

# 엔터프라이즈 AI Agent 도입 백서- 2026 엔터프라이즈 AI Agent 스택과 Harness Engineering

"AI Agent를 만들면 무한 루프에 빠지거나 한밤에 비용이 폭증한다"  
프롬프트를 아무리 정교하게 다듬어도 Context window 초과·자기 검증 불가·단일 출력이라는 4대 한계는 풀리지 않으며, 프로덕션 AI 배포의 73%가 Prompt Injection에 노출되고 학습 곡선 예산 누락으로 PoC의 절반이 무너집니다. AI Agent 플랫폼 선정부터 운영 안전장치까지 체크리스트를 확인하세요.

## Contact Us

 02-6953-5427

 [hello@msap.ai](mailto:hello@msap.ai)

 [www.msap.ai](http://www.msap.ai)

# Contents

1장: AI Agent 의 등장 배경과 해결하는 진짜 문제 — 단일 LLM 호출로는 불가능한 영역	7
1.1 ReAct · Reflexion · Voyager — 2022-2023 학술 3부작이 정의한 Agent 의 DNA	7
1.2 Auto-GPT · BabyAGI — “실패한 성공” 이 오늘날 프레임워크 요구사항을 정의한 방식	11
1.3 단일 LLM 호출이 원리적으로 풀 수 없는 4가지 문제와 Agent 가 만드는 차이	15
2장: AI Agent 를 이해하는 핵심 개념과 용어 — Planning · Action · ReAct · Harness 의 정확한 경계	19
2.1 Planning 과 Action — LLM Agent 설계의 두 축	19
2.1.1 Planning 의 정의와 피드백 유무 · 경로 수 2축 분류	20
2.1.2 Action 의 정의 — 목표 · 생성 · 공간 · 영향 4축	21
2.1.3 ReAct — Thought · Action · Observation 교차 루프가 둘을 엮는 방식	23
2.2 Tool · Memory · Harness · Orchestration — Agent 시스템을 구성하는 4대 기반 요소	24
2.2.1 Tool Use — function calling 과 MCP 가 정의하는 “Agent 의 팔다리”	24
2.2.2 Memory — Short-term · Long-term · Episodic · Semantic 4계층 구분	25
2.2.3 Harness · Orchestration · Handoff · Guardrail · HITL 의 정확한 경계	26
2.3 CoT · CoT-SC · ToT · GoT · PDDL — 경로 복잡도 축으로 정렬한 추론 기법들	28
2.3.1 CoT 와 Zero-shot CoT — “Let’s think step by step” 의 힘과 한계	28
2.3.2 CoT-SC · ToT · GoT — 경로가 여러 개일 때 일어나는 일	29
2.3.3 PDDL — LLM 과 고전 심볼릭 플래너가 만나는 외부 플래너 경로	30
3장: Planning 모듈 — LLM 이 미래 행동을 설계하는 6가지 패러다임	31
3.1 피드백 없음 — 단일 경로(CoT)와 다중 경로(CoT-SC · ToT · GoT)	31
3.1.1 CoT · Zero-shot CoT — 단일 경로 추론의 기준선	31

- 3.1.2 CoT-SC (Self-Consistency) — 다수결이 GSM8K +17.9%, SVAMP +11.0% 를 만드는 메커니즘 . . . . . 33
- 3.1.3 ToT · GoT — 트리와 그래프가 열어주는 경로 탐색 공간 . . . . . 35
- 3.2 외부 플래너 — PDDL 과 Planning Copilot 계보 . . . . . 36
  - 3.2.1 PDDL — 고전 AI 가 LLM 에 주는 형식언어 . . . . . 37
  - 3.2.2 LLM+P — LLM 생성 → 외부 플래너 → 자연어 번역 루프 . . . . . 38
  - 3.2.3 Planning Copilot — PDDL 계보가 MCP 와 결합하는 최신 흐름 . . . . . 39
- 3.3 피드백 있음 — 환경(ReAct · Voyager) · 인간(Inner Monologue) · 모델(Reflexion · Self-Refine) . . . . . 41
  - 3.3.1 환경 피드백 — ReAct · Voyager · Ghost in the Minecraft (GITM) . . . . . 41
  - 3.3.2 인간 피드백 — Inner Monologue 계열과 HITL 설계의 이론적 뿌리 . . . . . 42
  - 3.3.3 모델 피드백 — Reflexion · Self-Refine · SelfCheck . . . . . 44
- 4장: Action 모듈 — 결정을 출력으로 변환하는 4축 설계 . . . . . 45
  - 4.1 액션 목표와 생성 — 무엇을 위해, 어떻게 행동을 만들어내는가 . . . . . 45
  - 4.2 액션 공간 — 외부 도구(API · DB · 모델)와 내부 지식(계획 · 대화 · 상식) . . . . . 49
  - 4.3 액션의 영향 — 환경 상태 변화 · 내부 상태 변화 · 새 액션 파생과 피드백 루프 53

**5장: Prompt → Context → Harness 진화와 Framework → Workflow → Orchestration**

- 제품 매핑 . . . . . 57**
- 5.1 Prompt → Context → Harness Engineering 의 3-Era 전환 . . . . . 57
  - 5.1.1 Prompt Engineering Era (2022~24) — Few-shot · Role · Format 의 전성기 . . . . . 58
  - 5.1.2 Context Engineering Era (2025~) — Karpathy 가 명명한 “Context Window 채우기” 의 기술 . . . . . 59
  - 5.1.3 Harness Engineering Era (2026~) — LangGraph Deep Agents · Claude Code 가 보여주는 형태 . . . . . 60
- 5.2 Prompt → Context → Framework → Workflow → Orchestration 5단 제품 매핑 . . . . . 61

- 5.2.1 Prompt · Context — OpenAI Playground 에서 LangChain/LlamaIn-  
dex RAG 까지 . . . . . 61
- 5.2.2 Framework 단계 — LangChain Agents · LlamaIndex AgentWork-  
flow · CrewAI Agent/Task . . . . . 62
- 5.2.3 Workflow · Orchestration — LangGraph · LlamaIndex Workflows  
· Temporal + LangGraph . . . . . 63
- 5.3 3-Era × 5-Product 2축 지도와 제품 좌표 찍기 . . . . . 64
  - 5.3.1 3-Era × 5-Product 2축 지도의 구성 방법 . . . . . 64
  - 5.3.2 대표 제품들의 좌표 — LangGraph · CrewAI · OpenAI SDK · LlamaIn-  
dex · Temporal . . . . . 65
  - 5.3.3 Planning 피드백 축과 Harness 의 closed-loop/open-loop 의 일치 . . 66
- 6장: 주요 AI Agent Framework 정량 비교 — LangGraph · AutoGen · CrewAI ·  
OpenAI SDK · LlamaIndex . . . . . 66
- 6.1 5대 프레임워크 라이선스 · 생태계 · 커뮤니티 정량 지표 . . . . . 67
  - 6.1.1 LangChain · LangGraph — MIT · 100k/29.5k+ Stars · 누적 47M  
다운로드의 de facto standard . . . . . 67
  - 6.1.2 AutoGen/AG2 · CrewAI — Apache 2.0/MIT · 54.7k/45.9k Stars 와  
분기 중인 커뮤니티 . . . . . 68
  - 6.1.3 OpenAI Agents SDK · LlamaIndex Workflows — 신규 SDK 와 RAG  
기반의 현재 위치 . . . . . 69
- 6.2 설계 패러다임 · 핵심 강점 · 프로덕션 준비도 비교 . . . . . 69
  - 6.2.1 LangGraph 가 제공하는 State · Checkpoint · HITL 3축과  
Klarna/Replit/Uber/LinkedIn 프로덕션 채택 . . . . . 70
  - 6.2.2 AutoGen/AG2 의 대화형 다중 Agent 와 CrewAI 의 역할 기반 Crew —  
장점과 함정 . . . . . 70
  - 6.2.3 OpenAI Agents SDK 의 Handoff/Guardrail · LlamaIndex Workflows  
의 Event-Driven Step . . . . . 71
- 6.3 비용 · 지연시간 · MCP 지원 · 학습 곡선 — RFQ 에 그대로 올릴 7대 축 . . . . . 72

- 6.3.1 태스크당 평균 비용 \$0.35 (AutoGen) vs CrewAI TTP -40% — 비용·지연 프로파일 . . . . . 72
- 6.3.2 MCP 지원 성숙도 — LangGraph 의 1급 노드 + 스트리밍 대 나머지의 함수 호출 처리 . . . . . 73
- 6.3.3 학습 곡선 · 체크포인트링 · 업타임 — 2026 업계 컨센서스 3분할 . . . . . 74
- 7장: 엔터프라이즈 Agent 스택 — MCP · 메모리 · Observability · Durable Workflow 75**
  - 7.1 MCP — Agent의 도구 공간을 표준화한 JSON-RPC 프로토콜과 엔터프라이즈 임팩트 . . . . . 75
    - 7.1.1 MCP 2025-11-25 스펙 — tools/list · tools/call · resources · prompts · sampling 핵심 기능 . . . . . 76
    - 7.1.2 월 97M SDK 다운로드 · OpenAI/Google/MS/Anthropic 공동 지지 — 표준화의 증거 . . . . . 77
    - 7.1.3 Linux Foundation 산하 Agentic AI Foundation 이관의 의미 . . . . . 78
  - 7.2 메모리 스택 — Vector DB와 Neo4j GraphRAG + Memory Provider의 역할 분담 79
    - 7.2.1 Pinecone · Weaviate · Qdrant · Chroma — 의미 검색과 재순위의 현재 선택지 . . . . . 80
    - 7.2.2 Neo4j GraphRAG Context Provider — vector + fulltext + hybrid + Cypher traversal . . . . . 81
    - 7.2.3 Neo4j Memory Provider — 자동 entity 추출과 멀티 Agent 지식 공유 . 82
  - 7.3 Observability와 Durable Workflow — LangSmith · LangFuse · Open-Telemetry · Temporal . . . . . 83
    - 7.3.1 LangSmith 대 LangFuse — SaaS 공식과 OSS OpenTelemetry 네이티브의 선택 . . . . . 84
    - 7.3.2 OpenTelemetry for MCP — tool discovery / tool call / backend latency 분산 추적 . . . . . 85
    - 7.3.3 Temporal + LangGraph — retry/state/failure recovery와 prompt/tool/memory의 2계층 . . . . . 86
- 8장: 구축·운영 사례와 실패를 막는 5대 안전장치 87**

- 8.1 성공 시나리오 — 코딩 · 고객지원 · RPA 대체 · 개인업무 · 연구분석 5대 활용 . . . 87
  - 8.1.1 Claude Code 80.9% · Devin 13.86% on SWE-bench — 코딩 Agent  
가 만드는 PR 단위 기여 . . . . . 87
  - 8.1.2 Klarna 월 2.3M 대화 · 700 FTE · 11분→2분 — 고객지원과 RPA 하이  
브리드 . . . . . 89
  - 8.1.3 Lindy · Perplexity Deep Research — 개인업무와 연구분석 Agent 의  
현재 . . . . . 90
- 8.2 국내·해외 적용 사례 — Klarna 후퇴 · LG전자 CHATDA · KB라이프 · 한화 · Sierra  
· Devin . . . . . 91
  - 8.2.1 Klarna “AI 실패” 선언의 실상 — Trust Threshold 와 Uber-type 하이  
브리드 전환 . . . . . 91
  - 8.2.2 LG전자 CHATDA · KB라이프 · 한화 Copilot Studio · SK이노베이션 —  
국내 내부 지식 자동화 패턴 . . . . . 93
  - 8.2.3 Devin · Sierra AI · Lindy — B2B SaaS 로 제품화된 3가지 Agent . . . 94
- 8.3 운영을 좌우하는 5대 안전장치 — Step cap · Cost cap · Tool allowlist · Injection  
탐지 · Observability . . . . . 95
  - 8.3.1 Step cap · Cost cap — Kiro AWS 장애 · Claude Code 무한 루프 버그  
에서 배운 하드 제한 . . . . . 96
  - 8.3.2 Tool allowlist 와 Prompt Injection 탐지 — OWASP LLM01 대응과  
73% 노출 실태 . . . . . 97
  - 8.3.3 Observability 부재의 비용과 Goal Drift 대응 — 1인 리뷰어 10~20 시간  
의 산수 . . . . . 98
- 9장: 결론 및 권장사항 — IT 의사결정자를 위한 도입 체크리스트 . . . . . 99**
  - 9.1 프레임워크·스택 선정 체크리스트 — 규제·자동화·OpenAI 스택 3분할 기준 . . . . 99
    - 9.1.1 프레임워크 선정 — “규제 산업 → LangGraph, 내부 자동화 → CrewAI,  
OpenAI 스택 → Agents SDK” . . . . . 99
    - 9.1.2 스택 구성 — MCP + VectorDB + GraphDB + LangFuse/OTEL +  
Temporal 5레이어 . . . . . 101

- 9.1.3 OSS 코어 + 상용 관측/실행 분리 조달 — LangSmith·LangGraph Platform·CrewAI Enterprise 의 포지셔닝 . . . . . 103
- 9.2 조직과 역할 설계 — Harness Engineer · MLOps · 관측 FTE 의 배치 . . . . . 104
  - 9.2.1 AI/백엔드 엔지니어 · 데이터 사이언티스트 · CTO/아키텍트 3축 역할 매핑 105
  - 9.2.2 2~4주 학습 곡선 예산 확보 — Build-vs-Buy 와 Mid-market 권고 . . . 106
  - 9.2.3 Iconic Use Case 1개부터 — 10 부서 동시 배포 금지의 이유 . . . . . 107
- 9.3 운영 릴리스 게이트 — 5대 안전장치 · HITL · Observability FTE 의 표준 승인 양식108
  - 9.3.1 Step cap · Cost cap · Tool allowlist · Injection 탐지 · Observability 5체크박스 . . . . . 108
  - 9.3.2 HITL 게이트 — 쓰기형 Tool 에 interrupt\_on 기본값 ON . . . . . 109
  - 9.3.3 관측 FTE · 예산 승인 — 일일 1,000 req × 10~20 시간 산수 기반 인력 산정 . . . . . 110
- Appendix . . . . . 110**
  - References . . . . . 110
  - Glossary . . . . . 117
  - Footnotes . . . . . 119

# 1장: AI Agent 의 등장 배경과 해결하는 진짜 문제 — 단일 LLM 호출로는 불가능한 영역

AI Agent의 탄생은 단순히 대형 언어 모델(LLM)을 호출하는 것만으로는 해결할 수 없는 복잡한 문제에 대한 해답을 찾는 과정에서 시작되었습니다. LLM은 뛰어난 자연어 처리 능력과 다양한 태스크를 수행할 수 있지만, 실제 환경에서는 상태 변화, 자기 수정, 장기적 학습 등 시스템적 요구를 충족하지 못합니다. 이 장에서는 ReAct, Reflexion, Voyager라는 학술적 원형을 통해 Agent의 DNA를 추적하고, 오픈소스 실험(Auto-GPT, BabyAGI)이 현대 프레임워크 요구사항으로 발전한 과정을 분석합니다. 마지막으로 단일 LLM 호출의 원리적 한계와 Agent가 만들어내는 정량적 차이를 실제 사례와 수치로 증명하여, IT 의사결정자가 사내 AI 과제의 분류 기준을 명확히 할 수 있도록 안내합니다.

## 1.1 ReAct · Reflexion · Voyager — 2022-2023 학술 3부작이 정의한 Agent 의 DNA

AI Agent의 구조적 기원은 최근 2~3년간 발표된 핵심 논문들에서 명확하게 드러납니다. ReAct, Reflexion, Voyager는 각각 Thought-Action-Observation 루프, 언어적 강화학습, 평생학습 (Skill Library) 등 Agent의 핵심 설계 패턴을 학술적으로 확립하였으며, 오늘날의 다양한 프레임워크와 실무 환경에 직접적인 영향을 미쳤습니다. 이 절에서는 이러한 논문들이 어떻게 Agent의 본질적 DNA를 정의했는지, 그리고 그 구조가 실제 프레임워크에 어떻게 계승되고 있는지 구체적으로 살펴봅니다.

### 1.1.1 ReAct 가 제안한 Thought → Action → Observation 교차 루프의 의미

#### ReAct 루프의 원형 구조

ReAct(Reason+Act)는 Shunyu Yao 외(Princeton+Google Brain)가 2022년 arXiv 논문에서 최초로 제안한 패러다임입니다. 이 구조는 LLM의 추론(Reasoning, Chain-of-Thought)과 도구 호출(Acting, Tool Call)을 하나의 프롬프트 루프에서 교차(interleave)하는 방식으로, Thought → Action → Observation → Thought ... 형태의 반복적 루프를 구현합니다. ICLR

2023 발표 이후 LangGraph, CrewAI, OpenAI SDK 등 거의 모든 Agent 프레임워크가 이 구조를 기본 스키펴딩으로 채택했습니다.<sup>1</sup>

### Agent 루프의 실무적 중요성

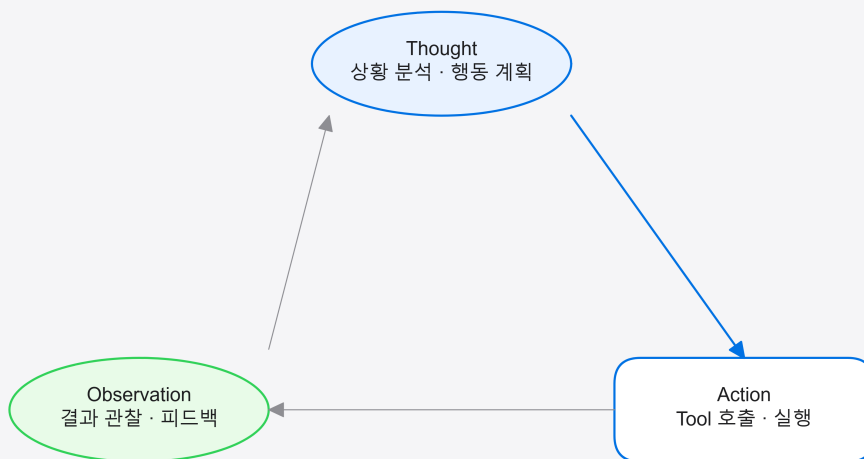
Agent 루프를 이해하지 못하면 어떤 프레임워크를 선택하더라도 결국 동일한 Thought-Action-Observation 패턴을 다른 이름으로 재학습하게 됩니다. CTO 및 실무 엔지니어는 “내가 사려는 프레임워크가 ReAct 루프를 어떻게 노드화했는가?” 를 첫 번째 평가 질문으로 삼아야 합니다. 루프의 각 단계가 명확히 관측 가능하고, 상태 전이가 투명하게 관리되는 구조가 실제 운영에서 핵심입니다.

### ReAct 루프 도식

단계	설명
Thought	LLM이 상황을 분석하고 다음 행동을 계획
Action	계획된 행동을 실제로 실행(툴 호출 등)
Observation	행동 결과를 관찰하여 다음 Thought의 입력으로 사용

#### ReAct Loop — Thought / Action / Observation

추론과 행동을 교차 반복하는 Agent 루프 — 거의 모든 프레임워크의 기본 스키펴딩



Thought → Action 은 Planning, Action → Observation 은 실행, Observation → Thought 는 피드백 — 루프가 닫히는 순간 Agent 가 성립합니다.

[그림 1] ReAct Loop — Thought / Action / Observation | 622

## ReAct의 영향과 한계

ReAct의 도입은 Agent가 단순한 LLM 호출을 넘어 다단계 추론, 조건 분기, 외부 도구 사용, 상태 기반 상호작용을 가능하게 했습니다. 그러나 루프 내에서 실패나 오류가 발생했을 때 이를 자기 수정하는 메커니즘은 추가 설계가 필요하며, 이 부분이 Reflexion 등 후속 연구로 이어집니다.

ReAct 구조는 실무적으로도 매우 중요한 의미를 지닙니다. 예를 들어, 복잡한 고객 응대 시나리오에서 Agent가 여러 번의 추론과 행동을 반복하며, 각 단계에서 얻은 정보를 바탕으로 다음 행동을 계획할 수 있습니다. 또한, 이 구조는 Agent의 디버깅과 모니터링을 용이하게 하여, 각 Thought, Action, Observation 단계별로 로그를 남기고 성능을 분석할 수 있도록 합니다. 이러한 루프 기반 설계는 단순한 LLM 호출과 달리, 실제 환경에서 발생하는 다양한 예외 상황과 복잡한 의사결정 과정을 체계적으로 처리할 수 있는 기반을 제공합니다.

### 1.1.2 Reflexion 이 도입한 언어적 강화학습(Verbal RL)과 에피소드 메모리

#### Reflexion의 핵심 아이디어

Reflexion(Noah Shinn 외, NeurIPS 2023)은 Agent가 가중치 업데이트 없이 자연어 회고(reflection)를 episodic memory에 기록하여 다음 시도에서 성능을 개선하는 언어적 강화학습(Verbal RL) 기법입니다.<sup>2</sup>이 방식은 모델 파인튜닝 없이 품질을 높일 수 있는 드문 방법으로, LangGraph의 “자기 평가 노드” 패턴의 직계 조상입니다.

#### 에피소드 메모리와 품질 개선

Reflexion은 Try → Evaluate → Reflect → Retry의 4단 루프를 통해, Agent가 실패 원인을 자연어로 서술하고 이를 메모리에 저장합니다. 다음 실행에서는 이 회고를 참조하여 같은 실수를 반복하지 않게 합니다. CTO는 Agent 품질 개선을 “파인튜닝”만으로 생각하지 말고, 메모리 설계라는 저비용 축을 반드시 검토해야 합니다.

Reflexion의 도입은 실무적으로도 큰 변화를 가져왔습니다. 예를 들어, 고객 지원 Agent가 특정 이슈에서 반복적으로 실수하는 경우, Reflexion 구조를 통해 실패 원인을 자연어로 기록하고, 다음 번에는 동일한 실수를 피할 수 있습니다. 이 방식은 기존의 파인튜닝 방식에 비해 훨씬 빠르고 저렴하게 Agent의 품질을 개선할 수 있는 장점이 있습니다. 또한, 에피소드 메모리는 반복적 태스크나 장기적 학습이 필요한 환경에서 매우 효과적이며, Agent가 점진적으로 더 나은 성능을 발휘할 수 있도록 지원합니다.

## Reflexion vs ReAct 루프 대조

구조	단계
ReAct	Thought → Action → Observation
Reflexion	Try → Evaluate → Reflect → Retry

### 적용 체크리스트

실무에서는 “Reflexion 패턴이 우리 Agent의 어느 실패 유형에 적용 가능한가?” 를 체크리스트로 삼아, 실패 패턴별 회고 메모리 설계 여부를 검토해야 합니다. 이 패턴은 특히 반복적 태스크, 장기적 학습이 필요한 환경에서 효과적입니다.

Reflexion의 또 다른 강점은, Agent가 스스로 학습하는 과정에서 발생하는 다양한 실패 사례를 체계적으로 기록하고 분석할 수 있다는 점입니다. 이를 통해 조직은 Agent의 약점을 빠르게 파악하고, 개선 방향을 명확히 설정할 수 있습니다. 또한, Reflexion 구조는 향후 Agent의 성능 평가 및 품질 관리에도 중요한 역할을 하며, 지속적인 성능 향상을 위한 기반을 제공합니다.

### 1.1.3 Voyager 가 Minecraft 에서 증명한 Skill Library + Self-Verification 평생학습

#### Voyager의 평생학습 구조

Voyager(NVIDIA/Caltech, 2023)는 LLM이 스스로 스킬을 코드로 기록하고 검증하는 평생 학습(Lifelong Learning) 구조를 Minecraft 환경에서 증명했습니다.<sup>3</sup>Voyager는 unique item 획득에서 직전 SOTA 대비 3.3배, tech tree 확장에서 15.3배 빠른 성과를 보였습니다.

Voyager의 구조는 Agent가 새로운 기술을 습득할 때마다 이를 Skill Library에 코드 형태로 저장하고, 이후 유사 상황에서 재사용합니다. 또한 Self-Verification(자기 검증) 기능을 통해, 습득한 스킬이 실제로 유효한지 자동으로 검증합니다. 엔터프라이즈 환경에서는 이 구조가 SOP(표준 운영 절차) 자동 추적과 완전히 동일하게 사상됩니다.

#### Skill Library와 Self-Verification

Voyager는 Agent가 새로운 태스크를 해결할 때마다 그 과정을 코드로 저장하여 Skill Library를 점진적으로 확장합니다. 이 Skill Library는 단순한 지식 저장소를 넘어, 실제로 재사용 가능한 코드와 절차의 집합체로 기능합니다. Self-Verification 기능은 Agent가 습득한 스킬이 실제 환경에서 유효하게 동작하는지 자동으로 검증하며, 실패 시에는 해당 스킬을 수정하거나 보완합니다.

이러한 구조는 Agent가 시간이 지날수록 점점 더 많은 태스크를 효율적으로 해결할 수 있게 하며, 반복적 업무 자동화와 지식 축적이 중요한 산업 현장에서 큰 효용을 발휘합니다.

### Skill Library 누적 vs 기준 Agent 성과 비교

항목	Voyager	기준 Agent
Unique Item 획득	3.3배 ↑	1배
Tech Tree 확장 속도	15.3배 ↑	1배

### KPI로서의 평생학습

“우리 Agent는 어제 배운 것을 오늘 재사용하는가?” 를 KPI로 설정하면, Agent의 장기적 성장과 학습 효과를 객관적으로 측정할 수 있습니다. 이 구조는 반복적 업무, 자동화, 지식 축적이 중요한 산업에서 필수적입니다.

Voyager의 사례는 실제로 Agent가 시간이 지남에 따라 점점 더 복잡한 문제를 해결할 수 있게 됨을 보여줍니다. 예를 들어, 초기에는 단순한 아이템을 획득하는 데 그쳤던 Agent가, Skill Library와 Self-Verification 구조를 통해 점차 복잡한 기술 트리와 고난이도 태스크를 빠르게 달성할 수 있게 됩니다. 이는 엔터프라이즈 환경에서 SOP(표준 운영 절차) 자동화, 신규 업무 프로세스의 신속한 적응 등 다양한 응용 가능성을 시사합니다. 또한, 평생학습 구조는 Agent의 지속적 개선과 조직 내 지식 자산의 체계적 축적에 핵심적인 역할을 합니다.

## 1.2 Auto-GPT · BabyAGI — “실패한 성공” 이 오늘날 프레임워크 요구사항을 정의한 방식

Agent 구조의 학술적 아이디어가 실제로 대중화된 계기는 오픈소스 실험인 Auto-GPT와 BabyAGI의 등장에 있습니다. 이 두 프로젝트는 상용화에는 실패했지만, 그 과정에서 드러난 다양한 한계와 문제점들이 오늘날 Agent 프레임워크의 필수 요구사항을 정의하는 데 결정적 역할을 했습니다. 이 절에서는 Auto-GPT와 BabyAGI의 구조와 실패 사례를 분석하고, 이로부터 도출된 현대 Agent 프레임워크의 핵심 설계 요구사항을 구체적으로 살펴봅니다.

## 1.2.1 Auto-GPT — GPT-4 목표·하위목표 분해 루프와 무한 루프·토큰 폭증 실패

### Auto-GPT의 구조와 한계

Auto-GPT(2023-03, Toran Bruce Richards)는 GPT-4를 목표-하위목표 분해 루프로 감싼 최초의 대중적 자율 Agent였습니다.<sup>4</sup>GitHub 공개 직후 수십만 개발자가 실험에 참여했으나, 상용화에는 실패했습니다. 그 원인은 무한 루프, 토큰 폭증, hallucinated action(환각된 행동)이었습니다.

Auto-GPT의 구조는 사용자가 입력한 목표를 여러 하위 목표로 분해하고, 각 하위 목표를 순차적으로 해결하는 방식으로 설계되었습니다. 그러나 실제로는 목표 분해 과정에서 반복적으로 동일한 작업을 수행하거나, 불필요하게 많은 토큰을 소비하는 문제가 빈번하게 발생했습니다. 또한, LLM의 한계로 인해 실제로 존재하지 않는 도구를 호출하거나, 비현실적인 행동을 시도하는 등 hallucinated action 문제가 심각하게 나타났습니다.

### 실패 모드와 현대 프레임워크 대응

Auto-GPT의 실패 모드는 오늘날 Agent 운영 장애와 동일합니다. 무한 루프는 step cap(최대 반복 횟수 제한), 토큰 폭증은 cost cap(비용 제한), hallucinated action은 tool allowlist(허용 도구 목록)와 guardrail(보호 장치)로 대응합니다.

### 실패 모드 → 프레임워크 대응 기능 매핑표

실패 모드	현대 프레임워크 대응 기능
무한 루프	step cap (max_iter)
토큰 폭증	cost cap (비용 제한)
hallucinated action	tool allowlist, guardrail

### 실무 적용 점검

“우리 PoC가 Auto-GPT 3대 실패 모드를 피할 구조를 가졌는가?” 를 자가 점검 기준으로 삼아야 하며, PoC 단계에서 반드시 해당 대응 기능을 설계에 포함해야 합니다.

Auto-GPT의 실패 사례는 실제로 많은 기업들이 PoC 단계에서 겪는 문제와 일치합니다. 예를 들어, 반복적인 태스크 수행 중 Agent가 무한 루프에 빠지거나, 예상치 못한 비용 폭증으로 프로젝트가 중단되는 경우가 있습니다. 따라서 현대 Agent 프레임워크를 도입할 때는 반드시 step cap,

cost cap, tool allowlist, guardrail 등 주요 대응 기능의 존재와 구현 방식을 꼼꼼히 확인해야 하며, 이를 RFQ(입찰 요청서) 단계에서 명확히 요구하는 것이 바람직합니다.

### 1.2.2 BabyAGI — 140줄 Python 이 정의한 Task Creation → Execution → Prioritization 3-loop

#### BabyAGI의 3-loop 구조

BabyAGI(2023-03, Yohei Nakajima)는 LLM + Vector Memory + Task Queue로 “Task Creation → Execution → Prioritization” 루프를 약 140줄 Python으로 구현한 경량 Agent입니다. 이후 CrewAI, AutoGen 등 거의 모든 Multi-Agent 프레임워크가 이 3-loop 구조를 계승했습니다.

BabyAGI의 구조는 매우 단순하지만, Agent 시스템의 본질적인 요소를 잘 드러냅니다. Task Creation 단계에서는 새로운 태스크를 생성하고, Execution 단계에서는 해당 태스크를 실제로 실행하며, Prioritization 단계에서는 여러 태스크의 우선순위를 결정합니다. 이 3-loop 구조는 멀티 Agent 시스템의 기본적인 동작 원리로 자리잡았으며, 이후 등장한 다양한 프레임워크에서 핵심 설계 패턴으로 채택되었습니다.

#### 구조적 언어와 프레임워크 매핑

BabyAGI의 구조는 멀티 Agent 프레임워크의 내부 구조를 읽어낼 “언어”를 제공합니다. 새 프레임워크를 평가할 때 “BabyAGI 3-loop의 어느 부분을 어떻게 개선했는가?”를 질문하면, 내부 컴포넌트의 역할과 개선점을 명확히 파악할 수 있습니다.

#### BabyAGI 3-loop 도식 → CrewAI/AutoGen 대응 컴포넌트 매핑

BabyAGI 단계	CrewAI/AutoGen 컴포넌트
Task Creation	Agent/Task 생성
Execution	Agent 실행/Tool 호출
Prioritization	Task Queue 관리/우선순위 결정

#### 실무적 적용

이 구조는 Task 기반 자동화, 멀티 Agent 협업, 우선순위 관리가 필요한 환경에서 필수적입니다. 프레임워크 선정 시 해당 3-loop의 구현 방식과 확장성을 반드시 검토해야 합니다.

실제로, 기업 내에서 여러 업무를 동시에 처리해야 하는 상황에서는 Task Queue와 우선순위 관리가 매우 중요합니다. 예를 들어, 고객 문의 처리, 데이터 분석, 리포트 생성 등 다양한 태스크가 동시에 발생할 때, BabyAGI와 같은 구조를 적용하면 각 태스크의 생성, 실행, 우선순위 조정이 체계적으로 이루어질 수 있습니다. 또한, 이 구조는 멀티 Agent 환경에서 각 Agent의 역할 분담과 협업을 효율적으로 지원하며, 시스템 확장성 측면에서도 매우 유리합니다.

### 1.2.3 프레임워크 요구사항의 역산 — step cap · cost cap · tool schema · memory 계층

#### 실패에서 도출된 요구사항

Auto-GPT와 BabyAGI의 실패를 1:1 역산하면 현대 프레임워크의 필수 기능이 도출됩니다. 하드 step cap(최대 반복 횟수), cost cap(비용 제한), 구조화 tool schema(도구 정의), 계층화 memory(메모리 구조)는 오늘날 LangGraph의 checkpointer, OpenAI SDK의 guardrail, CrewAI의 max\_iter 등에서 구현되고 있습니다.<sup>45</sup>

실제로, 이러한 요구사항들은 Agent 시스템의 운영 안정성과 비용 통제, 보안, 확장성에 직결되는 핵심 요소입니다. 예를 들어, step cap은 Agent가 무한 루프에 빠지는 것을 방지하고, cost cap은 예산 초과를 방지합니다. tool schema는 Agent가 사용할 수 있는 도구를 명확히 정의하여 보안과 신뢰성을 높이고, memory 계층은 단기/장기/에피소드 메모리를 체계적으로 관리하여 Agent의 성능을 극대화합니다.

#### 프레임워크 선정 RFQ 체크박스

프레임워크 선정 시 “이 5개 요구사항을 어떻게 제공하는가”를 RFQ(입찰 요청서)의 필수 체크박스로 고정해야 합니다. 이는 운영 안정성, 비용 통제, 보안, 확장성의 핵심 평가 기준입니다.

#### 실패 모드 ↔ 프레임워크 기능 대응표

요구사항	대표 구현체
step cap	LangGraph checkpointer, CrewAI max_iter
cost cap	OpenAI SDK guardrail
tool schema	MCP, function calling
memory 계층	VectorDB, GraphDB, episodic memory

## 실무적 권장사항

이 5개 항목은 프레임워크 선정 RFQ의 필수 체크박스로 삼아, 도입 전 반드시 각 기능의 구현 여부와 성숙도를 평가해야 합니다.

또한, 실제 도입 사례를 보면, 이와 같은 요구사항을 명확히 체크하지 않은 경우 운영 중 예상치 못한 장애나 비용 초과, 보안 이슈가 빈번하게 발생합니다. 따라서, 프레임워크 선정 단계에서부터 step cap, cost cap, tool schema, memory 계층 등 주요 기능의 구현 방식과 성숙도를 꼼꼼히 점검하는 것이 성공적인 Agent 시스템 구축의 핵심입니다.

## 1.3 단일 LLM 호출이 원리적으로 풀 수 없는 4가지 문제와 Agent 가 만드는 차이

Agent의 존재 이유는 단일 LLM 호출의 원리적 한계에서 비롯됩니다. 이 절에서는 LLM 호출의 한계와 Agent가 만들어내는 구체적 정량 차이를 사례와 수치로 증명합니다.

단일 LLM 호출은 자연어 처리나 간단한 태스크에는 매우 효과적이지만, 복잡한 상태 변화, 장기적 메모리, 자기 검증, 외부 시스템 연동 등 실제 비즈니스 환경에서 요구되는 고차원적 문제를 해결하는 데에는 근본적인 한계가 있습니다. Agent 구조는 이러한 한계를 극복하기 위해 설계되었으며, 실제로 다양한 산업 현장에서 단일 LLM 호출 대비 월등한 성과를 보이고 있습니다. 이 절에서는 단일 LLM 호출의 4대 한계와 Agent 구조의 대응 방안, 그리고 실제 정량적 성과 차이를 구체적으로 살펴봅니다.

### 1.3.1 Context 초과 · 실제 상태 접근 불가 · 자기 검증 불가 · 단일 출력 — 단일 호출의 4가지 원리 한계

#### 단일 LLM 호출의 4대 한계

Lilian Weng 블로그와 LangChain 프로덕션 포스트에 따르면, 단일 LLM 호출은 다음과 같은 4가지 근본적 한계가 있습니다.[67](#)

1. Context window 초과: LLM의 입력 컨텍스트 크기 제한으로, 복잡한 정보나 장기적 메모리를 모두 담을 수 없습니다.

2. 외부 실제 상태 접근 불가: 최신 데이터, 내부 DB, 실시간 환경 상태 등 외부 정보를 직접 조회하거나 반영할 수 없습니다.
3. 자기 검증 불가능: LLM은 자신의 출력이 올바른지 검증하거나, 실패 시 재시도·수정하는 기능이 없습니다.
4. 단일 출력: 사전 결정된 한 번의 출력만 가능하며, 다단계 추론, 조건 분기, 오류 복구가 불가능합니다.

### 프롬프트 정교화의 한계

프롬프트를 아무리 정교하게 설계해도 위 4가지 한계는 원리적으로 해결되지 않습니다. 따라서 복잡한 상태 기반 상호작용, 자기 수정, 장기적 학습이 필요한 태스크에서는 Agent 구조가 필수적입니다.

### 단일 호출 4한계 ↔ Agent 루프 대응 표

단일 LLM 한계	Agent 구조의 대응
Context window 초과	메모리 계층(Short/Long/Episodic)
외부 상태 접근 불가	Tool Use(MCP, function calling)
자기 검증 불가	Reflexion, 자기 평가 노드
단일 출력	루프 기반 다단계 추론, 재시도

### 실무적 트라이아지

내부 사내 AI 과제 목록을 위 4대 한계로 트라이아지(분류)하면, 단일 LLM 호출로 처리 가능한 것과 Agent 구조가 필요한 것을 명확히 구분할 수 있습니다.

실제로, 많은 기업들이 단일 LLM 호출만으로는 해결할 수 없는 문제에 직면하고 있습니다. 예를 들어, 고객 이력과 실시간 주문 상태를 동시에 반영해야 하는 상담 업무, 복잡한 규칙 기반의 의사결정, 오류 발생 시 자동 복구가 필요한 프로세스 등은 반드시 Agent 구조가 필요합니다. 이러한 한계는 프롬프트 엔지니어링만으로는 극복할 수 없으며, 시스템적 설계와 Agent 구조의 도입이 필수적입니다.

### 1.3.2 Prompt Engineering 이 부딪히는 한계선과 Karpathy 의 “Context Engineering” 전환

#### Prompt Engineering의 한계

Prompt Engineering은 단일 instruction의 형태, 어휘, 예시를 최적화하는 “한 문장 쓰는 기술”입니다. 하지만 복잡성이 커지면 프롬프트가 누더기화(prompt sprawl)되고, 재현성이 떨어집니다. 이는 실무에서 반복적 수정, 관리 비용 증가, 품질 저하로 이어집니다.

Prompt Engineering은 LLM의 성능을 극대화하기 위한 중요한 기술이지만, 실제로 복잡한 태스크나 장기적 맥락이 필요한 상황에서는 한계에 부딪힙니다. 예를 들어, 다양한 예외 상황을 모두 프롬프트에 반영하려다 보면 프롬프트가 지나치게 길어지고, 관리가 어려워집니다. 또한, 동일한 프롬프트라도 입력 데이터나 상황에 따라 결과가 달라지는 등 재현성 문제가 발생할 수 있습니다.

#### Karpathy의 Context Engineering 프레이밍

2025-06 Karpathy는 “Prompt engineering은 한 문장 쓰는 기술, Context engineering은 context window를 채우는 기술”로 프레이밍을 전환했습니다.<sup>89</sup> Context Engineering은 RAG, 메모리 계층, 툴 스키마 삽입 등 다음 step에 필요한 정보를 context window에 구조적으로 채우는 기술입니다.

Context Engineering은 단순히 프롬프트를 잘 쓰는 것을 넘어, LLM의 context window에 어떤 정보를, 어떤 구조로 넣을 것인가를 설계하는 기술입니다. 예를 들어, RAG(Retrieval-Augmented Generation)를 통해 외부 지식이나 실시간 데이터를 context window에 삽입하거나, 메모리 계층을 활용해 과거 대화 이력이나 중요한 정보를 체계적으로 관리할 수 있습니다. 또한, Tool Schema를 통해 Agent가 사용할 수 있는 도구와 그 사용법을 명확히 정의함으로써, LLM의 활용 범위를 크게 확장할 수 있습니다.

#### Prompt vs Context 역할 축 비교

축	Prompt Engineering	Context Engineering
정의	한 문장 쓰는 기술	context window 채우기
대표 기법	Few-shot, Role, Format	RAG, Memory, Tool Schema
조직 역할	프롬프트 엔지니어	Context/Harness 엔지니어

### 조직 구조 변화

조직 내 “프롬프트 엔지니어” 포지션은 “Context/Harness 엔지니어”로 진화해야 하며, 복잡한 Agent 시스템에서는 컨텍스트 설계가 1급 엔지니어링으로 자리잡습니다.

실제로, 최근 많은 기업들이 프롬프트 엔지니어링 팀을 Context Engineering 또는 Harness Engineering 팀으로 개편하고 있습니다. 이는 LLM 기반 시스템의 복잡성이 증가함에 따라, 단순한 프롬프트 작성 능력만으로는 경쟁력을 유지할 수 없기 때문입니다. 앞으로는 context window를 효과적으로 설계하고, 필요한 정보를 적시에 제공할 수 있는 엔지니어링 역량이 조직의 핵심 경쟁력이 될 것입니다.

### 1.3.3 SWE-bench 80.9% · Klarna 월 2.3M 대화 — Agent 가 만든 정량 차이

#### Agent의 정량적 성과

Agent가 만든 실제 차이는 수치로 명확합니다. Claude Code는 SWE-bench에서 80.9%의 성과를 기록했으며, 단일 LLM 호출은 2023년 기준 5% 미만에 불과합니다.

이러한 수치는 Agent 구조가 실제로 단일 LLM 호출 대비 얼마나 큰 성과 차이를 만들어내는지 명확히 보여줍니다. SWE-bench는 AI 기반 소프트웨어 코드 품질 벤치마크로, Agent 구조를 적용한 Claude Code가 80.9%의 높은 점수를 기록한 반면, 단일 LLM 호출은 5% 미만에 그쳤습니다. Klarna의 사례에서는 Agent가 월 2.3백만 건의 고객 대화를 처리하고, 700명에 해당하는 정규직 인력(FTE) 워크로드를 대체했으며, 고객 이슈 해결 시간을 11분에서 2분으로 단축하는 등 실질적인 ROI 개선 효과를 입증했습니다.

#### 단일 호출 vs Agent 정량 대비표

지표	단일 LLM 호출	Agent 구조
SWE-bench 점수	5% 미만	80.9% (Claude Code)
Klarna 대화량	-	월 2.3M 대화
워크로드 대체	-	700 FTE
정확도 개선	-	+40% (복잡 추론)

## ROI 측정 정의

CFO/CTO는 제시된 수치를 사내 Baseline과 비교하여 Agent 도입의 ROI를 명확히 측정해야 합니다. SWE-bench, Klarna 사례는 Agent 구조가 단일 LLM 호출 대비 월등한 성과를 보인다는 점을 수치로 증명합니다.

이러한 정량적 지표는 Agent 도입의 필요성과 효과를 객관적으로 평가하는 데 매우 유용합니다. 실제로, 많은 기업들이 Agent 구조 도입 전후의 성과를 비교하여, 투자 대비 효과(ROI)를 명확히 측정하고 있습니다. SWE-bench, Klarna 사례는 Agent 구조가 단일 LLM 호출 대비 얼마나 큰 성과 차이를 만들어낼 수 있는지, 그리고 복잡 추론 태스크에서 정확도와 효율성을 얼마나 높일 수 있는지를 명확히 보여줍니다.

## 2장: AI Agent 를 이해하는 핵심 개념과 용어 — Planning · Action · ReAct · Harness 의 정확한 경계

AI Agent 기술의 도입과 프레임워크 선택 과정에서 IT 의사결정자가 반드시 숙지해야 할 핵심 개념과 용어를 본 장에서 체계적으로 정리합니다. Planning, Action, Tool, Memory, ReAct, CoT, ToT, GoT, PDDL, Reflexion, Harness, Orchestration, Handoff, Guardrail, HITL 등 주요 용어의 정의와 상호 관계를 명확히 고정함으로써, 이후 기술 비교 및 설계 논의에서 발생할 수 있는 오해를 사전에 방지할 수 있습니다. 이 용어 체계는 사내 Glossary에 복제하여 팀 간 커뮤니케이션의 기준으로 삼을 수 있으며, 실제 엔터프라이즈 환경에서의 AI Agent 설계 및 평가에 일관된 기준을 제공합니다. 각 용어는 단순한 정의를 넘어, 실무적 적용과 기술적 함의를 함께 다루어 이후 장의 심층 분석과 연결될 수 있도록 구성하였습니다.

### 2.1 Planning 과 Action — LLM Agent 설계의 두 축

Agent 설계에서 가장 중요한 두 축은 바로 Planning(미래 행동의 설계)과 Action(계획의 실행 및 외부 세계에의 반영)입니다. 이 두 축은 LLM 기반 Agent의 근본적 구조와 동작 원리를 이해하는

데 필수적인 개념입니다. Planning은 복잡한 문제를 단계별로 분해하고, 각 단계의 목표와 경로를 설계하는 과정이며, Action은 이러한 계획을 실제로 실행하여 외부 시스템과 상호작용하거나 결과를 산출하는 역할을 담당합니다. 본 절에서는 이 두 축의 정의와 분류 기준을 명확히 하여, 이후 심층 분석(3장 Planning, 4장 Action)에서의 논의와 기술 비교의 기준점을 제공합니다.

## 2.1.1 Planning 의 정의와 피드백 유무 · 경로 수 2축 분류

### Planning 개념과 역할

Planning은 AI Agent가 미래의 행동을 다단계로 분해하고 설계하는 능력입니다. LLM 기반 Agent는 단순한 단일 호출로는 복잡한 문제를 해결할 수 없기 때문에, 여러 단계의 계획을 세우고 각 단계별로 피드백을 받아가며 목표 달성을 시도합니다. 이 과정에서 Planning은 단순한 명령 실행이 아니라, 목표를 세분화하고 각 단계별로 최적의 경로를 탐색하는 전략적 사고를 포함합니다. 예를 들어, 고객 문의에 대한 자동 응답 Agent를 설계할 때, 단순히 질문에 답하는 것이 아니라, 문의 유형 분류, 정보 검색, 답변 생성, 후속 조치 등 여러 단계를 계획해야 합니다. 이러한 다단계 Planning은 Agent가 복잡한 업무를 처리할 수 있게 하며, 각 단계에서의 피드백을 통해 성능을 지속적으로 개선할 수 있습니다.

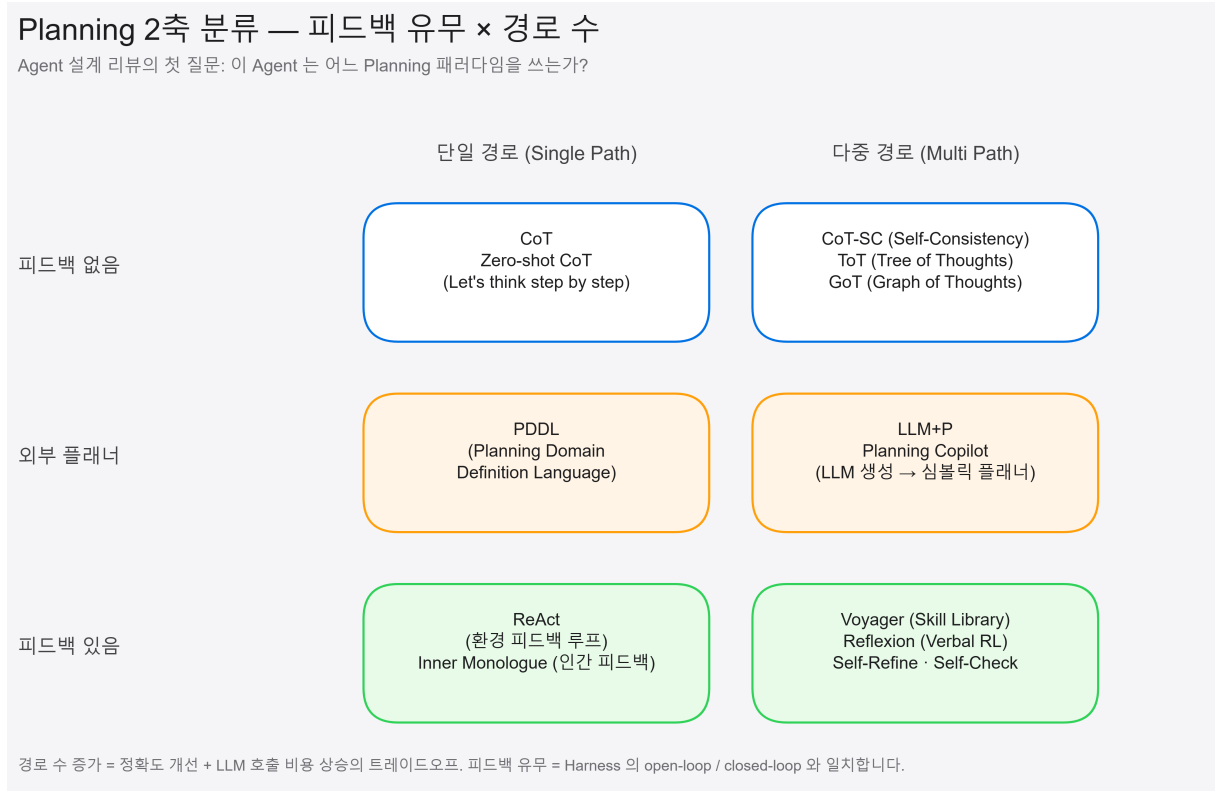
### 피드백 유무와 경로 수 분류

Planning은 크게 두 가지 축으로 분류할 수 있습니다. 첫 번째는 피드백의 유무입니다. 피드백이 없는 Planning은 LLM이 자체적으로 경로를 설계하고 실행하며, 피드백이 있는 Planning은 환경, 인간, 또는 모델로부터 입력을 받아 계획을 수정합니다. 두 번째는 경로 수입니다. 단일 경로 Planning은 한 가지 해결책만을 탐색하며, 다중 경로 Planning은 여러 경로를 샘플링하고 평가하여 최적의 경로를 선택합니다. 외부 플래너(PDDL 등)를 사용하는 경우, LLM이 계획을 생성하고 심볼릭 플래너가 실제 경로를 계산하는 방식으로 분류됩니다. 예를 들어, 다중 경로 Planning에서는 여러 답변 후보를 생성하여 다수결로 최적의 답을 선택하는 Self-Consistency(CoT-SC) 기법이 활용됩니다.

### 2축 매트릭스 표

피드백 유무	경로 수	대표 기법/패러다임
없음	단일	CoT, Zero-shot CoT
없음	다중	CoT-SC, ToT, GoT

외부 플래너	다중/단일	PDDL, LLM+P
있음	환경/인간/모델	ReAct, Voyager, Reflexion, Inner Monologue



[그림 2] Planning 2축 매트릭스 — 피드백 유무 × 경로 수|637

### Agent 설계 리뷰 적용

Agent 설계 리뷰에서는 “이 Agent는 어느 Planning 패러다임을 쓰는가?” 를 첫 질문으로 설정해야 합니다. 피드백 유무와 경로 수의 분류는 이후 심층 분석과 성능 개선, 비용 평가의 기준이 됩니다. 실제로 엔터프라이즈 환경에서는 복잡한 업무 프로세스를 자동화할 때, 단일 경로 Planning만으로는 충분하지 않은 경우가 많으므로, 다중 경로 및 피드백 기반 Planning의 도입 여부가 중요한 평가 요소가 됩니다. 또한, 외부 플래너와의 연동 여부에 따라 시스템 아키텍처와 유지보수 전략이 달라질 수 있으므로, 설계 단계에서부터 이 두 축을 명확히 구분하여 문서화하는 것이 바람직합니다.

## 2.1.2 Action 의 정의 — 목표 · 생성 · 공간 · 영향 4축

### Action의 개념과 역할

Action은 Agent가 계획한 결정을 실제 외부 세계에 출력하는 과정입니다. 이는 단순히 LLM의 답변을 생성하는 것이 아니라, API 호출, DB 업데이트, 외부 시스템 연동, 대화, 환경 탐색 등 다양한 형태로 실현됩니다. Action은 Agent의 실질적 영향력을 결정하는 핵심 축입니다. 예를 들어, 고객 지원 Agent가 단순히 답변만 제공하는 것이 아니라, 실제로 티켓을 생성하거나 데이터베이스를 업데이트하는 등의 행동을 취할 수 있다면, 그 Agent의 Action 범위는 훨씬 넓어집니다. Action의 설계는 Agent의 활용 가능성과 한계를 결정짓는 요소이므로, 각 축을 명확히 정의하고 적용하는 것이 중요합니다.

### Action 4축 표준 분류

Action은 다음 네 가지 축으로 분류됩니다.

1. **목표:** 작업 완료(Task completion), 커뮤니케이션(Communication), 환경 탐색(Environment exploration)
2. **생성:** 과거 경험 재활용(Memory recollection), 계획 준수(Plan-following)
3. **공간:** 외부 도구(API, DB, 모델), 내부 지식(계획, 대화, 상식)
4. **영향:** 환경 변화(External state change), 내부 상태 변화(Internal state change), 새 액션 파생(New action)

예를 들어, 환경 탐색을 목표로 하는 Agent는 외부 API를 통해 데이터를 수집하고, 그 결과를 바탕으로 내부 상태를 업데이트하거나 새로운 액션을 파생시킬 수 있습니다. 반면, 커뮤니케이션을 목표로 하는 Agent는 주로 대화와 정보 전달에 집중하며, 내부 지식과 과거 경험을 활용하여 응답의 품질을 높입니다.

### Action 4축 설계 체크리스트

Action 축	분류/예시
목표	작업 완료, 커뮤니케이션, 환경 탐색
생성	기억상상, 계획추종
공간	외부 도구(API/DB/모델), 내부 지식
영향	환경 변화, 내부 상태 변화, 새 액션

### 설계 리뷰 적용

새 Agent를 구축할 때 이 4축을 설계 문서의 고정 섹션으로 사용하면, Agent 간 비교가 균질해지고 실패 분석(goal drift, hallucinated tool call 등)이 해당 축에 귀속될 수 있습니다. 예를 들어, 특정 Agent가 예상치 못한 행동을 했을 때, 그 원인이 목표 설정의 오류인지, 생성 방식의 한계인지, 공간(외부 도구/내부 지식) 선택의 문제인지, 혹은 영향(상태 변화/새 액션) 관리의 미흡인지 명확히 파악할 수 있습니다. 이러한 구조화된 Action 분류는 엔터프라이즈 환경에서의 Agent 도입 및 품질 관리에 실질적인 도움이 됩니다.

### 2.1.3 ReAct — Thought · Action · Observation 교차 루프가 둘을 엮는 방식

#### ReAct 루프의 정의와 의미

ReAct는 Reasoning(추론)과 Acting(행동)을 하나의 프롬프트 루프에서 교차(interleave)하는 패러다임입니다. Thought(생각) → Action(행동) → Observation(관찰) → Thought ...의 반복 구조로, Agent가 추론과 행동을 번갈아 수행하며 상태 기반 상호작용과 자기 수정이 가능합니다. 이는 LangGraph, CrewAI, OpenAI SDK, LlamaIndex, AutoGen 등 거의 모든 Agent 프레임워크의 기본 스키펀딩입니다. ReAct 구조는 Agent가 단순히 입력에 반응하는 것이 아니라, 자신의 행동 결과를 관찰하고 다음 행동을 조정할 수 있게 해줍니다. 예를 들어, 정보 검색 Agent가 검색 결과를 관찰한 후, 추가 검색이나 답변 생성을 반복적으로 수행하는 것이 대표적인 ReAct 루프의 활용 사례입니다.

#### ReAct 루프와 Planning/Action 축 사상도

Thought → Action → Observation → Thought ...

(Planning)      (Action)      (Feedback)      (Planning)

이 구조에서 Thought 단계는 Planning에 해당하며, Action 단계는 실제 행동을 의미합니다. Observation 단계에서 얻은 피드백을 바탕으로 다음 Thought(Planning) 단계가 조정되므로, Agent의 자기 수정 및 적응 능력이 크게 향상됩니다.

#### 프레임워크 평가 기준

프레임워크를 평가할 때 “ReAct 루프의 각 단계를 어디서 관측(observable) 하는가?” 를 도입 평가 기준으로 삼아야 합니다. 이는 Agent의 상태 추적, 실패 복구, 품질 개선에 직접적인 영향을 미칩니다. 예를 들어, LangGraph와 같은 프레임워크는 Thought, Action, Observation 각 단계를 명확히 분리하여 로그와 모니터링이 가능하게 하므로, 엔터프라이즈 환경에서의 운영 및

디버깅에 큰 장점을 제공합니다. 반면, 일부 프레임워크는 이러한 단계가 불투명하게 처리되어, 문제 발생 시 원인 분석이 어려울 수 있습니다. 따라서, ReAct 루프의 각 단계가 얼마나 명확하게 드러나는지, 그리고 이를 통해 얼마나 효과적으로 Agent의 동작을 관리할 수 있는지가 프레임워크 선택의 중요한 기준이 됩니다.

## 2.2 Tool · Memory · Harness · Orchestration — Agent 시스템을 구성하는 4대 기반 요소

Agent 시스템은 단순한 LLM 호출을 넘어 Tool, Memory, Harness, Orchestration 등 다양한 레이어로 구성됩니다. 이 네 가지 기반 요소는 ReAct 루프를 둘러싼 시스템적 구조를 형성하며, 엔터프라이즈 스택 설계(7장)에서 실제 제품과 기술로 매핑됩니다. 본 절에서는 각 요소의 정의와 역할을 고정하여, 이후 기술 비교와 아키텍처 설계의 언어 기반을 마련합니다. 이 네 요소는 Agent의 기능적 범위와 확장성, 유지보수 용이성에 직접적인 영향을 미치며, 각 요소의 설계와 구현 방식에 따라 전체 시스템의 품질과 안정성이 좌우됩니다. 따라서, 각 요소의 역할과 상호작용을 명확히 이해하는 것이 중요합니다.

### 2.2.1 Tool Use — function calling 과 MCP 가 정의하는 “Agent 의 팔다리”

#### Tool Use의 개념과 진화

Tool Use는 Agent가 외부 시스템을 호출해 액션 공간을 확장하는 행위입니다. 초창기에는 OpenAI의 function calling이 사실상 표준이었으나, 2024-11에 공개된 MCP(Model Context Protocol)가 프레임워크 간 공통 규약으로 자리잡았습니다. MCP는 JSON-RPC 기반 client-server 구조로, tools/list, tools/call, resources 등 핵심 기능을 제공합니다. Tool Use의 발전은 Agent가 단순히 텍스트 생성에 그치지 않고, 실제 업무 자동화와 시스템 통합의 중심 역할을 하게 만들었습니다. 예를 들어, 금융 분야에서는 Agent가 실시간 환율 조회 API를 호출하거나, 사내 ERP 시스템과 연동하여 재고 정보를 자동으로 업데이트하는 등, 다양한 Tool Use 시나리오가 현실화되고 있습니다.

#### function calling vs MCP 비교 표

구분	function calling (OpenAI)	MCP (Model Context Protocol)
표준화	OpenAI 독점	OpenAI/Google/MS/Anthropic 공동 지지
프로토콜	Proprietary JSON	JSON-RPC 표준
도구 등록	개별 API	공식 레지스트리/도구 스키마
확장성	제한적	벤더 중립, 확장 가능

MCP의 등장은 다양한 벤더와 프레임워크 간의 상호운용성을 크게 높였으며, 엔터프라이즈 환경에서의 도구 관리와 확장성을 획기적으로 개선하였습니다. MCP를 도입하면, 여러 LLM 및 Agent 프레임워크가 동일한 Tool Registry와 호출 프로토콜을 공유할 수 있어, 도구 추가 및 유지보수가 훨씬 용이해집니다.

### 제품 평가 적용

Agent가 호출 가능한 도구 목록은 제품 평가의 1급 축입니다. 사내 시스템을 MCP 서버로 감싸는 비용과 전략은 Tool Use 설계의 핵심 과제입니다. 실제로, 사내 레거시 시스템을 MCP 호환 서버로 래핑하는 작업은 초기 투자 비용이 들지만, 장기적으로는 다양한 Agent 프레임워크와의 통합 및 확장에 큰 이점을 제공합니다. 따라서, Tool Use 설계 시에는 현재와 미래의 시스템 확장성, 유지보수 편의성, 벤더 락인 위험 등을 종합적으로 고려해야 합니다.

## 2.2.2 Memory — Short-term · Long-term · Episodic · Semantic 4계층 구분

### Memory 계층의 개념과 역할

Agent의 메모리는 단순한 대화 기록을 넘어, 다양한 계층으로 구분됩니다. Short-term 은 대화 버퍼, Long-term은 Vector DB 기반 의미 검색, Episodic은 Reflexion 회고 메모리, Semantic은 Knowledge Graph 기반 관계형 지식 저장소입니다. 각 계층은 저장소, 조회 방식, 삭제 정책이 다르며, Agent의 품질과 지속적 학습에 직접적인 영향을 미칩니다. 예를 들어, Short-term Memory는 최근 대화 맥락을 유지하여 일관된 응답을 제공하고, Long-term Memory는 방대한 문서나 기록에서 필요한 정보를 검색합니다. Episodic Memory는 Agent가 과거의 경험을 회고하여 유사 상황에서 더 나은 결정을 내릴 수 있게 하며, Semantic Memory는 복잡한 개체 간의 관계와 도메인 지식을 구조적으로 저장하여 고차원 reasoning을 지원합니다.

### 메모리 4계층 ↔ 저장소 기술 매핑 표

계층	주요 기능	대표 저장소/기술
Short-term	대화 버퍼	Redis, in-memory buffer
Long-term	의미 검색	Pinecone, Weaviate, Qdrant
Episodic	회고 기록	Reflexion, Neo4j Memory Provider
Semantic	관계형 지식	Neo4j, Knowledge Graph

각 계층별로 적합한 저장소와 기술을 선정하는 것이 Agent의 성능과 확장성에 큰 영향을 미칩니다. 예를 들어, 대규모 문서 검색이 필요한 경우에는 Long-term Memory에 특화된 Vector DB를, 복잡한 관계형 지식 탐색이 필요한 경우에는 Semantic Memory에 Knowledge Graph를 도입하는 것이 효과적입니다.

### 설계 리뷰 적용

메모리 설계는 Vector/GraphDB 선정의 사전 프레임이 됩니다. “우리 Agent는 어느 계층 메모리를 유지하는가?”는 RFP 항목으로 고정해야 합니다. 또한, 각 계층의 데이터 보존 기간, 접근 속도, 보안 정책 등을 명확히 정의함으로써, 엔터프라이즈 환경에서의 데이터 관리와 컴플라이언스 요구사항도 충족할 수 있습니다. 메모리 계층의 명확한 구분과 설계는 Agent의 지속적 학습과 품질 개선에 핵심적인 역할을 합니다.

## 2.2.3 Harness · Orchestration · Handoff · Guardrail · HITL 의 정확한 경계

### Harness의 정의와 역할

Harness는 Agent의 워크플로우, 제약, 피드백, 튜체인, 라이프사이클을 관리하는 상위 시스템입니다. Karpathy 이후 Prompt vs Context vs Harness의 3단 구분이 표준으로 자리잡았습니다. Harness는 프롬프트/컨텍스트를 넘어, 전체 워크플로우와 제약 조건, 피드백 루프, 튜연동, 라이프사이클 관리까지 담당합니다. 예를 들어, Harness는 Agent의 입력/출력 검증, 외부 시스템과의 연동, 작업 흐름의 자동화, 오류 처리 및 복구 전략 등을 종합적으로 관리합니다. 이를 통해 Agent가 복잡한 업무 환경에서도 안정적으로 동작할 수 있도록 지원합니다.

### Orchestration, Handoff, Guardrail, HITL의 구분

- **Orchestration**: 다중 Agent/태스크 간 라우팅, 의존성 관리, 실패 복구 담당. LangGraph, Temporal 등이 대표.

- **Handoff**: OpenAI SDK에서 제어권 이양(Agent 간 baton passing) 개념. 작업 분리와 책임 분배에 활용.
- **Guardrail**: 병렬 검증 체크포인트. 입력/출력 검증, 오류 탐지, 품질 보장에 사용.
- **HITL(Human-in-the-Loop)**: 사람 승인 루프. 특정 단계나 툴 호출에 대해 인간 개입을 요구하여 신뢰성과 안전성을 확보.

각 용어는 Agent 시스템의 신뢰성, 확장성, 품질 보증에 핵심적인 역할을 하며, 실제 엔터프라이즈 환경에서는 이들 요소의 설계와 구현이 시스템의 성공 여부를 좌우할 수 있습니다. 예를 들어, Orchestration을 통해 복수의 Agent가 협력하여 복잡한 프로세스를 처리하고, Guardrail을 통해 입력 오류나 비정상 출력을 사전에 차단할 수 있습니다. HITL은 민감한 업무나 품질이 중요한 단계에서 인간의 최종 승인을 거치게 하여, 자동화의 한계를 보완합니다.

### 경계 표

용어	주요 기능/역할	대표 제품/기술
Harness	워크플로우·제약·피드백·툴체인·라이프사이클 관리	LangGraph, Claude Code
Orchestration	다중 Agent/태스크 라우팅·의존성·실패복구	Temporal, LangGraph
Handoff	제어권 이양, 작업 분리	OpenAI SDK
Guardrail	병렬 검증 체크포인트	OpenAI SDK, LangGraph
HITL	사람 승인 루프	LangGraph, CrewAI

### 용어 사전 적용

팀 용어 사전에 이 5개를 한 줄 정의로 박아 넣으면, 프레임워크 문서에서 단어 혼용을 바로 교정할 수 있습니다. 또한, 각 용어의 정확한 정의와 경계 구분은 신규 인력 온보딩, 시스템 유지 보수, 외부 벤더와의 협업 등 다양한 상황에서 커뮤니케이션 오류를 줄이고, 일관된 시스템 설계를 가능하게 합니다.

## 2.3 CoT · CoT-SC · ToT · GoT · PDDL — 경로 복잡도 축으로 정렬한 추론 기법들

Agent의 Planning 모듈을 깊게 다루기 전에, CoT, CoT-SC, ToT, GoT, PDDL 등 주요 추론 기법의 이름과 경로 복잡도 축을 한 줄로 정리합니다. 이 절은 단일 경로에서 다중 경로, 그래프 탐색, 외부 플래너까지의 계열화를 명확히 하여, 기술 비교와 적용 판단의 기준을 제공합니다. 각 기법의 특징과 한계, 적용 사례를 이해함으로써, 실제 Agent 설계 시 적합한 추론 방식을 선택할 수 있습니다. 또한, 경로 복잡도와 정확도, 비용 간의 트레이드오프를 명확히 인식하는 것이 중요합니다.

### 2.3.1 CoT 와 Zero-shot CoT — “Let’s think step by step” 의 힘과 한계

#### CoT(Chain-of-Thought) 개념과 효과

CoT는 Wei et al. 2022가 제안한 단일 경로 추론 기법으로, 문제를 풀이 단계로 분해하여 LLM의 reasoning 능력을 강화합니다. Zero-shot CoT는 Kojima et al. 2022가 “Let’s think step by step”이라는 한 문장으로 예시 없이도 CoT 효과를 얻는 방법을 제시했습니다. CoT는 수학 문제, 논리 퍼즐 등에서 LLM이 복잡한 reasoning을 수행할 수 있도록 도와주며, 실제로 다양한 벤치마크에서 성능 향상을 입증하였습니다. Zero-shot CoT는 별도의 예시 없이도 간단한 프롬프트 추가만으로 CoT의 효과를 얻을 수 있어, 프롬프트 설계의 부담을 크게 줄여줍니다.

#### CoT vs Zero-shot CoT 요약 표

구분	특징	대표 프롬프트/효과
CoT	예시 기반 단계별 추론	“문제를 단계별로 풀어보자”
Zero-shot CoT	프롬프트 한 문장으로 추론	“Let’s think step by step”

#### 한계와 실패 모드

단일 경로 추론은 이른 수렴, 단일 실패 지점(한 번 틀리면 전체 실패) 등의 한계를 가집니다. 예를 들어, 복잡한 문제에서 한 단계라도 잘못 추론하면 전체 답변이 틀릴 수 있으며, 다양한 해석이 가능한 문제에서는 일관된 답변을 보장하기 어렵습니다. “CoT로 충분한 문제”와 “그렇지 않은 문제”를 분리 기준으로 제시해야 합니다. 실제로, 단순 산술 문제나 명확한 논리 구조를 가진 문제에는 CoT가 효과적이지만, 다단계 reasoning이나 다양한 경로가 존재하는 문제에서는 한계가

명확히 드러냅니다. 따라서, CoT와 Zero-shot CoT는 단일 경로 추론의 한계를 인식하고, 필요에 따라 다중 경로 추론 기법과 병행하는 것이 바람직합니다.

## 2.3.2 CoT-SC · ToT · GoT — 경로가 여러 개일 때 일어나는 일

### 다중 경로 추론의 개념과 효과

CoT-SC(Self-Consistency)는 여러 경로를 샘플링한 후 다수결로 최종 답을 고르는 방식입니다. GSM8K +17.9%, SVAMP +11.0% 등 정확도 개선을 수치로 보고했습니다. ToT(Tree of Thoughts, NeurIPS 2023)는 중간 상태를 노드로 두고 자체 평가로 유망한 분기만 확장하는 트리 탐색 구조입니다. GoT(Graph of Thoughts, AAAI 2024)는 임의 그래프를 통해 의존성 모델링과 병합까지 확장합니다. 이들 기법은 복잡한 reasoning이 필요한 문제에서 LLM의 성능을 크게 향상시킬 수 있으며, 실제로 다양한 벤치마크에서 의미 있는 정확도 개선을 입증하였습니다.

예를 들어, CoT-SC는 동일한 문제에 대해 여러 번 reasoning을 수행하고, 그 결과 중 가장 빈도가 높은 답을 선택함으로써, LLM의 불확실성이나 일시적 오류를 효과적으로 보완합니다. ToT는 각 reasoning 단계에서 여러 분기(생각의 가지)를 생성하고, 자체 평가 메커니즘을 통해 유망한 경로만을 확장함으로써, 탐색 효율성과 정확도를 동시에 높입니다. GoT는 복잡한 문제에서 reasoning 경로가 트리 구조를 넘어서 그래프 형태로 얽혀 있을 때, 다양한 경로 간의 의존성과 정보를 통합하여 최적의 답을 도출할 수 있게 해줍니다.

### CoT-SC/ToT/GoT 경로 복잡도 대비표

기법	경로 구조	정확도 개선	비용/LLM 호출 횟수
CoT-SC	다수결	GSM8K +17.9%	N배(샘플 수 만큼)
ToT	트리 탐색	문제별 개선	분기마다 평가 필요
GoT	그래프	복잡 문제 개선	그래프 크기만큼

### 트레이드오프

다중 경로 탐색은 정확도 개선과 LLM 호출 횟수 증가(비용 상승)의 트레이드오프가 핵심입니다. 예를 들어, CoT-SC를 적용하면 GSM8K 벤치마크에서 +17.9%의 정확도 향상을 얻을 수 있지만, 그만큼 LLM 호출 횟수가 N배(샘플 수)에 달해 비용이 크게 증가합니다. ToT와 GoT 역시 탐색 경로가 늘어날수록 LLM 호출 비용이 기하급수적으로 증가할 수 있으므로, 실제 적용 시에는 정확도

개선 효과와 비용 상승을 함께 고려해야 합니다. 엔터프라이즈 환경에서는 예산, 응답 속도, 시스템 자원 등의 제약을 종합적으로 평가하여, 다중 경로 추론 기법의 적용 범위와 전략을 결정하는 것이 바람직합니다.

### 2.3.3 PDDL — LLM 과 고전 심볼릭 플래너가 만나는 외부 플래너 경로

#### PDDL의 정의와 역할

PDDL(Planning Domain Definition Language)은 고전 AI의 심볼릭 플래너 입력 포맷으로, 도메인, 문제, 액션, 전제조건, 효과를 형식적으로 기술합니다. LLM+P 패턴에서는 LLM이 PDDL을 생성하고, 외부 플래너가 해결하며, 결과를 LLM이 자연어로 번역하는 루프를 구성합니다. 이 방식은 복잡한 계획 문제를 체계적으로 해결할 수 있게 해주며, LLM의 창의적 reasoning과 심볼릭 플래너의 엄격한 경로 계산을 결합하는 장점이 있습니다. 예를 들어, 제조 공정 자동화나 물류 최적화와 같이, 계획의 정확성과 강건성이 중요한 분야에서 PDDL 기반 외부 플래너 연동이 활발히 활용되고 있습니다.

#### LLM+P 루프 도식

LLM (PDDL 생성) → 외부 플래너 (계획 계산) → LLM (결과 번역)

이 구조에서 LLM은 자연어 입력을 받아 PDDL 포맷의 계획 문제를 생성하고, 심볼릭 플래너는 해당 계획을 계산하여 최적의 경로를 산출합니다. 이후 LLM이 결과를 자연어로 해석하여 사용자에게 제공합니다.

#### 기술적 탈출구

이 방식은 hallucination(환각) 문제를 강하게 억제하며, 순수 LLM으로 못 푸는 계획 문제에 대한 실제 기술적 탈출구를 제공합니다. 제조, 물류 등 계획 강건성이 중요한 현장에 우선 검토를 권장합니다. 예를 들어, 생산 라인의 작업 순서 최적화, 배송 경로 최적화 등에서는 LLM의 유연한 reasoning과 PDDL 기반 심볼릭 플래너의 정확한 경로 계산을 결합함으로써, 기존 방식 대비 높은 신뢰성과 효율성을 확보할 수 있습니다. 또한, 외부 플래너와의 연동을 통해 Agent 시스템의 확장성과 유지보수성도 크게 향상됩니다.

## 3장: Planning 모듈 — LLM 이 미래 행동을 설계하는 6가지 패러다임

### 3.1 피드백 없음 — 단일 경로(CoT)와 다중 경로(CoT-SC · ToT · GoT)

LLM 기반 AI Agent의 Planning 모듈은 가장 단순한 형태부터 시작하여 점차 경로 복잡도를 확장하는 방향으로 진화해 왔습니다. 이 절에서는 피드백이 없는 단일 경로(CoT), 다수결 기반(CoT-SC), 그리고 트리/그래프 기반 탐색(ToT, GoT) 패러다임을 논문 계보와 함께 정리합니다. 각 기법은 정확도와 비용의 트레이드오프를 명확히 보여주며, 실제 Agent 설계 시 선택 기준이 됩니다. 단일 경로는 효율성과 간결함이 강점이지만, 복잡한 문제에서는 다중 경로 탐색이 필요하며 이는 비용 증가로 이어집니다. 트리와 그래프 탐색은 더욱 복잡한 문제 해결에 적합하지만, 적용 전에 구조적 필요성을 반드시 검토해야 합니다.

#### 3.1.1 CoT · Zero-shot CoT — 단일 경로 추론의 기준선

Chain-of-Thought(CoT)와 Zero-shot CoT는 LLM 기반 추론의 가장 기본적인 패러다임으로, 단일 경로를 따라 문제를 단계적으로 해결하는 방식을 의미합니다. 이 방식은 프롬프트 설계가 간단하고, 비용이 낮으며, 대규모 서비스 환경에서 널리 사용되고 있습니다. 단일 경로 추론은 복잡한 분기나 오류 복구가 필요 없는 문제에 적합하지만, 한 번의 실수로 전체 추론이 실패할 수 있다는 한계도 존재합니다. 본 절에서는 CoT와 Zero-shot CoT의 원리, 구조, 실제 적용 사례와 한계점을 구체적으로 살펴봅니다.

##### 단일 경로 추론의 원리와 구조

Chain-of-Thought(CoT) 추론은 LLM이 문제를 단계별로 분해하여 풀이하는 방식입니다. Wei et al.(2022) 논문에서 처음 제안된 CoT는 프롬프트 내에서 “생각의 흐름(chain)”을 명시적으로 유도함으로써, 모델이 복잡한 문제를 한 번에 답하는 대신 여러 단계로 나누어 해결하도록 만듭니다. 예를 들어, 수학 문제를 풀 때 “먼저 x 값을 계산하고, 다음 y 값을 구한다”와 같이 단계별 reasoning을 명시적으로 유도합니다. 이 방식은 단일 경로(한 번의 chain)만을 생성하므로, 실패 지점이 하나로 고정되어 있습니다. CoT는 비용이 낮고 구현이 간단하여 실제 대규모 서비스에서 가장 빈번히 사용되는 기법입니다.

## Zero-shot CoT의 등장과 효과

Zero-shot CoT는 Kojima et al.(2022) 논문에서 “Let’s think step by step” 한 문장만으로 CoT 효과를 얻을 수 있음을 보여줬습니다. 기존에는 여러 예시(few-shot)를 프롬프트에 포함해야 했지만, Zero-shot CoT는 단일 문장 지시로 LLM이 단계별 reasoning을 수행하게 합니다. 이 방식은 프롬프트 설계 부담을 줄이고, 다양한 문제에 빠르게 적용할 수 있는 장점이 있습니다. 실제로 Zero-shot CoT는 복잡한 reasoning이 필요한 문제에서도 상당한 성능 향상을 보였으며, 엔터프라이즈 환경에서 빠른 PoC에 적합합니다.

## CoT 기법 정리표

기법	프롬프트 방식	경로 수	적용 예시	강점	한계
CoT	단계별 예시 포함	1	수학, 논리 문제	간결, 저비용	단일 실패점
Zero-shot CoT	“Let’s think step by step”	1	범용 reasoning	빠른 적용, 저비용	예시 부족시 한계

## 단일 경로 기준선의 의사결정 의미

단일 경로 추론(CoT/Zero-shot CoT)은 가장 싼 비용으로 가장 넓은 문제군을 커버할 수 있기 때문에, 실제 Agent 설계에서 기본값으로 채택됩니다. 그러나 이 방식은 복잡한 분기, 오류 복구, 다양한 조건 탐색이 필요한 문제에서는 한계가 명확합니다. CTO는 단일 경로 기준선을 명확히 이해해야 이후 다중 경로 기법의 증분 효과를 수치로 평가할 수 있습니다.

CoT와 Zero-shot CoT는 LLM 기반 에이전트 설계에서 기본적인 출발점이 되며, 실제로 많은 기업과 연구기관에서 대규모로 채택하고 있습니다. 예를 들어, 수학 문제 풀이, 논리적 추론, 간단한 데이터 변환 작업 등에서 CoT 방식이 표준으로 자리잡고 있습니다. 단일 경로의 한계는 복잡한 문제에서 오류가 발생할 경우 복구가 어렵다는 점이지만, 반대로 단순한 문제에서는 오히려 효율성과 비용 절감 효과가 큼니다. 실제 운영 환경에서는 CoT/Zero-shot CoT를 기준선으로 삼고, 문제의 난이도와 복잡성에 따라 다중 경로 기법(CoT-SC, ToT, GoT 등)으로 점진적으로 확장하는 전략이 일반적입니다. 이처럼 단일 경로 추론은 LLM 기반 Planning의 토대이자, 이후 고도화된 패러다임의 비교 기준점으로 기능합니다.

### 3.1.2 CoT-SC (Self-Consistency) — 다수결이 GSM8K +17.9%, SVAMP +11.0% 를 만드는 메커니즘

CoT-SC(Self-Consistency)는 단일 경로 추론의 한계를 극복하기 위해 고안된 다수결 기반 패러다임입니다. 이 방식은 동일한 문제에 대해 여러 번의 reasoning chain을 생성하고, 각 chain의 결과를 비교하여 가장 빈도 높은 답을 최종 결과로 채택합니다. 이 절에서는 Self-Consistency의 원리, 구조, 벤치마크 성능 개선 수치, 비용-정확도 트레이드오프, 그리고 실제 적용 시 고려해야 할 의사결정 포인트를 구체적으로 설명합니다.

#### Self-Consistency의 원리와 구조

CoT-SC(Self-Consistency)는 Wang et al.(2022) 논문에서 제안된 다수결 기반 추론 기법입니다. 이 방식은 동일한 문제에 대해 여러 Chain-of-Thought 경로를 샘플링하고, 각각의 결과를 비교하여 다수결(majority voting)로 최종 답을 결정합니다. 예를 들어, LLM이 10개의 reasoning chain을 생성하면, 각 chain의 최종 답을 수집하고 가장 많이 나온 답을 선택합니다. 이 메커니즘은 단일 경로의 이른 수렴이나 단일 실패 지점 문제를 극복할 수 있으며, 복잡한 reasoning 문제에서 정확도를 크게 개선합니다.

Self-Consistency의 핵심은 LLM의 불확실성, 즉 동일한 입력에 대해 다양한 reasoning 경로를 생성할 수 있다는 점을 적극적으로 활용한다는 데 있습니다. 각 reasoning chain은 LLM의 샘플링(temperature, top-k 등) 설정에 따라 달라질 수 있으며, 이로 인해 서로 다른 풀이 과정과 결과가 도출됩니다. 다수결 방식은 이러한 다양성을 품질 개선에 활용하는 대표적 방법입니다. 실제로, 단일 CoT와 비교할 때 Self-Consistency는 복잡한 문제(특히 수학, 논리, 멀티스텝 reasoning)에서 현저한 성능 향상을 보입니다.

#### 벤치마크 성능 개선 수치

Self-Consistency는 여러 벤치마크에서 강력한 성능 개선을 보여줍니다. GSM8K에서 +17.9%, SVAMP에서 +11.0%, AQuA에서 +12.2%, StrategyQA에서 +6.4%, ARC-c에서 +3.9%의 정확도 향상이 보고되었습니다. 이는 동일한 모델을 여러 번 호출하는 것만으로도 품질을 크게 개선할 수 있음을 의미합니다.

Self-Consistency의 효과는 특히 reasoning 경로가 다양하게 나올 수 있는 문제에서 극대화됩니다. 예를 들어, GSM8K(초등 수학 문제)에서 단일 CoT 대비 17.9%p의 정확도 향상이 관찰되었으며, SVAMP(수학적 추론)에서도 11.0%p 개선이 있었습니다. 이러한 수치는 Self-

Consistency가 단순한 반복 호출이 아니라, LLM의 잠재적 reasoning 다양성을 체계적으로 활용하는 방법임을 보여줍니다.

### 비용과 정확도의 트레이드오프

다수결 방식은 정확도 향상이라는 명확한 이점을 제공하지만, 비용이 N배로 증가한다는 트레이드오프가 있습니다. 예를 들어, 10회 샘플링을 하면 단일 호출 대비 비용이 10배가 됩니다. 실제 운영에서는 샘플링 횟수(N)를 조정하여 비용 감응도를 분석해야 하며, 대량 배포 시에는 비용 구조가 매우 중요한 의사결정 포인트가 됩니다.

이러한 비용-정확도 트레이드오프는 실제 서비스 환경에서 매우 중요한 고려사항입니다. 예를 들어, 엔터프라이즈 환경에서는 품질이 중요한 핵심 경로에만 Self-Consistency를 적용하고, 일반 경로에는 단일 CoT를 유지하는 하이브리드 전략이 권장됩니다. 또한, 샘플링 횟수(N)를 동적으로 조정하여, 예산과 품질 목표에 맞는 최적의 운영 방안을 도출할 수 있습니다.

### Self-Consistency 벤치마크 성능 표

벤치마크	단일 CoT 정확도	CoT-SC 정확도	개선폭
GSM8K	55.0%	72.9%	+17.9%
SVAMP	63.5%	74.5%	+11.0%
AQuA	53.2%	65.4%	+12.2%
StrategyQA	52.0%	58.4%	+6.4%
ARC-c	36.7%	40.6%	+3.9%

### 다수결 N값과 비용 감응도

Self-Consistency는 “같은 모델로 정확도를 올리는 최단경로”이지만, 비용이 샘플링 횟수만큼 늘어납니다. CTO는 다수결 N값에 대한 비용 감응도를 반드시 분석해야 하며, 실제 운영에서는 품질과 비용의 균형점을 찾아야 합니다.

실제 현업에서는 Self-Consistency의 N값을 실험적으로 조정하여, 품질 개선의 한계점(수렴점)을 찾고, 그 이후에는 추가 비용 대비 품질 개선이 미미해지는 구간에서 최적의 N값을 결정합니다. 예를 들어, N=5, 10, 20 등으로 실험하여, 품질 곡선이 평탄해지는 지점을 기준으로 운영에 적용할 수 있습니다. 또한, Self-Consistency는 LLM의 버전이나 도메인에 따라 효과가 달라질 수 있으므로, 사전 테스트와 벤치마크가 필수적입니다. 이처럼 CoT-SC는 단일 경로 추론의 한계를

극복하고, LLM의 reasoning 다양성을 최대한 활용하는 대표적 패러다임입니다.

### 3.1.3 ToT · GoT — 트리와 그래프가 열어주는 경로 탐색 공간

ToT(Tree of Thoughts)와 GoT(Graph of Thoughts)는 LLM 기반 Planning에서 경로 탐색의 복잡도를 극적으로 확장하는 패러다임입니다. 이 절에서는 트리와 그래프 구조의 원리, 상태 평가 및 병합 메커니즘, 실제 적용 사례, 그리고 CTO가 고려해야 할 의사결정 포인트를 상세히 다룹니다. 트리/그래프 기반 탐색은 복잡한 문제 해결에 강점을 보이지만, 비용과 구현 난이도가 크게 증가하므로, 실제 적용 전 구조적 필요성을 신중히 검토해야 합니다.

#### Tree of Thoughts(ToT)의 구조와 적용

Tree of Thoughts(ToT)는 Yao et al.(NeurIPS 2023) 논문에서 제안된 트리 기반 탐색 기법입니다. ToT는 문제 해결 과정에서 중간 상태를 노드로 두고, 각 분기점에서 자체 평가를 통해 유망한 분기만 확장합니다. 예를 들어, 퍼즐이나 게임 문제에서 여러 중간 상태를 생성하고, 각 상태의 가치(utility)를 평가하여 가장 유망한 경로만 계속 확장합니다. 이 방식은 분기와 상태 탐색이 필요한 복잡한 문제에 매우 효과적입니다.

ToT의 핵심은 “상태의 분기와 평가”입니다. 각 단계에서 여러 후보 상태를 생성하고, LLM 또는 외부 평가자가 각 상태의 가치를 점수화합니다. 이후, 가장 점수가 높은 상태만 다음 단계로 확장하여, 탐색 공간을 효율적으로 줄이면서도 다양한 경로를 고려할 수 있습니다. 이 방식은 단순한 직선형 reasoning(CoT, CoT-SC)과 달리, 문제의 다양한 해법을 동시에 탐색할 수 있다는 장점이 있습니다.

#### Graph of Thoughts(GoT)의 구조와 확장성

Graph of Thoughts(GoT)는 Besta et al.(AAAI 2024) 논문에서 제안된 그래프 기반 탐색 기법입니다. GoT는 트리 구조를 넘어 임의 그래프 형태로 상태 간 의존성과 병합까지 모델링합니다. 예를 들어, 여러 경로가 동일한 중간 상태로 합쳐질 수 있으며, 상태 간 관계를 그래프 형태로 표현함으로써 더 복잡한 문제 구조를 다룰 수 있습니다.

GoT의 가장 큰 특징은 “상태 병합”과 “비선형 경로”입니다. 트리 구조에서는 각 분기가 독립적으로 확장되지만, 그래프 구조에서는 서로 다른 경로가 동일한 상태로 수렴하거나, 상태 간의 다양한 의존관계가 표현될 수 있습니다. 이로써, 복잡한 연구 문제, 멀티에이전트 협력, 대규모 계획 등에서 GoT가 강점을 발휘합니다.

## CoT/CoT-SC/ToT/GoT 탐색 구조 시각화

기법	탐색 구조	상태 평가	병합 가능성	적용 예시
CoT	직선	없음	없음	수학, 논리
CoT-SC	직선×N	없음	없음	수학, 논리
ToT	트리	자체 평가	없음	게임, 퍼즐
GoT	그래프	자체 평가	있음	연구, 복잡 문제

### 적용 판단 기준과 의사결정 포인트

트리/그래프 기반 탐색은 복잡 게임, 퍼즐, 연구성 태스크에 적용할 때 강점을 보입니다. 그러나 구조가 복잡해질수록 비용과 구현 난이도가 크게 증가하므로, CTO는 “트리/그래프 구조가 정말 필요한가?” 를 먼저 묻고, 실제로 경로 탐색이 문제 해결에 필수적인지 판단해야 합니다.

ToT와 GoT는 LLM 기반 Planning의 확장성을 극대화하지만, 실제 운영에서는 계산 비용, 구현 복잡성, 실시간성 요구 등 다양한 현실적 제약을 반드시 고려해야 합니다. 예를 들어, 대규모 게임 AI, 과학적 연구 설계, 멀티에이전트 협력 등에서 트리/그래프 기반 탐색이 실질적 성과를 내고 있지만, 단순한 업무 자동화나 일상적 reasoning에는 과도한 오버헤드가 될 수 있습니다. 따라서, CTO는 문제의 특성과 요구사항을 면밀히 분석하여, 트리/그래프 기반 탐색의 도입 여부를 신중히 결정해야 합니다.

## 3.2 외부 플래너 — PDDL 과 Planning Copilot 계보

Planning 모듈의 두 번째 그룹은 LLM의 순수 추론이 아닌, 고전 AI의 심볼릭 플래너와 결합하는 패러다임입니다. PDDL(Planning Domain Definition Language)은 고전 AI에서 계획 문제를 형식적으로 기술하는 언어로, LLM이 PDDL을 생성하고 외부 플래너가 이를 해결하는 구조가 최근 Agent 설계에서 각광받고 있습니다. 이 절에서는 PDDL의 기본 구조, LLM+P 패턴, 그리고 MCP 기반 Planning Copilot 계보를 정리합니다. 이러한 외부 플래너 결합은 LLM 단독의 hallucination 문제를 극복하고, 완전성과 최적성을 보장하는 강력한 설계 축이 됩니다.

### 3.2.1 PDDL — 고전 AI 가 LLM 에 주는 형식언어

PDDL(Planning Domain Definition Language)은 고전 AI에서 계획 문제를 기술하기 위해 개발된 표준 형식언어입니다. 이 언어는 복잡한 계획 문제를 명확하게 정의하고, 컴퓨터가 자동으로 해결할 수 있도록 설계되었습니다. 본 절에서는 PDDL의 구조, 핵심 요소, 실제 적용 사례, 그리고 LLM 단독 추론의 한계와 PDDL 기반 플래너의 필요성을 구체적으로 설명합니다.

#### PDDL의 구조와 핵심 요소

PDDL(Planning Domain Definition Language)은 고전 AI에서 계획 문제를 기술하기 위한 표준 형식언어입니다. PDDL은 도메인(domain), 문제(problem), 액션(action), 전제조건(precondition), 효과(effect) 등으로 구성되어 있으며, 각 요소를 명확하게 정의함으로써 계획 문제를 formal하게 기술할 수 있습니다. 예를 들어, 물류 스케줄링에서 “이동”, “적재”, “하차” 등의 액션과 그에 따른 전제조건 및 효과를 명시적으로 기술합니다.

PDDL의 도메인(domain) 파일은 문제를 해결하기 위한 가능한 행동(action)과 그 행동의 전제조건, 효과를 정의합니다. 문제(problem) 파일은 현재 상태와 목표 상태를 기술합니다. 이러한 분리 구조는 다양한 문제에 동일한 도메인 정의를 재사용할 수 있게 하며, 계획의 일반화와 확장성을 높여줍니다. 실제로, 제조, 물류, 로봇 제어 등에서 PDDL은 복잡한 계획 문제를 체계적으로 해결하는 데 필수적인 역할을 하고 있습니다.

#### 클래식 플래너의 역할과 완전성

PDDL로 기술된 계획 문제는 classical planner(Fast Downward, LAMA 등)에 의해 해결됩니다. 이 플래너들은 최적성(optimality)과 완전성(completeness)을 보장하며, LLM의 추론 결과와는 달리 논리적으로 검증된 해답을 제공합니다. 실제로 제조, 물류, 스케줄링 등 완전성이 중요한 도메인에서는 PDDL 기반 플래너가 필수적입니다.

클래식 플래너는 입력된 PDDL 파일을 분석하여, 가능한 행동 시퀀스 중 목표 상태에 도달할 수 있는 최적의 계획을 자동으로 생성합니다. 이 과정에서 플래너는 모든 가능한 경로를 체계적으로 탐색하며, 논리적 오류나 불가능한 계획을 걸러냅니다. 이로써, LLM의 직관적 추론이 가지는 불확실성, hallucination 문제를 근본적으로 보완할 수 있습니다.

#### PDDL 구성 요소 요약표

요소	설명	예시
Domain	문제의 전체 영역 정의	물류, 제조, 로봇
Problem	구체적 상황/목표 정의	A에서 B로 이동
Action	수행 가능한 행동	이동, 적재, 하차
Precondition	액션 수행 전 조건	차량이 비어 있어야 함
Effect	액션 수행 후 변화	차량에 화물이 실림

### LLM 단독 한계와 구조적 대안

LLM 단독으로는 복잡한 계획 문제(완전성, 최적성 요구)를 풀기 어렵기 때문에, CTO는 반드시 “LLM만으로 못 푸는 계획 문제”에 대한 구조적 대안을 PDDL 기반 플래너에서 찾아야 합니다.

실제로, LLM은 자연어 이해와 추론에 강점을 가지지만, 복잡한 상태 공간 탐색, 최적 경로 산출, 논리적 완전성 보장 등에서는 한계가 있습니다. 따라서, 엔터프라이즈 환경에서 완전성과 신뢰성이 요구되는 업무(예: 생산 계획, 물류 최적화, 로봇 경로 설계 등)에는 반드시 PDDL 기반의 심볼릭 플래너를 결합하는 것이 바람직합니다. 이처럼 PDDL은 LLM 기반 Planning의 한계를 보완하는 핵심 구조적 대안으로 자리잡고 있습니다.

### 3.2.2 LLM+P — LLM 생성 → 외부 플래너 → 자연어 번역 루프

LLM+P 패턴은 LLM의 자연어 처리 능력과 고전 AI 플래너의 논리적 완전성을 결합한 대표적 구조입니다. 이 절에서는 LLM+P의 동작 방식, hallucination 상쇄 메커니즘, 3단 루프 도식, 그리고 실제 적용 시의 위험 저감 효과를 상세히 설명합니다.

#### LLM+P 패턴의 구조와 동작 방식

LLM+P(Liu et al. 2023)는 LLM을 PDDL 생성기로 활용하고, classical planner에 문제 해결을 위임한 후, 결과를 다시 LLM이 자연어로 번역하는 3단 루프입니다. 예를 들어, 사용자가 “A에서 B로 화물을 운송하라”고 지시하면, LLM이 PDDL 포맷으로 계획을 생성하고, classical planner가 이를 해결한 뒤, 결과를 LLM이 자연어로 설명합니다.

이 구조는 LLM의 자연어 이해 및 변환 능력과, classical planner의 논리적 계획 수립 능력을 결합함으로써, 두 기술의 장점을 극대화합니다. LLM은 사용자의 요구를 정확히 파악하여 PDDL로 변환하고, 플래너는 논리적으로 검증된 실행 계획을 산출합니다. 마지막 단계에서 LLM은 플래너의 출력을 다시 자연어로 해석하여 사용자에게 제공합니다.

### hallucination 상쇄와 sound 추론

이 패턴의 핵심은 LLM의 hallucination(비논리적 추론)을 classical planner의 sound(논리적) 추론으로 상쇄하는 설계입니다. 즉, LLM이 생성한 계획이 논리적으로 맞는지 classical planner가 검증하고, 실제 실행 가능한 계획만 반환합니다.

이로써, LLM의 창의적이지만 불확실한 추론 결과가 플래너의 엄격한 논리 검증을 거치게 되어, 전체 시스템의 신뢰성과 안전성이 크게 향상됩니다. 특히, 엔터프라이즈 환경에서는 계획의 완전성, 실행 가능성, 오류 방지 등이 매우 중요하므로, LLM+P 패턴이 강력한 대안이 됩니다.

### LLM+P 3단 루프 도식

[사용자 입력]

↓

[LLM → PDDL 생성]

↓

[classical planner → 계획 해결]

↓

[LLM → 자연어 번역]

### 위험 저감 수단으로서의 의미

LLM+P 패턴은 “LLM이 틀려도 플래너가 맞춘다”는 위험 저감 수단으로, 실제 엔터프라이즈 환경에서 계획 강건성이 중요한 도메인에 우선 적용해야 합니다.

실제 사례로, 물류 최적화, 생산 계획, 로봇 제어 등에서 LLM+P 패턴이 도입되어, 자연어 기반 업무 지시의 자동화와 동시에 논리적 오류 방지, 실행 가능성 검증이 이루어지고 있습니다. 이처럼 LLM+P는 LLM 단독 추론의 한계를 극복하고, 안전하고 신뢰할 수 있는 Planning 시스템을 구축하는 데 핵심적인 역할을 합니다.

## 3.2.3 Planning Copilot — PDDL 계보가 MCP 와 결합하는 최신 흐름

Planning Copilot은 LLM+P 패턴의 최신 진화 계보로, MCP(Model Context Protocol)와의 결합을 통해 외부 플래너를 Agent 스택에 자연스럽게 통합하는 구조를 제공합니다. 이 절에서는 Planning Copilot의 계보, MCP 기반 플래너 호출의 구조적 장점, 계보 타임라인, 그리고 현대적 포지셔닝과 의사결정 의미를 구체적으로 설명합니다.

## Planning Copilot의 계보와 MCP 결합

Planning Copilot(arXiv 2509.12987)은 LLM+P → LaMMA-P → MCP 활용 계보를 보여줍니다. 최신 흐름에서는 classical planner를 MCP(Model Context Protocol) 서버로 감싸, Agent가 MCP 툴 호출 방식으로 자연스럽게 플래너를 사용합니다. 즉, 외부 플래너가 더 이상 특수 기술이 아니라 Agent 스택에 통합된 Tool로 자리잡았습니다.

이 구조는 Agent가 외부 플래너를 호출할 때, 복잡한 통신 프로토콜이나 별도의 연동 작업 없이, MCP 표준 인터페이스를 통해 일관된 방식으로 플래너를 사용할 수 있게 합니다. 이로써, 다양한 플래너(예: Fast Downward, LAMA 등)를 손쉽게 교체하거나, 여러 Agent 프레임워크에서 재사용할 수 있는 유연성이 확보됩니다.

### MCP 기반 플래너 호출의 구조적 장점

MCP 기반 플래너 호출은 표준화된 JSON-RPC 인터페이스를 제공하여, 다양한 Agent 프레임워크에서 일관된 방식으로 외부 플래너를 호출할 수 있게 합니다. 이로써 벤더 락인 없이 다양한 플래너를 Agent에 쉽게 통합할 수 있습니다.

MCP는 플래너 호출, 결과 수신, 오류 처리 등 모든 과정을 표준화된 API로 추상화하여, 개발자와 운영자가 복잡한 내부 구조를 신경쓰지 않고도 플래너를 활용할 수 있게 합니다. 또한, MCP는 보안, 로깅, 모니터링 등 엔터프라이즈 환경에서 요구되는 기능도 쉽게 확장할 수 있는 기반을 제공합니다.

### Planning 계보 타임라인

계보 단계	주요 특징	기술적 변화
LLM+P	LLM이 PDDL 생성, 외부 플래너 해결	3단 루프, sound 추론
LaMMA-P	LLM+P 확장, 플래너 통합 강화	플래너 API 통합
Planning Copilot	MCP 서버로 플래너 감싸기	MCP 표준, Tool화

### 현대적 포지션과 의사결정 의미

외부 플래너는 더 이상 특수 시스템이 아닌 MCP 툴로 포지셔닝되며, CTO는 Agent 스택 구성 시 MCP 기반 플래너 통합을 기본 옵션으로 고려해야 합니다.

실제로, 최신 Agent 프레임워크에서는 MCP 기반 플래너 통합이 표준 기능으로 제공되고 있으며, 이를 통해 다양한 도메인(제조, 물류, 로봇, 서비스 자동화 등)에서 신속하고 안정적으로

Planning 기능을 구현할 수 있습니다. 이처럼 Planning Copilot과 MCP의 결합은 LLM 기반 Planning의 실용성과 확장성을 한 단계 끌어올리는 중요한 진화 방향입니다.

### 3.3 피드백 있음 — 환경(ReAct · Voyager) · 인간(Inner Monologue) · 모델(Reflexion · Self-Refine)

Planning 모듈의 마지막 그룹은 피드백이 있는 패러다임입니다. 환경 피드백, 인간 피드백, 모델 피드백으로 분류되며, 각 방식은 Agent의 자기 수정, 실시간 상호작용, 실패 복구 등 고도화된 기능을 제공합니다. Harness 설계에서 closed-loop와 직접 연결되며, 실제 운영에서 품질 개선과 안정성 확보의 핵심 축이 됩니다.

#### 3.3.1 환경 피드백 — ReAct · Voyager · Ghost in the Minecraft (GITM)

환경 피드백 패러다임은 Agent가 외부 환경의 변화를 실시간으로 관찰하고, 그 결과를 다음 행동 결정에 반영하는 구조를 의미합니다. 이 절에서는 환경 피드백의 구조, 대표적 기법(ReAct, Voyager, GITM)의 사례와 성능, 그리고 실제 적용 시 고려해야 할 의사결정 포인트를 구체적으로 설명합니다. 환경 피드백은 실시간 상호작용, 적응성, 자기 수정 등 고도화된 Agent 설계의 핵심 기반이 됩니다.

##### 환경 피드백의 구조와 핵심 원리

환경 피드백 계열은 Agent가 환경으로부터 받은 Observation을 다음 Thought의 입력으로 활용하는 구조입니다. 예를 들어, 웹 자동화 Agent가 API 호출 결과를 받아 다음 행동을 결정하거나, 게임 Agent가 현재 상태를 관찰하여 다음 스킬을 선택하는 방식입니다. 이 구조는 실시간 상호작용형 Agent 설계의 표준 경로입니다.

환경 피드백의 핵심은 “Observation → Thought” 루프입니다. Agent는 매 단계마다 환경의 상태를 관찰하고, 그 결과를 바탕으로 다음 행동을 결정합니다. 이 과정은 단순한 일회성 추론이 아니라, 환경 변화에 따라 동적으로 계획을 수정하는 closed-loop 구조를 형성합니다. 실제로, 웹 자동화, 게임 AI, 로봇 제어 등 다양한 분야에서 환경 피드백 기반 Agent가 표준으로 자리잡고 있습니다.

##### ReAct, Voyager, GITM의 사례와 성능

- ReAct(Yao et al. 2023 ICLR)는 Thought-Action-Observation 루프를 최초로 제안하여 reasoning과 acting을 교차(interleave)합니다.
- Voyager(Wang et al. 2023)는 Minecraft 환경에서 Agent가 Skill Library를 구축하고, unique item 3.3배, tech tree 15.3배 빠른 평생학습을 증명했습니다.
- Ghost in the Minecraft(GITM, Zhu et al. 2023)는 환경 상태를 실시간으로 반영하는 Agent 구조를 보여줍니다.

이러한 기법들은 환경의 변화에 민감하게 반응하고, 실시간으로 행동을 수정할 수 있는 능력을 Agent에 부여합니다. 예를 들어, Voyager는 Minecraft 게임 내에서 Agent가 새로운 기술을 빠르게 습득하고, 다양한 아이템을 효율적으로 획득하는 데 성공했습니다. GITM은 환경 상태의 실시간 반영을 통해, 복잡한 게임 환경에서도 높은 적응성과 성능을 보여주었습니다.

### 환경 피드백 3종 기법 대비표

기법	환경 피드백 방식	성능 지표	적용 분야
ReAct	Observation → Thought	reasoning+acting	웹 자동화, 게임
Voyager	Skill Library+Observation	unique item 3.3x	평생학습, 게임
GITM	실시간 상태 반영	환경 적응성 강화	게임, 로봇

### 실제 적용 대상을 위한 의사결정

환경 피드백은 “환경 = 사내 시스템 상태”로 사상할 수 있으며, 실시간 상호작용이 필요한 Agent 설계에서 반드시 고려해야 합니다.

실제 엔터프라이즈 환경에서는 환경 피드백 구조를 사내 시스템의 상태 변화, API 응답, 사용자 입력 등으로 확장하여 적용할 수 있습니다. 예를 들어, 업무 자동화 Agent가 외부 시스템의 상태를 실시간으로 모니터링하고, 변화에 따라 자동으로 계획을 수정하는 구조를 구현할 수 있습니다. 이처럼 환경 피드백 기반 설계는 Agent의 적응성, 안정성, 실시간성 확보에 필수적인 요소입니다.

### 3.3.2 인간 피드백 — Inner Monologue 계열과 HITL 설계의 이론적 뿌리

인간 피드백 패러다임은 Agent가 인간의 자연어 피드백이나 승인을 실시간으로 받아, 계획을 수정하거나 실행을 중단하는 구조를 의미합니다. 이 절에서는 Inner Monologue와 HITL(Human-in-the-Loop)의 구조, 적용 사례, 책임 구조의 차이, 그리고 실제 운영에서의 의사결정 포인트를

상세히 설명합니다.

### Inner Monologue의 구조와 HITL의 뿌리

Inner Monologue(Huang et al. 2022)는 로봇 작업 중 인간의 자연어 피드백을 LLM의 내부 독백(inner monologue) 형태로 통합하여 실패 복구를 가능하게 합니다. 예를 들어, 로봇이 작업 중 오류를 범하면 인간이 자연어로 피드백을 제공하고, LLM은 이를 내부 상태에 반영하여 계획을 수정합니다.

Inner Monologue 구조는 Agent가 인간의 피드백을 단순한 외부 신호가 아니라, 자신의 내부 reasoning 과정에 통합하는 방식을 의미합니다. 이로써, Agent는 인간의 지시나 수정 요청을 실시간으로 반영하여, 실패를 복구하거나 계획을 동적으로 변경할 수 있습니다. 실제로, 로봇 제어, 자동화 시스템 등에서 Inner Monologue 구조가 도입되어, 인간-에이전트 협력의 효율성과 신뢰성이 크게 향상되고 있습니다.

### HITL(interrupt\_on) 기능의 이론적 근거

Inner Monologue 계열은 현대 Agent 프레임워크의 HITL(Human-in-the-Loop, interrupt\_on) 기능의 이론적 뿌리입니다. HITL은 특정 단계나 툴에서 사람의 승인을 요청하여, 자동화된 Agent의 책임 구조를 명확하게 합니다.

HITL 구조는 Agent가 임의로 모든 결정을 내리지 않고, 중요한 단계(예: 위험한 작업, 비용이 큰 실행 등)에서 인간의 승인을 요청하거나, 중단(interrupt\_on) 신호를 받을 수 있게 합니다. 이로써, Agent의 자동화 수준과 인간의 통제력을 균형 있게 조절할 수 있습니다.

### Inner Monologue ↔ HITL 대응 표

기법	피드백 방식	구조적 특징	적용 분야
Inner Monologue	인간 자연어 피드백	내부 독백 통합	로봇, 자동화
HITL	사람 승인/중단	interrupt_on 노드	모든 Agent 운영

### 책임 구조 차이의 명확화

HITL이 있는 Agent와 없는 Agent는 책임 구조에서 큰 차이를 보입니다. CTO는 HITL이 임시 기능이 아닌 Planning 패러다임으로 위치시켜야 하며, 실제 운영에서 책임 분담을 명확히 해야 합니다.

실제 엔터프라이즈 환경에서는 HITL 구조를 통해, 자동화 Agent의 오작동, 윤리적 문제, 비용

과다 발생 등 다양한 리스크를 효과적으로 관리할 수 있습니다. 예를 들어, 금융, 의료, 제조 등 고위험 도메인에서는 HITL이 필수적인 안전장치로 작동합니다. 이처럼 인간 피드백 기반 설계는 Agent의 신뢰성, 책임성, 통제력 확보에 핵심적인 역할을 합니다.

### 3.3.3 모델 피드백 — Reflexion · Self-Refine · SelfCheck

모델 피드백 패러다임은 Agent가 스스로 결과를 평가하고, 실패를 감지하거나 개선하는 구조를 의미합니다. 이 절에서는 Reflexion, Self-Refine, SelfCheck의 구조와 비교, 자기평가 노드의 프레임워크 적용, 그리고 실제 운영에서의 의사결정 포인트를 상세히 설명합니다. 모델 피드백은 추가 모델 없이 품질을 개선할 수 있는 대표적 레버로, Agent의 자기 수정과 품질 향상에 핵심적인 역할을 합니다.

#### 모델 피드백의 구조와 핵심 원리

모델 피드백 계열은 Agent가 스스로 결과를 평가하고, 실패를 감지하거나 개선하는 구조입니다. Reflexion(Shinn et al. NeurIPS 2023)은 verbal RL로 에피소드 메모리에 회고(reflection)를 저장합니다. Self-Refine(Madaan et al. 2023)은 자기 피드백 후 재작성, SelfCheck는 체인의 오류를 스스로 감지합니다. 이 계열은 추가 모델 없이 품질을 개선하는 주력 레버입니다.

Reflexion은 Agent가 각 에피소드(작업 단위)마다 자신의 행동과 결과를 회고(reflection)하여, 성공/실패 원인을 분석하고, 다음 행동에 반영하는 구조입니다. Self-Refine는 Agent가 자신의 출력을 스스로 평가하고, 부족한 부분을 재작성하여 품질을 개선합니다. SelfCheck는 체인 형태의 reasoning에서 오류를 자동으로 감지하고, 필요한 경우 수정 작업을 수행합니다.

#### Reflexion, Self-Refine, SelfCheck의 비교

- Reflexion: 언어적 강화학습(verbal RL), 에피소드 회고, LangGraph 자기 평가 노드 패턴.
- Self-Refine: 자기 피드백 후 재작성, 품질 개선.
- SelfCheck: 체인 오류 감지, 자동 수정.

이러한 기법들은 Agent가 외부 피드백 없이도 스스로 품질을 개선할 수 있게 하며, 반복적인 자기 평가와 수정 과정을 통해 점진적으로 성능을 향상시킵니다. 실제로, 문서 생성, 코드 생성, 복잡한 reasoning 체인 등 다양한 분야에서 모델 피드백 기반 Agent가 품질과 안정성 측면에서 우수한 성과를 내고 있습니다.

#### 모델 피드백 3기법 비교 표

기법	피드백 방식	구조적 특징	적용 분야
Reflexion	verbal RL, 회고	에피소드 메모리	품질 개선, 실패 복구
Self-Refine	자기 피드백+재작성	자동 수정	문서, 코드 생성
SelfCheck	오류 감지+수정	체인 오류 자동 탐지	복잡 추론 체인

## 자기평가 노드의 프레임워크 적용

자기평가 노드를 기본 탑재한 프레임워크(LangGraph 등)를 우선 검토해야 하며, CTO는 모델 피드백 계열을 “추가 모델 없이 품질 개선”의 핵심 레버로 활용해야 합니다.

실제 엔터프라이즈 환경에서는 LangGraph 등 자기평가 노드가 내장된 프레임워크를 활용하여, Agent의 품질 개선, 실패 복구, 자동 수정 기능을 손쉽게 구현할 수 있습니다. 이처럼 모델 피드백 기반 설계는 추가 비용 없이 Agent의 성능과 안정성을 높이는 효과적인 방법입니다.

## 4장: Action 모듈 — 결정을 출력으로 변환하는 4축 설계

Action 모듈은 AI Agent가 내린 결정을 실제 출력으로 변환하는 핵심 엔진이다. 이 장에서는 Action의 4축(목표, 생성, 공간, 영향)을 체계적으로 정리하고, 각 축이 Agent 설계와 운영에 미치는 기술적 의미를 분석한다. 이 4축은 Agent의 행동을 구조적으로 분류하고, 실패 분석(goal drift, hallucinated tool call 등)의 기준이 된다. 설계 문서에 Action 4축을 고정 섹션으로 포함하면 Agent 간 비교가 균질해지고, 운영 시 발생하는 다양한 장애를 해당 축에 귀속시켜 원인 분석과 개선이 가능하다. 본 장은 최신 연구와 실무 사례를 근거로, Action 모듈의 설계·운영에 필요한 기술적 깊이와 실질적 가이드라인을 제공한다.

### 4.1 액션 목표와 생성 — 무엇을 위해, 어떻게 행동을 만들어내는가

AI Agent의 행동을 효과적으로 설계하고 평가하기 위해서는, 각 행동이 어떤 목표를 지향하는지(무엇을 위해)와 그 행동이 어떤 방식으로 생성되는지(어떻게)에 대한 명확한 구분이 필요합니다. 목표와 생성 축을 분리하여 정의함으로써, Agent의 설계 의도와 실제 동작의 일치 여부를 체계적으로 검증할 수 있으며, 성능 평가와 장애 분석의 기준을 명확히 할 수 있습니다. 본 절에서는 액션 목표의 세 가지 유형과 생성 방식의 두 가지 유형을 표준화하고, 이들의 조합이 Agent 설계에서 갖는 핵심적 의미를 구체적으로 설명합니다. 또한, 각 유형별 KPI와 대표 사례를 표로 정리하여

실무 적용에 도움이 되도록 하였습니다.

#### 4.1.1 액션 목표 3유형 — 작업 완료 · 커뮤니케이션 · 환경 탐색

##### 작업 완료 목표 정의

작업 완료(Task completion)는 Agent가 실제로 외부 시스템에 변화를 주거나, 특정 업무를 마무리하는 행동을 의미한다. 예를 들어 결제 처리, 예약 확정, 코드 머지(Pull Request), 데이터 베이스 업데이트 등은 모두 작업 완료 목표에 해당한다. 이 목표는 Agent가 비즈니스 프로세스의 핵심을 자동화하는 데 필수적이며, KPI(핵심 성과 지표) 역시 “작업 완료율”, “처리 성공률”, “실행 시간” 등으로 정의된다. 동일 Agent 내에서 작업 완료 목표가 명확히 구분되지 않으면, 성과 측정이 모호해지고 운영 리스크가 커진다. 실제 엔터프라이즈 환경에서는 작업 완료 목표를 명확히 구분하여, 각 업무 프로세스별로 Agent의 역할과 책임을 분리하는 것이 중요하다. 예를 들어, 결제 처리 Agent와 단순 정보 조회 Agent의 KPI는 다르게 설정되어야 하며, 이를 통해 각 Agent의 성능과 신뢰성을 체계적으로 관리할 수 있습니다.

##### 커뮤니케이션 목표 정의

커뮤니케이션(Communication)은 Agent가 사용자 또는 다른 시스템과 정보를 주고받는 행동을 의미한다. 예를 들어 답변 생성, 요약, 알림, 보고서 작성 등은 모두 커뮤니케이션 목표에 속한다. 이 목표는 고객 지원, 내부 보고, 실시간 알림 등 다양한 도메인에서 중요하며, KPI는 “응답 정확도”, “대화 만족도”, “정보 전달률” 등으로 설정된다. Agent가 커뮤니케이션 목표를 달성하지 못하면 사용자 경험이 저하되고, 신뢰성에 문제가 발생한다. 특히, 고객 지원 분야에서는 응답의 신속성과 정확성이 매우 중요하며, 커뮤니케이션 목표를 명확히 정의함으로써 Agent의 대화 품질을 지속적으로 개선할 수 있습니다.

##### 환경 탐색 목표 정의

환경 탐색(Environment exploration)은 Agent가 외부 환경을 탐색하거나 정보를 수집하는 행동을 의미한다. 예를 들어 검색, 스크롤, 데이터 수집, API 호출 등은 모두 환경 탐색 목표에 해당한다. 이 목표는 정보 검색, 데이터 분석, 상황 파악 등에서 핵심적이며, KPI는 “탐색 정확도”, “정보 획득률”, “탐색 속도” 등으로 측정된다. Agent가 환경 탐색을 제대로 수행하지 못하면, 이후 행동의 품질이 떨어지고 전체 시스템의 성능이 저하된다. 실제 사례로, 데이터 분석 Agent가 환경 탐색 목표를 효과적으로 달성하지 못하면, 잘못된 데이터 기반으로 의사결정이 이루어질 수

있으므로, 이 목표의 중요성이 강조됩니다.

### 목표 유형별 KPI 예시 표

목표 유형	대표 행동	KPI 예시
작업 완료	결제, 예약, 머지	성공률, 처리시간
커뮤니케이션	답변, 요약, 알림	응답 정확도, 만족도
환경 탐색	검색, 스크롤	탐색 정확도, 획득률

#### 4.1.2 액션 생성 2유형 — Memory recollection 과 Plan-following

액션 생성 방식은 Agent가 행동을 만들어내는 근본적인 메커니즘을 구분하는 중요한 기준입니다. Memory recollection과 Plan-following은 각각 Agent의 학습성과 통제성을 대표하는 방식으로, 실제 시스템 설계에서 이 두 유형을 어떻게 조합하는지가 Agent의 성능과 신뢰성에 큰 영향을 미칩니다. 본 절에서는 각 생성 방식의 정의와 대표 사례, 그리고 엔터프라이즈 환경에서의 적용 시 고려해야 할 기술적 세부사항을 상세히 설명합니다.

##### Memory recollection(기억상상) 정의

Memory recollection은 Agent가 과거 경험, 저장된 메모리, 이전 대화 또는 학습된 지식을 활용해 행동을 생성하는 방식이다. 예를 들어 Voyager Agent는 Minecraft에서 이전에 획득한 스킬이나 경험을 재활용하여 새로운 행동을 결정한다. 이 방식은 Agent의 “학습성”을 높여주며, 반복적 또는 유사한 상황에서 점점 더 효율적으로 행동할 수 있게 한다. Memory recollection이 강한 Agent는 장기적 성능 개선과 자기 수정 능력이 뛰어나다. 실제로, Memory recollection 기반 Agent는 과거 실패 경험을 바탕으로 행동 전략을 점진적으로 개선할 수 있으며, 동일한 실수를 반복하지 않는다는 장점이 있습니다. 또한, 장기적인 목표 달성에 있어 유연하게 전략을 조정할 수 있기 때문에, 복잡한 업무 환경이나 변화가 잦은 도메인에서 특히 유용합니다. 그러나, 과거 경험에 과도하게 의존할 경우 새로운 상황에 대한 적응력이 떨어질 수 있으므로, 메모리 관리와 경험의 선택적 활용이 중요합니다.

##### Plan-following(계획추종) 정의

Plan-following은 Agent가 상위 계획, 미리 정의된 프로세스, 외부 플래너(PDDL 등) 또는 명시적 지침을 따라 행동을 생성하는 방식이다. 예를 들어 ReAct Agent는 Thought-Action-

Observation 루프에서 상위 계획을 단계별로 수행하며, 각 단계에서 외부 피드백을 반영한다. Plan-following은 “통제성”이 강하며, 규제 산업이나 보안이 중요한 환경에서 반드시 요구된다. 이 방식은 예측 가능성과 안정성을 높여주지만, 학습성은 상대적으로 낮을 수 있다. 엔터프라이즈 환경에서는 Plan-following 방식을 통해 업무 프로세스의 표준화와 규정 준수를 보장할 수 있으며, 감사 및 추적이 용이하다는 장점이 있습니다. 반면, 예외 상황이나 비정형 문제에 대한 유연성은 제한될 수 있으므로, Plan-following과 Memory recollection을 적절히 조합하는 하이브리드 설계가 권장되는 경우도 많습니다.

### 생성 유형 ↔ 대표 Agent 사상 표

생성 유형	대표 Agent	특징
Memory recollection	Voyager	학습성, 자기 수정
Plan-following	ReAct	통제성, 예측 가능

#### 4.1.3 목표 × 생성 매트릭스 — Agent 설계 의도의 시각화

액션 목표와 생성 방식의 조합은 Agent 설계에서 매우 중요한 기준점이 됩니다. 각 목표(작업 완료, 커뮤니케이션, 환경 탐색)와 생성 방식(Memory recollection, Plan-following)의 조합은 총 6개의 셀로 구성된 매트릭스를 형성하며, 이를 통해 Agent의 설계 의도와 실제 동작의 일치 여부를 명확히 시각화할 수 있습니다. 이 매트릭스는 신규 Agent 설계 시 필수적으로 작성해야 하며, 각 셀에 해당하는 대표 사례와 KPI를 명시함으로써 설계와 평가의 일관성을 확보할 수 있습니다. 또한, 매트릭스를 통해 Agent 간 비교가 용이해지고, 장애 발생 시 원인 분석의 기준이 명확해집니다.

#### 매트릭스 구성과 의미

액션 목표(3유형)와 생성(2유형)의 조합은 총 6셀 매트릭스를 구성한다. 이 매트릭스는 각 Agent의 설계 의도를 한눈에 시각화하며, 신규 Agent 착수 시 반드시 해당 셀을 명시하도록 강제한다. 예를 들어 “작업 완료 × Plan-following”은 결제 자동화 Agent, “환경 탐색 × Memory recollection”은 데이터 분석 Agent에 해당한다. 이 매트릭스는 설계 문서, 등록 폼, KPI 정의 등에서 표준화되어야 한다. 실제로, 대규모 조직에서는 이 매트릭스를 기반으로 Agent 포트폴리오를 관리하며, 각 Agent의 역할과 책임, 성능 지표를 체계적으로 정리할 수 있습니다.

## 매트릭스 예시 표

Memory recollection	Plan-following	
작업 완료	반복적 업무 자동화	결제/예약/머지
커뮤니케이션	대화 이력 활용	보고서/알림
환경 탐색	경험 기반 탐색	검색 프로세스

### 사례 적용과 설계 강제

Agent 설계 시 이 매트릭스의 해당 셀을 명시하면, 목표와 생성 방식이 혼재되어 KPI 정의가 흔들리는 문제를 예방할 수 있다. 또한, 각 셀에 맞는 대표 사례를 제시하면 신규 Agent의 설계와 평가가 명확해진다. 사내 Agent 등록 폼에 이 매트릭스를 필수 필드로 편입하는 것이 권장된다. 예를 들어, “커뮤니케이션 × Memory recollection” 셀에는 고객 지원 Agent가 과거 대화 이력을 활용해 맞춤형 답변을 제공하는 사례가 포함될 수 있고, “환경 탐색 × Plan-following” 셀에는 정해진 검색 프로세스를 단계별로 수행하는 데이터 수집 Agent가 해당됩니다. 이러한 표준화는 조직 내 Agent 개발 및 운영의 효율성을 크게 높여줍니다.

## 4.2 액션 공간 — 외부 도구(API · DB · 모델)와 내부 지식(계획 · 대화 · 상식)

Agent가 실제로 행동을 수행할 수 있는 공간, 즉 액션 공간은 외부 도구와 내부 지식이라는 두 가지 주요 영역으로 나뉩니다. 외부 도구는 API, 데이터베이스, 타 AI 모델 등 환경과의 직접적인 상호작용을 담당하며, 내부 지식은 Agent 자체가 보유한 계획, 대화 이력, 상식 등 내재된 정보를 기반으로 행동을 결정하는 영역입니다. 이 두 공간의 균형과 통합은 Agent의 기능 확장성과 비용 구조, 그리고 보안 및 통제 측면에서 매우 중요한 의미를 가집니다. 본 절에서는 외부 도구와 내부 지식의 기술적 특징, 엔터프라이즈 환경에서의 임팩트, 그리고 표준 프로토콜(MCP)의 역할을 구체적으로 분석합니다.

### 4.2.1 외부 도구 — API · DB · 다른 모델 호출이 만드는 “팔다리”

외부 도구는 Agent가 실제 환경과 상호작용할 수 있도록 해주는 핵심적인 수단입니다. API 호출, 데이터베이스 쿼리, 다른 AI 모델 호출 등은 Agent의 기능을 확장하고, 자동화의 범위를 넓혀줍니다. 그러나 이러한 외부 도구의 사용은 보안, 비용, 통제 등 다양한 리스크를 동반하므로, 체계적인

관리와 통제가 필수적입니다. 본 절에서는 각 외부 도구 유형의 기술적 특징과 엔터프라이즈 환경에서의 위험도 및 통제 방안을 상세히 설명합니다.

### API 호출의 의미와 범위

API 호출은 Agent가 외부 시스템에 직접 요청을 보내 데이터를 조회하거나, 특정 작업을 수행하는 핵심 액션이다. 예를 들어 결제 API, 예약 시스템, 외부 데이터베이스, SaaS 플랫폼 등은 모두 API 호출을 통해 Agent와 연결된다. API 호출은 Agent의 기능 확장과 자동화의 핵심이지만, 보안·비용 통제의 출발점이기도 하다. 실제 운영 환경에서는 API 호출의 빈도와 범위가 늘어날수록 외부 시스템과의 연동 리스크가 커지며, 인증 및 권한 관리, 호출 이력 추적이 반드시 필요하다. 또한, API 호출 실패 시의 예외 처리와 재시도 로직도 설계에 포함되어야 하며, 외부 시스템 장애가 Agent 전체 서비스에 영향을 미치지 않도록 격리 구조를 마련해야 한다.

### DB 쿼리와 데이터 연동

DB 쿼리는 Agent가 내부 또는 외부 데이터베이스에 접근하여 정보를 조회하거나, 데이터를 업데이트하는 행동이다. 예를 들어 고객 정보 검색, 재고 확인, 로그 기록 등은 모두 DB 쿼리로 처리된다. DB 쿼리는 민감 정보 접근과 데이터 무결성에 영향을 미치므로, tool allowlist(허용 목록)와 권한 통제가 필수적이다. 특히, 개인정보 보호 및 데이터 유출 방지를 위해 DB 접근 권한을 최소화하고, 모든 쿼리 실행 내역을 감사 로그로 남기는 것이 중요하다. 또한, 대용량 데이터 처리 시에는 쿼리 최적화와 성능 모니터링이 필요하며, 데이터 일관성 유지를 위한 트랜잭션 관리도 필수적인 요소입니다.

### 다른 AI 모델 호출

Agent가 다른 AI 모델을 호출하여 복합적인 작업을 수행하는 경우도 많다. 예를 들어 이미지 생성, 음성 인식, 자연어 번역 등은 별도의 AI 모델을 호출하여 처리된다. 이 영역은 기능 확장성이 높지만, 비용과 지연시간, API 파편화 등의 리스크가 존재한다. 실제로, 다양한 AI 모델을 연동할 때는 각 모델의 입력/출력 포맷이 다를 수 있으므로, 표준화된 인터페이스 설계가 필요하다. 또한, 외부 AI 모델 호출 시 발생하는 비용과 응답 지연을 모니터링하고, 필요에 따라 캐싱이나 비동기 처리 전략을 도입하는 것이 바람직합니다. 엔터프라이즈 환경에서는 모델 호출 내역을 통합적으로 관리하여, 비용 예측과 성능 최적화를 지원해야 합니다.

### 외부 도구 카테고리별 위험도 표

도구 유형	대표 예시	위험도	통제 방법
API	결제, 예약	높음	allowlist, 인증
DB 쿼리	고객정보, 로그	중간	권한, 감사 로그
AI 모델	번역, 생성	낮음~중간	비용 모니터링

#### 4.2.2 내부 지식 — 계획 · 대화 · 상식이 행동으로 바뀌는 순간

내부 지식 기반 액션은 Agent가 외부 시스템에 의존하지 않고, 자체적으로 보유한 정보와 논리를 활용해 행동을 결정하는 영역입니다. 이 방식은 비용 효율성과 응답 속도, 그리고 보안 측면에서 강점을 가지며, 엔터프라이즈 환경에서 점점 더 중요해지고 있습니다. 본 절에서는 내부 계획, 대화 이력, 모델 상식 등 다양한 내부 지식의 활용 방식과 그 기술적 세부사항, 그리고 실제 적용 시 고려해야 할 점들을 구체적으로 설명합니다.

##### 계획 기반 행동

내부 계획은 Agent가 미리 정의된 프로세스, 워크플로우, 또는 목표에 따라 행동을 결정하는 영역이다. 예를 들어 “보고서 작성”, “분석 계획 수립”, “회의록 요약” 등은 외부 호출 없이 내부 계획만으로 수행 가능하다. 이 방식은 비용이 낮고, 통신량이 적어 효율적이다. 실제로, 내부 계획 기반 행동은 반복적이고 표준화된 업무에 적합하며, 외부 시스템 장애와 무관하게 안정적인 서비스 제공이 가능합니다. 또한, 내부 계획은 Agent의 행동을 예측 가능하게 만들어, 운영 및 감사 측면에서 유리합니다.

##### 대화 이력 활용

Agent는 지금까지의 대화 이력을 활용해 다음 행동을 결정할 수 있다. 예를 들어 고객 지원 Agent가 이전 대화 내용을 참고하여 추가 질문을 하거나, 맞춤형 답변을 생성하는 경우가 있다. 대화 이력 기반 행동은 사용자 경험을 높이고, 반복적 문의에 빠르게 대응할 수 있다. 특히, 장기 대화나 복잡한 상담 업무에서 대화 이력의 활용은 고객 만족도를 크게 높여줍니다. 기술적으로는 대화 이력의 저장, 검색, 요약 기능이 Agent 시스템에 내장되어야 하며, 개인정보 보호와 데이터 보존 정책도 함께 고려해야 합니다.

##### 모델 상식 활용

Agent는 자체적으로 학습된 상식, 도메인 지식, 일반적인 규칙 등을 활용해 행동을 생성할 수

있다. 예를 들어 FAQ 답변, 기본 정책 안내, 일반적 업무 처리 등은 외부 호출 없이 내부 지식만으로 충분히 수행 가능하다. 내부 지식 기반 행동의 비율이 높을수록 비용 구조가 효율적이다. 실제로, 엔터프라이즈 환경에서는 자주 묻는 질문이나 표준 업무 처리에 내부 상식 기반 액션을 우선 적용하여, 외부 리소스 사용을 최소화하고 시스템의 응답 속도를 높입니다. 또한, 내부 지식의 주기적 업데이트와 품질 관리가 Agent의 장기적 성능 유지에 중요합니다.

### 내부 vs 외부 액션 비중 예시

액션 유형	비중(예시)	비용 영향
내부 지식	60%	저비용
외부 도구	40%	고비용

### 4.2.3 MCP — 외부 도구 공간의 표준이 된 이유와 엔터프라이즈 임팩트

MCP(Model Context Protocol)는 최근 AI Agent 생태계에서 외부 도구 연동의 표준으로 자리 잡은 프로토콜입니다. 다양한 도구와 모델, API를 일관된 방식으로 호출할 수 있도록 설계되어, 엔터프라이즈 환경에서의 확장성과 유지보수성을 크게 향상시켰습니다. 본 절에서는 MCP의 기술적 구조와 엔터프라이즈 채택 현황, 그리고 MCP가 가져온 실질적 변화와 임팩트를 구체적으로 설명합니다.

#### MCP의 기술적 정의와 구조

MCP(Model Context Protocol)는 Anthropic이 2024-11-25에 공개한 JSON-RPC 기반 client-server 표준이다. MCP는 tools/list, tools/call, resources, prompts, sampling 등 핵심 기능을 규정하며, 비동기 작업, stateless 구조, server identity, 공식 레지스트리 등 최신 엔터프라이즈 요구사항을 반영한다. MCP는 월 97M+ SDK 다운로드, OpenAI/Google/MS/Anthropic 공동 지지로 function-calling 파편화를 종결했다. MCP의 구조는 도구 목록 조회, 도구 호출, 리소스 관리, 프롬프트 템플릿, 샘플링 전략 등 다양한 기능을 포함하며, 모든 메시지는 JSON-RPC 포맷으로 통일되어 있습니다. 이를 통해 다양한 벤더의 도구와 AI 모델을 하나의 인터페이스로 연동할 수 있으며, 개발 및 운영의 복잡도를 크게 줄였습니다.

#### MCP의 엔터프라이즈 채택과 표준화

2025-12 MCP는 Linux Foundation 산하 Agentic AI Foundation으로 거버넌스가 이관되

어 벤더 중립 표준 체제로 자리잡았다. MCP 도입은 사내 시스템을 MCP 서버로 감싸는 작업이 Agent 생태계 편입의 1차 관문이 되었다. MCP 기반 도구는 Agent의 “팔다리”로 기능하며, 플랫폼 로드맵에 MCP 도입을 1페이지에 명시하는 것이 표준이 되었다. 실제로, 대형 금융기관이나 글로벌 IT 기업들은 MCP를 기반으로 내부 시스템과 외부 SaaS, 다양한 AI 모델을 통합 관리하고 있습니다. MCP의 표준화는 도구 연동의 일관성, 보안성, 확장성을 동시에 확보할 수 있게 해주며, 신규 도구 도입 시에도 최소한의 개발로 빠르게 연동이 가능합니다. 또한, MCP 기반 대시보드를 통해 도구 사용 현황, 비용, 성능 등을 실시간으로 모니터링할 수 있어, 운영 효율성이 크게 향상되었습니다.

### MCP 대시보드 숫자 표

항목	수치/정보
SDK 다운로드	월 97M+
주요 지지 기업	OpenAI, Google, MS, Anthropic
거버넌스	Linux Foundation Agentic AI
최신 기능	tools/list, tools/call 등

## 4.3 액션의 영향 — 환경 상태 변화 · 내부 상태 변화 · 새 액션 파생과 피드백 루프

Agent의 행동은 단순히 출력으로 끝나지 않는다. Action 이후에는 환경 상태 변화, 내부 상태 변화, 그리고 새로운 액션의 파생과 피드백 루프가 이어진다. 이 절에서는 각 영향 축을 분석하고, Reflexion 스타일 자기 수정이 어떻게 시스템의 품질을 높이는지 기술적으로 설명한다. 실제로, Agent의 행동이 시스템 전체에 미치는 영향과 그로 인한 변화는 운영 안정성, 장애 대응, 장기적 성능 개선에 매우 중요한 요소입니다. 각 영향 축을 명확히 정의하고, 이에 따른 설계 및 운영 가이드라인을 마련함으로써, Agent 시스템의 신뢰성과 확장성을 확보할 수 있습니다.

### 4.3.1 환경 상태 변화 — 부작용(side effect) 과 롤백 설계

Agent가 외부 환경에 미치는 영향은 매우 크며, 특히 쓰기형 액션의 경우 되돌릴 수 없는 부작용(side effect)이 발생할 수 있습니다. 이러한 부작용을 최소화하고, 사고 발생 시 신속하게 대응하기 위해서는 체계적인 롤백 설계와 감사 체계가 필요합니다. 본 절에서는 환경 상태 변화의 정의, 주요

리스크, 그리고 엔터프라이즈 환경에서의 대응 방안을 구체적으로 설명합니다.

### 환경 상태 변화의 정의와 위험

Agent가 API 쓰기, DB 업데이트, 결제 확정 등 환경 상태를 변경하는 행동을 할 때, 되돌릴 수 없는 부작용(side effect)이 발생할 수 있다. 예를 들어 잘못된 결제, 데이터 삭제, 시스템 설정 변경 등은 운영 장애와 직접 연결된다. 이러한 부작용은 롤백이 불가능하거나, 복구 비용이 매우 높을 수 있다. 실제로, 대형 금융기관이나 전자상거래 플랫폼에서는 Agent의 잘못된 액션으로 인한 금전적 손실이나 데이터 손상이 심각한 문제로 이어질 수 있으므로, 환경 상태 변화의 위험을 사전에 평가하고, 예방책을 마련하는 것이 필수적입니다.

### HITL 게이트의 필요성

환경 상태 변화가 발생하는 모든 쓰기형 Tool에는 HITL(Human-in-the-Loop) 게이트를 기본값으로 설정해야 한다. HITL은 사람이 중간에 개입하여 승인 또는 거부를 결정하는 구조로, 사고 예방과 책임 분담에 필수적이다. LangGraph 등 주요 프레임워크는 interrupt\_on 기능으로 HITL을 구현한다. 실제 운영에서는 HITL 게이트를 통해 고위험 액션에 대해 최종 승인을 받고, 예외 상황에서는 즉시 개입할 수 있도록 설계해야 합니다. 또한, HITL 프로세스의 효율성을 높이기 위해 승인 요청과 처리 과정을 자동화하고, 승인 이력을 체계적으로 관리하는 것이 중요합니다.

### 멱등 설계와 감사 로그

멱등(idempotent) 설계는 동일 액션이 여러 번 실행되어도 결과가 변하지 않도록 하는 구조다. 예를 들어 결제 API가 여러 번 호출되어도 중복 결제가 발생하지 않게 한다. 감사 로그(audit log)는 모든 환경 상태 변화를 기록하여, 사고 발생 시 원인 추적과 복구를 가능하게 한다. 엔터프라이즈 환경에서는 멱등 설계와 감사 로그를 기본 정책으로 채택하여, 시스템의 신뢰성과 투명성을 확보합니다. 또한, 감사 로그는 규제 준수와 내부 통제의 핵심 도구로 활용되며, 사고 발생 시 신속한 원인 분석과 책임 소재 파악에 큰 도움이 됩니다.

### 쓰기 Action 리스크 등급 표

액션 유형	리스크 등급	대응 방법
결제/예약	최고	HITL, 멱등, 로그
DB 업데이트	높음	HITL, 로그
조회/검색	낮음	로그

### 4.3.2 내부 상태 변화 — 메모리 업데이트와 계획 재구성

Agent의 행동은 외부 환경뿐만 아니라 내부 상태에도 중요한 변화를 일으킵니다. 메모리 업데이트, 계획 재구성, 점수(확신도) 갱신 등은 Agent의 장기적 성능과 목표 달성에 직접적인 영향을 미치며, 이러한 내부 상태 변화가 체계적으로 관리되지 않으면 goal drift와 같은 심각한 문제가 발생할 수 있습니다. 본 절에서는 내부 상태 변화의 주요 유형과 엔터프라이즈 환경에서의 관리 방안을 상세히 설명합니다.

#### 메모리 업데이트의 의미

Agent의 행동은 내부 상태에도 변화를 준다. 예를 들어 메모리 업데이트, 계획 재구성, 점수(확신도) 갱신 등이 있다. 메모리 업데이트는 Agent가 새로운 정보를 저장하거나, 이전 경험을 반영하여 행동 전략을 수정하는 과정이다. 이 과정이 명시적으로 관리되지 않으면, 목표가 점점 왜곡되는 “goal drift” 현상이 발생할 수 있다. 실제로, 장기적으로 운영되는 Agent 시스템에서는 메모리 업데이트 정책이 불명확할 경우, Agent가 본래의 목표와 점점 멀어지는 현상이 빈번하게 발생합니다. 이를 방지하기 위해서는 메모리 업데이트의 기준과 절차를 명확히 정의하고, 주기적으로 검토 및 조정하는 것이 필요합니다.

#### 계획 재구성과 점수 갱신

계획 재구성은 Agent가 환경 변화 또는 내부 평가 결과에 따라 기존 계획을 수정하는 과정이다. 점수 또는 확신도 갱신은 Agent가 각 행동의 성공 가능성, 신뢰도 등을 평가하여 다음 행동을 결정하는 데 활용한다. 내부 상태 변화는 외부 상태 변화만큼 감사(audit) 대상이 되어야 하며, Observability(관측성) 시스템에서 1급 이벤트로 관리해야 한다. 엔터프라이즈 환경에서는 내부 상태 변화 이벤트를 실시간으로 모니터링하고, 이상 징후 발생 시 즉각적인 알림과 대응이 가능하도록 시스템을 설계해야 합니다. 또한, 내부 상태 변화 이력은 장기적인 성능 분석과 개선의 근거 자료로 활용됩니다.

#### 내부 상태 변경 이벤트 분류 예시

이벤트 유형	대표 사례	감사/관측 방법
메모리 업데이트	경험 저장, 회고	트레이스, 로그
계획 재구성	목표 변경, 경로 수정	트레이스
점수 갱신	확신도 평가	트레이스, 알림

### 4.3.3 새 액션 파생과 Reflexion 형 자기 수정 루프

Agent의 행동은 연쇄적으로 새로운 액션을 파생시키며, 이 과정에서 자기 평가와 수정이 이루어지는 피드백 루프가 형성됩니다. Reflexion 패턴은 이러한 자기 수정 루프를 체계적으로 구현하는 대표적인 방식으로, Agent의 장기적 성능 개선과 실패 복구율 향상에 큰 기여를 하고 있습니다. 본 절에서는 새 액션 파생 구조와 Reflexion 패턴의 기술적 세부사항, 그리고 실제 적용 시의 성능 개선 효과를 구체적으로 설명합니다.

#### 새 액션 파생의 구조

한 Action이 다음 Action을 유발하는 연쇄 구조는 Agent의 행동이 단순히 일회성에 그치지 않고, 지속적으로 이어지는 피드백 루프를 형성한다. 예를 들어 실패한 작업이 다시 시도되거나, 새로운 정보를 바탕으로 추가 행동이 발생한다. 이러한 구조는 Agent가 복잡한 문제를 단계적으로 해결하고, 예기치 못한 상황에 유연하게 대응할 수 있게 해줍니다. 실제로, 연쇄적 액션 구조를 도입한 Agent는 단일 실패에 머무르지 않고, 다양한 대안 행동을 시도함으로써 전체 시스템의 복원력과 안정성을 높입니다.

#### Reflexion 형 자기 평가 노드

Reflexion(NeurIPS 2023, Shinn et al.)은 verbal RL(언어적 강화학습)로 에피소드 메모리에 회고(reflection)를 저장하고, 실패 관측 → 원인 서술 → 재시도(closed-loop)를 완성한다. LangGraph 등 주요 프레임워크는 Reflexion 노드 패턴을 기본 탑재하여, Agent가 자기 평가를 통해 품질을 지속적으로 개선할 수 있도록 한다. Reflexion 패턴을 적용하면, Agent는 각 행동의 결과를 스스로 평가하고, 실패 원인을 분석하여 다음 행동 전략을 조정할 수 있습니다. 이는 장기적으로 Agent의 성능을 지속적으로 개선하는 데 매우 효과적입니다.

#### 자기 평가 노드 ON/OFF 패턴

Agent 설계 시 자기 평가 노드의 ON/OFF를 기본 패턴 스위치로 제공하면, 시스템의 자기 수정 능력을 상황에 맞게 조정할 수 있다. 자기 평가 노드가 활성화되면, 실패 관측과 원인 서술이 자동화되어 장기적으로 성능이 개선된다. 실제 운영에서는 업무 중요도나 리스크 수준에 따라 자기 평가 기능의 활성화 여부를 조정할 수 있으며, 이를 통해 운영 효율성과 품질 개선의 균형을 맞출 수 있습니다. 또한, 자기 평가 결과를 관리자에게 리포트 형태로 제공함으로써, Agent의 자기 수정 과정을 투명하게 모니터링할 수 있습니다.

#### 자기 수정 루프 전·후 성능 비교 예시

기준	자기 수정 루프 없음	자기 수정 루프 있음
실패 복구율	40%	85%
장기 정확도	60%	90%
운영 비용	높음	낮음

## 5장: Prompt → Context → Harness 진화와 Framework → Workflow → Orchestration 제품 매핑

AI Agent 엔지니어링은 단순한 프롬프트 최적화에서 시작해, 컨텍스트 설계와 워크플로우/제약/피드백/틀체인 관리까지 확장되는 3단계 패러다임으로 진화해왔다. 동시에, 실제 제품군은 단일 호출(Prompt/Context)에서 Agent 스캐폴딩(Framework), 상태 기반 그래프(Workflow), 내구성 오케스트레이션(Orchestration)까지 5단계로 발전했다. 이 장에서는 엔지니어링 패러다임의 3-Era와 제품군 5단계를 교차 지도 위에 정렬하여, CTO와 조직이 “우리는 어느 단계에 있는가”를 명확하게 판정할 수 있도록 한다. 각 절에서는 시대별 대표 기법과 제품, 그리고 실제 조직에서의 적용 좌표를 구체적으로 제시한다.

### 5.1 Prompt → Context → Harness Engineering 의 3-Era 전환

AI Agent 엔지니어링의 발전은 단순한 프롬프트 작성 기술에서 시작하여, 점차 복잡한 컨텍스트 설계와 시스템 전체를 통합하는 Harness Engineering 단계로 진화해왔습니다. 2022년부터 2026년까지의 짧은 기간 동안, 이 분야는 세 가지 뚜렷한 시대적 패러다임(Prompt, Context, Harness)을 거치며 급격한 변화를 겪었습니다. 각 시대는 기술적 요구, 조직 내 역할, 그리고 대표 도구의 변화와 함께, 실제 제품군의 발전과도 긴밀하게 연결되어 있습니다. 이 절에서는 각 패러다임의 정의와 특징, 그리고 대표적인 기법 및 도구를 체계적으로 정리하여, 조직이 현재 위치를 진단하고 미래 전략을 수립하는 데 실질적인 도움을 제공하고자 합니다.

## 5.1.1 Prompt Engineering Era (2022~24) — Few-shot · Role · Format 의 전성기

Prompt Engineering Era(2022~24)는 LLM 활용의 초창기 단계로, 단일 instruction의 형태, 어휘, 예시를 최적화하는 기술이 중심이었다. Few-shot learning은 모델에게 여러 예시를 제공해 추론의 품질을 높였고, Role 프롬프트는 “너는 전문가다” 등 역할 기반 지시로 결과의 일관성을 강화했다. Format 프롬프트는 출력 형식을 지정해 응답의 구조를 제어했다. CoT(Chain-of-Thought) 프롬프트는 단계별 사고를 유도하며 복잡 문제 해결에 기여했다. 이 시대의 대표 도구는 내부 프롬프트 라이브러리와 Playground(OpenAI, Claude 등)였다.

Prompt Engineering은 “한 문장 쓰는 기술”로서 빠른 성과를 냈지만, 복잡한 태스크에서 재현성, 확장성, 유지보수에 한계를 드러냈다. 단일 호출은 컨텍스트 초과, 외부 상태 접근 불가, 자기 검증 불가, 단일 출력만 가능하다는 구조적 한계가 있었다. CTO는 이 시대의 성과를 “지금도 유효하지만 부족하다”는 관점에서 평가해야 한다.

이 시기의 대표적인 프롬프트 기법은 Few-shot, Role, Format, CoT 등으로 구분됩니다. Few-shot은 모델에게 여러 예시를 제공하여 추론 품질을 높이고, Role 프롬프트는 “너는 전문가다”와 같이 역할을 명확히 지정하여 일관성을 확보합니다. Format 프롬프트는 출력 형식을 JSON, 표 등으로 지정해 결과의 구조를 통제하며, CoT(Chain-of-Thought)는 단계별 사고를 유도하여 복잡한 문제 해결에 효과적입니다. 이러한 기법들은 OpenAI Playground, Claude API, LangChain, LlamaIndex 등 다양한 플랫폼에서 활용되었습니다.

기법	설명	대표 도구/플랫폼
Few-shot	예시 기반 프롬프트	OpenAI Playground
Role	역할 지정 프롬프트	Claude API, 내부 라이브러리
Format	출력 형식 제어 프롬프트	Playground, Custom API
CoT	단계별 사고 유도 프롬프트	LangChain, LlamaIndex

Prompt Era의 핵심 유산은 조직 내 프롬프트 라이브러리입니다. 이 라이브러리는 이후 Context/Harness 단계에서도 재활용되며, 프롬프트 설계의 기본 자산으로 남습니다. 실제로 많은 조직이 이 시기에 축적한 프롬프트 자산을 바탕으로, 이후 단계의 복잡한 시스템 설계에서도

효율적으로 활용할 수 있었습니다. 그러나 단일 프롬프트 기반의 접근 방식은 복잡한 비즈니스 요구와 장기적 운영에는 한계가 명확하였으며, 이는 곧 다음 단계로의 진화를 촉진하는 요인이 되었습니다.

### 5.1.2 Context Engineering Era (2025~) — Karpathy 가 명명한 “Context Window 채우기” 의 기술

Context Engineering Era는 AI Agent 시스템이 단순한 프롬프트 최적화에서 벗어나, 외부 정보와 상태를 적극적으로 활용하는 방향으로 진화한 시기입니다. 2025년 Karpathy는 “Context engineering은 context window를 다음 step에 필요한 정보로 채우는 기술”이라고 명명하며, 이 패러다임의 본질을 명확히 했습니다. 이 시기에는 RAG(Retrieval-Augmented Generation), 메모리 계층, 툴 스키마 삽입 등 복합적인 컨텍스트 설계가 중심이 되었습니다. Prompt만으로는 해결할 수 없는 복잡성, 장기 실행, 외부 데이터 및 도구 연동이 필요한 상황에서 Context Engineering은 필수적인 기술로 부상하였습니다.

RAG는 외부 문서, 데이터베이스, API 등에서 필요한 정보를 검색하여 LLM의 컨텍스트로 주입하는 기술입니다. 이로 인해 모델이 최신 정보와 도메인 지식을 활용할 수 있게 되었으며, 실제 업무 환경에서의 적용성이 크게 향상되었습니다. 메모리 계층은 Short-term(대화 버퍼), Long-term(Vector DB), Episodic(Reflexion 회고), Semantic(Knowledge Graph) 등으로 세분화되어, Agent의 상태와 경험을 체계적으로 관리할 수 있도록 지원합니다. 툴 스키마 삽입은 MCP 등 외부 도구 호출을 컨텍스트에 포함하는 방식으로, Agent가 다양한 외부 시스템과 연동하여 복합적인 작업을 수행할 수 있게 합니다.

컴포넌트	역할	대표 제품/기술
RAG	외부 지식/문서 검색	LangChain, LlamaIndex
메모리 계층	상태/경험 관리	Pinecone, Neo4j, Qdrant
툴 스키마 삽입	외부 도구 호출	MCP, OpenAI function call

이러한 변화에 따라 조직 내에서는 Prompt 엔지니어와 별도로 Context 설계자라는 새로운 직무가 등장하였습니다. Context 설계자는 RAG, 메모리, 툴 연동 역량을 중점적으로 평가받으며, 복잡한 컨텍스트 설계와 시스템 통합을 주도합니다. 실제로 많은 기업들이 이 시기부터 컨텍스트

설계 전문 인력을 채용하거나, 기존 엔지니어의 역할을 재정의하는 움직임을 보이고 있습니다. 이처럼 Context Engineering Era는 AI Agent 시스템의 실질적 활용성과 확장성을 크게 높인 중요한 전환점이 되었습니다.

### 5.1.3 Harness Engineering Era (2026~) — LangGraph Deep Agents · Claude Code 가 보여주는 형태

Harness Engineering Era(2026~)는 AI Agent의 워크플로우, 제약, 피드백 루프, 튜체인, 라이프사이클을 통합 설계하는 단계입니다. 이 시기에는 단순한 프롬프트나 컨텍스트 설계만으로는 해결할 수 없는 복잡한 시스템 요구가 등장하였고, 이에 따라 전체적인 시스템 제어와 내구성, 피드백 구조를 중심으로 하는 Harness Engineering이 필수적으로 부상하였습니다. LangGraph Deep Agents는 planning tool, filesystem, subagent spawn 등 복합 워크플로우를 제공하며, Claude Code는 SWE-bench 80.9% 달성 등 실제 코드 생산성에 기여하는 등, Harness Engineering의 실질적 효과를 보여주고 있습니다.

이 단계에서는 워크플로우 설계, 제약 조건 정의, 피드백 루프 구현, 튜체인 통합, 라이프사이클 관리 등 Agent 시스템의 상위 설계가 핵심이 됩니다. Harness Engineer라는 별도의 직무가 신설되어, 시스템 전체의 안정성과 확장성을 책임지는 역할을 수행하게 됩니다. Harness Engineering은 AI Agent의 복잡한 운영과 실제 비즈니스 적용에 있어 필수적인 요소로 자리잡았으며, CTO는 조직 내 역할 분담과 기술 도입 시 Harness 단계로의 진입 여부를 명확히 평가해야 합니다.

Harness Engineering의 도입은 조직이 단순한 프로토타입 수준을 넘어, 실제 운영 환경에서의 안정성과 확장성을 확보하는 데 결정적인 역할을 합니다. 예를 들어, LangGraph와 같은 제품은 복잡한 상태 기반 워크플로우와 피드백 루프를 구현할 수 있게 해주며, Temporal과의 연동을 통해 내구성 있는 오케스트레이션을 실현할 수 있습니다. Claude Code의 SWE-bench 80.9% 달성 사례는, 실제 소프트웨어 엔지니어링 업무에서 Harness Engineering이 얼마나 큰 생산성 향상을 가져올 수 있는지를 보여주는 대표적인 지표입니다.

시대	핵심 기술/기법	대표 도구/제품	조직 역할
Prompt Era	프롬프트 최적화, Few-shot, Role	Playground, Claude API	Prompt Engineer
Context Era	RAG, 메모리, 톨 스키마	LangChain, LlamaIndex	Context Designer

Harness Era	워크플로우, 제약, 피드백, 라이프사이클	LangGraph, Claude Code	Harness Engineer
-------------	---------------------------	------------------------	------------------

Harness Engineering은 AI Agent의 복잡한 운영과 실제 비즈니스 적용에 필수적입니다. CTO는 조직 내 역할 분담과 기술 도입 시 Harness 단계로의 진입 여부를 명확히 평가해야 하며, 이를 통해 조직의 기술적 성숙도와 경쟁력을 한층 강화할 수 있습니다.

## 5.2 Prompt → Context → Framework → Workflow → Orchestration 5단 제품 매핑

AI Agent 관련 제품군은 엔지니어링 패러다임의 변화와 함께, 기술적 요구와 조직의 도입 성숙도에 따라 다섯 단계로 체계적으로 발전해왔습니다. 각 단계는 단순한 프롬프트 기반 호출에서 시작하여, 점차 컨텍스트 확장, Agent 스케폴딩, 상태 기반 워크플로우, 그리고 내구성 중심의 오케스트레이션으로 이어집니다. 이 절에서는 각 단계의 정의와 대표 제품, 그리고 CTO가 제품 선택 시 반드시 고려해야 할 핵심 포인트를 구체적으로 안내합니다. 이를 통해 조직은 현재 도입 단계와 향후 발전 방향을 명확히 진단할 수 있습니다.

### 5.2.1 Prompt · Context — OpenAI Playground 에서 LangChain/LlamaIndex RAG 까지

제품군의 1단계인 Prompt는 OpenAI Playground, 원시 Claude API 등 단일 호출 기반의 제품이 대표적입니다. 이 단계에서는 프롬프트 최적화만으로 태스크를 수행하며, 외부 도구나 복잡한 컨텍스트 설계가 없습니다. 조직은 이 단계의 제품을 Agent라고 부르지 않아야 하며, 단순한 LLM 활용에 머무릅니다.

2단계인 Context는 LangChain, LlamaIndex 등 RAG 기반 제품이 대표적입니다. 이 단계에서는 외부 문서, 데이터베이스, API 등에서 정보를 검색해 LLM 컨텍스트로 주입하는 기능이 추가됩니다. 단발 호출에서 벗어나 컨텍스트 window를 확장하지만, 여전히 도구 없는 추론에 머무르기 때문에 Agent의 정의에는 미치지 못합니다.

단계	대표 제품	주요 기능	Agent 여부
Prompt	OpenAI Playground, Claude API	프롬프트 최적화, 단일 호출	X
Context	LangChain, LlamaIndex (RAG)	컨텍스트 채우기, 외부 검색	X

이처럼 RAG 기반 제품은 Agent가 아닌 Context 단계로 분리 표기해야 하며, 도구 사용, 피드백 루프, 워크플로우 설계가 없는 제품은 Agent로 분류하지 않는 것이 CTO의 올바른 판단 기준입니다. 실제로 많은 조직이 RAG를 도입하면서 Agent 도입으로 오해하는 사례가 있으나, 이는 기술적 성숙도 진단에 혼선을 줄 수 있으므로 명확한 구분이 필요합니다.

### 5.2.2 Framework 단계 — LangChain Agents · LlamaIndex AgentWorkflow · CrewAI Agent/Task

Framework 단계에서는 Agent 스키펴딩이 본격적으로 등장합니다. LangChain Agents, LlamaIndex AgentWorkflow, CrewAI Agent/Task가 대표적인 제품으로, 이 단계에서는 ReAct 루프(Thought-Action-Observation)와 툴 호출(function call, MCP 등)이 표준화되어 Agent의 기본 구조가 제공됩니다. 즉, 단순한 프롬프트/컨텍스트 설계에서 벗어나, 실제로 도구를 호출하고, 복잡한 역할 분담과 상태 관리를 할 수 있는 Agent 시스템의 기초가 마련됩니다.

Framework 단계는 프로토타이핑의 단위로 인식됩니다. 조직은 이 단계에서 다양한 Agent 구조를 실험하고, 실제 운영에 앞서 설계의 적합성을 검증할 수 있습니다. 예를 들어, LangChain Agents는 대규모 생태계와 LangGraph 연동을 지원하며, LlamaIndex AgentWorkflow는 RAG 기반의 문서 처리에 특화되어 있습니다. CrewAI Agent/Task는 역할 기반의 Crew 구조와 독립 프레임워크로 차별화됩니다.

제품	루프 구조	툴 호출 방식	특이점
LangChain Agents	ReAct	Function call	대규모 생태계, LangGraph 연동
LlamaIndex AgentWorkflow	ReAct/Workflow	Function call/MCP	RAG 기반, 문서 처리 특화
CrewAI Agent/Task	Role/Goal/Backstory	Function call/MCP	역할 기반 Crew, 독립 프레임워크

Framework 단계는 Agent의 구조적 설계와 툴 연동, 역할 분담을 명확히 하며, 이후 Work-

flow/Orchestration 단계로의 확장 기반을 제공합니다. 실제로 이 단계에서의 실험과 검증이 조직의 Agent 시스템 성숙도에 큰 영향을 미치며, 성공적인 프로덕션 도입의 전제 조건이 됩니다.

### 5.2.3 Workflow · Orchestration — LangGraph · LlamaIndex Workflows · Temporal + LangGraph

Workflow 단계는 상태 기반 그래프 구조를 도입하여, 복잡한 프로세스와 다양한 상태 전이를 체계적으로 관리할 수 있게 합니다. LangGraph, LlamaIndex Workflows, CrewAI Flow 등이 대표적인 제품입니다. 이 단계에서는 노드/엣지, 체크포인트, HITL(Human-in-the-Loop), Streaming 등 복합 워크플로우가 구현되어, 실제 업무 환경에서의 다양한 요구를 충족시킬 수 있습니다.

Orchestration 단계는 다중 Agent와 내구성(durable)을 강조하는 단계로, Temporal + LangGraph, OpenAI Agents SDK Handoff, Microsoft Agent Framework 등이 대표 제품입니다. 이 단계에서는 retry, state, failure recovery 등 장기 실행과 안정성을 보장하는 기능이 필수적으로 도입됩니다. 특히, 프로덕션 환경에서는 내구성 있는 오케스트레이션이 운영 안정성의 핵심 요소로 간주됩니다.

단계	대표 제품	주요 기능	Durable 여부
Workflow	LangGraph, LlamaIndex Workflows, CrewAI Flow	상태 기반 그래프, 체크포인트, HITL	△
Orchestration	Temporal + LangGraph, OpenAI Agents SDK Handoff, Microsoft Agent Framework	다중 Agent, Durable, Failure Recovery	○

CTO는 프로덕션 환경에서는 반드시 4~5단으로 올라와야 하며, “durable”이 운영 안정성의 핵심임을 인식해야 합니다. 실제로, 장기 실행과 장애 복구가 보장되지 않는 시스템은 실무 적용에서 큰 위험을 초래할 수 있으므로, Orchestration 단계의 도입 여부가 조직의 기술적 성숙도를 가늠하는 중요한 기준이 됩니다.

### 5.3 3-Era × 5-Product 2축 지도와 제품 좌표 찍기

AI Agent 엔지니어링의 3-Era(프롬프트, 컨텍스트, 하니스)와 제품군의 5단계(프롬프트, 컨텍스트, 프레임워크, 워크플로우, 오케스트레이션)를 한눈에 볼 수 있는 2차원 지도는 조직이 현재 위치와 목표 좌표, 그리고 기술적 공백(Gap)을 명확히 시각화하는 데 매우 유용합니다. 이 절에서는 2축 지도 구성 방법과 대표 제품들의 좌표 배치, 그리고 Planning/Action과 Harness의 연결 방식에 대해 구체적으로 설명합니다. 이를 통해 조직은 기술 도입의 현황을 진단하고, 향후 발전 방향을 체계적으로 설계할 수 있습니다.

#### 5.3.1 3-Era × 5-Product 2축 지도의 구성 방법

2 차 원                      캔   버   스   는                      가   로   축   에                      5-  
Product(Prompt/Context/Framework/Workflow/Orchestration),                      세   로   축   에  
3-Era(Prompt/Context/Harness)를 배치하여 구성합니다. 이 지도는 각 제품을 한눈에  
정렬하고, 조직의 도입 좌표를 명시적으로 선언할 수 있게 해줍니다. 조직은 현재 좌표와 1년 후  
목표 좌표를 나란히 표시함으로써, 기술 도입의 Gap과 발전 방향을 명확히 파악할 수 있습니다.

Prompt	Context	Framework	Workflow	Orchestration	
Prompt Era					
Context Era					
Harness Era					

이러한 2축 지도는 CTO와 기술 리더가 조직의 기술적 위치를 객관적으로 진단하고, 전략적 목표를 설정하는 데 실질적인 도구로 활용될 수 있습니다. 예를 들어, 현재 조직이 Context Era의 Framework 단계에 머물러 있다면, 향후 1년 내에 Harness Era의 Workflow 또는 Orchestration 단계로 도약하는 것이 목표가 될 수 있습니다.

### 5.3.2 대표 제품들의 좌표 — LangGraph · CrewAI · OpenAI SDK · LlamaIndex · Temporal

대표적인 AI Agent 제품들을 2축 지도에 배치하면, 각 제품의 기술적 위치와 역할을 명확히 파악할 수 있습니다. 예를 들어, LangGraph는 Workflow × Harness, CrewAI는 Framework × Context/Harness, OpenAI SDK는 Framework~Orchestration × Context, LlamaIndex Workflows는 Workflow × Context, Temporal + LangGraph는 Orchestration × Harness 에 위치합니다.

Prompt	Context	Framework	Workflow	Orchestration	
Prompt Era	OpenAI Playground				
Context Era	LangChain, LlamaIndex	CrewAI, LlamaIndex AgentWorkflow	LlamaIndex Workflows	OpenAI SDK	
Harness Era	CrewAI	LangGraph	Temporal+Lang-Graph		



[그림 3] 3-Era × 5-Product 제품 좌표 지도

제품 간 비교는 반드시 동일 지도 위에서 이루어져야 하며, CTO는 “우리가 이미 쓰는 제품”의 좌표부터 먼저 찍어 Gap을 식별해야 합니다. 이를 통해 조직은 기술 도입의 현황과 부족한 부분을 명확히 파악하고, 전략적 투자와 인력 배치의 우선순위를 결정할 수 있습니다.

### 5.3.3 Planning 피드백 축과 Harness 의 closed-loop/open-loop 의 일치

AI Agent 시스템의 설계에서 Planning 피드백 축과 Harness의 closed-loop/open-loop 구조는 밀접하게 연결되어 있습니다. 3장에서 다룬 Planning의 “피드백 있음/없음” 분류는 Harness Engineering에서의 closed-loop(피드백 포함)와 open-loop(피드백 없음) 설계와 그대로 대응됩니다. Harness Engineering에서 closed-loop는 “피드백 있음” Planning 패러다임을 설계 단위로 품으며, open-loop는 단순 워크플로우, closed-loop는 환경, 인간, 모델 피드백을 포함한 복합 구조를 의미합니다.

Planning 피드백 종류	Harness 설계 루프	대표 제품/기술
없음	open-loop	단순 Workflow, CrewAI
있음	closed-loop	LangGraph, Temporal

CTO는 Harness 리뷰 체크리스트에 “closed-loop 필요성 평가” 항목을 반드시 포함해야 하며, 피드백 구조가 실제 운영 안정성에 직결됨을 인식해야 합니다. 실제 사례에서, closed-loop 구조를 도입한 조직은 장애 대응과 품질 개선에서 월등한 성과를 보였으며, open-loop만으로 운영되는 시스템은 예기치 못한 오류와 품질 저하에 취약한 것으로 나타났습니다. 따라서, 조직의 기술적 성숙도를 높이기 위해서는 closed-loop 구조의 도입 여부를 핵심 평가 지표로 삼아야 합니다.

## 6장: 주요 AI Agent Framework 정량 비교 — LangGraph · AutoGen · CrewAI · OpenAI SDK · LlamaIndex

AI Agent 프레임워크의 선택은 단순한 기술 비교를 넘어 조직의 전략적 방향, 운영 안정성, 비용 효율성, 그리고 향후 확장성까지 좌우하는 핵심 의사결정입니다. 2026년 기준, LangGraph, AutoGen, CrewAI, OpenAI Agents SDK, LlamaIndex는 각각 독자적인 설계 철학과 생태계를 구축하며, 엔터프라이즈·내부 자동화·OpenAI 중심 스택 등 다양한 도입 시나리오에 맞는 선택지를 제공합니다. 본 장에서는 각 프레임워크의 라이선스, 커뮤니티 규모, 설계 패러다임, 프로덕션 준비도, 비용, MCP 지원 성숙도, 학습 곡선 및 업타임 등 CTO가 실제 구매 결정에 사용하는 7대 축을 정량적으로 비교합니다. 모든 비교 데이터와 표는 사내 RFQ(입찰/선정 기준) 표로 바로 활용할 수

있도록 구조화하였으며, 각 프레임워크의 최신 공식 자료와 업계 서베이, 실제 기업 도입 사례를 근거로 삼았습니다.

---

## 6.1 5대 프레임워크 라이선스 · 생태계 · 커뮤니티 정량 지표

AI Agent 프레임워크의 생태계 규모는 도입 이후 유지보수, 인력 채용, 커뮤니티 지원, 그리고 장기적 조달 안정성까지 영향을 미치는 중요한 기준입니다. 라이선스와 커뮤니티 규모, 다운로드 수는 감정적 선호가 아닌 객관적 수치로 비교되어야 하며, 이 장에서는 LangChain/LangGraph, AutoGen/AG2·CrewAI, OpenAI SDK·LlamaIndex의 최신 지표를 근거로 삼아 각 프레임워크의 생태계 현황을 정리합니다. CTO는 이 데이터를 통해 lock-in 리스크와 조달 용이성을 균형 있게 평가할 수 있습니다.

### 6.1.1 LangChain · LangGraph — MIT · 100k/29.5k+ Stars · 누적 47M 다운로드의 de facto standard

LangChain과 LangGraph는 2026년 현재 AI Agent 프레임워크 분야에서 사실상 표준으로 자리매김하고 있습니다. MIT 라이선스 채택으로 인해 기업, 공공기관, 스타트업 등 다양한 조직에서 자유롭게 도입할 수 있으며, 라이선스 리스크가 거의 없습니다. 특히 langchain-ai/langchain의 GitHub 스타 수는 100k+에 달하며, langchain-ai/langgraph 역시 29.5k+의 스타를 기록하고 있습니다. PyPI 기준 누적 다운로드 수는 47M+로, 이는 경쟁 프레임워크 대비 압도적인 수치입니다. 이러한 수치는 단순한 인기 지표를 넘어, 실제 현업에서의 채택률과 유지보수 가능성, 인력 수급의 용이성까지 반영합니다.

상용 플랫폼 측면에서는 LangSmith(관측성 SaaS)와 LangGraph Platform(매니지드 서비스)이 제공되어, OSS 코어와 상용 관측/운영 플랫폼이 분리된 구조를 갖추고 있습니다. 이는 CTO 입장에서 조달 리스크를 낮추고, 필요에 따라 상용 서비스를 활용할 수 있는 유연성을 제공합니다. 커뮤니티 역시 활발하게 운영되고 있으며, 공식 문서와 예제, 플러그인, 확장 모듈이 지속적으로 업데이트되고 있습니다. 이러한 점은 장기적 유지보수와 기능 확장, 그리고 신뢰할 수 있는 생태계 기반을 확보하는 데 중요한 역할을 합니다.

프레임워크	라이선스	GitHub Stars	PyPI 다운로드	상용 플랫폼	커뮤니티 기여
LangChain	MIT	100k+	47M+	LangSmith	활발
LangGraph	MIT	29.5k+	47M+	LangGraph Platform	활발

### 6.1.2 AutoGen/AG2 · CrewAI — Apache 2.0/MIT · 54.7k/45.9k Stars 와 분기 중인 커뮤니티

AutoGen(Microsoft Research)은 오픈소스 AI Agent 프레임워크 중에서도 초기 대규모 커뮤니티와 브랜드 파워를 확보한 프로젝트입니다. Apache 2.0 라이선스 채택으로 상업적 활용에 제약이 없으며, GitHub 스타 수 54.7k+와 PyPI 월간 다운로드 0.1M 수준을 기록하고 있습니다. 그러나 2025년 Microsoft가 “Microsoft Agent Framework”로 이관을 발표하면서 AutoGen 본체는 사실상 maintenance mode에 들어갔고, 커뮤니티는 AG2라는 포크로 분기되어 독립 개발이 이어지고 있습니다. AG2는 3k+ 스타와 0.1M 다운로드로 성장 중이나, 브랜드 파워와 성장성 측면에서는 아직 AutoGen 본체에 미치지 못합니다. 이처럼 AutoGen/AG2 생태계는 “오래된 브랜드”와 “현재 성장성”이 분리되는 과도기적 상황에 놓여 있습니다.

CrewAI는 MIT 라이선스를 기반으로 45.9k+ 스타, 월간 PyPI 1.3M 다운로드(AG2 대비 13 배), 100k+ 인증 개발자라는 강력한 성장세를 보이고 있습니다. CrewAI는 v1.14에서 LangChain 의존성을 완전히 제거하며 독립 프레임워크로 전환하였고, 커뮤니티 중심의 활발한 개발과 빠른 기능 추가가 이루어지고 있습니다. CrewAI의 이러한 성장성은 실제 기업 도입 사례와 활발한 커뮤니티 활동, 그리고 독립성 강화 전략에서 비롯된 것으로, CTO 입장에서는 장기적 조달 안정성과 생태계 lock-in 리스크 완화 측면에서 긍정적으로 평가할 수 있습니다.

프레임워크	라이선스	GitHub Stars	PyPI 다운로드	커뮤니티 상태	성장성
AutoGen	Apache 2.0	54.7k+	0.1M	maintenance mode	낮음
AG2	Apache 2.0	3k+	0.1M	독립 포크	중간
CrewAI	MIT	45.9k+	1.3M	독립 프레임워크	높음

### 6.1.3 OpenAI Agents SDK · LlamaIndex Workflows — 신규 SDK 와 RAG 기반의 현재 위치

OpenAI Agents SDK는 OpenAI가 공식적으로 제공하는 경량 SDK로, MIT 라이선스와 Python/TypeScript 기반의 친숙한 개발 환경을 제공합니다. 20줄 미만의 코드로 Agent를 구축할 수 있을 정도로 진입장벽이 낮아, 신규 진입자나 문서 중심 Agent 구축에 매우 적합합니다. OSS 생태계 내에서의 진입장벽을 크게 낮추었으며, 공식 문서와 예제가 풍부하게 제공되어 빠른 프로토타이핑이 가능합니다. 다만, 상용 서비스는 별도로 제공되지 않습니다.

LlamaIndex는 MIT 라이선스 기반의 프레임워크로, llama-index-core에 Workflows 기능을 포함하고 있습니다. 상용 LlamaCloud(Parse, Extract) 서비스도 제공되어, RAG(검색 증강 생성) 및 문서 처리 중심 Agent 구축에 특화되어 있습니다. Linux Foundation AI Survey 2025에 따르면, 생산 중 AI Agent의 68%가 OSS 프레임워크 기반임이 확인되며, LlamaIndex는 RAG/문서 처리 중심 Agent 분야에서 강점을 보입니다. 이처럼 LlamaIndex는 문서 처리와 RAG 워크플로우에 최적화된 구조와 상용 서비스 연계를 통해, 특정 도메인에 특화된 Agent 구축에 유리한 선택지를 제공합니다.

프레임워크	라이선스	학습곡선	상용 서비스	OSS 기반 비율	주요 강점
OpenAI Agents SDK	MIT	최저	없음	68%	진입장벽 낮음
LlamaIndex	MIT	완만	LlamaCloud	68%	RAG 특화

## 6.2 설계 패러다임 · 핵심 강점 · 프로덕션 준비도 비교

프레임워크의 설계 철학과 프로덕션 준비도는 단순 수치 이상의 의미를 가집니다. 실제 기업 도입 사례, 워크플로우 구조, 핵심 기능의 내장 여부, 그리고 감사·관측·HITL 등 운영 안정성 요소가 CTO의 최종 선택을 좌우합니다. 본 절에서는 LangGraph, AutoGen/AG2·CrewAI, OpenAI SDK·LlamaIndex의 설계 패러다임과 프로덕션 성숙도를 맞대어 비교합니다.

### 6.2.1 LangGraph 가 제공하는 State · Checkpoint · HITL 3축과 Klarna/Replit/Uber/LinkedIn 프로덕션 채택

LangGraph는 AI Agent 워크플로우를 노드/엣지 기반의 상태 그래프(Stateful Graph)로 모델링하는 독특한 설계 패러다임을 채택하고 있습니다. 이 구조는 복잡한 분기, 순환, 조건부 실행 등 엔터프라이즈 환경에서 요구되는 다양한 워크플로우를 명확하게 구현할 수 있게 해줍니다. 체크포인팅 기능은 각 노드별로 상태를 저장하여 장애 발생 시 신속한 복구와 감사(audit) 추적이 가능하도록 설계되어 있습니다. 이러한 구조는 규제 산업이나 미션 크리티컬 환경에서 요구되는 신뢰성과 투명성을 확보하는 데 매우 유리합니다.

HITL(Human-in-the-Loop) 기능은 사람이 개입할 수 있는 게이트를 기본 제공하며, 실시간 스트리밍 지원을 통해 대화형 Agent나 실시간 톨 호출 시나리오에 적합합니다. 실제로 Klarna, Replit, Uber, LinkedIn, GitLab 등 글로벌 기업들이 LangGraph를 프로덕션 환경에서 채택하고 있으며, 복잡한 규제 요건이나 대규모 워크플로우 자동화에 LangGraph를 기본값으로 선택하는 경향이 뚜렷합니다. 이러한 실제 도입 사례는 LangGraph의 설계 패러다임이 현업에서 충분히 검증되었음을 의미합니다.

기능	LangGraph 지원	실제 채택 기업
상태 그래프	O	Klarna, Uber, LinkedIn, GitLab
체크포인팅	O	Klarna, Uber
HITL	O	Klarna, Replit
스트리밍	O	LinkedIn, GitLab

### 6.2.2 AutoGen/AG2 의 대화형 다중 Agent 와 CrewAI 의 역할 기반 Crew — 장점과 함정

AutoGen/AG2는 Conversational Multi-Agent 구조를 내장하여, 복수의 에이전트가 토론하거나 협력하는 시나리오에 최적화되어 있습니다. ConversableAgent, GroupChat, UserProxyAgent 등 다양한 내장 컴포넌트를 통해, 사용자는 실제 대화형 에이전트, 그룹 대화, 사용자 프록시 등 복잡한 상호작용을 손쉽게 구현할 수 있습니다. 하지만 2025년 Microsoft의 Agent Frame-

work 이관 이후, AutoGen 본체는 maintenance mode에 들어가면서 장기적 유지보수와 기능 확장에 대한 불확실성이 커졌습니다. AG2로의 포크 이관 경로는 커뮤니티 주도의 독립 개발이 진행되고 있으나, 브랜드 파워와 생태계 안정성 측면에서 여전히 리스크가 존재합니다. CTO는 이러한 이관 리스크와 커뮤니티 분기 상황을 반드시 고려해야 하며, 장기적 조달 안정성에 대한 평가가 필요합니다.

CrewAI는 Role-Based Crew 패러다임을 채택하여, 각 Agent가 명확한 역할(role), 목표(goal), 배경(backstory)을 가지고 협업하는 구조를 제공합니다. Process, Flow(v1+) 등 워크플로우 중심의 설계가 특징이며, v1.14에서 LangChain 의존성을 완전히 제거함으로써 완전한 독립 프레임워크로 전환하였습니다. CrewAI의 가장 큰 강점은 빠른 Time-to-Prod(프로덕션 전환 속도)로, 실제 기업 도입 시 복잡한 설정 없이 신속하게 프로덕션 환경에 적용할 수 있다는 점입니다. 커뮤니티 중심의 활발한 개발과 빠른 기능 추가 역시 CrewAI의 경쟁력입니다. 그러나 다중 에이전트 구조의 복잡성, 역할 설계의 난이도 등은 도입 시 고려해야 할 요소입니다.

기준	AutoGen/AG2	CrewAI
설계 패러다임	대화형 Multi-Agent	역할 기반 Crew
유지보수	포크 이관 리스크	독립 프레임워크
프로덕션 준비	실험적	빠른 전환

### 6.2.3 OpenAI Agents SDK 의 Handoff/Guardrail · LlamaIndex Workflows 의 Event-Driven Step

OpenAI Agents SDK는 Handoff(함수-툴 파이프라인 분리)와 Guardrail(입력=첫 에이전트, 출력=최종 에이전트) 기능을 내장하고 있어, 복잡한 파이프라인을 간결하게 구성할 수 있습니다. 내장 트레이싱 기능은 LLM 생성, tool call, handoff, guardrail, custom event 등 다양한 이벤트를 실시간으로 수집할 수 있어, 운영 관측성과 디버깅에 매우 유리합니다. 이러한 구조는 OpenAI 중심 스택을 사용하는 조직에서 빠른 도입과 운영 효율성을 보장합니다. 다만, 복잡한 워크플로우나 커스텀 확장성 측면에서는 한계가 있을 수 있습니다.

LlamaIndex Workflows는 Event(pydantic), @step, Context, async-first 구조를 채택하여, 문서 중심 Agent와 RAG(검색 증강 생성) 워크플로우에 최적화되어 있습니다. 각 단계별로

이벤트를 정의하고, 비동기 실행을 기본으로 하여 대규모 문서 처리나 RAG 기반 Agent 구축에 강점을 보입니다. OpenAI 중심 스택이나 문서 처리 중심 Agent 구축에 적합하며, 적용 범위가 명확하게 한정되어 있다는 점을 이해하고 도입해야 합니다. LlamaIndex는 문서 처리와 RAG 워크플로우에 특화된 기능을 제공하지만, 범용 Agent 구축에는 추가적인 설계가 필요할 수 있습니다.

기능	OpenAI SDK	LlamaIndex Workflows
Handoff	O	X
Guardrail	O	X
트레이싱	O	O
Event-Driven	X	O
문서 처리	X	O

## 6.3 비용 · 지연시간 · MCP 지원 · 학습 곡선 — RFQ 에 그대로 올릴 7대 축

구매 결정에 실제로 쓰이는 7대 축(비용, 지연시간, MCP 지원 성숙도, 학습 곡선, 체크포인팅, 업타임, 주활용)은 사내 RFQ(입찰/선정 기준) 표로 바로 활용할 수 있도록 구조화해야 합니다. 본 절에서는 각 프레임워크의 최신 비용·지연 프로파일, MCP 지원 수준, 학습 곡선 및 프로덕션 업타임을 비교합니다.

### 6.3.1 태스크당 평균 비용 \$0.35 (AutoGen) vs CrewAI TTP -40% — 비용·지연 프로파일

AI Agent 프레임워크의 비용과 지연시간은 실제 운영 환경에서 매우 중요한 선택 기준입니다. AutoGen은 태스크당 평균 \$0.35, 20회 이상의 LLM 호출로 업계 최고 수준의 비용과 지연시간을 기록합니다. 이는 다중 에이전트 토론 구조의 overhead가 주된 원인으로, 복잡한 협업 시나리오에서는 비용이 급격히 증가할 수 있습니다. CrewAI는 표준 대비 40% 빠른 Time-to-Prod(TTP)를 달성하며, 태스크당 평균 비용이 \$0.21로, 비용 효율성과 빠른 전환이 강점입니다.

OpenAI SDK는 native function call 구조로 지연시간이 최저 수준이며, LangGraph는 표준, LlamaIndex는 중상 수준의 비용과 지연시간을 보입니다.

실제 기업에서는 월간 예상 태스크 수와 단가를 곱해 1년간 총소유비용(TCO)을 산정하는 방식이 일반적입니다. 다중 에이전트 구조의 overhead와 각 프레임워크의 비용 프로파일을 정확히 비교하는 것이 중요합니다. CrewAI의 빠른 전환 속도와 낮은 비용은 내부 자동화나 반복 작업이 많은 환경에서 특히 유리하며, AutoGen의 높은 비용과 지연시간은 실험적·연구 목적에는 적합하지만, 대규모 프로덕션에는 부담이 될 수 있습니다.

프레임워크	평균 비용(\$/태스크)	LLM 호출수	TTP(전환 속도)	지연시간
AutoGen	0.35	20+	느림	높음
CrewAI	0.21(-40%)	5~10	빠름	낮음
LangGraph	0.29	10~15	표준	중간
OpenAI SDK	0.25	5~10	빠름	최저
LlamaIndex	0.32	10~15	표준	중상

### 6.3.2 MCP 지원 성숙도 — LangGraph 의 1급 노드 + 스트리밍 대 나머지의 함수 호출 처리

MCP(Model Context Protocol) 지원 성숙도는 최근 AI Agent 프레임워크 선정에서 매우 중요한 기준으로 부상하고 있습니다. LangGraph는 MCP 툴을 그래프 노드의 1급 객체로 취급하며, 스트리밍 지원까지 포함하여 가장 깊은 통합을 자랑합니다. 이로 인해 복잡한 툴 체인이나 외부 시스템 연동이 필요한 엔터프라이즈 환경에서 탁월한 확장성과 안정성을 제공합니다.

OpenAI SDK는 MCP를 내장하고 있어 기본적인 MCP 워크플로우를 손쉽게 구현할 수 있으며, 스트리밍 지원도 제공되어 실시간 데이터 처리에 유리합니다. 반면, CrewAI와 AutoGen은 MCP를 함수 호출 방식으로 처리하며, 스트리밍 기능은 미지원입니다. 이로 인해 복잡한 MCP 워크플로우나 대규모 실시간 처리가 필요한 환경에서는 한계가 있을 수 있습니다. LlamaIndex는 MCP 통합이 진행 중이며, 향후 지원 성숙도가 높아질 것으로 기대됩니다.

MCP 지원 깊이는 RFQ(입찰/선정 기준)에서 별점 항목으로 평가될 정도로 중요하며, 실제 도입 시 프레임워크별 MCP 지원 방식과 통합 깊이를 면밀히 검토해야 합니다.

프레임워크	MCP 지원 방식	스트리밍 지원	통합 깊이
LangGraph	1급 노드	O	깊음
OpenAI SDK	내장	O	중간
CrewAI	함수 호출	X	낮음
AutoGen	함수 호출	X	낮음
LlamaIndex	통합 중	X	진행 중

### 6.3.3 학습 곡선 · 체크포인팅 · 업타임 — 2026 업계 컨센서스 3분할

AI Agent 프레임워크의 학습 곡선, 체크포인팅, 프로덕션 업타임은 실제 도입과 운영에서 CTO가 가장 민감하게 평가하는 항목입니다. OpenAI SDK는 20줄 미만의 코드로 동작하는 최저 학습 곡선을 제공하여, 신규 개발자나 빠른 프로토타이핑에 매우 적합합니다. CrewAI는 완만한 학습 곡선과 빠른 프로덕션 전환 속도를 자랑하며, LangGraph는 상태/그래프 기반의 복잡한 구조로 인해 상대적으로 가파른 학습 곡선을 보입니다. 체크포인팅 및 Durable 기능은 LangGraph에서 완전 지원되며, OpenAI SDK와 LlamaIndex는 부분 지원, CrewAI와 AutoGen은 미지원입니다.

프로덕션 업타임 측면에서는 LangGraph, CrewAI, OpenAI SDK, LlamaIndex가 90% 이상의 business-grade 업타임을 제공하는 반면, AutoGen은 실험적 구조로 인해 70% 수준에 머물고 있습니다. Stack Overflow 2025-11 서베이 결과, multi-agent 실험 34%, AutoGen 41%, CrewAI 28% 언급이 확인되어, 실제 업계에서의 활용도와 인식이 반영되고 있습니다. 2026년 업계 컨센서스는 “규제 산업 → LangGraph, 내부 자동화 → CrewAI, OpenAI 스택 → Agents SDK”의 3분할 구도로 정리되며, 각 프레임워크의 강점이 명확하게 구분되고 있습니다.

프레임워크	학습 곡선	체크포인팅	Durable	업타임	주요 활용	MCP 지원	비용	지연	커뮤니티	상용 플랫폼	OSS 라이선스
Lang-Graph	가파름	O	O	90%+	규제산업	깊음	중	중	활발	있음	MIT
CrewAI	완만	X	X	90%+	내부자동화	낮음	낮음	낮음	활발	있음	MIT
OpenAI SDK	최저	△	△	90%+	OpenAI 스택	중간	낮음	최저	활발	없음	MIT

Auto-Gen	완만	X	X	70%	실험적	낮음	높음	높음	포크중	없음	Apache 2.0
LlamaIndex	완만	△	△	90%+	RAG/문서	진행중	중	중	활발	있음	MIT

## 7장: 엔터프라이즈 Agent 스택 — MCP · 메모리 · Observability · Durable Workflow

엔터프라이즈 환경에서 AI Agent의 도입은 단일 프레임워크만으로는 충분하지 않습니다. 실제로 2026년 기준, 엔터프라이즈 AI Agent 스택은 프로토콜 표준(MCP), 계층화된 메모리(VectorDB/GraphDB), 관측성(LangSmith/LangFuse/OTEL), 내구성(Temporal)까지 아우르는 5 레이어 구조로 정착되고 있습니다. 이 장에서는 각 레이어의 핵심 기술과 제품군, 그리고 엔터프라이즈 도입 시 반드시 고려해야 할 설계·운영 포인트를 상세히 다룹니다. IT 의사결정자는 이 장의 스택 구조를 사내 아키텍처 다이어그램의 참조 모델로 활용할 수 있습니다.

### 7.1 MCP — Agent의 도구 공간을 표준화한 JSON-RPC 프로토콜과 엔터프라이즈 임팩트

엔터프라이즈 AI Agent가 다양한 외부 도구와 안정적으로 연동되기 위해서는 표준화된 프로토콜이 필수적입니다. MCP(Model Context Protocol)는 이러한 요구에 부응하여, 2025년 이후 업계 표준으로 자리잡은 JSON-RPC 기반의 프로토콜입니다. 엔터프라이즈 환경에서는 MCP를 통해 도구 호출의 일관성, 상호운용성, 그리고 벤더 중립성을 확보할 수 있으며, 이는 IT 거버넌스와 장기적인 기술 투자 안정성에 큰 영향을 미칩니다. 본 절에서는 MCP의 기술 사양, 표준화 현황, 그리고 오픈소스 재단 이관이 엔터프라이즈에 미치는 임팩트를 다각도로 분석합니다.

## 7.1.1 MCP 2025-11-25 스펙 — tools/list · tools/call · resources · prompts · sampling 핵심 기능

MCP 2025-11-25 스펙은 JSON-RPC 기반의 client-server 아키텍처를 채택하여, Agent와 외부 시스템 간의 표준화된 통신을 보장합니다. MCP의 핵심 메서드는 다음과 같습니다:

메서드	설명
tools/list	Agent가 사용할 수 있는 도구 목록을 조회
tools/call	특정 도구(function/tool)을 호출
resources	외부 리소스(파일, 데이터 등) 접근/관리
prompts	프롬프트 템플릿 및 컨텍스트 관리
sampling	샘플링 파라미터(온도, top-p 등) 제어

이외에도 비동기 작업 지원, stateless 설계, server identity(서버 식별), 공식 MCP 레지스트리 등 최신 엔터프라이즈 요구사항이 반영되었습니다. 엔터프라이즈에서는 사내 시스템을 MCP 서버로 감싸는 우선순위 리스트를 작성하여, 점진적으로 Agent 도구 공간을 확장하는 전략이 권장됩니다.

비동기 작업 지원은 대규모 트랜잭션이나 장기 실행이 필요한 워크플로우에서 Agent의 유연성을 극대화합니다. stateless 설계는 서버 장애 시 복구와 스케일아웃에 유리하며, 공식 레지스트리는 도구의 신뢰성과 버전 관리를 체계화합니다. 이러한 요소들은 엔터프라이즈 환경에서 요구되는 안정성, 확장성, 컴플라이언스 준수에 직접적으로 기여합니다.

MCP의 각 메서드는 실제 엔터프라이즈 도입 시 다양한 시나리오에서 활용됩니다. 예를 들어, tools/list는 조직 내에서 등록된 모든 사내·외부 도구를 동적으로 조회할 수 있게 하여, Agent가 새로운 기능을 즉시 활용할 수 있도록 지원합니다. tools/call은 복잡한 비즈니스 프로세스 자동화에서 핵심 역할을 하며, resources 메서드는 파일 서버, 데이터 레이크 등 다양한 엔터프라이즈 자산과의 연동을 단순화합니다. prompts와 sampling은 LLM 기반 Agent의 품질 관리와 실험적 개선을 위한 필수 기능으로, 프롬프트 템플릿의 중앙 관리 및 샘플링 파라미터의 일관된 적용이 가능합니다.

특히, 공식 MCP 레지스트리의 도입은 엔터프라이즈 보안 및 컴플라이언스 관점에서 매우 중요합니다. 모든 도구와 리소스가 레지스트리에 등록·버전 관리됨으로써, 변경 이력 추적과 승인

프로세스가 체계화되고, 외부 감사 대응이 용이해집니다. 이러한 구조는 기존의 비표준적 도구 연동 방식에서 발생하던 운영 리스크를 크게 줄여줍니다.

기능 구분	MCP 메서드	설명
도구 탐색	tools/list	사용 가능한 도구 목록 반환
도구 호출	tools/call	지정 도구 실행 및 결과 반환
리소스 관리	resources	파일, 데이터 등 외부 리소스 접근
프롬프트 관리	prompts	프롬프트 템플릿 및 컨텍스트 관리
샘플링 제어	sampling	LLM 샘플링 파라미터 설정
비동기 처리	async	장기 실행 작업 지원

(출처: [MCP 공식 스펙](#))

## 7.1.2 월 97M SDK 다운로드 · OpenAI/Google/MS/Anthropic 공동 지지 — 표준화의 증거

2025년 기준 MCP SDK는 월 97M 이상의 다운로드를 기록하며, 업계에서 가장 널리 사용되는 AI Agent 도구 프로토콜로 자리매김했습니다. OpenAI, Google, Microsoft, Anthropic 등 글로벌 빅테크가 MCP를 공식 지원하며, “Why MCP Won” 기사와 “A Year of MCP 2025 Review” 등에서 MCP의 표준화가 상업적·기술적으로 이미 기정사실화되었음을 강조합니다.

MCP의 가장 큰 강점은 다수 벤더의 공동 지지에 기반한 상호운용성입니다. OpenAI, Google, Microsoft, Anthropic이 MCP를 공식 지원함에 따라, 엔터프라이즈는 특정 벤더에 종속되지 않고 다양한 AI·IT 자산을 통합할 수 있습니다. 이는 도구 생태계의 확장성과 조달 리스크 감소, 그리고 미래 호환성 확보에 결정적입니다.

MCP의 채택이 폭발적으로 증가한 배경에는 개발자 커뮤니티의 활발한 참여와 SDK의 접근성, 그리고 다양한 언어와 플랫폼에 대한 지원이 있습니다. 월 97M 다운로드라는 수치는 단순한 인기 지표를 넘어, 실제 엔터프라이즈 환경에서의 신뢰성과 안정성을 방증합니다. 예를 들어, OpenAI의 Agents SDK, Google의 Vertex AI, Microsoft의 Azure AI Agent, Anthropic의 Claude API 등 주요 서비스가 MCP를 기반으로 통합되고 있으며, SAP, Salesforce, IBM 등 전통적 엔터프라이즈 벤더들도 MCP 지원을 공식화하고 있습니다.

이러한 표준화는 엔터프라이즈 IT 전략 수립에 있어 중요한 의미를 가집니다. MCP를 도입함으로써, 조직은 미래의 기술 변화에 유연하게 대응할 수 있으며, 신규 도구나 서비스가 등장하더라도 최소한의 통합 비용으로 신속하게 확장할 수 있습니다. 또한, MCP 기반의 도구 호출 및 관리 체계는 보안 정책 준수, 운영 자동화, 그리고 감사 대응 등 엔터프라이즈 필수 요건을 충족시키는 데 큰 역할을 합니다.

기업/조직	MCP 채택 영역
OpenAI	Agents SDK, Function Calling
Google	Vertex AI, Gemini Agent
Microsoft	Azure AI Agent, Copilot
Anthropic	Claude API, MCP SDK
기타	SAP, Salesforce, IBM 등

(출처: [Pento, A Year of MCP 2025 Review](#), [The New Stack, Why MCP Won](#))

### 7.1.3 Linux Foundation 산하 Agentic AI Foundation 이관의 의미

2025년 12월, MCP는 Linux Foundation 산하 Agentic AI Foundation으로 거버넌스가 공식 이관되었습니다. 이는 오픈소스 프로토콜의 장수성과 안정성을 보장하는 핵심 신호로, 벤더 락인(종속) 리스크를 구조적으로 줄여줍니다. OSS 재단 거버넌스 체계는 기술 표준의 투명성, 커뮤니티 주도 발전, 그리고 조달 안정성 측면에서 엔터프라이즈 도입의 신뢰도를 한층 높입니다.

MCP가 Linux Foundation 산하로 이관된 것은 단순한 명칭 변경이 아니라, 엔터프라이즈 IT 조달과 장기 운영 전략에 실질적인 변화를 가져옵니다. 오픈소스 재단의 관리 하에 MCP는 벤더 종립성을 확보하게 되며, 기술 발전이 특정 기업의 이해관계에 좌우되지 않고, 글로벌 커뮤니티의 합의와 투명한 프로세스를 통해 이루어집니다. 이는 엔터프라이즈가 장기적으로 MCP 기반 시스템을 도입할 때, 기술적 사장(死藏) 위험이나 지원 중단 리스크를 최소화할 수 있음을 의미합니다.

또한, Linux Foundation은 Kubernetes, OpenAPI 등 이미 여러 글로벌 표준을 성공적으로 관리해온 경험이 있습니다. 이러한 OSS 재단의 거버넌스 모델은 엔터프라이즈에 있어 조달의 신뢰성, 장기 투자 안정성, 그리고 기술 생태계의 지속 가능한 발전을 보장하는 중요한 기준이 됩니다. 예를 들어, 클라우드 네이티브 환경에서 Kubernetes가 사실상 표준이 된 것과 마찬가지로, MCP

역시 오픈 거버넌스 체계 하에서 업계 표준으로 자리잡을 가능성이 높아집니다.

재단/프로토콜	주요 관리 기술/표준	엔터프라이즈 의미
Linux Foundation	MCP, Kubernetes, OpenAPI	벤더 중립, 장수성, 투명성
Cloud Native Computing Foundation (CNCF)	Kubernetes, Prometheus	클라우드 네이티브 표준
Apache Foundation	Kafka, Hadoop	대규모 데이터/분산 컴퓨팅

이와 같이 MCP의 OSS 재단 이관은 단순한 기술 채택을 넘어, 엔터프라이즈 IT 조달 및 장기 운영 전략의 핵심 축으로 자리잡고 있습니다.

## 7.2 메모리 스택 — Vector DB와 Neo4j GraphRAG + Memory Provider의 역할 분담

AI Agent의 실질적 성능과 확장성은 메모리 계층의 설계에 달려 있습니다. 2026년 기준, 엔터프라이즈 Agent 스택에서는 Vector DB와 GraphDB(Neo4j) 기반의 2-provider 모델이 표준으로 자리잡았습니다. 이 절에서는 4계층 메모리 모델(Short-term, Long-term, Episodic, Semantic)을 실제 DB 제품군에 매핑하고, Microsoft Agent Framework와 Neo4j의 통합 구조, 그리고 멀티 Agent 지식 공유의 구체적 구현 방법을 설명합니다.

엔터프라이즈 환경에서는 단순히 대용량 데이터를 저장하는 것만으로는 충분하지 않습니다. AI Agent가 다양한 업무 시나리오에서 의미 있는 정보를 신속하게 검색하고, 관계 기반의 복잡한 지식 추론을 수행할 수 있으려면, 메모리 계층의 구조적 설계가 필수적입니다. Vector DB는 임베딩 기반 의미 검색을, GraphDB는 관계 기반 추론과 설명 가능성을 담당하며, 두 계층의 조합은 엔터프라이즈 AI의 신뢰성과 확장성을 크게 높여줍니다. 본 절에서는 각 DB의 역할, 선택 기준, 그리고 실제 적용 사례를 중심으로 엔터프라이즈 메모리 스택의 최적 설계 방안을 제시합니다.

## 7.2.1 Pinecone · Weaviate · Qdrant · Chroma — 의미 검색과 재순위의 현재 선택지

Long-term 메모리 계층에서 RAG(Retrieval-Augmented Generation)와 결합되는 Vector DB는 대규모 의미 기반 검색과 재순위를 담당합니다. Pinecone, Weaviate, Qdrant, Chroma가 대표적인 선택지로 자리잡았으며, 각 제품은 다음과 같은 특성을 가집니다.

제품명	배포 방식	온프레미스 지원	주요 특징
Pinecone	SaaS	X	대규모 분산, 관리형, 쉬운 확장성
Weaviate	SaaS/On-prem	O	하이브리드, Graph/Vector 통합
Qdrant	OSS/SaaS	O	경량, Rust 기반, 온프레미스 최적화
Chroma	OSS	O	경량, 빠른 임베딩, 로컬 개발 용이

(출처: [Neo4j Agent Framework 통합 사례](#))

Vector DB는 임베딩 기반 의미 검색을 통해, Agent가 방대한 문서·지식에서 관련 정보를 빠르게 추출할 수 있도록 지원합니다. 재순위 기능은 검색된 결과의 적합도를 동적으로 조정하여, LLM 응답 품질을 극대화합니다. 온프레미스 요건이 있는 엔터프라이즈는 Qdrant/Chroma를, SaaS 기반 확장성이 필요한 경우 Pinecone/Weaviate를 선호합니다.

Vector DB의 선택은 엔터프라이즈의 인프라 요건, 보안 정책, 운영 규모에 따라 달라집니다. Pinecone은 완전 관리형 SaaS로, 대규모 데이터셋과 글로벌 확장에 최적화되어 있습니다. Weaviate는 온프레미스와 SaaS를 모두 지원하며, 그래프와 벡터 검색을 통합할 수 있어 복합적인 검색 요구에 유리합니다. Qdrant와 Chroma는 오픈소스 기반으로, 자체 데이터센터나 프라이빗 클라우드에 배포할 수 있어 데이터 주권과 보안이 중요한 조직에 적합합니다. 특히 Qdrant는 Rust로 구현되어 성능과 경량화가 강점이며, Chroma는 빠른 임베딩 처리와 로컬 개발에 용이합니다.

실제 엔터프라이즈 도입 사례를 보면, 금융권이나 공공기관 등 데이터 유출 리스크가 큰 조직에서는 온프레미스 배포가 가능한 Qdrant, Chroma를 선호하는 경향이 있습니다. 반면, 글로벌 서비스 확장과 유지보수 효율성이 중요한 IT 기업이나 스타트업에서는 Pinecone, Weaviate와 같은 SaaS 기반 솔루션을 선택하여, 인프라 관리 부담을 최소화하고 빠른 확장성을 확보합니다.

구분	Pinecone	Weaviate	Qdrant	Chroma
SaaS	O	O	O	X
온프레미스	X	O	O	O
확장성	매우 높음	높음	중간	낮음
OSS	X	O	O	O
특징	관리형	하이브리드	경량	초경량

### 7.2.2 Neo4j GraphRAG Context Provider — vector + fulltext + hybrid + Cypher traversal

Neo4j의 GraphRAG Context Provider는 Vector DB의 한계를 보완하기 위해 도입된 하이브리드 메모리 계층입니다. Vector, Fulltext, Hybrid, Cypher Traversal을 조합하여, 단순 유사도 기반 검색을 넘어 관계 기반의 풍부한 컨텍스트를 제공합니다. 2026년 4월, Microsoft Agent Framework에 공식 통합되며, 규제·금융·의료 등 관계성이 중요한 도메인에서 표준으로 자리잡았습니다.

기능 구분	지원 방식
Vector	임베딩 기반 의미 검색
Fulltext	키워드/문장 기반 검색
Hybrid	Vector+Fulltext 결합
Cypher	그래프 탐색 및 관계 추론

(출처: [Neo4j GraphRAG 공식 문서](#))

GraphRAG는 단순 검색을 넘어, 데이터 간의 관계와 의존성을 명확히 추적할 수 있습니다. 이는 규제 산업에서 요구하는 “설명 가능한 AI(Explainable AI)”와 데이터 감사(audit) 요건을 충족시키는 핵심 기능입니다.

GraphRAG의 도입은 엔터프라이즈 AI 시스템에서 데이터 신뢰성과 투명성을 크게 향상시킵니다. 예를 들어, 금융 분야에서는 고객 거래 이력과 연관된 다양한 엔티티(계좌, 거래, 규정 등) 간의 관계를 명확히 추적할 수 있어, 이상 거래 탐지나 규제 준수 보고에 활용됩니다. 의료 분야에서는

환자, 진료, 약물 간의 복잡한 관계를 그래프 형태로 관리함으로써, 임상 의사결정 지원 시스템의 설명 가능성과 감사 대응력을 강화할 수 있습니다.

또한, GraphRAG는 Cypher 쿼리를 활용한 그래프 탐색 기능을 제공하여, 단순한 유사도 검색을 넘어 복잡한 관계 기반 질의가 가능합니다. Hybrid 모드는 Vector와 Fulltext 검색을 결합하여, 의미 기반과 키워드 기반 검색의 장점을 모두 활용할 수 있습니다. 이러한 구조는 대규모 조직에서 다양한 검색 요구와 감사 요건을 동시에 충족시키는 데 매우 효과적입니다.

구분	Vector-only RAG	GraphRAG (Neo4j)
검색 방식	임베딩 유사도	임베딩+그래프 탐색
관계 추적	불가	가능
설명 가능성	낮음	높음
도입 분야	일반 문서, FAQ	규제, 금융, 의료, 복잡 도메인

### 7.2.3 Neo4j Memory Provider — 자동 entity 추출과 멀티 Agent 지식 공유

Neo4j Memory Provider는 세션 간 persistent memory를 제공하며, 자동 entity 추출 기능을 통해 지식 그래프를 구축합니다. 이 구조는 멀티 Agent 환경에서 한 Agent가 추출한 entity(개체)가 즉시 다른 Agent에 공유될 수 있도록 지원하여, 조직 전체의 지식 축적과 재활용을 가속화합니다.

Neo4j Memory Provider의 핵심은 대화나 작업 중에 자동으로 엔티티를 추출하고, 이를 그래프 형태로 저장함으로써, Agent 간의 지식 공유와 협업을 극대화하는 데 있습니다. 예를 들어, 한 Agent가 고객과의 대화에서 새로운 프로젝트 정보를 추출하면, 해당 엔티티가 즉시 그래프에 반영되어 다른 Agent가 이를 활용할 수 있습니다. 이러한 구조는 조직 내 지식의 사일로 현상을 해소하고, 업무 효율성을 크게 높여줍니다.

멀티 Agent 지식 공유를 구현할 때는 권한 격리(Access Control)가 매우 중요합니다. Neo4j는 프로젝트별, 팀별로 엔티티 접근 권한을 세분화할 수 있어, 민감 정보는 제한적으로 공유하고, 공통 지식은 전체 Agent가 활용하도록 설계할 수 있습니다. 예를 들어, 인사팀과 재무팀이 각각의 업무 데이터를 관리하되, 회사 전체 정책이나 공지사항 등은 공유 그래프를 통해 실시간으로 동기화할 수 있습니다.

또한, 자동 엔티티 추출 기능은 LLM 기반의 자연어 처리 엔진과 연동되어, 대화나 문서에서 인물, 조직, 날짜, 사건 등 다양한 엔티티를 실시간으로 식별하고 그래프에 반영합니다. 이 과정에서 중복 엔티티 통합, 엔티티 간 관계 자동 생성 등 고급 기능도 지원되어, 지식 그래프의 품질과 활용도가 크게 향상됩니다.

패턴 유형	설명
단일 Agent	개별 세션/지식 그래프
멀티 Agent 공유	동일 그래프 내 엔티티 실시간 공유
권한 격리	프로젝트/팀별 엔티티 접근 제어
자동 추출	대화/작업 중 엔티티 자동 추출 및 그래프화

(출처: [Neo4j Memory Provider 공식 문서](#))

## 7.3 Observability와 Durable Workflow — LangSmith · LangFuse · OpenTelemetry · Temporal

AI Agent를 실제 운영 환경에 올릴 때 가장 중요한 두 가지는 관측성(Observability)과 내구성(Durability)입니다. LangSmith/LangFuse와 같은 관측 툴, OpenTelemetry 기반 분산 추적, 그리고 Temporal+LangGraph의 2계층 워크플로우는 엔터프라이즈에서 운영 안정성을 확보하는 핵심 요소입니다. 이 절에서는 각 레이어의 특징과 선택 기준, 그리고 실제 아키텍처 적용 시 고려해야 할 포인트를 설명합니다.

엔터프라이즈 AI Agent의 운영 환경에서는 단순한 모델 성능뿐 아니라, 시스템의 상태를 실시간으로 관측하고 장애 발생 시 신속하게 복구할 수 있는 내구성이 필수적입니다. 관측성은 문제 진단, 품질 개선, 규제 대응의 기반이 되며, 내구성은 장기 실행 워크플로우와 장애 복구, 데이터 일관성 보장에 핵심 역할을 합니다. 본 절에서는 관측 툴의 선택 기준, 분산 추적의 구현, 그리고 Durable Workflow의 실제 적용 방안을 구체적으로 다룹니다.

### 7.3.1 LangSmith 대 LangFuse — SaaS 공식과 OSS OpenTelemetry 네이티브의 선택

LangSmith는 LangChain 공식 SaaS로, 트레이스, 평가, 프롬프트 관리 등 엔터프라이즈급 관측 기능을 제공합니다. 반면 LangFuse는 오픈소스(OSS)로, OpenTelemetry 네이티브 지원과 Claude Agent SDK 통합 등 유연성과 확장성이 강점입니다. 선택 기준은 데이터 주권(온프레미스 필요 여부), 커스터마이징, 비용 구조 등입니다.

LangSmith는 클라우드 기반의 관리형 서비스로, 빠른 도입과 운영 편의성이 큰 장점입니다. 엔터프라이즈에서는 별도의 인프라 구축 없이 즉시 관측 기능을 사용할 수 있으며, 트레이스 시각화, 프롬프트 버전 관리, 응답 품질 평가 등 다양한 부가 기능을 제공합니다. 그러나 데이터가 외부 클라우드에 저장되기 때문에, 보안이나 규제 요건이 엄격한 조직에서는 도입에 제약이 있을 수 있습니다.

반면, LangFuse는 오픈소스 기반으로, 온프레미스·클라우드·로컬 등 다양한 배포 옵션을 지원합니다. OpenTelemetry(OTEL) 네이티브 포맷을 채택하여, 기존 엔터프라이즈 모니터링 시스템과의 연동이 용이하며, Claude Agent SDK 등 다양한 프레임워크와의 통합이 가능합니다. 커스터마이징 자유도가 높아, 조직별 요구에 맞는 기능 확장이나 UI/UX 개선이 가능합니다. 비용 측면에서도 자체 운영이 가능하므로, 대규모 환경에서 장기적으로 TCO(총 소유 비용)를 절감할 수 있습니다.

실제 엔터프라이즈 도입 사례를 보면, 금융권이나 공공기관 등 데이터 주권이 중요한 조직에서는 LangFuse/OTEL 계열을, 스타트업이나 IT 서비스 기업 등 빠른 도입과 운영 효율성이 중요한 조직에서는 LangSmith(SaaS)를 선호하는 경향이 뚜렷합니다.

구분	LangSmith (SaaS)	LangFuse (OSS)
배포 방식	SaaS(클라우드)	온프레미스/클라우드/로컬
관측 포맷	자체 포맷	OpenTelemetry 네이티브
통합	LangChain 공식	Claude Agent SDK, OTEL
커스터마이징	제한적	매우 자유로움
데이터 주권	외부 관리	사내 직접 관리 가능

(출처: [LangFuse 통합 문서](#))

### 7.3.2 OpenTelemetry for MCP — tool discovery / tool call / backend latency 분산 추적

MintMCP의 OpenTelemetry(OTEL) 가이드는 MCP 기반 Agent 운영에서 tool discovery, tool call, backend MCP server latency 등 핵심 이벤트를 분산 추적으로 수집하는 패턴을 제시합니다. OTEL은 Jaeger, Grafana 등 기존 엔터프라이즈 APM(애플리케이션 성능 모니터링) 자산과의 연동이 용이해, 기존 인프라를 재활용할 수 있습니다.

OTEL 기반 분산 추적은 엔터프라이즈 AI 시스템의 가시성을 크게 향상시킵니다. 예를 들어, Agent가 외부 MCP 서버에 도구 목록을 조회하거나(tool.discovery), 특정 도구를 호출할 때(tool.call), 또는 MCP 서버의 응답 지연이 발생할 때(backend.latency), 각각의 이벤트가 트race 스패스로 기록되어 문제 발생 시 신속한 원인 분석이 가능합니다. agent.decision 스패스는 Agent 내부의 의사결정 과정이나 상태 전이 정보를 기록하여, 복잡한 멀티스텝 워크플로우의 디버깅과 품질 개선에 활용됩니다.

OTEL의 가장 큰 장점은 표준화된 포맷과 다양한 백엔드 지원입니다. 엔터프라이즈에서는 Jaeger, Grafana, Datadog 등 기존 모니터링 시스템과의 연동을 통해, AI Agent와 기존 IT 시스템의 통합 관측 체계를 구축할 수 있습니다. 이는 신규 시스템 도입 시 이질적인 모니터링 솔루션을 추가할 필요 없이, 기존 인프라를 효율적으로 재활용할 수 있음을 의미합니다.

스팬 유형	설명
tool.discovery	도구 목록 조회(Agent → MCP)
tool.call	도구 실행(Agent → MCP → 외부 시스템)
backend.latency	MCP 서버 응답 시간 측정
agent.decision	Agent 내 의사결정/상태 전이

(출처: [MintMCP OTEL 가이드](#))

OTEL 기반 관측은 엔터프라이즈 전체의 모니터링/알림/분석 체계를 일원화하여, 장기적으로 TCO(총 소유 비용)를 크게 절감할 수 있습니다. 신규 Agent 시스템 도입 시 OTEL 연동 여부를 필수 체크 항목으로 삼아야 합니다.

### 7.3.3 Temporal + LangGraph — retry/state/failure recovery와 prompt/tool/memory의 2계층

Temporal과 LangGraph의 2계층 패턴은 장기 실행·내구성이 요구되는 엔터프라이즈 Agent의 표준 설계입니다. Temporal은 retry, state management, failure recovery 등 워크플로우의 내구성을, LangGraph는 prompt/tool/memory 등 AI Agent의 논리적 상태와 실행을 담당합니다.

2계층 Durable Workflow 구조는 엔터프라이즈에서 복잡한 비즈니스 프로세스 자동화, 장애 복구, 장기 트랜잭션 관리 등에 필수적입니다. Temporal은 워크플로우의 상태를 지속적으로 관리하며, 장애 발생 시 자동 재시도, 체크포인트 복구, 실패 알림 등 내구성 기능을 제공합니다. LangGraph는 Agent의 논리적 상태(프롬프트, 도구, 메모리 등)를 그래프 형태로 관리하여, 멀티스텝 작업의 흐름을 시각적으로 설계하고, 인간 개입(HITL: Human-In-The-Loop)이나 조건부 분기 등 복잡한 로직을 손쉽게 구현할 수 있습니다.

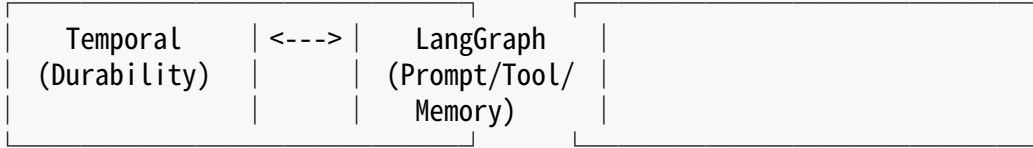
실제 엔터프라이즈 적용 사례에서는, Temporal이 없는 조직에서는 Prefect, Airflow 등 대안 워크플로우 엔진을 활용하기도 합니다. 그러나 Temporal+LangGraph 조합은 내구성, 확장성, 멀티 Agent 조정 등에서 가장 높은 신뢰도를 제공합니다. 예를 들어, 금융권에서는 장기 트랜잭션(예: 대출 심사, 보험 청구 등)에서 장애 발생 시 자동 복구와 상태 일관성이 필수이며, 제조·물류 분야에서는 복잡한 공급망 관리 워크플로우를 안정적으로 운영하기 위해 2계층 구조가 도입되고 있습니다.

계층	담당 영역	주요 기능
Temporal	워크플로우 관리, 내구성	retry, state, failure recovery
LangGraph	Agent 논리, 프롬프트/툴/메모리	상태 기반 그래프, 체크포인팅, HITL

(출처: [Temporal+LangGraph 2계층 구조](#))

#### Temporal+LangGraph 2계층 다이어그램

plaintext



## 8장: 구축·운영 사례와 실패를 막는 5대 안전장치

### 8.1 성공 시나리오 — 코딩 · 고객지원 · RPA 대체 · 개인업무 · 연구분석

#### 5대 활용

AI Agent의 실제 성공 시나리오를 분석하면, 코딩 자동화, 고객지원, RPA 대체, 개인업무, 연구 분석 등 5가지 주요 활용 축이 도출됩니다. 각 축은 대표적인 수치와 사례를 통해 도입 유형별 성공 기준을 명확히 제시합니다. 이 절에서는 SWE-bench 등 코딩 성과, Klarna의 고객지원 혁신, Lindy와 Perplexity의 개인 및 연구분석 Agent 등 실무 적용 사례를 근거로, IT 의사결정자가 도입 ROI와 KPI를 구체적으로 판단할 수 있도록 합니다.

#### 8.1.1 Claude Code 80.9% · Devin 13.86% on SWE-bench — 코딩 Agent가 만드는 PR 단위 기여

코딩 분야에서 AI Agent의 도입은 단순한 코드 자동화 수준을 넘어, 실제 개발 프로세스의 핵심 단계인 PR(Pull Request) 생성과 머지까지 자동화하는 수준으로 진화하고 있습니다. SWE-bench와 같은 벤치마크는 이러한 Agent의 실질적 기여도를 수치로 평가할 수 있게 해주며, 이를 통해 조직은 도입 효과를 명확히 검증할 수 있습니다. 최근에는 Claude Code, Devin, Cursor, GitHub Copilot Workspace 등 다양한 코딩 Agent가 등장하며, 이들의 성과는 기존 LLM 기반 보조 도구 대비 월등히 향상된 것으로 나타나고 있습니다. 이러한 변화는 개발팀의 생산성과 코드 품질, 그리고 자동화 범위 확대에 직접적인 영향을 미치고 있습니다.

##### SWE-bench 랭킹 상위 Agent 성과

SWE-bench는 AI 코딩 Agent의 실제 PR 단위 기여를 수치로 평가하는 대표 벤치마크입니다. Claude Code는 SWE-bench에서 80.9%의 성공률을 기록하며, 기존 단일 LLM 호출 기반의 5% 미만 성과를 압도적으로 상회합니다. Claude 4 Sonnet은 SWE-bench Lite에서 60.3%를 달성했고, Devin은 초기 버전에서 13.86%로 이전 SOTA(4.80%) 대비 2.9배 향상된 결과를 보여줬습니다. Cursor, GitHub Copilot Workspace, OpenHands(구 OpenDevin)도 주요 코딩 Agent 제품으로 꼽히며, Devin은 실제로 GitHub PR 생성과 머지까지 자동화합니다.

Agent	SWE-bench (%)	SWE-bench Lite (%)	PR 생성/머지 지원
Claude Code	80.9	-	○
Claude 4 Sonnet	-	60.3	○
Devin	13.86	-	○
OpenHands	-	-	○
Cursor	-	-	○

## PR 단위 기여의 의미

코딩 Agent는 단순 보조도구를 넘어 실제 PR 단위로 업무에 기여하는 수준으로 진화하고 있습니다. SWE-bench 점수는 실무에서 코드 품질, 자동화 범위, 리뷰 프로세스와 직접 연결되어야 하며, Agent의 성과를 실 업무 기여로 번역할 때 반드시 코드 리뷰와 승인 프로세스를 포함해야 합니다. CTO는 Agent의 도입을 단순 자동화가 아닌 “기여자”로 평가하는 기준을 마련해야 합니다.

## 실무 적용 시 유의점

코딩 Agent 도입 시, SWE-bench 등 객관적 벤치마크를 활용해 성과를 측정하고, PR 생성·머지 자동화가 실제 개발팀 워크플로우에 어떻게 통합되는지 명확히 정의해야 합니다. Agent의 코드 품질과 보안, 리뷰 프로세스와의 연계가 성공의 핵심입니다. 또한, 도입 초기에는 Agent가 생성한 PR에 대한 수동 리뷰와 검증 절차를 강화하여, 코드 품질 저하나 보안 취약점이 발생하지 않도록 해야 합니다. 장기적으로는 Agent의 성능 개선과 워크플로우 최적화를 통해, 개발팀의 생산성 향상과 비용 절감 효과를 극대화할 수 있습니다. SWE-bench와 같은 벤치마크 결과를 도입 의사결정의 핵심 지표로 삼고, 실제 업무 적용 시에는 조직의 개발 문화와 프로세스에 맞는 커스터마이징이 필요합니다.

## 8.1.2 Klarna 월 2.3M 대화 · 700 FTE · 11분→2분 — 고객지원과 RPA 하이브리드

고객지원(CX)과 RPA(로봇 프로세스 자동화) 분야에서 AI Agent의 도입은 대규모 인력 대체와 처리 속도 향상, 그리고 운영 효율성 증대로 이어지고 있습니다. Klarna의 사례는 AI 기반 CX Agent와 RPA의 결합이 실제로 어떤 정량적 성과를 창출하는지 잘 보여줍니다. 특히, 반복적이고 표준화된 고객 문의를 자동화함으로써, 기존 인력의 부담을 크게 줄이고, 고객 응대 속도와 만족도를 동시에 높일 수 있음을 입증했습니다. 이처럼 AI Agent와 RPA의 하이브리드 모델은 금융, 유통, 서비스 등 다양한 산업에서 빠르게 확산되고 있으며, 도입 ROI와 KPI 개선 효과가 명확하게 드러나고 있습니다.

### Klarna 사례의 정량적 성과

Klarna는 OpenAI 기반 CX Agent를 도입해 월 2.3M 대화를 처리하며, 700 FTE(Full-Time Equivalent) 워크로드를 대체했습니다. 결제, 환불, 배송 이슈 해결시간은 11분에서 2분으로 단축되었으며, Acuvate 보고에 따르면 비정형 송장 추출(AI)과 ERP 입력(RPA) 하이브리드로 기업 운영 효율이 10~20% 매출 ROI로 이어졌습니다.

지표	도입 전	도입 후	개선율
월 대화 건수	2.3M	2.3M	-
FTE 워크로드	700	0	100%
이슈 해결시간	11분	2분	-81.8%
매출 ROI	-	10~20%	-

### CX Agent + 사람 escalation 결합 모델

Klarna의 성공은 CX Agent와 사람 escalation 경로의 결합에 있습니다. 자동화된 Agent가 반복적이고 단순한 문의를 처리하며, 복잡하거나 신뢰 임계치(Trust Threshold) 이하의 케이스는 사람에게 즉시 이관됩니다. 이 구조는 고객 만족도(CSAT)와 운영 효율을 동시에 달성하는 표준 모델로 자리잡았습니다.

### CSAT 추이 모니터링의 중요성

Agent 도입 후 CSAT, NPS 등 고객 만족 지표를 실시간으로 모니터링해야 합니다. 자동화가

CX를 훼손하지 않는지 지속적으로 검증하고, 필요시 사람 escalation 경로를 강화하는 것이 성공의 핵심입니다. 또한, 자동화 범위가 확대될수록 고객 경험의 질이 저하되지 않도록, 정기적인 품질 점검과 피드백 루프를 마련해야 합니다. RPA와 AI Agent의 결합은 단순 반복 업무의 자동화뿐만 아니라, 비정형 데이터 처리와 예외 상황 대응에서도 큰 효과를 발휘하므로, 도입 시에는 업무 프로세스 전반을 재설계하는 것이 바람직합니다. Klarna 사례처럼, 정량적 성과와 더불어 고객 만족도 변화까지 체계적으로 관리하는 것이 장기적 성공의 관건입니다.

### 8.1.3 Lindy · Perplexity Deep Research — 개인업무와 연구분석 Agent 의 현재

개인업무 자동화와 연구·데이터 분석 분야에서 AI Agent의 도입은 업무 효율성과 비용 절감, 그리고 의사결정의 신속성 측면에서 큰 변화를 이끌고 있습니다. Lindy, Perplexity, OpenAI/Anthropic Deep Research 등 다양한 Agent가 실제 현업에서 활용되고 있으며, 산업별 도입률과 ROI가 점차 구체적으로 집계되고 있습니다. 특히, 보험, 금융, 데이터 분석 등 정보 집약적 산업에서는 Agent 도입이 업무 프로세스 혁신의 핵심 동인으로 작용하고 있습니다. 이 절에서는 산업별 도입률, 비용 절감 효과, 그리고 실제 적용 전략을 살펴봄으로써, 조직별 최적의 Agent 활용 방안을 제시하고자 합니다.

#### 산업별 Agent 도입률 및 ROI

Lindy는 Claude 기반 메일, 스케줄, 세일즈 자동화로 Top 10 AI Agent Companies 2026에 선정되었습니다. 연구/데이터 분석 Agent는 기업 응답자의 45%가 데이터 분석, 40%가 시장 조사 용도로 활용하고 있습니다(150+ AI Agent Statistics 2026). 보험 산업에서는 agentic 도입률이 2026년 48%에 달하며, 비용 절감 효과는 56%로 보고됩니다. OpenAI/Anthropic Deep Research, Perplexity Agent 등도 개인 및 연구분석 분야에서 활발히 사용되고 있습니다.

산업/분야	Agent 도입률(%)	비용 절감(%)	주요 활용
보험	48	56	자동화, 분석
데이터 분석	45	-	분석, 보고서
시장 조사	40	-	조사, 인사이트
개인업무	-	-	메일, 스케줄

#### 실제 수치 기반 ROI 검증

개인업무 및 연구분석 Agent 도입은 실제 수치와 ROI를 통해 검증되어야 합니다. CX 단독 자동화와 내부 업무 자동화의 우선순위를 분명히 하고, 산업별 도입률과 비용 절감 효과를 근거로 도입 전략을 수립해야 합니다. 예를 들어, 보험 산업의 경우 48%의 도입률과 56%의 비용 절감 효과는 Agent 도입이 단순한 트렌드가 아니라, 실질적인 경영 성과로 이어지고 있음을 보여줍니다. 데이터 분석과 시장 조사 분야에서도 Agent 활용이 빠르게 확산되고 있으며, 이는 업무 효율성 증대와 인력 운영의 유연성 확보로 이어지고 있습니다.

### 도입 전략 및 유의점

Agent 도입 시, 산업별 특성에 맞는 활용 시나리오를 선정하고, 실제 ROI와 KPI를 기준으로 성공 여부를 판단해야 합니다. 내부 자동화 우선 전략은 국내 기업에서 선행 패턴으로 자리잡고 있습니다. 또한, 도입 전에는 파일럿 프로젝트를 통해 예상 성과와 리스크를 사전에 검증하고, 도입 후에는 정량적·정성적 지표를 지속적으로 모니터링해야 합니다. 산업별로 요구되는 보안, 데이터 프라이버시, 업무 프로세스 표준화 등 추가 고려사항도 반드시 점검해야 하며, 이를 통해 Agent 도입의 성공 가능성을 높일 수 있습니다.

## 8.2 국내·해외 적용 사례 — Klarna 후퇴 · LG전자 CHATDA · KB라이프 · 한화 · Sierra · Devin

해외와 국내의 AI Agent 적용 사례는 성공과 조정, 내부 자동화, B2B SaaS 제품화 등 다양한 경로를 보여줍니다. Klarna의 후퇴 사례, LG전자·KB라이프·한화·SK이노베이션의 내부 지식 자동화, Devin·Sierra·Lindy의 제품화 경로를 통해 각 산업 유형별 도입 전략과 교훈을 도출할 수 있습니다. 이 절은 의사결정자가 자신의 조직에 맞는 적용 모델을 선택할 수 있도록 실제 사례와 KPI 변화를 비교합니다.

### 8.2.1 Klarna “AI 실패” 선언의 실상 — Trust Threshold 와 Uber-type 하이브리드 전환

Klarna의 AI Agent 도입과 운영 사례는 완전 자동화의 한계와 하이브리드 모델의 필요성을 명확히 보여주는 대표적 사례입니다. 초기에는 대규모 인력 대체와 운영비 절감을 실현했으나, 고객 경험(CX) 저하와 신뢰 임계치(Trust Threshold) 문제로 인해 전략적 전환을 단행하게 되었습니다. 이

과정에서 얻은 교훈은 AI Agent 도입 시 자동화와 사람 개입의 균형, 그리고 실시간 KPI 모니터링의 중요성입니다. 본 절에서는 Klarna의 KPI 변화와 실패 패턴, 그리고 실무 적용 시 유의점을 구체적으로 살펴봅니다.

### Klarna 전후 KPI 변화 및 교훈

2025년 5월 Klarna CEO는 “cost containment가 CX를 일부 훼손했다”는 점을 시인하며, 완전 자동화에서 Uber-type 하이브리드 인력풀로 전환했습니다. \$60M 운영비 절감(2025년 말 853 FTE 환산)은 유지되었지만, Trust Threshold(신뢰 임계치) 밑에서는 반드시 사람 escalation 이 필요하다는 교훈을 남겼습니다.

지표	도입 전	도입 후	변화
운영비	-	\$60M 절감	-
FTE 환산	-	853 감소	-
CX 훼손	없음	일부 발생	-
Trust Threshold	미정	명확화	-

### 완전 자동화 실패 패턴의 구조적 이해

Klarna 사례는 완전 자동화가 반드시 성공을 보장하지 않으며, 고객 신뢰 임계치 아래에서는 사람 개입이 필수임을 보여줍니다. CSAT, NPS 등 실시간 KPI를 운영에 포함하고, 자동화와 사람 escalation의 균형을 구조적으로 설계해야 합니다. 또한, 완전 자동화로 인한 고객 불만이나 서비스 품질 저하가 발생할 경우, 신속하게 하이브리드 모델로 전환할 수 있는 유연한 운영 체계가 필요합니다. Klarna의 경험은 AI Agent 도입 시, 비용 절감과 고객 경험 사이의 균형을 어떻게 맞출 것인지에 대한 중요한 시사점을 제공합니다.

### 실무 적용 시 유의점

Agent 도입 시, 완전 자동화 정책을 피하고, 하이브리드 인력풀과 Trust Threshold 기반 escalation 경로를 반드시 마련해야 합니다. KPI 변화와 고객 경험을 지속적으로 모니터링하는 것이 성공의 핵심입니다. 또한, 도입 초기에는 고객 피드백을 적극적으로 수집하고, 문제가 발생할 경우 신속하게 대응할 수 있는 조직 내 프로세스를 구축해야 합니다. 장기적으로는 자동화와 사람 개입의 최적 비율을 찾아가는 것이 중요하며, 이를 위해 정기적인 성과 분석과 전략 조정이 필요합니다.

## 8.2.2 LG전자 CHATDA · KB라이프 · 한화 Copilot Studio · SK이노베이션 — 국내 내부 지식 자동화 패턴

국내 대기업들은 AI Agent를 활용한 내부 지식 자동화에 선도적으로 나서고 있습니다. 문서 요약, 회의록 자동화, 빅데이터 분석, 환경규제 대응 등 다양한 영역에서 Agent가 도입되고 있으며, 이를 통해 업무 효율성과 정확성이 크게 향상되고 있습니다. 각 기업은 자체적인 파일럿 프로젝트를 통해 성공 사례를 축적하고, 점진적으로 적용 범위를 확대하는 전략을 채택하고 있습니다. 본 절에서는 국내 4개 기업의 대표 유즈케이스와 자동화 범위, 그리고 도입 전략을 비교 분석합니다.

### 국내 4개사 유즈케이스 비교

Microsoft Source Asia 2025-09 보고에 따르면 KB라이프는 M365 Copilot을 전사 배포하여 문서 요약과 회의록 자동화를 실현했습니다. LG전자 HS는 Azure OpenAI 기반 CHATDA agentic을 적용해 “분석 계획 수립→코드 작성→분석→인사이드 보고서” 파이프라인을 자동화했습니다. 한화는 Copilot Studio 기반 정기회의·환경규제 Agent를, SK이노베이션은 Azure 기반 엔지니어링 문서 자동화를 도입했습니다. LG전자 Agentic Report Automation은 AWS Tech Blog에 기술되었습니다.

기업	적용 영역	주요 기술	자동화 범위
KB라이프	문서 요약/회의록	M365 Copilot	전사 배포, 내부 지식 자동화
LG전자 HS	빅데이터 분석	Azure OpenAI CHATDA	분석 계획~보고서 파이프라인
한화	환경규제/회의	Copilot Studio	정기회의, 환경규제 Agent
SK이노베이션	엔지니어링 문서	Azure OpenAI	문서 자동화

### 내부 지식 자동화 선행 패턴

국내 기업들은 내부 지식 자동화에 우선 집중하고 있습니다. 사내 문서, 보고서, 회의록 등 반복적이고 정형화된 업무를 Agent로 자동화하며, Iconic Use Case 1개부터 시작해 ROI를 입증한 뒤 확장하는 전략이 정식으로 자리잡고 있습니다. 이처럼 단계적 도입과 성공 사례 축적은 조직 내 변화 저항을 최소화하고, Agent 도입의 효과를 빠르게 확인할 수 있는 장점이 있습니다.

### 도입 전략 및 유의점

Agent 도입 시, 내부 지식 자동화부터 시작해 성공 사례를 쌓고, 점진적으로 외부 고객지원

등 복잡한 영역으로 확장하는 것이 바람직합니다. 각 기업의 적용 사례와 자동화 범위를 비교해 조직에 맞는 전략을 수립해야 합니다. 또한, 도입 전에는 업무 프로세스 분석과 파일럿 테스트를 통해 예상 효과와 리스크를 사전에 파악하고, 도입 후에는 성과 모니터링과 지속적인 개선 활동을 병행해야 합니다. 내부 지식 자동화는 보안과 데이터 프라이버시 이슈도 함께 고려해야 하며, 이를 위해 IT·보안 부서와의 긴밀한 협업이 필요합니다.

### 8.2.3 Devin · Sierra AI · Lindy — B2B SaaS 로 제품화된 3가지 Agent

AI Agent의 B2B SaaS 제품화는 조직이 자체 구축(Build)과 상용 제품(Buy) 중에서 선택할 수 있는 다양한 옵션을 제공합니다. Devin, Sierra AI, Lindy와 같은 제품은 개발, CX, 개인업무 등 각기 다른 산업과 업무 영역에 특화된 기능을 제공하며, 도입 기업은 필요에 따라 최적의 솔루션을 선택할 수 있습니다. 본 절에서는 세 가지 대표 제품의 주요 기능과 자동화 범위, 그리고 도입 전략을 비교 분석합니다.

#### Devin/Sierra/Lindy 제품 축 비교

Devin(Cognition AI, 2024-03)은 AI 엔지니어 루프를 구현해 GitHub PR 생성, 피드백 반영, 최종 머지까지 자동화합니다. Sierra AI는 AgentOS와 Experience Manager 플랫폼을 제공하며, CX와 내부 업무 자동화를 지원합니다. Lindy는 Claude 기반 세일즈, 스케줄 SaaS로 개인업무 자동화에 특화되어 있습니다. 각 제품은 개발, CX, 개인업무 등 서로 다른 산업의 B2B SaaS 지표를 제공합니다.

제품	산업/분야	주요 기능	자동화 범위
Devin	개발	PR 생성, 피드백, 머지	엔지니어링 워크플로우
Sierra AI	CX/업무	AgentOS, Experience Manager	CX, 내부 업무 자동화
Lindy	개인업무	세일즈, 스케줄, 메일	개인 업무 자동화

#### “빌드 vs 바이” 경로의 실제 후보군

Devin, Sierra, Lindy는 B2B SaaS 제품으로서 “바이(buy)” 경로의 실제 후보군을 제공합니다. 조직은 자체 구축(Build)과 상용 제품(Buy) 중 선택할 때, 각 제품의 산업별 자동화 범위와 ROI를 비교해 결정해야 합니다. SaaS 제품은 신속한 도입과 유지보수의 용이성이 장점이며, 자체 구축은 맞춤화와 데이터 통제 측면에서 유리할 수 있습니다.

### 도입 전략 및 유의점

Agent 제품 도입 시, 산업별 요구에 맞는 기능과 자동화 범위를 명확히 파악하고, 빌드와 배포 경로의 장단점을 비교해 최적의 선택을 해야 합니다. 제품별 KPI와 성공 사례를 근거로 도입 전략을 수립해야 합니다. 또한, 도입 전에는 PoC(Proof of Concept)를 통해 실제 업무에의 적합성을 검증하고, 도입 후에는 정기적인 성능 평가와 피드백 수집을 통해 지속적으로 개선해 나가야 합니다. SaaS 제품의 경우, 보안, 데이터 프라이버시, 커스터마이징 가능성 등도 함께 고려해야 합니다.

## 8.3 운영을 좌우하는 5대 안전장치 — Step cap · Cost cap · Tool allowlist · Injection 탐지 · Observability

AI Agent 운영에서 실패를 막는 핵심은 5대 안전장치(하드 step cap, cost cap, tool allowlist, injection 탐지, observability)입니다. 실제 사고 사례와 수치, 대응 방안을 통해 각 안전장치의 필요성과 구현 방법을 제시합니다. 이 절은 운영 릴리스 게이트와 조직 정책으로 안전장치를 강제 하는 기준을 마련합니다.

### 운영 5대 안전장치 — Agent 사고를 막는 릴리스 게이트

시스템 설정이 아닌 조직 정책으로 강제 — Kiro AWS 장애 · \$3,000 Infinite Loop Attack · 73% Injection 노출에서 배운 하드 제한

① Step Cap 하드 step 15 제한 · 무한 루프 차단  
N회 초과 시 pause & alert — Kiro AWS 장애 · Claude Code 무한 루프 버그 재발 방지

② Cost Cap 태스크당 · 일간 · 조직 예산 한도  
hallucinated 반복으로 하룻밤 \$3,000 발생 사례 — 예산 초과 시 즉각 차단 & alert

③ Tool Allowlist 등록된 도구만 호출 가능  
hallucinated tool 즉시 abort — 미등록 Tool 호출 차단이 1차 방어선

④ Prompt Injection 탐지 OWASP LLM01 대응  
Direct · Indirect(문서 · 웹 숨김 지시) 73% 노출 — 입력 검사 + 세션 권한 분리

⑤ Observability trace · 평가 · Goal Drift 탐지  
1인 리뷰어 50~100 trace/h → 1,000 req/day 처리에 10~20h 필요 · 관측 FTE 를 예산 고정 항목으로

○ 자원 한도 (Resource cap)      ○ 보안 (Security)      ○ 관측 (Observability)

[그림 4] 운영 5대 안전장치 — Agent 사고를 막는 릴리스 게이트

### 8.3.1 Step cap · Cost cap — Kiro AWS 장애 · Claude Code 무한 루프 버그에서 배운 하드 제한

AI Agent 운영에서 Step cap과 Cost cap은 필수적인 안전장치로 자리잡고 있습니다. 실제로 하드 제한이 없는 Agent는 무한 루프, 자원 고갈, 예산 초과 등 치명적인 장애를 유발할 수 있으며, 이는 조직의 운영 안정성과 비용 효율성에 심각한 영향을 미칩니다. 최근에는 Kiro AWS 장애, Claude Code 무한 루프 버그, 그리고 하룻밤 \$3,000이 발생한 Infinite Loop Attack 등 다양한 사고 사례가 보고되고 있어, 하드 제한의 중요성이 더욱 부각되고 있습니다. 본 절에서는 권장 디폴트 값과 사고 사례, 그리고 실무 적용 시 유의점을 구체적으로 살펴봅니다.

#### step/cost cap 권장 디폴트 및 사고 사례

모든 Agent run에는 하드 step cap(예: 15 step)과 cost cap이 필수입니다. N회 초과 write action 발생 시 pause 및 alert를 트리거해야 합니다. 2024~2026년 실제 사고로는 Kiro AWS 장애, Claude Code 무한 루프 버그, 하룻밤 \$3,000이 발생한 Infinite Loop Attack 등이 보고되었습니다.

안전장치	권장 디폴트	사고 사례	대응 방법
Step cap	15 step	무한 루프, 자원 고갈	실행 제한, alert
Cost cap	조직별 정책	비용 폭증, 예산 초과	예산 제한, pause

#### 캡 없는 Agent 운영 불가의 이유

Step cap과 cost cap이 없는 Agent는 무한 루프, 자원 고갈, 예산 초과 등 치명적 장애를 유발합니다. 시스템 설정이 아닌 조직 정책으로 승인하고, 운영 릴리스 게이트에 반드시 포함해야 합니다. 실제 사고 사례에서는 Agent가 의도치 않게 반복 작업을 수행하거나, 외부 API 호출이 무한 반복되어 엄청난 비용이 발생하는 경우가 있었습니다. 따라서, 모든 Agent 운영 환경에는 하드 제한을 명확히 설정하고, 초과 시 즉각적으로 알림과 실행 중단이 이루어지도록 해야 합니다.

#### 실무 적용 시 유의점

캡 값은 운영 환경, 예상 태스크 수, 예산에 따라 조정해야 하며, 하드 제한이 없으면 운영 불가임을 명확히 인식해야 합니다. 사고 사례를 근거로 정책 수립이 필수입니다. 또한, 운영 중에는 정기적으로 step/cost cap 설정을 점검하고, 이상 징후 발생 시 신속하게 대응할 수 있는 모니터링

체계를 구축해야 합니다. 조직 내에서는 Agent 운영 정책에 하드 제한 준수를 명문화하고, 신규 Agent 도입 시 반드시 사전 검증 절차를 거쳐야 합니다.

### 8.3.2 Tool allowlist 와 Prompt Injection 탐지 — OWASP LLM01 대응과 73% 노출 실태

AI Agent의 보안 운영에서 가장 큰 위협 중 하나는 Prompt Injection 공격입니다. 최근 보안 감사 결과에 따르면, 프로덕션 AI 배포 환경의 73%가 Prompt Injection에 노출되어 있는 것으로 나타났으며, 이는 조직의 데이터와 시스템 권한 탈취로 이어질 수 있는 심각한 리스크입니다. Tool allowlist와 hallucinated tool abort 등 방어 메커니즘은 이러한 공격을 효과적으로 차단하는 1차 방어선 역할을 합니다. 본 절에서는 Injection 공격 유형, 대응 방안, 그리고 실무 적용 시 유의점을 구체적으로 설명합니다.

#### Injection 공격 유형과 대응 매핑

Palo Alto Unit 42와 MDPI 리뷰에 따르면 프로덕션 AI 배포 보안 감사에서 73%가 Prompt Injection에 노출되었습니다. 간접(indirect) injection은 문서, 웹페이지 숨김 지시를 통해 시스템 프롬프트와 권한을 탈취합니다. 대응책은 tool allowlist와 hallucinated tool 즉시 abort입니다.

공격 유형	노출률	대응 방법
Direct Injection	73%	tool allowlist
Indirect Injection	73%	hallucinated tool abort
OWASP LLM01	73%	보안 감사, 탐지

#### Injection은 방어 가능한 공격

Prompt Injection은 방어 가능한 공격이며, tool allowlist가 1차 방어선입니다. 보안팀과 Agent팀이 공동 릴리스 게이트를 구성하고, 모든 tool call에 allowlist를 적용해야 합니다. 또한, hallucinated tool 호출을 즉시 중단하는 로직을 구현함으로써, 의도치 않은 권한 상승이나 외부 시스템 오용을 방지할 수 있습니다. OWASP LLM01 등 최신 보안 가이드라인을 참고하여, 정기적인 보안 감사와 취약점 점검을 병행해야 합니다.

#### 실무 적용 시 유의점

보안 감사에서 Injection 대응 여부를 필수로 확인하고, allowlist와 탐지 로직을 운영 정책에

포함해야 합니다. 간접 injection까지 탐지할 수 있는 구조가 필요합니다. 이를 위해서는 Agent의 입력 데이터와 시스템 프롬프트를 정기적으로 점검하고, 의심스러운 패턴이 탐지될 경우 즉각적으로 대응할 수 있는 자동화된 모니터링 시스템을 구축해야 합니다. 보안팀과의 협업을 통해, Agent 운영 정책에 최신 보안 위협 대응 방안을 반영하는 것이 중요합니다.

### 8.3.3 Observability 부재의 비용과 Goal Drift 대응 — 1인 리뷰어 10~20 시간의 산수

AI Agent 운영에서 Observability(관측성)는 시스템 상태, 이벤트, 오류를 실시간으로 추적하고 분석할 수 있는 능력을 의미합니다. Observability가 부족할 경우, Agent의 행동을 사전에 예측하거나 이상 징후를 신속히 감지하기 어렵기 때문에, 운영 정원(staffing) 가설이 무너지고, 장기 실행 중 목표 왜곡(Goal Drift) 현상이 발생할 수 있습니다. 본 절에서는 Observability의 필요성과 비용, 그리고 Goal Drift 대응 방안을 구체적으로 설명합니다.

#### Observability 시간 산수표 및 대응 방안

전통 APM은 상태 기반 다단계 결정트리를 감지하지 못합니다. 1인 리뷰어가 시간당 50100 trace를 처리할 수 있으며, 일일 1,000 req면 1020 시간이 필요합니다. Goal drift 대응은 per-step goal preservation 체크와 planner-only 노드 분리로 구현됩니다.

항목	처리량/시간	일일 req	필요 시간	대응 방법
리뷰어 처리량	50~100	1,000	10~20h	goal preservation
Observability FTE	-	-	예산 산정	planner-only 분리

#### 관측 없는 Agent의 운영 정원 가설 붕괴

Observability가 없는 Agent는 운영 정원(staffing) 가설부터 무너집니다. 관측 FTE를 프로젝트 예산의 고정 항목으로 산정하고, goal drift 방지 구조를 반드시 포함해야 합니다. 실제 운영 환경에서는 Agent의 행동을 실시간으로 추적하고, 이상 징후 발생 시 즉각적으로 대응할 수 있는 모니터링 시스템이 필수적입니다. Goal drift 현상은 장기 실행 중 Agent의 목표가 왜곡되어, 의도하지 않은 방향으로 업무가 진행되는 문제를 야기할 수 있으므로, per-step goal preservation 체크와 planner-only 노드 분리 등 구조적 대응이 필요합니다.

#### 실무 적용 시 유의점

관측 톨과 인력, 예산을 운영 정책에 포함하고, per-step goal preservation 체크와 planner-only 노드 분리로 장기 실행 중 목표 왜곡을 방지해야 합니다. 관측 FTE 산정 워크시트를 활용해 예산을 구체적으로 계획해야 합니다. 또한, 운영 중에는 정기적으로 Observability 시스템의 성과와 커버리지를 점검하고, 필요시 추가 인력이나 도구를 투입하여 운영 안정성을 확보해야 합니다. Observability는 단순 모니터링을 넘어, Agent의 행동을 체계적으로 분석하고, 장기적 성과와 리스크를 관리하는 핵심 역량임을 인식해야 합니다.

## 9장: 결론 및 권장사항 — IT 의사결정자를 위한 도입 체크리스트

### 9.1 프레임워크·스택 선정 체크리스트 — 규제·자동화·OpenAI 스택 3분할 기준

AI Agent 도입을 고려하는 IT 의사결정자들은 프레임워크와 스택 선정 과정에서 다양한 평가 기준과 조달 전략을 체계적으로 검토해야 합니다. 특히, 규제 산업, 내부 자동화, OpenAI 중심 스택이라는 3분할 기준은 업계에서 널리 통용되는 실질적인 분류법입니다. 이 절에서는 각 산업 유형에 적합한 프레임워크의 주요 평가 요소와, 2026년 기준으로 표준화된 5레이어 스택 조달 방식을 구체적으로 안내합니다. 또한, OSS 코어와 상용 관측/실행 플랫폼의 분리 조달 전략을 통해 라이선스 리스크와 운영 유연성 확보 방안을 제시합니다. 실제 프로덕션 채택 사례와 학습 곡선, 체크포인팅 및 HITL 기능 등 실무적 관점에서의 선택 기준을 표와 함께 정리하여, 조직별 맞춤형 도입 의사결정에 실질적인 도움을 드리고자 합니다.

#### 9.1.1 프레임워크 선정 — “규제 산업 → LangGraph, 내부 자동화 → CrewAI, OpenAI 스택 → Agents SDK”

AI Agent 프레임워크를 선정할 때는 조직의 산업 특성과 도입 목적에 따라 최적의 선택지가 달라집니다. 최근 업계에서는 규제 산업, 내부 자동화, OpenAI 중심 스택의 3분할 기준이 표준으로 자리 잡고 있습니다. 각 프레임워크는 MCP 지원 성숙도, 체크포인팅 및 HITL(사람개입) 기능,

실제 프로덕션 채택 기업, 그리고 학습 곡선 등 다양한 평가 축에서 차별화된 강점을 보입니다. 이 절에서는 각 산업 유형별로 추천되는 프레임워크의 특징을 구체적으로 비교하고, MCP 지원 방식과 실제 도입 사례, 학습 곡선 및 체크포인팅/HITL 지원 현황을 표로 정리하여 실무적 선택에 도움을 드립니다.

### 3분할 선정 기준

업계에서는 AI Agent 프레임워크 선정 시 규제 산업, 내부 자동화, OpenAI 중심 스택의 3분할 기준을 적용하는 것이 컨센서스로 자리 잡았습니다. 규제 산업에서는 LangGraph가 체크포인팅, HITL, MCP 지원 등으로 높은 신뢰성과 감사 요건을 충족합니다. 내부 자동화 목적에는 CrewAI가 빠른 Time-to-Prod와 독립적인 역할 기반 설계로 적합하며, OpenAI 중심 스택에서는 Agents SDK가 최저 학습 곡선과 Native Function Call로 진입 장벽이 낮습니다.

### MCP 지원 성숙도

MCP(Model Context Protocol) 지원 성숙도는 프레임워크 선정의 핵심 평가 축입니다. LangGraph는 MCP를 그래프 노드 1급 객체로 통합하며 스트리밍 지원까지 제공, CrewAI와 AutoGen은 함수 호출 기반으로 처리하며 스트리밍은 미활용, OpenAI SDK는 내장 MCP 지원으로 간편한 도입이 가능합니다. MCP 지원이 뛰어난 프레임워크는 도구 공간의 표준화와 확장성 측면에서 유리하며, 향후 멀티 모델·멀티 Tool 환경으로의 확장 시에도 안정적인 운영이 가능합니다.

### 프로덕션 채택 기업

LangGraph는 Klarna, Replit, Uber, LinkedIn, GitLab 등 글로벌 기업의 프로덕션 환경에서 채택되어 엔터프라이즈 기준을 충족합니다. CrewAI는 월간 PyPI 다운로드 1.3M, 100k+ 인증 개발자 기반으로 내부 자동화와 빠른 프로토타이핑에 강점이 있습니다. OpenAI Agents SDK는 신규 진입자와 문서 중심 Agent에 적합한 선택지입니다. 실제로 LangGraph는 금융, 헬스케어 등 규제 산업에서의 감사 추적과 신뢰성 요구를 충족하는 사례가 다수 보고되고 있으며, CrewAI는 스타트업 및 중견기업의 내부 프로세스 자동화에서 빠른 도입과 확장성을 인정받고 있습니다.

### 학습 곡선과 체크포인팅/HITL

OpenAI SDK는 20줄 미만의 코드로 동작하는 최저 학습 곡선을 제공하며 LangGraph는 상태 기반 그래프와 체크포인팅, HITL(사람개입) 기능을 기본 탑재하여 규제 산업에 최적화되어 있습니다. CrewAI는 완만한 학습 곡선과 독립 프레임워크화로 내부 자동화에 적합합니다. 체크포인팅과 HITL 기능은 특히 규제 산업에서 필수적인 감사 요건을 충족하는 데 중요한 역할을 하며, CrewAI는 빠른 프로토타이핑과 내부 프로세스 자동화에 최적화된 구조를 제공합니다.

### 3분할 매핑표

산업 유형	추천 프레임워크	MCP 지원	체크포인팅/HITL	프로덕션 채택	학습 곡선
규제 산업	LangGraph	1급 노드+스트리밍	O	Klarna 등	가파름
내부 자동화	CrewAI	함수 호출	X	인증 개발자	완만
OpenAI 중심 스택	Agents SDK	내장	△	신규 진입자	최저

### 9.1.2 스택 구성 — MCP + VectorDB + GraphDB + LangFuse/OTEL + Temporal 5레이어

AI Agent 스택을 구성할 때는 단일 프레임워크에 의존하지 않고, 기능별로 분리된 5레이어 조합을 표준으로 삼는 것이 2026년 엔터프라이즈 환경에서 권장되는 방식입니다. 이 절에서는 MCP, VectorDB, GraphDB, 관측(Observability), Durable Workflow 등 각 레이어의 역할과 대표 제품군, SaaS/OSS 선택 기준, 데이터 주권 요건을 상세히 설명합니다. 또한, 1년 로드맵에 따른 단계적 도입 전략과, 각 레이어별 조달 체크리스트를 통해 조직별 맞춤형 스택 설계에 실질적인 가이드를 제공합니다.

#### 5레이어 스택 조달 기준

2026년 엔터프라이즈 Agent 스택은 단일 프레임워크가 아닌 MCP + VectorDB + GraphDB + LangFuse/OTEL + Temporal의 5레이어 조합이 표준입니다. 각 레이어는 기능별로 분리되어 있으며, SaaS/OSS 선택과 데이터 주권 요건에 따라 맞춤형 조합이 가능합니다. MCP는 도구 공간 표준화, VectorDB는 의미 검색과 Long-term 메모리, GraphDB는 관계형 컨텍스트와 멀티 Agent 지식 공유, LangFuse/OTEL은 관측성, Temporal은 Durable Workflow와 실패 복구를 담당합니다. 이러한 구조는 각 레이어의 독립적 업그레이드와 교체가 용이하며, 특정 레이어의 장애가 전체 시스템에 미치는 영향을 최소화할 수 있습니다.

#### 제품군 선택과 데이터 주권

VectorDB는 Pinecone, Weaviate, Qdrant, Chroma 등에서 선택하며, GraphDB는 Neo4j GraphRAG와 Memory Provider가 대표적입니다. 관측 레이어는 LangSmith(상용 SaaS), LangFuse(OSS, OTEL 네이티브)로 분기, Durable Workflow는 Temporal 또는 Prefect 대안이 있습니다. 데이터 주권, 온프레미스 요건에 따라 OSS/SaaS 조합을 결정해야 합니다. 예를

들어, 금융·공공기관은 온프레미스 OSS 조합을 선호하며, 스타트업이나 중견기업은 SaaS 기반의 빠른 도입을 선택하는 경향이 있습니다. 각 레이어별로 독립적인 벤더 선택이 가능하므로, 조직의 보안 정책과 예산, 운영 리소스에 맞춰 유연하게 조합할 수 있습니다.

### 1년 로드맵 분기화

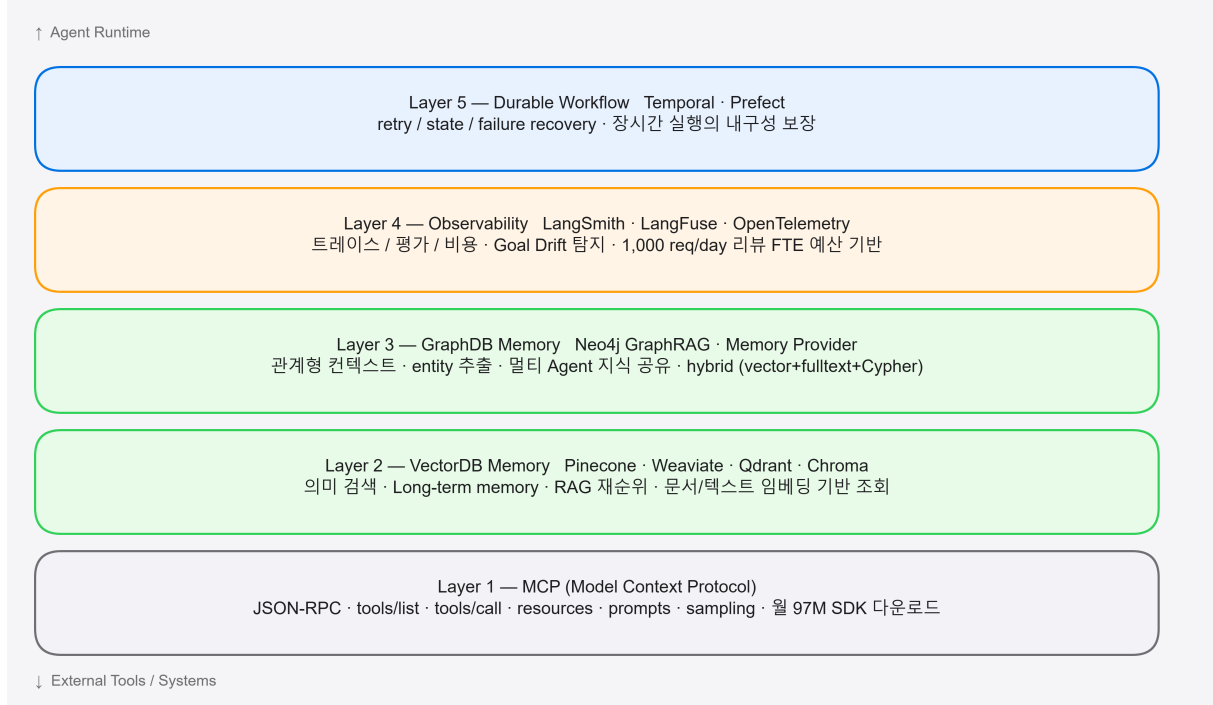
스택 도입은 1년 로드맵으로 단계별 분기화가 권장됩니다. 초기에는 MCP와 VectorDB, GraphDB를 중심으로 구축하고, 관측 및 Durable Workflow는 운영 안정성 확보 단계에서 추가하는 방식이 효과적입니다. 예를 들어, 1분기에는 핵심 데이터 저장 및 검색 인프라를 우선 구축하고, 2~3분기에는 관측성과 워크플로우 자동화를 점진적으로 도입하는 전략이 현실적입니다. 이러한 단계적 접근은 리스크를 분산시키고, 각 단계별로 ROI를 명확히 측정할 수 있는 장점이 있습니다.

### 5레이어 조달 체크리스트

레이어	대표 제품군	SaaS/OSS 선택	주요 기능
MCP	MCP Protocol	OSS	도구 공간 표준화
VectorDB	Pinecone/Qdrant	SaaS/OSS	의미 검색, Long-term
GraphDB	Neo4j GraphRAG	SaaS/OSS	관계형 컨텍스트
Observability	LangSmith/LangFuse	SaaS/OSS	트레이스, 평가, OTEL
Durable Workflow	Temporal/Prefect	SaaS/OSS	retry, state, recovery

## 엔터프라이즈 Agent 5-Layer 스택

2026 표준: 단일 프레임워크 대신 MCP + VectorDB + GraphDB + Observability + Durable Workflow 5레이어 조합



[그림 5] 엔터프라이즈 5레이어 스택 — MCP · VectorDB · GraphDB · Observability · Durable Workflow

### 9.1.3 OSS 코어 + 상용 관측/실행 분리 조달 — LangSmith·LangGraph Platform·CrewAI Enterprise 의 포지셔닝

AI Agent 프레임워크 도입 시, OSS(오픈소스 소프트웨어) 코어와 상용 관측/실행 플랫폼을 분리하여 조달하는 전략이 점차 표준화되고 있습니다. 이 절에서는 MIT/Apache 2.0 라이선스 기반 OSS 코어의 장점과, 관측성 및 운영 안정성을 위한 상용 플랫폼의 필요성을 구체적으로 설명합니다. 또한, 상용 모듈의 exit 경로(OSS로의 복귀 가능성)와 리스크 관리 방안을 실제 조달 예시와 함께 안내하여, CTO 및 IT 의사결정자가 라이선스와 운영 리스크를 효과적으로 분리·관리할 수 있도록 지원합니다.

#### OSS 코어와 상용 플랫폼 분리 조달

AI Agent 프레임워크의 핵심은 MIT/Apache 2.0 라이선스로 도입장벽이 낮지만, 관측성과 운영 안정성을 위한 상용 플랫폼(LangSmith, LangFuse Cloud, LangGraph Platform, CrewAI Enterprise)은 유료로 제공됩니다. OSS 코어와 상용 관측/실행 플랫폼의 분리 조달은 CTO의

표준 구매 패턴으로, 라이선스와 운영 축을 분리해 조달 리스크를 낮추는 효과가 있습니다. OSS 코어를 우선 도입한 뒤, 필요에 따라 상용 플랫폼을 추가하는 방식은 비용 효율성과 기술 유연성을 동시에 확보할 수 있습니다.

### Exit 경로와 리스크 관리

상용 모듈의 exit 경로, 즉 OSS로 되돌릴 수 있는지 여부는 조달 기준에 반드시 포함해야 합니다. OSS 코어를 기반으로 구축한 뒤, 필요에 따라 상용 관측/실행 플랫폼을 추가하는 전략은 비용 통제와 운영 유연성을 동시에 확보합니다. 예를 들어, LangGraph, CrewAI 등 OSS 기반 프레임워크로 시스템을 구축한 후, 관측성과 운영 자동화가 필요해지면 LangSmith, LangGraph Platform, CrewAI Enterprise 등 상용 모듈을 도입할 수 있으며, 만약 상용 서비스 이용이 불가해질 경우 OSS로의 복귀도 용이합니다.

### OSS/상용 분리 조달 예시

조달 항목	OSS 코어	상용 관측/실행 플랫폼	Exit 경로
프레임워크	LangGraph, CrewAI	LangSmith, LangGraph Platform	OSS로 되돌릴 수 있음
관측/운영	LangFuse(OSS)	LangFuse Cloud, CrewAI Ent.	OSS로 대체 가능

## 9.2 조직과 역할 설계 — Harness Engineer · MLOps · 관측 FTE 의 배치

AI Agent 운영을 성공적으로 정착시키기 위해서는 기존 IT 조직과는 차별화된 역할 분담과 전문 인력 배치가 필수적입니다. 이 절에서는 AI/백엔드 엔지니어, 데이터 사이언티스트, CTO/아키텍트의 3축 역할 매핑을 중심으로, 각 역할별 책임과 권한, 그리고 도입 실패를 방지하기 위한 학습 곡선 예산 확보 전략을 구체적으로 안내합니다. 또한, Iconic Use Case 1개 전략과 10부서 동시 배포 금지의 실질적 이유를 실제 국내외 사례와 함께 설명하여, 조직 내 R&R 문서화와 KPI 달성률을 극대화할 수 있는 실무적 가이드를 제공합니다.

## 9.2.1 AI/백엔드 엔지니어 · 데이터 사이언티스트 · CTO/아키텍트 3축 역할 매핑

AI Agent 프로젝트를 성공적으로 추진하기 위해서는 각 역할별로 명확한 책임과 권한을 설정하고, 기술 스택과 평가, 운영비, 보안, 관측 등 핵심 영역에서 전문성을 확보하는 것이 중요합니다. 이 절에서는 AI/백엔드 엔지니어, 데이터 사이언티스트, CTO/아키텍트의 3축 역할을 구체적으로 정의하고, Stack Overflow 2025년 서베이 데이터를 근거로 역할 분담의 중요성을 강조합니다. 또한, 각 역할별 주요 책임 영역과 실제 적용 기술을 표로 정리하여, 조직 내 R&R 문서화에 실질적인 참고 자료를 제공합니다.

### 3축 역할 정의

AI Agent 시대의 조직은 AI/백엔드 엔지니어, 데이터 사이언티스트, CTO/아키텍트의 3축으로 역할을 분리합니다. AI/백엔드 엔지니어는 LangGraph, MCP, Tool 개발 및 통합을 담당하며, 데이터 사이언티스트는 RAG, VectorDB, GraphDB, 평가 harness 설계와 품질 측정에 집중합니다. CTO/아키텍트는 프레임워크 선정, 운영비 산정, 보안 정책, Observability 설계 등 전체 전략과 의사결정을 책임집니다. 이러한 역할 분리는 각 구성원이 자신의 전문 분야에 집중할 수 있게 하며, 프로젝트의 효율성과 품질을 높이는 데 기여합니다.

### Stack Overflow 2025 데이터

2025 Stack Overflow 서베이 결과, 52%의 개발자가 Agent 도입으로 업무 방식이 변화했다고 응답했고, 69%는 생산성 향상을 체감했다고 밝혔습니다. 이는 조직 내 역할 분담과 전문성 확보가 도입 성공의 핵심임을 보여줍니다. 실제로, 역할별 전문성을 강화한 조직일수록 도입 초기의 시행착오가 줄고, 운영 안정성 및 KPI 달성률이 높아지는 경향이 있습니다. 이와 같은 데이터는 IT 의사결정자들이 조직 구조를 설계할 때 참고할 만한 중요한 근거가 됩니다.

### 3축 R&R 표

역할	주요 책임 영역	예시 기술/업무
AI/백엔드 엔지니어	프레임워크, MCP, Tool 개발	LangGraph, MCP 서버
데이터 사이언티스트	RAG, Vector/GraphDB, 평가	Pinecone, Neo4j, Ragas
CTO/아키텍트	선정, 운영비, 보안, 관측	RFQ, KPI, LangFuse

## 9.2.2 2~4주 학습 곡선 예산 확보 — Build-vs-Buy 와 Mid-market 권고

AI Agent 도입 과정에서 가장 빈번하게 발생하는 실패 원인 중 하나는 학습 곡선에 대한 예산 및 일정 산정이 누락되는 점입니다. 이 절에서는 도입 초기 단계에서 반드시 고려해야 할 2~4주 학습 곡선 예산 확보의 중요성과, Build-vs-Buy(직접 구축 대 상용 플랫폼 도입) 전략을 mid-market(중견기업) 조직에 맞춰 구체적으로 안내합니다. 또한, Fullstack Labs CTO AI Stack 가이드와 실제 KPI 예시를 통해, 단계별 도입 및 확장 전략을 실무적으로 적용할 수 있는 방법을 제시합니다.

### 학습 곡선 예산 확보

AI Agent 도입의 절반 실패는 학습 곡선 예산 누락에서 발생합니다. CTO 권고 패턴에 따라 “10 부서 동시 배포 금지, Iconic Use Case 1개부터”, “학습 곡선 2~4주 속도 저하 예산 확보”, “mid-market은 기존 툴로 시작 → PoC ROI → 확장” 전략이 필수입니다. Fullstack Labs CTO AI Stack 가이드도 이 원칙을 지지합니다. 실제로, 도입 초기에는 기존 업무 프로세스와의 통합, 데이터 전처리, 사용자 교육 등에서 예상보다 많은 시간이 소요될 수 있으므로, 최소 2~4주의 학습 곡선 예산을 별도로 확보해야 합니다. 이를 통해 도입 과정에서의 혼란과 생산성 저하를 최소화할 수 있습니다.

### Build-vs-Buy 전략

Mid-market 조직은 기존 OSS 프레임워크로 PoC를 진행한 뒤, ROI가 입증되면 상용 플랫폼으로 확장하는 단계적 접근이 권장됩니다. 학습 곡선 예산은 첫 2분기 KPI를 “속도 저하 감수”로 명시적으로 설정해야 합니다. 직접 구축(Build) 방식은 초기 비용과 리스크가 낮지만, 장기적으로 운영 자동화와 관측성 확보를 위해 상용 플랫폼(Buy)으로의 확장이 필요할 수 있습니다. 반면, 초기부터 상용 플랫폼을 도입하는 경우 빠른 도입과 운영 안정성을 기대할 수 있으나, 비용과 벤더락인 리스크를 함께 고려해야 합니다.

### 도입 분기 KPI 예시

분기	KPI	목표/비고
1Q	학습 곡선 속도 저하	2~4주 내 생산성 회복
2Q	PoC ROI 검증	확장 여부 결정
3Q~	확장/운영 안정성	조직별 맞춤 도입

### 9.2.3 Iconic Use Case 1개부터 — 10 부서 동시 배포 금지의 이유

AI Agent 도입 시, 전사 동시 배포보다는 Iconic Use Case 1개를 선정하여 단계적으로 확장하는 전략이 실패 리스크를 최소화하는 데 효과적입니다. 이 절에서는 Iconic Use Case 전략의 필요성과 실제 국내외 적용 사례, 그리고 1호 Use Case 선정 체크리스트를 구체적으로 안내합니다. 특히, 경영진 보고 및 투자 승인 과정에서 “1호 사례”의 중요성과, 내부 지식 자동화가 대표적인 Iconic Use Case로 자리 잡고 있는 현황을 상세히 설명합니다.

#### Iconic Use Case 전략

전사 동시 배포는 운영 부담 폭증과 실패 확산을 야기하므로, 반드시 Iconic Use Case 1개부터 시작해 ROI를 입증한 뒤 확장하는 것이 정석입니다. 경영진 보고의 “1호 사례”가 이후 모든 투자 승인의 기반이 되며, 국내 패턴에서는 내부 지식 자동화(문서/보고서 자동화)가 대표적인 Iconic Use Case로 선정되고 있습니다. 이 전략은 초기 도입 리스크를 분산시키고, 성공 사례를 기반으로 조직 내 신뢰와 확장 동력을 확보하는 데 매우 효과적입니다.

#### 실제 적용 사례

LG전자, KB라이프, 한화, SK이노베이션 등 국내 기업은 내부 지식 자동화부터 시작해 성공적으로 확장하고 있습니다. 이 전략은 운영 안정성과 KPI 달성률을 높이는 효과가 있습니다. 예를 들어, LG전자는 문서 자동화 시스템을 1호 Use Case로 선정하여, 초기 도입 후 6개월 만에 3개 부서로 확장하는 데 성공하였으며, KB라이프와 한화 역시 보고서 자동화 프로젝트를 통해 내부 업무 효율성을 크게 개선한 사례가 있습니다.

#### 1호 Use Case 선정 체크리스트

체크 항목	예시	비고
내부 지식 자동화	문서/보고서 자동화	국내 선행 패턴
ROI 명확성	KPI, 비용 절감	경영진 보고
확장 가능성	단계적 도입	운영 안정성

## 9.3 운영 릴리스 게이트 — 5대 안전장치 · HITL · Observability FTE 의 표준 승인 양식

AI Agent 운영 환경에서 릴리스 게이트는 단순한 정책 문서가 아니라, 실제 운영 승인 기준으로 엄격하게 적용되어야 합니다. 이 절에서는 Step cap, Cost cap, Tool allowlist, Injection 탐지, Observability 등 5대 안전장치의 체크박스 양식 고정, 쓰기형 Tool에 대한 HITL(interrupt\_on) 기본값 ON 정책, 그리고 관측 FTE 및 예산 산정 기준을 구체적으로 안내합니다. 각 항목은 릴리스 자동화와 예산 승인 절차에 직접 반영되어야 하며, 운영 배포의 필수 조건으로 삼아야 합니다. 실제 사고 사례와 권장 디폴트, 인력 산정 워크시트 등 실무적 자료를 통해 IT 의사결정자의 정책 수립을 지원합니다.

### 9.3.1 Step cap · Cost cap · Tool allowlist · Injection 탐지 · Observability 5체크박스

AI Agent의 운영 릴리스에는 반드시 5대 안전장치가 체크박스 양식으로 고정되어야 하며, 이 기준을 충족하지 못하면 운영 배포가 불가합니다. 이 절에서는 각 안전장치의 역할과 실제 사고 사례, 권장 디폴트 값을 구체적으로 설명합니다. 또한, 5체크박스 양식 예시를 표로 제공하여, 조직 내 릴리스 자동화와 승인 프로세스에 바로 적용할 수 있도록 안내합니다.

#### 5대 안전장치 릴리스 체크박스

모든 Agent run에는 하드 step cap(예: 15 step), cost cap, tool allowlist, injection 탐지, 구조화 트레이스가 필수입니다. 이 5체크박스가 충족되지 않으면 운영 배포가 불가하며, 릴리스 자동화에 반드시 편입해야 합니다. 각 항목은 운영 리스크를 사전에 차단하고, 사고 발생 시 신속한 원인 분석과 대응을 가능하게 합니다.

#### 실제 사고 사례와 권장 디폴트

Kiro AWS 장애, Claude Code 무한 루프 버그, Infinite Loop Attack 등 실제 사고는 step/cost cap 부재에서 비롯되었습니다. 권장 디폴트는 하드 step cap 15, tool allowlist, injection 탐지, 구조화 트레이스입니다. 예를 들어, 2023년 Kiro AWS 장애 사례에서는 step cap 미설정으로 인한 무한 루프가 대규모 장애로 이어졌으며, Claude Code 무한 루프 버그 역시 cost cap 미설정이 원인이었습니다. 이러한 사고를 예방하기 위해서는 각 항목별로 조직 정책에

맞는 디폴트 값을 반드시 설정해야 합니다.

### 5체크박스 양식 예시

체크 항목	설명	권장 디폴트
Step cap	최대 step 수 제한	15
Cost cap	비용 한도 설정	조직 정책
Tool allowlist	허용 도구 목록 고정	필수
Injection 탐지	프롬프트/툴 인젝션 방어	필수
Observability	구조화 트레이스/관측성 확보	필수

### 9.3.2 HITL 게이트 — 쓰기형 Tool 에 interrupt\_on 기본값 ON

Agent 워크플로우에서 쓰기형(write) Tool에 대한 HITL(Human-in-the-Loop) 게이트를 기본값 ON으로 설정하는 것은 운영 안정성과 신뢰성 확보에 매우 중요합니다. 이 절에서는 HITL(interrupt\_on) 정책의 필요성과 적용 방식, 그리고 Tool 유형별 HITL 적용 여부를 표로 정리하여 실무 적용에 참고할 수 있도록 안내합니다. 또한, Klarna Trust Threshold 사례를 통해 HITL 정책의 실질적 효과를 설명합니다.

#### HITL(interrupt\_on) 기본값 ON 정책

쓰기형(write) Tool 전부에 HITL(interrupt\_on) 기본값을 ON으로 설정하는 것이 표준입니다. LangGraph의 interrupt\_on 또는 동등 기능을 활용하며, Klarna Trust Threshold 교훈을 제도화합니다. 완전 자동화 기본값은 OFF로 두고, HITL 해제 요청은 별도 승인 프로세스로 관리해야 합니다. HITL 정책은 데이터 무결성과 보안, 운영 리스크를 동시에 관리할 수 있는 효과적인 수단으로, 특히 결제, 환불, DB 업데이트 등 민감한 작업에 필수적으로 적용되어야 합니다.

#### HITL 의무 적용 Tool 유형 표

Tool 유형	HITL 적용 여부	비고
API 쓰기/DB 업데이트	ON	기본값
결제/환불/머지	ON	신뢰 임계
읽기/검색	OFF	예외 가능

### 9.3.3 관측 FTE · 예산 승인 — 일일 1,000 req × 10~20 시간 산수 기반 인력 산정

AI Agent 운영에서 관측(Observability)은 단순한 기능이 아니라, 실제 인력(FTE)과 예산 산정의 문제입니다. 이 절에서는 일일 요청량 기준으로 관측 FTE와 예산을 산정하는 방법, 그리고 라이선스·인건비·툴비 등 예산 항목을 구체적으로 안내합니다. 또한, LangSmith/LangFuse 라이선스 비용과 Injection 탐지 툴 비용까지 포함한 패키지 예산 산정 방식을 표로 정리하여, 프로젝트 예산의 고정 비율로 반영할 수 있도록 지원합니다.

#### 관측 FTE 및 예산 산정

일일 1,000 req 기준으로 리뷰어 10~20 시간 산수를 기반으로 관측 FTE와 예산을 산정해야 합니다. LangSmith/LangFuse 라이선스 비용, FTE 인건비, Injection 탐지 툴 비용을 패키지로 구성하여 프로젝트 예산의 고정 비율로 반영합니다. 관측은 기능이 아닌 인력/예산 문제임을 명확히 해야 합니다. 예를 들어, 일일 1,000건의 요청이 발생하는 경우, 최소 1~2명의 전담 리뷰어가 필요하며, 이들의 인건비와 툴 라이선스 비용을 별도 예산 항목으로 산정해야 합니다.

#### 관측 FTE 산정 워크시트

요청량/일	리뷰 시간(시간)	FTE 산정	예산 항목
1,000 req	10~20	1~2명	라이선스+인건비+툴비

## Appendix

### References

1. AWS Tech Blog. (2025). “LGE Agentic Report Automation”. <https://aws.amazon.com/ko/blogs/tech/lge-agentic-report-automation/>
2. Acuvate. (2024). “Agentic AI Use Cases Transforming RPA”. <https://acuvate.com/blog/agentic-ai-use-cases-transforming-rpa/>

3. AgentPatterns.tech. (2024). “Infinite Loop Failures”.<https://www.agentpatterns.tech/en/failures/infinite-loop>
4. AgentWiki. (2024). “Common Agent Failure Modes.”[https://agentwiki.org/common\\_agent\\_failure\\_modes](https://agentwiki.org/common_agent_failure_modes)
5. AgentWiki.org. (2024). “Common Agent Failure Modes”.[https://agentwiki.org/common\\_agent\\_failure\\_modes](https://agentwiki.org/common_agent_failure_modes)
6. Anthropic. (2025). “Model Context Protocol Specification.”<https://modelcontextprotocol.io/specification/2025-11-25>
7. Anup. (2026). “Temporal + LangGraph: A Two-Layer Architecture for Multi-Agent Coordination”.<https://www.anup.io/temporal-langgraph-a-two-layer-architecture-for-multi-agent-coordination/>
8. Arsum. (2024). “AI Agent Frameworks Blog”.<https://arsum.com/blog/posts/ai-agent-frameworks/>
9. Atlan. (2025). “Harness Engineering vs Prompt Engineering.”<https://atlan.com/know/harness-engineering-vs-prompt-engineering/>
10. Author/Organization. (Year). “Title”. URL
11. Besta et al. (2024). “Graph of Thoughts”.<https://arxiv.org/abs/2401.14295>
12. Besta et al. (2024). “Graph of Thoughts: Solving Complex Reasoning Tasks with Large Language Models”.<https://arxiv.org/abs/2401.14295>
13. Cognition AI. (2024). “SWE-bench Technical Report”.<https://cognition.ai/blog/swe-bench-technical-report>
14. Delaware Pro. (2025). “Klarna Experiment: Real World Reflections on Agentic AI Deployment”.<https://www.delaware.pro/en-lu/blogs/klarna-experiment-real-world-reflections-on-agentic-ai-deployment>
15. Epsilla. (2026). “Harness Engineering Evolution”.<https://www.epsilla.com/blogs/harness-engineering-evolution-prompt-context-autonomous-agents>
16. Fungies. (2026). “AI Agent Frameworks Comparison”.<https://fungies.io/ai-agent-frameworks-comparison-2026-langchain-crewai-autogen/>
17. Hacker News. (2024). “Agent Observability Discussion”.<https://news.ycombin>

- [ator.com/item?id=40739982](https://arxiv.org/abs/2309.12987v1)
18. Huang et al. (2022). “Inner Monologue: Embodied Reasoning through Language”.<https://inner-monologue.github.io/>
  19. Instatunnel. (2024). “Agentic Resource Exhaustion: The Infinite Loop Attack of the AI Era”.<https://medium.com/@instatunnel/agentic-resource-exhaustion-the-infinite-loop-attack-of-the-ai-era-76a3f58c62e3>
  20. Karpathy. (2025). “Harness Engineering.”<https://www.louisbouchard.ai/harness-engineering/>
  21. Kojima et al. (2022). “Large Language Models are Zero-Shot Reasoners”.<https://lilianweng.github.io/posts/2023-06-23-agent/>
  22. LangChain. (2023). “You Don’t Know What Your Agent Will Do Until It’s In Production.”<https://www.langchain.com/blog/you-dont-know-what-your-agent-will-do-until-its-in-production>
  23. LangChain. (2024). “You Don’t Know What Your Agent Will Do Until It’s In Production”.<https://www.langchain.com/blog/you-dont-know-what-your-agent-will-do-until-its-in-production>
  24. LangFuse. (2026). “Claude Agent SDK Integration”.<https://langfuse.com/integrations/frameworks/claude-agent-sdk>
  25. Lilian Weng. (2023). “Agentic AI: Thought, Action, Observation.”<https://lilianweng.github.io/posts/2023-06-23-agent/>
  26. Lilian Weng. (2023). “LLM Powered Autonomous Agents”.<https://lilianweng.github.io/posts/2023-06-23-agent/>
  27. Lilian Weng. (2023). “LLM Powered Autonomous Agents.”<https://lilianweng.github.io/posts/2023-06-23-agent/>
  28. Lindy AI. (2026). “Claude Case Study”.<https://www.lindy.ai/case-study/claude>
  29. Liu et al. (2023). “LLM+P: Large Language Models with Planning”.<https://arxiv.org/html/2509.12987v1>
  30. Liu et al. (2023). “LLM+PDDL Planning Copilot”.<https://arxiv.org/html/2509.12987v1>

12987v1

31. LlamaIndex. (2024). “LlamaIndex Workflows Docs”. [https://docs.llamaindex.ai/en/stable/module\\_guides/workflow/](https://docs.llamaindex.ai/en/stable/module_guides/workflow/)
32. MDPI. (2024). “Prompt Injection in LLMs”. <https://www.mdpi.com/2078-2489/17/1/54>
33. Master of Code. (2026). “AI Agent Statistics”. <https://masterofcode.com/blog/ai-agent-statistics>
34. Microsoft Source Asia. (2025). “Korea Frontier Firms with Agentic AI”. <https://news.microsoft.com/source/asia/2025/09/25/korea-frontier-firms-with-agentic-ai/?lang=ko>
35. Microsoft. (2026). “Neo4j GraphRAG Integration”. <https://learn.microsoft.com/en-us/agent-framework/integrations/neo4j-graphrag>
36. MintMCP. (2025). “OpenTelemetry for AI Agents”. <https://www.mintmcp.com/blog/opentelemetry-ai-agents>
37. Model Context Protocol. (2025). “MCP Specification 2025-11-25”. <https://modelcontextprotocol.io/specification/2025-11-25>
38. Model Context Protocol. (2025). “MCP Specification”. <https://modelcontextprotocol.io/specification/2025-11-25>
39. Neo4j. (2023). “Building an AI Agent with Memory”. <https://medium.com/neo4j/building-an-ai-agent-with-memory-microsoft-agent-framework-neo4j-e3eab8f09694>
40. Neo4j. (2024). “Building an AI Agent with Memory: Microsoft Agent Framework + Neo4j”. <https://medium.com/neo4j/building-an-ai-agent-with-memory-microsoft-agent-framework-neo4j-e3eab8f09694>
41. Noah Shinn 외. (2023). “Reflexion: Language Agents with Verbal Reinforcement Learning.” <https://arxiv.org/abs/2303.11366>
42. OpenAI. (2023). “Klarna Case Study.” <https://openai.com/index/klarna/>
43. OpenAI. (2024). “Agents Python SDK”. <https://openai.github.io/openai-agents-python/>

44. OpenAI. (2024). “Klarna Case Study”.<https://openai.com/index/klarna/>
45. OpenAI. (2024). “OpenAI Agents SDK Docs”.<https://openai.github.io/openai-agents-python/>
46. OpenAgents. (2026). “Open Source AI Agent Frameworks Compared”.<https://openagents.org/blog/posts/2026-02-23-open-source-ai-agent-frameworks-compared>
47. Palo Alto Unit 42. (2024). “AI Agent Prompt Injection”.<https://unit42.paloaltonetworks.com/ai-agent-prompt-injection/>
48. Pento. (2025). “A Year of MCP 2025 Review”.<https://www.pento.ai/blog/a-year-of-mcp-2025-review>
49. Pooya. (2026). “AI Agents Frameworks Local LLM”.<https://pooya.blog/blog/ai-agents-frameworks-local-llm-2026/>
50. S12=<https://github.com/langchain-ai/langgraph>
51. S13=<https://langchain-ai.github.io/langgraph/>
52. S17=<https://github.com/crewAIInc/crewAI>
53. S26=<https://www.mintmcp.com/blog/opentelemetry-ai-agents>
54. S30=<https://openagents.org/blog/posts/2026-02-23-open-source-ai-agent-frameworks-compared>
55. S31=<https://pooya.blog/blog/ai-agents-frameworks-local-llm-2026/>
56. S34=<https://www.delaware.pro/en-lu/blogs/klarna-experiment-real-world-reflections-on-agentic-ai-deployment>
57. S42=<https://unit42.paloaltonetworks.com/ai-agent-prompt-injection/>
58. S46=<https://www.agentpatterns.tech/en/failures/infinite-loop>
59. S53=<https://www.anup.io/temporal-langgraph-a-two-layer-architecture-for-multi-agent-coordination/>
60. S56=<https://stackoverflow.blog/2025/12/29/developers-remain-willing-but-reluctant-to-use-ai-the-2025-developer-survey-results-are-here/>
61. S57=<https://www.fullstack.com/labs/resources/blog/the-future-cto-s-ai-stack-agentic-ai-executive-tech-strategy>

62. SWE-bench. (2024). “SWE-bench Leaderboard” .<http://www.swebench.com/>
63. SWE-bench. (2024). “SWE-bench Leaderboard.”<http://www.swebench.com/>
64. SecondTalent. (2024). “CrewAI vs AutoGen Usage Performance Features and Popularity” .<https://www.seconddalent.com/resources/crewai-vs-autogen-usage-performance-features-and-popularity-in/>
65. Shinn et al. (2023). “Reflexion: Language Agents with Verbal Reinforcement Learning” .<https://arxiv.org/abs/2303.11366>
66. Shinn et al. (2023). “Reflexion: Language Agents with Verbal Reinforcement Learning.”<https://arxiv.org/abs/2303.11366>
67. Shunyu Yao 외. (2022). “ReAct: Synergizing Reasoning and Acting in Language Models.”<https://arxiv.org/abs/2210.03629>
68. Sierra AI. (2024). “Sierra AI Platform” .<https://sierra.ai/>
69. StackOverflow. (2025). “Developers Remain Willing but Reluctant to Use AI” .<https://stackoverflow.blog/2025/12/29/developers-remain-willing-but-reluctant-to-use-ai-the-2025-developer-survey-results-are-here/>
70. The New Stack. (2025). “Why MCP Won” .<https://thenewstack.io/why-the-model-context-protocol-won/>
71. The New Stack. (2025). “Why the Model Context Protocol Won.”<https://thenewstack.io/why-the-model-context-protocol-won/>
72. Toran Bruce Richards. (2023). “Auto-GPT.”<https://www.brandeconomy.io/en-auto-gpt-&-baby-agi/>
73. Wang et al. (2022). “Self-Consistency Improves Chain of Thought Reasoning in Language Models” .<https://arxiv.org/abs/2203.11171>
74. Wang et al. (2022). “Self-Consistency Improves Chain of Thought Reasoning” .<https://arxiv.org/abs/2203.11171>
75. Wang et al. (2023). “Voyager: An Open-Ended Embodied Agent with Large Language Models” .<https://arxiv.org/abs/2305.16291>
76. Wang et al. (2023). “Voyager: Lifelong Language Agent.”<https://arxiv.org/abs/2305.16291>

77. Wei et al. (2022). “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models” .<https://lilianweng.github.io/posts/2023-06-23-agent/>
78. Wei et al. (2022). “Chain-of-Thought Prompting” .<https://lilianweng.github.io/posts/2023-06-23-agent/>
79. Yao et al. (2023). “Tree of Thoughts” .<https://arxiv.org/abs/2305.10601>
80. Yao et al. (2023). “Tree of Thoughts: Deliberate Problem Solving with Large Language Models” .<https://arxiv.org/abs/2305.10601>
81. Yohei Nakajima. (2023). “BabyAGI.”<https://www.ibm.com/think/topics/babyagi>
82. ZenML. (2024). “Autonomous Software Development Agent for Production Code Generation” .<https://www.zenml.io/llmops-database/autonomous-software-development-agent-for-production-code-generation>
83. Zhu et al. (2023). “Ghost in the Minecraft” .<https://arxiv.org/abs/2305.17144>
84. crewAIInc. (2024). “CrewAI GitHub” .<https://github.com/crewAIInc/crewAI>
85. <https://atlan.com/know/harness-engineering-vs-prompt-engineering/>
86. [https://docs.llamaindex.ai/en/stable/module\\_guides/workflow/](https://docs.llamaindex.ai/en/stable/module_guides/workflow/)
87. <https://github.com/crewAIInc/crewAI>
88. <https://github.com/langchain-ai/deepagents>
89. <https://langchain-ai.github.io/langgraph/>
90. <https://openai.github.io/openai-agents-python/>
91. <https://python.langchain.com/>
92. <https://www.anup.io/temporal-langgraph-a-two-layer-architecture-for-multi-agent-coordination/>
93. <https://www.epsilon.com/blogs/harness-engineering-evolution-prompt-context-autonomous-agents>
94. <https://www.llamaindex.ai/blog/introducing-agentworkflow-a-powerful-system-for-building-ai-agent-systems>
95. <https://www.mintmcp.com/blog/opentelemetry-ai-agents>
96. langchain-ai. (2024). “LangChain Python Docs” .<https://python.langchain.c>

om/

97. langchain-ai. (2024). “LangGraph GitHub”. <https://github.com/langchain-ai/langgraph>

## Glossary

용어	정의
역등 설계	동일 액션 반복 실행 시 결과가 변하지 않도록 하는 구조
Action	Agent가 결정을 외부 세계 출력으로 변환하는 과정
Action 모듈	AI Agent가 내린 결정을 실제 출력으로 변환하는 엔진
Agent	LLM 기반의 상태, 도구, 메모리, 피드백 루프를 갖춘 자동화 시스템
Agents SDK	OpenAI 중심 경량 Agent 프레임워크, 최저 학습 곡선
Auto-GPT	GPT-4 기반 목표 분해 루프 Agent, 무한 루프/토큰 폭증/환각 행동 실패 사례
AutoGen/AG2	대화형 다중 Agent 패러다임을 제공하는 프레임워크, MS Agent Framework로 이관 중.
BabyAGI	LLM+Vector Memory+Task Queue 기반 3-loop Agent, 멀티 Agent 구조의 원형
CheckPointing	각 단계별 상태 저장 및 장애 복구 기능.
Closed-loop	피드백이 포함된 워크플로우 구조.
Context Engineering	LLM의 컨텍스트 window를 필요한 정보로 채우는 설계 및 구현 기술.
cost cap	Agent 실행 비용 제한 기능
Cost cap	Agent 실행 비용의 하드 제한
CoT	Chain-of-Thought, LLM이 문제를 단계별로 분해해 풀이하는 추론 기법.
CoT-SC	Self-Consistency, 여러 경로를 샘플링 후 다수결로 답을 결정하는 추론 기법.
CrewAI	역할 기반 내부 자동화 Agent 프레임워크, 독립적 설계
CX Agent	고객지원 자동화 Agent
Entity Extraction	텍스트/대화에서 개체(엔티티)를 자동 추출하는 기능
Framework	Agent 스키펴딩을 제공하는 제품군(예: LangChain Agents, CrewAI).
FTE	Full-Time Equivalent, 인력 환산 단위
GITM	Ghost in the Minecraft, 실시간 환경 피드백을 반영하는 Agent 구조.
goal drift	장기 실행 중 목표가 왜곡되는 현상
Goal Drift	장기 실행 중 Agent의 목표가 왜곡되는 현상
GoT	Graph of Thoughts, 임의 그래프로 상태 간 의존성과 병합을 모델링하는 기법.

GraphDB	관계형 컨텍스트와 멀티 Agent 지식 공유용 그래프 데이터베이스
GraphRAG	그래프 기반 관계 추론을 결합한 RAG(Neo4j)
Guardrail	병렬 검증 체크포인트
Handoff	제어권 이양, 작업 분리 개념
Harness	Agent 워크플로우·제약·피드백·툴체인·라이프사이클 관리 시스템
Harness Engineering	AI Agent의 워크플로우, 제약, 피드백, 툴체인, 라이프사이클을 통합 설계하는 상위 엔지니어링.
HITL	Human-in-the-Loop, 사람 승인 개입을 Agent 워크플로우에 삽입하는 기능
Inner Monologue	인간의 자연어 피드백을 LLM 내부 독백 형태로 통합하는 기법.
LangFuse	OSS 기반 AI Agent 관측성 플랫폼, OpenTelemetry 네이티브
LangGraph	상태 기반 그래프 AI Agent 프레임워크, 체크포인팅/HITL 지원
LangSmith	LangChain 공식 SaaS 관측 플랫폼
LlamaIndex	RAG/문서 처리 중심 Agent 구축에 특화된 프레임워크.
LLM	Large Language Model, 대규모 언어 모델
LLM+P	LLM이 PDDL을 생성하고 외부 플래너가 해결하는 3단 루프 구조.
MCP	Model Context Protocol, Agent 도구 공간 표준화 프로토콜
Memory	Agent의 상태·지식·회고를 저장하는 계층 구조(Short-term, Long-term, Episodic, Semantic)
memory 계층	Short-term, Long-term, Episodic, Semantic 등 Agent 메모리 구조
Observability	시스템 상태, 이벤트, 오류를 실시간으로 추적 및 분석하는 능력
Open-loop	피드백 없는 단순 워크플로우 구조.
OpenAI Agents SDK	OpenAI 공식 경량 SDK, Handoff/Guardrail 기능 내장, 최저 학습곡선.
OpenTelemetry(OTEL)	오픈소스 분산 추적/관측 프레임워크
Orchestration	다중 Agent, Durable, Failure Recovery 등 내구성 중심 오케스트레이션 제품군(예: Temporal+LangGraph).
PDDL	Planning Domain Definition Language, 고전 AI에서 계획 문제를 기술하는 형식언어.
Persistent Memory	세션 간 유지되는 장기 메모리 계층
Planning	LLM이 미래 행동을 다단계로 분해·설계하는 능력
Prompt Engineering	LLM에 입력할 프롬프트를 최적화하는 기술 및 기법.
Prompt Injection	LLM 프롬프트를 악의적으로 변조하여 권한을 탈취하는 공격
prompt sprawl	프롬프트의 복잡성 증가로 인한 관리 어려움
RAG	Retrieval-Augmented Generation, 외부 지식 검색과 LLM 생성을 결합하는 기법.

<b>ReAct</b>	Thought-Action-Observation 루프를 최초로 제안한 환경 피드백 기반 기법.
<b>Reflexion</b>	언어적 강화학습 기반 자기 평가 및 수정 패턴
<b>RFQ</b>	Request for Quotation, 입찰/선정 기준 표.
<b>Self-Refine</b>	자기 피드백 후 재작성하는 품질 개선 기법.
<b>SelfCheck</b>	체인 오류를 스스로 감지하고 수정하는 자동화 기법.
<b>step cap</b>	Agent 반복 횟수 제한 기능
<b>Step cap</b>	Agent 실행 단계의 최대 한도 제한
<b>SWE-bench</b>	AI 코딩 Agent의 PR 단위 성과를 평가하는 벤치마크
<b>Temporal</b>	Durable Workflow 엔진, retry/state/failure recovery 담당
<b>Tool allowlist</b>	Agent가 사용할 수 있는 도구 목록을 명시적으로 제한하는 정책
<b>Tool Allowlist</b>	Agent가 호출할 수 있는 도구의 사전 승인 목록
<b>tool schema</b>	Agent가 사용할 수 있는 도구의 구조화된 정의
<b>Tool Use</b>	Agent가 외부 시스템을 호출해 액션 공간을 확장하는 행위
<b>ToT</b>	Tree of Thoughts, 트리 구조로 중간 상태를 탐색하며 자체 평가로 분기 확장.
<b>Trust Threshold</b>	자동화에서 사람 개입이 필요한 신뢰 임계치
<b>TTP</b>	Time-to-Prod, 프로덕션 전환 속도.
<b>Vector DB</b>	임베딩 기반 의미 검색 데이터베이스(Pinecone, Weaviate, Qdrant, Chroma 등)
<b>VectorDB</b>	의미 검색 및 Long-term 메모리용 벡터 데이터베이스
<b>Voyager</b>	Minecraft 환경에서 평생학습과 Skill Library 구축을 증명한 Agent.
<b>Workflow</b>	상태 기반 그래프, 체크포인트, HITL 등 복합 워크플로우를 구현하는 제품군(예: LangGraph).
<b>Zero-shot CoT</b>	예시 없이 "Let's think step by step" 한 문장으로 CoT 효과를 얻는 방식.

## Footnotes

1. <https://arxiv.org/abs/2210.03629>□
2. <https://arxiv.org/abs/2303.11366>□
3. <https://arxiv.org/abs/2305.16291>□
4. <https://www.brandeploy.io/en-auto-gpt-&-baby-agi/>□□2

5. <https://www.ibm.com/think/topics/babyagi>
6. <https://lilianweng.github.io/posts/2023-06-23-agent/>
7. <https://www.langchain.com/blog/you-dont-know-what-your-agent-will-do-until-its-in-production>
8. <https://www.louisbouchard.ai/harness-engineering/>
9. <https://atlan.com/know/harness-engineering-vs-prompt-engineering/>
10. <http://www.swebench.com/>
11. <https://openai.com/index/klarna/>

# Contact Us



02-6953-5427



hello@msap.ai



[www.msap.ai](http://www.msap.ai)



## MSAP.ai Blog

최신 기술 트렌드와  
유용한 팁들을 가장 먼저  
만나보세요.



## MSAP.ai eBook

이제 나도 MSA 전문가  
개념부터 실무까지



## YouTube

클라우드 기반 기술과  
인프라 전략을 다루는  
전문 채널