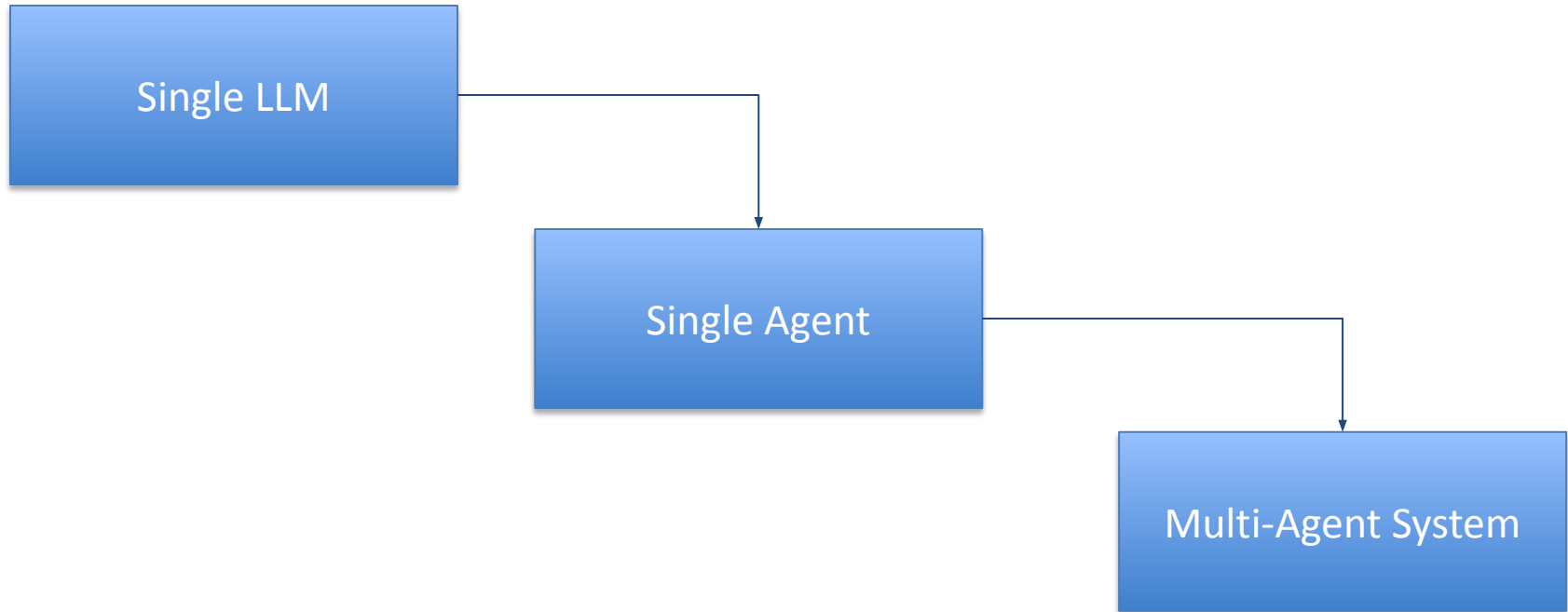


# Building Multi-Agent Systems with AutoGen

-Annu Sachan  
Principal Manager-AI/ML

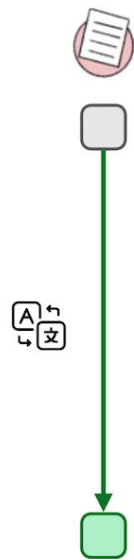
# Conceptual Progression



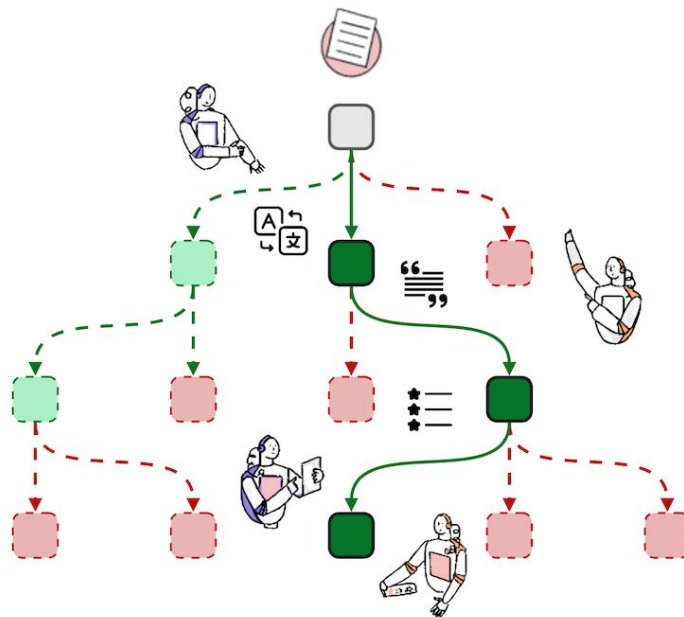
# Why Multi-Agent Systems?

- Single LLM limitations
- Role-based agents
- Coordination and specialization
- Scalable reasoning

# Single LLM vs Multi Agent



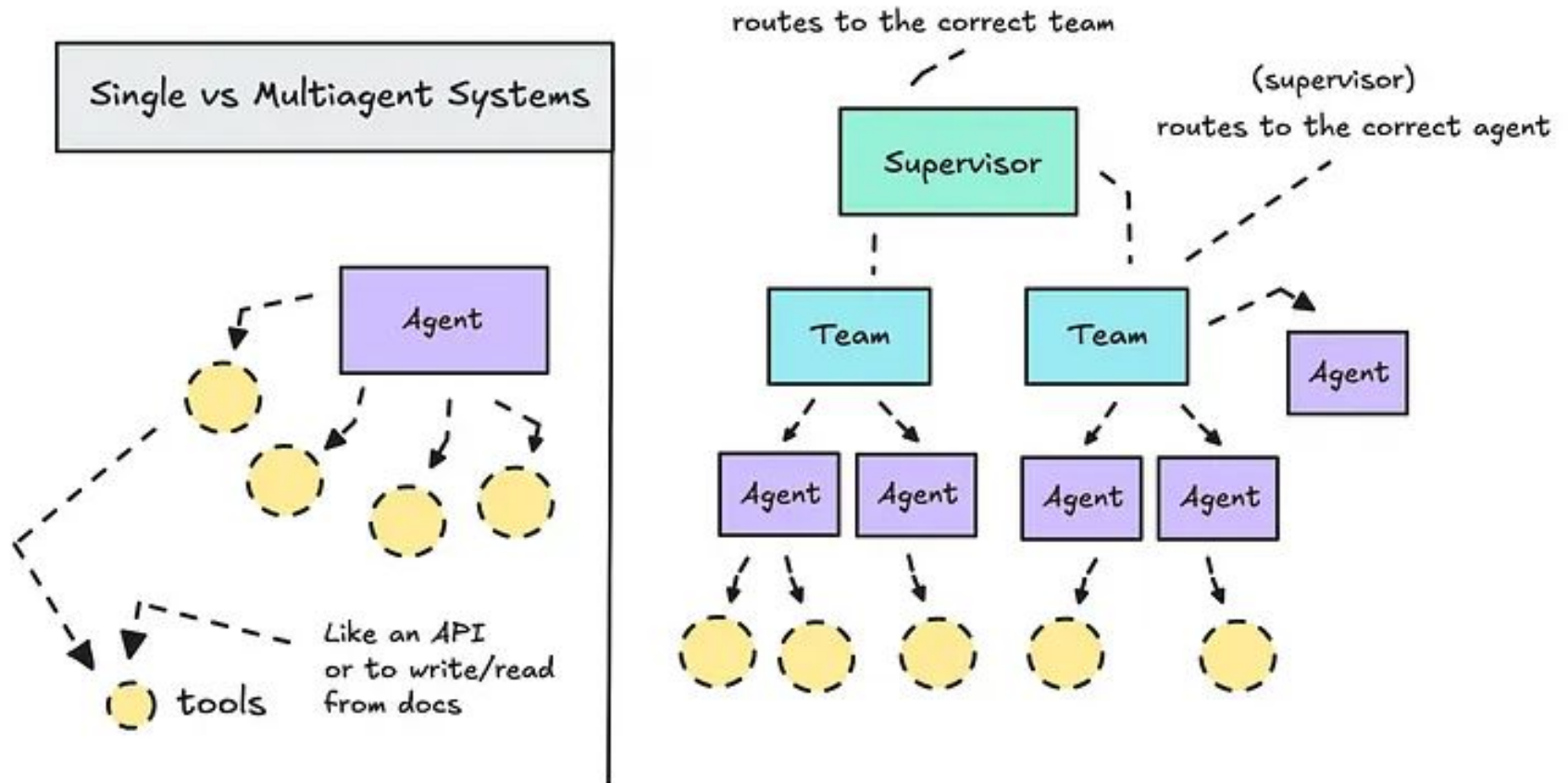
vs.



Single LLM Prompt Translation  
Standard AI Quality

 Multi-Agent Translation  
Best-in-class AI Quality

# Single Agent vs Multi-Agent



# Agent Framework Landscape

- AutoGen – conversation-first
- LangGraph – graph-based workflows
- CrewAI – role orchestration
- Google ADK – Governed enterprise agents
- OpenAI Swarm – Lightweight coordination

# Langgraph

**Best for:** Deterministic, stateful agent workflows

- DAG / state-machine based
- Explicit state persistence
- Strong retry & recovery semantics

 Less “free-form” agent chat

 Excellent for **production reliability**

# CrewAI

**Best for:** Role-based task execution (manager → workers)

- Very intuitive mental model
- Fast to prototype
- Clean “crew / role / task” abstraction



Less control over deep reasoning



Great for PoCs and product teams



# Google ADK

**Best for:** Enterprise-grade, deterministic agent workflows

- Gemini-first, plan-and-tool–driven agents
- Strictly typed tools and explicit execution flows
- Built-in observability, safety, and governance

 Less emergent agent behavior

 Excellent for regulated, production systems

# OpenAI Swarm

**Best for:** Lightweight coordination

- Minimal abstraction
- Simple handoffs between agents
- Very hackable

 No orchestration primitives

 Good for experiments and teaching

# Autogen

**Best for:** Research-grade, tool-heavy, human-in-the-loop systems

- Explicit agent conversations
- Strong tool orchestration
- Event-based logs (what you just saw)
- Supports group chats, selectors, supervisors



Steeper learning curve

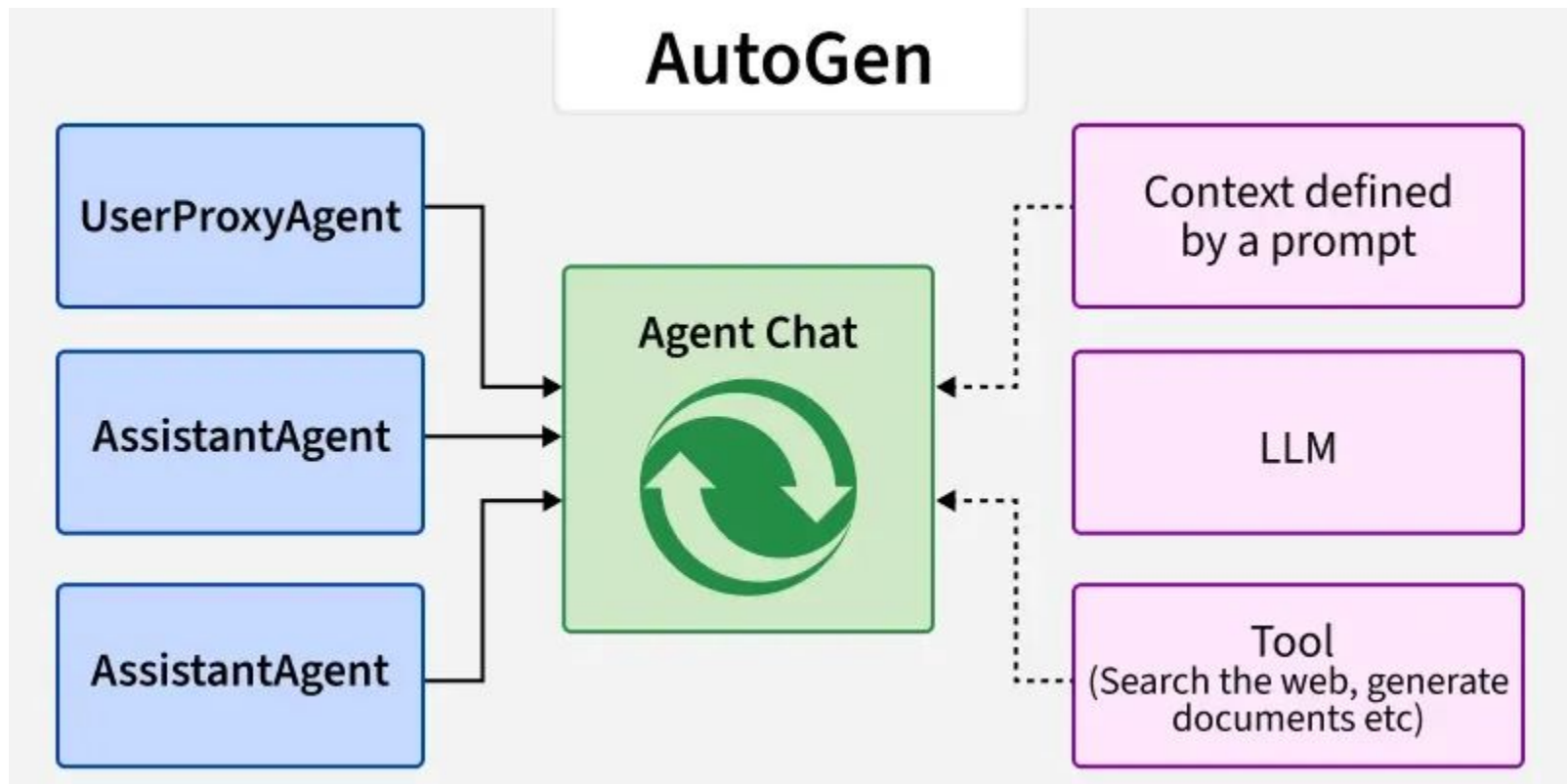


Best for *enterprise* and *complex workflows*

Dimension	AutoGen	LangGraph	CrewAI	Google Agent Development Kit (ADK)	OpenAI Swarm
Primary goal	Emergent multi-agent reasoning	Deterministic orchestration	Fast role-based execution	Governed enterprise agents	Lightweight coordination
Agent interaction	Conversational	Graph transitions	Manager-worker	Plan & tool execution	Handoff-based
Control model	Emergent	State machine / DAG	Semi-structured	Explicit & strict	Minimal
Determinism	Medium	High	Medium	Very high	Low
State persistence	Optional / manual	First-class	Limited	Scoped & governed	None
Tool integration	Flexible	Structured	Flexible	Strongly typed	Minimal

Dimension	AutoGen	LangGraph	CrewAI	Google Agent Development Kit (ADK)	OpenAI Swarm
<b>Observability</b>	Event logs	Graph traces	Basic logs	Enterprise-grade tracing	Minimal
<b>Failure handling</b>	Prompt-driven	Retries & recovery	Limited	Policy-driven	Manual
<b>Human-in-the-loop</b>	Native	Possible	Limited	Controlled	Manual
<b>Best LLM fit</b>	Model-agnostic	Model-agnostic	Model-agnostic	Gemini-first	OpenAI-first
<b>Setup complexity</b>	High	Medium	Low	Medium–High	Very low
<b>Production readiness</b>	Medium–High	High	Medium	Very high	Low
<b>Best use case</b>	Complex agent collaboration	Reliable production workflows	Rapid PoCs & automation	Regulated enterprise systems	Teaching & demos

# AutoGen Core Architecture



# High-Level Architecture

AutoGen's architecture revolves around the concept of *conversational agents*. These agents communicate with each other through messages, forming collaborative workflows. The core components are:

- **Agents:** Represent individual entities with specific roles, responsibilities and capabilities.
- **Teams:** A team is a group of agents that work together to achieve a common goal.
- **Messages:** The communication medium between agents, carrying information and instructions.
- **Conversations:** The dynamic exchange of messages between agents to achieve a common goal.
- **GroupChat:** A multi-agent conversation environment managed by a group chat manager.
- **Tool Calling:** The mechanism by which agents can invoke external functions or APIs.

# Agents

AutoGen AgentChat provides a set of preset Agents, each with variations in how an agent might respond to messages. All agents share the following attributes and methods:

- **name**: The unique name of the agent.
- **description**: The description of the agent in text.
- **run**: The method that runs the agent given a task as a string or a list of messages, and returns a **TaskResult**. Agents are expected to be stateful and this method is expected to be called with new messages, not complete history.
- **run\_stream**: Same as **run()** but returns an iterator of messages that subclass **BaseAgentEvent** or **BaseChatMessage** followed by a **TaskResult** as the last item



# Agent Type

AutoGen supports different types of agents, each with specific characteristics and roles:

- **AssistantAgent:** Designed to assist with specific tasks, leveraging LLMs for reasoning and problem-solving. They typically have pre-defined system prompts to guide their behavior.
- **UserProxyAgent:** Acts as a proxy for a human user, relaying messages and providing feedback. They can execute code and interact with the environment.

# Termination

We explored how to define agents, and organize them into teams that can solve tasks. However, a run can go on forever, and in many cases, we need to know *when* to stop them. This is the role of the termination condition.

AgentChat supports several termination condition by providing a base **TerminationCondition** class and several implementations that inherit from it.

Some examples of Built-In Termination Conditions:

1. **MaxMessageTermination**: Stops after a specified number of messages have been produced, including both agent and task messages.
2. **TextMentionTermination**: Stops when specific text or string is mentioned in a message (e.g., "TERMINATE").

For more refer to this link:

<https://microsoft.github.io/autogen/stable/user-guide/agentchat-user-guide/tutorial/termination.html>

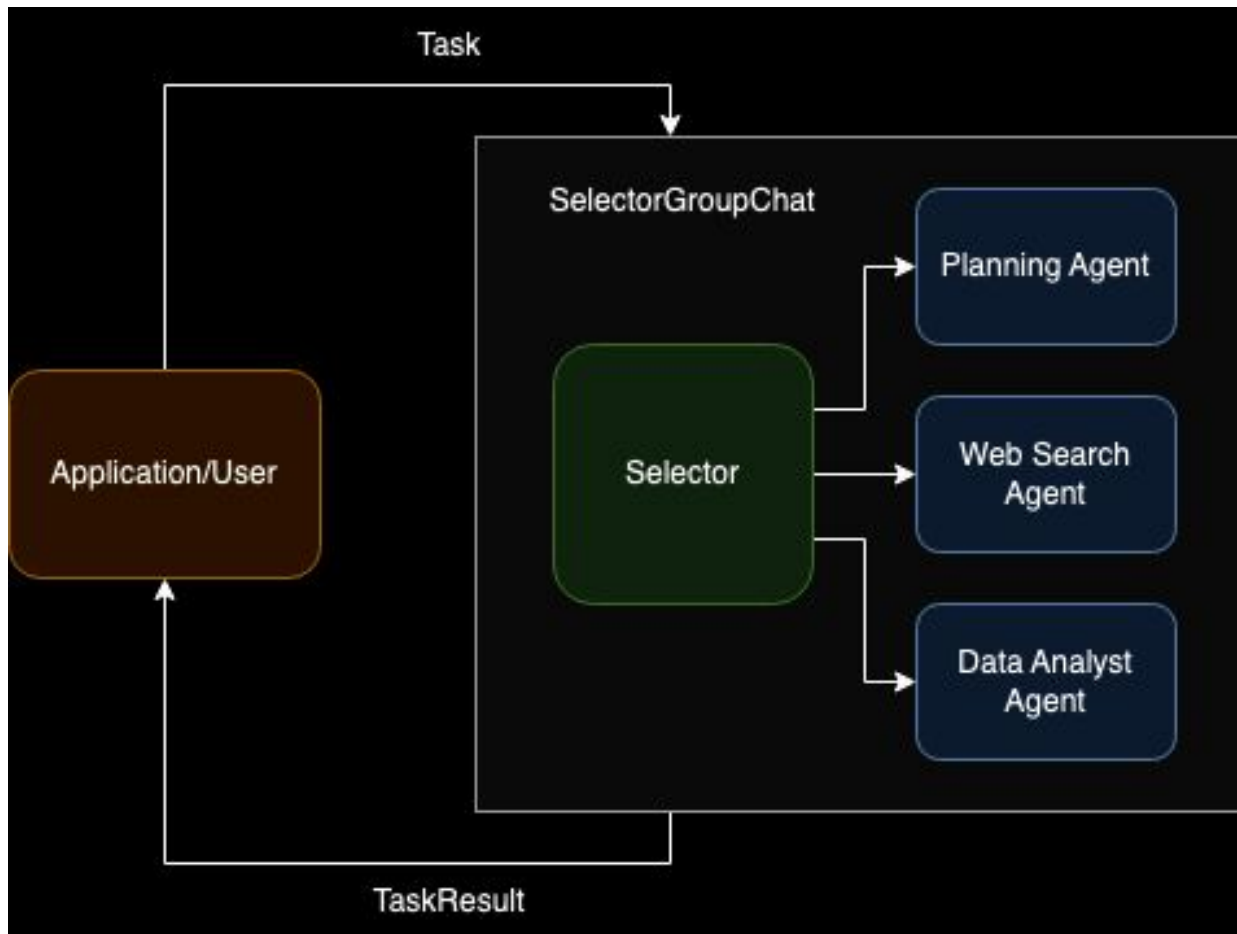
# Teams

A team is a group of agents that work together to achieve a common goal.

AgentChat supports several team presets:

- **RoundRobinGroupChat**: A team that runs a group chat with participants taking turns in a round-robin fashion (covered on this page). [Tutorial](#)
- **SelectorGroupChat**: A team that selects the next speaker using a ChatCompletion model after each message. [Tutorial](#)
- **MagenticOneGroupChat**: A generalist multi-agent system for solving open-ended web and file-based tasks across a variety of domains. [Tutorial](#)
- **Swarm**: A team that uses **HandoffMessage** to signal transitions between agents.

# Selector Group Chat



# Human-in-loop

TBD

# State Management

TBD

# Human + Agent Collaboration



Human

The diagram consists of three blue rectangular boxes arranged horizontally. The first box on the left is labeled 'Human', the middle box is labeled 'GroupChatManager', and the third box on the right is labeled 'Agents'. All boxes have a blue gradient and a thin black border.

GroupChatManager

Agents

# Persistence & Memory

- Stateless vs Stateful agents
- Short-term vs Long-term memory
- Persist decisions, not noise



# Persistence Architecture

Conversation

Memory Layer

Storage

Resume

# Chat UI Integration

- Streamlit
- Clear agent attribution

# End-to-End System

Chat UI

Backend

Multi-Agent System

# Wrap-Up & Next Steps

- Production considerations
- Evaluation & monitoring
- Deployment patterns