

Getting Started with Llama3.1

Lucas Soares

08-08-2024

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

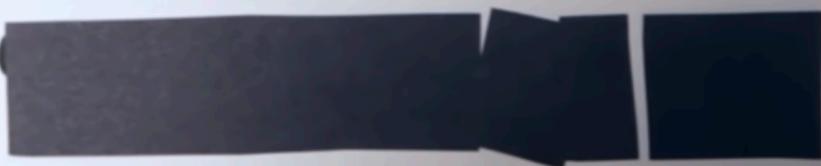
1. Presentation Block
2. Notebook Demo
3. Quick Q&A + Summary
4. Optional Exercise During Q&A

Methodology Notes

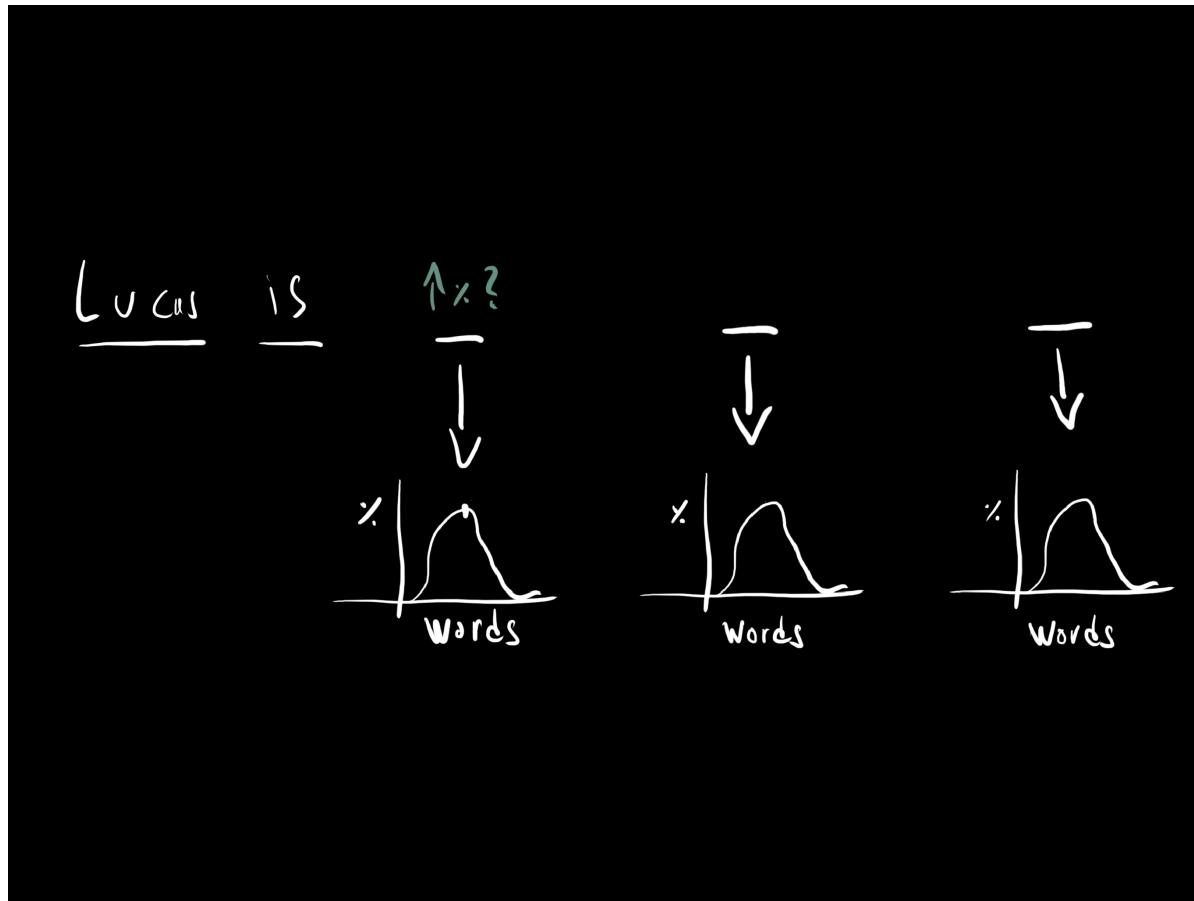
1. Presentation Block
2. Notebook Demo
3. Quick Q&A + Summary
4. Optional Exercise During Q&A
5. Repeat

LLMs Predict the Next Word

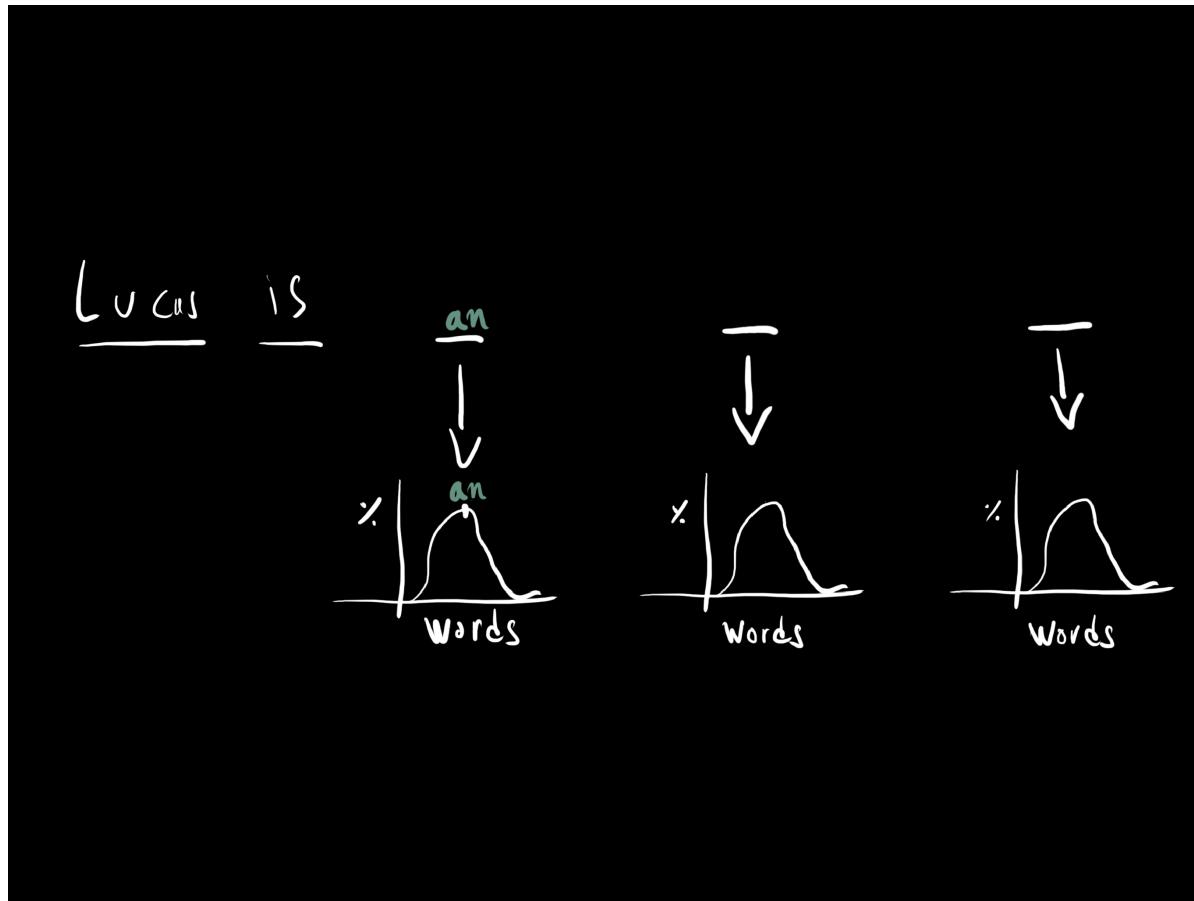
It is a thing you could not invent
with banks of



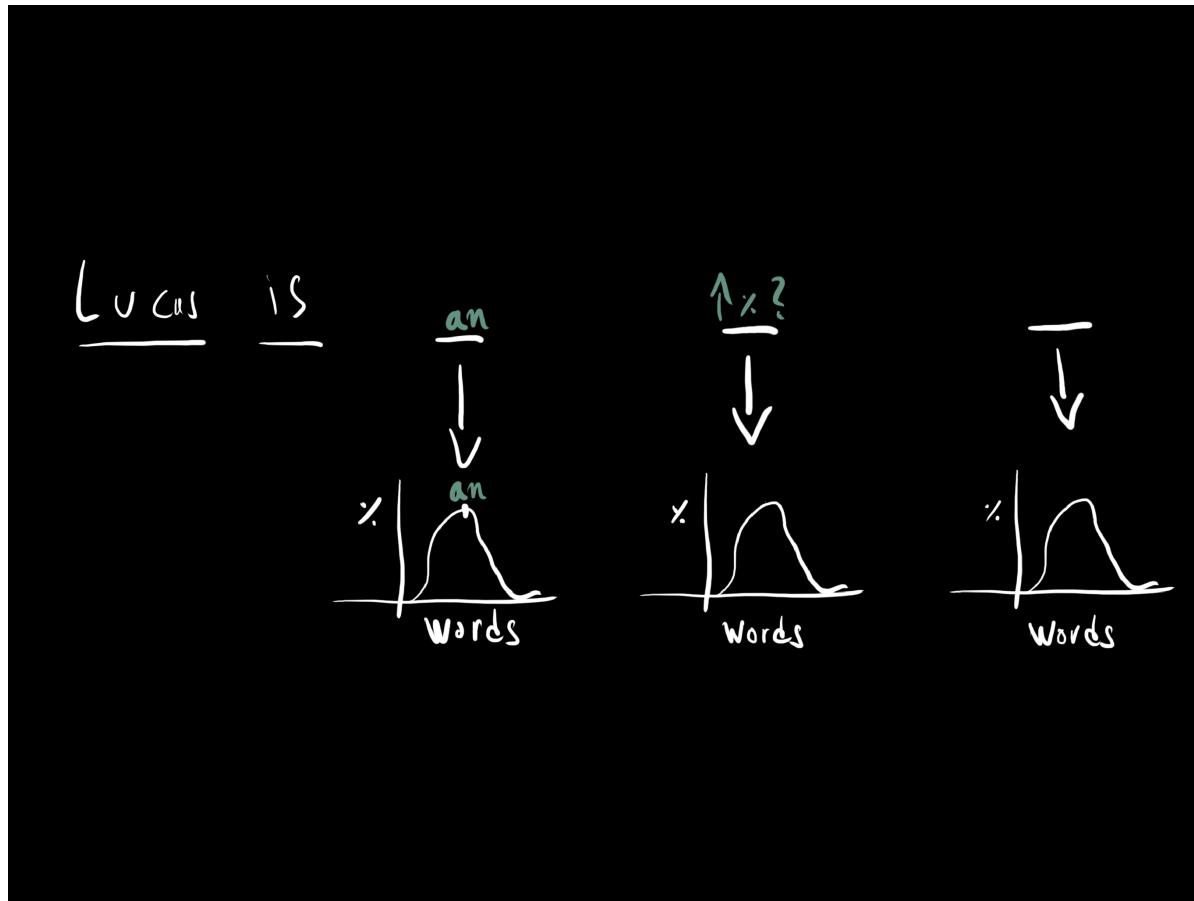
LLMs Predict the Next Word



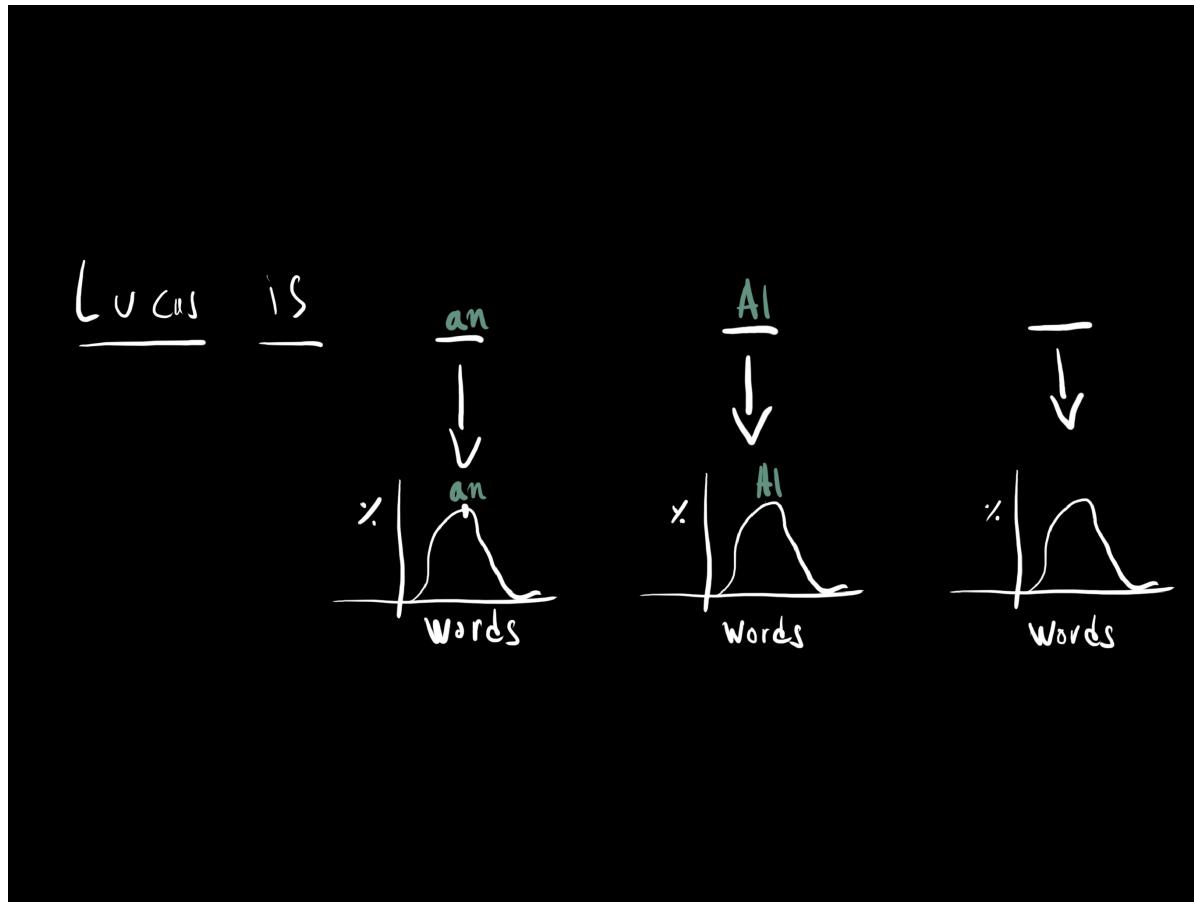
LLMs Predict the Next Word



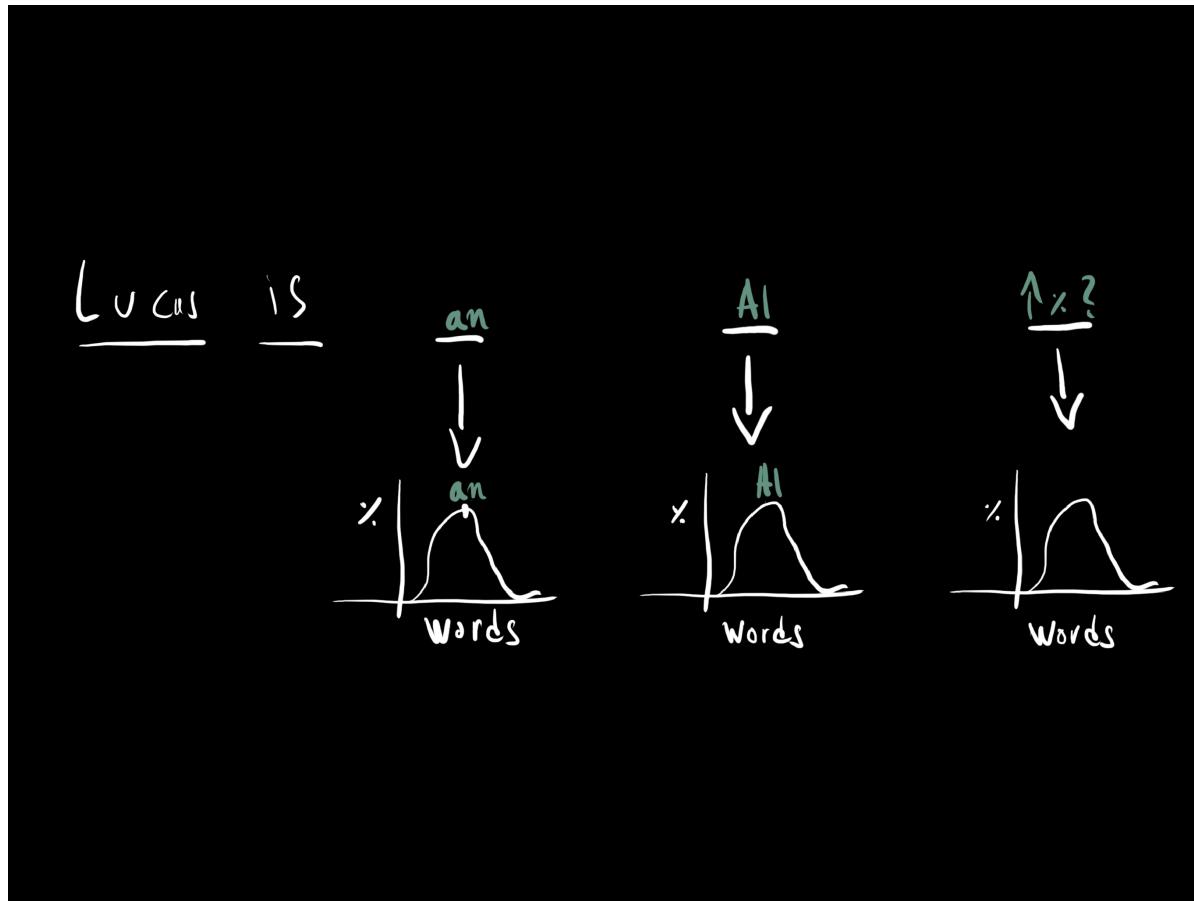
LLMs Predict the Next Word



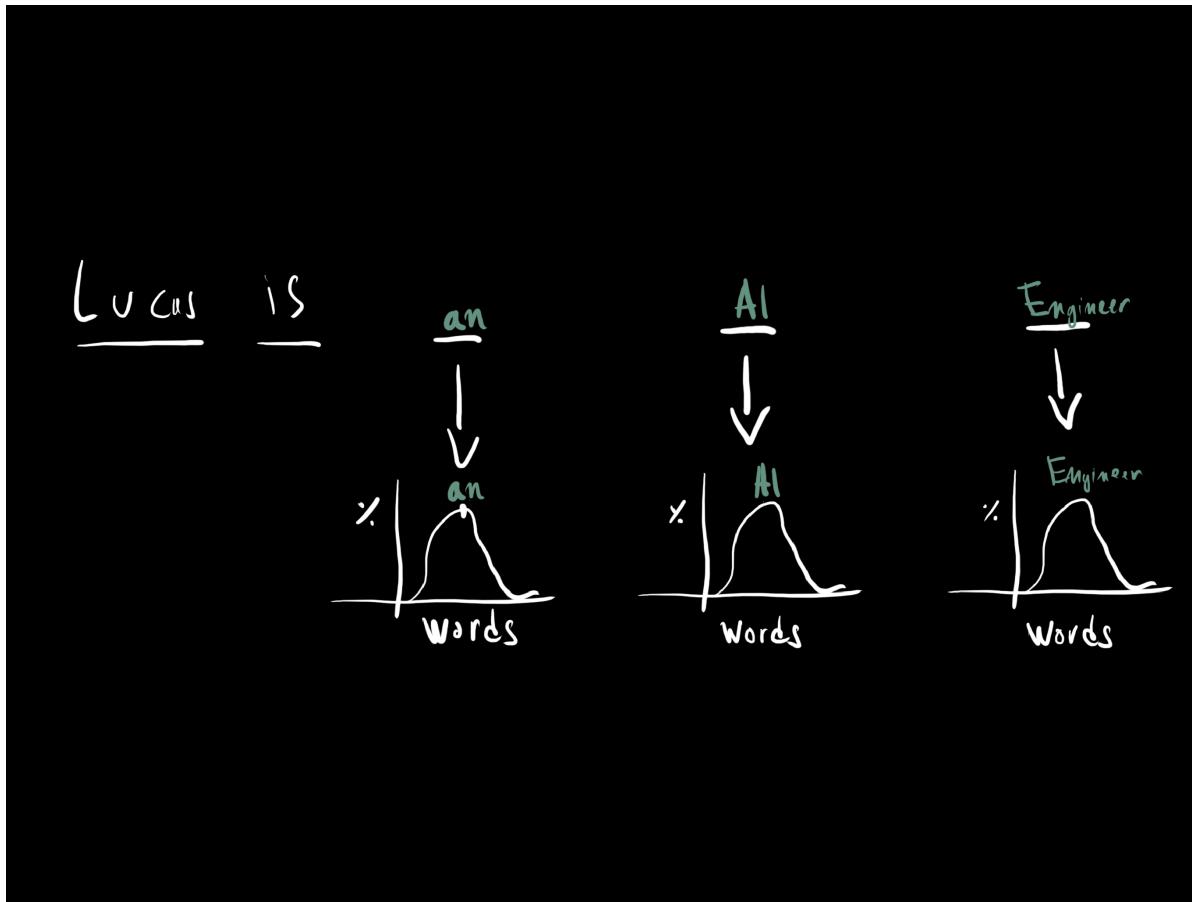
LLMs Predict the Next Word



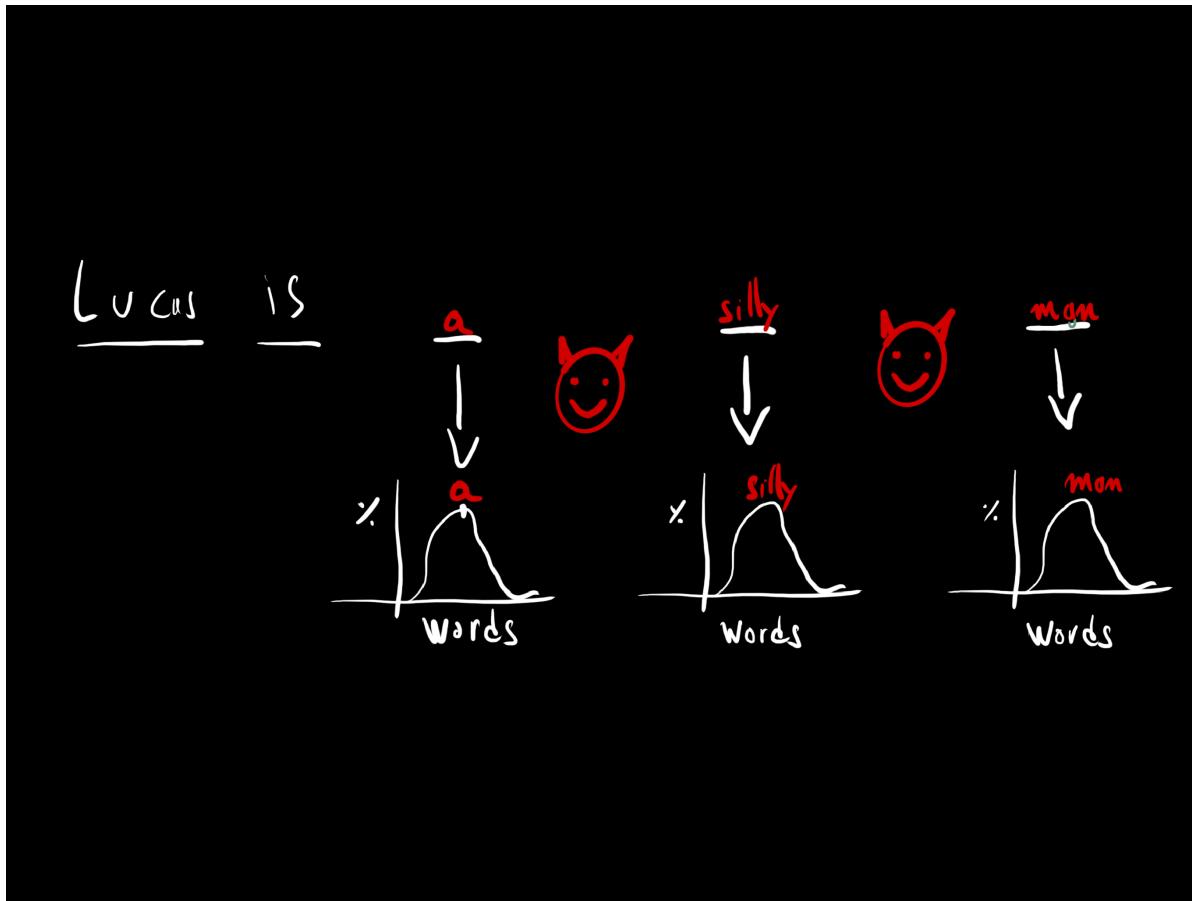
LLMs Predict the Next Word

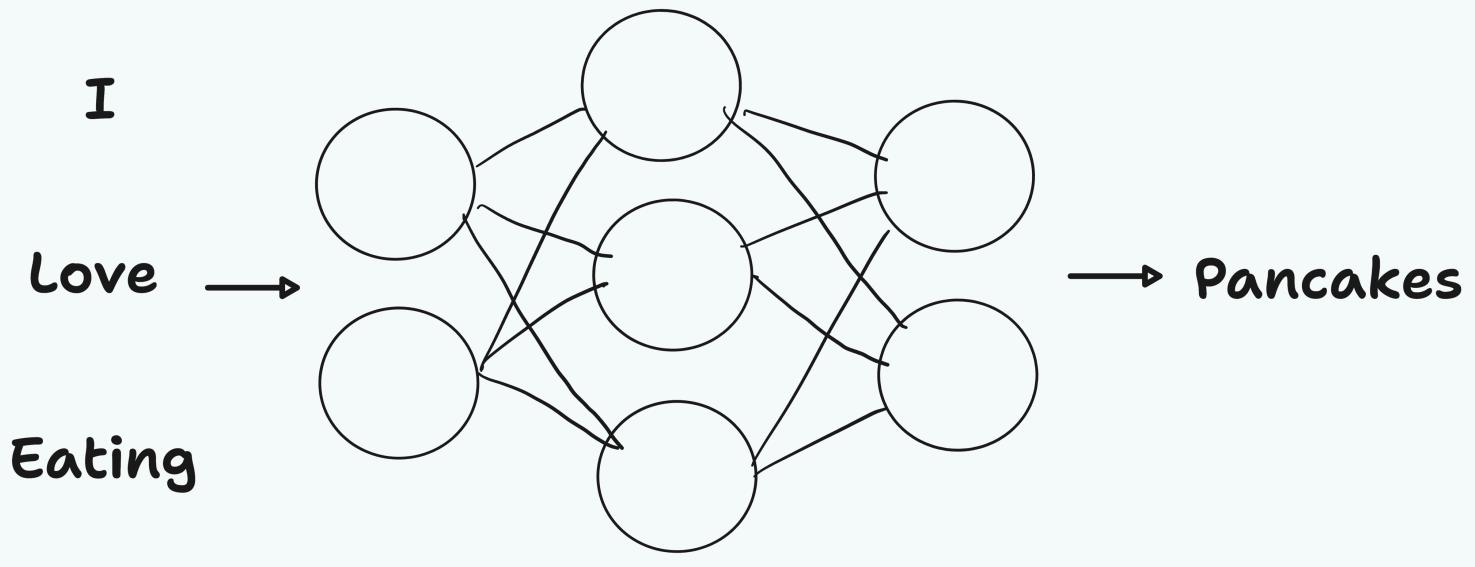


LLMs Predict the Next Word

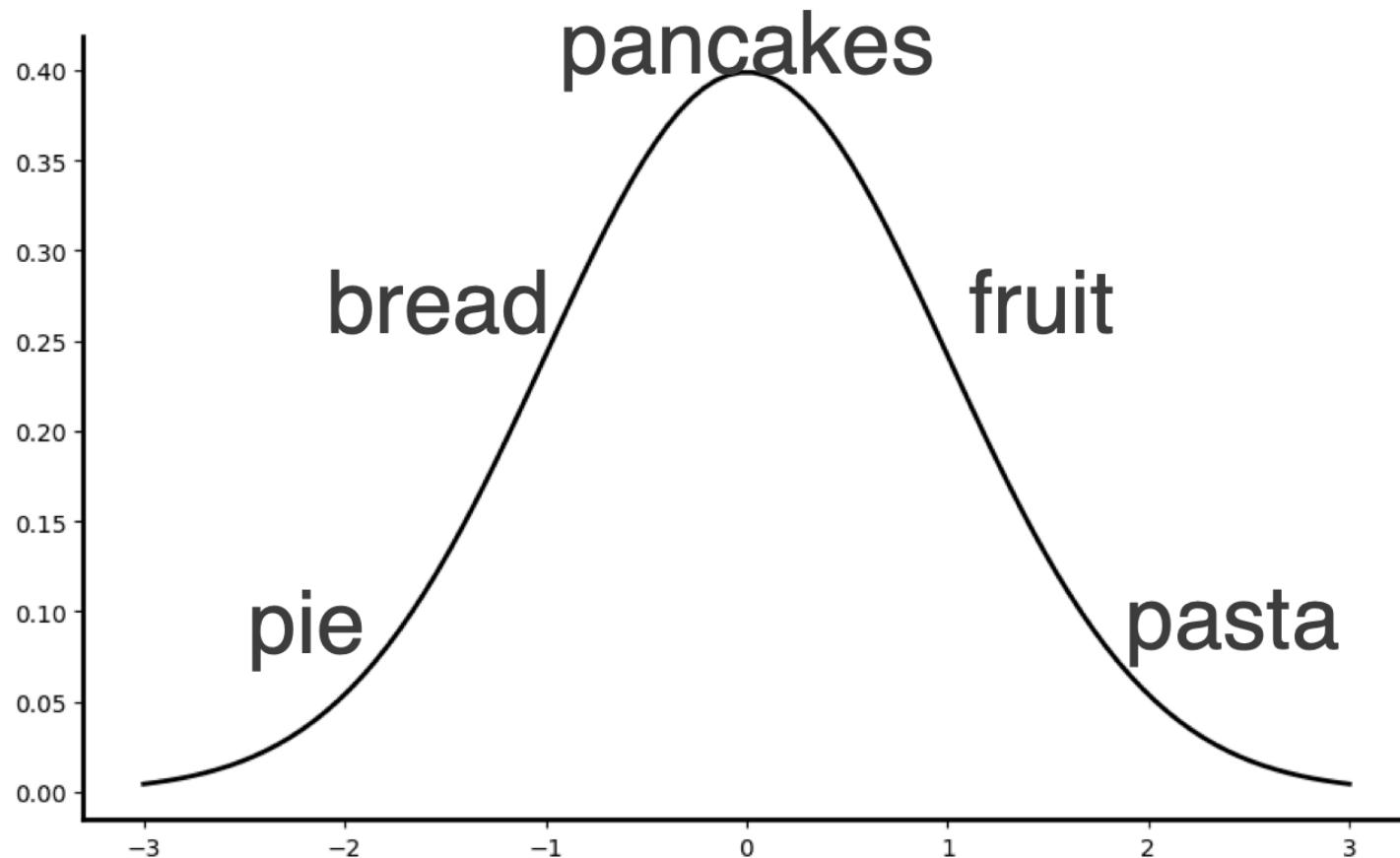


LLMs Predict the Next Word





Probability Distribution over the Next Word



Introduction to Llama3

405 B

Meet Llama 3.1

70 B

8 B

Llama3.1 Release

- LLM Released by Meta in July of 2024

Llama3.1 Release

- LLM Released by Meta in July of 2024
- Open source with a Commercial license (just like Llama2 & Llama3)

Llama3.1 Release

- LLM Released by Meta in July of 2024
- Open source with a Commercial license (just like Llama2 & Llama3)
- [Meta Llama3.1 Resources](#)

Incredible Performance in 2 sizes

Incredible Performance in 2 sizes

- Llama3.1 is OPEN SOURCE

Incredible Performance in 2 sizes

- Llama3.1 is OPEN SOURCE
- Released in 3 sizes: 8B, 70B & 405B parameters.

Incredible Performance in 2 sizes

- Llama3.1 is OPEN SOURCE
- Released in 3 sizes: 8B, 70B & 405B parameters.
- Incredible evaluation performances for all models.

Category Benchmark	Llama 3.1 8B	Gemma 2 9B IT	Mistral 7B Instruct	Llama 3.1 70B	Mixtral 8x22B Instruct	GPT 3.5 Turbo
General						
MMLU (0-shot, CoT)	73.0	72.3 (5-shot, non-CoT)	60.5	86.0	79.9	69.8
MMLU PRO (5-shot, CoT)	48.3	-	36.9	66.4	56.3	49.2
IFEval	80.4	73.6	57.6	87.5	72.7	69.9
Code						
HumanEval (0-shot)	72.6	54.3	40.2	80.5	75.6	68.0
MBPP EvalPlus (base) (0-shot)	72.8	71.7	49.5	86.0	78.6	82.0
Math						
GSMBK (8-shot, CoT)	84.5	76.7	53.2	95.1	88.2	81.6
MATH (0-shot, CoT)	51.9	44.3	13.0	68.0	54.1	43.1
Reasoning						
ARC Challenge (0-shot)	83.4	87.6	74.2	94.8	88.7	83.7
GPQA (0-shot, CoT)	32.8	-	28.8	46.7	33.3	30.8
Tool use						
BFCL	76.1	-	60.4	84.8	-	85.9
Nexus	38.5	30.0	24.7	56.7	48.5	37.2
Long context						
ZeroSCROLLS/QuALITY	81.0	-	-	90.5	-	-
InfiniteBench/En.MC	65.1	-	-	78.2	-	-
NIH/Multi-needle	98.8	-	-	97.5	-	-
Multilingual						
Multilingual MGSM (0-shot)	68.9	53.2	29.9	86.9	71.1	51.4

Llama 3.1 405B Rivals Closed Source Models

Category Benchmark	Llama 3.1 405B	Nemotron 4 340B Instruct	GPT-4 (0125)	GPT-4 Omni	Claude 3.5 Sonnet
General					
MMLU (0-shot, CoT)	88.6	78.7 (non-CoT)	85.4	88.7	88.3
MMLU PRO (5-shot, CoT)	73.3	62.7	64.8	74.0	77.0
IFEval	88.6	85.1	84.3	85.6	88.0
Code					
HumanEval (0-shot)	89.0	73.2	86.6	90.2	92.0
MBPP EvalPlus (base) (0-shot)	88.6	72.8	83.6	87.8	90.5
Math					
GSMBK (8-shot, CoT)	96.8	92.3 (0-shot)	94.2	96.1	96.4 (0-shot)
MATH (0-shot, CoT)	73.8	41.1	64.5	76.6	71.1
Reasoning					
ARC Challenge (0-shot)	96.9	94.6	96.4	96.7	96.7
GPQA (0-shot, CoT)	51.1	-	41.4	53.6	59.4
Tool use					
BFCL	88.5	86.5	88.3	80.5	90.2
Nexus	58.7	-	50.3	56.1	45.7
Long context					
ZeroSCROLLS/QuALITY	95.2	-	95.2	90.5	90.5
InfiniteBench/En.MC	83.4	-	72.1	82.5	-
NIH/Multi-needle	98.1	-	100.0	100.0	90.8
Multilingual					
Multilingual MGSM (0-shot)	91.6	-	85.9	90.5	91.6

Llama3.1 Technical Details

- Data: Trained on over 15 trillion tokens of text data

Llama3.1 Technical Details

- Data: Trained on over 15 trillion tokens of text data
- Context length: 128k tokens

Llama3.1 Technical Details

- Data: Trained on over 15 trillion tokens of text data
- Context length: 128k tokens
- System level approach for responsible development

[Meta AI Llama3.1 Release Blog Post](#)

Notebook Demo - Introduction to Llama3.1

Query Your Docs Locally with Llama3.1

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



Query Your Docs Locally with Llama3.1

- Privacy concern with closed source LLMs.



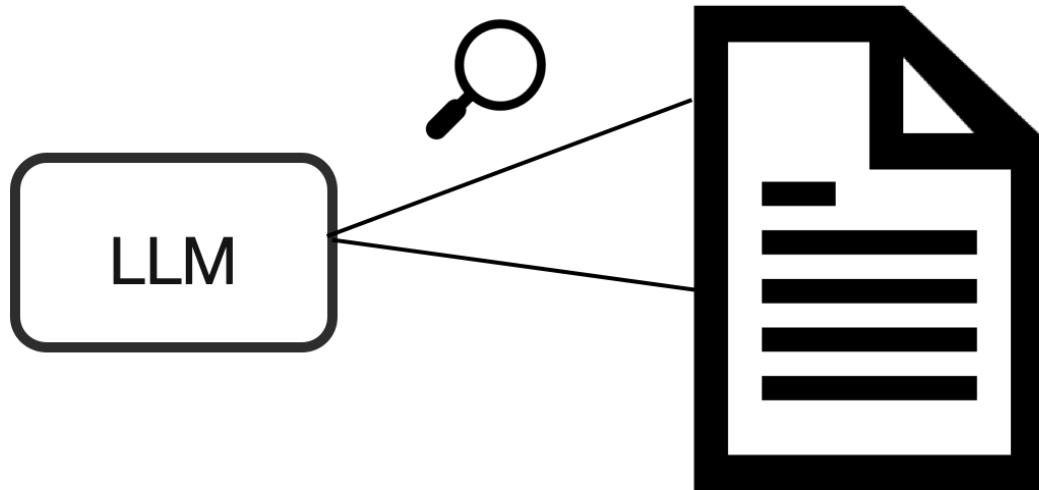
Query Your Docs Locally with Llama3.1

- Solution? Llama3.1!



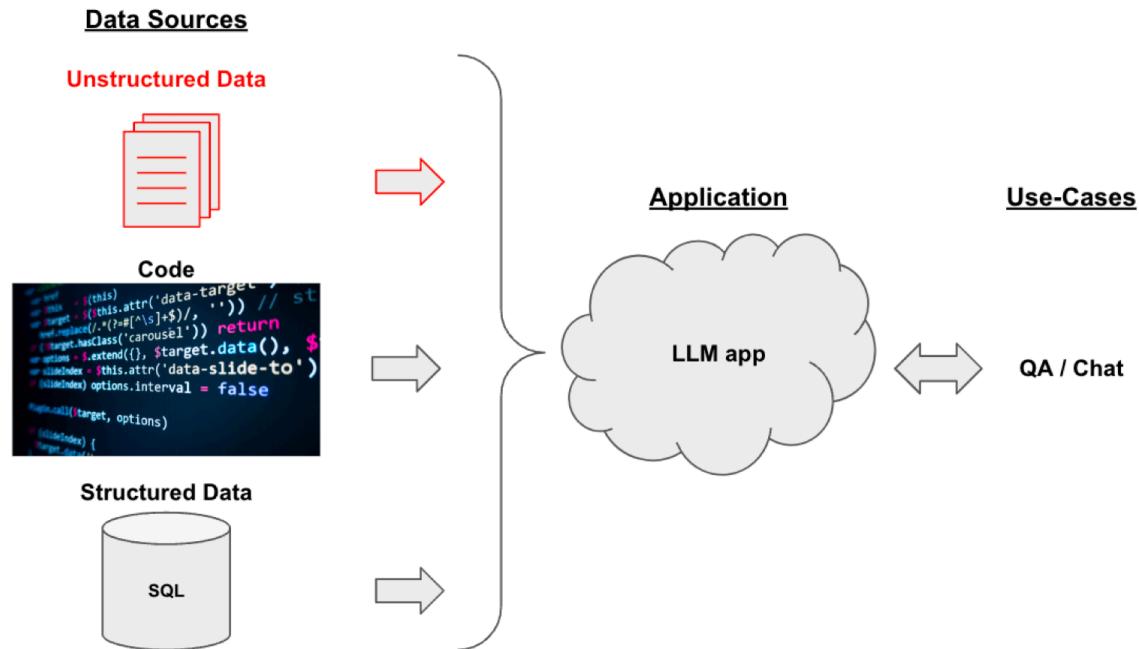
RAG with Llama3.1

- RAG - Retrieval Augmented Generation



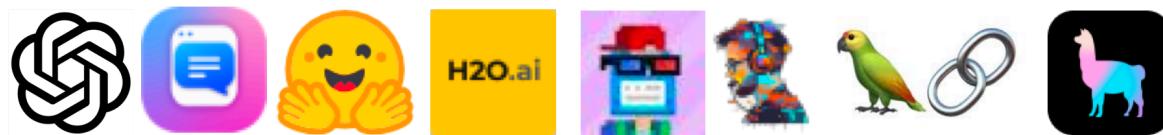
RAG with Llama3.1

- LLMs have a limited context length



Q&A RAG Tech Friction of Access

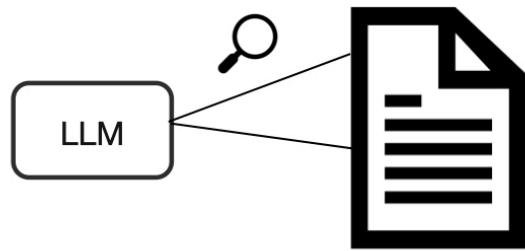
- Framework for RAG Systems
- Friction of Access



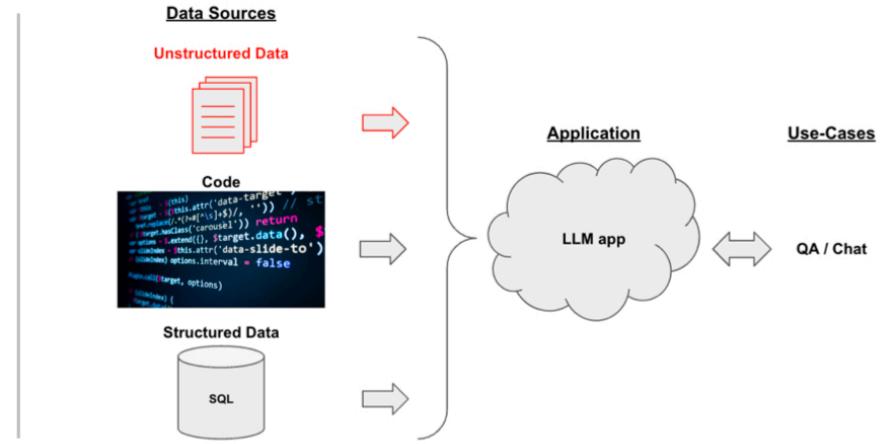
Friction of Access

[Langchain Docs](#)

RAG - Retrieval Augmented Generation



RAG - Retrieval Augmented Generation





RAG - Retrieval Augmented Generation

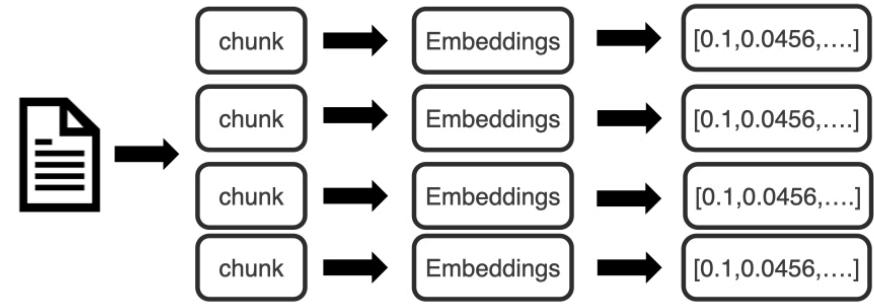


LLMs have a limited context length

RAG - Retrieval Augmented Generation

LLMs have a limited context length

Embeddings: capture content and meaning





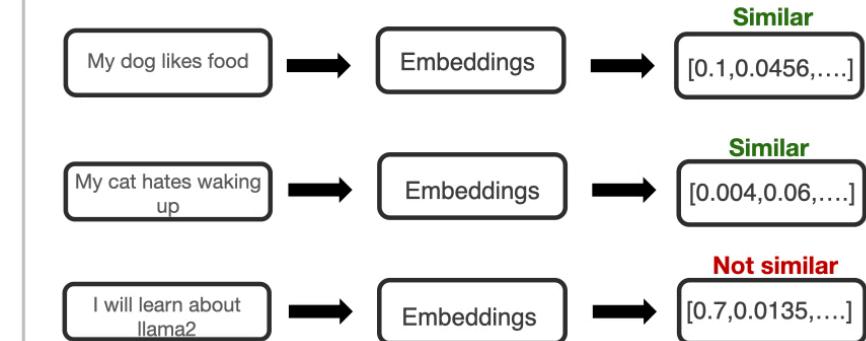
RAG - Retrieval Augmented Generation



LLMs have a limited context length



Embeddings: capture content and meaning

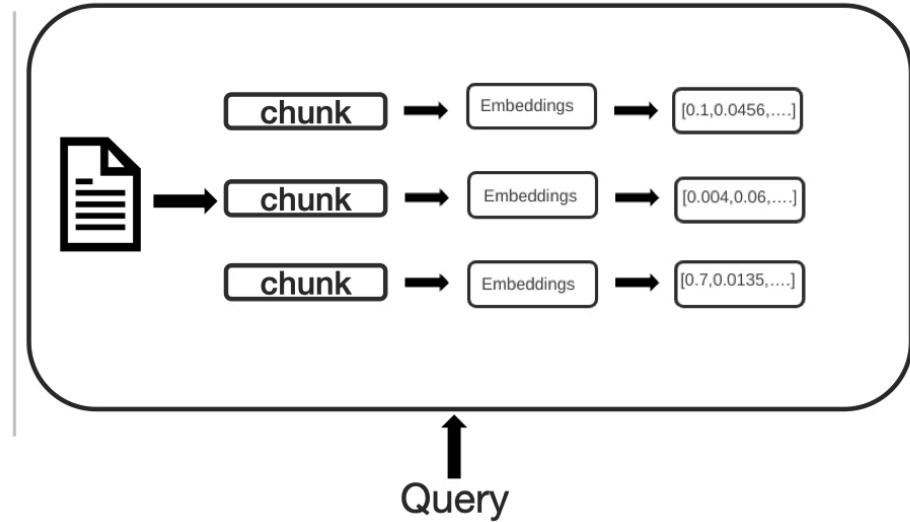


- RAG - Retrieval Augmented Generation

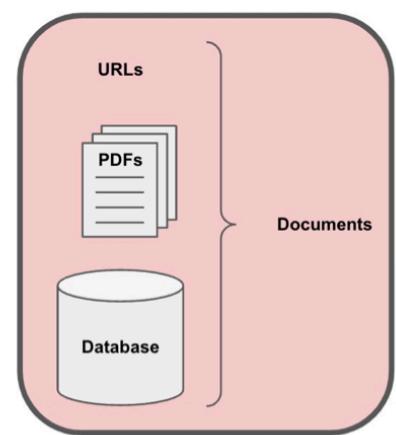
- LLMs have a limited context length

- **Embeddings:** capture content and meaning

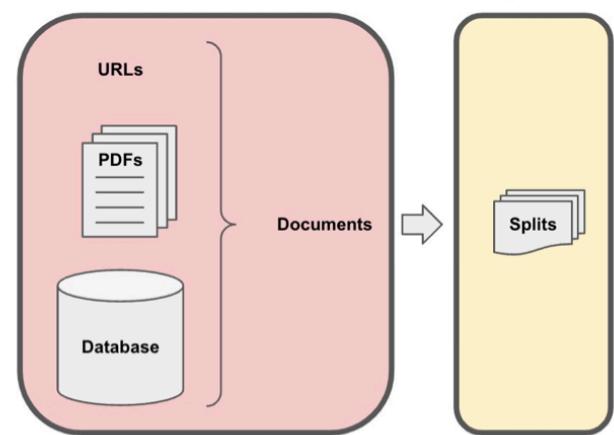
Vector Database



Document Loading

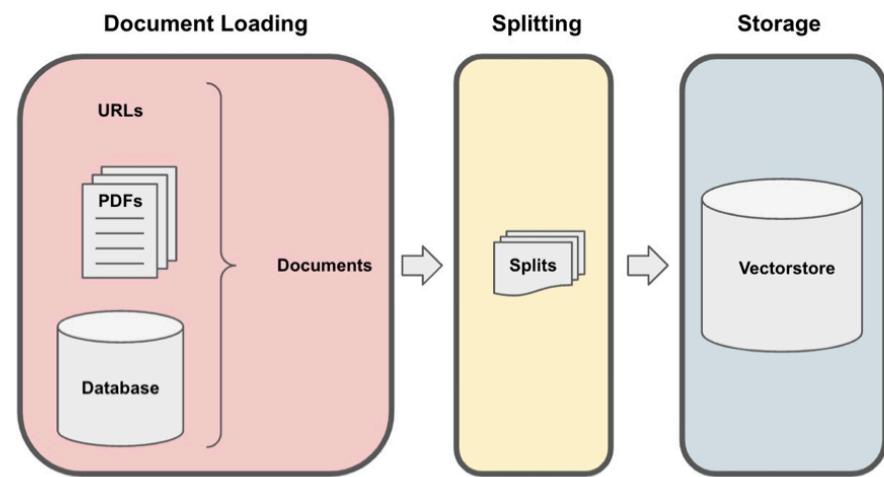


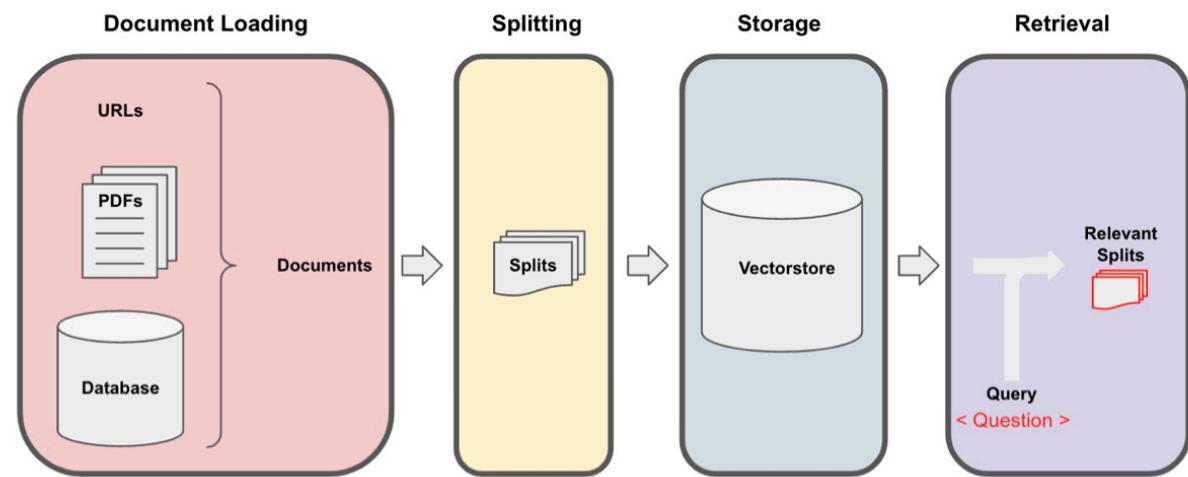
Document Loading

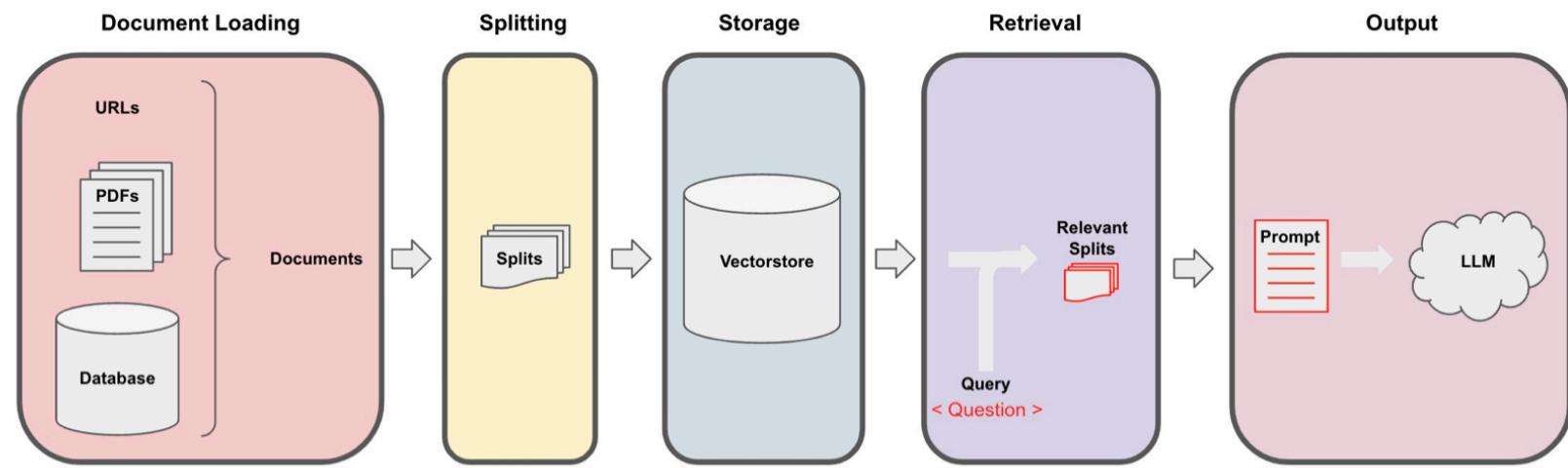


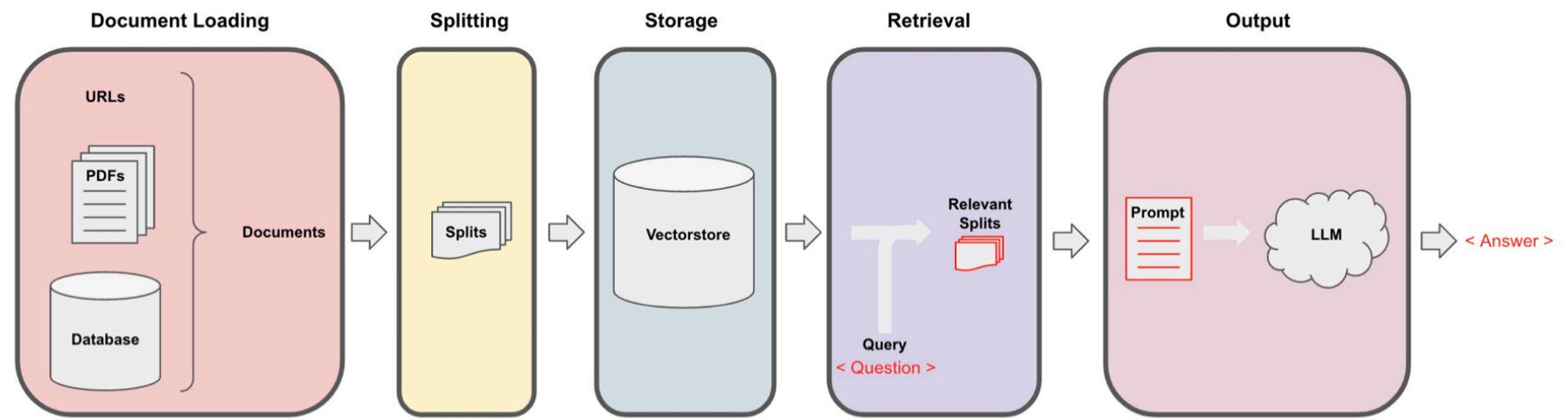
Splitting









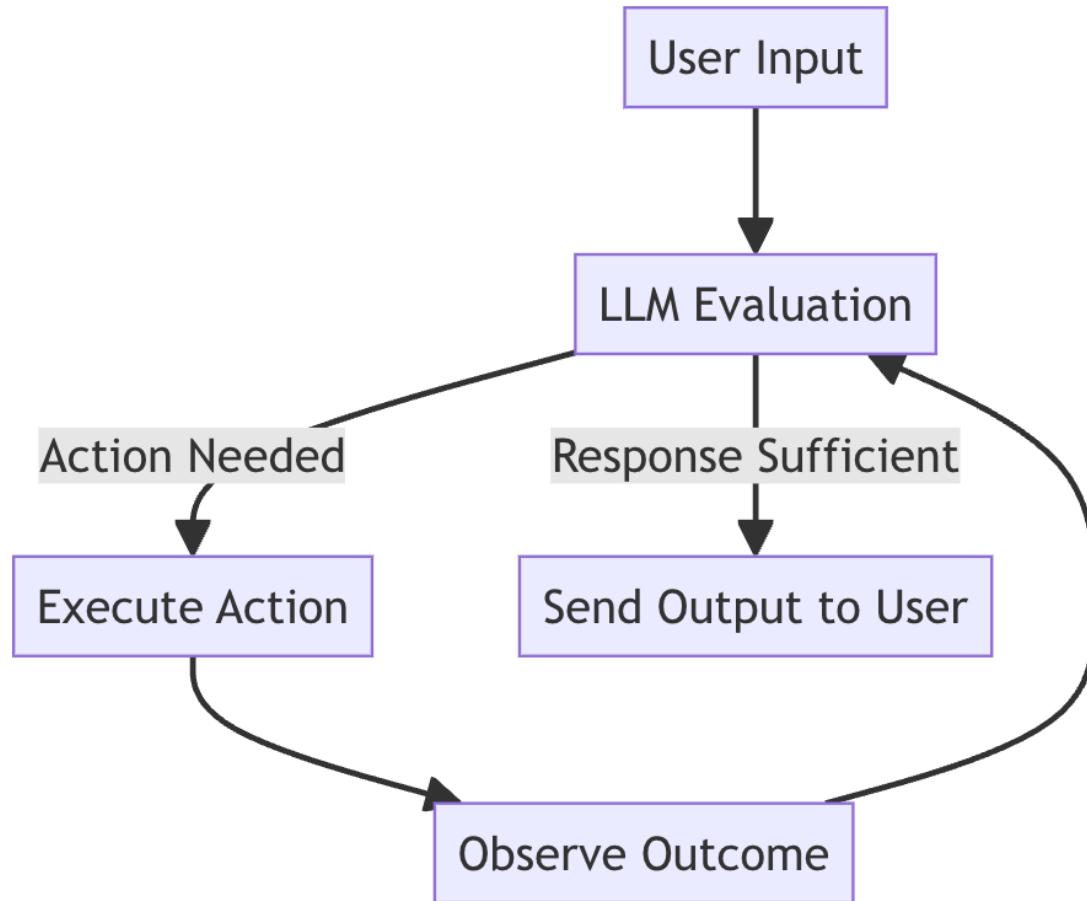


Notebook Demo - Local RAG with Llama 3.1

Q&A / Break

Local Agents with Llama 3.1

Local Agents with Llama 3.1



Explanation of the agent loop in cognitive architectures.

Practical Use Case: Customer Support Agent

Practical Use Case: Customer Support Agent

- **Scenario:** An LLM-powered customer support agent.

Practical Use Case: Customer Support Agent

- **Scenario:** An LLM-powered customer support agent.
- **User Input:** Customer asks about order status.

Practical Use Case: Customer Support Agent

- **Scenario:** An LLM-powered customer support agent.
- **User Input:** Customer asks about order status.
- **LLM Decision:** Determines if it can provide the status directly or if it needs to fetch data from the database.

Practical Use Case: Customer Support Agent

- **Scenario:** An LLM-powered customer support agent.
- **User Input:** Customer asks about order status.
- **LLM Decision:** Determines if it can provide the status directly or if it needs to fetch data from the database.
- **Action Taken:** If data fetch is needed, the agent queries the database and updates the user with the order status.

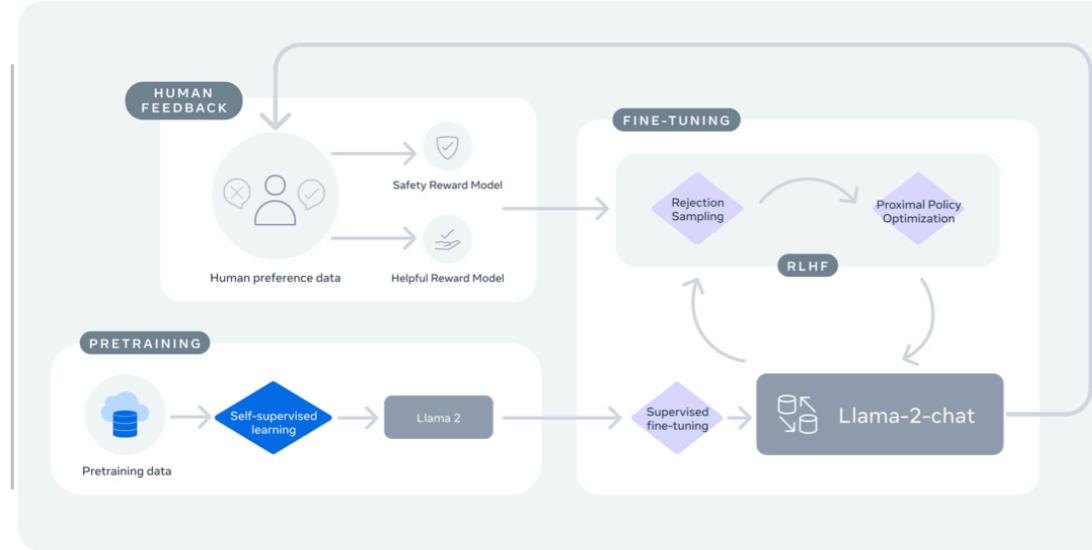
Practical use case of LLM agents in customer support.

Notebook Demo - Tool Calling and local agents with Llama 3.1

Q&A / Break

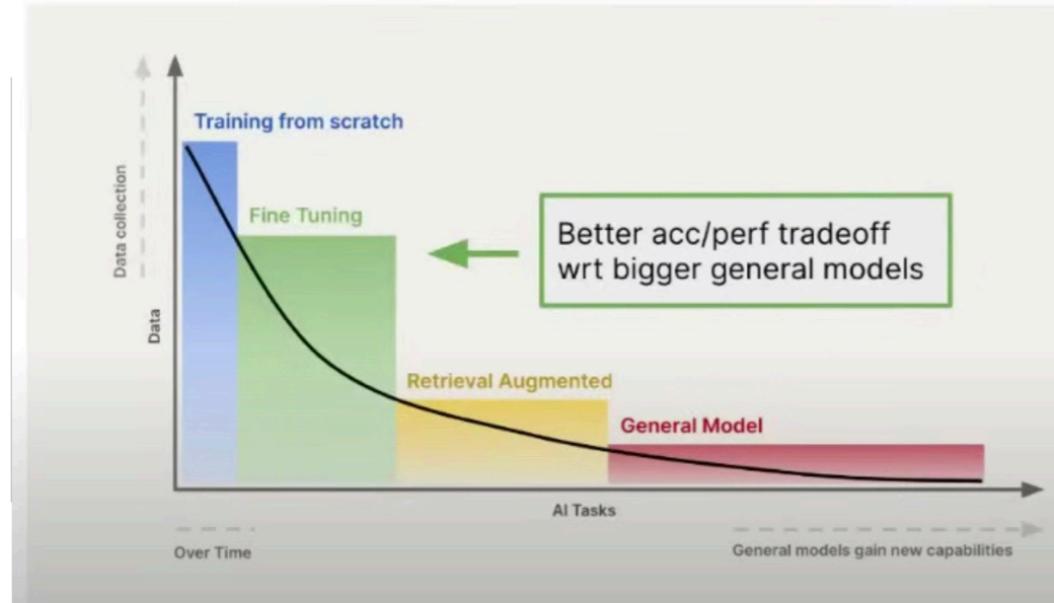
Fine Tuning Llama3.1

What is Fine Tuning?



What is Fine Tuning?

Why Fine Tune?



What is Fine Tuning?

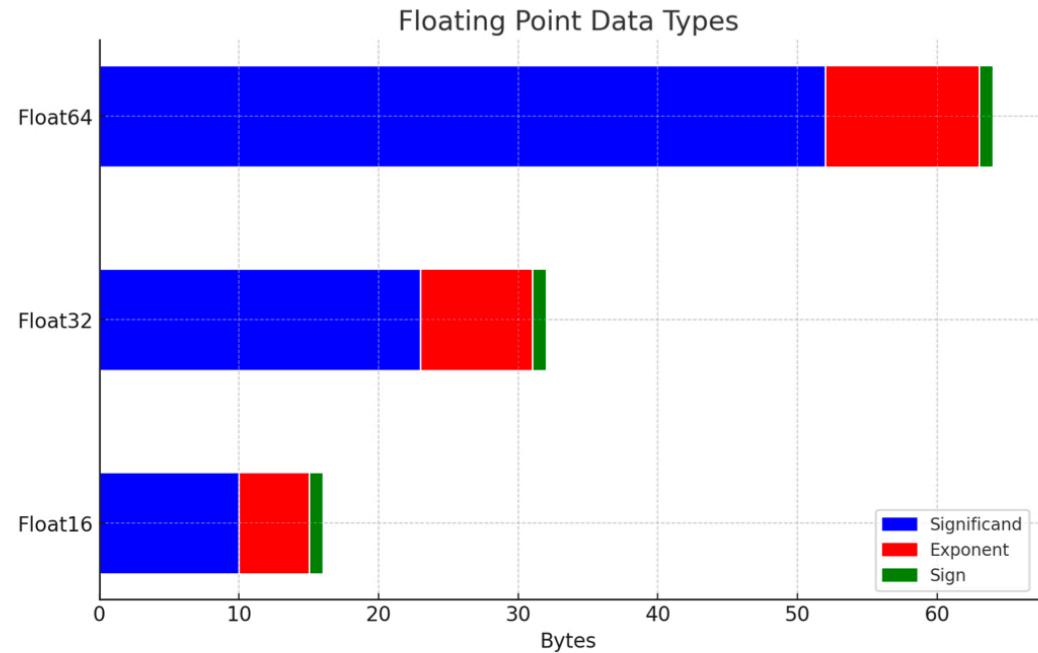
Why Fine Tune?

Memory cost of LLMs:
parameters, gradients,
optimiser states

The Memory Bottleneck: GPU comparison

GPU	Tier	\$ / hr (AWS)	VRAM (GiB)
H100	Enterprise	12.29	80
A100		5.12	80
V100		3.90	32
A10G		2.03	24
T4	Enterprise	0.98	16
RTX 4080	Consumer	N/A	16

Problem - Loading Params
Solution - Half Precision



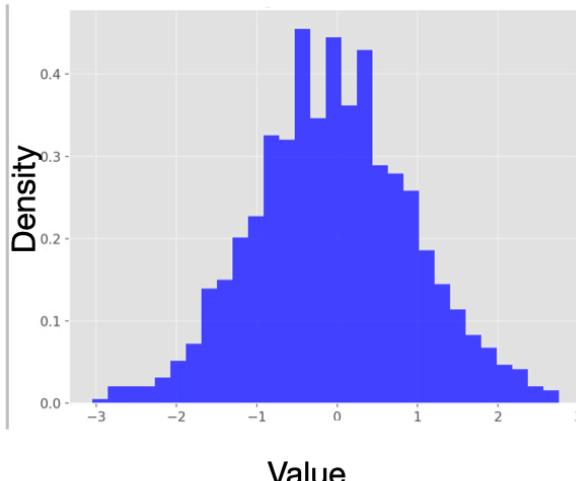
● Problem - Loading Params

Solution - Half Precision

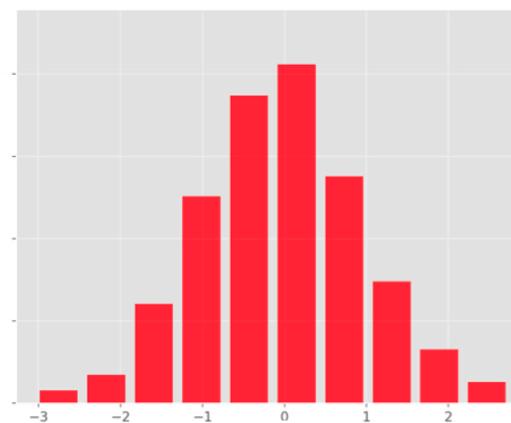
● Problem - Loading Gradients

Solution - Quantization

Original Distribution



Quantised Distribution



● Problem - Loading Params

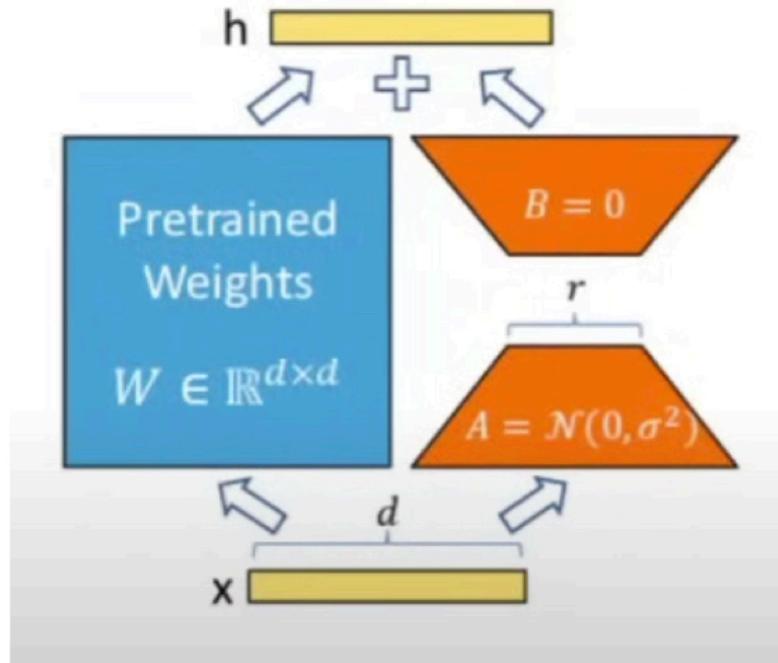
Solution - Half Precision

● Problem - Loading Gradients

Solution - Quantization

● Problem - Loading Optimizer
States

Solution - LoRA, QLora



Notebook Demo - Fine-Tuning Llama3.1 - Walkthrough