

Audio Summarization with Audio Features and Probability Distribution Divergence



Carlos-Emiliano González-Gallardo¹, Romain Deveaud¹, Eric SanJuan¹
and Juan-Manuel Torres-Moreno^{1,2}

¹LIA - Avignon Université, Avignon, France

²Département de GIGL, Polytechnique Montréal, Montréal, Canada

{carlos-emiliano.gonzalez-gallardo, eric.sanjuan, juan-manuel.torres}@univ-avignon.fr
romain.deveaud@gmail.com

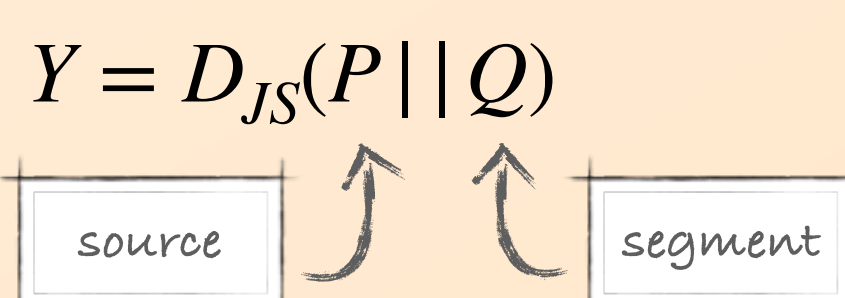


Introduction

- Audio summarization of massive online multimedia resources.
- Facilitate the understanding.
- Extractive audio summarization approaches:
 - textual methods
 - audio features
 - hybrid
- Our proposition:
 - Represent the information within the text in terms of its audio features.
 - Hybrid during training phase; text independent during summary creation.

Probability Distribution Divergence for Audio Summarization

- Training phase:
 - Audio features (275 MFCC + 2) & textual information.
 - Mapping between audio features X and an informativeness value Y .



$$Y = D_{JS}(P || Q)$$

- Audio summary creation:
 - Audio features (275 MFCC + 2) & textual information.
 - A score S_{Q_i} is computed to rank the pertinence of each segment $Q_1 \dots Q_k$
 - Segments with higher S_{Q_i} scores are chosen until θ is reached.

Training phase (Informativeness model)

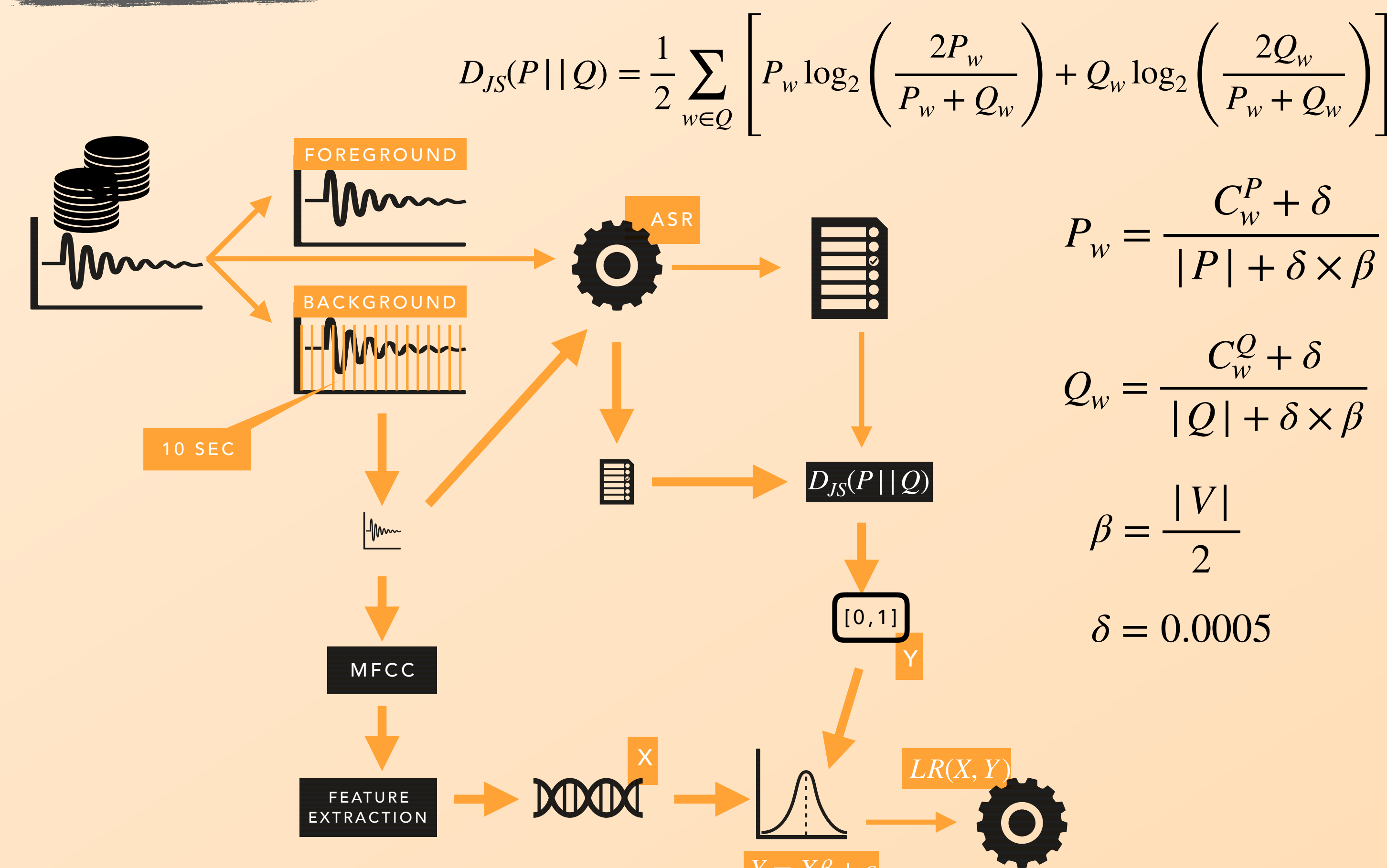


Fig. 1: Informativeness model scheme

Audio summary creation

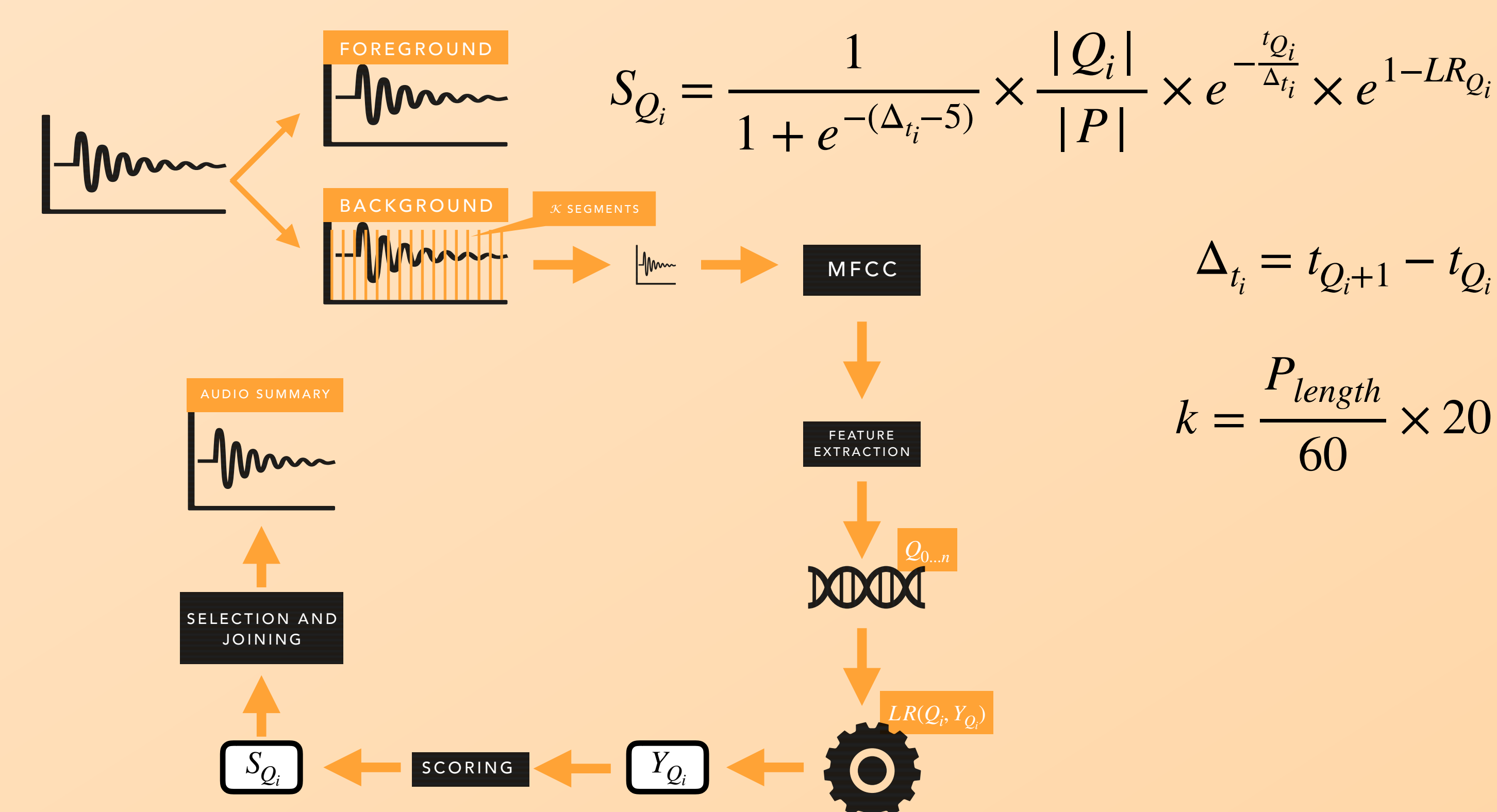


Fig. 2: Summary creation scheme

Experiments & Results

Training of Informativeness model

- 5,989 audio broadcasts (~310 hours) in French, English and Arabic.
- Automatically obtained transcripts are treated with a stemming process.
- A linear least squares regression model is trained to map the audio features X of ~111,600 training samples into an informativeness score Y .

Results

Score	Explanation	Sample	Length	Segments	Full Score	Average Score
1	Not informative	7	5m23s	8	3.20	3.75
2	Quite informative	6	9m45s	30	4.00	2.49
3	Half informative	5	8m47s	22	4.67	3.68
4	Mostly informative	4	1m42s	5	3.60	2.95
5	Full informative	3	2m47s	5	3.80	3.76
6	Quite informative	2	5m21s	13	3.50	2.78
7	Half informative	1	3m19s	8	4.20	2.90
8	Mostly informative	8	6m24s	20	3.75	2.84
9	Quite informative	9	7m35s	18	3.75	3.19
10	Not informative	10	2m01s	4	2.75	2.63

Table 1: Evaluation scale

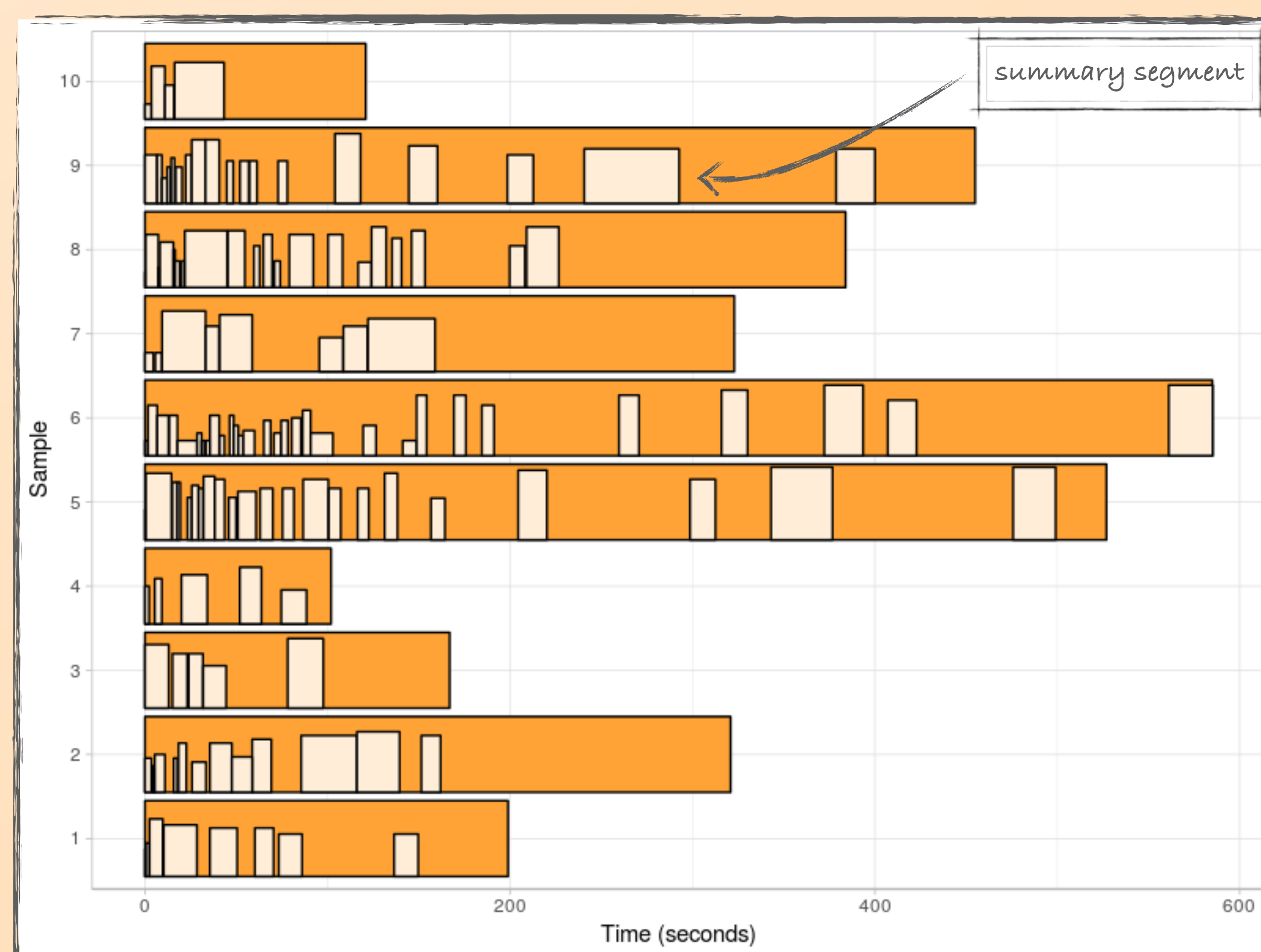


Fig. 3: Graphical representation of audio summarization performance

Conclusions & Perspectives

- Mapping informativeness from Mel-frequency Cepstral Coefficients (MFCC) features to their corresponding Jensen-Shannon (JS) divergence help to select those segments which are more relevant to the audio summary.
- Original approach; hybrid during training phase but text independent while creating summaries.
- It manages to generate at least half informative extractive summaries.
- Not a clear correlation between the quality of a summary and the quality of its parts.
- Future work will consider bigger evaluation datasets as well as French and Arabic summarization.

References

- Louis, A., Nenkova, A.: Automatically evaluating content selection in summarization without human models. In proceedings of EMNLP'09, (2009).
- Torres-Moreno, et al.: Summary evaluation without references. Polibits 42, 13-19 (2010).
- Raffii, Z., Pardo, B.: Music/voice separation using the similarity matrix. In proceedings of ISMIR'12, (2012).
- Jouvet, D., et al.: Adaptation of speech recognition vocabularies for improved transcription of youtube videos. ISGA 1(1), 1-9 (2018).

Acknowledgments

We would like to acknowledge the support of CHIST-ERA for funding this work through the Access Multilingual Information opinionS (AMIS) project (France - Europe).