Pierre Godard<sup>1</sup>, **Marcely Zanon Boito<sup>2</sup>**, Lucas Ondel<sup>3</sup>, Alexandre Berard<sup>2,4</sup>, François Yvon<sup>1</sup>, Aline Villavicencio<sup>5</sup> and Laurent Besacier<sup>2</sup>

<sup>1</sup>LIMSI, CNRS, Univ. Paris-Saclay, **France**<sup>2</sup>LIG, Univ. Grenoble Alpes, **France**<sup>3</sup>BUT, Brno, **Czech Republic**<sup>4</sup>CRISTAL, Univ. Lille, **France**<sup>5</sup>CSEE, Univ. of Essex, **UK** 

#### **OUTLINE**

- 01. MOTIVATION
- 02. APPROACH
- 03. RESULTS

### **MOTIVATION**

### **Computational Language Documentation**

- → 50 to 90% of the currently spoken language will go extinct before 2100\*
- Manually documenting all these languages is infeasible



### Computational Language Documentation (CLD)

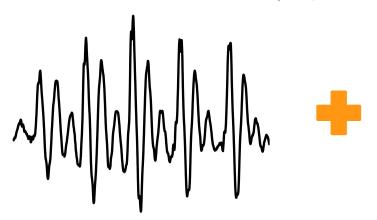
- → 50 to 90% of the currently spoken language will go extinct before 2100\*
- Manually documenting all these languages is infeasible



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

### **Endangered Languages Corpora**

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



#### **Translations**

to a well-documented language<sup>1</sup>

#### We focus on **UNSUPERVISED WORD SEGMENTATION**.

From speech

The system must output timestamps delimiting stretches of speech corresponding to real words in the language

#### We focus on **UNSUPERVISED WORD SEGMENTATION**.

→ From speech,

→ The system must output timestamps delimiting stretches of speech corresponding to real words in the language;

→ Slightly more favorable setup: the speech utterances are multilingually grounded (text translation in another language is available)

### THE TASK: Unsupervised Word Segmentation

from Speech using Attention

We focus on **UNSUPERVISED WORD SEGMENTATION**.

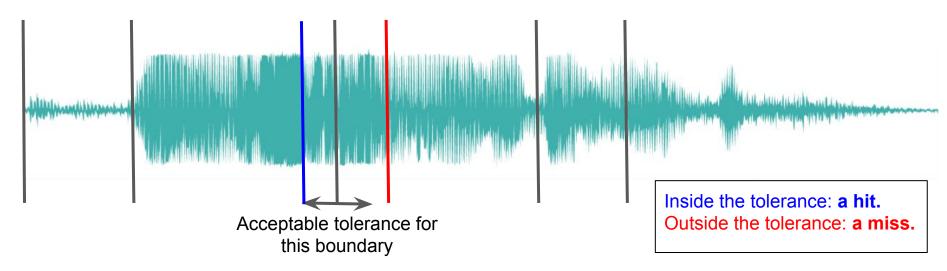


The tolerance window is defined on the **Zero Resource Challenge 2017 Track 2.** 

### THE TASK: Unsupervised Word Segmentation

from Speech using Attention

We focus on **UNSUPERVISED WORD SEGMENTATION**.



The tolerance window is defined on the Zero Resource Challenge 2017 Track 2.

#### **CONTRIBUTION:**

- First attempt of performing attentional (neural) word segmentation on speech
  - Previously proposed: a model working from symbolic level¹ (not speech)
- Low-resource setup, using only 5k sentences of the Mboshi-French parallel corpus<sup>2</sup>

### **OUR APPROACH**

#### **BACKGROUND:**

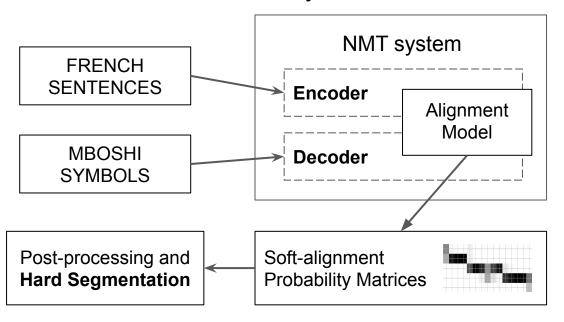
Attention-based encoder decoder models for Neural Machine Translation (NMT) are known to **jointly align and translate** a source into a target language<sup>1</sup>

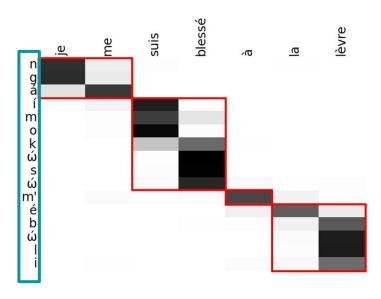
We use soft-alignment probability matrices learned during training to segment<sup>2</sup>

<sup>&</sup>lt;sup>1</sup> D. Bahdanau et al. **Neural Machine Translation by Jointly Learning to Align and Translate.** ICLR 2015.

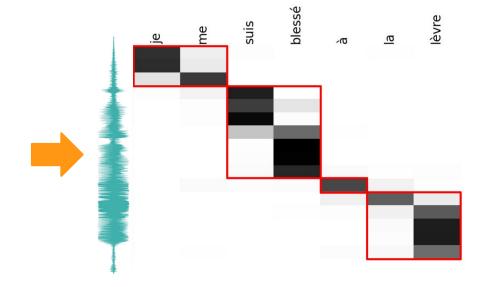
<sup>&</sup>lt;sup>2</sup> MZ Boito et al. Unwritten Languages Demand Attention too! Word Discovery Using Encoder-Decoder Models. ASRU 2017.

→ NMT systems¹ are trained with only 5k sentences





We would like to do the same, but from speech!

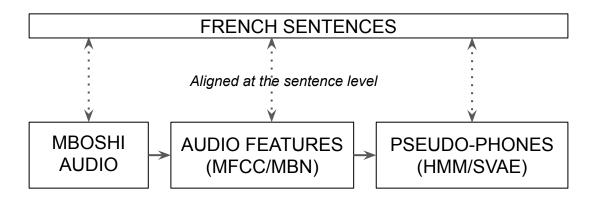


Problem: Infeasible training the model directly from speech with only 5k sentences

**Solution:** to extract pseudo-phones before training the network

We investigate Acoustic Unit Discovery (AUD) using two different audio feature extraction methods

### **Acoustic Unit Discovery (AUD)**

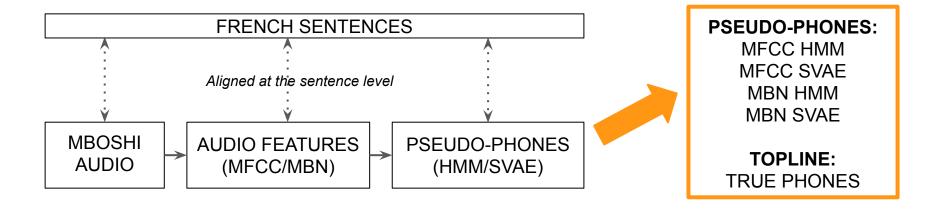


Two AUD models based on Bayesian Non-parametric HMM<sup>1</sup>

<sup>&</sup>lt;sup>1</sup> Implementation available a github.com/iondel/amdtk

Variational Inference for Acoustic Unit Discovery; L Ondel, L Burget, J Černocký; SLTU 2016.

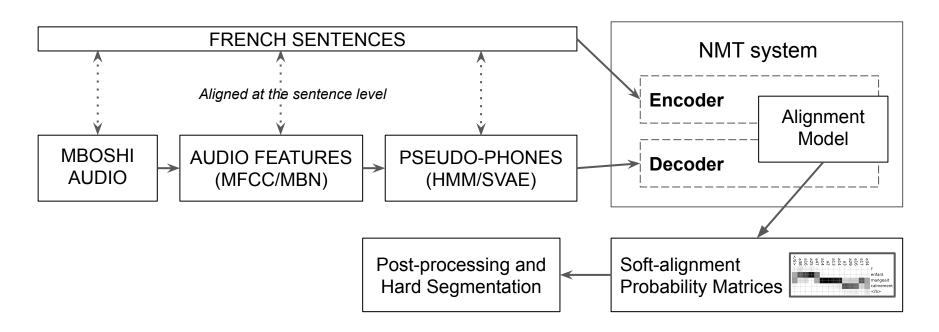
### **Acoustic Unit Discovery (AUD)**

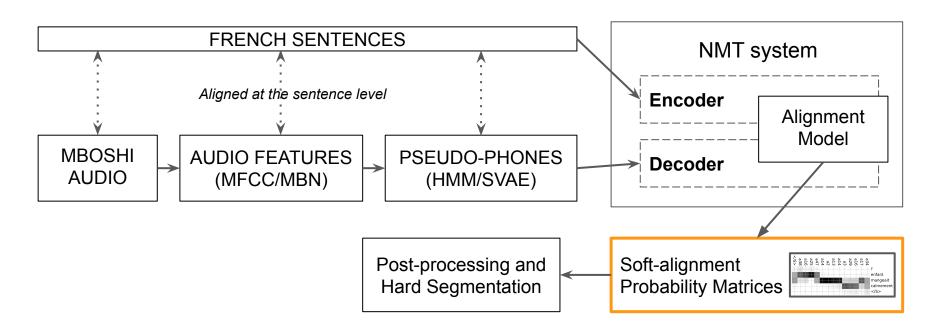


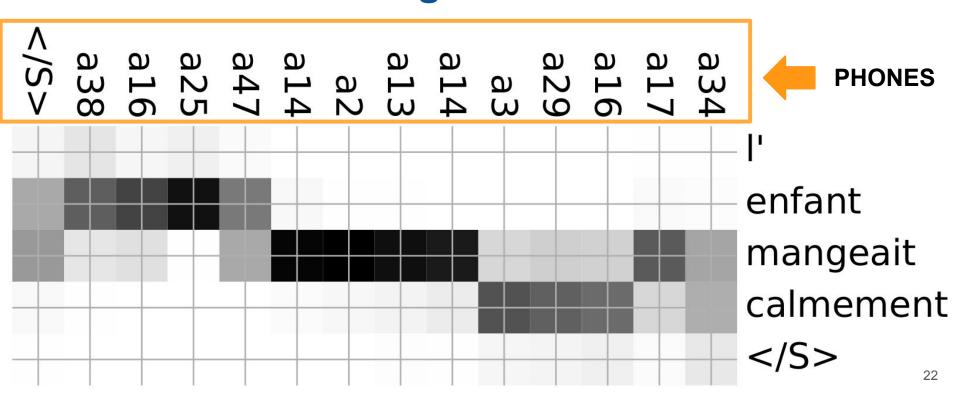
Two AUD models based on Bayesian Non-parametric HMM<sup>1</sup>

<sup>&</sup>lt;sup>1</sup> Implementation available a github.com/iondel/amdtk

Variational Inference for Acoustic Unit Discovery; L Ondel, L Burget, J Černocký; SLTU 2016.





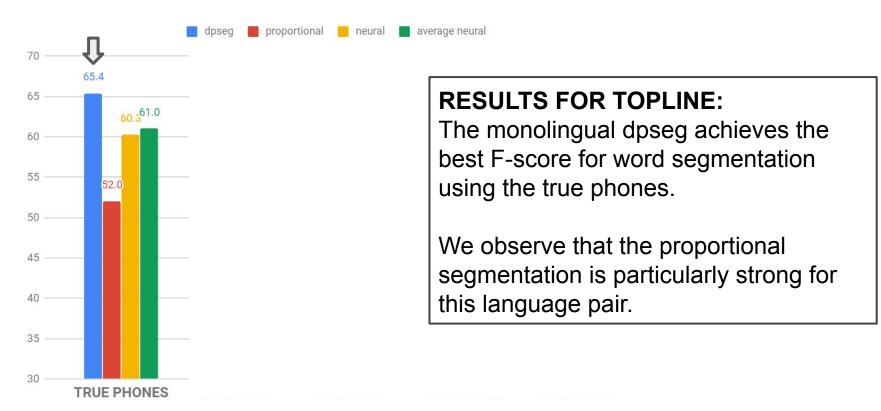


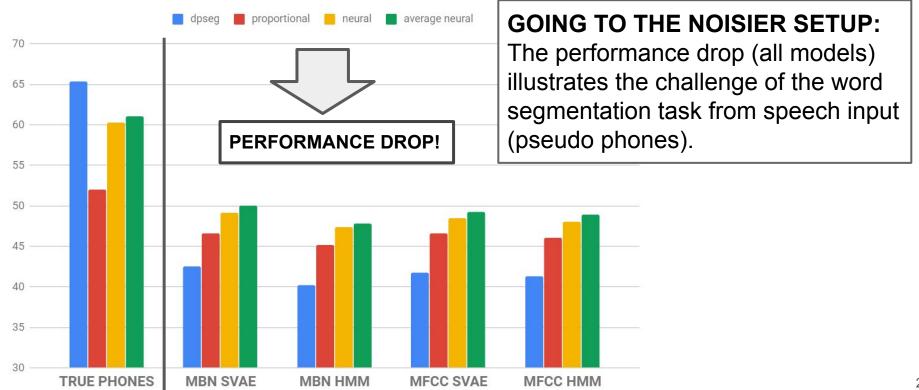
### RESULTS

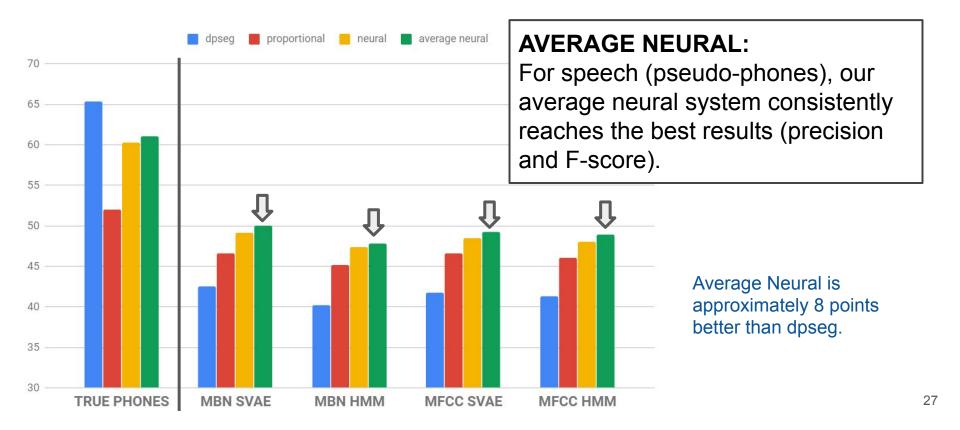
### **Baselines Comparison**

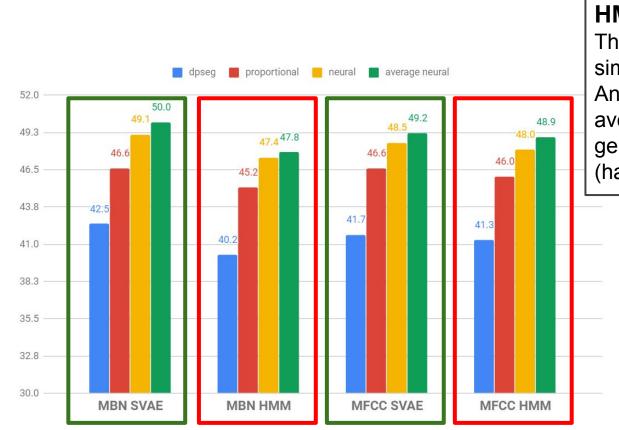
- Dpseg (Dirichlet process-based bigram LM¹) (monolingual)
- Proportional Segmentation (bilingual)
- → Neural Word Segmentation<sup>2</sup> (bilingual)
  Results are averaged over 5 runs with different splits.
- → Average Neural Word Segmentation (bilingual)
  Results are obtained through averaging 5 different soft-alignment matrices for each sentence.

S. Goldwater. A Bayesian framework for word segmentation: Exploring the effects of context. Cognition. 2009.





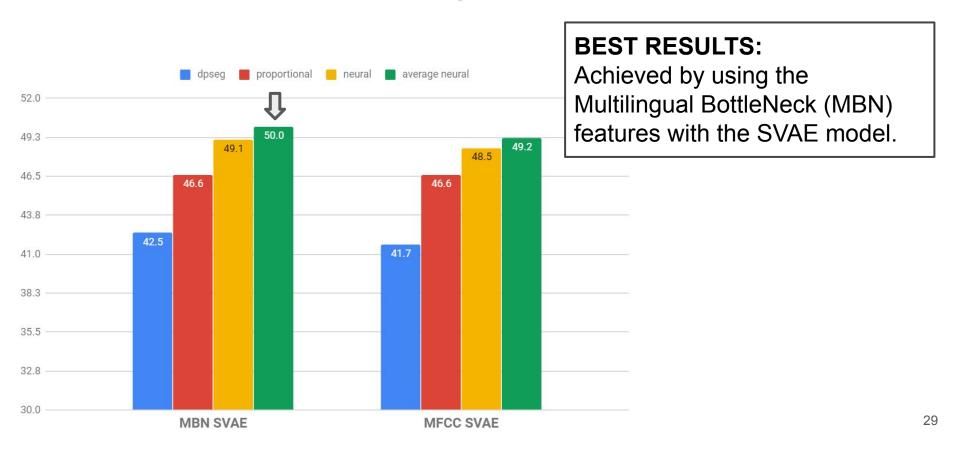




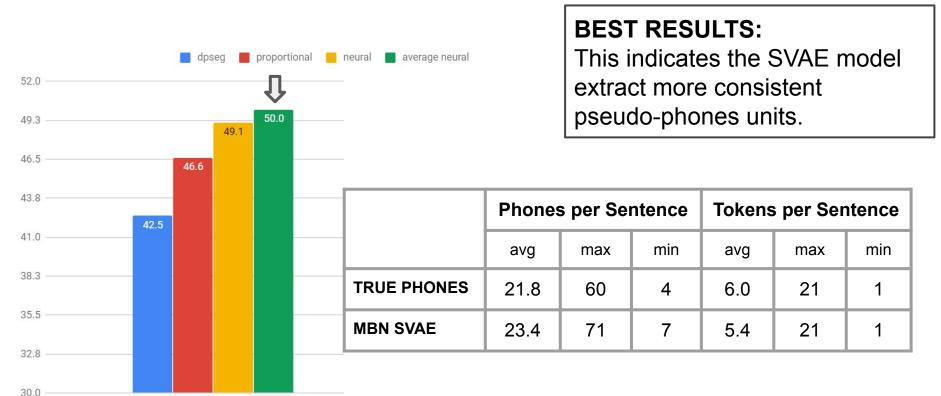
#### HMM RESULTS ARE WORSE:

This can be explained by the simplicity of the model.
Another explanation is the higher average number of pseudo-phones generated by sentence on HMMs (harder to segment).

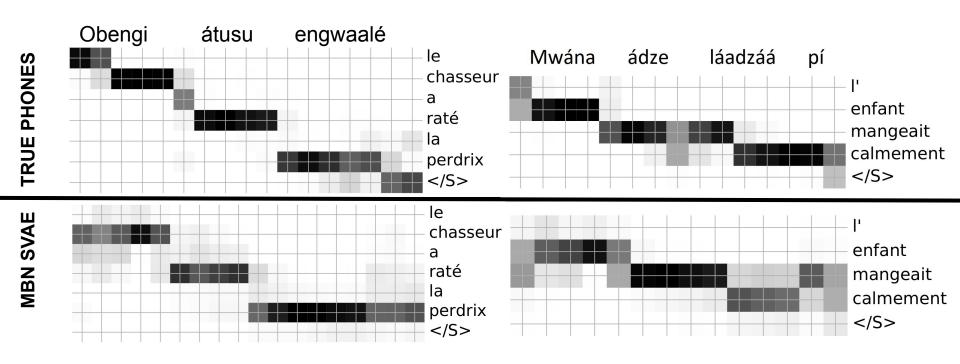
It's a 3 points difference between HMMs and SVAEs.



**MBN SVAE** 



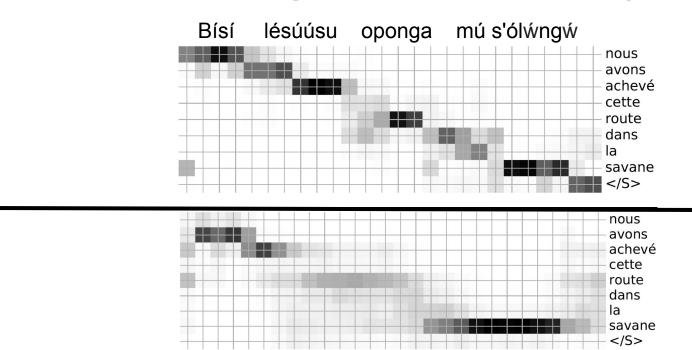
### **Example: Soft-alignment Probability Matrices**



**Figure:** Successful segmentation, examples of soft-alignment probability matrices. True Phones on top and MBN SVAE setup in the bottom.

31

### **Example: Soft-alignment Probability Matrices**



TRUE PHONES

**MBN SVAE** 

**Figure:** Failed segmentation, example of soft-alignment probability matrices. True Phones on top and MBN SVAE setup in the bottom.

### CONCLUSION

#### Conclusion

- Promising results for word segmentation from speech, outperforming two baselines in noisy (pseudo-phones) setup
- A deeper analysis of the world clusters obtained is needed to better understand how AUD affects the word discovery task
- → Word type results need improvement: 30.7% true phones, 14.1% best pseudo-phones setup

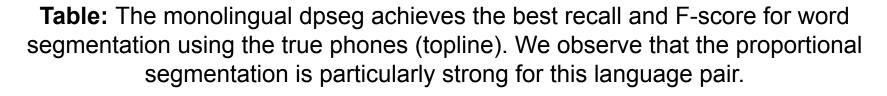
# Thank you! Questions?

Contact: marcely.zanon-boito@univ-grenoble-alpes.fr and pierre.godard@limsi.fr

Pierre Godard<sup>1</sup>, **Marcely Zanon Boito<sup>2</sup>**, Lucas Ondel<sup>3</sup>, Alexandre Berard<sup>2,4</sup>, François Yvon<sup>4</sup>, Aline Villavicencio<sup>5</sup> and Laurent Besacier<sup>2</sup>

<sup>1</sup>LIMSI, Univ. Paris-Saclay, **France**<sup>2</sup>LIG, Univ. Grenoble Alpes, **France**<sup>3</sup>BUT, Brno, **Czech Republic**<sup>4</sup>CRISTAL, Univ. Lille, **France**<sup>5</sup>CSEE, Univ. of Essex, **UK** 

AUD feat.	AUD	dpseg			proportional			neural			average neural		
	model	Р	R	F	Р	R	F	Р	R	F	Р	R	F
MFCC	HMM	27.9	80.2	41.3	42.6	49.9	46.0	51.6	44.9	48.0	55.5	43.7	48.9
MFCC	SVAE	29.8	69.1	41.7	42.2	51.9	46.6	52.7	45.0	48.5	55.7	44.1	49.2
MBN	HMM	27.8	72.6	40.2	42.5	48.1	45.2	50.8	44.5	47.4	54.1	42.9	47.8
MBN	SVAE	30.0	72.9	42.5	42.5	51.6	46.6	57.2	43.0	49.1	60.6	42.5	50.0
TRUE F	PHONES	53.8	83.5	65.4	44.5	62.6	52.0	60.5	59.9	60.3	62.8	59.3	61.0



AUD	AUD model	dpseg			proportional				neural		average neural		
feat.		Р	R	F	Р	R	F	Р	R	F	Р	R	F
MFCC	НММ	27.9	80.2	41.3	42.6	49.9	46.0	51.6	44.9	48.0	55.5	43.7	48.9
MFCC	SVAE	29.8	69.1	41.7	42.2	51.9	46.6	52.7	45.0	48.5	55.7	44.1	49.2
MBN	НММ	27.8	72.6	40.2	42.5	48.1	45.2	50.8	44.5	47.4	54.1	42.9	47.8
MBN	SVAE	30.0	72.9	42.5	42.5	51.6	46.6	57.2	43.0	49.1	60.6	42.5	50.0
TRUE F	PHONES	53.8	83.5	65.4	44.5	62.6	52.0	60.5	59.9	60.3	62.8	59.3	61.0

**Table:** The performance drop of all models illustrates the challenge of the word segmentation task from speech input (pseudo phones).

AUD feat.	AUD model	dpseg			proportional				neural		average neural		
		Р	R	F	Р	R	F	Р	R	F	Р	R	F
MFCC	НММ	27.9	80.2	41.3	42.6	49.9	46.0	51.6	44.9	48.0	55.5	43.7	48.9
MFCC	SVAE	29.8	69.1	41.7	42.2	51.9	46.6	52.7	45.0	48.5	55.7	44.1	49.2
MBN	НММ	27.8	72.6	40.2	42.5	48.1	45.2	50.8	44.5	47.4	54.1	42.9	47.8
MBN	SVAE	30.0	72.9	42.5	42.5	51.6	46.6	57.2	43.0	49.1	60.6	42.5	50.0
TRUE PHONES		53.8	83.5	65.4	44.5	62.6	52.0	60.5	59.9	60.3	62.8	59.3	61.0

**Table:** For speech (pseudo-phones), our average neural system consistently reaches the best results (precision and F-score).

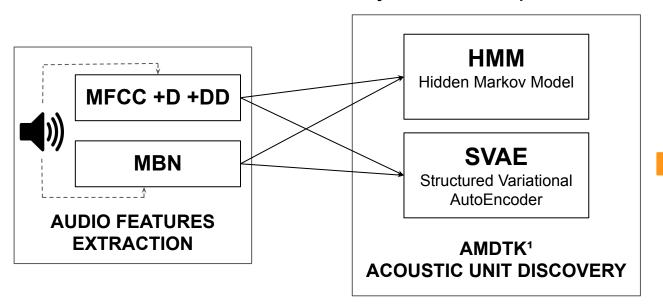
AUD feat.	AUD model	dpseg			proportional			neural			average neural		
		Р	R	F	Р	R	F	Р	R	F	Р	R	F
MFCC	НММ	27.9	80.2	41.3	42.6	49.9	46.0	51.6	44.9	48.0	55.5	43.7	48.9
MFCC	SVAE	29.8	69.1	41.7	42.2	51.9	46.6	52.7	45.0	48.5	55.7	44.1	49.2
MBN	нмм	27.8	72.6	40.2	42.5	48.1	45.2	50.8	44.5	47.4	54.1	42.9	47.8
MBN	SVAE	30.0	72.9	42.5	42.5	51.6	46.6	57.2	43.0	49.1	60.6	42.5	50.0
TRUE F	PHONES	53.8	83.5	65.4	44.5	62.6	52.0	60.5	59.9	60.3	62.8	59.3	61.0

**Table:** HMM models achieved worse results than SVAE models. This can be explained by the simplicity of the model. Another explanation would be the average number of pseudo-phones generated by sentence, which is higher on HMM models (harder to segment).

AUD	AUD model	dpseg			proportional			neural			average neural		
feat.		Р	R	F	Р	R	F	Р	R	F	Р	R	F
MFCC	НММ	27.9	80.2	41.3	42.6	49.9	46.0	51.6	44.9	48.0	55.5	43.7	48.9
MFCC	SVAE	29.8	69.1	41.7	42.2	51.9	46.6	52.7	45.0	48.5	55.7	44.1	49.2
MBN	нмм	27.8	72.6	40.2	42.5	48.1	45.2	50.8	44.5	47.4	54.1	42.9	47.8
MBN	SVAE	30.0	72.9	42.5	42.5	51.6	46.6	57.2	43.0	49.1	60.6	42.5	50.0
TRUE F	PHONES	53.8	83.5	65.4	44.5	62.6	52.0	60.5	59.9	60.3	62.8	59.3	61.0

**Table:** Best results were achieved by using the Multilingual BottleNeck (MBN) features with the SVAE model. This indicates this model extract more consistent pseudo-phones units.

**Two** AUD models based on Bayesian Non-parametric HMM:



MFCC HMM

**MFCC SVAE** 

**MBN HMM** 

MBN HMM

(These + translation are the NMT system input!) 42

<sup>&</sup>lt;sup>1</sup> Implementation available a github.com/iondel/amdtk

<sup>&</sup>lt;sup>2</sup> Variational Inference for Acoustic Unit Discovery; L Ondel, L Burget, J Černocký; 2016