# CS846
# Machine Learning for Software Engineering
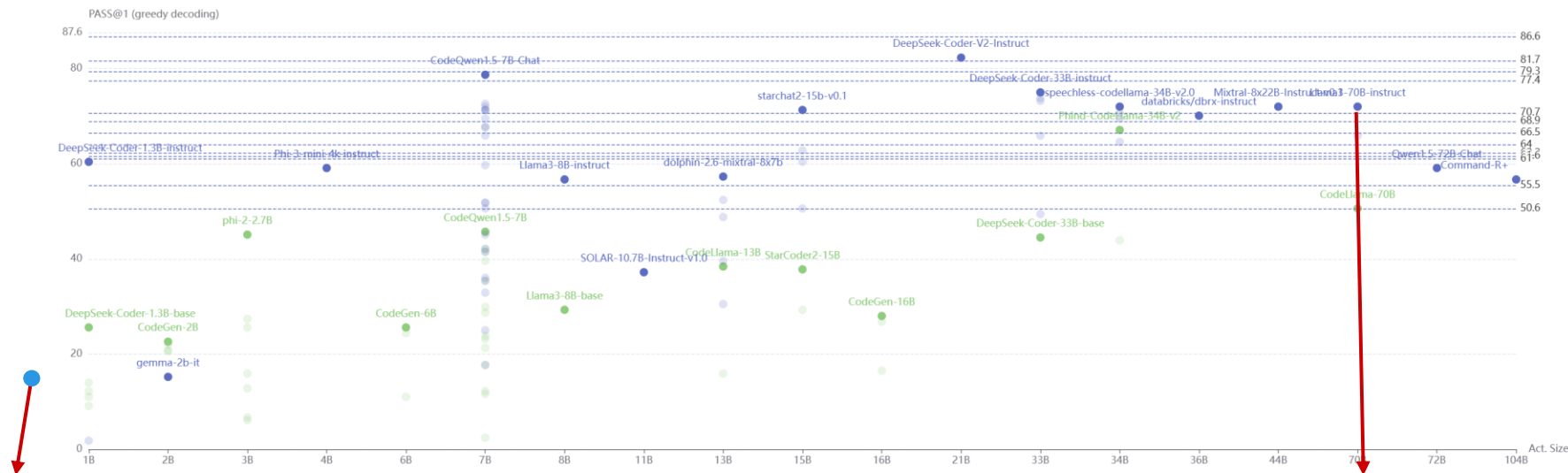
Pengyu Nie

# Large Language Models for Code

Training: pre-training, post-training

Inference: greedy/sampling, zero-shot/few-shot

# Large Language Models

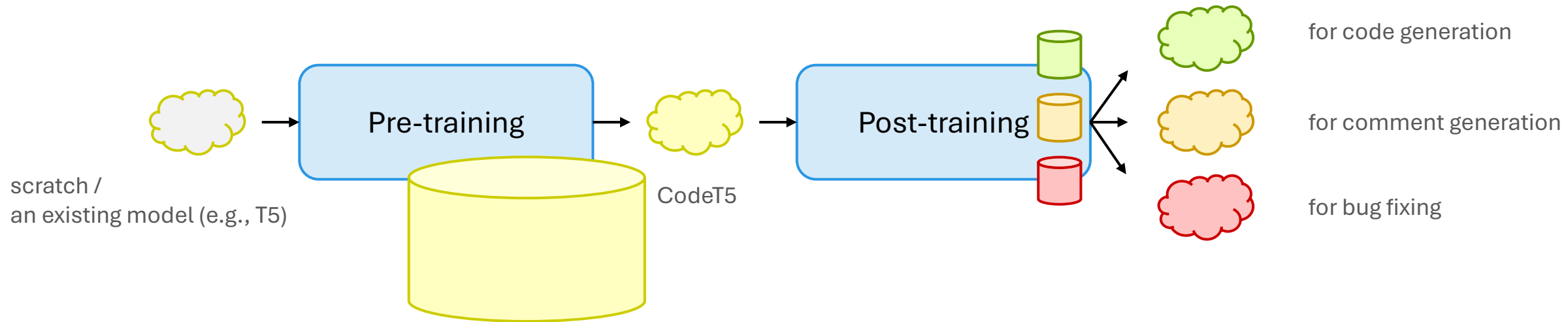- Large transformers trained with massive amount of data



CodeT5, 2021-09-02
- model size: 220M (2.2e8)
- training #tokens: ~4.1B (4.1e9)
- training cost: 12d on 16 x A100 GPUs
- context window: 512
- vocab size: 32K
- #layers: 12
- model dimension: 768
- attention head:

Llama3-70B-instruct, 2024-04-18
- model size: 70B (7e10)
- training #tokens: 15T (1.5e13)
- training cost: 54d on 16,384 x H100 GPUs
- context window: up to 128K
- vocab size: 128K
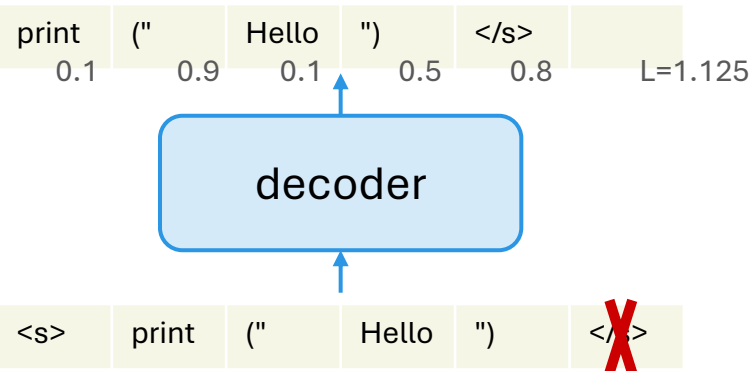- #layers: 80
- model dimension: 8K
- attention head: 64

# Training Overview



scratch /
an existing model (e.g., T5)

Pre-training

CodeT5

Post-training

for code generation

for comment generation

for bug fixing

# Pre-Training

- Massive dataset, self-supervised task(s)

next-token prediction
(aka casual language modeling)



| print | (" | Hello | ") | </s> | |
|-------|------|-------|------|------|--------|
| 0.1 | 0.9 | 0.1 | 0.5 | 0.8 | L=1.125 |

decoder

| <s> | print | (" | Hello | ") | </s> |
|-----|-------|-----|-------|-----|------|

$$L_{CE} = -\frac{1}{n}\sum \log P(y_i)$$

cross-entropy loss
stochastic gradient descent

(aka masked language modeling)



(a) Masked Span Prediction

(c) Masked Identifier Prediction

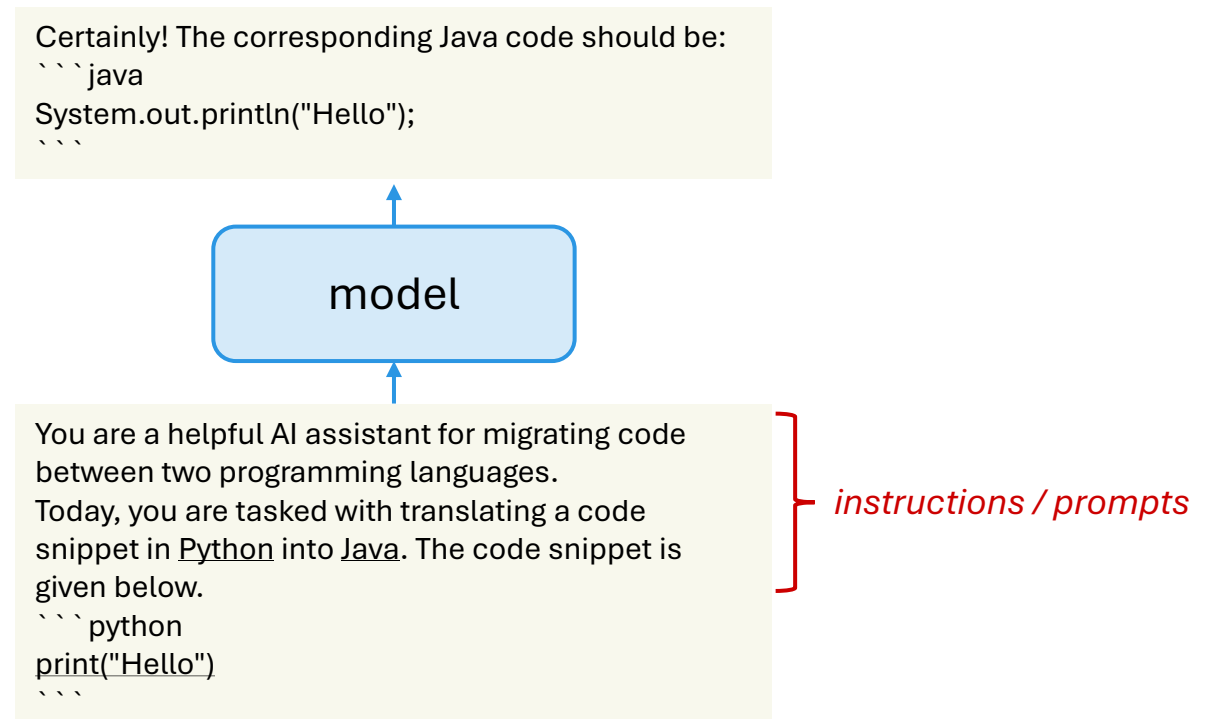(b) Identifier Tagging

(d) Bimodal Dual Generation

# Post-Training / Supervised Finetuning

- Smaller dataset, usually labelled by human
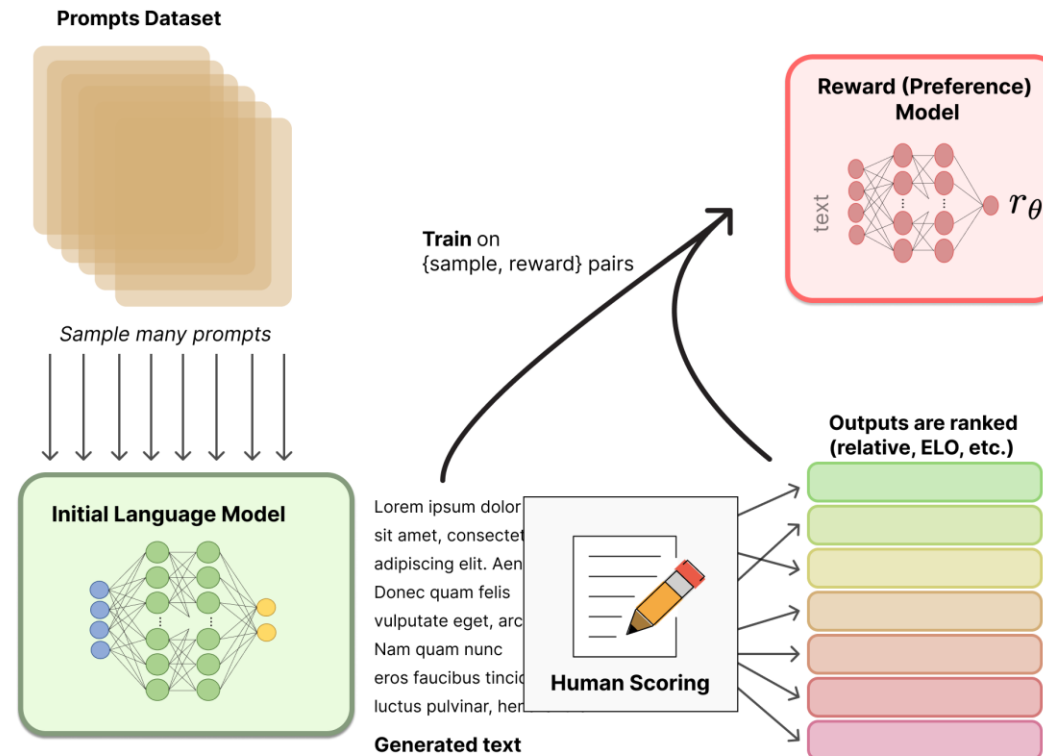- Cross-entropy loss + stochastic gradient descent (same as pre-training)



$$L_{CE} = -\frac{1}{n}\sum \log P(y_i)$$

Variant: instruction finetuning

Certainly! The corresponding Java code should be:
```java
System.out.println("Hello");
```



You are a helpful AI assistant for migrating code between two programming languages.
Today, you are tasked with translating a code snippet in Python into Java. The code snippet is given below.
```python
print("Hello")
```

*instructions / prompts*
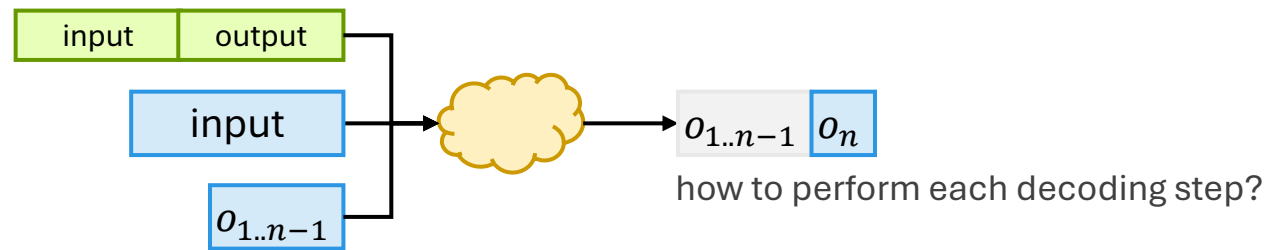
# Post-Training / Reinforcement Learning with Human Feedback

- For improving instruction-following capabilities "alignment" with human preferences
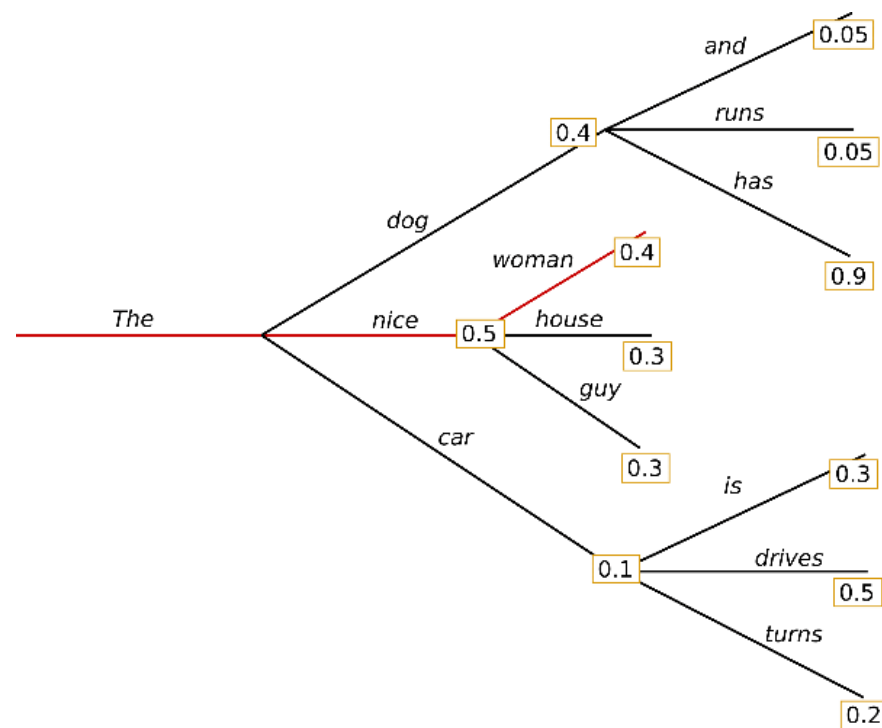
# Inference Overview



input → ☁ → output

provide input-output examples?

| input | output |
|---|---|

input

$o_{1..n-1}$

→ ☁ → $o_{1..n-1}$ $o_n$
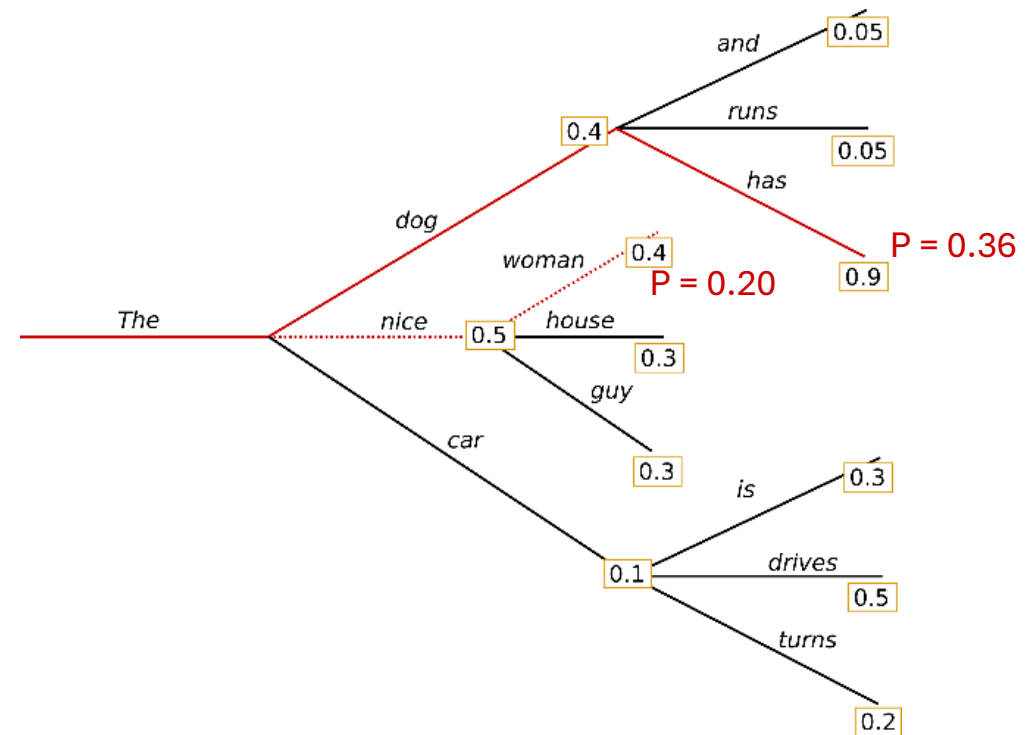
how to perform each decoding step?

# Greedy Decoding

- Choose the token with highest probability at each step
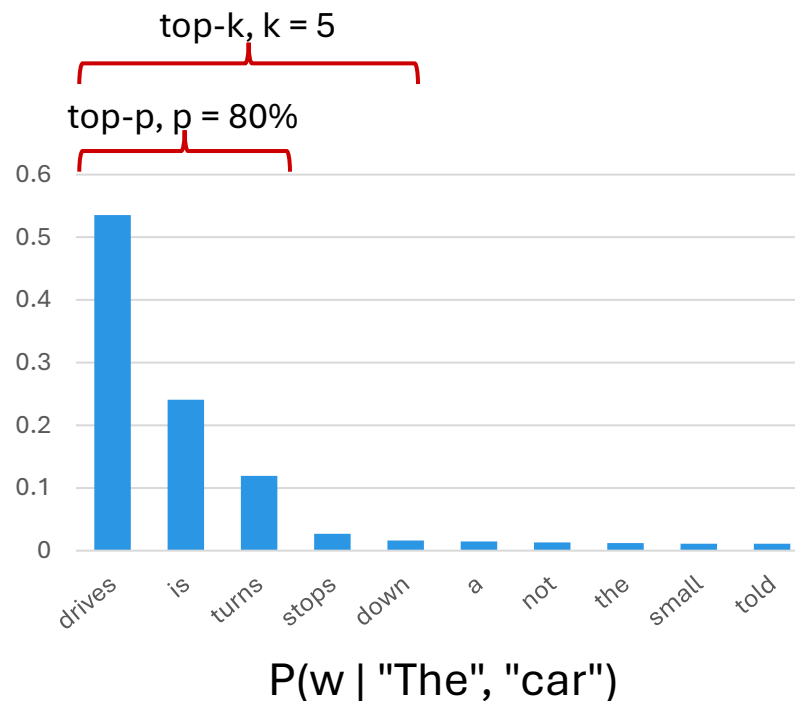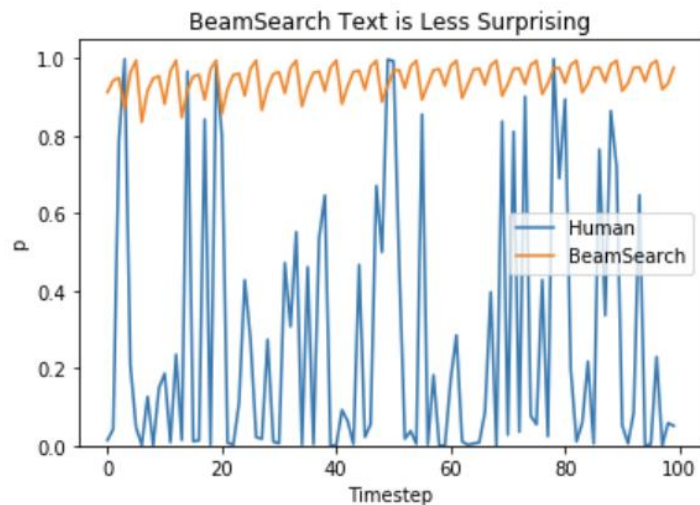- Local maximum
- Deterministic

# Beam Search Decoding

- Search for highest probability sequences,
  with keeping top-k most likely "hypotheses" at each time step

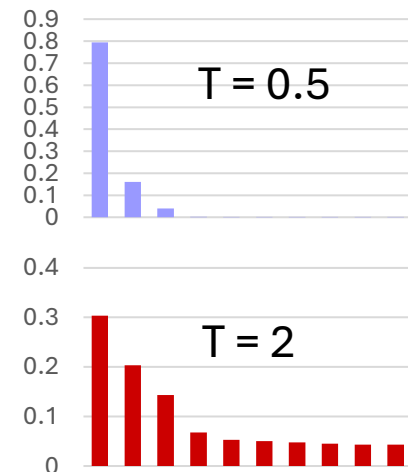- Closer to global maximum

- Deterministic

# Top-K / Top-P Sampling Decoding

- Beam search generated text/code does not look like human written ones, because they are "boring"

- Sample next token according to probability – constrained by top-k / top-p

- Random – randomness controlled by temperature

$$p_i = \text{softmax}\left(\frac{o_i}{T}\right) = \frac{\exp\frac{o_i}{T}}{\sum \exp\frac{o_i}{T}}$$



P(w | "The", "car")

# Zero-Shot, Few-Shot

- Zero-shot learning:
  ask the model to do something unseen during training or inference

- Few-shot learning (aka in-context learning):
  give the model a few input-output examples during inference

You are a helpful AI assistant for migrating code between two programming languages.
Today, you are tasked with translating a code snippet in Python into Java. The code snippet is given below.
```python
print("Hello")
```

Please output only the code in the target programming language and nothing else.

Certainly! The corresponding Java code should be:
```java
System.out.println("Hello");
```

You are a helpful AI assistant for migrating code between two programming languages.
Today, you are tasked with translating a code snippet in Python into Java. Please output only the code in the target programming language and nothing else.
For example:
Input:
```python
l = [1, 2, 3]
```

Output:
```java
List<Integer> l = List.of(1, 2, 3);
```

*few-shot examples / demos*

Now it is your turn
Input:
```python
print("Hello")
```

Output:

```java
System.out.println("Hello");
```