

# Discrete Audio Tokens for Multimodal LLMs

Mirco Ravanelli



# “Traditional” LLMs

- Language Models have a **long history**.
- **Goal:** Assign a probability to every sequence of words  $\mathbf{w}$

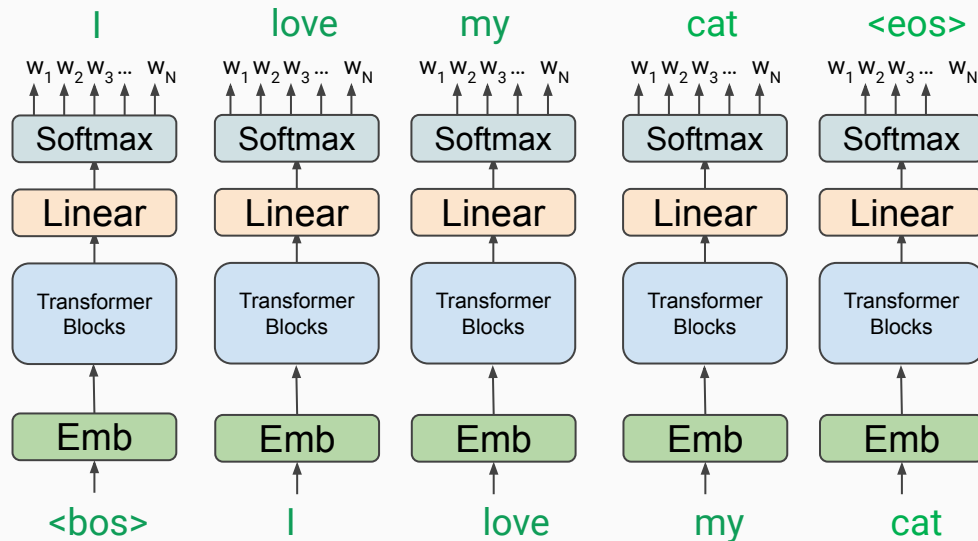
$$\begin{aligned} P(w_1, w_2, \dots, w_N) &= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)\dots P(w_N|w_1, \dots, w_{N-1}) \\ &= \prod_{i=1}^N P(w_i|w_1, \dots, w_{i-1}) \end{aligned}$$

- A language model can be trained to predict the **next word** given **all the previous ones**.



# Large Language Models

- **Dominant Approach:** autoregressive transformer trained to predict the next word.



**Brute Force:** Massively scaling up the model with more data and parameters results in impressive performance.



# Why Multimodality?

- A multimodal system can offer several advantages:

Audio



Video

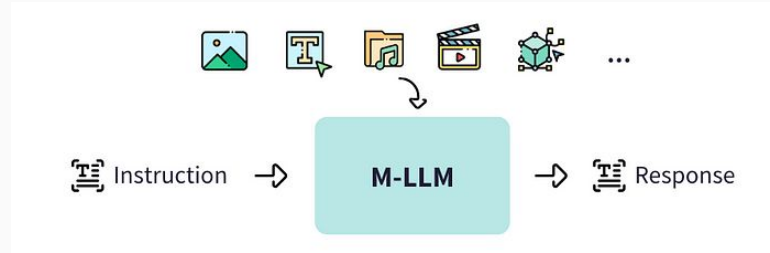


Image



Text

challenging to find hap  
vation that you don't  
'cause we're compar  
er people rather th  
it's really going on i  
ocus on the neo-ti



Deeper Understanding

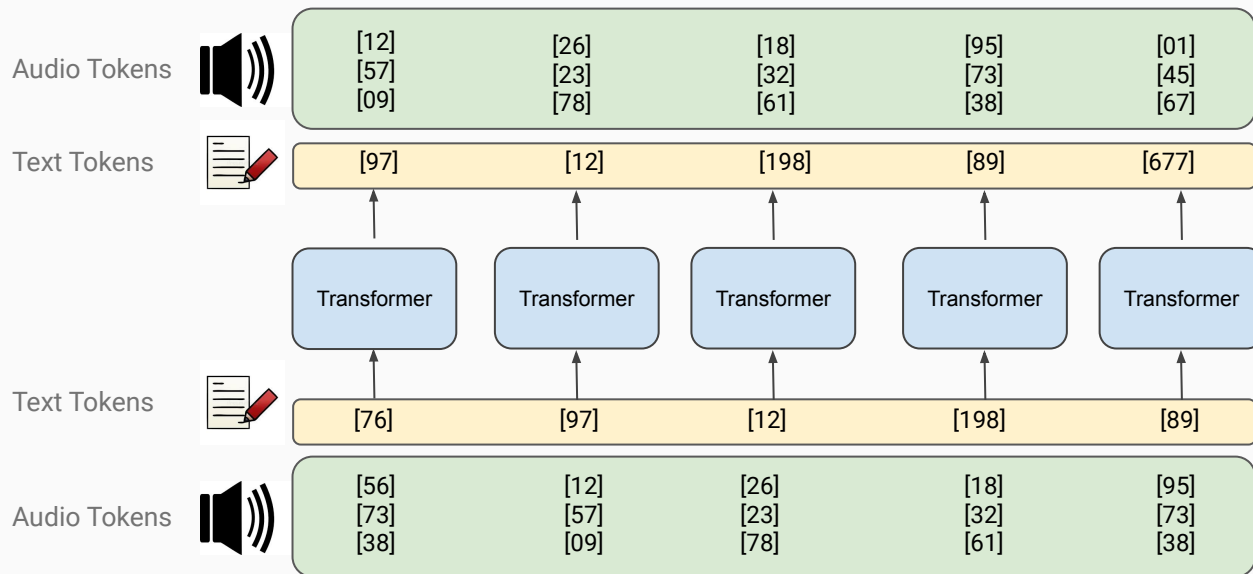
Better Generalization

Fewer Hallucinations

More Flexible  
Human-Machine Interaction

# Multimodal LLMs

- One way to train multimodal LLMs? **Tokenize all!**

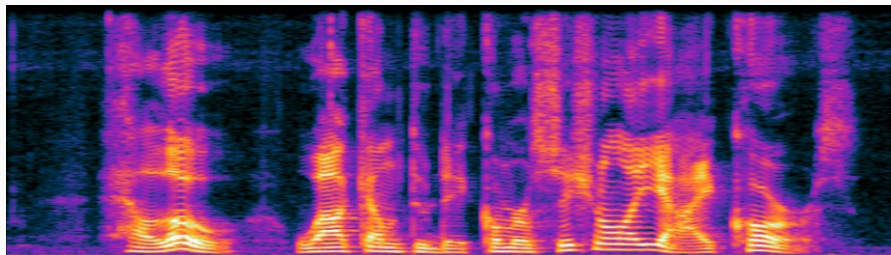


# How to Tokenize Audio?

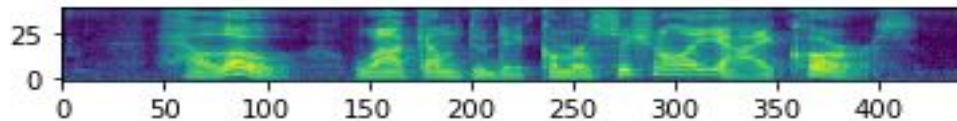
- Audio and Speech are **continuous**!



Waveform



Spectrogram



FBANKs

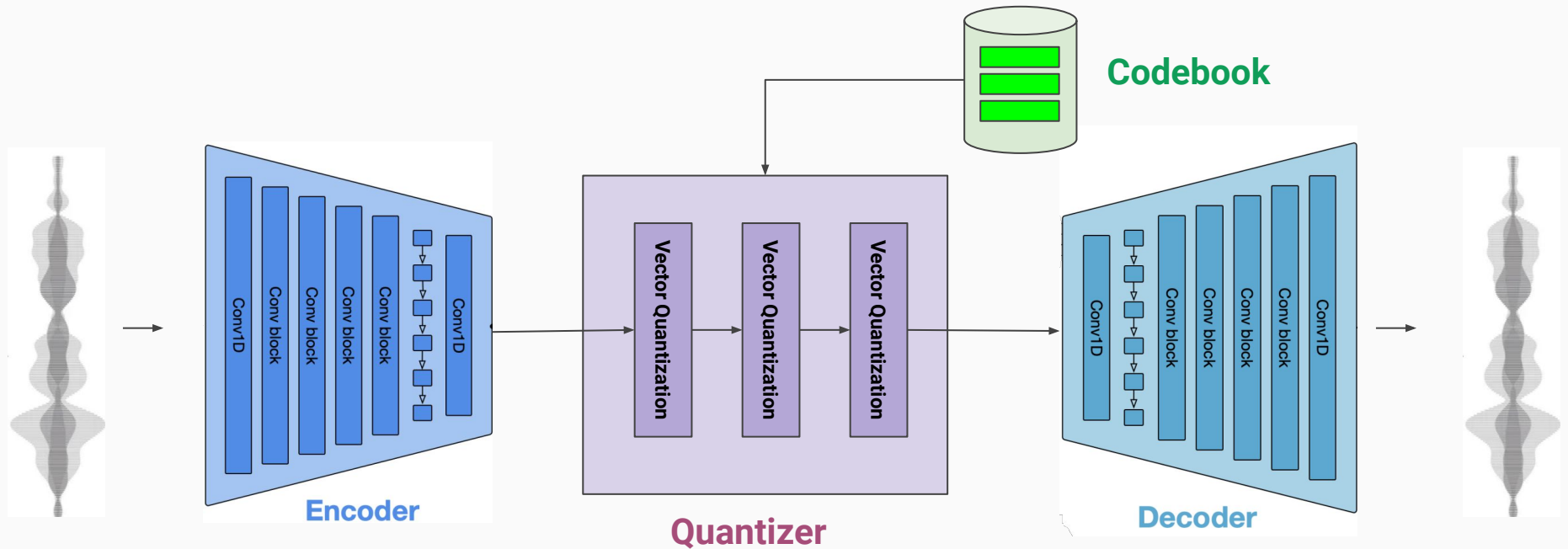


*How to Discretize Audio  
and Speech Without  
Losing Information?*

- SSL models (e.g., wav2vec, HuBERT, WavLM) also produce continuous representations!

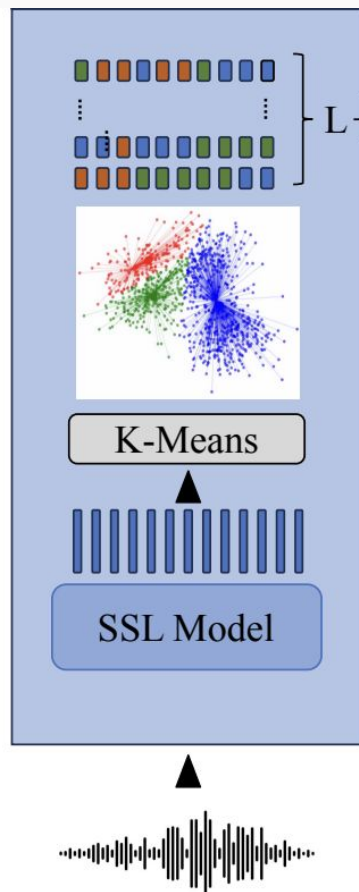
# Audio Codes

- **Compression tokens (Audio Codec)**



**Examples:** *SoundStream, Encodec*

# Semantic Tokens



## Semantic Tokens

- Large self-supervised models (e.g., *Wav2vec*, *HuBERT*, *WaLM*) achieve **state-of-the-art** performance.
- These models are based on **continuous** representations
- We convert continuous representations into discrete tokens through **clustering**.

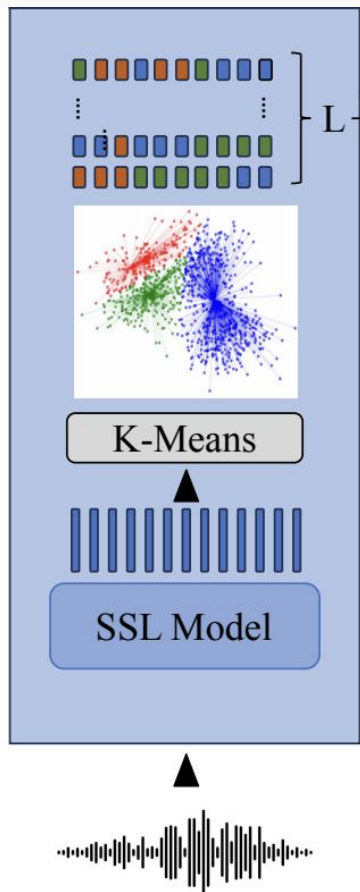


**How Should We Extract Discrete Audio Tokens from Self-Supervised Models?**

*Pooneh Mousavi<sup>1,2</sup>, Jarod Duret<sup>3</sup>, Salah Zaiem<sup>4</sup>, Luca Della Libera<sup>1,2</sup>, Artem Ploujnikov<sup>5,2</sup>, Cem Subakan<sup>6,2,1</sup>, Mirco Ravanelli<sup>1,2,5</sup>*

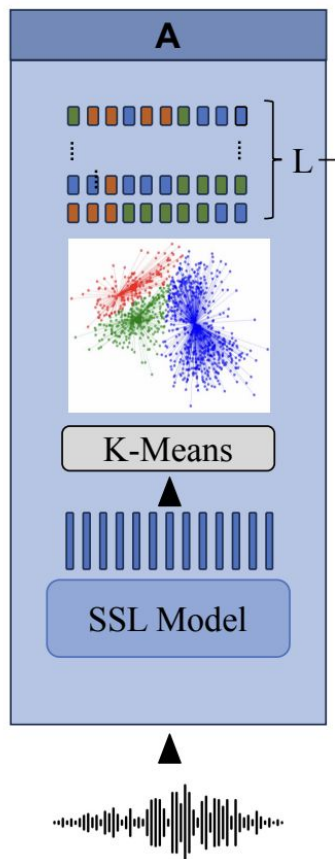


# Extracting Semantic Tokens



1. *Which layers should we cluster?*
2. *What is the optimal number of clusters?*
3. *Which datasets are we using for clustering?*
4. *What is the best approach to train the decoder (vocoder)?*
5. *How should we initialize the embeddings effectively?*
6. *Can we extract universal tokens for both discriminative and generative tasks?*

# Extracting Semantic Tokens



# Extracting Semantic Tokens

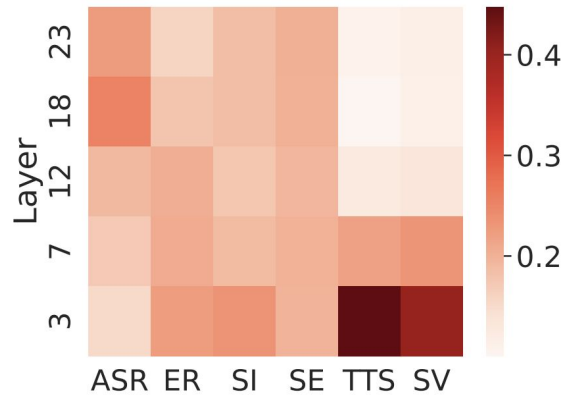


Figure 3: Attention analysis across various tasks and layers of the discrete WavLM model with in-domain tokenizers.

- Automatic Speech Recognition (ASR)
- Emotion Recognition (ER)
- Speaker Identification (SI)

- Speech Enhancement (SE)
- Text-to-Speech (TTS)
- Scalable Vocoder (SV)

# Extracting Semantic Tokens

Table 1: Assessing the impact of the number of clusters and embedding initialization on discrete WavLM-Large across different tasks.

Setting	ASR (EN)	ASR (FR)	SID	ER	SE		TTS		
	WER ↓	WER ↓	ACC ↑	ACC ↑	DNSMOS↑	dWER↓	UTMOS↑	WER ↓	
Effect of Number Of Clusters									
1000	7.15	34.61	79.0	61.8	3.93	6.75	3.65	5.76	
2000	6.96	32.94	79.5	67.2	3.93	6.58	3.55	5.62	
Effect of Embedding Initialization									
Random	6.96	32.94	81.0	67.2	3.93	6.75	3.65	5.76	
PreTrained & finetune	8.93	35.81	77.5	63.9	3.93	6.82	3.64	6.62	
PreTrained & freeze	9.26	35.12	73.1	67.0	3.93	6.98	3.66	6.42	

- The optimal number of clusters varies by task.
- Typically, using 1000 to 2000 clusters yields good performance.
- There is no advantage observed in initializing the embedding with pretrained centroid embeddings.

# Extracting Semantic Tokens

Table 2: *Out-of-domain and in-domain performance of discrete HuBERT and WavLM models across the downstream tasks.*

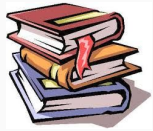
SSL Model	Tokenizer	ASR (EN)	ASR (FR)	SID	ER	SE		TTS		Vocoder	
		WER ↓	WER ↓	ACC ↑	ACC ↑	DNSMOS↑	dWER↓	UTMOS↑	WER ↓	UTMOS↑	dWER↓
HuBERT Large [4]	In-Domain	7.89	38.29	67.2	64.5	3.98	17.64	3.61	6.46	3.50	4.49
	Out-Of-Domain	N/A	39.50	67.8	61.7	3.95	15.92	3.54	5.45	3.48	2.92
WavLM Large [3]	In-Domain	6.96	32.94	81.0	67.2	3.93	6.75	3.65	5.76	3.49	2.98
	Out-Of-Domain	N/A	36.25	79.0	61.9	3.96	6.49	3.61	5.73	3.68	2.95

- As expected, the in-domain tokenizer outperforms its OOD counterpart.
- However, the performance drop is not always huge.

# Discrete Audio and Speech Benchmark (DASB)



*Are compression tokens better than semantic tokens?*



**Literature offers no clear answer**

---

**DASB - Discrete Audio and Speech Benchmark**

---

[DASB Code Repo](#)

Pooneh Mousavi<sup>1,2</sup>, Luca Della Libera<sup>1,2</sup>, Jarod Duret<sup>3</sup>, Artem Ploujnikov<sup>4,2</sup>,  
Cem Subakan<sup>5,2,1</sup>, Mirco Ravanelli<sup>1,2,4</sup>

<sup>1</sup>Concordia University <sup>2</sup>Mila - Quebec AI Institute <sup>3</sup>Avignon Université

<sup>4</sup>Université de Montréal <sup>5</sup>Université Laval

# Discrete Audio and Speech Benchmark (DASB)

Tokenizer	Type
Discrete HuBERT	Semantic
Discrete WavLM	Semantic
Discrete Wav2Vec2	Semantic
EnCodec	Compression
DAC	Compression
SpeechTokenizer	Hybrid

Task	Type
Automatic Speech Recognition (ASR)	Discriminative
Speaker Identification/Verification (SID, SV)	Discriminative
Emotion Recognition (ER)	Discriminative
Intent Classification (IC)	Discriminative
Keyword Spotting (KS)	Discriminative
Speech Enhancement (SE)	Generative
Speech Separation (SS)	Generative
Text-to-Speech (TTS)	Generative

# Discrete Audio and Speech Benchmark (DASB)

Table 2: Benchmarking results for discriminative tasks.

Models/Tasks	ASR-En		ASR-multiling		ER	IC	KS	SI	SV
	WER ↓		WER ↓		ACC ↑	ACC ↑	ACC ↑	ACC ↑	EER ↓
	Clean	Other	Welsh	Basque					
<i>Low Bitrate</i>									
Discrete Hubert	<b>8.99</b>	<b>21.14</b>	<b>58.50</b>	<b>26.83</b>	57.20	68.70	90.54	0.90	24.99
Discrete WavLM	11.72	27.56	60.37	28.63	<b>59.80</b>	73.40	<b>97.94</b>	0.70	26.02
Discrete Wav2Vec2	12.14	28.65	66.30	32.25	57.80	<b>74.10</b>	96.16	0.40	33.53
EnCodec	52.37	77.04	92.01	58.20	44.70	31.50	86.00	<b>58.30</b>	<b>17.40</b>
DAC	63.96	83.61	94.86	66.29	49.20	22.10	81.00	45.10	20.62
SpeechTokenizer	19.77	43.12	76.67	47.92	49.10	57.90	95.09	47.40	20.41
<i>Medium Bitrate</i>									
Discrete Hubert	<b>7.91</b>	<b>18.95</b>	54.77	23.63	<b>62.10</b>	70.50	94.69	67.40	15.71
Discrete WavLM	8.52	20.35	<b>54.22</b>	<b>22.06</b>	57.60	<b>78.00</b>	<b>98.09</b>	80.80	8.00
Discrete Wav2Vec2	8.76	21.32	60.39	26.64	59.10	75.10	96.64	65.47	17.64
EnCodec	46.80	74.24	91.23	47.95	51.30	31.40	88.70	<b>91.90</b>	<b>7.81</b>
DAC	59.54	81.48	97.43	56.16	45.80	18.90	76.60	83.80	11.78
SpeechTokenizer	18.32	41.21	75.17	38.94	52.10	57.80	94.86	91.40	7.88
<i>High Bitrate</i>									
EnCodec	<b>45.18</b>	<b>72.56</b>	<b>93.40</b>	<b>87.65</b>	46.40	<b>19.60</b>	<b>83.60</b>	<b>92.81</b>	<b>7.18</b>
DAC	99.53	99.38	99.40	99.68	<b>46.00</b>	15.70	75.20	85.61	10.89
<i>Continuous Baseline</i>									
SSL	3.370	7.04	41.77	14.32	63.10	86.10	99.00	99.70	2.10



- **Semantic tokens** outperform compression tokens in most discriminative tasks.
- The exception is **speaker recognition**, where EnCodec excels.
- **Big gap** compared to **continuous** baselines!



# Discrete Audio and Speech Benchmark (DASB)



Table 3: Benchmarking results for generative tasks. N.C. indicates “Not Converged”.

Models/Tasks	SE			SS			TTS	
	DNSMOS ↑	dWER ↓	SpkSim ↑	DNSMOS ↑	dWER ↓	SpkSim ↑	UTMOS ↑	dWER ↓
<i>Low Bitrate</i>								
Discrete HuBERT	3.33	<b>15.47</b>	0.824	3.52	80.86	0.840	3.24	<b>2.55</b>
Discrete WavLM	3.26	16.52	0.830	3.43	<b>62.34</b>	0.847	<b>3.84</b>	3.01
Discrete Wav2Vec2	<b>3.55</b>	18.86	0.779	<b>3.75</b>	96.70	0.787	3.32	3.45
EnCodec	3.15	34.35	0.852	3.11	83.55	<b>0.877</b>	1.46	8.85
DAC	3.30	57.41	0.853	3.01	102.00	0.854	1.97	10.68
SpeechTokenizer	3.18	30.13	<b>0.858</b>	3.13	85.25	0.874	2.51	3.69
<i>Medium Bitrate</i>								
Discrete HuBERT	3.48	12.62	0.875	3.70	66.29	0.891	3.80	3.40
Discrete WavLM	3.48	<b>10.18</b>	0.889	3.68	<b>34.03</b>	0.912	<b>3.82</b>	<b>2.45</b>
Discrete Wav2Vec2	<b>3.54</b>	17.60	0.858	<b>3.75</b>	78.42	0.866	3.68	2.89
EnCodec	3.10	19.07	0.885	3.09	48.57	0.906	1.50	94.6
DAC	3.49	31.14	<b>0.906</b>	3.26	55.43	<b>0.924</b>	1.71	71.26
SpeechTokenizer	3.49	23.44	0.876	3.42	60.75	0.906	1.96	53.26
<i>High Bitrate</i>								
EnCodec	2.87	68.22	0.814	<b>2.95</b>	<b>97.73</b>	<b>0.839</b>	N.C	N.C
DAC	<b>2.95</b>	<b>46.07</b>	<b>0.860</b>	2.53	208	0.784	N.C	N.C
<i>Continuous Baseline</i>								
SSL	3.49	4.92	0.928	3.68	9.97	0.939	3.71	2.94

- Semantic tokens show the best performance for generative tasks as well
- **Big gap** compared to **continuous** baselines

# Discrete Audio and Speech Benchmark (DASB)



TTS Samples



Ranking aggregation for models (medium bitrate)

Model	Disc.	Gen.	Comb.
Discrete HuBERT	2.66	3.62	3.11
Discrete WavLM	<b>2.00</b>	2.75	<b>1.94</b>
Discrete Wav2Vec2	3.33	<b>2.68</b>	3.41
EnCodec	4.11	3.93	4.23
DAC	5.55	4.06	4.64
SpeechTokenizer	3.44	3.81	3.64



# Discrete Audio and Speech Benchmark (DASB)

- We are still **far** from an ideal solution!



- Semantic tokens are computationally expensive
- They don't preserve speaker identities well
- Performance drops significantly compared to continuous representations

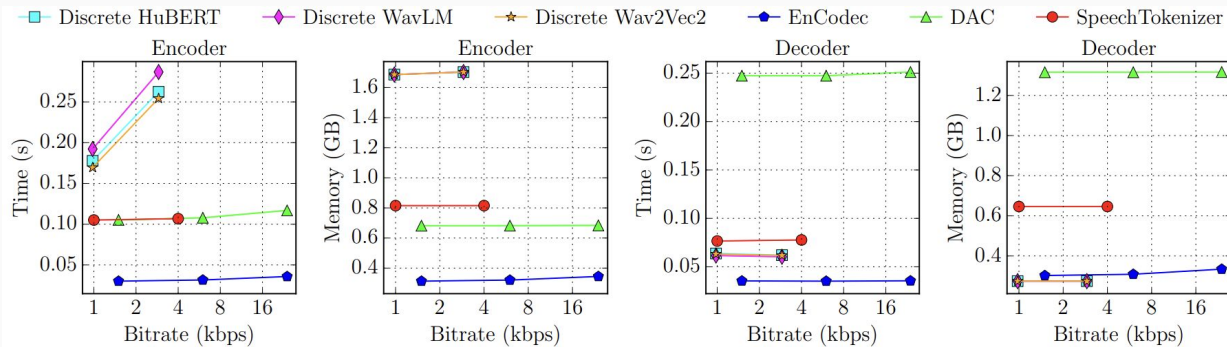


Figure 2: Time and memory required to process an utterance of 16 seconds for encoders and decoders of the considered audio tokenizers on an NVIDIA GeForce RTX 3070 GPU @ 8 GB.



# Future Directions

- Universal Audio Tokens

Emotion  
Recognition



Speaker  
Recognition

Speech  
Recognition

Audio Editing

Text-to-Speech

Enhancement

.....



## Some Ideas

- **Massive Multitask** learning (similar to [PASE](#))
- **Hierarchical** codebooks with **Dynamic Allocation** (more details for Music)
- **Perceptual Loss Optimization** (Similar to [MetricGAN](#))
- Better **Multi-scale** processing.

# Future Directions

- Can we learn “**interpretable**” audio tokens?

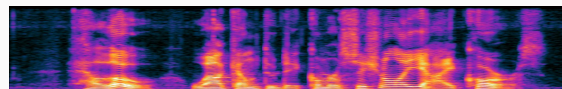
**Text:** “The City of Montréal”

**Tokenized:** “[The] [City] [of] [Mont] [ré] [a]”



Each token is **easily interpretable**, as it maps to a specific part of the text.

**Audio**



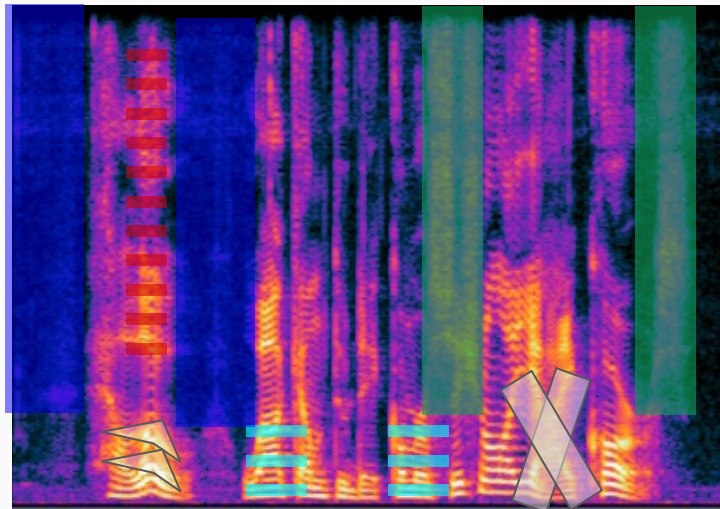
**Tokenized Audio**

[56]	[12]	[26]	[18]
[73]	[57]	[23]	[32]
[38]	[09]	[78]	[61]

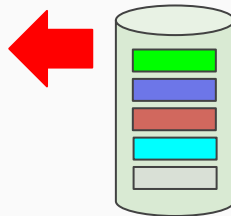


No clear mapping to the original signal

# Future Directions



Interpretable  
Codebook



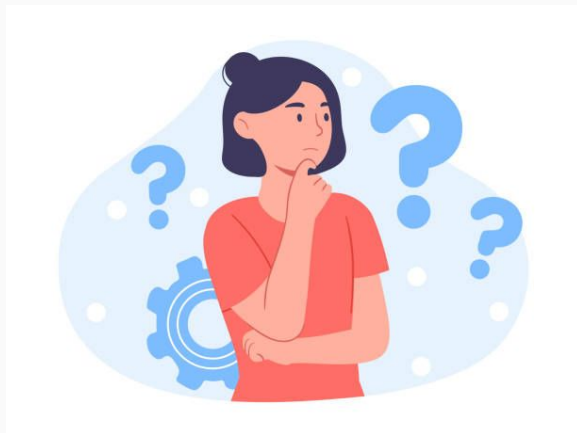
## Possible Advantages

- Increased **Transparency**
- Easier **Error Analyses**
- Interpretability might act as a **powerful regularizer**, enhancing both generalization and performance.

Our Related Paper: [NeurIPS 2024](#), [ICML 2024](#)

# Future Directions

- Is **audio tokenization** the best way to go?



- Survey Paper and Extended Benchmark (in Progress)

# Thank you!



Pooneh Mousavi



Artem Ploujnikov



Luca Della Libera



Jarod Duret



Salah Zaiem



Cem Subakan