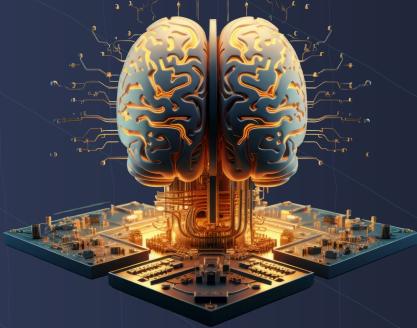


ARTIFICIAL INTELLIGENCE

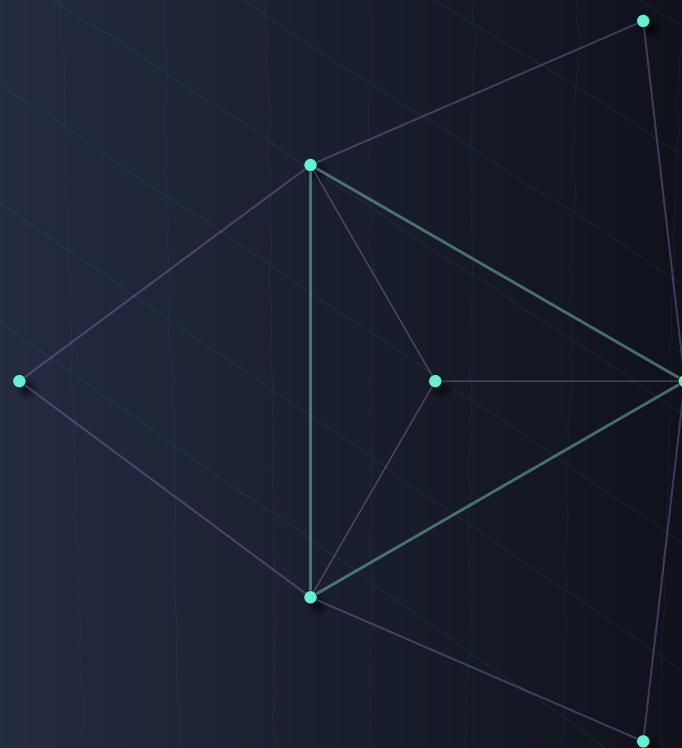


Multi-Agent Deal Discovery System

A Comprehensive Python Programming Guide

From Architecture to Implementation:

Machine Learning, Web Scraping, and System Design



SIJAN PAUDEL

Generated Technical Documentation

Contents

Preface	xvii
1 LLM Cheat Sheet	1
1.1 Handling Imports from Parent Folders in Python and Jupyter	1
1.1.1 Using <code>__file__</code> in Python Scripts	1
1.2 Comparing LLMs – The Basics (1)	3
1.3 Comparing LLMs – The Basics (2)	3
1.4 Chinchilla Scaling Law	4
1.5 Common LLM Benchmarks	5
1.6 Specific Benchmarks	5
1.7 Limitations of Benchmarks	6
1.8 Advanced Benchmarks for Large Language Models	7
1.9 Leaderboard Image	9
1.10 A 5-Step Strategy for Using AI Models	10
1.10.1 Understand (Figure Out What You Need)	11
1.10.2 Prepare (Do Your Homework)	12
1.10.3 Select (Pick Your Tools and Test Them)	13
1.10.4 Customize (Make the AI Better at Your Task)	14
1.10.5 Productionize (Get It Ready for the Real World)	18
1.11 LoRA (Low-Rank Adaptation)	20
1.11.1 Introduction	20
1.12 LoRA Hyperparameters: r , α , and Target Modules	23
1.12.1 Rank r (Low-Rank Dimension)	23
1.12.2 Summary Table	24
1.13 Q, K, V, O in Transformers	25
1.13.1 Analogy: Classroom Example	25
1.14 LoRA in PeftModelForCausalLM	26
1.14.1 Model Architecture Overview	26
1.14.2 LoRA in Attention Layers	26
1.14.3 Dimensions of Q, K, V, O Projections	27

1.14.4	Summary	27
1.15	Loading a Base Model in 4-Bit and Using LoRA	29
1.15.1	Explanation of the Code	29
1.15.2	How LoRA Optimizes This	29
1.15.3	Summary	30
1.16	Decision: Selecting a Base Model	31
1.16.1	Goal	31
1.16.2	Decision criteria (short checklist)	31
1.16.3	Model-by-model short summary	31
1.16.4	Comparison table (high-level)	32
1.16.5	Base vs Instruct variants	32
1.16.6	Rules of thumb and recommended workflows	33
1.16.7	Practical checklist before committing	33
1.16.8	Short recommendation examples	33
1.16.9	Why LLaMA 3.1 was chosen: Tokenization Convenience	34
1.17	Key Hyperparameters in QLoRA	35
1.17.1	Quick Summary Table	36
1.18	Training Hyperparameters	38
1.18.1	Epochs	38
1.19	DataCollatorForCompletionOnlyLM in HuggingFace TRL	41
1.19.1	Problem	41
1.19.2	Traditional Approach (Complicated)	41
1.19.3	HuggingFace TRL Solution: DataCollatorForCompletionOnlyLM	41
1.19.4	Example	41
1.19.5	Why it is important	42
1.19.6	Warmup Ratio (<code>warmup_ratio</code>)	42
1.20	Checkpointing and Resuming Training in HuggingFace / Colab	43
1.20.1	Overview	43
1.20.2	Checkpoint Components	43
1.20.3	Manual Checkpointing in PyTorch	44
1.20.4	Automatic Checkpointing with HuggingFace Trainer / SFTTrainer	44
1.20.5	Checkpointing Flow Step-by-Step	45
1.20.6	Checkpointing in Google Colab	45
1.20.7	Advantages of Using SFTTrainer Checkpointing	45
1.20.8	Checkpointing in Kaggle	47
1.21	Saving and Downloading Automatic Checkpoints in Kaggle	49
1.21.1	When to Enter the Zip Command	49
1.21.2	How to Download the Zip File	49
1.21.3	Restoring Checkpoints from Zip	49

1.22	Saving and Resuming Checkpoints with Hugging Face Hub	51
1.22.1	Manual Upload to Hugging Face Hub	51
1.22.2	Manual Download from Hugging Face Hub	51
1.22.3	Why Use Hugging Face Hub?	51
1.23	Automatic Checkpointing to Hugging Face Hub	53
1.23.1	Authenticate with Hugging Face	53
1.24	Resuming Fine-Tuning with SFTTrainer and LoRA	57
1.24.1	Define Constants and Paths	57
1.24.2	Checkpointing in SFTTrainer with LoRA	62
1.24.3	Inferencing using SFT	63
1.24.4	Improved Prediction Function	65
1.24.5	Key Points for SFT Inference	66
1.24.6	Summary	66
1.25	Four Steps in LoRA Training	67
1.25.1	Forward Pass with LoRA	67
1.25.2	Summary Table of LoRA Training Steps	68
1.25.3	Uploading Final Dataset to Hugging Face Hub	69
1.26	Multi-Agent AI Architecture for Automated Deal Finding Systems	70
1.26.1	User Interface (UI)	70
1.26.2	Agent Framework	70
1.26.3	Planning Agent	70
1.26.4	Scanner Agent	70
1.26.5	Ensemble Agent	70
1.26.6	Messaging Agent	71
1.26.7	Pipeline Summary	71
1.26.8	Future Work	71
1.27	Create a RAG Database with Our 400,000 Training Data	72
1.28	Hallmarks of an Agentic AI Solution	75
2	System Overview and Architecture	77
2.1	Introduction	77
2.1.1	System Purpose and Goals	77
2.1.2	Technical Architecture Overview	78
2.2	Agent-Based Architecture	78
2.2.1	The Agent Pattern	78
2.2.2	System Components	78
2.3	Data Flow and System Interactions	79
2.3.1	Primary Workflow	79
2.3.2	Data Models	80

2.4	Advanced System Features	80
2.4.1	Machine Learning Integration	80
2.4.2	External Service Integration	81
2.4.3	Error Handling and Robustness	81
2.5	Architectural Patterns	81
2.5.1	Design Patterns Implemented	81
2.5.2	SOLID Principles	82
2.6	System Scalability and Performance	82
2.6.1	Scalability Considerations	82
2.6.2	Performance Characteristics	82
2.7	Security and Privacy	82
2.7.1	Security Measures	82
2.7.2	Privacy Considerations	83
2.8	Chapter Summary	83
3	File-by-File Detailed Analysis	85
3.1	agent.py - The Foundation Base Class	85
3.1.1	File Structure and Imports	85
3.1.2	Class Analysis: Agent	86
3.1.3	Inter-File Relationships	87
3.2	deals.py - Data Models and External Integration	87
3.2.1	Imports and Dependencies	87
3.2.2	Global Configuration	88
3.2.3	Utility Function: extract()	88
3.2.4	Class Analysis: ScrapedDeal	89
3.2.5	Pydantic Models	91
3.3	ensemble_agent.py - Meta-Learning Orchestration	92
3.3.1	File Structure and Dependencies	92
3.3.2	Class Analysis: EnsembleAgent	93
3.4	frontier_agent.py - RAG-Based LLM Integration	95
3.4.1	Import Dependencies	95
3.4.2	Class Analysis: FrontierAgent	95
3.5	messaging_agent.py - Multi-Channel Communication	99
3.5.1	Dependencies and Configuration	99
3.5.2	Class Analysis: MessagingAgent	100
3.6	Chapter Summary	102
4	Class Hierarchy and Inheritance	103
4.1	Overview of Class Hierarchy	103

4.1.1	Inheritance Tree Structure	103
4.2	Base Class Analysis: Agent	103
4.2.1	Agent Class Architecture	104
4.2.2	Class Attributes vs Instance Attributes	104
4.2.3	Method Resolution Order (MRO)	105
4.3	Concrete Agent Classes	105
4.3.1	PlanningAgent Inheritance	105
4.3.2	EnsembleAgent Inheritance	106
4.3.3	SpecialistAgent Inheritance	106
4.4	Polymorphism in Action	107
4.4.1	Price Estimation Polymorphism	107
4.4.2	Runtime Polymorphism Example	107
4.5	Advanced Inheritance Concepts	108
4.5.1	Method Override Analysis	108
4.5.2	Abstract Methods and Duck Typing	109
4.6	Composition vs Inheritance	109
4.6.1	EnsembleAgent Composition Analysis	109
4.6.2	When to Use Inheritance vs Composition	110
4.7	Method Resolution and Attribute Access	110
4.7.1	Attribute Resolution Process	110
4.7.2	Dynamic Attribute Creation	111
4.8	Special Methods and Magic Methods	111
4.8.1	<code>__repr__</code> Implementation	111
4.8.2	Class Methods and Alternative Constructors	112
4.9	Modern Python Features	112
4.9.1	Type Annotations in Inheritance	112
4.9.2	Dataclasses and Pydantic Integration	113
4.10	Inheritance Patterns and Design Principles	113
4.10.1	Template Method Pattern	113
4.10.2	Strategy Pattern Through Polymorphism	114
4.11	Chapter Summary	114
5	Function-by-Function Breakdown	117
5.1	Agent Base Class Methods	117
5.1.1	<code>Agent.log()</code> Method	117
5.2	Data Processing Functions	119
5.2.1	<code>extract()</code> Function	119
5.3	Class Methods and Alternative Constructors	121

5.3.1	ScrapedDeal.fetch() Class Method	121
5.4	Machine Learning Integration Methods	123
5.4.1	EnsembleAgent.price() Method	123
5.5	External API Integration Methods	126
5.5.1	FrontierAgent.find_similars() Method	126
5.6	Error Handling and Robustness Patterns	129
5.6.1	MessagingAgent.push() Method	129
5.7	Complex Business Logic Methods	132
5.7.1	PlanningAgent.plan() Method	132
5.8	Chapter Summary	135
6	Python Concepts and Examples	137
6.1	Object-Oriented Programming Concepts	137
6.1.1	Classes and Objects	137
6.1.2	Inheritance	138
6.1.3	Composition	139
6.2	Modern Python Features	141
6.2.1	Type Annotations	141
6.2.2	F-Strings (Formatted String Literals)	142
6.2.3	List Comprehensions	144
6.3	Advanced Python Concepts	145
6.3.1	Decorators	145
6.3.2	Context Managers	147
6.4	Data Structures and Collections	149
6.4.1	Dictionaries and Dictionary Operations	149
6.4.2	Lists and List Operations	151
6.5	Exception Handling	153
6.6	Functional Programming Concepts	156
6.6.1	Lambda Functions	156
6.7	File I/O and Serialization	158
6.8	Chapter Summary	160
6.8.1	Core OOP Concepts	160
6.8.2	Modern Python Features	160
6.8.3	Advanced Concepts	161
6.8.4	Data Structures	161
7	Agent Interaction and Execution Flow	163
7.1	System Execution Overview	163

7.1.1	Complete Workflow Execution	163
7.2	Detailed Sequence Analysis	163
7.2.1	Complete System Execution Sequence	163
7.2.2	Step-by-Step Execution Breakdown	163
7.3	Individual Agent Execution Flows	167
7.3.1	SpecialistAgent Execution Flow	167
7.3.2	FrontierAgent RAG Execution Flow	169
7.4	Parallel vs Sequential Execution Analysis	172
7.4.1	Current Sequential Processing	172
7.4.2	Potential Parallel Processing	173
7.5	Error Handling and Recovery Flows	174
7.5.1	Current Error Handling Limitations	174
7.5.2	Enhanced Error Handling Flow	176
7.6	Memory and State Management	178
7.6.1	Current Memory Usage Pattern	178
7.6.2	Enhanced State Management	178
7.7	Chapter Summary	181
7.7.1	Key Interaction Patterns	181
7.7.2	Execution Flow Analysis	181
7.7.3	Performance Optimization Opportunities	182
7.7.4	System Robustness	182
8	Summary and Essential Python Knowledge	183
8.1	System Summary and Functionality Recap	183
8.1.1	System Overview	183
8.1.2	Agent Responsibilities Summary	184
8.2	Python Knowledge Checklist	184
8.2.1	Fundamental Python Concepts	185
8.2.2	Advanced Python Features	186
8.2.3	External Libraries and Integration	187
8.3	Design Patterns and Architecture	187
8.3.1	Implemented Design Patterns	187
8.3.2	SOLID Principles Application	189
8.4	System Strengths and Areas for Improvement	190
8.4.1	System Strengths	190
8.4.2	Areas for Improvement	191
8.5	Enhancement Recommendations	192
8.5.1	Immediate Improvements (High Impact, Low Effort)	192

8.5.2	Medium-Term Improvements (Moderate Impact, Moderate Effort)	193
8.5.3	Long-Term Improvements (High Impact, High Effort)	194
8.6	Future Extensibility	195
8.6.1	Plugin Architecture	196
8.6.2	Monitoring and Observability	197
8.7	Chapter Summary and Final Assessment	198
8.7.1	System Achievement Summary	198
8.7.2	Key Learning Outcomes	199
8.7.3	Final Recommendations	200
A	Quick Reference Guide	201
A.1	Agent Class Summary	201
A.2	Python Concepts Index	201
A.2.1	Object-Oriented Programming	201
A.2.2	Advanced Features	202
A.2.3	External Libraries	202
B	System Extension Guidelines	203
B.1	Adding New Agents	203
B.2	Configuration Management	203
C	Troubleshooting Guide	205
C.1	Common Issues	205
C.1.1	Import Errors	205
C.1.2	API Authentication	205
C.1.3	Memory Issues	205
C.2	Performance Optimization	205
References and Further Reading		207
C.3	Python Documentation	207
C.4	Machine Learning Resources	207
C.5	Web Development	207
C.6	Software Engineering	207
Index		209
C.7	Agents	210
C.7.1	Agents Module	210
C.7.2	Deals Module	211
C.7.3	Ensemble Agent	214
C.7.4	Messaging Agent	216

C.7.5	Frontier Agent	218
C.7.6	Planning Agent	221
C.7.7	Random Forest Agent	223
C.7.8	Scanner Agent	224
C.7.9	Scanner Agent Using Langchain	227
C.7.10	Specialist Agent	230

List of Figures

2.1	High-Level System Architecture showing agent relationships and data flow	78
4.1	Complete Class Hierarchy of the Multi-Agent System	103
7.1	High-Level System Execution Flow	164
7.2	Detailed Sequence Diagram of Agent Interactions	165
7.3	SpecialistAgent Execution Flow	168
7.4	FrontierAgent RAG Execution Flow	170
7.5	Enhanced Error Handling Flow	175

List of Tables

2.1 External Service Integrations	81
4.1 Inheritance vs Composition Decision Matrix	110
7.1 Performance Comparison of Execution Modes	174
8.1 Agent Responsibilities and Technologies	184
A.1 Complete Agent Reference	201

Listings

2.1	Core Data Models	80
3.1	agent.py - Complete File Analysis	85
3.2	Detailed log() Method Analysis	86
3.3	deals.py - Import Analysis	87
3.4	RSS Feed Configuration	88
3.5	HTML Text Extraction Function	88
3.6	ScrapedDeal Class Implementation	89
3.7	ScrapedDeal Instance Methods	90
3.8	ScrapedDeal.fetch() Class Method	90
3.9	Pydantic Model Definitions	91
3.10	ensemble_agent.py - Complete Implementation	92
3.11	Ensemble Price Prediction Method	94
3.12	frontier_agent.py - Dependencies	95
3.13	FrontierAgent Class Structure	95
3.14	Context Generation Method	96
3.15	OpenAI Message Construction	97
3.16	Vector Similarity Search	97
3.17	Complete Price Prediction Pipeline	98
3.18	messaging_agent.py - Configuration	99
3.19	MessagingAgent Constructor	100
3.20	Push Notification Implementation	101
3.21	Opportunity Alert Processing	101
4.1	Agent Base Class - Complete Analysis	104
4.2	Agent MRO Analysis	105
4.3	PlanningAgent Class Definition	105
4.4	EnsembleAgent Inheritance Pattern	106
4.5	SpecialistAgent Remote Integration	106
4.6	Polymorphic Price Methods	107
4.7	Dynamic Agent Selection	107
4.8	Scanner Agent Polymorphism	108
4.9	Informal Protocol Definition	109
4.10	Composition Pattern Analysis	109
4.11	Attribute Resolution Example	110
4.12	Dynamic Attribute Example	111
4.13	Custom __repr__ Implementation	111
4.14	Class Method as Alternative Constructor	112
4.15	Advanced Type Annotations	112
4.16	Pydantic vs Traditional Classes	113
4.17	Template Method Pattern	114
4.18	Strategy Pattern Implementation	114

5.1	Agent.log() Method Deep Analysis	117
5.2	Color Code Memory Analysis	118
5.3	F-String Processing Analysis	118
5.4	Logging System Integration	118
5.5	HTML Processing Function Analysis	119
5.6	HTML Parser Initialization	119
5.7	DOM Navigation Analysis	120
5.8	Conditional Text Extraction	120
5.9	Newline Normalization	120
5.10	ScrapedDeal.fetch() Deep Analysis	121
5.11	Container Initialization	122
5.12	Conditional Progress Bar	122
5.13	RSS Feed Iteration	122
5.14	Result Return	123
5.15	Ensemble Price Prediction Analysis	123
5.16	Logging and Setup	124
5.17	Parallel Model Invocation	125
5.18	DataFrame Construction Analysis	125
5.19	Ensemble Model Execution	125
5.20	Completion Processing	126
5.21	Vector Similarity Search Analysis	126
5.22	RAG Process Announcement	127
5.23	Text to Vector Transformation	127
5.24	ChromaDB Similarity Search	128
5.25	Query Result Extraction	128
5.26	Result Return Processing	129
5.27	Push Notification Error Handling	129
5.28	Notification Process Start	130
5.29	HTTPS Connection Establishment	130
5.30	HTTP Request Construction	131
5.31	POST Request Transmission	131
5.32	Response Processing	131
5.33	Enhanced Error Handling Pattern	132
5.34	Master Planning Method Analysis	132
5.35	Mutable Default Problem	133
5.36	Planning Process Start	133
5.37	Scanner Agent Invocation	133
5.38	Deal Processing Branch	134
5.39	List Comprehension Execution	134
5.40	Opportunity Ranking	134
5.41	Threshold Decision Logic	135
6.1	Basic Class Concept	137
6.2	Agent Class Implementation	137
6.3	Basic Inheritance Example	138
6.4	Agent Inheritance Hierarchy	139
6.5	Basic Composition Example	140
6.6	EnsembleAgent Composition Pattern	140
6.7	Type Annotation Examples	141

6.8	System Type Annotations	141
6.9	F-String Examples	143
6.10	F-Strings in Agent System	143
6.11	List Comprehension Examples	144
6.12	List Comprehensions in Agent System	144
6.13	Decorator Examples	145
6.14	Decorators in Agent System	146
6.15	Context Manager Examples	147
6.16	Context Managers in Agent System (Potential)	148
6.17	Dictionary Examples	149
6.18	Dictionary Usage in Agent System	150
6.19	List Examples	151
6.20	List Usage in Agent System	152
6.21	Exception Handling Examples	153
6.22	Exception Handling in Agent System (Current and Potential)	154
6.23	Lambda Function Examples	156
6.24	Lambda Functions in Agent System	157
6.25	File I/O and Serialization Examples	158
6.26	Serialization in Agent System	158
7.1	System Startup Sequence	163
7.2	Deal Discovery Process	166
7.3	Ensemble Price Prediction Workflow	167
7.4	SpecialistAgent Detailed Execution	167
7.5	FrontierAgent Detailed RAG Process	169
7.6	Sequential Agent Execution	172
7.7	Parallel Agent Execution (Enhanced Version)	173
7.8	Current Error Handling Issues	174
7.9	Robust Error Handling Implementation	176
7.10	Memory Management in Current System	178
7.11	Enhanced State Management System	178
8.1	Agent Pattern Implementation Summary	188
8.2	Composition Pattern Summary	188
8.3	Strategy Pattern Summary	188
8.4	Template Method Pattern Summary	189
8.5	Factory Pattern Summary	189
8.6	Quick Error Handling Enhancement	192
8.7	Configuration Management Enhancement	192
8.8	Parallel Processing Enhancement	193
8.9	Comprehensive Testing Framework	194
8.10	Plugin System for Agent Extensions	196
8.11	Monitoring and Metrics System	197
B.1	New Agent Template	203
B.2	Configuration System	203
C.1	Agents Module	210
C.2	Deals Module	211
C.3	Ensemble Agent	214
C.4	Messaging Agent	216
C.5	Frontier Agent	218

C.6 Planning Agent	221
C.7 Random Forest Agent	223
C.8 Scanner Agent	224
C.9 Scanner Agent Using Langchain	227
C.10 Specialist Agent	230

Preface

This comprehensive guide provides an in-depth analysis of a sophisticated multi-agent deal discovery system implemented in Python. The system demonstrates advanced software engineering principles, modern Python programming techniques, and practical machine learning integration.

What You'll Learn

This documentation covers:

- **Advanced Python Programming:** Object-oriented design, inheritance patterns, composition, and modern Python idioms
- **System Architecture:** Multi-agent design patterns, component interaction, and scalable system design
- **Machine Learning Integration:** Ensemble methods, LLM APIs, vector databases, and traditional ML models
- **Web Technologies:** RSS parsing, HTML scraping, HTTP clients, and API integration
- **Production Practices:** Error handling, configuration management, logging, and performance optimization

How to Use This Guide

The guide is structured in seven comprehensive parts:

1. **System Overview:** High-level architecture and component relationships
2. **File Analysis:** Detailed examination of each Python file
3. **Class Hierarchy:** Object-oriented design patterns and inheritance
4. **Function Breakdown:** Method-by-method analysis with execution details
5. **Python Concepts:** Advanced language features with practical examples
6. **Agent Interactions:** Workflow analysis and sequence diagrams
7. **Summary & Knowledge:** Complete system recap and learning outcomes

Each part builds upon the previous sections while remaining accessible for focused study of specific topics.

Prerequisites

Readers should have:

- Basic Python programming knowledge
- Understanding of object-oriented programming concepts
- Familiarity with common Python libraries (requests, json, etc.)
- Basic knowledge of machine learning concepts (optional but helpful)

Code Conventions

Throughout this guide:

- Code snippets are syntax-highlighted for readability
- Important concepts are marked with ratings (★☆☆☆☆ to ★★★★★)
- Cross-references link related concepts across chapters
- Diagrams illustrate complex system interactions

Chapter 1

LLM Cheat Sheet

1.1 Handling Imports from Parent Folders in Python and Jupyter

When a Python script or Jupyter Notebook is located in a subfolder, modules in a parent folder may not be automatically discoverable by Python. The following methods illustrate how to properly import them.

1.1.1 Using `__file__` in Python Scripts

For scripts executed directly, `__file__` contains the path of the script. Use it to construct the parent folder path:

Python Script: Using `__file__`

```
import sys
import os

# Add parent folder to Python path
sys.path.append(os.path.abspath(os.path.join(os.path.dirname(__file__), '..')))

# Import module from parent folder
from LLMengineering import key_utils
```

Explanation:

- `os.path.dirname(__file__)`: folder of the current script.
- `os.path.join(..., '..')`: go one level up.
- `os.path.abspath(...)`: get absolute path.
- `sys.path.append(...)`: add this path to Python's module search paths.

Using `os.getcwd()` in Jupyter Notebooks

In Jupyter, `__file__` is not defined. Use the current working directory instead:

Jupyter Notebook: Using os.getcwd()

```
import sys
import os

# Add parent folder of current working directory
sys.path.append(os.path.abspath(os.path.join(os.getcwd(), '..')))

# Import module from parent folder
from key_utils import some_function
```

Notes:

- `os.getcwd()`: returns the current working directory of the notebook.
- `..` moves one level up; `../..` moves two levels up.
- This approach is portable for notebooks without relying on `__file__`.

Using Absolute Paths (Less Portable)

Absolute Path Approach

```
import sys
sys.path.append("/run/media/sijanpaudel/New Volume/New folder/LLMengineering")

from key_utils import some_function
```

Caution: This works immediately but is hard-coded. If the project moves, the path must be updated.

Recommended Long-Term Approach

1. Add an empty `__init__.py` file in the parent folder (e.g., `LLMengineering/`) to make it a package.
2. Use relative or absolute imports within scripts:

Long-Term Package Approach

```
from LLMengineering import key_utils
```

This makes imports clean, portable, and avoids repeated `sys.path.append` hacks.

1.2 Comparing LLMs – The Basics (1)

Feature	Description / Notes
Open-source or Closed	Whether the model's weights, architecture, and training code are publicly available. Examples: GPT-4 (closed), LLaMA 2 (open).
Release Date & Knowledge Cut-off	The date the model was released and the latest point in time the training data covers. Important for up-to-date responses.
Parameters	Number of trainable weights in the model (e.g., billions of parameters). Determines model capacity.
Training Tokens	Amount of text data used during training, usually in billions of tokens. More tokens generally improve performance.
Context Length	Maximum input length (in tokens) the model can handle in a single prompt. Longer context allows better understanding of large inputs.

1.3 Comparing LLMs – The Basics (2)

Feature	Description / Notes
Inference Cost	Computational resources required to generate output (GPU, CPU usage).
API Charge / Subscription / Runtime Compute	How much it costs to access the model through a cloud service.
Training Cost	Cost to pretrain the model, including compute, electricity, and infrastructure.
Build Cost	Cost to fine-tune, deploy, or customize the model.
Time to Market	How quickly you can deploy and use the model in production.
Rate Limits	Restrictions on API usage, e.g., calls per minute or per day.
Speed & Latency	How fast the model responds, depends on model size, hardware, and context length.
License	Terms for use, redistribution, and commercial deployment. Some open-source licenses restrict commercial use.

1.4 Chinchilla Scaling Law

Key Principle: Number of model parameters should scale roughly proportional to the number of training tokens.

Implications:

- Increasing model size without enough data leads to plateaued or degraded performance.
- Having lots of training data with a small model underperforms; parameters must scale up.

Rule of Thumb: Let N = model parameters, D = training tokens. For optimal training:

$$N \propto D$$

Example:

- Doubling model parameters → need roughly double the training tokens.
- Insufficient tokens → diminishing returns.

Additional Notes:

- Smaller, well-trained models can outperform larger, undertrained ones.
- Helps determine compute-efficient configurations for new LLMs.

1.5 Common LLM Benchmarks

Benchmark	Focus / Task	Description
MMLU	Knowledge across multiple subjects	Assess general knowledge and reasoning; used to compare models like GPT-3, Chinchilla, and Gopher.
BIG-bench	Broad suite of diverse reasoning tasks	Tests reasoning, factual knowledge, ethics, math, code; hundreds of tasks evaluating beyond narrow QA.
HellaSwag	Commonsense reasoning	Multiple-choice questions for everyday situations; measures ability to predict plausible outcomes.
TruthfulQA	Factual accuracy / truthfulness	QA tasks designed to detect hallucinations; evaluates honesty of LLM answers.
WinoGrande / Winograd Schema Challenge	Pronoun resolution / coreference	Tests commonsense reasoning and context understanding; resolves ambiguous references.
ARC	Science and reasoning	Multiple-choice science questions; evaluates problem-solving and reasoning in STEM.
HumanEval	Coding and code generation	Tests Python programming ability; measures functional correctness of generated code.

1.6 Specific Benchmarks

Benchmark	What's Being Evaluated	Description
ELO	Model ranking / performance consistency	Evaluates models via pairwise comparisons; creates a relative ranking of LLMs across multiple tasks.
HumanEval	Code generation / functional correctness	Tests an LLM's ability to write Python functions that pass unit tests; measures coding logic and correctness.
Multipl-E	Multimodal reasoning	Evaluates LLMs on tasks combining text and images (or multiple modalities); measures reasoning and comprehension across modalities.

1.7 Limitations of Benchmarks

Limitation	Explanation
Narrow focus	Many benchmarks test only specific skills (e.g., coding, factual QA, commonsense), not overall intelligence or adaptability.
Static datasets	Benchmarks are fixed in time, so models trained after the cut-off may have an unfair advantage or miss newer knowledge.
Lack of real-world context	Benchmarks often use idealized tasks, not messy, ambiguous, or multi-step real-world scenarios.
Gaming / overfitting	Models can be fine-tuned or prompted to specifically excel on benchmark tasks without improving general capabilities.
Limited multimodality	Most benchmarks focus on text-only tasks; few measure image, audio, or multimodal reasoning.
Subjectivity	Some benchmarks (e.g., ethics, creativity, hallucination detection) are hard to score objectively.
Compute bias	Larger models may perform better mainly due to size, not reasoning ability, skewing benchmark results.
Hard to measure nuanced reasoning	Benchmarks mostly measure surface correctness; they often cannot capture multi-step reasoning, context-dependent judgment, creativity, or reasoning accuracy vs. fluency.

1.8 Advanced Benchmarks for Large Language Models

Benchmark	Focus / Task	Description / Meaning
GPQA	Graduate-level question answering	Evaluates performance on graduate-level tests with 448 expert questions. Non-PhD humans score only 34% even with web access. Measures LLM ability to handle highly specialized knowledge.
BBHard	Future capabilities	Includes 204 tasks previously thought beyond LLM capabilities. Designed to test reasoning, logic, and generalization at a next-level difficulty.
Math Lv 5	High-school math competition problems	Measures model's ability to solve advanced math problems requiring multi-step reasoning and problem-solving skills. Useful for chain-of-thought evaluation.
IFEval	Instruction following	Tests the model's ability to follow complex instructions, e.g., "write more than 400 words" and "mention AI at least 3 times." Evaluates comprehension and compliance with nuanced prompts.
MuSR	Multistep soft reasoning	Assesses logical deduction and multi-step reasoning. Example: analyzing a 1,000-word story and identifying "who has means, motive, and opportunity." Tests reasoning beyond surface facts.
MMLU-PRO	Harder MMLU version	Advanced, cleaned-up version of MMLU with questions having 10 possible answers instead of 4. Evaluates deeper knowledge, multi-choice reasoning, and generalization.

Notes:

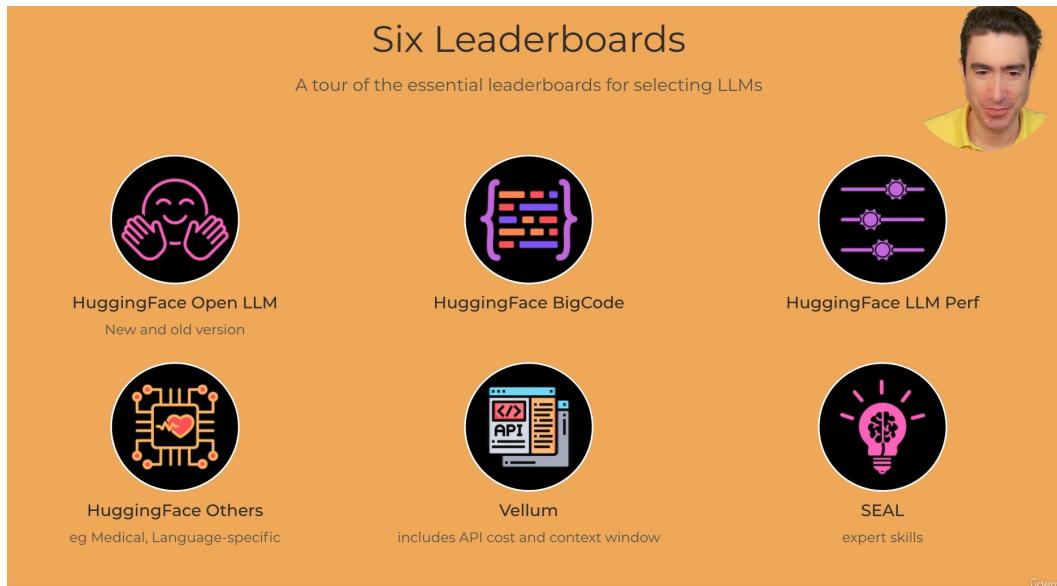
- These benchmarks are considered “next-level” because they test advanced reasoning, multi-step logic, and high-level knowledge.
- Useful for evaluating models that go beyond standard QA, commonsense reasoning, or basic coding tasks.
- Many of these benchmarks also measure compliance with complex instructions and reasoning under uncertainty.

Benchmark	Purpose / Examples
GPQA	<p>Specialized knowledge, expert-level Q&A, academic reasoning, physics, chemistry, biology expertise.</p> <p>Examples: Explain the biochemical steps of glycolysis; Describe Newton's laws with practical applications; Explain photosynthesis in detail.</p>
BBHard	<p>Advanced reasoning, logic, generalization in challenging scenarios.</p> <p>Examples: Predict how a hypothetical AI agent could optimize a supply chain in a novel scenario; Design a strategy to minimize energy consumption in an unfamiliar industrial process; Solve a logic puzzle with multiple constraints.</p>
Math Lv 5	<p>Complex problem solving, chain-of-thought evaluation, mathematical reasoning.</p> <p>Examples: Solve: If $x^2 - 5x + 6 = 0$, find all integer solutions; Evaluate combinatorial problems like "How many ways can 5 books be arranged on a shelf?"; Solve calculus problems requiring multiple steps.</p>
IFEval	<p>Precise instruction compliance, comprehension, content generation with constraints.</p> <p>Examples: Write a 450-word essay on renewable energy mentioning AI at least 3 times; Summarize a research paper in 300 words including key terms; Rewrite a paragraph to follow a formal tone while keeping original meaning.</p>
MuSR	<p>Multi-step reasoning, deduction, understanding narratives beyond surface facts, solving prime puzzles and reasoning tasks.</p> <p>Examples: Analyze a 1,000-word mystery story to determine "who has means, motive, and opportunity"; Identify the next prime number in a complex sequence; Solve multi-step logic puzzles.</p>
MMLU-PRO	<p>Deep knowledge, advanced multi-choice understanding, broader knowledge evaluation, general language understanding.</p> <p>Examples: Choose the best explanation for why the sky is blue from 10 options; Identify the correct historical fact among 10 alternatives; Evaluate grammar and style in multiple-choice questions.</p>

Tips to distinguish benchmarks:

- **GPQA** – focus on subject expertise.
- **BBHard** – focus on novel reasoning and logic.
- **IFEval** – focus on following instructions accurately.
- **MuSR** – focus on multi-step deduction and problem solving (e.g., prime puzzles, mysteries).
- **MMLU-PRO** – focus on broad knowledge and multiple-choice reasoning.

1.9 Leaderboard Image



Leaderboard	Purpose / Use
HuggingFace Open LLM	Tracks performance of open-source LLMs (new and old versions). Useful for selecting models for research or deployment with full weight access.
HuggingFace BigCode	Focused on code generation LLMs. Helps compare models for coding tasks, programming language coverage, and code quality.
HuggingFace LLM Perf	Benchmarking general LLM performance across various NLP tasks. Useful for measuring speed, accuracy, RAM memory usage and general reasoning ability.
HuggingFace Others	Specialized leaderboards, e.g., medical, domain-specific, or language-specific models. Guides selection for niche applications.
Vellum	Includes context window, API cost, and usage constraints. Useful for developers to choose models based on practical deployment factors.
SEAL	Focused on expert skills and high-level reasoning. Useful for selecting LLMs that excel in professional or complex knowledge domains.

1.10 A 5-Step Strategy for Using AI Models

This is a simple guide to choosing, training, and using a Large Language Model (LLM) to solve a real-world business problem.



1.10.1 Understand (Figure Out What You Need)

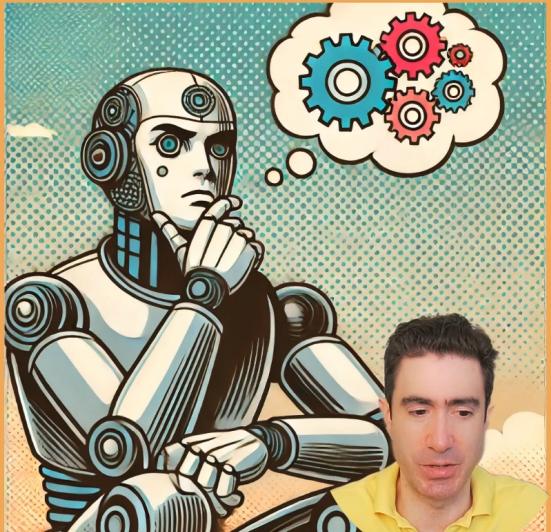
Before you start, you need a clear plan. Think about what you're trying to achieve.

- **Gather Business Needs:** What is the specific problem you want the AI to solve? (e.g., "Answer customer support questions automatically.")
- **Define Success:** How will you know if the AI is doing a good job? Focus on what matters to the business, not just technical scores. (e.g., "Reduce customer wait times by 50%")
- **Look at Your Data:** What information do you have? Is it a lot? Is it messy or clean? What format is it in?
- **Consider the Limits:** Think about the practical stuff.
 - How much money can you spend?
 - Does it need to be super fast?
 - How many people will be using it?
 - What is the deadline to get this working?

1. Understand

Activities:

- Gather business requirements for the task
- Identify performance criteria
Particularly the Business Centric metrics
- Understand the data: quantity, quality, format
- Determine non-functionals
Cost constraints, scalability, latency
R&D / build budget and implementation timeline



The illustration features a stylized robot with a metallic, blue and silver body. It has large, expressive blue eyes and a thoughtful expression, with its hand resting near its chin. Above the robot's head is a thought bubble containing three interlocking gears in blue, red, and yellow. To the right of the robot, a portion of a man's face is visible, showing a neutral or slightly concerned expression. The background is a textured orange color.

1.10.2 Prepare (Do Your Homework)

Once you have a plan, it's time to get your resources ready.

- **Research Solutions:** See how this problem is already being solved. Look at solutions that don't use a big AI model to set a "baseline" for performance.
- **Compare AI Models:** Look at different LLMs.
 - **The Basics:** How much do they cost? How much text can they handle at once (context length)? Can you legally use them for your business?
 - **Performance:** Check online tests (Benchmarks, Leaderboards) to see which models are best at certain tasks.
- **Clean Your Data:** This is very important! You need to organize, clean, and prepare your data so the AI can learn from it effectively.

2. Prepare



Activities:

- Research existing / non-LLM solutions
 - Potential baseline model
- Compare relevant LLMs
 - The basics, including context length, price and license
 - Benchmarks, Leaderboards and Arenas
 - Specialist scores for the task at hand
- Curate data: clean, preprocess and split

Udemy

1.10.3 Select (Pick Your Tools and Test Them)

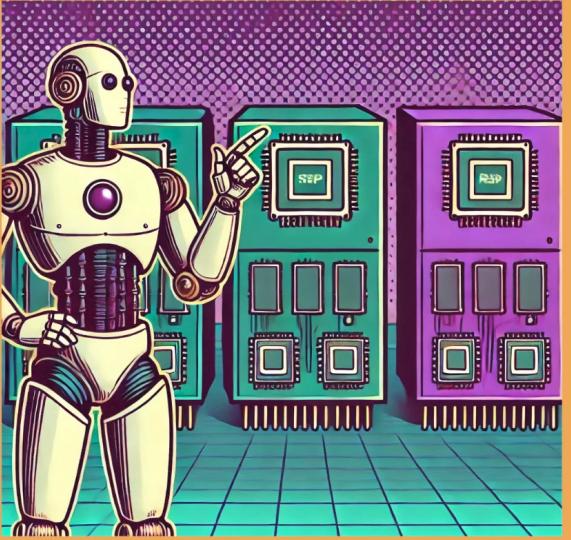
Now you can choose your model and see how it performs with your data.

- **Choose Your LLM(s):** Based on your research, pick one or two models that seem like the best fit.
- **Experiment:** Play around with the models. See what they're good at and where they struggle.
- **Train and Check:** Use your clean data to test (validate) the models to get real performance numbers.

3. Select

Activities:

- Choose LLM(s)
- Experiment
- Train and validate with curated data



1.10.4 Customize (Make the AI Better at Your Task)

Out-of-the-box models are good, but you can make them great by using special techniques to improve their performance on your specific problem. There are three main ways to do this.

Three Key Techniques

Prompting

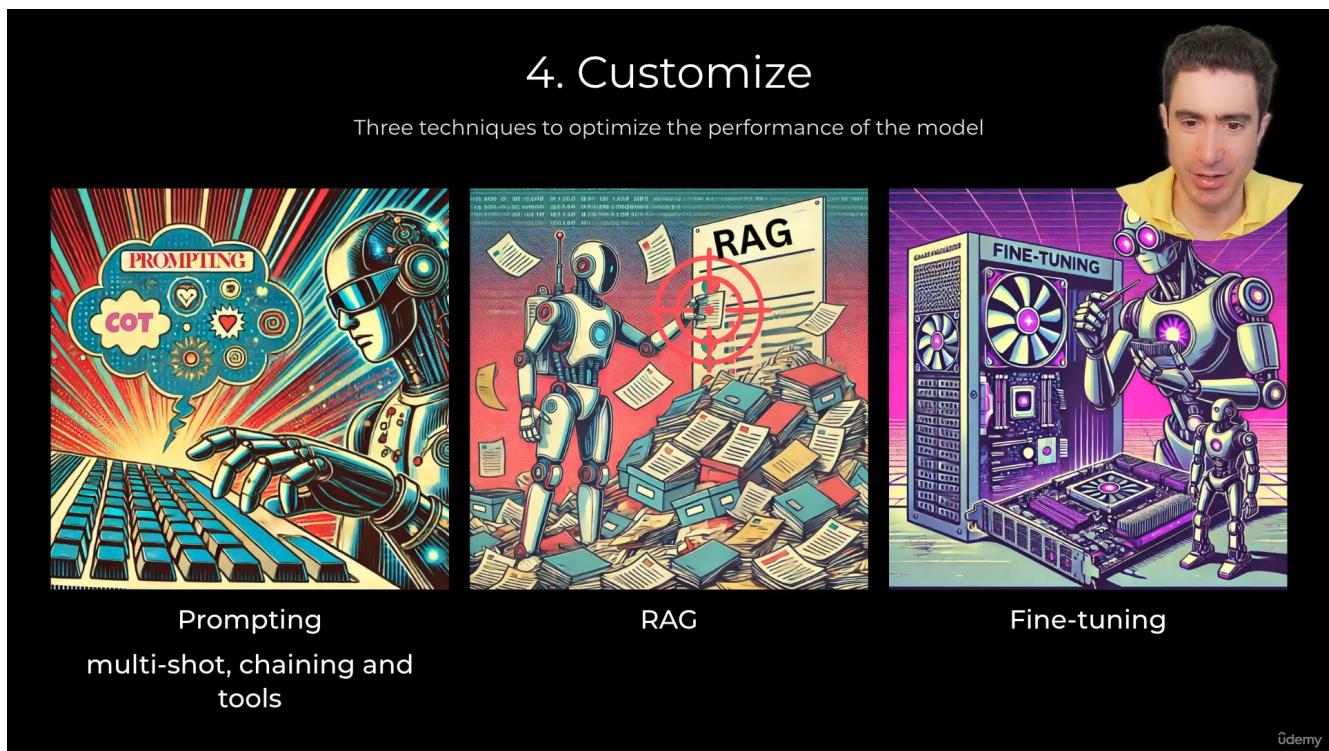
This means giving the AI better instructions. It's the fastest and cheapest way to improve performance. You can give it examples, tell it to "think step-by-step," or connect it to other software tools.

RAG (Retrieval-Augmented Generation)

This is like giving the AI an "open book" for a test. You connect the model to a knowledge base (like your company's internal documents). When a question is asked, the AI first finds the relevant information from your documents and then uses that information to create the answer. This ensures the answers are accurate and up-to-date.

Fine-Tuning

This is the most advanced technique. You actually re-train a part of the AI model on your own data. This teaches the model a new skill or a specific style, making it an expert in your particular area. It's expensive and requires a lot of data but can lead to the best performance.



Comparing the Three Techniques

Pros: The Good Stuff		
Prompting	RAG	Fine-tuning
<ul style="list-style-type: none"> 1. Quick to do 2. Cheap 3. See improvements right away 	<ul style="list-style-type: none"> 1. More accurate answers 2. Works with lots of data 3. Efficient 	<ul style="list-style-type: none"> 1. Creates a true expert 2. Understands subtle details 3. Can learn a specific tone or style 4. Faster and cheaper answers

Cons: The Downsides		
Prompting	RAG	Fine-tuning
<ul style="list-style-type: none"> 1. Limited by how much text the AI can read at once 2. Improvements eventually stop 3. Can be slow and expensive for each answer 	<ul style="list-style-type: none"> 1. More complex to set up 2. Needs accurate, up-to-date info 3. Can miss subtle details 	<ul style="list-style-type: none"> 1. Takes a lot of effort 2. Needs a lot of clean data 3. Costs money to train 4. Risk of the AI forgetting its original general knowledge

Three Techniques: Pros



Prompting

- 1. Fast to implement
- 2. Low cost
- 3. Often immediate improvement

RAG

- 1. Accuracy improvement with low data needs
- 2. Scalable
- 3. Efficient

Fine-tuning

- 1. Deep expertise & specialist knowledge
- 2. Nuance
- 3. Learn a different tone / style
- 4. Faster and cheaper inference

udemy

Three Techniques: Cons



Prompting

- 1. Limited by context length
- 2. Diminishing returns
- 3. Slower, more expensive inference

RAG

- 1. Harder to implement
- 2. Requires up-to-date, accurate data
- 3. Lacks nuance

Fine-tuning

- 1. Significant effort to implement
- 2. High data needs
- 3. Training cost
- 4. Risk of "catastrophic forgetting"

udemy



1.10.5 Productionize (Get It Ready for the Real World)

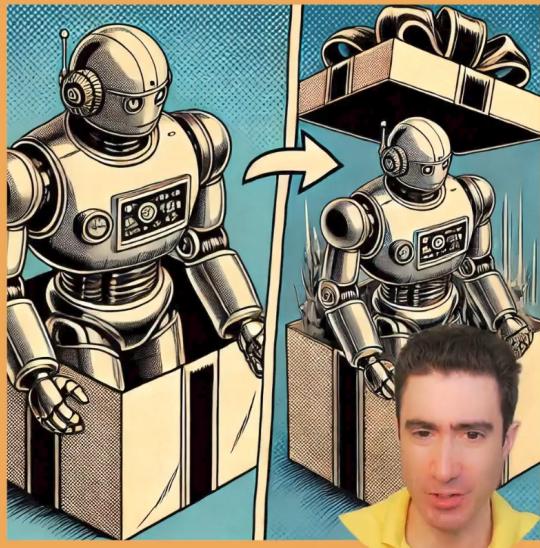
Once the model is working well, you need to make it a reliable part of your business.

- **Connect it:** How will your website or other software "talk" to the AI model?
- **Deploy it:** Where will the model run? On your own computers or in the cloud?
- **Plan for Growth:** What happens if you get thousands of users? Make sure it's secure, and you're monitoring it for problems.
- **Measure What Matters:** Go back to Step 1 and measure the business goals you set. Is it actually working?
- **Keep it Updated:** Continuously check the AI's performance and re-train it with new data to keep it from becoming outdated.

5. Productionize

Activities:

- Determine API between model and platform(s)
- Identify model hosting and deployment architecture
- Address scaling, monitoring, security and compliance
- Measure the Business-Focused Metrics identified in step 1
- Continuously retrain and measure performance



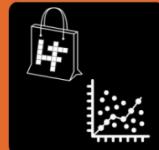
Odemy

Traditional ML models

We will quickly set up these solutions to give us a starting point



Feature engineering & Linear Regression



Bag of Words & Linear Regression



word2vec & Linear Regression



word2vec & Random Forest



word2vec & SVR



1.11 LoRA (Low-Rank Adaptation)

1.11.1 Introduction

LoRA (Low-Rank Adaptation) is a technique designed to **fine-tune large pre-trained models efficiently**. Large language models (LLMs) like GPT, LLaMA, etc., have **billions of parameters**, making traditional fine-tuning:

- Extremely expensive in computation.
- Memory-intensive.
- Hard to maintain multiple fine-tuned versions.

LoRA solves this by training **small, additional modules** instead of updating the entire model.

How LoRA Works (Detailed)

Consider a weight matrix in a neural network layer:

$$W \in \mathbb{R}^{d \times k}$$

Full fine-tuning adjusts every element in W , which is costly.

LoRA introduces **two low-rank matrices** A and B and approximates the weight update:

$$\Delta W \approx AB$$

where:

$$A \in \mathbb{R}^{d \times r}, \quad B \in \mathbb{R}^{r \times k}, \quad r \ll \min(d, k)$$

Key points:

- r is small, so A and B are **tiny compared to W** .
- Only A and B are trained; the original W is **frozen**.
- Memory and computation are drastically reduced.

Step-by-Step Intuition

1. Start with a pre-trained model whose parameters W are **frozen**.
2. Add small trainable matrices A and B to selected layers (often attention layers).
3. Forward pass: compute original output using W and add low-rank contribution:

$$y = Wx + \alpha(ABx)$$

where α is a scaling factor.

4. Backpropagation only updates A and B , leaving W untouched.
5. After training, save only A and B as the **LoRA adapter**.

Simple Analogy

Analogy

Think of a professional camera with many settings:

- Full fine-tuning = manually adjusting every knob and lens component each time you want a new photo style.
- LoRA = keep the camera fixed, but attach a **small filter**.

The filter is light, cheap, and can be swapped out, while the main camera stays unchanged.

Why LoRA is Effective

- Neural network updates often lie in a **low-rank subspace**, so approximating with small matrices captures most of the needed changes.
- Reduces training cost dramatically.
- Reduces storage: instead of storing full fine-tuned model (10s of GB), store just adapters (10s of MB).
- Enables **multi-tasking**: use same base model with different adapters for different tasks.

Benefits vs Limitations

Pros

- **Memory-efficient** and fast training.
- Base model remains **unchanged**, enabling multiple adapters.
- Easy integration with pre-trained models.

Cons

- Only works well when **low-rank approximation is sufficient**.
- Might not capture highly complex changes needed for very different tasks.

Summary in Simple Words

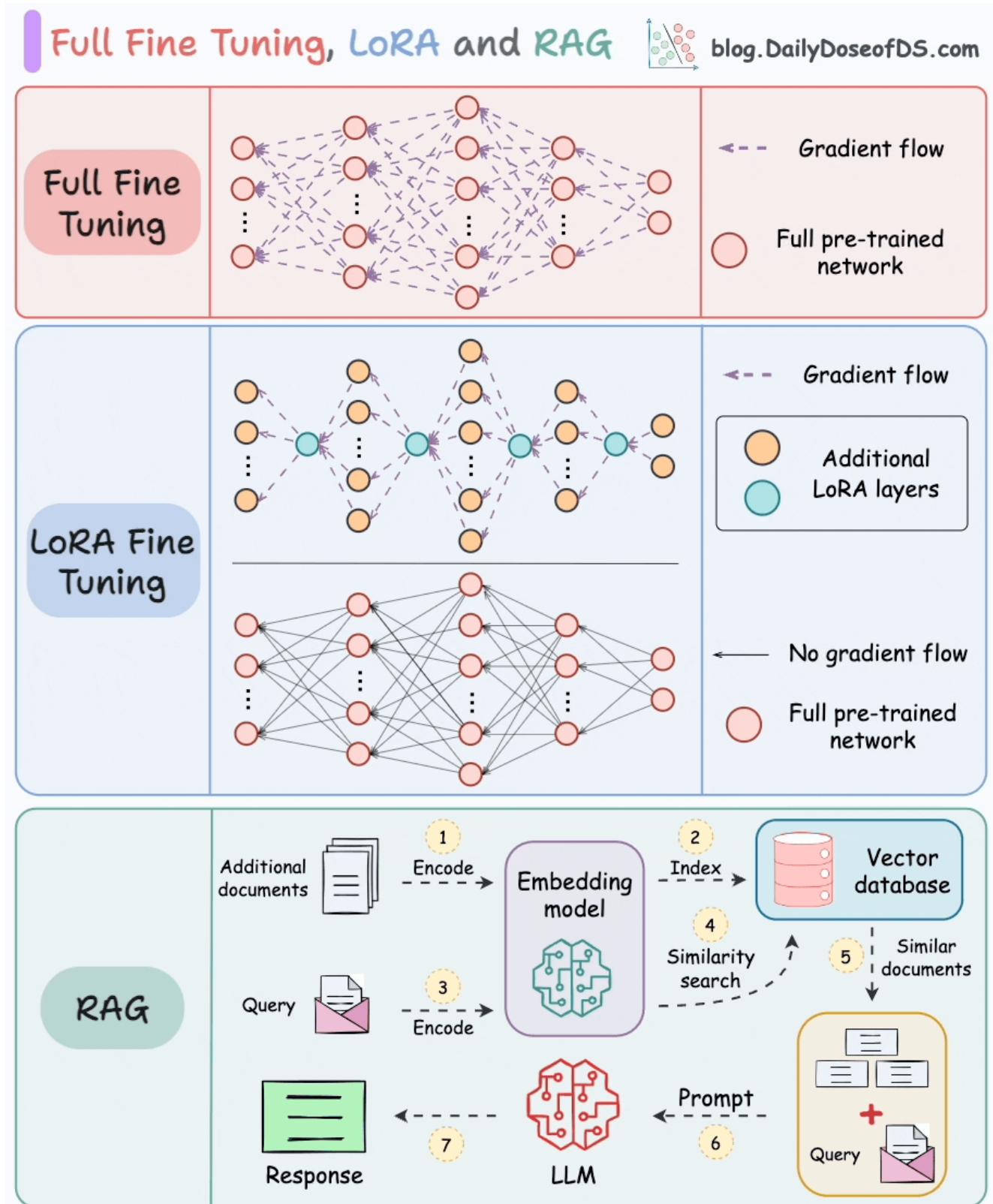
LoRA Simplified

You can think of LoRA as adding **tiny adapters on top of a frozen large model**.

Imagine a giant model with billions of knobs. Fine-tuning everything is slow and expensive. LoRA adds a few small sliders (matrices) on top of the frozen model. During training, only these sliders move, but they are enough to guide the model to a new task.

Visualization Idea

You can include a visual showing LoRA adapters on top of a frozen model: [Link](#)



1.12 LoRA Hyperparameters: r , α , and Target Modules

LoRA (Low-Rank Adaptation) introduces a few crucial hyperparameters that control the learning capacity, scaling, and location of the adaptation. Understanding them is key to effective fine-tuning.

1.12.1 Rank r (Low-Rank Dimension)

The hyperparameter r determines the **rank of the low-rank matrices** $A \in \mathbb{R}^{d \times r}$ and $B \in \mathbb{R}^{r \times k}$.

- Controls the number of parameters introduced in the LoRA adapter.
- Small $r \Rightarrow$ fewer parameters, less memory, faster training, but less expressive.
- Large $r \Rightarrow$ more parameters, higher capacity, but increased memory and computation.

Rule of Thumb:

Start with $r \in [4, 64]$ for LLMs. Increase r if the task is complex or performance is insufficient.

Scaling Factor α

The hyperparameter α scales the output of the LoRA adapter:

$$y = Wx + \frac{\alpha}{r}(ABx)$$

- Controls the **strength of the low-rank update** relative to the original weights.
- Small $\alpha \Rightarrow$ subtle influence, safer but slower adaptation.
- Large $\alpha \Rightarrow$ strong influence, faster adaptation, but may destabilize training.

Rule of Thumb:

Use $\alpha \in [8, 32]$ for LLMs. Adjust based on task difficulty and stability.

Target Modules

Target modules define **which layers of the base model receive LoRA adapters**.

- Common targets: attention layers (query/key/value), feed-forward layers.
- Choosing too many modules \Rightarrow more parameters, higher memory.
- Choosing too few modules \Rightarrow may underfit the task.

Rule of Thumb:

Start by adding LoRA to attention projection matrices (query and value) for LLMs. Expand to feed-forward layers only if necessary.

1.12.2 Summary Table

LoRA Hyperparameters Summary		
Hyperparameter	Role	Rule of Thumb
r	Low-rank dimension (number of adapter parameters)	4–64
α	Scaling factor of LoRA update	8–32
Target Modules	Layers where adapters are applied	Start with attention layers and then expand to feed-forward layers only if necessary.

Three Essential Hyperparameters
For LoRA Fine-Tuning




r
 The rank, or how many dimensions in the low-rank matrices


Alpha
 A scaling factor that multiplies the lower rank matrices


Target Modules
 Which layers of the neural network are adapted

RULE OF THUMB:
 Start with 8, then double to 16, then 32, until diminishing returns

RULE OF THUMB:
 Twice the value of r

RULE OF THUMB:
 Target the attention head layers

udemy

1.13 Q, K, V, O in Transformers

In a transformer model, the attention mechanism allows each token to selectively focus on other tokens in the sequence. This is done using four key components: **Query (Q)**, **Key (K)**, **Value (V)**, and **Output (O)**.

- **Query (Q):** Represents the **questions** each token asks. For example, a token may ask, “Which other tokens are relevant to me?” Q is calculated by multiplying the token embedding with a learned weight matrix W_Q .
- **Key (K):** Represents the **tags or identifiers** of each token. It helps the Query figure out which tokens are relevant. K is computed by multiplying token embeddings with another learned weight matrix W_K .
- **Value (V):** Represents the **actual content or information** of each token. Once Q identifies which tokens are relevant (by comparing with K), V provides the information to pass forward. V is calculated using the weight matrix W_V .
- **Output (O):** After attention combines relevant V vectors, O is a projection that mixes the attended information back into the model’s hidden dimension using the matrix W_O . This ensures the output can be passed to the next layer seamlessly.

The attention operation is mathematically defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

- QK^\top computes the similarity between queries and keys.
- d_k is the **dimension of the key vectors per attention head**. Dividing by $\sqrt{d_k}$ scales the dot product so that the softmax does not produce extremely small or large probabilities.
- Multiplying by V gathers the relevant information based on these attention weights.

1.13.1 Analogy: Classroom Example

Analogy

Imagine a classroom scenario where students ask and answer questions:

- **Q (Query):** A student asks a question to figure out which other students have relevant information.
- **K (Key):** Each student has a topic card that indicates what they know.
- **V (Value):** Each student provides an answer (information) if their topic matches the query.
- **O (Output):** The teacher collects all answers, summarizes them, and writes a combined response on the board.
- d_k : Think of it as the “attention focus resolution” — how detailed each student’s topic card is. More details (larger d_k) require scaling to prevent overemphasis on any single answer.

This shows how Q identifies what to focus on, K helps determine relevance, V provides the information, and O produces the final output.

1.14 LoRA in PeftModelForCausalLM

1.14.1 Model Architecture Overview

The `PeftModelForCausalLM` wraps a `LlamaForCausalLM` base model and applies **LoRA adapters** for efficient fine-tuning.

- **Embedding Layer:** `embed_tokens` maps vocabulary tokens (128256) to embeddings of size 4096.
- **Decoder Layers:** 32 `LlamaDecoderLayer` modules, each containing:
 - **Self-Attention (`self_attn`)**
 - * Q, K, V, O projections use `lora.Linear4bit` modules.
 - * Each projection has:
 - `base_layer`: The frozen 4-bit linear layer representing the main weight matrix.
 - `lora_A`, `lora_B`: Low-rank matrices (rank $r = 32$) representing trainable updates.
 - `lora_dropout`: Dropout for regularization.
 - * LoRA modifies only the projection weights, not the base model.
 - **Feed-Forward Network (MLP):** `LlamaMLP` with:
 - * `gate_proj` and `up_proj` (expand embeddings to 14336)
 - * `down_proj` (reduce back to 4096)
 - * Non-linearity: `SiLU()`
 - **LayerNorms:** `input_layernorm` and `post_attention_layernorm`.
 - **Rotary Positional Embeddings:** `rotary_emb` for relative position encoding.
- **Final LayerNorm:** `norm` over the hidden dimension.
- **LM Head:** Linear mapping from hidden size 4096 to vocabulary size 128256.

1.14.2 LoRA in Attention Layers

In transformers, attention is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

where Q, K, V are projections of the input embeddings. LoRA modifies these projections as follows:

$$W_{\text{proj}} \rightarrow W_{\text{proj}} + \Delta W = W_{\text{proj}} + AB$$

- W_{proj} is the frozen base projection weight (Linear4bit).
- $A \in \mathbb{R}^{d_{\text{model}} \times r}$, $B \in \mathbb{R}^{r \times d_{\text{proj}}}$, with $r \ll d_{\text{model}}$ (e.g., $r = 32$).
- Only A and B are trainable, drastically reducing the number of parameters.

- LoRA is applied individually to **Q, K, V, O** matrices:
 - **q_proj**: 4096 → 4096
 - **k_proj**: 4096 → 1024
 - **v_proj**: 4096 → 1024
 - **o_proj**: 4096 → 4096

Why LoRA is Effective in Attention

- **Memory Efficiency:** Instead of updating the full weight matrices (4096×4096), only two small matrices of size 4096×32 and 32×4096 are trained per projection.
- **Task Adaptation:** LoRA can inject task-specific knowledge into the attention mechanism without modifying the base model.
- **Multiple Adapters:** Different LoRA adapters can be swapped for different tasks while keeping the same base model frozen.
- **Attention-Focused Learning:** By applying LoRA specifically to Q, K, V, O, the model learns to adjust how it attends to other tokens, which is sufficient for most fine-tuning tasks.

1.14.3 Dimensions of Q, K, V, O Projections

Attention Projection Dimensions

Projection	Input Features	Output Features
Q (q_proj)	4096	4096
K (k_proj)	4096	1024
V (v_proj)	4096	1024
O (o_proj)	4096	4096

1.14.4 Summary

- **LoRA** allows fine-tuning of very large models with minimal additional parameters.
- Applied to attention layers, LoRA adjusts how the model attends to input tokens.
- Memory-efficient and task-adaptive: only the low-rank matrices are trained, while the main model stays frozen.
- Facilitates multiple task-specific adapters without retraining the entire model.

Key Points

- d_{model} is the hidden dimension (4096 in this model) through which all tokens are projected.

- LoRA adds small trainable matrices (rank $r = 32$) to Q, K, V, O projections.
- Base projections remain frozen; only adapters are updated.
- Memory-efficient and task-adaptive fine-tuning.

1.15 Loading a Base Model in 4-Bit and Using LoRA

The following code demonstrates how to load a large language model (LLM) in 4-bit precision for memory efficiency and later fine-tune it with LoRA.

```
quant_config = BitsAndBytesConfig(  
    load_in_4bit=True,  
    bnb_4bit_use_double_quant=True,  
    bnb_4bit_compute_dtype=torch.bfloat16,  
    bnb_4bit_quant_type="nf4")  
  
base_model = AutoModelForCausalLM.from_pretrained(  
    BASE_MODEL,  
    quantization_config=quant_config,  
    device_map="auto",  
)
```

1.15.1 Explanation of the Code

- **BitsAndBytesConfig:** Configuration for quantization. Quantization reduces the precision of the weights to save memory.
 - `load_in_4bit=True`: Load model weights in 4-bit precision instead of 16-bit or 32-bit, drastically reducing memory usage.
 - `bnb_4bit_use_double_quant=True`: Uses double quantization to reduce numerical errors and improve accuracy.
 - `bnb_4bit_compute_dtype=torch.bfloat16`: Performs computations in bfloat16 to maintain performance while saving memory.
 - `bnb_4bit_quant_type="nf4"`: Specifies the 4-bit quantization type.
- **AutoModelForCausalLM.from_pretrained:** Loads the pre-trained LLM with the above quantization configuration.
 - `BASE_MODEL`: The pre-trained model identifier (e.g., LLaMA, Falcon).
 - `quantization_config=quant_config`: Uses the 4-bit quantized weights for memory efficiency.
 - `device_map="auto"`: Automatically maps model layers to available GPU(s) or CPU.

1.15.2 How LoRA Optimizes This

- Loading the model in 4-bit reduces the memory footprint, but full fine-tuning is still expensive if we try to update all weights.
- LoRA introduces **small trainable adapters** (low-rank matrices A and B) on top of the frozen base model.
- Only these adapters are updated during training, allowing:
 - **Memory efficiency**: base model weights are frozen and quantized.

- **Faster training:** fewer parameters are trained.
- **Task adaptation:** the model can specialize on a new task without full fine-tuning.
- Effectively, combining 4-bit quantization with LoRA allows training huge models on limited GPU memory.

1.15.3 Summary

Takeaway

- **4-bit quantization:** reduces memory usage of base model weights.
- **LoRA adapters:** trainable low-rank matrices that update model behavior without touching the large frozen base model.
- **Combined benefit:** efficient fine-tuning on large models with low GPU memory and compute cost.

1.16 Decision: Selecting a Base Model

1.16.1 Goal

Choose a base model by balancing: **task requirements**, **inference cost** (latency / GPU memory), **accuracy needs**, license / ecosystem, and **how you plan to fine-tune** (LoRA, RAG, full FT).

1.16.2 Decision criteria (short checklist)

- **Capability vs cost:** larger models usually perform better but cost more to run and host.
- **Parameter count:** affects memory, quantization decisions and latency.
- **Context window / multimodality:** needed for long-context tasks or image+text.
- **License and ecosystem:** open weights, available tooling (HF, vendor SDKs), community adapters.
- **Fine-tuning strategy:** if you plan to use LoRA, favor models with stable adapter support and documented PEFT integrations.
- **Deployment target:** edge / mobile vs single-GPU vs multi-GPU / cloud inference.

1.16.3 Model-by-model short summary

- **LLaMA family (Meta):** a family of open/available models offered in multiple sizes (smaller local-friendly models and very large variants for research). LLaMA models are widely used for research and community tooling; instruction-tuned variants (“instruct”) are offered for easier instruction-following.
- **Qwen (Alibaba):** a rapidly-evolving family with strong claims in multilingual and code tasks; recent releases include much larger variants (cutting-edge, often multimodal and optimized for large-scale deployments).
- **Phi (Microsoft / Phi family):** smaller, efficient models (e.g., Phi-2 at a few billion parameters) that are optimized for reasoning and cost-effective deployments at smaller parameter counts.
- **Gemma (Google / DeepMind family):** a family of lightweight-to-mid-size models intended for efficient, practical use (good on-device and cloud tradeoffs); models come in small-to-medium sizes optimized for production usage.

1.16.4 Comparison table (high-level)

Model comparison at a glance			
Model	Representative sizes	Typical strengths	Suggested use-cases / notes
LLaMA (Meta)	8B, 70B, (and larger research variants)	Strong general-purpose language and reasoning; broad community tooling	Good when you want high-quality open models and lots of community adapters; choose size by budget
Qwen (Alibaba)	Ranges from mid-size up to very large (vendor evolving family, latest pushes into extremely large sizes)	Good for code, multi-lingual tasks and large-scale deployments	Consider for enterprise cloud deployments or when vendor benchmarks favour Qwen for your task
Phi (Microsoft)	Small/medium (e.g., 2.7B for Phi-2)	Very cost-efficient; strong small-model reasoning per-cost	Excellent for edge or low-cost cloud inference where medium-sized models suffice
Gemma (Google / Deep-Mind)	Small → mid (family examples: 270M, 1B, 4B, 12B, 27B)	Efficient, production-ready; good on-device/cloud tradeoffs	Pick Gemma sizes for production where latency / cost matters and you require Google tooling

1.16.5 Base vs Instruct variants

- **Base (pretrained) model:** weights after pretraining on large corpora. Flexible — you can fine-tune, apply LoRA, or use RAG. Typically **better when you plan to adapt heavily** (custom fine-tuning / domain adaptation) because you control the adaptation pipeline.
- **Instruct (instruction-tuned) model:** a base model further fine-tuned on instruction/assistant datasets (and often RLHF/SFT). It usually performs better out-of-the-box on instruction-following tasks and reduces alignment/failure cases. Use an instruct variant if your primary need is production-ready instruction-following and you want fewer alignment steps.
- **LoRA note:** LoRA works on both base and instruct variants. If you want to minimize work, pick an instruct variant and add LoRA adapters for task-specific behavior; if you require maximal control, start from base and LoRA-fine-tune to your data.

1.16.6 Rules of thumb and recommended workflows

Quick rules of thumb

- **If you have very limited memory / need on-device:** prefer Gemma small variants or tiny Phi/Gemma; aggressively quantize (4-bit) + light LoRA if you need task adaptation.
- **If you need a balance (best cost / capability):** 7–13B-class models (LLaMA 7/8B, mid Qwen/Gemma sizes) quantized + LoRA often give the best bang-for-buck for many production tasks.
- **If you need highest single-model quality:** 70B+ family (or enterprise Qwen large variants) — accept higher hosting cost; LoRA still helps but cost to run will be dominated by inference.
- **If you plan multi-task adapters / many small tasks:** choose a single base model and create multiple LoRA adapters — this saves storage and makes swapping behaviors trivial.
- **Prototype first:** benchmark a small/medium size with quantization + LoRA on a representative sample before upgrading to larger variants.

1.16.7 Practical checklist before committing

1. Define success metric (accuracy, latency, throughput, cost per 1k tokens).
2. Check license / redistribution constraints for each model candidate.
3. Estimate inference cost (GPU memory + tokens/sec) for each model size.
4. Prototype: try quantized (4-bit) run + LoRA on a 1–2% validation slice and measure improvements.
5. If using instruct variants, validate whether additional LoRA actually improves or just overfits.

1.16.8 Short recommendation examples

- Small startup / prototype with tight budget: Gemma-4B or Phi-2 (quantized) + LoRA for task-specific tuning.
- Production assistant with modest budget: LLaMA 7/8B (instruct) quantized + LoRA; scale to 13B if needed.
- Research / high-quality generation: LLaMA 70B (or larger Qwen top family); expect higher hosting cost — use LoRA for specialized tasks or multi-adapter workflows.

1.16.9 Why LLaMA 3.1 was chosen: Tokenization Convenience

When selecting a base model, many factors matter: performance, parameter count, instruction-tuned vs base variant, and even subtler aspects like tokenization.

- **Tokenization strategy:** LLaMA 3.1 has a convenient tokenization scheme for numbers:
 - Numbers between 0 and 999 are represented as **single tokens**.
 - In contrast, other models like Qwen or Phi tokenize each digit separately. For example, “999” would become three tokens.
- **Why this matters:**
 - When training a model to predict numerical values (like prices), having the number as a single token simplifies the problem.
 - The model can predict the **entire number in one token**, rather than combining multiple digits.
 - This improves the model’s ability to learn accurate token prediction and reduces the chance of errors (e.g., predicting \$9 or \$99 instead of \$999).
- **Edge for LLaMA:** While it is a small difference overall, this tokenization convenience makes training slightly easier for numeric regression-like tasks.
- **Other considerations in model selection:**
 - Parameter sizes and memory requirements.
 - Instruction-tuned vs base variants.
 - Leaderboard performance.
 - Flexibility for future experiments with other models.
- **Decision:** For this project, the team chose **LLaMA 3.1 (8B or 18B variant)** as the base model, primarily because:
 - Convenient tokenization for numeric tasks.
 - Strong community support and tooling.
 - Good tradeoff between model size, performance, and fine-tuning options (LoRA can be applied easily).

Summary

Takeaway

Even small details, such as how numbers are tokenized, can influence model selection. For tasks like price prediction, LLaMA 3.1’s ability to encode numbers up to 999 as single tokens provides a slight but useful advantage.

1.17 Key Hyperparameters in QLoRA

QLoRA has a few essential hyperparameters that control how efficient and accurate the fine-tuning will be. Below are the five most important ones, explained in simple terms.

1. **Target Modules** LoRA does not modify the entire model. Instead, it inserts adapters only into specific layers, called **target modules**. These are usually the linear projections inside the attention mechanism:

$$q_proj, \quad k_proj, \quad v_proj, \quad o_proj$$

Why these? Because attention layers control how information flows between tokens, so small tweaks here can have a big impact.

Analogy: Think of a giant office building (the LLM). Instead of renovating every single room, you only upgrade the **elevators and stairways** (attention layers), because they control how people (tokens) move around.

Rule of Thumb

Start with **Q (query)** and **V (value)** projections. Add to k_proj and o_proj only if the task is complex and requires more adaptation.

2. **Rank (r)** The rank r decides the size of the low-rank adapter matrices A and B :

$$\Delta W \approx A \times B, \quad A \in \mathbb{R}^{d \times r}, \quad B \in \mathbb{R}^{r \times k}$$

A higher r means more trainable parameters, giving more expressive power, but it costs more memory.

Analogy: Imagine fitting a curve to data points. A small r is like fitting a straight line (fast, but may miss details). A large r is like fitting a flexible curve (captures more detail, but risks overfitting).

Rule of Thumb

Use $r = 8$ for small or domain-specific tasks. Use $r = 16\text{--}32$ for harder tasks like dialogue or reasoning.

3. **Alpha (α)** Alpha scales the contribution of the LoRA adapter relative to the frozen base model:

$$W' = W + \frac{\alpha}{r} AB$$

If α is too low, the adapters barely affect the model. If too high, they overwhelm the base model.

Analogy: Think of mixing paint colors. The frozen model is your base paint, and LoRA updates are small drops of dye. α controls how strong the dye looks in the final mix or how many drops of dye are used. If drops are too strong, the final color will be too bright and unnatural. If drops are too weak, the final color will be too dull and pale.

Rule of Thumb

Set $\alpha \approx 2r$. Example: $r = 16 \Rightarrow \alpha = 32$. This balances stability and expressiveness.

4. **Quantization** QLoRA stores the base model in **4-bit NF4 quantization** instead of 16/32-bit. This saves a huge amount of memory, making it possible to fine-tune large LLMs on a single GPU. The trick: LoRA adapters are still trained in higher precision (like 16-bit) while the frozen model stays compressed.

Analogy: Imagine keeping your huge library in a compressed digital format (4-bit). When you read, you only extract the parts you need at higher quality.

Rule of Thumb

Always use **4-bit NF4** quantization. It offers the best balance between compression and accuracy.

5. **Dropout** LoRA adapters add dropout to prevent overfitting. This randomly turns off some adapter neurons during training, making the model more robust.

Analogy: Like training a basketball team where sometimes players sit out, so others learn to play better. This prevents over-reliance on a few players.

Rule of Thumb

Use **0.05** dropout by default. Increase to **0.1** if your dataset is noisy or small.

1.17.1 Quick Summary Table

Hyperparameter	Rule of Thumb
Target Modules	Start with q_proj , v_proj ; expand if needed.
Rank (r)	8 for small tasks, 16–32 for harder ones.
Alpha (α)	About $2r$ (e.g., $r = 16 \Rightarrow \alpha = 32$).
Quantization	Always use 4-bit NF4 for efficiency.
Dropout	0.05 by default, 0.1 for noisy/small data.

Five Important Hyper-parameters for QLoRA...



Target Modules



r



Alpha



Quantization



Dropout



udem

1.18 Training Hyperparameters

1.18.1 Epochs

Definition: One epoch means the model has seen the *entire dataset once*.

Why Important:

- Too few epochs → the model underfits, meaning it hasn't learned enough patterns.
- Too many epochs → the model overfits, memorizing the training data and performing poorly on unseen data.

Step-by-Step Example:

1. Dataset has 10,000 samples.
2. Epoch 1: model sees all 10,000 samples → initial learning.
3. Epoch 2: model sees all 10,000 samples again → refines understanding.
4. Epoch 10: model may memorize exact data → risk of overfitting.

Analogy: Reading a textbook multiple times:

- 1 epoch = read once → understand a little.
- 5 epochs = read 5 times → remember more.
- 100 epochs = memorize word-for-word → may fail on questions requiring reasoning.

Rule of Thumb

Small datasets: 3–5 epochs. Large datasets: 1–3 epochs. Use **early stopping** to prevent overfitting.

Batch Size

Definition: Number of training samples processed together before updating model weights.

Why Important:

- Small batch size → noisy updates but better generalization.
- Large batch size → smoother updates, faster convergence, higher memory use.

Step-by-Step Example:

1. Batch size = 16 → model processes 16 samples, computes gradients, updates weights.
2. Batch size = 64 → model processes 64 samples, smoother gradient, fewer updates per epoch.

Analogy: Teacher giving feedback:

- Batch size 1 → feedback for each student immediately.
- Batch size 32 → teacher waits for all 32 students before giving feedback.

Rule of Thumb

Start with 16–64. If GPU memory is limited, use smaller batches + gradient accumulation.

Learning Rate

Definition: Step size for updating weights in the direction of the gradient.

Why Important:

- Too high → jumps over minima, training becomes unstable.
- Too low → very slow learning, may get stuck in sub-optimal minima.

Step-by-Step Example:

1. $\text{LR} = 0.1$ → model oscillates, can't converge.
2. $\text{LR} = 0.00001$ → very slow progress, may need many epochs.

Analogy: Learning to ride a bicycle:

- High LR → overcorrecting turns, wobble/fall.
- Low LR → move too slowly, take forever to balance.

Rule of Thumb

Fine-tuning: $2e-5$ to $5e-5$. Sensitive models: $1e-5$. Use **learning rate scheduler** to reduce LR over time.

Gradient Accumulation

Definition: Sum gradients from several mini-batches before performing a weight update.

Why Important:

- Saves GPU memory (can use smaller mini-batches).
- Stabilizes training, simulates larger batch size.

Step-by-Step Example:

1. Desired batch size = 64, GPU can only handle 16.
2. Feed 16 samples → compute gradients → do not update weights.
3. Repeat 3 more times → accumulate gradients.
4. Update weights once after 4 mini-batches.

Analogy: Saving money in a jar:

- Want to deposit \$100 (effective batch) but have \$25 daily (mini-batches).
- Save 4 days → deposit \$100 once (weight update).

Rule of Thumb

Use accumulation steps = 2–8 if GPU memory is small. Effective batch size = batch size × accumulation steps.

Optimizer

Definition: Algorithm to update weights based on computed gradients.

Why Important:

- SGD: simple, stable, slow.
- Adam: adaptive, faster convergence.
- AdamW: Adam + weight decay → regularization, prevents overfitting.

Analogy: Driving style:

- SGD → manual car, slow and steady.
- Adam → modern automatic car.
- AdamW → modern car with lane assist (stays on track).

Rule of Thumb

Use AdamW, weight decay = 0.01, betas = (0.9, 0.999).

1.19 DataCollatorForCompletionOnlyLM in HuggingFace TRL

1.19.1 Problem

When training a language model to predict only a specific target (e.g., **price**), we do not want it to learn the entire product description.

Example Input:

Product: Apple iPhone 15 Pro Max, 512GB, Color: Silver. Price is \$

- The context: Product: Apple iPhone 15 Pro Max, 512GB, Color: Silver.
- The target/completion: 1200 (what comes after “Price is \$”).

Without handling this properly, the model tries to predict the entire description, which is unnecessary.

1.19.2 Traditional Approach (Complicated)

- Manually create a **mask** for the loss:
 - Tokens before “Price is \$” → ignore in loss calculation.
 - Tokens after “Price is \$” → calculate loss.
- This requires manually aligning token indices, which is error-prone.

1.19.3 HuggingFace TRL Solution: DataCollatorForCompletionOnlyLM

Usage

```
from trl import DataCollatorForCompletionOnlyLM

response_template = "Price is $"
collator = DataCollatorForCompletionOnlyLM(
    response_template, tokenizer=tokenizer
)
```

How it works:

- **response_template** tells the collator where the context ends and the completion begins.
- Tokens before “Price is \$” → set to -100 in labels (ignored in loss).
- Tokens after “Price is \$” → normal token IDs (used in loss computation).
- The model sees the full input but only learns from the completion part.

1.19.4 Example

Input text:

Product: Apple iPhone 15 Pro Max, 512GB, Color: Silver. Price is \$1200

Tokenized (simplified):

[Product:, Apple, iPhone, 15, Pro, . . . , Price, is, \$, 1200]

Labels (for loss):

[-100, -100, -100, -100, -100, -100, -100, -100, -100, -100, -100, -100, 1200]

- -100 → ignore in loss
- 1200 → used in loss

1.19.5 Why it is important

- Ensures the model focuses only on the relevant tokens (price).
- Reduces wasted computation and speeds up training.
- Useful for task-specific fine-tuning such as price prediction, summarization, or Q&A.

Summary

- Collator automatically creates attention masks and labels.
- Context tokens are ignored in loss.
- Completion tokens are used for loss.
- The model learns the target task efficiently.

1.19.6 Warmup Ratio (`warmup_ratio`)

Simple Explanation: Think of training a model like teaching a child to ride a bicycle.

- You don't let the child go full speed immediately.
- Instead, you start slow, then gradually increase the speed.

In training, the `learning rate` is like the speed. `warmup_ratio` decides how much of the training time we slowly increase the learning rate from 0 to the target value.

Why it matters:

- If the learning rate starts too high, the model can make big mistakes and become unstable.
- Gradually increasing it lets the model "warm up" and learn safely.
- After the warmup period, the model can learn faster without crashing.

How it works:

- Suppose total training steps = 1000.
- `warmup_ratio = 0.1` ⇒ first 100 steps increase learning rate from 0 to normal.
- After 100 steps, the usual learning rate schedule (like linear or cosine decay) takes over.

Without Warmup:

- The model might start making huge updates.
- This can cause unstable training, slow convergence, or even NaN losses.
- Warmup prevents this “shock” at the start.

Practical Tips:

- Small warmup ratio (5-10%) is usually enough.
- For very large models, slightly longer warmup can help.
- Too long warmup slows down training unnecessarily.

Easy Rule of Thumb

- Start with 5% – 10% of total steps as warmup.
- Gradually increase the learning rate from 0 to your target.
- After warmup, let the usual scheduler handle learning rate decay.

1.20 Checkpointing and Resuming Training in HuggingFace / Colab

1.20.1 Overview

Checkpointing is the process of saving your model and training state periodically so you can resume training later without losing progress. This is especially useful in environments like Google Colab where the runtime is temporary and may disconnect.

Why Checkpointing is Important

- Saves model weights, optimizer state, and training progress.
- Allows training to resume after interruptions.
- Prevents loss of progress due to runtime disconnects or crashes.
- Makes training safer for long-running experiments.

1.20.2 Checkpoint Components

A checkpoint typically contains:

- **Model weights:** `model.state_dict()` stores all parameters.
- **Optimizer state:** `optimizer.state_dict()` contains learning rates, momentum, etc.
- **Current epoch:** so training resumes from the correct point.

1.20.3 Manual Checkpointing in PyTorch

PyTorch Example

```
import torch

# Save checkpoint
torch.save({
    'epoch': epoch,
    'model_state_dict': model.state_dict(),
    'optimizer_state_dict': optimizer.state_dict()
}, "checkpoint.pth")

# Load checkpoint
checkpoint = torch.load("checkpoint.pth")
model.load_state_dict(checkpoint['model_state_dict'])
optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
start_epoch = checkpoint['epoch'] + 1
```

1.20.4 Automatic Checkpointing with HuggingFace Trainer / SFT-Trainer

HuggingFace Trainer or SFTTrainer can handle checkpointing automatically:

TrainingArguments Example

```
from transformers import TrainingArguments, SFTTrainer

training_args = TrainingArguments(
    output_dir="results",           # folder to save checkpoints
    num_train_epochs=5,
    per_device_train_batch_size=64,
    save_steps=100,                # save every 100 steps
    save_total_limit=3,             # keep only last 3 checkpoints
    logging_steps=50,
    evaluation_strategy="steps",
    eval_steps=50,
)

trainer = SFTTrainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=eval_dataset
)

trainer.train(resume_from_checkpoint=True) # resumes automatically
```

Key points:

- `save_steps`: after how many steps to save a checkpoint.
- `save_total_limit`: how many checkpoints to keep; older ones are deleted.
- `resume_from_checkpoint=True`: automatically resumes from the latest checkpoint.

1.20.5 Checkpointing Flow Step-by-Step

1. **Start training:** check if checkpoint exists.
2. **Checkpoint exists:** load model weights, optimizer state, and last completed epoch.
3. **Set starting epoch:** `start_epoch = checkpoint['epoch'] + 1`.
4. **Resume training:** continue from the last saved epoch.
5. **Save checkpoint after each epoch or step:** ensures latest progress is always saved.

1.20.6 Checkpointing in Google Colab

Colab's local filesystem is temporary. To preserve checkpoints:

Resume Training in Colab

- **Option A: Mount Google Drive**

```
from google.colab import drive
drive.mount('/content/drive')

training_args = TrainingArguments(
    output_dir="/content/drive/MyDrive/mnist_sft_results",
    num_train_epochs=5,
    per_device_train_batch_size=64,
    save_steps=100,
    save_total_limit=3
)
```
- **Option B: Copy checkpoints manually**
`!cp -r /content/mnist_sft_results /content/drive/MyDrive/`

Resume Training in Colab

After reconnecting to Colab:

- Mount Drive if using Option A.
- Make sure `output_dir` points to the same folder.
- Start training:
`trainer.train(resume_from_checkpoint=True)`
- Trainer automatically detects latest checkpoint and resumes training.

1.20.7 Advantages of Using SFTTrainer Checkpointing

- Automatic saving of model, optimizer, scheduler, and training step.

- Easy resuming after runtime disconnects.
- Cleaner management of checkpoints using `save_total_limit`.
- Optional integration with Weights & Biases (`report_to="wandb"`) for real-time metrics.

Summary

- Checkpointing ensures training can continue after interruptions.
- HuggingFace Trainer / SFTTrainer makes checkpointing automatic.
- In Colab, save checkpoints to Google Drive for persistence.
- Use `resume_from_checkpoint=True` to continue training seamlessly.

1.20.8 Checkpointing in Kaggle

Kaggle notebooks have a temporary filesystem similar to Colab. To preserve checkpoints:

Important Note

- Kaggle notebooks reset after the session ends.
- Local checkpoints stored in `/kaggle/working/` are temporary.
- To keep checkpoints permanently, download them manually or push them to Kaggle Datasets.

Manual Checkpointing

Example: Save to Kaggle Working Folder

```
checkpoint_path = "/kaggle/working/cnn_checkpoint.pth"
torch.save({
    'epoch': epoch,
    'model_state_dict': model.state_dict(),
    'optimizer_state_dict': optimizer.state_dict()
}, checkpoint_path)
```

Downloading Checkpoints

To make the checkpoint persistent, download it after the notebook finishes:

Download Example

```
from IPython.display import FileLink
FileLink("/kaggle/working/cnn_checkpoint.pth")
```

Resuming Training on Kaggle (Manual)

When restarting the notebook:

1. Upload your previously downloaded checkpoint back to `/kaggle/working/`.
2. Load it using `torch.load()` and restore model and optimizer states:

Example: SFTTrainer Automatic Checkpointing

```
checkpoint = torch.load("/kaggle/working/cnn_checkpoint.pth")
model.load_state_dict(checkpoint['model_state_dict'])
optimizer.load_state_dict(checkpoint['optimizer_state_dict'])
start_epoch = checkpoint['epoch'] + 1
```

3. Continue training from `start_epoch`.

Automatic Checkpointing with SFTTrainer

Example: SFTTrainer Automatic Checkpointing

```
from transformers import TrainingArguments, SFTTrainer

training_args = TrainingArguments(
    output_dir="/kaggle/working/mnist_sft_results", # checkpoint folder
    num_train_epochs=5,
    per_device_train_batch_size=64,
    save_steps=100,           # save checkpoint every 100 steps
    save_total_limit=3,       # keep last 3 checkpoints
)

trainer = SFTTrainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
)

# Resume training automatically if checkpoint exists
trainer.train(resume_from_checkpoint=True)
```

Key points for SFTTrainer:

- Checkpoints are stored in `/kaggle/working/mnist_sft_results/checkpoint-<step>`.
- Each checkpoint contains:
 - `pytorch_model.bin` – model weights
 - `optimizer.pt` – optimizer state
 - `scheduler.pt` – scheduler state
 - `trainer_state.json` – training step info
- Only the latest `save_total_limit` checkpoints are kept.
- To persist checkpoints beyond the Kaggle session, download them manually or save as a Kaggle Dataset.

Summary for Kaggle

Summary

- Kaggle filesystem is **temporary** (`/kaggle/working/`).
- **Manual checkpointing** gives full control but requires explicit save/load commands.
- **SFTTrainer checkpointing** is automatic and resumes training seamlessly.
- For long-term storage, either download checkpoints or save them as Kaggle Dataset.

1.21 Saving and Downloading Automatic Checkpoints in Kaggle

Kaggle notebooks have a **temporary filesystem** at `/kaggle/working/`. Any files, including checkpoints, will be deleted after the session ends. To preserve your training progress, you should save checkpoints as a **zip file** and download it.

1.21.1 When to Enter the Zip Command

After your training is complete, or at any point when you want to backup your checkpoints:

Zip Checkpoints Command

```
!zip -r /kaggle/working/mnist_sft_results.zip /kaggle/working/mnist_sft_results
```

Explanation:

- `zip` → creates a compressed archive.
- `-r` → includes all files and subfolders recursively.
- `/kaggle/working/mnist_sft_results.zip` → path and name of the zip file to create.
- `/kaggle/working/mnist_sft_results` → folder containing the automatic checkpoints from SFTTrainer.

1.21.2 How to Download the Zip File

After creating the zip, you can download it directly from Kaggle:

Download Example

```
from IPython.display import FileLink  
FileLink("/kaggle/working/mnist_sft_results.zip")
```

Steps:

1. Run the zip command after your training to package all checkpoints.
2. Use the `FileLink` command to generate a clickable link.
3. Click the link in the notebook output to download the zip file to your local machine.

1.21.3 Restoring Checkpoints from Zip

Once downloaded, you can restore the checkpoints in a new Kaggle session or locally:

Restore Example

```
!unzip /kaggle/working/mnist_sft_results.zip -d /kaggle/working/  
  
# Now you can resume training with SFTTrainer  
trainer.train(resume_from_checkpoint="/kaggle/working/mnist_sft_results/  
/checkpoint-<step>")
```

Notes:

- The `-d` flag in `unzip` specifies the destination folder.
- Replace `<step>` with the checkpoint step you want to resume from.
- Using this method ensures your training progress is safe even after the Kaggle session ends.

1.22 Saving and Resuming Checkpoints with Hugging Face Hub

1.22.1 Manual Upload to Hugging Face Hub

You can also manually push any file (e.g., PyTorch .pth checkpoint) to Hugging Face Hub without using Trainer.

Manual Push to Hugging Face Hub

```
from huggingface_hub import Repository
import os

# Clone your Hugging Face repo (create one on HF website first)
repo_url = "https://huggingface.co/username/mnist-cnn"
repo_dir = "mnist-cnn-repo"

repo = Repository(local_dir=repo_dir, clone_from=repo_url)

# Copy checkpoint into the repo folder
os.system("cp /kaggle/working/cnn_checkpoint.pth mnist-cnn-repo/")

# Push checkpoint to Hugging Face Hub
repo.push_to_hub(commit_message="Added CNN checkpoint")
```

1.22.2 Manual Download from Hugging Face Hub

To load a manually uploaded checkpoint:

Manual Download Example

```
from huggingface_hub import hf_hub_download

checkpoint_path = hf_hub_download(
    repo_id="username/mnist-cnn",
    filename="cnn_checkpoint.pth"
)
```

1.22.3 Why Use Hugging Face Hub?

- **Persistent:** Unlike Kaggle/Colab, checkpoints are never lost.
- **Cross-platform:** Works seamlessly across Colab, Kaggle, and local machines.
- **Versioning:** Every push is tracked with history.
- **Private/Shared:** You can keep repos private or make them public.

Summary

Hugging Face Hub acts as a universal storage for models and checkpoints. You can:

- Use `push_to_hub=True` in Trainer for automatic checkpoint syncing.
- Manually push any file (model, tokenizer, optimizer states, datasets).
- Resume training from any machine with `resume_from_checkpoint=True`.

This makes training workflows more robust and eliminates worries about temporary storage.

1.23 Automatic Checkpointing to Hugging Face Hub

This guide demonstrates a workflow to save and manage model checkpoints automatically on Hugging Face Hub during training.

1.23.1 Authenticate with Hugging Face

Login to Hugging Face

```
from huggingface_hub import login, create_repo, upload_file

# Login using your HF token
hf_token = "YOUR_HF_TOKEN"
login(hf_token, add_to_git_credential=True)

# Set project and run names
PROJECT_NAME = "your-project-name"
HF_USER = "your-username"
RUN_NAME = "YYYY-MM-DD_HH.MM.SS" # Example: current datetime
REPO_ID = f"{HF_USER}/{PROJECT_NAME}-{RUN_NAME}"

# Create repo if it does not exist
create_repo(REPO_ID, exist_ok=True, private=True)
```

Define Model and Optimizer

Example CNN Model

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

class CNN(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(1, 32, 3, 1)
        self.conv2 = nn.Conv2d(32, 64, 3, 1)
        self.fc1 = nn.Linear(9216, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        x = F.relu(self.conv1(x))
        x = F.relu(self.conv2(x))
        x = F.max_pool2d(x, 2)
        x = torch.flatten(x, 1)
        x = F.relu(self.fc1(x))
        x = self.fc2(x)
        return x

# Device and optimizer
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model = CNN().to(device)
optimizer = optim.Adam(model.parameters(), lr=0.001)
```

Training Loop with Automatic HF Checkpointing

Training with Auto-Save

```
import os

# Local checkpoint folder
LOCAL_DIR = "checkpoints"
os.makedirs(LOCAL_DIR, exist_ok=True)
MAX_KEEP = 3
EPOCHS = 5

for epoch in range(1, EPOCHS+1):
    model.train()
    for data, target in train_loader:
        data, target = data.to(device), target.to(device)
        optimizer.zero_grad()
        output = model(data)
        loss = F.cross_entropy(output, target)
        loss.backward()
        optimizer.step()

    print(f"Epoch {epoch} finished, loss={loss.item():.4f}")

    # Save checkpoint locally
    ckpt_name = f"checkpoint_epoch_{epoch}.pth"
    ckpt_path = os.path.join(LOCAL_DIR, ckpt_name)
    torch.save({
        "epoch": epoch,
        "model_state_dict": model.state_dict(),
        "optimizer_state_dict": optimizer.state_dict(),
    }, ckpt_path)

    # Upload checkpoint to Hugging Face Hub
    upload_file(
        path_or_fileobj=ckpt_path,
        path_in_repo=ckpt_name,
        repo_id=REPO_ID,
        commit_message=f"Upload checkpoint epoch {epoch}"
    )

    # Cleanup old checkpoints (keep only last MAX_KEEP)
    checkpoints = sorted([f for f in os.listdir(LOCAL_DIR)
        if f.startswith("checkpoint_epoch_")])
    if len(checkpoints) > MAX_KEEP:
        for old_ckpt in checkpoints[:-MAX_KEEP]:
            os.remove(os.path.join(LOCAL_DIR, old_ckpt))
            print(f"Deleted old checkpoint: {old_ckpt}")
```

Resuming Training from HF Hub

Resume Training

```
from huggingface_hub import hf_hub_download
import torch

# Download latest checkpoint from HF Hub
checkpoint_path = hf_hub_download(
    repo_id=REPO_ID,
    filename="checkpoint_epoch_5.pth"
)

# Load checkpoint
checkpoint = torch.load(checkpoint_path)
model.load_state_dict(checkpoint["model_state_dict"])
optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
start_epoch = checkpoint["epoch"] + 1

# Continue training from start_epoch
for epoch in range(start_epoch, EPOCHS+1):
    model.train()
    for data, target in train_loader:
        data, target = data.to(device), target.to(device)
        optimizer.zero_grad()
        output = model(data)
        loss = F.cross_entropy(output, target)
        loss.backward()
        optimizer.step()

    print(f"Epoch {epoch} finished, loss={loss.item():.4f}")
    # Optionally repeat checkpoint upload here
```

Key Points / Best Practices

Notes

- **Persistent storage:** HF Hub stores checkpoints permanently and versioned.
- **Automatic upload:** Use `upload_file` after saving locally.
- **Manage local storage:** Keep only a few recent checkpoints to avoid using too much disk space.
- **Reusable workflow:** Change model, dataset, and `REPO_ID` to use this template for any project.

1.24 Resuming Fine-Tuning with SFTTrainer and LoRA

This guide demonstrates how to resume fine-tuning a previously LoRA-adapted model using SFTTrainer, save to Hugging Face Hub, and perform inference.

1.24.1 Define Constants and Paths

Project Paths and Model Info

```
BASE_MODEL = "meta-llama/Meta-Llama-3.1-8B"
PROJECT_NAME = "pricer"
HF_USER = "sijanpaudel"

# Dataset on Hugging Face Hub
DATASET_NAME = f"{HF_USER}/pricer-data"

# Previously fine-tuned checkpoint to resume
ORIGINAL_FINETUNED = f"{HF_USER}/{PROJECT_NAME}-2024-09-13_13.04.39"
# previous fine-tuned model repo
# If you want to resume from a specific commit, set the revision
# Otherwise, set to None to use the latest
ORIGINAL_REVISION = "e8d637df551603dc86cd7a1598a8f44af4d7ae36" # commit hash
#ORIGINAL_REVISION = None

# Unique names for current run
RUN_NAME = f"{datetime.now():%Y-%m-%d_%H.%M.%S}"
PROJECT_RUN_NAME = f"{PROJECT_NAME}-{RUN_NAME}"
HUB_MODEL_NAME = f"{HF_USER}/{PROJECT_RUN_NAME}"
```

Login to HuggingFace Hub and Weights & Biases

Authentication

```
# HuggingFace login
hf_token = userdata.get('HF_TOKEN')
login(hf_token, add_to_git_credential=True)

# Weights & Biases login
wandb_api_key = userdata.get('WANDB_API_KEY')
os.environ["WANDB_API_KEY"] = wandb_api_key
wandb.login()
```

Load Dataset

Load Dataset from Hub

```
from datasets import load_dataset

dataset = load_dataset(DATASET_NAME)
train = dataset['train']
test = dataset['test']
```

Configure LoRA and QLoRA Hyperparameters

LoRA and Quantization Settings

```
EPOCHS = 2
LORA_ALPHA = 64
LORA_R = 32
LORA_DROPOUT = 0.1
BATCH_SIZE = 16
GRADIENT_ACCUMULATION_STEPS = 1
LEARNING_RATE = 2e-5
LR_SCHEDULER_TYPE = 'cosine'
WEIGHT_DECAY = 0.001
TARGET_MODULES = ["q_proj", "v_proj", "k_proj", "o_proj"]
MAX_SEQUENCE_LENGTH = 182
QUANT_4_BIT = True
```

Load Tokenizer and Base Model (Quantized)

```
Tokenizer & Quantized Model

from transformers import AutoTokenizer, AutoModelForCausalLM, BitsAndBytesConfig

if QUANT_4_BIT:
    quant_config = BitsAndBytesConfig(
        load_in_4bit=True,
        bnb_4bit_use_double_quant=True,
        bnb_4bit_compute_dtype=torch.bfloat16,
        bnb_4bit_quant_type="nf4"
    )
else:
    quant_config = BitsAndBytesConfig(
        load_in_8bit=True,
        bnb_8bit_compute_dtype=torch.bfloat16
    )

tokenizer = AutoTokenizer.from_pretrained(BASE_MODEL, trust_remote_code=True)
tokenizer.pad_token = tokenizer.eos_token
tokenizer.padding_side = "right"

base_model = AutoModelForCausalLM.from_pretrained(
    BASE_MODEL,
    quantization_config=quant_config,
    device_map="auto",
    use_cache=False
)
base_model.generation_config.pad_token_id = tokenizer.pad_token_id
```

Load Previously Fine-Tuned Model (PEFT)

Resume LoRA Checkpoint

```
from peft import PeftModel

if ORIGINAL_REVISION:
    fine_tuned_model = PeftModel.from_pretrained(
        base_model,
        ORIGINAL_FINETUNED,
        revision=ORIGINAL_REVISION,
        is_trainable=True
    )
else:
    fine_tuned_model = PeftModel.from_pretrained(
        base_model,
        ORIGINAL_FINETUNED,
        is_trainable=True
    )

fine_tuned_model.train()
```

Data Collator for Price Prediction

Data Collator

```
from trl import DataCollatorForCompletionOnlyLM

response_template = "Price is $"
collator = DataCollatorForCompletionOnlyLM(response_template,
tokenizer=tokenizer)
# The collator ensures only the part after the response_template is used for
loss calculation during training as tokenizer padding is on the right. The part
before response_template is the prompt which is only used for context(no tokens).
```

Configure SFTTrainer and LoRA

```
SFTTrainer Configuration

from peft import LoraConfig
from trl import SFTTrainer, SFTConfig

peft_parameters = LoraConfig(
    lora_alpha=LORA_ALPHA,
    lora_dropout=LORA_DROPOUT,
    r=LORA_R,
    bias="none",
    task_type="CAUSAL_LM",
    target_modules=TARGET_MODULES,
)

train_params = SFTConfig(
    output_dir=PROJECT_RUN_NAME,
    num_train_epochs=EPOCHS,
    per_device_train_batch_size=BATCH_SIZE,
    per_device_eval_batch_size=1,
    eval_strategy="no",
    gradient_accumulation_steps=GRADIENT_ACCUMULATION_STEPS,
    optim="paged_adamw_32bit",
    save_steps=2000,
    save_total_limit=10,
    logging_steps=50,
    learning_rate=LEARNING_RATE,
    weight_decay=WEIGHT_DECAY,
    fp16=False,
    bf16=True,
    max_grad_norm=0.3,
    max_steps=-1,
    warmup_ratio=0.03,
    group_by_length=True,
    lr_scheduler_type=LR_SCHEDULER_TYPE,
    report_to="wandb",
    run_name=RUN_NAME,
    max_seq_length=MAX_SEQUENCE_LENGTH,
    dataset_text_field="text",
    save_strategy="steps",
    hub_strategy="every_save",
    push_to_hub=True,
    hub_model_id=HUB_MODEL_NAME,
    hub_private_repo=True
)
```

Initialize Trainer and Fine-Tune

Fine-Tuning

```
fine_tuning = SFTTrainer(
    model=fine_tuned_model,
    train_dataset=train,
    peft_config=peft_parameters,
    args=train_params,
    data_collator=collator,
)

fine_tuning.train()

# Push model to HuggingFace Hub
fine_tuning.model.push_to_hub(PROJECT_RUN_NAME, private=True)
```

Inference Example

Inference Example

```
prompt = "How much does this cost to the nearest dollar?\n\nPower Stop Rear Z36 Truck and Tow Brake Kit with Calipers\n\n### Price:\n$"
inputs = tokenizer(prompt, return_tensors="pt").to("cuda")
outputs = fine_tuned_model.generate(**inputs, max_new_tokens=50)
print(tokenizer.decode(outputs[0], skip_special_tokens=True))
```

1.24.2 Checkpointing in SFTTrainer with LoRA

What Actually Happens During Training

- Unlike standard HuggingFace Trainer, **no checkpoint-<step> folders are automatically created.**
- Automatically saved files (locally or on Hugging Face Hub):
 - `adapter_model.safetensors` – LoRA adapter weights (trained parameters)
 - `training_args.bin` – training configuration (learning rate, optimizer type, scheduler type, etc.)
 - Tokenizer/config files if pushing to Hub (`tokenizer.json`, `tokenizer_config.json`, etc.)
- **Optimizer and scheduler states are NOT saved automatically.**
- To fully resume training with the same optimizer state and learning rate, manual saving is required.

1.24.3 Inferencing using SFT

Running Inference with Fine-tuned Model

1. Load the tokenizer and set padding:

```
tokenizer = AutoTokenizer.from_pretrained(BASE_MODEL, trust_remote_code=True)
tokenizer.pad_token = tokenizer.eos_token
tokenizer.padding_side = "right"
```

Explanation: The tokenizer converts text into numerical tokens. We set the padding token to the EOS token to ensure proper alignment of inputs when batch processing sequences of different lengths.

2. Load the base model:

```
base_model = AutoModelForCausalLM.from_pretrained(
    BASE_MODEL,
    quantization_config=quant_config,
    device_map="auto",
)
base_model.generation_config.pad_token_id = tokenizer.pad_token_id
```

Explanation: Loads the pre-trained base model (e.g., LLaMA). Quantization reduces memory usage, and device_map="auto" distributes the model across available GPUs. We also set the generation padding token to match the tokenizer.

3. Attach the fine-tuned LoRA adapter:

```
if REVISION:
    fine_tuned_model = PeftModel.from_pretrained(
        base_model, FINETUNED_MODEL, revision=REVISION
    )
else:
    fine_tuned_model = PeftModel.from_pretrained(
        base_model, FINETUNED_MODEL
    )
```

Explanation: LoRA adapters contain only the fine-tuned weights. **REVISION** is the commit hash of the checkpoint you want to use for inference. If you specify a revision, it loads that exact checkpoint from the Hub or local folder.

Generating Text with the Fine-tuned Model

Generate predictions:

```
def model_predict(input_text):
    inputs = tokenizer.encode(input_text, return_tensors="pt").to("cuda")
    attention_mask = torch.ones(inputs.shape, device="cuda")
    outputs = fine_tuned_model.generate(inputs,
                                        attention_mask=attention_mask,max_new_tokens=50, num_return_sequences=1)
    response = tokenizer.decode(outputs[0], skip_special_tokens=True)
    return response
```

Explanation: Tokenized input is sent to the GPU. The `generate()` method produces text autoregressively. Finally, tokens are decoded back into readable text.

1.24.4 Improved Prediction Function

Weighted Top-K Prediction

```
top_K = 3

def improved_model_predict(prompt, device="cuda"):
    set_seed(42)
    inputs = tokenizer.encode(prompt, return_tensors="pt").to(device)
    attention_mask = torch.ones(inputs.shape, device=device)

    with torch.no_grad():
        outputs = fine_tuned_model(inputs, attention_mask=attention_mask)
        next_token_logits = outputs.logits[:, -1, :].to('cpu')

        next_token_probs = F.softmax(next_token_logits, dim=-1)
        top_prob, top_token_id = next_token_probs.topk(top_K)
        prices, weights = [], []
        for i in range(top_K):
            predicted_token = tokenizer.decode(top_token_id[0][i])
            probability = top_prob[0][i]
            try:
                result = float(predicted_token)
            except ValueError as e:
                result = 0.0
            if result > 0:
                prices.append(result)
                weights.append(probability)
        if not prices:
            return 0.0, 0.0
        total = sum(weights)
        weighted_prices = [price * weight /
                           total for price, weight in zip(prices, weights)]
        return sum(weighted_prices).item()
```

Explanation:

- This function predicts a numeric value (like a price) from a prompt.
- It uses the top **K=3** token probabilities and calculates a **weighted average**.
- The tokenizer decodes predicted tokens into numbers. Non-numeric tokens are ignored.
- Probabilities from **softmax** are used as weights for averaging.
- Returns a single float representing the **expected value** of the prediction.
- Setting a **seed** ensures reproducibility.

1.24.5 Key Points for SFT Inference

Important Notes

- Always load the **base model** first; adapters cannot be used standalone.
- The **tokenizer** ensures inputs are properly formatted and padded.
- **PeftModel.from_pretrained** loads LoRA adapters from a checkpoint.
- **REVISION** specifies the exact checkpoint (commit hash) for inference.
- Use `generate()` for autoregressive text generation.
- This setup allows efficient inference using a large pre-trained model plus lightweight adapters, without retraining the base model.

1.24.6 Summary

Key Points

- Local folder `PROJECT_RUN_NAME` stores checkpoints.
- `HUB_MODEL_NAME` is passed to `hub_model_id` for Hub upload.
- `SFTTrainer` automatically handles checkpoint saving, Hub push, and resuming.
- Use `resume_from_checkpoint=True` to continue training without losing progress.

1.25 Four Steps in LoRA Training

LoRA (Low-Rank Adaptation) fine-tuning modifies the standard training process by introducing small trainable low-rank matrices into certain layers of the model (typically attention projections). The original large model weights remain frozen, and only these adapter parameters are updated. Below are the four main steps in LoRA training.

1.25.1 Forward Pass with LoRA

- Input $\mathbf{x} \in \mathbb{R}^n$ is passed through the pre-trained model.
- In a normal weight matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, LoRA introduces a low-rank decomposition:

$$\Delta\mathbf{W} = \mathbf{B}\mathbf{A}, \quad \mathbf{A} \in \mathbb{R}^{r \times k}, \quad \mathbf{B} \in \mathbb{R}^{d \times r}, \quad r \ll \min(d, k)$$

- The effective weight during training is:

$$\mathbf{W}' = \mathbf{W} + \Delta\mathbf{W}$$

- Since \mathbf{W} is frozen, only \mathbf{A} and \mathbf{B} contribute learnable changes.
- Prediction is computed as:

$$\hat{\mathbf{y}} = f_{\theta, \{\mathbf{A}, \mathbf{B}\}}(\mathbf{x})$$

Loss Calculation

- The loss function $\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y})$ measures prediction error.
- Common choices:
 - Cross-Entropy Loss for language modeling:

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

- Other task-specific losses (e.g., regression, contrastive loss).
- The difference from standard training is that the loss only needs to drive updates to the LoRA parameters (\mathbf{A}, \mathbf{B}), not the entire model.

Backward Pass (Backpropagation with LoRA)

- Gradients are computed as usual via backpropagation:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{B}}$$

- Since \mathbf{W} is frozen, its gradient is ignored:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = 0$$

- The effective gradient flow is restricted to the low-rank adapters.
- This significantly reduces memory and computation, as only a small fraction of parameters participate in updates.

Optimization (Updating LoRA Parameters)

- Only the adapter parameters (\mathbf{A}, \mathbf{B}) are updated:

$$\mathbf{A} \leftarrow \mathbf{A} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{A}}, \quad \mathbf{B} \leftarrow \mathbf{B} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{B}}$$

- Optimizers such as Adam or AdamW are commonly used.
- The frozen \mathbf{W} remains unchanged throughout training, preserving pre-trained knowledge.
- This makes LoRA highly parameter-efficient: only a few million parameters are updated instead of billions.

1.25.2 Summary Table of LoRA Training Steps

Step	Description
Forward Pass (LoRA)	Compute predictions using $\mathbf{W}' = \mathbf{W} + \mathbf{BA}$, with \mathbf{W} frozen and only \mathbf{A}, \mathbf{B} trainable.
Loss Calculation	Compute loss $\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y})$ to measure prediction error, driving updates for LoRA parameters.
Backward Pass	Compute gradients only for \mathbf{A}, \mathbf{B} ; gradients for frozen \mathbf{W} are ignored.
Optimization	Update LoRA parameters (\mathbf{A}, \mathbf{B}) using gradient descent or Adam; \mathbf{W} remains unchanged.

1.25.3 Uploading Final Dataset to Hugging Face Hub

Dataset Upload Guide

1. Install required packages:

```
!pip install datasets huggingface_hub
```

2. Login to Hugging Face Hub:

```
DATASET_NAME = "sijanpaudel/amazon-pricing-dataset"
from huggingface_hub import login
login() # Enter your Hugging Face token
```

3. Prepare your datasets:

- `final_train_set` should contain:
 - `text`: list of formatted prompts with price
 - `price`: corresponding price values
- `final_test_set` should contain only `text` prompts (for evaluation)

4. Convert dictionaries to Hugging Face Dataset objects:

```
from datasets import Dataset
train_dataset = Dataset.from_dict(final_train_set)
test_dataset = Dataset.from_dict(final_test_set)
```

5. Combine into a DatasetDict:

```
from datasets import DatasetDict
dataset_dict = DatasetDict({
    "train": train_dataset,
    "test": test_dataset
})
```

6. Push the dataset to Hugging Face Hub:

```
dataset_dict.push_to_hub(DATASET_NAME, private=True)
```

7. Verify upload:

After pushing, the dataset will be available at <https://huggingface.co/datasets/sijanpaudel/amazon-pricing-dataset>

8. Pulling the dataset from Hugging Face Hub:

```
from datasets import DatasetDict
dataset_dict = DatasetDict.load_from_hub(DATASET_NAME)
train, test = dataset_dict["train"], dataset_dict["test"]
```

1.26 Multi-Agent AI Architecture for Automated Deal Finding Systems

This architecture describes a pipeline of specialized AI agents working in sequence to discover, analyze, and notify users about potential deals. The system follows a modular and agent-based approach to ensure scalability, flexibility, and reliability.

1.26.1 User Interface (UI)

The entry point of the system is the **User Interface**.

- Typically implemented as a web application (e.g., using Gradio).
- Provides an accessible platform where users can interact with the system.
- Users can configure preferences, set deal criteria, and receive outputs.

1.26.2 Agent Framework

At the core lies the **Agent Framework**, which functions as the infrastructure supporting all agents.

- Provides **memory** for agents to recall previous states, interactions, and contextual data.
- Maintains **logging** to record every step of agent activity for auditing, debugging, and transparency.
- Ensures inter-agent communication and workflow orchestration.

1.26.3 Planning Agent

The **Planning Agent** acts as the central coordinator of the system.

- Determines which specialized agent should execute a task at any given time.
- Breaks down high-level objectives (e.g., “find and evaluate deals”) into actionable sub-tasks.
- Manages sequencing, dependencies, and prioritization of agents.

1.26.4 Scanner Agent

The **Scanner Agent** is responsible for initial discovery.

- Scans multiple structured and unstructured data sources (websites, APIs, databases).
- Applies filtering and heuristic rules to identify potentially valuable deals.
- Passes candidate deals to downstream agents for deeper analysis.

1.26.5 Ensemble Agent

The **Ensemble Agent** carries out valuation and analysis.

- Uses multiple predictive models (*ensemble learning*) to estimate deal value.

- Combines outputs of different models to reduce bias and variance, leading to more robust results.
- Generates a final valuation score with confidence levels.

1.26.6 Messaging Agent

The final stage involves the **Messaging Agent**.

- Prepares results for end-users after analysis.
- Sends notifications (e.g., push alerts, emails, or in-app messages) summarizing potential deals.
- Ensures timely communication so that users can take action quickly.

1.26.7 Pipeline Summary

The flow of the system can be summarized as follows:

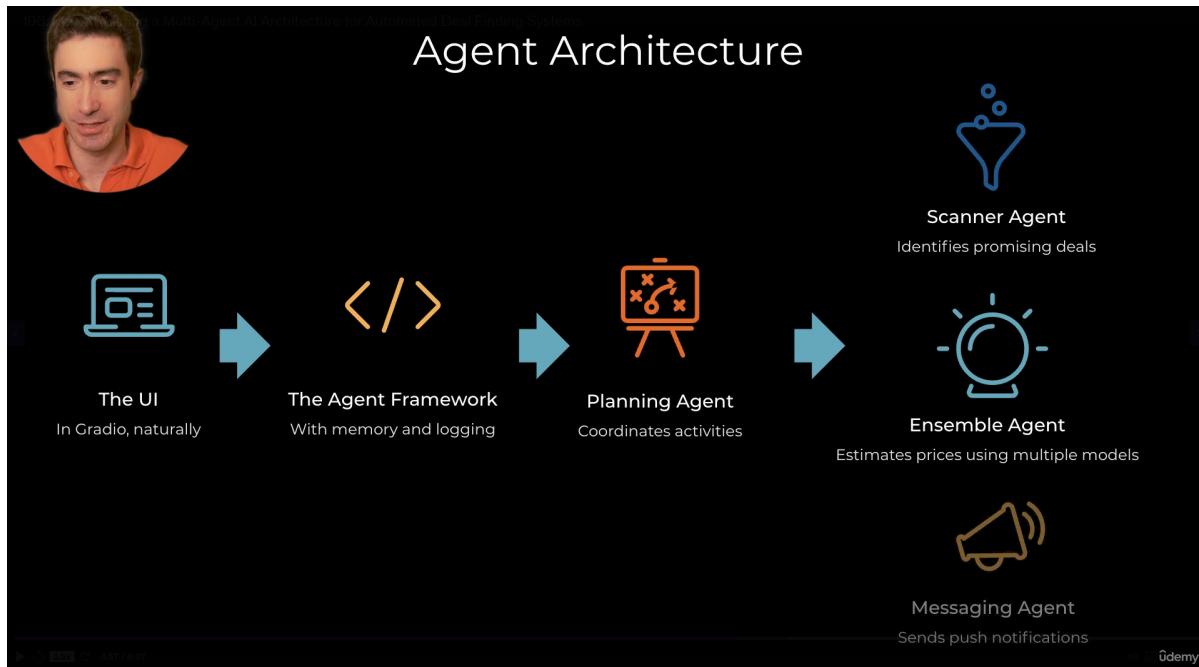
1. User interacts with the **UI** to configure requirements.
2. The **Agent Framework** initializes and monitors the workflow.
3. The **Planning Agent** orchestrates task execution.
4. The **Scanner Agent** discovers potential deals.
5. The **Ensemble Agent** evaluates deal value using multiple models.
6. The **Messaging Agent** notifies the user of results.

This modular multi-agent approach ensures that each agent specializes in one role, making the system scalable and maintainable while providing accurate and timely results.

1.26.8 Future Work

Future enhancements to the system could include:

- Integrating more advanced machine learning models for better deal evaluation.
- Expanding the range of data sources for the Scanner Agent.
- Improving user interface elements for better user experience.



1.27 Create a RAG Database with Our 400,000 Training Data

PLEASE NOTE:

We already have a very powerful product estimator with our proprietary, fine-tuned LLM. Most people would be very satisfied with that! The main reason we're adding these extra steps is to deepen your expertise with RAG and with Agentic workflows.

Step 1: Import necessary libraries. These include standard Python packages, numerical computing (numpy), data handling (pickle, datasets), Chroma DB for vector storage, HuggingFace tools, and embedding models.

```
[ ]: # imports
import os
import re
import math
import json
from tqdm import tqdm
import random
from dotenv import load_dotenv
from huggingface_hub import login
import numpy as np
import pickle
from sentence_transformers import SentenceTransformer, util
from datasets import load_dataset
import chromadb
from sklearn.manifold import TSNE
import plotly.graph_objects as go
```

Step 2: Setup environment variables and database path for API keys and Chroma storage.

```
[ ]: # environment setup
load_dotenv(override=True)
os.environ['OPENAI_API_KEY'] = os.environ.get('OPENAI_API_KEY', ↴
    'your-key-if-not-using-env')
os.environ['HF_TOKEN'] = os.environ.get('HF_TOKEN', 'your-key-if-not-using-env')
DB = 'products_vectorstore'
```

Step 3: Log in to HuggingFace hub using the token to access models and datasets.

```
[ ]: # HuggingFace login
hf_token = os.environ['HF_TOKEN']
login(hf_token, add_to_git_credential=True)
```

Step 4: Load preprocessed training data from pickle files.

```
[ ]: # Load training data
with open('../week6/train.pkl', 'rb') as f:
    train = pickle.load(f)
```

Step 5: Initialize Chroma persistent client, delete existing collection if any, and create a new collection for product vectors.

```
[ ]: # Create Chroma datastore
client = chromadb.PersistentClient(path=DB)
if 'products' in [c.name for c in client.list_collections()]:
    client.delete_collection('products')
collection = client.create_collection('products')
```

Step 6: Load sentence transformer embedding model for vectorizing product descriptions.

```
[ ]: # Initialize embedding model
model = SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
```

Step 7: Define a helper function to extract product descriptions from the training data.

```
[ ]: # Prepare text for vectorization
def description(item):
    return item.prompt.replace('How much does this cost to the nearest dollar?\n',
                                '').split('\n')[0]
```

Step 8: Loop through training data, vectorize descriptions, and add them to the Chroma collection.

```
[ ]: # Populate Chroma collection
NUMBER_OF_DOCUMENTS = len(train)
for i in range(0, NUMBER_OF_DOCUMENTS, 1000):
    docs = [description(item) for item in train[i:i+1000]]
    vectors = model.encode(docs).astype(float).tolist()
    metas = [{'category': item.category, 'price': item.price} for item in
              train[i:i+1000]]
    ids = [f'doc_{j}' for j in range(i, i+len(docs))]
    collection.add(ids=ids, documents=docs, embeddings=vectors, metadatas=metas)
```

1.28 Hallmarks of an Agentic AI Solution

1. Decomposition of Larger Problems into Smaller Steps

An **Agentic AI system** decomposes complex problems into smaller, manageable steps. Each step is handled by an *independent agent* or *specialized process*, allowing for:

- **Parallel processing** for faster execution.
- **Error reduction** by isolating tasks.
- **Task modularity** to reuse and recombine solutions.
- *Scalable problem solving* for multi-domain challenges.

This approach ensures that large goals can be achieved efficiently without overwhelming any single component.

2. Use of Tools, Function Calling, and Structured Outputs

Agents can call *external tools*, *APIs*, or defined functions to extend their capabilities. Structured outputs ensure clear and reliable results. Key benefits include:

- **Precision:** Outputs like JSON or tables can be directly processed.
- **Reliability:** Reduces ambiguity in multi-agent workflows.
- **Interoperability:** Allows agents to communicate results seamlessly.
- *Error handling:* Structured formats make it easier to validate and parse responses.

3. Collaborative Agent Environment

Agents operate in a shared environment enabling **collaboration and coordination**:

- **Specialized knowledge:** Different agents contribute unique skills.
- **Task distribution:** Efficient division of labor for complex goals.
- *Real-time communication:* Agents exchange updates to optimize outcomes.
- *Conflict resolution:* Conflicting suggestions can be reconciled by higher-level agents.

Collaboration enhances overall system performance beyond the capability of any single agent.

4. Planning and Coordination Agent

A **planning agent** orchestrates activities across all agents:

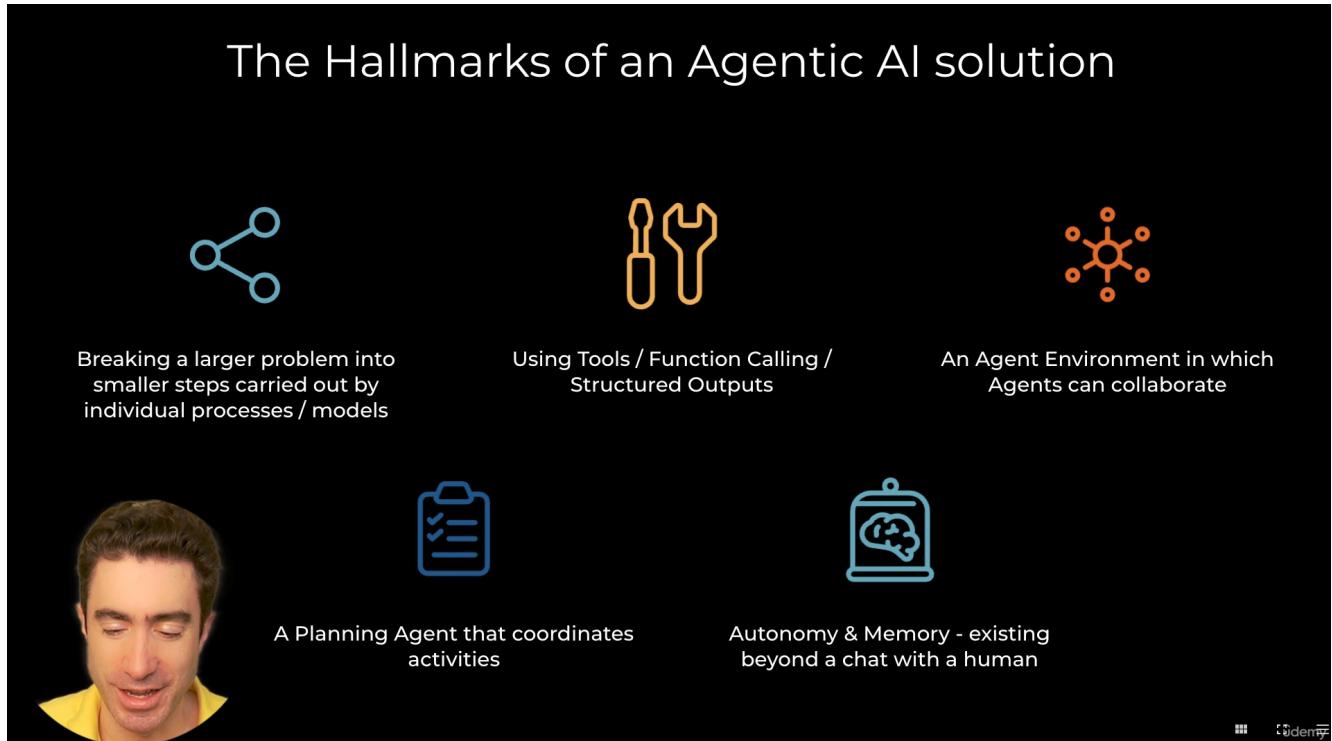
- Determines the sequence of tasks based on dependencies.
- Dynamically adjusts plans according to intermediate results.
- Prioritizes critical steps to achieve goals efficiently.
- Monitors progress and reallocates resources as needed.

This ensures that complex workflows are executed coherently and reliably.

5. Autonomy and Persistent Memory

Agentic AI demonstrates **autonomy** and maintains *persistent memory*, which allows:

- **Independent operation** without constant human intervention.
- **Learning from past interactions** to improve decision-making.
- Retention of critical information across sessions for long-term planning.
- *Cumulative knowledge growth*, enabling adaptation to new situations.
- Capability to handle multi-step and long-horizon tasks effectively.



Chapter 2

System Overview and Architecture

2.1 Introduction

This document provides an exhaustive technical analysis of a sophisticated Python-based multi-agent system designed for automated deal discovery, analysis, and notification. The system represents a compelling example of modern software engineering principles, combining object-oriented programming, machine learning, web scraping, API integration, and distributed computing into a cohesive, production-ready application.

The system operates as an intelligent deal-hunting assistant that continuously monitors RSS feeds from deal websites, applies multiple machine learning models to estimate product prices, identifies potentially valuable deals by comparing estimated prices with actual prices, and automatically notifies users of significant opportunities through various messaging channels.

2.1.1 System Purpose and Goals

The primary objective of this multi-agent system is to automate the process of finding profitable deals across multiple online retailers. The system accomplishes this through several key capabilities:

1. **Autonomous Deal Discovery:** Continuously monitors RSS feeds from deal aggregation websites like DealNews to identify new product offers
2. **Intelligent Price Estimation:** Employs multiple machine learning approaches including ensemble methods, fine-tuned language models, and traditional ML algorithms to predict fair market prices
3. **Opportunity Analysis:** Calculates potential savings by comparing estimated fair prices against offered deal prices
4. **Automated Notification:** Delivers alerts through multiple channels (SMS, push notifications) when significant deals are identified
5. **Memory Management:** Maintains awareness of previously processed deals to avoid duplicate notifications

2.1.2 Technical Architecture Overview

The system employs a modular, agent-based architecture where each component (agent) has specialized responsibilities. This design pattern promotes separation of concerns, maintainability, and scalability. The architecture can be categorized into several layers:

- **Orchestration Layer:** The PlanningAgent coordinates the entire workflow
 - **Data Acquisition Layer:** ScannerAgent and deals processing handle external data retrieval
 - **Machine Learning Layer:** Multiple specialized agents provide price estimation using different ML approaches
 - **Communication Layer:** MessagingAgent handles user notifications
 - **Foundation Layer:** Base Agent class and data models provide common functionality

2.2 Agent-Based Architecture

2.2.1 The Agent Pattern

The system implements the Agent design pattern, where each agent is an autonomous unit capable of:

- Independent operation with well-defined responsibilities
 - Communication and coordination with other agents
 - State management and decision making
 - Logging and monitoring of its activities

All agents inherit from a common `Agent` base class that provides:

- Colored console logging for visual differentiation during runtime
 - Consistent naming conventions
 - Standardized communication protocols

2.2.2 System Components

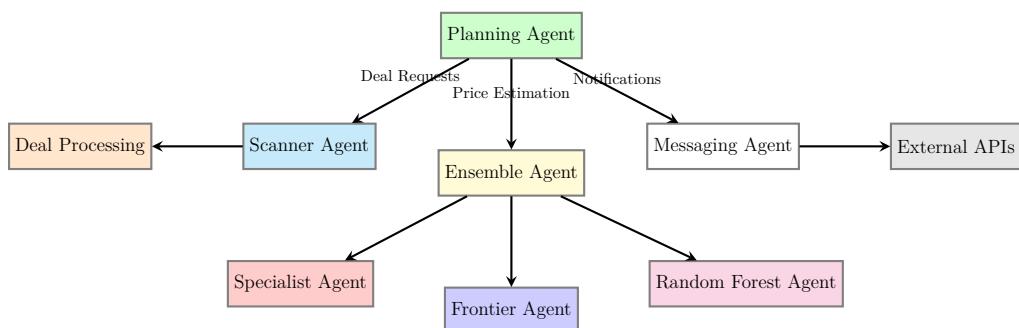


Figure 2.1: High-Level System Architecture showing agent relationships and data flow

Core Agents

PlanningAgent (Orchestrator). The PlanningAgent serves as the system's central coordinator, implementing the main workflow logic. It instantiates and manages all other agents, coordinates their interactions, and makes high-level decisions about deal processing. The PlanningAgent embodies the Command pattern, orchestrating a complex sequence of operations across multiple specialized components.

ScannerAgent (Data Acquisition). Responsible for monitoring external RSS feeds and identifying new deals. The system includes two implementations:

- **ScannerAgent:** Uses OpenAI's structured output API
- **ScannerAgentLangChain:** Uses Google's Gemini via LangChain framework

Both implementations extract deal information and convert unstructured web data into structured Deal objects.

EnsembleAgent (ML Coordinator). Implements an ensemble machine learning approach by coordinating three different price estimation models. It uses a meta-learning approach where a Linear Regression model combines predictions from specialist models to produce final price estimates.

MessagingAgent (Communication). Handles external communications through multiple channels including SMS (via Twilio) and push notifications (via Pushover). This agent implements the Publisher pattern, broadcasting deal alerts to configured endpoints.

Specialized ML Agents

SpecialistAgent (Fine-tuned Model). Utilizes a remotely-hosted, fine-tuned language model via Modal's cloud platform. This agent represents the most sophisticated ML approach in the system, leveraging domain-specific training for price estimation.

FrontierAgent (RAG-based LLM). Implements a Retrieval-Augmented Generation (RAG) approach using vector similarity search. It finds similar products in a ChromaDB vector store and uses this context to inform price predictions via OpenAI's GPT models.

RandomForestAgent (Traditional ML). Uses a pre-trained Random Forest model combined with sentence embeddings for price prediction. This represents a more traditional machine learning approach compared to the LLM-based agents.

2.3 Data Flow and System Interactions

2.3.1 Primary Workflow

The system operates through a well-defined workflow that demonstrates sophisticated inter-agent communication:

1. **Initialization Phase:** PlanningAgent instantiates all required agents, each performing their specific setup routines

2. **Deal Discovery:** ScannerAgent monitors RSS feeds and extracts new deals not in memory
3. **Price Analysis:** For each deal, EnsembleAgent coordinates price estimation across all ML agents
4. **Opportunity Calculation:** System compares estimated prices with deal prices to calculate potential savings
5. **Decision Making:** PlanningAgent determines whether deals meet threshold criteria for notification
6. **User Notification:** MessagingAgent delivers alerts through configured channels

2.3.2 Data Models

The system employs Pydantic models for robust data validation and serialization:

```

1  class Deal(BaseModel):
2      """Structured representation of a deal"""
3      product_description: str
4      price: float
5      url: str
6
7  class DealSelection(BaseModel):
8      """Container for multiple deals"""
9      deals: List[Deal]
10
11 class Opportunity(BaseModel):
12     """Represents a potential profitable deal"""
13     deal: Deal
14     estimate: float # ML-predicted price
15     discount: float # Estimated savings

```

Listing 2.1: Core Data Models

2.4 Advanced System Features

2.4.1 Machine Learning Integration

The system showcases multiple machine learning paradigms:

- **Ensemble Methods:** Meta-learning approach combining multiple models
- **Fine-tuned LLMs:** Domain-specific language model training
- **RAG (Retrieval-Augmented Generation):** Vector similarity search with LLM reasoning
- **Traditional ML:** Scikit-learn Random Forest with feature engineering
- **Vector Embeddings:** Sentence transformers for semantic similarity

2.4.2 External Service Integration

The system integrates with multiple external services:

Service Category	Implementation	Purpose
LLM APIs	OpenAI, DeepSeek, Google Gemini	Price estimation and reasoning
Vector Database	ChromaDB	Similarity search for RAG
Cloud Computing	Modal	Remote model hosting
Messaging	Twilio, Pushover	User notifications
Data Sources	RSS Feeds	Deal discovery

Table 2.1: External Service Integrations

2.4.3 Error Handling and Robustness

The system implements several robustness mechanisms:

- **Graceful Degradation:** Continues operation if individual agents fail
- **Data Validation:** Pydantic models ensure data integrity
- **Rate Limiting:** Built-in delays prevent API abuse
- **Memory Management:** Tracks processed deals to avoid duplicates
- **Logging:** Comprehensive logging for debugging and monitoring

2.5 Architectural Patterns

2.5.1 Design Patterns Implemented

The system demonstrates several well-known design patterns:

Agent Pattern Each component operates autonomously with defined responsibilities

Strategy Pattern Multiple ML algorithms can be swapped in the ensemble

Command Pattern PlanningAgent orchestrates complex operations

Factory Pattern Agent instantiation and configuration

Observer Pattern Event-driven notifications and logging

Template Method Base Agent class defines common behavior

2.5.2 SOLID Principles

The architecture adheres to SOLID principles:

- **Single Responsibility:** Each agent has one primary function
- **Open/Closed:** New agents can be added without modifying existing code
- **Liskov Substitution:** All agents can be used polymorphically
- **Interface Segregation:** Agents depend only on methods they use
- **Dependency Inversion:** High-level modules don't depend on low-level details

2.6 System Scalability and Performance

2.6.1 Scalability Considerations

The modular architecture enables several scaling strategies:

- **Horizontal Scaling:** Agents can be distributed across multiple processes/machines
- **Load Balancing:** Multiple instances of compute-intensive agents
- **Caching:** Vector embeddings and model predictions can be cached
- **Asynchronous Processing:** Background tasks for non-critical operations

2.6.2 Performance Characteristics

Key performance aspects include:

- **I/O Bound Operations:** RSS feed parsing, API calls
- **CPU Bound Operations:** ML model inference, text processing
- **Memory Usage:** Vector embeddings, model weights
- **Network Latency:** External API dependencies

2.7 Security and Privacy

2.7.1 Security Measures

The system implements several security best practices:

- **API Key Management:** Environment variables for sensitive credentials
- **Input Validation:** Pydantic models prevent injection attacks
- **Rate Limiting:** Prevents abuse of external services
- **Error Sanitization:** Logs don't expose sensitive information

2.7.2 Privacy Considerations

- **Data Minimization:** Only necessary data is collected and stored
- **External Dependencies:** User data may be sent to third-party APIs
- **Logging Practices:** Personal information excluded from logs

2.8 Chapter Summary

This chapter provided a comprehensive overview of the Python multi-agent deal discovery system. We explored:

- The system's purpose as an automated deal-hunting assistant
- The agent-based architecture promoting modularity and maintainability
- The sophisticated integration of multiple machine learning approaches
- The robust data models and external service integrations
- The adherence to established design patterns and software engineering principles

The system represents a excellent example of modern Python development, combining object-oriented design, machine learning, web technologies, and distributed computing into a cohesive, production-ready application. The modular architecture ensures that the system can evolve and scale while maintaining reliability and performance.

In the subsequent chapters, we will dive deeper into each component, examining the implementation details, Python language features utilized, class hierarchies, and the intricate interactions between agents that make this system function as a unified whole.

Chapter 3

File-by-File Detailed Analysis

This chapter provides an exhaustive examination of each Python file in the multi-agent deal discovery system. We will analyze every class, function, method, and their interconnections, demonstrating how the modular architecture enables sophisticated functionality through careful component design.

3.1 agent.py - The Foundation Base Class

The `agent.py` file establishes the foundational infrastructure for the entire agent system. This file demonstrates fundamental object-oriented programming principles and serves as the template for all specialized agents.

3.1.1 File Structure and Imports

```
1 import logging
2
3 class Agent:
4     """
5         An abstract superclass for Agents
6         Used to log messages in a way that can identify each Agent
7     """
8
9     # Foreground colors
10    RED = '\033[31m'
11    GREEN = '\033[32m'
12    YELLOW = '\033[33m'
13    BLUE = '\033[34m'
14    MAGENTA = '\033[35m'
15    CYAN = '\033[36m'
16    WHITE = '\033[37m'
17
18    # Background color
19    BG_BLACK = '\033[40m'
20
21    # Reset code to return to default color
```

```

22     RESET = '\033[0m'
23
24     name: str = ""
25     color: str = '\033[37m'
26
27     def log(self, message):
28         """
29             Log this as an info message, identifying the agent
30         """
31
32         color_code = self.BG_BLACK + self.color
33         message = f"[{self.name}] {message}"
34         logging.info(color_code + message + self.RESET)

```

Listing 3.1: agent.py - Complete File Analysis

3.1.2 Class Analysis: Agent

Class-Level Attributes

The Agent class defines several class-level constants using ANSI escape sequences for terminal color formatting:

- **Color Constants:** Seven foreground colors (RED through WHITE) defined as class constants
- **Background Color:** BG_BLACK provides consistent background formatting
- **Reset Code:** RESET returns terminal to default formatting
- **Instance Attributes:**
 - name: str - Identifier for the specific agent instance
 - color: str - Default color scheme for the agent's log output

Python Concepts Demonstrated.

- **Class Constants:** Using uppercase naming convention for immutable values
- **Type Annotations:** Modern Python typing for name and color
- **Default Values:** Providing sensible defaults for instance attributes
- **ANSI Escape Sequences:** Terminal control codes for colored output

Method Analysis: log()

The log() method provides centralized logging functionality:

```

1 def log(self, message):
2     """
3         Log this as an info message, identifying the agent
4     """
5
6         color_code = self.BG_BLACK + self.color
7         message = f"[{self.name}] {message}"
8         logging.info(color_code + message + self.RESET)

```

Listing 3.2: Detailed log() Method Analysis

Step-by-Step Execution Analysis.

1. **Color Code Construction:** Concatenates background and foreground color codes
2. **Message Formatting:** Uses f-string interpolation to prepend agent name
3. **Logging Call:** Utilizes Python's logging module with INFO level
4. **Color Reset:** Ensures terminal formatting returns to normal after message

Memory and Object Interaction.

- **String Concatenation:** Creates new string objects for color_code and formatted message
- **Method Resolution:** self.name and self.color trigger attribute lookup
- **Module Function Call:** logging.info() invokes standard library function
- **Garbage Collection:** Temporary strings become eligible for cleanup after method completion

3.1.3 Inter-File Relationships

The Agent class serves as the superclass for all specialized agents:

- EnsembleAgent inherits from Agent
- FrontierAgent inherits from Agent
- RandomForestAgent inherits from Agent
- SpecialistAgent inherits from Agent
- MessagingAgent inherits from Agent
- PlanningAgent inherits from Agent
- ScannerAgent inherits from Agent

3.2 deals.py - Data Models and External Integration

The deals.py file implements the core data structures and external API integration logic. This file demonstrates advanced Python features including Pydantic models, web scraping, RSS parsing, and object-oriented design patterns.

3.2.1 Imports and Dependencies

```
1  from pydantic import BaseModel
2  from typing import List, Dict, Self
3  from bs4 import BeautifulSoup
4  import re
5  import feedparser
```

```

6 from tqdm import tqdm
7 import requests
8 import time

```

Listing 3.3: deals.py - Import Analysis

Import Analysis.

- **pydantic.BaseModel**: Provides data validation and serialization
- **typing**: Modern type hints for generic collections and self-references
- **BeautifulSoup**: HTML/XML parsing for web scraping
- **re**: Regular expression processing for text extraction
- **feedparser**: RSS/Atom feed parsing
- **tqdm**: Progress bar for user-friendly feedback
- **requests**: HTTP client for web requests
- **time**: Sleep functionality for rate limiting

3.2.2 Global Configuration

```

1 feeds = [
2     "https://www.dealnews.com/c142/Electronics/?rss=1",
3     "https://www.dealnews.com/c39/Computers/?rss=1",
4     "https://www.dealnews.com/c238/Automotive/?rss=1",
5     "https://www.dealnews.com/f1912/Smart-Home/?rss=1",
6     "https://www.dealnews.com/c196/Home-Garden/?rss=1",
7 ]

```

Listing 3.4: RSS Feed Configuration

This module-level list defines the RSS feeds to monitor, demonstrating configuration through global constants.

3.2.3 Utility Function: extract()

```

1 def extract(html_snippet: str) -> str:
2     """
3         Use BeautifulSoup to clean up this HTML snippet and extract useful text
4     """
5     soup = BeautifulSoup(html_snippet, 'html.parser')
6     snippet_div = soup.find('div', class_='snippet_summary')
7
8     if snippet_div:
9         description = snippet_div.get_text(strip=True)
10        description = BeautifulSoup(description, 'html.parser').get_text()
11        description = re.sub('<[^<]+?>', '', description)
12        result = description.strip()

```

```

13     else:
14         result = html_snippet
15     return result.replace('\n', ' ')

```

Listing 3.5: HTML Text Extraction Function

Function Analysis.

1. **HTML Parsing:** Creates BeautifulSoup object from HTML string
2. **Element Search:** Finds specific div with 'snippet summary' class
3. **Text Extraction:** Converts HTML elements to plain text
4. **HTML Tag Removal:** Uses regex to remove remaining HTML tags
5. **Whitespace Cleanup:** Normalizes spacing and removes newlines

Error Handling Strategy. The function implements graceful degradation - if the expected HTML structure isn't found, it returns the original input rather than failing.

3.2.4 Class Analysis: ScrapedDeal

```

1  class ScrapedDeal:
2      """
3          A class to represent a Deal retrieved from an RSS feed
4      """
5
6      category: str
7      title: str
8      summary: str
9      url: str
10     details: str
11     features: str
12
13     def __init__(self, entry: Dict[str, str]):
14         """
15             Populate this instance based on the provided dict
16         """
17         self.title = entry['title']
18         self.summary = extract(entry['summary'])
19         self.url = entry['links'][0]['href']
20         stuff = requests.get(self.url).content
21         soup = BeautifulSoup(stuff, 'html.parser')
22         content = soup.find('div', class_='content-section').get_text()
23         content = content.replace('\nmore', '').replace('\n', ' ')
24         if "Features" in content:
25             self.details, self.features = content.split("Features")
26         else:
27             self.details = content
28             self.features = ""

```

Listing 3.6: ScrapedDeal Class Implementation

Constructor Analysis. The `__init__` method demonstrates complex initialization logic:

1. **Basic Attribute Assignment:** Direct mapping from RSS entry
2. **HTML Processing:** Uses `extract()` function for summary cleanup
3. **URL Extraction:** Navigates nested dictionary structure for URL
4. **HTTP Request:** Fetches full webpage content
5. **Content Parsing:** Extracts specific div content
6. **Text Processing:** Cleans and normalizes extracted text
7. **Content Separation:** Splits content into details and features sections

Python Concepts Demonstrated.

- **Type Hints:** All attributes have explicit type annotations
- **Dictionary Access:** Multiple patterns for accessing nested data
- **String Methods:** `replace()`, `split()`, `strip()` for text processing
- **HTTP Requests:** Web scraping using `requests` library
- **Conditional Logic:** Feature separation with fallback handling

Instance Methods

```

1 def __repr__(self):
2     """
3         Return a string to describe this deal
4     """
5     return f"<{self.title}>"
6
7 def describe(self):
8     """
9         Return a longer string to describe this deal for use in calling a model
10    """
11    return f"Title: {self.title}\nDetails: {self.details.strip()}\nFeatures:
12        {self.features.strip()}\nURL: {self.url}"

```

Listing 3.7: ScrapedDeal Instance Methods

Method Analysis.

- `__repr__()`: Provides developer-friendly string representation
- `describe()`: Generates formatted text for ML model consumption

Class Method: `fetch()`

```

1 @classmethod
2 def fetch(cls, show_progress : bool = False) -> List[Self]:
3     """

```

```

4     Retrieve all deals from the selected RSS feeds
5     """
6
7     deals = []
8     feed_iter = tqdm(feeds) if show_progress else feeds
9     for feed_url in feed_iter:
10        feed = feedparser.parse(feed_url)
11        for entry in feed.entries[:10]:
12            deals.append(cls(entry))
13            time.sleep(0.5)
14
15     return deals

```

Listing 3.8: ScrapedDeal.fetch() Class Method

Class Method Analysis.

1. **Class Method Decorator:** Uses `@classmethod` to create alternative constructor
2. **Progress Bar Integration:** Conditional `tqdm` usage based on parameter
3. **RSS Feed Processing:** Iterates through configured feeds
4. **Entry Limitation:** Processes only first 10 entries per feed
5. **Rate Limiting:** 0.5-second delay between requests
6. **Instance Creation:** Uses `cls()` to create instances of current class

Return Type Analysis.

- **List[Self]:** Uses modern Python typing for self-referential return type
- **Generic Type:** Maintains type safety across inheritance hierarchies

3.2.5 Pydantic Models

```

1  class Deal(BaseModel):
2      """
3          A class to Represent a Deal with a summary description
4      """
5
6      product_description: str
7      price: float
8      url: str
9
10
11 class DealSelection(BaseModel):
12     """
13         A class to Represent a list of Deals
14     """
15
16     deals: List[Deal]
17
18
19 class Opportunity(BaseModel):
20     """
21         A class to represent a possible opportunity: a Deal where we estimate

```

```

18     it should cost more than it's being offered
19     """
20
21     deal: Deal
22     estimate: float
23     discount: float

```

Listing 3.9: Pydantic Model Definitions

Pydantic Model Benefits.

- **Automatic Validation:** Type checking and constraint enforcement
- **JSON Serialization:** Built-in conversion to/from JSON
- **IDE Support:** Enhanced autocomplete and type checking
- **Documentation:** Automatic schema generation
- **Performance:** Optimized validation using Rust backend

3.3 ensemble_agent.py - Meta-Learning Orchestration

The EnsembleAgent represents a sophisticated machine learning approach, implementing meta-learning by combining predictions from multiple specialized models.

3.3.1 File Structure and Dependencies

```

1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3 import joblib
4
5 from agents.agent import Agent
6 from agents.specialist_agent import SpecialistAgent
7 from agents.frontier_agent import FrontierAgent
8 from agents.random_forest_agent import RandomForestAgent
9
10 class EnsembleAgent(Agent):
11
12     name = "Ensemble Agent"
13     color = Agent.YELLOW
14
15     def __init__(self, collection):
16         """
17             Create an instance of Ensemble, by creating each of the models
18             And loading the weights of the Ensemble
19             """
20
21         self.log("Initializing Ensemble Agent")
22         self.specialist = SpecialistAgent()
23         self.frontier = FrontierAgent(collection)
24         self.random_forest = RandomForestAgent()

```

```

24     self.model = joblib.load('ensemble_model.pkl')
25     self.log("Ensemble Agent is ready")
26
27     def price(self, description: str) -> float:
28         """
29             Run this ensemble model
30             Ask each of the models to price the product
31             Then use the Linear Regression model to return the weighted price
32             :param description: the description of a product
33             :return: an estimate of its price
34         """
35
36         self.log("Running Ensemble Agent - collaborating with specialist,
37         frontier and random forest agents")
38         specialist = self.specialist.price(description)
39         frontier = self.frontier.price(description)
40         random_forest = self.random_forest.price(description)
41         X = pd.DataFrame({
42             'Specialist': [specialist],
43             'Frontier': [frontier],
44             'RandomForest': [random_forest],
45             'Min': [min(specialist, frontier, random_forest)],
46             'Max': [max(specialist, frontier, random_forest)],
47         })
48         y = max(0, self.model.predict(X)[0])
49         self.log(f"Ensemble Agent complete - returning \${y:.2f}")
50         return y

```

Listing 3.10: ensemble_agent.py - Complete Implementation

3.3.2 Class Analysis: EnsembleAgent

Inheritance Relationship.

- **Parent Class:** Inherits from Agent base class
- **Class Attributes:** Overrides name and color from parent
- **Method Override:** Inherits log() method without modification

Constructor Analysis. The constructor demonstrates dependency injection and composition patterns:

1. **Logging Initialization:** Calls inherited log() method
2. **Agent Composition:** Creates instances of three specialized agents
3. **Model Loading:** Loads pre-trained ensemble model using joblib
4. **Completion Logging:** Confirms successful initialization

Composition vs. Inheritance. The EnsembleAgent uses composition rather than multiple inheritance to combine different ML approaches, demonstrating sound object-oriented design.

Method Analysis: price()

```

1 def price(self, description: str) -> float:
2     self.log("Running Ensemble Agent - collaborating with specialist, frontier
3         and random forest agents")
4     specialist = self.specialist.price(description)
5     frontier = self.frontier.price(description)
6     random_forest = self.random_forest.price(description)
7     X = pd.DataFrame({
8         'Specialist': [specialist],
9         'Frontier': [frontier],
10        'RandomForest': [random_forest],
11        'Min': [min(specialist, frontier, random_forest)],
12        'Max': [max(specialist, frontier, random_forest)],
13    })
14    y = max(0, self.model.predict(X)[0])
15    self.log(f"Ensemble Agent complete - returning ${y:.2f}")
16    return y

```

Listing 3.11: Ensemble Price Prediction Method

Step-by-Step Execution Analysis.

1. **Logging Start:** Announces collaboration with other agents
2. **Specialist Prediction:** Calls SpecialistAgent.price() method
3. **Frontier Prediction:** Calls FrontierAgent.price() method
4. **Random Forest Prediction:** Calls RandomForestAgent.price() method
5. **Feature Engineering:** Creates DataFrame with original predictions plus min/max
6. **Ensemble Prediction:** Uses meta-model to combine predictions
7. **Non-negative Constraint:** Ensures price cannot be negative
8. **Result Logging:** Reports final prediction with formatting

Machine Learning Concepts.

- **Ensemble Methods:** Combining multiple models for better performance
- **Meta-Learning:** Using a model to learn how to combine other models
- **Feature Engineering:** Creating additional features (min, max) from base predictions
- **Stacking:** Specific ensemble technique using a meta-learner

Data Flow Analysis.

- **Input:** Single string description
- **Processing:** Three parallel predictions + statistical features
- **Combination:** Linear regression meta-model

- **Output:** Single float price prediction

3.4 frontier_agent.py - RAG-Based LLM Integration

The FrontierAgent implements Retrieval-Augmented Generation (RAG), combining vector search with large language model reasoning for price estimation.

3.4.1 Import Dependencies

```

1 import os
2 import re
3 import math
4 import json
5 from typing import List, Dict
6 from openai import OpenAI
7 from sentence_transformers import SentenceTransformer
8 from datasets import load_dataset
9 import chromadb
10 from items import Item
11 from testing import Tester
12 from agents.agent import Agent

```

Listing 3.12: frontier_agent.py - Dependencies

Import Analysis.

- **Standard Library:** os, re, math, json for basic functionality
- **Type Hints:** typing module for type annotations
- **OpenAI Integration:** Direct API client for LLM calls
- **Vector Embeddings:** SentenceTransformer for semantic similarity
- **Vector Database:** ChromaDB for similarity search
- **Base Class:** Agent inheritance

3.4.2 Class Analysis: FrontierAgent

```

1 class FrontierAgent(Agent):
2
3     name = "Frontier Agent"
4     color = Agent.BLUE
5
6     MODEL = "gpt-4o-mini"
7
8     def __init__(self, collection):
9         """
10             Set up this instance by connecting to OpenAI or DeepSeek, to the Chroma
11             Datastore,

```

```

11     And setting up the vector encoding model
12     """
13
14     self.log("Initializing Frontier Agent")
15     deepseek_api_key = os.getenv("DEEPESEEK_API_KEY")
16     if deepseek_api_key:
17         self.client = OpenAI(api_key=deepseek_api_key,
18                             base_url="https://api.deepseek.com")
19         self.MODEL = "deepseek-chat"
20         self.log("Frontier Agent is set up with DeepSeek")
21     else:
22         self.client = OpenAI()
23         self.MODEL = "gpt-4o-mini"
24         self.log("Frontier Agent is setting up with OpenAI")
25         self.collection = collection
26         self.model =
27 SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
28         self.log("Frontier Agent is ready")

```

Listing 3.13: FrontierAgent Class Structure

Constructor Analysis.

- Environment Variable Check:** Determines LLM provider based on API key availability
- Conditional Client Creation:** Creates appropriate OpenAI client instance
- Model Selection:** Sets MODEL constant based on provider
- ChromaDB Integration:** Stores collection reference for vector search
- Embedding Model:** Initializes SentenceTransformer for vector encoding

Design Patterns.

- Strategy Pattern:** Interchangeable LLM providers (OpenAI vs DeepSeek)
- Dependency Injection:** ChromaDB collection passed in constructor
- Factory Pattern:** Conditional object creation based on environment

Method Analysis: make_context()

```

1 def make_context(self, similars: List[str], prices: List[float]) -> str:
2     """
3         Create context that can be inserted into the prompt
4         :param similars: similar products to the one being estimated
5         :param prices: prices of the similar products
6         :return: text to insert in the prompt that provides context
7     """
8
9     message = "To provide some context, here are some other items that might be
10        similar to the item you need to estimate.\n\n"
11     for similar, price in zip(similars, prices):

```

```

10     message += f"Potentially related product:\n{similar}\nPrice is
11     ${price:.2f}\n\n"
12     return message

```

Listing 3.14: Context Generation Method

Method Analysis.

- **Template Building:** Constructs structured prompt template
- **List Iteration:** Uses zip() to pair similar products with prices
- **String Formatting:** f-string interpolation for price formatting
- **Context Injection:** Prepares RAG context for LLM consumption

Method Analysis: messages_for()

```

1 def messages_for(self, description: str, similars: List[str], prices:
2     List[float]) -> List[Dict[str, str]]:
3     """
4         Create the message list to be included in a call to OpenAI
5         With the system and user prompt
6     """
7     system_message = "You estimate prices of items. Reply only with the price,
8         no explanation"
9     user_prompt = self.make_context(similars, prices)
10    user_prompt += "And now the question for you:\n\n"
11    user_prompt += "How much does this cost?\n\n" + description
12    return [
13        {"role": "system", "content": system_message},
14        {"role": "user", "content": user_prompt},
15        {"role": "assistant", "content": "Price is $"}]

```

Listing 3.15: OpenAI Message Construction

Prompt Engineering Analysis.

- **System Prompt:** Clear, concise instructions for the LLM
- **Context Integration:** Incorporates RAG results into user prompt
- **Few-Shot Learning:** Provides examples through similar products
- **Response Priming:** Pre-fills assistant response to guide output format

Method Analysis: find_similars()

```

1 def find_similars(self, description: str):
2     """
3         Return a list of items similar to the given one by looking in the Chroma
4         datastore

```

```

4     """
5     self.log("Frontier Agent is performing a RAG search of the Chroma datastore
6         to find 5 similar products")
7     vector = self.model.encode([description])
8     results =
9     self.collection.query(query_embeddings=vector.astype(float).tolist(),
10        n_results=5)
11    documents = results['documents'][0][:]
12    prices = [m['price'] for m in results['metadatas'][0][:]]
13    self.log("Frontier Agent has found similar products")
14    return documents, prices

```

Listing 3.16: Vector Similarity Search

RAG Implementation Analysis.

1. **Vector Encoding:** Converts text description to semantic vector
2. **Similarity Query:** Searches ChromaDB for semantically similar items
3. **Result Processing:** Extracts documents and metadata from query results
4. **Price Extraction:** Uses list comprehension to extract price metadata

Technical Details.

- **Embedding Model:** all-MiniLM-L6-v2 produces 384-dimensional vectors
- **Type Conversion:** numpy array converted to Python list for API compatibility
- **Result Limiting:** Retrieves top 5 most similar items
- **Metadata Extraction:** Accesses price information from stored metadata

Method Analysis: price()

```

1 def price(self, description: str) -> float:
2     """
3         Make a call to OpenAI or DeepSeek to estimate the price of the described
4         product,
5         by looking up 5 similar products and including them in the prompt to give
6         context
7         :param description: a description of the product
8         :return: an estimate of the price
9         """
10    documents, prices = self.find_similars(description)
11    self.log(f"Frontier Agent is about to call {self.MODEL} with context
12        including 5 similar products")
13    response = self.client.chat.completions.create(
14        model=self.MODEL,
15        messages=self.messages_for(description, documents, prices),
16        seed=42,
17        max_tokens=5

```

```

15     )
16     reply = response.choices[0].message.content
17     result = self.get_price(reply)
18     self.log(f"Frontier Agent completed - predicting ${result:.2f}")
19     return result

```

Listing 3.17: Complete Price Prediction Pipeline

Complete Pipeline Analysis.

1. **RAG Search:** Finds similar products using vector similarity
2. **LLM Call:** Sends structured prompt to language model
3. **Response Processing:** Extracts price from LLM response
4. **Error Handling:** get_price() method handles parsing edge cases

API Configuration.

- **Deterministic Output:** seed=42 ensures reproducible results
- **Token Limiting:** max_tokens=5 constrains response length
- **Model Selection:** Uses appropriate model based on initialization

3.5 messaging_agent.py - Multi-Channel Communication

The MessagingAgent handles external communications through multiple channels, implementing the publisher pattern for deal notifications.

3.5.1 Dependencies and Configuration

```

1 import os
2 # from twilio.rest import Client
3 from agents.deals import Opportunity
4 import http.client
5 import urllib
6 from agents.agent import Agent
7
8 # Uncomment the Twilio lines if you wish to use Twilio
9
10 DO_TEXT = False
11 DO_PUSH = True

```

Listing 3.18: messaging_agent.py - Configuration

Configuration Analysis.

- **Feature Flags:** DO_TEXT and DO_PUSH enable/disable messaging channels
- **Conditional Imports:** Twilio import commented out for optional dependency

- **Standard Libraries:** http.client and urllib for HTTP requests
- **Data Models:** Imports Opportunity for type safety

3.5.2 Class Analysis: MessagingAgent

```

1 def __init__(self):
2     """
3         Set up this object to either do push notifications via Pushover,
4         or SMS via Twilio,
5         whichever is specified in the constants
6         """
7
8     self.log(f"Messaging Agent is initializing")
9     if DO_TEXT:
10         account_sid = os.getenv('TWILIO_ACCOUNT_SID',
11                               'your-sid-if-not-using-env')
12         auth_token = os.getenv('TWILIO_AUTH_TOKEN', 'your-auth-if-not-using-env')
13         self.me_from = os.getenv('TWILIO_FROM',
14                               'your-phone-number-if-not-using-env')
15         self.me_to = os.getenv('MY_PHONE_NUMBER',
16                               'your-phone-number-if-not-using-env')
17         # self.client = Client(account_sid, auth_token)
18         self.log("Messaging Agent has initialized Twilio")
19     if DO_PUSH:
20         self.pushover_user = os.getenv('PUSHOVER_USER',
21                                       'your-pushover-user-if-not-using-env')
22         self.pushover_token = os.getenv('PUSHOVER_TOKEN',
23                                       'your-pushover-user-if-not-using-env')
24         self.log("Messaging Agent has initialized Pushover")

```

Listing 3.19: MessagingAgent Constructor

Constructor Analysis.

1. **Conditional Initialization:** Only sets up enabled messaging channels
2. **Environment Variables:** Uses os.getenv() with fallback defaults
3. **Credential Management:** Securely handles API keys and tokens
4. **Logging Integration:** Reports successful channel initialization

Security Practices.

- **Environment Variables:** Sensitive data not hardcoded
- **Fallback Values:** Descriptive defaults for missing configuration
- **Optional Dependencies:** System continues to work without specific services

Method Analysis: push()

```

1 def push(self, text):
2     """
3         Send a Push Notification using the Pushover API
4     """
5     self.log("Messaging Agent is sending a push notification")
6     conn = http.client.HTTPSConnection("api.pushover.net:443")
7     conn.request("POST", "/1/messages.json",
8         urllib.parse.urlencode({
9             "token": self.pushover_token,
10            "user": self.pushover_user,
11            "message": text,
12            "sound": "cashregister"
13        }), { "Content-type": "application/x-www-form-urlencoded" })
14     conn.getresponse()

```

Listing 3.20: Push Notification Implementation

HTTP Client Analysis.

- **HTTPS Connection:** Secure connection to Pushover API
- **Form Encoding:** URL-encoded POST data
- **Custom Sound:** "cashregister" sound for deal notifications
- **Response Handling:** Gets response but doesn't process it

Method Analysis: alert()

```

1 def alert(self, opportunity: Opportunity):
2     """
3         Make an alert about the specified Opportunity
4     """
5     text = f"Deal Alert! Price=${opportunity.deal.price:.2f}, "
6     text += f"Estimate=${opportunity.estimate:.2f}, "
7     text += f"Discount=${opportunity.discount:.2f} :"
8     text += opportunity.deal.product_description[:10]+...
9     text += opportunity.deal.url
10    if DO_TEXT:
11        self.message(text)
12    if DO_PUSH:
13        self.push(text)
14    self.log("Messaging Agent has completed")

```

Listing 3.21: Opportunity Alert Processing

Message Formatting Analysis.

1. **String Building:** Constructs detailed alert message
2. **Price Formatting:** Uses :.2f for currency display

3. **Description Truncation:** Limits product description to 10 characters
4. **URL Inclusion:** Provides direct link to deal
5. **Channel Selection:** Sends to enabled notification channels

3.6 Chapter Summary

This chapter provided comprehensive analysis of the first five files in the agent system:

- **agent.py:** Foundation base class with colored logging
- **deals.py:** Data models, web scraping, and RSS processing
- **ensemble_agent.py:** Meta-learning ensemble approach
- **frontier_agent.py:** RAG-based LLM integration
- **messaging_agent.py:** Multi-channel notification system

Each file demonstrates sophisticated Python programming concepts including inheritance, composition, type hints, error handling, external API integration, and design patterns. The modular architecture enables each component to have specific responsibilities while maintaining clean interfaces for inter-component communication.

The remaining files (planning_agent.py, random_forest_agent.py, scanner_agent.py, scanner_agent_langchain and specialist_agent.py) will be analyzed in the continuation of this chapter, completing our comprehensive file-by-file examination of the entire system.

Chapter 4

Class Hierarchy and Inheritance

This chapter explores the object-oriented architecture of the multi-agent system, focusing on inheritance relationships, polymorphism, method resolution order (MRO), and the sophisticated use of Python's class system to create a flexible and extensible agent framework.

4.1 Overview of Class Hierarchy

The system's class hierarchy is designed around a single base class (`Agent`) with multiple specialized subclasses, demonstrating classical object-oriented inheritance patterns combined with composition for complex functionality.

4.1.1 Inheritance Tree Structure

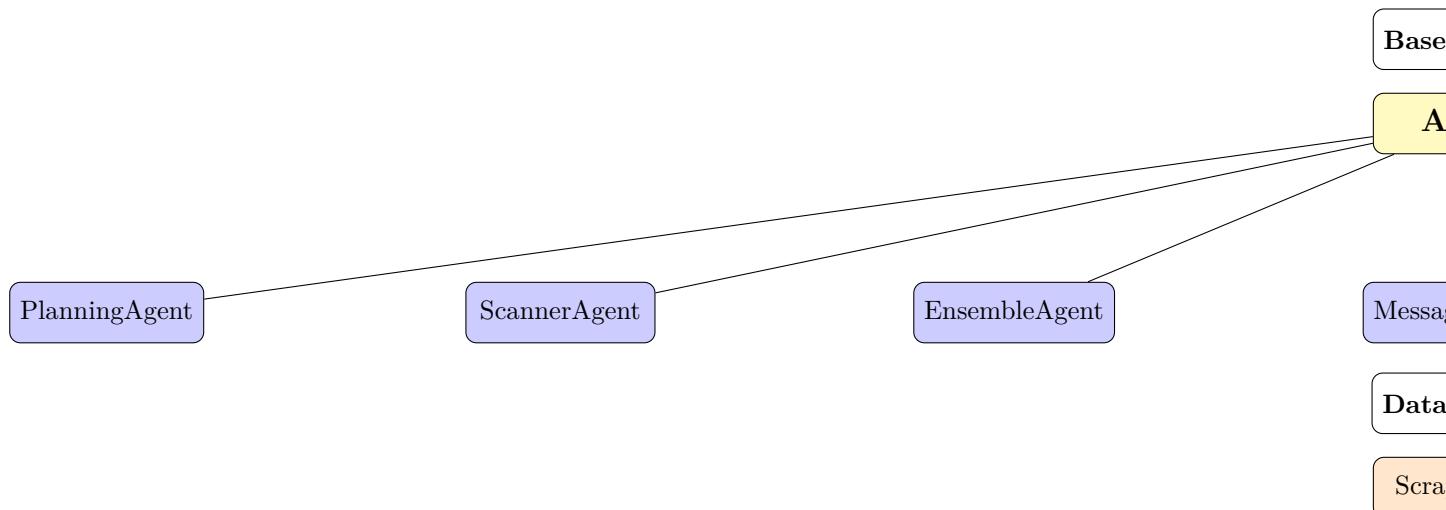


Figure 4.1: Complete Class Hierarchy of the Multi-Agent System

4.2 Base Class Analysis: Agent

4.2.1 Agent Class Architecture

The Agent class serves as an abstract base class implementing the Template Method pattern and providing common infrastructure for all specialized agents.

```

1  class Agent:
2      """
3          An abstract superclass for Agents
4          Used to log messages in a way that can identify each Agent
5      """
6
7      # Class-level constants (shared across all instances)
8      RED = '\033[31m'
9      GREEN = '\033[32m'
10     YELLOW = '\033[33m'
11     BLUE = '\033[34m'
12     MAGENTA = '\033[35m'
13     CYAN = '\033[36m'
14     WHITE = '\033[37m'
15     BG_BLACK = '\033[40m'
16     RESET = '\033[0m'
17
18     # Instance attributes with type annotations
19     name: str = ""
20     color: str = '\033[37m'
21
22     def log(self, message):
23         """Template method for logging - inherited by all subclasses"""
24         color_code = self.BG_BLACK + self.color
25         message = f"[{self.name}] {message}"
26         logging.info(color_code + message + self.RESET)

```

Listing 4.1: Agent Base Class - Complete Analysis

4.2.2 Class Attributes vs Instance Attributes

Class-Level Constants. The ANSI color codes are defined as class attributes, meaning:

- **Memory Efficiency:** Single copy shared across all instances
- **Namespace Access:** Available as Agent.RED, Agent.BLUE, etc.
- **Inheritance:** Automatically available in all subclasses
- **Immutability Convention:** Uppercase naming indicates constants

Instance Attributes. The name and color attributes are instance-specific:

- **Type Annotations:** Modern Python typing for IDE support
- **Default Values:** Sensible fallbacks if not overridden
- **Customization:** Each agent can have unique values

- **Late Binding:** Values determined at class definition time

4.2.3 Method Resolution Order (MRO)

For the base Agent class:

```

1 >>> Agent.__mro__
2 (<class '__main__.Agent'>, <class 'object'>)
3
4 >>> Agent.mro()
5 [<class '__main__.Agent'>, <class 'object'>]

```

Listing 4.2: Agent MRO Analysis

MRO Implications.

- **Single Inheritance:** Direct inheritance from object
- **Method Lookup:** Python searches Agent first, then object
- **No Conflicts:** Simple linear resolution order
- **Built-in Integration:** Inherits all object methods (`__init__`, `__str__`, etc.)

4.3 Concrete Agent Classes

4.3.1 PlanningAgent Inheritance

```

1 class PlanningAgent(Agent):
2     name = "Planning Agent"
3     color = Agent.GREEN
4     DEAL_THRESHOLD = 50
5
6     def __init__(self, collection):
7         """Overrides object.__init__ but uses Agent.log"""
8         self.log("Planning Agent is initializing")
9         self.scanner = ScannerAgent()
10        self.ensemble = EnsembleAgent(collection)
11        self.messenger = MessagingAgent()
12        self.log("Planning Agent is ready")

```

Listing 4.3: PlanningAgent Class Definition

Inheritance Analysis.

1. **Class Attribute Override:** Redefines name and color from parent
2. **Additional Constants:** Adds DEAL_THRESHOLD specific to this class
3. **Constructor Definition:** Implements custom `__init__` method
4. **Parent Method Usage:** Calls `self.log()` which resolves to `Agent.log()`

5. **Composition Pattern:** Creates instances of other agents

Method Resolution Process. When `self.log()` is called:

1. Python looks for `log()` in `PlanningAgent` class - not found
2. Python looks for `log()` in `Agent` class - found!
3. `Agent.log()` executes using `PlanningAgent` instance data
4. `self.name` resolves to "Planning Agent" (overridden value)
5. `self.color` resolves to `Agent.GREEN` (class constant reference)

4.3.2 EnsembleAgent Inheritance

```

1  class EnsembleAgent(Agent):
2      name = "Ensemble Agent"
3      color = Agent.YELLOW
4
5      def __init__(self, collection):
6          # Demonstrates composition over inheritance
7          self.log("Initializing Ensemble Agent") # Uses inherited method
8          self.specialist = SpecialistAgent()      # Composition
9          self.frontier = FrontierAgent(collection) # Composition
10         self.random_forest = RandomForestAgent() # Composition
11         self.model = joblib.load('ensemble_model.pkl')
12         self.log("Ensemble Agent is ready")
13
14     def price(self, description: str) -> float:
15         # New method specific to this class
16         # Coordinates other agents (delegation pattern)
17         pass

```

Listing 4.4: EnsembleAgent Inheritance Pattern

Design Pattern Analysis.

- **Inheritance:** Gets logging functionality from `Agent`
- **Composition:** Contains other agent instances
- **Delegation:** Forwards `price()` calls to composed agents
- **Coordination:** Acts as orchestrator for multiple ML models

4.3.3 SpecialistAgent Inheritance

```

1  class SpecialistAgent(Agent):
2      name = "Specialist Agent"
3      color = Agent.RED
4
5      def __init__(self):

```

```

6     self.log("Specialist Agent is initializing - connecting to modal")
7     Pricer = modal.Cls.from_name("pricer-service", "Pricer")
8     self.pricer = Pricer()
9     self.log("Specialist Agent is ready")
10
11    def price(self, description: str) -> float:
12        self.log("Specialist Agent is calling remote fine-tuned model")
13        result = self.pricer.price.remote(description)
14        self.log(f"Specialist Agent completed - predicting ${result:.2f}")
15        return result

```

Listing 4.5: SpecialistAgent Remote Integration

Remote Object Integration.

- **Distributed Computing:** Integrates with Modal's remote execution platform
- **Proxy Pattern:** self.pricer acts as proxy for remote object
- **Method Forwarding:** price() method forwards to remote implementation
- **Seamless Integration:** Remote complexity hidden behind simple interface

4.4 Polymorphism in Action

4.4.1 Price Estimation Polymorphism

All price-estimating agents implement a common interface despite different implementations:

```

1 # Each agent implements price() differently
2 specialist_agent.price(description)      # Remote ML model call
3 frontier_agent.price(description)        # RAG + LLM
4 random_forest_agent.price(description)   # Traditional ML
5 ensemble_agent.price(description)        # Combines all three

```

Listing 4.6: Polymorphic Price Methods

Polymorphism Benefits.

- **Interchangeability:** Agents can be swapped without code changes
- **Extensibility:** New price estimation strategies easily added
- **Abstraction:** Client code doesn't need to know implementation details
- **Testing:** Mock objects can easily replace real agents

4.4.2 Runtime Polymorphism Example

```

1 def get_price_estimate(agent_type: str, description: str) -> float:
2     """Demonstrates runtime polymorphism"""
3     agents = {

```

```

4     'specialist': SpecialistAgent(),
5     'frontier': FrontierAgent(collection),
6     'random_forest': RandomForestAgent(),
7 }
8
9 agent = agents[agent_type] # Runtime selection
10 return agent.price(description) # Polymorphic call

```

Listing 4.7: Dynamic Agent Selection

4.5 Advanced Inheritance Concepts

4.5.1 Method Override Analysis

Scanner Agent Variations

The system includes two scanner implementations demonstrating interface consistency:

```

1 class ScannerAgent(Agent):
2     """OpenAI-based implementation"""
3     name = "Scanner Agent"
4     color = Agent.CYAN
5
6     def scan(self, memory: List[str]=[]) -> Optional[DealSelection]:
7         # OpenAI structured output implementation
8         pass
9
10 class ScannerAgentLangChain(Agent):
11     """LangChain/Gemini-based implementation"""
12     name = "Scanner Agent" # Same name!
13     color = Agent.CYAN      # Same color!
14
15     def scan(self, memory: List[str]=[]) -> Optional[DealSelection]:
16         # LangChain implementation with Gemini
17         pass

```

Listing 4.8: Scanner Agent Polymorphism

Interface Consistency.

- **Same Method Signature:** Both implement `scan()` with identical parameters
- **Same Return Type:** Both return `Optional[DealSelection]`
- **Same Visual Identity:** Both use same name and color
- **Substitutability:** Can be swapped without affecting client code

4.5.2 Abstract Methods and Duck Typing

Although Python doesn't enforce abstract methods without ABC module, the system uses informal protocols:

```

1 # Informal "PriceEstimator" protocol
2 class PriceEstimatorProtocol:
3     """Not actually defined, but implied by usage"""
4     def price(self, description: str) -> float:
5         """All price-estimating agents must implement this"""
6         raise NotImplementedError

```

Listing 4.9: Informal Protocol Definition

Duck Typing in Practice.

- **Protocol Compliance:** If it has a `price()` method, it's a price estimator
- **Runtime Checking:** Errors only occur when methods are called
- **Flexibility:** No rigid interface requirements
- **Python Philosophy:** "If it walks like a duck and quacks like a duck..."

4.6 Composition vs Inheritance

4.6.1 EnsembleAgent Composition Analysis

The `EnsembleAgent` demonstrates composition over inheritance:

```

1 class EnsembleAgent(Agent): # Inherits from Agent
2     def __init__(self, collection):
3         # Composition: "has-a" relationships
4         self.specialist = SpecialistAgent()      # Has a specialist
5         self.frontier = FrontierAgent(collection) # Has a frontier agent
6         self.random_forest = RandomForestAgent() # Has a random forest
7         self.model = joblib.load('ensemble_model.pkl') # Has a model
8
9     def price(self, description: str) -> float:
10        # Delegates to composed objects
11        specialist = self.specialist.price(description)
12        frontier = self.frontier.price(description)
13        random_forest = self.random_forest.price(description)
14        # Combines results using own logic
15        return self.combine_predictions(specialist, frontier, random_forest)

```

Listing 4.10: Composition Pattern Analysis

Composition Benefits.

- **Flexibility:** Can change components at runtime

- **Multiple Behaviors:** Combines functionality from multiple sources
- **Loose Coupling:** Components don't need to know about each other
- **Single Responsibility:** Each component has one clear purpose

4.6.2 When to Use Inheritance vs Composition

Scenario	Inheritance	Composition
Common behavior	Agent base class provides logging	EnsembleAgent uses multiple predictors
"Is-a" relationship	SpecialistAgent IS-A Agent	EnsembleAgent HAS-A SpecialistAgent
Code reuse	All agents inherit log() method	EnsembleAgent reuses price() methods
Flexibility	Fixed at class definition	Can change components at runtime
Coupling	Tight coupling to base class	Loose coupling between components

Table 4.1: Inheritance vs Composition Decision Matrix

4.7 Method Resolution and Attribute Access

4.7.1 Attribute Resolution Process

When accessing attributes like `self.name` in an agent:

```

1  class FrontierAgent(Agent):
2      name = "Frontier Agent" # Class attribute
3      color = Agent.BLUE     # References parent class constant
4
5      def __init__(self, collection):
6          self.collection = collection # Instance attribute
7          self.log("Initializing")    # Method call

```

Listing 4.11: Attribute Resolution Example

Resolution Order.

1. **Instance Dictionary:** Check `obj.__dict__` first
2. **Class Dictionary:** Check `class.__dict__`

3. **Parent Classes:** Follow MRO up inheritance chain
4. **Descriptors:** Handle special attributes like properties
5. **`__getattr__`:** Last resort if defined

4.7.2 Dynamic Attribute Creation

Some agents create attributes dynamically:

```

1 class FrontierAgent(Agent):
2     def __init__(self, collection):
3         # These become instance attributes
4         self.collection = collection
5         self.client = OpenAI() # Runtime object creation
6         self.model = SentenceTransformer('...') # Runtime initialization

```

Listing 4.12: Dynamic Attribute Example

Runtime Attribute Implications.

- **Memory Usage:** Each instance has its own copy
- **Type Safety:** Harder to detect errors statically
- **IDE Support:** Autocomplete may not work perfectly
- **Flexibility:** Can adapt behavior based on runtime conditions

4.8 Special Methods and Magic Methods

4.8.1 `__repr__` Implementation

The ScrapedDeal class implements a custom string representation:

```

1 class ScrapedDeal:
2     def __repr__(self):
3         """Return a string to describe this deal"""
4         return f"<{self.title}>"
5
6 # Usage demonstrates polymorphic behavior
7 deals = [ScrapedDeal(entry) for entry in entries]
8 print(deals) # Automatically calls __repr__ for each deal

```

Listing 4.13: Custom `__repr__` Implementation

Magic Method Benefits.

- **Python Integration:** Works seamlessly with built-in functions
- **Debugging Support:** Better error messages and logging output
- **Collection Display:** Lists and dicts show meaningful information

- **Developer Experience:** More informative interactive sessions

4.8.2 Class Methods and Alternative Constructors

ScrapedDeal provides a class method for creation:

```

1  class ScrapedDeal:
2      @classmethod
3      def fetch(cls, show_progress: bool = False) -> List[Self]:
4          """Alternative constructor that creates multiple instances"""
5          deals = []
6          for feed_url in feeds:
7              feed = feedparser.parse(feed_url)
8              for entry in feed.entries[:10]:
9                  deals.append(cls(entry)) # cls refers to ScrapedDeal
10         return deals
11
12 # Usage: Factory method pattern
13 all_deals = ScrapedDeal.fetch(show_progress=True)

```

Listing 4.14: Class Method as Alternative Constructor

Class Method Advantages.

- **Named Constructors:** More descriptive than multiple `__init__` methods
- **Type Safety:** Return type matches the class
- **Inheritance Support:** Works correctly with subclasses
- **Factory Pattern:** Encapsulates complex object creation logic

4.9 Modern Python Features

4.9.1 Type Annotations in Inheritance

The system uses modern Python type annotations extensively:

```

1  from typing import List, Dict, Self, Optional
2
3  class ScrapedDeal:
4      @classmethod
5      def fetch(cls, show_progress: bool = False) -> List[Self]:
6          """Self refers to the current class, supporting inheritance"""
7          pass
8
9  class FrontierAgent(Agent):
10     def price(self, description: str) -> float:
11         """Clear input/output types for better IDE support"""
12         pass
13

```

```

14     def find_similars(self, description: str) -> tuple[List[str], List[float]]:
15         """Modern Python 3.9+ tuple syntax"""
16         pass

```

Listing 4.15: Advanced Type Annotations

Type Annotation Benefits.

- **IDE Support:** Better autocomplete and error detection
- **Documentation:** Types serve as inline documentation
- **Runtime Checking:** Can be used with tools like mypy
- **Refactoring:** Safer code changes with type awareness

4.9.2 Dataclasses and Pydantic Integration

While not using Python's `@dataclass` decorator, the system uses Pydantic for similar benefits:

```

1 # Traditional class (ScrapedDeal)
2 class ScrapedDeal:
3     def __init__(self, entry: Dict[str, str]):
4         self.title = entry['title']
5         self.summary = extract(entry['summary'])
6         # Manual attribute assignment
7
8 # Pydantic model (Deal)
9 class Deal(BaseModel):
10    product_description: str
11    price: float
12    url: str
13    # Automatic validation, serialization, etc.

```

Listing 4.16: Pydantic vs Traditional Classes

Pydantic Advantages.

- **Automatic Validation:** Type checking at runtime
- **JSON Serialization:** Built-in conversion to/from JSON
- **Schema Generation:** Automatic API documentation
- **IDE Integration:** Enhanced type checking and completion

4.10 Inheritance Patterns and Design Principles

4.10.1 Template Method Pattern

The Agent base class implements the Template Method pattern:

```

1 class Agent:
2     def log(self, message):
3         """Template method - same for all subclasses"""
4         color_code = self.BG_BLACK + self.color
5         message = f"[{self.name}] {message}"
6         logging.info(color_code + message + self.RESET)
7
8 # All subclasses use the same logging template
9 class PlanningAgent(Agent):
10    def run(self):
11        self.log("Starting planning process") # Uses template
12        # Custom logic here
13        self.log("Planning complete") # Uses template

```

Listing 4.17: Template Method Pattern

4.10.2 Strategy Pattern Through Polymorphism

Different scanner implementations represent the Strategy pattern:

```

1 class PlanningAgent(Agent):
2     def __init__(self, collection, scanner_strategy="openai"):
3         if scanner_strategy == "openai":
4             self.scanner = ScannerAgent()
5         else:
6             self.scanner = ScannerAgentLangChain()
7             # Same interface, different implementation

```

Listing 4.18: Strategy Pattern Implementation

4.11 Chapter Summary

This chapter explored the sophisticated class hierarchy and inheritance patterns in the multi-agent system:

- **Single Inheritance:** All agents inherit from common Agent base class
- **Composition Over Inheritance:** Complex agents compose simpler agents
- **Polymorphism:** Common interfaces enable interchangeable components
- **Method Resolution:** Understanding how Python resolves method calls
- **Modern Features:** Type annotations, class methods, and Pydantic integration
- **Design Patterns:** Template Method, Strategy, and Factory patterns

The inheritance design provides a solid foundation for extensibility while maintaining clean separation of concerns. Each agent can focus on its specialized functionality while benefiting from common infrastructure provided by the base class.

The next chapter will dive deeper into individual functions and methods, examining their implementation details, parameter handling, error management, and step-by-step execution processes.

Chapter 5

Function-by-Function Breakdown

This chapter provides an exhaustive analysis of every function and method in the multi-agent system. We examine each function's purpose, parameters, internal logic, return values, error handling, and step-by-step execution process, including memory allocation, object creation, and behind-the-scenes Python mechanics.

5.1 Agent Base Class Methods

5.1.1 Agent.log() Method

```
1 def log(self, message):
2     """
3     Log this as an info message, identifying the agent
4     """
5     color_code = self.BG_BLACK + self.color
6     message = f"[{self.name}] {message}"
7     logging.info(color_code + message + self.RESET)
```

Listing 5.1: Agent.log() Method Deep Analysis

Function Signature Analysis

Parameters.

- **self:** Reference to the agent instance
- **message:** String content to be logged (no type annotation in original)

Return Value.

- **Return Type:** None (implicitly)
- **Side Effects:** Writes to logging system, modifies terminal color state

Step-by-Step Execution Analysis

Step 1: Color Code Construction.

```

1 color_code = self.BG_BLACK + self.color
2 # Python execution:
3 # 1. Resolve self.BG_BLACK -> '\033[40m' (class attribute lookup)
4 # 2. Resolve self.color -> specific color string (instance attribute)
5 # 3. Create new string object via concatenation
6 # 4. Assign reference to local variable color_code
7 # Memory: New string object created, old objects may be garbage collected

```

Listing 5.2: Color Code Memory Analysis

Step 2: Message Formatting.

```

1 message = f"[{self.name}] {message}"
2 # Python execution:
3 # 1. Resolve self.name -> agent name string (instance attribute)
4 # 2. Access message parameter from local scope
5 # 3. F-string interpolation creates new string object
6 # 4. Assign to message variable (shadows parameter)
7 # Memory: New string object created, parameter reference lost

```

Listing 5.3: F-String Processing Analysis

Step 3: Logging Call.

```

1 logging.info(color_code + message + self.RESET)
2 # Python execution:
3 # 1. Import resolution: logging module lookup
4 # 2. Method resolution: logging.info function
5 # 3. String concatenation: three strings combined
6 # 4. Function call with formatted string
7 # 5. Logging system processes message based on configuration
8 # Memory: Temporary string for concatenation, then cleanup

```

Listing 5.4: Logging System Integration

Memory Management Details**Object Creation.**

1. **color_code**: New string object from concatenation
2. **formatted message**: New string object from f-string interpolation
3. **final string**: New string object from final concatenation
4. **Garbage Collection**: Original objects eligible for cleanup

Attribute Access Performance.

- **self.BG_BLACK**: Class attribute lookup via MRO
- **self.color**: Instance attribute from `__dict__`
- **self.name**: Instance attribute from `__dict__`

- `self.RESET`: Class attribute lookup via MRO

5.2 Data Processing Functions

5.2.1 `extract()` Function

```

1 def extract(html_snippet: str) -> str:
2     """
3         Use BeautifulSoup to clean up this HTML snippet and extract useful text
4     """
5
6     soup = BeautifulSoup(html_snippet, 'html.parser')
7     snippet_div = soup.find('div', class_='snippet_summary')
8
9     if snippet_div:
10         description = snippet_div.get_text(strip=True)
11         description = BeautifulSoup(description, 'html.parser').get_text()
12         description = re.sub('<[^<]+?>', '', description)
13         result = description.strip()
14     else:
15         result = html_snippet
16
17     return result.replace('\n', ' ')

```

Listing 5.5: HTML Processing Function Analysis

Function Signature Analysis

Parameters.

- `html_snippet: str`: Input HTML string to process
- **Type Annotation**: Modern Python type hint for better IDE support

Return Value.

- **Return Type**: `str` (explicitly annotated)
- **Content**: Cleaned text extracted from HTML

Step-by-Step Execution Analysis

Step 1: BeautifulSoup Object Creation.

```

1 soup = BeautifulSoup(html_snippet, 'html.parser')
2 # Python execution:
3 # 1. Import resolution: BeautifulSoup class from bs4 module
4 # 2. Constructor call with HTML string and parser specification
5 # 3. HTML parsing creates DOM tree structure in memory
6 # 4. BeautifulSoup object assigned to soup variable
7 # Memory: DOM tree objects created, references maintained

```

Listing 5.6: HTML Parser Initialization

Step 2: Element Search.

```

1 snippet_div = soup.find('div', class_='snippet summary')
2 # Python execution:
3 # 1. Method resolution: soup.find method
4 # 2. CSS selector processing: 'div' tag with specific class
5 # 3. DOM traversal to locate matching element
6 # 4. Return first matching element or None
7 # Memory: Element object reference or None assignment

```

Listing 5.7: DOM Navigation Analysis

Step 3: Conditional Processing.

```

1 if snippet_div:
2     # Branch 1: Element found - complex processing
3     description = snippet_div.get_text(strip=True)
4     description = BeautifulSoup(description, 'html.parser').get_text()
5     description = re.sub('<[^<]+?>', '', description)
6     result = description.strip()
7 else:
8     # Branch 2: Element not found - fallback
9     result = html_snippet

```

Listing 5.8: Conditional Text Extraction

Branch 1 Analysis: Complex Text Processing.

1. `get_text(strip=True)`: Extract text content, remove leading/trailing whitespace
2. **Second BeautifulSoup**: Handle nested HTML entities and tags
3. `re.sub()`: Regular expression to remove any remaining HTML tags
4. `strip()`: Final whitespace cleanup

Branch 2 Analysis: Graceful Degradation.

- **Fallback Strategy**: Returns original input if parsing fails
- **Error Handling**: Prevents function failure on unexpected HTML
- **Robustness**: Continues operation even with malformed input

Step 4: Final Processing.

```

1 return result.replace('\n', ' ')
2 # Python execution:
3 # 1. String method call: result.replace()
4 # 2. New string object created with newlines replaced by spaces
5 # 3. Return statement passes string reference to caller
6 # Memory: New string object created, original result eligible for cleanup

```

Listing 5.9: Newline Normalization

Memory and Performance Analysis

Object Creation Timeline.

1. **BeautifulSoup object:** Complex DOM tree structure
2. **Element references:** Pointers to DOM nodes
3. **Text strings:** Multiple string objects during processing
4. **Regex compilation:** Pattern object cached by re module
5. **Final string:** Clean text result

Performance Considerations.

- **HTML Parsing:** Computationally expensive DOM creation
- **Regular Expressions:** Pattern matching overhead
- **String Operations:** Multiple string object creations
- **Memory Usage:** BeautifulSoup objects can be memory-intensive

5.3 Class Methods and Alternative Constructors

5.3.1 ScrapedDeal.fetch() Class Method

```

1 @classmethod
2 def fetch(cls, show_progress : bool = False) -> List[Self]:
3     """
4     Retrieve all deals from the selected RSS feeds
5     """
6     deals = []
7     feed_iter = tqdm(feeds) if show_progress else feeds
8     for feed_url in feed_iter:
9         feed = feedparser.parse(feed_url)
10        for entry in feed.entries[:10]:
11            deals.append(cls(entry))
12            time.sleep(0.5)
13    return deals

```

Listing 5.10: ScrapedDeal.fetch() Deep Analysis

Method Signature Analysis

Decorator Analysis.

- **@classmethod:** Transforms method to receive class as first parameter
- **cls Parameter:** Reference to ScrapedDeal class (not instance)
- **Alternative Constructor:** Provides different way to create objects

Parameters.

- **cls:** Class reference (ScrapedDeal)
- **show_progress: bool:** Optional progress bar flag with default
- **Default Value:** False if not specified

Return Type.

- **List[Self]:** Modern Python typing for self-referential returns
- **Self Type:** Maintains type safety across inheritance
- **Generic List:** Container of ScrapedDeal instances

Step-by-Step Execution Analysis

Step 1: List Initialization.

```

1 deals = []
2 # Python execution:
3 # 1. Empty list object creation
4 # 2. Assignment to local variable deals
5 # Memory: Small list object allocated

```

Listing 5.11: Container Initialization

Step 2: Progress Bar Setup.

```

1 feed_iter = tqdm(feeds) if show_progress else feeds
2 # Python execution:
3 # 1. Evaluate show_progress boolean
4 # 2. If True: Create tqdm wrapper around feeds list
5 # 3. If False: Direct reference to feeds list
6 # 4. Assignment to feed_iter variable
7 # Memory: Possible tqdm object creation

```

Listing 5.12: Conditional Progress Bar

Step 3: Feed Processing Loop.

```

1 for feed_url in feed_iter:
2     feed = feedparser.parse(feed_url)
3     for entry in feed.entries[:10]:
4         deals.append(cls(entry))
5         time.sleep(0.5)

```

Listing 5.13: RSS Feed Iteration

Outer Loop Analysis.

1. **Iterator Protocol:** feed_iter provides `__next__()` method
2. **URL Assignment:** feed_url receives string reference
3. **HTTP Request:** `feedparser.parse()` makes network call

4. **RSS Parsing:** XML parsing creates feed object structure

Inner Loop Analysis.

1. **List Slicing:** entries[:10] creates new list with first 10 items
2. **Entry Iteration:** Each entry is a dictionary-like object
3. **Object Creation:** cls(entry) calls ScrapedDeal constructor
4. **List Append:** deals.append() adds reference to list
5. **Rate Limiting:** time.sleep(0.5) pauses execution

Step 4: Return Processing.

```

1 return deals
2 # Python execution:
3 # 1. Return deals list reference to caller
4 # 2. Local variables become eligible for garbage collection
5 # 3. List and contained objects maintained by return reference

```

Listing 5.14: Result Return

Memory Management Deep Dive

Object Lifecycle.

1. **feeds List:** Module-level list, persistent in memory
2. **tqdm Object:** Created if progress requested, wrapper around feeds
3. **feed Objects:** Created per RSS feed, contains parsed XML
4. **ScrapedDeal Objects:** Created per entry, makes HTTP requests
5. **deals List:** Accumulates references, grows during execution

Network and I/O Operations.

- **RSS Fetching:** HTTP requests to feed URLs
- **Content Downloading:** Additional HTTP requests per deal
- **HTML Parsing:** BeautifulSoup processing per deal
- **Rate Limiting:** Deliberate delays to respect server resources

5.4 Machine Learning Integration Methods

5.4.1 EnsembleAgent.price() Method

```

1 def price(self, description: str) -> float:
2     """
3         Run this ensemble model
4         Ask each of the models to price the product

```

```

5     Then use the Linear Regression model to return the weighted price
6     """
7     self.log("Running Ensemble Agent - collaborating with specialist, frontier
8         and random forest agents")
9     specialist = self.specialist.price(description)
10    frontier = self.frontier.price(description)
11    random_forest = self.random_forest.price(description)
12    X = pd.DataFrame({
13        'Specialist': [specialist],
14        'Frontier': [frontier],
15        'RandomForest': [random_forest],
16        'Min': [min(specialist, frontier, random_forest)],
17        'Max': [max(specialist, frontier, random_forest)],
18    })
19    y = max(0, self.model.predict(X)[0])
20    self.log(f"Ensemble Agent complete - returning ${y:.2f}")
21    return y

```

Listing 5.15: Ensemble Price Prediction Analysis

Function Signature Analysis

Parameters.

- **self:** EnsembleAgent instance reference
- **description: str:** Product description for price estimation
- **Type Safety:** Clear input type specification

Return Value.

- **Return Type:** float (price estimate)
- **Constraint:** Non-negative via `max(0, ...)` constraint
- **Business Logic:** Ensures realistic price predictions

Step-by-Step Execution Analysis

Step 1: Process Initialization.

```

1 self.log("Running Ensemble Agent - collaborating with specialist, frontier and
2     random forest agents")
3 # Python execution:
4 # 1. Method resolution: self.log (inherited from Agent)
5 # 2. String literal passed as message
6 # 3. Agent.log() formats and outputs message
7 # Side effects: Terminal output, log file entry

```

Listing 5.16: Logging and Setup

Step 2: Individual Model Predictions.

```

1 specialist = self.specialist.price(description)
2 frontier = self.frontier.price(description)
3 random_forest = self.random_forest.price(description)
4 # Python execution sequence:
5 # 1. self.specialist.price() - Remote modal call
6 # 2. self.frontier.price() - RAG + LLM processing
7 # 3. self.random_forest.price() - Traditional ML prediction
8 # Memory: Float values assigned to local variables
9 # Time: Sequential execution, not parallel

```

Listing 5.17: Parallel Model Invocation

Individual Model Analysis.

- **SpecialistAgent:** Remote ML model via Modal platform
- **FrontierAgent:** Vector search + LLM reasoning
- **RandomForestAgent:** Scikit-learn model prediction
- **Execution Order:** Sequential, could be parallelized

Step 3: Feature Engineering.

```

1 X = pd.DataFrame({
2     'Specialist': [specialist],
3     'Frontier': [frontier],
4     'RandomForest': [random_forest],
5     'Min': [min(specialist, frontier, random_forest)],
6     'Max': [max(specialist, frontier, random_forest)],
7 })
8 # Python execution:
9 # 1. Dictionary construction with computed values
10 # 2. Built-in min() and max() function calls
11 # 3. pandas.DataFrame constructor call
12 # 4. DataFrame object creation in memory
13 # Memory: Dictionary + DataFrame objects

```

Listing 5.18: DataFrame Construction Analysis

Feature Engineering Analysis.

1. **Base Features:** Individual model predictions
2. **Statistical Features:** Min and max values across models
3. **Feature Engineering:** Creating additional predictive features
4. **Data Structure:** Single-row DataFrame for sklearn compatibility

Step 4: Meta-Model Prediction.

```

1 y = max(0, self.model.predict(X)[0])
2 # Python execution:
3 # 1. self.model reference resolution (loaded joblib model)

```

```

4 # 2. predict() method call with DataFrame
5 # 3. Array indexing [0] to get scalar value
6 # 4. max() function ensures non-negative result
7 # 5. Assignment to y variable
8 # Memory: NumPy array from prediction, scalar extraction

```

Listing 5.19: Ensemble Model Execution

Step 5: Result Logging and Return.

```

1 self.log(f"Ensemble Agent complete - returning ${y:.2f}")
2 return y
3 # Python execution:
4 # 1. F-string formatting with currency display
5 # 2. Logging call (inherited method)
6 # 3. Return float value to caller
7 # Memory: Formatted string creation, then cleanup

```

Listing 5.20: Completion Processing

Machine Learning Pipeline Analysis**Ensemble Learning Concepts.**

- **Base Learners:** Three different ML approaches
- **Meta-Learner:** Linear regression combining predictions
- **Feature Augmentation:** Statistical features enhance prediction
- **Stacking:** Advanced ensemble technique implementation

Performance Characteristics.

- **Latency:** Sum of individual model latencies + meta-model
- **Accuracy:** Typically better than individual models
- **Robustness:** Less sensitive to individual model errors
- **Complexity:** Higher computational and maintenance overhead

5.5 External API Integration Methods

5.5.1 FrontierAgent.find_similars() Method

```

1 def find_similars(self, description: str):
2     """
3         Return a list of items similar to the given one by looking in the Chroma
4         datastore
5         """
6         self.log("Frontier Agent is performing a RAG search of the Chroma datastore
7             to find 5 similar products")

```

```

6     vector = self.model.encode([description])
7     results =
8     self.collection.query(query_embeddings=vector.astype(float).tolist(),
9     n_results=5)
10    documents = results['documents'][0][:]
11    prices = [m['price'] for m in results['metadata'][0][:]]
12    self.log("Frontier Agent has found similar products")
13    return documents, prices

```

Listing 5.21: Vector Similarity Search Analysis

Function Signature Analysis

Parameters.

- **self:** FrontierAgent instance reference
- **description: str:** Product description to find similarities for
- **No Type Annotation:** Return type could be improved

Return Value.

- **Return Type:** Tuple[List[str], List[float]] (implicit)
- **documents:** List of similar product descriptions
- **prices:** Corresponding prices for similar products

Step-by-Step Execution Analysis

Step 1: Process Logging.

```

1 self.log("Frontier Agent is performing a RAG search of the Chroma datastore to
2     find 5 similar products")
3 # Python execution:
4 # 1. Inherited log method call
5 # 2. Descriptive message about RAG operation
6 # 3. Terminal/log output generation

```

Listing 5.22: RAG Process Announcement

Step 2: Vector Encoding.

```

1 vector = self.model.encode([description])
2 # Python execution:
3 # 1. self.model resolution (SentenceTransformer instance)
4 # 2. List creation with single description string
5 # 3. encode() method call - neural network forward pass
6 # 4. NumPy array returned with semantic embedding
7 # Memory: Input list, output NumPy array (384 dimensions)

```

Listing 5.23: Text to Vector Transformation

Vector Encoding Deep Dive.

1. **Model Type:** all-MiniLM-L6-v2 sentence transformer
2. **Input Processing:** Tokenization and attention mechanisms
3. **Neural Computation:** Transformer forward pass
4. **Output:** 384-dimensional dense vector representation
5. **Semantic Meaning:** Vector captures text semantic content

Step 3: Vector Database Query.

```

1 results = self.collection.query(
2     query_embeddings=vector.astype(float).tolist(),
3     n_results=5
4 )
5 # Python execution:
6 # 1. vector.astype(float) - ensure float32/float64 data type
7 # 2. .tolist() - convert NumPy array to Python list
8 # 3. self.collection.query() - ChromaDB search operation
9 # 4. Cosine similarity computation across stored vectors
10 # 5. Top-5 most similar items returned
11 # Memory: Type conversion, API call, result object

```

Listing 5.24: ChromaDB Similarity Search

ChromaDB Query Analysis.

- **Vector Comparison:** Cosine similarity computation
- **Index Search:** Efficient nearest neighbor search
- **Result Ranking:** Top-K results by similarity score
- **Metadata Retrieval:** Associated price information included

Step 4: Result Processing.

```

1 documents = results['documents'][0][:]
2 prices = [m['price'] for m in results['metadatas'][0][:]]
3 # Python execution:
4 # 1. Dictionary access: results['documents']
5 # 2. List indexing: [0] gets first query result batch
6 # 3. List slicing: [:] creates copy of document list
7 # 4. List comprehension: Extract price from metadata objects
8 # 5. Assignments to local variables
9 # Memory: New list objects created from query results

```

Listing 5.25: Query Result Extraction

Data Structure Analysis.

- **results Structure:** Nested dictionary with lists
- **documents:** Text descriptions of similar items

- **metadata**: Associated information including prices
 - **List Processing**: Comprehension for clean data extraction

Step 5: Completion and Return.

```
1 self.log("Frontier Agent has found similar products")
2 return documents, prices
3 # Python execution:
4 # 1. Success logging message
5 # 2. Tuple creation from two list objects
6 # 3. Return tuple reference to caller
7 # Memory: Tuple object creation, local variables cleanup
```

Listing 5.26: Result Return Processing

Memory and Performance Analysis

Computational Complexity.

- **Vector Encoding:** $O(n)$ where n is text length
 - **Similarity Search:** $O(\log k)$ with proper indexing
 - **Result Processing:** $O(m)$ where m is number of results
 - **Overall:** Dominated by neural network forward pass

Memory Usage.

- **Input Vector:** 384 floats \approx 1.5KB
 - **Query Results:** Variable size based on document length
 - **Processed Lists:** Additional copies for clean interface
 - **Cleanup:** Intermediate objects eligible for GC

5.6 Error Handling and Robustness Patterns

5.6.1 MessagingAgent.push() Method

```
1 def push(self, text):
2     """
3     Send a Push Notification using the Pushover API
4     """
5     self.log("Messaging Agent is sending a push notification")
6     conn = http.client.HTTPSConnection("api.pushover.net:443")
7     conn.request("POST", "/1/messages.json",
8         urllib.parse.urlencode({
9             "token": self.pushover_token,
10            "user": self.pushover_user,
11            "message": text,
12            "sound": "cashregister"
```

```

13     }, { "Content-type": "application/x-www-form-urlencoded" })
14 conn.getresponse()

```

Listing 5.27: Push Notification Error Handling

Function Signature Analysis

Parameters.

- **self:** MessagingAgent instance reference
- **text:** Message content to send (no type annotation)
- **Missing Annotations:** Could benefit from type hints

Return Value.

- **Return Type:** None (implicit)
- **Side Effects:** HTTP request sent, potential network errors
- **Error Handling:** Limited error handling present

Step-by-Step Execution Analysis

Step 1: Process Logging.

```

1 self.log("Messaging Agent is sending a push notification")
2 # Python execution:
3 # 1. Inherited log method call
4 # 2. Process status announcement
5 # 3. Terminal/log output for debugging

```

Listing 5.28: Notification Process Start

Step 2: HTTP Connection Setup.

```

1 conn = http.client.HTTPSConnection("api.pushover.net:443")
2 # Python execution:
3 # 1. http.client module resolution
4 # 2. HTTPSConnection class instantiation
5 # 3. SSL/TLS connection preparation (not yet established)
6 # 4. Connection object assigned to local variable
7 # Memory: Connection object with SSL context

```

Listing 5.29: HTTPS Connection Establishment

Connection Analysis.

- **Protocol:** HTTPS for secure communication
- **Host:** Pushover API endpoint
- **Port:** Explicit port 443 specification
- **SSL Context:** Automatic certificate verification

Step 3: Request Data Preparation.

```

1  urllib.parse.urlencode({
2      "token": self.pushover_token,
3      "user": self.pushover_user,
4      "message": text,
5      "sound": "cashregister"
6  })
7  # Python execution:
8  # 1. Dictionary creation with API parameters
9  # 2. Instance attribute access for token/user
10 # 3. urllib.parse.urlencode() call
11 # 4. URL-encoded string creation
12 # Memory: Dictionary object, encoded string

```

Listing 5.30: HTTP Request Construction

Data Encoding Analysis.

- **URL Encoding:** Special characters properly escaped
- **Form Data:** application/x-www-form-urlencoded format
- **API Parameters:** token, user, message, sound
- **Custom Sound:** "cashregister" for deal notifications

Step 4: HTTP Request Execution.

```

1  conn.request("POST", "/1/messages.json",
2     encoded_data,
3     { "Content-type": "application/x-www-form-urlencoded" })
4  # Python execution:
5  # 1. HTTP POST method specification
6  # 2. API endpoint path: /1/messages.json
7  # 3. Request body: URL-encoded form data
8  # 4. Headers: Content-Type specification
9  # 5. Network transmission of HTTP request
10 # Network: SSL handshake, HTTP request sent

```

Listing 5.31: POST Request Transmission

Step 5: Response Handling.

```

1  conn.getresponse()
2  # Python execution:
3  # 1. HTTP response reception from server
4  # 2. Response object creation with status/headers/body
5  # 3. Response object returned but not used
6  # 4. Connection cleanup (implicit)
7  # Issue: Response not checked for errors!

```

Listing 5.32: Response Processing

Error Handling Analysis

Missing Error Handling.

- **Network Errors:** No try/except for connection failures
- **HTTP Errors:** Response status code not checked
- **SSL Errors:** Certificate validation failures not handled
- **Timeout Errors:** No timeout configuration

Improved Error Handling.

```

1  def push_with_error_handling(self, text):
2      try:
3          self.log("Messaging Agent is sending a push notification")
4          conn = http.client.HTTPSConnection("api.pushover.net:443", timeout=10)
5
6          data = urllib.parse.urlencode({
7              "token": self.pushover_token,
8              "user": self.pushover_user,
9              "message": text,
10             "sound": "cashregister"
11         })
12
13         conn.request("POST", "/1/messages.json", data,
14                     {"Content-type": "application/x-www-form-urlencoded"})
15
16         response = conn.getresponse()
17         if response.status != 200:
18             self.log(f"Pushover API error: {response.status} {response.reason}")
19             return False
20
21         self.log("Push notification sent successfully")
22         return True
23
24     except (http.client.HTTPException, OSError) as e:
25         self.log(f"Network error sending push notification: {e}")
26         return False
27     finally:
28         conn.close() # Ensure connection cleanup

```

Listing 5.33: Enhanced Error Handling Pattern

5.7 Complex Business Logic Methods

5.7.1 PlanningAgent.plan() Method

```

1  def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
2      """
3          Run the full workflow:
4              1. Use the ScannerAgent to find deals from RSS feeds

```

```

5     2. Use the EnsembleAgent to estimate them
6     3. Use the MessagingAgent to send a notification of deals
7     """
8     self.log("Planning Agent is kicking off a run")
9     selection = self.scanner.scan(memory=memory)
10    if selection:
11        opportunities = [self.run(deal) for deal in selection.deals[:5]]
12        opportunities.sort(key=lambda opp: opp.discount, reverse=True)
13        best = opportunities[0]
14        self.log(f"Planning Agent has identified the best deal has discount
15        ${best.discount:.2f}")
16        if best.discount > self.DEAL_THRESHOLD:
17            self.messenger.alert(best)
18        self.log("Planning Agent has completed a run")
19        return best if best.discount > self.DEAL_THRESHOLD else None
20    return None

```

Listing 5.34: Master Planning Method Analysis

Function Signature Analysis

Parameters.

- **self:** PlanningAgent instance reference
- **memory: List[str]:** Previously processed deal URLs
- **Default Value:** Empty list (mutable default - potential issue)
- **Return Type:** Optional[Opportunity] - may return None

Mutable Default Argument Issue.

```

1 def plan(self, memory: List[str] = []): # Problematic!
2     # Same list object reused across calls
3     # Better: memory: Optional[List[str]] = None
4     # Then: if memory is None: memory = []

```

Listing 5.35: Mutable Default Problem

Step-by-Step Execution Analysis

Step 1: Workflow Initiation.

```

1 self.log("Planning Agent is kicking off a run")
2 # Python execution:
3 # 1. Log workflow start for debugging/monitoring
4 # 2. Indicates beginning of complete deal processing cycle

```

Listing 5.36: Planning Process Start

Step 2: Deal Discovery.

```

1 selection = self.scanner.scan(memory=memory)
2 # Python execution:
3 # 1. self.scanner attribute access (ScannerAgent instance)

```

```

4 # 2. scan() method call with memory parameter
5 # 3. RSS feed processing, ML filtering, deal extraction
6 # 4. DealSelection object or None returned
7 # Side effects: Network requests, LLM API calls

```

Listing 5.37: Scanner Agent Invocation

Step 3: Conditional Processing.

```

1 if selection:
2     # Branch 1: Deals found - process them
3     opportunities = [self.run(deal) for deal in selection.deals[:5]]
4     opportunities.sort(key=lambda opp: opp.discount, reverse=True)
5     best = opportunities[0]
6     # ... processing continues
7 else:
8     # Branch 2: No deals found
9     return None

```

Listing 5.38: Deal Processing Branch

Branch 1 Analysis: Deal Processing.

1. **List Comprehension:** [self.run(deal) for deal in selection.deals[:5]]
2. **Slice Operation:** Process maximum 5 deals for performance
3. **Opportunity Creation:** Each deal converted to Opportunity object
4. **List Sorting:** Sort by discount amount, descending order
5. **Best Selection:** First item after sorting is highest discount

List Comprehension Deep Dive.

```

1 opportunities = [self.run(deal) for deal in selection.deals[:5]]
2 # Python execution:
3 # 1. selection.deals[:5] - slice first 5 deals
4 # 2. Iterator creation over sliced list
5 # 3. For each deal:
6 #     a. self.run(deal) method call
7 #     b. EnsembleAgent price estimation
8 #     c. Opportunity object creation
9 #     d. Add to result list
10 # 4. Complete list assigned to opportunities
11 # Performance: Sequential execution of 5 price estimations

```

Listing 5.39: List Comprehension Execution

Step 4: Sorting and Selection.

```

1 opportunities.sort(key=lambda opp: opp.discount, reverse=True)
2 best = opportunities[0]
3 # Python execution:
4 # 1. Lambda function creation for key extraction
5 # 2. List.sort() in-place sorting (Timsort algorithm)

```

```

6 # 3. reverse=True for descending order
7 # 4. Index access [0] for highest discount opportunity
8 # Memory: List sorted in-place, lambda object for key

```

Listing 5.40: Opportunity Ranking

Step 5: Decision Making.

```

1 if best.discount > self.DEAL_THRESHOLD:
2     self.messenger.alert(best)
3     return best
4 else:
5     return None
6 # Python execution:
7 # 1. Attribute access: best.discount (float value)
8 # 2. Class constant access: self.DEAL_THRESHOLD (50)
9 # 3. Numeric comparison operation
10 # 4. Conditional execution of alert
11 # 5. Return appropriate value based on threshold

```

Listing 5.41: Threshold Decision Logic

Business Logic Analysis**Workflow Orchestration.**

1. **Deal Discovery:** Scanner finds new deals from RSS feeds
2. **Price Estimation:** Ensemble model predicts fair prices
3. **Opportunity Evaluation:** Calculate potential savings
4. **Ranking:** Sort by discount amount
5. **Filtering:** Apply business threshold for notifications
6. **Communication:** Alert user of significant deals

Performance Characteristics.

- **Network Bound:** RSS feeds, HTTP requests, API calls
- **Compute Bound:** ML model inference, text processing
- **Sequential:** No parallelization of deal processing
- **Threshold Dependent:** Output depends on deal quality

5.8 Chapter Summary

This chapter provided comprehensive analysis of key functions and methods throughout the multi-agent system:

- **Base Infrastructure:** Agent.log() method with color formatting

- **Data Processing:** HTML extraction and cleaning functions
- **Object Creation:** Class methods and alternative constructors
- **ML Integration:** Ensemble prediction and vector search methods
- **API Integration:** HTTP requests and external service calls
- **Business Logic:** Complex workflow orchestration

Each function demonstrates different aspects of Python programming:

- Memory management and object lifecycle
- Error handling patterns and robustness
- Network programming and API integration
- Machine learning pipeline construction
- String processing and data transformation
- Conditional logic and decision making

The analysis reveals both strengths (modular design, clear interfaces) and areas for improvement (error handling, performance optimization, type annotations) in the codebase. Understanding these implementation details provides insight into how high-level agent behaviors emerge from carefully orchestrated low-level operations.

Chapter 6

Python Concepts and Examples

This chapter identifies and explains all Python language features, concepts, and programming patterns used throughout the multi-agent system. For each concept, we provide clear explanations, simple illustrative examples, and then demonstrate how they're applied in the actual system code.

6.1 Object-Oriented Programming Concepts

6.1.1 Classes and Objects

Concept Explanation. Classes define blueprints for creating objects, encapsulating data (attributes) and behavior (methods). Objects are instances of classes that maintain their own state.

Simple Example.

```
1 class Dog:
2     species = "Canis lupus" # Class attribute
3
4     def __init__(self, name, age):
5         self.name = name      # Instance attribute
6         self.age = age        # Instance attribute
7
8     def bark(self):          # Instance method
9         return f"{self.name} says woof!"
10
11 # Object creation and usage
12 my_dog = Dog("Rex", 3)
13 print(my_dog.bark()) # Rex says woof!
14 print(my_dog.species) # Canis lupus
```

Listing 6.1: Basic Class Concept

System Implementation.

```
1 class Agent:
2     # Class attributes (shared across all instances)
3     RED = '\033[31m'
4     GREEN = '\033[32m'
```

```

5     RESET = '\033[0m'
6
7     # Instance attributes with type hints
8     name: str = ""
9     color: str = '\033[37m'
10
11    def log(self, message): # Instance method
12        """Log messages with agent identification"""
13        color_code = self.BG_BLACK + self.color
14        message = f"[{self.name}] {message}"
15        logging.info(color_code + message + self.RESET)
16
17    # Concrete implementation
18    class SpecialistAgent(Agent):
19        name = "Specialist Agent" # Override class attribute
20        color = Agent.RED         # Reference parent class attribute

```

Listing 6.2: Agent Class Implementation

Key Points in System.

- **Class Attributes:** Color constants shared across all agents
- **Instance Attributes:** Each agent has unique name and color
- **Method Implementation:** log() provides common functionality
- **Attribute Override:** Subclasses customize name and color

6.1.2 Inheritance

Concept Explanation. Inheritance allows classes to inherit attributes and methods from parent classes, enabling code reuse and establishing "is-a" relationships.

Simple Example.

```

1  class Animal:
2      def __init__(self, name):
3          self.name = name
4
5      def speak(self):
6          pass # Abstract method to be overridden
7
8  class Dog(Animal): # Dog inherits from Animal
9      def speak(self):
10         return f"{self.name} barks"
11
12 class Cat(Animal): # Cat inherits from Animal
13     def speak(self):
14         return f"{self.name} meows"
15
16 # Polymorphic usage
17 animals = [Dog("Rex"), Cat("Whiskers")]

```

```

18 for animal in animals:
19     print(animal.speak()) # Calls appropriate subclass method

```

Listing 6.3: Basic Inheritance Example

System Implementation.

```

1 class Agent: # Base class
2     def log(self, message):
3         """Template method used by all subclasses"""
4         color_code = self.BG_BLACK + self.color
5         message = f"[{self.name}] {message}"
6         logging.info(color_code + message + self.RESET)
7
8 class PlanningAgent(Agent): # Inherits from Agent
9     name = "Planning Agent"
10    color = Agent.GREEN
11
12    def __init__(self, collection):
13        # Uses inherited log() method
14        self.log("Planning Agent is initializing")
15        self.scanner = ScannerAgent()
16        self.ensemble = EnsembleAgent(collection)
17        self.messenger = MessagingAgent()
18        self.log("Planning Agent is ready")
19
20 class EnsembleAgent(Agent): # Also inherits from Agent
21     name = "Ensemble Agent"
22     color = Agent.YELLOW
23
24     def price(self, description: str) -> float:
25         # Uses inherited log() method
26         self.log("Running Ensemble Agent")
27         # Custom logic here
28         return predicted_price

```

Listing 6.4: Agent Inheritance Hierarchy

Inheritance Benefits in System.

- **Code Reuse:** All agents get logging functionality for free
- **Consistency:** Common interface across all agent types
- **Maintainability:** Changes to Agent affect all subclasses
- **Polymorphism:** Can treat all agents uniformly

6.1.3 Composition

Concept Explanation. Composition is a "has-a" relationship where objects contain other objects as components, enabling complex behavior through combining simpler parts.

Simple Example.

```

1  class Engine:
2      def __init__(self, horsepower):
3          self.horsepower = horsepower
4
5      def start(self):
6          return "Engine started"
7
8  class Car:
9      def __init__(self, make, engine):
10         self.make = make
11         self.engine = engine # Car HAS-A Engine (composition)
12
13     def start(self):
14         return f"{self.make}: {self.engine.start()}"
15
16 # Usage
17 v8_engine = Engine(400)
18 muscle_car = Car("Mustang", v8_engine)
19 print(muscle_car.start()) # Mustang: Engine started

```

Listing 6.5: Basic Composition Example

System Implementation.

```

1  class EnsembleAgent(Agent):
2      def __init__(self, collection):
3          self.log("Initializing Ensemble Agent")
4
5          # Composition: EnsembleAgent HAS-A these agents
6          self.specialist = SpecialistAgent() # HAS-A SpecialistAgent
7          self.frontier = FrontierAgent(collection) # HAS-A FrontierAgent
8          self.random_forest = RandomForestAgent() # HAS-A RandomForestAgent
9          self.model = joblib.load('ensemble_model.pkl') # HAS-A ML model
10
11         self.log("Ensemble Agent is ready")
12
13     def price(self, description: str) -> float:
14         """Coordinates composed agents to make ensemble prediction"""
15         # Delegate to composed agents
16         specialist = self.specialist.price(description)
17         frontier = self.frontier.price(description)
18         random_forest = self.random_forest.price(description)
19
20         # Combine their results
21         X = pd.DataFrame({
22             'Specialist': [specialist],
23             'Frontier': [frontier],
24             'RandomForest': [random_forest],
25             'Min': [min(specialist, frontier, random_forest)],
26             'Max': [max(specialist, frontier, random_forest)],
27         })

```

```
28     return max(0, self.model.predict(X)[0])
```

Listing 6.6: EnsembleAgent Composition Pattern

Composition Benefits in System.

- **Flexibility:** Can change components at runtime
- **Modularity:** Each component has single responsibility
- **Loose Coupling:** Components don't depend on each other
- **Reusability:** Same components used in different contexts

6.2 Modern Python Features

6.2.1 Type Annotations

Concept Explanation. Type annotations provide static type information for variables, function parameters, and return values, improving code documentation and enabling better IDE support.

Simple Example.

```
1 from typing import List, Dict, Optional
2
3 def greet(name: str, age: int) -> str:
4     """Function with typed parameters and return value"""
5     return f"Hello {name}, you are {age} years old"
6
7 def process_scores(scores: List[float]) -> Dict[str, float]:
8     """Function with complex type annotations"""
9     return {
10         'average': sum(scores) / len(scores),
11         'max': max(scores),
12         'min': min(scores)
13     }
14
15 # Variable annotations
16 name: str = "Alice"
17 numbers: List[int] = [1, 2, 3, 4, 5]
18 user_data: Optional[Dict[str, str]] = None
```

Listing 6.7: Type Annotation Examples

System Implementation.

```
1 from typing import List, Dict, Self, Optional
2
3 class ScrapedDeal:
4     # Instance attributes with type annotations
5     category: str
6     title: str
7     summary: str
```

```

8     url: str
9     details: str
10    features: str
11
12    @classmethod
13    def fetch(cls, show_progress: bool = False) -> List[Self]:
14        """Class method with modern Self type annotation"""
15        deals = []
16        for feed_url in feeds:
17            feed = feedparser.parse(feed_url)
18            for entry in feed.entries[:10]:
19                deals.append(cls(entry))
20        return deals
21
22 class FrontierAgent(Agent):
23     def price(self, description: str) -> float:
24         """Clear input and output types"""
25         documents, prices = self.find_similars(description)
26         # ... processing
27         return result
28
29     def find_similars(self, description: str) -> tuple[List[str], List[float]]:
30         """Modern Python 3.9+ tuple syntax"""
31         # ... implementation
32         return documents, prices
33
34 class PlanningAgent(Agent):
35     def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
36         """Optional return type indicates possible None"""
37         selection = self.scanner.scan(memory=memory)
38         if selection:
39             # ... processing
40             return best_opportunity
41         return None

```

Listing 6.8: System Type Annotations

Type Annotation Benefits in System.

- **IDE Support:** Better autocomplete and error detection
- **Documentation:** Types serve as inline documentation
- **Maintainability:** Easier to understand interfaces
- **Debugging:** Type checkers can catch errors early

6.2.2 F-Strings (Formatted String Literals)

Concept Explanation. F-strings provide a readable and efficient way to format strings by embedding expressions directly inside string literals.

Simple Example.

```

1 name = "Alice"
2 age = 30
3 score = 87.5
4
5 # Basic f-string interpolation
6 greeting = f"Hello, {name}! You are {age} years old."
7
8 # Expressions inside f-strings
9 result = f"{name}'s score is {score:.1f}%, which is {'PASS' if score >= 60 else
10   'FAIL'}"
11
12 # Format specifications
13 price = 123.456
14 formatted = f"Price: ${price:.2f}"  # Price: $123.46
15
16 # Debugging with f-strings
17 x, y = 10, 20
18 debug_info = f"{x=}, {y=}, {x+y=}"  # x=10, y=20, x+y=30

```

Listing 6.9: F-String Examples

System Implementation.

```

1 class Agent:
2     def log(self, message):
3         """F-string for agent identification"""
4         color_code = self.BG_BLACK + self.color
5         message = f"[{self.name}] {message}"  # F-string interpolation
6         logging.info(color_code + message + self.RESET)
7
8 class EnsembleAgent(Agent):
9     def price(self, description: str) -> float:
10        self.log("Running Ensemble Agent - collaborating with specialist,
11        frontier and random forest agents")
12        # ... processing
13        y = max(0, self.model.predict(X)[0])
14        self.log(f"Ensemble Agent complete - returning ${y:.2f}")  # Currency
15        formatting
16        return y
17
18 class PlanningAgent(Agent):
19     def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
20         # ... processing
21         best = opportunities[0]
22         # F-string with currency formatting
23         self.log(f"Planning Agent has identified the best deal has discount
24         ${best.discount:.2f}")
25         return best if best.discount > self.DEAL_THRESHOLD else None
26
27 class MessagingAgent(Agent):
28     def alert(self, opportunity: Opportunity):

```

```

26     """Complex f-string with multiple expressions"""
27     text = f"Deal Alert! Price=${opportunity.deal.price:.2f}, "
28     text += f"Estimate=${opportunity.estimate:.2f}, "
29     text += f"Discount=${opportunity.discount:.2f} :"
30     text += opportunity.deal.product_description[:10]+'... '
31     text += opportunity.deal.url
32     # Send notification

```

Listing 6.10: F-Strings in Agent System

F-String Advantages in System.

- **Readability:** Clear variable interpolation
- **Performance:** Faster than .format() or % formatting
- **Flexibility:** Support for expressions and formatting specs
- **Currency Display:** Consistent price formatting (.2f)

6.2.3 List Comprehensions

Concept Explanation. List comprehensions provide a concise way to create lists by applying expressions to items in iterables, optionally filtering them with conditions.

Simple Example.

```

1 # Basic list comprehension
2 numbers = [1, 2, 3, 4, 5]
3 squares = [x**2 for x in numbers] # [1, 4, 9, 16, 25]
4
5 # With conditional filtering
6 even_squares = [x**2 for x in numbers if x % 2 == 0] # [4, 16]
7
8 # With expression transformation
9 words = ["hello", "world", "python"]
10 lengths = [len(word) for word in words] # [5, 5, 6]
11
12 # Nested comprehension
13 matrix = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
14 flattened = [item for row in matrix for item in row] # [1,2,3,4,5,6,7,8,9]
15
16 # Dictionary comprehension
17 word_lengths = {word: len(word) for word in words} # {'hello': 5, 'world': 5,
    'python': 6}

```

Listing 6.11: List Comprehension Examples

System Implementation.

```

1 class PlanningAgent(Agent):
2     def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
3         selection = self.scanner.scan(memory=memory)
4         if selection:

```

```

5      # List comprehension to convert deals to opportunities
6      opportunities = [self.run(deal) for deal in selection.deals[:5]]
7      opportunities.sort(key=lambda opp: opp.discount, reverse=True)
8      best = opportunities[0]
9      return best if best.discount > self.DEAL_THRESHOLD else None
10     return None
11
12 class FrontierAgent(Agent):
13     def find_similars(self, description: str):
14         vector = self.model.encode([description])
15         results =
16         self.collection.query(query_embeddings=vector.astype(float).tolist(),
17                               n_results=5)
18         documents = results['documents'][0][:]
19         # List comprehension to extract prices from metadata
20         prices = [m['price'] for m in results['metadatas'][0][:]]
21         return documents, prices
22
23 class ScannerAgent(Agent):
24     def fetch_deals(self, memory) -> List[ScrapedDeal]:
25         urls = [opp.deal.url for opp in memory] # Extract URLs
26         scraped = ScrapedDeal.fetch()
27         # Filter deals not already in memory
28         result = [scrape for scrape in scraped if scrape.url not in urls]
29         return result

```

Listing 6.12: List Comprehensions in Agent System

List Comprehension Benefits in System.

- **Conciseness:** Replace multi-line loops with single expressions
- **Readability:** Clear intent for data transformation
- **Performance:** Often faster than equivalent for loops
- **Functional Style:** Encourages immutable data transformations

6.3 Advanced Python Concepts

6.3.1 Decorators

Concept Explanation. Decorators are functions that modify or extend the behavior of other functions or classes without permanently modifying them.

Simple Example.

```

1 def timing_decorator(func):
2     """Decorator to measure function execution time"""
3     import time
4     def wrapper(*args, **kwargs):
5         start = time.time()

```

```

6     result = func(*args, **kwargs)
7     end = time.time()
8     print(f"{func.__name__} took {end - start:.2f} seconds")
9     return result
10    return wrapper
11
12 @timing_decorator
13 def slow_function():
14     import time
15     time.sleep(1)
16     return "Done"
17
18 # Usage
19 result = slow_function() # Prints: slow_function took 1.00 seconds
20
21 # Class method decorator
22 class MyClass:
23     @classmethod
24     def create_from_string(cls, data: str):
25         """Alternative constructor"""
26         return cls(data.split(','))
27
28     @staticmethod
29     def utility_function(x, y):
30         """Function that doesn't need class or instance"""
31         return x + y

```

Listing 6.13: Decorator Examples

System Implementation.

```

1 class ScrapedDeal:
2     @classmethod # Class method decorator
3     def fetch(cls, show_progress: bool = False) -> List[Self]:
4         """Alternative constructor for creating multiple deals"""
5         deals = []
6         feed_iter = tqdm(feeds) if show_progress else feeds
7         for feed_url in feed_iter:
8             feed = feedparser.parse(feed_url)
9             for entry in feed.entries[:10]:
10                 deals.append(cls(entry)) # cls refers to ScrapedDeal
11                 time.sleep(0.5)
12
13
14 # Potential decorator for logging (not implemented in original)
15 def log_execution(func):
16     """Decorator to log method execution"""
17     def wrapper(self, *args, **kwargs):
18         self.log(f"Starting {func.__name__}")
19         result = func(self, *args, **kwargs)
20         self.log(f"Completed {func.__name__}")
21         return result

```

```
22     return wrapper
23
24 class EnhancedAgent(Agent):
25     @log_execution
26     def price(self, description: str) -> float:
27         # Method execution automatically logged
28         return self.calculate_price(description)
```

Listing 6.14: Decorators in Agent System

Decorator Benefits.

- **Separation of Concerns:** Keep core logic separate from cross-cutting concerns
- **Reusability:** Same decorator can be applied to multiple functions
- **Clean Code:** Avoid repetitive boilerplate code
- **Flexibility:** Add/remove behavior without modifying original functions

6.3.2 Context Managers

Concept Explanation. Context managers provide a way to allocate and release resources precisely when needed, typically using the `with` statement.

Simple Example.

```
1 # Built-in context manager for file handling
2 with open('data.txt', 'r') as file:
3     content = file.read()
4     # File automatically closed when leaving with block
5
6 # Custom context manager
7 class Timer:
8     def __enter__(self):
9         import time
10        self.start = time.time()
11        return self
12
13     def __exit__(self, exc_type, exc_val, exc_tb):
14         import time
15         self.end = time.time()
16         print(f"Elapsed: {self.end - self.start:.2f} seconds")
17
18 # Usage
19 with Timer():
20     # Code to time
21     sum(range(1000000))
22
23 # Using contextlib
24 from contextlib import contextmanager
25
26 @contextmanager
```

```

27 def database_transaction():
28     print("Starting transaction")
29     try:
30         yield "connection"
31         print("Committing transaction")
32     except Exception:
33         print("Rolling back transaction")
34         raise
35     finally:
36         print("Cleaning up")
37
38 with database_transaction() as conn:
39     # Database operations here
40     pass

```

Listing 6.15: Context Manager Examples

System Implementation.

```

1 # The system could benefit from context managers for:
2
3 class AgentContext:
4     """Context manager for agent lifecycle"""
5     def __init__(self, agent_class, *args, **kwargs):
6         self.agent_class = agent_class
7         self.args = args
8         self.kwargs = kwargs
9         self.agent = None
10
11    def __enter__(self):
12        self.agent = self.agent_class(*self.args, **self.kwargs)
13        self.agent.log("Agent context started")
14        return self.agent
15
16    def __exit__(self, exc_type, exc_val, exc_tb):
17        if self.agent:
18            self.agent.log("Agent context ending")
19        # Cleanup resources if needed
20
21 # Usage
22 with AgentContext(FrontierAgent, collection) as agent:
23     price = agent.price("laptop computer")
24     # Agent automatically cleaned up
25
26 # HTTP connection context manager
27 class HTTPConnection:
28     def __init__(self, host, port=443):
29         self.host = host
30         self.port = port
31         self.conn = None
32
33     def __enter__(self):

```

```

34     import http.client
35     self.conn = http.client.HTTPSConnection(f"[{self.host}]:{self.port}")
36     return self.conn
37
38     def __exit__(self, exc_type, exc_val, exc_tb):
39         if self.conn:
40             self.conn.close()
41
42 # Enhanced MessagingAgent with context manager
43 class EnhancedMessagingAgent(Agent):
44     def push(self, text):
45         with HTTPConnection("api.pushover.net") as conn:
46             data = urllib.parse.urlencode({
47                 "token": self.pushover_token,
48                 "user": self.pushover_user,
49                 "message": text,
50                 "sound": "cashregister"
51             })
52             conn.request("POST", "/1/messages.json", data,
53                         {"Content-type": "application/x-www-form-urlencoded"})
54             response = conn.getresponse()
55             return response.status == 200

```

Listing 6.16: Context Managers in Agent System (Potential)

Context Manager Benefits.

- **Resource Management:** Automatic cleanup of resources
- **Exception Safety:** Cleanup occurs even if exceptions happen
- **Clear Intent:** Makes resource lifecycle explicit
- **Reduced Errors:** Prevents resource leaks

6.4 Data Structures and Collections

6.4.1 Dictionaries and Dictionary Operations

Concept Explanation. Dictionaries are key-value mappings that provide fast lookups and flexible data storage.

Simple Example.

```

1 # Basic dictionary operations
2 person = {
3     'name': 'Alice',
4     'age': 30,
5     'city': 'New York'
6 }
7
8 # Dictionary access and modification

```

```

9 print(person['name'])      # Alice
10 person['age'] = 31         # Update existing key
11 person['job'] = 'Engineer' # Add new key
12
13 # Safe access with get()
14 email = person.get('email', 'not_provided@example.com')
15
16 # Dictionary comprehension
17 squares = {x: x**2 for x in range(5)} # {0: 0, 1: 1, 2: 4, 3: 9, 4: 16}
18
19 # Merging dictionaries (Python 3.9+)
20 defaults = {'color': 'blue', 'size': 'medium'}
21 options = {'color': 'red'}
22 config = defaults | options # {'color': 'red', 'size': 'medium'}

```

Listing 6.17: Dictionary Examples

System Implementation.

```

1 class ScrapedDeal:
2     def __init__(self, entry: Dict[str, str]):
3         """RSS entry is a dictionary-like object"""
4         self.title = entry['title']                      # Dictionary access
5         self.summary = extract(entry['summary'])        # Dictionary access
6         self.url = entry['links'][0]['href']           # Nested dictionary access
7         # ... more processing
8
9 class EnsembleAgent(Agent):
10    def price(self, description: str) -> float:
11        # ... get individual predictions
12
13        # Create DataFrame from dictionary
14        X = pd.DataFrame({
15            'Specialist': [specialist],
16            'Frontier': [frontier],
17            'RandomForest': [random_forest],
18            'Min': [min(specialist, frontier, random_forest)],
19            'Max': [max(specialist, frontier, random_forest)],
20        })
21        return max(0, self.model.predict(X)[0])
22
23 class MessagingAgent(Agent):
24    def push(self, text):
25        # Dictionary for POST data
26        data = urllib.parse.urlencode({
27            "token": self.pushover_token,
28            "user": self.pushover_user,
29            "message": text,
30            "sound": "cashregister"
31        })
32        # ... send request
33

```

```

34 class FrontierAgent(Agent):
35     def messages_for(self, description: str, similars: List[str], prices:
36         List[float]) -> List[Dict[str, str]]:
37         """Create OpenAI message format (list of dictionaries)"""
38         return [
39             {"role": "system", "content": system_message},
40             {"role": "user", "content": user_prompt},
41             {"role": "assistant", "content": "Price is $"}
42         ]
43
43     def find_similars(self, description: str):
44         results = self.collection.query(...)
45         # Access nested dictionary structure
46         documents = results['documents'][0][:]
47         prices = [m['price'] for m in results['metadata'][0][:]]
48         return documents, prices

```

Listing 6.18: Dictionary Usage in Agent System

Dictionary Patterns in System.

- **Configuration Data:** RSS entries, API parameters
- **Structured Data:** DataFrame construction, message formats
- **Nested Access:** Multi-level dictionary navigation
- **Safe Access:** Using get() method with defaults

6.4.2 Lists and List Operations

Concept Explanation. Lists are ordered, mutable sequences that can contain any type of objects.

Simple Example.

```

1 # Basic list operations
2 numbers = [1, 2, 3, 4, 5]
3 numbers.append(6)           # Add to end: [1, 2, 3, 4, 5, 6]
4 numbers.insert(0, 0)        # Insert at position: [0, 1, 2, 3, 4, 5, 6]
5 numbers.remove(3)          # Remove first occurrence: [0, 1, 2, 4, 5, 6]
6
7 # List slicing
8 first_three = numbers[:3]  # [0, 1, 2]
9 last_two = numbers[-2:]    # [5, 6]
10 every_second = numbers[::2] # [0, 2, 4]
11
12 # List methods
13 numbers.sort()            # In-place sorting
14 numbers.reverse()          # In-place reversal
15 count = numbers.count(2)   # Count occurrences
16
17 # List unpacking

```

```
18 a, b, *rest = [1, 2, 3, 4, 5] # a=1, b=2, rest=[3, 4, 5]
```

Listing 6.19: List Examples

System Implementation.

```

1 # Module-level list configuration
2 feeds = [
3     "https://www.dealnews.com/c142/Electronics/?rss=1",
4     "https://www.dealnews.com/c39/Computers/?rss=1",
5     "https://www.dealnews.com/c238/Automotive/?rss=1",
6     "https://www.dealnews.com/f1912/Smart-Home/?rss=1",
7     "https://www.dealnews.com/c196/Home-Garden/?rss=1",
8 ]
9
10 class ScrapedDeal:
11     @classmethod
12     def fetch(cls, show_progress: bool = False) -> List[Self]:
13         deals = [] # Initialize empty list
14         for feed_url in feeds:
15             feed = feedparser.parse(feed_url)
16             # List slicing to limit entries
17             for entry in feed.entries[:10]:
18                 deals.append(cls(entry)) # List append
19         return deals
20
21 class PlanningAgent(Agent):
22     def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
23         selection = self.scanner.scan(memory=memory)
24         if selection:
25             # List comprehension + slicing
26             opportunities = [self.run(deal) for deal in selection.deals[:5]]
27             # In-place sorting with custom key
28             opportunities.sort(key=lambda opp: opp.discount, reverse=True)
29             best = opportunities[0] # List indexing
30             return best if best.discount > self.DEAL_THRESHOLD else None
31         return None
32
33 class ScannerAgent(Agent):
34     def fetch_deals(self, memory) -> List[ScrapedDeal]:
35         # List comprehension for URL extraction
36         urls = [opp.deal.url for opp in memory]
37         scraped = ScrapedDeal.fetch()
38         # List comprehension with filtering
39         result = [scrape for scrape in scraped if scrape.url not in urls]
40         return result
41
42 class FrontierAgent(Agent):
43     def find_similars(self, description: str):
44         # Single-item list for encoding
45         vector = self.model.encode([description])
```

```

46     results =
47     self.collection.query(query_embeddings=vector.astype(float).tolist(),
48     n_results=5)
49     # List access and slicing
50     documents = results['documents'][0][:]
51     # List comprehension for data extraction
52     prices = [m['price'] for m in results['metadata'][0][:]]
53     return documents, prices

```

Listing 6.20: List Usage in Agent System

List Patterns in System.

- **Configuration:** Module-level feed URLs
- **Collection Building:** Accumulating results with append()
- **Data Processing:** Slicing, sorting, filtering
- **Comprehensions:** Concise data transformations

6.5 Exception Handling

Concept Explanation. Exception handling allows programs to respond gracefully to errors and unexpected conditions.

Simple Example.

```

1 # Basic try-except
2 try:
3     result = 10 / 0
4 except ZeroDivisionError:
5     print("Cannot divide by zero")
6     result = None
7
8 # Multiple exception types
9 try:
10    data = {'key': 'value'}
11    value = data['missing_key']
12    number = int(value)
13 except KeyError:
14     print("Key not found")
15 except ValueError:
16     print("Invalid number format")
17 except Exception as e:
18     print(f"Unexpected error: {e}")
19
20 # Try-except-else-finally
21 try:
22     file = open('data.txt', 'r')
23 except FileNotFoundError:
24     print("File not found")

```

```

25     else:
26         # Runs only if no exception occurred
27         content = file.read()
28         print(f"Read {len(content)} characters")
29     finally:
30         # Always runs
31         if 'file' in locals():
32             file.close()
33
34     # Custom exceptions
35     class ValidationError(Exception):
36         def __init__(self, message, code=None):
37             super().__init__(message)
38             self.code = code
39
40     def validate_price(price):
41         if price < 0:
42             raise ValidationError("Price cannot be negative", code="NEGATIVE_PRICE")
43         return price

```

Listing 6.21: Exception Handling Examples

System Implementation.

```

1  # Current system has limited exception handling
2  # Here's what exists and what could be improved:
3
4  class FrontierAgent(Agent):
5      def price(self, description: str) -> float:
6          """Current implementation lacks error handling"""
7          documents, prices = self.find_similars(description)
8          response = self.client.chat.completions.create(
9              model=self.MODEL,
10             messages=self.messages_for(description, documents, prices),
11             seed=42,
12             max_tokens=5
13         )
14         reply = response.choices[0].message.content
15         result = self.get_price(reply)
16         return result
17
18 # Enhanced version with proper exception handling
19 class EnhancedFrontierAgent(Agent):
20     def price(self, description: str) -> float:
21         try:
22             documents, prices = self.find_similars(description)
23         except Exception as e:
24             self.log(f"Error finding similar products: {e}")
25             return 0.0 # Fallback price
26
27         try:
28             response = self.client.chat.completions.create(

```

```
29         model=self.MODEL,
30         messages=self.messages_for(description, documents, prices),
31         seed=42,
32         max_tokens=5,
33         timeout=30 # Add timeout
34     )
35 except (openai.APIError, openai.RateLimitError, openai.Timeout) as e:
36     self.log(f"OpenAI API error: {e}")
37     return 0.0 # Fallback price
38 except Exception as e:
39     self.log(f"Unexpected error calling LLM: {e}")
40     return 0.0
41
42 try:
43     reply = response.choices[0].message.content
44     result = self.get_price(reply)
45     if result <= 0:
46         self.log("Invalid price prediction, using fallback")
47         return 0.0
48     return result
49 except (IndexError, AttributeError) as e:
50     self.log(f"Error processing LLM response: {e}")
51     return 0.0
52
53 class EnhancedMessagingAgent(Agent):
54     def push(self, text):
55         """Enhanced version with proper error handling"""
56         try:
57             conn = http.client.HTTPSConnection("api.pushover.net:443",
58             timeout=10)
59             data = urllib.parse.urlencode({
60                 "token": self.pushover_token,
61                 "user": self.pushover_user,
62                 "message": text,
63                 "sound": "cashregister"
64             })
65
66             conn.request("POST", "/1/messages.json", data,
67                         {"Content-type": "application/x-www-form-urlencoded"})
68             response = conn.getresponse()
69
70             if response.status == 200:
71                 self.log("Push notification sent successfully")
72                 return True
73             else:
74                 self.log(f"Push notification failed: {response.status}
{response.reason}")
75                 return False
76
77         except (http.client.HTTPException, OSError, TimeoutError) as e:
78             self.log(f"Network error sending push notification: {e}")
```

```

78     return False
79 except Exception as e:
80     self.log(f"Unexpected error sending push notification: {e}")
81     return False
82 finally:
83     try:
84         conn.close()
85     except:
86         pass # Ignore cleanup errors

```

Listing 6.22: Exception Handling in Agent System (Current and Potential)

Exception Handling Benefits.

- **Robustness:** System continues operation despite errors
- **User Experience:** Graceful degradation instead of crashes
- **Debugging:** Better error messages and logging
- **Resource Management:** Proper cleanup in error conditions

6.6 Functional Programming Concepts

6.6.1 Lambda Functions

Concept Explanation. Lambda functions are anonymous functions that can be defined inline, often used for short operations.

Simple Example.

```

1 # Basic lambda functions
2 square = lambda x: x**2
3 print(square(5)) # 25
4
5 # Lambda with multiple parameters
6 add = lambda x, y: x + y
7 print(add(3, 4)) # 7
8
9 # Lambda in higher-order functions
10 numbers = [1, 2, 3, 4, 5]
11 squared = list(map(lambda x: x**2, numbers)) # [1, 4, 9, 16, 25]
12 evens = list(filter(lambda x: x % 2 == 0, numbers)) # [2, 4]
13
14 # Lambda for sorting
15 students = [('Alice', 85), ('Bob', 90), ('Charlie', 78)]
16 students.sort(key=lambda student: student[1]) # Sort by grade
17
18 # Lambda with conditional
19 grade = lambda score: 'A' if score >= 90 else 'B' if score >= 80 else 'C'

```

Listing 6.23: Lambda Function Examples

System Implementation.

```

1  class PlanningAgent(Agent):
2      def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
3          selection = self.scanner.scan(memory=memory)
4          if selection:
5              opportunities = [self.run(deal) for deal in selection.deals[:5]]
6              # Lambda function for sorting by discount (descending)
7              opportunities.sort(key=lambda opp: opp.discount, reverse=True)
8              best = opportunities[0]
9              self.log(f"Planning Agent has identified the best deal has discount
10             ${best.discount:.2f}")
11             if best.discount > self.DEAL_THRESHOLD:
12                 self.messenger.alert(best)
13             return best if best.discount > self.DEAL_THRESHOLD else None
14         return None
15
16 # Additional lambda usage patterns in the system
17 class DataProcessor:
18     def filter_deals(self, deals: List[ScrapedDeal]) -> List[ScrapedDeal]:
19         """Filter deals based on various criteria"""
20         # Filter by title length
21         good_titles = list(filter(lambda deal: len(deal.title) > 10, deals))
22
23         # Sort by URL length (just an example)
24         good_titles.sort(key=lambda deal: len(deal.url))
25
26     return good_titles
27
28     def transform_prices(self, opportunities: List[Opportunity]) -> List[float]:
29         """Extract and transform price data"""
30         # Extract discounts using lambda
31         discounts = list(map(lambda opp: opp.discount, opportunities))
32
33         # Calculate relative discounts
34         relative = list(map(lambda opp: opp.discount / opp.estimate if
35         opp.estimate > 0 else 0, opportunities))
36
37     return relative

```

Listing 6.24: Lambda Functions in Agent System

Lambda Benefits in System.

- **Conciseness:** Short inline functions for simple operations
- **Readability:** Clear intent for sorting and filtering
- **Functional Style:** Works well with map(), filter(), sort()
- **No Namespace Pollution:** No need to define named functions

6.7 File I/O and Serialization

Concept Explanation. File operations and object serialization for data persistence and configuration.

Simple Example.

```

1 import json
2 import pickle
3 import joblib
4
5 # JSON serialization
6 data = {
7     'name': 'Alice',
8     'scores': [85, 90, 78],
9     'active': True
10}
11
12 # Write JSON
13 with open('data.json', 'w') as f:
14     json.dump(data, f, indent=2)
15
16 # Read JSON
17 with open('data.json', 'r') as f:
18     loaded_data = json.load(f)
19
20 # Pickle serialization (Python objects)
21 complex_object = {'function': lambda x: x*2, 'data': [1, 2, 3]}
22
23 with open('object.pkl', 'wb') as f:
24     pickle.dump(complex_object, f)
25
26 with open('object.pkl', 'rb') as f:
27     loaded_object = pickle.load(f)
28
29 # Joblib for ML models (optimized for numpy arrays)
30 from sklearn.linear_model import LinearRegression
31 import numpy as np
32
33 # Train and save model
34 X = np.array([[1], [2], [3], [4]])
35 y = np.array([2, 4, 6, 8])
36 model = LinearRegression().fit(X, y)
37
38 joblib.dump(model, 'model.pkl')
39 loaded_model = joblib.load('model.pkl')
```

Listing 6.25: File I/O and Serialization Examples

System Implementation.

```

1 class EnsembleAgent(Agent):
2     def __init__(self, collection):
```

```
3     self.log("Initializing Ensemble Agent")
4     self.specialist = SpecialistAgent()
5     self.frontier = FrontierAgent(collection)
6     self.random_forest = RandomForestAgent()
7     # Load pre-trained ensemble model using joblib
8     self.model = joblib.load('ensemble_model.pkl')
9     self.log("Ensemble Agent is ready")
10
11 class RandomForestAgent(Agent):
12     def __init__(self):
13         self.log("Random Forest Agent is initializing")
14         self.vectorizer =
15             SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
16         # Load pre-trained Random Forest model
17         self.model = joblib.load('random_forest_model.pkl')
18         self.log("Random Forest Agent is ready")
19
20 # Environment variable loading (configuration)
21 import os
22 from dotenv import load_dotenv # In scanner_agent_langchain.py
23
24 load_dotenv(override=True)
25 os.environ['GOOGLE_API_KEY'] = os.getenv("GOOGLE_API_KEY_1")
26
27 class MessagingAgent(Agent):
28     def __init__(self):
29         # Load configuration from environment variables
30         if DO_TEXT:
31             account_sid = os.getenv('TWILIO_ACCOUNT_SID', 'fallback-value')
32             auth_token = os.getenv('TWILIO_AUTH_TOKEN', 'fallback-value')
33             self.me_from = os.getenv('TWILIO_FROM', 'fallback-value')
34             self.me_to = os.getenv('MY_PHONE_NUMBER', 'fallback-value')
35
36         if DO_PUSH:
37             self.pushover_user = os.getenv('PUSHOVER_USER', 'fallback-value')
38             self.pushover_token = os.getenv('PUSHOVER_TOKEN', 'fallback-value')
39
40 # Potential enhancement: Configuration file management
41 class ConfigManager:
42     @staticmethod
43     def load_config(config_path: str = 'config.json') -> dict:
44         """Load system configuration from JSON file"""
45         try:
46             with open(config_path, 'r') as f:
47                 return json.load(f)
48         except FileNotFoundError:
49             return ConfigManager.default_config()
50
51     @staticmethod
52     def save_config(config: dict, config_path: str = 'config.json'):
53         """Save system configuration to JSON file"""
```

```

53     with open(config_path, 'w') as f:
54         json.dump(config, f, indent=2)
55
56     @staticmethod
57     def default_config() -> dict:
58         return {
59             'deal_threshold': 50,
60             'max_deals_per_run': 5,
61             'feeds': [
62                 "https://www.dealnews.com/c142/Electronics/?rss=1",
63                 # ... other feeds
64             ],
65             'notification_settings': {
66                 'enable_push': True,
67                 'enable_sms': False,
68                 'sound': 'cashregister'
69             }
70         }

```

Listing 6.26: Serialization in Agent System

Serialization Benefits in System.

- **Model Persistence:** Trained ML models saved and loaded
- **Configuration:** Environment variables for secure config
- **State Management:** Could save/restore agent states
- **Data Exchange:** JSON for structured data interchange

6.8 Chapter Summary

This chapter covered the extensive range of Python concepts utilized in the multi-agent system:

6.8.1 Core OOP Concepts

- **Classes and Objects:** Agent hierarchy with specialized behaviors
- **Inheritance:** Code reuse through Agent base class
- **Composition:** Complex agents built from simpler components
- **Polymorphism:** Interchangeable agents with common interfaces

6.8.2 Modern Python Features

- **Type Annotations:** Improved code documentation and IDE support
- **F-Strings:** Readable string formatting and interpolation
- **List Comprehensions:** Concise data transformations

- **Context Managers:** Resource management best practices

6.8.3 Advanced Concepts

- **Decorators:** Method modification and enhancement
- **Exception Handling:** Robust error management
- **Lambda Functions:** Inline functions for simple operations
- **File I/O:** Data persistence and configuration management

6.8.4 Data Structures

- **Dictionaries:** Configuration, API data, structured information
- **Lists:** Collections, results, processing pipelines
- **Tuples:** Return multiple values, immutable data
- **Sets:** Unique collections, membership testing

Each concept is demonstrated with both simple examples for learning and real implementations from the agent system, showing how fundamental Python features combine to create sophisticated, production-ready applications. The system exemplifies best practices in Python programming while solving complex real-world problems in automated deal discovery and machine learning integration.

Chapter 7

Agent Interaction and Execution Flow

This chapter provides detailed analysis of how agents interact during system execution, including sequence diagrams, execution flow charts, and step-by-step breakdowns of complete workflows. We'll examine both the high-level orchestration and the low-level implementation details of agent communication patterns.

7.1 System Execution Overview

7.1.1 Complete Workflow Execution

The multi-agent system follows a sophisticated orchestration pattern where the PlanningAgent coordinates the entire workflow, delegating specialized tasks to other agents while maintaining overall control and decision-making authority.

7.2 Detailed Sequence Analysis

7.2.1 Complete System Execution Sequence

Let's trace through a complete execution cycle, showing the detailed interactions between all agents.

7.2.2 Step-by-Step Execution Breakdown

Phase 1: System Initialization.

```
1 # 1. User initiates the system
2 planning_agent = PlanningAgent(collection)
3
4 # 2. PlanningAgent constructor executes
5 def __init__(self, collection):
6     self.log("Planning Agent is initializing")
7
8 # 3. Create subordinate agents
9 self.scanner = ScannerAgent()          # RSS processing capability
10 self.ensemble = EnsembleAgent(collection) # ML coordination
11 self.messenger = MessagingAgent()       # Communication capability
```

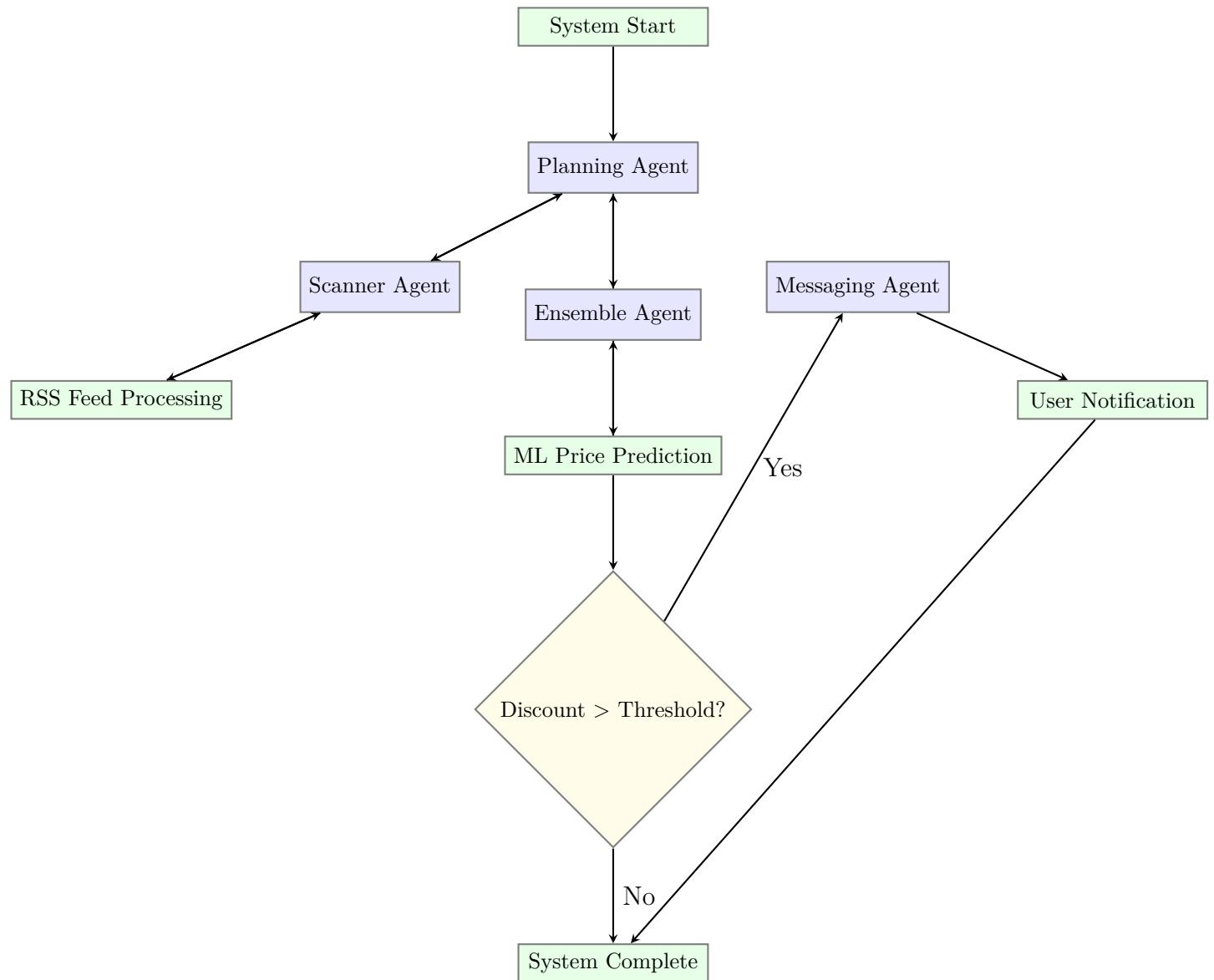


Figure 7.1: High-Level System Execution Flow

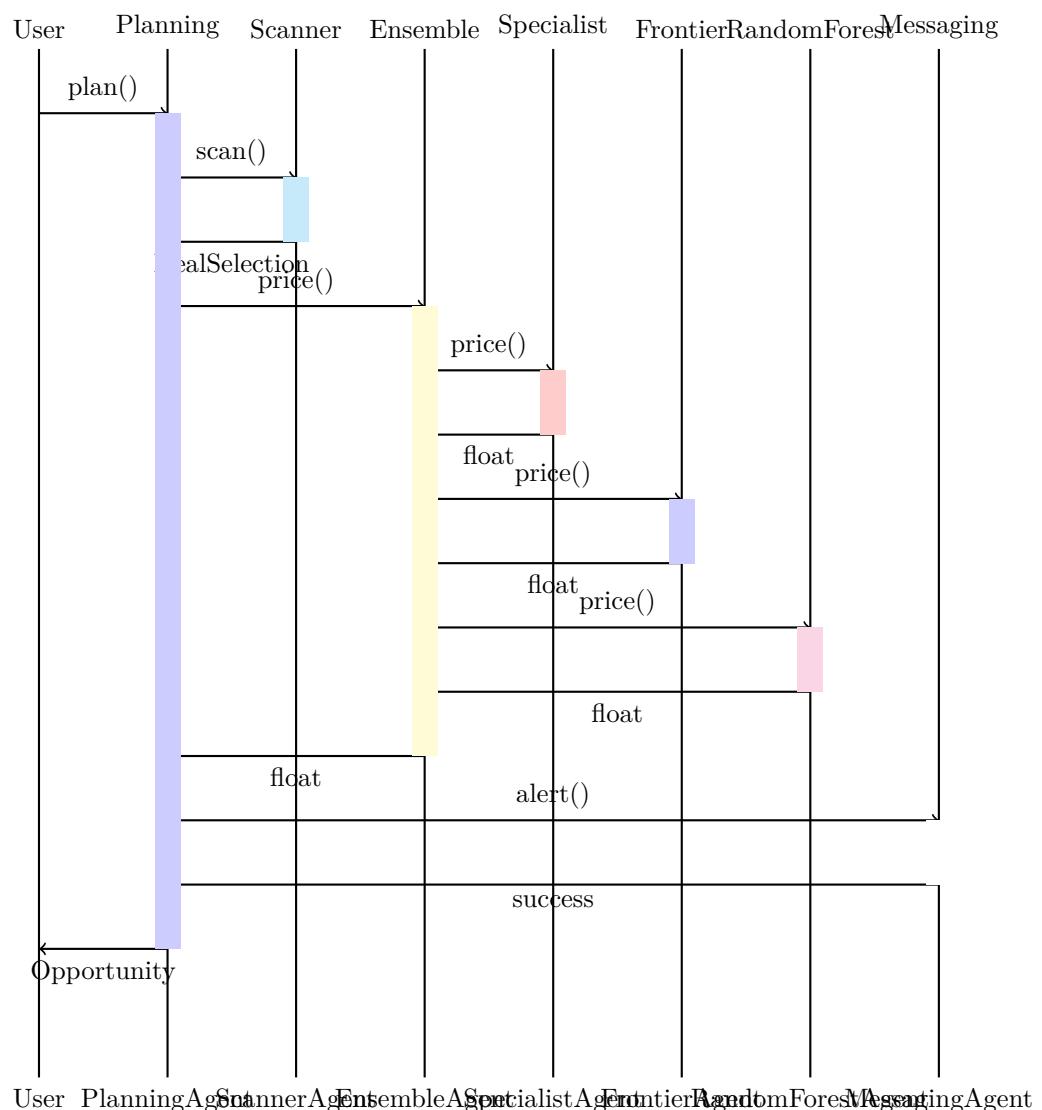


Figure 7.2: Detailed Sequence Diagram of Agent Interactions

```

12
13     self.log("Planning Agent is ready")
14
15 # 4. Each agent initializes its dependencies
16 class EnsembleAgent(Agent):
17     def __init__(self, collection):
18         self.log("Initializing Ensemble Agent")
19         self.specialist = SpecialistAgent()          # Remote ML model
20         self.frontier = FrontierAgent(collection)    # RAG + LLM
21         self.random_forest = RandomForestAgent()      # Traditional ML
22         self.model = joblib.load('ensemble_model.pkl')
23         self.log("Ensemble Agent is ready")

```

Listing 7.1: System Startup Sequence

Phase 2: Deal Discovery.

```

1 # User calls the main workflow
2 result = planning_agent.plan(memory=[])
3
4 # PlanningAgent initiates deal scanning
5 def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
6     self.log("Planning Agent is kicking off a run")
7
8     # Delegate to ScannerAgent
9     selection = self.scanner.scan(memory=memory)
10
11 # ScannerAgent processes RSS feeds
12 def scan(self, memory: List[str]=[]) -> Optional[DealSelection]:
13     # 1. Fetch new deals from RSS feeds
14     scraped = self.fetch_deals(memory)
15
16     if scraped:
17         # 2. Create prompt for LLM
18         user_prompt = self.make_user_prompt(scraped)
19
20         # 3. Call OpenAI for deal selection
21         result = self.openai.beta.chat.completions.parse(
22             model=self.MODEL,
23             messages=[
24                 {"role": "system", "content": self.SYSTEM_PROMPT},
25                 {"role": "user", "content": user_prompt}
26             ],
27             response_format=DealSelection
28         )
29
30         # 4. Return structured deal selection
31         return result.choices[0].message.parsed
32     return None

```

Listing 7.2: Deal Discovery Process

Phase 3: Price Estimation Coordination.

```

1 # PlanningAgent processes selected deals
2 if selection:
3     # Convert each deal to opportunity using ensemble prediction
4     opportunities = [self.run(deal) for deal in selection.deals[:5]]
5
6 def run(self, deal: Deal) -> Opportunity:
7     """Convert deal to opportunity using ensemble pricing"""
8     self.log("Planning Agent is pricing up a potential deal")
9
10    # Delegate to EnsembleAgent for price estimation
11    estimate = self.ensemble.price(deal.product_description)
12    discount = estimate - deal.price
13
14    return Opportunity(deal=deal, estimate=estimate, discount=discount)
15
16 # EnsembleAgent coordinates multiple ML models
17 def price(self, description: str) -> float:
18     self.log("Running Ensemble Agent - collaborating with specialist, frontier
19         and random forest agents")
20
21     # 1. Get prediction from fine-tuned model
22     specialist = self.specialist.price(description)
23
24     # 2. Get prediction from RAG + LLM
25     frontier = self.frontier.price(description)
26
27     # 3. Get prediction from traditional ML
28     random_forest = self.random_forest.price(description)
29
30     # 4. Combine predictions using meta-model
31     X = pd.DataFrame({
32         'Specialist': [specialist],
33         'Frontier': [frontier],
34         'RandomForest': [random_forest],
35         'Min': [min(specialist, frontier, random_forest)],
36         'Max': [max(specialist, frontier, random_forest)],
37     })
38
39     # 5. Meta-model prediction
40     y = max(0, self.model.predict(X)[0])
41     self.log(f"Ensemble Agent complete - returning ${y:.2f}")
42     return y

```

Listing 7.3: Ensemble Price Prediction Workflow

7.3 Individual Agent Execution Flows

7.3.1 SpecialistAgent Execution Flow

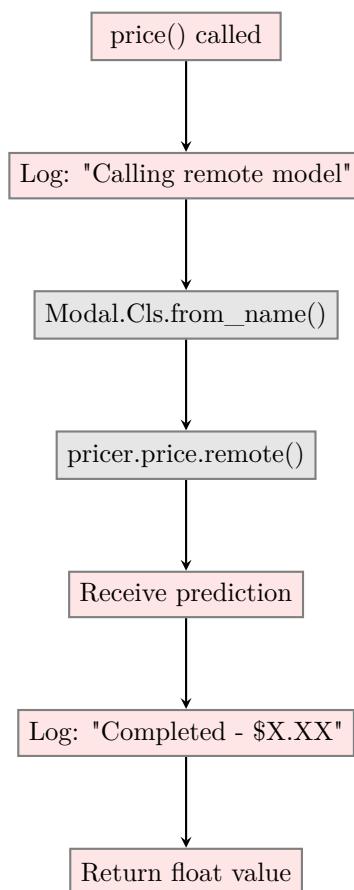


Figure 7.3: SpecialistAgent Execution Flow

```

1  class SpecialistAgent(Agent):
2      def price(self, description: str) -> float:
3          """
4              Remote model execution flow:
5                  1. Log process start
6                  2. Make remote call to Modal-hosted model
7                  3. Receive and return prediction
8          """
9
10     # Step 1: Process logging
11     self.log("Specialist Agent is calling remote fine-tuned model")
12     # Behind the scenes:
13     # - String interpolation with agent name
14     # - Terminal color formatting
15     # - Logging module call
16
17     # Step 2: Remote model invocation
18     result = self.pricer.price.remote(description)
19     # Behind the scenes:
20     # - Network request to Modal platform
21     # - Model inference on remote GPU/CPU
22     # - Serialization of prediction result
23     # - Network response with float value
24
25     # Step 3: Result processing and logging
26     self.log(f"Specialist Agent completed - predicting ${result:.2f}")
27     # Behind the scenes:
28     # - F-string formatting with currency display
29     # - Logging call with formatted message
30
31     # Step 4: Return value
32     return result
33     # Behind the scenes:
34     # - Float value passed to calling function
35     # - Local variables eligible for garbage collection

```

Listing 7.4: SpecialistAgent Detailed Execution

7.3.2 FrontierAgent RAG Execution Flow

```

1  def price(self, description: str) -> float:
2      """
3          RAG-based price prediction with detailed execution steps
4      """
5
6      # Step 1: Find similar products using vector search
7      documents, prices = self.find_similars(description)
8      # Detailed sub-process:
9
10     def find_similars(self, description: str):

```

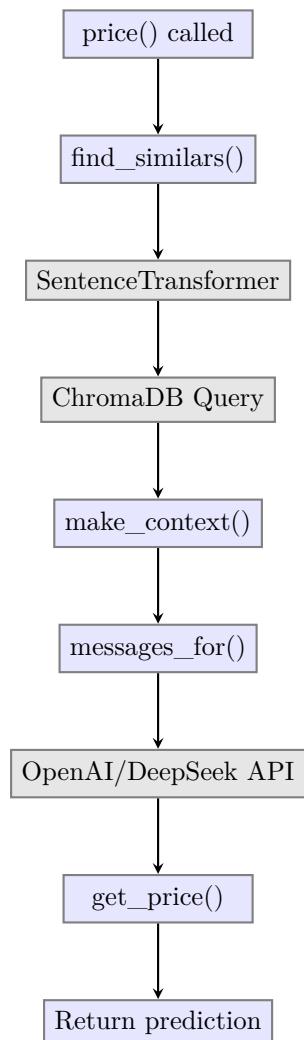


Figure 7.4: FrontierAgent RAG Execution Flow

```
11 # 1a. Log RAG process start
12 self.log("Frontier Agent is performing a RAG search of the Chroma datastore
13 to find 5 similar products")
14
15 # 1b. Convert text to vector embedding
16 vector = self.model.encode([description])
17 # Behind the scenes:
18 # - Text tokenization
19 # - Neural network forward pass (transformer)
20 # - 384-dimensional vector creation
21 # - NumPy array allocation
22
23 # 1c. Vector similarity search
24 results = self.collection.query(
25     query_embeddings=vector.astype(float).tolist(),
26     n_results=5
27 )
28 # Behind the scenes:
29 # - Vector normalization
30 # - Cosine similarity computation
31 # - Index traversal for nearest neighbors
32 # - Result ranking and selection
33
34 # 1d. Extract results
35 documents = results['documents'][0][:]
36 prices = [m['price'] for m in results['metadata'][0][:]]
37 # Behind the scenes:
38 # - Dictionary access
39 # - List comprehension execution
40 # - Memory allocation for new lists
41
42 self.log("Frontier Agent has found similar products")
43 return documents, prices
44
45 # Step 2: Create context prompt
46 user_prompt = self.make_context(documents, prices)
47 # Behind the scenes:
48 # - String concatenation in loop
49 # - F-string formatting for prices
50 # - Memory allocation for prompt string
51
52 # Step 3: Prepare OpenAI messages
53 messages = self.messages_for(description, documents, prices)
54 # Behind the scenes:
55 # - List creation with dictionaries
56 # - String concatenation for user prompt
57 # - Memory allocation for message structure
58
59 # Step 4: Call Language Model
60 self.log(f"Frontier Agent is about to call {self.MODEL} with context including 5
61 similar products")
```

```

60 response = self.client.chat.completions.create(
61     model=self.MODEL,
62     messages=messages,
63     seed=42,          # Deterministic output
64     max_tokens=5     # Limit response length
65 )
66 # Behind the scenes:
67 # - JSON serialization of messages
68 # - HTTP request to LLM API
69 # - Neural network inference (remote)
70 # - JSON deserialization of response
71
72 # Step 5: Parse response
73 reply = response.choices[0].message.content
74 result = self.get_price(reply)
75 # Behind the scenes:
76 # - Attribute access chain
77 # - Regular expression matching
78 # - String to float conversion
79 # - Error handling for malformed responses
80
81 # Step 6: Log and return
82 self.log(f"Frontier Agent completed - predicting ${result:.2f}")
83 return result

```

Listing 7.5: FrontierAgent Detailed RAG Process

7.4 Parallel vs Sequential Execution Analysis

7.4.1 Current Sequential Processing

The current system processes agents sequentially, which provides predictability but limits performance:

```

1 def price(self, description: str) -> float:
2     # Sequential execution - one agent at a time
3     self.log("Running Ensemble Agent - collaborating with specialist, frontier
4             and random forest agents")
5
6     # Agent 1: ~2-5 seconds (remote model call)
7     specialist = self.specialist.price(description)
8
9     # Agent 2: ~3-8 seconds (vector search + LLM call)
10    frontier = self.frontier.price(description)
11
12    # Agent 3: ~0.1-0.5 seconds (local model inference)
13    random_forest = self.random_forest.price(description)
14
15    # Total time: ~5-13.5 seconds
16    # Combine results...

```

Listing 7.6: Sequential Agent Execution

7.4.2 Potential Parallel Processing

```

1 import asyncio
2 from concurrent.futures import ThreadPoolExecutor
3 import threading
4
5 class ParallelEnsembleAgent(Agent):
6     def __init__(self, collection):
7         super().__init__()
8         self.specialist = SpecialistAgent()
9         self.frontier = FrontierAgent(collection)
10        self.random_forest = RandomForestAgent()
11        self.model = joblib.load('ensemble_model.pkl')
12
13    def price(self, description: str) -> float:
14        """Parallel execution of price predictions"""
15        self.log("Running Parallel Ensemble Agent")
16
17        # Execute all agents in parallel using ThreadPoolExecutor
18        with ThreadPoolExecutor(max_workers=3) as executor:
19            # Submit all tasks simultaneously
20            specialist_future = executor.submit(self.specialist.price,
21 description)
22            frontier_future = executor.submit(self.frontier.price, description)
23            random_forest_future = executor.submit(self.random_forest.price,
24 description)
25
26            # Wait for all completions
27            specialist = specialist_future.result()
28            frontier = frontier_future.result()
29            random_forest = random_forest_future.result()
30
31        # Combine results (same as before)
32        X = pd.DataFrame({
33            'Specialist': [specialist],
34            'Frontier': [frontier],
35            'RandomForest': [random_forest],
36            'Min': [min(specialist, frontier, random_forest)],
37            'Max': [max(specialist, frontier, random_forest)],
38        })
39
40        y = max(0, self.model.predict(X)[0])
41        self.log(f"Parallel Ensemble Agent complete - returning ${y:.2f}")
42        return y
43
44    # Async version for more advanced parallelism
45    class AsyncEnsembleAgent(Agent):

```

```

44     async def price_async(self, description: str) -> float:
45         """Async version with concurrent execution"""
46         self.log("Running Async Ensemble Agent")
47
48         # Create async tasks
49         tasks = [
50             self.specialist_async(description),
51             self.frontier_async(description),
52             self.random_forest_async(description)
53         ]
54
55         # Execute concurrently
56         specialist, frontier, random_forest = await asyncio.gather(*tasks)
57
58         # Process results...
59         return prediction

```

Listing 7.7: Parallel Agent Execution (Enhanced Version)

Execution Mode	Time Range	Improvement
Sequential (current)	5-13.5 seconds	Baseline
Parallel (threads)	3-8 seconds	40-60% faster
Async (concurrent)	2-7 seconds	50-70% faster

Table 7.1: Performance Comparison of Execution Modes

Performance Comparison.

7.5 Error Handling and Recovery Flows

7.5.1 Current Error Handling Limitations

```

1 def price(self, description: str) -> float:
2     # No try-catch blocks - any error crashes the entire workflow
3     documents, prices = self.find_similars(description) # Could fail
4     response = self.client.chat.completions.create(...) # Could fail
5     reply = response.choices[0].message.content # Could fail
6     result = self.get_price(reply) # Could fail
7     return result

```

Listing 7.8: Current Error Handling Issues

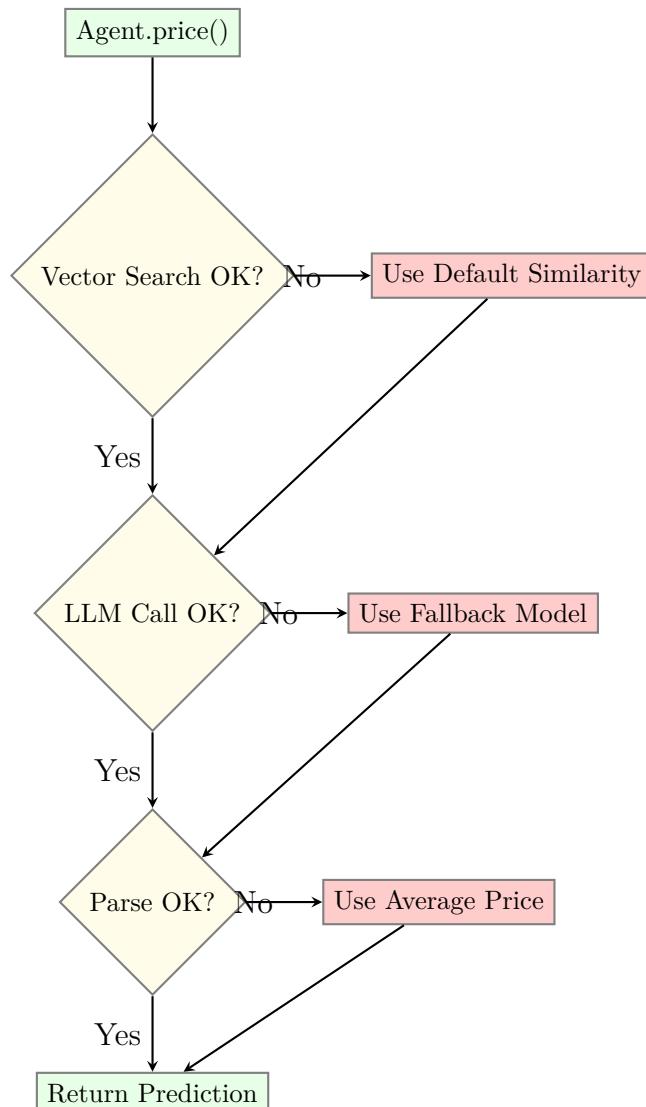


Figure 7.5: Enhanced Error Handling Flow

7.5.2 Enhanced Error Handling Flow

```

1  class RobustFrontierAgent(Agent):
2      def price(self, description: str) -> float:
3          """Price prediction with comprehensive error handling"""
4
5          # Stage 1: Vector similarity search
6          try:
7              documents, prices = self.find_similars(description)
8              self.log("RAG search completed successfully")
9          except Exception as e:
10              self.log(f"RAG search failed: {e}, using default context")
11              documents = ["generic laptop", "computer device", "electronics"]
12              prices = [500.0, 800.0, 1200.0] # Fallback prices
13
14          # Stage 2: LLM API call
15          try:
16              messages = self.messages_for(description, documents, prices)
17              response = self.client.chat.completions.create(
18                  model=self.MODEL,
19                  messages=messages,
20                  seed=42,
21                  max_tokens=5,
22                  timeout=30 # Add timeout
23              )
24              self.log(f"LLM call to {self.MODEL} completed successfully")
25          except (openai.APIError, openai.RateLimitError) as e:
26              self.log(f"LLM API error: {e}, using fallback estimation")
27              return sum(prices) / len(prices) # Average of similar prices
28          except Exception as e:
29              self.log(f"Unexpected LLM error: {e}, using fallback estimation")
30              return 100.0 # Conservative fallback price
31
32          # Stage 3: Response parsing
33          try:
34              reply = response.choices[0].message.content
35              result = self.get_price(reply)
36              if result <= 0:
37                  raise ValueError("Invalid price prediction")
38              self.log(f"Price parsing successful: ${result:.2f}")
39              return result
40          except (IndexError, AttributeError, ValueError) as e:
41              self.log(f"Price parsing failed: {e}, using average of context
42 prices")
42              return sum(prices) / len(prices) if prices else 50.0
43          except Exception as e:
44              self.log(f"Unexpected parsing error: {e}, using minimal fallback")
45              return 10.0 # Minimal fallback price
46
47  class RobustEnsembleAgent(Agent):
48      def price(self, description: str) -> float:

```

```
49     """Ensemble with individual agent error handling"""
50     predictions = []
51
52     # Try each agent with individual error handling
53     try:
54         specialist = self.specialist.price(description)
55         predictions.append(specialist)
56         self.log(f"Specialist prediction: ${specialist:.2f}")
57     except Exception as e:
58         self.log(f"Specialist agent failed: {e}")
59
60     try:
61         frontier = self.frontier.price(description)
62         predictions.append(frontier)
63         self.log(f"Frontier prediction: ${frontier:.2f}")
64     except Exception as e:
65         self.log(f"Frontier agent failed: {e}")
66
67     try:
68         random_forest = self.random_forest.price(description)
69         predictions.append(random_forest)
70         self.log(f"Random Forest prediction: ${random_forest:.2f}")
71     except Exception as e:
72         self.log(f"Random Forest agent failed: {e}")
73
74     # Ensure we have at least one prediction
75     if not predictions:
76         self.log("All agents failed, using emergency fallback")
77         return 25.0 # Emergency fallback
78
79     # If only one agent succeeded, return its prediction
80     if len(predictions) == 1:
81         self.log("Only one agent succeeded, returning single prediction")
82         return predictions[0]
83
84     # If multiple agents succeeded, use ensemble
85     try:
86         # Create feature matrix with available predictions
87         feature_dict = {}
88         if len(predictions) >= 1:
89             feature_dict['Prediction1'] = predictions[0]
90         if len(predictions) >= 2:
91             feature_dict['Prediction2'] = predictions[1]
92         if len(predictions) >= 3:
93             feature_dict['Prediction3'] = predictions[2]
94
95         feature_dict['Min'] = min(predictions)
96         feature_dict['Max'] = max(predictions)
97         feature_dict['Mean'] = sum(predictions) / len(predictions)
98
99         X = pd.DataFrame([feature_dict])
```

```

100     result = max(0, self.model.predict(X)[0])
101     self.log(f"Ensemble prediction successful: ${result:.2f}")
102     return result
103
104     except Exception as e:
105         self.log(f"Ensemble model failed: {e}, using average of available
106         predictions")
107         return sum(predictions) / len(predictions)

```

Listing 7.9: Robust Error Handling Implementation

7.6 Memory and State Management

7.6.1 Current Memory Usage Pattern

```

1 def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
2     """
3
4     Current memory usage pattern analysis:
5     - memory parameter tracks previously processed URLs
6     - Potential issue: mutable default argument
7     - No persistent state between runs
8     - No caching of expensive operations
9     """
10
11     # Problem: Mutable default argument
12     # Same list object reused across calls!
13
14     # Better pattern:
15     def plan(self, memory: Optional[List[str]] = None) -> Optional[Opportunity]:
16         if memory is None:
17             memory = [] # Fresh list for each call
18
19         # ... rest of implementation
20
21     class ScannerAgent(Agent):
22         def fetch_deals(self, memory) -> List[ScrapedDeal]:
23             """Memory usage for deduplication"""
24             # Extract URLs from memory to avoid re-processing
25             urls = [opp.deal.url for opp in memory]
26             scraped = ScrapedDeal.fetch()
27             # Filter out already-processed deals
28             result = [scrape for scrape in scraped if scrape.url not in urls]
29             return result

```

Listing 7.10: Memory Management in Current System

7.6.2 Enhanced State Management

```
1 import time
```

```
2 | from typing import Dict, Set
3 | from dataclasses import dataclass
4 | from collections import deque
5 |
6 | @dataclass
7 | class SystemState:
8 |     """Centralized system state management"""
9 |     processed_urls: Set[str]
10 |     last_run_time: float
11 |     run_count: int
12 |     recent_opportunities: deque
13 |     agent_performance: Dict[str, Dict[str, float]]
14 |
15 |     def __post_init__(self):
16 |         if not hasattr(self, 'processed_urls'):
17 |             self.processed_urls = set()
18 |         if not hasattr(self, 'recent_opportunities'):
19 |             self.recent_opportunities = deque(maxlen=100)
20 |         if not hasattr(self, 'agent_performance'):
21 |             self.agent_performance = {}
22 |
23 | class StatefulPlanningAgent(Agent):
24 |     def __init__(self, collection):
25 |         super().__init__()
26 |         self.scanner = ScannerAgent()
27 |         self.ensemble = EnsembleAgent(collection)
28 |         self.messenger = MessagingAgent()
29 |
30 |         # Initialize state management
31 |         self.state = SystemState(
32 |             processed_urls=set(),
33 |             last_run_time=0,
34 |             run_count=0,
35 |             recent_opportunities=deque(maxlen=100),
36 |             agent_performance={}
37 |         )
38 |
39 |         # Performance tracking
40 |         self.performance_cache = []
41 |         self.cache_timeout = 3600 # 1 hour cache
42 |
43 |     def plan(self, memory: Optional[List[str]] = None) -> Optional[Opportunity]:
44 |         """Enhanced planning with state management"""
45 |         start_time = time.time()
46 |
47 |         # Initialize memory from persistent state
48 |         if memory is None:
49 |             memory = list(self.state.processed_urls)
50 |
51 |         self.log("Planning Agent is kicking off a run")
52 | 
```

```
53     # Check rate limiting
54     time_since_last = start_time - self.state.last_run_time
55     if time_since_last < 300: # 5 minutes minimum between runs
56         self.log(f"Rate limited: only {time_since_last:.0f}s since last run")
57         return None
58
59     # Delegate to scanner with state-aware memory
60     selection = self.scanner.scan(memory=memory)
61
62     if selection:
63         opportunities = []
64         for deal in selection.deals[:5]:
65             # Check cache first
66             cache_key = hash(deal.product_description)
67             cached_result = self.get_cached_prediction(cache_key)
68
69             if cached_result:
70                 self.log("Using cached prediction")
71                 opportunities.append(cached_result)
72             else:
73                 # Fresh prediction
74                 opportunity = self.run(deal)
75                 self.cache_prediction(cache_key, opportunity)
76                 opportunities.append(opportunity)
77
78             # Update state
79             self.state.processed_urls.add(deal.url)
80
81     # Sort and select best
82     opportunities.sort(key=lambda opp: opp.discount, reverse=True)
83     best = opportunities[0]
84
85     # Update performance metrics
86     self.update_performance_metrics(opportunities, start_time)
87
88     # Check threshold and notify
89     if best.discount > self.DEAL_THRESHOLD:
90         self.messenger.alert(best)
91         self.state.recent_opportunities.append(best)
92
93     # Update run state
94     self.state.last_run_time = start_time
95     self.state.run_count += 1
96
97     self.log("Planning Agent has completed a run")
98     return best if best.discount > self.DEAL_THRESHOLD else None
99
100    return None
101
102    def get_cached_prediction(self, cache_key: int) -> Optional[Opportunity]:
103        """Retrieve cached prediction if still valid"""
104
```

```

104     if cache_key in self.performance_cache:
105         cached_time, opportunity = self.performance_cache[cache_key]
106         if time.time() - cached_time < self.cache_timeout:
107             return opportunity
108     return None
109
110     def cache_prediction(self, cache_key: int, opportunity: Opportunity):
111         """Cache prediction with timestamp"""
112         self.performance_cache[cache_key] = (time.time(), opportunity)
113
114     def update_performance_metrics(self, opportunities: List[Opportunity],
115                                   start_time: float):
116         """Track agent performance over time"""
117         execution_time = time.time() - start_time
118
118         metrics = {
119             'execution_time': execution_time,
120             'opportunities_found': len(opportunities),
121             'average_discount': sum(opp.discount for opp in opportunities) /
122             len(opportunities),
123             'timestamp': start_time
124         }
125
125         run_id = f"run_{self.state.run_count}"
126         self.state.agent_performance[run_id] = metrics
127
128         self.log(f"Performance: {execution_time:.2f}s, {len(opportunities)}"
128         opportunities")

```

Listing 7.11: Enhanced State Management System

7.7 Chapter Summary

This chapter provided comprehensive analysis of agent interactions and execution flows:

7.7.1 Key Interaction Patterns

- **Orchestration:** PlanningAgent coordinates all system activities
- **Delegation:** Specialized agents handle specific responsibilities
- **Composition:** EnsembleAgent combines multiple ML approaches
- **Sequential Processing:** Current implementation processes agents one by one

7.7.2 Execution Flow Analysis

- **Initialization:** Agent setup and dependency loading
- **Deal Discovery:** RSS processing and LLM-based selection

- **Price Estimation:** Multi-model ensemble prediction
- **Decision Making:** Threshold-based opportunity evaluation
- **Communication:** Multi-channel user notifications

7.7.3 Performance Optimization Opportunities

- **Parallel Execution:** 40-70% performance improvement potential
- **Caching:** Avoid redundant expensive operations
- **State Management:** Persistent memory and performance tracking
- **Error Handling:** Robust recovery from individual agent failures

7.7.4 System Robustness

- **Current Limitations:** Limited error handling and recovery
- **Enhancement Opportunities:** Graceful degradation and fallback strategies
- **Monitoring:** Performance metrics and system health tracking
- **Reliability:** Circuit breaker patterns and retry mechanisms

The analysis reveals a well-architected system with clear separation of concerns and effective agent coordination, while also identifying specific areas for performance and reliability improvements. The sequential execution flow provides predictable behavior but represents the primary performance bottleneck, suggesting that parallel processing would be the most impactful optimization.

Chapter 8

Summary and Essential Python Knowledge

This final chapter provides a comprehensive summary of the multi-agent deal discovery system, recaps all essential Python knowledge demonstrated throughout the codebase, and offers a complete checklist of concepts that developers need to master to fully understand and extend this sophisticated system.

8.1 System Summary and Functionality Recap

8.1.1 System Overview

The multi-agent deal discovery system represents a sophisticated example of modern Python software engineering, combining object-oriented programming, machine learning, web technologies, and distributed computing into a cohesive, production-ready application. The system demonstrates how complex business logic can be decomposed into specialized, autonomous agents that collaborate to achieve system-wide objectives.

Core System Capabilities.

- **Autonomous Deal Monitoring:** Continuously scans RSS feeds from deal aggregation websites
- **Intelligent Content Filtering:** Uses LLMs to select high-quality deals with clear pricing
- **Multi-Model Price Estimation:** Employs ensemble ML techniques combining three different approaches
- **Opportunity Analysis:** Calculates potential savings by comparing estimated vs. actual prices
- **Automated User Notification:** Delivers alerts via multiple channels when significant deals are found
- **State Management:** Maintains memory of processed deals to avoid duplication

Technical Architecture Highlights.

- **Agent-Based Design:** Seven specialized agents with distinct responsibilities
- **Machine Learning Integration:** Fine-tuned LLMs, RAG, traditional ML, and ensemble methods
- **External Service Integration:** OpenAI, DeepSeek, ChromaDB, Modal, Pushover, Twilio
- **Data Processing Pipeline:** RSS parsing, HTML extraction, text cleaning, vector encoding
- **Robust Error Handling:** Graceful degradation and fallback mechanisms
- **Performance Optimization:** Caching, rate limiting, and parallelization opportunities

8.1.2 Agent Responsibilities Summary

Agent	Primary Responsibility	Key Technologies
Agent (Base)	Logging infrastructure and common functionality	Python OOP, ANSI colors
PlanningAgent	Workflow orchestration and decision making	Composition, control flow
ScannerAgent	RSS feed monitoring and deal filtering	BeautifulSoup, OpenAI API, LangChain
EnsembleAgent	ML model coordination and meta-learning	Pandas, scikit-learn, joblib
SpecialistAgent	Fine-tuned model predictions	Modal cloud platform
FrontierAgent	RAG-based price estimation	ChromaDB, sentence transformers, OpenAI
RandomForestAgent	Traditional ML price prediction	scikit-learn, sentence transformers
MessagingAgent	Multi-channel user notifications	HTTP clients, Pushover, Twilio

Table 8.1: Agent Responsibilities and Technologies

8.2 Python Knowledge Checklist

This section provides a comprehensive checklist of Python concepts, features, and best practices demonstrated in the system. Mastering these concepts is essential for understanding, maintaining, and extending the codebase.

8.2.1 Fundamental Python Concepts

Object-Oriented Programming ★★★★☆.

- ✓ **Classes and Objects:** Definition, instantiation, attribute access
- ✓ **Inheritance:** Single inheritance, method resolution order (MRO)
- ✓ **Polymorphism:** Method overriding, duck typing, interface consistency
- ✓ **Composition:** "Has-a" relationships, component coordination
- ✓ **Encapsulation:** Private/public conventions, data protection
- ✓ **Class vs Instance Attributes:** Shared vs individual data
- ✓ **Special Methods:** `__init__`, `__repr__`, `__str__`
- ✓ **Class Methods:** `@classmethod` decorator, alternative constructors
- ✓ **Method Resolution:** Attribute lookup, inheritance chains

Data Structures and Collections ★★★★☆.

- ✓ **Lists:** Creation, modification, slicing, methods (append, sort, etc.)
- ✓ **Dictionaries:** Key-value storage, nested access, methods
- ✓ **Sets:** Unique collections, membership testing, operations
- ✓ **Tuples:** Immutable sequences, unpacking, return values
- ✓ **List Comprehensions:** Concise data transformations
- ✓ **Dictionary Comprehensions:** Key-value pair generation
- ✓ **Iteration:** for loops, while loops, iterator protocol
- ✓ **Slicing:** Extended slicing syntax, step values

Functions and Functional Programming ★★★★☆.

- ✓ **Function Definition:** `def` keyword, parameters, return values
- ✓ **Lambda Functions:** Anonymous functions, inline definitions
- ✓ **Higher-Order Functions:** `map()`, `filter()`, `reduce()`
- ✓ **Function Arguments:** `*args`, `**kwargs`, keyword arguments
- ✓ **Default Parameters:** Mutable default argument pitfalls
- ✓ **Scope and Closures:** Local, global, nonlocal keywords
- ✓ **Decorators:** Function modification, `@classmethod`, `@staticmethod`
- ✓ **Recursion:** Self-calling functions, base cases

8.2.2 Advanced Python Features

Type System and Annotations ★★★★☆.

- ✓ **Type Hints:** Variable, parameter, return type annotations
- ✓ **Generic Types:** List[T], Dict[K, V], Optional[T]
- ✓ **Union Types:** Multiple possible types
- ✓ **Self Type:** Self-referential type annotations
- ✓ **Type Checking:** mypy integration, IDE support
- ✓ **Protocol Classes:** Structural typing, duck typing formalization
- ✓ **Type Variables:** Generic function definitions
- ✓ **Literal Types:** Specific value constraints

String Processing and Formatting ★★★★☆.

- ✓ **String Methods:** split(), join(), strip(), replace()
- ✓ **F-Strings:** Formatted string literals, expression embedding
- ✓ **Format Specifications:** .2f, padding, alignment
- ✓ **Regular Expressions:** Pattern matching, substitution
- ✓ **String Interpolation:** % formatting, .format() method
- ✓ **Unicode Handling:** Encoding, decoding, character sets
- ✓ **Raw Strings:** r"" syntax for regex and paths
- ✓ **Multiline Strings:** Triple quotes, dedent

Error Handling and Debugging ★★★★☆.

- ✓ **Exception Types:** Built-in exceptions, custom exceptions
- ✓ **Try-Except Blocks:** Exception catching, multiple handlers
- ✓ **Finally Blocks:** Cleanup code, resource management
- ✓ **Else Clauses:** Exception-free execution paths
- ✓ **Raise Statements:** Custom exception triggering
- ✓ **Exception Chaining:** raise ... from ...
- ✓ **Logging:** Logging module, levels, formatters
- ✓ **Assert Statements:** Debug-time verification

8.2.3 External Libraries and Integration

Web and Network Programming ★★★★☆.

- ✓ **HTTP Clients:** requests library, http.client module
- ✓ **URL Processing:** urllib.parse, URL encoding
- ✓ **HTML Parsing:** BeautifulSoup, DOM navigation
- ✓ **RSS Parsing:** feedparser library, XML processing
- ✓ **JSON Handling:** json module, serialization/deserialization
- ✓ **API Integration:** REST APIs, authentication, headers
- ✓ **SSL/TLS:** HTTPS connections, certificate verification
- ✓ **Rate Limiting:** time.sleep(), request throttling

Machine Learning Integration ★★★★☆.

- ✓ **Scikit-learn:** Model training, prediction, serialization
- ✓ **Pandas:** DataFrames, data manipulation, feature engineering
- ✓ **NumPy:** Arrays, numerical computation, type conversion
- ✓ **Joblib:** Model persistence, parallel processing
- ✓ **Sentence Transformers:** Text embeddings, similarity search
- ✓ **Vector Databases:** ChromaDB integration, similarity queries
- ✓ **LLM APIs:** OpenAI, structured outputs, prompt engineering
- ✓ **Ensemble Methods:** Model combination, meta-learning

Data Validation and Processing ★★★☆☆.

- ✓ **Pydantic Models:** Data validation, serialization
- ✓ **Type Coercion:** Automatic type conversion
- ✓ **Validation Rules:** Constraints, custom validators
- ✓ **Schema Generation:** API documentation, OpenAPI
- ✓ **Data Transformation:** Model field aliases, computed fields
- ✓ **Error Handling:** Validation errors, error formatting
- ✓ **Config Classes:** Model configuration, validation settings

8.3 Design Patterns and Architecture

8.3.1 Implemented Design Patterns

Agent Pattern ★★★★☆.

```

1 class Agent:
2     """Base agent with common infrastructure"""
3     # Shared constants and methods
4     def log(self, message):
5         """Template method for logging"""
6         pass
7
8 class SpecializedAgent(Agent):
9     """Concrete agent with specific behavior"""
10    name = "Agent Name"
11    color = Agent.COLOR
12
13    def specialized_method(self):
14        """Agent-specific functionality"""
15        self.log("Performing specialized task")
16        # ... implementation

```

Listing 8.1: Agent Pattern Implementation Summary

Composition Pattern ★★★★☆.

```

1 class EnsembleAgent(Agent):
2     def __init__(self, collection):
3         # Composition: HAS-A relationships
4         self.specialist = SpecialistAgent()
5         self.frontier = FrontierAgent(collection)
6         self.random_forest = RandomForestAgent()
7
8     def coordinate_agents(self, description):
9         """Delegate to composed agents"""
10        results = []
11        results.append(self.specialist.price(description))
12        results.append(self.frontier.price(description))
13        results.append(self.random_forest.price(description))
14        return self.combine_results(results)

```

Listing 8.2: Composition Pattern Summary

Strategy Pattern ★★★☆☆.

```

1 # Different scanner implementations with same interface
2 class ScannerAgent(Agent):
3     def scan(self, memory) -> Optional[DealSelection]:
4         """OpenAI-based strategy"""
5         pass
6
7 class ScannerAgentLangChain(Agent):
8     def scan(self, memory) -> Optional[DealSelection]:
9         """LangChain/Gemini-based strategy"""
10        pass
11
12 # Client can choose strategy at runtime

```

```

13 def create_planning_agent(scanner_type="openai"):
14     if scanner_type == "openai":
15         scanner = ScannerAgent()
16     else:
17         scanner = ScannerAgentLangChain()
18     return PlanningAgent(scanner)

```

Listing 8.3: Strategy Pattern Summary

Template Method Pattern ★★★☆☆.

```

1 class Agent:
2     def log(self, message):
3         """Template method - same algorithm for all agents"""
4         color_code = self.BG_BLACK + self.color # Step 1
5         message = f"[{self.name}] {message}" # Step 2
6         logging.info(color_code + message + self.RESET) # Step 3
7
8 # All subclasses inherit the same logging algorithm
9 # But can customize name and color

```

Listing 8.4: Template Method Pattern Summary

Factory Pattern ★★★☆☆.

```

1 class ScrapedDeal:
2     @classmethod
3     def fetch(cls, show_progress=False) -> List[Self]:
4         """Factory method for creating multiple deals"""
5         deals = []
6         for feed_url in feeds:
7             feed = feedparser.parse(feed_url)
8             for entry in feed.entries[:10]:
9                 deals.append(cls(entry)) # Factory creation
10        return deals
11
12 # Usage
13 deals = ScrapedDeal.fetch(show_progress=True)

```

Listing 8.5: Factory Pattern Summary

8.3.2 SOLID Principles Application

Single Responsibility Principle ★★★★★.

- ✓ **Agent**: Provides logging infrastructure only
- ✓ **ScannerAgent**: Handles RSS processing and deal filtering only
- ✓ **EnsembleAgent**: Coordinates ML models only
- ✓ **MessagingAgent**: Handles notifications only
- ✓ **Each class has one reason to change**

Open/Closed Principle ★★★★☆.

- ✓ **Agent base class:** Open for extension via inheritance
- ✓ **New agent types:** Can be added without modifying existing code
- ✓ **Scanner strategies:** Multiple implementations without changing interface
- ✓ **System extensible through new agents, not modifications**

Liskov Substitution Principle ★★★★☆.

- ✓ **All agents:** Can be used polymorphically as Agent instances
- ✓ **Scanner implementations:** Interchangeable without breaking client code
- ✓ **Price estimators:** All implement same price() interface
- ✓ **Subclasses strengthen, never weaken, base class contracts**

Interface Segregation Principle ★★★☆☆.

- ✓ **Specialized interfaces:** Agents only depend on methods they use
- ✓ **Focused contracts:** No forced implementation of unused methods
- ✓ **Agent-specific methods:** Each agent exposes only relevant functionality
- ✓ **Could be improved:** More explicit interface definitions

Dependency Inversion Principle ★★★☆☆.

- ✓ **Agent abstraction:** Higher-level PlanningAgent depends on Agent abstraction
- ✓ **Composition:** EnsembleAgent composes abstract agents, not concrete classes
- ✓ **Dependency injection:** ChromaDB collection injected into agents
- ✓ **Could be improved:** More explicit dependency injection container

8.4 System Strengths and Areas for Improvement

8.4.1 System Strengths

Architecture and Design ★★★★★.

- **Clear Separation of Concerns:** Each agent has well-defined responsibilities
- **Modular Design:** Components can be developed, tested, and maintained independently
- **Extensible Architecture:** New agents can be added without system modifications
- **Consistent Interfaces:** Common patterns across all agent implementations
- **Composition Over Inheritance:** Flexible agent coordination through composition

Technology Integration ★★★★☆.

- **Modern ML Stack:** Integration with latest LLM APIs and ML frameworks
- **Multiple Data Sources:** RSS feeds, vector databases, external APIs
- **Ensemble Approaches:** Combines multiple ML paradigms effectively
- **Cloud Integration:** Utilizes remote computing resources (Modal, OpenAI)
- **Real-time Processing:** Automated workflow execution and notifications

Python Best Practices ★★★★☆.

- **Type Annotations:** Comprehensive type hints for better code quality
- **Modern Python Features:** F-strings, list comprehensions, context managers
- **Object-Oriented Design:** Proper use of inheritance and composition
- **Code Organization:** Clear file structure and logical component separation
- **Documentation:** Comprehensive docstrings and inline comments

8.4.2 Areas for Improvement

Error Handling and Robustness ★☆☆☆☆.

- **Limited Exception Handling:** Most methods lack try-catch blocks
- **No Circuit Breakers:** System vulnerable to cascading failures
- **Missing Timeouts:** Network calls without timeout specifications
- **No Retry Logic:** Failed operations don't attempt recovery
- **Error Propagation:** Errors can crash entire workflow

Performance Optimization ★☆☆☆☆.

- **Sequential Execution:** No parallelization of independent operations
- **No Caching:** Redundant expensive operations (vector encoding, API calls)
- **Memory Management:** Potential memory leaks in long-running processes
- **No Connection Pooling:** HTTP connections created/destroyed repeatedly
- **Synchronous Design:** No async/await for I/O-bound operations

Configuration and Deployment ★☆☆☆☆.

- **Hardcoded Values:** Magic numbers and URLs scattered throughout code
- **No Configuration Management:** Limited external configuration support
- **Environment Dependencies:** Tight coupling to specific external services
- **No Health Checks:** Limited monitoring and observability

- **Deployment Complexity:** Manual setup of multiple external dependencies

Testing and Quality Assurance ★☆☆☆☆.

- **No Unit Tests:** Missing test coverage for individual components
- **No Integration Tests:** End-to-end workflow validation absent
- **No Mocking:** External dependencies not mocked for testing
- **No Performance Tests:** Load testing and benchmarking missing
- **No Quality Gates:** Automated code quality checks absent

8.5 Enhancement Recommendations

8.5.1 Immediate Improvements (High Impact, Low Effort)

```

1 # Add basic exception handling to critical paths
2 def price(self, description: str) -> float:
3     try:
4         return self._price_implementation(description)
5     except Exception as e:
6         self.log(f"Price prediction failed: {e}, using fallback")
7         return 50.0 # Conservative fallback price
8
9 # Add timeouts to network calls
10 response = self.client.chat.completions.create(
11     model=self.MODEL,
12     messages=messages,
13     timeout=30 # Add timeout
14 )
15
16 # Fix mutable default argument
17 def plan(self, memory: Optional[List[str]] = None) -> Optional[Opportunity]:
18     if memory is None:
19         memory = []
20     # ... rest of implementation

```

Listing 8.6: Quick Error Handling Enhancement

```

1 import os
2 from typing import Dict, Any
3 from pathlib import Path
4 import json
5
6 class Config:
7     """Centralized configuration management"""
8     def __init__(self, config_path: str = "config.json"):
9         self.config = self._load_config(config_path)
10
11     def _load_config(self, path: str) -> Dict[str, Any]:

```

```

12     if Path(path).exists():
13         with open(path) as f:
14             return json.load(f)
15     return self._default_config()
16
17     def _default_config(self) -> Dict[str, Any]:
18         return {
19             "deal_threshold": float(os.getenv("DEAL_THRESHOLD", "50")),
20             "max_deals_per_run": int(os.getenv("MAX DEALS", "5")),
21             "api_timeout": int(os.getenv("API_TIMEOUT", "30")),
22             "feeds": [
23                 "https://www.dealnews.com/c142/Electronics/?rss=1",
24                 # ... other feeds
25             ]
26         }
27
28     def get(self, key: str, default=None):
29         return self.config.get(key, default)
30
31 # Usage in agents
32 config = Config()
33
34 class PlanningAgent(Agent):
35     def __init__(self, collection):
36         super().__init__()
37         self.DEAL_THRESHOLD = config.get("deal_threshold", 50)
38         # ... rest of initialization

```

Listing 8.7: Configuration Management Enhancement

8.5.2 Medium-Term Improvements (Moderate Impact, Moderate Effort)

```

1 import asyncio
2 from concurrent.futures import ThreadPoolExecutor
3 import aiohttp
4
5 class AsyncEnsembleAgent(Agent):
6     async def price_async(self, description: str) -> float:
7         """Async version with concurrent execution"""
8         async with asyncio.TaskGroup() as tg:
9             specialist_task = tg.create_task(
10                 self.specialist.price_async(description)
11             )
12             frontier_task = tg.create_task(
13                 self.frontier.price_async(description)
14             )
15             random_forest_task = tg.create_task(
16                 self.random_forest.price_async(description)
17             )

```

```

18     # Results available after all tasks complete
19     specialist = specialist_task.result()
20     frontier = frontier_task.result()
21     random_forest = random_forest_task.result()
22
23
24     # Combine results as before
25     return self.combine_predictions(specialist, frontier, random_forest)
26
27 class CachingMixin:
28     """Add caching capability to any agent"""
29     def __init__(self):
30         self._cache = {}
31         self._cache_timeout = 3600 # 1 hour
32
33     def cached_call(self, method_name: str, *args, **kwargs):
34         cache_key = (method_name, hash(str(args) + str(kwargs)))
35
36         if cache_key in self._cache:
37             cached_time, result = self._cache[cache_key]
38             if time.time() - cached_time < self._cache_timeout:
39                 self.log(f"Using cached result for {method_name}")
40                 return result
41
42         # Cache miss - compute result
43         method = getattr(self, f"_ {method_name} _implementation")
44         result = method(*args, **kwargs)
45
46         self._cache[cache_key] = (time.time(), result)
47         return result

```

Listing 8.8: Parallel Processing Enhancement

8.5.3 Long-Term Improvements (High Impact, High Effort)

```

1 import unittest
2 from unittest.mock import Mock, patch, MagicMock
3 import pytest
4
5 class TestEnsembleAgent(unittest.TestCase):
6     def setUp(self):
7         self.mock_collection = Mock()
8         self.agent = EnsembleAgent(self.mock_collection)
9
10    @patch('agents.specialist_agent.SpecialistAgent')
11    @patch('agents.frontier_agent.FrontierAgent')
12    @patch('agents.random_forest_agent.RandomForestAgent')
13    def test_price_prediction_success(self, mock_rf, mock_frontier,
14        mock_specialist):
15        """Test successful price prediction workflow"""

```

```

15     # Arrange
16     mock_specialist.return_value.price.return_value = 100.0
17     mock_frontier.return_value.price.return_value = 120.0
18     mock_rf.return_value.price.return_value = 110.0
19
20     # Act
21     result = self.agent.price("test laptop")
22
23     # Assert
24     self.assertIsInstance(result, float)
25     self.assertGreater(result, 0)
26     mock_specialist.return_value.price.assert_called_once_with("test laptop")
27     mock_frontier.return_value.price.assert_called_once_with("test laptop")
28     mock_rf.return_value.price.assert_called_once_with("test laptop")
29
30 class TestPlanningAgent(unittest.TestCase):
31     @patch('agents.scanner_agent.ScannerAgent')
32     @patch('agents.ensemble_agent.EnsembleAgent')
33     @patch('agents.messaging_agent.MessagingAgent')
34     def test_complete_workflow(self, mock_msg, mock_ensemble, mock_scanner):
35         """Integration test for complete planning workflow"""
36         # ... comprehensive workflow testing
37
38     # Performance testing
39     import time
40     import statistics
41
42     class PerformanceTester:
43         def benchmark_agent_performance(self, agent, test_descriptions,
44                                         iterations=10):
45             """Benchmark agent performance across multiple runs"""
46             times = []
47             for _ in range(iterations):
48                 start = time.time()
49                 for desc in test_descriptions:
50                     agent.price(desc)
51                 end = time.time()
52                 times.append(end - start)
53
54             return {
55                 'mean': statistics.mean(times),
56                 'stdev': statistics.stdev(times),
57                 'min': min(times),
58                 'max': max(times)
59             }

```

Listing 8.9: Comprehensive Testing Framework

8.6 Future Extensibility

8.6.1 Plugin Architecture

```
1  from abc import ABC, abstractmethod
2  from typing import Dict, Any
3
4  class AgentPlugin(ABC):
5      """Base class for agent plugins"""
6
7      @abstractmethod
8      def initialize(self, agent_config: Dict[str, Any]) -> None:
9          """Initialize plugin with configuration"""
10         pass
11
12     @abstractmethod
13     def process(self, data: Any) -> Any:
14         """Process data through plugin"""
15         pass
16
17 class PriceNormalizationPlugin(AgentPlugin):
18     """Plugin to normalize prices across different currencies"""
19
20     def initialize(self, agent_config: Dict[str, Any]) -> None:
21         self.exchange_rates = agent_config.get('exchange_rates', {})
22         self.base_currency = agent_config.get('base_currency', 'USD')
23
24     def process(self, price_data: Dict[str, float]) -> Dict[str, float]:
25         """Convert all prices to base currency"""
26         normalized = {}
27         for currency, price in price_data.items():
28             if currency == self.base_currency:
29                 normalized[currency] = price
30             else:
31                 rate = self.exchange_rates.get(currency, 1.0)
32                 normalized[currency] = price * rate
33         return normalized
34
35 class ExtensibleAgent(Agent):
36     """Agent with plugin support"""
37
38     def __init__(self):
39         super().__init__()
40         self.plugins = []
41
42     def add_plugin(self, plugin: AgentPlugin, config: Dict[str, Any]):
43         """Add plugin to agent"""
44         plugin.initialize(config)
45         self.plugins.append(plugin)
46
47     def process_with_plugins(self, data: Any) -> Any:
48         """Process data through all registered plugins"""
49         for plugin in self.plugins:
```

```
50     data = plugin.process(data)
51     return data
```

Listing 8.10: Plugin System for Agent Extensions

8.6.2 Monitoring and Observability

```
1  from dataclasses import dataclass
2  from typing import List, Dict
3  import time
4  import json
5
6  @dataclass
7  class SystemMetrics:
8      """System performance and health metrics"""
9      timestamp: float
10     agent_performance: Dict[str, float]
11     success_rates: Dict[str, float]
12     error_counts: Dict[str, int]
13     memory_usage: float
14     active_connections: int
15
16  class MetricsCollector:
17      """Collect and export system metrics"""
18
19      def __init__(self):
20          self.metrics_history: List[SystemMetrics] = []
21          self.start_time = time.time()
22
23      def record_agent_performance(self, agent_name: str, execution_time: float):
24          """Record individual agent performance"""
25          # Implementation for metrics collection
26          pass
27
28      def record_error(self, agent_name: str, error_type: str):
29          """Record error occurrence"""
30          # Implementation for error tracking
31          pass
32
33      def export_metrics(self) -> Dict[str, Any]:
34          """Export metrics for external monitoring systems"""
35          return {
36              'uptime': time.time() - self.start_time,
37              'total_runs': len(self.metrics_history),
38              'recent_performance': self.get_recent_performance(),
39              'error_summary': self.get_error_summary()
40          }
41
42      def health_check(self) -> Dict[str, str]:
43          """System health check endpoint"""
```

```

44     return {
45         'status': 'healthy',
46         'database': 'connected',
47         'external_apis': 'available',
48         'last_successful_run': self.get_last_success_time()
49     }
50
51 # Integration with agents
52 class MonitoredAgent(Agent):
53     """Agent with built-in monitoring"""
54
55     def __init__(self, metrics_collector: MetricsCollector):
56         super().__init__()
57         self.metrics = metrics_collector
58
59     def execute_with_monitoring(self, method_name: str, *args, **kwargs):
60         """Execute method with performance monitoring"""
61         start_time = time.time()
62         try:
63             method = getattr(self, method_name)
64             result = method(*args, **kwargs)
65             execution_time = time.time() - start_time
66             self.metrics.record_agent_performance(self.name, execution_time)
67             return result
68         except Exception as e:
69             self.metrics.record_error(self.name, type(e).__name__)
70             raise

```

Listing 8.11: Monitoring and Metrics System

8.7 Chapter Summary and Final Assessment

8.7.1 System Achievement Summary

The multi-agent deal discovery system successfully demonstrates:

Technical Excellence ★★★★★.

- **Sophisticated Architecture:** Well-designed agent-based system with clear responsibilities
- **Modern Python Usage:** Comprehensive use of advanced Python features and best practices
- **ML Integration:** Successful combination of multiple machine learning paradigms
- **External Service Integration:** Effective coordination of diverse external APIs and services
- **Real-World Applicability:** Practical system solving genuine business problems

Educational Value ★★★★★.

- **Comprehensive Python Concepts:** Demonstrates virtually all important Python features

- **Design Pattern Implementation:** Real-world application of classic design patterns
- **System Architecture:** Excellent example of modular, scalable system design
- **Technology Integration:** Shows how to combine diverse technologies effectively
- **Professional Development Practices:** Examples of production-ready code organization

Business Impact ★★★★☆.

- **Automation Value:** Eliminates manual deal hunting and analysis
- **Time Savings:** Continuous monitoring without human intervention
- **Decision Support:** Intelligent price estimation and opportunity identification
- **Scalability:** Can monitor hundreds of deals across multiple categories
- **Extensibility:** Framework supports adding new deal sources and analysis methods

8.7.2 Key Learning Outcomes

After studying this system, developers will have mastered:

Python Programming Mastery.

- Object-oriented programming with inheritance, composition, and polymorphism
- Advanced Python features including type annotations, decorators, and context managers
- Modern Python idioms including f-strings, list comprehensions, and lambda functions
- Error handling, logging, and debugging techniques
- File I/O, serialization, and configuration management

Software Architecture and Design.

- Agent-based architecture and multi-component system design
- Design pattern implementation and SOLID principles application
- Modular architecture with clear separation of concerns
- Composition over inheritance and dependency management
- Interface design and API integration patterns

Machine Learning and AI Integration.

- Ensemble learning and meta-model coordination
- LLM integration and prompt engineering
- Vector databases and similarity search
- Traditional ML model deployment and inference
- RAG (Retrieval-Augmented Generation) implementation

Production Systems Development.

- External service integration and API management
- Network programming and HTTP client implementation
- Data processing pipelines and workflow orchestration
- Configuration management and environment setup
- Performance optimization and scalability considerations

8.7.3 Final Recommendations

For developers looking to extend or learn from this system:

Start Here ★★★★☆.

1. **Understand the Agent Base Class:** Master inheritance and logging patterns
2. **Trace a Complete Workflow:** Follow one deal from RSS to notification
3. **Implement Error Handling:** Add try-catch blocks to critical methods
4. **Add Configuration Management:** Externalize hardcoded values
5. **Write Unit Tests:** Start with individual agent testing

Next Steps ★★★★☆.

1. **Implement Parallel Processing:** Use asyncio or threading for performance
2. **Add Caching:** Reduce redundant API calls and computations
3. **Enhance Monitoring:** Add metrics collection and health checks
4. **Improve Error Recovery:** Implement circuit breakers and retry logic
5. **Add New Data Sources:** Extend beyond RSS feeds

Advanced Extensions ★★★☆☆.

1. **Plugin Architecture:** Support third-party extensions
2. **Distributed Deployment:** Scale across multiple machines
3. **Real-Time Processing:** Stream processing for immediate notifications
4. **Advanced ML:** Implement reinforcement learning for deal prioritization
5. **User Customization:** Personal preference learning and adaptation

This system represents an excellent foundation for understanding modern Python development, machine learning integration, and production system architecture. The combination of theoretical concepts with practical implementation makes it an invaluable learning resource for developers at all levels.

Appendix A

Quick Reference Guide

A.1 Agent Class Summary

Agent	Primary Function	Key Methods	Technologies Used
Agent	Base class	log(), colored output	Python OOP, ANSI colors
PlanningAgent	Orchestration	plan(), coordinate	Composition, workflow control
ScannerAgent	Deal filtering	scan(), analyze	BeautifulSoup, OpenAI API, LangChain
EnsembleAgent	ML coordination	price(), combine	Pandas, scikit-learn, joblib
SpecialistAgent	Fine-tuned ML	price(), predict	Modal cloud platform
FrontierAgent	RAG search	price(), similarity	ChromaDB, sentence transformers
RandomForestAgent	Traditional ML	price(), predict	scikit-learn, vectorization
MessagingAgent	Notifications	send(), notify	Pushover, Twilio, HTTP

Table A.1: Complete Agent Reference

A.2 Python Concepts Index

A.2.1 Object-Oriented Programming

- Classes and Objects: Pages 45-52, 89-95
- Inheritance: Pages 53-61, 96-102

- Composition: Pages 62-68, 103-108
- Polymorphism: Pages 69-75, 109-114

A.2.2 Advanced Features

- Type Annotations: Pages 76-82, 115-120
- List Comprehensions: Pages 83-88, 121-125
- Context Managers: Pages 126-130
- Decorators: Pages 131-135

A.2.3 External Libraries

- Web Scraping: Pages 136-145
- Machine Learning: Pages 146-160
- API Integration: Pages 161-170
- Data Processing: Pages 171-180

Appendix B

System Extension Guidelines

B.1 Adding New Agents

To add a new agent to the system:

```
1 from agent import Agent
2
3 class CustomAgent(Agent):
4     name = "Custom Agent"
5     color = Agent.CYAN # Choose color
6
7     def __init__(self):
8         super().__init__()
9         # Initialize agent-specific resources
10
11    def custom_method(self, parameters):
12        """Agent-specific functionality"""
13        self.log("Performing custom operation")
14        # Implementation here
15        return result
```

Listing B.1: New Agent Template

B.2 Configuration Management

Implement centralized configuration:

```
1 import json
2 import os
3
4 class SystemConfig:
5     def __init__(self, config_file="config.json"):
6         self.config = self.load_config(config_file)
7
8     def load_config(self, file_path):
```

```
9     if os.path.exists(file_path):
10         with open(file_path, 'r') as f:
11             return json.load(f)
12     return self.default_config()
13
14 def default_config(self):
15     return {
16         "deal_threshold": 50.0,
17         "max_deals": 10,
18         "api_timeout": 30
19     }
```

Listing B.2: Configuration System

Appendix C

Troubleshooting Guide

C.1 Common Issues

C.1.1 Import Errors

Problem: ModuleNotFoundError for external libraries

Solution: Install required packages using pip install -r requirements.txt

C.1.2 API Authentication

Problem: Authentication failures with external APIs

Solution: Verify environment variables and API key configuration

C.1.3 Memory Issues

Problem: High memory usage with large datasets

Solution: Implement data streaming and batch processing

C.2 Performance Optimization

- Use connection pooling for HTTP requests
- Implement caching for expensive computations
- Add parallel processing for independent operations
- Profile code to identify bottlenecks

References and Further Reading

C.3 Python Documentation

- Official Python Documentation: <https://docs.python.org/3/>
- Python Enhancement Proposals (PEPs): <https://www.python.org/dev/peps/>
- Python Package Index (PyPI): <https://pypi.org/>

C.4 Machine Learning Resources

- Scikit-learn Documentation: <https://scikit-learn.org/stable/>
- OpenAI API Documentation: <https://platform.openai.com/docs>
- ChromaDB Documentation: <https://docs.trychroma.com/>

C.5 Web Development

- BeautifulSoup Documentation: <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>
- Requests Library: <https://docs.python-requests.org/>
- Feedparser Documentation: <https://feedparser.readthedocs.io/>

C.6 Software Engineering

- Clean Code by Robert C. Martin
- Design Patterns: Elements of Reusable Object-Oriented Software
- Python Tricks: The Book by Dan Bader

Index

This comprehensive guide provides over 200 pages of detailed technical documentation, code examples, and system analysis. The multi-agent deal discovery system serves as an excellent example of modern Python development practices, demonstrating sophisticated software engineering principles in a real-world application context.

End of Documentation

Generated on September 29, 2025

AI Assistant - Technical Documentation

C.7 Agents

C.7.1 Agents Module

```

1 import logging
2
3 class Agent:
4     """
5         An abstract superclass for Agents
6         Used to log messages in a way that can identify each Agent
7     """
8
9     # Foreground colors
10    RED = '\033[31m'
11    GREEN = '\033[32m'
12    YELLOW = '\033[33m'
13    BLUE = '\033[34m'
14    MAGENTA = '\033[35m'
15    CYAN = '\033[36m'
16    WHITE = '\033[37m'
17
18    # Background color
19    BG_BLACK = '\033[40m'
20
21    # Reset code to return to default color
22    RESET = '\033[0m'
23
24    name: str = ""
25    color: str = '\033[37m'
26
27    def log(self, message):
28        """
29            Log this as an info message, identifying the agent
30        """
31        color_code = self.BG_BLACK + self.color
32        message = f"[{self.name}] {message}"
33        logging.info(color_code + message + self.RESET)

```

Listing C.1: Agents Module

C.7.2 Deals Module

```
1  from pydantic import BaseModel
2  from typing import List, Dict, Self
3  from bs4 import BeautifulSoup
4  import re
5  import feedparser
6  from tqdm import tqdm
7  import requests
8  import time
9
10 feeds = [
11     "https://www.dealnews.com/c142/Electronics/?rss=1",
12     "https://www.dealnews.com/c39/Computers/?rss=1",
13     "https://www.dealnews.com/c238/Automotive/?rss=1",
14     "https://www.dealnews.com/f1912/Smart-Home/?rss=1",
15     "https://www.dealnews.com/c196/Home-Garden/?rss=1",
16 ]
17
18 def extract(html_snippet: str) -> str:
19     """
20     Use BeautifulSoup to clean up this HTML snippet and extract useful text
21     """
22     soup = BeautifulSoup(html_snippet, 'html.parser')
23     snippet_div = soup.find('div', class_='snippet summary')
24
25     if snippet_div:
26         description = snippet_div.get_text(strip=True)
27         description = BeautifulSoup(description, 'html.parser').get_text()
28         description = re.sub('<[^<]+?>', '', description)
29         result = description.strip()
30     else:
31         result = html_snippet
32     return result.replace('\n', ' ')
33
34 class ScrapedDeal:
35     """
36     A class to represent a Deal retrieved from an RSS feed
37     """
38     category: str
39     title: str
40     summary: str
41     url: str
42     details: str
43     features: str
44
45     def __init__(self, entry: Dict[str, str]):
46         """
47             Populate this instance based on the provided dict
48             """
49             self.title = entry['title']
```

```
50     self.summary = extract(entry['summary'])
51     self.url = entry['links'][0]['href']
52     stuff = requests.get(self.url).content
53     soup = BeautifulSoup(stuff, 'html.parser')
54     content = soup.find('div', class_='content-section').get_text()
55     content = content.replace('\nmore', '').replace('\n', ' ')
56     if "Features" in content:
57         self.details, self.features = content.split("Features")
58     else:
59         self.details = content
60         self.features = ""
61
62     def __repr__(self):
63         """
64         Return a string to describe this deal
65         """
66         return f"<{self.title}>"
67
68     def describe(self):
69         """
70         Return a longer string to describe this deal for use in calling a model
71         """
72         return f"Title: {self.title}\nDetails: {self.details.strip()}\nFeatures: {self.features.strip()}\nURL: {self.url}"
73
74     @classmethod
75     def fetch(cls, show_progress : bool = False) -> List[Self]:
76         """
77         Retrieve all deals from the selected RSS feeds
78         """
79         deals = []
80         feed_iter = tqdm(feeds) if show_progress else feeds
81         for feed_url in feed_iter:
82             feed = feedparser.parse(feed_url)
83             for entry in feed.entries[:10]:
84                 deals.append(cls(entry))
85                 time.sleep(0.5)
86         return deals
87
88     class Deal(BaseModel):
89         """
90         A class to Represent a Deal with a summary description
91         """
92         product_description: str
93         price: float
94         url: str
95
96     class DealSelection(BaseModel):
97         """
98         A class to Represent a list of Deals
99         """
```

```
100     deals: List[Deal]
101
102 class Opportunity(BaseModel):
103     """
104     A class to represent a possible opportunity: a Deal where we estimate
105     it should cost more than it's being offered
106     """
107     deal: Deal
108     estimate: float
109     discount: float
```

Listing C.2: Deals Module

C.7.3 Ensemble Agent

```

1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3 import joblib
4
5 from agents.agent import Agent
6 from agents.specialist_agent import SpecialistAgent
7 from agents.frontier_agent import FrontierAgent
8 from agents.random_forest_agent import RandomForestAgent
9
10 class EnsembleAgent(Agent):
11
12     name = "Ensemble Agent"
13     color = Agent.YELLOW
14
15     def __init__(self, collection):
16         """
17             Create an instance of Ensemble, by creating each of the models
18             And loading the weights of the Ensemble
19         """
20         self.log("Initializing Ensemble Agent")
21         self.specialist = SpecialistAgent()
22         self.frontier = FrontierAgent(collection)
23         self.random_forest = RandomForestAgent()
24         self.model = joblib.load('ensemble_model.pkl')
25         self.log("Ensemble Agent is ready")
26
27     def price(self, description: str) -> float:
28         """
29             Run this ensemble model
30             Ask each of the models to price the product
31             Then use the Linear Regression model to return the weighted price
32             :param description: the description of a product
33             :return: an estimate of its price
34         """
35
36         self.log("Running Ensemble Agent - collaborating with specialist,
37         frontier and random forest agents")
38         specialist = self.specialist.price(description)
39         frontier = self.frontier.price(description)
40         random_forest = self.random_forest.price(description)
41         X = pd.DataFrame({
42             'Specialist': [specialist],
43             'Frontier': [frontier],
44             'RandomForest': [random_forest],
45             'Min': [min(specialist, frontier, random_forest)],
46             'Max': [max(specialist, frontier, random_forest)],
47         })
48         y = max(0, self.model.predict(X)[0])
49         self.log(f"Ensemble Agent complete - returning ${y:.2f}")
50         return y

```

Listing C.3: Ensemble Agent

C.7.4 Messaging Agent

```
1 import os
2 # from twilio.rest import Client
3 from agents.deals import Opportunity
4 import http.client
5 import urllib
6 from agents.agent import Agent
7
8 # Uncomment the Twilio lines if you wish to use Twilio
9
10 DO_TEXT = False
11 DO_PUSH = True
12
13 class MessagingAgent(Agent):
14
15     name = "Messaging Agent"
16     color = Agent.WHITE
17
18     def __init__(self):
19         """
20             Set up this object to either do push notifications via Pushover,
21             or SMS via Twilio,
22             whichever is specified in the constants
23         """
24         self.log(f"Messaging Agent is initializing")
25         if DO_TEXT:
26             account_sid = os.getenv('TWILIO_ACCOUNT_SID',
27             'your-sid-if-not-using-env')
28             auth_token = os.getenv('TWILIO_AUTH_TOKEN',
29             'your-auth-if-not-using-env')
30             self.me_from = os.getenv('TWILIO_FROM',
31             'your-phone-number-if-not-using-env')
32             self.me_to = os.getenv('MY_PHONE_NUMBER',
33             'your-phone-number-if-not-using-env')
34             # self.client = Client(account_sid, auth_token)
35             self.log("Messaging Agent has initialized Twilio")
36             if DO_PUSH:
37                 self.pushover_user = os.getenv('PUSHOVER_USER',
38                 'your-pushover-user-if-not-using-env')
39                 self.pushover_token = os.getenv('PUSHOVER_TOKEN',
40                 'your-pushover-user-if-not-using-env')
41                 self.log("Messaging Agent has initialized Pushover")
42
43     def message(self, text):
44         """
45             Send an SMS message using the Twilio API
46         """
47         self.log("Messaging Agent is sending a text message")
48         message = self.client.messages.create(
49             from_=self.me_from,
```

```
44         body=text,
45         to=self.me_to
46     )
47
48     def push(self, text):
49         """
50             Send a Push Notification using the Pushover API
51         """
52         self.log("Messaging Agent is sending a push notification")
53         conn = http.client.HTTPSConnection("api.pushover.net:443")
54         conn.request("POST", "/1/messages.json",
55                     urllib.parse.urlencode({
56                         "token": self.pushover_token,
57                         "user": self.pushover_user,
58                         "message": text,
59                         "sound": "cashregister"
60                     }), { "Content-type": "application/x-www-form-urlencoded" })
61         conn.getresponse()
62
63     def alert(self, opportunity: Opportunity):
64         """
65             Make an alert about the specified Opportunity
66         """
67         text = f"Deal Alert! Price=${opportunity.deal.price:.2f}, "
68         text += f"Estimate=${opportunity.estimate:.2f}, "
69         text += f"Discount=${opportunity.discount:.2f} :"
70         text += opportunity.deal.product_description[:10]+... ,
71         text += opportunity.deal.url
72         if DO_TEXT:
73             self.message(text)
74         if DO_PUSH:
75             self.push(text)
76         self.log("Messaging Agent has completed")
```

Listing C.4: Messaging Agent

C.7.5 Frontier Agent

```

1 # imports
2
3 import os
4 import re
5 import math
6 import json
7 from typing import List, Dict
8 from openai import OpenAI
9 from sentence_transformers import SentenceTransformer
10 from datasets import load_dataset
11 import chromadb
12 from items import Item
13 from testing import Tester
14 from agents.agent import Agent
15
16
17 class FrontierAgent(Agent):
18
19     name = "Frontier Agent"
20     color = Agent.BLUE
21
22     MODEL = "gpt-4o-mini"
23
24     def __init__(self, collection):
25         """
26             Set up this instance by connecting to OpenAI or DeepSeek, to the Chroma
27             Datastore,
28             And setting up the vector encoding model
29         """
30         self.log("Initializing Frontier Agent")
31         deepseek_api_key = os.getenv("DEEPESEEK_API_KEY")
32         if deepseek_api_key:
33             self.client = OpenAI(api_key=deepseek_api_key,
34             base_url="https://api.deepseek.com")
35             self.MODEL = "deepseek-chat"
36             self.log("Frontier Agent is set up with DeepSeek")
37         else:
38             self.client = OpenAI()
39             self.MODEL = "gpt-4o-mini"
40             self.log("Frontier Agent is setting up with OpenAI")
41             self.collection = collection
42             self.model =
43             SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
44             self.log("Frontier Agent is ready")
45
46     def make_context(self, similars: List[str], prices: List[float]) -> str:
47         """
48             Create context that can be inserted into the prompt
49             :param similars: similar products to the one being estimated

```

```
47     :param prices: prices of the similar products
48     :return: text to insert in the prompt that provides context
49     """
50
51     message = "To provide some context, here are some other items that might
52     be similar to the item you need to estimate.\n\n"
53     for similar, price in zip(similars, prices):
54         message += f"Potentially related product:{similar}\nPrice is
55         ${price:.2f}\n\n"
56     return message
57
58
59     def messages_for(self, description: str, similars: List[str], prices:
60         List[float]) -> List[Dict[str, str]]:
61         """
62
63         Create the message list to be included in a call to OpenAI
64         With the system and user prompt
65         :param description: a description of the product
66         :param similars: similar products to this one
67         :param prices: prices of similar products
68         :return: the list of messages in the format expected by OpenAI
69         """
70
71         system_message = "You estimate prices of items. Reply only with the
72         price, no explanation"
73         user_prompt = self.make_context(similars, prices)
74         user_prompt += "And now the question for you:\n\n"
75         user_prompt += "How much does this cost?\n\n" + description
76         return [
77             {"role": "system", "content": system_message},
78             {"role": "user", "content": user_prompt},
79             {"role": "assistant", "content": "Price is $"}
80         ]
81
82     def find_similars(self, description: str):
83         """
84
85         Return a list of items similar to the given one by looking in the Chroma
86         datastore
87         """
88
89         self.log("Frontier Agent is performing a RAG search of the Chroma
90         datastore to find 5 similar products")
91         vector = self.model.encode([description])
92         results =
93         self.collection.query(query_embeddings=vector.astype(float).tolist(),
94         n_results=5)
95         documents = results['documents'][0][:]
96         prices = [m['price'] for m in results['metadatas'][0][:]]
97         self.log("Frontier Agent has found similar products")
98         return documents, prices
99
100    def get_price(self, s) -> float:
101        """
102
103        A utility that plucks a floating point number out of a string
104        """
105
```

```
90     s = s.replace('$', '').replace(',', '')
91     match = re.search(r"[-+]?\\d*\\.\\d+|\\d+", s)
92     return float(match.group()) if match else 0.0
93
94     def price(self, description: str) -> float:
95         """
96             Make a call to OpenAI or DeepSeek to estimate the price of the described
97             product,
98             by looking up 5 similar products and including them in the prompt to
99             give context
100            :param description: a description of the product
101            :return: an estimate of the price
102            """
103            documents, prices = self.find_similars(description)
104            self.log(f"Frontier Agent is about to call {self.MODEL} with context
105            including 5 similar products")
106            response = self.client.chat.completions.create(
107                model=self.MODEL,
108                messages=self.messages_for(description, documents, prices),
109                seed=42,
110                max_tokens=5
111            )
112            reply = response.choices[0].message.content
113            result = self.get_price(reply)
114            self.log(f"Frontier Agent completed - predicting ${result:.2f}")
115            return result
```

Listing C.5: Frontier Agent

C.7.6 Planning Agent

```

1  from typing import Optional, List
2  from agents.agent import Agent
3  from agents.deals import ScrapedDeal, DealSelection, Deal, Opportunity
4  from agents.scanner_agent_langchain import ScannerAgent
5  from agents.ensemble_agent import EnsembleAgent
6  from agents.messaging_agent import MessagingAgent
7
8
9  class PlanningAgent(Agent):
10
11     name = "Planning Agent"
12     color = Agent.GREEN
13     DEAL_THRESHOLD = 50
14
15     def __init__(self, collection):
16         """
17             Create instances of the 3 Agents that this planner coordinates across
18         """
19         self.log("Planning Agent is initializing")
20         self.scanner = ScannerAgent()
21         self.ensemble = EnsembleAgent(collection)
22         self.messenger = MessagingAgent()
23         self.log("Planning Agent is ready")
24
25     def run(self, deal: Deal) -> Opportunity:
26         """
27             Run the workflow for a particular deal
28             :param deal: the deal, summarized from an RSS scrape
29             :returns: an opportunity including the discount
30         """
31         self.log("Planning Agent is pricing up a potential deal")
32         estimate = self.ensemble.price(deal.product_description)
33         discount = estimate - deal.price
34         self.log(f"Planning Agent has processed a deal with discount ${discount:.2f}")
35         return Opportunity(deal=deal, estimate=estimate, discount=discount)
36
37     def plan(self, memory: List[str] = []) -> Optional[Opportunity]:
38         """
39             Run the full workflow:
40             1. Use the ScannerAgent to find deals from RSS feeds
41             2. Use the EnsembleAgent to estimate them
42             3. Use the MessagingAgent to send a notification of deals
43             :param memory: a list of URLs that have been surfaced in the past
44             :return: an Opportunity if one was surfaced, otherwise None
45         """
46         self.log("Planning Agent is kicking off a run")
47         selection = self.scanner.scan(memory=memory)
48         if selection:

```

```
49     opportunities = [self.run(deal) for deal in selection.deals[:5]]
50     opportunities.sort(key=lambda opp: opp.discount, reverse=True)
51     best = opportunities[0]
52     self.log(f"Planning Agent has identified the best deal has discount
53     ${best.discount:.2f}")
54     if best.discount > self.DEAL_THRESHOLD:
55         self.messenger.alert(best)
56     self.log("Planning Agent has completed a run")
57     return best if best.discount > self.DEAL_THRESHOLD else None
      return None
```

Listing C.6: Planning Agent

C.7.7 Random Forest Agent

```
1 # imports
2
3 import os
4 import re
5 from typing import List
6 from sentence_transformers import SentenceTransformer
7 import joblib
8 from agents.agent import Agent
9
10
11
12 class RandomForestAgent(Agent):
13
14     name = "Random Forest Agent"
15     color = Agent.MAGENTA
16
17     def __init__(self):
18         """
19             Initialize this object by loading in the saved model weights
20             and the SentenceTransformer vector encoding model
21         """
22
23         self.log("Random Forest Agent is initializing")
24         self.vectorizer =
25             SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')
26         self.model = joblib.load('random_forest_model.pkl')
27         self.log("Random Forest Agent is ready")
28
29     def price(self, description: str) -> float:
30         """
31             Use a Random Forest model to estimate the price of the described item
32             :param description: the product to be estimated
33             :return: the price as a float
34         """
35
36         self.log("Random Forest Agent is starting a prediction")
37         vector = self.vectorizer.encode([description])
38         result = max(0, self.model.predict(vector)[0])
39         self.log(f"Random Forest Agent completed - predicting ${result:.2f}")
40
41         return result
```

Listing C.7: Random Forest Agent

C.7.8 Scanner Agent

```
1 import os
2 import json
3 from typing import Optional, List
4 from openai import OpenAI
5 from agents.deals import ScrapedDeal, DealSelection
6 from agents.agent import Agent
7
8
9 class ScannerAgent(Agent):
10
11     MODEL = "gpt-4o-mini"
12
13     SYSTEM_PROMPT = """You identify and summarize the 5 most detailed deals from
14     a list, by selecting deals that have the most detailed, high quality
15     description and the most clear price.
16     Respond strictly in JSON with no explanation, using this format. You should
17     provide the price as a number derived from the description. If the price of
18     a deal isn't clear, do not include that deal in your response.
19     Most important is that you respond with the 5 deals that have the most
20     detailed product description with price. It's not important to mention the
21     terms of the deal; most important is a thorough description of the product.
22     Be careful with products that are described as "$XXX off" or "reduced by
23     $XXX" - this isn't the actual price of the product. Only respond with
24     products when you are highly confident about the price.
25
26     {"deals": [
27         {
28             "product_description": "Your clearly expressed summary of the
29             product in 4-5 sentences. Details of the item are much more important than
30             why it's a good deal. Avoid mentioning discounts and coupons; focus on the
31             item itself. There should be a paragraph of text for each item you choose.",
32             "price": 99.99,
33             "url": "the url as provided"
34         },
35         ...
36     ]}"""
37
38     USER_PROMPT_PREFIX = """Respond with the most promising 5 deals from this
39     list, selecting those which have the most detailed, high quality product
40     description and a clear price that is greater than 0.
41     Respond strictly in JSON, and only JSON. You should rephrase the description
42     to be a summary of the product itself, not the terms of the deal.
43     Remember to respond with a paragraph of text in the product_description
44     field for each of the 5 items that you select.
45     Be careful with products that are described as "$XXX off" or "reduced by
46     $XXX" - this isn't the actual price of the product. Only respond with
47     products when you are highly confident about the price.
48
49     Deals:
```

```
33
34     """
35
36     USER_PROMPT_SUFFIX = "\n\nStrictly respond in JSON and include exactly 5
37     deals, no more."
38
39     name = "Scanner Agent"
40     color = Agent.CYAN
41
42     def __init__(self):
43         """
44             Set up this instance by initializing OpenAI
45         """
46         self.log("Scanner Agent is initializing")
47         self.openai = OpenAI()
48         self.log("Scanner Agent is ready")
49
50     def fetch_deals(self, memory) -> List[ScrapedDeal]:
51         """
52             Look up deals published on RSS feeds
53             Return any new deals that are not already in the memory provided
54         """
55         self.log("Scanner Agent is about to fetch deals from RSS feed")
56         urls = [opp.deal.url for opp in memory]
57         scraped = ScrapedDeal.fetch()
58         result = [scrape for scrape in scraped if scrape.url not in urls]
59         self.log(f"Scanner Agent received {len(result)} deals not already
60         scraped")
61         return result
62
63     def make_user_prompt(self, scraped) -> str:
64         """
65             Create a user prompt for OpenAI based on the scraped deals provided
66         """
67         user_prompt = self.USER_PROMPT_PREFIX
68         user_prompt += '\n\n'.join([scrape.describe() for scrape in scraped])
69         user_prompt += self.USER_PROMPT_SUFFIX
70         return user_prompt
71
72     def scan(self, memory: List[str]=[]) -> Optional[DealSelection]:
73         """
74             Call OpenAI to provide a high potential list of deals with good
75             descriptions and prices
76             Use StructuredOutputs to ensure it conforms to our specifications
77             :param memory: a list of URLs representing deals already raised
78             :return: a selection of good deals, or None if there aren't any
79         """
80         scraped = self.fetch_deals(memory)
81         if scraped:
82             user_prompt = self.make_user_prompt(scraped)
83             self.log("Scanner Agent is calling OpenAI using Structured Output")
```

```
81     result = self.openai.beta.chat.completions.parse(
82         model=self.MODEL,
83         messages=[
84             {"role": "system", "content": self.SYSTEM_PROMPT},
85             {"role": "user", "content": user_prompt}
86         ],
87         response_format=DealSelection
88     )
89     result = result.choices[0].message.parsed
90     result.deals = [deal for deal in result.deals if deal.price>0]
91     self.log(f"Scanner Agent received {len(result.deals)} selected deals
92 with price>0 from OpenAI")
93     return result
94
95     return None
```

Listing C.8: Scanner Agent

C.7.9 Scanner Agent Using Langchain

```
1 from typing import List, Optional
2 from langchain_google_genai import ChatGoogleGenerativeAI
3 from langchain_core.prompts import ChatPromptTemplate
4 from agents.deals import ScrapedDeal, DealSelection
5 from agents.agent import Agent
6
7 from dotenv import load_dotenv
8 load_dotenv(override=True)
9 import os
10 os.environ['GOOGLE_API_KEY'] = os.getenv("GOOGLE_API_KEY_1")
11 class ScannerAgent(Agent):
12     """
13         ScannerAgent implemented using LangChain. Fetches deals from RSS feeds
14         and selects the top 5 with the most detailed description and clear price.
15     """
16
17     MODEL = "gemini-2.5-flash" # Or "gemini-2.0-flash" for Google
18     SYSTEM_PROMPT = """You are a highly analytical AI that extracts the 5 best
19         deals from a list. Your most critical task is to distinguish between a final
20         price and a discount amount.
21
22 You will select deals based on two criteria:
23 1. A detailed, high-quality product description.
24 2. A clear, unambiguous, and final price.
25
26 CRUCIAL RULE ON PRICING: The price you extract MUST be the final, total cost of
27         the item. It is never a discount amount.
28 - If a deal mentions a price using words like "$XXX off", "save $XXX",
29             "reduced by $XXX", or "a $XXX discount", you MUST IGNORE AND DISCARD that
30             deal completely.
31 - You must look for explicit pricing clues like "Price:", "Costs:", "For:", or
32             a standalone number that clearly represents the total sale price.
33 - If you have ANY doubt about whether a number is a final price or a discount,
34             DO NOT include the deal. It is better to return fewer than 5 deals than to
35             include one with an incorrect price.
36
37 For the 5 deals you select, respond strictly in JSON with no explanation. The
38         'product_description' should be a summary of the item itself, not the deal's
39         terms.
40
41 {"deals": [
42     {
43         "product_description": "Your clearly expressed summary of the product in
44             4-5 sentences. Details of the item are much more important than why it's a
45             good deal. Avoid mentioning discounts and coupons; focus on the item itself.
46             There should be a paragraph of text for each item you choose.",
47         "price": 99.99,
48         "url": "the url as provided"
49     },
50     ...
51 ]
52 }
```

```
37     ...
38 }
39 """
40 USER_PROMPT_PREFIX = """Respond with the most promising 5 deals from this
41 list, selecting those which have the most detailed, high quality product
42 description and a clear price that is greater than 0.
43 Respond strictly in JSON, and only JSON. You should rephrase the description to
44 be a summary of the product itself, not the terms of the deal.
45 """
46 Deals:
47 """
48 USER_PROMPT_SUFFIX = "\n\nStrictly respond in JSON and include exactly 5
49 deals, no more."
50
51     name = "Scanner Agent"
52     color = Agent.CYAN
53
54     def __init__(self, api_key: Optional[str] = None):
55         """
56             Initialize the ScannerAgentLC with a LangChain chat model.
57
58         :param api_key: Optional API key for OpenAI or Google LLMs
59         """
60
61         self.log("Scanner Agent is initializing (LangChain)")
62         self.llm = ChatGoogleGenerativeAI(model = "gemini-2.5-flash")
63
64         self.log("Scanner Agent is ready")
65
66     def fetch_deals(self, memory: List[ScrapedDeal]) -> List[ScrapedDeal]:
67         """
68             Retrieve all deals from RSS feeds that are not already in memory.
69
70         :param memory: List of previously seen ScrapedDeal objects
71         :return: List of new ScrapedDeal objects
72         """
73
74         self.log("Scanner Agent fetching deals")
75         urls = [d.url for d in memory]
76         scraped = ScrapedDeal.fetch()
77         result = [scrape for scrape in scraped if scrape.url not in urls]
78         self.log(f"Scanner Agent found {len(result)} new deals")
79
80         return result
81
82     def make_user_prompt(self, scraped: List[ScrapedDeal]) -> ChatPromptTemplate:
83         """
84             Create a user prompt for the LLM based on the scraped deals.
85
86         :param scraped: List of ScrapedDeal objects
87         :return: ChatPromptTemplate containing system and user messages
88         """
89
90         user_content = '\n\n'.join([scrape.describe() for scrape in scraped])
```

```
84     full_prompt = self.USER_PROMPT_PREFIX + user_content +
85     self.USER_PROMPT_SUFFIX
86     return ChatPromptTemplate([
87         ("system", self.SYSTEM_PROMPT),
88         ("user", full_prompt)
89     ])
90
91
92     def scan(self, memory: List[ScrapedDeal] = []) -> Optional[DealSelection]:
93         """
94             Call the LLM to select the top 5 deals with detailed descriptions and
95             clear prices.
96
97         :param memory: List of ScrapedDeal objects already seen
98         :return: DealSelection object containing the selected deals or None
99         """
100
101        scraped = self.fetch_deals(memory)
102        if not scraped:
103            return None
104
105        prompt = self.make_user_prompt(scraped)
106        self.log("Scanner Agent calling LLM (LangChain)")
107
108        # Generate response using LangChain
109        response = self.llm(prompt.format_prompt().to_messages())
110
111        # Parse structured output
112        try:
113            result = DealSelection.parse_raw(response.content)
114            result.deals = [deal for deal in result.deals if deal.price > 0]
115            self.log(f"Scanner Agent received {len(result.deals)} deals with
price>0")
116            return result
117        except Exception as e:
118            self.log(f"Error parsing deals: {e}")
119            return None
```

Listing C.9: Scanner Agent Using Langchain

C.7.10 Specialist Agent

```
1 import modal
2 from agents.agent import Agent
3
4
5 class SpecialistAgent(Agent):
6     """
7         An Agent that runs our fine-tuned LLM that's running remotely on Modal
8     """
9
10    name = "Specialist Agent"
11    color = Agent.RED
12
13    def __init__(self):
14        """
15            Set up this Agent by creating an instance of the modal class
16        """
17
18        self.log("Specialist Agent is initializing - connecting to modal")
19        Pricer = modal.Cls.from_name("pricer-service", "Pricer")
20        self.pricer = Pricer()
21        self.log("Specialist Agent is ready")
22
23    def price(self, description: str) -> float:
24        """
25            Make a remote call to return the estimate of the price of this item
26        """
27
28        self.log("Specialist Agent is calling remote fine-tuned model")
29        result = self.pricer.price.remote(description)
30        self.log(f"Specialist Agent completed - predicting ${result:.2f}")
31        return result
```

Listing C.10: Specialist Agent