

# Generative AI Innovative Applications



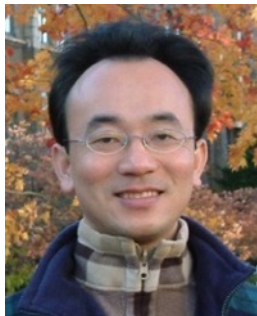
## Transformers for Natural Language Processing and Computer Vision; Large Language Models (LLMs)

1142GAIIA02

MBA, IM, NTPU (M6031) (Spring 2026)  
Tue 2, 3, 4 (9:10-12:00) (B3F17)



<https://meet.google.com/paj-zhhj-mya>



Min-Yuh Day, Ph.D.  
Professor and Director



Institute of Information Management, National Taipei University

<https://web.ntpu.edu.tw/~myday>



# Syllabus

Week	Date	Subject/Topics
1	2026/02/24	Introduction to Generative AI Innovative Applications
2	2026/03/03	Transformers for Natural Language Processing and Computer Vision; Large Language Models (LLMs)
3	2026/03/10	NVIDIA Building RAG Agents with LLMs Part I
4	2026/03/17	Case Study on Generative AI Innovative Applications I
5	2026/03/24	NVIDIA Building RAG Agents with LLMs Part II
6	2026/03/31	NVIDIA Building RAG Agents with LLMs Part III
7	2026/04/07	Make-up holiday for NTPU Sports Day (No Classes)
8	2026/04/14	Midterm Project Report

# Syllabus

**Week Date Subject/Topics**

**9 2026/04/21 Generative AI for Multimodal Information Generation**

**10 2026/04/28 NVIDIA Generative AI with Diffusion Models Part I**

**11 2026/05/05 NVIDIA Generative AI with Diffusion Models Part II**

**12 2026/05/12 Case Study on Generative AI Innovative Applications II**

**13 2026/05/19 NVIDIA Generative AI with Diffusion Models Part III**

**14 2026/05/26 Agentic AI and Large Multimodal Agents (LMAs)**

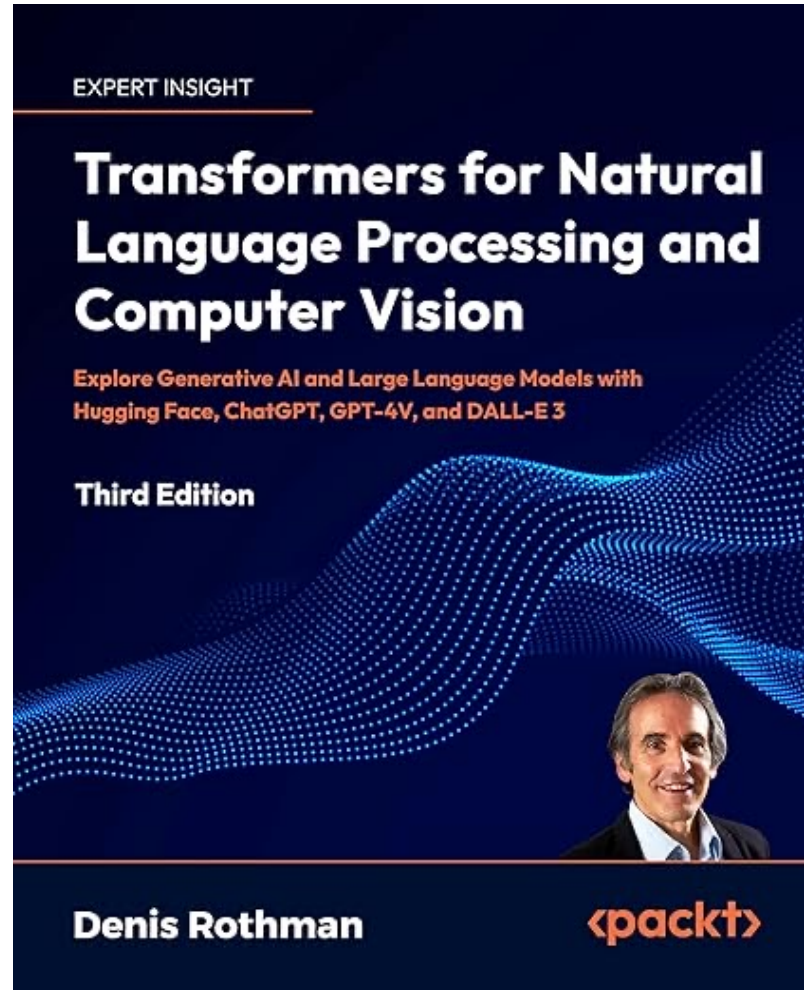
**15 2026/06/02 Final Project Report I**

**16 2026/06/09 Final Project Report II**

Denis Rothman (2024),

# Transformers for Natural Language Processing and Computer Vision:

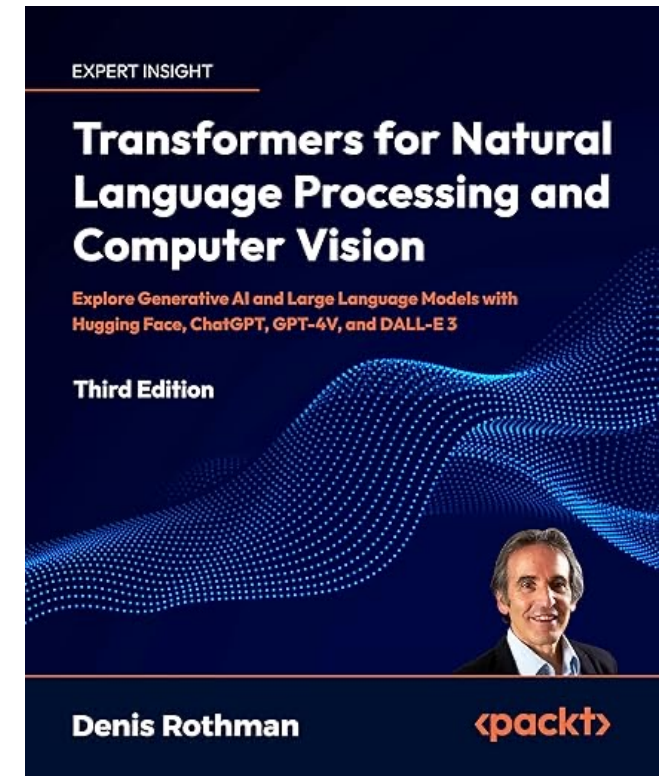
Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3,  
3rd Edition, Packt Publishing



Denis Rothman (2024),

# Transformers for Natural Language Processing and Computer Vision: Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3, 3rd Edition, Packt Publishing

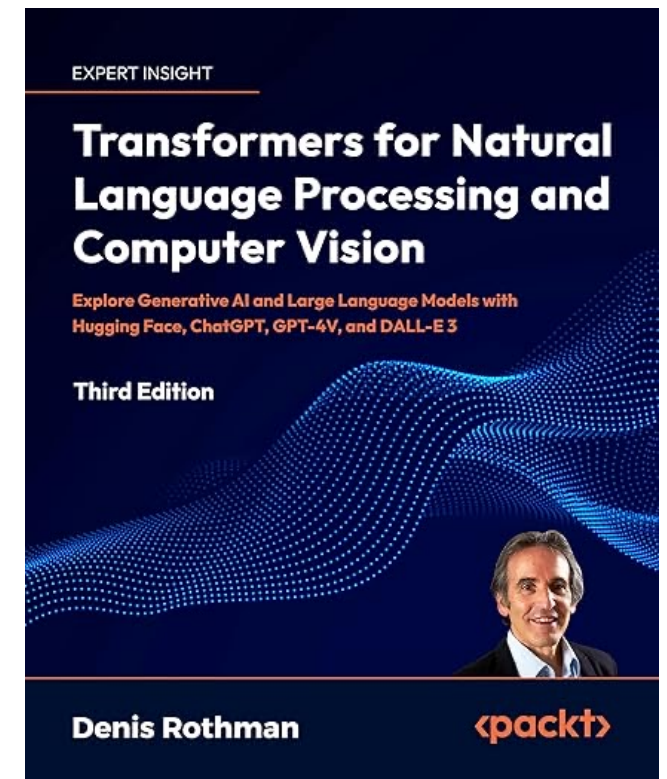
- 1.What Are Transformers?
- 2.Getting Started with the Architecture of the Transformer Model
- 3.Emergent vs Downstream Tasks: The Unseen Depths of Transformers
- 4.Advancements in Translations with Google Trax, Google Translate, and Gemini
- 5.Diving into Fine-Tuning through BERT
- 6.Pretraining a Transformer from Scratch through RoBERTa
- 7.The Generative AI Revolution with ChatGPT
- 8.Fine-Tuning OpenAI GPT Models
- 9.Shattering the Black Box with Interpretable Tools
- 10.Investigating the Role of Tokenizers in Shaping Transformer Models
- 11.Leveraging LLM Embeddings as an Alternative to Fine-Tuning
- 12.Toward Syntax-Free Semantic Role Labeling with ChatGPT and GPT-4
- 13.Summarization with T5 and ChatGPT
- 14.Exploring Cutting-Edge LLMs with Vertex AI and PaLM 2
- 15.Guarding the Giants: Mitigating Risks in Large Language Models
- 16.Beyond Text: Vision Transformers in the Dawn of Revolutionary AI
- 17.Transcending the Image-Text Boundary with Stable Diffusion
- 18.Hugging Face AutoTrain: Training Vision Models without Coding
- 19.On the Road to Functional AGI with HuggingGPT and its Peers
- 20.Beyond Human-Designed Prompts with Generative Ideation



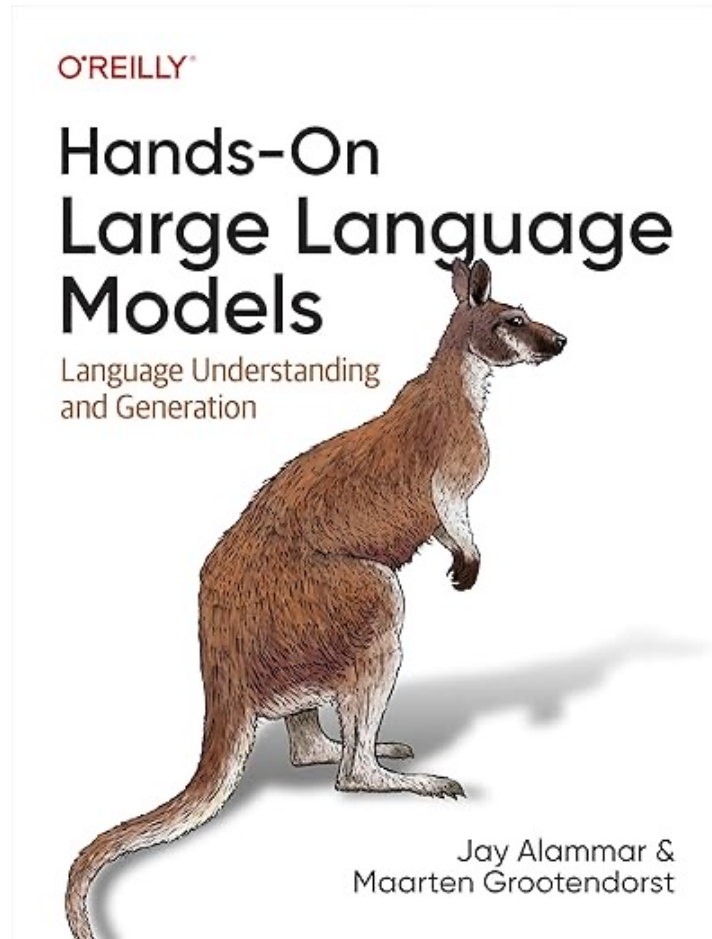
Denis Rothman (2024),

# Transformers for Natural Language Processing and Computer Vision: Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3, 3rd Edition, Packt Publishing

Chapter	Colab	Kaggle	Gradient	StudioLab
<b>Part I The Foundations of Transformer Models</b>				
<b>Chapter 1: What are Transformers?</b>				
<ul style="list-style-type: none"><li>• <code>XO_1_and_Accelerators.ipynb</code></li><li>• <code>ChatGPT_Plus_writes_and_explains_AI.ipynb</code></li></ul>	<a href="#">Open in Colab</a> <a href="#">Open in Colab</a>	<a href="#">Open in Kaggle</a> <a href="#">Open in Kaggle</a>	<a href="#">Run on Gradient</a> <a href="#">Run on Gradient</a>	<a href="#">Open Studio Lab</a> <a href="#">Open Studio Lab</a>
<b>Getting started with DeepSeek-R1 Reasoning models. Integrated into HuggingFace Hub and Together.</b>				
<ul style="list-style-type: none"><li>• <code>DeepSeek_Hugging_Face.ipynb</code></li></ul>	<a href="#">Open in Colab</a>	<a href="#">Open in Kaggle</a>	<a href="#">Run on Gradient</a>	<a href="#">Open Studio Lab</a>
<b>Chapter 2: Getting Started with the Architecture of the Transformer Model</b>				
<ul style="list-style-type: none"><li>• <code>Multi_Head_Attention_Sub_Layer.ipynb</code></li><li>• <code>positional_encoding.ipynb</code></li></ul>	<a href="#">Open in Colab</a> <a href="#">Open in Colab</a>	<a href="#">Open in Kaggle</a> <a href="#">Open in Kaggle</a>	<a href="#">Run on Gradient</a> <a href="#">Run on Gradient</a>	<a href="#">Open Studio Lab</a> <a href="#">Open Studio Lab</a>
<b>Chapter 3: Emergent vs Downstream Tasks: the Unseen Depths of Transformers</b>				
<ul style="list-style-type: none"><li>• <code>From_training_to_emergence.ipynb</code></li><li>• <code>Transformer_tasks_with_Hugging_Face.ipynb</code></li></ul>	<a href="#">Open in Colab</a> <a href="#">Open in Colab</a>	<a href="#">Open in Kaggle</a> <a href="#">Open in Kaggle</a>	<a href="#">Run on Gradient</a> <a href="#">Run on Gradient</a>	<a href="#">Open Studio Lab</a> <a href="#">Open Studio Lab</a>



Jay Alammar and Maarten Grootendorst (2024),  
**Hands-On Large Language Models:  
Language Understanding and Generation,**  
O'Reilly Media



Jay Alammar and Maarten Grootendorst (2024),  
**Hands-On Large Language Models:  
Language Understanding and Generation,**  
O'Reilly Media

Chapter 1: Introduction to Language Models

Chapter 2: Tokens and Embeddings

Chapter 3: Looking Inside Transformer LLMs

Chapter 4: Text Classification

Chapter 5: Text Clustering and Topic Modeling

Chapter 6: Prompt Engineering

Chapter 7: Advanced Text Generation Techniques and Tools

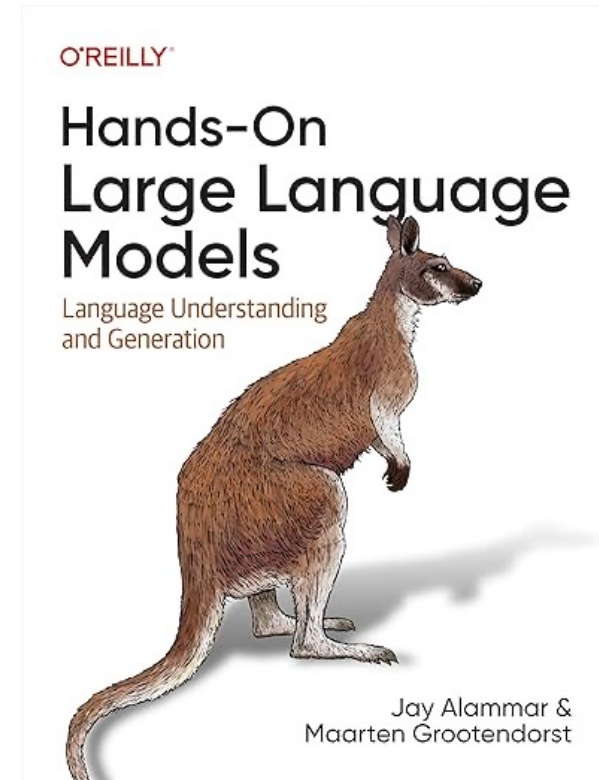
Chapter 8: Semantic Search and Retrieval-Augmented Generation

Chapter 9: Multimodal Large Language Models

Chapter 10: Creating Text Embedding Models

Chapter 11: Fine-tuning Representation Models for Classification

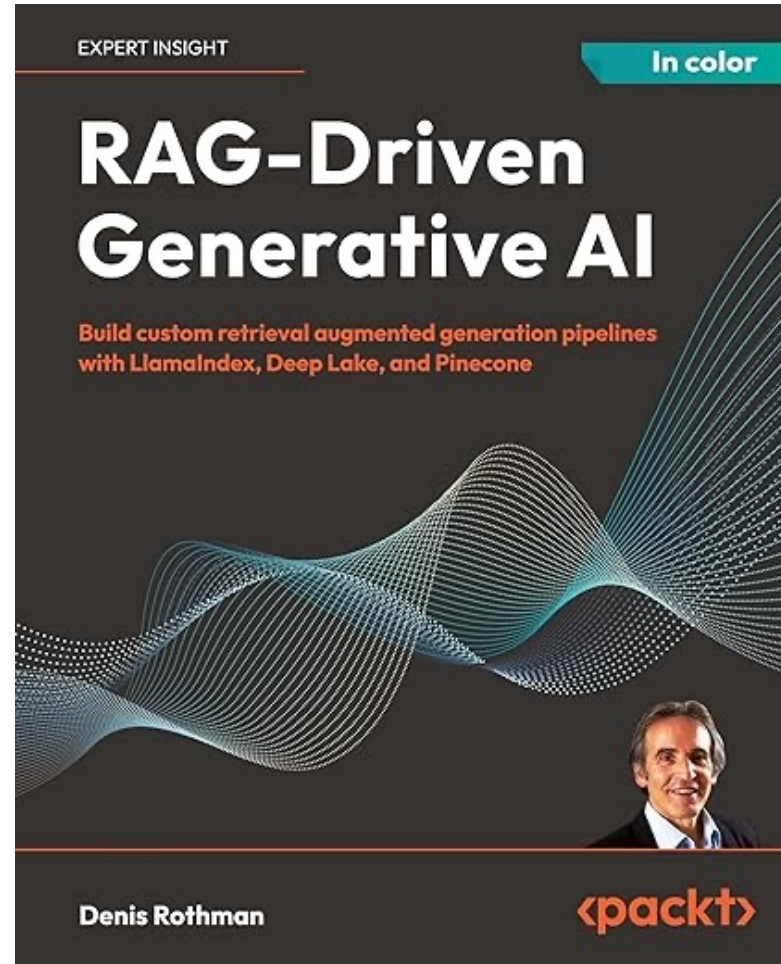
Chapter 12: Fine-tuning Generation Models



Denis Rothman (2024),

# RAG-Driven Generative AI:

Build custom retrieval augmented generation pipelines with LlamaIndex, Deep Lake, and Pinecone,  
Packt Publishing

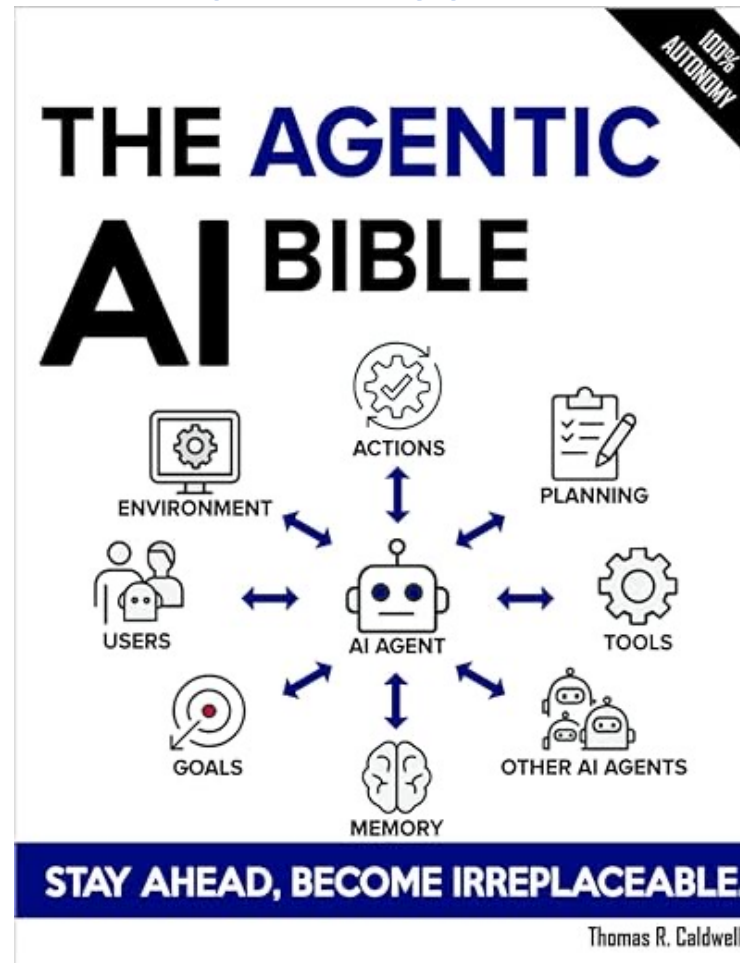


Thomas R. Caldwell (2025),

# The Agentic AI Bible:

The Complete and Up-to-Date Guide to Design, Build, and Scale Goal-Driven,  
LLM-Powered Agents that Think, Execute and Evolve,

Independently published



# **Generative AI**

# **Large Language Models**

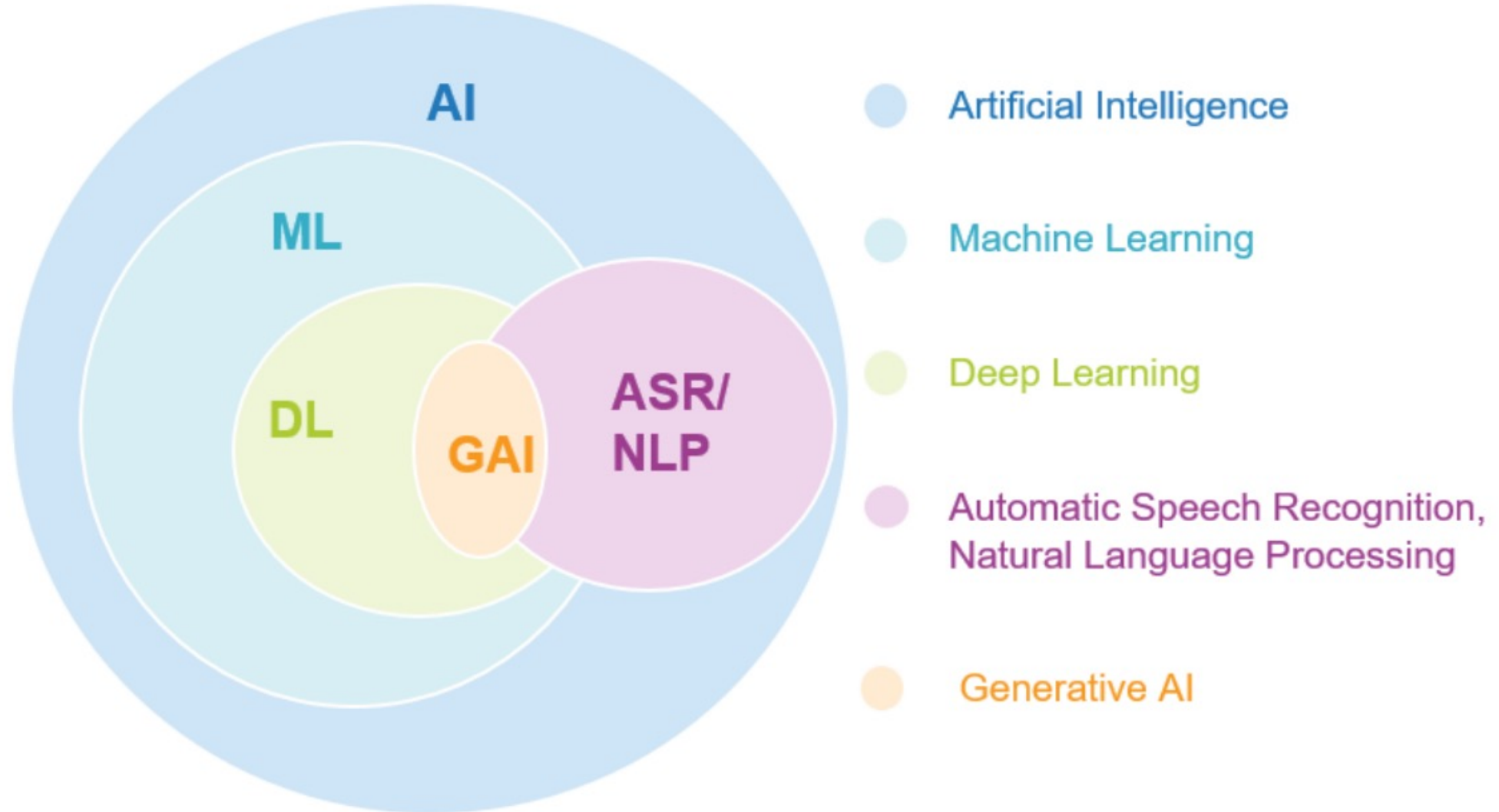
# **(LLMs)**

# **Foundation Models**

# Generative AI (Gen AI) AI Generated Content (AIGC)

# Artificial Intelligence (AI)

# AI, ML, DL, Generative AI



# Generative AI, Agentic AI, Physical AI

## Physical AI

Self-driving cars  
General robotics

## Agentic AI

Coding assistants  
Customer service  
Patient care

## Generative AI

Digital marketing  
Content creation

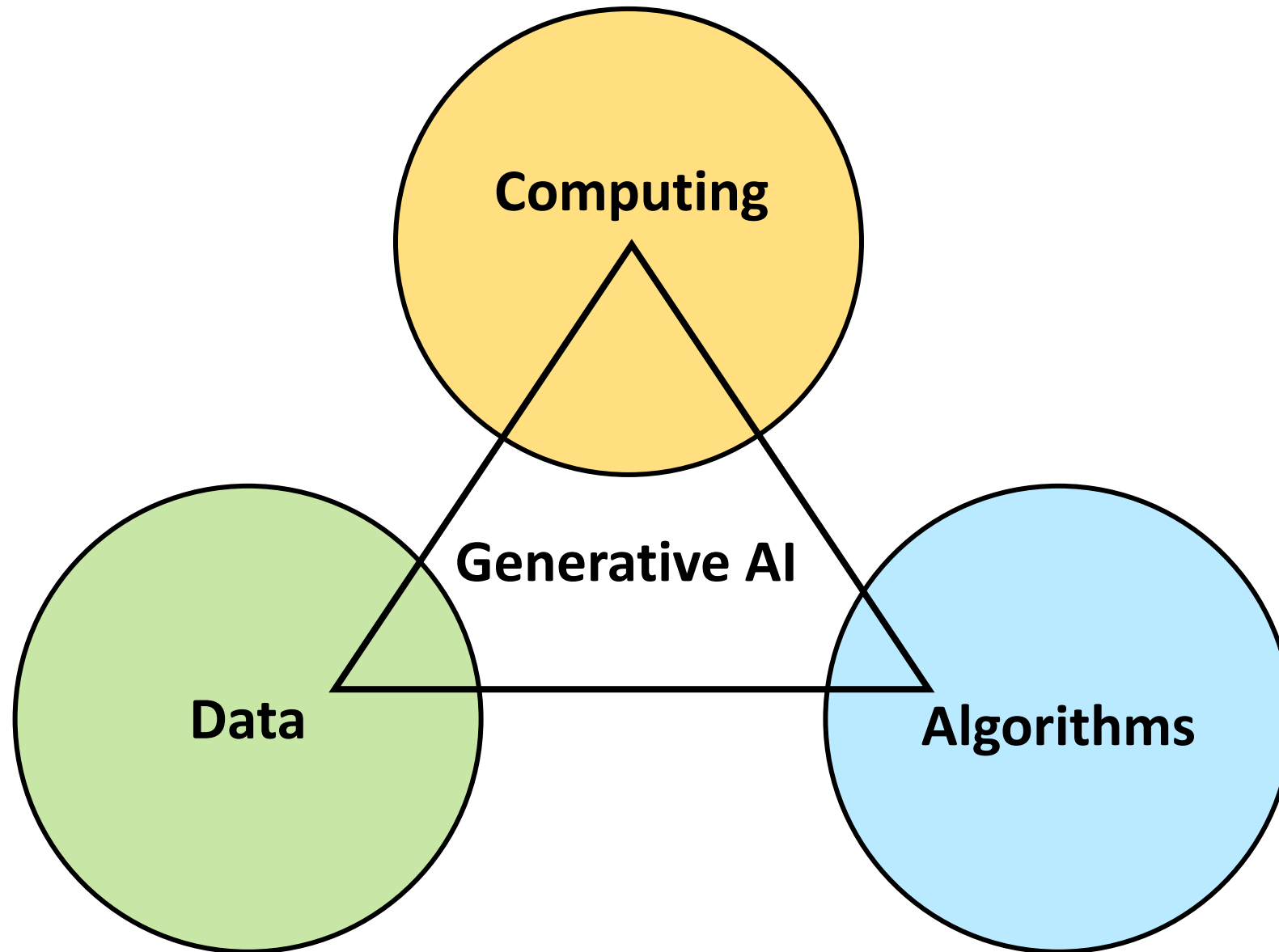
## Perception AI

Speech recognition  
Deep recommender systems  
Medical imaging

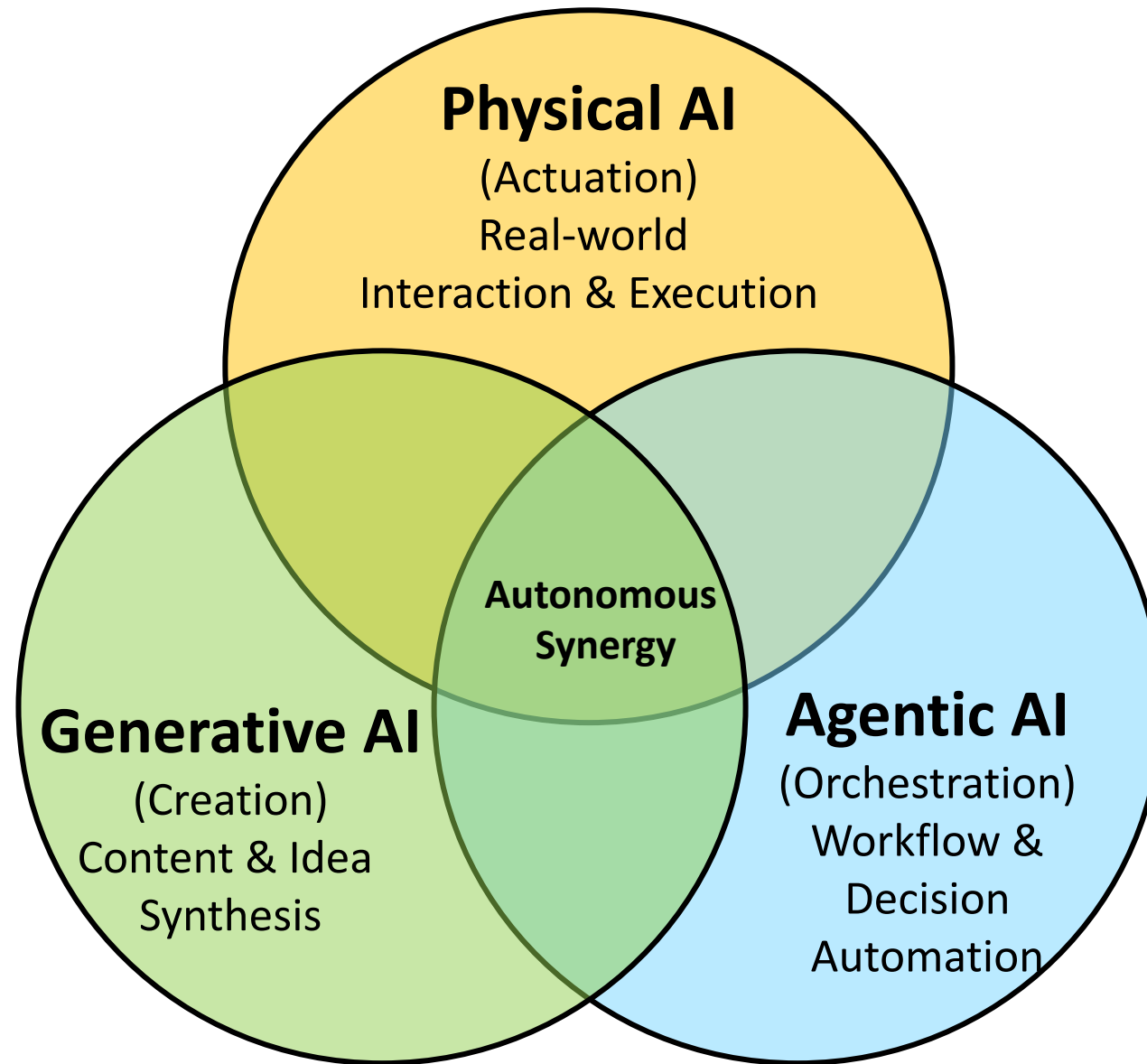
## 2012 AlexNet

Deep learning breakthrough

# Generative AI



# Generative AI, Agentic AI, Physical AI

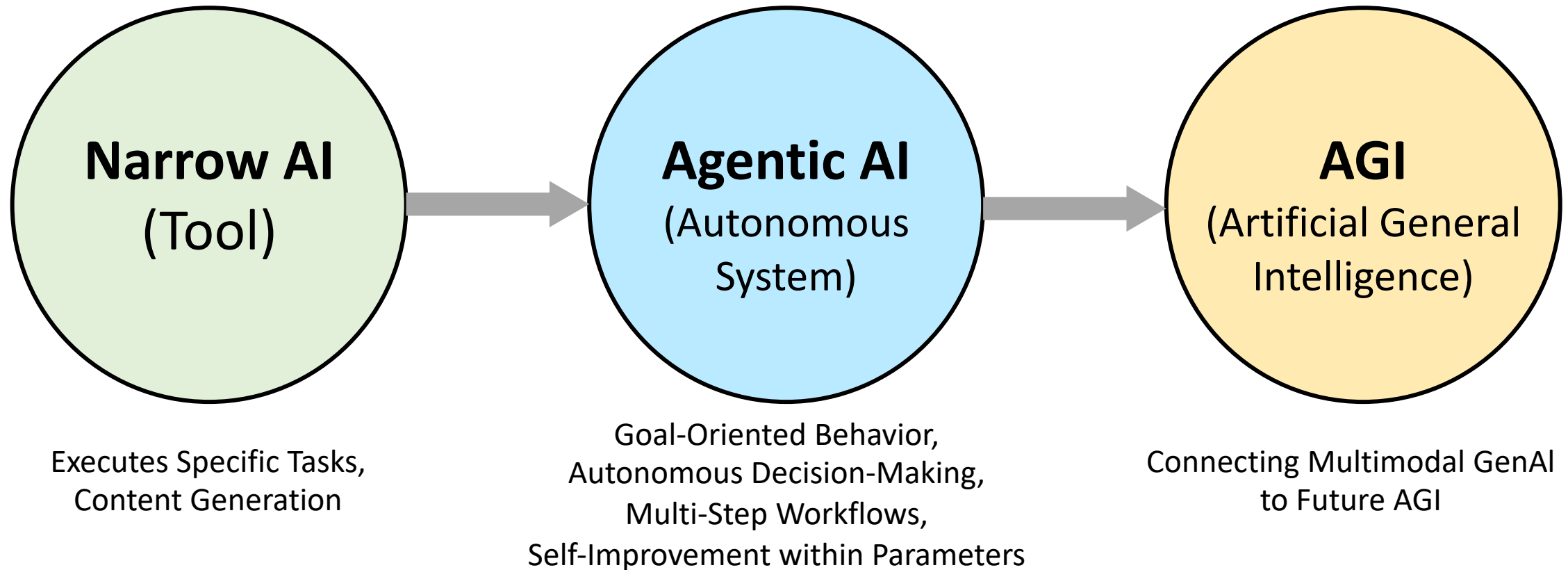


**New Economic  
Paradigm Shift:  
From Creation  
to Execution**

# The Future of AI

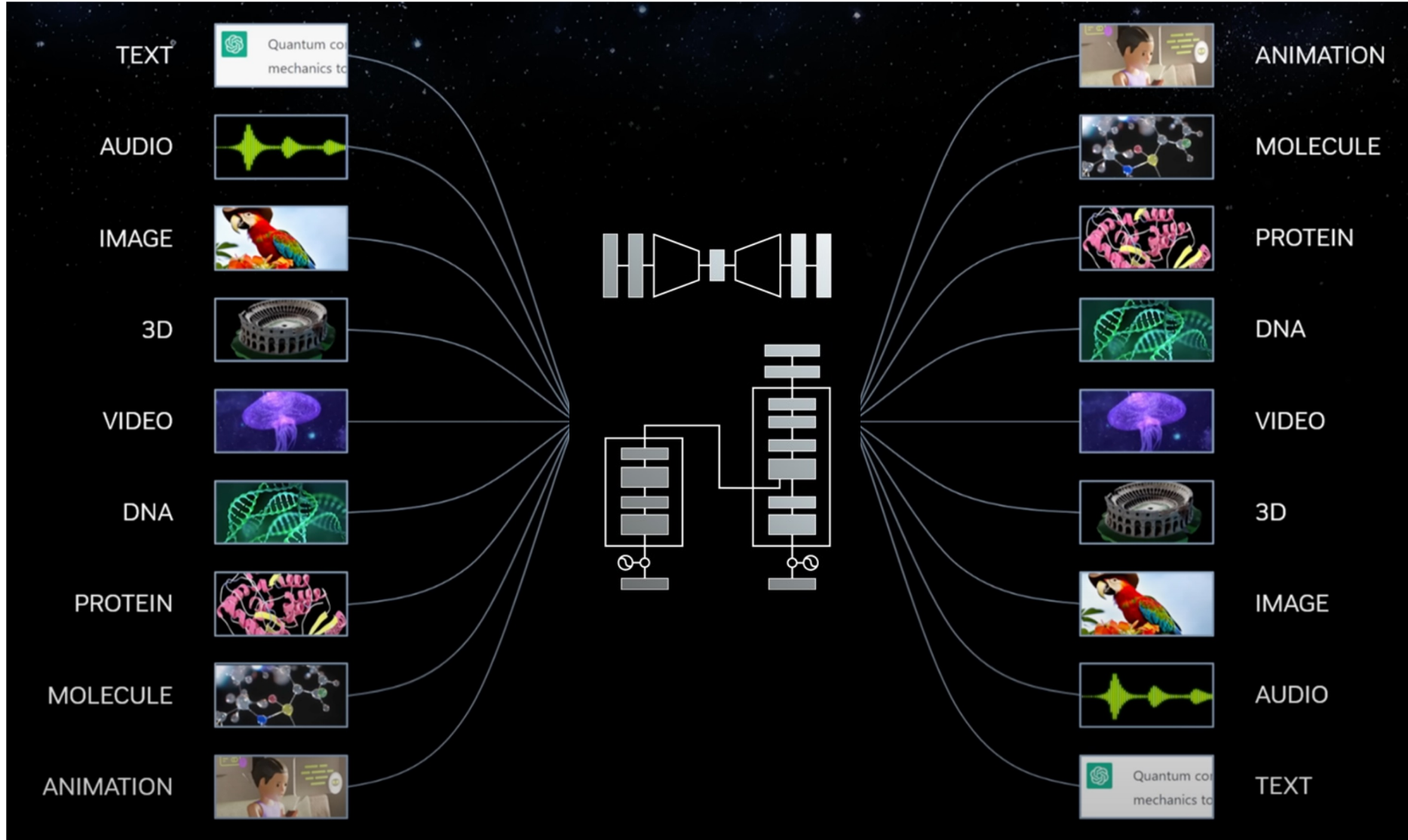
## From Tools to Agents:

### The Rise and Autonomy of Agentic AI

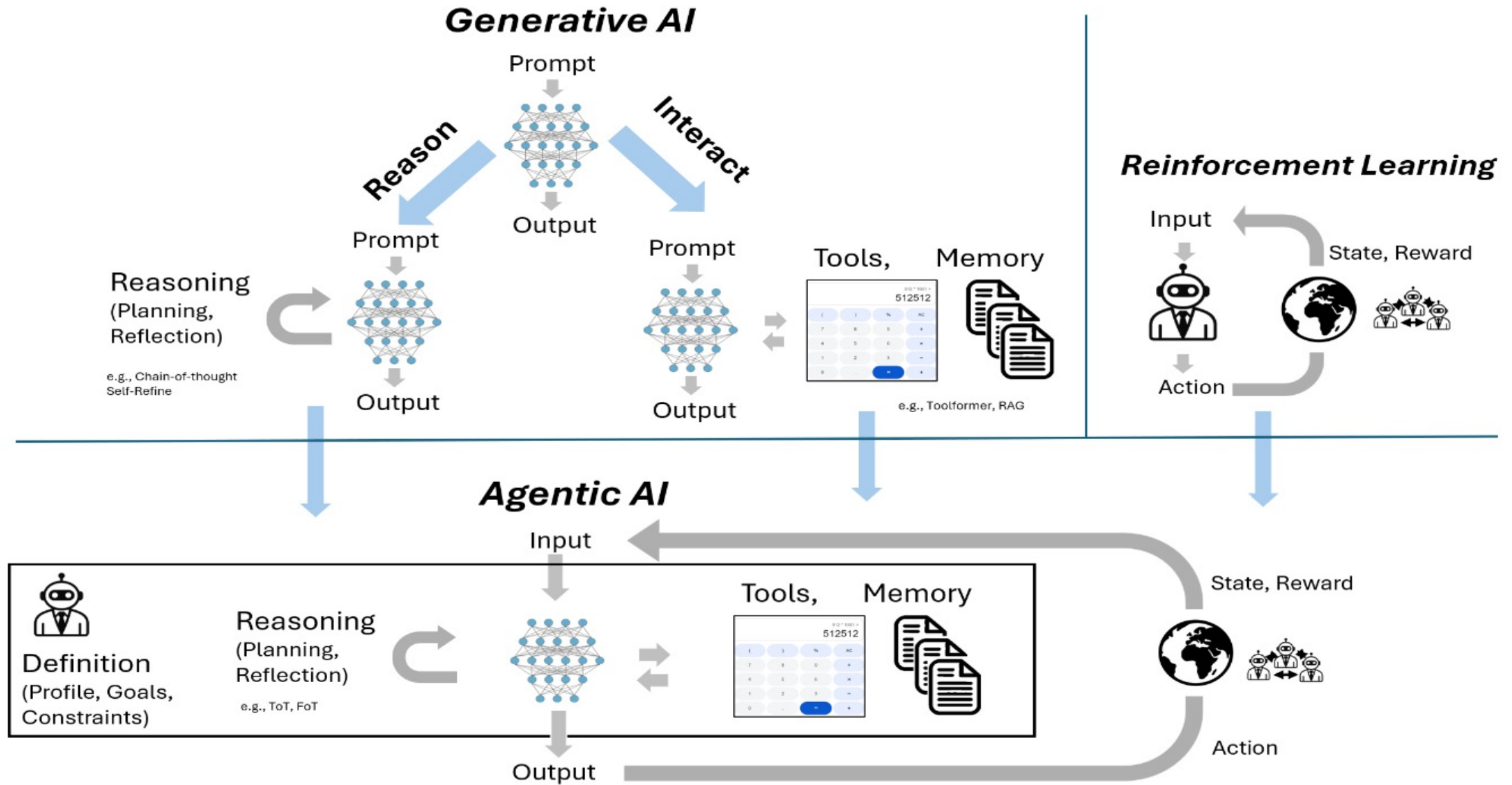


# Modular Modalities

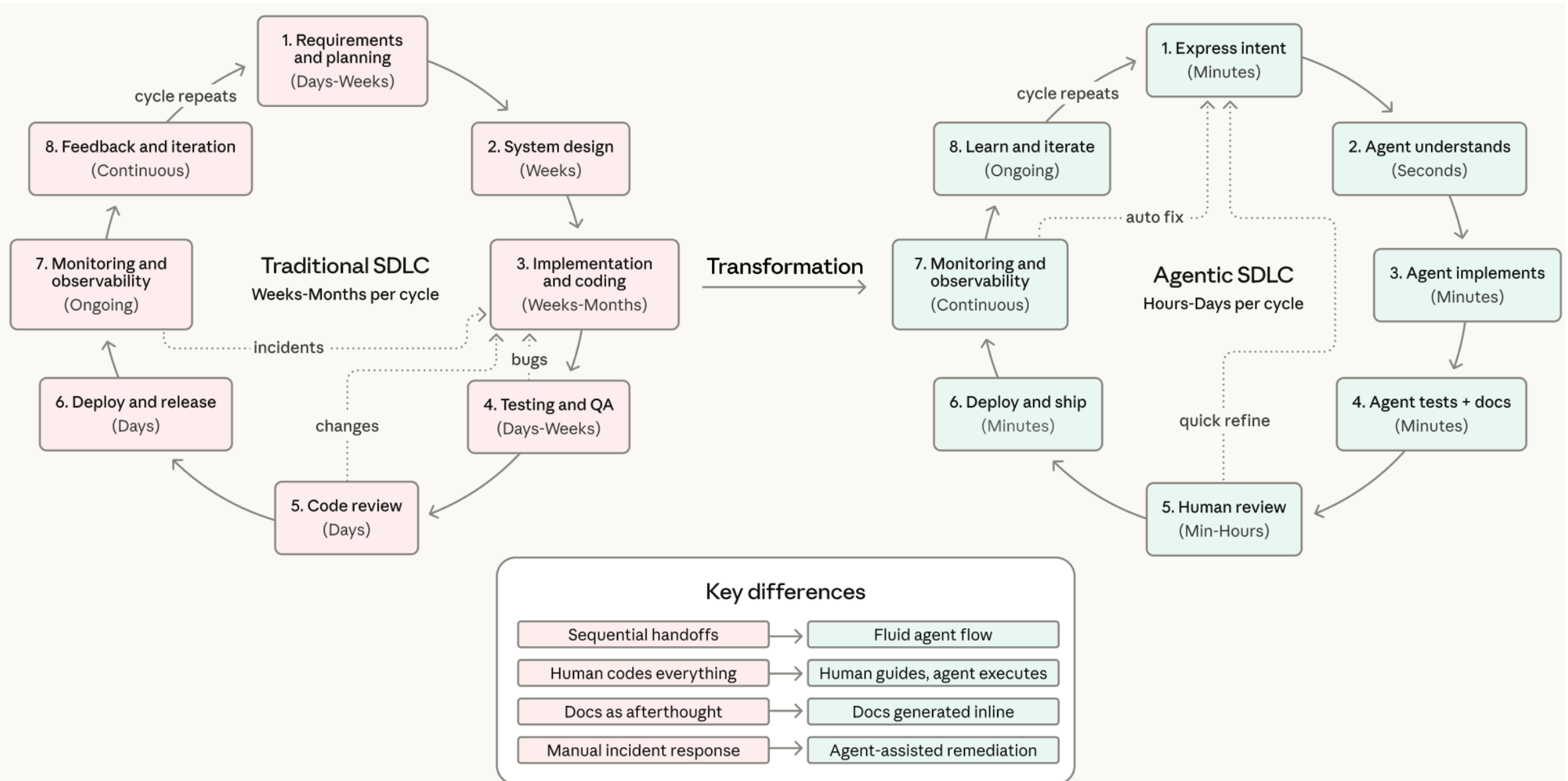
## Where Can The Transformer Fit?



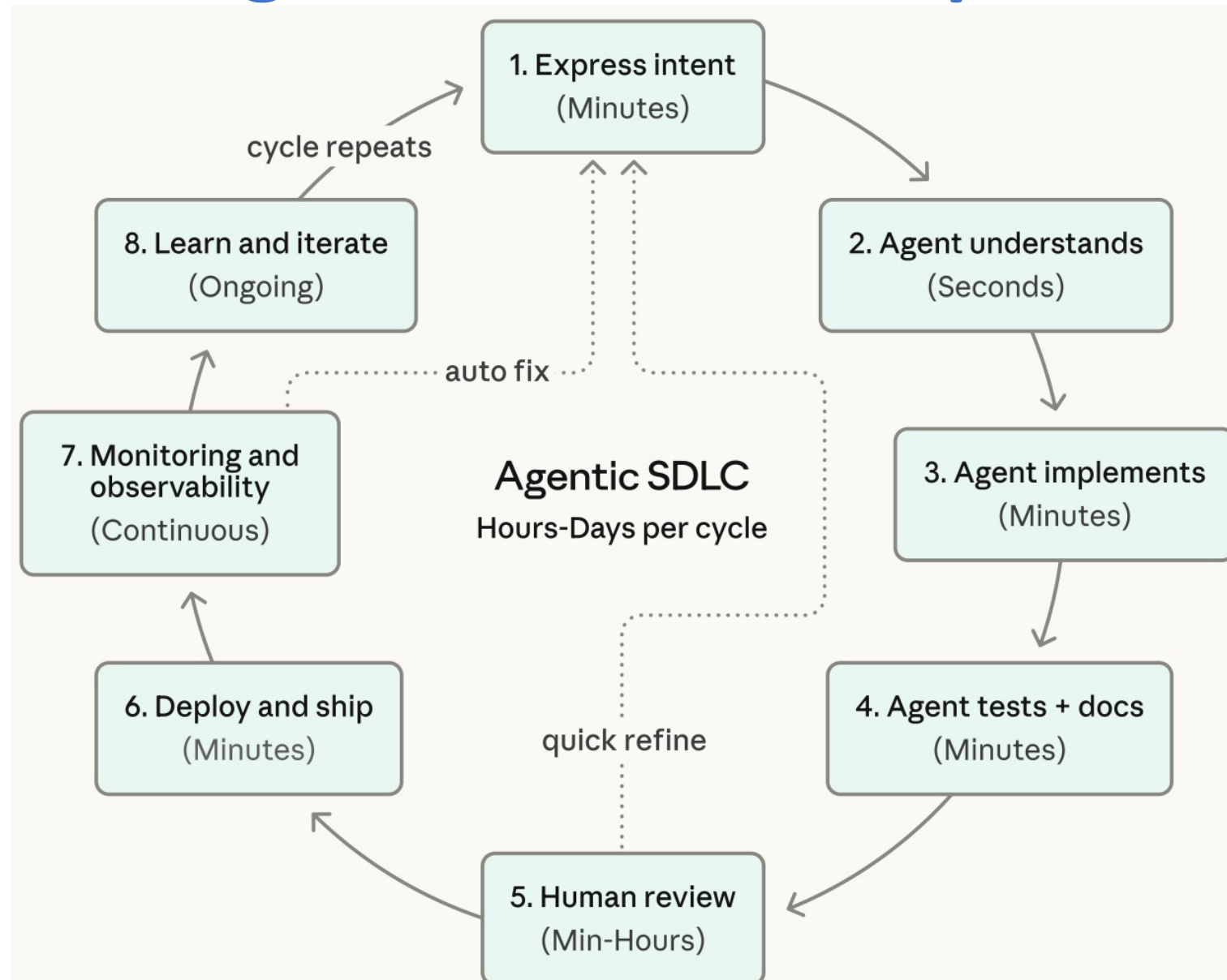
# From Generative AI to Agentic AI



# Agentic Coding Software Development Lifecycle

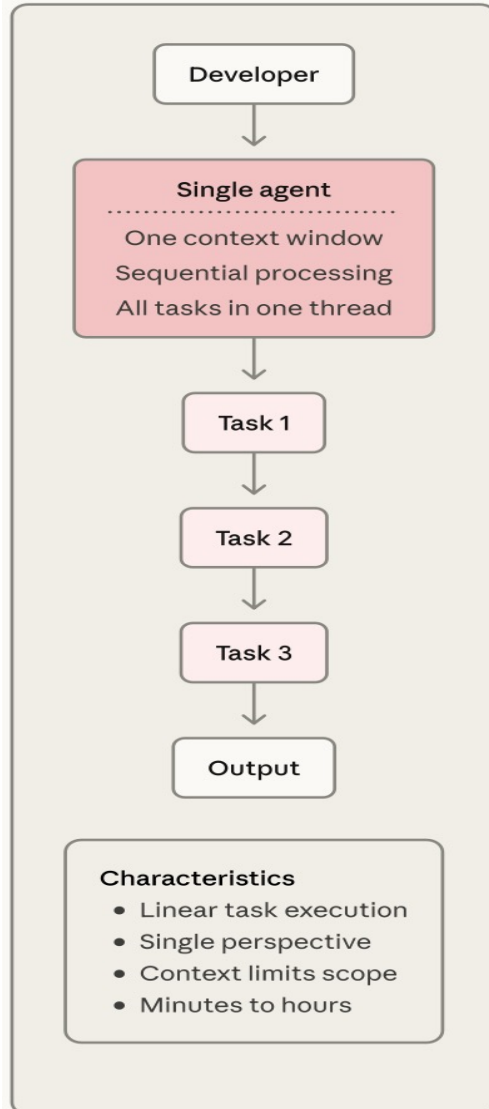


# Agentic Coding Software Development Lifecycle



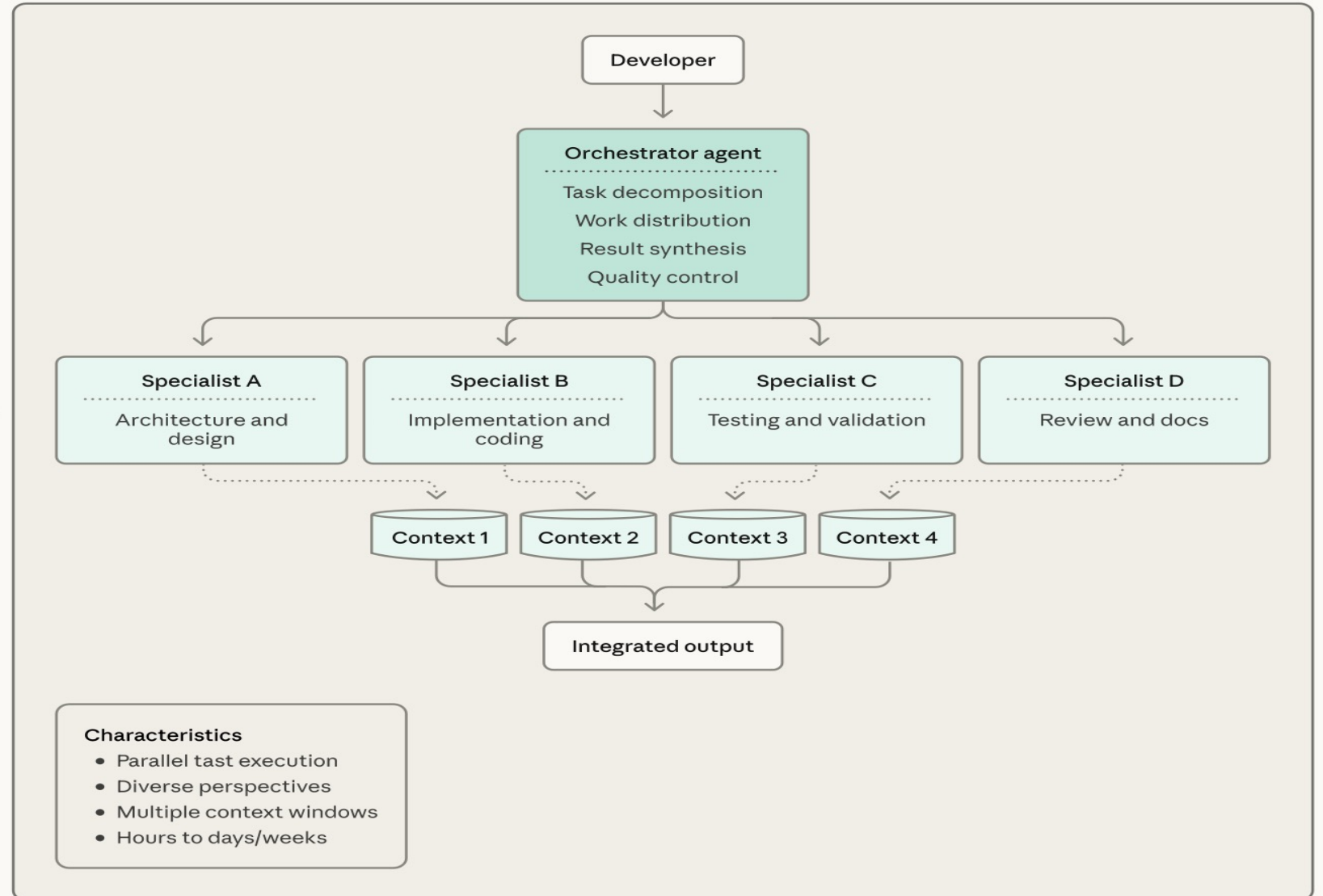
# Coding Agent Architectures: From Single Agents to Coordinated Teams

Single agent architecture

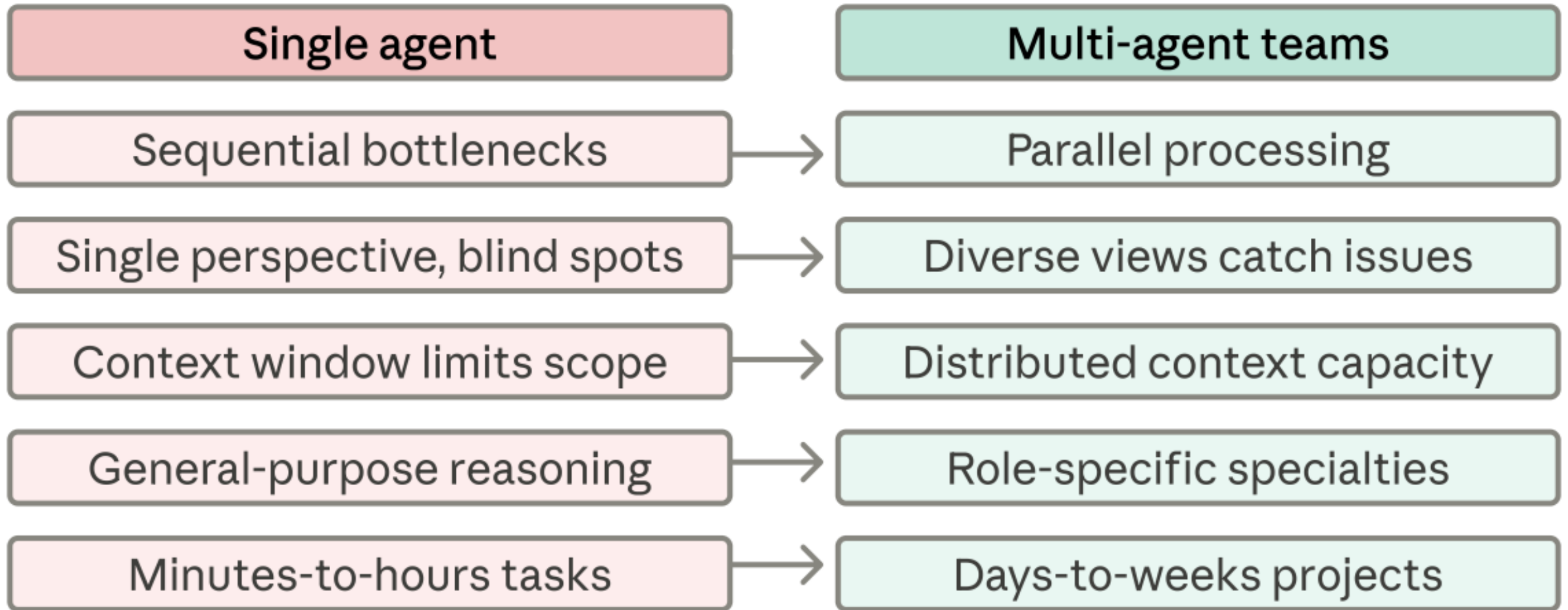


Evolution  
Task horizons  
expand from  
minutes to weeks

Multi-agent hierarchical architecture

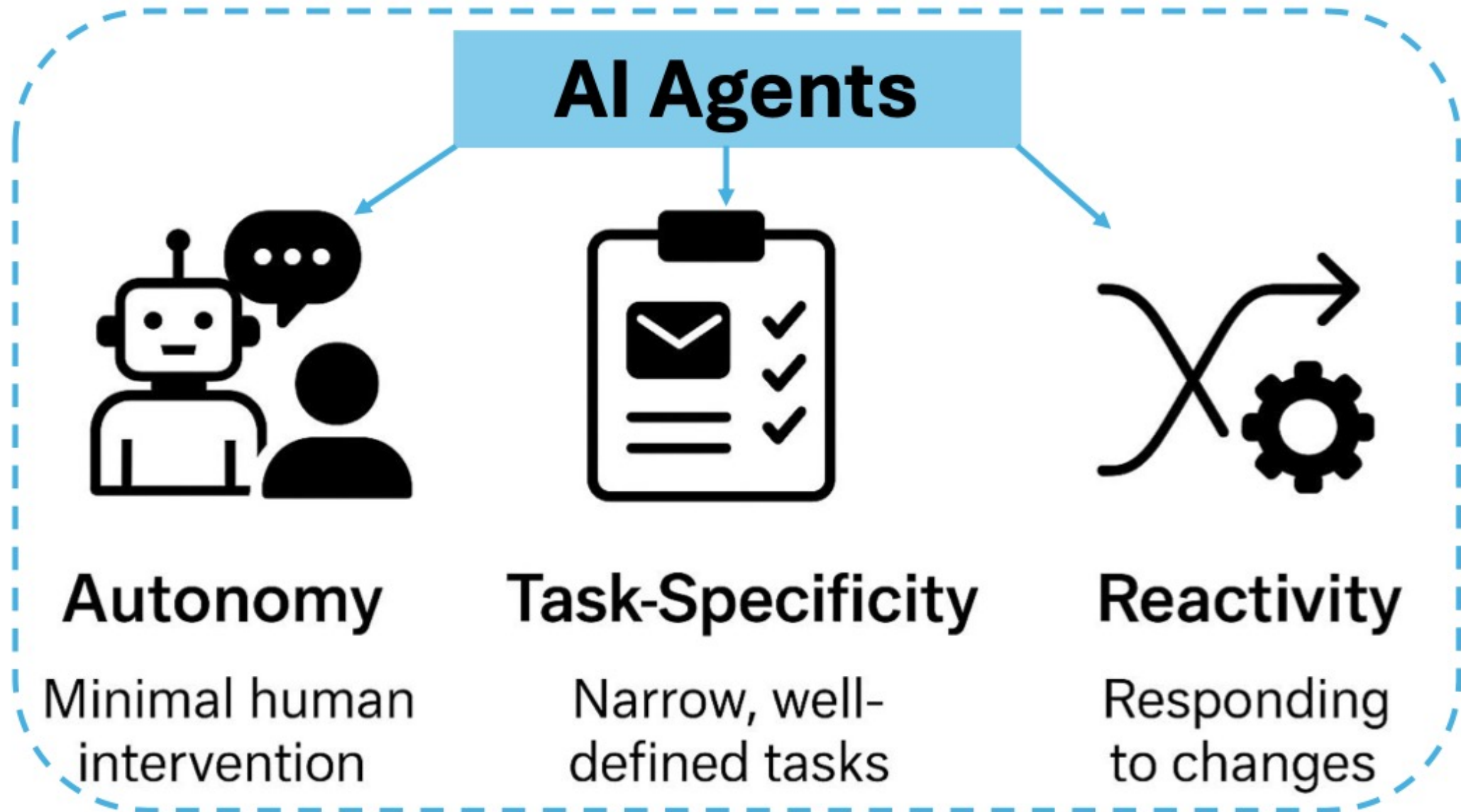


# Coding Agent Architectures: Single Agent to Multi-Agent Teams Performance Impact

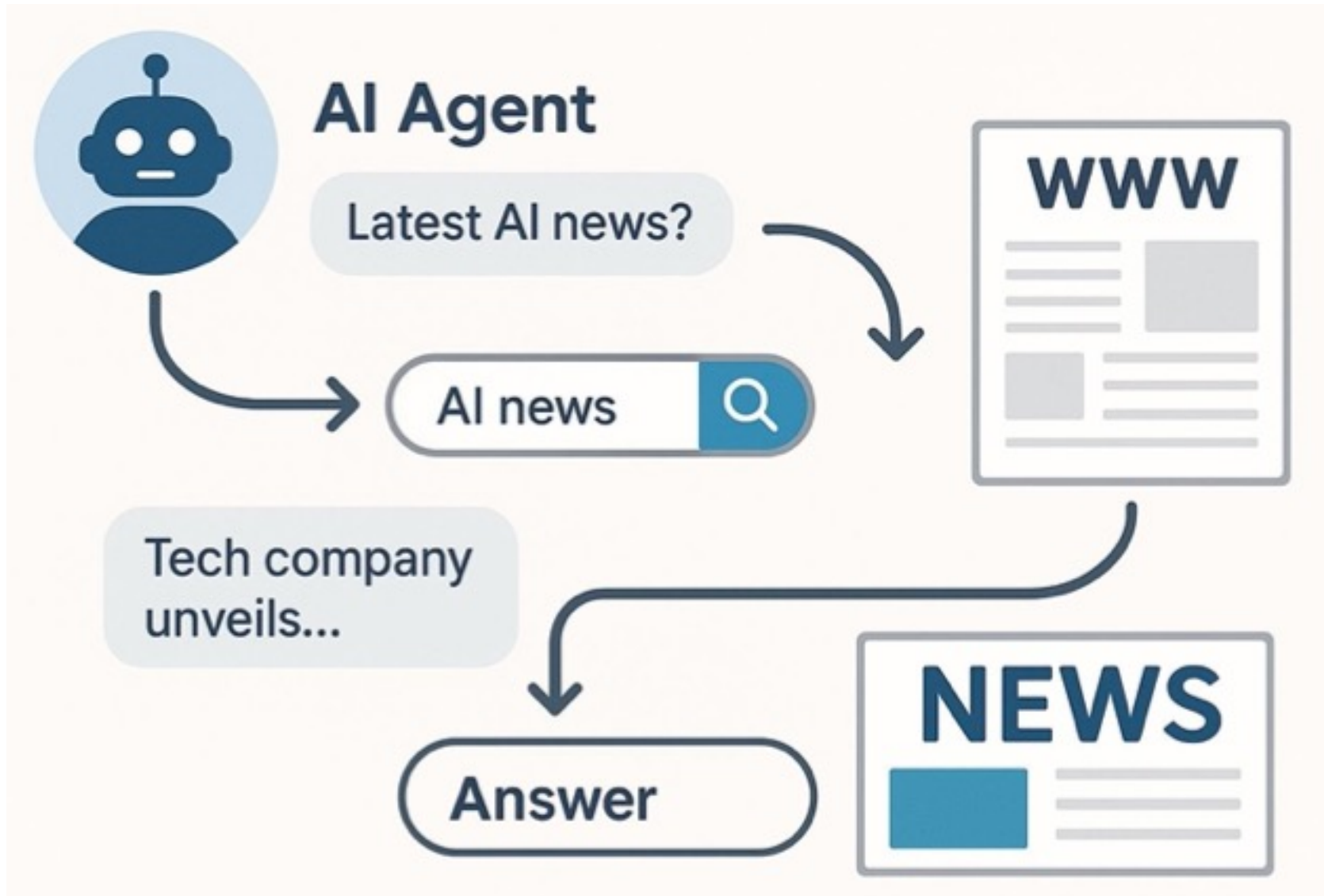


# Agentic AI

# AI Agents



# AI Agents



# Comparison of Generative AI and Traditional AI

Feature	Generative AI	Traditional AI
Output type	New content	Classification/Prediction
Creativity	High	Low
Interactivity	Usually more natural	Limited

# AI Agent / Agentic AI, Generative AI, Traditional AI

Feature	AI Agent / Agentic AI	Generative AI	Traditional AI
<b>Core Concept</b>	To autonomously perceive its environment, make decisions, and take actions to achieve specific goals.	To create new, original content (text, images, code, etc.) that resembles its training data.	To execute specific tasks based on pre-programmed rules or statistical patterns.
<b>Primary Function</b>	Action & Goal Achievement. Executes a series of tasks to complete an objective (e.g., "Book me a flight to Taipei next Tuesday.>").	Creation & Synthesis. Creates novel outputs in response to a prompt (e.g., "Write a poem about rain.>").	Classification & Prediction. Answers questions with a known range of outcomes (e.g., "Is this spam?").
<b>Decision Making</b>	Based on a continuous loop: Perceive -> Plan -> Act. It reasons about its goal, breaks it down, and executes steps.	Based on probabilistic patterns learned from massive, unstructured datasets. It predicts the next most likely word, pixel, or note.	Based on explicitly programmed logic (if-then rules) or learned patterns from structured data.
<b>Key Characteristic</b>	Autonomous & Goal-Oriented. Proactively takes steps and can adapt its plan based on new information.	Creative & Probabilistic. Can produce a wide variety of unique outputs from the same prompt.	Deterministic & Logic-Based. Given the same input, it will almost always produce the same output.

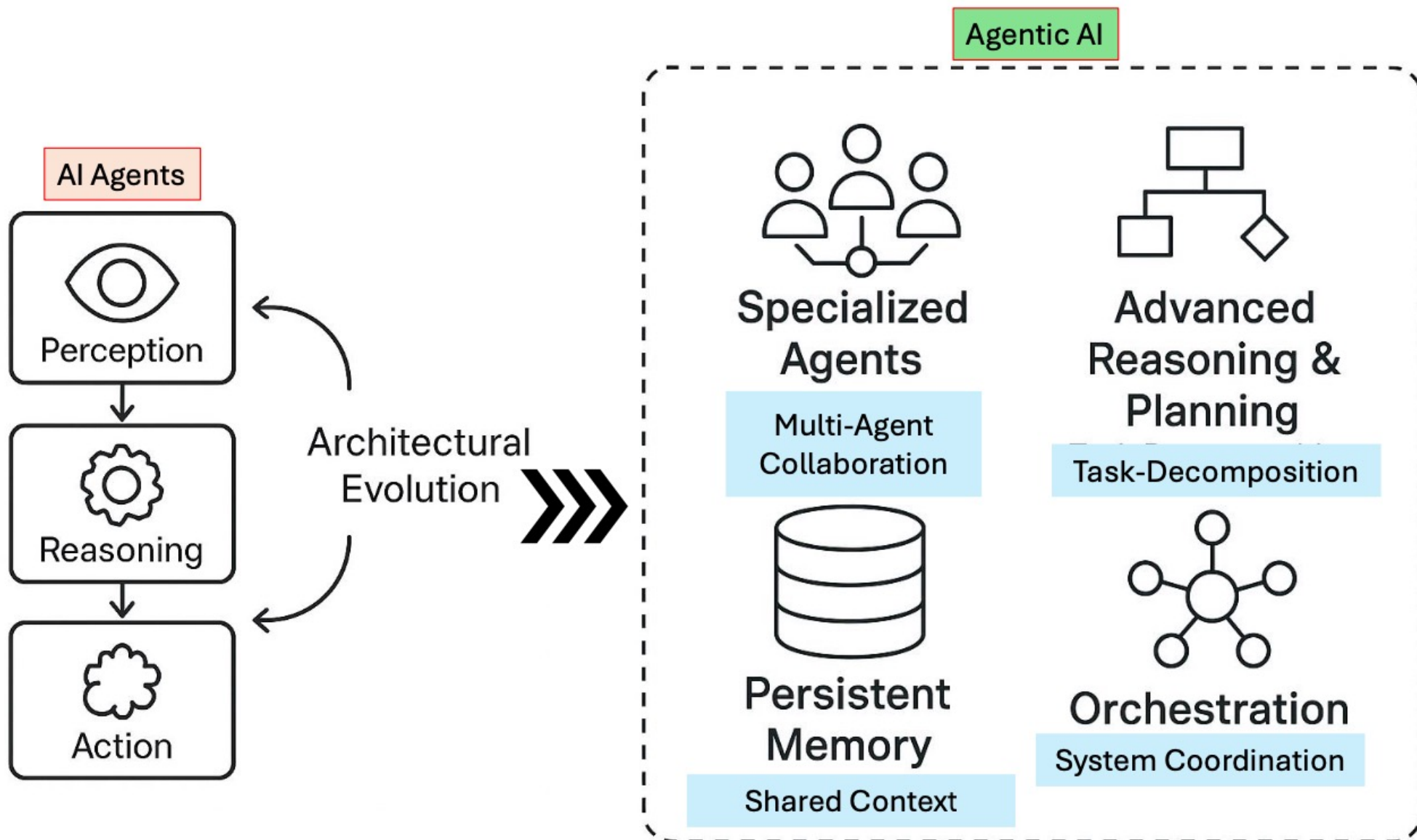
# AI Agent / Agentic AI, Generative AI, Traditional AI

Feature	AI Agent / Agentic AI	Generative AI	Traditional AI
<b>Interaction Model</b>	Proactive & Interactive. Actively observes its environment (digital or physical) and takes actions to change it.	Responsive. Engages in a dialogue or responds to a user's prompt to generate content.	Reactive. Responds to a direct input or query. It doesn't act on its own.
<b>Example Technologies</b>	Architectural frameworks like ReAct (Reason + Act), and systems that combine LLMs with tools and memory.	Large Language Models (LLMs) like GPT-4, Diffusion Models (for images), Generative Adversarial Networks (GANs).	Expert systems, decision trees, linear regression, traditional machine learning (ML) models.
<b>Common Use Cases</b>	Self-driving cars, autonomous trading bots, smart assistants that manage calendars, customer service agents that process refunds.	ChatGPT, Google Gemini, Midjourney (image generation), Copilot (code generation), music composition.	Spam filters, chess engines, recommendation systems (e.g., Netflix), credit scoring, medical diagnosis from scans.
<b>Relationship to Others</b>	An architecture or system that often uses Generative AI to reason and Traditional AI for specific sub-tasks to accomplish a goal.	Can serve as the "brain" or reasoning engine for an AI Agent, enabling it to understand, plan, and generate actions.	The foundation for modern AI. Its techniques can be components within larger AI systems.

# AI Agents vs Agentic AI

Feature	AI Agents	Agentic AI
<b>Definition</b>	Autonomous software programs that perform specific tasks.	Systems of multiple AI agents collaborating to achieve complex goals.
<b>Autonomy Level</b>	High autonomy within specific tasks.	Broad level of autonomy with the ability to manage multi-step, complex tasks and systems.
<b>Task Complexity</b>	Typically handle single, specific tasks.	Handle complex, multi-step tasks requiring coordination.
<b>Collaboration</b>	Operate independently.	Involve multi-agent information sharing, collaboration and cooperation.
<b>Learning and Adaptation</b>	Learn and adapt within their specific domain.	Learn and adapt across a wider range of tasks and environments.
<b>Applications</b>	Customer service chatbots, virtual assistants, automated workflows.	Supply chain management, business process optimization, virtual project managers.

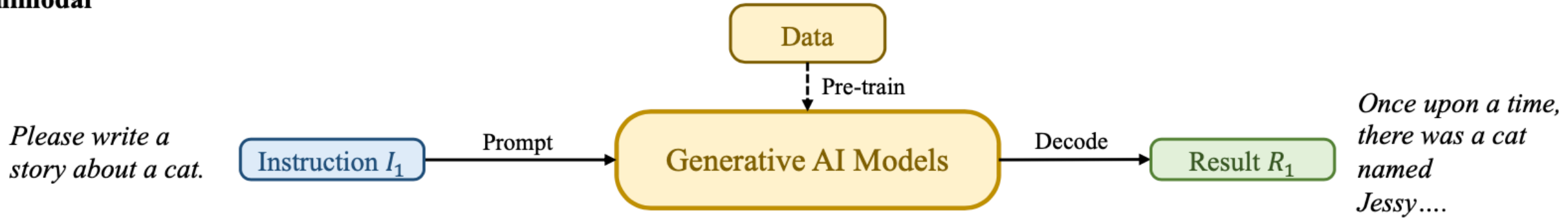
# AI Agents vs Agentic AI



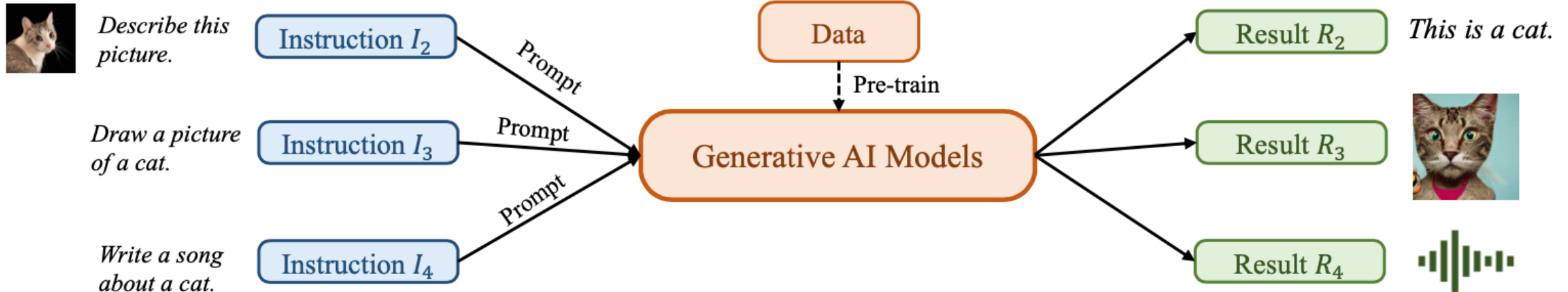
# Generative AI (Gen AI)

## AI Generated Content (AIGC)

### Unimodal

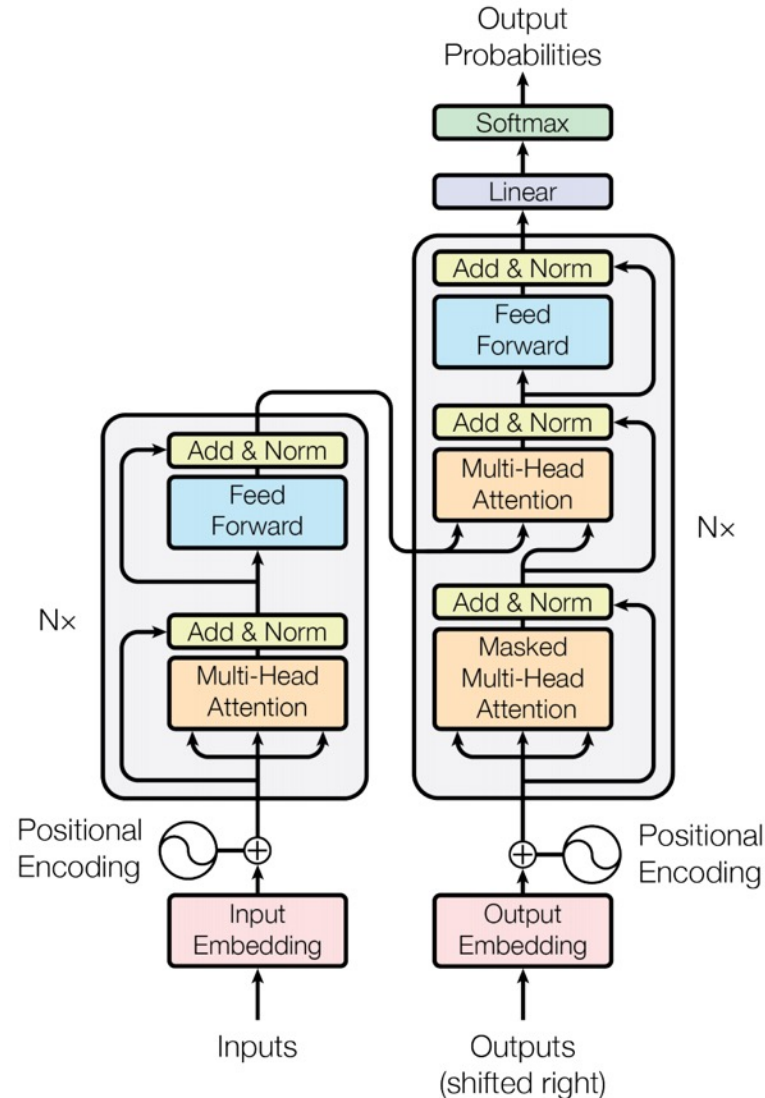


### Multimodal



# Transformer (Attention is All You Need)

(Vaswani et al., 2017)

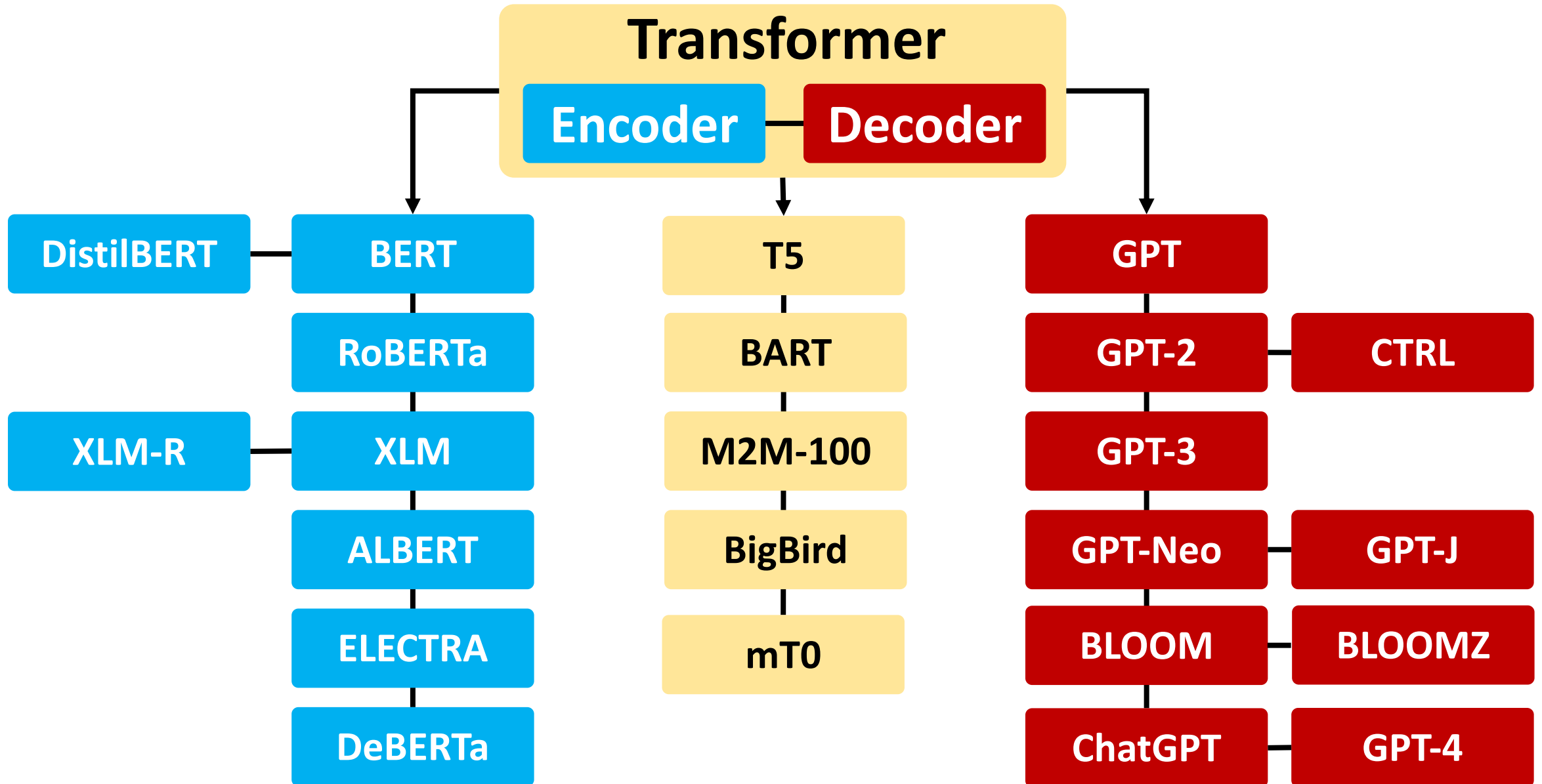


# Transformer (Attention is All You Need)

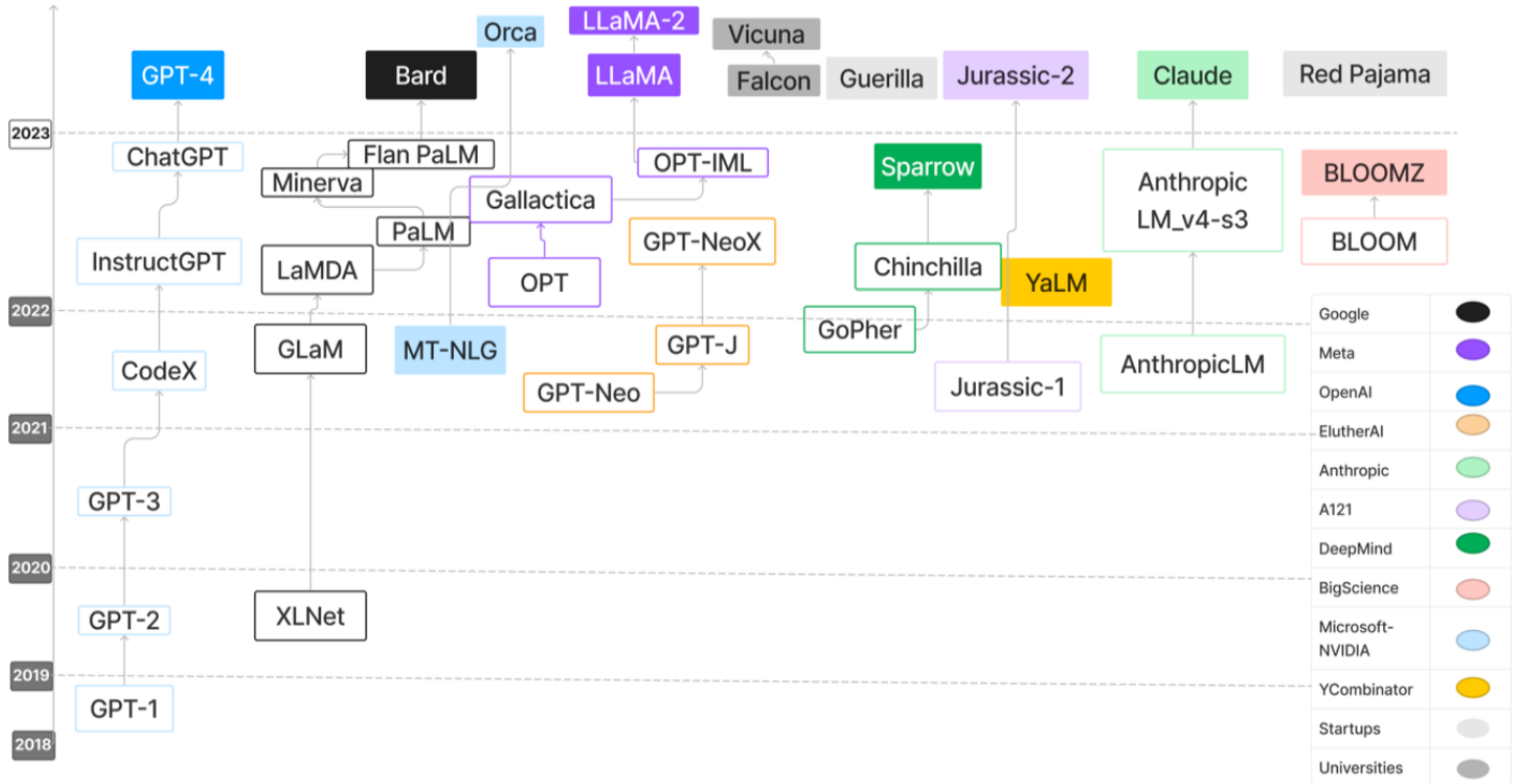
(Vaswani et al., 2017)

- A **Transformer** is a type of **deep learning model** introduced in the paper "**Attention Is All You Need**" (Vaswani et al., 2017).
- It revolutionized **Natural Language Processing (NLP)** by replacing traditional **sequence models** like **RNNs and LSTMs** with a **self-attention mechanism** that enables highly parallelizable training.

# Transformer Models



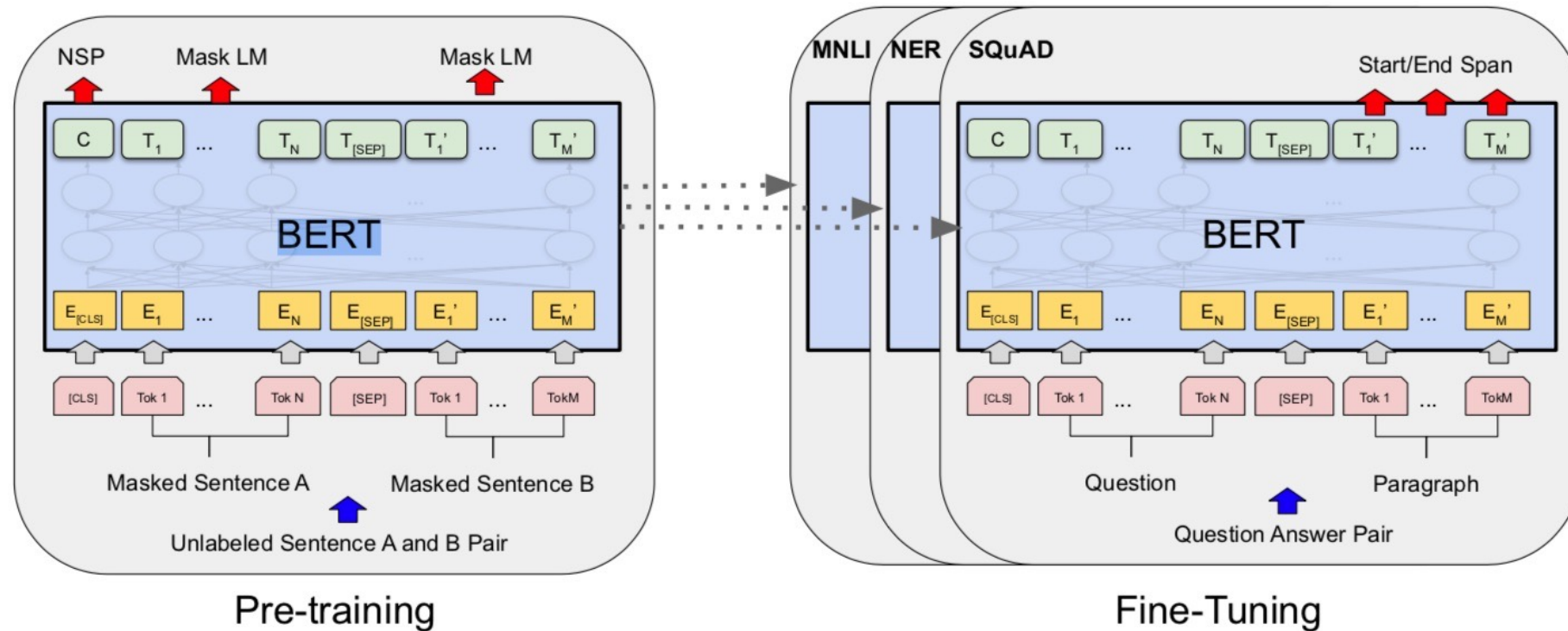
# Large Language Models (LLMs)



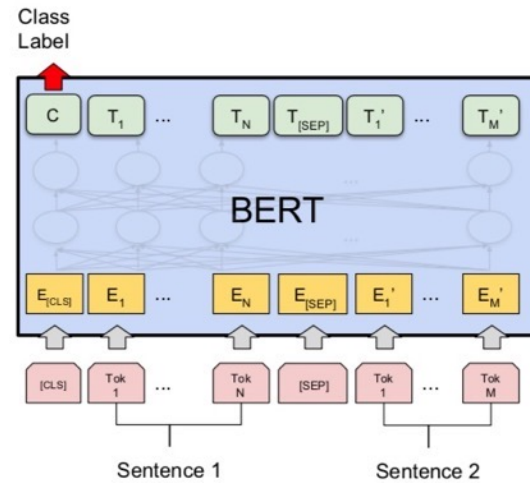
# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT (Bidirectional Encoder Representations from Transformers)

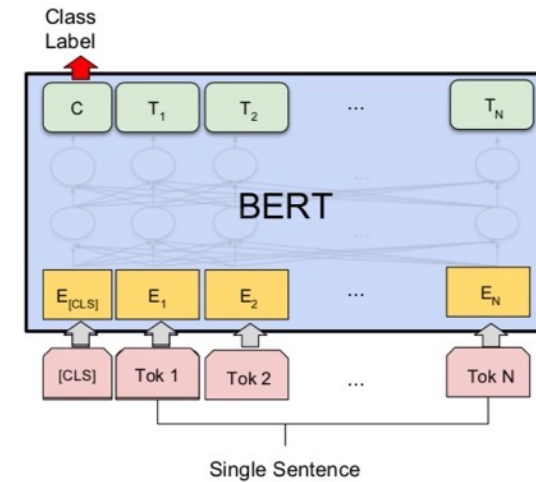
Overall pre-training and fine-tuning procedures for BERT



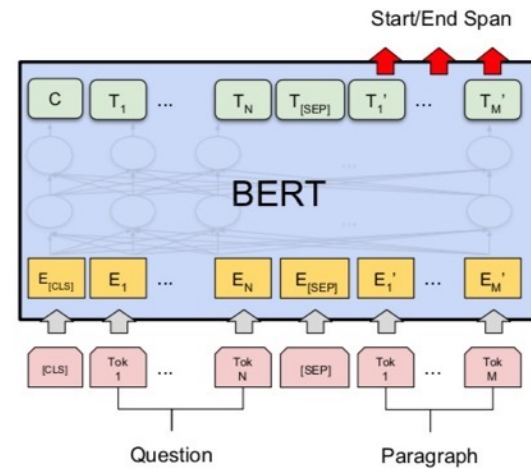
# Fine-tuning BERT on Different Tasks



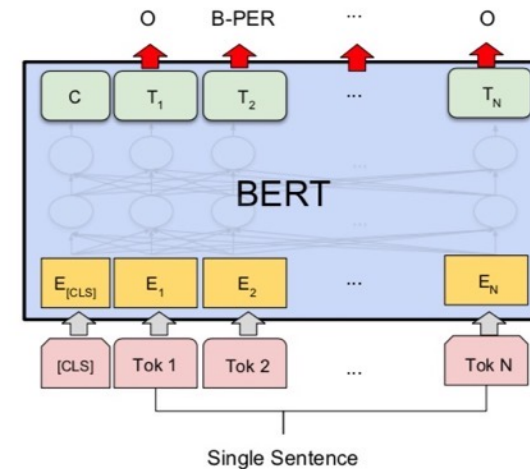
(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG



(b) Single Sentence Classification Tasks:  
SST-2, CoLA



(c) Question Answering Tasks:  
SQuAD v1.1

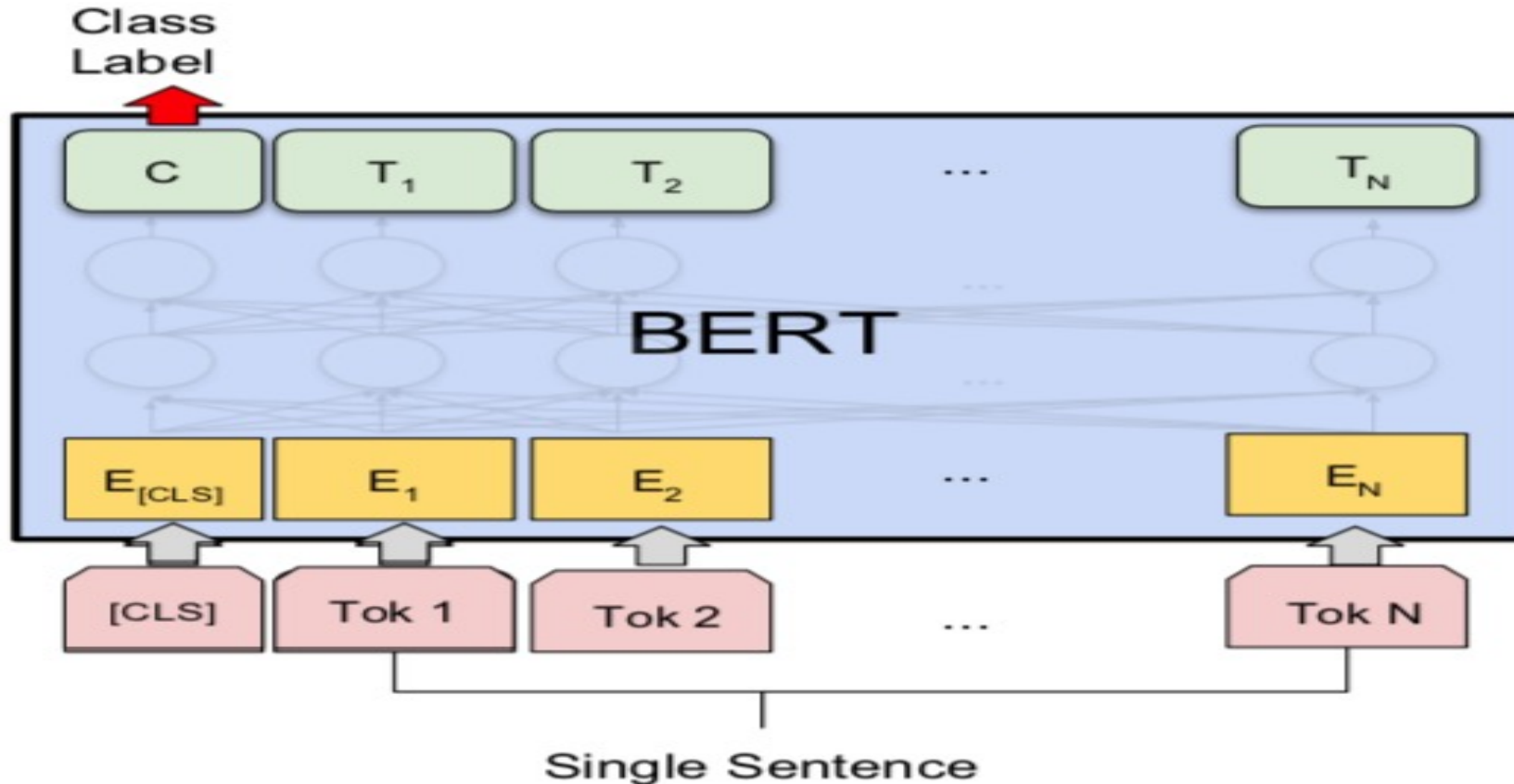


(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

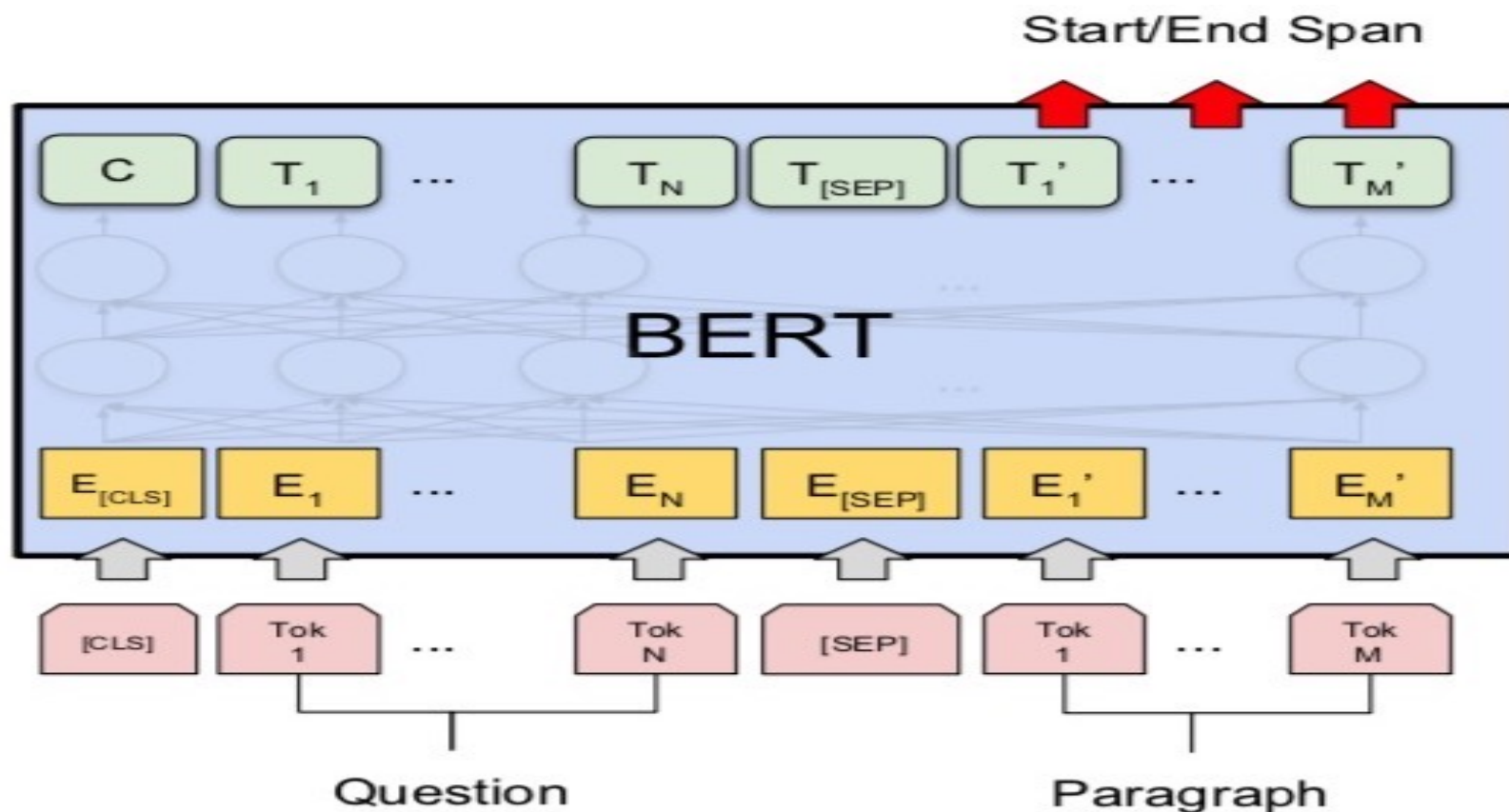
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Sentiment Analysis: Single Sentence Classification



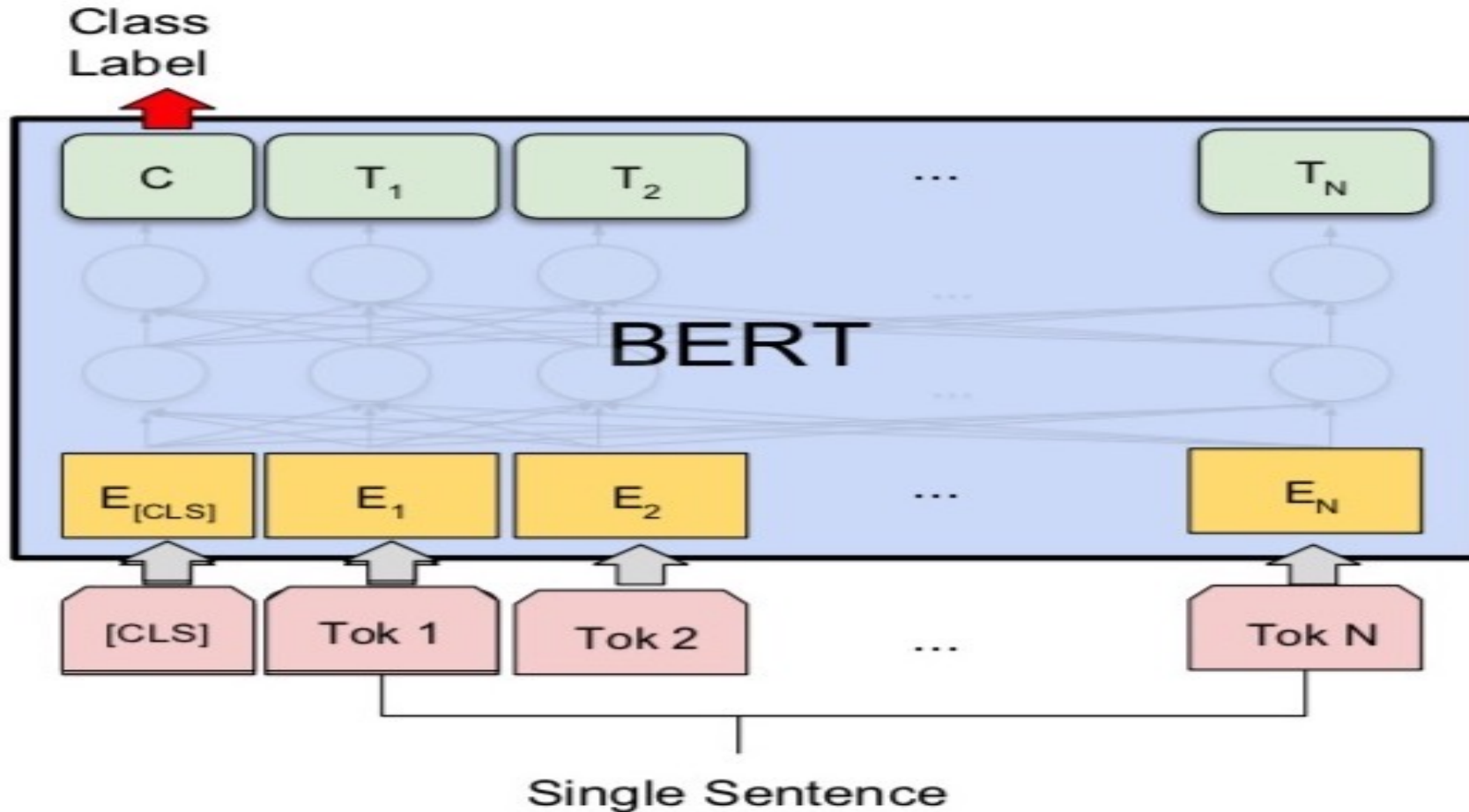
(b) Single Sentence Classification Tasks:  
SST-2, CoLA

# Fine-tuning BERT on Question Answering (QA)



(c) Question Answering Tasks:  
SQuAD v1.1

# Fine-tuning BERT on Dialogue Intent Detection (ID; Classification)



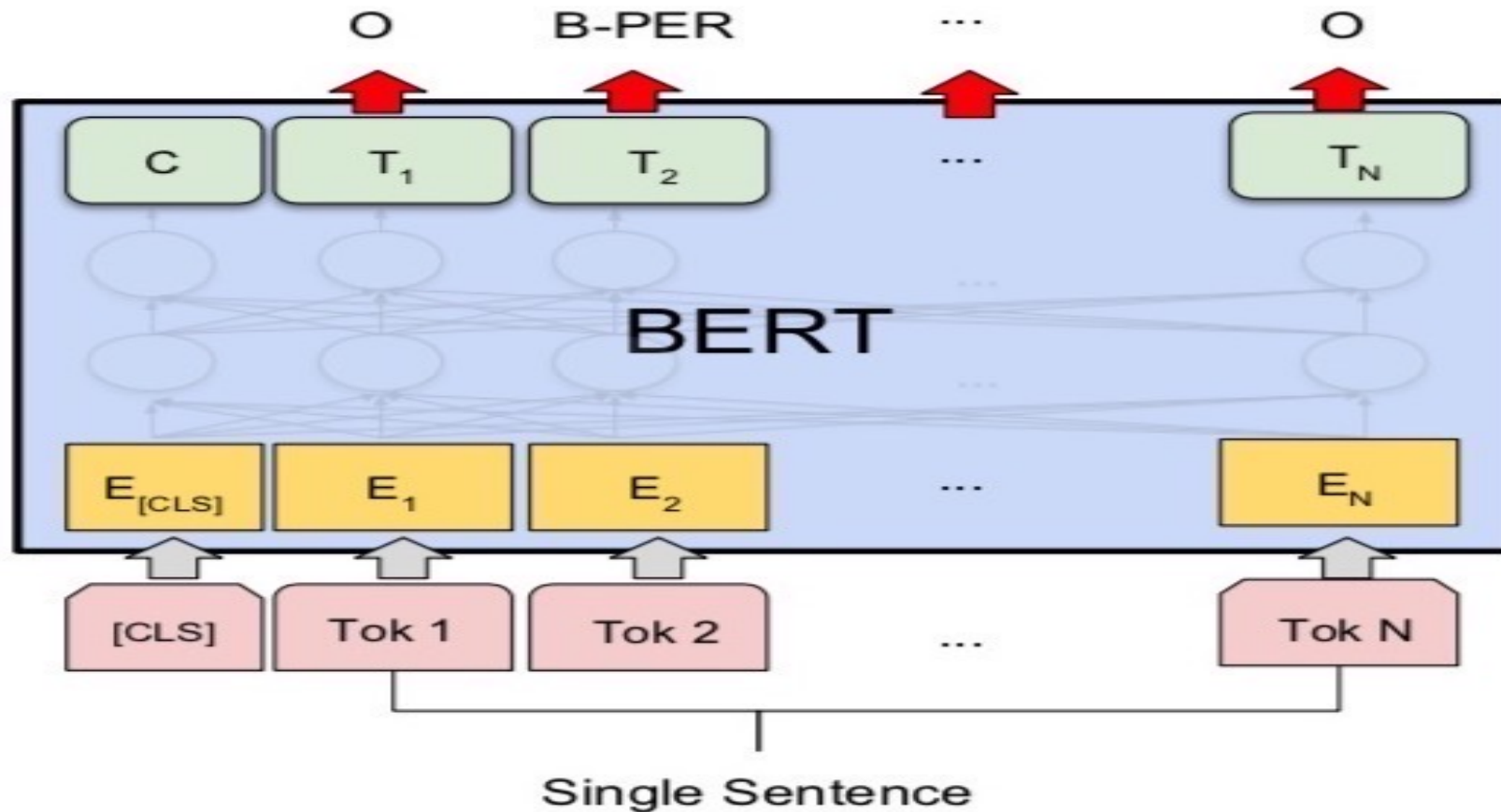
(b) Single Sentence Classification Tasks: SST-2, CoLA

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Fine-tuning BERT on Dialogue

## Slot Filling (SF)



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Key Features of the Transformer Model

- **Self-Attention Mechanism:** Allows the model to weigh the importance of different words in a sentence, regardless of their position.
- **Positional Encoding:** Since Transformers don't use recurrence (like RNNs), positional encodings are added to input embeddings to retain word order information.
- **Multi-Head Attention:** The model attends to different words simultaneously in multiple ways, capturing various relationships.
- **Feed-Forward Layers:** After attention, the output passes through dense layers for further transformation.
- **Layer Normalization & Residual Connections:** Improve gradient flow and training stability.
- **Encoder-Decoder Architecture:**
  - **Encoder:** Processes input text and converts it into contextual embeddings.
  - **Decoder:** Generates output text, often used in translation or text generation tasks.

# Popular Transformer-Based Models

- **BERT (Bidirectional Encoder Representations from Transformers)**
  - Used for tasks like classification and question answering.
- **GPT (Generative Pre-trained Transformer)**
  - Generates text based on input prompts.
- **T5 (Text-To-Text Transfer Transformer)**
  - Converts all NLP tasks into a text-to-text format.
- **ViT (Vision Transformer)**
  - Applies Transformer architecture to computer vision.

# Tokens in NLP (Text Processing)

Aspect	Description
Tokenization	Splitting text into meaningful units (words, subwords, or characters).
Subword Tokens	Methods like Byte Pair Encoding (BPE) and WordPiece break words into reusable subunits (e.g., "unhappiness" → ["un", "happiness"]).
Embeddings	Converts tokens into numerical vectors for processing.
Semantic Role Labeling (SRL)	Identifies sentence structure by assigning roles to tokens (e.g., subject, object).
Special Tokens	[CLS], [SEP]

# Tokens in CV (Image Processing)

Aspect	Description
<b>Patch Tokens</b>	Vision Transformers (ViT) split an image into small patches (e.g., 16x16 pixels), treating them as tokens.
<b>Positional Encoding</b>	Since images lack inherent sequence order like text, positional embeddings help ViTs understand spatial structure.
<b>Midjourney API Tokens</b>	Midjourney's AI processes text prompts into image tokens, converting descriptions into AI-generated art.
<b>CLIP Tokens</b>	OpenAI's CLIP model tokenizes both text and images, allowing cross-modal understanding (e.g., "dog" matches a picture of a dog).

# Tokens in NLP vs CV

Feature	NLP Tokens	CV Tokens
Basic Unit	Words, subwords, or characters	Image patches (e.g., 16×16 pixel grids)
Processing	Tokenized using BPE, WordPiece, SentencePiece	Tokenized as patch embeddings
Positional Encoding?	Needed to retain word order	Needed to retain spatial information
Example Models	BERT, GPT, T5, RAG	ViT, DINOv2, CLIP
Use Case	Text-based tasks (chatbots, summarization, RAG)	Vision-based tasks (image classification, AI art generation)

# Attention in NLP (Text Processing)

Aspect	Description
<b>Self-Attention</b>	Each token attends to every other token in a sentence, capturing dependencies across long text sequences.
<b>Scaled Dot-Product Attention</b>	Computes attention scores using query (Q), key (K), and value (V) vectors.
<b>Multi-Head Attention (MHA)</b>	Improves attention by using multiple attention heads that learn different relationships.
<b>Causal Attention (Decoder-Only Models)</b>	Restricts attention to past tokens only, enabling text generation without looking ahead.
<b>Cross-Attention</b>	The decoder attends to encoder outputs in seq-to-seq tasks like translation (e.g., T5, BART).
<b>Retrieval-Augmented Attention</b>	Fetches external knowledge before generating a response (RAG models).

# Attention in CV (Vision Processing)

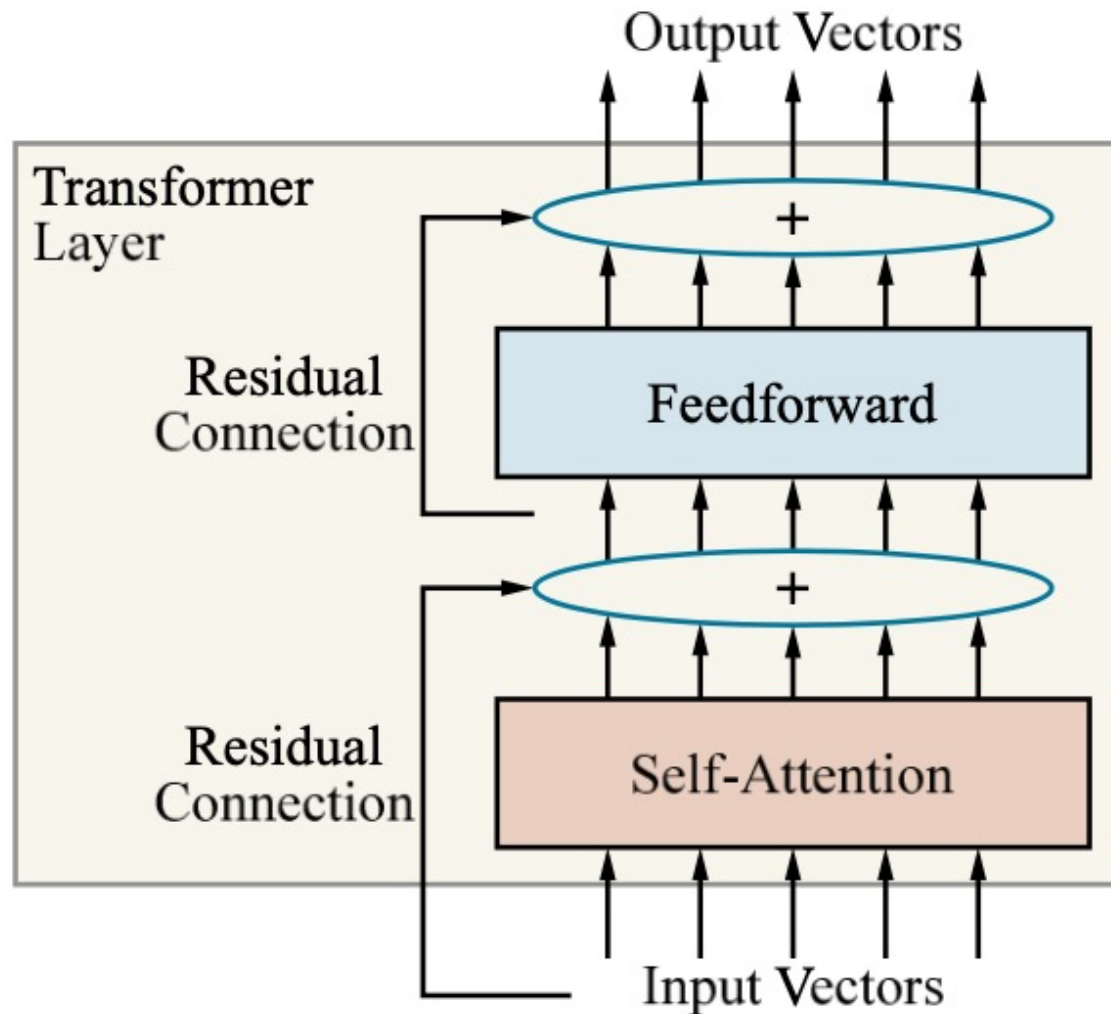
Aspect	Description
<b>Self-Attention in ViTs</b>	Treats images as a sequence of patches (like words in text) and applies attention to learn spatial relationships.
<b>Multi-Head Attention in ViTs</b>	Similar to NLP, multiple attention heads capture different visual features.
<b>Positional Encoding in Vision</b>	Since images lack inherent order, positional embeddings help maintain spatial structure.
<b>Cross-Attention in Multimodal AI</b>	Used in models like CLIP and Midjourney, where text descriptions attend to visual features.
<b>Attention Maps in Vision</b>	Heatmaps showing which image regions the model focuses on (e.g., for explainability).

# Attention in NLP vs CV

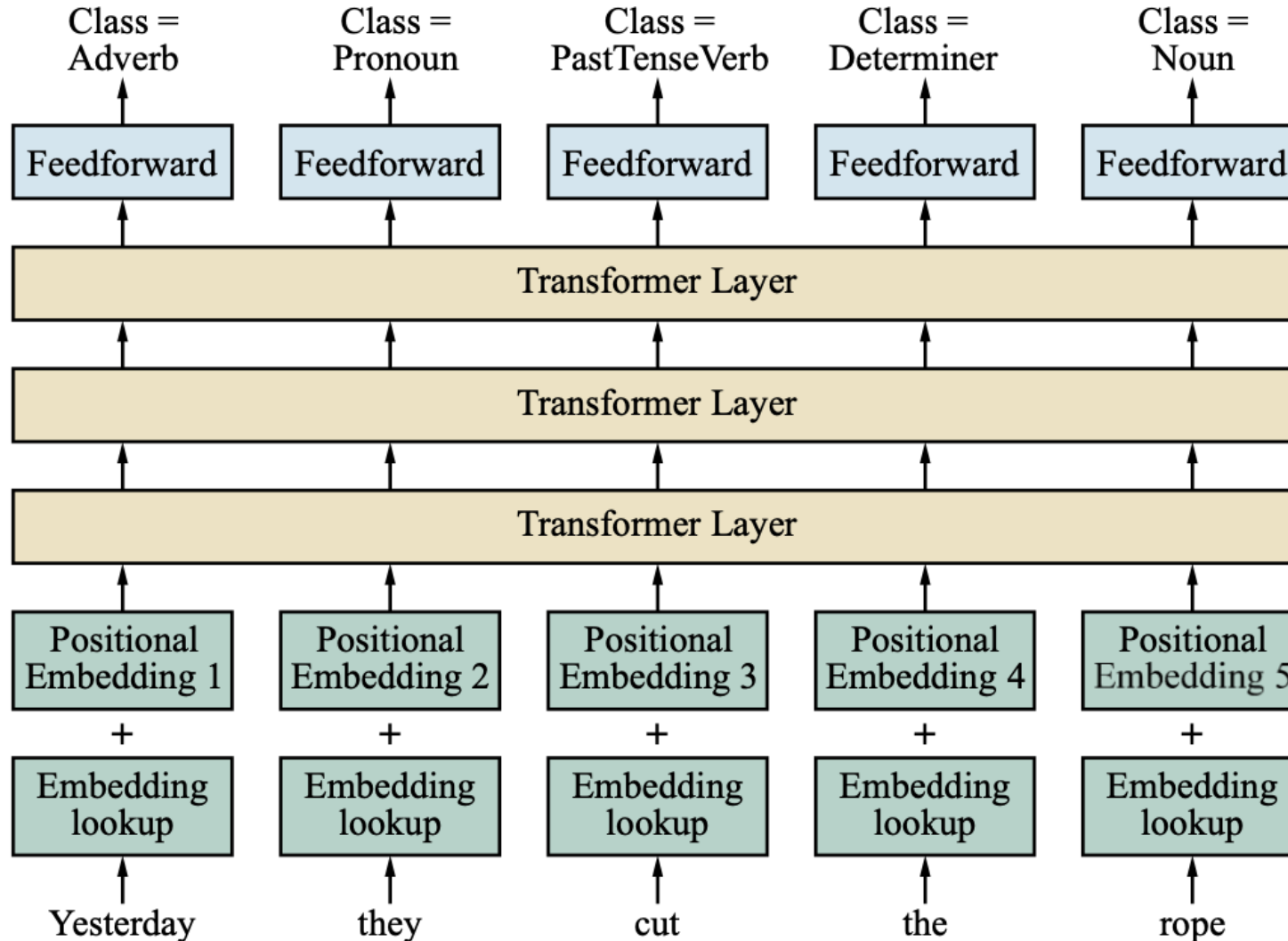
Feature	NLP Attention (Text)	CV Attention (Images)
Input Format	Tokenized text sequences	Image patches
Basic Unit	<b>Words/subwords</b>	<b>Pixels/patches</b>
Role of Attention	Captures long-range dependencies	Learns spatial & contextual relationships
Sequential Processing?	No, operates on full input (parallelizable)	No, processes patches like text
Model Examples	BERT, GPT, T5, RAG	ViT, CLIP, Midjourney

# Single-layer Transformer

consists of self-attention,  
a feedforward network, and residual connection

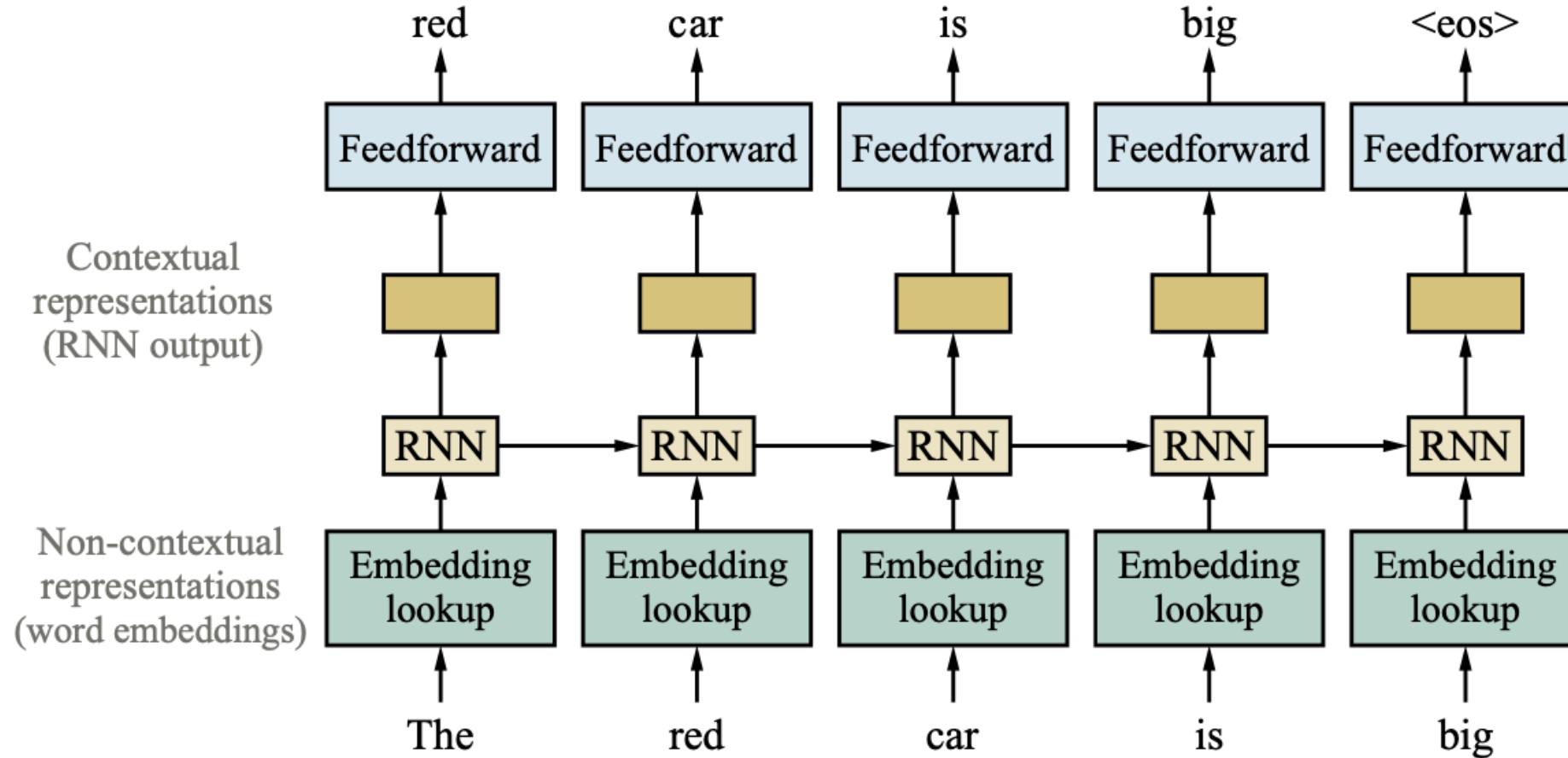


# Transformer Architecture for POS Tagging

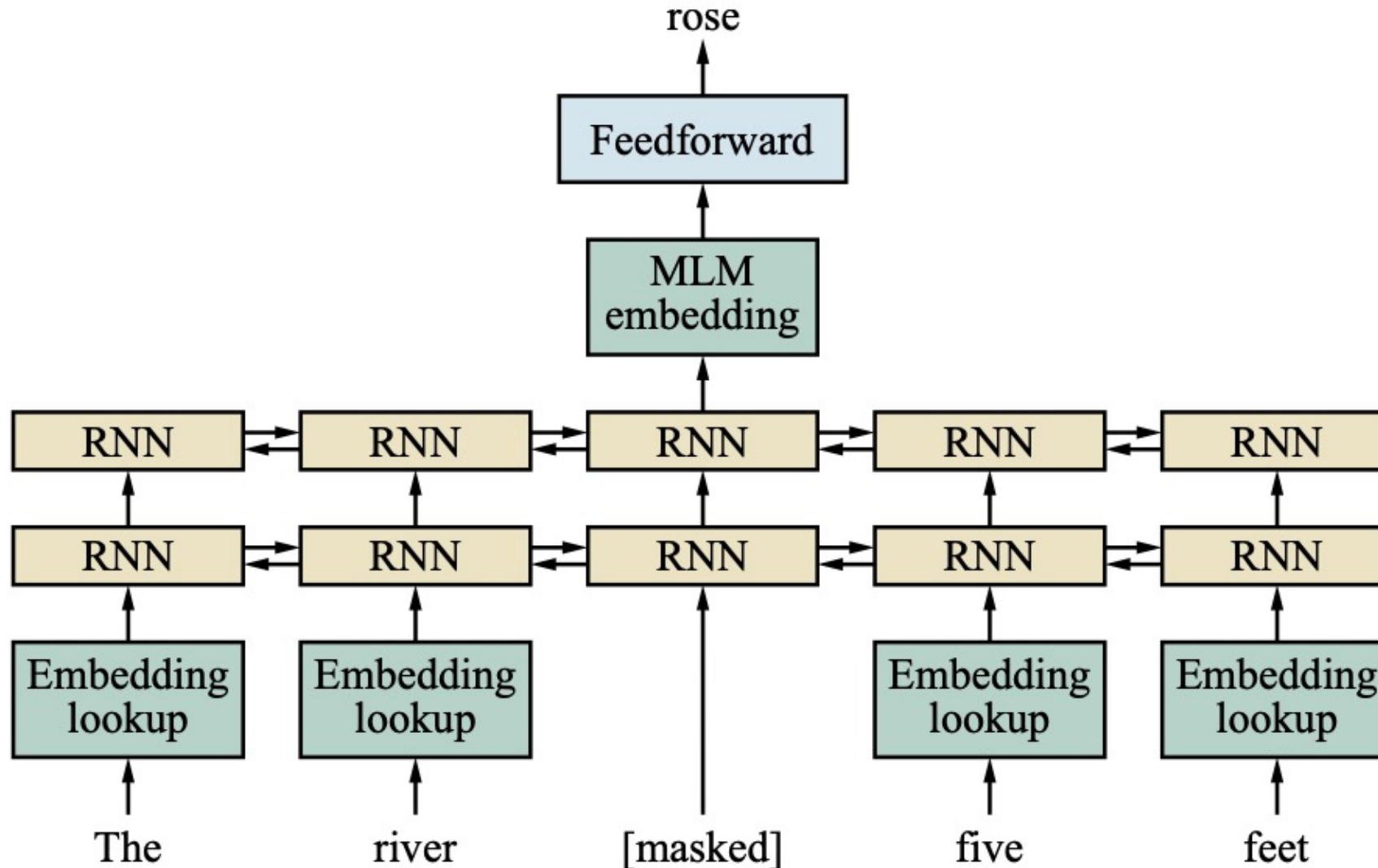


# Training Contextual Representations

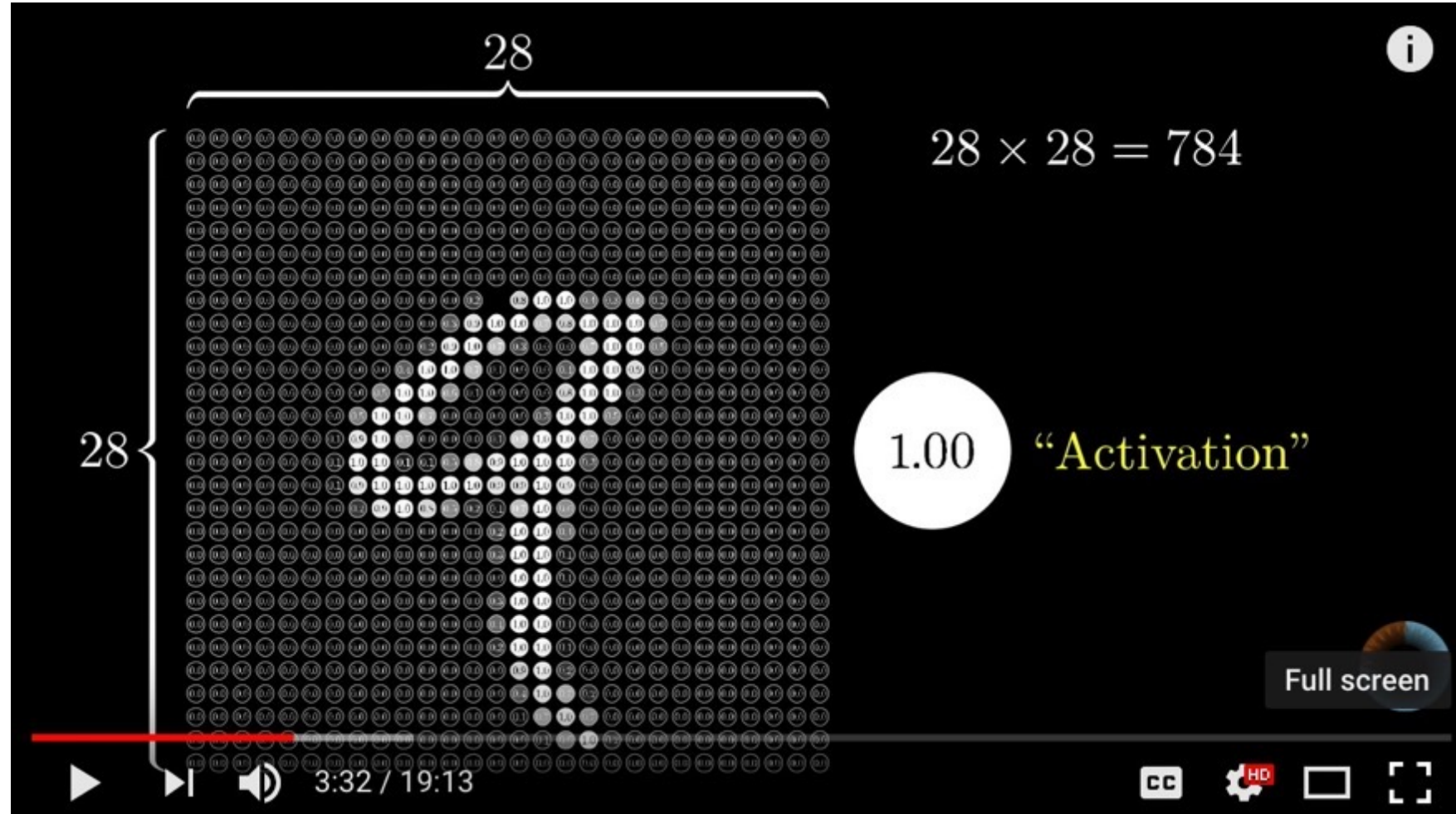
## using a left-to-right Language Model



# Masked Language Modeling: Pretrain a Bidirectional Model



# Neural Network and Deep Learning




Source: 3Blue1Brown (2017), But what \*is\* a Neural Network? | Chapter 1, deep learning,  
<https://www.youtube.com/watch?v=aircAruvnKk>

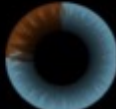
# Gradient Descent

## how neural networks learn

Average cost of all training data...

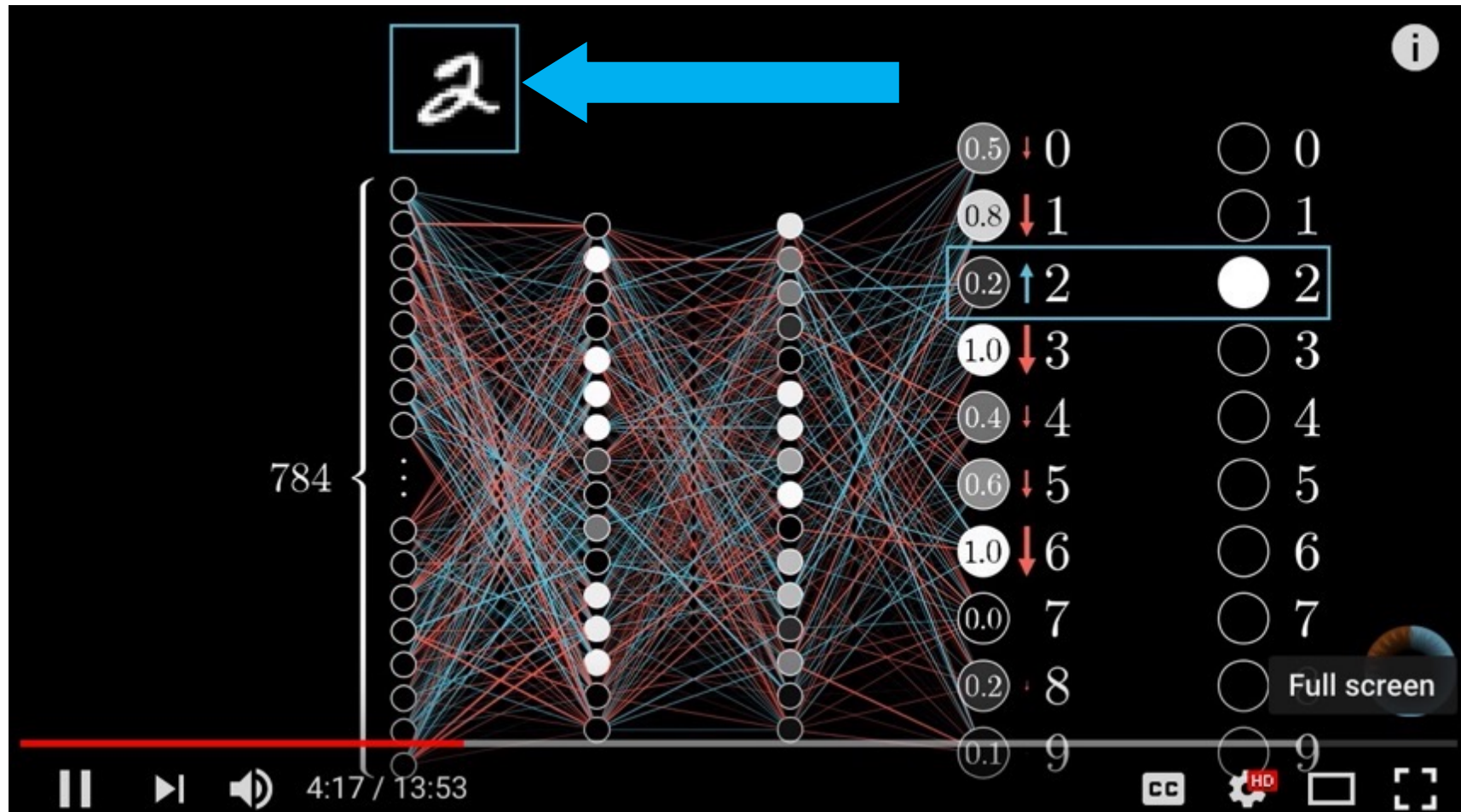
Cost of 

What's the "cost" of this difference?

Utter trash 

$(0.18 - 0.00)^2 +$	<input type="radio"/> 0	<input type="radio"/> 0
$(0.29 - 0.00)^2 +$	<input type="radio"/> 1	<input type="radio"/> 1
$(0.58 - 0.00)^2 +$	<input type="radio"/> 2	<input type="radio"/> 2
$(0.77 - 0.00)^2 +$	<input type="radio"/> 3	<input type="radio"/> 3
$(0.20 - 0.00)^2 +$	<input type="radio"/> 4	<input type="radio"/> 4
$(0.36 - 0.00)^2 +$	<input type="radio"/> 5	<input type="radio"/> 5
$(0.93 - 0.00)^2 +$	<input type="radio"/> 6	<input type="radio"/> 6
$(1.00 - 0.00)^2 +$	<input type="radio"/> 7	<input type="radio"/> 7
$(1.00 - 0.00)^2 +$	<input checked="" type="radio"/> 8	<input checked="" type="radio"/> 8
$(0.95 - 1.00)^2 +$	<input type="radio"/> 9	<input type="radio"/> 9
$(0.35 - 0.00)^2$		

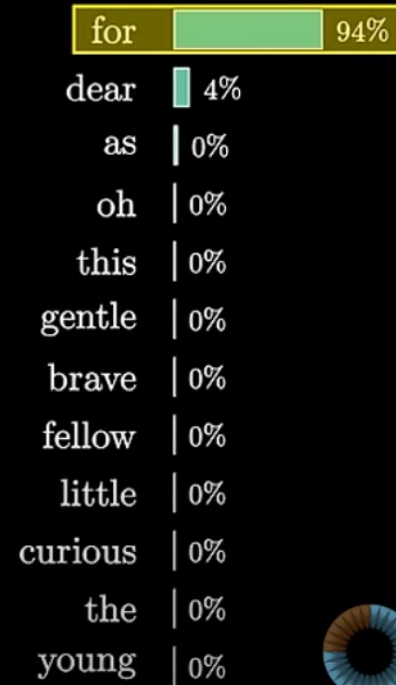
# Backpropagation



Source: 3Blue1Brown (2017), What is backpropagation really doing? | Chapter 3, deep learning, <https://www.youtube.com/watch?v=llg3gGewQ5U>

# Transformers (how LLMs work)

Behold, a wild pi creature, foraging in its native habitat of mathematical formulas and computer code! With its infinite digits and irrational tendencies, this strange creature is beloved by mathematicians and tech enthusiasts alike. Approach with caution, for attempting to calculate its exact value may lead to madness! But do not be afraid, **for**



# Attention in Transformers

**Query**  
1,572,864

$$\begin{bmatrix} -3.7 & +3.9 & -2.4 & -6.3 & -9.4 & -8.6 & +3.6 & -0.9 & \dots & +0.7 \\ +7.9 & +9.7 & -5.6 & +3.2 & -4.7 & -9.5 & +5.1 & -3.6 & \dots & -2.3 \\ +1.7 & +6.6 & +2.6 & +7.4 & -4.5 & +5.9 & -6.2 & +9.0 & \dots & +3.7 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ -5.6 & +8.9 & +4.6 & -4.9 & -5.7 & +0.4 & -9.4 & -5.8 & \dots & -1.5 \end{bmatrix}$$

**Key**  
1,572,864

$$\begin{bmatrix} -2.5 & -0.7 & -4.4 & +1.7 & +7.2 & -7.6 & +0.3 & -7.3 & \dots & +4.3 \\ -2.1 & +1.3 & -6.3 & -7.0 & -0.2 & -2.9 & +8.7 & +5.3 & \dots & +4.9 \\ +8.0 & -8.2 & +1.0 & +1.7 & +9.1 & -4.1 & -5.1 & -7.9 & \dots & -9.6 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ +8.5 & +3.4 & +5.6 & -4.3 & +1.7 & -8.6 & -0.3 & +9.5 & \dots & +7.5 \end{bmatrix}$$


**Value**  
 $12,288 \times 12,288 = 150,994,944$

12,288

$$\begin{bmatrix} -3.2 & +9.1 & -5.3 & +8.9 & +8.7 & +5.9 & +2.6 & +7.4 & \dots & -4.1 \\ +6.9 & +2.3 & -9.6 & -3.0 & -7.0 & +9.5 & -0.4 & -0.1 & \dots & +2.8 \\ -2.6 & -7.2 & +6.4 & -6.1 & +0.2 & -5.5 & -8.0 & +7.2 & \dots & +9.4 \\ +9.1 & +8.0 & +5.4 & -3.3 & -8.3 & -1.8 & -5.3 & -7.3 & \dots & -8.8 \\ +4.5 & -9.7 & +5.4 & -7.0 & -8.3 & -8.1 & +3.4 & -5.0 & \dots & -1.6 \\ +1.1 & +7.1 & +4.5 & -4.5 & -7.3 & -8.8 & -3.9 & -4.7 & \dots & -0.9 \\ +3.6 & +3.9 & -4.3 & -2.4 & -6.3 & +5.7 & -8.8 & +3.9 & \dots & +5.5 \\ +5.5 & -4.8 & -2.5 & +1.7 & -4.5 & -2.6 & -6.0 & -0.8 & \dots & -9.0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ +5.9 & -8.4 & +0.4 & -3.8 & +1.5 & +9.1 & +2.9 & -9.2 & \dots & -1.4 \end{bmatrix} \begin{bmatrix} +0.2 \\ +0.7 \\ +3.6 \\ -4.4 \\ -7.3 \\ -2.1 \\ +9.0 \\ -6.2 \\ \vdots \\ +0.9 \end{bmatrix} = \begin{bmatrix} -198.6 \\ +73.1 \\ -28.2 \\ +119.4 \\ -4.4 \\ +215.7 \\ +91.8 \\ -29.1 \\ -5.6 \\ \vdots \\ -5.1 \end{bmatrix}$$

16:48 / 26:09 • Counting parameters >

# How might LLMs store facts

 **GPT-3**

Total weights: **175,181,291,520** ⓘ

Embedding	$12,288 \times 50,257$ $d\_embed * n\_vocab$	$= 617,558,016$
Key	$128 \times 12,288 \times 96 \times 96$ $d\_query * d\_embed * n\_heads * n\_layers$	$= 14,495,514,624$
Query	$128 \times 12,288 \times 96 \times 96$ $d\_query * d\_embed * n\_heads * n\_layers$	$= 14,495,514,624$
Value	$128 \times 12,288 \times 96 \times 96$ $d\_value * d\_embed * n\_heads * n\_layers$	$= 14,495,514,624$
Output	$12,288 \times 128 \times 96 \times 96$ $d\_embed * d\_value * n\_heads * n\_layers$	$= 14,495,514,624$
Up-projection	$49,152 \times 12,288 \times 96$ $n\_neurons * d\_embed * n\_layers$	$= 57,982,058,496$
Down-projection	$12,288 \times 49,152 \times 96$ $d\_embed * n\_neurons * n\_layers$	$= 57,982,058,496$
<del>Unembedding</del>	<del><math>50,257 \times 12,288</math> <del><math>n\_vocab * d\_embed</math></del></del>	<del><math>= 617,558,016</math></del>

16:51 / 22:42 • Counting parameters >

# Large Language Models explained briefly

What follows is a conversation between a user and a helpful, very knowledgeable AI assistant.

User: Give me some ideas for what to do when visiting Santiago.

AI Assistant: Sure,

\_\_\_\_\_

Large Language Model

Token	Probability
,	53%
!	38%
thing	7%
.	0%
!	0%
-	0%
!	0%
!	0%
I	0%
thing	0%
-	0%
,I	0%
...	0%

1:49 / 8:47 · What are large language models? >

Source: 3Blue1Brown (2024), Large Language Models explained briefly, <https://www.youtube.com/watch?v=LPZh9BOjkQs>

**Generative AI,  
Agentic AI,  
AI Agent,  
RAG LLM  
for  
QA and Dialogue Systems**

**Chatbot**

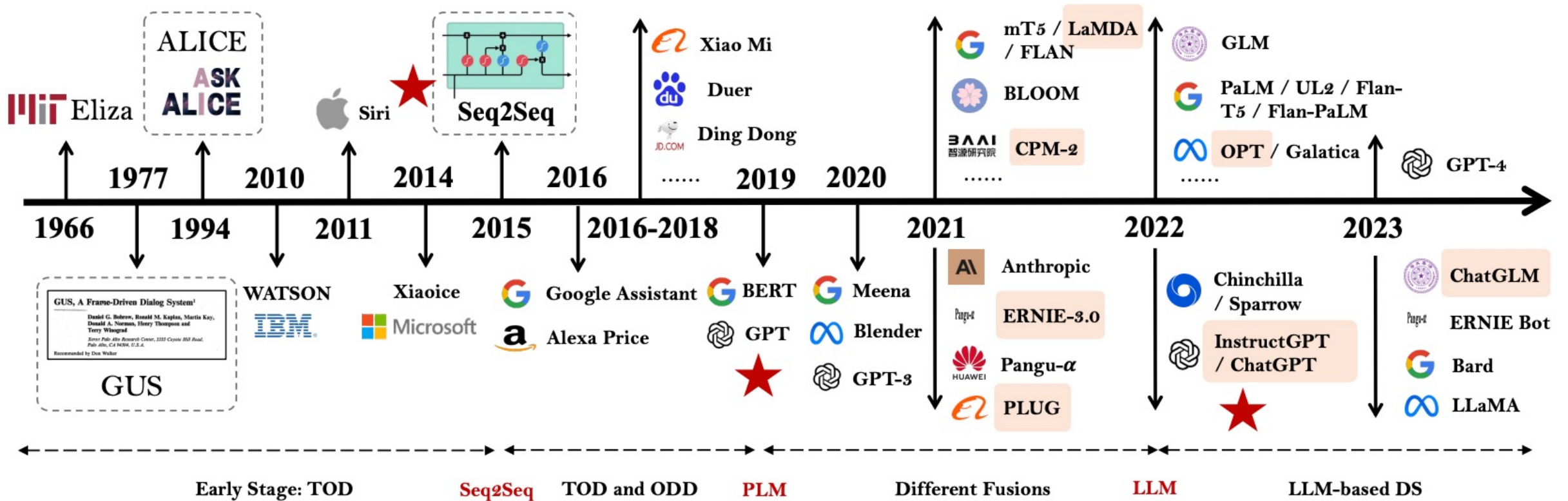
**Dialogue System**

**Intelligent Agent**

**Conversational AI**

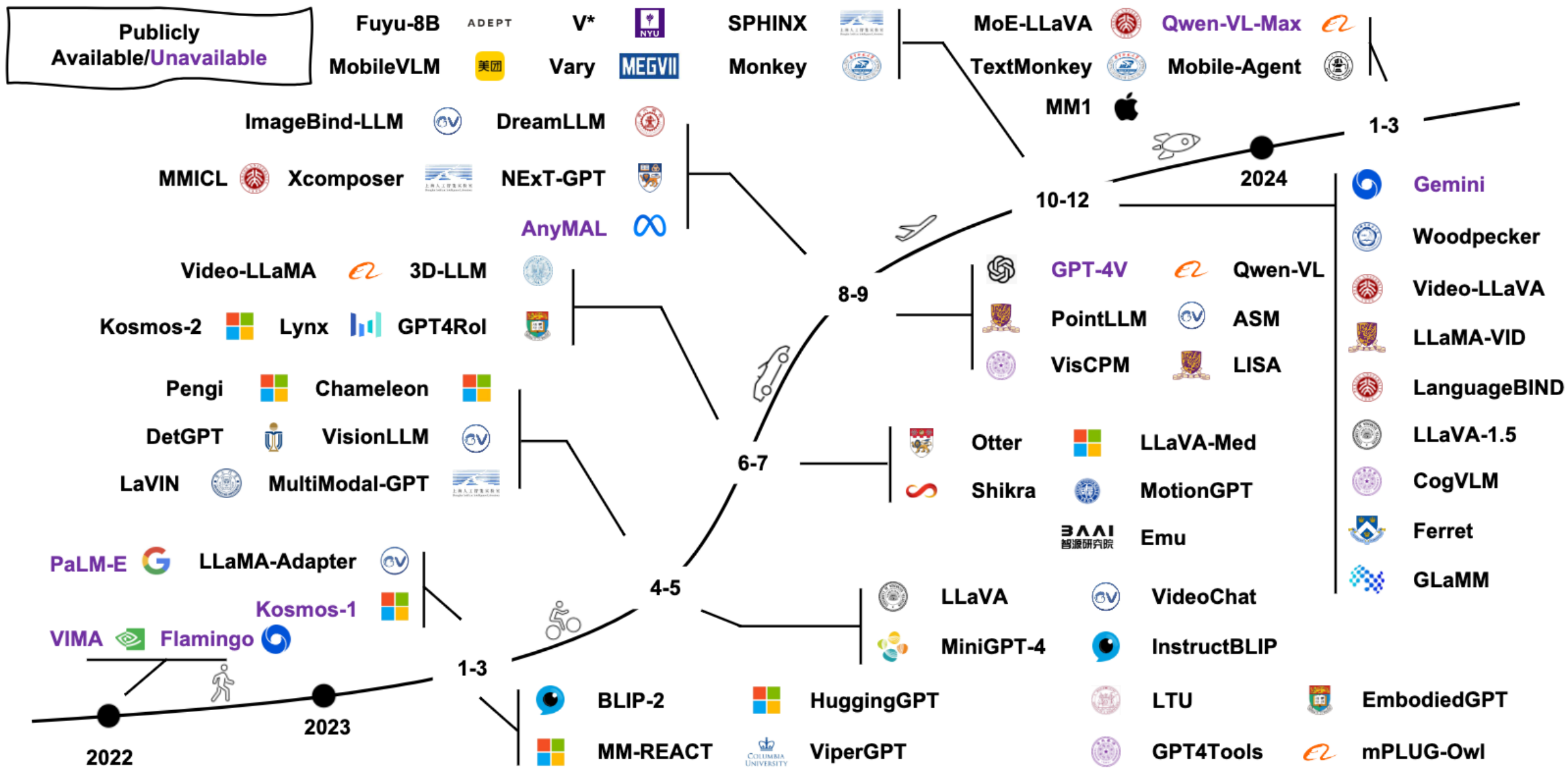
# The Development of LM-based Dialogue Systems

- 1) Early Stage (1966 - 2015)
- 2) The Independent Development of TOD and ODD (2015 - 2019)
- 3) Fusions of Dialogue Systems (2019 - 2022)
- 4) LLM-based DS (2022 - Now)

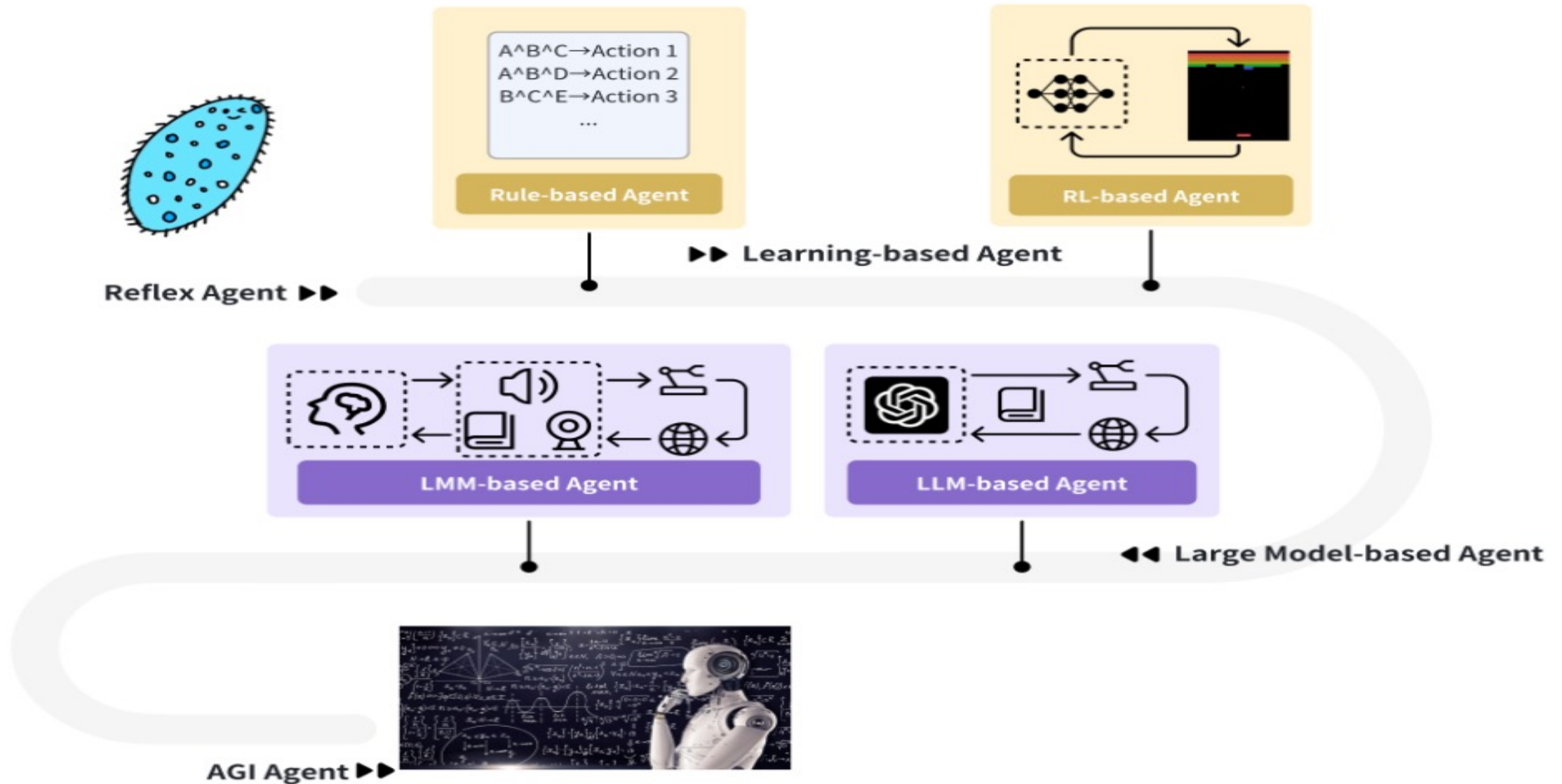


Task-oriented DS (TOD), Open-domain DS (ODD)

# Multimodal Large Language Models (MLLM)



# Intelligent Agents Roadmap



# AI Agents

- **Traditional AI Agents**

- **Simple reflex agents**
- **Model-based reflex agents**
- **Goal-based agents**
- **Utility-based agents**
- **Learning agents**

- **Evolution of AI Agents**

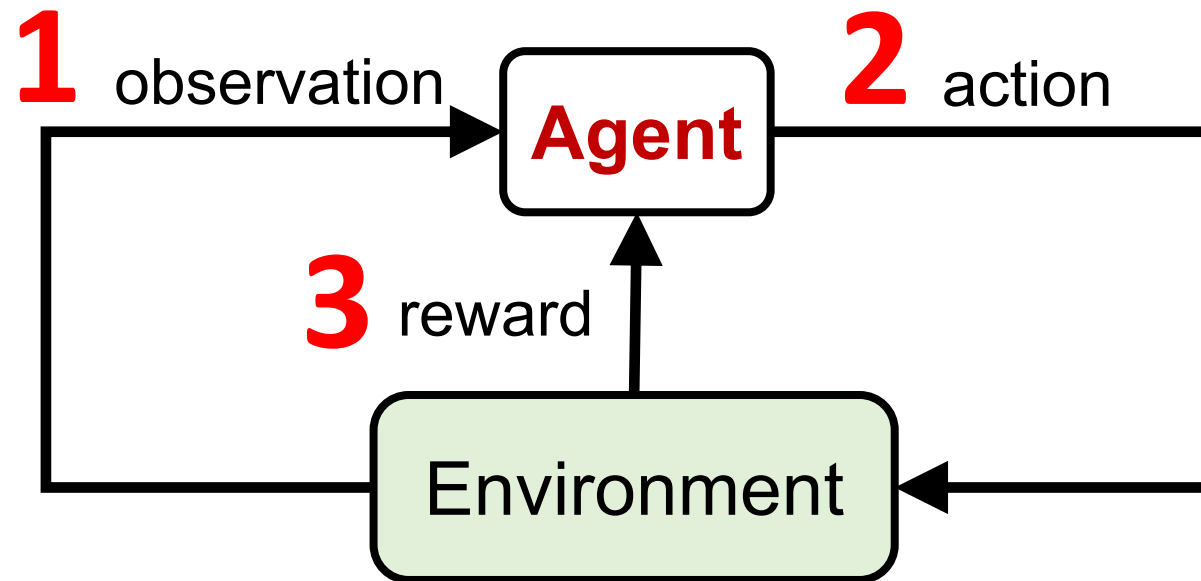
- **LLM-based Agents**
- **Multi-modal agents**
- **Embodied AI agents in virtual environments**
- **Collaborative AI agents**

# Reinforcement Learning (DL)

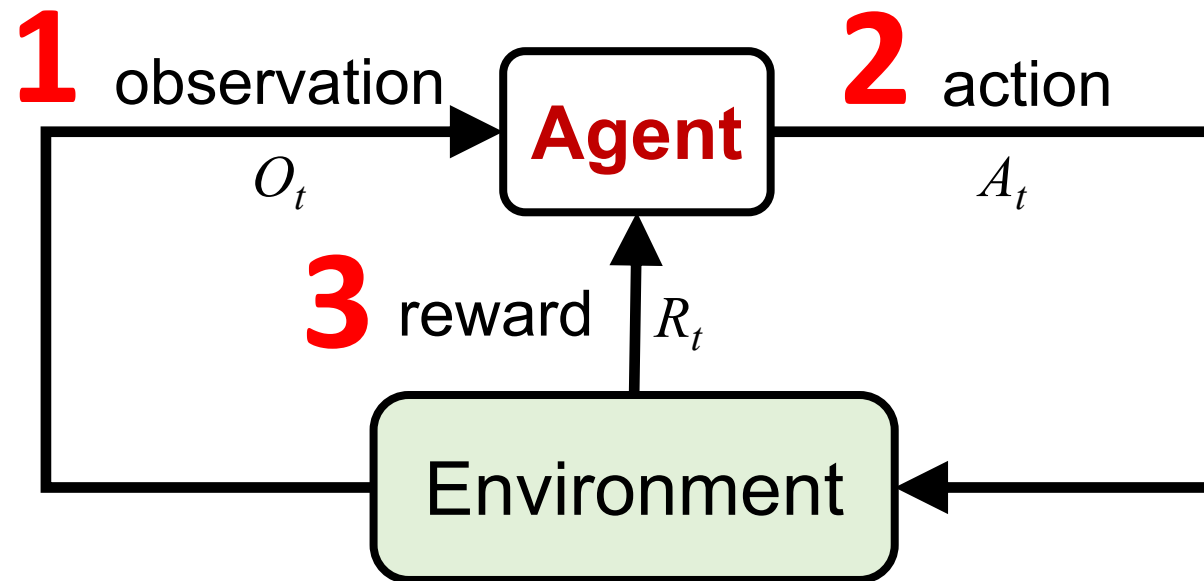
**Agent**

Environment

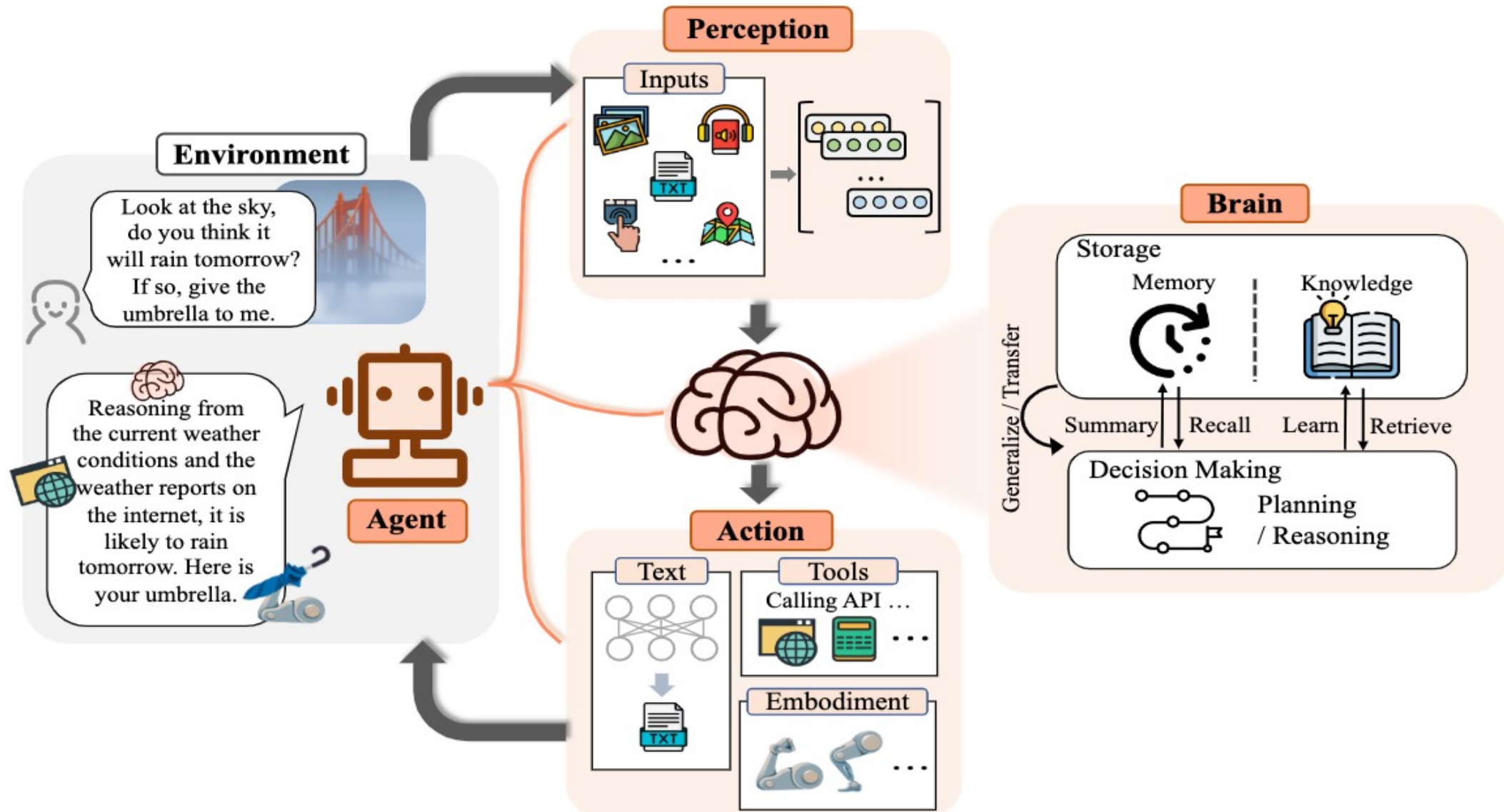
# Reinforcement Learning (DL)



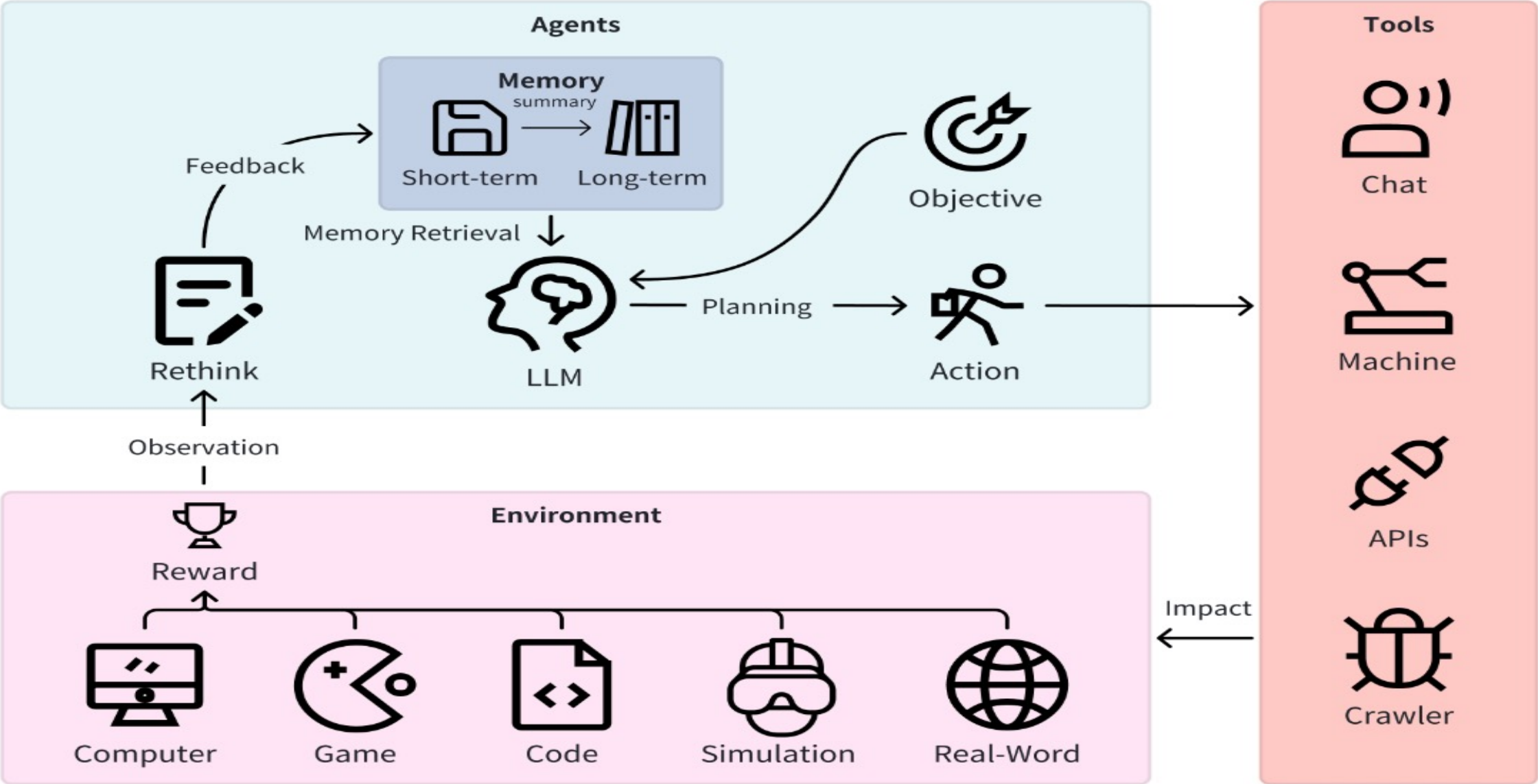
# Reinforcement Learning (DL)



# Large Language Model (LLM) based Agents

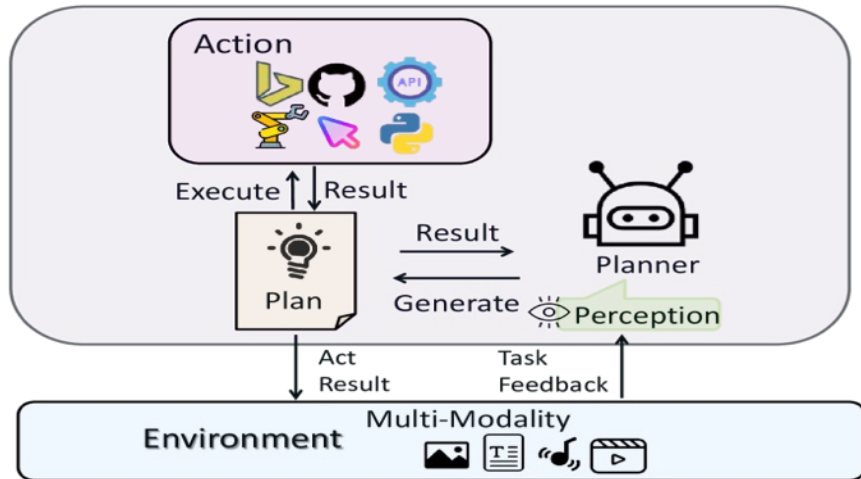


# LLM-based Agents

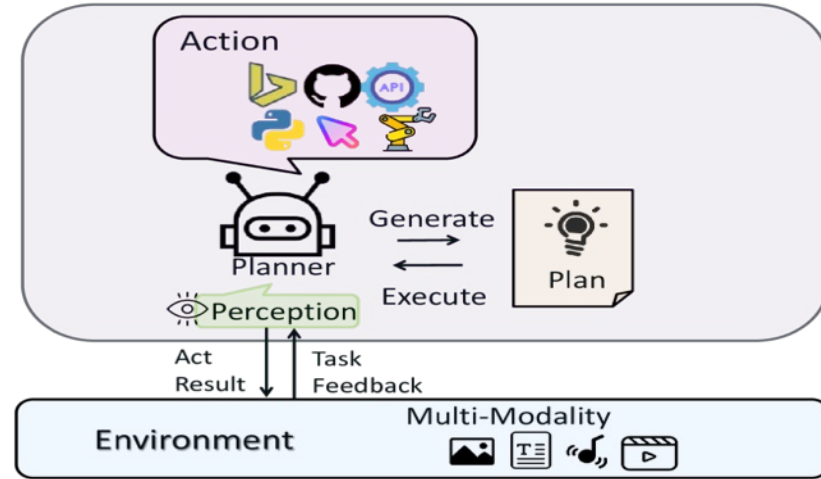


Source: Cheng, Yuheng, Ceyao Zhang, Zhengwen Zhang, Xiangrui Meng, Sirui Hong, Wenhao Li, Zihao Wang et al. "Exploring large language model based intelligent agents: Definitions, methods, and prospects." arXiv preprint arXiv:2401.03428 (2024).

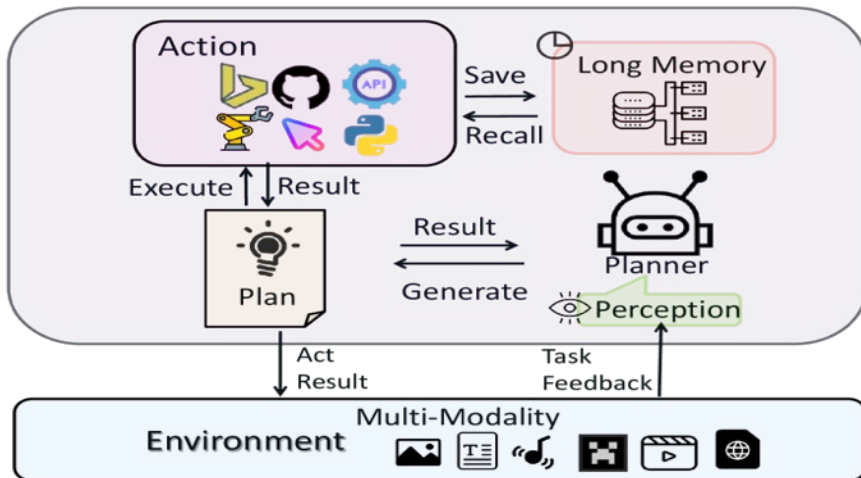
# Large Multimodal Agents (LMA)



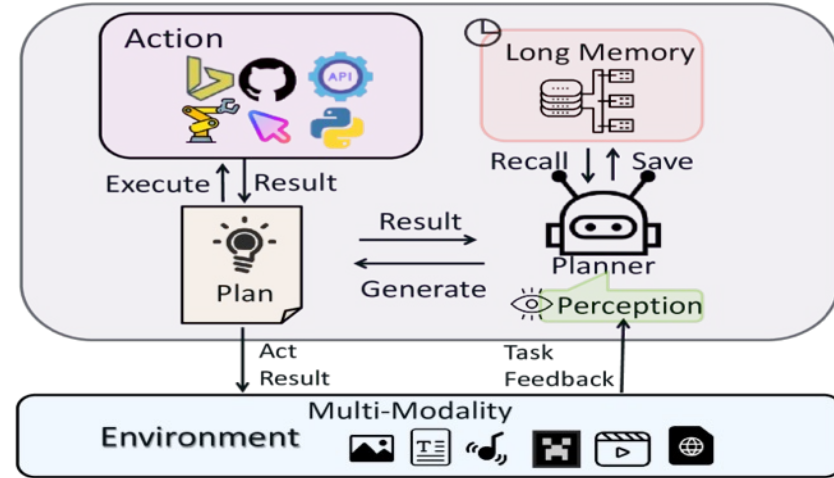
(a)



(b)

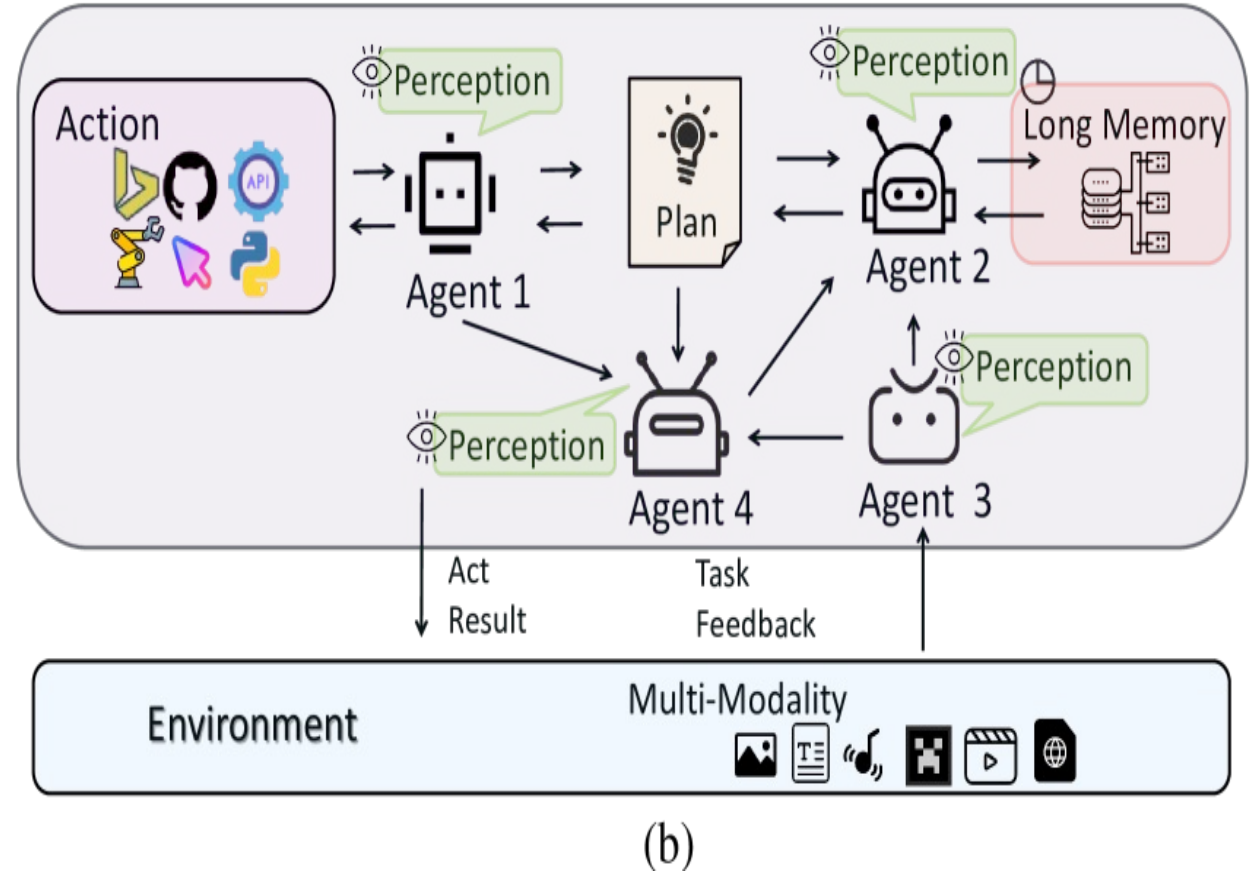
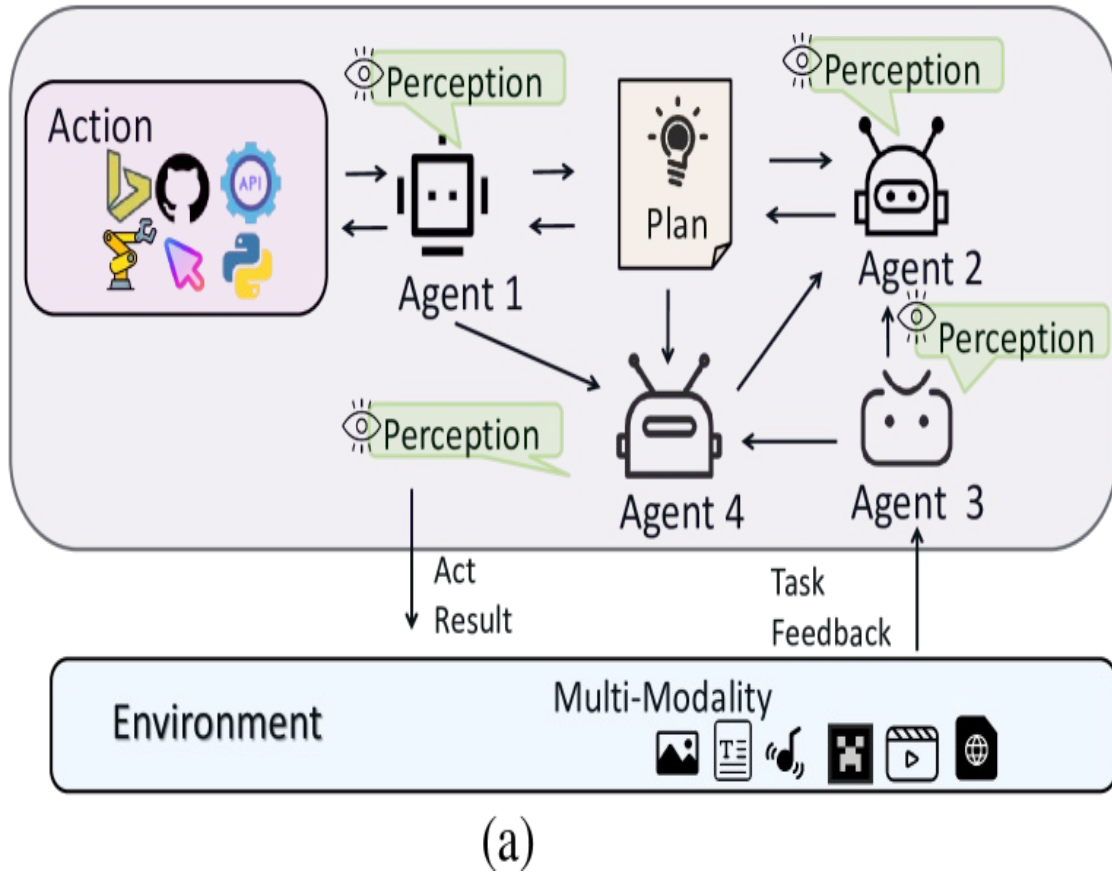


(c)

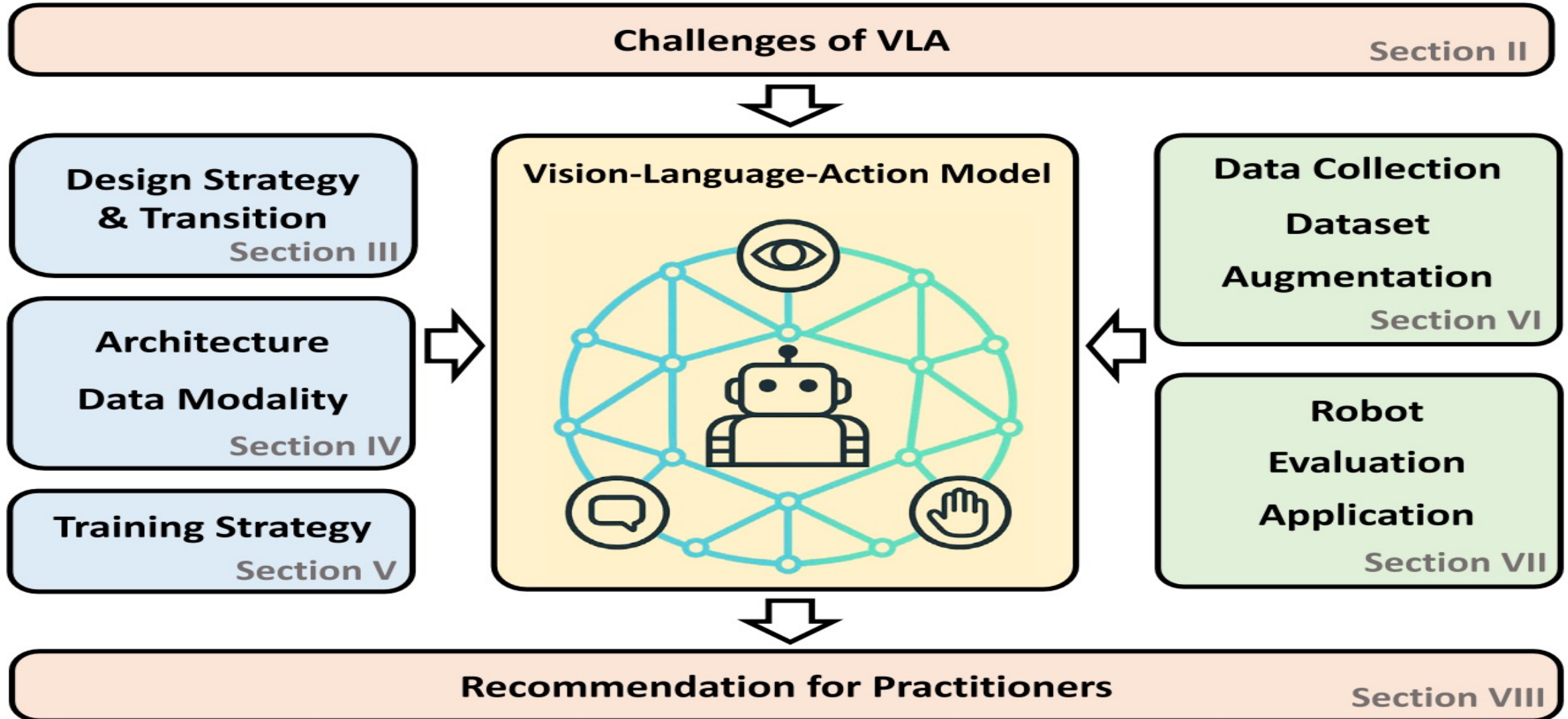


(d)

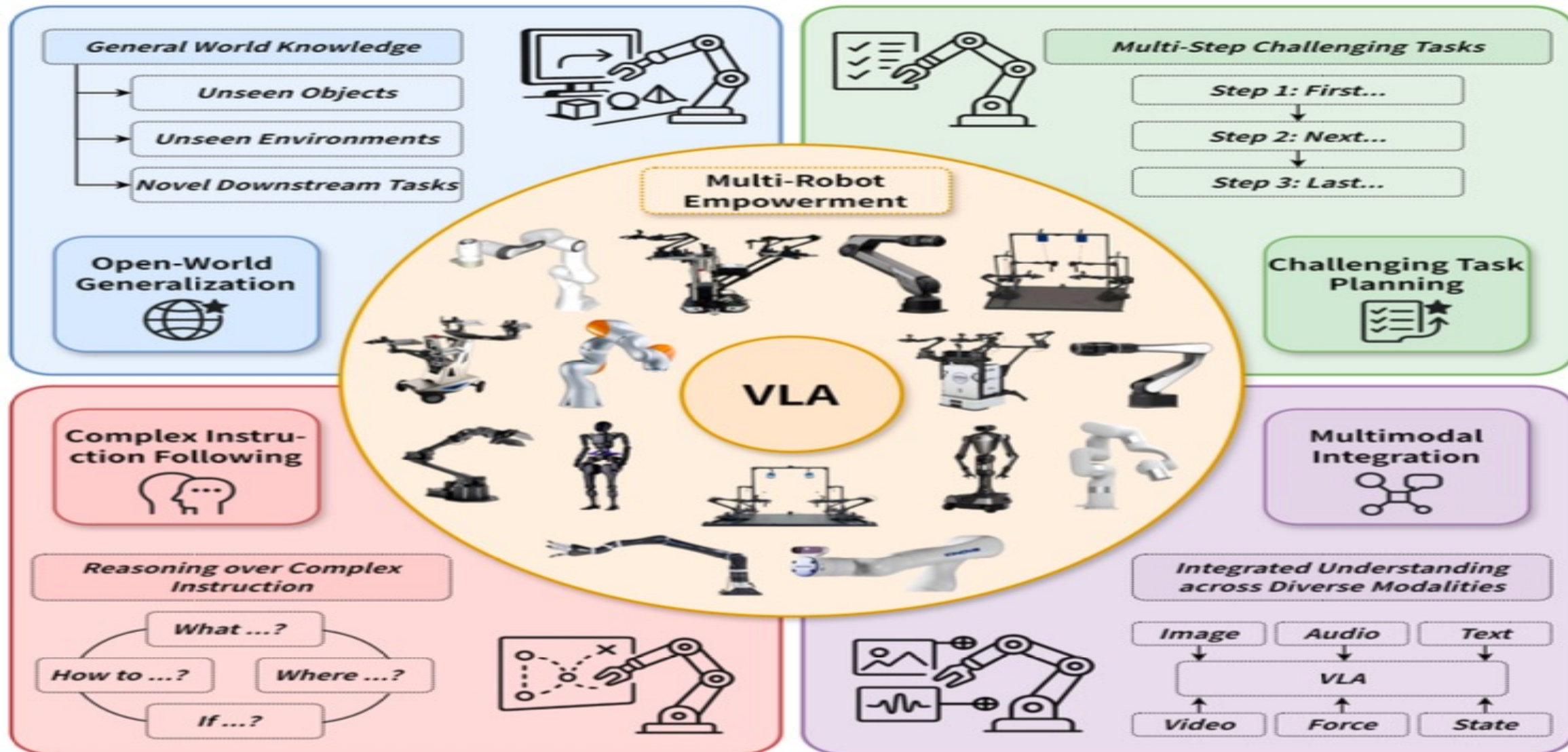
# Large Multimodal Agents (LMA)



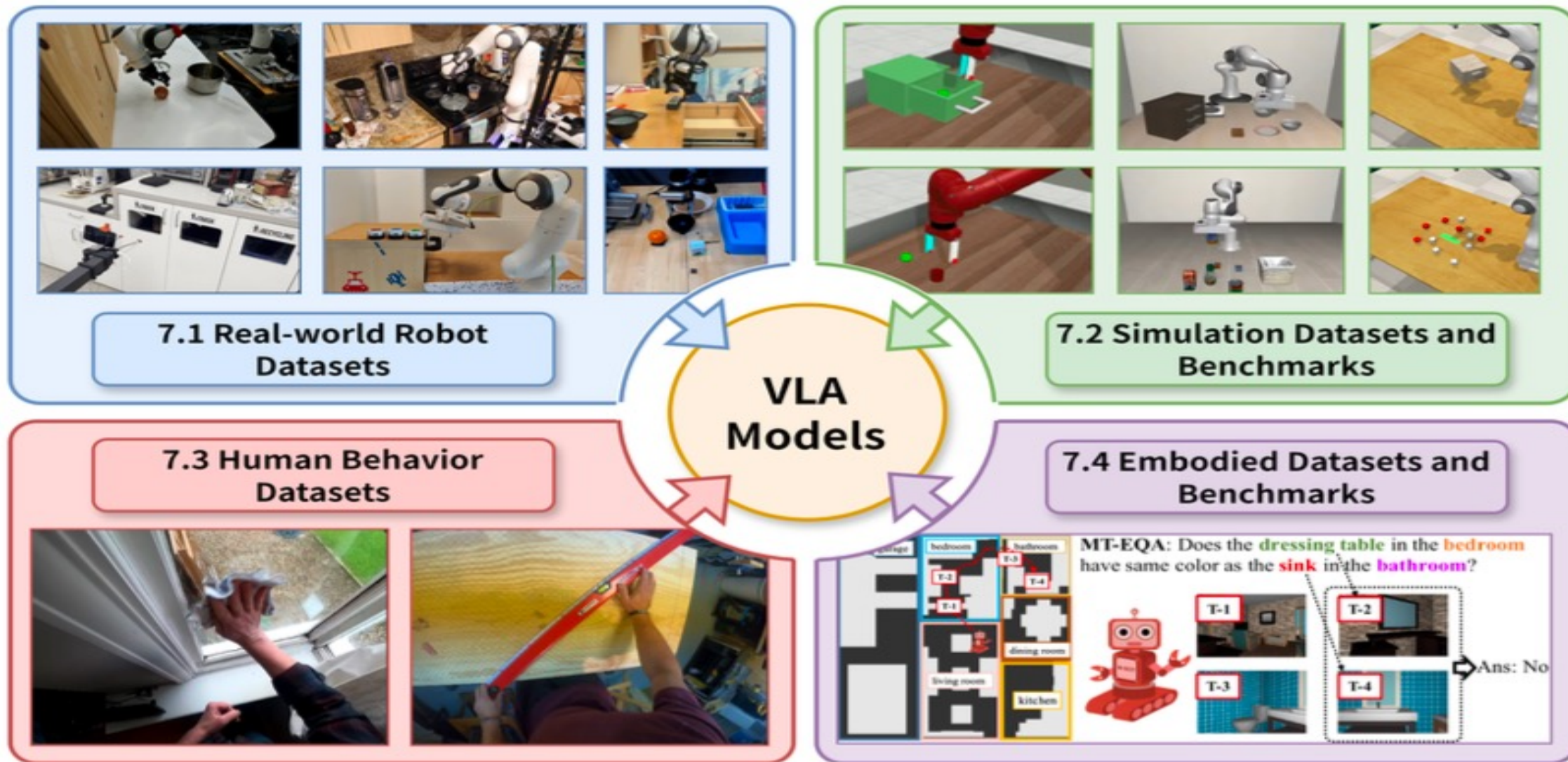
# Vision-Language-Action (VLA) Models for Robotics



# Large VLM-based Vision-Language-Action Models for Robotic Manipulation

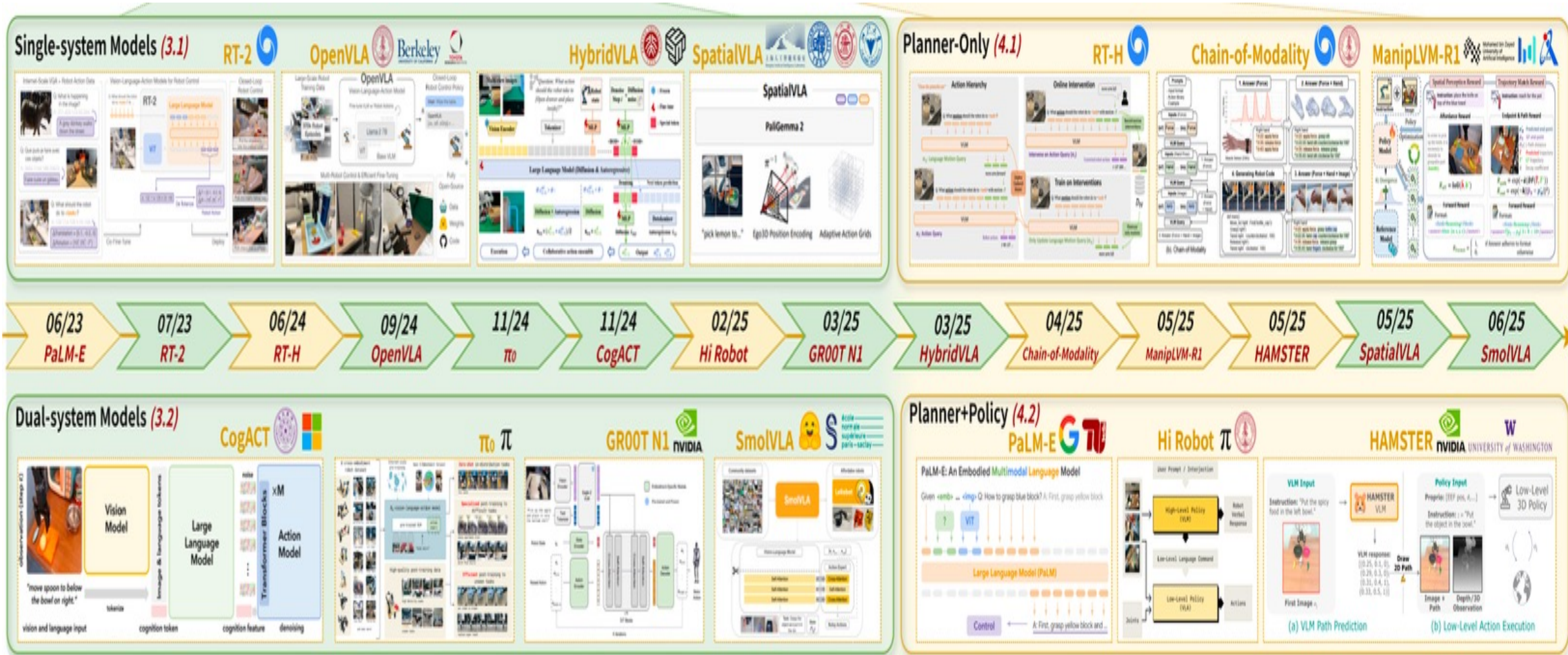


# Large VLM-based Vision-Language-Action Models for Robotic Manipulation



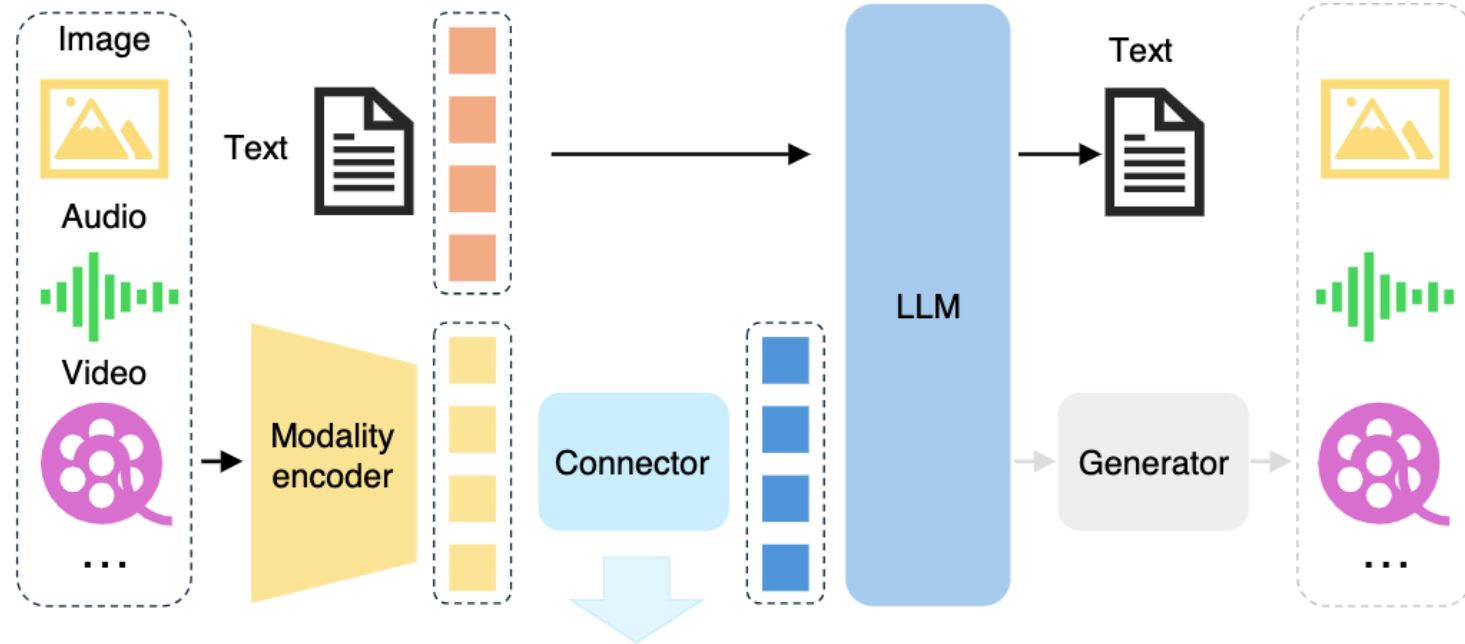
# Large VLM-based Vision-Language-Action Models for Robotic Manipulation (Timeline)

## Monolithic models and Hierarchical Models



Source: Rui Shao, Wei Li, Lingsen Zhang, Renshan Zhang, Zhiyang Liu, Ran Chen, and Liqiang Nie. (2025) "Large vlm-based vision-language-action models for robotic manipulation: A survey." arXiv preprint arXiv:2508.13073 (2025).

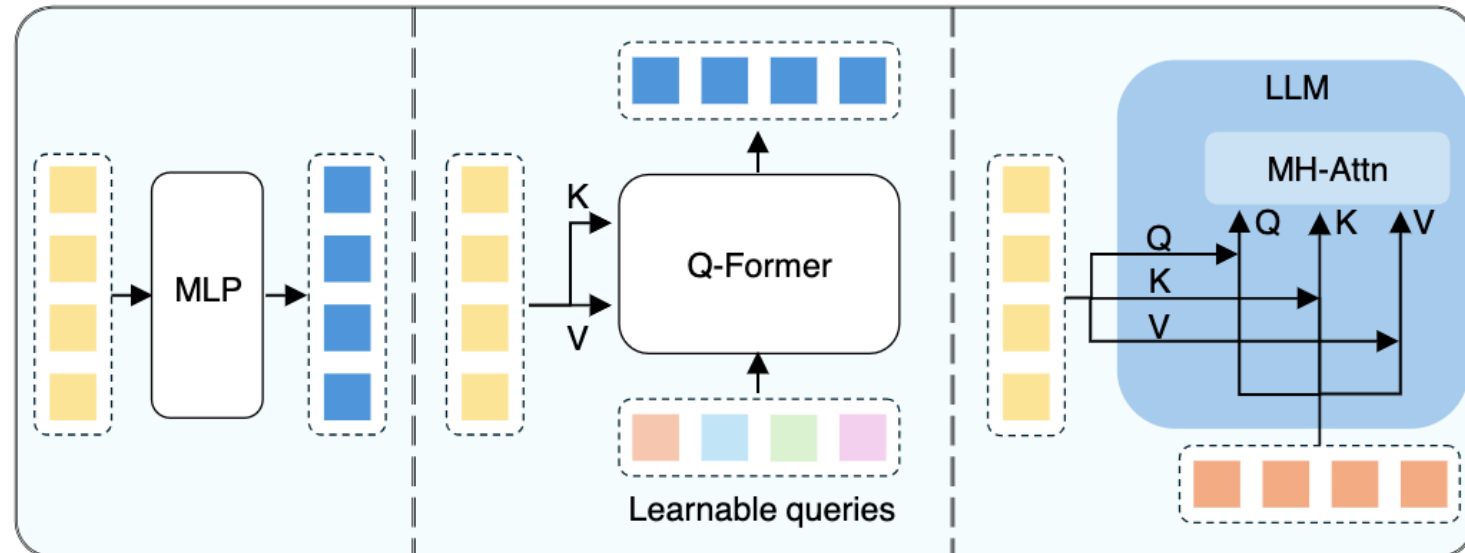
# Multimodal Large Language Models (MLLM)



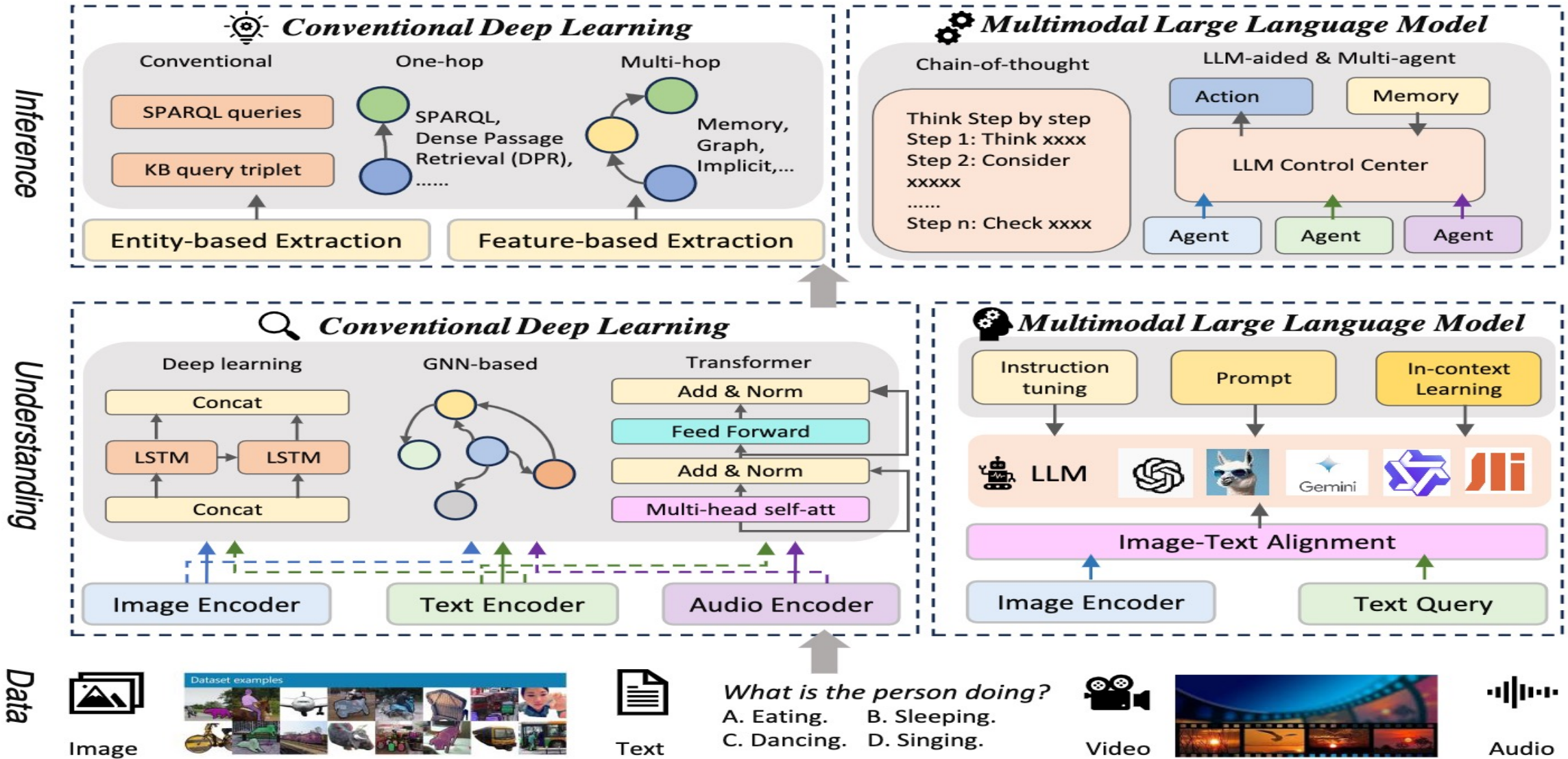
## Multimodal LLM

Three types of connectors:


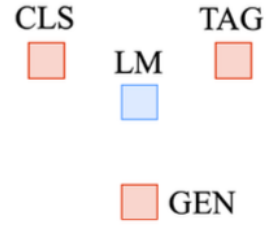
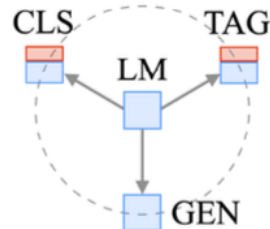
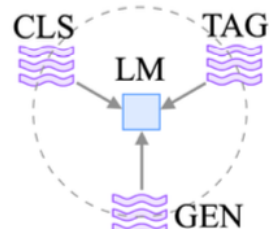
1. projection-based
2. query-based
3. fusion-based connectors



# Multimodal Large Language Model (MLLM) for Vision Question Answering

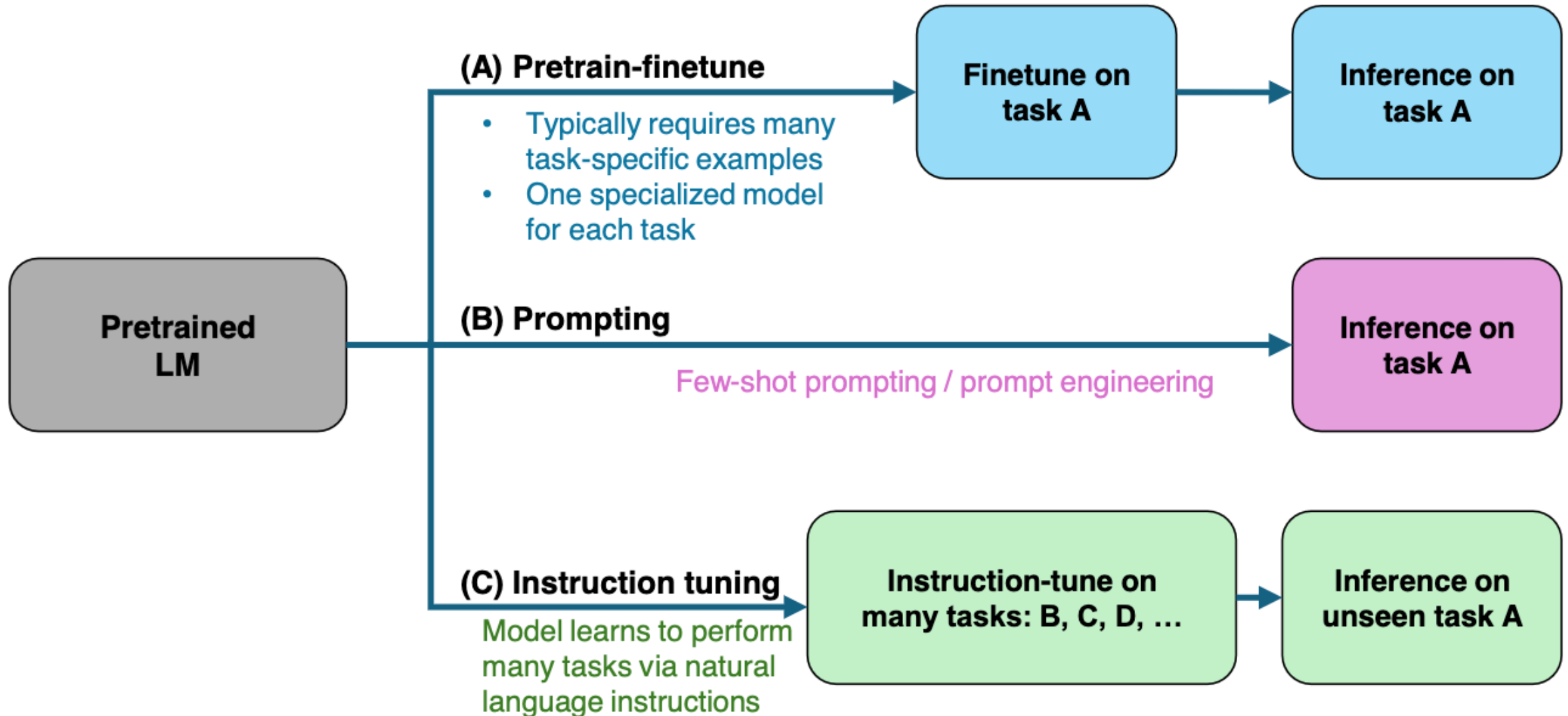


# Four Paradigms in NLP (LM)

Paradigm	Engineering	Task Relation
a. Fully Supervised Learning (Non-Neural Network)	Feature (e.g. word identity, part-of-speech, sentence length)	
b. Fully Supervised Learning (Neural Network)	Architecture (e.g. convolutional, recurrent, self-attentional)	
<b>Transfer Learning: Pre-training, Fine-Tuning (FT)</b>		
c. Pre-train, Fine-tune	Objective (e.g. masked language modeling, next sentence prediction)	
<b>GAI: Pre-train, Prompt, and Predict (Prompting)</b>		
d. Pre-train, Prompt, Predict	Prompt (e.g. cloze, prefix)	

# Large Language Models (LLM)

## Three typical learning paradigms



# Generative AI (Gen AI)

## AI Generated Content (AIGC)

### Image Generation

**Instruction 1:**

*An astronaut riding a horse in a photorealistic style.*

**Instruction 2:**

*Teddy bears working on new AI research on the moon in the 1980s.*

Figure 1



Figure 2

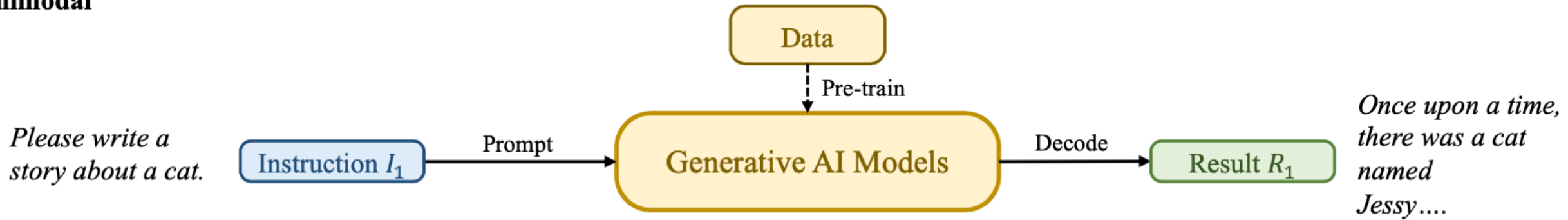


 **OpenAI DALL·E 2**

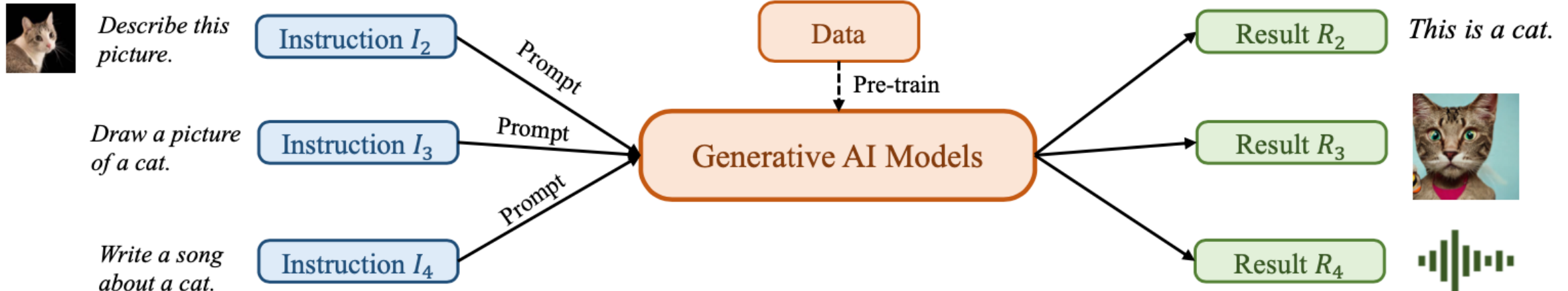
# Generative AI (Gen AI)

## AI Generated Content (AIGC)

### Unimodal

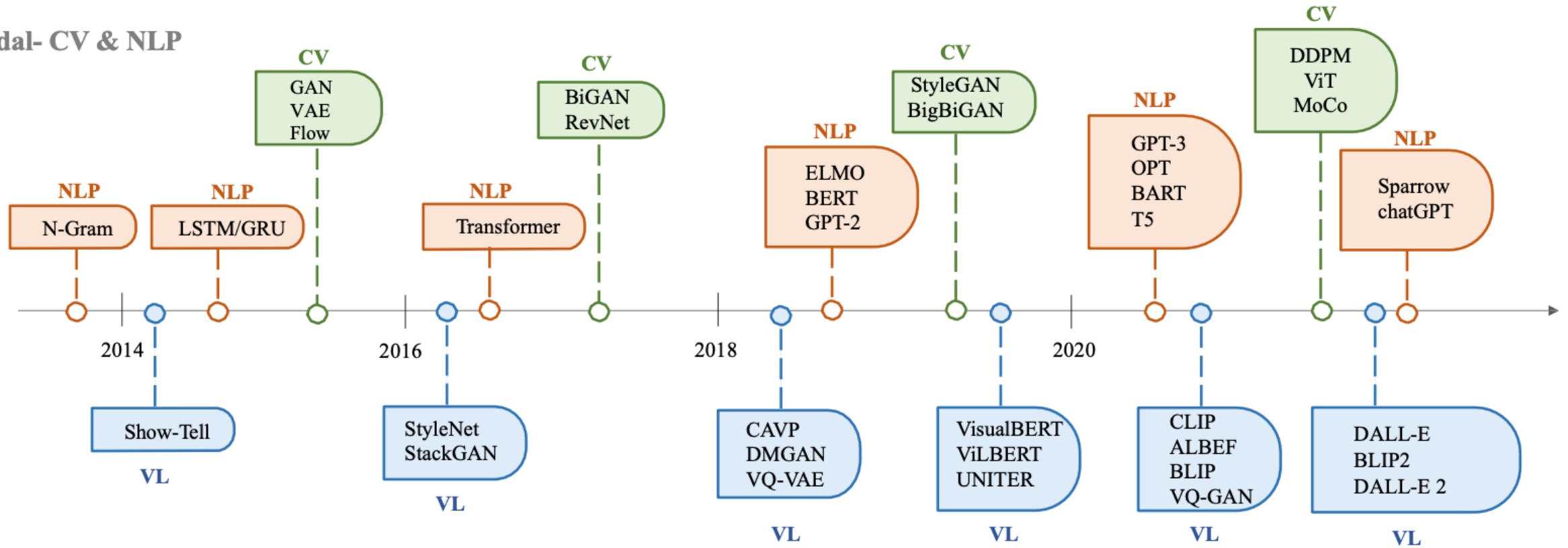


### Multimodal



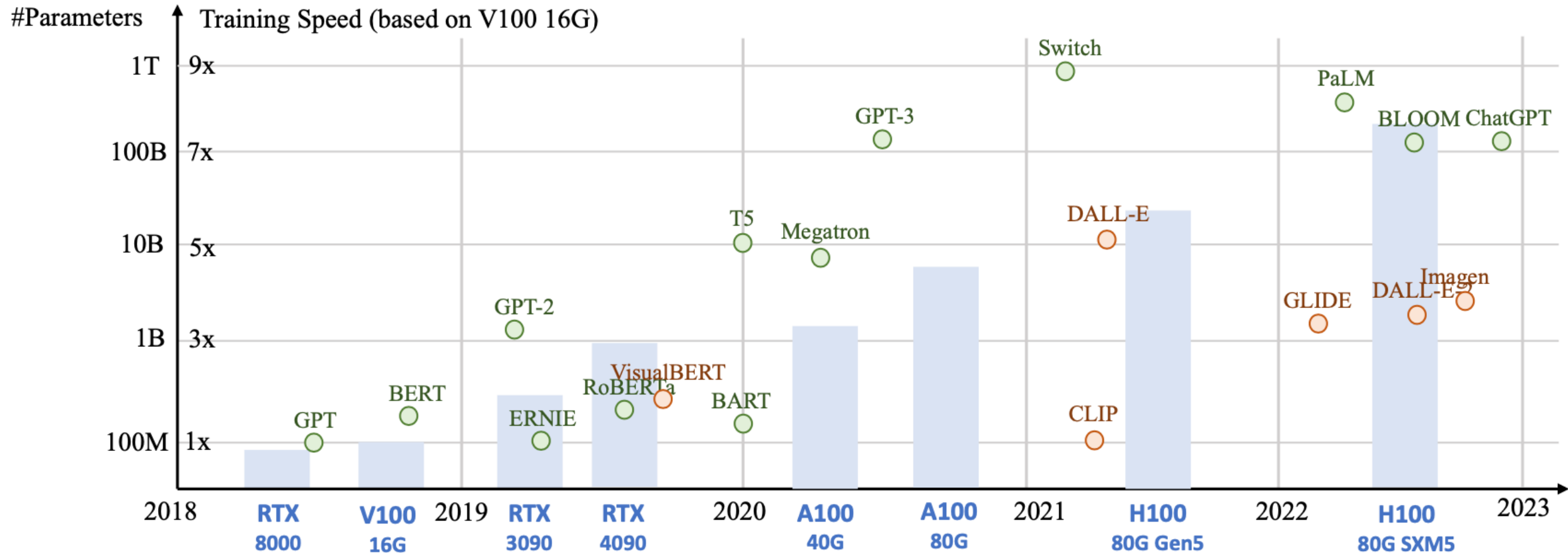
# The history of Generative AI in CV, NLP and VL

## Unimodal- CV & NLP

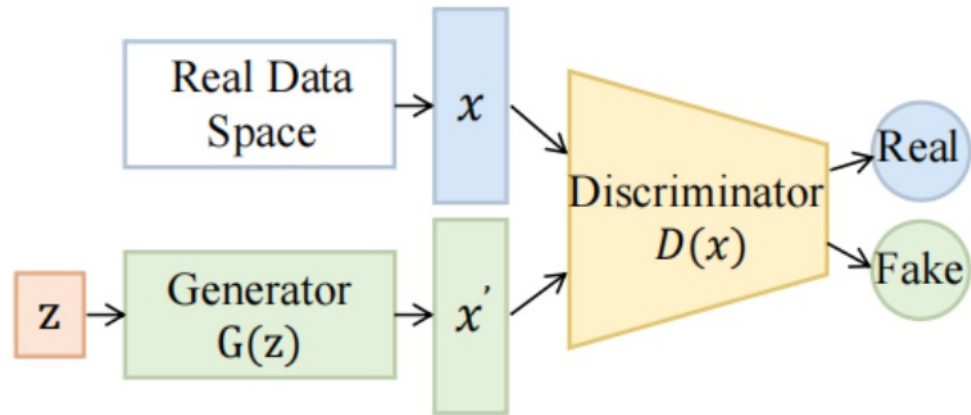


## Multimodal – Vision Language

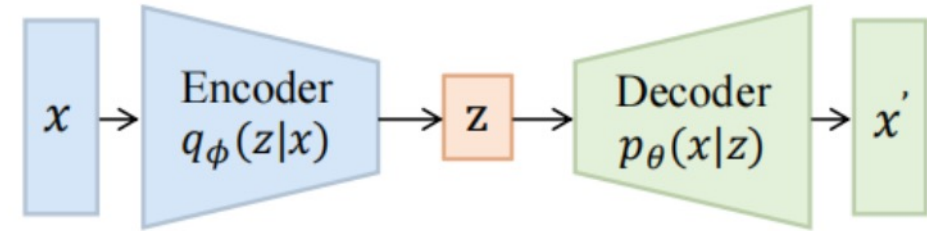
# Generative AI Foundation Models



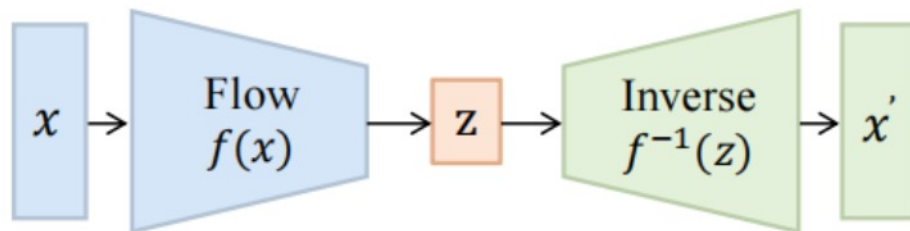
# Categories of Vision Generative Models



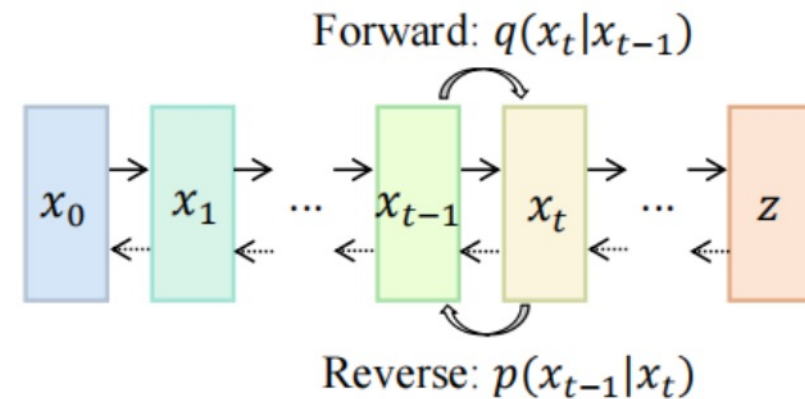
(1) Generative adversarial networks



(2) Variational autoencoders

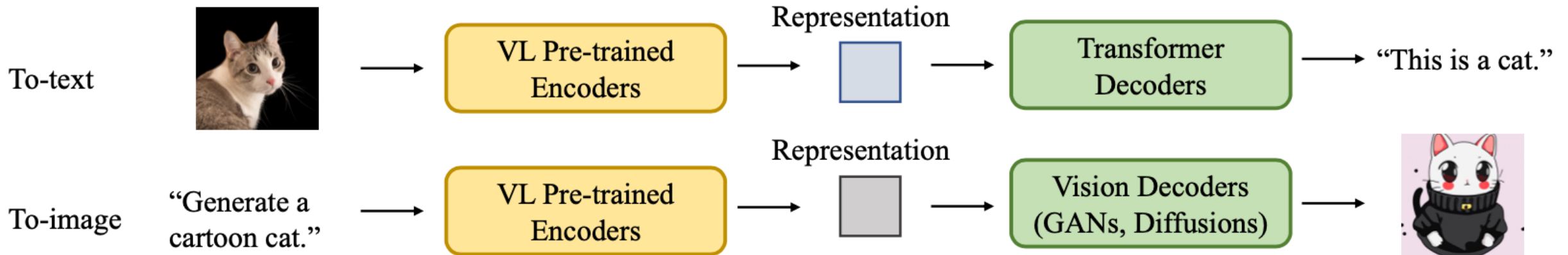
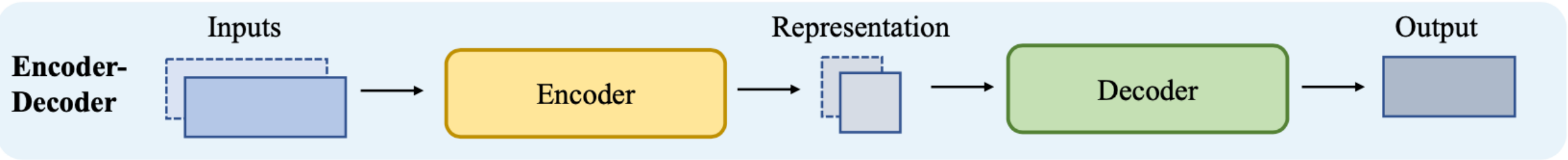


(3) Normalizing flows



(4) Diffusion models

# The General Structure of Generative Vision Language

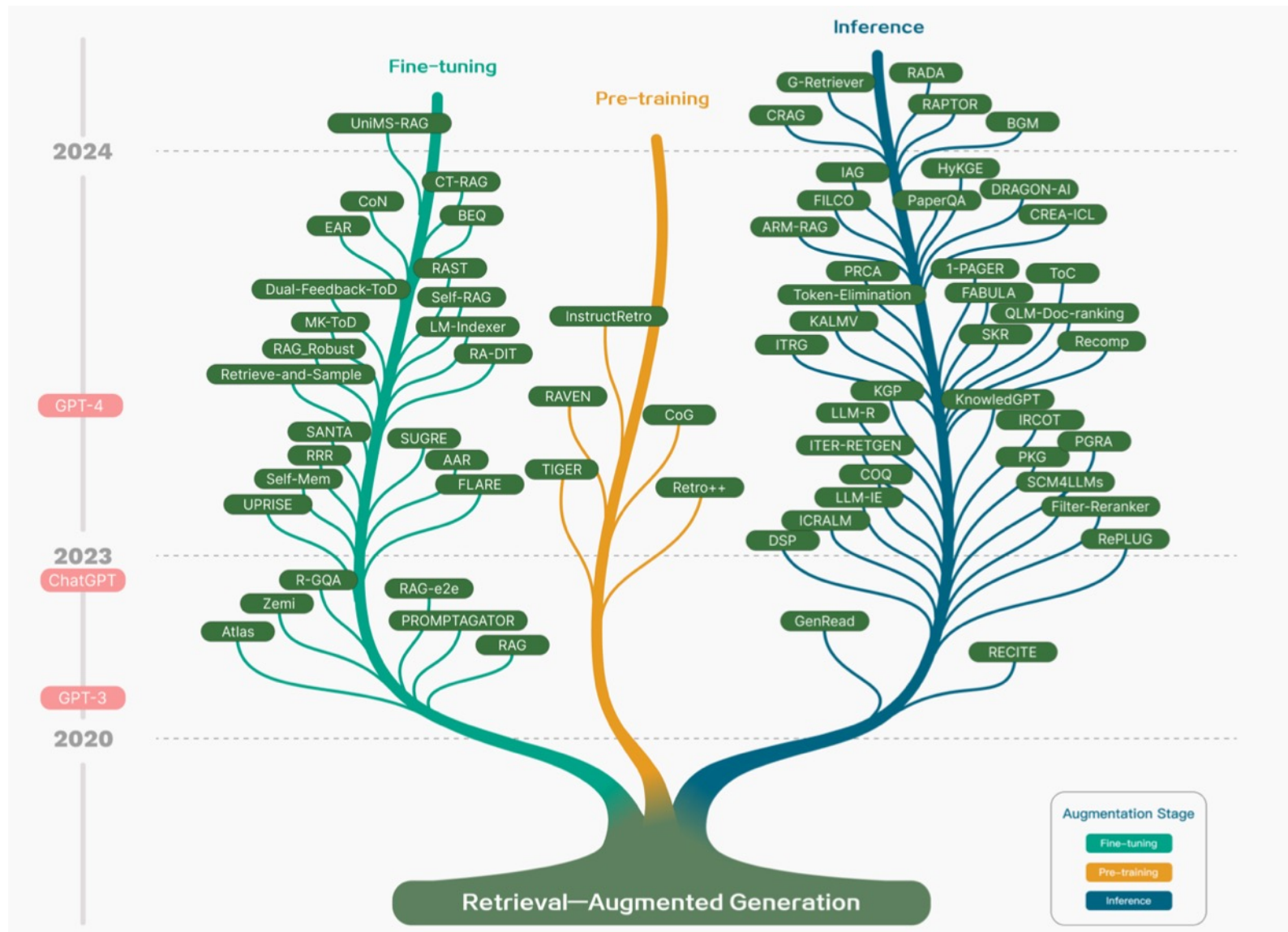


# **RAG LLM**

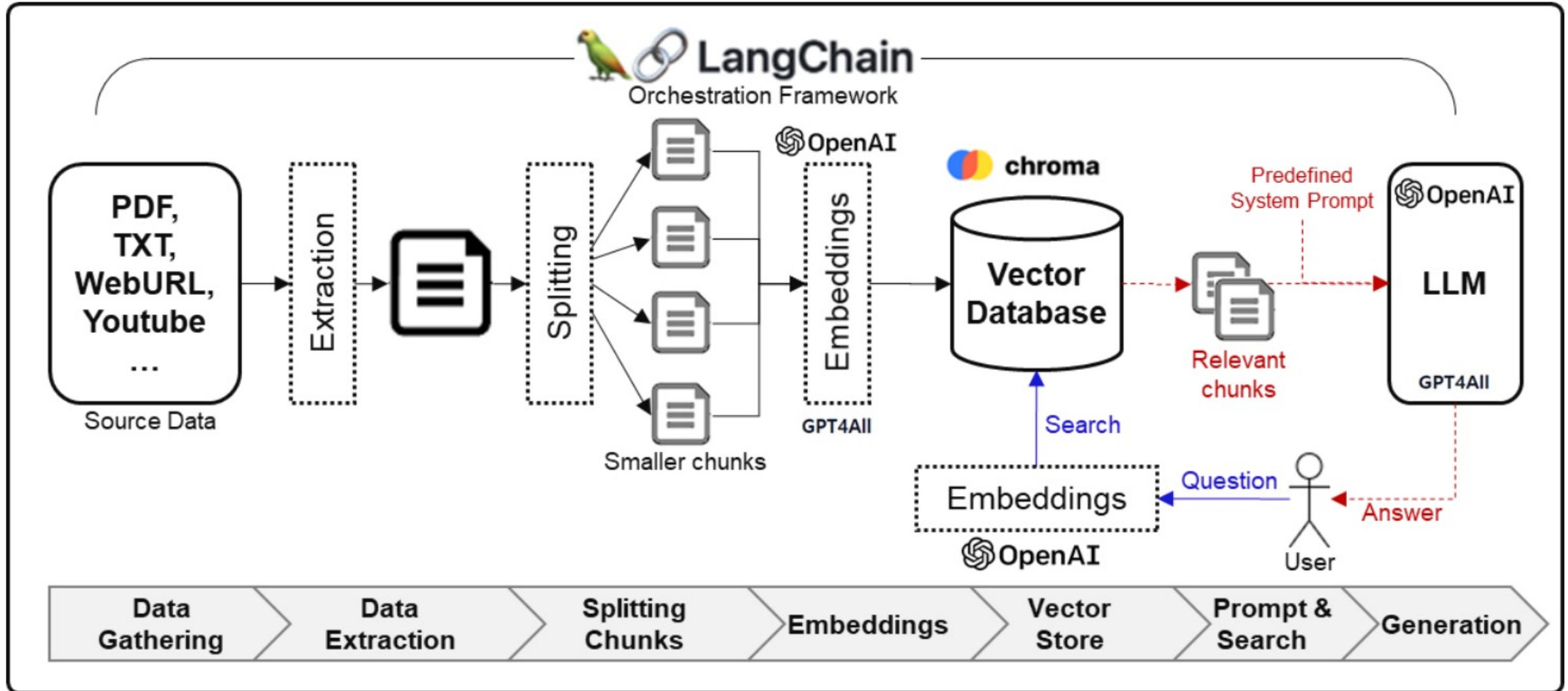
# **Dialogue Systems**

# Technology Tree of RAG Research

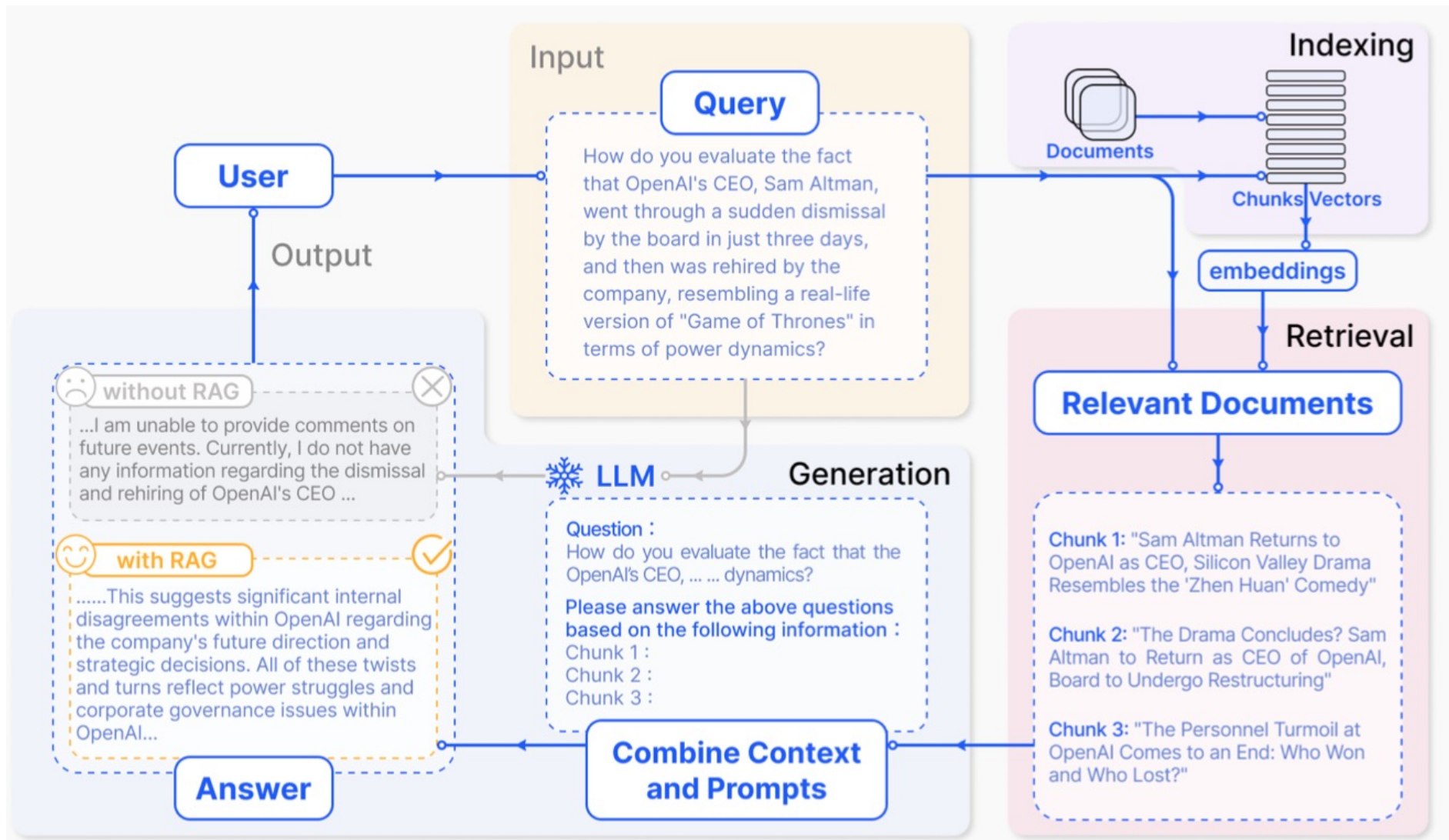
## Retrieval-Augmented Generation (RAG) for Large Language Models (LLMs)



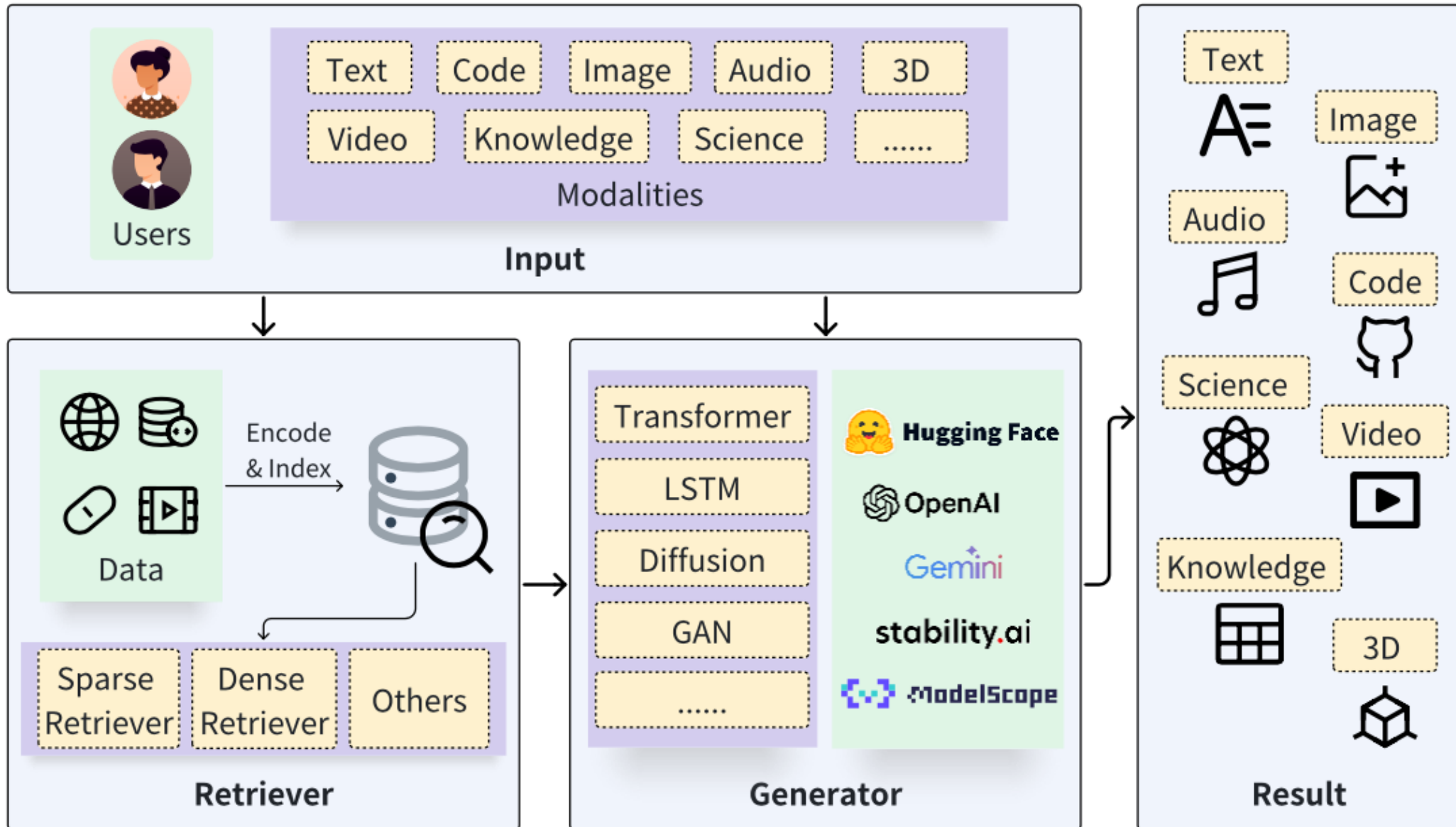
# Framework for Implementing Generative AI Services using RAG Model



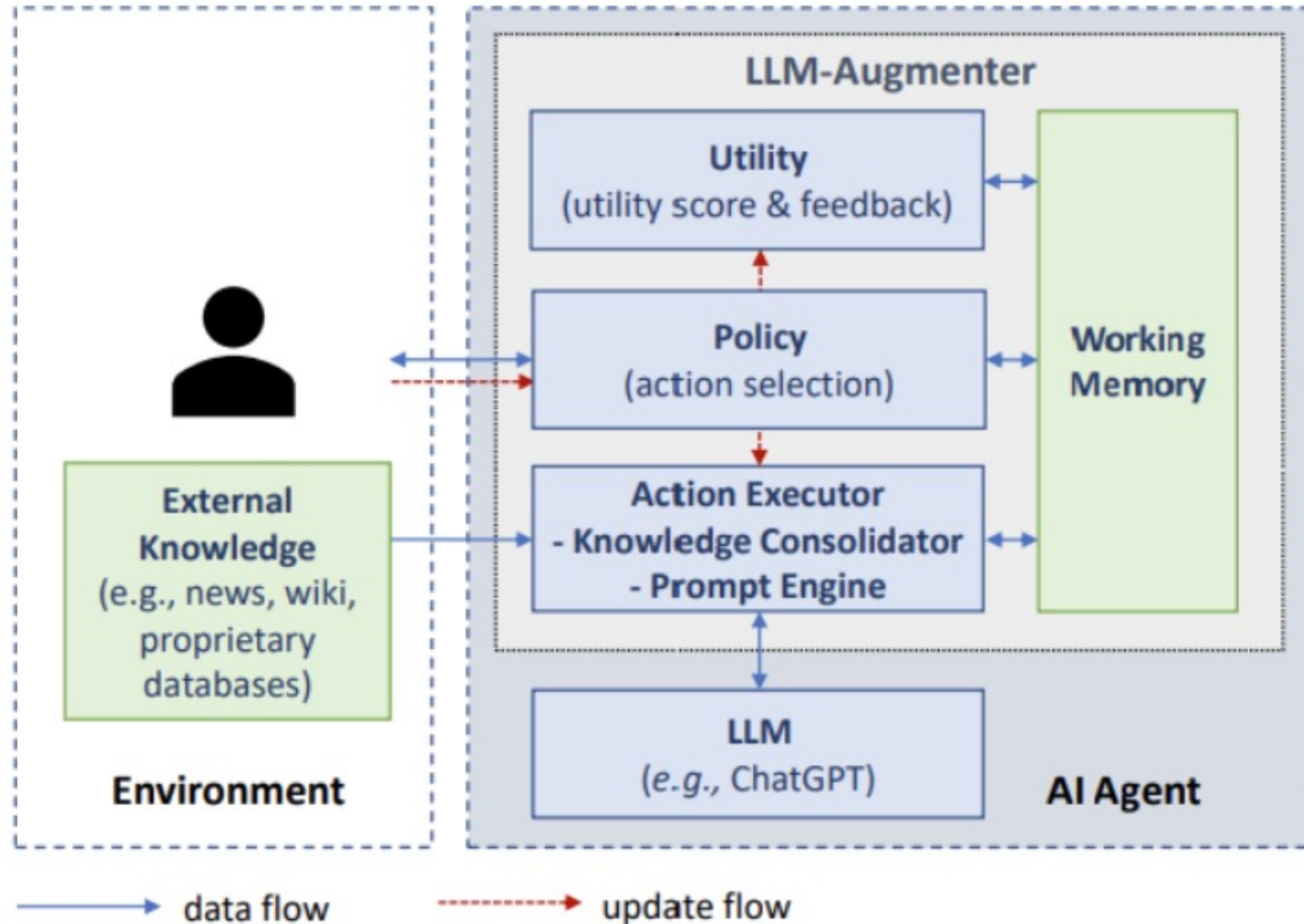
# Retrieval-Augmented Generation (RAG) for Large Language Models (LLMs)



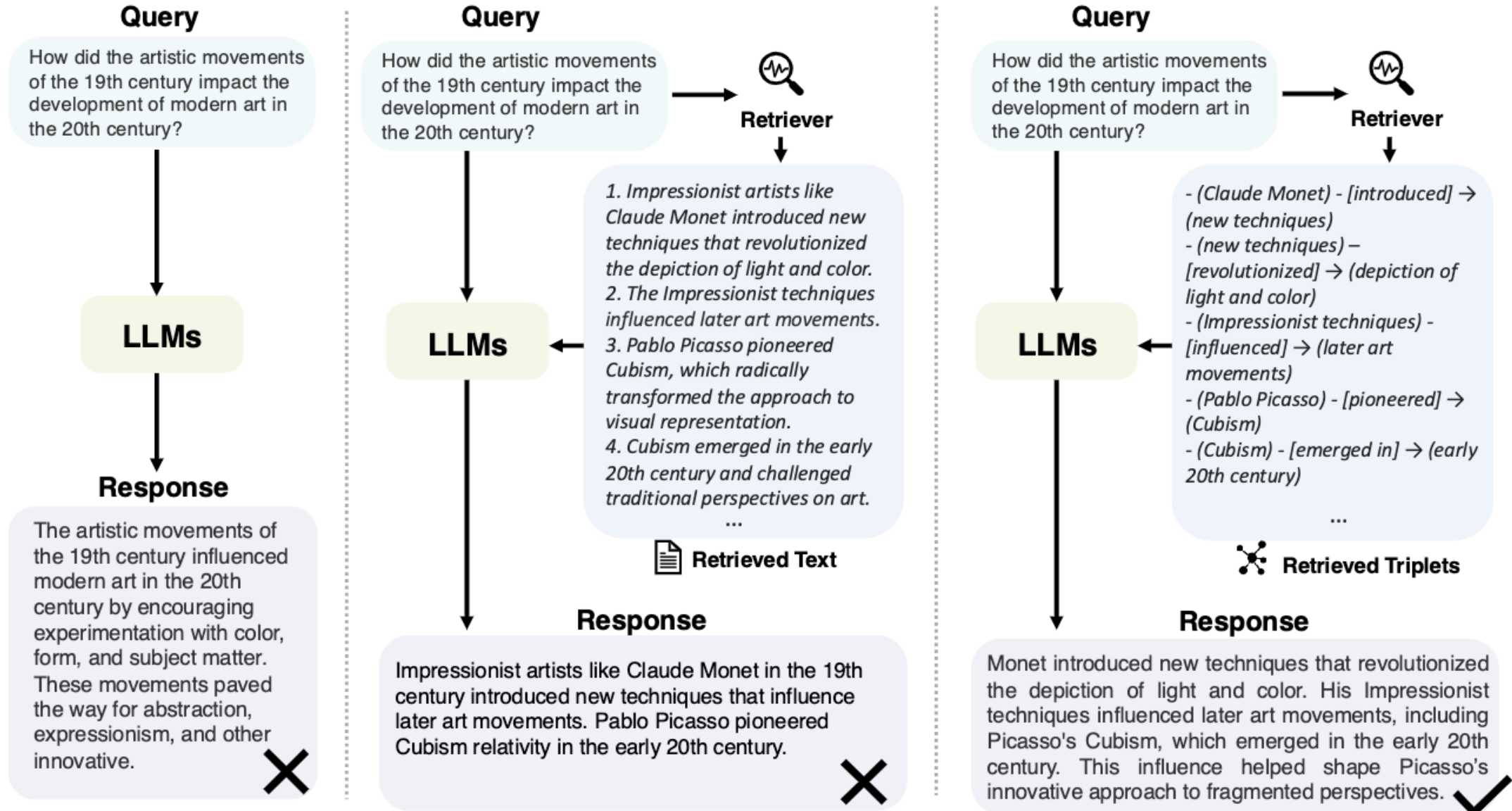
# Retrieval-Augmented Generation (RAG) Architecture



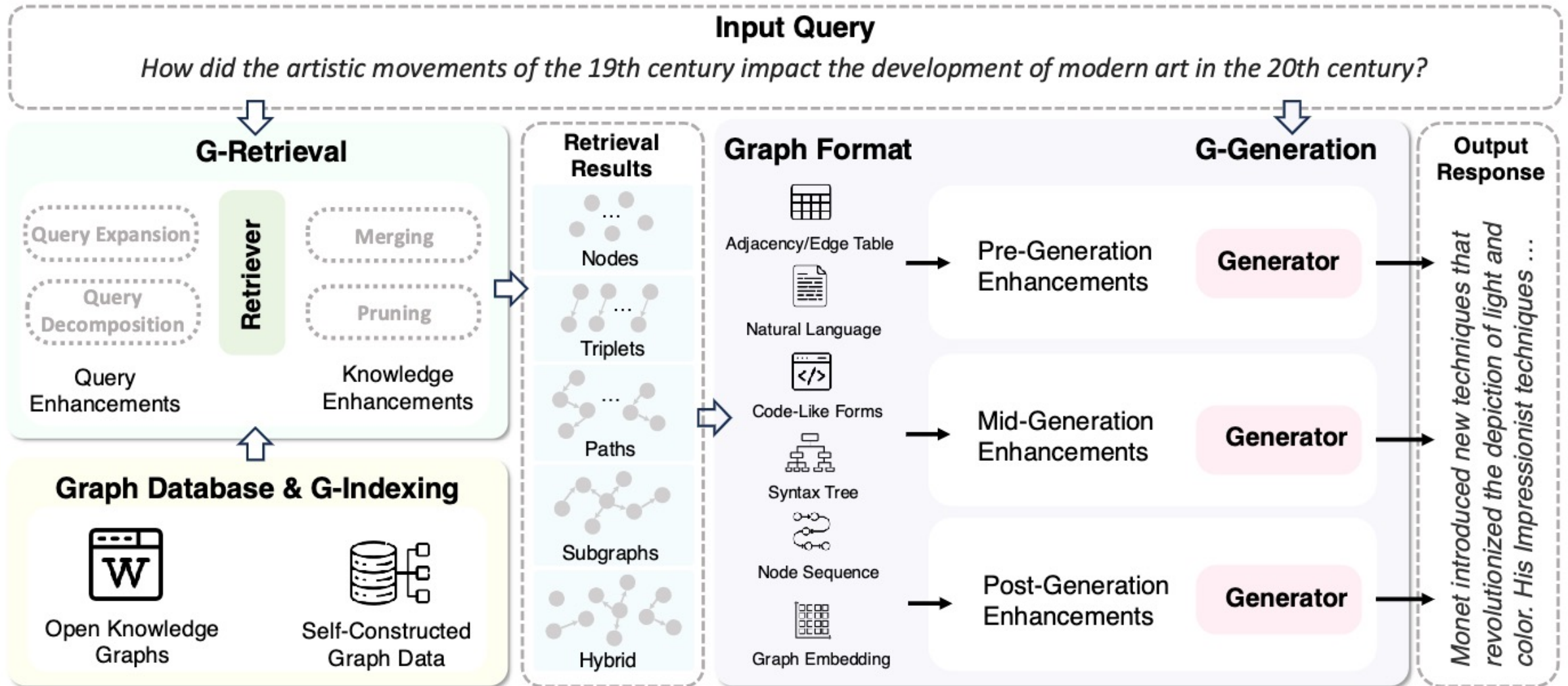
# A LLM-based Agent for Conversational Information Seeking



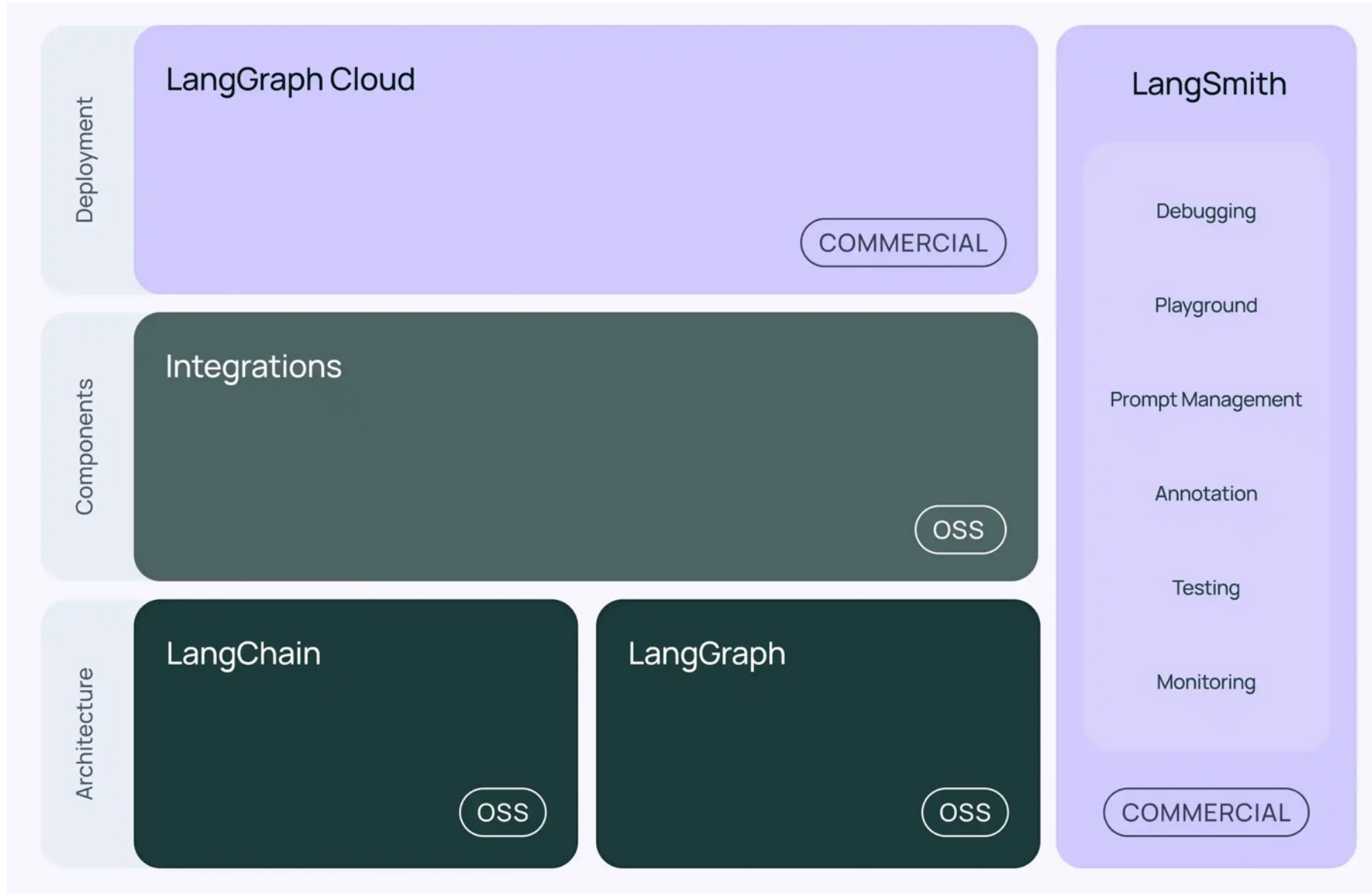
# Direct LLM, RAG, and GraphRAG



# GraphRAG Framework for Question Answering

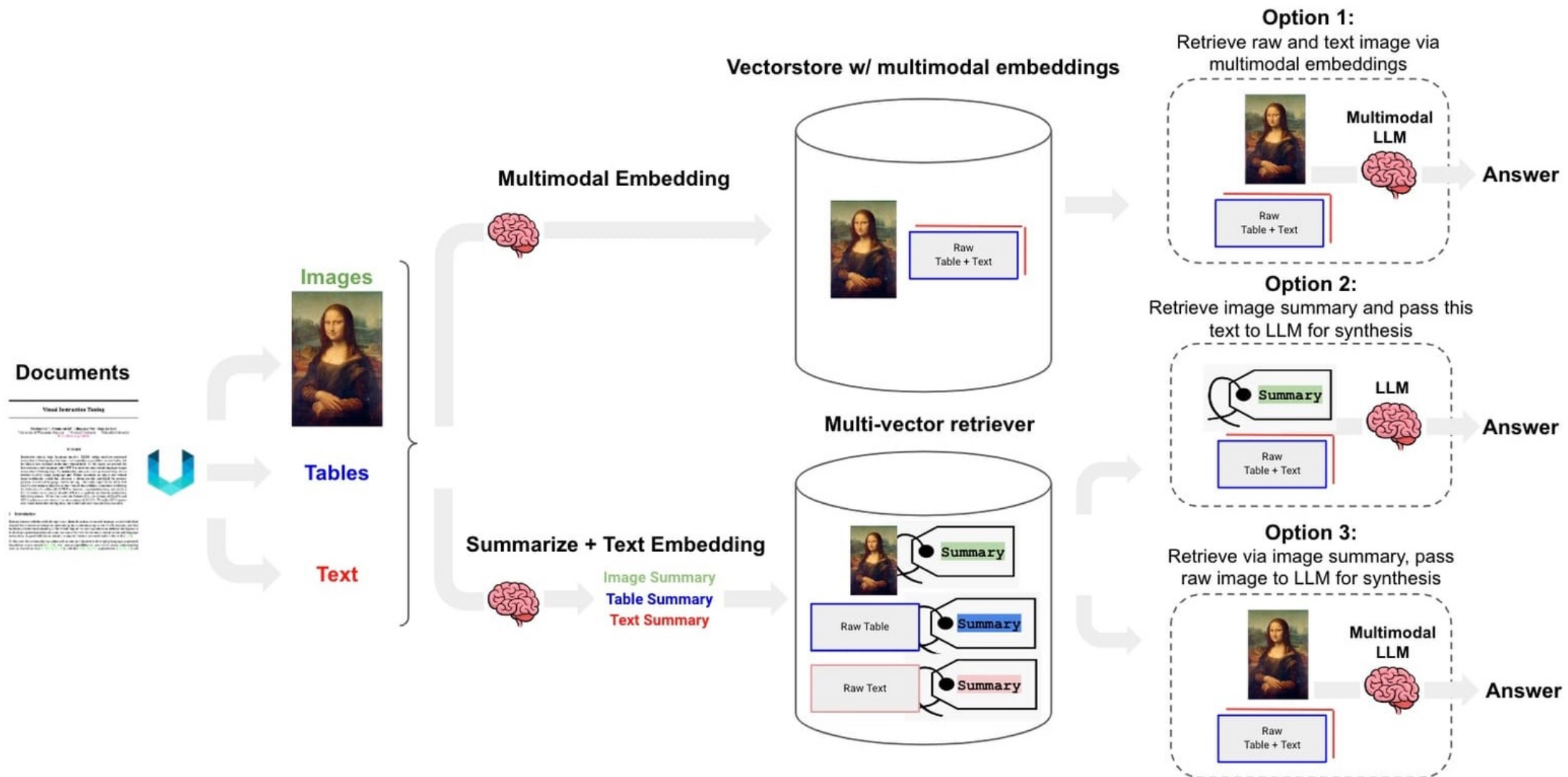


# LangChain Architecture



# Multimodal LLM RAG

## Multi-Vector Retriever for RAG



# Evaluating RAG with Ragas Metrics

## ragas score

generation

### faithfulness

how factually accurate is  
the generated answer

### answer relevancy

how relevant is the generated  
answer to the question

retrieval

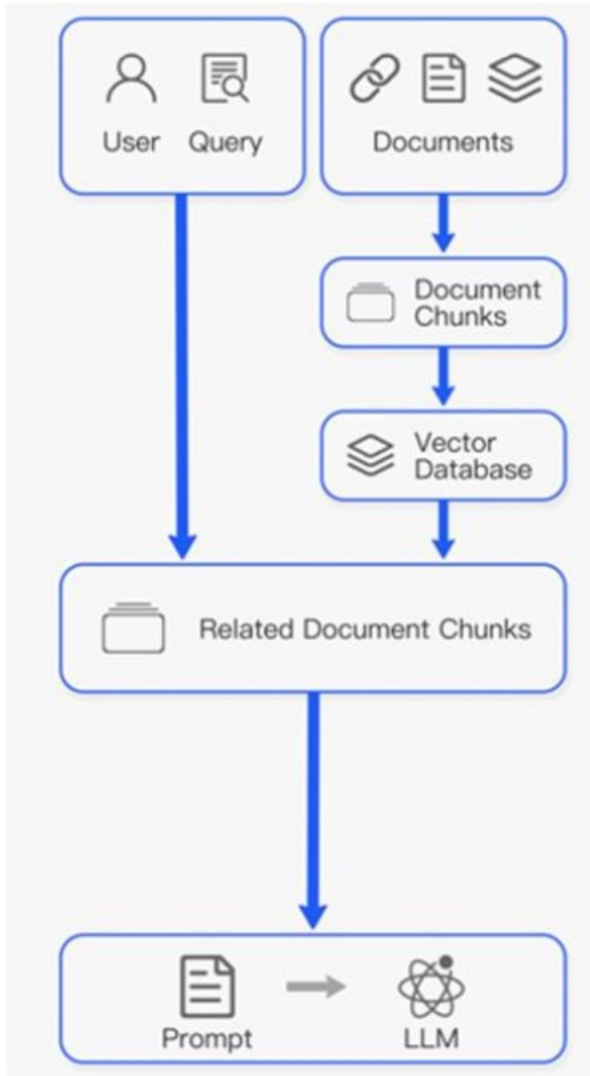
### context precision

the signal to noise ratio of retrieved  
context

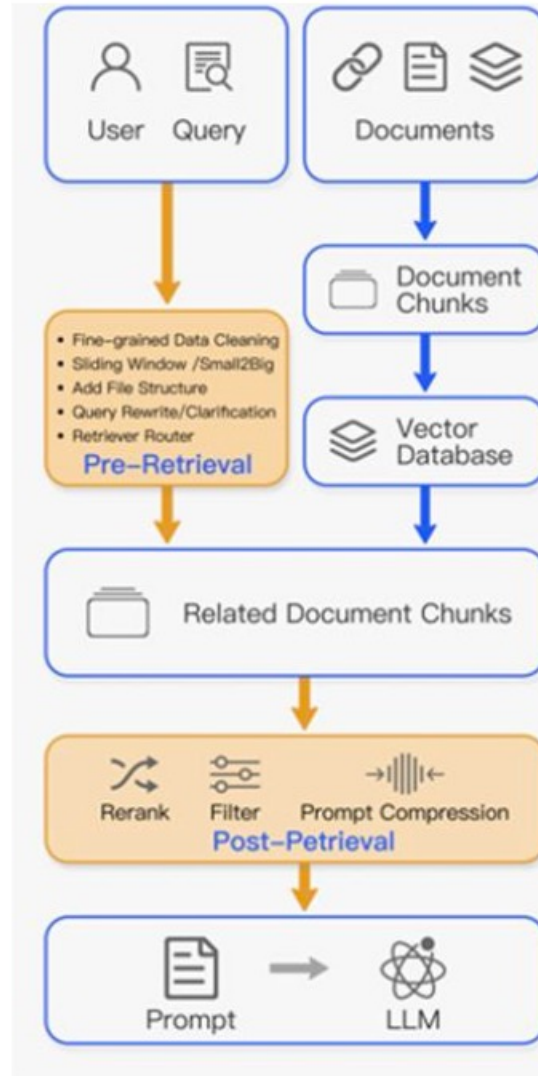
### context recall

can it retrieve all the relevant information  
required to answer the question

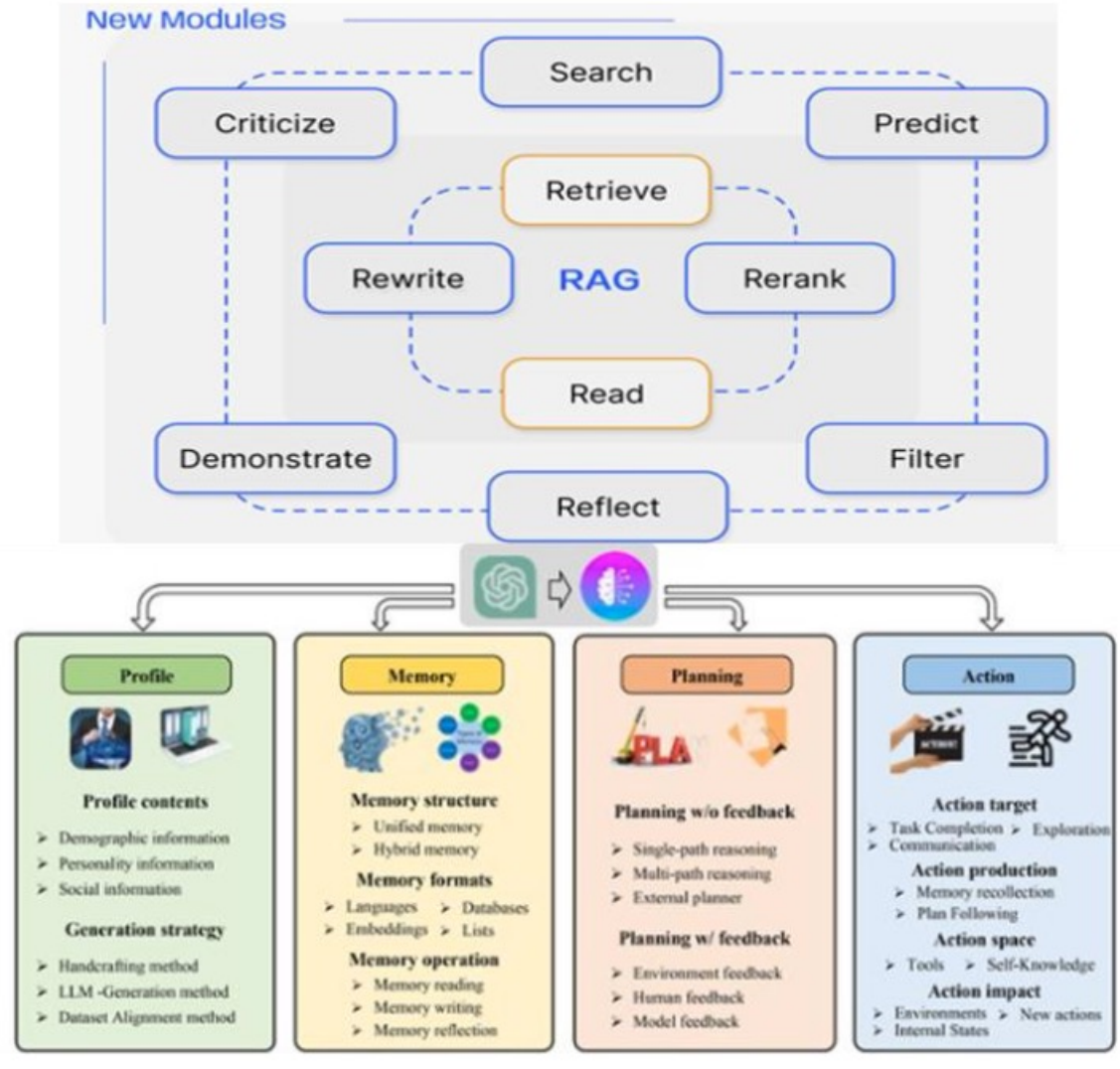
# Native RAG, Advanced RAG, Agentic RAG



Native RAG

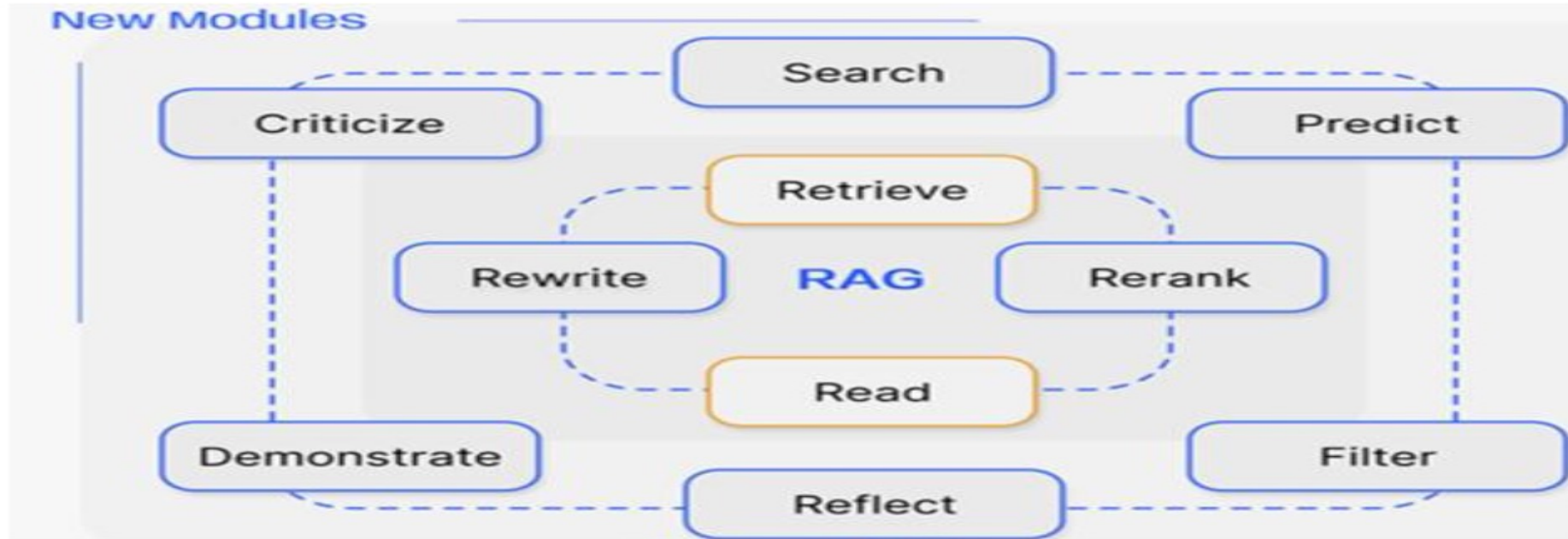


Advanced RAG

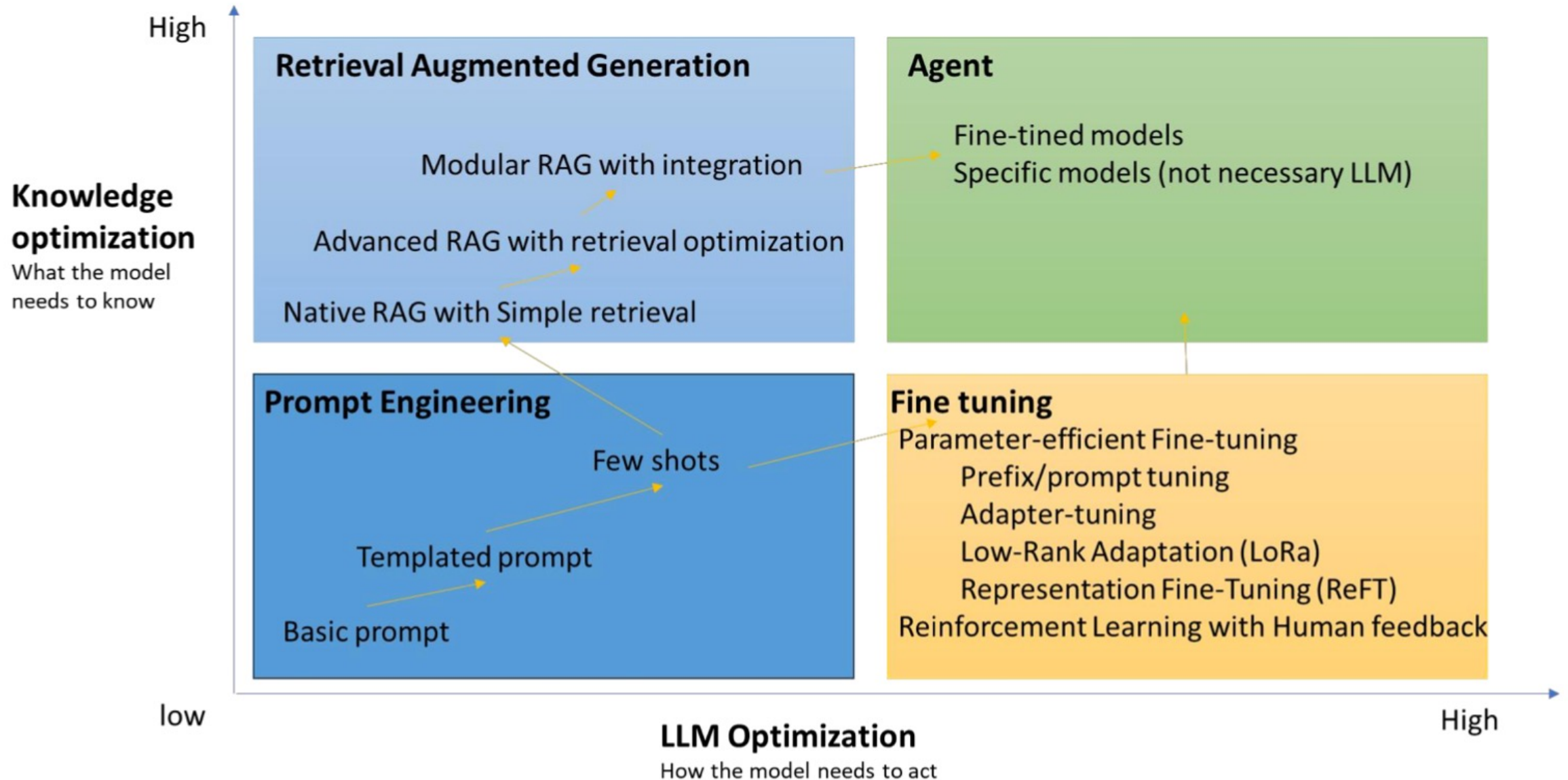


Modular + Agentic RAG

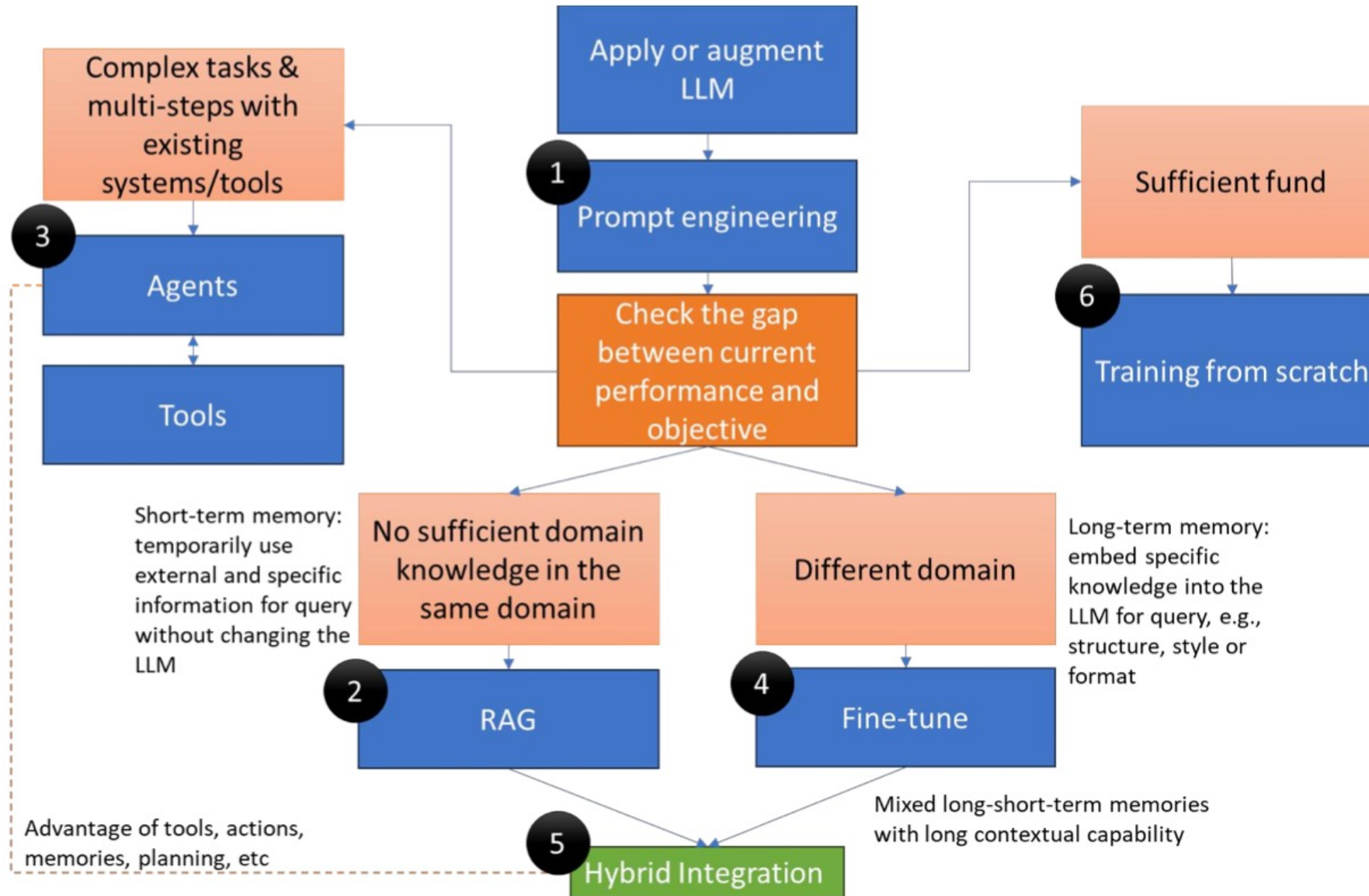
# Modular + Agentic RAG



# LLM Prompt, RAG, Agent, Fine tuning



# LLM Decision Path for Suitable Techniques



# LLM Stack

Layer	Sub-layer	
Application	Front-end service	Web services   Playground
	Back-end service	Input handler   output handler   Response formation
Orchestration		Deployment service   Governance service integration   Monitoring and reporting service   logging service   CI/CD pipeline
	LLM utility	Training   Fine-tuning   Prompt template   Agent service   Evaluation & Validation   Model hub   Guardrails
LLM	LLM API service	Embedding service   Generation service   Reasoning service   Planning service   LLM cache   LLM memory   Tooling API
	LLM API gateway	External LLM gateway   Internal LLM gateway   Application unified gateway
	Data storage	Data Analytics Store   KG / semantic layer   Vector database   Feature store   Application database
Data	Data engineering	Access control   Data processing   Quality control & review   Privacy & security   Data asset (active metadata)
	Data provision	Data lakehouse   Data lake   Data warehouse   In-memory storage
Infrastructure	Hardware	Cloud / on-premise / hybrid   Accelerator (GPU/TPU/FPGA/APU/ASIC/etc)   Quantum   Security   Network   Distributed/elastic / cluster

# LLM Functionalities

Service management

Scenario management	
Configuration	Task creation
Knowledge management	Application management
pipeline	Automated test

Model/data management		
Metadata / data annotation	Collaboration	Content moderation
Data mining & synthesis	Security and privacy: PII detection & masking	

Agent/prompt management		
Prompt studio / recipe	Agent studio / template	Tool management
NLP recipe	Safety and compliance	

Online services		
Security/personal Auth	Monitoring /Logging / cost reporting	audit trail
guardrails	Validation evaluation	
QC	API gateway	

Engineering

Prompt engineering	
Zero-short/ few-short	Chain/tree/graph-of-thought
Pattern-based template: co-star/automate	self-criticism / role-playing

RAG engineering		
Query rewriting	Routing	Rerank
Auto-merging	Recursive	Hybrid fusion
compression		Dense x

Agent engineering		
ReAct	self-criticism	memory
tooling	planning	decision
workflow	collaboration	

Fine-tuning engineering		
LoRa / QLoRa	quantization	SFT / RLHF
distillation	ReFT	Adaptor
PEFT	ZeRo /DeepSpeed	

Technology foundation

LLM (generation)		
OpenAI	Gemini / gemma	Claudia
ChatGLM	Llama	Mistral / Mixtral
Phi-MoE	Yi / Qwen	WizardLM
deepseek	Kimi	Blossom

Embedding	
BGE	FinGPT
Bert/ FinBert	M3E / ERNIE
Core algorithms	
Context fusion	memory
multimodal	Long context

Framework	
Langchain	Lammaindex
GPTflow	graphGPT
AutoAgent	HF Agents
Langroid	Haystack

Vector database		
pinecore	Chroma	pgvector
Faiss	Milvus	Qdrant

Feature store		
Feathr	Databricks	Feast
Vertex AI FS	Hopsworks	AWS FS

Source: Jun Xu (2024). "GenAI and LLM for Financial Institutions: A Corporate Strategic Survey." Available at SSRN 4988118.

# Domain-Specific Language Model (DSLML) for ESG

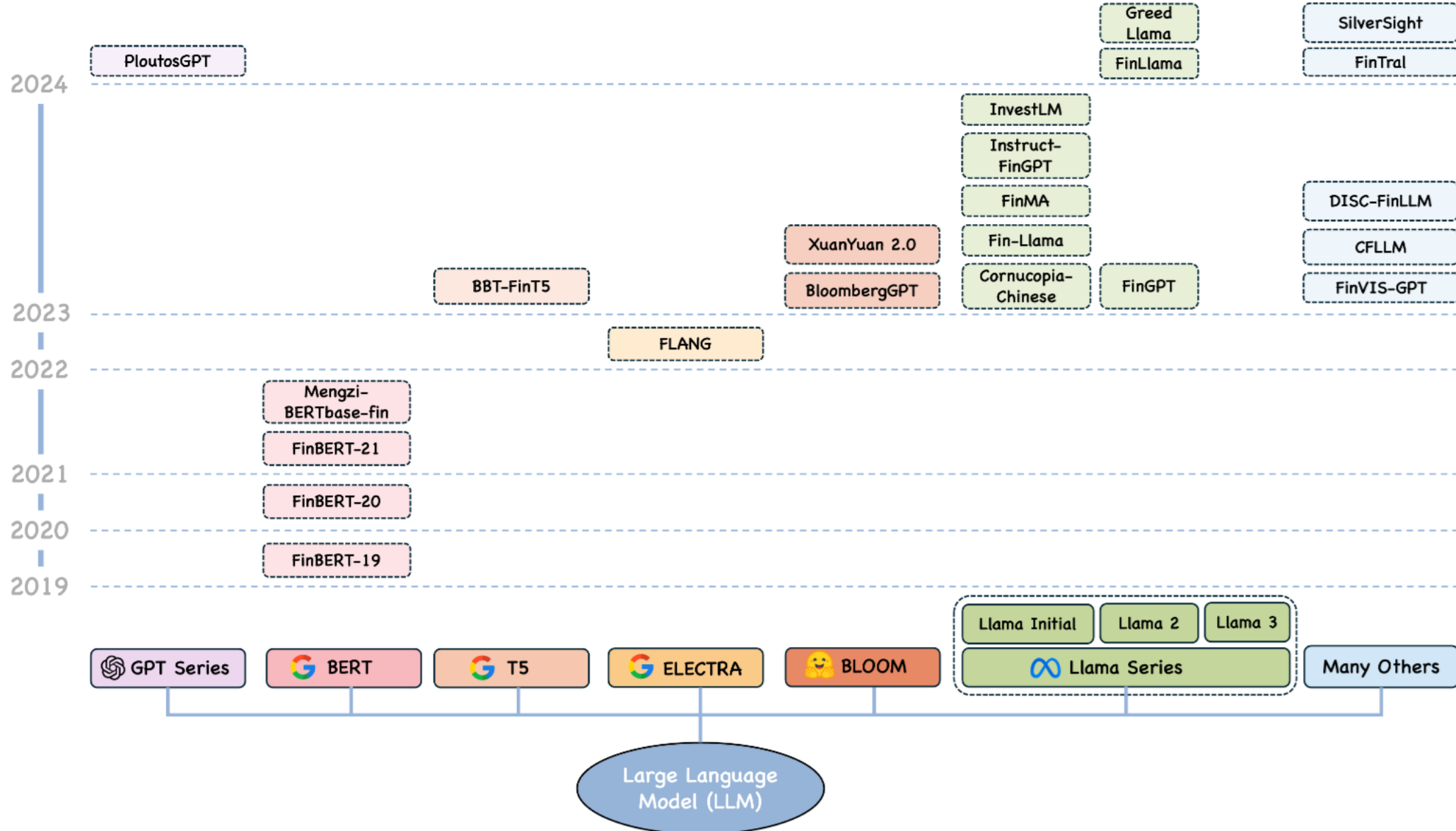
Model Architecture	Base Model	Parameter Count	Training Focus	Primary Use Case	Key Performance Metric
<b>ClimateBERT</b>	DistilRoBERTa	~82M	Climate news, science, & corporate reports	Fact-checking, Climate Specific Classification	48% improvement in domain MLM tasks (Webersinke et al., 2021)
<b>ESG-BERT</b>	BERT-Base	~110M	Sustainable investing corpora	Sentiment Analysis, Category Classification	F1-score of 0.90 on ESG text mining (Mehra et al., 2022)
<b>SusGen-GPT</b>	Llama/Mistral	7B - 8B	SusGen-30K (Financial + ESG tasks)	TCFD Report Generation, Relation Extraction	~2% performance gap to GPT-4 on specialized tasks (Wu et al., 2024)
<b>ClimateChat</b>	Llama 2 (tuned)	7B+	ClimateChat-Corpus (Q&A pairs)	Climate Science Q&A, Policy Analysis	Improved accuracy on scientific discovery tasks (ClimateChat Team, 2025)

# GenAI in ESG Datasets

Dataset Name	Size / Composition	Primary Application	Source
<b>SusGen-30K</b>	30,000 samples; 7 financial/ESG tasks. Sourced from TCFDHub & Hugging Face.	Training SusGen-GPT; Balancing financial NLP with ESG reporting tasks.	Wu et al. (2024)
<b>ClimateChat-Corpus</b>	Instruction-tuned Q&A pairs derived from scientific docs and web scraping.	Fine-tuning conversational agents for climate science.	ClimateChat Team (2025)
<b>A3CG</b>	1,679 sustainability reports (SGX companies); Annotated Aspect-Action pairs.	Greenwashing detection; Cross-category generalization testing.	Ong et al. (2025)
<b>ESG-Activities</b>	1,325 labeled text segments classified by EU ESG taxonomy.	Benchmarking LLM performance on granular ESG activity identification.	Intelligent ESG Evaluation Team (2025)
<b>FinRpt</b>	Chinese & English Equity Research Reports.	Evaluating automated generation of comprehensive financial reports.	Li et al. (2025)

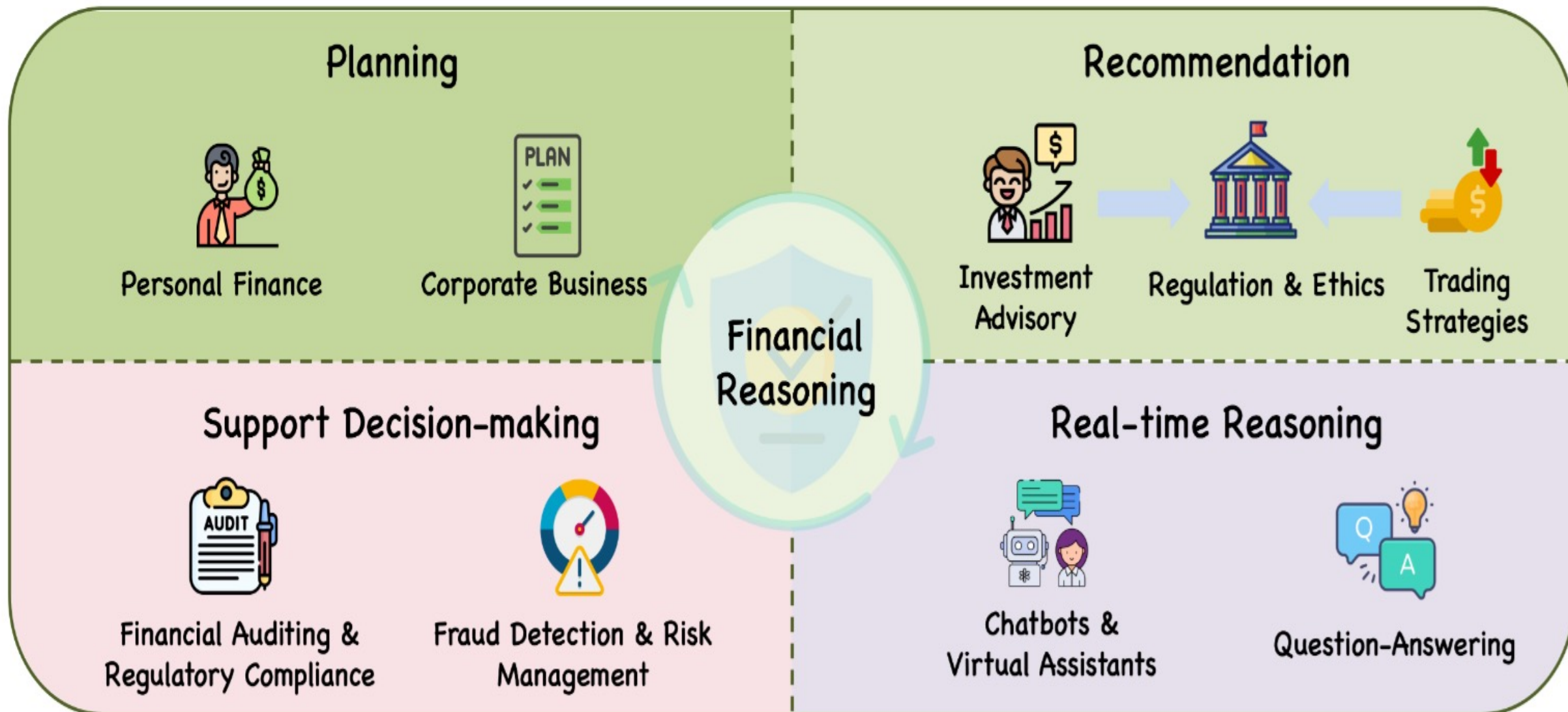


# Financial Large Language Models (Fin LLMs)

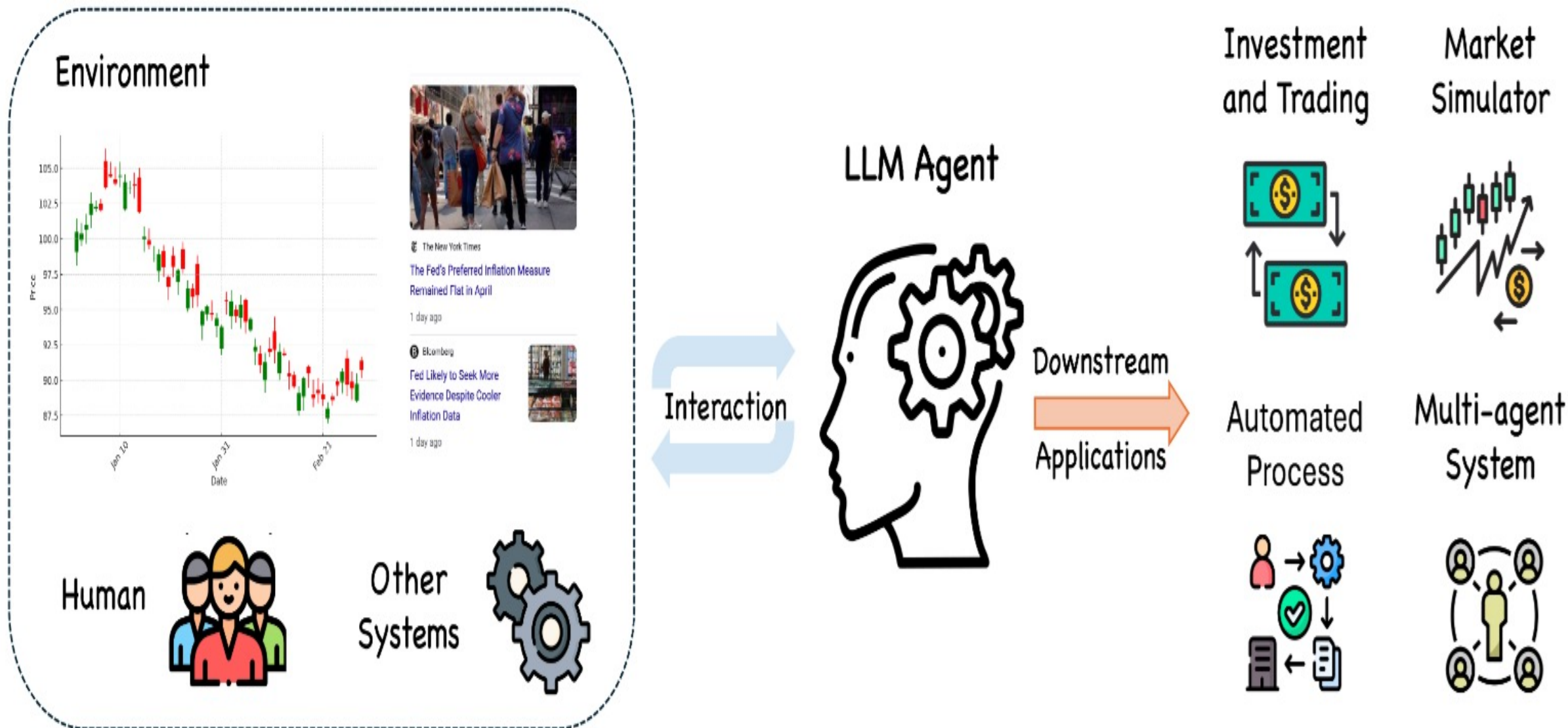


Source: Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M. Mulvey, H. Vincent Poor, Qingsong Wen, and Stefan Zohren (2024). "A Survey of Large Language Models for Financial Applications: Progress, Prospects and Challenges." arXiv preprint arXiv:2406.11903.

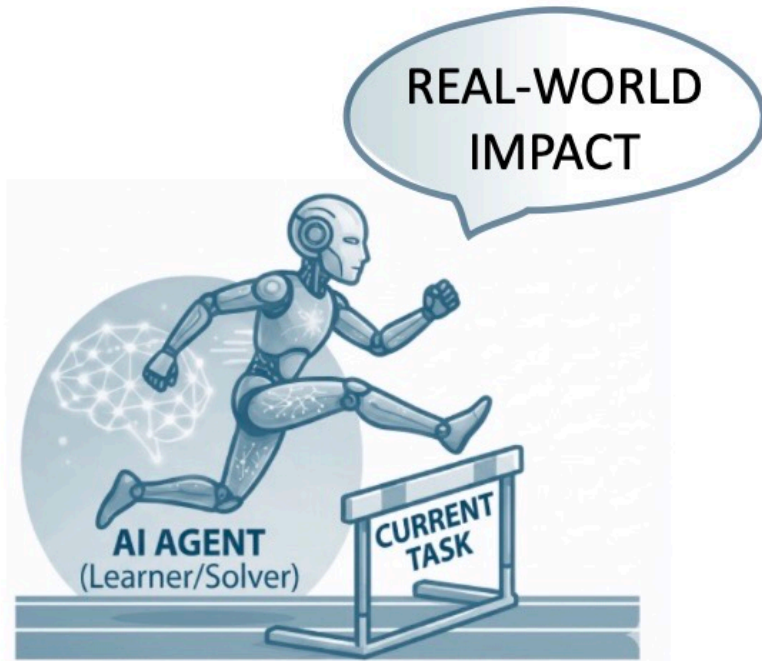
# Financial Reasoning Tasks



# Financial LLM Agent



# ESG Agent and ESG Benchmark



## ESG Agent



## ESG Benchmark



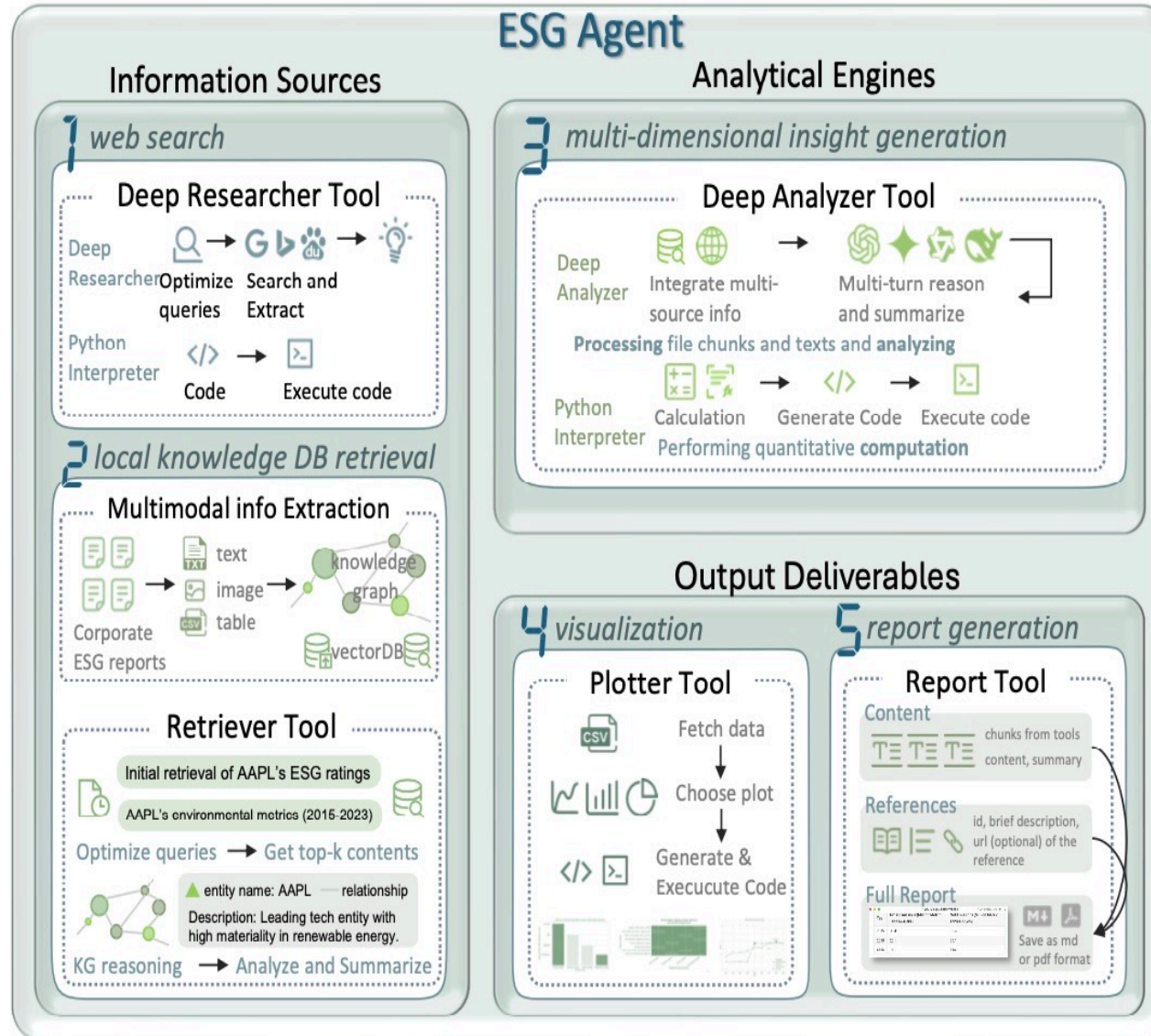
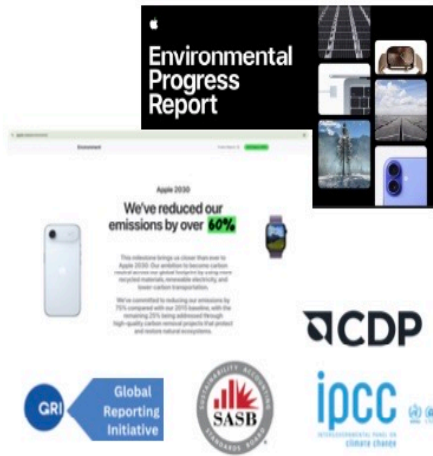
## Evaluation & Feedback



# ESG Agent Overall Architecture and Workflow

## Query / Task

Evaluate **AAPL's** environmental sustainability performance from 2015 to 2023 based on the provided corpus. Your response should include a quantitative trend of **ESG rating analysis** and at least one comparative insight regarding their **carbon neutrality goals**. Formulate your findings into a professional report structure.



Year	Gross Emissions (Million Metric Tons \$CO_2e\$)	Net Emissions (Million Metric Tons \$CO_2e\$)
2015	38.4	38.4
2019	25.1	25.1
2023	16.1	15.6

#### Apple 2030 vs. Industry Peers

Apple's 'Apple 2030' goal is to be carbon neutral across its entire value chain by 2030, which involves a 75% reduction in gross emissions from 2015 levels.

**Comparative Insight:** Compared to the IPCC recommendation of reaching net zero by 2050, Apple's target is 20 years ahead of global benchmarks. While peers like Microsoft have set ambitious 'carbon negative' targets by 2030, Apple's focus on a 75% absolute reduction in gross emissions highlights a rigorous approach to physical decarbonization rather than relying solely on offsets. [30, 31, 36]

#### 1.4 Visualizing the ESG Evolution

The following visualizations track Apple's ESG performance metrics and the strategic pivot in 2019 where environmental performance became a primary driver of the overall score.

Note: The chart highlights the significant improvement in Environmental and Social scores over the nine-year period, alongside the stability of Governance ratings. [36]

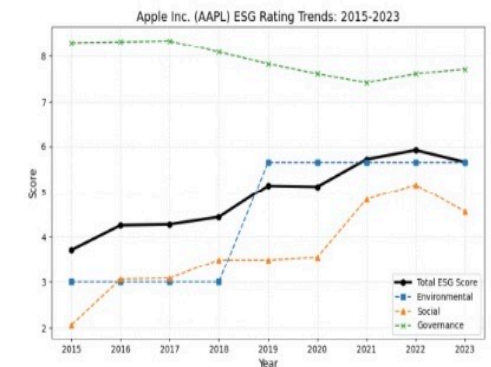


Chart: AAPL ESG Score Breakdown (2015-2023). Note the 2019 inflection point in Environmental scoring corresponding to 100% renewable energy achievement for global

# References

- Numa Dhamani and Maggie Engler (2024), Introduction to Generative AI, Manning
- Denis Rothman (2024), Transformers for Natural Language Processing and Computer Vision: Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3, 3rd Edition, Packt Publishing
- Thomas R. Caldwell (2025), The Agentic AI Bible: The Complete and Up-to-Date Guide to Design, Build, and Scale Goal-Driven, LLM-Powered Agents that Think, Execute and Evolve, Independently published
- Yilei Zhao, Wentao Zhang, Xiao Lei, Yandan Zheng, Mengpu Liu, and Wei Yang Bryan Lim. (2026) "Advancing ESG Intelligence: An Expert-level Agent and Comprehensive Benchmark for Sustainable Finance." arXiv preprint arXiv:2601.08676 (2026).
- Anthropic (2026), 2026 Agentic Coding Trends Report: How coding agents are reshaping software development, <https://resources.anthropic.com/hubfs/2026%20Agentic%20Coding%20Trends%20Report.pdf>
- NVIDIA DLI (2026), Building RAG Agents with LLMs, [https://learn.nvidia.com/courses/course-detail?course\\_id=course-v1:DLI+S-FX-15+V1](https://learn.nvidia.com/courses/course-detail?course_id=course-v1:DLI+S-FX-15+V1)
- NVIDIA DLI (2026), Generative AI with Diffusion Models, [https://learn.nvidia.com/courses/course-detail?course\\_id=course-v1:DLI+S-FX-14+V1](https://learn.nvidia.com/courses/course-detail?course_id=course-v1:DLI+S-FX-14+V1)
- Denis Rothman (2024), RAG-Driven Generative AI: Build custom retrieval augmented generation pipelines with LlamaIndex, Deep Lake, and Pinecone, Packt Publishing
- Jay Alammari and Maarten Grootendorst (2024), Hands-On Large Language Models: Language Understanding and Generation, O'Reilly Media
- Ben Auffarth (2023), Generative AI with LangChain: Build large language model (LLM) apps with Python, ChatGPT and other LLMs, Packt Publishing
- Chris Fregly, Antje Barth, and Shelbee Eigenbrode (2023), Generative AI on AWS: Building Context-Aware Multimodal Reasoning Applications, O'Reilly Media
- David Foster (2023), Generative Deep Learning: Teaching Machines to Paint, Write, Compose, and Play, 2nd Edition, O'Reilly & Associates Inc