# Explainability in LLMs

# TABLE OF CONTENTS

# 01
Motivation

# Motivation

Users need to understand model decisions, especially in high-stakes domains (e.g., healthcare, legal, finance).

Traditional ML models (e.g., trees, logistic regression) have established explainability methods, but LLMs require different approaches.

LLMs are rapidly complementing and even replacing traditional models, making it crucial to understand how they reach decisions.

Prompting strategies like CoT and ReAct enhance both performance and explainability.

Explainability builds trust and ensures reliable outcomes. Without clear reasoning behind LLM outputs, users may struggle to trust the model's decisions, particularly in sensitive contexts.

Poor explainability can lead to poor performance. LLMs often perform better when prompted in ways that reveal their reasoning process, improving accuracy and robustness.

# 02

Literature Review

# OUR TEAM

**SHAP**

Feature importance, explaining individual predictions, and model debugging
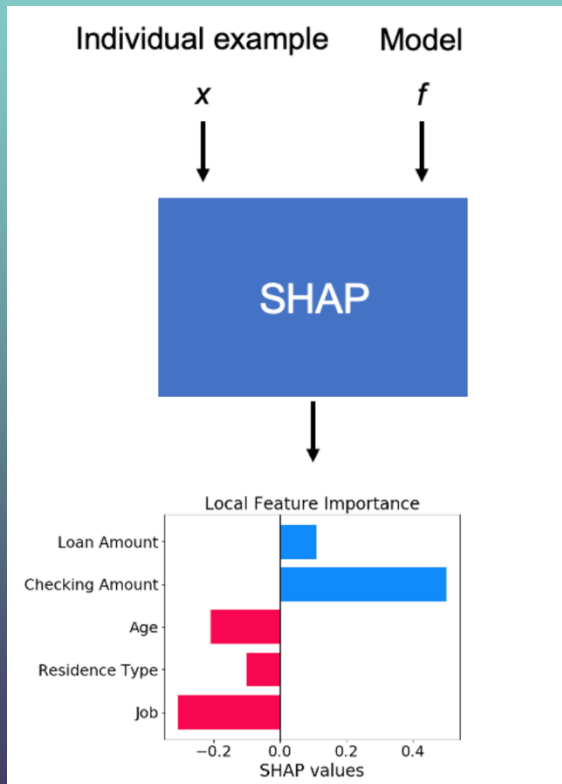
**LIME**

Intuitive for explaining tabular, text, and image data

**Less popular tools (ELI5)**

Quick model inspection, feature weight visualization, and debugging

# SHAP

# LIME

# ELI5



Impact of words (features) on prediction (=negative, in this case).

The words "movie" and "bad" contributed to the "negative" prediction positively and the word "was" contributed negatively.

# 03

## Approach Used

# Methods

| Method | DESCRIPTION |
|---|---|
| **Zero shot & Few shot prompting** | Zero shot : You provide the model with a prompt and expect it to generate the correct response without any examples. Few-shot is similar, but with examples. |
| **Chain of Thought prompting** | CoT encourages the model to explain its reasoning step by step, improving performance in complex reasoning tasks. |
| **Meta-prompting** | In meta-prompting, the prompt itself is designed to influence the structure or style of the model's response. |
| **Self-consistency prompting** | Generates multiple solutions to a given problem, then selects the most common or coherent result. |
| **Generated Knowledge prompting** | This technique instructs the model to first generate background knowledge or information before answering the main query. |
| **Prompt chaining** | Prompt chaining involves splitting complex tasks into multiple sequential prompts, where the output of one prompt becomes the input for the next. |
| **RAG** | RAG combines a language model with a retrieval system that fetches relevant documents or facts from an external database before generating a response. |
| **ReAct** | ReAct combines **reasoning** with **actions**, where the model iteratively reflects on information, interacts with tools or APIs, and refines its response. |

# Zero-shot & Few Shot



**Standard Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ❌

# Chain of Thought prompting



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✖

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔

# Meta-prompting



Focuses on the structural and syntactical aspects of tasks and problems rather than their specific content details

# Self-consistency prompting



Generates multiple solutions with different thought processes and selects answer with the most votes
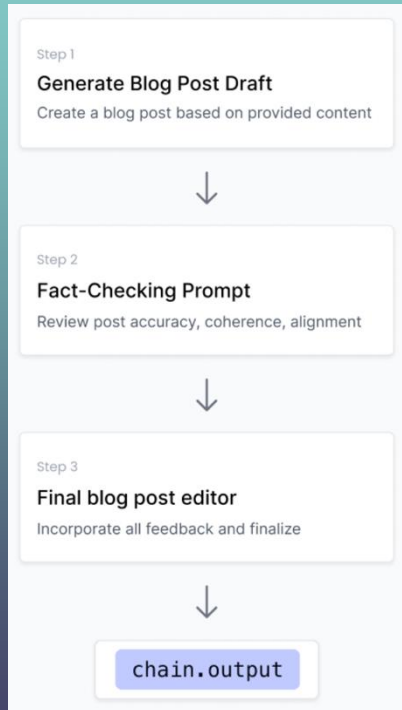
# Generated-knowledge prompting

Generate 4 facts about the Kermode bear, then use these facts to write a short blog post using the information:

The Kermode bear, also known as the spirit bear, is a rare subspecies of the American black bear that is found in the coastal regions of British Columbia, Canada.

1. The Kermode bear is a unique subspecies of the American black bear, with approximately 10% of the population having white fur.

2. The Kermode bear is a symbol of the First Nations people of British Columbia and is seen as a spirit animal.

3. The Kermode bear is an important part of the ecosystem, helping to disperse seeds and maintain the balance of the forest.

4. The Kermode bear is a protected species, and its population is monitored closely by the provincial government.
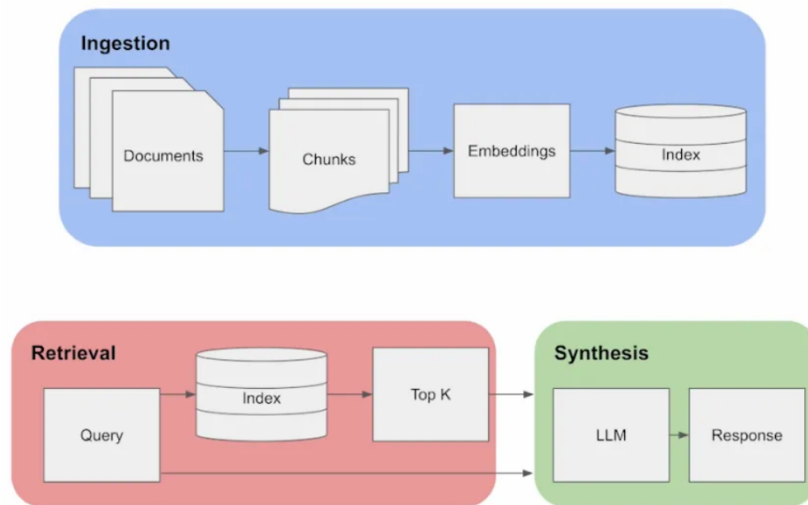
The Kermode bear is a unique and special animal that is found in the coastal regions of British Columbia, Canada. Approximately 10% of the population of this subspecies of the American black bear have white fur, making them a rare sight. The Kermode bear is a symbol of the First Nations people of British Columbia and is seen as a spirit animal. It is an important part of the ecosystem, helping to disperse seeds and maintain the balance of the forest. The Kermode bear is a protected species, and its population is monitored closely by the provincial government. It is a beautiful and majestic animal, and we must do our part to protect and conserve this species.
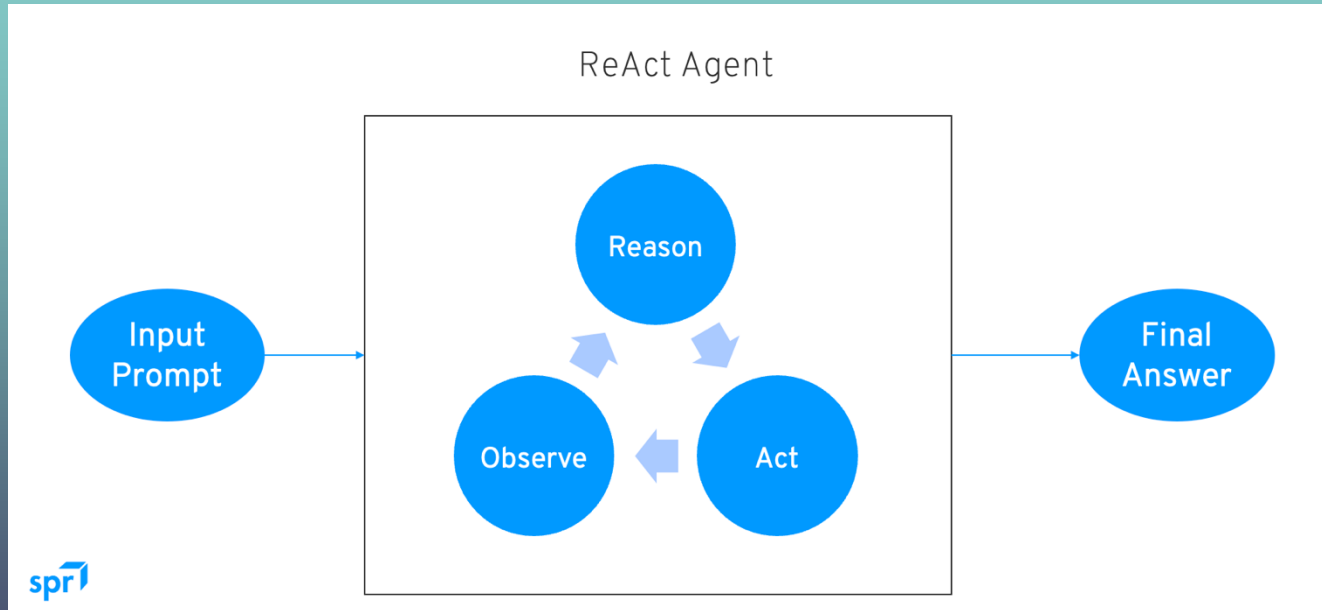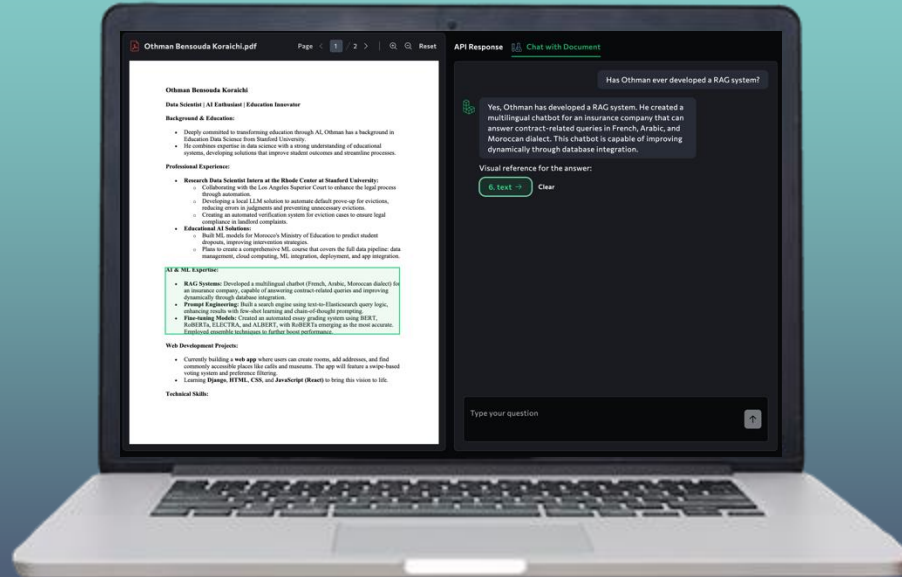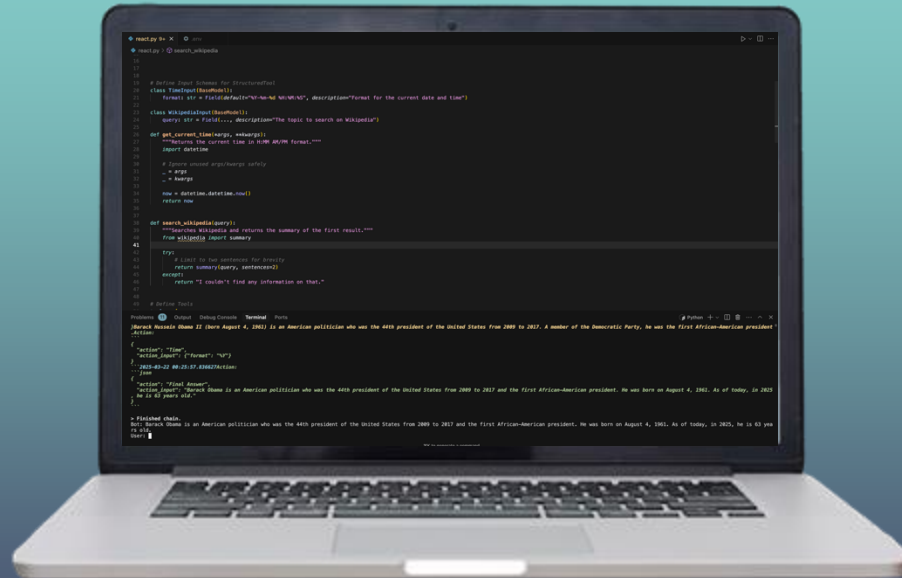
# Prompt chaining

# RAG

# ReAct

# 04

Demonstration

# DEMO (RAG)

DEMO (ReAct)

# References

- **SHAP :** Lundberg, S. M., & Lee, S.-I. (2017). *A Unified Approach to Interpreting Model Predictions*. Advances in Neural Information Processing Systems, 30.

- **LIME :** Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16), 1135-1144. https://doi

- **ELI5 :** https://eli5.readthedocs.io/en/latest/overview.html

- **Few-shot prompting :** Brown et al. (2020). *Language Models are Few-Shot Learners*. Advances in Neural Information Processing Systems, 33, 1877-1901.

- **CoT :** Wei et al. (2022). *Chain-of-thought prompting elicits reasoning in large language models*. In *Proceedings of the 36th International Conference on Neural Information Processing Systems (NeurIPS 2022)* (Article No. 1800, pp. 24824–24837).

- **Meta-prompting :** Zhang, Y., Yuan, Y., & Yao, A. C.-C. (2023). *Meta Prompting for AI Systems*. arXiv preprint arXiv:2311.11482. https://doi.org/10.48550/arXiv.2311.11482

- **Self-consistency :** Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., & Zhou, D. (2022). *Self-Consistency Improves Chain of Thought Reasoning in Language Models*. arXiv preprint arXiv:2203.11171. https://doi.org/10.48550/arXiv.2203.11171

- **Generated-Knowldege prompting :** Liu, J., Liu, A., Lu, X., Welleck, S., West, P., Le Bras, R., Choi, Y., & Hajishirzi, H. (2022). *Generated Knowledge Prompting for Commonsense Reasoning*. arXiv preprint arXiv:2110.08387. https://doi.org/10.48550/arXiv.2110.08387

- **Prompt chaining :** Sun, S., Yuan, R., Cao, Z., Li, W., & Liu, P. (2024). *Prompt Chaining or Stepwise Prompt? Refinement in Text Summarization*. arXiv preprint arXiv:2406.00507. https://doi.org/10.48550/arXiv.2406.00507

- RAG : Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. arXiv preprint arXiv:2005.11401. h

- **ReAct :** Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2022).*ReAct: Synergizing Reasoning and Acting in Language Models*. arXiv preprint arXiv:2210.03629. https://doi.org/10.48550/arXiv.2210.03629

- **Prompts :** https://www.promptingguide.ai/techniques/knowledge