

Decoding Strategies

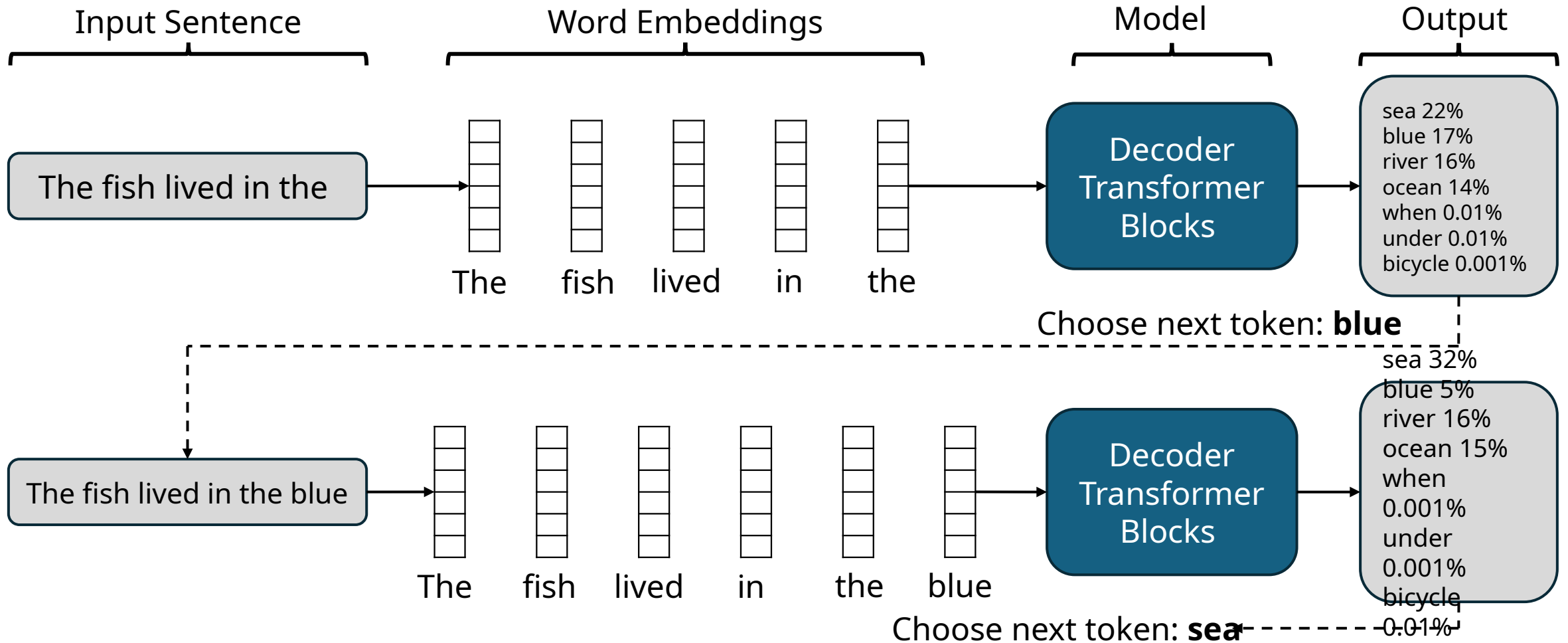
CS XXX: Introduction to Large Language Models

Contents

- Generating Text
- Greedy Decoding
- Beam Search Decoding
- Top k Sampling
- Top p (Nucleus) Sampling
- Generation Parameters

Generating Text

- Generate next token based on the current context



Generating Text

- More concretely, Once language model trained we generate the text in the following way
 - Sample a word in the output from the softmax distribution that results from using the beginning of sentence marker, <s>, as the first input.
 - Use the word embedding for that first word as the input to the network at the next time step, and then sample the next word in the same fashion.
 - Continue generating until the end of sentence marker, </s>, is sampled or a fixed length limit is reached.
- The goal of sequence generation is to produce a sequence of tokens that maximizes the overall probability of the complete sequence given the context.

Generating Text

- The goal of sequence generation is to produce a sequence of tokens that maximizes the overall probability of the complete sequence given the context.
- Consider the following distribution given the context “<s>The mouse”:

$$P(<s>, \text{The, mouse, ate, the, cheese, } </s>) = 0.02,$$

$$P(<s>, \text{The, cheese, ate, cheese, } </s>) = 0.01,$$

$$P(<s>, \text{mouse, the, the, cheese, ate, } </s>) = 0.0001.$$

⋮

then the generation of the model should result in the sequence with the highest probability

<s>The mouse ate the cheese</s>

Generating Text

- How should the model pick the next word at each time step?
- Consider the distribution after the word “ate” was generated

$$P(\text{cheese} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.30$$

$$P(\text{icecream} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.18$$

$$P(\text{the} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.25$$

$$P(\text{ate} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.12$$

$$P(\text{mouse} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.01$$

$$P(\text{The} \mid \langle s \rangle, \text{The, mouse, ate}) = 0.05$$

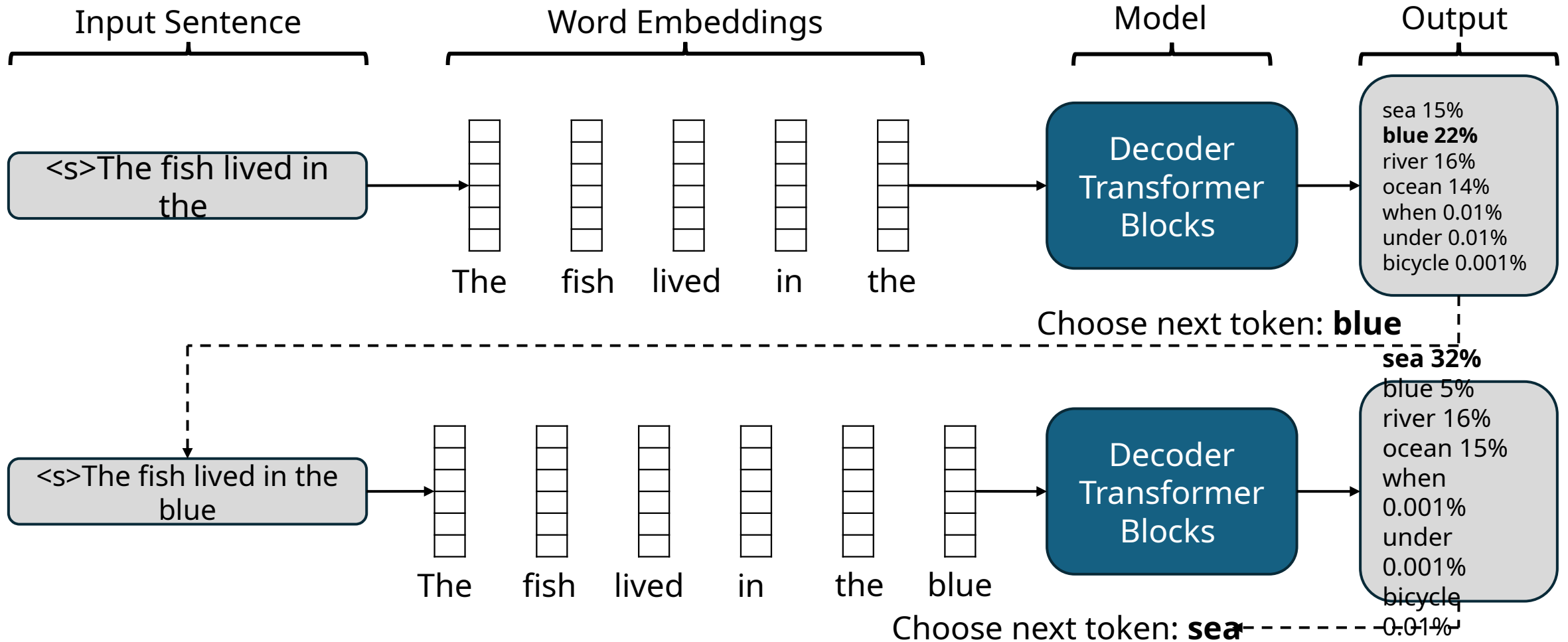
The next word generated should be “the” instead of “cheese” even though “cheese has a higher probability than “the” since predicting “the” will lead to the model generating the highest probability sequence “ $\langle s \rangle$ The mouse ate **the** cheese $\langle /s \rangle$ ”

Decoding Strategies

- Decoding strategies in Large Language Models (LLMs) refer to different methods used to generate text from a model's probability distribution over the vocabulary.
- These strategies determine how the next word (or token) is selected and can significantly impact the quality, coherence, and diversity of generated text.

Greedy Decoding

- **Greedy Decoding:** Choose the token with the greatest probability at each time step.



Greedy Decoding

- **Greedy Decoding:** Choose the token with the greatest probability at each time step.
- We terminate generation once the end of sequence token (`</s>`) is predicted or when the number of tokens predicted is equal to some max length.

Greedy Decoding

- **Greedy Decoding:** Choose the token with the greatest probability at each time step.
- We terminate generation once the end of sequence token (`</s>`) is predicted or when the number of tokens predicted is equal to some max length.
- **Problem:** Choosing the token with the highest probability may be locally optimal but not globally since optimal sequence generation requires considering the entire sequence probability, not just picking the most likely word at each step. Since greedy decoding does not explore multiple options, it often produces repetitive text. Greedy decoding always produces the same output for a given input (lack of diversity).

Beam Search

- Keep track of n most probable partial translations at each decoder step instead of just one!
- the beam size n is usually 5-10

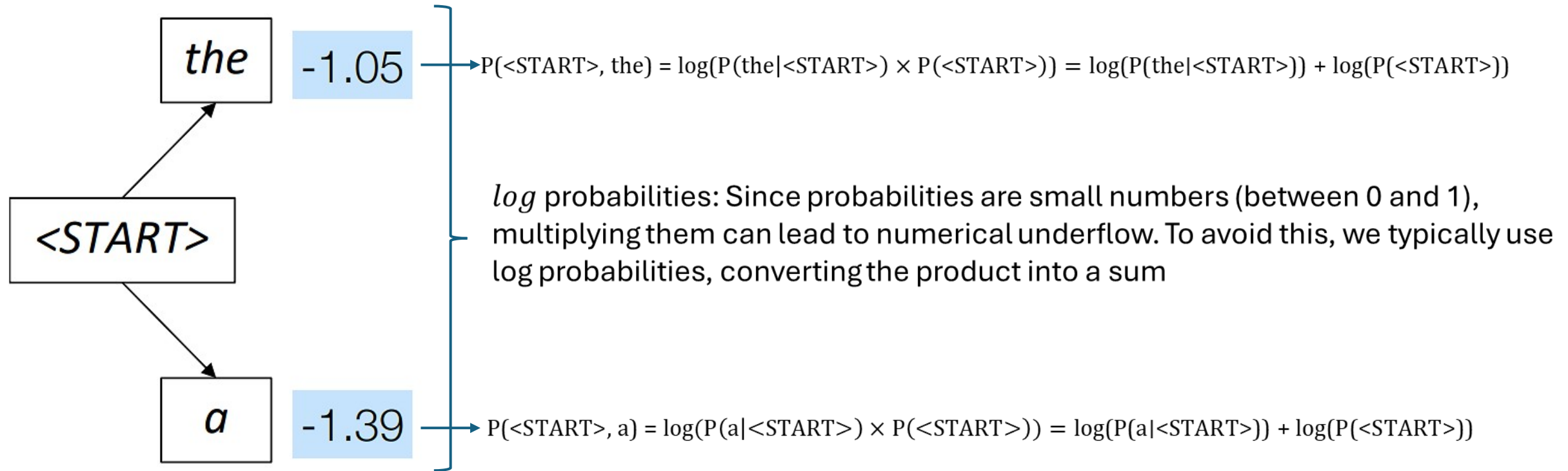
Beam Search

- Beam Size = 2

<START>

Beam Search

- Beam Size = 2
- Consider only the 2 (beam size) most likely words at the first time step



Note that $P(\langle \text{START} \rangle) = 1$, $\log(1) = 0$

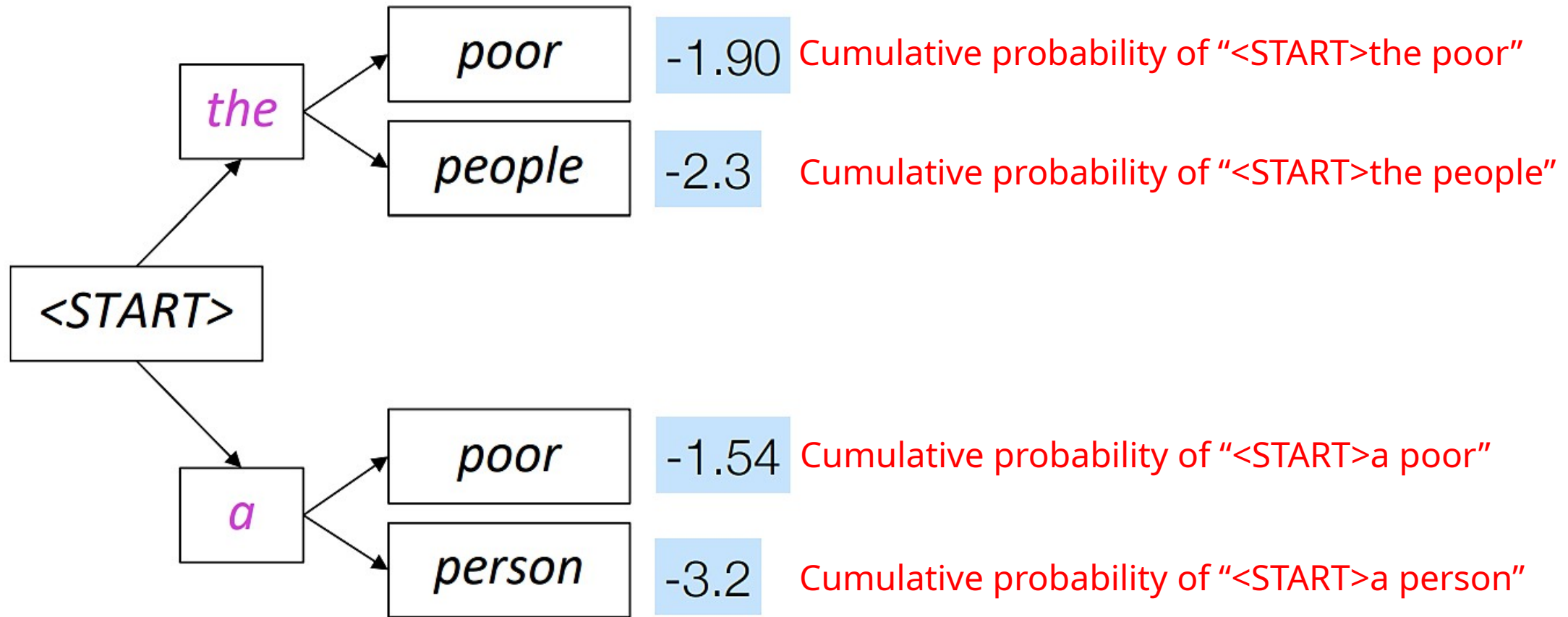
Beam Search

- Beam Size = 2
- Consider only the 2 (beam size) most likely words at the first time step



Beam Search

- Beam Size = 2



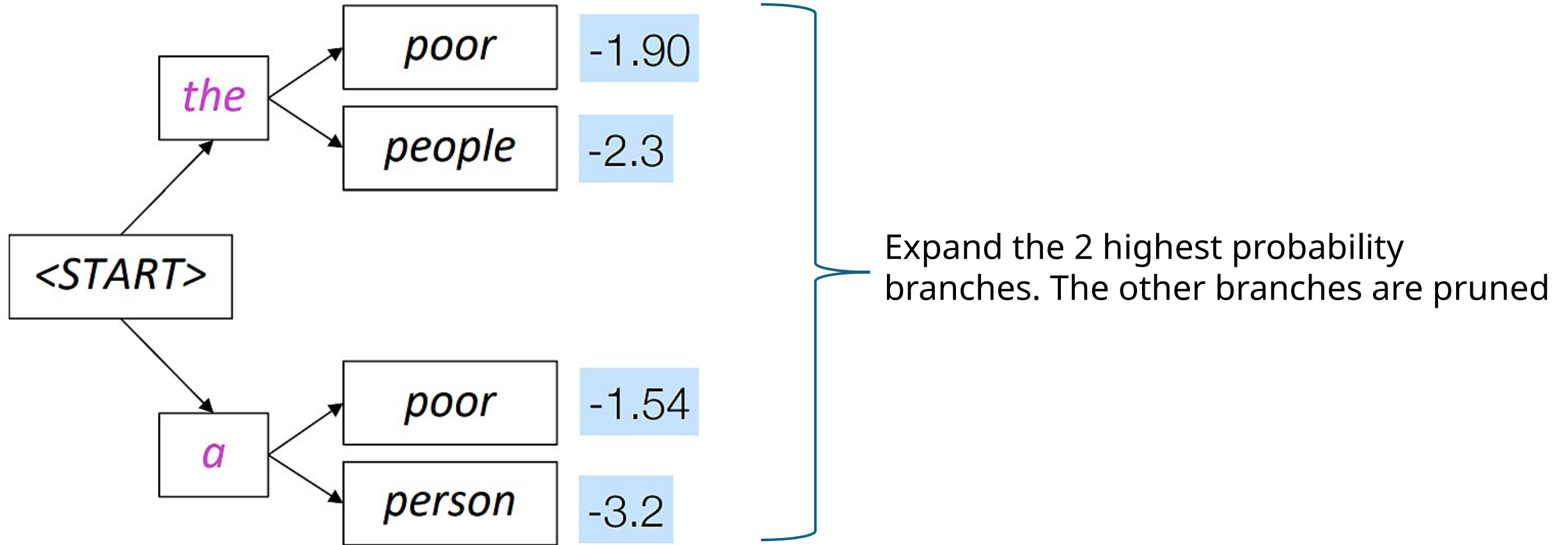
Beam Search

- Beam Size = 2



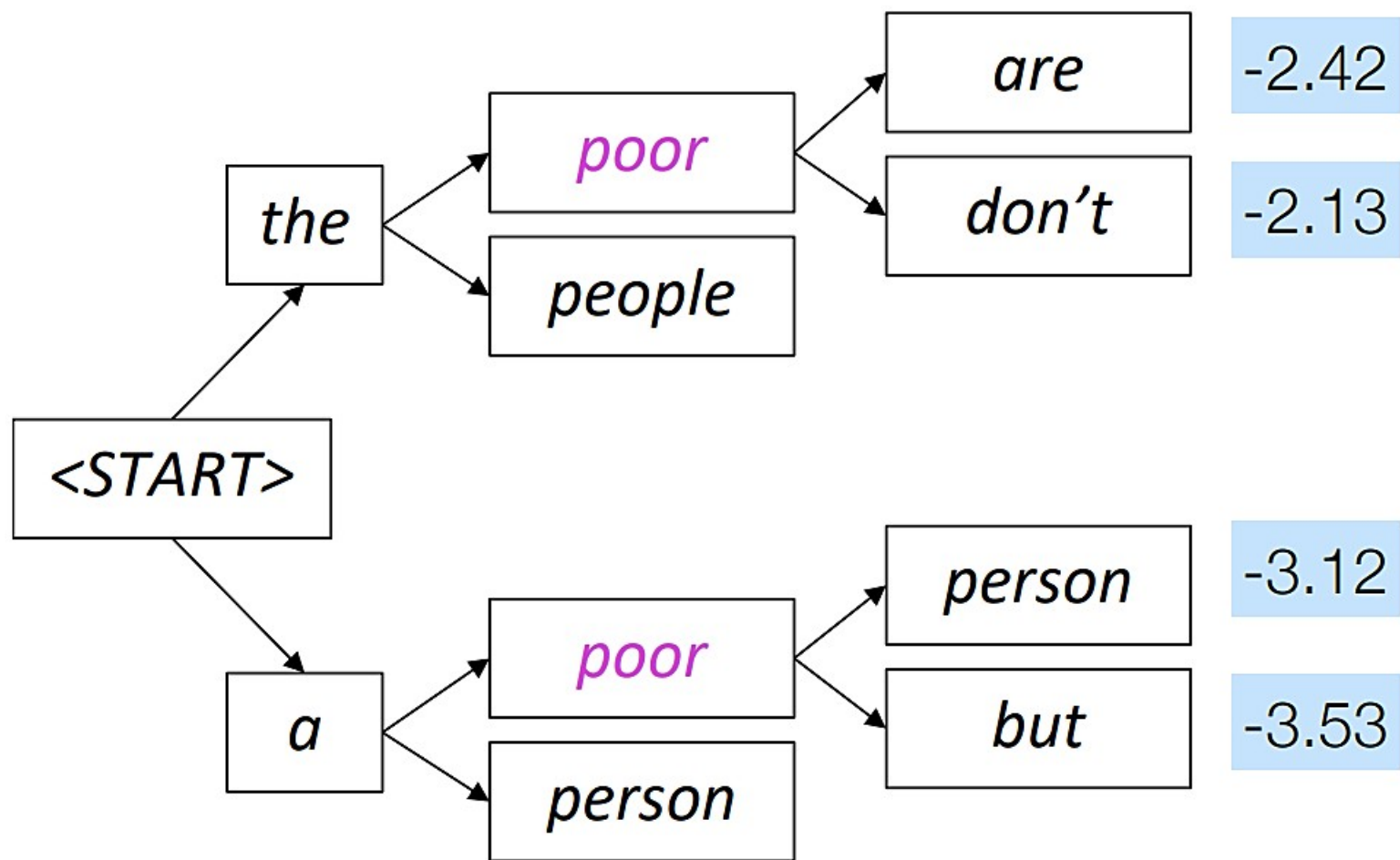
Beam Search

- Beam Size = 2

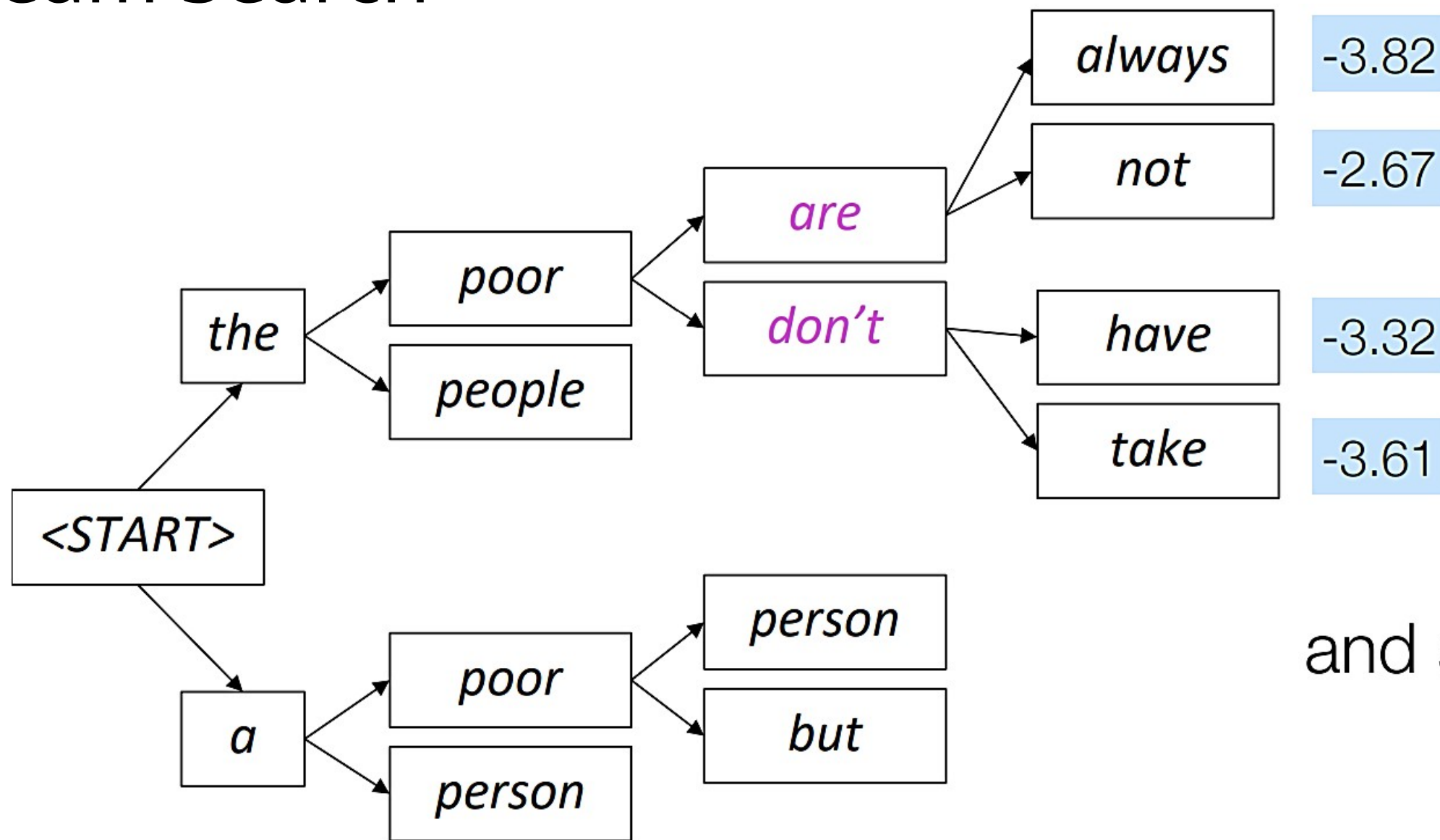


Beam Search

- Beam Size = 2

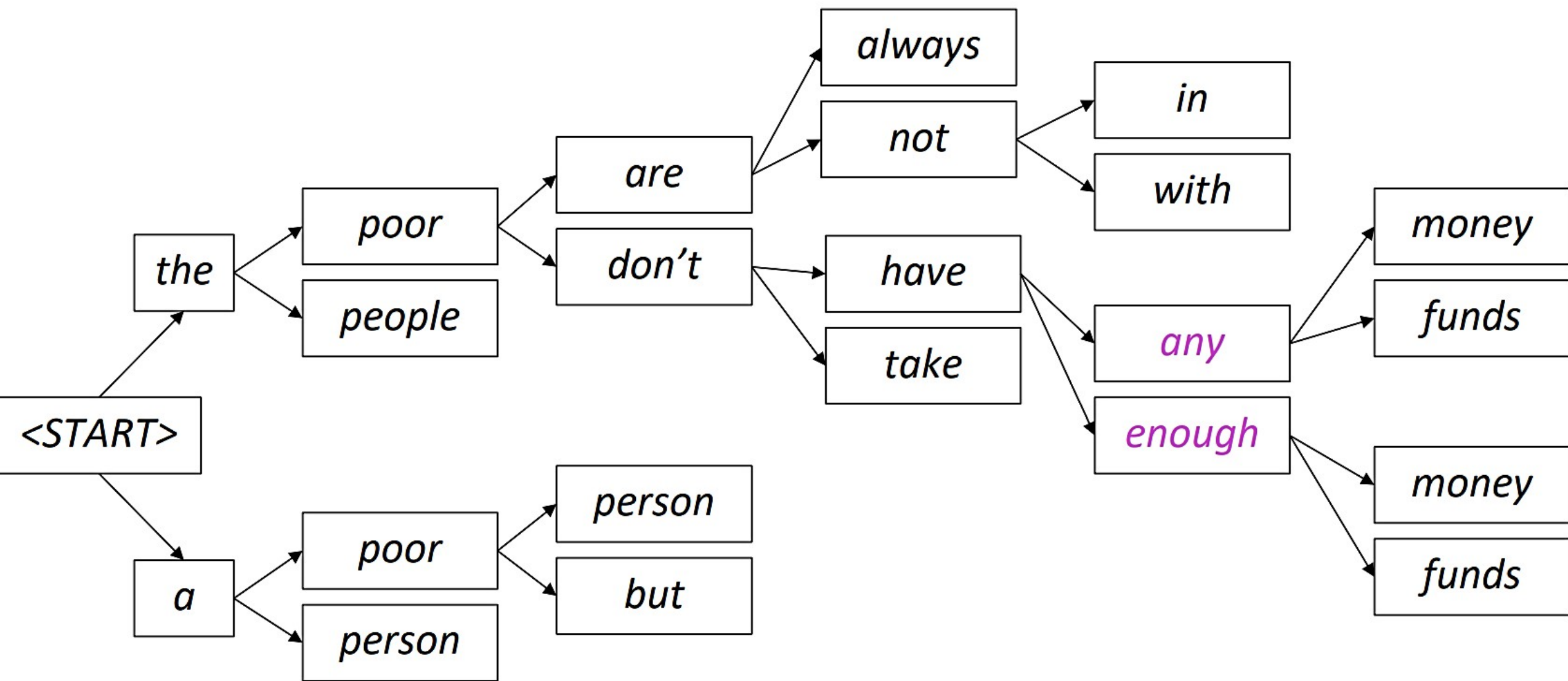


Beam Search



and so on...

Beam Search

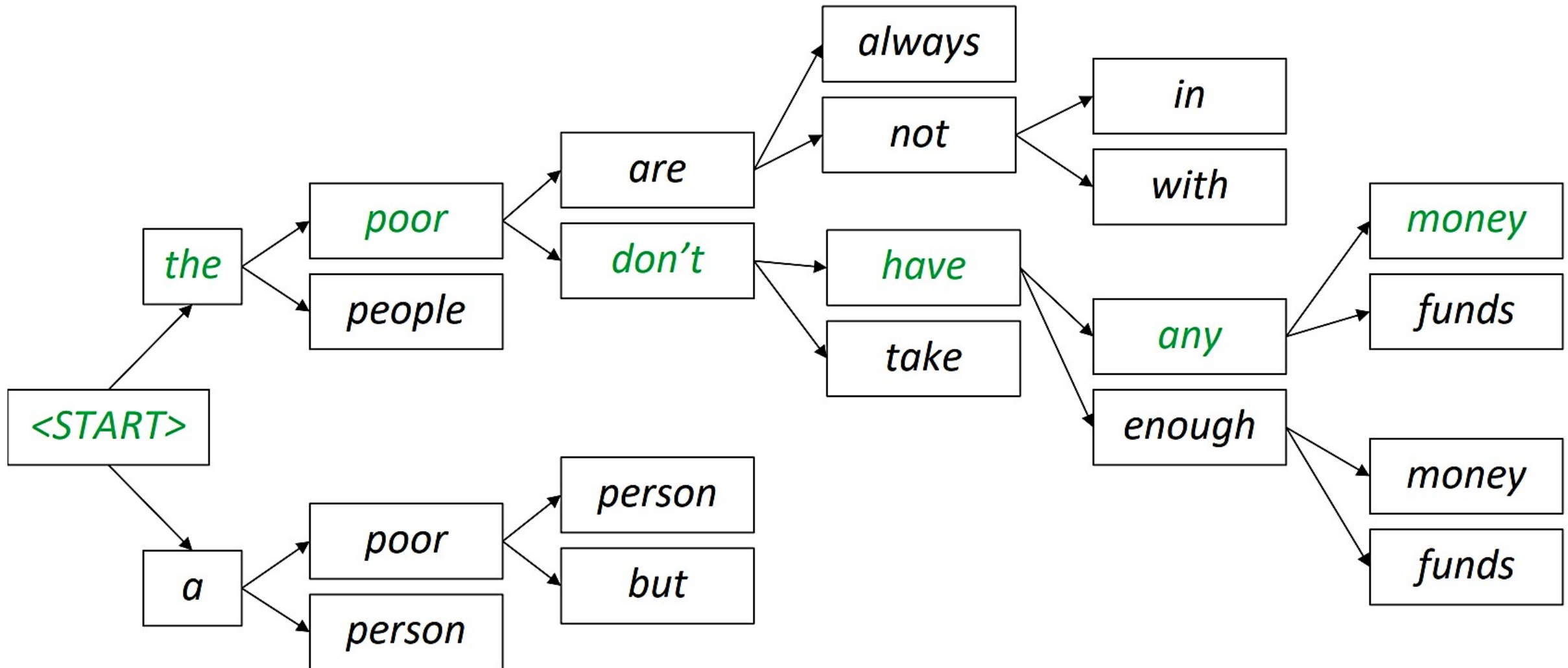


Beam Search

- Generation stops when all N (beam size) branches have either predicted the end-of-sequence (EOS e.g. </s>) token or reached the maximum allowed length.
- Once generation terminates, the sequence with the highest cumulative probability is selected as the final output.

Beam Search

- Once generation terminates, the sequence with the highest cumulative probability is selected as the final output.

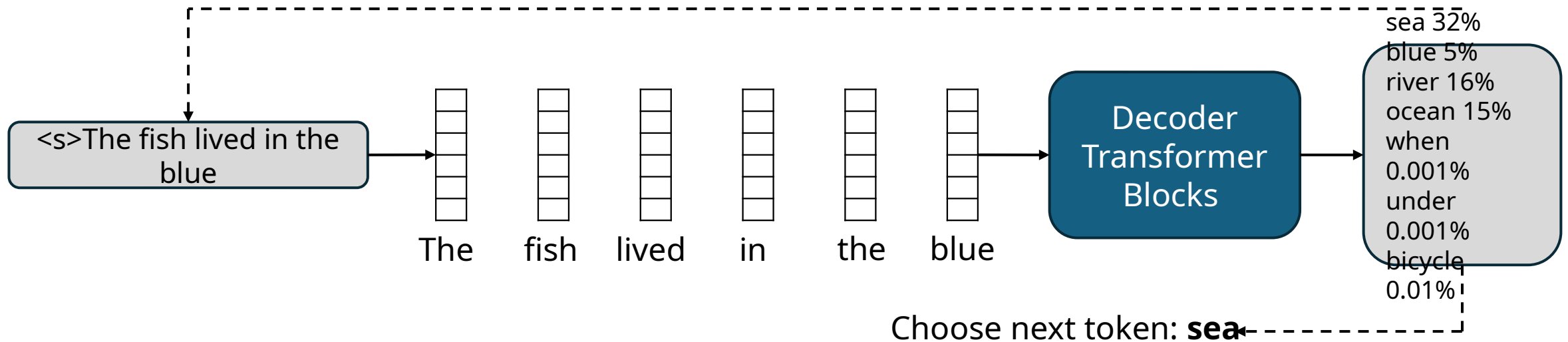


Beam Search

- Effect of beam size N
 - Small N means consider less hypotheses. The resulting generated sequence is likely to be ungrammatical, unnatural, nonsensical, or incorrect.
 - Large N means consider more hypotheses. The resulting generated sequence is likely to be more grammatical and natural, however larger N is more computationally expensive. A larger N may also make the generation more generic and less relevant.

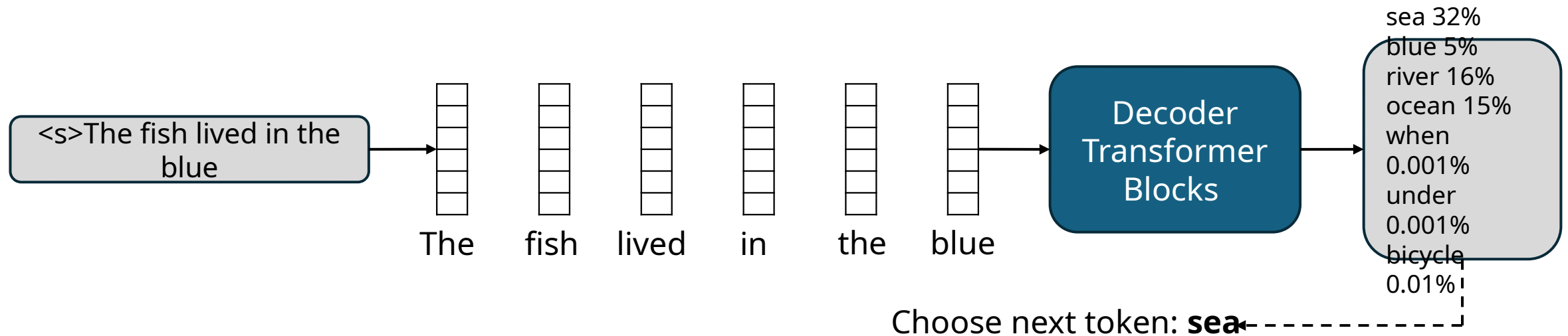
Top k Sampling

- On each step t , randomly sample from P_t , restricted to just the top-k most probable words
- Consider the probability distribution over the vocabulary in the following example



Top K Sampling

- Consider the probability distribution over the vocabulary in the following example



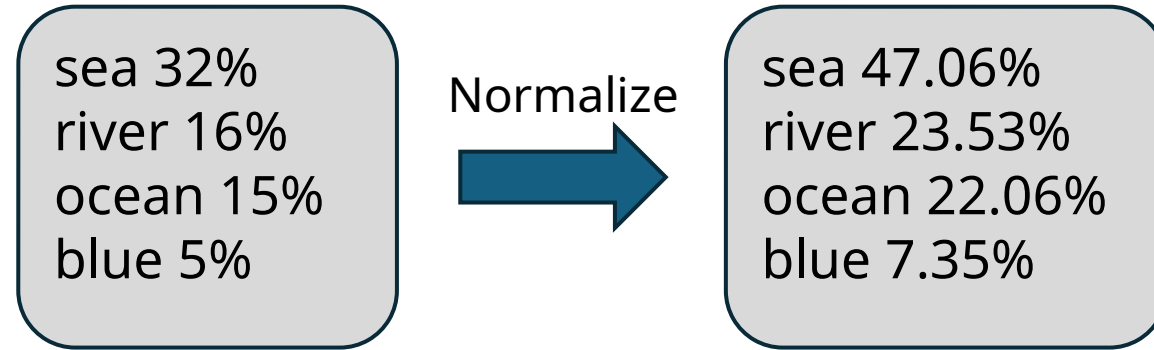
- If $k = 4$ the distribution is truncated to

sea 32%
river 16%
ocean 15%
blue 5%

Note that this truncation is
done at each time step

Top K Sampling

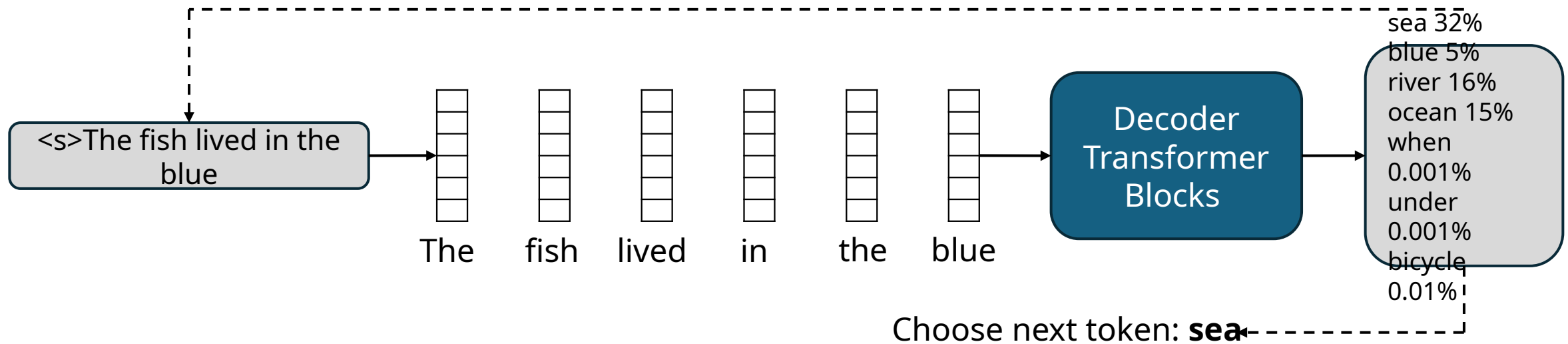
- If $k = 4$ we truncate the distribution to



- Now we randomly sample from the distribution. That is, 47.06% of the time we choose “sea” to be the next token, 23.53% of the time we choose “river” to be the next token, 22.06% of the time we choose “ocean” to be the next token, and 7.35% of the time we choose “blue” to be the next token.

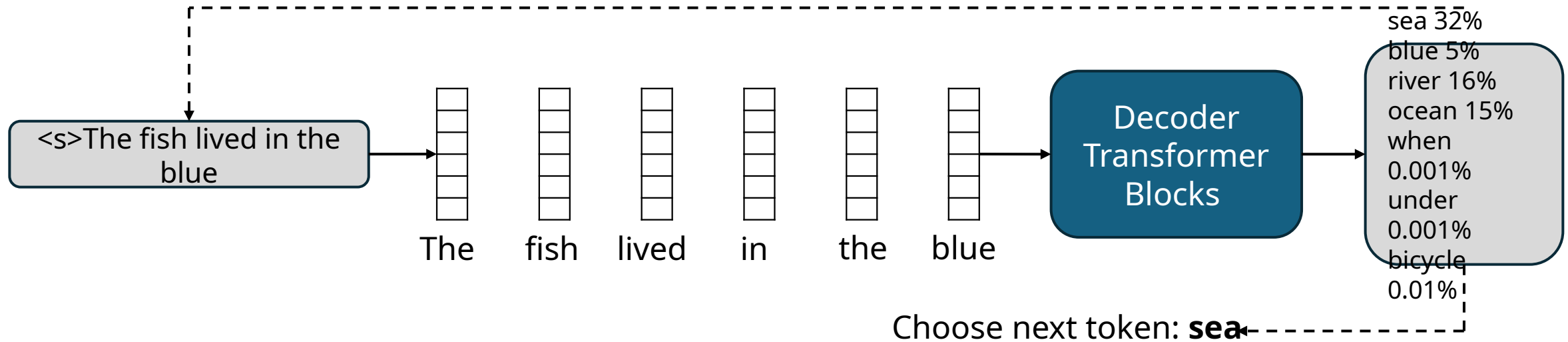
Top p (Nucleus) Sampling

- Truncate the distribution P_t , at every time step t , to the smallest set of words for which the cumulative probability just exceeds p
- Consider the probability distribution over the vocabulary in the following example



Top p (Nucleus) Sampling

- Consider the probability distribution over the vocabulary in the following example

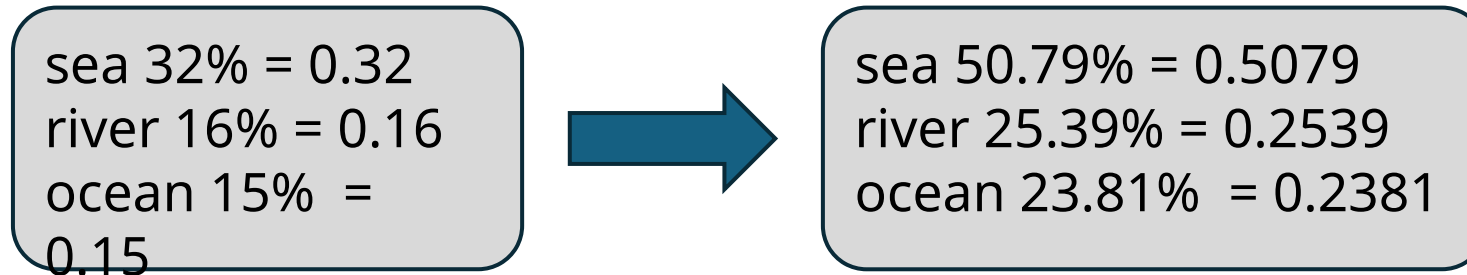


- If $p = 0.6$ the distribution is truncated to

sea 32% = 0.32
river 16% = 0.16
ocean 15% =
0.15

Note that this truncation is
done at each time step

Top p (Nucleus) Sampling



- Now we randomly sample from the distribution. That is, 50.79% of the time we choose “sea” to be the next token, 25.39% of the time we choose “river” to be the next token, and 23.81% of the time we choose “ocean” to be the next token.

Generation parameters

```
0 from transformers import pipeline
1
3 # Initialize the pipeline with a real model
4 gen = pipeline("text-generation", model="gpt2") # Using GPT-2 as an example
5 model
6
7 # Define generation parameters
8 generation_params = {
9     "max_length": 50,      # Maximum length of the output
10    "top_p": 0.9,          # Nucleus sampling (top_p)
11    "top_k": 50,           # Top-k sampling
12    "num_beams": 5,        # Beam search (set to 1 for greedy decoding)
13    "early_stopping": True, # Stop when at least `num_beams` sentences are
14    finished
15    "temperature": 0.7     # Controls randomness (lower = more deterministic)
16 }
17
18 # Use the pipeline
19 text = "The fish lived in"
20 output = gen(text, **generation_params)
21
22 # Print the output
23 print(output)
```

References

- Iyyer, M. (Instructor). (2022). Advanced Natural Language Processing (CS 685) UMass Amherst