

Getting Started with Llama3

Lucas Soares

12-03-2025

Methodology Notes

Methodology Notes

1. Presentation Block

Methodology Notes

1. Presentation Block

2. Notebook Demo

Methodology Notes

1. Presentation Block
2. Notebook Demo
3. Quick Q&A + Summary

Methodology Notes

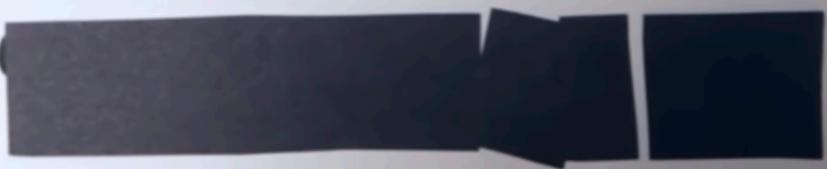
1. Presentation Block
2. Notebook Demo
3. Quick Q&A + Summary
4. Optional Exercise During Q&A

Methodology Notes

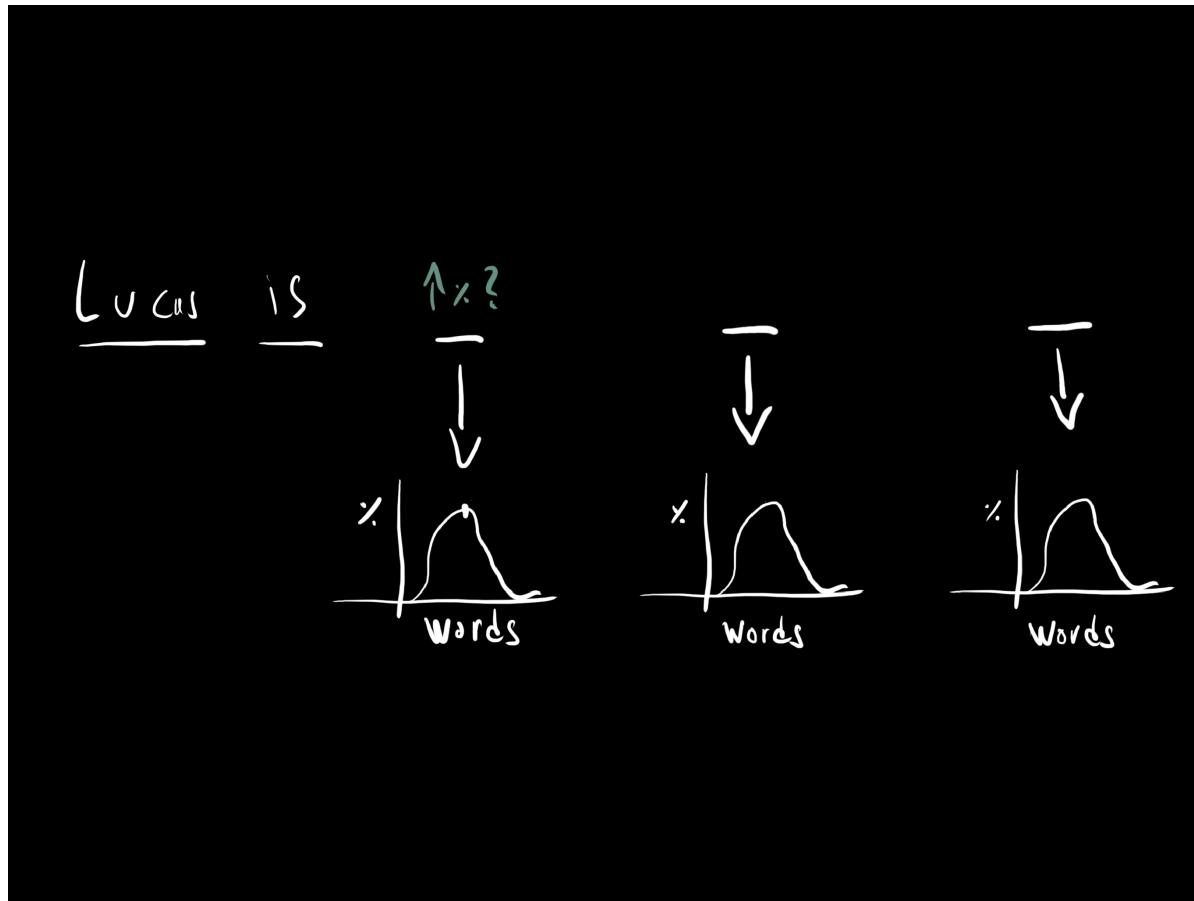
1. Presentation Block
2. Notebook Demo
3. Quick Q&A + Summary
4. Optional Exercise During Q&A
5. Repeat

LLMs Predict the Next Word

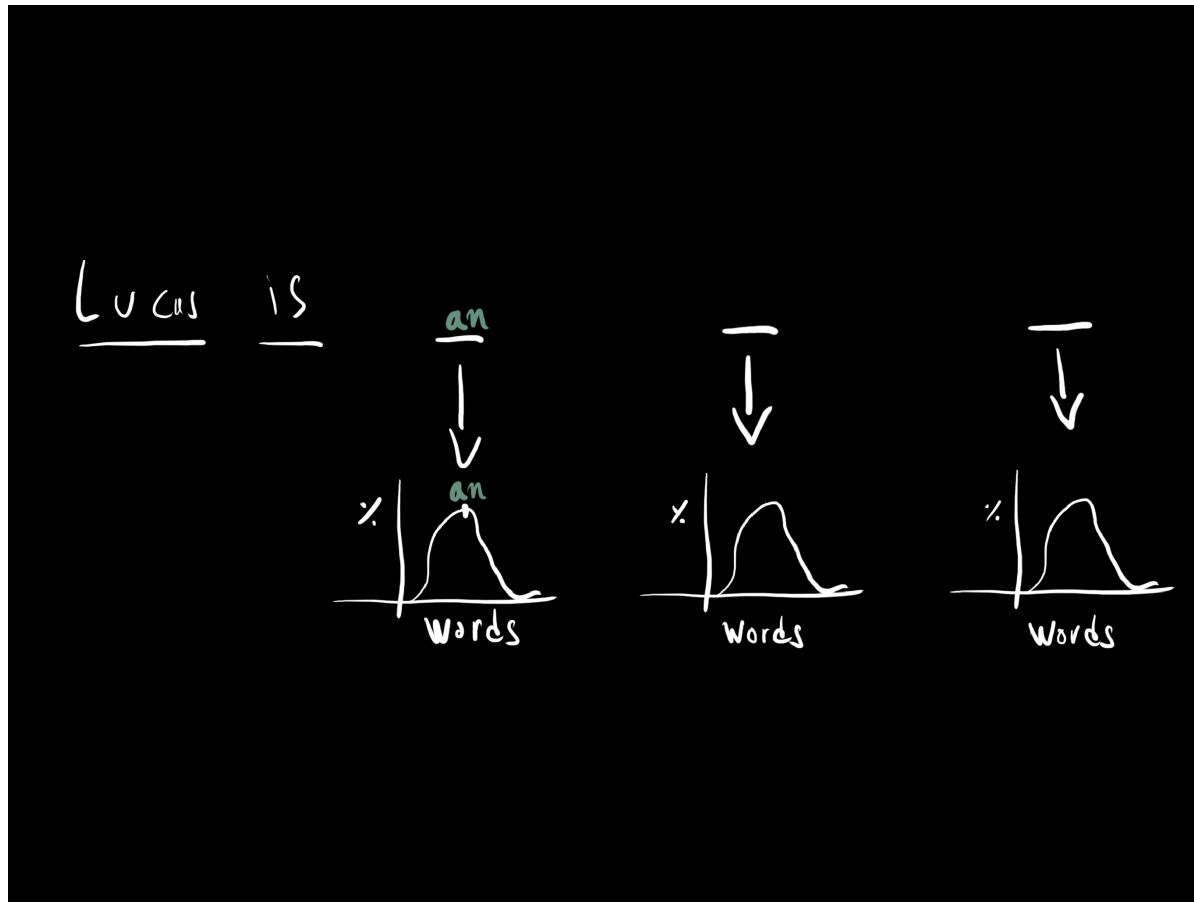
It is a thing you could not invent
with banks of



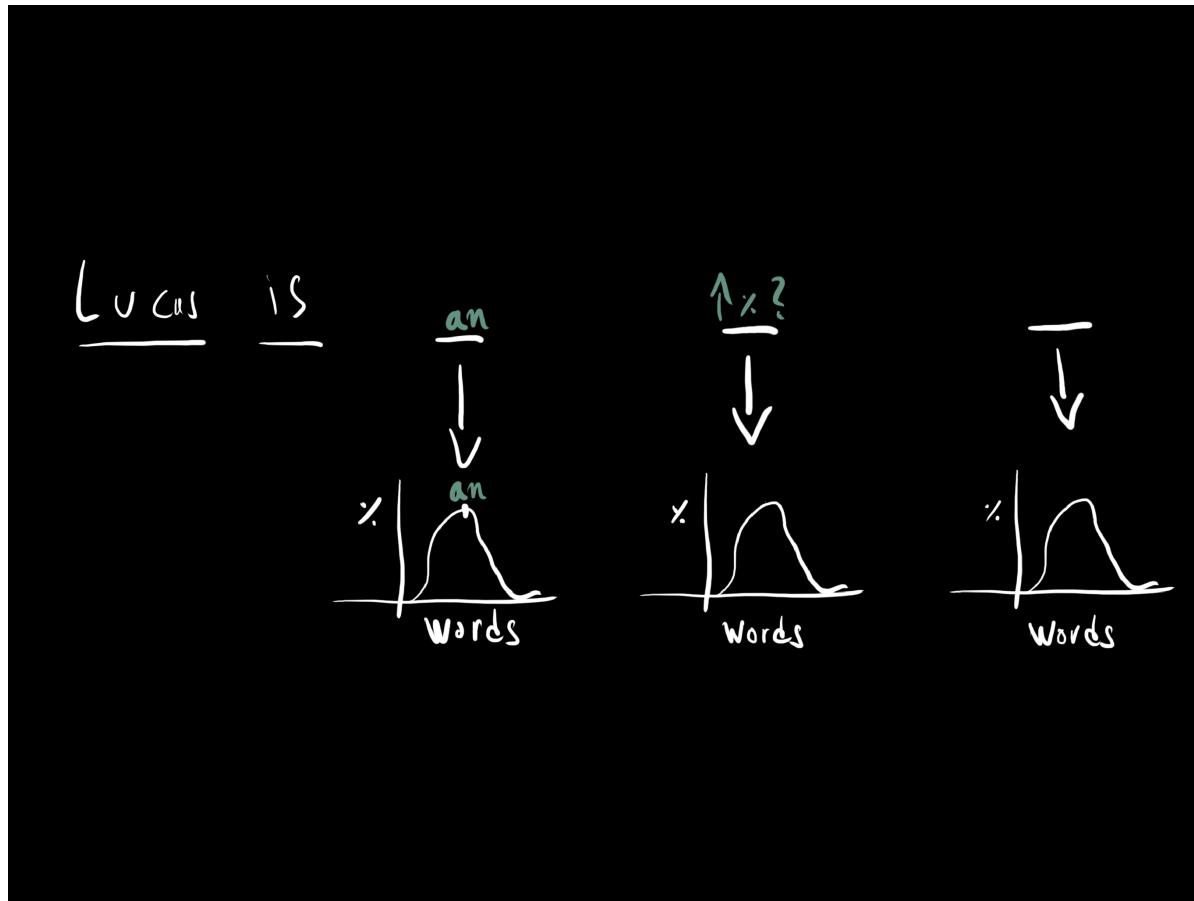
LLMs Predict the Next Word



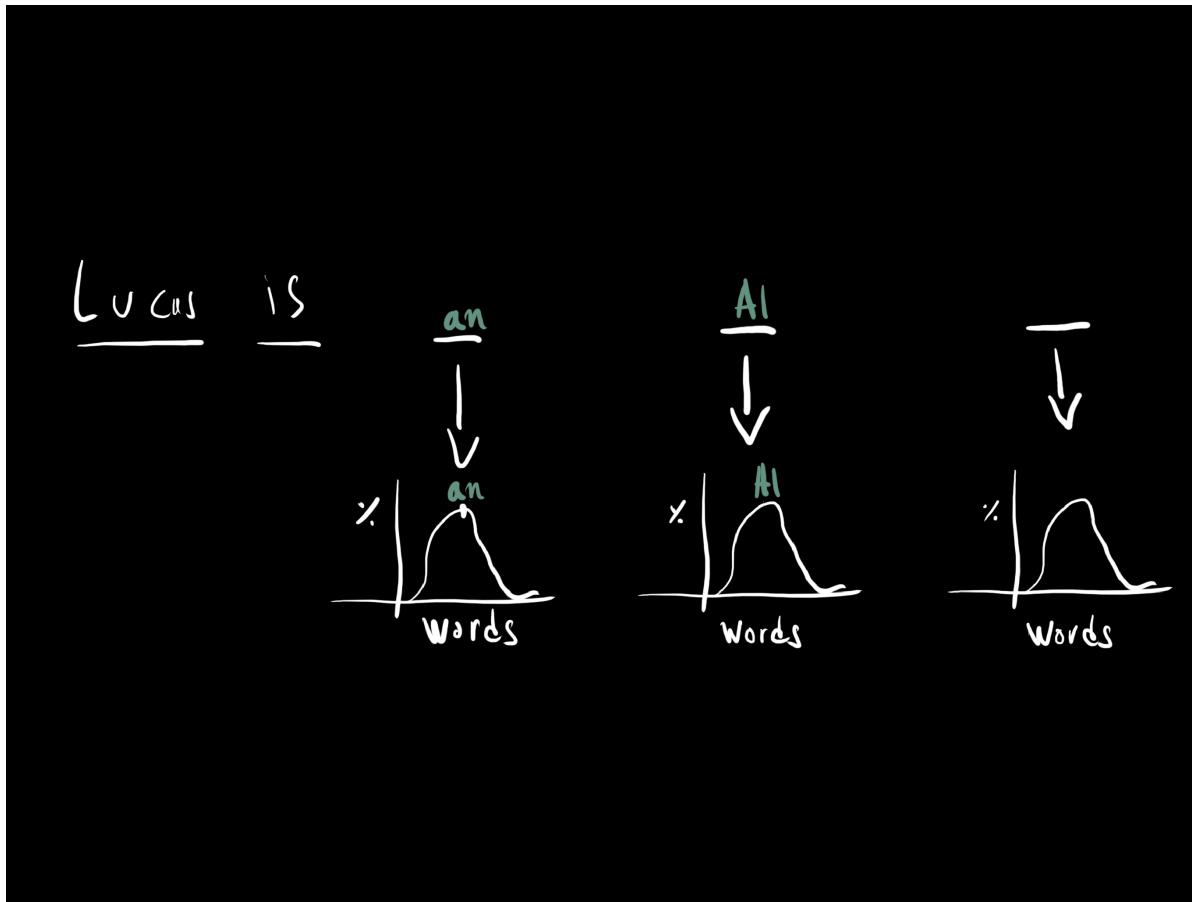
LLMs Predict the Next Word



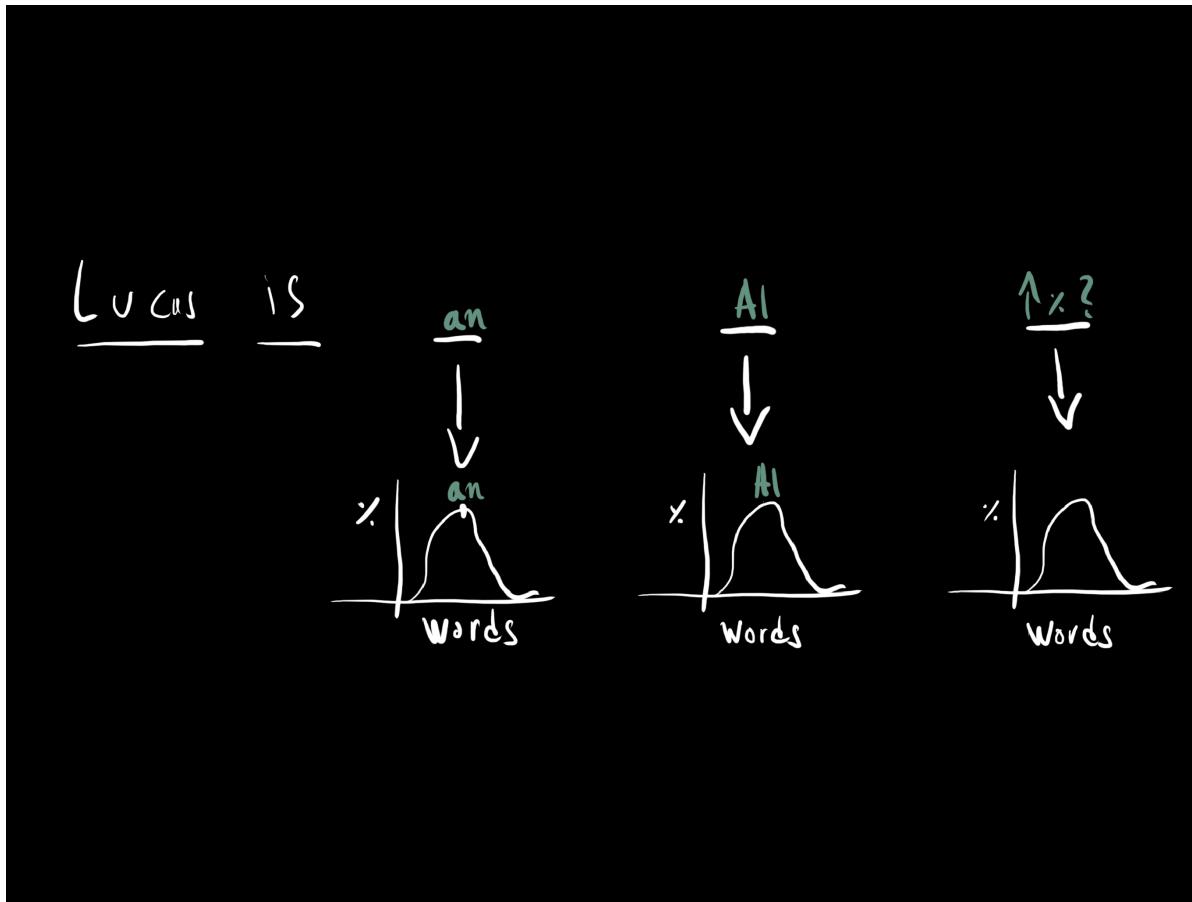
LLMs Predict the Next Word



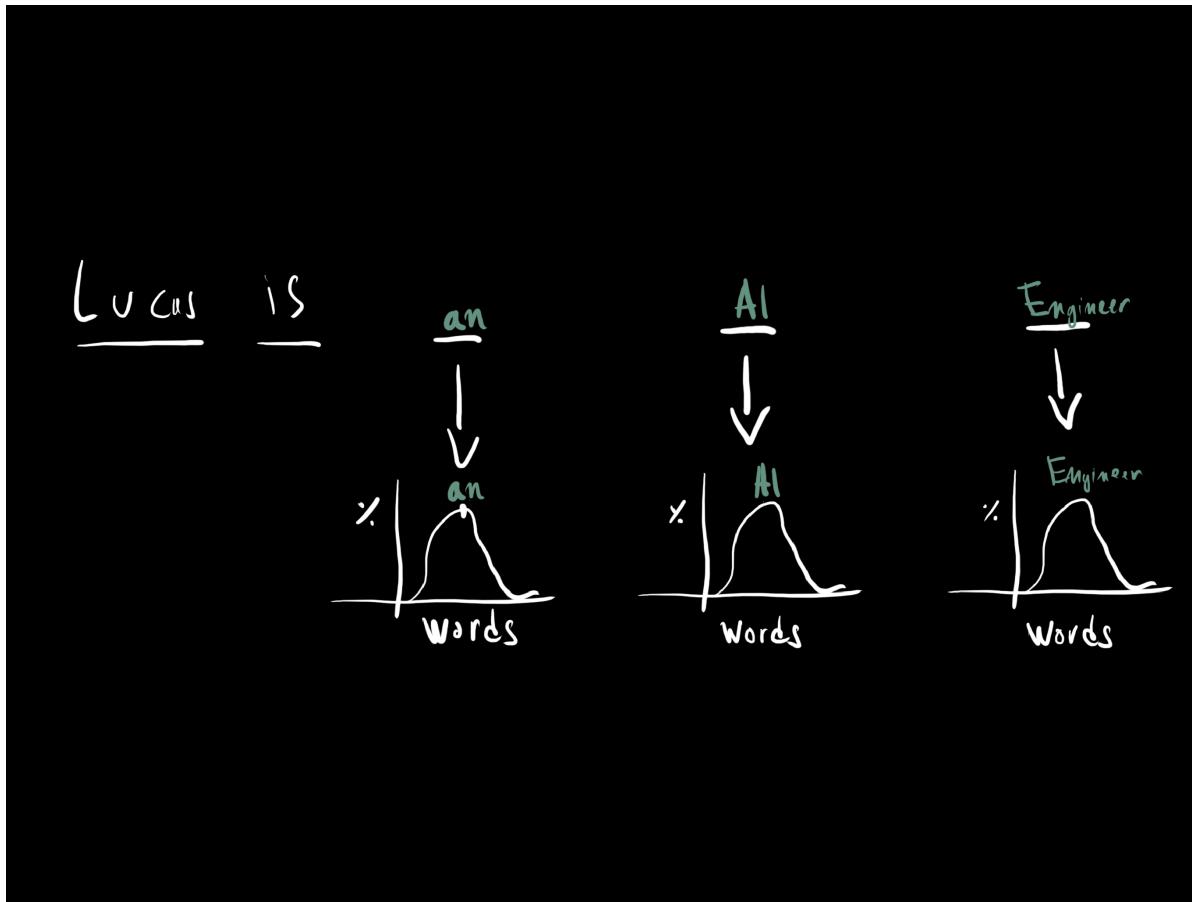
LLMs Predict the Next Word



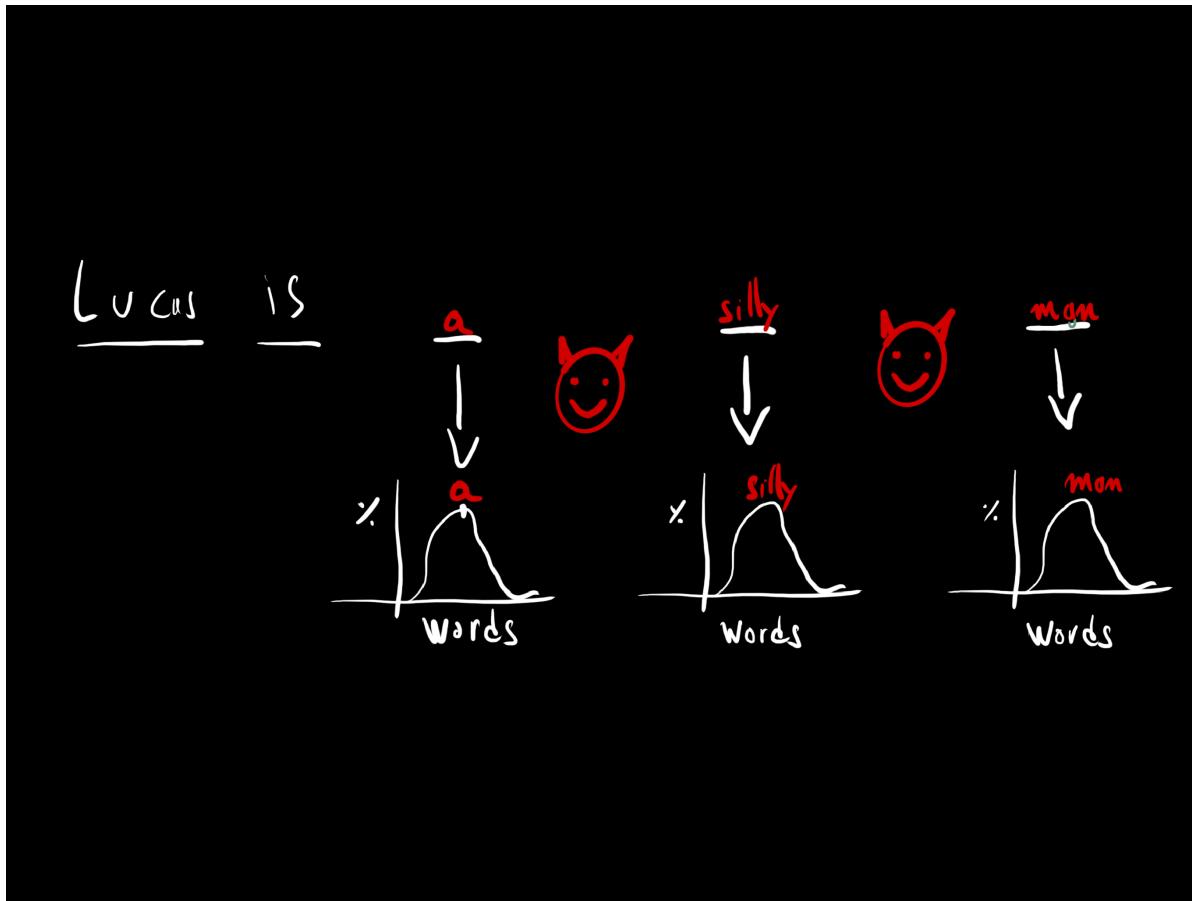
LLMs Predict the Next Word

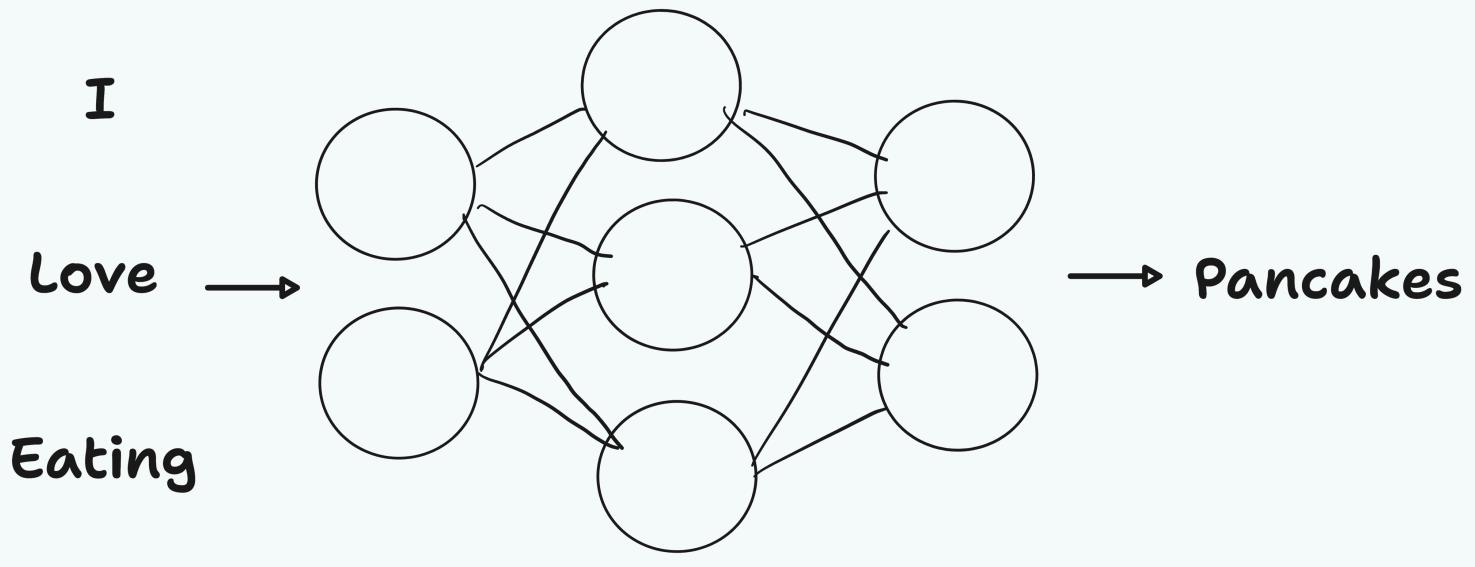


LLMs Predict the Next Word

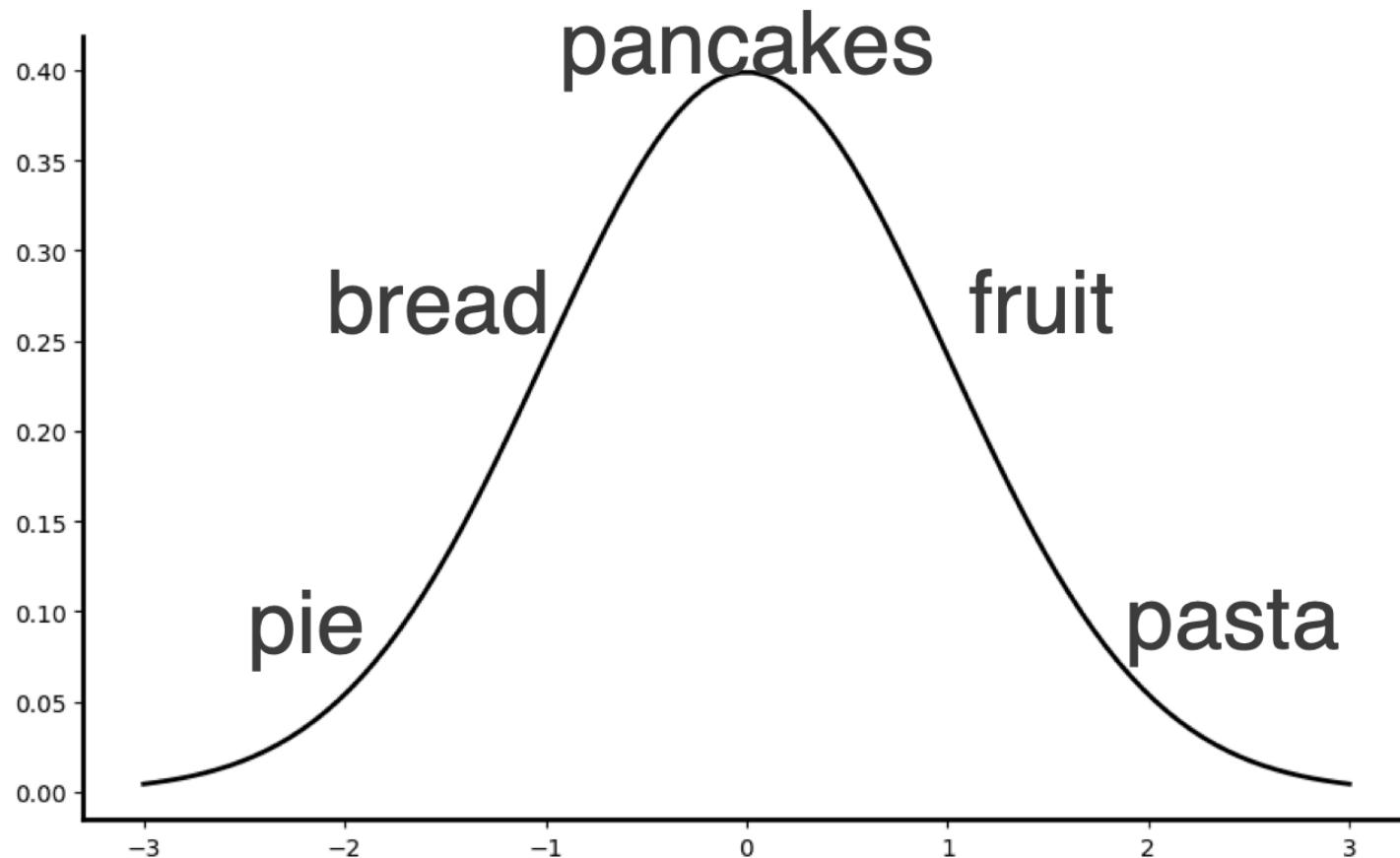


LLMs Predict the Next Word





Probability Distribution over the Next Word



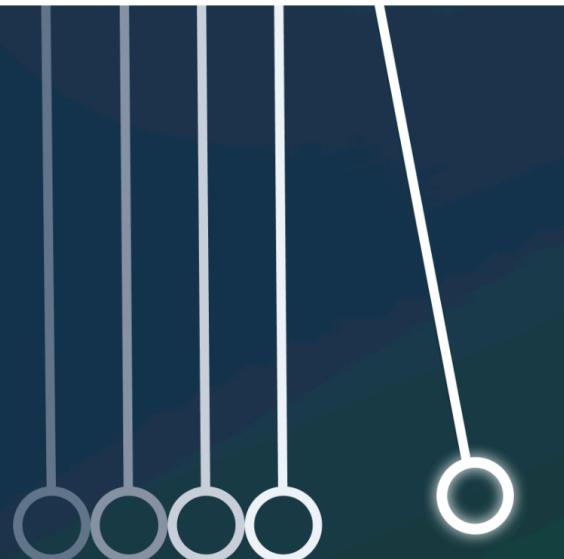
Introduction to Llama 3

The future of AI: Built with Llama

December 19, 2024 • 8 minute read

OPEN SOURCE

Looking at Llama's
impact in 2024 and
the path ahead



Meta

Llama3 Releases

- Meta Released Llama 3 in April of 2024
- Llama 3.1 was released in July of 2024
- Llama 3.2 was released in September of 2024
- Llama 3.3 was released in December of 2024

Llama 3 Series

- Open source with a Commercial license

Llama 3 Series

- Open source with a Commercial license
- The latest and greatest: Llama3.3

	Training Data	Params	Input modalities	Output modalities	Context length	GQA	Token count	Knowledge cutoff
Llama 3.3 (text only)	A new mix of publicly available online data.	70B	Multilingual Text	Multilingual Text and code	128k	Yes	15T+	December 2023

Llama 3.3 Evaluation Results

Category	Benchmark	# Shots	Metric	Llama 3.1 8B Instruct	Llama 3.1 70B Instruct	Llama-3.3 70B Instruct
	MMLU (CoT)	0	macro_avg/acc	73.0	86.0	86.0
	MMLU Pro (CoT)	5	macro_avg/acc	48.3	66.4	68.9
Steerability	IFEval			80.4	87.5	92.1
Reasoning	GPQA Diamond (CoT)	0	acc	31.8	48.0	50.5
Code	HumanEval	0	pass@1	72.6	80.5	88.4
	MBPP EvalPlus (base)	0	pass@1	72.8	86.0	87.6
Math	MATH (CoT)	0	sympy_intersection_score	51.9	68.0	77.0
Tool Use	BFCL v2	0	overall_ast_summary/macro_avg/valid	65.4	77.5	77.3
Multilingual	MGSM	0	em	68.9	86.9	91.1

Llama 3.3 Evaluation Results

Category	Benchmark	# Shots	Metric	Llama 3.1 8B Instruct	Llama 3.1 70B Instruct	Llama-3.3 70B Instruct
	MMLU (CoT)	0	macro_avg/acc	73.0	86.0	86.0
	MMLU Pro (CoT)	5	macro_avg/acc	48.3	66.4	68.9
Steerability	IFEval			80.4	87.5	92.1
Reasoning	GPQA Diamond (CoT)	0	acc	31.8	48.0	50.5
Code	HumanEval	0	pass@1	72.6	80.5	88.4
	MBPP EvalPlus (base)	0	pass@1	72.8	86.0	87.6
Math	MATH (CoT)	0	sympy_intersection_score	51.9	68.0	77.0
Tool Use	BFCL v2	0	overall_ast_summary/macro_avg/valid	65.4	77.5	77.3
Multilingual	MGSM	0	em	68.9	86.9	91.1

good at general problem solving

Llama 3.3 Evaluation Results

Category	Benchmark	# Shots	Metric	Llama 3.1 8B Instruct	Llama 3.1 70B Instruct	Llama-3.3 70B Instruct
	MMLU (CoT)	0	macro_avg/acc	73.0	86.0	86.0
	MMLU Pro (CoT)	5	macro_avg/acc	48.3	66.4	68.9
Steerability	IFEval			80.4	87.5	92.1
Reasoning	GPQA Diamond (CoT)	0	acc	31.8	48.0	50.5
Code	HumanEval	0	pass@1	72.6	80.5	88.4
	MBPP EvalPlus (base)	0	pass@1	72.8	86.0	87.6
Math	MATH (CoT)	0	sympy_intersection_score	51.9	68.0	77.0
Tool Use	BFCL v2	0	overall_ast_summary/macro_avg/valid	65.4	77.5	77.3
Multilingual	MGSM	0	em	68.9	86.9	91.1

follows
instructions
well

Llama 3.3 Evaluation Results

Category	Benchmark	# Shots	Metric	Llama 3.1 8B Instruct	Llama 3.1 70B Instruct	Llama-3.3 70B Instruct
	MMLU (CoT)	0	macro_avg/acc	73.0	86.0	86.0
	MMLU Pro (CoT)	5	macro_avg/acc	48.3	66.4	68.9
Steerability	IFEval			80.4	87.5	92.1
Reasoning	GPQA Diamond (CoT)	0	acc	31.8	48.0	50.5
Code	HumanEval	0	pass@1	72.6	80.5	88.4
	MBPP EvalPlus (base)	0	pass@1	72.8	86.0	87.6
Math	MATH (CoT)	0	sympy_intersection_score	51.9	68.0	77.0
Tool Use	BFCL v2	0	overall_ast_summary/macro_avg/valid	65.4	77.5	77.3
Multilingual	MGSM	0	em	68.9	86.9	91.1

good at thinking
GPT-4o is 54%

Llama 3.3 Evaluation Results

Category	Benchmark	# Shots	Metric	Llama 3.1 8B Instruct	Llama 3.1 70B Instruct	Llama-3.3 70B Instruct
	MMLU (CoT)	0	macro_avg/acc	73.0	86.0	86.0
	MMLU Pro (CoT)	5	macro_avg/acc	48.3	66.4	68.9
Steerability	IFEval			80.4	87.5	92.1
Reasoning	GPQA Diamond (CoT)	0	acc	31.8	48.0	50.5
Code	HumanEval	0	pass@1	72.6	80.5	88.4
	MBPP EvalPlus (base)	0	pass@1	72.8	86.0	87.6
Math	MATH (CoT)	0	sympy_intersection_score	51.9	68.0	77.0
Tool Use	BFCL v2	0	overall_ast_summary/macro_avg/valid	65.4	77.5	77.3
Multilingual	MGSM	0	em	68.9	86.9	91.1

good at coding

Whiteboard - Llama 3 Offering Breakdown

Q&A / Break

Notebook Demo - Introduction to Using The Llama Series Models

Q&A / Break

Query Your Docs Locally with Llama 3

- Need for LLMs with access to context-relevant data.



Query Your Docs Locally with Llama 3

- Privacy concern with closed source LLMs.



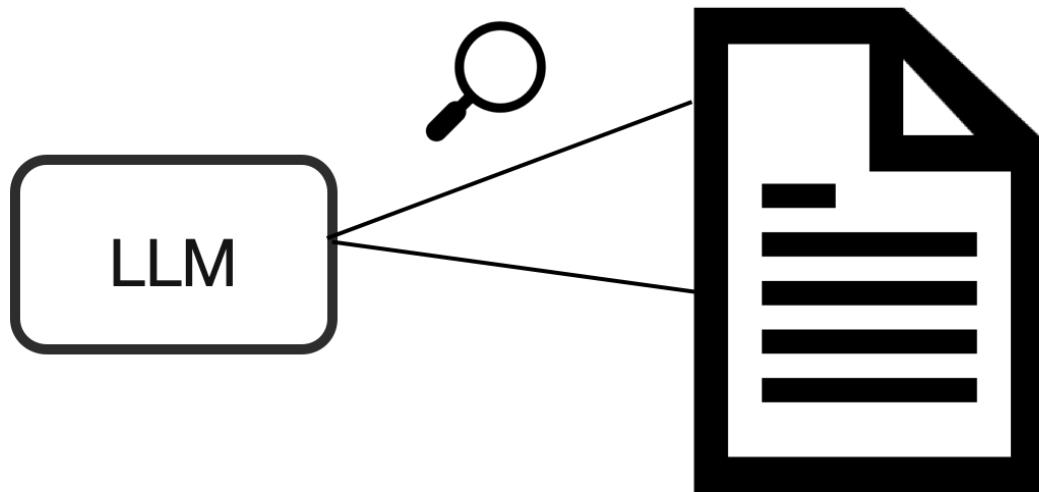
Query Your Docs Locally with Llama 3

- Solution? Llama 3!

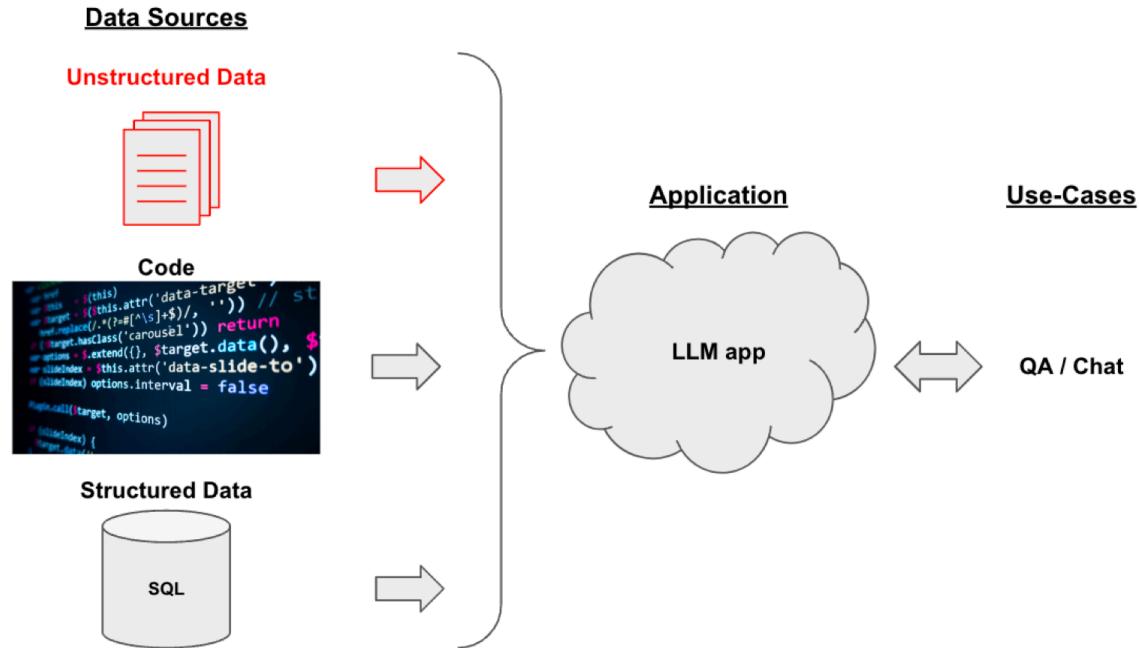


RAG with Llama 3

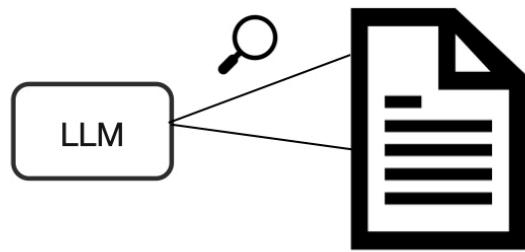
- RAG - Retrieval Augmented Generation



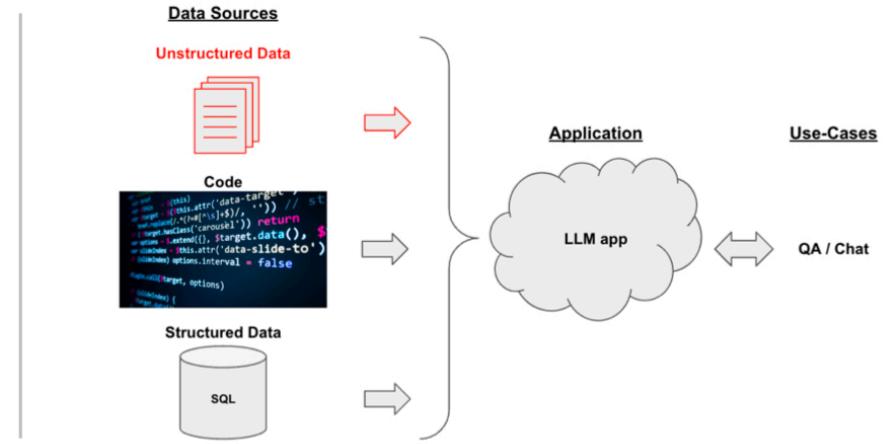
RAG with Llama 3



RAG - Retrieval Augmented Generation



RAG - Retrieval Augmented Generation





RAG - Retrieval Augmented Generation

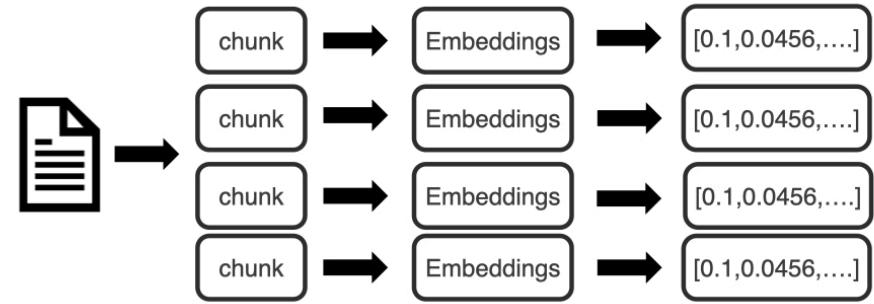


LLMs have a limited context length

RAG - Retrieval Augmented Generation

LLMs have a limited context length

Embeddings: capture content and meaning





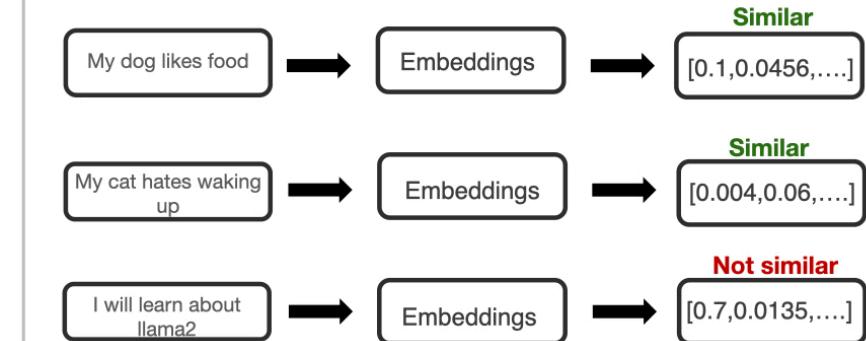
RAG - Retrieval Augmented Generation



LLMs have a limited context length



Embeddings: capture content and meaning

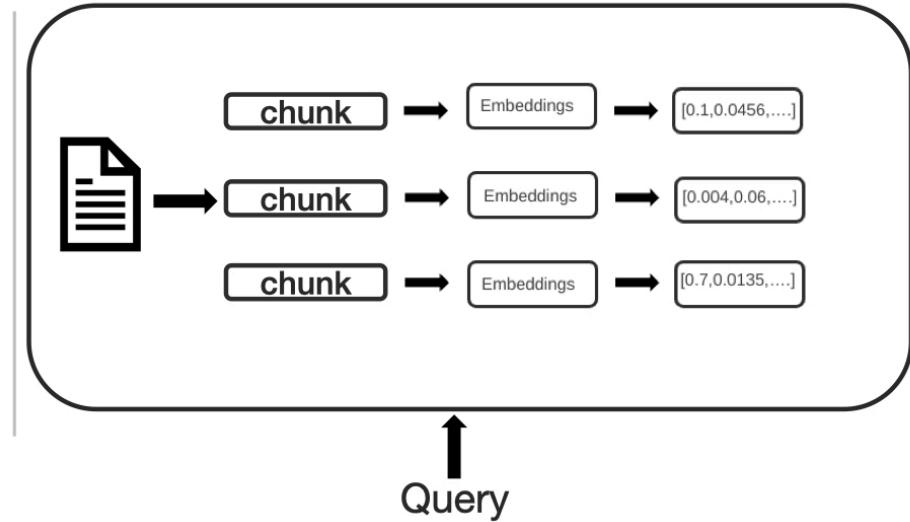


- RAG - Retrieval Augmented Generation

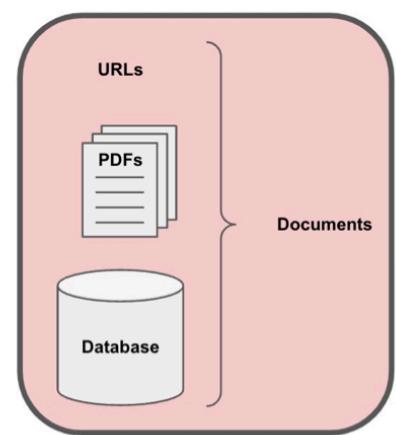
- LLMs have a limited context length

- **Embeddings:** capture content and meaning

Vector Database



Document Loading



Document Loading

URLs

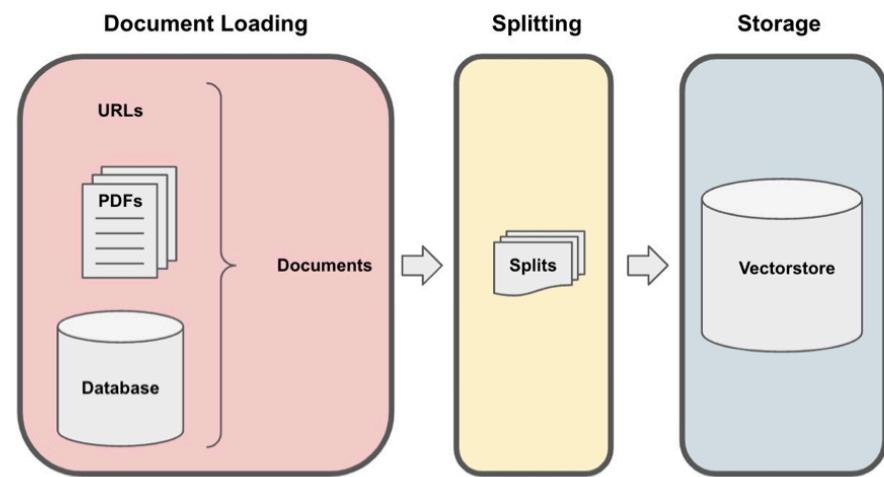


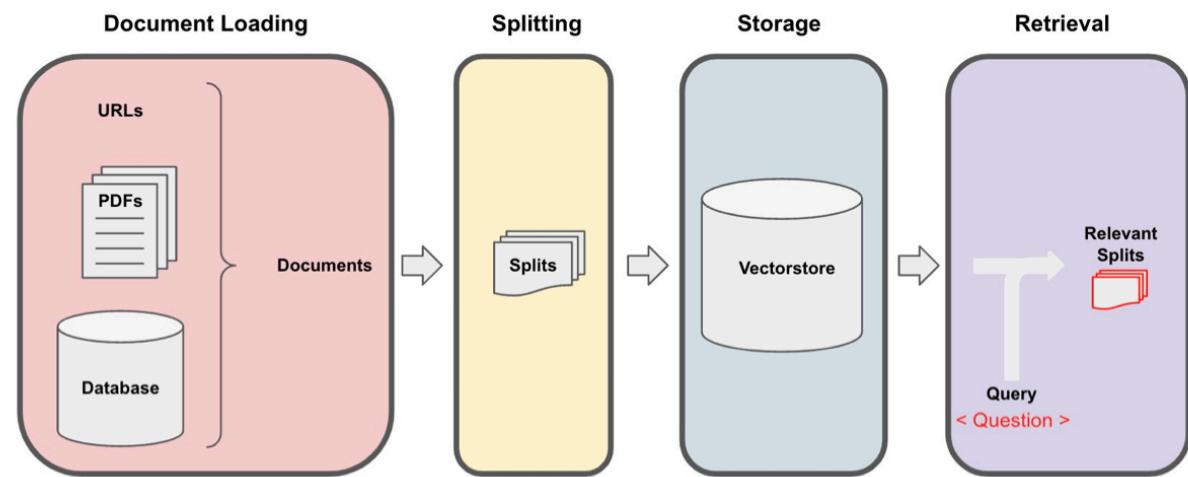
Database

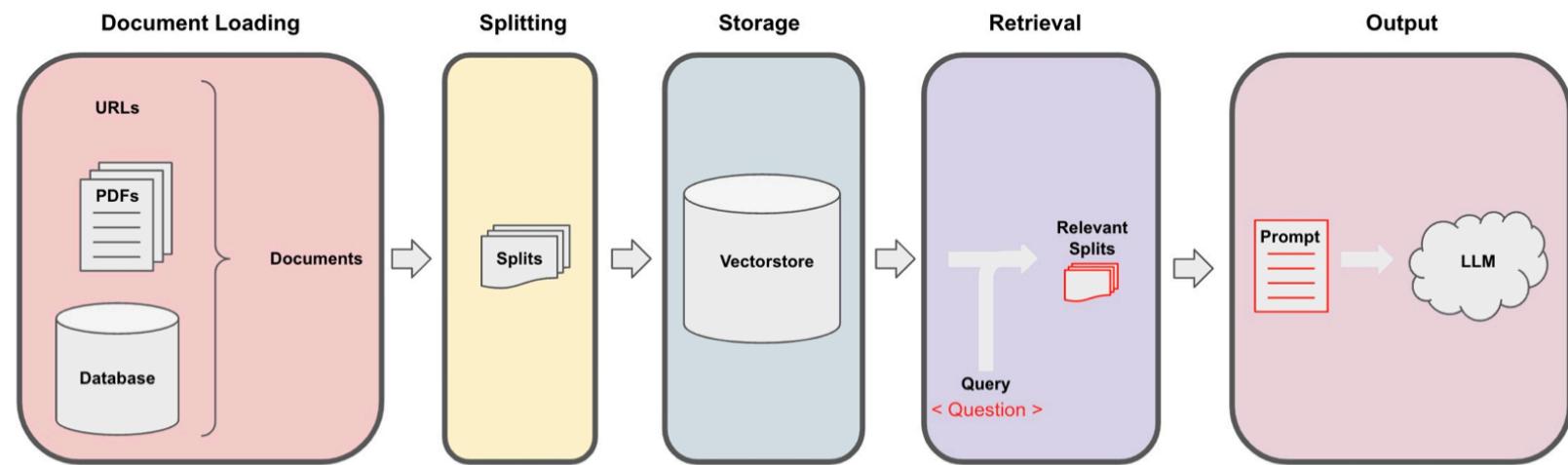
Splitting

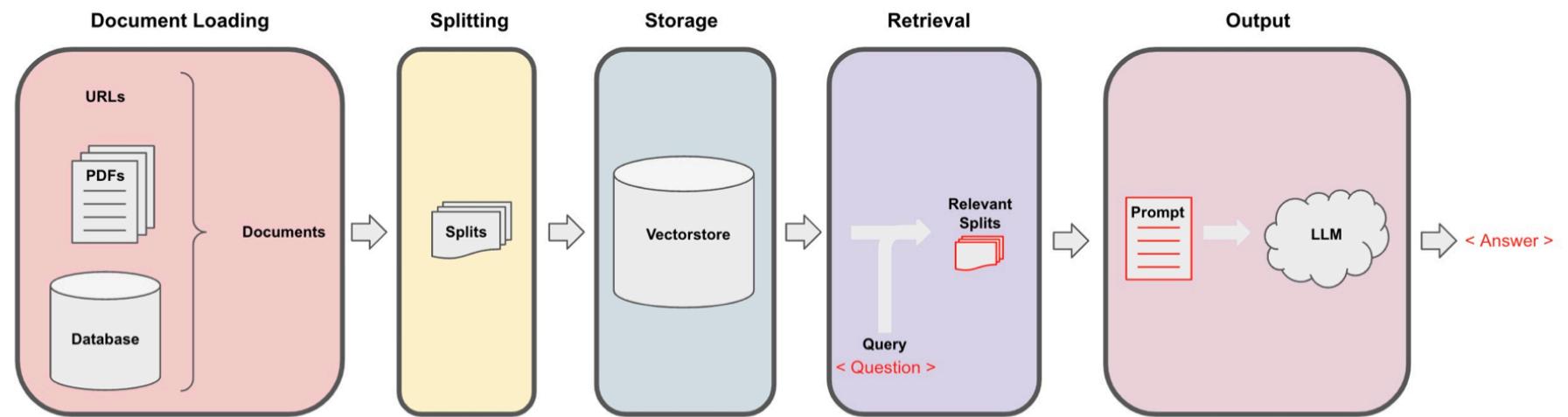
Documents









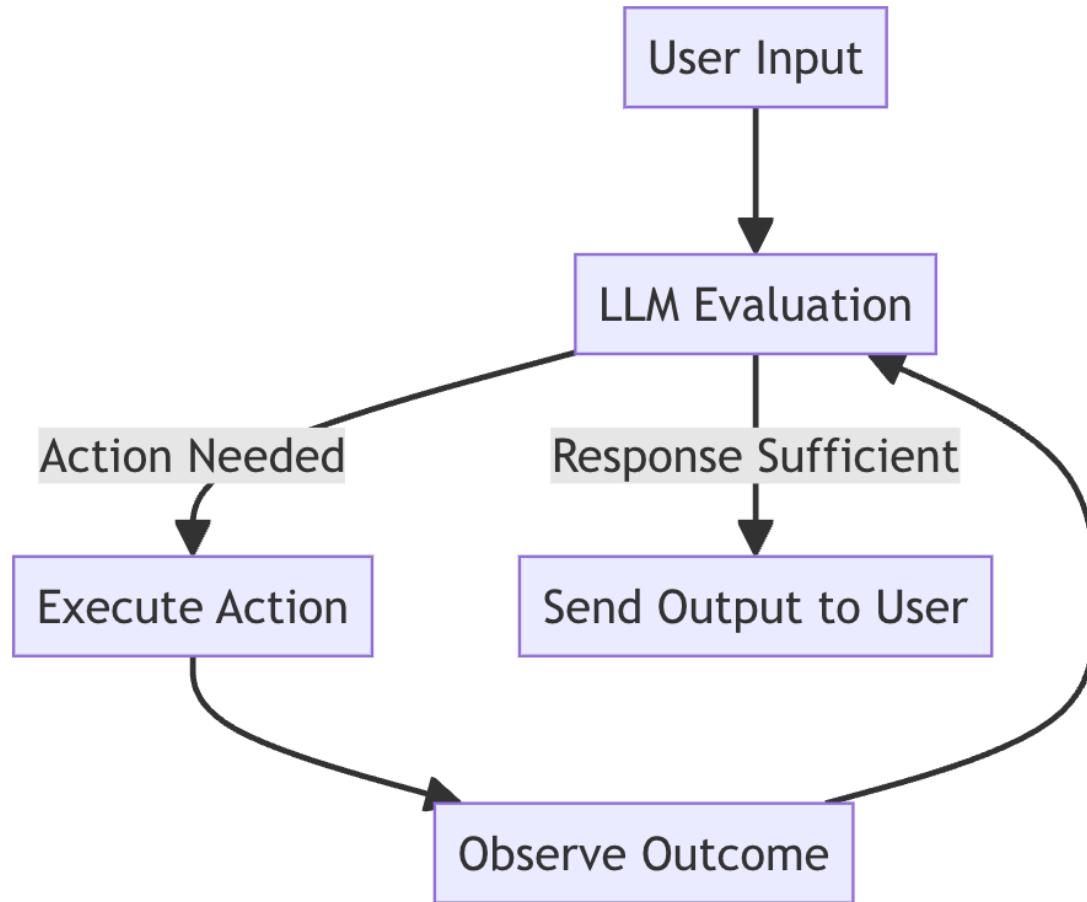


Notebook Demo - Local RAG with Llama 3

Q&A / Break

Local Agents with Llama 3

Local Agents with Llama 3



Explanation of the agent loop in cognitive architectures.

Practical Use Case: Private Docs Chatbot

Practical Use Case: Private Docs Chatbot

- **Scenario:** An LLM-powered chatbot designed to assist users by retrieving information from private documents.

Practical Use Case: Private Docs Chatbot

- **Scenario:** An LLM-powered chatbot designed to assist users by retrieving information from private documents.
- **User Input:** User asks about specific details from an internal document.

Practical Use Case: Private Docs Chatbot

- **Scenario:** An LLM-powered chatbot designed to assist users by retrieving information from private documents.
- **User Input:** User asks about specific details from an internal document.
- **LLM Decision:** Determines if the requested information requires search the document repository.

Practical Use Case: Private Docs Chatbot

- **Scenario:** An LLM-powered chatbot designed to assist users by retrieving information from private documents.
- **User Input:** User asks about specific details from an internal document.
- **LLM Decision:** Determines if the requested information requires search the document repository.
- **Action Taken:** If a search is required, the agent retrieves relevant document sections and provides a concise response based on the content.

Practical use case of LLM agents in customer support.

Notebook Demo - Tool Calling/Structured Outputs and local agents with Llama 3

Q&A / Break

Fine Tuning Llama 3

Whiteboard - Fine Tuning Llama 3: what, when, why and beyond

Fine Tuning Llama 3

- What is fine-tuning?

Fine Tuning Llama 3

- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.

Fine Tuning Llama 3

- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.
- What does it do?

Fine Tuning Llama 3

- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.
- What does it do?
- Fine-tuning transfers the pre-trained model's learned patterns and features to new tasks, improving performance and reducing training data needs.

Fine Tuning Llama 3

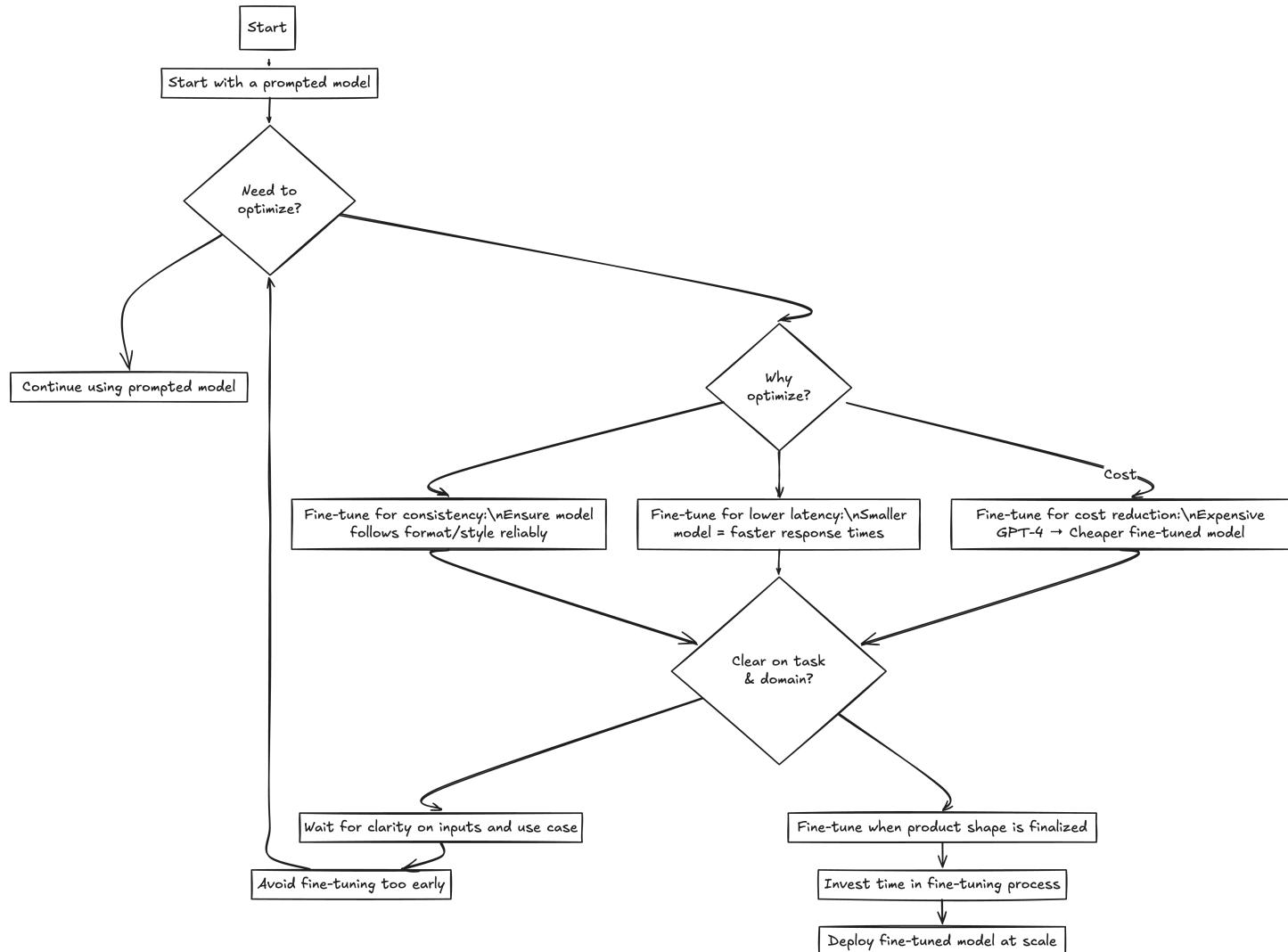
- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.
- What does it do?
- Fine-tuning transfers the pre-trained model's learned patterns and features to new tasks, improving performance and reducing training data needs.
- Why fine-tune?

Fine Tuning Llama 3

- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.
- What does it do?
- Fine-tuning transfers the pre-trained model's learned patterns and features to new tasks, improving performance and reducing training data needs.
- Why fine-tune?
- To improve the model's performance on a specific task.

Fine Tuning Llama 3

- What is fine-tuning?
- It's the process of training a smaller model on a specific task leveraging a domain-specific dataset.
- What does it do?
- Fine-tuning transfers the pre-trained model's learned patterns and features to new tasks, improving performance and reducing training data needs.
- Why fine-tune?
- To improve the model's performance on a specific task.
- When to fine-tune?



Notebook Demo - Fine-Tuning Llama 3 - Walkthrough

Q&A / Break

Notebook Demo - End To End Example with Llama 3

References

1. [Llama 3.3 Model Card](#)
2. [Fine Tuning Video](#)
3. [High Quality Structured Report Generation Demo](#)
4. [Llama Cookbook – Finetuning](#)
5. [Llama Recipes](#)
6. [Generating Synthetic Data](#)
7. [Llama Cookbook – End-to-End Use Cases](#)
8. [Structured Outputs with Ollama](#)
9. [Chunk Size Explanation](#)
10. [Embeddings + Llama by Simon Willison](#)
11. [Embeddings + Parquet for Small Projects](#)
12. [Awesome Local LLMs](#)
13. [LLM-Driven Data Engineering](#)
14. [Text Chunking in RAG](#)
15. [Llama Resources – AI.Meta](#)
16. [Llama 3 Paper \(arXiv\)](#)

