

# 메타-오케스트레이터 개선 계획 v3.0

일자: 2025-10-13

참고 기반 자료:

- 원본 계획: meta-orchestrator-improvement-analysis.md
- Anthropic 다중 에이전트 베스트프랙티스 (scalable.pdf)
- 한국 기술 감사 리포트: 멀티에이전트 아키텍처 전략 분석

분석 방법: 순차적 사고 (10단계) + 다중 문서 교차 검증

## 요약 (Executive Summary)

이 문서는 메타-오케스트레이터 프로젝트의 **생산 환경급 개선 계획**을 제시합니다. 원본 개선안이 에이전트 등록 문제와 SDK 패턴 준수 등을 정확히 짚어냈지만, **프로덕션급 멀티에이전트 시스템에 필수적인 몇 가지 원칙**이 누락되었습니다. 본 계획에서는 이를 보완하여 **구체적 구현 패턴**과 **우선순위별 실행 항목**을 제시합니다. 특히 강조되는 개선 사항은 다음과 같습니다:

1. **최소 권한 원칙 준수** - 각 에이전트에 필요한 최소한의 도구만 할당 (불필요한 웹검색 도구 제거 등) <sup>1</sup> <sup>2</sup>
2. **병렬 실행 패턴 도입** - 병렬 하위 에이전트 실행으로 지연 시간 최대 90% 감소 <sup>3</sup>
3. **컨텍스트 관리 전략** - 작업 단계별로 메모리 분류 및 주기적 요약/정리로 **컨텍스트 오염** 방지 <sup>4</sup> <sup>5</sup>
4. **에러 처리 및 복구 프로토콜** - 자동 재시도 제한, 3회 실패 시 인간 개입 등 **에스컬레이션 정책** 도입 <sup>6</sup> <sup>7</sup>
5. **도구 권한 위반 해결** - 예를 들어 knowledge-builder 에이전트가 웹검색 툴을 사용하던 문제 제거 등 보안 강화

이러한 개선안은 **Claude 공식 SDK 문서** <sup>1</sup> 와 **Anthropic 엔지니어링 블로그** <sup>3</sup> <sup>4</sup>, **Kenny Liao 튜토리얼** 등 권위있는 자료의 원칙을 토대로 합니다. 각 항목마다 구현 전후의 코드 예시와 **검증 기준**을 함께 제공하여 **즉시 실행 가능**한 형태로 제시합니다.

## Part 1: 주요 발견 사항 (원본 계획 외 추가 문제들)

### Issue #1: 도구 권한 위반 (최고 우선 순위)

**문제점:** knowledge-builder 에이전트가 "최소 권한 원칙(Least Privilege)"을 위반하여 불필요한 **웹 검색/문서 조회 도구**를 보유하고 있습니다 <sup>2</sup>. 각 에이전트는 담당 업무에 꼭 필요한 도구만 가져야 함에도, 현재 knowledge-builder가 인터넷 검색 관련 MCP 도구들을 포함하고 있습니다.

현재 상태 (코드):

```
# agents/knowledge_builder.py
tools = [
    'Read', 'Write', 'Edit', 'Grep', 'Glob',
    'TodoWrite',
    'mcp__brave-search__brave_web_search', # 웹검색 (불필요)
    'mcp__context7__resolve-library-id',   # 문서 ID 해석 (불필요)
```

```
'mcp__context7__get-library-docs',      # 문서 내용 조회 (불필요)
]
```

- 영향: "분리된 역할" 원칙 훼손 - 원래 research-agent만 웹 조사를 하고 knowledge-builder는 파일 생성에 집중해야 합니다. 권한이 겹치면 **공격 표면 증가** (예: knowledge-builder가 웹에서 프롬프트 인젝션될 위험) 및 **유지보수 복잡도**가 상승합니다.

개선 방안: knowledge-builder의 도구 목록에서 인터넷 검색 관련 MCP 도구를 제거하고 **파일 시스템 조작에 한정된 도구만** 유지합니다. Claude SDK에서는 `allowedTools` 설정을 통해 특정 도구만 허용하는 세밀한 제어가 가능하므로 이를 활용합니다 <sup>1</sup>. 수정 후 코드는 다음과 같습니다:

```
# agents/knowledge_builder.py (수정 후)
tools = [
    'Read', 'Write', 'Edit', 'Grep', 'Glob', # 파일 시스템 조작만 허용
    'TodoWrite',                             # 작업 플래닝 관련
    # 웹 검색 관련 도구 제거됨 (최소 권한 원칙 준수)
]
```

그리고 research\_agent.py에는 여전히 웹 검색 도구들을 포함시키고, knowledge-builder로 옮기지 않음을 명확히 합니다:

```
# agents/research_agent.py (변경 없음 - 원래 웹검색 도구 보유)
tools = [
    'mcp__brave-search__brave_web_search', # 웹 검색 (research-agent에만 허용)
    'mcp__context7__resolve-library-id',    # 문서 ID 해석 (research-agent 용)
    'mcp__context7__get-library-docs',      # 문서 내용 조회 (research-agent 용)
    'Read', 'Write',                       # 파일 입출력 (JSON 리포트 저장 등)
    'TodoWrite',
]
```

Claude SDK 가이드에 따르면 서브에이전트마다 필요한 도구만 허용함으로써 **의도치 않은 행동을 예방**할 수 있습니다 <sup>2</sup>. 이러한 권한 축소 조치는 보안 면에서 필수이며, 튜토리얼에서도 "에이전트에 꼭 필요한 도구만 허용하고 permissionMode를 적절히 설정"할 것을 권고합니다 <sup>8</sup>.

**워크플로우 변경:** 이제 메타-오케스트레이터는 다음 단계로 동작합니다: - ① research-agent가 필요한 조사를 수행하고 JSON 리포트를 생성 (`Write` 툴로 `/tmp/`에 저장). - ② knowledge-builder는 메타-오케스트레이터로부터 해당 JSON 데이터를 **전달받아** 마크다운 지식을 생성합니다 (파일 생성 전용 도구만 사용). - 직접 웹 검색 권한 없이 **Task** 도구로 리서치 결과를 전달받는 구조로 변경합니다 <sup>9</sup>. 이로써 각 에이전트의 역할이 분리되고 보안이 강화됩니다.

#### 검증 항목:

- [ ] knowledge-builder에서 `mcp__brave-search__*` 및 `mcp__context7__*` 도구 제거 확인
- [ ] knowledge-builder의 시스템 프롬프트를 "리서치 결과를 입력으로 받아 파일을 생성한다"는 내용으로 업데이트
- [ ] 전체 E2E 시나리오 테스트: research-agent → JSON 생성 → knowledge-builder에 Task 전달 → 마크다운 생성까지 정상 동작 확인 (웹 검색 시도 없는지 확인)

## Issue #2: 병렬 실행 패턴 미구현 (높은 우선 순위)

**문제점:** 원본 계획에서 "병렬 실행 패턴(concurrent pattern)"의 필요성을 언급했지만, **구체적인 구현 방법**은 제시되지 않았습니다. 현재 메타-오케스트레이터 프롬프트에 병렬 처리 가능성을 암시하는 내용은 있으나, 실제 코드에서는 여전히 **순차적(Task 호출 후 대기)**으로 에이전트를 실행하고 있습니다. 이로 인해 전체 수행 시간이 불필요하게 길어집니다.

**권위 자료 인사이트:** Anthropic 엔지니어링 블로그에 따르면, Claude Agent SDK의 **하위 에이전트(subagent)** 기능을 활용하면 여러 작업을 동시에 실행하여 **복잡한 작업의 전체 지연 시간을 극적으로 단축**할 수 있습니다 <sup>3</sup>. 실제로 3~5개의 에이전트를 병렬로 돌렸을 때 처리 시간이 최대 90%까지 감소한 사례가 보고되었습니다 <sup>3</sup>.

**현재 격차:** - 메타-오케스트레이터 설계 상 병렬 처리가 가능하다고 언급하지만 **"어떻게 구현하는지"** 구체적인 지침이나 코드 예시가 없습니다. - **배치 크기**나 동시 실행 에이전트 수에 대한 가이드라인이 없어, 개발자가 최적의 병렬화 전략을 알기 어렵습니다.

**개선 방안:** 메타-오케스트레이터의 시스템 프롬프트와 예제 코드에 **병렬 Task 호출** 방법을 추가 명시합니다. Claude Agent SDK에서는 하나의 메시지에 여러 `Task(...)` 호출을 나열하면 **자동으로 병렬 실행**됨이 문서화되어 있습니다 <sup>10</sup>. 이를 활용한 올바른/잘못된 구현 예시를 시스템 프롬프트에 포함시켜, 프롬프트 자체가 개발 가이드 역할을 하도록 합니다:

### ## 병렬 실행 패턴 가이드

#### **\*\*권장 구현 (병렬)\*\*:**

```
```python
# 여러 하위 작업을 한 메시지에서 실행 -> Claude가 병렬로 처리
Task(agent="research-agent", prompt="개념 A에 대한 조사")
Task(agent="research-agent", prompt="개념 B에 대한 조사")
Task(agent="research-agent", prompt="개념 C에 대한 조사")
# 위 모든 Task가 동시에 실행되며 결과를 기다린 후 한꺼번에 처리
```

#### **비권장 구현 (순차):**

```
# 하나의 작업이 끝난 후 다음 Task를 호출하는 방식 -> 비효율적
result1 = Task(agent="research-agent", prompt="개념 A 조사")
# 이 방식은 각 호출마다 기다리므로 병렬 실행이 안 됨
result2 = Task(agent="research-agent", prompt="개념 B 조사")
...
```

위와 같이 프롬프트에 명시하면 에이전트 스스로도 병렬 패턴을 따르도록 유도할 수 있습니다. 또한 **\*\*배치 크기 지침\*\***도 추가합니다:

- **\*\*권장 동시 작업 수\*\*:** 한 번에 3~5개의 하위 에이전트를 병렬 실행 (Anthropic 사례 기반 <sup>3</sup>)
- **\*\*예시\*\*:** 현재 57개의 토폴로지 개념을 순차 처리한다면, 이를 4~5개씩 \*12개의 배치\*로 나누어 각 배치를 병렬 실행합니다. 배치 간에는 순차 진행.

```
```plaintext
Batch 1: 개념 1~5 → 5개 하위 research-agent 병렬 실행
```

Batch 2: 개념 6~10 → 5개 병렬 실행  
...  
Batch 12: 개념 56~57 → 2개 병렬 실행 (마지막 배치)

이런 식으로 구현하면 원래 **5개 개념 순차 처리에 5분 걸리던 작업이 1분 내에 완료**될 수 있습니다 ③.

#### 검증 항목:

- [ ] 메타-오케스트레이터 시스템 프롬프트에 "병렬 실행 패턴" 섹션 추가 (위 코드 예시 포함)
- [ ] 5개의 개념에 대한 research-agent 병렬 실행 테스트: 순차 대비 지연 시간이 현저히 감소하는지 확인 (목표: >50% 단축)
- [ ] 동시 실행 에이전트 수 증가에 따른 Claude 응답 안정성 확인 (필요시 `--max-workers` 등 조정)

### Issue #3: 컨텍스트 오염 방지 전략 부재 (중간 우선 순위) ●

**문제점:** 메타-오케스트레이터가 오랜 시간 다단계로 동작할 경우, 대화/상태 컨텍스트가 점점 비대해지고 **관련 없는 정보가 누적**될 위험이 있습니다. 한국어 감사 리포트에서는 이를 "**컨텍스트 오염**" 현상으로 지적했는데, 시간이 지날수록 상태 객체에 불필요한 정보가 쌓여 후속 에이전트의 판단을 흐리는 문제를 말합니다. 현재 원본 계획에서는 memory-keeper 도구를 언급하지만, **구체적으로 무엇을 어떻게 저장/정리할지**에 대한 전략은 없습니다.

**Claude SDK 기준 개념:** SDK 공식 문서에서도 장시간 실행 시 **맥락 관리(context maintenance)**가 중요하다고 강조합니다. 이를 위해 **서브에이전트 별로 컨텍스트를 격리**하여 불필요한 정보 전파를 막고 ④, **자동 요약(compaction)** 기능을 통해 컨텍스트 길이가 한계에 달하면 이전 내용을 요약하도록 지원합니다 ⑤. 하지만 이러한 기능도 **체계적인 사용 전략**이 없으면 효과적이지 않습니다.

**현재 격차:** - 메모리 관리 도구(`mcp__memory-keeper__*`)는 존재하지만 **카테고리화나 주기적 정리 규칙**이 정의되어 있지 않습니다. - 각 단계별 어떤 정보를 저장하고, 언제 삭제하거나 요약할지 지침이 없어 컨텍스트가 기하급수적으로 늘어날 수 있습니다.

**개선 방안:** 메타-오케스트레이터의 프롬프트에 "**컨텍스트 관리 프로토콜**" 섹션을 추가하고, 다음과 같은 체계를 명시합니다:

#### ## 컨텍스트 관리 프로토콜

**\*\*문제\*\*:** 긴 워크플로우에서는 관련 없는 정보가 축적되어 이후 에이전트들을 혼란스럽게 할 수 있음.

**\*\*해결책\*\*:** 구조화된 memory-keeper 사용 및 주기적 요약

**\*\*1. 모든 컨텍스트 저장 항목에 카테고리 지정:\*\***

```
```python
```

```
# 현재 워크플로우 상태 저장 (session-state)
```

```
mcp__memory-keeper__context_save(  
    key="current-workflow-state",  
    value={"phase": "research", "done": 10, "remaining": 47},  
    category="session-state",  
    priority="high"
```

```
)
```

### # 에이전트 성능 기록 (agent-performance)

```
mcp__memory-keeper__context_save(  
    key=f"agent-{agent_name}-{timestamp}",  
    value={"duration": 45.3, "success": True, "quality": 8.5},  
    category="agent-performance",  
    priority="medium"  
)
```

### # 에러 발생 내역 (errors)

```
mcp__memory-keeper__context_save(  
    key=f"error-{timestamp}",  
    value={"agent": agent_name, "error": "API timeout", "retry": 1},  
    category="errors",  
    priority="high"  
)
```

## 2. 주기적 컨텍스트 정리:

- 예) 10회 에이전트 호출 후 실행: 1. category=="session-state"인 항목들을 모두 불러와 완료된 작업 요약
- 2. 개별 작업 기록들을 삭제하고 요약 결과만 남김 (중요 정보는 유지)
- 3. 새로운 요약 결과를 session-state에 저장 (과거 상태 압축)

## 3. 필요한 컨텍스트 선별 조회:

### # 현재 단계에 필요한 컨텍스트만 조회 (session-state의 최신 중요 정보 5개만)

```
mcp__memory-keeper__context_get(  
    category="session-state",  
    priorities=["high"],  
    limit=5  
)
```

위처럼 \*\*컨텍스트를 범주별로 관리\*\*하면, 불필요하게 커진 대화 기록을 단계별로 정리하여 \*\*컨텍스트 드리프트(Context Drift)\*\*를 막을 수 있습니다 <sup>4</sup> . 또한 compaction 기능(자동 요약)을 SDK에서 지원하므로 필요시 활성화하여 컨텍스트 한계 도달 전 요약을 수행합니다 <sup>5</sup> .

### \*\*검증 항목\*\*:

- [ ] 메타-오케스트레이터 프롬프트에 "컨텍스트 관리" 섹션 추가 및 위 가이드 적용
- [ ] 컨텍스트 카테고리 분류 체계 정의: 예) `session-state`, `agent-performance`, `errors`, `tasks`, `progress` 등
- [ ] 20회 이상의 에이전트 호출이 포함된 시나리오 실행 후 \*\*memory-keeper\*\*에 불필요한 오래된 항목이 남지 않았는지 확인 (요약 및 삭제가 이뤄졌는지)
- [ ] 장기간 실행 시에도 새로운 에이전트가 혼선 없이 필요한 정보만 받아가는지 확인 (컨텍스트 격리 효과 검증)

---

### Issue #4: 에러 핸들링 및 복구 프로토콜 부재 (높은 우선 순위)

**\*\*문제점\*\***: 현재 시스템에는 **\*\*에이전트 실패나 무한 루프 발생 시 대처 전략\*\***이 정의되어 있지 않습니다. 원본 계획에서 헬스체크 언급은 있었지만, 구체적인 재시도 횟수 제한이나 실패 누적에 따른 조치, 사용자 개입 절차 등이 없습니다. 그 결과 에이전트가 반복해서 실패하거나, 외부 API 오류로 멈추는 상황에 효과적으로 대응하기 어렵습니다.

**\*\*외부 인사이트\*\***: 한국 기술 리포트에서는 다음과 같은 실패 시나리오를 강조했습니다:

- **\*\*반복 오류 루프\*\***: 에이전트가 같은 오류를 계속 반복해 **\*\*무한 재시도\*\***하는 문제
- **\*\*도구 실행 실패\*\***: 외부 API 실패 등으로 인한 작업 중단
- **\*\*프롬프트 인젝션\*\***: 외부 데이터에 악의적인 지시가 포함되어 에이전트가 오동작

Claude SDK는 이러한 상황에 대비해 **\*\*빌트인 에러 핸들링과 세션 관리\*\***를 제공한다고 명시하고 있습니다<sup>6</sup>. 하지만 개발자가 명시적으로 정책을 지정해야 효과적으로 활용할 수 있습니다. 예컨대, Anthropic는 **\*\*고객 지원 에이전트\*\*** 예시에서 "필요시 인간에게 에스컬레이션"하는 것을 언급합니다<sup>7</sup>.

**\*\*개선 방안\*\***: 메타-오케스트레이터 프롬프트에 **\*\*"에러 처리 및 복구 프로토콜"\*\*** 섹션을 추가하고, 다음과 같은 정책을 명문화합니다:

```
```markdown
## 에러 처리 & 복구 프로토콜

**1. 재시도 횟수 추적:**

```python
# state 딕셔너리에 에러 횟수 기록
error_count = state.get('error_count', {})
agent_task_id = f"{agent_name}:{task_id}"

if agent_task_id not in error_count:
    error_count[agent_task_id] = 0

# 에이전트 호출 실패 시
error_count[agent_task_id] += 1

# memory-keeper에 누적 실패 횟수 저장
mcp__memory-keeper__context_save(
    key="error-tracking",
    value=error_count,
    category="errors",
    priority="high"
)
```

## 2. 에스컬레이션 정책:

- 1번째 실패: 즉시 같은 내용으로 자동 재시도
- 2번째 실패: 프롬프트를 보정하여 재시도 (예: 입력 단순화)
- 3번째 실패: 인간 사용자에게 에스컬레이션 (작업 중단 및 알림)

```

if error_count[agent_task_id] >= 3:
    print(f"⚠ 에이전트 {agent_name} 작업 {task_id} 3회 연속 실패 - 사용자 개입 필요")
    # 사용자에게 알림을 보내거나 대기 (Human-in-the-loop)
    Task(agent="notify-user", prompt=f"Agent {agent_name} needs assistance on task {task_id}")
    halt_further_actions()

```

### 3. 주요 실패 시나리오별 조치:

- **API 타임아웃**: (예: Brave Search나 Context7) → **지수 백오프** 적용해 1회 재시도. 그래도 실패하면 캐시된 결과 활용 또는 쿼리 단순화.
- **도구 출력 포맷 오류**: (예: research-agent의 JSON 파싱 실패) → quality-agent가 감지 시 해당 에이전트에 오류 피드백 전달 후 재시도.
- **리소스 고갈**: (예: 컨텍스트 창 한계 도달, 과도한 병렬 요청) → 워크플로우 일시 정지, `mcp__memory-keeper__compact` 또는 앞서 정의한 요약 절차 실행 후 재개.

### 4. 에러 로깅 (모든 오류에 대한 로그 기록):

```

mcp__memory-keeper__context_save(
    key=f"error-{timestamp}",
    value={
        "agent": agent_name,
        "task": current_task_desc,
        "error_type": str(type(error)),
        "message": str(error),
        "retry_count": error_count[agent_task_id],
        "context_snapshot": get_current_context()
    },
    category="errors",
    priority="high"
)

```

위와 같은 정책을 통해 **\*\*무한 오류 반복을 방지\*\***하고, 3회 이상 실패 시에는 **\*\*신속하게 인간 검토를 요청\*\***하는 등 안전 장치를 마련합니다 <sup>7</sup>. 또한 Claude SDK의 **\*\*Hook 기능\*\***을 활용하여 위험한 작업 (예: 파일 삭제, 이메일 발송 등)은 사전에 사용자 확인을 거치도록 구현할 것을 권고합니다 <sup>11</sup>. 모든 오류 및 결정은 memory-keeper에 구조화된 로그로 저장하여 사후 분석과 모니터링이 가능하도록 합니다.

**\*\*검증 항목\*\***:

- [ ] 메타-오케스트레이터 프롬프트에 "에러 처리" 섹션 추가 및 재시도/에스컬레이션 로직 구현
- [ ] 인위적인 오류 시나리오를 만들어 3회 실패 시 실제로 작업이 중단되고 **\*\*notify-user\*\*** 등의 경고가 발생하는지 확인
- [ ] memory-keeper에 오류 로그가 의도한 형식으로 기록되고 있는지 확인 (키, 에러유형, 메시지, 재시도 횟수 등)
- [ ] 프롬프트 인젝션 공격 시나리오를 시뮬레이션하여, 권한 없는 에이전트가 위험한 도구를 사용하려 하지 않는지 검사 (최소 권한 설정과 흑으로 대비)

---

## Part 2: 에이전트 기능 매트릭스 (도구 권한 맵)

**\*\*문제\*\***: 원본 계획에는 각 에이전트가 왜 특정 도구들을 가지고 있는지 **\*\*명시적 설명이나 문서화\*\***가 없습니다. 어떤 에이전트에 무슨 도구를 부여할지 명확한 기준이 없다면, 권한 부여에 일관성이 없어지고 **\*Least Privilege\*** 원칙도 흐려질 수 있습니다.

**\*\*해결\*\***: **\*\*에이전트 기능 매트릭스(Agent Capability Matrix)\*\***를 작성하여, 모든 에이전트의 도구 보유 현황과 그 **\*\*이유\*\***를 문서화합니다. 이는 권한 남용을 막고 신규 개발자도 구조를 이해하기 쉽게 합니다.

### ### 도구 카테고리 분류

<b>**카테고리**</b>	<b>**포함 도구**</b>	<b>**설명**</b>
`internet` (인터넷)	`mcp__brave-search__`, `mcp__context7__`	외부 웹검색, 문서 DB 조회
`filesystem` (파일시스템)	`Read`, `Write`, `Edit`, `Grep`, `Glob`	파일 읽기/쓰기/편집/검색 등
`validation` (검증)	`Read`, `Grep`, `Glob`	(**읽기 전용**) 품질 검증, 내용 검색 (파일 수정 불가)
`compute` (연산)	`Bash` (또는 Python 실행 MCP 등)	수학 계산, 코드 실행 등
`planning` (계획)	`TodoWrite`, `mcp__sequential-thinking__`	TODO 리스트 작성, 추론 보조
`memory` (메모리)	`mcp__memory-keeper__`	컨텍스트 영구 저장/조회
`delegation` (위임)	`Task`	하위 에이전트 호출/실행

위와 같이 도구를 카테고리로 묶으면 **\*\*각 에이전트가 어떤 영역의 작업을 하는지\*\*** 명확해집니다. 공식 문서에서도 유사한 권장 조합이 있으며, 예를 들어 **\*\*문서 검토 에이전트\*\***는 Read/Grep만 주어 **\*\*수정 권한을 제한\*\***하는 식으로 설계합니다 2 .

### ### 에이전트별 도구 할당 (최소 권한 적용)

- **\*\*메타-오케스트레이터 (Main Orchestrator)\*\***: `Task` (필수, 하위 호출), `Read`/`Write`/`Edit`/`Grep`/`Glob` (파일 조작), `TodoWrite` 및 `mcp\_\_sequential-thinking\_\_` (계획/추론), `mcp\_\_memory-keeper\_\_` (메모리 조작). **\*\*인터넷 도구는 직접 사용 안 함\*\*** (웹 조사는 research-agent에 위임).

- **\*\*research-agent (연구 에이전트)\*\***: `mcp\_\_brave-search\_\_`, `mcp\_\_context7\_\_` (외부 지식 검색), `Read` (입력 파일 읽기만), `Write` (조사 결과를 JSON으로 출력, 제한된 폴더에), `TodoWrite` (추가 조사 계획). **\*\*파일 수정 도구 없음\*\*** (읽기/출력만).

- **\*\*knowledge-builder (지식 생성 에이전트)\*\***: `Read`/`Write`/`Edit`/`Grep`/`Glob` (파일 시스템 접근), `TodoWrite` (계획). **\*\*인터넷 도구 없음\*\*** - 오직 전달받은 데이터로 지식 파일 생성 (앞서 Issue #1 개선 적용).

- **\*\*quality-agent (품질 검증 에이전트)\*\***: `Read`/`Grep`/`Glob` (파일을 읽어 검증만 함), `TodoWrite` (계획). **\*\*`Write`나 `Edit` 없음\*\*** - 보고만 하고 수정은 하지 않음 (사람 검토 또는 다른 에이전트에 수정 요청).

- **\*\*example-generator (예시 생성 에이전트)\*\***: `Read`/`Write`/`Edit` (파일 입출력), `TodoWrite` (계획), 필요시 `Bash` (Python/SymPy 연산 예시 생성 용). **\*\*인터넷 도구 없음\*\***.



- **dependency-mapper** (의존성 분석 에이전트): ``Read``,`Write``,`Edit``,`Grep``,`Glob`` (코드베이스 분석), ``TodoWrite``. **외부 접근 없음**.

- **socratic-planner** (문답형 플래너 에이전트): ``TodoWrite``, ``mcp__sequential-thinking``만. **파일/인터넷 도구 일절 없음** - 순수하게 대화 기반 계획 수립만 담당.

위 할당표를 통해 **각 에이전트의 역할과 권한 범위**를 명확히 했습니다. 신규 에이전트를 추가할 때에도 이 매트릭스를 참고하여 필요한 카테고리의 도구만 부여하면 됩니다. 또한 코드 상에서도 에이전트 정의 시 하위 클래스나 YAML 설정에서 해당 카테고리만 명시하도록 강제할 수 있습니다. (미래에는 아래 자동 할당 기능도 고려)

### 자동 도구 할당 (향후 개선 사항)

> **참고 코드 (의사 코드)** - Anthropic 커뮤니티 제안 <sup>12</sup> <sup>13</sup>

```
python
tool_by_category = {
    "internet": ["mcp__brave-search__", "mcp__context7__"],
    "filesystem": ["Read", "Write", "Edit", "Grep", "Glob"],
    "validation": ["Read", "Grep", "Glob"],
    "planning": ["TodoWrite", "mcp__sequential-thinking__"],
    "memory": ["mcp__memory-keeper__"],
    "delegation": ["Task"]
}

# 각 에이전트가 선언한 categories 기반으로 tools 자동 할당
for name, agent in agents_registry.items():
    tools = []
    for category in agent.categories:
        tools += tool_by_category.get(category, [])
    agent.tools = list(set(tools))
}
```

향후 에이전트 정의에 `categories = [...]`만 해주면, 위와 같은 로직으로 자동으로 최소 권한의 도구 세트를 할당하도록 구현할 수 있습니다.

## Part 3: 구현 로드맵 (우선순위별 단계)

각 개선 사항을 **4개의 단계(Phase)**로 나누어 우선순위에 따라 실행합니다. 각 단계마다 예상 소요 시간과 **검증 기준 (Definition of Done)**을 명시하여, 단계 완료 여부를 명확히 판단할 수 있도록 합니다.

### Phase 1: Critical Fixes (즉시 시행 - 이번 주)

**우선순위: 가장 높음 (Critical)** - 보안 및 기능상의 중대한 결함 해결이 목표입니다.

1. 도구 권한 위반 수정 (Issue #1 해결)
2. [ ] knowledge-builder에서 Brave Search MCP 제거

3. [ ] knowledge-builder에서 Context7 MCP 제거
4. [ ] knowledge-builder 시스템 프롬프트 수정 (데이터 입력 받아 파일 생성하도록)
5. [ ] 워크플로우 수정: research-agent가 JSON 생성 → knowledge-builder가 Task로 입력 받아 처리
6. **예상 소요:** 2시간
7. **검증:** E2E 테스트 (research-agent 웹검색 → JSON → knowledge-builder 파일 생성) 통과. knowledge-builder 실행 로그에 웹검색 시도 없음 확인.
8. **병렬 실행 패턴 추가 (Issue #2 해결)**
  9. [ ] 메타-오케스트레이터 프롬프트에 "병렬 실행 패턴" 섹션 추가
  10. [ ] Anthropic 사례의 3~5 에이전트 병렬 벤치마크 내용 포함 (90% 레이턴시 감소 언급)
  11. [ ] 병렬 Task 호출 코드 예시 추가 (올바른 vs 잘못된 방법)
  12. [ ] 배치 크기 가이드라인 추가
  13. **예상 소요:** 1시간
  14. **검증:** 한 번에 5개의 concept에 대해 research-agent 병렬 실행 실험 → 기존 순차 대비 응답 시간이 절반 이하로 감소하면 통과.
15. **에러 처리 프로토콜 구현 (Issue #4 해결)**
  16. [ ] 메타-오케스트레이터 프롬프트에 "에러 처리 & 복구" 섹션 추가
  17. [ ] 재시도 횟수 추적 로직 구현 (`error_count` 상태 활용)
  18. [ ] 3회 실패 시 human 에스컬레이션 로직 구현 (ex: `notify-user` Task 또는 경고 출력)
  19. [ ] memory-keeper 오류 로깅 추가 (에이전트, 에러메시지, 재시도횟수 등)
  20. **예상 소요:** 3시간
  21. **검증:** 강제로 실패를 유발하는 시나리오 실행 → 3번째 실패 후 실제로 human 개입 요청이 트리거되는지 확인. memory-keeper에 오류 로그가 저장되었는지 확인.
22. **컨텍스트 관리 전략 추가 (Issue #3 해결)**
  23. [ ] 메타-오케스트레이터 프롬프트에 "컨텍스트 관리 프로토콜" 섹션 추가
  24. [ ] 컨텍스트 카테고리 분류 체계 정의 (예: session-state, agent-performance, errors, tasks 등)
  25. [ ] 일정 호출 주기마다 메모리 정리/요약 가이드 추가 (예: 10회마다)
  26. **예상 소요:** 2시간
  27. **검증:** 20회 이상 에이전트 호출 후 memory-keeper에서 오래된 세션 항목이 요약/삭제되고 핵심만 남았는지 확인. 대화 컨텍스트가 불필요하게 커지지 않았는지 로그 점검.

#### Phase 1 완료 기준 (DoD):

- 전 구간 E2E 테스트 시나리오 모두 성공 (기존에 문제 없던 것 회귀 없음)
- knowledge-builder 에이전트가 더 이상 웹검색/문서조회 도구를 갖고 있지 않음
- 5개 에이전트 병렬 실행 시, 순차 대비 처리 시간이 **50% 이상 단축**됨 확인
- 에이전트 3회 연속 실패 상황에서도 시스템이 중단되지 않고 안전하게 사용자介入으로 전환됨 (크래시 없이 처리)

## Phase 2: Enhanced Capabilities (높은 우선 순위 - 다음 주) ●

**우선순위:** 높음 (High) - 1단계보다는 급하지 않으나, 시스템의 편의성과 성능 가시성을 향상시키는 작업들입니다.

### 1. 에이전트 버전 관리 및 변경 기록

2. [ ] 모든 에이전트 파일에 `VERSION` 필드 추가 (예: `__version__ = "1.0.0"`)
3. [ ] `LAST_UPDATED` 필드 추가 (마지막 수정 날짜)
4. [ ] 각 에이전트 코드 상단이나 주석에 **CHANGELOG** 섹션 추가 (주요 변경 내역 기록)
5. **예상 소요:** 1시간

6. **검증:** 7개 모든 에이전트 파일에 버전/날짜/변경내역 주석이 존재하는지 확인.

## 7. 에이전트 능력 매트릭스 문서화

8. [ ] `AGENT_CAPABILITIES.md` 신규 작성
9. [ ] 각 에이전트의 도구 카테고리 및 실제 도구 목록을 표로 정리 (Part 2 내용 활용)
10. [ ] 왜 그 도구들이 필요한지 각 항목에 justification 추가
11. **예상 소요:** 2시간
12. **검증:** 매트릭스 문서의 내용이 실제 코드의 에이전트 정의와 100% 일치하는지 크로스 체크.

## 13. 성능 모니터링 기능 강화

14. [ ] 에이전트 실행 성능을 structured logging (구조화된 로그)으로 기록
15. [ ] 예: 각 에이전트 완료 시 `mcp__memory-keeper__context_save(category="agent-performance", ...)` 로 duration, 성공여부, 품질점수 등을 저장
16. [ ] 추후 분석을 위해 로그를 한 곳에 모으는 전략 마련 (Phase 4에서 중앙 로깅 도입 예정)
17. **예상 소요:** 2시간
18. **검증:** 임의로 10회 정도 다양한 에이전트를 실행한 후 memory-keeper에서 **10개의 성능 레코드** (`agent-performance` 카테고리)에 기간, 성공여부 등이 기록되었는지 확인.

### Phase 2 완료 기준:

- 모든 에이전트 코드에 버전/변경기록이 포함됨 (개발자가 변경 시 누락하지 않도록 함)
- 능력 매트릭스 문서를 통해 모든 도구 할당의 배경을 이해할 수 있음
- 95% 이상의 에이전트 호출이 성능 로그를 남기고, memory-keeper에서 수집된 데이터를 확인 가능

## Phase 3: Scalability Improvements (중간 우선 순위 - 2주 이내) ●

**우선순위: 중간 (Medium)** - 시스템 확장성과 표준 준수를 높이는 개선입니다.

1. **동적 에이전트 로딩**
2. [ ] 에이전트 자동 발견 기능 구현 (예: `agents/__init__.py` 에서 `importlib` 로 `agents/` 디렉토리 내 모든 `.py` 파일 로드)
3. [ ] `main.py` 에서 개별 에이전트를 하드코딩 임포트하는 부분 제거
4. [ ] 새로운 에이전트 파일을 추가하면 코드 수정 없이도 인식되도록 개선
5. **예상 소요:** 4시간
6. **검증:** 임의로 `agents/test_agent.py` 파일을 만들고 main 코드 수정 없이 registry에 포함되는지 확인. (또는 시뮬레이션을 통해 신규 에이전트 추가 시간을 5분 이내로 달성)

## 7. MCP 표준화 정책 강화

8. [ ] 현재 사용하는 툴 중 **MCP 서버 연결 없이 로컬 구현된 툴**이 있는지 감사 (모든 tool 리스트 검토)
9. [ ] 가이드 문서 작성: "새로운 도구 추가 시 가능하면 MCP 서버 형태로 구현" (왜냐하면 권한 및 격리가 용이함을 명시) 14

10. [ ] 레거시 툴(예: 임시 Bash 스크립트 등)이 있다면 MCP 기반으로 전환 계획 수립

11. **예상 소요:** 3시간

12. **검증:** 문서에 MCP 사용 원칙이 추가되었는지 확인. 신규 툴 PR 시 체크리스트에 "MCP 사용 여부" 포함.

### 13. 헬스 체크 시스템 도입

- [ ] 각 사이클이나 주요 Task 실행 전에 핵심 MCP 툴 서버들의 **가용성/상태 체크**를 수행 (예: `mcp__brave-search__ping` 같은 헬스엔드포인트 호출)
- [ ] 체크 결과 실패시 해당 툴 사용을 건너뛰거나 대체 경로 사용, 로그 남기기
- [ ] memory-keeper에 각 에이전트 실행 시 사전 상태 ( `health-check` 카테고리) 기록
- **예상 소요:** 3시간
- **검증:** 의도적으로 Brave Search MCP를 종료하거나 네트워크 단절 시나리오 실행 -> 시스템이 이를 사전에 감지하여 실패를 예방하거나 적절히 우회하는지 확인. memory-keeper에 헬스체크 결과 (성공/실패)가 로그되었는지 확인.

#### Phase 3 완료 기준:

- 신규 에이전트를 파일 생성만으로 등록 가능 (메인 코드 수정 없이)
- "모델 컨텍스트 프로토콜(MCP)" 준수 원칙이 문서화되고 개발 workflow에 반영 (예: 코드 리뷰 체크포인트) <sup>14</sup>
- 주요 MCP 도구 장애 시, 시스템이 즉각 알리고 무조건 실패하지 않으며 (fail-safe), 헬스 상태가 기록되는 체계가 갖춰짐

## Phase 4: Observability & Maintainability (낮은 우선 순위 - 다음 달)

**우선순위: 낮음 (Low)** - 마지막으로 시스템 운영, 유지보수 편의 개선 사항들입니다.

### 1. 구조화된 로깅 인프라

- [ ] 모든 에이전트 호출 및 주요 이벤트에 대해 JSON 형식의 로그 출력 도입 (예: `{"timestamp": ..., "agent": "...", "event": "start", ...}`)
- [ ] 로그를 수집/조회할 수 있는 중앙 로그 저장 또는 대시보드 전략 마련 (예: ELK 스택, CloudWatch 등 - 선택사항)
- [ ] 외부 모니터링 시스템과 연계 (옵션: Sentry 등 에러 트래킹 연동)
- **예상 소요:** 4시간
- **검증:** 로그 출력 샘플이 JSON으로 잘 생성되고, 파서로 문제없이 읽히는지 확인. (외부 연계는 옵션이나 내부 테스트로 확인)

### 2. 평가 프레임워크 구축

- [ ] 현재 수동으로 수행하는 E2E 테스트를 자동화된 평가(test) 스위트로 전환
- [ ] 각 에이전트 유형별 **품질 지표** 정의 (예: research-agent - 검색 정확도, knowledge-builder - 마크다운 형식 준수 등)
- [ ] CI/CD 파이프라인에 통합하여 코드 변경 시마다 자동 평가 수행, **성능 회귀 검출** 설정
- **예상 소요:** 6시간
- **검증:** 테스트용으로 일부러 품질을 떨어뜨린 변경을 가한 후 CI 실행 -> 평가 스위트가 이를 감지하는지 확인. 1시간 내에 결과 리포트가 생성되는지 확인.

### 3. 환경설정/프롬프트 관리 개선

- [ ] 프롬프트 내용을 코드에서 분리하여 외부 YAML 또는 Markdown 설정 파일로 이전 (`configs/` 디렉토리 등 구성)
- [ ] `.claude/agents/` 폴더 구조를 재점검하여 에이전트 정의와 설정이 한 눈에 들어오도록 정리
- [ ] 동적으로 프롬프트/설정을 불러오도록 코드 수정 (예: `load_prompt(agent_name)` 함수로 파일 읽기)
- **예상 소요:** 5시간
- **검증:** 임의로 프롬프트 내용을 수정한 YAML을 배포 후 코드 변경 없이 적용되는지 테스트. 또한 설정 변경이 git diff로 용이하게 보이는지 확인.

#### Phase 4 완료 기준:

- 모든 에이전트 로그가 JSON 등 구조화된 포맷으로 출력되며, 필요시 필터링/검색이 가능
- 평가 스위트가 CI에서 실행되어 주요 회귀를 자동 탐지 (예: 중요 지표 악화 시 실패)
- 프롬프트 수정이 코드 수정 없이 가능하여, 운영 중에도 설정 튜닝을 신속히 적용 가능

## Part 4: 기존 계획 대비 핵심 개선사항

### 원본 계획에서 잘한 부분

- **에이전트 레지스트리 문제 인지:** 기존 개선안에서 에이전트 등록 방식의 문제점을 잘 짚고 있었습니다 (동적 로딩 등 Phase 3의 8번 항목과 연계됨).
- **Kenny Liao SDK 사용 준수 언급:** Claude Agent SDK의 모범 패턴을 따르는 것의 중요성을 인식하고 있었습니다.
- **메타-오케스트레이터의 조정자 역할 이해:** 메타-오케스트레이터가 하위 에이전트들을 어떻게 조율하는지 큰 그림은 정확히 파악하고 있었습니다.
- **memory-keeper 도구 활용 의식:** 컨텍스트 저장용 도구의 존재를 파악하고 필요성을 언급한 점은 좋았습니다.

### v3.0 계획에서 추가로 보완한 사항

1. **도구 권한 위반 문제 식별:** 원본 계획에서는 knowledge-builder가 불필요한 연구 도구를 갖고 있다는 점을 놓쳤으나, 이번 개선안에서 이를 **최우선 보안 이슈**로 지정하고 해결책을 제시했습니다.
2. **병렬 실행의 구체화:** 원본에서는 병렬 처리가 필요하다고만 했지만, v3.0에서는 **구현 코드 예시와 배치 전략**까지 명확히 하였습니다.
3. **에러 처리 전략 수립:** 이전에는 실패 대응이 추상적이었으나, 이번에는 **재시도 로직, 한계 횟수, human 에스컬레이션**까지 구체적인 프로토콜을 정의했습니다.
4. **컨텍스트 관리 방법론 제시:** memory-keeper 언급에서 더 나아가 **무엇을 어떻게 저장/정리할지** 카테고리 체계와 요약 전략을 도입했습니다.
5. **능력 매트릭스 문서화:** 기존에는 암묵적이던 에이전트-도구 대응을 명문화하여, 권한 부여의 **일관성과 투명성**을 확보했습니다.
6. **우선순위 및 일정 명시:** 기존 계획은 개선 항목별 중요도가 평면적이었지만, v3.0에서는 **단계별 우선순위, 예상 시간, DoD**를 지정하여 실행 계획의 현실성을 높였습니다.

### 프레임워크/패턴 적용 명확화

- **LangGraph vs Claude SDK:** 한국 감사 보고서는 Kenneth Liao의 예제 (LangGraph 기반)로 인해 본 시스템도 노드-엣지 기반이라고 간주했지만, 실제로는 **Claude Agent SDK (AgentDefinition, Task 등)**를 사용하고 있습니다. 이번 계획에서는 감사 보고서의 원칙들(최소 권한, 오류 처리, 상태 관리 등)을 수용하되, **구현은 Claude Agent SDK 패턴으로 수행함**을 분명히 했습니다 (예: Task 호출, MCP 톨 사용 등).

---

## Part 5: 성공 지표 (측정 가능 Metrics)

각 Phase 완료 후 달성해야 할 **객관적 성능/품질 지표**를 정의합니다. 이 지표들은 향후 회귀 테스트 및 지속적인 개선의 기준선이 됩니다.

### Phase 1 지표 (Critical)

- [] **E2E 테스트 통과율**: 100% (현재 기준 100% 미만이라면 완전 통과를 목표)
- [] **도구 권한 위반 건수**: 0건 (현재 1건: knowledge-builder의 인터넷 도구)
- [] **병렬 처리로 인한 지연 감소율**: 5개 에이전트 배치 기준 **50% 이상** (예: 5개 순차 300초 → 병렬 150초 이하)
- [] **에러 복구율**: 3회 연속 실패 상황에서도 **시스템 다운 없이 처리 전환** 100% (모든 치명적 실패를 graceful하게 처리)

### Phase 2 지표 (High Priority)

- [] **에이전트 문서화 수준**: 7/7 에이전트에 버전 및 변경이력 포함 (100%)
- [] **능력 매트릭스 커버리지**: 100% (모든 도구/에이전트 할당이 문서에 포함되고 타당한 justification 제공)
- [] **성능 로깅 비율**: 에이전트 호출 중 **95% 이상**에 대해 duration, success 여부 등의 메트릭이 기록됨 (일부 매우 짧은 trivial 호출 제외하고 거의 전체 로그 커버)

### Phase 3 지표 (Medium)

- [] **확장성**: 신규 에이전트 추가에 걸리는 시간 < **5분** (파일 작성 후 바로 시스템 인식, main 코드 수정 0)
- [] **MCP 준수율**: 새로 추가되는 도구의 **100%**가 MCP 방식으로 구현 또는 기존 라이브러리 MCP 활용 (내장 Bash 등 제외하고)
- [] **헬스체크 적용 범위**: 모든 에이전트 실행 전후로 주요 외부 의존 MCP에 대한 상태 점검이 1회 이상 수행 (로그에 남는 health-check 이벤트 100%)

### Phase 4 지표 (Low)

- [] **로그 커버리지**: **100%** 에이전트 이벤트가 구조화된 로그로 남음 (모든 start/finish, error 등이 JSON 등으로 기록)
- [] **평가 체계**: 코드 변경 후 **1시간 이내**에 자동 평가가 실행되고 결과 확인 (회귀 발견 시 알림)
- [] **설정 관리 효율**: 프롬프트 수정에 필요한 코드 배포 횟수 0회 (모든 프롬프트/시스템 설정 변경이 재배포 없이 가능, 설정 파일만 교체)

---

## Part 6: 위험 평가 및 대응

### 높은 위험 요소

#### 위험 1: 도구 제거로 인한 예기치 않은 부작용

- 영향: knowledge-builder에서 갑작스럽게 웹검색 도구를 제거하면, 혹시 해당 기능을 참조하던 흐름이 깨질 수 있습니다.
- 대응: - 변경 전에 관련 E2E 테스트 케이스를 모두 업데이트/추가하여 영향 범위를 파악합니다.
- 단계적 릴리스: 실제 운영에 바로 반영하기 전에 개발 환경에서 충분한 테스트.
- 만일 문제 발생 시 쉽게 원복할 수 있도록, 기존 knowledge-builder 코드를 백업/브랜치로 보존해 둡니다 (rollback 경로 확보).

### 위험 2: 병렬 실행 복잡성

- 영향: 병렬로 여러 에이전트를 띄우면 **경쟁 상태(race condition)**나 **자원 경합** 문제가 생길 수 있습니다 (예: 동시에 파일 쓰기 시 충돌).
- 대응: - 초기에 작은 배치 (예: 동시 2~3개)로 시험 가동하여 문제 발생 여부를 관찰합니다.
- Claude SDK의 `max_workers` 옵션 등 **타임아웃과 워커 제한** 설정을 활용해 과부하를 방지합니다.
- 프로덕션 투입 전 충분한 부하 테스트를 수행하고, 필요한 경우 동기화 락 또는 작업 큐 도입도 고려합니다.

### 위험 3: 컨텍스트 정리로 인한 정보 손실

- 영향: 주기적 컨텍스트 요약/삭제 과정에서 향후 필요할 정보를 잃어버릴 위험이 있습니다. 지나친 정리는 **필요한 컨텍스트마저 소실**시킬 수 있습니다.
- 대응: - 초기에는 **보수적인 정리 전략**을 채택합니다 (예: 최근 1-2회분만 삭제하고 나머지 유지).
- 자동 요약 내용은 사용자에게 가시적으로 보여줘 검토할 수 있게 합니다.
- 만약 중요한 정보가 사라졌다면 쉽게 복원할 수 있도록 memory-keeper에서 일정 기간 (예: 7일) 원본 보관 후 삭제를 실행합니다 (soft delete → hard delete).

## 중간 위험 요소 ●

### 위험 4: 버전 관리 규율 미준수

- 영향: 개발자들이 에이전트 수정 시 일일이 VERSION이나 CHANGELOG를 업데이트하지 않을 가능성. 시간이 지나면 버전 태그가 실제와 불일치하여 신뢰도 하락.
- 대응: Git pre-commit hook이나 CI에서 에이전트 파일 변경 시 버전 문자열 변경이 있는지 검사하는 스크립트를 추가합니다. 누락 시 PR을 막아 강제하도록 합니다.

### 위험 5: 동적 로딩 구현 버그

- 영향: importlib를 통한 자동 임포트 과정에서 예외가 발생하면, 최악의 경우 전체 메타-오케스트레이터 구동 실패 가능. 또는 특정 OS/환경에서 동작 안 할 수도 있음.
- 대응: - 다양한 환경(Ubuntu, Mac, Windows WSL 등)에서 테스트를 수행합니다.
- 문제가 발생할 경우를 대비해 `try-except`로 감싸고, 실패 시 기존 하드코딩 리스트로 fallback하도록 이중화합니다.
- 로깅을 통해 어떤 파일에서 로딩 실패했는지 추적 가능하게 합니다.

## Part 7: 향후 단계 (액션 아이템)

### 즉시 수행할 작업 (현재 세션)

1. 이 개선 계획을 **memory-keeper에 저장** - 우선 본 문서 내용을 시스템 내 장기 메모리에 기록해두어, 에이전트들이 참조할 수 있도록 합니다. 예를 들어:

```
mcp__memory-keeper__context_save(  
    key="improvement-plan-v3",  
    value=<이 문서 전체>,  
    category="architecture",  
    priority="high",  
    channel="main-workflow"  
)
```

이렇게 하면 차후 에이전트 대화나 계획 수립 시 architecture 카테고리에서 이 문서를 불러 참고할 수 있습니다.

2. **사용자 승인 받기** – 이 계획을 프로젝트 리더/팀에 공유하여 **Phase 1 항목들의 우선순위에** 이견이 없는지 확인합니다. 특히 병렬화나 에러 처리 정책은 시스템 동작에 큰 변화를 주므로, 사전에 팀의 동의를 구하고 세부 조정사항 (예: 재시도 횟수 3회 적절한지 등) 피드백을 반영합니다.

3. **Phase 1 구현 착수** – 승인되었다면 즉시 1단계 작업부터 시작합니다:

4. knowledge\_builder.py 수정하여 도구 리스트 정리 → 테스트

5. 메타-오케스트레이터 프롬프트 업데이트 (병렬, 컨텍스트, 에러 처리 섹션 추가) → 로컬에서 Claude SDK를 이용해 실행 테스트

6. 수정 후 통합 E2E 테스트 돌려보기 → 결과 공유 및 필요시 수정 반복

7. **향후 일정 조율** – Phase 2~4의 작업을 주간 스프린트 계획에 반영합니다. 우선 Phase 1 완료 후 결과를 리뷰하고, 문제 없으면 Phase 2 작업 (다음 주) 착수. 각 단계 후 성취도(metrics) 평가하여 다음 단계 계획을 조정합니다.

---

## Part 8: 결론

이번 개선 계획 v3.0은 메타-오케스트레이터를 **프로덕션 수준의 견고함과 확장성**으로 끌어올리기 위한 **종합 로드맵**입니다. 주요 내용을 요약하면 다음과 같습니다:

- **보안 강화:** 최소 권한 원칙을 철저히 적용하여 에이전트별 도구 권한을 세분화하고, 프롬프트 인젝션 등 공격 벡터를 줄였습니다 (예: knowledge-builder 웹접근 제거) <sup>2</sup>.
- **성능 최적화:** 병렬 하위 에이전트 실행으로 응답 지연을 대폭 감소시키고 <sup>3</sup>, 필요시 배치 실행 기법을 도입했습니다.
- **신뢰성 향상:** 컨텍스트 관리 및 에러 복구 프로토콜을 통해 장시간 동작 시 컨텍스트 혼잡을 예방하고, 오류 발생 시 자동대응 및 인간 개입 절차를 명확히 했습니다 <sup>4</sup> <sup>7</sup>.
- **운영 가시성:** 구조화된 로깅과 성능 모니터링으로 시스템 투명성을 높이고, 평가 스위트를 통해 지속적으로 품질을 관리합니다.
- **로드맵 실행 가능성:** 4단계에 걸친 구현 계획에는 각 단계별 구체적인 작업 내역, 예상 시간, 완료 기준이 포함되어 있어 팀이 **현실적인 일정 관리**를 할 수 있습니다.

모든 개선안에는 Claude 공식 SDK 문서와 실제 사례에서 검증된 패턴을 참고하여 신뢰도를 높였습니다. 예를 들어 **도구 권한 제어**는 SDK의 fine-grained permission 기능을 활용하고 <sup>1</sup>, **병렬/컨텍스트 관리**는 서브에이전트 및 compaction 기능에 근거했으며 <sup>3</sup> <sup>5</sup>, **에러 처리**는 Anthropic 권장 human-in-the-loop 전략을 채택했습니다 <sup>7</sup>.

이 계획이 승인되고 단계적으로 구현된다면, 메타-오케스트레이터 시스템은 **안전성, 효율성, 유지보수성 면에서 한층 성숙한 아키텍처**로 거듭날 것입니다. 이후 Phase 1부터 차근차근 실행에 옮겨 나가겠습니다.

---

## Part 9: 부록 (Appendices)

### A. References (공식 문서 및 자료)

#### Official Documentation:

- [Claude Agent SDK Overview – Fine-grained tool permissions, error handling, etc.] <sup>1</sup> <sup>6</sup>
- [Subagents in Claude SDK – Context isolation and parallelization] <sup>3</sup> <sup>4</sup>



- [Kenny Liao's Claude SDK Tutorial – Best practices \(GitHub\)](#)
- [Anthropic Engineering Blog – “Building agents with the Claude Agent SDK”] 3 5

#### Analysis Documents:

- Anthropic Community PDF: Multi-Agent Systems Best Practices (scalable.pdf)
- Korean Technical Audit Report: Strategic Analysis of Multi-Agent Architecture (2025)
- Original Improvement Plan: meta-orchestrator-improvement-analysis.md

### B. 도구 권한 위반 사례 - 수정 전후 코드

수정 전 (문제 상태): knowledge-builder가 불필요한 인터넷 도구 포함

```
# knowledge_builder.py (before)
tools = [
    'Read', 'Write', 'Edit', 'Grep', 'Glob',
    'TodoWrite',
    'mcp__brave-search__brave_web_search',    # 웹검색 도구 (권한 초과)
    'mcp__context7__resolve-library-id',      # 문서 ID 해석 (권한 초과)
    'mcp__context7__get-library-docs',        # 문서 내용 조회 (권한 초과)
]
```

수정 후 (개선 상태): knowledge-builder는 파일 작업 관련 도구만 보유

```
# knowledge_builder.py (after)
tools = [
    'Read', 'Write', 'Edit', 'Grep', 'Glob',    # 파일 시스템 작업
    'TodoWrite',                                # 일정/계획 관리
    # 웹검색/문서조회 도구 제거됨 (최소 권한 원칙 적용)
]
```

(참고: research\_agent.py는 여전히 인터넷 도구를 포함하며, knowledge-builder는 Task를 통해 연구 결과를 입력 받는 구조로 변경됨.)

### C. 병렬 실행 전후 비교

현재 (순차 실행 - 느림): 5개 개념에 대해 research-agent를 순차 호출한 예

```
시간 0초 → research-agent (개념1 조사 시작) ... [약 60초 소요]
시간 60초 → research-agent (개념2 조사 시작) ... [약 60초 소요]
시간 120초 → research-agent (개념3) ...
시간 180초 → research-agent (개념4) ...
시간 240초 → research-agent (개념5) ...
총 소요 시간 ≈ 300초 (5분)
```

개선 (병렬 실행 - 빠름): 동일 작업을 병렬 Task로 처리

시간 0초 → research-agent (개념1~5 동시 조사 시작, 병렬 실행)  
 시간 60초 → 모든 조사 완료  
 총 소요 시간 ≈ 60초 (1분) ← \*\*80%\*\* 이상 시간 단축!

(실제 성능은 각 에이전트 작업 시간 편차에 따라 다르지만, 병렬화로 큰 폭의 개선 가능 <sup>3</sup>.)

#### D. 컨텍스트 관리 카테고리 예시

아래는 제안된 memory-keeper 컨텍스트 분류 체계의 예시입니다:

카테고리	용도	우선순위	보존 기간
session-state	현재 워크플로우 단계 상태	높음 (High)	세션 종료 시 (또는 요약)
agent-performance	에이전트별 수행 메트릭 (시간, 품질)	중간 (Med)	최근 7일
errors	실패 내역, 재시도 횟수 등	높음 (High)	최근 30일
decisions	중요한 설계/결정 사항	높음 (High)	영구 보존
tasks	TODO 항목, 미완료 작업	중간 (Med)	완료 시까지
progress	완료된 작업, 마일스톤	중간 (Med)	최근 7일

위 표는 각 카테고리에 어떤 정보가 담기고 얼마나 유지될지를 정의한 것입니다. 이를 기반으로 memory-keeper 저장/정리 정책을 구현하면 체계적인 상태 관리가 가능합니다.

문서 버전: 3.0

작성자: Claude Sonnet 4.5 (AI)

상태: 대기 중 - 사용자 승인 및 Phase 1 실행 준비 완료

<sup>1</sup> <sup>6</sup> <sup>14</sup> Agent SDK overview - Claude Docs

<https://docs.claude.com/en/api/agent-sdk/overview>

<sup>2</sup> <sup>9</sup> <sup>10</sup> Subagents in the SDK - Claude Docs

<https://docs.claude.com/en/api/agent-sdk/subagents>

<sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>7</sup> Building agents with the Claude Agent SDK \ Anthropic

<https://www.anthropic.com/engineering/building-agents-with-the-claude-agent-sdk>

<sup>8</sup> <sup>11</sup> How to Create Agents with Claude Agents SDK - Bind AI IDE

<https://blog.getbind.co/2025/10/03/how-to-create-agents-with-claude-agents-sdk/>

<sup>12</sup> <sup>13</sup> meta-orchestrator-improvement-plan-v2.md

[file:///file\\_0000000061b861faa6f660c560569e9a](file:///file_0000000061b861faa6f660c560569e9a)