

# Les modèles pré-entraînés à l'épreuve des langues rares : expériences de reconnaissance de mots sur la langue japhug (sino-tibétain)

Séverine Guillaume, Guillaume Wisniewski, Cécile Macaire, Guillaume Jacques, Alexis Michaud, Benjamin Galliot, Maximin Coavoux, Solange Rossato, Minh-Châu Nguyễn, Maxime Fily

## Contexte / Objectifs à long terme

Faire bénéficier la documentation et l'étude des langues rares des outils de pointe en TAL  
Travail interdisciplinaire associant linguistes de terrain, chercheurs en informatique...  
Enjeu pour l'étude de la parole : disponibilité de données fiables et abondantes sur la diversité des langues

## Une approche de reconnaissance en deux étapes

- ▶ Entraînement non-supervisé par XLSR-53<sup>1</sup> : reconnaissance du signal de parole sur 53 langues (56 000 h)
    - ▷ Modèle de représentation du signal
  - ▶ Pré-entraînement sur un sous-ensemble transcrit du corpus de japhug
    - ▷ utilisation de données de terrain déposées dans la collection Pangloss (audio + transcription phonémique)
- Puis application à la reconnaissance des phonèmes (& des espaces) sur de l'audio.  
↳ Évaluation à l'aune de deux métriques classiques : CER (character error rate) et WER (word error rate)

## Résultats

### Audio en langue japhug transcrit automatiquement :

tçe kuaɕuŋɣu tçe iɕqha @mingchao(u→.) uraŋɣ nu-tçu pɣɣu tçendɣre iɕqha nɣki  
@yanguo kɣti ɣɣɣɣɣɣ ɣu nuɣɣɣɣ nu ku, iɕqha nu, iɕqha nu uftsa nuɣu  
ɣɣɣɣɣɣ lusuɣɣɣɣ pɣɣɣɣ. tçe nu ɣɣɣɣ lusuɣɣɣɣ pɣɣɣɣ tçe, tçendɣre nɣkinu,  
sɣtɕha ra toɣtɕoɣloɣnu zo ɕti tçe, tçendɣre iɕqha nu, @shandong nuɣu urmi  
@zhangxiaobing kurrmi ci tuɣtɕe ukuɣɣzu ci pɣɣu, tɣtçu. ɕendɣre urɕaɕ nu  
uskhru muɣɣɣɣdi ɣɣɣɣ ma muɣoɣzu ri [...]

**Sortie d'un modèle entraîné sans ponctuation.**  
**En rouge : corrections manuelles du linguiste**

tçe kuaɕuŋɣu tçe iɕqha; @mingchao uraŋ nuɣu pɣɣu; tçendɣre iɕqha, nɣki;  
@yanguo kɣti ɣɣɣɣɣɣ ɣu; nuɣɣɣɣ nu ku, iɕqha nu(.→.) iɕqha nu, uftsa  
nuɣu ɣɣɣɣɣɣ lusuɣɣɣɣ pɣɣɣɣ. tçe nu ɣɣɣɣ lusuɣɣɣɣ pɣɣɣɣ tçe, tçendɣre,  
nɣkinu, sɣtɕha ra toɣtɕoɣloɣnu zo ɕti tçe, tçendɣre iɕqha nu, @shandong nuɣu,  
urmi @zhangxiaobing kurrmi ci, tuɣtɕe ukuɣɣzu ci pɣɣu, tɣtçu. ɕendɣre urɕaɕ  
nu uskhru muɣɣɣɣd(er,→i) ɣɣɣɣ ma muɣoɣzu ri [...]

**Sortie d'un modèle entraîné avec ponctuation.**  
**En rouge : corrections manuelles du linguiste**

### De bons résultats dès 2h à 3h de données :

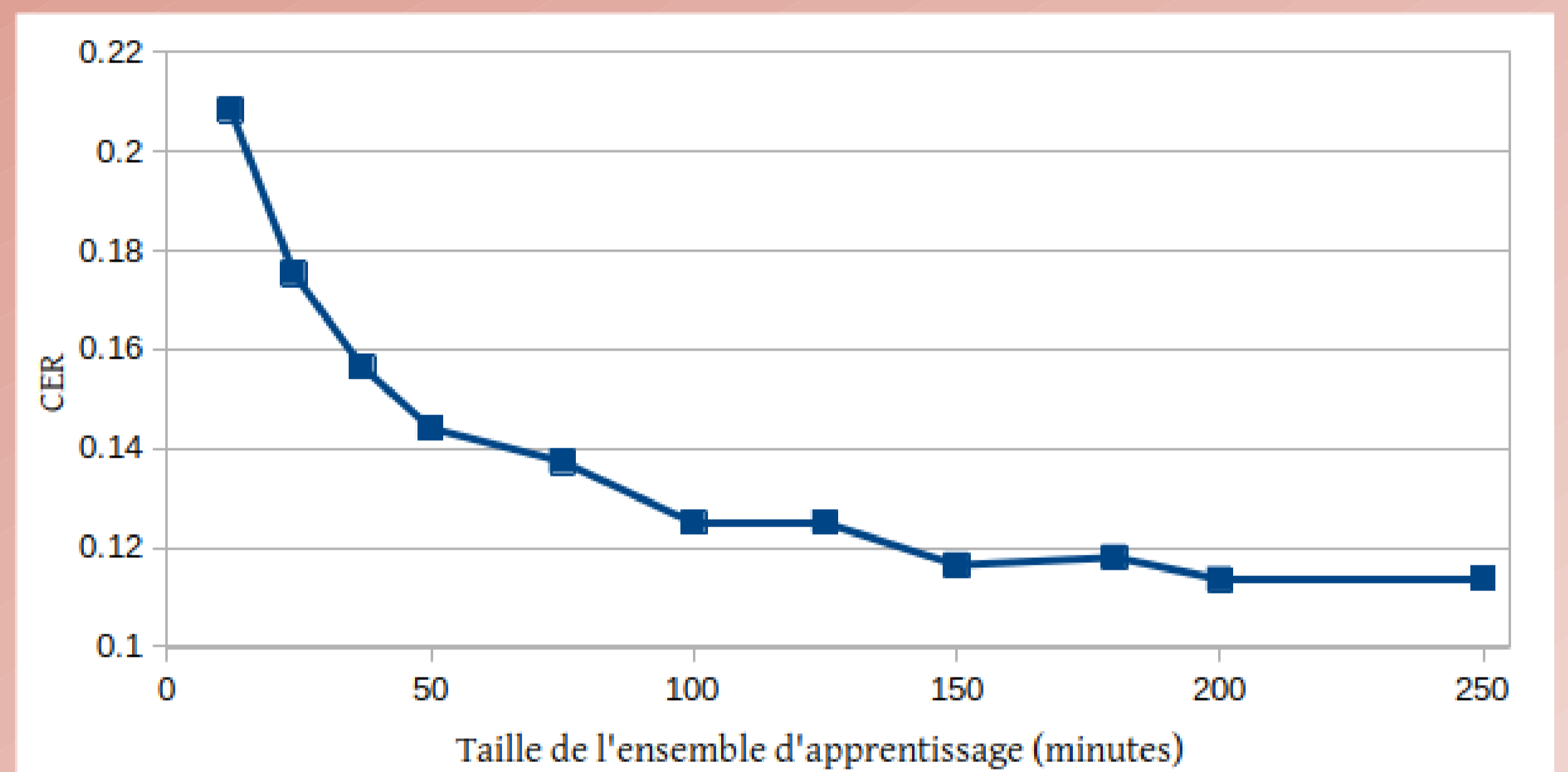


Figure 1 – Évolution des performances en fonction de la taille du corpus d'apprentissage

### Les mystères de l'apprentissage machine :

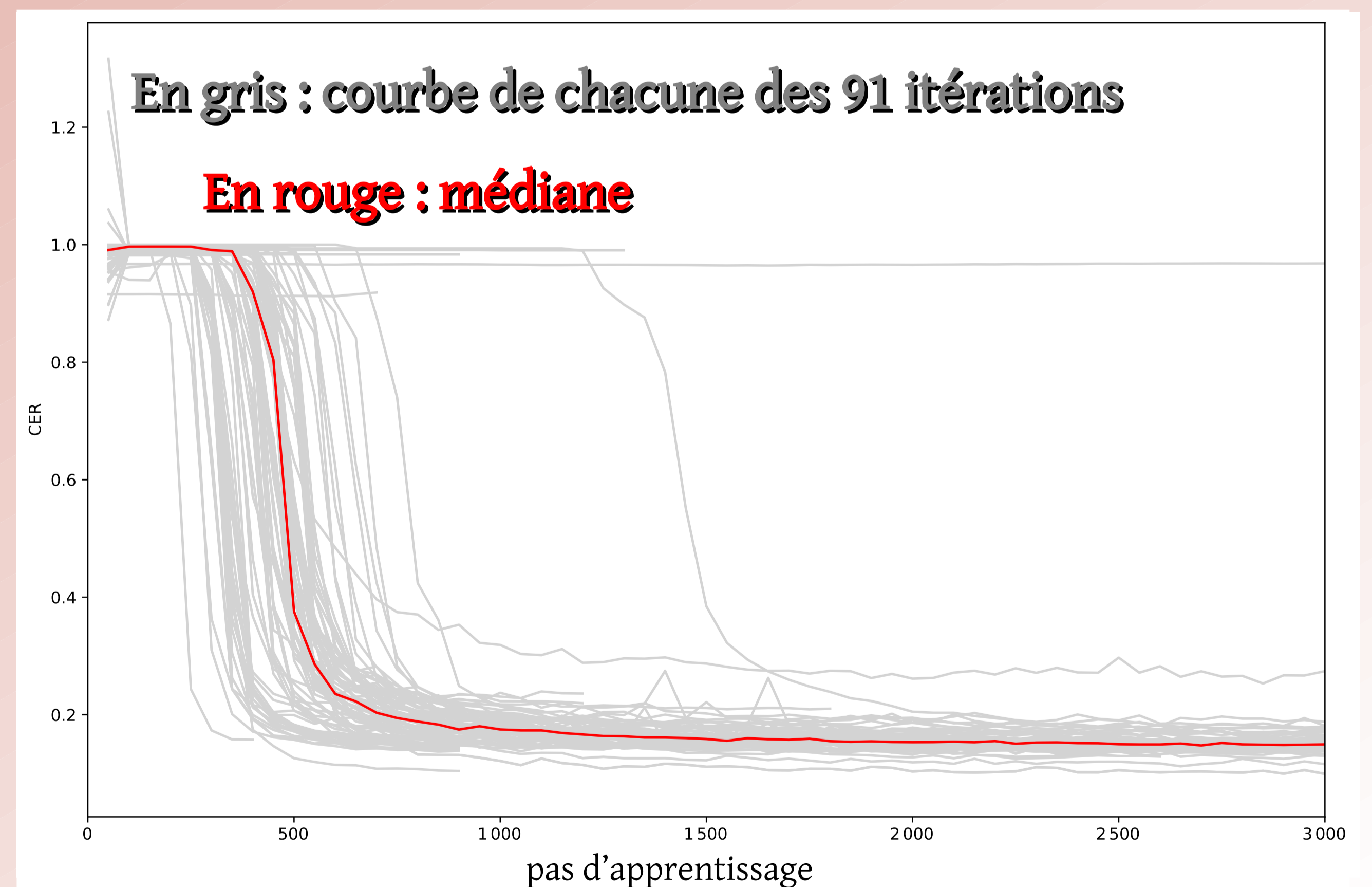


Figure 2 - CER sur l'ensemble de validation au cours de différentes itérations d'optimisation

## Perspectives

Gains en temps de transcription : après production "manuelle" d'un corpus de 2 à 3h, le modèle statistique produit une transcription qui sert de point de départ pour la suite.

Figure 3 – Extraits Pangloss du corpus Pangloss<sup>2</sup> de langue japhug de Guillaume Jacques. À gauche : un travail complet de transcription-glose-traduction. À droite : transcription seule..

## Conclusions

- Un résultat très satisfaisant au plan qualitatif
- Des perspectives applicatives prometteuses pour les langues peu dotées
- Des thèmes porteurs pour la recherche interdisciplinaire en parole

<sup>1</sup> [https://huggingface.co/ttransformers/model\\_doc/xlsr\\_wav2vec2.html](https://huggingface.co/ttransformers/model_doc/xlsr_wav2vec2.html)  
<sup>2</sup> <https://pangloss.cnrs.fr>