# Self-Consistency Pattern

Instructor

Dipanjan Sarkar

Head of Community & Principal AI Scientist at Analytics Vidhya
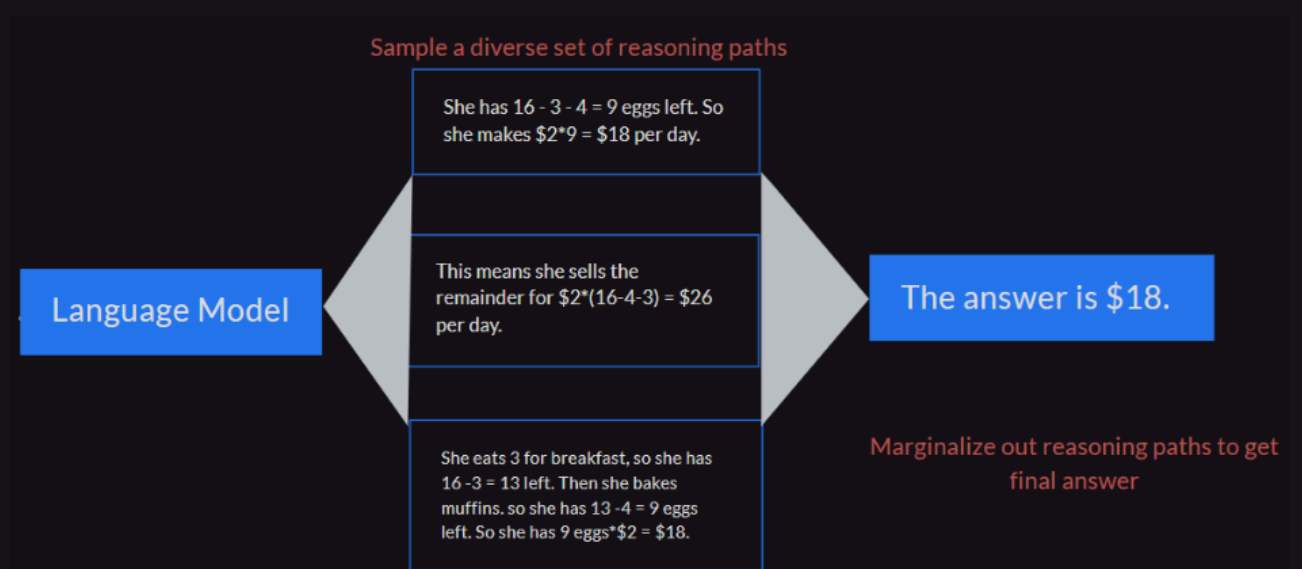
Google Developer Expert - ML & Cloud Champion Innovator
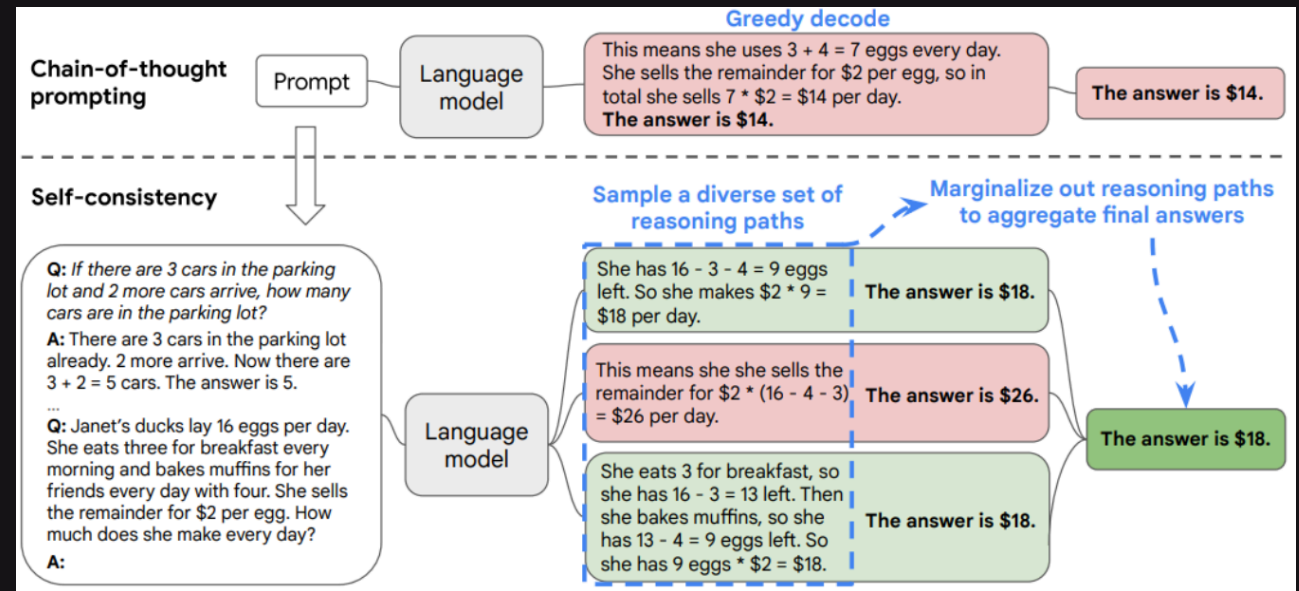
Published Author

# Self-Consistency Pattern

- Self-consistency prompting is a method that helps language models solve complex problems by considering multiple ways to think about a problem and then choosing the most common solution.

- Idea is to not rely on LLM greedy decoding which always selects the response tokens with the highest probability

- Instead, we make the LLM potentially explore diverse reasoning paths by sampling from potential response tokens when generating the response

Sample a diverse set of reasoning paths

She has 16 - 3 - 4 = 9 eggs left. So she makes $2*9 = $18 per day.

This means she sells the remainder for $2*(16-4-3) = $26 per day.

She eats 3 for breakfast, so she has 16 -3 = 13 left. Then she bakes muffins. so she has 13 -4 = 9 eggs left. So she has 9 eggs*$2 = $18.

Language Model

The answer is $18.

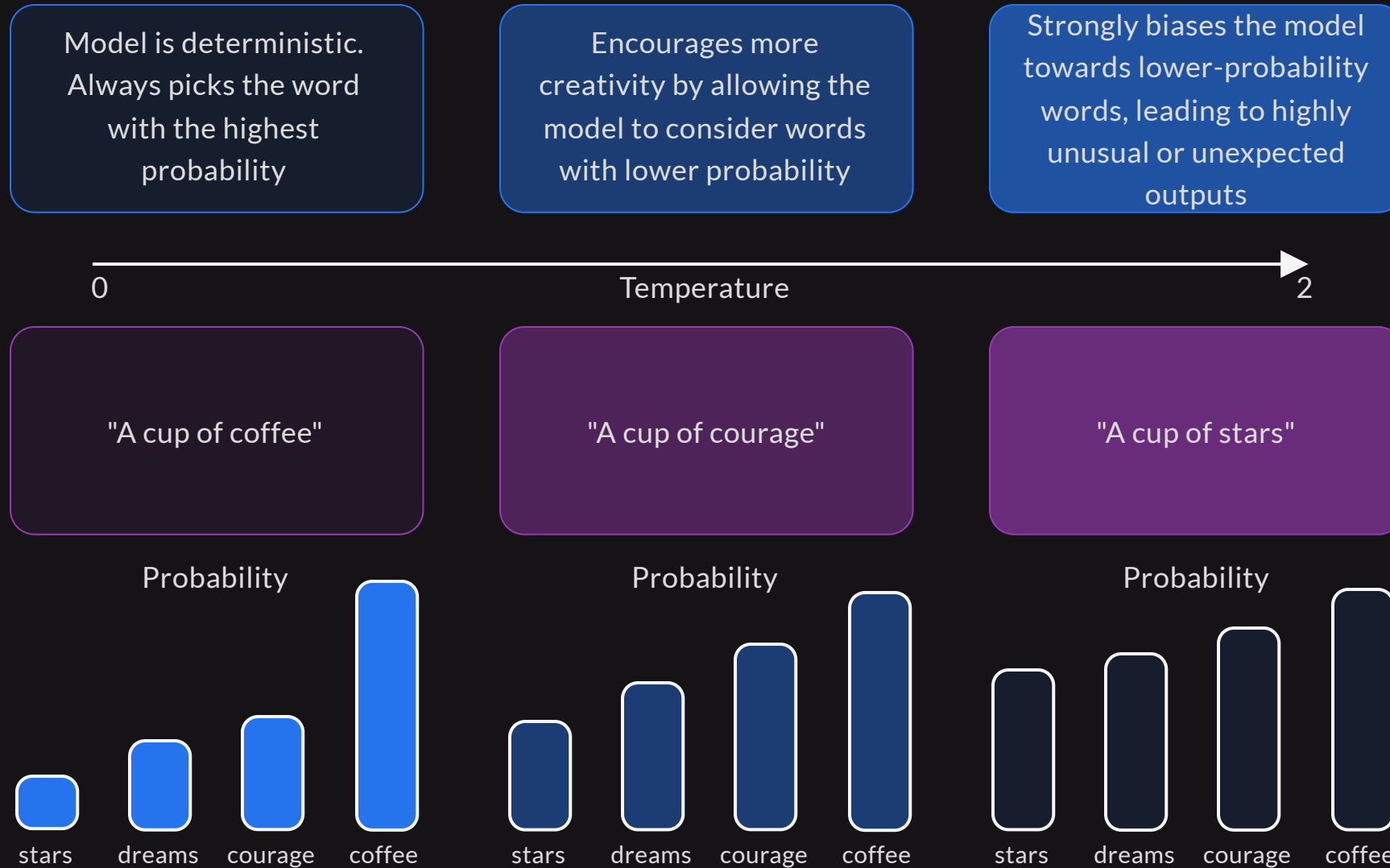Marginalize out reasoning paths to get final answer

# Self-Consistency Pattern

- Chain of Thought uses LLM greedy decoding to select the best response

- In Self-Consistency, we can make the LLM sample tokens and generate multiple variations of CoT responses to explore diverse reasoning paths and then take the majority vote to get to the answer

- Published in the paper, *Self-Consistency Improves Chain of Thought Reasoning in Language Models - March, 2023*

# How to Sample LLM Generation Tokens to get Different Responses?

Model is deterministic. Always picks the word with the highest probability

Encourages more creativity by allowing the model to consider words with lower probability

Strongly biases the model towards lower-probability words, leading to highly unusual or unexpected outputs

0           Temperature          2

"A cup of coffee"

"A cup of courage"

"A cup of stars"

Probability

Probability

Probability

stars   dreams   courage   coffee

stars   dreams   courage   coffee

stars   dreams   courage   coffee

# How to Sample LLM Generation Tokens to get Different Responses?

- In Greedy Decoding, LLMs generate the response tokens by always selecting the next probable token with the highest probability

- Setting the temperature of an LLM is a useful method to encourage the model to generate responses by including next probable tokens with lower probabilities.

- Other settings can also be leveraged like top_k or top_p to sample tokens randomly from a selection of top next tokens



Model is deterministic. Always picks the word with the highest probability

Encourages more creativity by allowing the model to consider words with lower probability

Strongly biases the model towards lower-probability words, leading to highly unusual or unexpected outputs

0     Temperature     2

"A cup of coffee"

"A cup of courage"

"A cup of stars"

Probability

Probability

Probability

stars   dreams   courage   coffee

stars   dreams   courage   coffee

stars   dreams   courage   coffee

Analytics Vidhya

# Thank You