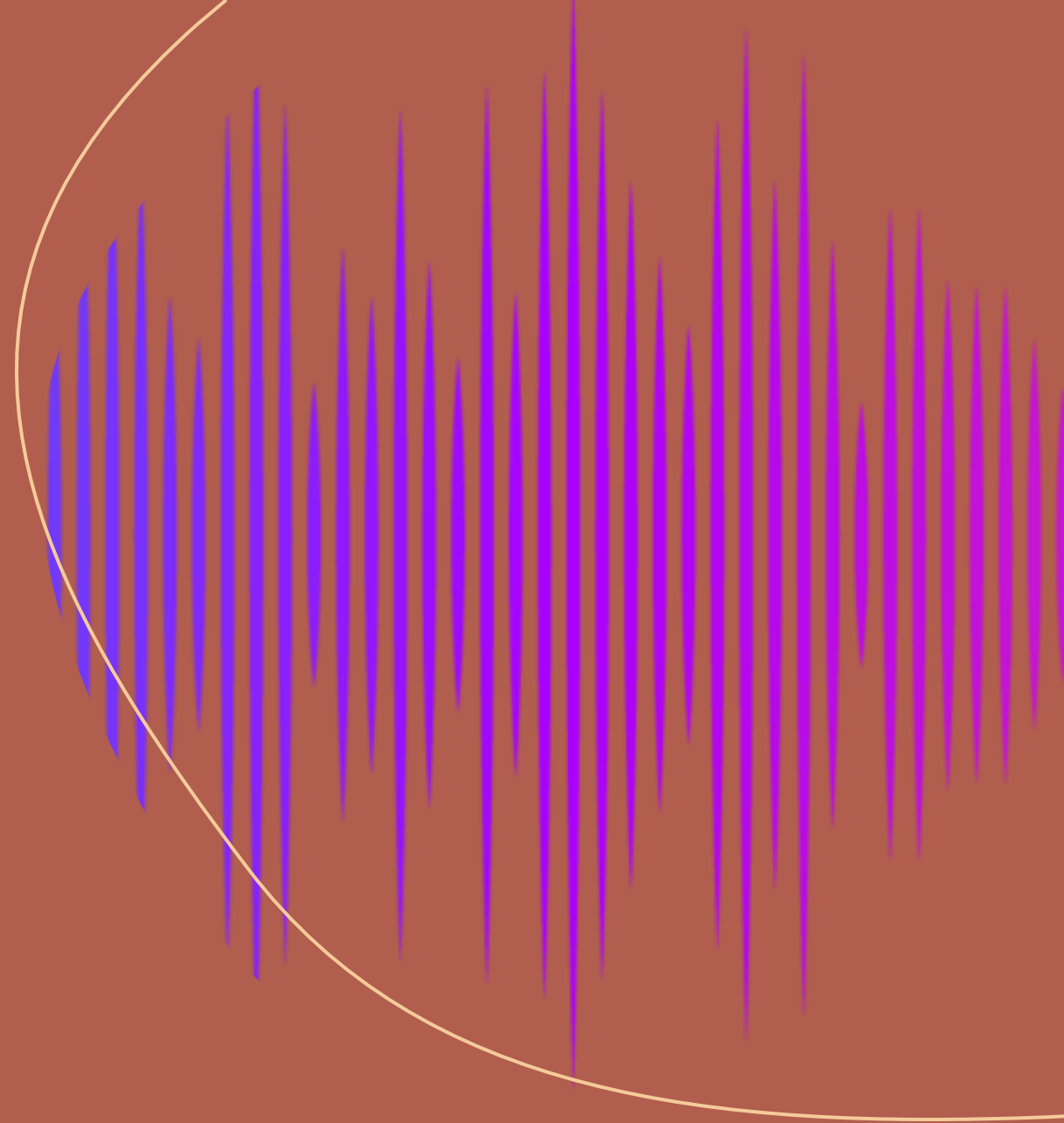




Natural Language Processing

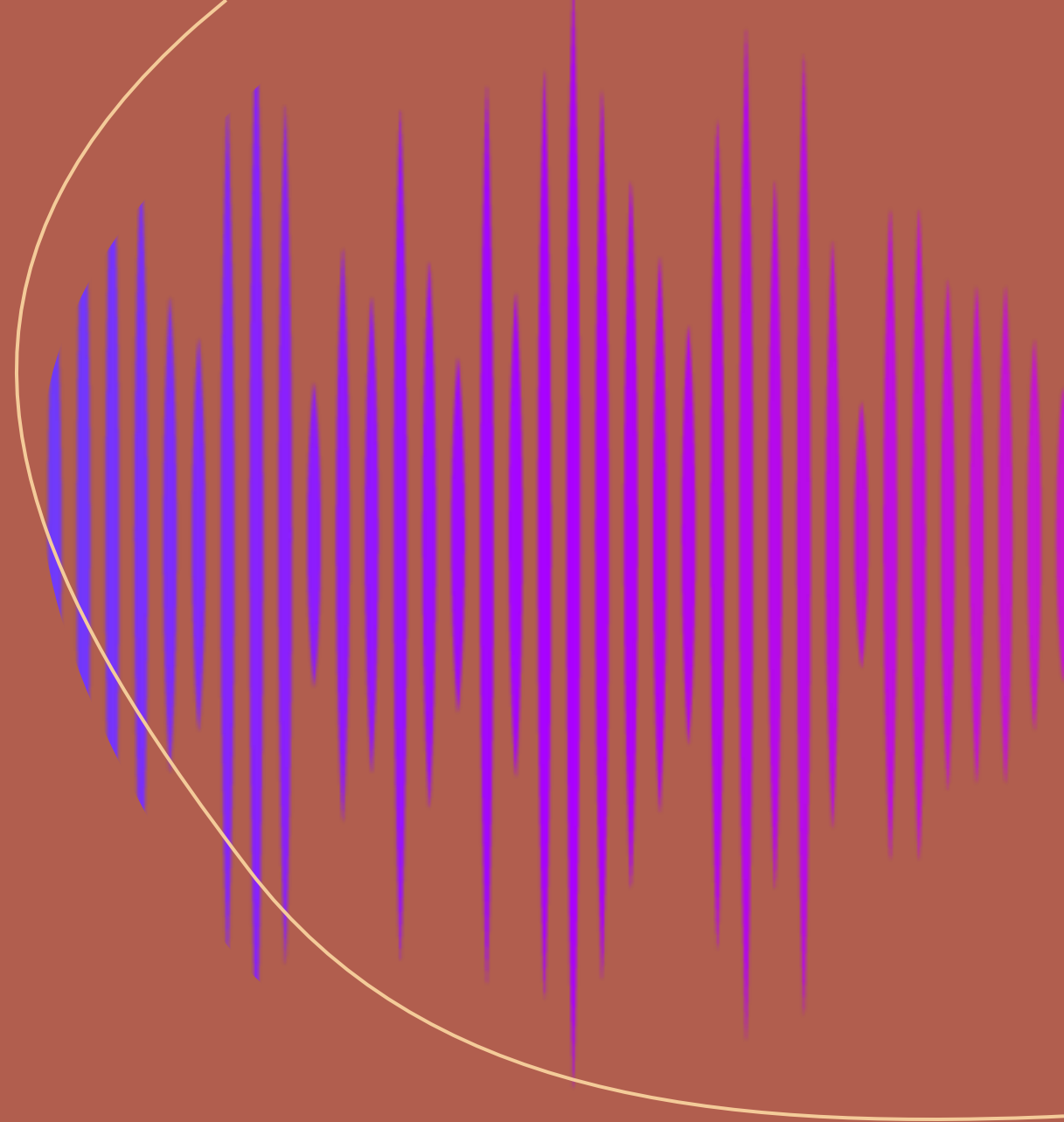
Яковенко Ольга

Задачи по работе с аудио



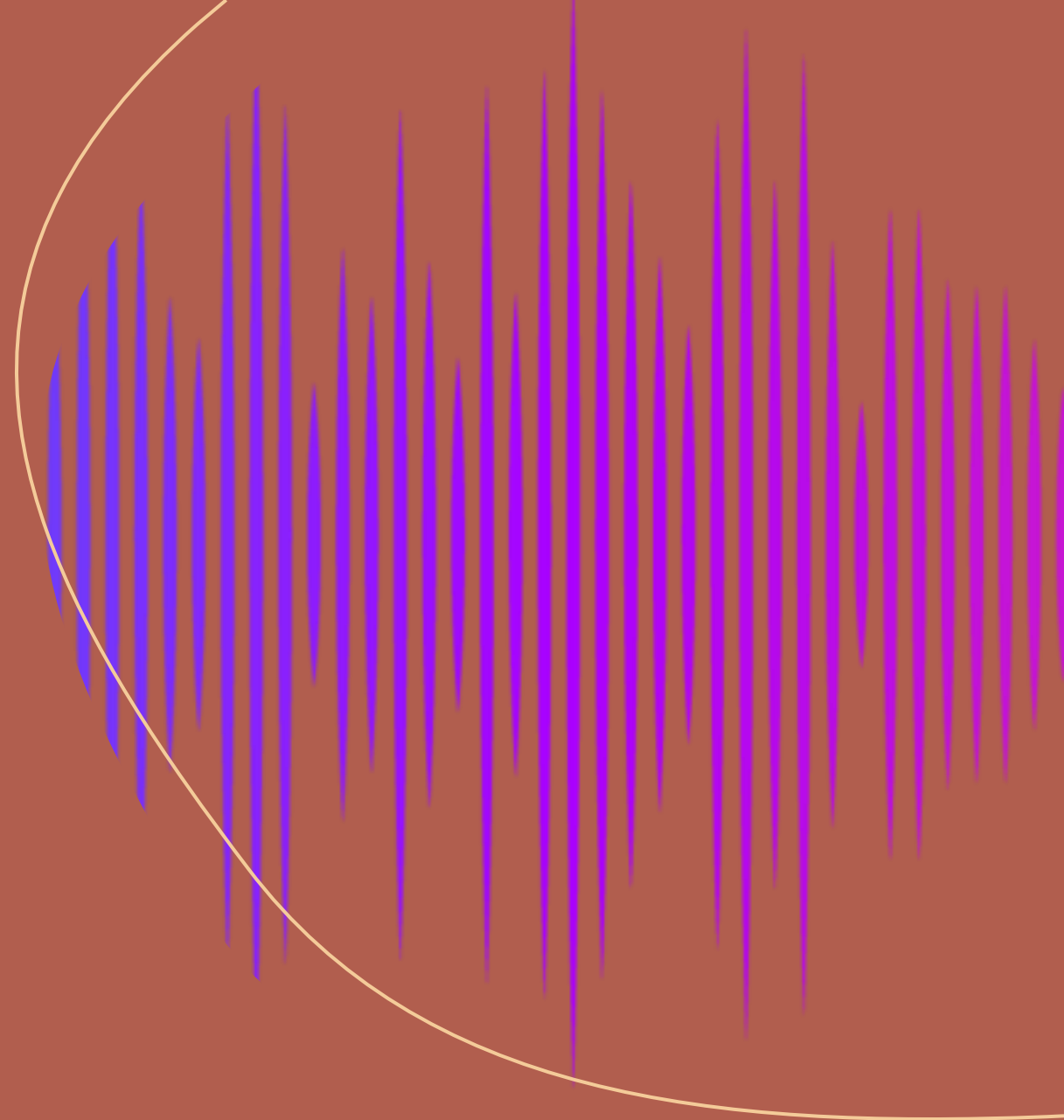
Задачи по работе с аудио

- Распознавание речи



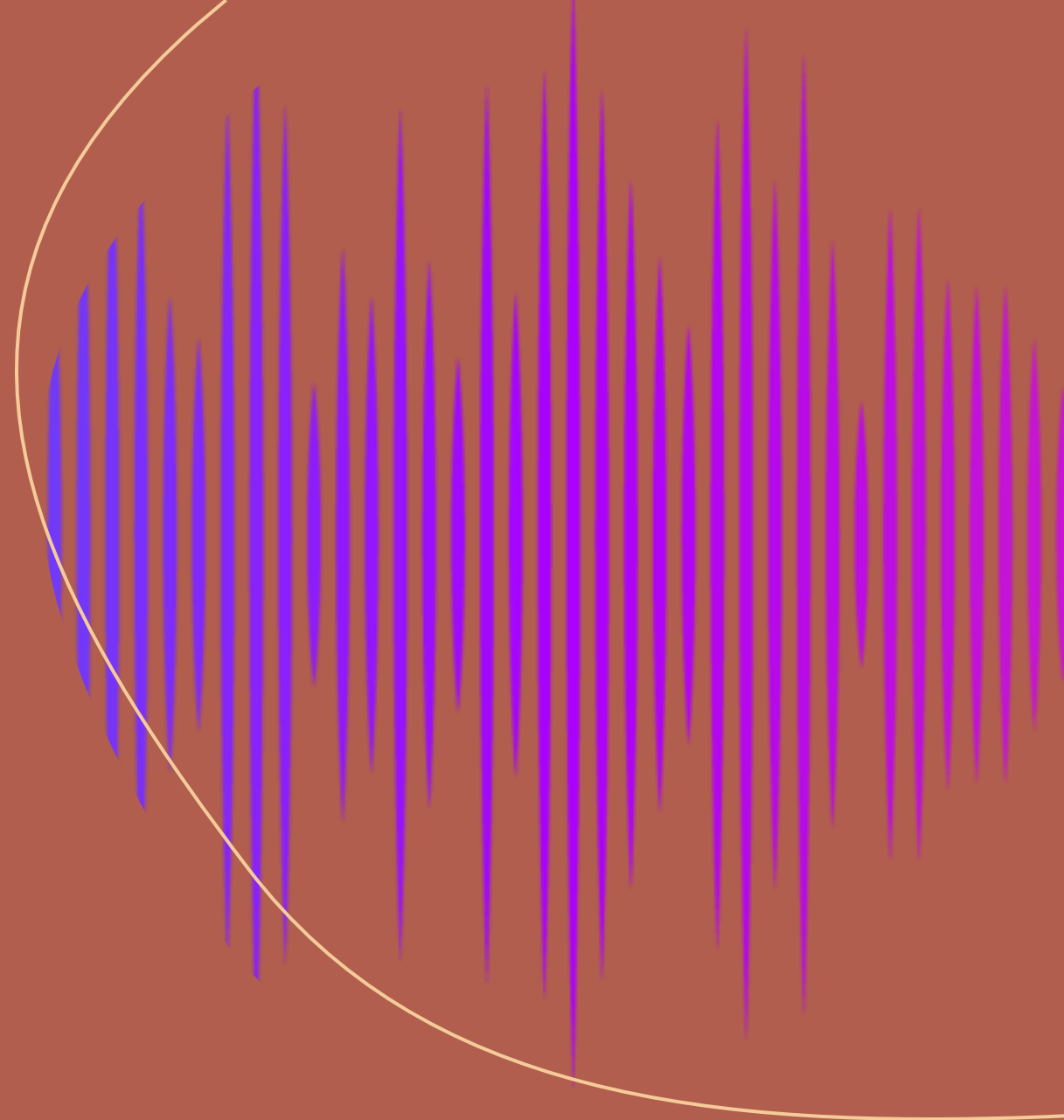
Задачи по работе с аудио

- Распознавание речи
- Синтез речи



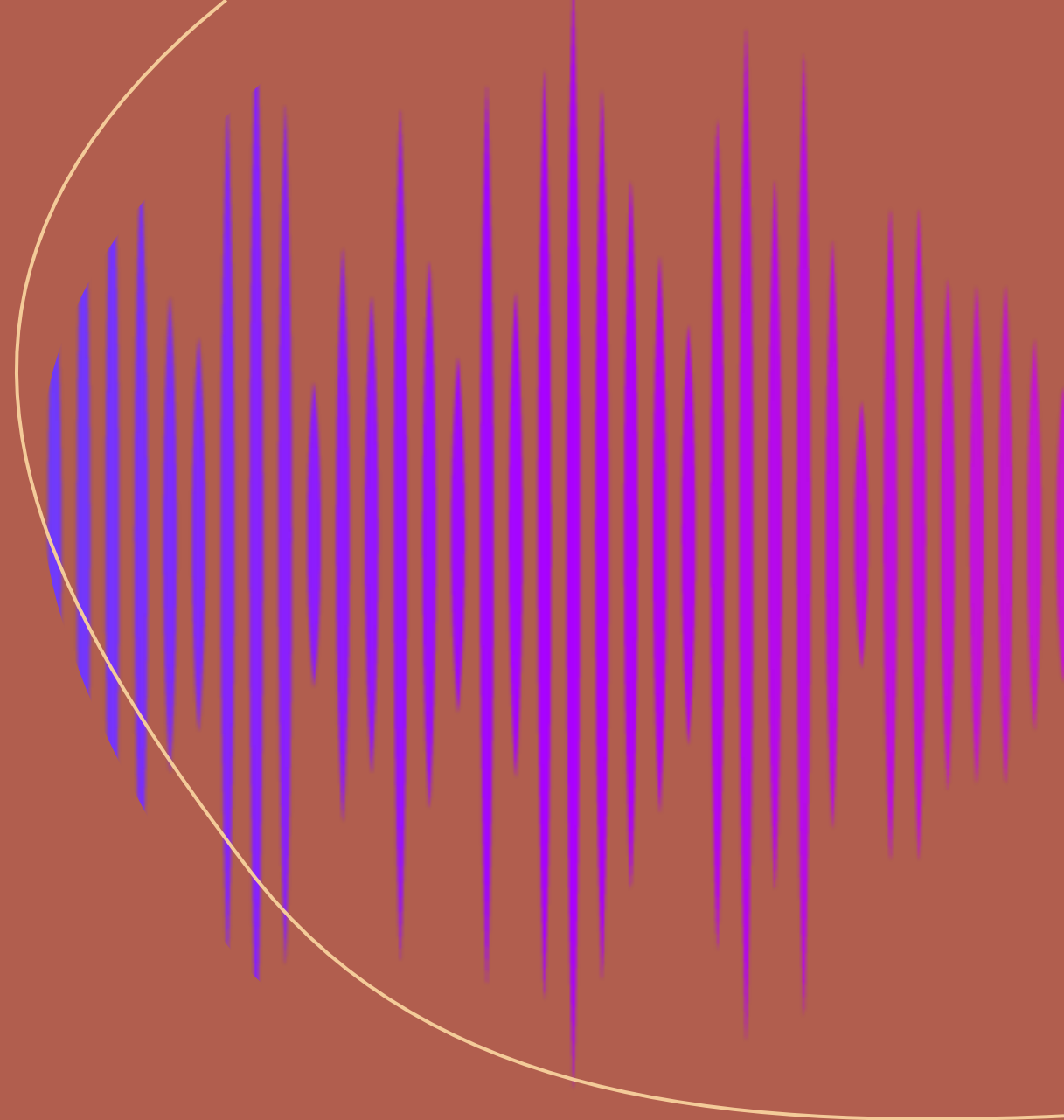
Задачи по работе с аудио

- Распознавание речи
- Синтез речи
- Верификация/идентификация диктора



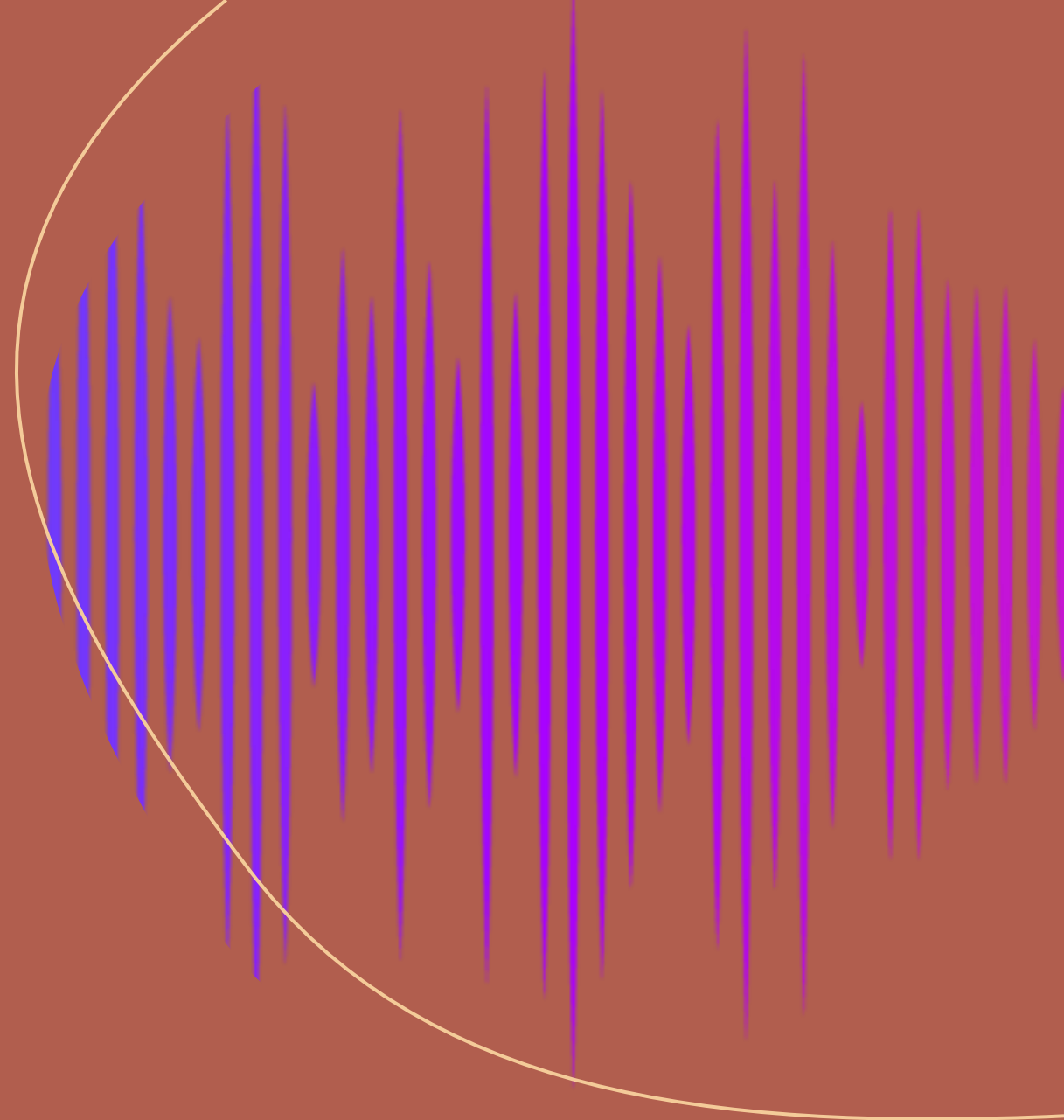
Задачи по работе с аудио

- Распознавание речи
- Синтез речи
- Верификация/идентификация диктора
- Поиск похожих композиций



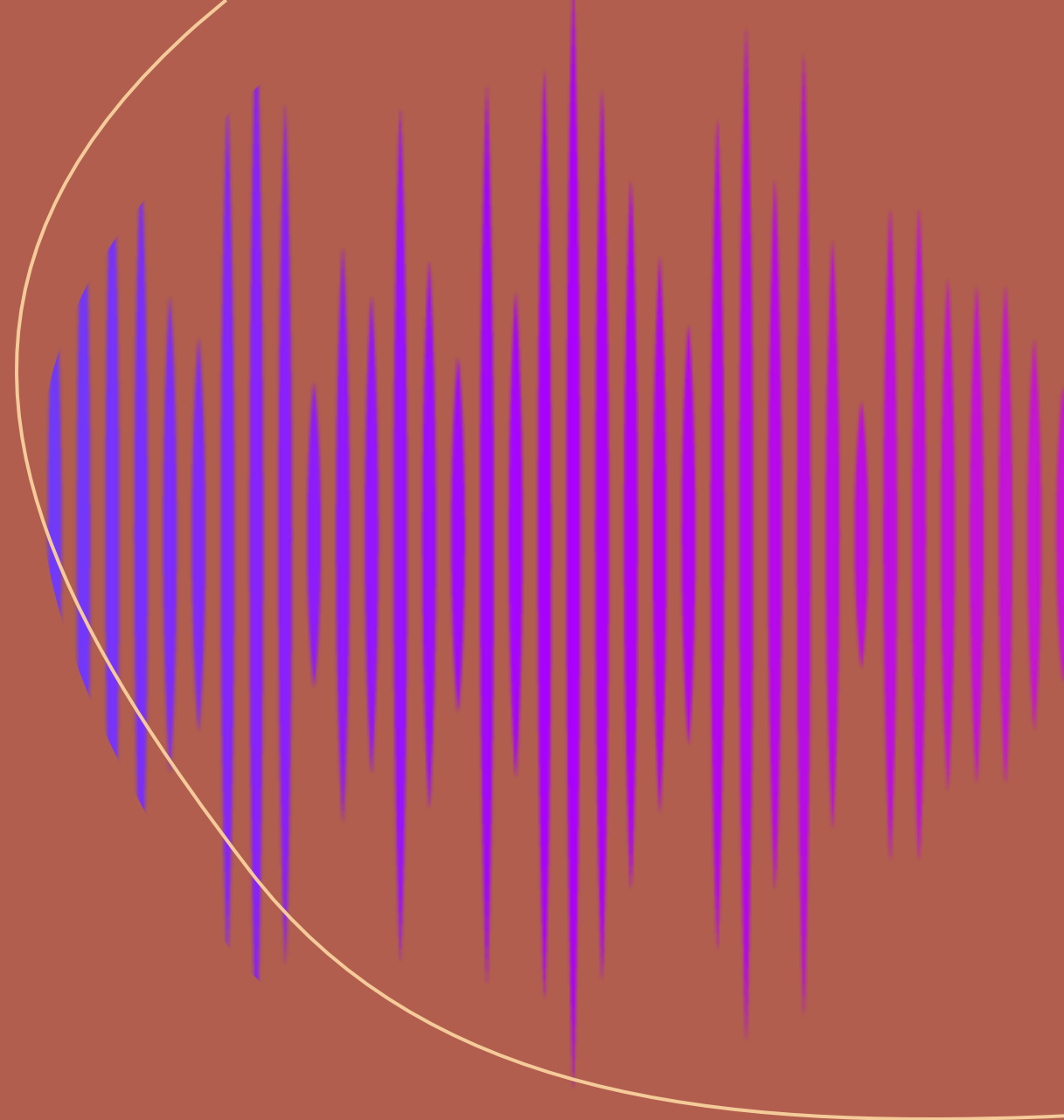
Задачи по работе с аудио

- Распознавание речи
- Синтез речи
- Верификация/идентификация диктора
- Поиск похожих композиций
- Шумоподавление



Задачи по работе с аудио

- Распознавание речи
- Синтез речи
- Верификация/идентификация диктора
- Поиск похожих композиций
- Шумоподавление



Распознавание речи



hello

Распознавание речи

Автоматическое обнаружение в
аудио произносимого человеком
текста



Распознавание речи

```
graph TD; A[Распознавание речи] --> B[Компонентный подход]; A --> C[End-to-End подход];
```

Компонентный
подход

End-to-End
подход

Компонентный подход

Входной аудио
поток

WAV файл

Извлечение
признаков

Мел-частотные
кепстральные
коэффициенты (MFCC)

I-vectors

Кепстральная
нормализация (CMVN)

Акустический
блок

Оконное
распознавание звуков

Определение
наиболее вероятной
цепочки звуков

Лингвистический
блок

Языковая модель

Выходной текст

Строка

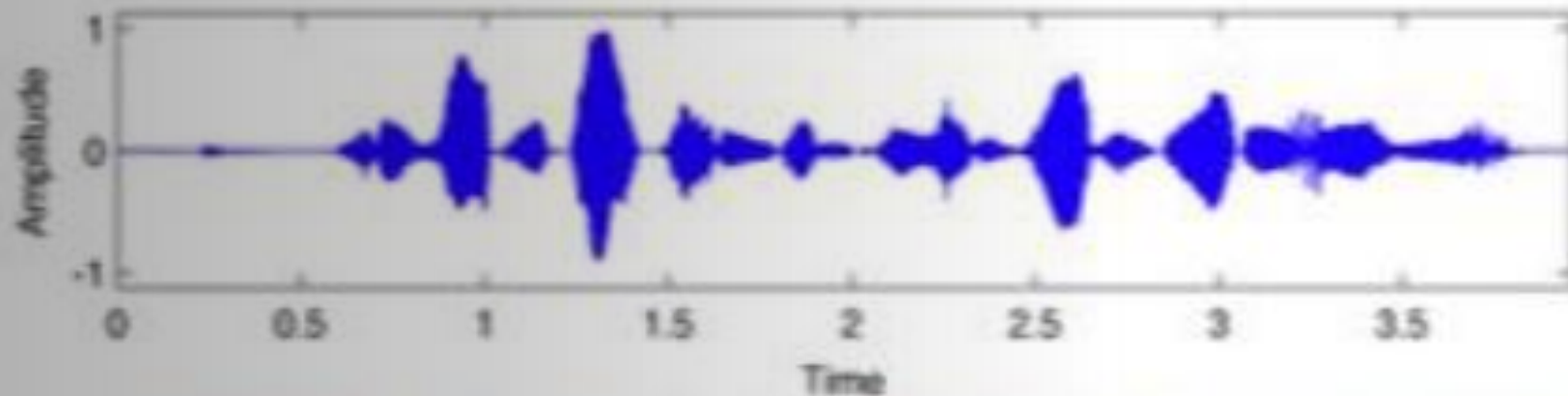
Входной поток аудио

- WAV формат:
- Частота дискретизации (8 кГц, 16 кГц, 44 кГц)
- Количество каналов
- Битовая глубина

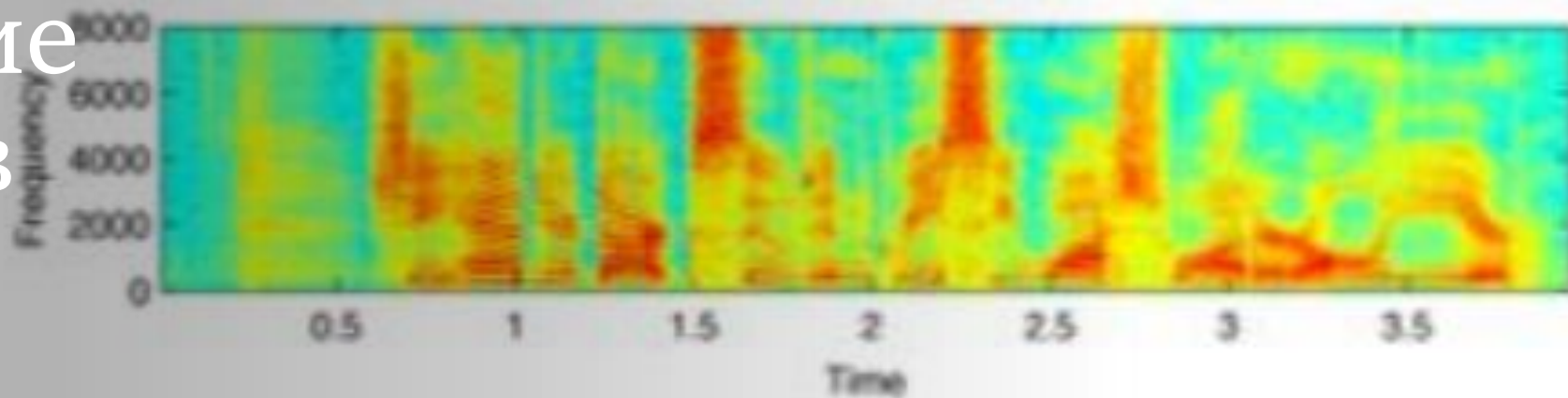
```
-bash-4.2$ soxi audio.wav
```

```
Input File      : 'audio.wav'  
Channels        : 1  
Sample Rate     : 8000  
Precision       : 16-bit  
Duration        : 01:31:33.97 = 43951752 samples ~ 412048 CDDA sectors  
File Size       : 87.9M  
Bit Rate        : 128k  
Sample Encoding: 16-bit Signed Integer PCM
```

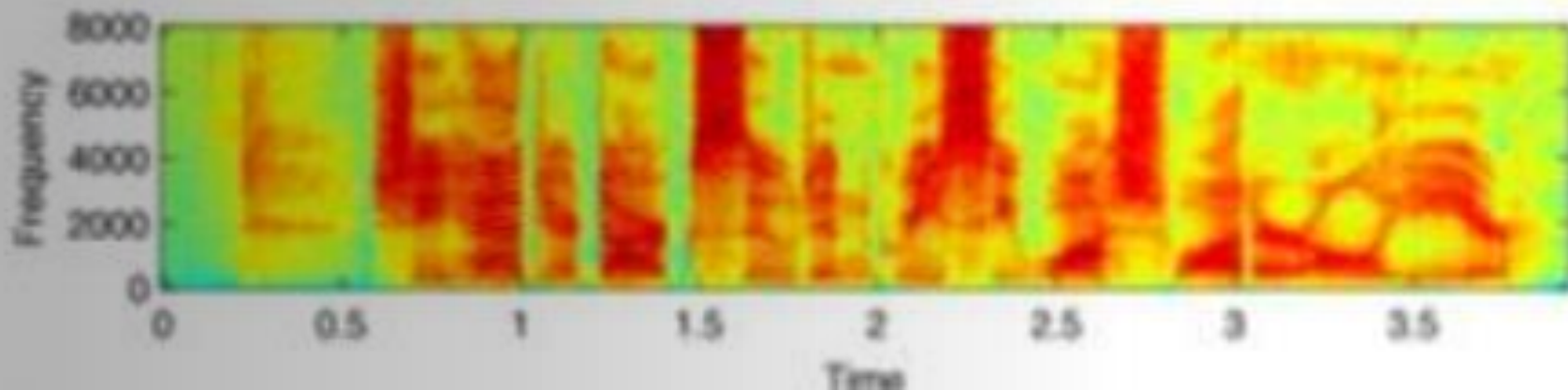
**Time Domain
Waveform**

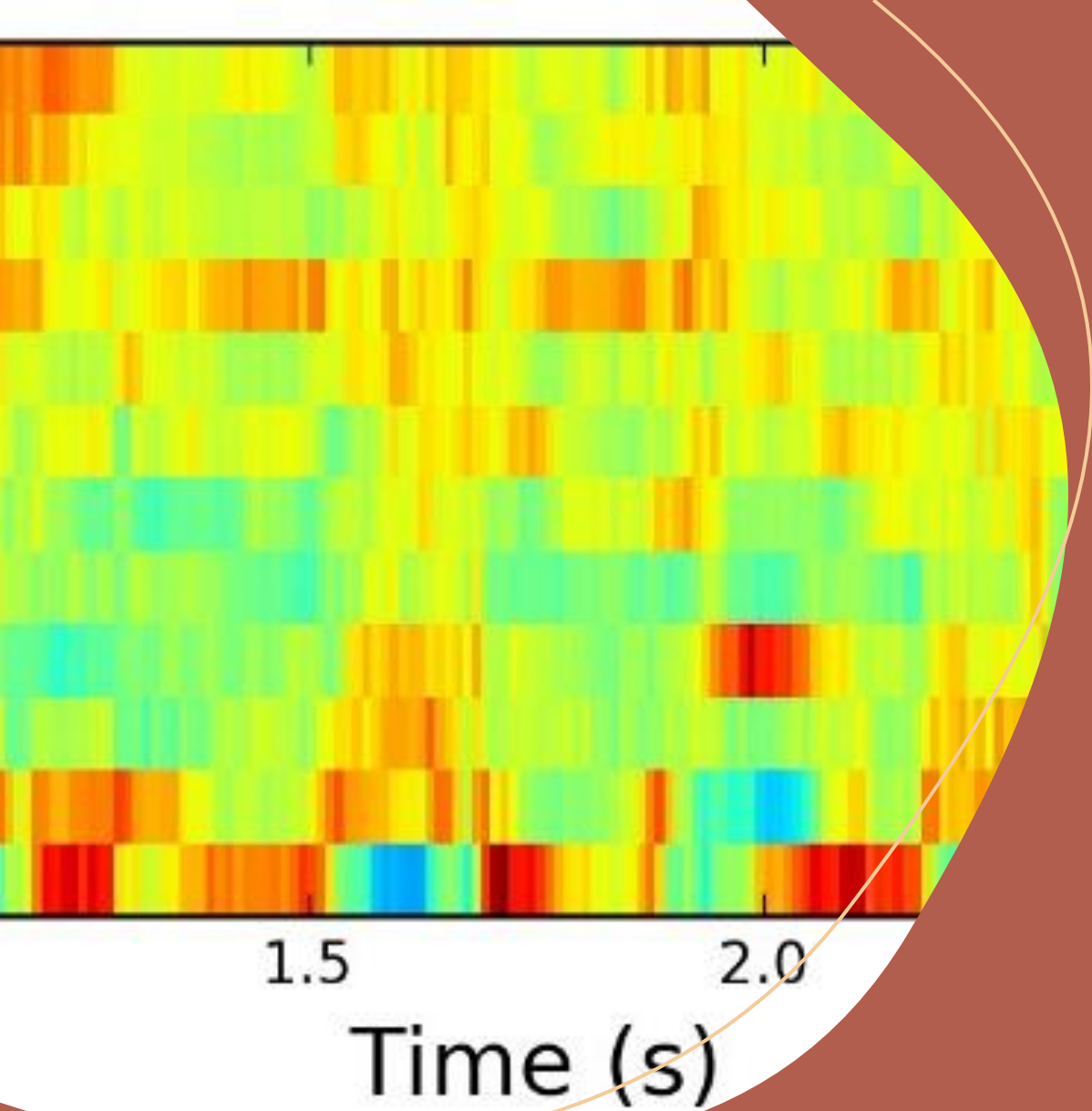


**Извлечение
признаков**
Spectrogram



**MFCC
Spectrogram**



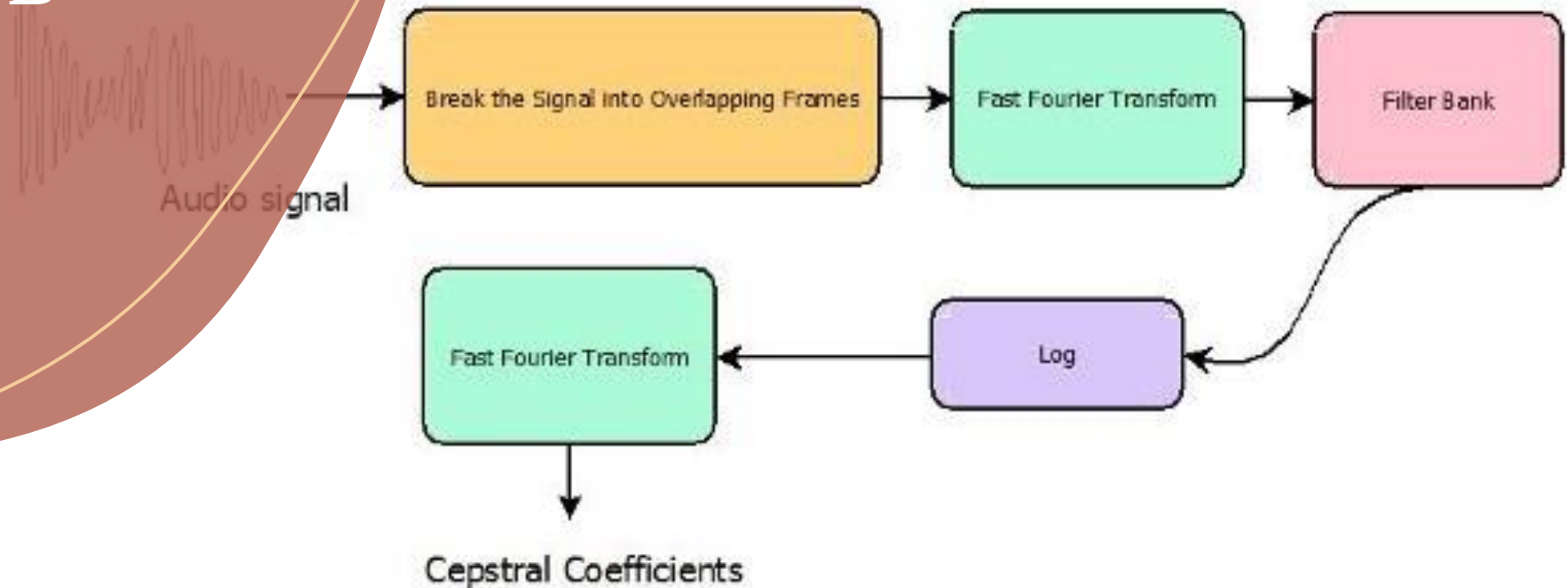
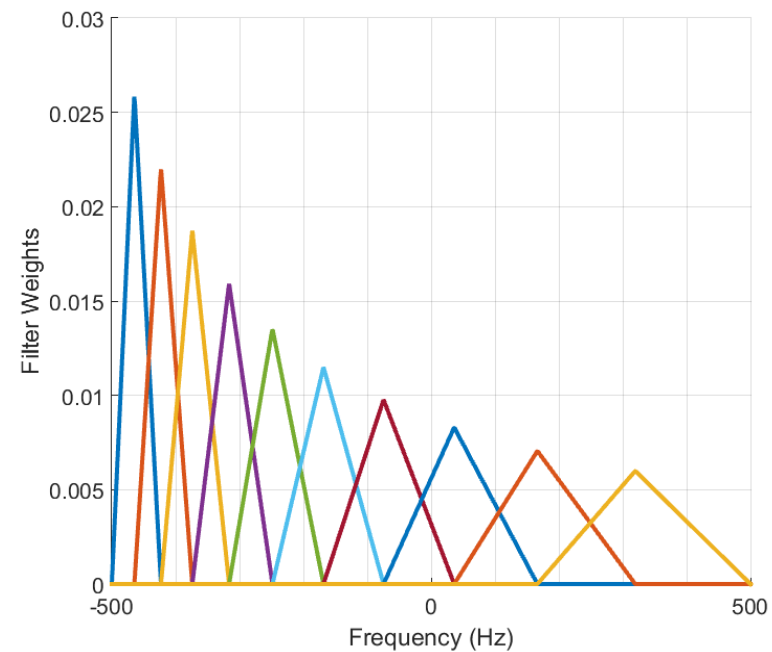


Извлечение признаков

- Мел-частотные кепстральные коэффициенты
- Mel-frequency cepstral coefficients
- MFCC

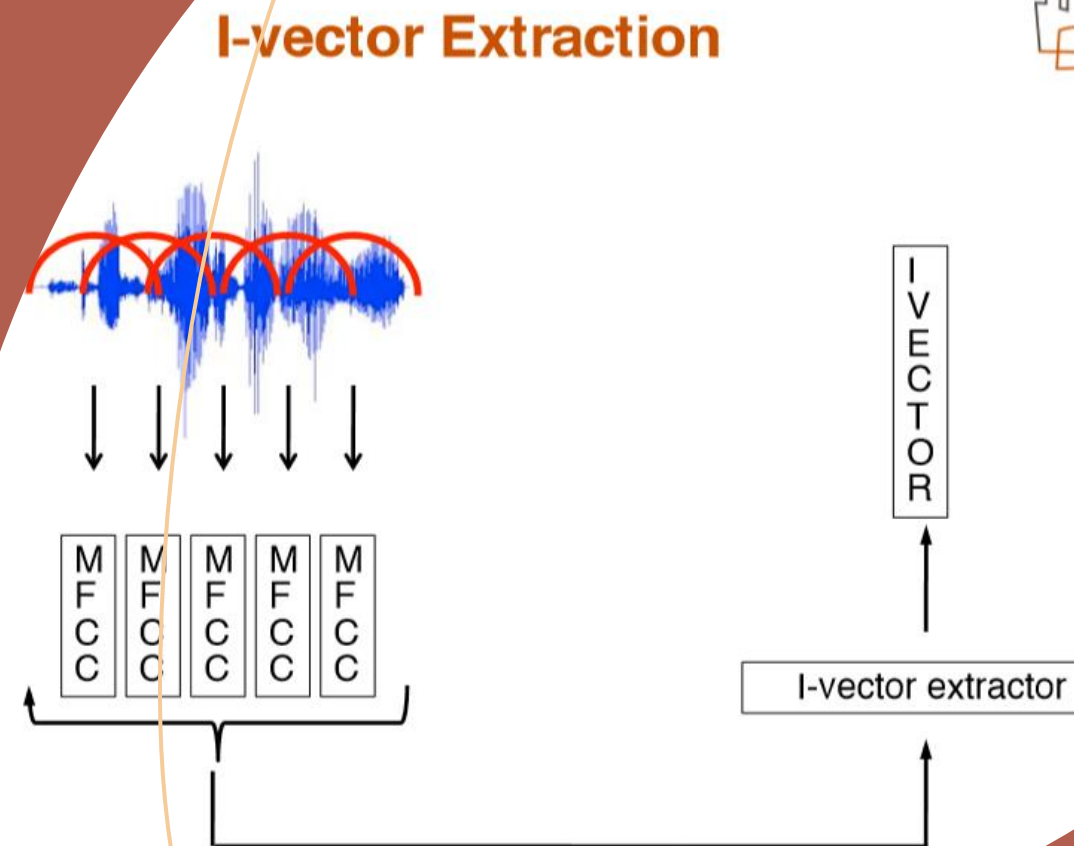
Извлечение признаков

MFCC



Извлечение признаков

- Identity vectors
- I-vectors



Language Systems, MIT CSAIL

Inter-session comparison

Акустический блок

Оконное распознавание

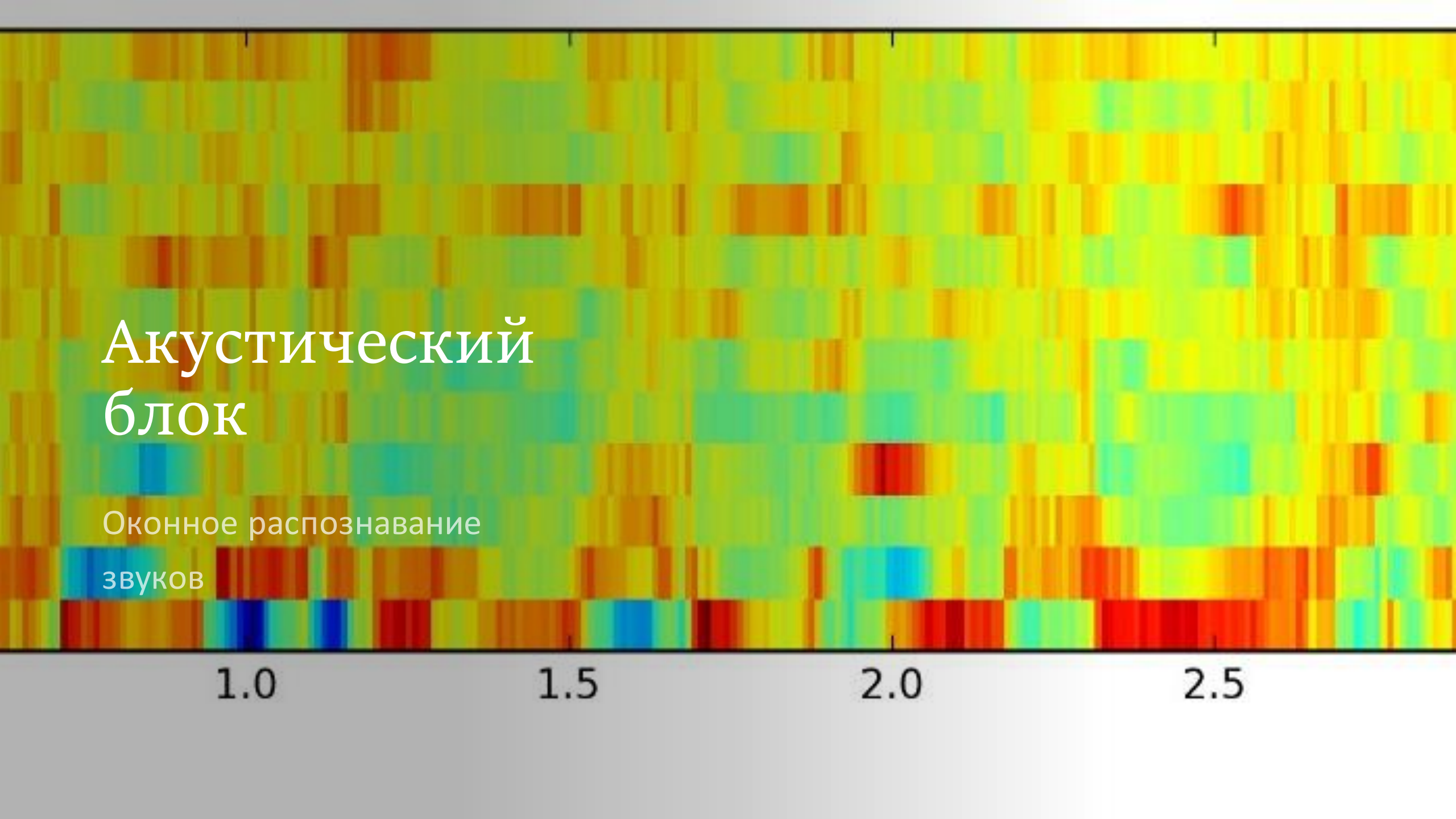
звуков

1.0

1.5

2.0

2.5



Акустический блок

Оконное распознавание
звуков

1.0

1.5

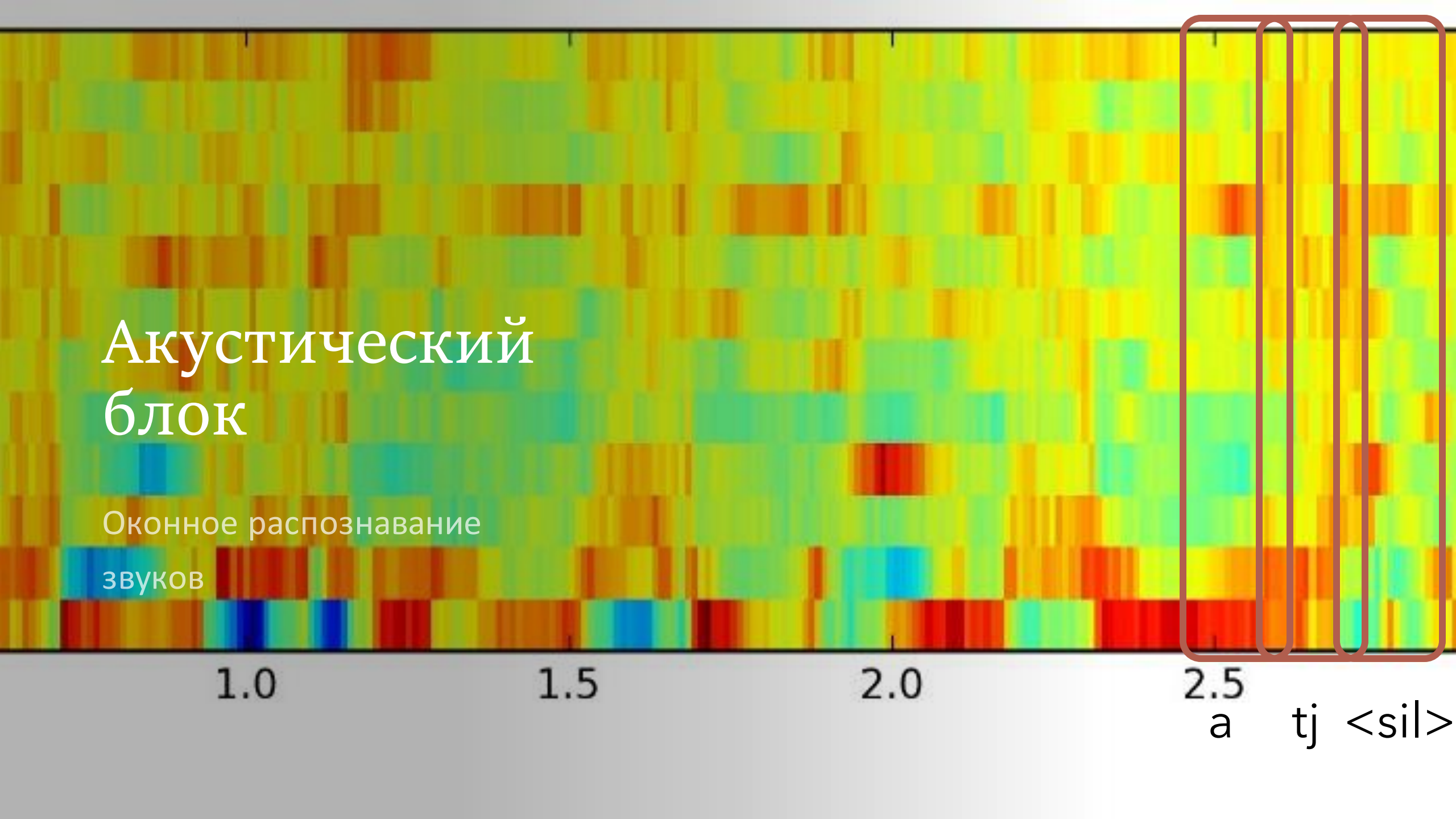
2.0

2.5

a

tj

<sil>



Акустический блок

Акустический блок

Оконное распознавание звуков

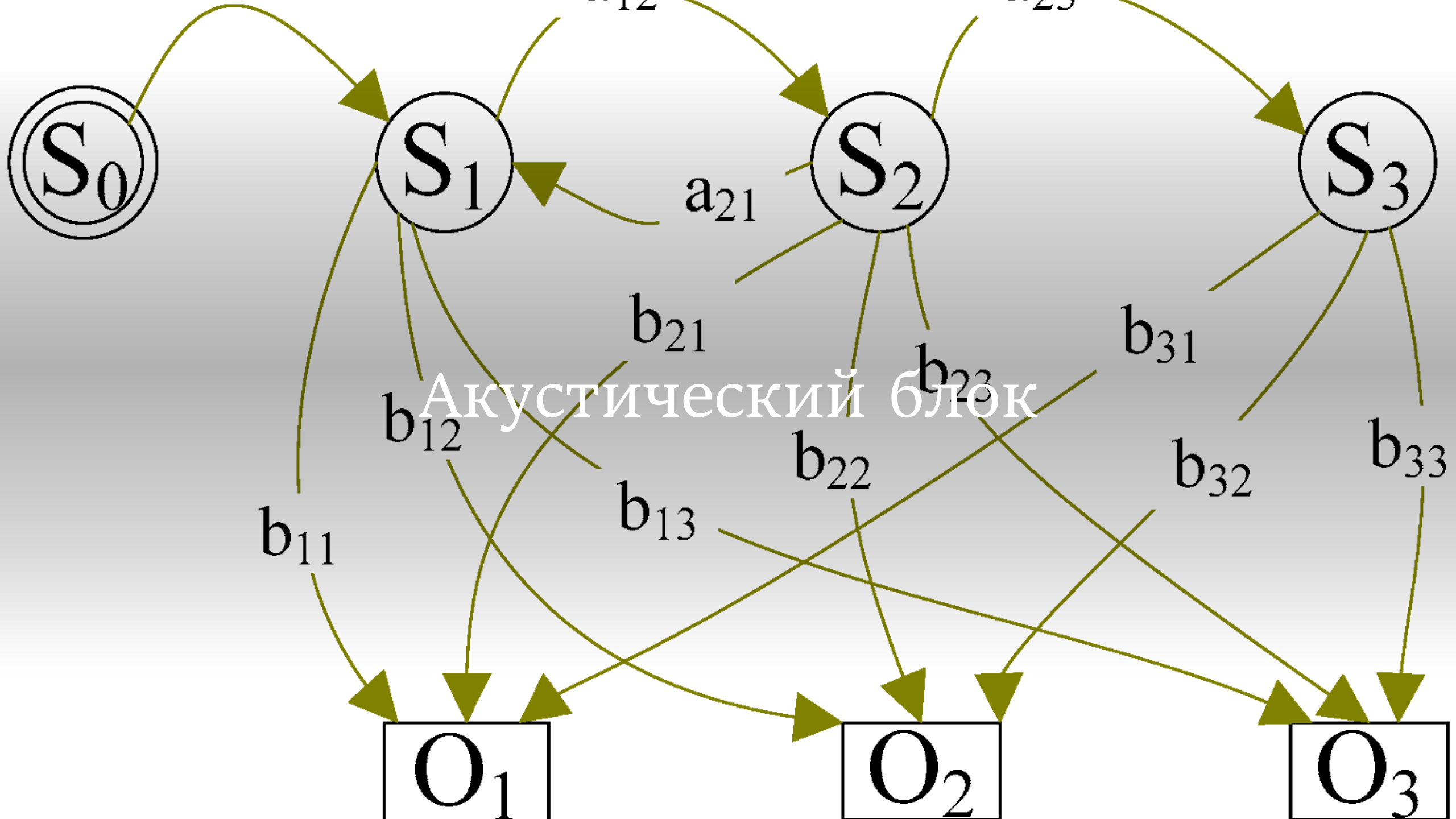
- Гауссовы смеси (Gaussian mixture models, GMM)
- Глубокие нейронные сети (Deep neural networks, DNN)

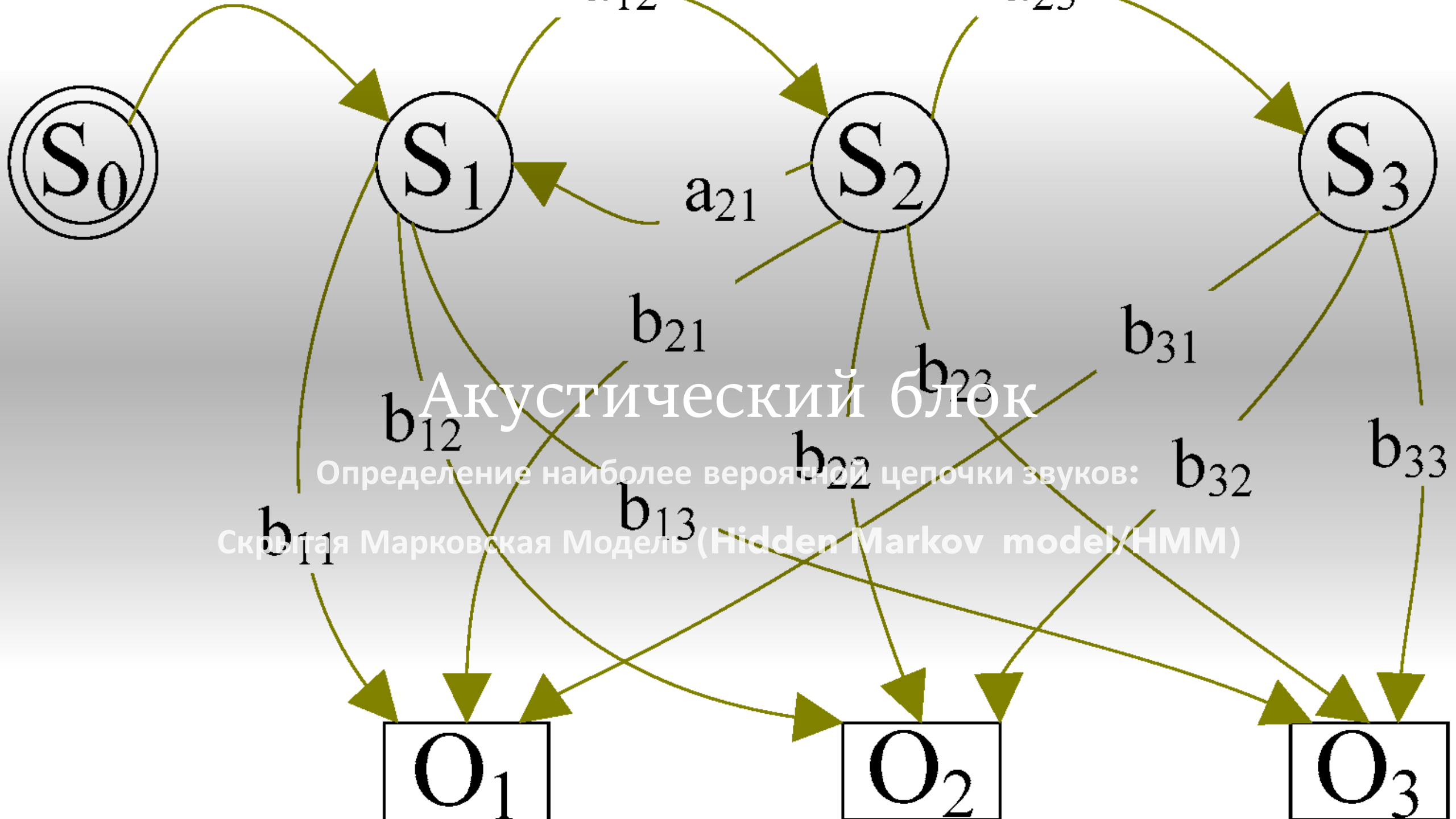
Акустический блок

Определение наиболее
вероятной цепочки звуков

Аудио признаки
+ вероятности
появления
звуков

Цепочка звуков







Лингвистический блок



Лингвистический блок

Языковая модель

- N-граммная языковая модель
- Рекуррентная нейронная сеть в качестве языковой модели (Recurrent neural network language model/RNNLM)

Лингвистический блок



Языковая модель

- N-граммная языковая модель
- Рекуррентная нейронная сеть в качестве языковой модели (Recurrent neural network language model/RNNLM)

Компонентный подход

- CMUSphinx
- HTK
- Kaldi

End-to-End подход

Входной аудио
поток

WAV файл

Извлечение
признаков

Мел-частотные
кепстральные
коэффициенты (MFCC)

I-vectors

Кепстральная
нормализация (CMVN)

Акустический
блок

Глубокие нейронные
сети (DNN)

Лингвистический
блок

Статистическая
языковая модель (n-
gram LM)

Выходной текст

Строка

Акустический блок

Определение наиболее
вероятной цепочки букв

Акустический блок

Определение наиболее
вероятной цепочки букв

Аудио
признаки

Буквы

End-to-End ПОДХОД

- Kaldi
- DeepSpeech
- Wav2letter

Практика

[https://github.com/DinoTheDinosaur/
FocusStart_NLP/blob/master/noteboo
ks/Speech_recognition_and_sentimen
t.ipynb](https://github.com/DinoTheDinosaur/FocusStart_NLP/blob/master/notebooks/Speech_recognition_and_sentiment.ipynb)