

# **Agentic System Design Course on Educative.io**

## **CHAPTER 1: AGENT DESIGN FUNDAMENTALS**

### **Lesson 1: Introduction to AI Agents**

Core Definition:

AI agents are systems that can perceive their environment, make decisions, and act autonomously to achieve specific goals following the 'perceive-reason-act' paradigm.

Three Fundamental Capabilities:

1. Perception: Sensing and extracting meaningful information from the environment (text, images, audio, sensor data)
2. Reasoning and Planning: Understanding context, selecting actions, planning steps toward goals using LLMs for language-based reasoning
3. Action Execution: Taking concrete actions (API calls, sending messages, controlling devices) based on reasoning

Key Characteristics of AI Agents:

<b>Feature</b>	<b>Description</b>
Autonomy	Acts independently without continuous human intervention

Goal-Oriented Behavior	Directs actions toward achieving predefined objectives
Perception & Feedback Loop	Observes environment and adjusts based on outcomes
Continuity	Maintains memory/context over time for multi-turn reasoning
Flexibility	Can revise plans when objectives or context change

### AI Models vs. AI Agents:

- AI Model: Passive, waits for input, computes output, stops. No goals, initiative, or context awareness beyond current input
- AI Agent: Active system using models as components. Has autonomy, memory, goals, and takes initiative. Operates over time in a perception-reason-action loop

### Three Categories of AI Agents:

1. Software-Based Agents (Sandbox Environment):
  - Operate in digital spaces (APIs, web interfaces, databases)
  - Examples: Chatbots, email triage agents, trading bots
  - No physical sensors/hardware
2. Physical Agents (Embodied Systems):
  - Interact with real world using sensors and actuators
  - Examples: Domestic robots, autonomous vehicles, factory robotic arms

- Handle uncertainty, delay, and safety in real environments
3. Hybrid and Adaptive Agents:
    - Combine digital and physical capabilities
    - Examples: AI traffic systems, healthcare assistants monitoring wearables, warehouse robots
    - Integrate multi-modal data and adapt through feedback loops

Evolution of Agents:

1. Rule-Based Agents (1970s-1980s): Fixed "if-then" logic, limited flexibility, no generalization (e.g., MYCIN, XCON)
2. Learning-Based Agents: Data-driven optimization through ML. Reinforcement Learning (RL) enables trial-and-error learning (e.g., AlphaGo), but requires large-scale training
3. LLM-Powered Agents (Current): Zero-shot generalization through massive pretraining. Key capabilities:
  - Tool use and selection
  - Dynamic task planning
  - Conversational memory
  - Generalization without retraining

When to Build an Agent:

Build agents when tasks involve:

1. Contextual decision-making: Nuanced judgment, ambiguous input, behavior adjustment based on prior interactions
  2. Complex rules: Too many conditional branches, exception cases making maintenance difficult
  3. Unstructured/natural language data: Parsing documents, emails, conversations requiring meaning extraction
- 

## Lesson 2: Agent Architecture - Core Agent

### Components

Three Core Components:

1. The Model (Agent's Brain):
  - Interprets inputs and reasons through decisions
  - LLMs work well because they're: flexible, compositional, adaptable
  - Model selection factors:
    - Task complexity (simple tasks = smaller models like Mistral/Gemma; complex = GPT-4/Claude Opus)
    - Latency requirements (fast responses need optimized models)
    - Cost considerations (mix models - small for routine, large when needed)
    - Context length (Claude 3 Sonnet, Gemini 1.5 for long conversations/documents)
2. Tools (Acting in Environment):
  - External functions, APIs, or environments for performing real-world tasks
  - Types: APIs, web search, code execution, databases, device interfaces
  - Enable agents to go beyond text generation to actual interaction
3. Instructions (Shaping Behavior):
  - Guide what agent should do, how it should behave, which goals to prioritize
  - Forms: Natural language prompts, system messages, examples, constraints, task definitions
  - Prompt engineering is critical for LLM-powered agents
  - Instructions:
    - Align output with expectations
    - Guide reasoning and tool use
    - Control safety and constraints

Component Interaction Flow:

1. Instructions set the stage (define role, goals, boundaries)
2. Model does the thinking (interprets, reasons, decides)
3. Tools handle the doing (retrieve data, send messages, trigger processes)
4. Steps may loop/repeat as needed

Why Modular Design?:

- Easily upgrade/swap models without system overhaul
- Integrate new tools to extend capabilities
- Adjust instructions for new domains/requirements

- Improve testing/debugging by isolating issues
- Enable fault isolation for system resilience
- Foster extensibility for diverse tasks

Example Workflow (Meeting Scheduling):

1. Perceive and plan (model + instructions)
  2. Gather information (Calendar API tool)
  3. Execute action (Email API tool + model for formatting)
  4. Confirm completion (Calendar API + model)
- 

## Lesson 3: Agent Architecture - Components Interaction and Agent Memory

The Architecture Blueprint - Four Primary Elements:

1. Input Interface: Receives information from environment/user (text, sensor data, API signals, multimodal inputs)
2. Reasoning Core (Model): LLM interprets input, plans steps, determines tool/action use
3. Tools and Actions: Based on reasoning, calls APIs, fetches documents, sends responses, triggers actions
4. Memory and Context: Stores/retrieves past information (conversational history, documents, user preferences)

The Agent Loop:

text

Input → Reasoning → Action → Memory → (loop continues)

This continuous loop is paramount for agentic system design, representing how agents function, learn, adapt, and evolve.

## Environmental Sensing (Perception):

- Gives agents context awareness, multimodal interaction, dynamic reactivity
- Examples of perception:
  - Virtual assistant: Speech-to-text + calendar reading
  - Customer support: Vision model/OCR for screenshot analysis
  - Factory monitoring: Sensor readings for real-time adjustment
  - Document agent: PDF/spreadsheet reading

## Perception in Workflow:

- Happens at beginning of loop
- Converts speech to text, extracts text from images (OCR), parses JSON, identifies intent
- Transforms raw inputs into meaningful signals

## Memory Systems - Three Types:

1. Short-Term Memory:
  - Holds recent context within single session
  - Tracks recent messages, follow-up questions
  - Implemented via in-session context windows or lightweight buffers
2. Long-Term Memory:
  - Persists across sessions
  - Stores: User preferences, past decisions, historical interactions
  - Requires explicit storage/retrieval (databases)
3. External/Knowledge Memory:
  - Looks up from outside sources (documents, knowledge bases, APIs)
  - Uses Retrieval-Augmented Generation (RAG) and vector databases
  - Vector databases: Store data as numerical embeddings for semantic search
  - Agent doesn't load entire knowledge base - performs efficient semantic search for relevant snippets

## Memory Operation Across Loop:

- Before reasoning: Retrieve relevant memories/facts
- During decision-making: Use memory as context
- After action: Store new information for future retrieval

Complete Agent Loop:

1. Perceive: Observe environment (message, file, sensor input)
  2. Recall: Pull relevant memories/background knowledge
  3. Reason and plan: Decide what to do (interpret intent, break down task, choose tools)
  4. Act: Execute actions (API calls, send messages, update files, physical actions)
  5. Store and learn: Save new information for future tasks
  6. Repeat
- 

## Lesson 4: Structuring Agent Behavior - Agent Orchestration Patterns

Orchestration Definition: How agents manage internal decision-making flow and how multiple agents coordinate across shared tasks.

Single-Agent Orchestration Patterns:

1. Tool Calling Loop:
  - Straightforward structure: receive task → decide tool → execute → observe → continue until goal met
  - Used in LangChain's AgentExecutor, OpenAI's tool calling
  - Simple but can lead to long, brittle chains
2. ReAct (Reasoning + Acting):
  - Alternates explicit reasoning steps with tool actions
  - Pattern: Think → Act → Observe → Think → Act
  - Improves transparency and traceability
  - Each thought/action logged for debugging/auditing
  - Example flow: "Think: I need to check flights" → "Act: Call flight API" → "Observe: See options" → "Think: Filter by price" → "Act: Return best result"
  - Used in OpenAI function-calling agents, LangChain, AutoGen
  - Verbosity can increase latency and token usage
3. Plan-and-Execute:

- Separates planning from execution
- First creates full plan (sequence of subtasks), then executes each step
- Useful for complex but known goals
- Components: Planning module (LLM call), execution module (agent/tool loop), memory
- Example: "Write Q2 report and email team" → Plan: [1. Retrieve data, 2. Create chart, 3. Write summary, 4. Send email] → Execute each
- Challenge: Plans may become outdated if environment changes

### Multi-Agent Coordination Patterns:

1. Manager-Worker Pattern:
  - Central manager oversees workflow
  - Manager: Plans, delegates to specialized workers, integrates results
  - Workers: Execute specific subtasks with distinct tools/reasoning
  - Coordination: Manager decomposes goal → assigns to workers → collects responses → finalizes output
  - Components: Central manager agent, multiple specialized worker agents, shared memory for task queues/results
  - Risk: Central point of failure, potential bottlenecks
2. Decentralized Handoff Pattern:
  - No central manager - agents operate as peers
  - Agents pass control based on task state or input type
  - Each agent maintains awareness of capabilities and context
  - Example: Travel planning - Agent A (query) → Agent B (flights) → Agent C (hotel)
  - More flexible and robust but requires careful design to avoid conflicts, deadlocks, missed responsibilities
  - Debugging more complex

### Choosing Orchestration Strategy:

#### Single vs. Multi-Agent Decision Factors:

- Task scope: Focused/narrow = single agent; multiple domains/subgoals = multi-agent
- Specialization: Distinct competencies needed = multi-agent
- Sequential vs. parallel execution: Linear = single; concurrent = multi-agent

- Observability: Single agent easier to monitor; multi-agent offers flexibility but needs coordination

#### Pattern Selection for Single-Agent:

- Plan-and-execute: Clean steps ahead of time (travel booking, multi-part documents)
- Tool calling loops: Dynamic decisions/exploration (debugging, search, knowledge synthesis)
- ReAct: Step-by-step transparency for monitoring/trust (regulated environments, educational tools)

#### Pattern Selection for Multi-Agent:

- Manager-Worker: Central planner can decompose and delegate (clearly defined roles)
- Decentralized handoffs: Independent operation, environmental changes, no single agent has full context (simulations, real-time collaboration, distributed monitoring)

#### Agent Orchestration Frameworks:

1. LangChain/LangGraph:
  - Agent executors for step-by-step tool selection
  - Multi-agent chains with logic-based routing
  - LangGraph: Framework for building agentic workflows as executable graphs
  - Features: Context memory, session state, persistent memory, workflow visualization, debugging tools
  - Highly extensible, integrates with other frameworks
2. AutoGen (Microsoft):
  - Built around multi-agent collaboration via dialogue
  - Agents as conversational entities with roles, goals, toolsets
  - Exchange structured messages (team discussions)
  - Supports arbitrary tool calling, chain-of-thought reasoning
  - Synchronous and asynchronous handoffs
  - Human-in-the-loop integration
  - Ideal for: Code writing, peer review, multi-step plans
3. CrewAI:
  - Designed around Manager-Worker architecture

- Evolved to support decentralized, hybrid, parallel execution
  - Components: Project manager (planner), workers (executors), shared task queue/memory
  - Supports human input/approval
  - Integrates with external LLMs, vector databases, APIs, other frameworks
- 

## Lesson 5: Building Trustworthy Agents - Guardrails and Human Oversight

Trust as Design Requirement: When agents take action (not just answer questions), trust becomes fundamental rather than optional.

Guardrails Definition: Safety mechanisms that wrap around agent behavior to ensure solutions stay within safe, expected limits.

Guardrail Implementation Points:

- Input guardrails: Before LLM receives input
- Output guardrails: After LLM generates response/plan
- Tool-use guardrails: Before tool execution

Types of Guardrails:

1. Contextual Grounding Checks:
  - Prevent drift off-topic or hallucinations
  - Ensure outputs remain focused and factually supported
  - Implementation:
    - Retrieval-based grounding (RAG) anchoring to source content
    - Prompt constraints to stay on topic
    - Post-response checks using classifiers/filters
    - Verify alignment and factual consistency
2. Safety and Moderation Filters:
  - Detect/block harmful, offensive, biased content

- Implementation:
  - Moderation APIs/classifiers for toxic language
  - System prompts avoiding specific topics
  - Blocked content types (hate speech, misinformation, sexual content)
  - Rate-limiting sensitive capabilities
- 3. PII and Data Protection Filters:
  - Detect, redact, restrict access to sensitive information
  - Implementation:
    - Regex/trained classifiers for PII detection
    - Redaction/masking before passing to model
    - Permission-based data access controls
    - Log/audit data access and responses
- 4. Tool Use Safeguards:
  - Restrict when, how, and under what conditions tools can be called
  - Implementation:
    - Allowlisting: Define permitted tools per context/role
    - Argument validation: Strict criteria before invocation
    - Preconditions: Logic checks before execution
    - Rate limits: Control usage frequency
    - Simulation: Generate plan, simulate without executing, then review
- 5. Rules-Based Output Validation:
  - Final check ensuring outputs meet defined rules
  - Implementation:
    - Regex checks for forbidden content
    - Format checks (JSON schema, email template)
    - Semantic checks with separate model/post-processing
    - Blocklists and allowlists for terms/phrases
- 6. Advanced Guardrails:
  - Critic/Evaluator/Optimizer: Separate agent/model reviews primary agent's outputs
  - Voting and Ensembling: Multiple LLMs/agents with consensus mechanism
  - Self-Reflection/Self-Critique: Agent reviews own actions against criteria
  - Multi-Agent Oversight: Mediate interactions between agents
  - Quantitative Monitoring: Continuous metrics tracking (risky tool calls, hallucination rate, PII detection)
  - Automated Rollback: Revert actions or restore safe state on errors

- Compliance and Explainability: Record and explain decisions for audit trails

#### Guardrail Frameworks:

- OpenAI Agents SDK: Input guardrails, safety checks
- Other libraries offer safety filters, PII detection, tool use safeguards

#### Limitations and Trade-offs:

- Over-blocking (false positives): Aggressive filters flag legitimate content
- Under-blocking (false negatives): Sophisticated attacks bypass filters
- Latency: Multiple checks increase processing time
- User friction: Excessive guardrails degrade experience

#### Best Practices:

- Layer defenses (multi-layered approach)
- Start with critical risks first
- Iterate and monitor continuously
- Maintain transparency with users
- Balance safety and utility

#### Controlling Agent Memory:

- Avoid retention of PII
- Filter inputs before storage
- Set memory expiration rules
- Restrict/permit topics

#### Human Oversight Types:

1. Human-in-the-Loop (HITL):
  - Humans directly involved in decision-making before action
  - Approval gates: Explicit review/approval before high-stakes actions
  - Escalation triggers: Low confidence/ambiguous input → defer to human expert
  - Example: "Do you want me to cancel this subscription now?"
2. Human-on-the-Loop (HOTL):

- Continuous monitoring with real-time intervention ability
  - Intervention interfaces: Dashboards allowing pause, override, edit plans in progress
  - Example: Supervisor monitoring agent via dashboard, can step in anytime
3. Human-above-the-Loop (HATL):
- Review behavior after completion for accountability/auditing
  - Audit trails and logging: Detailed logs of decisions, actions, outcomes
  - Vital for diagnosing failures, identifying patterns

Other Oversight Models:

- User-triggered oversight: End users flag unexpected behavior
- Crowdsourced/multi-reviewer: Multiple reviewers for nuanced judgment
- Automated oversight with escalation: Automated monitoring → human escalation on anomalies

Human Oversight Trade-offs:

- Latency (review delays)
  - Cost (human labor)
  - Scalability (bottlenecks)
  - Fatigue (over-escalation leads to desensitization)
- 

## **Lesson 6: Key Challenges and Design Strategies in Agentic AI Systems**

Key Challenges and Mitigation Strategies:

### **1. High Inference Latency**

Problem: LLMs are computationally intensive, causing delays especially in real-time applications

Design Strategies:

- Judicious model selection: Use smaller models for simple tasks, larger only for complex reasoning (routing pattern)
- Model optimization: Quantization, pruning, knowledge distillation
- Caching mechanisms: Cache frequent/deterministic outputs
- Parallelization: Concurrent LLM calls/tool invocations
- Asynchronous processing: Non-immediate tasks run async
- Hardware acceleration: GPUs, TPUs, edge computing

## 2. Output Uncertainty and Hallucination

Problem: LLMs generate plausible but incorrect/biased information

Design Strategies:

- RAG: Ground responses in verifiable knowledge bases
- Output validation guardrails: Rules-based checks for accuracy/safety
- Ensembling and voting: Multiple LLM instances with consensus
- Confidence scoring: Flag low-confidence responses for review
- Clear instructions: Precise prompts, explicit constraints

## 3. Memory Management and Consistency

Problem: Maintaining coherent, up-to-date context across interactions

Design Strategies:

- Structured memory systems: Clear short-term/long-term/external memory distinction using vector databases
- Explicit memory access policies: Define read/write/update permissions
- Selective memory retention: Filter what's stored and for how long
- Versioning and auditing: Track changes, ensure consistency

## 4. Scalability

Problem: System must handle increasing loads without performance degradation

Design Strategies:

- Modular architecture: Independent, scalable components
- Multi-agent orchestration: Distribute workload (Manager-Worker, decentralized)
- Efficient resource allocation: Dynamic allocation, load balancing, auto-scaling
- Optimized tool calls: Minimize redundancy, leverage batching

## 5. Integration Complexity

Problem: Interacting with diverse external tools, APIs, legacy systems

Design Strategies:

- Standardized tool definitions: Clear interfaces (JSON schemas)
- API abstraction layers: Separate agent logic from API specifics
- Robust error handling: Try-catch blocks, retry logic, fallback mechanisms
- Clear communication protocols: Explicit formats for multi-agent systems

## 6. Security and Privacy Vulnerabilities

Problem: Prompt injection, jailbreaking, data leakage

Design Strategies:

- Multi-layered guardrails: Input filtering, PII redaction, tool safeguards, output moderation
- Authentication and authorization: Strong access controls
- Secure deployment practices: Network security, encryption, audits
- Differential privacy: Protect individual data points (advanced)

## 7. Lack of Standardized Evaluation Metrics

Problem: Hard to measure dynamic, goal-oriented, interactive systems

Design Strategies:

- Task-oriented metrics: Goal achievement, problem-solving accuracy, interaction quality
- Simulation environments: Controlled, reproducible testing
- Human-centric evaluation: User surveys, A/B testing, expert reviews
- Observability and logging: Detailed interaction traces for analysis

## 8. Human-in-the-Loop Overhead

Problem: Human reviews introduce latency, cost, scalability issues

Design Strategies:

- Strategic placement: Only high-stakes/ambiguous situations
- Efficient interfaces: Intuitive dashboards with full context
- Automated summaries: Agent generates concise situation summary
- Continuous learning: Use human feedback to reduce future interventions (reflection mechanisms)

## 9. Fault Tolerance and Failure Recovery

Problem: Component failures can stall processes or crash system

Design Strategies:

- Graceful degradation: Continue with reduced functionality
- Retry mechanisms: Automatic retries with exponential backoff
- Checkpointing and state management: Save progress, resume from last good state
- Failure handoffs: Escalate to human or recovery agent
- Error boundaries: Isolate failures, prevent propagation

## Production Readiness Checklist:

### ✓ Performance and Latency:

- Response times within limits
- Tested under peak load
- Cost implications understood

### ✓ Reliability and Accuracy:

- Hallucination rates minimized
- Key decisions validated (RAG, guardrails, ensembling)
- Error handling and retries
- Graceful degradation

### ✓ Memory Management:

- Memory consistency maintained
- Sensitive data retention/expiration policies
- Scalable for long-term use

### ✓ Security and Privacy:

- Input/output filtering
- Tool access policies enforced
- Authentication/authorization robust

### ✓ Human Oversight:

- HITL points defined for high-risk actions
- Intuitive intervention interfaces
- Comprehensive audit trails

### ✓ Observability:

- Continuous monitoring and alerting
- Traceable reasoning and tool usage

✓ Scalability and Maintainability:

- Modular, extensible architecture
- Optimized deployment
- Clear documentation