# Introduction to AI Agents

January 29, 2026

Scalable Computational Intelligence Group (SCI)

Texas Advanced Computing Center (TACC)

# Outline

1. 1. Background on LLMs

2. 2. What even *is* an agent?

3. 3. Examples of agentic systems

# LLM is just text in and text out

**Gemma 2 - Instruct**

LLMs convert words or partial words to ids called

```
<start_of_turn>user
What is Cramer's Rule?<end_of_turn>
<start_of_turn>model
Cramer's Rule is ...<end_of_turn>
```
mum
t

LLM generation is recursive, so your output token amount is limited to:

Multi-turn conversations are just recorded as one long prompt with special tokens demarcating user and agent messages

The "Instruct" versions of models have been fine-tuned to understand this convention

$$\max output\ tokens = context\ window - input\ tokens$$

Input tokens

<u>Prompt</u>
Background Information:

Background Information:
My name is Dr. Orange

Instructions: Write a poem that exactly rhymes with my name

As I sat under0000000000000000000000000000000000000000
0

Context window
(e.g. 8k tokens)

LLM

under

# How fast can I run my model?

A little faster than you can read!

| Model Family | # of Parameters | iPhone 16 Pro Max tokens/s | Macbook Pro tokens/s | Vista GPU node tokens/s | Tejas tokens/s |
|---|---|---|---|---|---|
| Qwen3 | 1.7 billion | 13 | 73 | 313 | - |
| Llama3.1 | 8 billion | - | 25 | 173 | 1001 |
| Qwen3 | 30 billion | - | 36 | 130 | - |
| DeepSeek-R1 | 70 billion | - | - | 36 | - |
| DeepSeek-R1 | 671 billion | - | - | - | 135 |

Time to generate 100,000 tokens for agentic workloads:

| | 4 hours | 40 minutes | 10 minutes | 2 minutes |
|---|---|---|---|---|

Laptop is ~6x faster than phone
1 Vista GPU is ~4x faster than laptop
Tejas is ~5x faster than 1 Vista GPU

# Categories of tasks an LLM can accomplish

- Simple task examples
  - Summarize this text: {text}
  - Classify the sentiment of this statement as positive or negative: {statement}
  - Write a python script to implement a binary search of an array

- Complex task example
  - Reproduce the analytical model from the following paper, verify that it matches their results and then apply that model to the following data:

    **Paper Text**
    {paper_text}

    **Data**
    {data}

  - requires multiple steps
  - use of other software (web search, image analysis etc.)
  - math

[1] G. Mialon et al. "Gaia: a benchmark for general ai assistants" (2023) arXiv:2311.12983

# Outline

1. 1. Background on LLMs

2. 2. What even *is* an agent?

3. 3. Examples of agentic systems

# What are AI agents?

**AI Agent Definition:** "An artificial entity capable of perceiving its surroundings using sensors, making decisions, and then taking actions in response using actuators."

- The decision-making core of modern AI agents is either a reinforcement learning (RL) algorithm, or a LLM

- **RL-based agents** can struggle with complex tasks and their decision making can be inscrutable
  - e. g. in *long-horizon tasks* where the rewards during policy search are rare, the training times become extreme

- **LLM-based agents** utilize *goal decomposition to* break down goals into smaller subgoals
  - LLM can act as a planner
  - decides what action to take for a given subgoal
  - actions can be simple LLM tasks or external API calls to other software

Z. Xi et al. "The Rise and Potential of Large Language Model Based Agents: A Survey" *arXiv* (2023)

# Flow diagram of an agent: Perceive, Plan, and Act



Figure 5: Our generative agent architecture. Agents perceive their environment, and all perceptions are saved in a comprehensive record of the agent's experiences called the memory stream. Based on their perceptions, the architecture retrieves relevant memories and uses those retrieved actions to determine an action. These retrieved memories are also used to form longer-term plans and create higher-level reflections, both of which are entered into the memory stream for future use.

J. S. Park et al. "Generative Agents: Interactive Simulacra of Human Behavior" *UIST* (2023)

# Key functionalities an agent needs

1.  Planning and/or Reasoning through goal decomposition
    *   Goal decomposition: LLM writes a todo list
    *   ReAct: LLM iteratively decided next best step

2.  Acting through "tool" use
    *   LLM writes a command to trigger another piece of software

3.  Sensing current state of surroundings
    *   automatic snapshot of agent working area using software

4.  Memory retrieval
    *   Vector store/search
    *   Retrieval augmented generation (RAG)

# LLM agent - goal decomposition and tool use example



Goal

LLM Planner

Relevant Memory

User input text:
Look up the current weather at the Keck Observatory

Memories

List of available tools
[{'Tool name', 'tool description'}, ...]

List of past actions and results
[{'Action summary', 'result summary'}, ...]

# LLM agent - goal decomposition and tool use example



Goal → LLM Planner

Relevant Memory

Memories

**LLM prompt:**
Your goal is: {goal}

Relevant memories: {memories}

Break up your goal into smaller subgoals

# LLM agent - goal decomposition and tool use example



Goal

LLM Planner

Relevant Memory

Memories

Action Plan:
1. Use Tool: Search Engine to find information Keck Observatory location
2. Use Tool: Weather API to find weather at location

# LLM agent - goal decomposition and tool use example



Goal

LLM Planner

Relevant Memory

Memories

**LLM prompt:**
Your goal is: {goal}

Relevant memories: {memories}

Action plan: {action plan}

Perform step 1

# LLM agent - goal decomposition and tool use example

# LLM agent - goal decomposition and tool use example



Goal → LLM Planner

Memories → Relevant Memory → LLM Planner

**Final output text:**
Keck Observatory is located on the summit of Mauna Kea where the current weather is 16°C

# ReAct

- ReAct agents are AI systems that combine reasoning and action to solve problems by analyzing tasks and interacting with their environment.
- They operate in a loop of reasoning, acting, and observing, allowing for dynamic decision-making and real-time adjustments based on feedback.
- Contrasting ReAct agents with goal decomposition, ReAct agents only think 1 step ahead when reasoning. If they choose to execute an action they next ingest those results and reason on the next best step.

# LLM agent – ReAct and tool use example



Goal

LLM Brain

Relevant Memory

**User input text:**
Look up the current weather at the Keck Observatory

Memories

List of available tools
[{'Tool name', 'tool description'}, ...]

List of past actions and results
[{'Action summary', 'result summary'}, ...]

# LLM agent – ReAct and tool use example



Goal

LLM Brain

Relevant Memory

Memories

**LLM prompt:**
Your goal is: {goal}

Relevant memories: {memories}

Indicate what you should do next to solve this problem

# LLM agent – ReAct and tool use example



Goal

LLM Brain

Relevant Memory

Memories

Next Step:
1. Use Tool: Search Engine to find information Keck Observatory location

# LLM agent – ReAct and tool use example



Goal

LLM Brain

Relevant Memory

Memories

**LLM prompt:**
Your goal is: {goal}

Relevant memories: {memories}

Next step: {Next tool call}

Perform requested tool call.

# LLM agent – ReAct and tool use example

# LLM agent – ReAct and tool use example



Goal

LLM Brain

Relevant Memory

Memories

**LLM prompt:**
Your goal is: {goal}

Relevant memories: {memories… including keck observatory is located in Maunaki, Hawaii}

Indicate what you should do next to solve this problem

# LLM agent – ReAct and tool use example

# LLM agent – ReAct and tool use example



Goal → LLM brain → Tool

Relevant Memory

API call to get weather in Hawaii

Memories

Continue until LLM brain decides to send message back to user.

# LLM agent – ReAct and tool use example



Goal → LLM Brain → **Final output text:** Keck Observatory is located on the summit of Mauna Kea where the current weather is 16°C

Relevant Memory

Memories

# Outline

1. 1. Background on LLMs

2. 2. What even *is* an agent?

3. 3. Examples of agentic systems

# Comparing RL and LLM agents attempting the collect diamond challenge in Minecraft



Figure 2: **Comparison between RL-based method and our GITM.** RL agents try to map an complex goal directly to a sequence of low-level control signals, while our GITM leverages LLM to break down the goals and map them to structured actions for final control signals.

X. Zhu et al., "Ghost in the Minecraft" *arXiv* (2023)

# Extreme tool use example: HuggingGPT



An LLM-agent that looks up models on huggingface repo site

Attempts to solve sub goals via api calls to these models

Concatenates final results in the response

Y. Shen et al. "HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face" *arXiv* (2023)

# AI agents controlling multiple non-player characters in a "Sims-like" simulation

- Multi-agent simulation to study social dynamics and information transfer

- Each agent is initialized with a backstory in its memories

- Agents plan their schedules, move around, and interact with the environment and other characters through dialogs



Figure 1: Generative agents are believable simulacra of human behavior for interactive applications. In this work, we demonstrate generative agents by populating a sandbox environment, reminiscent of The Sims, with twenty-five agents. Users can observe and intervene as agents plan their days, share news, form relationships, and coordinate group activities.

J. S. Park et al. "Generative Agents: Interactive Simulacra of Human Behavior" *UIST* (2023)

apps
api
v2
src
lib
middleware
modules
api-keys
apps
atoms
auth
billing
booking-seat
conferencing
credentials
deployments
destination-calendars
controllers
inputs
outputs
services
TS destination-calendars.service.... 1
TS destination-calendars.module.ts
TS destination-calendars.repository.ts
email
event-types
jwt
kysely
memberships
oauth-clients
ooo
organizations
prisma
profiles
redis
router

TS destination-calendars.service.ts 1

nodules > destination-calendars > services > TS destination-calendars.service.ts >

```typescript
import { ConnectedCalendar, Calendar } from "@ee/calendars/c
import { CalendarsService } from "@ee/calendars/services/cal
import { DestinationCalendarsRepository } from "@modules/des
import { Injectable, NotFoundException } from "@nestjs/commor

@Injectable()
export class DestinationCalendarsService {
  constructor(
    private readonly calendarsService: CalendarsService,
    private readonly destinationCalendarsRepository: Destinat
  ) {}

  async updateDestinationCalendars(
    integration: string,
    externalId: string,
    userId: number,
    delegationCredentialId?: string
  ) {
    const userCalendars = await this.calendarsService.getCale
    const allCalendars: Calendar[] = userCalendars.connected
      .map((cal: ConnectedCalendar) => cal.calendars ?? [])
      .flat();
    const credentialId = allCalendars.find(
      (cal: Calendar) =>
        cal.externalId === externalId && cal.integration ===
    )?.credentialId;

    if (!delegationCredentialId && !credentialId) {
      throw new NotFoundException(`Could not find calendar ${
    }

    const delegatedCalendar = delegationCredentialId
      ? allCalendars.find(
        (cal: Calendar) =>
          cal.externalId === externalId &&
          cal.integration === integration &&
          cal.delegationCredentialId === delegationCredenti
      )
```

AUGMENT

Threads

New

0 files changed

Augment Memories    cal.com
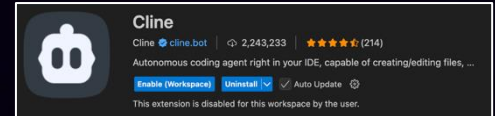
Ask or instruct Augment Agent

Agent    Auto    Default

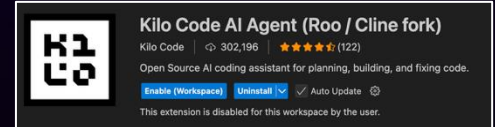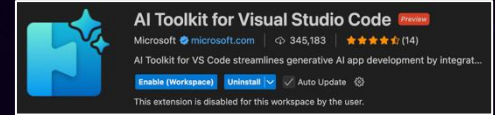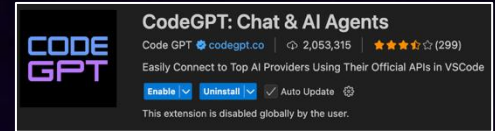main    1    0    Ln 1, Col 1    Spaces: 2    UTF-8    LF    TypeScript    Augment
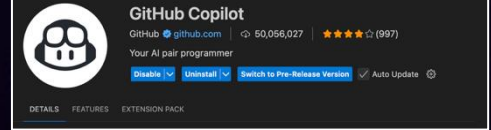
# AI Coding Tools

Deployed as:

- Full standalone IDEs
- Plugins for existing IDEs (VS Code etc.)

AI coding suggestions are provided at 3 different scales:

- Single line code completions
    - accepted by pressing Tab
- Highlighted code block rewrites
- Full codebase edits
    - Edits/additions across multiple files
    - Refactoring

# Today's Tutorial

- In today's tutorial we hope to apply the foundational knowledge just covered on AI Agents to use cases that cover different parts of researcher's scientific workflows.
- The remaining content of this tutorial will include the following labs:
  - ▶ Introduction to TACCs computing resources
  - ▶ Introduction to building AI agents
  - ▶ AI Agents for Simulations and Data Analysis
  - ▶ HPC integration of AI Agents