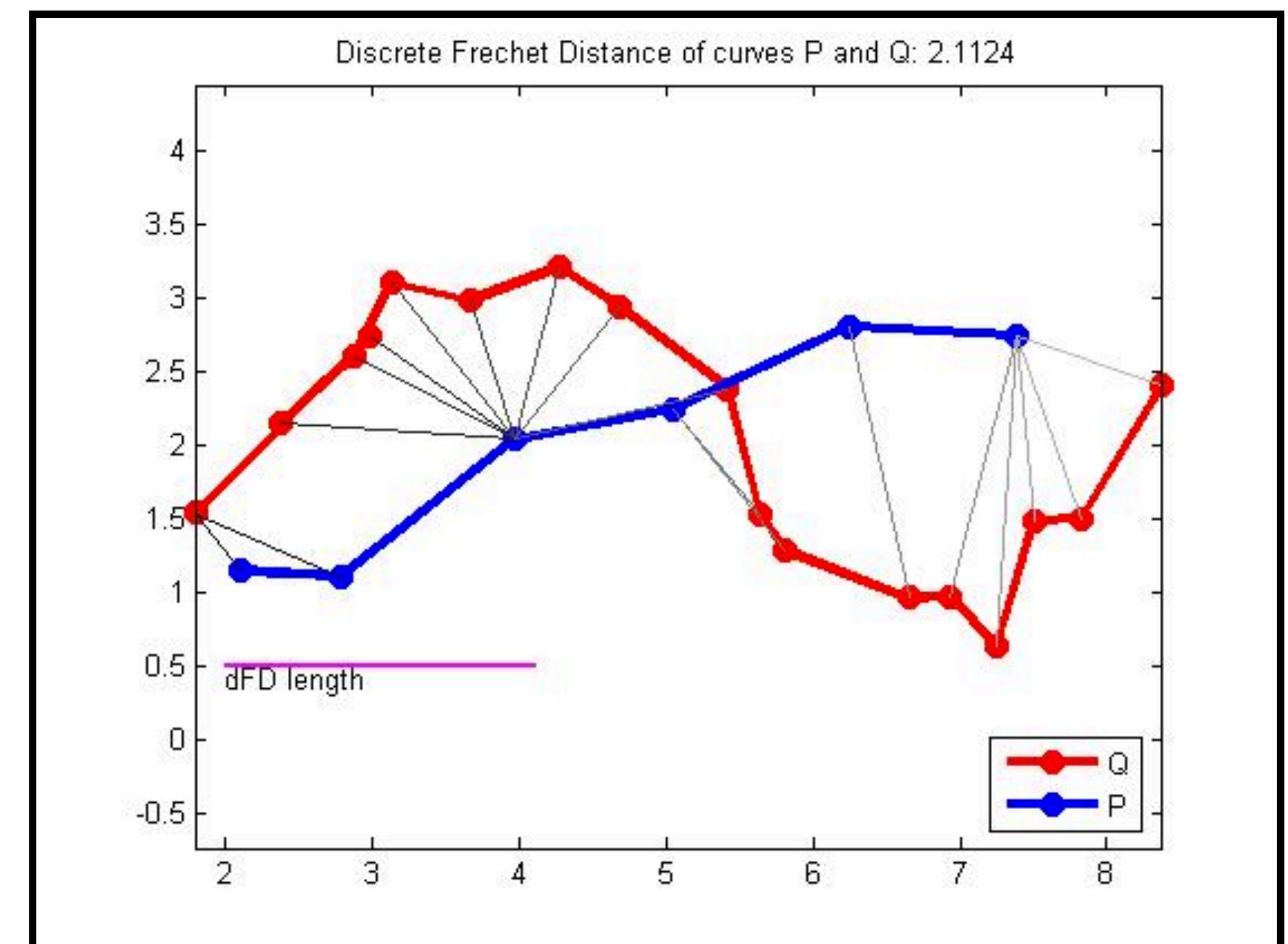
- Лишние паузы
- Не хватает пауз
- Неверное произношение
- Пропущены нужные / вставлены лишние звуки
- Некачественное аудио
- Неверная интонация
- Наличие заиканий
- Всё верно

Метрики качества

SNR:



Frechet distance:



Метрики качества

MCD:

$$\frac{10\sqrt{2}}{\ln 10} \frac{1}{T} \sum_{t=1}^T \sqrt{\sum_i (C_{ti} - \hat{C}_{ti})^2}.$$

nn-based:

MBNet: MOS Prediction for Synthesized Speech with Mean-Bias Network

[Yichong Leng](#), [Xu Tan](#), [Sheng Zhao](#), [Frank Soong](#), [Xiang-Yang Li](#), [Tao Qin](#)

Neural MOS Prediction for Synthesized Speech Using Multi-Task Learning With Spoofing Detection and Spoofing Type Classification

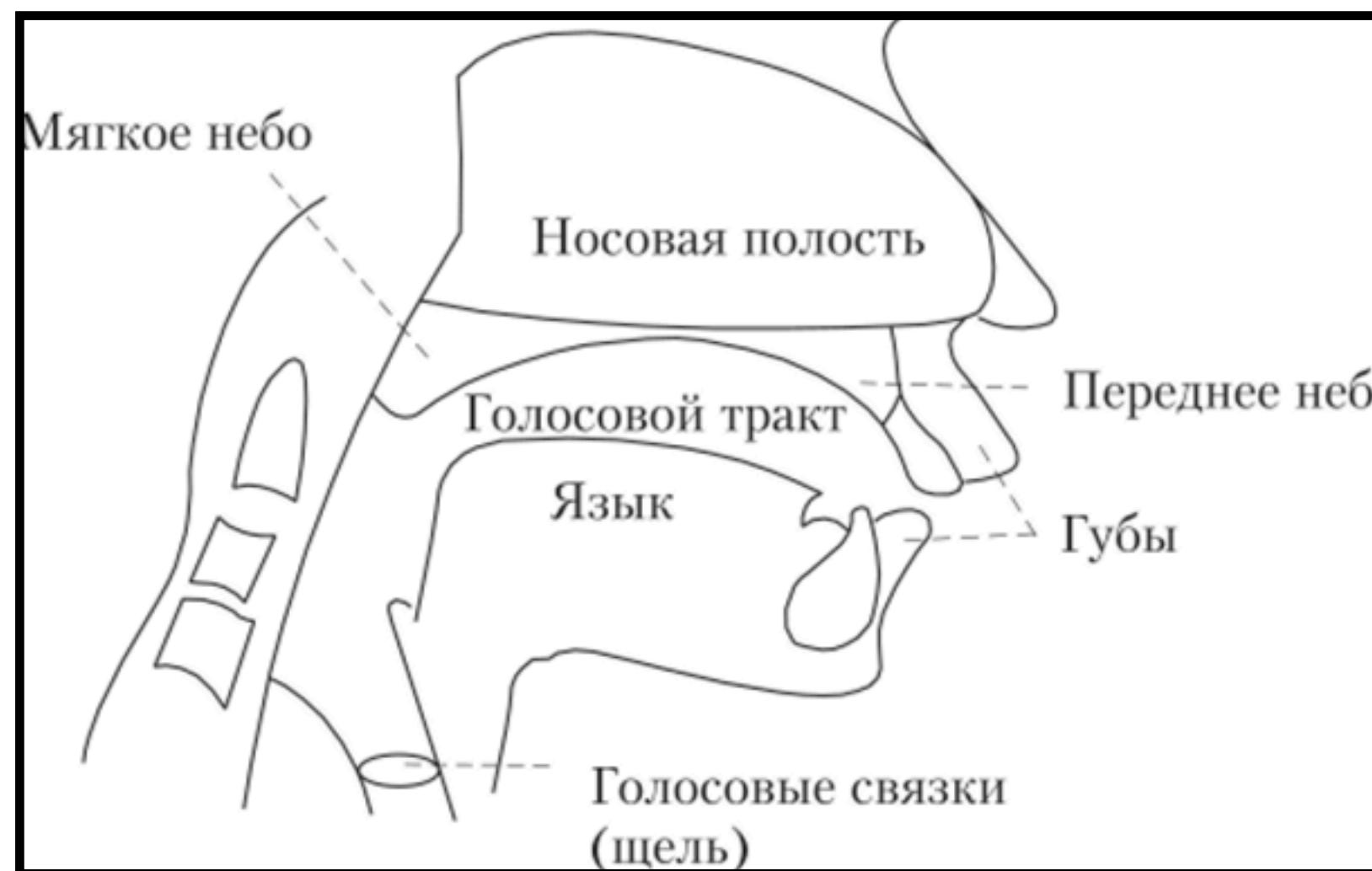
[Yeunju Choi](#), [Youngmoon Jung](#), [Hoirin Kim](#)

MOSNet: Deep Learning based Objective Assessment for Voice Conversion

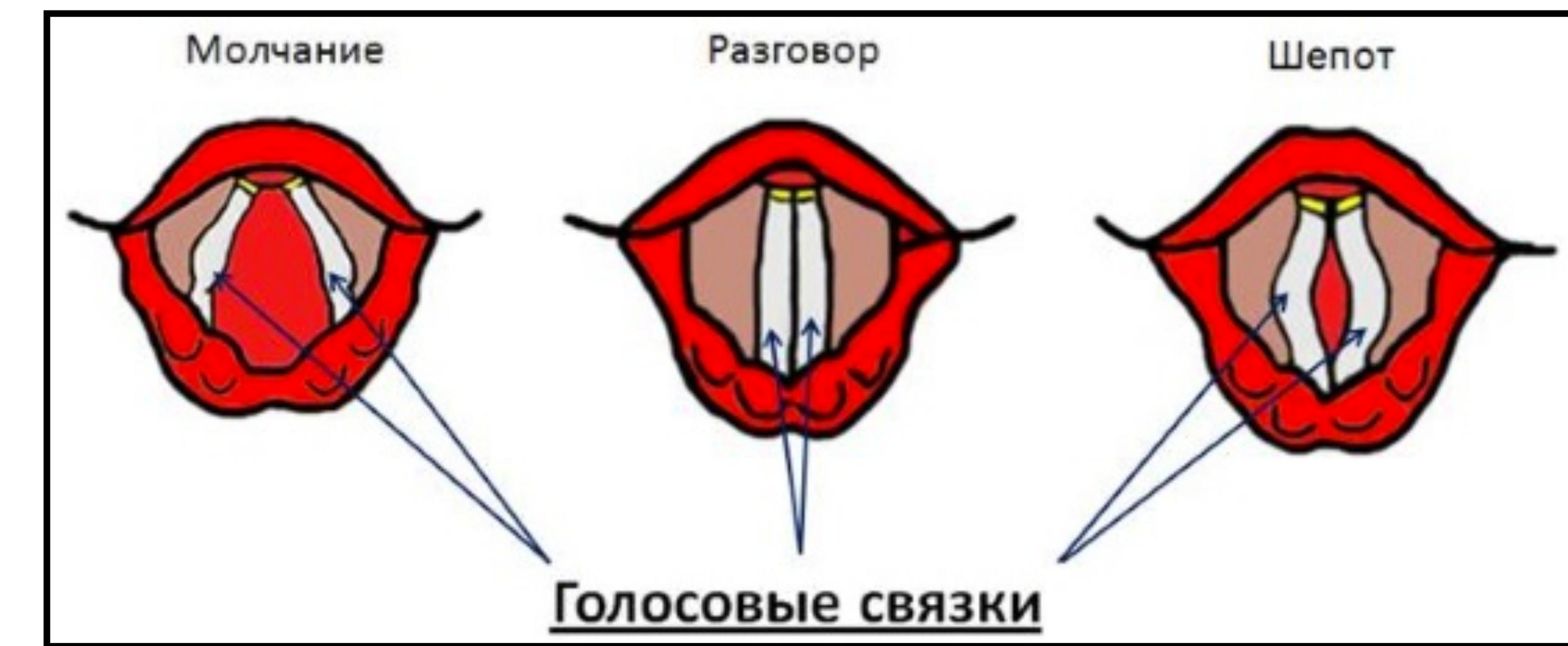
[Chen-Chou Lo](#), [Szu-Wei Fu](#), [Wen-Chin Huang](#), [Xin Wang](#), [Junichi Yamagishi](#), [Yu Tsao](#), [Hsin-Min Wang](#)

Что такое речь

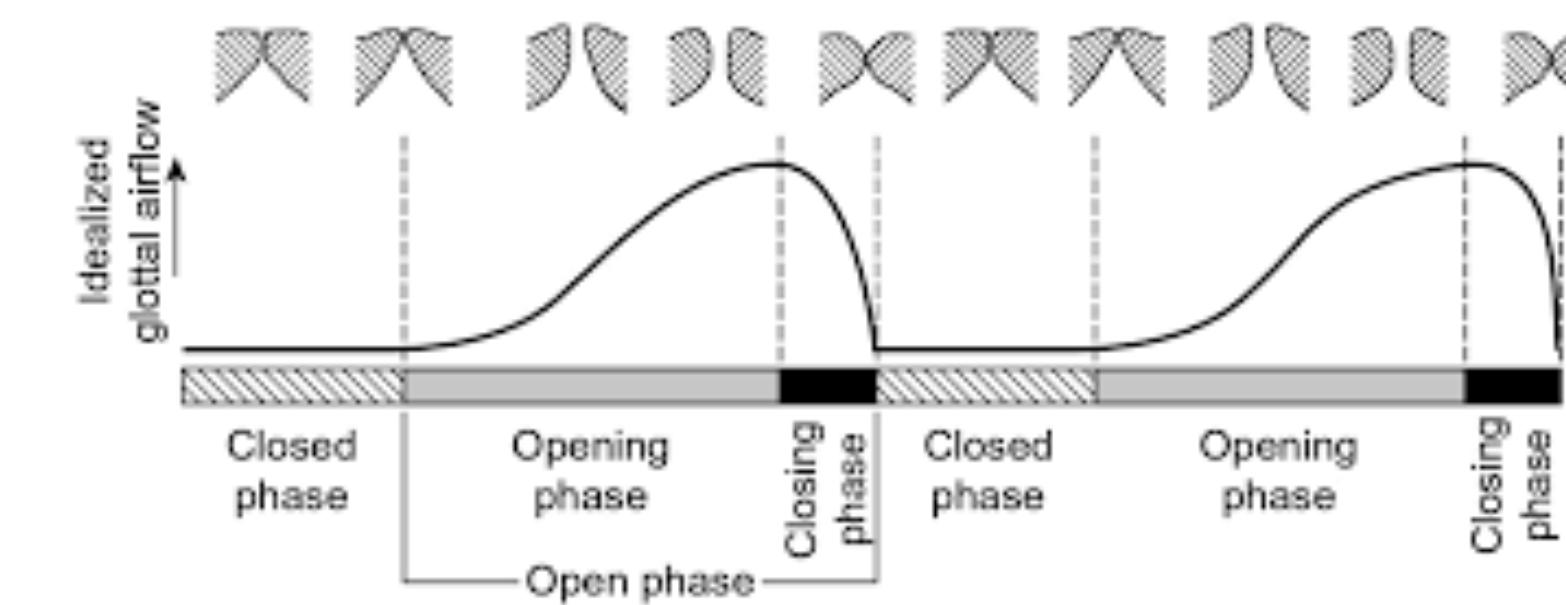
Речевой тракт человека



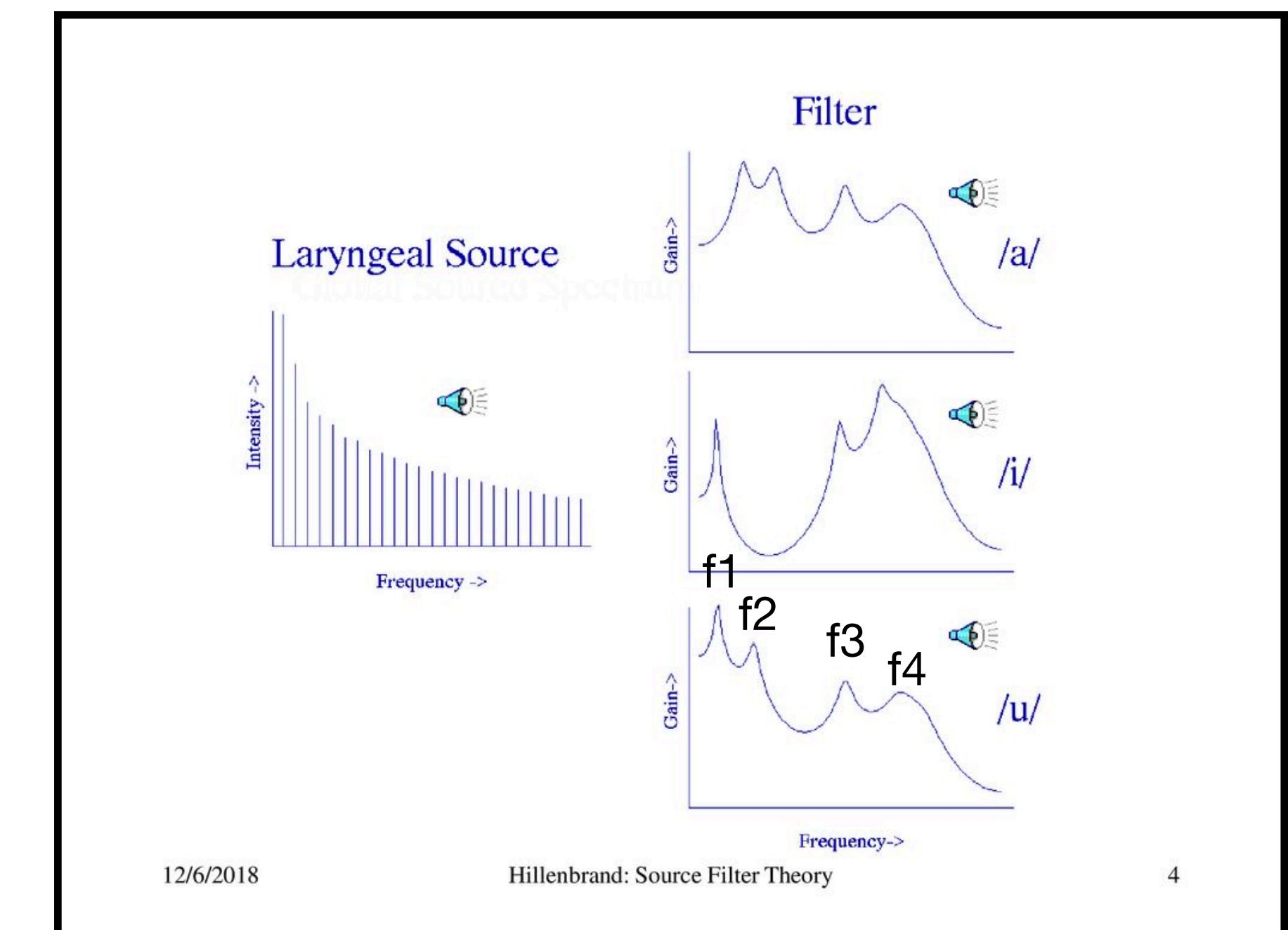
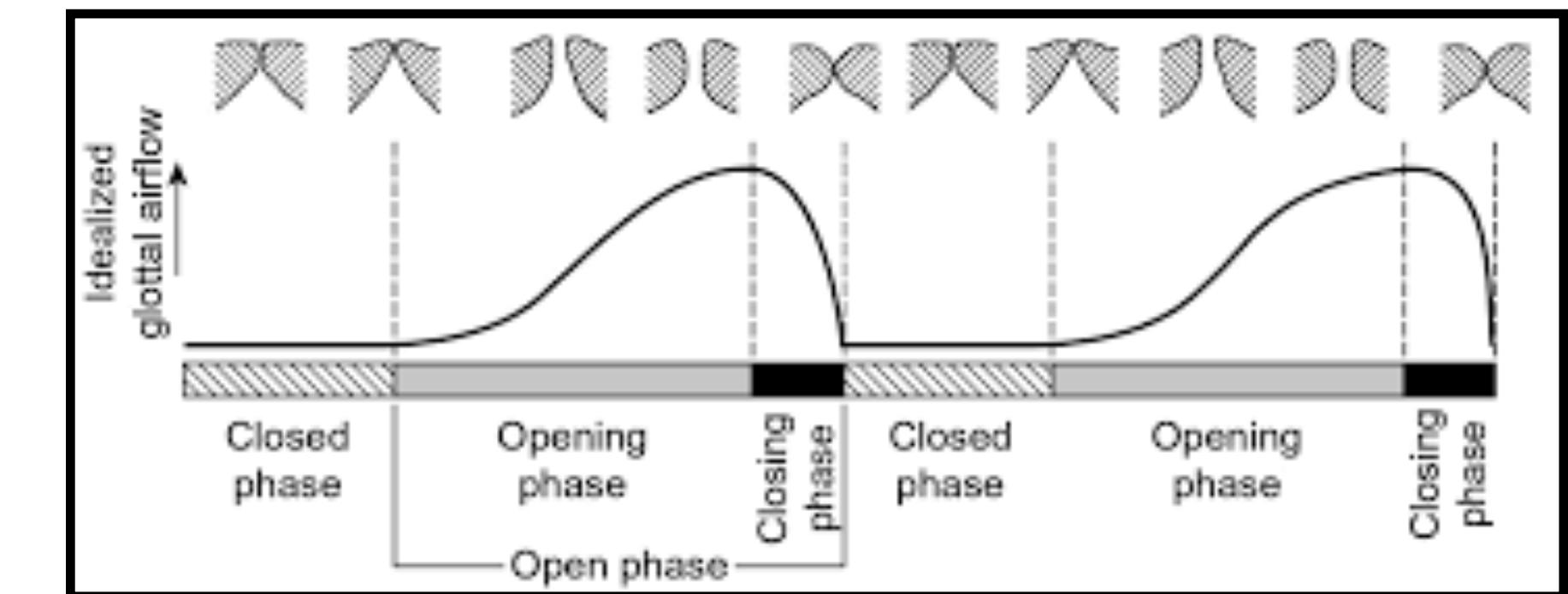
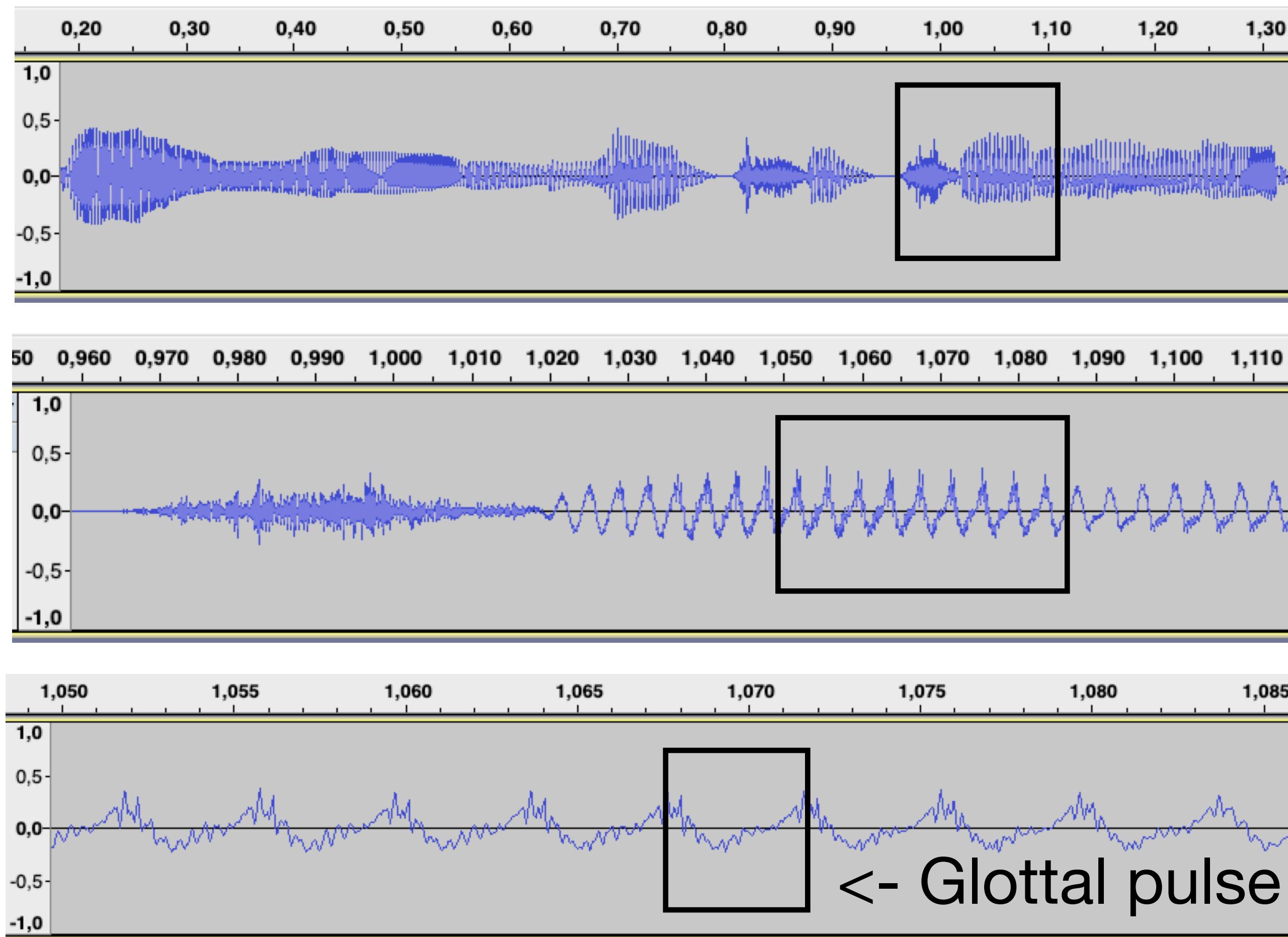
Голосовые связки



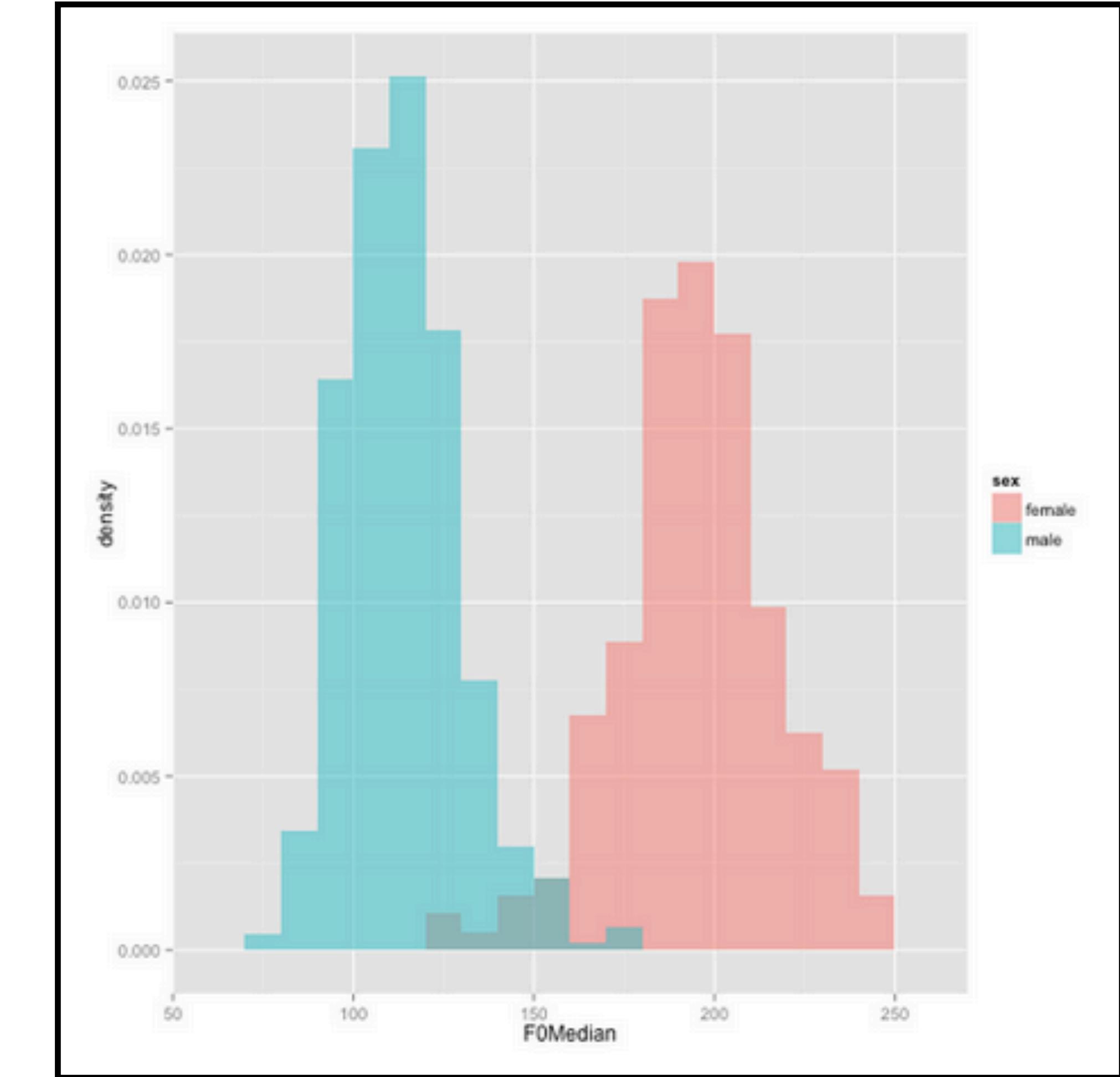
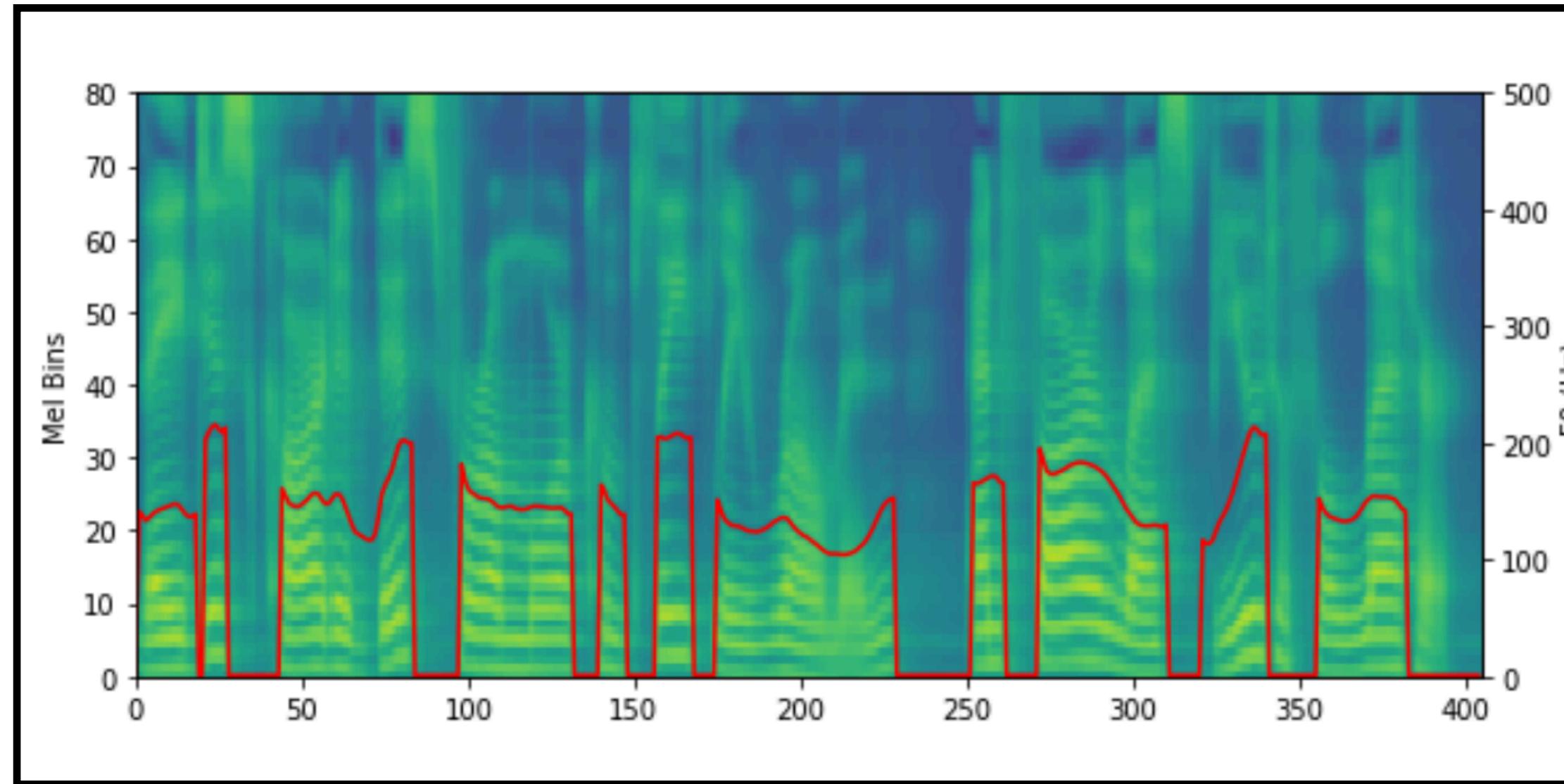
$$y(t) = h(t) * x(t),$$



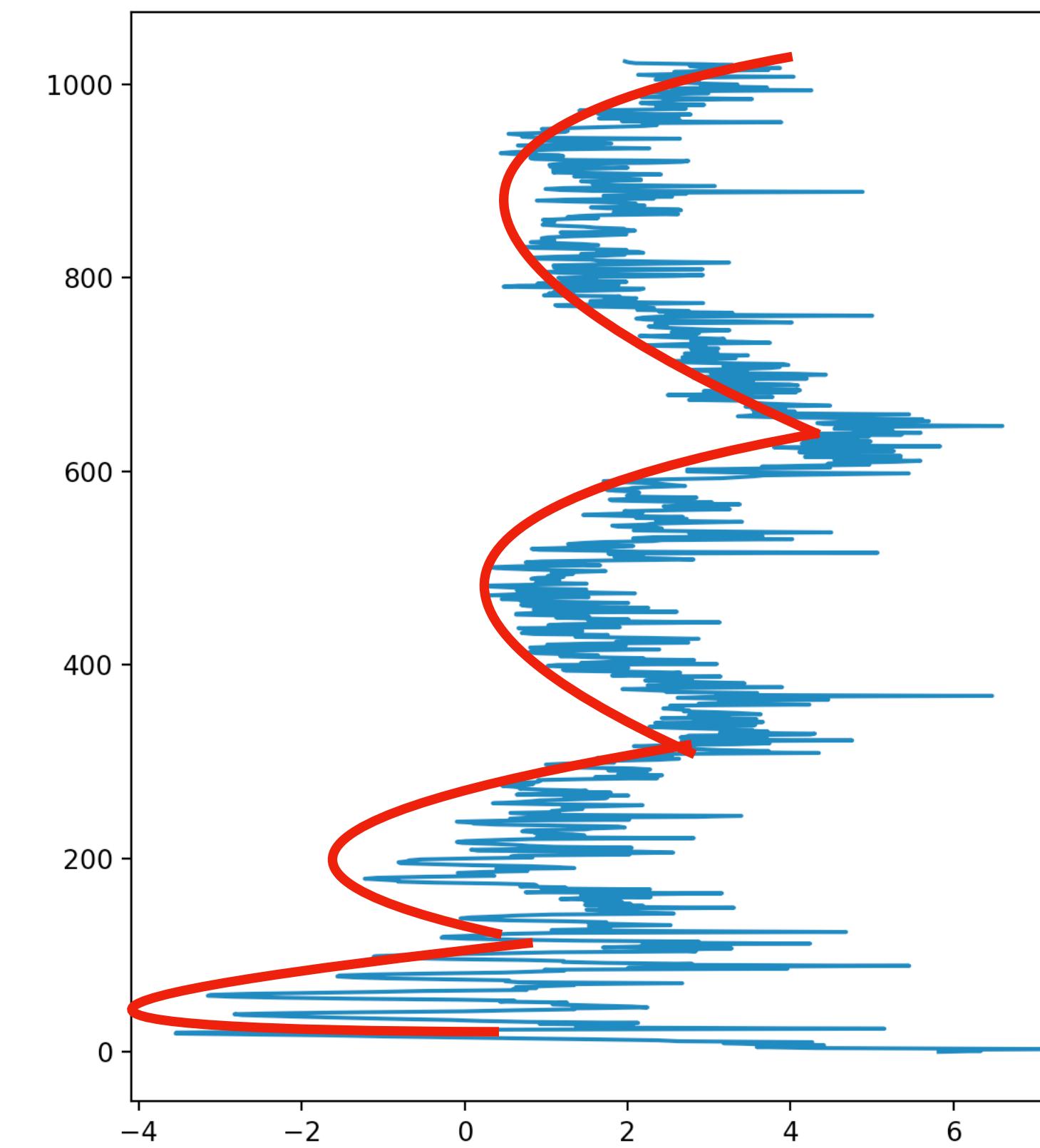
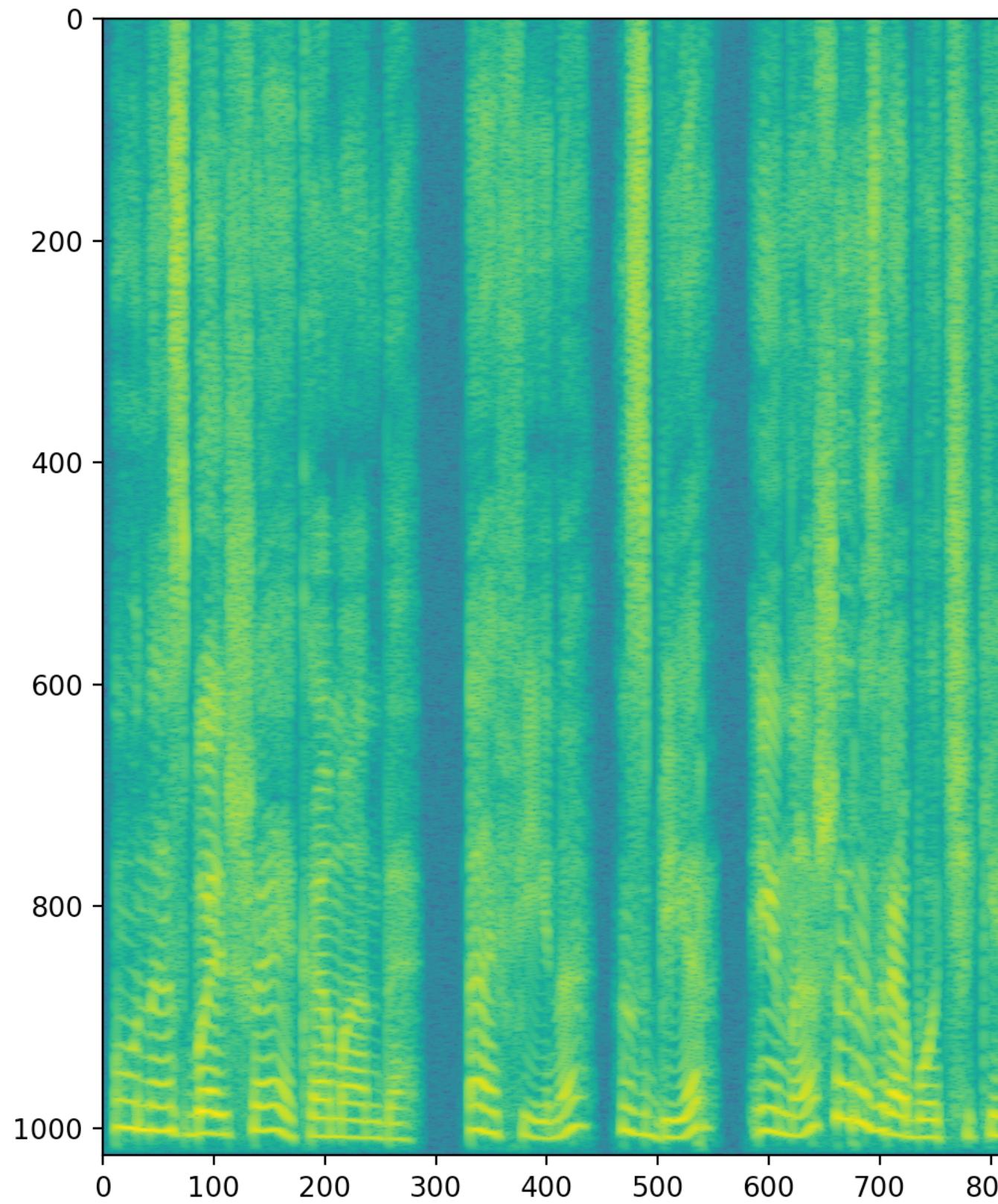
Что такое речь



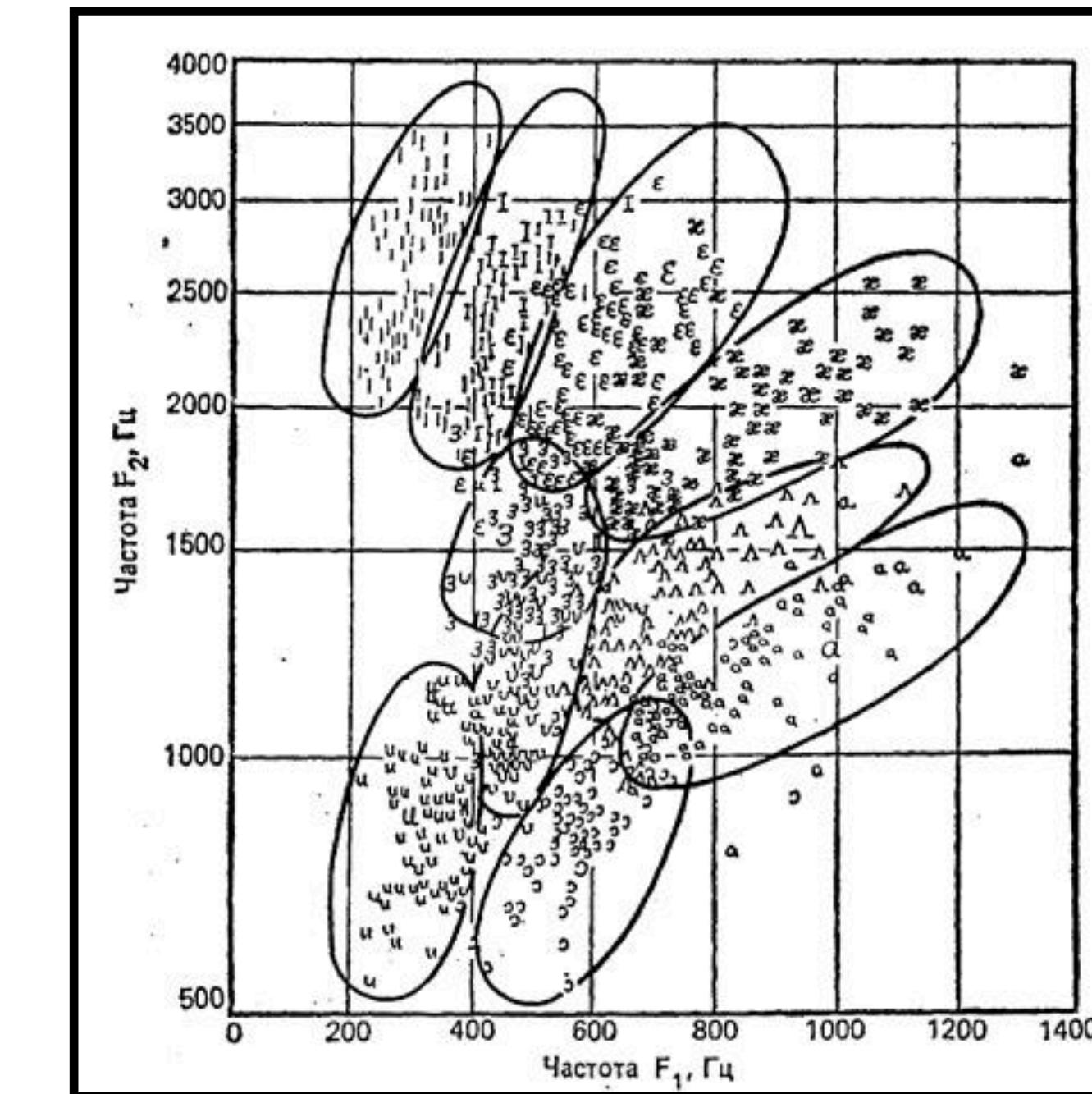
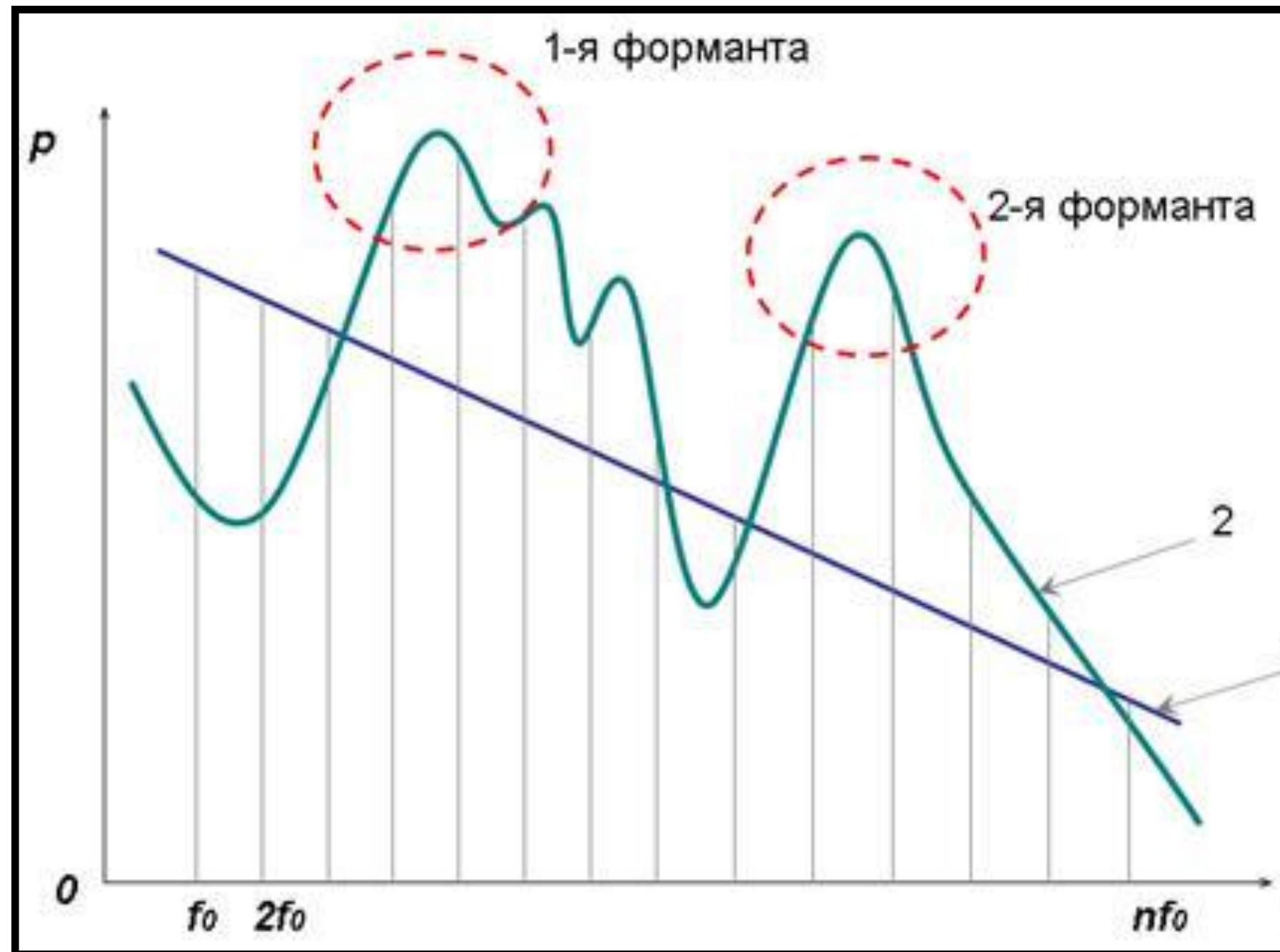
pitch (f0, частота основного тона)



Форманты



Форманты



Вокодеры

Не авторегрессионные:

- Griffin-Lim algorithm
- WORLD
- WaveGlow
- GAN-based

Авторегрессионные:

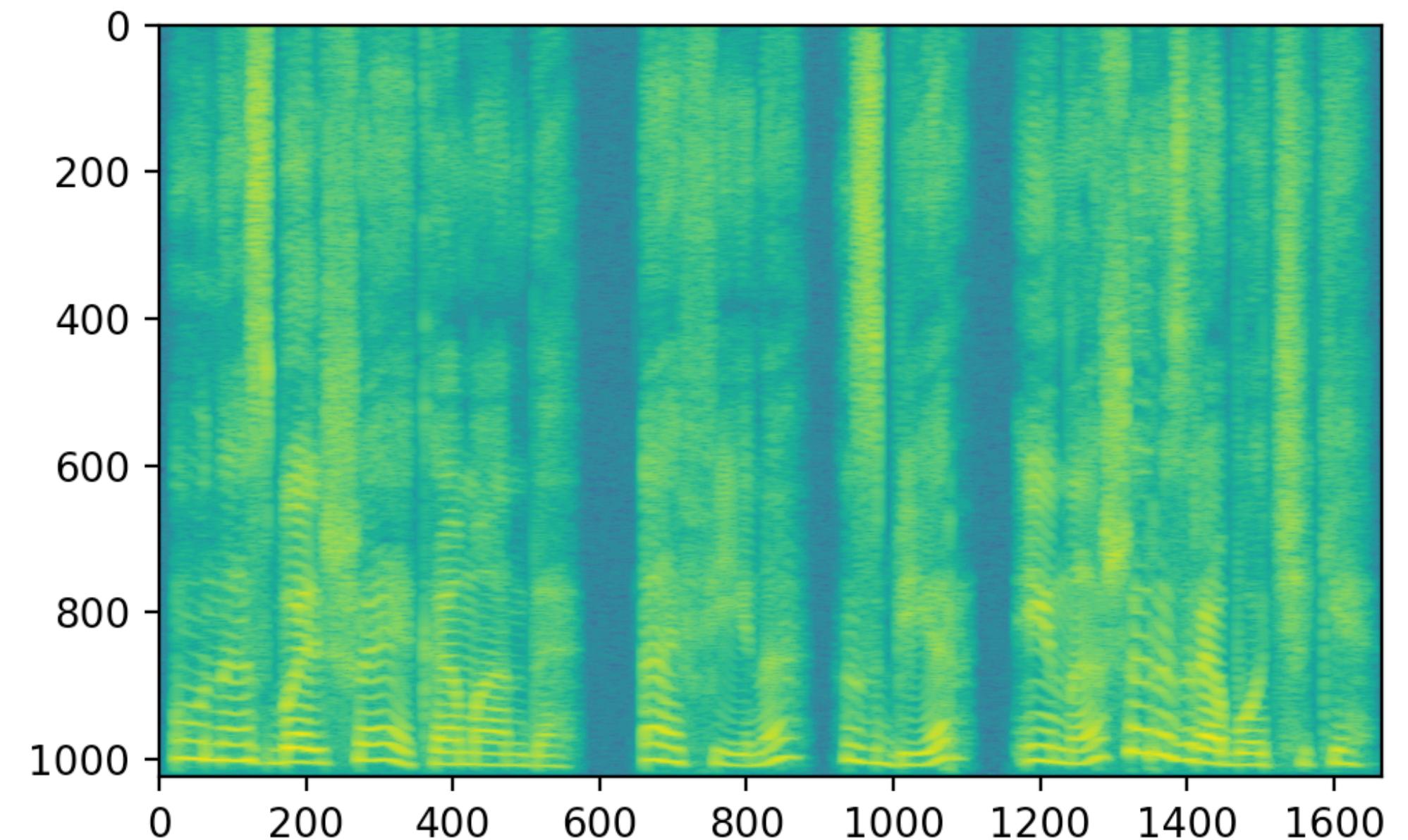
- WaveNet
- WaveRNN
- LPCNet

Griffin-Lim algorithm

Phase reconstruction:

a - искомый сигнал

$X = |\text{FFT}(a)|$ - спектrogramma

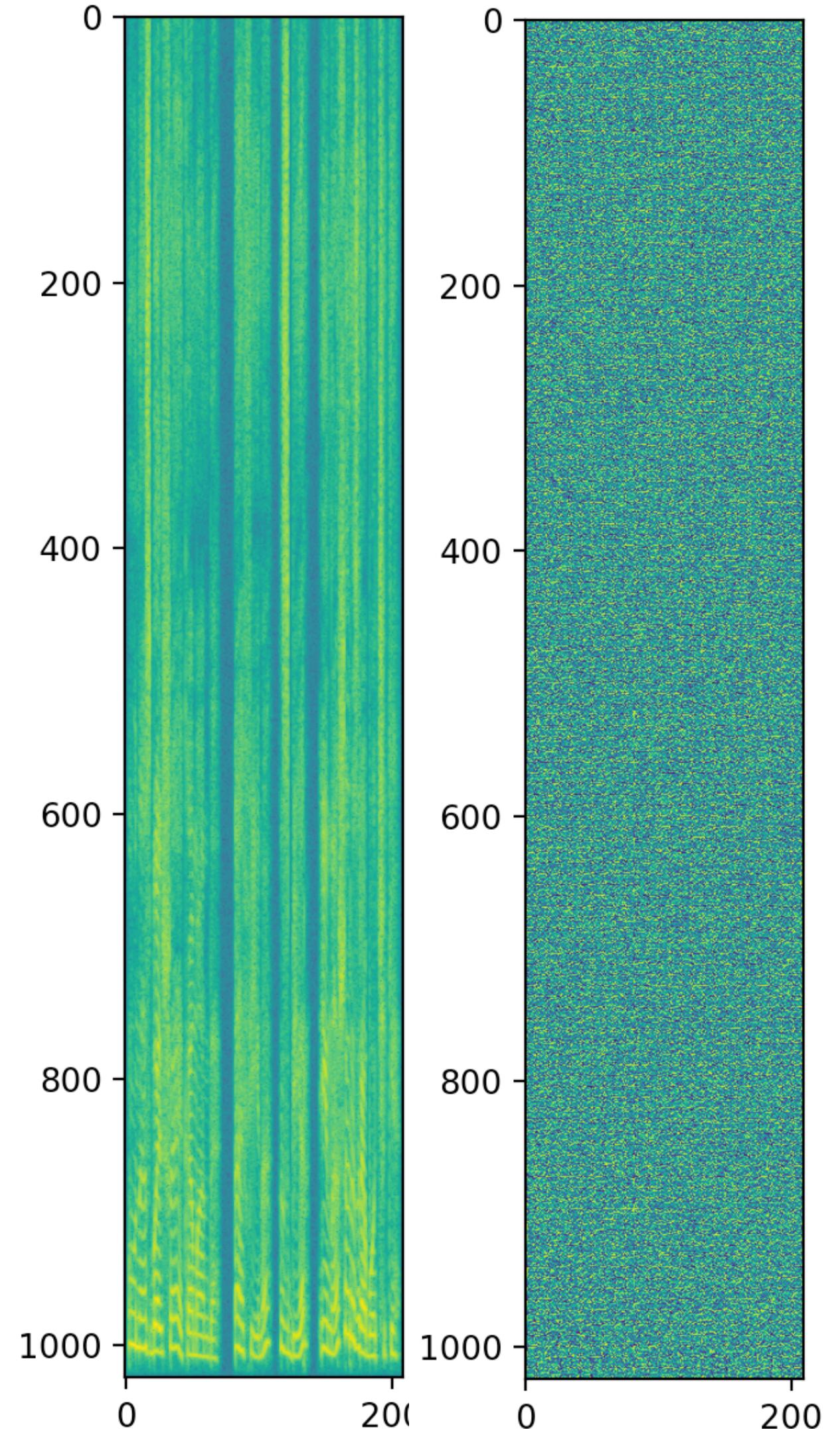
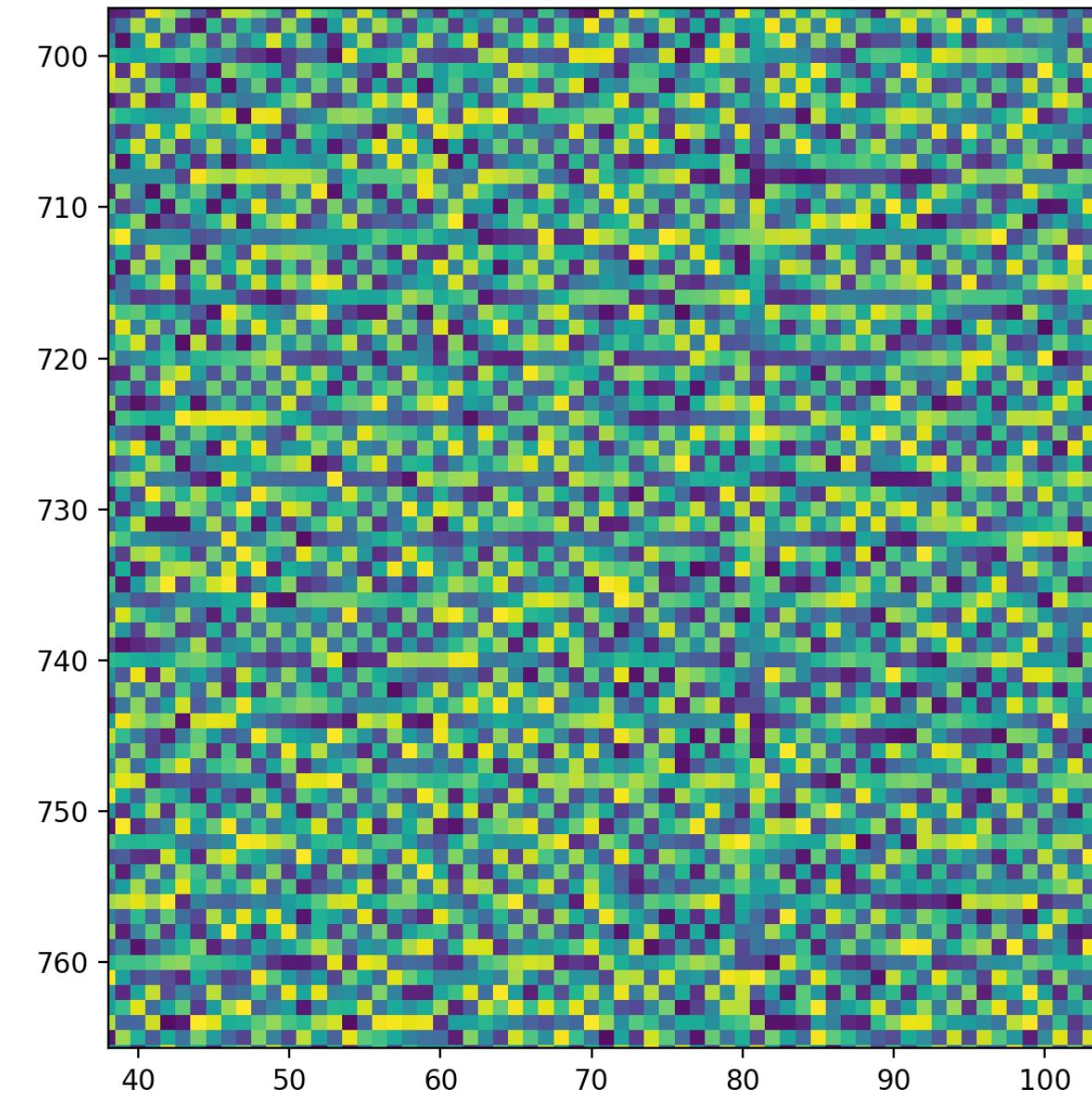


Griffin-Lim algorithm

Алгоритм:

1. Инициализируем фазу Y_0 случайно
2. $a_1 = \text{iFFT}(X e^{\wedge} iY_0)$
3. $Y_1 = \text{phase}(\text{FFT}(a_1))$
4. Итерируемся до сходимости
 $(a_j = \text{iFFT}(X e^{\wedge} iY_{j-1}))$

$$\min_{\mathbf{X}} \|\mathbf{X} - P_{\mathcal{C}}(\mathbf{X})\|_{\text{Fro}}^2 \text{ s.t. } \mathbf{X} \in \mathcal{A}$$



LPC (linear predictive coding)

Вейвформа - это свертка импульса с фильтром:

$$1. \quad y(t) = h(t) * x(t),$$

$$2. \quad x_t = \sum_{p=1}^P a_p x_{t-p} + \epsilon_t$$

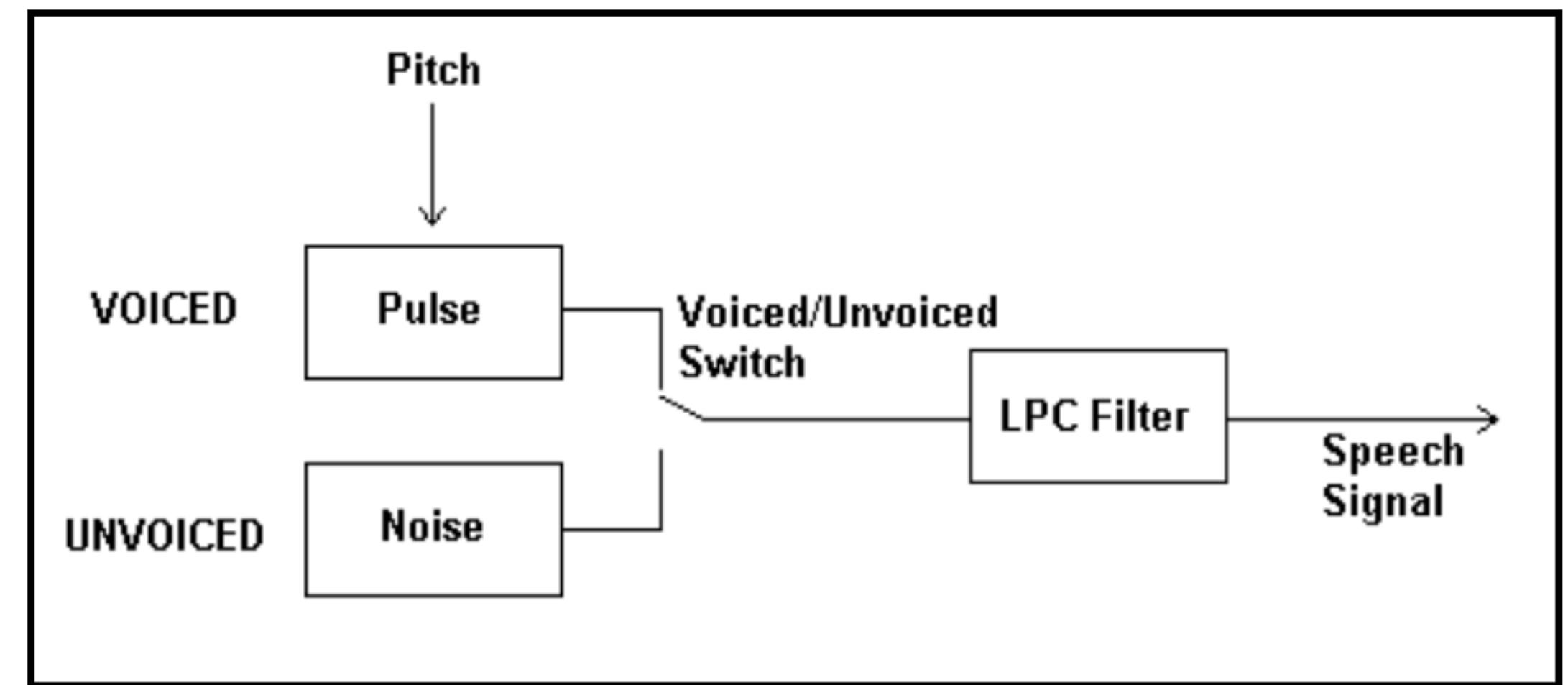


Figure 11 : LPC Decoder

Wavenet

$$2. \quad x_t = \sum_{p=1}^P a_p x_{t-p} + \epsilon_t$$

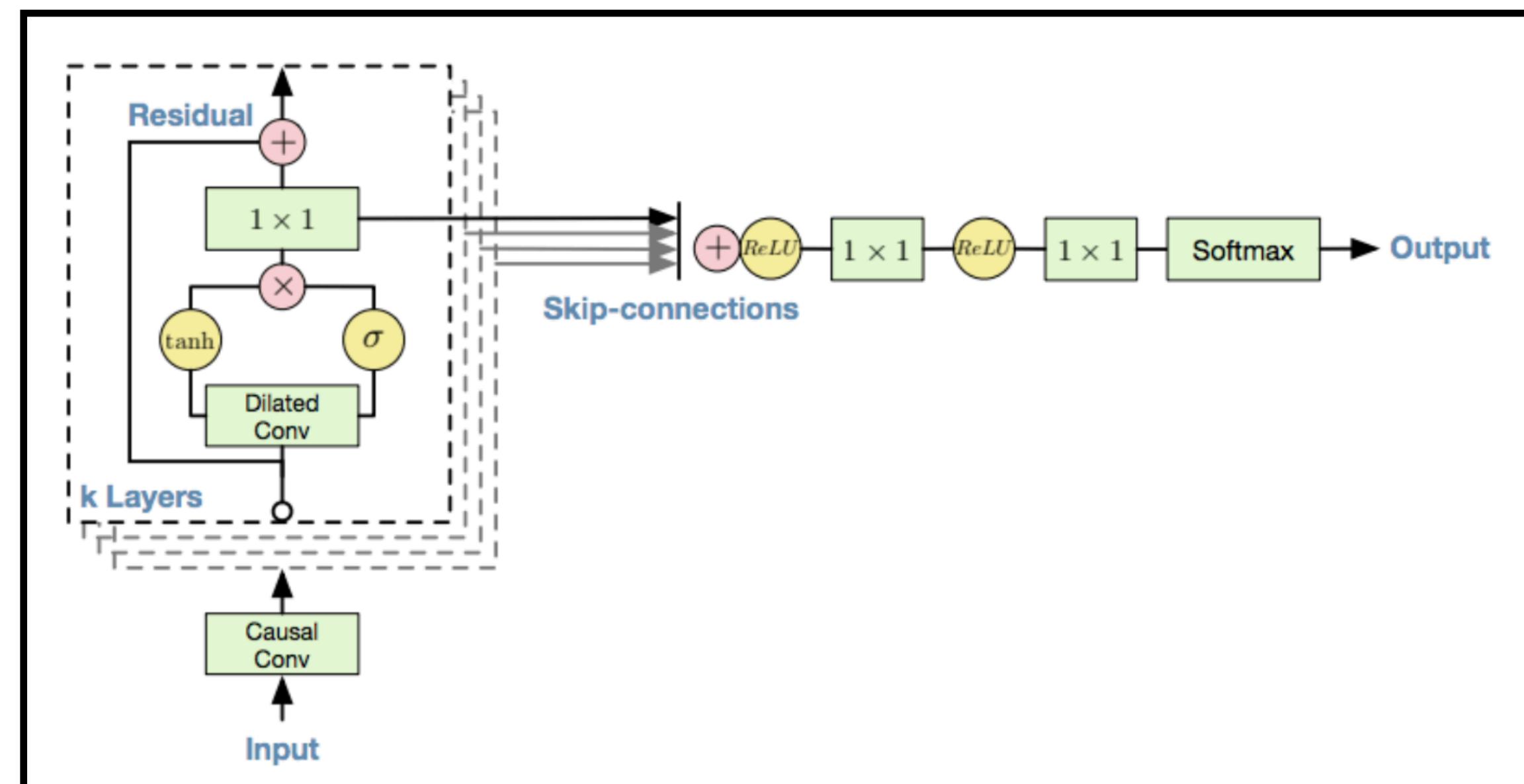
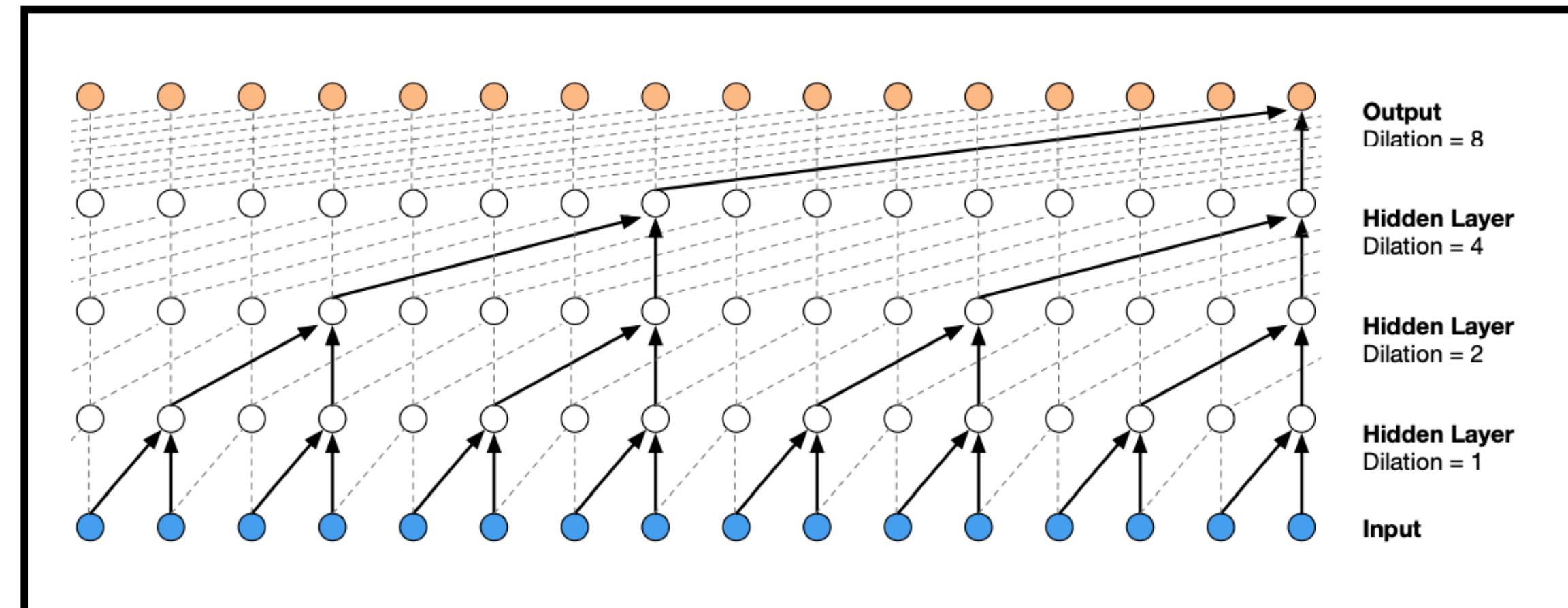
P = ?

eps = ?

$$3. \quad p(\mathbf{x} | \mathbf{h}) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}, \mathbf{h}).$$

$$4. \quad \mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k}^T \mathbf{h}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{h}).$$

\mathbf{h} - global conditioning:
speaker + phonemes



Wavenet

$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu |x_t|)}{\ln(1 + \mu)},$$

Inference:

