# CSE 256: NLP UCSD, Programming Assignment 4

## Text Decoding From GPT-2 using Beam Search (40 points)

### Due: Friday, Dec 2, 2024

IMPORTANT: After copying this notebook to your Google Drive, paste a link to it below. To get a publicly-accessible link, click the *Share* button at the top right, then click "Get shareable link" and copy the link.

Link: paste your link here:

https://colab.research.google.com/drive/1q_4G9CTdXao0DiFjsk5yqRrcbLfRh72I?usp=sharing

---

**Notes:**

Make sure to save the notebook as you go along.

Submission instructions are located at the bottom of the notebook.

## ⌄ Part 0: Setup

## ⌄ Adding a hardware accelerator

Go to the menu and add a GPU as follows:

```
Edit > Notebook Settings > Hardware accelerator > (GPU)
```

Run the following cell to confirm that the GPU is detected.

```
1 import torch
2
3 # Confirm that the GPU is detected
4 assert torch.cuda.is_available()
5
6 # Get the GPU device name.
7 device_name = torch.cuda.get_device_name()
8 n_gpu = torch.cuda.device_count()
9 print(f"Found device: {device_name}, n_gpu: {n_gpu}")
```

⇥  Found device: Tesla T4, n_gpu: 1

## ⌄ Installing Hugging Face's Transformers and Additional Libraries

We will use Hugging Face's Transformers (https://github.com/huggingface/transformers).

Run the following cell to install Hugging Face's Transformers library and some other useful tools.

```
1 pip install -q sentence-transformers==2.2.2 transformers==4.17.0 matplotlib
```

⇥      ──────────────────────────────────── 86.0/86.0 kB 5.0 MB/s eta 0:00:00
        Preparing metadata (setup.py) ... done
                        ──────────────────── 67.9/67.9 kB 4.8 MB/s eta 0:00:00
              ────────────────────────────── 3.8/3.8 MB 40.7 MB/s eta 0:00:00
              ────────────────────────────── 897.5/897.5 kB 28.1 MB/s eta 0:00:00
        Building wheel for sentence-transformers (setup.py) ... done

## ⌄ Part 1. Beam Search

We are going to explore decoding from a pretrained GPT-2 model using beam search. Run the below cell to set up some beam search utilities.

```
1 from transformers import GPT2LMHeadModel, GPT2Tokenizer
2
3 tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
4 model = GPT2LMHeadModel.from_pretrained("gpt2", pad_token_id=tokenizer.eos_token_id)
5
6 # Beam Search
7
```

```python
 8 def init_beam_search(model, input_ids, num_beams):
 9     assert len(input_ids.shape) == 2
10     beam_scores = torch.zeros(num_beams, dtype=torch.float32, device=model.device)
11     beam_scores[1:] = -1e9 # Break ties in first round.
12     new_input_ids = input_ids.repeat_interleave(num_beams, dim=0).to(model.device)
13     return new_input_ids, beam_scores
14
15
16 def run_beam_search_(model, tokenizer, input_text, num_beams=5, num_decode_steps=10, score_processors=[], to_cpu=Tru
17
18     input_ids = tokenizer.encode(input_text, return_tensors='pt')
19
20     input_ids, beam_scores = init_beam_search(model, input_ids, num_beams)
21
22     token_scores = beam_scores.clone().view(num_beams, 1)
23
24     model_kwargs = {}
25     for i in range(num_decode_steps):
26         model_inputs = model.prepare_inputs_for_generation(input_ids, **model_kwargs)
27         outputs = model(**model_inputs, return_dict=True)
28         next_token_logits = outputs.logits[:, -1, :]
29         vocab_size = next_token_logits.shape[-1]
30         this_token_scores = torch.log_softmax(next_token_logits, -1)
31
32         # Process token scores.
33         processed_token_scores = this_token_scores
34         for processor in score_processors:
35             processed_token_scores = processor(input_ids, processed_token_scores)
36
37         # Update beam scores.
38         next_token_scores = processed_token_scores + beam_scores.unsqueeze(-1)
39
40         # Reshape for beam-search.
41         next_token_scores = next_token_scores.view(num_beams * vocab_size)
42
43         # Find top-scoring beams.
44         next_token_scores, next_tokens = torch.topk(
45             next_token_scores, num_beams, dim=0, largest=True, sorted=True
46         )
47
48         # Transform tokens since we reshaped earlier.
49         next_indices = torch.div(next_tokens, vocab_size, rounding_mode="floor") # This is equivalent to `next_toker
50         next_tokens = next_tokens % vocab_size
51
52         # Update tokens.
53         input_ids = torch.cat([input_ids[next_indices, :], next_tokens.unsqueeze(-1)], dim=-1)
54
55         # Update beam scores.
56         beam_scores = next_token_scores
57
58         # Update token scores.
59
60         # UNCOMMENT: To use original scores instead.
61         # token_scores = torch.cat([token_scores[next_indices, :], this_token_scores[next_indices, next_tokens].unsc
62         token_scores = torch.cat([token_scores[next_indices, :], processed_token_scores[next_indices, next_tokens].u
63
64         # Update hidden state.
65         model_kwargs = model._update_model_kwargs_for_generation(outputs, model_kwargs, is_encoder_decoder=False)
66         model_kwargs["past"] = model._reorder_cache(model_kwargs["past"], next_indices)
67
68     def transfer(x):
69       return x.cpu() if to_cpu else x
70
71     return {
72         "output_ids": transfer(input_ids),
73         "beam_scores": transfer(beam_scores),
74         "token_scores": transfer(token_scores)
75     }
76
77
78 def run_beam_search(*args, **kwargs):
79     with torch.inference_mode():
80         return run_beam_search_(*args, **kwargs)
81
```

```
 82
 83 # Add support for colored printing and plotting.
 84
 85 from rich import print as rich_print
 86
 87 import numpy as np
 88
 89 import matplotlib
 90 from matplotlib import pyplot as plt
 91 from matplotlib import cm
 92
 93 RICH_x = np.linspace(0.0, 1.0, 50)
 94 RICH_rgb = (matplotlib.colormaps.get_cmap(plt.get_cmap('RdYlBu')))(RICH_x)[:, :3] * 255).astype(np.int32)[range(5, 45
 95
 96
 97 def print_with_probs(words, probs, prefix=None):
 98   def fmt(x, p, is_first=False):
 99     ix = int(p * RICH_rgb.shape[0])
100     r, g, b = RICH_rgb[ix]
101     if is_first:
102       return f'[bold rgb(0,0,0) on rgb({r},{g},{b})]{x}'
103     else:
104       return f'[bold rgb(0,0,0) on rgb({r},{g},{b})] {x}'
105   output = []
106   if prefix is not None:
107     output.append(prefix)
108   for i, (x, p) in enumerate(zip(words, probs)):
109     output.append(fmt(x, p, is_first=i == 0))
110   rich_print(''.join(output))
111
112 # DEMO
113
114 # Show range of colors.
115
116 for i in range(RICH_rgb.shape[0]):
117   r, g, b = RICH_rgb[i]
118   rich_print(f'[bold rgb(0,0,0) on rgb({r},{g},{b})]hello world rgb({r},{g},{b})')
119
120 # Example with words and probabilities.
121
122 words = ['the', 'brown', 'fox']
123 probs = [0.14, 0.83, 0.5]
124 print_with_probs(words, probs)
```

```
Downloading: 100%                                        0.99M/0.99M [00:00<00:00, 6.00MB/s]

Downloading: 100%                                        446k/446k [00:00<00:00, 588kB/s]

Downloading: 100%                                        26.0/26.0 [00:00<00:00, 648B/s]

Downloading: 100%                                        665/665 [00:00<00:00, 34.9kB/s]

Downloading: 100%                                        523M/523M [00:10<00:00, 70.0MB/s]

/usr/local/lib/python3.10/dist-packages/transformers/modeling_utils.py:1439: FutureWarning: You are using `torch.load` with `weights
  state_dict = torch.load(resolved_archive_file, map_location="cpu")
hello world rgb(215,49,39)
hello world rgb(244,111,68)
hello world rgb(253,176,99)
hello world rgb(254,226,147)
hello world rgb(251,253,196)
hello world rgb(217,239,246)
hello world rgb(163,210,229)
hello world rgb(108,164,204)
the brown fox
```

## Question 1.1 (5 points)

Run the cell below. It produces a sequence of tokens using beam search and the provided prefix.

```
1 num_beams = 5
2 num_decode_steps = 10
3 input_text = 'The brown fox jumps'
4
5 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=num_decode_steps)
6 for i, tokens in enumerate(beam_output['output_ids']):
```

```
7       score = beam_output['beam_scores'][i]
8       print(i, round(score.item() / tokens.shape[-1], 3), tokenizer.decode(tokens, skip_special_tokens=True))
```

```
0 -1.106 The brown fox jumps out of the fox's mouth, and the fox
1 -1.168 The brown fox jumps out of the fox's cage, and the fox
2 -1.182 The brown fox jumps out of the fox's mouth and starts to run
3 -1.192 The brown fox jumps out of the fox's mouth and begins to lick
4 -1.199 The brown fox jumps out of the fox's mouth and begins to bite
```

To get you more acquainted with the code, let's do a simple exercise first. Write your own code in the cell below to generate 3 tokens with a beam size of 4, and then print out the **third most probable** output sequence found during the search. Use the same prefix as above.

```
1 input_text = 'The brown fox jumps'
2
3 # WRITE YOUR CODE HERE!
4 num_beams = 4
5 num_decode_steps = 3
6
7 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=num_decode_steps)
8 print(tokenizer.decode(beam_output['output_ids'][2], skip_special_tokens=True))
```

```
The brown fox jumps up and down
```
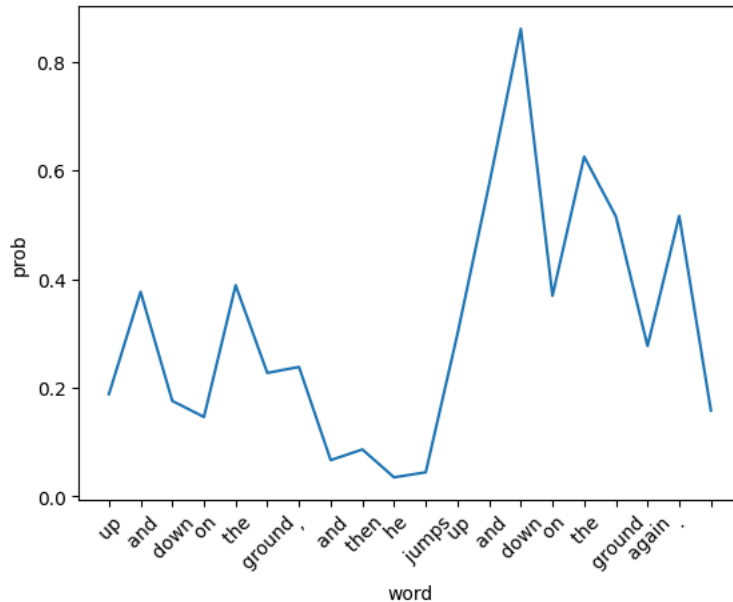
## ⌄ Question 1.2 (5 points)

Run the cell below to visualize the probabilities the model assigns for each generated word when using beam search with beam size 1 (i.e., greedy decoding).

```
1 input_text = 'The brown fox jumps'
2 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=1, num_decode_steps=20)
3 probs = beam_output['token_scores'][0, 1:].exp()
4 output_subwords = [tokenizer.decode(tok, skip_special_tokens=True) for tok in beam_output['output_ids'][0]]
5
6 print('Visualizeation with plot:')
7
8 fig, ax = plt.subplots()
9 plt.plot(range(len(probs)), probs)
10 ax.set_xticks(range(len(probs)))
11 ax.set_xticklabels(output_subwords[-len(probs):], rotation = 45)
12 plt.xlabel('word')
13 plt.ylabel('prob')
14 plt.show()
15
16 print('Visualization with colored text (red for lower probability, and blue for higher):')
17
18 print_with_probs(output_subwords[-len(probs):], probs, ' '.join(output_subwords[:-len(probs)]))
```

Visualizeation with plot:



Visualization with colored text (red for lower probability, and blue for higher):
The brown fox jumps **up and down on the ground , and then he jumps up and down on the ground again .**

Why does the model assign higher probability to tokens generated later than to tokens generated earlier?
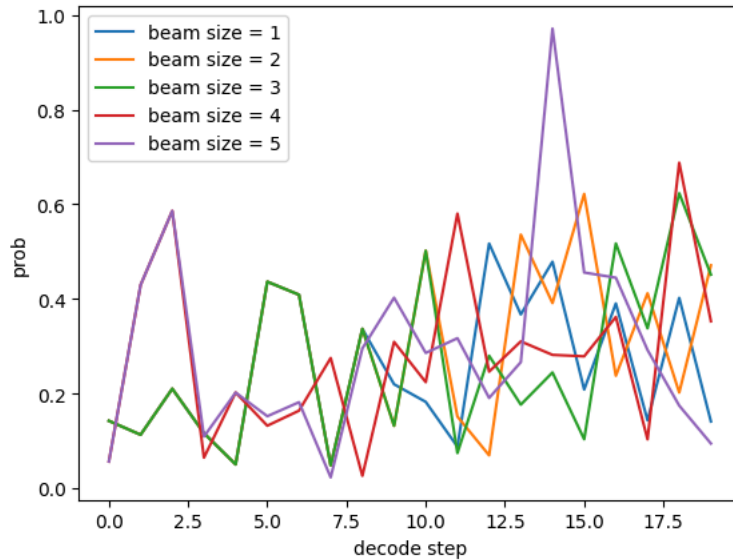
∨  Write your answer here

(1) Later tokens have more contextual information to inform their probability distribution, potentially making them more "confident".

(2) Later tokens are conditioned on an increasingly specific and refined context, which can lead to more focused and higher-probability predictions.

(3) The attention mechanism allows later tokens to have a more comprehensive view of the entire generated sequence, which wasn't possible for earlier tokens when they were first generated.

Run the cell below to visualize the word probabilities when using different beam sizes.

```
 1 input_text = 'Once upon a time, in a barn near a farm house,'
 2 num_decode_steps = 20
 3 model.cuda()
 4
 5 beam_size_list = [1, 2, 3, 4, 5]
 6 output_list = []
 7 probs_list = []
 8 for bm in beam_size_list:
 9   beam_output = run_beam_search(model, tokenizer, input_text, num_beams=bm, num_decode_steps=num_decode_steps)
10   output_list.append(beam_output)
11   probs = beam_output['token_scores'][0, 1:].exp()
12   probs_list.append((bm, probs))
13
14 print('Visualization with plot:')
15 fig, ax = plt.subplots()
16 for bm, probs in probs_list:
17   plt.plot(range(len(probs)), probs, label=f'beam size = {bm}')
18 plt.xlabel('decode step')
19 plt.ylabel('prob')
20 plt.legend(loc='best')
21 plt.show()
22
23 print('Model predictions:')
24 for bm, beam_output in zip(beam_size_list, output_list):
25   tokens = beam_output['output_ids'][0]
26   print(bm, beam_output['beam_scores'][0].item() / tokens.shape[-1], tokenizer.decode(tokens, skip_special_tokens=Tr
```

Visualization with plot:



## Question 1.3 (10 points)

Beam search often results in repetition in the predicted tokens. In the following cell we pass a score processor called `WordBlock` to `run_beam_search`. At each time step, it reduces the probability for any previously seen word so that it is not generated again.

Run the cell to see how the output of beam search changes with and without using `WordBlock`.

```
1 import collections
2
3 class WordBlock:
4     def __call__(self, input_ids, scores):
5         for batch_idx in range(input_ids.shape[0]):
6             for x in input_ids[batch_idx].tolist():
7                 scores[batch_idx, x] = -1e9
8         return scores
9
10 input_text = 'Once upon a time, in a barn near a farm house,'
11 num_beams = 1
12
13 print('Beam Search')
14 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=40, score_processo
15 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
16
17 print('Beam Search w/ Word Block')
18 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=40, score_processo
19 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
```

Model predictions:
1 -0.9706197796445905 Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He was playing with a sti
2 -0.9286185177889738 Once upon a time, in a barn near a farm house, a young boy was playing with a stick. The boy was playing with
3 -0.9597569667931759 Once upon a time, in a barn near a farm house, a young boy was playing with a stick. The boy, who had been pla
4 -0.9205132108746152 Once upon a time, in a barn near a farm house, there was a young girl who had been brought up by her mother. S
5 -0.9058780670166016 Once upon a time, in a barn near a farm house, there was a man who had been living in the house for a long tim
Beam Search
Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He was playing with a stick, and the boy was pl
Beam Search w/ Word Block
Once upon a time, in a barn near a farm house, the young girl was playing with her father's dog. She had been told that she would be

Is `WordBlock` a practical way to prevent repetition in beam search? What (if anything) could go wrong when using `WordBlock`?

**Write your answer here**

WordBlock is not a practical way to prevent repetition for following reasons:

(1) By completely blocking any previously seen token, the method can severely limit the model's generation capabilities, tokens such as he, she, it are required to appear multiple times in the sentence to make it more clear.

(2) This method does not distinguish between: meaningful repetitions, contextually appropriate repetitions or grammatically necessary repetitions.

(3) Forcibly removing tokens can lead to lower quality outputs and increased generation of less relevant alternative tokens.

(4) For each time step, the method iterates through all previously generated tokens, which can introduce computational complexity, especially for long sequences.

## ⌄ Question 1.4 (20 points)

Use the previous `WordBlock` example to write a new score processor called `BeamBlock`. Instead of uni-grams, your implementation should prevent tri-grams from appearing more than once in the sequence.

Note: This technique is called "beam blocking" and is described [here](#) (section 2.5). Also, for this assignment you do not need to re-normalize your output distribution after masking values, although typically re-normalization is done.

Write your code in the indicated section in the below cell.

```python
1 import collections
2
3 class BeamBlock:
4     def __call__(self, input_ids, scores):
5         batch_size, _ = input_ids.shape
6         vocab_size = scores.shape[-1]
7
8         for batch_idx in range(batch_size):
9             tokens = input_ids[batch_idx].tolist()
10            used_trigrams = set()
11            for i in range(len(tokens) - 2):
12                trigram = tuple(tokens[i:i + 3])
13                used_trigrams.add(trigram)
14
15            if len(tokens) >= 2:
16                last_two = tuple(tokens[-2:])
17                for next_token in range(vocab_size):
18                    potential_trigram = last_two + (next_token, )
19                    if potential_trigram in used_trigrams:
20                        scores[batch_idx, next_token] = -1e9
21        return scores
22
23 input_text = 'Once upon a time, in a barn near a farm house,'
24 num_beams = 1
25
26 print('Beam Search')
27 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=40, score_processc
28 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
29
30 print('Beam Search w/ Beam Block')
31 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_decode_steps=40, score_processc
32 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
```

```
Beam Search
Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He was playing with a stick, and the boy was pl
Beam Search w/ Beam Block
Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He was playing on the stick, and the boy was tr
```

## Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.
2. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
3. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.