

---

**THE GEORGE WASHINGTON UNIVERSITY**

---

WASHINGTON, DC

## Final Project Report

Chirag Lakhanpal

The George Washington University

DATS 6312: Natural Language Processing

Professor Amir Jafari

## Table of Contents

Introduction	2
Dataset	3
Individual Work (My contribution and Work Description)	4
1. Fine Tuning GPT2 Model	4
2. Data Preparation	4
3. Training Process and Hyper Parameter tuning	5
4. Streamlit Application Development	6
5. Things That Did Not Work	7
Results and Summary	8
1. ROUGE Scores	9
2. Training and Validation Loss	9
3. Perplexity	9
4. Observations from Extended Epoch Training	9
Best model Parameters	10
Conclusions	10
Code Contribution	11
References	12

## Introduction

This report covers a concentrated exploration in the topic of legal text analysis, especially in the domain of text production. It describes the application and fine-tuning of the GPT-2 model for producing legal texts using a subset of the legal\_contracts dataset that comprises 10% (65083 records) and 1% (6508 records). This technique tries to capitalize on GPT-2's sophisticated capabilities while adapting it to the complex and formal structure of legal language.

This project is significant because it has the potential to contribute legal document drafting and advisory creation through automated text generation. The project explores new possibilities in the efficiency and accuracy of legal text processing by fine-tuning GPT-2 using a properly selected legal dataset.

## Dataset

Link - [https://huggingface.co/datasets/albertvillanova/legal\\_contracts](https://huggingface.co/datasets/albertvillanova/legal_contracts)

For training the GPT-2 model in legal text generation, this project utilized a carefully selected subset of the albertvillanova/legal\_contracts dataset. From the original 650,833 records (33 GB), a 10% subset (65,083 records, 2.87 GB) was chosen for the main training phase, balancing data diversity with computational efficiency. Additionally, a 1% subset (6,508 records, 260 MB) was used specifically for hyperparameter tuning, optimizing model performance without excessive computational demand.

Quicklinks - Click here to rapidly navigate through this document

AMENDED AND RESTATED EMPLOYMENT AND NONCOMPETITION AGREEMENT

THIS AMENDED AND RESTATED EMPLOYMENT AND NONCOMPETITION AGREEMENT (the "Agreement") is made and entered into as of October 1, 2008, by and among Avocent Employment Services Co. (formerly known as Polycorn Investments, Inc.), a Texas corporation ("Employer"), Avocent Corporation, a Delaware corporation, and its direct or indirect subsidiary of Avocent Corporation engaged in the business of leasing employees to Avocent Corporation and its affiliates, including Apex, Inc. ("Apex") and Cyrex Computer Products Corporation ("Cyrex"); and the Employee, who is a direct or indirect subsidiary of Avocent Corporation and its affiliates (collectively referred to in this Agreement as "Avocent"). The Employee, Apex, and Cyrex are engaged in the business of designing, manufacturing, and selling stand-alone console/KVM switching systems, console/KVM remote access products, and integrated server cabinet solutions for the client/server computing market.

The Employee, Apex, and Cyrex entered into that certain Employment and Noncompetition Agreement dated July 1, 2008 (the "Original Employment Agreement"); and on March 1, 2008, Apex, Cyrex, and Avocent Corporation entered into the Reorganization Agreement dated March 1, 2008 (the "Reorganization Agreement"). Pursuant to the Reorganization Agreement, (i) Apex Acquisition Corp., a wholly-owned subsidiary of Avocent, merged with and into Apex on July 1, 2008 (the "Apex Merger"), and upon the Apex Merger, Apex became a wholly-owned subsidiary of Avocent; and (ii) Cyrex Acquisition Corp., a wholly-owned subsidiary of Avocent, merged with and into Cyrex (the "Cyrex Merger") on July 1, 2008, and upon the Cyrex Merger, Cyrex also became a wholly-owned subsidiary of Avocent; and in consideration of an increase in base pay, certain incentive bonus eligibility and awards, and an award of stock options that would not otherwise be made to Employee, Employee, Cyrex, and Avocent now wish to amend and restate the Original Employment Agreement with this Amended and Restated Employment and Noncompetition Agreement.

THE PARTIES HERETO AGREE AS FOLLOWS:

**1. DUTIES**

During the term of this Agreement, the Employee agrees to be employed by Employer and to serve Avocent as its Senior Vice President of Operations and Chief Operating Officer, and Employer agrees to employ the Employee and lease the Employee to Avocent to serve Avocent in such capacities. The Employee shall devote such of his business time, energy, and skill to the affairs of Avocent and Employer as shall be necessary to perform the duties of Senior Vice President of Operations and Chief Operating Officer. The Employee shall report to the President of the Employer, Cyrex, and Avocent Corporation and to the boards of directors of the Employer, Cyrex, and Avocent Corporation at all times during the term of this Agreement. The Employee shall have powers and duties at least commensurate with his position as Senior Vice President of Operations and Chief Operating Officer of Avocent Corporation.

**2. TERM OF EMPLOYMENT**

The Employee's employment with the Employer shall be for an indefinite period of time. For purposes of this Agreement the following terms shall have the following meanings:

**TERMINATION FOR CAUSE** shall mean termination by the Employer of the Employee's employment by the Employer for cause, including but not limited to, fraud upon, or deliberate injury to, the Employer or Avocent or its affiliates.

**TERMINATIONS OTHER THAN FOR CAUSE** shall mean termination by the Employer or Avocent Corporation of the Employee's employment by the Employer (other than in a Termination for Cause) and shall include any constructive termination of the Employee's employment by reason of material breach of this Agreement by the Employer or Avocent, such constructive termination to be effective upon thirty (30) days written notice from the Employer to the Employee of such constructive termination.

**VOLUNTARY TERMINATION** shall mean termination by the Employee of the Employee's employment by the Employer other than in a Termination for Cause and shall include any constructive termination of the Employee's employment by reason of material breach of this Agreement by the Employer or Avocent, such constructive termination to be effective upon thirty (30) days written notice from the Employee to the Employer of such constructive termination.

**CHANGE IN CONTROL** shall mean (i) a change in control as described in Section 2.1(b), and (ii) a termination by reason of the Employee's disability or death as described in Sections 2.5 and 2.6.

**TERMINATION UPON A CHANGE IN CONTROL** shall mean (i) a termination by the Employer of the Employee's employment with the Employer or services to Avocent within six (6) months following any "Change in Control" other than a "Termination for Cause" contemplated by or described in the Reorganization Agreement and/or resulting from the closing of the transactions described in the Reorganization Agreement including, without limitation, the Cyrex Merger, the Apex Merger, and the Merger (as such terms are defined in the Reorganization Agreement), or (ii) a termination by the Employer or Avocent Corporation of the Employee's employment by the Employer (other than a termination for cause) within eighteen (18) months following any "Change in Control" other than a "Change in Control" contemplated by or described in the Reorganization Agreement and/or resulting from the closing of the transactions described in the Reorganization Agreement including, without limitation, the Cyrex Merger, the Apex Merger, and the Merger (as such terms are defined in the Reorganization Agreement).

**CHANGE IN CONTROL** shall mean any one of the following events:

(a) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(b) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(c) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(d) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(e) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(f) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(g) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(h) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(i) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(j) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(k) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(l) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(m) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(n) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(o) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(p) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(q) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(r) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(s) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(t) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(u) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(v) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(w) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(x) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(y) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

(z) any person (other than Avocent) acquires beneficial ownership of ten percent (10%) or more of the combined voting power of Employer's, Cyrex's, or Avocent Corporation's then outstanding stock. For purposes of this Agreement, "beneficial ownership" shall be determined in accordance with the rules and regulations of the Securities Exchange Act of 1934, or any similar successor regulation or rule; and the term "person" shall include any natural person, corporation, partnership, trust or association, or any group or combination thereof, whose ownership of Employer's, Cyrex's, or Avocent Corporation's securities would be required to be reported under such regulation or rule.

Figure 1: First record for the dataset depicting Employment Non-Compete Agreement.

## Individual Work (My contribution and Work Description)

My contribution in this project was varied, covering the building of a Streamlit application to illustrate the model's capabilities in a real-world setting as well as the fine-tuning of the GPT-2 model for legal text synthesis.

### 1. Fine Tuning GPT2 Model

- **Model and Tokenizer Setup:** The project started with setting up the GPT-2 model and tokenizer. GPT-2, which is well-known for its powerful language production skills, was chosen for its ability to handle complicated language structures, such as those found in legal writings. I also tried BERT and DistilBERT but was unable to get them to function.

### 2. Data Preparation

- **Dataset Processing:** Normalization of text from the albertvillanova/legal\_contracts dataset was a vital step in the preprocessing procedure. To ensure text quality and uniformity, non-ASCII characters were removed throughout this procedure. To guarantee consistency, the text was transformed to lowercase.
- **Regular Expression Usage:** Regular expressions were extensively used for cleaning the data, including.
  - i. Replacing sequences of dashes with spaces
  - ii. Removing extraneous newline, carriage return, and tab characters
  - iii. Stripping excessive whitespace.
  - iv. Spell checking the words. Due to computation limitation this could not work.
  - v. NER on names, person, and organizations. Again, due to the sheer amount of the data, it got computationally expensive.

```

preprocess_function(examples):
    # Convert to lowercase and remove non-ASCII characters
    processed_text = [re.sub(r'[^\x00-\x7F]+', ' ', text).lower() for text in examples['text']]

    # Remove extra "-" characters
    processed_text = [re.sub(r'---+', ' ', text) for text in processed_text]

    # Remove newlines, tabs, and carriage returns
    processed_text = [text.replace('\n', ' ').replace('\r', ' ').replace('\t', ' ') for text in processed_text]

    # Remove extra whitespace
    processed_text = [re.sub(r'\s+', ' ', text).strip() for text in processed_text]

    # Remove numbers
    processed_text = [re.sub(r'\b\d+\b', '', text) for text in processed_text]

    # anonymized_text = []
    # for text in processed_text:
    #     doc = nlp(text)
    #     for ent in doc.ents:
    #         if ent.label_ in ["PERSON", "ORG", "GPE"]:
    #             text = text.replace(ent.text, f'[{ent.label_}]')
    #     anonymized_text.append(text)

    # corrected_text = []
    # for text in anonymized_text:
    #     blob = TextBlob(text)
    #     corrected_text.append(blob.correct().string)

    return tokenizer(processed_text, truncation=True, padding='max_length', max_length=512)

```

Figure 2: Python function for preprocessing text data, showing various steps including normalization, cleaning, and tokenization.

### 3. Training Process and Hyper Parameter tuning

- **Epoch Configuration:** For the smaller dataset, epochs were set at **1, 10, and 50**. For the larger dataset, **25 epochs were used**, which took around ~45 mins per epoch, balancing depth of learning with computational efficiency.
- **Batch Size Variation:** Two different batch sizes, 8 and 16, were experimented with. I tried with bigger batch sizes starting from 128 and reducing it to 32 in the power of 2, however I encountered “CUDA out of memory” error.
- **Optimizer Selection:** For the training used two different optimizers- **AdamW and SGD (Stochastic Gradient Descent)**. **AdamW**.
- **Learning Rate Experimentation:** Multiple learning rates were tested: **5e-5, 1e-5, 5e-4, 5e-3, and 5e-2**. This range allowed for observing how the model responded to both subtle and significant changes in learning rate, influencing the speed and stability of the learning process.
- **Gradient Clipping:** To address the potential issue of exploding gradients, gradient clipping was incorporated. This technique involved capping the gradients during backpropagation to a predefined range, ensuring they did not exceed manageable levels.

## 4. Streamlit Application Development

- The Streamlit application for this project is structured into two primary scripts: one for the chat functionality and the other for the main application, featuring a landing page for task selection. Here's a brief overview of each aspect:
- Landing Page and Navigation: Users are first directed to a landing page offering task selections such as document summarization, classification, or initiating a legal chat. This is facilitated through a tabbed layout for easy navigation.
- Custom Styling: I applied custom CSS styling, enhancing the visual appeal and user experience.
- Input Sanitization: A crucial aspect of the app is the sanitization of user inputs using regular expressions, ensuring the inputs are clean and concise for optimal model response.

## 5. Things That Did Not Work

- In this project, an initial attempt to create a Q&A style chatbot using the `nguha/legalbench` dataset, which contains 162 specialized tasks, faced significant challenges.
- **Computational Power Limitations:** Initially, I tried to process the entire dataset in large batches for model training. This approach quickly led to out-of-memory errors, indicating that the dataset's size and complexity were too demanding for the available computational resources.
- **Modular Approach Issues:** Shifting to a modular strategy, where the dataset was fed one task at a time, did not resolve these memory issues. It became clear that both the batch size and the individual task complexity contributed to the computational challenges.
- **Dataset Specificity:** 'Yes' or 'no' responses to specific legal questions made up most of the sample. Because of its narrow focus, the chatbot was less useful to a broader audience because it was unable to adequately respond to a larger variety of legal inquiries.

```

[ec2-user@ip-10-0-0-187 NLP]$ /opt/conda/envs/pytorch/bin/python "/home/ec2-user/NLP/Fine Tuning GPT2.py"
2023-12-09 02:53:08.870338: E external/local_xla/xla/stream_executor/cuda/cuda_dnn.cc:9261] Unable to register cuDNN factory: Attempting to register factory for plugin cuDNN when one has already been registered
2023-12-09 02:53:08.870389: E external/local_xla/xla/stream_executor/cuda/cuda_fft.cc:607] Unable to register cuFFT factory: Attempting to register factory for plugin cuFFT when one has already been registered
2023-12-09 02:53:08.871207: E external/local_xla/xla/stream_executor/cuda/cuda_blas.cc:1515] Unable to register cuBLAS factory: Attempting to register factory for plugin cuBLAS when one has already been registered
2023-12-09 02:53:08.877060: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2023-12-09 02:53:09.639536: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
Training on subset: canada_tax_court_outcomes
11%|██████████          | 1/9 [00:00<00:05, 1.33it/s]
Traceback (most recent call last):
  File "/home/ec2-user/NLP/Fine Tuning GPT2.py", line 102, in <module>
    train_and_evaluate()
  File "/home/ec2-user/NLP/Fine Tuning GPT2.py", line 89, in train_and_evaluate
    trainer.train()
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer.py", line 1530, in train
    return inner_training_loop(
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer.py", line 1919, in _inner_training_loop
    self._maybe_log_save_evaluate(tr_loss, model, trial, epoch, ignore_keys_for_eval)
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer.py", line 2253, in _maybe_log_save_evaluate
    metrics = self.evaluate(ignore_keys=ignore_keys_for_eval)
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer.py", line 2983, in evaluate
    output = eval_loop()
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer.py", line 3198, in evaluation_loop
    preds_host = logits if preds_host is None else nested_concat(preds_host, logits, padding_index=-100)
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer_pt_utils.py", line 123, in nested_concat
    return torch_pad_and_concatenate(tensors, new_tensors, padding_index=padding_index)
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/transformers/trainer_pt_utils.py", line 82, in torch_pad_and_concatenate
    return torch.cat((tensors, new_tensors), dim=0)
  File "/opt/conda/envs/pytorch/lib/python3.10/site-packages/torch/cuda/__init__.py", line 222, in cat
    return torch.cat(tensors, dim)
torch.cuda.OutOfMemoryError: CUDA out of memory. Tried to allocate 9.28 GiB (GPU 0; 22.19 GiB total capacity; 9.69 GiB already allocated; 9.05 GiB free; 12.82 GiB reserved in total by PyTorch. If reserved memory is >> allocated memory try setting max_split_size_mb to avoid fragmentation. See documentation for Memory Management and PYTORCH_CUDA_ALLOC_CONF)
3/9 [00:05<00:10, 1.79s/it]

```

Figure 3: Screenshot of a CUDA out-of-memory error encountered during model training.



## Results and Summary

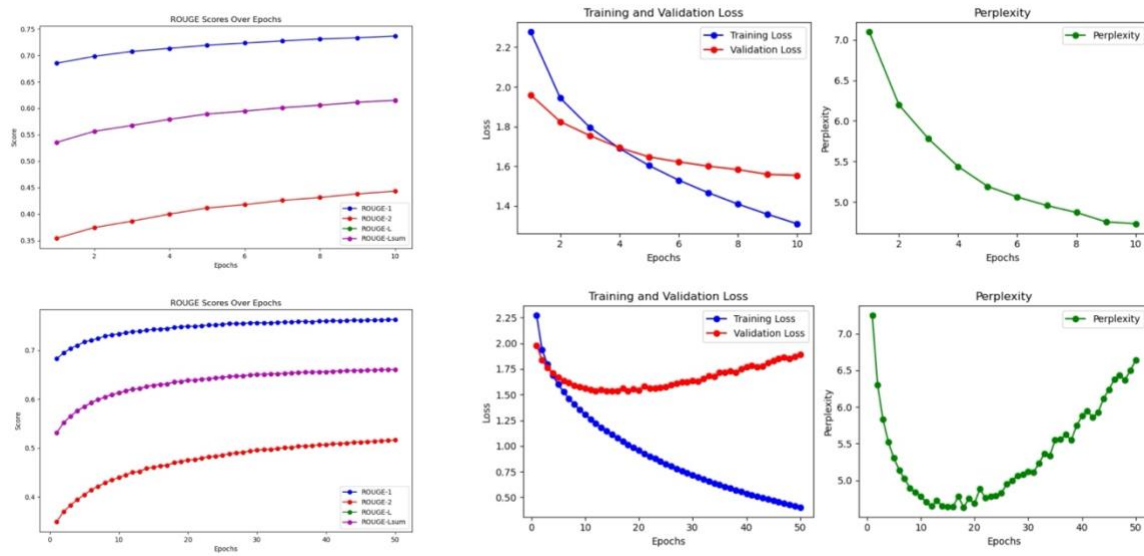


Figure 4: Comparative analysis of ROUGE scores (left), training and validation losses, and perplexity (right) over different numbers of training epochs on a small dataset (~300 MB).

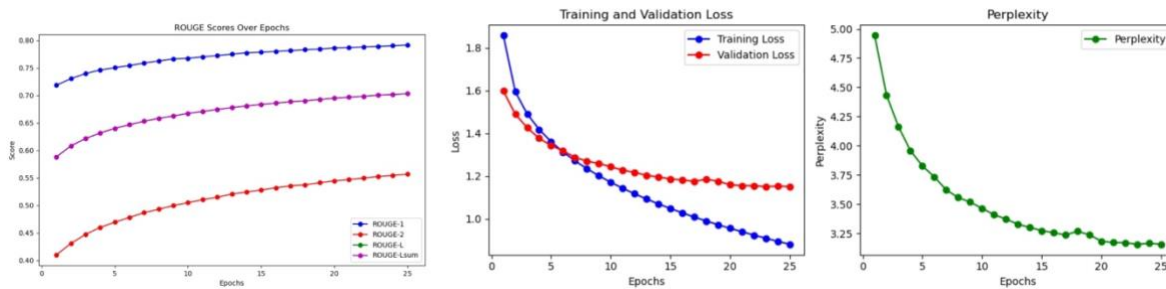


Figure 5: Comparative analysis of ROUGE scores (left), training and validation losses, and perplexity (right) over different numbers of training epochs on a large dataset (~3 GB).

The figures displayed represent the results of training a text generation model across different epochs, as evidenced by ROUGE scores, training and validation losses, and perplexity metrics.

## 1. ROUGE Scores

- ROUGE (Recall-Oriented Understudy for Gisting Evaluation) scores are a set of metrics for evaluating automatic summarization and machine translation software. They compare the overlap of n-grams, word sequences, and word pairs between the computer-generated output and a set of reference texts. Higher scores indicate better performance, with ROUGE-1, ROUGE-2, and ROUGE-L representing the overlap of unigrams, bigrams, and the longest common subsequence, respectively.
- Top Graphs: The top left graph illustrates an improvement in ROUGE scores across all three metrics over 10 epochs. The consistent upward trend across ROUGE-1, ROUGE-2, and ROUGE-L sum indicates the model's increasing proficiency in generating text closely matching the reference texts in terms of content and structure.

## 2. Training and Validation Loss

- Loss is a measure of how far the model's predictions are from the actual outcomes. Lower loss values correspond to better model performance.
- Top Middle Graphs: The training loss (blue) and validation loss (red) both show a downward trend over 10 epochs, suggesting that the model is learning effectively and generalizing well to unseen data.

## 3. Perplexity

- Perplexity is a measurement of how well a probability distribution or probability model predicts a sample. In the context of language models, lower perplexity indicates a better predictive performance.
- Top Right Graph: A decreasing perplexity trend over 10 epochs is observed, implying the model's improving ability to predict the next word in the sequence.

## 4. Observations from Extended Epoch Training

- When extending the training to 50 epochs, the graphs show distinct trends.
- Bottom Graphs: The ROUGE scores continue to improve, although they begin to plateau, indicating that the model might be approaching its optimal performance. The training and validation loss graphs depict a notable divergence after around 20 epochs, potentially suggesting overfitting where the model learns the training data too

well and may not perform as effectively on new, unseen data. The perplexity initially decreases, reflecting improved model predictions, but then increases dramatically, further implying that the model's performance on the validation set is worsening, possibly due to overfitting.

## Best model Parameters

Parameter	Best Value	Experimental Values
Learning rate	5e-5	5e-5, 1e-5, 5e-4, 5e-3, and 5e-2
Optimizer	AdamW	AdamW and SGD
Batch Size	8	8, and 16 (Working)
Momentum (For SGD)	0.9	0.9
Epochs	10	1,10,25,50
Number of Workers (Data Loading)	8	4, 8, and 16
Dataset Size	-	10% (65083 records) and 1% (6508 records)

```
(pytorch) [ec2-user@ip-10-0-0-187 NLP] $ /opt/conda/envs/pytorch/bin/python "/home/ec2-user/NLP/Fine Tuning GPT2 Inference.py"
Enter a prompt: Draft a sales agreement contract that outlines terms and conditions for the sale of goods between two parties.
Draft a sales agreement contract that outlines terms and conditions for the sale of goods between two parties, this contract is entered into by and between the following parties on the 1st day of January, (hereinafter referred to as the "effective date"): party a: taiyuan putal business consulting co., ltd. legal address: no. xuefu street, shangdi, haidian district, beijing legal representative: mr. qingjie sheng party b: shanxi puda resources international, inc., a company incorporated under the laws of the province of british columbia, canada, and having its principal place of business at suite, west broadway, provo, utah, v6c 2j1 party c: zhao ming party d: xin jia party e: shenzhen hong party f: li shao party g: jianquan li whereas: (a) the parties have agreed to establish a joint venture company in the people's republic of china (the "prc") in accordance with the relevant laws and regulations of prc; and (b) it is the intention of both parties that this agreement shall regulate their relations and undertakings. now, therefore, for good and valuable consideration, the receipt and sufficiency of which is hereby acknowledged by each party, they hereby agree as follows: definitions. "contract" shall mean this sales contract and all exhibits and schedules attached hereto and made a part hereof, together with all amendments, modifications, supplements and extensions thereto and any exhibits or schedules to any of them which may be executed and/or delivered hereunder by either party and are incorporated herein by reference; "party a's address" for notices shall be at the address set forth above or at such other place or to such party as may from time to time be notified to the other party by written notice in writing. product means any natural, chemical, oil, gas, hydrocarbon substances and other petroleum products, including but not limited to oil and gas liquids, electricity, steam, air, water and riparian waste, crops, timber for crop nutrition, livestock, poultry, wild animals and wild game, all types of wild and domestications as well as plants and trees, shrubbery, branches, garnishments, solids, scales, emblems, badges, letters of identification, customs, patents, patent applications, trade names, copyrights, rights to sue and recover for past infringement or misappropriation of any patent, trademark or other intellectual property rights
(pytorch) [ec2-user@ip-10-0-0-187 NLP] $
```

Figure 6: Inference on custom model trained on large dataset.

## Conclusions

- ROUGE scores improved consistently across 10 epochs, indicating better alignment with reference texts in the model's output.
- Both training and validation loss decreased over 10 epochs, suggesting effective learning and generalization to unseen data.
- Perplexity metrics showed a downward trend over 10 epochs, reflecting the model's improved predictive accuracy.
- In extended training up to 50 epochs, ROUGE scores began to plateau, suggesting a limit to the model's improvement in text generation.

- Divergence between training and validation loss after 20 epochs indicated potential overfitting, where the model excessively learned from the training data at the expense of generalization.
- A dramatic increase in perplexity after an initial decrease during extended training suggested diminishing returns in predictive performance, possibly due to overfitting.

## **Code Contribution**

Total lines of code = 427 (Streamlit + train). Modified lines = ~80. Added lines = ~40.

Percentage = 89.66%

## References

1. Alammar, J. (n.d.). The Illustrated GPT-2 (Visualizing Transformer Language Models). Retrieved from <https://jalammar.github.io/illustrated-gpt2/>
2. Streamlit Docs. (n.d.). Build a ChatGPT-like app. Retrieved from <https://docs.streamlit.io/knowledge-base/tutorials/build-conversational-apps#build-a-chatgpt-like-app>
3. OpenAI. (n.d.). ChatGPT.
4. Papers with Code. (n.d.). GPT-2. Retrieved from <https://paperswithcode.com/method/gpt-2>
5. Image References:
  - a. ResearchGate. (n.d.). Structure of the applied GPT-2 medium architecture. Retrieved from [https://www.researchgate.net/figure/Structure-of-the-applied-GPT-2-medium-architecture\\_fig2\\_365625866](https://www.researchgate.net/figure/Structure-of-the-applied-GPT-2-medium-architecture_fig2_365625866)
  - b. Alammar, J. (n.d.). The Illustrated GPT-2 (Visualizing Transformer Language Models). Retrieved from <https://jalammar.github.io/illustrated-gpt2/>