



wav2tok: Deep Sequence Tokenizer for Audio Retrieval



Adhiraj Banerjee



Vipul Arora

Department of Electrical Engineering
Indian Institute of Technology Kanpur, India
{adhiraj, vipular}@iitk.ac.in



Introduction

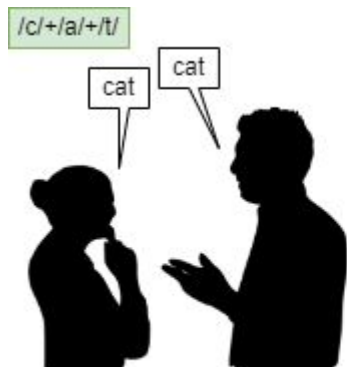
Sequence Retrieval aims at retrieving sequences similar to a query sequence, with the constraint that

an ordered alignment exists between the query and the target sequence.

Introduction

Sequence Retrieval aims at retrieving sequences similar to a query sequence, with the constraint that

an ordered alignment exists between the query and the target sequence.

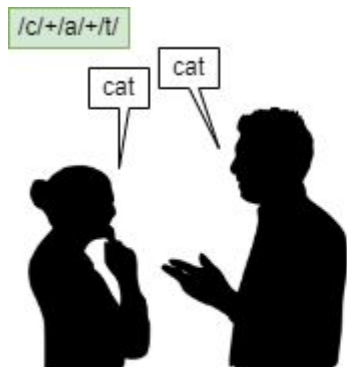


Speech Search

Introduction

Sequence Retrieval aims at retrieving sequences similar to a query sequence, with the constraint that

an ordered alignment exists between the query and the target sequence.



Speech Search

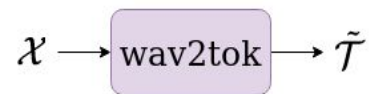


Music Search

Introduction

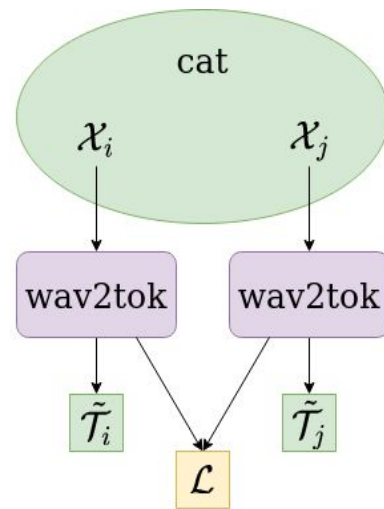
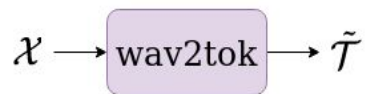


A model mapping audio
 \mathcal{X} to discrete tokens $\tilde{\mathcal{T}}$



Introduction

A model mapping audio \mathcal{X} to discrete tokens $\tilde{\mathcal{T}}$

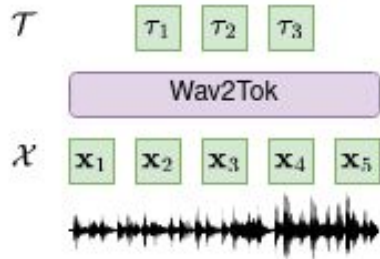


Model learns the tokens un-supervised from pairs of similar audio

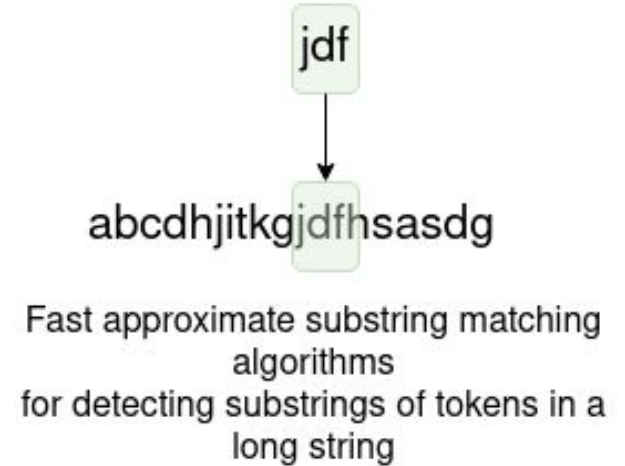
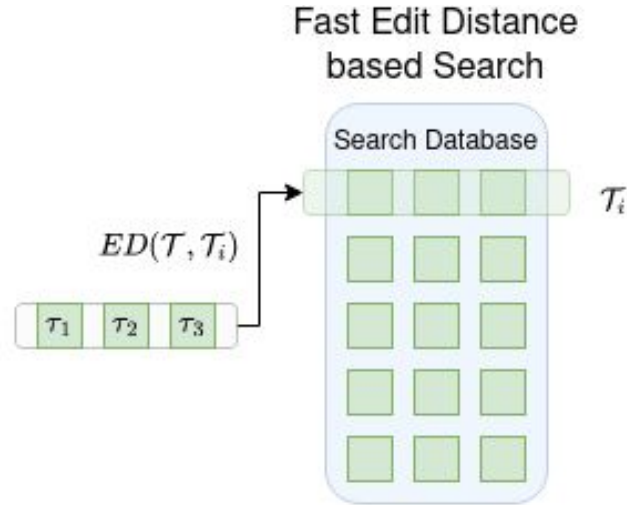
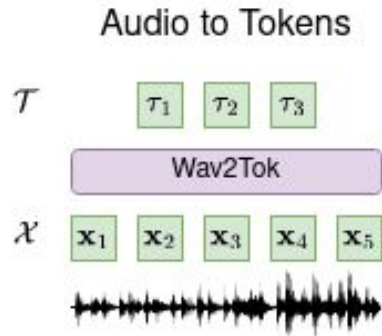


Motivation

Audio to Tokens



Motivation



Motivation



Preserving Languages with no written form

Bo, (Tribal Language in Andaman
Islands, extinct February 2010) [1]



Khoisan, (Kalahari Desert,
Tanzania) [2]



Baghati (Solan, Himachal Pradesh)
[1]



Sfyria (Greece) [2]



[1] <https://www.tribuneindia.com/news/archive/features/when-a-language-faces-extinction-584691>

[2] <https://peakd.com/top10languages/@calmbrain/10-extraordinary-languages-that-do-not-involve-speaking>

Motivation

Bird Language Acquisition

Distress
Call ?

Hungry?



Better monitoring for wildlife conservation

Preserving Languages with no written form

Bo, (Tribal Language in Andaman
Islands, extinct February 2010) [1]



Khoisan, (Kalahari Desert,
Tanzania) [2]



Baghati (Solan, Himachal Pradesh)
[1]



Sfyria (Greece) [2]



[1] <https://www.tribuneindia.com/news/archive/features/when-a-language-faces-extinction-584691>

[2] <https://peakd.com/top10languages/@calmbrain/10-extraordinary-languages-that-do-not-involve-speaking>

wav2tok: A Framework for learning discrete tokens



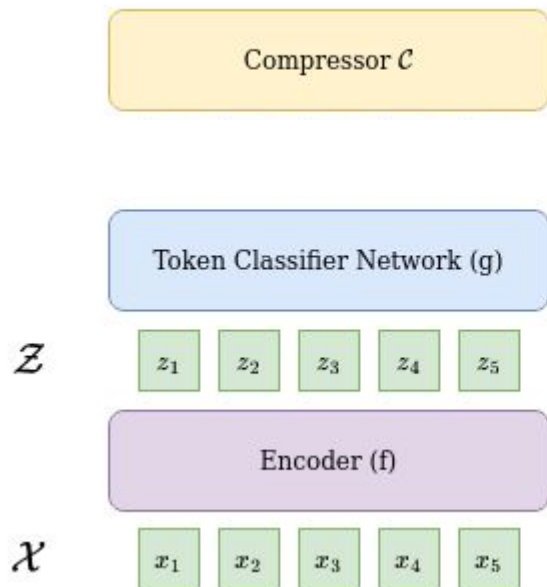
Compressor C

Token Classifier Network (g)

Encoder (f)

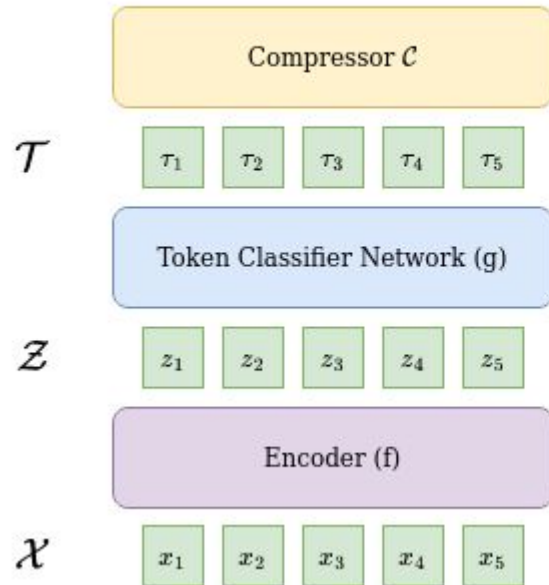


wav2tok: A Framework for learning discrete tokens



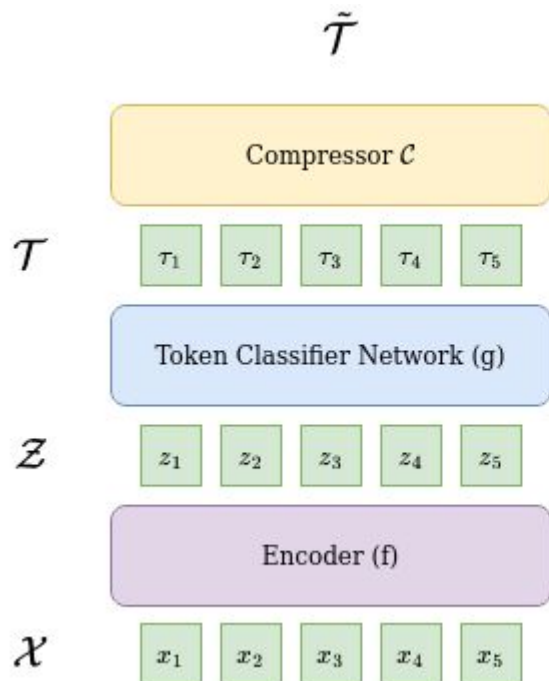


wav2tok: A Framework for learning discrete tokens



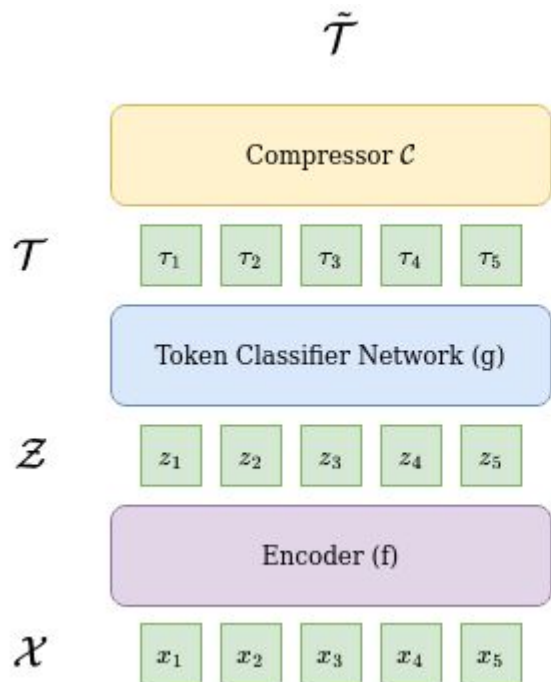


wav2tok: A Framework for learning discrete tokens





wav2tok: A Framework for learning discrete tokens



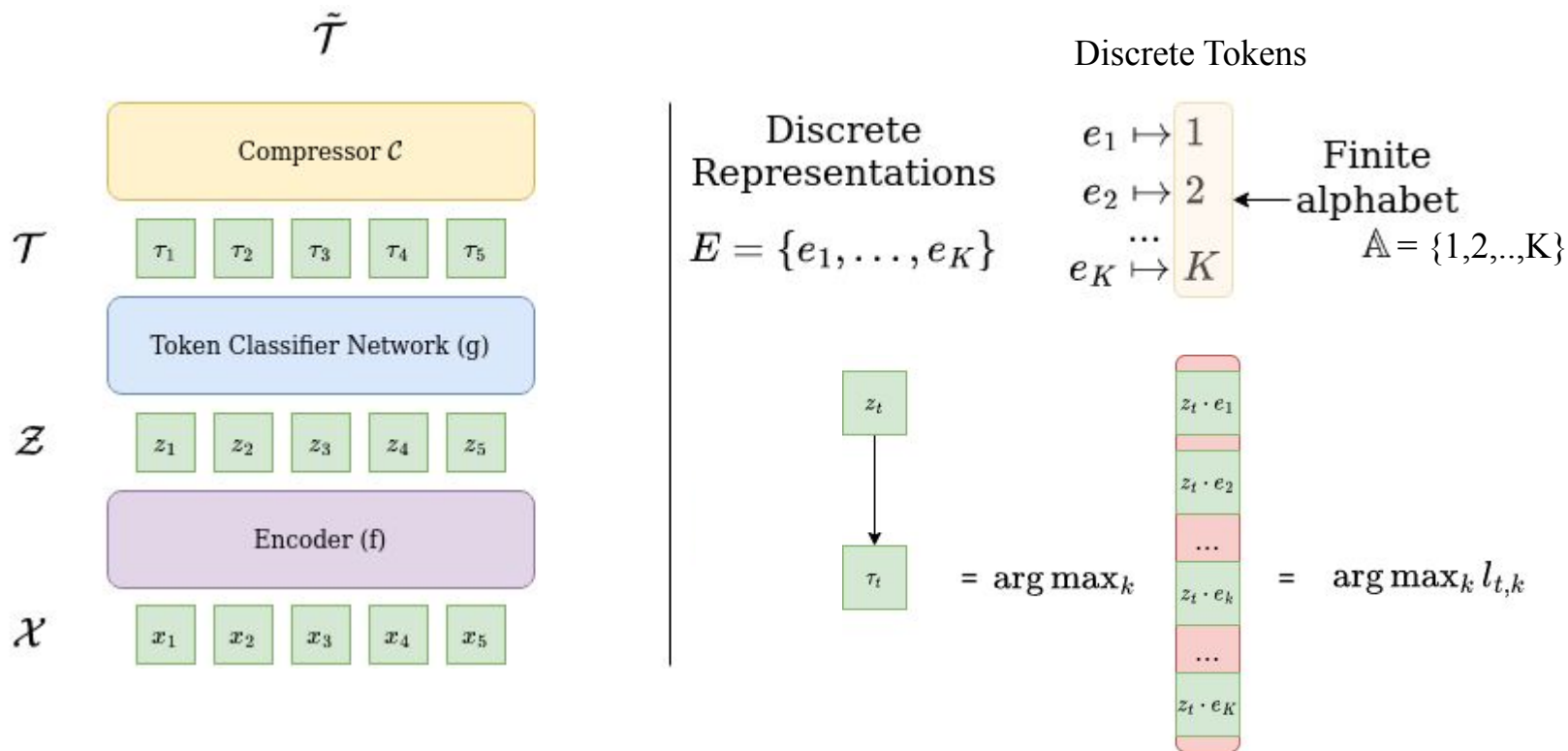
Discrete
Representations

$$E = \{e_1, \dots, e_K\}$$

Discrete Tokens

$$\begin{array}{l} e_1 \mapsto 1 \\ e_2 \mapsto 2 \\ \vdots \\ e_K \mapsto K \end{array} \leftarrow \text{Finite alphabet } \mathbb{A} = \{1, 2, \dots, K\}$$

wav2tok: A Framework for learning discrete tokens



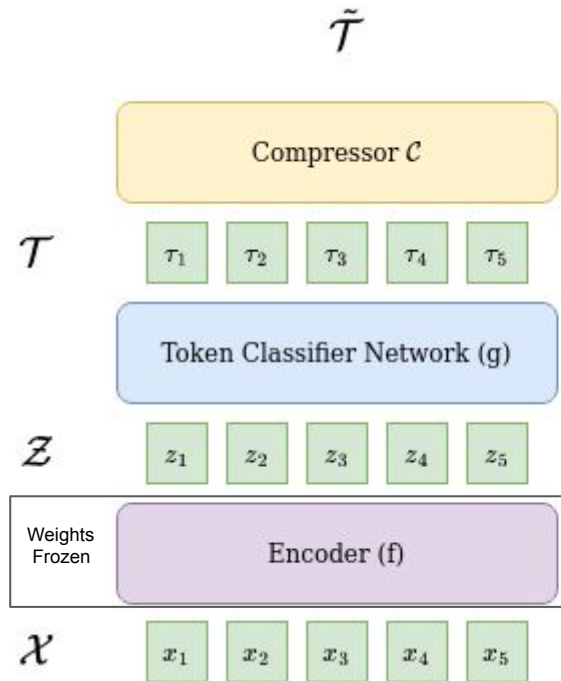


Training

Trained on pairs of similar sequences $(\mathcal{X}, \mathcal{X}')$

Expectation Maximization (EM) algorithm like training,

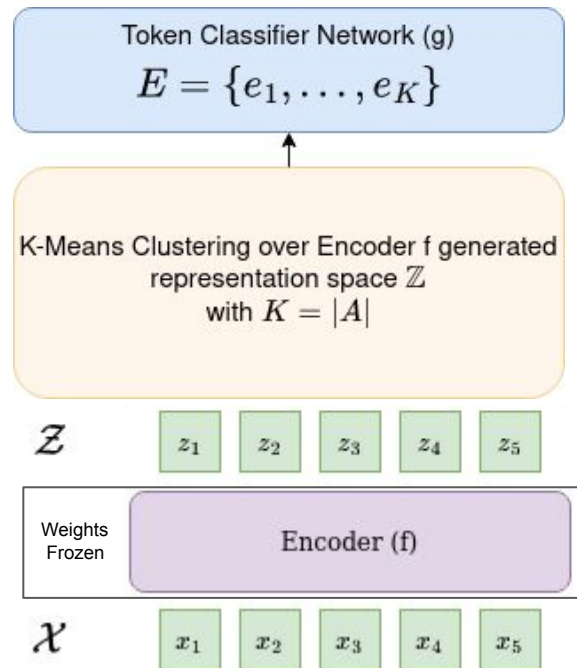
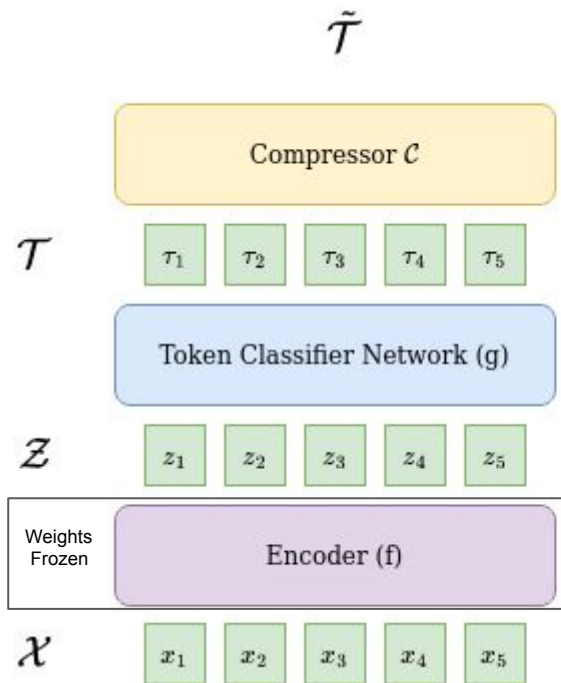
E-Step



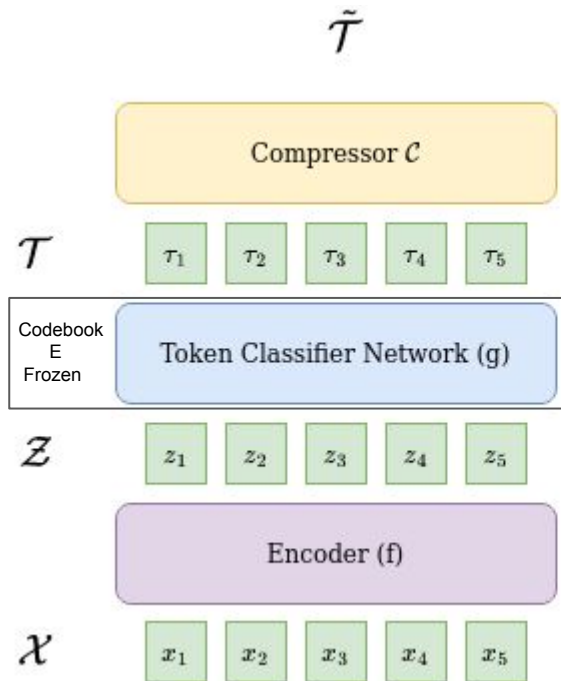
Token Classifier Network (g)

$$E = \{e_1, \dots, e_K\}$$

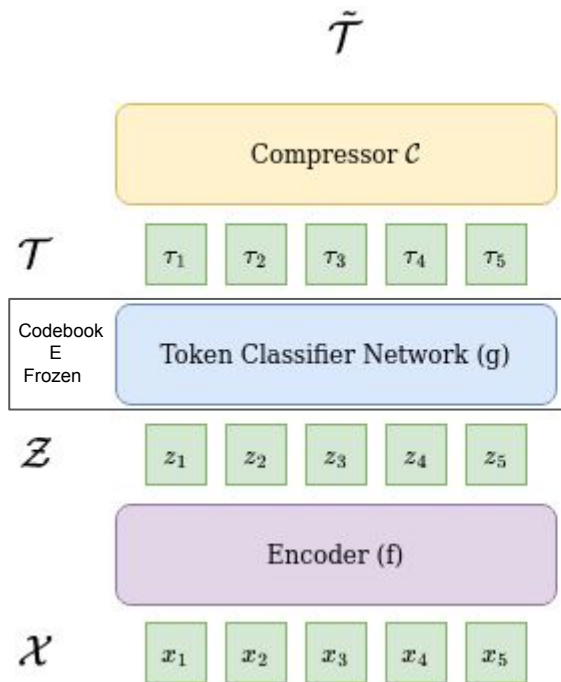
E-Step



M-Step



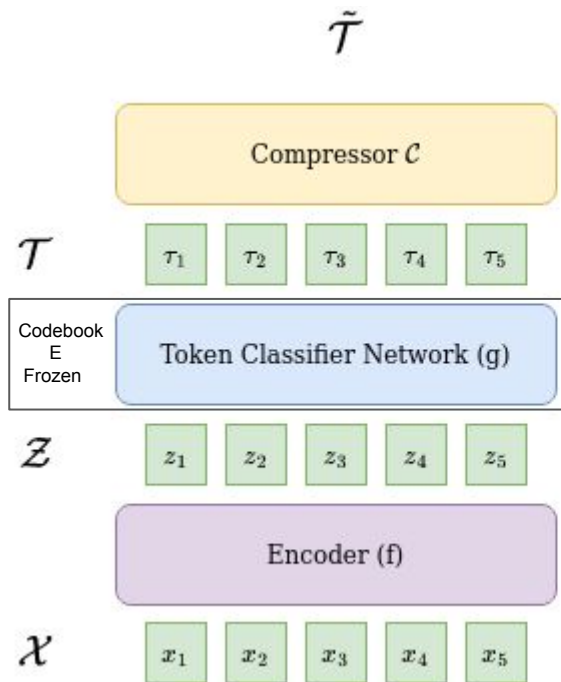
M-Step



$$\mathcal{X} \mapsto \tilde{\mathcal{T}}, \mathcal{X}' \mapsto \tilde{\mathcal{T}'}$$

$$\mathcal{L} = \mathcal{L}_m(\mathcal{X}, \mathcal{X}') + \alpha \mathcal{L}_{ctc}(\mathcal{X}, \tilde{\mathcal{T}'}) + \beta \mathcal{L}_{ctc}(\mathcal{X}', \tilde{\mathcal{T}})$$

M-Step



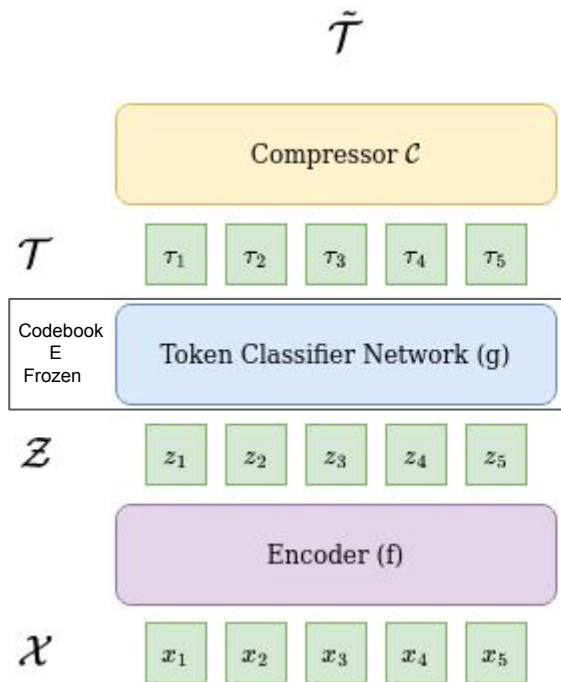
$$\mathcal{X} \mapsto \tilde{\mathcal{T}}, \mathcal{X}' \mapsto \tilde{\mathcal{T}}'$$

$$\mathcal{L} = \mathcal{L}_m(\mathcal{X}, \mathcal{X}') + \alpha \mathcal{L}_{ctc}(\mathcal{X}, \tilde{\mathcal{T}}') + \beta \mathcal{L}_{ctc}(\mathcal{X}', \tilde{\mathcal{T}})$$

Likelihood Loss



M-Step



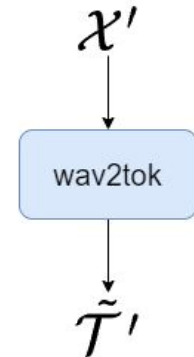
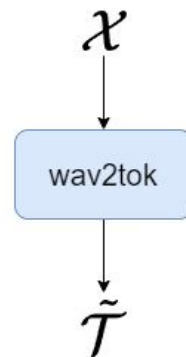
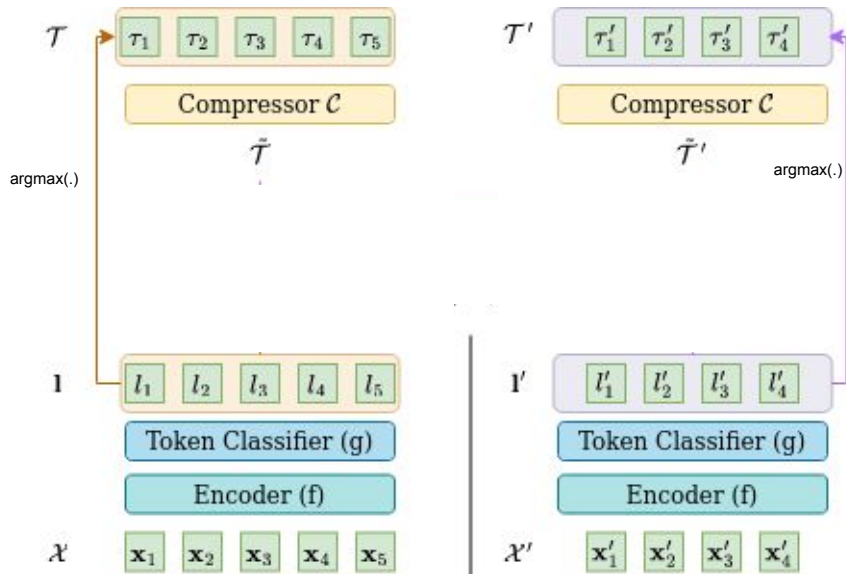
$$\mathcal{X} \mapsto \tilde{\mathcal{T}}, \mathcal{X}' \mapsto \tilde{\mathcal{T}'}$$

$$\mathcal{L} = \mathcal{L}_m(\mathcal{X}, \mathcal{X}') + \alpha \mathcal{L}_{ctc}(\mathcal{X}, \tilde{\mathcal{T}'}) + \beta \mathcal{L}_{ctc}(\mathcal{X}', \tilde{\mathcal{T}})$$

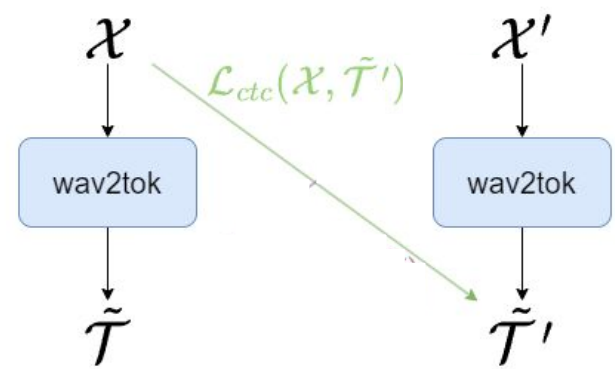
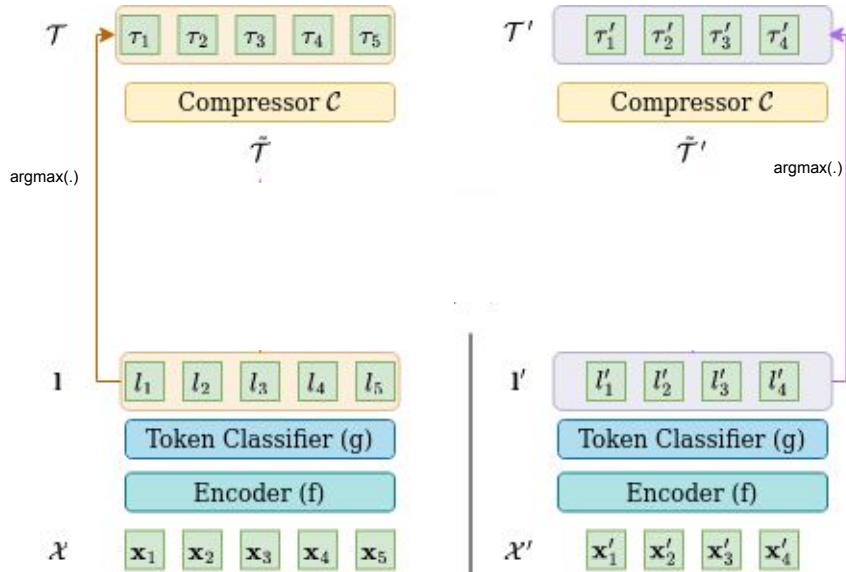
Contrastive Loss

Likelihood Loss

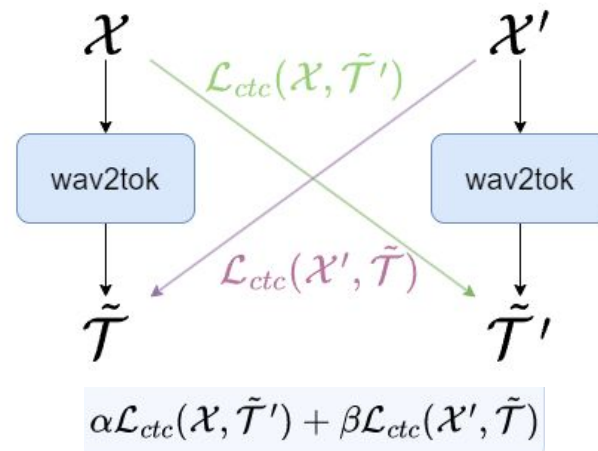
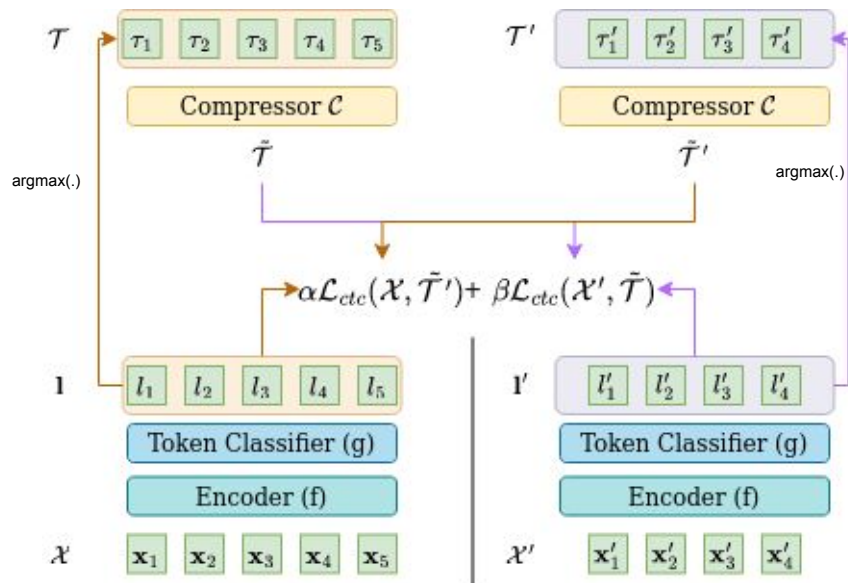
Likelihood Loss



Likelihood Loss

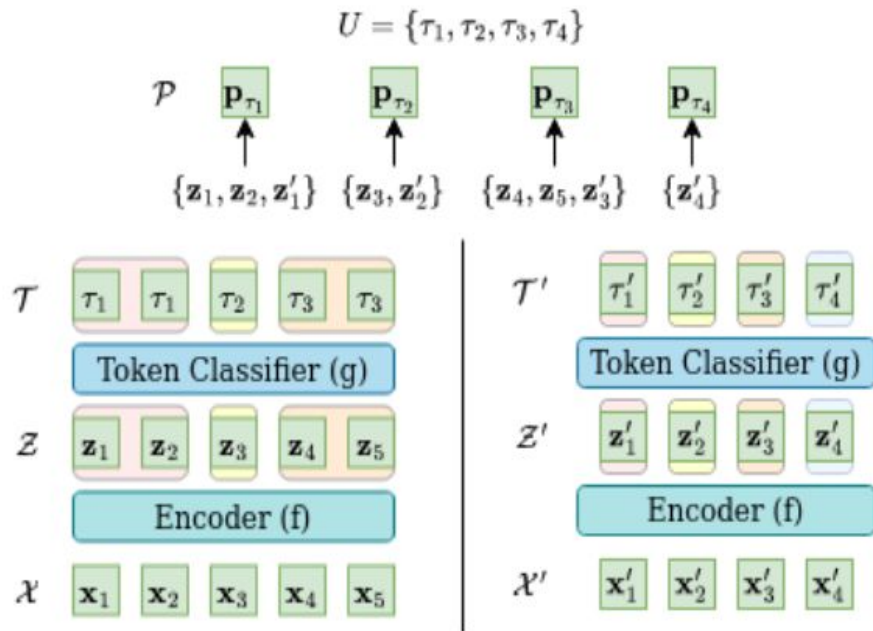


Likelihood Loss

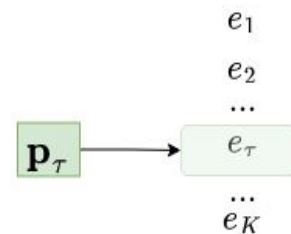


- We use the CTC forward backward algorithm [3]

Contrastive Loss



Contrastive Task



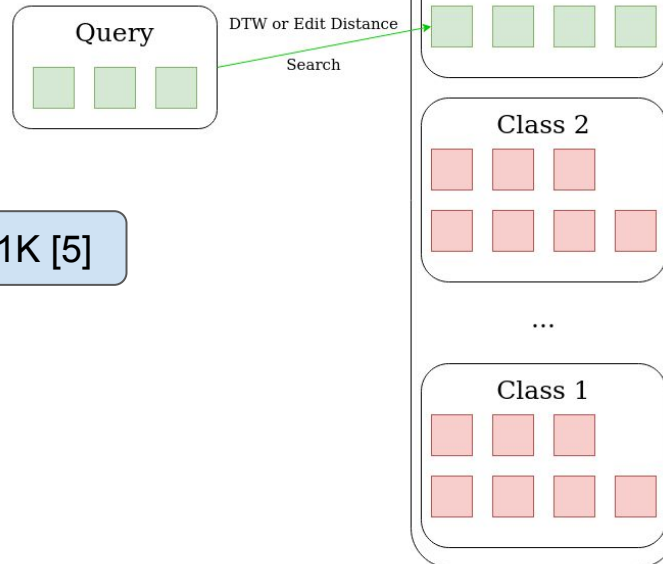
Experiments



Experiments

Melody Retrieval

Evaluation of tokens learnt
from Music audio



Training:

MIR-QbSH [4]

MIR-1K [5]

Testing:

MIR-QbSH [4]

[4] https://www.music-ir.org/mirex/wiki/2020:Query_by_Singing/Humming

[5] <http://mirlab.org/dataset/public/>



Results (Music tokens)

Table 1: Quality of Tokenization

Model	Vanilla query ✓ (MRR)	Time Stretched query TS (MRR)	Pitch Shifted query PS (MRR)	Search Time (s)	Infer (s)
MIDI ED	0.75	0.64	0.72	3.84	0.62
Relative Note DTW	0.84	0.74	0.8	0.02	0.62
wav2vec2-O ED	0.72	0.72	0.71	0.01	0.43
wav2vec2-Multi ED	0.82	0.82	0.82		1.2
wav2tok ED	0.84	0.84	0.84	0.04	0.14



Results (Music tokens)

Audio to MIDI with SOTA melody extraction algorithm [6]

Table 1: Quality of Tokenization

Model	Vanilla query	Time Stretched query	Pitch Shifted query	Search Time (s)	Infer (s)
	V	TS	PS		
MIDI ED	0.75	0.64	0.72	3.84	0.62
Relative Note DTW	0.84	0.74	0.8	0.02	0.62
wav2vec2-O ED	0.72	0.72	0.71	0.01	0.43
wav2vec2-Multi ED	0.82	0.82	0.82		1.2
wav2tok ED	0.84	0.84	0.84	0.04	0.14

[6] Justin Salamon and Emilia Gomez. Melody extraction from polyphonic music signals using pitch contour characteristics. IEEE transactions on audio, speech, and language processing, 20(6): 1759–1770, 2012.



Results (Music tokens)

Audio to MIDI with SOTA melody extraction algorithm [6]

MIDI [6] to Relative Note sequence [7]

MIDI : 55,...,55,56,...,78
to
Notes: (55,0.1s), (56,0.3s), ...
to
Relative notes: (0,0.1s),(+3, 0.3s),...

Table 1: Quality of Tokenization

	Vanilla query	Time Stretched query	Pitch Shifted query		
Model	V	TS	PS	Search Time	Infer
	(MRR)	(MRR)	(MRR)	(s)	(s)
MIDI ED	0.75	0.64	0.72	3.84	0.62
Relative Note DTW	0.84	0.74	0.8	0.02	0.62
wav2vec2-O ED	0.72	0.72	0.71	0.01	0.43
wav2vec2-Multi ED	0.82	0.82	0.82		1.2
wav2tok ED	0.84	0.84	0.84	0.04	0.14

[6] Justin Salamon and Emilia Gomez. Melody extraction from polyphonic music signals using pitch contour characteristics. IEEE transactions on audio, speech, and language processing, 20(6): 1759–1770, 2012.

[7] Naziba Mostafa and Pascale Fung. A note based query by humming system using convolutional neural network. In INTERSPEECH, pp. 3102–3106, 2017.



Results (Music tokens)

Audio to MIDI with SOTA melody extraction algorithm [6]

MIDI [6] to Relative Note sequence [7]

MIDI : 55,...,55,56,...,78
to
Notes: (55,0.1s), (56,0.3s), ...
to
Relative notes: (0,0.1s),(+3, 0.3s),...

wav2vec 2.0 base model pretrained on 960h LibriSpeech and Finetuned on TIMIT [8]

Table 1: Quality of Tokenization

	Vanilla query	Time Stretched query	Pitch Shifted query		
Model	V	TS	PS	Search Time	Infer
	(MRR)	(MRR)	(MRR)	(s)	(s)
MIDI ED	0.75	0.64	0.72	3.84	0.62
Relative Note DTW	0.84	0.74	0.8	0.02	0.62
wav2vec2-O ED	0.72	0.72	0.71	0.01	0.43
wav2vec2-Multi ED	0.82	0.82	0.82		1.2
wav2tok ED	0.84	0.84	0.84	0.04	0.14

[6] Justin Salamon and Emilia Gomez. Melody extraction from polyphonic music signals using pitch contour characteristics. IEEE transactions on audio, speech, and language processing, 20(6): 1759–1770, 2012.

[7] Naziba Mostafa and Pascale Fung. A note based query by humming system using convolutional neural network. In INTERSPEECH, pp. 3102–3106, 2017.

[8] [wav2vec 2.0: A Framework for Self-Supervised ... - NIPS papers](#)



Results (Music tokens)

Audio to MIDI with SOTA melody extraction algorithm [6]

MIDI [6] to Relative Note sequence [7]

MIDI : 55,...,55,56,...,78
to
Notes: (55,0.1s), (56,0.3s), ...
to
Relative notes: (0,0.1s),(+3, 0.3s),...

wav2vec 2.0 base model pretrained on 960h LibriSpeech and Finetuned on TIMIT [8]

Multi-Lingually (53 languages) pretrained wav2vec 2.0 large model and Finetuned on Common Voice [9]

Table 1: Quality of Tokenization

	Vanilla query	Time Stretched query	Pitch Shifted query		
Model	V	TS	PS	Search Time	Infer
	(MRR)	(MRR)	(MRR)	(s)	(s)
MIDI ED	0.75	0.64	0.72	3.84	0.62
Relative Note DTW	0.84	0.74	0.8	0.02	0.62
wav2vec2-O ED	0.72	0.72	0.71	0.01	0.43
wav2vec2-Multi ED	0.82	0.82	0.82		1.2
wav2tok ED	0.84	0.84	0.84	0.04	0.14

[6] Justin Salamon and Emilia Gomez. Melody extraction from polyphonic music signals using pitch contour characteristics. IEEE transactions on audio, speech, and language processing, 20(6): 1759–1770, 2012.

[7] Naziba Mostafa and Pascale Fung. A note based query by humming system using convolutional neural network. In INTERSPEECH, pp. 3102–3106, 2017.

[8] [wav2vec 2.0: A Framework for Self-Supervised ... - NIPS papers](#)

[9] Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli. Unsupervised cross-lingual representation learning for speech recognition. CoRR, abs/2006.13979, 2020. URL <https://arxiv.org/abs/2006.13979>



Results (Music tokens)

Table 2: Some Variations

	Vanilla query	Time Stretched query	Pitch Shifted query
Model	V (MRR)	TS (MRR)	PS (MRR)
log-mel DTW	0.72	0.7	0.67
vq-log-mel ED	0.71	0.6	0.62
wav2tok+MIR1K ED	0.72	0.64	0.67
wav2tok ED	0.84	0.84	0.84



Results (Music tokens)

Table 2: Some Variations

Model	Vanilla query	Time Stretched query	Pitch Shifted query
	V (MRR)	TS (MRR)	PS (MRR)
log-mel DTW	0.72	0.7	0.67
vq-log-mel ED	0.71	0.6	0.62
wav2tok+MIR1K ED	0.72	0.64	0.67
wav2tok ED	0.84	0.84	0.84

K-Means over log-mel features



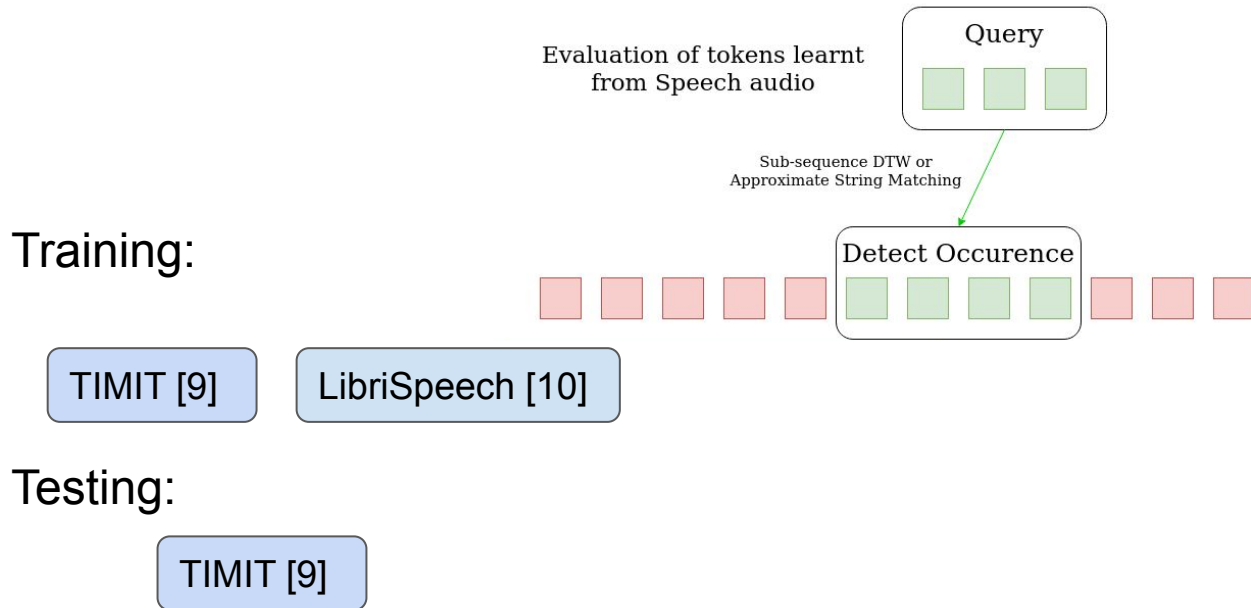
Results (Music tokens)

Table 2: Some Variations

	Vanilla query	Time Stretched query	Pitch Shifted query
Model	V (MRR)	TS (MRR)	PS (MRR)
K-Means over log-mel features	0.72	0.7	0.67
vq-log-mel ED	0.71	0.6	0.62
wav2tok trained on MIR-1K polyphonic music audio	0.72	0.64	0.67
wav2tok ED	0.84	0.84	0.84

Experiments

Spoken Term Detection



[9] Garofolo, John S., et al. TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download. Philadelphia: Linguistic Data Consortium, 1993.

[10] V. Panayotov, G. Chen, D. Povey and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), South Brisbane, QLD, Australia, 2015, pp. 5206-5210, doi: 10.1109/ICASSP.2015.7178964.



Results (Speech tokens)

Table 4: Spoken Term Detection

Model	ED (F1)	Search Time (s)	DTW (F1)	Search Time (s)
log-mel DTW	-	-	0.41	0.003
wav2vec2-O	0.61	0.29	0.43	0.23
wav2vec2-Multi	0.63	0.72	0.48	0.66
wav2tok	0.65	0.064	0.52	0.09
wav2tok+Libri	0.63		0.44	0.1



Results (Speech tokens)

Table 4: Spoken Term Detection

Spoken Term Detection system based on [11]

S-DTW over Posteriorograms generated by SOTA ASRs

Model	ED (F1)	Search Time (s)	DTW (F1)	Search Time (s)
log-mel DTW	-	-	0.41	0.003
wav2vec2-O	0.61	0.29	0.43	0.23
wav2vec2-Multi	0.63	0.72	0.48	0.66
wav2tok	0.65	0.064	0.52	0.09
wav2tok+Libri	0.63		0.44	0.1



Thank You!