

Unsupervised Symbol Discovery

by applying language constraints on sequential data

Charly Lamothe¹ Thierry Artières² Ricard Marxer³

¹Aix-Marseille Université
QARMA - équipe d'Apprentissage de Marseille

²ECM - École Centrale de Marseille
QARMA - équipe d'Apprentissage de Marseille
LIS - Laboratoire d'Informatique et Systèmes

³Université de Toulon
Aix-Marseille Université
LIS - Laboratoire d'Informatique et Systèmes

September 06, 2019

Summary

- 1** Introduction
- 2** Preliminary notions
- 3** Work achieved
- 4** Evaluation and results
- 5** Conclusion and perspectives

Introduction

Sequential data

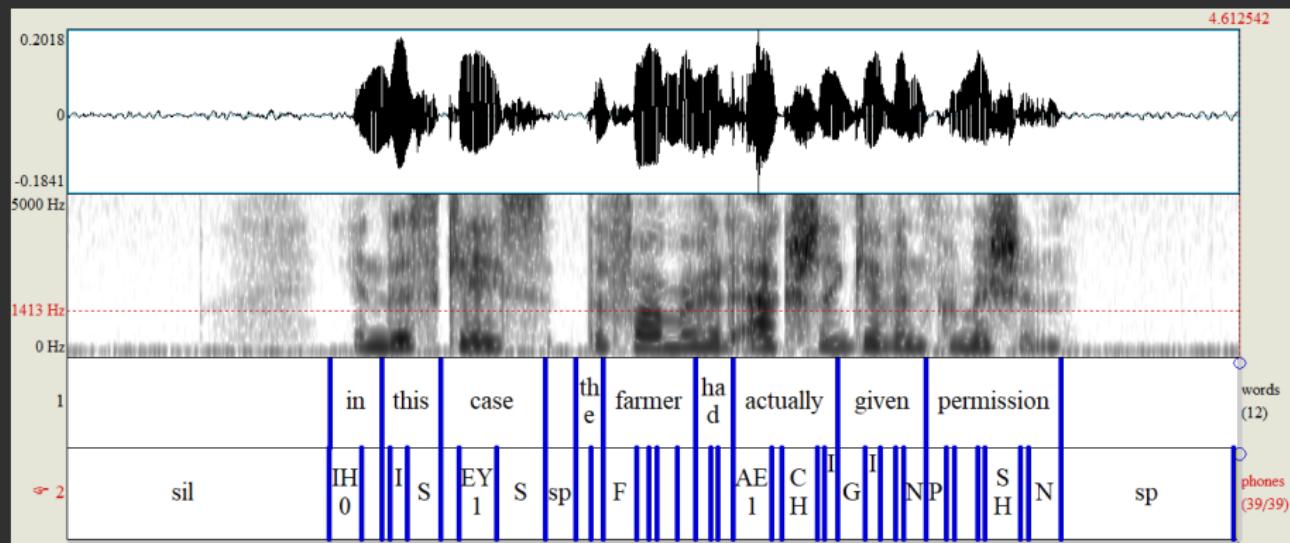


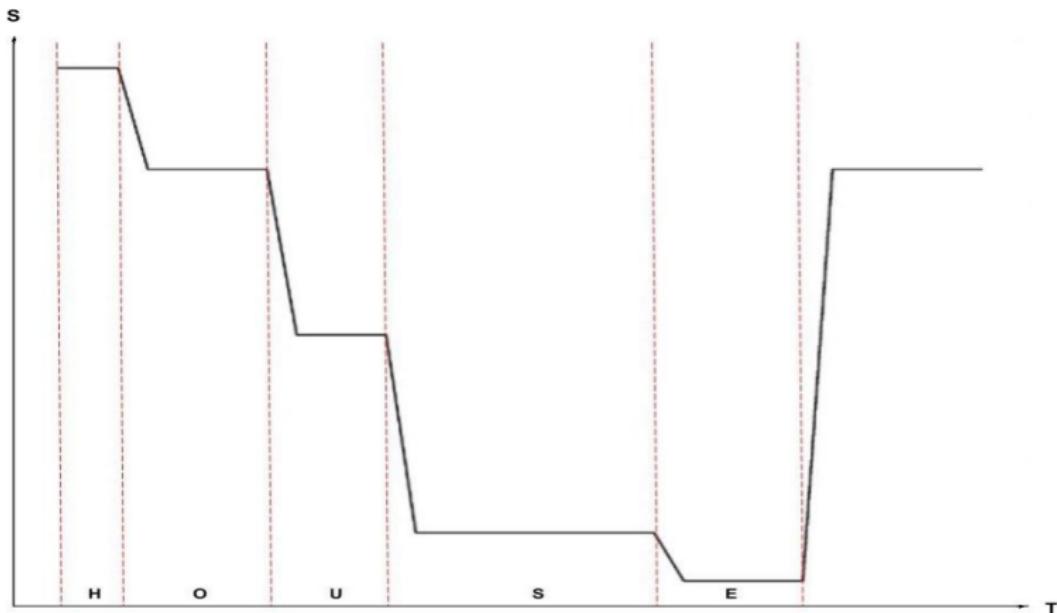
Figure: Example of speech signal with its aligned transcription.

Sequential data

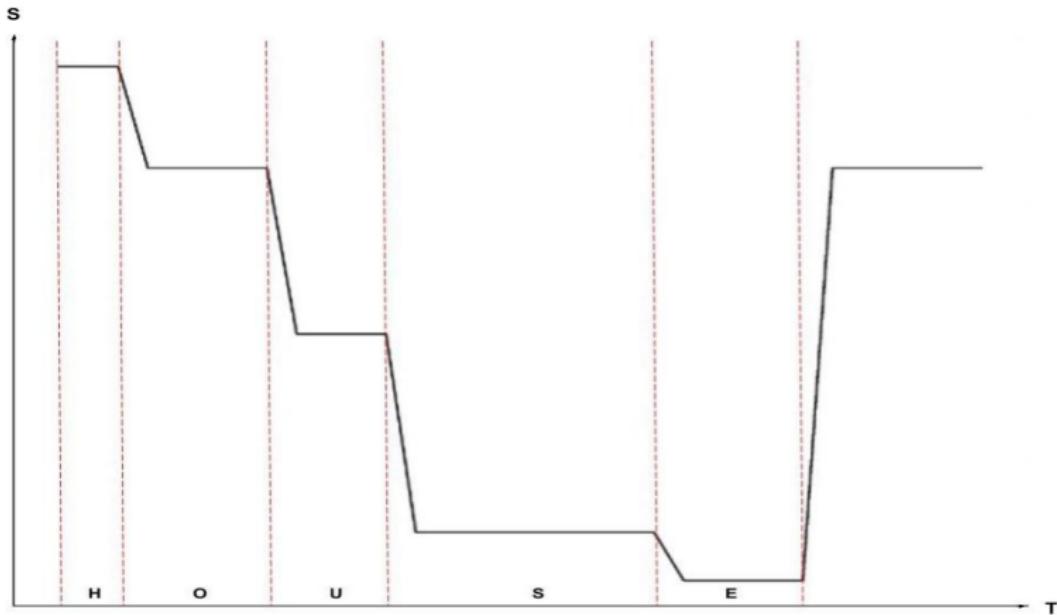


Figure: Example of handwritten text image with its aligned transcription.

Information in sequential structure



Information in sequential structure



$$S_{\text{house}} = - \sum_i P_i \log P_i \quad \forall_i, P := \{P(H), P(O), P(U), P(S), P(E)\}$$

Hypotheses

- (1) The existing statistical regularities of all languages L could be exploited to learn better representation of high level semantic content within sequential data;
- (2) Different languages L share some equivalent statistical regularities;
- (3) These regularities are shared across datasets S of different modalities M .

Hypotheses

- (1) The existing statistical regularities of all languages L could be exploited to learn better representation of high level semantic content within sequential data;
- (2) Different languages L share some equivalent statistical regularities;
- (3) These regularities are shared across datasets S of different modalities M .

Hypotheses

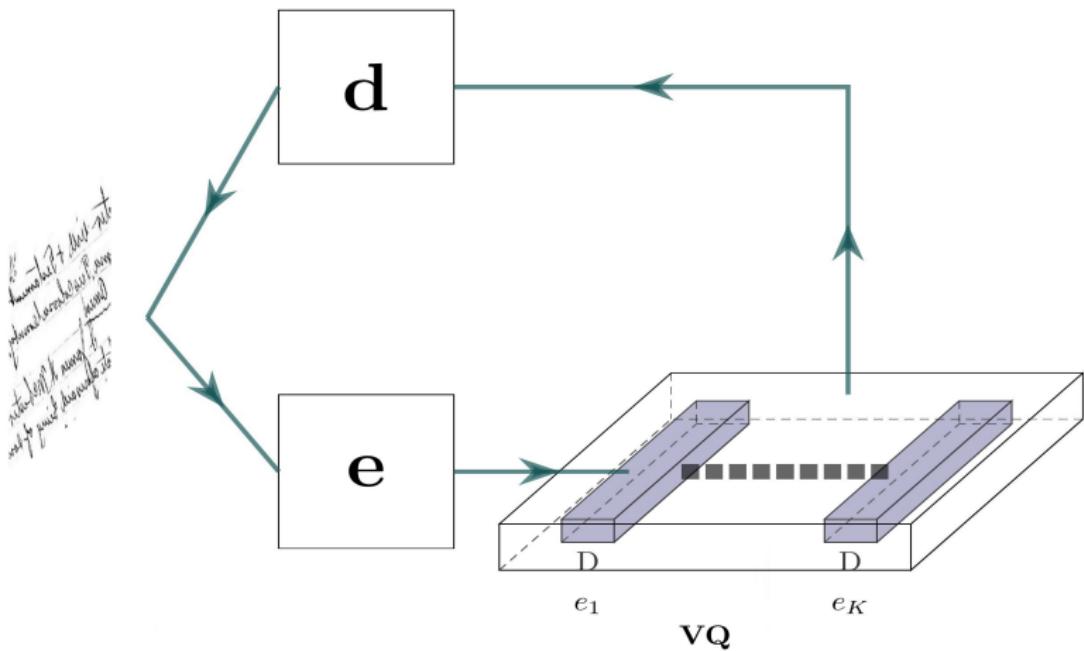
- (1) The existing statistical regularities of all languages L could be exploited to learn better representation of high level semantic content within sequential data;
- (2) Different languages L share some equivalent statistical regularities;
- (3) These regularities are shared across datasets S of different modalities M .

Preliminary notions

$\mathcal{S} := \{(x_i^{(t)}, y_i^{(t)})\}_{i=1}^m$ a dataset with $(x_i^{(t)}, y_i^{(t)}) \in \mathcal{X} \times \mathcal{Y}$;

$\forall y_i^{(t)} \in \mathcal{Y}$ an aligned transcription of $x_i^{(t)} \in \mathcal{X}$;

$\mathcal{S}_{\text{train}} \in \mathcal{S} \setminus \{y^{(t)}\}_{i=1}^m$, $\mathcal{S}_{\text{eval}} \in \mathcal{S}$.



\mathbf{e} encodes an input sequential data $x^{(t)}$ at timestep t such that $\hat{x}^{(t)} = \mathbf{d}(\mathbf{e}(x^{(t)}))$, where $z_{\mathbf{e}(x^{(t)})} = \mathbf{e}(x^{(t)})$ is a lower-dimensional space that keeps the most important features of $x^{(t)}$.

Vector Quantized-Variational Autoencoder (VQ-VAE)

Let $\{e_i \in \mathbb{R}^D, i = 1, \dots, K\}$ be the K prototype vectors of dimension D ;

$$q(x^{(t)}) = \operatorname{argmin}_i \|z_{\mathbf{e}(x^{(t)})} - e_i\|_2^2 \quad \forall i \in \{1, \dots, K\};$$

$$z_{q(x^{(t)})} = e_{q(x^{(t)})};$$

$$\hat{x}^{(t)} = \mathbf{d}(z_{q(x^{(t)})});$$

$$\mathcal{L}_{\text{VQ-VAE}} = \mathcal{L}_{\text{err}} + \| \operatorname{sg}(z_{\mathbf{e}(x^{(t)})}) - e_{q(x^{(t)})} \|_2^2 + \gamma \| z_{\mathbf{e}(x^{(t)})} - \operatorname{sg}(e_{q(x^{(t)})}) \|_2^2$$

Vector Quantized-Variational Autoencoder (VQ-VAE)

Let $\{e_i \in \mathbb{R}^D, i = 1, \dots, K\}$ be the K prototype vectors of dimension D ;

$$q(x^{(t)}) = \operatorname{argmin}_i \|z_{\mathbf{e}(x^{(t)})} - e_i\|_2^2 \quad \forall i \in \{1, \dots, K\};$$

$$\hat{z}_{q(x^{(t)})} = e_{q(x^{(t)})};$$

$$\hat{x}^{(t)} = \mathbf{d}(z_{q(x^{(t)})});$$

$$\mathcal{L}_{\text{VQ-VAE}} = \mathcal{L}_{\text{err}} + \| \operatorname{sg}(z_{\mathbf{e}(x^{(t)})}) - e_{q(x^{(t)})} \|_2^2 + \gamma \| z_{\mathbf{e}(x^{(t)})} - \operatorname{sg}(e_{q(x^{(t)})}) \|_2^2$$

Vector Quantized-Variational Autoencoder (VQ-VAE)

Let $\{e_i \in \mathbb{R}^D, i = 1, \dots, K\}$ be the K prototype vectors of dimension D ;

$$q(x^{(t)}) = \operatorname{argmin}_i \|z_{\mathbf{e}(x^{(t)})} - e_i\|_2^2 \quad \forall i \in \{1, \dots, K\};$$

$$z_{q(x^{(t)})} = e_{q(x^{(t)})};$$

$$\hat{x}^{(t)} = \mathbf{d}(z_{q(x^{(t)})});$$

$$\mathcal{L}_{\text{VQ-VAE}} = \mathcal{L}_{\text{err}} + \| \text{sg}(z_{\mathbf{e}(x^{(t)})}) - e_{q(x^{(t)})} \|_2^2 + \gamma \| z_{\mathbf{e}(x^{(t)})} - \text{sg}(e_{q(x^{(t)})}) \|_2^2$$

Vector Quantized-Variational Autoencoder (VQ-VAE)

Let $\{e_i \in \mathbb{R}^D, i = 1, \dots, K\}$ be the K prototype vectors of dimension D ;

$$q(x^{(t)}) = \operatorname{argmin}_i \|z_{\mathbf{e}(x^{(t)})} - e_i\|_2^2 \quad \forall i \in \{1, \dots, K\};$$

$$z_{q(x^{(t)})} = e_{q(x^{(t)})};$$

$$\hat{x}^{(t)} = \mathbf{d}(z_{q(x^{(t)})});$$

$$\mathcal{L}_{\text{VQ-VAE}} = \mathcal{L}_{\text{err}} + \| \text{sg}(z_{\mathbf{e}(x^{(t)})}) - e_{q(x^{(t)})} \|_2^2 + \gamma \| z_{\mathbf{e}(x^{(t)})} - \text{sg}(e_{q(x^{(t)})}) \|_2^2$$

Vector Quantized-Variational Autoencoder (VQ-VAE)

Let $\{e_i \in \mathbb{R}^D, i = 1, \dots, K\}$ be the K prototype vectors of dimension D ;

$$q(x^{(t)}) = \operatorname{argmin}_i \|z_{\mathbf{e}(x^{(t)})} - e_i\|_2^2 \quad \forall i \in \{1, \dots, K\};$$

$$z_{q(x^{(t)})} = e_{q(x^{(t)})};$$

$$\hat{x}^{(t)} = \mathbf{d}(z_{q(x^{(t)})});$$

$$\mathcal{L}_{\text{VQ-VAE}} = \mathcal{L}_{\text{err}} + \| \mathsf{sg}(z_{\mathbf{e}(x^{(t)})}) - e_{q(x^{(t)})} \|_2^2 + \gamma \| z_{\mathbf{e}(x^{(t)})} - \mathsf{sg}(e_{q(x^{(t)})}) \|_2^2$$

Example of quantized embedding space

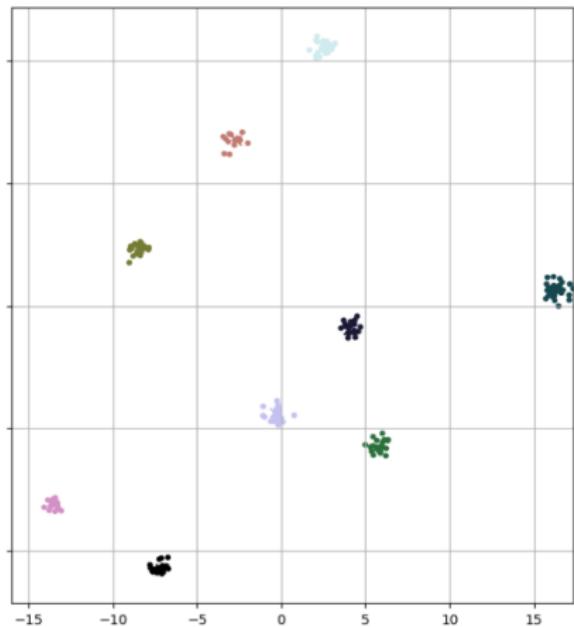
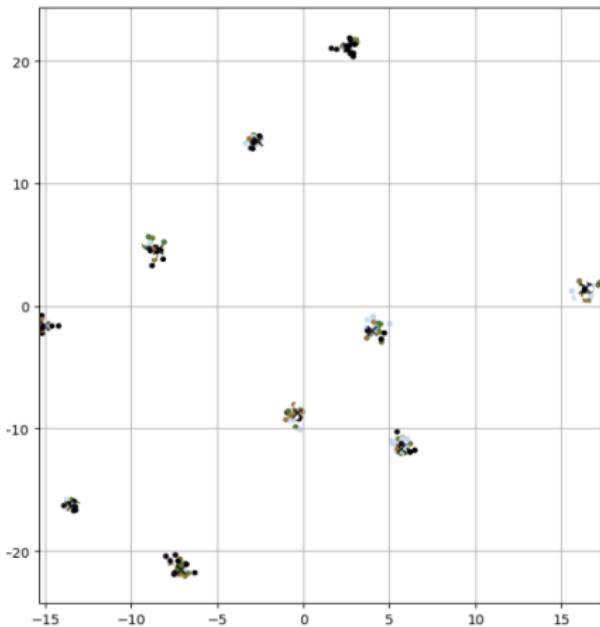
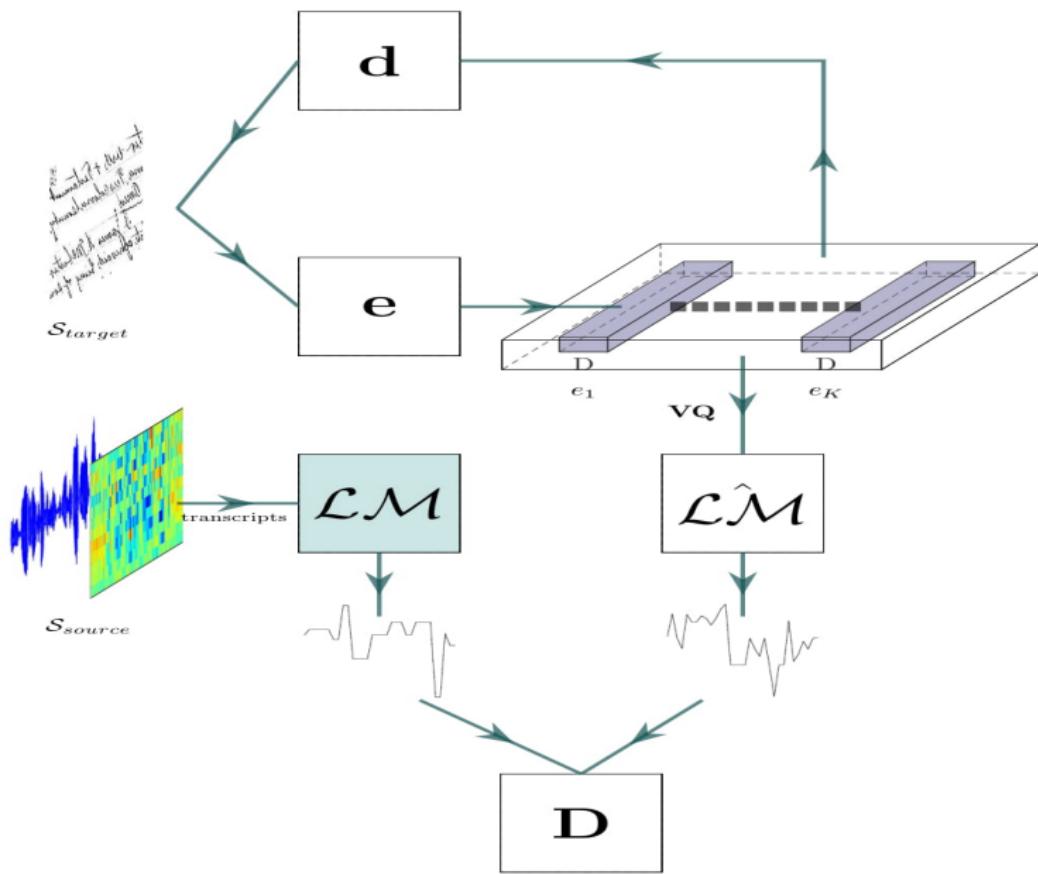


Figure: *Left* Audio frame points colored by speaker id and the embedding marks colored in black. *Right* Aaudio frame points colored by encoding indices and the embedding marks colored by number of embedding vectors.

Work achieved



Experiments

$$\mathbf{L}_a, \mathbf{M}_a \rightarrow \mathbf{L}_a, \mathbf{M}_a$$

$$\mathbf{L}_a, \mathbf{M}_a \rightarrow \mathbf{L}_a, \mathbf{M}_b$$

Both train a Language Model \mathcal{LM} computed from a dataset $\mathcal{S}_{\text{source}} \sim \mathbf{L}, \mathbf{M}$ and a Language Model $\hat{\mathcal{LM}}$ computed from a dataset $\mathcal{S}_{\text{target}} \sim \mathbf{L}, \mathbf{M}$ and use them on the main model training with a discriminator \mathbf{D} .

n-gram Language Model

Compute the n-gram matrix \mathcal{G} from the source dataset S_{source} ;

Compute the entropy vector \mathcal{G}_H using \mathcal{G} ;

Compute the entropy dataset $S_{\mathcal{G}_H}$ using the entropy vector \mathcal{G}_H using \mathcal{G} and the source dataset S_{source} ;

Do the same steps for $\hat{\mathcal{G}}$ and $S_{\hat{\mathcal{G}}_H}$ on S_{target} .

n-gram Language Model

Compute the n-gram matrix \mathcal{G} from the source dataset $\mathcal{S}_{\text{source}}$;

Compute the entropy vector \mathcal{G}_H using \mathcal{G} ;

Compute the entropy dataset $\mathcal{S}_{\mathcal{G}_H}$ using the entropy vector \mathcal{G}_H using \mathcal{G} and the source dataset $\mathcal{S}_{\text{source}}$;

Do the same steps for $\hat{\mathcal{G}}$ and $\mathcal{S}_{\hat{\mathcal{G}}_H}$ on $\mathcal{S}_{\text{target}}$.

n-gram Language Model

Compute the n-gram matrix \mathcal{G} from the source dataset $\mathcal{S}_{\text{source}}$;

Compute the entropy vector \mathcal{G}_H using \mathcal{G} ;

Compute the entropy dataset $\mathcal{S}_{\mathcal{G}_H}$ using the entropy vector \mathcal{G}_H using \mathcal{G} and the source dataset $\mathcal{S}_{\text{source}}$;

Do the same steps for $\hat{\mathcal{G}}$ and $\mathcal{S}_{\hat{\mathcal{G}}_H}$ on $\mathcal{S}_{\text{target}}$.

n-gram Language Model

Compute the n-gram matrix \mathcal{G} from the source dataset $\mathcal{S}_{\text{source}}$;

Compute the entropy vector \mathcal{G}_H using \mathcal{G} ;

Compute the entropy dataset $\mathcal{S}_{\mathcal{G}_H}$ using the entropy vector \mathcal{G}_H using \mathcal{G} and the source dataset $\mathcal{S}_{\text{source}}$;

Do the same steps for $\hat{\mathcal{G}}$ and $\mathcal{S}_{\hat{\mathcal{G}}_H}$ on $\mathcal{S}_{\text{target}}$.

Example of bi-gram Language Model

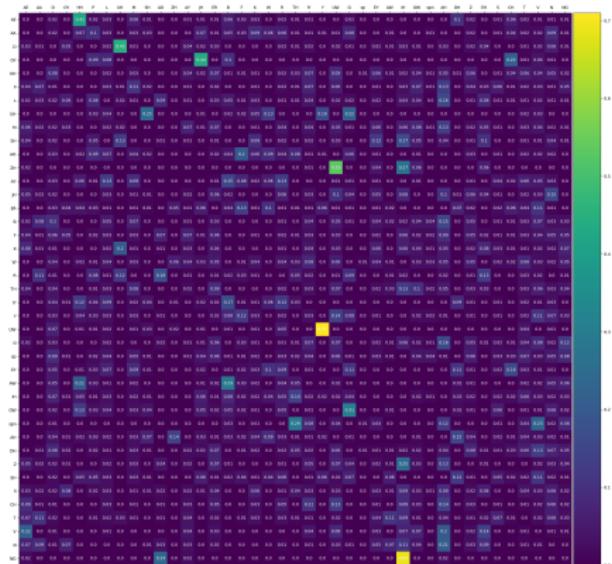
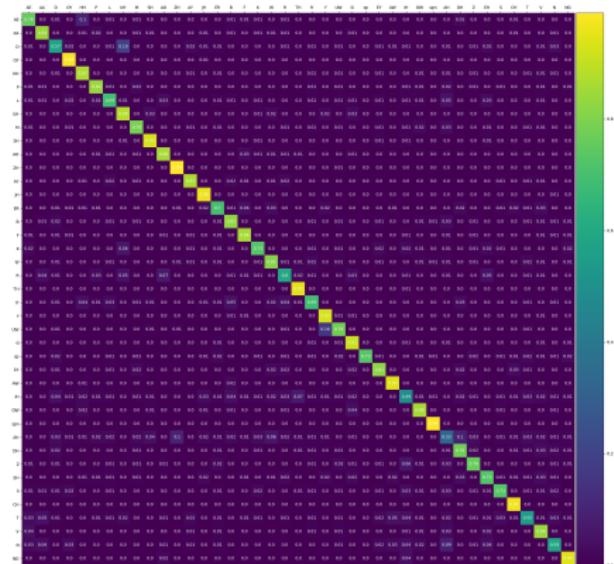


Figure: The bi-gram matrix of the groundtruth alignments of the validation subset of VCTK dataset, with a step of 20ms.

RNN Language Model

Train iteratively the model \mathcal{R} on $\mathcal{S}_{\text{source}}$;

Compute the source entropy dataset $\mathcal{S}_{\hat{\mathcal{R}}_H}$ from \mathcal{R} ;

Do the same steps for $\hat{\mathcal{R}}$ and $\mathcal{S}_{\hat{\mathcal{R}}_H}$ on $\mathcal{S}_{\text{target}}$.

RNN Language Model

Train iteratively the model \mathcal{R} on $\mathcal{S}_{\text{source}}$;

Compute the source entropy dataset $\mathcal{S}_{\mathcal{R}_H}$ from \mathcal{R} ;

Do the same steps for $\hat{\mathcal{R}}$ and $\mathcal{S}_{\hat{\mathcal{R}}_H}$ on $\mathcal{S}_{\text{target}}$.

RNN Language Model

Train iteratively the model \mathcal{R} on $\mathcal{S}_{\text{source}}$;

Compute the source entropy dataset $\mathcal{S}_{\mathcal{R}_H}$ from \mathcal{R} ;

Do the same steps for $\hat{\mathcal{R}}$ and $\mathcal{S}_{\hat{\mathcal{R}}_H}$ on $\mathcal{S}_{\text{target}}$.

Adversarial training

Fooling loss: $\text{BCE}(\mathbf{D}(\hat{y}_{\mathcal{G}_H}^{(t)}), 1);$

Autoencoder training objective using ngram LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{ngram}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}}$$

Autoencoder training objective using RNN LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{rnn}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}} + \mathcal{L}_{\hat{\mathcal{R}}}$$

Discriminator training objective for both case:

$$\max_{\mathbf{D}} \mathcal{L}_{\mathbf{D}} = 0.5 \times (\text{BCE}(\mathbf{D}(y_{\mathcal{G}_H}^{(t)} \in \mathcal{S}_{\mathcal{G}_H}), 1) + \text{BCE}(\mathbf{D}(sg(\hat{y}_{\mathcal{G}_H}^{(t)})), 0))$$

Adversarial training

Fooling loss: $\text{BCE}(\mathbf{D}(\hat{y}_{\mathcal{G}_H}^{(t)}), 1);$

Autoencoder training objective using ngram LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{ngram}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}}$$

Autoencoder training objective using RNN LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{rnn}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}} + \mathcal{L}_{\hat{\mathcal{R}}}$$

Discriminator training objective for both case:

$$\max_{\mathbf{D}} \mathcal{L}_{\mathbf{D}} = 0.5 \times (\text{BCE}(\mathbf{D}(y_{\mathcal{G}_H}^{(t)}) \in \mathcal{S}_{\mathcal{G}_H}), 1) + \text{BCE}(\mathbf{D}(sg(\hat{y}_{\mathcal{G}_H}^{(t)})), 0))$$

Adversarial training

Fooling loss: $\text{BCE}(\mathbf{D}(\hat{y}_{\mathcal{G}_H}^{(t)}), 1);$

Autoencoder training objective using ngram LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{ngram}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}}$$

Autoencoder training objective using RNN LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{rnn}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}} + \mathcal{L}_{\hat{\mathcal{R}}}$$

Discriminator training objective for both case:

$$\max_{\mathbf{D}} \mathcal{L}_{\mathbf{D}} = 0.5 \times (\text{BCE}(\mathbf{D}(y_{\mathcal{G}_H}^{(t)}) \in \mathcal{S}_{\mathcal{G}_H}), 1) + \text{BCE}(\mathbf{D}(sg(\hat{y}_{\mathcal{G}_H}^{(t)})), 0))$$

Adversarial training

Fooling loss: $\text{BCE}(\mathbf{D}(\hat{y}_{\mathcal{G}_H}^{(t)}), 1);$

Autoencoder training objective using ngram LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{ngram}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}}$$

Autoencoder training objective using RNN LM:

$$\min_{\mathbf{e}, \mathbf{d}} \mathcal{L}_{\text{rnn}} = \mathcal{L}_{\text{VQ-VAE}} + \mathcal{L}_{\text{fool}} + \mathcal{L}_{\hat{\mathcal{R}}}$$

Discriminator training objective for both case:

$$\max_{\mathbf{D}} \mathcal{L}_{\mathbf{D}} = 0.5 \times (\text{BCE}(\mathbf{D}(y_{\mathcal{G}_H}^{(t)} \in \mathcal{S}_{\mathcal{G}_H}), 1) + \text{BCE}(\mathbf{D}(sg(\hat{y}_{\mathcal{G}_H}^{(t)})), 0))$$

Evaluation and results

Evaluation

Let $u_1 = \{1, 1, 1, 2, 2, 3, 3\}$ and $v_1 = \{4, 4, 4, 1, 1, 3, 3\}$ be two samples of respectively subsets U and V ;

$s = ARI(., .)$ a clustering metric;

$$ARI(u_1, v_1) = ARI(v_1, u_1) = 1.0.$$

Evaluation

Let $u_1 = \{1, 1, 1, 2, 2, 3, 3\}$ and $v_1 = \{4, 4, 4, 1, 1, 3, 3\}$ be two samples of respectively subsets U and V ;

$s = ARI(., .)$ a clustering metric;

$$ARI(u_1, v_1) = ARI(v_1, u_1) = 1.0.$$

Evaluation

Let $u_1 = \{1, 1, 1, 2, 2, 3, 3\}$ and $v_1 = \{4, 4, 4, 1, 1, 3, 3\}$ be two samples of respectively subsets U and V ;

$s = ARI(., .)$ a clustering metric;

$$ARI(u_1, v_1) = ARI(v_1, u_1) = 1.0.$$

Results table

| id | Architecture | ARI | AMI | Bottleneck accuracy | Perplexity |
|-----------|---|---------------|---------------|----------------------------|-------------------|
| (i) | <i>basic</i> , random | 0.0062 | 0.0028 | n/a | 12.8577 |
| (ii) | <i>basic</i> + $J_p = 0.12$ | 0.0433 | 0.0203 | n/a | 37.9220 |
| (iii) | <i>basic</i> + $\mathcal{S}_{\mathcal{R}_H^2}$ + $J_p = 0.12$, aligned, M_a | 0.0288 | 0.0244 | n/a | n/a |
| (iv) | <i>pixcnn</i> + $J_p = 0.12$ | 0.1588 | 0.2785 | 0.6066 | 185.9 |
| (v) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^2}$, $J_p = 0.12$, aligned, M_a | 0.1825 | 0.2742 | 0.5767 | 63.7217 |
| (vi) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^2}$, $J_p = 0.12$, unaligned, M_a | 0.1784 | 0.2643 | 0.5670 | 59.53 |
| (vii) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^2}$, $J_p = 0.12$, unaligned, M_b | 0.1706 | 0.2704 | 0.5709 | 65.74 |
| (viii) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^3}$, $J_p = 0.12$, aligned, M_a | 0.1785 | 0.2805 | 0.5897 | 71.7067 |
| (ix) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^3}$, $J_p = 0.12$, unaligned, M_a | 0.1585 | 0.2645 | 0.5704 | 64.9483 |
| (x) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{G}_H^3}$, $J_p = 0.12$, unaligned, M_b | 0.1836 | 0.2682 | 0.5777 | 67.6567 |
| (xi) | <i>pixcnn</i> + $\mathcal{S}_{\mathcal{R}_H^4}$, $J_p = 0.25$, unaligned, M_b | 0.1893 | 0.2583 | 0.4940 | 51.5633 |

TABLE 1 – This figure shows the results obtained by the experiments described in the present section. J_p stands for the probability of jittering, if a jitter layer is used (cf. Section 3.3). $\mathcal{S}_{\mathcal{R}_H^4}$ stands for an entropy dataset computed from a RNN (cf. Section 3.6). M_a and M_b means the source and the target datasets have two distinct modalities. ARI, AMI, are respectively the Adjusted Rand Index and the Adjusted Mutual Information scores (cf. Section 3.7 and Appendix F). The language L of the datasets are always the same. Here *basic* uses 44 tokens and *pixcnn* 128 tokens in the bottleneck.

Results summary

Using VCTK (speech) as the source dataset and Scribblelens (handwritten) as the target dataset:

Best ARI score using $\mathcal{S}_{\mathcal{R}_H^4}$;

Best AMI score using $\mathcal{S}_{\mathcal{G}_H^3}$;

These scores are slightly better than the state-of-the-art baseline.

Conclusion and perspectives

The results found on the ARI metric are consistent with our hypotheses (1) and (3);

The RNN has not reached its full capacity yet.

The results found on the ARI metric are consistent with our hypotheses (1) and (3);

The RNN has not reached its full capacity yet.

Hyperparameters tuning;

Repeat the previous experiments on multiple datasets;

Increase the size of k-seed;

Hyperparameters tuning;

Repeat the previous experiments on multiple datasets;

Increase the size of k-seed;

Hyperparameters tuning;

Repeat the previous experiments on multiple datasets;

Increase the size of k-seed;

Thank you for listening!

Example of obtained entropy curves

