

# WIKIPEDIA SIMPLIFIER (SUMMARY)

## 1. INTRODUCTION:

The simple English wikipedia is an English language version of wikipedia. An automated translation from one readability level into another is generally difficult. Transforming main wikipedia articles into simple English includes a system that can separate language attributes from the content. Natural language does not exhibit an obvious style - content separability. We can discern four known approaches to deal with this issue in the context of our task and the known NLP techniques.

First, a pretraining step can be employed to reduce transformation task to translation. Secondly, articles in simple English can be seen as a special case of Content Summarization. It does not rely on compression of knowledge. Next, delete - retrieve - generate (DRG) technique. DRG workflow intends to identify key phrases responsive to changes and perform interventions on them. In simplification context, this might help with abatement of uncommon words and phrases. Finally, a growing number of publications take the holistic approach to language transformation and rely on pretrained generators to produce samples in target style

## 2. APPROACH

The main idea is to exploit the fact that a text generation model can produce samples aligned with the seed. This seed 'activates' associations acquired in training, and repeated

Sampling can produce a wide variety of responses in Simple English. A pipeline comparing these samples to target content in meaning can reject samples that fail to match the embedding style of the source, effectively manipulating the output by means of a rejection-pass filter (RPF). Latent generative distribution with style  $z$  and content  $c$ :

$$\hat{x} \sim G(z, c) = \pi P(x_t | x_{t-1}, \dots, 0, c)$$

Here we can formulate the problem of style transformation for source lexeme  $s$  as finding such sequence  $x_t$  that the score difference norm  $\|P(s) - P(x_t | x_{t-1}, \dots, 0, c)\|_1 \leq \tau$ . More

Precisely, since a trained generative model already includes style  $z$ , we can control generator output simply by means of repeated sampling admitted on the condition of satisfying a similarity threshold  $\tau$ . This process of threshold-based filtering is recurring and self-adjusting.

### 3. EVALUATION CRITERIA, BASELINE AND

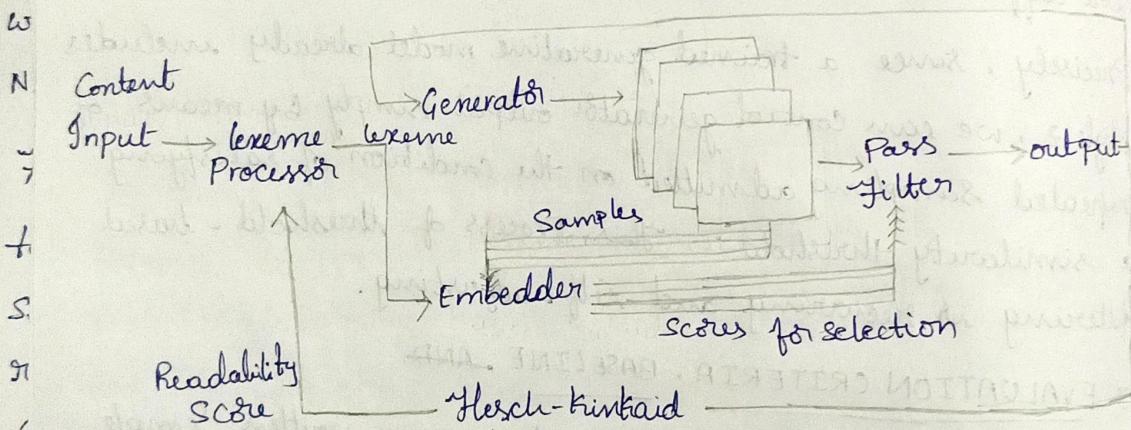
Given our aspiration to supplement human-written Simple English articles by automation, we propose employing the readability score of the former as a baseline. The Average Flesch Score of current simple English Wikipedia is 61.46, which is more than the flesch score of main Wikipedia.

The Flesch-Kincaid formula is not perfect as it only accounts for average length of sentences and words while ignoring the actual frequency distribution over the entire vocabulary. Here lexemes of sentences are subjected to potential replacements from the text generator trained on the simple English corpus. After this, if the content similarity score is

Less than the threshold  $\tau$ , it will be rejected and original lexeme is taken. And in the next step we use the lexeme as a seed to generate samples for replacement of the other.

#### 4 FRAME WORK ARCHITECTURE

We are using a GPT-2 transformer-based text generation model and a Google Sentence Encoder as the building blocks of the architecture. The pipelining layer that takes the source content, parses it into lexeme feeds into the generator, and applies the rejection-Pass Filter is all-original.



The hyperparameters define the work of the generators and lexeme processor, and can be critical to the final output quality.

RPF Hyperparameters :  $\tau$  is a rejection threshold, {MIN, MINSOFT, MAX} are the unconditional minimum, Punctuation-marked minimum, and the maximum lexeme size in words respectively, SEED is the number of lexemes provided as context, NSAMPLES stand for the number of candidates produced by the generator per each lexeme, and  $t^0$  is the generation temperature.

## RPF architecture algorithm:

Input: Content Sources

Output: Simplified text  $St = ""$

repeat

    Identify next lexeme  $L = \text{process}(s)$

    for  $i=1$  to  $n$  samples do

        Generate next sample  $x_i = \text{generate}(s)$

        Save embedding score  $e(cx_i)$

    end for

    if  $\max(e(cx_i) - e(L)) > T$  then

$St = \arg \max(e(cx_i) - e(L)) + St$

    else

$St = L + St$

    end if

until source  $s$  is done

## 5. DATASET AND IMPLEMENTATION DETAILS.

Our lexeme Processor is a heuristics engine implemented in Python. It splits the incoming sentence into lexemes preferring terminal punctuations, and falling back to other punctuation markers if no terminal boundaries are found within a reasonable span. The Configurable hyperparameters provide the compromise between the content preservation and simplification. If we wikipedia simplification works best when the threshold is conservative, and the generator operates at a lowered temperature.

## 6. EXPERIMENTS

The lexeme replacement rates can vary based on the article. The rates are driven by the amount of material available for training (Simple Wikipedia corpus) adjacent to the articles

for transformation. "Personal" Category that describe biographies are seen to have lowest replacement rates whereas the highest replacement rates are seen when there is rich and factual context to draw from. Such as in articles on popular religions or cultural phenomena.

The changes in readability scores are invariably positive, with a qualifier that such a change might be trivial if achieved with means simpler than a neural net. A trivial change can lead to large differences in the Flesch-Kincaid score. A more complicated intervention would involve a restructuring of the whole phrase. Not every replacement results in a coherent modification.

## 7. ERROR ANALYSIS

The formal features of transformed text alone can not adequately describe the transformation because they do not estimate the qualitative language features such as preservation of content and text fluidity.

7.1 phrase encoder errors: This project employs Google sentence Encoder for calculation of candidate acceptance score. This is a state-of-the-art encoder, but sometimes it fails to indicate a mismatch. Although the sentences are very similar according to their scores they can communicate entirely different meaning. To combat such blunders we can raise the acceptance threshold  $\tau$ , but this may lower the effective replacement rates.

7.2 Suboptimal treatment of name entities and rare words:  
The framework with names and words that were not in the training set which promote RPF to produce credible-looking replacements with some "hard to spot" errors. In all these cases the network does a good effort to comply with a seemingly impossible requirement to produce words not seen in training.

#### 7.3 oversimplification:

It seems like some of the simple English constructs picked from the training set appear to be too simple. In these cases the RPF may pick unwanted language biases from training set.

#### 7.4 Non-Verbatim quotes:

Another interesting artifact is seen when rendering quotes. The RPF may alter the quote in some cases. The preservation of verbatim quotes is not learned in training, so this artifact highlights one gap in our framework.

### 8. FUTURE WORK:

There are several limitations to the proposed framework. By doing a second pass over the generated text going in the opposite direction can alleviate low replacement rates. Better training method can be used when dealing with Name Entities and rare words. A brute-force solution to a large number of trials for finding a match which ~~slow~~ slows the generation process would require faster GPU, but a more intelligent solution might require stopping sampling early if there are enough high-quality candidates.