

Assignment 5: Natural Language Generation

Instructor: Luke Zettlemoyer

CSED503

Due: 11:59pm PT, 8th June 2026

In this assignment, you will learn about generating text from neural language models using different decoding algorithms. You will also learn how to fine-tune generative models, showcased through knowledge distillation: using a bigger teacher language model to generate data, which is used to fine-tune a smaller language model and improve its performance.

You will submit both your **code** and **writeup** (as PDF) via Gradescope.

Required Deliverables

- **Code Notebook:** This assignment has an associated Jupyter notebook (CSED503_Assignment5.ipynb). You need to submit the notebook with your solutions. On Google Colab you can do so by **File** → **Download** → **Download .ipynb**. **Please comment out any additional code you had written to solve the write-up exercises before submitting on gradescope to avoid timeouts.**
- **Write-up:** For written answers and open-ended reports, produce a single PDF for §1-2 and submit it in Gradescope. We recommend using Overleaf to typeset your answers in L^AT_EX, but other legible typed formats are acceptable.

Recommended Reading

The homework is based on lectures on language generation, decoding algorithms, and knowledge distillation, so the lecture slides should be your best resource. For more detailed reading we recommend checking Chapters 9 and 10 of [Jurafsky and Martin](#). We also recommend checking Patrick von Platen's great blog post on [decoding algorithms](#).

Required Compute

Except §1.1, all of the exercises will require you to use a GPU to run your code. We have tested the reference implementations on the free tier T4 GPU on Colab and you should be able to use it to solve the exercises. The most compute intensive exercise is the Knowledge Distillation section. If you run into any issues with the compute please contact the course staff.

Acknowledgement

This assignment is adapted by Hamish Ivison, designed by Kabir Ahuja with help from Sofia Serrano and Nikita Haduong in ironing out issues and improving clarity of the assignments. Melanie Sclar, Khushi Khandelwal, Melissa Mitchell, and Kavel Rao provided invaluable feedback and helped in testing the notebooks. §1 of this homework is adapted from assignments created by Yegor Kuznetsov, Liwei Jiang, Jaehun Jung, and Gary Jiacheng Liu.

1 Natural Language Generation (100 pts)

In this part of the homework, you will learn about generating text from neural language models using different decoding algorithms. You will also learn how to fine-tune generative models, showcased through knowledge distillation.

1.1 Decoding Algorithms (60 pts)

You will implement and experiment with various decoding algorithms for language generation. In particular, we will focus on basic decoding techniques like greedy decoding, random sampling, temperature sampling, and top-p/top-k sampling. While these algorithms are conceptually simple, these have become ubiquitous in the era of modern LLMs and if you have used any modern LLM systems like ChatGPT, these are the decoding algorithms that these models use to generate text!

1.1.0 Set Up Evaluation Metrics

Dataset In this assignment, we focus on the open-ended story generation task (data available [here](#)). This dataset contains *prompts* for story generation, modified from the [ROCStories dataset](#).

Evaluation Metrics :

- **Fluency:** [The CoLA classifier](#) is a RoBERTa-large classifier trained on the CoLA corpus (Warstadt et al., 2019), which contains sentences paired with grammatical acceptability judgments. We will use this model to evaluate fluency of generated sentences.
- **Diversity:** **The Count of Unique N-grams** is used to measure the diversity of the generated sentences.
- **Naturalness:** **The Perplexity** of generated sentences under the language model is used to measure the naturalness of language. You can directly use [the perplexity function from HuggingFace evaluate-metric package](#) for this assignment.

1.1.1 Greedy Decoding

The idea of greedy decoding is simple: select the next token as the one that receives the highest probability. **Implement the greedy() function that processes tokens in batch.** Its input argument `next_token_logits` is a 2-D FloatTensor where the first dimension is batch size and the second dimension is the vocabulary size, and you should output `next_tokens` which is a 1-D LongTensor where the first dimension is the batch size.

The softmax function is monotonic—in the same vector of logits, if one logit is higher than the other, then the post-softmax probability corresponding to the former is higher than that corresponding to the latter. Therefore, for greedy decoding you won't need to actually compute the softmax.

1.1.2 Vanilla Sampling, Temperature Sampling

To get more diverse generations, you can randomly sample the next token from the distribution implied by the logits. This decoding is called sampling, or vanilla sampling (since we will see more variations of sampling). Formally, the probability of for each candidate token w is

$$p(w) = \frac{\exp z(w)}{\sum_{w' \in V} \exp z(w')}$$

where $z(w)$ is the logit for token w , and V is the vocabulary. This probability on all tokens can be derived at once by running the softmax function on vector \mathbf{z} .

Temperature sampling controls the randomness of generation by applying a temperature t when computing the probabilities. Formally,

$$p(w) = \frac{\exp(z(w)/t)}{\sum_{w' \in V} \exp(z(w')/t)}$$

where t is a hyper-parameter.

Implement the `sample()` and `temperature()` functions. When testing the code we will use $t = 0.8$, but your implementation should support arbitrary $t \in (0, \infty)$.

1.1.3 Top- k Sampling

Top- k sampling decides the next token by randomly sampling among the k candidate tokens that receive the highest probability in the vocabulary, where k is a hyper-parameter. The sampling probability among these k candidate tokens should be proportional to their original probability implied by the logits, while summing up to 1 to form a valid distribution.

Implement the `topk()` function that achieves this goal. When testing the code we will use $k = 20$, but your implementation should support arbitrary $k \in [1, |V|]$.

1.1.4 Top- p Sampling

Top- p sampling, or nucleus sampling, is a bit more complicated. It considers the smallest set of top candidate tokens such that their cumulative probability is greater than or equal to a threshold p , where $p \in [0, 1]$ is a hyper-parameter. In practice, you can keep picking candidate tokens in descending order of their probability, until the cumulative probability is greater than or equal to p (though there's more efficient implementations). You can view top- p sampling as a variation of top- k sampling, where the value of k varies case-by-case depending on what the distribution looks like. Similar to top- k sampling, the sampling probability among these picked candidate tokens should be proportional to their original probability implied by the logits, while summing up to 1 to form a valid distribution.

Implement the `topp()` function that achieves this goal. When testing the code we will use $p = 0.7$, but your implementation should support arbitrary $p \in [0, 1]$.

1.1.5 Evaluation

Run the evaluation cell. This will use the first 10 prompts of the test set, and generate 10 continuations for each prompt with each of the above decoding methods. Each decoding method will output its overall evaluation metrics: perplexity, fluency, and diversity.

Deliverables:

- Code (40 pts, 8 pts for each decoding algorithm):** Implement code blocks denoted by YOUR CODE HERE: in Section 1 of CSED503_Assignment5.ipynb.
- Write-up (20 pts):** Answer the following questions in your write-up:
 - Q1:** In greedy decoding, what do you observe when generating 10 times from the test prompt? (2 pts)
 - Q2:** In vanilla sampling, what do you observe when generating 10 times from the test prompt? (2 pts)
 - Q3:** In temperature sampling, play around with the value of temperature t . Which value of t makes it equivalent to greedy decoding? Which value of t makes it equivalent to vanilla sampling? (4 pts)

- **Q4:** In top- k sampling, play around with the value of k . Which value of k makes it equivalent to greedy decoding? Which value of k makes it equivalent to vanilla sampling? (4 pts)
- **Q5:** In top- p sampling, play around with the value of p . Which value of p makes it equivalent to greedy decoding? Which value of p makes it equivalent to vanilla sampling? (4 pts)
- **Q6:** Is there a decoding method that wins over all others on all three evaluation metrics? If not, which method strikes the best balance in your opinion? (There is no single correct answer here, any reasoning faithful to your experimental results will receive full credit.) (4 pts)

How do I know if my code is working correctly? It is hard to automate the evaluation of the decoding algorithms due to the issues with reproducibility during sampling. We recommend two ways to check the correctness of your code. First, you can run the evaluation cell and check if you get numbers close to the reference values that we provide for each decoding algorithm. Another way we recommend is to go through the write up questions, think of the answers that you expect for these questions and see if your implementation of the decoding algorithms behaves accordingly. E.g. for Q4, from your understanding of top- k sampling you should be able to guess what value of k makes the algorithm equivalent to greedy decoding. When you choose that value of k , does your implementation return the output which is the same as the output you get when generating using the **greedy** method?

1.2 Knowledge Distillation (40 pts)

In this part of the homework, you will learn how we can use knowledge distillation from a larger teacher model to a smaller student model. Particularly, we will be focusing on the task of text summarization and using the [CNN/Daily Mail dataset](#). We will use Qwen2.5-1.5B-Instruct as our teacher model, which is a 1.5B parameter decoder-only model pre-trained on 18T tokens of data and then further fine-tuned to follow instructions to perform different tasks (similar to something like ChatGPT). You can read more about Qwen2.5 models [here](#). For the student model, we will be using the GPT-2 small model, which is a 124M parameter model.

1.2.1 Background.

As discussed in lectures, knowledge distillation (KD) is the process of transforming large models into smaller ones. Why we might need to do something like this is because for many practical scenarios it might be impossible to serve large models as those will have high latency and inference costs. One of the reasons why high performing language models are so large is because they are supposed to be general purpose models with a wide range of capabilities. However, if for an application at hand, we only need one specific capability of the large model, e.g. summarization, we can use knowledge distillation to specialise a much smaller language model towards that particular task.

KD is not a recent idea and dates back to at least [Hinton 2015](#). While there are many flavors to how to distill knowledge from a large neural network (teacher model) to a smaller network (student), we will focus on the synthetic data approach, which has become very common in the LLM era due to their generative nature. The idea is very simple, we start with a teacher model and use it to generate data for the task which we want the student model to specialise towards. For e.g. if we want to specialise a small model to do better summarization, we will use a large teacher model and generate summaries of a bunch of articles using this model. The generated data is then used to train / fine-tune the smaller student model. It can be useful to filter the synthetic data generated by the teacher model before using it to train a smaller model to get rid of low-quality samples, see – [West et al. 2022](#), [Sclar et al. 2022](#) and [Wang et al. 2023](#). However, for the purposes of this homework we will simply train the student model without any filtering.

1.2.2 Implementing Knowledge Distillation for Text Summarization.

Step 1: Set up Student Model (10 pts)

- `prepare_articles_for_student_model()` (4 pts): [Implement this function](#).

In this function you format and tokenize the data so that it can be used for summarization using the student model i.e. GPT-2. Note that GPT-2 is a language model and inherently a language model's job is to predict continuations of a sequence by predicting one token at a time. To perform specific tasks like summarization using language models, we need to prepare the data in such a format such that the possible continuation of the sequence is the output we want i.e. in this case the summary of the article. The GPT-2 paper found adding a TL;DR to the end of the article helps the model in generating better summaries. Post formatting, you should then tokenize the formatted articles, which means breaking the article into a list of (sub-)words and converting them into token ids corresponding to the indices of words in the language model's vocabulary. Both of these steps can be conveniently done using a single line of code by calling the pre-trained tokenizer from huggingface: `tokenizer()`.

- `summarize_wth_student_model()` (6 pts): [Implement this function](#).

In this function you implement the code for generating summaries using the student model by first formatting and tokenizing the articles by calling the above function and then feeding the tokenized inputs to the student model to generate summaries. We will be using top-p / nucleus sampling for generation. As with the tokenization, generation is also very convenient using the pre-trained models from huggingface and can be done by simply calling `model.generate()`. To use top-p sampling, simply provide the argument `top_p = <p>` to the `generate` method.

Step 2: Set up Teacher Model (10 pts)

- `prepare_articles_teacher()` (4 pts): [Implement this function.](#)

Similar to the student model, we will need to format and tokenize data for the teacher model. Our teacher model i.e. Qwen2.5-1.5B-Instruct is something we call an instruction tuned language model. What it means is that in addition to being trained on next-word prediction on a large text corpora, the model was further fine-tuned to follow instructions for a wide range of problems (e.g. different NLP tasks, chat bot queries like write me an email). Please check [Ouyang et al.](#) if you are interested to learn more about instruction tuning, as instruction tuning has been one of the key ideas that has led to the success of modern LLMs. Coming back to the function implementation, you will need to format your prompt in way that is appropriate for instruction following rather than text completion. We will do this by adding an instruction to the beginning of each article, i.e., “Summarize the following article.” Further, we will also instruct the model to output the summary in a specific format by appending a suffix to the end of each article, i.e., “Start your summary with 'TL;DR:'”. This will help us easily extract the summary from the generated response of the model. We also add something called a *System Prompt* at the beginning of each input, which is useful to ground the model towards a particular role or persona. Like for this problem we use the system prompt: “You are a helpful assistant and an expert at summarizing articles.” (we add the system prompt for you, you don’t need to add that on your own).

- `summarize_with_teacher_model()` (6 pts): [Implement this function.](#)

Similar to `summarize_wth_student_model()`, just uses the teacher model to summarise the articles.

- `generate_synthetic_data_for_distillation()`: [You do NOT need to implement this function.](#) Calls `summarize_with_teacher_model()` with the articles in the training data and generates summaries using the teacher model.

Step 3: Fine-tuning Student Model on Synthetic Summaries (10 pts)

- `prepare_data_for_distillation()` (10 pts): [Implement this function.](#)

This function formats the data in a specific way so that it can be used to fine-tune the student model. You will follow pretty much the same process as you did for the student model in `prepare_articles_for_student_model` with a few changes.

1. First we will include the summaries in the input text along with the articles. This is done because we are now training the student model to generate summaries from the articles. Hence the format of the input text will be `<article>\nTL;DR:<summary>`.
2. In the dictionary returned by the tokenizer, we now need to add a new key, “`labels`”, which contains the labels to train the language model. For language models, the labels are the same as the input IDs since the model is expected to generate the next word in the sequence. However, while fine-tuning, we want the model to learn how to generate the summaries from the articles and we do not care about the model learning to predict tokens in the original articles. Therefore, we replace the labels for the prompt tokens with -100, which is a special token id that is used to signal the loss function to ignore the loss for those tokens.

This process of fine-tuning a language model to generate output text conditioned on an input is commonly referred to as *Supervised Fine-tuning*, which is one of the simplest yet effective forms of instruction tuning.

- `fine_tune_student_model()`: [You do NOT need to implement this function.](#)

This function fine-tunes the student model using the `Trainer` API from huggingface. *Fine-tuning takes roughly 5 minutes on Google Colab T4 GPU.*

Deliverables:

1. **Code (30 pts):** Implement code blocks denoted by `YOUR CODE HERE:` in Section 2 of `CSED503_Assignment5.ipynb`.
2. **Write-up (10 pts): Effectiveness of Synthetic Data in Comparison to Human Data.** Note that the CNN/Daily Mail dataset that we are using does have human written summaries of the articles in the training data. For this exercise, we ask you to train on those summaries and compare the performance with the student model fine-tuned using synthetic data generated from the teacher model. You can use just the first 1000 article, summary pairs from the training data. Report a plot with ROUGE scores comparing GPT-2 (no fine-tuning), GPT-2 (KD fine-tuning), GPT-2 (real-data fine-tuning), and Qwen 1.5B-Instruct. Also, explain the trends in 2-3 lines in the writeup.

You can use the following code to prepare the human annotated training data for fine-tuning:

```
# Select just first 1000 examples
train_data_og = cnn_dm_cse447_dataset["train"].select(range(1000))

train_tokenized_data_og = train_data_og.map(
    lambda example: prepare_data_for_distillation(
        example["article"],
        example["summary"],
        student_tokenizer,
        max_length=1024,
    ),
    batched=False,
    remove_columns=cnn_dm_cse447_dataset["train"].column_names,
)
# You can now run fine_tune_student_model() with train_tokenized_data_og
```