

AutoGen

A framework for building AI agents and applications



0.4 version

What is autogen?

a framework for building multi agents
AI system.



Why does it matter?

Autogen lets multiple agents specialize
and work together as a unit

team of superheroes

Autogen v0.4

async foundation

faster and more flexible
for complete tasks.

layered API design

Agent chot API

Quick, task
focused API

Better observability

Core API

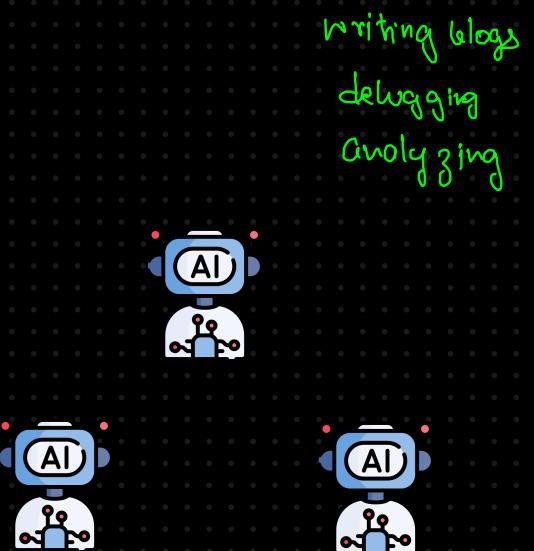
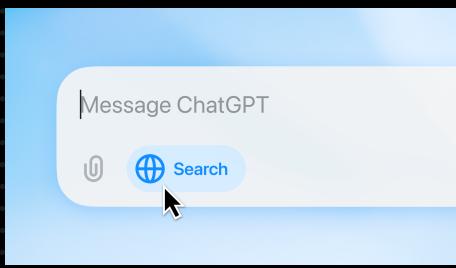
Low level
control

AutoGen

A framework for building AI agents and applications



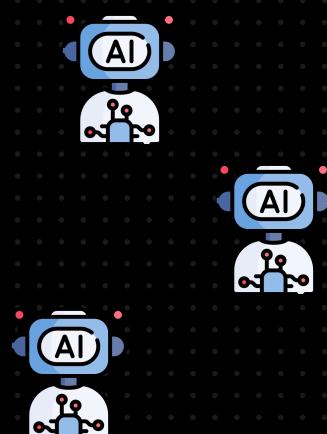
AutoGen is a comprehensive framework developed by Microsoft for building AI agents and applications. It facilitates the creation of both single and multi-agent systems capable of complex interactions and workflows.



Autogen is a framework developed by microsoft to create
multi-agent AI systems.

2023
|
2025 v0.4

Single AI
Lone worker



Why it matters?

Smarter & scalable systems

flexibility ✓

Microsoft AI ecosystem

↳ Azure
↳ Semantic
Kernel

Autogen v0.4

early 2025

1. async foundation

2. Layered

design

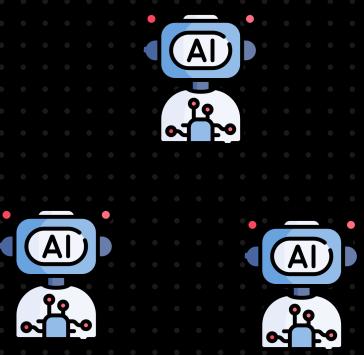


Core

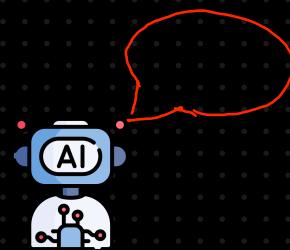
Agent Chat



3. Better Observability



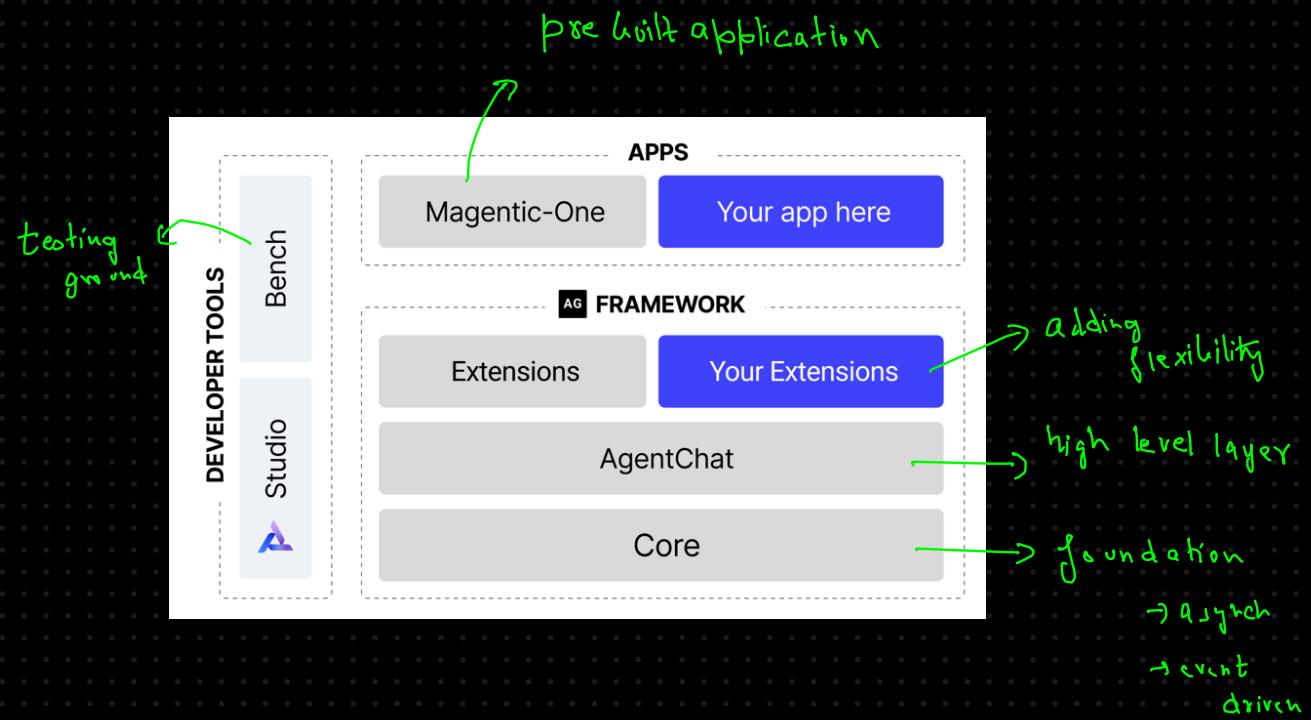
Customer Support



Autogen vs AG2

Setting Up Microsoft Autogen in VS Code

Autogen Architecture



Autogen v0.4 vs Other AI Frameworks

LangChain (rewAI, microsoft Kernel, llmqa index)

Autogen AG

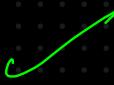
Architecture

event-driven
async

multi agent
capabilities



Observability



Scalability

async foundation handling
complex agent network
efficiently

Use case
suitability

collaborative
AI / team

LangChain



chain-based sequential
execution

better for a single
agent workflow

depth missing

synchron. design struggles
with large scale

data processing /
Q-A tools

Autogen AG

Crew AI



Flexibility

Core

Agent-chat

predefined workflows,
less flexible
for customization

Customization

Event system
& we can do
tweaks

template-based
approach

Maturity

matured

newer &
growing

Learning
Curve

requires bit
more learning

Simplex

Use Case

complex projects

straightforward
team task

Autogen

AG

Semantic Kernel



Focus

multi agent
collaboration

integrate AI into existing
apps

Integration

emphasizes stand-alone
agent network

embedding AI into
systems

Programming Model

event - driven

function - composition
paradigm.

Memory management

context management
is sophisticated

skill based memory

Use Case

multi agents hinge

in-app specific
AI enhancement



Primary Purpose

Agents for task
& conversation

Knowledge retrieval
and document integration

Data Handling

Focused on agent
interaction

Excels with structured
or unstructured data sources

Use Case

Multi-agent or
dynamic agent
teams

Llama index knowledge
based AI

Async Functionality in Python



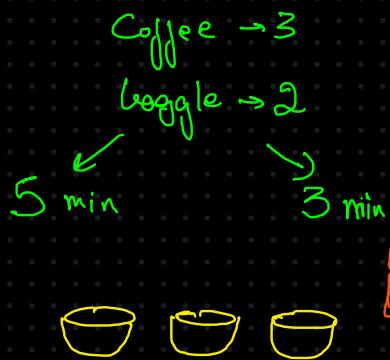
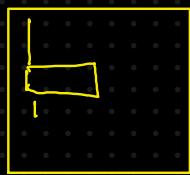
→ API

→ Autogen uses async functionality for real-time agent interactions.

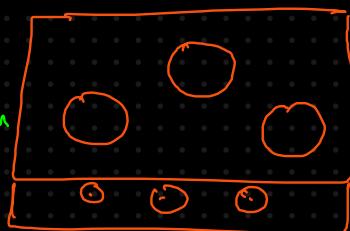
Synchronous
programming



Microwave

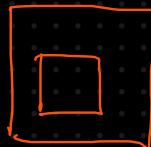


Asynchronous
programming



Subroutine

```
def getData(a,b):
```



```
def main(a,b):
```

getData

main
↓
getData
↓ completed
main

Co-routine

```
async def getData
```



```
def main():
```

getData

Parallelism

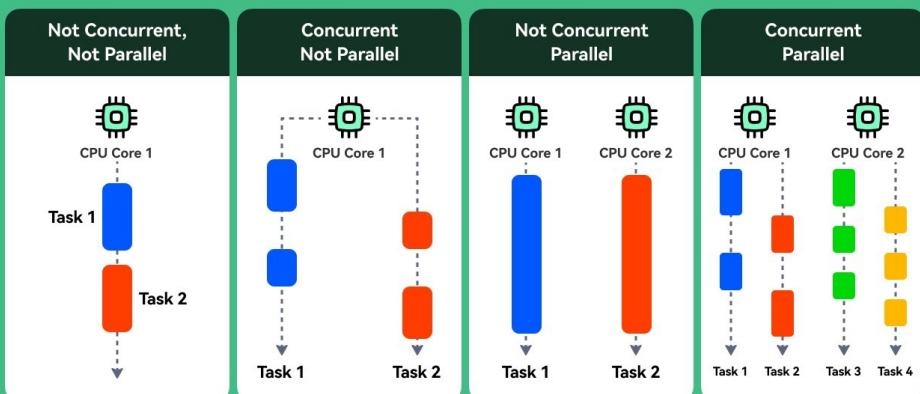
Running multiple tasks simultaneously using multiple threads or processes.



Concurrency

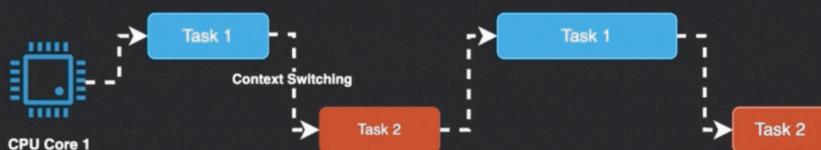
Managing multiple tasks that can start, run & finish with overlapping times.

Concurrency is NOT Parallelism

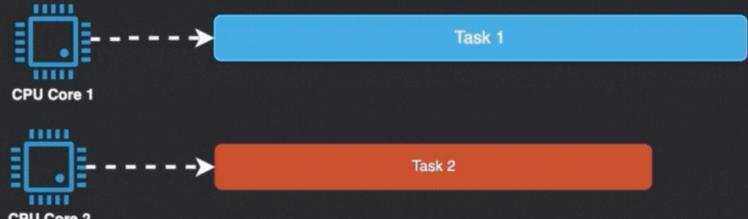


CONCURRENCY VS PARALLELISM

Concurrency

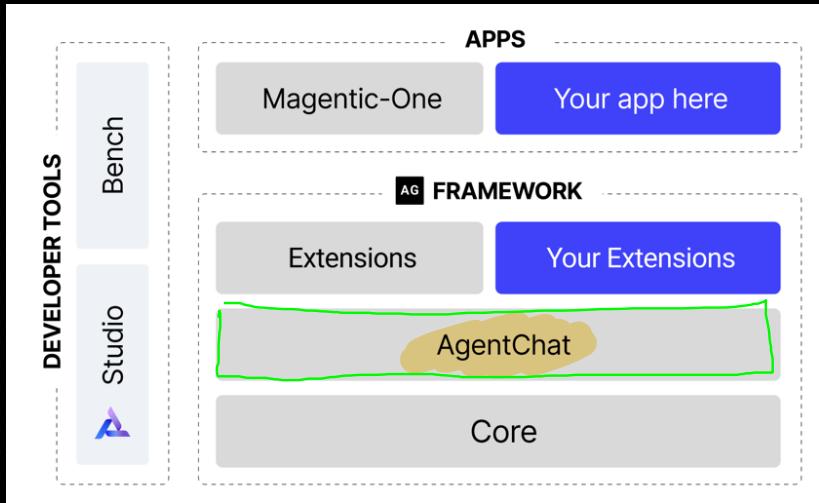


Parallelism



Async Functionality in Python - Practical

Introduction to Agent Chat API



Agent chat provides a modern, async, powered way to handle agent conversation

What it does?

- define agents
- enable them to chat, collaborate and solve problems.
- event driven i.e. response to task as they come

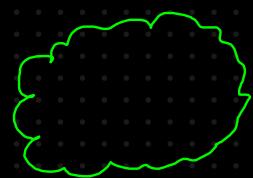
Key Features

- customization with prompts
- connect it to LLMs

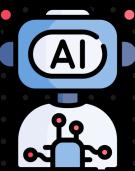
Why it Matters

- makes dev fast and flexible

Configuring LLMs (OpenAI, Gemini, llama)



OpenAI
Gemini
llama



import/install library

```
from autogen_ext.models.openai import OpenAIChatCompletionClient  
  
openai_model_client = OpenAIChatCompletionClient(  
    model="gpt-4o-2024-08-06", get 4o, gpt 3.5-t  
    # api key="sk-...", # Optional if you have an OPENAI_API_KEY  
)
```

class in our lib
which connects
to openAI
models

to authorize & use

Gemini

Google Cloud account → 3 months
\$300 credits

Google AI Studio → Gemini Key

Ollama → local model

