

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 final report for the APS360 Final Project. Total Pages: 12

1 INTRODUCTION

“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 a gated recurrent unit (GRU) 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.



Figure 1: Basic Model Pipeline

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.

2 ILLUSTRATION

The high-level description of the model can be seen below in 2 . The model consists of Gated Recurrent Unit that is trained on poems sourced from the web. These poems also build up the dictionary of words used by the model to generate haikus. Finally, the diagram also shows the expected usage of this model. A user would prompt the model with the starting words of a poem, press enter, and the model would generate a haiku with the prompt as its starting point.

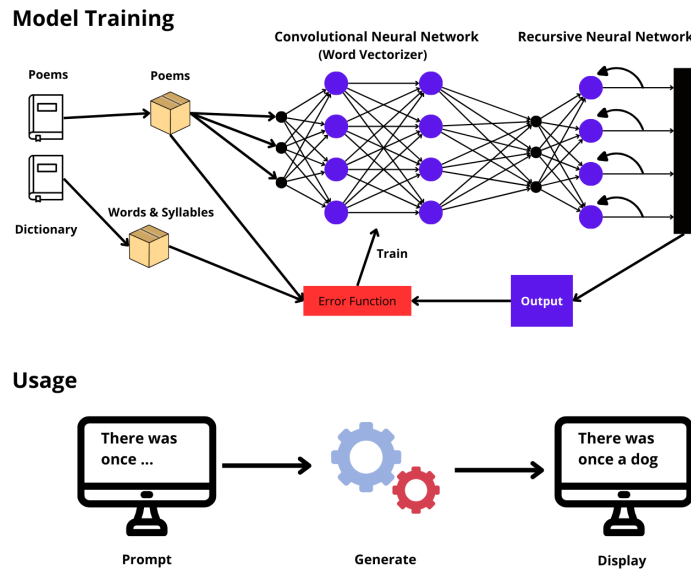


Figure 2: Top-Level Model Description

3 BACKGROUND

Previous poetry generation models utilized Recurrent Neural Networks (RNN) with some form of memory limiting to generate poems off a large dataset.

One model, created by researchers at the University of Saarland, utilized encoder decoder systems, as well as long short-term memory (LSTM) and Gated Recurrent Units (GRU) to address vanishing gradients, (Novikova et al., 2022).

Another model, developed at the University of Edinburgh utilized multiple neural networks chained together, with a convolutional neural network (CNN) for encoding, one RNN for creating hidden context parameters, and another RNN with a separate weight vector for modifying the final output (Zhang & Lapata, 2014). The output of this model was particularly compelling, with human ratings for the generated poems being over 75% of the ratings given to human written poems.

Another model, built at City, University of London, used syllables as the inputs and outputs of their RNN instead of words (Lewis et al., 2021). This model's poems were incredibly convincing, with 56% of polled individuals believing the AI generated poems were written by humans.

When focusing specifically on haikus, even a relatively basic model built by students at Stanford University was able to score well (Aguilar & Liao, 2017). For this model human evaluators gave the AI generated haikus an average quality score that was over 80% of the average score given to human poems. This model only used a single RNN with LSTM. Granted, the evaluations were only done by 8 individuals, but the strength of RNNs for this kind of task is evident.

Finally, a model presented at the 13th conference on language resources and evaluation found a large

degree of success in constraining the output of neural networks to generate poems (Popescu-Belis et al., 2022). This model forced the AI to repeatedly output lines of a poem until certain syntax requirements were met. It also enabled the user to specify emotional tone by encoding category preferences (happy, sad, etc.) into a weighted vector that was then used after generating a line to swap out words that didn't fit the categories well with words that did.

This was most effectively done when words were chosen based on word frequency from other poems that fit the specified category. This model also utilized a softmax function that truncates the least likely words from the sample space during generation, increasing the probability of likely words and drastically improving output.

Each one of these model's utilized different techniques. Whether the differences were in the way they encoded and decoded the inputs and outputs, or the or the specific method they used to manage hidden states. Whether they used syllables as the their input output token, or whether they used words. Whether they forced the model to regenerate lines that didn't fit certain syntax requirements, used algorithms to swap out words after generating the lines, or simply left the lines as is. Each model had their own unique way of addressing the somewhat existential task of creating poetry. Even the specific methods used to normalize their logits or sample from their word space was varied. most of the models weren't even trying to write the same kind of poem.

However, every model utilized some form of encoding and decoding the text, some RNN architecture with a hidden state to efficiently calculate gradients, some variant of softmax to normalize their output, and some probabilistic sampling to choose the next word. Many of the poems also implemented some algorithm to prevent individual lines of each poem from being incoherent. These principles are what guided our design as we developed our model throughout the semester.

In the end we mostly ended up utilizing the work out of the University of Saarland, and the 13th conference on language resources and evaluation. Our design was largely based on the frameworks used in their models, specifically in terms of utilizing a GRU, usinig individual words as outputs, and processing lines after generating them to fit syntax requirements. However, just like every other model, we used the main encoding-RNN with hidden state-decoding framework to generate individual lines.

4 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 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.

4.1 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.

4.2 TEAM CONTRIBUTIONS AND COMMUNICATION

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

5 NOTABLE CONTRIBUTION

5.1 DATA PROCESSING

The main data used for this project consists of the Gutenberg Poetry Corpus dataset in github (biglam, 2022). This dataset is publicly available and comes from a web scraper run in 2018. The web scraper collected poems from the Gutenberg Project, and online repository for public domain literature. The data comes in the form of a JSON file which has the following fields: "s", a line of poetry, and "gid" an id of the book in the gutenberg project database. To clean the data we stripped every line of the poem, and segmented them into "poems" of 10 lines. We decided to split the data every 10 lines since the main reason for using this data set is to draw from more historical and poetic language. Therefore, the actual length of the text we trained our model on is not as relevant since we work out the structure of the haiku through the training.

Another data source we used is the CMU pronouncing dictionary (University, 2014). 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 line of the JSON file, and removes uncommon characters such as the return-carriage character (`\r`), tabs (`\t`), and em-dashes (`—`). Next it converts all the characters to lowercase. And finally, it appends the haiku to a text file in which the end of a poem section 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):

"gid": "442" }

```
{ "s": "That was our bench the time you said to  
me", "gid": "442" }  
{ "s": "The long new poem—but how different  
now,", "gid": "442" }  
{ "s": "How eerie with the curtain of the fog",  
"gid": "442" }  
{ "s": "Making it strange to all the friendly  
trees!", "gid": "442" }  
{ "s": "There is no wind, and yet great curving  
scrolls", "gid": "442" }  
{ "s": "Carve themselves, ever changing, in the  
mist.", "gid": "442" }  
{ "s": "Walk on a little, let me stand here watch-  
ing", "gid": "442" }  
{ "s": "To see you, too, grown strange to me and  
far. . .", "gid": "442" }  
{ "s": "I used to wonder how the park would  
be", "gid": "442" }  
{ "s": "If one night we could have it all alone—",
```

Finalized (After cleaning):

That was our bench the time you said to me
 The long new poem but how different now
 How eerie with the curtain of the fog
 Making it strange to all the friendly trees

There is no wind, and yet great curving scrolls
 Carve themselves, ever changing, in the mist
 Walk on a little, let me stand here watching
 To see you, too, grown strange to me and far
 I used to wonder how the park would be
 If one night we could have it all alone

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, testing new data simply means prompting the model differently. In addition to being the largest collection of poems we could find, this data set provides the additional benefit of having a wide range of styles. This allows us to make more diverse poems as seen below in the qualitative results section.

6 ARCHITECTURE

For our final model, we switched to a neural network centred around the gated recurrent unit (GRU) architecture, which replaced our previous long-short term memory (LSTM) architecture. The gated recurrent unit has been shown to be similar to the LSTM architecture on performance and quality of output, while requiring less compute resources during training (LST, 2020), which appeared to hold true when testing our model and comparing to our previous design.

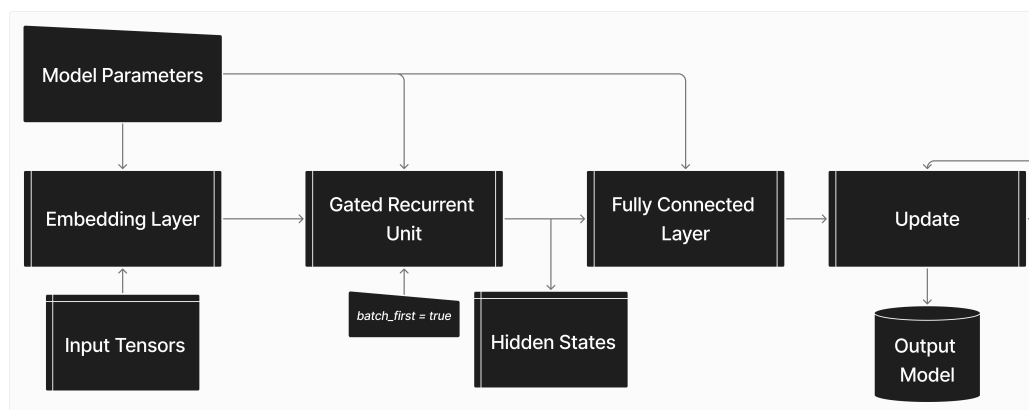


Figure 3: Model Architecture

The architecture of the model is as follows, as described in the Figure 3 flowchart.

Firstly, the model parameters are manually configured. These include the following, with the hyper-parameters we found to work best with our dataset:

1. The size of the vocabulary dictionary - which was 109976 distinct words for our data
2. Embedding dimensions - 128 was found to work well
3. The index of the padding token within the vocabulary dictionary
4. The number of hidden dimensions - set to 256
5. The number of GRU layers - set to 3
6. How much dropout to apply between layers - we did not use any
7. If the GRU will perform bidirectional passes - we did not perform these passes

When data is passed to the model, it first goes through the embedding layer. This layer serves to transform this data into vector representations that capture the relationships between words through-

out the data. The layer takes in the number of unique words in the dataset, the size of the embedding dimension, and the padding token index so that it does not contribute to gradients when training.

The output of that layer is then passed to the Gated Recurrent Unit layer, which serves to process the sequence of embeddings to capture contextual information across words. It takes in the embedding dimensions, the hidden dimensions, the number of GRU layers to have, batch_first (which should always be set to true for our data), how much dropout to apply, and if the GRU will perform bidirectional passes. The resulting output of this layer are two tensors, one with the output features from the GRU, and another that has the hidden states of the GRU for each layer.

Finally, the output features are passed to a linear layer, which maps them as logits to each word in the vocabulary set. The input size to the layer is the number of input features, which is the size of the hidden dimension, and the output size is the number of output features, equal to the size of the vocabulary set.

It should be noted that if the GRU is to perform bidirectional passes, the number of hidden dimensions must be multiplied by two in order to align with the required tensor size.

This model can then be passed to train, for which we used Cross Entropy Loss criterion with the padding index passed as a parameter to ensure it does not contribute to gradients, and the Adaptive Movement Estimation (Adam) optimizer. Additionally, we used a learning rate of 0.001, a batch size of 64, and 50 epochs.

7 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:

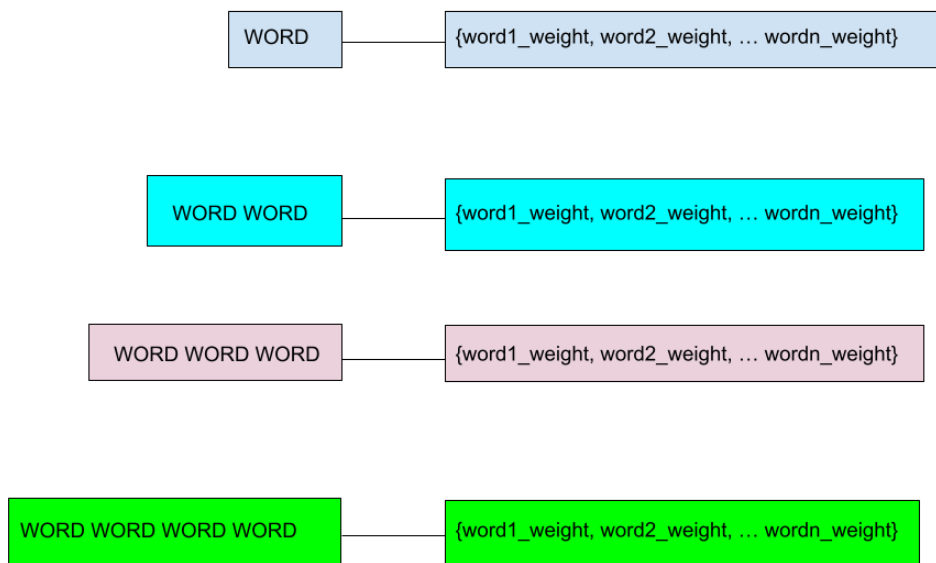


Figure 4: general weighting scheme for markov chain

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`.

8 QUANTITATIVE RESULTS

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

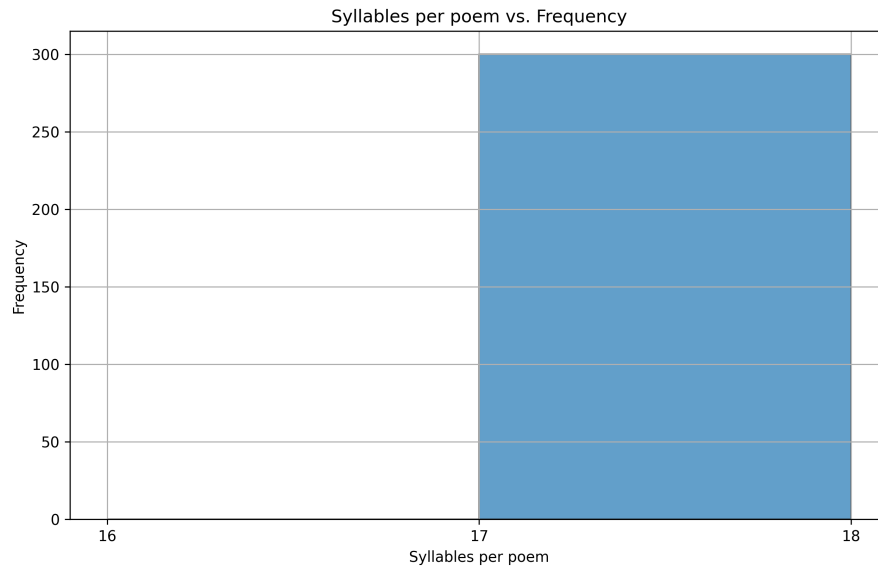


Figure 5: Syllables per haiku from 300 outputs

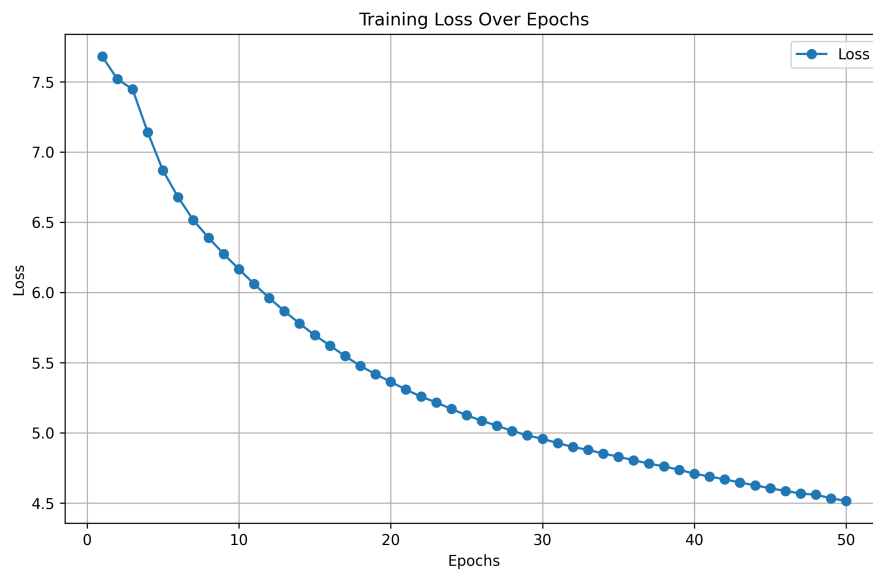


Figure 6: Training Loss over 50 Epochs

As shown in Figure 7, all haikus that were generated by the model meet the 17 syllable count, demonstrating the model's ability to regenerate words dynamically in order to maintain the traditional haiku structure - which was one of the targets from our previous model iteration.

During training, the loss and the rate of loss steadily decreases over the epochs as seen in Figure 6, showing our model is continuing to learn the meanings of words relative to their neighbours in sentences. The loss is calculated as the rate of accuracy between the model's prediction of the next word in a sequence and the original word in that sequence, which reached 4.5 after 50 epochs.

9 QUALITATIVE RESULTS

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

- *"river as thou fledst
she glows within thy dreams where
babe sad favorite"*
- *"water from the breeze
there love dejected lies not
here and could he dream"*
- *"grass and the moon is
the smoke of a broken tree
is great for the sun"*
- *"ocean says she has
to stir stainless nightingale
but well they judge me"*
- *"car wolves reach of sighs
and coat an hurrah like strife
truest mother round"*
- *"long does fairer by
their heart yet cracks and half so
late they bent him off"*
- *"when the death is no
knight ere for my dower and the
vine and hints of that"*
- *"wake when spite of love
when that their light are red and
many a star sweet"*

As we can see, the model is able to generate haikus that could pass for human-written, however, it still generates incoherent haikus. While the decision to make train the haikus in poetry rather than haikus resulted in more poetic language, it also increased the frequency of rarely used words.

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 (University, 2014) 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 opted for heuristic algorithms to fix this, although this is still not perfect in practice, it significantly reduced the number of times the 5-7-5 meter was broken.

Despite gaining more resources for training through Vertex AI, this is not a sustainable way of training the model as it quickly runs out of free compute credits. However, using even just the free Tier drastically increased the speed at which we were able to train the model.

10 EVALUATION ON NEW DATA

To evaluate our model, we generated dozens of new poems from varying inputs and chose the five that we believed were the most coherent to represent the model. We then created a survey with 15 poems. These consisted of the 5 made by our model, 5 chosen from the more coherent outputs of the baseline markov chain, and 5 chosen from the repository of human written poems that we used in training. We had 24 people rate the quality of each poem from 1 to 10, and averaged all the responses. The results are shown in the figure below:



Figure 7: Syllables per haiku from 300 outputs

the five poems used in the survey from our model are:

- *"winter chill with hopes
of peace they see the home to
sleep and we are heard"*
- *"new grass like diamonds
blowing through the darker sounds
of heaven to flame"*
- *"the mightiest makes
the memory of some three
shepherds from the glow"*
- *"the charm of his heart
the graceful hills of the brow
and the black essence"*
- *the ravenous sand
the graveyard for the raindrops
and the gothic feet"*

11 ETHICAL CONSIDERATIONS

Despite AI being around for several decades, it was only recently that computational power necessary to generate human-level works of art became widely available. As impressive as they are, increasingly improving models that generate images, written stories, and poems also raise questions about how these models get and utilize their training data. While this is a very trending topic, legislation protecting individuals from having their written works used to train AI models is still in early

stages. Because of this we have decided to only use works in the public domain. The works we use in this model can all be found in the Gutenberg Project digital library meaning none of them are held under exclusive rights. Finally, this model is a purely academic experiment, and we do not seek to benefit from the material it generates. Therefore, the literary works we will deal with in this project and our use of them fall under fair dealings as per Canadian law (Branch, 1985) .

We have discussed the steps that we will take to prevent misuse of our training data, now let's dig deeper into how someone with harmful intentions might use our project. The biggest cause of concern comes from poem-writing competitions. Instead of having to naturally come up with poems, someone might take advantage of a deep learning model to easily generate several entries, increasing their chances of winning or simply overloading the contest's submission portal. While many competitions may not explicitly ban AI-generated poems, it is not hard to see how a generated poem winning over many human authors may be seen as controversial. Our system somewhat counteracts this by having the user be vital in the poem creation as they must prompt the model in order for it to generate anything. Additionally, our limited scope of haikus would make it harder for someone to take advantage of the model on a massive scale as haikus are notoriously short. With these ethical considerations, we can confidently say that misuse of this model is highly unlikely.

REFERENCES

- Lstm and gru neural network performance comparison study: Taking yelp review dataset as an example. In *2020 13th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, pp. 655–660. IEEE, 2020. doi: 10.1109/CISP-BMEI51763.2020.9221727. URL <https://ieeexplore.ieee.org/document/9221727>.
- Rui Aguiar and Kevin Liao. Haiku generation. GitHub, 2017. URL https://github.com/raguiar2/haiku_generation/tree/master.
- biglam. gutenberg-poetry-corpus, 2022. URL <https://huggingface.co/datasets/biglam/gutenberg-poetry-corpus>.
- Legislative Services Branch. Consolidated federal laws of canada, copyright act, 1985. URL <https://laws-lois.justice.gc.ca/eng/acts/c-42/page-6.html#docCont>.
- Danielle Lewis, Andrea Zugarini, and Eduardo Alonso. Syllable neural language models for english poem generation. In *12th International Conference on Computational Creativity (ICCC’21)*, pp. 350–356. Association for Computational Creativity, 2021.
- Svetlana Novikova, Sangeet Sagar, Peilu Lin, Meng Li, and Pavle Markovic. English and chinese poetry generation software project: Deep learning for the processing and interpretation of literary texts. 2022.
- Andrei Popescu-Belis, Alex R Atrio, Valentin Minder, Aris Xanthos, Gabriel Luthier, Simon Mattei, and Antonio Rodriguez. Constrained language models for interactive poem generation. In *Proceedings of the 13th Conference on Language Resources and Evaluation (LREC 2022)*. 20-25 June 2022, 2022.
- Carnegie Mellon University. The carnegie mellon university pronouncing dictionary, 2014. URL <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>.
- Xingxing Zhang and Mirella Lapata. Chinese poetry generation with recurrent neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 670–680, 2014.