# Optimize LLMs and build generative AI applications
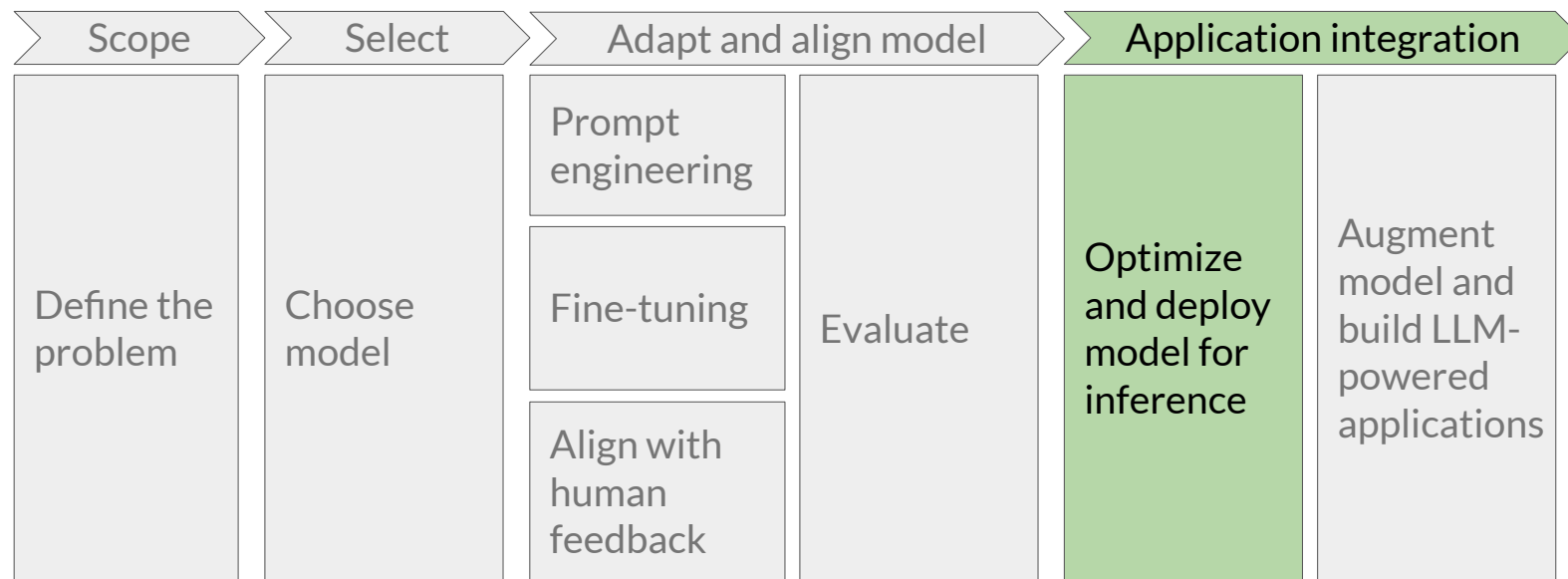
# Generative AI project lifecycle



| Scope | Select | Adapt and align model | | Application integration | |
|---|---|---|---|---|---|
| Define the problem | Choose model | Prompt engineering | Evaluate | Optimize and deploy model for inference | Augment model and build LLM-powered applications |
| | | Fine-tuning | | | |
| | | Align with human feedback | | | |

DeepLearning.AI                    aws

# Generative AI project lifecycle

# Generative AI project lifecycle

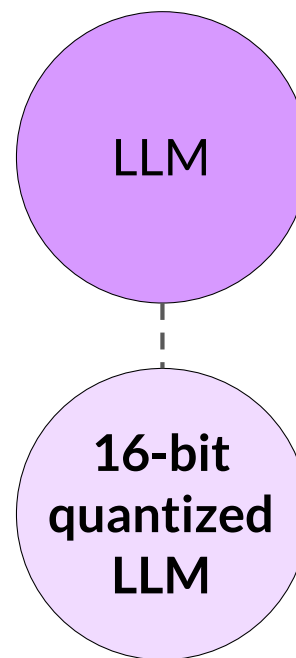| Scope | Select | Adapt and align model | | **Application integration** | |
|---|---|---|---|---|---|
| Define the problem | Choose model | Prompt engineering<br><br>Fine-tuning<br><br>Align with human feedback | Evaluate | **Optimize and deploy model for inference** | Augment model and build LLM-powered applications |

# Model optimizations to improve application performance

# LLM optimization techniques

**Distillation**

# Distillation

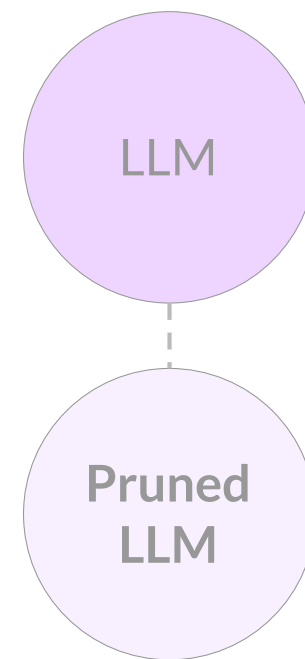Train a smaller student model from a larger teacher model

# Distillation

Train a smaller student model from a larger teacher model

# Distillation
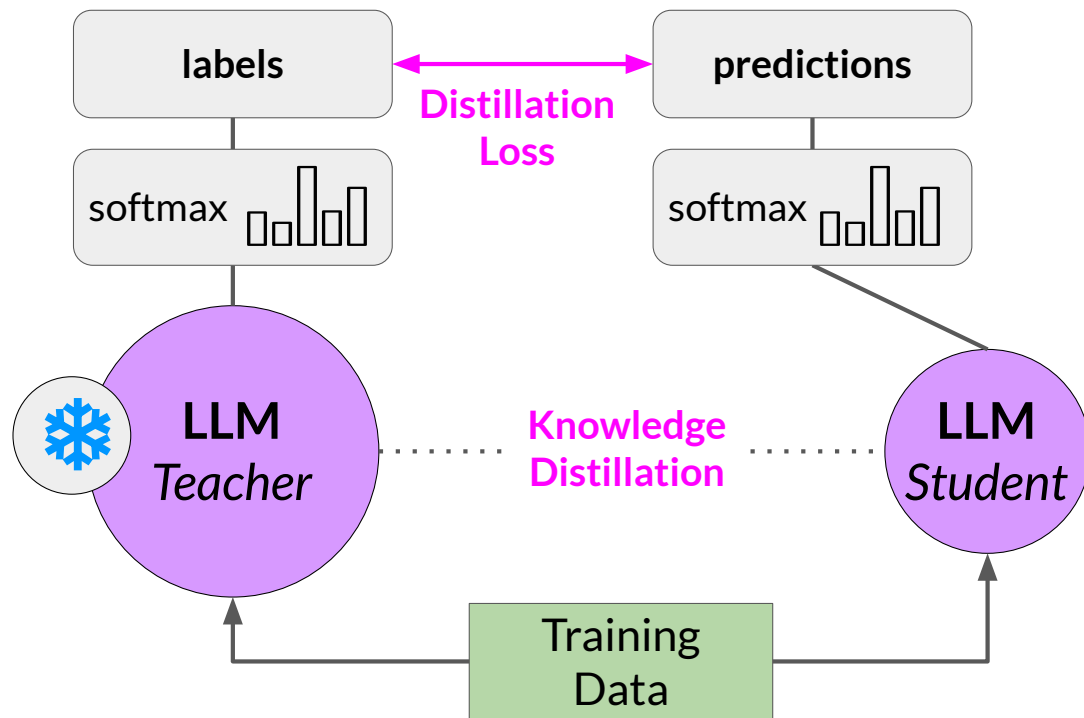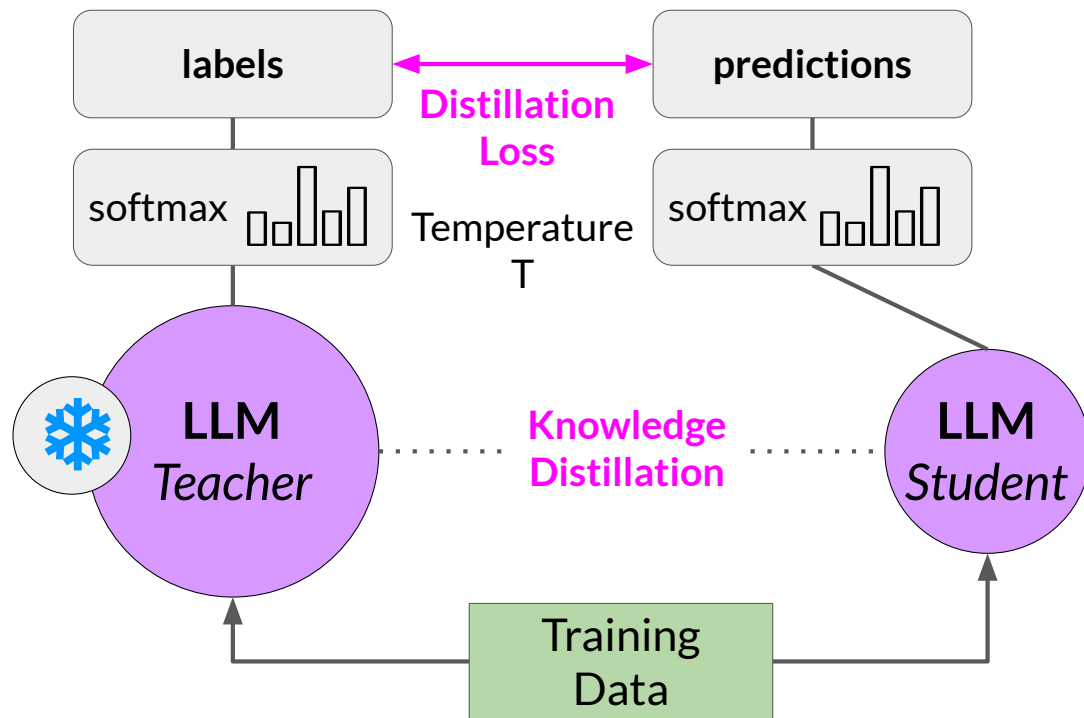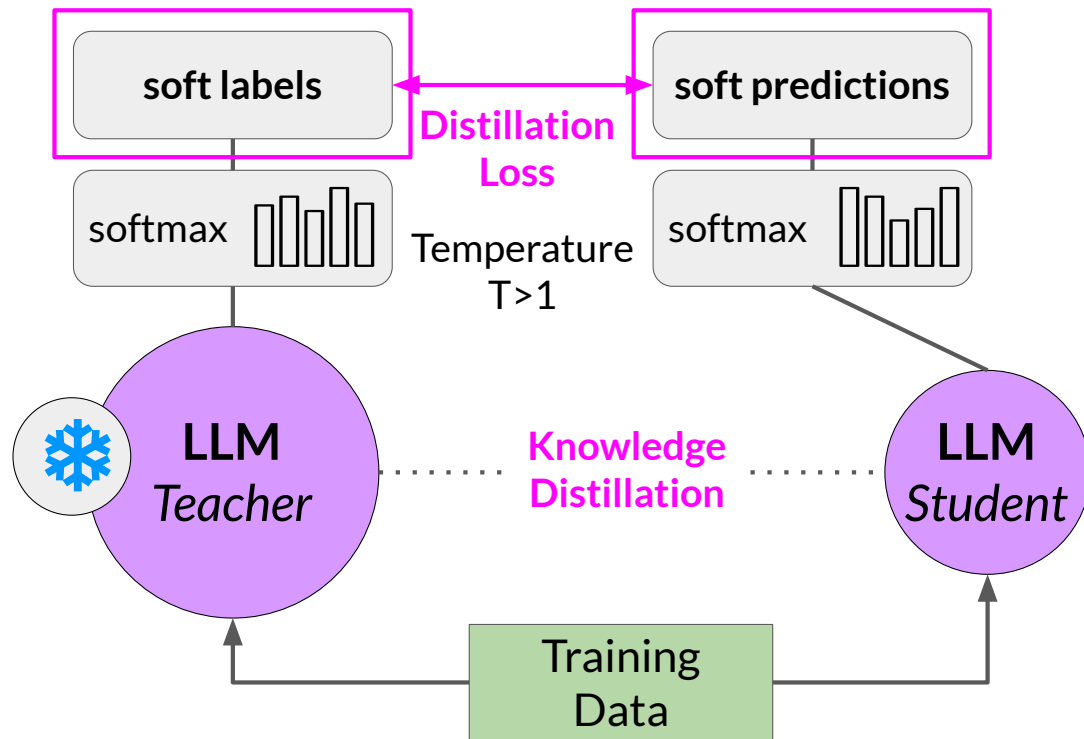
Train a smaller student model from a larger teacher model

# Distillation

Train a smaller student model from a larger teacher model

# Distillation

Train a smaller student model from a larger teacher model

# Post-Training Quantization (PTQ)

Reduce precision of model weights



LLM

8-bit quantized LLM

**MIN ~3e⁻³⁸**   0.0   **MAX ~3e³⁸**

**FP32**
32-bit floating point

- **Applied to model weights** (and/or activations)

- **Requires calibration** to capture dynamic range

?   0   ?

**FP16 | BFLOAT16 | INT8**
16-bit floating point | 8-bit integer

DeepLearning.AI                                    aws

# Pruning

Remove model weights with values close or equal to zero



- Pruning methods
  - Full model re-training
  - PEFT/LoRA
  - Post-training

- In theory, reduces model size and improves performance

- In practice, only small % in LLMs are zero-weights

# Cheat Sheet - Time and effort in the lifecycle

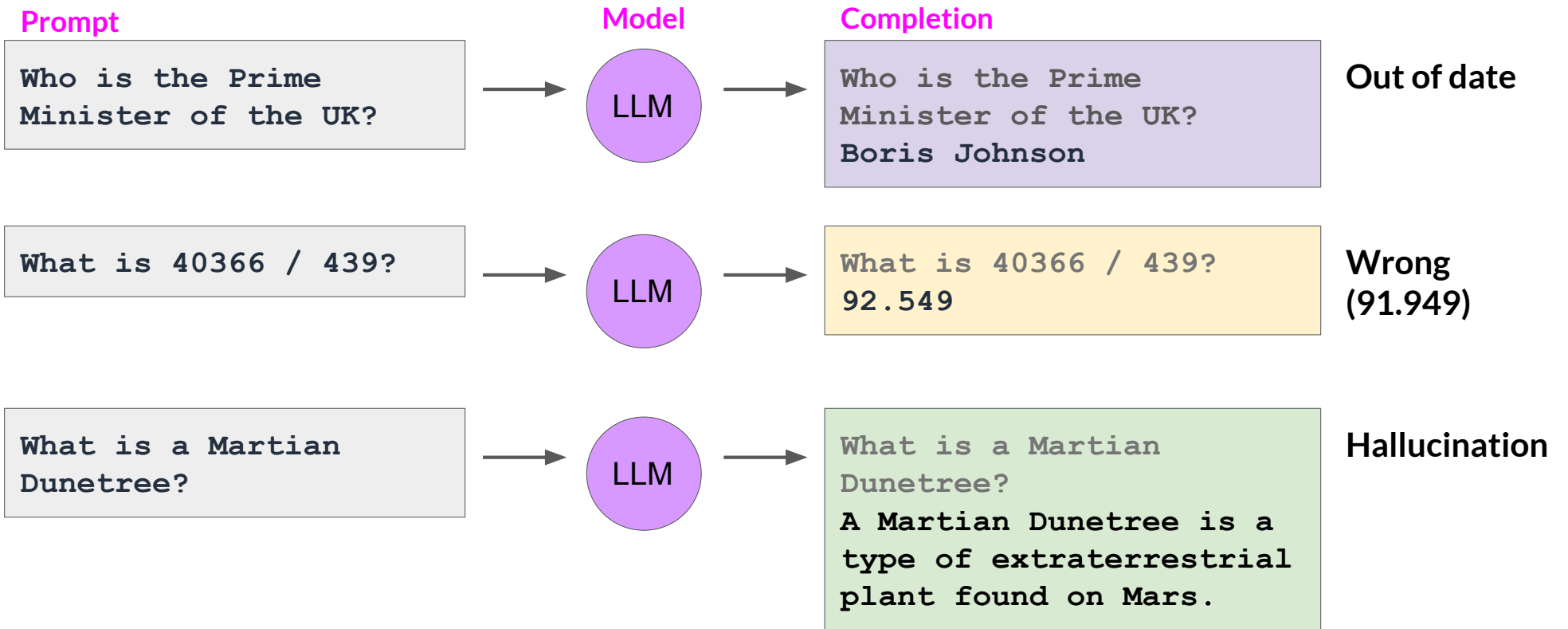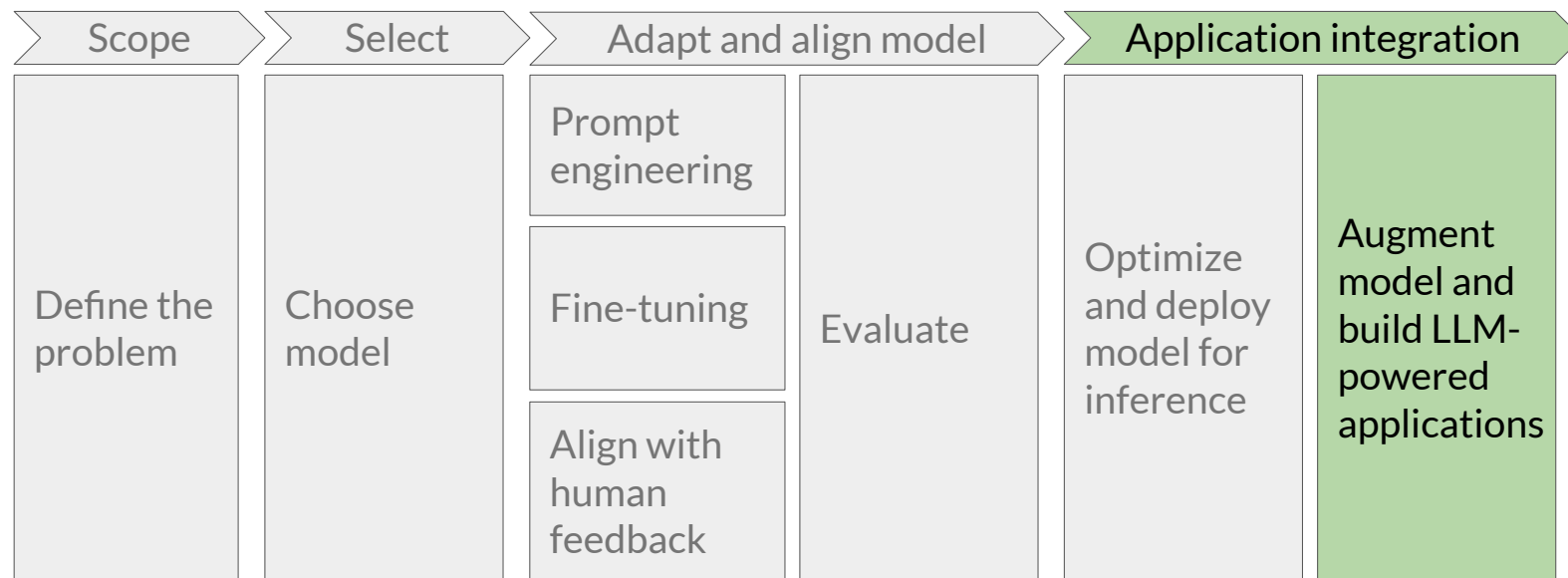| | Pre-training | Prompt engineering | Prompt tuning and fine-tuning | Reinforcement learning/human feedback | Compression/ optimization/ deployment |
|---|---|---|---|---|---|
| Training duration | Days to weeks to months | Not required | Minutes to hours | Minutes to hours similar to fine-tuning | Minutes to hours |
| Customization | Determine model architecture, size and tokenizer.<br><br>Choose vocabulary size and # of tokens for input/context<br><br>Large amount of domain training data | No model weights<br><br>Only prompt customization | Tune for specific tasks<br><br>Add domain-specific data<br><br>Update LLM model or adapter weights | Need separate reward model to align with human goals (helpful, honest, harmless)<br><br>Update LLM model or adapter weights | Reduce model size through model pruning, weight quantization, distillation<br><br>Smaller size, faster inference |
| Objective | Next-token prediction | Increase task performance | Increase task performance | Increase alignment with human preferences | Increase inference performance |
| Expertise | High | Low | Medium | Medium-High | Medium |

# Using the LLM in applications

# Models having difficulty

**Prompt:**
```
Who is the Prime
Minister of the UK?
```

LLM

**Completion:**
```
Who is the Prime
Minister of the UK?
Boris Johnson
```

**Out of date**

**Prompt:**
```
What is 40366 / 439?
```

LLM

**Completion:**
```
What is 40366 / 439?
92.549
```

**Wrong
(91.949)**

**Prompt:**
```
What is a Martian
Dunetree?
```

LLM

**Completion:**
```
What is a Martian
Dunetree?
A Martian Dunetree is a
type of extraterrestrial
plant found on Mars.
```

**Hallucination**

DeepLearning.AI

aws

# Generative AI project lifecycle

# LLM-powered applications

# LLM-powered applications

# Knowledge cut-offs in LLMs

**Prompt**

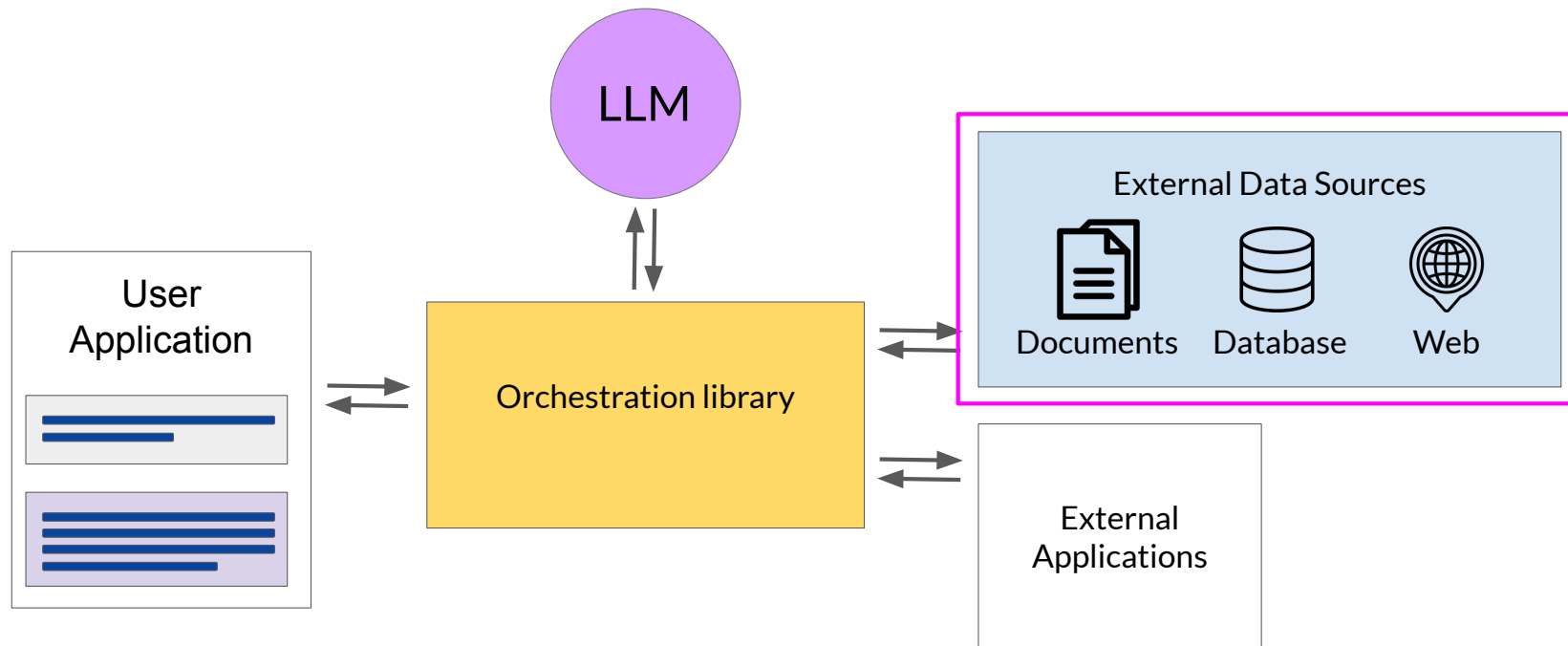Who is the
current Prime
Minister of the
United Kingdom?

**Model**

LLM

**Completion**

Who is the
current Prime
Minister of the
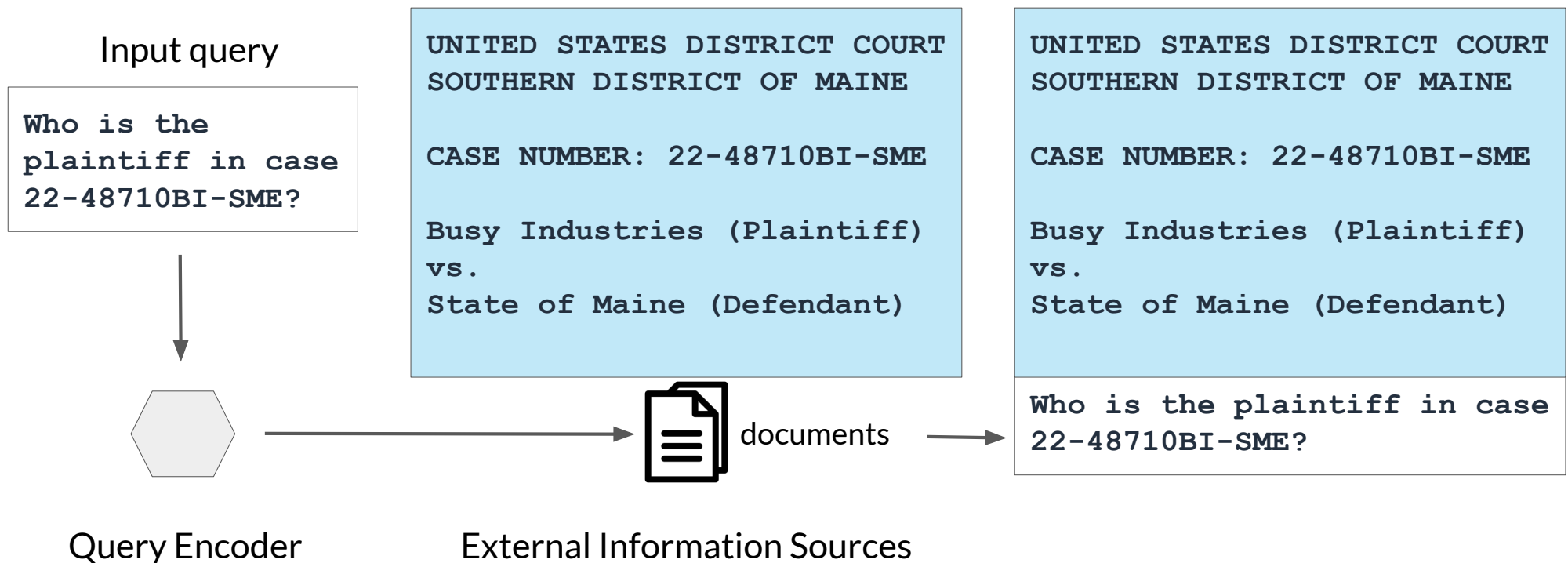United Kingdom?

Boris Johnson

DeepLearning.AI

aws

# LLM-powered applications

# Retrieval Augmented Generation (RAG)



Retriever

Lewis et al. 2020 "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"

DeepLearning.AI    aws

# Example: Searching legal documents

Input query

```
Who is the
plaintiff in case
22-48710BI-SME?
```

Query Encoder

External Information Sources

documents

```
UNITED STATES DISTRICT COURT
SOUTHERN DISTRICT OF MAINE


CASE NUMBER: 22-48710BI-SME


Busy Industries (Plaintiff)
vs.
State of Maine (Defendant)
```

```
UNITED STATES DISTRICT COURT
SOUTHERN DISTRICT OF MAINE


CASE NUMBER: 22-48710BI-SME


Busy Industries (Plaintiff)
vs.
State of Maine (Defendant)
```

```
Who is the plaintiff in case
22-48710BI-SME?
```

# Example: Searching legal documents

# RAG integrates with many types of data sources



Query encoder → External information sources
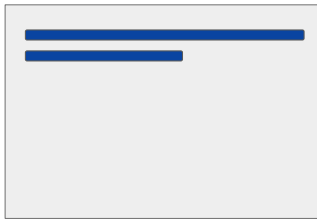
Retriever

External Information Sources
- Documents
- Wikis
- Expert Systems
- Web pages
- Databases
- Vector Store

# Data preparation for vector store for RAG

Two considerations for using external data in RAG:

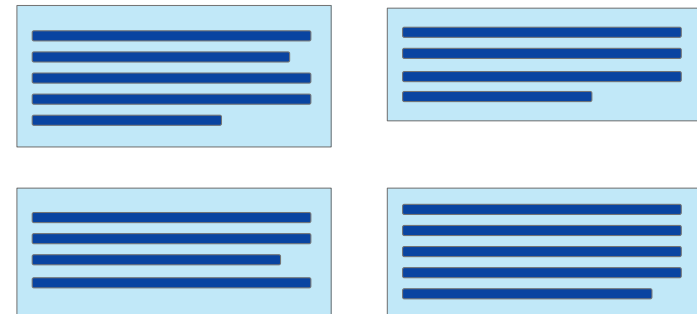1. Data must fit inside context window

Prompt context limit few 1000 tokens

Single document too large to fit in window

Split long sources into short chunks

# Data preparation for RAG

Two considerations for using external data in RAG:

1. Data must fit inside context window
2. Data must be in format that allows its relevance to be assessed at inference time: **Embedding vectors**

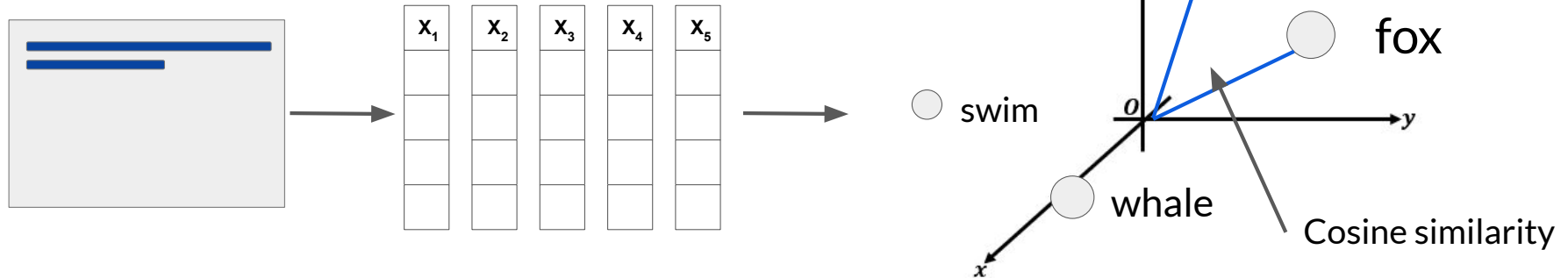Prompt text converted to embedding vectors

# Data preparation for RAG
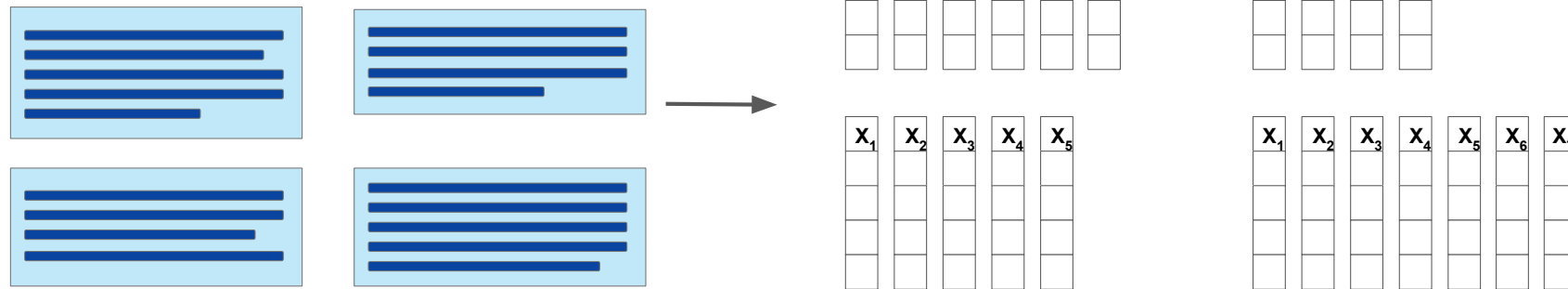
Two considerations for using external data in RAG:

1. Data must fit inside context window
2. Data must be in format that allows its relevance to be assessed at inference time: **Embedding vectors**

Process each chunk with LLM
to produce embedding vectors

# Vector database search



- Each text in vector store is identified by a key
- Enables a **citation** to be included in completion

# Enabling interactions with external applications

# Having an LLM initiate a clothing return

# Having an LLM initiate a clothing return

Lookup with RAG

**ShopBot**

Ok, I've found your order.

Do you want to return any other items from that order?

No, only the jeans.

API call

Ok great. Let me get a return label from our shipping partner.  Can you remind me of the email you used to order?

# Having an LLM initiate a clothing return



API call to the shipper

**ShopBot**

Sure, its tim.b@email.net

Thank you! I've just sent the shipping label to your email address. Please return your jeans within the next 5 days.

Great, thank you!

You're welcome!

DeepLearning.AI

aws

# LLM-powered applications

# Requirements for using LLMs to power applications

**Plan actions**

Steps to process return:
**Step 1:** Check order ID
**Step 2:** Request label
**Step 3:** Verify user email
**Step 4:** Email user label

**Format outputs**

SQL Query:
**SELECT COUNT(*)**
**FROM orders**
**WHERE order_id = 21104**

**Validate actions**

Collect required user information and make sure it is in the completion

User email:
tim.b@email.net

Prompt structure is important!

Helping LLMs reason and
plan with Chain-of-Thought
Prompting

# LLMs can struggle with complex reasoning problems

**Prompt**

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**

LLM

**Completion**

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?

A: The answer is 27. ❌

# Humans take a step-by-step approach to solving complex problems

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?

Start: Roger started with 5 balls.
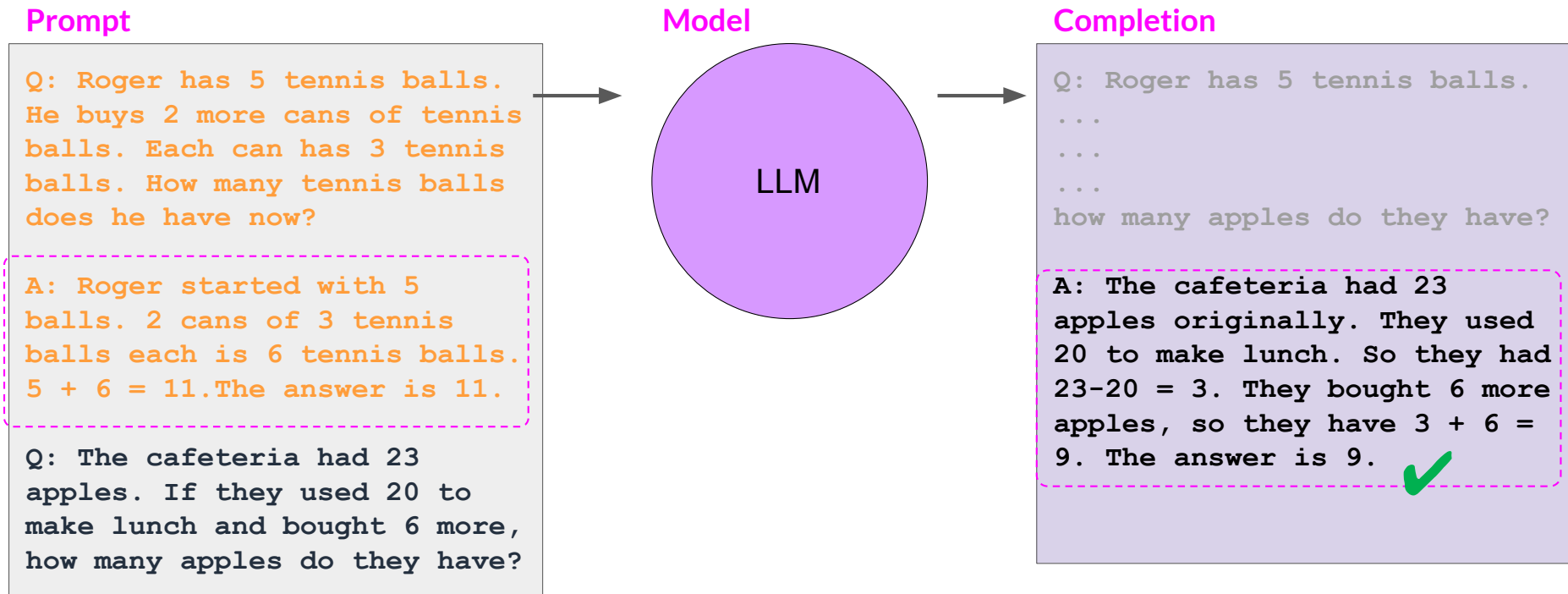Step 1: 2 cans of 3 tennis balls each is 6 tennis balls.
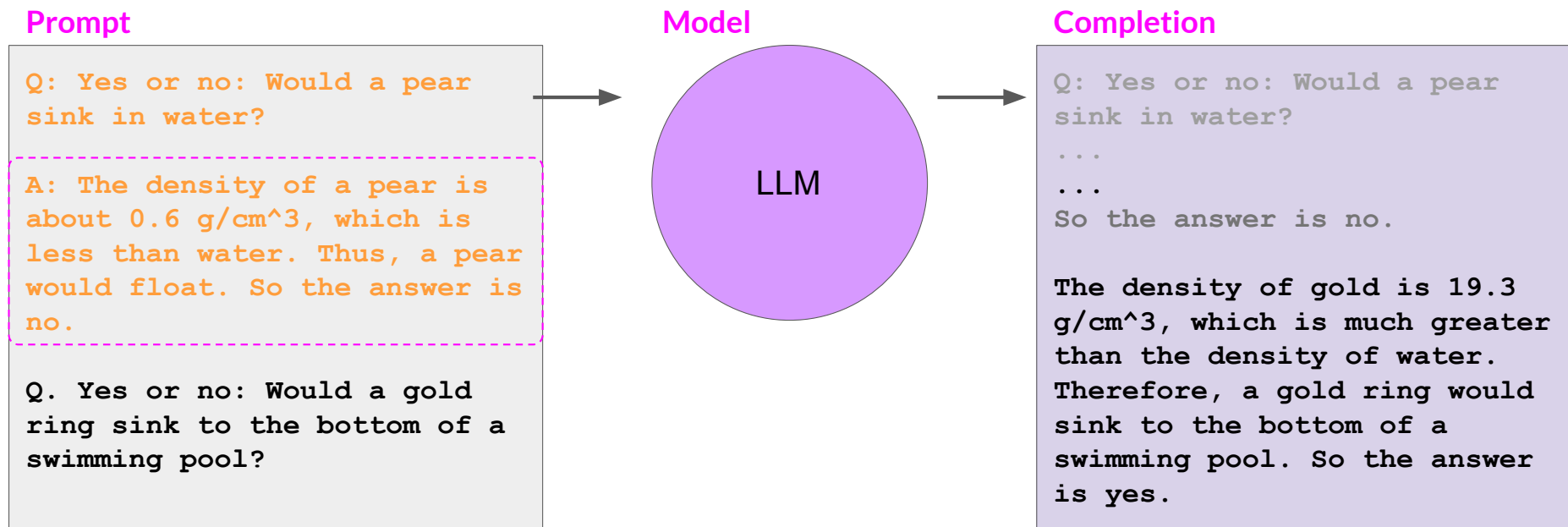Step 2: 5 + 6 = 11
End: The answer is 11

Reasoning steps

"Chain of thought"

# Chain-of-Thought Prompting can help LLMs reason

**Prompt**

**Model**

**Completion**

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?

LLM

Q: Roger has 5 tennis balls.
...
...
...
how many apples do they have?

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. ✔

Source: Wei et al. 2022, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models"

DeepLearning.AI

aws

# Chain-of-Thought Prompting can help LLMs reason

**Prompt**

Q: Yes or no: Would a pear sink in water?

A: The density of a pear is about 0.6 g/cm^3, which is less than water. Thus, a pear would float. So the answer is no.

Q. Yes or no: Would a gold ring sink to the bottom of a swimming pool?

**Model**

LLM

**Completion**

Q: Yes or no: Would a pear sink in water?
...
...
So the answer is no.

The density of gold is 19.3 g/cm^3, which is much greater than the density of water. Therefore, a gold ring would sink to the bottom of a swimming pool. So the answer is yes.

Source: Wei et al. 2022, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models"

DeepLearning.AI

aws

# Program-aided Language Models

# LLMs can struggle with mathematics

```
What is 40366 / 439?
```

LLM

```
What is 40366 / 439?
92.549
```

# Program-aided language (PAL) models



Source: Gao et al. 2022, "PAL: Program-aided Language Models"

# PAL example

**Prompt with one-shot example**

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?

Answer:
```
# Roger started with 5 tennis balls
tennis_balls = 5
# 2 cans of tennis balls each is
bought_balls = 2 * 3
# tennis balls. The answer is
answer = tennis_balls + bought_balls
```

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves did they have left?

# PAL example

**Prompt with one-shot example**

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?

Answer:
```
# Roger started with 5 tennis balls
tennis_balls = 5
# 2 cans of tennis balls each is
bought_balls = 2 * 3
# tennis balls. The answer is
answer = tennis_balls + bought_balls
```

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves did they have left?

**Completion, CoT reasoning (blue) , and PAL execution (pink)**

Answer:
```
# The bakers started with 200 loaves
loaves_baked = 200          ←
# They sold 93 in the morning and 39 in the
afternoon
loaves_sold_morning = 93    ←
loaves_sold_afternoon = 39  ←
# The grocery store returned 6 loaves.
loaves_returned = 6         ←
# The answer is
answer = loaves_baked
   - loaves_sold_morning
   - loaves_sold_afternoon
   + loaves_returned
```

# Program-aided language (PAL) models



PAL prompt template

question

PAL formatted prompt

```
def solution:

return answer
```

Python script

Python interpreter

answer = 74

LLM