

# 人工智慧新趨勢：AI 代理

## New Trends in Artificial Intelligence: AI Agent

Time: 2025/7/18 (五) 13:40-16:30

Place: E504電腦教室, 臺北市文山區萬美街二段21巷20號

Organizer: 臺北市政府公務人員訓練處 綜合企劃組 劉雨青



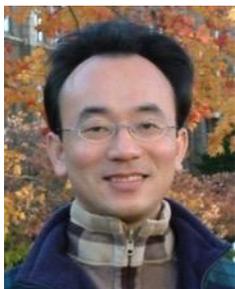
**戴敏育 教授 (Prof. Min-Yuh Day)**

國立臺北大學 資訊管理研究所 教授

金融科技暨綠色金融研究中心 主任

永續辦公室 永續發展組 組長





# 戴敏育 教授

## Prof. Min-Yuh Day



**Professor, Information Management, NTPU**

**Director, Intelligent Financial Innovation Technology, IFIT Lab, IM, NTPU**

**Director, Fintech and Green Finance Center (FGFC), NTPU**

**Division Director, Sustainable Development, Sustainability Office, NTPU**

**Visiting Scholar, IIS, Academia Sinica**

**Ph.D., Information Management, NTU**

Publications Co-Chairs, International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2013- )

Program Co-Chair, IEEE International Workshop on Empirical Methods for Recognizing Inference in Text (IEEE EM-RITE 2012- )

Publications Chair, The IEEE International Conference on Information Reuse and Integration for Data Science (IEEE IRI 2007- )



2020 Cohort



# 人工智慧新趨勢 (AI Agent)

課程目標：

透過智慧代理技術革新臺北市政府服務

- 深度探索 AI Agent 創新應用
- 聚焦政府實務操作與服務提升
- 整合臺北智慧城市創新實證計畫
- 為成果發表提供專業基礎

# Outline

1. AI Agent 基礎與最新趨勢
2. AI Agent 政府應用與智慧城市實作
3. AI Agent 實施策略與治理架構

# Innovative Agentic AI Technology for Autonomous ESG Report Generation

Industrial Technology Research Institute (ITRI),  
Fintech and Green Finance Center (FGFC, NTPU),  
NTPU-114A513E01, 2025/03/01~2025/12/31

# Generative AI-Driven ESG Report Generation Technology

Industrial Technology Research Institute (ITRI),  
Fintech and Green Finance Center (FGFC, NTPU),  
NTPU-113A513E01, 2024/03/01~2024/12/31

# Agentic AI Powering Digital Sustainability Innovation

# 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

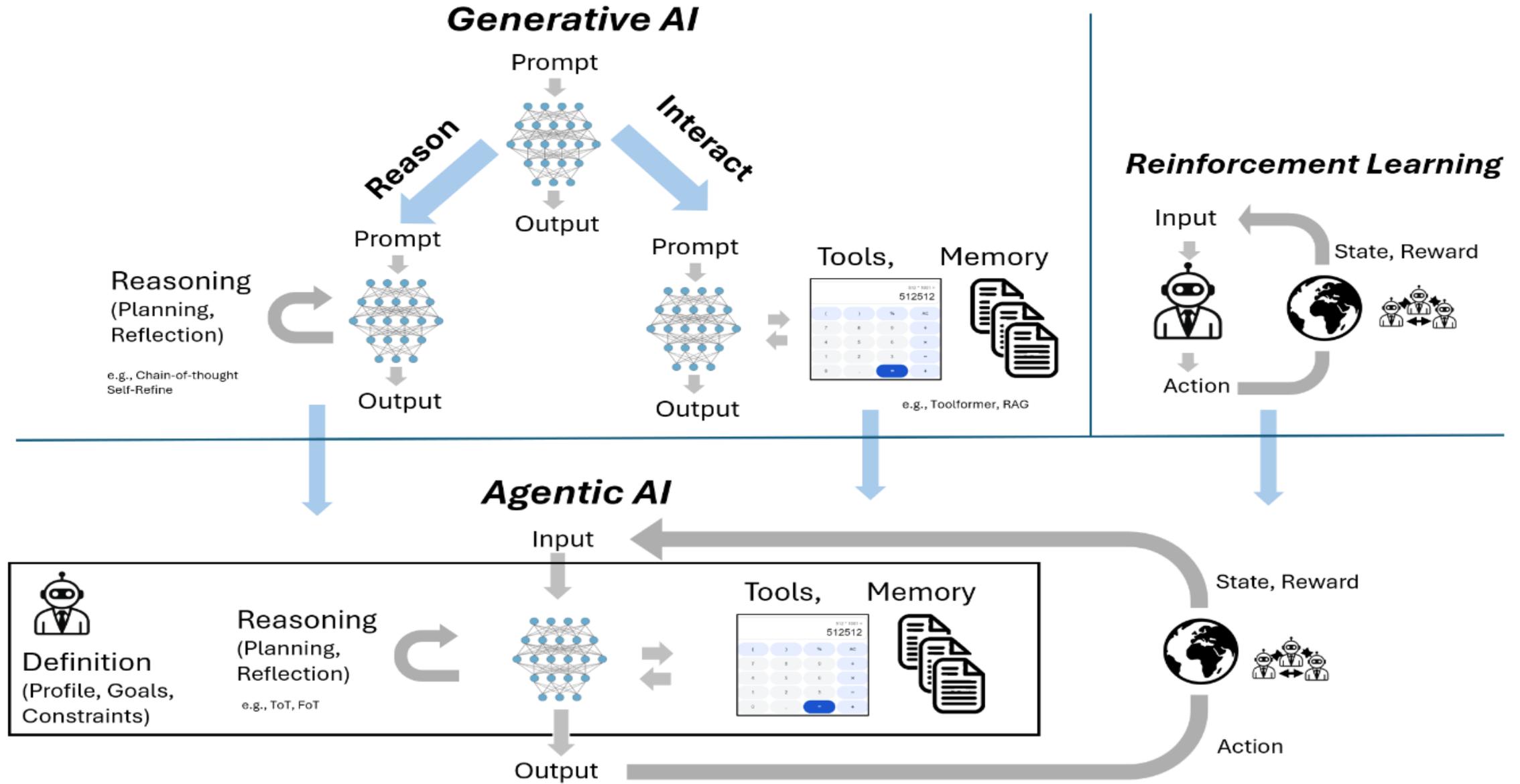
Deep learning breakthrough

# AI Agent 定義

## 從感知到行動的自主決策系統

- **AI Agent 定義**
  - 能夠感知環境、自主決策並採取行動以達成特定目標的智慧系統 (IBM, 2024)
- **AI Agent 與傳統軟體差異：**
  - 從固定規則轉向動態學習適應 (Ahmed, 2025)

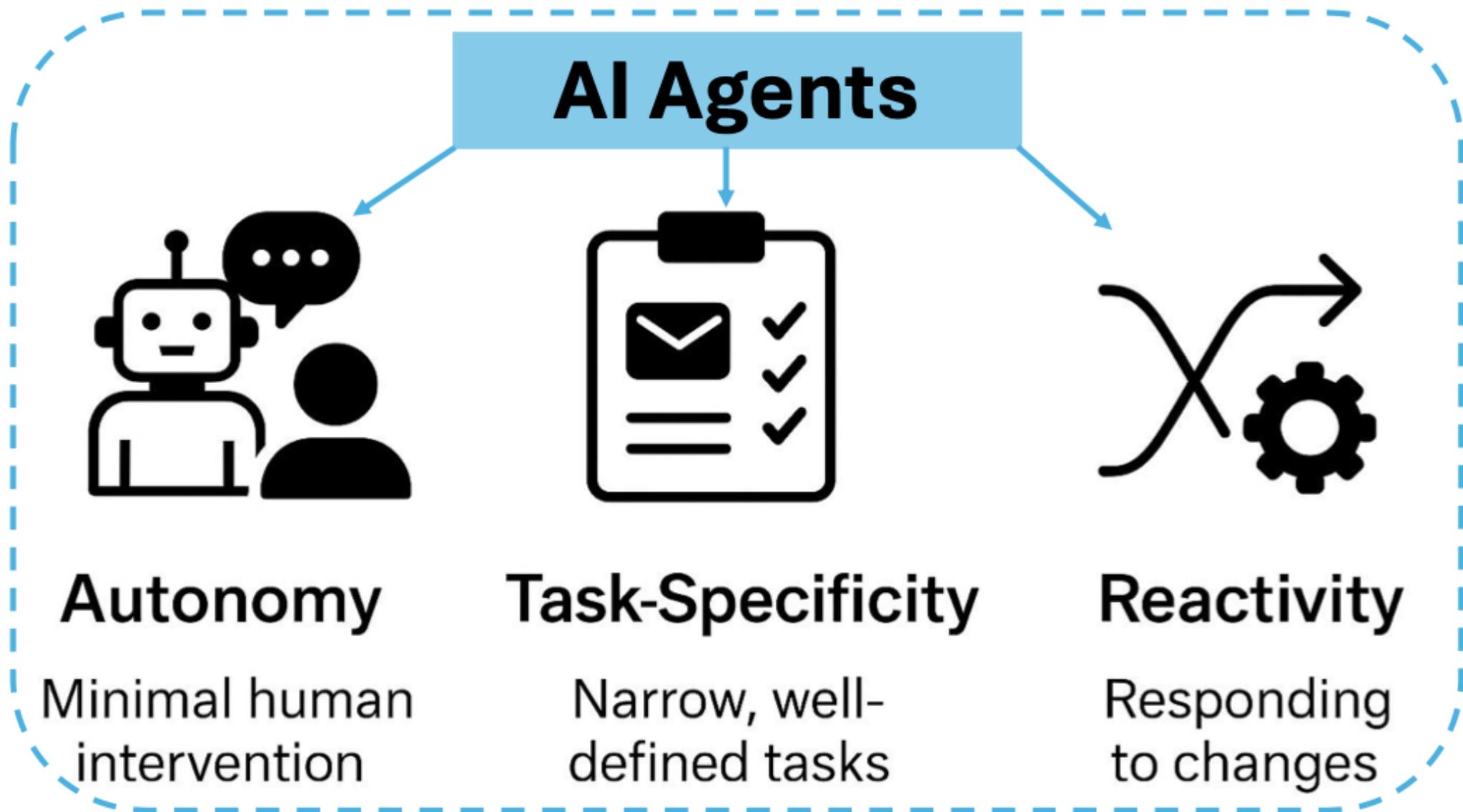
# From Generative AI to Agentic AI



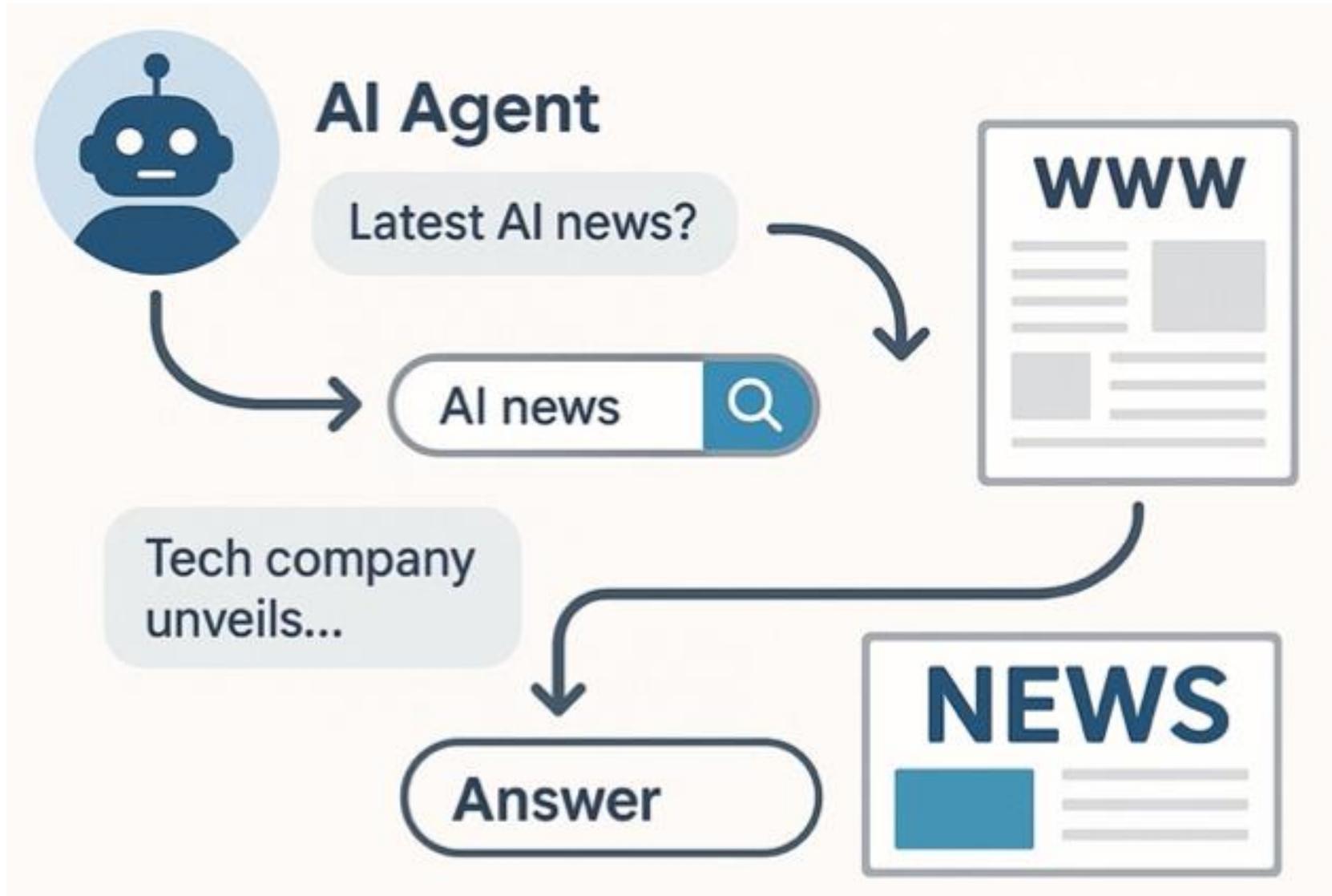
# AI Agent 核心能力

- **感知能力**：自然語言理解、多模態輸入處理
- **知識基礎**：混合式符號-分散式知識表徵系統
- **推理能力**：多步驟邏輯推論與因果分析
- **行動能力**：工具使用、API 整合、即時適應

# 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.

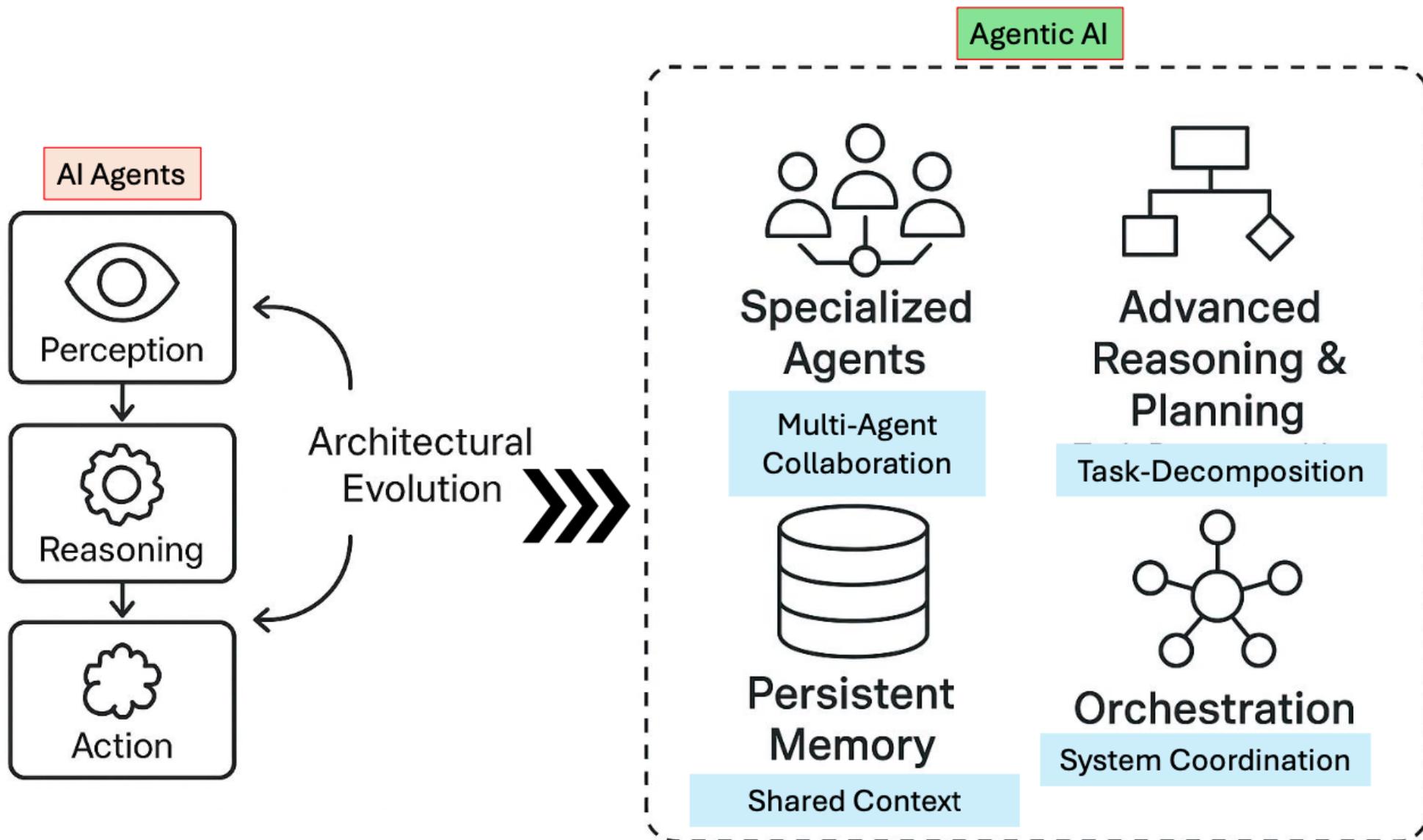
# AI Agent / Agentic AI, Generative AI, Traditional AI

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

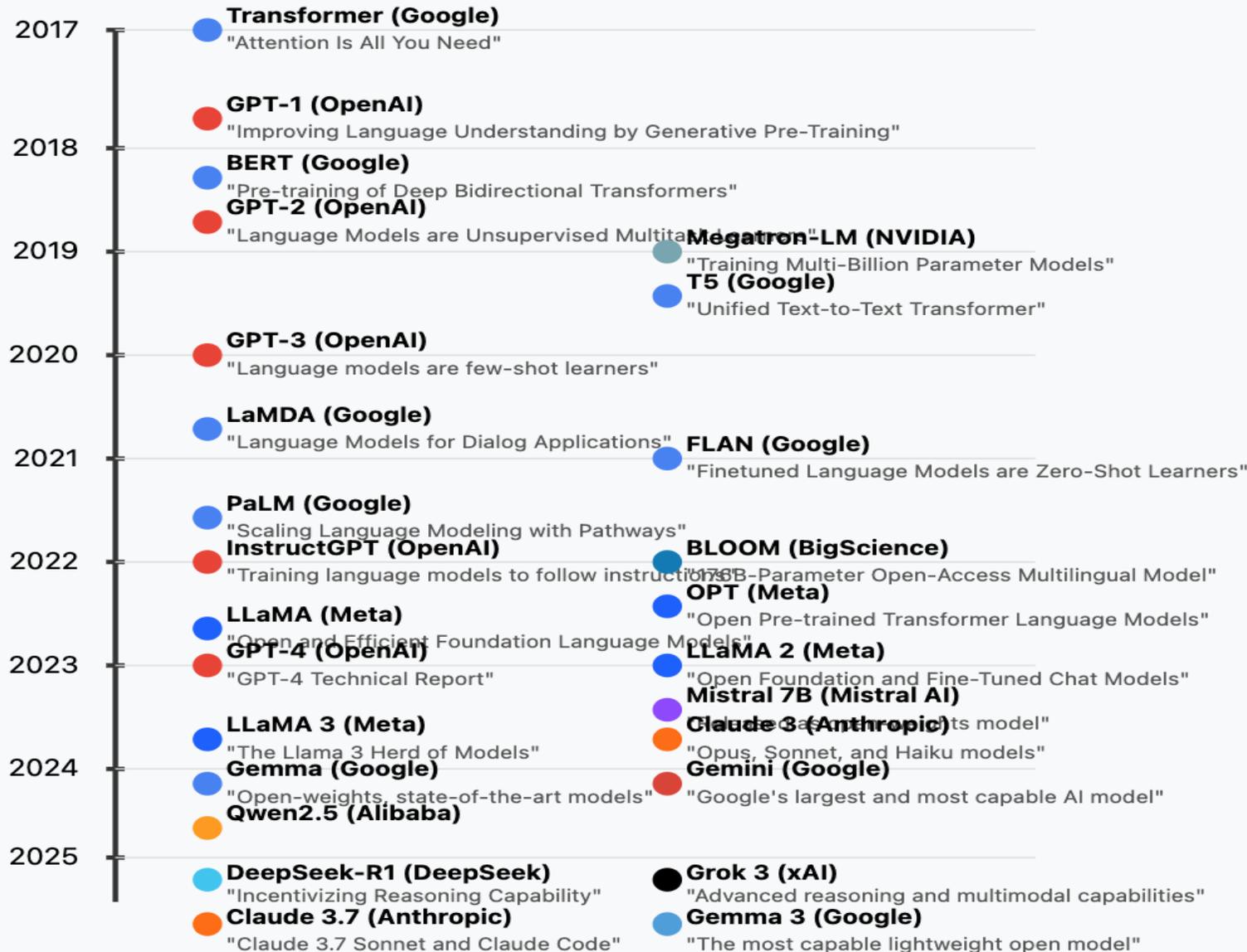
Feature	AI Agents	Agentic AI
<b>Definition</b>	Autonomous software programs that perform specific tasks.	Systems of multiple AI agents collaborating to achieve complex goals.
<b>Autonomy Level</b>	High autonomy within specific tasks.	Broad level of autonomy with the ability to manage multi-step, complex 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



# **AI Agents and Large Multimodal Agents (LMAs)**

# Generative AI LLMs (2017-2025)



### Key Organizations

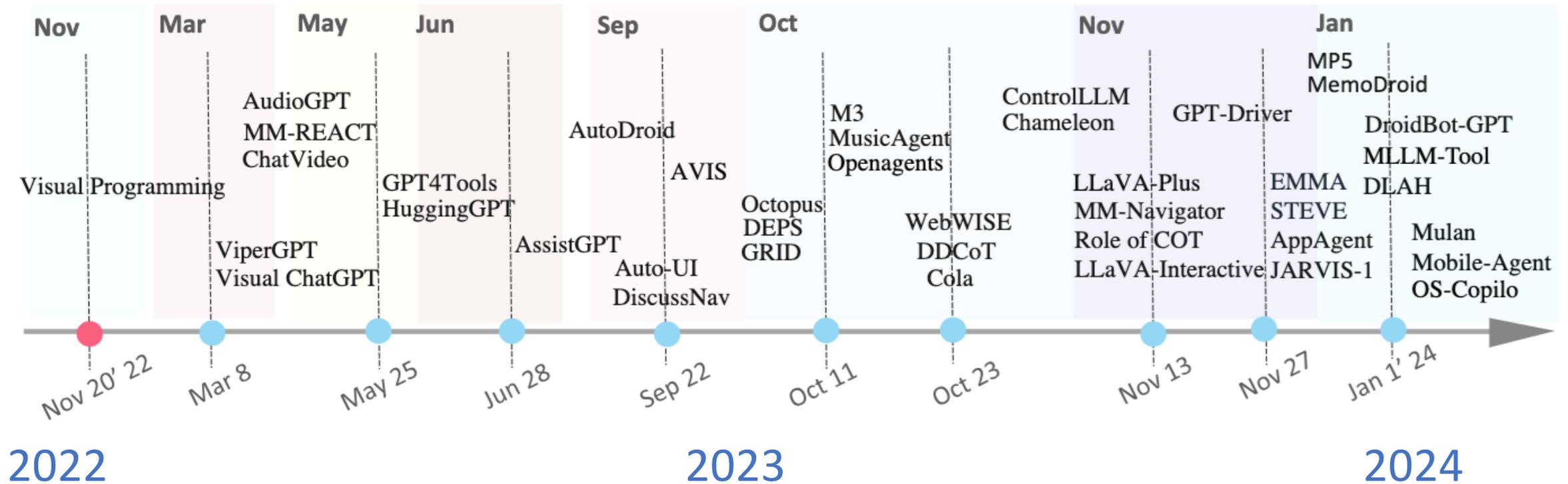
- Google
- OpenAI
- Meta
- Mistral AI
- Alibaba
- xAI
- Anthropic
- NVIDIA
- BigScience

### Key Milestones

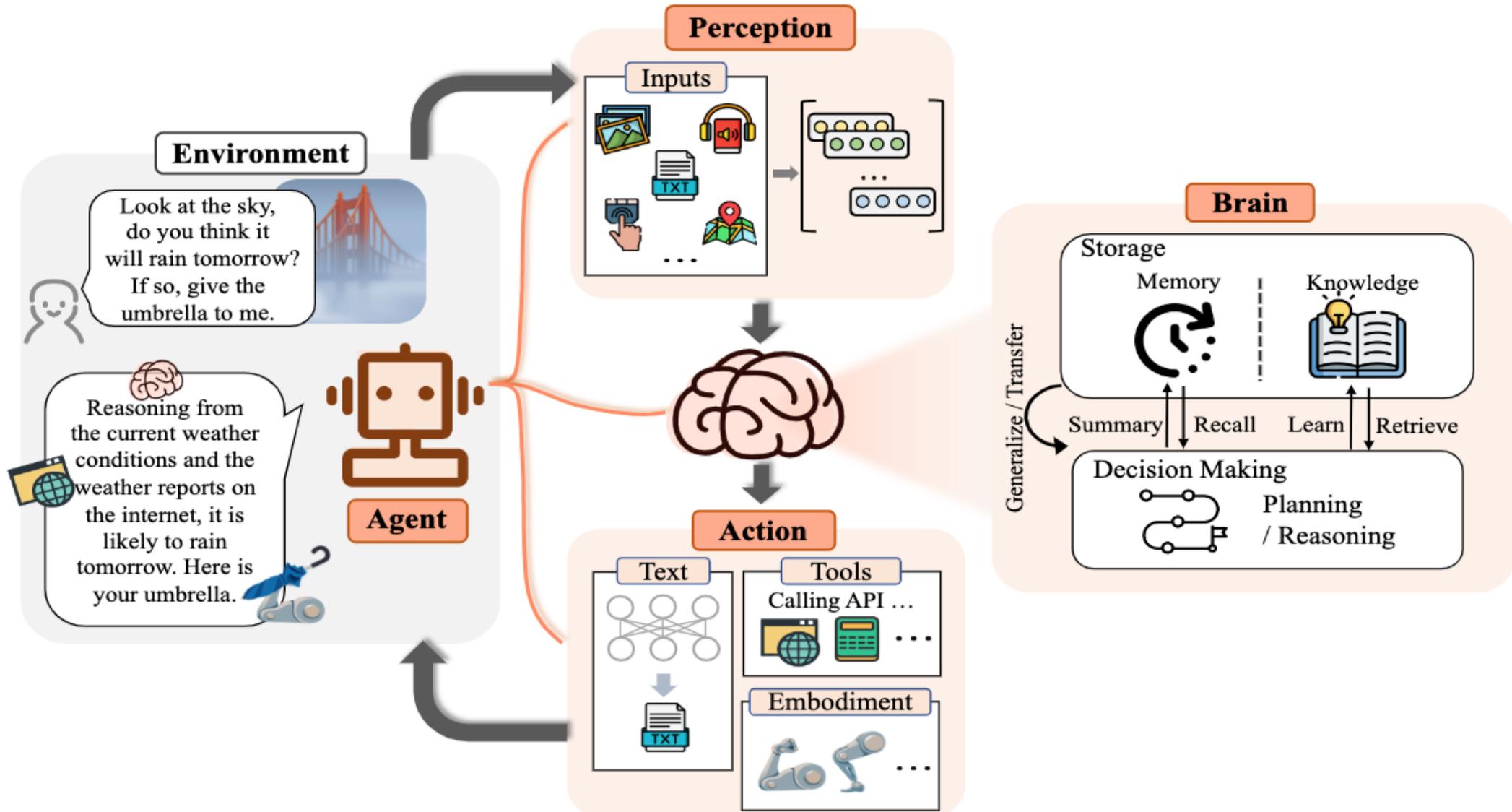
- 2017:** Transformer architecture
- 2018:** First-gen GPT, BERT
- 2020:** GPT-3 (175B parameters)
- 2022:** Emergent abilities, instruction tuning
- 2023:** GPT-4, multimodal models
- 2024:** Open-weights race, Mamba2
- 2025:** DeepSeek-R1, Grok 3  
Claude 3.7, Gemma 3

# LLM-powered Multimodal Agents

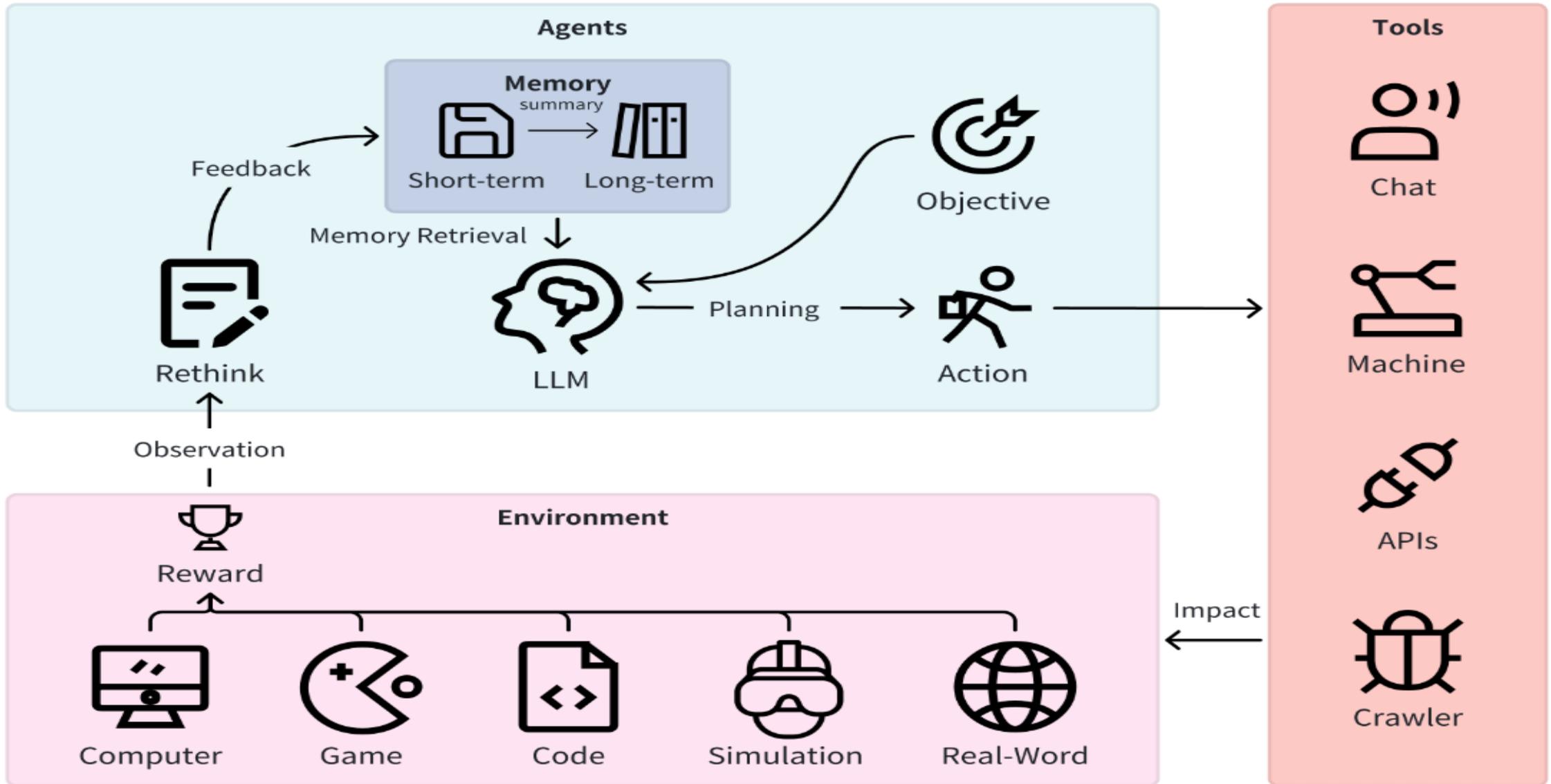
## Large Multimodal Agents (LMAs)



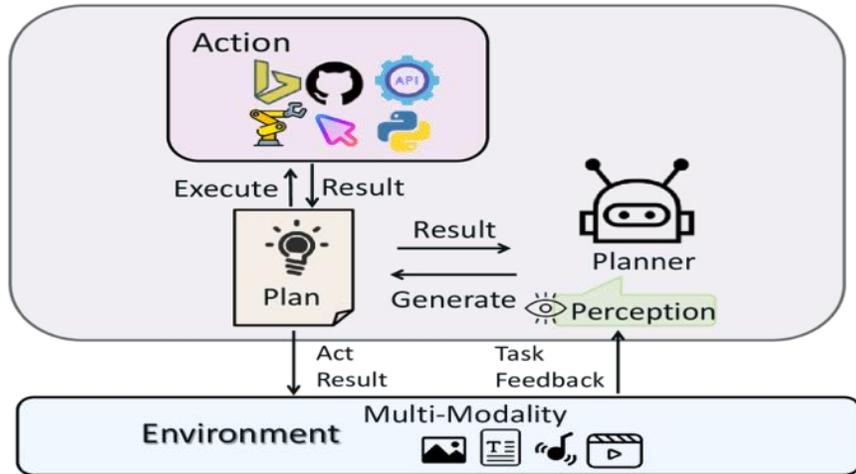
# Large Language Model (LLM) based Agents



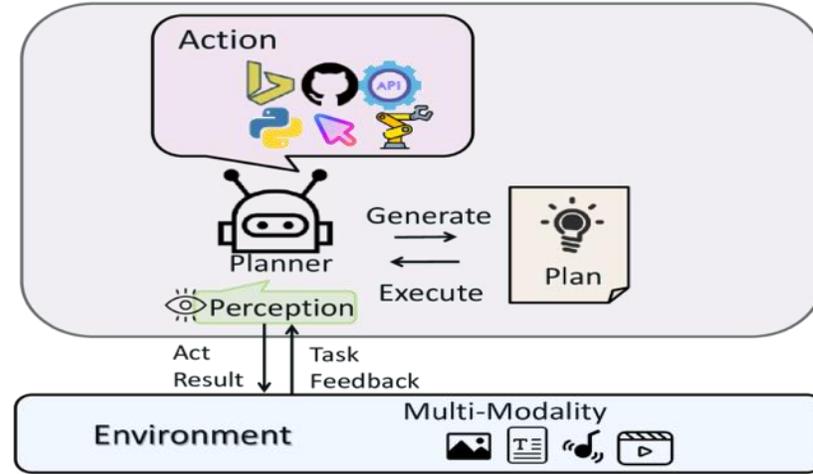
# LLM-based Agents



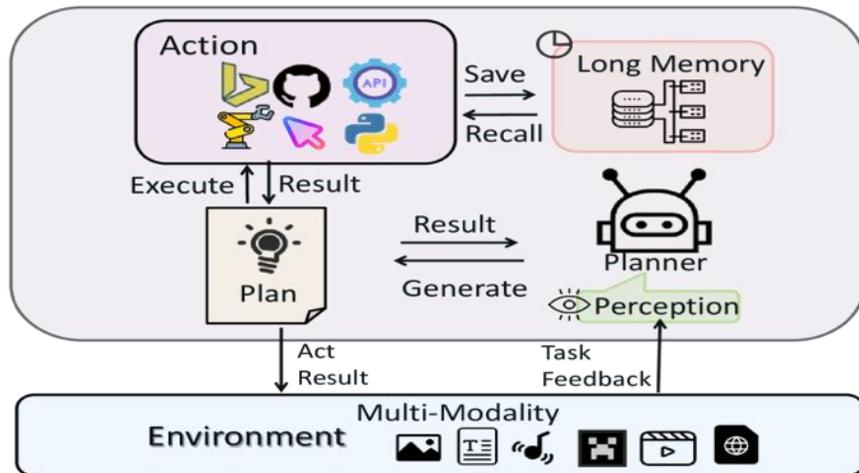
# Large Multimodal Agents (LMA)



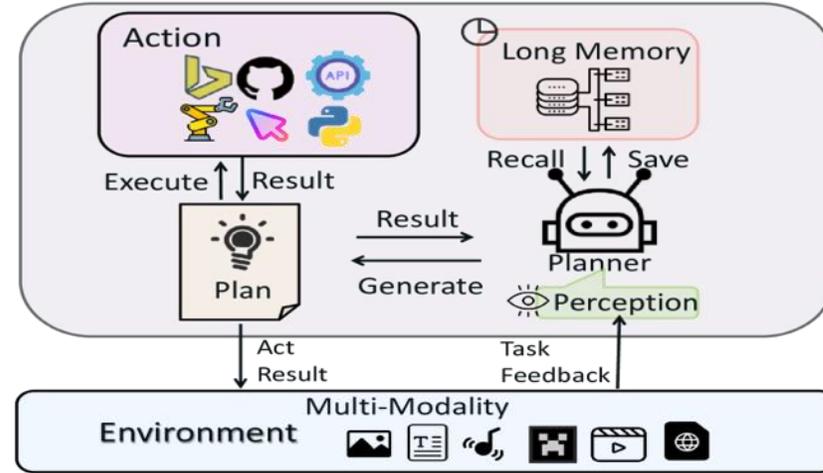
(a)



(b)

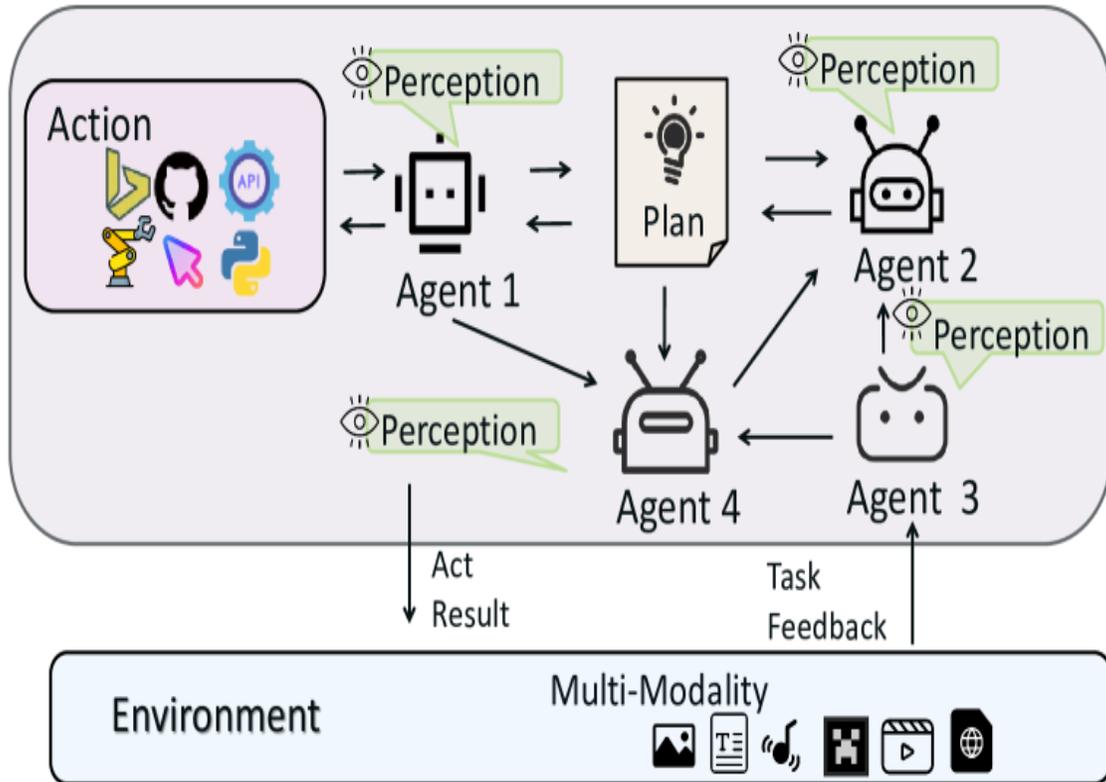


(c)

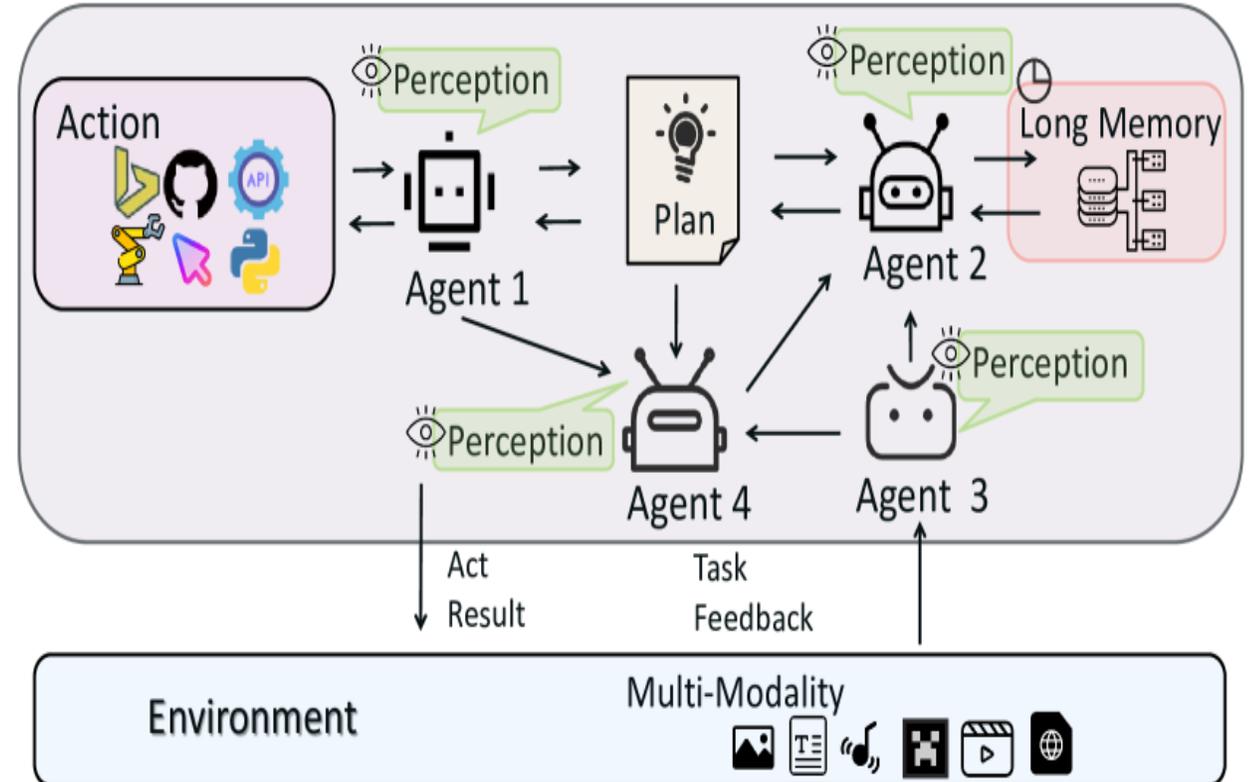


(d)

# Large Multimodal Agents (LMA)



(a)



(b)

# Agentic AI Cloud Architecture

## Microservices and Serverless Architecture

Containers (Docker, Kubernetes)

Serverless platforms (AWS Lambda, Google Cloud Functions)

## APIs and Tooling Integration via MCP

Agents access tools (e.g., databases, APIs, CRMs, payment gateways)  
using Model Context Protocol (MCP)

Enhances tool-using behavior of LLM agents

## Tools and Frameworks

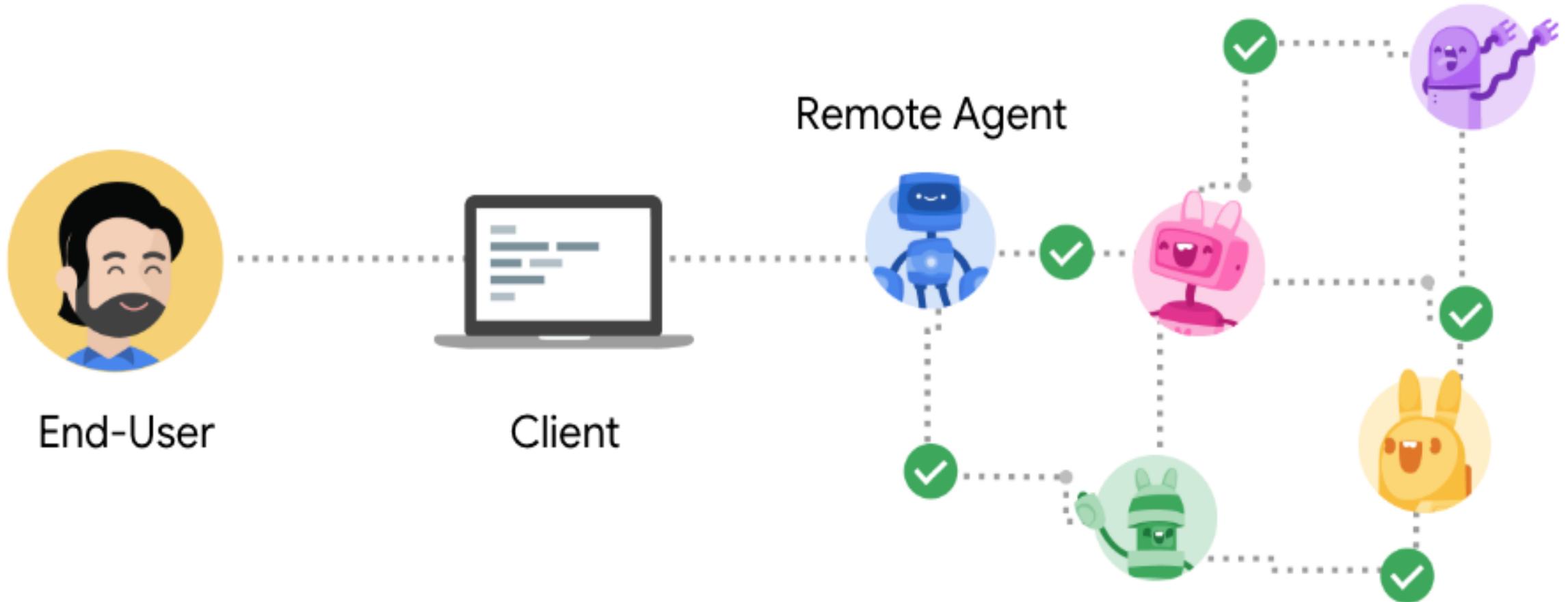
LangChain, AutoGen, CrewAI: for orchestrating LLM agents

Anthropic's MCP, Google's A2A: communication protocols

Vector DBs (Pinecone, Weaviate): for agent memory

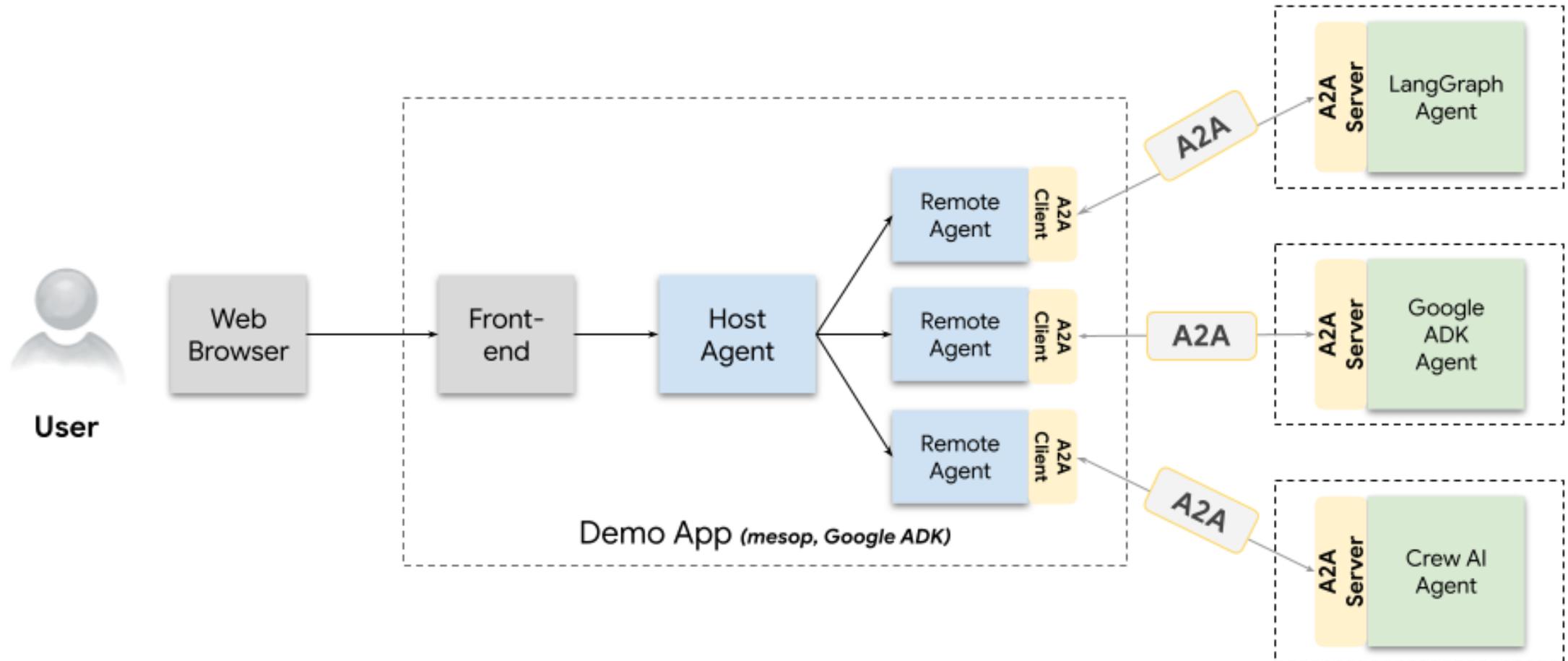
# Agent2Agent Protocol (A2A)

An open protocol enabling Agent-to-Agent interoperability, bridging the gap between opaque agentic systems



# A2A Demo Web App

Agents talking to other agents over A2A



# A2A

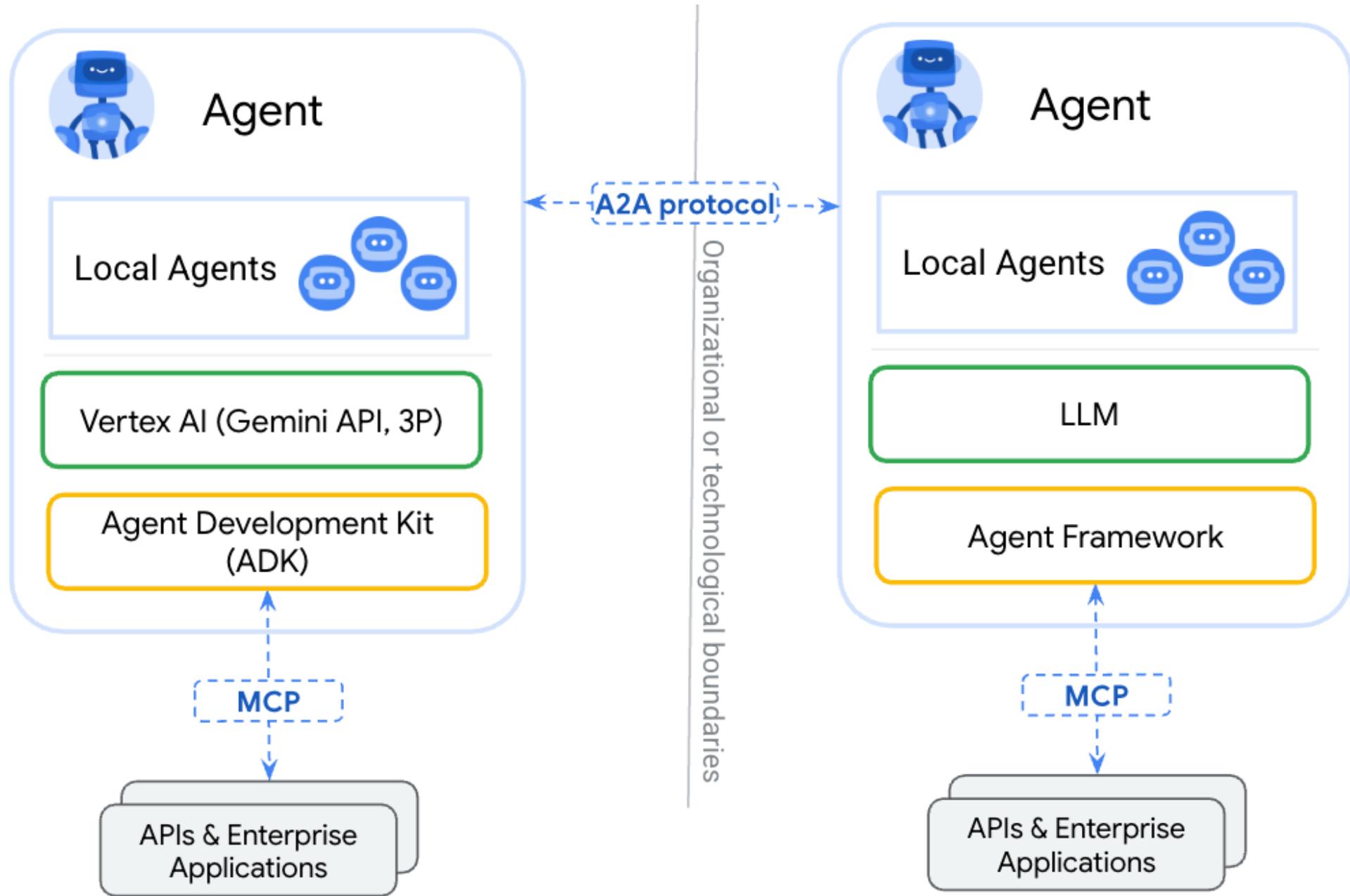
(Agent2Agent Protocol)

for agent-agent collaboration

# MCP

(Model Context Protocol)

for tools and resources

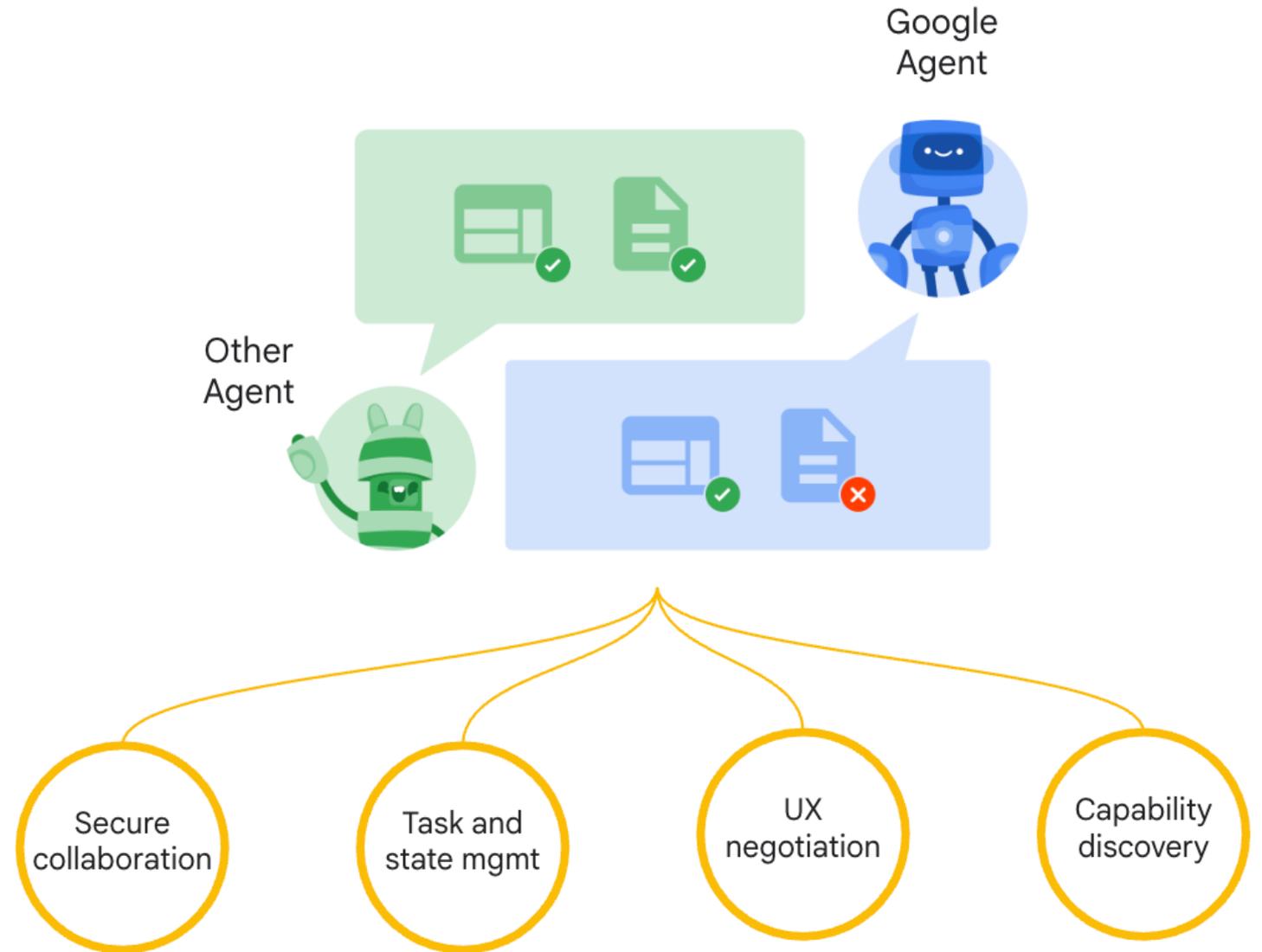


# Google A2A (Agent2Agent Protocol)

Seamless Agent Collaboration

Simplifies Enterprise Agent Integration

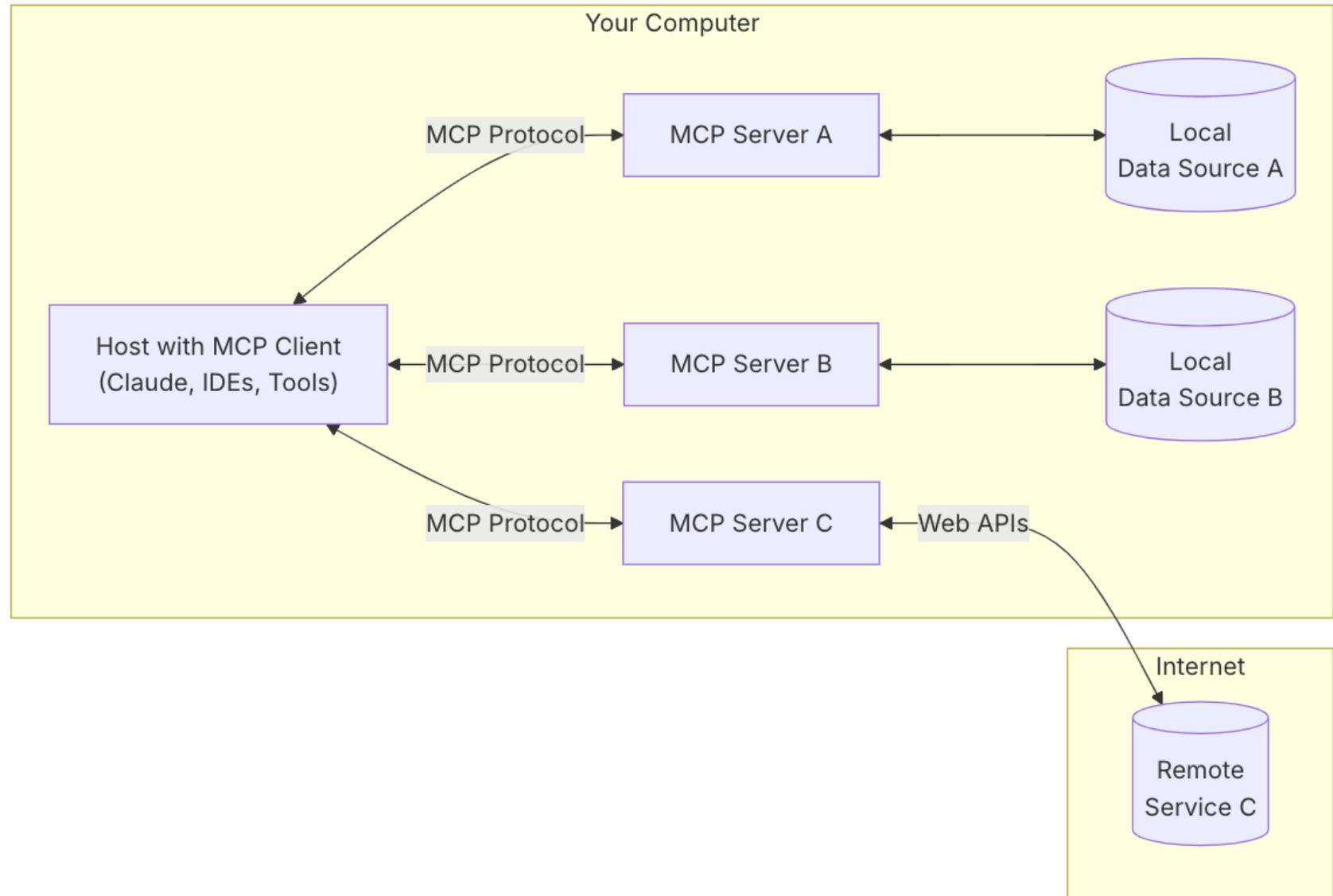
Supports Key Enterprise Requirements



# MCP (Model Context Protocol)

**MCP is an open protocol that standardizes how applications provide context to LLMs.**

**MCP: USB-C port for AI applications.**



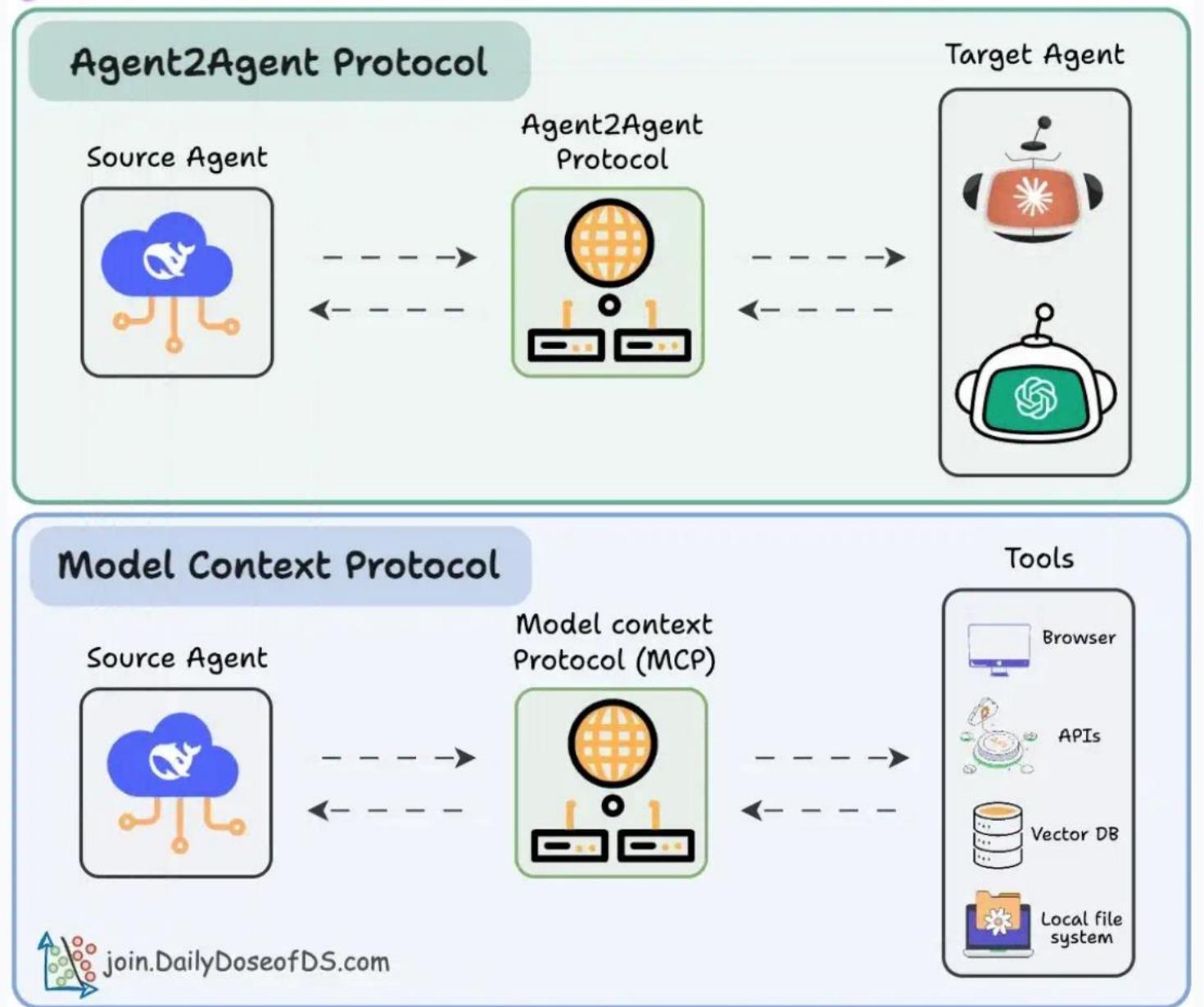
# MCP and A2A

- **MCP (Model Context Protocol) for tools and resources**
  - **Connect agents to tools, APIs, and resources with structured inputs/outputs.**
  - **Google ADK supports MCP tools. Enabling wide range of MCP servers to be used with agents.**
- **A2A (Agent2Agent Protocol) for agent-agent collaboration**
  - **Dynamic, multimodal communication between different agents without sharing memory, resources, and tools**
  - **Open standard driven by community.**
  - **Samples available using Google ADK, LangGraph, Crew.AI**

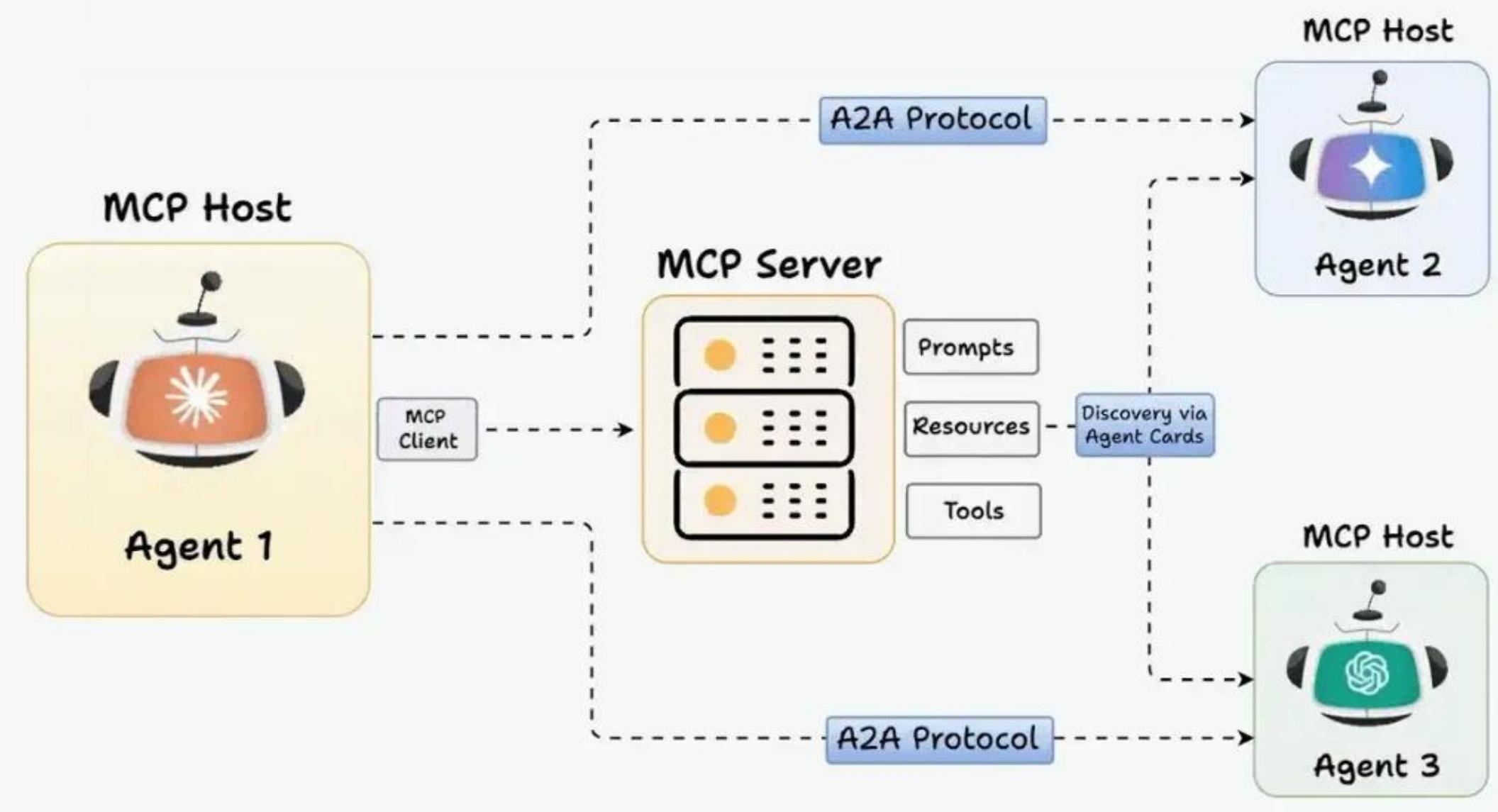
# Agentic applications require both A2A and MCP

A2A allows agents to connect with other agents and collaborate in teams.

MCP provides agents with access to tools

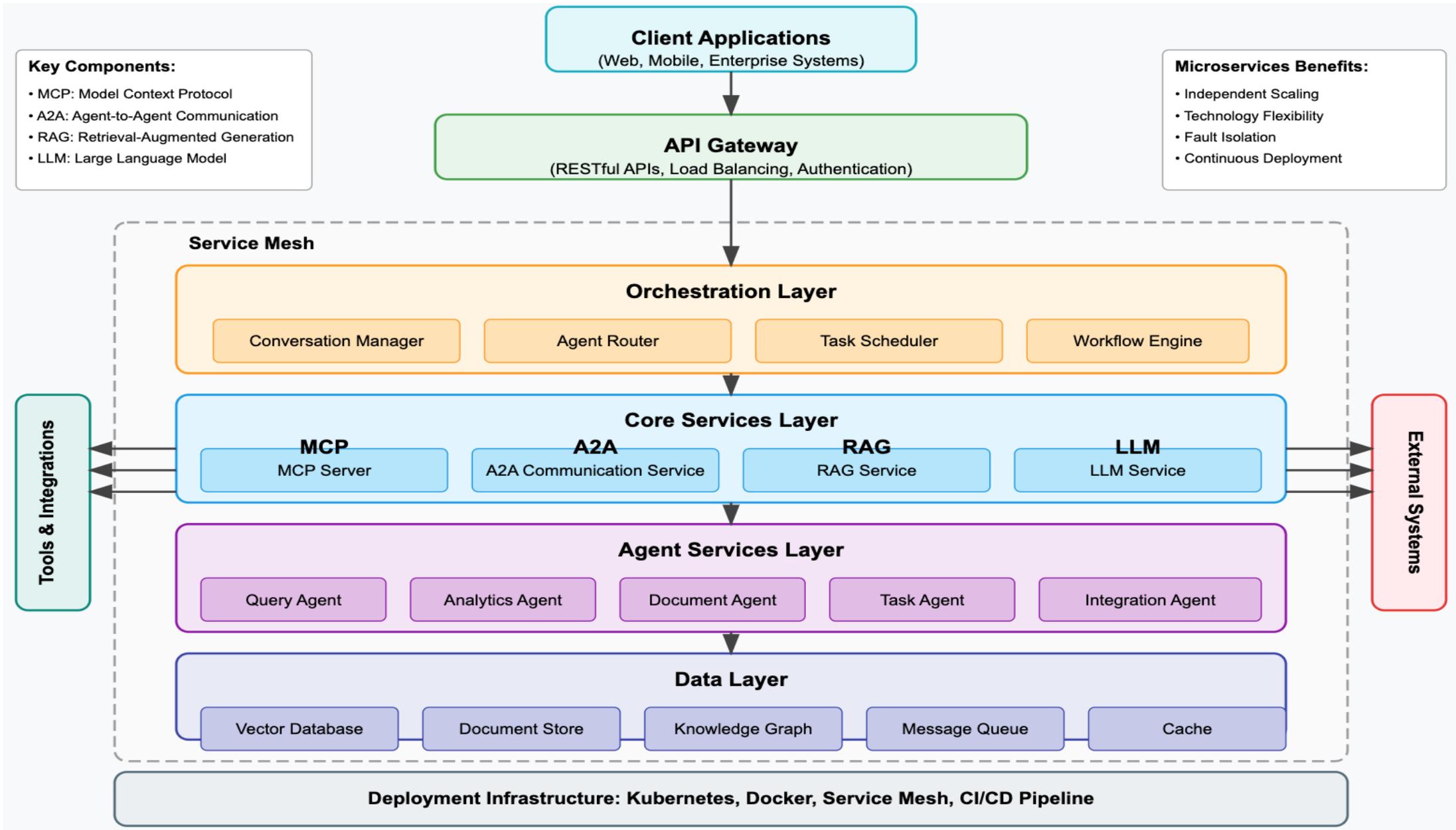


# MCP and A2A Protocol for AI Agents

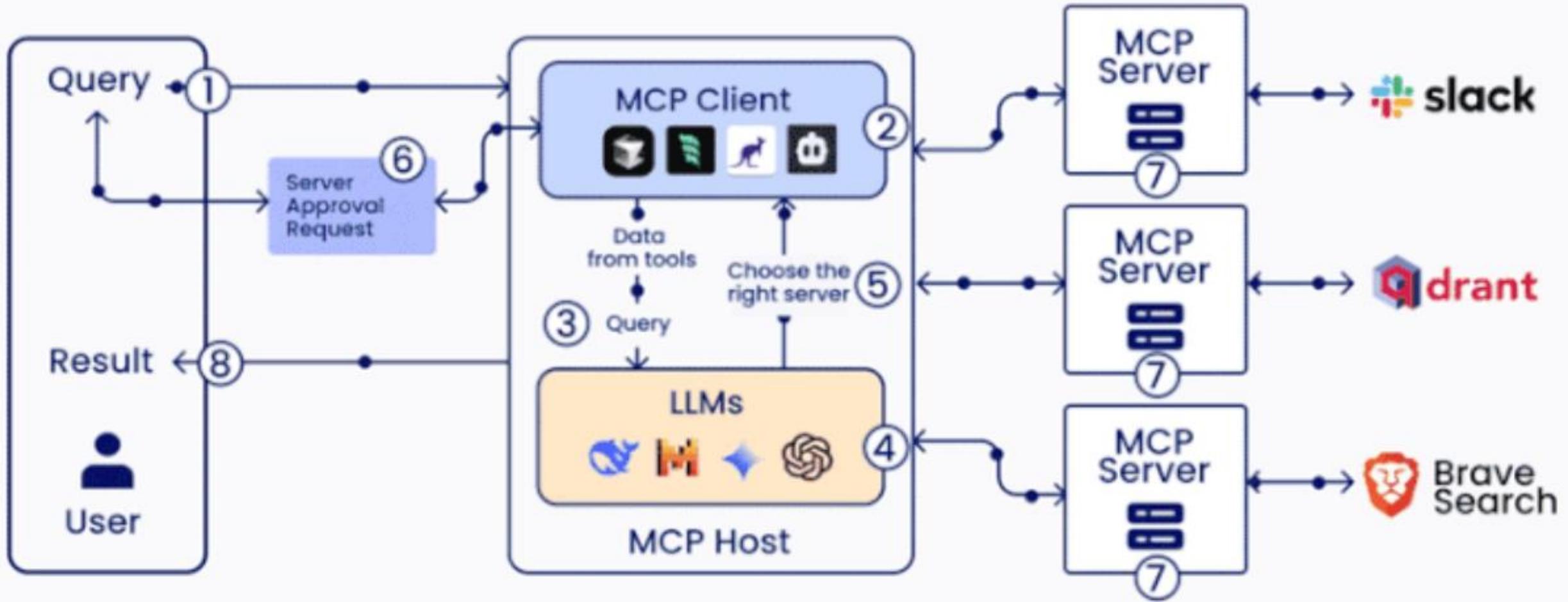


Source: <https://blog.dailydoseofds.com/p/a-visual-guide-to-agent2agent-a2a>

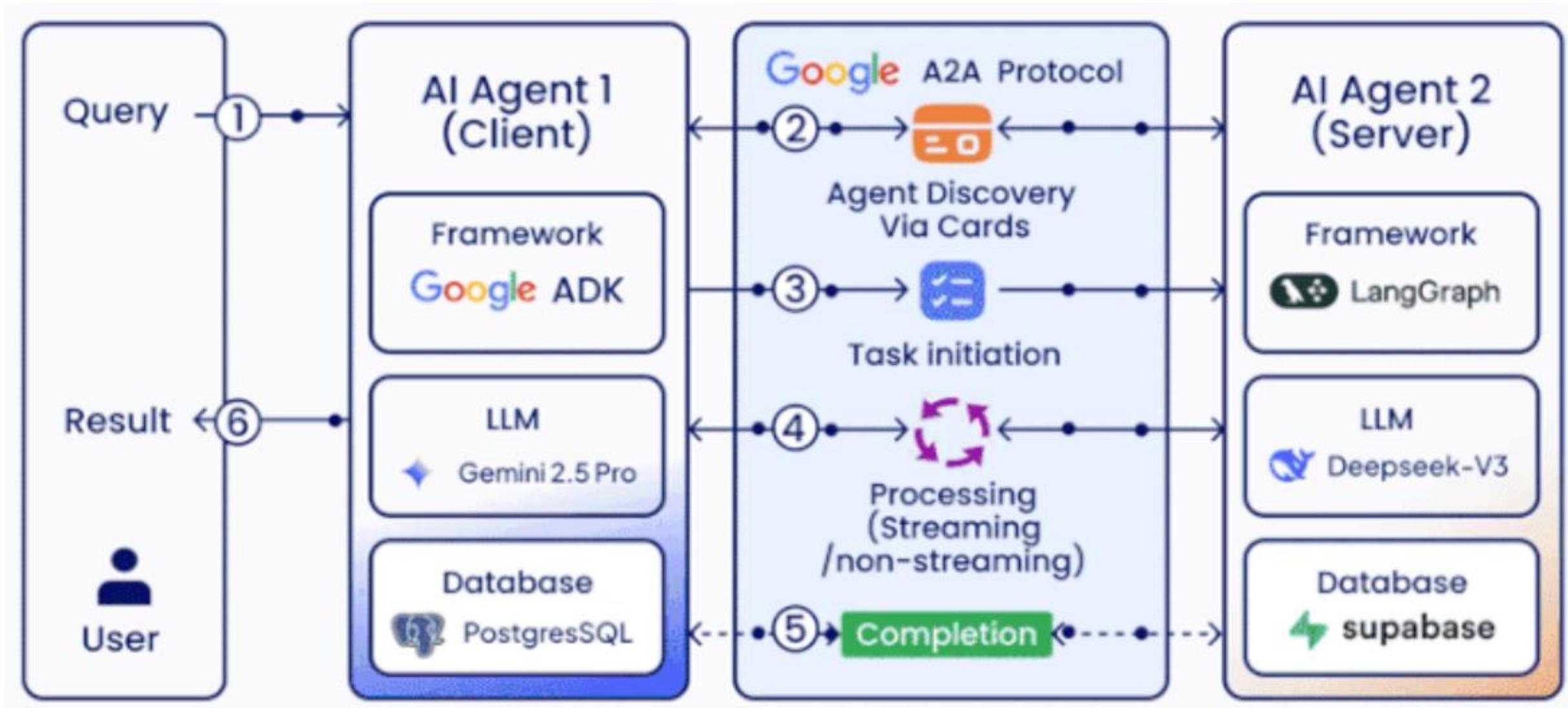
# Agentic AI System with Microservices Architecture



# MCP (Model Context Protocol)



# A2A (Agent2Agent Protocol)



**Generative AI**  
**Large Language Models**  
**(LLMs)**  
**Foundation Models**

# Language Models

# Text Image

# Speech

# Video

# Models

- Text To Image
- Speech To Text
- Text To Speech
- Speech To Speech
- Video Generation

## Text To Image

Image generation models and API providers

ALL MODELS	IMAGE ARENA
 METHODOLOGY	 DALLE
 Stable Diffusion	 Midjourney
 Playground	 Amazon Titan
 Ideogram	 Google Imagen
 Leonardo.Ai Phoenix	 Recraft
 Janus Pro	 Luma Labs
 Infinity	 MiniMax
 Gemini	 OpenAI GPT
 Reve	 FLUX
 SANA-Sprint	 HiDream

# 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

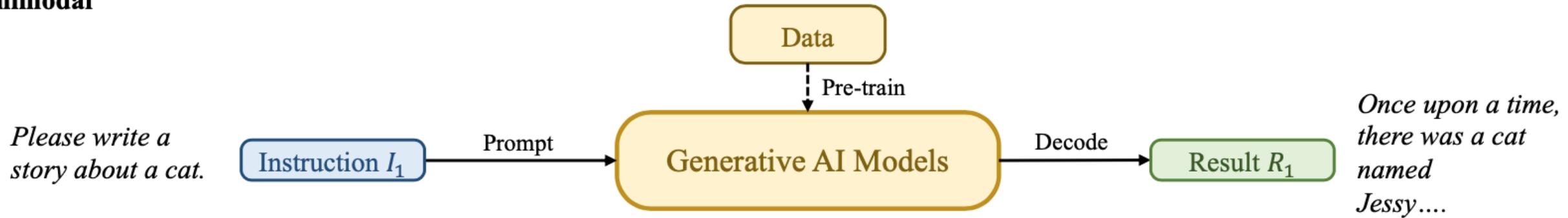


 **OpenAI DALL·E 2**

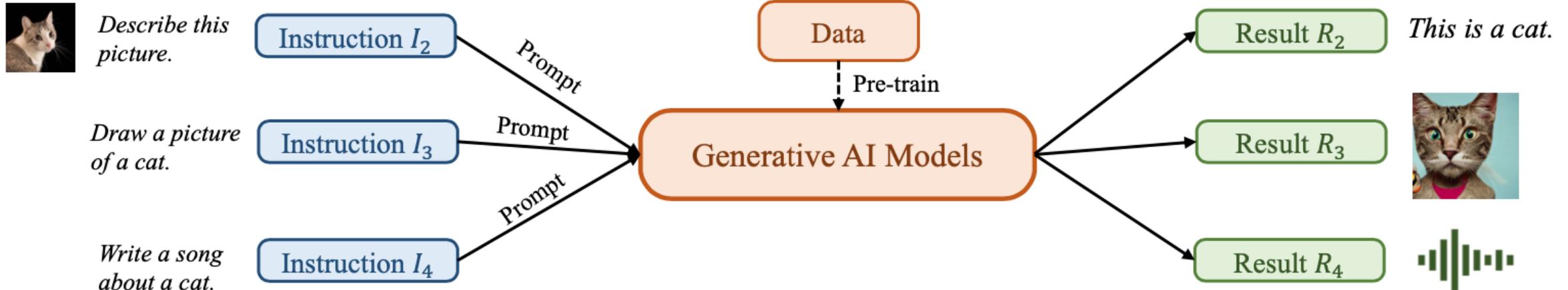
# Generative AI (Gen AI)

## AI Generated Content (AIGC)

### Unimodal

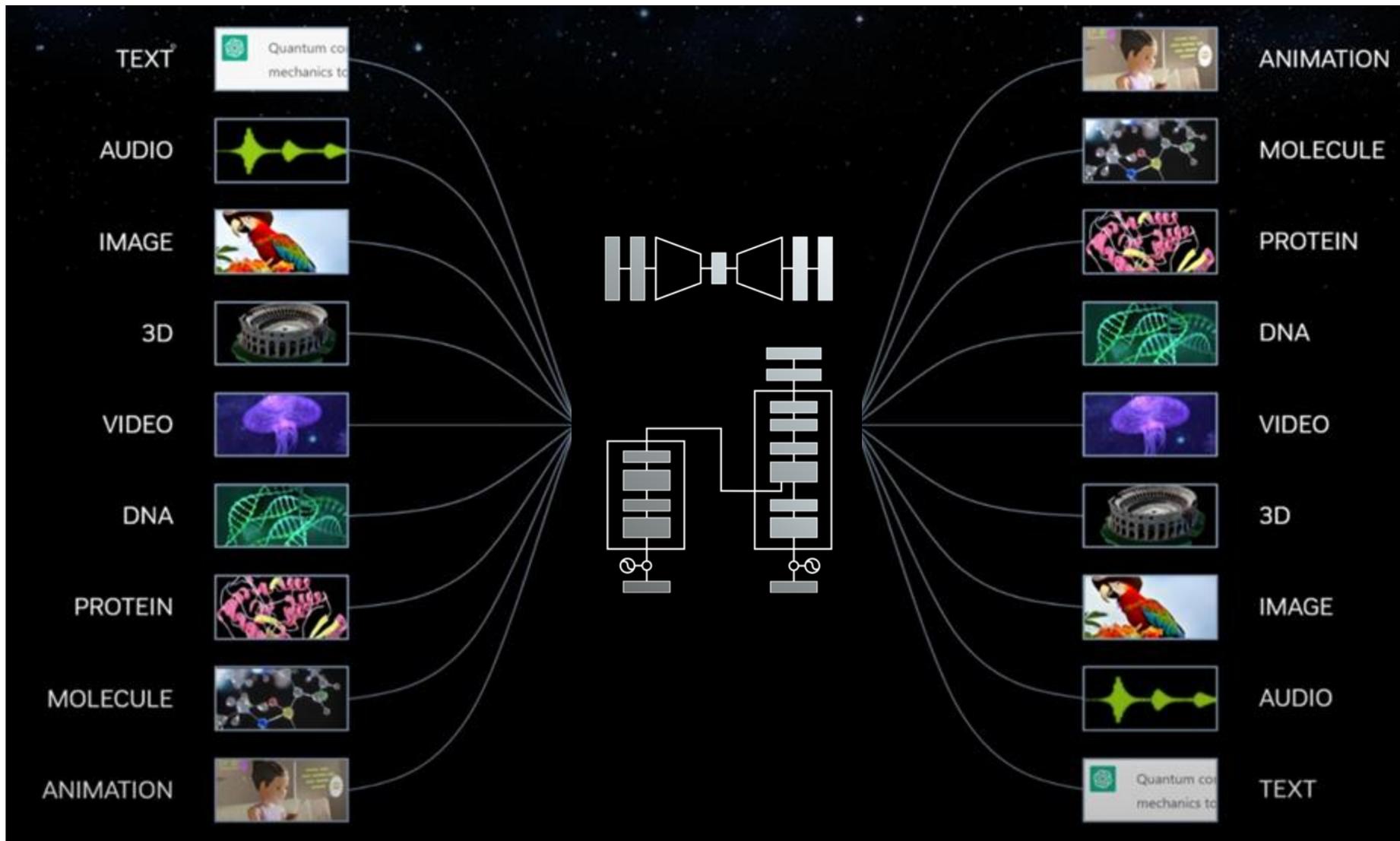


### Multimodal



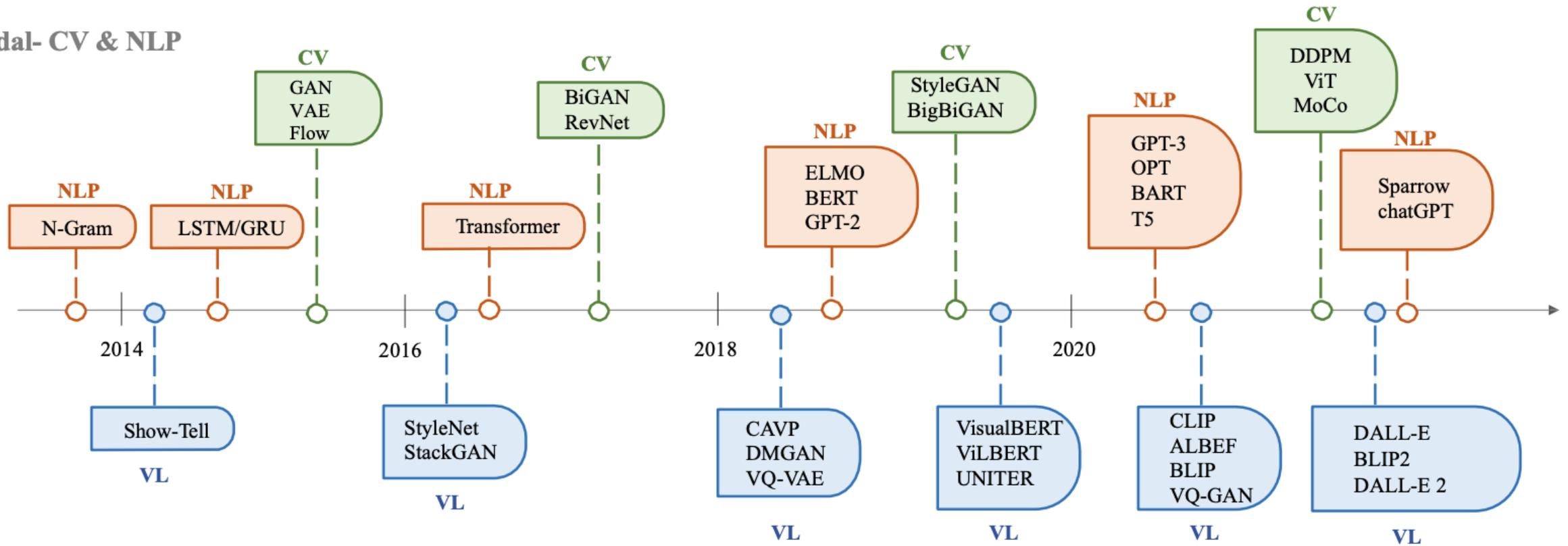
# Modular Modalities

## Where Can The Transformer Fit?



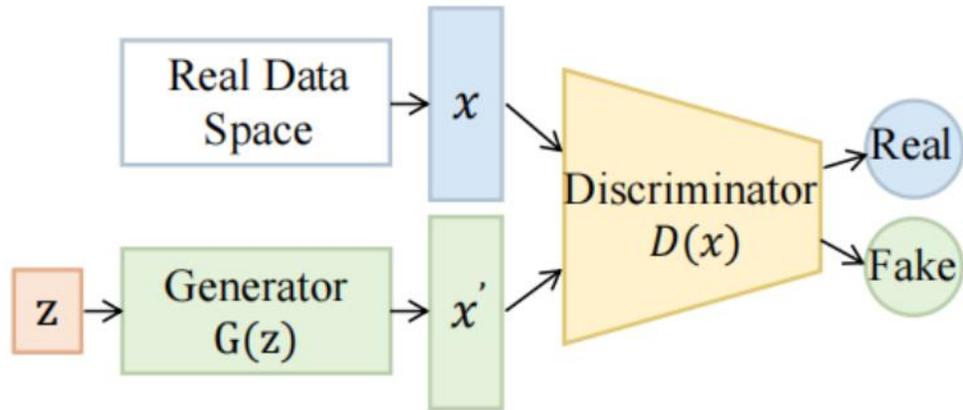
# The history of Generative AI in CV, NLP and VL

## Unimodal- CV & NLP

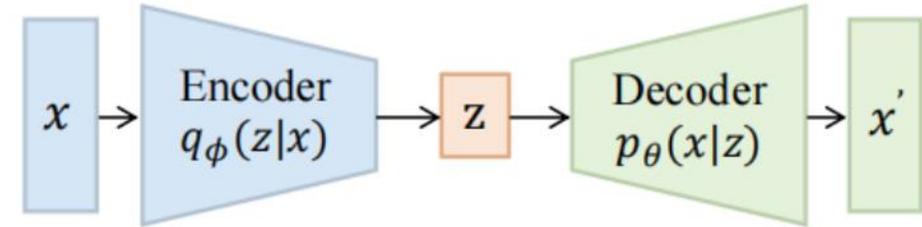


## Multimodal – Vision Language

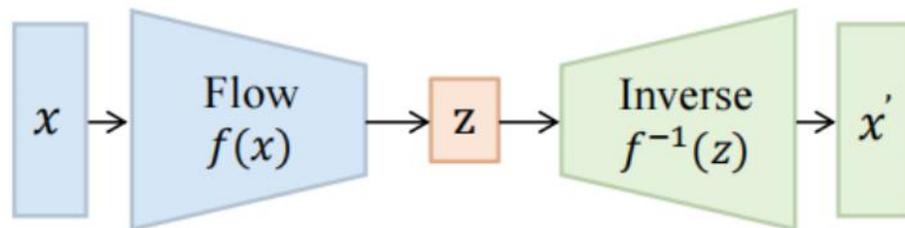
# Categories of Vision Generative Models



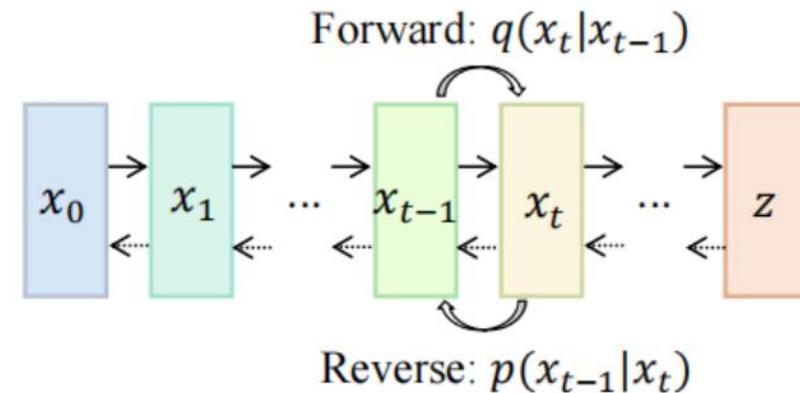
(1) Generative adversarial networks



(2) Variational autoencoders

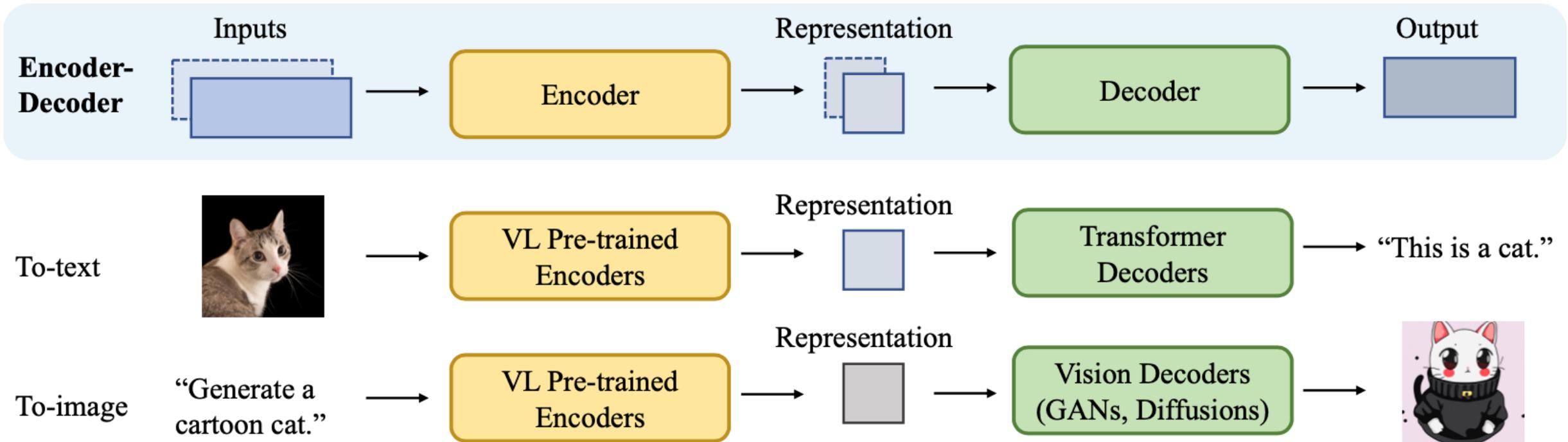


(3) Normalizing flows



(4) Diffusion models

# The General Structure of Generative Vision Language



# Artificial Analysis **Text to Image Arena**

Artificial Analysis LANGUAGE MODELS ▾ SPEECH, IMAGE & VIDEO MODELS ▾ LEADERBOARDS ▾ ARENAS ▾ ABOUT ▾ Newsletter **Subscribe**

Arena    Leaderboard    Personal Leaderboard

TEXT TO IMAGE ARENA 🚩 + Submit prompt

13/30 to view your model preferences 🙄 Try the new 🗣️ Speech Arena

**Which image best reflects this prompt?**

Clumsy robot trying to cook in a cartoon kitchen



♥ Prefer (← Key)



♥ Prefer (→ Key)

# Artificial Analysis **Text to Speech Arena**

Arena

Leaderboard

Personal Leaderboard

TEXT TO SPEECH ARENA 🚩

5/30 to view your model preferences 🧠

API Performance & Price Analysis 📄

### Which do you prefer?

Imagine this voice as a conversational AI assistant, customer support system or reading you an email

The ISS travels at approximately 17,500 miles per hour, orbiting Earth every 90 minutes and experiencing 16 sunrises and sunsets each day.

▶ 0:11 / 0:11



▶ 0:10 / 0:10



Playing

♥ Prefer (← Key)

♥ Prefer (→ Key)

#### Notes:

**Models compared:** TTS-1, TTS-1 HD, Studio, Journey, Neural2, WaveNet, Standard, Polly Long-Form, Polly Neural, Polly Standard, Azure Neural, MetaVoice v1, XTTS v2, StyleTTS 2, OpenVoice v2, Sonic English (Oct '24), Turbo v2.5, Multilingual v2, GPT-4o Realtime Preview, 3.0 mini, T2A-01-HD, T2A-01-Turbo, Zonos-v0.1, Kokoro 82M v1.0, Polly Generative, Flash v2.5, Fish Speech 1.5, Dialog, GPT-4o mini TTS, LMNT

**Methodology:** For further details, see our [Speech to Text methodology page](#).

**Other notable links:** See also [TTS-Arena](#) on Hugging Face for another arena which includes more open-source models.

# Artificial Analysis Video Generation Model Arena

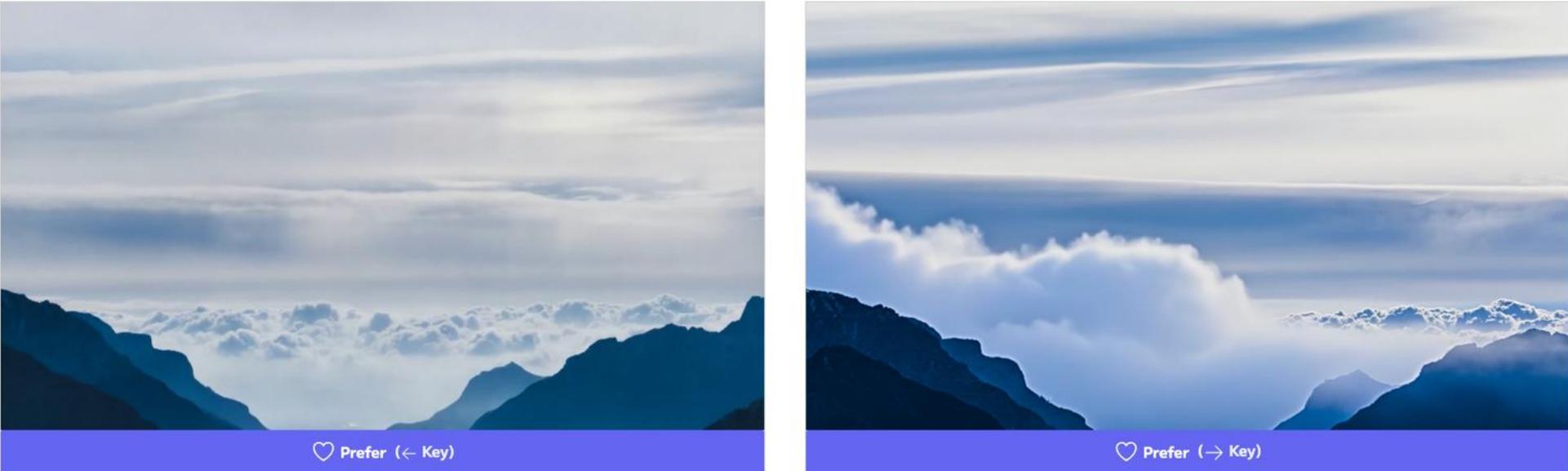
Artificial Analysis LANGUAGE MODELS SPEECH, IMAGE & VIDEO MODELS LEADERBOARDS ARENAS ABOUT Newsletter **Subscribe**

Arena Leaderboard Personal Leaderboard

VIDEO GENERATION MODEL ARENA  0/30 to view your model preferences  + Submit prompt  
Try the new  Speech Arena

Which video best reflects this prompt?

Clouds flow gently through the mountain valley, billowing and expanding as they move from left to right.



 Prefer (← Key)

 Prefer (→ Key)

# Artificial Analysis **Text to Image Leaderboard**

## Text to Image AI Model & Provider Leaderboard

Analysis and comparison of Text to Image generation models & API providers. Artificial Analysis has analyzed text to image models and hosting providers across quality, generation time, and price. For further details, see our methodology page.

**Image Arena**  
 Contribute to the Quality ELO score and see your personal model ranking

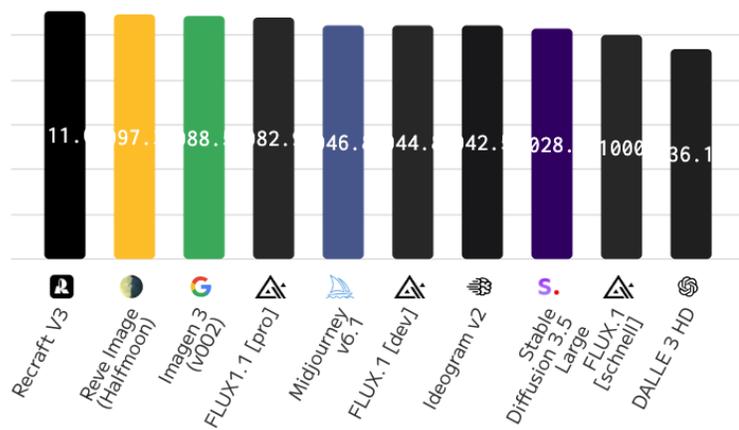
[Image Arena](#)

**Text to image models & providers compared:** Phoenix 0.9 Ultra, Playground v2.5, Stable Diffusion 3 Medium, Stable Diffusion XL 1.0, SDXL Lightning, Stable Diffusion 1.5, Stable Diffusion 2.1, Amazon Titan G1 (Standard), DALLE 2, DALLE 3 HD, DALLE 3, Midjourney v6, Stable Diffusion 1.6, Stable Diffusion 3 Large Turbo, Stable Diffusion 3 Large, Midjourney v6.1, Amazon Titan G1 v2 (Standard), Playground v3 (beta), Ideogram v2, FLUX.1 [pro], FLUX.1 [dev], Stable Diffusion 3.5 Medium, Ideogram v2 Turbo, Ideogram v1, FLUX1.1 [pro], Recraft 20B, FLUX.1 [schnell], Stable Diffusion 3.5 Large, Stable Diffusion 3.5 Large Turbo, Recraft V3, Luma Photon Flash, Adobe Firefly 3, GPT-4o, Janus Pro, Luma Photon, Lumina Image v2, Phoenix 1.0 Fast, Phoenix 1.0 Ultra, Image-01, Gemini 2.0 Flash Experimental, Reve Image (Halfmoon), Ideogram v2a, Ideogram v2a Turbo, Imagen 3 (v002), Ideogram 3.0, Midjourney v7 Alpha, Sana Sprint 1.6B, HiDream-I1-Dev, and Grok 2.

### Highlights

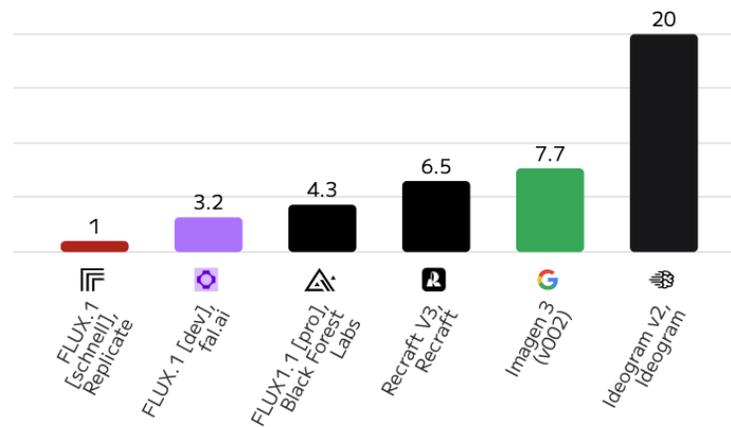
#### QUALITY ELO

ELO score in Artificial Analysis Image Arena (relative metric of image generation quality), Higher is better



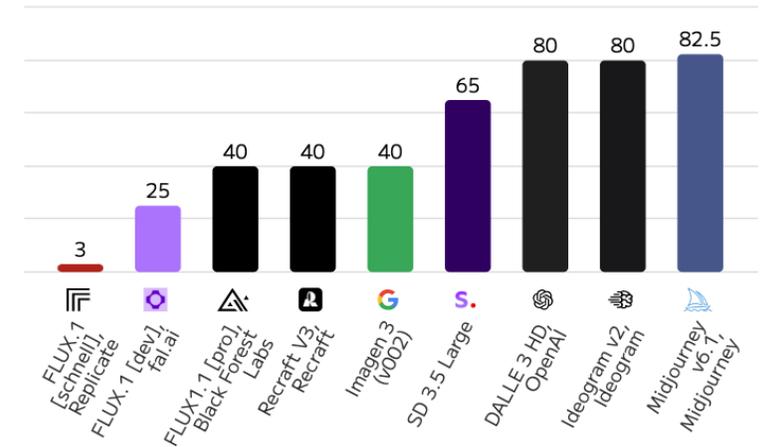
#### GENERATION TIME

Generation time: Seconds to generate 1 image, Lower is better



#### PRICE

Price: USD per 1000 image generations, Lower is better

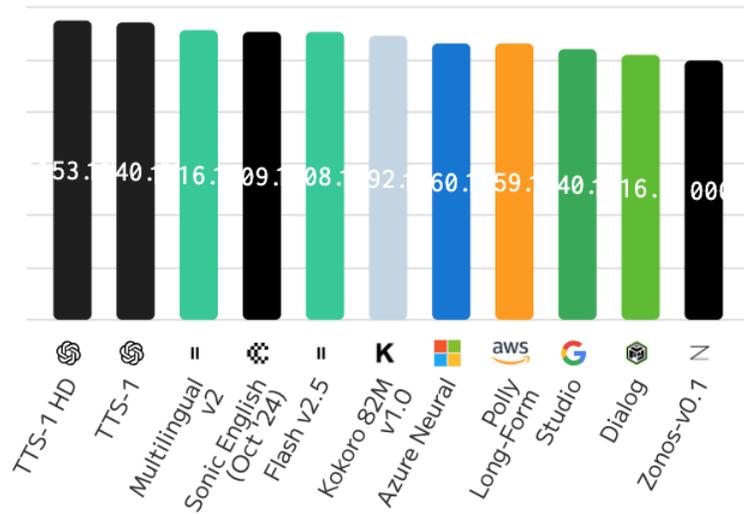


# Text to Speech (TTS) AI Model & Provider Leaderboard

Text to speech models & providers compared: TTS-1, TTS-1 HD, Studio, Journey, Neural2, WaveNet, Standard, Polly Long-Form, Polly Neural, Polly Standard, Azure Neural, MetaVoice v1, XTTS v2, StyleTTS 2, OpenVoice v2, Sonic English (Oct '24), 3.0 mini, Turbo v2.5, Multilingual v2, T2A-01-HD, T2A-01-Turbo, Zonos-v0.1, Kokoro 82M v1.0, Polly Generative, Flash v2.5, Dialog, Murf Speech Gen 2, and Step TTS Mini.

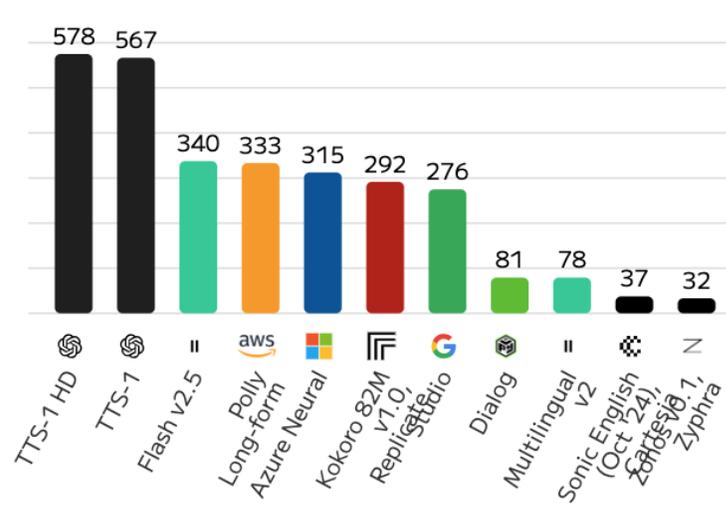
## QUALITY ELO

Arena ELO: Average ELO rating of the model, Higher is better



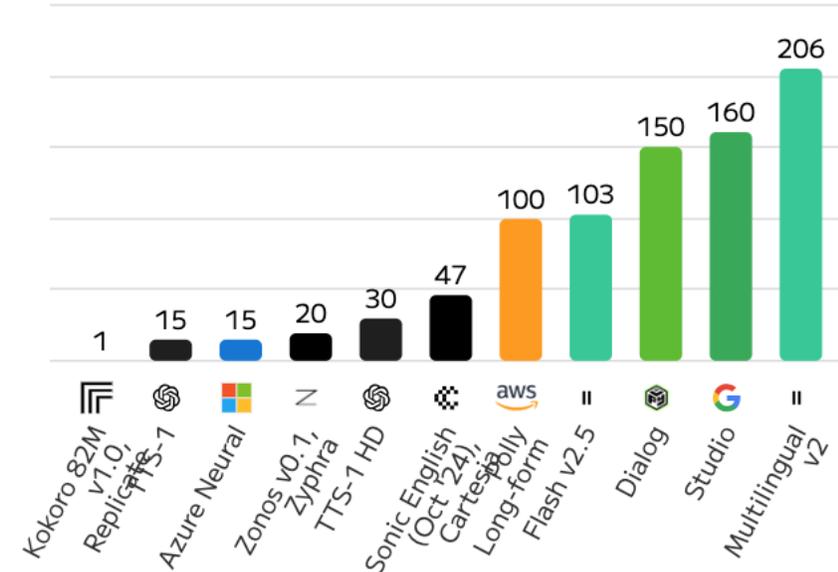
## CHARACTERS PER SECOND

Characters processed per second: # of characters per second of generation time, Higher is better



## PRICE

Price: USD per 1M characters of text, Lower is better



# Text to Speech (TTS) AI Model & Provider Leaderboard

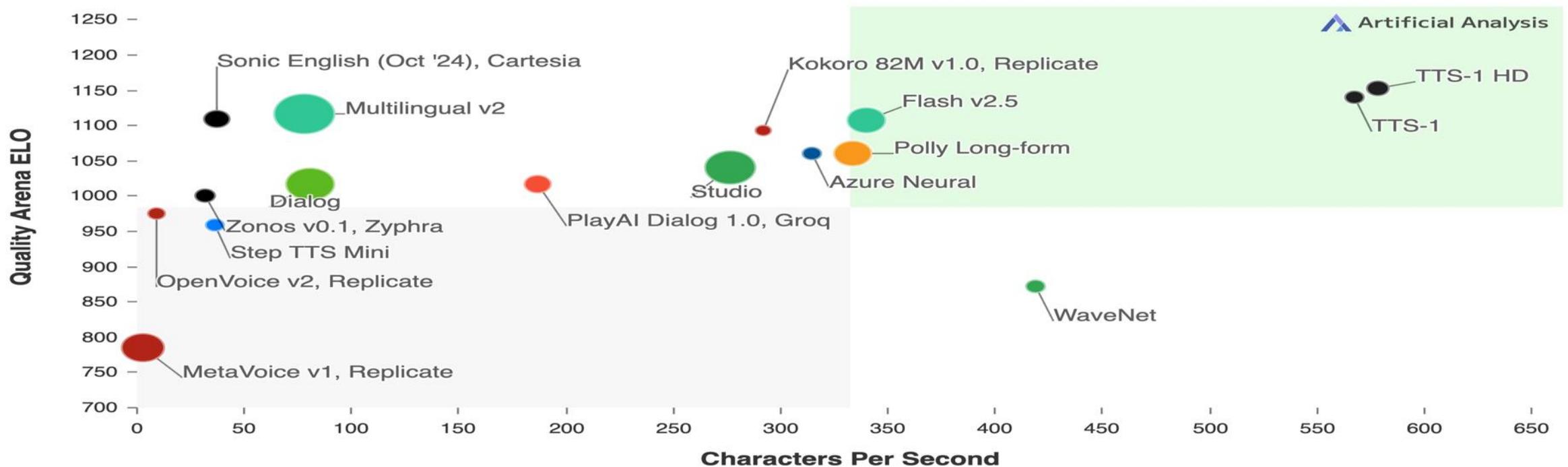
## Quality vs. Speed

Arena ELO: Average ELO rating of the model, Characters processed per second: # of characters per second of generation time

Most attractive quadrant

Size represents Price: USD per 1M characters of text

- TTS-1
- TTS-1 HD
- Studio
- WaveNet
- Polly Long-form
- Azure Neural
- MetaVoice v1, Replicate
- OpenVoice v2, Replicate
- Sonic English (Oct '24), Cartesia
- Multilingual v2
- Zonos v0.1, Zyphra
- Kokoro 82M v1.0, Replicate
- Flash v2.5
- Dialog
- Step TTS Mini
- PlayAI Dialog 1.0, Groq



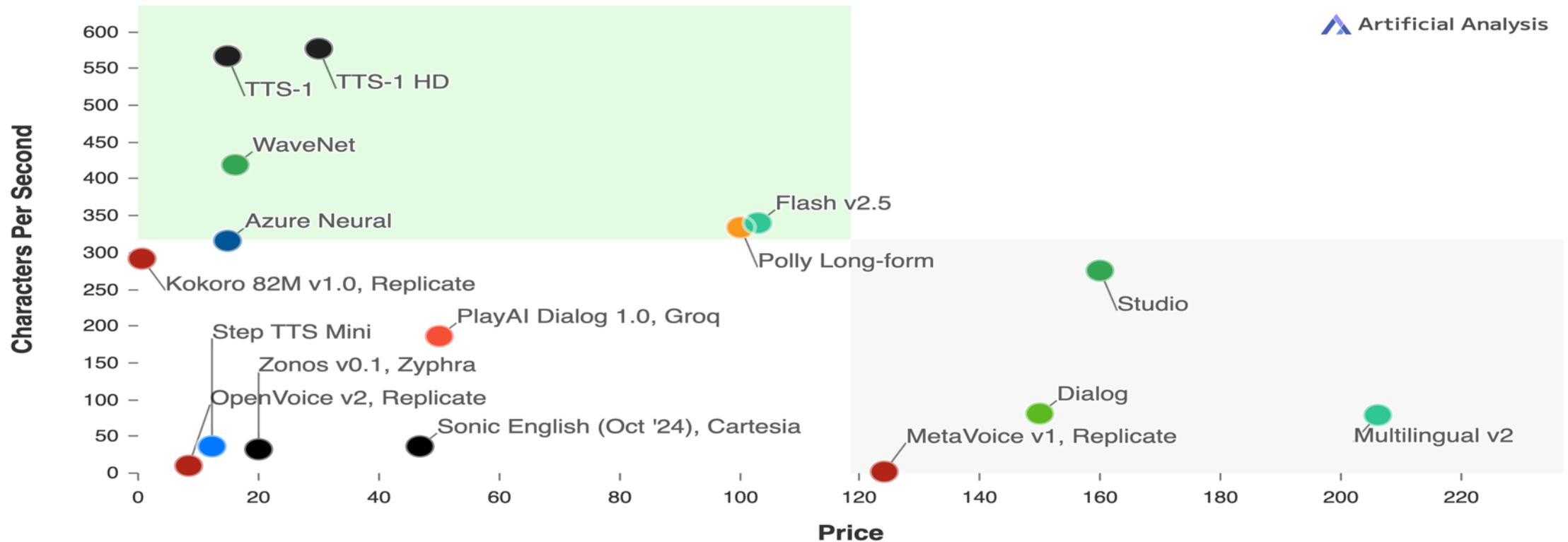
# Text to Speech (TTS) AI Model & Provider Leaderboard

## Speed vs. Price

Characters processed per second: # of characters per second of generation time, Price: USD per 1M characters of text

Most attractive quadrant

- TTS-1
- TTS-1 HD
- Studio
- WaveNet
- Polly Long-form
- Azure Neural
- MetaVoice v1, Replicate
- OpenVoice v2, Replicate
- Sonic English (Oct '24), Cartesia
- Multilingual v2
- Zonos v0.1, Zyphra
- Kokoro 82M v1.0, Replicate
- Flash v2.5
- Dialog
- Step TTS Mini
- PlayAI Dialog 1.0, Groq



# Text to Speech (TTS) AI Model & Provider Leaderboard

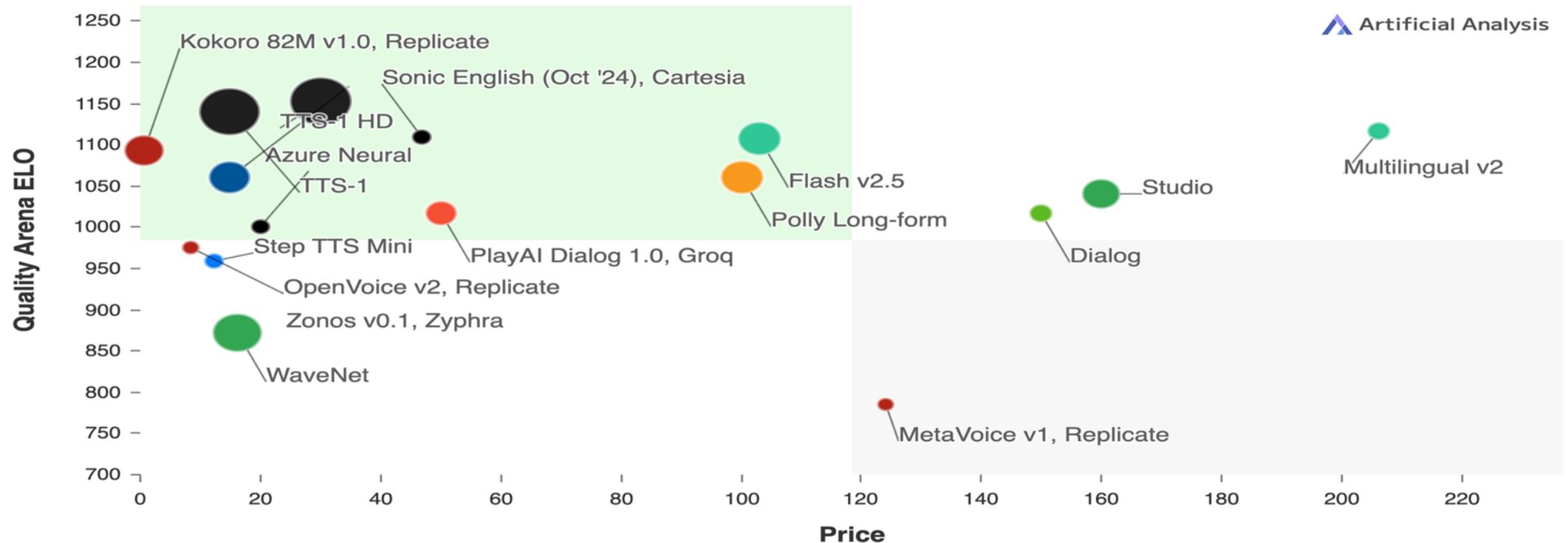
## Quality vs. Price

Arena ELO: Average ELO rating of the model, Price: USD per 1M characters of text

Most attractive quadrant

Size represents Characters processed per second: # of characters per second of generation time

- TTS-1
- TTS-1 HD
- Studio
- WaveNet
- Polly Long-form
- Azure Neural
- MetaVoice v1, Replicate
- OpenVoice v2, Replicate
- Sonic English (Oct '24), Cartesia
- Multilingual v2
- Zonos v0.1, Zyphra
- Kokoro 82M v1.0, Replicate
- Flash v2.5
- Dialog
- Step TTS Mini
- PlayAI Dialog 1.0, Groq

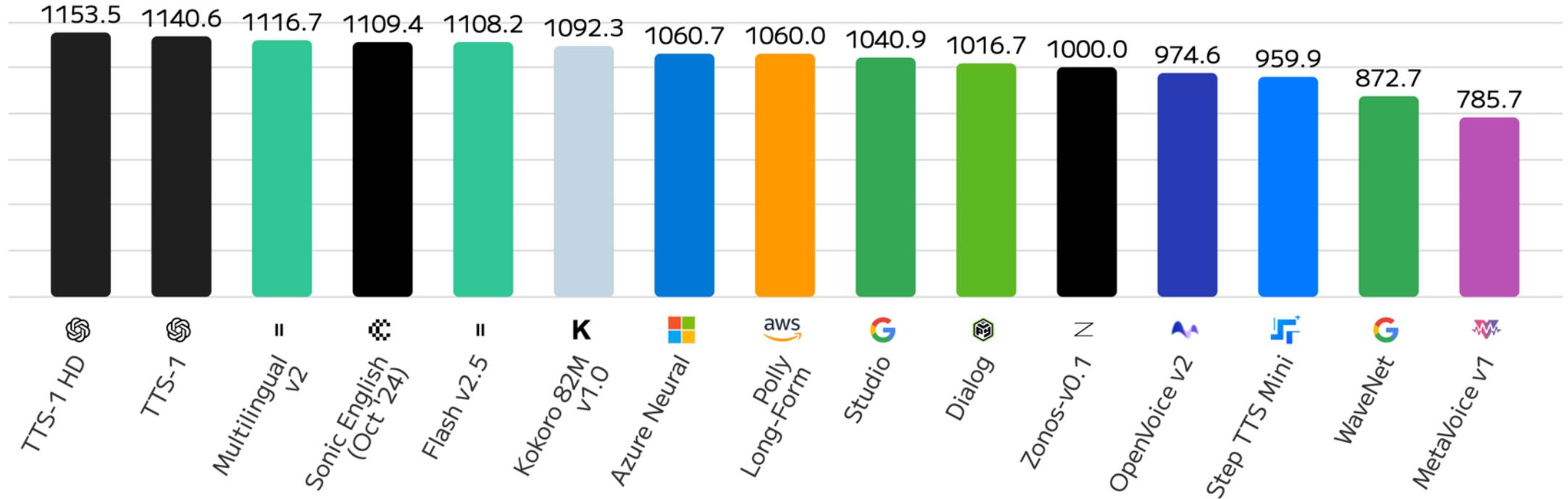


# Text to Speech (TTS) AI Model & Provider Leaderboard

## Quality Arena ELO (Text to Speech Arena)

Arena ELO: Average ELO rating of the model, Higher is better

Artificial Analysis



# Speech to Text (STT) AI Model & Provider Leaderboard

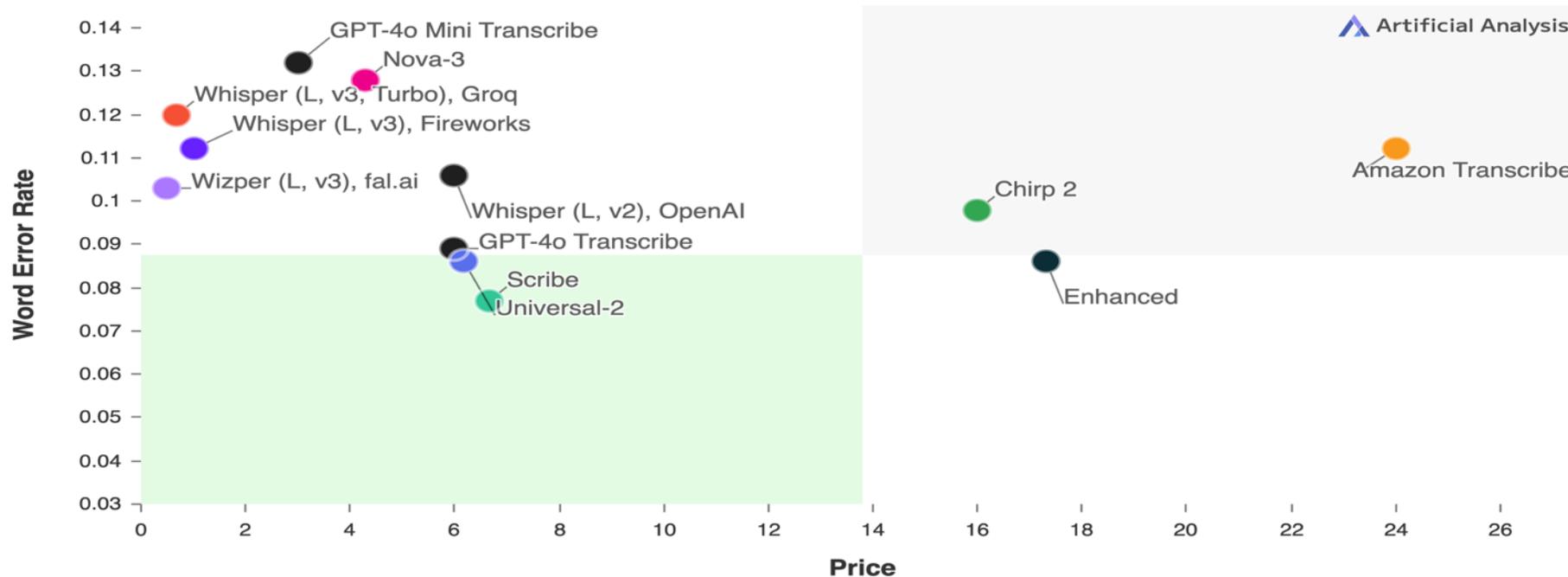
Speech-to-text models & providers compared: Whisper (L, v2), OpenAI, Universal-1, Standard, Whisper (L, v2), Azure, Enhanced, Nano, Wizper (L, v3), fal.ai, Incredibly Fast Whisper, Replicate, Nova-2, Whisper (L, v2), Replicate, Whisper (L, v3), Replicate, Base, WhisperX, Replicate, Whisper (L v2), Deepgram, Gladia, Whisper (L, v3), Groq, Distil-Whisper, Groq, Whisper (L, v3), fal.ai, Whisper (L, v3), Deepinfra, Whisper (L, v3, Turbo), Groq, Whisper (L, v3), Fireworks, Whisper (L, v3, Turbo), Fireworks, Universal-2, Amazon Transcribe, Fish Speech to Text, Nova-3, Chirp, Chirp 2, Scribe, GPT-4o Transcribe, and GPT-4o Mini Transcribe.

## Word Error Rate vs. Price

Word error rate: % of words transcribed incorrectly, Price: USD per 1000 minutes of audio

Most attractive quadrant

- Whisper (L, v2), OpenAI
Enhanced
Wizper (L, v3), fal.ai
Whisper (L, v3, Turbo), Groq
- Whisper (L, v3), Fireworks
Universal-2
Amazon Transcribe
Nova-3
Chirp 2
Scribe
- GPT-4o Transcribe
GPT-4o Mini Transcribe



# Speech to Text (STT) AI Model & Provider Leaderboard

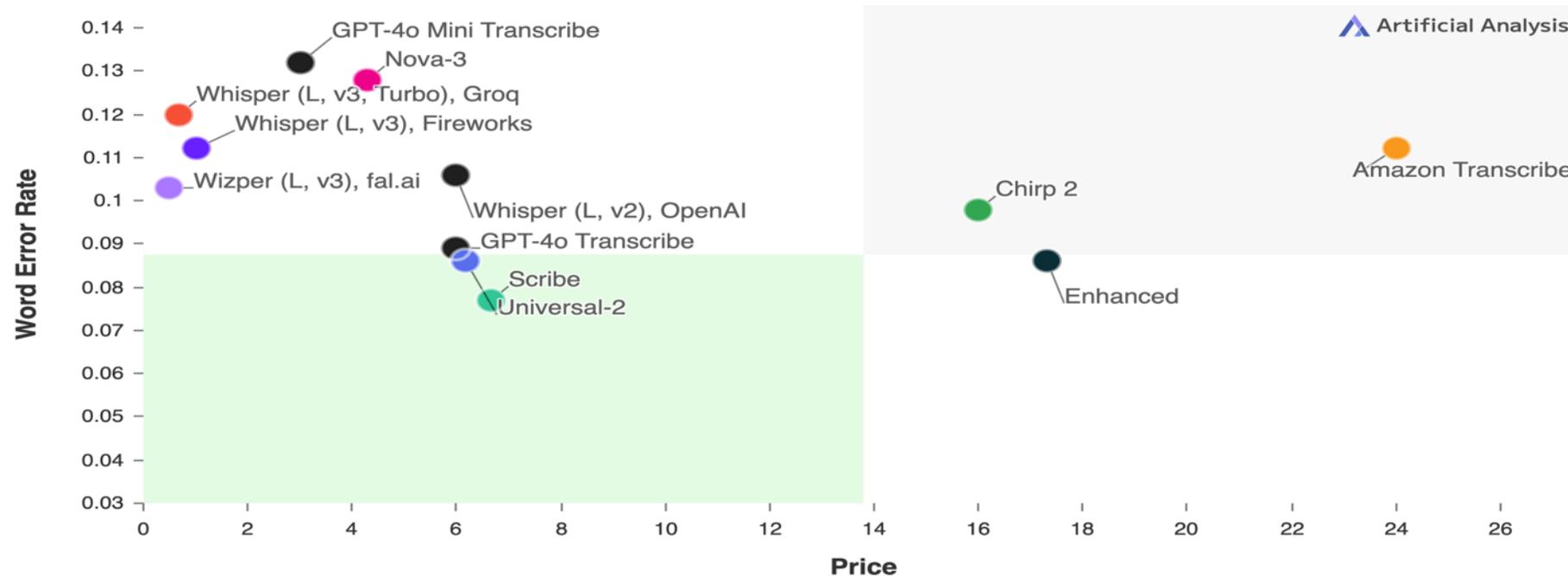
Speech-to-text models & providers compared: Whisper (L, v2), OpenAI, Universal-1, Standard, Whisper (L, v2), Azure, Enhanced, Nano, Wizper (L, v3), fal.ai, Incredibly Fast Whisper, Replicate, Nova-2, Whisper (L, v2), Replicate, Whisper (L, v3), Replicate, Base, WhisperX, Replicate, Whisper (L v2), Deepgram, Gladia, Whisper (L, v3), Groq, Distil-Whisper, Groq, Whisper (L, v3), fal.ai, Whisper (L, v3), Deepinfra, Whisper (L, v3, Turbo), Groq, Whisper (L, v3), Fireworks, Whisper (L, v3, Turbo), Fireworks, Universal-2, Amazon Transcribe, Fish Speech to Text, Nova-3, Chirp, Chirp 2, Scribe, GPT-4o Transcribe, and GPT-4o Mini Transcribe.

## Word Error Rate vs. Price

Word error rate: % of words transcribed incorrectly, Price: USD per 1000 minutes of audio

Most attractive quadrant

- Whisper (L, v2), OpenAI
- Enhanced
- Wizper (L, v3), fal.ai
- Whisper (L, v3, Turbo), Groq
- Whisper (L, v3), Fireworks
- Universal-2
- Amazon Transcribe
- Nova-3
- Chirp 2
- Scribe
- GPT-4o Transcribe
- GPT-4o Mini Transcribe



# Artificial Analysis **Text to Video** Leaderboard

Text to Video		Image to Video		
CREATOR	NAME	ARENA ELO	95% CI	# APPEARANCES
 Google	<b>Veo 2</b>	1124	-10/+10	6,452
 Kuaishou	<b>Kling 1.5 (Pro)</b>	1053	-6/+6	20,631
 OpenAI	<b>OpenAI Sora</b>	1049	-5/+5	23,649
 MiniMax	<b>T2V-01</b>	1039	-4/+4	43,450
 Pika Art	<b>Pika 2.0</b>	1038	-6/+6	20,432
 Kuaishou	<b>Kling 1.6 (Standard)</b>	1029	-7/+6	13,607
 MiniMax	<b>T2V-01-Director</b>	1022	-9/+9	7,765

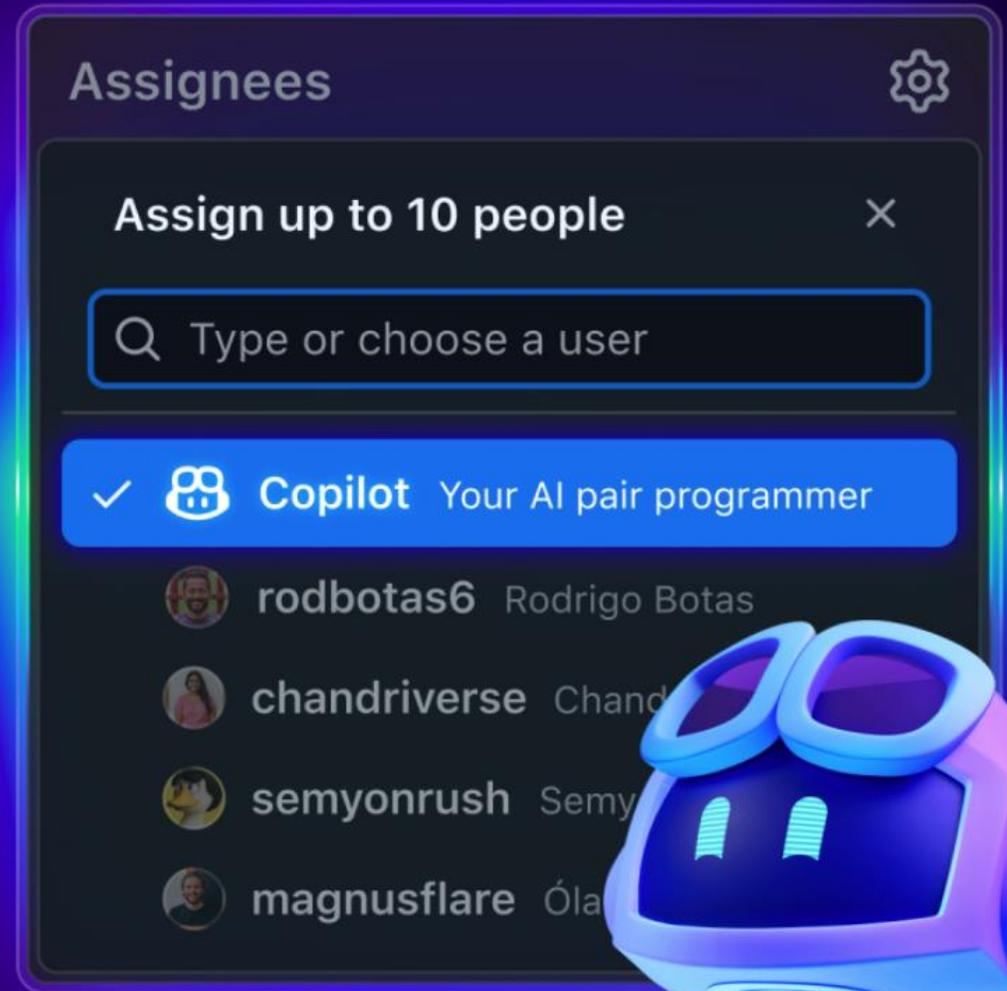
# Artificial Analysis **Image to Video** Leaderboard

Text to Video		Image to Video		
CREATOR	NAME	ARENA ELO	95% CI	# APPEARANCES
 Kuaishou	<b>Kling 1.6 (Pro)</b>	1121	-17/+18	2,748
 Runway	<b>Runway Gen 4</b>	1115	-14/+15	7,314
 Google	<b>Veo 2</b>	1113	-18/+17	2,770
 MiniMax	<b>I2V-01-Director</b>	1031	-15/+15	7,407
 Pika Art	<b>Pika 2.2</b>	1001	-19/+17	2,740
 Alibaba	<b>Wan 2.1 14B</b>	1000	+0/+0	2,700
 Runway	<b>Runway Gen 3 Alpha Turbo</b>	992	-15/+14	7,420
 Runway	<b>Runway Gen 3 Alpha</b>	971	-18/+16	2,558
 OpenAI	<b>OpenAI Sora</b>	960	-19/+18	2,552
 Tencent	<b>Hunyuan Video</b>	922	-18/+17	2,535

Source: [https://artificialanalysis.ai/text-to-video/arena?tab=Leaderboard&leaderboard\\_tab=t2v](https://artificialanalysis.ai/text-to-video/arena?tab=Leaderboard&leaderboard_tab=t2v)

# GitHub Copilot CODING AGENT

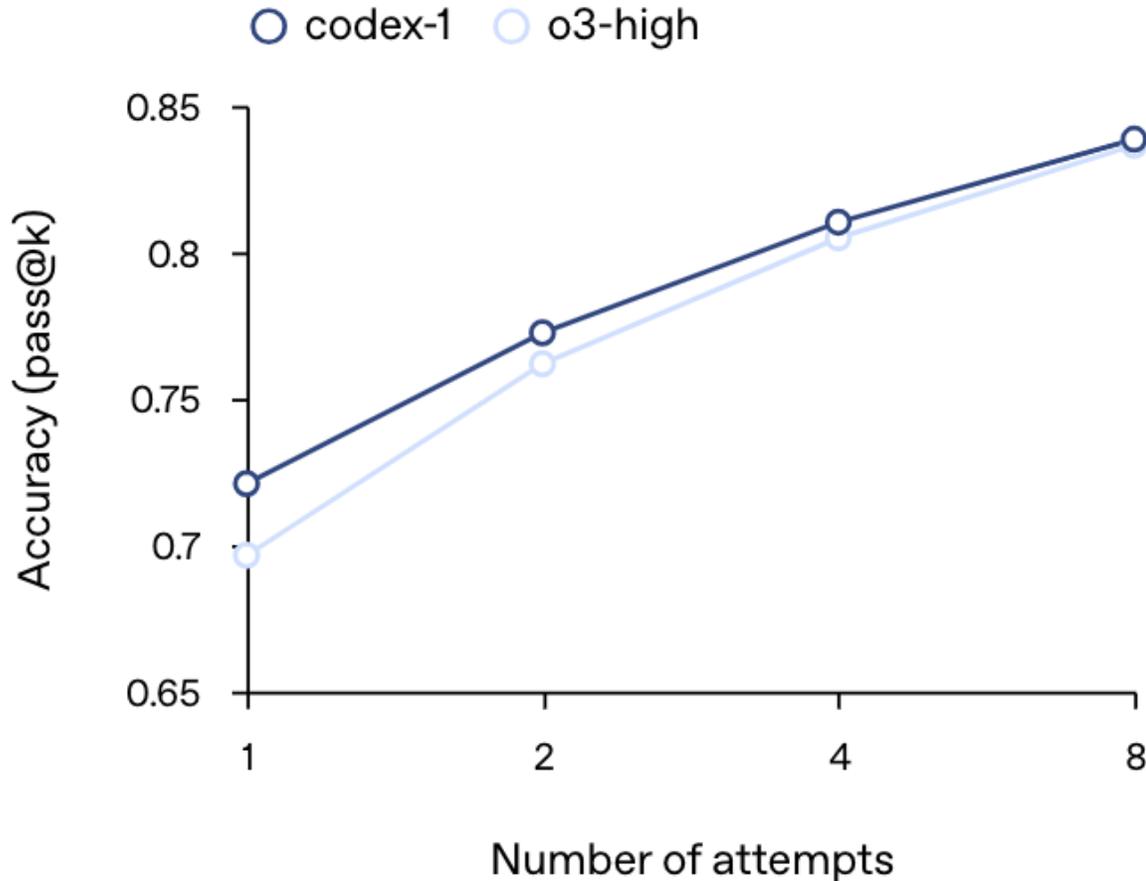
 **GitHub Copilot**  
**CODING**  
**AGENT**



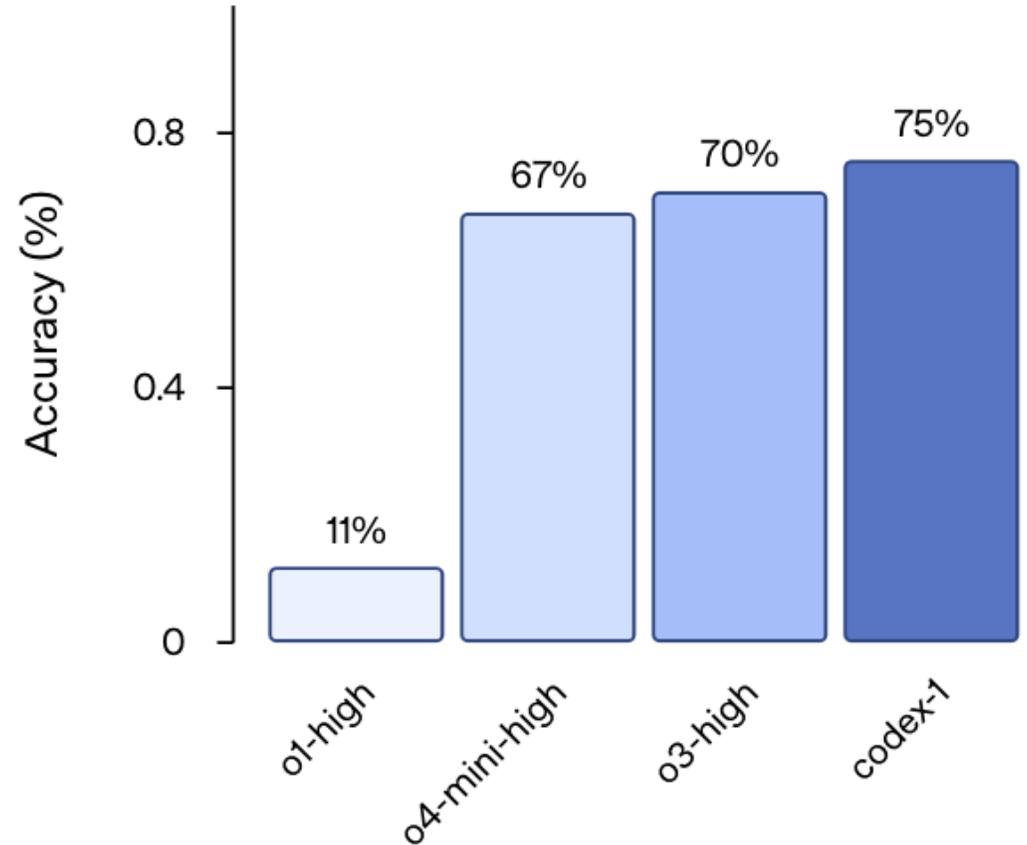
# OpenAI Codex

A cloud-based software engineering agent that can work on many tasks in parallel, powered by codex-1

SWE-Bench Verified



OpenAI Internal SWE tasks



# SWE-bench Verified Leaderboard

AI coding agents ranking by software engineering problems resolved rate

Model	% Resolved	Org	Date
  OpenHands	65.80		2025-04-15
Augment Agent v0	65.40		2025-03-16
 Amazon Q Developer Agent (v20250405-dev)	65.40		2025-04-05
W&B Programmer 01 crosscheck5	64.60		2025-01-17
 PatchPilot-v1.1	64.60		2025-05-03
AgentScope	63.40		2025-02-06
Tools + Claude 3.7 Sonnet (2025-02-24)	63.20		2025-02-24
Blackbox AI Agent	62.80	-	2025-01-10
EPAM AI/Run Developer Agent v20250219 + Anthopic Claude 3.5 Sonnet	62.80		2025-02-28
 SWE-agent + Claude 3.7 Sonnet w/ Review Heavy	62.40		2025-02-25

Source: <https://www.swebench.com/>

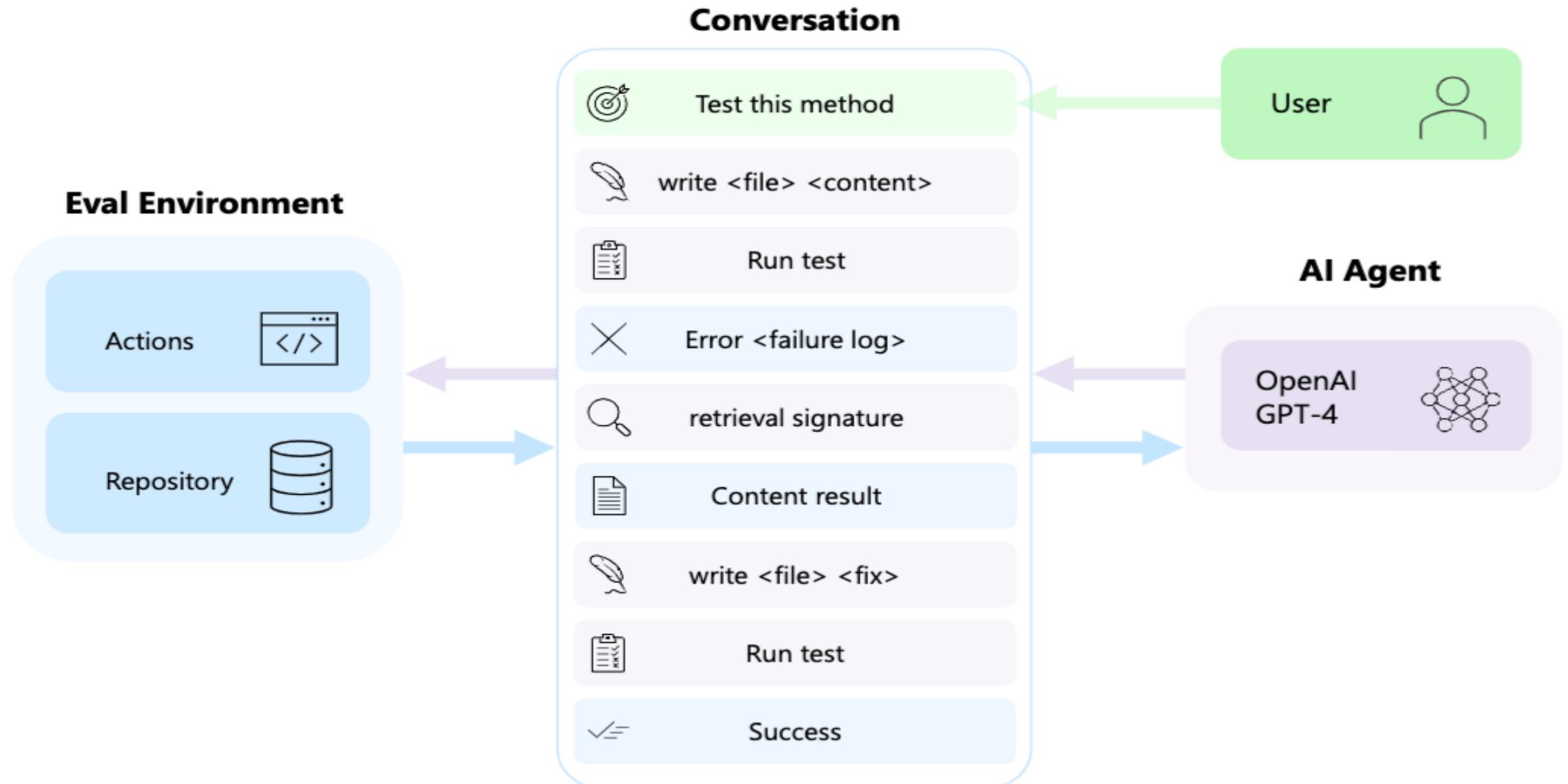
# WebDev Arena Leaderboard

## AI coding competition in web development challenges

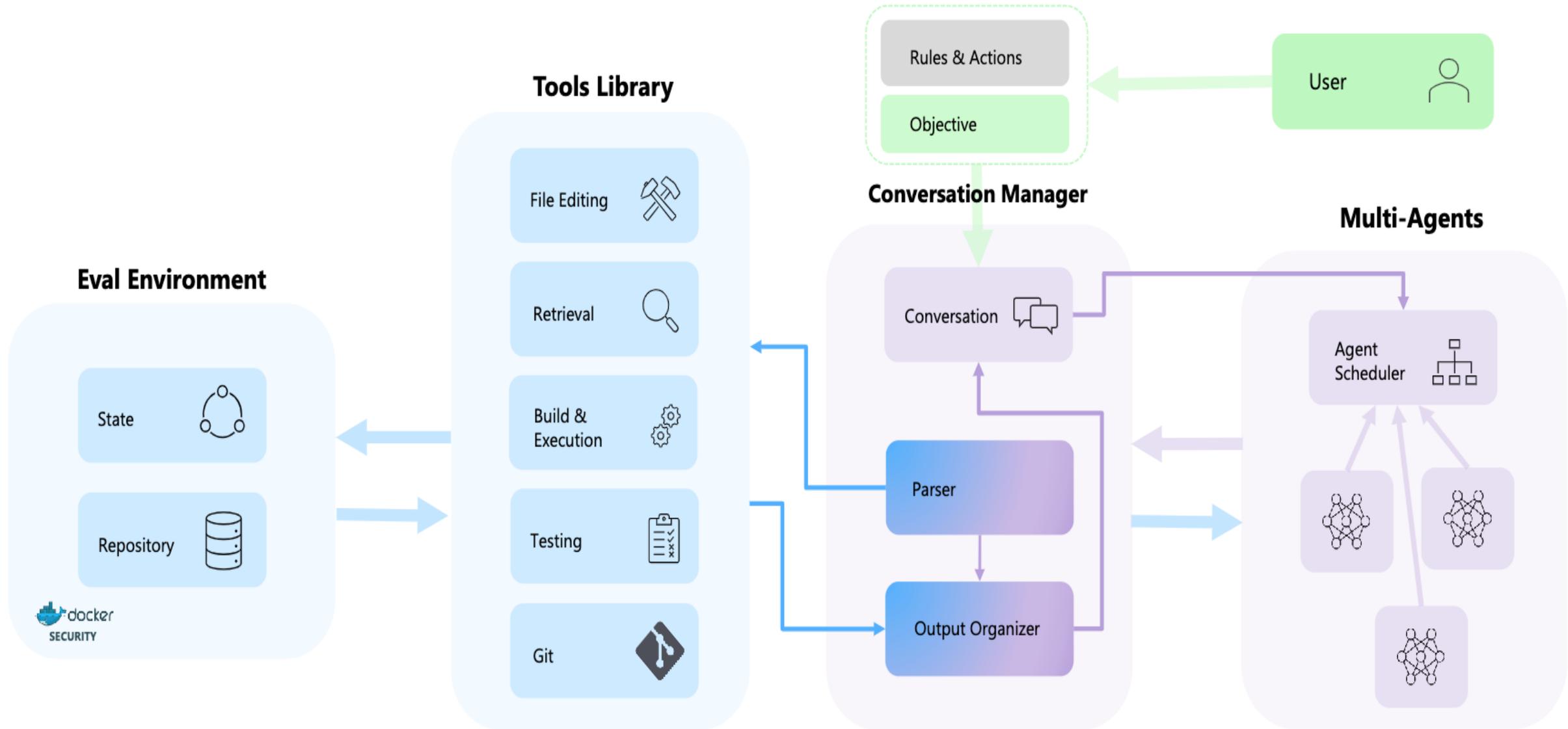
Rank (UB) ⓘ	Model	Arena Score	95% CI ⓘ	Votes	Organization	License
1	 Gemini-2.5-Pro-Preview-05-06	1414.64	+13.57 / -15.14	3,464	Google	Proprietary
2	 Claude 3.7 Sonnet (20250219)	1357.05	+9.03 / -6.72	7,481	Anthropic	Proprietary
3	 Gemini-2.5-Flash-Preview-05-20	1310.42	+19.10 / -21.23	981	Google	Proprietary
4	 GPT-4.1-2025-04-14	1257.20	+9.49 / -8.26	4,880	OpenAI	Proprietary
5	 Claude 3.5 Sonnet (20241022)	1237.74	+4.15 / -4.66	26,338	Anthropic	Proprietary
6	 DeepSeek-V3-0324	1206.67	+20.99 / -20.94	1,097	DeepSeek	MIT
6	 DeepSeek-R1	1198.68	+10.46 / -8.71	3,760	DeepSeek	MIT
6	 o3-2025-04-16	1190.47	+10.43 / -9.42	3,617	OpenAI	Proprietary
6	 GPT-4.1-mini-2025-04-14	1185.05	+10.34 / -10.54	2,995	OpenAI	Proprietary
6	 Qwen3-235B-A22B	1177.78	+13.42 / -16.64	2,024	Alibaba	Apache 2.0

Source: <https://web.lmarena.ai/leaderboard>

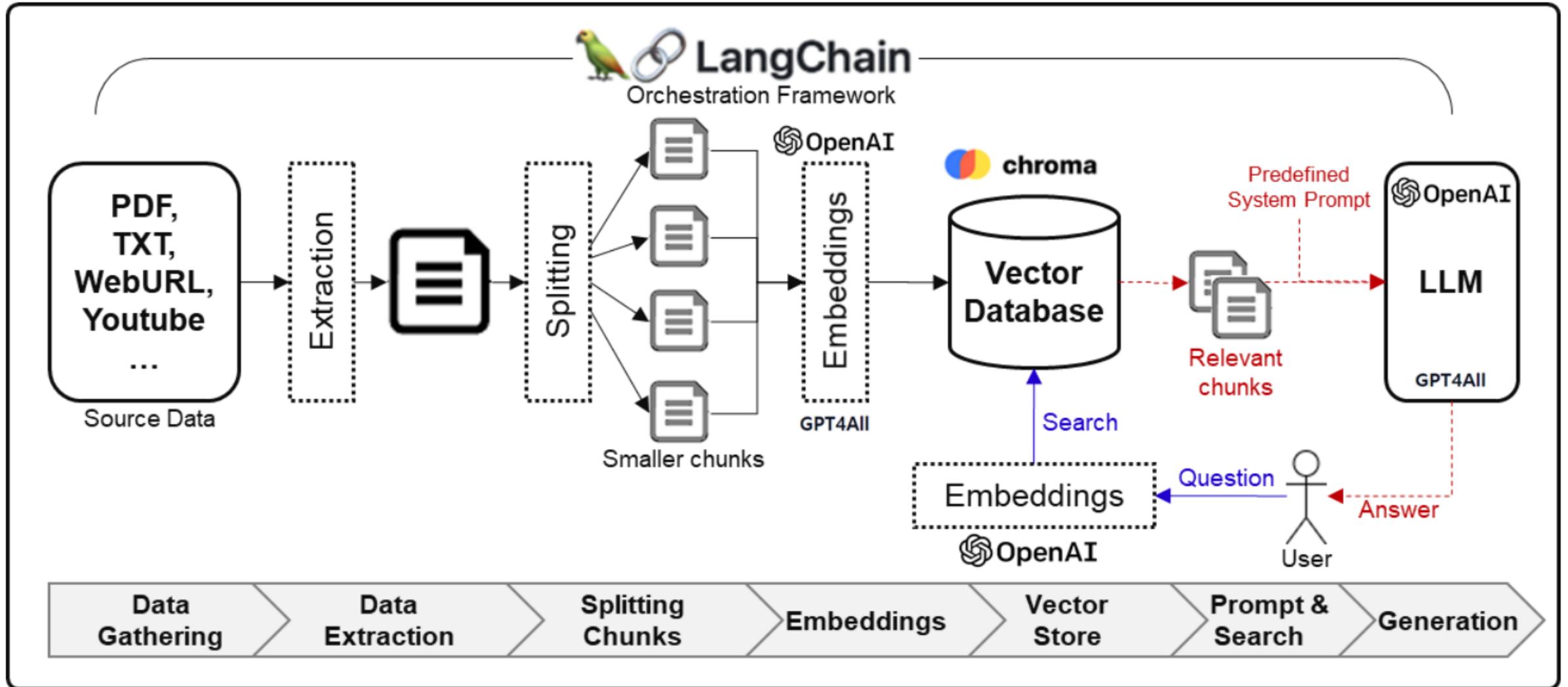
# AutoDev: Automated AI-Driven Development



# AutoDev: Automated AI-Driven Development



# Framework for Implementing Generative AI Services using RAG Model



Spring 2025

**Generative AI**  
**Innovative Applications**



## University Ambassador



This certificate acknowledges that

# Min-Yuh Day

has been certified to deliver NVIDIA instructor-led workshop for  
academia

A handwritten signature in black ink, appearing to read "Greg Estes", written over a horizontal line.

**Greg Estes**

Vice President, NVIDIA

Issue Date: : March 7, 2025

Ambassador Certification ID: cCFh1ZWWTvqKTq7dcKkEWw



## Certified Instructor



This certificate acknowledges that

# Min-Yuh Day

has been certified to deliver the instructor-led workshop

## Building RAG Agents with LLMs

A handwritten signature in black ink, appearing to read "Greg Estes".

**Greg Estes**

Vice President, NVIDIA

Issue Date: : March 7, 2025

Certification ID: OVmqY4cSSya0BdMQBWHxzw

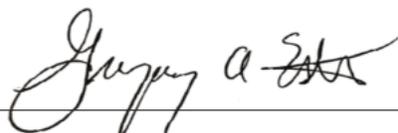
## Certificate of Completion

This certificate is awarded to

# Min-Yuh Day

for successfully completing

### Building RAG Agents with LLMs



**Greg Estes**

Vice President, NVIDIA

Issue Date: : December 8, 2024

Certification ID: ed-qOCIMQatzU8SNUNxgw |

[https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw/courses/course?course\\_id=course-v1:DLI+S-FX-15+V1](https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw/courses/course?course_id=course-v1:DLI+S-FX-15+V1)

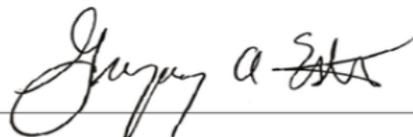
<https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw>

## Certificate of Competency

This certificate is awarded to

# Min-Yuh Day

for demonstrating competence in the completion of  
**Generative AI with Diffusion Models**



**Greg Estes**

Vice President, NVIDIA

Issue Date: : February 28, 2025

Certification ID: q3oo-oBhTQKtyCCote2E-Q

# NVIDIA Certified Instructors (12 from Taiwan)

The screenshot shows the NVIDIA Deep Learning Institute website's instructor directory search results. The page is titled "Deep Learning Institute" and features a navigation bar with links for "Find Training", "Self Paced Courses", "Instructor-Led Workshops", and "Educator Programs". A search bar at the top right contains the text "Search for instructors by name. Minimum". Below the search bar, there are filter options for "Workshop Certification", "Location [1]", "Organization", and "Specialization", along with a "Reset filters" button. The results are sorted by "Name A-Z".

Name	Organization	Featured Workshop
Chi-Hung Chuang	Chung Yuan Christian University	Fundamentals of Deep Learning, Taiwanese
Chia Yu Hsu	National Taiwan University of Science and Technology (NTUST)	Applications of AI for Predictive Maintenance, Taiwanese
Chien-Yu Chen	National Taiwan University (NTU)	Fundamentals of Deep Learning, Taiwanese
Chun-Yi Lee	National Taiwan University (NTU)	Fundamentals of Deep Learning, Taiwanese
David Tseng	Cavedu	Getting Started with AI on Jetson Nano, Taiwanese
Hsinmin Lu	National Taiwan University (NTU)	Fundamentals of Deep Learning, English
Min-Yuh Day	National Taipei University (NTPU)	Building RAG Agents with LLMs, English
Ming-Che Chen	Southern Taiwan University of Science and Technology (STUST)	Getting Started with AI on Jetson Nano, Taiwanese
MingChe Hu	Chung Yuan Christian University	Fundamentals of Deep Learning, Taiwanese
Ping-Chun Hsieh	National Yang Ming Chiao Tung University (NYCU)	Fundamentals of Deep Learning, English
Po-Chih Kuo	National Tsing Hua University (NTHU)	Fundamentals of Deep Learning, English
Shu-Kai Hsieh	National Taiwan University (NTU)	Fundamentals of Deep Learning, Taiwanese

# NVIDIA Developer Program

<https://developer.nvidia.com/join-nvidia-developer-program>

## NVIDIA

## Deep Learning Institute (DLI)

<https://learn.nvidia.com/>

# Get NVIDIA DLI Certificate

- **Step 1. Join NVIDIA Developer Program (Free)**  
<https://developer.nvidia.com/join-nvidia-developer-program>
- **Step 2. Visit NVIDIA Deep Learning Institute (DLI)**  
<https://learn.nvidia.com/>
- **Step 3. Enroll "Generative AI with Diffusion Models"**  
Self-Paced Course (\$90)  
[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)

# Join the NVIDIA Developer Program

take one of the  
complimentary  
technical self-  
paced courses  
(worth up to \$90)

Generative AI and LLMs   Graphics and Simulation   Accelerated Computing   Data Science   Deep Learning

<p>8 hours</p> <h3>Getting Started With Deep Learning</h3> <p>Explore the fundamentals of deep learning by training neural networks and using results to improve performance and capabilities.</p>	<p>2 hours</p> <h3>Modeling Time-Series Data With Recurrent Neural Networks in Keras</h3> <p>Explore how to classify and forecast time-series data using recurrent neural networks (RNNs), such as modeling a patient's health over time.</p>	<p>4 hours</p> <h3>Deploying a Model for Inference at Production Scale</h3> <p>Learn how to deploy your own machine learning models on a GPU server.</p>
<p>8 hours</p> <h3>Building Real-Time Video AI Applications</h3> <p>Gain the knowledge and skills needed to enable the real-time transformation of raw video data from widely deployed camera sensors into deep learning-based insights.</p>	<p>2 hours</p> <h3>Introduction to Graph Neural Networks</h3> <p>Learn the basic concepts, models, and applications of graph neural networks.</p>	<p>4 hours</p> <h3>Introduction to Physics-Informed Machine Learning With Modulus</h3> <p>Learn the various building blocks of NVIDIA Modulus, which turbocharges use cases by building physics-based deep learning models that are 100,000X faster than traditional methods and offers high-fidelity simulation results.</p>
<p>2 hours</p> <h3>Get Started With Highly Accurate Custom ASR for Speech AI</h3> <p>Learn to build, train, fine-tune, and deploy a GPU-accelerated automatic speech recognition (ASR) service with NVIDIA® Riva that includes customized features.</p>	<p>2 hours</p> <h3>Integrating Sensors With NVIDIA DRIVE</h3> <p>Find out how to integrate automotive sensors into your applications using NVIDIA DRIVE®.</p>	

<https://developer.nvidia.com/join-nvidia-developer-program>

# NVIDIA Deep Learning Institute (DLI)

Self-Paced Course

**Generative AI Explained**

Free  
2 hours

Self-Paced Course

**Getting Started With Deep Learning**

Certificate available  
\$90  
8 hours

Instructor-Led Workshop

**Fundamentals of Deep Learning**

Certificate available  
\$500  
8 hours

Self-Paced Course

**Introduction to Transformer-Based Natural Language Processing**

Certificate available  
\$30  
6 hours

Self-Paced Course

**Building RAG Agents With LLMs**

Certificate available  
Free  
8 hours

Instructor-Led Workshop

**Building RAG Agents With LLMs**

Certificate available  
\$500  
8 hours

Self-Paced Course

**Generative AI with Diffusion Models**

Certificate available  
\$90  
8 hours

Instructor-Led Workshop

**Generative AI with Diffusion Models**

Certificate available  
\$500  
8 hours

## What do you want to learn today?

### Filters

Level +

Format +

Topics -

- Deep Learning
- Accelerated Computing
- Generative AI/LLM
- Graphics and Simulation
- OpenUSD
- Data Science
- NIMS
- NIM
- RAPIDS

Free / Paid +

Language +



Sort by: -- ▾

Showing 19 results

**Generative AI** x

# Generative AI

## All Courses

<b>Self-paced</b> <b>Generative AI Explained</b>  Free 02:00	<b>Self-paced</b> <b>Generative AI with Diffusion Models</b>  \$90 08:00	<b>Instructor-Led</b> <b>Generative AI with Diffusion Models</b>  08:00
<b>Self-paced</b> <b>Augment your LLM Using</b>	<b>Self-paced</b> <b>Introduction to Transformer-</b>	<b>Instructor-Led</b> <b>Rapid Application</b>

Self-paced Course

## Generative AI Explained

In this no-coding course, learn Generative AI concepts and applications, as well as the challenges and opportunities in this exciting field.

[About Course](#) [Objectives](#) [Topics Covered](#) [Course Outline](#) [Stay Informed](#) [Contact Us](#)

[Continue Learning](#)

### About this Course

Generative AI describes technologies that are used to generate new content based on a variety of inputs. In recent time, Generative AI involves the use of neural networks to identify patterns and structures within existing data to generate new content. In this course, you will learn Generative AI concepts, applications, as well as the challenges and opportunities in this exciting field.

### Learning Objectives

Upon completion, you will have a basic understanding of Generative AI and be able to more effectively use the various tools built on this

### Course Details

**Duration:** 02:00

**Price:** Free

**Level:** Technical - Beginner

**Subject:** Generative AI/LLM

**Language:** English

[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)

Self-paced Course

## Building RAG Agents with LLMs

Agents powered by large language models (LLMs) have shown great retrieval capability for using tools, looking at documents, and plan their approaches. This course will show you how to deploy an agent system in practice with the flexibility to scale up your system to meet the demands of users and customers.



[Continue Learning](#)

## About this Course

**This course is free for a limited time.**

The evolution and adoption of large language models (LLMs) have been nothing short of revolutionary, with retrieval-based systems at the forefront of this technological leap. These models are not just tools for automation; they are partners in enhancing productivity, capable of holding informed conversations by interacting with a vast array of tools and documents. This course is designed for those eager to explore the potential of these systems, focusing on practical deployment and the efficient implementation required to manage the considerable demands of both users and deep learning models. As we delve into the intricacies of LLMs, participants will gain insights into advanced orchestration techniques that include internal reasoning, dialog management, and effective tooling strategies.

## Course Details

**Duration:** 08:00

**Price:** Free

**Level:** Technical - Intermediate

**Subject:** Generative AI/LLM

**Language:** English

**Course Prerequisites:**

Introductory deep learning knowledge, with comfort

[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)

Self-paced Course

## Generative AI with Diffusion Models

Take a deeper dive into denoising diffusion models, which are a popular choice for text-to-image pipelines, with applications in creative content generation, data augmentation, simulation and planning, anomaly detection, drug discovery, personalized recommendations, and more.



[About Course](#) [Objectives](#) [Topics Covered](#) [Course Outline](#) [Stay Informed](#) [Contact Us](#)

[Continue Learning](#)

### About this Course

Thanks to improvements in computing power and scientific theory, generative AI is more accessible than ever before. Generative AI plays a significant role across industries due to its numerous applications, such as creative content generation, data augmentation, simulation and planning, anomaly detection, drug discovery, personalized recommendations, and more. In this course, learners will take a deeper dive into denoising diffusion models, which are a popular choice for text-to-image pipelines.

### Learning Objectives

### Course Details

**Duration:** 08:00

**Price:** \$90

**Subject:** Generative AI/LLM

**Language:** English

**Course Prerequisites:**

A basic understanding of [Deep Learning Concepts](#).

[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)

Self-paced Course

## Rapid Application Development with Large Language Models (LLMs)

Get started quickly in developing LLM-based applications by exploring the open-sourced ecosystem including pretrained LLMs.

Self-paced courses are temporarily unavailable for purchase outside the USA as we transition to a new ecommerce system. We apologize for any inconvenience. **Free courses** remain available for enrollment.

About Course Objectives Topics Covered Course Outline Stay Informed Contact Us

Buy Now Redeem Code

### About this Course

Recent advancements in both the techniques and accessibility of large language models (LLMs) have opened up unprecedented opportunities to help businesses streamline their operations, decrease expenses, and increase productivity at scale. Additionally, enterprises can use LLM-powered apps to provide innovative and improved services to clients or strengthen customer relationships. For example, enterprises could provide customer support via AI companions or use sentiment analysis apps to extract valuable customer insights. In this course you will gain a strong understanding and practical knowledge of LLM application development by exploring the open-sourced ecosystem including pretrained LLMs, enabling you to get started quickly in developing LLM-based applications.

### Learning Objectives

By participating in this course, you will:

- Find, pull in, and experiment with the HuggingFace model repository and Transformers API.
- Use encoder models for tasks like semantic analysis, embedding, question-answering, and zero-shot classification.
- Work with conditioned decoder-style models to take in and generate interesting data formats, styles, and modalities.
- Kickstart and guide generative AI solutions for safe, effective, and scalable natural data tasks.
- Explore the use of LangChain for orchestrating data pipelines and environment-enabled agents.

### Course Details

**Duration:** 08:00

**Price:** \$90

**Level:** Technical - Beginner

**Subject:** Generative AI/LLM

**Language:** English

**Course Prerequisites:**

Introductory deep learning, with comfort with PyTorch and transfer learning preferred. Content covered by [DLI's Getting Started with Deep Learning](#) or [Fundamentals of Deep Learning](#) courses, or similar experience is sufficient.

Intermediate Python experience, including object-oriented programming and libraries. Content covered by

# Rapid Application Development with Large Language Models (LLMs)

Search



## Monthly Activity

Skill Points	0
Time Spent	
Courses in Progress	16
Courses Completed	12
Watched Videos	
Assessments	

## Skills

## Certificates



Introduction to Transformer-Based Natural Language Processing



Building RAG Agents with LLMs



Building RAG Agents with LLMs



Accelerating End-to-End Data Science Workflows



Generative AI with Diffusion Models



Building Agentic AI Applications with LLMs

## Completed Courses

View more < >

Self-paced

Sizing LLM Inference Systems

100% Completed  
03:00

Self-paced

Augment your LLM Using Retrieval Augmented Generation

100% Completed  
01:00

Self-paced

Building RAG Agents with LLMs

100% Completed  
08:00

Self-paced

Generative AI Explained

100% Completed  
02:00

Self-paced

Introduction to Transform Based Natural Language Processing

100% Completed  
06:00

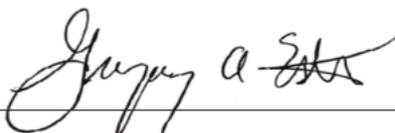
## Certificate of Completion

This certificate is awarded to

# Min-Yuh Day

for successfully completing

### Building RAG Agents with LLMs



**Greg Estes**

Vice President, NVIDIA

Issue Date: : December 8, 2024

Certification ID: ed-qOCIMQatzU8SNUNxgw |

[https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw/courses/course?course\\_id=course-v1:DLI+S-FX-15+V1](https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw/courses/course?course_id=course-v1:DLI+S-FX-15+V1)

<https://learn.nvidia.com/certificates?id=ed-qOCIMQatzU8SNUNxgw>

## All Self-Paced Courses

Accelerated Computing Data Science Deep Learning **Generative AI/LLM** Graphics and Simulation Infrastructure

[Share Generative AI/LLM Courses](#)

<p>Self-paced</p> <p><b>Generative AI Explained</b></p> <p>Free 02:00</p>	<p>Self-paced</p> <p><b>Introduction to NVIDIA NIM™ Microservices</b></p> <p>Free 02:00</p>	<p>Self-paced</p> <p><b>Introduction to Deploying RAG Pipelines for Production at Scale</b></p> <p>\$90 03:00</p>	<p>Self-paced</p> <p><b>Generative AI with Diffusion Models</b></p> <p>\$90 08:00</p>
<p>Self-paced</p> <p><b>Techniques for Improving the Effectiveness of RAG Systems</b></p> <p>\$30 03:00</p>	<p>Self-paced</p> <p><b>Introduction to Transformer- Based Natural Language Processing</b></p> <p>\$30 06:00</p>	<p>Self-paced</p> <p><b>Building LLM Applications With Prompt Engineering</b></p> <p>\$90 08:00</p>	<p>Self-paced</p> <p><b>Synthetic Tabular Data Generation Using Transformers</b></p> <p>\$30 04:00</p>
<p>Self-paced</p> <p><b>Sizing LLM Inference Systems</b></p> <p>Free 03:00</p>	<p>Self-paced</p> <p><b>Building RAG Agents with LLMs</b></p> <p>Free 08:00</p>	<p>Self-paced</p> <p><b>Augment your LLM Using Retrieval Augmented Generation</b></p> <p>Free 01:00</p>	

# NVIDIA Deep Learning Institute (DLI)

**Deep Learning Institute** Find Training Self Paced Courses Instructor-Led Workshops Educator Programs Enterprise Solutions Certification Resources

## Building RAG Agents with LLMs

Course Progress Bookmarks Updates

Building RAG Agents with LLMs Introduction Introduction

Building RAG Agents with LLMs

Introduction

Environment and LLMs

LangChain

Documents and Embeddings

Retrieval-Augmented Generation

Next Steps

Feedback

Previous

Next



## Building RAG Agents with LLMs

Introduction



# Building RAG Agents with LLMs

## Building RAG Agents with LLMs

[Course](#)
[Progress](#)
[Bookmarks](#)
[Updates](#)

[Building RAG Agents with LLMs](#)
[Environment and LLMs](#)
[Environment \[0, 1, 2\]](#)

Building RAG Agents with LLMs

Introduction

Introduction

Course Slides

Environment and LLMs

Environment [0, 1, 2]

Part 1: Course Environment

Part 2: LLM Services

LangChain

Environment [3, 4]

Part 3: LangChain

Previous

Next

welcome to **BUILDING RAG AGENTS WITH LLMs**. In this first section, we will get introduced to the overall course environment, LLM services, and recommended workflows!

**This tab contains the course environment for this section, which will contain the notebooks for the next two videos!** Please click through the videos in the remaining tabs to watch the material and work through the exercises!

**Please click the "Start" button to start up your own private server for hands-on coding practice.** It will take a few minutes to start up, so go ahead and click it now and then proceed to the next video! After a few minutes when the server has loaded, click "Launch" to access the code labs.



This Lab 0:01:06 / 2:00:00

Course 13:45:51 / 32:00:00



# Building RAG Agents with LLMs

Building RAG Agents with LLMs

Introduction

Environment and LLMs

LangChain

Documents and Embeddings

Environment [5, 6]

Part 5: Documents

Part 6: Embeddings

Retrieval-Augmented Generation

Environment [7, 8, Assessment]

Part 7: Vector Stores

Part 8: Evaluation

Next Steps

Previous

Next

[Bookmark this page](#)

**In this section, we will combine all of our prior efforts to integrate and evaluate retrieval-augmented generation pipelines!** Along the way, you will also get the opportunity to work through the assessment, which will involve Gradio, LangServe, FAISS, RAG, and Evaluation! **Good Luck!**

**Please click the "Start" button to start up your own private server for hands-on coding practice.** It will take a few minutes to start up, so go ahead and click it now and then proceed to the next video! After a few minutes when the server has loaded, click "Launch" to access the code labs.



This Lab 0 : 15 : 39 / 4 : 00 : 00



Course 14 : 12 : 18 / 32 : 00 : 00

LAUNCH

STOP TASK

ASSESS TASK

# Building RAG Agents with LLMs

← → ↻ Not Secure 34.227.20.149/lab/lab/tree/08\_evaluation.ipynb ☆ 📄 🔍 ⋮

File Edit View Run Kernel Tabs Settings Help

Name	Last Modified
chatbot	yesterday
composer	yesterday
docker_router	yesterday
frontend	yesterday
imgs	yesterday
llm_client	yesterday
slides	yesterday
solutions	yesterday
00_jupyterlab.ipynb	yesterday
01_microservices.ipynb	yesterday
02_llms.ipynb	yesterday
03_langchain_intro.ipynb	yesterday
04_running_state.ipynb	yesterday
05_documents.ipynb	yesterday
06_embeddings.ipynb	yesterday
07_vectorstores.ipynb	yesterday
08_evaluation.ipynb	yesterday
09_langserve.ipynb	yesterday
64_guardrails.ipynb	yesterday
99_table_of_contents.ipynb	yesterday

## DEEP LEARNING INSTITUTE

## Notebook 8 [Assessment]: RAG Evaluation

Welcome to the last notebook of the course! In the previous notebook, you integrated a vector store solution into a RAG pipeline! In this notebook, you will take that same pipeline and evaluate it using numerical RAG evaluation techniques incorporating LLM-as-a-Judge metrics!

**Learning Objectives:**

- Learn how to integrate the techniques from prior notebooks to numerically approximate the goodness of your RAG pipeline.

Simple 0 \$ 12 Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 08\_evaluation.ipynb 1 🔔

# AI Agent 對政府的核心價值

## 全球智慧城市投資趨勢與效益分析

- **市場規模**：AI Agent 市場從 51 億美元成長至 2030 年 471 億美元 (Alvarez & Marsal, 2025)
- **政府投資**：全球 66% 城市大力投資 AI 技術 (UN-Habitat, 2024)
- **投資效益**：成功實作案例顯示 3:1 至 5:1 的投資報酬率 (Coolfire Solutions, 2024)
- **實際案例**：
  - **新加坡**：AI 交通管理年省 10 億美元 (Singapore Smart Nation, 2024)
  - **巴塞隆納**：智慧城市 AI 綜合年省 1.45 億美元 (Barcelona City Council, 2024)

# 多智慧代理系統的技術突破

## 從文字生成到整合推理系統的演進

- **Agent-Native 基礎模型（2024-2025）：**
  - **OpenAI o1/o3：**內建推理與規劃能力
  - **Gemini 2.0 Flash：**原生多模態介面互動
  - **Microsoft Copilot Agents：**角色專業化與階層協調
- **技術突破：**從文字生成器演進為具備原生代理能力的整合推理系統 (Ahmed, 2025)
- **協作模式：**專業化 Agent 間的無縫協調與資訊交換

# 大型語言模型整合與檢索增強生成技術

- **RAG 架構組件：**
  - **檢索器：**向量搜尋結合 BM25 與密集嵌入
  - **生成器：**LLM 基礎的連貫回應生成
  - **整合層：**內容排序與信心評分

# 大型語言模型整合與檢索增強生成技術

- **政府服務優勢：**
  - 即時知識更新，無需重新訓練模型
  - 透過驗證資料降低幻覺問題
  - 擴展處理大型政府文件的能力

# 政府適用的 AI Agent 開發框架

- **AutoGen (Microsoft)**
  - 企業級可靠性、 Docker 容器支援
- **CrewAI**
  - 快速原型開發、 角色導向設計
- **LangGraph**
  - 圖形化工作流程控制， 適合複雜流程
- **n8n**
  - AI Agent 視覺化工作流程自動化平台
  - 自建或雲端部署選項

# AI Agent 技術選擇與評估標準

- 安全功能與合規要求
- 擴展性與效能表現
- 整合能力與 API 支援
- 技術支援生態系統

# Outline

1. AI Agent 基礎與最新趨勢
2. AI Agent 政府應用與智慧城市實作
3. AI Agent 實施策略與治理架構

# 全球智慧城市 AI 實作成功案例

## 量化效益與關鍵成功因素分析

- **新加坡智慧國家計畫：**  
**年省 10 億美元的 AI 交通管理** (Singapore Smart Nation, 2024)
  - **1,500+ 路口中 80% 由 AI 管理**
  - **尖峰時段延誤減少 20%**
  - **碳排放減少 15% (年減 50 萬噸)**

# 全球智慧城市 AI 實作成功案例

## 量化效益與關鍵成功因素分析

- **巴塞隆納智慧城市：**  
**年省 1.45 億美元的綜合效益**

(Barcelona City Council, 2024; Harvard Kennedy School, 2024)

- **智慧水資源管理：節省 5,800 萬美元**
- **智慧停車收入：增加 5,000 萬美元**
- **智慧照明系統：節省 3,700 萬美元**

# 首爾市 AI 基礎行政體系構建計畫

2026年前投資2,064億韓元的全面 AI 轉型

- **投資規模：2,064億韓元（約1.55億美元） 至2026年** (Seoul Metropolitan Government, 2024)
- **核心應用：**
  - 獨居老人智慧安心服務
  - 119緊急電話 AI 語音分析
  - 智慧路口交通最佳化
  - 1,300位公務員生成式 AI 培訓

# 首爾市 AI 基礎行政體系構建計畫

2026年前投資2,064億韓元的全面 AI 轉型

- **治理架構**：數位政策室統籌協調，確保跨部門整合 (Seoul Metropolitan Government, 2024)
- **成功關鍵**：高層承諾與系統性人才培育

# 智慧城市 AI 投資報酬率分析框架

## 量化效益評估與成本效益計算

- **量化效益** (Thales Group, 2024; Minnovation Technologies, 2024) :
  - **交通管理**：旅行時間減少 12-25%
  - **行政效率**：處理速度提升最高 50%
  - **能源管理**：消耗量減少 10-25%
  - **市民滿意度**：服務品質提升 35%
- **投資回報**：智慧感測器 5 年期投資報酬率達 24:1  
(Coolfire Solutions, 2024)

# 智慧城市 AI 投資報酬率分析框架

## 量化效益評估與成本效益計算

- **成本節約：**

- **Barcelona 智慧城市實作年省 1.45 億美元**

(Harvard Kennedy School, 2024)

- **評估標準：**

- **技術效益、營運效率、市民滿意度、環境影響**

# 全球經驗應用於臺北智慧城市試辦計畫

## 具體實作策略與預期效益分析

- **第一組 - 無障礙出行導航：**
  - **借鑑新加坡 ITS 模式：即時無障礙路線規劃**  
(Singapore Smart Nation, 2024)
  - **預期效益：年省 50 萬美元支援服務成本**  
(Coolfire Solutions, 2024)

# 全球經驗應用於臺北智慧城市試辦計畫

## 具體實作策略與預期效益分析

- **第二組 - 青年心理支持：**
  - **參考阿姆斯特丹數位參與模式  
與首爾社會服務整合**  
(Seoul Metropolitan Government, 2024)
  - **預期效益：服務回應時間改善 40%**

# 全球經驗應用於臺北智慧城市試辦計畫

## 具體實作策略與預期效益分析

- **第三組 - 1999 服務品質：**
  - **結合新加坡 SingPass 數位身分整合  
與巴塞隆納市民服務**  
(Singapore Smart Nation, 2024; Barcelona City Council, 2024)
  - **預期效益：通話處理時間減少 60%**

# 全球經驗應用於臺北智慧城市試辦計畫

## 具體實作策略與預期效益分析

- **第四組 - 兒少心理健康：**
  - **適用首爾老人照護 AI 監測模式**  
(Seoul Metropolitan Government, 2024)
  - **預期效益：提早 30% 識別介入需求**

# 三階段 AI Agent 實施策略框架

## 從基礎建設到全面部署的系統性方法

- **第一階段：基礎建構（6-12個月）**
- **第二階段：擴展整合（12-18個月）**
- **第三階段：全面部署（持續進行）**

# 三階段 AI Agent 實施策略框架

## 從基礎建設到全面部署的系統性方法

- **第一階段：基礎建構（6-12個月）**

(Chan, 2023)

- **API 優先架構連接既有系統**
- **主資料管理統一治理**
- **特定部門試辦計畫**

# 三階段 AI Agent 實施策略框架

## 從基礎建設到全面部署的系統性方法

- **第二階段：擴展整合（12-18個月）**

(Singapore Smart Nation, 2024)

- 微服務架構模組化部署
- 即時資料串流提供即時洞察
- 跨部門協作機制建立

# 三階段 AI Agent 實施策略框架

## 從基礎建設到全面部署的系統性方法

- **第三階段：全面部署（持續進行）**

(Barcelona City Council, 2024)

- 雲原生解決方案確保擴展性
- 持續監控與改善機制
- 全市協調與最佳化

# AI Agent 實施風險評估與緩解策略

## 技術、組織與法規風險的全方位管理

- **技術風險** (DataGuard, 2024; IBM, 2024) :
  - **資料品質**：透過驗證框架確保
  - **系統可靠性**：透過冗餘機制保障
  - **資訊安全**：透過加密與存取控制防護

# AI Agent 實施風險評估與緩解策略

## 技術、組織與法規風險的全方位管理

- **組織風險** (Chan, 2023) :
  - **變革管理**：透過培訓計畫推動
  - **利害關係人參與**：透過透明溝通建立
  - **治理架構**：透過明確問責結構確立
- **評估矩陣**：
  - **高、中、低風險分級與對應緩解措施**  
(European Union, 2024)

# AI Agent 實作 (n8n + LangChain 實作)

## 1999 服務 AI 客服代理建構實例

- AI Agent 架構：
  - 市民來電 →
  - n8n 工作流程 →
  - LangChain 處理 →
  - 知識庫查詢 →
  - 回應生成 →
  - 系統記錄

# AI Agent 實作 (n8n + LangChain 實作)

## 1999 服務 AI 客服代理建構實例

- **n8n 工作流程節點：**
  - **Webhook 接收器：**接收語音轉文字結果
  - **條件判斷：**識別問題類型  
(一般諮詢/緊急事件/轉接需求)
  - **LangChain 節點：**調用 AI 模型進行語意理解
  - **RAG 檢索：**從政府 FAQ 知識庫獲取相關資訊
  - **回應合成：**生成個人化回應
  - **資料記錄：**存儲對話記錄供後續分析

# AI Agent 實作 (n8n + LangChain 實作)

## 1999 服務 AI 客服代理建構實例

- **LangChain 關鍵組件：**
  - **ConversationMemory**：維護對話歷史
  - **VectorStore**：政府政策文件向量化儲存
  - **RetrievalQA Chain**：實現 RAG 檢索問答
  - **OutputParser**：確保回應格式符合政府標準

# AI Agent 實作 (n8n + LangChain 實作)

## 實作練習建議

- **第一組：**  
**建立無障礙路線查詢 Agent，整合交通即時資訊 API**
- **第二組：**  
**設計青年諮詢 Agent，具備情緒識別和資源推薦功能**
- **第三組：**  
**完整 1999 客服 Agent，包含多語言支援和智慧分流**
- **第四組：**  
**兒少關懷 Agent，具備風險評估和轉介機制**

# Outline

1. AI Agent 基礎與最新趨勢
2. AI Agent 政府應用與智慧城市實作
3. AI Agent 實施策略與治理架構

# 政府 AI 應用倫理框架與國際標準

## UNESCO 與 OECD 指導原則的實務應用

- **UNESCO 核心價值 (UNESCO, 2021) :**
  1. **人權與人性尊嚴**
  2. **和平、公正與互聯社會**
  3. **多元性與包容性**
  4. **環境與生態系統繁榮**

# 政府 AI 應用倫理框架與國際標準

## UNESCO 與 OECD 指導原則的實務應用

- **OECD AI 原則 (OECD, 2024) :**
  1. **包容性成長與福祉**
  2. **以人為本的價值與公平性**
  3. **透明度與可解釋性**
  4. **穩健性、安全性與保障**
  5. **問責制**

# AI 治理結構設計與組織架構

## 多層級決策機制與監督體系

- **AI 治理委員會**

(Seoul Metropolitan Government, 2024; Singapore Smart Nation, 2024) •

- 多方利害關係人組成
- 策略監督與政策制定
- 定期框架檢討

# AI 治理結構設計與組織架構

## 多層級決策機制與監督體系

- **AI 倫理委員會** (MIT Technology Review, 2025) :
  - 獨立諮詢機構
  - 高風險應用評估
  - 持續倫理監控

# AI 治理結構設計與組織架構

## 多層級決策機制與監督體系

- **AI 長角色** (Chan, 2023) :
  - 行政層級 AI 策略
  - 跨部門協調
  - 外部利害關係人介面

# 基於風險評估的 AI 實施分級管理

## 從最小風險到不可接受風險的管理框架

- **風險分級 (European Union, 2024) :**
  - **最小風險**：標準合規與文件記錄
  - **有限風險**：透明度要求與 AI 標示
  - **高風險**：全面評估、人工監督、詳細稽核
  - **不可接受風險**：禁止應用（社會評分、操弄等）

# 基於風險評估的 AI 實施分級管理

## 從最小風險到不可接受風險的管理框架

- **政府特定風險**

(MIT Technology Review, 2025; DataGuard, 2024) •

- 市民服務中的演算法偏見
- 資料處理中的隱私侵犯
- 關鍵系統的安全漏洞
- 自動化決策的問責空隙

# AI 系統生命週期治理框架

## 從設計開發到部署營運的全程管控

- **第一階段：設計與開發**

(UNESCO, 2021; OECD, 2024)

- 倫理影響評估
- 偏見測試與緩解
- 隱私設計原則實施
- 人權影響評估

# AI 系統生命週期治理框架

## 從設計開發到部署營運的全程管控

- **第二階段：測試與驗證**

(MIT Technology Review, 2025)

- 全面測試協定
- 獨立驗證程序
- 利害關係人諮詢
- 受影響社群試辦測試

# AI 系統生命週期治理框架

## 從設計開發到部署營運的全程管控

- **第三階段：部署與營運**

(Barcelona City Council, 2024)

- 分階段推出與監控
- 人員培訓計畫
- 公眾溝通策略
- 持續效能監控

# AI 系統效能評估與成效測量框架

## Kirkpatrick 四層次模型在政府 AI 的應用

- **Kirkpatrick 四層次模型**

(Kirkpatrick & Kirkpatrick, 2016; Docebo, 2025) :

- **第一層次：反應**（滿意度、參與度）
- **第二層次：學習**（知識、技能）
- **第三層次：行為**（應用、效能）
- **第四層次：結果**（效率、成果）

# AI 系統效能評估與成效測量框架

## Kirkpatrick 四層次模型在政府 AI 的應用

- **政府專用指標** (SafetyCulture, 2024) :
  - 市民滿意度改善
  - 服務提供效率提升
  - 成本節約與投資報酬率測量 (Coolfire Solutions, 2024)
  - 政策合規與效力評估

# 建立公眾信任的透明度機制

## 資訊公開與市民參與的最佳實務

- **透明度要求** (Inter-Parliamentary Union, 2025) :
  - 公開 AI 系統登記與清單
  - 定期透明度報告發布
  - 主動揭露系統限制
  - 清楚說明 AI 決策過程

# 建立公眾信任的透明度機制

## 資訊公開與市民參與的最佳實務

- **市民參與策略** (Barcelona City Council, 2024; bee smart city GmbH, 2024) •
  - 參與式設計流程
  - 定期公眾諮詢
  - 回饋機制與申訴管道
  - 社群基礎監督機制

# AI Agent 實施三階段檢核表

## 確保成功部署的實務指引

- **即時行動（0-6個月）**

(Chan, 2023) :

- **建立治理結構**
- **進行 AI 系統盤點**
- **發展風險評估能力**
- **啟動利害關係人參與**

# AI Agent 實施三階段檢核表

## 確保成功部署的實務指引

- **短期行動（6-18個月）**

(European Union, 2024; Singapore Smart Nation, 2024) •

- **實施核心框架**
- **部署試辦計畫**
- **建立監控程序**
- **進行合規評估**

# AI Agent 實施三階段檢核表

## 確保成功部署的實務指引

- **長期行動（18個月以上）** (Barcelona City Council, 2024) :
  - 擴展成功實施
  - 持續改善流程
  - 國際合作交流
  - 進階能力發展

# 成果發表會簡報架構

## 展現專業能力與治理智慧的關鍵要素

- **問題定義：清楚闡述市民需求**
- **AI 解決方案設計：技術方法與能力**
- **實施策略：時程、資源、風險**
- **預期成果：可測量效益與成功指標**
- **治理計畫：倫理考量與監督機制**

# 成果發表會簡報標準

## 展現專業能力與治理智慧的關鍵要素

- 技術可行性與創新性
- 政府營運實務適用性
- 倫理考量與風險緩解
- 利害關係人參與與溝通
- 可測量成果與永續性

# 成果發表實作展示指引

## n8n + LangChain 原型展示標準

- **技術展示：**
  - **工作流程視覺化：**完整的 n8n 流程圖展示
  - **實際操作演示：**現場模擬市民服務情境
  - **效能指標呈現：**回應時間、準確率、用戶滿意度
  - **錯誤處理機制：**異常情況的處理流程

# 成果發表實作展示指引

## n8n + LangChain 原型展示標準

- **原型功能：**
  - **基本對話能力：**理解常見市民問題
  - **知識庫整合：**連接相關政策文件資料庫
  - **多輪對話支援：**維持對話上下文
  - **人工轉接機制：**複雜問題的升級處理

# 成果發表實作展示指引

## n8n + LangChain 原型展示標準

- **評估標準：**
  - **技術實現度 (30%)：** 原型完整性與功能性
  - **創新應用 (25%)：** 解決方案的創意與實用性
  - **使用者體驗 (20%)：** 介面友善度與易用性
  - **系統整合 (15%)：** 與現有政府系統的相容性
  - **擴展潛力 (10%)：** 未來發展與規模化可能性

# 成果發表實作展示指引

## n8n + LangChain 原型展示標準

- 實作交付物清單：
  - 1. n8n 工作流程檔案 (.json)
  - 2. LangChain 程式碼 (.py)
  - 3. 知識庫設計文件
  - 4. 測試案例與結果報告
  - 5. 部署說明文件
  - 6. 5 分鐘實作展示影片

# 第一組 - 無障礙出行導航 AI Agent :

- **技術要求：**
  - 使用 n8n 建立工作流程，整合臺北市即時交通 API
  - 實作 LangChain RAG 系統，知識庫包含無障礙設施資訊
  - 建立路線最佳化演算法，考慮無障礙需求

# 第一組 - 無障礙出行導航 AI Agent :

- **核心功能：**
  1. 接收使用者位置和目的地
  2. 識別無障礙需求類型（輪椅、視障、聽障等）
  3. 即時查詢交通狀況和設施可用性
  4. 生成個人化無障礙路線建議
  5. 提供語音和視覺雙重輸出

# 第一組 - 無障礙出行導航 AI Agent :

- **展示情境：**
  - **輪椅使用者從臺北車站到市政府的最佳路線**
  - **視障者搭乘捷運的詳細導引**
  - **臨時電梯故障時的替代方案推薦**

# 第二組 - 青年心理支持 AI Agent :

- **技術要求 :**
  - **n8n 工作流程整合情緒分析 API**
  - **LangChain 對話系統具備心理諮商知識庫**
  - **實作風險評估與轉介機制**

# 第二組 - 青年心理支持 AI Agent :

- **核心功能：**
  1. **自然語言情緒狀態分析**
  2. **基於心理學原理的對話引導**
  3. **危機識別與緊急轉介**
  4. **資源推薦與追蹤關懷**
  5. **隱私保護與資料加密**

# 第二組 - 青年心理支持 AI Agent :

- **展示情境：**
  - **青年表達學業壓力的對話處理**
  - **識別自殺傾向並啟動緊急程序**
  - **提供客製化心理健康資源推薦**

# 第三組 - 1999 服務品質 AI Agent :

- **技術要求 :**
  - **n8n 多語言語音轉文字工作流程**
  - **LangChain RAG 整合完整政府 FAQ 資料庫**
  - **智慧分流與人工轉接機制**

# 第三組 - 1999 服務品質 AI Agent :

- **核心功能：**
  1. **多語言語音識別與處理**
  2. **意圖識別與問題分類**
  3. **知識庫檢索與答案生成**
  4. **複雜問題智慧分流**
  5. **服務品質即時監控**

# 第三組 - 1999 服務品質 AI Agent :

- **展示情境：**
  - **處理繁體中文、英文、台語混合查詢**
  - **複雜稅務問題的專業轉接**
  - **緊急案件的快速升級處理**

# 第四組 - 兒少心理健康 AI Agent :

- **技術要求 :**
  - **n8n 多源資料整合工作流程**
  - **LangChain 兒童心理發展知識系統**
  - **跨部門協作通報機制**

# 第四組 - 兒少心理健康 AI Agent :

- **核心功能：**
  1. **多模態兒童行為分析**
  2. **發展里程碑追蹤評估**
  3. **風險因子預警系統**
  4. **跨部門資源整合推薦**
  5. **家長與教師協作平台**

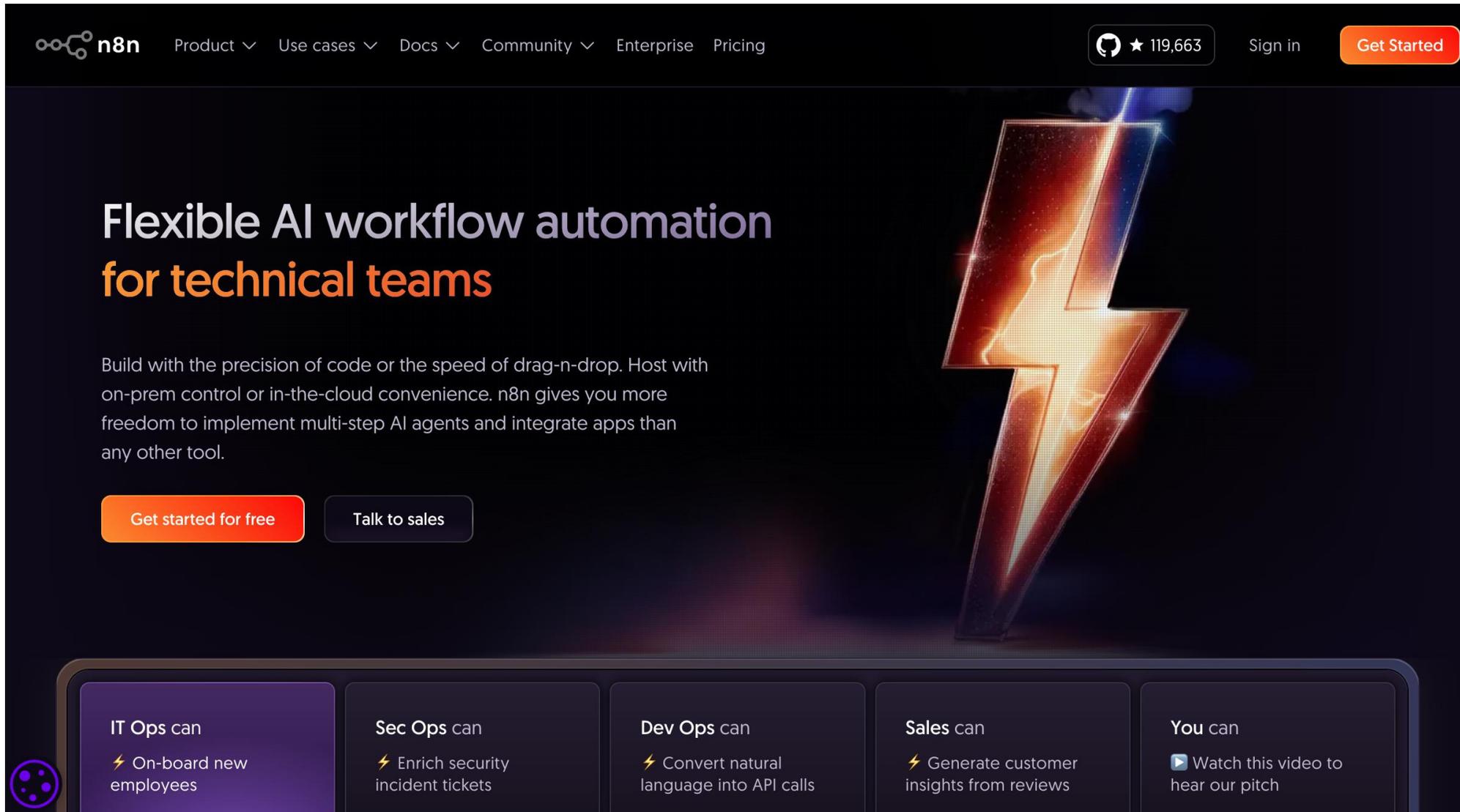
# 第四組 - 兒少心理健康 AI Agent :

- **展示情境：**
  - **分析兒童繪畫作品識別情緒問題**
  - **學校通報系統的自動風險評估**
  - **整合社會局、教育局、衛生局資源推薦**

# 開發環境設置

- **n8n 本地安裝指引與雲端部署選項**
- **LangChain 開發環境配置**
- **臺北市政府 API 存取權限申請**
- **測試資料集與模擬環境**

# AI Agent with n8n (nodemation)



The screenshot shows the n8n website homepage. At the top, there is a navigation bar with the n8n logo, menu items for Product, Use cases, Docs, Community, Enterprise, and Pricing, a GitHub repository star count of 119,663, a Sign in link, and a Get Started button. The main content area features a large, glowing lightning bolt graphic on the right. The headline reads "Flexible AI workflow automation for technical teams". Below this, a paragraph describes the tool's flexibility in building AI agents. Two buttons, "Get started for free" and "Talk to sales", are positioned below the text. At the bottom, a horizontal row of five cards lists use cases: IT Ops (onboarding employees), Sec Ops (enriching tickets), Dev Ops (converting natural language to API calls), Sales (generating insights from reviews), and a general "You can" card with a video pitch link.

n8n Product Use cases Docs Community Enterprise Pricing 119,663 Sign in Get Started

## Flexible AI workflow automation for technical teams

Build with the precision of code or the speed of drag-n-drop. Host with on-prem control or in-the-cloud convenience. n8n gives you more freedom to implement multi-step AI agents and integrate apps than any other tool.

Get started for free Talk to sales

- IT Ops can**
  - ⚡ On-board new employees
- Sec Ops can**
  - ⚡ Enrich security incident tickets
- Dev Ops can**
  - ⚡ Convert natural language into API calls
- Sales can**
  - ⚡ Generate customer insights from reviews
- You can**
  - ▶ Watch this video to hear our pitch

<https://n8n.io/>

# AI Agent with n8n Workflow Automation Templates

The screenshot shows the n8n.io/workflows/ website. The browser address bar displays 'n8n.io/workflows/'. The website header includes the n8n logo, navigation links for Product, Use cases, Docs, Community, Enterprise, and Pricing, a GitHub repository link with 119,663 stars, a Sign in button, and a Get Started button. The main content area features a large heading '3501 Workflow Automation Templates' and a search bar with the placeholder text 'Search apps, roles, usecases...'. Below the search bar are several category buttons: AI, Sales, IT Ops, Marketing, Document Ops, Other, and Support. At the bottom left, there is a purple circular icon with four dots and the text 'Newcomer essentials: learn by doing'.

<https://n8n.io/workflows/>



# AI Agent with n8n

## Workflow Automation Templates

The screenshot displays the n8n web interface. On the left, a workflow template titled "AI Timesheet Generator with Gmail, Calendar & GitHub to Google Sheets" is visible, with a "Use for free" button. On the right, a "Use template" modal is open, providing options to import the template to an n8n destination. The modal includes a close button (X) in the top right corner.

**Use template**

import to an n8n destination

- Import template to **myday** cloud workspace
- Copy template to clipboard (JSON)

Get started with n8n

- Get started free with n8n cloud
- Open self-hosting installation docs

<https://n8n.io/workflows/>

# AI Agent with n8n Workflow Automation Templates

The screenshot shows the n8n website interface. At the top, there is a navigation bar with the n8n logo, links for Product, Use cases, Docs, Community, Enterprise, and Pricing, a star rating of 119,663, a Sign in button, and a Get Started button. Below the navigation bar, there is a left sidebar with a 'Back to Templates' link, a menu icon, a pencil icon, a Telegram icon, and a '+6' button. The main content area features the title 'Angie, Personal AI Assistant with Telegram Voice and Text' and a 'Use for free' button. The central focus is a workflow diagram titled 'Process Telegram Request'. The workflow starts with 'Listen for Incoming events' (Update message), followed by 'Voice or Text' (Manual), an 'If' condition, and 'Get Voice File' (get file). It then proceeds to 'Speech to Text' (Transcribe Recording), 'Angie, AI Assistant' (Chat GPT), and 'Telegram' (sendMessage message). Below the main workflow, there are several sub-workflows connected to the 'Angie, AI Assistant' node, including 'OpenAI Chat ModelWindow Buffer Memory', 'Memory', 'Get Email' (google message), 'Google Calendar' (calendar event), 'Tasks' (update task), and 'Contacts' (get contacts). At the bottom of the workflow diagram, there are four icons: a square with a circle, a magnifying glass, a magnifying glass with a minus sign, and a circular arrow.

<https://n8n.io/workflows/2462-angie-personal-ai-assistant-with-telegram-voice-and-text/>

# AI Agent with n8n (nodemation)

<https://n8n.io>

n8n



Personal / Angie, Personal AI Assistant... + Add tag

Inactive

Share

Save

Star 119,777

Overview

Personal

Projects

+ Add project

Admin Panel

Templates

Variables

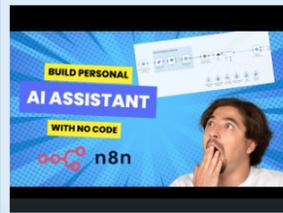
Insights

Help

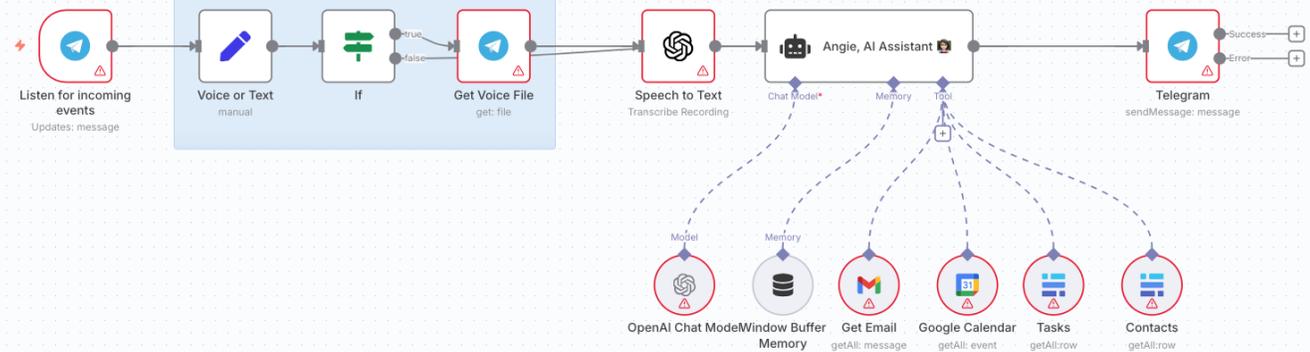
What's New

Set up template

Start here: Step-by Step Youtube Tutorial



Process Telegram Request



Execute workflow

# AI Agent with n8n (nodemation)

<https://n8n.io>

The screenshot displays the n8n workflow editor interface. The browser address bar shows the URL: `imyday.app.n8n.cloud/workflow/new?projectId=rKk0Z2QTLwuwRqTm`. The interface includes a left sidebar with navigation options: Overview, Personal, Projects (with an 'Add project' button), Admin Panel, Templates, Variables, Insights, Help, and What's New. The main workspace shows a workflow on a grid background. The workflow starts with a trigger node labeled 'When chat message received'. This is followed by an 'AI Agent' node, which is highlighted with a red border. Below the AI Agent node are three sub-nodes: 'Chat Model\*', 'Memory', and 'Tool', each with a plus sign icon for adding more instances. The workflow ends with a plus sign icon. At the top right of the editor, there are controls for 'Inactive' (a toggle switch), 'Share', 'Save', and a 'Star' button with a count of 119,739. Below the workflow, there is a toolbar with icons for full screen, zoom in, zoom out, and a red 'Open chat' button. The bottom status bar shows 'Chat', 'Session: 57a31...', and 'Logs'.

**結語：**

# 人工智慧新趨勢（AI Agent） 從學習到實踐的轉換

- **技術實踐**：開發真正解決市民問題的 AI 系統
- **治理創新**：建立負責任的 AI 治理機制
- **經驗分享**：為全球智慧城市發展貢獻台北經驗
- **持續學習**：在快速變化的 AI 世界中保持領先

# Summary

1. AI Agent 基礎與最新趨勢
2. AI Agent 政府應用與智慧城市實作
3. AI Agent 實施策略與治理架構

# References

- Ahmed, S. (2025, January 15). AI agents in 2025: A comprehensive review and future outlook. Medium. <https://medium.com/@sahin.samia/current-trends-in-ai-agents-use-cases-and-the-future-ahead-1026c4d753fd>
- Alvarez & Marsal. (2025). Demystifying AI agents in 2025: Separating hype from reality and navigating market outlook. <https://www.alvarezandmarsal.com/thought-leadership/demystifying-ai-agents-in-2025-separating-hype-from-reality-and-navigating-market-outlook>
- Barcelona City Council. (2024). Smart city Barcelona: IoT integration and citizen engagement. Smart City Hub. <https://smartcityhub.com/technology-innovation/barcelona-showcase-smart-city-dynamics/>
- California Department of Education. (2024). Learning with AI, learning about AI: Professional learning framework. <https://www.cde.ca.gov/ci/pl/aiincalifornia.asp>
- Chan, C. K. Y. (2023). A comprehensive AI policy education framework for university teaching and learning. International Journal of Educational Technology in Higher Education, 20(1), 38. <https://doi.org/10.1186/s41239-023-00408-3>
- Coolfire Solutions. (2024). Smart city initiatives: Where's the ROI? <https://coolfiresolutions.com/blog/measuring-smart-city-roi/>
- DataGuard. (2024). The growing data privacy concerns with AI: What you need to know. <https://www.dataguard.com/blog/growing-data-privacy-concerns-ai/>
- Dibia, V. (2024, December 28). AI agents 2024 rewind: A year of building and learning. Medium. <https://medium.com/@victor.dibia/ai-agents-2024-rewind-a-year-of-building-and-learning-fc6dd490bce2>
- Docebo. (2025). How to measure training effectiveness: The 2025 framework. <https://www.docebo.com/learning-network/blog/how-to-measure-training-effectiveness/>
- Emerging Tech Brew. (2022, September 21). What cities can learn from Sidewalk and Toronto's failed city of the future. <https://www.emergingtechbrew.com/stories/2022/09/21/how-miscommunication-derailed-sidewalk-s-usd1-3-billion-city-of-the-future>
- European Union. (2024). Regulation (EU) 2024/1689 of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act). <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>
- Frontiers. (2022). Social acceptance of smart city projects: Focus on the Sidewalk Toronto case. Frontiers in Environmental Science, 10. <https://doi.org/10.3389/fenvs.2022.898922>
- FutureAGI. (2025). RAG architecture: Enhancing LLM agents with AI retrieval. <https://futureagi.com/blogs/rag-architecture-llm-2025>
- Google for Developers. (2025, January). A2A: A new era of agent interoperability. <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/>
- Harrison, C. (2024). LangChain documentation: Building applications with LLMs through composability. LangChain Inc. <https://python.langchain.com/docs/>
- Harvard Kennedy School. (2024). How smart city Barcelona brought the Internet of Things to life. Data-Smart City Solutions. <https://datasmart.hks.harvard.edu/news/article/how-smart-city-barcelona-brought-the-internet-of-things-to-life-789>
- IBM. (2024). What is AI agent evaluation? <https://www.ibm.com/think/topics/ai-agent-evaluation>
- Inter-Parliamentary Union. (2025, February). What AI means for public engagement with parliaments. <https://www.ipu.org/news/case-studies/2025-02/what-ai-means-public-engagement-with-parliaments>
- Kirkpatrick, D. L., & Kirkpatrick, J. D. (2016). Evaluating training programs: The four levels (3rd ed.). Berrett-Koehler Publishers.
- Minnovation Technologies. (2024). Example of a smart city: A case study into Barcelona. <https://minnovation.com.au/smart-cities-2/example-of-a-smart-city-a-case-study-into-barcelona/>

# References

- MIT Technology Review. (2022, June 29). Toronto wants to kill the smart city forever. <https://www.technologyreview.com/2022/06/29/1054005/toronto-kill-the-smart-city/>
- MIT Technology Review. (2025, June 11). Inside Amsterdam's high-stakes experiment to create fair welfare AI. <https://www.technologyreview.com/2025/06/11/1118233/amsterdam-fair-welfare-ai-discriminatory-algorithms-failure/>
- n8n.io. (2024). n8n documentation: Workflow automation for technical people. n8n GmbH. <https://docs.n8n.io/>
- Organisation for Economic Co-operation and Development. (2024). How countries are implementing the OECD principles for trustworthy AI. OECD.AI. <https://oecd.ai/en/wonk/national-policies-2>
- Policy Options. (2020, May). Sidewalk Labs' city-of-the future in Toronto was a stress test we needed. Institute for Research on Public Policy. <https://policyoptions.irpp.org/magazines/may-2020/sidewalk-labs-city-of-the-future-in-toronto-was-a-stress-test-we-needed/>
- SafetyCulture. (2024). 10 training evaluation methods. SC Training. <https://training.safetyculture.com/blog/training-evaluation-methods/>
- Seoul Metropolitan Government. (2024). Seoul aims to become the best city in AI utilization: AI-based administration to enhance well-being of citizens. <https://english.seoul.go.kr/seoul-aims-to-become-the-best-city-in-ai-utilization-ai-based-administration-to-enhance-well-being-of-citizens/>
- Singapore Smart Nation and Digital Government Office. (2024). National artificial intelligence strategy. <https://www.smartnation.gov.sg/initiatives/strategic-national-projects/>
- Thales Group. (2024). Singapore: The world's smartest city. <https://www.thalesgroup.com/en/worldwide-digital-identity-and-security/iot/magazine/singapore-worlds-smartest-city>
- UNESCO. (2021). Recommendation on the ethics of artificial intelligence. <https://www.unesco.org/en/artificial-intelligence/recommendation-ethics>
- UN-Habitat. (2024). Global assessment of responsible AI in cities: Research and recommendations to leverage AI for people-centred smart cities. <https://unhabitat.org/global-assessment-of-responsible-ai-in-cities>
- Victordibia. (2024). AI agents 2024 rewind: A year of building and learning. <https://newsletter.victordibia.com/p/ai-agents-2024-rewind-a-year-of-building>
- Wiggins, G., & McTighe, J. (2005). Understanding by design (Expanded 2nd ed.). Association for Supervision and Curriculum Development.

## Q & A

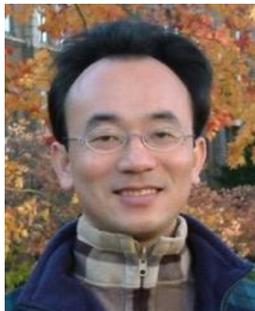
# 人工智慧新趨勢：AI 代理

## New Trends in Artificial Intelligence: AI Agent

Time: 2025/7/18 (五) 13:40-16:30

Place: E504電腦教室, 臺北市文山區萬美街二段21巷20號

Organizer: 臺北市政府公務人員訓練處 綜合企劃組 劉雨青



## 戴敏育 教授 (Prof. Min-Yuh Day)

國立臺北大學 資訊管理研究所 教授

金融科技暨綠色金融研究中心 主任

永續辦公室 永續發展組 組長

