# Artificial Intelligence

# Generative AI, Agentic AI, and Physical AI
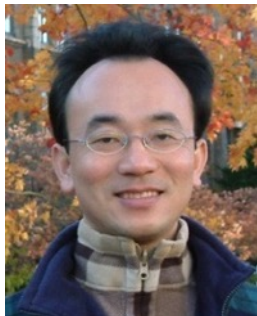
1141AI09
MBA, IM, NTPU (M5276) (Fall 2025)
Tue 2, 3, 4 (9:10-12:00) (B3F17)

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

**Institute of Information Management**, **National Taipei University**

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

2025-12-02

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

# Syllabus

Week   Date   Subject/Topics

1 2025/09/09 Introduction to Artificial Intelligence

2 2025/09/16 Artificial Intelligence and Intelligent Agents;
Problem Solving

3 2025/09/23 Knowledge, Reasoning and Knowledge Representation;
Uncertain Knowledge and Reasoning

4 2025/09/30 Case Study on Artificial Intelligence I

5 2025/10/07 Machine Learning: Supervised and Unsupervised Learning;
The Theory of Learning and Ensemble Learning

# Syllabus

**Week   Date   Subject/Topics**

**6 2025/10/14 NVIDIA Fundamentals of Deep Learning I:**
              **Deep Learning; Neural Networks**

**7 2025/10/21 NVIDIA Fundamentals of Deep Learning II:**
              **Convolutional Neural Networks;**
              **Data Augmentation and Deployment**

**8 2025/10/28 Self-Learning**

**9 2025/11/04 Midterm Project Report**

**10 2025/11/11 NVIDIA Fundamentals of Deep Learning III:**
               **Pre-trained Models; Natural Language Processing**

# Syllabus

Week   Date   Subject/Topics

11 2025/11/18 Case Study on Artificial Intelligence II

12 2025/11/25 Computer Vision and Robotics

13 2025/12/02 Generative AI, Agentic AI, and Physical AI

14 2025/12/09 Philosophy and Ethics of AI and the Future of AI

15 2025/12/16 Final Project Report I

16 2025/12/23 Final Project Report II

# Generative AI, Agentic AI, and Physical AI

# Outline

- **Generative AI**

- **Agentic AI**

- **Physical AI (Robotics)**

# Generative AI, Agentic AI, Physical AI

**Physical AI**
Self-driving cars
General robotics

**Agentic AI**
Coding assistants
Customer service
Patient care

**Generative AI**
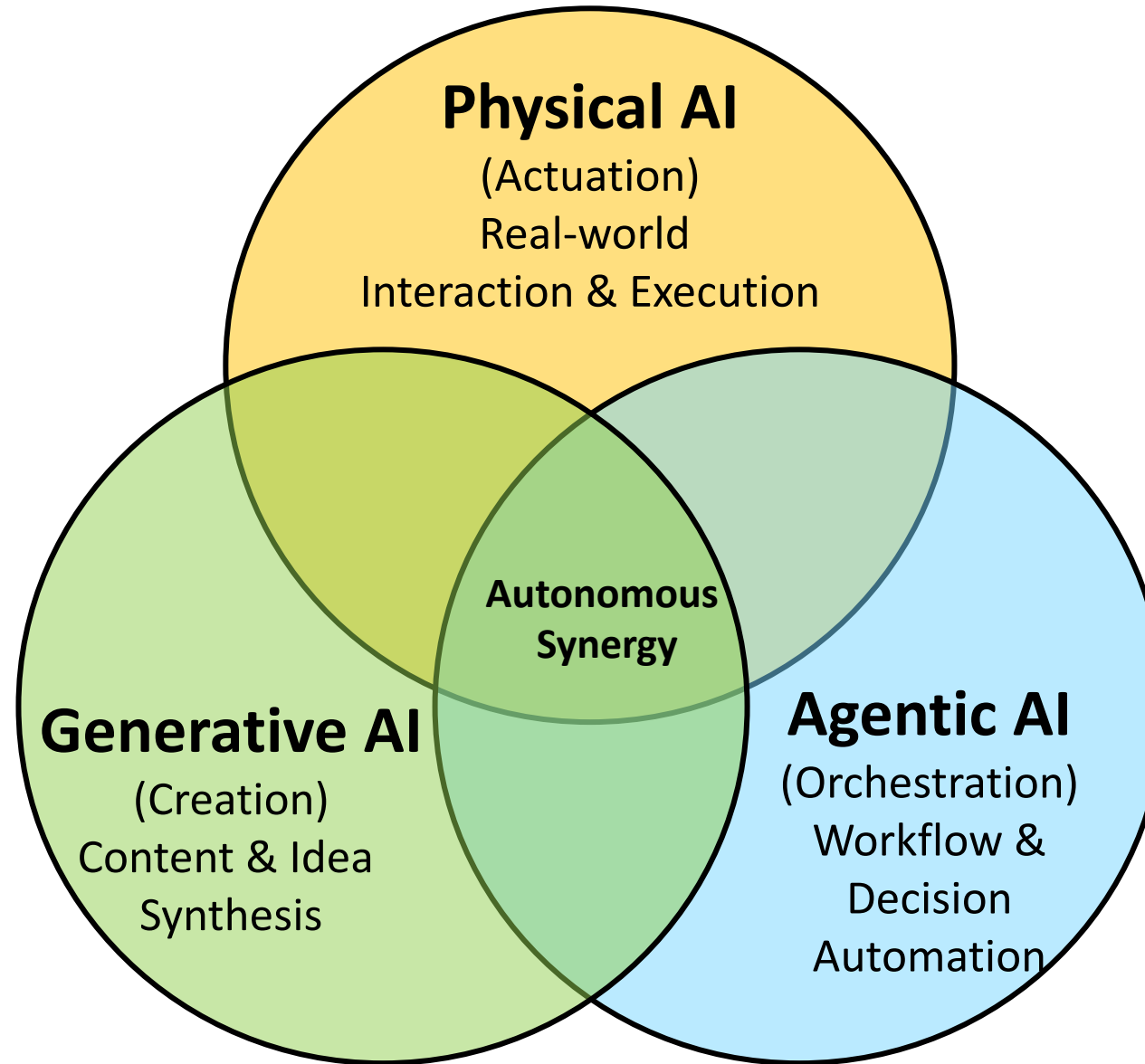Digital marketing
Content creation

**Perception AI**
Speech recognition
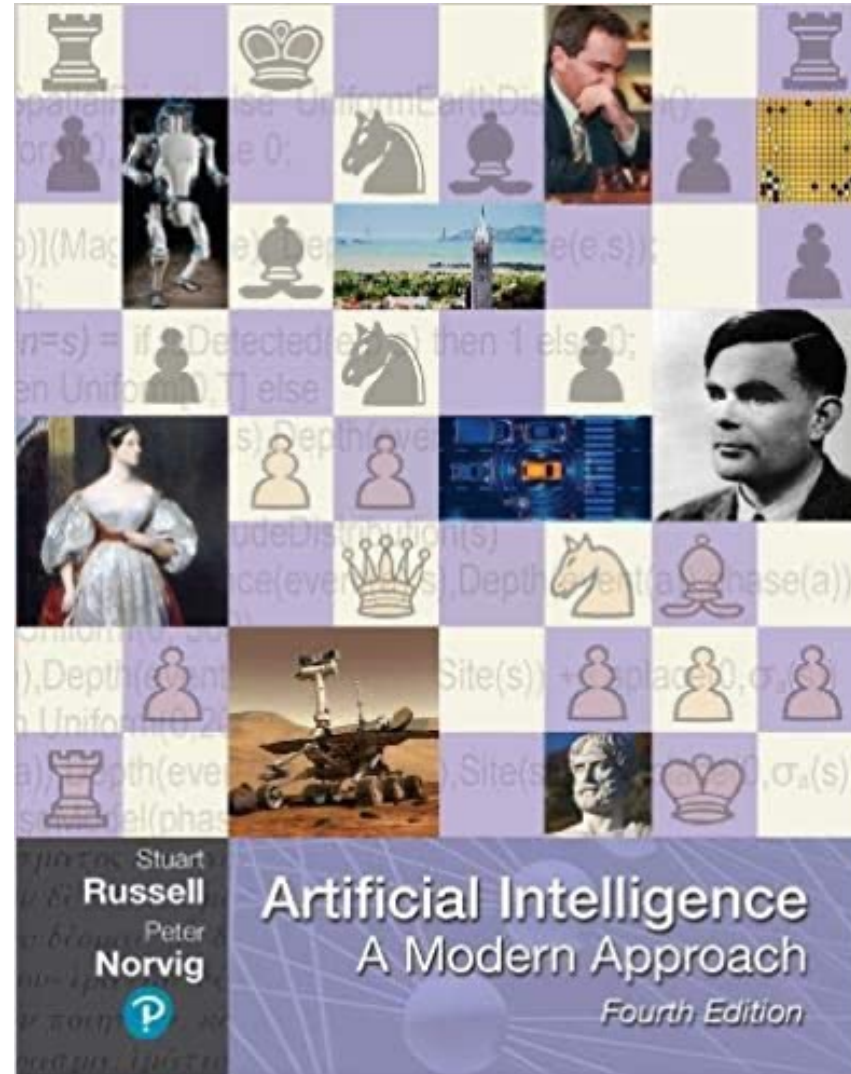Deep recommender systems
Medical imaging

**2012 AlexNet**
Deep learning breakthrough

# Generative AI, Agentic AI, Physical AI



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

**Physical AI**
(Actuation)
Real-world
Interaction & Execution

**Autonomous Synergy**

**Generative AI**
(Creation)
Content & Idea
Synthesis

**Agentic AI**
(Orchestration)
Workflow &
Decision
Automation

# Stuart Russell and Peter Norvig (2020),
# Artificial Intelligence: A Modern Approach,
## 4th Edition, Pearson



Source: Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson

# Artificial Intelligence:
# A Modern Approach

1. **Artificial Intelligence**

2. **Problem Solving**

3. **Knowledge and Reasoning**

4. **Uncertain Knowledge and Reasoning**

5. **Machine Learning**

6. **Communicating, Perceiving, and Acting**

7. **Philosophy and Ethics of AI**

# Artificial Intelligence: Communicating, perceiving, and acting

# Artificial Intelligence:
## 6. Communicating, Perceiving, and Acting

- **Natural Language Processing**

- **Deep Learning for Natural Language Processing**
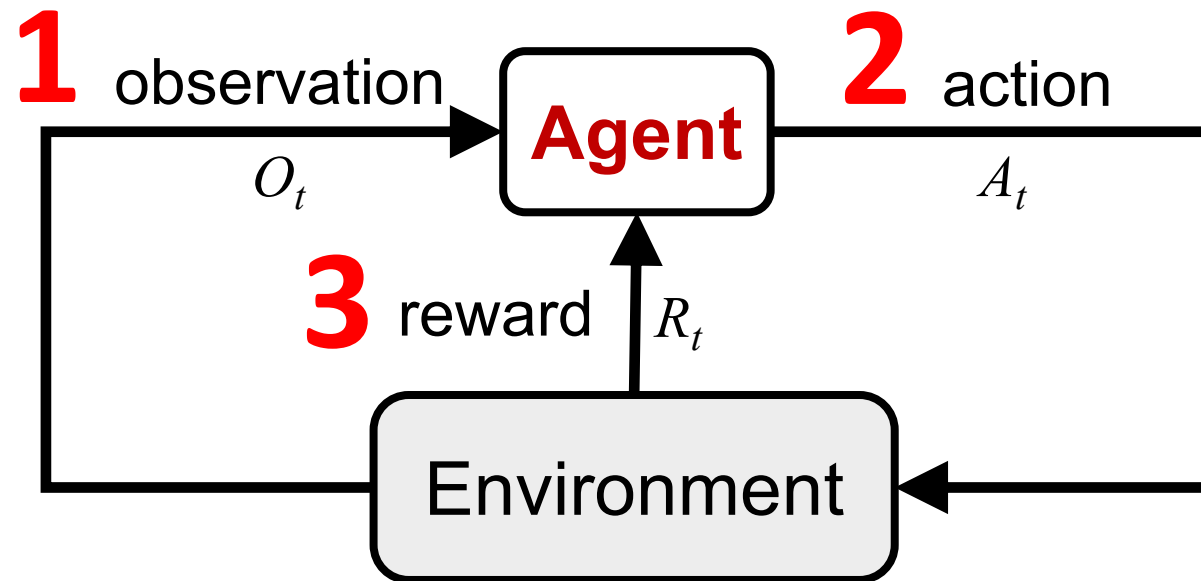
- **Computer Vision**

- **Robotics**

# Reinforcement Learning (DL)

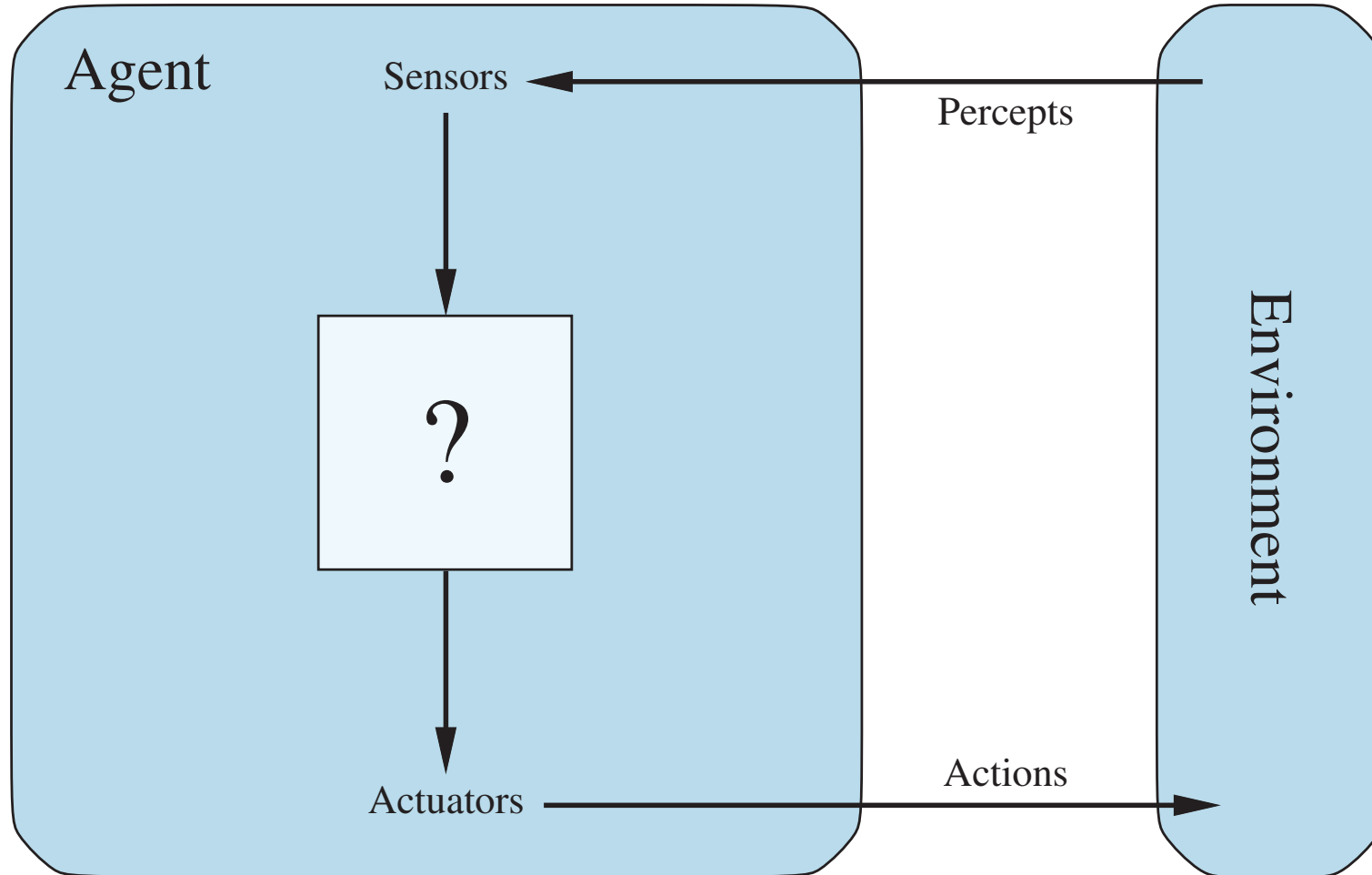# Reinforcement Learning (DL)

# Reinforcement Learning (DL)

# Agents interact with environments through sensors and actuators

# AI Acting Humanly:
# The Turing Test Approach
## (Alan Turing, 1950)

- **Knowledge Representation**

- **Automated Reasoning**

- **Machine Learning (ML)**

  - **Deep Learning (DL)**

- **Computer Vision (Image, Video)**

- **Natural Language Processing (NLP)**
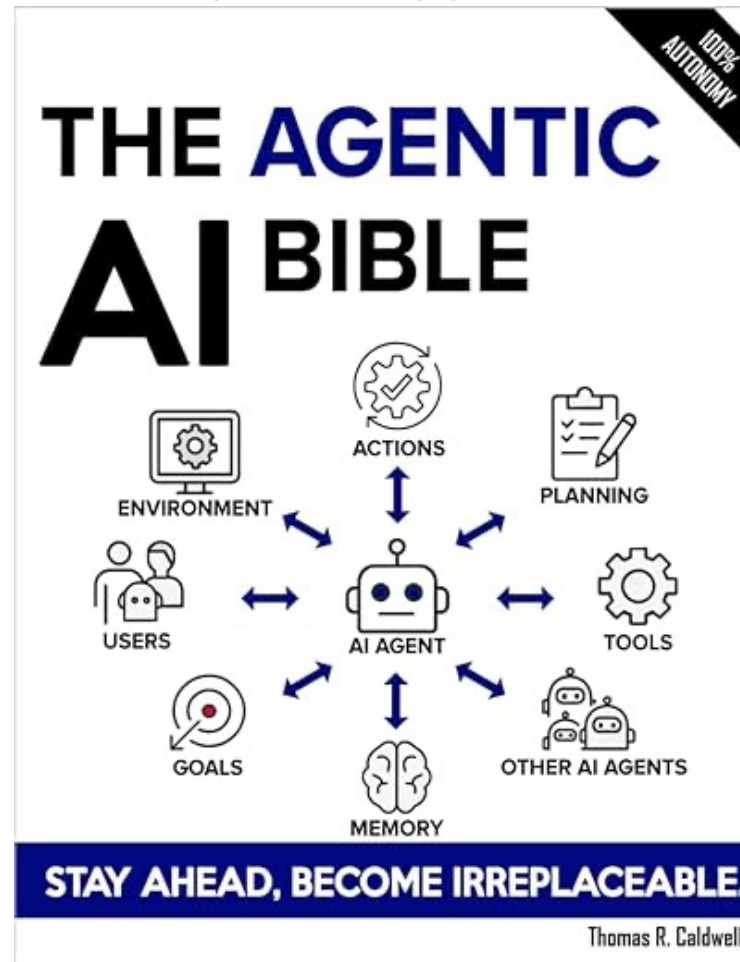
- **Robotics**

# 4 Approaches of AI

| | |
|---|---|
| **2.**<br>**Thinking Humanly:**<br>**The Cognitive Modeling Approach** | **3.**<br>**Thinking Rationally:**<br>**The "Laws of Thought" Approach** |
| **1.**<br>**Acting Humanly:**<br>**The Turing Test Approach** **(1950)** | **4.**<br>**Acting Rationally:**<br>**The Rational Agent Approach** |

**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**

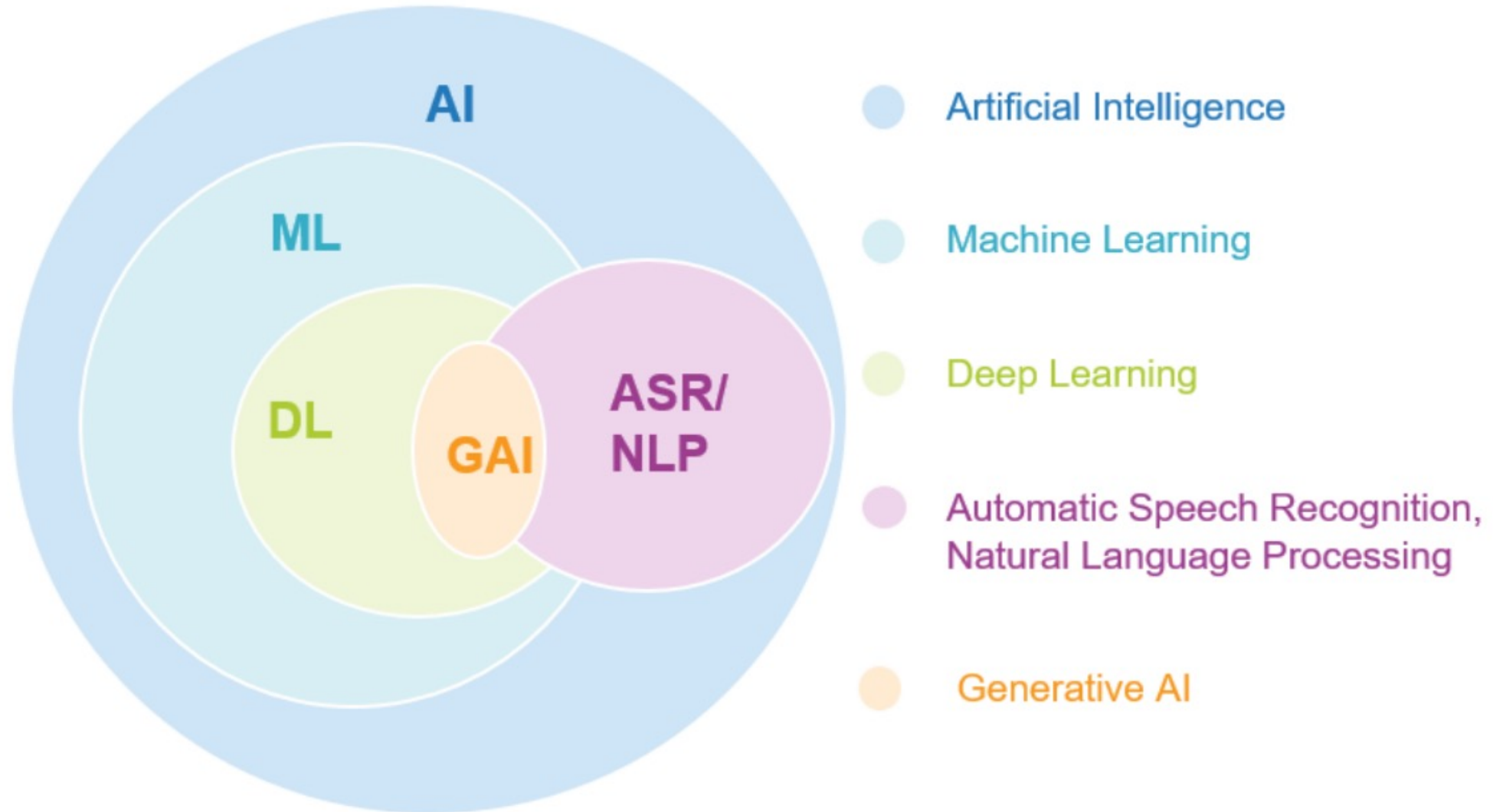# Generative AI-Driven ESG Report Generation Technology

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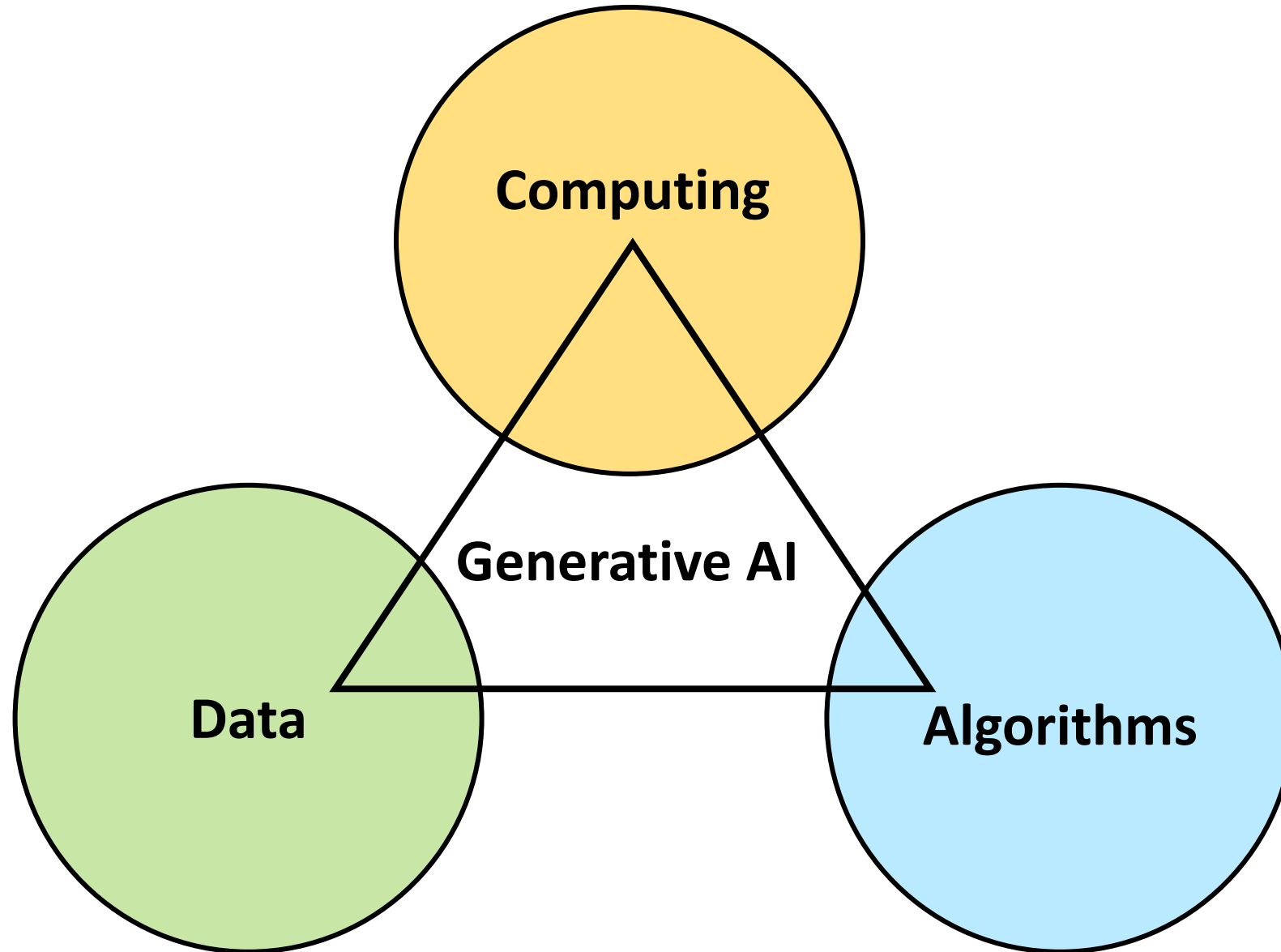NTPU-113A513E01, 2024/03/01~2024/12/31

# Innovative Agentic AI Technology for Autonomous ESG Report Generation

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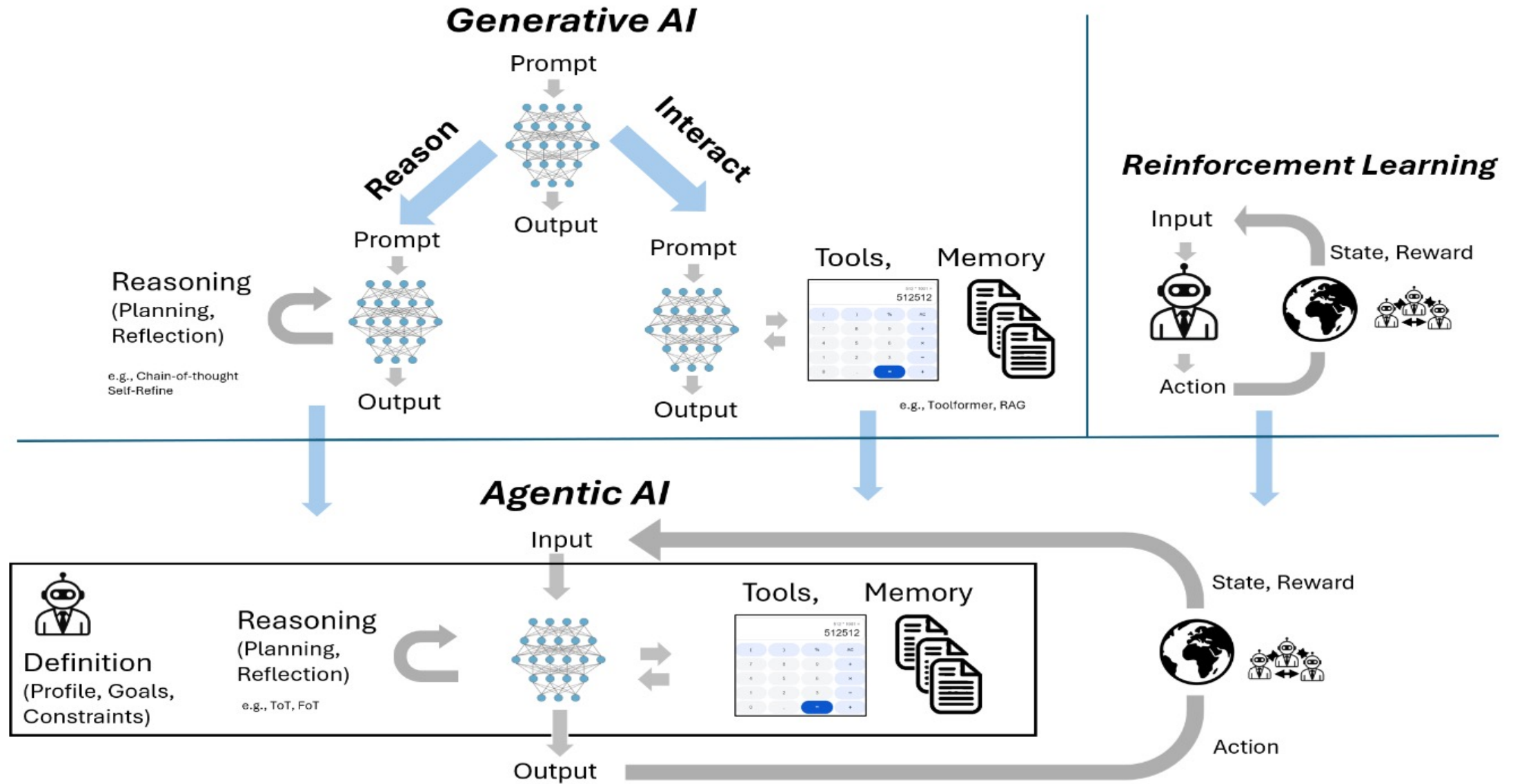NTPU-114A513E01, 2025/03/01~2025/12/31

# AI, ML, DL, Generative AI



Source: Jeong, Cheonsu. "A Study on the Implementation of Generative AI Services Using an Enterprise Data-Based LLM Application Architecture." arXiv preprint arXiv:2309.01105 (2023).
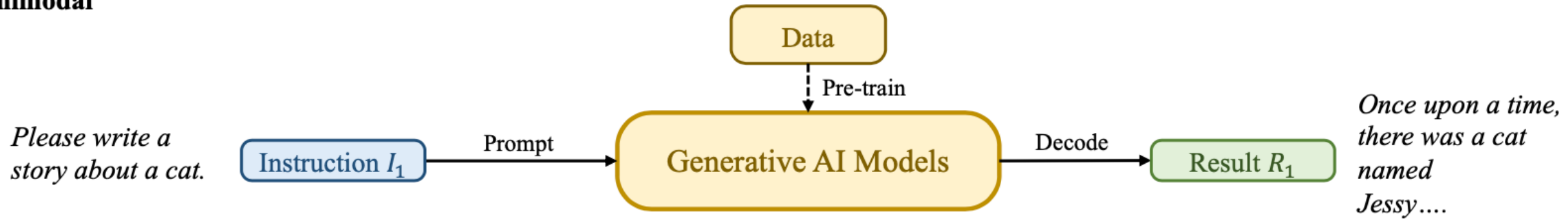
22

# Generative AI

# From Generative AI to Agentic AI

# Generative AI
# Text, Image, Video, Audio Applications
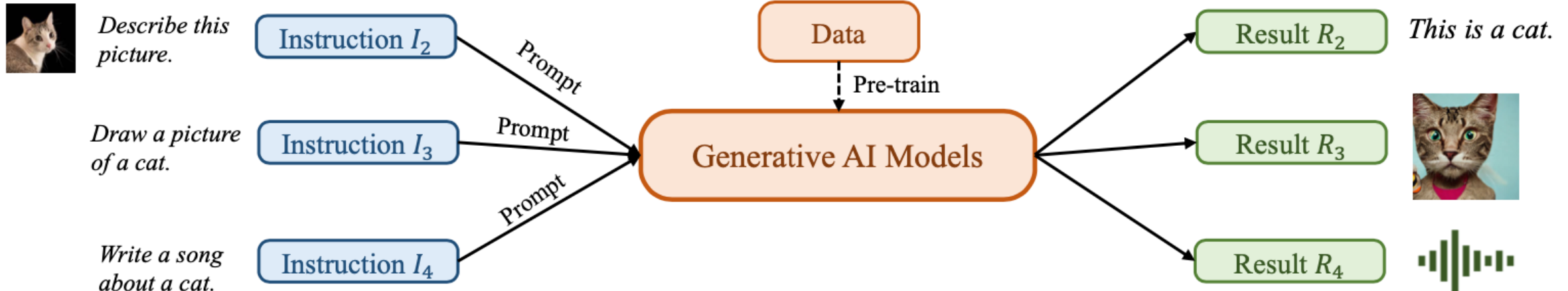
# Generative AI (Gen AI)
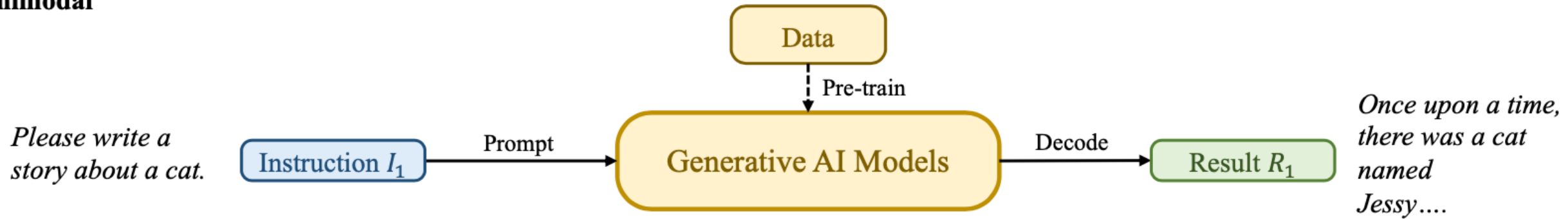## AI Generated Content (AIGC)
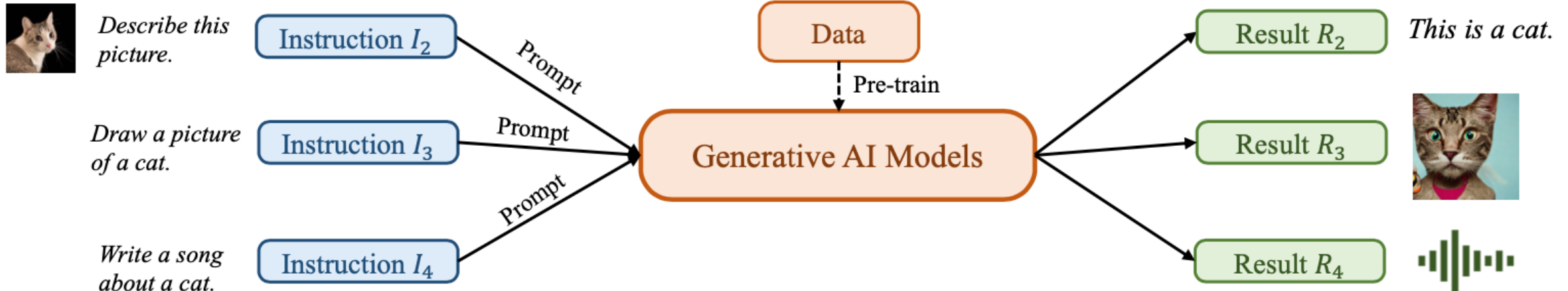
# Generative AI (Gen AI)
## AI Generated Content (AIGC)
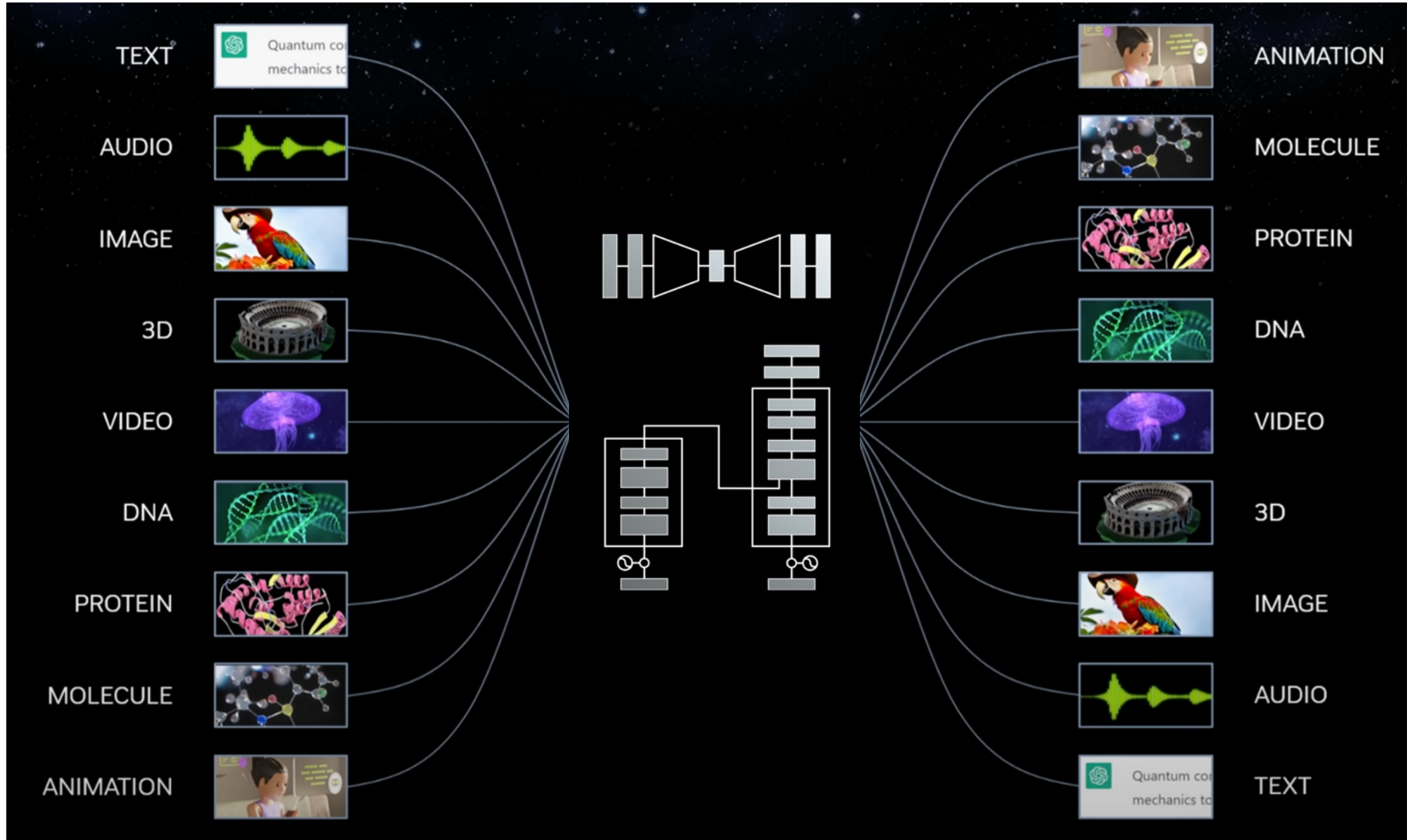
# Modular Modalities
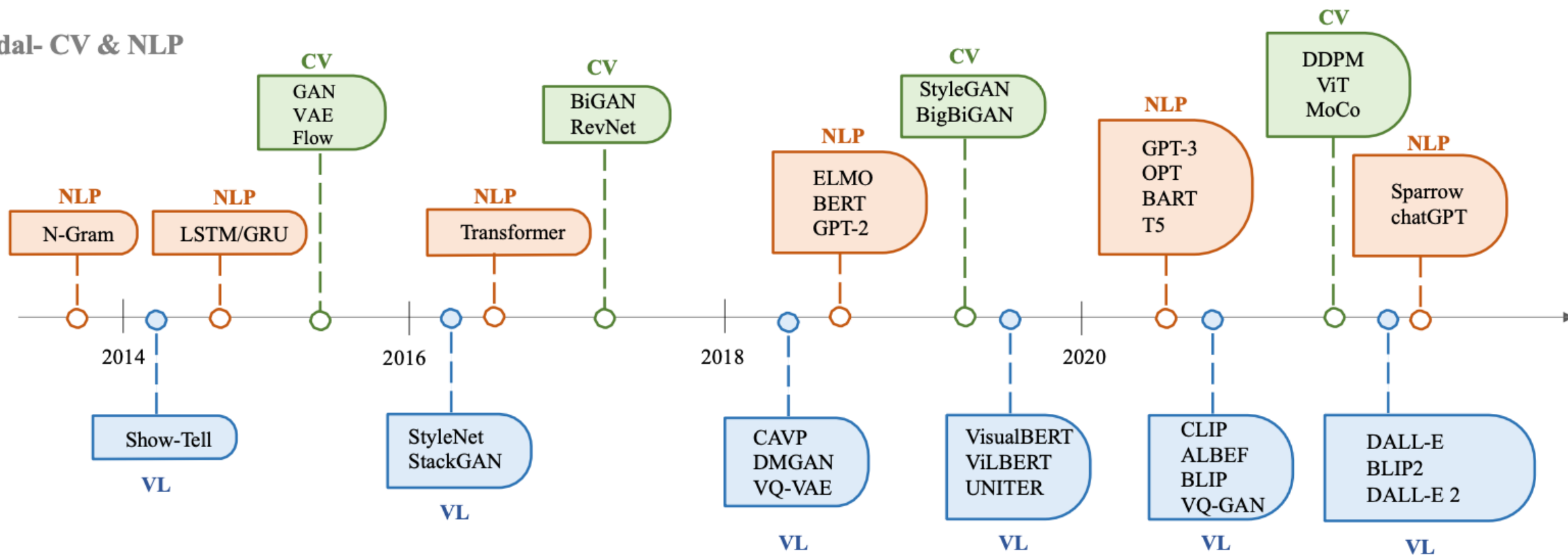## Where Can The Transformer Fit?

# 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

# Transformer Models

32

# Large Language Models (LLMs)

Source: Hadi, Muhammad Usman, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, and Seyedali Mirjalili. "Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects." Authorea Preprints (2023).

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

# 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 |

# Generative AI
# Text, Image, Video, Audio Applications

# 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)



Multimodall LLM
Three types of connectors:
1. projection-based
2. query-based
3. fusion-based connectors

Source: Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. (2024) "A survey on multimodal large language models." National Science Review (2024): nwae403.

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

# Large Language Models (LLM)
# Three typical learning paradigms



**(A) Pretrain-finetune**
- Typically requires many task-specific examples
- One specialized model for each task

**Finetune on task A** → **Inference on task A**

**Pretrained LM**

**(B) Prompting**

Few-shot prompting / prompt engineering

**Inference on task A**

**(C) Instruction tuning**

Model learns to perform many tasks via natural language instructions

**Instruction-tune on many tasks: B, C, D, …** → **Inference on unseen task A**

Source: Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. (2024) "A survey on multimodal large language models." National Science Review (2024): nwae403.

40

# Conversational AI
## to deliver contextual and personal experience to users



Source: Huynh-The, Thien, Quoc-Viet Pham, Xuan-Qui Pham, Thanh Thi Nguyen, Zhu Han, and Dong-Seong Kim (2022). "Artificial Intelligence for the Metaverse: A Survey." arXiv preprint arXiv:2202.10336.

41

# Technological Integration for Multimodal AI

# AutoDev: Automated AI-Driven Development

Source: Tufano, Michele, Anisha Agarwal, Jinu Jang, Roshanak Zilouchian Moghaddam, and Neel Sundaresan. (2024) "AutoDev: Automated AI-Driven Development." arXiv preprint arXiv:2403.08299 (2024).

# Framework for Implementing Generative AI Services using RAG Model

# 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

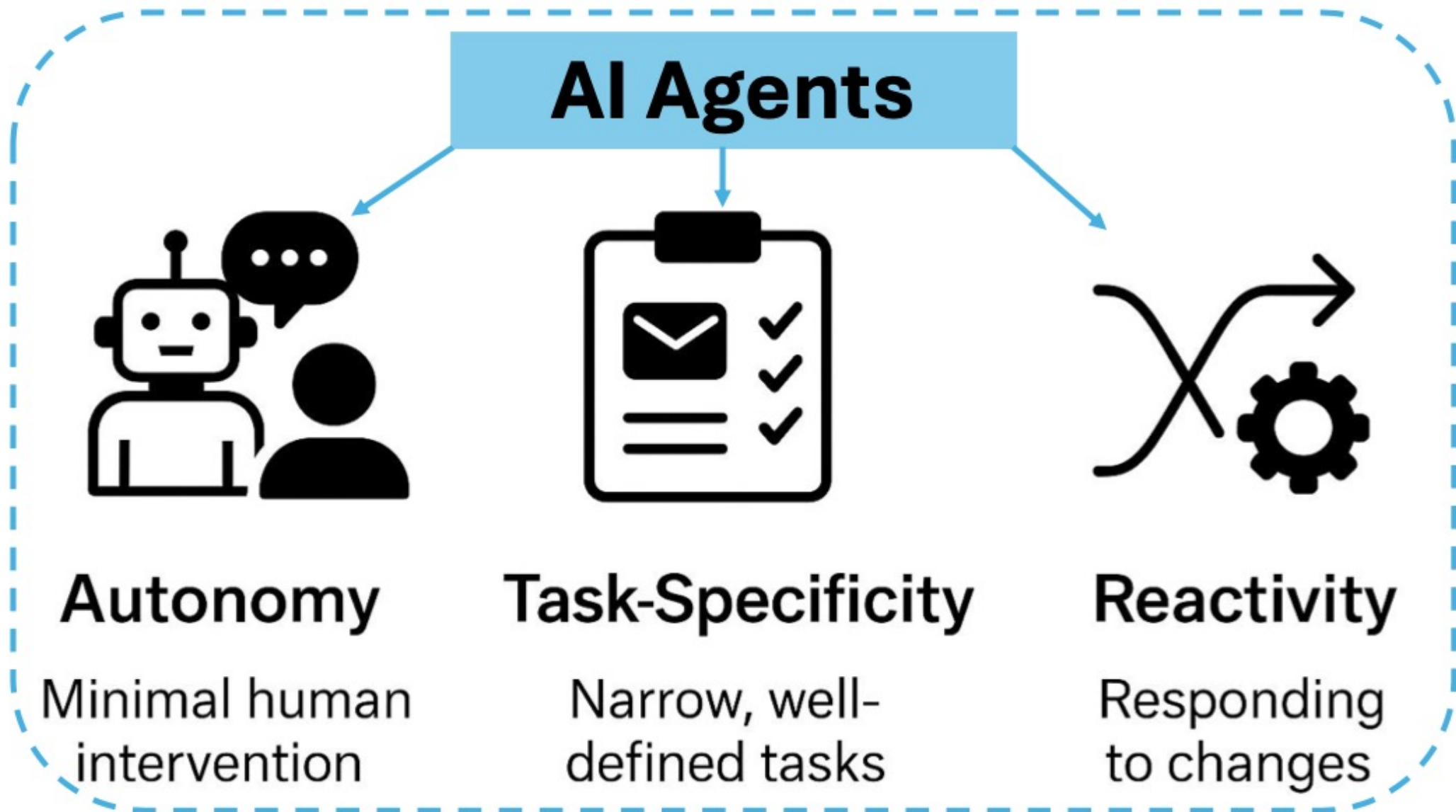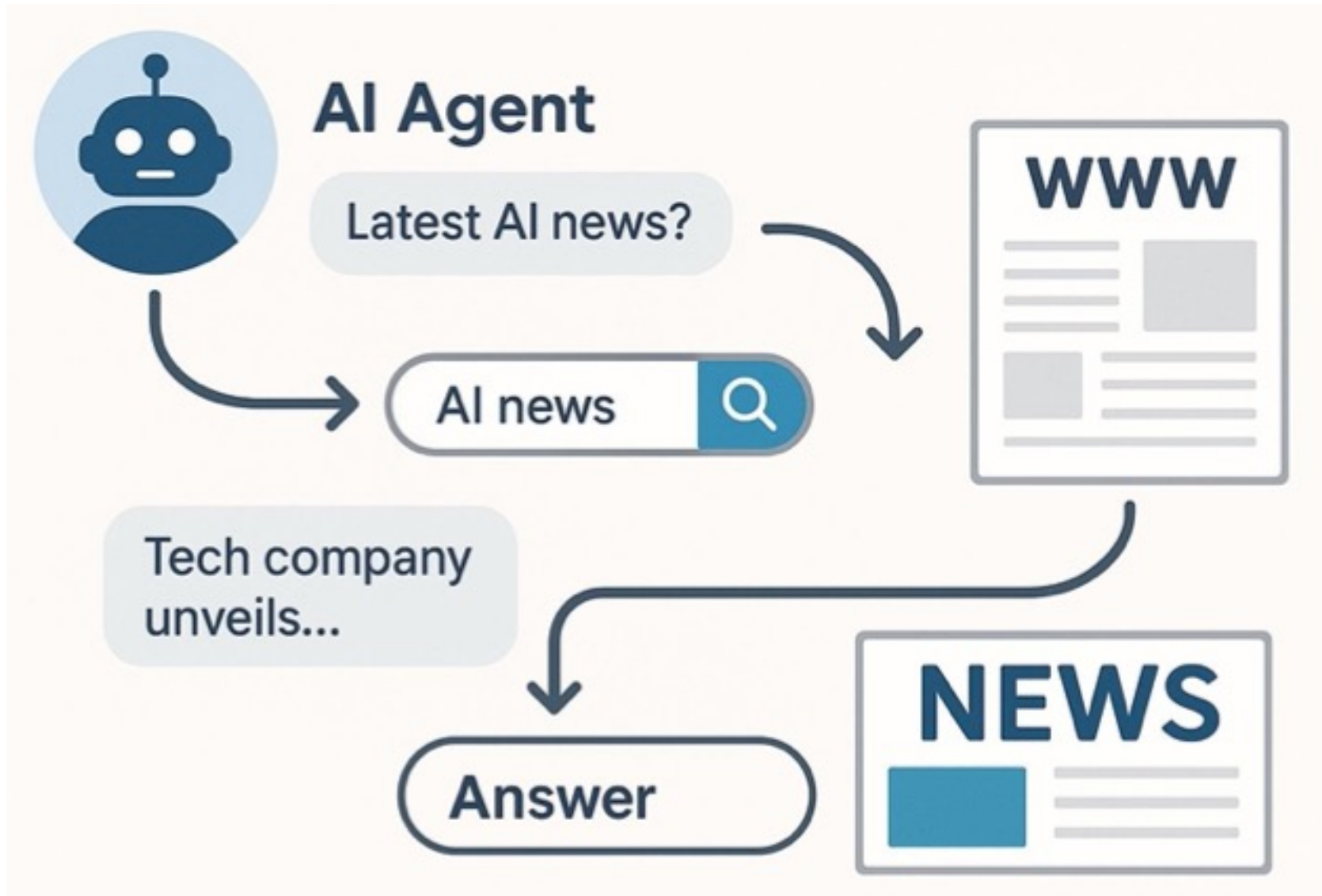| Feature | AI Agent / Agentic AI | Generative AI | Traditional AI |
|---|---|---|---|
| Core Concept | To autonomously perceive its environment, make decisions, and take actions to achieve specific goals. | To create new, original content (text, images, code, etc.) that resembles its training data. | To execute specific tasks based on pre-programmed rules or statistical patterns. |
| Primary Function | 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?"). |
| Decision Making | 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. |
| Key Characteristic | 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. |

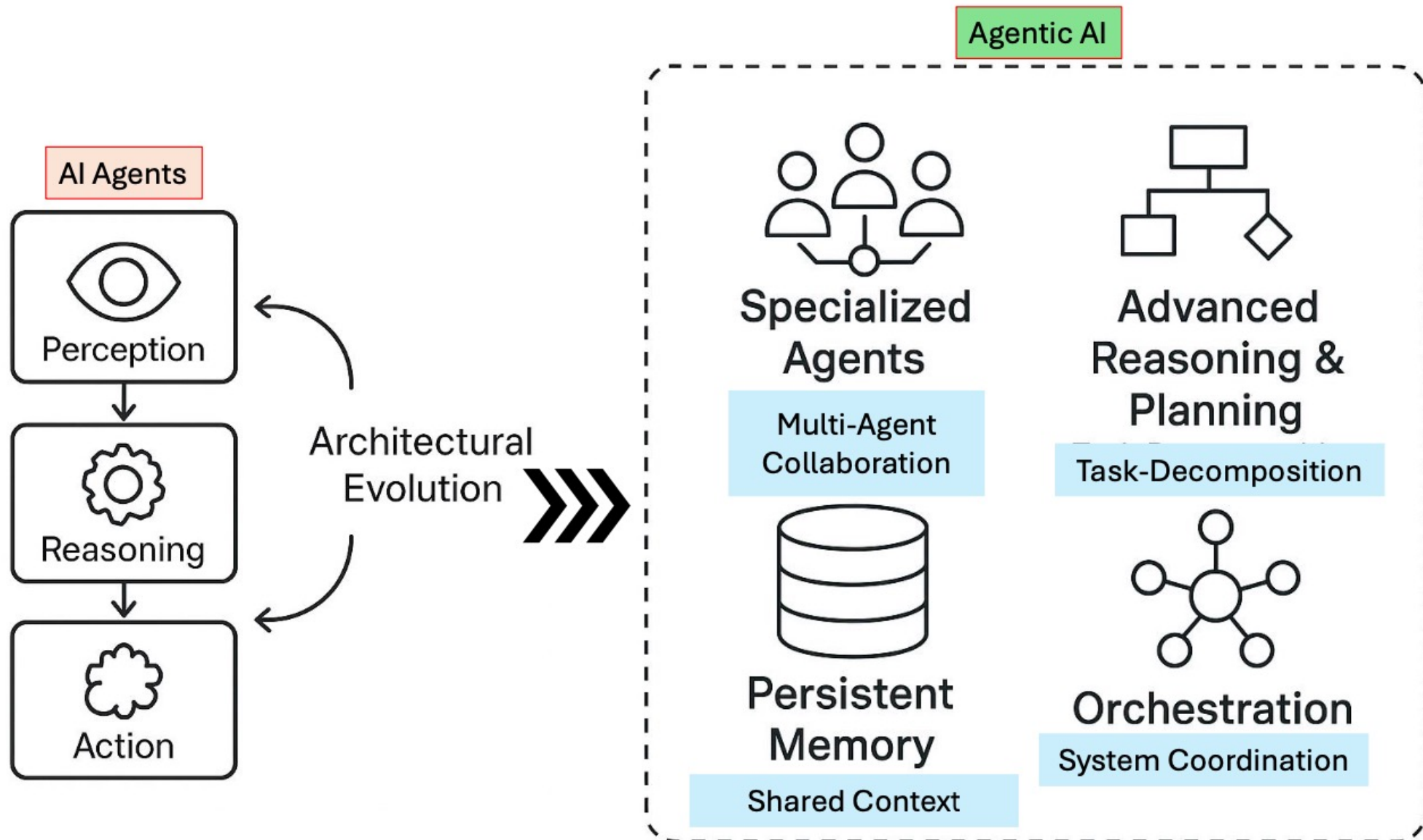# AI Agent / Agentic AI, Generative AI, Traditional AI

| Feature | AI Agent / Agentic AI | Generative AI | Traditional AI |
|---|---|---|---|
| Interaction Model | 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. |
| Example Technologies | 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. |
| Common Use Cases | 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. |
| Relationship to Others | 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

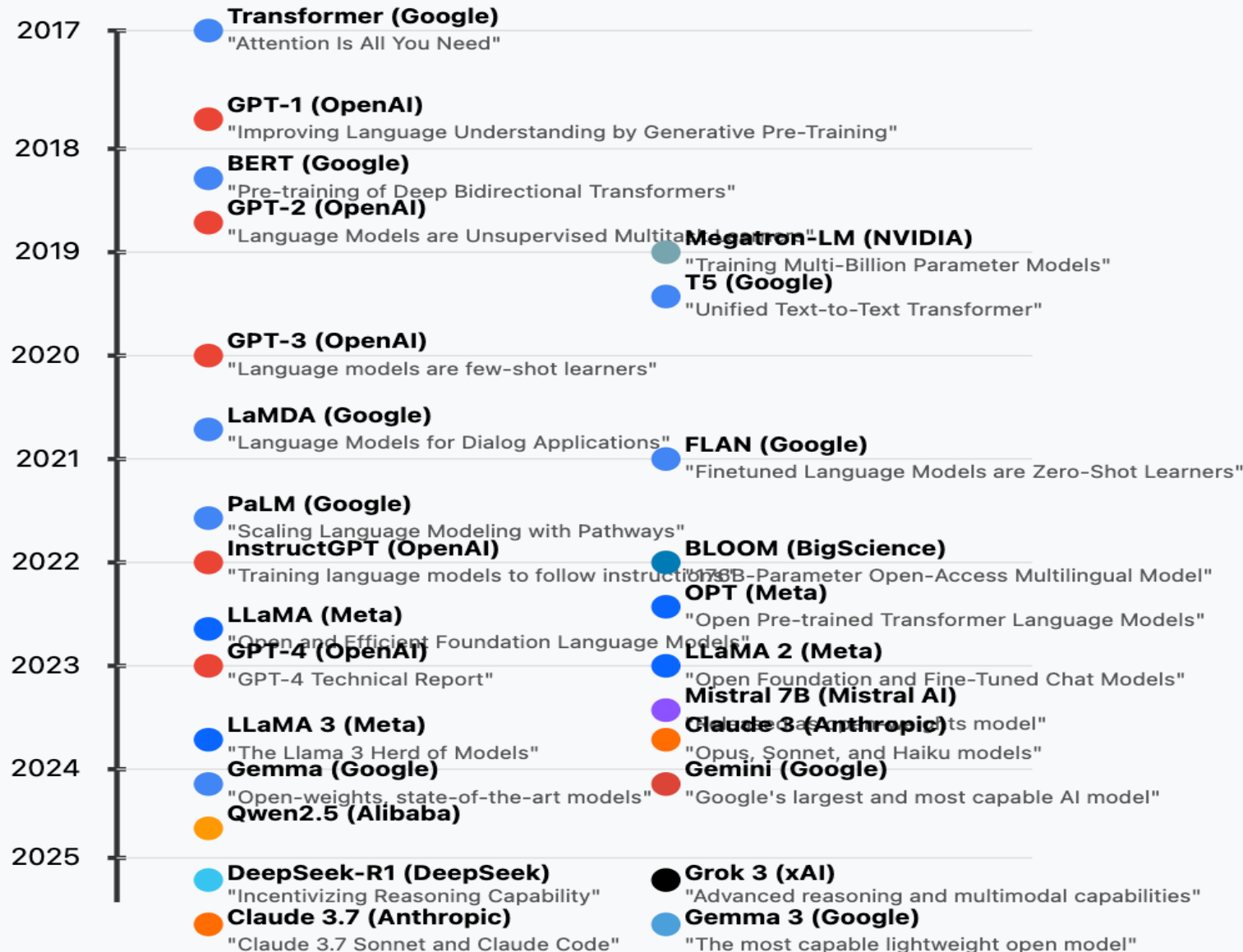| Feature | AI Agents | Agentic AI |
|---|---|---|
| **Definition** | Autonomous software programs that perform specific tasks. | Systems of multiple AI agents collaborating to achieve complex goals. |
| **Autonomy Level** | High autonomy within specific tasks. | Broad level of autonomy with the ability to manage multi-step, complex tasks and systems. |
| **Task Complexity** | Typically handle single, specific tasks. | Handle complex, multi-step tasks requiring coordination. |
| **Collaboration** | Operate independently. | Involve multi-agent information sharing, collaboration and cooperation. |
| **Learning and Adaptation** | Learn and adapt within their specific domain. | Learn and adapt across a wider range of tasks and environments. |
| **Applications** | 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)



**Timeline**

**2017**
- **Transformer (Google)** — "Attention Is All You Need"

**2018**
- **GPT-1 (OpenAI)** — "Improving Language Understanding by Generative Pre-Training"
- **BERT (Google)** — "Pre-training of Deep Bidirectional Transformers"
- **GPT-2 (OpenAI)** — "Language Models are Unsupervised Multitask Learners"

**2019**
- **Megatron-LM (NVIDIA)** — "Training Multi-Billion Parameter Models"
- **T5 (Google)** — "Unified Text-to-Text Transformer"

**2020**
- **GPT-3 (OpenAI)** — "Language models are few-shot learners"

**2021**
- **LaMDA (Google)** — "Language Models for Dialog Applications"
- **FLAN (Google)** — "Finetuned Language Models are Zero-Shot Learners"

**2022**
- **PaLM (Google)** — "Scaling Language Modeling with Pathways"
- **InstructGPT (OpenAI)** — "Training language models to follow instructions"
- **BLOOM (BigScience)** — "176B-Parameter Open-Access Multilingual Model"
- **OPT (Meta)** — "Open Pre-trained Transformer Language Models"

**2023**
- **LLaMA (Meta)** — "Open and Efficient Foundation Language Models"
- **GPT-4 (OpenAI)** — "GPT-4 Technical Report"
- **LLaMA 2 (Meta)** — "Open Foundation and Fine-Tuned Chat Models"
- **Mistral 7B (Mistral AI)** — "Grouped-query attention model"
- **Claude 3 (Anthropic)** — "Opus, Sonnet, and Haiku models"

**2024**
- **LLaMA 3 (Meta)** — "The Llama 3 Herd of Models"
- **Gemma (Google)** — "Open-weights, state-of-the-art models"
- **Gemini (Google)** — "Google's largest and most capable AI model"
- **Qwen2.5 (Alibaba)**

**2025**
- **DeepSeek-R1 (DeepSeek)** — "Incentivizing Reasoning Capability"
- **Grok 3 (xAI)** — "Advanced reasoning and multimodal capabilities"
- **Claude 3.7 (Anthropic)** — "Claude 3.7 Sonnet and Claude Code"
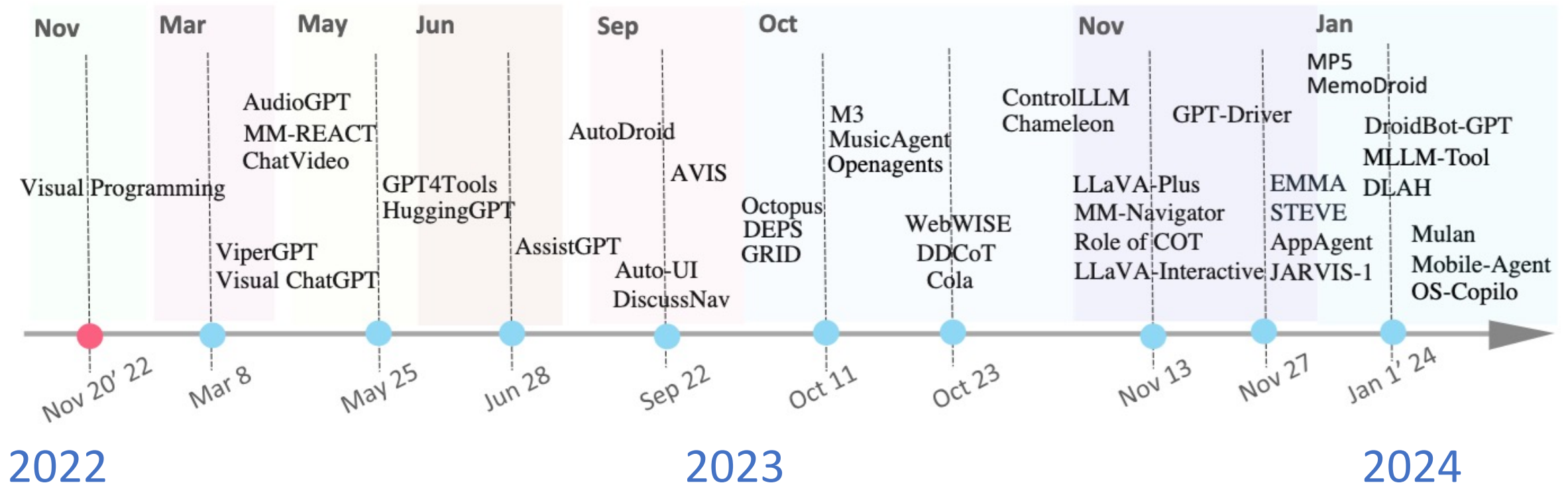- **Gemma 3 (Google)** — "The most capable lightweight open model"

**Key Organizations**
- Google
- OpenAI
- Meta
- Mistral AI
- Alibaba
- xAI
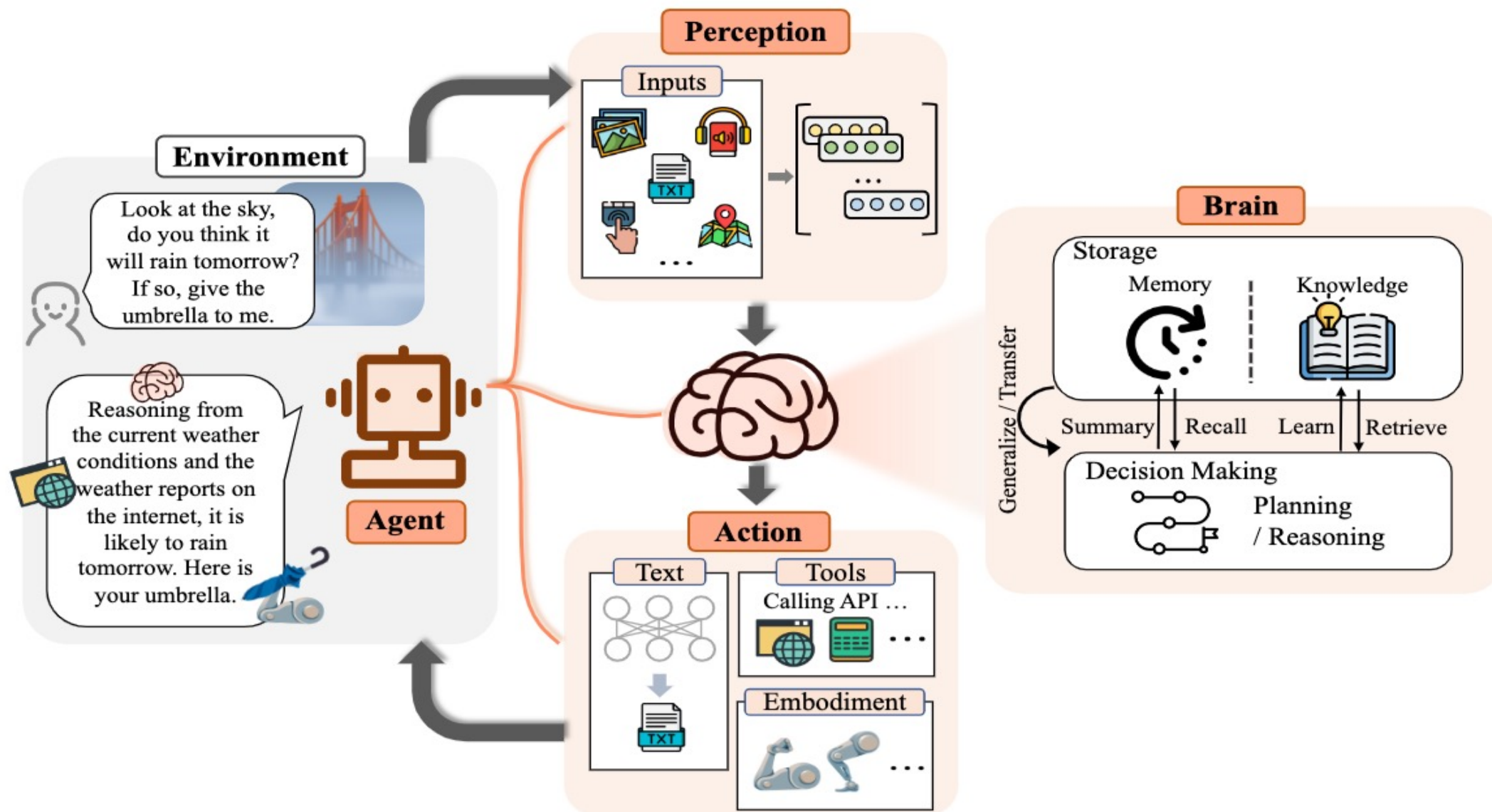- Anthropic
- NVIDIA
- BigScience

**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
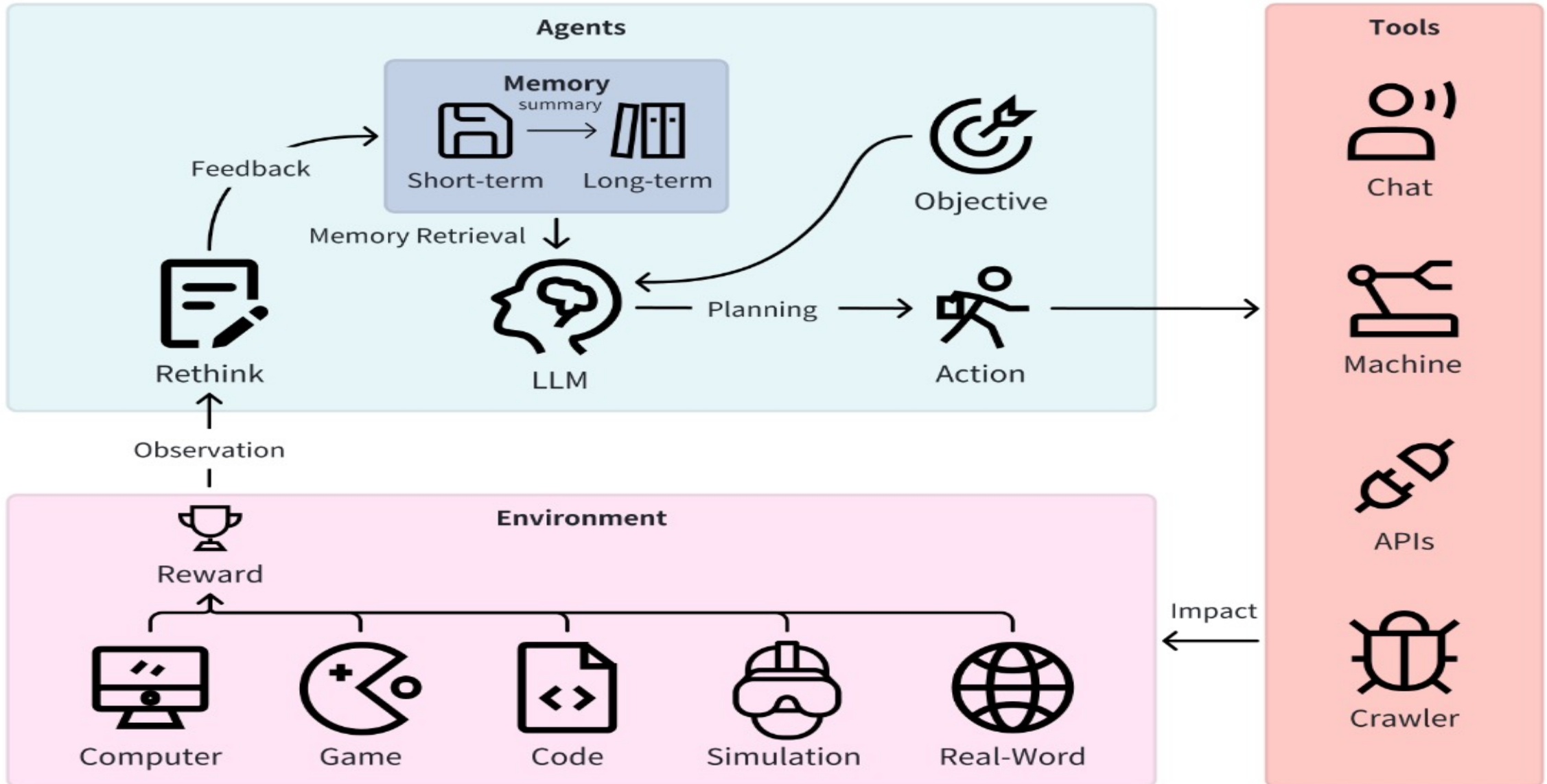# Large Multimodal Agents (LMAs)

Source: Xie, Junlin, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. "Large Multimodal Agents: A Survey." arXiv preprint arXiv:2402.15116 (2024).

# Large Language Model (LLM) based Agents
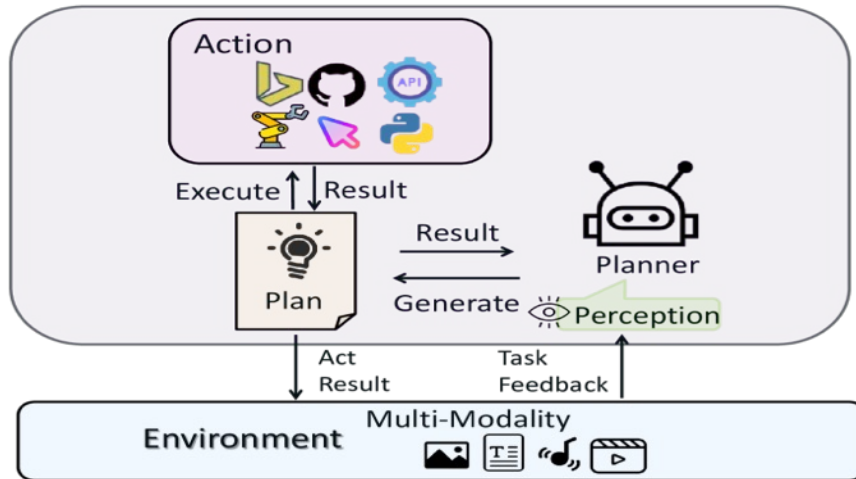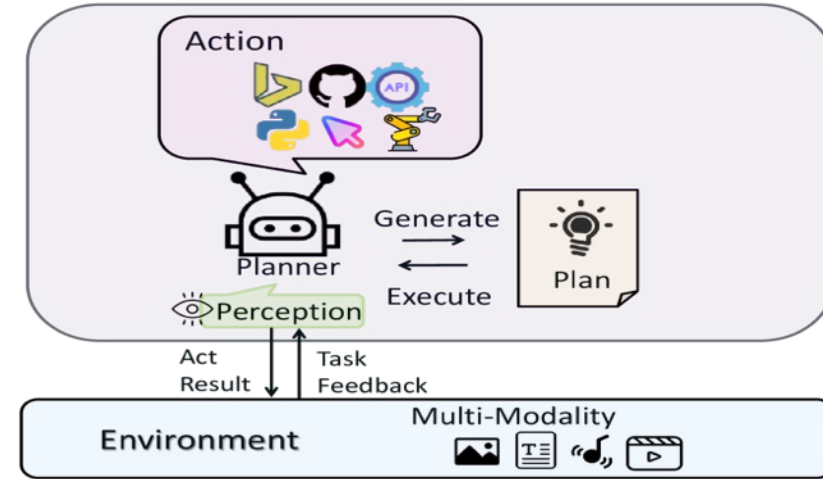
# LLM-based Agents

- **Definition: <span style="color:red">AI agents</span> that use <span style="color:red">Large Language Models</span> as their <span style="color:red">core decision-making</span> mechanism**

- **Key Features:**

  - **Natural language interface**

  - **Vast knowledge base**

  - **Ability to understand context and nuance**

  - **Generalize to new tasks with minimal additional training**

# LLM-based Agents

# Large Multimodal Agents (LMA)



(a)  (b)  (c)  (d)

Source: Xie, J., Chen, Z., Zhang, R., Wan, X., & Li, G. (2024). Large Multimodal Agents: A Survey. ArXiv, abs/2402.15116.

# 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
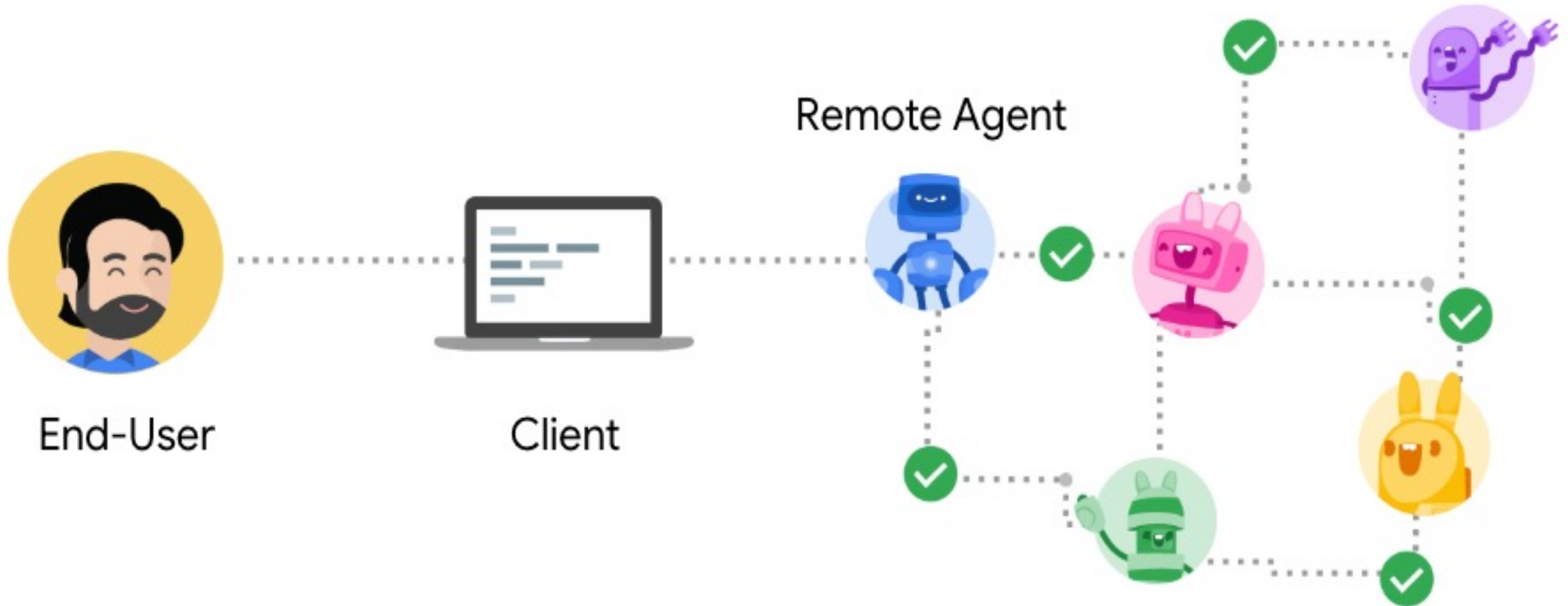
## Tools and Frameworks
LangChain, AutoGen, CrewAI: for orchestrating LLM agents
Anthropic's MCP, Google's A2A: communication protocols
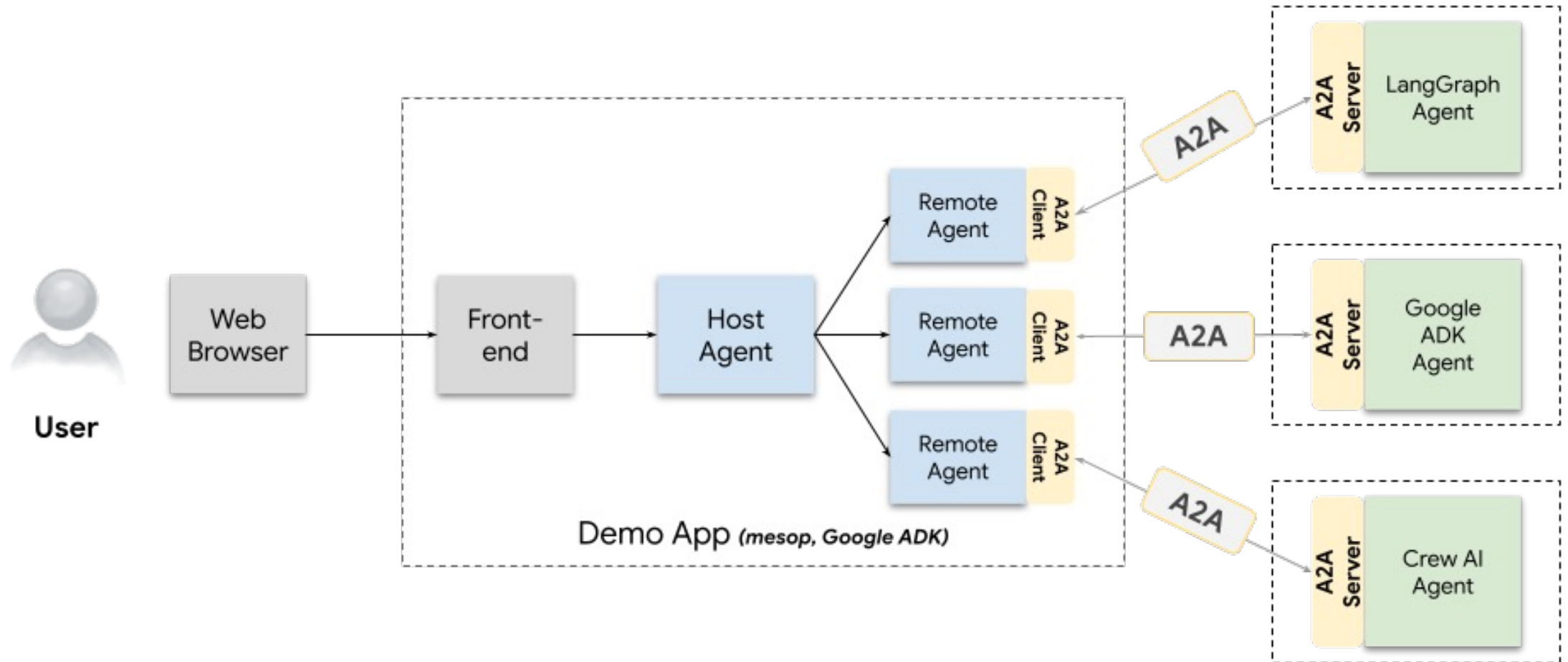Vector DBs (Pinecone, Weaviate): for agent memory

61

# Agent2Agent Protocol (A2A)

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

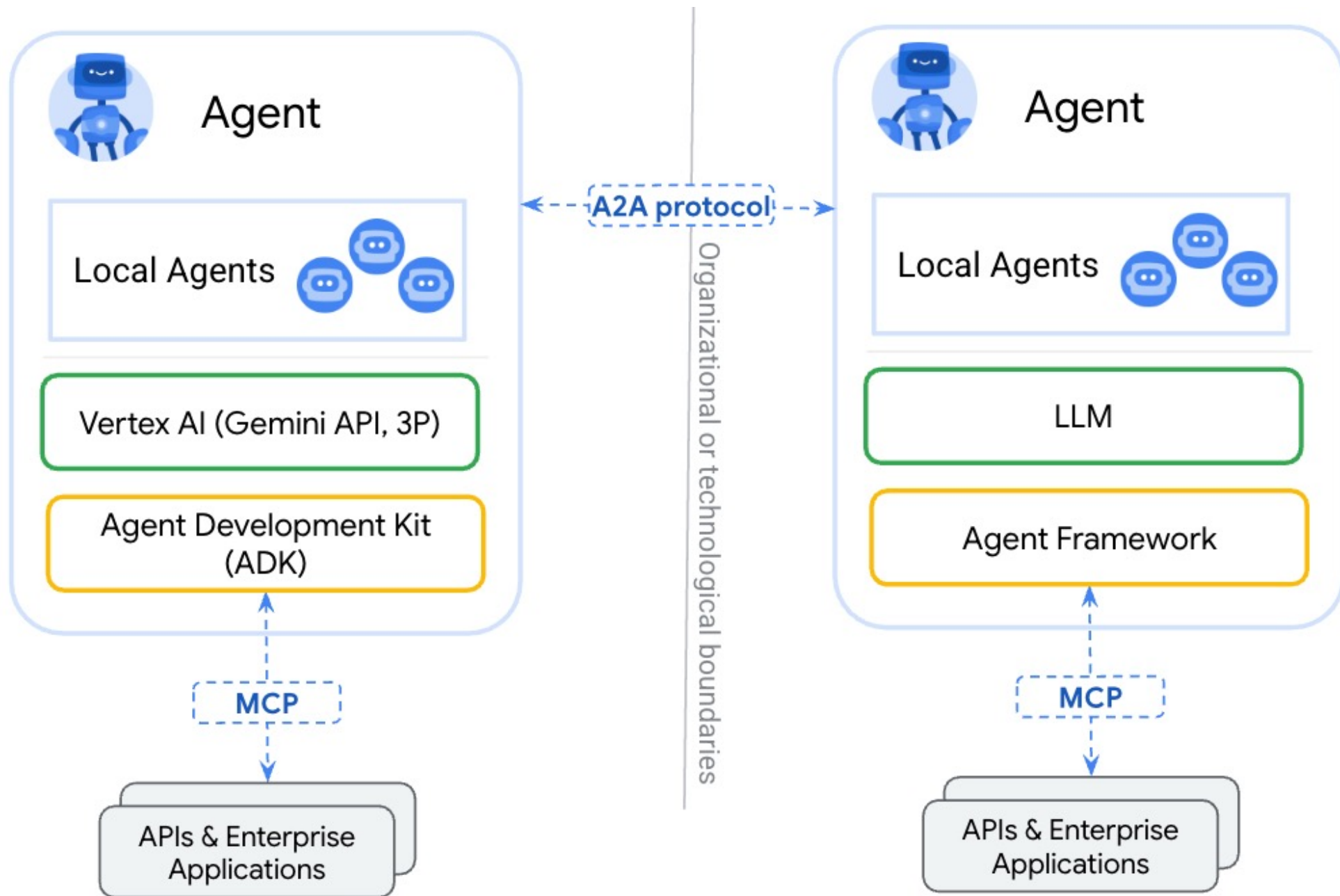# A2A Demo Web App
## Agents talking to other agents over A2A

**A2A**
**(Agent2Agent Protocol)**
for agent-agent collaboration

**MCP**
**(Model Context Protocol)**
for tools and resources
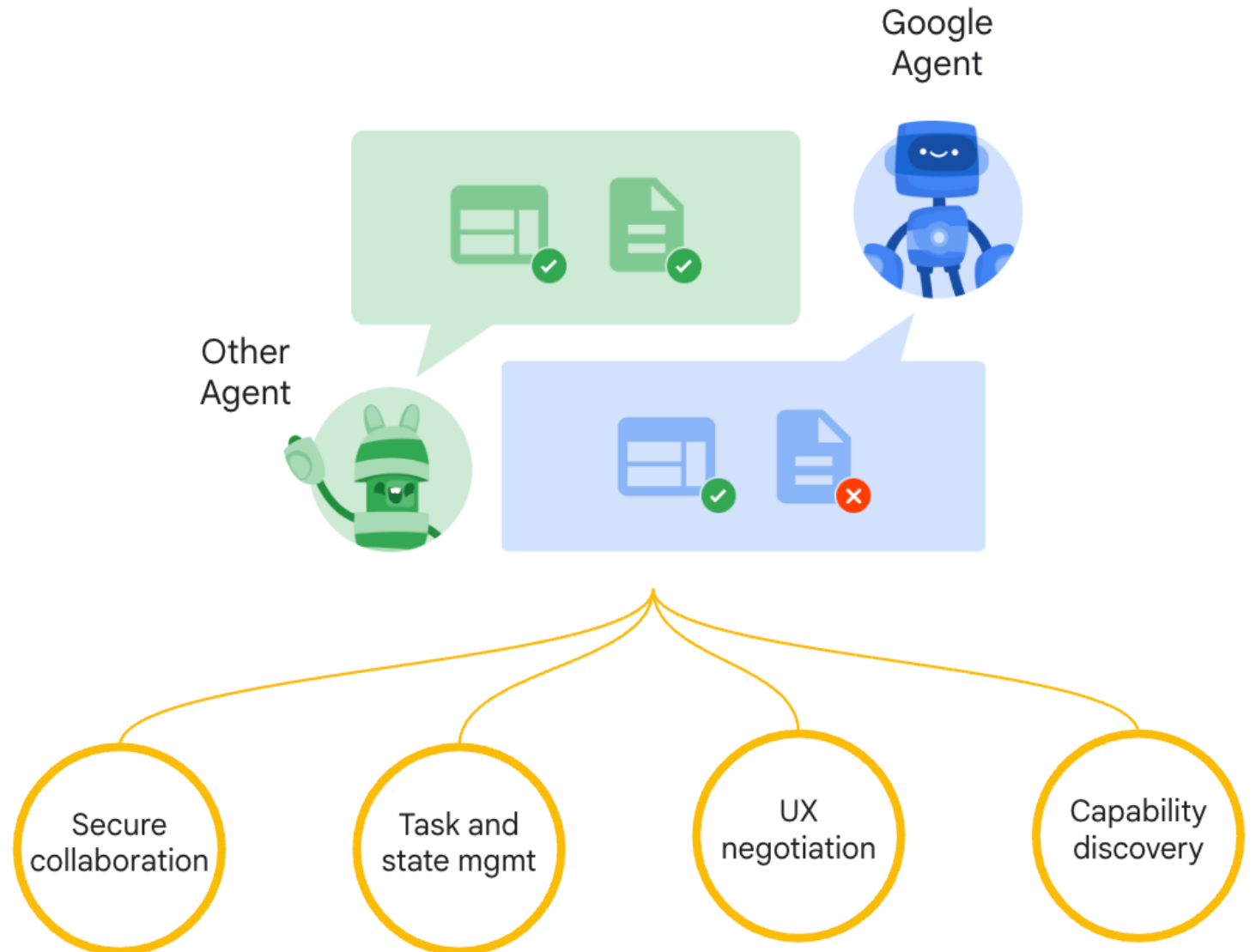
Source: https://google.github.io/A2A/

# Google A2A (Agent2Agent Protocol)

**Seamless Agent Collaboration**
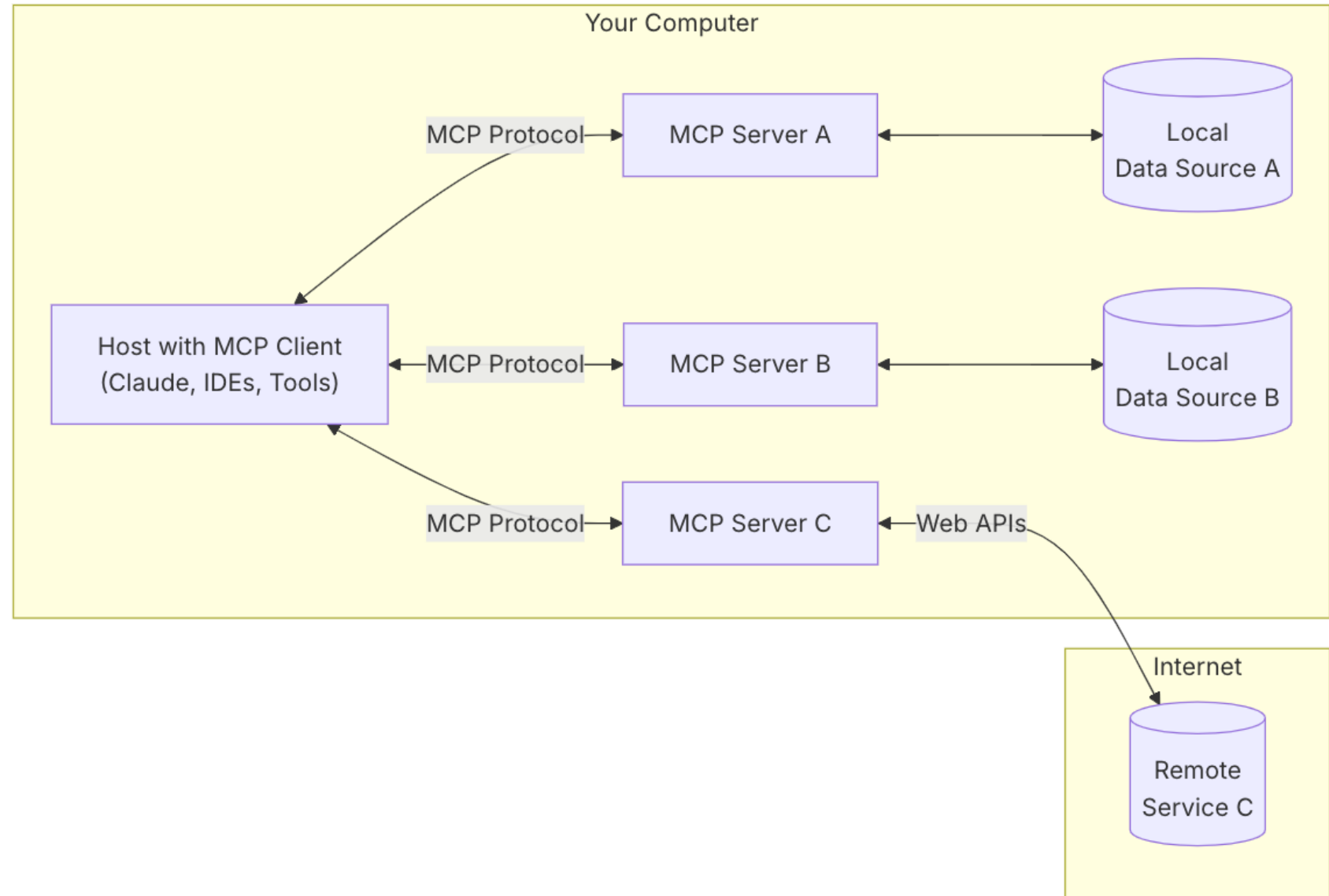
**Simplifies Enterprise Agent Integration**

**Supports Key Enterprise Requirements**

# MCP (Model Context Protocol)

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

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



Your Computer

MCP Protocol → MCP Server A ↔ Local Data Source A

Host with MCP Client (Claude, IDEs, Tools) ← MCP Protocol → MCP Server B ↔ Local Data Source B

MCP Protocol → MCP Server C ← Web APIs

Internet

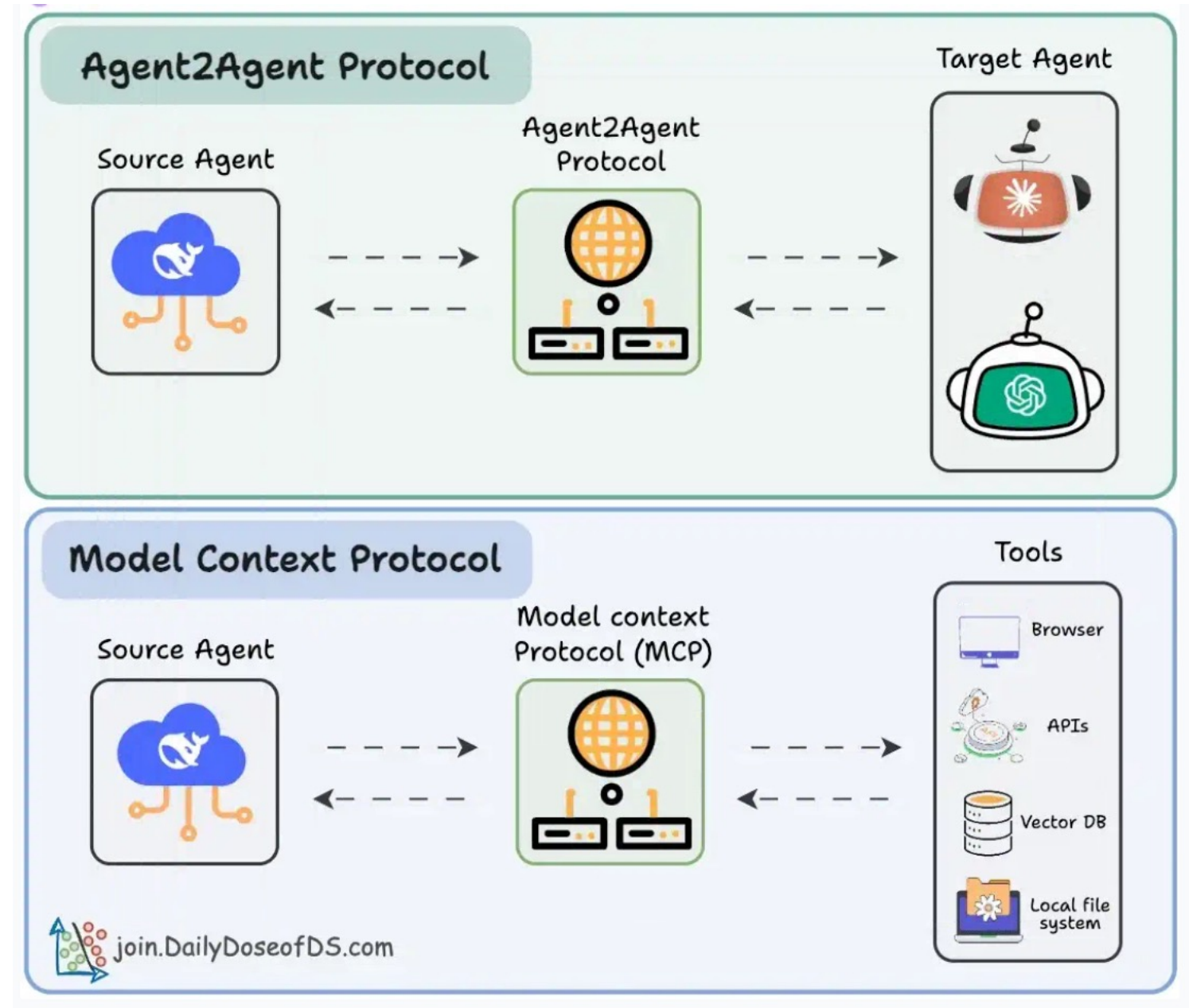Remote Service C

# 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**
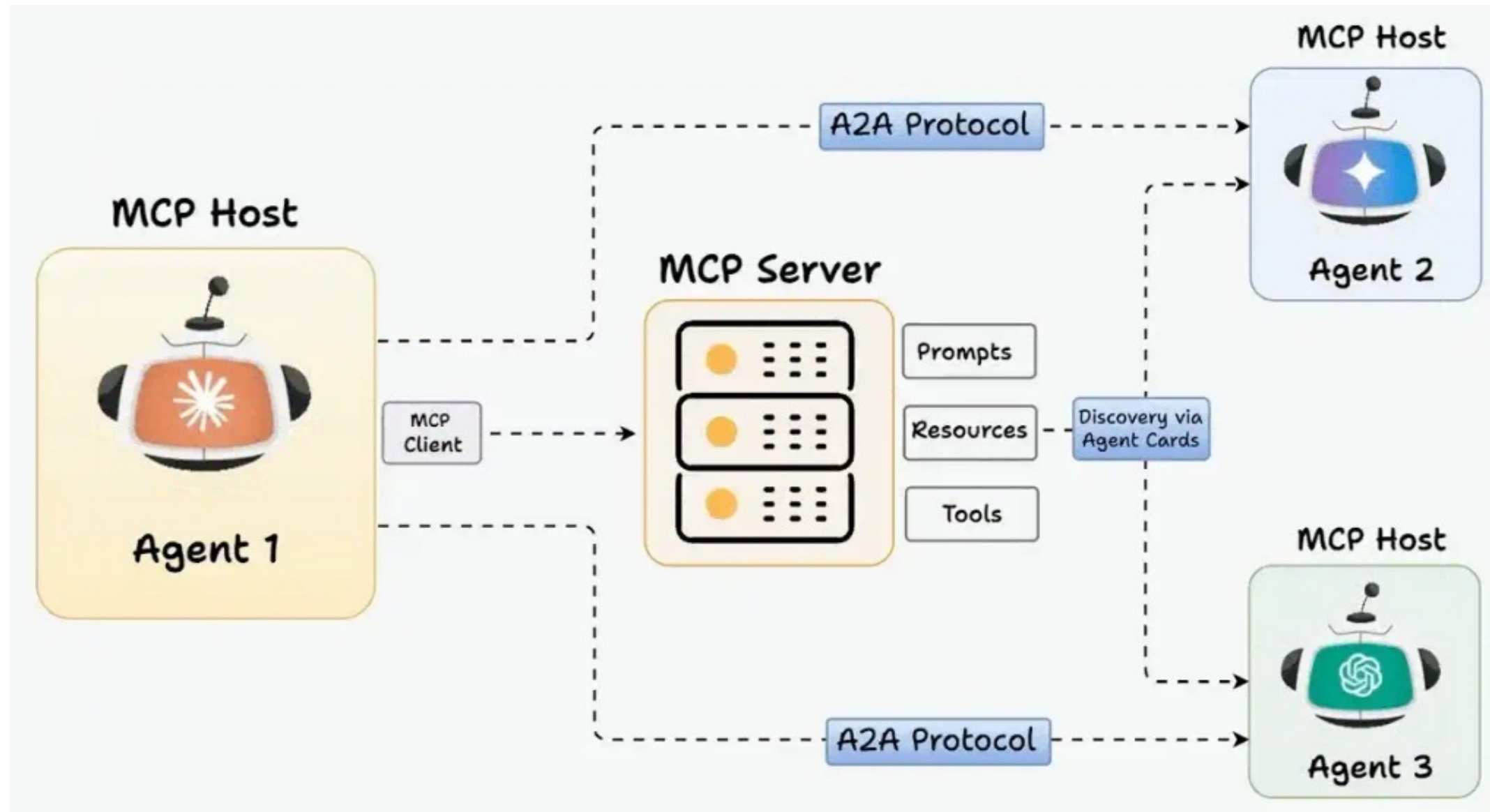
# Agentic applications require both A2A and MCP

A2A allows agents to connect with other agents and collaborate in teams.
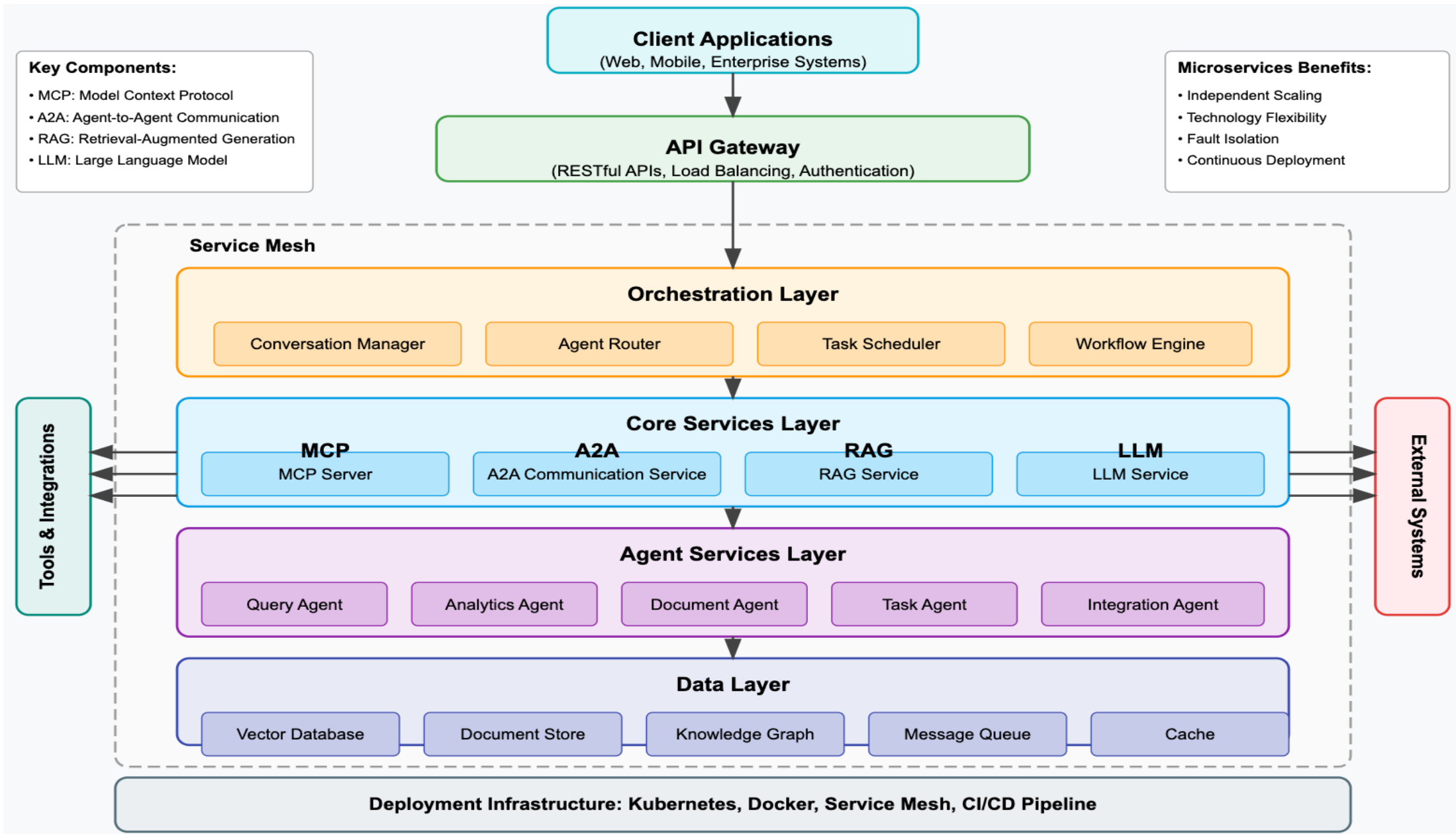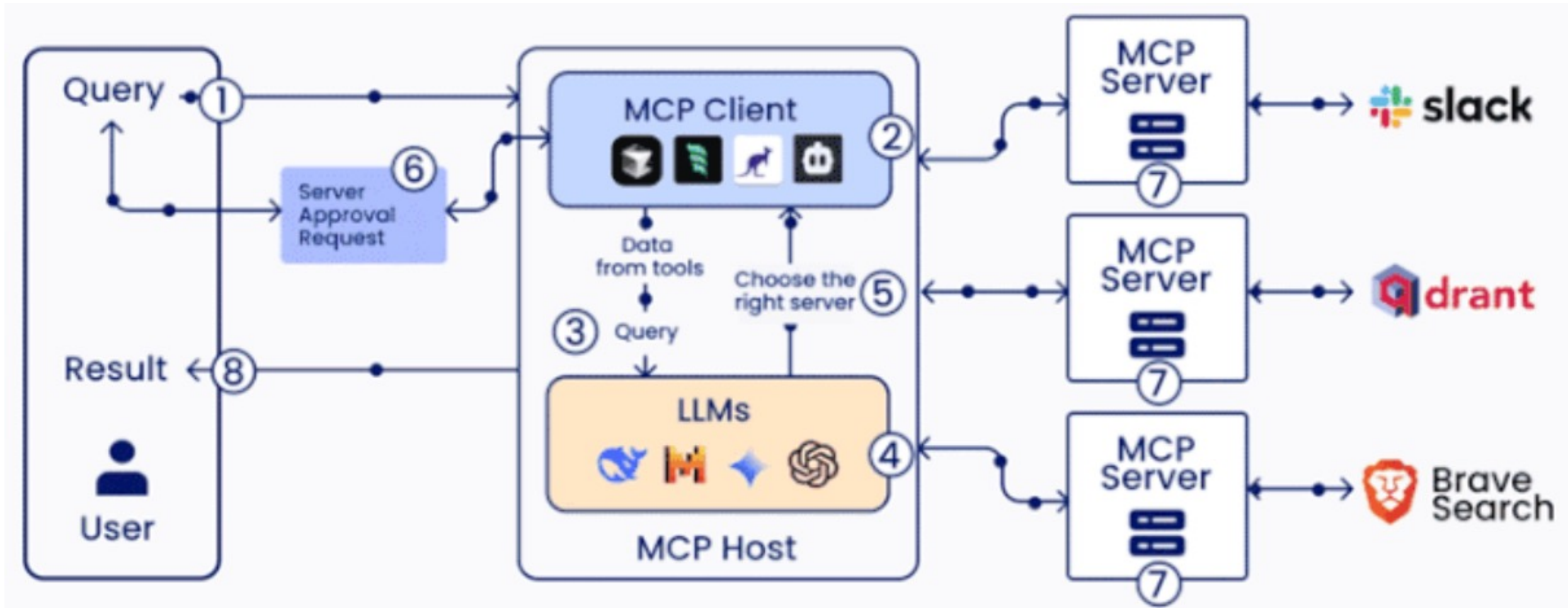
MCP provides agents with access to tools
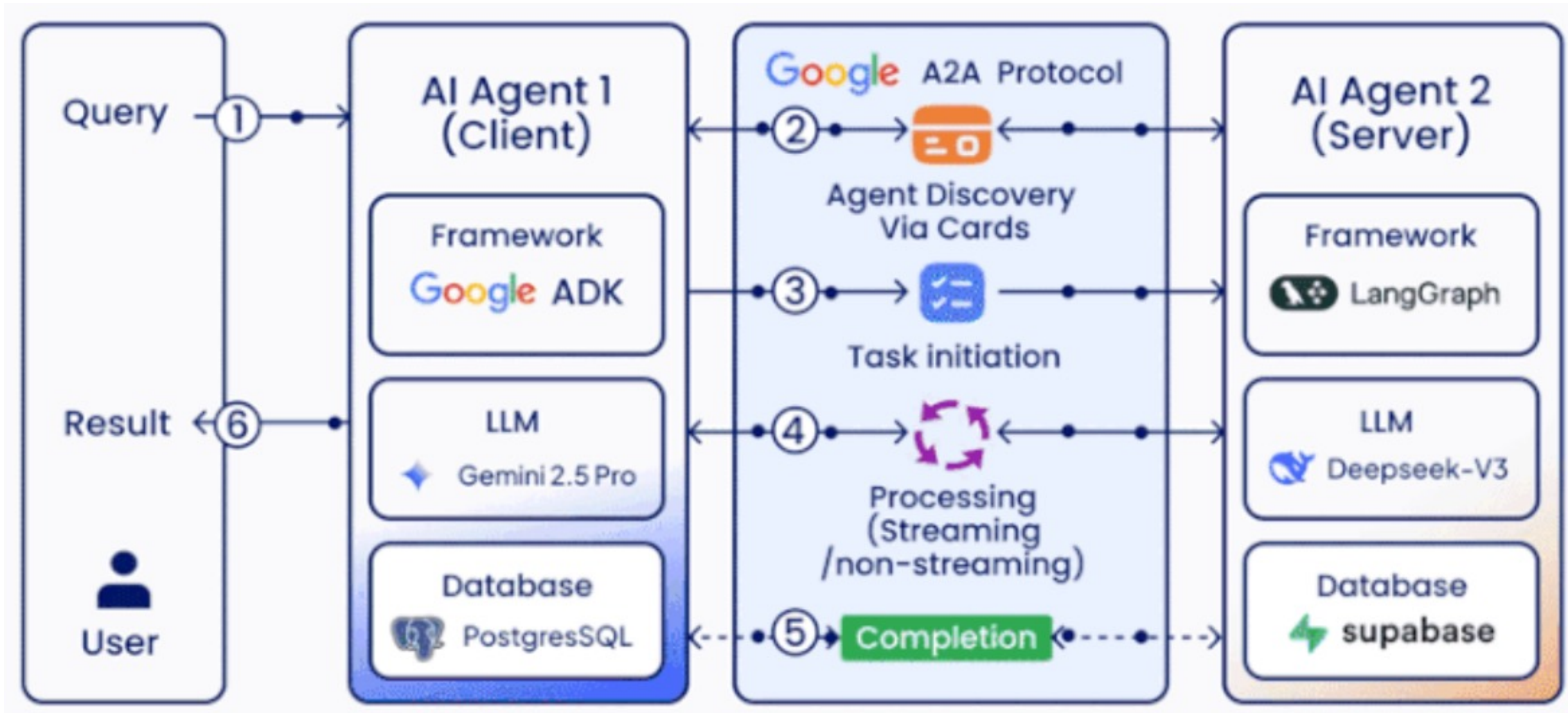
# MCP and A2A Protocol for AI Agents

# Agentic AI System with Microservices Architecture



**Client Applications**
(Web, Mobile, Enterprise Systems)

**Key Components:**
- MCP: Model Context Protocol
- A2A: Agent-to-Agent Communication
- RAG: Retrieval-Augmented Generation
- LLM: Large Language Model

**Microservices Benefits:**
- Independent Scaling
- Technology Flexibility
- Fault Isolation
- Continuous Deployment

**API Gateway**
(RESTful APIs, Load Balancing, Authentication)

**Service Mesh**

**Orchestration Layer**

| Conversation Manager | Agent Router | Task Scheduler | Workflow Engine |

**Core Services Layer**

**Tools & Integrations**

**External Systems**

| MCP | A2A | RAG | LLM |
| MCP Server | A2A Communication Service | RAG Service | LLM Service |

**Agent Services Layer**

| Query Agent | Analytics Agent | Document Agent | Task Agent | Integration Agent |

**Data Layer**

| Vector Database | Document Store | Knowledge Graph | Message Queue | Cache |

**Deployment Infrastructure: Kubernetes, Docker, Service Mesh, CI/CD Pipeline**

70

# MCP (Model Context Protocol)

Source: Houssem Eddine Lassoued (2025), How MCP and A2A Protocols Are Revolutionizing AI Agent Development, https://medium.com/@houssemeddinelassoued/

# A2A (Agent2Agent Protocol)

72

# Agentic AI and World Model for Edge General Intelligence

Source: Changyuan Zhao, Guangyuan Liu, Ruichen Zhang, Yinqiu Liu, Jiacheng Wang, Jiawen Kang, Dusit Niyato et al (2025), "Edge general intelligence through world models and agentic AI: Fundamentals, solutions, and challenges." arXiv preprint arXiv:2508.09561.

# **Physical AI (Robotics)**

# Framework of the Embodied Agent based on MLMs and WMs

# Embodied AI



Embodied AI: Aligning Cyber Space with Physical World

**Embodied Robots**

Fixed-base Robots,
Wheeled &Tracked Robots,
Quadruped Robots,
Humanoid Robots,
Biomimetic Robots,
...

**Simulators**

General Simulator,
Real-Scene based
Simulator,
...

**Embodied Perception**

Active Visual Exploration,
3D Visual Grounding,
Visual-Language Navigation,
Non-Visual Perception,
...

**Embodied Interaction**

Embodied Question
Answering,
Embodied Grasping,
...

**Embodied Agent**

Embodied Multi-modal
Foundation Model,
Embodied Task Planning,
...

**Sim-to-Real Adaptation**

Embodied World Model, Data Collection and Training, Embodied Control

Applications | Robotics | Autonomous Driving | Healthcare | Domestic Assistance | Industrial Automation | Search and Rescue

# Embodied Agents



MLM: Multimodal Language Model, which directly perceive the world and control the embodiment

VLM: Visual-Language Model with the outer policy models

LLM + VLM: LLM-based agent that perceives the world utilizing the VLM, and LLM
means the Large-Language Model with visual context and outer policy models.

# Boston Dynamics: Spot

## Automate sensing and inspection, capture limitless data, and explore without boundaries.

# Boston Dynamics: Atlas
## The world's most dynamic humanoid robot

Atlas is a research platform designed to push the limits of whole-body mobility

# Boston Dynamics: Atlas Goes Hands On

Atlas uses a machine learning (ML) vision model to detect and localize the environment fixtures and individual bins.

The robot uses a specialized grasping policy and continuously estimates the state of manipulated objects to achieve the task.



Fully Autonomous

# Boston Dynamics: Atlas



#13 ON TRENDING
What's new, Atlas?

https://www.youtube.com/watch?v=fRj34o4hN4I

# Humanoid Robot: Sophia



https://www.youtube.com/watch?v=S5t6K9iwcdw

# Can a robot pass a university entrance exam?
## Noriko Arai at TED2017

# Embodied Robots



(a) Fixed-base Robots (Franka Emika Panda)

(b) Wheeled Robots (Jackal robot)

(c) Tracked Robots (iRobot PackBot)

(d) Quadruped Robots (Boston Dynamics Spot)

(e) Humanoid Robots (Tesla Optimus)

(f) Biomimetic Robots

# Gemini Robotics:
# Bringing AI into the Physical World
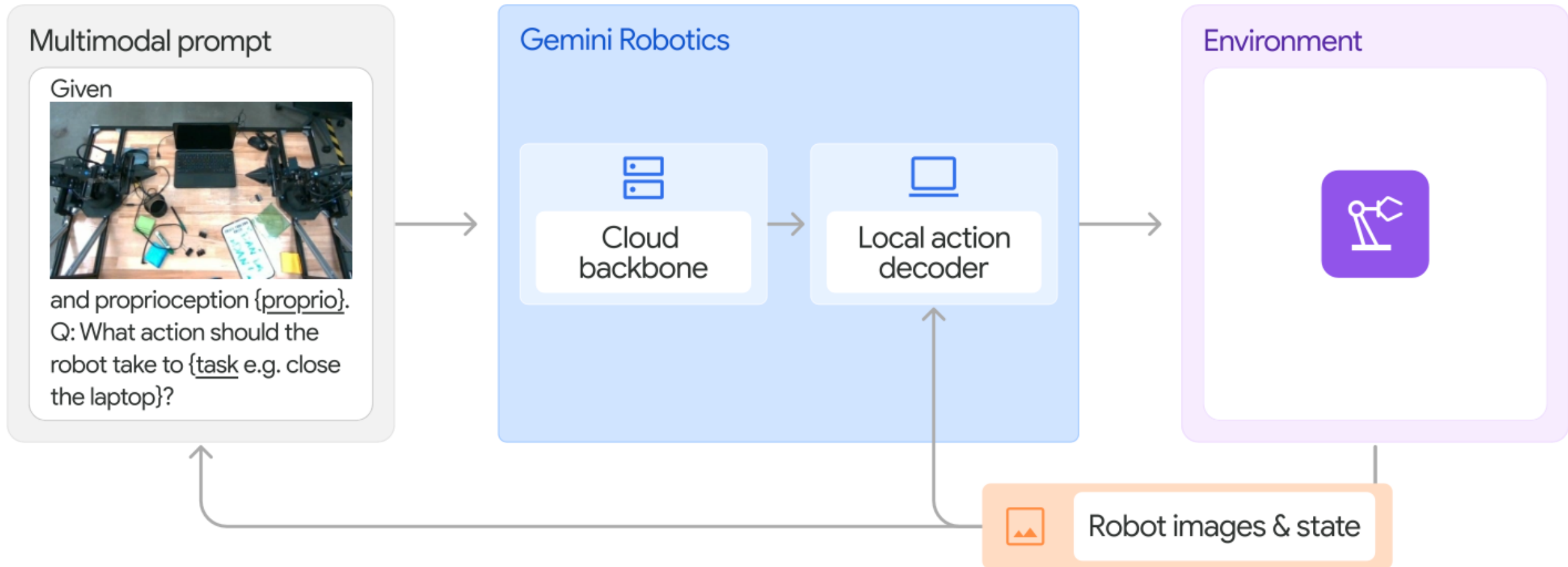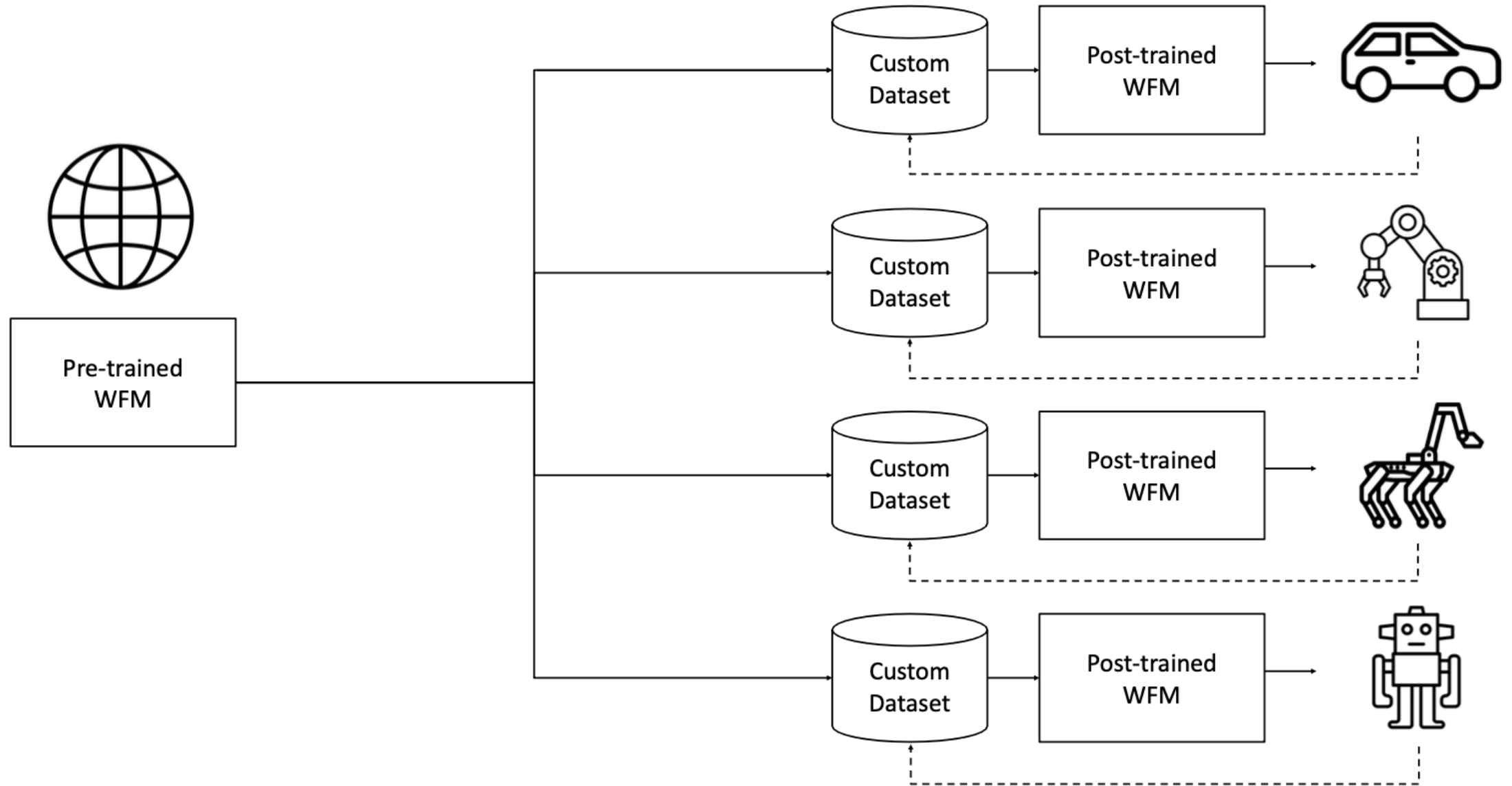


Source: Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna et al.(2025)
"Gemini robotics: Bringing ai into the physical world." arXiv preprint arXiv:2503.20020 (2025).

# Gemini Robotics Models:
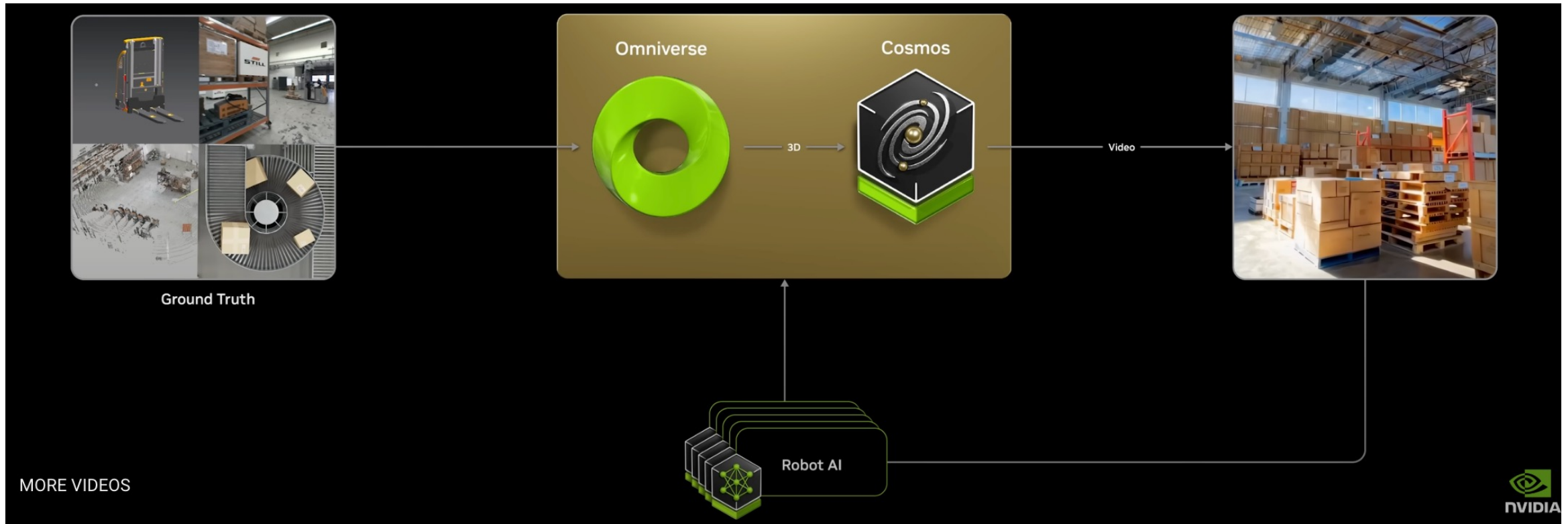# Architecture, Input and Output

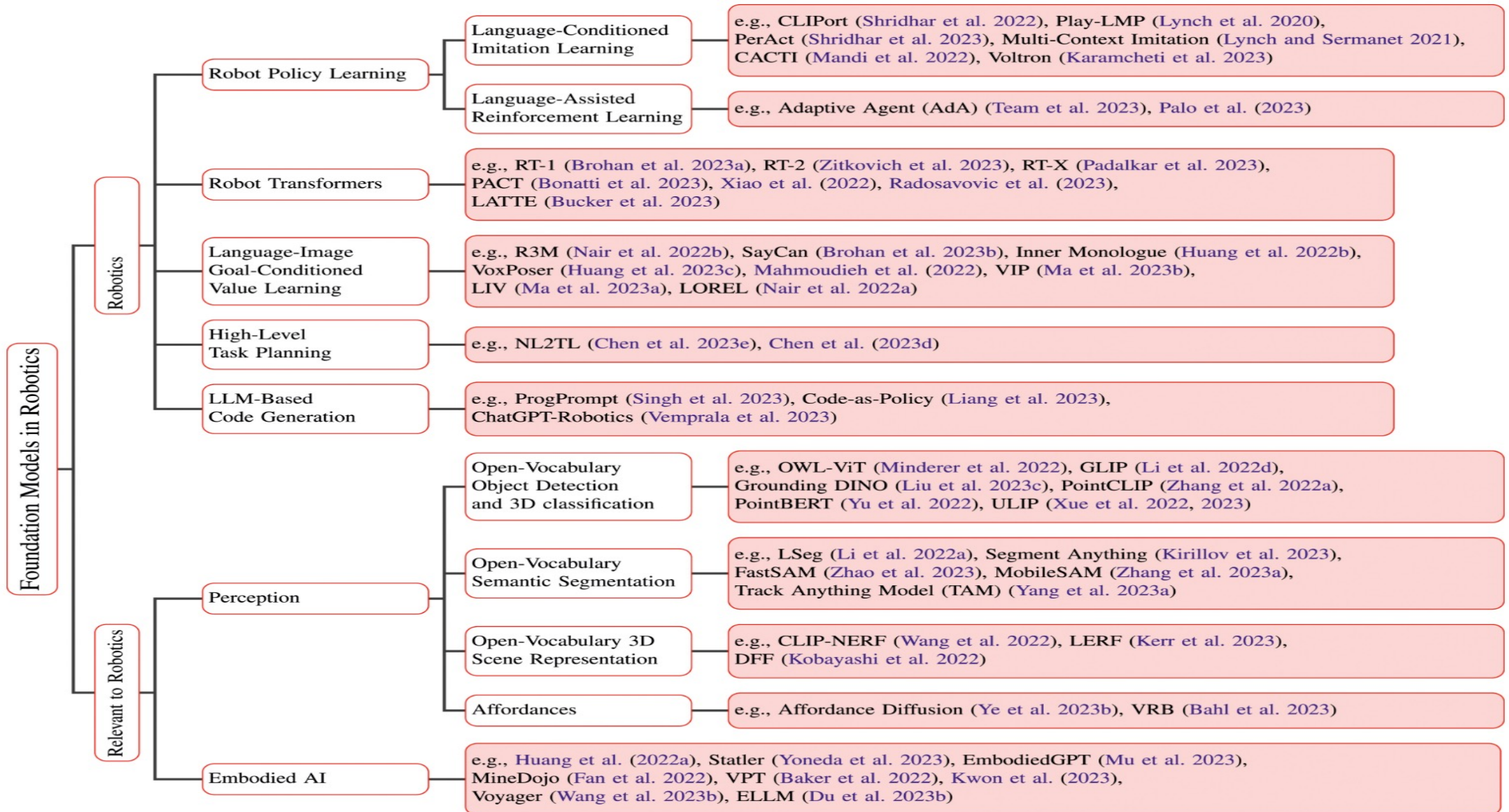# World Foundation Model Platform for Physical AI

# NVIDIA Cosmos
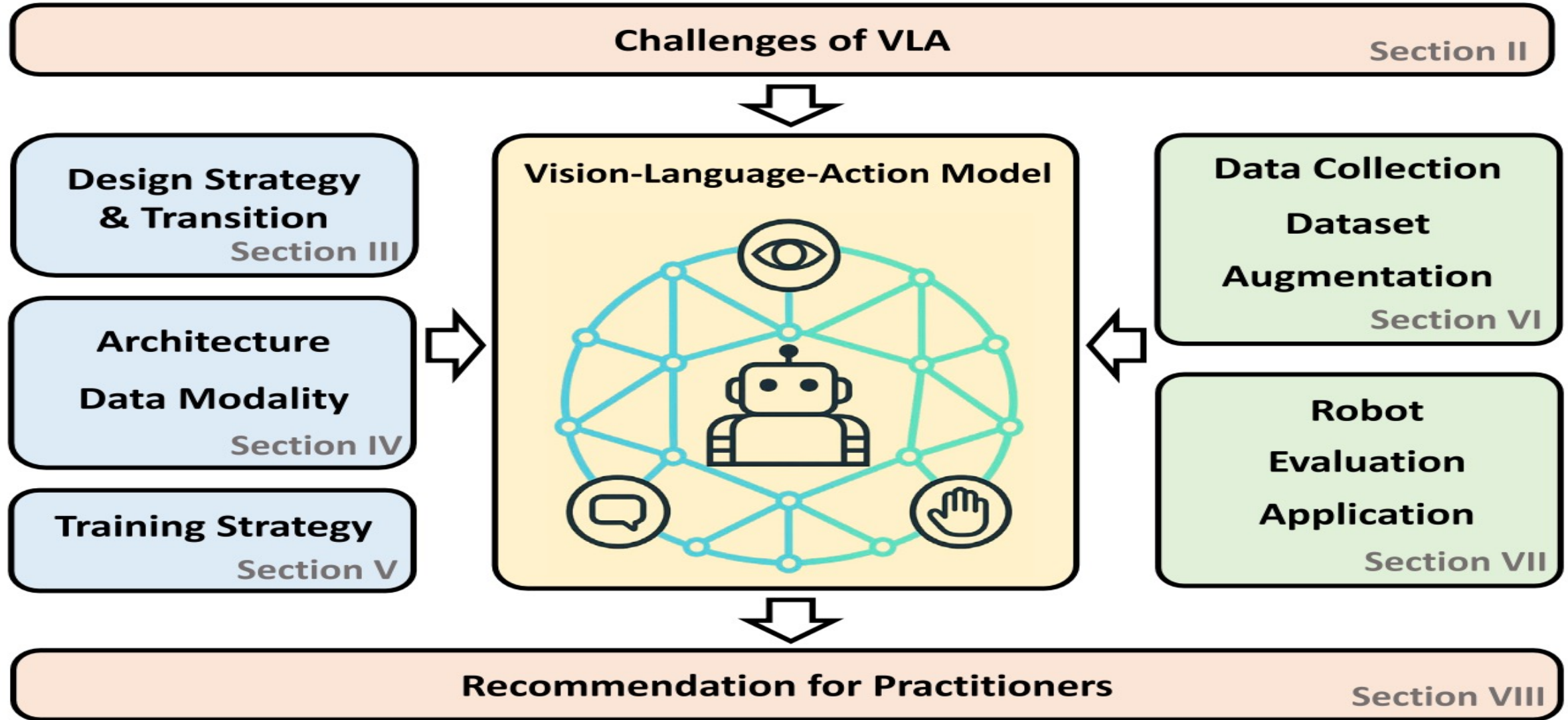# World Foundation Model Platform for Physical AI

Source: Agarwal, Niket, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay et al. "Cosmos world foundation model platform for physical ai."
arXiv preprint arXiv:2501.03575 (2025)., https://research.nvidia.com/labs/dir/cosmos-predict1/
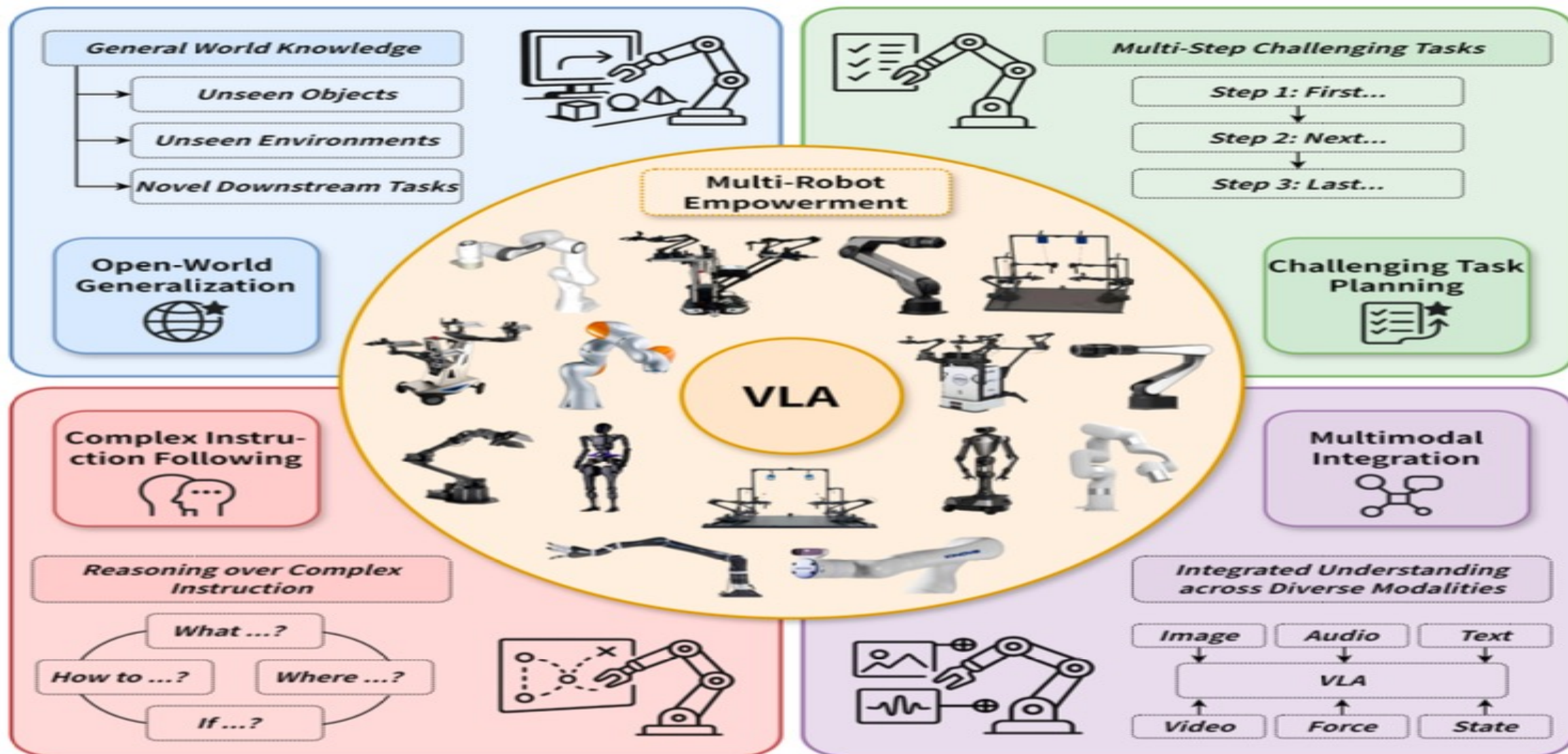
# Foundation Models in Robotics

89

# Vision Language Action (VLA) Models for Robotics

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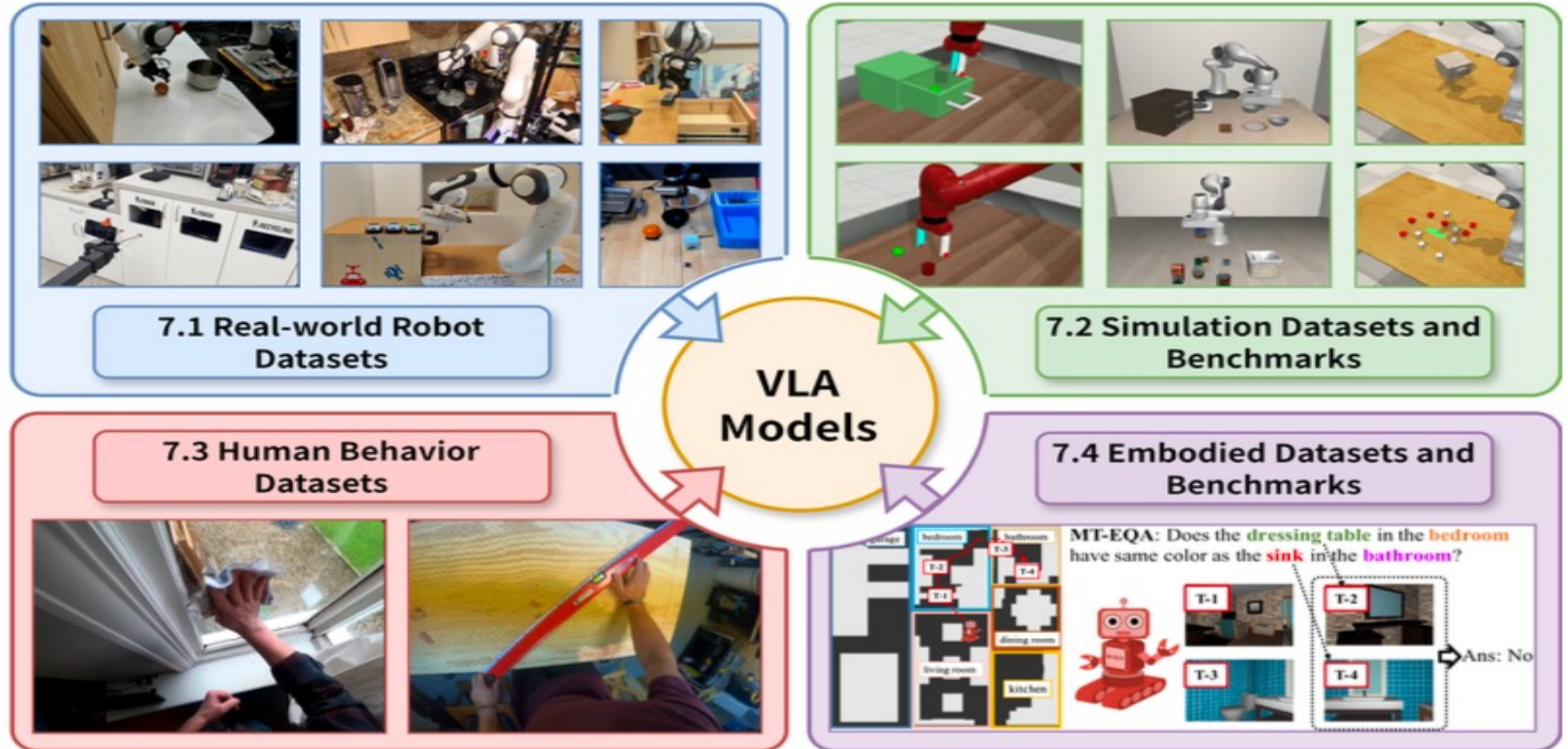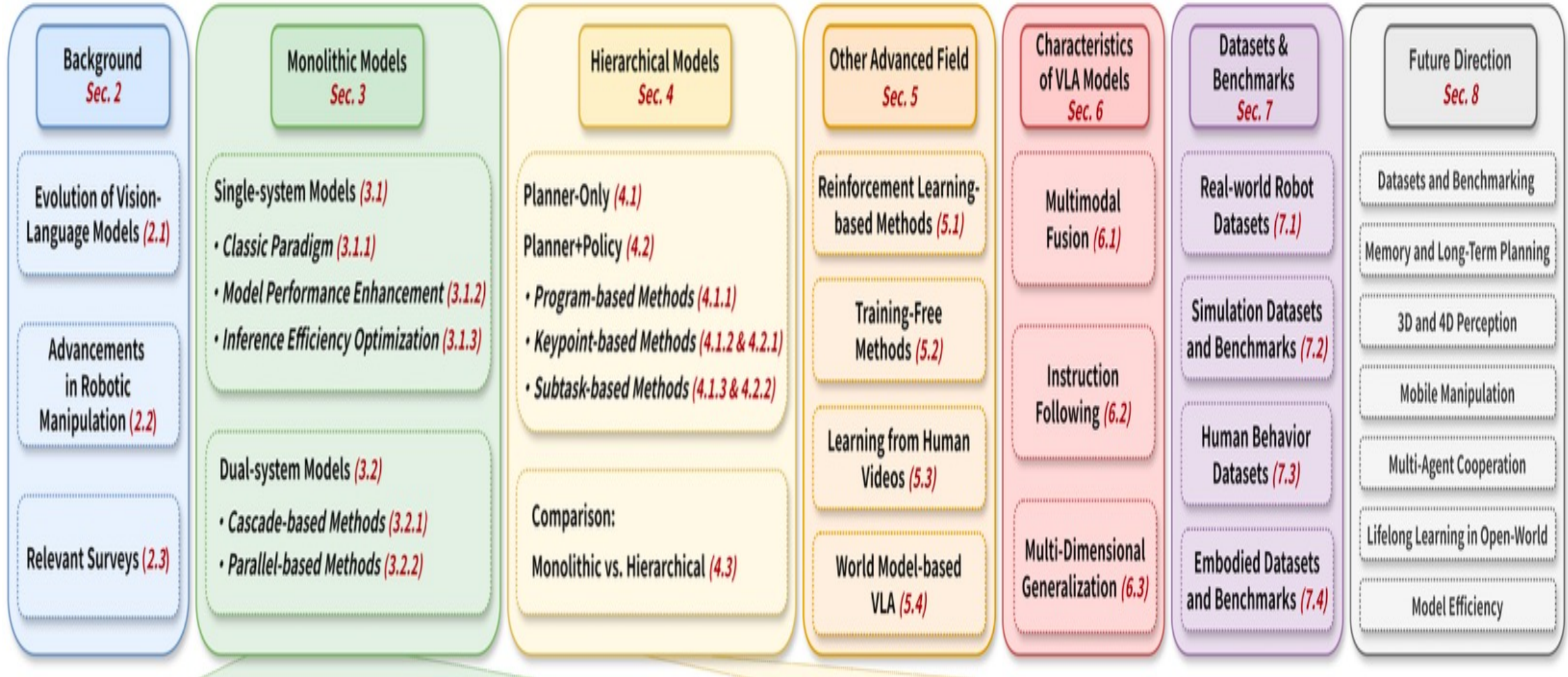
# Large VLM-based Vision-Language-Action Models
## for Robotic Manipulation

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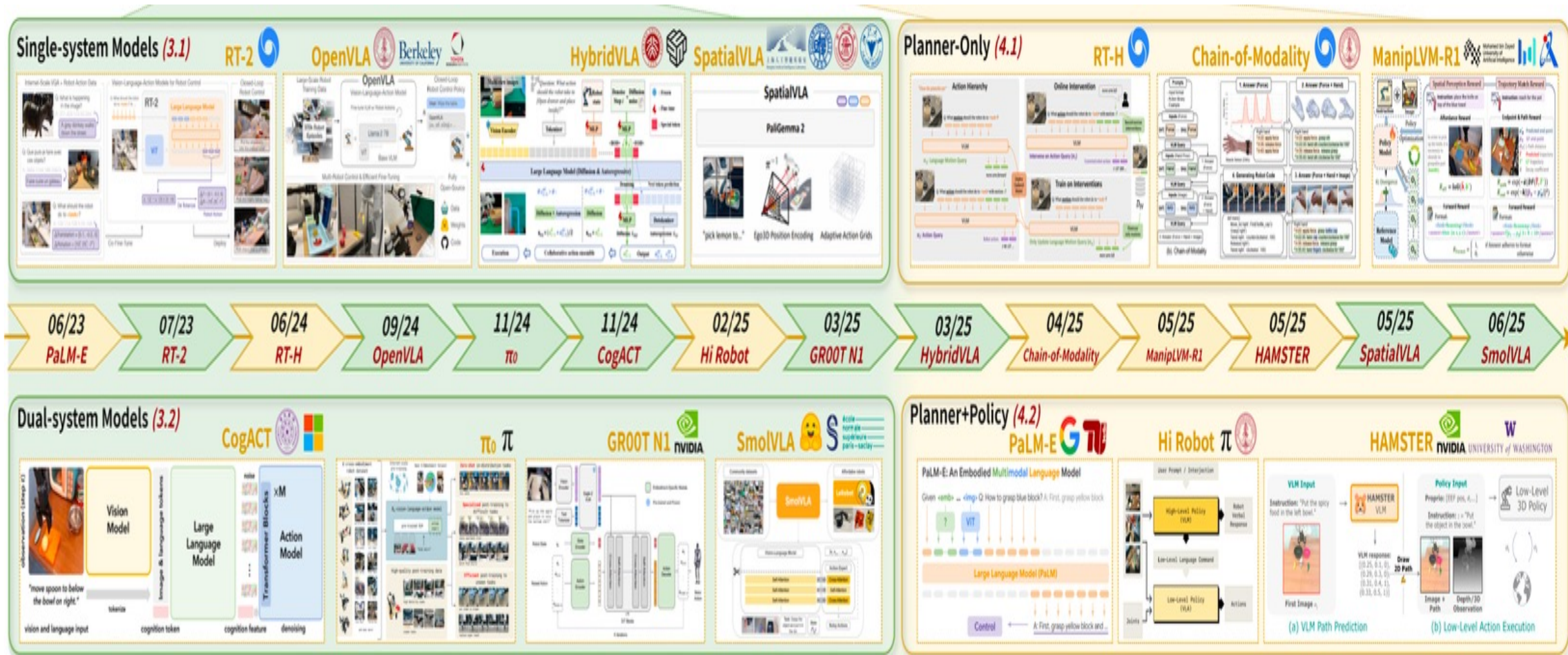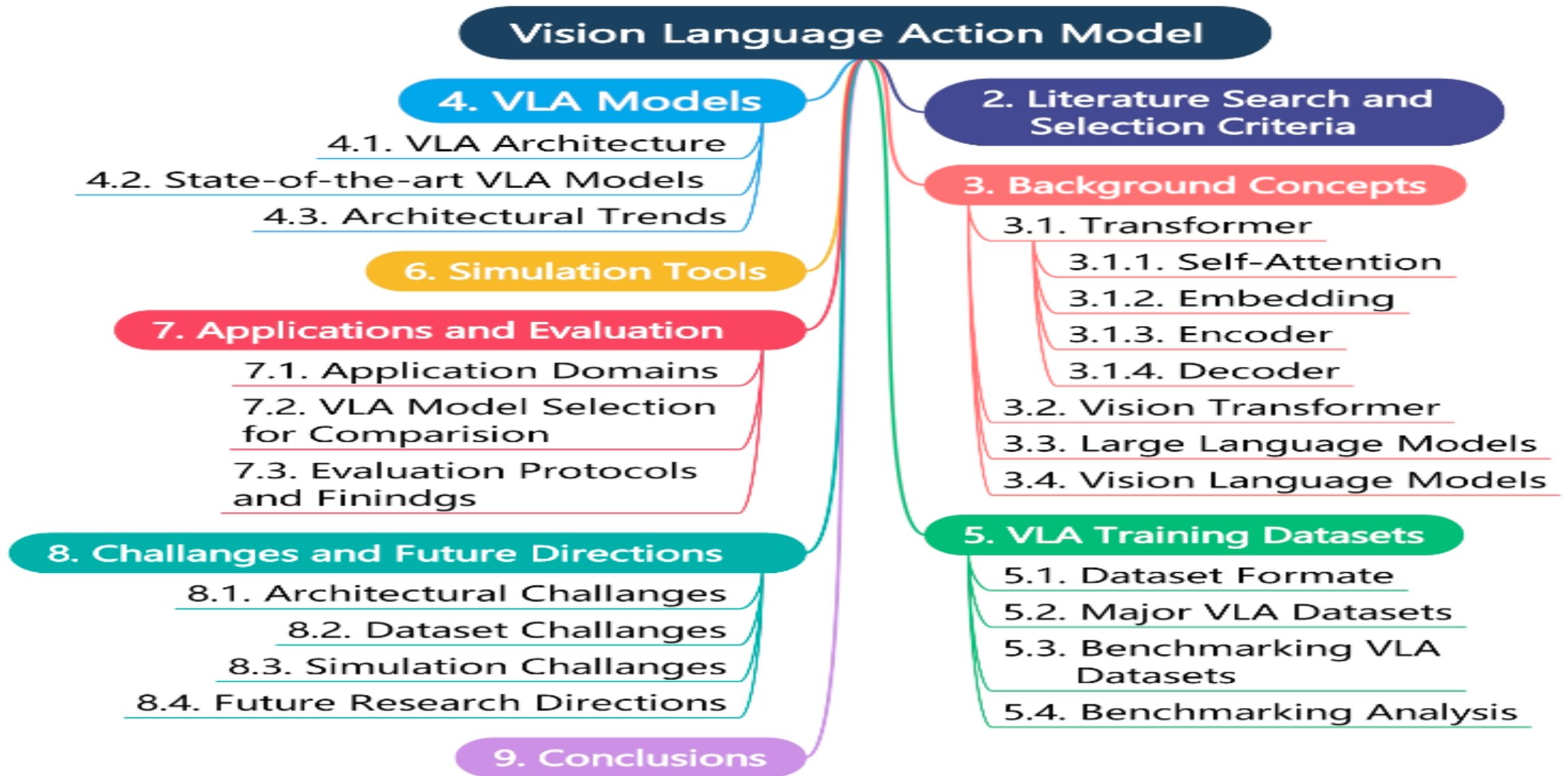


Source: Rui Shao, Wei Li, Lingsen Zhang, Renshan Zhang, Zhiyang Liu, Ran Chen, and Liqiang Nie. (2025) "Large vlm-based vision-language-action models for robotic manipulation: A survey." arXiv preprint arXiv:2508.13073 (2025).

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

**Background** *Sec. 2*

- Evolution of Vision-Language Models *(2.1)*
- Advancements in Robotic Manipulation *(2.2)*
- Relevant Surveys *(2.3)*

**Monolithic Models** *Sec. 3*

Single-system Models *(3.1)*
- *Classic Paradigm (3.1.1)*
- *Model Performance Enhancement (3.1.2)*
- *Inference Efficiency Optimization (3.1.3)*

Dual-system Models *(3.2)*
- *Cascade-based Methods (3.2.1)*
- *Parallel-based Methods (3.2.2)*

**Hierarchical Models** *Sec. 4*

Planner-Only *(4.1)*

Planner+Policy *(4.2)*
- *Program-based Methods (4.1.1)*
- *Keypoint-based Methods (4.1.2 & 4.2.1)*
- *Subtask-based Methods (4.1.3 & 4.2.2)*

Comparison:
Monolithic vs. Hierarchical *(4.3)*

**Other Advanced Field** *Sec. 5*

- Reinforcement Learning-based Methods *(5.1)*
- Training-Free Methods *(5.2)*
- Learning from Human Videos *(5.3)*
- World Model-based VLA *(5.4)*

**Characteristics of VLA Models** *Sec. 6*

- Multimodal Fusion *(6.1)*
- Instruction Following *(6.2)*
- Multi-Dimensional Generalization *(6.3)*

**Datasets & Benchmarks** *Sec. 7*

- Real-world Robot Datasets *(7.1)*
- Simulation Datasets and Benchmarks *(7.2)*
- Human Behavior Datasets *(7.3)*
- Embodied Datasets and Benchmarks *(7.4)*

**Future Direction** *Sec. 8*

- Datasets and Benchmarking
- Memory and Long-Term Planning
- 3D and 4D Perception
- Mobile Manipulation
- Multi-Agent Cooperation
- Lifelong Learning in Open-World
- Model Efficiency

# Large VLM-based Vision-Language-Action Models
## for Robotic Manipulation (Timeline)
### Monolithic models and Hierarchical Models



Timeline: 06/23 PaLM-E · 07/23 RT-2 · 06/24 RT-H · 09/24 OpenVLA · 11/24 π0 · 11/24 CogACT · 02/25 Hi Robot · 03/25 GR00T N1 · 03/25 HybridVLA · 04/25 Chain-of-Modality · 05/25 ManipLVM-R1 · 05/25 HAMSTER · 05/25 SpatialVLA · 06/25 SmolVLA

# Vision Language Action Models in Robotic Manipulation

# Vision Language Action (VLA) Models, Datasets

**Contributing institutions: Academic** (e.g., CMU, CNRS, UC, Peking Uni)
**Industrial Labs** (e.g., Google, NVIDIA, Microsoft)

# VLA Models and Foundational VLA Datasets
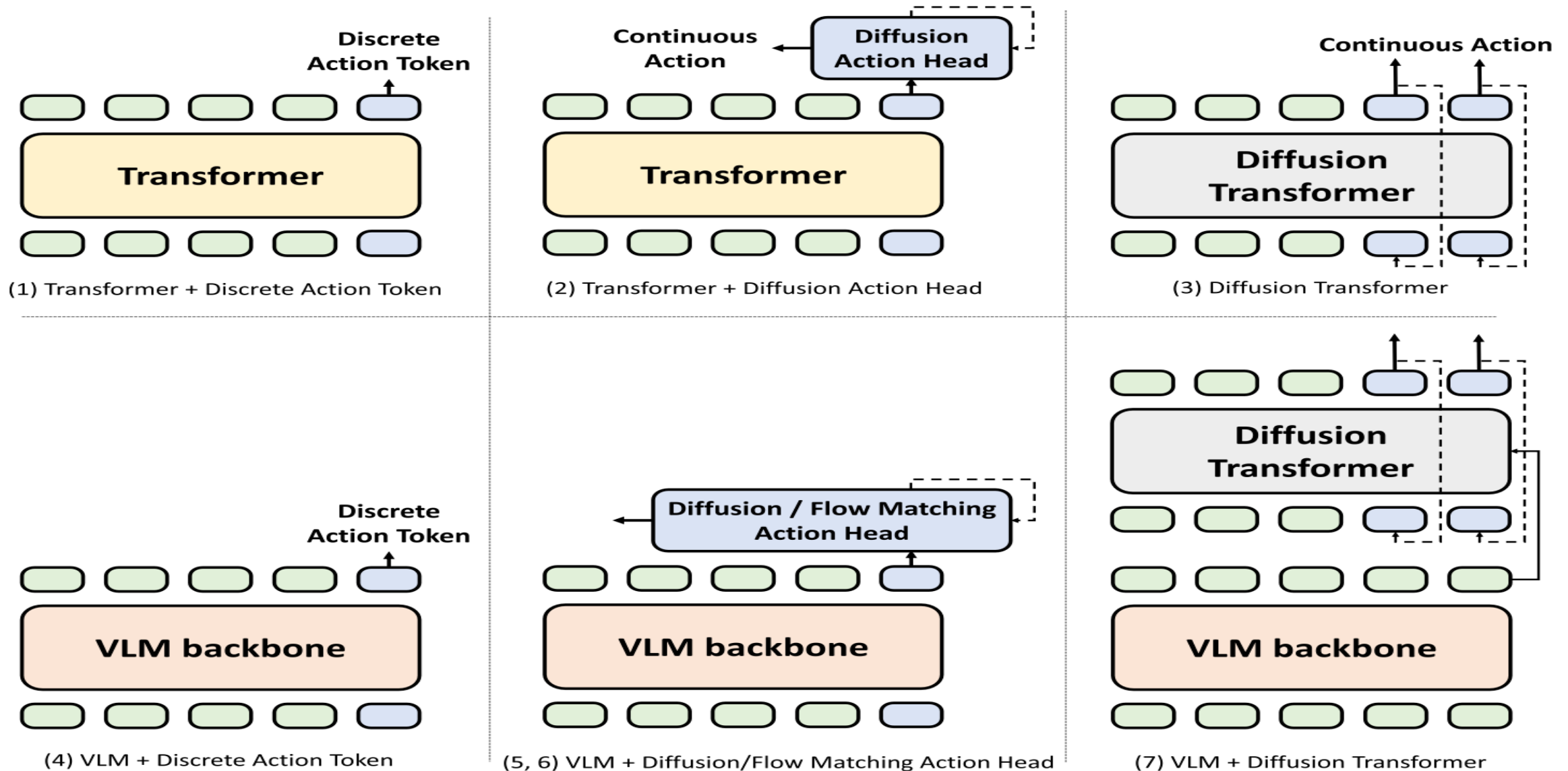
# Timeline of Vision-Language-Action (VLA) Models

# VLA Model Components and Training Paradigms

# Architecture of Sensorimotor Models for VLA



(1) Transformer + Discrete Action Token

(2) Transformer + Diffusion Action Head

(3) Diffusion Transformer

(4) VLM + Discrete Action Token

(5, 6) VLM + Diffusion/Flow Matching Action Head

(7) VLM + Diffusion Transformer

Source: Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. (2025) "Vision-language-action models for robotics: A review towards real-world applications." IEEE Access (2025).

# Transformer Architecture:

## Encoder-decoder structure and the internal mechanism of multi-head attention
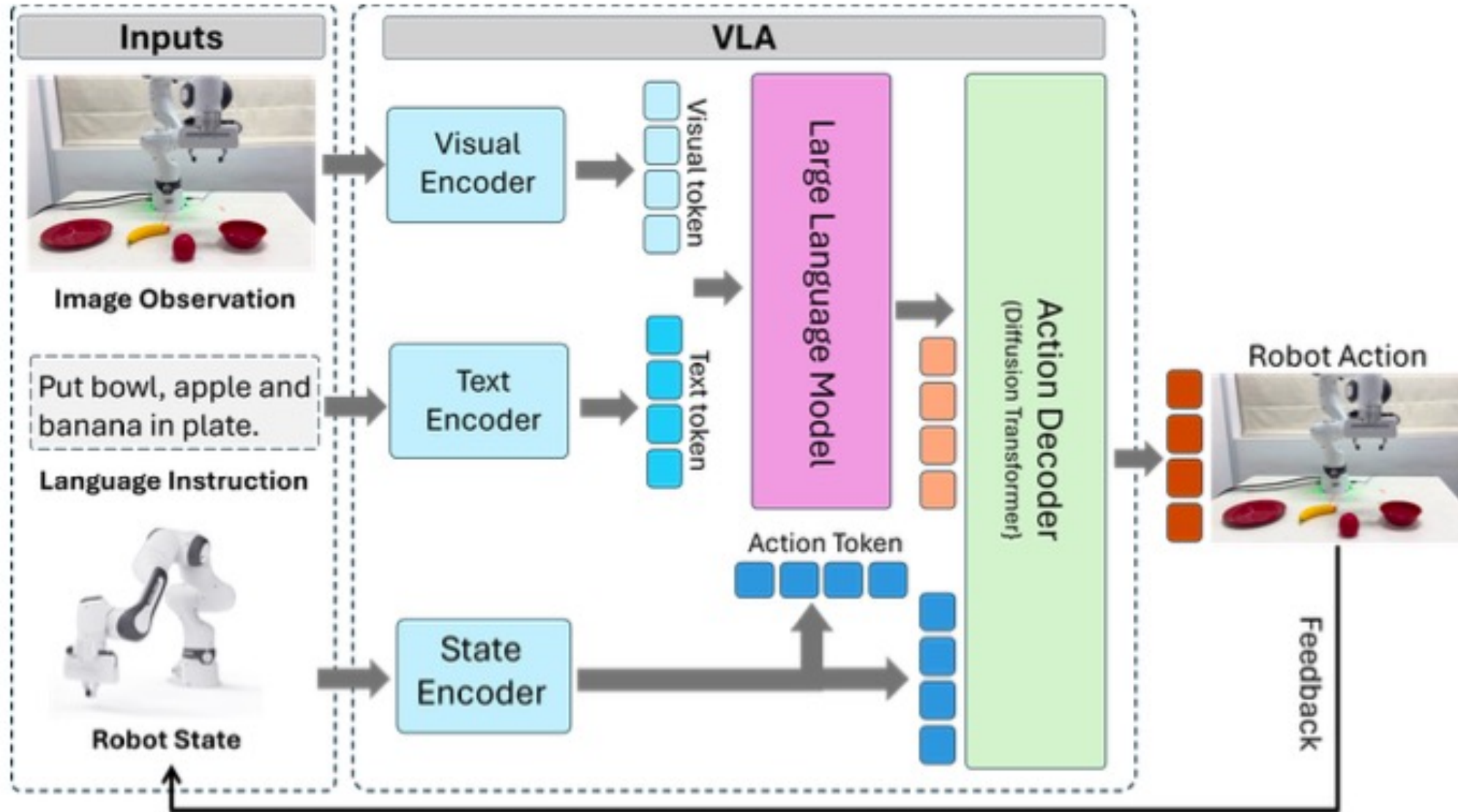
# Architecture of the ViT

# Architecture of VLM
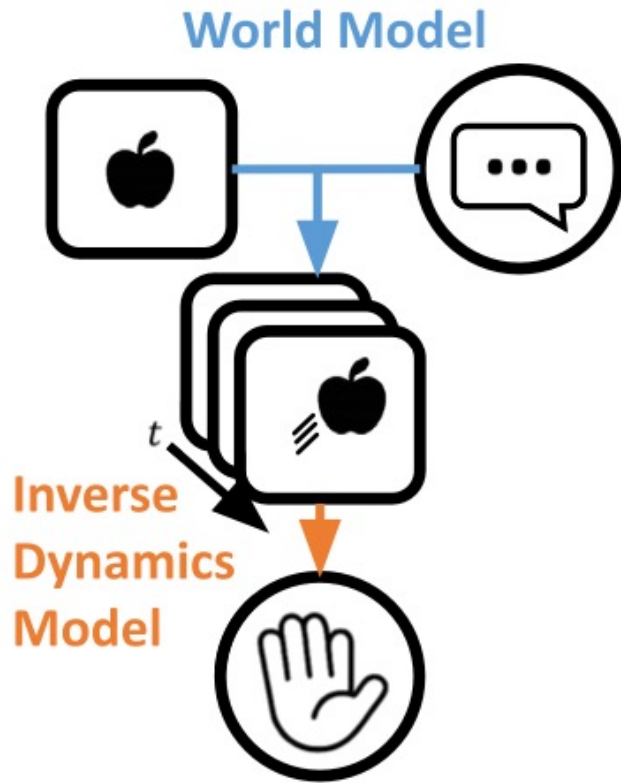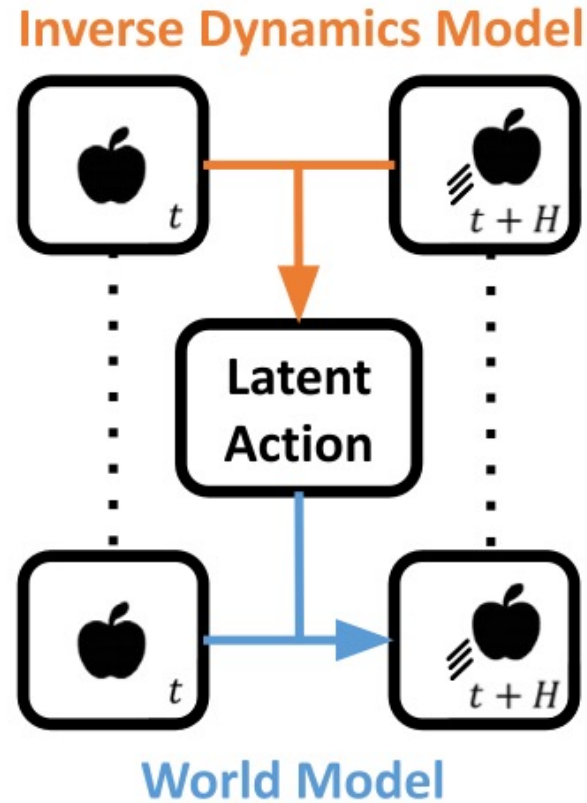# for Image Captioning and Semantic Understanding

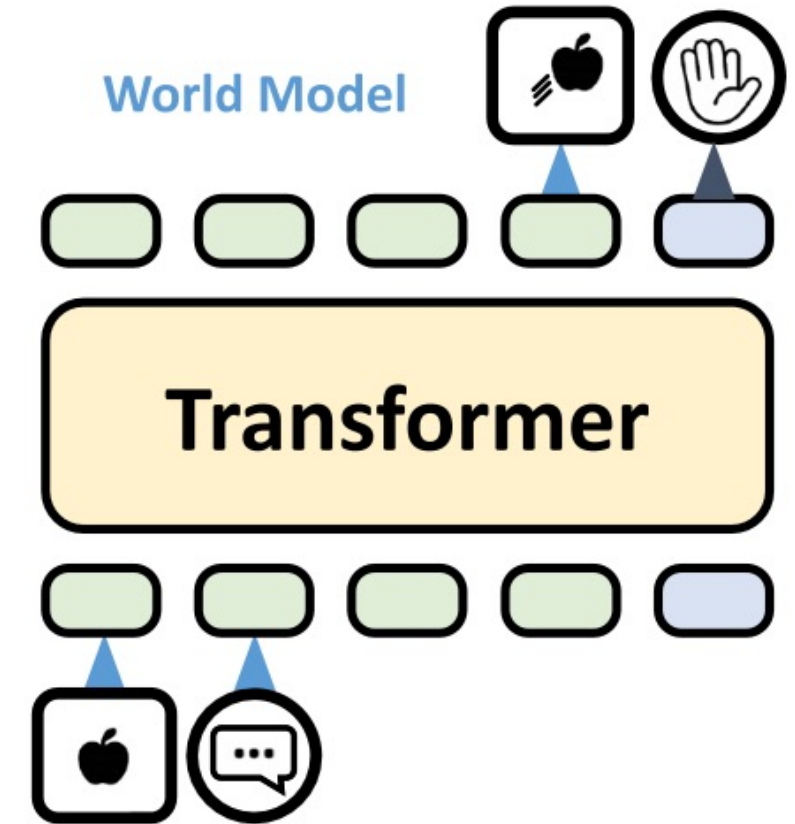# Architecture of a VLA System for Robotic Manipulation

# Design Patterns for Incorporating World Models in VLA



(1) Action generation in world models

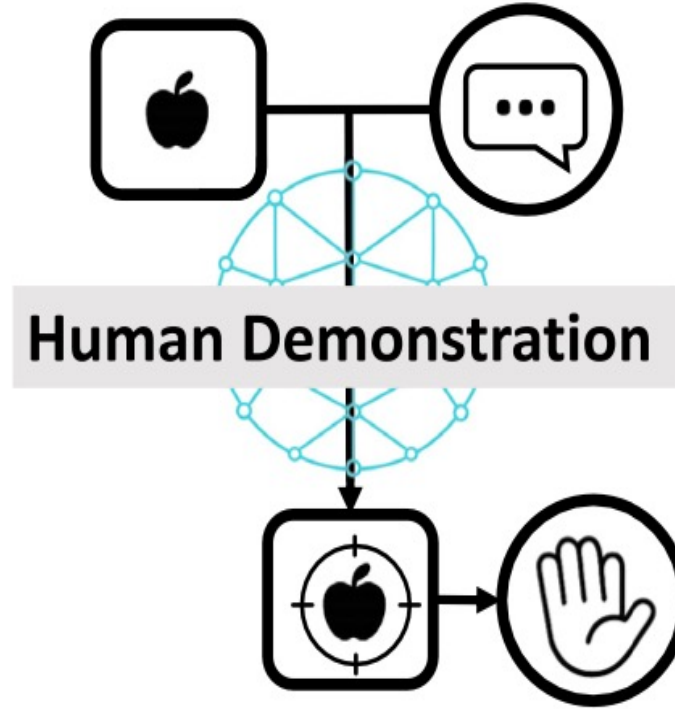(2) Latent action generation via world models

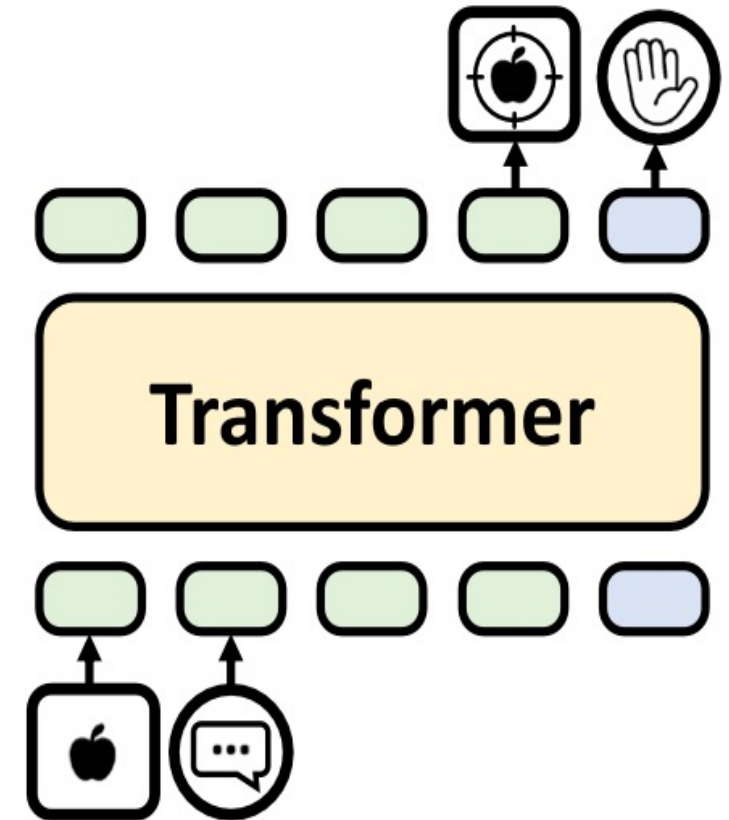(3) Sensorimotor models with implicit world models

# Design Patterns for Incorporating Affordance-based Models in VLA



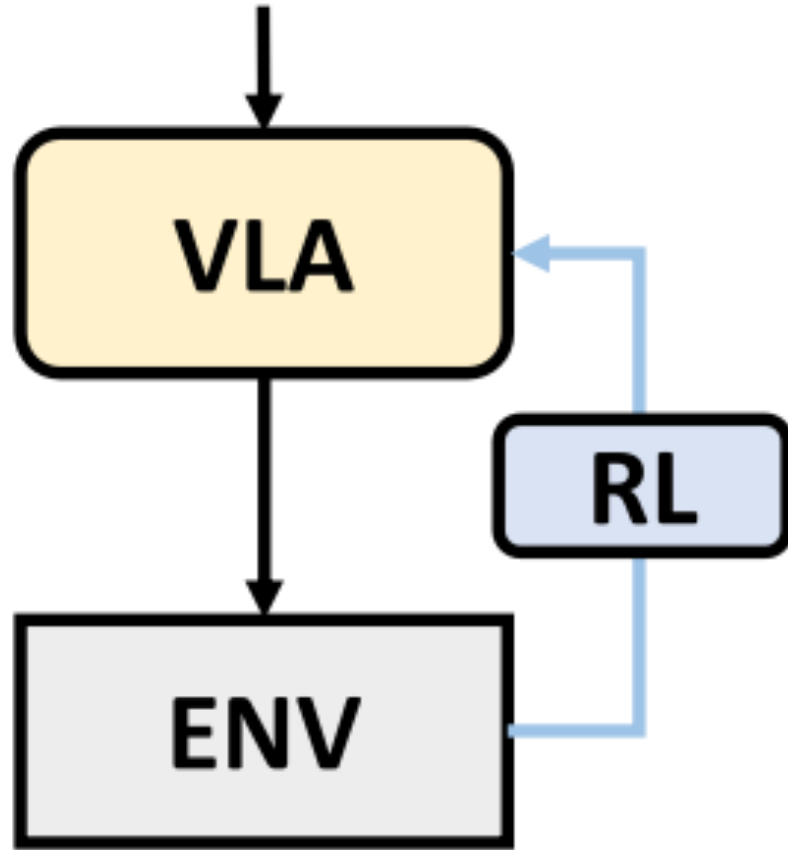(1) Affordance prediction and action generation using VLMs

(2) Affordance extraction from human datasets
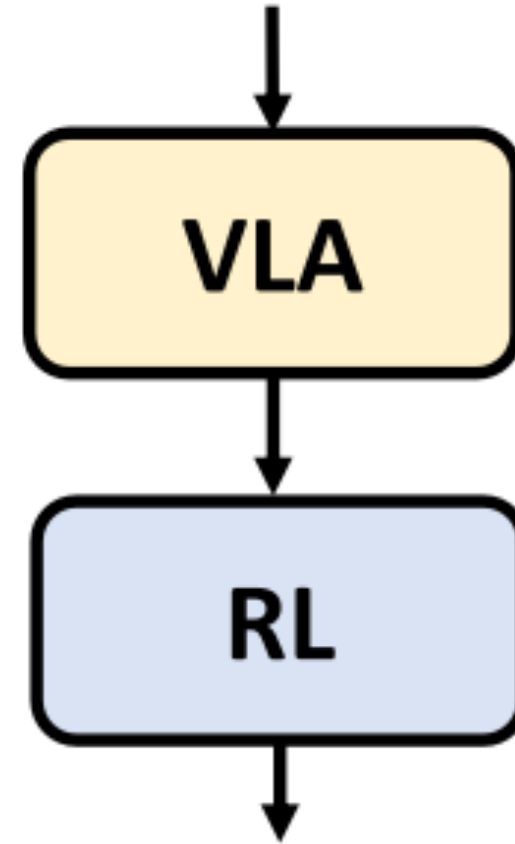
(3) Integration of sensorimotor models and affordance-based models

# Integrating RL with VLA Models



(1) Improving VLA using RL

(2) Using VLAs as high-level policies and RL for low-level control

Source: Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. (2025) "Vision-language-action models for robotics: A review towards real-world applications." IEEE Access (2025).

# Robots Used in VLA Research



**Robot**
- Manipulator
- Hand/Gripper
- Mobile Robot
- Quadruped Robot
- Humanoid Robot

**Data Collection**
- Teleoperation
- Proxy Devices
- Human Data Collection

**Dataset**

Human Video Datasets
- Ego4D  Ego-Exo4D
- HOI4D  ARCTIC

Simulation Datasets
- RoboTurk  MimicGen

Real Robot Datasets
- QT-Opt  RT-X
- BC-Z  DROID  ...

**Augmentation**
- Vision
- Language
- Action

**Evaluation**

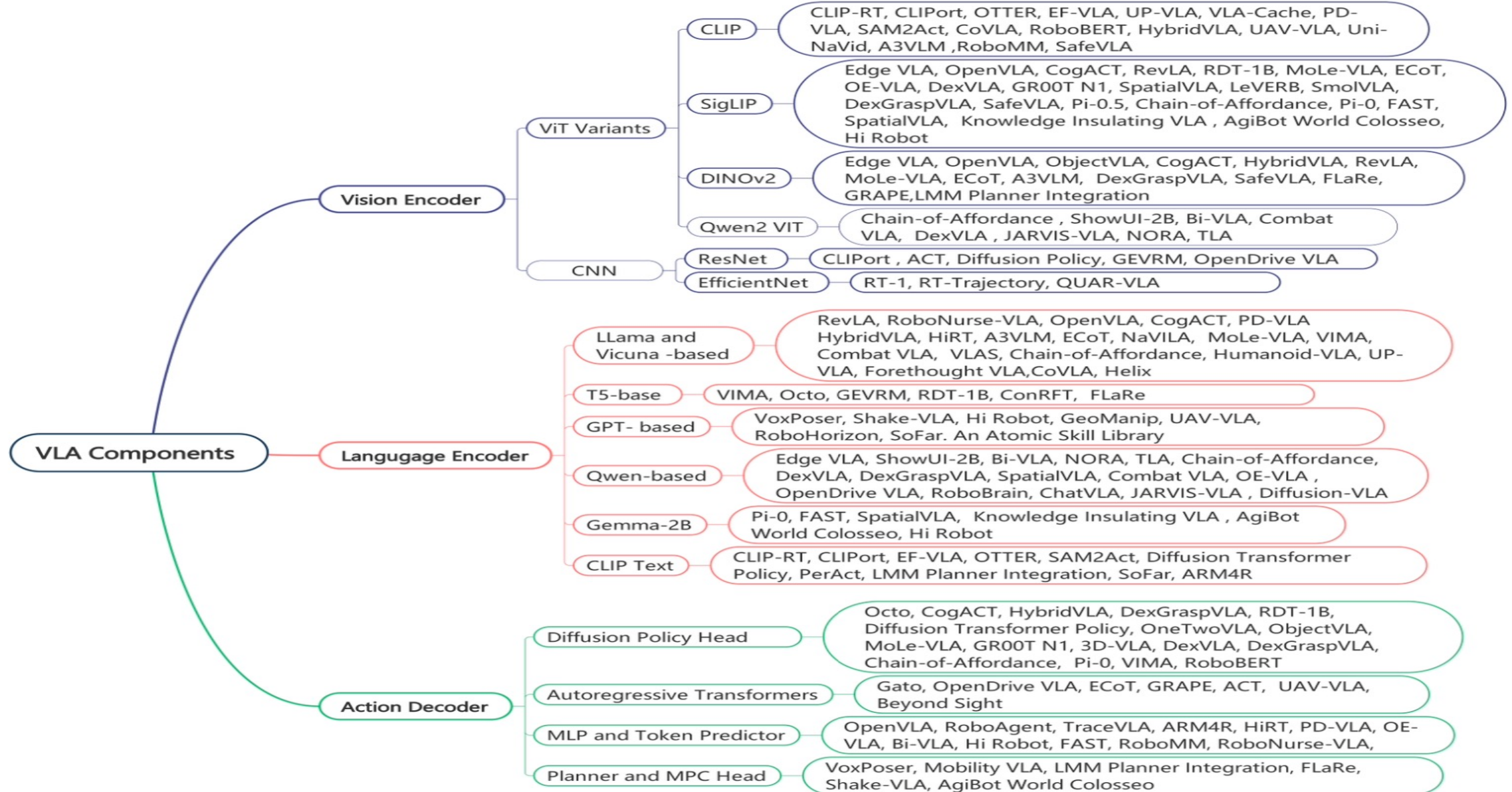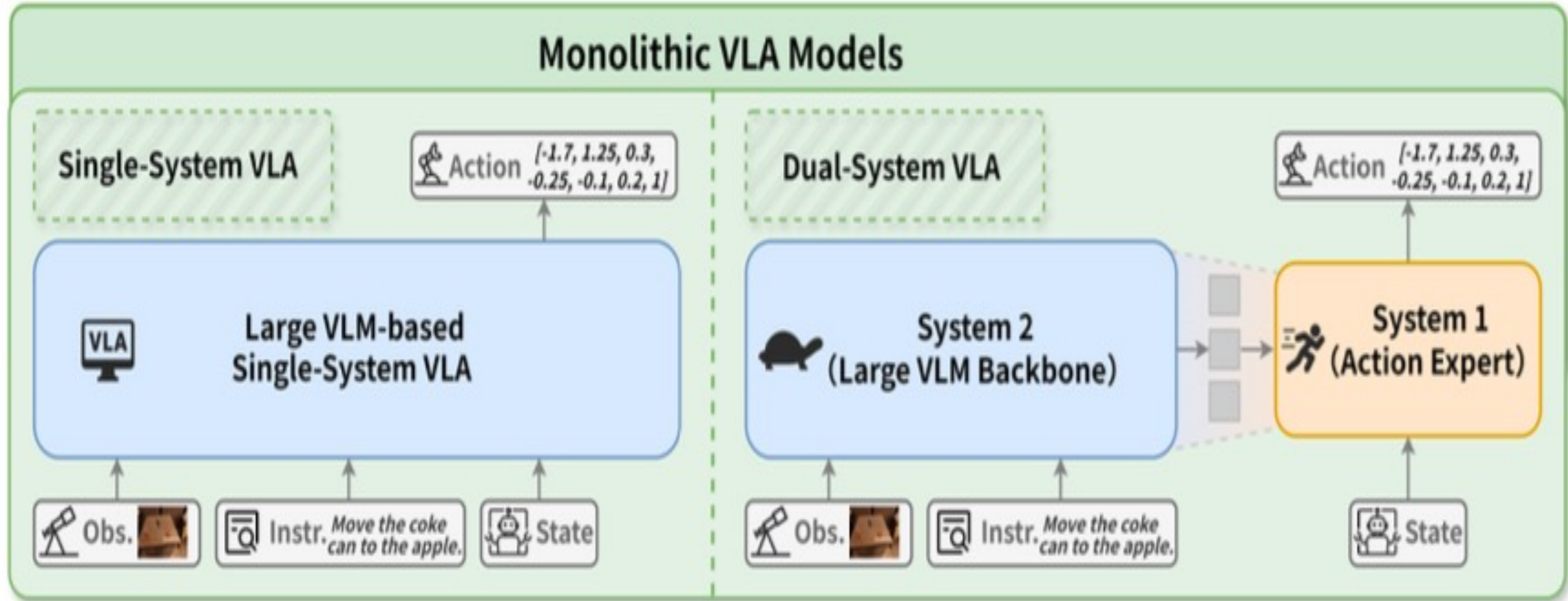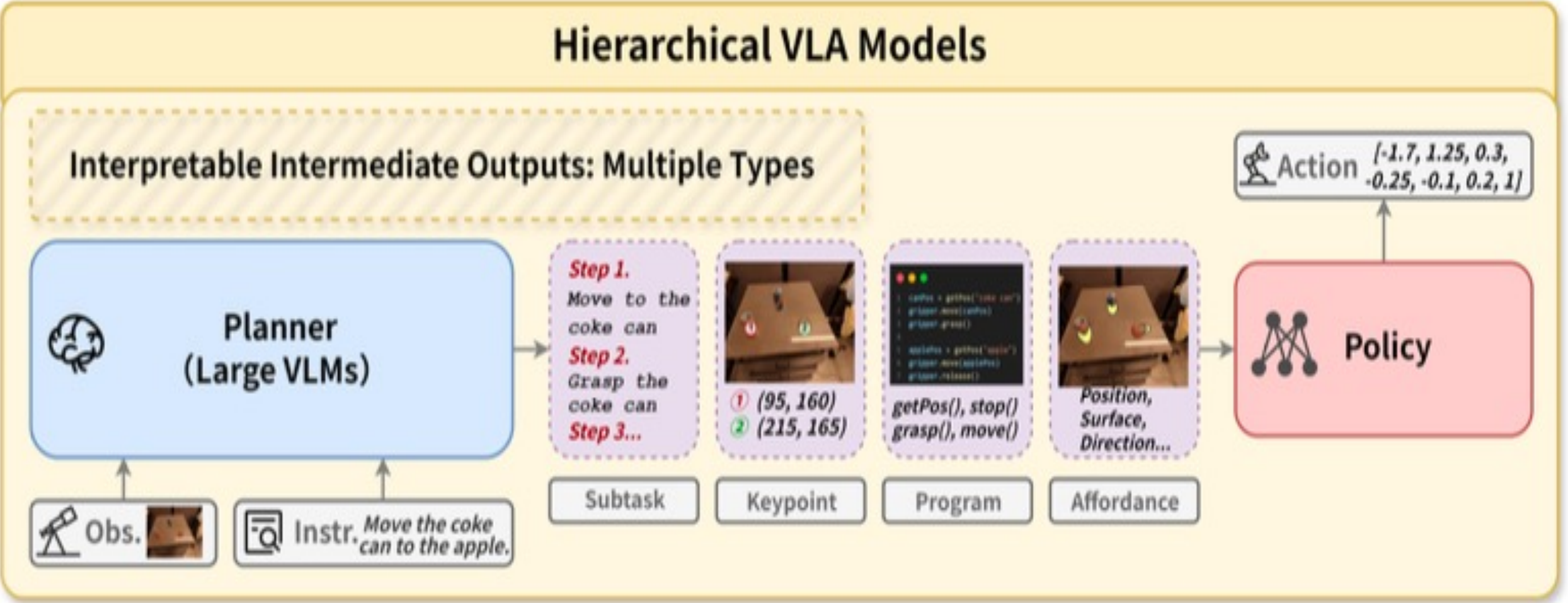| CALVIN | Habitat | robosuite | ManiSkill | RLBench |
| AI2-THOR | Habitat 2.0 | robomimic | ManiSkill 2 | COLOSSEUM |
| Meta-World | Habitat 3.0 | RoboCasa | ManiSkill 3 | SIMPLER |
|  |  | LIBERO | ManiSkill-HAB | RoboArena |

# Vision Language Action (VLA) Components

# Large VLM-based Vision-Language-Action Models
# Monolithic VLA Models

# Large VLM-based Vision-Language-Action Models
# Hierarchical VLA Models

# Large VLM-based Vision-Language-Action Models Paradigms in Monolithic Single-system 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).
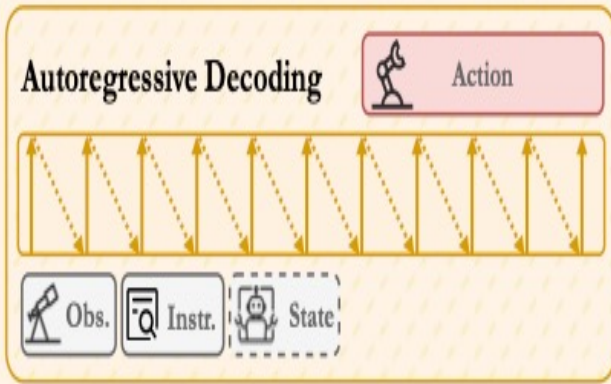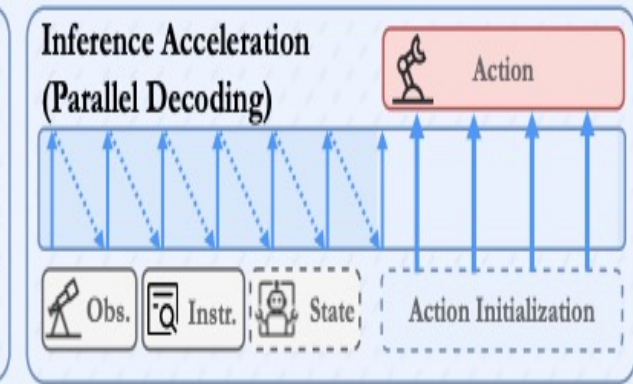
# Large VLM-based Vision-Language-Action Models Paradigms in Monolithic Dual-system Models



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).

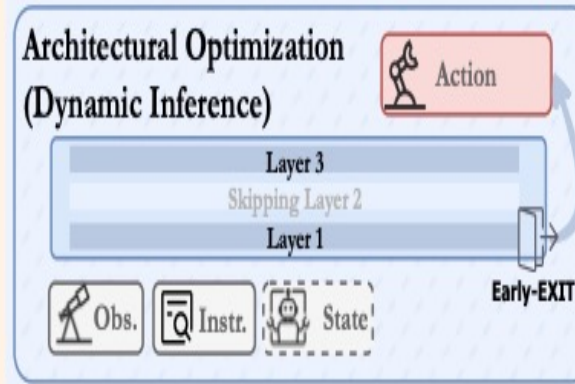# Large VLM-based Vision-Language-Action Models
## Hierarchical Models

# Schematic of the Unified VLA Training Data Format

# Benchmarking VLA Datasets
## by Task Complexity and Modality Richness

# Real-World Robot Datasets for VLA Research

| Name | Episodes | Skill | Task | Modality | Embodiment | Collection |
|---|---|---|---|---|---|---|
| QT-Opt | 580K | 1 (Pick) | NA | RGB | KUKA LBR iiwa | Learned |
| MT-Opt | 800K | 2 | 12 | RGB, L | 7 robots | Scripted, Learned |
| RoboNet | 162K | NA | NA | RGB | 7 robots | Scripted |
| BridgeData | 7.2K | 4 | 71 | RGB, L | WidowX 250 | Teleop |
| BridgeData V2 | 60.1K | 13 | NA | RGB-D, L | WidowX 250 | Teleop |
| BC-Z | 26.0K | 3 | 100 | RGB, L | Google EDR | Teleop |
| Language Table | 413K | 1 (Push) | NA | RGB, L | xArm | Teleop |
| RH20T | 110K | 42 | 147 | RGB-D, L, F, A | 4 robots | Teleop |
| RT-1 | 130K | 12 | 700+ | RGB, L | Google EDR | Teleop |
| OXE | 1.4M | 527 | 160,266 | RGB-D, L | 22 robots | Mixed |
| DROID | 76K | 86 | NA | RGB-D, L | Franka | Teleop |
| FuSe | 27K | 2 | 3 | RGB, L, T, A | WidowX 250 | Teleop |
| RoboMIND | 107K | 38 | 479 | RGB-D, L | 4 robots | Teleop |
| AgiBot World | 94K | 87 | 217 | RGB-D, L | AgiBot G1 | Teleop |

# Benchmarks for Vision-Language-Action Evaluation

**Simulation Environments:** Navigation (Nav), Manipulation (Manip), and Whole-Body Control (WBC)

| Name | Task | Scenes / Objects | Observation | Physics | Built Upon | Description |
|---|---|---|---|---|---|---|
| robosuite | Manip | NA / 10 | RGB-D, S | MuJoCo | NA | Modular framework, 11 tasks |
| robomimic | Manip | NA / NA | RGB | MuJoCo | robosuite | Offline learning, 8 tasks |
| RoboCasa | Manip | 120 / 2.5K | RGB | MuJoCo | robosuite | 100 kitchen tasks, photorealistic |
| LIBERO | Manip | NA / NA | RGB | MuJoCo | robosuite | 130 tasks in 4 task suites |
| Meta-World | Manip | 1 / 80 | Pose | MuJoCo | NA | 50 Manip tasks for Meta-RL |
| LeVERB-Bench | Nav, WBC | 4 / NA | RGB | PhysX | Isaac Sim | Humanoid control |

Source: Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. (2025) "Vision-language-action models for robotics: A review towards real-world applications." IEEE Access (2025).

# Benchmarks for Vision-Language-Action Evaluation

## Simulation Environments: Navigation (Nav), Manipulation (Manip), and Whole-Body Control (WBC)

| Name | Task | Scenes / Objects | Observation | Physics | Built Upon | Description |
|---|---|---|---|---|---|---|
| ManiSkill | Manip | NA / 162 | RGB-D, PC, S | PhysX | SAPIEN | 4 tasks, 36K demos |
| ManiSkill 2 | Manip | NA / 2.1K | RGB-D, PC | PhysX | ManiSkill | Extended task diversity |
| ManiSkill 3 | Nav, Manip, WBC | NA / NA | RGB-D, PC, S | PhysX | ManiSkill 2 | GPU-parallelized simulation |
| ManiSkill-HAB | Manip | 105 / 92 | RGB-D | PhysX | ManiSkill 3, Habitat 2.0 | HAB tasks from Habitat 2.0 |
| RoboTwin | Manip | NA / 731 | RGB-D | PhysX | SAPIEN | Dual-arm tasks |
| Ravens | Manip | NA / NA | RGB-D | PyBullet | NA | 10 tabletop tasks |
| VIMA-BENCH | Manip | NA / 29 | RGB, S | PyBullet | Ravens | 17 multimodal prompt tasks |
| LoHoRavens | Manip | 1 / 3 | RGB-D | PyBullet | Ravens | Long-horizon planning |
| CALVIN | Manip | 4 / 7 | RGB-D | PyBullet | NA | Long-horizon lang-cond tasks |

# Benchmarks for Vision-Language-Action Evaluation

**Simulation Environments:** Navigation (Nav), Manipulation (Manip), and Whole-Body Control (WBC)

| Name | Task | Scenes / Objects | Observation | Physics | Built Upon | Description |
|------|------|------------------|-------------|---------|------------|-------------|
| Habitat | Nav | 185 / NA | RGB-D, S | Bullet | NA | Fast, Nav only |
| Habitat 2.0 | Nav, Manip | 105 / 92 | RGB-D | Bullet | Habitat | Mobile manipulation (HAB) |
| Habitat 3.0 | Nav, Manip | 211 / 18K | RGB-D | Bullet | Habitat 2.0 | Human avatars support |
| RLBench | Manip | 1 / 28 | RGB-D, S | PyBullet | V-REP | Tiered task difficulty |
| THE COLOSSEUM | Manip | 1 / 107 | RGB-D | PyBullet | RLBench | 20 tasks, 14 env variations |
| AI2-THOR | Nav, Manip | NA / 118 | RGB-D, S | Unity | NA | Object states, task planning |
| CHORES | Nav | 191K / 40K | RGB | Unity | AI2-THOR | Shortest-path planning |
| SIMPLER | Manip | 4 / 17 | RGB | PhysX | SAPIEN, Isaac Sim | Real-to-sim evaluation |
| RoboArena | Manip | NA / NA | RGB | Real | NA | Distributed real-world evaluation |

Source: Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. (2025) "Vision-language-action models for robotics: A review towards real-world applications." IEEE Access (2025).

# Summary

- **Generative AI**

- **Agentic AI**

- **Physical AI (Robotics)**

# References

- Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson.
- Denis Rothman (2024), Transformers for Natural Language Processing and Computer Vision - Third Edition: Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3, 3rd ed. Edition, Packt Publishing
- Aurélien Géron (2022), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 3rd Edition, O'Reilly Media.
- Steven D'Ascoli (2022), Artificial Intelligence and Deep Learning with Python: Every Line of Code Explained For Readers New to AI and New to Python, Independently published.
- Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. (2022) "YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors." arXiv preprint arXiv:2207.02696.
- Chien-Yao Wang, I-Hau Yeh, and Hong-Yuan Mark Liao. (2025) "Yolov9: Learning what you want to learn using programmable gradient information." In European Conference on Computer Vision, pp. 1-21. Springer, Cham.
- Nidhal Jegham, Chan Young Koh, Marwan Abdelatti, and Abdeltawab Hendawi. (2024) "Evaluating the Evolution of YOLO (You Only Look Once) Models: A Comprehensive Benchmark Study of YOLO11 and Its Predecessors." arXiv preprint arXiv:2411.00201.
- Ranjan Sapkota, and Manoj Karkee. (2024) "Yolo11 and vision transformers based 3d pose estimation of immature green fruits in commercial apple orchards for robotic thinning." arXiv preprint arXiv:2410.19846.
- Yang Liu, Weixing Chen, Yongjie Bai, Xiaodan Liang, Guanbin Li, Wen Gao, and Liang Lin. (2024) "Aligning cyber space with physical world: A comprehensive survey on embodied ai." arXiv preprint arXiv:2407.06886.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry et al. (2021) "Learning transferable visual models from natural language supervision." In International Conference on Machine Learning, pp. 8748-8763. PMLR.
- Wonjae Kim, Bokyung Son, and Ildoo Kim. (2021) "Vilt: Vision-and-language transformer without convolution or region supervision."  In International Conference on Machine Learning, pp. 5583-5594. PMLR.
- Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R. Martin, Ming-Ming Cheng, and Shi-Min Hu. (2022) "Attention mechanisms in computer vision: A survey." Computational Visual Media ,:1-38.
- Valentin Bazarevsky, Ivan Grishchenko, Karthik Raveendran, Tyler Zhu, Fan Zhang, and Matthias Grundmann.(2020) "Blazepose: On-device real-time body pose tracking." arXiv preprint arXiv:2006.10204.
- Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna et al.(2025) "Gemini robotics: Bringing ai into the physical world." arXiv preprint arXiv:2503.20020 (2025).
- Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay et al. (2025) "Cosmos world foundation model platform for physical ai." arXiv preprint arXiv:2501.03575 (2025).
- Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke Zhu et al. (2025) "Foundation models in robotics: Applications, challenges, and the future." The International Journal of Robotics Research 44, no. 5 (2025): 701-739.
- Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Posner, and Yuke Zhu. (2025) "Vision-language-action models for robotics: A review towards real-world applications." IEEE Access (2025).
- Muhayy Ud Din, Waseem Akram, Lyes Saad Saoud, Jan Rosell, and Irfan Hussain. (2025) "Vision language action models in robotic manipulation: A systematic review." arXiv preprint arXiv:2507.10672 (2025)
- 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)..
- Min-Yuh Day (2025), Python 101, https://tinyurl.com/aintpupython101