

# ① Deep Agents from Scratch (by LangChain Academy)

- Earlier LLMs were good for generating text and could reason about information presented in their input context
  - single pass operations like RAG or Chat
- Today LLMs can be trained to generate structured output, such as tool calls
  - interact with traditional computing infrastructure
  - build an agent (e.g. ReAct agent that reasons and acts)
    - Operates in a loop (tool call → feedback → repeat if necessary).
- Long-running agents: Claude Code, Deep Research, Manus
  - planning tool
  - file system
  - subagents
  - prompting (detailed)
- Earlier, lots of logic needed to be implemented
- Now, the logic has moved into the prompt and the machinery surrounding the LLM has become simpler
- Deep agent:
  - long-running
  - highly capable } agent
  - <https://github.com/langchain-ai/deepagents>
    - Planning tool
    - File system
    - sub agents
    - system prompt
    - pip install deepagents (or uv pip install deepagents)
    - LangGraph graph (streaming, HTML, memory, studio)
    - tools (req.), instructions (req.), subagents (opt.), model (opt.), builtin\_tools (opt.)
    - Supports: async, MCP, configurations
- Python 3.11+ w/ uv package manager

②

- `uv venv --python 3.11`
  - `uv sync`
  - `uv run jupyter notebook` (localhost:8888)
- Environment variables: LangSmith, OpenAI/Anthropic, Tailly
- LangChain Academy courses: <https://academy.langchain.com/courses/deep-agents-with-langgraph>
- Common deep agent approaches
  - Task planning (e.g. todo), often with recitation
  - Context offloading to file systems
  - Context isolation through sub-agent delegation

## Module 1: Course Overview

- Agents are working on
  - 1) more general tasks
  - 2) over longer time horizons
- Deep agent may do tens of tool calls or turns
- 4 central principles common for deep agents:
  - 1) Planning (todo list, updates), approval for plan, start work
  - 2) Offload context (file system), read/write, ~~long-term memory~~ long-term memory
  - 3) Task delegation (sub-agents), context isolation, careful w/ conflicts, parallelization + summary
  - 4) Careful/extensive prompt engineering, also ok
- Deep agents on GitHub: <https://github.com/langchain-ai/deepagents>

## Module 2: Create Agent

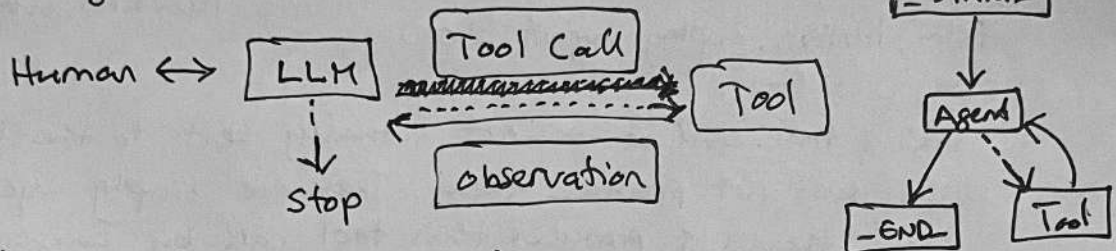
- Create ReAct agent (pre-built)
- `D_create_agent.ipynb`
- ReAct agent is an AI agent that uses the "Reasoning and Acting" (ReAct) framework to combine chain-of-thought (CoT) reasoning with external tool use.
- Three components (this agent): LLM, tools, prompt (context)
- `create_react_agent` → `create_agent`: key capabilities
  - Memory integration (short-term & long-term)
  - Human-in-the-loop control: interrupt execution (human review)
  - Streaming support: agent state, model tokens, tool outputs, combined stream
  - Deployment tooling: LangGraph Platform supports testing, debugging, and deployment
- Studio (UI), LangSmith (observability),



③

- Simple agent: agent - tool loop
- Complex agent: — " — + customized with pre/post hooks and/or structured outputs

• ReAct agent



• @tool decorator for def calculator(...):

• from langgraph.prebuilt import create\_react\_agent

• from langchain.chat\_models import init\_chat\_model

model = init\_chat\_model(model="openai:gpt-4o-mini",  
temperature=0.0)  
tools = [calculator]

agent = create\_react\_agent(model, tools, prompt=SYSTEM\_PROMPT)  
• with\_config({"recursion\_limit": 20})

from IPython.display import Image, display

display(Image(agent.get\_graph(xray=True).draw\_mermaid\_png()))

result1 = agent.invoke({"messages": [{"role": "user",  
"content": "what is 3.1 \* 4.2?"}]})

from utils import format\_messages

format\_messages(result1["messages"])

HumanMessage → AI Message → Tool Message → AI Message  
(tool call) (tool output)

• LangGraph state can be shared between Nodes in the graph

AgentState (default) - remaining\_steps: NotRequired[RemainingSteps] ← recursion limit  
- messages: Annotated[Sequences[BaseMessage], add\_messages]

• Custom state: CalcState(AgentState):

ops: Annotated[List[str], reduce\_list]

- extend calculator's AgentState with "ops"

- injected state: not sent to LLM

def calculator\_wstate(operation, a, b,

state: Annotated[CalcState, InjectedState]

tool\_call\_id: Annotated[str, InjectedToolCallId])

- update state: return Command(update={"ops": ops,  
"messages": [ToolMessage(...)]})

④

```
- agent = create_react_agent(
    model,
    tools,
    prompt = SYSTEM_PROMPT,
    !!! state_schema = CalcState ), with_config(
    { "recursion_limit": 203 }
```

- from IPython.display import JSON  
JSON(result2)

- state & tool-call-id are not actually sent to the LLM.  
LLM does not produce them. They are simply injected  
after the LLM produces this tool call by LangGraph.  
This gives tool node access to them when it's actually  
invoking the tool but doesn't burden the LLM itself  
with having to produce those arguments.

- Command: goto & update simultaneously

• create\_react\_agent(  
model, tools, \*,  
Optional: prompt, response\_format, pre\_model\_hook,  
post\_model\_hook, state\_schema, context\_schema,  
Checkpointter, store, interrupt\_before,  
interrupt\_after, debug, version, name,  
\*\*deprecated\_kwargs) → CompiledStateGraph:

Annotations:  
- "before END structured output" → "structured response"  
- "called before agent node" → "e.g. summarize messages"  
- "call after LLM call but before tool call" → "e.g. HTTP, post-processing"

### Module 3: Todo Lists

- planning component / tool
- Steer long-running, complex tasks → structured todo lists before executing tasks
- Basic pre-built AgentState has messages  
→ add todos (context + status) and files (virtual file system)
- A long-running agent can use a todo list to stay on task  
→ tools: write\_todo, read\_todo
- List contains individual tasks.  
→ update with a full rewrite of the list  
→ allows LLM to reconsider tasks as it makes progress
- from deep-agents\_from\_scratch.prompts import WRITE\_TODOS\_DESCRIPTION  
(system prompt on how to manage todos)  
(create and)
- @tool (description = WRITE\_TODOS\_DESCRIPTION)  
def write\_todos(todos: list[Todo], todo\_call\_id: Annotated[str, injected\_tool\_call\_id]) → Command



⑤

```
return Command(update = {
```

```
  "todos": todos,
```

```
  "messages": ToolMessage(f"Updated todo list to {todos}",
    tool_call_id = tool_call_id)
```

```
@tool(parse_docstring=True)
```

```
def read_todos(state: Annotated[DeepAgentState, InjectedState],
    tool_call_id: Annotated[str, InjectedToolCallId]) -> str:
```

```
    ...
    return result(string)
```

→ ~~read~~ Create ReAct agent will package that string into a tool message and update the messages field in state within the tool observation

- Follow Marcus approach of having a todo recitation after each task from deep-agents-from-scratch. prompts input TODO\_USAGE\_INSTRUCTIONS (system prompt)

```
...
tools = [write_todos, web_search, read_todos]
```

```
prompt = TODO_USAGE_INSTRUCTIONS ... SIMPLE_RESEARCH_INSTRUCTIONS
```

- write\_todos tool call contains todos list which will be stored to the agent's state

- web\_search tool call contains result to a web query

- read\_todos tool call lists all todos in the state

- write\_todos tool call will finally update the todos list by marking pending items as completed

- Finally AI is ready to provide the final result/output to the user

- LangGraph state is a typed data structure

→ available to each node of the graph for its duration  
(also possible to store for long-term storage)

→ data types & reducer function (how info is added to the element)

→ especially useful when a task is mapped to multiple nodes, which are executed in parallel and update state simultaneously

- typing. Annotated allows to attach arbitrary metadata to a type hint  
Annotated[Type, metadata1, metadata2, ...]

- remaining\_steps tracks steps in a graph (recursion\_limit)

- create\_react\_agent

- pre-hook, post-hook

- response\_format

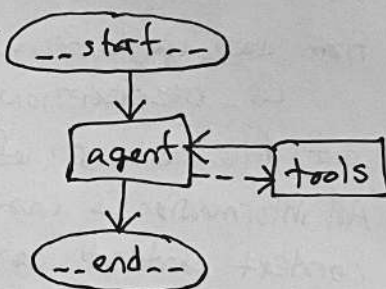
e.g.

summarize (node before)

e.g. HTTP, guardrails, validation, postprocessing

(node after)

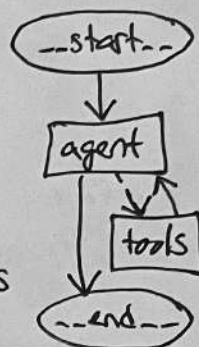
structured\_output (structured response)



- ⑥
- InjectedState provides the tool access to the graph state
  - Command updates values in state + redirect to next node
  - Prompt for using todos (for the agent)

## Module 4: Files

- Agent context window can grow rapidly during complex tasks
  - offload context through file system operations
  - save information to files and fetch it as-needed
- Mock file system: key = file path, value = file content
  - short-term, thread-wise persistence
  - single-agent conversation
  - not suited for information that needs to persist across different conversation threads



- tools: ls, read-file, write-file
- Write data to a file, read some or all from the file later
- from deep-agents-from-scratch.prompts import (
  - LS-DESCRIPTION, READ-FILE-DESCRIPTION, WRITE-FILE-DESCRIPTION)
  - def file\_reducer(left, right):
- All information is easily kept in context, but for long-running agents context content can be compressed or eliminated. Storing information prior to compression and retrieving it when it's needed is smart context engineering.
- <sup>system</sup> prompt for the agent: you have three tools: xxx to manage a virtual file system
- Built on top of a mock file system in LangGraph state

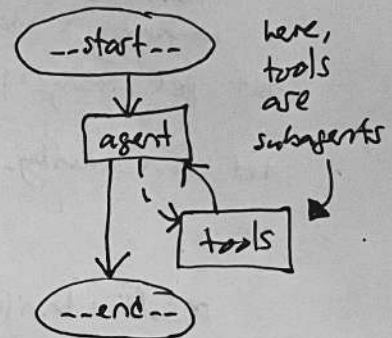
## Module 5: Subagents

- Specialized agents → context isolation
- Issues with long context
  - context clash
  - context confusion
- Context isolation provides an effective solution by delegating tasks to specialized sub-agents, each operating within their own isolated context window.
  - prevents context clashes, confusion, poisoning, and dilution



⑦

- enables focused, specialized task execution
- Subagents with different tool sets tailored for specific tasks
  - registry dictionary with subagent\_type as the key
  - subagent operates in complete isolation from the parent's context → clean separation of concerns
- SubAgent(TypedDict):
  - name: str
  - description: str
  - prompt: str
  - tools: NotRequired[List[str]]
- Task delegation tool for context isolation through sub-agents.
  - ~~can~~ generate sub-agents as tools
  - subagents act as both tools and agents
- Supervisor agent
  - subagent 1
  - subagent 2
  - ...
- supervisor agent's prompt
  - task: coordinate to subagents
  - available\_tools: task, think
  - hard\_limits: iterations, focus
  - scaling\_rules: instructions for use ~~with~~ examples
- Do not give a single agent a lot of tools → gets confused!



## Module 6: Full Agent

- LangGraph > create\_react\_agent → LangChain (v1.0) > create\_agent
- save raw tool observation to reproduce, use summary in tool calls (token heavy)
- Tamily: search engine that is well designed for AI applications
  - 1) search execution using Tamily API
  - 2) context summarization using GPT-4o-mini
  - 3) Result processing using markdownify to convert ~~to~~ HTML to markdown
  - 4) context offloading saves to file
  - 5) strategic thinking provides structured reflection
- search tool that offloads raw contents to files and returns only a summary to the agent

⑧ → solves token efficiency problem by storing detailed search results in files while keeping the agent's working context minimal and focused.

```
• import httpx
from markdownify import markdownify
from pydantic import BaseModel, Field
from taily import TailyClient
from deep-agents-from-scratch.prompts import SUMMARIZE_WEB_SEARCH
from deep-agents-from-scratch.state import DeepAgentState

summarization_model = init_chat_model(model="gpt-4o-mini")
taily_client = TailyClient()
class Summary(BaseModel):
    filename: str = ...
    summary: str = ...
def get_today_str() → str:
    ...
def run_taily_search(search_query: str, max_results: int = 1,
                    topic: Literal["general", "news", "finance"] = "general",
                    include_raw_content: bool = True) → dict:
    result = taily_client.search(search_query, max_results=max_results,
                                include_raw_content=include_raw_content,
                                topic=topic)
    return result
def summarize_webpage_content(webpage_content: str) → Summary:
    structured_model = summarization_model.with_structured_output(Summary)
    summary_obj, filename = structured_model.invoke(
        [HumanMessage(content=SUMMARIZE_WEB_SEARCH.format(
            webpage_content=webpage_content,
            date=get_today_str()))])
    return summary_obj, filename
def process_search_results(results: dict) → list(dict):
    processed_results = []
    HTTPX_CLIENT = httpx.Client()
    for result in results.get('results', []):
        processed_results.append({
            'url': result['url'],
            'title': result['title'],
            'summary': summary_obj.summary,
            'filename': summary_obj.filename,
            'raw_content': raw_content })
    return processed_results
```



9

@tool (parse\_docstring=True)

def tavily\_search(query: str, state: ...) → Command:

# search web and save detailed results to files while  
# returning minimal context. (context offloading)

search\_results = run\_tavily\_search(...)

processed\_results = process\_search\_results(search\_results)

# save each result to a file and prepare summary

# create minimal summary for tool message

return Command(update = {"files": files, "messages": [ToolMessage(summary\_text, tool\_call\_id = tool\_call\_id)]})

FINAL TOOL

@tool (parse\_docstring=True)

def think\_tool(reflection: str) → str:

return f"Reflection recorded: {reflection}"

used after each search to analyze results and plan next steps systematically.

STOP & REFLECT FOR UH

• Tavily search → HTTPX.get(url) → markdownify(response.text)  
→ summarize\_webpage\_content(raw\_content)

• Only summaries are provided to messages, offloading raw content to files  
→ use read\_file() to access full details when needed

• Deep Agent

model = init\_chat\_model(model="anthropic:claudesonnet-4-20250514", temperature=0.0)

max\_concurrent\_research\_units = 3

max\_researcher\_iterations = 3

sub\_agent\_tools = [tavily\_search, think\_tool]

built\_in\_tools = [ls, read\_file, write\_file, write\_todos, read\_todos, think\_tool]

research\_sub\_agent = {"name": "research-agent",

"description": "Delegate research to a subagent researcher. Only give this researcher one topic at a time.",

"prompt": RESEARCHER\_INSTRUCTIONS.format(date=get\_today\_str()),

"tools": ["tavily\_search", "think\_tool"]}

task\_tool = \_create\_task\_tool(sub\_agent\_tools, [research\_sub\_agent], model, DeepAgentState)

delegation\_tools = [task\_tool]

all\_tools = sub\_agent\_tools + built\_in\_tools + delegation\_tools

①

```
SUBAGENT_INSTRUCTIONS = SUBAGENT_USAGE_INSTRUCTIONS.format(  
    max_concurrent_research_units = max_concurrent_research_units,  
    max_researcher_iterations = max_researcher_iterations,  
    date = datetime.now().strftime("%a %b %-d, %Y"))
```

- Research Agent Prompt (instructions)
  - 1) Role
  - 2) task
  - 3) available tools
  - 4) instructions
  - 5) hard limits
  - 6) show your thinking

RESEARCHER\_INSTRUCTIONS

- Supervisor (Deep Agent) (instructions)

- 1) Todo management
- 2) File system usage
- 3) Subagent delegation
  - 3.1) Task
  - 3.2) Available tools
  - 3.3) Hard limits
  - 3.4) Scaling rules

INSTRUCTIONS

- agent = create\_react\_agent(model, all\_tools, prompt = INSTRUCTIONS,  
 state\_schema = DeepAgentState)

```
result = agent.invoke({ "messages": [{ "role": "user",  
    "content": "Give me an overview  
of MCP" } ] })
```

- deepagents package: <https://github.com/langchain-ai/deepagents>

- file system tools (context offloading)
- todo tool (planning)
- task tool

→ only supply subagent and any tools you want to use

```
agent = create_deep_agent(subagents_tools, INSTRUCTIONS,  
    subagents = [research_sub_agent],  
    model = model)  
result = agent.invoke(...)
```

- Deep agents

- planning: prior to task execution
- computer access: shell & file system
- subagent delegation: isolated task execution





## ⑪ · Deep Agents

- 1) Planning Tool
- 2) File System
- 3) subagents
- 4) Prompts (System)

### · Deep Agents UI

- UI for interacting with deep agents

## Module 7: Conclusion