

Spring-AI Intro



Kim Horn

Version 1.0

22 March 2025

Note: the slides when presented were animated, these animations are lost in the PDF

LLM = Large
Language
Model

*These can't
learn or
improve
themselves*

AI
(1956)

Symbolic AI

Machine Learning

Differentiators

Our
Knowledge

Knowledge
Graphs

Our Data

Our data, our
systems, and
other APIs.

Context

Context

Neural Nets
(1943)

Deep Learning
(2006)

Gen AI

CA SB 1047: Safe and Secure Innovation
for Frontier Artificial Intelligence
Models Act.

INTEGRATION

LLM

Proprietary / Open Source

Generative AI (GenAI) refers to Deep Learning models (LLMs) that can generate coherent human like high-quality text, images and other content based on the data they were trained on.

There are many issues:

- don't know about our business, our domain, where's the differentiation
- not trained on our data
- training data is incomplete
- trained at a past date
- input context size is limited
- not connected to the world, e.g. time
- stateless, no memory
- no structured output (e.g. Java Object)
- hallucinate, indecisive, can give biased, immoral, unethical and illegal advice
- hard to test and validate, catch 22
- security, privacy, and cost issues
- legal problems, different Laws, Standards for Safety First and Responsible AI, in every country.

Spring AI



Spring AI is an application framework for AI engineering. At its core, Spring AI addresses the fundamental challenge of AI integration: Connecting your enterprise Data and APIs with the AI Models.

Spring AI is a new framework. The project was started back in 2023. The first publicly available version, the 0.8, was released in February 2024. Now its at 1.0.0.M6. The 1st release is coming....

I could never find good working examples.

In Feb 2024 I started to build a set of 40+ Examples, off Snapshot build.....

Java is one of the top languages used with Azure OpenAI



Gen AI app

App framework



Firebase



spring

LLM



Gemini

Transactional database



AlloyDB

Scalable infrastructure



GKE

FULL SIZED ADULTS!



GOOGLE CLOUD NEXT, 2024

GOTO Conferences

POST <https://api.openai.com/v1/chat/completions>

```
{
  "model": "gpt-4o",
  "messages": [
    {
      "role": "system",
      "content": "You are a helpful polite assistant."
    },
    {
      "role": "user",
      "content": "What is a dog ?"
    }
  ]
}
```

Request

Try This:

```
"role": "system",
"content": "You are a bad pirate."
```

```
{
  "id": "chatcmpl-A8j01LYAMpbOdH0VL1scVjjbuJ78",
  "object": "chat.completion",
  "created": 1726642965,
  "model": "gpt-4o-2024-05-13",
  "choices": [
    {
      "index": 0,
      "message": {
        "role": "assistant",
        "content": "A domesticated mammal and a subspecies of grey wolf.  
Canis lupus familiaris....."
      },
      "logprobs": null,
      "finish_reason": "stop"
    }
  ],
  "usage": {
    "prompt_tokens": 23,
    "completion_tokens": 11,
    "total_tokens": 34,
    "completion_tokens_details": {
      "reasoning_tokens": 0
    }
  },
  "system_fingerprint": "fp_a5d11b2ef2"
}
```

Response

Every LLM vendor has a different API !!!

The prompt, our user's **message** = **"What is a dog ?"**

ChatClient (fluent)

```
ChatClient chatClient = chatClientBuilder.build();

String answer = chatClient
    .prompt()
    .user(message)
    .call()
    .content();
```

ChatModel (lower level)

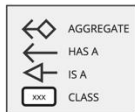
```
OpenAiChatModel chatModel;

ChatResponse response =
    chatModel.call(new Prompt(message));

String answer = response.getResult()
    .getOutput()
    .getText();
```

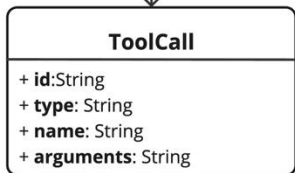
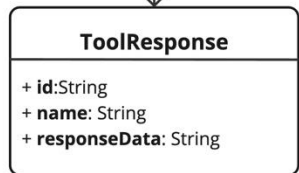
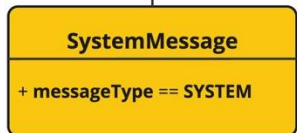
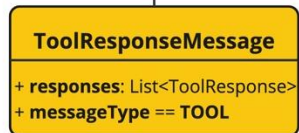
- Our Open AI key is automatically taken from our environment variables
- Pom specifies the specific vendor api: *spring-ai-openai-spring-boot-starter*

Spring AI Message API

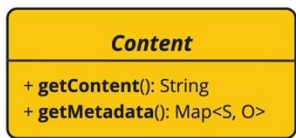
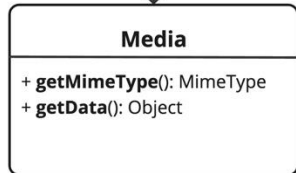


Role

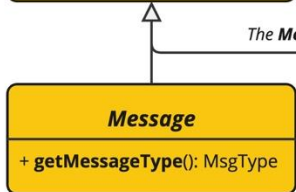
Prompt



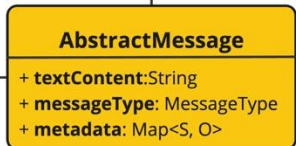
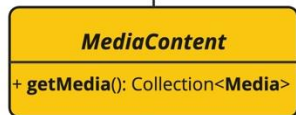
The data of a media attachment in a user message. It consists of a MIME type and either the raw media data or an URI to the data.



The **Content** field contains the primary, textual content along with optional **Metadata**.



The **Media** fields contains the multimodal inputs.



The Message interface encapsulates a **Prompt** textual content, a collection of metadata attributes, and a **MessageType**.

The collection specifies the **instructions** to the LLM.

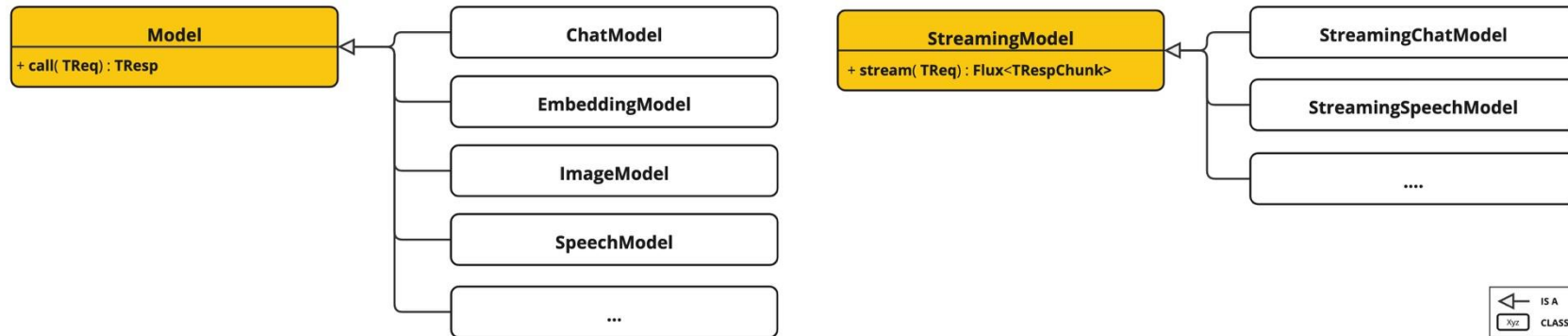
```
public class Prompt implements
ModelRequest<List<Message>>
{
    private final List<Message> messages;
    private ChatOptions chatOptions;
}
```

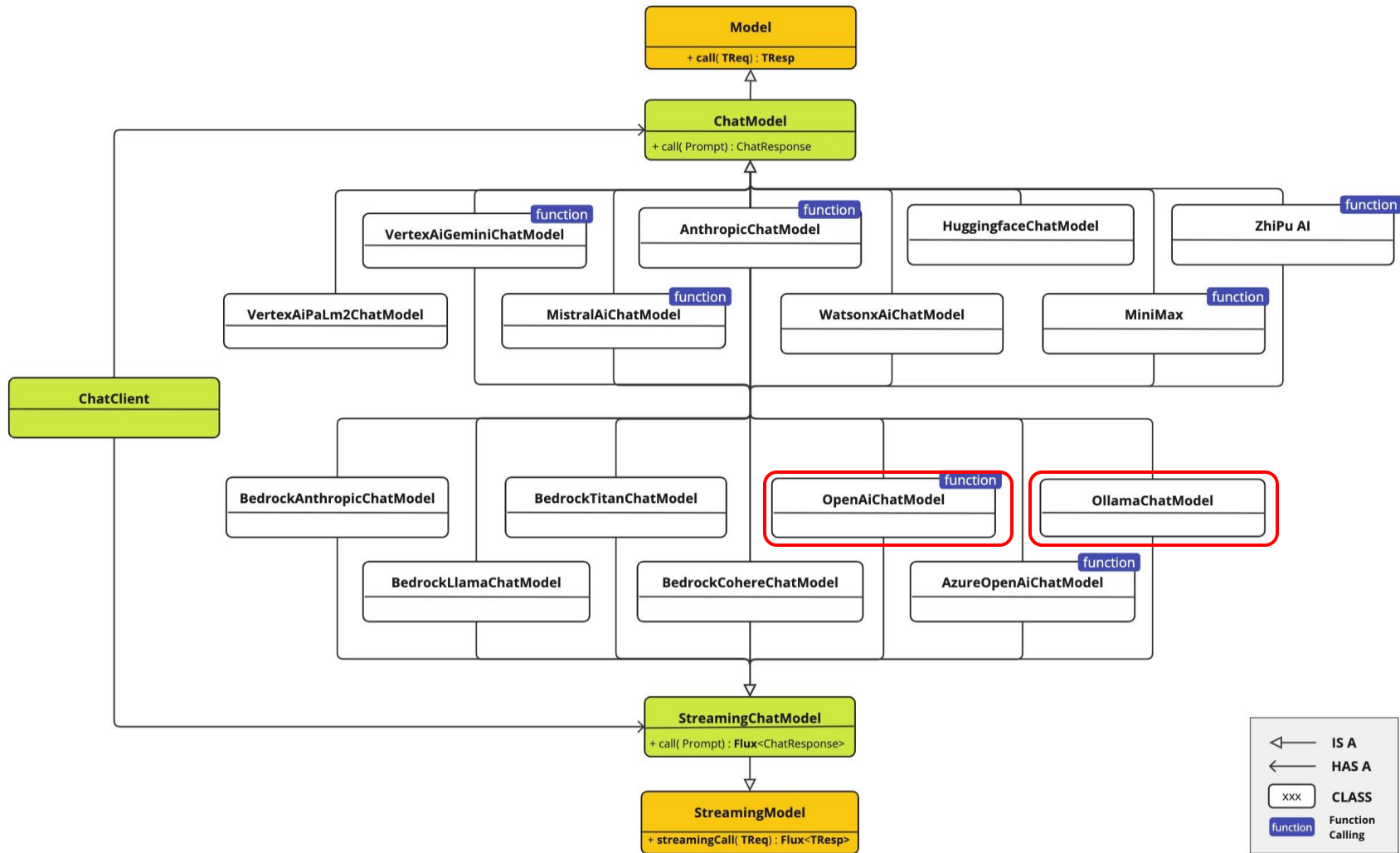
```
public interface Content {
    String getContent();
    Map<String, Object> getMetadata();
}
```

```
public interface Message extends Content {
    MessageType getMessageType();
}
```

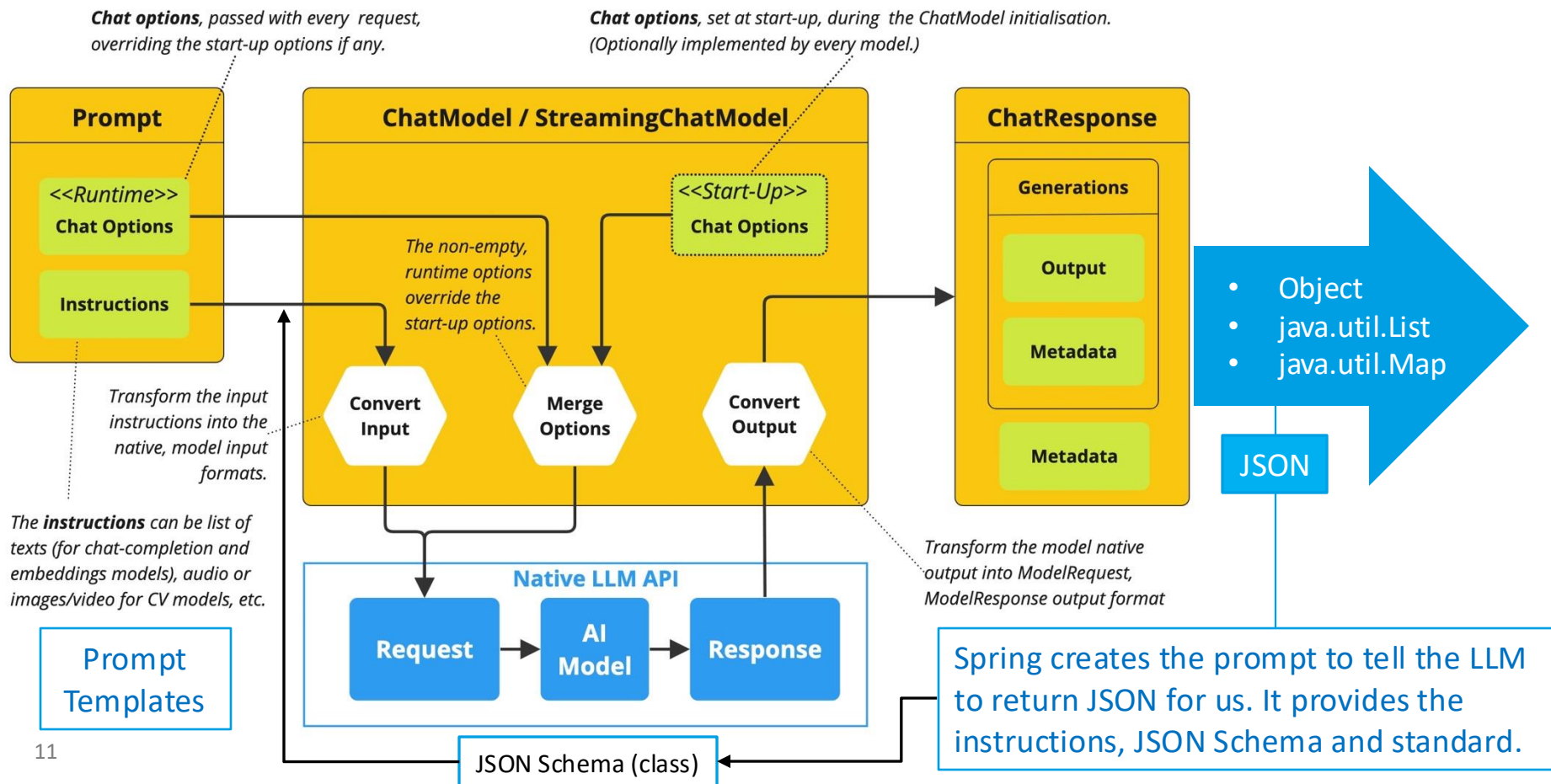


Client	Description
ChatClient / ChatModel	The core interface for <u>text-based interactions</u> . This is used for simple requests, prompt-based requests, and requests requiring a response in a specific format (such as JSON, XML etc).
ImageModel	A client for calling the vendor-specific <u>image generation</u> API. When we send some text as a request using this client, we get the URL of the generated image in the response.
SpeechModel / StreamingSpeechModel	A client for calling the vendor-specific <u>text-to-speech generation</u> API (TTS-1). When we send some text as a request using this client, we get the speech file path in the response.
AudioTranscriptionModel	A client for calling the vendor-specific audio transcription (<u>speech-to-text</u>) models. Currently, only OpenAI's 'whisper-1' is supported.
EmbeddingModel	A client for calling the vendor-specific <u>vector embedding</u> models, commonly used for making similarity searches in vector stores.

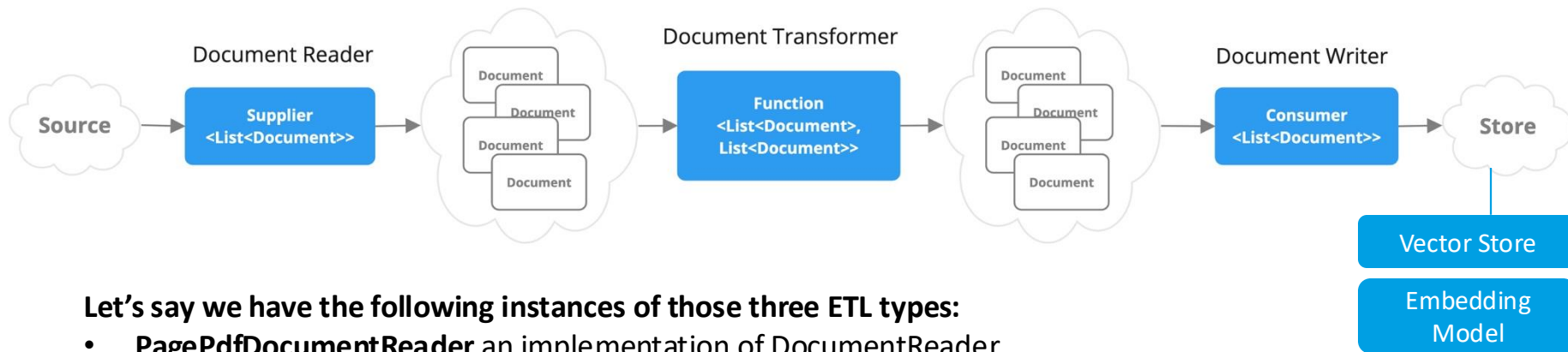




Components



How to we get our information (data, documents), into LLM ?



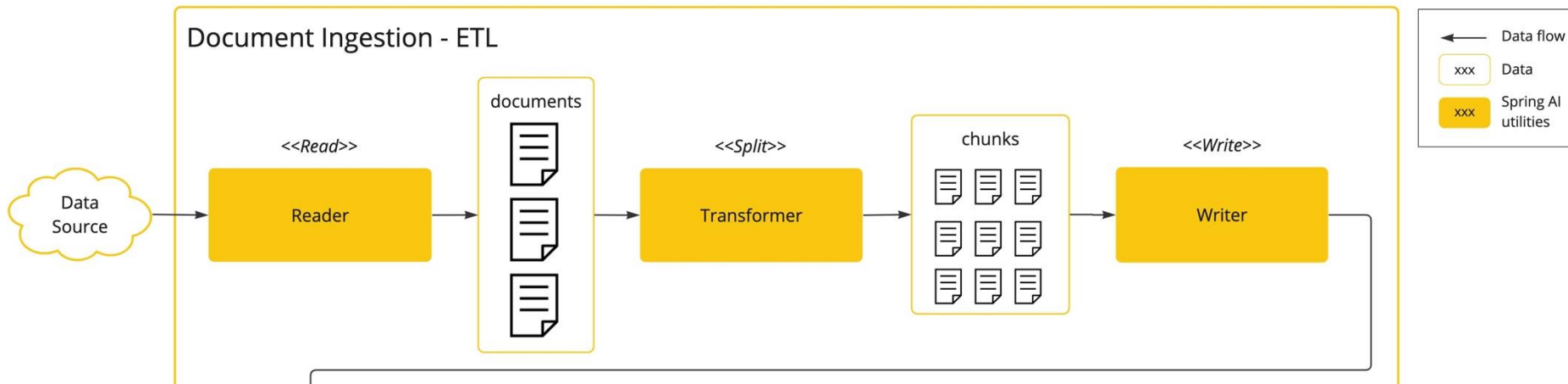
Let's say we have the following instances of those three ETL types:

- **PagePdfDocumentReader** an implementation of DocumentReader
- **TokenTextSplitter** an implementation of DocumentTransformer, to split into chunks
- **VectorStore** an implementation of DocumentWriter
- Many other built-in ETL types.

We can now read data (documents), split them into chunks, then load into a Vector Store, via an Embedding model, with the following code:

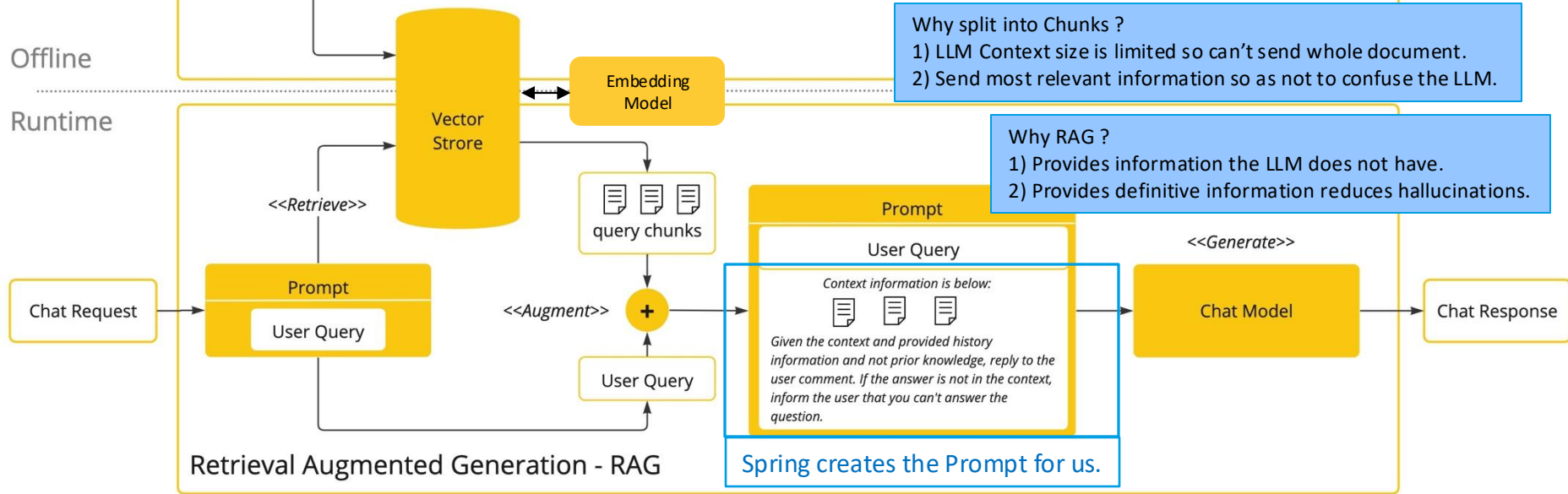
```
vectorStore.add(tokenTextSplitter.apply(pdfReader.get()));
```

Document Ingestion - ETL



Offline

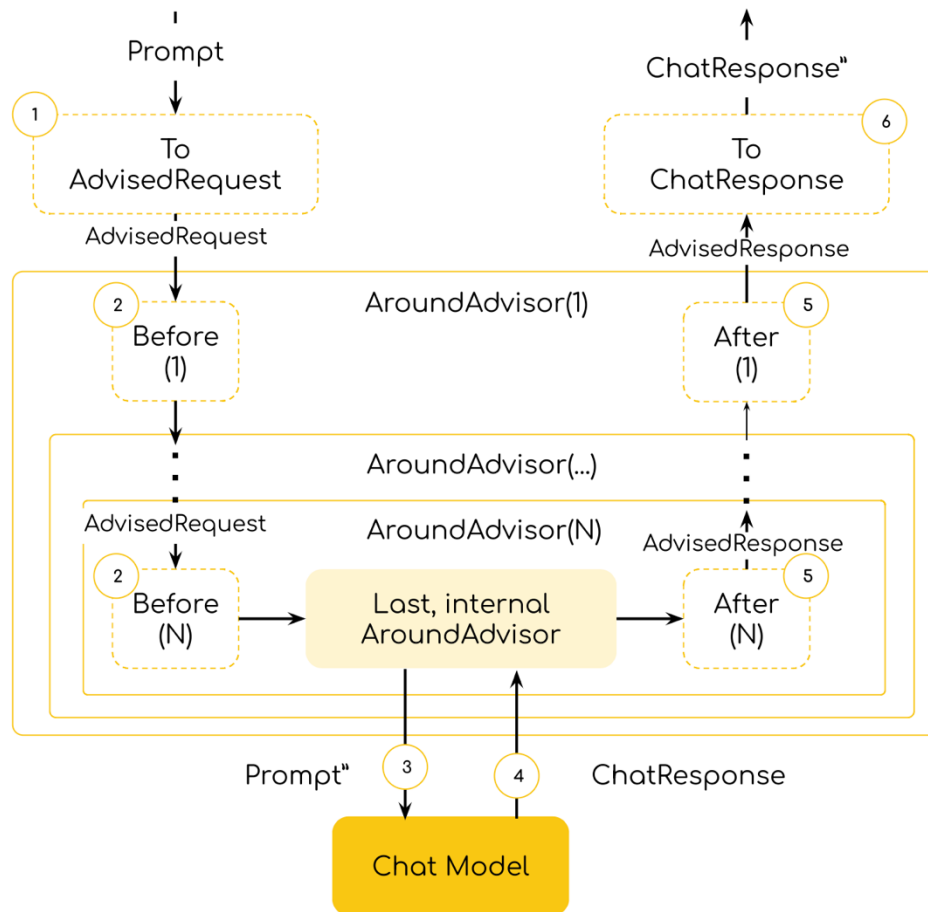
Runtime



The Advisors API provides a flexible and powerful way to intercept, modify, and enhance interactions with LLMs.

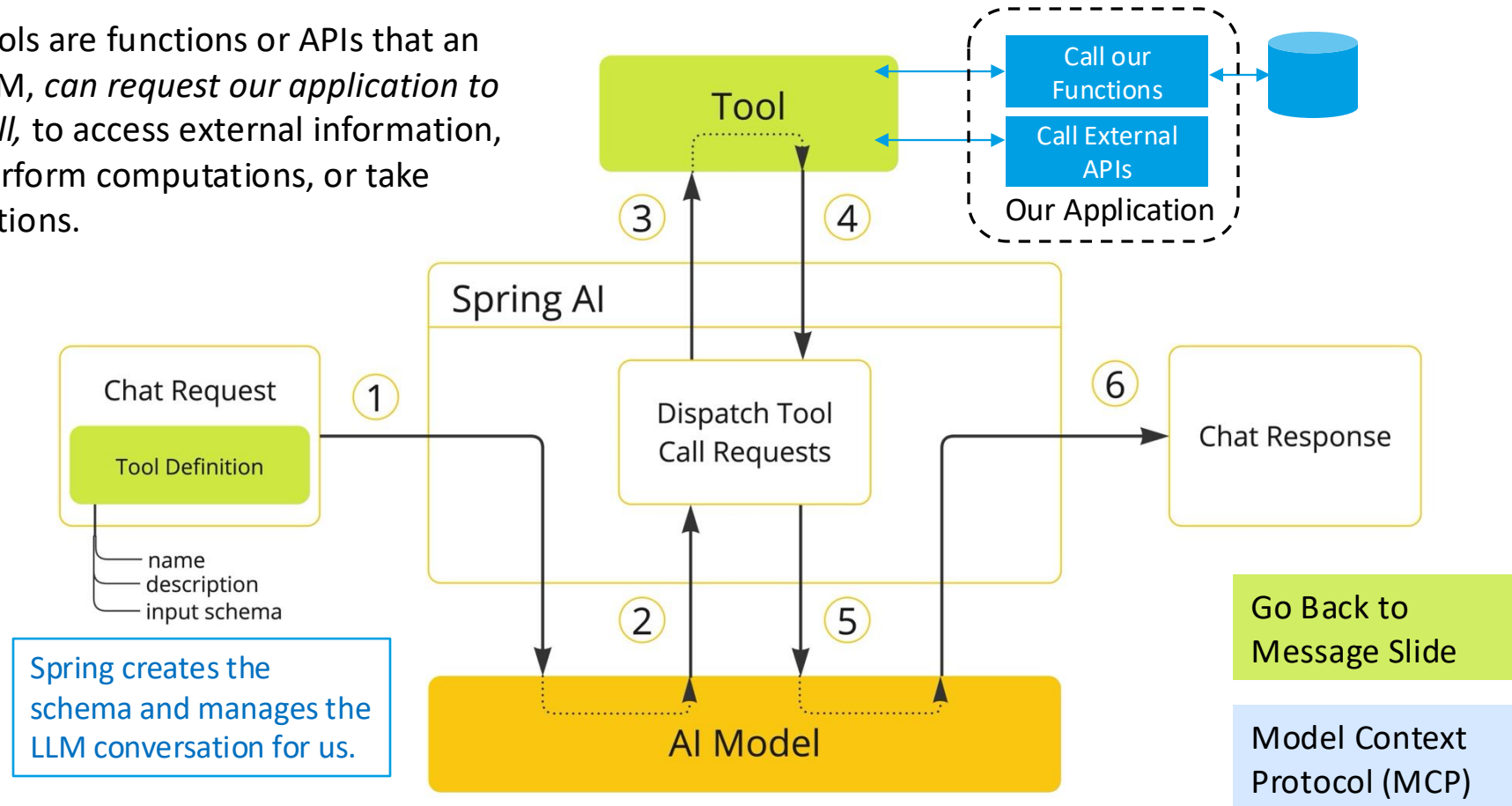
Types:

1. MemoryAdvisor
2. QuestionAnswerAdvisor – RAG
3. SafeGuardAdvisor
4. SimpleLoggingAdvisor
5.
6.
7. Roll your own....



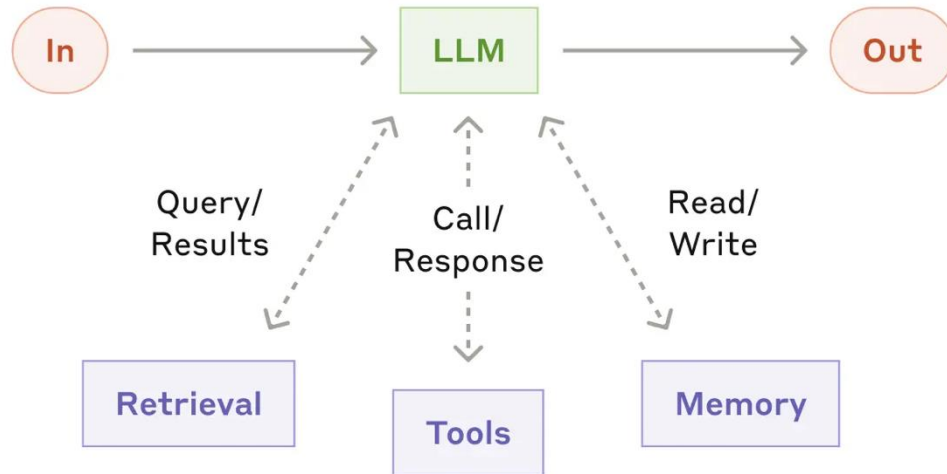
Tools (functions) – Perform actions, retrieve our information

Tools are functions or APIs that an LLM, *can request our application to call*, to access external information, perform computations, or take actions.



The basic feature of Spring-AI now provide all the building blocks for complex agentic systems, and its easy to implement their patterns.

<https://www.anthropic.com/research/building-effective-agents>



Anthropic Patterns :

- Chaining
- Routing
- Parallelization
- Orchestration
- Evaluator-Optimizer
- Agent*

*A definition is not really agreed yet

Examples

1) Chat Bot – RAG :

- talk to an Invoice PDF
- Run the PDF ETL, Embeds in Vector store, search, calls LLM with right Context
- Uses OpenAI – GPT-4o

2) Support Application/Agent:

- Routing Agent Pattern - to pass support tickets to the right expert team.
- Uses Local Ollama
 - Spring loads llama3.2 for us

About this Invoice ?

The Big Company

27 Little Street

Sydney, 2000

Phone: (061) 123-4567

INVOICE



INVOICE #

2034

DATE

21/2/2018

BILL TO

Fred Smith

Small Company

40A Big Street

Nice Ville, 3456

0437 856342

jimy@small.com.au

CUSTOMER ID

564

TERMS

Due Upon Receipt

DESCRIPTION	QTY	UNIT PRICE	AMOUNT
Service fee	1	200.00	200.00
Labor: 5 hours at \$75/hr	5	75.00	375.00
New client discount		(50.00)	(50.00)
Part Item 4523 - Big Nut	3	45.00	135.00
Part Item 3423 - Flexible cable	1	12.00	12.00
			-
			-
			-
			-
			-
			-
			-
			-
			-
			-
			-

Thank you for your business!

SUBTOTAL

672.00

TAX RATE

4.250%

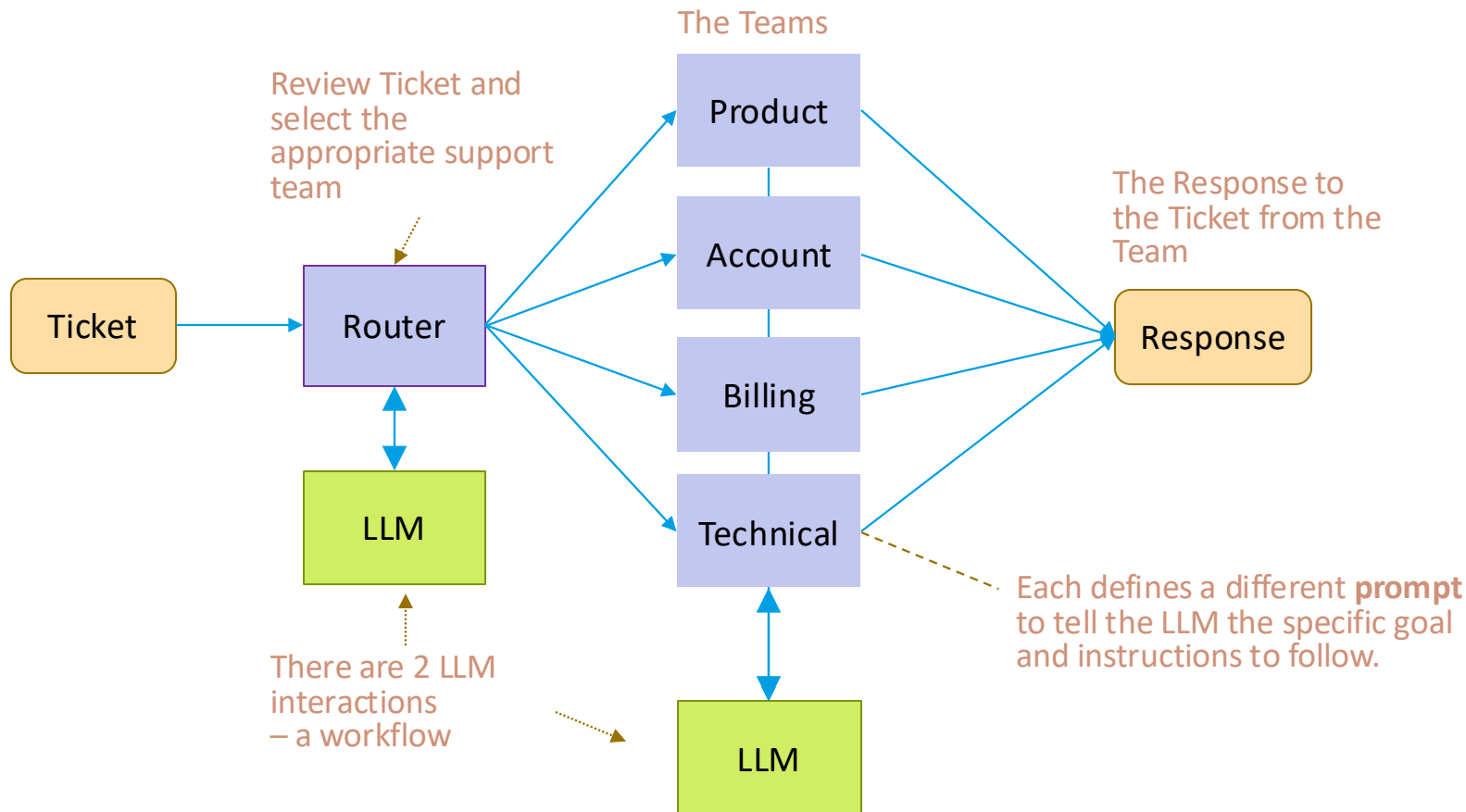
TAX

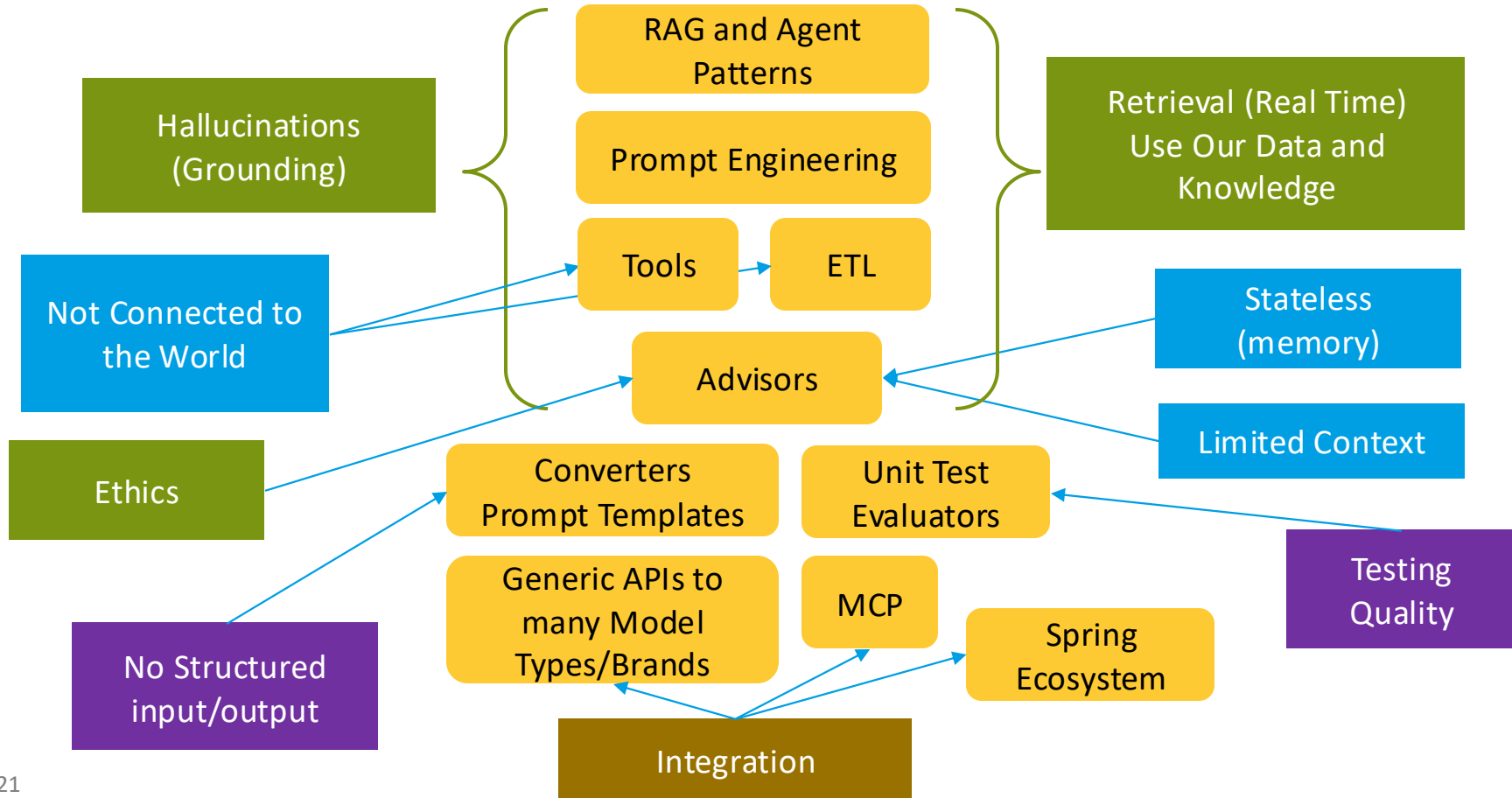
28.56

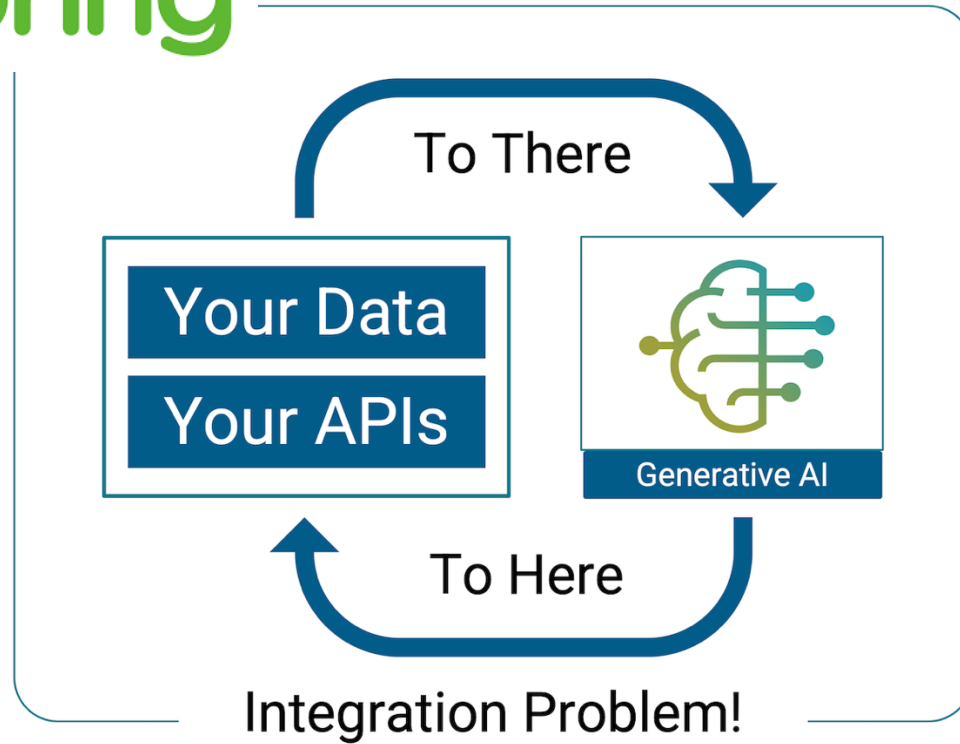
TOTAL

\$

700.56







The End