

## GPU 수가 아닌 MCP 수가 AI 경쟁력: 기업 AI 성숙도의 새로운 기준

기업의 AI 경쟁력을 보유한 GPU 수가 아닌, AI가 실제 업무와 얼마나 연결되었는지를 나타내는 'MCP(Model Context Protocol) 수'로 재정의해야 합니다. 이를 위해 MCP를 표준 인터페이스로 삼아 AI를 '일하는 에이전트'로 만들고, 기존 시스템을 AI 자원으로 전환하여 실질적인 비즈니스 가치를 창출하는 전략을 제시합니다.

#### **Contact Us**



02-6953-5427



hello@msap.ai



www.msap.ai



#### Contents

제1장 서론 - GPU에서 MCP 수 중심의 AI 경쟁력으로	7
1.1 문제 제기: GPU·모델 중심 사고의 구조적 한계	7
1.1.1 도입: 현재의 AI 인프라 논의 분석	7
1.1.2 컴퓨팅 파워 집착이 초래한 AI 인프라 논의의 왜곡	7
1.1.3 GPU·파라미터 수와 실제 업무 생산성 간의 불일치	8
1.1.4 "GPU 몇 장인가?"에서 "AI로 실행되는 업무 수는 얼마인가?"로	8
1.2 Al 경쟁력의 새로운 정의: '업무 접점 수(MCP 수)'	8
1.2.1 도입: 새로운 경쟁력 지표 제시	8
1.2.2 MCP 수 정의: LLM이 직접 호출할 수 있는 프로덕션 업무 기능 수	9
1.2.3 GPU 수 vs MCP 수: 필요 조건과 충분 조건의 구분	10
1.2.4 "모델 보유"에서 "업무 접점 수"로 이동하는 성숙도 관점	11
1.3 본 백서의 목적, 대상 독자, 범위	11
1.3.1 도입: 백서의 가이드 역할 정의	11
1.3.2 MSA·쿠버네티스·클라우드 네이티브에 대한 최소 전제	11
1.3.3 공공·금융·엔터프라이즈 IT 의사결정자가 이 백서를 읽어야 하는 이유	12
1.3.4 이 백서에서 다루지 않는 범위(모델 연구·알고리즘 레벨 논의 등)	12
1.4 백서 구성과 활용 방법	12
1.4.1 도입: 백서의 실용적 가치 강조	12
1.4.2 온프레미스 AI 플랫폼 → MCP 기술 → 성숙도 모델 → 기존 시스템 MCP	
화 → MSAP.ai 순서	13
1.4.3 각 장별 산출물: 평가 지표, 설계 패턴, 참조 아키텍처	13
1.4.4 "나중에"가 아닌 "지금" 무엇을 시작할 것인가에 대한 가이드	13
제2장: 엔터프라이즈 AI 인프라와 온프레미스 AI 플랫폼 전략	14
2.1 엔터프라이즈 AI 인프라 4계층 재정의	15
2.1.1 연산 계층: GPU/TPU, 노드·클러스터 구성	15



2.1.2 모델 계층: 파운데이션 모델·도메인 LLM·서빙 스택	16
2.1.3 컨텍스트 계층: MCP·RAG·데이터 커넥터 계층	16
2.1.4 애플리케이션·업무 계층: 챗봇, 에이전트, 업무 시스템	17
2.2 왜 공공·엔터프라이즈에는 온프레미스 AI 플랫폼이 필요한가	18
2.2.1 AI를 "내부 자산(업무 시스템·데이터)을 지능화하는 플랫폼"으로 보는 관점	18
2.2.2 내부 시스템·데이터 근접성, 지연·가용성·통제 측면의 이점	18
2.2.3 조직 안의 IT 자산을 어떻게 엮느냐가 GPU·LLM 보유보다 중요한 이유	19
2.3 퍼블릭 클라우드 LLM에 내부 데이터를 보내는 리스크	20
2.3.1 데이터 주권·레지던시·규제(공공·금융·개인정보) 관점의 제약	20
2.3.2 벡터DB·로그·모델 피드백에 남는 민감 정보 이슈	20
2.3.3 단일 질의가 조직 전체 규제 위반으로 이어지는 전형적인 패턴들	21
2.4 온프레미스·하이브리드 AI 플랫폼 설계 원칙	21
2.4.1 온프레미스 LLM 서빙 + MCP + RAG + Observability 통합 구조	21
2.4.2 퍼블릭 LLM 활용 범위: 마스킹·프록시·샌드박스 전략	22
2.4.3 온프레미스·프라이빗 클라우드·하이브리드 도입 패턴 비교	22
2.5 "나중에"가 아닌 "지금" 시작해야 하는 이유	24
2.5.1 GPU 가격·기술·규제를 기다리는 전략이 위험한 이유	24
2.5.2 지금 당장 할 수 있는 일: 자산 목록화, 후보 업무 선정, PoC 범위 정의	25
2.5.3 "먼저 MCP 인벤토리를 가진 조직이 AI 성숙도를 선점한다"는 관점	25
제3장 전략적 선택: 왜 '모델 학습'이 아닌 '추론 활용'에 집중해야 하는가	26
3.1 모델 학습(Training)과 추론(Inference)의 기술적 차이	26
3.1.1 Pre-training, Fine-tuning, Instruction Tuning 개념 정리	27
3.1.2 학습 파이프라인과 추론 파이프라인의 비교	27
3.1.3 배치 학습과 온라인 추론의 자원 사용 패턴	28
3.2 파운데이션 모델 학습의 경제학	29
3.2.1 글로벌 빅테크의 LLM 학습 비용과 자본 구조	29
3.2.2 데이터·인력·인프라를 포함한 총소유비용(TCO) 요소	29
3.2.3 공공기관·중견기업이 독자 LLM 개발에서 실패하는 구조적 이유	30



3.3 GPU 증설 중심 접근의 한계와 Idle 자원 문제	30
3.3.1 GPU 수 증가가 곧 업무 자동화 증가로 이어지지 않는 이유	31
3.3.2 활용되지 못하는 연산 자원과 조직 내부 프로세스의 미정비	31
3.3.3 GPU 투자와 MCP·RAG 투자 간의 우선순위 재조정	31
3.4 추론 중심 전략과 MCP 수의 연계	32
3.4.1 MCP 수 = 실제 추론이 개입하는 업무 수	32
3.4.2 MCP 수를 기준으로 한 투자 우선순위: 어떤 업무부터 AI화할 것인가	33
3.4.3 모델 학습에 투자해야 하는 극소수 사례를 구분하는 기준	34
제4장 MCP(Model Context Protocol)의 기술적 정의와 표준 생태계	34
4.1 MCP의 기본 개념과 아키텍처	35
4.1.1 MCP의 목적: LLM과 외부 시스템을 연결하는 개방형 표준	36
4.1.2 MCP의 핵심 구성: Host, Client, Server, Tools, Resources, JSON-RPC	36
4.1.3 요청 라우팅·모델 오케스트레이션·컨텍스트 주입 흐름	37
4.2 기존 연동 방식과 MCP의 차이	38
4.2.1 API·Webhook·플러그인 방식의 N×M 연동 문제	38
4.2.2 MCP 도구 호출 모델이 연동 복잡도를 N+M으로 줄이는 방식	39
4.2.3 단순 HTTP API 호출과 컨텍스트 인식 프로토콜로서 MCP의 차별점	40
4.3 MCP 생태계와 글로벌 표준 동향	41
4.3.1 Anthropic·OpenAl·Microsoft·GitHub 등 주요 벤더의 MCP 채택 현황 .	41
4.3.2 OS·IDE·브라우저·에이전트 런타임 등 플랫폼 레벨 통합 흐름	42
4.3.3 MCP 생태계와 글로벌 표준 동향	43
4.4 엔터프라이즈 도입 시 고려사항	43
4.4.1 인증·인가·감사·네트워크 구조(Host 중심 보안 모델)	44
4.4.2 프롬프트 인젝션·데이터 유출 등 MCP 특유의 보안 위협	44
4.4.3 MCP 버전 전략·거버넌스: 등록·검증·배포 정책	45
제5장 AI 성숙도 모델과 MCP 수 기반 평가 프레임워크	46
5.1 Al 성숙도의 새로운 지표: '업무 접점 수'	46
5.1.1 모델·GPU 보유 중심 성숙도 평가의 한계	47



5.1.2 MCP 수로 측정하는 AI 활용 범위(업무 도메인·기능 수)	47
5.1.3 MCP 수와 자동화율·의사결정 지원 비율의 상관관계	48
5.2 MCP 수 측정 방법론	49
5.2.1 도메인·시스템별 MCP 인벤토리 구축 절차	49
5.2.2 MCP 카테고리 분류(조회·작성·승인·배포·모니터링 등)	50
5.2.3 AI 노출 범위(Coverage)·사용 빈도·성공률 지표 설계	51
5.3 MCP 기반 AI 성숙도 모델	52
5.3.1 Level 0: 실험 단계 - PoC 챗봇·파일 업로드 수준	52
5.3.2 Level 1~2: 일부 도메인 MCP화 - 특정 업무 도메인 자동화 단계	52
5.3.3 Level 3~4: 전사 MCP 포털·AI 운영 내재화 단계	53
5.4 조직별 성숙도 진단 기준과 로드맵	53
5.4.1 공공기관(민원·행정·통계)의 성숙도 진단 항목	54
5.4.2 금융·서비스·제조 분야의 성숙도 진단 항목	54
5.4.3 1년·3년 단위 MCP 로드맵 수립 예시	55
제6장 기존 시스템의 지능형 MCP화 전략 – MSA 관점 레거시 재설계	56
제6장 기존 시스템의 지능형 MCP화 전략 - MSA 관점 레거시 재설계 6.1 '교체'가 아닌 '진화' 전략	56 57
	57
6.1 '교체'가 아닌 '진화' 전략	57 57 58
6.1 '교체'가 아닌 '진화' 전략	57 57 58 59
6.1 '교체'가 아닌 '진화' 전략	57 57 58 59
6.1 '교체'가 아닌 '진화' 전략	57 57 58 59 59
6.1 '교체'가 아닌 '진화' 전략	57 57 58 59 59 60
6.1 '교체'가 아닌 '진화' 전략	57 58 59 59 60 60 61
6.1 '교체'가 아닌 '진화' 전략	57 58 59 59 60 60 61
6.1 '교체'가 아닌 '진화' 전략          6.1.1 모놀리식·3-Tier 레거시 시스템 구조 분석          6.1.2 "버리는 것이 아니라 감싸되, 다시 설계해서 감싸는" 방식의 전환          6.1.3 기존 비즈니스 로직을 AI 추론 호출 대상으로 승격시키는 방법          6.2 기존 자산을 MCP로 노출하는 패턴          6.2.1 REST·gRPC·SOAP API의 MCP 서버화 패턴          6.2.2 배치·스케줄러·RPA 기능의 MCP 도구 추상화          6.2.3 DB 조회·리포트·통계 쿼리의 MCP화 방법          6.3 MSA(Microservice Architecture) 관점 MCP 설계 원칙	57 58 59 59 60 61 61
6.1 '교체'가 아닌 '진화' 전략          6.1.1 모놀리식·3-Tier 레거시 시스템 구조 분석          6.1.2 "버리는 것이 아니라 감싸되, 다시 설계해서 감싸는" 방식의 전환          6.1.3 기존 비즈니스 로직을 AI 추론 호출 대상으로 승격시키는 방법          6.2 기존 자산을 MCP로 노출하는 패턴          6.2.1 REST·gRPC·SOAP API의 MCP 서버화 패턴          6.2.2 배치·스케줄러·RPA 기능의 MCP 도구 추상화          6.2.3 DB 조회·리포트·통계 쿼리의 MCP화 방법          6.3 MSA(Microservice Architecture) 관점 MCP 설계 원칙          6.3.1 DDD 관점에서 기존 API를 MCP로 재설계하는 방법론	57 58 59 59 60 61 61 62 62



6.4.1 RBAC 기반 권한·역할 모델 설계	64
6.4.2 호출 로그·추적·감사를 포함한 Observability 전략	64
6.4.3 쿠버네티스 기반 Auto-scaling·Service Discovery·Resilience 확보	65
6.5 단계별 MCP 전환 로드맵	65
6.5.1 3개월 내 구현 가능한 MCP 후보 도출	66
6.5.2 핵심 도메인(민원·운영·개발 생산성 등) 우선 전환 전략	66
6.5.3 PoC에서 전사 MCP 카탈로그로 확장하는 조직·프로세스 변화	67
제7장: MCP와 RAG로 내부 데이터를 AI 자산으로 전환하는 방법	68
7.1 두 개의 뇌를 가진 AI: 행동을 담당하는 MCP, 지식을 담당하는 RAG	68
7.1.1 AI의 그럴듯한 거짓말, '할루시네이션'과 신뢰의 열쇠	69
7.1.2 RAG 개념과 전형적인 아키텍처(인덱싱·검색·생성)	69
7.1.3 MCP 도구로서 RAG 인덱싱·질의 기능을 노출하는 방식	71
7.2 기존 데이터를 AI가 활용할 수 있게 만드는 절차	71
7.2.1 문서·메일·업무 로그의 정제·구조화·분류	72
7.2.2 임베딩 전략·메타데이터 설계와 보안 태깅	72
7.2.3 하나의 인덱스를 여러 에이전트·업무에서 재사용하는 패턴	73
7.3 공공·엔터프라이즈 데이터와 규제·보안	74
7.3.1 개인정보·기밀 정보·규제 데이터의 취급 기준	74
7.3.2 중앙집중형 벡터DB vs 소스 시스템 실시간 조회(에이전트형 아키텍처)	75
7.3.3 RAG·MCP를 활용한 보안·감사·접근제어 전략	77
7.4 MCP·RAG 통합 워크플로와 도입 단계	77
7.4.1 질의 → 검색 → MCP 툴 호출 → 결합 응답 워크플로	78
7.4.2 민원 응대, 청약·통계 질의, 운영 로그 분석 등 도메인별 활용 예시	78
7.4.3 RAG·MCP 결합 효과 측정 지표(정확도·처리 시간·재작업률)	79
제 8장 MSAP.ai란 무엇일까요? Al Native 플랫폼	80
8.1. 왜 'AI 플랫폼'이 필요한가요?	80
8.1.1. MSAP.ai의 정체: 'AI 발전소'를 위한 '스마트 전력망'	80
8 1 2   소프트웨어 개발이 모드 과정을 ΔI로 스마트하게!	ຂ1



8.1.3. 공공 및 엔터프라이즈 도입 모델	82
8.2 기존 API JSON 스키마의 MCP 자동 변환 기능	82
8.2.1. API-MCP 매핑: 설계 패러다임의 전환	83
8.2.2. 대규모 API 자산의 체계적 전환 프로세스	83
8.2.3. 변경에 대응하는 자동 갱신 구조	84
8.3 MCP·RAG·Observability 통합 아키텍처	84
8.3.1. 통합 인프라 구성: MSA 기반의 AI 서비스망	85
8.3.2. RAG-MCP 통합 구조: 지식과 행동의 결합	85
8.3.3. VibeOps/AIOps 시나리오: MSAP APM & Observability 기반의 지능	
형 자율 운영	86
8.4 프롬프트 중심 업무 UX와 Widget 기반 통합	87
8.4.1. 핵심 기능: Prompts와 Elicitation	87
8.4.2. 컨텍스트 전환 없는 업무 연속성 유지	88
8.4.3. 업무 집중도와 생산성에 미치는 영향	88
8.5 단계별 도입 로드맵 및 실행 체크리스트	89
1단계: 개념 증명 (PoC, Proof of Concept)	89
2단계: 핵심 업무 확장	90
3단계: 플랫폼 내재화	90
제9장: References & Links	91



#### 제1장 서론 - GPU에서 MCP 수 중심의 AI 경쟁력으로

본 백서는 기업 AI 전략의 핵심적인 패러다임 전환을 제안합니다. 지금까지 AI 경쟁력의 척도는 GPU와 같은 고가의 컴퓨팅 인프라 보유량으로 평가되어 왔습니다. 그러나 이제 그 무게 중심은 실제 업무 현장에서 가치를 창출하는 '업무 접점의 수', 즉 MCP(Model Context Protocol) 수로 빠르게 이동하고 있습니다. 이 전환은 단순한 기술적 변화를 넘어, 기업의 AI 투자 효율성과 비즈니스 성과를 좌우하는 가장 중요한 전략적 변수입니다.

본 장은 GPU 숫자 세기에 몰두하는 경쟁에서 벗어나, AI를 실제 업무 자동화와 연결하여 측정 가능한 가치를 창출하는 MCP 기반 전략으로 전환하지 않으면 안 되는 이유를 명확히 제시합니다.

#### 1.1 문제 제기: GPU·모델 중심 사고의 구조적 한계

#### 1.1.1 도입: 현재의 AI 인프라 논의 분석

현재 업계 전반의 AI 인프라 관련 논의는 컴퓨팅 파워의 양적 측면에 과도하게 집중되어 있습니다. 이러한 경향은 AI의 실질적인 비즈니스 가치 평가를 왜곡하고 있으며, 막대한 초기 투자에도 불구하고 기대에 미치지 못하는 ROI(투자수익률)의 근본적인 원인이 되고 있습니다. 따라서 기업의 AI 투자가 실질적인 성과로 이어지기 위해서는 인프라에 대한 관점과 투자 방향성을 시급히 재설정해야 할 필요가 있습니다.

#### 1.1.2 컴퓨팅 파워 집착이 초래한 AI 인프라 논의의 왜곡

- 1. 성능 지표에 매몰된 논의: 현재 AI 인프라 논의는 GPU의 성능 지표에 집중되어 있습니다. 예를 들어, FP8 TensorCore: 최대 1,979TFL0PS와 같은 구체적인 수치는 기술 사양서의 핵심 항목으로 다뤄집니다. 이는 AI 경쟁력에 대한 논의가 TFLOPS와 같은 원초적인 연산 능력 확보에 국한되어 있음을 보여줍니다.
- 2. '생산'과 '활용'의 혼동: 이러한 막대한 컴퓨팅 파워는 주로 거대 언어 모델(LLM)을 처음부터 학습시키는 LLM 모델 학습을 위해 요구됩니다. 이는 AI 가치를 '생산'하는 단계에 해당합니다. 그러나 실제 비즈니스 가치는 이미 검증된 모델을 업무 프로세스에 연결하여 '활용'하는 추론(Inference) 단계에서 발생합니다. 이는 AI 가치 사슬의 마지막 단계인 '활용'을



무시하고, 오직 '생산' 단계에만 자원을 집중시키는 구조적 왜곡을 초래한다. 이는 발전소만 짓고 전력망 투자는 외면하는 것과 같다.

#### 1.1.3 GPU·파라미터 수와 실제 업무 생산성 간의 불일치

GPU 보유 수량이나 AI 모델의 파라미터 수가 실제 업무 생산성 향상으로 직결된다는 보장은 없습니다. 이 관계를 '발전소'와 '전력망'의 비유로 설명할 수 있다. GPU 클러스터는 전력을 생산하는 '발전소(Power Plant)'에 해당한다. 하지만 아무리 강력한 발전소가 있어도, 생산된 전력을 각가정과 공장, 즉 실제 업무 시스템으로 전달하는 '전력망(Power Grid)'이 없다면 아무런 가치를 창출할 수 없습니다. 이 전력망의 역할을 하는 것이 바로 MCP입니다. MCP가 부재한 상황에서 GPU에 대한 투자는 결국 활용되지 못하는 유휴 자원(Idle Asset) 으로 전락할 위험이 큽니다.

#### 1.1.4 "GPU 몇 장인가?"에서 "AI로 실행되는 업무 수는 얼마인가?"로

앞선 논의를 종합해 볼 때, 우리는 AI 경쟁력을 평가하는 질문을 근본적으로 바꾸어야 합니다.

"우리 조직은 GPU를 몇 장 보유하고 있는가?"

라는 기존의 질문에서,

"우리 조직의 AI는 몇 종류의 핵심 업무를 실제로 실행할 수 있는가?"

라는 새로운 질문으로의 전환이 시급합니다.

이 질문의 전환은 AI 투자의 초점을 단순 비용(Capex) 중심에서 실질적인 비즈니스 가치 (Value) 중심으로 옮기는 핵심적인 과정입니다. 그리고 이 새로운 질문에 답하기 위해서는, AI 경 쟁력을 측정할 새로운 지표가 필요합니다.

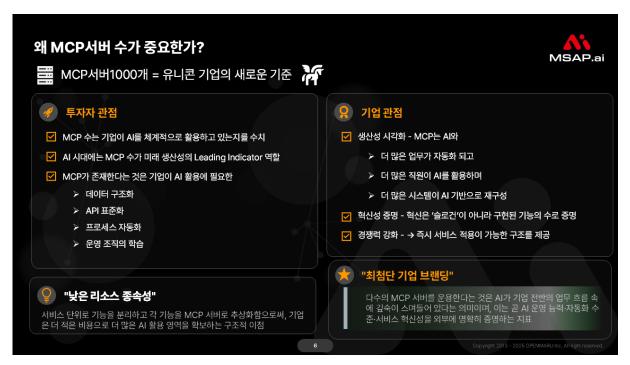
#### 1.2 AI 경쟁력의 새로운 정의: '업무 접점 수(MCP 수)'

#### 1.2.1 도입: 새로운 경쟁력 지표 제시

AI 시대의 진정한 경쟁력은 인프라 보유량이 아닌, AI가 실제 업무 프로세스와 얼마나 많이, 그리고 깊게 연결되어 있는지를 측정하는 지표로 정의되어야 합니다. 우리는 이 새로운 지표를 '업무접점 수', 즉 'MCP 수' 라고 명명합니다. 이 지표는 기업의 AI 성숙도와 실제 비즈니스 영향력을



가장 정확하게 반영합니다. AI가 더 많은 업무 접점을 가질수록, 조직 전체에 지능화된 자동화가 더 넓고 깊게 확산되고 있음을 의미하기 때문입니다.



[그림 1] 왜 MCP 구현 개수가 중요한가?

1.2.2 MCP 수 정의: LLM이 직접 호출할 수 있는 프로덕션 업무 기능 수 'MCP 수'의 명확한 정의는 다음과 같습니다.

LLM이 Model Context Protocol을 통해 직접 호출하여 실행할 수 있는, 실제 운영 환경 (Production)에 배포된 고유한 비즈니스 기능(업무)의 총 수

이 정의는 추상적인 개념이 아니다. 이미 대한민국 유수의 공공 및 금융 기관들이 제안요청서 (RFP)를 통해 요구하는 핵심 업무 기능들이 바로 미래의 MCP 후보 자산들이다. 다음은 실제 프로젝트 요구사항에서 발췌한 구체적인 예시다.

#### • 재무/총무:

- 가상계좌 생성 및 조회
- 지급의뢰서 일괄 출력
- 지출결의서 이체확인증 출력



#### • 인사관리:

- 재직증명서 및 연구실적증명서 일괄 발급
- 채용공고 외부 시스템 연계

#### • 재난대응:

- 소방안전지도 GIS 기능 조회
- 차량동태 관제
- 기상정보 연계

#### • 공공서비스:

- 공연/숙박시설 예약 기능
- 환급금 신청 및 조회

#### 1.2.3 GPU 수 vs MCP 수: 필요 조건과 충분 조건의 구분

GPU와 MCP의 관계는 비즈니스 가치 창출 관점에서 '필요 조건'과 '충분 조건'으로 명확히 구분할 수 있습니다.

구분	지표	역할 및 의미	비유 (발전소/전력망)
필요 조건	GPU 수	AI 모델을 구동하기 위	발전소
		한 최소한의 연산 능력	
		을 확보하는 단계. 이	
		것만으로는 비즈니스	
		가치가 창출되지 않음.	
충분 조건	MCP 수	확보된 연산 능력을 실	전력망
		제 업무 프로세스에 연	
		결하여 자동화와 지능	
		화를 구현하고, 실질적	
		인 생산성 향상과 가치	
		를 창출하는 단계.	



#### 1.2.4 "모델 보유"에서 "업무 접점 수"로 이동하는 성숙도 관점

기업의 AI 성숙도는 두 단계로 명확히 구분된다.

- 초기 단계(Immature Stage) 는
  - '모델 및 인프라 확보'에 집착합니다.
  - 이 단계의 조직은 "어떤 LLM을 도입했는가?" 혹은 "GPU를 몇 대 보유했는가?"를 묻는다.
- 고도화 단계(Mature Stage) 는
  - '업무 접점 중심'으로 전환한다.
  - 이 단계의 성숙한 조직은 "우리의 AI는 몇 개의 핵심 업무를 수행하는가?" 즉, "MCP를 몇 개나 확보했는가?"를 기준으로 성공을 측정합니다.

#### 1.3 본 백서의 목적, 대상 독자, 범위

#### 1.3.1 도입: 백서의 가이드 역할 정의

본 백서는 단순한 기술 소개를 넘어, 독자들이 MCP 기반의 AI 전략을 수립하고 실행하는 데 필요한 명확한 청사진을 제공하는 것을 목표로 합니다. 이를 위해 백서가 추구하는 목적, 소통하고자하는 핵심 독자층, 그리고 다루고자 하는 내용의 경계를 명확히 설정함으로써 실질적인 가이드를 제공하고자 합니다.

#### 1.3.2 MSA·쿠버네티스·클라우드 네이티브에 대한 최소 전제

본 백서는 독자가 마이크로서비스 아키텍처(MSA), 쿠버네티스, 클라우드 네이티브와 같은 현대적인 IT 아키텍처에 대한 기본적인 이해를 갖추고 있음을 전제로 합니다. 다수의 공공 및 금융기관에서 이러한 기술들은 이미 핵심적인 시스템 전환 요건으로 반복적으로 등장하고 있습니다. MCP는 이러한 검증된 기술 스택 위에서 AI의 잠재력을 비즈니스 가치로 변환하는 논리적이고 필연적인 아키텍처 상위 계층입니다.



#### 1.3.3 공공·금융·엔터프라이즈 IT 의사결정자가 이 백서를 읽어야 하는 이유

AI 투자의 패러다임이 변하고 있습니다. 이제 경쟁의 핵심은 막대한 초기 비용이 드는 GPU 인프라 확충이 아니라, 기존 IT 자산을 재활용하고 실제 업무 효율을 높이는 MCP 구축 역량에 있습니다.

본 백서는 이러한 전환의 최전선에 있는 공공, 금융, 엔터프라이즈 IT 의사결정자들을 위한 전략 지침서입니다. MCP 중심의 접근법이 어떻게 AI 프로젝트의 ROI(투자수익률)를 극대화하고, 실질적인 비즈니스 목표를 달성하는 가장 현실적인 경로인지를 명확히 제시합니다. 이러한 프로젝트의 공통점은 특정 LLM 모델 도입 자체가 목표가 아니라, '연구 지식 관리', '항만 운영'이라는 명확한 비즈니스 도메인의 지능화를 목표로 한다는 점이다. 이는 성공적인 AI 도입이 GPU 인프라 규모가 아닌, 얼마나 많은 핵심 업무(MCP)를 AI와 연결했는지에 따라 결정됨을 명백히 보여줍니다.

#### 1.3.4 이 백서에서 다루지 않는 범위(모델 연구·알고리즘 레벨 논의 등)

이 백서는 LLM 모델 자체를 개발하거나, 특정 Al 알고리즘의 성능을 비교 분석하는 학술적 논의는 다루지 않습니다.

대신, 이미 시장에서 검증된 강력한 LLM을 어떻게 기업의 특수한 업무 환경에 효과적으로 '통합하고 활용'할 것인가에 대한 아키텍처와 엔지니어링 전략에 집중합니다. 우리의 목표는 학술적탐구가 아니라, AI를 당신의 대차대조표에 긍정적인 영향을 미치는 실질적인 자산으로 전환하는엔지니어링 청사진을 제공하는 것입니다.

#### 1.4 백서 구성과 활용 방법

#### 1.4.1 도입: 백서의 실용적 가치 강조

본 백서는 단순한 정보 전달을 넘어 독자가 실제 행동 계획을 수립할 수 있도록 구성된 실용적인 가이드입니다. 각 장은 독립적인 주제를 다루면서도 유기적으로 연결되어, 독자가 자신의 조직 상황에 맞는 전략을 도출할 수 있도록 설계되었습니다. 백서의 전체적인 흐름과 각 장에서 얻을 수있는 구체적인 산출물을 이해하면, 본 백서를 더욱 전략적으로 활용할 수 있습니다.



1.4.2 온프레미스 AI 플랫폼 → MCP 기술 → 성숙도 모델 → 기존 시스템 MCP 화 → MSAP.ai 순서

백서 전체는 다음과 같은 논리적 흐름에 따라 구성되어 있습니다.

- 1. 온프레미스 AI 플랫폼 구축: AI 전략의 기반이 되는 인프라 환경을 정의합니다.
- 2. MCP 기술 심층 분석: AI와 업무 시스템을 연결하는 핵심 기술인 MCP를 상세히 다룹니다.
- 3. AI 성숙도 모델: 조직의 현재 위치를 진단하고 목표를 설정하는 프레임워크를 제시합니다.
- 4. 기존 시스템의 MCP화: 레거시 자산을 AI가 활용 가능한 도구로 전환하는 구체적인 패턴을 설명합니다.
- 5. MSAP.ai 솔루션: 위 모든 과정을 가속화하는 통합 플랫폼을 소개합니다.

#### 1.4.3 각 장별 산출물: 평가 지표, 설계 패턴, 참조 아키텍처

독자는 백서의 각 장을 통해 다음과 같은 구체적이고 실용적인 결과물을 얻을 수 있습니다.

주제 영역	핵심 산출물	설명
Al 성숙도 모델	평가 지표	조직의 AI 역량을 정량적으로
		측정하고 개선 영역을 식별할
		수 있는 체크리스트
기존 시스템 MCP화	설계 패턴	레거시 API, DB 쿼리, 배치 작
		업 등을 MCP 서버로 전환하는
		재사용 가능한 아키텍처 솔루
		션
MSAP.ai	참조 아키텍처	실제 엔터프라이즈 환경에 적
		용 가능한 MCP 기반 AI 플랫
		폼의 구체적인 시스템 구성도

## 1.4.4 "나중에"가 아닌 "지금" 무엇을 시작할 것인가에 대한 가이드이 장을 닫기 전에 독자에게는 명확한 실행 계획이 손에 쥐어져야 한다.



MCP 중심의 AI 전략은 거대한 마스터플랜이 아니라, 오늘 당장 시작할 수 있는 작고 빠른 실행의 반복으로 완성된다. 다음 4단계 접근법은 단순한 로드맵이 아니라, 지금 즉시 조직의 AI 투자 방향을 재설정하기 위한 첫 번째 실행 과제다.

- 1. 1단계 자산 파악: 조직 내 시스템들을 기능 단위로 분석하여 "AI가 호출하면 유의미한" 기능 후보 목록을 작성합니다.
- 2. 2단계 최소 MCP 도메인 선정: 리스크가 낮고 성공 효과가 빠르게 나타날 수 있는 업무 도메인(예: 내부 규정 질의, 고객 문의 분류)을 우선적으로 선정합니다.
- 3. 3단계 MCP 기반 PoC 설계: 선정된 도메인의 핵심 기능 3~5개를 MCP 도구로 구현하는 소규모 검증 프로젝트를 설계합니다.
- 4. 4단계 효과 측정 체계 정의: MCP 수, 업무 처리 시간 단축률 등을 핵심 성과 지표(KPI)로 정의하고, 향후 AI 투자의 우선순위를 이 지표에 기반하여 결정하는 체계를 수립합니다.

### 제2장: 엔터프라이즈 AI 인프라와 온프레미스 AI 플랫폼 전략

엔터프라이즈 환경에서 인공지능(AI)의 성공적인 도입은 단순히 강력한 언어 모델(LLM)을 확보하는 것을 넘어, 전체 인프라 스택에 대한 깊이 있는 이해와 전략적 접근을 요구하는 복합적인 과제가 되었습니다. 많은 조직이 AI 모델의 잠재력에 집중하는 동안, 그 모델이 실제 비즈니스 가치를 창출하기 위해 필요한 견고한 기반, 즉 인프라의 중요성을 간과하는 경우가 많습니다. AI 기술이기업의 핵심 데이터 및 업무 시스템과 안전하고 효율적으로 통합되지 않는다면, 아무리 뛰어난 모델도 고립된 기술 시연에 그칠 뿐입니다.

본 장에서는 기업이 진정한 AI 경쟁력을 확보하기 위해 반드시 고려해야 할 인프라 전략의 새로운 설계도를 제시합니다. 먼저, 기존 IT 인프라 모델을 넘어선 새로운 엔터프라이즈 AI 인프라 4계층 모델(연산, 모델, 컨텍스트, 애플리케이션)을 재정의하고 각 계층의 역할을 분석합니다. 이어서 데이터 주권과 내부 자산 활용이 중요한 공공 및 엔터프라이즈 환경에서 왜 온프레미스 (On-premise) AI 플랫폼이 전략적 필수 요건이 되는지 심층적으로 논의합니다.



마지막으로, AI 모델과 기업의 다양한 IT 자산을 연결하는 표준화된 신경망으로서 MCP(Model Context Protocol)가 가지는 핵심적인 역할과 그 전략적 가치를 조명할 것입니다.

#### 2.1 엔터프라이즈 AI 인프라 4계층 재정의

AI 기술을 기업의 핵심 자산과 안전하게 통합하고 비즈니스 가치를 극대화하기 위해서는, 기존의 IT 인프라 모델을 넘어서는 새로운 아키텍처적 접근이 필수적입니다. 본 절에서는 엔터프라이즈 AI 인프라를 연산(Computation), 모델(Model), 컨텍스트(Context), 애플리케이션(Application)의 4개 계층으로 재정의합니다. 이 4계층 모델은 AI 도입을 위한 기술 요소를 논리적으로 구분하여, 각 계층의 역할을 명확히 하고 이들이 어떻게 유기적으로 상호작용하며 최종적인 비즈니스 가치를 창출하는지에 대한 전략적 청사진을 제공합니다.

#### 2.1.1 연산 계층: GPU/TPU, 노드·클러스터 구성

연산 계층은 AI 모델의 학습과 추론을 위한 물리적 기반이자 모든 AI 워크로드의 엔진 역할을 담당합니다. 이 계층의 핵심은 GPU(Graphics Processing Unit)와 같은 고성능 병렬 처리 하드웨어입니다. 예를 들어, 공공기관의 LLM 도입 사업에서는 다음과 같은 엔터프라이즈급 GPU 서버사양이 요구됩니다.

• GPU: NVIDIA HGX-H200 \* 8ea 이상

• 성능: FP8 TensorCore: 최대 1,979TFLOPS

이러한 고성능 연산 자원은 거대 언어 모델의 빠른 연구 개발 주기 확보와 대규모 AI 모델 학습 및 추론 실험을 가능하게 합니다. 그러나 단순히 고사양 하드웨어를 보유하는 것만으로는 충분하지 않습니다. 이 자원들을 효과적으로 활용하기 위해서는 여러 GPU 서버를 하나의 논리적 단위로 묶는 노드(Node) 및 클러스터(Cluster) 구성이 필수적입니다. 잘 설계된 클러스터는 자원의독점적 사용을 보장하고, 보안 및 관리 부담을 최소화하며, 안정적인 AI 서비스를 제공하는 기반이됩니다. 결론적으로, 연산 계층의 성능은 모델 계층의 잠재력을 결정하고, 컨텍스트 계층의 응답속도에 직접적인 영향을 미치는 AI 인프라의 물리적 심장입니다.



#### 2.1.2 모델 계층: 파운데이션 모델·도메인 LLM·서빙 스택

모델 계층은 AI의 '두뇌' 역할을 하는 소프트웨어 스택으로, 실제 지능을 구현하는 핵심 요소입니다. 이 계층은 크게 두 부분으로 구성됩니다.

첫째는 방대한 데이터로 사전 학습된 파운데이션 모델(Foundation Model)이며,

둘째는 이를 특정 산업이나 업무 도메인(예: 법률, 금융, 의료) 데이터로 미세 조정하여 전문성을 강화한 도메인 특화 LLM입니다.

이렇게 준비된 모델을 실제 서비스로 제공하기 위해서는 안정적인 서빙 스택(Serving Stack)이 필요합니다.

서빙 스택은 모델을 메모리에 로드하고, 사용자 요청에 따라 추론을 수행하며, 그 결과를 애플리케이션에 전달하는 일련의 과정을 관리합니다. "LLM 기반 지능형 정보시스템 구축"과 같은 실제 사업은 이 모델 계층이 어떻게 특정 비즈니스 문제 해결에 직접적으로 관여하는 지를 명확히 보여주는 사례입니다. 결국 모델 계층은 연산 계층의 자원을 소모하여 지능을 생성하고, 컨텍스트 계층을 통해 조직의 실제 문제와 연결되어 애플리케이션 계층에서 가치를 발휘하는 Al의 두뇌라 할수 있습니다.

#### 2.1.3 컨텍스트 계층: MCP·RAG·데이터 커넥터 계층

컨텍스트 계층은 모델 계층의 AI가 기업 내부의 자산(데이터, 시스템, API)과 소통하고 상호작용할 수 있도록 연결하는 '신경망'에 비유할 수 있습니다. AI 모델이 아무리 뛰어나더라도 조직의 맥락(Context)을 이해하지 못하면 실질적인 가치를 창출하기 어렵습니다. 이 계층의 핵심 기술은 MCP와 RAG입니다.

- MCP (Model Context Protocol): MCP는 'LLM과 외부 시스템을 연결하는 개방형 표준 프로토콜'입니다. 기업 내 수많은 AI 에이전트(M)와 다양한 업무용 도구(N)를 개별적으로 연동할 때 발생하는 'N x M 통합 문제'를 해결하는 핵심적인 역할을 합니다. MCP는 모든 시스템과 모델이 단일한 표준 인터페이스를 통해 통신하게 함으로써, 복잡한 개별 연동 작업을 제거하고 확장성 있는 AI 생태계를 구축할 수 있도록 지원합니다.
- RAG (Retrieval-Augmented Generation): RAG는 LLM이 답변을 생성할 때, 고정된 내부 지식에만 의존하는 것이 아니라 기업의 최신 내부 데이터를 실시간으로 검색하고 참조



하게 하는 기술입니다. 이를 위해 내부 문서들은 의미 단위로 분할(Chunking)되고 벡터 임 베딩으로 변환되어 벡터 DB에 저장되며, 사용자 질의와 가장 관련성 높은 문서 조각을 실시 간으로 찾아내 LLM의 컨텍스트에 주입하는 방식으로 동작합니다. 예를 들어, 챗봇 상담 시스템 구축 시 "다양한 유형의 문서 파일(hwp, pdf, txt, doc 등) 수집 및 처리"를 통해 RAG를 구현하면, 최신 주거 정책이나 내부 규정이 반영된 정확한 답변을 사용자에게 제공할 수 있습니다.

이처럼 컨텍스트 계층은 연산 계층의 속도와 모델 계층의 지능을 애플리케이션 계층의 구체적인 업무와 연결시키는, 가치 창출의 핵심 통로입니다.

#### 2.1.4 애플리케이션·업무 계층: 챗봇, 에이전트, 업무 시스템

애플리케이션·업무 계층은 앞선 세 개의 계층을 기반으로 AI 기술이 최종적으로 비즈니스 가치를 창출하는 접점입니다. 이 계층에서는 AI가 기존 업무 시스템에 내장(Embedding)되거나, 새로운 지능형 애플리케이션(챗봇, AI 에이전트 등)의 형태로 최종 사용자에게 제공됩니다. 다수의 공공 기관 제안요청서(RFP)에 명시된 실제 업무 시스템들은 AI가 어떻게 기존 업무를 지능화하는지 구체적으로 보여줍니다.

- 인사관리: 참여 연구진의 재직증명서 및 연구실적증명서를 일괄 발급하고, 이메일·카카오톡 인증 등 채용관리 시스템을 개선하여 프로세스를 효율화합니다.
- 총무관리: 전자증빙시스템을 활용해 지급의뢰서를 일괄 출력하고, 결의서 반려 프로세스를 자동화하여 행정 업무 부담을 경감합니다.
- 재난대응: 소방안전지도 시스템이 재난 상황, 대응 전략 등 다양한 정보를 빅데이터로 관리하고 데이터 기반의 효율적인 현장 대응을 강화합니다.
- 민원서비스: LLM 기반 대화형 챗봇을 통해 주거정책 및 제도 관련 상담을 24시간 자동화하고, 사용자의 자연어 질의에 실시간으로 응답합니다.

결론적으로, 애플리케이션 계층은 연산, 모델, 컨텍스트 계층의 모든 기술적 투자가 집약되어 최종 비즈니스 가치로 발현되는 최상위 인터페이스입니다. 이 4계층 모델은 AI 도입을 위한 기술적 구성요소를 체계적으로 이해하는 틀을 제공합니다. 특히 기업의 핵심 자산인 데이터와 업무 시스템과의 안전하고 효율적인 연동을 위해서는, 외부 클라우드 서비스에 의존하기보다 내부 통제권



을 확보할 수 있는 온프레미스 AI 플랫폼이 왜 중요한 전략적 선택지가 되는지를 명확히 보여줍니다.

#### 2.2 왜 공공·엔터프라이즈에는 온프레미스 AI 플랫폼이 필요한가

퍼블릭 클라우드 AI 서비스가 제공하는 신속성과 편리성에도 불구하고, 데이터 주권, 보안 규제, 그리고 내부 자산의 전략적 활용이 최우선 과제인 공공 및 엔터프라이즈 환경에서는 온프레미스 AI 플랫폼 구축이 단순한 선택이 아닌 전략적 필수 요건으로 부상하고 있습니다. 이는 AI를 외부의 지능을 빌려오는 도구로 보는 시각에서 벗어나, 조직의 핵심 자산을 지능화하는 내재된 플랫폼으로 인식하는 관점의 전환을 요구합니다.

2.2.1 AI를 "내부 자산(업무 시스템·데이터)을 지능화하는 플랫폼"으로 보는 관점

엔터프라이즈 AI의 본질은 외부의 범용 모델을 사용하는 것이 아니라, 조직이 이미 보유한 가장 가치 있는 자산인 '업무 시스템과 데이터'의 잠재력을 극대화하는 데 있습니다. 수많은 공공기관이 기존 핵심 시스템의 고도화를 목표로 하는 것은 이러한 관점을 명확히 보여줍니다. AI는 이들 시스템에 내장되어 프로세스를 자동화하고, 데이터 기반의 의사결정을 지원하며, 서비스의 질을 향상시키는 역할을 합니다. 이처럼 AI를 내부 자산을 지능화하는 플랫폼으로 바라볼 때, AI 인프라는 자연스럽게 내부 시스템과 가장 가까운 곳, 즉 온프레미스 환경에 위치해야 한다는 결론에 도달하게 됩니다.

2.2.2 내부 시스템·데이터 근접성, 지연·가용성·통제 측면의 이점 온프레미스 플랫폼은 기술적으로도 명확한 이점을 제공합니다.

• 근접성 및 지연(Latency): AI 추론이 내부 데이터 및 시스템과 물리적으로 가까운 위치에서 수행될 때, 네트워크 지연 시간이 최소화됩니다. 이는 실시간 응답이 중요한 대국민 서비스 나 신속한 의사결정이 필요한 내부 업무 시스템에서 결정적인 경쟁 우위가 됩니다. 예를 들어, '서울소방재난본부'의 소방안전지도 시스템과 같이 재난 현장의 실시간 정보 공유가 인



명 구조와 직결되는 미션 크리티컬한 업무에서, 외부 클라우드 왕복으로 인한 수백 밀리초 (ms)의 지연은 의사결정의 골든타임을 놓치게 할 수 있습니다.

• 가용성 및 통제(Availability & Control): 외부 클라우드 서비스의 장애나 정책 변경에 종속되지 않고 독자적인 서비스 연속성을 확보할 수 있습니다. 시스템 장애 발생 시 자체적인 복구 및 통제가 가능하다는 점은 공공 서비스의 안정성을 위한 핵심 요건입니다. 온프레미스 플랫폼은 조직에 완전한 통제권을 부여하여 시스템의 안정성과 가용성을 최고 수준으로 유지할 수 있게 합니다.

# 2.2.3 조직 안의 IT 자산을 어떻게 엮느냐가 GPU·LLM 보유보다 중요한 이유 AI 시대의 진정한 경쟁력은 단순히 고가의 GPU를 얼마나 많이 보유했는지, 혹은 최신 LLM을 도입했는지에 따라 결정되지 않습니다. 핵심은 조직 내부에 흩어져 있는 IT 자산들을 얼마나 효과적으로 엮어내어 AI가 활용할 수 있도록 만드느냐에 있습니다. 다음과 같은 비유는 이 점을 명확하게 설명합니다.

- GPU/LLM 클러스터는 전기를 생산하는 '발전소' 입니다.
- Al 기능(MCP, RAG 등)은 발전소에서 생산된 전력을 공장과 가정으로 전달하는 '전력망' 입니다.
- 조직의 업무 시스템과 담당자는 이 전력을 사용하는 '소비자' 입니다.

아무리 강력한 발전소를 보유하고 있더라도, 생산된 전력을 각 가정과 공장(소비자)까지 안전하고 효율적으로 전달하는 '전력망'이 없다면 발전소는 무용지물입니다. 엔터프라이즈 AI 환경에서 이 '전력망'의 역할을 하는 것이 바로 MCP와 같은 표준화된 연결 체계입니다. 조직의 IT 자산을 MCP를 통해 엮어내는 작업이야말로, GPU 투자 가치를 실현하고 AI를 조직 전체에 확산시키는 가장 중요한 전제 조건입니다.

결국 온프레미스 플랫폼의 필요성은 데이터와 시스템에 대한 완전한 통제권 확보와 직결됩니다. 이는 단순히 기술적 우위를 넘어, 기업의 가장 민감한 자산을 외부 클라우드로 전송할 때 발생하는 심각한 리스크를 근본적으로 차단하기 위한 전략적 결정입니다.



#### 2.3 퍼블릭 클라우드 LLM에 내부 데이터를 보내는 리스크

기업의 가장 민감하고 전략적인 자산인 내부 데이터를 외부 퍼블릭 클라우드 LLM으로 전송하는 행위는 단순한 기술적 API 호출을 넘어, 예측하기 어려운 규제 및 보안 리스크를 조직에 전가할 수 있습니다. 편리성과 비용 효율성이라는 장점 이면에 숨겨진 이러한 리스크를 정확히 인지하지 못한다면, 단 한 번의 잘못된 데이터 전송이 조직 전체를 위험에 빠뜨릴 수 있습니다.

#### 2.3.1 데이터 주권·레지던시·규제(공공·금융·개인정보) 관점의 제약

공공기관 및 금융기관의 제안요청서에서 반복적으로 강조되는 것은 관련 법규 및 규정 준수 의무입니다. 특히 "개인정보보호법", "국가정보보안기본지침" 등은 민감한 데이터의 외부 전송을 엄격히 통제합니다. 공공, 금융, 개인정보와 같이 고도로 민감한 데이터를 다루는 기관에게 데이터의물리적 위치를 의미하는 데이터 레지던시(Data Residency)와 데이터에 대한 통제권을 의미하는데이터 주권(Data Sovereignty)은 타협할 수 없는 원칙입니다. 내부 데이터를 국외에 위치한 퍼블릭 클라우드 데이터센터로 보내는 것은 이러한 규정을 위반할 소지가 매우 크며, 이는 법적 제재로 이어질 수 있습니다.

#### 2.3.2 벡터DB·로그·모델 피드백에 남는 민감 정보 이슈

데이터가 한번 조직의 네트워크 경계를 벗어나면, 그 흔적은 다양한 형태로 외부 시스템에 남게 됩니다.

- 사용자 질의: "특정 고객의 민원 처리 내역 요약"과 같은 질의 자체에 민감 정보가 포함될 수 있습니다.
- RAG 프로세스: 내부 문서를 기반으로 답변을 생성하는 과정에서 문서의 일부 또는 전체가 외부 벡터 DB에 임베딩되어 저장될 수 있습니다. 외부 벡터 DB에 저장되는 임베딩 자체는 원문 텍스트의 의미적 정보를 압축하여 포함하므로, 민감한 원본 데이터가 삭제된 후에도 데이터의 본질적인 내용이 유추되거나 복원될 잠재적 위험을 내포합니다.
- 시스템 로그 및 피드백: API 호출 기록, 오류 로그, 모델 성능 개선을 위한 피드백 데이터 등 에도 민감 정보가 포함되어 외부 서비스 제공자의 서버에 장기간 보관될 수 있습니다.



이러한 정보의 외부 잔존은 "보안서약서" 및 "누출금지 대상 정보" 조항에 명시된 기밀유지 의무를 위반하는 행위가 될 수 있으며, 계약 위반에 따른 심각한 법적 책임을 초래할 수 있습니다.

#### 2.3.3 단일 질의가 조직 전체 규제 위반으로 이어지는 전형적인 패턴들

조직의 보안 정책이 아무리 견고하더라도, 직원의 단순한 질의 하나가 조직 전체를 규제 위반의 위험에 빠뜨릴 수 있습니다. 예를 들어, 한 직원이 업무 편의를 위해 "특정 민원인의 개인정보가 포함된 내부 문서를 요약해줘"라는 질의를 외부 퍼블릭 LLM 기반 챗봇에 입력했다고 가정해 보겠습니다. 이 순간, 해당 민원인의 개인정보는 조직의 통제 범위를 벗어나 외부로 전송되며, 이는 명백한 개인정보보호법 위반에 해당합니다. 이러한 보안 위규 사항은 "[별첨 2] 보안 위약금 부과 기준"에 따라 총 사업비의 일정 비율에 해당하는 위약금 부과나 부정당업자 등록과 같은 심각한 제재로 이어질 수 있습니다.

이러한 리스크를 근본적으로 회피하기 위해서는 데이터를 외부로 보내지 않는, 즉 온프레미스 및 하이브리드 환경에 최적화된 기술적, 아키텍처적 대응 방안이 반드시 필요합니다.

#### 2.4 온프레미스·하이브리드 AI 플랫폼 설계 원칙

앞서 논의된 퍼블릭 클라우드 LLM 활용의 리스크를 최소화하고, 엔터프라이즈 내부 자산의 활용을 극대화하기 위해서는 명확한 설계 원칙에 기반한 온프레미스 및 하이브리드 AI 플랫폼 구축이 필요합니다. 이는 단순히 인프라를 내부에 두는 것을 넘어, 보안, 통합, 관찰 가능성을 모두 고려한 체계적인 아키텍처 설계를 의미합니다.

2.4.1 온프레미스 LLM 서빙 + MCP + RAG + Observability 통합 구조 이상적인 온프레미스 AI 플랫폼은 다음과 같은 핵심 구성 요소를 통합한 아키텍처를 가집니다.

- 온프레미스 LLM 서빙: 내부 GPU 클러스터를 활용하여 자체적으로 LLM을 운영합니다. 이를 통해 모델 추론 과정에서 데이터가 외부로 유출되는 것을 원천적으로 차단하고, 모델에 대한 완전한 통제권을 확보합니다.
- MCP (Model Context Protocol): 모든 내부 업무 시스템과 AI 모델을 연결하는 표준 인 터페이스 역할을 합니다. MCP를 통해 각 시스템은 표준화된 '도구(Tool)'로 추상화되어,



AI가 일관된 방식으로 업무를 호출하고 실행할 수 있게 됩니다.

- RAG (Retrieval-Augmented Generation): 외부 데이터 유출 없이 최신 내부 지식 자산을 AI가 활용할 수 있게 하는 핵심 메커니즘입니다. 모든 데이터 수집, 정제, 임베딩 과정이 내부에서 이루어지므로 데이터 주권을 완벽하게 지킬 수 있습니다.
- Observability (관찰 가능성): 시스템의 상태, 성능, 비용, 보안 감사를 위한 필수 요소입니다. 여러 제안요청서에서 요구하는 "모니터링 강화" 및 "중앙집중식 로깅" 요건을 충족해야합니다. Open Telemetry 기반의 중앙 집중식 로그 수집 및 분석 체계와, Observability를 활용한 시스템 메트릭 모니터링 체계 구축을 의미합니다. 이를 통해 장애 발생 시 신속한원인 분석과 선제적인 성능 관리가 가능해집니다.

#### 2.4.2 퍼블릭 LLM 활용 범위: 마스킹·프록시·샌드박스 전략

온프레미스를 중심으로 구축하되, 특정 비즈니스 요구에 따라 퍼블릭 LLM의 장점을 안전하게 활용하기 위한 하이브리드 전략도 고려할 수 있습니다. 이 경우, 강력한 보안 통제가 전제되어야 합니다. "용역업체 사용 전산망과 기관 전산망의 분리"나 "내부 정보시스템 접근 통제"와 같은 보안요구사항 원칙을 적용하여 다음과 같은 전략을 구현할 수 있습니다.

- 데이터 마스킹(Data Masking): 외부 LLM에 질의를 보내기 전, 개인정보나 기업 기밀과 같은 민감 정보를 식별 불가능한 형태로 변환(마스킹)합니다.
- 보안 프록시(Secure Proxy): 모든 외부 API 호출이 중앙의 보안 프록시를 통과하도록 강제합니다. 이 프록시는 호출 내용 로깅, 데이터 마스킹 적용, 비인가된 API 호출 차단 등의 보안 정책을 일괄적으로 수행합니다.
- 샌드박스(Sandbox): 민감 데이터와 무관한 일반적인 정보 검색이나 콘텐츠 생성과 같이 리스크가 낮은 업무에 한해서만 제한된 환경(샌드박스)에서 외부 모델을 활용하도록 허용합니다.

#### 2.4.3 온프레미스·프라이빗 클라우드·하이브리드 도입 패턴 비교

조직의 보안 요구 수준, IT 역량, 예산, 확장성 요구사항에 따라 최적의 도입 모델은 달라질 수 있습니다. 각 패턴의 특징을 비교하면 다음과 같습니다.



도입 패턴	핵심 특징	장점	단점	적합한 조직
완전 온프레미스	조직이 소유한 데	- 데이터 주권 및		
	이터센터 내에 모	보안 통제 수준		
	든 하드웨어	극대화		
	(GPU 서버 등)와			
	소프트웨어를 직			
	접 구축 및 운영			
- 내부 시스템과	- 높은 초기 투자			
초저지연 연동 가	비용			
<b>⊣</b> 0				

- 인프라 운영 및 관리 복잡성
- 확장성 확보에 제약 | 데이터 주권과 물리적 통제가 최우선인 공공, 국방, 금융 기관 및 '국 가정보보안기본지침'의 엄격한 적용을 받는 조직 | | 프라이빗 클라우드 | 내부 또는 지정된 데이터센터에 클라우드 기술(가상화, 컨테이너 등)을 적용하여 구축한 전용 클라우드 | - 온 프레미스 수준의 보안 및 통제
- 클라우드 기술을 통한 자원 효율성 및 민첩성 확보 | 클라우드 플랫폼 구축 및 운영을 위한 높은 기술 역량 필요
- 퍼블릭 클라우드 대비 신기술 도입 속도가 느릴 수 있음 | IT 거버넌스가 확립되어 있고, 클라우드 네이티브 기술 내재화를 목표로 하는 대기업 | | 하이브리드 AI | 온프레미스 플랫폼을 중심으로, 민감하지 않은 워크로드는 퍼블릭 클라우드 서비스를 함께 활용 | 워크로드 특성에 따른 최적의 인프라 선택
- 비용 효율성과 보안 통제의 균형
- 퍼블릭의 최신 AI 기술을 안전하게 활용 가능 | 온프레미스와 클라우드 간 복잡한 네트워크 및 데이터 관리 필요
- 일관된 보안 정책 적용 및 거버넌스 수립이 어려움 | 신속한 서비스 개발과 보안 통제의 균형이 필요한 디지털 전환 선도 기관 |

이러한 설계 원칙과 도입 패턴을 이해했다면, 다음 질문은 '언제 이 여정을 시작해야 하는가'입



니다. 많은 조직이 기술의 안정화나 비용 하락을 기다리지만, AI 시대에는 그 기다림 자체가 가장 큰 리스크가 될 수 있습니다.

#### 2.5 "나중에"가 아닌 "지금" 시작해야 하는 이유

많은 의사결정자들이 GPU 가격 하락, LLM 기술의 성숙, 혹은 관련 규제의 명확화를 기다리며 Al 플랫폼 구축을 '나중'의 과제로 미루는 경향이 있습니다. 그러나 급변하는 Al 시대에 이러한 '관 망' 전략은 혁신의 기회를 놓치는 것을 넘어, 따라잡기 힘든 기술 격차와 더 큰 기회비용을 초래하는 위험한 선택일 수 있습니다.

#### 2.5.1 GPU 가격·기술·규제를 기다리는 전략이 위험한 이유

기다리는 전략은 세 가지 측면에서 조직을 위험에 빠뜨립니다.

- 기회비용의 발생: 선도적인 공공기관과 기업들은 이미 '24년, '25년 사업으로 LLM, RAG, 클라우드 네이티브 전환 프로젝트를 발주하며 발 빠르게 움직이고 있습니다. 경쟁자들이 AI를 통해 업무 효율성을 높이고 새로운 서비스를 창출하는 동안 관망하는 것은 경쟁 우위를 스스로 포기하는 것과 같습니다.
- 내부 역량 축적의 실패: AI 플랫폼을 성공적으로 구축하고 운영하는 경험과 노하우는 단기 간에 축적되지 않습니다. 지금 작게라도 시작하여 기술적 시행착오를 겪고 운영 역량을 내 재화하지 않으면, 향후 기술 격차는 더욱 벌어져 후발주자로서 시장을 따라잡는 것이 거의 불가능해질 수 있습니다.
- MCP 생태계 선점 기회 상실: MCP 서버를 구축하는 것은 단순히 기술을 도입하는 행위가 아닙니다. 이는 조직의 핵심 업무 프로세스와 데이터를 AI가 이해하고 활용할 수 있는 표준화된 '디지털 자산'으로 만드는 과정입니다. 이 디지털 자산을 먼저, 그리고 더 많이 구축하는 조직이 미래 AI 기반 비즈니스 생태계에서 주도권을 선점하게 될 것입니다.



2.5.2 지금 당장 할 수 있는 일: 자산 목록화, 후보 업무 선정, PoC 범위 정의 거창한 마스터 플랜이 없더라도 즉시 시작할 수 있는 구체적이고 실용적인 첫걸음들이 있습니다. 실제 사업 수행 절차를 참고하여 다음과 같은 3단계 접근을 제안합니다.

- 1. 자산 목록화: 조직이 현재 보유한 내부 시스템, 데이터베이스, API를 기능 단위로 목록화합니다. 이 중에서 "AI와 연동했을 때 비즈니스 가치가 있을 것"으로 판단되는 자산을 식별하여 잠정적인 MCP 후보 리스트를 작성합니다. (RFP의 현황분석 단계에 해당)
- 2. 후보 업무 선정: 전사적인 영향이 크면서도 초기 도입 리스크가 비교적 낮은 업무를 파일럿 대상으로 선정합니다. 예를 들어, 내부 규정 질의응답, 회의록 요약, 보고서 초안 생성과 같은 업무는 좋은 시작점이 될 수 있습니다. (RFP의 **ISP** 또는 **상세설계** 단계에 해당)
- 3. PoC 범위 정의: 선정된 업무에 대해 3~5개의 핵심 기능을 MCP 도구(Tool)로 구현하고, 관련 내부 문서를 RAG 기술과 결합하는 소규모 개념증명(PoC)을 설계하고 실행합니다. 이를 통해 최소한의 자원으로 AI 도입의 실질적인 효과와 기술적 과제를 검증할 수 있습니다.

#### 2.5.3 "먼저 MCP 인벤토리를 가진 조직이 AI 성숙도를 선점한다"는 관점

결론적으로, 미래 기업의 AI 경쟁력을 측정하는 새로운 핵심 지표는 보유한 GPU의 수가 아니라, 실제 업무 프로세스와 유기적으로 연결된 'MCP 인벤토리의 수'가 될 것입니다. MCP 인벤토리가 풍부하다는 것은 AI가 조직의 더 많은 업무 영역에 깊숙이 관여하고 있다는 의미이며, 이는 곧 조직의 AI 성숙도가 높고 비즈니스 민첩성과 혁신 잠재력이 크다는 것을 방증합니다. 따라서 진정한 AI 전략가의 첫 질문은 'GPU를 몇 장 구매할 것인가?'가 아니라, '어떤 핵심 업무부터 MCP 인벤토리로 전환하여 디지털 자산화할 것인가?'가 되어야 한다.



# 제3장 전략적 선택: 왜 '모델 학습'이 아닌 '추론 활용'에 집 중해야 하는가

기업의 인공지능(AI) 경쟁력을 가늠하는 척도가 근본적으로 변하고 있습니다. 과거에는 거대 언어모델(LLM)을 직접 '생산'할 수 있는 역량이 핵심으로 여겨졌지만, 이제는 검증된 상용 모델의 강력한 기능을 기존 업무 시스템에 얼마나 효과적으로 '유통'하고 '활용'하는지가 성패를 가르는 시대가 되었습니다. 이 패러다임의 전환을 이해하기 위해, 우리는 AI 인프라를 전력 시스템에 비유할 수 있습니다. 고성능 GPU 클러스터는 전기를 생산하는 '발전소'에, 이를 통해 생성되는 MCP나 RAG는 '전력망'에, 그리고 실제 가치를 창출하는 각 부서의 업무 시스템은 최종 '소비자'에 해당합니다.

IT 의사결정자는 "우리도 자체 LLM을 만들어야 하는가?"라는 질문에서 벗어나, "어떻게 하면 가장 효율적인 AI 전력망을 구축하여 조직 전체에 AI라는 전력을 공급할 것인가?"라는 전략적 질문에 집중해야 합니다.

본 장에서는 이러한 관점을 바탕으로 AI 투자 전략의 핵심을 논증하고자 합니다. 먼저 모델 '학습(Training)'과 '추론(Inference)'의 기술적·경제적 차이를 명확히 분석하여, 왜 대부분의 기업에 독자적인 모델 학습이 비현실적인 선택인지를 설명합니다. 이어서 GPU 증설만으로는 AI 도입의효과를 거두기 어려운 구조적 한계를 지적하고, 그 대안으로 추론 중심 전략의 필요성을 역설할 것입니다. 최종적으로, 기업의 실질적인 AI 성숙도를 측정하는 핵심 지표가 GPU 보유량이 아닌, 실제 업무와 연결된 'MCP의 수'가 되어야 하는 이유를 구체적인 사례와 함께 증명하겠습니다.

#### 3.1 모델 학습(Training)과 추론(Inference)의 기술적 차이

성공적인 AI 전략을 수립하기 위한 첫걸음은 기술 스택의 가장 기본적인 두 축인 '학습'과 '추론'의 차이를 명확히 이해하는 것입니다. 이 두 개념은 목표, 프로세스, 그리고 요구되는 자원과 인프라 구조 측면에서 근본적으로 다릅니다. 모델을 '만드는' 과정인 학습과, 만들어진 모델을 '사용하는' 과정인 추론을 구분하는 것은 한정된 예산과 인력을 어디에 집중해야 할지 결정하는 합리적인 AI 투자 전략의 명확한 출발점이 됩니다.



#### 3.1.1 Pre-training, Fine-tuning, Instruction Tuning 개념 정리

'모델 학습'은 특정 목적을 수행할 수 있는 AI 모델을 생성하는 모든 활동을 포괄하며, 주로 다음과 같은 세부 과정으로 구성됩니다. IT 의사결정자는 각 과정의 목적과 자원 소모 수준을 이해해야 합 니다.

- 사전 학습 (Pre-training): 대규모의 정제되지 않은 데이터셋을 기반으로 언어의 일반적인 패턴, 문법, 사실적 지식 등을 모델에 내재화하는 과정입니다. 이는 파운데이션 모델의 근간을 이루는 단계로, 천문학적인 규모의 연산 자원과 데이터를 필요로 합니다.
- 미세 조정 (Fine-tuning): 사전 학습된 모델을 특정 도메인이나 과제에 특화시키기 위해, 상대적으로 작지만 잘 정제된 데이터셋으로 추가 학습을 진행하는 과정입니다.
- 명령어 튜닝 (Instruction Tuning): 모델이 사용자의 다양한 지시(Instruction)를 더 잘 이해하고 따르도록, '질문-답변' 쌍으로 구성된 데이터셋을 활용해 학습하는 단계입니다.

이러한 과정들은 모두 막대한 연산 자원을 요구하는 R&D 중심의 활동입니다. 이는 모델 학습이 상용 서비스 운영보다는, 특정 가설을 검증하고 연구개발 사이클을 단축하기 위한 실험적 성격이 강한 활동임을 보여줍니다.

#### 3.1.2 학습 파이프라인과 추론 파이프라인의 비교

모델 학습과 추론은 파이프라인 구성에서도 뚜렷한 차이를 보입니다. 학습 파이프라인이 모델을 '제조'하는 공정에 가깝다면, 추론 파이프라인은 완성된 제품을 소비자에게 '배송'하고 안정적으로 '운영'하는 물류 및 서비스망에 비유할 수 있습니다.

	학습 파이프라인 (Training	추론 파이프라인 (Inference
구분	Pipeline)	Pipeline)
핵심 목표	모델 생성 및 성능 최적화	AI 기능의 안정적인 제공 및 확
		장
주요 활동	데이터 수집·준비, 분산 학습,	모델 서빙, API 엔드포인트 노
	모델 검증, 하이퍼파라미터 튜	출, 요청 처리, 실시간 모니터
	JO	링, 오토스케일링



	학습 파이프라인 (Training	추론 파이프라인 (Inference
구분	Pipeline)	Pipeline)
관련 기술	대규모 데이터 처리 프레임워	MSA, API 게이트웨이, CI/
	크, 분산 컴퓨팅 라이브러리	CD, 서비스 메시, 쿠버네티스
		기반 오토스케일링
성격	R&D, 제조(Build)	운영(Operation), 서비스 제공
		(Serve)

주목할 점은 추론 파이프라인을 구성하는 핵심 기술들은 이미 만들어진 AI 모델의 기능을 민원 서비스나 재난 대응과 같은 실제 업무에 안정적으로 제공하고, 사용자 요청에 따라 유연하게 확장하는 운영 중심의 기술 요소들입니다. 이는 많은 공공 및 민간 조직의 당면 과제가 새로운 모델을 '만드는 것'이 아니라, 기존 시스템과 AI를 안정적으로 '연결하고 운영하는 것'에 있음을 시사합니다.

#### 3.1.3 배치 학습과 온라인 추론의 자원 사용 패턴

자원 사용 패턴의 차이는 두 활동의 경제성을 결정하는 중요한 요소입니다.

- 배치 학습 (Batch Training): 고성능 GPU 자원을 수 주에서 수 개월에 걸쳐 100%에 가깝게 사용하는 패턴을 보입니다. 이는 마치 최대 출력으로 장기간 가동되는 '발전소'와 유사하며, 막대한 초기 투자 비용과 지속적인 운영 비용을 발생시킵니다.
- 온라인 추론 (Online Inference): 반면, '대표홈페이지 민원서비스'와 같은 실제 업무 환경에서의 추론은 사용자의 요청이 발생할 때마다 필요한 만큼만 GPU 자원을 순간적으로 사용합니다. 이는 일상적인 '전력 소비' 패턴과 같습니다. 이러한 가변적인 자원 사용 패턴은 『Model Context Protocol (MCP): An Al for FinOps Use Case』에서 강조하는 비용 최적화 및 거버넌스와 직접적으로 연결됩니다. 필요한 만큼만 자원을 할당하고 비용을 통제하는 것이 핵심이며, 이는 추론 중심 아키텍처의 중요한 경제적 이점입니다.

이처럼 모델 학습과 추론은 기술적 구성부터 자원 활용 패턴에 이르기까지 명확한 차이를 보입니다. 이러한 기술적 차이는 자연스럽게 각 활동의 경제적 타당성 분석으로 이어지며, 왜 대다수 기업이 추론 중심 전략을 선택해야 하는지에 대한 근거를 제공합니다.



#### 3.2 파운데이션 모델 학습의 경제학

자체 파운데이션 모델을 학습하겠다는 결정은 기술적 도전을 넘어 막대한 자본 투자를 요구하는 경제적 선택입니다. 따라서 대부분의 기업에게 이는 경제적으로 비현실적인 경로일 가능성이 높습니다. 이 섹션에서는 글로벌 빅테크 기업과 일반 기업 및 공공기관 간의 근본적인 자본 구조 차이를 분석하고, 인프라, 인력, 데이터 운영을 포함한 총소유비용(TCO)의 구체적인 구성 요소를 살펴봅니다. 이를 통해 왜 많은 중견기업과 공공기관이 독자 LLM 개발에서 구조적으로 실패할 수밖에 없는지를 심층적으로 논증하고자 합니다.

#### 3.2.1 글로벌 빅테크의 LLM 학습 비용과 자본 구조

업계에 알려진 바와 같이, 파운데이션 모델 하나를 학습하는 데에는 수천억 원에 달하는 막대한 비용이 소요됩니다. 이러한 투자는 단순히 비용의 문제를 넘어, 실패를 용인하고 장기적인 R&D를 지속할 수 있는 소수의 글로벌 기업만이 감당할 수 있는 자본 구조를 전제로 합니다. 이들 기업은 시 모델 개발을 통해 파생되는 클라우드 서비스, 광고, 소프트웨어 생태계 등 다각화된 사업 모델로 투자 비용을 회수할 수 있는 구조를 갖추고 있습니다. 반면, 대부분의 공공기관이나 중견기업에게 이러한 규모의 투자는 핵심 사업의 존속을 위협할 수 있는 비현실적인 전략입니다.

#### 3.2.2 데이터·인력·인프라를 포함한 총소유비용(TCO) 요소

AI 모델 개발의 총소유비용(TCO)은 단순히 GPU 서버 구매 비용으로 끝나지 않습니다. 성공적인 모델 운영을 위해서는 다음과 같은 지속적인 투자가 필수적입니다.

- 인프라 비용: 초기 투자 비용이 가장 가시적인 부분입니다. NVIDIA HGX-H200 8개를 탑 재한 서버와 같은 고사양 장비 도입은 수억 원에 달하는 막대한 초기 비용을 요구합니다.
- 전문 인력 비용: AI 모델을 개발하고 유지보수하기 위해서는 높은 수준의 전문 인력이 필요하며, 이는 지속적인 비용을 발생시킵니다. '응용SW 개발자 평균 임금'데이터를 참고하면, 고급 AI 전문가 확보 및 유지에 상당한 인건비가 소요됨을 알 수 있습니다.



• 데이터 및 운영 비용: TCO에서 간과하기 쉽지만 가장 지속적으로 발생하는 비용입니다. '데이터 수집, 정제, 프롬프트 설계' 및 '지속적인 데이터 학습 및 성능 개선 체계 마련' 등의 활동은 모델의 생명주기 동안 끊임없이 투입되어야 하는 운영 비용입니다.

이처럼 TCO는 일회성 투자가 아닌, 인프라, 인력, 데이터 운영이 맞물려 지속적으로 발생하는 구조를 가지고 있어 신중한 재무적 검토가 반드시 필요합니다.

#### 3.2.3 공공기관·중견기업이 독자 LLM 개발에서 실패하는 구조적 이유

앞서 분석한 막대한 TCO와 자본 구조의 차이 외에도, 국내 다수의 공공기관 사업 입찰 조건 자체가 구조적 한계를 명확히 보여줍니다. 공공기관의 AI 사업의 예산과 규모가 파운데이션 모델 개발이 아닌, 기존 기술을 활용한 시스템 통합 및 응용 서비스 개발에 초점을 맞추고 있음을 방증합니다. 이러한 현실적인 제약 조건하에서 중소·중견기업이 선택할 수 있는 가장 합리적인 AI 도입 전략은 독자적인 모델 개발이 아닌, 검증된 외부 모델을 효율적으로 '활용'하는 데 집중하는 것입니다.

결론적으로, 파운데이션 모델 학습은 소수의 글로벌 기업을 제외한 대부분의 조직에게 경제적 장벽이 매우 높은 선택지입니다. 그렇다면 대안은 무엇일까요? 단순히 GPU 인프라를 증설하는 것만으로 AI 시대의 경쟁력을 확보할 수 있을까요? 다음 섹션에서는 이러한 접근법이 왜 충분하지 않은지에 대한 문제를 제기하고 새로운 해법을 모색하겠습니다.

#### 3.3 GPU 증설 중심 접근의 한계와 Idle 자원 문제

많은 조직이 AI 도입을 서두르며 GPU 서버 증설에 집중하지만, 이러한 투자가 기대했던 업무 혁신으로 곧바로 이어지지 않는 경우가 많습니다. 이는 AI 도입의 본질을 오해한 결과입니다. AI 역량은 단순히 연산 능력의 총량, 즉 '발전소'의 규모만으로 결정되지 않습니다. 생산된 AI라는 '전력'을 조직 내 모든 업무 현장까지 안전하고 효율적으로 전달할 '전력망'이 부재하다면, 아무리 강력한 발전소도 무용지물입니다. 본 섹션에서는 GPU 증설 중심 접근법의 근본적인 한계와 이로 인해 발생하는 유휴(Idle) 자원 문제를 논증하겠습니다.



#### 3.3.1 GPU 수 증가가 곧 업무 자동화 증가로 이어지지 않는 이유

GPU는 '발전소'이고, MCP는 '전력망'입니다. GPU 서버 수를 늘리는 것은 발전소의 총 발전 용량을 키우는 것에 불과합니다. 그러나 이 전력을 각 업무 시스템이라는 '소비자'에게 안정적으로 연결하고, 필요에 따라 전력량을 제어하며, 사용량을 측정할 표준화된 인터페이스, 즉 MCP(Model Context Protocol)가 없다면 AI 기능은 실제 업무 프로세스에 적용될 수 없습니다. GPU 클러스터와 업무 시스템 사이에 놓인 이 거대한 단절이 바로 GPU 투자 효과가 나타나지 않는 근본적인원인입니다.

#### 3.3.2 활용되지 못하는 연산 자원과 조직 내부 프로세스의 미정비

결과적으로, 전력망(MCP)이 부재한 상황에서 증설된 GPU는 '활용되지 못하는 Idle 자원'으로 전락할 수밖에 없습니다. 이는 막대한 초기 투자 비용이 회수되지 못하고 감가상각으로 사라지는 심각한 자원 낭비를 초래합니다.

더 나아가, 기술적 연결의 부재뿐만 아니라 조직 내부 프로세스의 미정비 또한 AI 도입의 큰 걸림돌이 됩니다. AI와 같은 혁신적인 신기술을 효과적으로 도입하기 위해서는, 기술 자체의 도입에 앞서 이를 수용할 수 있는 조직 내부 프로세스의 현대화와 유연성 확보가 반드시 선행되어야 합니다. 낡은 수도관에 강력한 수압의 물을 흘려보낼 수 없듯, 경직된 레거시 프로세스에 AI를 접목하는 것은 비효율과 실패를 낳을 뿐입니다.

#### 3.3.3 GPU 투자와 MCP·RAG 투자 간의 우선순위 재조정

따라서 합리적인 AI 투자 전략은 우선순위의 재조정을 요구합니다. GPU 증설과 같은 '생산 설비' 투자에서, 조직의 AI 활용 능력을 극대화하는 '활용 인프라' 투자로 무게 중심을 옮겨야 합니다. 여기서 활용 인프라의 두 핵심 축은 다음과 같습니다.

- 1. MCP (Model Context Protocol): 시스템과 AI를 연결하는 '전력망'
- 2. RAG (Retrieval-Augmented Generation): 조직의 지식 자산과 AI를 연결하는 '지식 파이프라인'

MCP가 AI가 수행할 수 있는 '동사(Actions)'를 제공한다면, RAG는 그 행동의 근거가 되는



'명사(Knowledge)'를 제공합니다. 진정으로 지능적인 시스템은 아는 능력과 실행하는 능력을 모두 필요로 합니다.

새로운 LLM을 만드는 대신, 기존의 방대한 문헌 정보를 RAG 구조로 설계하여 LLM과 연결함으로써 질병과 유전자 간의 관계를 추론하는 것을 목표로 합니다. 이는 새로운 모델을 개발하는 것보다, 조직이 이미 보유한 지식 자산을 AI가 활용할 수 있도록 '연결'하는 투자가 훨씬 더 즉각적이고 높은 투자수익률(ROI)을 제공함을 명확히 증명하는 사례입니다.

GPU 중심의 접근법은 명백한 한계를 가집니다. 이를 극복하기 위한 대안은 AI의 '생산'이 아닌 '활용'에 초점을 맞춘 추론 중심 전략이며, 그 전략의 기술적 핵심에는 바로 MCP가 자리하고 있습니다. 다음 섹션에서는 이 추론 중심 전략과 MCP의 역할에 대해 더 깊이 논의하겠습니다.

#### 3.4 추론 중심 전략과 MCP 수의 연계

AI 투자의 초점을 모델 학습에서 추론 활용으로 전환할 때, 우리는 새로운 질문에 직면하게 됩니다: 기업의 실질적인 AI 도입 성과를 어떻게 측정하고 관리할 것인가? GPU 보유량이나 모델 파라미터 수가 더 이상 유효한 지표가 아니라면, 무엇이 AI 성숙도를 가늠하는 새로운 척도가 되어야하는가? 이 질문에 대한 해답은 추론 중심 전략의 핵심인 MCP에 있습니다. 본 섹션에서는 'MCP의 수'가 기업의 실질적인 AI 적용 범위와 깊이를 나타내는 핵심 성과 지표(KPI)가 되는 이유를 구체적인 사례를 통해 증명하겠습니다.

#### 3.4.1 MCP 수 = 실제 추론이 개입하는 업무 수

"기업이 보유하고 운영하는 MCP 서버의 수는 곧 AI를 통해 자동화하거나 지능화한 실제 업무의 범위와 비례한다."이 주장은 MCP가 단순한 기술 프로토콜이 아니라, AI와 실제 비즈니스 프로 세스가 만나는 '업무 접점' 그 자체라는 사실에 기반합니다.

'LLM 기반 지능형 정보시스템 구축' 제안요청서의 기능 요구사항은 이를 명확히 보여주는 예시입니다.

• 요구사항 FUN-009 (총무관리): '지급의뢰서 일괄 출력', '반려 신청 및 자동화 처리' 등 다수의 기능



• 요구사항 FUN-010 (인사관리): '참여연구진 전원의 재직증명서와 연구실적증명서 일괄 발급', '채용관리 시스템 개선' 기능

이 각각의 기능들은 AI 에이전트가 호출할 수 있는 별도의 MCP '도구(Tool)'로 구현될 수 있습니다. 즉, '지급의뢰서 출력 도구', '재직증명서 발급 도구', '채용 공고 연계 도구' 등이 각각 하나의 MCP 인터페이스를 통해 제공됩니다. 이처럼 세분화된 도구 기반 접근법은 각 도구가 독립적으로 배포 및 관리가 가능한 자체 포함 비즈니스 역량이라는 마이크로서비스 아키텍처(MSA) 원칙과 정확히 일치합니다. 이는 추상적인 AI 역량이 아닌, AI가 직접 개입하여 처리하는 구체적인업무의 수를 의미합니다. 따라서 조직이 보유한 MCP의 수는 AI가 얼마나 많은 실제 업무에 통합되었는지를 나타내는 가장 직접적이고 정량적인 지표가 됩니다.

# 3.4.2 MCP 수를 기준으로 한 투자 우선순위: 어떤 업무부터 AI화할 것인가 MCP 인벤토리를 구축하는 것은 AI 도입을 위한 체계적인 로드맵을 수립하는 가장 효과적인 방법입니다. 업무 도메인과 MCP화 대상 업무를 매핑하면 다음과 같은 인벤토리를 구성할 수 있습니다.

업무 도메인	관련 정보 시스템	MCP화 대상 (업무 접점)
연구관리	연구정보지식시스템, 한국연구	LLM 기반 제규정 자연어 질
	자정보(KRI)	의, 연구업적 정보 통합 관
		리
재난대응	소방안전지도 시스템, 재난 정	공간정보(GIS) 기반 작전도
	보 관리	관리, 타 기관 재난상황 정
		보 연동
민원/행정	대표홈페이지, The건강보험	환급금(지원금) 신청·조회,
	(APP)	임신·출산 진료비 신청·조
		회

이러한 인벤토리를 바탕으로, 조직은 리스크가 낮고 효과가 빠른 도메인부터 MCP 기반의 개념 증명(PoC)을 설계하고 점진적으로 확장해 나가야 합니다. 이 과정에서 MCP 수, 각 MCP의 호



출 빈도, 업무 처리 시간 단축률 등은 AI 투자의 효과를 측정하고 다음 투자 우선순위를 결정하는 핵심 지표로 활용되어야 합니다. 이는 막연한 기대가 아닌, 데이터 기반의 합리적인 AI 투자 의사결정을 가능하게 합니다.

#### 3.4.3 모델 학습에 투자해야 하는 극소수 사례를 구분하는 기준

기존 시스템과 데이터를 '활용'하는 추론 중심 전략을 지지하고 있습니다. 이는 대부분의 비즈니스 문제가 검증된 파운데이션 모델의 추론 능력과 조직 내부 데이터를 결합하는 RAG 및 MCP 아키텍처를 통해 해결될 수 있음을 시사합니다.

따라서 자체 모델 학습은 극히 예외적인 전략으로 간주되어야 합니다. 이는 기존의 검증된 파운데이션 모델과 RAG/MCP 조합으로도 도저히 해결할 수 없는, 매우 특수하고 미션 크리티컬한 요구사항이 존재하며, 해당 모델이 시장에서 대체 불가능한 독점적인 경쟁 우위를 제공할 수 있을 때에만 제한적으로 고려되어야 합니다. 예를 들어, 특정 분야의 독점 데이터를 기반으로 한 초고도의 전문성이 요구되는 영역이 이에 해당할 수 있습니다.

결론적으로, AI 시대의 진정한 경쟁력은 더 많은 GPU를 소유하는 하드웨어의 경쟁이 아닙니다. 그것은 조직의 핵심 업무 프로세스 곳곳에 얼마나 더 많은 MCP를 구축하여 AI 추론 능력을 깊숙이, 그리고 광범위하게 통합하는지에 달려 있습니다. AI를 단순한 기술 전시가 아닌 실질적인 비즈니스 동력으로 전환하고자 하는 조직이라면, 이제 GPU 수를 세는 하드웨어 중심의 자본적 지출(CapEx) 관점에서 벗어나, MCP 수를 통해 실제 업무 통합 가치를 측정하는 운영 효율성(OpEx) 관점으로 전환해야 할 때입니다.

## 제4장 MCP(Model Context Protocol)의 기술적 정의와 표준 생태계

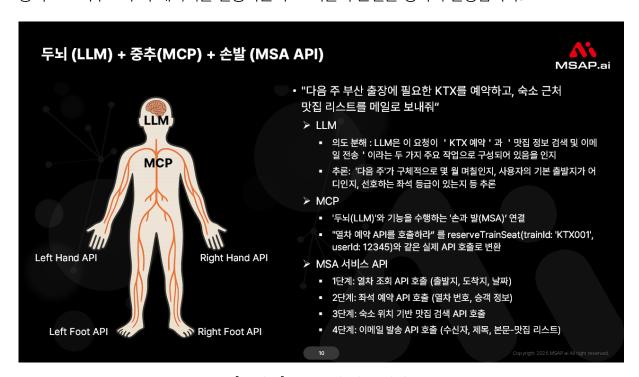
엔터프라이즈 AI 전략의 성공은 단순히 강력한 거대 언어 모델(LLM)을 도입하는 것을 넘어, 이 AI를 기업 내외부의 데이터 및 시스템과 얼마나 유기적으로 연결하여 실제 업무 프로세스에 녹여낼수 있는지에 달려있습니다. MCP는 단순한 연동 기술이 아닌, AI가 기업의 IT 자산을 이해하고 활



용하는 방식을 근본적으로 바꾸는 개방형 표준 프로토콜입니다. 본 장에서는 MCP의 기술적 정의와 아키텍처, 기존 연동 방식과의 차이점을 분석하고, 빠르게 성장하는 글로벌 생태계와 엔터프라이즈 도입 시 반드시 고려해야 할 보안 및 거버넌스 사항을 심도 있게 다룹니다. 이를 통해 MCP에 대한 깊이 있는 이해가 어떻게 기업의 AI 활용 패러다임을 혁신하고, 궁극적으로 비즈니스 경쟁력과 직결되는지를 명확히 제시하고자 합니다.

#### 4.1 MCP의 기본 개념과 아키텍처

Model Context Protocol(MCP)은 AI 에이전트와 외부 세계 간의 상호작용을 표준화하는 신경 망과 같습니다. LLM이 인간의 언어를 이해하는 두뇌라면, MCP는 외부 시스템이라는 신체의 각부분과 신호를 주고받으며 실제 행동을 이끌어내는 표준화된 신경계에 비유할 수 있습니다. 본 섹션에서는 MCP 아키텍처를 구성하는 핵심 요소와 그 작동 방식을 상세히 분석하여, AI가 어떻게 동적으로 외부 도구와 데이터를 활용하는지 그 기술적 본질을 명확히 설명합니다.



[그림 2] MCP 의 기본 개념

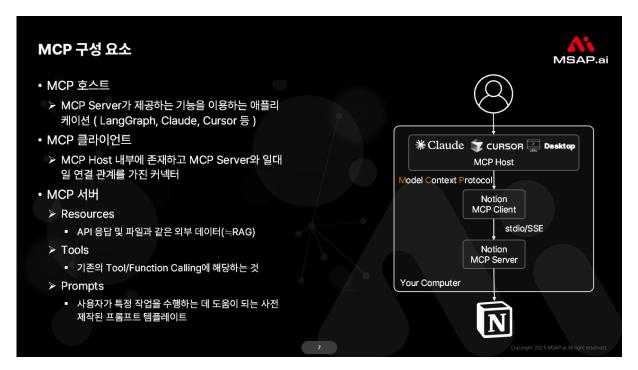


### 4.1.1 MCP의 목적: LLM과 외부 시스템을 연결하는 개방형 표준

MCP(Model Context Protocol)의 근본적인 목적은 LLM 기반 애플리케이션이 외부 데이터 소스 및 도구(Tools)와 원활하게 통합될 수 있도록 설계된 '개방형 상호운용성 프로토콜'을 제공하는 것입니다. 과거 개발자 도구가 각기 다른 언어를 지원하는 복잡성을 해결하기 위해 LSP(Language Server Protocol)라는 표준을 만들어 생태계를 구축하고 혁신을 이끌었던 것에서 영감을 받았습니다. MCP 역시 다양한 AI 에이전트, LLM, 그리고 기업용 도구들이 서로 다른 기술 스택에 구애받지 않고 안전하게 통신하고 구조화된 작업을 호출할 수 있는 표준 인터페이스 역할을 수행합니다. 이를 통해 개발자들은 복잡하고 파편화된 개별 연동 작업에서 벗어나 AI의 핵심 기능 구현에 집중할 수 있으며, AI 통합의 복잡성을 해결하고 표준화를 통해 혁신을 가속하는 것을 목표로 합니다.

4.1.2 MCP의 핵심 구성: Host, Client, Server, Tools, Resources, JSON-RPC

MCP 아키텍처는 명확한 역할을 가진 여러 구성 요소의 유기적인 결합으로 이루어집니다. 각 요소는 AI 시스템의 안전하고 효율적인 외부 자원 활용을 보장하도록 설계되었습니다.



[그림 3] MCP 구성 요소



- Host: LLM을 내장하고 사용자 요청을 중재하는 전체 AI 애플리케이션 또는 환경을 의미합니다. VS Code와 같은 IDE나 Claude Desktop이 대표적인 예시이며, AI가 접근할 수 있는 외부 자원을 결정하고 모든 통신을 통제하는 중앙 게이트키퍼(Gatekeeper)이자 보안 브로커(Security Broker) 역할을 수행합니다.
- Client: Host 내부에 존재하며, LLM의 요청을 MCP 표준 메시지 형식으로 변환하여 적절한 Server에 전달하는 통신 오케스트레이터입니다. Host 내부에 존재하는 Client는 특정 MCP Server와의 통신을 전담하는 구성 요소입니다. 즉, Host는 여러 Client를 통해 동시에 다수의 Server와 상호작용할 수 있습니다.
- Server: LLM에 필요한 컨텍스트, 데이터, 또는 실행 가능한 기능을 제공하는 외부 서비스를 의미합니다. GitHub, 데이터베이스, 슬랙(Slack) 등 기업의 내외부 시스템이 MCP 서버로 구현될 수 있습니다.
- Tools: Server가 제공하는 기능 중, AI 모델이 특정 작업을 수행하기 위해 호출할 수 있는 함수(functions)를 의미합니다. 예를 들어, '사용자 정보 조회'나 '파일 쓰기'와 같은 구체적 인 액션이 도구에 해당합니다.
- Resources: Server가 제공하는 데이터로, AI 모델이나 사용자에게 풍부한 컨텍스트를 제공하는 역할을 합니다. 파일 내용, 데이터베이스 스키마, Git 히스토리 등이 리소스의 예시입니다.
- JSON-RPC: Client와 Server 간의 모든 구조화된 데이터 교환에 사용되는 표준 통신 메시지 형식입니다. 이를 통해 언어나 플랫폼에 구애받지 않는 상호운용성을 확보합니다.

### 4.1.3 요청 라우팅·모델 오케스트레이션·컨텍스트 주입 흐름

MCP 기반의 모델 오케스트레이션은 사전에 정의된 규칙이 아닌, LLM의 실시간 판단에 따라 작업 흐름이 동적으로 결정된다는 점에서 기존의 자동화와 차별화됩니다. LLM이 스스로 문제를 분석하고 최적의 도구를 선택하여 해결하는 과정은 다음과 같은 단계로 이루어집니다.

1. 도구 탐색: MCP Client가 연결된 Server로부터 현재 사용 가능한 도구(Tools)의 전체 목록과 각 도구가 어떤 기능을 수행하는지에 대한 설명(description)을 받아옵니다.



- 2. 쿼리 및 컨텍스트 전송: 사용자의 요청(쿼리)이 1단계에서 확보한 도구 설명과 함께 LLM(예: Claude)에게 컨텍스트로 전달됩니다.
- 3. 도구 사용 판단: LLM은 자신의 사전 지식, 사용자의 쿼리 내용, 그리고 제공된 도구 설명을 종합적으로 분석하여 문제 해결에 가장 적합한 도구가 있는지 스스로 판단합니다.
- 4. 도구 실행 요청: LLM이 특정 도구를 사용하기로 결정하면, 이 결정은 MCP Client에게 전 달되고, Client는 해당 도구의 실행을 Server에 공식적으로 요청합니다.
- 5. 결과 반환 및 컨텍스트화: Server에서 도구를 실행한 결과(예: API 응답, 데이터베이스 조회 결과)는 다시 LLM에게 추가적인 컨텍스트로 전달됩니다.
- 6. 최종 응답 생성: LLM은 도구 실행 결과를 바탕으로 사용자의 초기 질문에 대한 최종적인 자연어 응답을 생성하여 제시합니다.

이처럼 LLM의 추론 능력을 중심으로 외부 도구를 동적으로 선택하고 실행하는 아키텍처는, 모든 연동 로직이 사전에 고정되어야 했던 기존의 정적인 연동 방식과는 근본적으로 다른 차원의 유연성과 확장성을 제공합니다.

### 4.2 기존 연동 방식과 MCP의 차이

과거 AI 시스템의 확장을 가로막았던 가장 큰 장벽은 '통합의 복잡성' 문제였습니다. 새로운 AI 모델이나 외부 도구가 추가될 때마다 연결 지점이 기하급수적으로 늘어나 시스템 전체가 경직되고 유지보수가 불가능해지는 현상이 빈번했습니다. MCP는 이러한 문제를 구조적으로 해결하기 위해 등장한 표준 프로토콜입니다.

### 4.2.1 API·Webhook·플러그인 방식의 N×M 연동 문제

기존의 점대점(Point-to-Point) API 연동 방식은 'N×M 통합 문제'라는 고질적인 복잡도를 야기합니다. 예를 들어, 3개의 서로 다른 LLM(N=3)을 기업 내 5개의 주요 도구(M=5) - CRM, 캘린더, 고객 지원 시스템, 내부 지식 베이스, 문서 분석기 - 와 연동해야 하는 상황을 가정해 보겠습니다. 이 경우, 각 LLM은 5개의 도구와 개별적으로 연동되어야 하므로 총  $3 \times 5 = 15$ 개의 맞춤형 연결(Custom Connection)이 필요합니다.





[그림 4] 기존 LLM 연동 방식과 MCP 의 차이점

이러한 아키텍처는 초기 구축 비용도 높지만, 더 큰 문제는 유지보수와 확장성입니다. 새로운 도구 하나가 추가되면 기존 3개의 LLM과 모두 연동해야 하므로 3개의 연결이 추가로 필요하며, 새로운 LLM이 도입되면 기존 5개의 도구와 모두 연결해야 하므로 5개의 연결이 추가로 필요합니다. 이처럼 시스템의 규모가 커질수록 연결 지점은 기하급수적으로 증가하여 결국에는 관리가 불가능하고 변화에 대응할 수 없는 취약한 구조를 만들게 됩니다.

### 4.2.2 MCP 도구 호출 모델이 연동 복잡도를 N+M으로 줄이는 방식

MCP는 N×M 연동 문제를 N+M 수준으로 획기적으로 단순화합니다. 이는 MCP가 AI 모델과 외부 도구 사이에 '통합된 프로토콜 계층(Unified Protocol Layer)' 역할을 수행하기 때문에 가능합니다.

MCP 아키텍처에서는 모든 도구(N개)가 표준 MCP 서버 인터페이스를 구현하고, 모든 LLM 애플리케이션(M개)은 표준 MCP 클라이언트 인터페이스만 구현하면 됩니다. 즉, 각 구성 요소는 MCP라는 단일 표준에 맞춰 한 번만 구현하면 생태계 내의 다른 모든 구성 요소와 상호작용할 수 있습니다.

따라서 3개의 LLM과 5개의 도구를 연동하는 시나리오에서 필요한 연결은 각 LLM이 MCP



클라이언트를 구현하는 3개와 각 도구가 MCP 서버를 구현하는 5개를 더한, 총 3 + 5 = 8개의 연결만으로 충분합니다. 새로운 도구가 추가되면 해당 도구의 MCP 서버 하나만 추가하면 되고, 새로운 LLM이 추가되어도 해당 LLM의 MCP 클라이언트 하나만 추가하면 됩니다. 이처럼 전체 시스템의 연결 복잡도가 기하급수적이 아닌 선형적으로 증가하므로, 확장성과 유지보수성이 극대화됩니다.

### 4.2.3 단순 HTTP API 호출과 컨텍스트 인식 프로토콜로서 MCP의 차별점

MCP는 단순히 기존의 REST API를 대체하는 기술이 아닙니다. 두 프로토콜은 설계 철학과 핵심 목적에서 근본적인 차이를 보이며, MCP는 AI와의 상호작용에 최적화된 새로운 패러다임을 제시 합니다.

		MCP (Model Context
구분	기존 REST API	Protocol)
설계 패러다임	자원 중심	역량 중심
	(Resource-Oriented): 특정	(Capability-Oriented): AI7
	자원(데이터)을 생성, 조회, 수	'무엇을 할 수 있는지'에 대한
	정, 삭제하는 데 중점을 둡니	역량을 노출하고 실행하는 데
	다.	중점을 둡니다.
상호작용 방식	정적 호출: 사전에 정의된 엔드	동적 탐색 및 오케스트레이션:
	포인트를 직접 호출하는 방식	LLM이 실시간으로 사용 가능
	으로, 상호작용이 고정적입니	한 도구를 탐색하고, 상황에 맞
	다.	게 최적의 도구를 선택하여 호
		출 흐름을 동적으로 결정합니
		다.



		MCP (Model Context
구분	기존 REST API	Protocol)
핵심 목적	데이터 교환: 애플리케이션 간	컨텍스트 인식 기반 실행:
	의 데이터 검색 및 조작이 주된	LLM에게 풍부한 컨텍스트를
	목적입니다.	제공하고, 그 판단에 따라 외부
		시스템에서 실제 작업을 수행
		하게 하는 것이 주된 목적입니
		다.

결론적으로, MCP의 진정한 차별점은 정적인 데이터 호출을 넘어선 '컨텍스트 인식'과 '역량 중심' 설계에 있습니다. AI가 단순히 데이터를 요청하는 것을 넘어, "이 문제를 해결하기 위해 어떤 도구를 사용해야 하는가?"를 스스로 판단하고 실행하게 만드는 이 능력이 바로 MCP 생태계 확장의 핵심 원동력입니다.

### 4.3 MCP 생태계와 글로벌 표준 동향

MCP는 특정 기업의 독점 기술을 넘어, 주요 빅테크 기업들이 적극적으로 참여하고 플랫폼 레벨에 통합하면서 AI 연동의 산업 표준으로 빠르게 자리 잡고 있습니다. 이는 MCP가 제공하는 상호 운용성과 확장성이 개별 기업의 이해관계를 초월하는 가치를 지니고 있음을 증명하며, AI 시대의 필수적인 기반 기술로 부상하고 있음을 보여줍니다.

4.3.1 Anthropic·OpenAl·Microsoft·GitHub 등 주요 벤더의 MCP 채택 현황

글로벌 기술 시장을 선도하는 주요 벤더들이 MCP를 자사의 핵심 AI 제품 및 플랫폼에 통합하며 생태계 확장을 주도하고 있습니다.

• Anthropic: MCP 개념을 최초로 제안했으며, 자사의 Claude.ai 및 Claude Desktop을 통해 MCP 생태계를 선도하는 핵심적인 역할을 수행하고 있습니다.



- Microsoft: 개발자 생태계에서의 영향력을 바탕으로 Copilot Studio와 VS Code에 MCP 지원을 공식적으로 통합하여, 수많은 개발자들이 AI 에이전트를 쉽게 구축하고 활용할 수 있는 기반을 마련했습니다.
- OpenAI: 경쟁사의 프로토콜임에도 불구하고 자사의 AI 에이전트가 MCP를 지원하도록 함으로써, 특정 모델에 종속되지 않는 개방형 프로토콜로서의 MCP의 위상을 강화하고 상호운용성을 높이는 데 기여했습니다.
- GitHub: GitHub Copilot이 MCP 클라이언트로 작동하여, 코드베이스 분석, 풀 리퀘스트(PR) 검토, 이슈 추적 등 개발과 관련된 다양한 작업을 외부 MCP 서버를 통해 수행할 수 있도록 지원합니다.
- AWS: 클라우드 환경과의 긴밀한 통합을 주도하며, 코드 어시스턴트를 위한 AWS MCP Servers 제품군(예: AWS CDK 서버, 비용 분석 서버 등)을 공식적으로 출시하여 개발자들이 클라우드 자원을 AI 에이전트를 통해 효과적으로 관리할 수 있도록 지원합니다.

이러한 채택 현황은 MCP가 특정 영역에 국한되지 않음을 보여준다. 마이크로소프트와 GitHub는 개발자 워크플로우를, AWS는 클라우드 자원 관리를, 그리고 Anthropic과 OpenAl는 핵심 AI 에이전트 기능을 중심으로 생태계를 확장하며, MCP가 엔터프라이즈 AI의 다양한 접점에서 표준 인터페이스로 기능할 수 있는 잠재력을 입증하고 있다.

### 4.3.2 OS·IDE·브라우저·에이전트 런타임 등 플랫폼 레벨 통합 흐름

MCP의 영향력은 개별 애플리케이션을 넘어 운영체제(OS), 통합 개발 환경(IDE) 등 플랫폼 수준으로 깊숙이 확장되고 있습니다. 이는 AI 기능이 특정 앱의 부가 기능이 아닌, 플랫폼 자체의 핵심역량으로 자리 잡고 있음을 의미합니다.

- IDE 및 코드 에디터: VS Code, Zed, Cursor 등 다수의 최신 개발 환경이 MCP를 내장하여, 코드 분석, 파일 시스템 접근, 버전 관리 등 복잡한 개발 워크플로우를 AI 에이전트가 직접 지원할 수 있게 합니다.
- 데스크톱 애플리케이션: Claude Desktop과 같은 네이티브 애플리케이션이 MCP 호스트 역할을 수행하며, 로컬 파일 시스템이나 설치된 프로그램 등 개인화된 시스템 자원과의 안 전하고 강력한 연동을 제공합니다.



• 에이전트 프레임워크: Firebase Genkit, LangChain 어댑터 등 차세대 AI 에이전트 개발 프레임워크들이 MCP를 외부 도구와의 표준 연동 방식으로 채택하고 있어, 개발자들이 보다 쉽게 강력한 에이전트를 구축할 수 있는 기반을 제공합니다.

### 4.3.3 MCP 생태계와 글로벌 표준 동향

MCP 생태계는 자발적인 참여와 투명한 거버넌스를 통해 빠르게 성장하고 있으며, 개발자들의 참여를 촉진하는 다양한 지원 체계를 갖추고 있습니다.

- 레퍼런스 서버(Reference Servers): MCP에 기여하는 핵심 개발자들이 직접 유지보수하는 공식 예제 서버들입니다. 이는 새로운 개발자들이 프로토콜의 핵심 기능을 이해하고 자체 서버를 구현할 때 참고할 수 있는 가장 신뢰도 높은 기준점 역할을 합니다.
- 공식 SDK (Software Development Kits): TypeScript, Python, Go, Kotlin 등 주요 프로그래밍 언어별로 공식 SDK를 제공합니다. 이를 통해 개발자들은 프로토콜의 세부 사항을 모두 이해하지 않아도 쉽고 안정적으로 MCP 클라이언트와 서버를 구축할 수 있습니다.
- 커뮤니티 주도 로드맵: MCP의 발전 방향은 특정 기업이 아닌 커뮤니티에 의해 투명하게 결정됩니다. 현재 로드맵에서는 장시간 소요되는 작업을 위한 비동기 작업(Asynchronous Operations) 지원, 대규모 엔터프라이즈 환경을 위한 무상태(Statelessness) 및 확장성 개선 등이 활발하게 논의되고 있으며, 이는 프로토콜이 실제 산업 현장의 요구사항을 반영하며 지속적으로 발전하고 있음을 보여줍니다.

이러한 표준화 흐름과 강력한 생태계는 MCP 도입을 고려하는 기업에게 안정성과 미래 확장 성을 보장합니다. 이제 기업들은 이러한 흐름에 동참하기 위해 실질적인 도입 전략과 고려사항을 검토해야 할 시점입니다.

### 4.4 엔터프라이즈 도입 시 고려사항

MCP 도입은 단순히 새로운 기술을 선택하는 것을 넘어, 기업의 민감한 데이터와 핵심 시스템에 AI 에이전트의 접근을 허용하는 중대한 전략적 결정입니다. 따라서 성공적인 도입을 위해서는 보안, 거버넌스, 그리고 MCP 환경에서 발생할 수 있는 잠재적 위협에 대한 깊은 이해가 반드시 필요합니다.



### 4.4.1 인증·인가·감사·네트워크 구조(Host 중심 보안 모델)

MCP 보안의 핵심은 '중재 접근 패턴(Mediated Access Pattern)'에 기반한 Host 중심 보안 모델입니다. 이 모델에서 MCP Host는 Al 모델과 외부 자원 사이에서 모든 상호작용을 통제하고 검증하는 '보안 브로커(Security Broker)' 역할을 수행합니다. 이는 Al가 통제되지 않은 상태로 외부 시스템에 직접 접근하는 것을 원천적으로 차단하고, 모든 활동을 감사할 수 있는 단일 지점을확보하는 가장 중요한 보안 원칙입니다.

엔터프라이즈 환경에서 이 모델을 성공적으로 구현하기 위해서는 다음과 같은 보안 통제 방안 이 필수적입니다.

- 인증 및 인가 (Authentication & Authorization): 모든 MCP 서버는 **0Auth 2.1**과 같은 표준 인증 프로토콜을 의무적으로 구현해야 합니다. 특히, 인증을 처리하는 서버와 실제 자원을 제공하는 리소스 서버는 반드시 분리되어야 합니다. 이를 통해 중앙화된 ID 관리 시스템(IdP)과의 연동이 가능해지며, 안전한 토큰 기반의 인증 및 인가 체계를 구축할 수 있습니다.
- 세분화된 권한 관리 (Fine-Grained Permissions): 최소 권한 원칙(Principle of Least Privilege)에 따라, 접근 제어 목록(ACLs)을 사용하여 AI 에이전트가 특정 작업을 수행하는 데 필요한 최소한의 도구와 데이터에만 접근하도록 엄격히 통제해야 합니다. 예를 들어, '읽기 전용' 에이전트에게는 파일 쓰기 도구에 대한 접근 권한을 부여해서는 안 됩니다.
- Host의 통제 역할: 모든 요청은 반드시 Host를 경유하도록 네트워크 아키텍처를 설계해야합니다. Al 모델이 임의의 서버에 직접 네트워크 요청을 보내는 것을 원천적으로 차단함으로써, Host가 모든 상호작용을 감사하고, 비정상적인 활동을 탐지하며, 정책에 위반되는 요청을 차단하는 중앙 통제 지점으로서의 역할을 완벽히 수행할 수 있도록 보장해야합니다.

### 4.4.2 프롬프트 인젝션·데이터 유출 등 MCP 특유의 보안 위협

MCP는 강력한 기능을 제공하는 만큼, 이를 악용하려는 새로운 유형의 보안 위협에 노출될 수 있습니다. 기업은 이러한 위협을 사전에 인지하고 방어 전략을 수립해야 합니다.



보안 위협 유형	설명
프롬프트 인젝션 (Prompt Injection)	악의적인 사용자가 교묘하게 조작된 입력(프롬
	프트)을 통해 LLM을 속여, 개발자가 의도하지
	않은 MCP 도구를 실행하게 만들거나 민감한
	정보를 노출하도록 유도하는 공격입니다.
데이터 유출 및 비인가 접근	권한이 부적절하게 설정된 MCP 도구를 통해
	AI 에이전트가 자신의 접근 권한 범위를 벗어나
	는 데이터(파일, 데이터베이스 레코드 등)를 읽
	거나 외부로 유출시키는 심각한 위협입니다.
잘못된 구성 (Poor Configuration)	인증/인가 범위(Scope)를 지나치게 광범위하
	게 설정하거나 보안 설정을 잘못 구성하여, 해
	킹 또는 내부자 위협 발생 시 피해 범위(Blast
	Radius)를 크게 확대시키는 인적 오류 기반의
	위협입니다.

### 4.4.3 MCP 버전 전략·거버넌스: 등록·검증·배포 정책

기업이 장기적인 관점에서 MCP를 신뢰하고 도입하기 위해서는 프로토콜 자체가 안정적이고 투명한 거버넌스 체계 위에서 발전해야 합니다. MCP는 특정 벤더에 종속되지 않는 커뮤니티 주도 거버넌스 모델을 채택하여 이를 보장합니다.

- Working Groups 및 Interest Groups: MCP 커뮤니티 내에는 보안, 확장성, 특정 도메인 적용 등 특정 주제를 심도 있게 논의하고 해결책을 모색하는 공식적인 그룹들이 활동하고 있습니다. 이를 통해 산업 현장의 다양한 요구사항이 프로토콜 발전에 체계적으로 반영됩니다.
- SEP (Specification Enhancement Proposal) 워크플로우: 프로토콜에 대한 모든 중요한 변경 제안은 'SEP'라는 공식 제안서 양식을 통해 제출되어야 합니다. 제출된 제안서는 커뮤니티의 충분한 공개 토론과 기술적 검토, 그리고 합의 과정을 거쳐 투명하게 채택 여부가 결정됩니다. 이러한 공식 절차는 MCP가 특정 벤더의 이해관계에 따라 급진적으로 변경



되는 것을 방지하고, 안정적으로 발전할 수 있는 강력한 기반이 됩니다.

결론적으로, 본 장에서 분석한 MCP는 단순한 기술 프로토콜을 넘어 엔터프라이즈 AI 아키텍처의 근본적인 진화를 의미합니다. MCP는 N×M 통합 문제를 N+M 수준으로 해결하는 구조적 우아함을 제공할 뿐만 아니라, 정적인 자원 중심 API에서 동적인 역량 중심 프로토콜로의 패러다임 전환을 이끈합니다. Anthropic, Microsoft, AWS 등 경쟁사를 아우르는 광범위한 벤더들의 채택으로 빠르게 성장하는 벤더 중립적 생태계는 기업에게 기술 종속성을 피하고 장기적인 확장성을 확보할 수 있는 전략적 이점을 제공합니다. 마지막으로, Host를 보안 브로커로 삼는 중재 접근패턴과 투명한 거버년스를 기반으로 한 보안 우선주의적 접근 방식은, MCP의 강력한 기능을 엔터프라이즈 환경에서 안전하게 활용하기 위한 필수 전제 조건입니다. 이 모든 요소를 종합적으로이해하고 내재화하는 것이야말로 AI 시대를 선도하는 기업의 핵심 경쟁력이 될 것입니다.

### 제5장 AI 성숙도 모델과 MCP 수 기반 평가 프레임워크

### 5.1 AI 성숙도의 새로운 지표: '업무 접점 수'

기업의 AI 경쟁력을 평가하는 전통적인 방식은 보유한 거대 언어 모델(LLM)이나 그래픽 처리 장치 (GPU) 인프라 규모에 초점을 맞춰왔습니다. 그러나 이러한 지표들은 잠재력의 후행 지표(lagging indicator)에 불과합니다. AI 기술의 진정한 가치는 보유가 아닌 '활용'에서 비롯되며, 실제 비즈 니스 프로세스와의 통합 수준이야말로 경쟁 우위의 핵심입니다.

본 장에서는 AI의 실현된 가치와 경쟁적 민첩성을 측정하는 선행 지표(leading indicator)로 서 '업무 접점 수', 즉 MCP(Model Context Protocol) 수를 제안합니다. MCP는 AI 모델이 기업의 내부 시스템과 외부 도구에 연결되는 표준화된 통로이며, 그 수는 AI가 조직의 비즈니스에 얼마나 깊숙이 통합되어 있는지를 나타내는 가장 직접적인 척도입니다. 이 프레임워크는 AI 경쟁력 확보를 위해 모든 조직이 채택해야 할 필수 전략이며, AI 투자의 패러다임을 '인프라 비축'에서 '가치 창출을 위한 연결 확대'로 전환할 것을 촉구합니다.



### 5.1.1 모델·GPU 보유 중심 성숙도 평가의 한계

GPU 클러스터나 LLM 모델을 소유하는 것만으로는 비즈니스 가치를 창출할 수 없습니다. 이는 발전소와 전력망의 관계에 비유할 수 있습니다.

- 발전소 (GPU·LLM): AI 기능을 생성하는 핵심 인프라입니다. 높은 컴퓨팅 파워를 제공하여 복잡한 연산을 처리합니다.
- 전력망 (MCP): 발전소에서 생산된 전력(Al 기능)을 실제 소비처(업무 시스템, 사용자)까지 전달하는 표준화된 네트워크입니다.

최첨단 발전소(GPU)를 수백 기 보유하고 있더라도, 각 가정과 공장으로 전력을 실어 나를 전력망(MCP)이 없다면 발전소는 가동되지 않는 유휴 자산에 불과합니다. 마찬가지로, 기업이 막대한 비용을 들여 GPU 클러스터를 구축해도 생성된 AI 기능을 실제 재무, 인사, 고객 관리 시스템과 같은 업무 현장으로 전달할 MCP가 없다면 해당 투자는 가치를 증명할 수 없습니다. 결론적으로, GPU 투자의 ROI는 MCP라는 전력망의 구축 밀도와 범위에 의해 결정됩니다. 전력망 없이는 발전소 증설이 무의미한 지출에 불과한 것과 같습니다.

결국, GPU 보유량이나 모델의 파라미터 수는 AI의 '잠재력'을 보여줄 뿐, 실제 '경쟁력'을 대변하지 못합니다. 진정한 AI 경쟁력은 "조직 전체에 AI 기능이 얼마나 원활하게 흘러가고 있는 가?"라는 질문에 답할 수 있어야 하며, 그 해답은 MCP의 구축 수준에 달려 있습니다.

### 5.1.2 MCP 수로 측정하는 AI 활용 범위(업무 도메인·기능 수)

MCP는 LLM이 기업의 다양한 외부 시스템, 즉 데이터 소스, 내부 애플리케이션, 외부 API 등을 호출하고 상호작용할 수 있도록 연결하는 표준화된 인터페이스입니다. 기존의 API가 인간 개발자가 사용할 정적인 엔드포인트를 제공하는 것에 가깝다면, MCP는 AI 모델이 동적으로 발견하고 상황에 맞게 사용할 수 있는 '능력(Capability)'을 노출하는 프로토콜이라는 점에서 근본적인 차이가 있습니다. 따라서 조직 내에 구현된 MCP의 수는 AI가 관여하고 있는 업무 도메인 및 기능의 수를 직접적으로 반영합니다.

요구하는 기능들을 MCP화하는 시나리오를 통해 AI의 활용 범위를 구체적으로 파악할 수 있습니다.



업무 도메인	세부 기능	잠재적 MCP 명칭 (예시)
인사 관리	각종 증명서(재직, 연구실적	hr.issueCertificateBatch
	등) 일괄 발급	
총무 관리	전자증빙시스템을 활용한 지급	ga.printPaymentRequestBatch
	의뢰서 일괄 출력	
재난 대응	GIS 기반 차량 동태 실시간 관	disaster.getVehicleStatus
	리	
재난 대응	기상청 외부 데이터 연계	disaster.getExternalWeather
서비스 (예약)	온라인 공연 예매 기능 구현	service.reservePerformanceT
서비스 (예약)	예매 내역 확인 및 환불 기능	service.processRefund
	구현	
지식 관리	내부 규정 및 매뉴얼에 대한 자	knowledge.queryInternalDocs
	연어 질의응답	

위 표에서 볼 수 있듯, MCP가 하나씩 추가될 때마다 AI가 수행할 수 있는 업무의 범위는 인사, 총무, 재난 대응, 고객 서비스 등 다양한 영역으로 명확하게 확장됩니다. 따라서 MCP의 총 수는 조직의 AI 활용 범위를 정량적으로 측정하는 가장 효과적인 지표가 됩니다.

### 5.1.3 MCP 수와 자동화율·의사결정 지원 비율의 상관관계

MCP 수의 증가는 AI가 관여하는 업무 프로세스의 확장을 의미하며, 이는 곧 조직의 자동화 수준과 직결되는 핵심 지표입니다. MCP 수가 많다는 것은 AI가 개입하여 스스로 처리하거나, 인간의 판단을 지능적으로 지원할 수 있는 업무 접점이 많다는 것을 의미하기 때문입니다. 이러한 관계는 기업의 핵심 성과와 직접적으로 연결됩니다.

- RAG 기반 질의응답 시스템: 내부 지식 베이스를 검색하는 MCP가 구현되면, 직원이 규정이 이나 매뉴얼을 찾기 위해 소비하던 시간이 AI 기반 자동응답으로 대체됩니다. 이는 직접적인 업무 자동화를 통해 직원 생산성을 향상시키고 운영 비용을 절감하는 사례입니다.
- 지출결의 자동화 처리: 재무 시스템의 '지출결의 자동화 처리' 기능이 MCP로 노출되면, Al 에이전트는 특정 조건에 따라 결재 문서를 자동으로 생성하거나 반려할 수 있습니다. 이는



복잡한 의사결정 과정을 지원하고 가속화하여 비즈니스 사이클을 단축시킵니다.

이처럼 MCP는 AI의 지능을 실제 업무 프로세스에 주입하는 통로입니다. 따라서 높은 MCP수는 감소된 운영 비용, 더 빠른 의사결정 주기, 향상된 직원 생산성으로 직접 변환되며, C-레벨경영진이 주목해야 할 핵심 KPI(핵심성과지표)입니다.

결론적으로, AI의 실제 가치를 평가하기 위해서는 모델의 성능이나 인프라 규모를 넘어, AI가 비즈니스와 얼마나 많은 '접점'을 가지고 있는지를 측정해야 합니다. MCP 수는 이러한 활용 범위를 객관적으로 나타내는 핵심 지표이며, 이를 체계적으로 측정하고 관리하기 위한 구체적인 방법론이 필요합니다.

### 5.2 MCP 수 측정 방법론

MCP 수를 조직의 AI 성숙도 핵심 지표로 활용하기 위해서는, 이를 체계적이고 일관되게 측정할수 있는 방법론이 필수적입니다. 주관적인 판단이나 개별 프로젝트 단위의 집계만으로는 전사적인 AI 도입 현황을 정확히 파악하기 어렵습니다.

본 섹션에서는 조직 내 모든 MCP를 식별하여 인벤토리를 구축하는 절차부터, 기능적 특성에 따라 분류하고, 나아가 AI 활용 성과를 측정하기 위한 핵심 지표를 설계하는 구체적인 방법론을 제시합니다. 이 방법론은 기업이 AI 도입 현황을 정량적으로 파악하고, 데이터에 기반하여 향후 투자 및 개선 방향을 도출할 수 있도록 지원하는 데 그 목적이 있습니다.

### 5.2.1 도메인·시스템별 MCP 인벤토리 구축 절차

조직의 모든 MCP 자산을 체계적으로 목록화하는 인벤토리 구축은 AI 활용 현황을 파악하는 첫걸음입니다. 다음 4단계 접근법을 통해 잠재적 MCP 후보를 발굴하고 실질적인 자산으로 전환할 수있습니다.

- 1. 자산 파악 조직이 보유한 모든 IT 자산(레거시 시스템, API 등)을 기능 단위로 분해하여 AI가 호출할 수 있는 잠재적 MCP 후보 목록을 작성합니다.
- 2. 최소 MCP 도메인 선정 처음부터 모든 시스템을 MCP화하는 것은 비효율적입니다. 리스크가 비교적 낮고 AI 도입 효과를 빠르게 검증할 수 있는 도메인(예: 내부 규정 질의, 고객 문의)을 우선적으로 선정합니다.



- 3. MCP 기반 PoC 설계 선정된 도메인의 가장 핵심적인 기능 3~5개를 MCP 도구로 구현하는 소규모 개념 증명(PoC)을 설계합니다. 이때, 내부 문서 검색을 위한 RAG(Retrieval-Augmented Generation) 인덱스와 결합하여 '지식 검색'과 '기능 실행'이 하나의 시나리오 안에서 이루어지도록 구성하는 것이 효과적입니다.
- 4. MCP 수 및 효과 측정 체계 정의 PoC 단계부터 성과를 측정할 지표를 명확히 정의합니다. 정량적 지표로는 '구현된 MCP 수', '적용된 업무 도메인 범위', 'MCP 호출 빈도' 등을, 정 성적 지표로는 '업무 처리 시간 단축률', '사용자 만족도' 등을 설정하여 AI 투자의 효과를 지속적으로 추적합니다.

### 5.2.2 MCP 카테고리 분류(조회·작성·승인·배포·모니터링 등)

구축된 MCP 인벤토리는 기능적 특성에 따라 분류해야 관리 및 활용이 용이합니다. 다양한 시스템 요구사항 분석을 통해 도출된 MCP를 다음과 같은 카테고리로 분류할 수 있습니다.

- 조회 (Inquiry) AI가 시스템의 데이터를 읽어 사용자에게 정보를 제공하거나 다음 행동을 결정하는 데 사용합니다.
  - 예시: 재규정 자연어 질의, 과거 민원접수 현황 조회, 실시간 재난 현황 검색, 가상계 좌 조회.
- 작성 (Creation) AI가 새로운 데이터나 콘텐츠를 생성하여 시스템에 기록하는 작업을 수행합니다.
  - 예시: 여입결의서 초안 자동 작성, 회의록 기반 보고서 초안 생성, 공지사항 게시물 등록.
- 실행/처리 (Execution/Processing) 단순 조회를 넘어 시스템의 특정 프로세스를 실행하 거나 데이터를 변경하는 작업을 포함합니다.
  - 예시: 지급의뢰서 일괄 출력, 데이터 이관 및 검증 작업 실행, 온라인 공연 예매 및 환급금 신청 처리.
- 승인 (Approval) 워크플로우와 관련된 승인, 반려 등의 의사결정 프로세스를 AI가 수행하거나 지원합니다.



- 예시: 지출결의 자동화 처리(반려/승인), 시스템 접근 권한 신청 자동 승인, 완료 문서 반려 프로세스 실행.
- 모니터링 (Monitoring) 시스템의 상태, 성능, 로그 등을 주기적으로 확인하여 이상 징후를 감지하고 보고하는 작업을 수행합니다.
  - 예시: 연계 시스템 API 상태 모니터링, 인프라 자원(CPU, Memory) 사용량 수집, 실시간 챗봇 상담 진행 상황 확인.

다수의 공공 및 서비스 분야 RFP를 분석한 결과, 초기 AI 도입은 '조회' 유형의 MCP에 집중되지만, 성숙도가 높아질수록 '작성' 및 '실행/처리'와 같이 시스템 상태를 변경하는 능동적인 MCP에 대한 요구사항이 증가하는 경향을 보입니다.

### 5.2.3 AI 노출 범위(Coverage) 사용 빈도·성공률 지표 설계

MCP 인벤토리와 분류 체계가 마련되었다면, 이를 기반으로 AI 활용 성과를 측정할 세 가지 핵심지표를 설계해야 합니다. 이 지표들은 AI가 조직에 미치는 영향력을 다각도로 분석하는 기준이 됩니다.

지표	정의	측정 방법 예시
- Al 노출 범위 (Coverage)	조직의 전체 업무 도메인 중,	(MCP가 구현된 업무 도메인
	AI가 MCP를 통해 접근 가능	수 / 전체 업무 도메인 수) *
	한 업무 도메인의 비율	100
사용 빈도 (Usage	특정 기간 동안 각 MCP 도구	중앙 로그 관리 시스템(ELK
Frequency)	가 AI 에이전트에 의해 호출된	등)을 통해 API Gateway의
	총 횟수	호출 로그를 집계하여 도구별/
		일별/월별 호출 수 분석
성공률 (Success Rate)	MCP 도구 호출 시도가 기술적	(성공한 호출 수 / 전체 호출
	오류 없이 성공적으로 완료된	시도 수) * 100. API 응답 코
	비율	드(예: HTTP 2xx)를 기준으
		로 산정



이러한 정량적 측정 체계는 AI 도입의 현주소를 명확히 보여주는 나침반 역할을 합니다. 이제이 나침반을 활용하여, 조직이 AI 여정의 어느 단계에 있는지 진단하고 다음 목적지로 나아갈 항로를 그릴 수 있는 구체적인 성숙도 모델을 정의할 차례입니다.

### 5.3 MCP 기반 AI 성숙도 모델

조직의 AI 도입 수준을 단순히 '도입했다' 또는 '안했다'의 이분법으로 나누는 것은 무의미합니다. AI의 가치는 활용의 깊이와 범위에 따라 달라지기 때문입니다. 따라서 MCP(Model Context Protocol)의 활용 수준을 기준으로 조직의 AI 성숙도를 체계적인 단계로 구분하는 모델이 필요합니다.

본 섹션에서는 0단계부터 4단계까지, 총 다섯 단계로 구성된 AI 성숙도 모델을 제시합니다. 각 단계는 MCP의 구현 범위, 시스템 통합 수준, 운영 내재화 정도를 기준으로 정의됩니다. 이 모델은 기업이 현재 자신의 위치를 객관적으로 진단하고, 다음 단계로 나아가기 위한 명확한 방향성을 설정하는 데 실질적인 도움을 줄 것입니다.

### 5.3.1 Level 0: 실험 단계 - PoC 챗봇·파일 업로드 수준

특징: 이 단계는 특정 비즈니스 문제 해결을 위해 매우 제한된 범위에서 AI 기술의 가능성을 탐색하는 개념 증명(Proof of Concept, PoC) 수준입니다. AI 기능은 대부분 전사 시스템과 통합되지 않은 독립적인 환경에서 운영되며, 특정 부서의 특정 업무에 한정되어 실험적으로 사용됩니다.

사례 분석: LLM 기반 대화형 챗봇 상담 시스템 구축 사업이 대표적인 예입니다. 이 프로젝트는 내부 문서(hwp, pdf 등)를 기반으로 지식 베이스를 구축하고, RAG(Retrieval-Augmented Generation) 기술을 활용하여 사용자의 질문에 답변하는 챗봇을 만드는 것을 목표로 합니다. 이는 전사적 데이터 연동 없이 특정 상담 업무에 대한 AI 도입의 기술적 타당성과 효과를 검증하는 과정으로, Level 0의 전형적인 특징을 보여줍니다. 이 단계에서는 MCP가 공식적으로 정의되지 않았을 수 있으며, 파일 업로드나 간단한 API 호출 형태를 띱니다.

### 5.3.2 Level 1~2: 일부 도메인 MCP화 - 특정 업무 도메인 자동화 단계

특징: 실험 단계를 넘어, AI의 가치를 실질적인 업무 효율화로 연결하는 단계입니다. 인사, 총무, 고객 서비스 등 특정 업무 도메인의 반복적인 작업을 자동화하기 위해 기존 레거시 시스템이나 내



부 데이터베이스의 기능들을 표준화된 MCP로 노출하여 AI와 본격적으로 연동하기 시작합니다. AI는 더 이상 독립된 챗봇이 아니라, 기간계 시스템의 데이터를 조회하고 특정 기능을 실행하는 주체로 기능합니다.

사례 분석: '인사관리 데이터 정비 및 기능 개선(FUN-010)' 요구사항은 이 단계의 좋은 예시 입니다. 예를 들어, "참여연구진 전원의 재직증명서와 연구실적증명서를 일괄 발급"하는 기능을 hr. issueCertificateBatch라는 MCP 도구로 구현하는 시나리오를 상상할 수 있습니다. AI 에 이전트가 이 MCP를 호출하면, 인사 시스템은 여러 명의 증명서를 자동으로 생성하고 출력하여 기존의 수작업을 대체합니다. 이처럼 명확한 비즈니스 가치를 지닌 특정 업무들이 MCP를 통해 자동화되는 것이 Level 1~2의 핵심입니다.

### 5.3.3 Level 3~4: 전사 MCP 포털·AI 운영 내재화 단계

특징: 개별 업무 도메인에서의 성공적인 MCP 연동 경험을 바탕으로, 이를 전사적으로 확산하고 체계적으로 관리하는 플랫폼화 단계입니다. 수십, 수백 개로 늘어난 MCP 서버를 중앙에서 관리, 배포, 모니터링하고 검색할 수 있는 전사 MCP 포털이 구축됩니다. 이 단계에서 AI 기능의 개발부터 배포, 운영, 보안에 이르는 전 과정(AI Ops 또는 DevSecOps)이 내재화되고 자동화됩니다.

근거: 조직이 Level 3에 도달하면 MCP의 수는 수백 개로 증가하여 거대한 소프트웨어 제공 및 운영상의 과제를 야기합니다. 이러한 복잡성으로 인해, DevSecOps 체계와 중앙 집중식 로깅 및 모니터링 시스템( MSAP Observability)은 더 이상 '좋은 시도'가 아닌, AI를 대규모로 관리하기 위한 '필수 전제조건'이 됩니다. 이는 단순히 AI 기능을 개발하는 것을 넘어, 안정적이고 확장가능한 AI 서비스를 전사적으로 운영할 수 있는 역량을 확보했음을 의미하며, AI가 조직의 핵심 인 프라로 자리 잡았음을 증명합니다.

이처럼 MCP 기반 성숙도 모델은 조직의 AI 여정을 명확한 단계로 제시합니다. 그러나 이 모델을 실제 조직에 효과적으로 적용하기 위해서는 각 조직의 특성에 맞는 구체적인 진단 기준과 실현 가능한 목표를 담은 로드맵 수립 방안이 필요합니다.

### 5.4 조직별 성숙도 진단 기준과 로드맵

AI 성숙도 모델은 이론적 청사진에 불과합니다. 실제 가치는 이 청사진을 조직의 고유한 비즈니스 환경, 규제 제약, 기술 부채라는 현실에 맞춰 실행 가능한 건축 설계도로 변환하는 데서 나옵니다.



이 섹션에서는 바로 그 변환 과정을 다룹니다.

모든 조직에 동일한 AI 성숙도 기준을 일괄적으로 적용하는 것은 비효율적입니다. 산업의 특성, 규제 환경, 기존 IT 인프라 수준을 고려한 맞춤형 진단 기준과 실현 가능한 로드맵 수립이 무엇보다 중요합니다. 공공, 금융, 서비스 등 주요 분야별 특성을 반영한 진단 항목 예시를 통해 조직의 현 위치를 진단하고, 이를 바탕으로 단기 및 중장기 목표를 설정하는 로드맵 수립 방안을 제시하여 이론적 모델을 실질적인 실행 계획으로 전환하는 과정을 안내합니다.

### 5.4.1 공공기관(민원·행정·통계)의 성숙도 진단 항목

공공기관은 대국민 서비스, 내부 행정 효율화, 재난 대응, 방대한 지식 자산 관리 등 고유한 업무특성을 가집니다. 이러한 특성을 반영한 MCP 기반 AI 성숙도 진단 항목은 다음과 같습니다.

진단 영역	진단 항목 (예시)	성숙도 단계 (0-4)
민원 서비스	환급금(지원금) 신청/조회 기	
	능이 MCP 도구로 제공되는	
	가?	
내부 행정	전자결재 시스템의 문서 반려/	
	승인 프로세스가 AI 에이전트	
	를 통해 호출 가능한가?	
재난 대응	GIS 기능(지도제어, 레이어 관	
	리)이 MCP로 노출되어 타 시	
	스템과 연동되는가?	
지식 관리	기관 내/외부 지식(법령, 매뉴	
	얼) 검색이 RAG와 결합된	
	MCP 도구로 구현되었는가?	

### 5.4.2 금융·서비스·제조 분야의 성숙도 진단 항목

금융, 서비스, 제조 분야는 신속한 고객 응대, 거래 처리의 정확성, 생산 데이터 분석 등이 핵심 경쟁력입니다. 각 분야의 특성을 반영한 진단 항목은 다음과 같이 구성할 수 있습니다.



진단 영역	진단 항목 (예시)	성숙도 단계 (0-4)
서비스 (예약/판매)	온라인 예약/예매/환불 기능이	
	MCP 도구로 구현되어 있는	
	가?	
금융 (결제/정산)	다양한 온라인 결제 수단(신용	
	카드, 계좌이체) 연동이 MCP	
	화 되어 있는가?	
공통 (고객관리)	CMS를 통한 고객 정보 조회	
	및 관리가 AI 챗봇을 통해 가능	
	한가?	
제조 (생산/재고)	생산 계획 데이터 조회가 MCP	
	리소스로 제공되는가?	

### 5.4.3 1년·3년 단위 MCP 로드맵 수립 예시

성숙도 진단 결과를 바탕으로, 조직은 구체적인 실행 계획인 로드맵을 수립해야 합니다. 아래는 가상의 공공기관을 위한 1년, 3년 단위의 MCP 로드맵 예시입니다.

- 1년 로드맵 (단기 목표: 파일럿 성공 및 기반 마련)
  - 목표: 성숙도 Level 1 달성
  - 주요 과제:
    - \* '대국민 챗봇 상담 시스템'을 파일럿 프로젝트로 선정하고 RAG 기반 지식 검색 및 자주 묻는 질문 자동응답 MCP 도구 개발.
    - \* MCP 개발 및 운영을 위한 최소한의 보안 가이드라인 초안 수립.
  - 예상 도전 과제:
    - \* 내부 데이터(HWP, PDF 등 비정형 문서)의 정제 및 AI 학습 데이터로의 전환 어려움.
  - 성공 전제 조건:



- \* 프로젝트 목표에 대한 경영진의 명확하고 지속적인 후원.
- \* 파일럿 도메인의 현업 부서(예: 민원 상담팀)의 적극적인 참여와 피드백.
- 3년 로드맵 (중장기 목표: 전사 확산 및 플랫폼 내재화)
  - 목표: 성숙도 Level 3 달성
  - 주요 과제:
    - \* 파일럿 프로젝트 성공을 바탕으로 인사, 재무 등 핵심 내부 행정 도메인으로 MCP 적용 확대.
    - \* 전사 MCP 포털 구축 프로젝트를 착수하여 MCP 서버의 중앙 관리, 검색, 모니터 링 기능 구현.
    - \* Al 기능의 안정적 운영을 위해 선진 사례에서 요구하는 DevSecOps 체계 및 중 앙 집중식 로깅 및 모니터링 시스템(Observability)을 지원하는 플랫폼 구축. Al 운영 전담팀 구성.
  - 예상 도전 과제:
    - \* 부서별로 고립된 데이터 사일로(silo)를 해소하고 데이터를 연동하는 과정의 복잡성.
    - \* 전사 데이터 거버넌스 및 표준화 정책 수립에 대한 부서 간 이견 조율.
  - 성공 전제 조건:
    - \* 전사 AI 전략을 총괄하는 중앙 거버넌스 조직(AICC, AI Center of Excellence 등)의 확립.
    - \* 단기 성과를 넘어 장기적인 플랫폼 구축을 위한 전사적 투자 결정 및 예산 확보.

### 제6장 기존 시스템의 지능형 MCP화 전략 - MSA 관점 레 거시 재설계

기업이 AI 시대의 경쟁력을 확보하기 위해 풀어야 할 가장 큰 숙제 중 하나는 지난 수십 년간 축적해 온 막대한 IT 자산, 즉 레거시 시스템을 어떻게 활용할 것인가입니다. 이 시스템들은 현재 비즈



니스의 근간을 이루고 있지만, 동시에 기술적 부채와 경직된 구조로 인해 새로운 AI 기술의 신속한 도입을 저해하는 양날의 검과 같습니다.

이러한 상황에서 기존 시스템을 전면 교체(Rip-and-Replace)하는 방식은 막대한 비용과 예측 불가능한 리스크를 동반하기에 현실적인 대안이 되기 어렵습니다. 보다 현명하고 실용적인 접근법은 Model Context Protocol(MCP)을 활용하여 기존 시스템을 점진적으로 '진화'시키는 것입니다. MCP는 LLM(거대 언어 모델)과 기업의 다양한 외부 시스템을 연결하는 개방형 표준 프로토콜로서, 레거시 시스템의 안정적인 기능은 유지하면서 AI가 이해하고 활용할 수 있는 새로운인터페이스를 제공하는 교두보 역할을 합니다.

본 장에서는 MCP를 중심으로 레거시 시스템을 '교체'가 아닌 '진화'시키는 구체적인 전략과 기술적 패턴, 그리고 장기적인 확장성을 보장하기 위한 마이크로서비스 아키텍처(MSA) 기반의 설계 원칙을 심도 있게 다룰 것입니다. 이 전략은 Al 모델이라는 강력한 '두뇌'와 실제 업무 프로세스라는 '신경망'을 연결하여, 기업의 IT 자산을 진정한 의미의 지능형 자산으로 전환하는 핵심 경로를 제시합니다.

### 6.1 '교체'가 아닌 '진화' 전략

AI 시대에 레거시 시스템을 현대화하는 과제는 '전면 교체'가 아닌 '점진적 진화'의 관점에서 접근해야 합니다. 예를 들어, 구축된 지 7년 이상 경과하여 사소한 기능 변경에도 많은 테스트와 검증이 필요한 시스템과 같은 경우, 전면 교체는 비즈니스 연속성에 심각한 위협이 될 수 있습니다. 진화적 접근법은 기존 시스템이 제공하는 안정성과 핵심 비즈니스 로직의 가치를 인정하고, 이를 유지하면서 AI와 같은 새로운 기술을 유연하게 접목시키는 전략적 선택입니다. 이는 위험을 최소화하며 변화에 점진적으로 적응하고, 투자를 보호하며 AI 시대의 경쟁력을 확보하는 가장 현실적인 경로입니다.

### 6.1.1 모놀리식·3-Tier 레거시 시스템 구조 분석

대부분의 기업 환경에 존재하는 레거시 시스템은 모놀리식(Monolithic) 또는 3-Tier 아키텍처로 구성되어 있습니다. 이러한 구조는 모든 기능이 하나의 거대한 단위로 긴밀하게 결합되어 있어, 작은 변경 사항 하나가 시스템 전체에 예기치 않은 영향을 미칠 수 있습니다.



- '소방방제시스템'의 경우, 상용 DBMS에 대한 높은 종속성으로 인해 GIS와 같은 핵심 기능의 유연한 확장에 어려움을 겪고 있습니다. 특정 기술에 묶여 있어 새로운 데이터 소스를 연동하거나 기능을 개선하는 데 제약이 따릅니다.
- '통합예약발권시스템'역시 사소한 기능을 추가하거나 변경할 때마다 시스템 전반에 걸친 광범위한 테스트와 검증 과정이 필요하여 신속한 업데이트가 불가능한 구조적 한계를 보여 줍니다.

이처럼 기능 간 결합도가 높은 구조는 독립적인 개발과 배포를 어렵게 만들어, 빠르게 변화하는 AI 기술을 유연하게 적용하는 데 큰 걸림돌이 됩니다. 새로운 AI 모델을 특정 업무에 적용하고 자 해도, 시스템 전체를 수정해야 하는 부담 때문에 시도조차 하기 어려운 것이 현실입니다.

### 6.1.2 "버리는 것이 아니라 감싸되. 다시 설계해서 감싸는" 방식의 전환

레거시 시스템을 진화시키는 핵심 전략은 '스트랭글러 피그(Strangler Fig)' 또는 '랩(Wrap)' 패턴을 적용하는 것입니다. 이는 기존 시스템의 코드를 직접 수정하는 대신, 그 외부를 새로운 인터페이스로 감싸 점진적으로 대체해 나가는 방식입니다.

여기서 중요한 점은 단순히 기존의 기술 중심적인 API를 그대로 노출하는 것이 아니라, AI가 이해하고 활용하기 쉽도록 비즈니스 도메인 관점에서 다시 설계하여 감싸야 한다는 것입니다. 예를 들어, 내부적으로 복잡한 파라미터를 요구하는 함수가 있더라도, AI 에이전트에게는 '고객 등급별 할인율 조회'와 같이 명확한 비즈니스 목적을 가진 인터페이스로 제공해야 합니다.

이 방식은 다음과 같은 명확한 장점을 가집니다.

- 위험 최소화: 기존 시스템의 핵심 로직은 변경 없이 그대로 유지되므로 안정성과 비즈니스 연속성을 보장할 수 있습니다.
- 점진적 현대화: 전체 시스템을 한 번에 바꾸는 대신, 우선순위가 높은 기능부터 하나씩 새로 운 인터페이스로 감싸며 점진적으로 시스템을 현대화할 수 있습니다.
- AI 친화적 설계: AI 모델이 쉽게 호출하고 그 결과를 해석할 수 있는 능력 지향적(Capability-oriented) 인터페이스를 제공하여 AI 통합의 복잡성을 크게 낮춥니다.



### 6.1.3 기존 비즈니스 로직을 AI 추론 호출 대상으로 승격시키는 방법

'랩' 패턴의 최종 목표는 레거시 시스템 내부에 캡슐화되어 있던 비즈니스 로직을 LLM이 직접 호출하고 활용할 수 있는 독립된 'MCP 도구(Tool)'로 승격시키는 것입니다.

MCP 명세에 따르면 '도구'란 AI 모델이 특정 작업을 수행하기 위해 호출하는 함수(Function)를 의미합니다. 기존 비즈니스 로직을 MCP 도구로 승격시키는 기술적 과정은 다음과 같습니다.

- 1. API화: 레거시 시스템의 특정 기능(예: 재고 조회 로직)을 호출할 수 있는 내부 API(REST, gRPC 등)를 만듭니다.
- 2. MCP 서버 노출: 이 API를 호출하는 MCP 서버를 구현합니다. MCP 서버는 LLM으로부터 tools/call 요청을 받아 해당 API를 실행하고 결과를 반환하는 어댑터 역할을 합니다.
- 3. 도구로 승격: LLM은 이제 '재고 조회'라는 명확한 이름과 inputSchema를 가진 신뢰할 수 있는 도구를 직접 호출할 수 있게 됩니다. LLM은 더 이상 복잡한 내부 구현을 알 필요 없이, 비즈니스 목적에 따라 필요한 도구를 동적으로 선택하고 실행할 수 있습니다.

이러한 진화 전략을 성공적으로 구현하기 위해서는 기업이 보유한 자산의 유형, 즉 API, 배치 작업, DB 쿼리 등의 특성에 따라 각기 다른 MCP 노출 패턴을 적용해야 합니다.

### 6.2 기존 자산을 MCP로 노출하는 패턴

기업 내부에는 REST API, 주기적으로 실행되는 배치 작업, 데이터베이스 쿼리 등 수많은 형태의 IT 자산이 존재합니다. 이러한 자산들을 AI가 일관된 방식으로 활용하기 위해서는 표준화된 인터 페이스가 필수적이며, MCP가 바로 그 역할을 수행합니다. MCP는 M개의 AI 모델과 N개의 외부 도구 간의 통합 복잡성이 기하급수적으로 증가하는 'M x N 통합 문제'를 해결하는 표준 프로토콜 입니다. 각 자산의 기술적 특성에 맞는 최적의 MCP화 패턴을 적용함으로써, 이기종 시스템의 장 벽을 허물고 모든 IT 자산을 AI가 활용 가능한 '도구'의 생태계로 통합할 수 있습니다.



### 6.2.1 REST gRPC SOAP API의 MCP 서버화 패턴

기업 내에 이미 존재하는 RESTful, SOAP, gRPC 등의 API는 MCP화를 위한 가장 좋은 출발점입니다. 이 패턴에서 MCP 서버는 기존 API와 AI 에이전트 사이에서 '어댑터(Adapter)' 또는 '프록시(Proxy)' 역할을 수행합니다.

프로세스는 다음과 같이 진행됩니다.

- 1. AI 에이전트가 특정 작업 수행을 위해 MCP 클라이언트를 통해 tools/call 요청을 보냅니다.
- 2. MCP 서버는 이 요청을 수신하여, 미리 정의된 매핑 규칙에 따라 기존 백엔드 API(예: REST API)에 대한 표준 HTTP 요청으로 변환하여 전송합니다.
- 3. 백엔드 API로부터 응답을 받으면, MCP 서버는 이를 AI 모델이 이해하기 쉬운 JSON 형식의 텍스트 콘텐츠로 가공합니다.
- 4. 가공된 결과를 AI 에이전트에게 최종적으로 반환합니다.

이 패턴을 통해 기존 API의 수정 없이도 레거시 자산을 AI 생태계에 빠르고 안전하게 편입시킬 수 있습니다.

### 6.2.2 배치·스케줄러·RPA 기능의 MCP 도구 추상화

'참여연구진 증명서 일괄 발급', '지급의뢰서 일괄 출력'과 같이 실시간 응답이 아닌 백그라운드에서 장시간 수행되는 배치(Batch) 및 스케줄러 작업 역시 MCP 도구로 추상화할 수 있습니다. 이러한 작업들은 복잡한 내부 실행 과정을 가지고 있지만, AI 에이전트에게는 단순한 도구 호출로 캡슐화하여 제공할 수 있습니다.

예를 들어, LLM이 start\_batch\_job(job\_name='issue\_all\_certificates')와 같은 도 구를 호출하면, MCP 서버는 이를 받아 내부 스케줄링 시스템에 해당 배치 작업을 트리거하는 요 청을 보냅니다. 작업의 시작 여부나 작업 ID를 즉시 반환하고, 실제 작업은 비동기적으로 백그라 운드에서 수행됩니다. 이를 통해 AI 에이전트는 복잡한 일괄 처리 프로세스를 직접 제어할 필요 없이, 간단한 명령만으로 강력한 백엔드 기능을 활용할 수 있게 됩니다.



### 6.2.3 DB 조회·리포트·통계 쿼리의 MCP화 방법

'버스정보시스템(BIS) DB 연계'나 '내부결재 현황 조회'와 같이 데이터베이스에서 특정 정보를 조회하거나 리포트를 생성하는 업무는 AI 기반 의사결정 지원에 매우 중요한 자산입니다. 이 패턴은 복잡한 데이터베이스 스키마나 SQL을 AI 에이전트로부터 숨기고, 의미 있는 정보 요청을 처리하는 도구를 제공합니다.

가령, 사용자가 "지난 분기 재무 보고서를 보여줘"라고 자연어로 질의하면, LLM은 이를 get\_financial\_report(start\_date='2024-01-01', end\_date='2024-03-31') 형태의 도구 호출로 변환합니다. MCP 서버는 이 호출을 받아, 미리 정의되고 검증된 SQL 쿼리를 데이터 베이스에 실행합니다. 쿼리 실행 결과로 얻은 데이터를 가공하여 AI 에이전트에게 반환하면, AI는 이를 바탕으로 사용자에게 표나 그래프 형태의 보고서를 생성해 줄 수 있습니다. 이 방식은 데이터 접근을 표준화하고 보안을 강화하면서, AI가 기업의 핵심 데이터에 안전하게 접근하여 인사이트를 도출할 수 있도록 합니다.

이처럼 다양한 패턴을 통해 기존 자산을 기술적으로 노출하는 것을 넘어, 장기적인 확장성과 유지보수성을 확보하기 위해서는 마이크로서비스 아키텍처(MSA)의 설계 원칙을 MCP 설계에 적 용하는 것이 필수적입니다.

### 6.3 MSA(Microservice Architecture) 관점 MCP 설계 원칙

MCP 서버들을 단순히 기능별로 무분별하게 구현할 경우, 서비스 간의 의존성이 복잡하게 얽힌 또 다른 형태의 기술 부채인 '분산 모놀리식(Distributed Monolith)'으로 전략할 위험이 있습니다. 이는 각 서비스가 독립적으로 배포되거나 확장될 수 없어, 마이크로서비스의 장점을 전혀 살리지 못하는 결과를 초래합니다. 진정한 민첩성과 회복탄력성을 확보하기 위해서는 마이크로서비스 아키텍처(MSA)의 핵심 원칙, 즉 '높은 응집도(High Cohesion)'와 '느슨한 결합(Loose Coupling)'을 MCP 설계에 반드시 적용해야 합니다. 이는 각 MCP 서버가 명확한 책임 단위를 가지며, 다른 서버에 대한 의존성을 최소화하여 독립적인 진화가 가능하도록 만드는 핵심 철학입니다.



### 6.3.1 DDD 관점에서 기존 API를 MCP로 재설계하는 방법론

기술 중심의 레거시 API를 비즈니스 중심의 MCP 도구로 성공적으로 재설계하기 위해서는 도메인 주도 설계(DDD, Domain-Driven Design) 방법론이 매우 효과적입니다. DDD의 핵심은 현업 전문가와 개발자가 '유비쿼터스 언어(Ubiquitous Language)'라는 공통의 언어를 사용하여 비즈니스 도메인을 함께 이해하고 모델링하는 것입니다.

이 과정에서 '바운디드 컨텍스트(Bounded Context)'라는 논리적 경계를 식별하게 되는데, 이는 특정 비즈니스 도메인이 책임지는 명확한 범위(예: 인사, 회계, 주문)를 의미합니다. 기술적 관점이 아닌 비즈니스 관점에서 식별된 이 경계를 기준으로 MCP 서버의 책임 단위를 설계하면, 자연스럽게 비즈니스와 기술 아키텍처가 일치하게 됩니다. 예를 들어, getUserInfo, getEmpDetails 등 기술적으로 파편화된 API들을 '인사관리'라는 바운디드 컨텍스트 안에서 '직원 정보 조회'라는 일관된 MCP 도구로 재설계할 수 있습니다.

#### 6.3.2 마이크로서비스 경계와 MCP 툴 단위 정렬

잘 설계된 마이크로서비스의 경계는 앞서 DDD를 통해 식별한 바운디드 컨텍스트와 일치해야 합니다. 이 원칙을 MCP 설계에도 그대로 적용하여, 마이크로서비스의 책임 범위와 MCP 서버가 제공하는 도구의 범위를 정렬하는 것이 중요합니다.

예를 들어, '인사관리'라는 바운디드 컨텍스트가 하나의 마이크로서비스로 구현되었다면, 관련 MCP 서버 역시 '증명서발급', '채용지원자조회', '휴가신청' 등 오직 인사관리 도메인에 속하는 도구들로만 구성되어야 합니다. '주문처리'나 '결제'와 관련된 도구가 이 서버에 포함되어서는 안 됩니다. 이러한 명확한 책임 분리는 다음과 같은 효과를 가져옵니다.

- 독립적 배포: 인사관리 기능의 변경이 다른 시스템에 영향을 주지 않고 독립적으로 배포될 수 있습니다.
- 자율적 팀 운영: 인사관리 마이크로서비스와 MCP 서버를 담당하는 팀이 다른 팀과의 불필 요한 조율 없이 자율적으로 개발하고 운영할 수 있습니다.
- 선택적 확장: 채용 시즌에 지원자 조회 요청이 급증할 경우, 다른 서비스에 영향을 주지 않고 인사관리 MCP 서버만 선택적으로 확장(scale-out)할 수 있습니다.



### 6.3.3 재사용성과 확장성을 고려한 MCP 스키마 설계

MCP 도구의 재사용성과 확장성은 inputSchema를 어떻게 설계하는지에 따라 크게 좌우됩니다. 스키마를 설계할 때는 특정 AI 에이전트나 단일 사용 사례에 종속되지 않고, 여러 컨텍스트에서 재 사용될 수 있도록 일반화하는 것이 핵심 원칙입니다.

예를 들어, '민원 상태 조회' 도구의 inputSchema를 재사용성과 확장성을 고려하여 다음과 같이 설계할 수 있습니다.

```
{
  "type": "object",
  "properties": {
    "complaint_id": { "type": "string", "description": "특정 민원 번호" },
    "complaint_type": { "type": "string", "description": "조회할 민원 유형 (예: 환경, 교통)" },
    "start_date": { "type": "string", "format": "date", "description": "조회 시작일" },
    "end_date": { "type": "string", "format": "date", "description": "조회 종료일" }
  },
   "required": []
}
```

이처럼 다양한 조회 시나리오에 대응할 수 있도록 옵션 파라미터를 포함하면, '특정 민원 조회', '특정 유형의 민원 목록 조회', '특정 기간 동안의 민원 조회' 등 여러 AI 에이전트가 다양한 맥락에서 이 도구를 재사용할 수 있게 되어 시스템 전체의 효율성이 극대화됩니다.

이와 같이 MSA 원칙에 따라 잘 설계된 MCP 서버들을 실제 운영 환경에서 안정적이고 안전하게 관리하기 위해서는, 다음으로 운영, 보안, 그리고 시스템의 상태를 깊이 있게 파악할 수 있는 관측성(Observability)에 대한 설계가 반드시 뒤따라야 합니다.

### 6.4 MCP 운영·보안·관측 설계

수많은 MCP 서버가 분산된 환경에서 AI 에이전트에 의해 동적으로 호출되는 현대적 아키텍처에서는, 전통적인 시스템 모니터링 및 관리 방식만으로는 안정성과 보안을 보장하기 어렵습니다. 예측 불가능한 호출 패턴과 서비스 간의 복잡한 상호작용을 효과적으로 통제하고 문제를 신속하게 해결하기 위해서는, 새로운 차원의 운영, 보안, 그리고 관측성(Observability) 설계가 필수적입니



다. 이는 단순히 시스템을 지켜보는 것을 넘어, 시스템 내부 상태에 대해 질문하고 답을 얻을 수 있는 능력을 확보하는 것을 의미합니다.

### 6.4.1 RBAC 기반 권한·역할 모델 설계

MCP 환경의 핵심 보안 요구사항은 '누가 무엇을 할 수 있는가'를 정밀하게 통제하는 것입니다. 이를 위해 역할 기반 접근 제어(RBAC, Role-Based Access Control) 모델을 설계해야 합니다. 이는 AI 에이전트나 사용자에게 '역할(Role)'을 부여하고, 각 역할이 호출할 수 있는 MCP 도구에 대한 권한을 정의하는 방식입니다.

구체적인 시나리오는 다음과 같습니다.

- '인사담당자' 역할을 가진 AI 에이전트는 '인사정보조회', '증명서발급' 도구를 호출할 수 있는 권한을 가집니다.
- '재무담당' 역할을 가진 AI 에이전트는 '지급의뢰서생성' 도구는 호출할 수 있지만, '인사정 보조회' 도구는 호출할 수 없습니다.

이를 기술적으로 구현하기 위해 MCP 명세는 OAuth 2.1과 같은 표준 프로토콜 활용을 강조합니다. 중앙 인증 서버(Authorization Server)를 통해 인증을 수행하고, 각 AI 에이전트에게는 자신이 수행할 작업에 필요한 최소한의 권한만을 담은 범위가 지정된 토큰(Scoped Token)을 발급합니다. MCP 서버는 도구 호출 요청을 받을 때마다 이 토큰을 검증하여, 요청자가 해당 작업을 수행할 정당한 권한을 가졌는지 확인한 후에만 요청을 처리해야 합니다.

### 6.4.2 호출 로그·추적·감사를 포함한 Observability 전략

신뢰할 수 있는 MCP 생태계를 구축하기 위해서는 모든 도구 호출에 대한 가시성을 확보하는 것이 중요합니다. 이를 위한 핵심 전략은 로깅, 추적, 감사입니다.

• 로깅 및 감사(Logging & Audit): 모든 MCP 도구 호출에 대해 '누가, 언제, 무엇을, 어떤 파라미터로 호출했고, 그 결과는 어떠했는가'를 명확히 알 수 있도록 구조화된 로그를 생성하고 중앙 집중식으로 수집해야 합니다. 이는 장애 발생 시 원인을 분석하는 첫 단서가 되며, 보안 감사 및 규제 준수 요구사항을 충족하는 필수 요소입니다.



• 분산 추적(Distributed Tracing): AI 에이전트의 최초 요청부터 시작하여 API 게이트웨이, 여러 MCP 서버, 그리고 최종적으로 레거시 시스템에 도달하기까지의 전체 호출 흐름을 하나의 트랜잭션으로 묶어 추적하고 시각화합니다. 이를 통해 전체 서비스 중 어느 구간에서 병목이 발생하는지, 어떤 서비스의 실패가 다른 서비스에 영향을 미쳤는지를 직관적으로 파악하여 신속한 문제 해결을 가능하게 합니다.

6.4.3 쿠버네티스 기반 Auto-scaling·Service Discovery·Resilience 확보 쿠버네티스(Kubernetes) 기반의 클라우드 네이티브 환경은 변동성이 큰 MCP 서버 운영에 최적의 솔루션을 제공합니다.

- 자동 확장(Auto-scaling): 특정 민원 신청(예: 환급금 신청)이 특정 시간대에 폭주할 경우, 쿠버네티스는 해당 MCP 서버를 실행하는 컨테이너의 수를 자동으로 늘려 부하를 분산시키고 안정적인 서비스를 보장합니다. 반대로 트래픽이 줄어들면 자원을 회수하여 비용 효율성을 높입니다.
- 서비스 디스커버리(Service Discovery): 새로운 MCP 서버가 배포되거나 기존 서버가 확장될 때, AI 에이전트가 해당 서버의 네트워크 위치(IP 주소, 포트)를 동적으로 찾아 연결할수 있도록 지원합니다. 이를 통해 수동 설정 없이도 유연한 서비스 확장이 가능합니다.
- 회복탄력성(Resilience): 서킷 브레이커(Circuit Breaker)와 같은 패턴을 적용하여, 특정 레거시 시스템의 응답이 지연되거나 장애가 발생했을 때 해당 시스템을 호출하는 MCP 서 버로의 요청을 일시적으로 차단합니다. 이는 일부 시스템의 장애가 전체 AI 서비스로 전파되는 것을 막아 시스템 전체의 안정성을 유지하는 핵심적인 방어 메커니즘입니다.

이와 같은 기술적, 운영적 전환은 단번에 이루어질 수 없으며, 전사적인 관점에서 체계적인 로드맵에 따라 위험을 관리하며 단계적으로 추진되어야 합니다.

### 6.5 단계별 MCP 전환 로드맵

전사적인 MCP 전환을 '빅뱅(Big-bang)' 방식으로 한 번에 시도하는 것은 높은 실패 위험을 안고 있습니다. 대신, 작고 가시적인 성공을 통해 기술의 효용성을 입증하고 전사적 공감대를 형성한 후



점진적으로 확대해 나가는 단계별 로드맵이 반드시 필요합니다. 이 접근법은 초기 단계에서 학습과 경험을 축적하고, 이를 바탕으로 더 복잡하고 중요한 영역으로 안정적으로 확장할 수 있는 기반을 마련해 줍니다.

### 6.5.1 3개월 내 구현 가능한 MCP 후보 도출

성공적인 첫걸음을 위해서는 단기간(예: 3개월) 내에 가시적인 성과를 낼 수 있는 '빠른 성공 (Quick-win)' 과제를 신중하게 선정해야 합니다. 후보 과제 선정 기준은 다음과 같습니다.

- 기술적 단순성: 외부 시스템 의존성이 적고, 기존 로직을 API로 노출하는 과정이 비교적 간단한 업무를 우선으로 합니다.
- 명확한 효과: 현업 부서 담당자들이 AI 도입으로 인한 업무 시간 단축이나 편의성 증대를 명확히 체감할 수 있어야 합니다.
- 낮은 리스크: 핵심적인 거래 처리나 민감 데이터를 다루는 업무보다는, 실패 시 영향이 적은 내부 지원 업무가 적합합니다.

이러한 기준에 부합하는 후보로는 '내부 규정 및 매뉴얼에 대한 질의응답 시스템'이나, 반복적인 '특정 양식 데이터 자동 입력'과 같은 업무를 고려할 수 있습니다.

### 6.5.2 핵심 도메인(민원·운영·개발 생산성 등) 우선 전환 전략

초기 PoC(Proof of Concept)가 성공적으로 완료되면, 비즈니스 가치와 파급 효과가 큰 핵심 도메인을 중심으로 MCP 전환을 확대해야 합니다. 다수의 공공기관 제안요청서에서 반복적으로 언급되는 다음과 같은 영역들이 우선 고려 대상이 될 수 있습니다.

- 대국민 민원 서비스: '환급금 신청/조회', '증명서 발급' 등 반복적이고 표준화된 민원 업무를 MCP화하여 24시간 응대 가능한 지능형 서비스를 제공하고, 상담원의 업무 부담을 경감시킵니다.
- 시스템 운영 자동화: '서버 상태 모니터링', '로그 분석을 통한 이상 징후 탐지', '장애 발생 시 초기 대응' 등의 운영 업무를 MCP 도구로 만들어 AI 에이전트가 수행하도록 함으로써, 운영 효율성과 안정성을 극대화합니다.



• 개발 생산성 향상: '코드 리뷰 요청', '테스트 환경 자동 구성', '배포 파이프라인 실행' 등의 개발 관련 작업을 MCP 도구로 제공하여, 개발자들이 반복적인 업무에서 벗어나 핵심적인 개발 활동에 집중할 수 있도록 지원합니다.

### 6.5.3 PoC에서 전사 MCP 카탈로그로 확장하는 조직·프로세스 변화

소규모 PoC를 넘어 전사적으로 재사용 가능한 MCP 자산을 체계적으로 축적하고 관리하기 위해서는 기술적 노력과 더불어 조직 및 프로세스의 변화가 반드시 수반되어야 합니다. 그 중심에는 '전사 MCP 카탈로그(또는 MCP 포털)' 구축이 있습니다. 이는 등록된 모든 MCP 서버와 도구의명세, 사용법, 담당자, 보안 등급 등을 관리하고, 개발자들이 쉽게 검색하여 재사용할 수 있도록 지원하는 중앙 저장소입니다.

이러한 MCP 카탈로그를 성공적으로 운영하기 위해서는 다음과 같은 조직적, 프로세스적 변화가 필수적입니다.

- 중앙 거버넌스 조직(CoE, Center of Excellence) 설립: 전사 MCP 전략을 수립하고, 개발 표준 및 보안 정책을 정의하며, 재사용성을 감독하는 중앙 전담 조직을 구성합니다.
- 표준화된 개발 및 관리 프로세스 수립: MCP 서버의 설계, 개발, 테스트, 배포, 버저닝에 대한 표준화된 가이드라인을 제공하여 품질의 일관성을 확보합니다.
- 성과 측정 및 재사용 촉진: MCP 도구의 재사용률, 호출 빈도 등을 측정하고, 재사용 성과 가 높은 팀이나 개인에게 인센티브를 제공하여 자발적인 참여와 공유 문화를 확산시킵니다.

결론적으로, 성공적인 MCP 전환은 단순히 기술을 도입하는 것을 넘어, Al와 협업하는 새로운 업무 방식을 조직 문화에 내재화하는 전사적인 혁신 과정입니다. 이는 본 장에서 제시한 '진화' 전략의 최종 목표이며, 기업의 모든 IT 자산을 실질적인 비즈니스 가치를 창출하는 지능형 자산으로 승격시키는 유일한 경로입니다.



# 제7장: MCP와 RAG로 내부 데이터를 AI 자산으로 전환하는 방법

이 장의 핵심 목표는 기업 내부에 축적된 방대한 비정형 데이터를 단순한 정보 저장소를 넘어, AI가 실질적인 가치를 창출하는 '지능형 자산'으로 전환하는 구체적인 기술 아키텍처와 전략을 제시하는 것입니다. 많은 기업이 생성형 AI의 잠재력을 인지하고 있지만, 그 힘을 내부 데이터와 안전하게 결합하여 실제 업무 프로세스를 혁신하는 단계에는 이르지 못하고 있습니다. 본 장에서는 Model Context Protocol(MCP)을 통한 '기능적 지능'과 Retrieval-Augmented Generation(RAG)을 통한 '맥락적 지능'의 결합이 어떻게 이러한 혁신을 가능하게 하는지 전략적 중요성을 강조하여 서술합니다. 이 두 기술을 체계적으로 통합함으로써, AI는 비로소 기업의 고유한 지식과 프로세스를 이해하고 실질적인 '행동'을 수행하는 신뢰할 수 있는 파트너로 거듭날 수 있습니다.

## 7.1 두 개의 뇌를 가진 AI: 행동을 담당하는 MCP, 지식을 담당하는 RAG

AI 시스템을 성공적으로 구축하기 위한 첫걸음은 시스템의 역할을 '기능(Function)'과 '맥락 (Context)'으로 명확히 분리하여 설계하는 것입니다. 이는 시스템 아키텍처의 유연성, 확장성, 그리고 무엇보다 신뢰도를 극대화하는 핵심 전략입니다. 이 개념을 숙련된 요리사와 레시피북의 관계에 비유할 수 있습니다.

요리사는 재료를 다듬고, 불을 조절하며, 요리를 완성하는 실질적인 '행동'을 수행합니다. 이는 외부 시스템과 연동하여 작업을 처리하는 MCP(기능적 지능)에 해당합니다. 반면, 레시피북은 어떤 요리를 어떤 순서로 만들어야 할지에 대한 정확한 '지식'과 '맥락'을 제공하며, 이는 RAG(맥락적 지능)의 역할과 같습니다. 요리사의 손기술(MCP)이 아무리 뛰어나도, 검증된 레시피(RAG) 없이는 일관된 최고의 결과물을 낼 수 없습니다. 이처럼 각 요소를 분리하면, 새로운 레시피를 추가하거나(RAG 업데이트) 요리사의 기술을 향상시키는(MCP 개선) 작업을 독립적으로 수행할 수 있어 전체 시스템이 훨씬 유연하고 강력해집니다.

즉 이 모델에서 MCP는 외부 시스템과 연동하여 데이터 조회, 프로세스 실행 등 실질적인 '행동'을 수행하는 기능적 지능을 담당합니다. 반면, RAG는 기업 내부의 방대한 문서와 데이터에 기



반하여 신뢰도 높은 '지식'을 제공하는 맥락적 지능의 역할을 맡습니다. 이처럼 두 지능을 분리함으로써, 우리는 복잡한 AI 시스템을 더 작고 관리하기 쉬운 단위로 나누고, 각 구성 요소의 역할을 명확히 하여 독립적인 개발과 개선을 가능하게 할 수 있습니다.

### 7.1.1 AI의 그럴듯한 거짓말, '할루시네이션'과 신뢰의 열쇠

최신 거대 언어 모델(LLM)은 방대한 지식을 바탕으로 인간과 유사한 대화를 나누는 놀라운 능력을 보여줍니다. 하지만 이 똑똑해 보이는 AI에게는 아직 넘지 못한 거대한 벽이 존재하는데, 바로일본의 AI 권위자 마츠오 유타카 교수가 언급한 '문맥 이해의 벽(依然として口る文脈理解の壁)'입니다. 이는 AI가 세상의 모든 지식을 학습했더라도, 특정 조직이나 상황이라는 '특수한 맥락' 완벽히 이해하지 못하는 본질적 한계를 의미합니다.

이 한계 때문에 발생하는 가장 치명적인 문제가 바로 '할루시네이션(Hallucination, 환각)' 현상입니다. 쉽게 말해 AI가 그럴듯한 거짓말을 지어내는 것입니다.

할루시네이션은 왜 발생할까요?

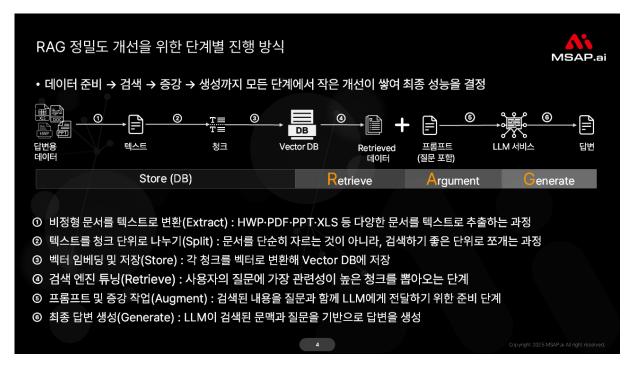
LLM의 작동 원리는 '진실'을 찾는 것이 아니라, 주어진 질문에 대해 '가장 그럴듯한 다음 단어'를 확률적으로 예측하여 문장을 완성하는 방식이기 때문입니다. 만약 기업 내부 정보처럼 LLM이 전혀 학습한 적 없는 '맥락'에 대해 질문을 받으면, LLM은 "모르겠습니다"라고 답하는 대신, 자신이 아는 온갖 일반 지식을 조합해 가장 자연스럽고 설득력 있게 들리는 답변을 창작해냅니다.

마치 "눈치 빠른 신입사원에게 우리 회사 내부 규정에 대해 물었더니, 잘 모름에도 불구하고 어디선가 들어본 다른 회사 이야기들을 섞어 자신감 있게 대답하는 상황"과 같습니다. 이 답변은 겉보기에는 매우 전문적이지만 실제로는 사실이 아니며, 기업 환경에서 이러한 오류는 잘못된 의사결정, 고객 신뢰 하락 등 심각한 비즈니스 리스크를 초래할 수 있습니다.

### 7.1.2 RAG 개념과 전형적인 아키텍처(인덱싱·검색·생성)

RAG(Retrieval-Augmented Generation)는 LLM에 외부 지식 데이터베이스를 연결하여 답변의 정확성과 신뢰도를 획기적으로 높이는 기술입니다. LLM이 가진 내부 지식의 한계를 극복하고 할루시네이션을 방지하기 위해, RAG는 답변 생성에 필요한 근거를 실시간으로 외부 데이터 소스에서 찾아와 LLM에 제공합니다. 다수의 공공 및 민간 프로젝트 제안요청서에서 RAG 도입을 필수 요건으로 명시하는 것은 이러한 신뢰성 확보의 중요성을 방증합니다.





[그림 5] RAG 아키텍처와 동작 방법

RAG의 전형적인 아키텍처는 다음과 같은 3단계로 구성됩니다.

- 1. 인덱싱 (Indexing): 이 단계에서는 기업 내부에 흩어져 있는 다양한 비정형 데이터(hwp, pdf, 이메일, 업무 로그 등)를 수집하고 정제합니다. 이후, 정제된 텍스트를 AI가 의미적으로 이해할 수 있는 숫자 형태의 벡터(Vector)로 변환(임베딩)하고, 이를 신속하게 검색할 수 있는 벡터 데이터베이스(Vector DB)에 저장합니다. 이 과정은 내부 지식을 AI가 활용할 수 있는 형태로 '자산화'하는 첫 단계입니다.
- 2. 검색 (Retrieval): 사용자의 질의가 입력되면, 시스템은 먼저 이 질의를 벡터로 변환합니다. 그 다음, 인덱싱된 벡터 데이터베이스에서 사용자의 질의 벡터와 의미적으로 가장 유사한 데이터 조각(Chunk)들을 찾아냅니다. 이 단계의 정확성이 RAG 전체 성능의 핵심을 좌우합니다.
- 3. 생성 (Generation): 마지막으로, 검색 단계에서 찾아낸 근거 데이터(Context)를 사용자의원본 질의와 함께 LLM에 전달합니다. LLM은 이 구체적인 근거를 바탕으로 최종 답변을 생성하게 됩니다. 이 방식을 통해 LLM은 막연한 추측이 아닌, 검증 가능한 내부 데이터에 기반한 정확하고 신뢰도 높은 답변을 제공할 수 있습니다.



### 7.1.3 MCP 도구로서 RAG 인덱싱·질의 기능을 노출하는 방식

앞서 설명한 기능과 맥락의 분리 원칙을 아키텍처로 구현하는 가장 효과적인 방법은 RAG의 검색기능을 MCP '도구(Tool)'로 추상화하여 노출하는 것입니다. 이는 RAG 인덱싱 파이프라인 전체를 하나의 표준화된 인터페이스 뒤로 캡슐화하는 전략입니다.

이 아키텍처에서 AI 에이전트(LLM)는 내부 지식이 필요할 때, 직접 RAG 시스템의 복잡한 구조를 이해할 필요 없이, 사전에 등록된 MCP 도구를 호출하기만 하면 됩니다. 예를 들어, 에이전트는 다음과 같은 간단한 명령어를 실행합니다.

search\_internal\_documents(query: "지난 분기 데이터베이스 연결 오류의 주요 원인과 해결책을 요약해줘.")

이 search\_internal\_documents 도구 호출이 MCP Host를 통해 전달되면, 내부적으로는 앞서 설명한 RAG의 검색(Retrieval) 파이프라인이 실행됩니다.

RAG 시스템은 벡터 데이터베이스에서 가장 관련성 높은 문서를 찾아내고, 그 내용을 다시 LLM 에이전트에게 구조화된 형태로 반환합니다.

이 방식은 RAG라는 복잡한 기술 스택을 표준화된 인터페이스 뒤로 숨김으로써, 여러 AI 에이전트가 일관된 방식으로 내부 지식 베이스에 접근할 수 있게 만듭니다.

'인사 규정 챗봇'과 '기술 장애 분석 에이전트'는 서로 다른 목적을 가졌지만, 동일한 search\_internal\_documents 도구를 호출하여 필요한 지식을 얻을 수 있습니다.

이러한 아키텍처적 분리는 단순히 기술적 우아함을 넘어 비즈니스 필수 요건입니다.

이를 통해 기업은 지식 소스(RAG)와 행동 역량(MCP)을 독립적으로, 그리고 병렬적으로 개선 하여 AI 기반 혁신의 속도를 가속화할 수 있습니다.

다음 절에서는 이 구조를 기반으로 기존 데이터를 AI가 실제로 활용할 수 있는 자산으로 만드는 구체적인 절차를 살펴보겠습니다.

### 7.2 기존 데이터를 AI가 활용할 수 있게 만드는 절차

여기는 AI 프로젝트가 혁신적인 비즈니스 자산이 되느냐, 아니면 값비싼 과학 실험으로 남느냐가 결정되는 중대한 갈림길입니다. 기업 내부에 다양한 형태로 축적된 데이터를 AI가 '이해하고 활용할 수 있는' 형태로 변환하는 체계적인 절차는 실제 업무 프로세스에 적용 가능한 AI를 구축하기



위한 필수 조건입니다. 이 과정에 적용되는 엔지니어링의 엄격함이 전체 AI 이니셔티브의 투자 대비 효과(ROI)를 직접적으로 결정합니다. 이 과정은 단순히 파일 형식을 바꾸는 데이터 변환을 넘어, 흩어져 있던 정보의 가치를 재발견하고 AI 시스템의 최종 성능을 결정짓는 핵심적인 데이터 엔지니어링 활동입니다.

#### 7.2.1 문서·메일·업무 로그의 정제·구조화·분류

'LLM 데이터 활용을 위한 데이터 정비 및 기능 개선'을 요구하는 것은, 원본 데이터의 품질이 Al의 성능과 직결되기 때문입니다. RAG 시스템이 데이터를 효과적으로 활용하기 위해서는 다음과 같은 전처리 과정을 반드시 거쳐야 합니다.

- 데이터 정제: AI의 이해를 방해하는 노이즈를 제거하는 과정입니다. 스캔된 문서 이미지에서 되는 보다 이미지에서 부모를 추출하기 위해 OCR(광학 문자 인식) 기술을 적용하고, 웹 페이지에서 수집한 데이터의 불필요한 HTML 태그를 제거하며, 다양한 형식의 특수문자를 정규화하는 작업이 포함됩니다. 깨끗한 데이터는 정확한 임베딩의 전제 조건입니다.
- 데이터 구조화: 정제된 비정형 텍스트를 의미 있는 단위(Chunking)로 분할하고, 각 데이터 조각에 풍부한 맥락을 부여하는 과정입니다. 단순히 텍스트를 자르는 것을 넘어, 각 데이터 에 출처(예: '2024년 1분기 재무 보고서'), 작성일, 작성 부서, 문서 버전 등의 메타데이터 를 포함시켜야 합니다. 이는 이후 검색 정확도를 높이는 데 결정적인 역할을 합니다.
- 데이터 분류: 데이터를 업무 도메인에 따라 체계적으로 분류하는 것은 검색 효율성을 높이는 데 매우 중요합니다. 예를 들어, 모든 문서를 하나의 인덱스에 저장하는 대신 '인사 규정', '재무 보고서', '기술 장애 보고' 등과 같이 분류된 인덱스를 구성하면, AI 에이전트는 특정 도메인에 국한된 검색을 수행하여 더 빠르고 정확한 결과를 얻을 수 있습니다.

#### 7.2.2 임베딩 전략·메타데이터 설계와 보안 태깅

임베딩 전략: 텍스트를 AI가 이해할 수 있는 벡터로 변환하는 '임베딩(Embedding)' 과정은 RAG의 검색 성능을 좌우하는 핵심 요소입니다. 짧은 질의응답 데이터, 긴 보고서, 소스 코드 등 문서의 종류와 길이에 따라 최적의 성능을 내는 임베딩 모델이 다르므로, 데이터 특성에 맞는 모델을 선택하는 전략이 필요합니다.



메타데이터 설계: 잘 설계된 메타데이터는 RAG 검색의 정확도를 극적으로 향상시킵니다. "문서 메타데이터를 활용하여 검색 범위를 축소하고 정확도를 향상"시키는 것이 많은 과업의 핵심 요구사항입니다. 예를 들어, 사용자가 "최신 연차 규정"을 질문했을 때, 메타데이터에 '문서 종류: 규정', '생성일: 2024-01-15', '소관 부서: 인사팀'과 같은 정보가 있다면, AI는 수많은 문서 중가장 관련성 높은 최신 인사 규정을 정확하게 필터링하여 찾아낼 수 있습니다.

보안 태깅: 개인정보와 기업 기밀을 다루는 엔터프라이즈 환경에서는 데이터 보안이 무엇보다 중요합니다. RAG 인덱싱 과정에서 데이터의 민감도에 따라 '개인정보', '대외비', '기밀' 등과 같은 보안 등급을 반드시 태깅해야 합니다. 이 태깅 정보는 이후 AI 에이전트가 데이터에 접근할 때 사용자의 권한을 확인하고 접근을 제어하는 보안 로직의 기반이 됩니다.

#### 7.2.3 하나의 인덱스를 여러 에이전트·업무에서 재사용하는 패턴

이 패턴은 중복되고 사일로화된 AI 솔루션을 만들려는 조직적 경향에 직접적으로 맞서는 방법입니다. 이는 기술 부채와 일관성 없는 사용자 경험의 흔한 원인이 됩니다. 잘 구축된 RAG 인덱스(지식 베이스)는 특정 챗봇이나 단일 업무에 종속되는 일회성 자산이 아닙니다. MCP 전략을 통해 이지식 베이스는 여러 AI 에이전트와 업무에서 반복적으로 재사용될 수 있는 전사적인 AI 자산으로 거듭납니다.

하나의 고품질 인덱스를 중앙에서 관리하고, 다양한 AI 에이전트가 MCP 도구를 통해 이 인덱스를 공유하여 활용하는 패턴은 다음과 같이 매우 효율적입니다.

RAG 인덱스 (지식 베이스)	활용하는 AI 에이전트/업무	주요 질의 내용
전사 규정 인덱스	1. 신입사원 온보딩 챗봇	"연차 사용 규정에 대해 알려
		줘."
	2. 감사팀 규정준수 검토 에이	"프로젝트 하도급 계약 시 필수
	전트	검토사항은?"
기술 장애 보고서 인덱스	1. SRE 장애 대응 지원 에이	"지난 분기 데이터베이스 연결
	전트	오류의 주요 원인과 해결책을
		요약해줘."



RAG 인덱스 (지식 베이스)	활용하는 AI 에이전트/업무	주요 질의 내용
	2. 고객 지원 챗봇	"00 서비스 접속 불가 시 사
		용자가 직접 확인할 사항은?"

이러한 재사용 패턴은 부서별로 유사한 지식 베이스를 중복해서 구축하고 관리하는 비효율을 방지합니다. 또한, 모든 AI 에이전트가 동일한 중앙 지식 소스를 참조함으로써 전사적으로 일관된 정보에 기반한 AI 응답 품질을 보장하는 최적의 아키텍처입니다. 그러나 데이터의 활용 범위가 넓어질수록, 데이터 보안 및 규제 준수의 중요성은 더욱 커지게 되며, 이는 다음 절에서 심도 있게 다룰 주제입니다.

#### 7.3 공공·엔터프라이즈 데이터와 규제·보안

이 섹션은 엔터프라이즈 AI의 근본적인 긴장 관계를 탐색합니다. 즉, LLM의 무한한 잠재력을 활용하면서도 개인정보 보호법과 같은 데이터 프라이버시 법률과 엄격한 내부 보안 정책이라는 경직되고 협상 불가능한 경계 내에서 운영해야 하는 과제입니다. 공공 및 엔터프라이즈 환경에서 AI를 성공적으로 도입하기 위해서는 기술적 성능만큼, 혹은 그 이상으로 중요한 것이 데이터 규제 준수와 철저한 보안 체계 구축입니다. AI가 내부의 민감한 개인정보와 기업 기밀을 다루게 되는 만큼, RAG 및 MCP 아키텍처를 설계하는 초기 단계부터 강력한 보안 원칙과 전략을 반드시 통합해야 합니다.

#### 7.3.1 개인정보·기밀 정보·규제 데이터의 취급 기준

개인정보 보호법 준수 및 '누출금지 대상정보'의 철저한 관리를 공통적으로 요구하고 있습니다. 민감 데이터를 취급하는 AI 시스템은 다음 원칙들을 반드시 준수해야 합니다.

• 비식별화 조치: RAG 인덱싱 과정에서 주민등록번호, 연락처, 주소 등 개인을 식별할 수 있는 정보는 원칙적으로 사전에 마스킹하거나 삭제하는 등 비식별화 처리를 거쳐야 합니다. LLM에 개인정보가 직접 노출되는 것을 원천적으로 차단하는 것이 가장 안전한 접근 방식입니다.



- 암호화 저장 및 전송: 벡터 데이터베이스에 저장되는 모든 민감 데이터와 메타데이터는 반드시 암호화되어야 합니다. 또한, MCP를 통해 AI 에이전트와 MCP 서버 간에 데이터가 전송되거나, LLM으로 근거 데이터가 전달되는 모든 통신 구간은 SSL/TLS와 같은 표준 프로토콜을 사용하여 암호화해야 합니다.
- 접근 기록 관리: 접근 기록은 최소 1년 이상 보관해야 하며, 특히 5만 명 이상의 개인정보 또는 고유식별정보를 처리하는 시스템의 경우 관련 규정에 따라 2년 이상 보관하여 감사 추 적성을 확보해야 합니다. 이를 통해 누가, 언제, 어떤 데이터에 접근했는지에 대한 모든 로 그를 상세히 기록하여 보안 사고 발생 시 원인 분석 및 책임 추적의 근거로 활용할 수 있어 야 합니다.

#### 7.3.2 중앙집중형 벡터DB vs 소스 시스템 실시간 조회(에이전트형 아키텍처)

기업 데이터를 AI에 연동하는 방식은 크게 두 가지 아키텍처로 나눌 수 있으며, 각각 장단점이 뚜렷합니다.

아키텍처 유형	설명	장점	단점
중앙집중형 벡터DB	모든 소스 시스템의 데	- 여러 소스의 정보를	
	이터를 사전에 ETL(추	통합하여 답변 생성이	
	출, 변환, 적재)하여 하	가능하고 검색 속도가	
	나의 대규모 벡터 데이	빠릅니다.	
	터베이스에 인덱싱하		
	고, RAG는 이 통합		
	DB만을 조회하는 방		
	식입니다.		
- 복잡한 쿼리에 대해	- 원본 데이터와 벡터		
일관된 성능을 보장합	DB 간 데이터 동기화		
니다.	지연 문제가 발생할 수		
	있습니다.		



 아키텍처 유형	서며	자저	 다저
이기택시 ㅠ엉 	설명 	장점 	단점 
- 대규모 DB를 구축			
하고 유지보수하는 데			
상당한 비용과 노력이			
필요합니다.			
소스 시스템 실시간 조	각 소스 시스템(데이터	- 데이터 동기화 없이	
회 (에이전트형)	베이스, 내부 API 등)	항상 최신 데이터에 접	
	을 직접 호출하는	근할 수 있습니다.	
	MCP 도구를 구현합		
	니다. AI 에이전트가		
	필요에 따라 실시간으		
	로 해당 도구를 호출하		
	여 데이터를 조회하는		
	'Mediated Access		
	Pattern'을 따릅니다.		
- 기존 시스템을 그대	- 각 소스 시스템에 대		
로 활용하므로 초기 구	한 실시간 조회로 인해		
축 비용을 절감할 수	응답 지연이 발생할 수		
있습니다.	있습니다.		
- AI 에이전트의 동시			
호출이 많아질 경우 각			
소스 시스템의 부하가			
증가할 수 있습니다.			

이 하이브리드 모델은 타협이 아니라 전략적 필수입니다. '전사 규정 인덱스'(7.2.3절)와 같이 정적이고 기초적인 지식의 경우, 중앙집중형 벡터DB는 고속 검색을 보장합니다. 반대로, '지난 주말 모구리야영장 예약률'(20230638109...pdf 참조)과 같이 동적이고 운영적인 데이터의 경우, 실시간 MCP 도구 호출이 정확성을 보장하는 유일한 방법입니다. 성공적인 아키텍처는 이 두 가



지 요구사항을 모두 수용해야 합니다.

#### 7.3.3 RAG·MCP를 활용한 보안·감사·접근제어 전략

AI의 강력한 기능을 안전하게 활용하기 위해서는 MCP와 RAG를 기반으로 한 다계층 보안 전략이 필수적입니다. MCP 아키텍처는 AI의 모든 외부 접근을 중앙에서 통제할 수 있는 구조를 제공하므로, 이를 활용하여 강력한 'Mandatory Security Controls'를 구현할 수 있습니다.

- 1. MCP Host를 통한 중앙 통제 (Security Broker): MCP Host는 모든 AI 데이터 트래픽에 대한 단일하고 강력하게 요새화된 검문소 역할을 합니다. 어떠한 요청도 사용자를 인증하고, 권한을 확인하며, 전체 트랜잭션을 기록하는 이 중앙 보안 브로커를 통과하지 않고는 내부 시스템이나 지식 베이스에 도달할 수 없습니다. 이를 통해 AI 모델이 임의의 시스템에 직접 접근하는 것을 원천적으로 차단하고 모든 상호작용을 통제된 채널로 유도합니다.
- 2. 역할 기반 접근 제어 (RBAC): 사용자의 역할(예: 인사팀, 재무팀, 개발팀)에 따라 접근 가능한 RAG 인덱스의 범위와 호출할 수 있는 MCP 도구를 엄격히 제한하는 정책을 수립하고 적용해야 합니다. 예를 들어, 인사팀 사용자는 '인사 규정 인덱스'에만 접근할 수 있고, '재무 시스템 조회' MCP 도구는 호출할 수 없도록 통제해야 합니다.
- 3. 호출 로그 기반 감사 및 모니터링: 모든 RAG 질의와 MCP 도구 호출 내역(사용자, 시간, 입력 파라미터, 반환 결과)을 상세히 로깅해야 합니다. 이 로그들을 중앙 Observability 플랫폼으로 전송하여 분석함으로써 비정상적인 접근 시도나 데이터 유출 징후를 실시간으로 탐지하고, 보안 감사 시 신뢰할 수 있는 증적 자료로 활용하는 체계를 구축해야 합니다.

이러한 보안 아키텍처는 Al라는 강력한 엔진에 안전벨트와 브레이크를 장착하는 것과 같습니다. 이는 Al의 도입을 가속화하면서도 기업의 가장 중요한 자산인 데이터를 보호하기 위한 필수 전제조건이며, 다음 절에서는 이를 바탕으로 한 구체적인 도입 워크플로를 살펴보겠습니다.

#### 7.4 MCP·RAG 통합 워크플로와 도입 단계

지금까지 논의된 RAG의 맥락적 지능, MCP의 기능적 지능, 그리고 이를 뒷받침하는 보안 전략을 통합하여, 실제 업무에 적용할 수 있는 구체적인 워크플로와 단계별 도입 방안을 제시합니다. 본



섹션의 목표는 이론적 논의를 넘어 실용적인 청사진을 제공함으로써, 기업의 IT 의사결정자들이 AI 도입의 막연한 구상 단계를 지나 실질적인 첫걸음을 내디딜 수 있도록 안내하는 것입니다.

#### 7.4.1 질의 → 검색 → MCP 툴 호출 → 결합 응답 워크플로

사용자의 복합적인 요청이 처리되는 과정을 구체적인 워크플로를 통해 살펴보겠습니다. 이 워크플로는 RAG를 통한 지식 검색과 MCP를 통한 실시간 데이터 조회가 결합되어 시너지를 창출하는 과정을 명확하게 보여줍니다.

워크플로: "지난 달 우리 부서의 클라우드 비용 상위 3개 서비스를 찾아보고, 비용 절감 가이드를 요약해줘."

- 1. [질의 분석]: LLM 에이전트가 사용자의 자연어 질의를 분석하여 두 가지 과업(① 비용 조회, ② 문서 검색 및 요약)으로 분해합니다.
- 2. [MCP 도구 호출]: 첫 번째 과업을 해결하기 위해, LLM은 등록된 MCP 도구 목록에서 '비용 조회' 기능에 가장 적합한 get\_cloud\_cost(team: "우리 부서", period: "지난 달") 도구를 선택하여 호출합니다.
- 3. [실시간 데이터 획득]: MCP 서버는 이 호출을 받아 내부 비용 관리 시스템을 실시간으로 조회하고, 비용 상위 3개 서비스 목록([서비스A, 서비스B, 서비스C])을 구조화된 데이터 (JSON)로 반환합니다.
- 4. [RAG 기반 검색]: LLM은 두 번째 과업을 위해, 획득한 서비스 목록과 '비용 절감 가이드' 라는 키워드를 결합하여, RAG 검색 MCP 도구 search\_internal\_documents(query: "서비스 A, B, C 비용 절감 가이드")를 호출합니다.
- 5. [근거 데이터 확보]: RAG 시스템은 내부 지식 베이스에서 관련 가이드 문서를 찾아 그 내용을 LLM에 컨텍스트로 제공합니다.
- 6. [결합 응답 생성]: LLM은 실시간 조회 결과(3단계)와 검색된 문서 내용(5단계)을 종합하여, 사용자에게 최종적인 자연어 보고서를 생성하여 전달합니다.

7.4.2 민원 응대, 청약·통계 질의, 운영 로그 분석 등 도메인별 활용 예시 MCP와 RAG 통합 아키텍처는 다양한 산업 및 업무 도메인에서 구체적인 가치를 창출할 수 있습니다. 여러 요구사항을 바탕으로 한 실제 적용 시나리오는 다음과 같습니다.



업무 도메인	활용 시나리오
공공 민원 응대	먼저, RAG가 "임신·출산 진료비 신청 방법을
	알려줘"라는 질문에 대해 관련 규정 문서를 검
	색하여 정확한 정보를 제공합니다 (맥락적 지
	능). 이어서 사용자가 "바로 신청할래"라고 하
	면, MCP 도구가 실제 신청 프로세스를 실행합
	니다 (기능적 지능).
청약·통계 질의	예약 시스템 관리자가 "지난 주말 모구리야영장
	예약률과 주요 고객 연령대 통계 알려줘"라고
	질의하면, MCP 도구가 예약 DB를 실시간으로
	쿼리하여 통계 데이터를 추출하고, LLM이 이를
	분석하여 보고서 형태로 생성합니다.
운영 로그 분석	시스템 운영자가 "어젯밤 10시경 발생한 서버
	접속 지연의 원인을 관련 로그를 기반으로 분석
	해줘"라고 요청하면, RAG가 중앙 로그 관리 시
	스템(ELK)에서 관련 로그를 검색·추출하고,
	LLM이 로그 패턴을 분석하여 잠재적인 원인을
	제시합니다.

#### 7.4.3 RAG·MCP 결합 효과 측정 지표(정확도·처리 시간·재작업률)

성공적인 AI 시스템 도입은 구축으로 끝나지 않습니다. 그 성과를 객관적으로 측정하고 지속적으로 개선하기 위한 핵심 성과 지표(KPI)를 정의하고 관리하는 것이 중요합니다.

- 답변 정확도 (Accuracy): RAG를 통해 생성된 답변이 내부 근거 데이터와 얼마나 일치하는 지를 측정하는 지표입니다. (예: LLM기반 대화형 챗봇 구축 과업에서 요구한 바와 같이, 발주자가 제공한 테스트케이스 100개 중 95개 이상 정답 목표) 이는 할루시네이션 발생률을 역으로 추적하여 관리함으로써 시스템의 신뢰도를 나타냅니다.
- 업무 처리 시간 (Processing Time): 사용자가 질의를 시작한 순간부터 최종 답변(또는



MCP 도구 실행 완료)을 받기까지 걸리는 시간입니다. AI 도입 전후의 동일 업무 처리 시간을 비교하여 생산성 향상(시간 단축률)을 측정할 수 있습니다.

• 재작업률 (Rework Rate): AI가 제공한 답변이나 처리 결과가 부정확하여 사용자가 추가 질문을 하거나 수동으로 수정해야 하는 비율입니다. 이 지표의 감소는 AI 시스템의 신뢰도 와 업무 효율성 증가를 의미합니다.

궁극적으로 이러한 KPI는 단순한 측정 기준을 넘어 비즈니스 가치의 언어입니다. 정확성, 효율성, 신뢰성을 엄격하게 추적함으로써 IT 리더들은 AI에 대한 논의를 비용 중심에서 수익 동력으로 전환하고, 잘 설계된 시스템이 지속적인 혁신과 시장 리더십의 기반임을 증명할 수 있습니다.

## 제 8장 MSAP.ai란 무엇일까요? Al Native 플랫폼

#### 8.1. 왜 'AI 플랫폼'이 필요한가요?

과거에는 AI를 도입한다고 하면, 특정 팀에서 단발성으로 진행하는 '신기한 기술 실험'에 그치는 경우가 많았습니다. 이런 방식으로는 AI를 회사 전체의 경쟁력으로 만들 수 없습니다.

이제 AI는 일부 부서의 전유물이 아니라, 회사 전체가 함께 사용하고 비즈니스에 녹여내야 하는 핵심 기술이 되었습니다. 그러기 위해서는 AI 기술을 전사적으로 확산시키고, 누구나 쉽게 활용할 수 있도록 도와주는 'AI 플랫폼'이 반드시 필요합니다.

MSAP.ai는 바로 이러한 역할을 수행하는 AI 플랫폼입니다. 단순히 기술을 모아놓은 것이 아니라, AI라는 강력한 엔진을 회사의 모든 업무 시스템에 효과적으로 연결하고 공급하는 '중앙 통로'라고 할 수 있습니다.

#### 8.1.1. MSAP.ai의 정체: 'AI 발전소'를 위한 '스마트 전력망'

MSAP.ai를 쉽게 이해하기 위해 '발전소'와 '전력망'에 비유해 보겠습니다.

• 발전소 (거대 언어 모델, GPU 등): AI 기술의 핵심입니다. 엄청난 양의 '전기(AI 능력)'를 만들어냅니다.



- 전력망 (MSAP.ai 플랫폼): 발전소에서 만든 전기를 전선, 변압기 등을 통해 각 가정과 공장 까지 안전하고 효율적으로 보내주는 역할을 합니다.
- 공장이나 가정의 전기 제품 (Al Native 애플리케이션): 전기를 이용하는 도구 들입니다. 예를 들어, 챗봇의 답변 생성, 데이터 분석 및 요약 같은 기능들이죠.

아무리 뛰어난 발전소가 있어도 전기를 집까지 보내주는 전력망이 없다면 아무 소용이 없겠죠? MSAP.ai는 바로 이 '전력망'의 역할을 합니다. Al라는 강력한 에너지를 회사의 모든 업무 현장에 안정적으로 공급하고, 누구나 쉽게 사용할 수 있도록 돕는 핵심 기반 시설입니다.

이러한 '스마트 전력망'은 '마이크로서비스(Microservice)'라는 기술로 만들어집니다. AI의 다양한 기능들을 마치 레고 블록처럼 각각 독립된 부품으로 만들어, 필요한 기능만 빠르게 개발하고 수정하며 추가할 수 있습니다. 덕분에 시장 변화나 새로운 비즈니스 요구에 훨씬 민첩하게 대응할 수 있습니다.

#### 8.1.2. 소프트웨어 개발의 모든 과정을 AI로 스마트하게!

MSAP.ai는 소프트웨어 생애주기의 설계, 개발, 배포, 운영 전반에 걸쳐 AI를 통합하여 생산성과 안정성을 극대화합니다. 이는 DevSecOps 및 CI/CD 파이프라인 자동화 개념을 통해 구체화됩니다.

모든 현대 소프트웨어에 DevSecOps가 중요하지만, '코드'에 해당하는 모델과 도구, '데이터'에 해당하는 RAG 소스가 끊임없이 변화하는 AI 플랫폼에서는 그 역할이 더욱 증폭됩니다.

견고하고 자동화된 파이프라인은 이렇게 동적인 환경을 통제하고 새로운 AI 역량이 안전하고 신속하게 전달되도록 보장하는 유일한 방법입니다. 예를 들어, 개발 단계에서는 GitHub와 같은 코드 관리 시스템과 연동하여 AI가 코드 리뷰를 수행하고, 커밋 메시지를 제안하며, 잠재적인 보안 취약점을 사전에 경고하는 등 GitOps 자동화 시나리오를 지원합니다. 이를 통해 개발자는 반복적인 작업에서 벗어나 더 창의적인 문제 해결에 집중할 수 있으며, 조직은 표준화된 고품질 코드를 유지할 수 있습니다. 이는 개발 생산성을 획기적으로 향상시키는 동시에, 안정적인 시스템 운영의 기반을 마련합니다.



#### 8.1.3. 공공 및 엔터프라이즈 도입 모델

MSAP.ai는 데이터 주권과 보안이 중요한 공공 및 엔터프라이즈 환경에 최적화된 도입 모델을 제공합니다.

도입 모델	특징 및 고려사항
온프레미스 (On-premise)	데이터 주권 및 강력한 보안:

조직의 데이터센터 내에 직접 인프라를 구축하여 모든 데이터와 AI 모델을 내부에서 통제합니다. 특히 데이터 외부 유출에 민감한 공공, 금융, 의료 분야에 적합합니다. 고려사항: GPU 서버 도입 및 직접적인 인프라 운영, 관리에 대한 전문 인력과 투자가 필요합니다. | | 프라이빗 클라우드 (Private Cloud) | 유연성과 통제력의 균형: 조직이 이미 보유하거나 계약한 프라이빗 클라우드 환경을 활용하여 인프라를 구축합니다. 기존 클라우드 자원을 활용하면서도 온프레미스와 유사한 수준의 보안 및 통제력을 확보할 수 있습니다. 고려사항: 클라우드 환경에 대한 이해와 운영역량이 필요합니다. |

MSAP.ai는 AI 기술을 조직의 핵심 역량으로 전환시키는 전략적 기반을 제공합니다. 이는 단순히 새로운 기능을 추가하는 것을 넘어, 기존의 IT 자산을 AI와 연결하여 그 가치를 극대화하는 과정에서 진정한 힘을 발휘합니다. 다음 섹션에서는 이러한 플랫폼이 기업이 보유한 방대한 API 자산을 어떻게 MCP(Model Context Protocol)로 변환하여 AI와 연결하는지에 대한 구체적인 기술을 심도 있게 살펴보겠습니다.

#### 8.2 기존 API JSON 스키마의 MCP 자동 변환 기능

기업이 수 년간 축적해 온 방대한 기존 API 자산은 AI 시대의 걸림돌이 아니라, 가장 중요한 핵심 자산으로 전환될 수 있습니다. 이 혁신적인 전환의 중심에는 기존 API를 AI가 이해하고 자율적으로 호출할 수 있는 표준화된 '도구(Tool)'로 변환하는 MCP의 역할이 있습니다. 이는 레거시 시스템의 기능을 AI 에이전트의 팔과 다리로 만들어주는 것과 같습니다.



#### 8.2.1. API-MCP 매핑: 설계 패러다임의 전환

기존 REST API의 JSON 스키마를 MCP 도구로 매핑하는 과정은 단순한 기술적 변환을 넘어 설계 패러다임의 전환을 요구합니다.

REST API는 본질적으로 '리소스 중심(Resource-oriented)'으로 설계되어 특정 데이터(리소스)에 대한 생성, 조회, 수정, 삭제(CRUD)에 초점을 맞춥니다. 반면, MCP 도구는 LLM이 특정 목표를 달성하기 위해 사용하는 '역량 중심(Capability-oriented)'으로 설계됩니다.

이러한 본질적 차이 때문에 OpenAPI 명세서 등을 이용한 단순 자동 변환은 한계가 명확합니다. LLM이 도구의 목적과 사용법을 명확히 이해하고 혼동 없이 사용하기 위해서는 도구의 이름, 설명, 파라미터를 LLM 친화적으로 재정의하는 과정이 필수적입니다. 예를 들어, 시스템 내에 get\_user\_info, fetch\_user\_profile, lookup\_user와 같이 유사한 기능을 수행하지만 이름이 다른 API들이 산재해 있다면, LLM은 어떤 도구를 사용해야 할지 혼란을 겪을 수 있습니다. 따라서 이러한 API들을 '사용자 프로필 조회'라는 명확한 단일 역량으로 추상화하여 하나의 MCP도구로 정의하는 설계적 판단이 필요합니다.

이러한 추상화는 단순히 인간의 편의를 위한 것이 아니라, LLM 자체를 위한 결정적인 최적화 과정입니다. get\_user\_info와 fetch\_user\_profile 같은 도구들 사이의 모호성을 제거함으로 써 우리는 모델에 가해지는 인지적 부하를 극적으로 줄일 수 있습니다. 이는 곧 낭비되는 도구 선택 토큰의 감소, 추론 비용 절감, 그리고 환각이나 부정확한 도구 사용의 위험성을 현저히 낮추는 결과로 이어집니다.

#### 8.2.2. 대규모 API 자산의 체계적 전환 프로세스

조직 내 수백, 수천 개의 API를 MCP로 효과적으로 전환하기 위해서는 체계적인 접근이 필요합니다.

- 1. 1단계: MCP 인벤토리 구축 조직의 업무 도메인별로 정보 시스템과 API를 목록화하여 'MCP화 대상' 후보군을 정의합니다. "AI가 호출할 수 있으면 유의미한 기능은 무엇인 가?"라는 관점에서 자산을 파악하고 우선순위를 정하는 단계입니다.
- 2. 2단계: MCP 서버 전략 수립 모든 API를 단일 MCP 서버로 노출하는 것은 비효율적이며 LLM의 성능을 저하시킬 수 있습니다. 특정 MCP 서버가 전체 API(107개)를 제공하는 엔



드포인트와 핵심 기능 중심의 API(38개)를 제공하는 엔드포인트를 별도로 운영하는 사례처럼, 사용 사례나 업무 도메인에 맞춰 도구 집합(Tool Set)을 그룹화하여 여러 MCP 서버로 제공하는 전략이 효과적입니다.

3. 3단계: 자동화 및 프록시 적용 전환 작업의 효율성을 높이기 위해 OpenAPI 명세서로부터 MCP 서버 코드의 기본 골격(boilerplate)을 자동 생성하는 도구를 활용합니다. 개발자는 자동 생성된 코드를 기반으로, 2단계에서 수립된 전략에 따라 LLM에 최적화된 도구 이름과 설명을 추가하고 파라미터를 재정의하는 작업에 집중합니다.

#### 8.2.3. 변경에 대응하는 자동 갱신 구조

API는 비즈니스 요구사항에 따라 빈번하게 변경될 수 있습니다. API 명세가 변경될 때마다 MCP 정의가 수동으로 갱신되어야 한다면, 이는 또 다른 기술 부채가 될 것입니다. 따라서 자동 갱신 구조를 CI/CD 파이프라인에 통합하는 DevSecOps 관점의 접근이 중요합니다.

API 명세서의 변경이 감지되면 자동으로 MCP 서버를 재배포하는 파이프라인을 구축할 수 있습니다. 더 나아가, 컴파일이나 재배포 없이 실시간으로 API 변경을 MCP 정의에 반영할 수 있는 '라이브 프록시 서버(Live proxy server)' 아키텍처를 도입하면, 빈번하게 변경되는 API에 대해서도 최고의 민첩성을 확보할 수 있습니다.

기존 API 자산을 MCP라는 표준화된 인터페이스로 전환하는 것은 단순히 개별 기능을 AI와 연결하는 것을 넘어섭니다. 이는 조직의 모든 디지털 자산을 AI가 활용할 수 있는 거대한 도구 상자로 만드는 과정입니다. 다음 섹션에서는 이렇게 준비된 MCP 도구들이 RAG, Observability와 같은 다른 AI 핵심 기술과 어떻게 통합되어 강력한 시너지를 창출하는지 구체적인 아키텍처를 통해 살펴보겠습니다.

#### 8.3 MCP·RAG·Observability 통합 아키텍처

현대의 AI 플랫폼은 단일 기술의 합이 아닌, 여러 핵심 기술 요소가 유기적으로 결합된 통합 아 키텍처 위에서 그 진정한 가치를 발휘합니다. 특히 AI의 '행동'을 담당하는 MCP, '지식'을 제공하는 RAG(Retrieval-Augmented Generation), 그리고 시스템의 '상태'를 감시하고 분석하는 Observability는 분리된 기술이 아니라 하나의 아키텍처로 융합될 때 비로소 지능적이고 자율적인 시스템을 구현할 수 있습니다. 이 세 가지 요소의 통합은 AI 플랫폼의 두뇌, 지식, 그리고 신경



계를 완성하는 것과 같습니다.

#### 8.3.1. 통합 인프라 구성: MSA 기반의 AI 서비스망

MSAP.ai 기반의 통합 아키텍처는 클라우드 네이티브의 핵심 요소들로 구성됩니다. 먼저, API 게이트웨이가 외부의 모든 요청에 대한 단일 진입점(Single Entry Point) 역할을 수행하며, 인증, 라우팅, 사용량 제어 등 공통 기능을 처리합니다.

게이트웨이를 통과한 요청은 서비스 메시(Service Mesh) 에 의해 관리되는 마이크로서비스 환경으로 전달됩니다.

서비스 메시는 마이크로서비스 간의 통신을 안정적이고 안전하게 제어하며, 트래픽 관리, 서비스 디스커버리, 서킷 브레이커와 같은 기능을 제공합니다.

이 환경에서 각 MCP 서버는 특정 비즈니스 도메인(예: 인사, 재무, 재고 관리)에 대한 책임을 지는 독립적인 마이크로서비스로 컨테이너화되어 배포됩니다. 쿠버네티스는 이러한 수많은 컨테이너화된 MCP 서버들을 자동으로 배포, 확장, 관리하는 오케스트레이션 역할을 수행하여, 전체시스템의 탄력성과 회복탄력성을 보장합니다.

#### 8.3.2. RAG-MCP 통합 구조: 지식과 행동의 결합

RAG와 MCP의 통합은 AI 에이전트가 단순히 주어진 지식을 검색하는 것을 넘어, 검색된 지식을 바탕으로 실질적인 행동을 수행하게 만드는 핵심적인 구조입니다. 이 통합 워크플로우는 다음과 같이 진행됩니다.

사용자가 "최근 체결된 '클라우드 네이티브 전환 사업' 계약의 주요 내용을 요약하고, 관련 담당자에게 검토 요청 메일을 보내줘"와 같은 복합적인 질의를 하면, AI 에이전트는 먼저 RAG 파이프라인을 호출하는 '내부 문서 검색 MCP 도구(Tool #1)'를 실행합니다. 이 도구는 내부 지식 베이스에서 해당 계약 문서를 찾아 핵심 내용을 추출합니다. 그 다음, 에이전트는 추출된 정보를 바탕으로 조직의 인사 시스템과 연동된 '담당자 조회 및 메일 발송 MCP 도구(Tool #2)'를 호출하여 담당자에게 요약된 내용과 함께 검토 요청을 보냅니다. 최종적으로 이 모든 과정을 종합하여 사용자에게 "계약서 요약과 함께 담당자에게 메일 발송을 완료했습니다."라는 최종 답변을 생성합니다.

더 나아가, 'RAG 인덱싱 파이프라인의 MCP 도구화'라는 개념을 적용할 수 있습니다. 이는 문서 수집, 정제, 임베딩, 인덱싱에 이르는 RAG의 데이터 준비 과정 자체를 AI 에이전트가 직접



MCP 도구로 호출할 수 있게 만드는 것입니다. 이를 통해 에이전트는 필요에 따라 최신 정보를 동적으로 학습하고 지식 베이스를 스스로 갱신하는, 한 차원 높은 자율성을 갖게 됩니다.

# 8.3.3. VibeOps/AlOps 시나리오: MSAP APM & Observability 기반의 지능형 자율 운영

MSAP APM과 MSAP Observability를 통해 수집된 실시간 성능 데이터(Metrics, Traces, Logs)는 MSPA.ai의 AI 플랫폼과 결합하여, 복잡한 시스템 운영을 자동화하고 지능화하는 AlOps, 즉 MSPA.ai가 지향하는 VibeOps를 실현합니다.

예를 들어, 시스템 운영자가 "특정 서비스에서 지난 1시간 동안 에러율이 급증한 원인을 분석하고 관련 담당자에게 알림을 보내줘"라는 자연어 프롬프트를 입력했다고 가정해 봅시다.

MSAP VibeOps AI 에이전트는 다음과 같은 일련의 작업을 자율적으로 수행합니다.

- 1. [로그 분석] 먼저 MSAP Observability의 통합 로그 관리 기능에 접근합니다. '로그 분석 MCP 도구'를 호출하여 해당 서비스의 에러 로그를 실시간으로 수집하고, AI 기반으로 이상 패턴을 즉시 식별합니다. 더 이상 별도의 로그 시스템을 조회할 필요 없이, 단일 플랫폼 내에서 분석이 완료됩니다.
- 2. [메트릭-트레이스 연관 분석] 다음으로, 에러 발생 시점의 상세한 컨텍스트를 파악하기 위해 MSAP APM의 트랜잭션 데이터와 MSAP Observability의 인프라 메트릭(CPU, Memory, Network 등)을 동시에 조회합니다. '메트릭-트레이스 연관 분석 MCP 도구'를 호출하여 급증한 에러와 특정 리소스 병목 현상, 비정상적인 트랜잭션 간의 인과관계를 심층 분석합니다.
- 3. [인시던트 생성 및 조치] 분석된 근본 원인(Root Cause)과 핵심 지표를 바탕으로, Jira, Slack, PagerDuty 등 외부 시스템과 연동된 '인시던트 생성 및 협업 MCP 도구'를 호출합니다. 단순히 알림을 보내는 것을 넘어, 원인 분석 리포트와 조치 권고안을 포함한 인시던트를 자동 생성하여 담당자에게 전달하고 협업을 유도합니다.

이처럼 MSPA.ai의 AI 플랫폼(MCP, RAG)과 MSAP APM, MSAP Observability의 유기적인 통합은 단순한 모니터링을 넘어, 시스템이 스스로 문제를 진단하고, 해결 방안을 제시하며, 지속적으로 학습하는 지능형 자율 운영(Intelligent Autonomous Operations) 체계를 완성합니다. 이는 IT 운영의 패러다임을 바꾸는 VibeOps의 핵심입니다.



MCP, RAG, Observability의 유기적인 통합은 지능적이고 자율적인 시스템 운영의 미래를 제시합니다. 그러나 이러한 강력한 기술 통합의 최종적인 가치는 결국 사용자의 업무 경험(UX)을 어떻게 혁신하는가에 달려 있습니다. 다음 섹션에서는 이러한 기술들이 어떻게 최종 사용자의 업무 방식을 근본적으로 바꾸는지 논의해 보겠습니다.

#### 8.4 프롬프트 중심 업무 UX와 Widget 기반 통합

기술의 발전은 궁극적으로 사용자의 경험(UX)을 어떻게 혁신하는가에 따라 그 가치가 평가됩니다. 더 이상 사용자가 복잡한 GUI 화면의 메뉴와 버튼을 찾아 여러 화면을 탐색할 필요 없이, 일상 언어로 된 프롬프트 하나만으로 원하는 업무를 즉시 처리하는 '프롬프트 중심 UX'가 바로 그 혁신의 정점에 있습니다. Model Context Protocol(MCP)은 이러한 혁신적인 경험을 현실로 만드는 핵심 기술입니다.

#### 8.4.1. 핵심 기능: Prompts와 Elicitation

프롬프트 창과 위젯(Widget) 기반의 대화형 UX는 MCP가 제공하는 두 가지 핵심 기능, 'Prompts'와 'Elicitation'을 통해 기술적으로 구현됩니다. 이 두 기능은 AI와 사용자 간의 소통을 더욱 빠르고 직관적으로 만들어 줍니다.

• MCP Prompts: 자주 쓰는 업무를 예약하는 '단축키'

자주 사용하는 복잡한 업무 절차를 재사용 가능한 템플릿으로 정의하는 기능입니다. 예를 들어, /보고서작성, /휴가신청, /법인카드결제 와 같은 '슬래시 커맨드' 형태로 프롬프트를 미리 정의해 둘 수 있습니다. 사용자는 간단한 명령어만 입력하면 여러 단계로 구성된 복잡한 워크플로우를 즉시 시작할 수 있습니다. 이는 사용자가 반복적인 업무 절차를 일일이 기억할 필요 없이, 일관되고 효율적인 방식으로 작업을 수행하도록 돕는 강력한 단축키 역할을 합니다.

• MCP Elicitation: AI가 먼저 물어보는 '눈치 빠른 대화'

AI 에이전트가 업무를 수행하는 도중 추가 정보나 사용자의 결정이 필요할 때, 대화의 흐름을 끊지 않고 사용자에게 먼저 말을 걸어 필요한 것을 요청하는 지능적인 상호작용 방식입니다.



쉽게 말해, AI가 똑똑한 비서처럼 다음 단계를 예측하고 사용자에게 필요한 질문을 먼저 던지는 것입니다. 예를 들어, 사용자가 /휴가신청 프롬프트를 실행하면, AI는 가만히 기다리는 대신 "휴가 날짜를 선택해주세요"라는 메시지와 함께 달력 위젯을 띄워줍니다. 보고서제출 전에는 '승인/반려' 버튼을 보여주며 최종 의사를 묻습니다.

이처럼 Elicitation 기능은 사용자가 모든 정보를 한 번에 입력해야 하는 부담을 덜어주고, AI가 필요한 정보를 능동적으로 이끌어내어(Elicit) 막힘없는 업무 처리를 가능하게 합니다. 덕분에 사용자는 자연스러운 대화 속에서 더 쉽고 정확하게 결정을 내릴 수 있습니다.

#### 8.4.2. 컨텍스트 전환 없는 업무 연속성 유지

프롬프트 중심 UX의 가장 큰 장점은 사용자가 여러 웹 화면이나 애플리케이션을 전환하지 않고, 단일 인터페이스 내에서 모든 업무 흐름을 완결할 수 있다는 점입니다.

예를 들어, 개발자가 IDE나 협업 툴의 채팅창에서 "다음 주 부산 출장 계획을 세우고 보고서를 작성해줘"라고 입력했다고 상상해 보십시오. AI 에이전트는 이 프롬프트 하나를 받아 백그라운드에서 다음과 같은 작업을 MCP 도구들을 통해 순차적으로 또는 병렬적으로 처리합니다.

- 1. 항공편 조회 MCP 도구를 호출하여 최적의 항공편을 검색합니다.
- 2. 숙소 예약 MCP 도구를 호출하여 예약 가능한 호텔을 찾습니다.
- 3. 재무 시스템과 연동된 MCP 도구를 통해 법인카드 한도를 확인합니다.
- 4. 모든 예약이 완료되면, 사용자의 캘린더에 일정을 등록하는 MCP 도구를 실행합니다.
- 5. 마지막으로, 출장 계획을 바탕으로 보고서 초안을 작성하는 MCP 도구를 호출합니다.

이 모든 과정이 완료되면, AI 에이전트는 최종 결과물인 '출장 계획 요약 및 보고서 초안'만을 사용자에게 보여줍니다. 사용자는 항공사 웹사이트, 호텔 예약 사이트, 사내 재무 시스템, 캘린더, 문서 작성 툴을 일일이 오갈 필요 없이, 원래 작업하던 화면에서 모든 업무를 완료할 수 있습니다.

#### 8.4.3. 업무 집중도와 생산성에 미치는 영향

이러한 프롬프트 중심 UX는 개인과 조직의 생산성에 지대한 긍정적 영향을 미칩니다.

• 집중도 향상: 불필요한 애플리케이션 간의 '컨텍스트 스위칭(Context Switching)' 비용을



극적으로 줄여줍니다. 사용자는 여러 도구를 사용하는 방법을 배우거나 기억할 필요 없이, 오직 자신의 핵심 업무 목표에만 깊이 몰입할 수 있습니다.

• 생산성 증대: Gartner의 분석에 따르면, MCP 도입은 모델 정확도를 25% 향상시키고 시 솔루션의 배포 속도를 50% 단축시키는 효과가 있습니다. 이는 개발자뿐만 아니라 일반 현업 사용자에게도 유사한 수준의 생산성 향상으로 이어질 수 있음을 시사합니다. 복잡한 시스템 조작에 소요되던 시간이 단축되고, 자연어 기반의 직관적인 상호작용이 가능해지면서 모든 사용자가 전문가 수준의 업무 처리 속도와 정확도를 달성할 수 있게 됩니다.

결론적으로, 프롬프트 중심 UX는 MCP와 같은 AI 통합 기술의 최종 가치를 사용자에게 전달하는 핵심 통로입니다. 이는 단순한 인터페이스의 변화를 넘어, 일하는 방식 자체를 근본적으로 혁신하는 패러다임의 전환입니다. 다음 섹션에서는 이러한 혁신적인 경험을 조직에 성공적으로 도입하기 위한 단계별 로드맵과 구체적인 실행 체크리스트를 제시하겠습니다.

## 8.5 단계별 도입 로드맵 및 실행 체크리스트

성공적인 AI 플랫폼 도입은 하룻밤에 이루어지는 '빅뱅(Big-bang)' 방식이 아닌, 측정 가능하고 반복 가능한 단계적 접근을 통해 이루어져야 합니다. MCP 기반 AI 플랫폼을 조직에 성공적으로 안착시키기 위한 현실적인 3단계 로드맵과 각 단계별 핵심 실행 과제를 제시하여, 의사결정자와 실무자 모두에게 명확한 가이드를 제공하고자 합니다.

1단계: 개념 증명 (PoC, Proof of Concept)

이 단계의 목표는 1~2개의 핵심 도메인을 선정하여 MCP와 RAG의 기술적 타당성을 검증하고, 측정 가능한 비즈니스 가치를 입증하는 파일럿을 성공적으로 수행하는 것입니다.

실행 체크리스트:

- 1. AI 연계 자산 식별 및 목록화: 현재 시스템 기능 중 AI가 호출했을 때 유의미한 가치를 창출할 수 있는 기능 후보 목록을 작성하여 잠정 MCP 후보군을 정의하십시오.
- 2. 최초 실행 도메인 선정: 리스크가 낮고 성공 시 효과가 명확하게 드러나는 도메인(예: 내부 규정 및 매뉴얼 질의응답)을 우선적으로 선정하십시오.



- 3. PoC 시나리오 설계 및 구현: 선정된 도메인 내 3~5개 핵심 기능을 MCP 도구로 구현하고, 관련 문서를 RAG 인덱스로 구축하여 '자연어 질의-실행' 통합 시나리오를 구성하십시오.
- 4. 성공 측정 지표(KPI) 정의 및 확립: PoC의 성공을 객관적으로 판단할 명확한 KPI(예: 구현된 MCP 도구 수, API 호출 빈도, 업무 처리 시간 단축률)를 사전에 정의하십시오.

2단계: 핵심 업무 확장

PoC의 성공을 바탕으로, MCP 적용 범위를 조직의 핵심 업무 도메인으로 확대하고, 전사적인 MCP 자산을 체계적으로 관리하기 위한 카탈로그와 거버넌스 기반을 구축하는 단계입니다.

실행 체크리스트:

- 전사 MCP 개발 표준 수립 및 강제: 기술 부채를 방지하고 상호운용성을 보장하기 위해, 도구 명명 규칙, 데이터 연동 가이드라인, 버전 관리 정책을 포함한 전사적 MCP 개발 표준을 수립하고 준수를 의무화하십시오.
- 2. 전사 보안 정책 정의 및 적용: OAuth 2.1 기반의 표준 인증 프로토콜을 적용하고, 역할 기반 접근 제어(ACL) 정책을 수립하여 MCP 도구별 접근 권한을 통제하는 전사 보안 표준을 확립하십시오.
- 3. 중앙 MCP 카탈로그(포털) 구축 의무화: 모든 사업부에서 자산을 재사용하고 개발을 가속 화할 수 있도록, 등록된 MCP 서버와 도구를 검색하고 재사용할 수 있는 중앙 집중식 카탈 로그 구축을 지시하십시오.

3단계: 플랫폼 내재화

외부 기술 및 인력에 대한 의존도를 줄이고, 자체적인 AI 플랫폼 개발 및 운영 역량을 확보하여 AI 기술을 조직의 핵심 DNA로 완전히 내재화하고, 지속 가능한 운영 거버넌스를 정착시키는 단계입니다.

실행 체크리스트:

1. Al 운영(DevSecOps) 전담 조직 구성: MCP 서버의 개발, 배포, 모니터링, 보안 및 전체 생애주기를 전담하는 중앙 거버넌스 조직을 구성하여 플랫폼의 안정적인 운영과 지속적인 발전을 책임지게 하십시오.



- 2. Observability 체계 고도화: 모든 MCP 도구 호출에 대한 로그, 메트릭, 트레이스 데이터 를 수집하고, 이를 중앙에서 통합 분석하여 성능 병목을 식별하고 비용을 최적화하는 고도 화된 모니터링 시스템을 구축하십시오.
- 3. 자체 성숙도 평가 및 발전 로드맵 수립: 정기적으로 MCP 성숙도 모델에 기반하여 조직의 AI 활용 수준을 자체 평가하고, 그 결과를 다음 단계의 기술 투자 및 인력 양성 계획에 반영하는 선순환 구조를 만드십시오.

# 제9장: References & Links

- Microservices Martin Fowler (https://martinfowler.com/articles/microservices.html)
- Building Microservices: Designing Fine-Grained Systems Sam Newman, O'Reilly Media (https://www.oreilly.com/library/view/building-microservices-2nd/9781492034018/)
- Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks Patrick Lewis, et al. (Facebook AI) – (https://arxiv.org/abs/2005.11401)
- What is RAG? Pinecone (https://www.pinecone.io/learn/retrieval-augmented-generation/)
- Kubernetes Official Documentation (https://kubernetes.io/docs/home/)
- Cloud Native Trail Map Cloud Native Computing Foundation (CNCF) (https://raw.githubusercontent.com/cncf/trailmap/master/CNCF\_TrailMap\_latest.png)
- The Economics of Large Language Models SemiAnalysis (https://www.semianalysis.com/p/the-economics-of-large-language)
- Trends in the dollar training cost of machine learning systems Epoch Al
   (https://epochai.org/blog/trends-in-the-dollar-training-cost-of-machine-learning-systems)
- Bringing AI to the Edge: A New Architecture for a Distributed World –
   Andreessen Horowitz (a16z) (https://a16z.com/ 2020/02/16/ ai-edge computing-new-architecture/)



- 9.1.1 Model Context Protocol Specification https://modelcontextprotocol.io/ specification/2025-06-18 Model Context Protocol
- 9.1.2 Model Context Protocol (MCP) GitHub https://github.com/modelcontextprotocol/modelcontextprotocol GitHub
- 9.1.3 Introducing the Model Context Protocol https://www.anthropic.com/ news/model-context-protocol Anthropic
- 9.1.4 Model Context Protocol Wikipedia https://en.wikipedia.org/wiki/Model\_Context\_Protocol 위키백과
- 9.1.5 Catalog of official Microsoft MCP servers https://github.com/microsoft/mcp GitHub
- 9.2.1 About Model Context Protocol (MCP) GitHub Copilot https://docs.github.com/en/copilot/concepts/context/mcp GitHub Docs
- 9.2.2 Using Model Context Protocol GitHub Copilot https://docs.github.com/ copilot/customizing-copilot/using-model-context-protocol GitHub Docs
- 9.2.3 Windows is getting support for the "USB-C of Al apps" https://www.theverge.com/ news/ 669298/ microsoft-windows-ai-foundry-mcp-support The Verge
- 9.2.4 What is Model Context Protocol (MCP)? https://www.itpro.com/tech-nology/artificial-intelligence/what-is-model-context-protocol-mcp IT Pro
- 9.2.5 Model Context Protocol Curriculum for Beginners https://github.com/microsoft/mcp-for-beginners GitHub
- 9.3.1 What is Retrieval-Augmented Generation (RAG)? AWS https:// aws.amazon.com/ what-is/ retrieval-augmented-generation/ Amazon Web Services, Inc.
- 9.3.2 Retrieval Augmented Generation (RAG) overview Azure Al Search
   https://learn.microsoft.com/azure/search/retrieval-augmented-generation-overview Microsoft Learn
- 9.3.3 What is Retrieval-Augmented Generation (RAG)? Google Cloud https://cloud.google.com/use-cases/retrieval-augmented-generation Google Cloud



- 9.3.4 What is retrieval-augmented generation (RAG)? McKinsey https://www.mckinsey.com/ featured-insights/ mckinsey-explainers/ what-is-retrieval-augmented-generation-rag McKinsey & Company
- 9.3.5 What Is Retrieval-Augmented Generation (RAG)? NVIDIA Blog https:// blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/ NVIDIA Blog
- 9.4.1 The Compute Divide: The Cost of Training Large Language Models arXiv:2404.09356 arXiv
- 9.4.2 Training large language models more efficiently Amazon Science
   https://www.amazon.science/blog/training-large-language-models-more-efficiently amazon.science
- 9.4.3 Trends in AI inference energy consumption Energy & AI (ScienceDirect)
   ScienceDirect
- 9.4.4 Densing law of LLMs Nature Machine Intelligence https://www.nature.com/articles/s42256-025-01137-0 Nature
- 9.4.5 China's DeepSeek says its hit Al model cost just \$294,000 to train –
   Reuters Reuters
- 9.5.1 MSAP.ai 공식 홈페이지 https://www.msap.ai/ MSAP
- 9.5.2 MSAP.ai Product MSAP.ai Platform & AI https://www.msap.ai/product/msap-ai/ MSAP
- 9.5.3 클라우드 네이티브 MSA 인프라 자동화 솔루션 MSAP.ai Platform https://www.msap.ai/product/msap-ai/platform/ MSAP
- 9.5.4 [백서] 차세대 지능형 애플리케이션 플랫폼: AI와 MSA 시대를 위한 MSAP.ai 기술 백서 https://www.msap.ai/resource/ai-msa-platform/ MSAP
- 9.5.5 투라인클라우드, MSAP.ai·컨설팅서비스 조달청 디지털서비스몰 등록 주요 기사 URL (예: ITDaily, DigitalToday 등) https://www.digitaltoday.co.kr/news/articleView.html?idxno=600509&utm\_source=chatgpt.com



# **Contact Us**



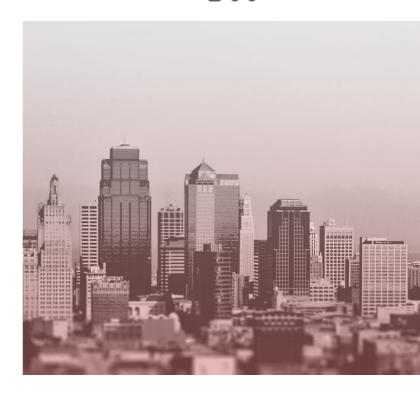
02-6953-5427



hello@msap.ai



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



엠에스에이피닷에이아이 | MSAP.ai

전화 : (02) 6953 - 5427 팩스 : (02) 469 - 7247 메일 : hello@msap.ai