

Empirical Evaluation of Sequence-to-Sequence Models for Word Discovery in Low-resource Settings

Marcely Zanon Boito¹, Aline Villavicencio² and Laurent Besacier¹

¹LIG, Univ. Grenoble Alpes, **France** ²DCS, Univ. of Sheffield, **UK** ³INF, UFRGS, **Brazil**

Motivation



Computational Language Documentation (CLD)

- → 50 to 90% of the currently spoken languages will go extinct before 2100¹
- Manually documenting all these languages is infeasible



Computational Language Documentation (CLD)

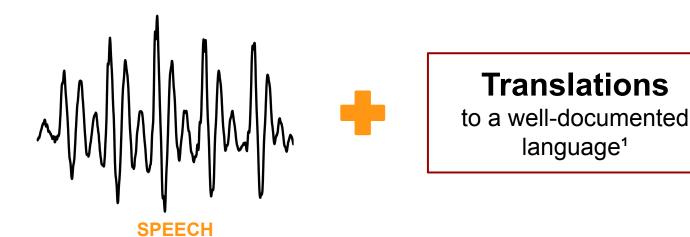
- → 50 to 90% of the currently spoken languages will go extinct before 2100¹
- Manually documenting all these languages is infeasible



GOAL: to automatically retrieve information about language structures to speed up language documentation

Documentation Corpora

- Small size (difficult to collect)
- Often lack written form (oral-tradition languages)
- → Parallel information (translations instead of transcriptions)



Documentation Corpora

- Small size (difficult to collect)
- → Often lack written form (oral-tradition languages)
- Parallel information (translations instead of transcriptions)

CLD approaches

- Deal with speech
- 2. Incorporate bilingual (or multilingual) annotations
- 3. Robust to low-resource

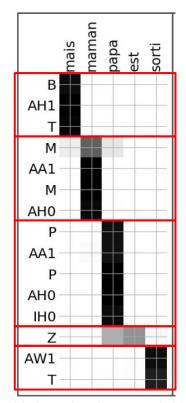


Unsupervised Word Segmentation from Speech with Attention¹

Low-resource speech word segmentation

We use **Neural Machine Translation** (NMT) model for generating **soft-alignment probability** matrices

Use of the RNN model from Bahdanau et. 2015²



But since then...

SOTA performance with different neural architectures (Transformer¹, convS2S², pervasive attention³, etc) using different attention mechanisms

But since then...

SOTA performance with different neural architectures (Transformer¹, convS2S², pervasive attention³, etc) using different attention mechanisms

 New discussion about the interpretability of the attention mechanism^{4,5,6,7}

Contribution

Empirical evaluation of 3 attention-based seq2seq models

Using different architectures, we want to investigate:

- Their capacity of generating exploitable alignment
- Their robustness to data scarcity



Experimental Settings



Task Pipeline

Encoder:

word1,word2,word3,word4...

Decoder:

phn1,phn2,phn3,phn4,phn10,phn1...

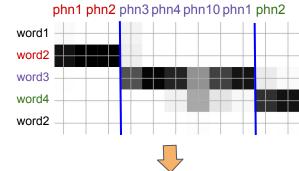


Unit Discovery System



seq2seq system





Segmentation:

phn1phn2, phn3phn4phn10phn1, phn2

Alignment:

(phn1phn2, word2); (phn3phn4phn10phn1, word3); (phn2, word4)

Task Pipeline: Data

Encoder:

word1,word2,word3,word4...

Decoder:

phn1,phn2,phn3,phn4,phn10,phn1...

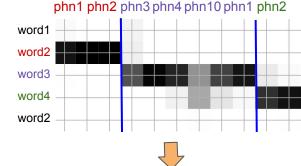


Unit Discovery System



seq2seq system





Segmentation:

phn1phn2, phn3phn4phn10phn1, phn2

Alignment:

(phn1phn2, word2); (phn3phn4phn10phn1, word3); (phn2, word4)

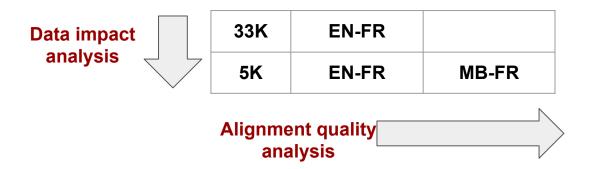
Corpora (3)

(MB-FR) Mboshi-French parallel corpus.

documentation dataset; tailored sentences

(EN-FR) English-French parallel corpus.

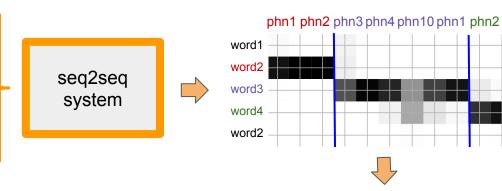
librispeech augmentation in French; noisy aligned information (filtered)





Task Pipeline: Models

Encoder: word1,word2,word3,word4... **Decoder:** phn1,phn2,phn3,phn4,phn10,phn1... **Unit Discovery** System



Segmentation:

phn1phn2, phn3phn4phn10phn1, phn2

Alignment:

(phn1phn2, word2); (phn3phn4phn10phn1, word3); (phn2, word4)

Sequence-to-Sequence (seq2seq) Models

RNN-based

Global attention.

(Bahdanau et al. 2015)

Attention mechanism creates context vectors.

$$c_t = Att(H, s_{t-1}) = \sum_{i=1}^{A} \alpha_i^t h_i$$

$$\alpha_i^t = \operatorname{softmax}(align(h_i, s_{t-1}))$$

Sequence-to-Sequence (seq2seq) Models

RNN-based

Global attention.

(Bahdanau et al. 2015)

Attention mechanism creates context vectors.

$$c_t = Att(H, s_{t-1}) = \sum_{i=1}^{A} \alpha_i^t h_i$$

$$\alpha_i^t = \operatorname{softmax}(align(h_i, s_{t-1}))$$

Transformer

Multi-head attention using Scaled dot-product attention.

(Vaswani et al. 2017)

Attention is a mapping problem between key-value and query vectors.

$$MultiHead(V,K,Q) = f(Concat(H))$$

$$h_i = Att(f_i(V), f_i(K), f_i(Q))$$

$$Att(V, K, Q) = softmax(\frac{QK^{T}}{\sqrt{d_k}})V$$

Sequence-to-Sequence (seq2seq) Models

RNN-based

Global attention.

(Bahdanau et al. 2015)

Attention mechanism creates context vectors.

$$c_t = Att(H, s_{t-1}) = \sum_{i=1}^{A} \alpha_i^t h_i$$

$$\alpha_i^t = \operatorname{softmax}(align(h_i, s_{t-1}))$$

Transformer

Multi-head attention using Scaled dot-product attention.

(Vaswani et al. 2017)

Attention is a mapping problem between key-value and query vectors.

MultiHead(V,K,Q) = f(Concat(H))

$$h_i = Att(f_i(V), f_i(K), f_i(Q))$$

$$Att(V, K, Q) = softmax(\frac{QK^{T}}{\sqrt{d_k}})V$$

2D CNN-based

Pervasive attention.

(Elbayad et al. 2018)

Source elements are re-encoded considering the generated output.

$$\alpha = \operatorname{softmax}(W_1 \tanh(H_L W_2))$$

$$H_L^{\rm Att} = \alpha H_L$$

Task Pipeline: Evaluation (1)

Encoder:

word1,word2,word3,word4...

Decoder:

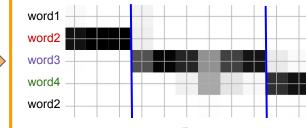
phn1,phn2,phn3,phn4,phn10,phn1...



Unit Discovery System



seq2seq system



Segmentation:

phn1 phn2 phn3 phn4 phn10 phn1 phn2

phn1phn2, phn3phn4phn10phn1, phn2

Alignment:

(phn1phn2, word2); (phn3phn4phn10phn1, word3); (phn2, word4)

Task Pipeline: Evaluation (2)

Encoder:

word1,word2,word3,word4...

Decoder:

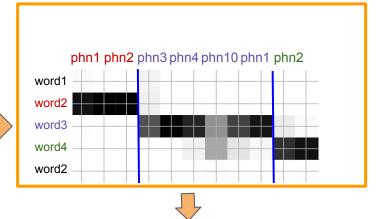
phn1,phn2,phn3,phn4,phn10,phn1...



Unit Discovery System



seq2seq system



Segmentation:

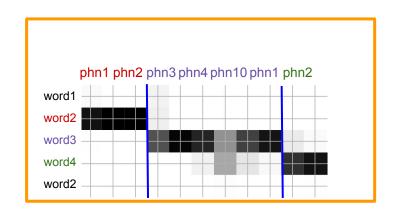
phn1phn2, phn3phn4phn10phn1, phn2

Alignment:

(phn1phn2, word2); (phn3phn4phn10phn1, word3); (phn2, word4)

Intrinsic Evaluation

How do we evaluate alignment quality without having gold (word-level) alignment information?

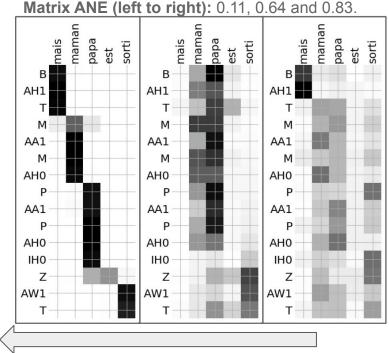


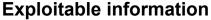
We use **Average Normalized Entropy** for assessing *alignment confidence*.



Intuition: sharper alignments are more informative.

Soft-alignment probability matrix: one probability distribution per line (target symbol)

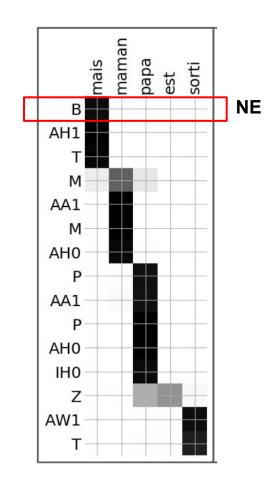






For every line in the matrix we compute normalized entropy (NE).

$$NE(t_i, s) = -\sum_{j=1}^{|s|} P(t_i, s_j) \cdot \log_{|s|} (P(t_i, s_j))$$

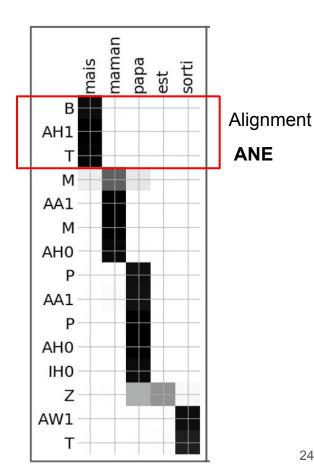




For every line in the matrix we compute normalized entropy (NE). We average over sets of distributions.

$$NE(t_i, s) = -\sum_{j=1}^{|s|} P(t_i, s_j) \cdot \log_{|s|} (P(t_i, s_j))$$

$$ANE(t, s) = \frac{\sum_{i=1}^{|t|} NE(t_i, s)}{|t|}$$

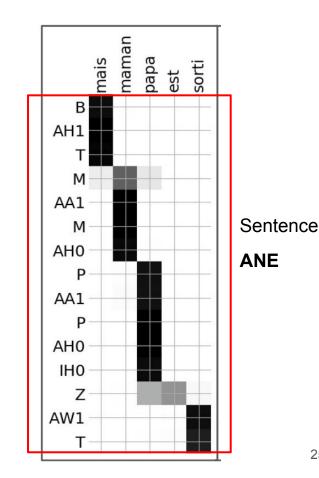




For every line in the matrix we compute normalized entropy (NE). We average over sets of distributions.

$$NE(t_i, s) = -\sum_{j=1}^{|s|} P(t_i, s_j) \cdot \log_{|s|} (P(t_i, s_j))$$

$$ANE(t, s) = \frac{\sum_{i=1}^{|t|} NE(t_i, s)}{|t|}$$

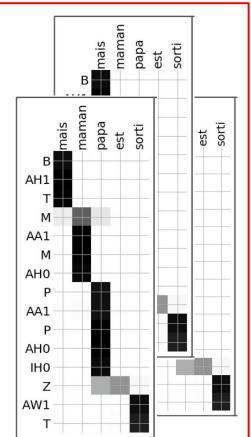




For every line in the matrix we compute normalized entropy (NE). We average over sets of distributions.

$$NE(t_i, s) = -\sum_{j=1}^{|s|} P(t_i, s_j) \cdot \log_{|s|} (P(t_i, s_j))$$

$$ANE(t, s) = \frac{\sum_{i=1}^{|t|} NE(t_i, s)}{|t|}$$



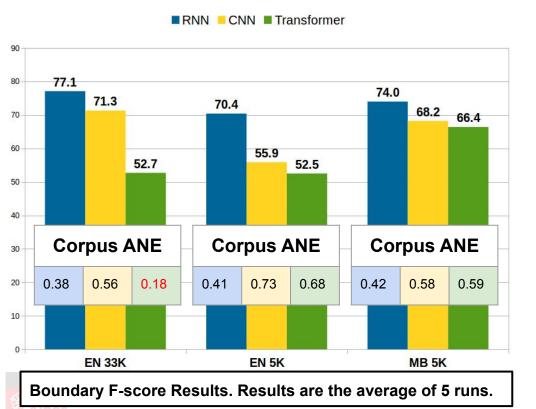
Corpus

ANE

Results



Unsupervised Word Segmentation Results



- RNN-based model performed the best
- Models with lower Corpus ANE reached better segmentation results (strong negative Pearson's correlation verified for all models)

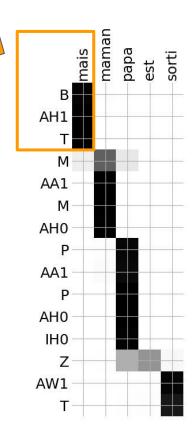
Experimental settings and corpora: https://gitlab.com/mzboito/attention_study

ANE Application: Exploiting the Alignments

Aligned pair

We accumulate ANE for all the (discovered type, aligned information) pairs discovered by our best 5K models in the whole corpus

This allow us to rank discovered alignments by their confidence.





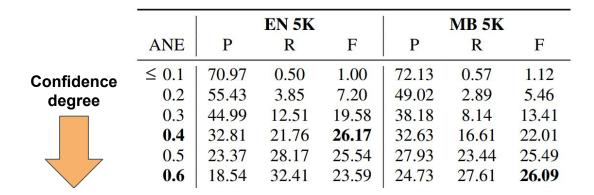
ANE over (discovered type, aligned information) pairs for the entire dataset



 High-confidence alignments cover a small portion of the corpus, but have high precision



ANE over (discovered type, aligned information) pairs for the entire dataset



 Accepting a wider confidence window, we decrease precision results, but increase coverage



Alignment ANE can be used for filtering the resulting lexicon, increasing type discovery results

(4)							
		EN 5K			MB 5K		
	ANE	P	R	F	P	R	F
2	≤ 0.1	70.97	0.50	1.00	72.13	0.57	1.12
	0.2	55.43	3.85	7.20	49.02	2.89	5.46
	0.3	44.99	12.51	19.58	38.18	8.14	13.41
	0.4	32.81	21.76	26.17	32.63	16.61	22.01
	0.5	23.37	28.17	25.54	27.93	23.44	25.49
	0.6	18.54	32.41	23.59	24.73	27.61	26.09
	0.7	16.23	34.34	22.04	23.00	30.12	26.08
	0.8	15.21	35.16	21.23	22.17	30.95	25.84
	0.9	15.01	35.31	21.06	22.06	31.05	25.80
	all	15.01	35.34	21.07	22.06	31.05	25.80



- → Low ANE: more frequently correct types, good alignment
- → High ANE: more frequently incorrect types and alignments artifacts

	Phoneme Sequence	Grapheme	Aligned Information
1	SER1	sir	
2	HHAH1SH	hush	chut
3	FIH1SHER0	fisher	fisher
4	KLER1K	clerk	clerc
5	KIH1S	kiss	embrasse
6	GRIH1LD	grilled	grilled
7	WUH1D	would	m'ennuierais
8	HHEH1LP	help	aidez
9	DOW1DOW0	dodo	dodo
10	KRAE1BZ	crabs	crabes

	Phoneme Sequence	Grapheme	Aligned Information	
1	AH0	а	convenablement	
2	IH1	Not a word	ah	
3	D	Not a word	riant	
4	N	Not a word	obéit	
5	YUW1	you	diable	
6	IH1	Not a word	qu'en	
7	AE1T	at	laquelle	
8	Z	Not a word	bas	
9	YUW1P	Not a word		
10	L	Not a word	parfaitement	

Conclusion



Concluding...

For low-resource,

- RNN-based model outputs the most easily exploitable soft-alignments
- → Transformer soft-alignment exploitation remains a challenge (what are the heads learning?¹,²)

Concluding...

For low-resource,

- → RNN-based model outputs the most easily exploitable soft-alignments
- → Transformer soft-alignment exploitation remains a challenge (what are the heads learning?^{1,2})

We introduced Average Normalized Entropy (ANE) showing:

- Its correlation to the segmentation scores
- Its usefulness for filtering alignments





Thank you! Questions?

Contact: marcely.zanon-boito@univ-grenoble-alpes.fr



Experimental settings and corpora: https://gitlab.com/mzboito/attention_study



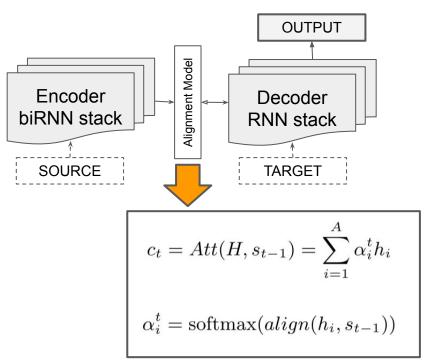
Empirical Evaluation of Sequence-to-Sequence Models for Word Discovery in Low-resource Settings

Marcely Zanon Boito¹, Aline Villavicencio² and Laurent Besacier¹

¹LIG, Univ. Grenoble Alpes, **France** ²CSEE, Univ. of Essex, **UK** ³INF, UFRGS, **Brazil**

Models: RNN

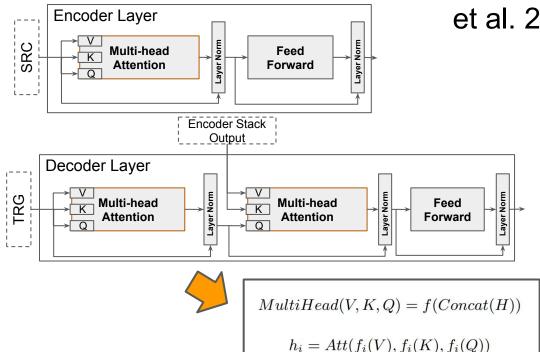
Global attention from Bahdanau et al. 2015



Attention appears at the form of a context vector for a decoder step *t*.

It is computed using the set of source annotations H and the last state of the decoder network s_{t-1} (translation context).

Models: Transformer



 $Att(V, K, Q) = softmax(\frac{QK^T}{\sqrt{d_*}})$

Multi-head attention from Vaswani et al. 2017

Attention is a mapping problem: given a pair of key-value vectors and a query vector, the goal is to compute the weighted sum of the key-values.

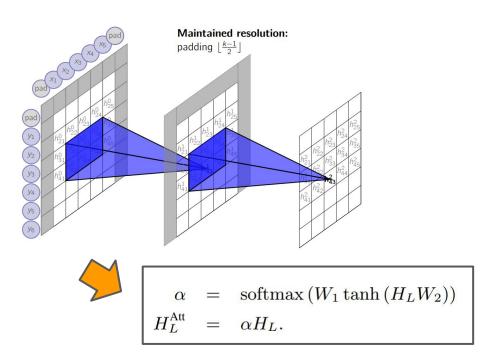
Weights are learned by compatibility functions between key-query pairs using Scaled dot-product (SDP) for each head in H.

Multi-head attention: SDP for several heads, which are then concatenated and projected again to yield the final result.



Models: CNN

Pervasive attention form Elbayad et al. 2018



Source and target are encoded jointly, what acts an an attention-like mechanism since individual source elements are re-encoded as the output is generated.

Attention weight tensor α is computed from the last activation tensor H_L , to pool the elements of the same tensor along the source dimension.