# Progress Update

**Yu-Hsiang Lin**
Language Technologies Institute
Carnegie Mellon University
yuhsianl@andrew.cmu.edu

## Abstract

Place holder [**?** ].

## 1 TTR using Language-independent tokenization

Tokenizer: SentencePiece (https://github.com/google/sentencepiece)

Data: Spanish–English (/home/gneubig/exp/transfer-exp/data/spa_eng)

For debugging purposes, we first consider English, and use only the training set (ted-train.mtok.eng) to train the tokenizer model. The number of sentences in this training set is 196,036.

Model: The model type is unigram, vocabulary size is 8,000, character coverage is 100%.

If we run tokenizer on the same training dataset using this model, the TTR result is:

distinct token count = 7,989

token count = 5,386,583

TTR = 0.001483

Note that a few sample high-/low-frequency tokens are shown in Figures 1 and 2. Actually when it performs the tokenization, the UNK, BOS, EOS have been removed (Figure 3).

Should probably further remove things like "quot", "_&", and punctuations.

If I increase the vocabulary size to 16,000, the result is:

distinct token count = 15,981

token count = 5,232,195

TTR = 0.003054

Give more example of vocabularies in the middle in Figure 4.

(It seems that increasing vocabulary size of model will increase TTR, which indicates that the previous vocabulary size may be too small to be representative.)

vocabulary size = 32,000

distinct token count = 29,765

token count = 5,177,707

TTR = 0.005749

Vocabulary size larger than 32,000 seems to be too large. Encounter error during training: the training ends early and reports error

RuntimeError: Internal: /Users/travis/build/google/sentencepiece/src/trainer_interface.cc(343) [(trainer_spec_.vocab_size()) == (model_proto->pieces_size())]

Figure 1: Top 20.



Figure 2: Tail 20.



Figure 3: Example.

```
Yu-Hsiangs-MBP:ttr yuhsianglin$ head -6000 spa_eng.eng.vocab | tail -100
_Min     -11.895
_critic -11.8952
_Shi     -11.8961
lah      -11.8973
_remotely        -11.8975
_Jas     -11.8981
ough     -11.8981
_shap    -11.8982
_disc    -11.8988
_tablet -11.899
atory    -11.8992
_atmospher       -11.8993
du       -11.9005
_Harr    -11.9023
_controvers      -11.9026
_obsole -11.9027
_Elizabeth       -11.9027
_Michigan        -11.9027
_Public -11.9027
_Sylvi  -11.9027
_accompani       -11.9027
_calendar        -11.9027
_circular        -11.9027
_commodity       -11.9027
_competitor      -11.9027
_jungle -11.9027
_snail  -11.9027
_sulfide         -11.9027
_triumph         -11.9027
_unfair -11.9027
_virgin -11.9027
_cream  -11.9027
_manuscript      -11.9027
```

Figure 4: Middle (6000+).

3