



**HAL**  
open science

## ASR performance prediction on unseen broadcast programs using convolutional neural networks

Zied Elloumi, Laurent Besacier, Olivier Galibert, Juliette Kahn, Benjamin Lecouteux

► **To cite this version:**

Zied Elloumi, Laurent Besacier, Olivier Galibert, Juliette Kahn, Benjamin Lecouteux. ASR performance prediction on unseen broadcast programs using convolutional neural networks. ICASSP, 2018, Alberta, Canada. hal-02088829

**HAL Id: hal-02088829**

**<https://hal.science/hal-02088829v1>**

Submitted on 3 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ASR Performance Prediction on Unseen Broadcast Programs using Convolutional Neural Networks

Zied Elloumi, Laurent Besacier, Olivier Galibert, Juliette Kahn, Benjamin Lecouteux

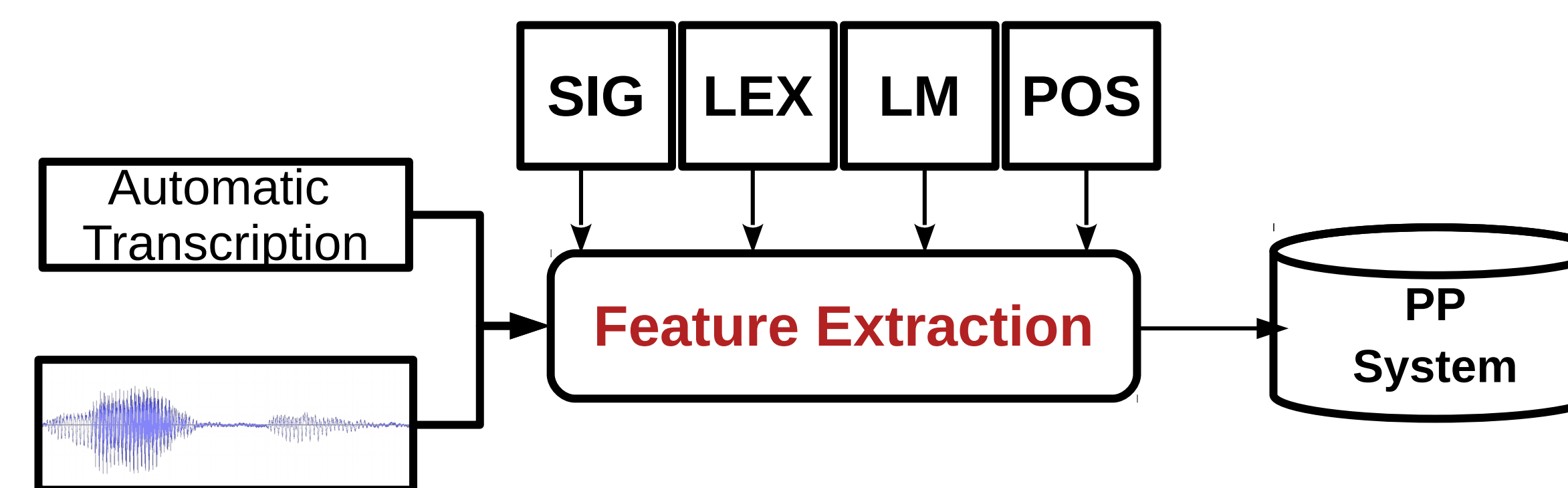


## Abstract

- Propose an heterogeneous French corpus dedicated to performance prediction task
- Compare two prediction approaches: regression (**engineered features**) based vs a new strategy based on convolutional neural networks (**learnt features**).
- The joint use of textual and signal features did not work for the regression baseline while the combination of inputs for CNNs leads to the best WER prediction performance.
- CNN prediction predicts the WER distribution on a collection of speech recordings

## Regression Baseline

- We used an existing tool named TransRater for baseline regression approach
- We adapted the TransRater tool from English to French that requires **engineered features** to predict the WER performance



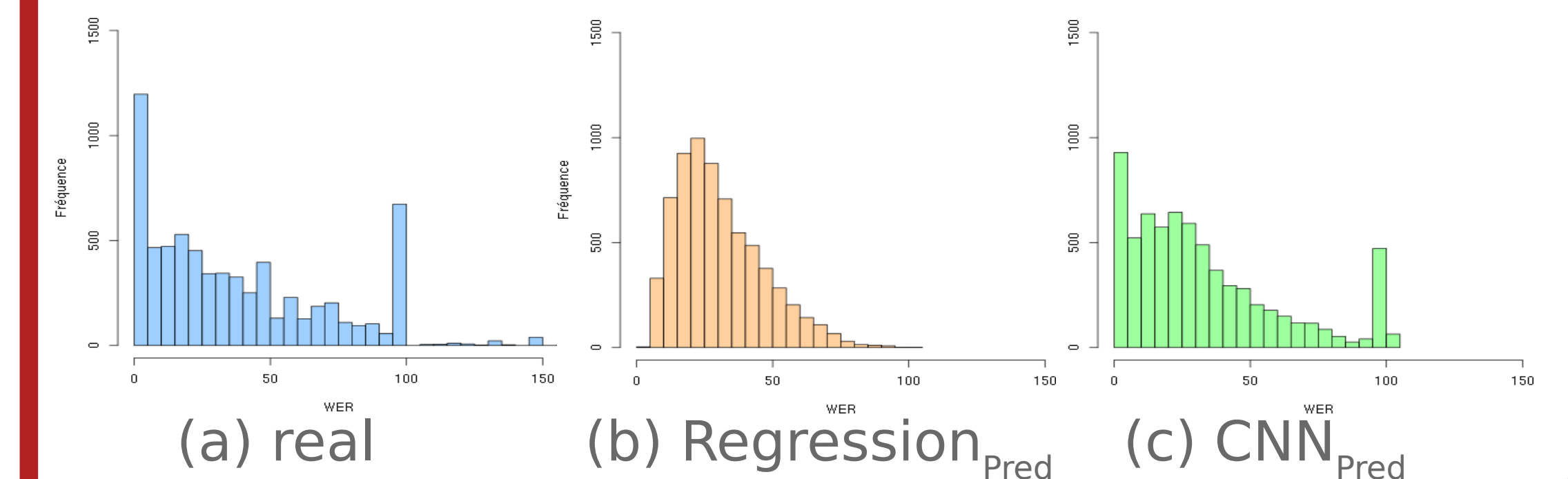
## Results & Analysis

Model	Input	MAE	$\tau$
<b>Textual features</b>			
<b>Regression</b>	POS+LEX+LM	22.01	<b>44.16</b>
<b>CNN<sub>Softmax</sub></b>	EMBED	<b>21.48</b>	38.91
<b>CNN<sub>ReLU</sub></b>	EMBED	22.30	38.13
<b>Signal features</b>			
<b>Regression</b>	SIG	25.86	23.36
<b>CNN<sub>Softmax</sub></b>	RAW-SIG	25.97	23.61
<b>CNN<sub>ReLU</sub></b>	RAW-SIG	26.90	21.26
<b>CNN<sub>Softmax</sub></b>	MEL-SPEC	29.11	19.76
<b>CNN<sub>ReLU</sub></b>	MEL-SPEC	26.07	24.29
<b>CNN<sub>Softmax</sub></b>	MFCC	<b>25.52</b>	<b>26.63</b>
<b>CNN<sub>ReLU</sub></b>	MFCC	26.17	25.41
<b>Textual and Signal features</b>			
<b>Regression</b>	POS+LEX+LM+SIG	21.99	45.82
<b>CNN<sub>Softmax</sub></b>	EMBED+RAW-SIG	<b>19.24</b>	<b>46.83</b>
<b>CNN<sub>ReLU</sub></b>	EMBED+RAW-SIG	20.56	45.01
<b>CNN<sub>Softmax</sub></b>	EMBED+MEL-SPEC	20.93	40.96
<b>CNN<sub>ReLU</sub></b>	EMBED+MEL-SPEC	20.93	44.38
<b>CNN<sub>Softmax</sub></b>	EMBED+MFCC	19.97	44.71
<b>CNN<sub>ReLU</sub></b>	EMBED+MFCC	20.32	45.52

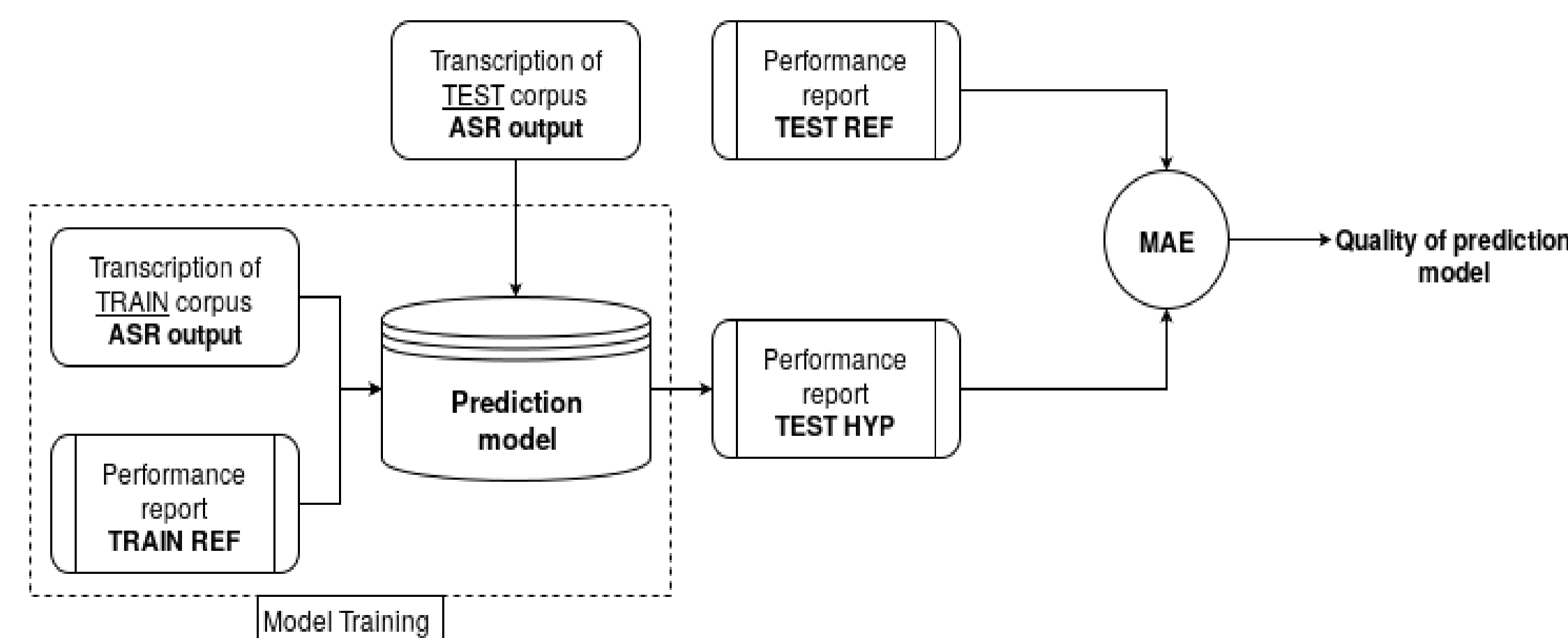
- Analysis of predicted WERs

	NS	S	NS + S
WER <sub>REF</sub>	21.47	38.83	31.20
WER <sub>Pred</sub> Regression	<b>22.08</b>	28.72	25.82
WER <sub>Pred</sub> CNN <sub>Softmax</sub>	18.93	<b>33.99</b>	<b>27.37</b>
#Utterances	3,1k	3,7k	6,8k
#Words <sub>REF</sub>	49.8k	63.3k	113,1k

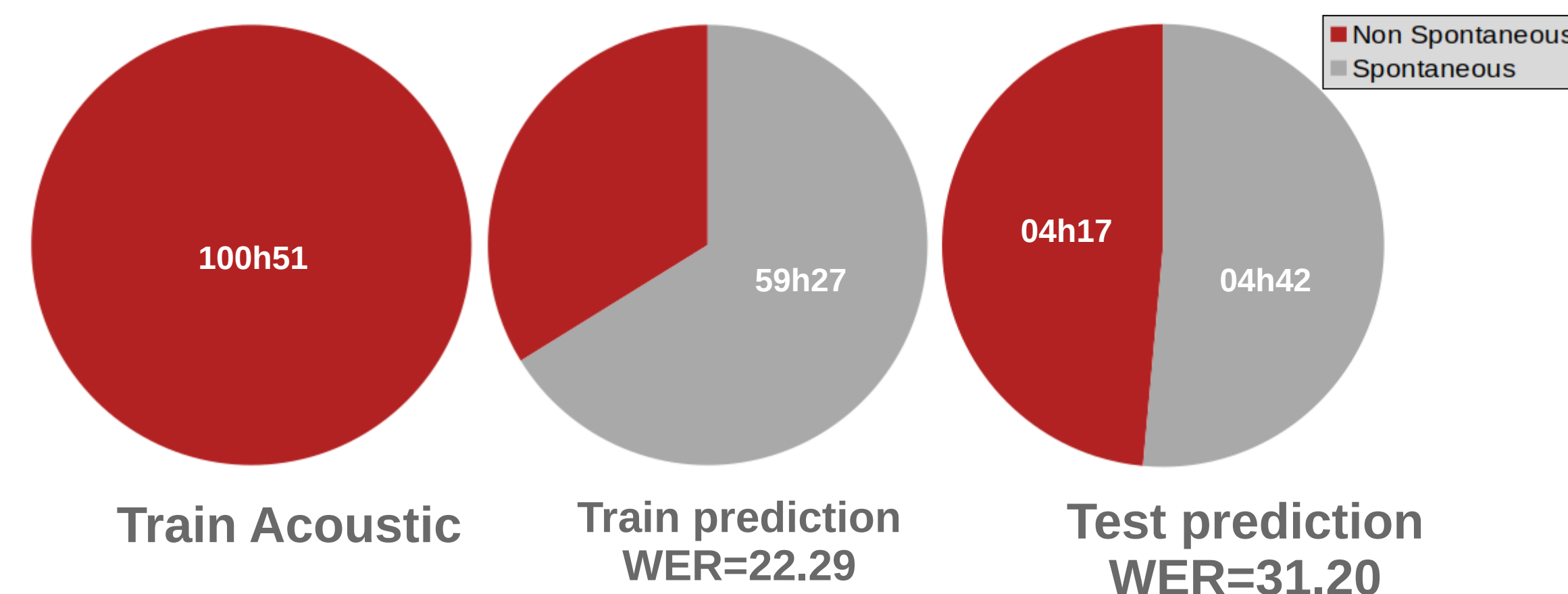
- Distribution of speech turns according to their WER



## Evaluation framework

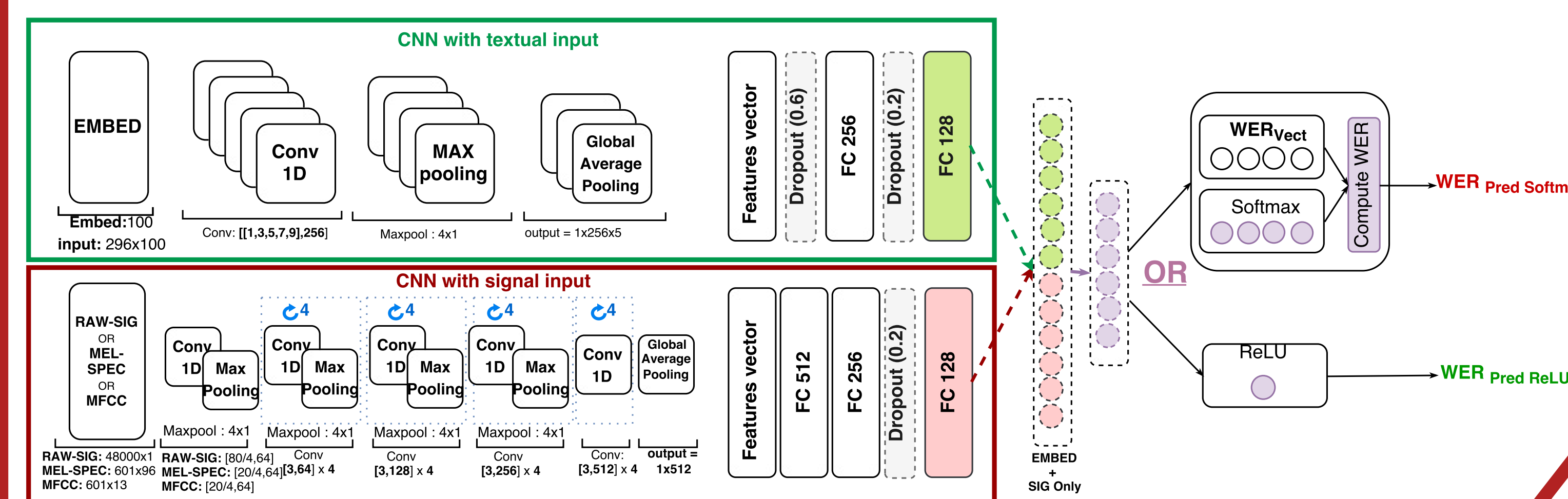


- We built our own French ASR system based on the KALDI toolkit
- We used a French data from different broadcast collections: ESTER, ETAPE, REPERE, Quaero



## Our proposed approach

- We proposed a new End-to-End performance prediction system based on CNNs.
- Several approaches to encode speech signal
- Input**: textual, speech signal or the both textual+ speech signal
- Output**: two proposed methods to predict a continuous value ( CNN<sub>Softmax</sub> and CNN<sub>ReLU</sub> )



## Conclusion

- We presented an evaluation framework for evaluating ASR performance prediction on unseen broadcast programs
- CNNs were very efficient encoding both textual (ASR transcript) and signal to predict WER.
- Future work**: Analyze the learnt representations of our performance prediction system



Contact :  
Zied.elloumi@lne.fr

