

Speech emotion recognition

Hamed Vaheb, Elnaz Khaveh

Introduction

Models

Conclusion



Introduction

Applications

Data set



Introduction

Importance:

Applications

Data set



Introduction

Importance:

- The most natural way to express ourselves

Applications

Data set



Introduction

Importance:

- The most natural way to express ourselves

The Evolution of Human Vocal Emotion
Gregory A. Bryant¹
First Published June 24, 2020 | Research Article |  <https://doi.org/10.1177/1754073920930791>
[Article information](#) ▾

Applications

Data set

Introduction

Importance:

- The most natural way to express ourselves
- Influence over judgement
 - "go home!"

The Evolution of Human Vocal Emotion
Gregory A. Bryant¹
First Published June 24, 2020 | Research Article |  <https://doi.org/10.1177/1754073920930791>
[Article information](#) ▾

Applications

Data set



Introduction

Importance:

- The most natural way to express ourselves
- Influence over judgement
 - "go home!"
- Human-computer interaction (HCI)
 - Personal Assistants
 - Next Generation AI

The Evolution of Human Vocal Emotion
Gregory A. Bryant
First Published June 24, 2020 | Research Article |  <https://doi.org/10.1177/1754073920930791>
[Article information](#) ▾

Applications

Data set

Applications

Applications

- Psychological consultation



Applications

- Psychological consultation



Applications

- Psychological consultation



Predicting couple therapy outcomes based on speech acoustic features

Md Nasir, Brian Robert Baucom, Panayiotis Georgiou, Shrikanth Narayanan

Published: September 21, 2017 • <https://doi.org/10.1371/journal.pone.0185123>

Applications

- Psychological consultation



Predicting couple therapy outcomes based on speech acoustic features

Md Nasir, Brian Robert Baucom, Panayiotis Georgiou, Shrikanth Narayanan

Published: September 21, 2017 • <https://doi.org/10.1371/journal.pone.0185123>

Interspeech. 2018 Sep; 2018: 3413–3417.

doi: [10.21437/interspeech.2018-2551](https://doi.org/10.21437/interspeech.2018-2551)

PMID:

Using Prosodic and Lexical Information for Learning Utterance-level Behaviors in Psychotherapy

Karan Singla,¹ Zhuohao Chen,¹ Nikolaos Flemotomos,¹ James Gibson,¹ Dogan Can,¹ David C. Atkins,² and Shrikanth Narayanan¹

Applications



Predicting couple therapy outcomes based on speech acoustic features

Md Nasir, Brian Robert Baucom, Panayiotis Georgiou, Shrikanth Narayanan

Published: September 21, 2017 • <https://doi.org/10.1371/journal.pone.0185123>

Interspeech. 2018 Sep; 2018: 3413–3417.

doi: [10.21437/interspeech.2018-2551](https://doi.org/10.21437/interspeech.2018-2551)

PMID:

Using Prosodic and Lexical Information for Learning Utterance-level Behaviors in Psychotherapy

Karan Singla,¹ Zhuohao Chen,¹ Nikolaos Flemotomos,¹ James Gibson,¹ Dogan Can,¹ David C. Atkins,² and Shrikanth Narayanan¹

Applications



Interspeech. 2018 Sep; 2018: 3413–3417.
doi: [10.21437/interspeech.2018-2551](https://doi.org/10.21437/interspeech.2018-2551)

PMID:

Using Prosodic and Lexical Information for Learning Utterance-level Behaviors in Psychotherapy

Karan Singla,¹ Zhuohao Chen,¹ Nikolaos Flemotomos,¹ James Gibson,¹ Dogan Can,¹ David C. Atkins,² and Shrikanth Narayanan¹

Applications

Interspeech. 2018 Sep; 2018: 3413–3417.
doi: [10.21437/interspeech.2018-2551](https://doi.org/10.21437/interspeech.2018-2551)

PMID:

Using Prosodic and Lexical Information for Learning Utterance-level Behaviors in Psychotherapy

Karan Singla,¹ Zhuohao Chen,¹ Nikolaos Flemotomos,¹ James Gibson,¹ Dogan Can,¹
David C. Atkins,² and Shrikanth Narayanan¹

Applications

Applications

- Call centers

Applications

- Call centers

Recognizing low/high anger in speech for call centers

Authors:  Fu-Ming Lee,  Li-Hua Li,  Ru-Yi Huang [Authors Info & Claims](#)

ISPRA'08: Proceedings of the 7th WSEAS International Conference on Signal Processing, Robotics and Automation • February 2008 • Pages 171–176

Introduction

Importance:

- The most natural way to express ourselves
- Influence over judgement
 - "go home!"
- Human-computer interaction (HCI)
 - Personal Assistants
 - Next Generation AI

The Evolution of Human Vocal Emotion
Gregory A. Bryant¹
First Published June 24, 2020 | Research Article |  <https://doi.org/10.1177/1754073920930791>
[Article information](#) ▾

Applications

Data set



RAVDESS

Plots

Processing

Feature
extraction

RAVDESS

- A validated multimodal database

Plots

Processing

Feature extraction

RAVDESS

- A validated multimodal database
- 24 professional actors

Plots

Processing

Feature extraction

RAVDESS

- A validated multimodal database
- 24 professional actors
- Gender-balanced

Plots

Processing

Feature extraction

RAVDESS

- 24 professional actors
- Gender-balanced

Plots

Processing

Feature extraction

RAVDESS

- Gender-balanced

Plots

Processing

Feature extraction

RAVDESS

Plots

Processing

Feature extraction

RAVDESS

- 8 classes :

Plots

Processing

Feature extraction

RAVDESS

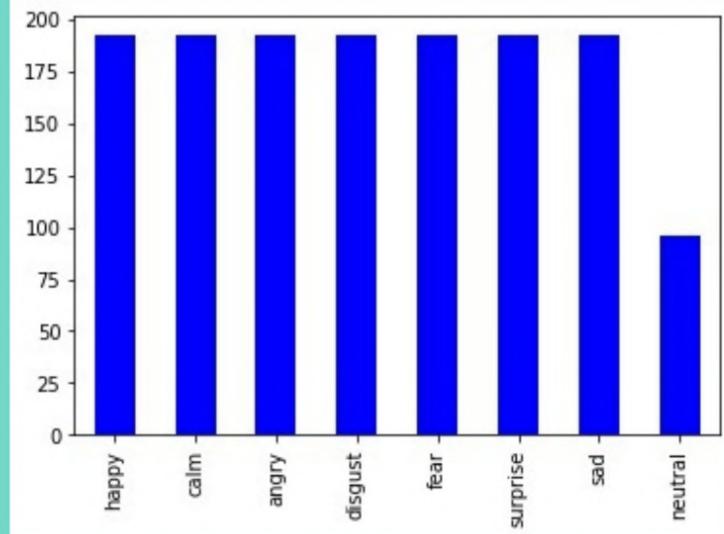
- 8 classes :
Neutral, calm, happy, sad, angry, fearful, disgust, surprise

Plots

Processing

Feature extraction

RAVDESS



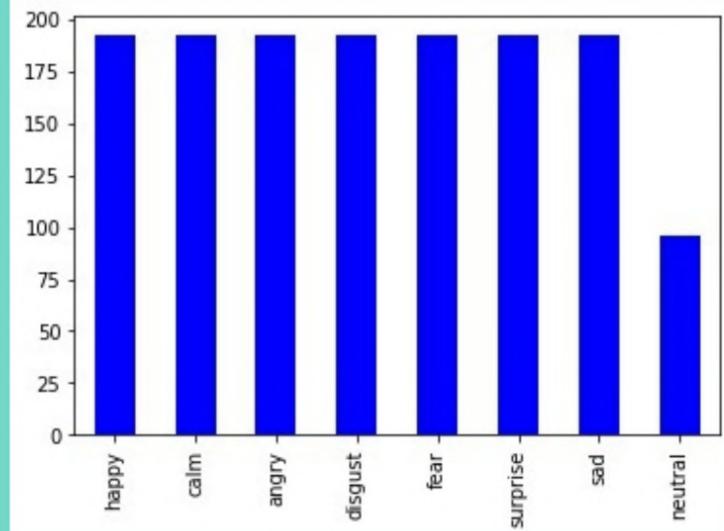
- 8 classes :
Neutral, calm, happy, sad, angry, fearful, disgust,
surprise

Plots

Processing

Feature
extraction

RAVDESS



Neutral, calm, happy, sad, angry, fearful, disgusted,
surprised

Plots

Processing

Feature extraction

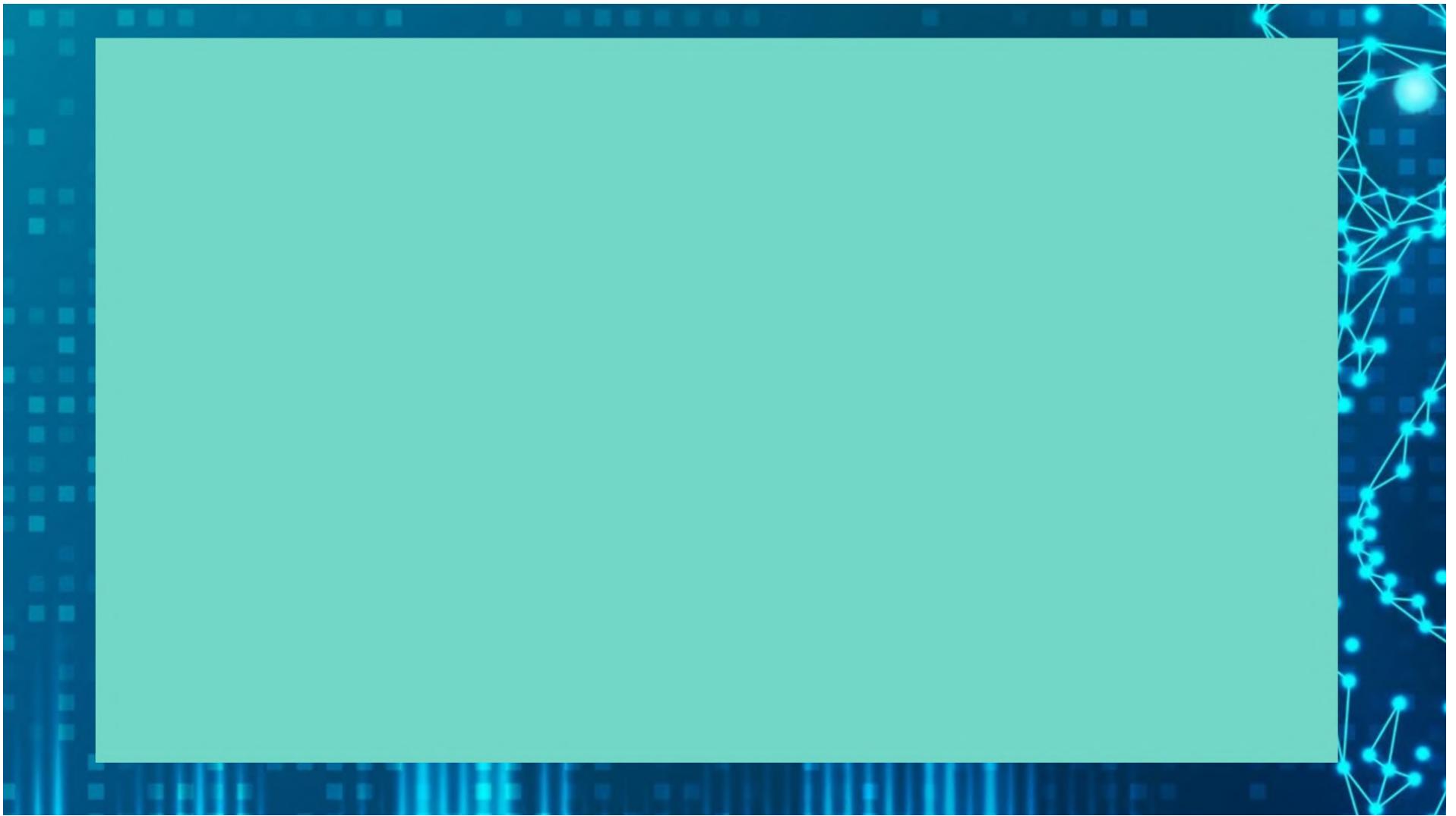
RAVDESS

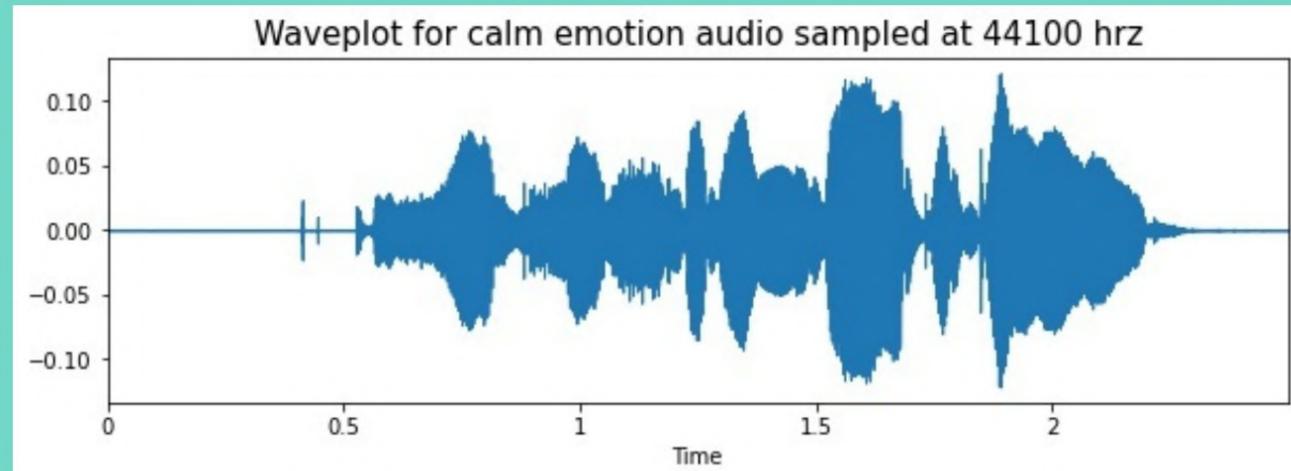
Neutral, calm, happy, sad, angry, fearful, disgust, surprise

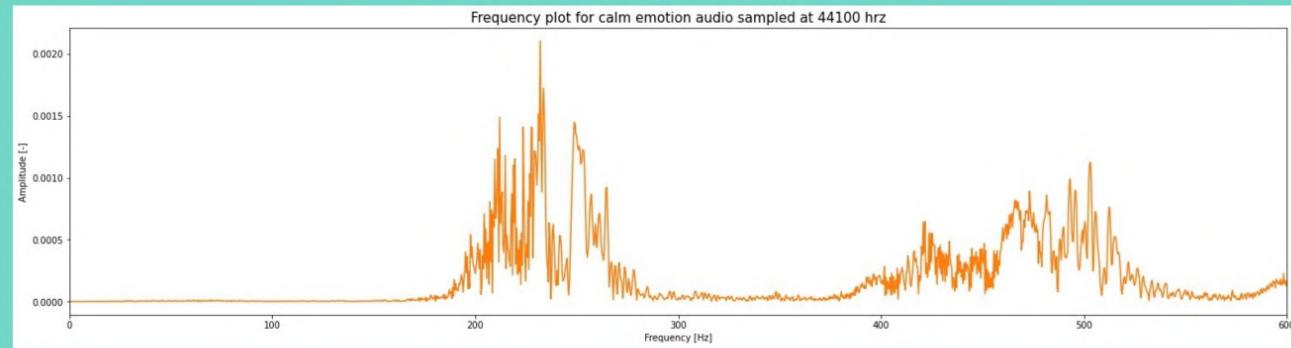
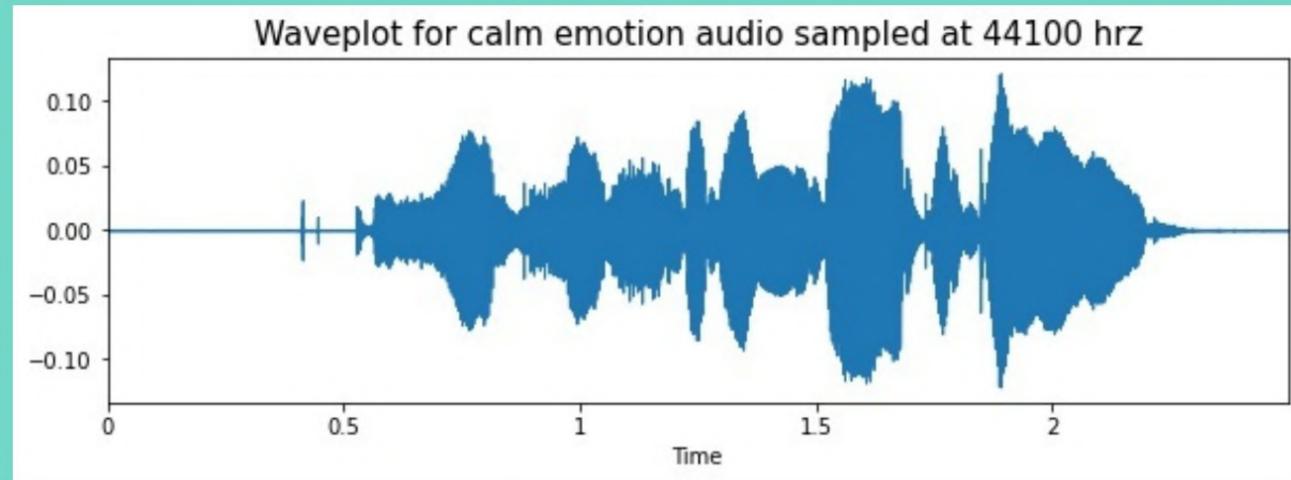
Plots

Processing

Feature extraction







RAVDESS

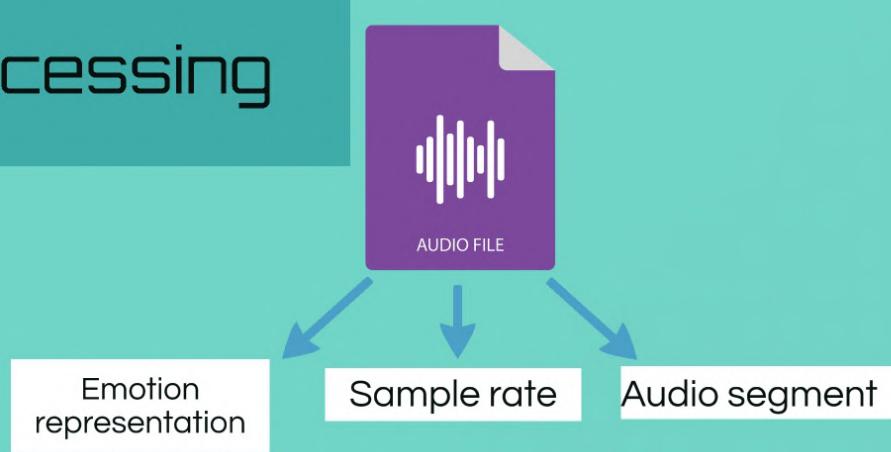
Neutral, calm, happy, sad, angry, fearful, disgust, surprise

Plots

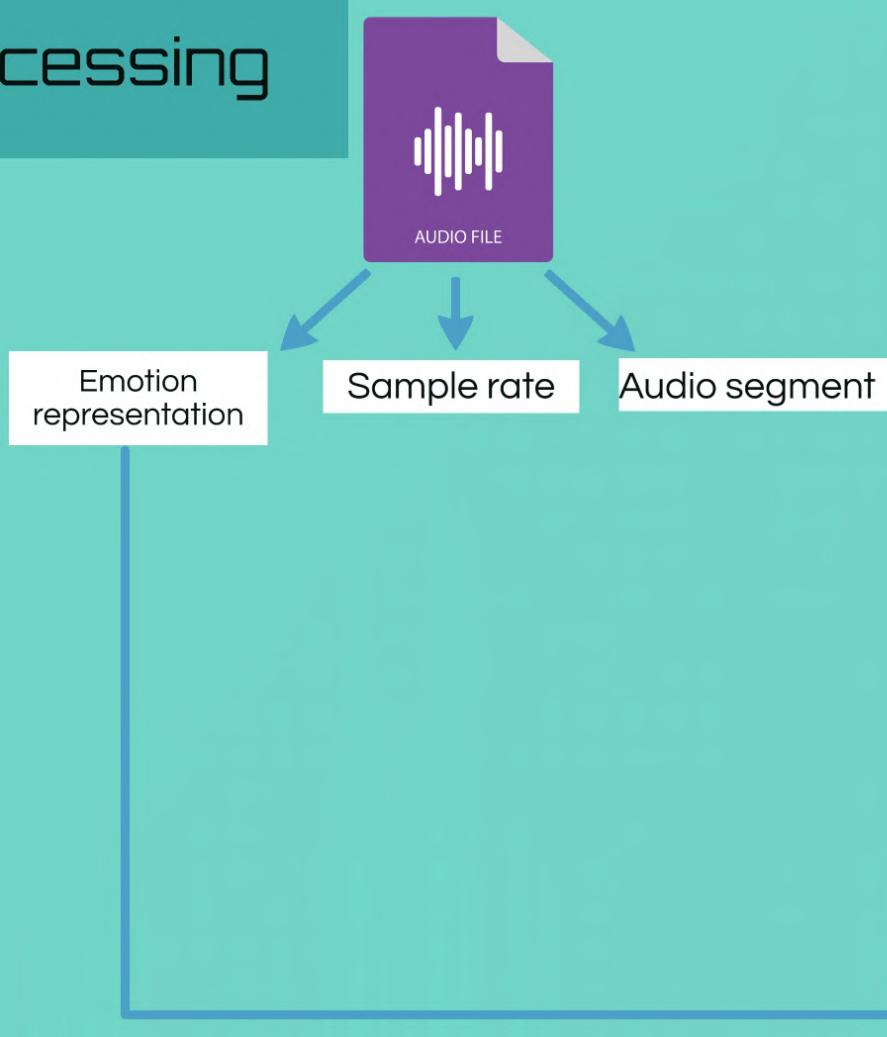
Processing

Feature extraction

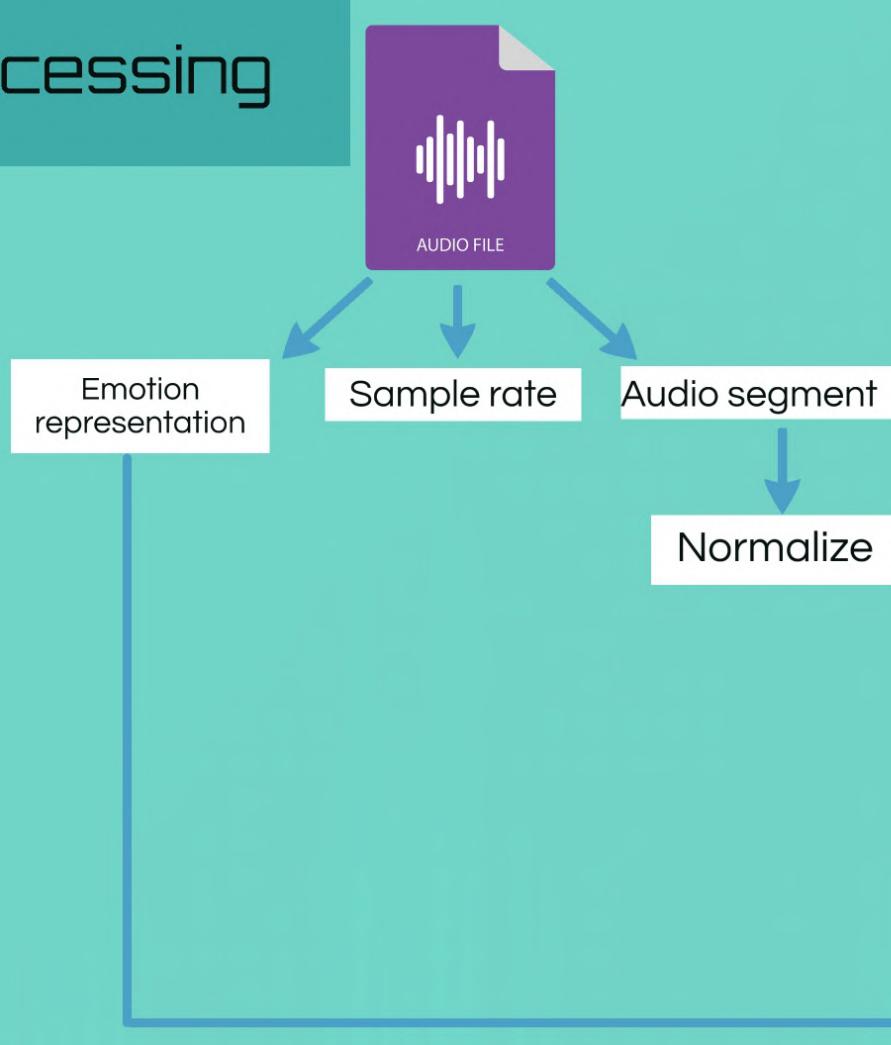
Processing



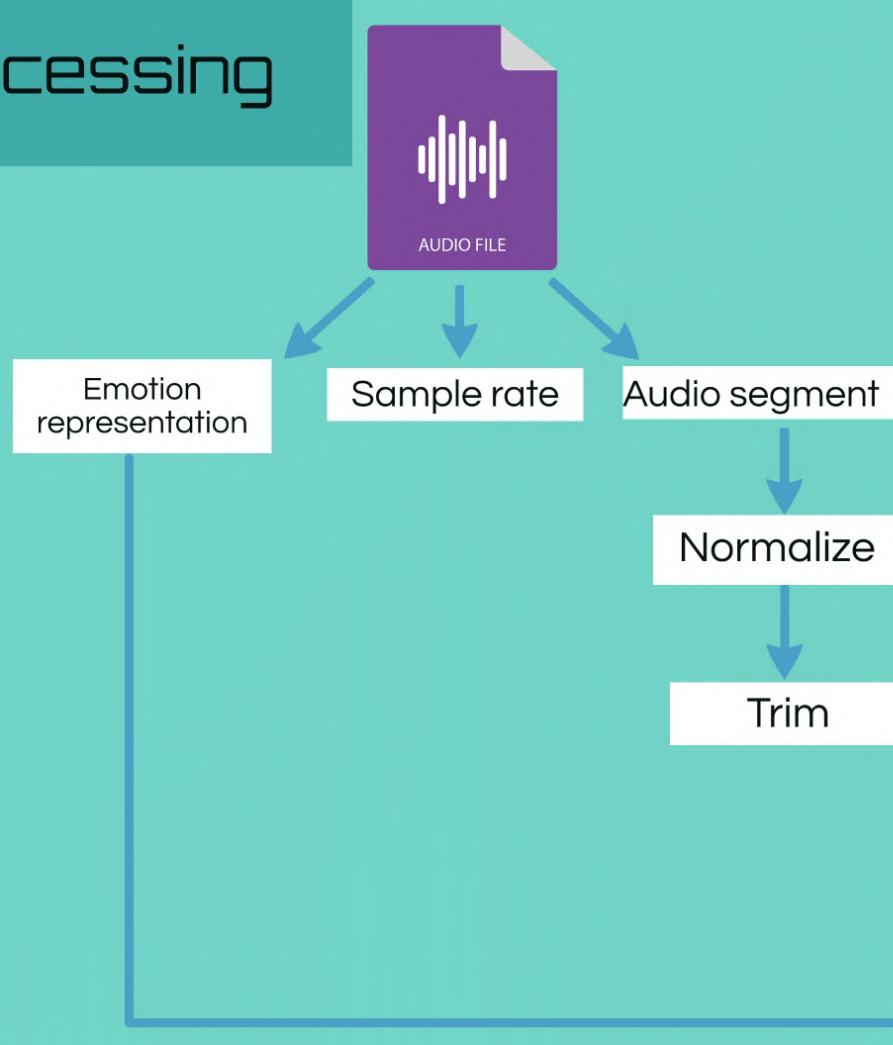
Processing



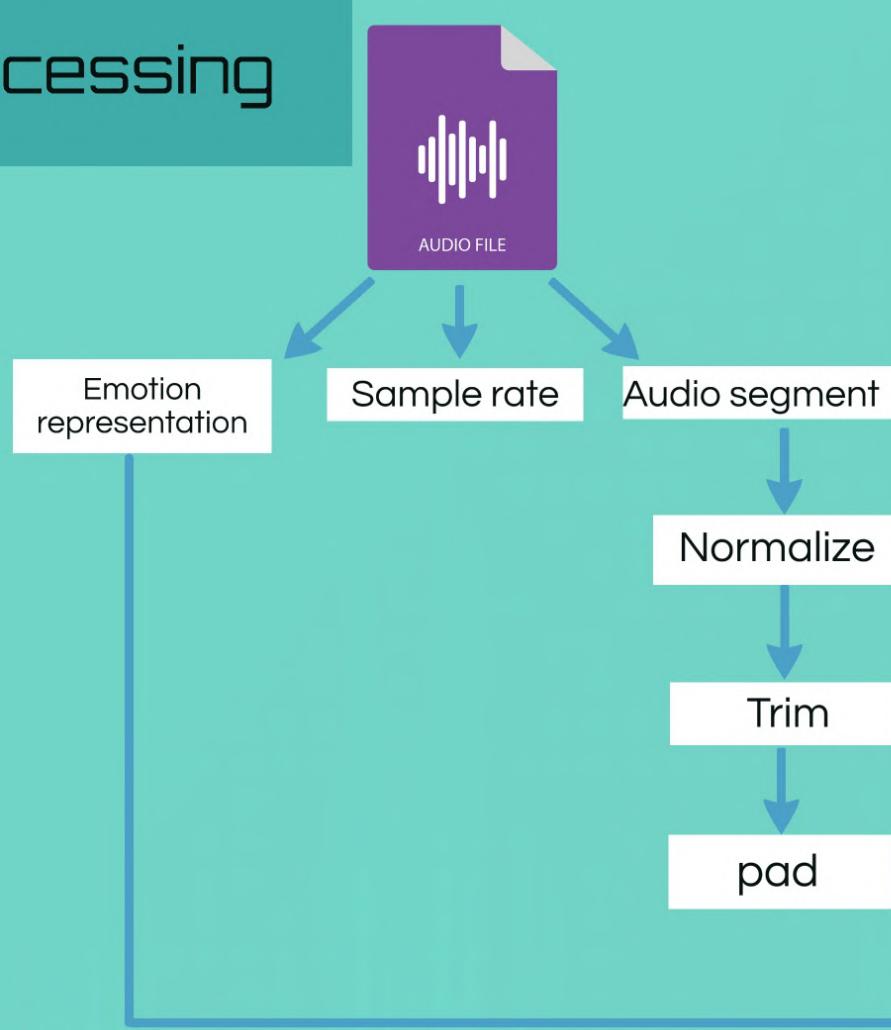
Processing

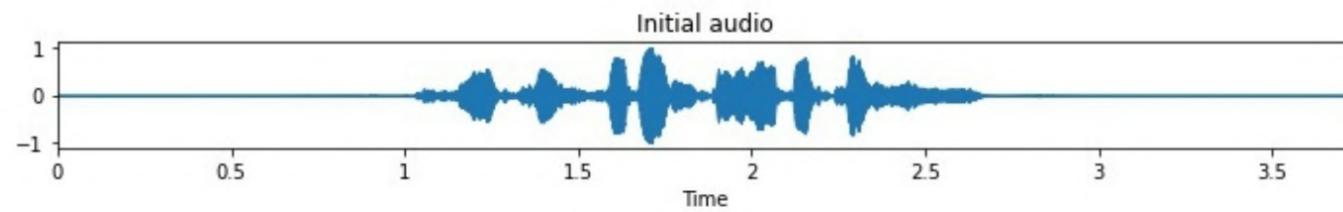


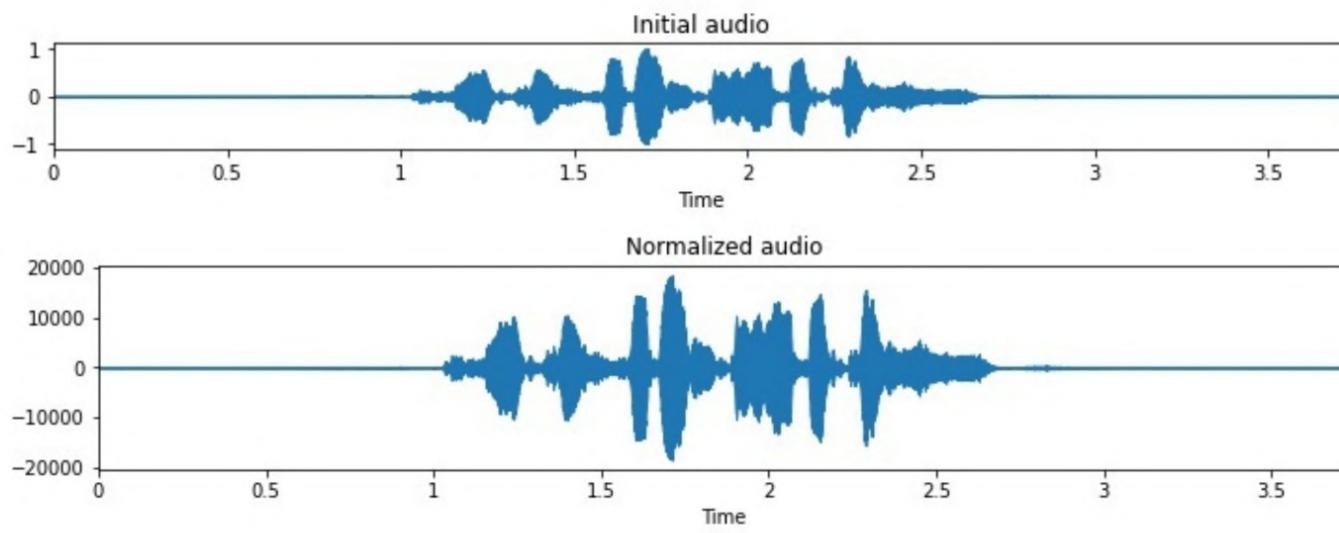
Processing

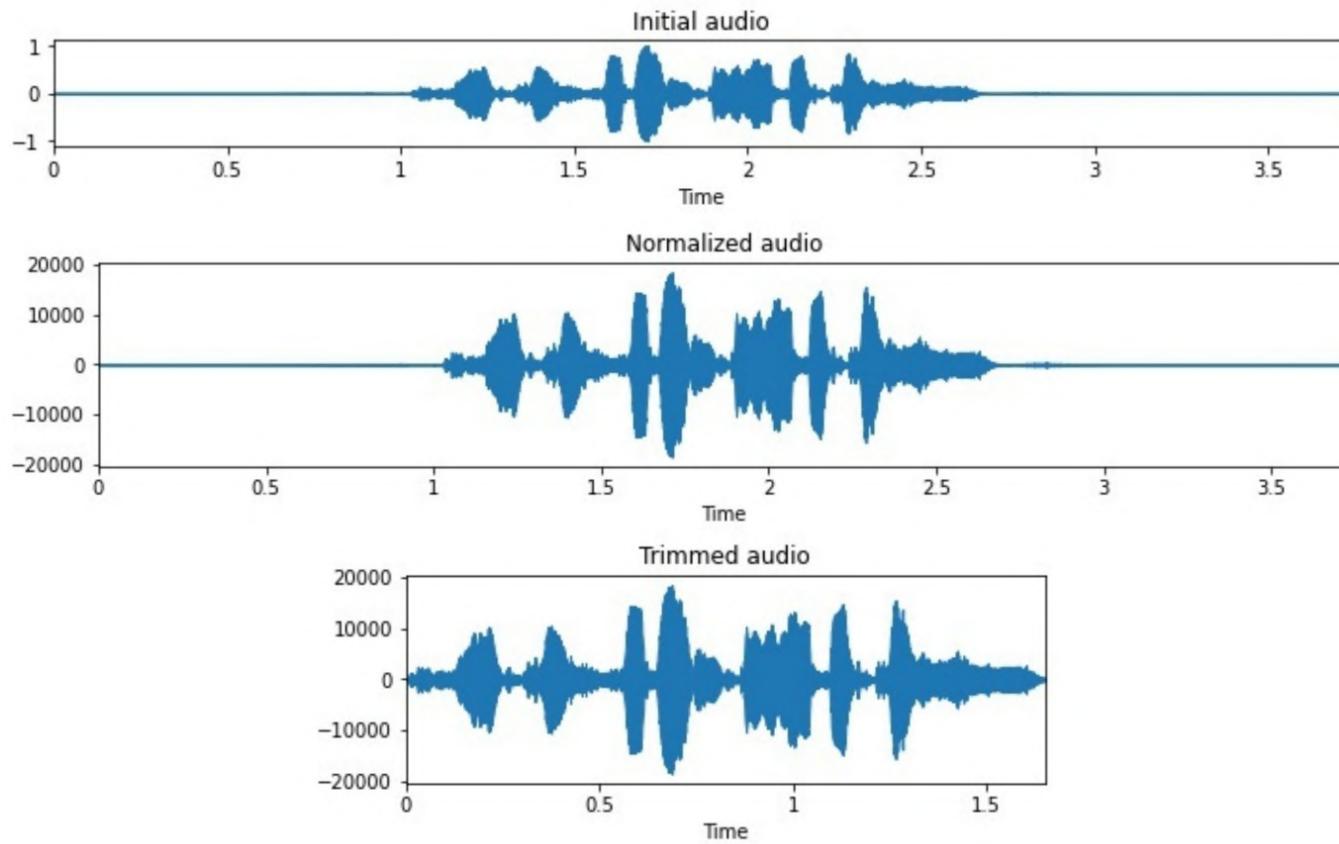


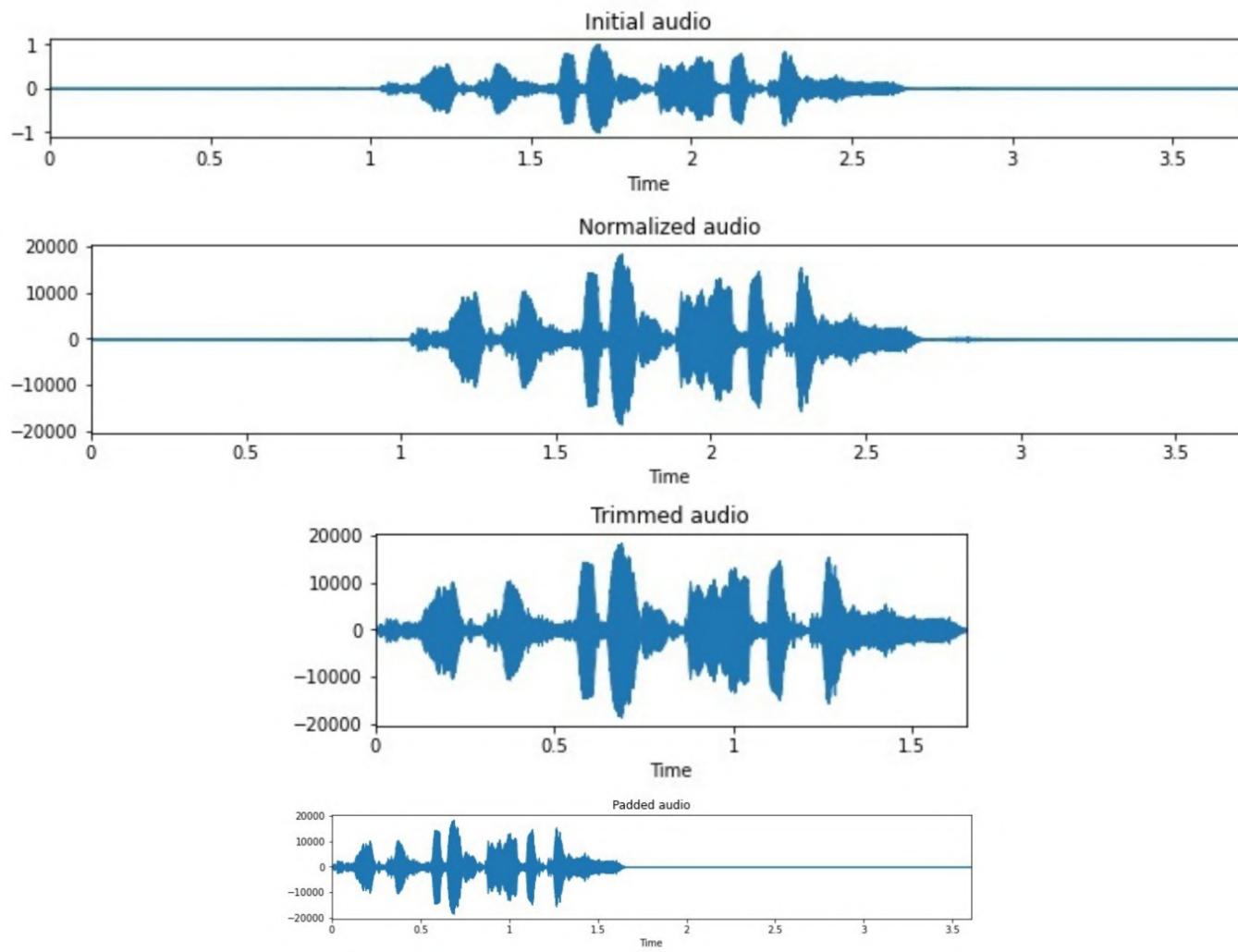
Processing



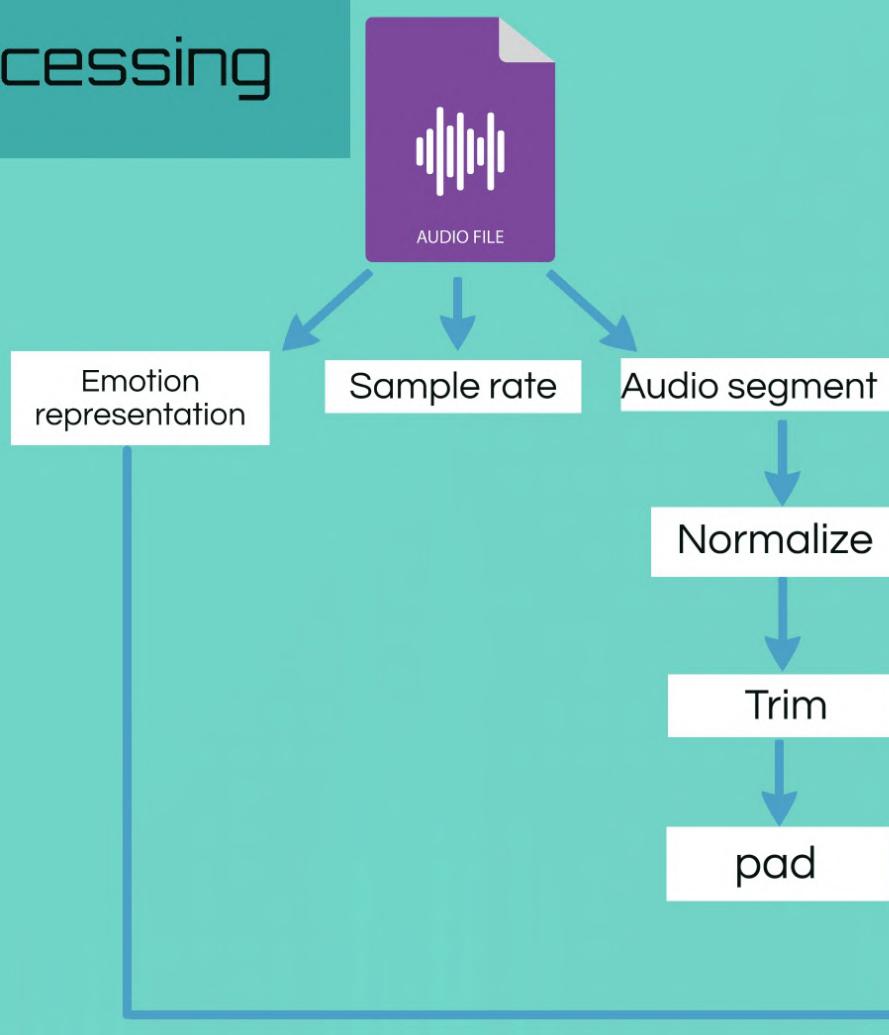




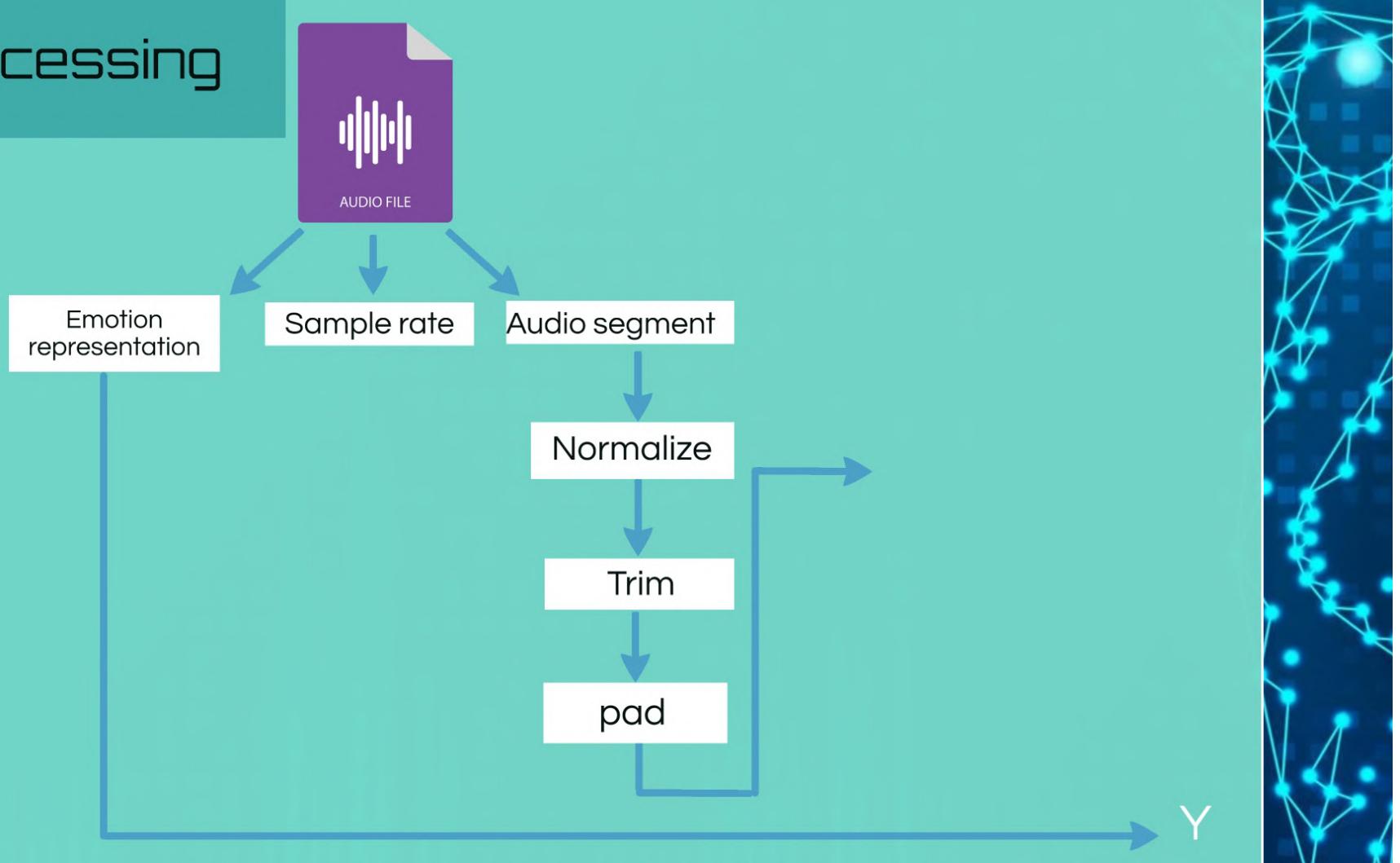




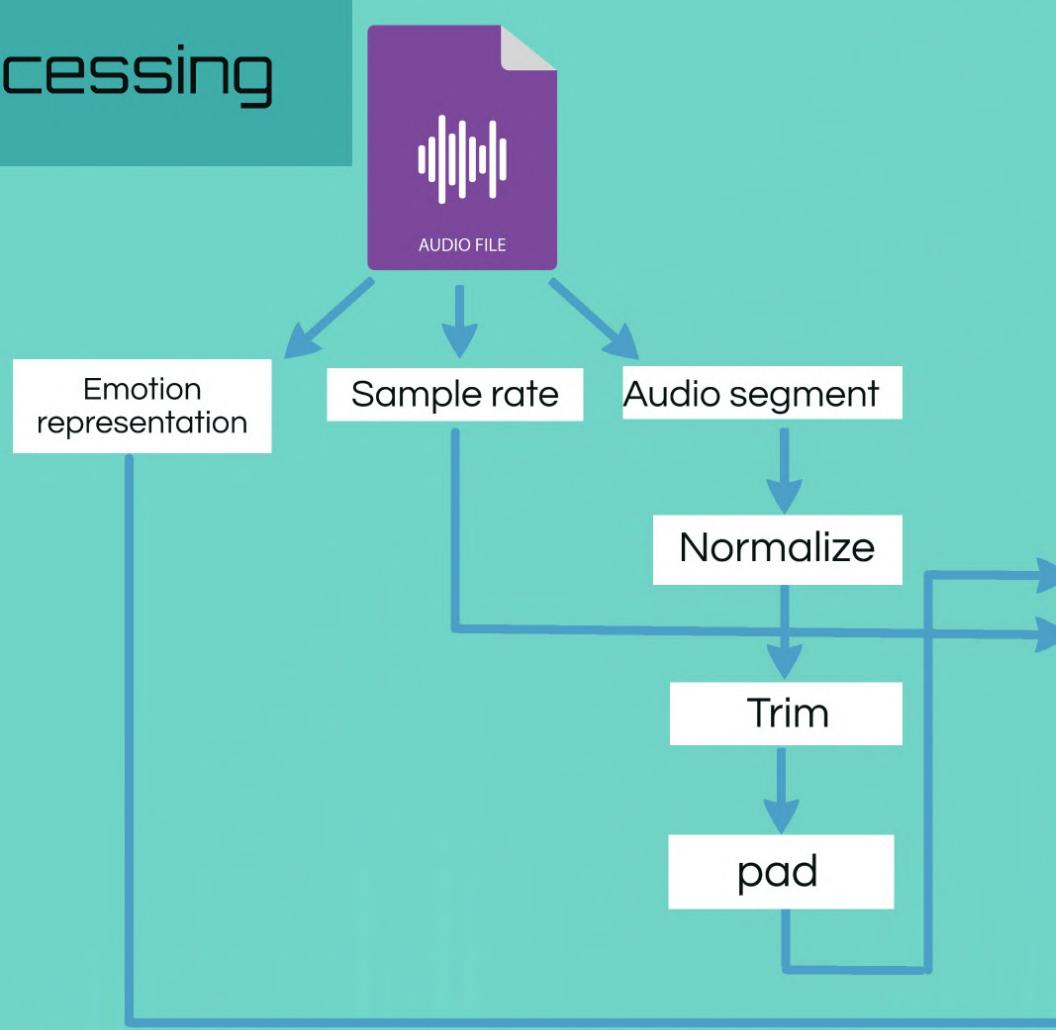
Processing



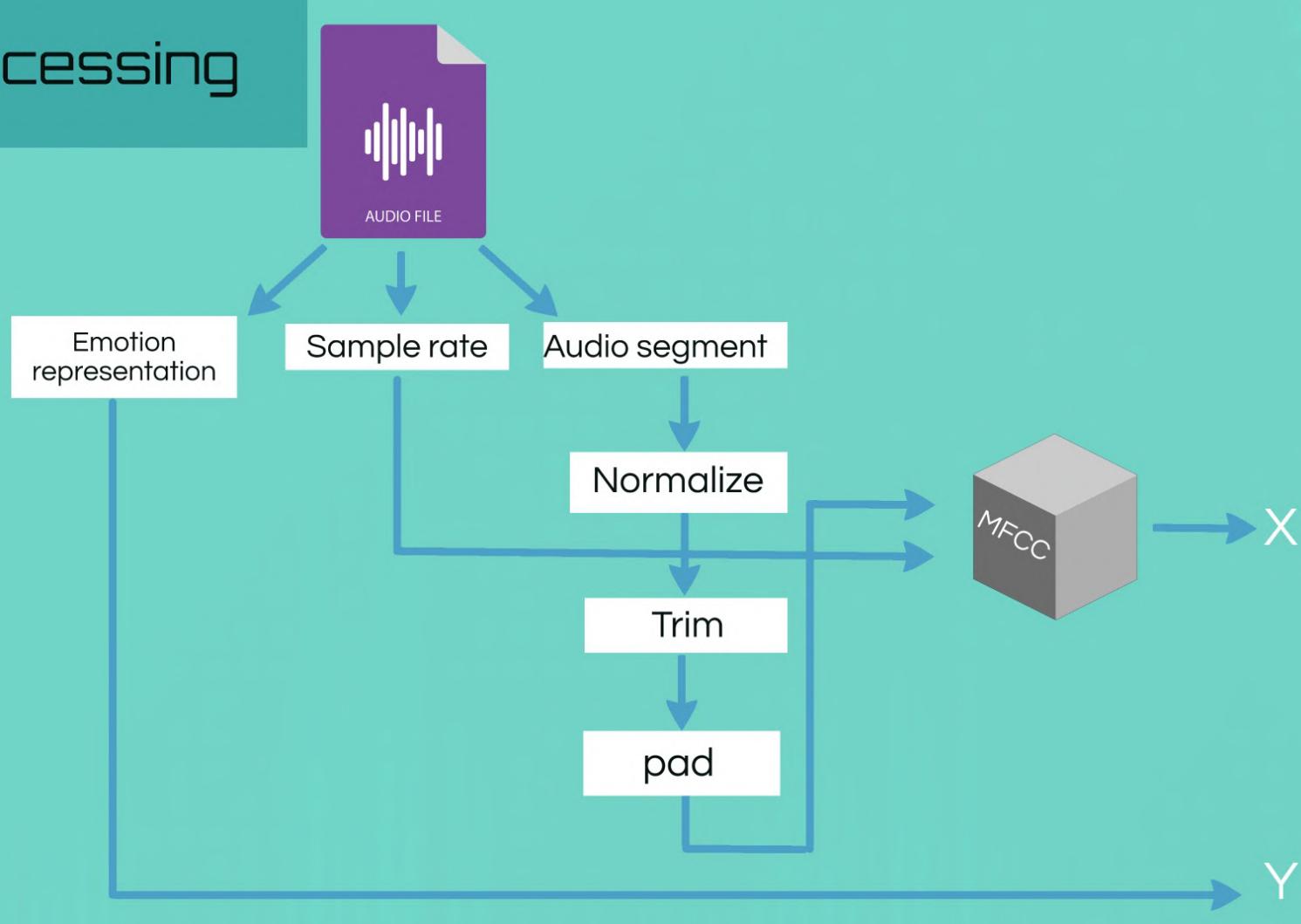
Processing



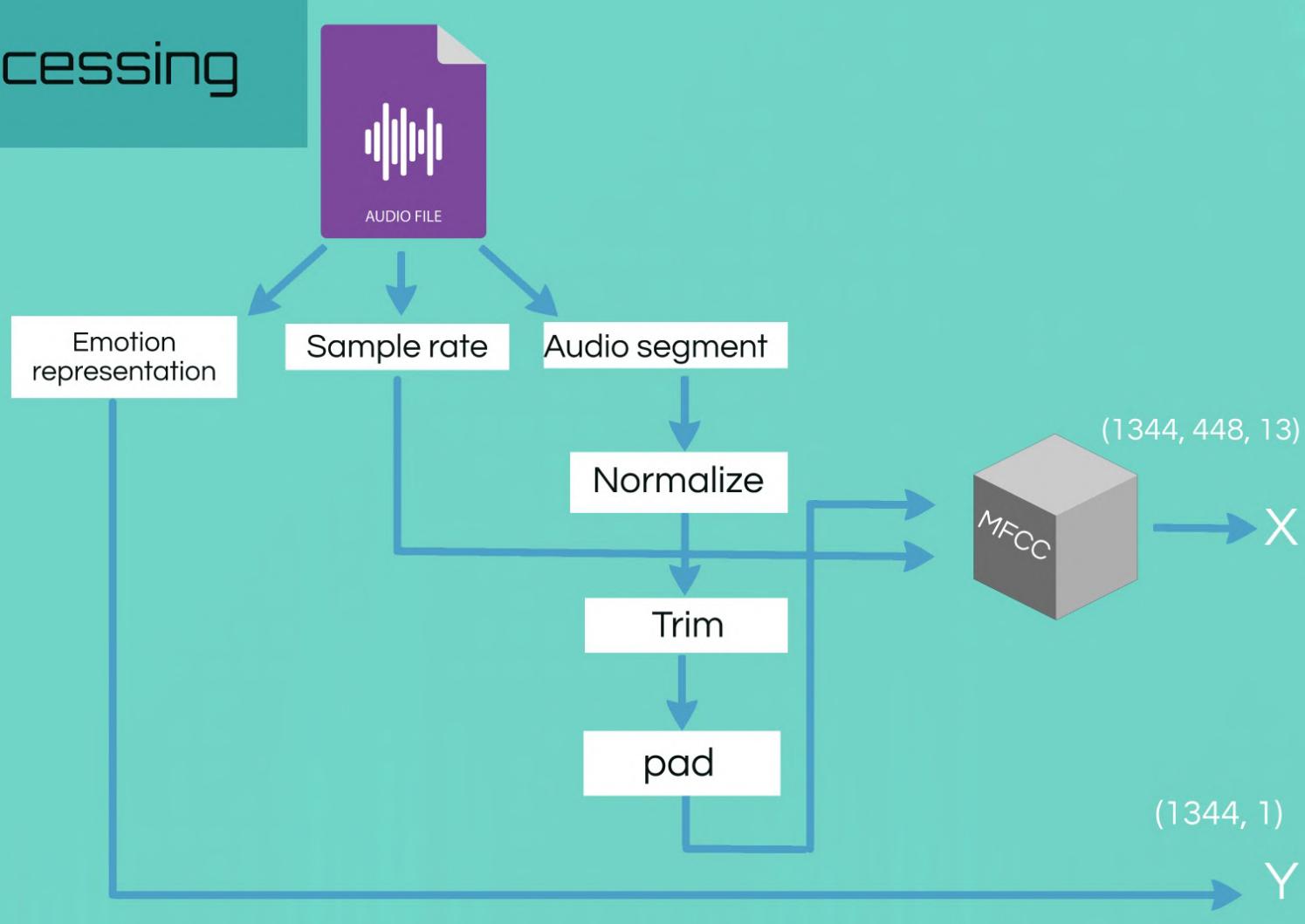
Processing



Processing



Processing



RAVDESS

Neutral, calm, happy, sad, angry, fearful, disgust, surprise

Plots

Processing

Feature extraction

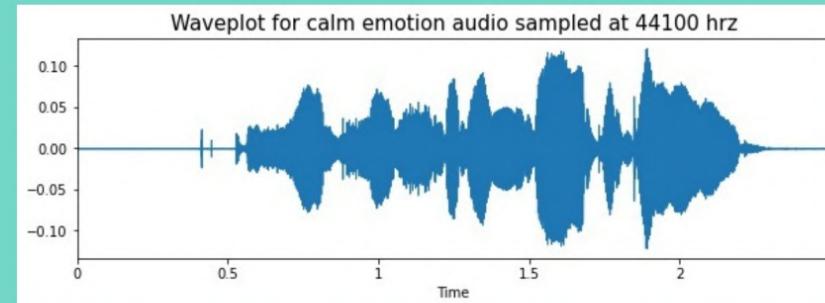
Feature extraction

Mel-Spectrogram

MFCC

Feature extraction

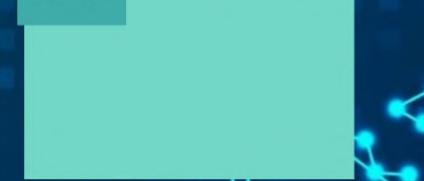
- Time-Domain
 - RMSE, ZCR



Mel-Spectrogram

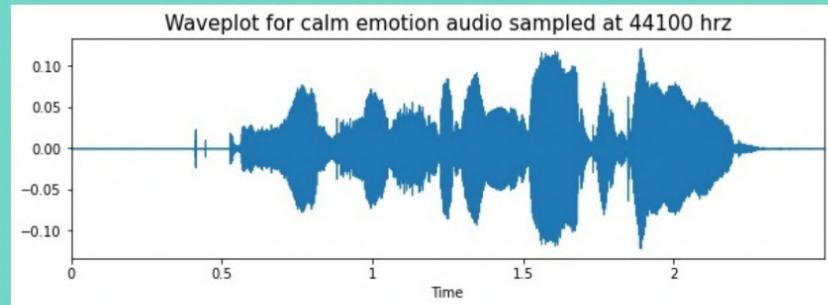


MFCC

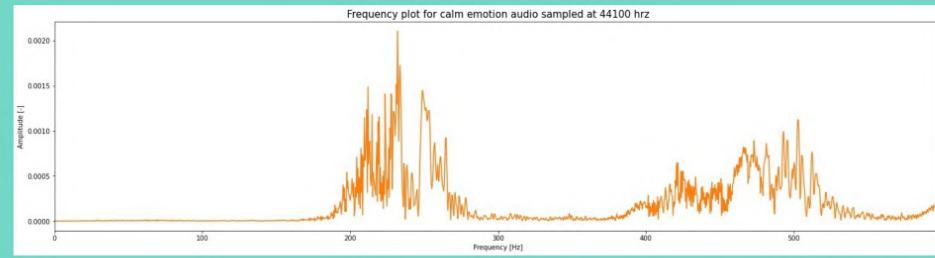


Feature extraction

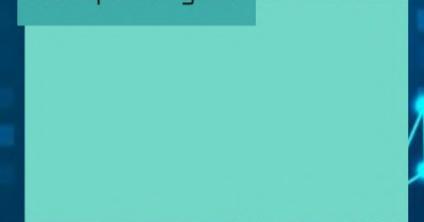
- Time-Domain
 - RMSE, ZCR



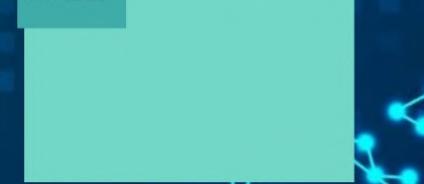
- Frequency Domain
 - Spectral Centroid



Mel-Spectrogram



MFCC



Mel Spectrogram

Mel Spectrogram

1. Extract STFT

Mel Spectrogram

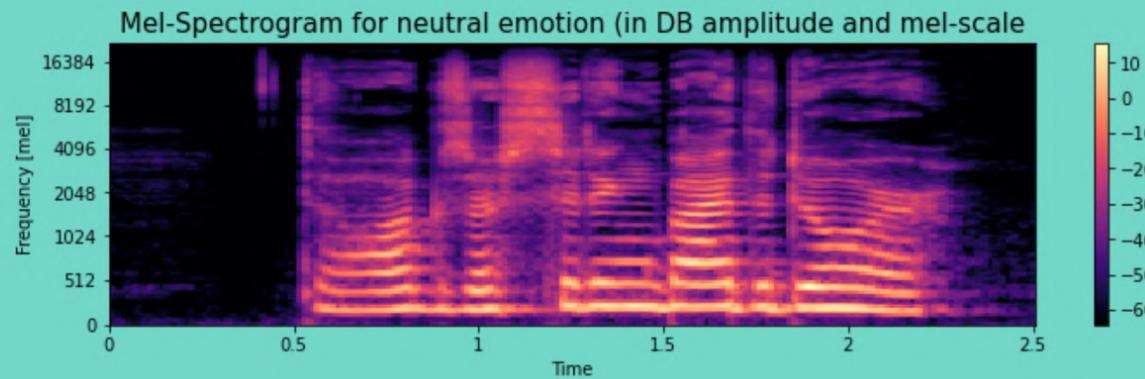
1. Extract STFT
2. Convert amplitude to DBs

Mel Spectrogram

1. Extract STFT
2. Convert amplitude to DBs
3. Convert frequencies to Mel Scale

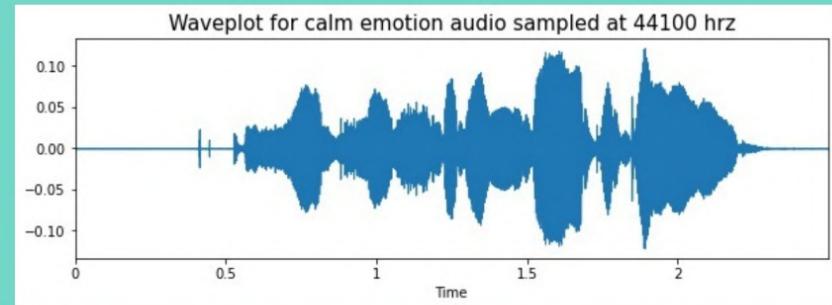
Mel Spectrogram

1. Extract STFT
2. Convert amplitude to DBs
3. Convert frequencies to Mel Scale

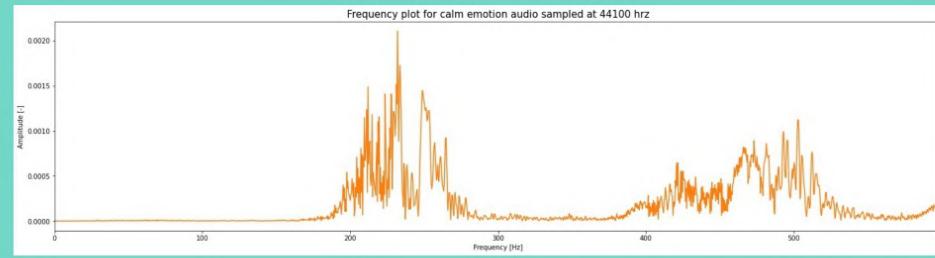


Feature extraction

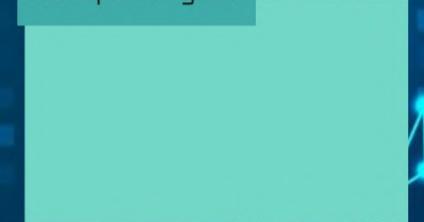
- Time-Domain
 - RMSE, ZCR



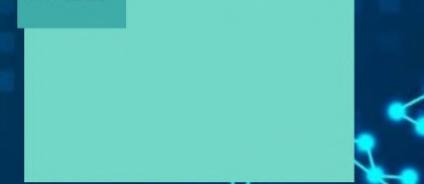
- Frequency Domain
 - Spectral Centroid



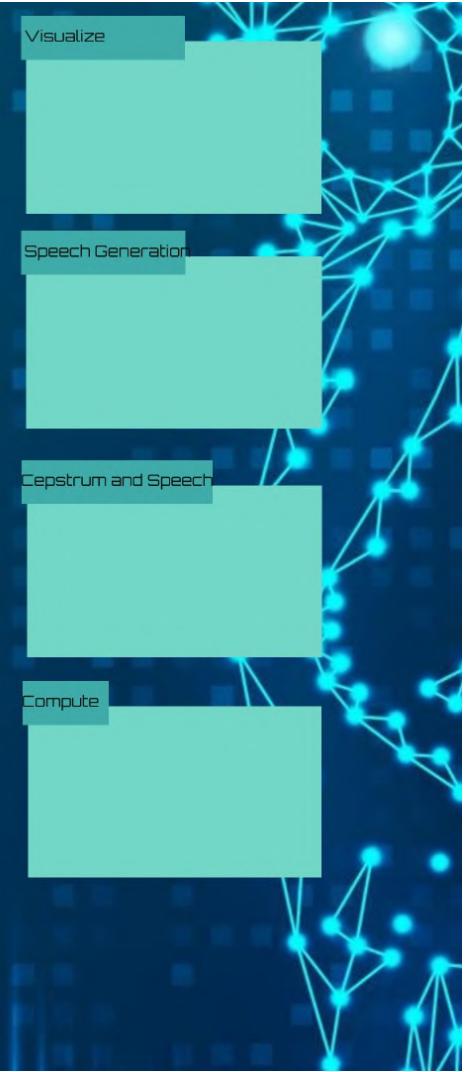
Mel-Spectrogram



MFCC



Mel-Frequency Cepstrum Coefficients (MFCCs)



Mel-Frequency Cepstrum Coefficients (MFCCs)

Cepstrum

Time-domain signal

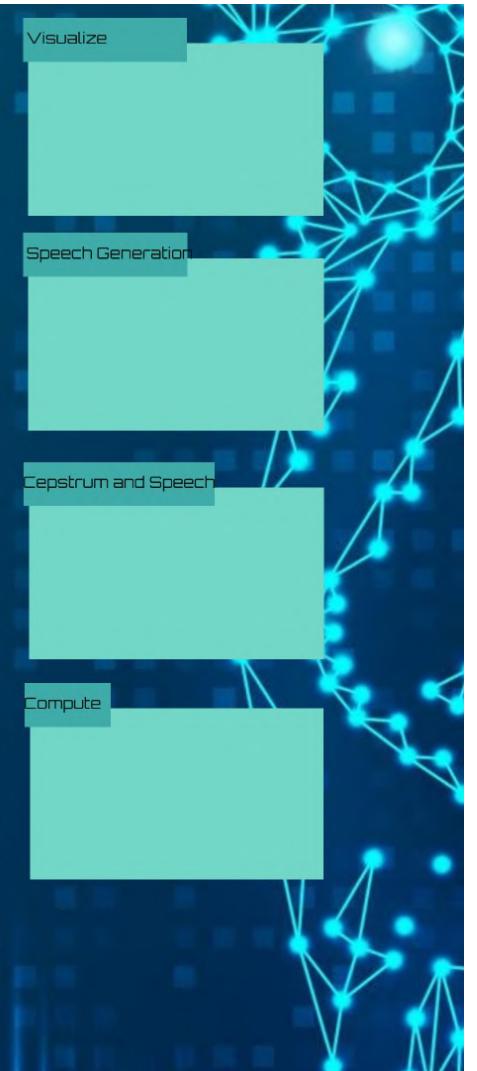
$$C(x(t)) = F^{-1} [\log(F[x(t)])]$$

Spectrum

Log spectrum

Cepstrum

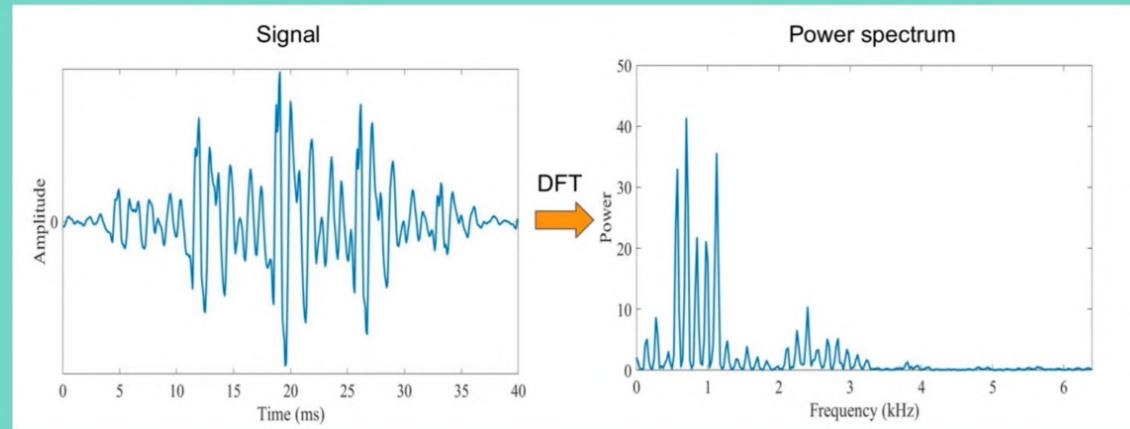
The diagram shows the mathematical formula for computing the Cepstrum. A pink box highlights the term $x(t)$ in the argument of the Fourier transform F . Another pink box highlights the term $x(t)$ in the argument of the logarithm \log . A green box encloses the entire expression $F^{-1} [\log(F[x(t)])]$. Inside this green box, an orange box highlights the term $F[x(t)]$ in the argument of the logarithm. Labels indicate the 'Time-domain signal' at the top left, 'Spectrum' above the inner orange box, 'Log spectrum' below it, and 'Cepstrum' at the bottom.





$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

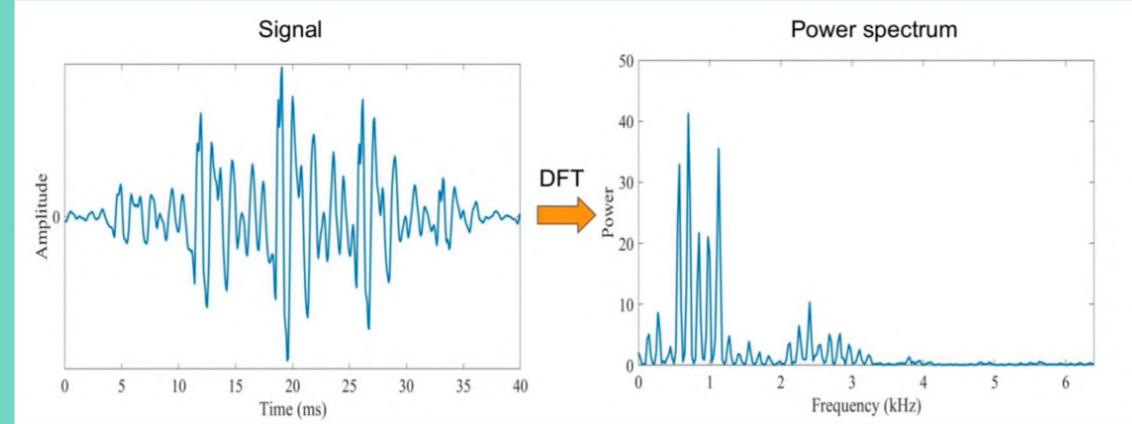
$$C(x(t)) = F^{-1}[\log(F[x(t)])] \quad x(t) \quad F[x(t)]$$



$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

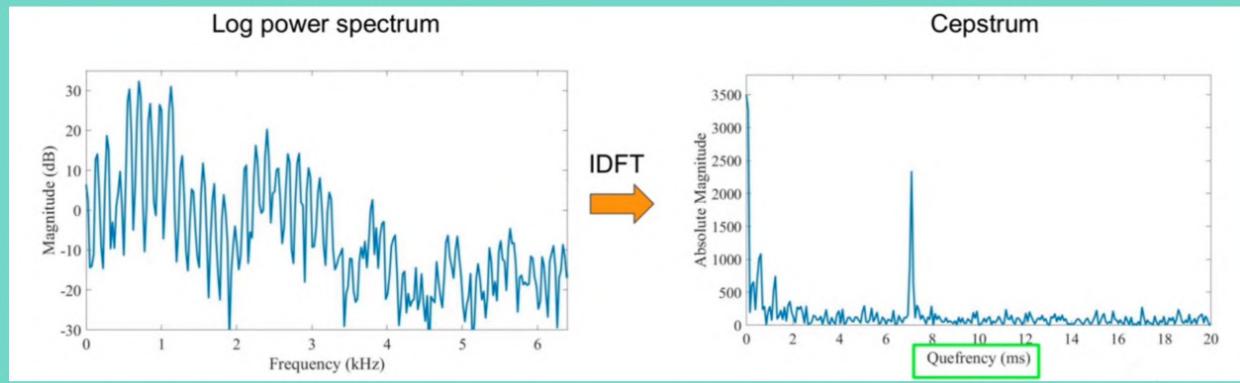
$x(t)$

$F[x(t)]$



$$\log(F[x(t)])$$

$$F^{-1}$$



$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

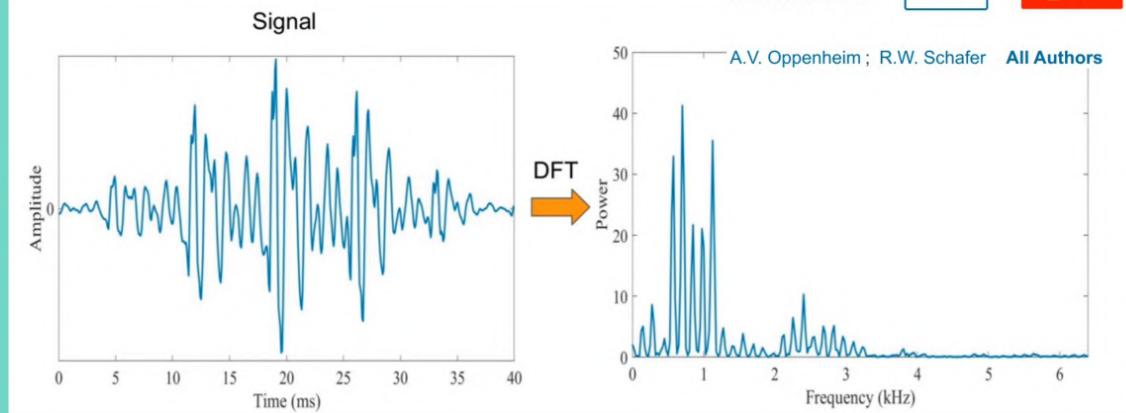
$x(t)$

From frequency to quefrency: a history of the cepstrum

Publisher: IEEE

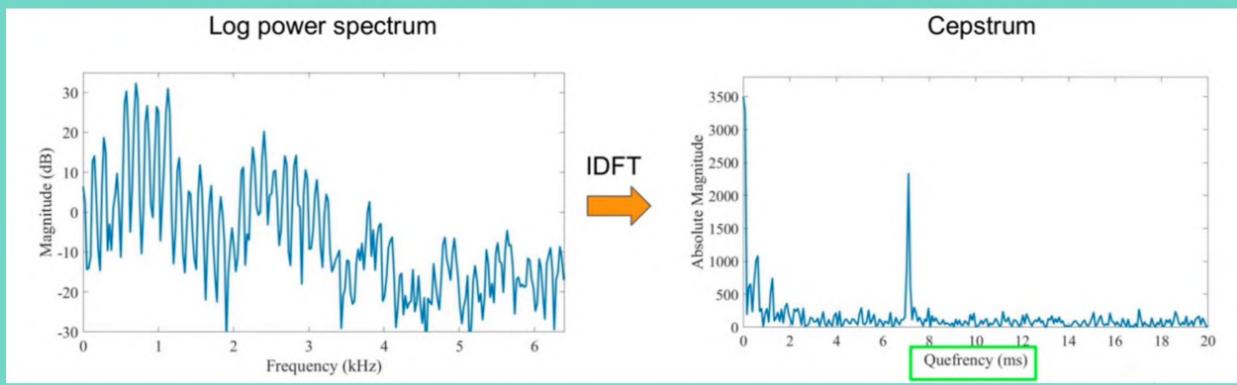
Cite This

PDF



$$\log(F[x(t)])$$

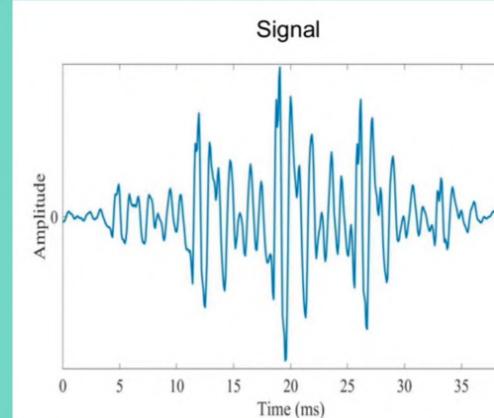
$$F^{-1}$$



$$C(x(t)) = F^{-1}[\log(F[x(t)])]$$

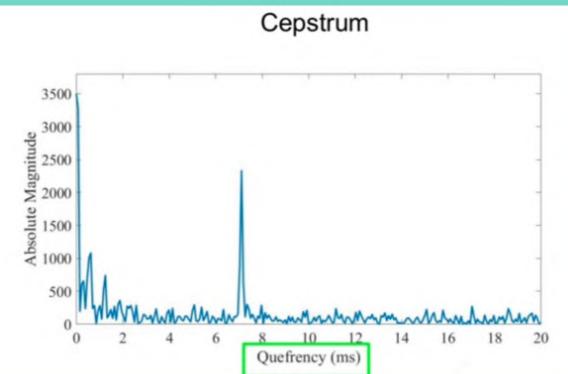
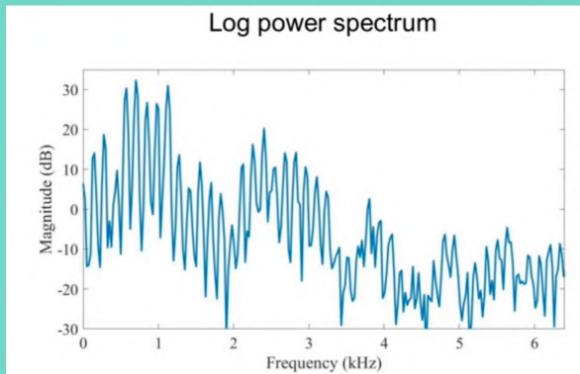
$x(t)$

$F[x(t)]$



$$\log(F[x(t)])$$

$$F^{-1}$$



Mel-Frequency Cepstrum Coefficients (MFCCs)

Cepstrum

Time-domain signal

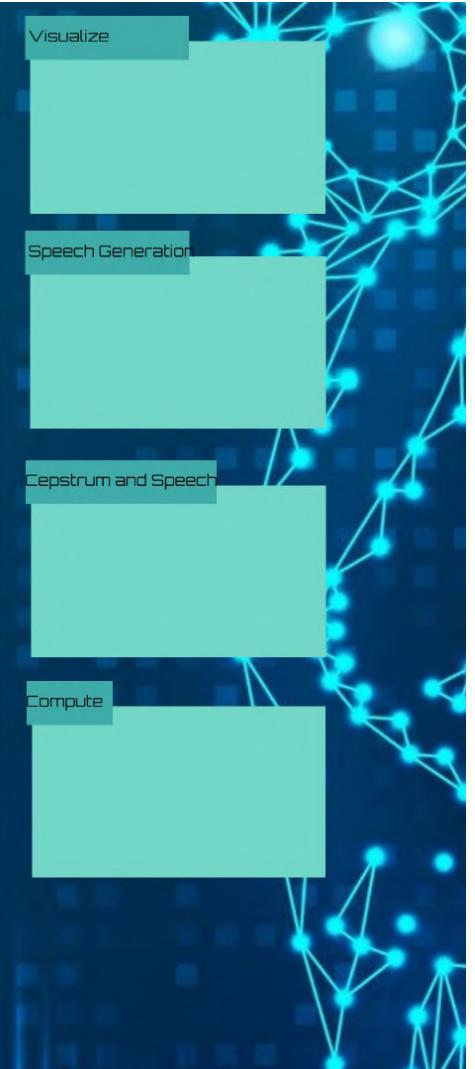
$$C(x(t)) = F^{-1} [\log(F[x(t)])]$$

Spectrum

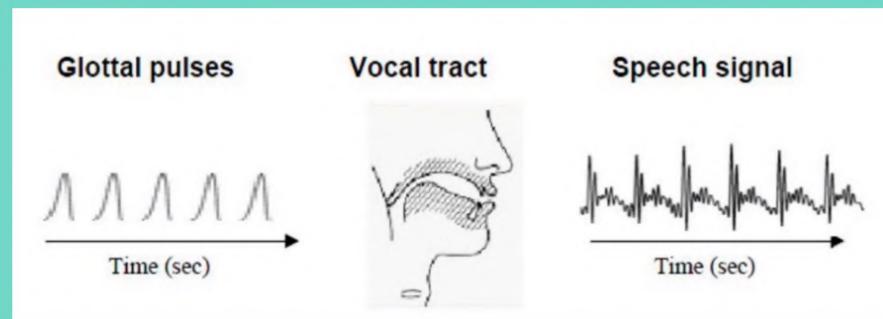
Log spectrum

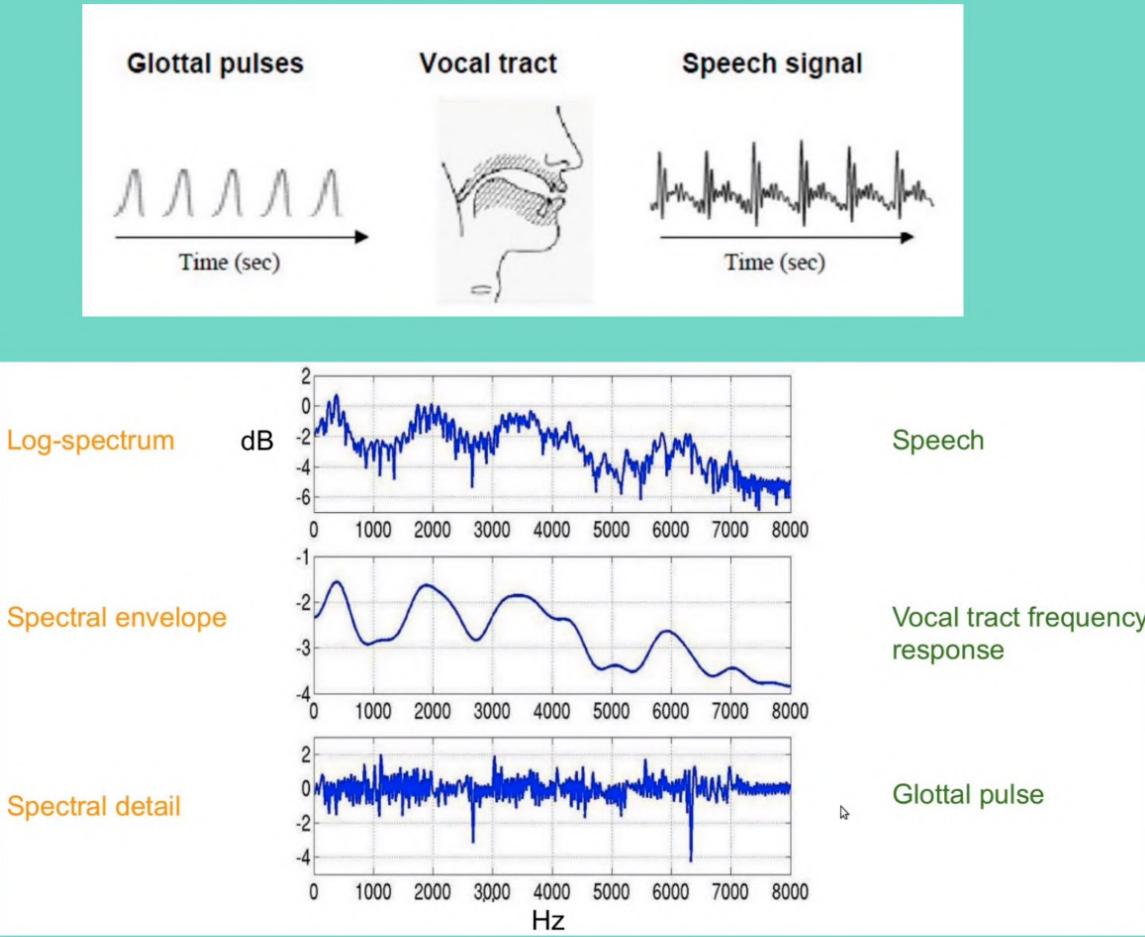
Cepstrum

The diagram illustrates the mathematical formula for computing the Cepstrum. It shows a green box containing the inverse Fourier transform (F^{-1}) of a log spectrum. This log spectrum is contained within an orange box, which is itself derived from the time-domain signal $x(t)$, indicated by a pink box.









Mel-Frequency Cepstrum Coefficients (MFCCs)

Cepstrum

Time-domain signal

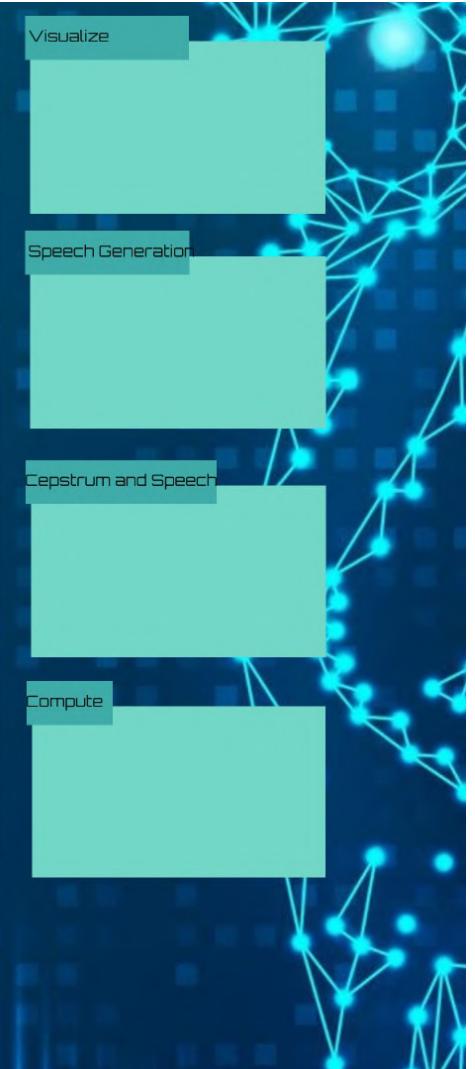
$$C(x(t)) = F^{-1} [\log(F[x(t)])]$$

Spectrum

Log spectrum

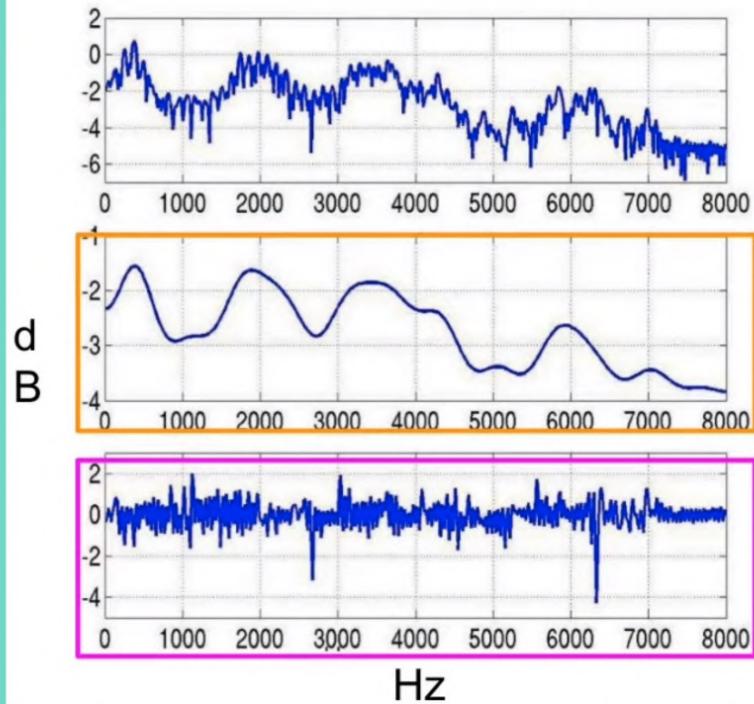
Cepstrum

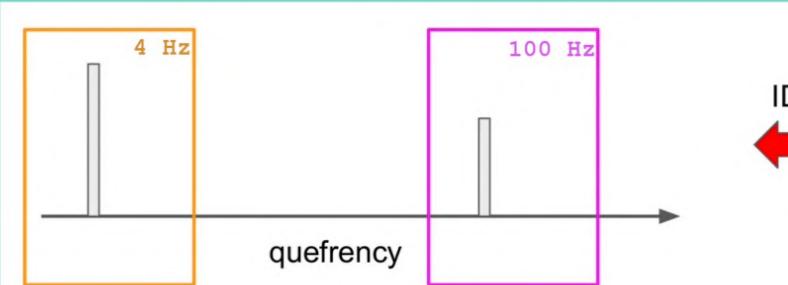
The diagram illustrates the mathematical formula for the Cepstrum. It shows a time-domain signal $x(t)$ being processed through a Fourier Transform F to produce a spectrum. This spectrum is then converted into a log spectrum. Finally, the inverse Fourier Transform F^{-1} is applied to the log spectrum to yield the Cepstrum. The components are color-coded: the time-domain signal is pink, the spectrum is orange, the log spectrum is blue, and the final cepstrum is green.





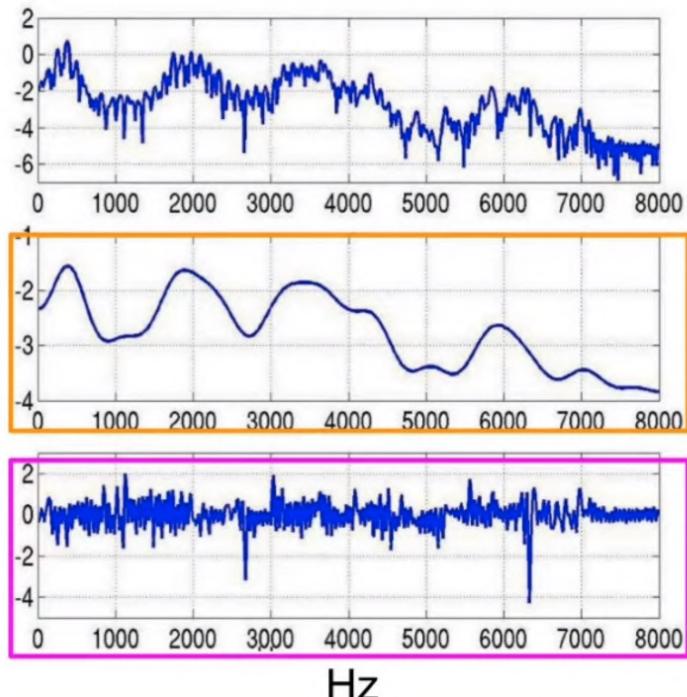
$$\log(X(t)) = \log(E(t)) + \log(H(t))$$





$$X(t) = \boxed{E(t)} \cdot \boxed{H(t)}$$

$$\log(X(t)) = \log(E(t)) + \log(H(t))$$



Mel-Frequency Cepstrum Coefficients (MFCCs)

Cepstrum

Time-domain signal

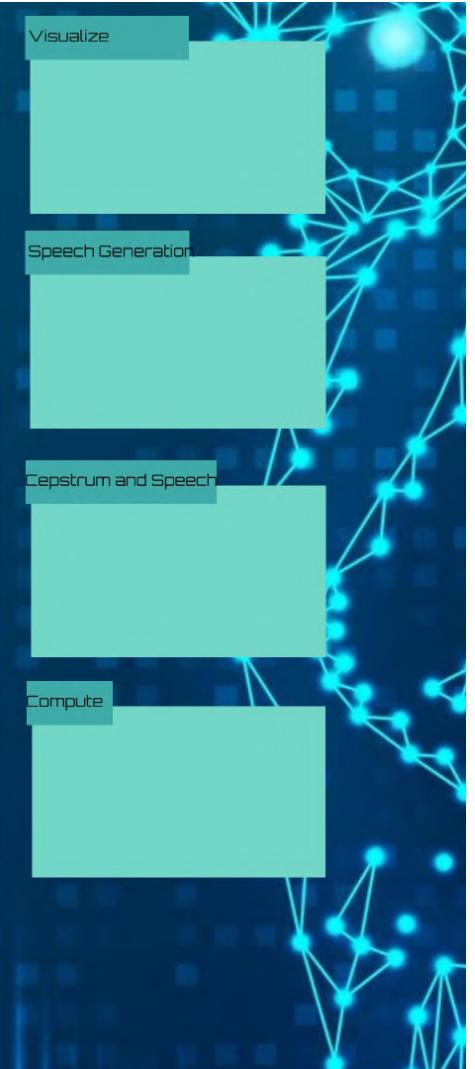
$$C(x(t)) = F^{-1} [\log(F[x(t)])]$$

Spectrum

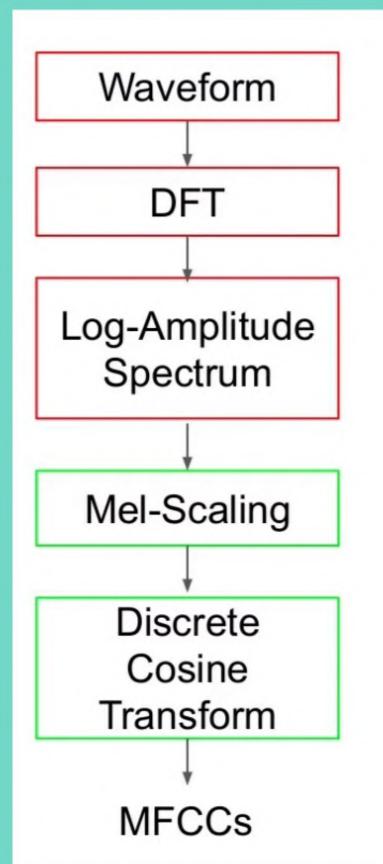
Log spectrum

Cepstrum

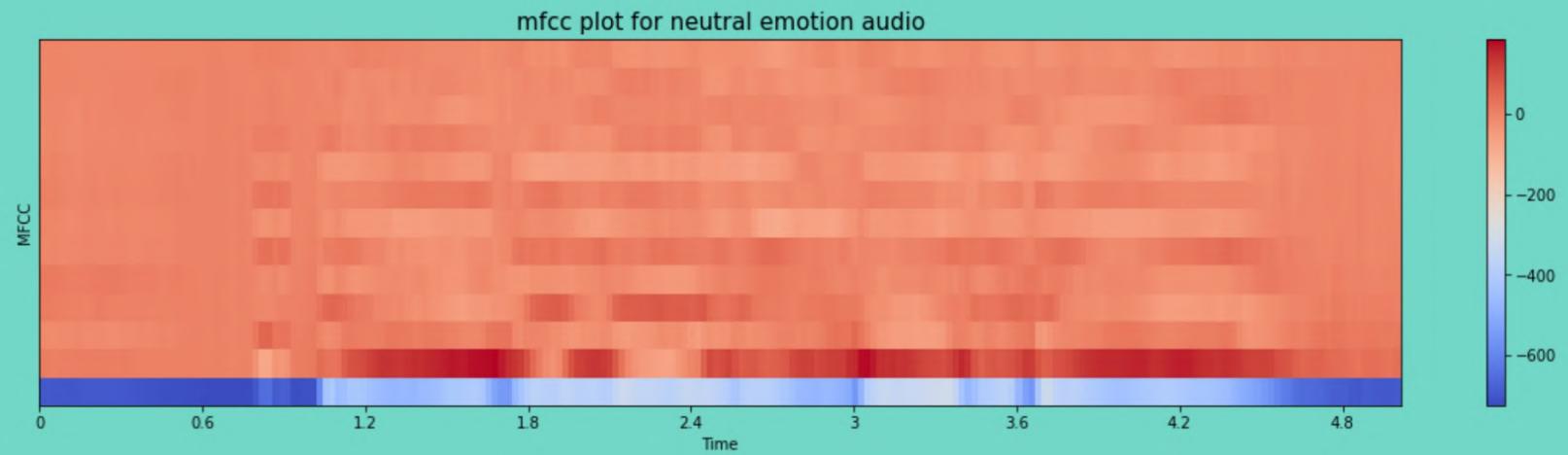
The diagram shows the mathematical formula for calculating the Cepstrum. A pink box highlights the term $x(t)$ in the argument of the Fourier transform F . Another pink box highlights the term $x(t)$ in the argument of the logarithm \log . A green box encloses the entire expression $F^{-1} [\log(F[x(t)])]$. Inside this green box, an orange box highlights the term $F[x(t)]$ in the argument of the logarithm. Labels indicate the 'Time-domain signal' at the top left, 'Spectrum' above the inner orange box, 'Log spectrum' below it, and 'Cepstrum' at the bottom.





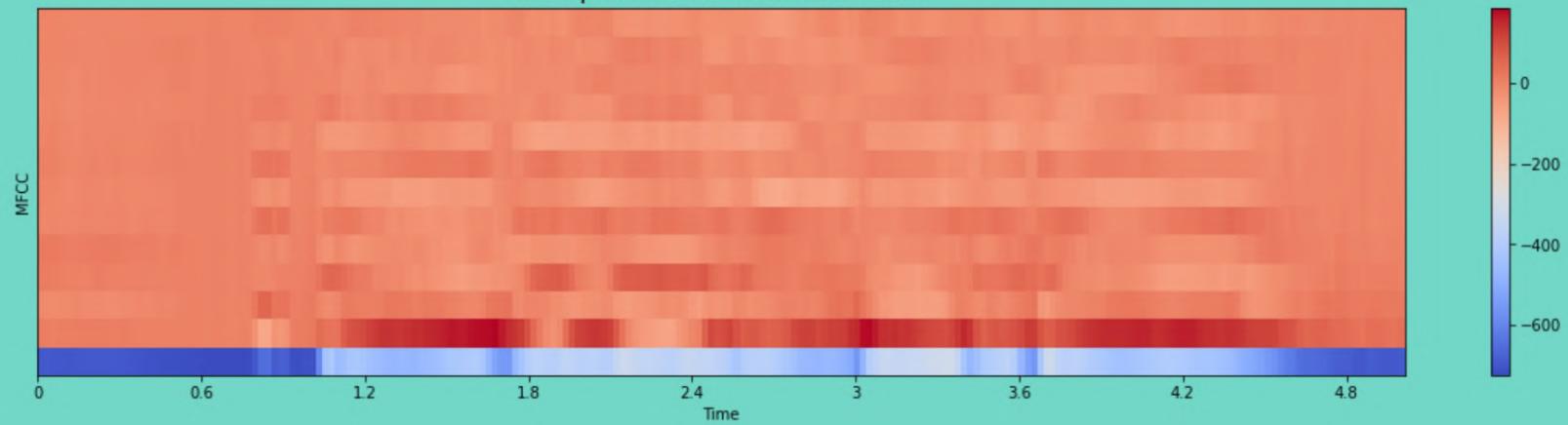






X

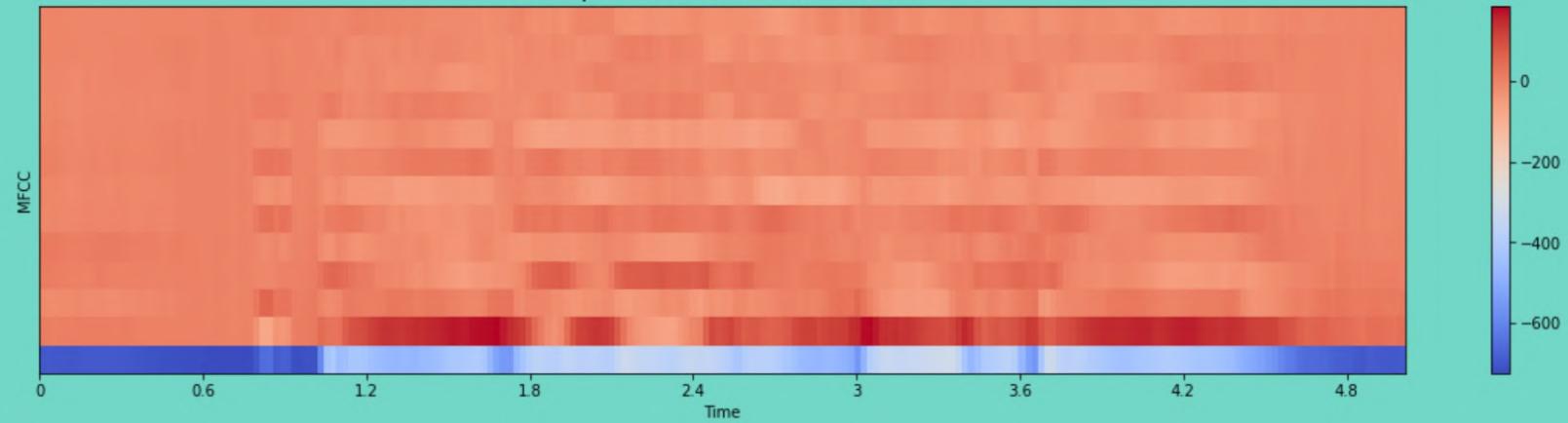
mfcc plot for neutral emotion audio



X

(1344, 448, 13)

mfcc plot for neutral emotion audio



X

(1344, 448, 13)

Mel-Frequency Cepstrum Coefficients (MFCCs)

Cepstrum

Time-domain signal

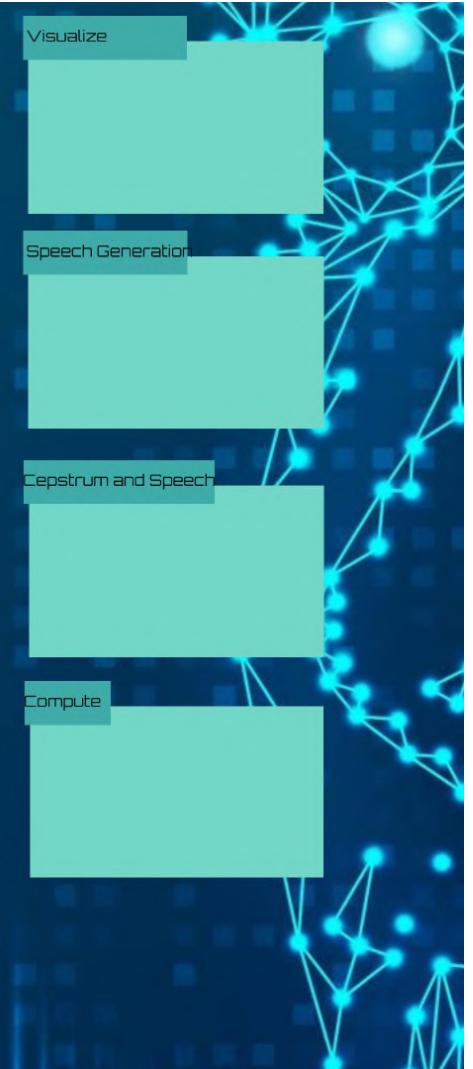
$$C(x(t)) = F^{-1} [\log(F[x(t)])]$$

Spectrum

Log spectrum

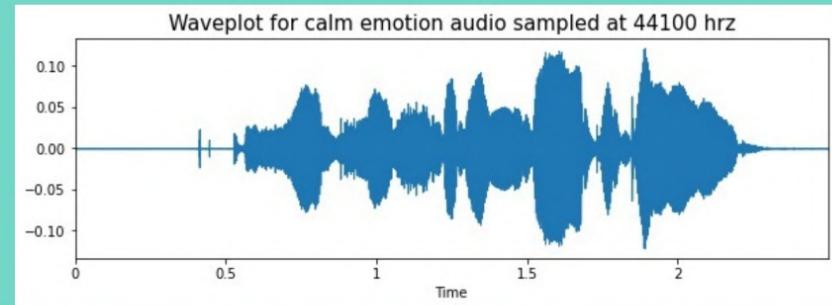
Cepstrum

The diagram illustrates the mathematical formula for computing the Cepstrum. It shows a green box containing the inverse Fourier transform (F^{-1}) of a log spectrum. This log spectrum is contained within an orange box, which is itself derived from the time-domain signal $x(t)$, indicated by a pink box.

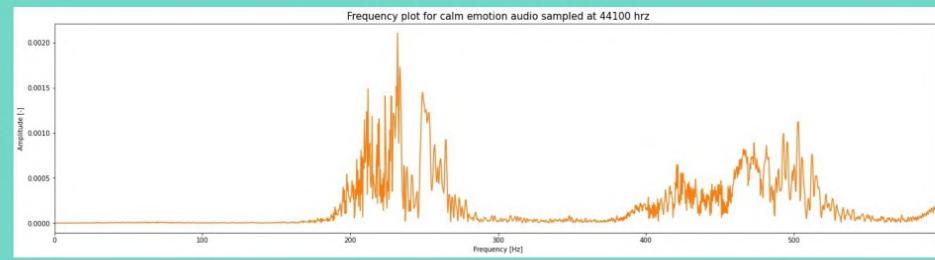


Feature extraction

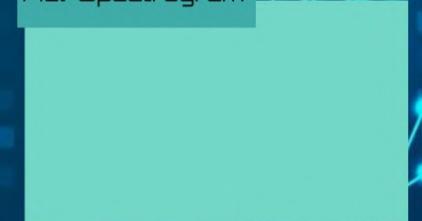
- Time-Domain
 - RMSE, ZCR



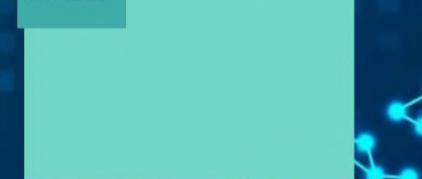
- Frequency Domain
 - Spectral Centroid



Mel-Spectrogram



MFCC



RAVDESS

Neutral, calm, happy, sad, angry, fearful, disgust, surprise

Plots

Processing

Feature extraction

Introduction

Importance:

- The most natural way to express ourselves
- Influence over judgement
 - "go home!"
- Human-computer interaction (HCI)
 - Personal Assistants
 - Next Generation AI

The Evolution of Human Vocal Emotion
Gregory A. Bryant¹
First Published June 24, 2020 | Research Article |  <https://doi.org/10.1177/1754073920930791>
[Article information](#) ▾

Applications

Data set



Speech emotion recognition

Hamed Vaheb, Elnaz Khaveh

Introduction

Models

Conclusion



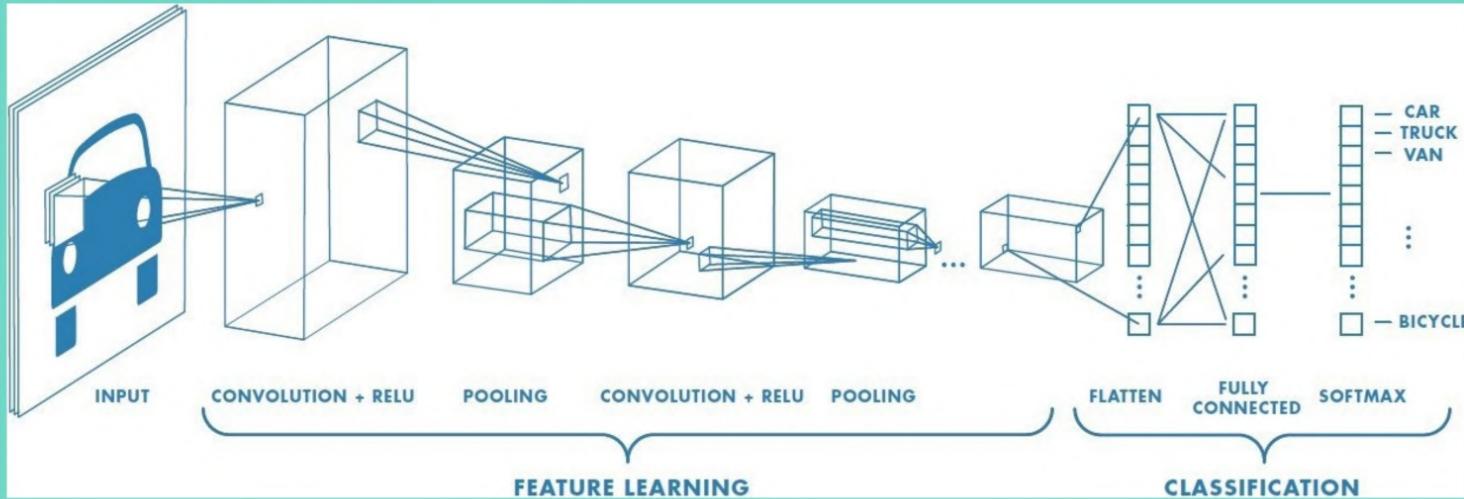
Models

CNN

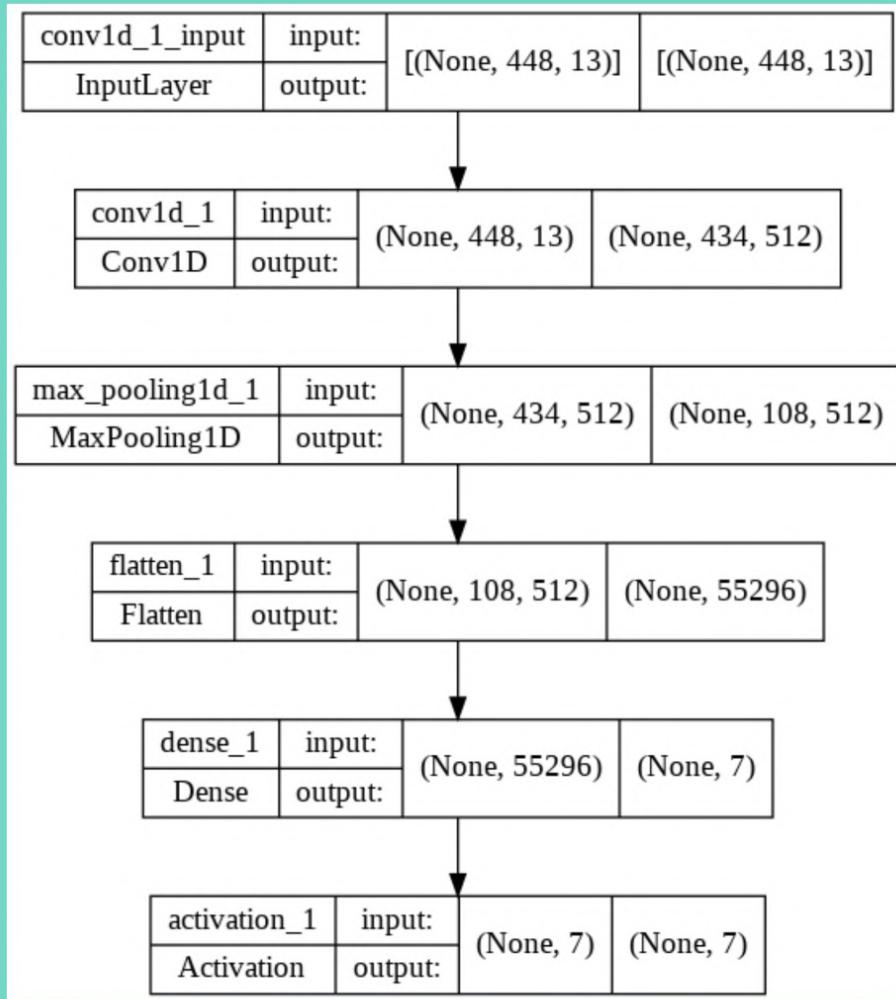
RNN

Results

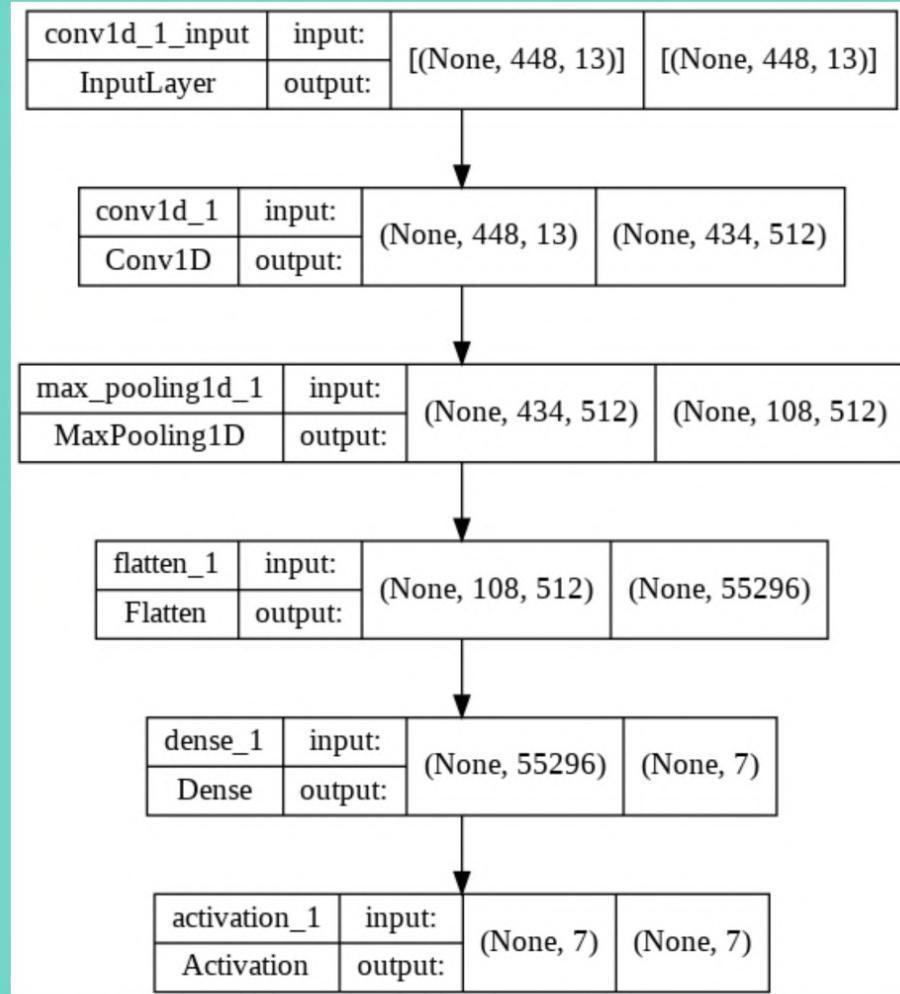
CNN



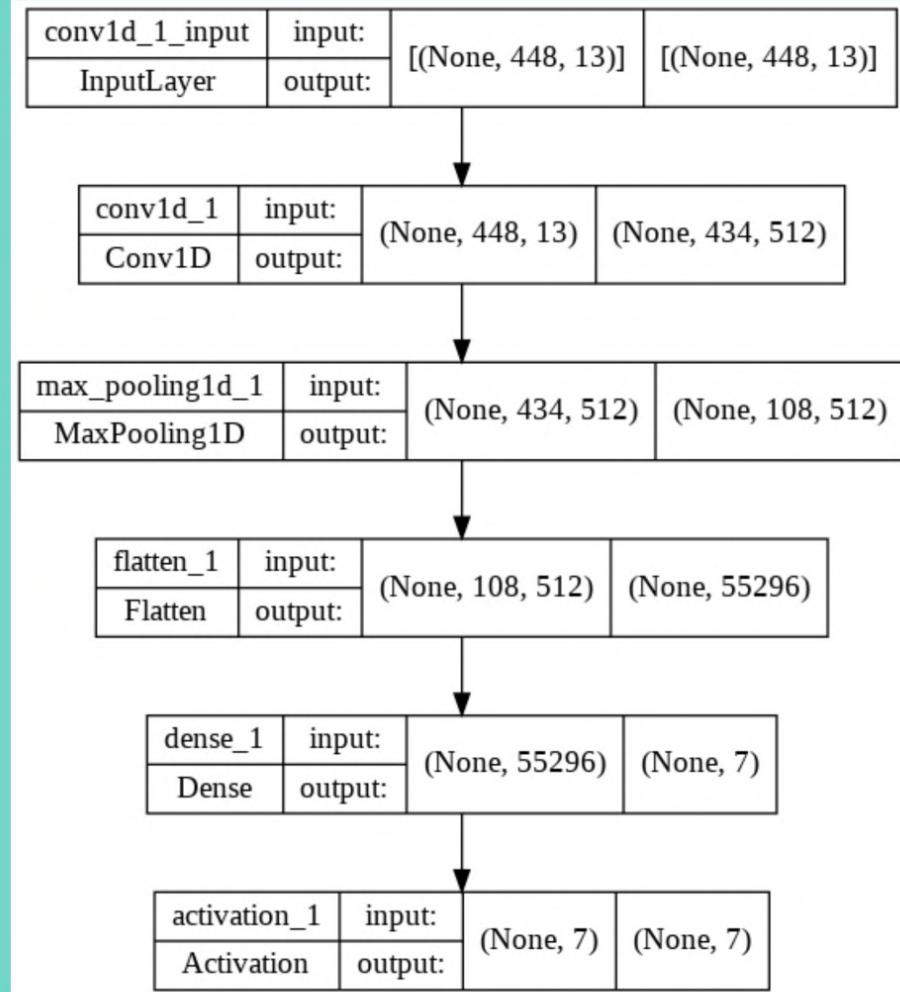
architecture



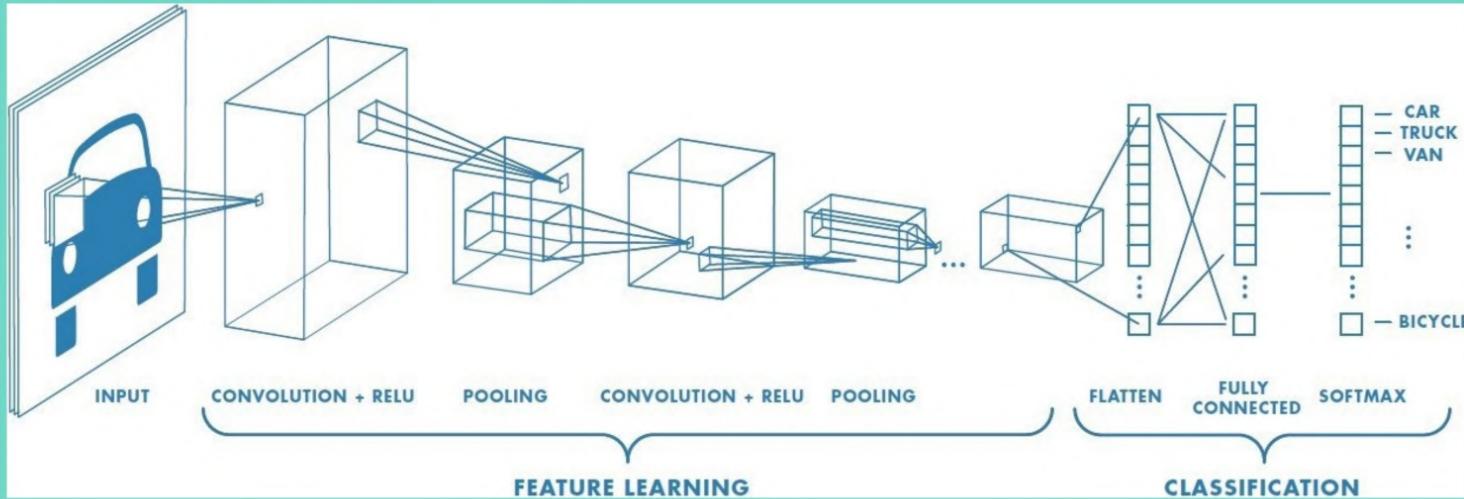
X



X
(1344, 448, 13)



CNN



architecture

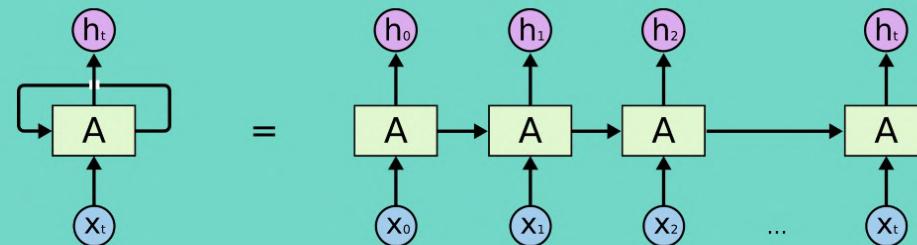
Models

CNN

RNN

Results

RNN



GRU

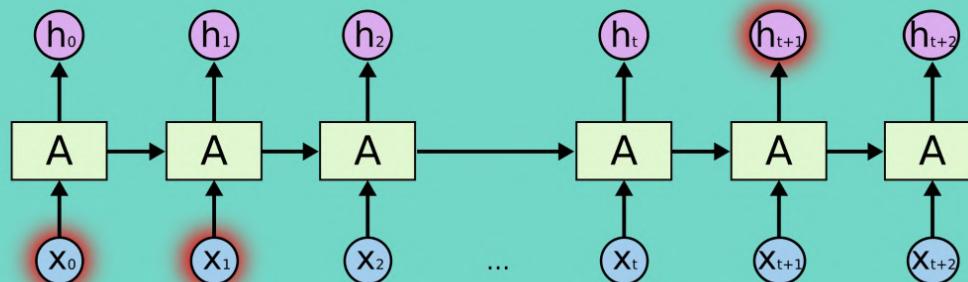
architecture

RNN

GRU

architecture

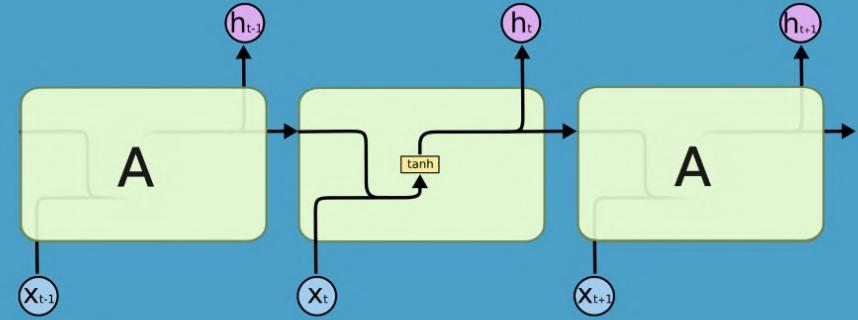
RNN

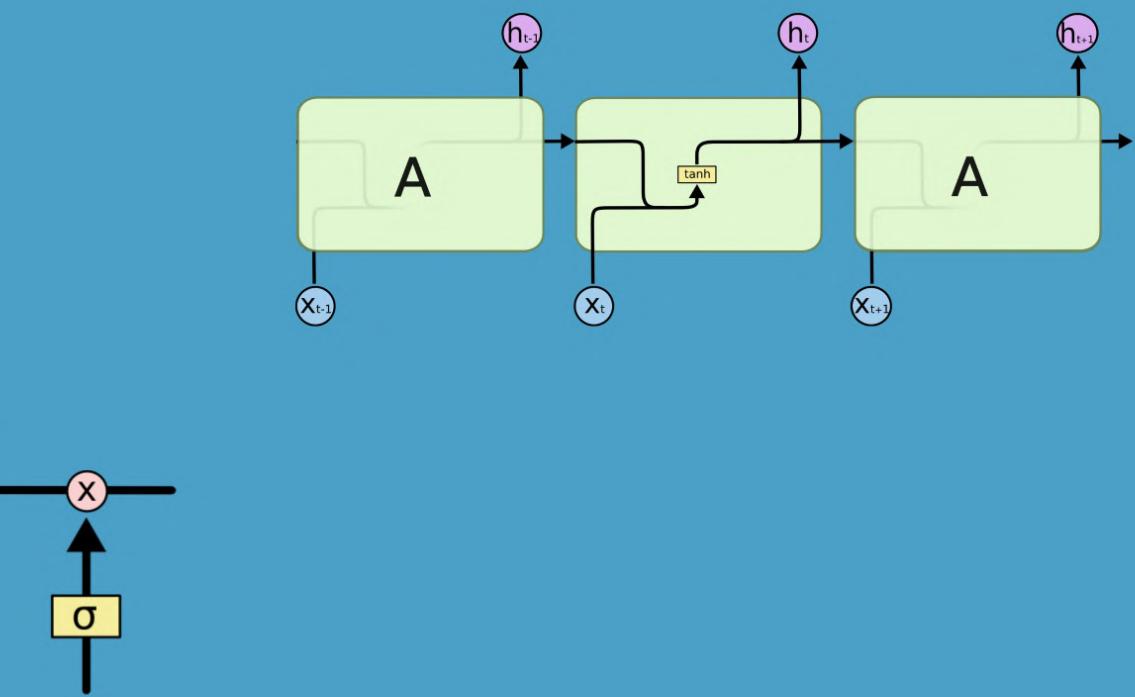


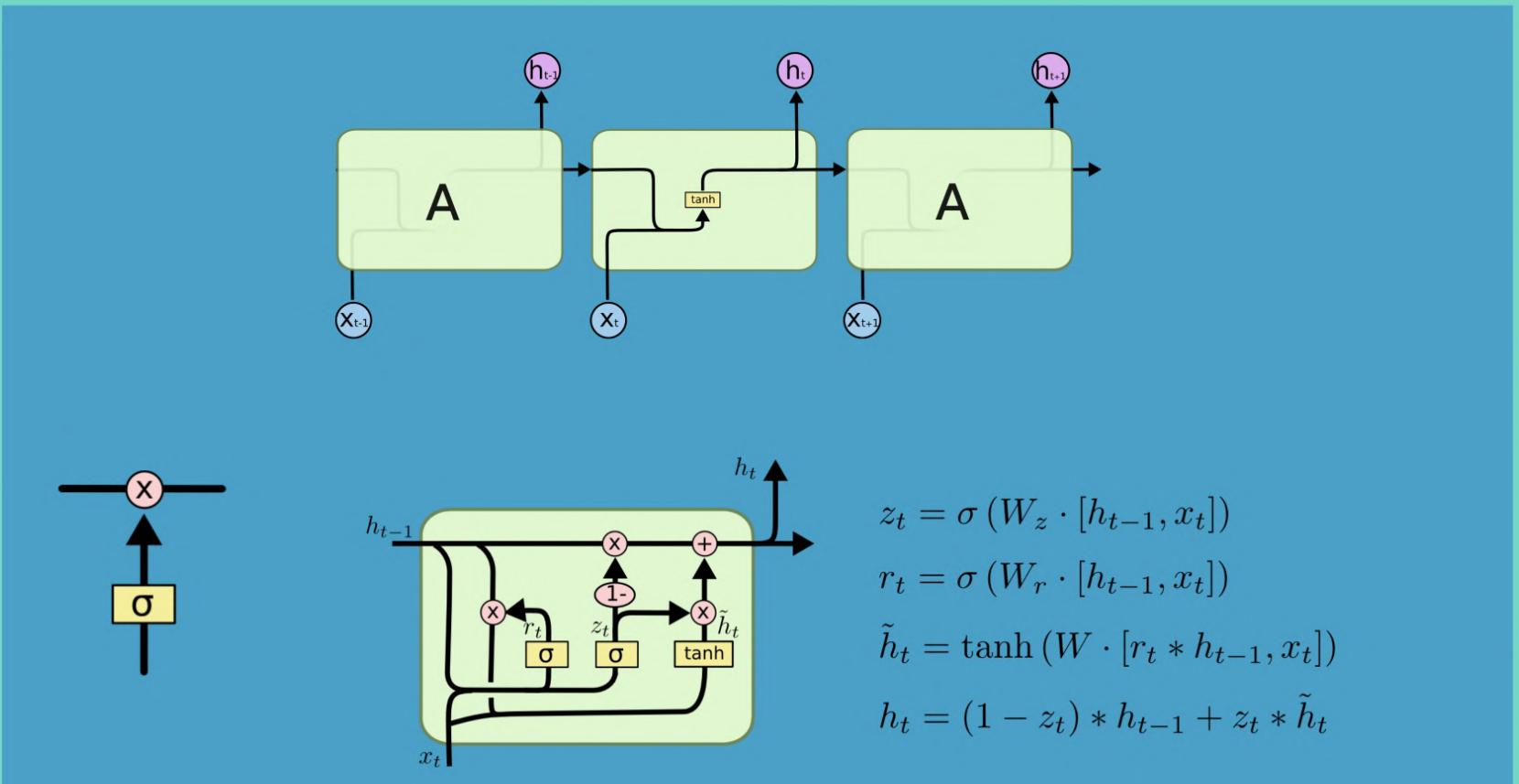
GRU

architecture

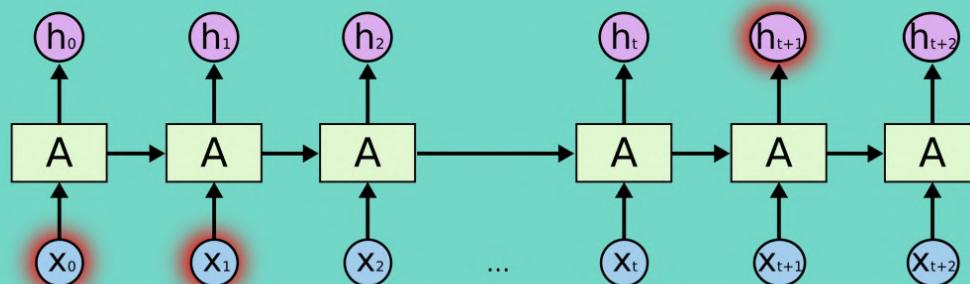






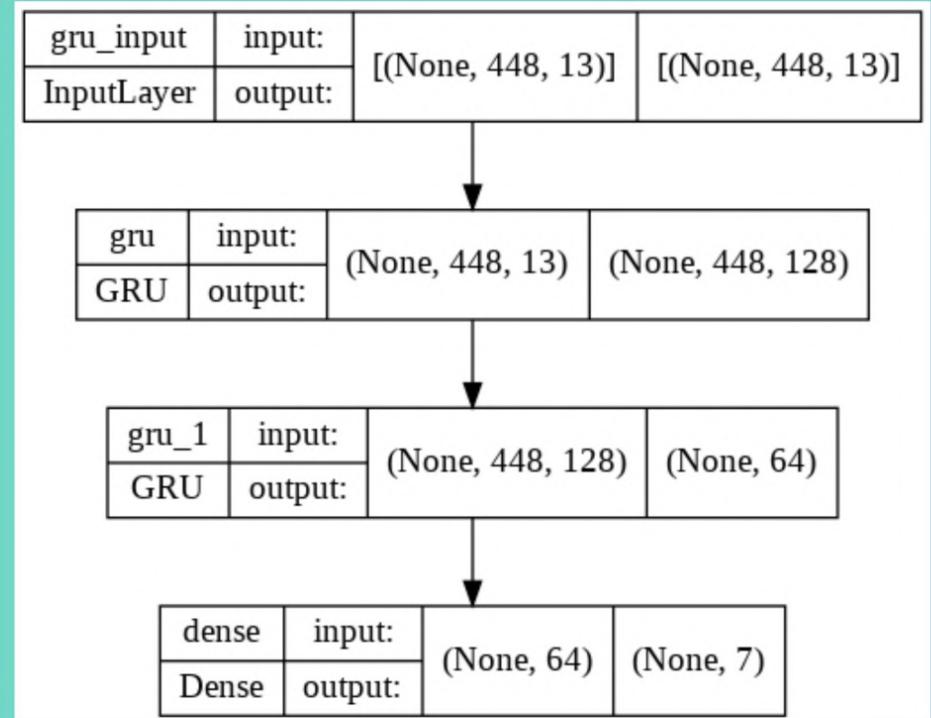


RNN

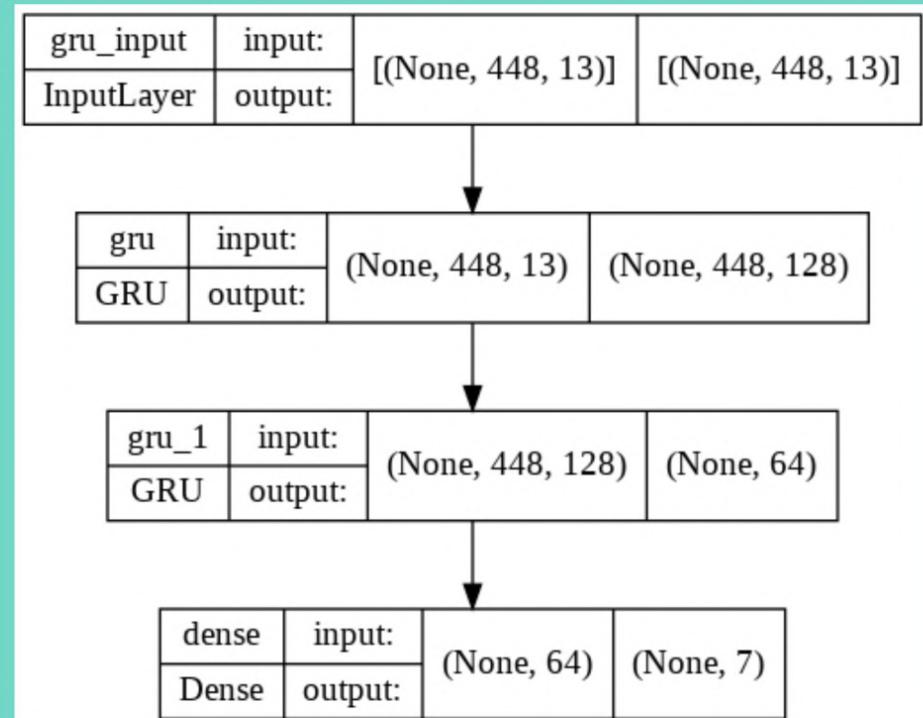


GRU

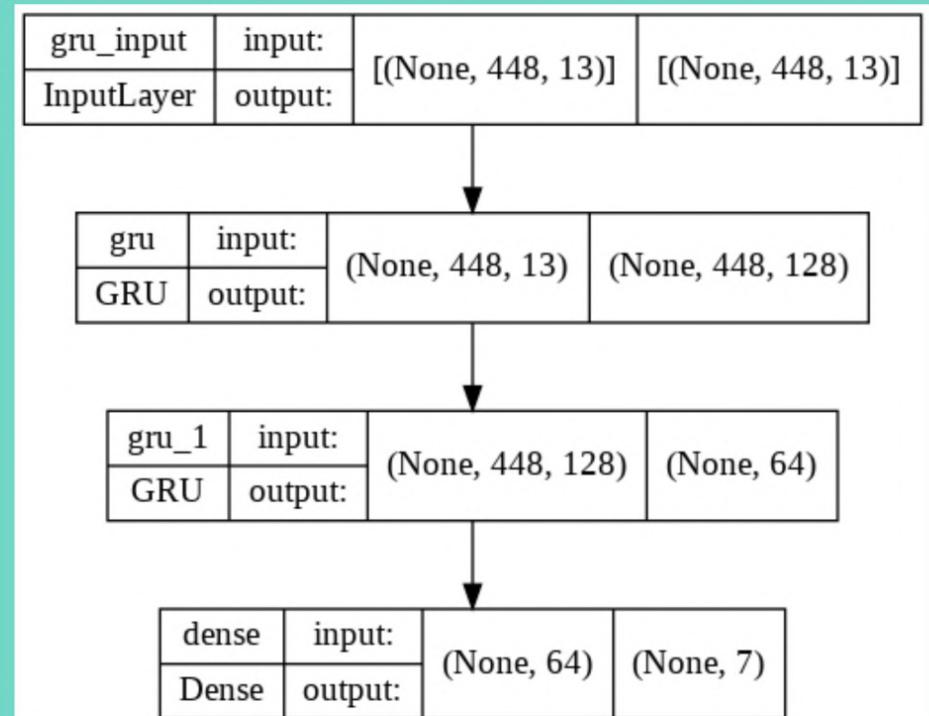
architecture



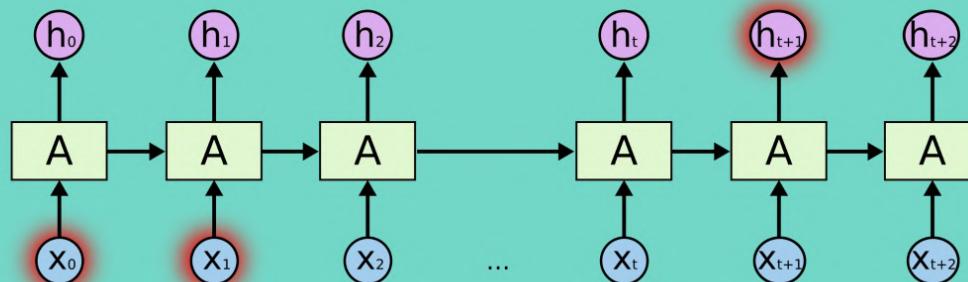
X



X
(1344, 448, 13)



RNN



GRU

architecture

Models

CNN

RNN

Results

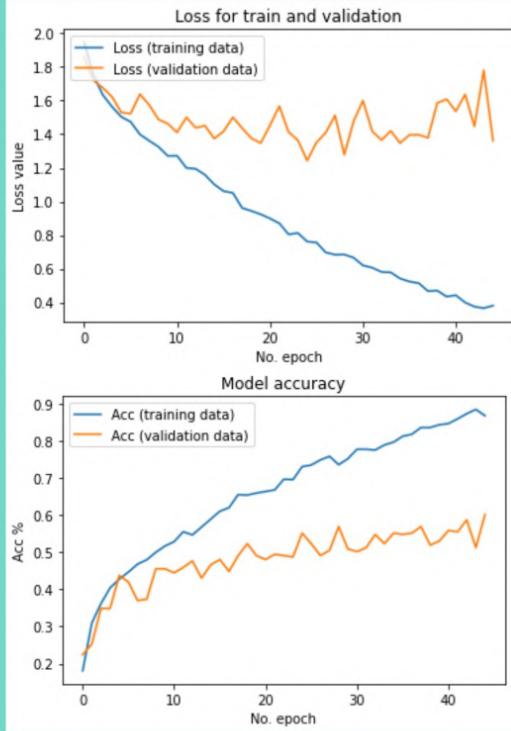
Compare the results

Loss-accuracy

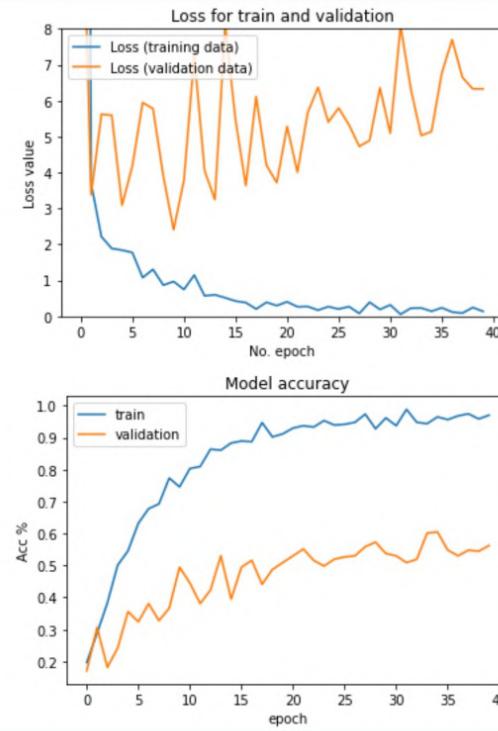
Table

Confusion matrix

RNN



CNN



Compare the results

Loss-accuracy

Table

Confusion matrix

	RNN	CNN
Accuracy	60%	61%
Precision	63%	64%
Recall	62%	62%
F1 score	61%	60%
loss	1.39	6.58

The training time for RNN: 18 minutes

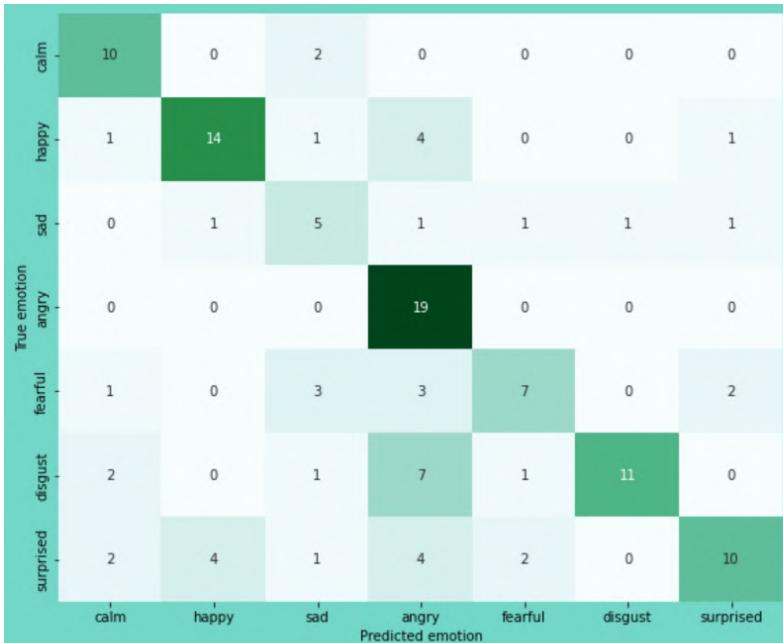
The training time for CNN: 3 minutes

Compare the results

Loss-accuracy

Table

Confusion matrix



CNN



RNN

Compare the results

Loss-accuracy

Table

Confusion matrix

Models

CNN

RNN

Results

Speech emotion recognition

Hamed Vaheb, Elnaz Khaveh

Introduction

Models

Conclusion



Conclusion

- CNN and RNN perform almost the same
- Time complexity of CNN is slightly better
- MFCC was the most determinant feature

Future work

Future work

- Real-time
- continuous representation of emotions
- Small data set
- hyper-parameter
- ensemble

Future work

- Real-time
- continuous representation of emotions
- Small data set
- hyper-parameter
- ensemble

[https://github.com/
berserkhmdvhb/Project-SER](https://github.com/berserkhmdvhb/Project-SER)

Future work

- Real-time
- continuous representation of emotions
- Small data set
- hyper-parameter
- ensemble

Thanks for your attention !

Future work

- Real-time
- continuous representation of emotions
- Small data set
- hyper-parameter
- ensemble

Conclusion

- CNN and RNN perform almost the same
- Time complexity of CNN is slightly better
- MFCC was the most determinant feature

Future work

Speech emotion recognition

Hamed Vaheb, Elnaz Khaveh

Introduction

Models

Conclusion

