

WHAT ARE AI AGENTS?

Yogesh Haribhau Kulkarni

YHK

Outline

1 INTRO

2 ADK

3 REFS

About Me

Yogesh Haribhau Kulkarni

Bio:

- ▶ 20+ years in CAD/Engineering software development
- ▶ Got Bachelors, Masters and Doctoral degrees in Mechanical Engineering (specialization: Geometric Modeling Algorithms).
- ▶ Currently doing Coaching in fields such as Data Science, Artificial Intelligence Machine-Deep Learning (ML/DL) and Natural Language Processing (NLP).
- ▶ Feel free to follow me at:
 - ▶ Github (github.com/yogeshhk)
 - ▶ LinkedIn (www.linkedin.com/in/yogeshkulkarni/)
 - ▶ Medium (yogeshharibhaukulkarni.medium.com)
 - ▶ Send email to [yogeshkulkarni at yahoo dot com](mailto:yogeshkulkarni@yahoo.com)



Office Hours:
Saturdays, 2 to 5pm
(IST); Free-Open to all;
email for appointment.

Introduction

Classical idea of Agents

- ▶ Agents are autonomous
- ▶ Can look at their environment and analyze the situation
- ▶ Make comprehensive plans to achieve specific goals
- ▶ Actually take action to execute those plans
- ▶ Agents bridge the gap between answering and doing

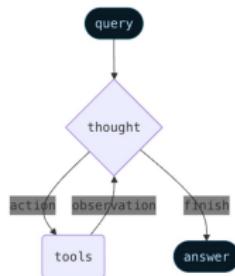
Welcome to AI Agents

- ▶ AI agents represent one of the most exciting frontiers in AI
- ▶ Unlike LLMs that just respond to prompts, agents do things
- ▶ Not just everyday chatbots, systems that reason, plan, and take action
- ▶ Can take on complex multi-step tasks autonomously or with human-in-loop
- ▶ Technology is advancing rapidly from conversational to agentic AI
- ▶ 2025 is the year of AI agents.

The Impossible Job

- ▶ If your boss tasks you with planning a massive get-together
- ▶ Must prepare a guest list and plan fancy menu
- ▶ Needs to find entertainment for the event
- ▶ A simple chatbot (with answering LLM, not reasoning LLM) cannot handle these complex requirements
- ▶ You need an AI agent, not just responses, but actions
- ▶ Agents have a Reasoning LLM as brain and tools as its hands-legs.
- ▶ AI agents improve with reasoning, acting, and memory components.
(ReAct = Reasoning + Acting)

ReAct (Reasoning + Acting) Agent



A ReAct agent selects tools by thinking about the request, selecting a tool, observing its result, and iterating until the request is fulfilled.

Based on paper:
[ReAct: Synergizing Reasoning and Acting in Language Models](#)
<https://arxiv.org/abs/2210.03629>

Diagram from Langchain docs: <https://docs.langchain.com/oss/python/langchain-agents>

(Ref: Python + AI Agents - Microsoft)

The Evolution of AI Capabilities

- ▶ **Traditional Programming:** Needed code to operate
- ▶ **Traditional ML:** Needed feature engineering
- ▶ **Deep Learning:** ML with Neural networks, no feature engineering
- ▶ **Foundational Models:** Many tasks single model
- ▶ **Agents (2024 . . .):** Can actually **do things**, not just talk

Why Does “Taking Action” Matter?

- ▶ In 2022, ChatGPT was revolutionary because AI felt conversational
- ▶ By 2024, people wanted more than conversation, they wanted **execution**
- ▶ Examples of what users now expect:
 - ▶ Instead of listing leads ? **email them directly**
 - ▶ Instead of summarizing docs ? **file and create workflow tasks**
 - ▶ Instead of suggesting products ? **customize landing pages**
- ▶ This shift from **information** to **action** defines the agent era

How Agents Work?

- ▶ Agent acts, take you from one state to the other state(ReAct paper: Reasoning and Action),
- ▶ It can plan and make decisions, provides value by workflow automation.
- ▶ Agents have access to tools (ToolFormer paper) e.g. Search APIs, booking, send email etc.
- ▶ Interacting of external environment and other Agents, etc.
- ▶ Memory to keep the history of conversations/actions done so far.
- ▶ May have human-in-loop to keep it sane in the wild-world.
- ▶ Agents were there from 1950's but they are effective because of LLMs.
- ▶ Agents are systems where LLMs dynamically direct their own paths and tool usage
- ▶ Agents can have autonomous-ness or predefined workflow paths
- ▶ Essential component in modern AI systems with varying degrees of autonomy

The Agent's Fundamental Game Loop

- ▶ Not a one-and-done action, but a continuous reasoning loop
- ▶ Similar to a programming while loop that keeps iterating
- ▶ **Thought:** Analyzes situation and plans next step
- ▶ **Action:** Calls specific tools to execute the plan
- ▶ **Observation:** Examines results of the action taken
- ▶ Cycle repeats: Thought → Action → Observation
- ▶ Continues until the task is completely accomplished
- ▶ This loop enables continuous adaptation and problem-solving

Agent's Inner Monologue

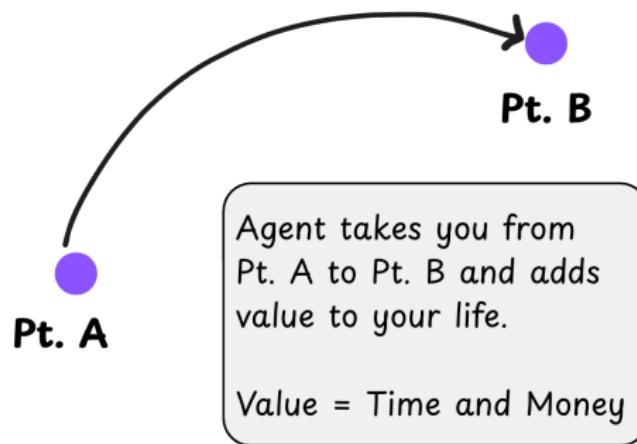
- ▶ Agents have visible thought processes before taking action
- ▶ Example: "User wants weather in New York. I have a tool for that".
- ▶ "My first move is to call the weather API".
- ▶ Internal planning step makes agents more than reactive programs
- ▶ Reasoning through problems before execution
- ▶ This deliberation distinguishes agents from simple scripts
- ▶ Shows intelligent decision-making rather than blind execution

How Do Agents Take Action?

- ▶ The magic lies in **tools** and **function calling**
- ▶ Agents are paired with APIs, plugins, or external systems
- ▶ Instead of just text responses, LLMs output structured commands:
 - ▶ “Call the send_email() function with these inputs...”
 - ▶ “Fetch records from CRM using this query...”
 - ▶ “Schedule a meeting for Tuesday at 2PM...”
- ▶ **Mental model:** LLM = brain, Tools = hands
- ▶ Without tools, agents just talk. With tools, they act.

Defining AI Agents with an example

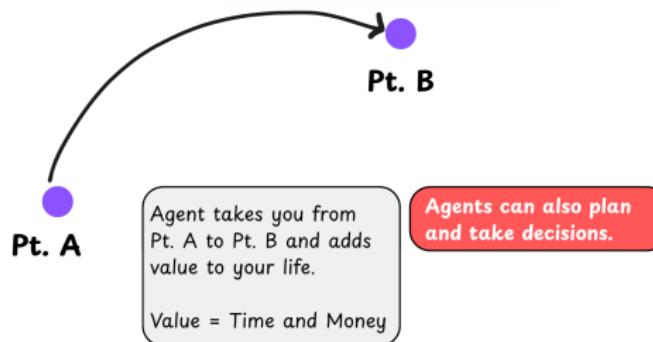
- ▶ Planning a trip involves many complex tasks
- ▶ Point A: Just discussing the trip
- ▶ Point B: All bookings and itinerary ready
- ▶ AI Agents aim to take you from A to B
- ▶ First requirement: Agent adds value by saving time/money



(Ref: Vizuara AI Agents Bootcamp)

Evolving Definition of Agents

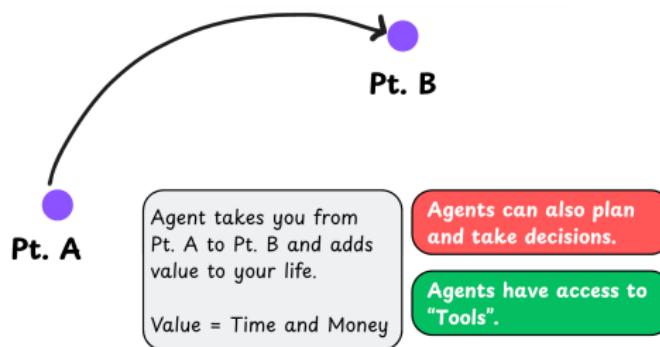
- ▶ Not all tools from A to B are agents (e.g., cars)
- ▶ Agents must plan and make decisions
- ▶ Example: Choosing flights based on budget
- ▶ Planning daily itinerary needs contextual judgment
- ▶ Second requirement: Agent includes decision-making ability



(Ref: Vizuara AI Agents Bootcamp)

Agents Need Tools

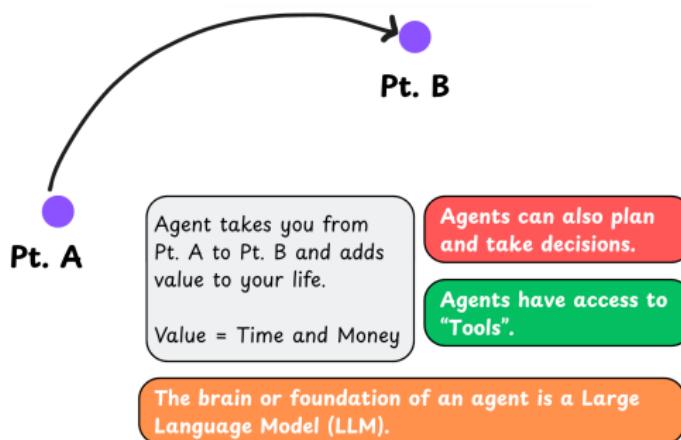
- ▶ Even self-driving cars plan but are not agents
- ▶ Agents need access to external tools
- ▶ Tools = Access to services (e.g., Gmail, Booking)
- ▶ Agents perform tasks using these tools
- ▶ Third requirement: Agent adds tool access to capabilities



(Ref: Vizuara AI Agents Bootcamp)

Rise of LLMs in Agents

- ▶ Transformers (2017) enabled powerful LLMs
- ▶ LLMs understand and generate human language
- ▶ Agents use LLMs for reasoning and planning
- ▶ LLMs enable understanding of webpages and writing emails
- ▶ Fourth requirement: Agents are LLMs with tools and planning ability



(Ref: Vizuara AI Agents Bootcamp)

What Is an Agent? (Technical Definition)

- ▶ Agent acts and takes you from one state to another, providing value through workflow automation
- ▶ Based on ReAct paradigm: **Reasoning + Acting**
- ▶ Key capabilities:
 - ▶ Can plan and make decisions
 - ▶ Has access to tools (search APIs, booking, email, etc.)
 - ▶ Interacts with external environments and other agents
 - ▶ Maintains memory of conversations and actions
 - ▶ May include human-in-the-loop for safety
- ▶ Agents existed since the 1950s but are now effective because of LLMs



Two Ways to Define Agents

Technical View:

- ▶ LLM (brain)
- ▶ + Tools (hands)
- ▶ + Planning (strategy)
- ▶ + Memory (context)
- ▶ + State management

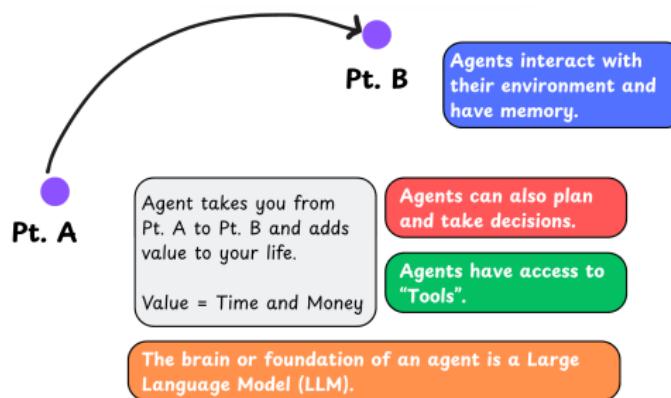
Business View:

- ▶ Systems that complete tasks end-to-end
- ▶ Focus on outcomes, not components
- ▶ Solve real-world problems
- ▶ Provide measurable value

Important: Today's agents are **engineering wrappers** around AI models, the intelligence comes from the LLMs, agents help act on that intelligence.

Final Definition of Agents

- ▶ Agents can learn from feedback and environment
- ▶ Agents interact with tools, humans, and websites
- ▶ They improve with experience (memory)
- ▶ Agents evolve over time via memory and feedback
- ▶ Fifth requirement: LLMs + Tools + Planning + Learning



(Ref: Vizuara AI Agents Bootcamp)

Understanding Agency

- ▶ Agency = Level of autonomy an agent has
- ▶ Low agency → less value
- ▶ High agency → high value
- ▶ More autonomous agents can handle complex tasks
- ▶ Agency is key to measuring agent usefulness

Agency Level	Description	Name	Example	🔗
●○○○○	Agent does not influence what happens next	Simple Processor	Grammar checker that rewrites sentences	
●●○○○	Agent determines basic control flow	Router	Customer Query → Tech Support or Sales	
●●●○○	Agent determines function execution	Tool Caller	A smart calendar assistant that spots "let's meet on Tuesday" and books the meeting	
●●●●○	Lays out a short plan and carries it step by step	Multi-step Agent	Personal travel planner that gathers flight options, hotels, local activities	
●●●●●	One agentic workflow starts another agentic workflow	Multi-agent	Travel planner agent → Booking agent ← Email agent	

(Ref: Vizuara AI Agents Bootcamp)



Tools: The Agent's Hands

- ▶ LLM is the agent's brain; tools are its hands
- ▶ Tools are functions agents call to interact with the world
- ▶ Can search web, run calculations, or query databases
- ▶ Bridge between thinking and doing in the real world
- ▶ Enable agents to move from planning to execution
- ▶ Without tools, agent thoughts would be useless
- ▶ Tools provide the interface to external systems and data

Two Ways Agents Use Tools

- ▶ **JSON Agent:** Writes structured work orders for other systems
- ▶ JSON approach requires external system to read and execute
- ▶ **Code Agent:** Directly writes and runs code blocks
- ▶ Code approach is more direct and powerful
- ▶ Code is naturally more expressive than JSON
- ▶ Can handle complex logic like loops and conditionals
- ▶ Modular, easier to debug, and taps into existing libraries
- ▶ Code agents can access thousands of APIs directly

Code Agent in Action

- ▶ Alfred needs a gala menu - agent has “suggest_menu” tool
- ▶ Agent doesn't just make up suggestions randomly
- ▶ Generates and runs actual code to call the specific tool
- ▶ Gets real results from the tool execution
- ▶ Super direct, efficient, and powerful way to take action
- ▶ Code generation enables precise tool interaction
- ▶ Results are based on actual tool capabilities, not hallucination

Advanced Pattern: Agentic RAG

- ▶ Traditional RAG: Retrieval Augmented Generation fetches info before answering
- ▶ Agentic RAG supercharges this with intelligent multi-step processes
- ▶ Turns retrieval itself into an agent-driven task
- ▶ Like having a master researcher on staff
- ▶ Doesn't just do one search - runs complete research processes
- ▶ Rewrites queries for better results and runs multiple searches
- ▶ Uses findings to inform next searches and validates accuracy
- ▶ Pulls from both private data and public web sources



Multi-Agent Systems: Digital Teams

- ▶ Complex problems like finding the missing Batmobile need teams
- ▶ Single agents can't handle web searches, calculations, and visualization
- ▶ Solution: Build teams of specialized agents
- ▶ Manager agent acts as project lead breaking down big tasks
- ▶ Delegates work to specialist agents with specific skills
- ▶ Web agent handles online searching while manager coordinates
- ▶ Manager focuses on big picture and final integration
- ▶ Digital division of labor for complex problem solving

The GAIA Benchmark Reality Check

- ▶ GAIA LLM benchmark tests real-world multi-step problems
- ▶ Measures how well systems handle tricky, complex tasks
- ▶ Results are eye-opening and show current limitations
- ▶ Humans solve these tasks with 92% accuracy
- ▶ Today's most advanced AI models: only 15% accuracy
- ▶ Massive 77% gap between human and AI performance
- ▶ This gap is exactly what agentic systems aim to close
- ▶ Shows the enormous potential for improvement



The Fundamental Shift

- ▶ Moving from conversational AI to agentic AI era
- ▶ Old paradigm: Ask questions, get answers
- ▶ New paradigm: State goals, systems plan and accomplish them
- ▶ Represents fundamental change in human-computer interaction
- ▶ Technology for building personal AI assistants advancing rapidly
- ▶ Not a question of “if” but “when” this becomes reality
- ▶ The future Alfred is closer than we think
- ▶ Prepare for AI that can handle “impossible” complex tasks

Papers that Shaped AI Agents

- ▶ Core research papers laid the foundation
- ▶ Introduced key frameworks and architectures
- ▶ Sparked recent boom in agent development
- ▶ Include Transformer and Agentic frameworks
- ▶ Major driving force in LLM-based agent systems

Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

Jason Wei Xuezhi Wang Dale Schuurmans Maarten Bosma

Brian Ichter Fei Xia Ed H. Chi Quoc V. Le Denny Zhou

Google Research, Brain Team
{jasonwei,dennyyzhou}@google.com

REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS

Shunyu Yan¹, Jeffrey Zhao², Dian Yu², Nan Du², Izhak Shafran², Karthik Narasimhan¹, Yuan Cao²

¹Department of Computer Science, Princeton University

²Google Research, Brain team

{shunuyu,karthikn}@princeton.edu
{jeffreyzhao,diyanyu,dunan,izhak,yuancao}@google.com

Toolformer: Language Models Can Teach Themselves to Use Tools

Timo Schick Jane Dwivedi-Yu Roberto Dessì[†] Roberta Raileanu
Maria Lomeli Luke Zettlemoyer Nicola Cancedda Thomas Scialom

Meta AI Research [†]Universitat Pompeu Fabra

Generative Agents: Interactive Simulacra of Human Behavior

Joon Sung Park
Stanford University
Stanford, USA
joonspk@stanford.edu

Joseph C. O'Brien
Stanford University
Stanford, USA
jobjrnc3@stanford.edu

Carrie J. Cai
Google Research
Mountain View, CA, USA
cja@ai.google.com

Meredith Ringel Morris
Google DeepMind
Seattle, WA, USA
mmemo@google.com

Percy Liang
Stanford University
Stanford, USA
pliang@cs.stanford.edu

Michael S. Bernstein
Stanford University
Stanford, USA
mbs@cs.stanford.edu

(Ref: Vizuara AI Agents Bootcamp)



Frameworks

Python AI Agent Frameworks Overview

- ▶ Python offers several frameworks for building single and multi-agent AI systems.
- ▶ These frameworks support workflows, tool use, reasoning, and human-in-the-loop capabilities.
- ▶ Common applications include autonomous assistants, research copilots, and workflow orchestration.

LangChain v1

- ▶ **Type:** Multi-agent, graph-based framework.
- ▶ **Built on:** LangGraph for agent state management and orchestration.
- ▶ **Features:**
 - ▶ Agent-centric architecture.
 - ▶ Workflow graphs for reasoning and tool execution.
 - ▶ Human-in-the-loop integration via LangSmith.
 - ▶ Supports OpenAI, Anthropic, Gemini, and local models.



Pydantic-AI

- ▶ **Type:** Lightweight single/multi-agent framework.
- ▶ **Focus:** Strong type safety using Pydantic for structured inputs/outputs.
- ▶ **Features:**
 - ▶ Declarative agent definition with validation.
 - ▶ Easy LLM orchestration and chaining.
 - ▶ Built for developers needing data integrity.
- ▶ **Ideal for:** AI applications requiring schema consistency and typed reasoning.

Crew AI

- ▶ **Type:** Multi-agent coordination framework.
- ▶ **Concept:** “Crew” of specialized agents working collaboratively.
- ▶ **Features:**
 - ▶ Role-based agent interactions.
 - ▶ Shared memory and task assignment.
 - ▶ Supports autonomous workflow creation.
 - ▶ Integrates with LangChain, OpenAI, and custom tools.
- ▶ **Use Case:** Complex multi-role task solving (research, planning, content generation).



Autogen

- ▶ **Developed by:** Microsoft Research.
- ▶ **Type:** Multi-agent conversational framework.
- ▶ **Features:**
 - ▶ Agent-to-agent and human-in-loop communication.
 - ▶ Supports function calling and code execution.
 - ▶ Built-in memory and orchestration logic.
 - ▶ Extendable with custom agents and tools.
- ▶ **Example Code:**

```
1 from autogen import AssistantAgent, UserProxyAgent  
2  
3 assistant = AssistantAgent("assistant", llm="gpt-4")  
4 user = UserProxyAgent("user", code_execution=True)  
5  
6 user.initiate_chat(assistant, message="Build a weather app.")
```

OpenAI Agents

- ▶ **Type:** Hosted agent platform by OpenAI (2024+).
- ▶ **Features:**
 - ▶ Create, configure, and deploy GPT-based agents.
 - ▶ Access via OpenAI API and GPTs interface.
 - ▶ Integrates with files, APIs, tools, and memory.
 - ▶ Supports human feedback and evaluation.
- ▶ **Example Use:** Deploy a custom assistant with knowledge base and tools.

Google ADK (Agent Development Kit)

- ▶ **Developed by:** Google DeepMind / Gemini ecosystem.
- ▶ **Type:** Multi-agent and workflow orchestration framework.
- ▶ **Features:**
 - ▶ Integrates Gemini 1.5 models with tool-use capabilities.
 - ▶ Supports perception, memory, planning, and dialogue.
 - ▶ Human-in-the-loop and autonomous modes.
 - ▶ Built for scalable, production-grade AI agents.
- ▶ **Use Case:** Building Gemini-powered assistants and reasoning systems.



Comparison Summary

- ▶ **LangChain v1:** Graph-based, modular, integrates with LangSmith.
- ▶ **Pydantic-AI:** Strong typing, data-safe orchestration.
- ▶ **Crew AI:** Role-based multi-agent collaboration.
- ▶ **Autogen:** Conversational multi-agent system by Microsoft.
- ▶ **OpenAI Agents:** Hosted GPT-based no-code agent deployment.
- ▶ **Google ADK:** Gemini-native multi-agent orchestration suite.



Final Thoughts

Why Agents?



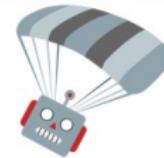
Flexible

Can **reason** based on personalized inputs, handle unexpected **edge cases**, and respond in **natural language**. **Outcome-driven** rather than path-driven.



Easy to Prototype

Reduced dev time **integrating** with external APIs and databases, and in building deterministic, “if this then that” systems.



Proactive

Can run an agent in response to a **user prompt**, on a **trigger** or **schedule** (autonomously).

(Ref: Google Cloud Labs H2 India Agents for All)

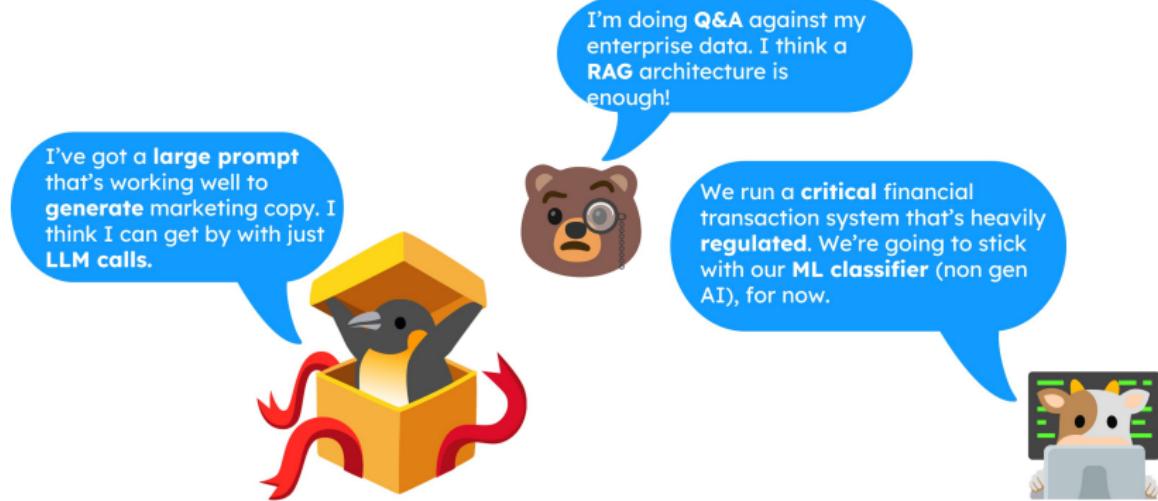
When to Use Agents?

- ▶ Best suited for tasks requiring flexibility and model-driven decision-making
- ▶ Consider tradeoffs: agents increase latency and cost for better task performance
- ▶ Recommended for open-ended problems with unpredictable steps
- ▶ Simple solutions preferred - single LLM calls with retrieval often sufficient



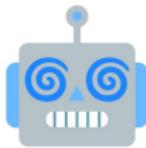
(Ref: Google Cloud Labs H2 India Agents for All)

When NOT to Use Agents?



(Ref: Google Cloud Labs H2 India Agents for All)

Agent Challenges



Predictability

Agents use LLMs to reason and respond.
LLMs are nondeterministic.



Stability

Complexity of using many tools, **error-handling, security, privacy...**

+ Framework + protocol landscape is **evolving rapidly**.



Operations

Tracking calls through a **distributed system** is hard, debugging is tricky, tracking **token throughput** is important (**costs**), **quota**, LLM reasoning can be **opaque**...

(Ref: Google Cloud Labs H2 India Agents for All)

Mitigating Predictability challenges

- ▶ choose reasoning models, when possible. (eg. gemini 2.5 flash)
- ▶ ground LLM in trusted data to reduce hallucinations (RAG, search)
- ▶ guide your agent to think step by step (chain of thought or few-shot prompting)
- ▶ use state management and memory.
- ▶ implement checks in the agent's workflow: tool-call request validation, safety filters, logic checks...
- ▶ test and evaluate your agent.

(Ref: Google Cloud Labgs H2 India Agents for All)



Mitigating Stability challenges

- ▶ use consistent tool abstractions via Model Context Protocol (MCP)
- ▶ when writing your agent's system prompt, be as detailed as possible. (tool descriptions)
- ▶ implement or leverage error-handling in your agent (circuit-breakers, retries, hand off to a human in the loop...)
- ▶ utilize state management, especially the workflow status field (recover from interruptions, prevent infinite loops)
- ▶ be willing to refactor and adjust with new framework features, tool interfaces...

(Ref: Google Cloud Labs H2 India Agents for All)



Mitigating Operations challenges

- ▶ in the dev phase, log verbosely.
- ▶ build robust test and evaluation datasets. automate these tests on every commit.
- ▶ gather necessary metrics, like token throughput to your agent's models, and response code distribution. create alerts if these trip certain thresholds.
- ▶ leverage agent framework's tracing features to track model and tool calls, find + address latency bottlenecks.

(Ref: Google Cloud Labgs H2 India Agents for All)



Future AI Applications

- ▶ What are future AI applications like?
 - ▶ **Generative:** Generate content like text and images
 - ▶ **Agentic:** Execute complex tasks on behalf of humans
- ▶ How do we empower every developer to build them?
 - ▶ **Co-Pilots:** Human-AI collaboration
 - ▶ **Autonomous:** Independent task execution
- ▶ 2024 is expected to be the year of AI agents



The Big Question

- ▶ Agentic AI technology is moving incredibly fast
- ▶ Personal AI assistants will soon be capable of complex tasks
- ▶ Think beyond simple queries to multi-step accomplishments
- ▶ Consider what “impossible” tasks you want to delegate
- ▶ What complex, party-of-the-century level challenge will you tackle?
- ▶ The era of AI that truly does rather than just discusses
- ▶ Prepare for AI assistants that can handle your biggest challenges

Google ADK (Agent Development Kit)

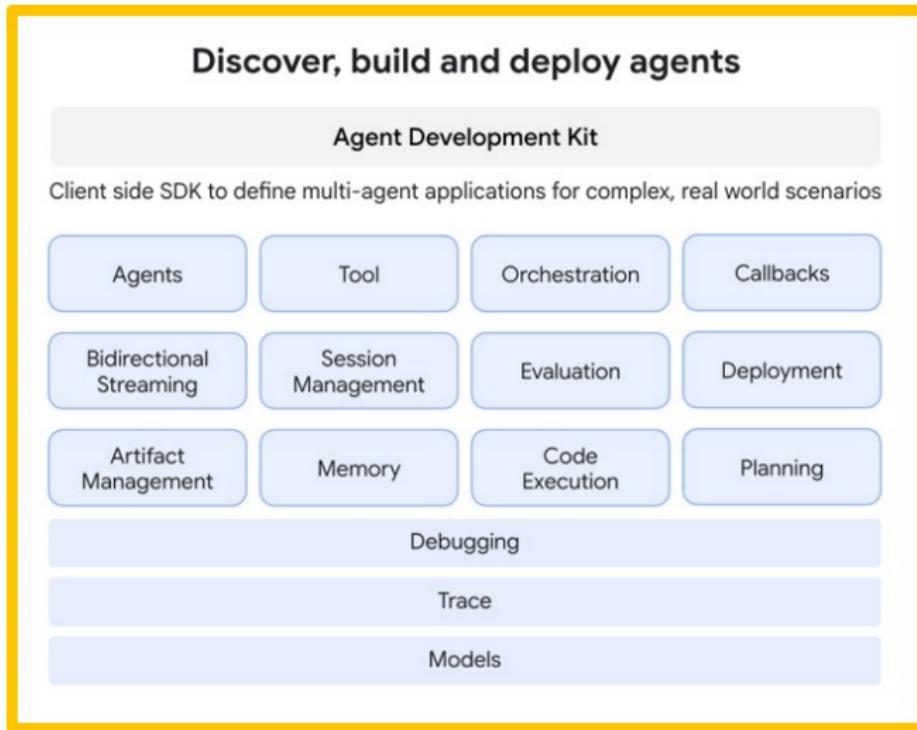
YHK

Introduction to ADK Framework

- ▶ ADK (Agent Development Kit) is Google's framework for building sophisticated AI agents and agentic applications
- ▶ Official documentation available at:
<https://google.github.io/adk-docs/>
- ▶ Designed for production-ready agentic systems with focus on scalability and reliability
- ▶ Supports multi-agent systems, tool integration, and flexible deployment options
- ▶ Built on Google's Gemini models with support for multiple LLM providers
- ▶ Framework emphasizes type safety, async operations, and enterprise-grade features



Key Components



(Ref: Google Cloud Labs H2 India Agents for All)

Key Components

- ▶ **Agents:** Core building blocks with instructions, tools, and model configuration
- ▶ **Tools:** Functions that agents can call to interact with external systems and APIs
- ▶ **Sessions:** Manage conversation context and state across multiple interactions
- ▶ **Runners:** Execute agent logic with support for streaming and async operations
- ▶ **Memory:** Store and retrieve conversation history and context
- ▶ **Multi-Agent Support:** Coordinate multiple specialized agents working together
- ▶ **Deployment:** Built-in support for FastAPI and serverless deployment



ADK Architecture Levels

- ▶ **Basic Agents:** Single agent with tools and instructions
- ▶ **Stateful Agents:** Agents with session management and memory
- ▶ **Multi-Agent Systems:** Multiple specialized agents collaborating
- ▶ **Agentic Workflows:** Complex orchestration with control flow
- ▶ **Production Systems:** Deployed services with monitoring and scaling
- ▶ Each level provides additional capabilities for building sophisticated AI applications

Key Features - Model Support & Flexibility

- ▶ **Gemini Integration:** Native support for Google's Gemini models (1.5 Pro, 2.0 Flash)
- ▶ **Multi-Provider:** Works with OpenAI, Anthropic, and other LLM providers
- ▶ **Type Safety:** Full TypeScript/Python type annotations for reliability
- ▶ **Async-First:** Built on async/await for high-performance concurrent operations
- ▶ **Streaming Support:** Real-time response streaming for better UX
- ▶ **Tool Calling:** Native function calling with structured input/output
- ▶ **Production Ready:** Built-in FastAPI integration and deployment patterns



Advanced Capabilities - Grounding & Multi-Modal

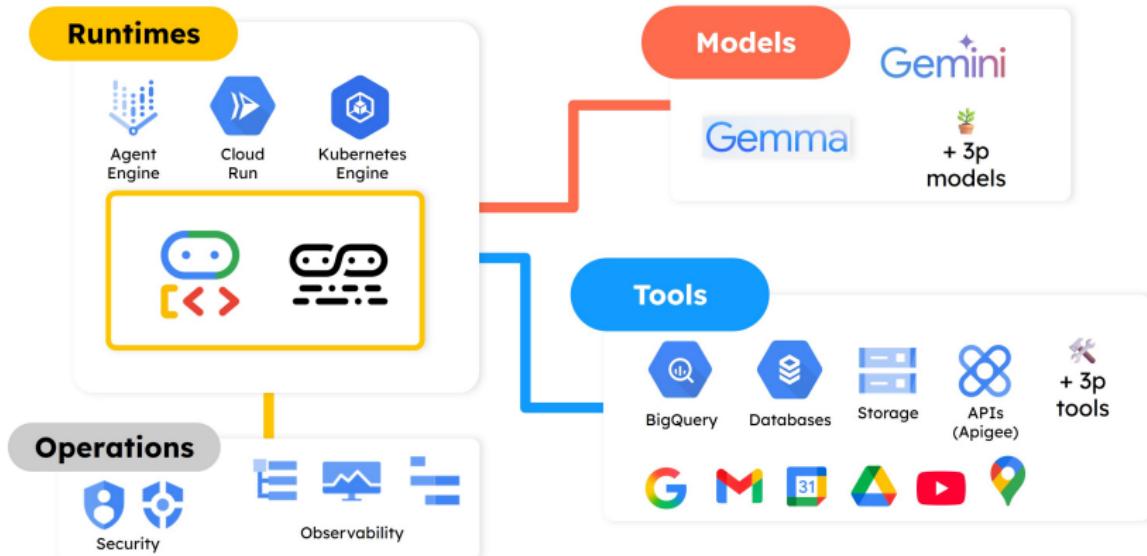
- ▶ **Google Search Grounding:** Connect agents to real-time web search results
- ▶ **Multi-Modal Input:** Support for text, images, audio, and video
- ▶ **Code Execution:** Built-in code interpreter for dynamic computation
- ▶ **Vertex AI Integration:** Enterprise features like RAG and vector search
- ▶ **Function Calling:** Structured tool execution with automatic parameter extraction
- ▶ **Safety Filters:** Built-in content safety and harm prevention



Enterprise Features - State & Deployment

- ▶ **Session Management:** Persistent conversation state across interactions
- ▶ **Cloud Integration:** Native support for Google Cloud services
- ▶ **Structured Output:** Pydantic models for type-safe responses
- ▶ **Error Handling:** Robust retry logic and error recovery
- ▶ **Monitoring:** Integration with Cloud Logging and Tracing
- ▶ **Authentication:** Built-in OAuth and API key management
- ▶ **Scalability:** Designed for high-throughput production workloads

Build Agents with Google Cloud-ADK



(Ref: Google Cloud Labs H2 India Agents for All)

Installation & Setup

- ▶ **Simple Installation:** Install via pip with minimal dependencies
- ▶ **CLI Tools:** Built-in commands for project creation and management
- ▶ **Project Generation:** Automatic scaffolding with best practices
- ▶ **API Key Creation:** Get the API key from
<https://aistudio.google.com/api-keys>
- ▶ **API Key Setup:** Configure Gemini API key via environment variable
- ▶ **Quick Start:** Start building agents immediately after installation

```
1 pip install google-adk
2 pip install google-adk --use-deprecated=legacy-resolver
3
4 # Create a new agent project
5 adk create my_agent
6
7 # Set API key in .env file
8 echo 'GOOGLE_API_KEY="YOUR_API_KEY"' > my_agent/.env # write your API key into
   an .env
```

Use CLI not VS code

Use Anaconda prompt with Admin permissions.

```
(google-adk) D:\Yogesh\GitHub\TeachingDataScience\Code\google-adk>adk create my_agent
Choose a model for the root agent:
1. gemini-2.5-flash
2. Other models (fill later)
Choose model (1, 2): 1
1. Google AI
2. Vertex AI
Choose a backend (1, 2): 1

Don't have API Key? Create one in AI Studio: https://aistudio.google.com/apikey

Enter Google API key [AIzaSyCwIflyXG0VC0l3W5Cj49M-hVajZDu7X0c]:
Agent created in D:\Yogesh\GitHub\TeachingDataScience\Code\google-adk\my_agent:
- .env
- __init__.py
- agent.py
```

Explore the agent project

- ▶ The created agent project has the following structure, with the agent.py file containing the main control code for the agent.
- ▶ The agent.py file contains a root_agent definition which is the only required element of an ADK agent. You can also define tools for the agent to use. Update the generated agent.py
- ▶ Run your agent using the adk run command-line tool.
- ▶ The ADK framework provides web interface you can use to test and interact with your agent.

```
my_agent/
2     agent.py      # main agent code
3     .env         # API keys or project IDs
4     __init__.py
5
6 adk run my_agent
7
8 adk web --port 8000 my_agent
```

Wrong time!!

```
[google-adk] D:\Yogesh\GitHub\TeachingDataScience\Code\google-adk>adk run my_agent
.log setup complete: C:\Users\yoges\AppData\Local\Temp\agents_log\agent.20251027_181028.log
to access latest log: tail -F C:\Users\yoges\AppData\Local\Temp\agents_log\agent.latest.log
D:\Yogesh\anaconda3\envs\google-adk\lib\site-packages\google\adk\cli\cli.py:154: UserWarning: [EXPERIMENTAL] InMemoryC
redentialService: This feature is experimental and may change or be removed in future versions without notice. It may in
duce breaking changes at any time.
    credential_service = InMemoryCredentialService()
D:\Yogesh\anaconda3\envs\google-adk\lib\site-packages\google\adk\auth\credential_service\in_memory_credential_service.
:33: UserWarning: [EXPERIMENTAL] BaseCredentialsService: This feature is experimental and may change or be removed in f
uture versions without notice. It may introduce breaking changes at any time.
    super().__init__()
running agent root_agent, type exit to exit.
D:\Yogesh\anaconda3\envs\google-adk\lib\site-packages\google\adk\cli\cli.py:98: UserWarning: [EXPERIMENTAL] App: This
feature is experimental and may change or be removed in future versions without notice. It may introduce breaking change
at any time.
    else App(name=session.app_name, root_agent=root_agent_or_app)
[user]: what can i do?
[root_agent]: I can tell you the current time in a specified city.
[user]: What is the current time in Pune?
[root_agent]: The current time in Pune is 10:30 AM.
[user]: bye
[root_agent]: Goodbye!
[user]:
[aborted]
```



Basic Agent Example - Simple Setup

- ▶ Create a basic agent with Gemini 2.0 Flash model
- ▶ Define tools using Python functions with type hints
- ▶ Agent automatically handles function calling and response generation
- ▶ Clean separation between tool definition and agent logic
- ▶ Synchronous execution for simple use cases

```
from adk import Agent
from adk.models import GeminiModel
import yfinance as yf
def get_stock_price(symbol: str) -> dict:
    """Get current stock price for a given symbol."""
    # Implementation using yfinance or similar
    stock = yf.Ticker(symbol)
    return {"symbol": symbol, "price": stock.info.get('currentPrice')}
agent = Agent(
    model=GeminiModel(model_name="gemini-2.0-flash-exp"),
    tools=[get_stock_price],
    instructions="You are a helpful financial assistant. Use tools when needed."
)
response = agent.run("What is NVIDIA's current stock price?")
print(response.text)
```



Agent with Multiple Tools

- ▶ Define multiple specialized tools for the agent with docstring
- ▶ Agent selects appropriate tools based on docstring and user query

```
from adk import Agent
from adk.models import GeminiModel
import yfinance as yf
2
def get_stock_price(symbol: str) -> float:
    """Get current stock price."""
    return yf.Ticker(symbol).info.get('currentPrice', 0)
4
def get_company_info(symbol: str) -> dict:
    """Get company information."""
    ticker = yf.Ticker(symbol)
10
    return { "name": ticker.info.get('longName'),
             "sector": ticker.info.get('sector'),
12
             "summary": ticker.info.get('longBusinessSummary')}
14
def get_analyst_recommendations(symbol: str) -> str:
    """Get analyst recommendations."""
16
    return yf.Ticker(symbol).recommendations.to_string()
18
agent = Agent(model=GeminiModel(model.name="gemini-2.0-flash-exp"),
              tools=[get_stock_price, get_company_info, get_analyst_recommendations],
20
              instructions="Provide detailed financial analysis using available tools.")
22
for chunk in agent.run_stream("Write a report on NVIDIA stock"):
    print(chunk.text, end="", flush=True)
24
26
```



Agent with Session Management

- ▶ Sessions maintain conversation history and context
- ▶ Enable multi-turn conversations with memory
- ▶ Store and retrieve previous interactions
- ▶ Support for persistent storage backends

```
1 from adk import Agent, Session
2 from adk.models import GeminiModel
3 from adk.storage import InMemoryStorage
4 import yfinance as yf
5
6 def get_stock_price(symbol: str) -> float:
7     """Get current stock price."""
8     return yf.Ticker(symbol).info.get('currentPrice', 0)
9 def get_company_info(symbol: str) -> dict:
10    """Get company information."""
11    ticker = yf.Ticker(symbol)
12    return { "name": ticker.info.get('longName'),
13            "sector": ticker.info.get('sector'),
14            "summary": ticker.info.get('longBusinessSummary') }
15
16 # Create storage for session history
17 storage = InMemoryStorage()
18 agent = Agent(
19     model=GeminiModel(model_name="gemini-2.0-flash-exp"),
20     tools=[get_stock_price, get_company_info],
21     instructions="You are a financial advisor. Remember previous conversations.")
```



Agent with Session Management

```
# Create a session for this conversation
2 session = Session(
    agent=agent,
4     storage=storage,
    session_id="user-123"
6 )
8 # Multi-turn conversation
response1 = session.run("What's NVIDIA's stock price?")
10 print(response1.text)

12 # Agent remembers previous context
response2 = session.run("How about their competitor AMD?")
14 print(response2.text)

16 # Session history is maintained
response3 = session.run("Compare both companies")
18 print(response3.text)
```



Agent with Google Search Grounding

- ▶ Google Search grounding provides real-time web information
- ▶ Automatically cites sources in responses
- ▶ Reduces hallucination with factual grounding
- ▶ Requires Vertex AI setup for production use

```
from adk import Agent
from adk.models import GeminiModel
from adk.extensions import GoogleSearchGrounding
agent = Agent(
    model=GeminiModel(
        model_name="gemini-2.0-flash-exp",
        extensions=[GoogleSearchGrounding()])
),
instructions="Provide well-researched answers with citations."
)
response = agent.run(
    "What are the latest developments in AI semiconductor technology?"
)
print(response.text)
# Response will include citations from search results
```



Structured Output with Pydantic

- ▶ Define expected output schema to ensure type-safe and structured responses
- ▶ Automatic validation ideal for integration with downstream systems

```
1 from adk import Agent
2 from adk.models import GeminiModel
3 from pydantic import BaseModel, Field
4 def get_stock_price(symbol: str) -> float:
5     :
6 def get_company_info(symbol: str) -> dict:
7     :
8 class StockAnalysis(BaseModel):
9     symbol: str = Field(description="Stock ticker symbol")
10    current_price: float = Field(description="Current stock price")
11    recommendation: str = Field(description="Buy/Hold/Sell recommendation")
12    reasoning: str = Field(description="Analysis reasoning")
13    agent = Agent(
14        model=GeminiModel(model_name="gemini-2.0-flash-exp"),
15        tools=[get_stock_price, get_company_info],
16        output_schema=StockAnalysis,
17        instructions="Analyze the stock and provide structured recommendation.")
18    result: StockAnalysis = agent.run("Analyze NVIDIA stock")
19    print(f"Symbol: {result.symbol}")
20    print(f"Price: ${result.current_price}")
21    print(f"Recommendation: {result.recommendation}")
22    print(f"Reasoning: {result.reasoning}")
23
```



Multi-Agent System - Architecture

- ▶ **Specialized Agents:** Each agent focuses on specific domain expertise
- ▶ **Orchestration:** Coordinator agent manages task delegation
- ▶ **Parallel Execution:** Multiple agents can work simultaneously
- ▶ **Result Aggregation:** Combine outputs from multiple agents
- ▶ **Scalable Design:** Handle complex workflows through collaboration

Multi-Agent Implementation - Part 1

```
from adk import Agent
from adk.models import GeminiModel

def web_search(query: str) -> str:
    """Search the web for information."""
    # Implementation using search API
    pass

def get_stock_data(symbol: str) -> dict:
    """Get comprehensive stock data."""
    import yfinance as yf
    ticker = yf.Ticker(symbol)
    return {
        "price": ticker.info.get('currentPrice'),
        "recommendations": ticker.recommendations.tail(5).to_dict()
    }

web_agent = Agent(
    name="Web Agent",
    model=GeminiModel(model_name="gemini-2.0-flash-exp"),
    tools=[web_search],
    instructions="Search the web and provide sourced information."
)

finance_agent = Agent(
    name="Finance Agent",
    model=GeminiModel(model_name="gemini-2.0-flash-exp"),
    tools=[get_stock_data],
    instructions="Analyze financial data and present in clear tables."
)
```

Multi-Agent Implementation - Part 2

- ▶ Coordinator agent orchestrates multiple specialized agents
- ▶ Delegates tasks to appropriate agents based on requirements
- ▶ Aggregates and synthesizes results from multiple sources

```
from adk import Agent, MultiAgentOrchestrator
2 from adk.models import GeminiModel

4 orchestrator = MultiAgentOrchestrator(
    agents=[web_agent, finance_agent],
6     coordinator=Agent(
        model=GeminiModel(model_name="gemini-2.0-flash-exp"),
8     instructions="""You coordinate a team of specialized agents.
        — Use Web Agent for market news and trends
        — Use Finance Agent for stock data and analysis
        — Synthesize their outputs into comprehensive reports."""))
10
12 # Run multi-agent workflow
14 result = orchestrator.run(
    "Provide a comprehensive analysis of AI semiconductor companies "
16     "including market outlook and financial performance")
18 print(result.text)

20 # Access individual agent outputs
for agent_name, agent_result in result.agent_outputs.items():
22     print(f"\n{agent_name} output:")
23     print(agent_result.text)
24
```

YHK

FastAPI Deployment

Deploy agents as REST APIs using FastAPI

```
1 from fastapi import FastAPI, HTTPException
2 from fastapi.responses import StreamingResponse
3 from adk import Agent, Session
4 from adk.models import GeminiModel
5 from pydantic import BaseModel
6 import uvicorn
7
8 app = FastAPI(title="ADK Agent API")
9
10 agent = Agent(
11     model=GeminiModel(model_name="gemini-2.0-flash-exp"),
12     tools=[get_stock_price, get_company_info],
13     instructions="You are a financial assistant API."
14 )
15
16 class QueryRequest(BaseModel):
17     message: str
18     session_id: str
19
```

FastAPI Deployment

```
1  @app.post("/chat")
2  async def chat(request: QueryRequest):
3      session = Session(agent=agent, session.id=request.session.id)
4      response = await session.run.async(request.message)
5      return {"response": response.text}
6
7  @app.post("/chat/stream")
8  async def chat_stream(request: QueryRequest):
9      session = Session(agent=agent, session.id=request.session.id)
10
11     async def generate():
12         async for chunk in session.run.stream.async(request.message):
13             yield chunk.text
14
15     return StreamingResponse(generate(), media_type="text/plain")
16
17 if __name__ == "__main__":
18     uvicorn.run(app, host="0.0.0.0", port=8000)
```

Vertex AI Integration

```
1 from adk import Agent
2 from adk.models import VertexAIModel
3 from adk.extensions import VertexAIVectorSearch
4 from google.cloud import aiplatform
5
6 # Initialize Vertex AI
7 aiplatform.init(project="your-project-id", location="us-central1")
8
9 # Create agent with Vertex AI backend
10 agent = Agent(
11     model=VertexAIModel(
12         model_name="gemini-2.0-flash-exp",
13         project="your-project-id",
14         location="us-central1"
15     ),
16     extensions=[
17         VertexAIVectorSearch(
18             index_endpoint="your-index-endpoint",
19             deployed_index_id="your-index-id"
20         )
21     ],
22     instructions="Answer questions using both your knowledge and the vector store."
23 )
24
25 response = agent.run("What are our company policies on remote work?")
26 print(response.text)
27
```



Error Handling & Reliability

```
from adk import Agent
from adk.models import GeminiModel
from adk.exceptions import ToolExecutionError, ModelError
import logging
import yfinance as yf
logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(__name__)

def safe_tool_wrapper(func):
    """Decorator for safe tool execution with fallback."""
    def wrapper(*args, **kwargs):
        try:
            return func(*args, **kwargs)
        except Exception as e:
            logger.error(f"Tool {func.__name__} failed: {e}")
            return {"error": str(e), "fallback": True}
    return wrapper

@safe_tool_wrapper
def get_stock_price(symbol: str) -> dict:
    """Get stock price with error handling."""
    return {"price": yf.Ticker(symbol).info['currentPrice']}
```

Error Handling & Reliability

```
1 agent = Agent(  
2     model=GeminiModel(  
3         model.name="gemini-2.0-flash-exp",  
4         max_retries=3,  
5         timeout=30  
6     ),  
7     tools=[get_stock_price],  
8     instructions="Handle errors gracefully and inform users."  
9 )  
10  
11 try:  
12     response = agent.run("What's the price of NVDA?")  
13     print(response.text)  
14 except ModelError as e:  
15     logger.error(f"Model error: {e}")  
16     print("Sorry, I'm having trouble processing your request.")  
17
```



Framework Comparison: LangGraph vs Agno vs ADK

Feature	LangGraph	Agno	ADK
Primary Focus	Graph-based workflows	Multi-agent systems	Enterprise AI agents
Developer	LangChain AI	Agno (ex-phidata)	Google
Architecture	State machines & graphs	Agent teams	Agent orchestration
Model Support	100+ providers	23+ providers	Gemini + multi-provider
Learning Curve	Steep (graph concepts)	Moderate	Moderate
Setup Time	Complex	Minutes	Minutes
Performance	Good	3μs instantiation	Optimized for Gemini
State Management	Built-in (checkpoints)	Session storage	Session & Cloud
Memory	Memory modules	Built-in drivers	In-memory & Cloud
Multi-Agent	Via subgraphs	Native teams	Orchestrator pattern
Workflow Control	Explicit graphs	Coordinate mode	Coordinator agent
Streaming	Yes	Yes	Yes (native async)
RAG Support	Vector stores	20+ vector DBs	Vertex AI Search
Deployment	Custom	FastAPI built-in	FastAPI + Cloud Run
Monitoring	LangSmith	agno.com	Cloud Console
Production Ready	Yes	Yes	Yes (enterprise)
Best For	Complex workflows	Fast prototyping	Google Cloud users
Unique Feature	Graph visualization	6.5KB/agent memory	Search grounding

Best Practices

- ▶ **Start Simple:** Begin with basic agents before adding complexity
- ▶ **Clear Instructions:** Provide specific, actionable instructions to agents
- ▶ **Tool Design:** Keep tools focused and well-documented with type hints
- ▶ **Error Handling:** Always implement proper error handling and retries
- ▶ **Testing:** Test agents thoroughly before production deployment
- ▶ **Monitoring:** Implement logging and monitoring for production systems
- ▶ **Security:** Validate inputs and sanitize outputs
- ▶ **Cost Management:** Monitor API usage and implement rate limiting

Getting Started - Next Steps

- ▶ **Documentation:** Explore comprehensive guides at
<https://google.github.io/adk-docs/>
- ▶ **Examples Repository:** Check out example projects and templates
- ▶ **Community:** Join Google Cloud community for support
- ▶ **Start Building:** Begin with simple agents and iterate
- ▶ **Cloud Integration:** Leverage Google Cloud services for production
- ▶ **Vertex AI:** Explore enterprise features for advanced use cases
- ▶ **Monitoring:** Use Cloud Console for production monitoring

Resources & References

- ▶ **Official Documentation:** <https://google.github.io/adk-docs/>
- ▶ **GitHub Repository:** <https://github.com/google/adk>
- ▶ **Gemini API:** <https://ai.google.dev/>
- ▶ **Vertex AI:** <https://cloud.google.com/vertex-ai>
- ▶ **Google Cloud:** <https://cloud.google.com/>
- ▶ **Community Support:** Google Cloud Community forums

References

- ▶ CS 194/294-196 (LLM Agents) - Lecture 3, Chi Wang and Jerry Liu
- ▶ LLM Powered Autonomous Agents Lil'Log
- ▶ Power of Autonomous AI Agents - Yogesh Kulkarni
- ▶ Microsoft AutoGen- Yogesh Kulkarni
- ▶ Microsoft AutoGen using Open Source Models- Yogesh Kulkarni
- ▶ A CAMEL ride - Yogesh Kulkarni
- ▶ Autonomous AI Agents (LLM, VLM, VLA) - Code Your Own AI
- ▶ Awesome LLM-Powered Agent
<https://github.com/hyp1231/awesome-lilm-powered-agent>
- ▶ Autonomous Agents (LLMs). Updated daily
<https://github.com/tmgthb/Autonomous-Agents>



Thanks ...

- ▶ Search "**Yogesh Haribhau Kulkarni**" on Google and follow me on LinkedIn, GitHub, Medium
- ▶ Office Hours: Saturdays, 2 to 3 pm (IST); Free-Open to all; email for appointment.
- ▶ yogeshkulkarni at yahoo dot com
- ▶ Call + 9 1 9 8 9 0 2 5 1 4 0 6



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://medium.com/@yogeshharibhaukularkarni>)



(<https://www.github.com/yogeshhk/>)

YHK

Pune AI Community (PAIC)

- ▶ Two-way communication:

- ▶ Website puneaicommunity dot org
- ▶ Email puneaicommunity at gmail dot com
- ▶ Call + 9 1 9 8 9 0 2 5 1 4 0 6
- ▶ LinkedIn:
<https://linkedin.com/company/pune-ai-community>

- ▶ One-way Announcements:

- ▶ Twitter (X) @puneaicommunity
- ▶ Instagram @puneaicommunity
- ▶ WhatsApp Community: Invitation Link
<https://chat.whatsapp.com/LluOrhyEzuQLDr25ixZ>
- ▶ Luma Event Calendar: puneaicommunity



Website

- ▶ Contribution Channels:

- ▶ GitHub: Pune-AI-Community and puneaicommunity
- ▶ Medium: pune-ai-community
- ▶ YouTube: @puneaicommunity

Pune AI Community (PAIC) QR codes



Website



Medium Blogs



Twitter-X



LinkedIn Page



Github Repository



WhatsApp Invite



Luma Events



YouTube Videos



Instagram

