

# 세션 PromptOps : 동시접속자 데이터 기반 증설 의사결정 프레임워크

"세션이 10만인데 실제 동시접속자는 얼마죠?"— 운영 현장에서 매일 반복되는 이 질문에 정확히 답할 수 있는 기업은 드뭅니다. 기존 WAS 모니터링 도구는 활성 세션 수만 보여줄 뿐 실제 동시접속자를 구분하지 못하며, 장애 발생 시 세션·트랜잭션·인프라 로그를 수작업으로 교차 분석하는 데 수 시간이 소모됩니다. 본 백서는 16KB 메모리로 수억 사용자를 오차율 0.81% 이내로 집계하는 HyperLogLog 기반 측정, IMDG 세션 클러스터링, 그리고 자연어 한 줄로 장애 원인을 분석하는 PromptOps 방법론을 세션-트랜잭션-AI 3계층 통합 아키텍처로 제시합니다.



 [hello@cncf.co.kr](mailto:hello@cncf.co.kr)

 02-469-5426

 [www.cncf.co.kr](http://www.cncf.co.kr)

# Contents

<b>1장. WAS 세션 관리의 과제와 PromptOps의 탄생</b>	<b>4</b>
1.1 엔터프라이즈 WAS 환경의 세션 관리 과제 . . . . .	4
1.1.1 “활성 세션 수”와 “실제 동시 사용자 수”의 괴리 . . . . .	4
1.1.2 수작업 운영 통계의 한계와 Seasonality 대응 부재 . . . . .	5
1.1.3 장애 원인 추적의 컨텍스트 단절 . . . . .	6
1.2 솔루션 기업의 접근: 세션-트랜잭션-AI 통합 플랫폼 . . . . .	7
1.2.1 세션 클러스터링 솔루션의 시작과 배경 . . . . .	7
1.2.2 PromptOps(Prompt Operations)의 정의 . . . . .	8
1.2.3 VibeOps 패러다임: AI 코파일럿 기반 운영 체계로의 전환 . . . . .	9
<b>2장: 세션-트랜잭션-AI 3계층 통합 아키텍처</b>	<b>10</b>
2.1 세션 클러스터링 솔루션: IMDG 기반 세션 계층 . . . . .	10
2.1.1 IMDG 기반 세션 저장소와 Failover 메커니즘 . . . . .	10
2.1.2 이기종 WAS 간 세션 공유와 중복 로그인 방지 . . . . .	12
2.2 APM 솔루션: HyperLogLog 기반 트랜잭션 계층 . . . . .	13
2.2.1 HyperLogLog(HLL) 기반 동시접속자 집계 원리 . . . . .	14
2.2.2 시간 단위 롤업 구조와 3가지 사용자 식별 모드 . . . . .	15
2.2.3 100% 실시간 트랜잭션 모니터링과 OpenTelemetry 통합 . . . . .	17
2.3 CogentAI / PromptOps: LLM+RAG+MCP 기반 AI 계층 . . . . .	18
2.3.1 하이브리드 LLM과 이중 데이터 소스 아키텍처 . . . . .	19
2.3.2 MCP 프로토콜 기반 시스템 연동과 개인정보 보호 . . . . .	20
2.4 경쟁 제품 대비 아키텍처 차별성 . . . . .	22
2.4.1 Redis+별도APM vs Hazelcast+Datadog vs 솔루션 기업 통합 플랫폼 비교	22
<b>3장: PromptOps 핵심 활용 시나리오</b>	<b>23</b>
3.1 동시접속자 기반 서버 사이징 의사결정 . . . . .	24
3.1.1 실시간 동시접속자 조회와 과거 동일 시간대 비교 . . . . .	24

- 3.1.2 서버 증설·감소 의사결정을 위한 데이터 기반 근거 . . . . . 26
- 3.1.3 Kubernetes HPA 연동을 통한 자동 스케일링 . . . . . 27
- 3.2 Seasonality 패턴 분석과 장애 예방 . . . . . 29
  - 3.2.1 시간축 롤업 데이터 기반 주기적 패턴 탐지 . . . . . 29
  - 3.2.2 과거 동일 시기 대비 접속 패턴 비교와 사전 대응 . . . . . 30
  - 3.2.3 장애 근본 원인 분석(RCA): 사용자-요청-인프라 연결 추적 . . . . . 32
- 3.3 IT 의사결정자를 위한 보고서 자동 생성 . . . . . 33
  - 3.3.1 자연어 질의 기반 동시접속자 보고서 생성 . . . . . 33
  - 3.3.2 사용자 행동 분석과 비정상 접속 탐지 . . . . . 34
- 3.4 사용 시 주의사항과 데이터 정확도 . . . . . 36
  - 3.4.1 HyperLogLog 오차율과 사용자 식별 모드별 주의점 . . . . . 36
  - 3.4.2 LLM 할루시네이션 대응과 운영 데이터 신뢰성 . . . . . 37
- 4장: 도입 사례와 사용자 그룹별 활용 가치 . . . . . 38**
  - 4.1 PromptOps 적용 사례 . . . . . 39
    - 4.1.1 세션-트랜잭션-LLM 통합 운영 시나리오 . . . . . 39
    - 4.1.2 WAS OOM 장애의 AI 자동 분석 사례 . . . . . 40
    - 4.1.3 공공기관 도입 사례와 운영 효율화 성과 . . . . . 41
  - 4.2 사용자 그룹별 PromptOps 활용 가치 . . . . . 42
    - 4.2.1 CTO·IT Director: 서버 투자 의사결정 근거 확보 . . . . . 42
    - 4.2.2 기술 기획팀장·WAS 운영자: Seasonality 기반 리소스 계획과 장애 분석 . . . . . 44
    - 4.2.3 DevOps 엔지니어·보안 관리자: 자동 스케일링과 비정상 접속 탐지 . . . . . 45
- 5장: PromptOps 적용 가이드 . . . . . 46**
  - 5.1 3단계 점진적 도입 경로 . . . . . 46
    - 5.1.1 Phase 1: 세션 클러스터링 솔루션 도입 — 세션 클러스터링 기반 마련 . . . . . 46
    - 5.1.2 Phase 2: APM 솔루션 연동 — 트랜잭션 모니터링과 동시접속자 집계 . . . . . 48
    - 5.1.3 Phase 3: CogentAI/PromptOps 적용 — AI 기반 자연어 운영 활성화 . . . . . 49
  - 5.2 기술 연동과 확장 . . . . . 51
    - 5.2.1 Kubernetes·OpenTelemetry·REST API 연동 . . . . . 51

5.2.2 MCP 기반 사내 시스템 연동: ERP·그룹웨어·DB . . . . .	53
5.3 라이선스와 인프라 요구사항 . . . . .	54
5.3.1 제품별 라이선스 구조와 배포 옵션 . . . . .	54
5.3.2 인프라 요구사항과 GPU 서버 무상 임대 프로그램 . . . . .	56
<b>Appendix</b>	<b>57</b>
References . . . . .	57
Glossary . . . . .	59

# 1장. WAS 세션 관리의 과제와 PromptOps의 탄생

## 1.1 엔터프라이즈 WAS 환경의 세션 관리 과제

엔터프라이즈 환경에서 Web Application Server(WAS)는 대규모 트래픽과 다양한 사용자 요구를 실시간으로 처리해야 하며, 이 과정에서 세션 관리가 시스템의 안정성과 확장성을 결정짓는 핵심 요소로 작용합니다. 하지만 실제 현장에서는 활성 세션 수와 동시 사용자 수의 괴리, 수작업에 의존하는 운영 통계, 장애 원인 추적의 어려움 등 다양한 문제가 발생하고 있습니다. 이러한 과제들은 단순히 기술적인 문제를 넘어 IT 운영의 효율성과 정확성, 그리고 비즈니스 민첩성까지 저해하는 요인으로 작용합니다. 본 절에서는 엔터프라이즈 WAS 환경에서 세션 관리가 직면하는 주요 과제와 그 배경, 그리고 이로 인해 발생하는 운영상의 복잡성을 심층적으로 분석합니다.

### 1.1.1 “활성 세션 수”와 “실제 동시 사용자 수”의 괴리

엔터프라이즈 WAS 환경에서 활성 세션 수는 일반적으로 서버의 세션 관리 모듈이 유지하는 세션 객체의 총합을 의미합니다. 하지만 실제 동시 사용자 수와는 큰 차이가 발생할 수 있습니다. 예를 들어, 사용자가 로그아웃하지 않고 브라우저를 닫거나 네트워크 장애로 세션이 비정상적으로 종료되는 경우, WAS는 여전히 해당 세션을 활성 상태로 유지합니다. 반대로, 짧은 시간 동안 여러 사용자가 접속했다가 빠르게 이탈하는 경우 실제 동시 접속자는 많지만 세션 수는 적을 수 있습니다. 이러한 괴리는 서버 사이징이나 리소스 할당에 있어 정확한 판단을 어렵게 만듭니다.

Sticky Session 방식은 로드 밸런서가 사용자의 요청을 항상 동일한 WAS 인스턴스로 전달하는 구조입니다. 이 방식은 세션 데이터의 일관성을 보장하지만, 특정 서버에 장애가 발생하면 해당 사용자의 세션 데이터가 유실되고 서비스가 중단되는 단일 장애점(SPOF, Single Point of Failure) 문제가 발생합니다. 이로 인해 고가용성(HA) 환경에서는 Sticky Session을 피하거나 외부 세션 저장소(IMDG 등)를 도입해야 합니다.

세션 데이터 유실은 서비스 품질 저하와 사용자 불만으로 이어질 수 있습니다. 특히 장애 발생 시 세션 클러스터링이 되어 있지 않은 환경에서는 복구가 어렵습니다. 서버 사이징 의사결정 역시 활성 세션 수에만 의존하면 실제 부하와 맞지 않는 결과가 나오기 쉽습니다. 운영자는 경험이나 감에 의존하여 서버 증설 또는 감소를 결정하게 되며, 이는 비용 효율성과 안정성 모두에서 불확실성을

높입니다.

실제 현장에서는 이러한 괴리가 IT 인프라의 과도한 리소스 할당 또는 부족 현상으로 이어질 수 있습니다. 예를 들어, 한 대형 금융기관에서는 활성 세션 수를 기준으로 서버 용량을 산정했으나, 실제로는 사용자 접속 패턴이 시간대별로 크게 달라져 불필요한 리소스가 장기간 유휴 상태로 남는 문제가 발생했습니다. 반대로, 이벤트나 프로모션 등으로 단기간에 동시 접속자가 급증할 때는 활성 세션 수가 실제 부하를 따라가지 못해 장애가 발생하는 사례도 있습니다. 또한, 세션 타임아웃 설정이 비효율적으로 관리될 경우, 불필요하게 오래 유지되는 세션이 시스템 자원을 점유하여 전체 성능 저하로 이어질 수 있습니다. 이처럼 활성 세션 수와 실제 동시 사용자 수의 괴리는 단순한 숫자 차이가 아니라, IT 운영의 전략과 비용, 서비스 품질에 직접적인 영향을 미치는 중요한 이슈임을 알 수 있습니다.

### 1.1.2 수작업 운영 통계의 한계와 Seasonality 대응 부재

기존 WAS 환경에서는 운영자가 직접 SQL 쿼리를 작성하거나, 로그 파일을 수집·분석하여 세션 데이터를 추출하는 방식이 일반적입니다. 이러한 수작업 기반의 운영 통계는 많은 시간과 인적 자원을 소모하며, 쿼리 오류나 데이터 해석 실수로 인해 정확한 통계 산출이 어렵습니다. 특히, 시스템이 복잡해질수록 데이터 소스가 분산되고, 일관된 통계 기준을 유지하기 힘듭니다. 이로 인해 운영팀은 반복적인 데이터 추출과 가공, 보고서 작성에 많은 노력을 들이게 되며, 실시간성이 떨어지는 결과물이 만들어집니다.

운영자가 과거 데이터와 현재 데이터를 수동으로 비교 분석하는 경우, 반복적으로 발생하는 Seasonality(계절성) 패턴을 놓치기 쉽습니다. 예를 들어, 특정 요일이나 시간대에 트래픽이 급증하는 현상, 연말·연시 또는 이벤트 기간에 반복되는 부하 상황 등을 체계적으로 파악하지 못하면 장애 대응이나 리소스 계획이 부정확해집니다. 실제로, 대형 이커머스 기업에서는 블랙프라이데이, 연말 세일 등 특정 시즌에 트래픽이 급증하는데, 과거 데이터를 체계적으로 분석하지 못해 서버 증설 타이밍을 놓치거나, 반대로 불필요하게 과도한 리소스를 할당하는 사례가 빈번하게 발생합니다.

경영진이나 의사결정권자를 위한 보고서(동시접속자 추이, 서버별 부하 등)를 매번 수작업으로 생성해야 하는 현실도 큰 비효율을 초래합니다. 데이터 추출, 가공, 시각화 과정이 반복되면서 운영 팀의 업무 부담이 증가하고, 실시간성이 떨어지는 보고서가 만들어집니다. 이는 신속한 의사결정과 전략 수립에 장애가 됩니다. 또한, 수작업 통계는 데이터의 신뢰성과 일관성을 보장하기 어렵기

때문에, 경영진이 잘못된 정보를 바탕으로 의사결정을 내릴 위험도 존재합니다.

이처럼 수작업 기반의 운영 통계와 Seasonality 대응 부재는 IT 운영의 민첩성과 효율성을 저해하는 주요 원인입니다. 자동화된 통계 수집 및 분석, 그리고 계절성 패턴의 체계적 관리가 이루어지지 않으면, 기업은 예측 불가능한 트래픽 변화에 효과적으로 대응할 수 없으며, 이는 곧 서비스 품질 저하와 비용 증가로 이어질 수 있습니다. 따라서 엔터프라이즈 WAS 환경에서는 운영 통계의 자동화와 Seasonality 분석 역량의 확보가 필수적입니다.

### 1.1.3 장애 원인 추적의 컨텍스트 단절

엔터프라이즈 IT 환경에서는 장애 발생 시 신속하고 정확한 원인 분석이 매우 중요합니다. 그러나 기존 모니터링 시스템에서는 “CPU 부하 급증”, “메모리 사용량 초과” 등 인프라 레벨의 알림만 제공되는 경우가 많습니다. 하지만 이러한 알림만으로는 실제 장애의 근본 원인, 즉 어떤 사용자가 어떤 요청을 했는지 파악하기 어렵습니다. 이는 장애 대응을 위한 정확한 컨텍스트 제공에 한계가 있음을 의미합니다.

세션(사용자 정보), 트랜잭션(요청 내역), 인프라(리소스 사용량) 데이터가 각각 별도의 시스템에서 관리되는 경우, 장애 분석(RCA, Root Cause Analysis) 과정에서 데이터 연계가 단절됩니다. 예를 들어, 특정 사용자가 대량의 데이터를 요청하여 CPU 부하가 발생했다라도, 이를 세션 데이터와 트랜잭션 로그, 인프라 메트릭을 통합하여 분석하기 어렵습니다.

컨텍스트 단절로 인해 장애 발생 시 근본 원인 분석(RCA)이 지연되고, 반복적인 장애가 발생할 수 있습니다. 운영자는 다양한 시스템을 오가며 데이터를 수집하고, 수작업으로 연관성을 찾아야 하므로 대응 시간과 정확성이 떨어집니다. 이는 IT 운영의 효율성 저하와 서비스 품질 악화로 이어집니다.

실제 사례로, 한 금융권 고객사는 대규모 트랜잭션이 발생한 후 시스템 장애가 반복적으로 발생했으나, 인프라 모니터링 시스템에서는 단순히 CPU 사용량 초과만을 알릴 뿐, 어떤 사용자의 어떤 요청이 문제를 일으켰는지 파악하지 못해 근본 원인 분석이 장기간 지연되었습니다. 이 과정에서 운영자는 세션 로그, 트랜잭션 로그, 인프라 메트릭을 각각 별도로 추출하여 수작업으로 매칭해야 했고, 이로 인해 장애 대응 시간이 수 시간에서 수일로 늘어나는 문제가 발생했습니다. 또한, 컨텍스트가 단절된 환경에서는 반복적인 장애가 발생할 때마다 동일한 분석 과정을 반복해야 하므로, 운영팀의 업무 부담이 가중되고 서비스 신뢰도도 저하됩니다.

따라서, 세션, 트랜잭션, 인프라 데이터를 통합적으로 연계하여 장애 원인을 신속하게 분석할 수 있는 체계가 마련되지 않는 한, 엔터프라이즈 WAS 환경에서의 장애 대응과 서비스 품질 개선은 한계에 부딪힐 수밖에 없습니다.

## 1.2 솔루션 기업의 접근: 세션-트랜잭션-AI 통합 플랫폼

솔루션 기업은 엔터프라이즈 WAS 환경에서 발생하는 다양한 세션 관리 과제를 해결하기 위해, 세션-트랜잭션-AI 통합 플랫폼이라는 혁신적인 접근 방식을 도입하였습니다. 이 플랫폼은 IMDG 기반 세션 클러스터링, HyperLogLog 기반 트랜잭션 집계, 그리고 LLM+RAG+MCP 통합 AI 계층 등 최신 기술을 결합하여, 고가용성, 확장성, 실시간 분석, 자연어 질의 기반 운영 자동화 등 다양한 혁신을 실현합니다. 솔루션 기업의 통합 플랫폼은 단순한 세션 관리 도구를 넘어, 엔터프라이즈 IT 운영의 패러다임을 변화시키는 핵심 솔루션으로 자리매김하고 있습니다. 본 절에서는 솔루션 기업의 기술적 접근 방식과 주요 플랫폼 구성 요소, 그리고 이를 통해 실현되는 운영 혁신의 구체적인 내용을 상세히 설명합니다.

### 1.2.1 세션 클러스터링 솔루션의 시작과 배경

오픈소스 미들웨어 기술 지원 경험을 바탕으로, 엔터프라이즈 WAS 환경에서 발생하는 세션 관리 문제를 해결하고자 IMDG(In-Memory Data Grid) 기반 세션 클러스터링 솔루션을 개발했습니다. 기존 WAS 내부 세션 저장 방식은 장애 발생 시 세션 데이터 유실 위험이 크고, 확장성에 한계가 있었습니다. 이에 WAS 외부에 세션을 저장하는 아키텍처를 도입하여 고가용성과 확장성을 확보했습니다.

IMDG는 분산 메모리 구조를 통해 여러 노드에 세션 데이터를 저장하고, 장애 발생 시 자동으로 Failover를 지원합니다. 이를 통해 단일 장애점(SPOF)을 제거하고, 서버 증설·감소 시에도 세션 데이터가 안전하게 유지됩니다. 또한, 세션 클러스터링은 이기종 WAS 간 세션 공유를 가능하게 하여, JBoss EAP, Tomcat, WebLogic 등 다양한 환경에서 통합된 세션 관리가 가능합니다.

세션 클러스터링 솔루션은 WAS 외부에 세션을 저장함으로써, 서비스 중단 없이 서버 리소스 증설·감소가 가능하고, 장애 발생 시에도 사용자 경험을 보장합니다. 세션 클러스터링은 엔터프라이즈 환경에서 필수적인 고가용성, 확장성, 장애 복구 능력을 제공합니다.

IMDG 기반 세션 클러스터링의 도입은 실제 현장에서 다양한 이점을 제공합니다. 예를 들어, 한

대형 제조기업에서는 기존에 WAS 서버 장애 시 세션 데이터가 유실되어 사용자 불만이 빈번하게 발생했으나, 세션 클러스터링 솔루션 도입 이후에는 장애 발생 시에도 세션 데이터가 자동으로 다른 노드로 이관되어 서비스 연속성이 보장되었습니다. 또한, IMDG는 노드 간 데이터 복제를 통해 데이터 일관성을 유지하며, 서버 증설이나 축소 시에도 세션 데이터의 무중단 이전이 가능합니다. 이기종 WAS 환경에서도 세션 클러스터링을 통해 통합 관리가 가능하므로, 기업은 다양한 애플리케이션 서버를 혼합 운영하면서도 일관된 세션 정책을 적용할 수 있습니다. 이러한 아키텍처적 결정은 단순히 기술적 안정성분만 아니라, IT 운영의 민첩성과 비용 효율성, 그리고 서비스 품질 향상까지 동시에 실현할 수 있는 기반을 제공합니다.

## 1.2.2 PromptOps(Prompt Operations)의 정의

PromptOps는 세션 클러스터링 솔루션의 세션 데이터와 APM 솔루션의 트랜잭션 데이터를 LLM(Large Language Model)과 통합하여, 운영자가 자연어로 세션 정보를 조회하고 운영 의사 결정을 지원받는 차세대 운영 패러다임입니다. 기존의 수작업 기반 데이터 추출·분석을 AI 기반 자동화로 전환하여, 운영 효율성과 정확성을 극대화합니다.

PromptOps는 특허 NP25073-KR에 기반한 이중 데이터 소스 구조를 채택합니다. 현재 상태 데이터는 세션 서버(IMDG 기반)에서 실시간으로 조회하고, 과거 이력 데이터는 APM(애플리케이션 성능 모니터링) 시스템에서 분석합니다. 이 구조는 실시간 운영과 과거 트렌드 분석을 동시에 지원하며, 자연어 질의로 “현재 동시접속자 수”, “지난 주 같은 시간대 접속자 수”, “1개월간 통계” 등 다양한 정보를 조회할 수 있습니다.

PromptOps는 LLM과 RAG(Retrieval-Augmented Generation) 기술을 활용하여, 운영자가 복잡한 SQL이나 스크립트 없이 자연어로 세션·트랜잭션 정보를 질의할 수 있습니다. 예를 들어, “서버 부하율 대비 증설이 필요한가?”, “장시간 접속 사용자 목록을 보여줘” 등 다양한 운영 시나리오에 대해 AI가 데이터 기반 답변을 제공합니다.

PromptOps의 가장 큰 특징은 운영 자동화와 의사결정 지원의 혁신입니다. 예를 들어, 운영자는 “최근 1시간 동안 동시접속자 수가 급증한 서버를 알려줘”와 같은 자연어 질의를 통해, 복잡한 쿼리 작성 없이 실시간 데이터를 즉시 확인할 수 있습니다. 또한, PromptOps는 운영 히스토리와 트렌드 데이터를 결합하여, “과거 동일 이벤트 기간과 비교해 올해 트래픽이 얼마나 증가했는가?”와 같은 고차원 분석도 지원합니다. 이중 데이터 소스 구조는 실시간성과 이력성, 두 가지 요구를 모두

충족시키며, LLM+RAG 기술은 방대한 데이터에서 핵심 정보를 추출해 운영자에게 직관적으로 제공합니다. 실제로, PromptOps 방식을 적용한 한 금융기관에서는 운영 통계 생성 시간이 기존 수시간에서 수분 이내로 단축되었으며, 장애 대응 시나리오별로 AI가 자동으로 원인 후보를 제시해 운영자의 업무 효율성이 크게 향상되었습니다. 이처럼 PromptOps는 엔터프라이즈 IT 운영의 패러다임을 자연어 기반 AI 자동화로 전환하는 운영 방법론입니다.

### 1.2.3 VibeOps 패러다임: AI 코파일럿 기반 운영 체계로의 전환

VibeOps는 기존 AIOps(Artificial Intelligence for IT Operations)를 넘어 AI 코파일럿 기반의 지능형 운영 패러다임입니다. AI 코파일럿은 운영자의 질의에 대해 실시간 데이터 분석, 원인 추적, 조치 제안 등 다양한 역할을 수행하며, 운영자는 운영 기준과 자동·승인 대응 경계를 설계하는 “운영 체계 설계자”로 역할이 전환됩니다.

기존에는 운영자가 장애 대응의 최전선에서 직접 데이터를 분석하고 조치를 취해야 했으나, VibeOps 환경에서는 AI가 장애 원인 분석과 대응 방안을 제시합니다. 운영자는 AI가 제안한 조치의 승인 여부, 자동화 정책, SLA 기준 등을 설계·관리하며, 운영 체계의 품질과 효율성을 높입니다.

VibeOps는 세션, 트랜잭션, 인프라 데이터를 통합하여 AI 기반 RCA(Root Cause Analysis)를 자동화합니다. 운영자는 자동 대응 정책과 승인 대응 경계를 설계하고, AI가 반복적인 장애 대응을 자동화합니다. 이는 운영 효율성, 장애 대응 속도, 서비스 품질 모두에서 혁신을 실현합니다.

VibeOps 패러다임의 도입은 운영자의 역할을 단순한 장애 대응자에서 전략적 운영 설계자로 변화시킵니다. 예를 들어, 기존에는 장애 발생 시 운영자가 로그를 직접 분석하고, 각종 데이터를 수집·매칭하여 원인을 추적해야 했으나, VibeOps 환경에서는 AI 코파일럿이 실시간으로 데이터를 분석하여 “특정 서버의 CPU 부하가 급증한 원인은 대량 트랜잭션을 발생시킨 사용자 A의 요청 때문입니다”와 같이 구체적인 원인과 대응 방안을 제시합니다. 운영자는 이러한 AI의 제안을 바탕으로 자동화 정책을 설계하거나, 승인 기반의 대응 경계를 설정할 수 있습니다. 또한, VibeOps는 반복적으로 발생하는 장애 유형에 대해 자동 대응 시나리오를 학습하여, 동일한 문제가 재발할 경우 AI가 즉시 조치를 취하도록 할 수 있습니다. 이로 인해 장애 대응 속도가 획기적으로 단축되고, 운영팀의 업무 부담도 크게 줄어듭니다. 더불어, 서비스 품질과 신뢰성이 향상되어, 기업은 IT 운영의 전략적 가치를 극대화할 수 있습니다. VibeOps는 단순한 자동화를 넘어, AI와 인간 운영자의 협업을 통한 지능형 운영 체계로의 전환을 의미합니다.

## 2장: 세션-트랜잭션-AI 3계층 통합 아키텍처

세션-트랜잭션-AI 3계층 통합 아키텍처는 엔터프라이즈 WAS 환경에서 사용자 세션 관리, 트랜잭션 집계, 그리고 AI 기반 운영 자동화까지 하나의 플랫폼에서 통합적으로 제공하는 구조를 의미합니다. 세션 클러스터링 솔루션은 IMDG 기반 세션 계층을 담당하며, APM 솔루션은 HyperLogLog를 활용한 트랜잭션 계층을, CogentAI/PromptOps는 LLM+RAG+MCP 기반 AI 계층을 구현합니다. 이 장에서는 각 계층의 기술적 원리와 아키텍처, 그리고 경쟁 제품 대비 차별성을 심층적으로 다룹니다.

### 2.1 세션 클러스터링 솔루션: IMDG 기반 세션 계층

세션 클러스터링 솔루션은 In-Memory Data Grid(IMDG) 기반의 세션 저장소를 통해 WAS 환경에서 세션의 고가용성과 확장성을 보장합니다. 세션 데이터를 WAS 외부에 저장함으로써 장애 발생 시에도 세션 유실 없이 Failover가 가능하며, 이기종 WAS 간 세션 공유와 중복 로그인 방지 등 엔터프라이즈 환경에 최적화된 기능을 제공합니다. 또한, 세션 클러스터링 솔루션은 다양한 WAS 환경에서의 유연한 연동성과 운영 효율성을 높이기 위한 다양한 기능을 내장하고 있습니다. 이 계층은 대규모 트래픽 환경에서의 안정적인 세션 관리와 장애 대응, 그리고 운영 편의성 측면에서 기존 솔루션 대비 차별화된 경쟁력을 제공합니다.

#### 2.1.1 IMDG 기반 세션 저장소와 Failover 메커니즘

##### IMDG 아키텍처의 핵심 구조

IMDG(In-Memory Data Grid)는 WAS 환경에서 세션 데이터를 메모리 기반 분산 저장소에 관리하는 방식입니다. 세션 클러스터링 솔루션은 Red Hat Data Grid, Hazelcast, Apache Ignite, Infinispan 등 오픈소스 IMDG 엔진을 활용하거나 자체 엔진을 적용하여, WAS 인스턴스가 직접 세션을 저장하지 않고 외부 IMDG 노드에 세션 정보를 위임합니다. 이 구조는 WAS 서버 장애 시에도 세션 데이터가 유지되어, 사용자의 로그인 상태나 트랜잭션 컨텍스트가 끊임없이 이어질 수 있습니다.

IMDG는 각 노드가 메모리 내에 데이터를 저장하고, 네트워크를 통해 데이터를 실시간으로 동기화합니다. 이를 통해 데이터 접근 속도가 매우 빠르며, 대규모 분산 환경에서도 일관된 세션

관리를 보장할 수 있습니다. 세션 클러스터링 솔루션은 IMDG의 이러한 특성을 활용하여, WAS 서버의 확장이나 축소, 장애 발생 시에도 세션 데이터의 일관성과 가용성을 유지합니다. 또한, IMDG는 데이터 파티셔닝과 복제 기능을 통해 데이터 손실 위험을 최소화하며, 분산 환경에서의 데이터 일관성 정책(Strong Consistency, Eventual Consistency 등)을 유연하게 설정할 수 있습니다.

### Failover 메커니즘과 고가용성 구현

Failover 메커니즘은 IMDG 노드 간의 데이터 복제와 자동 장애 조치(Auto Failover)를 통해 실현됩니다. 세션 데이터는 최소 2개 이상의 노드에 복제되며, 한 노드가 장애를 일으켜도 다른 노드가 즉시 세션을 제공하여 서비스 연속성을 보장합니다. WAS 서버가 재시작되거나 장애가 발생해도, IMDG에 저장된 세션을 불러와 사용자가 끊임 없이 서비스를 이용할 수 있습니다. 이는 Sticky Session 방식의 단일 장애점(SPOF) 문제를 근본적으로 해결합니다.

세션 클러스터링 솔루션의 Failover는 노드 상태를 지속적으로 감시하는 헬스 체크(Health Check) 기능과 연동되어, 장애 발생 시 자동으로 복제된 데이터를 활성 노드로 승격시킵니다. 복구 시간(RTO, Recovery Time Objective)이 짧아, 실제 운영 환경에서 장애로 인한 서비스 중단을 최소화할 수 있습니다. 또한, 데이터 복제 주기와 복제 방식(동기/비동기)을 환경에 맞게 조정할 수 있어, 성능과 데이터 안전성 간의 균형을 맞출 수 있습니다.

### 세션 MBean 모니터링 기능

세션 클러스터링 솔루션은 세션 관리의 관찰 가능성을 강화하기 위해 MBean 기반 모니터링을 제공합니다. 주요 모니터링 항목은 Active 세션 수, 세션 생성/소멸 수, 세션별 메모리 사용량 등입니다. 운영자는 JMX 콘솔 또는 Grafana 대시보드를 통해 실시간으로 세션 상태를 파악할 수 있으며, 장애 조치나 서버 증설·감소 결정에 필요한 데이터를 즉시 확보할 수 있습니다.

MBean은 Java Management Extensions(JMX) 표준을 따르기 때문에, 다양한 운영 도구와 쉽게 연동할 수 있습니다. 예를 들어, 운영자는 특정 시간대의 세션 폭주 현상을 실시간으로 감지하고, 자동 알림을 받아 신속하게 대응할 수 있습니다. 또한, 세션별 상세 정보(예: 사용자별 세션 유지 시간, 세션 속성 값 등)를 조회하여, 비정상적인 세션 증가나 메모리 누수 현상을 조기에 탐지할 수 있습니다. 이러한 모니터링 기능은 운영 효율성뿐 아니라, 장애 예방과 용량 계획(Capacity Planning) 측면에서도 큰 가치를 제공합니다.

### 운영상의 장점과 유의점

IMDG 기반 세션 저장소는 서버 증설·감소에 따른 세션 일관성 문제를 해결하며, 대규모 트래

픽 환경에서 확장성과 안정성을 동시에 제공합니다. 단, IMDG 노드의 메모리 용량과 네트워크 대역폭, 장애 복구 시 데이터 일관성 정책(Strong vs Eventual Consistency) 설정에 유의해야 합니다.

특히, IMDG 노드의 메모리 부족이나 네트워크 장애가 발생할 경우, 세션 데이터의 일부 손실이나 지연이 발생할 수 있으므로, 운영 환경에 맞는 리소스 할당과 네트워크 품질 관리가 중요합니다. 또한, IMDG의 데이터 복제 정책을 적절히 설정하지 않으면, 장애 복구 시 데이터 불일치가 발생할 수 있으므로, 운영 목적에 따라 Strong Consistency와 Eventual Consistency 중 적합한 정책을 선택해야 합니다. 마지막으로, IMDG 클러스터의 노드 수와 복제 인자를 조정하여, 장애 발생 시에도 충분한 데이터 가용성을 확보할 수 있도록 설계하는 것이 중요합니다.

## 2.1.2 이기종 WAS 간 세션 공유와 중복 로그인 방지

### 이기종 WAS 환경에서의 세션 공유

세션 클러스터링 솔루션은 JBoss EAP, Tomcat, WebLogic, Jeus, WebSphere, Resin 등 다양한 WAS 제품을 지원하며, 이기종 WAS 간에도 동일한 IMDG 세션 저장소를 공유할 수 있습니다. 이를 통해 여러 WAS 인스턴스 또는 서로 다른 WAS 제품에서 사용자 세션을 일관되게 관리할 수 있으며, 시스템 확장 또는 마이그레이션 시에도 세션 데이터의 일관성이 유지됩니다.

이러한 이기종 WAS 간 세션 공유 기능은 엔터프라이즈 환경에서 시스템 이중화, 단계적 마이그레이션, 멀티벤더 전략 등에 매우 유용하게 활용됩니다. 예를 들어, 기존에 Tomcat 기반으로 운영하던 시스템을 WebLogic으로 점진적으로 이전할 때, 세션 클러스터링 솔루션을 통해 세션 데이터의 일관성을 유지하면서 무중단 전환이 가능합니다. 또한, 서로 다른 WAS 환경에서 동일한 사용자 경험을 제공할 수 있으므로, 대규모 시스템 통합 프로젝트에서도 안정적인 세션 관리가 가능합니다.

### 중복 로그인 방지 기능

중복 로그인 방지 기능은 동일한 사용자 ID로 여러 곳에서 동시 로그인 시 기존 세션을 자동으로 종료하거나, 새로운 로그인 시도를 차단하는 방식으로 구현됩니다. IMDG 세션 저장소는 사용자별 세션 상태를 실시간으로 확인하여 중복 로그인 탐지 및 제어를 수행합니다. 이는 보안 강화와 라이선스 관리 측면에서 중요한 역할을 합니다.

구체적으로, 세션 클러스터링 솔루션은 세션 생성 시 사용자 ID를 키로 하여 세션 정보를 저

장하고, 새로운 로그인 요청이 들어올 경우 기존 세션의 존재 여부를 검사합니다. 이미 활성화된 세션이 있을 경우, 정책에 따라 기존 세션을 강제로 종료하거나, 새로운 로그인 시도를 거부할 수 있습니다. 이 기능은 금융, 공공, 교육 등 보안 및 라이선스 준수가 중요한 산업에서 널리 활용되며, 불필요한 리소스 소모와 보안 위협을 효과적으로 차단할 수 있습니다.

### 세션 생성 필터링과 불필요한 세션 관리

세션 클러스터링 솔루션은 정적 콘텐츠(예: 이미지, CSS, JS)에 대한 세션 생성을 필터링하여 불필요한 세션 데이터가 생성되는 것을 방지합니다. Servlet 2.5 이상을 지원하는 모든 WAS 환경에서 세션 생성 조건을 세밀하게 설정할 수 있으며, 세션 관리 정책을 통해 메모리 사용량을 최적화할 수 있습니다.

운영자는 세션 생성 필터를 통해, 로그인, 장바구니, 결제 등 실제 사용자 활동에만 세션이 생성되도록 정책을 설정할 수 있습니다. 이를 통해 불필요한 세션 데이터로 인한 메모리 낭비를 줄이고, 시스템의 전체적인 성능을 향상시킬 수 있습니다. 또한, 세션 만료 정책, 세션 갱신 주기, 비정상 세션 탐지 등 다양한 세션 관리 기능을 통해, 대규모 트래픽 환경에서도 안정적인 세션 운영이 가능합니다.

### 지원 범위와 확장성

Servlet 2.5 이상을 지원하는 WAS라면 세션 클러스터링 솔루션과 연동이 가능하며, 신규 WAS 제품이나 클라우드 환경에서도 IMDG 세션 저장소를 활용할 수 있습니다. 이는 엔터프라이즈 환경에서의 유연한 확장과 장기적인 투자 보호를 보장합니다.

특히, 클라우드 네이티브 환경에서는 컨테이너 기반 WAS 인스턴스가 동적으로 생성·소멸될 수 있으므로, 외부 IMDG 세션 저장소를 통한 세션 일관성 유지가 필수적입니다. 세션 클러스터링 솔루션은 Kubernetes, Docker 등 컨테이너 오케스트레이션 환경과도 쉽게 연동할 수 있으며, Auto Scaling, Blue-Green Deployment 등 최신 운영 패턴에서도 안정적인 세션 관리를 제공합니다. 이러한 확장성은 엔터프라이즈 고객이 미래의 IT 환경 변화에 유연하게 대응할 수 있도록 지원합니다.

## 2.2 APM 솔루션: HyperLogLog 기반 트랜잭션 계층

APM 솔루션은 HyperLogLog(HLL) 알고리즘을 활용하여 대규모 트랜잭션 환경에서 동시접속자 집계와 시간 단위 롤업 데이터 관리, 실시간 트랜잭션 모니터링을 제공합니다. Seasonality 패턴

분석과 서버 사이징 의사결정에 필요한 정확한 데이터 기반을 구축합니다. 이 계층은 기존의 단순 카운트 방식이나 샘플링 기반 모니터링과 달리, 대규모 분산 시스템 환경에서도 높은 정확도와 메모리 효율성을 동시에 제공합니다. 또한, 다양한 사용자 식별 모드와 시간 단위 롤업 구조를 통해, 운영자는 장기적인 트래픽 패턴 분석과 실시간 장애 대응을 모두 손쉽게 수행할 수 있습니다.

## 2.2.1 HyperLogLog(HLL) 기반 동시접속자 집계 원리

### HyperLogLog 알고리즘의 원리

HyperLogLog(HLL)는 고유 사용자 수(동시접속자)를 효율적으로 추정하는 확률적 데이터 구조입니다. 사용자 ID(예: IP, 세션ID, 쿠키)를 해시 함수로 이진수 변환한 뒤, 변환 결과에서 선행 0(leading zeros)의 개수를 기록합니다. 여러 해시값에서 가장 긴 선행 0의 개수를 집계하면, 전체 고유 사용자 수를 근사적으로 계산할 수 있습니다.

HLL은 기존의 HashSet이나 Bloom Filter와 달리, 집합의 크기가 매우 커져도 메모리 사용량이 거의 증가하지 않는다는 장점이 있습니다. 예를 들어, 16KB의 메모리만으로도 수억 명의 고유 사용자를 집계할 수 있으며, 이는 대규모 서비스 환경에서 매우 중요한 특성입니다. HLL의 내부 구조는 여러 개의 레지스터로 구성되어 있으며, 각 레지스터는 해시값의 특정 비트를 기준으로 분할되어 데이터를 저장합니다. 이 방식은 메모리 효율성과 집계 속도 모두에서 뛰어난 성능을 보장합니다.

#### 메모리 효율성과 대규모 집계

HLL의 가장 큰 장점은 메모리 사용량이 매우 적다는 점입니다. 16KB 메모리만으로도 수억 명의 고유 사용자 집계가 가능하며, 표준 오차율은 약 0.81%로 서버 사이징 의사결정에는 충분한 정확도를 제공합니다. 이는 기존의 HashSet이나 Bloom Filter 대비 월등한 효율성을 의미합니다.

실제로, HashSet을 사용하여 수백만 명의 고유 사용자를 집계하려면 수백 MB의 메모리가 필요하지만, HLL은 동일한 작업을 수십 KB로 처리할 수 있습니다. 이로 인해, 분산 환경에서 각 WAS 인스턴스가 별도의 HLL 스케치를 유지하더라도, 전체 시스템의 메모리 부담이 거의 발생하지 않습니다. 또한, HLL은 병합(Merge) 연산이 매우 간단하여, 여러 인스턴스의 데이터를 중앙 서버에서 손쉽게 통합할 수 있습니다. 이러한 특성은 대규모 분산 시스템에서의 실시간 집계와 장애 대응에 매우 적합합니다.

#### 분산 환경에서의 HLL 스케치 병합

APM 솔루션은 각 WAS 인스턴스에 HLL 스케치를 생성하고, 분산 환경에서는 이 스케치들을 병합하여 전체 시스템의 중복 사용자 집계를 제거합니다. 병합된 HLL 스케치는 전체 시스템의 동시접속자 수를 정확하게 추정하며, 서버별, 인스턴스별, 전체 집계 모두가 가능합니다.

HLL 스케치의 병합은 각 레지스터의 최대값을 비교하여 새로운 스케치를 생성하는 방식으로 이루어집니다. 이 과정은 연산 비용이 매우 낮아, 실시간 집계가 가능합니다. 예를 들어, 100대 이상의 WAS 인스턴스가 동시에 운영되는 환경에서도, 각 인스턴스의 HLL 데이터를 주기적으로 중앙 서버에 전송하고, 중앙 서버는 이를 병합하여 전체 동시접속자 수를 빠르게 산출할 수 있습니다. 이 방식은 분산 환경에서의 확장성과 집계 정확도를 모두 만족시킵니다.

### 운영상의 활용과 한계

HLL 기반 집계는 실시간 서버 부하 판단, 서버 증설·감소 의사결정, Seasonality 패턴 분석 등 다양한 운영 시나리오에 활용됩니다. 다만, 정밀 과금이나 법적 인증이 필요한 경우에는 오차율을 고려한 추가 검증이 필요합니다.

예를 들어, 대규모 쇼핑몰에서 특정 시간대의 동시접속자 수를 실시간으로 모니터링하여, 트래픽 폭주 시 자동으로 서버를 증설하거나, 장애 발생 시 신속하게 원인 분석을 수행할 수 있습니다. 하지만, HLL은 확률적 추정 방식이기 때문에, 0.81% 내외의 오차가 존재합니다. 따라서, 과금이나 법적 증빙 등 정확한 숫자가 필요한 업무에는 HLL 결과를 참고값으로 활용하고, 추가적인 로그 분석이나 데이터 검증 절차를 병행하는 것이 바람직합니다.

## 2.2.2 시간 단위 롤업 구조와 3가지 사용자 식별 모드

### 시간 단위 롤업 데이터 구조

APM 솔루션은 트랜잭션 데이터를 2초→1분→5분→1시간 단위로 롤업(roll-up)하여 저장합니다. 이 구조는 대용량 데이터의 장기 보관과 Seasonality 패턴 분석에 최적화되어 있습니다. 운영자는 특정 시간대의 동시접속자 수, 트랜잭션 부하, 장애 패턴을 장기간에 걸쳐 분석할 수 있습니다.

롤업 구조는 원시 데이터의 저장 공간을 효율적으로 관리하면서, 장기적인 트렌드 분석과 피크 타임 예측에 매우 유리합니다. 예를 들어, 실시간 모니터링을 위해 2초 단위의 세밀한 데이터를 수집하되, 일정 시간이 지나면 1분, 5분, 1시간 단위로 데이터를 집계하여 저장함으로써, 저장 공간을 절약하고 분석 효율성을 높일 수 있습니다. 이러한 롤업 데이터는 운영자가 월간, 분기별,

연간 트래픽 패턴을 손쉽게 파악하고, 리소스 계획이나 장애 예방 전략을 수립하는 데 큰 도움이 됩니다.

### 3가지 사용자 식별 모드의 특성

APM 솔루션은 사용자 식별을 위해 3가지 모드를 제공합니다:

- Mode 0: IP 주소 기반 식별(간단하지만 NAT/프록시 환경에서 정확도 저하)
- Mode 1: JSESSIONID 기반 식별(세션 만료 시 중복 가능)
- Mode 2: KHANUSER 쿠키 기반 식별(쿠키 비활성화 시 제한)

각 모드는 환경에 따라 정확도와 활용성이 다르며, 운영 목적에 따라 적합한 모드를 선택해야 합니다.

Mode 0(IP 기반)은 설정이 간단하고 별도의 추가 작업 없이 바로 적용할 수 있지만, 여러 사용자가 동일한 NAT 또는 프록시 서버를 통해 접속하는 환경에서는 실제 사용자 수보다 적게 집계될 수 있습니다. Mode 1(JSESSIONID 기반)은 세션 단위로 사용자를 식별하므로, 세션이 만료되거나 브라우저가 재시작될 경우 동일 사용자가 중복 집계될 수 있습니다. Mode 2(KHANUSER 쿠키 기반)은 고유 쿠키를 활용하여 가장 정확한 사용자 식별이 가능하지만, 일부 사용자가 쿠키를 비활성화한 경우에는 집계에서 누락될 수 있습니다.

### 정확도 차이와 선택 기준

IP 기반은 빠르고 간단하지만 NAT, 프록시 환경에서는 여러 사용자가 동일 IP로 집계될 수 있습니다. JSESSIONID 기반은 세션 만료 시 중복이 발생할 수 있으며, KHANUSER 쿠키 기반은 가장 높은 정확도를 제공하지만, 쿠키 비활성화 환경에서는 제한이 있습니다. 운영 목적에 따라 모드를 조합하거나, 보안 정책과 사용자 환경을 고려해 선택해야 합니다.

예를 들어, 사내망이나 폐쇄망 환경에서는 IP 기반 식별만으로도 충분할 수 있지만, 외부 고객이 다수 접속하는 서비스에서는 KHANUSER 쿠키 기반 식별을 권장합니다. 또한, 운영자는 모드별 집계 결과를 비교 분석하여, 실제 사용자 수와 추정치 간의 차이를 파악하고, 환경 변화에 따라 식별 정책을 유연하게 조정할 수 있습니다. 필요에 따라, 여러 모드를 동시에 적용하여 복수의 집계 결과를 병렬로 관리하는 것도 가능합니다.

### 롤업 데이터의 운영적 가치

시간 단위 롤업 데이터는 주기적 피크 타임, 장애 패턴, 리소스 계획 등 Seasonality 대응에 필수적인 정보를 제공합니다. 운영자는 자연어 질의로 “1개월간 1시간 기준 통계 정보”를 조회하며,

데이터 기반 의사결정을 수행할 수 있습니다.

이러한 롤업 데이터는 단순한 트래픽 집계뿐만 아니라, 장애 발생 시점과 패턴 분석, 리소스 증설·감소 시뮬레이션, SLA(서비스 수준 협약) 준수 여부 확인 등 다양한 운영 시나리오에 활용됩니다. 예를 들어, 특정 요일이나 시간대에 반복적으로 트래픽이 급증하는 현상을 롤업 데이터를 통해 사전에 감지하고, 자동 스케일링 정책을 적용하여 장애를 예방할 수 있습니다. 또한, 장기적인 데이터 분석을 통해, 서비스 성장 추세나 계절별 트래픽 변동 등 비즈니스 인사이트를 도출할 수 있습니다.

### 2.2.3 100% 실시간 트랜잭션 모니터링과 OpenTelemetry 통합

#### 실시간 트랜잭션 모니터링 구조

APM 솔루션은 모든 트랜잭션을 100% 실시간으로 모니터링합니다. WAS별 에이전트가 트랜잭션 데이터를 실시간으로 수집하며, 장애 발생 시 즉시 원인 분석과 대응이 가능합니다. 기존 샘플링 방식 대비 완전한 데이터 수집으로 장애 근본 원인 분석(RCA)과 운영 보고서 생성에 강점을 갖습니다.

실시간 모니터링은 트랜잭션의 시작부터 종료까지의 모든 이벤트를 기록하며, 각 트랜잭션의 처리 시간, 오류 발생 여부, 호출 경로 등을 상세히 추적합니다. 이를 통해, 운영자는 장애 발생 시점과 원인을 신속하게 파악할 수 있으며, 서비스 성능 저하나 비정상 트랜잭션을 조기에 탐지하여 선제적으로 대응할 수 있습니다. 또한, 모든 트랜잭션 데이터를 저장하므로, 사후 분석이나 감사(Audit) 용도로도 활용이 가능합니다.

#### OpenTelemetry 표준 기반 통합 관측성

APM 솔루션은 OpenTelemetry 표준을 기반으로 세션, 트랜잭션, 분산 트레이스, 로그 데이터를 하나의 통합 스키마로 관리합니다. 이는 Prometheus, Grafana, Jaeger 등 기존 모니터링 스택과의 연동을 가능하게 하며, 클라우드 네이티브 환경에서의 관찰 가능성을 극대화합니다.

OpenTelemetry는 다양한 언어와 플랫폼을 지원하는 오픈소스 관측성 표준으로, 분산 시스템에서의 트레이스, 메트릭, 로그 데이터를 일관된 방식으로 수집하고 처리할 수 있습니다. APM 솔루션은 OpenTelemetry Collector와 연동하여, 외부 시스템과의 데이터 교환, 커스텀 메트릭 수집, 알림 연동 등 다양한 확장 기능을 제공합니다. 이를 통해, 운영자는 단일 대시보드에서 전체 시스템의 상태를 종합적으로 모니터링할 수 있습니다.

### 세션-트랜잭션-로그 통합 아키텍처

세션, 트랜잭션, 로그 데이터는 통합된 데이터 모델로 관리되어, 운영자는 자연어 질의로 “어떤 사용자가 어떤 URL을 호출했는가?” 를 즉시 분석할 수 있습니다. 장애 대응, 보안 분석, SLA 준수 여부 확인 등 다양한 운영 시나리오에 활용됩니다.

이 통합 아키텍처는 데이터 간의 상관관계를 손쉽게 분석할 수 있도록 설계되어 있습니다. 예를 들어, 특정 사용자의 세션에서 발생한 모든 트랜잭션과 로그 이벤트를 한눈에 파악할 수 있으며, 장애 발생 시점의 전체 시스템 상태를 신속하게 재구성할 수 있습니다. 또한, 자연어 질의 인터페이스를 통해, 비전문가도 복잡한 데이터 분석을 손쉽게 수행할 수 있습니다. 이는 운영 효율성과 데이터 기반 의사결정 품질을 크게 향상시킵니다.

### 운영상의 장점과 확장성

실시간 트랜잭션 모니터링과 OpenTelemetry 통합은 엔터프라이즈 환경에서 장애 대응 속도와 데이터 신뢰성을 크게 향상시킵니다. 클라우드, 온프레미스, 하이브리드 환경 모두에서 확장성이 뛰어나며, 기존 모니터링 시스템과의 공존이 가능합니다.

APM 솔루션은 다양한 배포 환경에서 유연하게 적용할 수 있으며, 기존의 Prometheus, Datadog, New Relic 등과도 연동이 가능합니다. 또한, API 기반 확장 기능을 통해, 맞춤형 대시보드, 자동화된 알림, 외부 시스템 연동 등 다양한 운영 요구사항을 충족할 수 있습니다. 이러한 확장성은 엔터프라이즈 고객이 변화하는 IT 환경에 신속하게 대응할 수 있도록 지원합니다.

## 2.3 CogentAI / PromptOps: LLM+RAG+MCP 기반 AI 계층

CogentAI(Cowork Agent AI)는 기업이나 기관에서 함께 일하는 AI Agent 제품에 대한 브랜드입니다. LLM, RAG, MCP를 결합하여 IT 운영자가 자연어 대화를 통해 시스템 진단 및 자동화된 조치를 수행할 수 있도록 지원합니다.

CogentAI와 PromptOps는 LLM(대형 언어 모델), RAG(검색 증강 생성), MCP(모델 컨텍스트 프로토콜) 통합 아키텍처를 통해 운영자가 자연어로 세션·트랜잭션 정보를 조회하고, AI 기반 운영 자동화를 실현합니다. 실시간 데이터와 과거 이력 분석을 동시에 제공하며, 개인정보 보호와 할루시네이션 최소화 기능을 내장합니다. 이 계층은 운영자의 업무 효율성을 극대화하고, 데이터 기반 의사결정의 신뢰도를 높이며, 엔터프라이즈 환경에서 요구되는 보안 및 컴플라이언스 요건을 충족할 수 있도록 설계되어 있습니다.

## 2.3.1 하이브리드 LLM과 이중 데이터 소스 아키텍처

### 하이브리드 LLM 동적 선택 메커니즘

CogentAI는 온프레미스 LLM, GPT-4, Claude, Gemma3 등 다양한 LLM을 하이브리드로 동적 선택합니다. 질의의 복잡도, 응답 품질, 비용, 언어 특성에 따라 최적의 LLM을 자동으로 선택하여 자연어 질의에 응답합니다. 이는 운영자의 요구에 맞는 맞춤형 AI 코파일럿을 제공하며, 복수 LLM 조합으로 할루시네이션 위험을 최소화합니다.

하이브리드 LLM 구조는 단일 모델의 한계를 극복하고, 다양한 업무 시나리오에 최적화된 응답을 제공할 수 있도록 설계되었습니다. 예를 들어, 기술적 질의에는 GPT-4를, 비용이 중요한 단순 질의에는 Gemma3를, 한국어 특화 질의에는 KoGPT와 같은 언어 특화 모델을 선택적으로 활용할 수 있습니다. 이 과정은 AI 오케스트레이터가 자동으로 질의 유형을 분석하고, 각 LLM의 강점과 약점을 고려하여 최적의 조합을 결정합니다. 또한, 복수 LLM의 응답을 교차 검증하거나, 앙상블 방식으로 결합하여 응답 품질과 신뢰성을 더욱 높일 수 있습니다.

### 이중 데이터 소스 구조의 구현

CogentAI는 특허 기반 이중 데이터 소스 구조를 적용합니다. 현재 상태는 세션서버에서 실시간으로 조회하고, 과거 이력은 APM에서 시계열 데이터를 분석합니다. 이 구조는 실시간 운영 데이터와 장기 이력 분석을 동시에 제공하여, 장애 대응, 리소스 계획, 보고서 생성 등 다양한 운영 시나리오에 최적화되어 있습니다.

이중 데이터 소스 구조는 데이터 신뢰성과 분석 품질을 동시에 확보할 수 있는 핵심 아키텍처입니다. 실시간 데이터 소스는 현재 시스템 상태, 동시접속자 수, 트랜잭션 부하 등 즉각적인 정보 제공에 최적화되어 있으며, 이력 데이터 소스는 장기적인 트렌드 분석, 장애 패턴 탐지, 리소스 증설·감소 의사결정 등에 활용됩니다. CogentAI는 두 데이터 소스를 자동으로 결합하여, 운영자가 자연어로 요청한 복합 질의(예: “현재와 과거의 트래픽 패턴 비교”)에 대해 신속하고 정확한 답변을 제공합니다.

### 실시간 조회와 이력 분석의 통합

운영자는 자연어 질의로 “현재 동시접속자 수와 지난 주 같은 요일, 같은 시간대의 접속자 수를 비교해줘”와 같은 요청을 할 수 있으며, CogentAI는 세션서버와 APM의 데이터를 결합하여 즉시 응답합니다. 이중 데이터 소스 구조는 데이터 신뢰성과 분석 품질을 크게 향상시킵니다.

이 통합 구조는 데이터 소스 간의 시계열 동기화, 데이터 정합성 검증, 응답 시간 최적화 등 다양

한 기술적 과제를 해결하여, 운영자가 복잡한 데이터 분석을 손쉽게 수행할 수 있도록 지원합니다. 또한, 실시간 데이터와 이력 데이터를 동시에 활용함으로써, 장애 조기 탐지, 리소스 최적화, SLA 준수 여부 확인 등 다양한 운영 시나리오에서 높은 실효성을 발휘합니다.

### 운영 자동화와 AI 코파일럿의 가치

하이브리드 LLM과 이중 데이터 소스 구조는 운영자가 장애 대응, 리소스 증설·감소, 보고서 생성 등 운영 의사결정을 자연어 한 줄로 수행할 수 있게 하며, AI가 복잡한 데이터 분석과 조치 제안을 자동으로 제공합니다.

AI 코파일럿은 단순 질의 응답을 넘어, 운영자가 직면한 문제 상황을 분석하고, 최적의 대응 방안을 제안하거나, 반복적인 운영 업무를 자동화할 수 있습니다. 예를 들어, 트래픽 급증 시 자동으로 서버 증설을 제안하거나, 장애 발생 시 원인 분석과 조치 절차를 단계별로 안내할 수 있습니다. 이러한 AI 기반 운영 자동화는 운영자의 업무 부담을 줄이고, 서비스 안정성과 효율성을 동시에 향상시킵니다.

## 2.3.2 MCP 프로토콜 기반 시스템 연동과 개인정보 보호

### MCP(Model Context Protocol) 표준 연동 구조

MCP는 파일 시스템, DB, ERP, 그룹웨어 등 다양한 사내 시스템과 CogentAI를 표준화된 방식으로 연동하는 프로토콜입니다. 각 시스템의 데이터 컨텍스트를 정의하고, LLM이 자연어 질의에 필요한 데이터를 자동으로 참조할 수 있게 합니다. REST API, OpenTelemetry, JDBC 등 다양한 인터페이스와 호환됩니다.

MCP는 데이터 소스별로 메타데이터와 접근 권한, 데이터 스키마를 정의하여, LLM이 질의 목적에 맞는 데이터를 안전하게 참조할 수 있도록 지원합니다. 예를 들어, ERP 시스템의 인사 데이터, DB의 트랜잭션 로그, 그룹웨어의 일정 정보 등을 MCP 표준에 따라 연동하면, 운영자는 자연어로 “최근 1주일간 장애 발생과 인력 투입 현황을 분석해줘”와 같은 복합 질의도 손쉽게 수행할 수 있습니다. MCP는 데이터 보안과 접근 제어 기능도 내장하고 있어, 민감 정보의 무단 접근을 방지할 수 있습니다.

### 한국어 개인정보 자동 탐지 및 마스킹

CogentAI는 한국어 개인정보(이름, 주민번호, 전화번호 등)를 자동으로 탐지하고, 마스킹 기능을 제공합니다. 운영자는 개인정보 보호 정책에 따라 데이터 출력 시 자동으로 민감 정보를

숨길 수 있으며, GDPR, 개인정보보호법 등 컴플라이언스 요구사항을 충족합니다.

이 기능은 LLM이 자연어 질의 결과를 생성할 때, 개인정보 패턴을 실시간으로 탐지하고, 사전에 정의된 마스킹 정책(예: 이름의 일부만 표시, 주민번호 뒷자리 숨김 등)을 적용합니다. 운영자는 마스킹 수준을 정책에 따라 세밀하게 조정할 수 있으며, 로그 데이터, 트랜잭션 내역 등 다양한 데이터 소스에 일관되게 적용할 수 있습니다. 이를 통해, 개인정보 유출 위험을 최소화하고, 외부 감사나 법적 요구사항에도 효과적으로 대응할 수 있습니다.

### RAG 기반 할루시네이션 최소화 메커니즘

CogentAI는 RAG(검색 증강 생성) 기술을 활용하여 내부 문서, 운영 데이터, 시스템 로그를 실시간으로 참조합니다. LLM의 할루시네이션(허위 생성) 위험을 최소화하며, 운영자는 AI 답변의 근거 데이터를 즉시 확인할 수 있습니다. 완전한 제거는 불가능하지만, 데이터 기반 교차 검증과 운영 기준 설정으로 신뢰성을 높입니다.

RAG는 LLM이 답변을 생성할 때, 사전에 인덱싱된 신뢰할 수 있는 데이터 소스에서 관련 정보를 검색하여, 답변의 근거로 활용합니다. 이 과정에서, LLM이 근거 없는 정보를 임의로 생성하는 할루시네이션 현상을 줄이고, 실제 운영 데이터에 기반한 신뢰성 높은 응답을 제공합니다. 운영자는 AI가 참조한 데이터 소스, 문서, 로그 내역 등을 함께 확인할 수 있어, 답변의 신뢰도를 직접 검증할 수 있습니다. 또한, 운영 기준에 따라 RAG의 검색 범위와 우선순위를 조정하여, 업무 목적에 맞는 응답 품질을 확보할 수 있습니다.

### 운영 데이터 신뢰성과 확장성

MCP 기반 연동은 다양한 사내 시스템과의 통합을 가능하게 하며, 운영 데이터의 신뢰성과 확장성을 동시에 제공합니다. AI 기반 운영 자동화, 보고서 생성, 장애 분석 등 모든 운영 시나리오에서 개인정보 보호와 데이터 신뢰성이 보장됩니다.

CogentAI는 MCP를 통해 신규 시스템이나 외부 데이터 소스와의 연동도 손쉽게 확장할 수 있으며, 데이터 소스 추가 시에도 일관된 보안 정책과 데이터 정합성을 유지할 수 있습니다. 이러한 확장성은 엔터프라이즈 환경에서 다양한 비즈니스 요구와 IT 인프라 변화에 유연하게 대응할 수 있도록 지원합니다. 또한, 운영자는 AI 기반 자동화 기능을 통해, 반복적인 보고서 생성, 장애 분석, 보안 점검 등 다양한 업무를 효율적으로 수행할 수 있습니다.

## 2.4 경쟁 제품 대비 아키텍처 차별성

솔루션 기업 통합 플랫폼은 세션-트랜잭션-AI 네이티브 통합 구조와 HyperLogLog 기반 동시접속자 집계, 자연어 세션 질의(PromptOps) 등 고유 기능으로 경쟁 제품 대비 차별성을 확보합니다. Redis+별도APM+ChatGPT, Hazelcast+Datadog 등 기존 조합과 비교하여, 통합성, 확장성, AI 운영 지원 측면에서 우위를 설명합니다. 이 절에서는 주요 경쟁 조합과 솔루션 기업 통합 플랫폼의 아키텍처적 차별성을 구체적으로 비교하여, 엔터프라이즈 환경에서의 실질적인 운영 효율성과 미래 확장성 측면에서 솔루션 기업의 강점을 강조합니다.

### 2.4.1 Redis+별도APM vs Hazelcast+Datadog vs 솔루션 기업 통합 플랫폼 비교

#### 세션 클러스터링 아키텍처 비교

Redis 기반 세션 클러스터링은 단순하고 빠르지만, 장애 복구와 데이터 일관성 측면에서 한계가 있습니다. Hazelcast는 IMDG 기반으로 세션 복제와 장애 조치가 가능하지만, 별도의 APM과 연동이 필요합니다. 솔루션 기업은 IMDG 기반 세션 저장소와 트랜잭션 집계, AI 운영 지원까지 네이티브 통합 플랫폼으로 제공하여, 장애 대응과 확장성에서 우위를 점합니다.

Redis는 단일 마스터 구조에서 장애 발생 시 복구 시간이 길어질 수 있으며, 세션 데이터의 일관성 보장을 위해 별도의 복제 및 장애 조치 구성이 필요합니다. Hazelcast는 IMDG의 특성상 데이터 복제와 장애 조치가 내장되어 있으나, 트랜잭션 집계와 AI 연동을 위해서는 별도의 외부 솔루션이 필요합니다. 반면, 솔루션 기업 통합 플랫폼은 세션, 트랜잭션, AI 계층이 하나의 통합 아키텍처로 설계되어, 장애 발생 시 자동 복구와 데이터 일관성 유지, 운영 자동화까지 원스톱으로 제공합니다.

#### APM 통합과 동시접속자 집계 방식

Redis+별도APM 조합은 동시접속자 집계에 HashSet 또는 단순 카운트 방식을 사용하며, 대규모 환경에서 메모리 사용량과 정확도 문제가 발생합니다. Hazelcast+Datadog 조합은 분산 트레이스와 모니터링은 강점이 있으나, HyperLogLog 기반 동시접속자 집계는 제공하지 않습니다. 솔루션 기업은 HyperLogLog 알고리즘으로 수억 명의 동시접속자를 16KB 메모리로 집계하며, 분산 환경에서 중복 제거와 시간 단위 롤업 데이터까지 제공합니다.

솔루션 기업의 HyperLogLog 기반 집계는 대규모 트래픽 환경에서도 높은 정확도와 메모리

효율성을 보장하며, 분산 환경에서의 데이터 병합과 중복 제거가 자동으로 이루어집니다. 이는 기존 HashSet 방식의 메모리 한계, 단순 카운트 방식의 중복 집계 문제를 근본적으로 해결합니다. 또한, 시간 단위 롤업 구조를 통해 장기적인 트래픽 패턴 분석과 리소스 계획에도 최적화되어 있습니다.

### AI 운영 지원과 자연어 세션 질의

ChatGPT 연동 방식은 세션 데이터와 트랜잭션 데이터를 별도로 관리하며, 자연어 질의의 품질과 데이터 신뢰성에 한계가 있습니다. PromptOps는 세션-트랜잭션-AI 통합 구조에서 자연어 질의로 세션 정보, 동시접속자 수, 장애 원인 분석, 보고서 생성까지 자동화합니다. 이는 운영자의 업무 효율성과 데이터 기반 의사결정 품질을 크게 향상시킵니다.

PromptOps는 LLM, RAG, MCP 기반의 통합 아키텍처를 통해, 운영자가 복잡한 데이터 분석을 자연어 한 줄로 수행할 수 있도록 지원합니다. 또한, 실시간 데이터와 이력 데이터를 결합하여, 장애 대응, 리소스 최적화, 보안 분석 등 다양한 운영 시나리오에서 높은 신뢰성과 효율성을 제공합니다. 경쟁 조합의 경우, 데이터 소스 간의 연동과 자연어 질의 품질이 제한적이기 때문에, 엔터프라이즈 환경에서의 운영 자동화 수준이 솔루션 기업에 비해 낮을 수밖에 없습니다.

### 통합성, 확장성, 운영 자동화 측면의 차별성

솔루션 기업 통합 플랫폼은 세션 클러스터링, 트랜잭션 집계, AI 운영 자동화까지 하나의 플랫폼에서 제공하며, Kubernetes, OpenTelemetry, REST API 등 클라우드 네이티브 환경과 완벽하게 연동됩니다. 경쟁 조합 대비 통합성과 확장성, 운영 자동화 품질에서 명확한 차별성을 갖습니다.

솔루션 기업은 단일 플랫폼 내에서 모든 계층의 데이터와 기능을 통합 관리할 수 있으므로, 운영자는 복잡한 시스템 연동이나 데이터 정합성 문제 없이, 신속하게 새로운 기능을 도입하고 확장할 수 있습니다. 또한, 클라우드 네이티브 환경에서의 자동 확장, 멀티 클러스터 관리, API 기반 연동 등 최신 IT 트렌드에 부합하는 아키텍처를 제공합니다. 이러한 통합성과 확장성, 운영 자동화 품질은 엔터프라이즈 고객이 미래의 비즈니스 변화에 신속하게 대응할 수 있도록 지원하는 핵심 경쟁력입니다.

## 3장: PromptOps 핵심 활용 시나리오

PromptOps는 솔루션 기업의 세션-트랜잭션-AI 통합 플랫폼에서 자연어 기반 질의와 AI 분석을 통해 WAS 환경의 운영 효율성을 극대화하는 핵심 도구입니다. 이 장에서는 PromptOps를 활용

한 동시접속자 기반 서버 사이징, Seasonality 패턴 분석, 장애 예방, 보고서 자동 생성, 그리고 데이터 정확도와 운영 시 주의사항까지 다양한 실무 시나리오를 다룹니다. 실제 운영 환경에서 PromptOps가 어떻게 서버 증설·감소 의사결정, 주기적 부하 대응, 장애 근본 원인 분석, IT 의사결정자 보고서 자동화, 그리고 데이터 신뢰성 확보에 기여하는지 구체적으로 설명합니다.

## 3.1 동시접속자 기반 서버 사이징 의사결정

동시접속자 기반 서버 사이징은 엔터프라이즈 WAS 환경에서 가장 중요한 운영 의사결정 중 하나입니다. PromptOps는 실시간 동시접속자 수와 과거 동일 시간대의 접속자 수 비교, 서버 증설·감소 판단, 그리고 Kubernetes HPA 연동 자동 스케일링까지 데이터 기반의 의사결정을 지원합니다. 이 섹션에서는 자연어 질의를 통한 실시간 서버 부하 판단, 데이터 기반 증설 판단, 그리고 자동화된 스케일링 구현 방법을 상세히 다룹니다. 또한, 각 기능이 실제 운영 환경에서 어떻게 적용되고, 기존 방식 대비 어떤 장점과 유의점이 있는지 구체적으로 설명하여, 운영자가 실질적으로 활용할 수 있는 실무적 통찰을 제공합니다.

### 3.1.1 실시간 동시접속자 조회와 과거 동일 시간대 비교

#### 실시간 동시접속자 수 조회 방법

PromptOps를 활용하면 운영자는 “현재 동시접속자 수를 알려줘”와 같은 자연어 질의로 실시간 서버 부하를 즉시 파악할 수 있습니다. APM 솔루션은 HyperLogLog(HLL) 알고리즘을 기반으로 각 WAS 인스턴스의 고유 사용자 수를 빠르고 정확하게 집계합니다. 이 데이터는 IMDG 기반 세션 클러스터와 연동되어 서버 장애 시에도 세션 데이터를 잃지 않고, 실시간 집계가 가능합니다. 운영자는 웹 대시보드 혹은 CLI에서 PromptOps 프롬프트를 입력하여 현재 동시접속자 수, 서버별 분포, 메모리 사용량 등 다양한 지표를 즉시 확인할 수 있습니다.

실시간 동시접속자 조회는 운영자가 시스템의 현재 상태를 빠르게 파악하고, 예기치 않은 부하 급증이나 장애 징후를 조기에 감지하는 데 매우 효과적입니다. 예를 들어, 대형 이벤트나 프로모션이 진행되는 시점에서 예상보다 많은 사용자가 접속할 경우, PromptOps를 통해 즉각적으로 동시접속자 수를 확인하고, 필요 시 서버 증설이나 리소스 재분배 등의 조치를 신속하게 내릴 수 있습니다. 또한, 실시간 데이터는 운영팀의 의사소통에도 중요한 역할을 하며, 각 부서 간 협업 시 신뢰할 수 있는 근거 자료로 활용됩니다.

## 과거 동일 시간대 접속자 수 비교

PromptOps는 “지난 주 같은 요일, 같은 시간대의 동시접속자 수와 비교해줘”와 같은 질의도 지원합니다. APM 솔루션의 시계열 데이터와 시간 단위 롤업(2초→1분→5분→1시간) 구조를 활용하여, 과거 데이터와 현재 데이터를 자동으로 비교 분석합니다. 이 기능은 Seasonality 패턴 분석과 장애 예방에 필수적입니다. 운영자는 24시간 동안의 동시접속자 추이 그래프, 피크 타임, 평균 부하 등 다양한 통계 정보를 자연어로 요청할 수 있으며, AI가 자동으로 시각화된 결과를 제공합니다.

과거와 현재의 동시접속자 수를 비교함으로써, 운영자는 시스템 부하의 정상 범위와 이상 징후를 명확하게 구분할 수 있습니다. 예를 들어, 평소보다 20% 이상 높은 부하가 감지되면, 이는 이벤트 효과일 수도 있지만, 비정상적인 트래픽이나 공격의 신호일 수도 있습니다. PromptOps는 이러한 비교 결과를 표, 그래프 등으로 시각화하여 한눈에 파악할 수 있도록 하며, 운영자는 데이터 기반으로 신속하게 원인 분석과 대응 방안을 마련할 수 있습니다.

## 프롬프트 예시와 결과 해석

실무에서는 다음과 같은 프롬프트가 활용됩니다:

- “24시간 동안 동시접속자 추이 그래프를 보여줘”
- “현재 부하가 지난 주 같은 시간대 대비 얼마나 증가했는지 알려줘”

PromptOps는 결과를 표, 그래프, 요약 통계 등 다양한 형태로 반환하며, 운영자는 이를 바탕으로 서버 증설·감소, 장애 대응, SLA 준수 여부 등을 판단할 수 있습니다. 결과 해석 시에는 HyperLogLog의 오차율(약 0.81%)을 감안하여 데이터 기반 의사결정이 이루어집니다.

실제 운영 사례에서는, 예를 들어 월요일 오전 9시에 동시접속자 수가 평소 대비 30% 이상 증가한 경우, PromptOps의 비교 분석 결과를 근거로 즉각적인 서버 증설이나 트래픽 분산 조치를 결정할 수 있습니다. 또한, AI가 제공하는 요약 통계와 시각화 자료는 경영진이나 IT 의사결정자에게 신속하게 보고할 수 있는 자료로도 활용됩니다.

## 장점과 유의사항

PromptOps의 실시간 동시접속자 조회 기능은 기존의 수작업 SQL 쿼리, 로그 분석 대비 압도적으로 빠르고 정확합니다. 하지만 NAT, 프록시 환경에서는 IP 기반 집계 오차가 발생할 수 있으므로, 사용자 식별 모드(JSESSIONID, KHANUSER 등)를 적절히 선택해야 합니다. 또한,

과거 데이터와의 비교 시에는 시스템 환경 변화(서버 증설, 로드밸런서 변경 등)를 고려하여 해석해야 합니다.

운영자는 실시간 데이터의 신뢰성을 높이기 위해, 집계 방식과 사용자 식별 모드의 특성을 충분히 이해하고, 필요 시 추가적인 로그 분석이나 트래픽 샘플링을 병행해야 합니다. 특히, 대규모 이벤트나 시스템 구조 변경이 있었던 경우에는 단순 수치 비교보다는 맥락을 고려한 해석이 필요 합니다. 이러한 유의사항을 준수하면, PromptOps를 통한 동시접속자 기반 서버 사이징과 장애 예방이 더욱 효과적으로 이루어질 수 있습니다.

### 3.1.2 서버 증설·감소 의사결정을 위한 데이터 기반 근거

#### HyperLogLog 기반 동시접속자 집계

APM 솔루션은 HyperLogLog(HLL) 알고리즘을 활용하여 동시접속자 수를 메모리 효율적으로 집계합니다. HLL은 사용자 ID를 해시로 변환하고, 선행 0의 개수로 고유 사용자 수를 추정합니다. 16KB 메모리로 수억 명의 사용자를 0.81% 오차로 집계할 수 있으며, 각 WAS 인스턴스의 HLL 스케치를 병합하여 전체 시스템 중복을 제거합니다. 이 방식은 감에 의존하던 기존 서버 사이징 의사결정에서 데이터 기반 판단으로 전환하는 핵심 기술입니다.

HyperLogLog의 도입으로 인해, 기존에는 운영자의 경험이나 예측에 의존하던 서버 증설·감소 판단이 과학적이고 객관적인 데이터에 기반하게 되었습니다. 예를 들어, 대규모 쇼핑몰에서는 이벤트 기간 동안 동시접속자 수가 급격히 증가할 수 있는데, HLL 기반 집계는 이러한 변화에 신속하게 대응할 수 있는 근거 데이터를 제공합니다. 또한, 분산 환경에서 여러 서버의 데이터를 통합하여 중복 없이 집계할 수 있기 때문에, 전체 시스템의 실제 부하를 정확하게 파악할 수 있습니다.

#### 자연어 질의와 AI 답변 프로세스

운영자는 “서버 3대 기준, 현재 부하율 대비 증설이 필요한가?”와 같은 질의를 PromptOps에 입력할 수 있습니다. AI는 실시간 동시접속자 수, 서버별 부하율, CPU/메모리 사용량, 과거 피크 타임 데이터를 종합 분석하여 증설·감소 필요성을 데이터 기반으로 답변합니다. AI는 서버별 세션 분포, 부하 임계치, SLA 준수 여부 등을 자동으로 계산하여 운영자에게 근거 있는 의사결정 정보를 제공합니다.

이 과정에서 AI는 단순히 현재 수치만을 보여주는 것이 아니라, 과거 트렌드와 예측 모델을 결합하여 향후 부하 변화까지 예측할 수 있습니다. 예를 들어, “현재 부하가 곧 임계치에 도달할

것으로 예상되니, 서버를 2대 증설하는 것이 적절하다”는 식의 구체적인 조언을 제공합니다. 이러한 AI 기반 분석은 운영자의 의사결정을 신속하고 정확하게 만들어주며, 불필요한 리소스 낭비나 장애 위험을 최소화할 수 있습니다.

### 데이터 기반 의사결정의 장점

데이터 기반 서버 증설·감소 의사결정은 감이나 경험에 의존하던 기존 방식 대비 객관적이고 신뢰성 높은 결과를 제공합니다. HyperLogLog의 오차율이 서버 사이징에는 미미하므로, 실시간 부하 변화에 따른 신속한 대응이 가능합니다. 운영자는 AI가 제시하는 근거와 통계 정보를 바탕으로 서버 증설, 감소, 리소스 재분배 등 다양한 조치를 효율적으로 수행할 수 있습니다.

실제 사례로, 금융권에서는 월말 급여 이체나 대규모 결제 이벤트 시, PromptOps의 데이터 기반 분석을 통해 사전에 서버 증설을 결정하고, 장애 없이 트래픽을 처리한 경험이 있습니다. 이처럼 데이터 기반 의사결정은 SLA 준수와 고객 만족도 향상에 직접적으로 기여합니다.

### 유의사항과 한계

HyperLogLog 기반 집계는 과금 시스템 등 정밀한 사용자 수 집계에는 부적합할 수 있습니다. 또한, 사용자 식별 모드(IP, JSESSIONID, KHANUSER 등)의 선택에 따라 집계 정확도가 달라지므로, 운영 환경에 맞는 모드 선택이 중요합니다. AI 답변은 데이터 기반이지만, 시스템 환경 변화나 특이 패턴(예: 대량 로그인 이벤트) 발생 시에는 추가 검증이 필요합니다.

운영자는 데이터 기반 의사결정의 한계를 명확히 인식하고, 필요 시 추가적인 로그 분석, 트래픽 샘플링, 전문가의 경험적 판단을 병행해야 합니다. 특히, 예외적 상황이나 시스템 구조 변경이 있었던 경우에는 AI의 분석 결과를 맹신하지 않고, 교차 검증을 통해 최종 결정을 내려야 합니다.

## 3.1.3 Kubernetes HPA 연동을 통한 자동 스케일링

### HPA Custom Metric 연동 구조

APM 솔루션은 Kubernetes Horizontal Pod Autoscaler(HPA)와 연동하여 동시접속자 수 기반 자동 스케일링 정책을 구현할 수 있습니다. APM에서 집계한 동시접속자 수를 Custom Metric으로 Kubernetes에 전달하면, HPA가 실시간 부하에 따라 Pod 수를 자동으로 조정합니다. 이 구조는 기존의 수동 서버 사이징에서 완전 자동화된 스케일링으로 발전하는 핵심 경로입니다.

Kubernetes HPA는 기본적으로 CPU, 메모리 사용량을 기준으로 Pod 수를 조정하지만, APM 솔루션과의 연동을 통해 실제 동시접속자 수와 같은 비표준(Custom) 메트릭을 활용할 수

있습니다. 이를 통해, 단순 리소스 사용량이 아닌 실제 사용자 부하에 맞춘 스케일링 정책을 설계할 수 있으며, 이는 서비스 품질 유지와 비용 효율성 측면에서 큰 이점을 제공합니다.

### 스케일링 정책 설계 방법

운영자는 PromptOps에서 “동시접속자 수가 1,000명을 넘으면 Pod를 2개 추가해줘”와 같은 자연어 질의로 스케일링 정책을 설계할 수 있습니다. AI는 과거 부하 패턴, 현재 트렌드, 서버별 리소스 사용량을 분석하여 최적의 스케일링 임계치와 정책을 제안합니다. HPA는 CPU, 메모리뿐 아니라 동시접속자 Custom Metric을 기준으로 Pod 수를 조정하므로, 실제 사용자 부하에 맞는 리소스 할당이 가능합니다.

실제 적용 시에는, 예를 들어 “동시접속자 수가 2,000명을 초과하면 Pod를 3개까지 확장하고, 500명 미만으로 감소하면 1개로 축소”와 같은 세부 정책을 설계할 수 있습니다. AI는 과거 이벤트, 트래픽 변화, 장애 이력 등을 참고하여, 불필요한 스케일링이나 리소스 낭비를 최소화하는 최적의 정책을 추천합니다.

### 자동화의 장점과 실무 적용 사례

자동 스케일링은 피크 타임, 이벤트, 장애 상황에서 신속한 리소스 확장·축소를 가능하게 하며, 운영자의 수동 개입을 최소화합니다. APM 솔루션과 Kubernetes HPA 연동은 대형 금융, 공공 기관, 이커머스 환경에서 이미 적용되고 있으며, SLA 준수, 비용 효율성, 장애 예방 측면에서 높은 효과를 보이고 있습니다.

예를 들어, 대형 이커머스 사이트에서는 블랙프라이데이와 같은 대규모 트래픽 이벤트 시, 자동 스케일링 정책을 통해 수십 대의 Pod가 자동으로 증설되어 장애 없이 트래픽을 처리한 사례가 있습니다. 이처럼 자동화된 스케일링은 운영자의 부담을 줄이고, 예기치 않은 트래픽 급증에도 유연하게 대응할 수 있도록 지원합니다.

### 유의사항과 확장성

HPA 연동 시에는 Custom Metric의 집계 주기, 오차율, 사용자 식별 모드에 따라 스케일링 반응성이 달라질 수 있습니다. 또한, Pod 수 증설 시 세션 클러스터링(세션 클러스터링 솔루션)과 연동하여 세션 데이터 유실 방지, 중복 로그인 방지 등 추가 설정이 필요합니다. 운영자는 PromptOps와 APM, Kubernetes의 통합 구조를 이해하고, 실무 환경에 맞게 정책을 지속적으로 개선해야 합니다.

특히, 스케일링 정책의 임계치 설정이 지나치게 민감하거나 둔감할 경우, 불필요한 리소스 증설 또는 서비스 품질 저하가 발생할 수 있습니다. 따라서, 실제 트래픽 패턴과 비즈니스 요구사항을

충분히 반영하여 정책을 설계하고, 운영 중에도 지속적으로 모니터링 및 튜닝을 수행하는 것이 중요합니다.

## 3.2 Seasonality 패턴 분석과 장애 예방

Seasonality 패턴 분석과 장애 예방은 반복적 부하, 장애 상황을 사전에 예측하고 대응하는 운영의 핵심입니다. PromptOps는 시간축 롤업 데이터 기반 패턴 탐지, 과거 동일 시기 대비 접속 패턴 비교, 그리고 사용자-요청-인프라 연결 추적을 통한 장애 근본 원인 분석(RCA)을 지원합니다. 이 섹션에서는 AI와 시계열 데이터, 트랜잭션 분석을 활용한 실무 시나리오를 상세히 설명합니다. 또한, 실제 운영 환경에서 Seasonality 패턴이 어떻게 나타나며, 이를 통해 장애를 예방하고 리소스를 효율적으로 관리할 수 있는지 구체적인 사례와 방법론을 중심으로 다룹니다.

### 3.2.1 시간축 롤업 데이터 기반 주기적 패턴 탐지

#### 시간 단위 롤업 데이터 구조

APM 솔루션은 트랜잭션 데이터를 2초→1분→5분→1시간 단위로 롤업하여 저장합니다. 이 구조는 대량의 세션·트랜잭션 데이터를 효율적으로 집계하고, 장기간의 부하 패턴 분석에 최적화되어 있습니다. PromptOps는 “1개월간 1시간 기준 통계 정보 조회”와 같은 자연어 질의를 통해 주기적 피크 타임, 평균 부하, 이상 패턴 등을 자동으로 탐지합니다.

시간축 롤업 데이터 구조의 가장 큰 장점은, 대규모 트랜잭션 데이터도 장기간에 걸쳐 손쉽게 분석할 수 있다는 점입니다. 예를 들어, 1년치 트랜잭션 로그를 1시간 단위로 롤업하면, 수십억 건의 데이터를 수천 건으로 요약할 수 있어, 운영자는 장기적인 트렌드와 반복 패턴을 빠르게 파악할 수 있습니다. 또한, 롤업 데이터는 저장 공간을 크게 절약하면서도, 주요 이상 패턴이나 피크 타임을 놓치지 않고 탐지할 수 있도록 도와줍니다.

#### Seasonality 패턴 분석 방법

운영자는 PromptOps를 활용하여 “최근 1개월간 피크 타임을 알려줘”, “평일과 주말의 부하 패턴을 비교해줘”와 같은 질의를 입력할 수 있습니다. AI는 롤업된 시계열 데이터를 분석하여 반복적으로 발생하는 부하, 장애, 이벤트 패턴을 자동으로 시각화하고, 예측 정보를 제공합니다. 이 기능은 금융, 교육, 이커머스 등 반복적 이벤트가 많은 환경에서 장애 예방과 리소스 확보에 필수적입니다.

예를 들어, 교육기관에서는 매주 월요일 오전에 수강신청 트래픽이 반복적으로 증가하는 패턴이 관찰될 수 있습니다. PromptOps는 이러한 Seasonality 패턴을 자동으로 탐지하여, 운영자가 사전에 서버 증설이나 트래픽 분산 정책을 준비할 수 있도록 지원합니다. 또한, AI는 과거 데이터와 현재 트렌드를 결합하여, 향후 발생 가능한 피크 타임이나 이상 부하를 예측하는 데도 도움을 줍니다.

### 패턴 탐지의 실무 활용

실무에서는 주기적 피크 타임(예: 점심시간, 야간, 이벤트 기간) 탐지, 이상 부하 패턴(예: 갑작스런 트래픽 폭증) 분석, SLA 준수 여부 모니터링 등에 활용됩니다. PromptOps는 운영자가 반복되는 Seasonality 패턴을 놓치지 않고, 사전에 리소스를 확보하거나 장애 대응 계획을 수립할 수 있도록 지원합니다.

특히, 이벤트나 프로모션 기간과 같이 트래픽이 급증하는 시점에는, 과거 패턴을 참고하여 적절한 리소스 확보와 장애 예방 조치를 미리 준비할 수 있습니다. 또한, 반복적으로 발생하는 이상 부하나 장애 패턴을 조기에 탐지함으로써, 서비스 중단이나 SLA 위반을 효과적으로 방지할 수 있습니다.

### 장점과 유의사항

시간축 롤업 데이터 기반 패턴 분석은 수작업 로그 분석 대비 압도적으로 빠르고 정확합니다. 하지만 롤업 주기 설정, 데이터 집계 방식에 따라 세밀한 패턴 탐지가 제한될 수 있으므로, 운영 환경에 맞는 롤업 구간 선택이 중요합니다. 또한, 이벤트 발생 시 과거 데이터와의 비교를 통해 장애 예방 조치를 신속하게 수행해야 합니다.

운영자는 롤업 데이터의 한계를 인식하고, 필요 시 원본 데이터나 세부 로그를 추가로 분석하여, 세밀한 이상 패턴이나 예외 상황을 놓치지 않도록 해야 합니다. 또한, 롤업 주기와 집계 방식이 실제 비즈니스 요구와 일치하는지 주기적으로 점검하는 것이 중요합니다.

## 3.2.2 과거 동일 시기 대비 접속 패턴 비교와 사전 대응

### 과거-현재 비교 분석 구조

PromptOps는 “작년 같은 시기(블랙프라이데이, 수강신청) 대비 현재 접속 패턴을 비교해줘”와 같은 자연어 질의를 지원합니다. APM 솔루션의 시계열 데이터와 CogentAI의 LLM+RAG 기반 분석 기능을 결합하여, 과거와 현재의 접속 패턴, 부하, 장애 발생 빈도 등을 자동으로 비교합니다.

AI는 반복되는 이벤트, 장애 패턴을 탐지하고, 사전 대응 방안을 제안합니다.

이러한 비교 분석 구조는, 단순히 현재 트래픽만을 모니터링하는 것이 아니라, 과거 유사 이벤트와의 차이점을 명확하게 파악할 수 있도록 도와줍니다. 예를 들어, 블랙프라이데이와 같은 대규모 이벤트에서, 작년과 올해의 트래픽 패턴, 장애 발생 시점, 리소스 사용량 등을 비교함으로써, 올해 발생할 수 있는 위험 요인을 사전에 식별할 수 있습니다.

### 사전 대응 시나리오

운영자는 PromptOps를 통해 “작년 수강신청 기간과 올해의 부하 패턴을 비교해줘”, “블랙프라이데이 이벤트 대비 리소스 확보가 충분한지 알려줘”와 같은 질의를 입력할 수 있습니다. AI는 과거 데이터와 현재 트렌드를 분석하여, 리소스 증설, 장애 예방, SLA 준수 여부 등 사전 대응 조치를 자동으로 제안합니다.

실제 적용 사례로, 대학교의 수강신청 시스템에서는 매년 반복되는 트래픽 급증 시점을 미리 예측하고, PromptOps의 비교 분석 결과를 바탕으로 서버 증설과 네트워크 용량 확보를 사전에 완료하여, 장애 없이 서비스를 제공한 경험이 있습니다. 이처럼 과거-현재 비교 분석은 반복적 이벤트 대응에 매우 효과적입니다.

### 실무 적용과 장점

과거-현재 비교 분석은 반복적 장애, 부하 상황을 사전에 예측하고 대응하는 데 매우 효과적입니다. PromptOps와 CogentAI의 통합 분석은 운영자가 수작업으로 과거 로그, 통계 데이터를 비교하던 비효율을 완전히 제거하며, 신속한 의사결정과 리소스 확보에 기여합니다.

특히, 반복적으로 발생하는 장애나 트래픽 급증 이벤트에 대해, 사전에 충분한 대비책을 마련할 수 있으므로, 서비스 품질 유지와 SLA 준수에 큰 도움이 됩니다. 또한, AI가 자동으로 분석 결과를 시각화하고, 주요 차이점과 위험 요인을 요약해주기 때문에, 운영자의 업무 효율성이 크게 향상됩니다.

### 유의사항과 한계

과거-현재 비교 시 시스템 환경 변화(서버 증설, 정책 변경 등)를 반드시 고려해야 하며, 데이터의 정확도와 집계 방식에 따라 분석 결과가 달라질 수 있습니다. AI의 자동 분석 결과는 운영자의 경험과 교차 검증을 통해 최종 의사결정에 활용해야 합니다.

운영자는 단순 수치 비교에 의존하지 않고, 시스템 구조, 정책, 사용자 행동 패턴 등 다양한 맥락을 함께 고려해야 하며, 필요 시 추가적인 로그 분석이나 전문가의 판단을 병행해야 합니다. 또한, 데이터 품질과 집계 방식의 일관성을 유지하는 것이 중요합니다.

### 3.2.3 장애 근본 원인 분석(RCA): 사용자-요청-인프라 연결 추적

#### 세션-트랜잭션-인프라 연동 구조

PromptOps는 CPU 부하 급증, 메모리 사용량 폭증 등 인프라 알림 발생 시 “어떤 사용자가 어떤 URL을 호출했는가?” 를 자연어 질의로 추적할 수 있습니다. 세션 클러스터링 솔루션의 세션 데이터, APM의 트랜잭션 데이터, 인프라 리소스 모니터링 정보를 하나의 연속된 맥락으로 연결하여, 장애 근본 원인 분석(RCA)을 지원합니다.

이 연동 구조를 통해, 기존에는 각각 분리되어 있던 세션, 트랜잭션, 인프라 데이터를 통합적으로 분석할 수 있게 되었습니다. 예를 들어, 특정 시간대에 CPU 사용량이 급증한 원인을 파악할 때, 해당 시간대에 어떤 사용자가 어떤 요청을 보냈는지, 그 요청이 어떤 트랜잭션으로 이어졌는지, 그리고 그 트랜잭션이 인프라 자원에 어떤 영향을 미쳤는지를 한눈에 추적할 수 있습니다.

#### RCA 프로세스와 AI 분석

운영자는 “CPU 부하가 급증한 시간대에 가장 많이 호출된 URL과 사용자 목록을 알려줘” 와 같은 질의를 PromptOps에 입력합니다. AI는 세션(사용자)↔트랜잭션(요청)↔인프라(리소스) 데이터를 결합하여, 장애 발생 시점의 주요 트랜잭션, 사용자 행동, 리소스 사용 패턴을 자동 분석합니다. RCA 결과는 표, 그래프, 요약 통계 등 다양한 형태로 제공되며, 운영자는 근본 원인에 신속하게 접근할 수 있습니다.

실제 사례로, 대형 공공기관의 민원 시스템에서 특정 시간대에 CPU 부하가 급증한 원인을 분석한 결과, 특정 사용자가 반복적으로 대용량 파일 다운로드 요청을 보낸 것이 원인임을 PromptOps의 RCA 분석을 통해 신속하게 파악할 수 있었습니다. 이를 바탕으로, 해당 사용자의 접근을 제한하고, 시스템 설정을 조정하여 장애를 예방할 수 있었습니다.

#### 실무 적용과 장점

RCA 기능은 기존의 분리된 모니터링 시스템(IP, 로그, 트랜잭션, 인프라 등) 대비 압도적으로 빠르고 정확한 장애 분석을 제공합니다. PromptOps와 CogentAI의 통합 분석은 장애 발생 시 원인 사용자, 요청 URL, 리소스 사용 패턴을 자동으로 추적하여, 신속한 대응과 조치에 기여합니다.

이러한 통합 분석은 장애 발생 시 운영자의 스트레스를 크게 줄여주며, 반복적 장애의 근본 원인을 빠르게 파악하여 재발 방지 대책을 수립하는 데도 큰 도움이 됩니다. 또한, RCA 결과는 IT 의사결정자에게 신뢰할 수 있는 근거 자료로 활용될 수 있습니다.

#### 유의사항과 한계

RCA 분석은 데이터의 정확도, 집계 방식, 사용자 식별 모드에 따라 결과가 달라질 수 있습니다. 또한, AI 분석 결과는 운영자의 경험과 추가 검증을 통해 최종 조치에 활용해야 하며, 시스템 환경 변화(정책 변경, 서버 증설 등)를 반드시 고려해야 합니다.

운영자는 RCA 결과를 맹신하지 않고, 필요 시 원본 로그나 추가 데이터를 참고하여 교차 검증을 수행해야 하며, 시스템 구조나 정책 변경이 있었던 경우에는 분석 결과의 해석에 각별히 주의해야 합니다.

### 3.3 IT 의사결정자를 위한 보고서 자동 생성

PromptOps는 IT 의사결정자가 자연어 질의 한 줄로 동시접속자 보고서, 서버별 세션 분포, SLA 준수 여부 등 다양한 운영 보고서를 자동 생성할 수 있도록 지원합니다. 이 섹션에서는 보고서 자동 생성 시나리오와 사용자 행동 분석, 비정상 접속 탐지 등 보안 관점의 활용까지 상세히 설명합니다. 또한, 실제 IT 의사결정자들이 PromptOps를 통해 어떻게 신속하고 정확한 보고서를 받아보고, 이를 기반으로 전략적 의사결정을 내릴 수 있는지 구체적인 활용 사례와 주의사항을 함께 다룹니다.

#### 3.3.1 자연어 질의 기반 동시접속자 보고서 생성

##### 보고서 자동 생성 구조

PromptOps는 “이번 주 동시접속자 최대/최소/평균을 보고서 형태로 만들어줘”와 같은 자연어 질의를 지원합니다. AI는 APM 솔루션의 시계열 데이터, 서버별 세션 분포, SLA 준수 여부 등을 자동으로 집계하여, 주간/월간 동시접속자 추이 보고서, 서버별 세션 분포 현황, SLA 모니터링 보고서를 생성합니다.

보고서 자동 생성 구조는 IT 의사결정자가 복잡한 데이터 추출이나 수작업 보고서 작성 없이, 자연어 질의 한 줄만으로 원하는 형태의 보고서를 신속하게 받아볼 수 있도록 설계되었습니다. 예를 들어, “지난 달 서버별 세션 분포와 SLA 준수 현황을 표와 그래프로 요약해줘”와 같은 요청을 하면, AI가 자동으로 데이터를 집계하고, 시각화 자료와 함께 보고서를 완성해줍니다. 이 과정에서 PDF, Excel, HTML 등 다양한 포맷으로 결과를 받을 수 있어, 내부 회의나 외부 보고에 즉시 활용할 수 있습니다.

##### 보고서 생성 프롬프트 예시

실무에서는 다음과 같은 프롬프트가 활용됩니다:

- “월간 동시접속자 추이 그래프와 최대/최소/평균 값을 보고서로 만들어줘”
- “서버별 세션 분포 현황을 표로 보여줘”
- “SLA 준수 여부를 모니터링 보고서로 만들어줘”

PromptOps는 결과를 PDF, Excel, HTML 등 다양한 형태로 자동 생성하며, IT 의사결정자는 신속하게 보고서를 받아볼 수 있습니다.

이러한 프롬프트는 반복적인 보고서 작성 업무를 획기적으로 단축시켜주며, 운영팀과 경영진 간의 커뮤니케이션을 원활하게 해줍니다. 또한, AI가 자동으로 데이터의 이상치나 특이점도 함께 분석하여 보고서에 포함시킬 수 있으므로, 의사결정자는 단순 수치뿐만 아니라 인사이트까지 함께 얻을 수 있습니다.

### 장점과 실무 적용

보고서 자동 생성 기능은 기존의 “운영팀 요청→데이터 추출→보고서 작성” 프로세스를 자연어 한 줄로 단축시킵니다. IT 의사결정자는 실시간 데이터 기반 보고서를 즉시 받아볼 수 있으며, 서버 투자, 리소스 확보, SLA 준수 등 다양한 의사결정에 활용할 수 있습니다.

실제 사례로, 대기업의 IT 부서에서는 매주 월요일 오전마다 PromptOps를 통해 주간 동시접속자 추이와 SLA 준수 현황을 자동으로 보고서로 받아, 경영진 회의에 즉시 활용하고 있습니다. 이처럼 자동화된 보고서 생성은 업무 효율성 향상과 신속한 의사결정에 큰 기여를 하고 있습니다.

### 유의사항과 한계

보고서 자동 생성 시 데이터의 정확도, 집계 방식, 사용자 식별 모드에 따라 결과가 달라질 수 있습니다. AI가 자동 생성한 보고서는 운영자의 경험과 추가 검증을 통해 최종 의사결정에 활용해야 하며, 시스템 환경 변화(서버 증설, 정책 변경 등)를 반드시 고려해야 합니다.

운영자는 보고서의 주요 수치와 통계가 실제 시스템 상태를 정확히 반영하는지 주기적으로 검증해야 하며, 필요 시 추가적인 데이터 분석이나 전문가의 의견을 참고해야 합니다.

## 3.3.2 사용자 행동 분석과 비정상 접속 탐지

### 사용자 행동 분석 구조

PromptOps는 “현재 접속한 사용자 아이디 목록”, “guest1로 접속한 사용자가 조회한 URL 리스트”, “1시간 이상 접속 중인 사용자 목록” 등 다양한 자연어 질의를 지원합니다. AI는 세션 데이터, 트랜잭션 로그, 접속 시간, URL 패턴 등을 자동 분석하여, 사용자 행동, 장시간 접속,

비정상 접속 패턴을 탐지합니다.

사용자 행동 분석 구조는 단순히 접속자 수를 집계하는 것을 넘어, 각 사용자의 세부 행동 패턴까지 심층적으로 분석할 수 있도록 설계되었습니다. 예를 들어, 특정 사용자가 비정상적으로 많은 요청을 보내거나, 장시간 동안 동일한 세션을 유지하는 경우, PromptOps는 이를 자동으로 탐지하여 운영자에게 경고를 제공합니다. 또한, URL 패턴 분석을 통해, 정상적인 사용 흐름과 다른 비정상 요청이나 잠재적 공격 시도를 신속하게 파악할 수 있습니다.

### 비정상 접속 탐지 시나리오

운영자는 PromptOps를 통해 “비정상 URL 패턴을 탐지해줘”, “장시간 접속 중인 사용자 목록을 알려줘” 등 보안 관점의 질의를 입력할 수 있습니다. AI는 트랜잭션 로그, 세션 데이터, 접속 시간, URL 패턴을 분석하여, 비정상 접속, 장시간 접속, 대량 로그인 등 다양한 보안 이벤트를 자동으로 탐지합니다.

예를 들어, 대형 금융기관에서는 PromptOps를 활용하여, 단기간 내에 동일 IP에서 다수의 로그인 시도가 발생하는 경우를 자동으로 탐지하고, 잠재적 보안 위협에 신속하게 대응할 수 있었습니다. 또한, 장시간 접속이나 비정상적인 URL 접근이 반복되는 사용자를 조기에 식별하여, 계정 도용이나 데이터 유출 위험을 사전에 차단할 수 있습니다.

### 실무 적용과 장점

사용자 행동 분석과 비정상 접속 탐지 기능은 보안 관리자, DevOps 엔지니어가 실시간으로 이상 패턴을 탐지하고, 신속하게 대응할 수 있도록 지원합니다. PromptOps와 CogentAI의 통합 분석은 기존의 수작업 로그 분석 대비 압도적으로 빠르고 정확한 보안 이벤트 탐지를 제공합니다.

특히, 대규모 시스템에서는 수작업으로 모든 로그를 분석하는 것이 사실상 불가능하므로, AI 기반 자동 분석이 필수적입니다. PromptOps는 실시간 알림, 시각화, 자동 보고서 생성 등 다양한 기능을 통해, 운영자의 업무 효율성과 보안 수준을 동시에 향상시켜줍니다.

### 유의사항과 한계

사용자 행동 분석, 비정상 접속 탐지 시 데이터의 정확도, 집계 방식, 사용자 식별 모드에 따라 결과가 달라질 수 있습니다. AI 분석 결과는 운영자의 경험과 추가 검증을 통해 최종 조치에 활용해야 하며, 시스템 환경 변화(정책 변경, 서버 증설 등)를 반드시 고려해야 합니다.

운영자는 AI가 탐지한 이상 패턴이나 보안 이벤트가 실제 위협인지, 혹은 정상적인 사용자의 오탐지인지 신중하게 판단해야 하며, 필요 시 추가적인 로그 분석이나 보안 전문가의 의견을 참고해야 합니다.

## 3.4 사용 시 주의사항과 데이터 정확도

PromptOps와 APM 솔루션, CogentAI를 활용할 때는 HyperLogLog 오차율, 사용자 식별 모드, LLM 할루시네이션, 운영 데이터 신뢰성 등 다양한 주의사항을 반드시 고려해야 합니다. 이 섹션에서는 데이터 정확도와 운영 관점의 주요 주의사항을 상세히 설명합니다. 또한, 실제 운영 환경에서 발생할 수 있는 데이터 품질 이슈와 AI 분석의 한계, 그리고 이를 극복하기 위한 실무적 대응 방안을 구체적으로 안내합니다.

### 3.4.1 HyperLogLog 오차율과 사용자 식별 모드별 주의점

#### HyperLogLog 오차율 특성

HyperLogLog(HLL)는 약 0.81%의 표준 오차율로 동시접속자 집계를 수행합니다. 이 오차율은 서버 사이징, 장애 예방 등 운영 의사결정에는 미미하지만, 정밀한 과금 시스템이나 법적 사용자 수 집계에는 부적합할 수 있습니다. 운영자는 HLL의 오차율 특성을 이해하고, 데이터 기반 의사결정에 적절히 활용해야 합니다.

실제 운영 환경에서는, 예를 들어 10만 명의 동시접속자를 집계할 때 약 800명의 오차가 발생할 수 있습니다. 이는 서버 증설이나 장애 예방에는 큰 영향을 미치지 않지만, 과금이나 법적 기준이 엄격한 환경에서는 추가적인 보정이나 다른 집계 방식을 병행해야 할 수 있습니다.

#### 사용자 식별 모드별 주의사항

APM 솔루션은 3가지 사용자 식별 모드(Mode 0: IP 주소, Mode 1: JSESSIONID, Mode 2: KHANUSER 쿠키)를 제공합니다.

- IP 주소 모드는 NAT, 프록시 환경에서 중복 집계 오차가 발생할 수 있습니다.
- JSESSIONID 모드는 세션 만료 시 중복 집계가 발생할 수 있습니다.
- KHANUSER 쿠키 모드는 쿠키 비활성화 환경에서 집계가 제한될 수 있습니다.

운영자는 환경에 맞는 모드 선택과 오차 보정 방법을 반드시 고려해야 합니다.

예를 들어, 대기업의 사내망이나 공공기관처럼 NAT 환경이 많은 경우에는 IP 기반 집계보다는 JSESSIONID나 KHANUSER 쿠키 기반 집계를 사용하는 것이 더 정확할 수 있습니다. 반대로, 쿠키 사용이 제한된 환경에서는 IP 기반 집계가 불가피할 수 있으므로, 각 환경에 맞는 최적의 식별 모드를 선택해야 합니다.

### 실무 적용과 장점

HyperLogLog 기반 집계는 대규모 분산 환경에서 메모리 효율적이고 빠른 집계를 가능하게 하며, 운영자의 감이나 경험에 의존하던 기존 방식 대비 신뢰성 높은 데이터 기반 의사결정을 제공합니다.

실제 사례로, 대형 이커머스 사이트에서는 수백만 명의 동시접속자를 HyperLogLog 기반으로 실시간 집계하여, 서버 증설 및 장애 예방에 성공적으로 활용하고 있습니다. 이처럼 HLL은 대규모 환경에서 특히 강력한 성능을 발휘합니다.

### 유의사항과 한계

오차율, 사용자 식별 모드의 한계는 운영 환경에 따라 달라질 수 있습니다. 운영자는 PromptOps와 APM의 집계 방식, 오차 특성을 이해하고, 최종 의사결정에 추가 검증을 반드시 수행해야 합니다.

특히, 법적 책임이 따르거나 과금 기준이 엄격한 환경에서는, HyperLogLog 집계 결과만을 근거로 삼지 말고, 필요 시 원본 로그나 추가적인 집계 방식을 병행해야 합니다. 또한, 운영 환경 변화(정책 변경, 서버 증설 등)가 있을 때마다 집계 방식의 적합성을 재점검하는 것이 중요합니다.

## 3.4.2 LLM 할루시네이션 대응과 운영 데이터 신뢰성

### RAG 기반 할루시네이션 최소화

CogentAI는 RAG(Retrieval-Augmented Generation) 기술을 활용하여 내부 문서를 실시간 참조하고, LLM 할루시네이션(근거 없는 답변)을 최소화합니다. RAG는 LLM이 실시간 데이터, 문서, 로그를 참조하여 정확한 답변을 생성하도록 지원하며, 운영 데이터 신뢰성을 높입니다.

RAG 기반 분석은 AI가 단순히 학습된 지식에만 의존하지 않고, 최신 운영 데이터와 문서를 실시간으로 참고하므로, 답변의 정확도와 신뢰성이 크게 향상됩니다. 예를 들어, 시스템 구조가 변경되거나 새로운 정책이 도입된 경우에도, RAG는 최신 문서를 반영하여 올바른 답변을 제공할 수 있습니다.

### 할루시네이션 대응 방법

할루시네이션 완전 제거는 불가능하므로, 운영자는 AI 답변을 교차 검증하고, 주요 의사결정에는 데이터 기반 근거와 추가 검증을 반드시 수행해야 합니다. PromptOps는 운영자가 AI 답변의 근거 데이터, 출처, 통계 정보를 확인할 수 있도록 지원하며, 신뢰성 높은 운영 환경을 제공합니다.

실제 운영에서는, AI가 제공한 답변이 실제 시스템 상태와 일치하는지, 혹은 근거가 명확한지 반드시 확인해야 하며, 필요 시 원본 데이터나 전문가의 의견을 참고해야 합니다. 또한, PromptOps는 AI 답변의 근거가 되는 데이터와 문서의 출처를 함께 제공하여, 운영자가 신뢰도를 직접 판단할 수 있도록 도와줍니다.

### 메모리 관리와 세션 생성 필터링

운영 데이터 신뢰성 확보를 위해 불필요한 세션 생성 필터링, 메모리 관리, 세션 클러스터링 등 다양한 운영 관점의 주의사항을 반드시 고려해야 합니다. PromptOps는 세션 생성 필터링, 중복 로그인 방지, 메모리 사용량 모니터링 등 다양한 기능을 제공하여, 운영 데이터의 품질과 신뢰성을 높입니다.

예를 들어, 불필요한 세션이 반복적으로 생성되는 경우, 메모리 사용량이 급증하여 시스템 성능 저하나 장애로 이어질 수 있습니다. PromptOps는 이러한 세션 생성 패턴을 자동으로 탐지하고, 운영자에게 알림을 제공하여, 사전에 문제를 예방할 수 있도록 지원합니다.

### 유의사항과 한계

AI 답변, RAG 기반 분석 결과는 운영자의 경험과 추가 검증을 통해 최종 의사결정에 활용해야 하며, 시스템 환경 변화(정책 변경, 서버 증설 등)를 반드시 고려해야 합니다. 운영 데이터 신뢰성 확보를 위해 지속적인 모니터링, 검증, 정책 개선이 필요합니다.

운영자는 AI와 RAG의 한계를 명확히 인식하고, 필요 시 수동 검증이나 추가 데이터 분석을 병행해야 하며, 시스템 구조나 정책 변경이 있을 때마다 데이터 품질과 신뢰성을 재점검해야 합니다. 또한, 운영 데이터의 신뢰성을 높이기 위해, 주기적인 모니터링과 정책 개선을 지속적으로 수행하는 것이 중요합니다.

## 4장: 도입 사례와 사용자 그룹별 활용 가치

본 장에서는 PromptOps의 실제 적용 사례와 다양한 사용자 그룹별 활용 가치를 심층적으로 분석한다. 엔터프라이즈 환경에서 PromptOps가 어떻게 세션-트랜잭션-LLM 통합 운영을 실현하며, 장애 대응과 리소스 최적화, 보고서 자동화 등 실무적 효익을 제공하는지 구체적으로 설명한다. 또한 CTO, IT Director, 기술 기획팀장, WAS 운영자, DevOps 엔지니어, 보안 관리자 등 각 역할별로 PromptOps 방식이 제공하는 핵심 가치를 다각도로 조명한다.

## 4.1 PromptOps 적용 사례

PromptOps는 다양한 산업군에서 실질적인 운영 효율성과 장애 대응 능력을 입증해왔습니다. 실제 현장에서는 세션-트랜잭션-LLM 통합 운영 시나리오를 통해 복잡한 장애 원인을 신속하게 분석하고, AI 기반 자동화 도구로 장애 대응 시간을 크게 단축하는 등 가시적인 효과를 경험하고 있습니다. 특히 공공기관과 대규모 엔터프라이즈 환경에서 PromptOps의 적용은 운영 자동화와 보고서 생성, 보안 강화 등 여러 측면에서 높은 만족도를 이끌어내고 있습니다. 이 섹션에서는 세션-트랜잭션-LLM 통합 운영 시나리오, AI 기반 장애 분석, 그리고 공공기관 도입 및 고객 만족도 사례를 중심으로 PromptOps의 실무 적용 효과를 구체적으로 소개합니다.

### 4.1.1 세션-트랜잭션-LLM 통합 운영 시나리오

#### CPU 부하 급증 추적 사례

엔터프라이즈 WAS 환경에서 CPU 부하가 급증하는 상황은 빈번하게 발생한다. 기존 모니터링 시스템은 인프라 지표만을 제공해 장애 원인 추적이 어렵지만, PromptOps는 세션-트랜잭션-LLM 통합 구조를 통해 근본 원인 분석을 혁신한다. 예를 들어, CPU 부하 알림이 발생했을 때 운영자는 자연어로 “CPU 부하 급증 시 어떤 사용자가 어떤 URL을 호출했는지 알려줘”라고 질의할 수 있다. PromptOps는 세션 데이터와 트랜잭션 로그를 실시간으로 연계 분석하여, 해당 시간대에 특정 사용자가 반복적으로 특정 URL을 호출한 패턴을 즉시 제시한다. 이 과정에서 LLM이 운영 데이터의 맥락을 이해하고, RCA(Root Cause Analysis) 보고서를 자동 생성한다. 이러한 사례는 <https://www.openmaru.io/session-llm/> 백서에서 상세히 소개되고 있다.

#### 세션-트랜잭션 연계 분석 구조

PromptOps는 세션(사용자)↔트랜잭션(요청)↔인프라(리소스)를 하나의 연속된 맥락으로 연결한다. 예를 들어, 특정 사용자가 로그인 후 대량의 데이터 조회 요청을 반복하여 CPU 부하를 유발하는 경우, 트랜잭션 로그와 세션 이력을 결합해 해당 사용자의 행동 패턴을 시각화한다. LLM은 이 데이터를 기반으로 “이 사용자가 지난 1주일간 동일한 URL을 몇 회 호출했는지”, “해당 트랜잭션이 전체 CPU 사용량의 몇 %를 차지하는지” 등 정밀한 분석 결과를 제공한다. 이로써 운영자는 장애 대응뿐 아니라, 비정상 행동 탐지와 시스템 튜닝에도 활용할 수 있다.

#### LLM 기반 RCA 자동화의 장점

LLM이 RCA 프로세스를 자동화함으로써 운영자의 분석 부담이 크게 줄어든다. 기존에는 로그 파일을 일일이 검색하고 SQL 쿼리를 작성해야 했지만, PromptOps에서는 자연어 한 줄로 복합적인 분석이 가능하다. 또한, LLM은 과거 유사 장애 사례와 비교 분석을 제시하여, 반복되는 장애 패턴의 사전 예방까지 지원한다. 이처럼 세션-트랜잭션-LLM 통합 운영은 엔터프라이즈 환경에서 장애 대응의 패러다임을 근본적으로 전환한다.

세션-트랜잭션-LLM 통합 운영 시나리오의 도입은 단순한 장애 대응을 넘어, 운영 데이터의 맥락적 이해와 자동화된 분석을 가능하게 합니다. 예를 들어, 운영자는 특정 기간 동안의 트랜잭션 패턴을 자연어로 질의하고, LLM이 이를 분석하여 비정상 트래픽이나 잠재적 장애 요인을 사전에 식별할 수 있습니다. 또한, LLM은 대량의 로그 데이터와 세션 정보를 빠르게 종합하여, 운영자가 놓치기 쉬운 미세한 이상 징후까지 탐지할 수 있습니다. 이러한 자동화된 분석과 보고는 운영팀의 숙련도에 관계없이 일관된 품질의 장애 대응을 보장하며, 신규 인력의 온보딩에도 큰 도움이 됩니다. 실제로, PromptOps 방식을 적용한 기업들은 장애 대응 시간의 단축뿐 아니라, 운영 효율성 및 서비스 안정성 향상이라는 장기적 효과도 경험하고 있습니다. 이처럼 세션-트랜잭션-LLM 통합 운영 시나리오는 엔터프라이즈 IT 운영의 새로운 표준으로 자리잡고 있습니다.

## 4.1.2 WAS OOM 장애의 AI 자동 분석 사례

### Java heap space 장애 자동 분석

WAS 환경에서 Out Of Memory(OOM) 장애는 서비스 중단을 초래하는 치명적 이슈다. 기존에는 운영자가 heap dump를 분석하고, 트랜잭션/쿼리 로그를 수작업으로 검토해야 했다. PromptOps는 CogentAI와 연동하여, Java heap space 장애 발생 시 트랜잭션 및 쿼리 데이터를 자동 분석한다. 예를 들어, “최근 OOM 장애의 근본 원인을 트랜잭션 관점에서 분석해줘”라는 질의에 대해, CogentAI는 메모리 사용량 급증 트랜잭션, 대용량 쿼리 실행, 세션 생성 폭증 등 다양한 원인 후보를 신속히 제시한다.

### AI 기반 장애 분석의 효율성

AI 자동 분석은 장애 원인 추적 시간을 획기적으로 단축한다. 기존 수동 분석은 수시간~수일이 소요되지만, PromptOps는 수분 내에 RCA 보고서를 생성한다. CogentAI는 트랜잭션별 메모리 사용량, 쿼리 실행 빈도, 세션 생성 패턴을 시각화하여, 운영자가 즉시 조치할 수 있는 근거를 제공한다. 특히, 반복되는 쿼리나 대용량 데이터 처리 트랜잭션이 OOM의 주 원인임을 데이터

기반으로 명확히 제시한다.

### 실제 도입 효과와 개선점

실제 도입 사례에서는 AI 분석 결과를 바탕으로 쿼리 튜닝, 캐시 정책 변경, 세션 생성 제한 등 조치가 신속히 이루어졌다. 기존의 수동 분석 대비 장애 대응 시간이 70% 이상 단축되었으며, 운영자의 피로도도 실수 가능성이 크게 감소했다. 또한, AI가 과거 유사 장애와 비교 분석을 제공함으로써, 재발 방지와 시스템 안정성 향상에 기여했다.

PromptOps의 AI 기반 OOM 장애 분석은 단순한 로그 분석을 넘어, 복합적인 데이터 상관관계를 자동으로 파악하는 데 큰 강점이 있습니다. 예를 들어, Heap dump 내 객체 분포와 트랜잭션별 메모리 증감 추이를 연계하여, 특정 쿼리나 서비스 호출이 메모리 누수의 원인임을 명확히 밝힐 수 있습니다. 또한, CogentAI는 과거 장애 이력과 비교하여, 동일한 패턴이 반복되는지 여부를 자동으로 탐지합니다. 이 과정에서 운영자는 자연어로 “최근 1개월간 OOM 장애가 발생한 트랜잭션 유형을 비교해줘”와 같은 질의를 할 수 있으며, LLM은 표와 그래프 형태로 결과를 제공합니다. 이러한 자동화는 장애 대응의 신속성뿐만 아니라, 장애 재발 방지와 시스템 구조 개선에도 중요한 인사이트를 제공합니다. 실제 현장에서는 AI 분석 결과를 기반으로 메모리 관리 정책을 개선하고, 반복적 장애의 근본 원인을 제거하는 등 지속적인 품질 향상 효과를 보고 있습니다. 이처럼 PromptOps 기반의 AI 자동 분석은 엔터프라이즈 IT 운영에서 필수적인 방법론으로 자리매김하고 있습니다.

## 4.1.3 공공기관 도입 사례와 운영 효율화 성과

### 공공기관 도입 사례

통합 애플리케이션 플랫폼(iAP)은 디지털서비스몰에 등록되어, 공공기관에서도 PromptOps 방식을 손쉽게 적용할 수 있다. 실제로 여러 공공기관에서 APM 솔루션과 PromptOps 방식을 적용하여, 동시접속자 기반 서버 사이징, 장애 대응, 보고서 자동화 등 다양한 운영 효율성을 경험하고 있다. 공공기관은 보안과 신뢰성이 중요한데, PromptOps는 개인정보 자동 마스킹, 실시간 모니터링, AI 기반 분석 기능을 통해 높은 신뢰도를 확보한다.

### 운영 효율화 성과

다수의 공공기관 및 기업 고객이 지속적으로 유지보수 계약을 갱신하고 있으며, 높은 고객 만족도를 유지하고 있다. 이는 PromptOps 방식이 제공하는 자연어 질의 기반 운영 자동화, 실시간

장애 대응, 보고서 자동 생성 등 실무적 효익이 실제 현장에서 인정받고 있음을 의미한다. 고객들은 “운영팀의 업무 부담이 획기적으로 줄었다”, “장애 대응 시간이 크게 단축되었다”, “보고서 작성이 자동화되어 경영진 보고가 신속해졌다” 등 긍정적인 피드백을 지속적으로 제공하고 있다.

### 신뢰성과 확장성의 시사점

공공기관 도입 사례는 PromptOps의 신뢰성과 확장성을 입증한다. 디지털서비스몰 등록으로 정부기관, 교육기관 등 다양한 분야에서 활용이 가능하며, ON-PREM, 클라우드, 하이브리드 환경 모두 지원한다. 고객 만족도가 높은 이유는, 솔루션 기업의 기술적 완성도와 PromptOps의 실무적 가치가 결합된 결과임을 알 수 있다.

공공기관에서의 PromptOps 적용은 단순한 시스템 도입을 넘어, 실제 운영 환경에서의 신뢰성과 확장성을 동시에 검증하는 중요한 사례로 평가받고 있습니다. 예를 들어, 교육청, 지방자치단체, 공공 연구기관 등 다양한 조직에서 PromptOps 방식을 적용하여, 각기 다른 IT 인프라 환경에서도 일관된 성능과 보안 수준을 유지하고 있습니다. 특히, 개인정보 보호와 관련된 엄격한 규정을 충족하기 위해, PromptOps는 데이터 마스킹, 접근 제어, 감사 로그 등 다양한 보안 기능을 제공하고 있습니다. 또한, 클라우드와 온프레미스 환경을 모두 지원함으로써, 기관별 IT 전략에 유연하게 대응할 수 있습니다. 실제로 한 공공기관에서는 PromptOps 적용 이후 장애 대응 시간이 60% 이상 단축되었고, 운영팀의 보고서 작성 업무가 자동화되어 행정 효율성이 크게 향상되었습니다. 이러한 성공 사례는 다른 공공기관 및 민간 기업에도 긍정적인 영향을 미치고 있으며, PromptOps 방식에 대한 높은 고객 유지율의 주요 요인으로 작용하고 있습니다. 이처럼 공공기관 도입 사례는 PromptOps의 실질적 가치와 시장에서의 신뢰도를 명확히 보여줍니다.

## 4.2 사용자 그룹별 PromptOps 활용 가치

PromptOps는 다양한 사용자 그룹에게 맞춤형 운영 효익을 제공한다. CTO, IT Director, 기술 기획팀장, WAS 운영자, DevOps 엔지니어, 보안 관리자 등 각 역할별로 동시접속자 기반 의사결정, Seasonality 분석, 자동 스케일링, 비정상 접속 탐지 등 핵심 가치를 실현할 수 있다. 이 섹션에서는 각 그룹별 활용 시나리오와 구체적 장점을 심층적으로 설명한다.

### 4.2.1 CTO·IT Director: 서버 투자 의사결정 근거 확보

#### 자연어 질의 기반 데이터 접근

CTO와 IT Director는 서버 투자 의사결정에 있어 정확한 동시접속자 트렌드 데이터가 필수적이다. PromptOps는 “이번 주 동시접속자 최대/최소/평균을 알려줘”, “지난달 대비 서버 부하가 얼마나 증가했는지 보고서로 만들어줘” 와 같은 자연어 질의를 통해 실시간 데이터에 직접 접근할 수 있다. 기존에는 운영팀이 데이터를 추출하고 보고서를 작성하는 복잡한 프로세스가 필요했지만, PromptOps의 LLM 기반 질의 응답으로 한 줄의 질의만으로 즉시 결과를 확인할 수 있다.

### 의사결정 프로세스의 혁신

PromptOps 적용으로 CTO/IT Director의 의사결정 프로세스가 획기적으로 단축된다. 서버 증설·감소 판단, SLA 준수 여부 확인, 리소스 최적화 등 다양한 의사결정이 데이터 기반으로 이루어진다. 특히, HyperLogLog 기반 동시접속자 집계는 감에 의존하던 기존 방식과 달리, 0.81% 오차율의 정밀한 집계를 제공하여 투자 판단의 근거를 명확히 한다.

### 경영진 보고 자동화의 장점

PromptOps는 주간/월간 동시접속자 추이 보고서, 서버별 세션 분포 현황, SLA 준수 여부 등 다양한 경영진 보고서를 자동으로 생성한다. CTO/IT Director는 보고서 작성에 소요되는 시간을 크게 줄이고, 신속한 의사결정과 경영진 커뮤니케이션을 실현할 수 있다. 이는 조직 전체의 운영 효율성과 경쟁력을 높이는 핵심 요소다.

CTO와 IT Director는 조직의 IT 인프라 투자와 전략적 의사결정에 있어 신속하고 정확한 데이터가 필수적입니다. PromptOps는 기존의 수작업 데이터 집계와 보고서 작성 과정을 자동화하여, 경영진이 실시간으로 의사결정에 필요한 정보를 얻을 수 있도록 지원합니다. 예를 들어, 서버 증설이 필요한 시점을 HyperLogLog 기반 동시접속자 집계로 과학적으로 산출할 수 있으며, 투자 대비 효과를 수치로 분석하여 경영진에게 명확한 근거를 제시할 수 있습니다. 또한, PromptOps의 자연어 질의 기능은 IT 비전문가인 경영진도 손쉽게 운영 데이터를 조회하고, 주요 지표를 실시간으로 파악할 수 있게 해줍니다. 이로 인해, IT 부서와 경영진 간의 커뮤니케이션이 원활해지고, 전략적 의사결정의 속도와 품질이 크게 향상됩니다. 실제로 PromptOps 방식을 적용한 기업에서는 서버 투자 비용의 최적화, 불필요한 리소스 낭비 방지, SLA 준수율 향상 등 다양한 실질적 효과를 경험하고 있습니다. 이러한 변화는 조직 전체의 경쟁력 강화로 이어지며, PromptOps가 IT 의사결정의 핵심 도구로 자리잡는 이유가 되고 있습니다.

## 4.2.2 기술 기획팀장·WAS 운영자: Seasonality 기반 리소스 계획과 장애 분석

### Seasonality 패턴 기반 리소스 계획

기술 기획팀장은 반복되는 피크 타임과 Seasonality 패턴을 분석하여 리소스 확보 계획을 수립해야 한다. PromptOps는 “1개월간 1시간 기준 통계 정보 조회”, “작년 블랙프라이데이 대비 현재 접속 패턴 비교” 등 자연어 질의를 통해 시간 단위 롤업 데이터와 시계열 분석 결과를 제공한다. 이를 바탕으로 피크 타임 예측, 서버 증설·감소 계획, 장애 예방 전략을 수립할 수 있다.

### 실시간 장애 분석과 신속 대응

WAS 운영자는 장애 발생 시 실시간으로 원인을 파악해야 한다. PromptOps는 “CPU 부하 급증 시 어떤 사용자가 어떤 URL을 호출했는지 알려줘”, “1시간 이상 접속 중인 사용자 목록을 보여줘” 등 자연어 질의를 통해 세션-트랜잭션 데이터를 즉시 분석한다. LLM과 RAG 기술이 결합되어, 과거 유사 장애와 비교 분석, RCA 자동화, 비정상 행동 탐지 등 신속한 장애 대응이 가능하다.

### 운영 효율성과 업무 부담 감소

PromptOps 적용으로 기술 기획팀장과 WAS 운영자의 업무 부담이 크게 줄어든다. 반복적인 데이터 추출, 보고서 작성, 장애 분석 등 수작업이 자동화되어, 핵심 업무에 집중할 수 있다. 또한, 운영 데이터의 신뢰성과 정확도가 향상되어, 리소스 계획과 장애 대응의 품질이 높아진다.

기술 기획팀장과 WAS 운영자는 IT 서비스의 안정성과 효율성을 동시에 책임지는 중요한 역할을 맡고 있습니다. PromptOps는 Seasonality 분석을 통해 반복적인 트래픽 패턴과 예외적 피크 현상을 쉽게 파악할 수 있도록 지원합니다. 예를 들어, 연말 쇼핑 시즌이나 특정 이벤트 기간의 접속자 급증 패턴을 사전에 예측하여, 서버 리소스를 미리 확보하거나 자동 스케일링 정책을 적용할 수 있습니다. 또한, PromptOps의 실시간 장애 분석 기능은 장애 발생 시 즉각적으로 원인을 파악하고, 과거 유사 사례와 비교하여 최적의 대응 방안을 제시합니다. 이 과정에서 LLM과 RAG 기술이 결합되어, 방대한 로그와 트랜잭션 데이터를 신속하게 분석하고, 운영자가 놓치기 쉬운 비정상 행동이나 잠재적 위험 신호까지 탐지할 수 있습니다. 실제 현장에서는 PromptOps 적용 이후, 장애 대응 시간이 절반 이하로 단축되고, 리소스 과다 할당이나 불필요한 서버 증설이 줄어드는 등 운영 효율성이 크게 향상되었습니다. 이러한 변화는 기술 기획팀과 운영팀의 업무 부담을 줄이고, 조직 전체의 IT 서비스 품질을 한 단계 높이는 데 기여하고 있습니다.

### 4.2.3 DevOps 엔지니어·보안 관리자: 자동 스케일링과 비정상 접속 탐지

#### Kubernetes HPA 연동 자동 스케일링

DevOps 엔지니어는 Kubernetes HPA(Custom Metric) 연동을 통해 동시접속자 수 기반 자동 스케일링 정책을 구현할 수 있다. PromptOps는 APM 솔루션과 연계하여, 실시간 동시접속자 메트릭을 Kubernetes에 제공한다. 예를 들어, “현재 동시접속자 수가 임계치 이상일 때 자동으로 서버를 증설해줘”와 같은 정책이 자연어로 정의되고, 실제 운영에 적용된다. 이는 기존의 수동 서버 사이징 의사결정에서 자동화된 스케일링으로 발전하는 경로를 제시한다.

#### 비정상 접속 탐지와 보안 강화

보안 관리자는 장시간 접속 사용자, 비정상 URL 패턴, 중복 로그인 등 다양한 보안 위협을 탐지해야 한다. PromptOps는 “현재 접속한 사용자 아이디 목록”, “guest1로 접속한 사용자가 조회한 URL 리스트”, “1시간 이상 접속 중인 사용자 목록” 등 자연어 질의를 통해 실시간 행동 분석을 수행한다. LLM과 RAG 기술이 결합되어, 비정상 접속 패턴을 신속히 탐지하고, 보안 정책을 강화할 수 있다.

#### 운영 자동화와 신뢰성 확보

DevOps 엔지니어와 보안 관리자는 PromptOps를 통해 운영 자동화와 신뢰성을 동시에 확보한다. 자동 스케일링, 비정상 접속 탐지, 실시간 모니터링 등 다양한 기능이 자연어 질의로 구현되어, 운영 효율성과 보안 수준이 크게 향상된다. 또한, PromptOps는 기존 Grafana/Prometheus 스택과 공존하며, 확장성과 통합성이 뛰어나다.

DevOps 엔지니어와 보안 관리자는 현대 IT 인프라 운영에서 자동화와 보안의 균형을 맞추는데 중요한 역할을 수행합니다. PromptOps는 Kubernetes HPA와의 연동을 통해, 동시접속자 수와 같은 실시간 메트릭을 기반으로 자동 스케일링 정책을 손쉽게 구현할 수 있도록 지원합니다. 예를 들어, 특정 임계치 이상의 접속자가 발생할 경우, 자동으로 Pod를 증설하거나 리소스를 재분배하여 서비스 안정성을 보장할 수 있습니다. 이 과정에서 운영자는 복잡한 스크립트나 수동 설정 없이, 자연어로 정책을 정의하고 적용할 수 있어 운영 효율성이 크게 향상됩니다. 또한, 보안 관리자는 PromptOps의 실시간 행동 분석 기능을 활용하여, 장시간 접속, 비정상 URL 접근, 중복 로그인 등 다양한 보안 위협을 신속하게 탐지할 수 있습니다. LLM과 RAG 기술이 결합되어, 대규모 로그 데이터 속에서도 이상 징후를 빠르게 식별하고, 필요시 즉각적인 대응 조치를 취할 수 있습니다. 실제로 PromptOps 방식을 적용한 조직에서는 운영 자동화와 보안 강화가 동시에

실현되어, 인프라 관리의 신뢰성과 효율성이 모두 향상되는 효과를 경험하고 있습니다. 또한, 기존의 Grafana/Prometheus와 같은 오픈소스 모니터링 도구와의 통합이 용이하여, 기존 인프라 환경에 무리 없이 적용할 수 있다는 점도 큰 장점입니다. 이처럼 PromptOps는 DevOps와 보안 관리자의 업무를 혁신적으로 변화시키는 핵심 솔루션으로 자리잡고 있습니다.

## 5장: PromptOps 적용 가이드

PromptOps는 엔터프라이즈 WAS 환경에 세션-트랜잭션-AI 통합 운영을 제공하는 차세대 플랫폼입니다. 본 장에서는 PromptOps를 효과적으로 도입하기 위한 단계별 경로, 기술 연동 및 확장 방안, 그리고 라이선스 및 인프라 요구사항을 상세히 안내합니다. 도입 과정은 기존 시스템의 변경 부담을 최소화하면서 점진적으로 AI 기반 자연어 운영 자동화까지 확장할 수 있도록 설계되어 있습니다. 각 단계별 주요 기술적 고려사항과 실무 적용 시 유의점, 그리고 확장성·비용·인프라 측면의 핵심 정보를 제공합니다.

### 5.1 3단계 점진적 도입 경로

PromptOps의 적용은 단순한 솔루션 적용을 넘어, 기존 WAS 환경에 AI 기반의 자연어 운영 자동화를 점진적으로 확장하는 전략적 접근이 필요합니다. 본 절에서는 세션 클러스터링 기반 마련, 트랜잭션 모니터링 및 동시접속자 집계, 그리고 AI 기반 자연어 운영 활성화의 3단계 도입 경로를 구체적으로 안내합니다. 각 단계는 시스템의 안정성과 확장성을 보장하면서, 운영 효율성과 자동화 수준을 단계적으로 높일 수 있도록 설계되어 있습니다. 실무 환경에서의 적용 시 고려해야 할 핵심 기술 요소와 도입 효과, 그리고 단계별 유의사항을 함께 다룹니다.

#### 5.1.1 Phase 1: 세션 클러스터링 솔루션 도입 — 세션 클러스터링 기반 마련

##### 기존 애플리케이션 수정 불필요

세션 클러스터링 솔루션은 기존 WAS 환경에 추가적으로 설치할 수 있으며, 애플리케이션 코드 수정 없이 세션 클러스터링을 적용할 수 있습니다. 이는 기존의 Sticky Session 방식에서 벗어나, IMDG(In-Memory Data Grid) 기반의 외부 세션 저장소를 활용함으로써 서버 장애 시에도 세션

데이터가 유지되는 고가용성 구조를 제공합니다. 설치 과정은 WAS별 플러그인 또는 설정 파일 변경만으로 진행되며, 실제 서비스 중단 없이 점진적 전환이 가능합니다.

세션 클러스터링 솔루션의 가장 큰 장점 중 하나는 기존 애플리케이션의 소스 코드나 아키텍처를 변경할 필요가 없다는 점입니다. 대부분의 WAS 환경에서는 세션 관리 방식의 변경이 서비스 중단이나 대규모 코드 수정으로 이어질 수 있으나, 세션 클러스터링 솔루션은 WAS의 세션 저장소 설정만 외부 IMDG로 전환하면 되므로, 운영 중에도 무중단으로 적용할 수 있습니다. 실제로 여러 금융권, 공공기관, 대형 이커머스 사이트에서 애플리케이션 수정 없이 세션 클러스터링 솔루션을 도입하여, 서비스 안정성과 확장성을 동시에 확보한 사례가 있습니다. 이러한 접근 방식은 도입 초기의 리스크를 최소화하고, 운영팀의 부담을 크게 줄여줍니다.

### Sticky Session에서 IMDG 기반으로 전환

Sticky Session 방식은 로드 밸런서가 특정 사용자의 세션을 하나의 WAS 인스턴스에 고정시키는 구조로, 단일 장애점(SPOF)과 세션 데이터 유실 위험이 존재합니다. 세션 클러스터링 솔루션은 IMDG 기반 세션 저장소를 도입하여, 모든 WAS 인스턴스가 외부 세션 저장소를 공유하도록 설계합니다. 이를 통해 WAS 장애나 재시작 시에도 세션 데이터가 안전하게 유지되며, 클러스터 확장/축소에 따라 세션 데이터의 일관성이 보장됩니다.

IMDG 기반 세션 관리 방식은 기존 Sticky Session의 한계를 극복하는 데 매우 효과적입니다. Sticky Session은 로드 밸런서가 세션을 특정 WAS 인스턴스에 고정시키므로, 해당 인스턴스 장애 시 세션 데이터가 소실될 수 있습니다. 반면, IMDG는 모든 세션 데이터를 메모리 기반 분산 저장소에 저장하여, 여러 WAS 인스턴스가 동일한 세션 데이터에 접근할 수 있도록 합니다. 이로 인해 서버 장애, 재시작, 확장/축소 등 다양한 상황에서 세션 데이터의 일관성과 가용성이 유지됩니다. 또한, IMDG는 세션 데이터의 복제와 자동 Failover 기능을 제공하여, 장애 발생 시에도 서비스 연속성을 보장합니다. 실제로 IMDG 기반 세션 클러스터링을 도입한 기업들은 서비스 장애율 감소, 운영 효율성 향상, 그리고 확장성 측면에서 큰 이점을 경험하고 있습니다.

### 로드 밸런서 설정 변경 없이 적용

세션 클러스터링 솔루션은 기존 로드 밸런서 설정을 변경하지 않고도 적용할 수 있습니다. Sticky Session을 사용하던 환경에서는 세션 저장소를 외부로 이동시키는 것만으로도 세션 클러스터링이 활성화됩니다. 로드 밸런서의 라운드로빈, Least Connection 등 다양한 분산 전략과 호환되며, 서비스 확장 시에도 별도의 네트워크 설정 변경이 필요 없습니다. 이는 운영자의 부담을 줄이고, 도입 초기의 리스크를 최소화하는 중요한 장점입니다.

실제 도입 사례를 살펴보면, 로드 밸런서의 세션 고정 설정을 해제하지 않고도 세션 클러스터링 솔루션을 적용하여, 운영 환경의 네트워크 구조를 그대로 유지하면서 세션 클러스터링 효과를 얻을 수 있었습니다. 또한, 다양한 로드 밸런서(예: F5, NGINX, HAProxy 등)와의 호환성이 검증되어 있어, 기존 인프라를 변경하지 않고도 도입이 가능합니다. 이러한 특성 덕분에, 운영팀은 네트워크 구조나 보안 정책을 변경하지 않고도 세션 클러스터링의 이점을 누릴 수 있으며, 신규 서비스 확장이나 클라우드 전환 시에도 유연하게 대응할 수 있습니다.

## 5.1.2 Phase 2: APM 솔루션 연동 — 트랜잭션 모니터링과 동시접속자 집계

### APM 에이전트 설치와 수집 서버 구성

APM 솔루션은 WAS별 에이전트를 설치하여 트랜잭션 데이터를 실시간으로 수집합니다. 에이전트는 JBoss EAP, Tomcat, WebLogic 등 주요 WAS 환경에서 동작하며, 세션 클러스터와 연계되어 세션-트랜잭션 데이터의 통합 분석 기반을 제공합니다. 수집 서버는 중앙에서 트랜잭션 데이터를 집계·분석하며, 대규모 환경에서는 분산 수집 서버 구조로 확장할 수 있습니다.

APM 에이전트는 WAS 프로세스에 경량으로 삽입되어, 애플리케이션의 성능 저하 없이 트랜잭션, 세션, 시스템 메트릭을 실시간으로 수집합니다. 설치는 단순히 WAS의 설정 파일에 에이전트 라이브러리를 추가하는 방식으로 진행되며, 운영 중에도 적용이 가능합니다. 중앙 수집 서버는 수집된 데이터를 실시간으로 집계·분석하며, 장애 탐지, 성능 병목 분석, 용량 계획 등에 활용됩니다. 대규모 환경에서는 여러 대의 수집 서버를 클러스터로 구성하여, 데이터 처리량과 가용성을 높일 수 있습니다. 또한, 수집 서버는 REST API, Prometheus Exporter 등 다양한 인터페이스를 제공하여, 외부 시스템과의 연동도 용이합니다.

### HyperLogLog 기반 동시접속자 집계 활성화

동시접속자 집계는 HyperLogLog(HLL) 알고리즘을 활용하여, 사용자 ID를 해시 변환 후 선행 0 개수로 고유 사용자 수를 추정합니다. HLL은 16KB 메모리로 수억 명의 사용자를 0.81% 오차로 집계할 수 있으며, 분산 WAS 환경에서 각 인스턴스의 HLL 스케치를 병합하여 전체 시스템 중복을 제거합니다. 이 방식은 기존의 감이나 경험에 의존한 서버 사이징에서 벗어나, 데이터 기반 의사결정을 가능하게 합니다.

HyperLogLog는 대규모 트래픽 환경에서 동시접속자 수를 정확하고 효율적으로 집계하는 데 최적화된 알고리즘입니다. 기존에는 단순 카운팅이나 세션 테이블 스캔 방식이 주로 사용되었으나,

이는 메모리 사용량이 급증하거나 실시간 집계가 어려운 단점이 있었습니다. HLL은 해시 기반의 확률적 추정 방식을 사용하여, 극히 적은 메모리로도 수백만~수억 명의 고유 사용자를 실시간으로 추정할 수 있습니다. APM 솔루션은 각 WAS 인스턴스에서 생성된 HLL 스케치를 중앙 수집 서버에서 병합하여, 전체 시스템의 중복을 제거한 정확한 동시접속자 수를 제공합니다. 이를 통해 IT 운영자는 서버 증설, 용량 계획, 장애 대응 등 다양한 의사결정에 있어 신뢰할 수 있는 데이터를 기반으로 판단할 수 있습니다. 실제로 HLL 기반 집계 도입 후, 서버 오버프로비저닝이나 과소평가로 인한 장애 발생이 크게 감소한 사례가 다수 보고되고 있습니다.

### 시간 단위 롤업 데이터 축적

APM은 트랜잭션 데이터를 2초, 1분, 5분, 1시간 단위로 롤업하여 저장합니다. 시간축 롤업 구조는 Seasonality(주기적 부하 패턴) 분석과 장애 예방에 필수적이며, 주간/월간 동시접속자 추이 분석, 피크 타임 탐지, 과거-현재 비교 등 다양한 운영 보고서 생성의 기반이 됩니다. 롤업 데이터는 Grafana, Prometheus 등 외부 모니터링 시스템과 연동하여 시각화할 수 있습니다.

시간 단위 롤업은 대용량 트랜잭션 데이터를 효율적으로 저장하고, 장기적인 트렌드 분석을 가능하게 합니다. 예를 들어, 2초 단위의 세밀한 데이터는 실시간 장애 탐지와 원인 분석에 활용되고, 1분/5분/1시간 단위의 롤업 데이터는 장기적인 부하 패턴, 피크 타임, 계절성 분석 등에 사용됩니다. 이러한 구조는 데이터 저장소의 용량 부담을 줄이면서도, 필요한 수준의 세밀한 분석을 지원합니다. 또한, 롤업 데이터는 Grafana, Prometheus와 같은 시각화 도구와 연동되어, 운영자와 IT 의사결정자가 직관적으로 시스템 상태를 파악할 수 있도록 도와줍니다. 실제로 시간 단위 롤업 구조를 도입한 기업들은 장애 사전 예방, 용량 계획, SLA 준수 등 다양한 운영 목표를 효과적으로 달성하고 있습니다.

## 5.1.3 Phase 3: CogentAI/PromptOps 적용 — AI 기반 자연어 운영 활성화

### GPU 서버 요구사항과 무상 임대 프로그램

CogentAI 및 PromptOps는 LLM(대형 언어 모델) 기반의 자연어 질의 기능을 제공하며, 운영 자동화와 보고서 생성에 활용됩니다. LLM 모델 서빙에는 GPU 서버가 필요하며, 솔루션 기업에서는 NVIDIA GPU 서버 무상 임대 프로그램을 통해 초기 도입 부담을 완화할 수 있습니다. GPU 서버는 모델 서빙(vLLM, TGI 등)에 최적화되어 있으며, 온프레미스 또는 클라우드 환경에서 유연하게 배포할 수 있습니다.

GPU 서버는 LLM 모델의 대규모 연산을 실시간으로 처리하기 위해 필수적인 인프라입니다. NVIDIA A100, H100 등 최신 GPU는 수천 개의 연산 유닛과 대용량 메모리를 제공하여, 자연어 질의에 대한 빠른 응답과 대규모 동시 사용자 지원이 가능합니다. 솔루션 기업의 무상 임대 프로그램을 활용하면, 초기 투자 비용 없이 GPU 서버를 임시로 도입하여, PoC(개념 검증) 및 초기 운영 자동화 프로젝트를 빠르게 시작할 수 있습니다. 실제로 여러 대기업과 공공기관이 무상 임대 프로그램을 통해 AI 기반 운영 자동화의 효과를 사전에 검증하고, 본격적인 도입을 결정한 사례가 있습니다. 또한, GPU 서버는 온프레미스와 클라우드 환경 모두에서 배포가 가능하므로, 보안 정책이나 데이터 주권 요구사항에 따라 유연하게 선택할 수 있습니다.

### PromptOps 초기 설정과 사용자 교육

PromptOps 적용 시 초기 설정은 세션 데이터와 트랜잭션 데이터의 연결, LLM 모델 선택 (GPT-4, Claude, Gemma3 등), MCP 프로토콜 연동, 자연어 프롬프트 템플릿 구성 등으로 이루어집니다. 운영자 및 IT 의사결정자를 위한 사용자 교육 프로그램이 제공되며, 실제 운영 시나리오 기반의 프롬프트 활용 예제와 자동화 정책 설계 방법을 안내합니다. 초기 도입 후에는 지속적인 운영 데이터 축적과 AI 모델 개선이 이루어집니다.

PromptOps의 초기 설정 과정은 비교적 직관적이지만, 각 기업의 IT 환경과 운영 정책에 맞는 최적화가 필요합니다. 세션 및 트랜잭션 데이터의 연결은 세션 클러스터링 솔루션과 APM에서 수집된 데이터를 통합하여, LLM이 실시간 운영 현황을 이해할 수 있도록 하는 핵심 단계입니다. LLM 모델은 기업의 데이터 보안 정책, 응답 속도, 비용 등을 고려하여 선택할 수 있으며, 필요에 따라 사내 전용 모델을 구축할 수도 있습니다. MCP 프로토콜 연동을 통해 ERP, 그룹웨어, DB 등 다양한 사내 시스템과의 데이터 통합이 가능해집니다. 자연어 프롬프트 템플릿은 운영자들이 자주 사용하는 질의 유형을 미리 정의하여, 반복적인 업무 자동화와 보고서 생성을 손쉽게 할 수 있도록 지원합니다. 사용자 교육은 실제 운영 시나리오를 기반으로 진행되며, 운영자들이 PromptOps의 다양한 기능을 효과적으로 활용할 수 있도록 실습 예제와 가이드가 제공됩니다. 도입 이후에는 운영 데이터가 지속적으로 축적되고, LLM 모델의 성능 개선 및 프롬프트 템플릿의 최적화가 이루어져, 점진적으로 운영 자동화의 수준이 높아집니다.

## 5.2 기술 연동과 확장

PromptOps는 다양한 IT 인프라와의 연동을 통해 운영 자동화와 통합 관측성, 외부 시스템과의 데이터 연계 등 폭넓은 확장성을 제공합니다. Kubernetes, OpenTelemetry, REST API, MCP 등 표준 기술과의 연동을 통해 기존 모니터링 스택과의 공존, 솔루션 기업 생태계와의 시너지, 그리고 사내 시스템 통합을 실현할 수 있습니다. 본 절에서는 PromptOps의 주요 연동 기술과 실무 적용 시 고려해야 할 확장 방안에 대해 구체적으로 설명합니다.

### 5.2.1 Kubernetes·OpenTelemetry·REST API 연동

#### Kubernetes HPA Custom Metric 연동

PromptOps와 APM 솔루션은 Kubernetes Horizontal Pod Autoscaler(HPA)와 연동하여 동시접속자 수 기반 자동 스케일링 정책을 구현할 수 있습니다. APM에서 수집한 동시접속자 메트릭을 Custom Metric API로 Kubernetes에 전달하면, HPA가 실시간 트래픽 변화에 따라 Pod 수를 자동으로 조정합니다. 이는 서버 증설·감소 의사결정을 자동화하고, Peak 타임 대응을 효율적으로 지원합니다.

Kubernetes 환경에서의 자동 스케일링은 서비스의 가용성과 비용 효율성을 동시에 확보하는데 매우 중요합니다. APM 솔루션이 제공하는 동시접속자 수, 트랜잭션 부하 등 실시간 메트릭을 HPA의 Custom Metric API로 연동하면, 단순 CPU/메모리 사용량 기반이 아닌 실제 서비스 부하에 맞춘 지능형 스케일링 정책을 구현할 수 있습니다. 예를 들어, 동시접속자 수가 특정 임계치를 초과하면 Pod가 자동으로 증설되고, 부하 감소 시에는 자동으로 축소되어 리소스 낭비를 방지할 수 있습니다. 이러한 구조는 대규모 이벤트, 프로모션, 예기치 못한 트래픽 급증 등 다양한 상황에서 안정적인 서비스 제공을 가능하게 합니다. 실제로 여러 클라우드 네이티브 기업들이 APM 솔루션과 Kubernetes HPA 연동을 통해 운영 효율성과 비용 절감 효과를 동시에 달성하고 있습니다.

#### OpenTelemetry 표준 기반 통합 관측성

APM 솔루션은 OpenTelemetry 표준을 지원하여, 세션·트랜잭션·분산 트레이스·로그 데이터를 하나의 스키마로 통합합니다. OpenTelemetry Collector와 연동하면, Prometheus, Grafana, Jaeger 등 다양한 관측성 도구와 호환되며, IT 운영자는 기존 모니터링 스택과 PromptOps를 병행하여 사용할 수 있습니다. 표준 기반 통합은 장애 분석, SLA 준수 모니터링, 보고서

자동 생성 등 다양한 운영 시나리오에 활용됩니다.

OpenTelemetry는 클라우드 네이티브 환경에서 관측성 데이터의 표준화를 목표로 하는 오픈소스 프로젝트입니다. APM 솔루션은 OpenTelemetry Collector와의 연동을 통해, 세션, 트랜잭션, 로그, 분산 트레이스 등 다양한 데이터를 하나의 표준 스키마로 변환하여 외부 시스템에 전달합니다. 이를 통해 Prometheus, Grafana, Jaeger, Zipkin 등 다양한 오픈소스 및 상용 관측성 도구와의 호환성이 보장됩니다. 운영자는 기존에 사용하던 모니터링 대시보드와 PromptOps의 AI 기반 자연어 질의 기능을 병행하여, 장애 분석, 용량 계획, SLA 준수 등 다양한 운영 업무를 효율적으로 수행할 수 있습니다. 또한, 표준 기반 통합은 신규 시스템 도입, 클라우드 전환, 멀티클러스터 운영 등 다양한 확장 시나리오에서 유연성을 제공합니다.

### REST API를 통한 외부 시스템 연동

PromptOps는 REST API를 제공하여 외부 시스템(ERP, 그룹웨어, DB 등)과 연동할 수 있습니다. RESTful 인터페이스를 통해 세션/트랜잭션 데이터 조회, 보고서 생성, 자동화 정책 실행 등이 가능하며, 기존 시스템과의 통합을 위한 커스텀 개발이 용이합니다. API 연동은 사내 IT 생태계와의 시너지를 극대화하며, 운영 효율성을 높이는 핵심 요소입니다.

REST API는 시스템 간 연동의 표준 인터페이스로, 다양한 언어와 플랫폼에서 손쉽게 활용할 수 있습니다. PromptOps의 REST API는 인증, 권한 관리, 데이터 필터링 등 보안과 확장성을 고려하여 설계되어 있습니다. 이를 활용하면 ERP 시스템에서 실시간 운영 데이터를 조회하거나, 그룹웨어와 연동하여 자동화된 보고서를 생성하는 등 다양한 업무 자동화 시나리오를 구현할 수 있습니다. 또한, REST API는 외부 개발자나 파트너사가 자체 애플리케이션과 PromptOps를 연동할 때도 활용도가 높아, 기업 전체의 IT 생태계 통합을 촉진합니다. 실제로 REST API 기반 연동을 통해 운영 자동화, 데이터 집계, 실시간 알림 등 다양한 혁신 사례가 보고되고 있습니다.

### Grafana/Prometheus와의 공존 방안

PromptOps는 기존 Grafana/Prometheus 모니터링 스택과 공존할 수 있도록 설계되었습니다. APM에서 수집한 메트릭을 Prometheus Exporter로 제공하고, Grafana 대시보드에서 시각화할 수 있습니다. 기존 모니터링과 AI 기반 자연어 질의 기능을 병행하여, 운영자와 IT 의사 결정자가 각자의 업무에 맞는 도구를 선택적으로 활용할 수 있습니다.

Grafana와 Prometheus는 오픈소스 모니터링 분야에서 널리 사용되는 도구로, 시각화와 실시간 알림 기능이 강점입니다. PromptOps는 이러한 기존 모니터링 스택과의 호환성을 보장하여, 기존 운영팀이 익숙한 환경을 그대로 유지하면서도, AI 기반 자연어 질의 및 자동화 기능을 추가로

활용할 수 있도록 지원합니다. 예를 들어, 운영자는 Grafana 대시보드에서 실시간 메트릭을 모니터링하면서, PromptOps를 통해 자연어로 장애 원인 분석이나 보고서 생성을 요청할 수 있습니다. 이처럼 두 시스템의 공존은 IT 운영의 유연성과 효율성을 극대화하며, 신규 기능 도입에 따른 교육 부담도 최소화합니다.

## 5.2.2 MCP 기반 사내 시스템 연동: ERP·그룹웨어·DB

### MCP 프로토콜의 표준화 연동 구조

MCP(Model Context Protocol)는 파일 시스템, DB, ERP, 그룹웨어 등 다양한 사내 시스템과 CogentAI를 연동하는 표준 프로토콜입니다. MCP는 데이터 구조와 컨텍스트를 LLM에 전달하여, 자연어 질의와 AI 분석이 사내 시스템 데이터에 직접 적용될 수 있도록 지원합니다. 표준화된 연동 구조는 확장성과 보안, 데이터 일관성을 보장합니다.

MCP는 사내 시스템의 데이터 구조와 맥락 정보를 LLM에 전달하는 데 최적화된 프로토콜로, 데이터의 스키마, 권한, 접근 정책 등을 표준화된 방식으로 정의할 수 있습니다. 이를 통해 ERP, 그룹웨어, DB 등 이기종 시스템 간의 데이터 통합이 용이해지며, LLM이 각 시스템의 맥락을 이해하고 자연어 질의에 적합한 답변을 생성할 수 있습니다. MCP는 TLS 기반 암호화, 인증/인가 정책, 감사 로그 등 보안 기능도 내장하고 있어, 기업의 데이터 보안 요구사항을 충족할 수 있습니다. 실제로 MCP를 도입한 기업들은 사내 시스템 간 데이터 연동의 복잡성을 크게 줄이고, AI 기반 업무 자동화의 범위를 확장하는 데 성공하고 있습니다.

### 솔루션 기업 생태계와의 시너지

솔루션 기업 Dashboard, COP(Cloud Operations Platform), Observability, iAP(통합 애플리케이션 플랫폼) 등 솔루션 기업 제품 생태계와 PromptOps의 연동은 운영 효율성과 데이터 활용도를 극대화합니다. MCP를 통해 각 시스템의 데이터와 운영 정책을 AI 기반으로 통합 관리할 수 있으며, 보고서 자동 생성, 장애 분석, 리소스 계획 등 다양한 시나리오에서 시너지가 발생합니다.

솔루션 기업 생태계는 다양한 운영 도구와 플랫폼을 통합적으로 제공하여, 기업 IT 환경의 복잡성을 줄이고 데이터 활용도를 높입니다. PromptOps는 MCP를 통해 이들 시스템과 데이터를 실시간으로 연동하여, 운영자와 IT 의사결정자가 자연어로 다양한 업무를 자동화할 수 있도록 지원합니다. 예를 들어, 솔루션 기업 Dashboard에서 실시간 모니터링 정보를 PromptOps에 연동하여, “이번 주 장애 발생 원인 분석 보고서를 생성해줘”와 같은 자연어 요청을 처리할 수

있습니다. 또한, COP, iAP 등과의 연동을 통해 리소스 계획, 배포 자동화, SLA 준수 모니터링 등 다양한 운영 시나리오에서 AI 기반 자동화의 효과를 극대화할 수 있습니다. 이러한 시너지는 기업의 IT 운영 효율성, 데이터 기반 의사결정, 업무 자동화 수준을 한 단계 끌어올립니다.

### 사내 시스템 통합과 데이터 활용

MCP 기반 연동은 사내 시스템의 데이터(ERP 인사정보, 그룹웨어 일정, DB 트랜잭션 등)를 PromptOps에서 자연어 질의로 조회·분석할 수 있게 합니다. 운영자는 복잡한 SQL이나 API 호출 없이, “지난달 인사 이동자 목록을 보고서로 만들어줘”와 같은 프롬프트로 자동화된 결과를 얻을 수 있습니다. 이는 IT 운영의 효율성과 정확도를 크게 높이는 혁신적 변화입니다.

사내 시스템 통합은 기존에는 복잡한 커스텀 개발, API 연동, 데이터 변환 작업이 필수적이었으나, MCP와 PromptOps를 활용하면 자연어 프롬프트만으로 다양한 데이터를 조회·분석할 수 있습니다. 예를 들어, ERP 시스템의 인사정보, 그룹웨어의 일정 데이터, DB의 트랜잭션 기록 등 다양한 정보를 LLM이 이해할 수 있는 형태로 통합하여, 운영자가 원하는 보고서나 분석 결과를 자동으로 생성할 수 있습니다. 이러한 방식은 반복적인 데이터 집계, 보고서 작성, 이상 탐지 등 업무의 자동화 수준을 크게 높여주며, 데이터 활용의 정확도와 신뢰성도 함께 향상시킵니다. 실제로 MCP 기반 통합을 도입한 조직에서는 IT 운영팀의 업무 부담이 크게 줄고, 데이터 기반 의사결정의 속도와 품질이 향상된 사례가 다수 보고되고 있습니다.

## 5.3 라이선스와 인프라 요구사항

PromptOps 적용 시에는 제품별 라이선스 정책, 배포 옵션, 그리고 인프라 요구사항을 명확히 이해하는 것이 매우 중요합니다. 각 솔루션 기업 제품군은 에이전트 수 기반의 가격 정책과 다양한 배포 옵션(클라우드, 온프레미스 등)을 제공하며, AI 기반 운영 자동화를 위한 GPU 서버 도입 시에는 무상 임대 프로그램을 통해 초기 비용 부담을 줄일 수 있습니다. 본 절에서는 실무 환경에서 PromptOps 방식을 적용할 때 반드시 고려해야 할 라이선스 구조와 인프라 요구사항을 상세히 안내합니다.

### 5.3.1 제품별 라이선스 구조와 배포 옵션

#### 에이전트 수 기반 가격 정책

APM 솔루션, 세션 클러스터링 솔루션, CogentAI는 에이전트 수(설치된 WAS 인스턴스 수)

기반의 가격 정책을 적용합니다. 각 WAS별로 설치된 에이전트 수에 따라 라이선스 비용이 산정되며, 대규모 환경에서는 볼륨 할인 및 연간 유지보수 옵션이 제공됩니다. 이는 투명한 비용 산정과 예측 가능한 운영비용 관리에 도움이 됩니다.

에이전트 수 기반 라이선스 정책은 기업의 인프라 규모와 운영 환경에 따라 유연하게 적용할 수 있는 장점이 있습니다. 예를 들어, 초기 도입 시에는 소규모로 시작하여 점진적으로 에이전트 수를 늘릴 수 있으며, 대규모 확장 시에는 볼륨 할인 혜택을 받을 수 있습니다. 연간 유지보수 옵션을 선택하면, 정기적인 소프트웨어 업데이트, 기술 지원, 장애 대응 등 다양한 서비스를 안정적으로 제공받을 수 있습니다. 또한, 라이선스 비용 산정이 명확하여 예산 계획과 비용 관리가 용이하며, 도입 후에도 운영 환경 변화에 따라 라이선스 규모를 조정할 수 있습니다. 실제로 에이전트 수 기반 정책을 도입한 기업들은 예산 초과나 불필요한 비용 지출 없이, 효율적으로 IT 운영비용을 관리하고 있습니다.

### 클라우드/온프레미스 지원 범위

솔루션 기업 제품군은 AWS, Azure, GCP 등 퍼블릭 클라우드 환경과 온프레미스 데이터센터 모두에서 배포가 가능합니다. 클라우드 SaaS, 프라이빗 클라우드, 온프레미스 설치 등 다양한 배포 옵션이 제공되며, 고객 환경에 맞게 유연하게 선택할 수 있습니다. 배포 방식에 따라 라이선스 정책 및 지원 범위가 달라질 수 있으므로, 도입 전 상세한 상담이 필요합니다.

클라우드 환경에서는 솔루션 기업의 SaaS 서비스를 통해 신속한 도입과 운영이 가능하며, 인프라 관리 부담이 크게 줄어듭니다. 프라이빗 클라우드나 온프레미스 환경에서는 보안, 데이터 주권, 커스터마이징 등 기업별 요구사항에 맞는 맞춤형 배포가 가능합니다. 각 배포 옵션은 네트워크 구성, 보안 정책, 데이터 저장 위치 등 다양한 요소에 영향을 미치므로, 도입 전 IT 인프라팀과의 긴밀한 협의가 필요합니다. 또한, 배포 방식에 따라 기술 지원 범위, 업데이트 정책, 장애 대응 방식 등이 달라질 수 있으므로, 솔루션 기업과의 사전 상담을 통해 최적의 옵션을 선택하는 것이 중요합니다.

### CogentAI의 유연한 배포 옵션

CogentAI는 클라우드 SaaS(솔루션 기업 Cloud), 프라이빗 클라우드, 온프레미스 등 다양한 배포 옵션을 제공합니다. GPU 서버 요구사항에 따라 클라우드 기반 GPU 인스턴스 또는 온프레미스 GPU 서버를 선택할 수 있으며, 초기 도입 시 무상 임대 프로그램을 활용할 수 있습니다. 배포 옵션은 보안, 데이터 주권, 운영 효율성 등 고객 요구에 따라 맞춤형으로 설계됩니다.

CogentAI의 유연한 배포 옵션은 기업의 IT 전략과 보안 정책에 따라 다양한 선택지를 제공함

니다. 클라우드 SaaS는 빠른 도입과 확장, 낮은 초기 투자 비용이 장점이며, 온프레미스 배포는 데이터 보안, 내부 정책 준수, 커스터마이징 등에서 강점을 가집니다. GPU 서버는 솔루션 기업의 무상 임대 프로그램을 통해 초기 도입 부담을 줄일 수 있으며, 클라우드 기반 GPU 인스턴스와 온프레미스 GPU 서버 중에서 선택이 가능합니다. 실제로 보안이 중요한 금융권, 공공기관 등에서는 온프레미스 배포를 선호하며, 스타트업이나 중소기업은 클라우드 SaaS를 통해 빠르고 유연하게 AI 기반 운영 자동화를 도입하고 있습니다. 각 배포 옵션은 솔루션 기업의 컨설팅을 통해 기업별 요구에 맞게 최적화할 수 있습니다.

### 5.3.2 인프라 요구사항과 GPU 서버 무상 임대 프로그램

#### IMDG 노드용 메모리 중심 서버

세션 클러스터링 솔루션 도입 시 IMDG 노드(세션 저장소)는 메모리 중심 서버가 필요합니다. 세션 데이터의 고속 저장과 Failover를 위해 충분한 RAM(최소 32GB 이상)이 권장되며, CPU와 네트워크 대역폭도 클러스터 규모에 따라 적절히 설계해야 합니다. 장애 대비를 위해 노드 간 복제와 자동 Failover 기능을 활성화해야 합니다.

IMDG 노드는 세션 데이터의 실시간 저장과 장애 복구를 담당하므로, 메모리 용량과 네트워크 성능이 매우 중요합니다. 대규모 트래픽 환경에서는 노드 간 데이터 복제, 분산 저장, 자동 Failover 기능을 통해 서비스 연속성을 보장해야 하며, 이를 위해 최소 32GB 이상의 RAM과 고성능 CPU, 10Gbps 이상의 네트워크 대역폭이 권장됩니다. 또한, IMDG 소프트웨어의 설정을 통해 노드 간 데이터 복제 정책, 장애 발생 시 자동 전환(Failover) 기능을 활성화해야 합니다. 실제로 IMDG 노드의 메모리 용량 부족이나 네트워크 병목 현상은 세션 데이터 유실, 서비스 지연 등 심각한 장애로 이어질 수 있으므로, 인프라 설계 단계에서 충분한 용량과 성능을 확보하는 것이 매우 중요합니다.

#### APM 에이전트+수집 서버 요구사항

APM 솔루션은 WAS별 에이전트와 중앙 수집 서버로 구성됩니다. 에이전트는 WAS 서버에 경량으로 설치되며, 수집 서버는 트랜잭션 데이터 집계와 분석을 담당합니다. 수집 서버는 CPU와 I/O 성능이 중요하며, 대규모 환경에서는 분산 수집 서버 구조로 확장할 수 있습니다. 네트워크 안정성과 데이터 저장소(SSD 등)도 고려해야 합니다.

APM 에이전트는 WAS 서버의 리소스를 최소한으로 사용하도록 설계되어 있어, 애플리케이션 성능에 미치는 영향이 거의 없습니다. 중앙 수집 서버는 실시간으로 대량의 트랜잭션 데이터를

집계·분석하므로, 고성능 CPU와 빠른 디스크 I/O(SSD 등)가 필수적입니다. 대규모 환경에서는 수집 서버를 여러 대로 분산하여, 데이터 처리량과 가용성을 높일 수 있습니다. 또한, 네트워크 안정성은 데이터 유실 방지와 실시간 분석 정확도에 직접적인 영향을 미치므로, 전용 네트워크 또는 고속 네트워크 환경을 구축하는 것이 권장됩니다. 데이터 저장소는 SSD 기반으로 구성하여, 대용량 트랜잭션 데이터의 빠른 저장과 조회를 지원해야 합니다. 실제로 이러한 인프라 요구사항을 충족한 기업들은 트랜잭션 모니터링, 장애 분석, 용량 계획 등 다양한 운영 업무를 안정적으로 수행하고 있습니다.

### GPU 서버 요구사항과 무상 임대 안내

CogentAI 및 PromptOps 적용 시 GPU 서버가 필요합니다. NVIDIA A100, H100 등 최신 GPU가 권장되며, 모델 서빙(vLLM, TGI 등)에 최적화된 환경이 필요합니다. 솔루션 기업에서는 AI 기반 운영을 도입하는 기업을 위해 NVIDIA GPU 서버 무상 임대 프로그램을 제공하며, 초기 도입 부담을 줄이고 AI 운영 자동화를 빠르게 실현할 수 있습니다. GPU 서버는 클라우드 또는 온프레미스 환경에서 선택적으로 배포할 수 있습니다.

GPU 서버는 LLM 모델의 대규모 연산을 실시간으로 처리하기 위해 필수적인 인프라입니다. NVIDIA A100, H100 등 최신 GPU는 수천 개의 연산 유닛과 대용량 메모리를 제공하여, 자연어 질의에 대한 빠른 응답과 대규모 동시 사용자 지원이 가능합니다. 솔루션 기업의 무상 임대 프로그램을 활용하면, 초기 투자 비용 없이 GPU 서버를 임시로 도입하여, PoC(개념 검증) 및 초기 운영 자동화 프로젝트를 빠르게 시작할 수 있습니다. 실제로 여러 대기업과 공공기관이 무상 임대 프로그램을 통해 AI 기반 운영 자동화의 효과를 사전에 검증하고, 본격적인 도입을 결정한 사례가 있습니다. 또한, GPU 서버는 온프레미스와 클라우드 환경 모두에서 배포가 가능하므로, 보안 정책이나 데이터 주권 요구사항에 따라 유연하게 선택할 수 있습니다.

---

## Appendix

### References

1. “HyperLogLog 설명 및 활용 사례”. [출처 링크](#)

2. “Kubernetes Horizontal Pod Autoscaler 공식 문서”. [출처 링크](#)
3. “OpenTelemetry 공식 문서”. [출처 링크](#)
4. “통합 애플리케이션 플랫폼(iAP) 디지털서비스몰 등록 사례”. [출처 링크](#)
5. “솔루션 기업 세션-트랜잭션-LLM 통합 운영 백서”. [출처 링크](#)
6. “PromptOps 공식 소개 및 사례”. [출처 링크](#)
7. “AIOps와 VibeOps 비교”. [출처 링크](#)
8. “CogentAI Technical Overview”. [출처 링크](#)
9. “CogentAI 공식 문서”. [출처 링크](#)
10. “CogentAI 공식 소개”. [출처 링크](#)
11. “Hazelcast IMDG 공식 문서”. [출처 링크](#)
12. “HyperLogLog Algorithm”. [출처 링크](#)
13. “HyperLogLog 알고리즘 설명”. [출처 링크](#)
14. “IMDG 기반 세션 클러스터링 아키텍처”. [출처 링크](#)
15. “JBoss EAP 공식 문서”. [출처 링크](#)
16. “Kubernetes HPA Documentation”. [출처 링크](#)
17. “Kubernetes 공식 문서”. [출처 링크](#)
18. “LLM+RAG 기반 운영 자동화”. [출처 링크](#)
19. “NVIDIA GPU 서버 안내”. [출처 링크](#)
20. “APM 솔루션 Documentation”. [출처 링크](#)
21. “APM 솔루션 공식 문서”. [출처 링크](#)
22. “APM 솔루션 공식 문서”. [출처 링크](#)
23. “세션 클러스터링 솔루션 Documentation”. [출처 링크](#)
24. “세션 클러스터링 솔루션 공식 문서”. [출처 링크](#)
25. “세션 클러스터링 솔루션 공식 문서”. [출처 링크](#)
26. “솔루션 기업 Session-LLM 백서”. [출처 링크](#)
27. “OpenTelemetry 공식 문서”. [출처 링크](#)
28. “PromptOps 공식 소개”. [출처 링크](#)
29. “PromptOps 적용 안내”. [출처 링크](#)
30. “PromptOps 특허 NP25073-KR”. [출처 링크](#)

31. “Redis 공식 문서”. [출처 링크](#)

## Glossary

용어	정의
롤업(roll-up)	데이터의 시간 단위 집계 및 축약
세션 클러스터링	여러 WAS 인스턴스 간 세션 데이터 공유 및 장애 복구를 위한 분산 저장 기술
할루시네이션	LLM이 근거 없는 허위 정보를 생성하는 현상
AIOps	Artificial Intelligence for IT Operations, AI 기반 IT 운영 자동화
APM	Application Performance Monitoring, 애플리케이션 성능 모니터링 시스템
CogentAI	LLM+RAG+MCP 기반 AI 운영 자동화 솔루션.
Failover	장애 발생 시 자동으로 대체 노드로 전환하는 메커니즘
HPA	Horizontal Pod Autoscaler, Kubernetes의 자동 스케일링 기능.
HyperLogLog	고유 사용자 수를 효율적으로 추정하는 알고리즘.
HyperLogLog(HLL)	메모리 효율적 고유 사용자 수 추정 알고리즘.
IMDG	In-Memory Data Grid, 메모리 기반 분산 데이터 저장소.
Kubernetes HPA	Horizontal Pod Autoscaler, Kubernetes에서 자동 스케일링을 담당하는 컴포넌트
LLM	Large Language Model, 대형 언어 모델.
LLM 할루시네이션	근거 없는 AI 답변 생성 현상.
MBean	Java Management Extensions에서 관리 객체를 의미하는 용어
MCP	Model Context Protocol, 표준화된 시스템 연동 프로토콜.
OOM	Out Of Memory, 메모리 부족으로 인한 시스템 장애
APM 솔루션	트랜잭션 모니터링과 동시접속자 집계 기능을 제공하는 애플리케이션 성능 관리 솔루션.
세션 클러스터링 솔루션	IMDG 기반 세션 클러스터링 솔루션으로 WAS 환경의 고가용성 세션 관리 제공.
OpenTelemetry	표준 기반 통합 관측성 프레임워크.
PromptOps	자연어 질의와 AI 분석을 통해 세션·트랜잭션 정보를 운영 자동화하는 방법론.
RAG	Retrieval-Augmented Generation, LLM의 실시간 데이터 참조 기술.
RCA	Root Cause Analysis, 장애 근본 원인 분석 프로세스
REST API	Representational State Transfer Application Programming Interface, 시스템 간 연동을 위한 표준 인터페이스.
Seasonality	주기적 패턴, 반복되는 트래픽/부하 현상
SPOF	Single Point of Failure, 단일 장애점

<b>Sticky Session</b>	로드 밸런서가 세션을 특정 WAS 인스턴스에 고정시키는 방식.
<b>VibeOps</b>	AI 코파일럿 기반 지능형 IT 운영 패러다임
<b>WAS</b>	Web Application Server, 웹 애플리케이션을 실행하는 서버 소프트웨어

---

# Contact Us

 [hello@cncf.co.kr](mailto:hello@cncf.co.kr)

 02-469-5426

 [www.cncf.co.kr](http://www.cncf.co.kr)

## CNF Blog

다양한 콘텐츠와 전문 지식을 통해 더 나은 경험을 제공합니다.

## CNF eBook

이제 나도 클라우드 네이티브 전문가  
쿠버네티스 구축부터 운영 완전 정복

## CNF Resource

Community Solution의 최신 정보와  
유용한 자료를 만나보세요.

