

# Character-Based Text Generation

Davide Giuseppe Griffon

## Abstract

This document serves as the report for the fourth task in the "Natural Language Processing" course completed by student Davide Giuseppe Griffon at Vilnius University as part of the Master's program in Data Science.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Text generation</b>	<b>2</b>
2.1	Loading the Data . . . . .	2
2.2	One-hot Encoding and Decoding . . . . .	2
2.3	Building the Training Dataset . . . . .	3
2.4	LSTM Architecture . . . . .	3
2.5	Performance Evaluation . . . . .	4
2.6	Challenges and Possible Improvements . . . . .	5
<b>3</b>	<b>Appendix - Code</b>	<b>6</b>

# 1 Introduction

In this project, I developed a character-based text generation system using Long Short-Term Memory (LSTM) networks. The model was trained on a collection of works by G.K. Chesterton from the Gutenberg corpus, with the goal of generating new text that mimics the author’s writing style.

A text corpus is a large, structured collection of texts that serves as a foundational resource in natural language processing (NLP). Text corpora provide the raw material necessary for analyzing language patterns, training models, and evaluating their performance. In NLP applications, corpora are essential because they offer authentic examples of language usage, allowing models to learn naturally occurring patterns rather than relying on manually coded rules.

The remainder of this document is organized as follows. First, I describe the data loading and preprocessing steps, including the creation of character encodings. Then, I detail the architecture of the LSTM model and the training process. Finally, I present the results of text generation experiments and discuss potential improvements to the system.

## 2 Text generation

### 2.1 Loading the Data

For this project, I utilized Chesterton’s works available in the NLTK Gutenberg corpus, which comprises three major books: “The Ball and the Cross”, “The Wisdom of Father Brown”, and “The Man Who Was Thursday”. To facilitate data loading and preprocessing, I implemented a `GutenbergLoader` class with various utility functions.

The corpus analysis revealed the following statistics for each work:

Work	Characters	Words	Lines
The Ball and the Cross	457,450	81,598	9,548
Father Brown	406,629	71,626	7,654
The Man Who Was Thursday	320,525	57,955	6,793
<b>Total</b>	1,184,604	211,179	23,995

Table 1: Statistics of Chesterton’s works in the corpus

The combined corpus consists of 1,184,604 characters in total, with a vocabulary of 91 unique characters. The complete character set includes:

! " \$ % ' ( ) \* + , - . / 0 1 2 3 4 5 6 7 8 9 : ; < > ? @  
A B C D E F G H I J K L M N O P Q R S T U V W X Y Z [ ] \_ ‘  
a b c d e f g h i j k l m n o p q r s t u v w x y z ~ ` é ê ì

### 2.2 One-hot Encoding and Decoding

For this purpose, I wrote the `CharacterEncoder` class which is responsible for encoding and decoding characters using one-hot encoding. The method `fit` creates bidirectional mappings through two simple dictionary comprehensions:

```

1 self.char_to_index = {
2     char: idx for idx, char in enumerate(unique_chars)
3 }
4 self.index_to_char = {
5     idx: char for idx, char in enumerate(unique_chars)
6 }

```

The first line creates a mapping from characters to indices, while the second creates the reverse mapping. Using Python's `enumerate` function, each unique character is assigned a consecutive integer index, making the encoding and decoding process straightforward.

Furthermore, methods `save_mappings` and `load_mappings` are responsible for saving and loading these mappings from a JSON file, making the loading process more efficient during subsequent steps.

## 2.3 Building the Training Dataset

To train the character-based language model, I needed to create overlapping sequences from the input text. For this purpose, I implemented the `TextDataset` class, which inherits from PyTorch's `Dataset` class. Each sequence in the training set consists of 40 characters (`self.seq_length = 40`), and the target is the character that follows this sequence.

The core functionality resides in the `__getitem__` method:

```

1 def __getitem__(self, idx: int) -> tuple[torch.Tensor, torch.Tensor]:
2     # Get the sequence and target
3     sequence = self.text_indices[idx : idx + self.seq_length]
4     target = self.text_indices[idx + self.seq_length]
5
6     # Create one-hot encoding for input sequence
7     X = torch.zeros(
8         self.seq_length, self.char_encoder.vocab_size, dtype=torch.
9         float32
10    )
11    for t, char_idx in enumerate(sequence):
12        X[t, char_idx] = 1
13
14    # Return target as a single integer (not one-hot encoded)
15    return X, torch.tensor(target, dtype=torch.long)

```

This method creates overlapping sequences by sliding a window of 40 characters over the text. For each sequence, it returns:

- A one-hot encoded tensor of the input sequence with shape `(40, vocab_size)`
- The index of the target character (the 41st character) as a single integer

## 2.4 LSTM Architecture

For this project, I implemented a simple LSTM (Long Short-Term Memory) network using PyTorch. The model architecture consists of two main components:

```

1 class LSTM(nn.Module):
2     def __init__(self, input_size: int, hidden_size: int = 128):
3         super(LSTM, self).__init__()
4         self.hidden_size = hidden_size

```

```

5         self.lstm = nn.LSTM(input_size, hidden_size, num_layers=1,
6                               batch_first=True)
7         self.fc = nn.Linear(hidden_size, input_size)
8
9         def forward(self, x: torch.Tensor) -> torch.Tensor:
10             x, _ = self.lstm(x)
11             # Return logits (softmax will be applied in the generation step)
12             return self.fc(x[:, -1, :])

```

The architecture consists of:

- A single LSTM layer with a hidden size of 128 units
- A fully connected layer that maps the LSTM's output back to the vocabulary size

The model takes a sequence of one-hot encoded characters as input and outputs raw logits over the possible next characters. The softmax function is not applied in the forward pass but rather during the text generation phase to convert these logits into probabilities.

## 2.5 Performance Evaluation

The generation of text is performed using a proxy class `TextGenerator` that loads the trained model and generates text using the `generate_text` method. The diversity parameter controls the randomness in the text generation process by scaling the logits before applying the softmax function.

The code for text generation is as follows:

```

1 seed_text = (
2     "There was an instant of rigid silence, and then Syme in his turn "
3     "fell furiously on the other, filled with a flaming curiosity."
4 )[:40]
5 # Note that the truncated text is:
6 # "There was an instant of rigid silence, a"
7 for diversity in [0.2, 0.5, 1.0]:
8     print(f"\nDiversity: {diversity}")
9     generated = generator.generate_text(
10         seed_text=seed_text, length=100, diversity=diversity
11     )
12     print(generated)

```

Here are examples of generated text with different diversity levels:

- **Diversity: 0.2**  
 "There was an instant of rigid silence, and the beard and the black strange of the streets of the street and face and strong and streets that"
- **Diversity: 0.5**  
 "There was an instant of rigid silence, and in the two some light and sunset the chair and sat desportably entirely like a strong face and be"
- **Diversity: 1.0**  
 "There was an instant of rigid silence, and the writity-slack't read and gatening silks, by might began tomprecy that they have then heally e"

The results demonstrate several limitations of the character-level model:

- Despite occasionally generating existing words, the overall text lacks coherence and semantic meaning
- The model shows no understanding of broader context or grammar
- Higher diversity levels (`diversity = 1.0`) lead to the generation of non-existent words and more random text sequences
- Lower diversity levels (`diversity = 0.2`) produce more repetitive patterns but still lack meaningful structure

These limitations are inherent to the character-level approach, as the model operates without any understanding of word-level semantics or linguistic structure. While it can learn character patterns and combinations common in English text, it cannot capture higher-level language features that would be necessary for generating coherent narratives.

## 2.6 Challenges and Possible Improvements

One of the main challenges I encountered during this project was the significant computational time required for the model training phase. This slowness stems from the inherent complexity of processing large volumes of character-level data combined with the computational demands of the model architecture. While I utilized the MPS (Metal Performance Shaders) device available on my Apple machine, which provided better performance than CPU processing, the lack of GPU (CUDA from NVIDIA) access remained a limiting factor in achieving optimal training speeds. I attempted to address these performance issues by experimenting with a reduced dataset, using only one book from Chesterton, but this approach resulted in notably worse generation quality. Increasing the batch size provided some improvement in training speed, though the overall process remained time-consuming.

The limitations in the generated text quality can be attributed primarily to the character-level approach of the model. By processing text one character at a time and only considering the previous 40 characters for context, the model struggles to grasp higher-level linguistic concepts such as word meanings, grammatical structures, and broader contextual relationships. This limitation manifests in the generated text as a lack of coherence and meaningful narrative flow, even though the model can sometimes produce valid English words.

Looking toward potential improvements, several things could enhance both the model's performance and the quality of generated text. The most straightforward enhancement would be implementing the training process on a GPU, which would significantly reduce training time and allow to train more epochs in fewer time. From an architectural perspective, transitioning to a word-level model could provide better semantic understanding and more coherent text generation. Additionally, implementing a transformer architecture, which has demonstrated impressive results in various text generation tasks, could significantly improve the quality of the generated content. The model could also benefit from an expanded training dataset incorporating more works from Chesterton and/or other authors, providing a richer foundation for learning language patterns and improving the overall generation capabilities.

### 3 Appendix - Code

All the code is available in the GitHub repository <https://github.com/Griffosx/nlp> under the src/task\_4 folder. For completeness, I include here the code for the main files used in this project.

File loader.py

```
1 import nltk
2 from nltk.corpus import gutenberg
3
4
5 class GutenbergLoader:
6     def __init__(self):
7         # Ensure the Gutenberg dataset is downloaded
8         nltk.download("gutenberg")
9
10    def list_available_works(self) -> list[str]:
11        """List all available works in the Gutenberg corpus."""
12        return gutenberg.fileids()
13
14    def list_author_works(self, author: str) -> list[str]:
15        """List all works by a specific author."""
16        return [
17            work
18            for work in self.list_available_works()
19            if author.lower() in work.lower()
20        ]
21
22    def load_chesterton_works(self) -> dict[str, str]:
23        """Load Chesterton's major works from the Gutenberg corpus."""
24        chesterton_works = {
25            "ball": "chesterton-ball.txt",
26            "brown": "chesterton-brown.txt",
27            "thursday": "chesterton-thursday.txt",
28        }
29
30        loaded_works = {}
31        for title, filename in chesterton_works.items():
32            try:
33                text = gutenberg.raw(filename)
34                loaded_works[title] = text
35                print(f"Successfully loaded: {filename}")
36            except Exception as e:
37                print(f"Error loading {filename}: {str(e)}")
38
39        return loaded_works
40
41    def get_combined_text(self, works: dict[str, str]) -> str:
42        """Combine all loaded works into a single text."""
43        return "\n\n".join(works.values())
44
45    def get_chesterton_combined_text(self) -> str:
46        """Load and combine Chesterton's major works."""
47        works = self.load_chesterton_works()
48        return self.get_combined_text(works)
49
50    def print_corpus_stats(self, works: dict[str, str]) -> None:
```

```

51     """Print basic statistics about the loaded works."""
52     print("\nCorpus Statistics:")
53     total_chars = 0
54     total_words = 0
55     total_lines = 0
56
57     for title, text in works.items():
58         num_chars = len(text)
59         num_words = len(text.split())
60         num_lines = len(text.splitlines())
61
62         total_chars += num_chars
63         total_words += num_words
64         total_lines += num_lines
65
66         print(f"\n{title.title():}")
67         print(f"Characters: {num_chars:,}")
68         print(f"Words: {num_words:,}")
69         print(f"Lines: {num_lines:,}")
70
71     print("\nTotal Statistics:")
72     print(f"Total Characters: {total_chars:,}")
73     print(f"Total Words: {total_words:,}")
74     print(f"Total Lines: {total_lines:,}")
75
76     def show_unique_characters(self, text: str) -> None:
77         """Display all unique characters in the text."""
78         unique_chars = sorted(list(set(text)))
79         print("\nUnique Characters in Corpus:")
80         print(f"Number of unique characters: {len(unique_chars)}")
81         print("Characters: ", " ".join(unique_chars))
82
83
84     def print_chesterton_info():
85         # Initialize the loader
86         loader = GutenbergLoader()
87
88         # List all Chesterton works available
89         print("Available Chesterton works in Gutenberg:")
90         chesterton_works = loader.list_author_works("chesterton")
91         for work in chesterton_works:
92             print(f"- {work}")
93
94         # Load Chesterton's works
95         print("\nLoading Chesterton's works...")
96         works = loader.load_chesterton_works()
97
98         # Print statistics for each work
99         loader.print_corpus_stats(works)
100
101         # Get combined text
102         combined_text = loader.get_combined_text(works)
103
104         # Show unique characters
105         loader.show_unique_characters(combined_text)

```

File encoder.py

```

1 import json

```

```

2 from loader import GutenbergLoader
3
4
5 class CharacterEncoder:
6     def __init__(self):
7         self.char_to_index: dict[str, int] = {}
8         self.index_to_char: dict[int, str] = {}
9         self.vocab_size: int = 0
10
11     def fit(self, text: str) -> None:
12         """Create character to index mappings from text."""
13         # Get unique characters and sort them for consistency
14         unique_chars = sorted(list(set(text)))
15         self.vocab_size = len(unique_chars)
16
17         # Create bidirectional mappings
18         self.char_to_index = {char: idx for idx, char in enumerate(
19             unique_chars)}
20         self.index_to_char = {idx: char for idx, char in enumerate(
21             unique_chars)}
22
23     def save_mappings(self, filename: str) -> None:
24         """Save character mappings to a JSON file."""
25         mappings = {
26             "char_to_index": self.char_to_index,
27             "index_to_char": {
28                 str(k): v for k, v in self.index_to_char.items()
29             }, # Convert int keys to str for JSON
30             "vocab_size": self.vocab_size,
31         }
32         with open(filename, "w", encoding="utf-8") as f:
33             json.dump(mappings, f, ensure_ascii=False, indent=2)
34
35     def load_mappings(self, filename: str) -> None:
36         """Load character mappings from a JSON file."""
37         with open(filename, "r", encoding="utf-8") as f:
38             mappings = json.load(f)
39         self.char_to_index = mappings["char_to_index"]
40         self.index_to_char = {
41             int(k): v for k, v in mappings["index_to_char"].items()
42         } # Convert str keys back to int
43         self.vocab_size = mappings["vocab_size"]
44
45     def print_mappings(self) -> None:
46         """Print character mappings in a readable format."""
47         print("\nCharacter to Index Mappings:")
48         print("-" * 30)
49         for char, idx in sorted(self.char_to_index.items(), key=lambda
50             x: x[1]):
51             if char.isspace():
52                 char_display = f"[space-{ord(char)}]"
53             elif char == "\n":
54                 char_display = "[newline]"
55             elif char == "\t":
56                 char_display = "[tab]"
57             else:
58                 char_display = char
59             print(f"'{char_display}' -> {idx}")

```



```
57 print(f"\nTotal vocabulary size: {self.vocab_size}")
```

# File generator.py

```
1 import numpy as np
2 import torch
3 import torch.nn as nn
4 from torch.utils.data import Dataset, DataLoader
5 from task_4.loader import GutenbergLoader
6 from task_4.encoder import CharacterEncoder
7
8
9 class TextDataset(Dataset):
10     def __init__(self, text: str, char_encoder: CharacterEncoder,
11                  seq_length: int):
12         self.text = text
13         self.char_encoder = char_encoder
14         self.seq_length = seq_length
15         self.text_indices = [char_encoder.char_to_index[char] for char
16                             in text]
17         self.num_sequences = len(self.text_indices) - self.seq_length -
18                             1
19
20     def __len__(self) -> int:
21         return self.num_sequences
22
23     def __getitem__(self, idx: int) -> tuple[torch.Tensor, torch.Tensor
24 ]:
25         # Get the sequence and target
26         sequence = self.text_indices[idx : idx + self.seq_length]
27         target = self.text_indices[idx + self.seq_length]
28
29         # Create one-hot encoding for input sequence
30         X = torch.zeros(
31             self.seq_length, self.char_encoder.vocab_size, dtype=torch.
32             float32
33         )
34         for t, char_idx in enumerate(sequence):
35             X[t, char_idx] = 1
36
37         # Return target as a single integer (not one-hot encoded)
38         return X, torch.tensor(target, dtype=torch.long)
39
40
41 class LSTM(nn.Module):
42     def __init__(self, input_size: int, hidden_size: int = 128):
43         super(LSTM, self).__init__()
44         self.hidden_size = hidden_size
45         self.lstm = nn.LSTM(input_size, hidden_size, num_layers=1,
46                             batch_first=True)
47         self.fc = nn.Linear(hidden_size, input_size)
48
49     def forward(self, x: torch.Tensor) -> torch.Tensor:
50         x, _ = self.lstm(x)
51         # Return logits (softmax will be applied in the generation step
52         )
53         return self.fc(x[:, -1, :])
```

```

49 class TextGenerator:
50     def __init__(
51         self,
52         seq_length: int = 40,
53         device: str = "mps",
54     ):
55         self.seq_length = seq_length
56         self.device = device
57         self.model = None
58
59     def build_model(self, vocab_size: int) -> None:
60         """Construct the LSTM model."""
61         self.model = LSTM(vocab_size).to(self.device)
62
63     def train(
64         self,
65         text: str,
66         char_encoder: CharacterEncoder,
67         model_save_path: str = "model.pth",
68         encoder_save_path: str = "char_encoder.json",
69         epochs: int = 10,
70         batch_size: int = 32,
71         learning_rate: float = 0.001,
72     ) -> list[float]:
73         """Train the model and save it to disk."""
74         if self.model is None:
75             self.build_model(char_encoder.vocab_size)
76
77         dataset = TextDataset(text, char_encoder, self.seq_length)
78         dataloader = DataLoader(dataset, batch_size=batch_size, shuffle
                                =True)
79
80         # CrossEntropyLoss expects raw logits and target indices
81         criterion = nn.CrossEntropyLoss()
82         optimizer = torch.optim.RMSprop(self.model.parameters(), lr=
            learning_rate)
83
84         losses = []
85         self.model.train()
86
87         for epoch in range(epochs):
88             epoch_loss = 0
89             batch_count = 0
90
91             for batch_X, batch_y in dataloader:
92                 batch_X = batch_X.to(self.device)
93                 batch_y = batch_y.to(self.device)
94
95                 optimizer.zero_grad()
96                 # Get logits from model
97                 logits = self.model(batch_X)
98                 # CrossEntropyLoss expects logits and target class
99                 # indices
100                 loss = criterion(logits, batch_y)
101
102                 loss.backward()
103                 optimizer.step()

```

```

104         epoch_loss += loss.item()
105         batch_count += 1
106
107         del batch_X, batch_y, logits, loss
108
109         if batch_count % 100 == 0:
110             print(f"Epoch {epoch+1}, Batch {batch_count}/{len(
111                 dataloader)}")
112
113             avg_loss = epoch_loss / batch_count
114             losses.append(avg_loss)
115             print(f"Epoch [{epoch+1}/{epochs}], Loss: {avg_loss:.4f}")
116
117             # Save model and encoder
118             torch.save(self.model.state_dict(), model_save_path)
119             char_encoder.save_mappings(encoder_save_path)
120             print(f"Model saved to {model_save_path}")
121             print(f"Character encoder saved to {encoder_save_path}")
122
123             return losses
124
125     def generate_text(
126         self,
127         seed_text: str,
128         model_load_path: str = "model.pth",
129         encoder_load_path: str = "char_encoder.json",
130         length: int = 100,
131         diversity: float = 0.5,
132     ) -> str:
133         """Load model and generate text."""
134         char_encoder = CharacterEncoder()
135         char_encoder.load_mappings(encoder_load_path)
136
137         if self.model is None:
138             self.build_model(char_encoder.vocab_size)
139             self.model.load_state_dict(
140                 torch.load(model_load_path, map_location=self.device,
141                     weights_only=True)
142             )
143
144         if len(seed_text) != self.seq_length:
145             raise ValueError(f"Seed text must be {self.seq_length}
146                 characters long")
147
148         self.model.eval()
149         current_sequence = seed_text
150         generated_text = seed_text
151
152         with torch.no_grad():
153             for _ in range(length):
154                 # Prepare input tensor
155                 x_pred = torch.zeros(
156                     (1, self.seq_length, char_encoder.vocab_size),
157                     dtype=torch.float32
158                 )
159                 for t, char in enumerate(current_sequence):
160                     if char in char_encoder.char_to_index:
161                         x_pred[0, t, char_encoder.char_to_index[char]]

```

```

158                                     = 1
159
160     # Generate prediction
161     x_pred = x_pred.to(self.device)
162     # Get logits from model
163     logits = self.model(x_pred)
164
165     # Apply softmax and temperature scaling for generation
166     preds = torch.softmax(logits / diversity, dim=-1).cpu()
167         .numpy()[0]
168
169     # Sample next character
170     next_index = np.random.choice(len(preds), p=preds)
171     next_char = char_encoder.index_to_char[next_index]
172
173     # Update sequences
174     generated_text += next_char
175     current_sequence = current_sequence[1:] + next_char
176
177     del x_pred, logits
178
179     return generated_text
180
181 def train():
182     loader = GutenbergLoader()
183     encoder = CharacterEncoder()
184     generator = TextGenerator(seq_length=40)
185
186     combined_text = loader.get_chesterton_combined_text()
187     encoder.fit(combined_text)
188
189     epochs = 10
190     batct_size = 512
191     losses = generator.train(
192         text=combined_text, char_encoder=encoder, epochs=epochs,
193         batch_size=batct_size
194     )
195     print("Training completed!")
196
197 def generate():
198     generator = TextGenerator(seq_length=40)
199
200     seed_text = (
201         "There was an instant of rigid silence, and then Syme in his
202         turn fell "
203         "furiously on the other, filled with a flaming curiosity."
204     )[:40]
205     print(f"Seed text: {seed_text}")
206
207     print("\nSeed text:")
208     print(seed_text)
209     print("\nGenerated texts with different diversity levels:")
210
211     for diversity in [0.2, 0.5, 1.0]:
212         print(f"\nDiversity: {diversity}")
213         generated = generator.generate_text(

```

```
212         seed_text=seed_text, length=100, diversity=diversity
213     )
214     print(generated)
```