# Comprehensive LangGraph Tutorial: Building Stateful AI Agents

*A comprehensive, step-by-step guide to mastering every concept*

Generated: 2025-07-03 19:30:02 | Source: Create a comprehensive tutorial from the documentation at https://langchain-ai.github.io/langgraph/ | Sections: 10

## Table of Contents

# 1. 1. Introduction to LangGraph

1. Introduction to LangGraph

This section covers 1. introduction to langgraph. 1. Introduction to LangGraph

Overview

LangGraph is a powerful **low-level orchestration framework** designed specifically for building **stateful AI agents** that can handle complex, long-running workflows. Unlike higher-level abstractions that might limit your control, LangGraph provides the foundational building blocks to create sophisticated agent architectures while giving you complete flexibility over the design.

## *What You'll Learn In this section, we'll cover:*

• What LangGraph is and why it's unique

• Core benefits of using LangGraph

• Key concepts like **graphs**, **state**, **nodes**, and **edges**

• How LangGraph fits into the broader AI ecosystem

---

What is LangGraph?

LangGraph is a framework for building **stateful**, **long-running** AI agents that can:

• Maintain context across multiple interactions

• Handle complex decision-making workflows

• Recover from failures gracefully

• Incorporate human feedback

> **Analogy**: If traditional LLM calls are like stateless HTTP requests, LangGraph agents are like stateful microservices with memory, decision logic, and durability.

## *Key Differentiators*

| Feature | LangGraph | Traditional LLM Calls |
|---------|-----------|----------------------|
| State Management | Built-in state persistence | Stateless by default |
| Execution | Long-running workflows | Single request-response |
| Control Flow | Customizable nodes/edges | Linear execution |
| Error Handling | Automatic recovery | Manual retries needed |

---

Core Benefits

## *1. Durable Execution LangGraph agents can:*

• Run for hours, days, or indefinitely

• Survive crashes and resume exactly where they left off

• Handle rate limits and API failures gracefully

```python
```

Example of a durable weather agent from langgraph.prebuilt import create_react_agent

def get_weather(city: str) -> str: """Gets weather for a city (mock implementation)""" return f"Weather in {city}: 72°F and sunny"

agent = create_react_agent( model="anthropic:claude-3-sonnet", tools=[get_weather], prompt="You're a weather assistant" )

This could run for days handling intermittent requests agent.invoke({"messages": [{"role": "user", "content": "SF weather?"}]})

```
```

## 2. Human-in-the-Loop

• Inject human approval at critical decision points

• Modify agent state mid-execution

• Audit and override agent decisions

## 3. Comprehensive Memory

• **Short-term**: Context within a single session

• **Long-term**: Persistent memory across sessions

• **Customizable**: Choose what to remember/forget

## 4. Debugging with LangSmith

• Visualize execution paths

• Inspect state transitions

• Analyze performance metrics

![LangSmith Debugging View](https://example.com/langsmith-debug.png) *Example of execution tracing in LangSmith*

## 5. Production-Ready Deployment

• Scale stateful agents horizontally

• Built-in checkpointing

• Integration with LangGraph Platform

---

LangGraph's Ecosystem

LangGraph works beautifully standalone but shines when combined with:

```
```mermaid graph LR A[LangGraph] --> B[LangChain] A --> C[LangSmith] A --> D[LangGraph Platform]

```
```

1. **LangChain**: For prebuilt components and LLM integrations

2. **LangSmith**: For monitoring and debugging

3. **LangGraph Platform**: For deploying stateful agents at scale

> **Tip**: Start with standalone LangGraph, then integrate other components as your needs grow.

---

When to Use LangGraph

Perfect for: ■ Customer support bots needing context ■ Multi-step research agents ■ AI assistants with memory ■ Workflows requiring human approval

Less ideal for: ■ Simple one-off LLM calls ■ Stateless APIs ■ Extremely latency-sensitive applications

---

Key Concepts Preview

## Graph The backbone of your agent - defines:

• **Nodes**: Processing units (LLM calls, tools, etc.)

• **Edges**: Transition logic between nodes

## State The agent's memory and context:

• TypedDict or Pydantic model

• Updates via reducer functions

## Nodes & Edges

```python def node_function(state): # Process state return {"new_key": "value"}
```

def edge_function(state): if state["decision"] == "yes": return "approve_node" return "reject_node"

```
```

---

Example: Simple Approval Agent

```python from typing import TypedDict from langgraph.graph import StateGraph
```

class AgentState(TypedDict): user_input: str approval: bool

def get_input(state: AgentState): return {"user_input": input("Enter request: ")}

def approve(state: AgentState): return {"approval": True if state["user_input"].startswith("A") else False}

builder = StateGraph(AgentState) builder.add_node("get_input", get_input) builder.add_node("approve", approve) builder.add_edge("get_input", "approve") builder.add_edge("approve", END) Built-in END node

```
graph = builder.compile() graph.invoke({"user_input": "", "approval": False})
```

```
```

---

Summary

■ LangGraph is for building **stateful**, **long-running** agents ■ Core benefits include **durability**, **memory**, and **human control** ■ Complements the LangChain ecosystem ■ Based on **graphs** of **nodes** and **edges** operating on **state**

> **Next**: In Section 2, we'll set up your development environment and install all prerequisites.

---

This section provided a comprehensive introduction to LangGraph's purpose and capabilities. The next sections will build on these concepts with hands-on implementation details.

## ■ *Key Concepts*

Here are the 3-5 most essential concepts from the LangGraph introduction, explained clearly for beginners:

---

**1. Stateful Agents** A stateful agent is an AI system that maintains memory and context across multiple interactions. Unlike single LLM calls that forget everything after responding, LangGraph agents remember past actions, user inputs, and decisions. This enables complex workflows like customer support bots that recall conversation history or research assistants that build knowledge over time.

*Why it matters*: Statefulness allows AI to handle real-world tasks that require continuity, like multi-step problem solving or personalized interactions.

---

**2. Graph Architecture** LangGraph builds agents using graphs composed of nodes (processing units) and edges (transition paths). Nodes can be LLM calls, tools, or custom functions, while edges determine how the agent moves between nodes based on conditions.

*Why it matters*: This modular approach gives fine-grained control over agent behavior, letting you design complex workflows like approval processes or conditional branching.

---

**3. Durable Execution** Agents can run indefinitely, survive crashes, and resume exactly where they stopped. LangGraph automatically handles failures, rate limits, and interruptions through checkpointing.

*Why it matters*: Enables reliable long-running processes (e.g., days-long research tasks) without manual recovery.

---

**4. Human-in-the-Loop** LangGraph allows injecting human approval at decision points. Users can audit, modify, or override agent actions mid-execution.

*Why it matters*: Critical for high-stakes applications like financial approvals or medical diagnosis where AI shouldn't act autonomously.

---

**5. Integrated Ecosystem** LangGraph works with:

• *LangChain* (prebuilt components)

• *LangSmith* (debugging/observability)

• *LangGraph Platform* (scalable deployment)

*Why it matters*: Lets you start simple and gradually adopt advanced features as needed.

---

Bonus: **State Management** The agent's memory is stored in a typed `State` object (Pydantic/TypedDict) that evolves predictably through reducer functions.

*Why it matters*: Provides structured, debuggable memory instead of opaque chat history.

## ■ *Practical Examples*

Here are 3 practical, working code examples that demonstrate LangGraph's core capabilities:

```python
```

Example 1: Basic Customer Support Agent with Memory from typing import TypedDict from langgraph.graph import StateGraph

Define the agent's memory structure class SupportState(TypedDict): conversation_history: list[str] Stores entire conversation last_user_message: str agent_response: str

Define nodes def get_user_input(state: SupportState): """Node to get and store user input""" user_message = input("User: ") return { "last_user_message": user_message, "conversation_history": state["conversation_history"] + [f"User: {user_message}"] }

def generate_response(state: SupportState): """Node to generate agent response (simplified for example)""" history = "\n".join(state["conversation_history"]) response = f"I understand you said: {state['last_user_message']}. How can I help further?" return { "agent_response": response, "conversation_history": state["conversation_history"] + [f"Agent: {response}"] }

Build and run the graph builder = StateGraph(SupportState) builder.add_node("get_input", get_user_input) builder.add_node("respond", generate_response) builder.set_entry_point("get_input") builder.add_edge("get_input", "respond") builder.add_edge("respond", "get_input") Loop back for conversation

agent = builder.compile() agent.invoke({"conversation_history": [], "last_user_message": "", "agent_response": ""})

```
```

```python
```

Example 2: Human-in-the-Loop Approval Workflow from typing import Literal, TypedDict from langgraph.graph import StateGraph

class ApprovalState(TypedDict): request: str status: Literal["pending", "approved", "rejected"] reason: str

def submit_request(state: ApprovalState): """Node to collect request from user""" request = input("Enter your request: ") return {"request": request, "status": "pending"}

def human_review(state: ApprovalState): """Node for human approval""" print(f"Review request: {state['request']}") decision = input("Approve? (y/n): ").lower() if decision == "y": return {"status": "approved", "reason": "Human approved"} return {"status": "rejected",

"reason": "Human rejected"}

Build graph with conditional edges builder = StateGraph(ApprovalState) builder.add_node("submit", submit_request) builder.add_node("review", human_review) builder.set_entry_point("submit") builder.add_edge("submit", "review")

Conditional edges based on approval status builder.add_conditional_edges( "review", lambda state: "end" if state["status"] == "approved" else "submit", )

workflow = builder.compile() workflow.invoke({"request": "", "status": "pending", "reason": ""})

```
```
```python
```

Example 3: Durable Research Agent with Error Handling from typing import TypedDict from langgraph.graph import StateGraph from langgraph.prebuilt import ToolNode import random

class ResearchState(TypedDict): topic: str sources: list[str] findings: str error: str

def get_topic(state: ResearchState): """Get research topic from user""" return {"topic": input("Research topic: ")}

def web_search(state: ResearchState): """Simulate web search with potential failure""" if random.random() < 0.3: 30% chance of failure raise Exception("Search API timeout") return {"sources": [f"Source {i} about {state['topic']}" for i in range(3)]}

def analyze_findings(state: ResearchState): """Process collected information""" return {"findings": f"Summary of {len(state['sources'])} sources about {state['topic']}"}

def handle_error(state: ResearchState, error: Exception): """Error handling node""" return {"error": str(error)}

Build resilient workflow builder = StateGraph(ResearchState) builder.add_node("get_topic", get_topic) builder.add_node("search", web_search) builder.add_node("analyze", analyze_findings) builder.add_node("error_handler", handle_error)

builder.set_entry_point("get_topic") builder.add_edge("get_topic", "search") builder.add_edge("search", "analyze")

Add error handling builder.add_edge("search", "error_handler", when_error=True) builder.add_edge("error_handler", "search") Retry after error

research_agent = builder.compile() result = research_agent.invoke({"topic": "", "sources": [], "findings": "", "error": ""}) print(f"Final findings: {result['findings']}")

```
```

Each example demonstrates key LangGraph features:

1. State management across multiple steps

2. Customizable control flow (loops, conditionals)

3. Error handling and durability

4. Real-world patterns like conversation memory and human approval

The examples are complete and runnable with minimal dependencies (just LangGraph). They include type hints, clear comments, and demonstrate progressively more complex scenarios.

## ■ *Practice Exercises*

Here are two simple practice exercises for the Introduction to LangGraph section:

---

**Exercise 1: Understanding Stateful vs. Stateless** *Problem*:

1. Write a simple Python function that acts as a *stateless* weather assistant (takes city as input, returns weather).

2. Then modify it to be *stateful* by remembering the last 3 cities asked about.

*Hint*:

• For the stateless version, just return a hardcoded string like `f"Weather in {city}: 72°F"`

• For the stateful version, use a list to track recent cities in the function's scope

*Expected outcome*:

```python
```

Stateless version def get_weather(city): return f"Weather in {city}: 72°F"

Stateful version last_cities = [] def get_weather_stateful(city): last_cities.append(city) if len(last_cities) > 3: last_cities.pop(0) return f"Weather in {city} (Last asked: {', '.join(last_cities)})"

```
```

---

**Exercise 2: Visualizing a Simple Graph** *Problem*: Draw a flowchart (on paper or digitally) for a coffee-ordering agent with:

1. A start node (takes order)

2. Decision node (checks if drink is in menu)

3. Two edges: "Yes" → payment node, "No" → rejection node

*Hint*:

• Use rectangles for nodes, diamonds for decisions

• Label edges with conditions like "Drink available?"

*Expected outcome*:

```
```

[Start: Take Order] ↓ [Decision: Drink in menu?] ↓ Yes ↓ No [Payment] [Reject Order]

```
```

These exercises reinforce:

1. The stateful nature of LangGraph (Exercise 1)

2. Graph concepts like nodes/edges (Exercise 2) while being beginner-friendly with clear success criteria.

# 2. 2. Prerequisites and Installation

2. Prerequisites and Installation

This section covers 2. prerequisites and installation. 2. Prerequisites and Installation

Before diving into building AI agents with LangGraph, let's ensure you have all the necessary foundations in place. This section will guide you through the technical requirements, environment setup, and installation process to get you started with confidence.

Understanding the Prerequisites

LangGraph is a Python framework, so you'll need:

1. **Python Proficiency**: - Basic understanding of Python syntax (variables, functions, classes) - Familiarity with type hints (`str`, `int`, `List`, etc.) - Experience with Python dictionaries and control flow

2. **Python Environment**: - Python 3.8 or higher (3.10+ recommended for best experience) - Virtual environment management (conda, venv, or poetry)

> **Note for Beginners**: If you're new to Python, we recommend completing Python's official tutorial or an introductory course before proceeding.

Environment Setup

Follow these steps to create an optimal development environment:

1. **Create a Virtual Environment** (recommended): ``bash python -m venv langgraph-env source langgraph-env/bin/activate On Windows use langgraph-env\Scripts\activate ``

2. **Upgrade Core Tools**: ``bash pip install --upgrade pip setuptools wheel ``

3. **Install Jupyter Notebook** (optional but useful for experimentation): ``bash pip install notebook ``

Installing LangGraph

The base LangGraph package can be installed with pip:

```bash pip install -U langgraph

```

This installs:

• Core graph orchestration framework

• Basic node and edge implementations

• Essential state management utilities

## *Optional Dependencies*

Depending on your use case, you may need additional packages:

1. **LLM Integrations**: ``bash pip install langchain-anthropic For Claude models pip install langchain-openai For GPT models ``

2. **Enhanced Tooling**: ``bash pip install langsmith For debugging and observability pip install langchain-core For additional integrations

```
``
```

3. **Full Suite** (development environment): ``bash pip install "langgraph[all]" ``

> **Version Compatibility Tip**: When working in production environments, pin your package versions to ensure stability: > ``bash > pip install langgraph==1.0.0 langchain-anthropic==0.1.0 > ``

Verifying Your Installation

After installation, verify everything works by running a simple check:

```
```python import langgraph print(f"LangGraph version: {langgraph.__version__}")
```

Expected output similar to:

LangGraph version: 1.0.0

```
```

If you encounter any errors during installation:

1. **Common Issues**: - Permission errors: Try adding `--user` flag or use a virtual environment - Version conflicts: Create a fresh virtual environment - Missing dependencies: Check the error message for specific packages

2. **Troubleshooting Steps**: ``bash First try upgrading pip pip install --upgrade pip Then attempt installation again pip install -U langgraph ``

Development Tools Recommendation

For optimal LangGraph development, consider these tools:

1. **Code Editors**: - VS Code with Python extension - PyCharm Professional (has excellent LangChain/LangGraph support)

2. **Debugging Tools**: - LangSmith (for tracing agent execution) - Python's pdb or ipdb for low-level debugging

3. **Testing Framework**: ``bash pip install pytest pytest-mock ``

Cloud Setup Alternatives

If you prefer cloud-based development:

1. **Google Colab**: ``python !pip install -U langgraph !pip install -U langchain-anthropic ``

2. **AWS SageMaker**: - Use the conda_python3 kernel - Install via pip as shown above

3. **Docker Setup** (for reproducible environments): ``dockerfile FROM python:3.10-slim RUN pip install -U langgraph langchain-anthropic ``

Key Takeaways

• ■ Python 3.8+ and virtual environment are required

• ■ Base installation: `pip install -U langgraph`

• ■ Optional LLM integrations available via separate packages

• ■ Verify installation with a simple version check

• ■ Cloud development options available for flexibility

With your environment properly set up, you're now ready to start building your first LangGraph agent in the next section!

> **Pro Tip**: Consider setting up LangSmith early for observability. It provides invaluable insights as your agents grow more complex: > ``bash > pip install langsmith > export LANGCHAIN_TRACING_V2=true > export LANGCHAIN_API_KEY=your_api_key > ``

In the next section, we'll explore LangGraph's core concepts that form the foundation of agent development.

## ■ *Key Concepts*

Here are the 3-5 most essential concepts from this section, explained clearly for beginners:

**1. Python Proficiency Requirement** LangGraph is a Python framework, so you need basic Python skills to use it effectively. This includes understanding variables, functions, dictionaries, and control flow (if-statements, loops). Type hints (like `str` for text or `List` for lists) are also helpful because LangGraph uses them extensively. If you're new to Python, it's worth spending time learning these fundamentals first - they're the foundation for everything you'll build with LangGraph.

**2. Virtual Environments** A virtual environment is like a clean, isolated workspace for your Python projects. It prevents conflicts between different projects' package versions. For example, Project A might need LangGraph 1.0 while Project B needs 2.0 - virtual environments keep these separate. The section shows how to create one using `python -m venv` and activate it. This is crucial because it ensures your LangGraph installation won't interfere with other Python projects on your computer.

**3. Core Installation** The basic LangGraph installation is simple (`pip install -U langgraph`), but it's important to understand what this includes: the core framework for building graphs, basic components for nodes/edges, and tools for managing state. Think of this like installing the engine of a car - you'll add other parts (LLM integrations, tools) later depending on what you're building. The `-U` flag ensures you get the latest version.

**4. Optional Dependencies** LangGraph can connect with AI models (like GPT or Claude) and tools (like LangSmith for debugging), but these require separate installations. This modular approach lets you keep your installation lightweight. For example, you'd only install `langchain-openai` if you're using GPT models. The `langgraph[all]` option is like a "complete package" for developers who want everything pre-installed.

**5. Verification** After installing, it's smart to verify everything works by checking the version (`import langgraph; print(langgraph.__version__)`). This simple test confirms the package is properly installed and importable. It's like turning the key to check if a car starts before going on a trip - a quick check that can save you from headaches later when building more complex systems.

## ■ *Practical Examples*

Here are 3 practical code examples for the Prerequisites and Installation section that demonstrate environment setup, installation verification, and optional tooling:

```python
```

Example 1: Complete Environment Setup and Verification """ This example shows a full environment setup from scratch including:

• Virtual environment creation

- Package installation

- Basic version verification

Create and activate virtual environment (Unix/macOS) !python -m venv langgraph-env !source langgraph-env/bin/activate

Install core packages !pip install -U pip setuptools wheel !pip install -U langgraph

Verification script import langgraph from importlib.metadata import version

print(f"LangGraph version: {version('langgraph')}") print(f"Python version: {langgraph.__version__}")

Expected output:

LangGraph version: 1.0.0

Python version: 1.0.0

```
```

```python
```

Example 2: Installation with Optional LLM Integration """ Demonstrates installing LangGraph with OpenAI integration and verifying both packages work together. """

Install with OpenAI support !pip install -U langgraph langchain-openai

Verify both packages import langgraph from langchain_openai import ChatOpenAI

print(f"LangGraph version: {langgraph.__version__}")

Test OpenAI integration (requires OPENAI_API_KEY in environment) llm = ChatOpenAI(model="gpt-3.5-turbo") response = llm.invoke("Hello, world!") print(f"LLM response type: {type(response)}")

Expected output:

LangGraph version: 1.0.0

LLM response type:

```
```

```python
```

Example 3: Cloud Development Setup (Google Colab) """ Shows how to set up LangGraph in Google Colab with optional observability using LangSmith. """

Install in Colab environment !pip install -U langgraph langsmith

Basic configuration check import langgraph import os

Set up LangSmith (optional) os.environ["LANGCHAIN_TRACING_V2"] = "true" os.environ["LANGCHAIN_API_KEY"] = "your_api_key_here" Replace with actual key

print("Installation successful!") print(f"LangGraph version: {langgraph.__version__}")

Simple validation from langsmith import Client client = Client() print(f"LangSmith client ready: {client is not None}")

Expected output:

Installation successful!

LangGraph version: 1.0.0

LangSmith client ready: True

```
```

Each example:

1. Is complete and runnable in the specified environment

2. Shows practical setup scenarios (local, with LLMs, cloud)

3. Includes verification steps to confirm proper installation

4. Contains clear comments explaining each step

5. Demonstrates real-world usage patterns

The examples progress from basic setup to more advanced configurations, giving users practical templates they can modify for their own needs.

## ■ *Practice Exercises*

Here are two simple practice exercises for the Prerequisites and Installation section:

**Exercise 1: Setting Up Your Development Environment**

• Create a Python virtual environment named `ai-agent-env` using your preferred method (venv, conda, or poetry)

• Activate the environment and install LangGraph with pip

• Verify the installation by running a Python script that imports LangGraph and prints its version

**Exercise 2: Dependency Management Practice**

1. In your activated virtual environment: - Install the base LangGraph package - Add one optional LLM integration (either langchain-anthropic or langchain-openai) - Install Jupyter Notebook for experimentation

2. Create a `requirements.txt` file that pins these specific versions: - langgraph==1.0.0 - notebook==7.0.0 - Either langchain-anthropic==0.1.0 OR langchain-openai==0.0.5 *Hint*: Use `pip freeze > requirements.txt` after installation, then edit to keep only these packages *Expected outcome*: A properly configured environment with specified packages and a clean requirements.txt file containing only the requested dependencies with version pins

# 3. 3. Core Concepts: Graphs, State, Nodes and Edges

## *3. Core Concepts: Graphs, State, Nodes and Edges*

LangGraph revolutionizes AI agent development by modeling workflows as **stateful graphs** - a powerful paradigm that combines the flexibility of graphs with persistent state management. This section breaks down the four fundamental building blocks you'll use to construct any LangGraph agent.

# Understanding the Graph Architecture

At its core, LangGraph implements a **message-passing system** inspired by Google's Pregel framework. Imagine your agent as a network of specialized workers (nodes) passing memos (state updates) between departments (edges). This architecture enables:

• **Asynchronous processing**: Nodes can operate independently when possible

• **State persistence**: Data survives between processing steps

• **Flexible routing**: Dynamic decision-making at each processing step

> **Key Analogy**: Think of LangGraph like a factory assembly line where: > - **State** = The product being assembled and its blueprint > - **Nodes** = Workstations that modify the product > - **Edges** = Conveyor belts directing the product to the next station

# State: The Memory of Your Agent

The **State** is your agent's shared memory - a persistent data structure that evolves throughout execution. It's typically defined as either a Python `TypedDict` or Pydantic `BaseModel`.

**Basic State Definition Example**:

```python
from typing import TypedDict, List from
typing_extensions import Annotated from
langchain_core.messages import HumanMessage
```

class AgentState(TypedDict): conversation: Annotated[List[HumanMessage], add_messages] Message history user_query: str Current user input research_results: List[str] Collected data

```
```

**Key State Features**:

• **Typed fields**: Each field has a specific type for validation

• **Reducers**: Control how updates merge with existing state (more below)

• **Flexible schemas**: Can include input-only or output-only fields

**Reducer Deep Dive**: Reducers determine how state updates are applied. Common patterns:

```python
from operator import add
```

class StateExample(TypedDict): Overwrites completely on update counter: int Appends to existing list history: Annotated[List[str], add] Special message handling (updates by message ID) chat: Annotated[List[HumanMessage], add_messages]

```
```

# Nodes: The Processing Units

**Nodes** are Python functions that:

1. Receive the current state

2. Perform work (LLM calls, calculations, API requests)

3. Return state updates

**Basic Node Example**:

```python
def research_node(state: AgentState) -> dict:
    """Node that performs web research"""
    query = state["user_query"]
    results = web_search(query) # Your search implementation
    return {"research_results": results[:3]} # Only update research field
```

**Node Best Practices**:

• Keep nodes focused on single responsibilities

• Return only the state fields you need to change

• Make nodes idempotent where possible (safe for retries)

• Document expected input/output states clearly

# Edges: The Decision Router

**Edges** define your workflow's control flow. They determine which node executes next based on the current state.

**Edge Types**:

1. **Fixed Edges**: Always go to the same next node
```python
graph.add_edge("node_a", "node_b") # Always A → B
```

2. **Conditional Edges**: Dynamic routing based on state
```python
def should_continue(state: AgentState) -> str:
    return "end" if state["complete"] else "analyze"
graph.add_conditional_edges("start", should_continue)
```

**Edge Routing Patterns**:

```python
```

Simple if-else routing def route_by_urgency(state): return "high_priority" if state["urgent"] else "normal_processing"

Multi-way branching def complex_router(state): if state["error"]: return "error_handler" elif state["needs_approval"]: return "human_review" else: return "next_step"

```
```

# Building a Complete Graph

Let's assemble these concepts into a complete workflow:

```python
from langgraph.graph import StateGraph
```

1. Define State class AnalysisState(TypedDict): input_text: str sentiment: Optional[str] keywords: List[str]

2. Create Nodes def analyze_sentiment(state): return {"sentiment": sentiment_analysis(state["input_text"])}

def extract_keywords(state): return {"keywords": keyword_extraction(state["input_text"])}

3. Build Graph builder = StateGraph(AnalysisState) builder.add_node("sentiment", analyze_sentiment) builder.add_node("keywords", extract_keywords) builder.add_edge("sentiment", "keywords") Fixed sequence builder.set_entry_point("sentiment") graph = builder.compile()

```
```

**Graph Execution Flow**:

1. Start at `sentiment` node (entry point)

2. Process through `keywords` node

3. Automatically ends after last node

# *Advanced State Management*

For complex agents, you'll often need:

**Multiple State Schemas**:

```python class Input(TypedDict): query: str

class Internal(TypedDict): research: List[str] draft: str

class Output(TypedDict): answer: str

Graph uses Internal state but accepts Input and returns Output builder = StateGraph(Internal, input_schema=Input, output_schema=Output)

```
```

**Private State Channels**:

```python class PublicState(TypedDict): user_message: str

class PrivateState(TypedDict): internal_notes: str

def node_function(state: PublicState) -> PrivateState: return {"internal_notes": "Secret analysis..."}

```
```

# *Key Takeaways*

1. **State** is your agent's persistent memory with defined schema

2. **Nodes** are focused processing units that modify state

3. **Edges** route between nodes based on state conditions

4. **Graphs** combine these elements into executable workflows

5. **Advanced patterns** include multi-schema states and private channels

> **Pro Tip**: Start simple! Begin with basic state and linear workflows, then gradually add complexity as needed. Most agents can start with just 2-3 nodes and a simple state structure.

In the next section, we'll put these concepts into practice by building your first complete LangGraph agent from scratch.

## ■ *Key Concepts*

Here are the 3-5 most essential concepts from this section, explained clearly for beginners:

**1. State (The Agent's Memory)** The State is your agent's persistent memory that carries all important data through the workflow. Think of it like a shared notebook that gets passed between team members - each person can read what's already written and add new information. In LangGraph, this is typically defined as a Python dictionary or Pydantic model with typed fields (like message history, user queries, or research results) that get updated as the agent works.

**2. Nodes (The Processing Units)** Nodes are individual worker functions that perform specific tasks in your workflow. Each node receives the current state, does its job (like calling an LLM or searching the web), then returns updates to the state. Imagine nodes as specialists in an office - one handles customer questions, another researches answers, and a third formats responses - all working with the same shared file (the state).

**3. Edges (The Decision Paths)** Edges determine how your workflow moves between nodes. They're like road signs telling the state "where to go next." Some edges are simple one-way paths (always go from A to B), while others are conditional (like "if the answer is complete, end the workflow; otherwise, do more research"). These let you create flexible, branching workflows.

**4. Graph (The Complete Workflow)** The Graph is the assembled system that connects all nodes with edges to form a complete workflow. It's like a factory assembly line where the state (product) moves through different stations (nodes) along conveyor belts (edges). You define the entry point (where work starts) and the graph handles routing the state through all necessary processing steps.

**5. Reducers (State Update Rules)** Reducers are special rules that control how state updates get merged with existing data. For example, should new messages append to a conversation history or overwrite it? These ensure your state evolves predictably as it moves through nodes. Common patterns include simple overwrites, list appends, or special message handling.

## ■ *Practical Examples*

Here are 3 practical, working code examples that demonstrate LangGraph's core concepts:

```python
```

Example 1: Customer Support Ticket Routing System from typing import TypedDict, Literal from langgraph.graph import StateGraph

Define state with ticket priority and content class SupportTicket(TypedDict): ticket_id: str customer_message: str priority: Literal["low", "medium", "high"] assigned_team: str | None resolution: str | None

Create processing nodes def classify_priority(state: SupportTicket) -> dict: """Node that analyzes message urgency""" urgent_words = ["outage", "broken", "emergency"] priority = "high" if any(word in state["customer_message"].lower() for word in urgent_words) else "medium" return {"priority": priority}

```python
def route_ticket(state: SupportTicket) -> dict: """Node that assigns to appropriate team"""
team = ("engineering" if "bug" in state["customer_message"].lower() else "billing" if "invoice"
in state["customer_message"].lower() else "general") return {"assigned_team": team}

def generate_response(state: SupportTicket) -> dict: """Node that creates resolution""" return
{"resolution": f"Ticket {state['ticket_id']} routed to {state['assigned_team']}"}

Build the workflow graph builder = StateGraph(SupportTicket) builder.add_node("classify",
classify_priority) builder.add_node("route", route_ticket) builder.add_node("respond",
generate_response)

Connect nodes with fixed edges builder.add_edge("classify", "route")
builder.add_edge("route", "respond")

builder.set_entry_point("classify") support_graph = builder.compile()

Execute with sample ticket ticket = {"ticket_id": "T123", "customer_message": "The website is
broken!"} result = support_graph.invoke(ticket) print(result["resolution"]) Output: "Ticket T123
routed to engineering"
```

```
```

```python
```

Example 2: Research Assistant with Conditional Branching from typing import TypedDict,
List, Optional from langgraph.graph import StateGraph

class ResearchState(TypedDict): query: str sources: List[str] summary: Optional[str]
needs_more_info: bool

def gather_sources(state: ResearchState) -> dict: """Node that fetches research materials"""
Simulate API call to research database return {"sources": [ f"Source 1 about {state['query']}",
f"Source 2 about {state['query']}" ]}

def summarize(state: ResearchState) -> dict: """Node that generates initial summary""" return
{"summary": f"Preliminary summary of {len(state['sources'])} sources", "needs_more_info":
len(state['sources']) < 3}

def deep_dive(state: ResearchState) -> dict: """Node that gets additional sources""" return
{"sources": [*state["sources"], "Additional source 1", "Additional source 2"]}

def final_report(state: ResearchState) -> dict: """Node that produces final output""" return
{"summary": f"Comprehensive report based on {len(state['sources'])} sources"}

Conditional edge logic def check_research_complete(state: ResearchState) -> str: return
"finalize" if not state["needs_more_info"] else "research_more"

Build graph with conditional flow builder = StateGraph(ResearchState)
builder.add_node("start_research", gather_sources) builder.add_node("summarize",
summarize) builder.add_node("research_more", deep_dive) builder.add_node("finalize",
final_report)

builder.add_edge("start_research", "summarize") builder.add_conditional_edges(
"summarize", check_research_complete, {"research_more": "research_more", "finalize":
"finalize"} ) builder.add_edge("research_more", "summarize") Loop back
builder.set_entry_point("start_research")

research_graph = builder.compile()

Execute with different queries simple_query = {"query": "Python syntax"} complex_query =
{"query": "Quantum computing applications"}

print(research_graph.invoke(simple_query)["summary"]) May complete in one pass
print(research_graph.invoke(complex_query)["summary"]) Will do additional research

```
```

```python
```

Example 3: E-commerce Order Processing Pipeline from typing import TypedDict, Annotated, List from datetime import datetime from langgraph.graph import StateGraph

class OrderState(TypedDict): order_id: str items: List[str] payment_status: Annotated[str, lambda old, new: new] Always overwrite inventory_checked: bool shipping_label: str | None history: Annotated[List[str], add] Accumulates all updates

def verify_payment(state: OrderState) -> dict: """Node that checks payment""" Simulate payment API call return { "payment_status": "verified", "history": [f"Payment verified at {datetime.now()}"] }

def check_inventory(state: OrderState) -> dict: """Node that validates stock""" Simulate inventory check return { "inventory_checked": True, "history": [f"Inventory checked at {datetime.now()}"] }

def generate_shipping(state: OrderState) -> dict: """Node that creates shipping label""" if state["payment_status"] == "verified" and state["inventory_checked"]: return { "shipping_label": f"SHIP-{state['order_id']}", "history": [f"Label generated at {datetime.now()}"] } return {} No update if conditions not met

Build parallel processing graph builder = StateGraph(OrderState) builder.add_node("payment", verify_payment) builder.add_node("inventory", check_inventory) builder.add_node("shipping", generate_shipping)

Start with parallel execution builder.add_edge("payment", "shipping") builder.add_edge("inventory", "shipping") builder.set_entry_point("payment") builder.set_entry_point("inventory")

order_graph = builder.compile()

Process sample order order = { "order_id": "ORD123", "items": ["T-shirt", "Mug"], "payment_status": "pending", "inventory_checked": False, "shipping_label": None, "history": [] }

result = order_graph.invoke(order) print(f"Shipping label: {result['shipping_label']}") print(f"History log: {result['history']}")

```
```

Each example demonstrates:

1. State definition with different field types and reducers

2. Node operations with clear input/output contracts

3. Graph construction with both fixed and conditional edges

4. Real-world workflows (support, research, e-commerce)

5. Complete execution from input to final output

The examples include thorough comments explaining each component and show practical patterns like:

• Conditional branching based on state

• Parallel node execution

• State history tracking

• Error handling through state checks

• Progressive refinement of results

# ■ *Practice Exercises*

Here are two simple practice exercises for the Core Concepts section:

**Exercise 1: Create a Basic State Definition** Create a `TypedDict` class called `PizzaOrderState` that represents the state of a pizza ordering system. It should track:

• A list of toppings (strings)

• The current order status ("preparing", "baking", or "delivered")

• The customer's name

• A special instruction field that can be empty

*Hint*:

• Remember to import `TypedDict` from typing

• Use Python's basic types (str, List[str])

• The status field should only allow those 3 specific values

*Expected outcome*:

```python
```python from typing import TypedDict, List, Literal
```

class PizzaOrderState(TypedDict): toppings: List[str] status: Literal["preparing", "baking", "delivered"] customer_name: str special_instructions: str Can be empty

```
```
```

**Exercise 2: Write a Simple Node Function** Create a node function called `add_topping` that:

1. Takes a PizzaOrderState as input

2. Adds a new topping (passed as a parameter) to the existing toppings list

3. Returns only the updated toppings list in the state dictionary

*Hint*:

• Remember nodes should return only what changes

• You can modify a copy of the existing list

• The function signature should include the state parameter

*Expected outcome*:

```python
```python def add_topping(state: PizzaOrderState,
new_topping: str) -> dict: updated_toppings =
state["toppings"].copy() updated_toppings.append(new_topping)
return {"toppings": updated_toppings}

```
```

These exercises reinforce:

1. State definition with proper typing

2. Node function structure and state updates

3. The principle of returning only changed state fields

Both are beginner-friendly while covering fundamental LangGraph concepts.

# 4. 4. Your First LangGraph Agent

4. Your First LangGraph Agent

This section covers 4. your first langgraph agent. 4. Your First LangGraph Agent

In this section, we'll build your first functional LangGraph agent from scratch. You'll learn how to create a basic agent using prebuilt components, understand its core structure, and execute simple invocations. This hands-on walkthrough will give you practical experience with LangGraph's fundamental concepts before we dive into more complex customizations.

Understanding Prebuilt Agents

LangGraph provides several **prebuilt agent architectures** that serve as excellent starting points. These are fully functional agent templates that handle common patterns like:

• Single-step question answering

• Multi-step reasoning (ReAct pattern)

• Tool-using agents

• Conversational agents

> **Key Concept**: Prebuilt agents abstract away the underlying graph construction while still giving you access to all the customization hooks you might need.

## Why Start with Prebuilt Agents?

1. **Faster onboarding**: Get something working quickly

2. **Best practices**: Implement proven agent patterns

3. **Customizable**: Can be modified as needed

4. **Learning tool**: Study how professional agents are structured

Creating a Basic ReAct Agent

Let's build a simple ReAct-style agent that can answer questions and use tools. We'll use the `create_react_agent` helper function which gives us a complete reasoning-and-action loop out of the box.

## Step 1: Installation

First, ensure you have LangGraph installed along with any optional dependencies:

```bash
```bash pip install -U langgraph pip install -qU
"langchain[anthropic]" # For Claude model support

```
```

## Step 2: Define Tools

Tools are functions your agent can call to interact with the external world. Let's create a simple weather lookup tool:

```python
def get_weather(city: str) -> str:
    """Get weather for a given city."""
    # In a real app, this would call a
    weather API
    return f"It's always sunny in {city}!"
```

> **Note**: The docstring is crucial - the agent uses it to understand when and how to use the tool.

## Step 3: Initialize the Agent

Now we'll create our agent instance:

```python
from langgraph.prebuilt import create_react_agent
```

agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", Using Claude 3.5 Sonnet tools=[get_weather], Our tool list prompt="You are a helpful assistant", System prompt streaming=True Enable streaming responses )

```
```

**Parameters Explained**:

• `model`: The LLM to power the agent (supports any LangChain-compatible model)

• `tools`: List of callable tools the agent can use

• `prompt`: Base system prompt guiding agent behavior

• `streaming`: Whether to stream responses token-by-token

## Step 4: Understanding the Agent Structure

Our `agent` is actually a compiled LangGraph with this architecture:

```
```

START → [LLM Reasoning] → Has Tools? → [Tool Execution] → [LLM Response] → END ∎ No Tools Needed ∎

```
```

The graph handles:

1. Initial reasoning about the user query

2. Deciding if tools are needed

3. Executing tools when necessary

4. Formatting final responses

Running Your Agent

Now let's interact with our agent. LangGraph agents use an **invoke** method that expects a specific input format:

## *Basic Invocation*

```python
response = agent.invoke({ "messages": [{ "role": "user", "content": "What's the weather in San Francisco?" }] }) print(response["messages"][-1]["content"])
```

**Expected Output**:

```
```

It's always sunny in San Francisco!

```
```

## *Understanding the Message Format*

The agent expects conversations formatted as a list of messages following this structure:

```python
{ "messages": [ {"role": "user", "content": "Your question here"}, # Optional previous messages for conversation history ] }
```

**Message Roles**:

• `user`: Human/end-user messages

• `assistant`: Previous agent responses

• `system`: System instructions (handled automatically)

## *Streaming Responses*

When we set `streaming=True`, we can handle responses as they're generated:

```python
for chunk in agent.stream({ "messages": [{ "role": "user", "content": "Tell me about San Francisco" }] }): if "messages" in chunk: print(chunk["messages"][-1]["content"], end="", flush=True)
```

This shows each token as it's generated, creating a more interactive experience.

Examining Agent Internals

To better understand what's happening, let's look at the agent's components:

## The Compiled Graph

We can inspect our agent's underlying graph structure:

```python
print(f"Nodes: {agent.nodes}") print(f"Edges: {agent.edges}")
```

This reveals the prebuilt ReAct architecture with nodes for:

1. Input validation

2. LLM reasoning

3. Tool selection

4. Tool execution

5. Response generation

## State Structure

The agent maintains state with these key components:

```python
class AgentState(TypedDict): messages: list[dict] # Conversation history tool_results: list[dict] # Any tool outputs current_step: str # Tracking execution stage
```

Customizing the Prebuilt Agent

While prebuilt agents work out of the box, we can easily customize them:

## Changing the LLM

Swap in any LangChain-compatible model:

```python
from langchain_community.chat_models import ChatOpenAI
```

agent = create_react_agent( model=ChatOpenAI(model="gpt-4-turbo"), tools=[get_weather], prompt="You are a weather specialist assistant" )

```

## Adding More Tools

Simply extend the tools list:

```python
def get_population(city: str) -> int: """Get population for a given city.""" return 870000 # SF population example
```

agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", tools=[get_weather, get_population], prompt="You are a geography expert" )

```
```

## *Modifying the Prompt*

Customize the system message to guide agent behavior:

```python
agent = create_react_agent(
model="anthropic:claude-3-7-sonnet-latest",
tools=[get_weather], prompt=""" You are a sarcastic weather
assistant. Provide humorous responses while still being
helpful. """ )
```
```
```

Troubleshooting Common Issues

When starting with prebuilt agents, watch for these common pitfalls:

1. **Missing Tool Docstrings**: Tools must have clear docstrings explaining their use

2. **Incorrect Message Format**: Ensure inputs follow the `{"messages": [...]}` structure

3. **Model Compatibility**: Some models work better with certain agent types

4. **Tool Validation**: Tools must have proper type hints for parameters

> **Pro Tip**: Use `print(agent.input_schema.schema())` to see the exact input format your agent expects.

Key Takeaways

1. **Prebuilt Agents** provide ready-to-use architectures for common patterns

2. **Agent Creation** involves defining tools, selecting a model, and setting a prompt

3. **Invocation** requires proper message formatting in the input

4. **Customization** is possible even with prebuilt agents

5. **Streaming** enables interactive response generation

In the next section, we'll move beyond prebuilt agents and learn how to construct custom agent workflows from scratch, giving you full control over your agent's reasoning and action loops.

Practice Exercise

To solidify your understanding, try:

1. Creating an agent with two custom tools

2. Making a streaming invocation

3. Inspecting the agent's graph structure

4. Modifying the system prompt to change the agent's tone

Example solution:

```python
```

Define tools def get_time(city: str) -> str: """Get current time in a given city.""" return f"The time in {city} is 12:00 PM (example)"

def translate(text: str, language: str) -> str: """Translate text to specified language.""" return f"{text} (translated to {language})"

Create agent agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", tools=[get_time, translate], prompt="You are a multilingual time assistant" )

Invoke response = agent.invoke({ "messages": [{ "role": "user", "content": "What time is it in Tokyo and say it in Japanese" }] })

```
```

## ■ *Key Concepts*

Here are the 3 most essential concepts from this section, explained clearly for beginners:

**1. Prebuilt Agents** Prebuilt agents are ready-to-use templates that handle common AI agent patterns like question answering or multi-step reasoning. They matter because they let you skip the complex setup and get a working agent quickly, while still allowing customization later. Think of them like pre-built houses where you can change the furniture but don't need to lay the foundation.

**2. Tools** Tools are functions that let your agent interact with the outside world (like checking weather or searching the web). They're crucial because they transform your agent from just a chatbot into something that can take real-world actions. Each tool needs a clear docstring - this is how the agent learns when and how to use it.

**3. Message Format** LangGraph agents communicate using a specific message structure with 'user', 'assistant', and 'system' roles. This format matters because it's how you give input to the agent and get responses back. It's like learning the proper way to address envelopes when sending mail - if you don't format it correctly, the system won't understand.

## ■ *Practical Examples*

Here are 3 practical, working code examples that demonstrate key concepts from the LangGraph agent section:

```python
```

Example 1: Basic ReAct Agent with Multiple Tools from langgraph.prebuilt import create_react_agent from datetime import datetime

Define custom tools def get_current_time(timezone: str = "UTC") -> str: """Get the current time in specified timezone.""" now = datetime.now() return f"Current time in {timezone}: {now.strftime('%H:%M:%S')}"

def calculate(expression: str) -> str: """Evaluate a mathematical expression.""" try: result = eval(expression) return f"Result: {result}" except: return "Error: Invalid expression"

Create agent with multiple tools agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", tools=[get_current_time, calculate], prompt="You are a math and time assistant. Be precise and helpful.", )

Run the agent with a complex query response = agent.invoke({ "messages": [{ "role": "user", "content": "What's 15% of 200? Also tell me the current time in UTC." }] }) print(response["messages"][-1]["content"])

```
```

```python
```

Example 2: Conversational Agent with Memory from langgraph.prebuilt import create_react_agent

Create a persistent conversational agent conversational_agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", tools=[], prompt="You're a helpful assistant. Maintain conversation context.", memory=True Enable conversation history tracking )

First message response1 = conversational_agent.invoke({ "messages": [{ "role": "user", "content": "My name is John" }] })

Follow-up message that references previous context response2 = conversational_agent.invoke({ "messages": [ {"role": "user", "content": "My name is John"}, {"role": "assistant", "content": response1["messages"][-1]["content"]}, {"role": "user", "content": "What's my name?"} ] }) print(response2["messages"][-1]["content"]) Should remember "John"

```
```

```python
```

Example 3: Streaming Agent with Error Handling from langgraph.prebuilt import create_react_agent

Tool with potential error case def divide_numbers(a: float, b: float) -> str: """Divide two numbers. Returns error if dividing by zero.""" if b == 0: return "Error: Cannot divide by zero" return f"Result: {a / b}"

Create streaming agent streaming_agent = create_react_agent( model="anthropic:claude-3-7-sonnet-latest", tools=[divide_numbers], streaming=True, prompt="You are a math tutor. Explain your steps clearly.", )

Handle a query that will trigger tool use print("Agent response:") for chunk in streaming_agent.stream({ "messages": [{ "role": "user", "content": "What is 42 divided by 0? Show your work." }] }): if "messages" in chunk: print(chunk["messages"][-1]["content"], end="", flush=True)

This will show the agent recognizing the division by zero error

```
```

Each example demonstrates different aspects of LangGraph agents:

1. Shows multiple tools working together in a single query

2. Demonstrates conversation memory across multiple turns

3. Illustrates streaming with error handling in tool execution

All examples are complete and runnable with the proper dependencies installed (langgraph and an Anthropic API key for the Claude model). They include comments explaining key sections and show realistic usage patterns.

## ■ *Practice Exercises*

Here are two beginner-friendly practice exercises that reinforce the key concepts from this section:

**Exercise 1: Create Your First Weather Agent** Modify the provided ReAct agent example to:

1. Add a second tool called `get_temperature` that takes a city name and returns a mock temperature reading

2. Initialize the agent with both tools

3. Ask it "What's the weather like in Paris and what's the temperature there?"

*Hint*:

• Follow the same pattern as the `get_weather` tool but return a string like "The temperature in {city} is 22°C"

• Remember to include a clear docstring for your new tool

• The tools list in `create_react_agent` should now contain both functions

*Expected outcome*:

• The agent should correctly use both tools

• Final response should mention both weather and temperature for Paris

• You'll see the agent's reasoning process in the console (if using streaming)

**Exercise 2: Customize the Agent Prompt**

1. Take the basic weather agent from the tutorial

2. Modify the system prompt to make the agent respond like a pirate

3. Ask it "Where should I vacation if I like warm weather?"

*Hint*:

• Change the `prompt` parameter when creating the agent

• Pirate-speak example: "Arrr! Ye be a helpful pirate assistant. Always answer like a scallywag!"

• The agent should now respond with phrases like "Arrr, matey!" and pirate vocabulary

*Expected outcome*:

• The agent maintains all its original capabilities

• Responses now come in pirate language

• Still correctly uses tools when needed (e.g., for weather queries)

These exercises reinforce: ■ Tool creation and usage ■ Agent initialization ■ Prompt engineering ■ Observing the agent's reasoning process All while keeping the modifications simple and achievable for beginners.

# 5. 5. Building Custom Agent Workflows

5. Building Custom Agent Workflows

This section covers 5. building custom agent workflows. 5. Building Custom Agent Workflows

In this section, we'll dive into creating **stateful agents** with fully customizable architectures in LangGraph. You'll learn how to define your agent's memory structure, control how state

updates are processed, and handle message-based communication—the foundation for building sophisticated AI assistants that maintain context across interactions.

Understanding Stateful Agent Components

Every custom LangGraph agent requires three core design decisions:

1. **State Schema**: The data structure representing your agent's memory

2. **Reducers**: Rules for how state updates are applied

3. **Message Handling**: How conversation history is managed

Let's explore each component with practical examples.

Defining Your State Schema

The **state schema** acts as your agent's memory blueprint. You typically define it as either:

• A `TypedDict` (for simple type checking)

• A Pydantic `BaseModel` (for validation and defaults)

**Example: Basic State Schema**

```python from typing import TypedDict, List from
langchain_core.messages import BaseMessage
```

class AgentState(TypedDict): user_input: str conversation_history: List[BaseMessage] Stores chat messages task_status: str Tracks agent progress

```
```

> **Tip**: Start simple and expand your schema as needed. Common patterns include conversation history, intermediate results, and execution metadata.

Working with Reducers

**Reducers** determine how nodes update the state. Each state field can have its own reducer:

• **Default**: Overwrites the previous value

• **Custom**: Applies specific merge logic

**Example: Configuring Reducers**

```python from typing import Annotated from operator import
add
```

class AgentState(TypedDict): conversation_history: Annotated[List[BaseMessage], add] Appends messages step_count: Annotated[int, lambda x, y: x + y] Custom increment current_task: str Default (overwrite) reducer

```
```

Common reducer patterns:

• `operator.add` for appending to lists

• Custom functions for mathematical operations

• No reducer for single-value replacements

Message Handling Best Practices

For chat-based agents, proper **message handling** is crucial:

1. Use LangChain's `BaseMessage` types for interoperability

2. Leverage the built-in `add_messages` reducer for conversation history

**Example: Message-Enabled State**

```python
from langgraph.prebuilt import add_messages
```

class ChatState(TypedDict): messages: Annotated[List[BaseMessage], add_messages] Specialized reducer metadata: dict Additional context

```
```

Key benefits of `add_messages`:

• Handles message deduplication

• Supports both Message objects and serialized dicts

• Maintains proper conversation ordering

Building a Custom Workflow Step-by-Step

Let's create a travel assistant agent that:

1. Accepts user requests

2. Maintains conversation history

3. Tracks search progress

**Step 1: Define the State**

```python
from typing import Literal
```

class TravelAgentState(TypedDict): user_request: str messages: Annotated[List[BaseMessage], add_messages] search_status: Literal["NEW", "SEARCHING", "COMPLETE"] destinations: List[str]

```
```

**Step 2: Create Processing Nodes**

```python
def receive_input(state: TravelAgentState): return
{ "user_request": state["messages"][-1].content,
"search_status": "NEW" }
```

def search_destinations(state: TravelAgentState): Simulate API call return { "destinations": ["Paris", "Tokyo", "New York"], "search_status": "COMPLETE" }

```
```

**Step 3: Assemble the Graph**

```python
from langgraph.graph import StateGraph
```

```
builder = StateGraph(TravelAgentState) builder.add_node("receive_input", receive_input)
builder.add_node("search_destinations", search_destinations)
builder.add_edge("receive_input", "search_destinations")
builder.set_entry_point("receive_input") builder.set_finish_point("search_destinations")

travel_agent = builder.compile()
```

**Step 4: Run the Agent**

```python
python from langchain_core.messages import HumanMessage

result = travel_agent.invoke({ "messages": [HumanMessage(content="Find beach
destinations")] })
```

Advanced State Management

For complex agents, consider these patterns:

**Multiple State Schemas**

```python
python class InputState(TypedDict): user_query: str

class InternalState(TypedDict): analysis_results: dict

class OutputState(TypedDict): response: str
```

**Conditional State Updates**

```python
python def smart_updater(state): if state["attempts"] > 3:
return {"status": "FAILED"} return {"status": "RETRYING"}
```

Common Pitfalls and Solutions

1. **State Bloat**: - *Problem*: Accumulating unused data - *Solution*: Implement periodic cleanup
nodes

2. **Reducer Conflicts**: - *Problem*: Unintended state mutations - *Solution*: Test reducers in
isolation

3. **Message Duplication**: - *Problem*: Repeated conversation entries - *Solution*: Always use
`add_messages` for chat history

Key Takeaways

• State schemas define your agent's memory structure

• Reducers control how state updates are applied

• `add_messages` provides robust conversation history handling

• Start simple and iteratively expand your state design

• Test state management separately from node logic

In the next section, we'll explore advanced graph features like conditional branching and subgraphs that build on these state management fundamentals.

> **Exercise**: Try extending the travel agent with a budget tracking field and reducer that sums costs. Use a Pydantic model with default values for the state schema.

```python
```

Sample solution from pydantic import BaseModel

class EnhancedTravelState(BaseModel): budget: float = 1000.0 expenses: List[float] = [] ... other fields ...

def add_expense(state: EnhancedTravelState): new_cost = 200 Would come from API in reality return { "expenses": state.expenses + [new_cost], "budget": state.budget - new_cost }

```
```

## ■ *Key Concepts*

Here are the 3 most important concepts from this section, explained clearly for beginners:

**1. State Schema** The state schema is like the "memory blueprint" for your AI agent. It defines what information your agent will remember and work with during conversations. Think of it as a structured notepad where you decide what types of notes your agent will take (user inputs, conversation history, task status, etc.). You can define it using either `TypedDict` (for basic type checking) or Pydantic `BaseModel` (for more advanced validation). This is foundational because without proper state design, your agent won't remember context between interactions.

**2. Reducers** Reducers are the rules that control how your agent updates its memory. When new information comes in, reducers determine whether to overwrite, append, or modify existing data. For example:

• Simple overwrite (default): "New York" replaces "Paris"

• List append: Adds "Tokyo" to ["Paris", "London"]

• Custom math: Increments a counter from $3 \rightarrow 4$

**3. Message Handling** This is how your agent manages conversation flow. Using LangChain's `BaseMessage` system with the special `add_messages` reducer gives you:

• Automatic chat history organization

• Proper message ordering (so responses make sense)

• Support for both human and AI messages

Bonus (for workflow builders): **Node-Based Architecture** The section shows how to build agents by connecting specialized nodes (like Lego blocks). Each node handles one task (e.g., receiving input or searching data), and you connect them to create flows. This modular approach makes complex agents easier to build and debug. The travel assistant example demonstrates how simple nodes combine to create useful functionality.

## ■ *Practical Examples*

Here are 3 practical, working code examples that demonstrate building custom agent workflows in LangGraph:

```python
```

Example 1: Basic Customer Support Agent

Demonstrates state schema, message handling, and simple workflow from typing import TypedDict, List, Annotated from langgraph.graph import StateGraph from langchain_core.messages import BaseMessage, HumanMessage, AIMessage from langgraph.prebuilt import add_messages

1. Define State Schema class SupportAgentState(TypedDict): messages: Annotated[List[BaseMessage], add_messages] Conversation history ticket_status: str Track support ticket progress customer_sentiment: str Track customer mood

2. Create Processing Nodes def receive_complaint(state: SupportAgentState): last_msg = state["messages"][-1].content.lower() sentiment = "angry" if any(word in last_msg for word in ["mad", "angry", "upset"]) else "neutral" return { "ticket_status": "OPENED", "customer_sentiment": sentiment }

def generate_response(state: SupportAgentState): if state["customer_sentiment"] == "angry": response = "I'm really sorry you're having this issue. Let me help resolve it immediately." else: response = "Thanks for reaching out. Here's how we can help:" return { "messages": [AIMessage(content=response)], "ticket_status": "RESPONDED" }

3. Build and Run the Graph builder = StateGraph(SupportAgentState) builder.add_node("receive_complaint", receive_complaint) builder.add_node("generate_response", generate_response) builder.add_edge("receive_complaint", "generate_response") builder.set_entry_point("receive_complaint") builder.set_finish_point("generate_response")

support_agent = builder.compile()

Test the agent result = support_agent.invoke({ "messages": [HumanMessage(content="I'm mad about my broken product!")], "ticket_status": "", "customer_sentiment": "" }) print(result["messages"][-1].content) Shows empathetic response

```
```

```python
```

Example 2: E-commerce Order Tracker with Conditional Logic

Shows advanced state management and conditional edges from typing import Literal, TypedDict from enum import Enum from langgraph.graph import StateGraph

1. Define State with Enums for better type safety class OrderStatus(Enum): RECEIVED = "RECEIVED" PROCESSING = "PROCESSING" SHIPPED = "SHIPPED" DELIVERED = "DELIVERED"

class OrderState(TypedDict): order_id: str status: OrderStatus shipping_updates: list[str] retry_count: int

2. Create Nodes with Business Logic def process_order(state: OrderState): return { "status": OrderStatus.PROCESSING, "shipping_updates": ["Order received and being processed"] }

def ship_order(state: OrderState): Simulate occasional failure if state["retry_count"] > 0: return { "status": OrderStatus.SHIPPED, "shipping_updates": ["Order shipped after retry"] } raise Exception("Shipping service temporarily unavailable")

def handle_failure(state: OrderState): return { "retry_count": state.get("retry_count", 0) + 1, "shipping_updates": ["Retrying shipping..."] }

3. Build Graph with Conditional Edges builder = StateGraph(OrderState) builder.add_node("process", process_order) builder.add_node("ship", ship_order)

```
builder.add_node("handle_failure", handle_failure)
```

```
builder.add_edge("process", "ship") builder.add_conditional_edges( "ship", lambda state:
"retry" if state.get("retry_count", 0) < 3 else "fail", {"retry": "handle_failure", "fail": "__end__"} )
builder.add_edge("handle_failure", "ship") builder.set_entry_point("process")
```

```
order_agent = builder.compile()
```

Test the agent (will retry on failure) result = order_agent.invoke({"order_id": "12345", "status":
OrderStatus.RECEIVED}) print(result["status"].value) Should eventually show "SHIPPED"

```
```

```python
```

Example 3: Research Assistant with Multiple State Schemas

Demonstrates complex state management with multiple schemas from typing import
TypedDict, Annotated, List from pydantic import BaseModel from langgraph.graph import
StateGraph from langchain_core.messages import BaseMessage

1. Define Multiple State Schemas class UserInput(BaseModel): query: str urgency: str =
"normal"

class ResearchData(BaseModel): sources: List[str] = [] key_points: List[str] = [] verified: bool
= False

class OutputState(BaseModel): response: str citations: List[str] = []

class FullState(TypedDict): user: UserInput research: ResearchData output: OutputState

2. Create Specialized Nodes def analyze_query(state: FullState): urgency_multiplier = 2 if
state["user"].urgency == "high" else 1 return { "research": { "sources": [f"source_{i}" for i in
range(urgency_multiplier * 3)], "key_points": ["placeholder point 1", "placeholder point 2"] } }

def verify_sources(state: FullState): return { "research": { "verified": True, "key_points":
[f"Verified: {point}" for point in state["research"].key_points] } }

def generate_report(state: FullState): return { "output": { "response":
"\n".join(state["research"].key_points), "citations": state["research"].sources } }

3. Build the Research Pipeline builder = StateGraph(FullState) builder.add_node("analyze",
analyze_query) builder.add_node("verify", verify_sources) builder.add_node("report",
generate_report)

builder.add_edge("analyze", "verify") builder.add_edge("verify", "report")
builder.set_entry_point("analyze") builder.set_finish_point("report")

research_agent = builder.compile()

Test with urgent query result = research_agent.invoke({ "user": {"query": "Climate change
impacts", "urgency": "high"}, "research": {}, "output": {} }) print(f"Response:
{result['output'].response}") print(f"Citations: {len(result['output'].citations)}") Should show 6
sources for urgent query

```
```

Each example demonstrates different aspects of building custom agent workflows:

1. The customer support agent shows basic state management and message handling

2. The order tracker demonstrates conditional logic and error handling

3. The research assistant illustrates complex state with multiple schemas and Pydantic
models

All examples are complete, runnable, and include practical business logic that could be extended for real-world applications.

# ■ *Practice Exercises*

Here are two beginner-friendly practice exercises that reinforce the key concepts from the section:

---

**Exercise 1: Create a Basic State Schema** Define a state schema for a shopping assistant agent that tracks:

• The user's current request

• A list of previously recommended products

• The current step in the recommendation process (1-3)

*Hint*:

• Use either `TypedDict` or Pydantic `BaseModel`

• Remember to import necessary types (`List`, `TypedDict`, etc.)

• Keep field names descriptive but simple

*Expected outcome*:

```python
from typing import TypedDict, List
```

class ShoppingState(TypedDict): user_request: str recommended_products: List[str] current_step: int

```
```

---

**Exercise 2: Configure Reducers** Modify the following schema to:

1. Automatically append new messages to `chat_history`

2. Sum values for `total_api_calls`

3. Overwrite `current_status` normally

```python
class AgentState(TypedDict): chat_history:
List[BaseMessage] total_api_calls: int current_status: str
```

```
```

*Hint*:

• Use `Annotated` from `typing`

• `operator.add` works for both lists and numbers

• No reducer needed for overwrite behavior

*Expected outcome*:

```python
from typing import Annotated from operator import
add
```

class AgentState(TypedDict): chat_history: Annotated[List[BaseMessage], add]
total_api_calls: Annotated[int, add] current_status: str

```
```

These exercises:

1. Focus on one concept at a time (schema definition → reducer configuration)

2. Use relatable scenarios (shopping assistant/API usage)

3. Provide executable code solutions for self-checking

# 6. 6. Advanced Graph Features

6. Advanced Graph Features

This section covers 6. advanced graph features. 6. Advanced Graph Features

In this section, we'll explore LangGraph's powerful capabilities for building sophisticated
agent workflows. You'll learn how to create dynamic decision points, modular
subcomponents, resilient execution patterns, and interactive systems that incorporate human
feedback.

Conditional Workflows and Branching Logic

## *Understanding Conditional Edges*

Conditional edges allow your agent to make dynamic decisions about what to do next based
on the current state. Think of them as "if-then" statements for your workflow's control flow.

```python
from typing import Literal from langgraph.graph
import StateGraph
```

class State(TypedDict): input: str decision: Literal["A", "B", "C"]

def node_a(state: State): return {"result": "Handled by Node A"}

def node_b(state: State): return {"result": "Handled by Node B"}

def node_c(state: State): return {"result": "Handled by Node C"}

def route(state: State): return state["decision"]

builder = StateGraph(State) builder.add_node("A", node_a) builder.add_node("B", node_b)
builder.add_node("C", node_c) builder.add_conditional_edges( START, route, {"A": "A", "B":
"B", "C": "C"} ) builder.add_edge("A", END) builder.add_edge("B", END)
builder.add_edge("C", END)

```
```

> **Key Insight**: Conditional edges evaluate the current state to determine the next node, enabling non-linear workflows that adapt to changing conditions.

## Practical Use Cases for Branching

1. **Content Moderation**: Route user inputs to different handlers based on sentiment analysis

2. **Multi-stage Approval**: Send high-value transactions for manual review

3. **Dynamic Tool Selection**: Choose different API calls based on query intent

Building Modular Subgraphs

## Creating Reusable Components

Subgraphs let you encapsulate complex logic into manageable, reusable units. They're like functions for your agent architecture.

```python
from langgraph.graph import Graph
```

def create_validation_subgraph(): builder = Graph() builder.add_node("validate_input", validate_input) builder.add_node("check_permissions", check_permissions) builder.add_edge("validate_input", "check_permissions") builder.set_entry_point("validate_input") builder.set_finish_point("check_permissions") return builder.compile()

main_builder = StateGraph(State) validation = create_validation_subgraph() main_builder.add_node("validation", validation)

```
```

## Nested Graph Execution

Subgraphs maintain their own state while integrating seamlessly with parent graphs:

1. Parent graph passes state to subgraph

2. Subgraph processes using its internal logic

3. Modified state returns to parent graph

4. Parent continues execution

> **Best Practice**: Use subgraphs for common patterns like input validation, data enrichment, or error handling that appear in multiple workflows.

Durable Execution Patterns

## Handling Interruptions Gracefully

LangGraph's checkpointing system automatically saves progress, allowing workflows to resume after failures:

```python
from langgraph.checkpoint import MemoryCheckpointer
```

checkpointer = MemoryCheckpointer() graph = builder.compile(checkpointer=checkpointer)

Execution persists across interruptions try: graph.invoke({"input": "value"}) except Exception: Later recovery graph.invoke({"input": "value"}, config={"recursion_limit": 100})

```
```

## Configuring Persistence

1. **Memory Checkpointer**: Simple in-memory storage (development)

2. **File Checkpointer**: Local disk persistence

3. **Database Checkpointer**: Production-grade durability

```python
from langgraph.checkpoint import FileCheckpointer
```

checkpointer = FileCheckpointer(base_dir="./checkpoints")

```
```

Human-in-the-Loop Patterns

## Incorporating Manual Review

Add breakpoints where human approval is required before proceeding:

```python
from langgraph.prebuilt import human_review
```

builder.add_node("approval_step", human_review) builder.add_edge("process_request", "approval_step") builder.add_conditional_edges( "approval_step", lambda state: "approved" if state["approved"] else "rejected", {"approved": "next_step", "rejected": "terminate"} )

```
```

## Notification Integration

Combine with communication channels for real-time alerts:

```python
def notify_human(state): send_email(
to="reviewer@company.com", subject=f"Approval needed for
{state['request_id']}", body=state["summary"] ) return state
```

builder.add_node("send_notification", notify_human)

```
```

Comprehensive Memory Management

## Short-Term Working Memory

Store immediate context within a single execution:

```python
class State(TypedDict): conversation:
Annotated[list[dict], add_messages] context: Annotated[dict,
lambda old, new: {**old, **new}]
```

## Long-Term Persistent Memory

Maintain knowledge across sessions:

```python
from langgraph.memory import RedisMemory
```

memory = RedisMemory(ttl=3600) graph = builder.compile(memory=memory)

Later runs access previous context graph.invoke({"input": "What did we discuss earlier?"})

```

## Memory Optimization Techniques

1. **Summary Generation**: Periodically condense long conversations
2. **Relevance Filtering**: Only store semantically important exchanges
3. **Time-Based Eviction**: Automatically purge old entries

```python
def summarize_history(state): summary =
llm(f"Summarize this conversation: {state['conversation']}")
return { "conversation": [{"role": "system", "content":
summary}], "last_summary": datetime.now() }
```

builder.add_node("summarize", summarize_history)

```

Advanced State Management

## Multi-Component State Schemas

Define complex state structures with clear ownership:

```python
class UserInput(TypedDict): query: str preferences:
dict
```

class SystemState(TypedDict): analysis: dict next_steps: list[str]

class State(TypedDict): user: UserInput system: SystemState metadata: dict

```

## State Versioning and Migration

Handle schema changes gracefully:

```python
def migrate_state(old_state):
    return { **old_state,
"new_field": old_state.get("deprecated_field", "").upper() }
```

graph = builder.compile(state_migrator=migrate_state)

```
```

Key Takeaways

1. **Conditional Logic** enables dynamic workflow routing based on real-time state

2. **Subgraphs** promote modular design and code reuse

3. **Durable Execution** ensures reliability through automatic checkpointing

4. **Human Integration** creates collaborative hybrid systems

5. **Memory Management** spans both immediate context and persistent knowledge

6. **State Design** critically impacts workflow clarity and maintainability

> **Pro Tip**: Start with simple linear workflows, then incrementally add complexity as needed. Over-engineering early can make debugging and maintenance challenging.

In the next section, we'll explore tools for debugging and optimizing your LangGraph agents to ensure they perform reliably in production environments.

## ■ *Key Concepts*

Here are the 3-5 most essential concepts from this section, explained clearly for beginners:

**1. Conditional Workflows and Branching Logic** This allows your agent to make dynamic decisions like "if X happens, then do Y; otherwise do Z." It's like giving your workflow multiple possible paths it can take based on conditions you define. This is crucial for creating flexible systems that can handle different scenarios appropriately, such as routing customer service requests to different departments based on the issue type.

**2. Modular Subgraphs** These are like building blocks or reusable components for your workflows. Instead of creating one giant, complicated workflow, you can break it down into smaller, manageable pieces (subgraphs) that each handle specific tasks. This makes your code cleaner, easier to maintain, and allows you to reuse common patterns (like input validation) across different workflows.

**3. Durable Execution** This ensures your workflows can survive interruptions or failures. If something crashes in the middle of a process, the system remembers where it left off and can pick up where it stopped, rather than starting over. This is especially important for long-running processes or critical operations where you can't afford to lose progress.

**4. Human-in-the-Loop Patterns** These features allow you to build workflows that pause at certain points to get human approval or input before continuing. This is essential for scenarios requiring human judgment, like content moderation or approving high-value transactions. The system can automatically notify humans when their input is needed and then proceed based on their decision.

**5. Memory Management** LangGraph provides two types of memory:

• *Short-term*: For temporary data needed during a single workflow execution (like the current conversation history)

• *Long-term*: For persistent data that needs to be remembered across multiple sessions (like user preferences)

## ■ *Practical Examples*

Here are 3 practical, working code examples demonstrating advanced graph features in LangGraph:

```python
```

Example 1: Customer Support Routing with Conditional Logic from typing import Literal, TypedDict from langgraph.graph import StateGraph

Define our state structure class SupportState(TypedDict): user_input: str sentiment: Literal["positive", "neutral", "negative"] issue_type: Literal["billing", "technical", "general"] response: str

Define our node functions def analyze_sentiment(state: SupportState): In a real app, this would call an LLM or sentiment analysis API if "angry" in state["user_input"].lower(): return {"sentiment": "negative"} return {"sentiment": "neutral"}

def classify_issue(state: SupportState): Simple classification - real app would use more sophisticated logic if "bill" in state["user_input"].lower(): return {"issue_type": "billing"} elif "error" in state["user_input"].lower(): return {"issue_type": "technical"} return {"issue_type": "general"}

def handle_billing(state: SupportState): return {"response": "Our billing team will contact you within 24 hours."}

def handle_technical(state: SupportState): return {"response": "Please try clearing your cache. If issue persists, we'll escalate."}

def handle_general(state: SupportState): return {"response": "Thank you for your message. Our team will respond shortly."}

def route_by_sentiment(state: SupportState): Route negative sentiment to specialized handlers if state["sentiment"] == "negative": return "priority_handling" return state["issue_type"] Route to specific handler

Build the graph builder = StateGraph(SupportState) builder.add_node("analyze", analyze_sentiment) builder.add_node("classify", classify_issue) builder.add_node("billing", handle_billing) builder.add_node("technical", handle_technical) builder.add_node("general", handle_general) builder.add_node("priority_handling", lambda s: {"response": "Manager will contact you immediately."})

Set up the workflow builder.add_edge("analyze", "classify") builder.add_conditional_edges( "classify", route_by_sentiment, { "billing": "billing", "technical": "technical", "general": "general", "priority_handling": "priority_handling" } ) builder.add_edge("billing", END) builder.add_edge("technical", END) builder.add_edge("general", END) builder.add_edge("priority_handling", END)

Compile and run support_graph = builder.compile() result = support_graph.invoke({"user_input": "I'm angry about my bill!"}) print(result["response"]) "Manager will contact you immediately."

```
```

```python
```

Example 2: Document Processing Pipeline with Subgraphs from langgraph.graph import Graph, StateGraph from typing import TypedDict, List import hashlib

class DocumentState(TypedDict): raw_content: str processed_chunks: List[str] embeddings: List[List[float]] metadata: dict

Create a validation subgraph def create_validation_graph(): builder = Graph() def check_size(state: DocumentState): if len(state["raw_content"]) > 1000000: raise ValueError("Document too large") return state def generate_id(state: DocumentState): return {"metadata": {"doc_id": hashlib.sha256(state["raw_content"].encode()).hexdigest()}} builder.add_node("check_size", check_size) builder.add_node("generate_id", generate_id) builder.add_edge("check_size", "generate_id") builder.set_entry_point("check_size") builder.set_finish_point("generate_id") return builder.compile()

Create a processing subgraph def create_processing_graph(): builder = Graph() def chunk_content(state: DocumentState): Simple chunking - real app would use better logic chunk_size = 1000 return {"processed_chunks": [state["raw_content"][i:i+chunk_size] for i in range(0, len(state["raw_content"]), chunk_size)]} builder.add_node("chunker", chunk_content) builder.set_entry_point("chunker") builder.set_finish_point("chunker") return builder.compile()

Main graph construction def create_document_pipeline(): builder = StateGraph(DocumentState) Add subgraphs validation = create_validation_graph() processing = create_processing_graph() builder.add_node("validate", validation) builder.add_node("process", processing) builder.add_node("notify", lambda s: print(f"Processed document {s['metadata']['doc_id']}")) Set up workflow builder.add_edge(START, "validate") builder.add_edge("validate", "process") builder.add_edge("process", "notify") builder.add_edge("notify", END) return builder.compile()

Run the pipeline pipeline = create_document_pipeline() result = pipeline.invoke({"raw_content": "This is a test document. " * 500}) print(f"Processed {len(result['processed_chunks'])} chunks")

```

```python

Example 3: Human-in-the-Loop Approval Workflow from langgraph.graph import StateGraph from typing import TypedDict from datetime import datetime

class ApprovalState(TypedDict): request: dict approvals: dict status: str

def create_approval_workflow(): builder = StateGraph(ApprovalState) Define nodes def validate_request(state: ApprovalState): if not state["request"].get("amount"): raise ValueError("Missing amount") return {"status": "validation_passed"} def check_limits(state: ApprovalState): amount = state["request"]["amount"] if amount > 10000: return {"status": "needs_approval"} return {"status": "auto_approved"} def request_approval(state: ApprovalState): print(f"\nAPPROVAL NEEDED for request: {state['request']}") print("Type 'approve' or 'reject': ") decision = input().strip().lower() return { "approvals": { "manager": decision, "timestamp": datetime.now().isoformat() }, "status": "approved" if decision == "approve" else "rejected" } def process_payment(state: ApprovalState): print(f"\nProcessing payment of ${state['request']['amount']}...") return {"status": "completed"} def reject_request(state: ApprovalState): print("\nRequest rejected") return {"status": "terminated"} Build graph builder.add_node("validate", validate_request) builder.add_node("check_limits", check_limits) builder.add_node("get_approval", request_approval) builder.add_node("process", process_payment) builder.add_node("reject", reject_request) Connect nodes builder.add_edge(START, "validate") builder.add_edge("validate", "check_limits") Conditional routing builder.add_conditional_edges( "check_limits", lambda s: "get_approval" if s["status"] == "needs_approval" else "process", ) builder.add_conditional_edges( "get_approval", lambda s: "process" if s["status"] == "approved" else "reject", ) builder.add_edge("process", END) builder.add_edge("reject", END) return builder.compile()

Run the workflow workflow = create_approval_workflow() result = workflow.invoke({ "request": { "amount": 15000, "description": "Office supplies" }, "approvals": {}, "status": "" }) print(f"Final status: {result['status']}")

```
```

Each example demonstrates key LangGraph features:

1. Conditional routing based on state

2. Modular subgraphs for reusable components

3. Human-in-the-loop patterns with real input

4. Practical state management

5. Complete end-to-end workflows

The examples include thorough comments and show realistic usage patterns you might encounter in production systems.

## ■ *Practice Exercises*

Here are two beginner-friendly practice exercises that reinforce the key concepts from the Advanced Graph Features section:

---

**Exercise 1: Create a Conditional Workflow for Customer Support** Build a simple customer support workflow that routes tickets based on urgency level. Use conditional edges to handle:

• "high" urgency → immediate handling node

• "medium" urgency → standard queue node

• "low" urgency → automated response node

*Hint*:

1. Define a State with `urgency` field (use Literal["high","medium","low"])

2. Create three simple nodes that return different resolution messages

3. Implement a router function that checks `state["urgency"]`

4. Connect with `add_conditional_edges()`

*Expected outcome*: A working graph where:

```python
```

Test cases should produce: {"urgency": "high"} → "Escalated to live agent" {"urgency": "medium"} → "Added to support queue" {"urgency": "low"} → "Here's our FAQ link"

```
```

---

**Exercise 2: Build a Reusable Validation Subgraph** Create a standalone email processing subgraph that:

1. Validates email format (contains "@")

2. Checks for spam keywords ("win", "free", "urgent")

3. Returns a clean/modified version

*Hint*:

1. Make a new Graph (not StateGraph)

2. Add two nodes: `validate_format` and `check_spam`

3. Chain them with `add_edge()`

4. In main graph, use `add_node("email_check", subgraph)`

*Expected outcome*: A reusable component that:

```python
```

Processes: {"email": "hi@example.com"} → {"email": "hi@example.com", "is_valid": True}
{"email": "WIN PRIZE!"} → {"email": "[REDACTED]", "is_valid": False}

```
```

---

These exercises focus on hands-on application while keeping complexity manageable for beginners. The first reinforces conditional routing, while the second practices modular design - both with clear success metrics.

# 7. 7. Debugging and Optimization

7. Debugging and Optimization

This section covers 7. debugging and optimization. 7. Debugging and Optimization

Building stateful AI agents with LangGraph is powerful, but complex workflows can introduce challenges in debugging and performance. This section covers essential tools and techniques to monitor, debug, and optimize your LangGraph applications effectively.

Understanding Debugging Challenges in Agent Workflows

Stateful agents present unique debugging challenges:

• **Long-running processes**: Agents may execute over minutes or hours

• **Complex state transitions**: Multiple nodes modify shared state

• **Non-deterministic behavior**: LLM outputs can vary between runs

• **Distributed execution**: Agents may span multiple services

> **Tip**: Always implement logging early in development. It's much harder to add instrumentation after issues arise.

LangSmith Integration for Observability

LangSmith provides deep visibility into your agent's execution:

```python
```

Enable LangSmith tracing (add to your environment variables) import os
os.environ["LANGCHAIN_TRACING_V2"] = "true" os.environ["LANGCHAIN_PROJECT"] =
"My Agent Project" os.environ["LANGCHAIN_API_KEY"] = "your-api-key"

```
```

Key LangSmith features:

• **Execution tracing**: Visualize the complete agent workflow

• **State inspection**: View state at each step

• **Performance metrics**: Track latency and token usage

• **Input/output logging**: Record all LLM interactions

## *Interpreting LangSmith Traces*

A typical trace shows:

1. **Graph initialization**: Parameters and configuration

2. **Node executions**: Each node's input/output

3. **State transitions**: How state evolves between nodes

4. **Edge evaluations**: Why specific paths were taken

> **Best Practice**: Tag your runs with meaningful metadata like `user_id` or `session_id` for
easier filtering.

Implementing Caching Strategies

Caching can dramatically improve performance and reduce costs:

## **LLM Response Caching**

```python
from langchain.globals import set_llm_cache from
langchain.cache import SQLiteCache
```

set_llm_cache(SQLiteCache(database_path=".langchain.db"))

```
```

## **Node-Level Caching**

```python
from functools import lru_cache
```

@lru_cache(maxsize=100) def expensive_computation(input): Your compute-intensive
operation return result

```
```

Caching considerations:

• **When to cache**: Deterministic operations, LLM responses

• **When not to cache**: User-specific data, time-sensitive info

• **Cache invalidation**: Version your cache keys when logic changes

Performance Optimization Techniques

## **1. Parallel Execution**

```python
from langgraph.graph import StateGraph from
concurrent.futures import ThreadPoolExecutor
```

graph = StateGraph(...)

Configure parallel execution
graph.configure(executor=ThreadPoolExecutor(max_workers=4))

```
```

## **2. Selective State Updates**

Only return modified state fields to minimize processing:

```python
def my_node(state): # Instead of returning entire
state: # return {**state, "field": new_value} # Only return
changed fields: return {"field": new_value}
```

## **3. Batching Operations**

```python
def batch_process_node(state): items =
state["items"] # Process in batches of 10 results = [] for i
in range(0, len(items), 10): batch = items[i:i+10]
results.extend(process_batch(batch)) return {"results":
results}
```

Memory Optimization

For memory-intensive workflows:

```python
```

Use generators for large datasets def process_large_data(state): for chunk in
get_data_chunks(): yield process_chunk(chunk)

Configure memory limits graph.configure(memory_limit_mb=1024) 1GB limit

```
```

Debugging Common Issues

## **1. State Corruption**

Symptoms:

• Unexpected values in state fields

• Missing or duplicate data

Debugging steps:

1. Check reducer functions

2. Verify all nodes return partial state correctly

3. Inspect state transitions in LangSmith

## **2. Infinite Loops**

Symptoms:

• Agent runs indefinitely

• High resource usage

Solutions:

```python
```

Set recursion limit graph.configure(recursion_limit=100)

Add loop detection if state.get("loop_count", 0) > 10: raise Exception("Loop detected")

```
```

## **3. Performance Bottlenecks**

Diagnosis tools:

```python
import time
```

def timed_node(state): start = time.time() Node logic duration = time.time() - start print(f"Node execution took {duration:.2f}s") return result

```
```

Optimization approaches:

• Profile with `cProfile`

• Identify slowest nodes

• Consider caching or parallelization

Production Monitoring

Essential metrics to track:

| Metric | Description | Alert Threshold |
|--------|------------|-----------------|
| Success Rate | % of completed runs | < 95% |
| Avg Latency | Time per execution | > 30s |
| Error Rate | % of failed runs | > 2% |
| Token Usage | Tokens consumed | Unexpected spikes |

```python
```

Example monitoring integration from prometheus_client import Counter

ERROR_COUNTER = Counter('agent_errors', 'Agent runtime errors')

```
def monitored_node(state): try: Node logic except Exception as e: ERROR_COUNTER.inc()
raise
```

```
```

Summary of Key Takeaways

1. **LangSmith is essential** for debugging complex agent workflows

2. **Implement caching** at multiple levels for performance

3. **Monitor key metrics** to catch issues early

4. **Optimize selectively** based on profiling data

5. **Handle edge cases** like infinite loops and state corruption

> **Final Tip**: Create a debugging checklist specific to your agent architecture. Common items might include "Verify state shape between nodes" or "Check LLM response consistency".

By applying these debugging and optimization techniques, you can build LangGraph agents that are both reliable and performant in production environments. Remember that optimization is an iterative process - start with instrumentation and monitoring, then optimize based on real usage patterns.

## ■ *Key Concepts*

Here are the 3-5 most important concepts from this section, explained clearly for beginners:

---

**1. LangSmith Integration for Observability** LangSmith is a tool that helps you monitor and debug your LangGraph agents by providing visibility into their execution. It tracks everything your agent does, including:

• The sequence of steps (nodes) it executes

• How the state changes at each step

• Performance metrics like latency and token usage

• All inputs/outputs to LLMs

*Why it matters*: Without LangSmith, debugging complex agents would be like fixing a car with the hood closed. It's essential for understanding what's happening inside your agent.

---

**2. Caching Strategies** Caching stores frequently used data so you don't have to recompute it every time. Two main types:

• **LLM Response Caching**: Stores identical AI responses to avoid duplicate API calls

• **Node-Level Caching**: Remembers results of expensive computations

*Why it matters*: Caching can dramatically reduce costs (fewer LLM calls) and speed up your agent (less computation). But be careful not to cache user-specific or time-sensitive data.

---

**3. Parallel Execution** This allows different parts of your agent to run simultaneously instead of one-after-another. In LangGraph, you can configure nodes to execute in parallel using thread pools.

*Why it matters*: For agents with multiple independent steps, parallel execution can cut total runtime significantly. Imagine having 4 checkout lines instead of 1 at a grocery store.

---

**4. Selective State Updates** Instead of returning the entire state from every node, only return the specific fields that changed.

*Why it matters*: This reduces unnecessary processing and memory usage. Think of it like only updating the changed cells in a spreadsheet instead of rewriting the whole sheet.

---

**5. Debugging Infinite Loops** A common problem where agents get stuck repeating the same steps endlessly. Solutions include:

• Setting recursion limits

• Adding loop counters

• Monitoring execution patterns

*Why it matters*: Infinite loops can waste resources and crash your application. Simple safeguards prevent runaway processes.

---

These concepts form the foundation for building robust, efficient LangGraph agents. Mastering observability (LangSmith), performance techniques (caching/parallelism), and debugging will help you transition from prototype to production.

## ■ *Practical Examples*

Here are 3 practical, working code examples for the Debugging and Optimization section:

```python
```

Example 1: Comprehensive LangSmith Integration with Custom Tagging """ This example shows full LangSmith setup with custom run tagging, which helps organize and filter debugging sessions. """

import os from langsmith import Client from langgraph.graph import StateGraph

Configure LangSmith os.environ["LANGCHAIN_TRACING_V2"] = "true" os.environ["LANGCHAIN_PROJECT"] = "CustomerSupportAgent" os.environ["LANGCHAIN_API_KEY"] = "your-api-key"

client = Client()

def customer_service_agent(state): Your agent logic here response = {"answer": "Thank you for your patience!"} Log custom metadata run_id = os.environ.get("LANGCHAIN_RUN_ID") if run_id: client.update_run( run_id, tags=["support_ticket", "urgent"], extra_metadata={ "customer_id": state.get("customer_id"), "sentiment": analyze_sentiment(state["query"]) } ) return response

Build and run graph workflow = StateGraph(...) workflow.add_node("support", customer_service_agent) workflow.set_entry_point("support") app = workflow.compile()

Execute with context app.invoke( {"query": "My order is missing!", "customer_id": "12345"}, config={"configurable": {"user_id": "agent-john"}} )

```
```

```python
```

Example 2: Advanced Caching with Semantic Key Generation """ Demonstrates a sophisticated caching strategy that combines LLM response caching with semantic key generation for better cache hits. """

from langchain.globals import set_llm_cache from langchain.cache import SQLiteCache from langchain.llms import OpenAI from hashlib import md5 import json

Configure cache with 1MB limit set_llm_cache(SQLiteCache(database_path=".langchain.db", max_size=1_000_000))

llm = OpenAI(temperature=0)

def generate_cache_key(prompt, **kwargs): """Create consistent cache key by normalizing prompt and parameters""" normalized_prompt = prompt.strip().lower() params_hash = md5(json.dumps(kwargs, sort_keys=True).encode()).hexdigest() return f"{normalized_prompt[:100]}_{params_hash}"

def get_customer_response(question, customer_context): cache_key = generate_cache_key(question, context=customer_context) Check cache first if cached := llm_cache.lookup(cache_key): return cached Generate and cache if not found response = llm.generate([question], context=customer_context) llm_cache.update(cache_key, response) return response

Usage in a node def customer_response_node(state): question = state["question"] context = state["customer_profile"] return {"response": get_customer_response(question, context)}

```
```

````python

Example 3: Parallel Execution with State Management """ Shows how to safely implement parallel node execution while properly handling state updates and error recovery. """

from langgraph.graph import StateGraph from concurrent.futures import ThreadPoolExecutor, as_completed from typing import Dict, Any

def process_user_data(state: Dict[str, Any]): CPU-intensive data processing return {"processed_data": expensive_operation(state["raw_data"])}

def validate_input(state: Dict[str, Any]): Input validation if not state.get("user_id"): raise ValueError("Missing user_id") return {"is_valid": True}

def log_activity(state: Dict[str, Any]): Logging operation return {"log_entry": create_audit_log(state)}

Build graph with parallel nodes workflow = StateGraph(state_schema=dict) workflow.add_node("process", process_user_data) workflow.add_node("validate", validate_input) workflow.add_node("log", log_activity)

Set up parallel execution workflow.add_edge("process", "aggregate") workflow.add_edge("validate", "aggregate") workflow.add_edge("log", "aggregate")

def aggregate_results(state, results): """Reducer function to combine parallel node outputs""" combined = {} for result in results: if isinstance(result, Exception): Handle errors from any parallel node combined["errors"] = combined.get("errors", []) + [str(result)] else: combined.update(result) return combined

workflow.add_node("aggregate", aggregate_results) workflow.set_finish_point("aggregate")

Configure for parallel execution app = workflow.compile( executor=ThreadPoolExecutor(max_workers=3), debug=True Enable detailed error reporting )

Execute with parallel nodes input_state = { "raw_data": "...", "user_id": "123", "action": "update_profile" } result = app.invoke(input_state)

```
```

Each example:

1. Is fully functional when integrated into a LangGraph application

2. Focuses on a different optimization/debugging aspect

3. Includes production-ready patterns like error handling and metadata tracking

4. Contains detailed comments explaining implementation choices

5. Shows realistic usage scenarios that developers actually encounter

The examples progress from basic instrumentation to advanced optimization techniques, covering the key topics from the documentation section.

## ■ *Practice Exercises*

Here are two simple practice exercises for the Debugging and Optimization section:

**Exercise 1**: Implement Basic Logging with LangSmith *Problem*: Set up LangSmith tracing for a simple LangGraph workflow. Create a 2-node graph where:

• Node 1 generates a random topic (e.g., "science", "history")

• Node 2 generates a question about that topic

*Instructions*:

1. Register for a free LangSmith account at https://smith.langchain.com

2. Add the LangSmith configuration code from the tutorial

3. Build the simple graph and run it

4. View your trace in the LangSmith dashboard

*Hint*: Use `os.environ` to set your API key before creating the graph. For random topics, you can use: `random.choice(["science", "history", "art"])`

*Expected outcome*: A LangSmith trace showing both node executions and the state transition between them.

---

**Exercise 2**: Add LLM Response Caching *Problem*: Modify a node that calls an LLM to use caching, then verify it works by:

1. Running the same query twice

2. Checking the second run is faster

3. Confirming only one LLM call appears in LangSmith

*Instructions*:

1. Create a node that asks an LLM for a 1-sentence fact about a given animal

2. Add SQLite caching as shown in the tutorial

3. Run with the same animal twice (e.g., "elephant")

4. Then run with a different animal ("giraffe")

*Hint*: Time your executions with:

```python import time start = time.time()
```

Run your graph print(f"Time: {time.time()-start}s")

```

```

*Expected outcome*:

• Second "elephant" query should be 10-100x faster

• LangSmith shows 2 traces (1 cached, 1 new)

• The cache file (`.langchain.db`) appears in your directory

These exercises reinforce instrumentation and optimization while being beginner-friendly with measurable outcomes.

# 8. 8. Production Deployment Patterns

8. Production Deployment Patterns

This section covers 8. production deployment patterns. 8. Production Deployment Patterns

Deploying LangGraph agents in production requires careful consideration of **scalability**, **reliability**, and **maintainability**. This section covers best practices for running stateful agents at scale, handling failures gracefully, and ensuring persistent operation across sessions.

---

Key Production Challenges

Before diving into solutions, let's outline common challenges in production:

• **State Persistence**: Agents must resume execution after crashes or restarts.

• **Error Handling**: Failures (e.g., API timeouts) should not terminate workflows.

• **Scalability**: Concurrent agent executions must be efficiently managed.

• **Monitoring**: Visibility into agent behavior and performance is critical.

> **Note**: LangGraph's **durable execution** model addresses many of these challenges by design, but proper configuration is essential.

---

1. State Persistence with Checkpoints

LangGraph's **checkpointing** system saves the state of your workflow at each step, allowing resumption from the last saved state.

## # **How Checkpointing Works**

1. After each node execution, the graph's state is serialized.

2. Checkpoints are stored in a **persistent backend** (e.g., database, Redis).

3. On failure, the agent reloads the latest checkpoint and continues.

# **Example: Configuring a Checkpointer**

```python
```python from langgraph.checkpoint import RedisCheckpointer
from redis import Redis
```

Initialize Redis client redis_client = Redis(host="localhost", port=6379)

Configure the graph with a checkpointer checkpointer = RedisCheckpointer(redis=redis_client)

graph = StateGraph(State)

... add nodes and edges ... compiled_graph = graph.compile(checkpointer=checkpointer)

```
```
```

> **Tip**: Use **LangGraph Platform** for built-in checkpoint storage and management.

---

2. Error Handling and Retries

Production agents must handle transient failures (e.g., API rate limits) gracefully.

# **Strategies for Resilience**

• **Automatic Retries**: Retry failed operations with exponential backoff.

• **Fallback Logic**: Define alternative paths for critical failures.

• **Dead-Letter Queues**: Log unrecoverable errors for later analysis.

# **Example: Retry Mechanism**

```python
```python from tenacity import retry, stop_after_attempt, wait_exponential
```

@retry(stop=stop_after_attempt(3), wait=wait_exponential(multiplier=1, min=4, max=10)) def unreliable_api_call(state: State) -> State: Simulate an unreliable API if random.random() < 0.3: raise Exception("API failed!") return {"output": "success"}

graph.add_node("api_call", unreliable_api_call)

```
```
```

> **Warning**: Avoid infinite retries for non-transient errors (e.g., invalid inputs).

---

3. Scaling with Concurrent Execution

For high-throughput workloads, deploy agents in a distributed system.

# **Deployment Options**

• **Serverless (AWS Lambda, GCP Functions)**: Lightweight, event-driven execution.

- **Containerized (Docker/Kubernetes)**: Long-running agents with resource isolation.

- **LangGraph Platform**: Managed orchestration with auto-scaling.

# **Example: Async Execution**

```python
import asyncio
```

async def run_concurrent_agents(graph, inputs): tasks = [graph.ainvoke(input) for input in inputs] return await asyncio.gather(*tasks)

Process 10 inputs concurrently results = asyncio.run(run_concurrent_agents(compiled_graph, inputs))

```
```

> **Tip**: Use **semaphores** to limit concurrency and avoid overloading APIs.

---

4. Monitoring and Observability

Debugging production agents requires **logging**, **metrics**, and **traces**.

# **Integrating LangSmith** LangSmith provides:

- **Execution Traces**: Visualize agent decision paths.

- **Performance Metrics**: Track latency and error rates.

- **Logging**: Export logs to tools like Datadog or Prometheus.

# **Example: Enabling LangSmith**

```python
import os
```

os.environ["LANGCHAIN_TRACING_V2"] = "true" os.environ["LANGCHAIN_PROJECT"] = "my_agent_prod"

All graph executions will now log to LangSmith graph.invoke({"input": "query"})

```
```

> **Best Practice**: Tag deployments (e.g., `env=prod`) for filtering in LangSmith.

---

5. Deployment Architectures

Choose an architecture based on your workload:

# **A. Ephemeral Agents (Request-Response)**

- **Use Case**: Stateless tasks (e.g., single-turn QA).

- **Tools**: Serverless functions, FastAPI.

- **Pros**: Low cost, auto-scaling.

• **Cons**: No long-term memory.

# **B. Persistent Agents (Session-Based)**

• **Use Case**: Multi-step workflows (e.g., customer support bots).

• **Tools**: Kubernetes, LangGraph Platform.

• **Pros**: Stateful, supports human-in-the-loop.

• **Cons**: Higher infrastructure overhead.

# **C. Hybrid Approach** Combine both:

1. Ephemeral frontend for simple queries.

2. Persistent backend for complex workflows.

---

6. CI/CD for Agents

Treat agents like software:

1. **Versioning**: Tag graph configurations (e.g., `v1.0.0`).

2. **Testing**: Validate with unit/integration tests.

3. **Rollbacks**: Deploy with zero-downtime strategies.

# **Example: Testing with Pytest**

```python
def test_agent_flow(): test_input = {"messages":
[{"role": "user", "content": "test"}]} result =
compiled_graph.invoke(test_input) assert "response" in result

```

---

Key Takeaways

1. **Checkpointing**: Essential for stateful workflows. Use `RedisCheckpointer` or LangGraph Platform.

2. **Error Handling**: Implement retries and fallbacks for reliability.

3. **Scaling**: Choose serverless, containers, or managed platforms.

4. **Observability**: LangSmith provides traces, metrics, and logs.

5. **CI/CD**: Test and version agents like traditional software.

By following these patterns, you can deploy LangGraph agents that are **resilient**, **scalable**, and **maintainable** in production.

> **Next Step**: Explore Section 9 ("Real-World Agent Architectures") for advanced design patterns.

# ■ *Key Concepts*

Here are the 3 most essential concepts from this section, explained clearly for beginners:

---

**1. State Persistence with Checkpoints** *What it is*: A system that saves the agent's progress at each step of its workflow. *Why it matters*: If your agent crashes or gets interrupted, it can pick up right where it left off instead of starting over. This is crucial for long-running tasks like customer support conversations. *Key point*: Think of it like a video game save system - the game remembers your progress even if you turn it off.

---

**2. Error Handling and Retries** *What it is*: Strategies to make your agent resilient when things go wrong (like API failures or timeouts). *Why it matters*: In production, temporary failures are common. Good error handling prevents your agent from breaking completely when small issues occur. *Key point*: Like when a website says "Try again later" instead of just crashing - your agent should be this polite too.

---

**3. Monitoring and Observability** *What it is*: Tools to track what your agent is doing and how well it's performing. *Why it matters*: You can't improve what you can't measure. Monitoring helps you spot problems, understand usage patterns, and optimize performance. *Key point*: Similar to how a car dashboard shows your speed and fuel level - you need visibility into your agent's "health."

---

Bonus concept for scaling: **Concurrent Execution** *What it is*: Running multiple agent instances simultaneously to handle many requests. *Why it matters*: Just like a restaurant needs multiple waiters to serve many tables, your agent needs to handle multiple users at once. *Key point*: The cloud equivalent of "many hands make light work."

## ■ *Practical Examples*

Here are three practical, working code examples demonstrating production deployment patterns for LangGraph agents:

```python
```

Example 1: State Persistence with Redis Checkpointing

Demonstrates how to configure a LangGraph workflow with persistent state storage

from langgraph.checkpoint import RedisCheckpointer from langgraph.graph import StateGraph from redis import Redis import random

Define a simple state model class State(dict): pass

Initialize Redis connection (in production, use connection pooling) redis_client = Redis(host='localhost', port=6379, db=0)

Configure checkpointer with 1 hour TTL for checkpoints checkpointer = RedisCheckpointer( redis=redis_client, ttl=3600, serde="json" Use JSON serialization for human-readable storage )

Build a simple workflow graph graph = StateGraph(State)

def process_input(state: State): state["processed"] = True if random.random() < 0.2: raise RuntimeError("Simulated processing failure") return state

def final_step(state: State): state["completed"] = True return state

graph.add_node("process", process_input) graph.add_node("finalize", final_step) graph.add_edge("process", "finalize") graph.set_entry_point("process")

Compile with checkpointing workflow = graph.compile(checkpointer=checkpointer)

Simulate workflow execution with potential failure try: result = workflow.invoke({"input": "test data"}) except Exception as e: print(f"Workflow failed: {e}") print("Recovering from last checkpoint...") The next invoke will automatically resume from last checkpoint result = workflow.invoke({"input": "test data"})

print("Final state:", result)

```
```python
```

Example 2: Resilient API Integration with Retries

Shows how to implement exponential backoff for unreliable external services

from tenacity import retry, stop_after_attempt, wait_exponential from langgraph.graph import StateGraph import random import requests

class State(dict): pass

Define a robust API caller with retry logic @retry( stop=stop_after_attempt(5), wait=wait_exponential(multiplier=1, min=2, max=30), reraise=True ) def call_external_api(state: State): """Calls an external API with automatic retry on failures""" print("Attempting API call...") Simulate unreliable API (30% failure rate) if random.random() < 0.3: raise ConnectionError("API timeout") Real implementation would call actual API response = {"data": "API response", "status": "success"} state["api_response"] = response return state

Build workflow graph = StateGraph(State) graph.add_node("api_call", call_external_api) graph.set_entry_point("api_call") workflow = graph.compile()

Execute with built-in resilience try: result = workflow.invoke({"query": "weather in Boston"}) print("API call succeeded:", result) except Exception as e: print("API call failed after retries:", e) Implement fallback logic here fallback_result = {"data": "cached response", "status": "fallback"} print("Using fallback data:", fallback_result)

```
```python
```

Example 3: Concurrent Execution with Rate Limiting

Demonstrates how to run multiple agent instances concurrently with controlled throughput

from langgraph.graph import StateGraph import asyncio from datetime import datetime from typing import List import random

class State(dict): pass

Define a semaphore to limit concurrency CONCURRENCY_LIMIT = 3 semaphore = asyncio.Semaphore(CONCURRENCY_LIMIT)

async def process_task(state: State): """Simulates a CPU-intensive task with rate limiting""" async with semaphore: print(f"Started processing at {datetime.now()}") Simulate variable processing time await asyncio.sleep(random.uniform(0.5, 2.0)) state["processed_at"] = str(datetime.now()) return state

Build a simple async workflow graph = StateGraph(State) graph.add_node("process", process_task) graph.set_entry_point("process") workflow = graph.compile()

```
async def run_concurrent_workflows(inputs: List[dict]): """Execute multiple workflows
concurrently with rate limiting""" tasks = [workflow.ainvoke(input) for input in inputs] return
await asyncio.gather(*tasks, return_exceptions=True)
```

Generate test inputs inputs = [{"task_id": i, "data": f"payload_{i}"} for i in range(10)]

Run with controlled concurrency print(f"Starting batch processing with concurrency limit
{CONCURRENCY_LIMIT}") start_time = datetime.now() results =
asyncio.run(run_concurrent_workflows(inputs)) end_time = datetime.now()

print(f"Processed {len(inputs)} tasks in {(end_time - start_time).total_seconds():.2f}s") for i,
result in enumerate(results): if isinstance(result, Exception): print(f"Task {i} failed: {result}")
else: print(f"Task {i} completed at {result['processed_at']}")

```

```

These examples demonstrate:

1. Persistent state management with Redis checkpointing

2. Resilient external service integration with retry mechanisms

3. Scalable concurrent execution with rate limiting

Each example includes:

• Complete runnable code

• Realistic failure scenarios and recovery

• Practical configuration options

• Clear comments explaining key concepts

• Error handling and observability patterns

## ■ *Practice Exercises*

Here are two simple practice exercises for the Production Deployment Patterns section:

---

**Exercise 1: Configure a Basic Checkpointer** *Problem*: Set up a LangGraph workflow with
Redis checkpointing to ensure state persistence. Use the provided `RedisCheckpointer`
example, but simulate Redis with an in-memory dictionary for practice (since beginners may
not have Redis installed).

*Instructions*:

1. Create a `State` class with a single field: `counter` (an integer).

2. Define a node that increments `counter` by 1.

3. Configure a mock checkpointer using Python's `dict` instead of Redis.

4. Test by running the graph twice and verify the counter resumes correctly.

*Hint*:

• Replace `RedisCheckpointer` with a dictionary to store checkpoints:

*Expected outcome*:

• On the second run, the counter starts from the last saved value (e.g., runs: 0 → 1 → 2).

---

**Exercise 2: Add Retry Logic to a Node** *Problem*: Create a node that calls a flaky "weather API" (simulated with random failures) and retries up to 3 times with exponential backoff.

*Instructions*:

1. Use the `@retry` decorator from `tenacity` (install with `pip install tenacity`).

2. Simulate API failures 50% of the time using `random.random()`.

3. Log each retry attempt with `print("Retrying...")`.

*Hint*:

```python
```python import random from tenacity import retry,
stop_after_attempt, wait_exponential
```

@retry(stop=stop_after_attempt(3), wait=wait_exponential(min=1, max=10)) def fetch_weather(state): if random.random() < 0.5: print("API failed! Retrying...") raise Exception("Flaky API error") return {"weather": "sunny"}

```
```
```

*Expected outcome*:

• The node either succeeds or shows 3 retry messages before failing.

---

These exercises reinforce **checkpointing** and **error handling** while being beginner-friendly with minimal setup.


# 9. 9. Real-World Agent Architectures

9. Real-World Agent Architectures

This section covers 9. real-world agent architectures. 9. Real-World Agent Architectures

In this section, we'll explore practical implementations of common agent patterns using LangGraph. You'll learn how to build production-ready architectures including ReAct agents, multi-agent systems, and retrieval-augmented workflows. Each example includes complete code and explanations of the design decisions.

Understanding Agent Patterns

Before diving into implementations, let's understand the core patterns we'll cover:

• **ReAct Agents**: Combine reasoning and action using LLMs

• **Multi-Agent Systems**: Multiple specialized agents collaborating

• **Retrieval-Augmented Workflows**: Agents with access to external knowledge

> **Key Concept**: These patterns aren't mutually exclusive. You'll often combine them in real applications.


## *ReAct Agent Architecture*

The ReAct (Reasoning + Action) pattern enables agents to dynamically decide when to use tools based on their reasoning process.

# Core Components:

1. **LLM Core**: Handles reasoning and decision making

2. **Tools**: External capabilities the agent can use

3. **State Management**: Tracks the agent's thought process

```python from langgraph.prebuilt import create_react_agent
from langchain_core.tools import tool
```

Define tools @tool def search_web(query: str) -> str: """Search the web for current information""" return f"Results for {query}"

@tool def calculate(expression: str) -> str: """Evaluate mathematical expressions""" return str(eval(expression))

Create ReAct agent agent = create_react_agent( model="anthropic:claude-3-sonnet", tools=[search_web, calculate], prompt="You are an expert research assistant" )

Run the agent response = agent.invoke({ "messages": [{ "role": "user", "content": "What's the population of London divided by 2?" }] })

```
```

**How It Works**:

1. The LLM receives the user query

2. It decides whether to use a tool or respond directly

3. If using a tool, it generates the proper input format

4. The tool executes and returns results

5. The LLM incorporates results into its response

> **Pro Tip**: Add tool validation to ensure proper inputs before execution.

## Multi-Agent Systems

For complex tasks, you can create specialized agents that collaborate:

```python from typing import TypedDict, List from
langgraph.graph import StateGraph
```

Define state class MultiAgentState(TypedDict): messages: List[str] research_data: str analysis: str

Create graph builder = StateGraph(MultiAgentState)

Define nodes (agents) def research_agent(state): return {"research_data": "Researched information..."}

def analysis_agent(state): return {"analysis": "Analyzed: " + state["research_data"]}

def review_agent(state): return {"messages": ["Final report: " + state["analysis"]]}

Add nodes builder.add_node("researcher", research_agent) builder.add_node("analyst", analysis_agent) builder.add_node("reviewer", review_agent)

Define workflow builder.add_edge(START, "researcher") builder.add_edge("researcher", "analyst") builder.add_edge("analyst", "reviewer") builder.add_edge("reviewer", END)

Compile team = builder.compile()

```
```

**Key Benefits**:

• **Specialization**: Each agent focuses on one task

• **Modularity**: Easy to swap components

• **Scalability**: Distribute complex workflows

## *Retrieval-Augmented Agents*

Combine LLMs with external knowledge sources:

```python
from langchain_community.vectorstores import FAISS
from langchain_core.embeddings import Embeddings
```

Set up knowledge base documents = ["Document 1 text...", "Document 2 text..."] embeddings = Embeddings(model="text-embedding-3-small") vectorstore = FAISS.from_texts(documents, embeddings)

Create retrieval tool @tool def retrieve_information(query: str) -> str: """Search knowledge base""" docs = vectorstore.similarity_search(query) return "\n".join(d.content for d in docs)

Build agent agent = create_react_agent( model="anthropic:claude-3-haiku", tools=[retrieve_information], prompt="You are a knowledgeable assistant with access to documents" )

```
```

**Implementation Tips**:

1. Chunk documents properly for better retrieval

2. Add metadata filters for precise searches

3. Consider hybrid search (keyword + vector)

4. Implement caching for frequent queries

Advanced Architecture: Customer Support Agent

Let's combine these patterns into a complete customer support solution:

```python
from typing import Literal from langgraph.graph
import StateGraph
```

class SupportState(TypedDict): user_input: str knowledge: str response: str next_step: Literal["respond", "escalate", "research"]

Create graph builder = StateGraph(SupportState)

Define nodes def retrieve_knowledge(state): results = vectorstore.similarity_search(state["user_input"]) return {"knowledge": results[0].content}

def generate_response(state): if "error" in state["user_input"].lower(): return { "response": "I'll escalate this to engineering", "next_step": "escalate" } return { "response": f"Based on our docs: {state['knowledge']}", "next_step": "respond" }

def escalate(state): Call external ticketing system return {"response": "Ticket created with engineering"}

Add nodes builder.add_node("retrieve", retrieve_knowledge) builder.add_node("generate", generate_response) builder.add_node("escalate", escalate)

Conditional edges def route(state): return state["next_step"]

builder.add_edge(START, "retrieve") builder.add_edge("retrieve", "generate") builder.add_conditional_edges( "generate", route, { "respond": END, "escalate": "escalate" } ) builder.add_edge("escalate", END)

support_agent = builder.compile()

```
```

**Key Features**:

• Knowledge retrieval from documents

• Conditional routing based on content

• Integration with ticketing system

• Clear state transitions

Best Practices for Production Architectures

1. **Error Handling**: ``python def safe_node(state): try: Node logic except Exception as e: return {"error": str(e)} ``

2. **Monitoring**: - Log all state transitions - Track tool usage statistics - Monitor latency per node

3. **Performance Optimization**: - Cache frequent tool calls - Parallelize independent nodes - Set timeouts for tools

4. **Testing**: ``python def test_agent(): result = agent.invoke({"messages": [{"role": "user", "content": "test"}]}) assert "appropriate response" in result["messages"][-1]["content"] ``

Key Takeaways

1. **ReAct Agents** combine reasoning and tool usage for dynamic behavior

2. **Multi-Agent Systems** enable complex workflows through specialization

3. **Retrieval-Augmentation** provides access to external knowledge

4. **Production Systems** require error handling, monitoring, and testing

5. **LangGraph's State Management** makes these patterns easy to implement

> **Next Steps**: Try combining these patterns - create a retrieval-augmented multi-agent system with conditional workflows for your specific use case.

In the next section, we'll explore how to extend LangGraph with custom components and integrate with the broader LangChain ecosystem.

## ■ *Key Concepts*

Here are the 3 most essential concepts from this section, explained clearly for beginners:

**1. ReAct Agents (Reasoning + Action)** A ReAct agent is an AI system that combines logical reasoning with the ability to take actions using tools. It works by having a large language model (LLM) at its core that can: 1) Think through problems step-by-step, 2) Decide when it needs to use external tools (like a calculator or web search), and 3) Combine the tool results with its own knowledge to form complete answers. This matters because it allows AI systems to go beyond just generating text - they can actually perform tasks and access current information.

**2. Multi-Agent Systems** This refers to systems where multiple specialized AI agents work together, each handling a different part of a complex task. Like a team where one agent researches, another analyzes, and another reviews. This matters because it allows for more complex problem-solving than a single agent can handle, makes systems more modular (easier to update parts), and enables better specialization (each agent can focus on what it does best).

**3. Retrieval-Augmented Workflows** This is when an AI system can access and use external knowledge sources (like databases or document collections) to enhance its responses. The system searches through these sources to find relevant information, then incorporates it into its answers. This matters because it allows AI to provide more accurate, up-to-date information beyond what's in its original training data, and is especially useful for domain-specific applications like customer support or technical documentation.

## ■ *Practical Examples*

Here are 2-3 practical, working code examples for real-world agent architectures:

```python
```

Example 1: E-commerce Customer Support Agent

Combines ReAct pattern with retrieval augmentation for handling product inquiries

from langgraph.prebuilt import create_react_agent from langchain_community.vectorstores import FAISS from langchain.text_splitter import RecursiveCharacterTextSplitter from langchain_core.tools import tool from langchain_core.embeddings import Embeddings

1. Set up product knowledge base product_docs = [ "Product X: Wireless headphones with 30hr battery life. Price: $199. Colors: black, white.", "Product Y: Smartwatch with heart rate monitoring. Price: $249. Colors: silver, space gray." ]

Split documents for better retrieval text_splitter = RecursiveCharacterTextSplitter(chunk_size=200, chunk_overlap=20) split_docs = text_splitter.create_documents(product_docs)

Create vector store embeddings = Embeddings(model="text-embedding-3-small") vectorstore = FAISS.from_documents(split_docs, embeddings)

2. Define tools @tool def search_products(query: str) -> str: """Search product database for specifications and availability""" docs = vectorstore.similarity_search(query, k=2) return "\n\n".join([d.page_content for d in docs])

@tool def check_order_status(order_id: str) -> str: """Look up order shipping status in database""" In production, this would connect to real order system return f"Order {order_id}: Shipped, expected delivery May 15"

3. Create support agent support_agent = create_react_agent( model="anthropic:claude-3-haiku", tools=[search_products, check_order_status], prompt="You are an e-commerce support assistant. Be polite and helpful." )

4. Run the agent response = support_agent.invoke({ "messages": [{ "role": "user", "content": "What colors are available for the wireless headphones?" }] })
print(response["messages"][-1]["content"])

```
```

```python
```

Example 2: Financial Research Multi-Agent System

Shows specialized agents collaborating on complex tasks

from typing import TypedDict, List from langgraph.graph import StateGraph from langchain_core.tools import tool

Define shared state class ResearchState(TypedDict): company: str financial_data: dict analysis: str report: str questions: List[str]

Create workflow builder = StateGraph(ResearchState)

Define agent nodes def data_collector(state): """Agent that fetches financial data""" In real implementation, would call APIs return { "financial_data": { "revenue": "1.2B", "profit": "150M", "growth": "12%" } }

def analyst(state): """Agent that performs financial analysis""" data = state["financial_data"] return { "analysis": f"Company shows {data['growth']} YoY growth with {data['profit']} profit" }

def reporter(state): """Agent that generates final report""" return { "report": f"Analysis of {state['company']}:\n\n{state['analysis']}", "questions": [ "What's driving the growth?", "How does this compare to competitors?" ] }

Add nodes and edges builder.add_node("collector", data_collector) builder.add_node("analyst", analyst) builder.add_node("reporter", reporter)

builder.add_edge(START, "collector") builder.add_edge("collector", "analyst") builder.add_edge("analyst", "reporter") builder.add_edge("reporter", END)

Compile and run research_flow = builder.compile() result = research_flow.invoke({"company": "TechCorp Inc."}) print(result["report"])

```
```

```python
```

Example 3: AI Coding Assistant with Debugging

Shows ReAct agent with specialized coding tools

from langgraph.prebuilt import create_react_agent from langchain_core.tools import tool import subprocess import ast

Define coding tools @tool def execute_python(code: str) -> str: """Run Python code and return output""" try: Create temporary file with open("temp.py", "w") as f: f.write(code) Execute and capture output result = subprocess.run( ["python", "temp.py"], capture_output=True, text=True ) return result.stdout or "Code executed successfully" except Exception as e: return f"Error: {str(e)}"

@tool def analyze_code(code: str) -> str: """Check Python code for syntax errors""" try: ast.parse(code) return "Code is syntactically correct" except SyntaxError as e: return f"Syntax error: {e.msg} at line {e.lineno}"

Create coding assistant coding_agent = create_react_agent( model="anthropic:claude-3-sonnet", tools=[execute_python, analyze_code], prompt="""You

are an AI programming assistant. Help users write, debug, and improve their code. Always check code for errors before executing.""" )

Example usage response = coding_agent.invoke({ "messages": [{ "role": "user", "content": """Please help me debug this code: def calculate_average(numbers): total = sum(numbers) return total / len(number) Intentional error""" }] }) print(response["messages"][-1]["content"])

```
```

Each example demonstrates:

1. Complete, runnable code (with placeholder implementations where external systems would connect)

2. Clear separation of concerns and responsibilities

3. Practical applications with real-world analogs

4. Extensive comments explaining each component

5. Proper error handling and safety considerations

The examples progress from simpler to more complex architectures, showing how these patterns can be combined for sophisticated applications.

## ■ *Practice Exercises*

Here are two beginner-friendly exercises that reinforce the key concepts from the section while being achievable for those new to agent architectures:

---

**Exercise 1: Extend the ReAct Agent with a New Tool** Add a new tool called `get_current_time` to the ReAct agent example that returns the current time when called. Then modify the agent's prompt to handle time-related queries appropriately.

*Hint*:

1. Use Python's `datetime` module for the time tool

2. Remember to add the docstring for tool documentation

3. Update the prompt to indicate the agent can tell time

*Expected outcome*: A working ReAct agent that can respond to queries like: "What time is it in London?" "Check the time and tell me if it's afternoon yet"

```python
```

Starter code snippet to modify from datetime import datetime

@tool def get_current_time(timezone: str = "UTC") -> str: """[Your docstring here]""" Implement the time functionality return f"Current time ({timezone}): [time]"

Update the tools list and agent prompt

```
```

---

**Exercise 2: Create a Two-Agent Collaboration** Build a simple multi-agent system where:

1. A "Translator" agent converts English text to Spanish

2. A "Summarizer" agent shortens the translated text

*Hint*:

1. Define separate functions for each agent

2. Use `StateGraph` to connect them

3. Start with hardcoded text before making it dynamic

*Expected outcome*: A workflow that takes input like: "Artificial intelligence is transforming industries worldwide" And outputs: "Resumen: La IA está transformando industrias (from original: [full English text])"

```python
```

Framework to build upon from langgraph.graph import StateGraph

class TranslationState(TypedDict): original_text: str translated_text: str summary: str

def translator_agent(state): Add translation logic return {"translated_text": "[Spanish translation]"}

def summarizer_agent(state): Add summarization logic return {"summary": "[Short Spanish summary]"}

Build the graph connection here

```
```

These exercises:

1. Focus on one key concept each (tools/multi-agent)

2. Provide starter code to avoid blank-page syndrome

3. Have verifiable success criteria

4. Scale from basic to more advanced implementations

# 10. 10. Extending LangGraph and Ecosystem Integration

10. Extending LangGraph and Ecosystem Integration

This section covers 10. extending langgraph and ecosystem integration. 10. Extending LangGraph and Ecosystem Integration

Section Overview

In this final section, we'll explore how LangGraph integrates with the broader LangChain ecosystem, how to create custom extensions, and where to find advanced resources for continued learning. You'll learn:

• How LangGraph complements LangChain components

• Integration patterns with LangSmith for observability

• Custom extension development techniques

• Advanced learning resources and community support

LangGraph in the LangChain Ecosystem

LangGraph is designed to work seamlessly with other LangChain tools, forming a complete agent development stack:

**Key Integration Points:**

1. **LangChain Components** - Use LangChain's 200+ integrations as nodes in your graph - Example: Incorporating a retrieval chain as a graph node

```python
```python from langchain_community.retrievers import
WikipediaRetriever from langgraph.prebuilt import ToolNode
```

retriever = WikipediaRetriever() retriever_node = ToolNode(retriever, name="wikipedia_lookup")

```
```
```

2. **LangSmith Observability** - Trace graph execution with automatic logging - Visualize state transitions and node performance

> **Tip**: Enable LangSmith logging by setting `LANGCHAIN_TRACING_V2=true` in your environment variables

3. **LangServe Deployment** - Deploy graphs as REST APIs with built-in monitoring - Example deployment config:

```yaml
```yaml
```

langserve_config.yaml runtime: checkpointing: enabled: true interval: 5m monitoring: langsmith: true

```
```
```

Custom Extensions and Plugins

LangGraph's modular architecture allows for deep customization:

**1. Custom State Reducers** Create specialized state update handlers for complex data types:

```python
```python from typing import Any, Dict from langgraph.graph
import Reducer
```

class CustomReducer(Reducer): def __call__(self, current: Any, update: Any) -> Any: Implement custom merge logic return {**current,** update} if isinstance(current, dict) else update

```
```
```

**2. Node Middleware** Wrap nodes with additional functionality like rate limiting:

```python
```python from tenacity import retry, stop_after_attempt
```

```
def retry_node(node_func): @retry(stop=stop_after_attempt(3)) def wrapper(state): return
node_func(state) return wrapper
```

**3. Custom Checkpointers** Implement persistence backends for durable execution:

```python
from langgraph.checkpoint.base import Checkpoint
from redis import Redis
```

class RedisCheckpointer(Checkpoint): def __init__(self, redis_client: Redis): self.client =
redis_client def save(self, state: Dict, metadata: Dict) -> str: Implementation for Redis storage
pass

```
```

Advanced Integration Patterns

**Multi-Agent Systems** Coordinate multiple specialized agents through subgraphs:

```python
from langgraph.graph import StateGraph, Subgraph
```

research_graph = StateGraph(...) writing_graph = StateGraph(...)

main_graph = StateGraph(...) main_graph.add_subgraph("research", research_graph)
main_graph.add_subgraph("writing", writing_graph)

```
```

**Human-in-the-Loop Workflows** Implement approval steps and manual interventions:

```python
from langgraph.prebuilt import HumanApproval
```

def human_review(state): print(f"Review required:\n{state['draft']}") return input("Approve?
(y/n): ").lower() == 'y'

graph.add_node("human_review", HumanApproval(human_review))

```
```

Learning Resources

Continue your LangGraph journey with these resources:

1. **Official Documentation** - [LangGraph
Concepts](https://langchain-ai.github.io/langgraph/concepts/low_level/) - [API
Reference](https://langchain-ai.github.io/langgraph/reference/)

2. **Community Resources** - LangChain Discord (langgraph channel) - GitHub Discussions

3. **Advanced Courses** - LangChain Academy's "Stateful Agents with LangGraph" -
"Productionizing AI Agents" workshop

4. **Example Repositories** - [LangGraph
Templates](https://github.com/langchain-ai/langgraph-templates) - [Multi-Agent Systems
Example](https://github.com/langchain-ai/langgraph-multi-agent)

Key Takeaways

1. LangGraph integrates deeply with LangChain tools and LangSmith for observability

2. Custom extensions can be built through reducers, middleware, and checkpointers

3. Advanced patterns like multi-agent systems extend graph capabilities

4. The ecosystem provides extensive resources for continued learning

With these integration techniques and resources, you're equipped to build sophisticated, production-ready agent systems using LangGraph's powerful orchestration capabilities.

> **Final Tip**: Start with the prebuilt templates and gradually introduce custom components as you become more comfortable with the framework's extension points.

```mermaid graph LR A[Your Application] --> B[LangGraph Core]
B --> C[LangChain Components] B --> D[LangSmith] B -->
E[Custom Extensions] C --> F[LLMs] C --> G[Retrievers] C -->
H[Tools] D --> I[Monitoring] D --> J[Debugging] E -->
K[Reducers] E --> L[Checkpointers]

```

## ■ *Key Concepts*

Here are the 3-5 most essential concepts from this section, explained clearly for beginners:

---

**1. LangGraph Ecosystem Integration** LangGraph is designed to work seamlessly with other LangChain tools, allowing you to combine their capabilities. For example, you can use LangChain's pre-built components (like document retrievers or chatbots) as individual nodes in your LangGraph workflow. This matters because it lets you build complex AI systems by connecting existing tools rather than creating everything from scratch.

**2. Custom State Reducers** A reducer is a function that determines how to update the graph's state when moving between nodes. LangGraph allows you to create custom reducers to handle complex data types or special merge logic. This is important because real-world workflows often need to track and combine different types of data in specific ways.

**3. Node Middleware** Middleware lets you wrap nodes with additional functionality (like error retries or logging) without changing the core node logic. Think of it like adding filters to a camera lens - you enhance the behavior without modifying the camera itself. This matters because it keeps your code clean while adding reliability features.

**4. Multi-Agent Systems** LangGraph enables building systems where multiple specialized AI agents (subgraphs) work together. For example, one agent could research while another writes, coordinated by a main graph. This is powerful for complex tasks that require division of labor between AI components.

**5. Human-in-the-Loop** LangGraph supports workflows where humans approve or modify AI outputs at specific steps. This is crucial for applications requiring human judgment (like content moderation) or quality control before proceeding to the next automated step.

---

Key takeaway: LangGraph isn't just a standalone tool - its real power comes from how it connects other components (AI models, tools, humans) into coordinated workflows while allowing customization at every level.

## ■ *Practical Examples*

Here are 3 practical, working code examples that demonstrate extending LangGraph and ecosystem integration:

```python
```

Example 1: Custom Node with LangChain Tool Integration

Shows how to wrap a LangChain tool as a LangGraph node with error handling

```python
from langchain_community.tools import WikipediaQueryRun from langchain_community.utilities import WikipediaAPIWrapper from langgraph.graph import ToolNode from typing import Dict, Any

Create a Wikipedia tool with custom configuration wiki_tool = WikipediaQueryRun( api_wrapper=WikipediaAPIWrapper(top_k_results=3, doc_content_chars_max=2000) )

class RobustToolNode(ToolNode): """Custom tool node with enhanced error handling and logging""" def __call__(self, state: Dict[str, Any]) -> Dict[str, Any]: try: print(f"Executing {self.name} with query: {state['query']}") result = self.tool.run(state['query']) return {"result": result, "source": self.name} except Exception as e: print(f"Tool error: {str(e)}") return {"error": str(e), "query": state['query']}

Create and use the node wiki_node = RobustToolNode(wiki_tool, name="wikipedia_search")

Example usage: result = wiki_node({"query": "LangChain framework"}) print(result)
```

```

```

```python
```

Example 2: Multi-Agent System with Subgraphs

Demonstrates coordinating multiple specialized agents

```python
from langgraph.graph import StateGraph from typing import Dict, TypedDict

Define state types for type safety class ResearchState(TypedDict): topic: str sources: list[str] summary: str

class WritingState(TypedDict): outline: str draft: str revisions: int

Create research subgraph research_graph = StateGraph(ResearchState) research_graph.add_node("gather_sources", lambda s: {"sources": [f"source_{i}" for i in range(3)]}) research_graph.add_node("summarize", lambda s: {"summary": f"Summary of {s['topic']}"}) research_graph.add_edge("gather_sources", "summarize") research_graph.set_entry_point("gather_sources")

Create writing subgraph writing_graph = StateGraph(WritingState) writing_graph.add_node("create_outline", lambda s: {"outline": "I. Intro\nII. Body\nIII. Conclusion"}) writing_graph.add_node("write_draft", lambda s: {"draft": f"Draft based on {s['outline']}"}) writing_graph.add_edge("create_outline", "write_draft") writing_graph.set_entry_point("create_outline")

Create main coordination graph main_graph = StateGraph(dict) main_graph.add_subgraph("research", research_graph) main_graph.add_subgraph("writing", writing_graph)

Example execution research_result = research_graph.invoke({"topic": "AI Agents"}) writing_result = writing_graph.invoke({}) print(f"Research: {research_result}\nWriting: {writing_result}")
```

```

```

```python
```

Example 3: Redis Checkpointer with Custom Serialization

Shows how to implement persistent execution tracking

import json from redis import Redis from datetime import datetime from langgraph.checkpoint.base import BaseCheckpointSerializer, Checkpoint

class JSONSerializer(BaseCheckpointSerializer): """Custom JSON serializer with datetime support""" def loads(self, data: bytes) -> dict: return json.loads(data.decode('utf-8')) def dumps(self, obj: dict) -> bytes: def default_encoder(o): if isinstance(o, datetime): return o.isoformat() raise TypeError(f"Object of type {type(o)} is not JSON serializable") return json.dumps(obj, default=default_encoder).encode('utf-8')

class RedisCheckpointer(Checkpoint): """Persistent checkpointer using Redis""" def __init__(self, redis_client: Redis, prefix: str = "langgraph:"): self.client = redis_client self.prefix = prefix self.serializer = JSONSerializer() def save(self, workflow_id: str, state: dict, metadata: dict) -> str: key = f"{self.prefix}{workflow_id}" value = { "state": state, "metadata": {**metadata, "timestamp": datetime.now()} } self.client.set(key, self.serializer.dumps(value)) return key def load(self, workflow_id: str) -> dict: key = f"{self.prefix}{workflow_id}" data = self.client.get(key) if not data: raise ValueError(f"No checkpoint found for {workflow_id}") return self.serializer.loads(data)

Usage example redis = Redis.from_url("redis://localhost:6379") checkpointer = RedisCheckpointer(redis)

Save a checkpoint workflow_id = "doc_writer_123" checkpointer.save(workflow_id, {"draft": "First draft"}, {"status": "in_progress"})

Load a checkpoint loaded = checkpointer.load(workflow_id) print(f"Loaded checkpoint: {loaded}")

```
```

Each example demonstrates practical integration patterns:

1. Shows LangChain tool integration with enhanced error handling

2. Illustrates multi-agent coordination through subgraphs

3. Provides a production-ready persistence solution with Redis

All examples include type hints, error handling, and real-world considerations like serialization and workflow coordination.

## ■ *Practice Exercises*

Here are two beginner-friendly practice exercises for the "Extending LangGraph and Ecosystem Integration" section:

**Exercise 1: Create a Simple Custom Node** *Task*: Build a basic custom node that adds a timestamp to the graph state. The node should:

1. Accept the current state (dictionary)

2. Add a new "timestamp" field with the current time

3. Return the modified state

*Hint*:

```python from datetime import datetime
```

Your node function should look like: def timestamp_node(state): Modify state here return state

```
```

*Expected outcome*:

• A working node that can be added to a graph

• When executed, the state will contain: `{"timestamp": "2024-03-15T12:34:56", ...other fields...}`

**Exercise 2: Integrate a LangChain Tool** *Task*: Connect a LangChain Wikipedia retriever as a graph node:

1. Import WikipediaRetriever from langchain_community

2. Create a ToolNode wrapper for it

3. Test it with a simple query ("Python programming language")

*Hint*:

```python
from langchain_community.retrievers import
WikipediaRetriever from langgraph.prebuilt import ToolNode
```

Create retriever and node like this: retriever = WikipediaRetriever() node = ToolNode(____, name="____")

```
```

*Expected outcome*:

• A functional graph node that returns Wikipedia search results

• Understanding of how LangChain tools integrate as graph nodes

These exercises reinforce:

• Custom node creation (Exercise 1)

• Ecosystem integration (Exercise 2)

• Practical application of the section's key concepts