

POETRY THROUGH PROPOGATION, GENERATING HAIKUS WITH DEEP LEARNING RECURRENT NEURAL NETWORKS

Evan Banerjee

Student# 1009682309

evan.banerjee@mail.utoronto.ca

Diego Ciudad Real Escalante

Student# 1009345308

diego.ciudadrealescalante@mail.utoronto.ca

Noah Monti

Student# 1009452398

noah.monti@mail.utoronto.ca

Ji Hong Sayo

Student# 1007314728

ji.sayo@mail.utoronto.ca

ABSTRACT

This is our progress report for the APS360 Final Project. Total Pages: 9

1 INTRODUCTION

2 ILLUSTRATION

3 BACKGROUND

4 PROJECT DESCRIPTION

Project Description “Genuine poetry can communicate before it is understood” —T.S. Eliot In this project we are working to build a deep learning model to generate short poems of various topics. Our model is a haiku generator which can be prompted with the first line of a haiku and generate the content of subsequent lines, built on an LSTM (long-short-term-memory) architecture. Poetry has long been regarded as a deep expression of emotion and human experience. Given that the model we construct will possess neither emotions nor intuition, we are interested to see whether our construction will be able to imitate human poetry, and if so to what degree. If the model is capable of producing compelling poems as output, this would contradict the idea that poetry requires emotion to produce, with implications for our understanding of the nature of creative work in general.

Several ML approaches can be taken to text-generation. For instance, we compare our model to markov chains as a baseline. However, deep learning offers advantages that other statistical approaches cannot. As compared to simpler statistical models, neural nets can learn more abstract and general characteristics in the poems used for training. Given that a hallmark of good poetry is abstraction and unity around a sustaining theme, we believe deep learning is uniquely well-suited for the problem of poetry generation.

5 INDIVIDUAL CONTRIBUTIONS AND RESPONSIBILITIES

Thus far our team has been performing effectively—we have been able to meet both internal and external deadlines. Currently, we’ve been collaborating over github, with most team members working in jupyter notebooks we upload to the shared github page for review. Team communication takes place over whatsapp, and both systems have been working well thus far.

On a technical level, we are quite pleased with our progress thus far, which has been facilitated by the fact that training on text is less memory-intensive than training on images, which allows

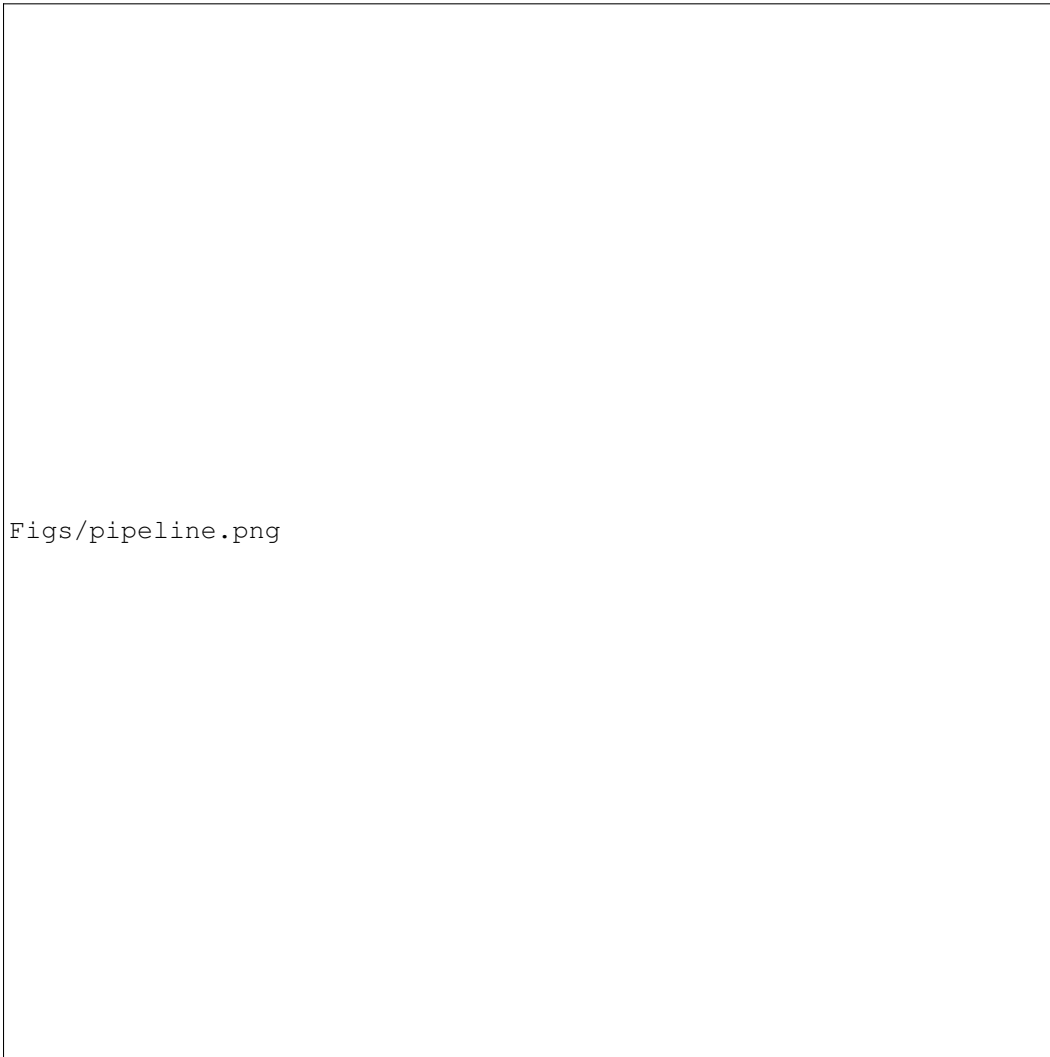


Figure 1: Basic Model Pipeline

for more rapid prototyping without consuming very large amounts of computing power. We have already implemented both 'baseline' and prototype models. What remains is hyperparameter tuning, improvements to the model and testing, all of which we feel will be completed well within the timeframe of the project. In summary, we are on track to meet the goals outlined in our proposal, and may even implement some additional testing and functionality beyond what was outlined in our initial proposal.

5.1 REDUNDANCIES

The most crucial remaining aspect of our project is training the finalized model with a consolidated dataset. This is also a potentially sensitive step, as connection losses or malfunctions could result in significant investments of time and computational resources being wasted. Given this, we plan to have two team members responsible for training. If the primary team member runs into difficulties, they will quickly alert the team via WhatsApp, so that the secondary group member can take over and begin training on a separate system, thus reducing the risk of a technical or personal problem derailing the project.

Several of our group members have separately implemented or are working to implement different model architectures. While we have a consensus on which architecture we plan to use for the

final model, this experimentation provides a strong element of redundancy. If unexpected issues emerge with the architecture we've agreed on, we can discuss as a group and fall back on one of the experimental models.

5.2 TEAM CONTRIBUTIONS AND COMMUNICATION

As shown in the task tracker (Figure 2) and the communications table (Figure 3), all team members have been responsive and met the necessary deadlines for their work.

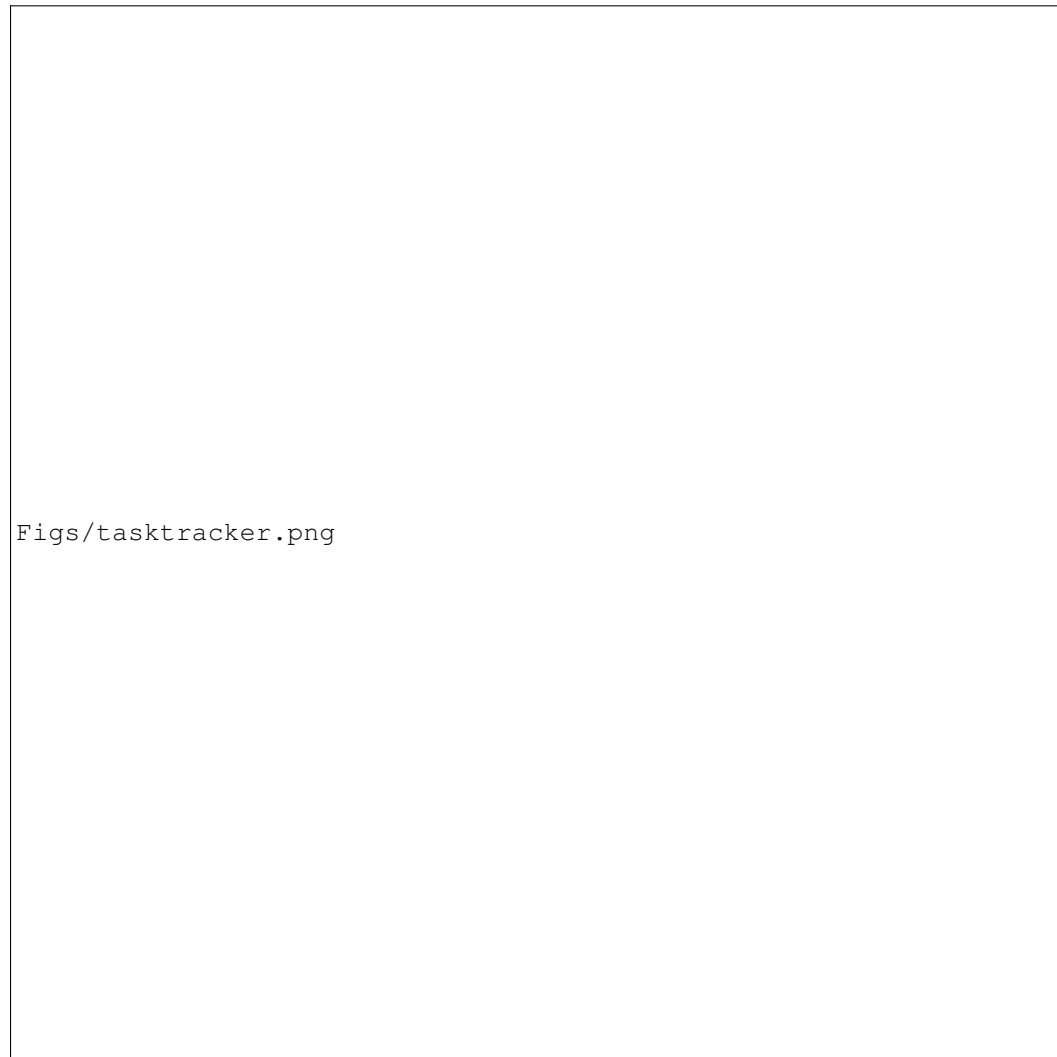


Figure 2: Updated Task Tracker

6 NOTABLE CONTRIBUTION

6.1 DATA PROCESSING

The main data used for this project consists of the haiku4u dataset in kaggle (?). This dataset is publicly available and comes from a web scraper run in October 2023. The data comes in the form of a csv file which has the following fields: Haiku, a field containing the entire haiku with lines separated by a "|" character; Domain, the name of the web page where it was originally found; Domains, the amount of times this specific haiku was encountered while scraping the web; URL,

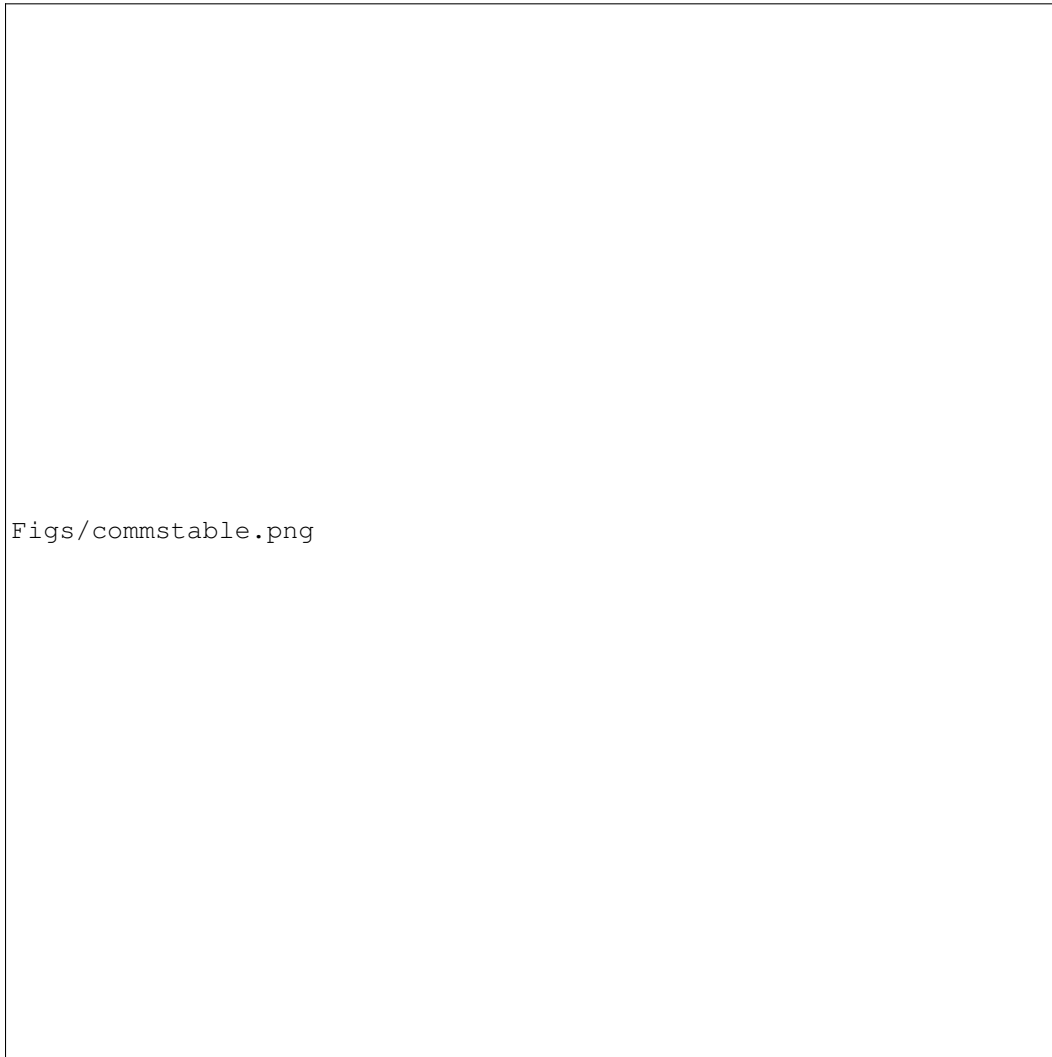


Figure 3: Team communication and contributions

the specific url where the poem was found; URLs, the number of urls that were found in in the same domain. Siblings, the number of other Poems that were found in that same URL. Another data source we used is the CMU pronouncing dictionary (?). We use this dictionary to count the amount of syllables in each line of the haiku. Given that this is an open source project with a python library, data cleanup was minimal for this dataset.

After downloading the data, the cleaning process for the haiku dataset consists is run by a python script. This python script takes in each row of the csv file, and removes all the rows that do not contain a haiku. Then, it splits the haiku into its three lines based on where the | delimiter can be found. Next it converts all the characters to lowercase. And finally, it appends the haiku to a text file in which the end of a haiku is denoted by two new line characters. The text file is then used at training when loading the data. Below is an example of the raw input data, and its processed form.

Original (Raw text):

Yesterday it worked.|Today it is not
working.|Windows is like that.,richardneill.org,27,
<https://richardneill.org/humour.php>,27,19

Finalized (After cleaning):

yesterday it worked.
today it is not working.
windows is like that.

Another use of the data is in building the vocabulary before training. To do this we run the entirety of the text file word-by-word into a python dictionary that maps each word to an index. This lets us represent the words encountered in a numeric form that is easier to deal with. Finally, we must say that because the nature of this project is generative, we testing new data simply means prompting the model differently. However, we are actively looking into ways of increasing the training data we have access to. Some examples of how we could do this are concatenating this data set with others like it and scraping the web ourselves in the hopes we find many more haikus that were not seen originally in 2023.

6.2 BASELINE MODEL

Our first main baseline model is a markov chain of depth four. This was trained on the same data used in the primary model. The main structure is shown in figure 0:

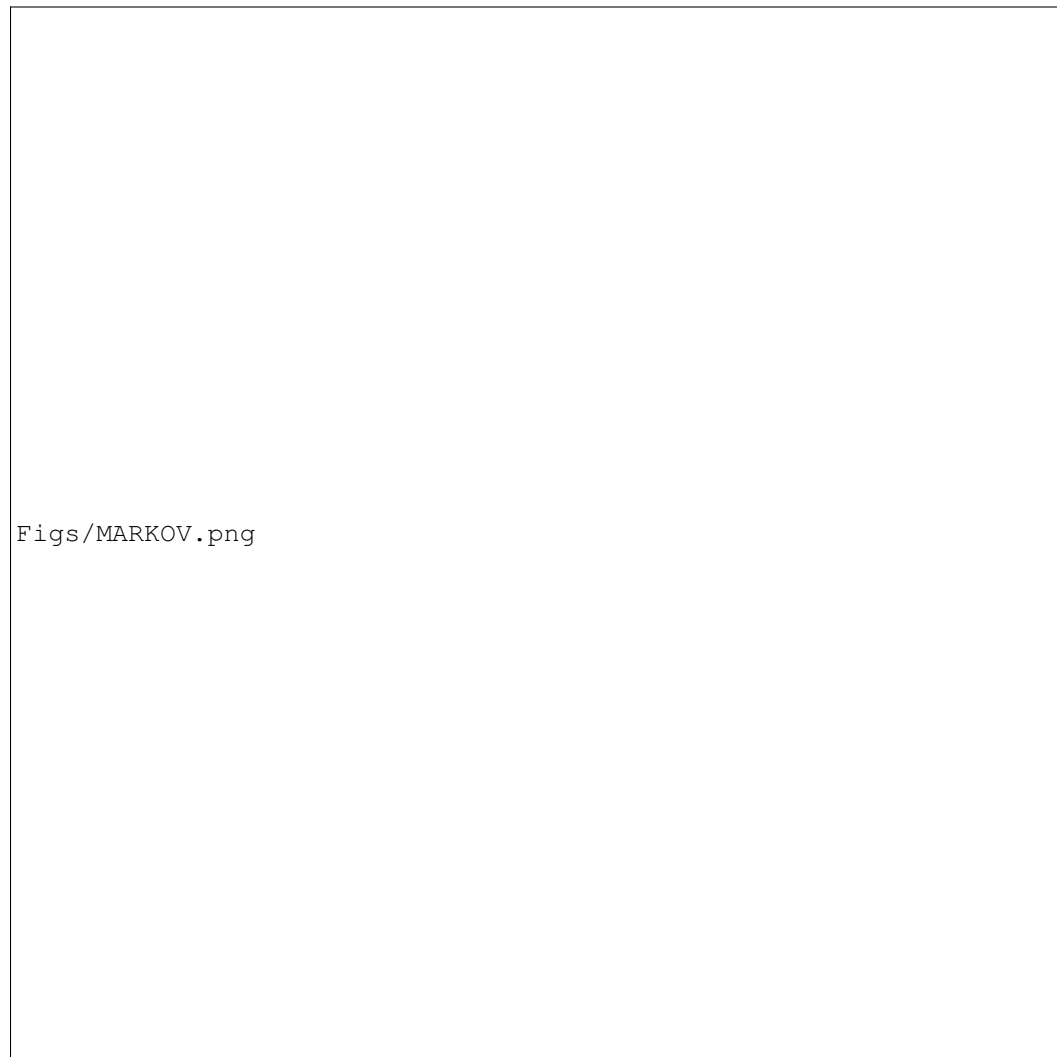


Figure 4: general weighting scheme for markov chain

Each individual word in the corpora is assigned a weighting scheme for every word (including itself) based on how frequently those other words come after the first word.
Then every pair of words is assigned a weighting scheme for every word in the corpora based on how likely the word appears after the pair.
Then every triple of words is assigned a similar weighting scheme.

This repeats until you have a weighting scheme for every n-gram sequence of words up to the depth provided by the function.

Throughout this process we treated newline characters and "end of poem" lines (EOP) as distinct words. The model would take in a starting string of words and would try to autocomplete the poem. It chose words by taking a weighted random choice based on the weights calculated for the markov datastructure

We chose to use a depth of 4 because lower depths had far more incoherent outputs while higher depths seemed to either return no output or just directly copy poems for the most part. Here are some example inputs, and their corresponding outputs:

Show me:

*show me a garden that's
bursting into life*

Give:

*give your
puppies a little extra
thankful that god took me out
this funny as hell
steve's hair evolving*

Teach:

*teach people
how to get folder icons
to show in snowboard*

It should be noted that these poems were hand-chosen from the markov chain for being more coherent. Often, the markov chain will simply copy off another poem, return no output, or only return a couple of words, such as the input "teach" generating the output "teach great place to live".

Even in the best case, this model struggles to correctly write lines with the correct number of syllables.

We created another baseline model from a very simple LSTM RNN. the structure for this one was a much simpler version of our primary model.

It had a 128 dimensional embedding scheme with 2 hidden layers of dimension 256 and a learning rate of 0.001. It was trained with a batch size of 64 for 30 epochs. The output was mostly incoherent so we decided to use the Markov chain as a baseline. All the source code for the first RNN and the Markov chain is on github, under the jupyter notebook `RNN_Model_1_1`.

6.3 PRIMARY MODEL

The overall model architecture is described in the Figure 5 flowchart.

For our haiku generation model, we continued with a Long Short-Term Memory (LSTM) architecture to capture the contextual nuances, and move towards generating coherent and accurate haikus.

The architecture of the model is as follows:

First, the input data is passed through an embedding layer. It transforms this data into vector representations that capture the relationships between words throughout the data. The layer takes in the number of unique words in the dataset, including our special tokens. These special tokens consist of the following:



Figure 5: Model Architecture

- Padding token - used to make sequences uniform in length within a batch.
- Unknown token - represents a word that is not in the vocabulary.
- End-of-Sequence token - signifies where the model should stop generating text.

The embedding dimensions are also inputs to this layer, which are the size of each word in the embedding vector, as is the padding index - the index reserved for our padding tokens that ensure padded positions don't contribute to the learning process.

Next, the output of the embedding layer is passed through the LSTM layer. This layer serves to process the sequence of embeddings to capture contextual information across words. It takes in the numbers of features in the hidden states of the LSTM, the number of stacked LSTM layers, and the batch size. The resulting output of this layer is a tensor that has the output features from the LSTM, as well as the hidden states of the LSTM for each layer.

Finally, the output of the LSTM layer is passed through a fully-connected layer. This maps the LSTM outputs to our vocabulary set, and produces logits for each word in this vocabulary. The input features are equal to the LSTM's hidden dimensions, and the output features are equal to the size of the vocabulary.

In our current configuration, we used 128 embedding dimensions, 256 hidden dimensions, and 2 LSTM layers.

With our current training data, we have a vocabulary size of 50554 words.

The embedding layer has $50554 \cdot 128 = 6470912$ parameters.

Next the the LSTM layer has $2 \cdot 4 \cdot (256 \cdot (128 + 256) + 256) = 788480$ parameters. The four is to account for the input gate, forget gate, cell gate, and output gate in an LSTM.

Finally, the fully-connected layer has $256 \cdot 50554 + 50554 = 12992378$ parameters.

Therefore, the total number of parameters in the current model is 20251770.

As for our training hyperparameters, we chose to train over 50 epochs, with a learning rate of 0.001, a batch size of 64, an Adam optimizer¹, and a Cross Entropy Loss criterion².

Quantitative Results:

We tested out model with 300 one or two word prompts to the model, and counted the syllables in the output haiku.

As we can see in Figure 6, the vast majority of generated haikus are between 17 and 19 syllables, with a few outliers. These can be explained by either the model incorrectly counting syllables from words that are not in the CMU Pronouncing Dictionary (?), or the syllable counting algorithm miscalculating the output. We are aiming to have our model only output 17 syllable haikus where possible - there are more methods we will be integrating into the model.

Qualitative Results:

The model generates a mix of interesting or funny haikus, and incoherent outputs:

- *“river flows wildly
air patches on the surface
dissipates in life”*
- *“rainfall through the roof
dancers from serenity
as his life goes by”*
- *“grass on the wayside
looking for the light to come
warm clouds and weather”*
- *“bloom gives way to me
for i am not your hero
i am george carlin”*
- *“pond ripples to steal
winters for having their gaze
beg for hope bears now”*
- *“branch and wood filters
of night bright and sunny day
hot lemonade waits”*
- *“stream sits and pure fish
in between and i am glad
with what we i”*
- *“ancient tree freedom
anxiety creation
ignites and section”*

As we can see, the model is able to generate haikus that could pass for human-written, however, it still generates incoherent haikus. Moving towards generating more coherent outputs will be a focus of ours as we improve to model.

¹Chosen from findings in (?)

²Chosen from approach in (?)

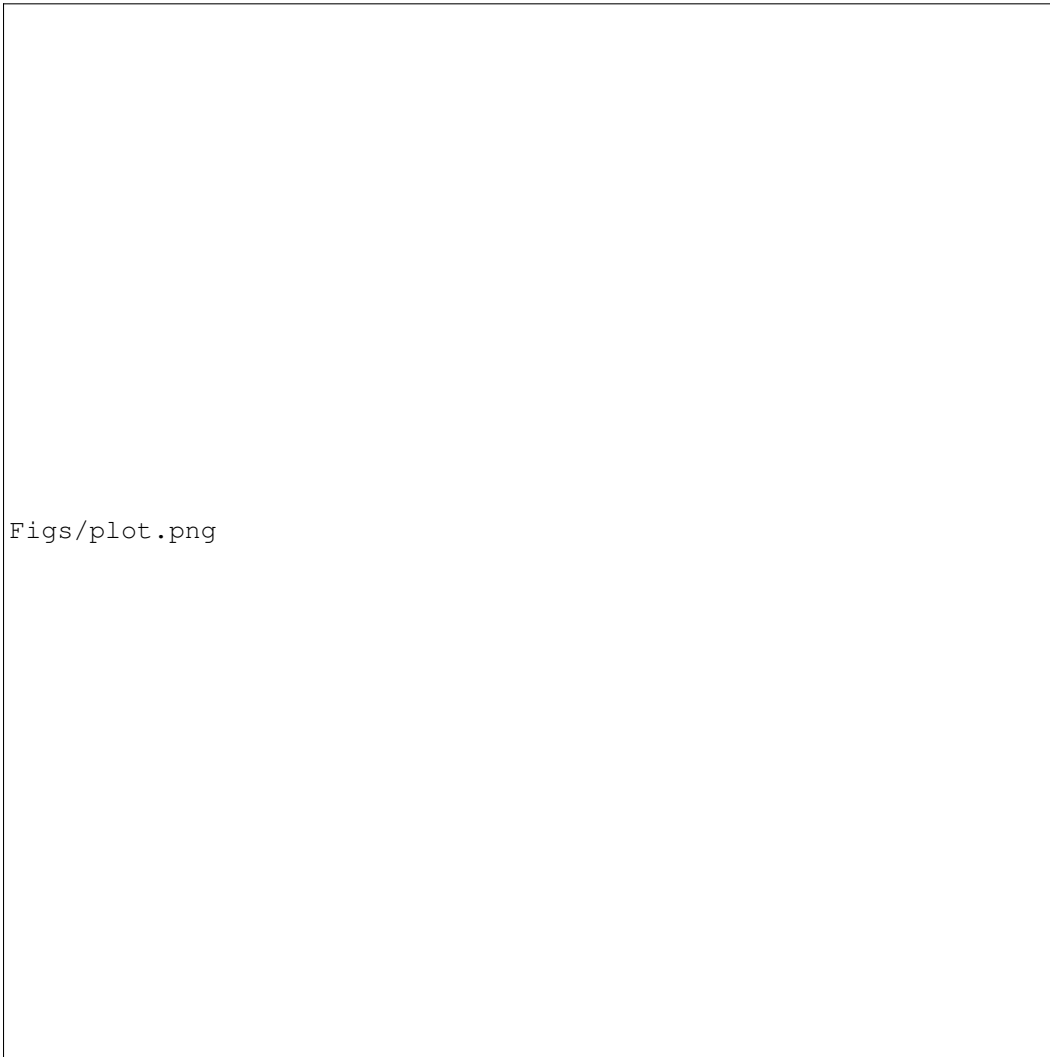


Figure 6: Syllables per haiku from 300 outputs

The main challenges we faced making this model were related to ensuring syllable structure and was maintained and the syllable counts were accurate, and slow training times.

Because of irregularities in the English Language, we found it difficult to come up with a reliable method to count the number of syllables in a given word. Our current solution involves referencing the CMU Pronouncing Dictionary (?) to get the syllable count of a given word. However, this method is not perfect, as the CMU Pronouncing Dictionary does not contain all words that may show up in our vocabulary. This causes issues when a word is not in the dictionary, so we are looking to integrate algorithms that can help resolve this, although this still only provides an approximation which can lead to slight issues during generation.

The slow training times were due to the large amount of data and limited access to compute resources. We tested differing batch sizes to reach a reasonable training time given our resources, however we are looking into ways to gain access to more powerful GPUs to speed up training.

REFERENCES

Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. MIT Press, 2007.

Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.

Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.