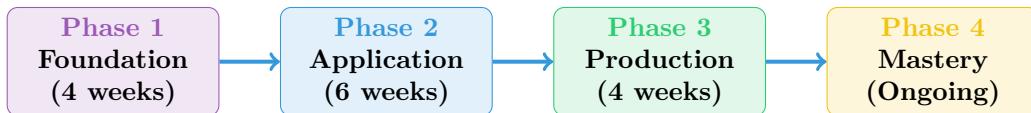# COMPREHENSIVE PROGRAM ROADMAP

# Prompt Engineering & LLM Application Development

A Structured Learning Journey from
Foundation to Expert-Level Mastery

*Technology Agnostic • Provider Independent • Industry Aligned*

| **Phase 1** Foundation (4 weeks) | → | **Phase 2** Application (6 weeks) | → | **Phase 3** Production (4 weeks) | → | **Phase 4** Mastery (Ongoing) |
|---|---|---|---|---|---|---|

| **14–18 Weeks** Core Program | **200+ Hours** Learning Time | **4 Capstones** Major Projects | **50+ Skills** Competencies |
|---|---|---|---|

# Contents

# Part I

# Program Overview and Navigation

# Chapter 1

# Executive Roadmap Summary

## 1.1 Program Vision

This comprehensive roadmap provides a structured, technology-agnostic pathway for mastering prompt engineering and building production-ready LLM-powered applications. The program is designed around the principle that effective LLM engineering requires a combination of foundational prompt crafting skills, software engineering best practices, and operational excellence.

The roadmap emphasizes practical application over theoretical knowledge, with each phase building upon the previous one through hands-on projects, iterative skill development, and progressive complexity. Learners emerge not just as prompt engineers, but as complete AI application developers capable of designing, building, deploying, and maintaining sophisticated LLM systems.

## 1.2 Program Structure Overview

| Phase 1<br>**Foundation**<br>4 weeks | Phase 2<br>**Application**<br>6 weeks | Phase 3<br>**Production**<br>4 weeks | Phase 4<br>**Mastery**<br>Ongoing |
|---|---|---|---|
| Prompting Fundamentals | Architecture Patterns | System Design | Advanced Techniques |
| Advanced Patterns | Tool Integration | Evaluation at Scale | Research Frontiers |
| Evaluation Methods | RAG Systems | LLMOps | Specialization |
| Specialized Domains | Conversational AI | Safety & Security | Leadership |

## 1.3 Learning Outcomes by Phase

### 1.3.1 Phase 1: Foundation (Weeks 1–4)

Upon completing Phase 1, learners will demonstrate mastery of core prompting principles applicable across any LLM platform:

- **Prompt Construction:** Design clear, effective prompts that consistently produce high-quality outputs regardless of the underlying model

- **Pattern Recognition:** Identify and apply 25+ distinct prompting patterns to appropriate use cases

- **Reasoning Enhancement:** Implement chain-of-thought, self-consistency, and other reasoning amplification techniques

- **Quality Evaluation:** Systematically assess and iterate on prompt performance using both qualitative and quantitative methods

- **Multi-Modal Competence:** Extend prompting skills to code generation, structured outputs, and vision-language tasks

### 1.3.2 Phase 2: Application Development (Weeks 5–10)

Phase 2 transforms prompt engineers into application developers:

- **Architecture Design:** Select and implement appropriate architectural patterns for LLM applications

- **Tool Integration:** Enable LLMs to interact with external systems through function calling and tool use

- **Knowledge Augmentation:** Build retrieval-augmented generation (RAG) systems for knowledge-intensive applications

- **Conversational Systems:** Design stateful agents with memory, context management, and persona consistency

- **Provider Abstraction:** Create applications that work across multiple LLM providers without code changes

### 1.3.3 Phase 3: Production Engineering (Weeks 11–14)

Phase 3 focuses on operational excellence and production readiness:

- **System Design:** Make informed trade-offs between prompting, fine-tuning, and RAG approaches

- **Quality Assurance:** Implement comprehensive evaluation pipelines with automated and human-in-the-loop assessment

- **Observability:** Deploy monitoring, logging, and alerting systems for LLM applications

- **Cost Optimization:** Manage inference costs through caching, model selection, and prompt optimization

- **Security Hardening:** Protect against prompt injection, data leakage, and other LLM-specific vulnerabilities

### 1.3.4 Phase 4: Mastery (Ongoing)

Phase 4 represents continuous professional growth:

- **Advanced Techniques:** Master emerging patterns like constitutional AI, recursive self-improvement, and multi-agent orchestration

- **Research Integration:** Stay current with academic developments and integrate new techniques into practice

- **Domain Specialization:** Develop deep expertise in specific verticals (healthcare, finance, legal, etc.)

- **Technical Leadership:** Guide teams and organizations in LLM adoption and best practices

## 1.4   Target Audience Profiles

This roadmap accommodates multiple learner profiles with tailored pathways:

| Profile | Background & Goals | Recommended Track |
|---|---|---|
| **Software Engineer** | Experienced developer seeking to integrate LLM capabilities into existing applications | Standard (14 weeks) |
| **Data Scientist** | ML practitioner transitioning to LLM-focused roles; familiar with model evaluation | Accelerated (10 weeks) |
| **Product Manager** | Technical PM requiring deep understanding of LLM capabilities and limitations | Foundation+ (8 weeks) |
| **Technical Writer** | Content creator using generative AI for documentation and communication | Custom (Phase 1 + select modules) |
| **Career Changer** | Professional from another field entering AI/ML; strong analytical skills | Extended (18+ weeks) |

# Chapter 2

# Prerequisites and Preparation

## 2.1 Required Knowledge Assessment

Before beginning the program, learners should honestly assess their current competencies against the following requirements. This self-assessment determines your recommended starting point and pace.

> ⚠️ **Prerequisites**
>
> **Essential Prerequisites (Must Have):**
>
> - Proficiency in at least one programming language (Python strongly recommended)
>
> - Understanding of basic software development concepts (version control, APIs, data structures)
>
> - Familiarity with command-line interfaces and development environments
>
> - Access to API credits from at least one major LLM provider
>
> - Dedicated study time of 8–15 hours per week

### 2.1.1 Technical Skills Matrix

Rate your current proficiency (1 = Beginner, 5 = Expert) for each skill:

| Skill Area | 1 | 2 | 3 | 4 | 5 | Required |
|---|---|---|---|---|---|---|
| Python programming | ☐ | ☐ | ☐ | ☐ | ☐ | 3+ |
| REST API consumption | ☐ | ☐ | ☐ | ☐ | ☐ | 2+ |
| JSON/YAML data formats | ☐ | ☐ | ☐ | ☐ | ☐ | 3+ |
| Git version control | ☐ | ☐ | ☐ | ☐ | ☐ | 2+ |
| Database fundamentals (SQL) | ☐ | ☐ | ☐ | ☐ | ☐ | 2+ |
| Web development basics | ☐ | ☐ | ☐ | ☐ | ☐ | 1+ |
| Machine learning concepts | ☐ | ☐ | ☐ | ☐ | ☐ | 1+ |
| Cloud services familiarity | ☐ | ☐ | ☐ | ☐ | ☐ | 1+ |

### 2.1.2 Prerequisite Gap Remediation

If you identified gaps in the essential prerequisites, the following resources can help you prepare:

| Gap Area | Recommended Resource | Time Investment |
|---|---|---|
| Python Programming | Python official tutorial; Codecademy Python course | 20–40 hours |
| REST APIs | RESTful API tutorials; Postman learning center | 5–10 hours |
| Git/Version Control | Pro Git book (free online); GitHub Learning Lab | 5–10 hours |
| Command Line | Linux Command Line Basics; terminal tutorials | 5–8 hours |
| JSON/Data Formats | JSON/YAML tutorials; practice with real APIs | 2–4 hours |
| ML Concepts | 3Blue1Brown neural networks; fast.ai intro | 10–15 hours |

## 2.2   Environment Setup Checklist

Complete this checklist before beginning Week 1:

---

### ⋮≣ Self-Evaluation Checklist

**Development Environment:**

- ☐ Python 3.9+ installed with pip/conda
- ☐ Code editor or IDE configured (VS Code, PyCharm, etc.)
- ☐ Git installed and GitHub account created
- ☐ Virtual environment manager available (venv, conda, poetry)
- ☐ Jupyter Notebook or JupyterLab installed

**API Access:**

- ☐ Primary LLM provider account and API key (any major provider)
- ☐ Secondary provider account (recommended for comparison)
- ☐ API credit budget allocated ($50–100 recommended for full program)
- ☐ Environment variables configured for secure key storage

**Learning Infrastructure:**

- ☐ Note-taking system established (Notion, Obsidian, etc.)
- ☐ Calendar blocked for study sessions
- ☐ Community access (Discord, Slack, or study group)
- ☐ Core textbooks acquired (see Resources chapter)

---

## 2.3   Mindset Preparation

Success in this program requires embracing several key mindsets:

### 2.3.1   Empirical Iteration

LLM behavior is inherently probabilistic and sometimes surprising. Effective prompt engineers adopt an experimental mindset, systematically testing hypotheses rather than assuming outcomes. Every prompt is a hypothesis to be validated, not a certainty.

### 2.3.2   Provider Agnosticism

While you may develop on one platform, the principles taught here transfer across all major providers. Avoid over-fitting your mental models to any single provider's quirks or features. The goal is portable expertise.

### 2.3.3   Continuous Learning

The LLM landscape evolves rapidly. Techniques that are state-of-the-art today may be superseded within months. This program provides foundational skills and frameworks for ongoing adaptation, not a static body of knowledge.

### 2.3.4   Production Orientation

From Day 1, think about how techniques would work in production environments. Consider edge cases, failure modes, costs, and user experience even during learning exercises.

# Part II

# Detailed Learning Path

# Chapter 3

# Phase 1: Foundation Mastery

> **⏱ Time Investment**
>
> **Duration:** 4 weeks
> **Weekly Commitment:** 10–12 hours
> **Total Hours:** 40–48 hours
> **Delivery:** Personal Prompt Cookbook

## 3.1 Phase 1 Learning Path Visualization

| Week 1 Fundamentals | Week 2 Patterns | Week 3 Evaluation | Week 4 Specialization |
| --- | --- | --- | --- |
| LLM Mechanics | Chain-of-Thought | Metrics Design | Code Generation |
| Prompt Anatomy | Few-Shot Learning | Test Set Creation | Structured Output |
| Core Techniques | Self-Consistency | A/B Testing | Multi-Modal |

Checkpoint: First 10 Patterns

Capstone: Prompt Cookbook

## 3.2 Week 1: Prompt Engineering Fundamentals

### 3.2.1 Learning Objectives

By the end of Week 1, you will:

- Understand how language models process and generate text at a conceptual level

- Identify the components of effective prompts and their purposes

- Apply fundamental prompting techniques with confidence

- Configure generation parameters appropriately for different tasks

- Recognize the capabilities and limitations of current LLM technology

### 3.2.2 Day-by-Day Breakdown

**Days 1–2: Understanding Language Models**

**Core Concepts:**

- Tokenization: How text becomes numbers; why token boundaries matter

- Context windows: Working within limits; strategies for context management

- Probabilistic generation: Outputs as samples from distributions

- Instruction tuning: How models learn to follow directions (RLHF, Constitutional AI concepts)

- Model families: Understanding capability tiers without vendor lock-in

**Practical Activities:**

1. Experiment with a tokenizer to understand token boundaries

2. Test identical prompts across different model sizes, noting capability differences

3. Explore generation parameters (temperature, top-p) with the same prompt

**Time Investment:** 4–5 hours

**Days 3–4: Anatomy of Effective Prompts**

**Core Concepts:**

- Prompt structure: Role, context, instruction, format, constraints

- Clarity techniques: Specificity, explicit requirements, unambiguous language

- Output formatting: JSON, XML, Markdown, custom structures

- Positive and negative constraints: What to include versus avoid

- Delimiters and structure: Organizing complex prompts

**Practical Activities:**

1. Transform 5 vague prompts into precise, well-structured versions

2. Create a "prompt anatomy" template for future use

3. Build a prompt that produces valid JSON output reliably

**Time Investment:** 4–5 hours

**Days 5–7: Core Prompting Techniques**

**Techniques to Master:**

| Technique | Implementation Focus |
| --- | --- |
| Role Prompting | Assign personas to influence response style, depth, and perspective |
| Context Setting | Provide relevant background to ground the model's responses |
| Instruction Clarity | Write unambiguous, specific, actionable instructions |
| Constraint Specification | Define boundaries: length, scope, format, style |
| Output Formatting | Specify exact output structure for reliable parsing |
| Task Decomposition | Break complex requests into sequential steps |

**Practical Activities:**

1. Apply each technique to 3 different task types

2. Compare outputs with and without each technique

3. Document which techniques work best for which scenarios

**Time Investment:** 3–4 hours

**⋮≡ Self-Evaluation Checklist**

**Week 1 Self-Assessment:**

☐ I can explain why the same prompt might produce different outputs

☐ I can identify at least 5 components of a well-structured prompt

☐ I can transform a vague request into a precise prompt

☐ I can reliably produce structured output (JSON) from prompts

☐ I understand when to use high vs. low temperature settings

☐ I have documented at least 6 fundamental prompting techniques

## 3.3   Week 2: Advanced Prompting Patterns

### 3.3.1   Learning Objectives

By the end of Week 2, you will:

- Implement chain-of-thought prompting and its variants effectively

- Design few-shot examples that maximize model performance

- Apply self-consistency techniques for high-stakes decisions

- Use critique-and-revise patterns to improve output quality

- Understand when to apply each advanced pattern

### 3.3.2  Pattern Deep Dives

**Chain-of-Thought (CoT) Prompting**

Chain-of-thought prompting enables complex reasoning by having the model generate intermediate steps before reaching conclusions.

| CoT Variant | Trigger Phrase | Best Use Case |
|---|---|---|
| Zero-Shot CoT | "Let's think step by step" | Quick reasoning enhancement without examples |
| Few-Shot CoT | Provide reasoning examples | Complex problems requiring specific reasoning style |
| Self-Consistency | Multiple CoT paths | High-stakes decisions requiring confidence |
| Tree-of-Thought | Branching exploration | Problems with multiple valid approaches |
| Program-Aided | Include code execution | Mathematical or logical verification needed |

> **Pro Tip**
>
> Chain-of-thought is most effective for multi-step reasoning problems. For simple factual recall or creative writing, CoT may actually reduce quality by introducing unnecessary verbosity. Match the technique to the task.

**Few-Shot Learning Mastery**

| Aspect | Best Practice |
|---|---|
| Example Count | Start with 3–5 examples; add more only if needed for consistency |
| Example Selection | Choose diverse, representative cases; include edge cases sparingly |
| Example Ordering | Place most similar examples closest to the actual task |
| Format Consistency | Maintain identical structure across all examples |
| Negative Examples | Include sparingly when boundary clarification is essential |
| Dynamic Selection | Consider retrieval-based example selection for varied inputs |

**Iterative Refinement Patterns**

| Pattern | Implementation | When to Use |
|---|---|---|
| Critique-Revise | Generate, self-critique, improve | Quality-critical outputs |

| Pattern | Implementation | When to Use |
|---------|----------------|-------------|
| Multi-Turn Refinement | Progressive improvement via conversation | Complex, evolving requirements |
| Self-Evaluation | Assess output against criteria | Before final delivery |
| Adversarial Testing | Probe for weaknesses | Security-sensitive applications |
| Constitutional Revision | Apply principle-based evaluation | Alignment-critical outputs |

> **❗ Common Pitfall**
>
> Common Pitfall: Over-engineering prompts with too many patterns at once. Start simple and add complexity only when simpler approaches fail. Each additional technique increases token usage and potential failure points.

### 3.3.3 Week 2 Practical Exercises

> **🖧 Project Deliverable**
>
> **Exercise 2.1: CoT Implementation Study**
> Select three distinct problem types:
>
> 1. Mathematical word problem
>
> 2. Logical reasoning puzzle
>
> 3. Code debugging scenario
>
> For each, implement both zero-shot and few-shot CoT. Document:
>
> - Accuracy comparison
>
> - Reasoning quality assessment
>
> - Token usage difference
>
> - Recommendation for when to use each
>
> **Deliverable:** Comparative analysis document with all prompts and outputs.

> **🖧 Project Deliverable**
>
> **Exercise 2.2: Few-Shot Engineering**
> Create a classification task (sentiment, intent, or category assignment):
>
> 1. Build few-shot prompts with 1, 3, 5, and 10 examples
>
> 2. Test each variant on 20 test cases
>
> 3. Measure accuracy, consistency, and token cost
>
> 4. Vary example ordering and measure impact
>
> **Deliverable:** Performance analysis with optimal configuration recommendation.

> **⋮☰ Self-Evaluation Checklist**
>
> **Week 2 Self-Assessment:**
>
> ☐ I can implement zero-shot and few-shot chain-of-thought prompting
>
> ☐ I can design effective few-shot examples for classification tasks
>
> ☐ I understand when self-consistency improves reliability
>
> ☐ I can implement a critique-and-revise pipeline
>
> ☐ I have documented at least 10 distinct prompting patterns
>
> ☐ I can articulate the trade-offs between different patterns

## 3.4 Week 3: Systematic Evaluation

### 3.4.1 Learning Objectives

By the end of Week 3, you will:

- Design appropriate evaluation metrics for different task types

- Create comprehensive test sets including edge cases

- Implement automated evaluation pipelines

- Use LLM-as-judge techniques effectively

- Establish baselines and track improvements systematically

### 3.4.2 Evaluation Framework

**Metric Selection by Task Type**

| Task Type | Primary Metrics | Implementation Notes |
|---|---|---|
| Classification | Accuracy, F1, Precision, Recall | Use stratified test sets; report per-class metrics |
| Extraction | Exact Match, Token F1 | Define matching criteria carefully; handle partial matches |
| Generation | BLEU, ROUGE, BERTScore | Choose based on what aspect matters (fluency, coverage, similarity) |
| Summarization | ROUGE-L, factual accuracy | Include human evaluation for faithfulness |
| Code Generation | Execution pass rate, correctness | Automated testing against test cases |
| Open-Ended | LLM-as-judge, human rating | Define rubrics explicitly; ensure consistency |

**Test Set Construction**

A high-quality test set includes:

- **Representative Cases (60%):** Typical inputs the system will encounter

- **Edge Cases (20%):** Boundary conditions, unusual inputs, extreme values

- **Adversarial Cases (10%):** Inputs designed to confuse or break the system

- **Regression Cases (10%):** Previously failed inputs that have been fixed

**LLM-as-Judge Implementation**

> **Pro Tip**
>
> When using LLMs to evaluate other LLM outputs, provide explicit rubrics with scoring criteria. Include examples of outputs at each score level. Consider using a different model family for evaluation to avoid systematic biases.

### 3.4.3 Week 3 Practical Exercises

> **Project Deliverable**
>
> **Exercise 3.1: Golden Test Set Creation**
> For a prompt developed in Weeks 1–2:
>
> 1. Create 30+ test cases with expected outputs
>
> 2. Categorize by case type (representative, edge, adversarial)
>
> 3. Implement automated scoring for applicable metrics
>
> 4. Run baseline evaluation and document results
>
> **Deliverable:** Test set file (JSON/YAML) plus evaluation script with baseline results.

> **Project Deliverable**
>
> **Exercise 3.2: LLM-as-Judge Pipeline**
> Build an evaluation system that:
>
> 1. Defines clear scoring rubrics (1–5 scale with descriptions)
>
> 2. Implements LLM-based evaluation with consistent prompts
>
> 3. Calculates inter-rater reliability (if using multiple judge prompts)
>
> 4. Generates summary reports with score distributions
>
> **Deliverable:** Working LLM-as-judge implementation with documentation.

## 3.5 Week 4: Specialized Prompting Domains

### 3.5.1 Learning Objectives

By the end of Week 4, you will:

- Apply prompting techniques to code generation tasks effectively

- Produce reliable structured outputs (JSON, XML, YAML)

- Work with vision-language models for multi-modal tasks

- Adapt prompting strategies for domain-specific applications

- Complete the Phase 1 capstone project

### 3.5.2  Code Generation Patterns

| Pattern | Implementation |
| --- | --- |
| Specification-First | Provide detailed requirements including edge cases before requesting code |
| Test-Driven | Supply test cases; ask for implementation that passes them |
| Incremental Development | Build functionality step-by-step, verifying each step |
| Pseudo-code Scaffolding | Outline structure first, then request implementation |
| Code Review Prompting | Analyze and improve existing code systematically |
| Documentation Generation | Generate docstrings, comments, and READMEs from code |

### 3.5.3  Structured Output Reliability

- **Schema Definition:** Always provide explicit JSON schemas or format examples

- **Validation:** Implement parsing with error handling; retry on malformed output

- **Mode Selection:** Use provider-specific structured output modes when available

- **Fallback Strategies:** Define behavior when structured output fails

### 3.5.4  Multi-Modal Prompting

When working with vision-language models:

- **Image Description:** Be specific about what aspects of the image to analyze

- **Grounding:** Reference image regions explicitly ("in the upper left," "the red object")

- **Task Clarity:** Distinguish between description, analysis, and generation tasks

- **Limitations:** Understand model limitations with text in images, fine details, spatial reasoning

## 3.6   Phase 1 Capstone: Personal Prompt Cookbook

> **⚐ Milestone**
>
> **Deliverable: Personal Prompt Cookbook**
> Compile a comprehensive prompt cookbook containing:
> **Part 1: Pattern Library (20–30 patterns)**
>
> - Pattern name and category
> - Description and use cases
> - Template with placeholders
> - 2–3 concrete examples with outputs
> - Known limitations and failure modes
> - Provider-specific notes (if any)
>
> **Part 2: Evaluation Framework**
>
> - Test cases for representative patterns
> - Automated evaluation scripts
> - Quality metrics and thresholds
> - LLM-as-judge rubrics
>
> **Part 3: Personal Insights**
>
> - Best practices discovered through experimentation
> - Anti-patterns to avoid
> - Decision framework for pattern selection
>
> **Format:** GitHub repository with Markdown documentation and supporting code.

> **☑ Assessment Checkpoint**
>
> **Phase 1 Assessment Criteria:**
>
> | Criterion | Weight |
> | --- | --- |
> | Pattern variety and coverage (20+ distinct patterns) | 25% |
> | Quality of documentation and examples | 25% |
> | Evaluation implementation and rigor | 20% |
> | Personal insights and original contributions | 15% |
> | Code quality and organization | 15% |
>
> **Passing Threshold:** 70% overall, with no criterion below 50%.

# Chapter 4

# Phase 2: Application Development

> **🕐 Time Investment**
>
> **Duration:** 6 weeks
> **Weekly Commitment:** 10–12 hours
> **Total Hours:** 60–72 hours
> **Delivery:** 2–3 Portfolio Applications

## 4.1 Phase 2 Learning Path Visualization

| Week 5 Architecture | Week 6 Tools | Week 7 RAG I | Week 8 RAG II | Week 9 Agents | Week 10 Projects |
|---|---|---|---|---|---|
| Patterns | Function Call | Chunking | Retrieval | Memory | Integration |
| State Mgmt | Tool Design | Embedding | Reranking | Planning | Deployment |

Project 1: Tool Agent  Project 2: RAG System  Project 3: Full App

## 4.2 Week 5: LLM Application Architecture

### 4.2.1 Learning Objectives

- Identify and select appropriate architectural patterns for LLM applications

- Design component interactions for maintainability and scalability

- Implement proper state management and conversation handling

- Apply caching strategies for cost and latency optimization

- Handle errors gracefully in LLM-powered systems

### 4.2.2 Architecture Pattern Catalog

| Pattern | Description | Complexity | Use Cases |
|---|---|---|---|
| Direct API | Simple request-response | ★☆☆☆ | Single-turn tasks, MVP |

| Pattern | Description | Complexity | Use Cases |
|---------|-------------|------------|-----------|
| Prompt Chain | Sequential prompts | ★★☆☆ | Multi-step workflows |
| Router | Classify and delegate | ★★☆☆ | Multi-intent systems |
| RAG Pipeline | Retrieve then generate | ★★★☆ | Knowledge-based QA |
| Agent Loop | Plan-act-observe cycle | ★★★☆ | Autonomous tasks |
| Multi-Agent | Coordinated specialists | ★★★★ | Complex workflows |
| Human-in-Loop | LLM + human verification | ★★★☆ | High-stakes decisions |

### 4.2.3 Component Design Principles

**Prompt Management**

- Store prompts in version-controlled templates, not hard-coded strings

- Implement parameterization for dynamic content injection

- Maintain separate prompts for development, testing, and production

- Document prompt purpose, expected inputs, and output format

**Context Management**

- Implement context window tracking to prevent truncation

- Design summarization strategies for long conversations

- Consider sliding window, summarization, or hybrid approaches

- Handle context overflow gracefully with clear user feedback

**Response Handling**

- Parse and validate outputs before use

- Implement retry logic with exponential backoff

- Define fallback behaviors for malformed responses

- Log all requests and responses for debugging

> **⚎ Project Deliverable**
>
> **Exercise 5.1: Architecture Design Document**
> Design an LLM application (choose: customer support bot, code review assistant, or research summarizer):
>
> 1. Create system context and component diagrams
>
> 2. Document data flows and API contracts
>
> 3. Define error handling and fallback strategies
>
> 4. Estimate costs at different usage levels
>
> 5. Plan for horizontal scaling
>
> **Deliverable:** Architecture design document with diagrams.

## 4.3 Week 6: Tool Integration and Function Calling

### 4.3.1 Learning Objectives

- Implement function calling across different provider APIs

- Design tool schemas that guide model behavior effectively

- Handle parallel and sequential tool execution

- Manage errors and edge cases in tool chains

- Apply security best practices for tool access

### 4.3.2 Tool Design Best Practices

| Aspect | Best Practice |
|---|---|
| Naming Convention | Use clear, action-oriented names: `search_database`, `send_email`, `calculate_total` |
| Description Quality | Provide detailed descriptions of when, why, and how to use each tool |
| Parameter Design | Well-typed parameters with clear descriptions, constraints, and examples |
| Return Format | Consistent, parseable return formats; include both data and status |
| Error Messages | Informative errors the LLM can interpret and respond to appropriately |
| Scope Limitation | Minimal necessary permissions; explicit boundaries on what tools can do |
| Confirmation Gates | Require confirmation for destructive or expensive operations |

> **❗ Common Pitfall**
>
> Security Warning: Tools that access external systems (databases, APIs, file systems) must validate all LLM-generated parameters. Never pass LLM output directly to system commands or database queries without sanitization. Implement rate limiting and audit logging for all tool invocations.

> **⌁ Project Deliverable**
>
> **Project 1: Multi-Tool Agent**
> Build an agent with access to 4+ tools:
>
> - Information retrieval tool (web search or database)
>
> - Computation tool (calculator, data analysis)
>
> - Communication tool (email draft, notification)
>
> - File/data tool (read/write structured data)
>
> Requirements:
>
> - Proper tool selection based on user intent
>
> - Sequential and parallel tool execution
>
> - Error handling and recovery
>
> - Conversation history integration
>
> - Security validation on all inputs
>
> **Deliverable:** Working agent with documentation and test suite.

## 4.4 Weeks 7–8: Retrieval-Augmented Generation (RAG)

### 4.4.1 Learning Objectives

- Design and implement complete RAG pipelines

- Apply appropriate document chunking strategies

- Select and configure embedding models and vector stores

- Implement retrieval optimization techniques

- Evaluate and iterate on RAG system quality

### 4.4.2 RAG Pipeline Architecture

Ingestion Pipeline

```
Documents → Chunking → Embedding → Vector Store
```

Query Pipeline

```
User Query → Query Embedding → Retrieval → Reranking

Response ← LLM Generation
```

### 4.4.3 Chunking Strategies

| Strategy | Description | Best For |
|----------|-------------|----------|
| Fixed Size | Split at token/character count | Simple implementation, uniform retrieval |
| Recursive | Split hierarchically (paragraph, sentence, word) | Varied document structures |
| Semantic | Split at topic boundaries | Topic-coherent retrieval |
| Document-Aware | Respect document structure (headers, sections) | Structured documents |
| Sliding Window | Overlapping chunks | Context preservation |
| Parent-Child | Small chunks linked to larger context | Precise retrieval with broad context |

### 4.4.4 Retrieval Optimization

| Technique | Implementation |
|-----------|----------------|
| Hybrid Search | Combine dense (embedding) and sparse (BM25) retrieval |
| Query Expansion | Generate multiple query variants to improve recall |
| HyDE | Generate hypothetical answer, embed it for retrieval |
| Reranking | Use cross-encoder to reorder initial retrieval results |
| Metadata Filtering | Pre-filter by date, source, category before semantic search |
| Multi-Query | Generate diverse queries from single user input |

> **⛕ Project Deliverable**
>
> **Project 2: RAG-Based Document Q&A**
> Build a complete RAG system:
>
> 1. Document ingestion pipeline (PDF, Markdown, HTML)
>
> 2. Configurable chunking with multiple strategies
>
> 3. Vector storage with any provider-agnostic approach
>
> 4. Conversational interface with context tracking
>
> 5. Source citation in responses
>
> 6. Evaluation suite with retrieval and generation metrics
>
> Optimization Phase:
>
> - Compare chunking strategies on same document set
>
> - Implement and compare reranking approaches
>
> - Test hybrid retrieval vs. dense-only
>
> - Document performance trade-offs
>
> **Deliverable:** Working RAG application with optimization analysis report.

## 4.5   Week 9: Conversational Agents and Memory

### 4.5.1   Learning Objectives

- Design memory architectures for different conversation types

- Implement context summarization and compression

- Build agents with consistent personas across sessions

- Handle entity tracking and reference resolution

- Manage long-term memory persistence

### 4.5.2   Memory Architecture Comparison

| Type | Implementation | Pros | Cons |
| --- | --- | --- | --- |
| Buffer | Full history in context | Complete context | Context limit hit quickly |
| Window | Last N turns only | Predictable tokens | Loses early context |
| Summary | Compressed history | Handles long chats | Information loss |
| Entity | Track entities/facts | Efficient storage | Complex implementation |

| Type | Implementation | Pros | Cons |
|---|---|---|---|
| Vector | Semantic retrieval | Relevant recall | Retrieval quality varies |
| Hybrid | Combined approaches | Balanced trade-offs | System complexity |

### 4.5.3 Persona Consistency Techniques

- **System Prompt Anchoring:** Define persona characteristics in system prompt
- **Fact Sheets:** Maintain explicit lists of persona facts for reference
- **Consistency Checking:** Validate responses against persona constraints
- **Memory Integration:** Include persona facts in memory retrieval

## 4.6 Week 10: Integration and Capstone Projects

### 4.6.1 Portfolio Project Requirements

Complete at least two portfolio-quality applications:

---

**🤖 Track A: Conversational Agent**

- Multi-turn chatbot with memory
- Consistent persona across sessions
- Tool integration for actions
- Conversation analytics

---

**🔍 Track B: Knowledge System**

- Document ingestion pipeline
- Semantic search interface
- Answer with citations
- Evaluation dashboard

---

**🗄 Track C: Data Interface**

- Natural language to structured queries
- Result visualization
- Query explanation
- Multi-source integration

> **⚑ Milestone**
>
> **Phase 2 Capstone Deliverables:**
> For each project:
>
> - Working application deployed (local or cloud)
>
> - Source code in version control
>
> - Architecture documentation with diagrams
>
> - README with setup and usage instructions
>
> - Evaluation suite with test cases
>
> - Performance metrics and analysis
>
> - Provider abstraction (works with multiple LLMs)

> **☑ Assessment Checkpoint**
>
> **Phase 2 Assessment Criteria:**
>
> | Criterion | Weight |
> | --- | --- |
> | Application functionality and completeness | 30% |
> | Code quality, architecture, and maintainability | 25% |
> | Documentation quality (README, architecture docs) | 15% |
> | Evaluation implementation and coverage | 15% |
> | Provider abstraction and portability | 15% |
>
> **Passing Threshold:** 70% overall, with working functionality required.

# Chapter 5

# Phase 3: Production Engineering

> **🕐 Time Investment**
>
> **Duration:** 4 weeks
> **Weekly Commitment:** 10–12 hours
> **Total Hours:** 40–48 hours
> **Delivery:** Production Playbook

## 5.1 Phase 3 Learning Path Visualization



| **Week 11** System Design | **Week 12** Evaluation | **Week 13** Operations | **Week 14** Security |
|---|---|---|---|
| Decision Framework | Eval Pipelines | Observability | Prompt Injection |
| Cost Modeling | LLM-as-Judge | Incident Response | Data Privacy |
| Scaling Strategy | A/B Testing | Feedback Loops | Guardrails |

**Production Readiness Review**

## 5.2 Week 11: System Design for Production

### 5.2.1 Learning Objectives

- Apply decision frameworks for prompting vs. fine-tuning vs. RAG

- Model costs accurately across different usage scenarios

- Design systems for latency, reliability, and scalability requirements

- Select appropriate models based on comprehensive criteria

- Plan multi-provider strategies for resilience

### 5.2.2 Decision Framework: Prompting vs. Fine-Tuning vs. RAG

| Factor | Better Prompting | Fine-Tuning | RAG |
|---|---|---|---|
| Knowledge Type | General/static | Style/format | Dynamic/external |
| Data Volume | Any | 100s–1000s examples | Any corpus size |
| Update Frequency | Immediate | Requires retraining | Real-time possible |
| Iteration Speed | Very fast | Days to weeks | Fast (retrieval tuning) |
| Cost Structure | Per-token | Training + inference | Storage + retrieval + inference |
| Attribution | None | None | Source citation possible |
| Hallucination | Model-dependent | Reduced for trained domain | Reduced with good retrieval |

### 5.2.3  Cost Modeling Framework

- **Input Costs:** System prompt + context + user input tokens

- **Output Costs:** Generated response tokens

- **RAG Costs:** Embedding generation + vector storage + retrieval queries

- **Tool Costs:** External API calls, compute for tool execution

- **Hidden Costs:** Retry attempts, evaluation runs, logging storage

> **💡 Pro Tip**
>
> Cost Optimization Strategies:
>
> 1. Cache frequent queries and their responses
>
> 2. Use smaller models for classification/routing, larger for generation
>
> 3. Compress system prompts; eliminate redundancy
>
> 4. Implement early stopping for streaming responses when appropriate
>
> 5. Batch requests when real-time response isn't required

## 5.3  Week 12: Evaluation at Scale

### 5.3.1  Learning Objectives

- Build comprehensive evaluation pipelines for production systems

- Implement continuous evaluation with automated alerting

- Design and execute A/B tests for prompt changes

- Detect and respond to model drift and regression

- Balance automated and human evaluation appropriately

### 5.3.2   Multi-Level Evaluation Framework

| Level | What to Measure | How to Measure | Frequency |
|---|---|---|---|
| Component | Individual prompt quality | Unit tests, golden sets | Every change |
| Pipeline | End-to-end correctness | Integration tests, traces | Daily |
| System | User-facing quality | A/B tests, user feedback | Weekly |
| Business | Impact on KPIs | Analytics, conversion | Monthly |

### 5.3.3   LLM-as-Judge Best Practices

- **Explicit Rubrics:** Define 1–5 scale with detailed descriptions for each level

- **Calibration Examples:** Include examples of outputs at each score level

- **Multiple Dimensions:** Evaluate relevance, accuracy, completeness, safety separately

- **Judge Diversity:** Use multiple judge prompts or models to reduce bias

- **Human Calibration:** Regularly compare LLM scores to human ratings

> **⛋ Project Deliverable**
>
> **Exercise 12.1: Production Evaluation Pipeline**
> Build an evaluation system that:
>
> 1. Runs scheduled evaluations against versioned test sets
>
> 2. Implements multiple evaluation strategies (exact match, semantic, LLM-judge)
>
> 3. Tracks metrics over time with visualization
>
> 4. Alerts on regression beyond thresholds
>
> 5. Integrates with your Phase 2 application
>
> **Deliverable:** Working evaluation pipeline with dashboard and alerts.

## 5.4   Week 13: LLMOps and Production Operations

### 5.4.1   Learning Objectives

- Implement comprehensive observability for LLM systems

- Design prompt versioning and deployment workflows

- Build feedback collection and processing pipelines

- Establish incident response procedures for LLM failures

- Create continuous improvement loops from production data

### 5.4.2   Observability Stack Components

| Component | Purpose | Key Metrics |
|---|---|---|
| Request Logging | Audit trail, debugging | Request/response content, latency, tokens |
| Distributed Tracing | Multi-step visibility | Span durations, tool calls, retrieval steps |
| Performance Metrics | System health | Latency percentiles, throughput, error rates |
| Cost Tracking | Budget management | Token usage, cost per request, cost trends |
| Quality Monitoring | Output health | Evaluation scores, feedback signals |

### 5.4.3   Prompt Version Management

- **Version Control:** Treat prompts as code; store in Git

- **Change Documentation:** Record rationale for each change

- **Staged Rollout:** Deploy to canary, then percentage, then full

- **Rollback Plan:** Maintain ability to instantly revert

- **A/B Testing Integration:** Support running multiple versions simultaneously

## 5.5   Week 14: Safety, Security, and Ethics

### 5.5.1   Learning Objectives

- Identify and mitigate LLM-specific security vulnerabilities

- Implement input validation and output filtering

- Design systems that protect user privacy

- Apply responsible AI principles in practice

- Navigate regulatory requirements (GDPR, AI Act concepts)

### 5.5.2   Security Vulnerability Catalog

| Vulnerability | Description | Mitigation Strategies |
|---|---|---|
| Prompt Injection | Malicious input overrides system instructions | Input validation, instruction isolation, output monitoring |
| Indirect Injection | Malicious content in retrieved documents | Content sanitization, source validation |

| Vulnerability | Description | Mitigation Strategies |
|---|---|---|
| Data Exfiltration | Extracting training data or user data | Output filtering, PII detection, rate limiting |
| Jailbreaking | Bypassing safety guardrails | Layered defenses, output classification |
| Excessive Agency | Unintended autonomous actions | Permission boundaries, human-in-loop |
| Model Extraction | Stealing model via API queries | Rate limiting, query pattern detection |

### 5.5.3   Defense-in-Depth Strategy

1. **Input Layer:** Validate, sanitize, and classify incoming requests

2. **Prompt Layer:** Isolate system instructions from user content

3. **Model Layer:** Use models with safety training; consider guardrails

4. **Output Layer:** Filter, classify, and validate generated content

5. **Monitoring Layer:** Detect anomalies, log for audit, alert on violations

> ❗ **Common Pitfall**
>
> Never trust LLM output for security-sensitive operations. Always validate:
>
> - SQL queries before execution
>
> - API calls before making requests
>
> - File paths before access
>
> - Shell commands before running
>
> - Any output that will be rendered as HTML

## 5.6   Phase 3 Capstone: Production Playbook

> 🏁 **Milestone**
>
> **Deliverable: Production Playbook**
> Create a comprehensive operational playbook containing:
> **Part 1: Decision Frameworks**
>
> - When to use prompting vs. RAG vs. fine-tuning decision tree
>
> - Model selection criteria and comparison matrix
>
> - Architecture pattern selection guide
>
> - Cost estimation templates and calculators
>
> **Part 2: Quality Management**
>
> - Evaluation strategy templates by task type
>
> - Test set creation checklists
>
> - LLM-as-judge rubric templates
>
> - Quality metrics definitions and thresholds
>
> **Part 3: Operations**
>
> - Observability setup guides and dashboards
>
> - Incident response procedures for common failures
>
> - On-call runbooks for LLM-specific issues
>
> - Feedback processing workflows
>
> - Cost monitoring and optimization playbooks
>
> **Part 4: Safety and Compliance**
>
> - Security hardening checklists
>
> - Input validation and output filtering guidelines
>
> - Privacy protection procedures
>
> - Responsible AI review checklist
>
> **Format:** Wiki, Notion, or documentation site with templates and checklists.

☑ **Assessment Checkpoint**

**Phase 3 Assessment Criteria:**

| Criterion | Weight |
| --- | --- |
| Playbook completeness and practicality | 30% |
| Decision framework clarity and usefulness | 20% |
| Operational procedures quality | 20% |
| Security and safety coverage | 15% |
| Production enhancement of Phase 2 project | 15% |

**Passing Threshold:** 70% overall.

# Chapter 6

# Phase 4: Continuous Mastery

> **⏱ Time Investment**
>
> **Duration:** Ongoing
> **Weekly Commitment:** 3–5 hours
> **Focus:** Continuous professional development

## 6.1 Mastery Development Framework

Phase 4 represents ongoing professional growth beyond the structured curriculum. Success requires establishing sustainable learning habits and contributing back to the community.

### 6.1.1 Advanced Technique Areas

| Area | Topics to Explore |
| --- | --- |
| Multi-Agent Systems | Agent coordination, task decomposition, specialized roles, consensus mechanisms |
| Advanced Reasoning | Constitutional AI, recursive self-improvement, formal verification of outputs |
| Efficiency Optimization | Prompt compression, speculative decoding, distillation techniques |
| Multimodal Integration | Vision-language-action models, embodied AI, cross-modal reasoning |
| Domain Specialization | Vertical-specific techniques (healthcare, legal, finance, code) |
| Research Frontiers | Following and implementing academic advances |

### 6.1.2 Specialization Tracks

> **‹/› Code Generation Specialist**
>
> Deep expertise in code generation, completion, review, and debugging. Focus on IDE integration, multi-file context, and repository-level understanding.

> **🔍 Knowledge Systems Architect**
>
> Advanced RAG systems, knowledge graphs, multi-hop reasoning, and enterprise search. Focus on accuracy, attribution, and scale.

> **💬 Conversational AI Expert**
>
> Advanced dialogue systems, persona consistency, emotional intelligence, and multi-party conversation. Focus on user experience and engagement.

### 6.1.3 Continuous Learning Practices

- **Weekly:** Read 2–3 new papers or blog posts on LLM techniques

- **Monthly:** Implement one new technique from recent research

- **Quarterly:** Complete a substantial side project exploring new areas

- **Annually:** Contribute to open source, write technical content, or present at meetups

### 6.1.4 Community Engagement

- Participate in prompt engineering communities (Discord, Reddit, forums)

- Share learnings through blog posts, talks, or tutorials

- Contribute to open-source tools and frameworks

- Mentor others earlier in their learning journey

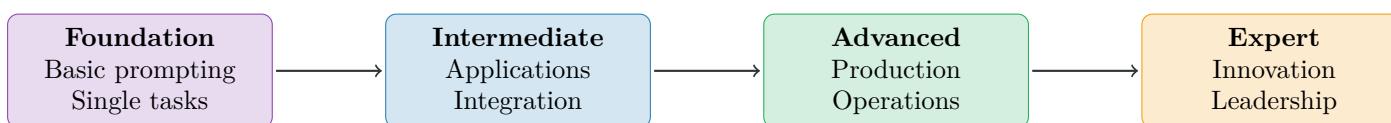- Attend conferences and local AI/ML meetups

# Part III

# Supporting Materials

# Chapter 7

# Competency Framework

## 7.1 Skill Progression Model

This framework defines four competency levels, with specific observable behaviors for each.

| Foundation<br>Basic prompting<br>Single tasks | → | Intermediate<br>Applications<br>Integration | → | Advanced<br>Production<br>Operations | → | Expert<br>Innovation<br>Leadership |

## 7.2 Detailed Competency Matrix

### 7.2.1 Prompt Engineering Competencies

| Skill | Foundation | Intermediate | Advanced | Expert |
|---|---|---|---|---|
| Basic Prompting | Write clear prompts; get useful outputs | Consistent quality across tasks | Optimize for edge cases | Teach others; create frameworks |
| Pattern Application | Use 5–10 patterns | Master 20+ patterns; select appropriately | Design custom patterns | Pioneer new patterns |
| Evaluation | Manual assessment | Automated metrics | Full evaluation pipelines | Design eval frameworks |
| Multi-Modal | Basic image+text | Complex multi-modal workflows | Production multi-modal systems | Multimodal architecture design |

### 7.2.2 Application Development Competencies

| Skill | Foundation | Intermediate | Advanced | Expert |
|---|---|---|---|---|
| Architecture | Understand patterns | Implement patterns | Design novel architectures | Set organizational standards |
| Tool Integration | Basic function calls | Multi-tool agents | Complex tool orchestration | Tool framework design |

| Skill | Foundation | Intermediate | Advanced | Expert |
|---|---|---|---|---|
| RAG Systems | Basic retrieval | Optimized pipelines | Production RAG at scale | RAG architecture innovation |
| Agents | Simple agents | Stateful conversational | Multi-agent systems | Agent framework design |

### 7.2.3  Production Engineering Competencies

| Skill | Foundation | Intermediate | Advanced | Expert |
|---|---|---|---|---|
| System Design | Understand trade-offs | Make good decisions | Optimize complex systems | Define best practices |
| Operations | Basic logging | Full observability | Incident response | Operations strategy |
| Security | Aware of risks | Implement protections | Comprehensive hardening | Security architecture |
| Cost Management | Track costs | Optimize costs | Cost-efficient design | Cost strategy leadership |

# Chapter 8

# Self-Assessment Tools

## 8.1 Phase Readiness Checklists

### 8.1.1 Ready for Phase 2?

> **⋮≡ Self-Evaluation Checklist**
>
> **Phase 1 Completion Checklist:**
> **Knowledge (must answer "yes" to all):**
>
> ☐ Can explain how tokenization affects prompt design
>
> ☐ Can describe at least 5 prompting patterns and when to use them
>
> ☐ Can implement chain-of-thought prompting effectively
>
> ☐ Can design few-shot examples for classification tasks
>
> ☐ Can create evaluation metrics appropriate to different task types
>
> **Skills (must answer "yes" to all):**
>
> ☐ Have documented 20+ prompt patterns in personal cookbook
>
> ☐ Have implemented automated evaluation for at least one prompt
>
> ☐ Have successfully generated structured output (JSON) reliably
>
> ☐ Have compared prompt performance across different models
>
> **Artifacts:**
>
> ☐ Personal Prompt Cookbook complete and documented
>
> ☐ Evaluation scripts functional and tested

### 8.1.2   Ready for Phase 3?

> **☷ Self-Evaluation Checklist**
>
> **Phase 2 Completion Checklist:**
> **Knowledge (must answer "yes" to all):**
>
> ☐ Can explain different LLM application architecture patterns
>
> ☐ Can design tool schemas for function calling
>
> ☐ Can describe RAG pipeline components and trade-offs
>
> ☐ Can explain memory architectures for conversational agents
>
> ☐ Can discuss provider abstraction strategies
>
> **Skills (must answer "yes" to all):**
>
> ☐ Have built working application with tool integration
>
> ☐ Have implemented RAG system with retrieval optimization
>
> ☐ Have built conversational agent with memory
>
> ☐ Applications work with at least two different LLM providers
>
> **Artifacts:**
>
> ☐ 2–3 portfolio applications complete and documented
>
> ☐ Architecture documentation with diagrams
>
> ☐ Evaluation suites for each application

### 8.1.3   Ready for Phase 4 (Mastery)?

> **⽇ Self-Evaluation Checklist**
>
> **Phase 3 Completion Checklist:**
> **Knowledge (must answer "yes" to all):**
>
> ☐ Can apply decision frameworks for prompting vs. fine-tuning vs. RAG
>
> ☐ Can design comprehensive evaluation pipelines
>
> ☐ Can implement observability for LLM systems
>
> ☐ Can identify and mitigate LLM security vulnerabilities
>
> ☐ Can discuss responsible AI principles and their application
>
> **Skills (must answer "yes" to all):**
>
> ☐ Have built production evaluation pipeline with alerting
>
> ☐ Have implemented security hardening for LLM application
>
> ☐ Have created operational runbooks for incident response
>
> ☐ Have optimized application for cost and latency
>
> **Artifacts:**
>
> ☐ Production Playbook complete with decision frameworks
>
> ☐ Phase 2 application enhanced for production
>
> ☐ Operational documentation complete

# Chapter 9

# Resources and References

## 9.1 Core Reading List

### 9.1.1 Primary Textbooks

| Title | Focus Area | Best For |
|---|---|---|
| *Prompt Engineering for Generative AI* | Core prompting skills | Phase 1, structured playbook |
| *AI Engineering: Building Applications with Foundation Models* | Full lifecycle | Phase 2–3, production focus |
| *LLM Engineer's Handbook* | Operations and infrastructure | Phase 3, production systems |
| *Building LLM Powered Applications* | Hands-on development | Phase 2, practical projects |
| *Prompt Engineering for LLMs* | Application building | Phase 1–2, comprehensive |

### 9.1.2 Official Documentation (Provider-Agnostic Approach)

Study documentation from multiple providers to understand both common principles and provider-specific features:

- Major LLM provider documentation (prompt engineering guides)

- API references for function calling and tool use

- Best practices and cookbook examples

- Safety and usage policies

## 9.2 Frameworks and Tools

| Category | Examples | Use Case |
| --- | --- | --- |
| Orchestration | LangChain, LlamaIndex, Haystack | Building LLM applications |
| Vector Stores | Pinecone, Weaviate, Chroma, pgvector | RAG systems |
| Observability | LangSmith, Phoenix, Helicone | Monitoring and tracing |
| Evaluation | OpenAI Evals, RAGAS, custom | Quality assessment |
| Safety | NeMo Guardrails, custom filters | Content safety |

## 9.3  Key Research Papers

1. "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models" (Wei et al.)

2. "Self-Consistency Improves Chain of Thought Reasoning" (Wang et al.)

3. "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks" (Lewis et al.)

4. "Constitutional AI: Harmlessness from AI Feedback" (Bai et al.)

5. "Tree of Thoughts: Deliberate Problem Solving with LLMs" (Yao et al.)

6. "ReAct: Synergizing Reasoning and Acting in Language Models" (Yao et al.)

7. "Toolformer: Language Models Can Teach Themselves to Use Tools" (Schick et al.)

8. "The Shift from Models to Compound AI Systems" (Berkeley AI Research)

# Chapter 10

# Schedule Options and Pacing Guides

## 10.1 Standard Track (14 Weeks)

| Week | Phase | Focus | Hours |
|------|-------|-------|-------|
| 1 | | Fundamentals | 10–12 |
| 2 | Phase 1 | Advanced Patterns | 10–12 |
| 3 | | Evaluation | 10–12 |
| 4 | | Specialization + Capstone | 10–12 |
| 5 | | Architecture | 10–12 |
| 6 | | Tool Integration | 10–12 |
| 7 | Phase 2 | RAG Fundamentals | 10–12 |
| 8 | | RAG Optimization | 10–12 |
| 9 | | Agents & Memory | 10–12 |
| 10 | | Projects + Capstone | 10–12 |
| 11 | | System Design | 10–12 |
| 12 | Phase 3 | Evaluation at Scale | 10–12 |
| 13 | | Operations | 10–12 |
| 14 | | Safety + Capstone | 10–12 |

## 10.2 Accelerated Track (10 Weeks)

For experienced developers with 15–20 hours per week:

| Week | Phase | Focus | Hours |
|------|-------|-------|-------|
| 1 | Phase 1 | Fundamentals + Patterns | 15–20 |
| 2 | | Evaluation + Specialization + Capstone | 15–20 |
| 3 | | Architecture + Tools | 15–20 |
| 4 | | RAG Complete | 15–20 |
| 5 | Phase 2 | Agents + Memory | 15–20 |
| 6 | | Project 1 | 15–20 |
| 7 | | Projects 2–3 | 15–20 |
| 8 | | System Design + Evaluation | 15–20 |
| 9 | Phase 3 | Operations + Security | 15–20 |
| 10 | | Production Playbook | 15–20 |

## 10.3   Extended Track (18+ Weeks)

For career changers or those with limited weekly time (5–8 hours):

- **Weeks 1–6:** Phase 1 (spread across 6 weeks)

- **Weeks 7–14:** Phase 2 (8 weeks for deeper practice)

- **Weeks 15–18:** Phase 3 (4 weeks at comfortable pace)

- Include additional review weeks as needed

## 10.4   Self-Paced Guidelines

- **Minimum Commitment:** 5 hours per week to maintain momentum

- **Phase Completion:** Complete each phase before starting the next

- **Project-First Option:** Start projects early, learn concepts as needed

- **Milestone Cadence:** Set weekly goals and track progress

- **Community Support:** Join study groups for accountability

# Appendix A

# Quick Reference Cards

## A.1 Prompting Pattern Quick Reference

| Pattern | When to Use | Template Snippet |
|---|---|---|
| Role Prompting | Need expert perspective | "You are a [role] with expertise in [domain]..." |
| Few-Shot | Need consistent format | "Examples: [Ex1] [Ex2] Now: [Task]" |
| Chain-of-Thought | Complex reasoning | "Let's think through this step by step..." |
| Output Format | Structured output needed | "Respond in JSON format: {...}" |
| Constraints | Specific requirements | "Requirements: 1. ... 2. ... 3. ..." |
| Self-Consistency | High-stakes decisions | Generate N responses, select most common |
| Critique-Revise | Quality improvement | "Critique the above, then improve it" |
| Decomposition | Complex multi-part task | "Break this into steps: [Task]" |

## A.2 Architecture Decision Quick Reference

| Requirement | Recommended Pattern | Key Considerations |
|---|---|---|
| Simple Q&A | Direct API | Lowest complexity, highest latency sensitivity |
| Multi-step workflow | Prompt Chain | Define clear handoffs between steps |
| Multiple intents | Router Pattern | Classifier accuracy is critical |
| External knowledge | RAG Pipeline | Chunk size and retrieval quality matter |
| Tool execution | Agent Loop | Tool design and error handling critical |

| Requirement | Recommended Pattern | Key Considerations |
| --- | --- | --- |
| High stakes | Human-in-Loop | Define clear escalation criteria |

## A.3   Security Checklist

**⬚ Self-Evaluation Checklist**

**Pre-Production Security Checklist:**

☐ Input validation implemented for all user inputs

☐ System prompt isolated from user content

☐ Output filtering for PII and sensitive content

☐ Rate limiting configured

☐ Audit logging enabled

☐ Prompt injection test suite passing

☐ Tool permissions minimized (least privilege)

☐ Human confirmation for destructive actions

☐ Error messages don't leak system details

☐ API keys properly secured (not in code)

# Appendix B

# Glossary

| Term | Definition |
| --- | --- |
| Chain-of-Thought (CoT) | Prompting technique that encourages step-by-step reasoning |
| Context Window | Maximum number of tokens a model can process in a single request |
| Embedding | Dense vector representation of text for semantic similarity |
| Few-Shot Learning | Providing examples in the prompt to guide model behavior |
| Fine-Tuning | Training a model on custom data to adapt its behavior |
| Function Calling | Enabling models to invoke external functions/tools |
| Guardrails | Safety mechanisms to constrain model outputs |
| Hallucination | Model generating plausible but factually incorrect information |
| HyDE | Hypothetical Document Embedding; generating answer for retrieval |
| LLM | Large Language Model |
| LLMOps | Operations practices specific to LLM systems |
| Prompt Injection | Attack where malicious input overrides system instructions |
| RAG | Retrieval-Augmented Generation; combining retrieval with generation |
| Reranking | Re-ordering retrieved results using a more sophisticated model |
| Self-Consistency | Generating multiple responses and selecting most common answer |
| System Prompt | Instructions provided to the model that persist across turns |
| Temperature | Parameter controlling randomness in model outputs |
| Token | Basic unit of text processing in LLMs |
| Top-p (Nucleus) | Sampling parameter limiting to most probable tokens |
| Vector Store | Database optimized for storing and querying embeddings |

| Term | Definition |
| --- | --- |
| Zero-Shot | Prompting without providing examples |