

Report on Automatic Speech Recognition Experiment for Yoruba and French Languages Using 1 Hour of Data

Group Members:

Sewade Olaolu Ogun (Yoruba Male Speaker)

Tatiana Moteu (French Female Speaker)

Introduction

Automatic speech recognition (ASR) for high-resource languages is commonplace, but low-resource ASR is gaining interest in the community as humans (infants) require few examples to learn speech and language. Low-resource ASR requires learning with few examples which is an interesting area of research. This project aims to explore transfer learning for low-resource ASR by utilizing a pretrained Contrastive Predictive Coding (CPC) model and fine-tuning it on Yoruba and French Language ASRs given only 1 hour of data for each language.

Data Preparation

Recording was done using the Lig-Aikuma Android app. It is an easy-to-use app with

a good interface for recording and elicitation. In the elicitation mode, a text is displayed on screen while the speaker reads out the text carefully. Dataset for the experiments are hosted at [1] and [2] for both Yoruba and French datasets.

Preprocessing

Raw speech dataset is split into;

- 1 hour of speech data (split into train and val (80%-20% split))
- 1 hour of speech as test set

The audio files are then preprocessed with their respective text files. Text files are linked to the audio files and converted to character strings as required by the CPC model.

Table 1: Training Results for Yoruba Language Experiments

		Train	Val		Test
Model	Parameters	Loss	Loss	CER	CER
CPC trained from scratch	$lr = 2e-4$ epochs=80, Early stopping, BeamSearch: beam_width=20, cutoff_top_n=20	1.8845	2.0444	0.6053	0.6258
CPC Finetuned without LM	$lr = 2e-4$ epochs=80, Early stopping, BeamSearch: beam_width=20, cutoff_top_n=20	0.4719	1.1641	0.3040	0.3476
CPC with 4-Gram LM	$lr = 2e-4$ epochs=80, Early stopping, BeamSearch: beam_width=20, cutoff_top_n=20	0.4885	1.2059	0.3072	0.3395
CPC with 5-Gram LM	$lr = 2e-4$ epochs=80, Early stopping, BeamSearch: beam_width=20, cutoff_top_n=20	0.3208	1.2234	0.3040	0.3270

Table 2: Training Results for French Language Experiments

		Train	Val		Test
Model	Parameters	Loss	Loss	CER	CER
CPC trained from scratch	lr = 2e-4 epochs=80, Early stopping, BeamSearch: beam_width=20,cutoff_top_n=20	2.6973	2.6705	0.7108	0.7177
CPC Finetuned without LM	lr = 2e-4 epochs=80, Early stopping, BeamSearch: beam_width=20,cutoff_top_n=20	0.7225	1.933	0.4792	0.4792

Summary of Results

The CPC model was initially trained without starting with pretrained weights to see how it will perform with few examples. Character Error Rates (CER) are measured for all the experiments as well training/validation losses. The pretrained CPC model is gotten from the CPC audio library from Facebook research [2]. Table 3 summarizes results of experiments comparing both languages when trained from scratch and when fine-tuned using a pretrained model. After 80 epochs, the CER was still very high without pretraining. By fine-tuning a pretrained model for the same experiment, we were able to reduce the word-error-rate by a factor of 2. For the yoruba experiments, the model indicates it can get better but will require more training time to converge to a reasonable result.

Training loss plateau for up to 20-30 epochs for most experiments before decreasing considerably for the remaining results show that language models can serve as a good prior to the beam search. 5-gram language models performed better than 4-gram language models on the test set while further performing better than models not using any language model. Table 2 shows the results of experiments

epochs. The ReduceOnPlateau learning rate scheduler helped in ensuring training continues.

Table 3: Comparison of Results for Yoruba and French Models

	Yoruba	French
Model	WER (Test)	WER (Test)
CPC trained from scratch	0.6258	0.7177
CPC Finetuned without LM	0.3476	0.4792

Furthermore, we experimented with 4-gram and 5-gram language models for the yoruba language CPC model. The language models were trained using the kenlm library following the procedure in [4]. Table 1 shows the results of all experiments performed for yoruba. The

performed on the french dataset. In comparison with the yoruba models, the WER on the french dataset was not as low as that of yoruba language on the same task. This might indicate that yoruba language is closer to english in language proximity than French.

Future Directions

Some other interesting ideas which we could not explore at the moment include;

- Increase the training data in steps and observing the WER due to data size
- Multilingual finetuning using both french and yoruba dataset

Challenges

The major challenge faced was in creating and preparing the datasets for this experiment. Getting the data in the format required for the dataloaders require that scripts had to be written to place the data in the right format.

Conclusion

In this project, we built an Automatic speech recognition system (character level) for both yoruba language and french language with only 1 hour of labeled data. The results show that pretrained models can be explored for transfer learning in speech.

- Use RNN language model or Transformer-based language model as language prior.
- Multilingual pretraining (using CPC).

Other experimental figures and plots are included in the appendix.

References

- [1] https://github.com/ogunlao/yoruba_speech_project
- [2] https://github.com/TatianaMoteuN/french_speech_rec
- [3] https://github.com/facebookresearch/CPC_audio
- [4] <https://kheafield.com/code/kenlm/>

Appendix Yoruba Experiment Plots

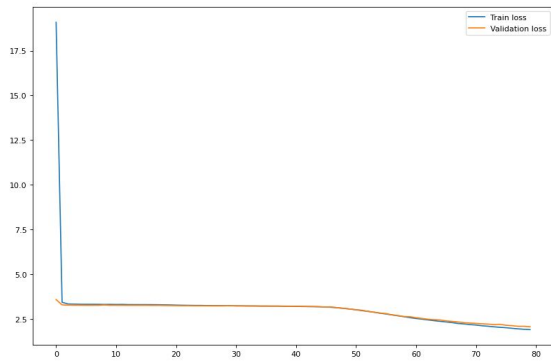


Fig 1: Train and Validation Loss for CPC Trained from Scratch and No Language Model

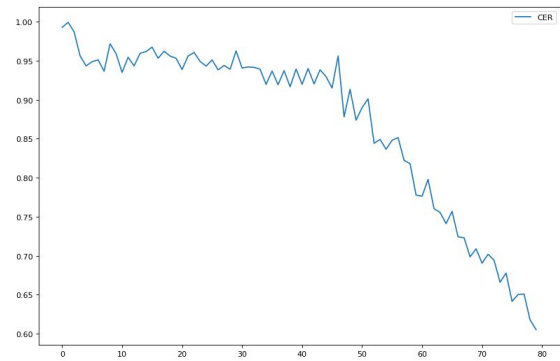


Fig 5: CER Plot for CPC Trained From Scratch

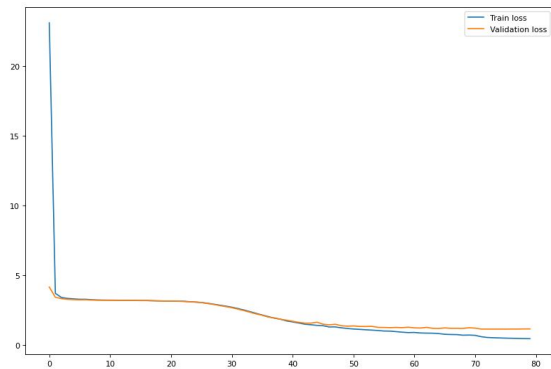


Fig 2: Train and Validation Loss for CPC with Fine Tuning and No Language Model

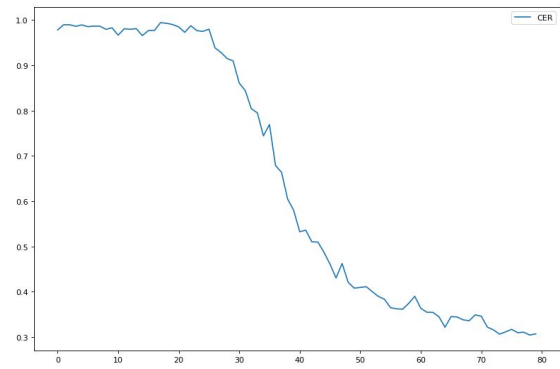


Fig 6: CER Plot for CPC with No Language Model

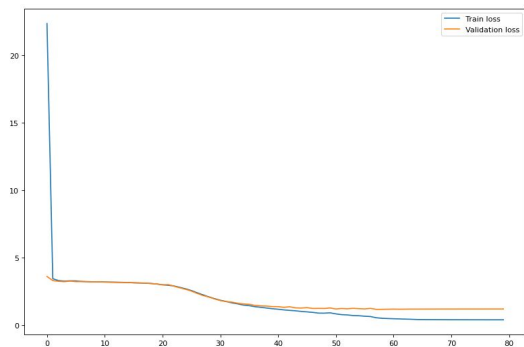


Fig 3: Train and Validation Loss for CER with Finetuning and 4-Gram Language Model

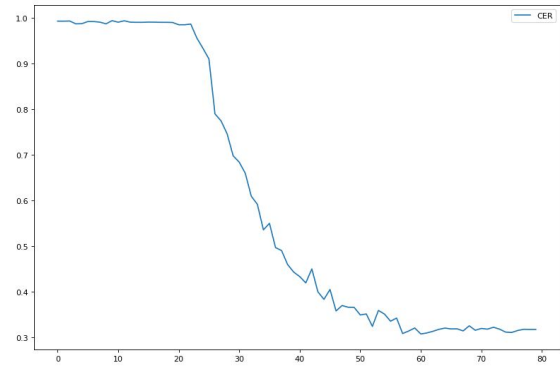


Fig 7: CER Curve for CPC with 4-Gram Language Model

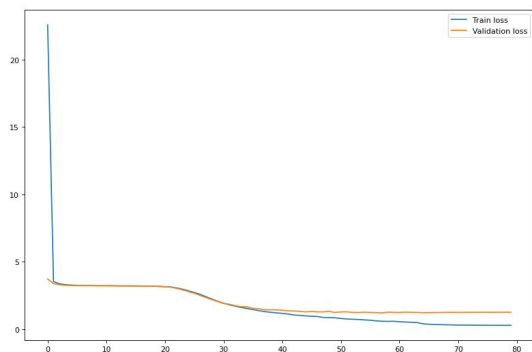


Fig 4: Train and Validation Loss for CER with Finetuning and 5-Gram Language Model

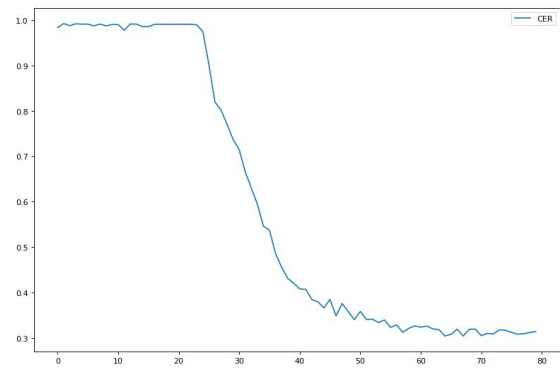


Fig 8: CER Curve for CPC with 5-Gram Language Model

French Experiment Plots

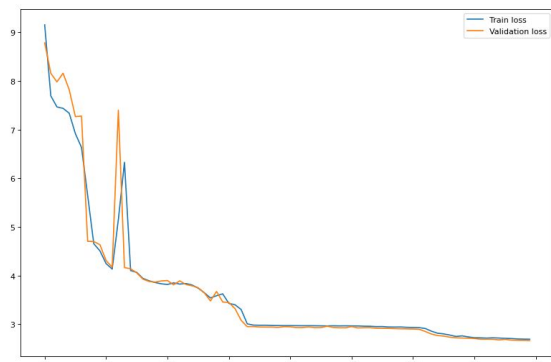


Fig 9: Train and Validation Loss for CPC From Scratch - No Language Model

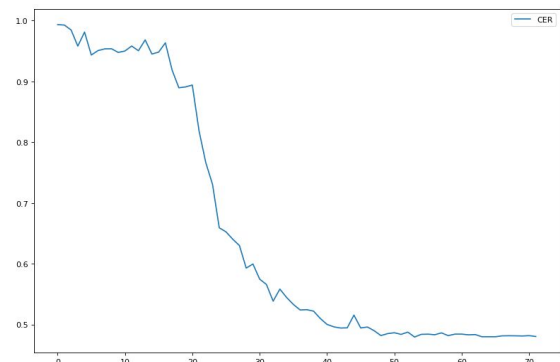


Fig 12: CER Plot - CPC with No Language Model

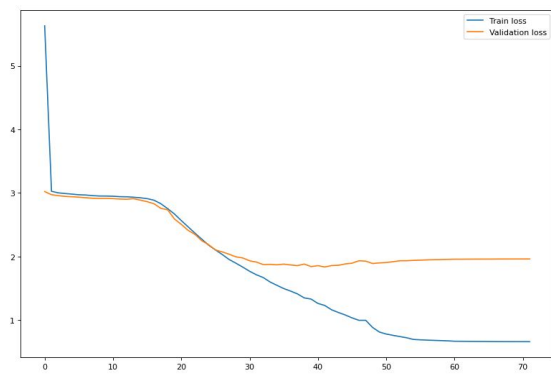


Fig 10: Train and Validation Loss for CPC with Finetuning - No Language Model

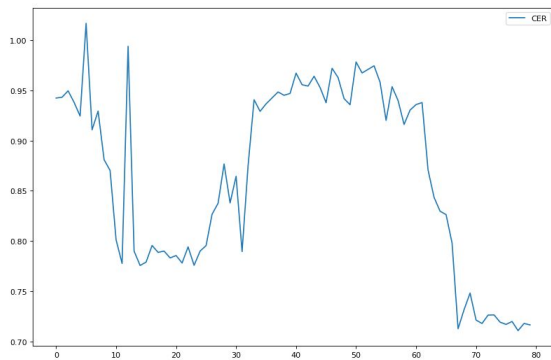


Fig 11: CER Plot - CPC from Scratch with No Language Model