

WaveNet: A Generative Model for Raw Audio

Nadezda Kondrateva, 2023

План

Звук

Звук — колебательное движение частиц упругой среды, распространяющееся в виде волн в газообразной, жидкой или твёрдой средах

Звуковая волна - разрежение и сгущение звука

отвечает за
громкость

Амплитуда

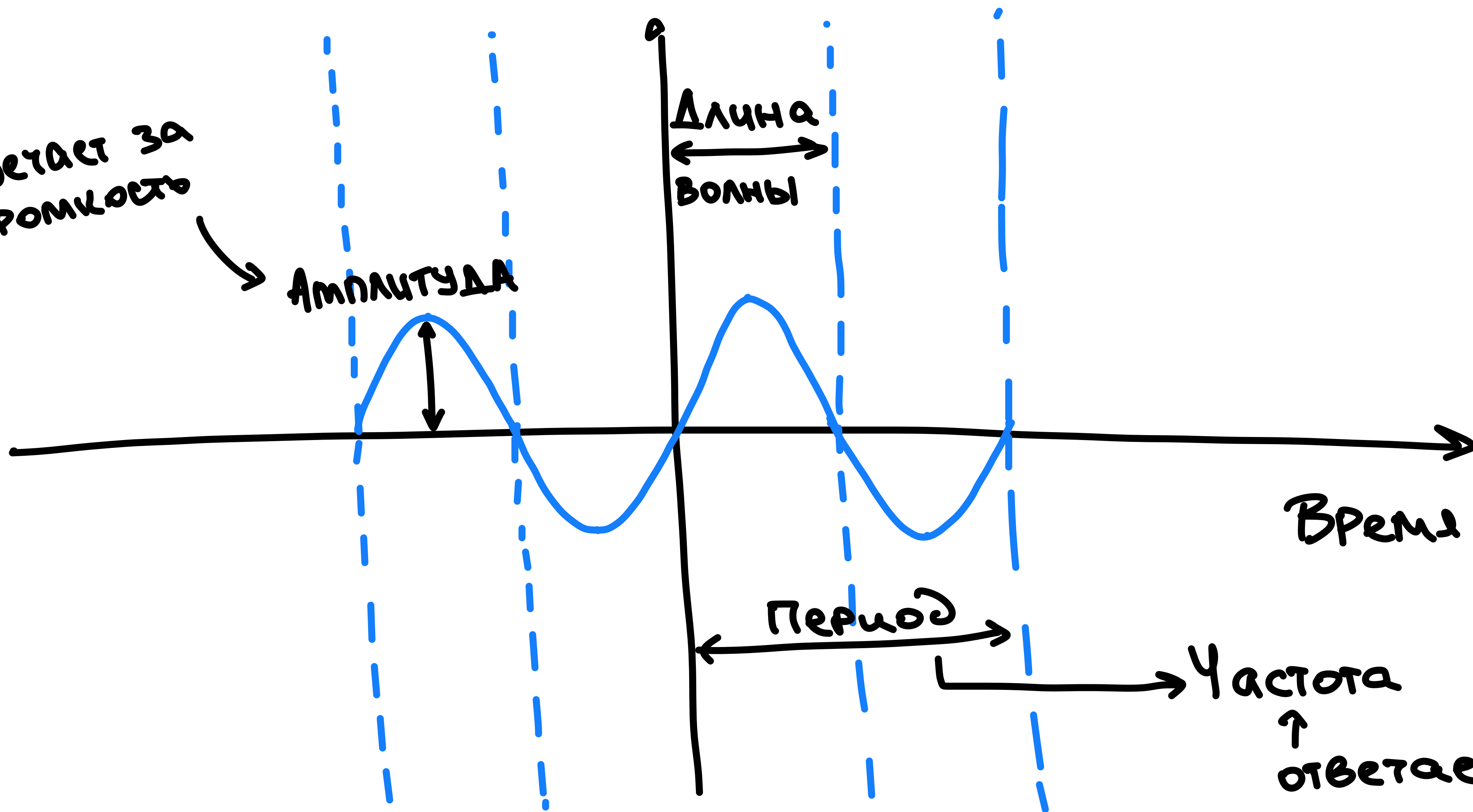
Длина
волны

Время

Период

Частота

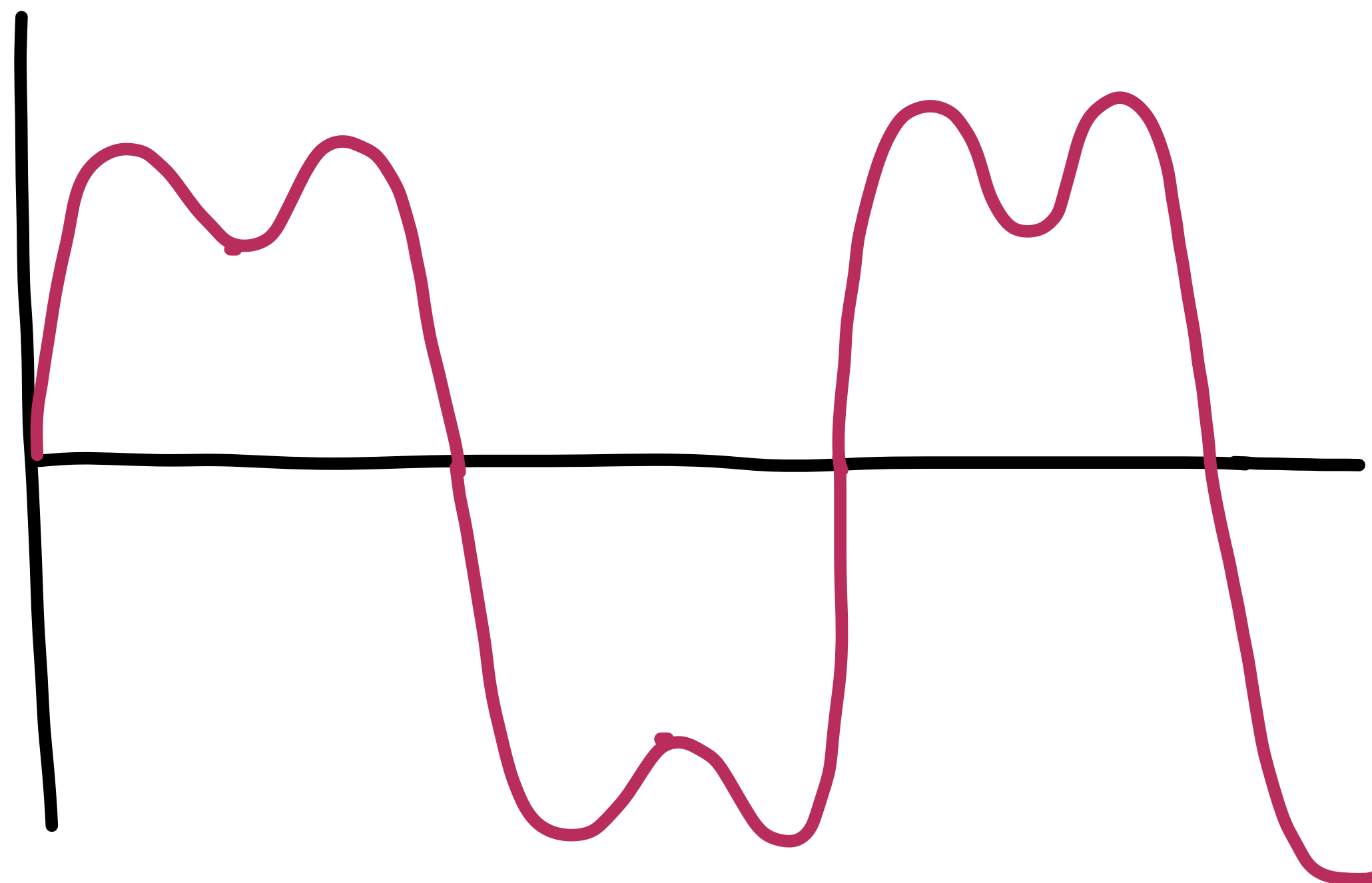
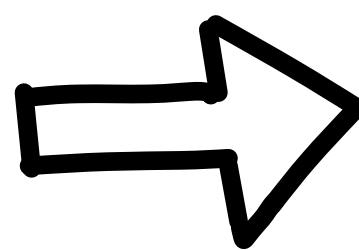
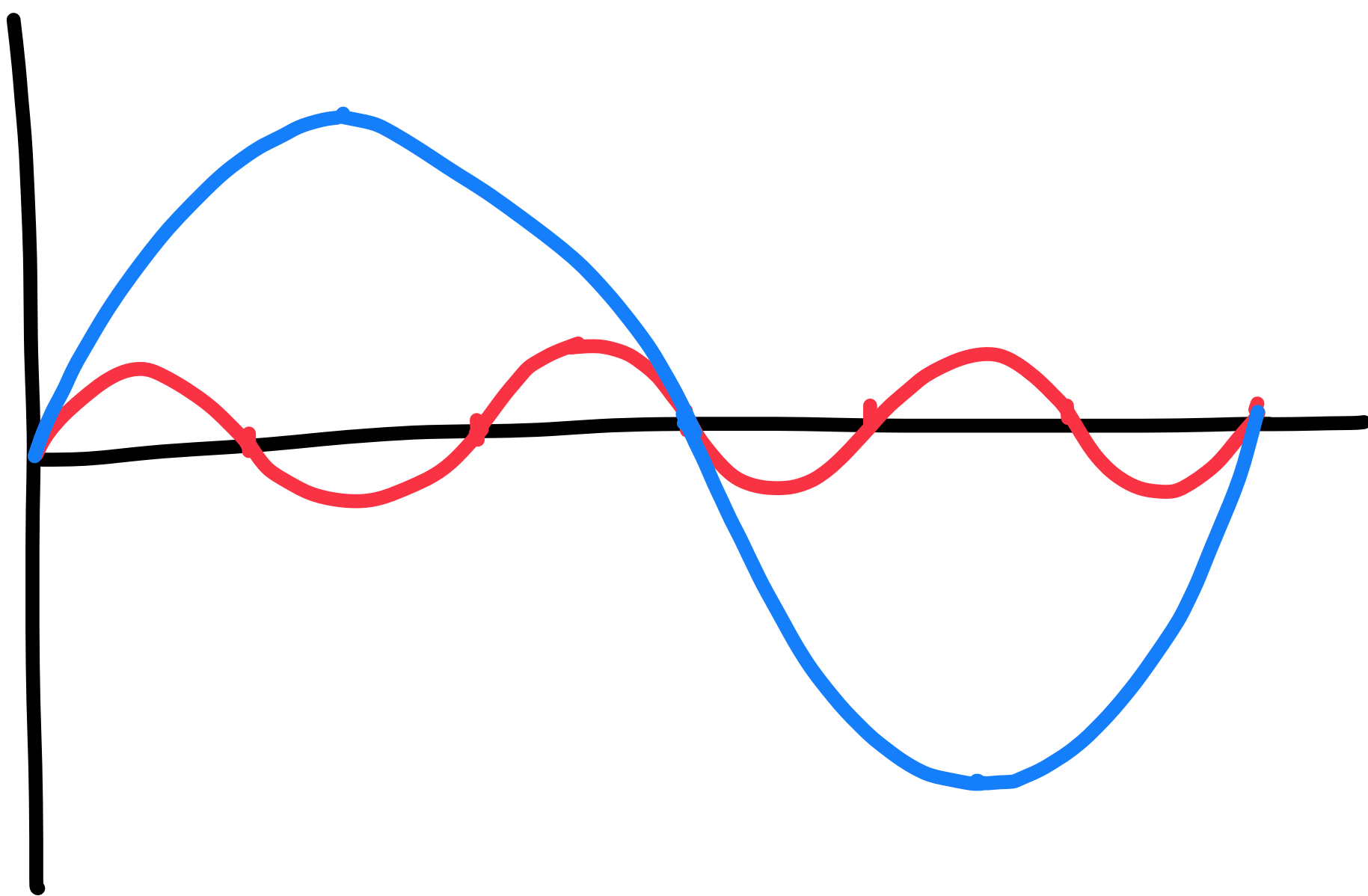
↑
отвечает
за тон





1 Second





vibrations in the air
are analog signal

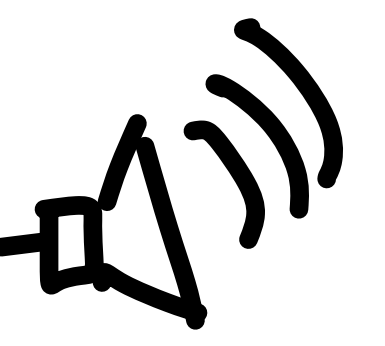


Inside a
computer

Digital to analog

DAC

Amplifier

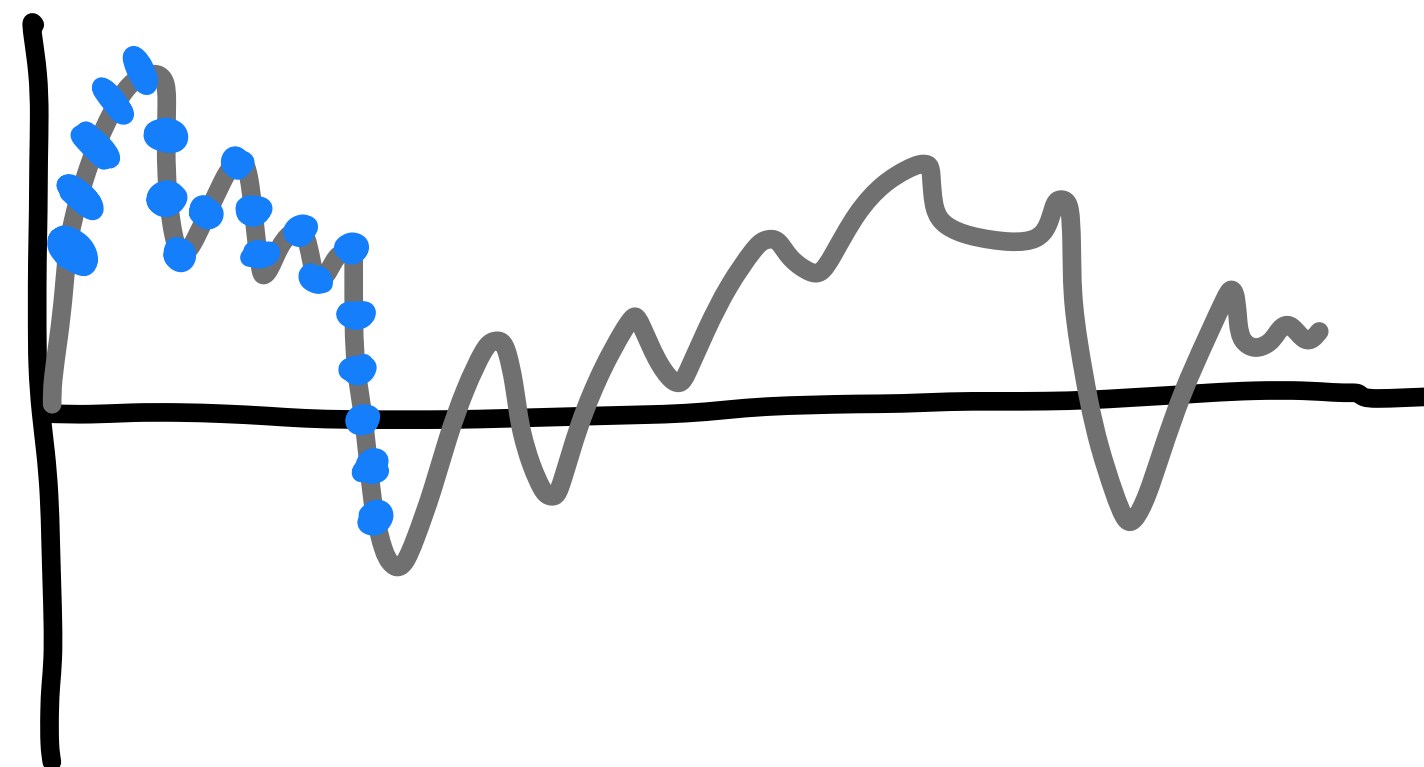


Binary signal

Storage
device

ADC

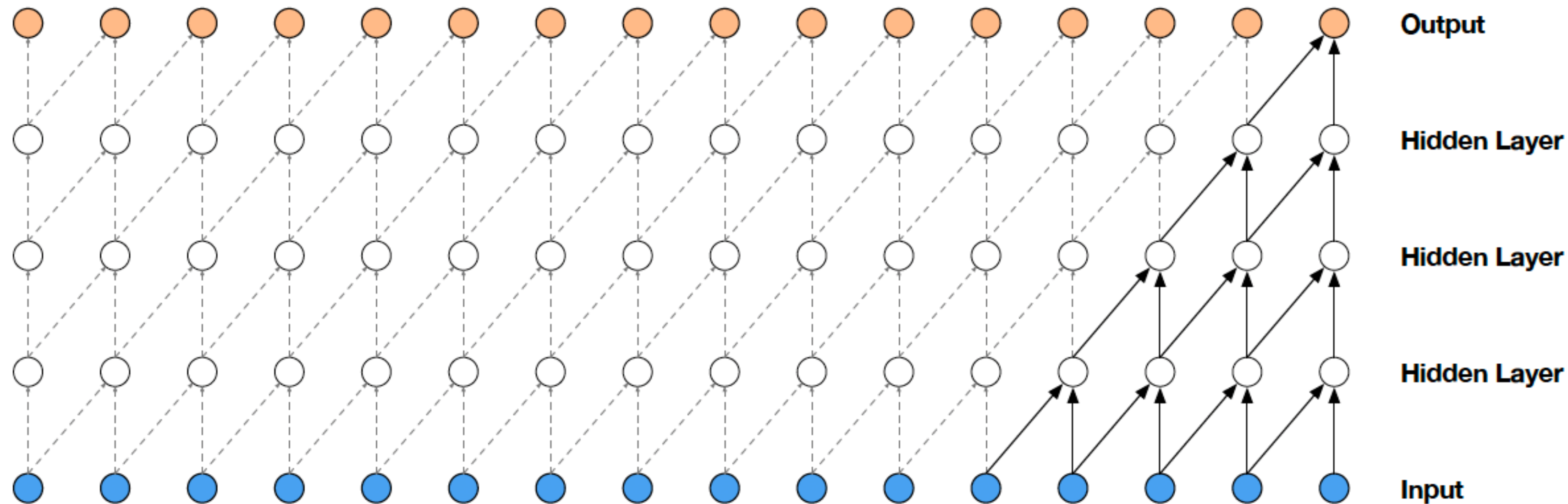
Analog to
digital converter
measures it all at
regular intervals



represent each
data point
(Delta-Sigma
modulation)
with 16 bits,
≥ 44000 rec. per sec.

WaveNet

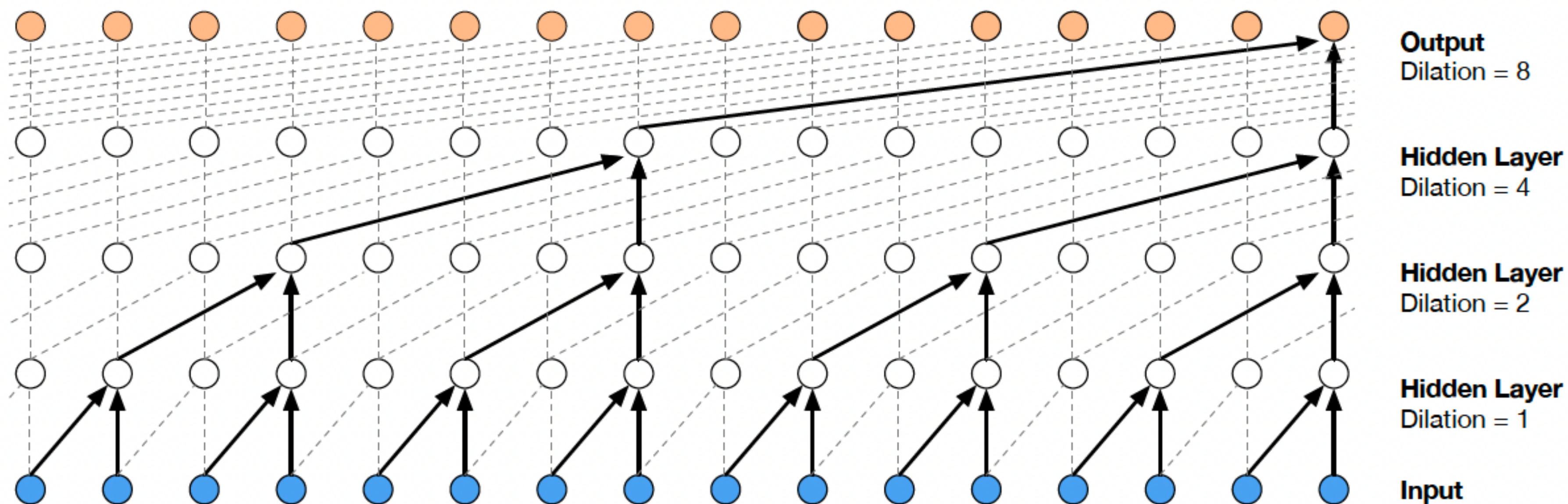
Causal Convolutions



Проблема: нужно много слоев или большая длина фильтра для достаточно большого рецептивного поля

Решение: Dilated Casual Convolutions

Dilated Casual Convolutions



При обучении шаг фильтра удваивается для каждого слоя

Softmax Distributions

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \text{for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

Проблема: 65,636 вероятностей для каждого промежутка

Решение: сжать значения амплитуд:

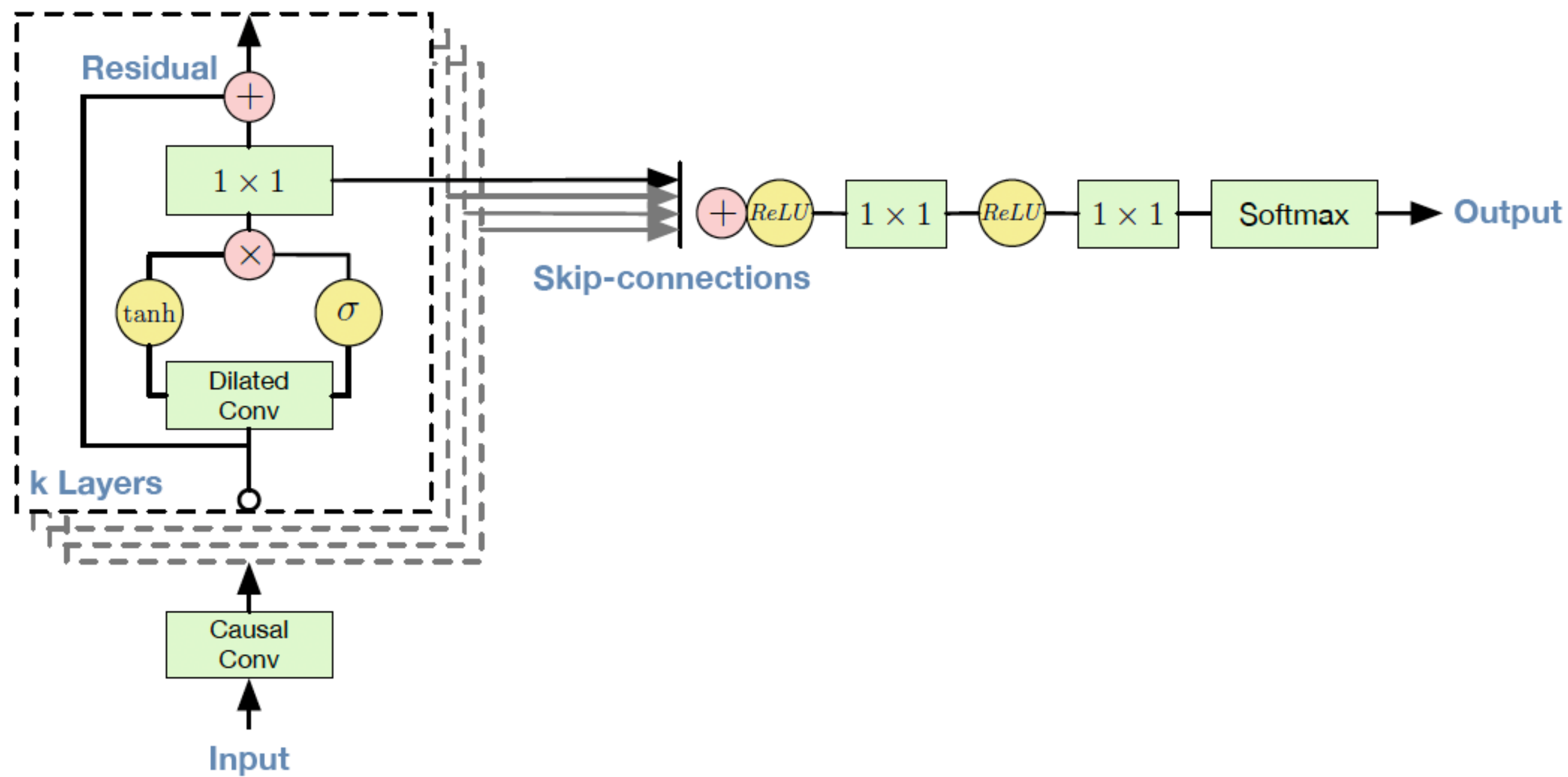
$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu |x_t|)}{\ln(1 + \mu)}, \quad -1 < x_t < 1, \quad \mu = 255$$

Gated Activation Units

↙ οδυσταστική φωνή

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x}) \odot \sigma(W_{g,k} * \mathbf{x})$$

Residual and Skip Connections



А если хотим аудио с определенными характеристиками?

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

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

Эксперименты

Multi-speaker Speech Generation

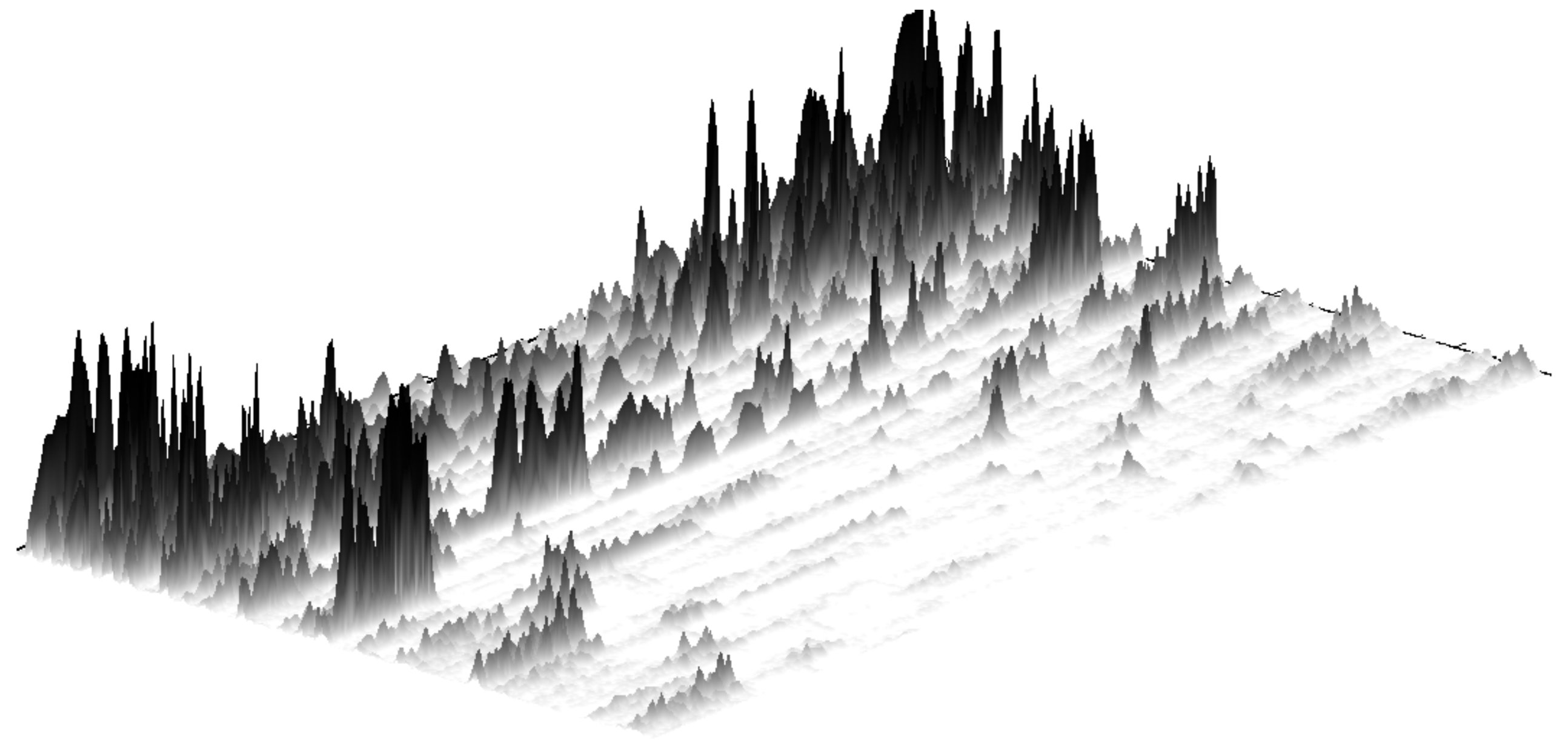
- **Параметризация: ID спикера**
- **Генерирует звуки, похожие на человеческие, дополнительные характеристики воспроизводит**
- **Речь несвязная из-за размера рецептивного поля**

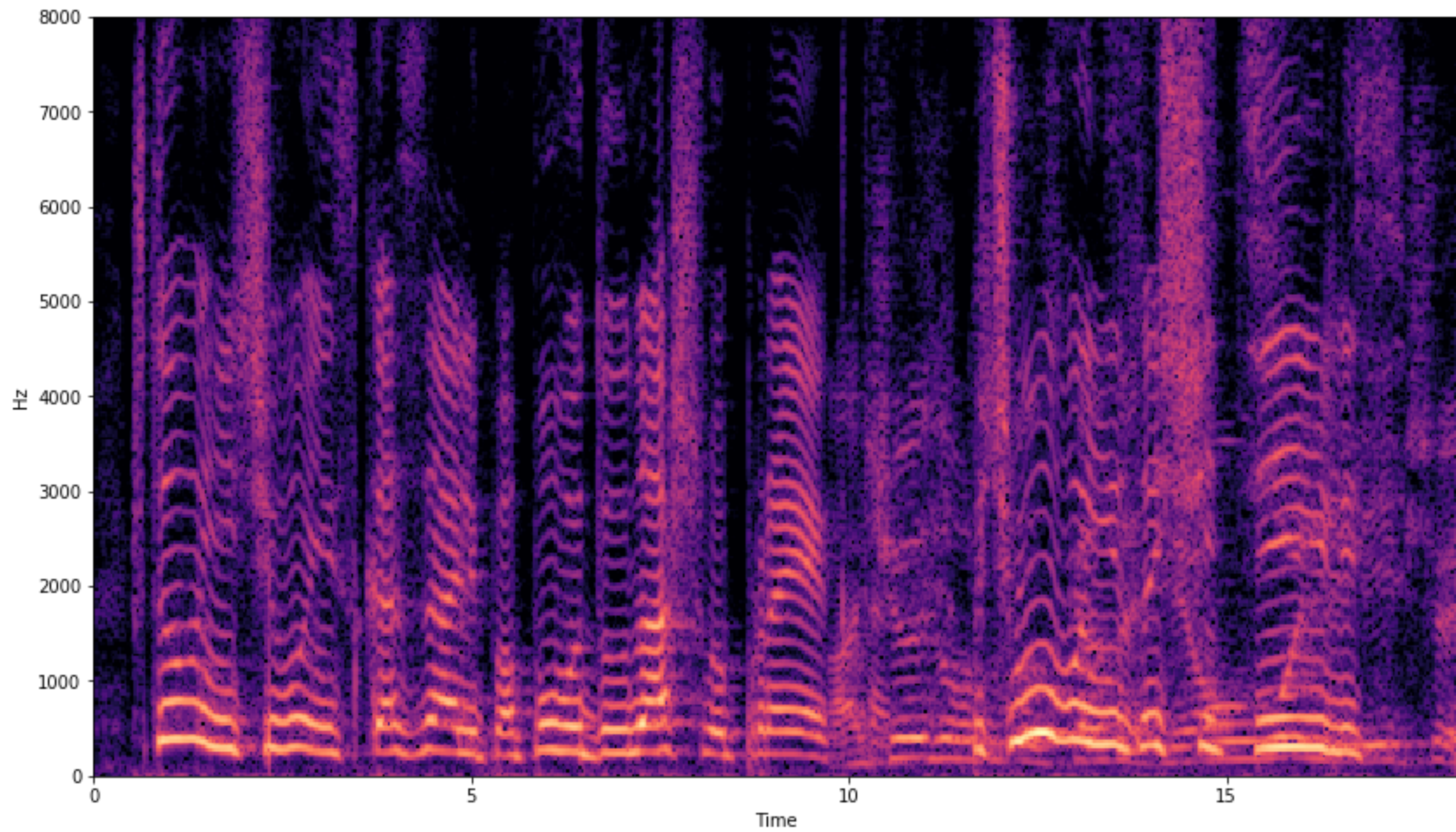
Спектрограмма и мел-спектрограмма

**Спектрограмма показывает
зависимость амплитуды от времени и
частоты**

Построение

- Применяется преобразование Фурье для каждого фрагмента звукового сигнала
- По-сути, берет некоторый входной сигнал и выводит соответствующие веса всех частот, составляющих сигнал



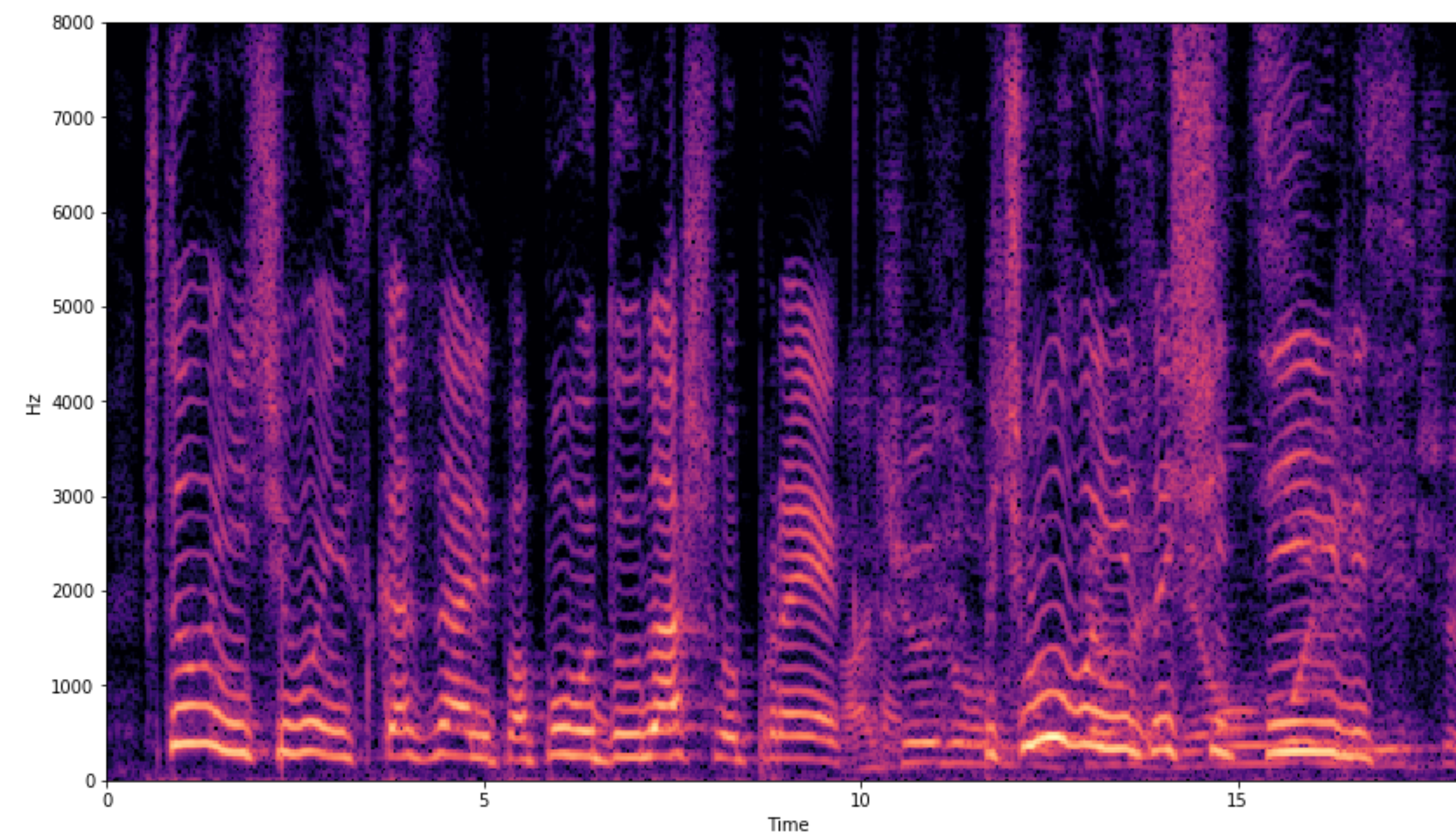


Проблема!

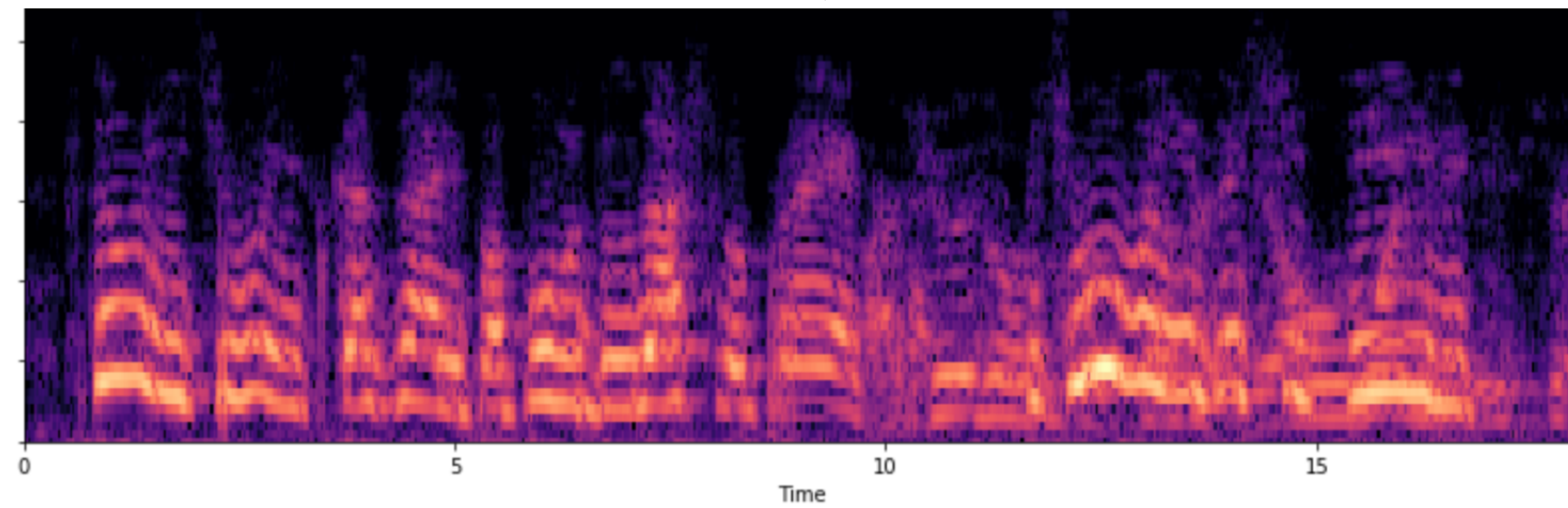
**Разные частоты воспринимаются человеком по-разному
Хотим, чтобы одинаковые изменения по шкале одинаково
воспринимались для человека**

Мел-спектрограмма

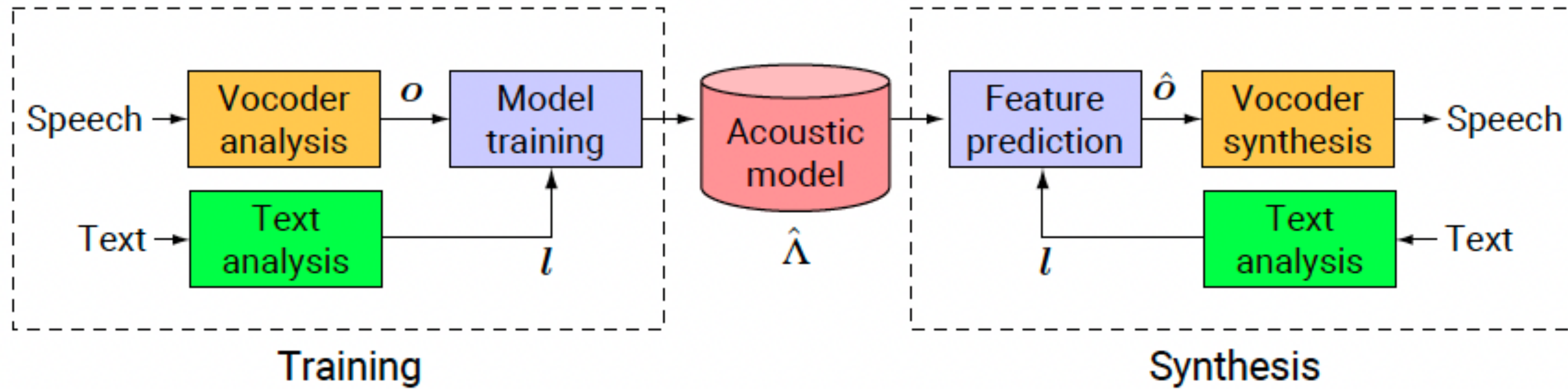
Значения из низких частот спектрограммы остаются практически неизменными на мел-спектре, а в высоких частотах происходит усреднение значений из более широкого диапазона



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



Text-To-Speech



Music

Speech recognition

Итоги

Конец?