# Recap

## Objectives
- Improve customer/employee experience by delivering better search relevance
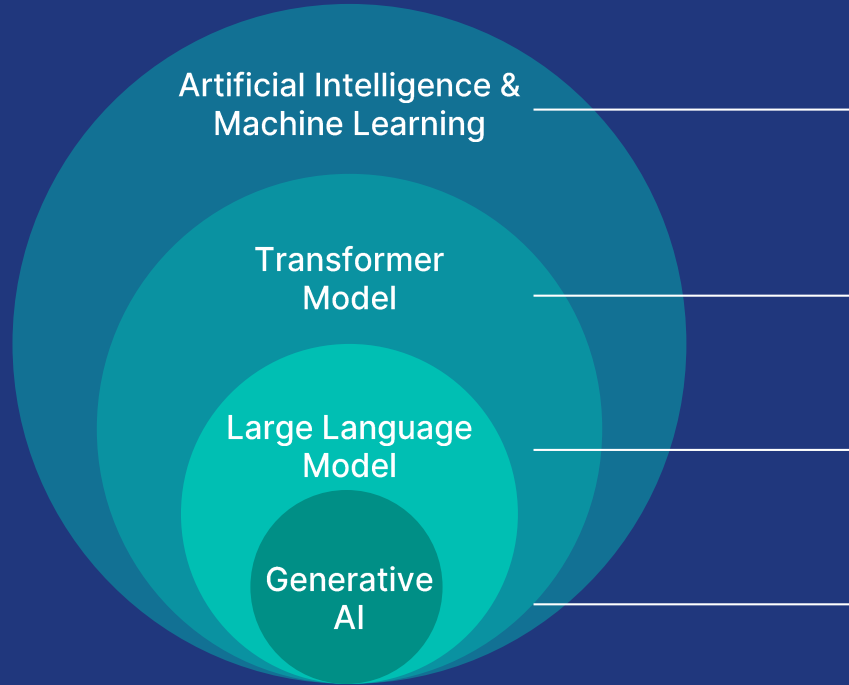- Minimize time to implement and fine tune relevance

## Challenges
- Getting the right information to the right audience
- Supporting legacy technology takes away from innovation and strategic impact
- Slow, unreliable search performance

# Introduction

Introduction to Elasticsearch and Generative AI

elastic

# ML, Transformer, and Large Language Model basics

**What is it?**

**Artificial Intelligence & Machine Learning**

The science of teaching computers to think, learn, and improve on their own.

**Transformer Model**

A Neural Network architecture that considers word relationships and context.

**Large Language Model**

An AI model that uses massive data to generate human-like text and perform language tasks exceptionally well

**Generative AI**

A large language model trained to compose content and responses to human prompts.

elastic

# Like any new **revolutionary** tech, Generative AI comes with a new set of strengths and limitations
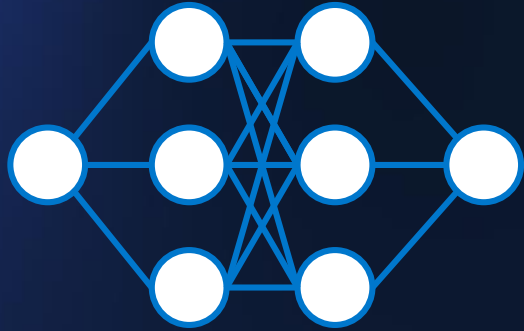
## Excels at human-like, iterative content creation

Natural language processing of large data sets used for creating human-like conversations, writing content, and providing code examples

## Limitations of Large Language Models

- Base **models** are trained on public data
- Data is frozen in time after training and fine-tuning
- Non-deterministic results and Hallucinations
- Cost & Privacy concerns for large scale use

**elastic**

A large language model is not a database

# Searching YOUR Data

How to search your data:
BM25, Embeddings, Vector Similarity, Retrieval Strategies

elastic

# Elasticsearch is a search engine - Inverted indices

## Indexing text fields in documents

✓ **Inverted indices:** tokens, count of frequency, and documents containing them

✓ **Freq and positions:** frequency and token offset for scoring, phrase queries and more

✓ **BM25:** A sparse, unsupervised model for lexical search, improvement over TF/IDF

✓ **Speed and Space Optimizations:** Block Max WAND, query short circuiting, match_only_text, and more!



| term | freq | documents |
|------|------|-----------|
| choice | 1 | 3 |
| coming | 1 | 1 |
| fury | 1 | 2 |
| is | 3 | 1, 2, 3 |
| ours | 1 | 2 |
| the | 2 | 2, 3 |
| winter | 1 | 1 |
| yours | 1 | 3 |

1: Winter is coming.
2: Ours is the fury.
3: The choice is yours.

Dictionary          Postings

elastic

How do we get the most relevant context
to answer the user's question,
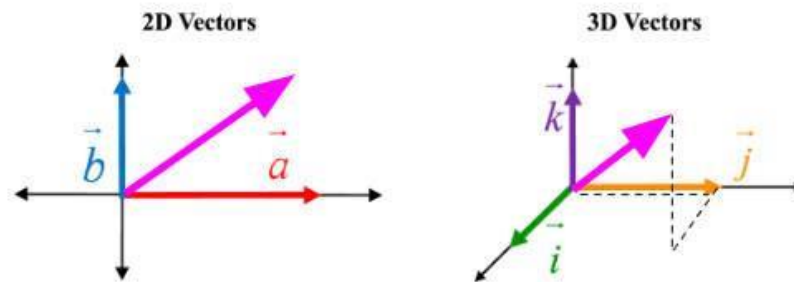in a natural way?

Beyond the BM25 "bag of words"

# Remember Vectors?



2D Vectors    3D Vectors

v = [1.0, 0.5, -2.1]

You likely learned the math for 3D vectors.

Some may have learned the math for n-dimensional vectors later in school.

NLP uses vectors with hundreds to thousands of dimensions (not pictured)

*What math textbooks looked like before common core (US centric joke, sorry)*

elastic

# Meaning can be encoded as a high dimensional Vector

**Text Embedding Model**

Chunking -> Transformer or LLM -> vector

**Image Classification (i.e. Google image search)**

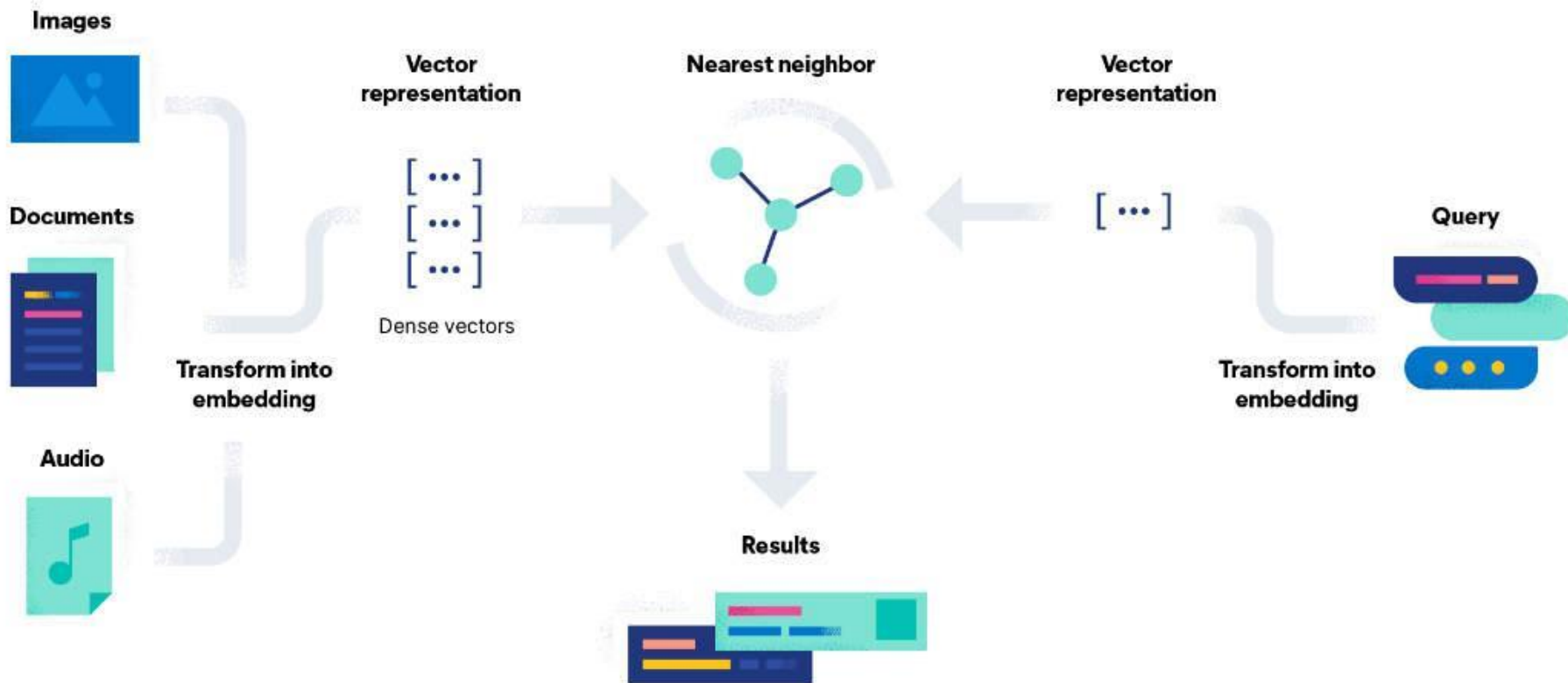Convolutional Neural Network -> Trained model -> vector

**Soundwave Classification**

FFT + time -> Spectrograph -> Image (see above)

# Vector search conceptual architecture
*Use vector nearest neighbor to generate a search ranking*

**Images**

**Documents**

**Transform into embedding**

**Audio**

**Vector representation**

$$[ \ ... \ ]$$
$$[ \ ... \ ]$$
$$[ \ ... \ ]$$

Dense vectors

**Nearest neighbor**

**Vector representation**

$$[ \ ... \ ]$$

**Query**

**Transform into embedding**

**Results**

elastic

# Data Ingestion and Embedding Generation

POST /_doc

```
{
    "_id":"product-1234",
    "product_name":"Summer Dress",
    "description":"Our best-selling…",
    "Price": 118,
    "color":"blue",
    "fabric":"cotton",
} "desc_embedding":[0.452,0.3242,…]
}
```

You asked, we answered: Our best-selling classic wrap dress now comes in a cotton poplin that's wear-all-day perfect. Bonus: stripes (our favorite).

FIT
• 39" from high point of shoulder

DETAILS
• Cotton.
• Lined.
• Machine wash.
• Import.

**Source data**

POST /_doc

## 🔘 ML Inference pipelines                    ⊕ Add inference pipeline

Inference pipelines will be run as processors from the Enterprise Search Ingest Pipeline

**ml-inference-embedding-generation**                    Actions ⠿

● Deployed    pytorch    text_embedding

**ml-inference-emotional-analysis**                    Actions ⠿

● Deployed    pytorch    text_classification

Learn more about deploying ML models in Elastic ↗

elastic

# Vector Query

**summer clothes**

**Transformer model**

```
GET product-catalog/_search

{
    "knn": {
        "field": "desc_embbeding",
        "k": 5,
        "num_candidates": 50,
        "query_vector_builder": {
                "text_embedding": {
                    "model_text": "summer clothes",
                    "model_id": <text-embedding-model>
        },

        "filter": {
            "term": {
                "department": "women"
            }
        }
    },
    "size": 10
}
```
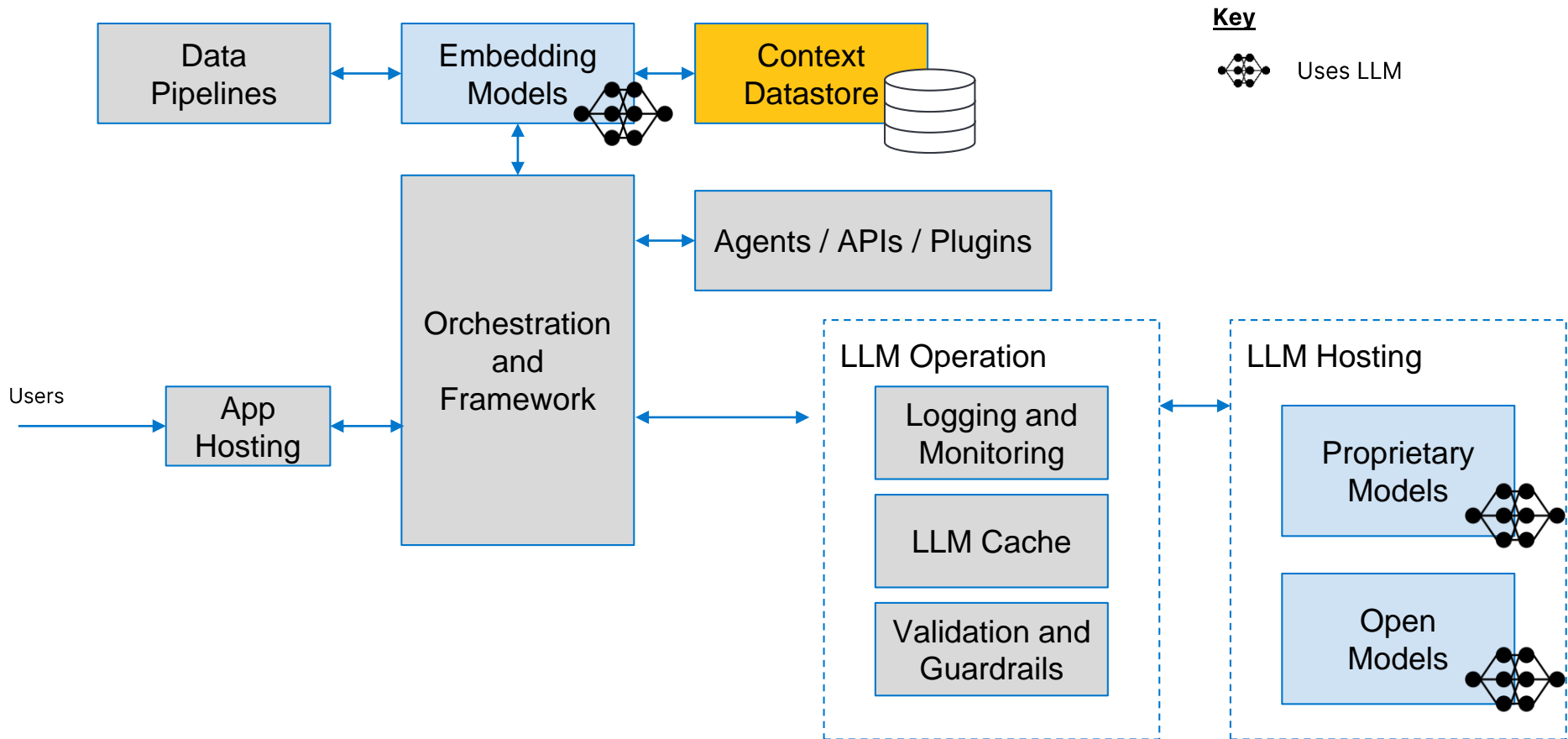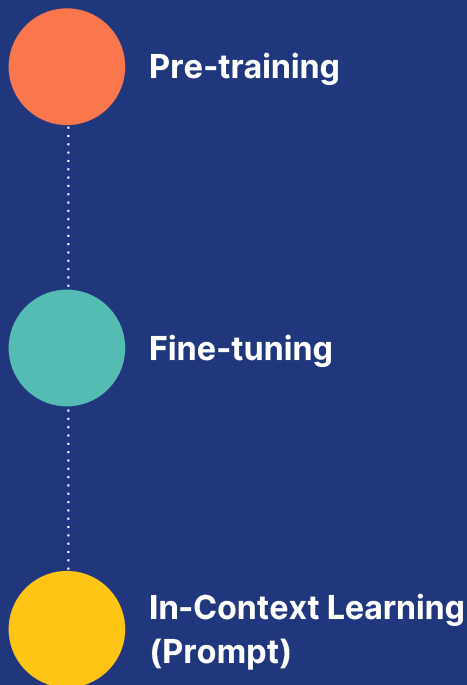
elastic

# RAG Applications

RAG = Retrieval Augmented Generation

Developing AI Applications with Context

elastic

# Emerging Generative AI App Architecture

# The 3 ways LLMs get 'smart'

**Pre-training**

*Foundational* or *Base* model training costs tens to hundreds of millions of $USD. LLMs learn **language** and **knowledge** from massive public data sets.

**Fine-tuning**

- **Task specific training** (classification, etc)
- Improve **quality of responses** in a domain
- Add **knowledge** from a specific data source
- **Alignment** with safeguards and ethical limits

**In-Context Learning (Prompt)**

- **Prompt engineering** techniques
  - **In-context learning and instruction**
- **Retrieval Augmented Generation**
  - Include **knowledge** in prompt

elastic

# 'Prompt Engineering'

Prompt Engineering is the art and science of designing effective prompts to guide the responses of Large Language Models.

a.k.a. Coding in 'Natural Language"

Prompts can ask the LLM to:

- Complete a task like summarization
- Follow provided context
- Make step by step plans or instructions
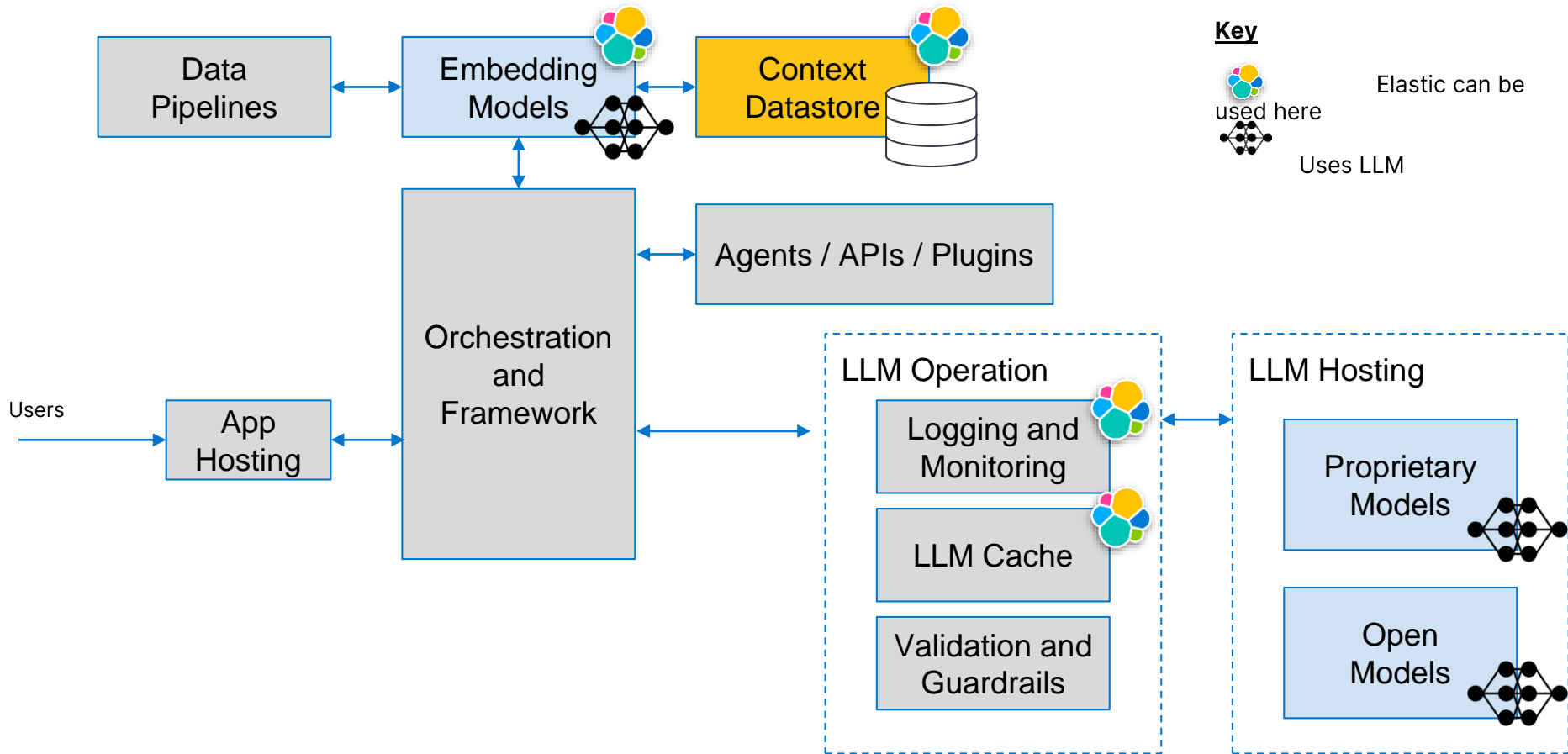- Format output in specific useful ways



elastic

# Anatomy of a Prompt

You are a helpful AI assistant who answers questions.

System Prompt

<<Additional Context>>

Opportunity for additional instructions / Context

User: <<User Supplied Input>>
AI:

User Input

Context Window

Completion

*Not technically part of the prompt*

elastic

# Retrieval Augmented Generation (RAG)

You are a helpful AI assistant who answers questions using the following supplied context. If you can't answer the question using this context say "I don't know"

System Prompt

Context Window

Context: The color of the sky is purple today

Supplied Context

User: What does the sky look like today?
AI:

User Input

Completion

elastic

RAG uses semantic search techniques like those in Elasticsearch to act as the bridge between private data and Generative AI

PII

Proprietary Information

Customer Case History

Private Knowledge Base

elastic

# Emerging Generative AI App Architecture

# Question Answering + Context Retrieval Workflow

**What is the policy of the company about X?**

original question

Retrieved context

question as search query

**Elastic Relevance Engine (ESRE)**

search results

Contextual data store 1

Contextual data store 2

Contextual data store 3

elastic

# Use Cases

# Generative AI is evolving within enterprises: Legal

**TODAY**

🔍 Work contract for California

**TOMORROW**

🔍 What are the main labor and employment law requirements at our California office?

elastic

# Generative AI is evolving within enterprises: Customer Success



**TODAY**

🔍 **Customer shopping locations**

**TOMORROW**

🔍 **Are my customers in Dallas buying products at locations most convenient for them, and with the deepest discounts?**

elastic

# Thank You

Elastic is a SEARCH Company

www.elastic.co