

DSPy: Not Your Average Prompt Engineering

WHY DONT YOU LIKE PROMPT ENGINEERING

- It is simple, effective and cost-efficient.

→ At almost no cost on GPU

→ Take you 5 mins to validate the effectiveness via LLM UI/API

→ Most tricks can be explained in one or few-sentences, instead of 8-page

- It is brittle and lacks of a systematic way to improve it.

DSPy:
Not Your Average Prompt Engineering

Basic Prompt Engineering Modules

ZERO-SHOT

```

Zero-shot

1 Text: Oh, great. Another meeting. Just what I needed to make my day even more exciting.
2 Sentiment:
  
```



FEW-SHOT

Few-shot examples

```

1 Text: Wow, thanks for leaving your dirty dishes in the sink for the tenth time this week. It really adds to the ambiance of our kitchen.
2 Sentiment: negative
3 Text: Oh, wonderful. Another software update. I'm thrilled to see what new bugs and glitches it brings.
4 Sentiment: negative
5 Text: Wow, you're so good at interrupting people. I wish I had that skill too.
6 Sentiment: negative
7 Text: Oh, fantastic! The traffic is absolutely amazing today. I can't wait to spend another hour sitting in my car.
8 Sentiment: negative
9 Text: I had an amazing time on my vacation. The sights, the food, and the people were all incredible.
10 Sentiment:

```

Training data

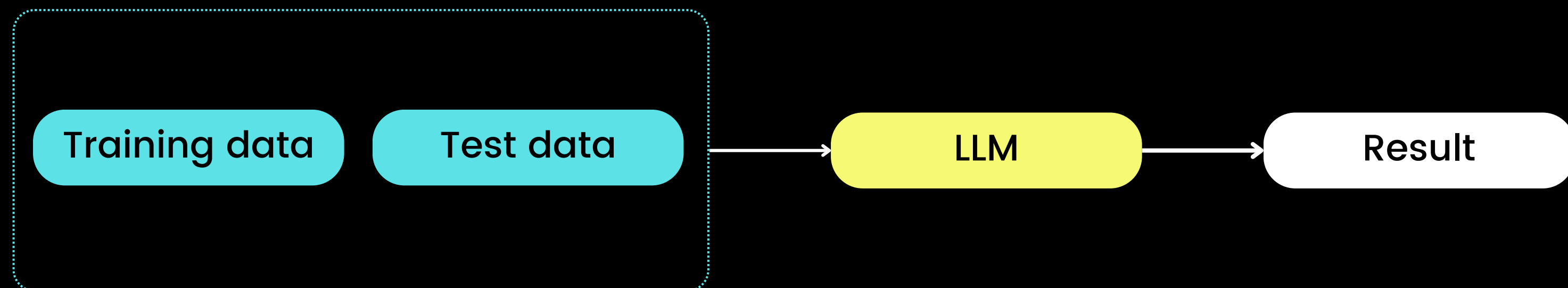
Test data

LLM

Result

FEW-SHOT

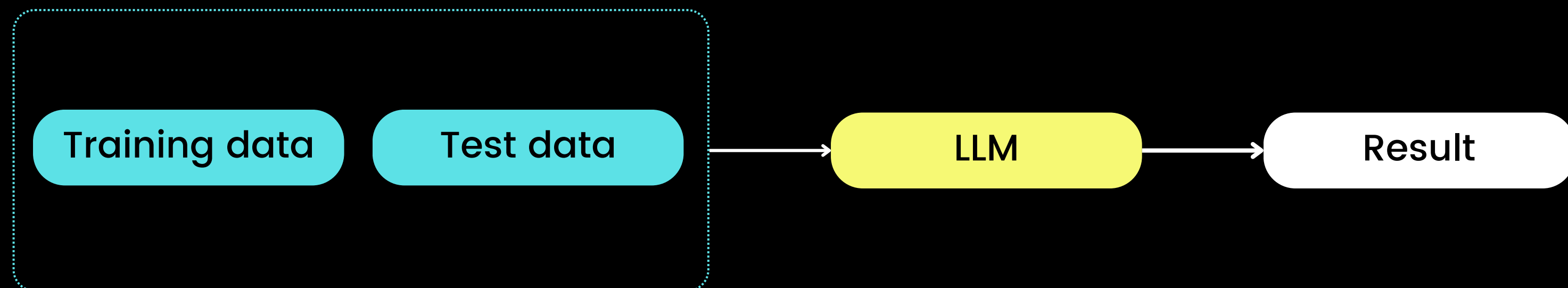
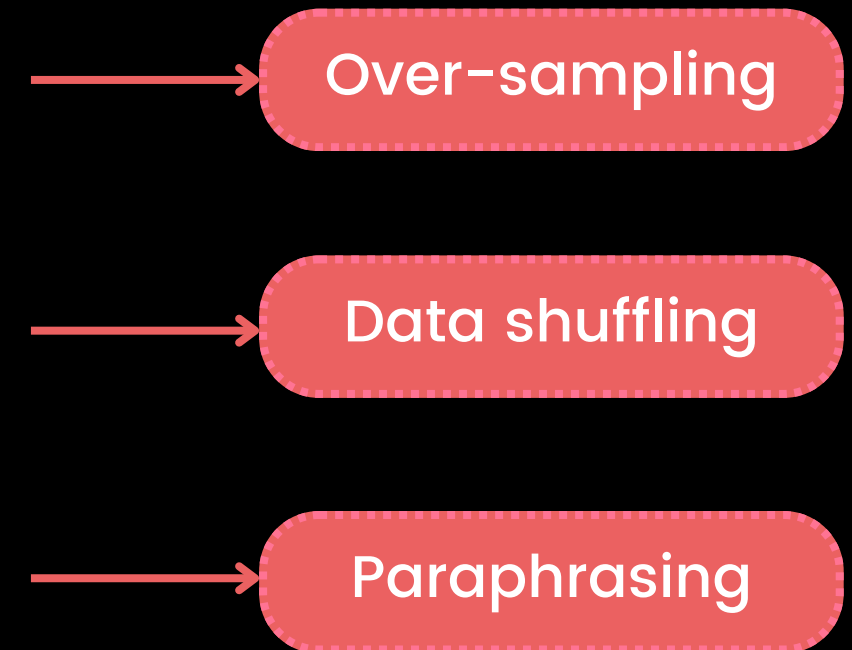
- Presents a set of high-quality demonstrations, each consisting of both input and desired output, on the target task.
- Performance is influenced by
 - Training examples, and the order of the examples
 - Example string template



FEW-SHOT BIASES

- **Majority label bias** exists if distribution of labels among the examples is unbalanced;
- **Recency bias** refers to the tendency where the model may repeat the label at the end;
- **Common token bias** indicates that LLM tends to produce common tokens more often than rare tokens.

In classic ML

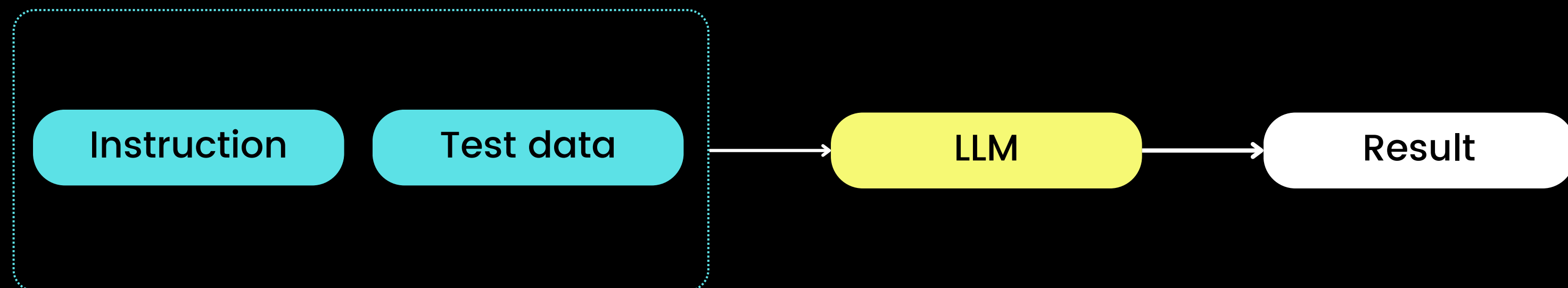


INSTRUCTION PROMPT



Instruction prompt

```
1 Please label sentiment of the sentence below into "positive", "negative" and "neutral".  
2 Text: I had an amazing time on my vacation. The sights, the food, and the people were  
   all incredible.  
3 Sentiment:
```

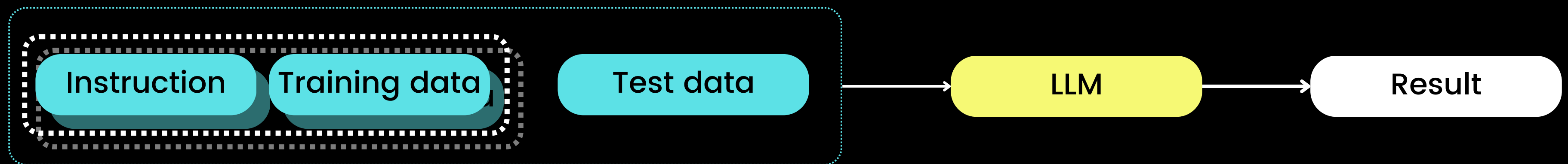


INSTRUCTION PROMPT

- Explicitly telling model what to do, instead of showing a set of demonstrations (i.e. few-shot) and let model immitate.



Instruction few-shot prompting



Few-shot instruction prompting

CHAIN-OF-THOUGHT (COT)

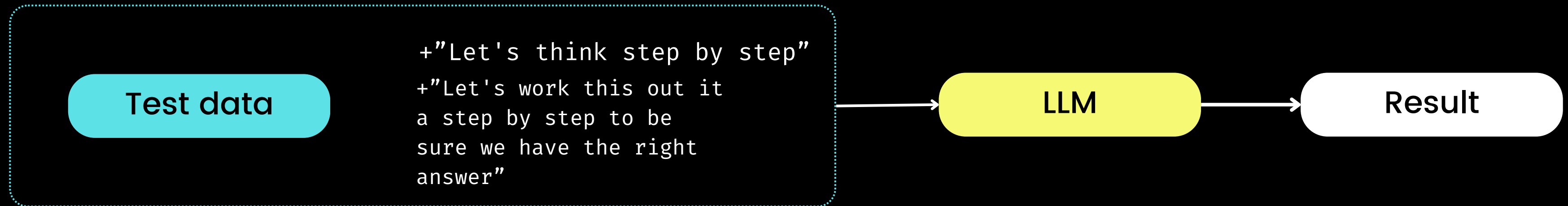


Chain of thoughts prompt

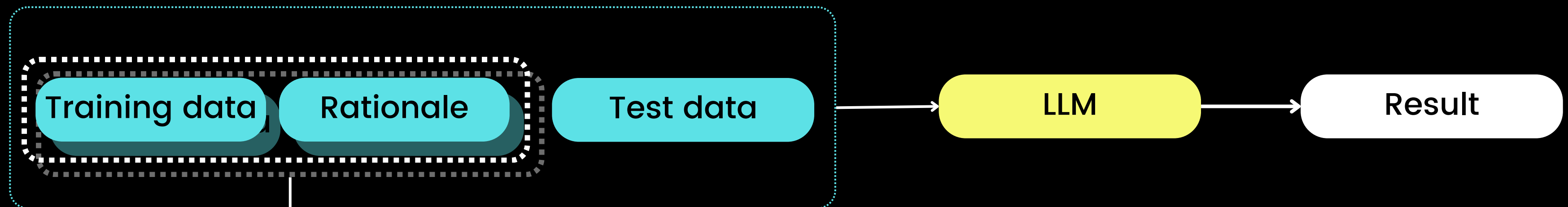
- 1 Question: Tom and Elizabeth have a competition to climb a hill. Elizabeth takes 30 minutes to climb the hill. Tom takes four times as long as Elizabeth does to climb the hill. How many hours does it take Tom to climb up the hill?
- 2 Answer: It takes Tom $30 \times 4 = \langle\langle 30 \times 4 = 120 \rangle\rangle 120$ minutes to climb the hill.
- 3 It takes Tom $120 / 60 = \langle\langle 120 / 60 = 2 \rangle\rangle 2$ hours to climb the hill.
- 4 So the answer is 2.
- 5 \equiv
- 6 Question: Jack is a soccer player. He needs to buy two pairs of socks and a pair of soccer shoes. Each pair of socks cost \$9.50, and the shoes cost \$92. Jack has \$40. How much more money does Jack need?
- 7 Answer: The total cost of two pairs of socks is $\$9.50 \times 2 = \$\langle\langle 9.5 \times 2 = 19 \rangle\rangle 19$.
- 8 The total cost of the socks and the shoes is $\$19 + \$92 = \$\langle\langle 19 + 92 = 111 \rangle\rangle 111$.
- 9 Jack need $\$111 - \$40 = \$\langle\langle 111 - 40 = 71 \rangle\rangle 71$ more.
- 10 So the answer is 71.
- 11 \equiv
- 12 **Question: There are three birds on the tree, shot one down, how many are left on the tree?**
- 13 **Answer:**

CHAIN-OF-THOUGHT (CoT)

Zero-shot CoT



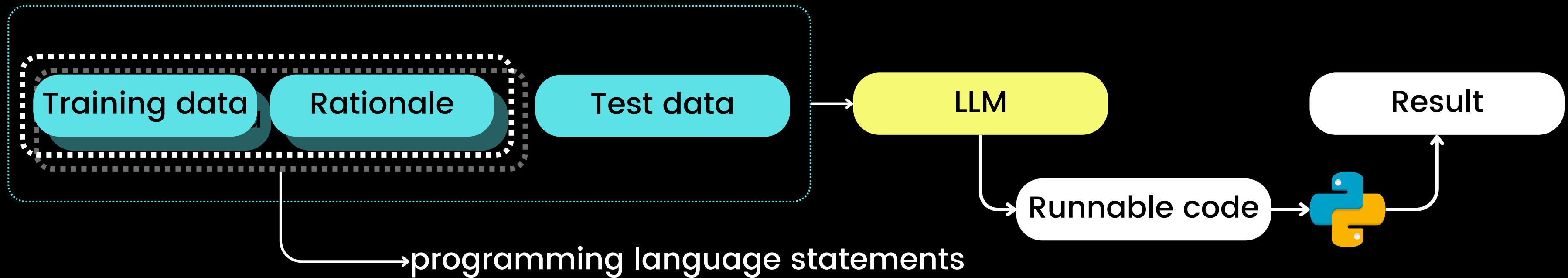
Few-shot CoT



→ manually written/model-generated high-quality reasoning chains.

PROGRAM-OF-THOUGHT (POT)

Few-shot PoT



Question: In Fibonacci sequence, it follows the rule that each number is equal to the sum of the preceding two numbers. Assuming the first two numbers are 0 and 1, what is the 50th number in Fibonacci sequence?

The first number is 0, the second number is 1, therefore, the third number is 0+1=1. The fourth number is 1+1=2. The fifth number is 1+2=3. The sixth number is 2+3=5. The seventh number is 3+5=8. The eighth number is 5+8=13.
..... (Skip 1000 tokens)
The 50th number is 32,432,268,459.

CoT

32,432,268,459



```
length_of_fibonacci_sequence = 50
fibonacci_sequence = np.zeros(length_of_)
fibonacci_sequence[0] = 0
fibonacci_sequence[1] = 1
For i in range(3, length_of_fibonacci_sequence):
    fibonacci_sequence[i] = fibonacci_sequence[i-1] +
    fibonacci_sequence[i-2]
ans = fibonacci_sequence[-1]
```

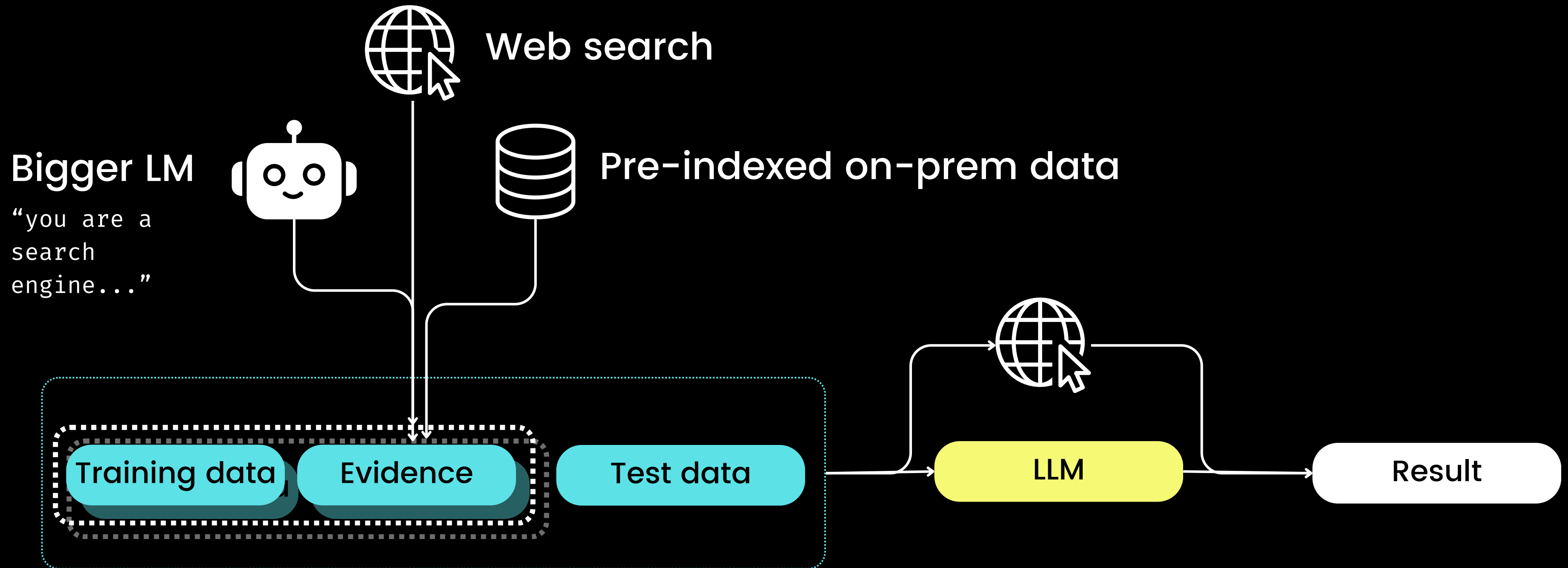
PoT



12,586,269,025



RETRIEVER



READER

Convert any URL to an LLM-friendly input with a simple prefix `https://r.jina.ai/`



READER



Enter your URL

https://arxiv.org/abs/2310.19923

Click below to fetch the source code of the page directly

Reader URL

https://r.jina.ai/https://arxiv.org/abs/2310.19923

Stream Mode

Click below to obtain the content through our Reader API

↓

FETCH CONTENT

<?xml version="1.0" encoding="UTF-8"?>

<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/DTD,

<html xmlns="http://www.w3.org/1999/xhtml" lang="en" xml:lang="en">

<head> <title>[2310.19923] Jina Embeddings 2: 8192-Token General-Purpose Text Embeddings for Lo

<meta name="viewport" content="width=device-width, initial-scale=1">

<link rel="apple-touch-icon" sizes="180x180" href="/static/browse/0.3.4/images/icons/apple-to

<link rel="icon" type="image/png" sizes="32x32" href="/static/browse/0.3.4/images/icons/favico

<link rel="icon" type="image/png" sizes="16x16" href="/static/browse/0.3.4/images/icons/favico

<link rel="manifest" href="/static/browse/0.3.4/images/icons/site.webmanifest">

<link rel="mask-icon" href="/static/browse/0.3.4/images/icons/safari-pinned-tab.svg" color="#!

<meta name="msapplication-TileColor" content="#da532c">

<meta name="theme-color" content="#ffffff">

<link rel="stylesheet" type="text/css" media="screen" href="/static/browse/0.3.4/css/arXiv.cs

<link rel="stylesheet" type="text/css" media="print" href="/static/browse/0.3.4/css/arXiv-pri

<link rel="stylesheet" type="text/css" media="screen" href="/static/browse/0.3.4/css/browse_sc

</head>

<body>

</body>

</html>

Title: Jina Embeddings 2: 8192-Token General-Purpose Text Embeddings for Long Documents

URL Source: https://arxiv.org/abs/2310.19923

Markdown Content:

Authors: [Michael Günther](https://arxiv.org/search/cs?searchtype=author&query=G%C3%BCnther,+M),

[View PDF](https://arxiv.org/pdf/2310.19923) [HTML (experimental)](https://arxiv.org/html/2310.1

> Abstract:Text embedding models have emerged as powerful tools for transforming sentences into

> To address these challenges, we introduce Jina Embeddings 2, an open-source text embedding mo

Submission history

From: Han Xiao \[[view email](https://arxiv.org/show-email/11a0a836/2310.19923)\]

**P.5.4.77(4)https://arxiv.org/abs/2310.19923v1[3] Fri, 20 Oct 2023 10:35:30 UTC (307 KB)

Pose a Question

Summarize the content in two sentences.

Input a question and combine it with the fetched content for LLM to generate an answer

The content discusses the challenges faced by existing open-source text embedding models in representing lengthy documents and introduces Jina Embeddings 2 as a solution to address these challenges.

The paper introduces Jina Embeddings 2, an open-source text embedding model that can handle up to 8192 tokens, addressing the limitation of existing models in representing long documents. The model achieves state-of-the-art performance on various embedding-related tasks and matches the performance of OpenAI's ada-002 model, with experiments showing that extended context improves performance in tasks like NarrativeQA.

ALL BASIC MODULES

Prompt

“Parameter” of the prompt

Zero-shot

Test data

LLM

Result

Few-shot

Training data

Test data

LLM

Result

Instruction
few-shot

Instruction

Training data

Test data

LLM

Result

Chain of
Thought

Training data

Rationale

Test data

LLM

Result

Retrieval

Training data

Evidence

Test data

LLM

Result



DSPy:
Not Your Average Prompt Engineering


What is DSPy


DSPY


- **D**eclarative **S**elf-improving Language **P**rograms, pythonically.
- DSPy is a framework for algorithmically optimizing prompts and LM weights, especially in a prompt pipeline.
- However, it is hard to learn.
 - *"Yeah man, I have been seeing DSPy everywhere but haven't found time to check it out yet"* - almost everyone I talk to about the project.


stanfordnlp/dspy


DSPy: The framework for programming—not prompting—foundation models





 124
Contributors

 105
Used by

 21
Discussions


 10k
Stars

 689
Forks



stanfordnlp/dspy: DSPy: The framework for programming—not prompting—foundation models

DSPy: The framework for programming—not prompting—foundation models - stanfordnlp/dspy

 GitHub

UNDERSTANDING DSPY

- DSPy closes the loop of prompt engineering;
- DSPy separates the logic (what) from textual representation (how).

UNDERSTANDING DSPY

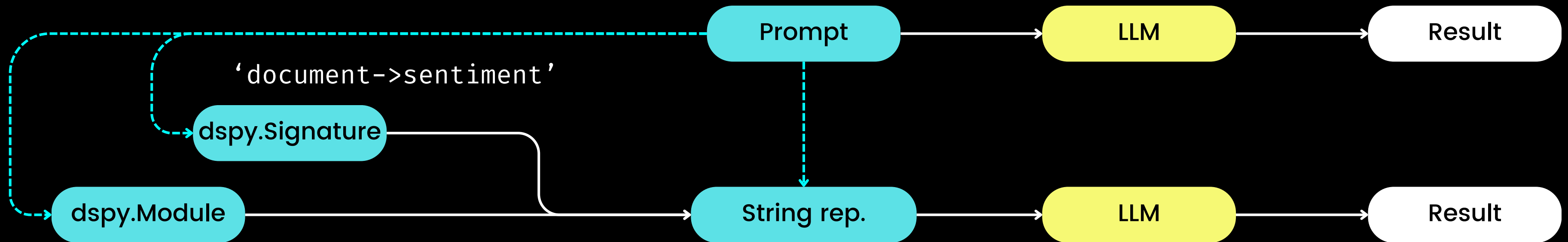
- DSPy closes the loop of prompt engineering;

Transforming prompt engineering from what is often a *manual, handcrafted process* into a *structured, well-defined machine learning workflow*: i.e. preparing datasets, defining the model, training, evaluating, and testing. In my opinion, this is the most revolutionary aspect of DSPy.

UNDERSTANDING DSPY

- DSPy separates the logic (what) from textual representation (how).

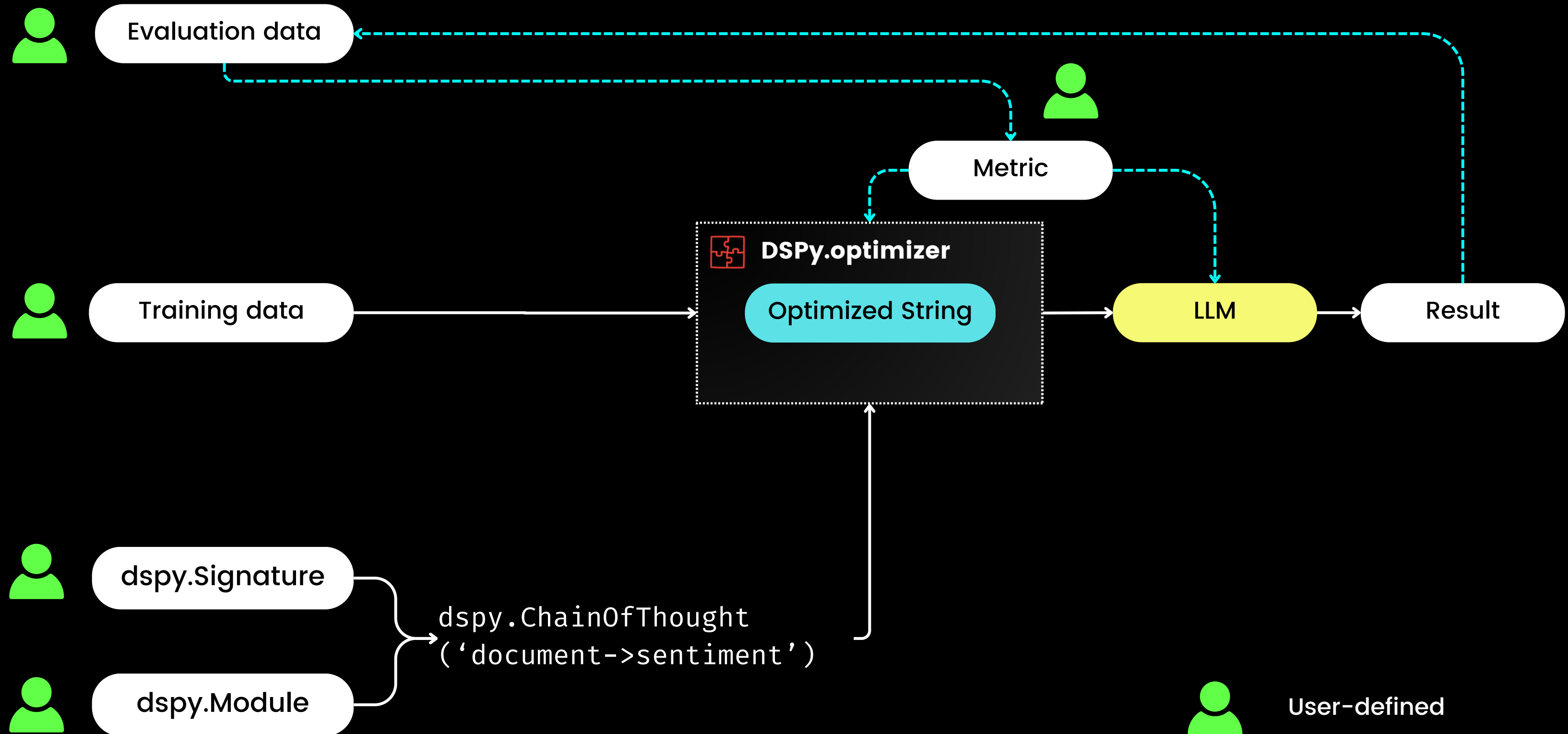
“This is important to me, I will lose my job if I can’t get the sentiment classification correct ...”



Predict
ChainOfThought
ReAct
ProgramOfThought

“... get the sentiment classification correct ... important ... lose my job ...”

DSPY.OPTIMIZER.COMPILE



WHAT EXACTLY DSPY.COMPILE OPTIMIZE

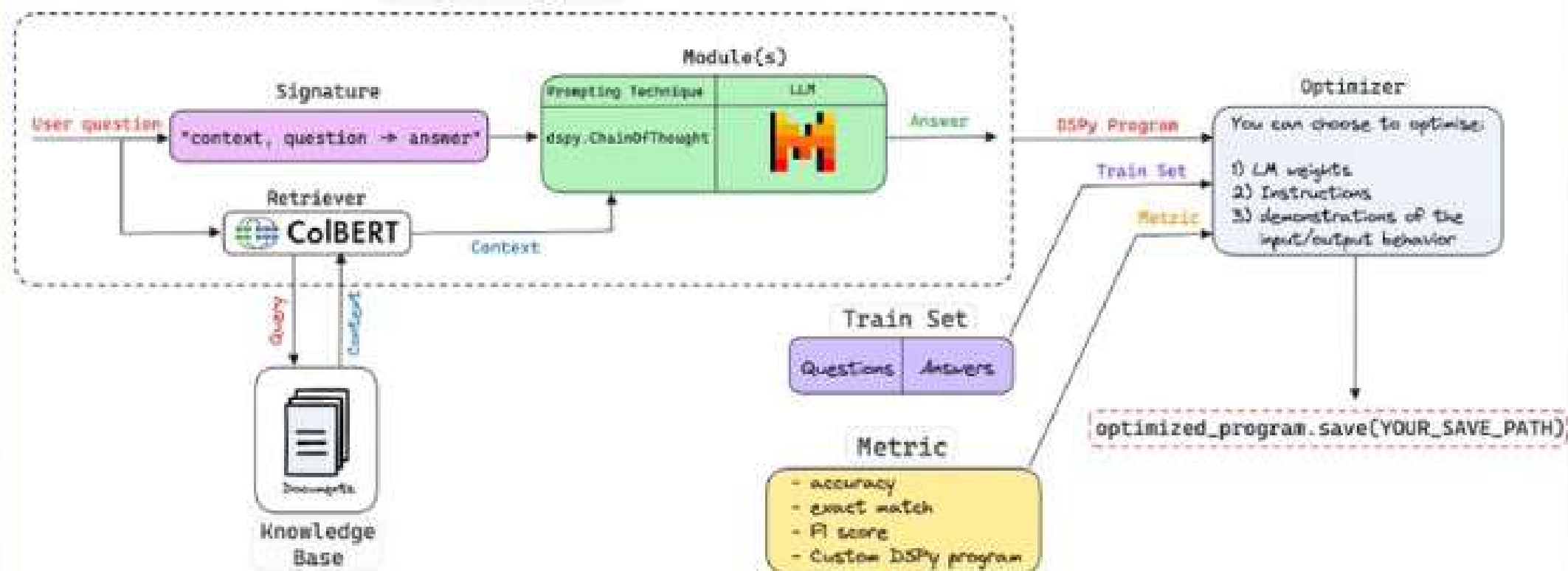
The compile function acts at the heart of this optimizer, akin to calling `optimizer.optimize()`. Think of it as the DSPy equivalent of training. This `compile()` process aims to tune:

- the few-shot demonstrations
- the instructions
- the LLM weights

You can imagine DSPy as a toolbox of **discrete optimization methods**.

DSPy: Programming not prompting LMs

DSPy Program



A blue wireframe sphere in the top right corner of the slide.

DSPy:
Not Your Average Prompt Engineering

Demo