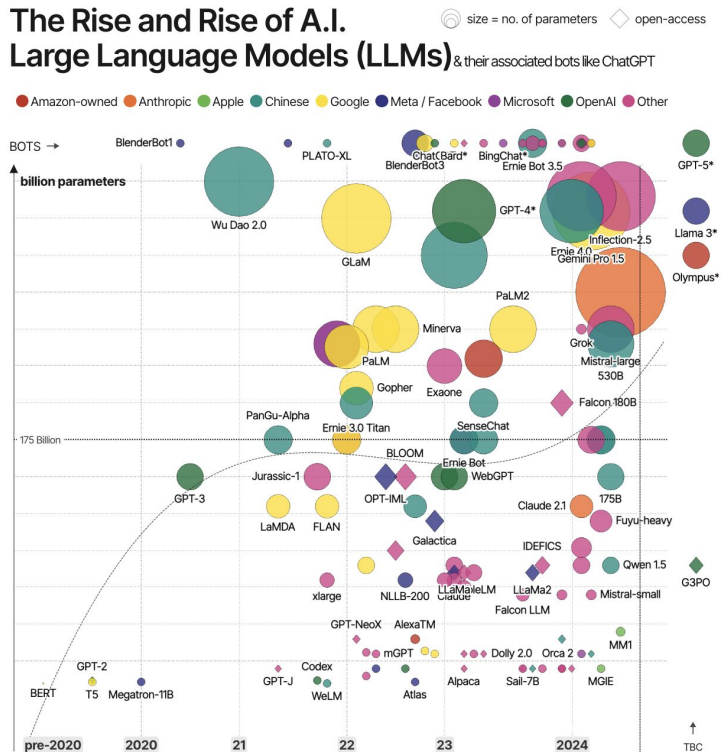# CS 294/194-196:
# Large Language Model Agents

# Teaching Staff

- **Instructor: Prof. Dawn Song**
- **(guest) Co-instructor: Dr. Xinyun Chen**
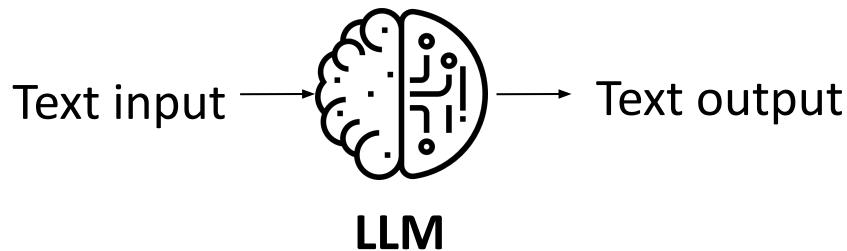- GSIs: Alex Pan & Sehoon Kim
- Readers: Tara Pande & Ashwin Dara

# Accelerated development of large language models (LLMs)



The Rise and Rise of A.I.
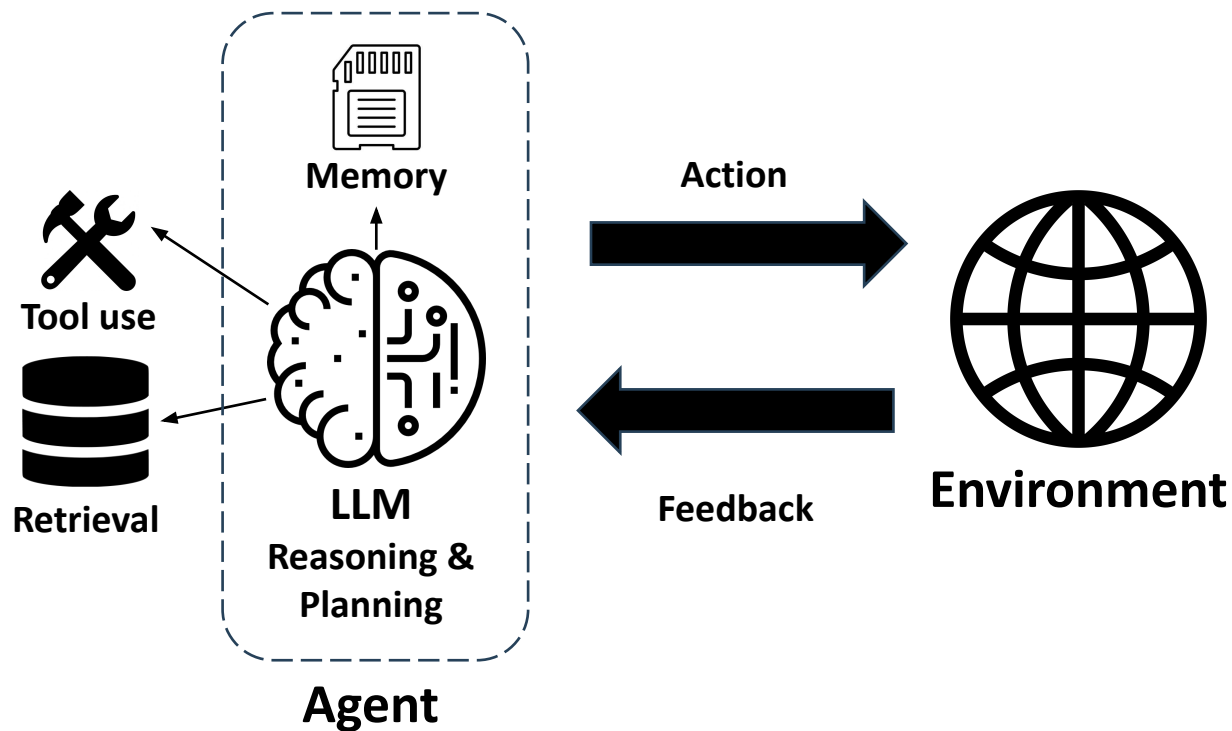Large Language Models (LLMs) & their associated bots like ChatGPT

size = no. of parameters    open-access

Amazon-owned   Anthropic   Apple   Chinese   Google   Meta / Facebook   Microsoft   OpenAI   Other

Text input → LLM → Text output

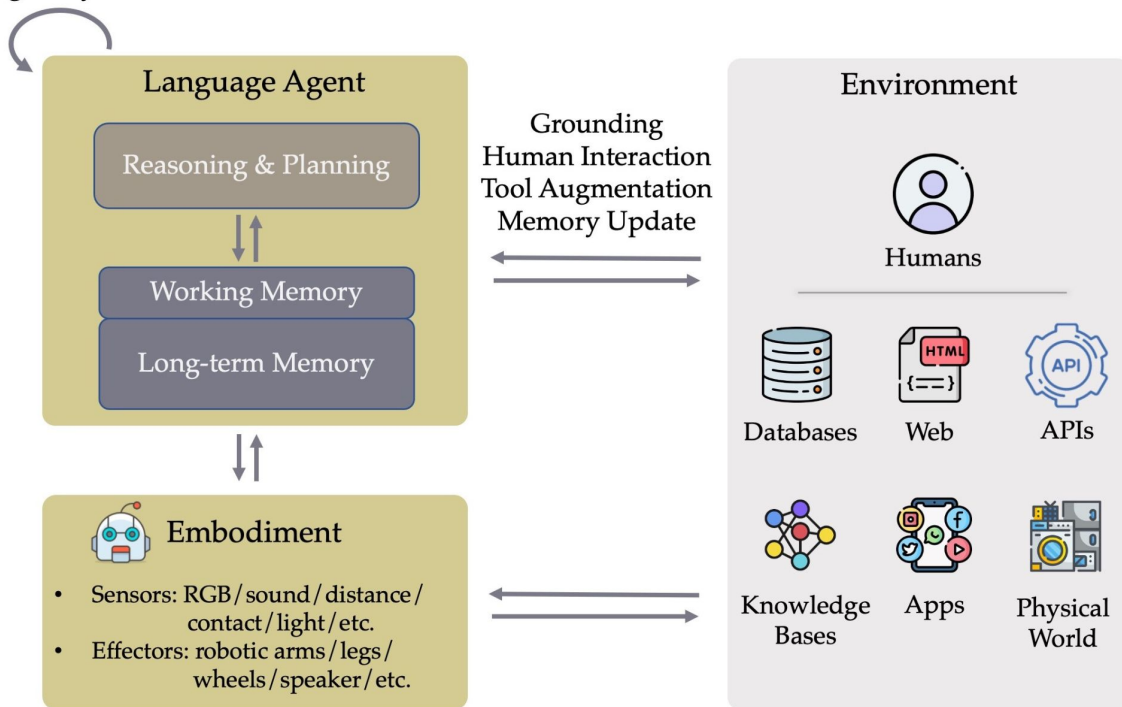# LLM agents: enabling LLMs to interact with the environment

# LLM Agents in Diverse Environments

# Multi-agent collaboration: division of labor for complex tasks



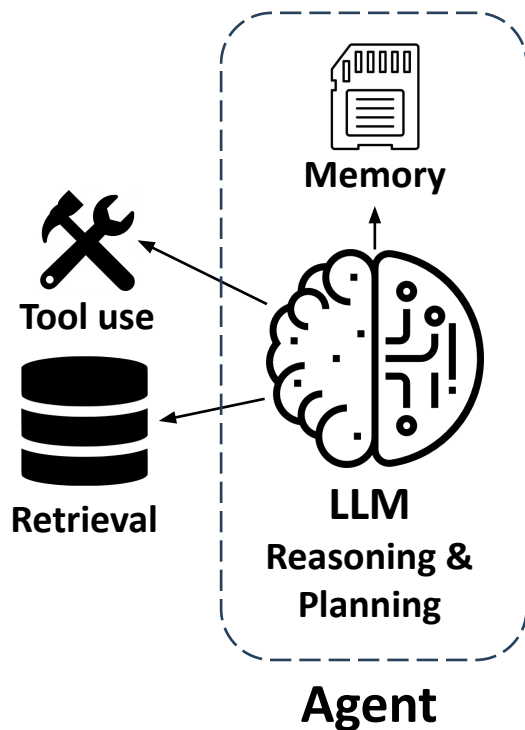**Specialized agents for different subtasks**
Autogen, CrewAI, CAMEL, Mixture-of-Agents,...
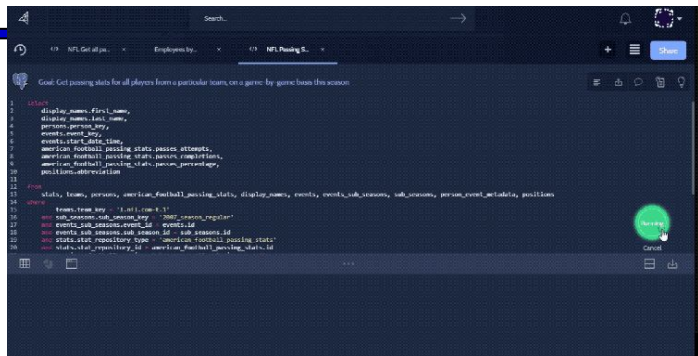
**Emergence of social behaviors with role-play LLMs**
Generative agents, Project Sid,...

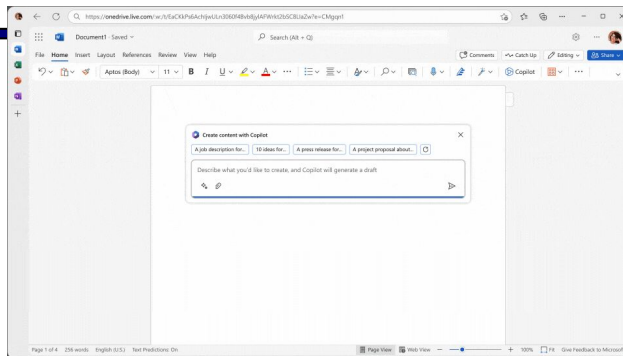# Why empowering LLMs with the agent framework



**Agent**

- Solving real-world tasks typically involves a trial-and-error process

- Leveraging external tools and retrieving from external knowledge expand LLM's capabilities

- Agent workflow facilitates complex tasks
    - Task decomposition
    - Allocation of subtasks to specialized modules
    - Division of labor for project collaboration
    - Multi-agent generation inspires better responses
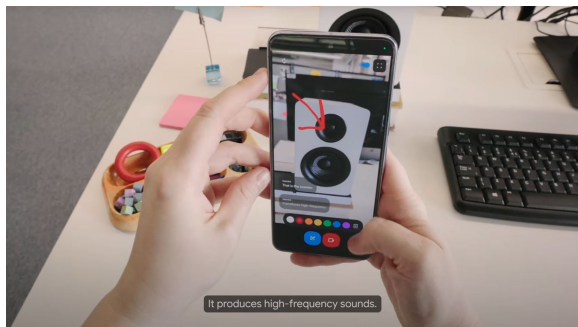
# LLM agents transformed various applications



**Code generation**
Cursor, GitHub Copilot, Devin, Replit,…



**Workflow automation**
Microsoft Copilot, Multi-On,…



**Personal assistant**
Google Astra, OpenAI GPT-4o,…



**Robotics**
Figure AI, Tesla Optimus,…

- Education
- Law
- Finance
- Healthcare
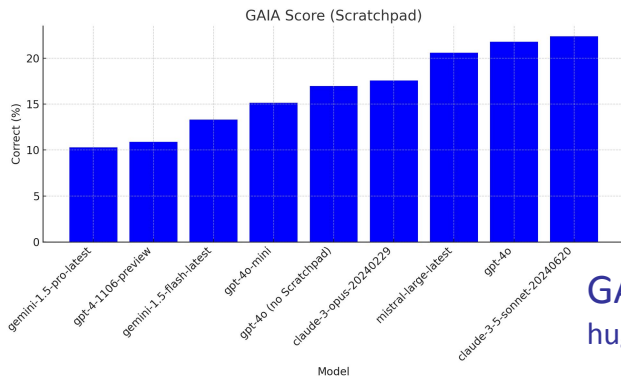- Cybersecurity

…

# LLM agents are improving



SWE-Bench (Jimenez*, Yang*, et al.)
swebench.com

GAIA (Mialon et al.)
huggingface.co/gaia-benchmark

WebArena
(Zhou et al.)
webarena.dev

# Challenges for LLM agent deployment in the wild

- Reasoning and planning
  - LLM agents tend to make mistakes when performing complex tasks end-to-end
- Embodiment and learning from environment feedback
  - LLM agents are not yet efficient at recovering from mistakes for long-horizon tasks
  - Continuous learning, self-improvement
  - Multimodal understanding, grounding and world models
- Multi-agent learning, theory of mind
- Safety and privacy
  - LLMs are susceptible to adversarial attacks, can emit harmful messages and leak private data
- Human-agent interaction, ethics
  - How to effectively control the LLM agent behavior, and design the interaction mode between humans and LLM agents

# Topics covered in this course

- Model core capabilities
  - Reasoning
  - Planning
  - Multimodal understanding
- LLM agent frameworks
  - Workflow design
  - Tool use
  - Retrieval-augmented generation
  - Multi-agent systems
- Applications
  - Software development
  - Workflow automation
  - Multimodal applications
  - Enterprise applications
- Safety and ethics

# Large Language Model Agents MOOC