

가상화 vs 클라우드 네이티브 백서- 2026 클라우드 네이티브는 AI 시대의 인프라 표준

"탈VMware를 한다는데, 다른 가상화로 바꾸기만 하면 정말 끝나는 걸까"
Broadcom 인수 후 라이선스가 2~10배 오르고 CentOS EOL로 OS
구독비는 TCO의 30~50%를 잠식하지만, 대안 IaaS로의 VM→VM 이주는
OS만 다시 설치 할 뿐이며 게스트 OS·정적 할당·명령형 운영이라는 3중
한계는 그대로 남습니다. MSAP.ai는 영구 라이선스 하나에
웹서버·WAS·APM·Pod Cluster·AI·MSA Accelerator를 단일 번들로 통합한
국산 PaaS로, 해외 상용 K8s의 반복 구독 구조를 1회 원화 자산화로
재정의합니다. 이 백서에서 VM 대비 3~10배 집적률·하드웨어 30~50%
절감의 회계적 근거와, Booking.com·Adidas·국내 금융권 사례를 기반으로
한 12개월 3분류 현대화 실행 로드맵, 그리고 이사회 설득용 3대
KPI(TCO·집적률·배포 리드타임)를 확인하실 수 있습니다.

Contact Us

 02-6953-5427

 hello@msap.ai

 www.msap.ai

Contents

- 1.1 하이퍼바이저 가상화의 탄생: 서버 하드웨어 낭비에서 출발한 1999~2010년대 IaaS 시대 7
 - 1.1.1 VMware·OpenStack·Nutanix·CloudStack의 발원과 목적 차이 7
 - 1.1.2 VMWare 를 Broadcom 이 인수 후 상용 IaaS 라이선스 지형 변화(1차 정리 지점) 9
- 1.2 컨테이너·쿠버네티스의 탄생: Google Borg에서 CNCF v1.0 정의까지 11
 - 1.2.1 FreeBSD jails → LXC → Docker로 이어진 Linux 컨테이너 계보와 2013년 UX 혁신 11
 - 1.2.2 Borg(2003)→Kubernetes(2014)와 CNCF 5대 정의(선언형·불변 인프라) 13
- 2장: IaaS와 PaaS의 핵심 개념·용어와 운영 모델의 구조적 차이** 16
 - 2.1 NIST 정의로 정렬한 IaaS·PaaS·SaaS와 VM·Container·Pod의 격리 원리 16
 - 2.1.1 NIST SP 800-145 기준 IaaS와 PaaS의 경계 16
 - 2.1.2 Linux namespace·cgroups·OverlayFS로 본 컨테이너 격리의 커널 공유 구조 17
 - 2.1.3 VM·Container·Pod 3계층의 격리 경계와 자원 할당 단위 비교 20
 - 2.2 명령형 vs 선언형, 가변 vs 불변, 정적 vs 동적 — 운영 모델 3대 축 22
 - 2.2.1 kubectl apply와 컨트롤러 루프로 구현된 멱등성의 의미 22
 - 2.2.2 불변 인프라(Immutable Infrastructure)와 이미지 레이어링의 결합 23
 - 2.2.3 Request/Limit·HPA·VPA·KEDA·Cluster Autoscaler의 동적 자원 할당 25
- 3장: VM→VM 이주가 중복 투자인 구조적 이유 — OS tax·정적 할당·명령형의 3중 한계** 27
 - 3.1 게스트 OS 중복(OS tax)이 TCO의 숨은 50%를 차지하는 메커니즘 27
 - 3.1.1 VM 1대당 수백 MB~수 GB 고정 점유와 100 VM=OS 100개 패치 구조 27
 - 3.1.2 금융·공공 필수 OS 종속 솔루션(라이선스·접근제어·SMS·EDR·접속 프로그램)의 VM별 중복 비용 구조 28

- 3.1.3 컨테이너에는 Guest OS가 존재하지 않는다 — VM과의 본질적 구조 차이
 - 와 OS 종속 비용 전 카테고리의 구조적 제거 30
- 3.2 정적 자원 할당과 명령형 배포의 한계 31
 - 3.2.1 vCPU·메모리 사전 할당과 피크 아닌 시점의 유휴 자원 누적 32
 - 3.2.2 이미지 클론·구성관리 방식의 비역등성과 드리프트 누적 33
 - 3.2.3 VM 부팅 지연 vs 컨테이너 기동의 복구·확장 탄력성 격차 34
- 3.3 엔터프라이즈 Linux 구독 비용 상승·CentOS EOL 압력에 대한 근본 해결 — 다른 Linux로의 이전이 아니라 OS 없는 시스템 환경으로의 전환 35
 - 3.3.1 CentOS EOL과 글로벌 엔터프라이즈 Linux 구독 비용 상승 — VM 운영이
 - 떠안게 된 2020년대 OS 리스크 35
 - 3.3.2 Rocky Linux·AlmaLinux 등 대체 Linux 이전의 한계 — “OS 갱신 비용”
 - 이 반복되는 리프트앤시프트 36
 - 3.3.3 근본 해결 경로 — 컨테이너 전환으로 OS 자체가 없는 시스템 환경 구축 . 37
- 4장: 하이퍼바이저·IaaS 진영 경쟁 비교 (VMware·OpenStack·Nutanix·CloudStack) 38**
 - 4.1 IaaS 4종 비교 매트릭스 — 라이선스·커뮤니티·에코시스템의 구조적 차이 38
 - 4.1.1 VMware vSphere/VCF — IaaS 4종 중 유일한 독점 상용 라이선스 포지션 38
 - 4.1.2 OpenStack — Apache 2.0 IaaS와 디스트로 파편화(RHOSP→RHOSO
 - 이행) 39
 - 4.1.3 Nutanix AHV·Apache CloudStack — HCI와 통신사 중심 IaaS 포지셔닝 40
 - 4.2 IaaS 4종이 공유하는 “VM 단위 운영”의 한계 41
 - 4.2.1 IaaS 4종의 K8s 지원은 모두 “VM 위 K8s” 형태라는 사실 41
 - 4.2.2 RHOSP→RHOSO 이행과 OpenStack 워크로드의 K8s 통합 — 2024년
 - IaaS 벤더 지형 변화 42
 - 4.3 IaaS 진영의 Cloud Native 응답 — KubeVirt는 브리지(Bridge) 기술이지 근본
 - 해결책이 아니다 43
 - 4.3.1 KubeVirt의 포지셔닝 — K8s 위에서 VM을 파드 오브젝트로 다루는 브리지
 - 기술 43

4.3.2 브리지의 본질적 한계 — OS tax·정적 할당·명령형 운영이 K8s 표면으로
 옮겨질 뿐 44

**5장: 클라우드 네이티브 진영 경쟁 비교 — MSAP.ai 기준 (vs Kubernetes 업스트림·해외
 상용 엔터프라이즈 K8s 디스트리뷰션 A·B·C) 46**

5.1 MSAP.ai 기준 비교 매트릭스 — 경쟁 PaaS(K8s 업스트림·해외 상용 디스트리뷰션
 A·B·C) 대비 라이선스·구성·지원 차이 46

5.1.1 MSAP.ai — PaaS 기준 제품 46

5.1.2 Kubernetes 업스트림 (Apache 2.0, CNCF Graduated) 48

5.1.3 해외 상용 엔터프라이즈 K8s 디스트리뷰션 A —
 Route·Build·OAuth·SCC·멀티클러스터 관리 결합형 (경쟁 제품) . . . 49

5.1.4 해외 상용 디스트리뷰션 B·C — 멀티클러스터 중심형과 가상화 벤더 계열
 결합형 (경쟁 제품) 50

5.2 공공·금융·AI 관점 제품 선택 기준 51

5.2.1 베어메탈 지원 범위와 GPU Operator 공식 지원 수준 비교 51

5.2.2 국내외 공공·금융 채택 현황과 엔터프라이즈 지원 모델 52

6장: VM 위 K8s vs 베어메탈 K8s — 3중 스택 과도기 타협과 목표 아키텍처 54

6.1 VM 위 K8s가 만들어내는 3중 스택의 구조적 비효율 54

6.1.1 하이퍼바이저·게스트 OS·컨테이너 런타임 3계층 누적 구조 54

6.1.2 TKG·Karbon·Magnum·CKS 제공 형태와 벤더 락인 경로 56

6.2 베어메탈 K8s 목표 아키텍처와 필수 구성 요소 58

6.2.1 MetalLB·Rook-Ceph·Longhorn·cert-manager·External DNS 자체
 구성 58

6.2.2 Booking.com·CERN·Adidas 등 대규모 베어메탈 K8s 운영 사례 60

7장: 라이선스·전력·인건비·Auto Scale 사이징 4대 TCO 축 비교 62

7.1 라이선스 축 — 구(舊) 상용 IaaS 라이선스와 상용 오픈소스 기반 MSAP.ai의 회계적
 비교 63

7.1.1 구(舊) 상용 IaaS 구독 라이선스의 3년 예산 영향 — 구조적 요인만 압축 서술 63

- 7.1.2 MSAP.ai — 영구 라이선스·번들 플랫폼 기반의 비용 효율적 국내 라이선스
 - 구조 64
- 7.2 전력·자원 집적도·인건비 축 66
 - 7.2.1 동일 H/W에서 VM 대비 컨테이너 3~10배 집적률의 자원·전력 효과 . . . 67
 - 7.2.2 VM 운영팀 + K8s 운영팀 이중화에서 단일 K8s 표면으로의 FTE 수렴 . . 68
 - 7.2.3 “예시적(illustrative)” TCO 시뮬레이션 제시 원칙과 변동성 요인 69
- 7.3 Auto Scale 기반 하드웨어 사이징 절감 축 — VM 정적 할당 vs K8s 동적 스케일 70
 - 7.3.1 VM 정적 할당이 강제하는 구조적 유헴 — 자원 사용률이 낮을 수밖에 없는
 - 이유 71
 - 7.3.2 HPA·KEDA·Cluster Autoscaler의 0→N→0 동적 스케일 메커니즘 . . . 72
 - 7.3.3 동일 워크로드 기준 하드웨어 사이징 30~50% 절감(예시적) — 산출 전제와
 - 한계 73
- 8장: AI 스택(모델·Vector DB·RAG·API·MCP·Agent)의 K8s 적합성과 VM 비경제성 74**
 - 8.1 AI 스택 구성 요소의 K8s 기본 기능 의존성 74
 - 8.1.1 KServe·vLLM·Triton을 통한 선언형 모델 서빙과 자동 확장 75
 - 8.1.2 Milvus·Weaviate·Qdrant Vector DB의 Operator·Helm 배포 표준 . . 76
 - 8.1.3 LangChain·LlamaIndex·MCP 서버·Agent의 컨테이너 배포 패턴 . . . 78
 - 8.2 GPU 공유·MIG·Time-slicing으로 본 VM 비경제성의 기술적 근거 80
 - 8.2.1 NVIDIA GPU Operator·Device Plugin·MIG·Time-slicing의 다중 파드
 - 공유 80
 - 8.2.2 VM GPU 패스스루의 전용 점유 구조와 AI 인프라 단가 영향 81
- 9장: 관리 포인트 이중화 비용과 단일 K8s 표면으로의 수렴 전략 83**
 - 9.1 VM과 컨테이너가 공존할 때 발생하는 이중 운영 비용 83
 - 9.1.1 FTE·모니터링·보안 정책의 3중 중복 경로 83
 - 9.1.2 CI/CD·백업·DR 파이프라인의 분기 유지 비용 85
 - 9.2 단일 K8s 표면으로 통합하는 운영 아키텍처 86
 - 9.2.1 Prometheus·OpenTelemetry·Falco·Trivy·Cosign으로 통합되는 관측·
 - 보안 87

9.2.2 Argo CD·Flux·Tekton·Crossplane으로 통합되는 GitOps·IaC·DR . . . 88

10장: 세계적 클라우드 네이티브 전환 트렌드 — CNCF 공식 정의와 글로벌 레퍼런스로 본

표준화 흐름 **90**

10.1 CNCF 공식 정의와 세계적 기술 표준의 방향 90

 10.1.1 “VM 이관이 아닌 컨테이너·K8s 기반 전환”이 세계적 주류 경로로 확립된
 배경 90

 10.1.2 AI·데이터 공유 플랫폼 기반으로서의 클라우드 네이티브 — 세계적 트렌드 92

10.2 글로벌 공공·엔터프라이즈 사례와 국내 확산 동향 93

 10.2.1 CERN·Spotify·Booking.com·Adidas 등 글로벌 공공·엔터프라이즈 클
 라우드 네이티브 전환 사례 94

 10.2.2 국내 금융·공공의 클라우드 네이티브 확산 동향과 MSAP.ai 채택 흐름 . 95

11장: VM이 여전히 유효한 예외 케이스와 KubeVirt 흡수 전략 **97**

11.1 VM을 유지해야 하는 기술·규제 조건 97

 11.1.1 상이 커널 요구·레거시 드라이버·실시간 OS·Windows 워크로드 97

 11.1.2 하드 멀티테넌시 규제와 Kata Containers·gVisor 대체 조건 99

11.2 KubeVirt로 잔존 VM을 K8s 오브젝트로 흡수하는 아키텍처 102

 11.2.1 KubeVirt IOPS·네트워크 오버헤드 벤치마크와 적합 워크로드 판단 . . . 102

 11.2.2 Pod와 VM을 동일 API·동일 GitOps 파이프라인에서 관리하는 패턴 . . . 104

12장: 결론 — IT 의사결정자 체크리스트와 운영 모델 선택 권장사항 **106**

12.1 4대 운영 모델 체크리스트 기반 현 상태 진단 106

 12.1.1 운영 모델 체크리스트 — 선언형·불변·자동확장·단일 표면 4항 107

 12.1.2 IaaS 전략 가치 체크리스트 — 특수 워크로드 %·3년 라이선스 전
 망·AI/MSA 전제 108

 12.1.3 관리 포인트·목표 아키텍처 체크리스트 — FTE·vMotion 중복·KubeVirt
 흡수 109

12.2 12개월 단위 실행 권장 순서와 결정 프레임 110

 12.2.1 현대화 대상·KubeVirt 흡수·VM 유지 3분류 결정 프레임 111

12.2.2 이사회 제시용 3대 KPI — TCO·집적률·배포 리드타임	112
부록: 전체 출처 목록	113
부록: 전체 출처 목록	113
S01~S10: IaaS·가상화·클라우드 네이티브 기원 및 역사	113
S11~S20: 라이선스·벤더 정책 변화·VM→K8s 브리지 기술	115
S21~S30: 운영 모델·관측성·자동확장·컨테이너 네이티브 기능	116
S31~S40: Operator·Service Mesh·정책 엔진·관측성·보안	117
S41~S50: Kubernetes 릴리스·Kubeflow·Booking.com·CERN·Adi- das·Windows 컨테이너	118
S51~S60: 설치 가이드·AI 스택·관측성·보안·연례 조사·시장 동향	119
S61~S64: 벤더 제품·문서·국산 PaaS·운영 사례	121
Appendix	121
References	121
Glossary	129

1.1 하이퍼바이저 가상화의 탄생: 서버 하드웨어 낭비에서 출발한 1999~2010년대 IaaS 시대

1999년부터 2010년대 초반까지 IT 인프라 환경은 서버 하드웨어의 활용 효율성에 대한 고민이 가장 큰 이슈였습니다. 기업들은 각 애플리케이션마다 별도의 물리 서버를 할당하는 방식으로 시스템을 운영했기 때문에, 실제로는 많은 자원이 유휴 상태로 방치되는 경우가 많았습니다. 이러한 비효율성은 비용 증가와 관리 복잡성으로 이어졌으며, 이를 해결하기 위한 기술적 혁신의 필요성이 대두되었습니다. 이에 따라 하이퍼바이저 기반 가상화 기술이 등장하게 되었고, 이는 곧 IaaS(Infra as a Service) 플랫폼의 탄생으로 이어졌습니다. 이 절에서는 VMware, OpenStack, Nutanix, CloudStack 등 대표적인 IaaS 기술의 등장 배경과 각자의 목적, 그리고 최근 Broadcom의 VMware 인수 이후 변화된 라이선스 정책까지 살펴봅니다. 이를 통해 VM 기반 IaaS의 한계와, 단순한 VM→VM 이주가 문제의 본질을 해결하지 못한다는 점을 기술적·회계적 관점에서 명확히 짚어봅니다.

1.1.1 VMware·OpenStack·Nutanix·CloudStack의 발원과 목적 차이

VMware, OpenStack, Nutanix, CloudStack는 모두 IaaS 시대를 대표하는 플랫폼으로, 각기 다른 배경과 기술적 목표를 가지고 출발하였습니다. 이들의 등장은 서버 하드웨어 자원의 낭비를 해결하고, 데이터센터의 효율성을 극대화하기 위한 공통된 필요에서 비롯되었습니다. 그러나 각 플랫폼은 설계 철학, 운영 방식, 그리고 시장에서의 포지셔닝에 있어 뚜렷한 차이를 보입니다. 이러한 차이점을 이해하는 것은 오늘날의 인프라 전략 수립에 있어 매우 중요합니다.

VMware 창립과 ESX 하이퍼바이저의 등장

VMware는 1998년 Diane Greene, Mendel Rosenblum 등 핵심 인물에 의해 설립되었으며, 2001년 최초의 x86 하이퍼바이저인 ESX를 출시했다. 이 기술은 하나의 물리 서버에 여러 VM을 올려 하드웨어 자원을 효율적으로 분할하는 혁신을 가져왔다. 2003년 VMotion 기능이 등장하면서 VM의 실시간 이동이 가능해져, 데이터센터 운영의 유연성이 크게 향상되었다. VMware는 이후 vSphere, vCenter 등으로 엔터프라이즈 시장을 선점했고, 2023년 Broadcom에 690억 달러에 인수되며 라이선스 모델 변화의 분기점을 맞았다.

VMware의 등장은 단순한 가상화 기술의 도입을 넘어, 데이터센터 운영의 패러다임을 바꾸는

계기가 되었습니다. 하이퍼바이저를 통해 여러 VM을 한 대의 서버에서 운영할 수 있게 됨으로써, 자원 활용률이 극적으로 개선되었습니다. 또한, VMotion과 같은 기능은 무중단 서비스 제공과 장애 대응의 유연성을 크게 높였습니다. 이러한 기술적 진보는 기업 IT 인프라의 표준으로 자리잡았으며, 이후 시장의 경쟁 구도를 형성하는 데 결정적인 역할을 하였습니다.

OpenStack의 오픈소스 IaaS 계보

OpenStack은 NASA Nebula 프로젝트와 Rackspace Cloud Files를 기반으로 2010년 출범했다. Nova(컴퓨트), Neutron(네트워크), Cinder(스토리지), Keystone(인증) 등 20여 개의 컴포넌트로 구성되어, 오픈소스 IaaS 플랫폼의 표준을 제시했다. Red Hat, Canonical, Mirantis 등 다양한 디스트리뷰션이 등장하며, 벤더 종속을 피하려는 대형 조직의 선택지로 자리잡았다.

OpenStack의 가장 큰 특징은 오픈소스 생태계를 기반으로 한다는 점입니다. 이는 기업들이 특정 벤더에 종속되지 않고, 자체적으로 인프라를 구축·운영할 수 있는 자유를 제공합니다. 다양한 컴포넌트 구조는 유연한 확장성과 맞춤형 아키텍처 구성을 가능하게 하였으며, 대규모 클라우드 환경을 필요로 하는 기업과 공공기관에서 널리 채택되었습니다. 오픈소스라는 특성상 활발한 커뮤니티와 빠른 기술 발전이 이루어졌으나, 반면에 복잡한 설치 및 운영, 전문 인력 확보의 어려움 등도 함께 존재하였습니다.

Nutanix의 HCI(하이퍼컨버지드 인프라) 모델

Nutanix는 2009년 창업하여 서버·스토리지·네트워크를 단일 어플라이언스로 통합하는 HCI 모델을 제시했다. Nutanix AHV(독점 하이퍼바이저)와 Prism 관리 콘솔을 통해 VM 단위 운영을 강화했지만, 컨테이너 네이티브와는 구별되는 IaaS 범주에 속한다.

Nutanix의 하이퍼컨버지드 인프라(HCI) 모델은 전통적인 3계층 아키텍처(서버, 스토리지, 네트워크)를 하나의 통합된 솔루션으로 단순화시켰습니다. 이를 통해 인프라 구축과 확장이 훨씬 용이해졌으며, 관리 복잡성도 크게 줄어들었습니다. Nutanix AHV는 자체 개발한 하이퍼바이저로, VMware와의 라이선스 경쟁에서 차별화 포인트를 제공하였습니다. 그러나 Nutanix 역시 VM 단위의 자원 운영에 초점을 맞추고 있어, 컨테이너 기반의 클라우드 네이티브와는 본질적으로 구분됩니다.

CloudStack의 통신사·호스팅 중심 확산

Citrix는 2012년 cloud.com을 인수한 후 CloudStack을 Apache 재단에 기부했다. CloudStack은 통신사와 호스팅 사업자 중심으로 확산되었으며, VM 단위로 클라우드 서비스를 제공하는 구조를 유지하고 있다.

CloudStack은 비교적 단순한 구조와 빠른 배포가 가능하다는 장점으로, 통신사와 호스팅 사업자 등 대규모 멀티테넌트 환경에서 주로 채택되었습니다. 오픈소스 프로젝트로서의 유연성과, 다양한 하이퍼바이저 지원 등으로 인해 특정 시장에서 강점을 보였으나, 글로벌 대형 퍼블릭 클라우드와의 경쟁에서는 다소 제한적인 영향력을 보여주었습니다.

laaS 4종 비교 표

제품명	창립/출범	최초 공개 시점	운영 단위
VMware	1998	ESX: 2001	VM
OpenStack	2010	2010	VM
Nutanix	2009	AHV: 2015	VM
CloudStack	2012	2012	VM

VM 단위 운영의 공통 전제

이들 laaS 4종은 모두 VM 단위 운영을 전제로 하며, 서버 하드웨어 낭비를 해결하는 1세대 laaS 계보에 속한다. 즉, “대안 laaS”를 선택해도 해결하는 문제의 본질은 동일하며, 이후 6장에서 논의할 목표 아키텍처(베어메탈 K8s 등)와는 근본적으로 다른 총위임을 명확히 해야 한다.

대안 laaS의 시대적 문제 인식

IT 의사결정자는 “우리가 검토 중인 대안 laaS가 실제로는 같은 시대, 같은 문제를 푸는 사촌일 수 있다”는 경각심을 가져야 한다. VM→VM 이주는 하이퍼바이저만 교체할 뿐, OS 종속성·라이선스·운영 포인트는 그대로 유지된다.

1.1.2 VMWare 를 Broadcom 이 인수 후 상용 laaS 라이선스 지형 변화(1차 정리 지점)

2023년 Broadcom의 VMware 인수는 laaS 시장에 큰 충격을 주었습니다. 기존의 라이선스 정책이 대폭 변경되면서, 많은 기업과 기관들이 비용 부담과 운영 전략 재검토에 직면하게 되었습니다. 이 절에서는 Broadcom 인수 이후 VMware의 라이선스 모델 변화와 그에 따른 시장의 반응, 그리고 이러한 변화가 실제로 IT 인프라 운영과 비용 구조에 어떤 영향을 미치는지 살펴봅니다. 또한, 단순히 하이퍼바이저를 교체하는 VM→VM 이주가 근본적 문제 해결로 이어지지 않는다는 점을 회계적·운영적 관점에서 강조합니다.

2023년 Broadcom의 VMware 인수와 라이선스 변화

2023년 Broadcom은 VMware를 690억 달러에 인수하며, 2024년부터 라이선스 정책을 대대적으로 변경했다. 기존의 영구 라이선스 및 독립 제품 판매를 중단하고, 번들 구독 강제와 코어당 과금 모델로 전환했다. 이로 인해 기존 고객들은 2~10배의 인상률을 경험하고 있으며, AT&T 등 대형 통신사는 라이선스 변화에 반발해 소송까지 제기하는 등 시장 혼란이 지속되고 있다.

Broadcom의 라이선스 정책 변화는 단순한 가격 인상에 그치지 않고, 고객의 선택권을 크게 제한하는 방향으로 진행되었습니다. 영구 라이선스의 폐지와 번들 구독 강제는 기존에 개별적으로 필요한 제품만 도입하던 기업들에게 불리하게 작용하였으며, 코어당 과금 체계는 서버 집적률이 높은 환경에서 비용 부담을 더욱 가중시켰습니다. 이러한 변화는 단기적으로는 탈VMware 흐름을 촉진시키는 요인으로 작용하고 있으나, 실제로는 대안 IaaS 역시 유사한 라이선스 구조와 운영 모델을 가지고 있어, 근본적인 비용 절감이나 운영 효율성 개선으로 이어지지 않는 경우가 많습니다.

과금 구조 변화 요약 표

연도	라이선스 모델	과금 단위	주요 변화점
2023	영구 라이선스	서버/소켓/vCPU	독립 제품 판매, 영구 사용
2024	번들 구독(강제)	코어당 과금	구독 번들, 코어당 최소 단위

탈VMware 흐름과 회계적 강제

이러한 변화는 “탈VMware” 흐름을 시장에 회계적으로 강제하고 있으며, IaaS 간 이주가 단순히 라이선스 모델의 변화만을 가져올 뿐, 운영 모델이나 TCO(총소유비용)의 근본적 개선으로 이어지지 않는다는 점을 강조한다. 실제 논의의 무게중심은 이후 장에서 다루는 운영 모델 전환과 TCO 4대 축(라이선스·전력·인건비·Auto Scale 사이징)으로 옮겨가야 한다.

Broadcom의 라이선스 정책 변화는 단순히 비용 문제를 넘어, IT 인프라 운영의 전략적 방향성에도 큰 영향을 미치고 있습니다. 많은 기업들이 대안 IaaS로의 전환을 모색하고 있지만, 실제로는 VM 단위의 운영 모델과 라이선스 구조가 유사하기 때문에, 기대했던 수준의 비용 절감이나 운영 혁신을 경험하지 못하는 경우가 많습니다. 따라서, 진정한 혁신을 위해서는 VM→VM 이주가 아닌, 운영 모델 자체의 근본적 변화를 고민해야 하며, 이는 이후 장에서 다루는 베어메탈 K8s와 같은 현대적 아키텍처로의 전환 논의로 이어집니다.

1.2 컨테이너·쿠버네티스의 탄생: Google Borg에서 CNCF v1.0 정의까지

2010년대 중반 이후 IT 인프라 환경은 하드웨어 자원 활용의 효율성에서 벗어나, 애플리케이션의 빠른 배포와 확장, 그리고 자동화된 운영에 초점을 맞추게 되었습니다. 이러한 변화의 중심에는 컨테이너와 쿠버네티스(Kubernetes)가 자리잡고 있습니다. 이 절에서는 Linux 컨테이너 기술의 발전 과정과 Docker의 혁신, 그리고 Google Borg에서 시작된 대규모 오케스트레이션 시스템이 Kubernetes와 CNCF(Cloud Native Computing Foundation) 표준으로 이어지는 과정을 살펴 봅니다. 이를 통해 VM 기반 인프라의 한계를 넘어, 베어메탈 K8s와 같은 현대적 아키텍처로의 전환 필요성을 기술적·운영적 관점에서 강조하고자 합니다.

1.2.1 FreeBSD jails → LXC → Docker로 이어진 Linux 컨테이너 계보와 2013년 UX 혁신

컨테이너 기술은 2000년대 초반부터 다양한 형태로 발전해왔으며, 그 중심에는 프로세스와 파일 시스템의 격리, 자원 제한, 그리고 운영의 자동화라는 공통된 목표가 있었습니다. 초기에는 FreeBSD jails와 Solaris Zones와 같은 운영체제 수준의 격리 기술이 등장하였고, 이후 Linux 커널의 cgroups와 네임스페이스 기술을 활용한 LXC가 등장하면서 컨테이너의 경량화와 집적률이 크게 향상되었습니다. 하지만, 이러한 기술들은 여전히 복잡한 설정과 관리의 어려움이 존재하였고, 개발자와 운영자의 생산성을 높이기에는 한계가 있었습니다.

FreeBSD jails와 Solaris Zones의 초기 격리 모델

2000년 FreeBSD jails는 프로세스와 파일 시스템을 격리하는 최초의 OS 레벨 가상화 기술로, 서버 내에서 여러 독립된 환경을 운영할 수 있게 했다. 2004년 Solaris Zones는 네임스페이스와 리소스 제한을 결합하여, 더욱 강력한 격리와 자원 관리 기능을 제공했다.

FreeBSD jails는 서버 내에서 여러 개의 독립된 환경을 운영할 수 있도록 하여, 보안성과 자원 격리 측면에서 큰 진전을 이루었습니다. Solaris Zones는 네임스페이스와 리소스 제한 기능을 결합함으로써, 한 단계 더 발전된 격리와 자원 관리 능력을 제공하였습니다. 이러한 기술들은 이후 컨테이너 기술의 토대가 되었으며, 대규모 서버 환경에서의 효율적 운영을 가능하게 하였습니다.

Linux 컨테이너(LXC)의 등장과 기술적 진화

2008년에는 Linux 컨테이너(LXC)가 cgroups와 네임스페이스 기술을 결합하여, 프로세스·네트워크·파일 시스템·사용자 등 다양한 자원을 격리하고 제한할 수 있게 되었다. LXC는 컨테이너의 경량성과 집적률을 크게 높였지만, 개발자 입장에서는 여전히 복잡한 설정과 관리가 요구되었다.

LXC는 Linux 커널의 네임스페이스와 cgroups 기능을 활용하여, 프로세스, 네트워크, 파일 시스템, 사용자 등 다양한 자원을 격리하고 제한할 수 있게 하였습니다. 이를 통해 한 대의 서버에서 수십, 수백 개의 컨테이너를 효율적으로 운영할 수 있게 되었으며, 자원 활용률이 크게 향상되었습니다. 그러나 LXC는 설정 파일의 복잡성과 관리의 어려움, 표준화된 배포 방식의 부재 등으로 인해, 개발자와 운영자 모두에게 진입 장벽이 높았습니다.

Docker의 이미지 레이어링과 UX 혁신

2013년 Docker(dotCloud, Solomon Hykes)는 LXC를 기반으로 이미지 레이어링과 레지스트리 관리, 표준화된 CLI를 결합하여 개발자 UX를 혁신했다. Docker는 “가상화의 대안”이 아니라, 애플리케이션 패키징·실행 단위의 혁신을 목표로 설계되었다. 이미지 레이어링은 애플리케이션 배포와 롤백, 확장성을 극적으로 개선했으며, 컨테이너 오케스트레이션(K8s)으로 이어지는 기반이 되었다.

Docker의 등장은 컨테이너 기술의 대중화에 결정적인 역할을 하였습니다. Docker는 LXC 기반의 컨테이너를 표준화된 이미지로 패키징하고, 레이어 구조를 통해 효율적인 배포와 롤백을 가능하게 하였습니다. 또한, Docker Hub와 같은 중앙 레지스트리를 통해 이미지의 공유와 재사용이 쉬워졌으며, 표준화된 CLI를 도입함으로써 개발자와 운영자의 생산성을 크게 높였습니다. 이러한 혁신은 컨테이너 오케스트레이션 시스템(Kubernetes 등)의 발전으로 이어졌으며, 오늘날 클라우드 네이티브 아키텍처의 핵심 기반이 되었습니다.

컨테이너 계보 연대표 및 능력 비교 표

연도	기술	주요 능력
2000	FreeBSD jails	프로세스·파일 격리
2004	Solaris Zones	네임스페이스·자원 제한
2008	LXC	cgroups·네임스페이스
2013	Docker	이미지 레이어링·UX 혁신

Docker는 가상화가 아닌 패키징 혁신

컨테이너의 설계 목표는 하이퍼바이저 대체가 아니라, 애플리케이션 배포 단위의 재정의였다. Docker는 VM과 달리 OS 전체를 복제하지 않고, 호스트 커널을 공유하며 애플리케이션 실행 환경만을 분리한다. “Docker는 가상화인가?” 라는 오해는, 기술적 구조와 설계 목표의 차이에서 명확히 해소된다.

Docker의 가장 큰 혁신은 애플리케이션의 배포와 실행 환경을 완전히 표준화하고, 개발과 운영의 경계를 허물었다는 점입니다. VM은 OS 전체를 복제하여 무겁고 느린 반면, Docker 컨테이너는 필요한 실행 환경만을 패키징하여 경량화와 빠른 배포가 가능합니다. 또한, 동일한 이미지를 어디서나 실행할 수 있기 때문에, 개발 환경과 운영 환경의 불일치 문제를 근본적으로 해결할 수 있습니다. 이러한 특성은 DevOps, CI/CD, 마이크로서비스 등 현대적 소프트웨어 개발·운영 방식의 기반이 되었으며, 클라우드 네이티브 패러다임의 핵심으로 자리잡았습니다.

1.2.2 Borg(2003)→Kubernetes(2014)와 CNCF 5대 정의(선언형·불변 인프라)

컨테이너 기술의 발전과 함께, 대규모 분산 환경에서 수많은 컨테이너를 효율적으로 관리하고 자동화하는 오케스트레이션 시스템의 필요성이 대두되었습니다. Google은 일찍이 내부적으로 Borg와 Omega와 같은 대규모 오케스트레이션 시스템을 개발하여, 수십만 노드의 클러스터에서 컨테이너 기반 워크로드를 자동화하는 데 성공하였습니다. 이러한 경험은 오픈소스 프로젝트인 Kubernetes의 탄생으로 이어졌고, 이후 CNCF(Cloud Native Computing Foundation)가 클라우드 네이티브의 표준을 정립하는 데 결정적인 역할을 하였습니다. 이 절에서는 Borg에서 Kubernetes, 그리고 CNCF v1.0 정의에 이르는 사상적·기술적 계보를 살펴보고, 선언형·불변 인프라의 중요성을 강조합니다.

Google Borg와 Omega의 내부 오케스트레이션

Google은 2003년부터 Borg라는 내부 오케스트레이션 시스템을 운영하며, 수십만 노드의 클러스터에서 컨테이너 기반 워크로드를 자동화했다. Borg는 선언형 API, 자동 스케일, 불변 인프라, 서비스 메시 등 현대적 클라우드 네이티브의 핵심 개념을 실전에서 구현했다. 2013년 Omega 프로젝트를 통해 스케줄링·리소스 관리의 유연성을 추가하며, Kubernetes 탄생의 기반을 마련했다.

Borg는 Google의 대규모 인프라 운영 경험에서 탄생한 시스템으로, 컨테이너 기반 워크로드의 배치, 확장, 장애 복구, 자원 최적화 등을 자동화하였습니다. 선언형 API를 통해 원하는 상태를

정의하면, 시스템이 실제 상태를 자동으로 수렴시키는 구조를 가지고 있습니다. Omega는 Borg의 한계를 보완하기 위해 개발된 차세대 오케스트레이터로, 스케줄링과 리소스 관리의 유연성을 강화하였습니다. 이러한 경험은 이후 Kubernetes의 설계 철학에 그대로 반영되었습니다.

Kubernetes의 공개와 CNCF 설립

2014년 Joe Beda, Brendan Burns, Craig McLuckie 등 Google 엔지니어가 Kubernetes를 공개했다. Kubernetes는 Borg의 핵심 개념을 오픈소스화하며, 컨테이너 오케스트레이션의 표준으로 자리잡았다. 2015년 CNCF(Cloud Native Computing Foundation)가 설립되어, Kubernetes를 중심으로 클라우드 네이티브 생태계를 확장했다.

Kubernetes는 컨테이너의 배포, 확장, 복구, 네트워크, 스토리지 등 인프라 운영의 모든 측면을 선언적이고 자동화된 방식으로 관리할 수 있도록 설계되었습니다. 오픈소스 프로젝트로서 빠른 생태계 확장과 다양한 벤더의 참여가 이루어졌으며, 현재는 클라우드 네이티브 인프라의 사실상 표준으로 자리잡았습니다. CNCF는 Kubernetes를 중심으로 다양한 오픈소스 프로젝트를 통합·표준화하며, 클라우드 네이티브 생태계의 발전을 이끌고 있습니다.

CNCF v1.0 공식 정의(2018년)와 5대 요소

2018년 CNCF는 v1.0 공식 정의를 발표하며, 클라우드 네이티브의 표준을 다음 5대 요소로 규정했다:

1. 컨테이너(Container)
2. 서비스 메시(Service Mesh)
3. 마이크로서비스(Microservices)
4. 불변 인프라(Immutable Infrastructure)
5. 선언형 API(Declarative API)

이 5대 요소는 VM 중심 운영과 근본적으로 다르며, “VM 위 K8s”라는 과도기 구조가 아닌, 선언형·불변 인프라를 표준으로 제시한다.

CNCF의 5대 요소는 기존 VM 기반 인프라와는 완전히 다른 운영 모델을 지향합니다. 컨테이너는 OS 전체를 복제하지 않고, 필요한 실행 환경만을 경량화하여 제공합니다. 서비스 메시는 마이크로서비스 간의 트래픽 관리, 보안, 관측을 자동화하며, VM 기반 네트워크의 한계를 극복합니다. 마이크로서비스는 각 서비스가 독립적으로 배포·확장될 수 있도록 하여, 대규모 시스템의 유연성과 확장성을 극대화합니다. 불변 인프라는 운영 중 변경을 금지하고, 빌드 시점에 고정된

아티팩트만을 배포함으로써, 환경 드리프트와 장애의 원인을 근본적으로 차단합니다. 마지막으로 선언형 API는 원하는 상태를 선언하면 시스템이 자동으로 실제 상태를 맞추어 주는 방식으로, 수동 관리와 명령형 운영의 한계를 극복합니다.

CNCF v1.0 5대 요소 × VM 중심 운영 대응 표

요소	VM 중심 운영 대응 불가 이유
컨테이너	OS 전체 복제 대신 커널 공유·경량화
서비스 메시	VM은 네트워크 정책·트래픽 분산 제한
마이크로서비스	VM은 단일 서비스·확장성 한계
불변 인프라	VM은 운영 중 변경·드리프트 누적
선언형 API	VM은 명령형·수동 관리에 의존

쿠버네티스는 도구가 아니라 운영 모델

Kubernetes는 단순한 도구가 아니라, 선언형·불변 인프라 운영 모델의 표준이다. VM→VM 이주는 과도기적 타협일 뿐, 운영 모델의 근본적 전환을 위해서는 컨테이너·K8s 기반 현대화가 필수적이다.

Kubernetes의 도입은 단순히 새로운 소프트웨어를 사용하는 것이 아니라, 인프라 운영의 철학과 방식을 근본적으로 바꾸는 일입니다. VM 기반 인프라에서 벗어나, 선언형·불변 인프라를 중심으로 한 자동화·확장성·유연성의 시대가 도래하였으며, 이는 곧 기업의 경쟁력과 혁신 역량을 좌우하는 핵심 요소가 되었습니다. 따라서, VM→VM 이주에 머무르지 않고, 클라우드 네이티브 아키텍처로의 전환을 적극적으로 모색해야 할 시점입니다.

2장: IaaS와 PaaS의 핵심 개념·용어와 운영 모델의 구조적 차이

2.1 NIST 정의로 정렬한 IaaS·PaaS·SaaS와 VM·Container·Pod의 격리 원리

클라우드 인프라 환경에서 IaaS와 PaaS, 그리고 SaaS는 각기 다른 책임 경계와 운영 단위를 정의하며, 실제 서비스 도입 및 운영 전략에 직접적인 영향을 미칩니다. 이 절에서는 NIST SP 800-145의 공식 정의를 바탕으로, 주요 클라우드 서비스 모델의 구분과 각 계층별 대표 제품을 명확히 정리합니다. 또한 VM, 컨테이너, Pod 등 3계층의 격리 원리를 기술적으로 분석하여, 조직이 인프라 현대화 과정에서 반드시 짚고 넘어가야 할 운영 단위와 관리 책임 경계의 본질을 제시합니다. 이를 통해, 용어 혼재와 벤더별 마케팅 주장에 흔들리지 않고 객관적으로 아키텍처를 평가할 수 있는 기준을 마련할 수 있습니다.

2.1.1 NIST SP 800-145 기준 IaaS와 PaaS의 경계

NIST 정의와 용어 정렬

NIST SP 800-145는 클라우드 서비스 모델을 IaaS, PaaS, SaaS 세 가지로 정의합니다. IaaS(Infrastructure as a Service)는 “프로세싱, 스토리지, 네트워크 등 기본 자원을 제공” 하며, 사용자가 OS와 애플리케이션을 직접 설치·운영합니다. PaaS(Platform as a Service)는 “런타임, 미들웨어, 개발·배포 환경까지 포함된 플랫폼을 제공” 하여 사용자는 애플리케이션 개발·배포에 집중할 수 있습니다. SaaS(Software as a Service)는 “완성된 애플리케이션을 서비스 형태로 제공” 하는 모델입니다. OpenStack과 VMware는 IaaS에 해당하며, Kubernetes 업스트림과 해외 상용 엔터프라이즈 K8s 디스트리뷰션(Red Hat OpenShift, Suse Rancher, VMWare Tanzoo), 그리고 MSAP.ai(투라인클라우드)는 PaaS 층위에 가까운 구조를 갖습니다.

대표 제품 매칭과 관리 책임 경계

IaaS와 PaaS의 경계를 명확히 하기 위해 대표 제품을 NIST 정의에 따라 매칭하면 다음과 같습니다. IaaS는 OpenStack, VMware vSphere, Nutanix AHV 등 서버·스토리지·네트워크

인프라를 VM 단위로 제공하는 제품이 해당됩니다. PaaS는 Kubernetes 업스트림, 해외 상용 엔터프라이즈 k8s 상용 제품, MSAP.ai 처럼 컨테이너·Pod·서비스 오브젝트 단위로 배포·운영이 가능한 플랫폼이 해당됩니다. 관리 책임 경계는 IaaS에서는 OS·미들웨어·애플리케이션까지 사용자가 직접 관리하지만, PaaS에서는 OS·런타임·미들웨어까지 플랫폼이 책임지고 사용자는 애플리케이션 개발·배포에 집중합니다.

IaaS 벤더 마케팅 용어의 검증

많은 IaaS 벤더가 “PaaS 제공”을 주장하지만, 실제로는 VM 단위 인프라 제공에 머무르는 경우가 많습니다. NIST 정의를 기준으로 벤더 제안서의 용어를 검증하면, “PaaS”라는 단어가 등장해도 실제로는 IaaS 범주에 속하는 경우가 많으므로, 반드시 NIST 정의와 관리 책임 경계로 재확인해야 합니다.

NIST 3층 정의 × 대표 제품 매칭 표

서비스 모델	NIST 정의	대표 제품	관리 책임 경계
IaaS	프로세싱·스토리지·네트워크 자원 제공	OpenStack, VMware vSphere, Nutanix AHV	OS·미들웨어·애플리케이션 직접 관리
PaaS	런타임·미들웨어 포함 배포 환경 제공	Kubernetes 업스트림, MSAP.ai, 해외 엔터프라이즈 K8s	애플리케이션 개발·배포 집중, OS·런타임 플랫폼 책임
SaaS	완성 애플리케이션 제공	Salesforce, Google Workspace	사용자: 서비스 이용, 벤더: 전체 운영

실무적 판단 기준

IT 의사결정자는 벤더 제안서에서 “PaaS”라는 용어가 등장할 때, 반드시 NIST 정의와 관리 책임 경계를 기준으로 검증해야 합니다. 실제로는 VM 단위 인프라 제공에 머무르는 경우가 많으므로, 운영 단위와 책임 경계가 명확히 PaaS에 부합하는지 체크리스트로 확인하는 것이 중요합니다.

2.1.2 Linux namespace·cgroups·OverlayFS로 본 컨테이너 격리의 커널 공유 구조

컨테이너 기술은 기존의 가상머신(VM)과는 달리, 운영체제 커널을 공유하면서도 각 컨테이너가 독립된 환경처럼 동작할 수 있도록 다양한 커널 레벨의 격리 메커니즘을 활용합니다. 이 절에서는 Linux 커널의 namespace, cgroups, OverlayFS가 어떻게 컨테이너의 격리와 자원 제한, 이미

지 관리에 기여하는지 구체적으로 설명합니다. 또한 이러한 커널 공유 구조가 가져오는 경량성과 보안상의 트레이드오프, 그리고 이를 보완하기 위한 최신 기술 동향까지 실무적으로 짚어봅니다.

컨테이너 격리의 기술적 구조

컨테이너는 Linux 커널의 namespace와 cgroups, 그리고 OverlayFS 이미지 레이어링을 기반으로 격리와 자원 제한을 구현합니다. namespace는 PID, NET, MNT, UTS, IPC, USER 등 다양한 시스템 리소스를 분리하여 각 컨테이너가 독립된 환경처럼 동작하게 합니다. 예를 들어, PID namespace는 각 컨테이너가 별도의 프로세스 트리를 갖도록 하며, NET namespace는 네트워크 인터페이스와 라우팅 테이블을 분리합니다. MNT namespace는 파일 시스템 마운트 포인트를 분리하여, 컨테이너마다 서로 다른 파일 시스템 뷰를 제공합니다. UTS namespace는 호스트 이름과 도메인 이름을 분리하고, IPC namespace는 프로세스 간 통신 리소스를 분리합니다. USER namespace는 사용자 및 그룹 ID를 분리하여, 컨테이너 내부의 root 권한이 호스트에는 영향을 미치지 않도록 합니다.

cgroups(Control Groups)는 CPU, 메모리, IO 등 자원 사용량을 제한하여 컨테이너 간 자원 경쟁을 제어합니다. 예를 들어, 특정 컨테이너에 CPU 사용량을 1코어로 제한하거나, 메모리 사용량을 512MB로 제한할 수 있습니다. 이를 통해 하나의 컨테이너가 과도하게 자원을 점유하여 다른 컨테이너의 성능에 영향을 주는 것을 방지할 수 있습니다.

OverlayFS는 여러 이미지 레이어를 중첩하여 컨테이너 이미지를 효율적으로 관리하고, 불변성을 강화합니다. 기본 이미지는 읽기 전용으로 유지되며, 컨테이너 실행 시 변경 사항은 별도의 쓰기 레이어에 저장됩니다. 이를 통해 동일한 베이스 이미지를 여러 컨테이너가 공유하면서도, 각 컨테이너의 변경 사항은 독립적으로 관리할 수 있습니다.

커널 공유의 경량성과 보안적 함의

컨테이너는 VM과 달리 OS 커널을 공유합니다. 이 구조 덕분에 VM 대비 컨테이너의 기동 시간은 수백 밀리초수 초로 매우 빠르며, 집적률도 310배까지 높아질 수 있습니다. 컨테이너는 하이퍼바이저를 거치지 않고, 호스트 OS 위에서 직접 실행되므로 오버헤드가 적고, 리소스 효율성이 극대화됩니다. 실제로 대규모 마이크로서비스 환경에서는 수백~수천 개의 컨테이너가 단일 호스트에서 동시에 실행될 수 있습니다.

그러나 커널 공유는 보안 격리 측면에서 VM보다 약할 수 있으며, 커널 취약점이 전체 컨테이너에 영향을 미칠 수 있다는 트레이드오프가 존재합니다. 예를 들어, 컨테이너 내부에서 커널 취약점이 악용될 경우, 호스트 전체 또는 다른 컨테이너로 공격이 확산될 위험이 있습니다. 이를

보완하기 위해 Kata Containers, gVisor 등 하드 멀티테넌시 옵션이 개발되고 있습니다. Kata Containers는 경량 VM을 활용하여 컨테이너마다 별도의 커널을 제공하고, gVisor는 사용자 공간에서 커널 호출을 가로채는 샌드박스 계층을 추가하여 보안성을 강화합니다.

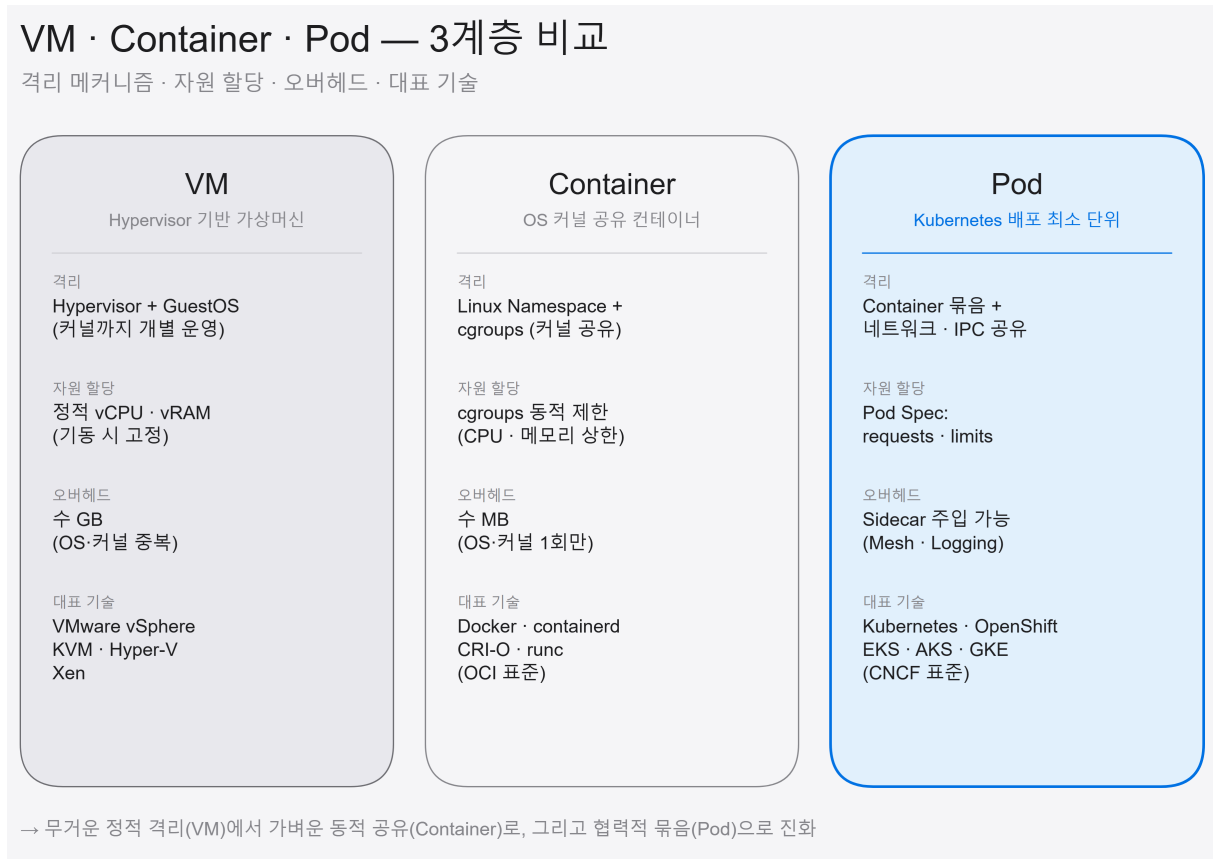
컨테이너 격리 3요소 표

격리 요소	역할	구현 방식	보안 영향
namespace	리소스 분리	PID, NET, MNT, UTS, IPC, USER	컨테이너 간 프로세스·네트워크 완전 분리
cgroups	자원 제한	CPU, 메모리, IO 제한	자원 경쟁·폭주 방지, 과다 사용 차단
OverlayFS	이미지 레이어링	여러 레이어 중첩	불변성 강화, 이미지 변조 방지

경량성과 트레이드오프

컨테이너의 경량성은 커널 공유 구조에서 비롯되며, VM 대비 자원 효율성이 매우 높습니다. 그러나 보안 격리 수준은 VM보다 낮을 수 있으므로, 금융·공공 등 규제 환경에서는 추가적인 격리 옵션(Kata Containers 등)을 검토해야 합니다. 실제로 금융권이나 공공기관에서는 컨테이너 환경 도입 시, 커널 취약점에 대한 대응 방안과 멀티테넌시 보안 정책을 별도로 마련하는 것이 필수적입니다. 또한, SELinux, AppArmor, seccomp 등 리눅스 보안 모듈을 적극적으로 활용하여 컨테이너 내부의 권한 상승이나 시스템 콜 제한을 강화하는 것도 실무에서 중요한 고려사항입니다.

2.1.3 VM·Container·Pod 3계층의 격리 경계와 자원 할당 단위 비교



[그림 1] VM·Container·Pod 3계층 격리 경계와 자원 할당 단위 비교

클라우드 네이티브 환경에서 VM, 컨테이너, 그리고 Pod는 각각 다른 수준의 격리와 자원 할당 단위를 제공합니다. 이 절에서는 각 계층의 기술적 구조와 운영상의 차이, 그리고 조직이 VM 중심에서 Pod 중심으로 전환할 때 마주치는 심리적·운영적 장벽을 구체적으로 분석합니다. 이를 통해 조직이 배포 단위와 격리 경계의 변화에 어떻게 대응해야 하는지 실무적 시사점을 제시합니다.

3계층 구조의 기술적 차이

VM, Container, Pod는 각각 다른 수준의 격리와 자원 할당 단위를 제공합니다. VM은 하이퍼바이저 위에 OS 전체를 포함하여 완전한 격리를 제공하며, 기동 시간은 수십 초~수 분이 소요됩니다. VM은 하드웨어 가상화를 통해 각 인스턴스가 독립된 커널과 사용자 공간을 가지므로, 보안 격리와 장애 격리 측면에서 가장 강력합니다. Container는 OS 커널을 공유하며, 프로세스 수준에서 격리됩니다. 컨테이너는 하이퍼바이저 없이 호스트 OS 위에서 동작하므로, 기동 속도가 빠르고 리소스 효율성이 높지만, 커널 취약점에 대한 보안 이슈가 존재합니다. Pod는 Kubernetes의 최소

배포 단위로, 복수 컨테이너가 네트워크와 스토리지 네임스페이스를 공유합니다. Pod 내 컨테이너들은 localhost 네트워크와 볼륨을 공유하며, 하나의 논리적 애플리케이션 단위로 동작합니다.

배포 단위 전환의 심리적 벽

조직이 “VM 1대=서비스 1개”에서 “Pod=배포 단위”로 전환할 때, 배포 단위와 격리 경계의 변화가 심리적·운영적 저항을 유발할 수 있습니다. 기존 VM 기반 운영에서는 각 서비스가 별도의 VM에서 독립적으로 실행되므로, 장애 발생 시 영향 범위가 명확하고, 운영팀의 책임 경계도 분명합니다. 그러나 Kubernetes 기반 환경에서는 여러 컨테이너가 하나의 Pod로 묶여 배포되고, Pod가 노드 간에 동적으로 이동할 수 있으므로, 장애 분석이나 자원 할당 방식이 달라집니다. 운영팀은 Pod의 라이프사이클, 스케줄링, 오토스케일링 등 새로운 개념을 학습해야 하며, 배포 자동화와 롤백, 롤링 업데이트 등 DevOps 문화와도 밀접하게 연계됩니다.

기술적으로는 Pod가 OS 커널을 공유하면서도 네트워크·스토리지 네임스페이스를 분리하여 충분한 격리와 유연성을 제공합니다. 예를 들어, Pod 내 컨테이너들은 서로 localhost로 통신할 수 있지만, 외부 네트워크와는 별도의 네임스페이스로 분리되어 있습니다. 또한, PersistentVolume을 활용하면 Pod가 재기동되더라도 데이터 일관성을 유지할 수 있습니다.

VM·Container·Pod 3계층 비교 매트릭스

계층	격리 경계	기동 시간	자원 할당 단위	가변성
VM	OS 전체(커널 포함)	수십 초~수 분	VM(고정)	낮음(정적)
Container	OS 커널 공유, 프로세스 수준	수백 ms~수 초	컨테이너(동적)	중간(동적)
Pod	네트워크·스토리지 네임스페이스 공유	수백 ms~수 초	Pod(복수 컨테이너)	높음(동적)

실무적 체크리스트

IT 의사결정자는 “배포 단위가 무엇인가?”를 조직 내에서 체크리스트로 진단하여, VM 중심 운영에서 Pod 중심 운영으로 전환할 때 발생하는 기술적·운영적 변화에 대한 준비 상태를 점검해야 합니다. 예를 들어, 장애 발생 시 영향 범위, 자원 할당 정책, 보안 격리 수준, 배포 자동화 프로세스 등 각 항목에 대해 기존 VM 환경과 Kubernetes 환경을 비교 분석하고, 전환에 따른 조직 내 교육과 프로세스 개선이 필요합니다. 또한, Pod 기반 운영에서는 서비스 디스커버리, 로드밸런싱, 오토스케일링 등 클라우드 네이티브 특화 기능을 적극적으로 활용할 수 있으므로, 기존 운영 방식과의 차이를 명확히 인식하고 단계적으로 도입 전략을 수립해야 합니다.

2.2 명령형 vs 선언형, 가변 vs 불변, 정적 vs 동적 — 운영 모델 3대 축

클라우드 인프라의 운영 모델은 단순히 기술 스택의 변화에 그치지 않고, 시스템 구성과 관리 방식의 근본적인 패러다임 전환을 요구합니다. 명령형과 선언형, 가변과 불변, 정적과 동적이라는 세 가지 축은 현대 인프라 운영의 핵심 원리를 이루며, 각 축마다 VM 기반 IaaS와 컨테이너 기반 PaaS의 구조적 차이가 뚜렷하게 드러납니다. 이 절에서는 Kubernetes가 기본적으로 제공하는 운영 모델의 구현체와, 전통적 VM 환경과의 본질적 차이를 심층적으로 분석하여, 조직이 인프라 현대화의 방향성을 명확히 진단할 수 있도록 안내합니다.

2.2.1 kubectl apply와 컨트롤러 루프로 구현된 멱등성의 의미

Kubernetes의 선언형 운영 모델은 현대 인프라 관리의 복잡성을 획기적으로 줄여주는 핵심 원리입니다. 이 절에서는 kubectl apply와 컨트롤러 루프가 어떻게 멱등성을 보장하는지, 그리고 명령형 운영과의 구조적 차이가 무엇인지 구체적으로 설명합니다. 또한, 실무에서 선언형 운영이 가져오는 효과와 KPI 측정 방법까지 다각도로 분석합니다.

선언형 운영과 멱등성

Kubernetes의 핵심은 선언형 운영 모델입니다. 사용자는 `kubectl apply -f` 명령을 통해 원하는 상태를 YAML 파일로 선언하면, Kubernetes 컨트롤러가 실제 상태와 원하는 상태를 지속적으로 비교(reconcile)하여 수렴시킵니다. 이 구조 덕분에 같은 YAML 파일을 여러 번 적용해도 결과는 항상 동일하며, 이를 멱등성(idempotency)이라고 부릅니다. 멱등성은 운영 중 드리프트(구성 불일치) 사고를 방지하는 핵심 원리입니다.

Kubernetes의 컨트롤러 루프는 지속적으로 클러스터 상태를 모니터링하고, 선언된 스펙과 실제 상태가 일치하지 않을 경우 자동으로 조정 작업을 수행합니다. 예를 들어, 사용자가 3개의 파드를 선언했는데, 장애로 인해 1개가 사라지면 컨트롤러가 즉시 새로운 파드를 생성하여 원하는 상태를 복구합니다. 이러한 자동화 덕분에 운영자는 수동으로 장애를 복구하거나, 반복적인 명령을 실행할 필요가 없습니다.

명령형 vs 선언형 운영 비교

VM 기반 IaaS에서는 이미지 클론, Ansible/Puppet 등 명령형 구성관리 방식이 주로 사용됩니다. 이 방식은 같은 명령을 여러 번 실행할 때 결과가 달라질 수 있으며, 수동 변경이나 스크립트

오류로 인해 드리프트가 누적될 위험이 큼니다. 명령형 방식에서는 작업 순서, 실행 환경, 중간 상태에 따라 결과가 달라질 수 있으므로, 일관성 있는 운영이 어렵습니다. 반면, 선언형 운영은 상태 수렴을 자동으로 보장하므로, 운영 사고의 빈도를 획기적으로 줄일 수 있습니다. 선언형 모델에서는 인프라의 '최종 상태'만 정의하면 되고, 그 상태에 도달하는 과정은 시스템이 자동으로 처리합니다.

명령형 vs 선언형 운영 예시 표

운영 방식	동일 명령 n회 실행 결과	드리프트 발생 가능성
명령형	매번 다를 수 있음	높음
선언형	항상 동일	낮음

실무적 KPI

운영 사고의 재발 빈도, 드리프트 발생 건수 등을 KPI로 삼으면 선언형 운영의 가치가 즉시 드러납니다. 조직은 선언형 운영 모델로 전환하여 관리 효율성과 안정성을 동시에 확보할 수 있습니다. 실제로 대규모 클러스터 환경에서는 선언형 운영을 통해 배포 자동화, 롤백, 롤링 업데이트 등 다양한 운영 시나리오에서 일관성과 신뢰성을 확보할 수 있습니다. 또한, GitOps와 같은 선언형 운영 확장 모델을 도입하면, 코드 기반 인프라 관리와 변경 이력 추적이 가능해져, 감사와 규정 준수 측면에서도 큰 이점을 얻을 수 있습니다.

2.2.2 불변 인프라(Immutable Infrastructure)와 이미지 레이어링의 결합

불변 인프라(Immutable Infrastructure) 원칙은 현대 클라우드 운영에서 신뢰성과 보안성을 높이는 핵심 전략입니다. 이 절에서는 컨테이너 이미지의 불변성, OverlayFS 기반 이미지 레이어링, 그리고 이미지 서명·런타임 변조 탐지 등 공급망 보안 기술이 어떻게 결합되어 운영 관리의 투명성과 신뢰성을 강화하는지 상세히 설명합니다. 또한, 가변 인프라와의 비교를 통해 불변 인프라의 실무적 이점을 구체적으로 제시합니다.

불변 인프라 원칙과 이미지 레이어링

컨테이너 이미지는 빌드 시점에 고정된 아티팩트로, 실행 중에는 수정이 금지되는 불변성(immutable) 원칙을 따릅니다. OverlayFS 레이어링은 여러 이미지 레이어를 중첩하여, 변경이 필요한 경우 새로운 레이어를 추가하는 방식으로 관리합니다. 이 구조는 패치·감사·롤백 등 운영 관리의 투명성과 신뢰성을 높여줍니다.

불변 인프라에서는 운영 중 직접 패치나 수동 수정이 허용되지 않으며, 모든 변경은 소스 코드 또는 빌드 파이프라인을 통해 새로운 이미지를 생성하고, 이를 재배포하는 방식으로 이루어집니다. OverlayFS는 베이스 이미지(예: OS, 런타임)와 애플리케이션 레이어, 설정 파일 레이어 등 여러 계층을 분리하여 관리합니다. 이로 인해 동일한 베이스 이미지를 여러 서비스가 공유할 수 있고, 변경 이력 추적과 롤백이 용이해집니다.

이미지 서명·런타임 변조 탐지

Sigstore/Cosign 같은 이미지 서명·검증 도구는 공급망 보안을 강화하며, Falco는 런타임 변조 탐지 기능을 제공하여 실행 중 위협을 실시간으로 감지합니다. 이미지 서명은 컨테이너 이미지가 신뢰할 수 있는 빌드 파이프라인에서 생성되었음을 증명하고, 배포 전 검증을 통해 무단 변조나 악성 코드 삽입을 차단합니다. Falco와 같은 런타임 보안 도구는 컨테이너 내부에서 비정상적인 시스템 콜이나 파일 접근, 네트워크 활동 등을 실시간으로 모니터링하여, 공격 징후를 조기에 탐지하고 대응할 수 있도록 지원합니다.

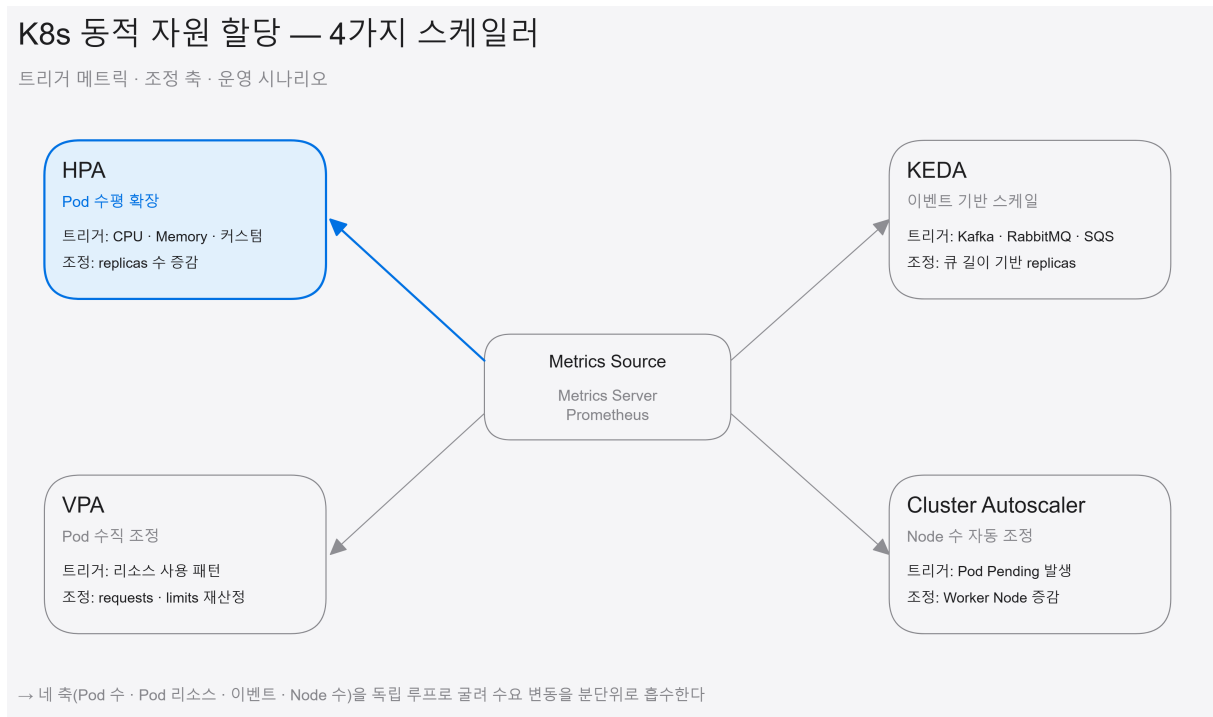
가변 인프라 vs 불변 인프라 비교 표

인프라 유형	패치 방식	감사 투명성	롤백 용이성
가변 인프라	운영 중 직접 수정	낮음	어려움
불변 인프라	이미지 재빌드·재배포	높음	쉬움

감사·운영 관점 이점

불변 인프라는 운영 중 긴급 수정 관행을 빌드 파이프라인으로 우회시키며, 감사·재현·롤백이 간편해집니다. 모든 변경 사항은 코드와 이미지 레이어로 추적 가능하므로, 보안 사고 발생 시 원인 분석과 재현이 용이합니다. 조직은 패치·감사·롤백의 모든 과정에서 투명성과 신뢰성을 확보할 수 있습니다. 또한, 불변 인프라 모델은 DevSecOps와 같은 자동화된 보안·운영 프로세스와도 자연스럽게 결합되어, 인프라 전체의 품질과 신뢰성을 한층 높여줍니다.

2.2.3 Request/Limit·HPA·VPA·KEDA·Cluster Autoscaler의 동적 자원 할당



[그림 2] K8s 동적 자원 할당 — 4종 스케일러 독립 루프|620

Kubernetes의 동적 자원 할당 체계는 전통적 VM 환경의 정적 할당 모델과 근본적으로 다릅니다. 이 절에서는 Request/Limit, HPA, VPA, KEDA, Cluster Autoscaler 등 Kubernetes의 4단계 자동 확장 메커니즘을 상세히 설명하고, VM 환경과의 자원 활용도 차이를 실무적 관점에서 비교합니다. 이를 통해 조직이 클라우드 네이티브 운영의 자원 효율성과 비용 절감 효과를 명확히 이해할 수 있도록 안내합니다.

동적 자원 할당의 4단계 체계

Kubernetes는 컨테이너 Request/Limit, HPA(Horizontal Pod Autoscaler), VPA(Vertical Pod Autoscaler), KEDA(Event-driven Autoscaler), Cluster Autoscaler 등 4단계 자동 확장 체계를 제공합니다. Request/Limit은 컨테이너별 최소·최대 자원 할당을 정의하며, HPA는 CPU·메모리 등 지표 기반으로 파드 수를 자동 조정합니다. VPA는 파드의 리소스 크기를 동적으로 조정하며, KEDA는 메시지 큐·HTTP 요청 등 이벤트 기반 확장을 지원합니다. Cluster Autoscaler는 노드 풀 자체를 확장·축소하여 전체 클러스터의 자원 활용도를 극대화합니다.

Request/Limit은 각 컨테이너가 사용할 수 있는 최소/최대 리소스를 명확히 지정함으로써, 과도한 자원 점유나 오버커밋에 따른 장애를 예방합니다. HPA는 실시간 리소스 사용량을 모니터

링하여, 트래픽 증가 시 자동으로 파드 수를 늘리고, 부하 감소 시 줄입니다. VPA는 워크로드 특성에 따라 파드의 CPU/메모리 크기를 동적으로 조정하여, 리소스 낭비를 최소화합니다. KEDA는 메시지 큐, 이벤트, 외부 트리거 등 다양한 신호에 따라 파드 수를 자동으로 조정할 수 있어, 비동기 처리나 이벤트 드리븐 아키텍처에 적합합니다. Cluster Autoscaler는 파드가 스케줄되지 못할 때 자동으로 노드를 추가하거나, 유휴 노드를 감지해 축소함으로써, 클러스터 전체의 자원 효율성을 극대화합니다.

VM의 정적 할당과 자원 유휴

VM 환경에서는 vCPU·메모리를 피크 부하 기준으로 사전 고정 할당하므로, 피크 아닌 시점에는 자원이 유휴화되어 평균 이용률이 낮아집니다. 예를 들어, 트래픽이 몰리는 특정 시간대에 맞춰 VM을 과도하게 할당하면, 나머지 시간에는 대부분의 자원이 놀게 됩니다. 반면, Kubernetes의 동적 할당 모델은 실시간 부하 변화에 따라 자원을 자동으로 조정하므로, 동일 워크로드에서 평균 이용률을 획기적으로 높일 수 있습니다. 실제로 대규모 서비스에서는 VM 기반 환경의 평균 자원 이용률이 1020%에 불과한 반면, Kubernetes 환경에서는 5070%까지 끌어올릴 수 있습니다.

자동 확장 4단계 매트릭스

자동 확장 방식	트리거	대상	적합 시나리오
HPA	CPU·메모리	파드 수	웹서비스·API 서버
VPA	리소스 사용	파드 리소스 크기	데이터 처리·배치 작업
KEDA	이벤트	파드 수	메시지 큐·비동기 처리
Cluster Autoscaler	파드 스케줄 실패	노드 풀	대규모 확장·축소

실무적 적용과 비교

조직은 “우리 VM 평균 이용률이 10% 미만”이라는 현실을 출발점으로, K8s 동적 할당 모델을 도입하면 자원 효율성과 비용 절감 효과를 실질적으로 체감할 수 있습니다. 자동 확장 체계는 클라우드 네이티브 운영의 핵심이며, VM 기반 모델과의 구조적 차이를 수치로 보여줍니다. 예를 들어, 대규모 이벤트 트래픽이 발생하는 전자상거래 서비스에서는 KEDA와 HPA를 결합하여, 주문량 급증 시 파드와 노드를 자동으로 확장하고, 이벤트 종료 후에는 자원을 신속하게 회수할 수 있습니다. 이처럼 Kubernetes의 동적 자원 할당 모델은 비용 최적화, 운영 자동화, 서비스 안정성 등 다양한 측면에서 조직에 실질적인 경쟁력을 제공합니다.

3장: VM→VM 이주가 중복 투자인 구조적 이유 — OS tax·정적 할당·명령형의 3중 한계

3.1 게스트 OS 중복(OS tax)이 TCO의 숨은 50%를 차지하는 메커니즘

VM 환경에서 발생하는 OS tax는 단순히 라이선스 비용에 국한되지 않고, 운영 인건비, 보안 솔루션, 관리 에이전트 등 다양한 영역에서 중복 비용을 야기합니다. 이 절에서는 VM 기반 인프라에서 Guest OS가 반복적으로 탑재됨으로써 발생하는 비용 구조를 세밀하게 분석합니다. 또한, 컨테이너 환경에서는 Guest OS 자체가 사라지면서 OS 종속 비용 카테고리들이 어떻게 근본적으로 제거되는지 기술적으로 설명하고, VM과 컨테이너 환경의 비용 구조를 정량적으로 비교합니다. 이를 통해 하이퍼바이저나 OS 배포판 교체만으로는 해결할 수 없는 구조적 한계를 드러내고, 컨테이너 전환의 근본적인 비용 절감 효과를 증명합니다.

3.1.1 VM 1대당 수백 MB~수 GB 고정 점유와 100 VM=OS 100개 패치 구조

VM 환경에서는 각 VM이 독립적으로 Guest OS를 탑재합니다. 일반적인 엔터프라이즈 Linux 배포판은 최소 수백 MB에서 수 GB의 디스크 공간을 점유하며, 메모리 역시 OS 커널 및 기본 서비스가 수백 MB 이상을 고정적으로 사용합니다. 예를 들어, 100 VM을 운영하는 조직은 100개의 OS 패치, 100개의 보안 관리, 100개의 라이선스 비용을 동시에 부담해야 합니다. 이 구조는 VM 수가 곧 OS 수로 직결되어, 관리 포인트와 비용이 선형적으로 증가합니다.

컨테이너 환경에서는 호스트 커널을 공유하고, 각 컨테이너는 네임스페이스와 cgroups로 격리됩니다. 이로 인해 동일 하드웨어에서 VM 대비 3~10배의 워크로드 집적률을 달성할 수 있습니다(Google SRE, Borg 사례 참고). 컨테이너는 Guest OS를 반복 탑재하지 않으므로, 디스크 및 메모리 자원 낭비가 대폭 줄어듭니다. 실제 절감 효과는 서버 대수 절감이 아니라 OS 수 감소에서 발생합니다.

많은 조직이 “서버 대수를 줄이면 비용이 절감된다”고 착각하지만, 실제로는 OS 수가 비용의 핵심 드라이버입니다. OS 패치, 보안 솔루션, 라이선스, 인건비 등 모든 운영 비용이 OS 수에 비례해 누적되기 때문입니다. 컨테이너 전환은 OS 수를 획기적으로 줄여 TCO 절감의 근본 원인을 제공합니다.

VM 환경에서 Guest OS가 반복적으로 탑재되는 구조는 단순히 디스크와 메모리 자원 낭비에 그치지 않습니다. 각 Guest OS마다 개별적으로 보안 패치, 모니터링, 접근제어, 백업 등 다양한 관리 작업이 필요하며, 이로 인해 운영팀의 업무 부담과 복잡성이 기하급수적으로 증가합니다. 특히 대규모 엔터프라이즈 환경에서는 OS별로 취약점 관리와 컴플라이언스 대응이 필수적이기 때문에, Guest OS 수가 늘어날수록 인력 투입과 관리 비용이 선형적으로 상승합니다. 반면, 컨테이너 환경에서는 이러한 관리 포인트가 호스트 OS 수준으로 통합되어, 운영 효율성과 보안 대응 속도가 크게 향상됩니다. 실제로 Google, Netflix 등 글로벌 클라우드 네이티브 기업들은 컨테이너 기반 인프라로의 전환을 통해 수천~수만 개 워크로드를 소수의 호스트 OS로 집적하여, 관리 효율성과 비용 절감 효과를 극대화하고 있습니다.

환경	집적률(워크로드 수)	OS 점유(디스크/메모리)
VM	1x	수백 MB~수 GB/VM
컨테이너	3~10x	호스트 커널 공유

3.1.2 금융·공공 필수 OS 종속 솔루션(라이선스·접근제어·SMS·EDR·접속 프로그램)의 VM별 중복 비용 구조

금융 및 공공기관에서는 보안과 내부통제 요건으로 인해 각 서버 OS에 다양한 솔루션을 의무적으로 설치해야 합니다. 이 절에서는 OS 구독, 접근제어(PAM), SMS 모니터링, EDR/백신, 보안 접속 프로그램, 백업·패치 관리, OS 패치 인건비 등 7가지 카테고리가 VM OS 수에 선형적으로 누적되는 구조를 구체적으로 설명합니다. 컨테이너 환경에서는 이 과금 단위가 호스트 노드 수로 재정의되거나 K8s 네이티브 기능으로 대체되어, 구조적으로 비용이 축소되는 원리를 사례와 함께 분석합니다.

금융·공공 분야의 VM 환경에서는 각 VM마다 OS 구독과 유지보수 라이선스가 별도로 필요합니다. 예를 들어, CentOS EOL 이후 엔터프라이즈 Linux 구독 의존도가 높아지면서, VM 100대 운영 시 OS 구독 100개가 필요하게 됩니다. 이는 단순히 라이선스 비용뿐만 아니라, 접근제어(PAM) 솔루션, SMS 모니터링 에이전트, EDR/백신, 보안 접속 프로그램, 백업 및 패치 관리 에이전트 등 다양한 보안 및 관리 솔루션의 라이선스와 운영 비용이 VM OS 수에 따라 선형적으로 증가함을 의미합니다. 각 솔루션은 서버 OS 단위로 에이전트가 설치되고, 이에 따라 관리 포인트와

비용이 중복됩니다.

컨테이너 환경에서는 이러한 솔루션들이 호스트 노드 단위로 통합되거나, Kubernetes의 네이티브 기능(RBAC, ServiceAccount, Audit Log, Prometheus, Falco 등)으로 대체됩니다. 예를 들어, 접근제어는 K8s의 RBAC와 OIDC로, 모니터링은 Prometheus와 Grafana로, EDR/백신은 Falco와 Trivy로, 백업은 Velero로 대체할 수 있습니다. 이로 인해 과금 단위가 VM OS 수에서 호스트 노드 수(수~수십 개)로 대폭 축소되며, 관리 효율성과 보안 수준이 동시에 향상됩니다.

또한, OS 패치 및 업그레이드에 투입되는 인건비 역시 VM 환경에서는 OS 수 × 패치 윈도우당 공수로 누적되지만, 컨테이너 전환 시에는 이미지 재빌드와 롤링 업데이트를 통한 자동화로 운영 부담이 크게 줄어듭니다. 실제로 대형 금융기관의 컨테이너 전환 사례에서는, 기존 100대 VM의 OS 패치 및 보안 관리에 투입되던 인력이 컨테이너 환경에서는 1/10 수준으로 감소하였으며, 보안 사고 대응 속도와 컴플라이언스 준수율이 크게 향상된 것으로 보고되고 있습니다.

과금 단위	VM 환경	컨테이너 환경
OS 구독	VM OS 수 × 단가	노드 수 × 단가

요건	VM OS 단위 PAM	컨테이너 노드 단위	K8s 네이티브 대체
권한 분리·감사 추적	PAM 라이선스	RBAC/ServiceAccount	K8s Audit Log

과금 단위	VM OS 수	컨테이너 노드 수	K8s 네이티브 스택
SMS 에이전트	VM OS 단위	노드 단위	Prometheus 등

요건	VM OS 단위 EDR·백신	컨테이너 이미지 스캔·런타임	K8s 정책(OPA)
취약점 스캔·런타임 탐지	EDR·백신	Trivy·Falco	OPA/Gatekeeper

요건	VM OS 단위 보안 접속	컨테이너 환경 GitOps·kubectl 감사	K8s 네이티브 대체
원격 셸·세션 녹화	보안 접속 프로그램	GitOps·kubectl 감사 로그	RBAC·OIDC

요건	VM OS 단위 에이전트	컨테이너 이미지 재빌드·롤링	K8s Velero·롤링 업데이트
백업·패치	OS 에이전트	이미지 재빌드·롤링	Velero·롤링 업데이트

인건비	VM OS 수 × 공수	컨테이너 이미지 재빌드 자동화	K8s 최소 OS 호스트 노드
패치·업그레이드	OS 수 × 월·분기 공수	자동화 공수	노드 수 × 공수

카테고리	VM OS 수	컨테이너 노드 수	K8s 네이티브 대체
OS 구독	100	10	-
PAM	100	10	RBAC/OIDC
SMS	100	10	Prometheus
EDR/백신	100	10	Falco/Trivy
보안 접속	100	10	GitOps
백업/패치	100	10	Velero
인건비	100 × 공수	10 × 공수	자동화

3.1.3 컨테이너에는 Guest OS가 존재하지 않는다 — VM과의 본질적 구조 차이와 OS 종속 비용 전 카테고리의 구조적 제거

VM과 컨테이너 환경의 가장 본질적인 차이점 중 하나는 Guest OS의 존재 여부입니다. 이 절에서는 컨테이너에는 Guest OS가 존재하지 않는다는 기술적 근거와, 그로 인해 OS 종속 비용 카테고리들이 구조적으로 제거되는 원리를 상세히 설명합니다. 또한, “컨테이너도 OS가 있다”는 오해를 차단하고, 실제로 컨테이너 이미지의 distroless/minimal base가 무엇인지, 그리고 VM Guest OS와 컨테이너 구조가 어떻게 다른지 표와 함께 비교합니다.

VM 환경에서는 하이퍼바이저 위에서 각 VM이 독립적으로 Guest OS(커널, systemd, 패키지 전체)를 구동합니다. 이로 인해 각 VM마다 OS 패치, 보안 에이전트, 접근제어, 모니터링 등 다양한 관리 작업이 반복적으로 필요합니다. 반면, 컨테이너 환경에서는 호스트 커널을 공유하고, 네임스페이스와 cgroups로 격리된 프로세스 단위로 동작합니다. 컨테이너 이미지에 포함된 distroless/minimal base는 시스템 OS가 아니라 애플리케이션 실행 파일과 런타임만을 포함합

니다. 즉, 컨테이너에는 Guest OS가 존재하지 않으며, OS 종속 비용이 구조적으로 사라집니다.

3.1.2에서 분해한 6종 OS 종속 비용 카테고리는 컨테이너 환경에서는 노드 OS(호스트) 수십 개로 과금 단위가 재정의되거나 K8s 네이티브 메커니즘(RBAC, NetworkPolicy, Falco, Trivy, Cosign, OpenTelemetry, Prometheus)으로 대체되어 카테고리 자체가 소멸합니다. 실제로 VM Guest OS 스택(커널, systemd, 패키지, 에이전트들)과 컨테이너(호스트 커널 공유, distroless 이미지, K8s 네이티브 제어)는 기술적 동치가 아닙니다.

“컨테이너도 OS가 있다?”는 오해가 종종 있지만, 컨테이너 이미지에 포함된 distroless/minimal base는 OS가 아니라 애플리케이션 실행 환경에 불과합니다. 예를 들어, distroless 이미지에는 셸이나 패키지 매니저가 없고, 오직 실행 파일과 필수 라이브러리만 포함되어 있습니다. 이로 인해 보안 공격 표면이 최소화되고, 관리 포인트가 대폭 축소됩니다. 실제로 Google, Netflix 등은 distroless 이미지를 표준화하여 운영 효율성과 보안성을 동시에 확보하고 있습니다.

구성 요소	VM Guest OS 스택	컨테이너 구조
커널	VM마다 독립	호스트 커널 공유
systemd/패키지	VM마다 반복	없음, distroless/minimal base
에이전트	OS 단위 설치	K8s 네이티브 기능 대체

카테고리	VM OS 수 과금	컨테이너 환경 대체 경로
OS 구독	OS 수	노드 수
PAM	OS 수	RBAC/OIDC
SMS	OS 수	Prometheus
EDR/백신	OS 수	Falco/Trivy
보안 접속	OS 수	GitOps
백업/패치	OS 수	Velero/롤링 업데이트

3.2 정적 자원 할당과 명령형 배포의 한계

OS tax 외에도 VM 환경은 정적 자원 할당과 명령형 배포라는 두 가지 구조적 한계를 내포하고 있습니다. 이 절에서는 VM 기반 인프라가 왜 평균 자원 이용률이 낮고, 운영 드리프트가 누적되는

지, 그리고 복구·확장 탄력성에서 컨테이너와의 격차가 어떻게 발생하는지 기술적으로 분석합니다. 이를 통해 VM 환경의 비효율성과 컨테이너 환경의 우수성을 구체적으로 비교하고, 실제 운영 사례와 기술적 세부사항을 바탕으로 설명합니다.

3.2.1 vCPU·메모리 사전 할당과 피크 아닌 시점의 유휴 자원 누적

VM 환경에서는 vCPU와 메모리를 사전 고정 할당하는 것이 기본입니다. 이 구조는 피크 부하를 기준으로 자원을 산정하게 하며, 실제 평균 이용률은 10~20% 수준으로 낮게 고착됩니다. 피크 아닌 시점에는 유휴 자원이 누적되지만, 운영 정책상 자원 재할당이나 축소가 어렵습니다.

Kubernetes는 Request/Limit과 Horizontal Pod Autoscaler(HPA)를 통해 자원을 동적으로 할당합니다. 워크로드의 부하에 따라 파드 수와 자원 할당량이 자동으로 조정되어 피크 시점이 아닌 평균 부하에 맞는 자원 사용이 가능합니다.

VM 환경에서 유휴 자원은 실제로 비용 절감에 반영되지 않으며, 구매 승인이나 자원 재할당이 어렵습니다. K8s 환경에서는 자원 활용률이 높아져 데이터센터 공간, 전력, 냉각 비용까지 절감 효과가 발생합니다.

VM 환경의 정적 자원 할당 구조는 실제로 많은 비효율을 초래합니다. 예를 들어, 금융기관의 핵심 업무 시스템은 일일 트래픽 피크에 맞춰 VM 자원을 할당받지만, 대부분의 시간 동안 평균 부하는 피크 대비 10~20%에 불과합니다. 이로 인해 전체 인프라의 80~90% 자원이 유휴 상태로 남게 되며, 이는 전력, 냉각, 하드웨어 감가상각 등 모든 TCO 항목에 불필요한 비용을 발생시킵니다. 반면, 컨테이너 환경에서는 HPA와 같은 자동 확장 기능을 통해 워크로드 부하에 따라 파드 수와 자원 할당량이 실시간으로 조정됩니다. 실제로 대형 이커머스 기업의 사례에서는, 컨테이너 기반 인프라 도입 후 평균 자원 활용률이 20%에서 60% 이상으로 상승하였고, 데이터센터 운영 비용이 30% 이상 절감된 것으로 보고되었습니다.

또한, VM 환경에서는 자원 축소나 재할당이 기술적으로 어렵고, 운영 정책상 승인 절차가 복잡하여 유휴 자원이 장기간 방치되는 경우가 많습니다. 반면, K8s 환경에서는 자원 할당이 선언적으로 관리되며, 필요에 따라 자동으로 확장·축소가 이루어져 운영 효율성이 크게 향상됩니다.

환경	피크 자원	평균 자원	유휴 자원
VM	정적 할당	낮음	많음
K8s	동적 할당	높음	적음

3.2.2 이미지 클론·구성관리 방식의 비역등성과 드리프트 누적

VM 환경에서의 배포 방식은 주로 이미지 클론과 Ansible, Puppet 등 구성관리 도구를 활용한 명령형 배포에 의존합니다. 이 절에서는 명령형 배포가 어떻게 운영 드리프트를 누적시키고, 결과적으로 운영 사고의 구조적 원인이 되는지 분석합니다. 또한, Kubernetes의 선언형 배포와 GitOps 모델이 어떻게 역등성을 보장하고, 운영 효율성과 신뢰성을 높이는지 구체적으로 설명합니다.

명령형 배포 방식에서는 동일한 스크립트를 여러 번 실행해도 결과가 항상 일치하지 않을 수 있습니다. 예를 들어, 운영 중에 수동으로 설정 파일이 변경되거나, 예외적인 패치가 적용된 경우, 구성관리 도구가 이를 감지하지 못하고 일관성 없는 상태가 누적됩니다. 이러한 드리프트는 시간이 지날수록 시스템 간 설정 불일치, 보안 취약점, 예기치 않은 장애로 이어질 수 있습니다. 실제로 대형 금융기관에서 발생한 배포 실패 사례 중 상당수가 “수동 변경 추정”에 기인하며, 이는 명령형 배포의 구조적 한계에서 비롯됩니다.

반면, Kubernetes는 YAML 선언형 배포와 컨트롤러 루프를 통해 역등성을 보장합니다. GitOps(Argo CD, Flux)는 Git을 단일 진실 공급원으로 삼아 클러스터와 동기화하며, 모든 변경 이력이 Git에 기록되어 배포 실패의 원인을 명확하게 추적할 수 있습니다. 롤백 역시 Git 커밋 단위로 손쉽게 수행할 수 있어, 운영 사고 발생 시 신속한 복구가 가능합니다.

또한, 선언형 운영 모델은 변경 관리와 컴플라이언스 준수에 있어 큰 장점을 제공합니다. 모든 배포 및 설정 변경이 코드로 관리되므로, 외부 감사나 보안 심사 시 변경 이력을 투명하게 제출할 수 있습니다. 실제로 글로벌 금융기관들은 GitOps 도입 이후 운영 드리프트 사고 빈도가 급감하였으며, 배포 실패율이 50% 이상 감소한 것으로 보고하고 있습니다.

운영 방식	배포 과정	결과 일관성	드리프트 위험
명령형	이미지 클론+CM	낮음	높음
선언형(GitOps)	YAML+GitOps+컨트롤러	높음	낮음

3.2.3 VM 부팅 지연 vs 컨테이너 기동의 복구·확장 탄력성 격차

VM 환경과 컨테이너 환경의 복구 및 확장 탄력성에는 근본적인 차이가 존재합니다. 이 절에서는 VM 부팅 지연과 컨테이너 기동 시간의 차이가 실제 운영에서 어떤 영향을 미치는지, 그리고 장애 복구(MTTR)와 자동 확장(스케일 아웃)에서 컨테이너 환경이 가지는 우위를 구체적으로 설명합니다.

VM은 부팅에 수십 초에서 수 분이 소요됩니다. OS 커널, 서비스, 각종 에이전트가 모두 초기화되어야 하며, 장애 복구나 확장 시에도 동일한 지연이 발생합니다. 예를 들어, 대규모 트래픽이 몰리는 이벤트성 서비스에서 VM 기반 확장 전략을 사용할 경우, 실시간 트래픽 증가에 대응하지 못하고 서비스 지연이나 장애가 발생할 수 있습니다. 반면, 컨테이너는 수백 밀리초에서 수 초 내에 기동이 완료됩니다. 이는 이미지 프리풀링, 네임스페이스 격리, 런타임 최적화 등 컨테이너 기술의 발전 덕분에 가능해졌습니다.

복구·확장 탄력성의 구조적 차이도 큼니다. VM 환경에서는 장애 복구(MTTR)와 스케일 아웃 응답성이 부팅 지연에 의해 제한됩니다. 실제로 장애 발생 시 신규 VM을 기동하는 데 수 분이 소요되어, 서비스 복구가 지연되는 사례가 빈번합니다. K8s 환경에서는 장애 발생 시 새로운 파드가 즉시 투입되어 MTTR이 대폭 단축되고, HPA 및 Cluster Autoscaler를 통한 자동 확장이 실시간으로 이루어집니다. 예를 들어, 글로벌 이커머스 기업의 블랙프라이데이 트래픽 대응 사례에서는, 컨테이너 기반 인프라가 실시간 트래픽 변화에 맞춰 수십 초 내에 수백 개 파드를 자동 확장하여, 무중단 서비스를 제공한 바 있습니다.

자동 확장이 이론상 가능하다는 것과 실제 요청 피크에 맞춰 확장이 작동한다는 것은 기동 시간의 차이에서 결정됩니다. 컨테이너 환경은 실제 운영 효과가 뛰어나며, 장애 복구와 확장 탄력성 측면에서 VM 환경 대비 월등한 경쟁력을 보유하고 있습니다.

환경	기동 시간	MTTR 영향	스케일 응답성
VM	수십 초~수 분	느림	제한적
컨테이너	수백 ms~수 초	빠름	실시간

3.3 엔터프라이즈 Linux 구독 비용 상승·CentOS EOL 압력에 대한 근본 해결 — 다른 Linux로의 이전이 아니라 OS 없는 시스템 환경으로의 전환

최근 CentOS EOL과 글로벌 엔터프라이즈 Linux 구독 비용 상승은 VM 기반 운영 환경에 심각한 구조적 리스크를 초래하고 있습니다. 이 절에서는 단순히 OS 배포판을 교체하는 방식이 반복적인 리프트앤시프트 구조의 한계에 머무를 수밖에 없는 이유를 분석합니다. 또한, 컨테이너 전환을 통한 OS 자체 제거가 어떻게 근본적인 해결책이 되는지, 그리고 VM→VM 이주와 OS 교체가 반복될 때와 컨테이너 기반 아키텍처로 전환할 때의 효과를 정량적으로 비교합니다.

3.3.1 CentOS EOL과 글로벌 엔터프라이즈 Linux 구독 비용 상승 — VM 운영이 떠안게 된 2020년대 OS 리스크

2020년 말 CentOS 8 EOL 및 CentOS Stream 전환 발표로 무상 엔터프라이즈 호환 Linux 트랙이 사실상 종료되었습니다. 이후 글로벌 엔터프라이즈 Linux 벤더들은 구독 모델을 강화하고, 가격 인상을 단행하며 소스 접근 정책을 변경하고 구독 단위를 재정의했습니다. VM 기반 운영에서는 OS 수 × 구독 × 인건비가 모두 이 시장 변화에 직격됩니다.

OS 선택의 자유도가 줄어들고, 구독 비용이 OS 수에 비례해 상승합니다. 3년 예산 전망에 이 리스크를 반영해야 하며, VM 환경에서는 OS 구독 비용이 전체 인프라 비용의 상당 부분을 차지하게 됩니다.

CentOS EOL 이후 등장한 Rocky Linux, AlmaLinux 등 대체 배포판은 일시적으로 무상 엔터프라이즈 호환성을 제공하지만, 글로벌 벤더의 구독 정책 강화와 가격 인상 추세는 계속되고 있습니다. 실제로 Red Hat, SUSE 등 주요 벤더들은 소스 접근 정책을 변경하고, 구독 단위를 vCPU, 소켓, VM 단위로 세분화하여 라이선스 비용을 높이고 있습니다. 이러한 시장 변화는 VM 기반 운영 환경에서 OS 수가 많을수록 비용 부담이 기하급수적으로 늘어나는 구조적 리스크를 내포합니다. 특히 금융·공공기관과 같이 대규모 VM 인프라를 운영하는 조직에서는, OS 구독 비용이 전체 인프라 TCO의 30~50%를 차지하는 사례도 보고되고 있습니다.

연도	주요 사건
2020	CentOS 8 EOL, CentOS Stream 전환

2021	Rocky/AlmaLinux 등장
2022~	글로벌 벤더 구독 정책 강화

3.3.2 Rocky Linux·AlmaLinux 등 대체 Linux 이전의 한계 — “OS 갱신 비용”이 반복되는 리프트앤시프트

CentOS EOL 대응으로 Rocky Linux, AlmaLinux, Oracle Linux 등 대체 배포판 이전이 단기적으로는 무상 엔터프라이즈 호환성을 복원합니다. 하지만 VM 전수 대상 OS 교체, 검증, 호환성 테스트 인건비가 반복되고, 기존 OS 종속 에이전트(3.1.2의 6종)의 호환성·재인증 비용, 배포판 지속가능성·보안 대응 속도 리스크, 금융·공공 내부통제 심사 갱신 비용이 중복 투자로 반복됩니다.

OS만 바꾸면 된다는 임시방편은 3년 뒤 같은 이슈를 재현시키는 구조적 반복입니다. 하이퍼바이저 교체만으로 3중 한계를 해결할 수 없는 것과 동일하게, OS 교체만으로 OS tax를 해결할 수 없습니다.

대부분의 다른 Linux 로 이전하는 프로젝트는 실제로 많은 숨은 비용과 리스크를 내포하고 있습니다. 예를 들어, 대규모 VM 인프라를 운영하는 금융기관이 CentOS에서 Rocky Linux로 이전할 경우, 수백~수천 대 VM의 OS 교체와 검증, 호환성 테스트, 내부통제 심사 갱신에 막대한 인건비와 시간이 소요됩니다. 또한, 기존에 설치된 접근제어, 모니터링, 백신 등 OS 종속 에이전트의 호환성 재인증과 라이선스 갱신이 필수적으로 요구됩니다. 만약 대체 배포판의 커뮤니티 지원이 약화되거나, 보안 패치가 지연될 경우, 전체 인프라의 보안 리스크가 급격히 증가할 수 있습니다. 이러한 구조적 한계로 인해, OS 교체만으로는 근본적인 비용 절감이나 운영 효율성 개선을 기대하기 어렵습니다.

항목	대체 Linux 이전 프로젝트	컨테이너 전환 프로젝트
OS 교체·검증 인건비	반복	최초 1회
에이전트 재인증	반복	대체·축소
내부통제 갱신	반복	K8s 네이티브 대체

3.3.3 근본 해결 경로 — 컨테이너 전환으로 OS 자체가 없는 시스템 환경 구축

컨테이너 전환은 단순한 OS 교체와는 차원이 다른 근본적인 효과를 제공합니다. 이 절에서는 3.1.3에서 확인한 “컨테이너에는 Guest OS가 없다”는 구조적 사실을 3.3.1~3.3.2의 시장 리스크 대응 논리와 결합하여, 엔터프라이즈 Linux 구독 비용 상승 압력에 대한 근본 해결책이 무엇인지 설명합니다. 즉, “어떤 Linux로 갈 것인가”가 아니라 “애플리케이션을 컨테이너로 옮겨 Guest OS가 존재하지 않는 운영 단위로 재정의” 하는 것이 핵심임을 강조합니다.

이행 아키텍처로는 K8s 호스트 노드 OS(수십 개 규모, 최소 설치 또는 K8s 전용 OS)만 유지하고, 워크로드는 distroless/minimal base 컨테이너 이미지로 빌드하는 표준 패턴을 제시합니다. 호스트 노드 OS 자체도 향후 Bottlerocket, Talos, Flatcar 같은 K8s 최소 OS로 수렴 가능합니다.

컨테이너 전환은 단순히 OS tax를 제거하는 것에 그치지 않고, 운영 효율성, 보안성, 자동화 수준을 획기적으로 개선합니다. 예를 들어, 글로벌 클라우드 사업자들은 이미 K8s 최소 OS(Bottlerocket, Talos 등)를 도입하여, 호스트 OS 관리 포인트를 최소화하고, 워크로드는 distroless 이미지로 표준화하여 운영 효율성과 보안성을 동시에 확보하고 있습니다. 이러한 구조적 전환은 OS 교체나 구독 갱신과 달리, 3년, 5년, 10년 후에도 반복 투자와 리스크가 재현되지 않는 근본적인 해결책임을 실제 사례와 함께 입증할 수 있습니다.

경로	3년 TCO	재현 리스크	호환성 보증
구독 갱신	높음	반복	벤더 종속
대체 Linux 이전	중간	반복	불확실
컨테이너 전환	낮음	소멸	K8s 표준

단계	OS 유형
일반 Linux	Ubuntu/CentOS/RHEL
K8s 최소 OS	Bottlerocket/Talos/Flatcar

4장: 하이퍼바이저·IaaS 진영 경쟁 비교 (VMware·OpenStack·Nutanix·CloudStack)

4.1 IaaS 4종 비교 매트릭스 — 라이선스·커뮤니티·에코시스템의 구조적 차이

IaaS(Infra as a Service) 4종(VMware, OpenStack, Nutanix, CloudStack)은 모두 VM 단위 운영이라는 공통된 전제를 갖고 있으나, 각 플랫폼은 라이선스 모델, 커뮤니티 구조, 시장 포지셔닝, 그리고 Kubernetes(K8s) 지원 방식에서 뚜렷한 차이를 보입니다. 이 절에서는 각 진영의 라이선스 유형과 운영 단위, K8s 지원 방식, 시장 동향을 비교 매트릭스 형태로 정렬하여, 기술·비용·생태계 관점에서 의사결정자가 각 플랫폼의 본질적 차이를 명확히 인지할 수 있도록 합니다. 특히, “탈VMware” 대안 검토 시 단순 IaaS 재선택이 근본적 문제 해결이 아님을 강조하며, 이후 장에서 다루는 목표 아키텍처와 연결되는 논리적 기반을 제공합니다. 각 플랫폼의 구조적 특성과 시장 내 위치를 종합적으로 비교함으로써, 조직의 인프라 전략 수립에 실질적인 참고 자료를 제공합니다.

4.1.1 VMware vSphere/VCF — IaaS 4종 중 유일한 독점 상용 라이선스 포지션

VMware는 IaaS 4종 중 유일하게 독점 상용 라이선스 모델을 채택하고 있습니다. VMware vSphere 및 VMware Cloud Foundation(VCF)은 코어당 과금, 번들 구독, 그리고 라이선스 갱신이 필수적인 구조를 갖고 있습니다. 이 모델은 OpenStack, CloudStack, Nutanix와 달리 Apache 2.0 오픈소스 라이선스나 HCI 어플라이언스 기반 라이선스와는 근본적으로 다릅니다. VMware의 라이선스 정책 변화(2024년 Broadcom 인수 이후 구독 번들 강제, 코어당 과금 전환 등)는 IaaS 시장에서 가격 구조 및 운영 모델에 중대한 영향을 미치고 있습니다.

VMware의 핵심 운영 단위는 VM이며, 하이퍼바이저 기반의 인프라 관리와 고가의 상용 지원 체계가 특징입니다. VMware는 엔터프라이즈 시장에서 높은 신뢰성과 안정성을 제공하지만, 최근 라이선스 정책 변화로 인해 비용 부담이 크게 증가하였고, 탈VMware 대안 검토가 확산되고 있습니다. 이 진영은 라이선스 모델 자체가 다른 IaaS와 구조적으로 다르다는 점을 반드시 인식해야 합니다.

VMware는 Tanzu Kubernetes Grid(TKG) 등 K8s 지원 솔루션을 제공하지만, 기본적으로 VM 위에서 K8s를 배포하는 구조를 유지합니다. 이는 기존 인프라 환경을 크게 변경하지 않고 컨테이너 기반 워크로드를 도입할 수 있다는 장점이 있지만, VM 단위 자원 할당과 라이선스 비용, 그리고 복잡한 관리 구조로 인해 전체적인 TCO가 높게 형성되는 경향이 있습니다. 특히, VMware의 K8s 지원은 하이퍼바이저 기반 인프라의 안정성을 유지하면서도, 현대적인 클라우드 네이티브 워크로드를 수용할 수 있도록 설계되어 있습니다. 그러나 이러한 구조는 오히려 컨테이너의 집적률과 자동화 이점을 충분히 활용하지 못하는 한계도 내포하고 있습니다.

IaaS 플랫폼	라이선스 유형	운영 단위	K8s 지원 방식
VMware	상용 구독	VM	VM 위 K8s
OpenStack	Apache 2.0	VM	VM 위 K8s
Nutanix AHV	HCI 어플라이언스	VM	VM 위 K8s
CloudStack	Apache 2.0	VM	VM 위 K8s

4.1.2 OpenStack — Apache 2.0 IaaS와 디스트로 파편화(RHOSP→RHOSO 이행)

OpenStack은 Apache 2.0 라이선스 기반의 오픈소스 IaaS 플랫폼으로, Nova(컴퓨트), Neutron(네트워크), Cinder(스토리지), Keystone(인증) 등 20개 이상의 핵심 컴포넌트로 구성되어 있습니다. 오픈소스이므로 벤더 종속이 없다는 명제가 있지만, 현실에서는 Red Hat, Canonical, Mirantis 등 주요 디스트리뷰터 중심으로 파편화가 발생하며, 각 디스트로마다 지원 범위와 업그레이드 경로가 달라지는 문제가 있습니다.

2024년 Red Hat은 RHOSP(Red Hat OpenStack Platform) 신규 판매를 중단하고 RHOSO(Red Hat OpenShift Service on OpenStack)로 이행하는 전략을 발표했습니다. 이는 OpenStack 워크로드조차 K8s 기반 운영 모델로 수렴하고 있다는 세계적 트렌드의 신호로 해석할 수 있습니다. 디스트로 파편화와 지원 정책 변화는 조직의 장기적 인프라 전략에 중요한 영향을 미치므로, OpenStack 도입 시 디스트로 선택과 이행 경로를 명확히 이해해야 합니다.

OpenStack의 디스트로 파편화는 실제 운영 환경에서 다양한 문제를 야기할 수 있습니다. 예를 들어, Red Hat, Canonical(Ubuntu), Mirantis 등 각 벤더가 제공하는 OpenStack 배포판은

지원 주기, 기능 통합, 보안 패치 제공 방식, 업그레이드 경로 등이 상이합니다. 이로 인해 동일한 OpenStack 기반이라 하더라도, 실제 운영 및 유지보수에 필요한 역량과 비용이 크게 달라질 수 있습니다. 특히, RHOSP에서 RHOSO로의 이행은 단순한 제품 변경이 아니라, OpenStack 워크로드를 K8s 기반으로 통합하는 전략적 전환을 의미합니다. 이는 기존 VM 기반 워크로드와 컨테이너 기반 워크로드의 경계가 점차 사라지고, 인프라 운영의 패러다임이 변화하고 있음을 보여줍니다. 따라서 OpenStack을 도입하거나 운영 중인 조직은 디스트로 선택뿐만 아니라, 장기적인 기술 로드맵과 벤더 지원 정책의 변화를 면밀히 모니터링해야 합니다.

컴포넌트 수	주요 디스트로	2024년 상태 변화
20+	Red Hat, Canonical, Mirantis	RHOSP 신규 판매 중단, RHOSO로 이행

4.1.3 Nutanix AHV·Apache CloudStack — HCI와 통신사 중심 IaaS 포지셔닝

Nutanix는 2009년 창업 이후 서버, 스토리지, 네트워크를 단일 어플라이언스로 통합한 HCI(Hyper-Converged Infrastructure) 모델을 제공하고 있습니다. Nutanix AHV는 독점 하이퍼바이저 라이선스 기반이며, VM 단위 운영을 전제로 합니다. Nutanix는 엔터프라이즈 및 데이터센터 시장에서 높은 확장성과 통합 관리 기능을 내세우지만, 클라우드 네이티브와는 개념적으로 구분됩니다.

CloudStack은 Citrix가 cloud.com을 인수한 후 2012년 Apache 재단에 기부된 오픈소스 IaaS 플랫폼입니다. 주로 통신사, 호스팅 사업자, 대규모 서비스 사업자 중심으로 활용되며, VM 단위 운영이 기본입니다. Kubernetes 지원은 VM 위에서 K8s를 배포하는 형태로 제공됩니다.

Nutanix AHV는 하드웨어와 소프트웨어가 긴밀하게 통합된 HCI 어플라이언스 형태로 제공되어, 초기 구축 및 확장 과정에서의 복잡성을 크게 줄여줍니다. Nutanix Prism과 같은 통합 관리 도구를 통해 인프라 운영의 효율성을 높이고, 장애 대응 및 리소스 확장도 자동화할 수 있습니다. 그러나 Nutanix의 VM 중심 운영 모델은 컨테이너 네이티브 환경과는 다소 거리가 있으며, K8s 지원 역시 VM 기반 위에서 구현됩니다. 반면, Apache CloudStack은 오픈소스 커뮤니티 주도로 발전해 왔으며, 대규모 멀티테넌트 환경에 적합한 네트워크 및 스토리지 가상화 기능을 제공합니다. 특히 통신사와 호스팅 사업자들이 대규모 VM 인프라를 효율적으로 운영하기 위해 CloudStack을 채택하는 사례가 많습니다. 하지만, 두 플랫폼 모두 VM 단위 자원 할당과 관리 구조를 유지하고

있어, 클라우드 네이티브의 자동화·집적률·운영 효율성 측면에서는 한계가 존재합니다.

플랫폼	대상 시장	운영 단위	K8s 지원 형태
Nutanix AHV	엔터프라이즈	VM	VM 위 K8s
CloudStack	통신사·호스팅	VM	VM 위 K8s

4.2 IaaS 4종이 공유하는 “VM 단위 운영”의 한계

IaaS 4종은 모두 VM 단위 운영을 전제로 하며, 하이퍼바이저, 게스트 OS, 컨테이너 런타임의 3중 스택 구조를 반복합니다. 이 절에서는 IaaS 간 이주가 근본적 문제 해결이 아닌, 동일한 구조적 한계를 유지하는 이유를 기술적으로 분석합니다. 최근 시장 동향과 벤더 전략 변화 역시 “운영 모델 교체(VM→K8s)”로 수렴하고 있음을 확인할 수 있습니다. VM 기반 인프라의 구조적 한계와, 이를 극복하기 위한 업계의 변화 방향을 구체적으로 살펴봅니다.

4.2.1 IaaS 4종의 K8s 지원은 모두 “VM 위 K8s” 형태라는 사실

VMware Tanzu(TKG), Nutanix Karbon, OpenStack Magnum, Apache CloudStack CKS 등 IaaS 4종은 모두 VM 위에서 K8s를 배포하는 형태를 기본으로 제공합니다. 이 방식은 하이퍼바이저, 게스트 OS, 컨테이너 런타임이 3중으로 쌓이는 구조를 반복하여, VM 단위 자원 할당, OS 패치·라이선스 부담, 그리고 컨테이너 집적률 저하 등 구조적 비효율을 야기합니다.

이러한 구조적 한계는 VM 단위 운영의 본질적인 제약에서 비롯됩니다. VM 기반 인프라에서는 각 VM마다 별도의 OS가 필요하며, 이로 인해 패치, 보안, 라이선스 관리 등 운영 부담이 누적됩니다. 또한, VM 단위로 자원을 할당하다 보니, 실제 워크로드의 요구와 무관하게 과도한 리소스 예약이 발생할 수 있습니다. 컨테이너 기반 환경에서는 워크로드의 집적률을 높이고, 자원 활용도를 극대화할 수 있지만, VM 위 K8s 구조에서는 이러한 이점을 온전히 누리기 어렵습니다. 특히, 하이퍼바이저와 게스트 OS가 중첩됨으로써, 네트워크 지연 및 스토리지 I/O 오버헤드가 발생할 수 있으며, 이는 대규모 서비스 환경에서 성능 저하로 이어질 수 있습니다. 결과적으로, IaaS 플랫폼 간의 단순 이주만으로는 이러한 구조적 한계를 극복할 수 없으며, 베어메탈 K8s와 같은 새로운 운영 모델로의 전환이 요구됩니다.

IaaS 플랫폼	K8s 지원 제품명	배치 형태	운영 단위
VMware	Tanzu TKG	VM 위 K8s	VM
OpenStack	Magnum	VM 위 K8s	VM
Nutanix AHV	Karbon	VM 위 K8s	VM
CloudStack	CKS	VM 위 K8s	VM

4.2.2 RHOSP→RHOSO 이행과 OpenStack 워크로드의 K8s 통합 — 2024년 IaaS 벤더 지형 변화

2024년 Red Hat은 RHOSP(Red Hat OpenStack Platform) 신규 판매를 중단하고 RHOSO(OpenShift 기반)로 이행하는 전략을 발표했습니다. 이는 OpenStack 워크로드조차 K8s 기반 운영 모델로 수렴하고 있다는 세계적 트렌드의 지표입니다. AT&T 등 대형 통신사들도 OpenStack Helm + K8s 조합으로 핵심 네트워크 워크로드를 이전하는 사례가 늘고 있습니다.

이러한 변화는 단순히 제품 라인업의 변경에 그치지 않고, 인프라 운영 패러다임의 근본적 전환을 의미합니다. OpenStack의 전통적인 VM 기반 워크로드가 점차 컨테이너 기반으로 이전되고 있으며, 이는 운영 자동화, 자원 집적률, 민첩성 측면에서 큰 이점을 제공합니다. 실제로, AT&T와 같은 글로벌 통신사는 OpenStack 기반 네트워크 기능 가상화(NFV) 환경을 K8s 기반으로 재구성하여, 네트워크 서비스의 배포 속도와 유연성을 크게 향상시키고 있습니다. 동시에, VMware 진영에서도 2024년 Broadcom 인수 이후 라이선스 정책이 급격히 변화하면서, 기존 IaaS 고객들이 운영 모델 전환을 고민하게 되었습니다. 이러한 시장 변화는 단순한 IaaS 플랫폼 교체가 아니라, VM 중심에서 K8s 중심의 클라우드 네이티브 운영 모델로의 전환을 가속화하고 있습니다. 즉, IaaS 벤더의 전략 변화는 곧 조직의 인프라 전략과 직결되며, 장기적으로는 K8s 기반의 선언형·자동화 운영 모델이 표준으로 자리잡을 가능성이 높아지고 있습니다.

분기	주요 변화 신호
Q1	RHOSP 신규 판매 중단, RHOSO 이행
Q2	AT&T 등 통신사 K8s 통합 사례 증가
Q3	VMware 라이선스 구조 변화
Q4	K8s 기반 운영 모델 수렴 가속화

4.3 IaaS 진영의 Cloud Native 응답 — KubeVirt는 브리지(Bridge) 기술이지 근본 해결책이 아니다

IaaS 진영은 컨테이너 시대에 대응하기 위해 “K8s 위 VM(KubeVirt)” 라는 브리지 기술을 내놓고 있습니다. 이 절에서는 KubeVirt의 포지셔닝과 기술적 한계, 그리고 과도기 브리지로서의 역할을 명확히 구분하여, 운영 모델 전환의 근본적 해답이 아님을 기술적으로 논증합니다. KubeVirt가 제공하는 기능과 한계를 종합적으로 분석함으로써, 조직이 미래 인프라 전략을 수립할 때 유의해야 할 점을 제시합니다.

4.3.1 KubeVirt의 포지셔닝 — K8s 위에서 VM을 파드 오브젝트로 다루는 브리지 기술

KubeVirt는 Kubernetes CRD(Custom Resource Definition)로 VM을 파드처럼 관리하는 구조를 제공합니다. VM을 K8s 오브젝트로 선언형 YAML로 배포하고, K8s의 API·컨트롤러 루프·GitOps 파이프라인과 연동할 수 있습니다. Red Hat과 CNCF 생태계에서 Early Majority 단계에 진입했으며, VM을 K8s 표면으로 흡수하는 과도기 브리지 역할을 수행합니다.

KubeVirt의 가장 큰 특징은 기존 VM 워크로드를 K8s 환경으로 자연스럽게 이전할 수 있도록 지원한다는 점입니다. 이를 통해 조직은 기존 VM 자산을 보호하면서, 점진적으로 컨테이너 기반 운영 모델로 전환할 수 있습니다. KubeVirt는 VM을 K8s 리소스로 추상화하여, 네이티브 K8s 오브젝트와 동일한 방식으로 배포·운영·모니터링이 가능합니다. 예를 들어, VM의 라이프사이클 관리, 네트워킹, 스토리지 연결 등이 K8s의 선언형 API와 통합되어, DevOps 및 GitOps 파이프라인 내에서 자동화할 수 있습니다. 또한, KubeVirt는 OpenShift, Rancher, Upstream K8s 등 다양한 배포판과 호환성을 제공하며, VM과 컨테이너 워크로드를 단일 클러스터 내에서 혼합 운영할 수 있는 유연성을 제공합니다.

IaaS 4종이 제공하는 “VM 위 K8s”와는 반대로, KubeVirt는 “K8s 위 VM”이라는 방향성을 갖습니다. 이는 기존 VM 자산을 보호하면서 점진적으로 컨테이너 기반 운영 모델로 전환할 수

있는 전략적 다리입니다. 실제로, KubeVirt는 금융, 공공, 통신 등 레거시 VM 워크로드가 많은 산업군에서 과도기적 브리지로서의 역할을 수행하고 있습니다.

Red Hat 공식 벤치마크 기준, KubeVirt의 IOPS는 네이티브 KVM 대비 ~3% 내, 네트워크 지연은 ~5% 내에서 오버헤드가 발생합니다. 이는 기술적으로 성숙 단계에 도달했음을 의미하지만, VM 단위 운영의 한계를 근본적으로 해결하지는 못합니다. 성능 측면에서는 대다수 엔터프라이즈 워크로드에 무리가 없으나, 초고성능이 요구되는 환경에서는 네이티브 VM 대비 미세한 성능 저하가 발생할 수 있습니다. 또한, 운영 자동화와 통합 관점에서는 K8s 네이티브 워크로드와 동일한 수준의 관리 편의성을 제공하지만, VM 특유의 자원 할당 및 OS 관리 부담은 여전히 남아 있습니다.

방향성 비교	laaS 위 K8s (VM 위 컨테이너)	KubeVirt (K8s 위 VM)
구조	하이퍼바이저 위 VM 위 K8s	K8s 위 VM CRD

성능 비교	KVM vs KubeVirt
IOPS	~3% 내 오버헤드
네트워크	~5% 내 오버헤드

4.3.2 브리지의 본질적 한계 — OS tax·정적 할당·명령형 운영이 K8s 표면으로 옮겨질 뿐

KubeVirt로 VM을 K8s에 흡수해도, 게스트 OS tax(패치·라이선스·보안 이미지 관리 부담)는 그대로 유지됩니다. VM 단위 OS 관리, 라이선스 과금, 패치 인건비 등 구조적 비용 항목이 컨테이너 환경에서 제거되는 것과 달리, KubeVirt는 이 부담을 K8s 표면으로 옮길 뿐 근본적으로 해소하지 못합니다.

KubeVirt 기반 VM은 VM 단위 정적 자원 할당을 전제로 하므로, K8s의 HPA(수평 자동확장), KEDA(이벤트 기반 확장) 등 동적 스케일의 경제성을 잃게 됩니다. VM 이미지·스냅샷 기반의 명령형 운영은 선언형 오브젝트 운영 모델의 명등성과 자동화 장점을 실질적으로 깨뜨립니다.

KubeVirt는 하이퍼바이저(KVM), 게스트 OS, K8s 런타임의 3중 스택을 형태만 바꿔 잔존 시킵니다. 결과적으로 “laaS 교체” 를 KubeVirt로 우회해도 TCO(총소유비용), 운영 복잡도, AI

워크로드 정합성 문제는 해결되지 않습니다. KubeVirt 일정에는 반드시 “컨테이너 전환 목표 비율·시점”이 붙어야 하며, 브리지 위에 눌러앉지 않도록 출구 지표를 명확히 해야 합니다.

KubeVirt의 한계는 실제 운영 환경에서 더욱 뚜렷하게 드러납니다. 예를 들어, VM 기반 워크로드는 여전히 OS 패치, 보안 업데이트, 라이선스 관리 등 반복적인 운영 작업이 필요하며, 이는 컨테이너 기반 환경에서 기대하는 자동화 및 경량화와는 거리가 있습니다. 또한, VM 단위로 자원을 할당하는 구조는 워크로드의 동적 확장 및 축소에 비효율적일 수밖에 없습니다. K8s의 HPA나 KEDA와 같은 자동 확장 기능을 적용할 수 없으므로, 리소스 낭비와 관리 복잡도가 증가합니다. 명령형 운영 방식 역시 VM 이미지를 기반으로 한 수동 배포, 스냅샷, 롤백 등 전통적인 관리 패턴을 그대로 유지하게 만듭니다. 이는 선언형 인프라와 멍등성 기반의 자동화 파이프라인 구축에 장애물이 될 수 있습니다. 특히, AI·머신러닝 워크로드와 같이 대규모 자원 집약적 작업에서는 3중 스택 구조로 인한 성능 저하와 복잡성이 더욱 두드러집니다. 따라서, KubeVirt는 레거시 VM 자산을 점진적으로 컨테이너 환경으로 이전하기 위한 과도기적 솔루션일 뿐, 장기적인 인프라 혁신의 궁극적 해답은 될 수 없습니다. 조직은 반드시 컨테이너 전환 목표와 일정을 명확히 설정하고, 브리지 기술에 머무르지 않도록 전략적 로드맵을 수립해야 합니다.

근본 한계 비교	해결 여부
OS tax	잔존
정적 할당	잔존
명령형 운영	잔존
3중 스택	잔존

관점 차이 요약	11장(예외 흡수 전략)	4.3(브리지 한계 비판)
목적	잔존 VM 흡수	근본 한계 비판

5장: 클라우드 네이티브 진영 경쟁 비교 — MSAP.ai 기준 (vs Kubernetes 업스트림·해외 상용 엔터프라이즈 K8s 디스트리뷰션 A·B·C)

클라우드 네이티브 전환은 단순히 VM을 대체하는 수준을 넘어서, 선언형 오브젝트 기반의 운영 모델로 조직의 인프라 구조를 근본적으로 변화시키는 과정이다. 이 장에서는 국내 공공·금융 규제에 맞춘 MSAP.ai(오픈마루·투라인클라우드 공동개발 PaaS 기준 제품)를 중심으로 삼아, Kubernetes 업스트림과 해외 상용 엔터프라이즈 K8s 디스트리뷰션 A·B·C와의 비교를 통해 실제 선택 기준과 운영 모델의 구조적 차이를 분석한다. 특히 기능 유무가 아니라 국내 지원 경계, 한글화, 국산 인증, 총소유비용, AI 네이티브 통합 등 실무적 의사결정 포인트를 중심으로 경쟁 매트릭스를 제시한다. 독자는 이 장을 통해 조직에 가장 적합한 클라우드 네이티브 플랫폼을 5가지 기준으로 자가 진단할 수 있다.

5.1 MSAP.ai 기준 비교 매트릭스 — 경쟁 PaaS(K8s 업스트림·해외 상용 디스트리뷰션 A·B·C) 대비 라이선스·구성·지원 차이

클라우드 네이티브 플랫폼을 선택할 때 단순히 기능의 유무나 오픈소스 여부만으로 결정하는 것은 한계가 있습니다. 실제로는 라이선스 구조, 운영 관리 방식, 멀티클러스터 지원, 국내 지원 체계, 한글화, AI 네이티브 통합 등 다양한 실무적 요소가 복합적으로 작용합니다. 본 절에서는 MSAP.ai를 기준 축으로 삼아 Kubernetes 업스트림, 해외 상용 엔터프라이즈 K8s 디스트리뷰션 A, 그리고 해외 상용 디스트리뷰션 B·C와의 차이를 매트릭스 형태로 구조적으로 비교합니다. 이를 통해 국내 조직이 실제로 고려해야 할 선택 포인트를 명확히 제시하고, 각 플랫폼의 실질적인 경쟁력과 한계점을 구체적으로 파악할 수 있도록 안내합니다.

5.1.1 MSAP.ai — PaaS 기준 제품

MSAP.ai는 오픈마루와 투라인클라우드가 공동 개발한 국내 PaaS 기준 제품으로, Kubernetes 기반 클라우드 네이티브 오케스트레이션을 중심으로 설계되었습니다. 이 제품은 국내 공공·금융

규제에 맞춘 한글화 UI/문서, 국산 인증, 국내 지원 체계를 기본 제공하며, AI 워크로드에 특화된 스택을 통합하여 국내 환경에 최적화된 플랫폼입니다.

MSAP.ai의 구성 요소는 컨테이너 오케스트레이션(K8s), AI 스택(모델 서빙, Vector DB, RAG, MCP, Agent), 운영 관리(관측성, 보안, CI/CD, 백업/DR), 국내 지원(한글화, SLA, 온사이트 대응), CNCF Certified Kubernetes Conformance Program 등재를 통한 표준 API·기능 호환성을 모두 갖추고 있습니다. 특히 AI 네이티브 통합(한글 임베딩, sLLM 추론, GPU MIG 분할 등)은 국내 조직의 AI 파일럿·실무 적용에 중요한 경쟁력입니다. 예를 들어, 한글 자연어 처리에 특화된 임베딩 모델이나, 국내 데이터 규제 환경에 맞춘 데이터 거버넌스 기능은 해외 제품에서 쉽게 제공되지 않는 차별점입니다.

공공·금융의 규제·운영 요건에 맞춘 한글화 UI/문서, 국산 인증(정보보호, 금융 규제 등), 국내 지원 체계(SLA, 온사이트, 한글 운영 매뉴얼)는 해외 상용 디스트리뷰션과 차별화되는 핵심 요소입니다. 실제로 국내 금융권 및 공공기관에서 MSAP.ai를 도입하는 사례가 증가하고 있으며, 규제 대응과 인력 수급 측면에서 국산 옵션의 가치를 높이고 있습니다. 예를 들어, 국내 모 은행은 해외 상용 제품의 한글화 미비와 느린 장애 대응으로 인해 MSAP.ai로 전환하여 운영 효율성을 크게 개선한 사례가 있습니다.

MSAP.ai는 CNCF Certified Kubernetes Conformance Program에 등재된 공식 인증 디스트리뷰션으로, 표준 K8s API 및 기능 호환성을 CNCF가 공식 검증합니다. 이는 “국산 제품은 표준 호환성이 약하다”는 선입견을 완전히 차단하며, 업스트림 K8s와 동일한 기술적 기반을 보증합니다. 덕분에 오픈소스 생태계와의 연동, 최신 기능 도입, 서드파티 솔루션 통합 등에서도 글로벌 표준과 동일한 수준의 확장성을 확보할 수 있습니다.

MSAP.ai는 해외 상용 디스트리뷰션 대비 총소유비용(TCO) 절감, 국내 지원 접근성, AI 네이티브 통합 등에서 실무적 경쟁력을 갖습니다. 특히 AI 스택의 K8s 정합성, GPU Operator 공식 지원, 멀티클러스터 관리, 베어메탈 지원 등은 국내 조직의 실제 운영 환경에 최적화된 선택지입니다. 예를 들어, GPU 리소스의 세밀한 분할 및 할당, AI 모델의 신속한 배포와 운영, 그리고 장애 발생 시의 신속한 현장 대응은 국내 고객에게 매우 중요한 요소입니다.

아래 포지셔닝 카드 예시는 각 구성 요소별로 MSAP.ai와 경쟁 제품의 커버리지를 한눈에 비교할 수 있도록 정리한 것입니다.

구성 요소	MSAP.ai 커버리지	경쟁 제품 커버리지
컨테이너 오케스트레이션	O	O
AI 스택	O	△ (부분)
운영 관리	O	O
국내 지원	O	X
CNCF Certified K8s	O	O

5.1.2 Kubernetes 업스트림 (Apache 2.0, CNCF Graduated)

Kubernetes 업스트림은 Apache License 2.0 하에서 완전한 오픈소스로 제공되며, 상용 제한이 없습니다. CNCF Graduated 프로젝트로서 GitHub Stars 10만+, 기여자 수만 명, 활발한 릴리스 주기와 커뮤니티 지원을 갖추고 있습니다. 모든 기능은 오픈소스 커뮤니티에서 개발·검증되며, 글로벌 표준으로 자리 잡았습니다.

업스트림 K8s만으로도 상용 제약은 전혀 없으며, 기능적으로는 모든 클라우드 네이티브 워크로드를 지원합니다. 다만 운영 관리, 국내 지원, AI 네이티브 통합 등 실무적 요소는 조직이 직접 구축하거나 별도 지원 계약을 통해 보완해야 합니다. 예를 들어, 업스트림 K8s는 기본적으로 관리 UI, 모니터링, 보안, 백업 등 부가 기능이 포함되어 있지 않으므로, 이를 위해서는 Prometheus, Grafana, Argo CD, Velero 등 오픈소스 툴을 별도로 설치·운영해야 합니다. 또한, 장애 발생 시 공식적인 엔터프라이즈 지원을 받으려면 별도의 컨설팅이나 서드파티 벤더와의 계약이 필요합니다.

MSAP.ai는 업스트림 K8s를 기반으로 엔터프라이즈 기능, AI 스택, 국내 지원을 결합한 제품입니다. 따라서 업스트림 호환성을 유지하면서 국내 환경에 최적화된 부가 기능과 지원을 제공합니다. 조직은 업스트림으로 시작하고 필요 시 MSAP.ai 등 국내 지원 계약을 붙이는 하이브리드 경로를 선택할 수 있습니다. 실제로 일부 공공기관은 개발 환경은 업스트림 K8s로 운영하고, 서비스 환경은 MSAP.ai로 전환하여 안정성과 지원 체계를 강화하는 전략을 채택하고 있습니다.

아래 프로젝트 지표 카드 예시는 업스트림 K8s의 글로벌 영향력과 MSAP.ai의 호환성을 함께 보여줍니다.

항목	Kubernetes 업스트림	MSAP.ai 호환성
GitHub Stars	100,000+	O

기여자 수	수만 명	0
릴리스 주기	3~4개월	0

5.1.3 해외 상용 엔터프라이즈 K8s 디스트리뷰션 A — Route·Build·OAuth·SCC·멀티클러스터 관리 결합형 (경쟁 제품)

해외 상용 엔터프라이즈 K8s 디스트리뷰션 A는 K8s에 Route(네트워크 라우팅), Build(이미지 빌드 파이프라인), ImageStream(이미지 관리), OAuth(인증 연동), SCC(보안 컨텍스트 제약), 멀티클러스터 관리 등 엔터프라이즈 기능을 결합한 제품입니다. 글로벌 벤더의 엔터프라이즈 지원 체계와 함께 구독 기반(소켓/코어/노드 단위 과금)으로 제공됩니다. 이러한 결합형 구조는 대규모 조직에서 복잡한 워크로드와 다양한 보안 요구사항을 충족시키는 데 효과적입니다.

글로벌 벤더는 글로벌 엔터프라이즈 지원(SLA, 온사이트, 다국어 문서)을 제공하지만, 국내 장애 대응 SLA, 한글화, 규제 대응 등은 제한적일 수 있습니다. 예를 들어, 장애 발생 시 영어 기반의 글로벌 지원 센터를 통해 티켓을 접수해야 하며, 국내 현장 대응이나 한글 문서 제공은 추가 비용이 발생하거나 지원이 제한될 수 있습니다. 과금 구조는 소켓/코어/노드 단위로, 대규모 환경에서는 비용이 빠르게 누적될 수 있습니다. 실제로 국내 대형 금융기관의 경우, 노드 수가 증가함에 따라 연간 유지비용이 수억 원에 달하는 사례도 있습니다.

MSAP.ai는 동일한 K8s 기반 엔터프라이즈 기능 커버리지를 국내 지원 채널(한글화, SLA, 온사이트)로 제공하며, 총소유비용(TCO) 절감 효과가 큼니다. 금융·공공의 컴플라이언스 요건이 강한 조직은 해외 상용 지원과 국내 지원 체계의 책임 경계를 계약 단계에서 명확히 비교해야 합니다. 예를 들어, 개인정보보호법, 전자금융감독규정 등 국내 규제 준수 여부와 장애 발생 시의 신속한 대응 가능성은 실제 도입 결정에 큰 영향을 미칩니다.

아래 기능 매핑 표는 주요 엔터프라이즈 기능과 국내 지원 요소를 MSAP.ai와 디스트리뷰션 A 간에 비교한 것입니다.

기능	디스트리뷰션 A	MSAP.ai
Route	0	0
Build	0	0

ImageStream	O	O
OAuth	O	O
SCC	O	O
멀티클러스터 관리	O	O
한글화	X	O
국내 SLA	X	O

5.1.4 해외 상용 디스트리뷰션 B·C — 멀티클러스터 중심형과 가상화 벤더 계열 결합형 (경쟁 제품)

해외 상용 디스트리뷰션 B는 Apache 2.0 오픈소스 기반에 엔터프라이즈 지원을 별도로 결합하며, 멀티클러스터 관리 중심의 아키텍처를 제공합니다. 이 제품은 베어메탈 지원이 완전하며, GPU Operator도 공식 또는 부분 지원합니다. 멀티클러스터 관리 기능은 여러 데이터센터 또는 하이브리드 클라우드 환경에서 단일 콘솔로 여러 K8s 클러스터를 통합 관리할 수 있게 해주어, 대규모 분산 인프라를 운영하는 조직에 적합합니다.

디스트리뷰션 C는 가상화 벤더(vSphere) 계열로, K8s on vSphere를 전제로 독점 라이선스 모델을 채택합니다. VM 위 K8s가 기본 전제로 베어메탈 지원이 제한적이며, GPU 공유도 제약이 있습니다. 즉, 가상화 환경에 최적화되어 있지만, AI·데이터 워크로드와 같이 높은 성능과 리소스 집적률이 요구되는 환경에서는 한계가 있습니다. 예를 들어, GPU를 여러 워크로드에 효율적으로 분할하거나, 베어메탈 환경에서 직접적인 리소스 제어가 필요한 경우에는 디스트리뷰션 C의 구조적 제약이 뚜렷하게 드러납니다.

MSAP.ai는 베어메탈·VM 위 K8s 모두를 지원하며, 국내 AI 워크로드(한글 임베딩, sLLM, GPU MIG 분할 등) 최적화 번들을 기본 포함합니다. 기존 가상화 고객이 디스트리뷰션 C를 기본 경로로 선택할 경우 3중 스택(하이퍼바이저·게스트 OS·컨테이너 런타임) 구조적 한계가 강화됩니다. MSAP.ai는 베어메탈 전환의 국내 표준 옵션으로, 멀티클러스터 관리와 AI 네이티브 통합에서 경쟁력을 갖습니다. 실제로, 베어메탈 환경에서 GPU 리소스의 효율적 분할 및 AI 워크로드의 고성능 처리가 필요한 연구기관이나 대형 데이터센터에서는 MSAP.ai가 선호되고 있습니다.

아래 비교 표는 각 디스트리뷰션의 라이선스, 베어메탈 지원, 멀티클러스터, 가상화 결합, AI

최적화 측면에서의 차이를 요약한 것입니다.

항목	디스트리뷰션 B	디스트리뷰션 C	MSAP.ai
라이선스	Apache 2.0	독점	Apache 2.0
베어메탈 지원	O	X	O
멀티클러스터	O	△	O
기존 가상화 결합	X	O	O
AI 최적화	△	X	O

5.2 공공·금융·AI 관점 제품 선택 기준

클라우드 네이티브 플랫폼 선택에서 베어메탈 지원, GPU Operator 공식 지원, 공공·금융 채택 현황 등은 기술적 기능 이상의 실무적 기준이 됩니다. 본 절에서는 K8s 5종(업스트림, 해외 디스트리뷰션 A·B·C, MSAP.ai)의 베어메탈·GPU 지원 범위와 국내외 공공·금융 채택 현황, 엔터프라이즈 지원 모델을 비교하여 실제 선택 가이드를 제공합니다. AI·데이터 워크로드 비중이 높은 조직에는 GPU 공식 지원이 사실상 1차 제외 기준이 되며, 규제 대응과 인력 수급도 선택의 핵심 변수임을 강조합니다. 또한, 각 플랫폼별로 실제 도입 시 고려해야 할 실무적 이슈와 국내외 사례를 통해 조직별 맞춤형 선택 전략을 안내합니다.

5.2.1 베어메탈 지원 범위와 GPU Operator 공식 지원 수준 비교

베어메탈 지원은 클라우드 네이티브 플랫폼 선택에서 매우 중요한 요소입니다. 업스트림 K8s, 해외 상용 디스트리뷰션 A·B, MSAP.ai는 모두 베어메탈을 완전 지원합니다. 이는 VM 위 K8s의 3중 스택 구조적 한계를 피하고, AI·데이터 워크로드의 집적률과 확장성을 극대화할 수 있는 기반입니다. 예를 들어, 베어메탈 환경에서는 하이퍼바이저 오버헤드 없이 GPU, 네트워크, 스토리지 등의 리소스를 직접 제어할 수 있어, AI 모델 학습 및 대용량 데이터 처리에 최적화된 성능을 발휘할 수 있습니다. 반면, 디스트리뷰션 C(가상화 벤더 계열)는 VM 전제를 기본으로 하며, 베어메탈 지원이 제한적입니다. 이로 인해 리소스 집적률이 낮아지고, 고성능 워크로드에 불리한 구조적 한계가 발생합니다.

NVIDIA GPU Operator 공식 지원 여부는 AI 워크로드의 단가와 성능을 결정하는 핵심 선택

기준입니다. GPU Operator는 Kubernetes 환경에서 GPU 리소스를 자동으로 관리하고, 드라이버·컨테이너 런타임·모니터링 등을 통합 제공합니다. 디스트리뷰션 A와 MSAP.ai는 GPU Operator를 공식 지원하며, GPU MIG 분할, Time-slicing 등 최신 기능을 기본 제공하여, 여러 AI 워크로드가 하나의 GPU를 효율적으로 공유할 수 있습니다. 디스트리뷰션 B는 부분 지원, 디스트리뷰션 C는 제한적 지원에 머뭅니다. 예를 들어, 디스트리뷰션 C에서는 GPU 할당이 제한적이거나, 공식적인 지원이 없어 직접 커스텀 작업이 필요할 수 있습니다. MSAP.ai는 국내 AI 워크로드(한글 임베딩, sLLM 추론, GPU MIG 분할 등) 최적화 번들을 기본 포함하여 실제 파일럿·실무 적용에 강점을 가집니다. 실제로 국내 AI 스타트업이나 연구기관에서는 MSAP.ai의 GPU Operator 공식 지원과 한글 특화 AI 기능을 활용해, 빠른 프로토타이핑과 대규모 모델 운영을 성공적으로 수행한 사례가 있습니다.

AI·데이터 워크로드 비중이 높은 조직에는 GPU 공식 지원이 사실상 1차 제외 기준입니다. MSAP.ai는 국내 AI 워크로드에 최적화된 번들을 기본 제공하며, GPU Operator 공식 지원, 베어메탈·VM 위 K8s 모두 지원, 멀티클러스터 관리 등에서 실무적 경쟁력을 갖습니다. 예를 들어, 대형 병원이나 금융기관에서 AI 기반 데이터 분석을 위해 GPU 집적 환경을 구축할 때, MSAP.ai의 공식 지원과 베어메탈 최적화 기능이 도입 성공의 핵심 요인으로 작용하고 있습니다.

아래 비교 표는 각 플랫폼의 베어메탈 지원, GPU Operator 공식 지원, MIG 지원 여부를 정리한 것입니다.

플랫폼	베어메탈 지원	GPU Operator 공식 지원	MIG 지원
업스트림 K8s	O	O	O
디스트리뷰션 A	O	O	O
디스트리뷰션 B	O	△	O
디스트리뷰션 C	X	X	X
MSAP.ai	O	O	O

5.2.2 국내외 공공·금융 채택 현황과 엔터프라이즈 지원 모델

글로벌 금융·공공 채택 현황을 살펴보면, 해외 상용 디스트리뷰션 A·B는 글로벌 금융·공공에서 상위 채택률을 기록하며, 엔터프라이즈 지원(SLA, 온사이트, 다국어 문서)이 규제 컴플라이언스에

중요한 가치를 제공합니다. 예를 들어, 미국, 유럽, 일본 등 주요 금융기관과 정부기관에서는 글로벌 벤더의 엔터프라이즈 K8s 디스트리뷰션을 통해 안정적인 운영과 규제 준수를 달성하고 있습니다. CNCF 연례 조사에서도 엔터프라이즈 K8s 채택률이 세계적으로 주류화되고 있음을 확인할 수 있습니다.

국내 금융권에서도 대형 은행의 해외 상용 디스트리뷰션 A 채택, 기존 금융기관의 컨테이너 전환, 신규 디지털 금융사의 K8s 기반 시작 등 다양한 사례가 확산되고 있습니다. 예를 들어, 국내 모 대형 은행은 2022년 디스트리뷰션 A를 도입하여 전체 인프라를 컨테이너 기반으로 전환하였고, 신규 금융사는 AI 파일럿 프로젝트에 MSAP.ai를 도입하여 한글 특화 AI 서비스와 신속한 장애 대응을 실현하였습니다. MSAP.ai는 오픈마루·투라인클라우드의 국내 지원 체계와 한글 운영 매뉴얼, 규제 대응 레퍼런스를 바탕으로 국내 공공·금융에서 채택을 확대하고 있는 국산 PaaS 기준 제품입니다. 실제로, 공공기관이나 금융권에서는 국산 인증, 한글화, 국내 SLA, 온사이트 지원 등 국내 환경에 최적화된 요소가 도입 결정에 큰 영향을 미치고 있습니다.

엔터프라이즈 지원 모델 비교 측면에서, 규제 대응과 인력 수급은 기술 선정 기준의 상수이며, 국내 공공·금융은 해외 벤더 종속 리스크를 완화할 수 있는 국산 옵션(MSAP.ai)을 동시에 검토해야 합니다. 해외 벤더 대비 국내 지원 응답성(SLA, 온사이트, 한글 문서)은 실제 장애 대응, 운영 효율, 규제 심사에서 큰 차이를 만듭니다. 예를 들어, 장애 발생 시 국내 엔지니어가 신속하게 현장 대응을 할 수 있는지, 한글 운영 매뉴얼이 제공되어 내부 인력 교육이 용이한지, 규제 심사 시 국산 인증이 유리한지 등이 실제 도입 성공의 핵심 요인입니다.

아래 채택 현황 비교 표는 실제 국내외 조직의 제품 채택 사례와 범위를 정리한 것입니다.

조직	채택 제품	채택 시점	범위
대형 은행	디스트리뷰션 A	2022년	전체 인프라
신규 금융사	MSAP.ai	2023년	AI 파일럿
기존 금융기관	디스트리뷰션 B	2021년	일부 서비스
공공기관	업스트림 K8s	2020년	개발 환경

6장: VM 위 K8s vs 베어메탈 K8s — 3중 스택 과도기 타협과 목표 아키텍처

클라우드 네이티브 전환이 가속화됨에 따라 인프라 운영 모델의 선택은 단순히 하이퍼바이저 종류를 고르는 문제가 아닌, 전체 아키텍처의 효율과 미래 확장성을 결정하는 핵심 의사결정이 되었다. 본 장에서는 “VM 위 K8s”와 “베어메탈 K8s”라는 두 가지 배치 모델이 만들어내는 구조적 차이와 과도기적 타협의 현실, 그리고 목표 아키텍처로서 베어메탈 K8s가 갖는 기술적·운영적 장점을 분석한다. 특히 VM 위 K8s가 기존 VMware 등 가상화 자산을 보호하는 과도기 전략임을 명확히 하고, 궁극적으로는 베어메탈 위 K8s에 KubeVirt 등 브리지 기술로 잔존 VM을 흡수하는 방향이 세계적 표준임을 강조한다. 독자는 자기 조직의 K8s 클러스터가 VM 위인지, 베어메탈 위인지, 그리고 각 모델이 실제로 어떤 비용·복잡성을 야기하는지 한 페이지 다이어그램으로 매핑할 수 있어야 한다.

6.1 VM 위 K8s가 만들어내는 3중 스택의 구조적 비효율

VM 위 K8s는 기존 가상화 인프라의 자산을 보호하면서 컨테이너 기반 운영 모델로 전환하는 과도기적 전략으로 널리 채택되어 왔습니다. 하지만 이 방식은 하이퍼바이저, 게스트 OS, 컨테이너 런타임이라는 세 계층이 중첩되어 운영되기 때문에, 실제 현장에서는 기술적·운영적 비효율이 누적되는 문제가 발생합니다. 본 절에서는 VM 위 K8s가 만들어내는 구조적 부담과 기능 중복, 그리고 벤더별 제품(TKG, Karbon, Magnum, CKS)이 어떻게 락인과 이식성 제한을 유발하는지 구체적으로 분석합니다. 이 분석을 통해 VM 위 K8s 모델이 단기적 타협책임을 이해하고, 장기적으로는 구조적 한계를 극복할 필요성을 인식할 수 있습니다.

6.1.1 하이퍼바이저·게스트 OS·컨테이너 런타임 3계층 누적 구조

VM 위 K8s 환경은 하이퍼바이저(예: VMware vSphere), 게스트 OS(예: RHEL, Ubuntu 등), 그리고 컨테이너 런타임(Docker, containerd, CRI-O)이라는 세 계층이 중첩되어 운영된다. 이 구조는 각 계층마다 별도의 라이선스, 패치, 운영 관리 포인트를 요구하며, 실제로는 단일 워크로드를 위해 세 번의 관리·감사·보안 작업이 중복된다. 하이퍼바이저는 물리 서버의 자원을 분할하고,

게스트 OS는 각 VM마다 독립적으로 운영되며, 그 위에 K8s가 컨테이너를 배포한다. 이 구조는 VM 한 대당 OS 패치·보안·라이선스 비용이 선형적으로 누적되고, 컨테이너 집적률의 이점을 상쇄한다.

이러한 3중 스택 구조는 실제 운영 환경에서 여러 가지 문제를 야기합니다. 첫째, 각 계층별로 관리 포인트가 분리되어 있기 때문에 장애 발생 시 원인 분석이 복잡해집니다. 예를 들어, 컨테이너에서 발생한 네트워크 지연이 하이퍼바이저, 게스트 OS, 컨테이너 런타임 중 어디에서 비롯된 것인지 추적하는 데 많은 시간이 소요될 수 있습니다. 둘째, 보안 측면에서도 각 계층별로 별도의 취약점 관리와 패치가 필요하므로, 전체 시스템의 보안 유지에 더 많은 리소스가 투입됩니다. 실제로 하이퍼바이저와 게스트 OS, 컨테이너 런타임 각각의 보안 패치 주기가 다르기 때문에, 운영팀은 상시적으로 여러 계층의 보안 이슈를 모니터링해야 하며, 이로 인해 운영 복잡도가 크게 증가합니다.

기능 중복과 TCO 낭비

하이퍼바이저 계층에서는 vMotion, HA(High Availability), DRS(Distributed Resource Scheduler) 등 고가의 기능을 제공하지만, K8s 역시 Pod 재스케줄, ReplicaSet, StatefulSet, 자동 복구 등 유사 기능을 네이티브로 제공한다. 두 계층이 동일한 복구·확장 기능을 중복으로 제공하는 것은 명백한 TCO(총소유비용) 낭비이며, 실제로는 장애 복구 시 두 계층의 정책이 충돌하거나 불필요하게 복잡해지는 사례가 빈번하다. 예를 들어, VM HA와 K8s ReplicaSet이 동시에 장애 복구를 시도하면, 복구 순서와 책임 경계가 모호해져 운영 사고가 발생할 수 있다.

이와 같은 기능 중복은 단순히 비용 문제를 넘어, 실제 장애 상황에서 복구 프로세스가 꼬이거나, 불필요한 리소스 소모로 이어질 수 있습니다. 예를 들어, VM HA가 장애를 감지해 VM을 재시작하는 동시에, K8s ReplicaSet이 동일한 파드의 복구를 시도할 경우, 일시적으로 리소스가 과도하게 할당되거나, 복구가 중복 실행되어 오히려 서비스 가용성이 저하될 수 있습니다. 이로 인해 운영팀은 장애 복구 정책을 계층별로 조정하거나, 특정 기능을 비활성화하는 등 추가적인 운영 정책 수립이 요구됩니다.

운영 부담 요소 표

계층	부담 요소
하이퍼바이저	라이선스, 패치, HA 기능 중복
게스트 OS	패치, 보안, 라이선스, 에이전트
컨테이너 런타임	배포, 스케일, 네트워크 관리

동일 기능 중복 구매의 경고

실제 현장에서는 하이퍼바이저와 K8s가 각각 HA 기능을 제공함에도 불구하고, 두 계층 모두에 비용을 지불하고 있다. “동일 기능을 두 번 구매하고 있지 않은가?” 라는 질문은 VM 위 K8s의 구조적 비효율을 직관적으로 드러내며, IT 의사결정자는 이 중복이 조직의 예산과 운영 복잡도에 어떤 영향을 미치는지 반드시 점검해야 한다.

이러한 구조적 비효율은 단기적으로는 기존 자산 보호와 전환 리스크 완화라는 장점이 있지만, 장기적으로는 조직의 IT 비용 구조와 운영 효율성에 부정적인 영향을 미칠 수 있습니다. 특히 대규모 인프라를 운영하는 조직일수록, 이러한 중복 비용과 복잡성은 누적되어 경쟁력 저하로 이어질 수 있으므로, 전략적 관점에서 VM 위 K8s의 한계를 명확히 인식해야 합니다.

6.1.2 TKG·Karbon·Magnum·CKS 제공 형태와 벤더 락인 경로

VM 위 K8s 환경을 지원하는 대표적인 상용 및 오픈소스 제품으로는 VMware Tanzu Kubernetes Grid(TKG), Nutanix Karbon, OpenStack Magnum, Apache CloudStack CKS 등이 있습니다. 이들 제품은 각자의 IaaS 플랫폼 위에서 K8s 클러스터를 손쉽게 배포하고 관리할 수 있도록 다양한 자동화 기능과 관리 콘솔을 제공합니다. 하지만 이러한 편의성 이면에는 벤더 종속성, 즉 락인(Lock-in) 리스크가 내재되어 있습니다. 본 절에서는 각 제품별 배치 전제와 네트워킹, 스토리지, 업그레이드 경로의 벤더 종속성, 그리고 실제 이식성의 한계를 구체적으로 살펴봅니다.

VM 위 K8s 제품별 배치 전제

VMware Tanzu Kubernetes Grid(TKG), Nutanix Karbon, OpenStack Magnum, Apache CloudStack CKS 등 주요 IaaS 벤더들은 VM 위에서 K8s를 배포하는 제품을 제공한다. 이들 제품은 K8s를 VM 기반 인프라에 맞게 재포장하여, 설치·운영·업그레이드·네트워킹·스토리지 경로를 벤더 특화 방식으로 구성한다. 예를 들어, TKG는 vSphere 환경에 최적화된 K8s 클러스터를 생성하며, Karbon은 Nutanix AHV 위에서 K8s를 배포한다. Magnum은 OpenStack Nova VM 위에 K8s를 올리고, CKS는 CloudStack VM 위에 K8s를 배치한다.

이러한 제품들은 각기 다른 하이퍼바이저 및 VM 관리 플랫폼에 최적화되어 있어, 초기 도입 시에는 빠른 배포와 손쉬운 관리가 가능하다는 장점이 있습니다. 그러나 이 과정에서 네트워킹,

스토리지, 업그레이드 경로 등 핵심 인프라 요소가 벤더 특화 방식으로 구성되기 때문에, 표준 K8s 환경과의 호환성이나 이식성에 한계가 발생할 수 있습니다.

벤더 종속적 네트워킹·스토리지·업그레이드 경로

이러한 제품들은 네트워킹과 스토리지 구성에서 벤더 종속성을 강화한다. 예를 들어, TKG는 NSX-T와 통합된 네트워크를 사용하며, Karbon은 Nutanix Files/Volumes와 연동된다. 업그레이드 경로 역시 벤더의 관리 콘솔이나 API를 통해 이루어지므로, 표준 K8s 업스트림과는 다르게 이식성 테스트가 필수적이다. 실제로는 클러스터를 다른 벤더 환경으로 이전하거나 멀티클러스터를 구성할 때 제약이 발생한다.

실제 운영 환경에서는 네트워크 플러그인, 스토리지 드라이버, 인증 및 보안 정책 등 다양한 영역에서 벤더 특화 기능이 적용됩니다. 예를 들어, TKG의 경우 NSX-T와의 통합을 통해 네트워크 세분화 및 보안 정책을 구현할 수 있지만, 이러한 구성이 vSphere 환경에 종속되기 때문에, 동일한 클러스터를 다른 퍼블릭 클라우드나 온프레미스 환경으로 이전할 때 추가적인 마이그레이션 작업이 필요합니다. 마찬가지로, Karbon은 Nutanix Files/Volumes와의 연동을 전제로 하므로, Nutanix 외의 스토리지 솔루션으로의 전환이 어렵습니다.

멀티클러스터·이식성 제한

“VM 위 K8s” 제품들은 멀티클러스터 관리 기능을 제공하지만, 벤더 특화 API와 관리 콘솔에 의존하는 경우가 많아 표준 K8s API와 완전히 호환되지 않는 사례가 존재한다. 이식성 테스트(PoC)는 반드시 배치 전제, 네트워킹, 스토리지, 업그레이드 경로를 포함하여 진행해야 하며, 실제로는 벤더 락인(locked-in) 리스크가 내재되어 있다. “K8s는 표준이므로 락인 없음”이라는 가정은 VM 위 K8s 제품에서는 부분적으로만 성립한다.

이러한 제한은 멀티클러스터 환경에서 더욱 두드러집니다. 예를 들어, 여러 벤더의 VM 위 K8s 클러스터를 통합 관리하거나, 클러스터 간 워크로드 이동을 시도할 때, 벤더 특화 API와 관리 콘솔의 차이로 인해 표준 K8s 방식만으로는 모든 기능을 동일하게 구현하기 어렵습니다. 실제로 멀티클러스터 관리 솔루션을 도입하더라도, 벤더별로 추가적인 커스텀 연동 작업이 필요하며, 이로 인해 운영 복잡성과 유지보수 비용이 증가합니다.

제품 비교 표

제품	배치 전제	네트워킹	스토리지	업그레이드 경로
TKG	vSphere VM	NSX-T	vSAN	vSphere 관리 콘솔

Karbon	AHV VM	Nutanix 네트워크	Nutanix Files/Volumes	Prism 관리 콘솔
Magnum	Nova VM	OpenStack Neutron	Cinder/Swift	OpenStack CLI/API
CKS	CloudStack VM	CloudStack 네트워크	CloudStack 스토리지	CloudStack 관리 콘솔

이 표는 각 제품이 VM 위 K8s를 배치할 때 어떤 경로로 구성되는지, 그리고 표준 K8s와의 이식성/확장성에 어떤 제약이 있는지 한눈에 보여준다.

결론적으로, VM 위 K8s 제품들은 기존 인프라 자산을 활용하는 데에는 효과적이지만, 장기적으로는 벤더 락인과 이식성 제한이라는 구조적 한계를 내포하고 있습니다. 따라서 조직은 도입 시 단기적 편의성뿐만 아니라, 미래의 확장성, 멀티클라우드 전략, 벤더 종속성 해소 방안까지 종합적으로 고려해야 합니다.

6.2 베어메탈 K8s 목표 아키텍처와 필수 구성 요소

베어메탈 K8s는 VM 계층을 제거하고 물리 서버 위에 직접 K8s 클러스터를 배포하는 방식으로, 3중 스택의 오버헤드와 기능 중복을 해소합니다. 이 모델은 자원 집적률, 운영 효율, 장애 복구 속도에서 구조적으로 유리하며, 대규모 글로벌 조직에서 이미 현실적으로 운영되고 있습니다. 본 절에서는 베어메탈 K8s의 필수 구성 요소와 실제 대규모 운영 사례를 통해 현실 가능성과 성공 요인을 분석합니다. 특히, 베어메탈 환경에서 요구되는 오픈소스 솔루션의 직접 구축과 운영 난이도, 그리고 대규모 조직에서의 실전 경험을 바탕으로, 베어메탈 K8s가 왜 목표 아키텍처로 자리잡고 있는지 구체적으로 설명합니다.

6.2.1 MetalLB·Rook-Ceph·Longhorn·cert-manager·External DNS 자체 구성

베어메탈 환경에서 K8s를 운영하기 위해서는 클라우드 환경에서 자동으로 제공되는 여러 인프라 기능을 직접 구축해야 합니다. 대표적으로 로드밸런서, 스토리지, PKI, DNS, Control Plane HA, 백업·복구 등 핵심 인프라 요소를 오픈소스 솔루션이나 상용 제품으로 직접 구성하고 운영해야 하므로, 운영팀의 기술적 역량과 책임이 크게 증가합니다. 이 절에서는 베어메탈 K8s에서 반드시 고려해야 할 필수 구성 요소와 각 요소별 대표 오픈소스, 운영 난이도, 대체 벤더를 구체적으로

살펴봅니다.

로드밸런서(MetalLB) 구성

베어메탈 환경에서는 클라우드 제공 로드밸런서가 없으므로 MetalLB 같은 오픈소스 로드밸런서를 직접 구성해야 한다. MetalLB는 Layer 2/3 네트워크에서 IP 할당과 트래픽 분산을 담당하며, K8s 서비스에 외부 접근을 가능하게 한다. 운영자는 네트워크 아키텍처와 IP 풀 관리, 장애 대응에 대한 역량을 갖춰야 한다.

MetalLB를 도입할 때는 네트워크 토폴로지, BGP 또는 L2 모드 선택, IP 주소 풀 관리 등 세부적인 네트워크 설계가 필요합니다. 또한, 장애 발생 시 IP failover, 네트워크 충돌, 라우팅 정책 등 다양한 상황에 대한 대응 시나리오를 사전에 마련해야 하므로, 네트워크 엔지니어와 K8s 운영팀 간의 긴밀한 협업이 요구됩니다.

스토리지(Rook-Ceph·Longhorn) 구성

블록/오브젝트 스토리지 역시 클라우드의 managed 서비스 대신 Rook-Ceph(오브젝트·블록·파일 스토리지), Longhorn(경량 블록 스토리지) 등 오픈소스 솔루션을 자체 배포해야 한다. Rook-Ceph는 복제, 스냅샷, 확장성을 제공하며, Longhorn은 노드 장애 시 자동 복구와 볼륨 관리에 강점을 가진다. 운영자는 스토리지 클러스터의 HA, 백업, 복구 전략을 직접 설계해야 한다.

Rook-Ceph의 경우, Ceph 클러스터의 설계(모니터 노드, OSD, MDS 등), 데이터 복제 정책, 장애 복구, 성능 튜닝 등 고난도의 스토리지 운영 역량이 필요합니다. Longhorn은 상대적으로 경량이지만, 노드 장애 및 데이터 일관성 보장, 볼륨 스냅샷 및 복구 정책을 직접 관리해야 하므로, 스토리지에 대한 기본적인 이해와 운영 경험이 요구됩니다.

PKI(cert-manager)와 DNS(External DNS) 구성

cert-manager는 클러스터 내 PKI 인프라를 자동화하며, TLS 인증서 발급·갱신을 관리한다. External DNS는 DNS 레코드 자동화와 외부 DNS 연동을 담당한다. 이 두 구성 요소는 서비스 노출과 보안 인증에 필수적이며, 운영자는 DNS·PKI 정책을 직접 관리해야 한다.

cert-manager를 통해 ACME(예: Let's Encrypt) 기반 인증서 자동화, 사내 CA 연동, 인증서 만료 모니터링 등 다양한 보안 정책을 구현할 수 있습니다. External DNS는 Route53, Infoblox 등 외부 DNS 서비스와 연동하여, K8s 서비스의 외부 노출을 자동화합니다. 이 과정에서 DNS TTL, 레코드 동기화, 권한 관리 등 세부 정책을 직접 설계하고 운영해야 하므로, 보안 및 네트워크에 대한 폭넓은 이해가 필요합니다.

etcd HA 및 백업·복구

K8s Control Plane의 HA는 etcd 3노드 이상 구성, API server 다중화, 정기 백업·복구 전략이 필수다. 운영자는 etcd 데이터 손실·장애 대응 시나리오를 사전에 설계해야 하며, 복구 자동화 도구를 활용할 수 있어야 한다.

etcd 클러스터의 안정적 운영을 위해서는 노드 간 네트워크 레이턴시, 데이터 일관성, 스냅샷 백업 주기, 장애 발생 시 롤백 및 복구 프로세스 등 다양한 운영 정책을 수립해야 합니다. Velero 등 백업 솔루션을 활용해 K8s 리소스와 데이터를 주기적으로 백업하고, 실제 장애 상황에서 신속하게 복구할 수 있는 자동화 스크립트와 절차를 마련하는 것이 중요합니다.

필수 구성 요소 표

구성 요소	대표 오픈소스	운영 난이도	대체 벤더
로드밸런서	MetalLB	중	F5, NGINX
스토리지	Rook-Ceph, Longhorn	높음	NetApp, Dell EMC
PKI	cert-manager	중	Venafi, Let's Encrypt
DNS	External DNS	중	Infoblox, Route53
Control Plane HA	etcd	높음	-
백업·복구	Velero	중	Veeam, Commvault

이 표는 베어메탈 K8s에서 반드시 직접 구성해야 하는 6개 주요 요소와 대표 오픈소스, 운영 난이도, 대체 벤더를 보여준다. “베어메탈은 어렵다”라는 인식은 실제로 이 구성 요소들의 직접 운영 책임에서 비롯되며, 인력·교육 계획을 산정할 때 체크리스트로 활용할 수 있다.

결국 베어메탈 K8s의 성공적인 운영을 위해서는 각 인프라 요소별 오픈소스 솔루션의 특성과 운영 난이도를 충분히 이해하고, 조직 내에 필요한 기술 역량을 확보하는 것이 필수적입니다. 이를 통해 클라우드 네이티브 환경의 자율성과 확장성을 극대화할 수 있습니다.

6.2.2 Booking.com·CERN·Adidas 등 대규모 베어메탈 K8s 운영 사례

베어메탈 K8s가 실제로 대규모 환경에서 안정적으로 운영되고 있는 대표적인 사례로는 Booking.com, CERN, Adidas 등이 있습니다. 이들 조직은 VM 계층을 제거하고 물리 서버 위에 K8s를 직접 배포함으로써, 자원 집적률과 운영 효율성을 극대화하고, 장애 복구 및 배포 속도를 획기적으로 개선하였습니다. 본 절에서는 각 조직의 구체적인 운영 사례와 그 성공 요인을 살펴봅니다.

Booking.com 베어메탈 K8s 대규모 운영

Booking.com은 KubeCon EU 2022에서 수천 노드 베어메탈 환경에서 VM 없이 K8s 클러스터를 직접 운영한다고 발표했다. 이 조직은 VM 계층을 제거함으로써 자원 집적률을 극대화하고, 장애 복구·배포 속도를 크게 개선했다. 운영 조직 규모는 수십 명 수준이며, 배포 리드타임은 기존 VM 기반 대비 3~4배 단축되었다.

Booking.com의 사례는 베어메탈 K8s가 단순히 이론적 아키텍처가 아니라, 실제 대규모 서비스 환경에서 충분히 실현 가능성을 보여줍니다. 이들은 자체적으로 네트워크, 스토리지, 보안 등 모든 인프라 요소를 직접 설계하고 운영하며, 자동화된 배포 파이프라인과 모니터링 시스템을 구축하여 운영 효율성을 극대화하였습니다. 또한, 장애 발생 시 빠른 복구를 위해 멱등성 배포, 무중단 롤링 업데이트, 자동화된 장애 감지 및 복구 시스템을 도입하였으며, 이를 통해 서비스 가용성과 신뢰성을 크게 향상시켰습니다.

CERN의 베어메탈+OpenStack+K8s 혼합 운영

CERN은 LHC 데이터 처리와 대규모 과학 연산을 위해 베어메탈, OpenStack, K8s를 혼합 운영한다. VM과 컨테이너를 병행하지만, 핵심 워크로드는 베어메탈 위 K8s로 전환하고 있다. 이로써 데이터 처리 효율과 확장성을 확보했으며, 운영 조직은 수백 명 규모로 복잡한 인프라를 안정적으로 관리하고 있다.

CERN은 베어메탈 서버와 OpenStack 기반 VM, 그리고 K8s 클러스터를 유연하게 조합하여, 워크로드 특성에 따라 최적의 배치 모델을 선택합니다. 특히, 대규모 데이터 분석 및 과학 연산 워크로드는 베어메탈 K8s 위에서 실행하여, 최대의 성능과 확장성을 확보하고 있습니다. 이 과정에서 자체 개발한 인프라 자동화 도구와 모니터링 시스템을 활용하여, 수천~만 노드 규모의 인프라를 안정적으로 운영하고 있습니다.

Adidas의 K8s 전면 전환

Adidas는 2019년부터 4000+ 파드, 200+ 서비스 규모의 K8s 클러스터를 베어메탈 환경에서 운영하고 있다. 배포 속도는 VM 기반 대비 3~4배 개선되었으며, 장애 복구 시간(MTTR)은 50% 이상 단축되었다. 운영 조직은 DevOps·SRE 팀이 통합적으로 관리하며, 멱등성·자동확장·불변 인프라 원칙을 실현하고 있다.

Adidas의 사례는 베어메탈 K8s가 단순히 대규모 IT 조직에만 적합한 것이 아니라, 애플리케이션 개발과 운영의 민첩성을 추구하는 다양한 산업군에서도 충분히 적용 가능성을 보여줍니다. 이들은 DevOps와 SRE 팀이 협업하여, 인프라 코드화, 자동화된 배포, 실시간 모니터링, 장애

대응 프로세스를 구축하였으며, 이를 통해 서비스 품질과 운영 효율성을 동시에 달성하였습니다.

운영 사례 표

조직	노드 규모	배포 지표(속도/MTTR)	운영 조직 규모
Booking.com	수천 노드	배포 속도 3~4배 개선	수십 명
CERN	수천~만 노드	데이터 처리 효율 극대화	수백 명
Adidas	4000+ 파드, 200+ 서비스	MTTR 50% 단축, 배포 속도 3~4배 개선	DevOps·SRE 통합

이 표는 베어메탈 K8s가 “비현실”이라는 인식을 실제 글로벌 사례로 반박하며, 대규모 운영이 이미 현실적으로 가능함을 입증한다. “우리 규모 대비 3~10배 큰 조직이 이미 베어메탈로 돌리고 있다”는 근거는 리스크 완화 도구로 활용할 수 있다.

이러한 사례들은 베어메탈 K8s가 단순히 이론적 이상향이 아니라, 실제로 대규모 조직에서 성공적으로 운영되고 있음을 보여줍니다. 조직의 규모와 무관하게, 기술적 역량과 체계적인 운영 프로세스를 갖춘다면 베어메탈 K8s로의 전환이 충분히 가능하며, 이는 장기적으로 IT 비용 절감과 서비스 품질 향상이라는 두 마리 토끼를 모두 잡을 수 있는 전략적 선택임을 시사합니다.

7장: 라이선스·전력·인건비·Auto Scale 사이징 4대 TCO

축 비교

클라우드 네이티브 전환의 회계적 근거를 명확히 제시하는 본 장에서는 라이선스, 전력·자원 집적도, 인건비, 그리고 Auto Scale 기반 하드웨어 사이징 절감이라는 4대 TCO(Total Cost of Ownership) 축을 중심으로 IaaS와 PaaS(클라우드 네이티브) 환경의 구조적 비용 차이를 분석합니다. 최근 시장에서 라이선스 모델의 급격한 변화와 데이터센터 전력·공간·인력 비용의 압박, 그리고 자동 확장(Autoscale) 기술이 실제 하드웨어 구매·운영에 미치는 영향까지, 각 축별로 조직이 직접 자기 수치를 대입해 3년 전망을 산출할 수 있도록 근거와 예시적 시나리오를 제공합니다. 특히 VM의 정적 할당과 K8s의 동적 스케일이 하드웨어 사이징 자체를 어떻게 변화시키는지 회계적 관점에서 분석하는 것이 본 장의 핵심입니다.

7.1 라이선스 축 — 구(舊) 상용 IaaS 라이선스와 상용 오픈소스 기반 MSAP.ai의 회계적 비교

최근 클라우드 인프라 시장에서는 라이선스 정책의 변화가 조직의 IT 예산에 미치는 영향이 점점 커지고 있습니다. 특히 구(舊) 상용 IaaS 벤더와 외산 상용 쿠버네티스 디스트리뷰션이 영구 라이선스에서 고가의 번들 구독·코어당 과금 모델로 전환하면서, 라이선스 비용이 예측 불가능하게 증가하는 사례가 늘고 있습니다. 이에 반해, 국내 PaaS 제품인 MSAP.ai는 자체 영구 라이선스 체계를 유지하면서 웹서버·WAS·APM·Pod Cluster·AI·MSA Accelerator를 단일 라이선스에 번들 포함한 혁신적 가격 정책을 채택하여, 외산 상용 쿠버네티스의 고가 서브스크립션 대비 획기적인 비용 절감을 제공합니다. 본 절에서는 구 상용 IaaS·외산 상용 K8s의 구독 라이선스와 MSAP.ai의 영구 라이선스·번들 구조를 회계적 관점에서 비교하고, 실제 조직이 3년 예산 시나리오를 어떻게 구성할 수 있는지 구체적으로 설명합니다. 결론적으로, MSAP.ai는 구 라이선스 구조가 야기하는 반복 구독 비용 팽창을 제거하면서 핵심 플랫폼 구성 요소를 단일 번들로 제공하는 비용 효율적 국내 라이선스 구조라는 점을 회계적 근거와 함께 제시합니다.

7.1.1 구(舊) 상용 IaaS 구독 라이선스의 3년 예산 영향 — 구조적 요인만 압축 서술 구 라이선스의 구독 전환과 예산 영향

구 상용 IaaS 진영에서는 2024년을 기점으로 라이선스 모델이 영구 라이선스에서 번들 구독 및 코어당 최소 단위 과금으로 전환되었습니다. 이 구조적 변화는 기존의 “서버 단위 라이선스”에서 “코어 단위 구독”으로 비용 산정 방식이 바뀌어, 대규모 환경에서는 2~10배의 인상률이 발생할 수 있습니다. 예를 들어, VMware의 라이선스 정책 변경은 AT&T 등 대형 통신사의 소송으로 이어졌으며, 실제 예산 산정 시 3년 단위로 갱신 주기와 과금 단위가 조직의 TCO에 직접적인 영향을 미칩니다. 이러한 변화는 단순히 가격 인상 문제가 아니라, 구 라이선스 계약의 구조적 요인이 예산 계획에 반영되어야 함을 의미합니다.

예시적(illustrative) 시나리오

예산 시뮬레이션은 반드시 “예시적(illustrative)”임을 명시해야 하며, 실제 수치는 워크로드, 라이선스 계약, 인건비 등 조직별 변수에 따라 크게 달라집니다. 예를 들어, 현행 유지(구 상용 IaaS 구독 지속), 부분 전환(일부 워크로드를 MSAP.ai로 이전), 전면 전환(전체 워크로드를 MSAP.ai로

이전) 세 가지 시나리오를 비교하면, 라이선스·국내 지원 계약·마이그레이션 비용이 각각 다르게 산출됩니다. 현행 유지의 경우 구독 비용이 지속적으로 누적되고, 부분 전환은 마이그레이션 비용과 국내 지원 계약이 혼합되며, 전면 전환은 **상용 오픈소스 기반 국내 지원 계약만 남아 라이선스 비용이 급감하는 구조를** 갖습니다.

구조적 요인과 회계적 판단

TCO 논쟁은 감정이 아닌 계약·과금 단위·갱신 주기의 산술로 내려져야 합니다. 조직은 반드시 2024년 라이선스 지형 변화의 구조적 요인을 예산 계획에 반영하고, 구체 수치가 필요하면 1.1.2의 라이선스 지형 변화 표를 참조해야 합니다. 구 라이선스 측은 단순히 비용 항목이 아니라, 운영 모델 전환의 회계적 근거로 활용되어야 하며, 구독 전환·번들화·코어당 과금의 구조적 특성을 명확히 이해해야 합니다. 특히 구 라이선스는 환율 변동과 해외 본사 정책 변경에 예산이 직접 노출되어 있어, **원화 기반 국내 계약으로의 전환 효과는 예산 안정성 측면에서도 유의미합니다.**

구 상용 IaaS 구독 모델의 실제 사례와 추가적 고려사항

실제 국내외 대기업에서는 구 상용 IaaS 구독 모델 전환에 따라 예산 계획을 재조정하는 사례가 늘고 있습니다. 기존에 영구 라이선스와 연간 유지보수만으로 운영하던 환경이 구독 기반으로 전환되면, 3년 단위로 전체 라이선스 비용이 재산정되고, 갱신 시점마다 예산이 대폭 증가할 수 있습니다. 또한, 코어 수 증가나 신규 워크로드 도입 시 추가 구독이 필수화되어, 예산 예측이 더욱 어려워집니다. 이와 달리, 일부 워크로드를 MSAP.ai로 이전하는 부분 전환 시나리오에서는 마이그레이션 비용이 일시적으로 증가하지만, 장기적으로는 라이선스 비용이 점진적으로 감소하고 **원화 기반 국내 계약으로 예산 예측 가능성이 높아지는 구조를** 갖게 됩니다. 전면 전환의 경우에는 초기 마이그레이션과 MSAP.ai 국내 지원 계약 비용만 남아, 라이선스 비용 부담이 크게 줄어드는 효과를 볼 수 있습니다. 조직은 이러한 구조적 변화를 미리 예측하고, 각 시나리오별로 3년 단위 예산을 시뮬레이션하여 최적의 전략을 수립해야 합니다.

7.1.2 MSAP.ai — 영구 라이선스·번들 플랫폼 기반의 비용 효율적 국내 라이선스 구조

외산 상용 K8s의 구독 라이선스와 MSAP.ai 영구 라이선스의 구조적 차이

외산 상용 쿠버네티스 디스트리뷰션은 “라이선스=연간 구독=갱신 시점마다 재지불”이라는 반복 지불 구조를 강제합니다. 즉, 계약 기간이 종료되면 동일 기능을 계속 사용하기 위해 매년

혹은 3년 단위로 고가의 구독료를 재지불해야 하며, 갱신 시점의 본사 정책·환율 변동에 예산이 직접 노출됩니다. 반면 MSAP.ai는 국내 PaaS 제품으로서 자체 영구 라이선스 체계를 채택하여, 한 번의 라이선스 구매로 기간 제한 없는 사용권을 확보할 수 있습니다. 이 구조는 외산 상용 K8s의 “사용권 = 매년 지출되는 반복 비용” 공식을 해체하고, “사용권 = 1회 지불 자산화”로 재정적함으로써 조직의 3년·5년 다년도 예산 예측 가능성을 근본적으로 회복시킵니다.

단일 라이선스 번들에 포함되는 혁신적 구성 — 웹서버·WAS·APM·Pod Cluster·AI·MSA Accelerator

MSAP.ai의 혁신적 가격 정책의 핵심은 단일 라이선스 안에 기업 애플리케이션 운영에 필요한 핵심 플랫폼 구성 요소가 모두 번들로 포함된다는 점입니다. 구체적으로 웹서버, WAS(Web Application Server), APM(Application Performance Management), Pod Cluster(쿠버네티스 기반 컨테이너 오케스트레이션), AI 스택, MSA(Microservice Architecture) Accelerator가 단일 MSAP.ai 라이선스에 통합 제공됩니다. 외산 환경에서는 이 구성 요소들을 개별 벤더·개별 구독 계약으로 각각 도입해야 하며, 각 품목마다 코어당 구독료·지원 계약·갱신 인상률이 누적되어 총비용이 기하급수적으로 증가합니다. MSAP.ai는 이를 하나의 영구 라이선스 번들로 제공하여, 조직이 중복 구독 비용 없이 통합 플랫폼을 운영할 수 있도록 합니다.

MSAP.ai 국내 라이선스 구조의 비용 효율성

MSAP.ai의 국내 라이선스 구조가 비용 효율적인 이유는 네 가지로 요약됩니다. 첫째, 영구 라이선스 체계 — 한 번의 구매로 기간 제한 없이 사용 가능하며, 외산 상용 K8s의 매년 반복 구독료 대비 3년·5년 누적 TCO가 획기적으로 낮아집니다. 둘째, 번들 통합 가격 — 웹서버·WAS·APM·Pod Cluster·AI·MSA Accelerator가 단일 라이선스에 포함되어, 개별 품목을 외산으로 도입할 때 발생하는 중복 구독료가 원천 제거됩니다. 셋째, 원화 기반 계약 — 달러 환율 변동과 해외 본사 가격 인상 정책에 예산이 노출되지 않아 다년도 예산 예측이 정확해집니다. 넷째, 국내 지원 체계 포함 — 한글 문서·국내 SLA·온사이트 대응이 라이선스 구조 내에서 제공되어, 외산 벤더에서 별도 청구되던 프리미엄 지원 항목의 추가 비용이 발생하지 않습니다.

외산 상용 K8s 대비 MSAP.ai의 예산 안정성

외산 상용 쿠버네티스 라이선스는 3년 갱신 주기마다 코어당 단가 인상, 번들 강제 편입, 본사 정책 변경에 예산이 직접 노출됩니다. MSAP.ai는 영구 라이선스 체계이므로 동일 구성 요소의 지속 사용에 대해 갱신 인상 리스크가 구조적으로 존재하지 않으며, 유지보수 계약만 연 단위 원화로 관리되어 환율·해외 본사 정책과 독립적으로 운영됩니다. 특히 공공·금융과 같이 다년도 예산

편성의 정확성이 요구되는 조직에서, **사용권은 1회 영구 자산화·유지보수만 연 계약**이라는 국내 라이선스 구조는 외산 상용 K8s 대비 3년·5년 TCO의 상방 리스크를 획기적으로 축소합니다.

외산 상용 K8s vs MSAP.ai 라이선스 구조 비교 표

비교 항목	외산 상용 쿠버네티스	MSAP.ai (국내 PaaS)
라이선스 모델	고가 연간/3년 구독	자체 영구 라이선스
반복 지불 구조	매년/갱신 시점 재지불	1회 구매 후 영구 사용
번들 포함 범위	K8s 중심 (웹서버·WAS·APM·AI 등 별도 구독)	웹서버·WAS·APM·Pod Cluster·AI·MSA Accelerator 단일 번들
과금 통화	USD 기반 (환율 노출)	원화 기반 (환율 독립)
국내 SLA·한글화	제한적·추가 비용	기본 포함
3년·5년 누적 TCO	높음 (구독 반복)	낮음 (영구 라이선스)
갱신 인상 리스크	높음 (본사 정책·환율)	없음 (사용권 영구)

실제 적용 시 예산·운영 효과

국내 금융권·공공기관에서는 외산 상용 쿠버네티스의 3년 갱신 시점에 구독 단가가 2~10배까지 상승한 사례가 보고되고 있으며, 이러한 예산 충격을 완화하기 위해 MSAP.ai로 전면 또는 부분 전환하는 사례가 증가하고 있습니다. 외산 구독에서 MSAP.ai 영구 라이선스로 전환할 경우, **매년 반복되던 고가 구독료가 1회 영구 라이선스 비용으로 치환**되고, 개별 품목(웹서버·WAS·APM·AI 등)으로 분리 구매해야 했던 플랫폼 구성 요소가 단일 번들로 통합되기 때문에, 총소유비용(TCO)의 상당 부분을 차지하던 해외 구독료 누적이 사라지는 회계적 효과가 발생합니다. 또한 한글 문서·국내 SLA·온사이트 대응이 기본 포함되어, 외산 벤더에서 별도 청구되던 프리미엄 지원 항목이 라이선스 구조 내에 흡수되는 구조적 이점이 있습니다. 결과적으로 MSAP.ai는 라이선스 축에서 국내 조직이 선택할 수 있는 가장 비용 효율적인 국내 라이선스 구조로 평가되며, 외산 상용 쿠버네티스 구독 모델의 구조적 예산 리스크를 근본적으로 해소하는 **영구 라이선스·번들 통합형 회계적 대안**이 됩니다.

7.2 전력·자원 집적도·인건비 축

전력·자원 집적도·인건비 축은 라이선스 외에 실제 데이터센터 운영 비용에 직접적으로 영향을 미치는 항목입니다. 컨테이너 기반 K8s 환경은 VM 대비 3~10배의 집적률을 제공하며, 동일 하

드웨어에서 더 많은 워크로드를 수용할 수 있습니다. 이로 인해 전력·공간·냉각 비용이 감소하며, 운영팀(FTE)·도구·정책의 이중화가 단일 K8s 표면으로 수렴될 때 인건비 절감 효과가 발생합니다. 본 절에서는 집적률과 전력, FTE 수렴, 그리고 예시적 TCO 시뮬레이션의 변동성 요인을 구체적으로 설명합니다. 특히, 실제 데이터센터 운영에서 집적률 향상과 인건비 절감이 어떻게 실질적인 비용 절감으로 이어지는지, 그리고 각 조직이 이를 어떻게 측정하고 적용할 수 있는지에 대한 실무적 가이드라인을 제공합니다.

7.2.1 동일 H/W에서 VM 대비 컨테이너 3~10배 집적률의 자원·전력 효과

컨테이너 집적률의 구조적 우위

Google SRE 책과 Borg 사례에 따르면, 동일 하드웨어에서 VM 대비 컨테이너가 3~10배 더 많은 워크로드를 수용할 수 있습니다. 이는 VM이 각 인스턴스마다 게스트 OS를 포함해 자원 오버헤드가 발생하는 반면, 컨테이너는 호스트 커널을 공유하며 경량화된 프로세스 수준 격리로 집적률이 극대화되기 때문입니다. IBM 리포트에 따르면, 컨테이너는 베어메탈 대비 CPU·메모리 오버헤드가 0% 수준이며, VM(KVM)은 25% CPU, I/O는 최대 40% 오버헤드가 발생합니다.

전력·공간·냉각 비용 절감

집적률이 높아지면 동일한 데이터센터 공간에서 더 많은 워크로드를 처리할 수 있으며, 서버 대수를 줄이면 전력·공간·냉각 비용이 직접적으로 감소합니다. ESG(환경·사회·지배구조) 관점에서도 데이터센터 전력 예산 절감은 중요한 회계적 근거가 됩니다. 특히 공간 부족 조직에서는 집적률 향상이 공간 비용 절감으로 직결되므로, K8s 환경으로의 전환이 단순한 기술적 선택이 아닌 회계적·환경적 전략임을 강조해야 합니다.

집적률·CPU 오버헤드·I/O 오버헤드 비교 표

항목	VM 환경	컨테이너 환경
집적률	1x	3~10x
CPU 오버헤드	2~5%	~0%
I/O 오버헤드	최대 40%	~0%

실제 데이터센터 적용 사례와 추가 분석

실제 글로벌 기업의 데이터센터에서는 컨테이너 기반 집적률 향상으로 서버 대수를 30~70%

까지 줄인 사례가 보고되고 있습니다. 예를 들어, Adidas의 CNCF 사례에서는 컨테이너 도입 후 동일한 하드웨어에서 처리 가능한 워크로드가 5배 이상 증가하였고, 이에 따라 전력 및 냉각 비용이 연간 수억 원 단위로 절감되었습니다. 또한, 서버 대수 감소는 랙 공간, 네트워크 스위치, UPS 등 부대 인프라 비용까지 줄여주어, 전체 TCO 절감 효과가 더욱 극대화됩니다. 이러한 효과는 단순히 기술적 효율성에 그치지 않고, 조직의 ESG 목표 달성 및 친환경 경영에도 크게 기여할 수 있습니다. 조직은 실제 집적률, 전력 소비량, 공간 활용률 등을 주기적으로 모니터링하여, 컨테이너 기반 전환의 실질적 효과를 정량적으로 평가하고, 추가적인 최적화 방안을 도출해야 합니다.

7.2.2 VM 운영팀 + K8s 운영팀 이중화에서 단일 K8s 표면으로의 FTE 수렴

이중 운영의 구조적 비용

VM 환경(하이퍼바이저·OS 중심)과 K8s 환경(컨테이너 중심)이 공존할 경우, 운영팀(FTE), 도구, 정책이 2중으로 유지되어 중복 비용이 발생합니다. VM 운영팀은 하이퍼바이저 관리, OS 패치, 라이선스 관리에 집중하며, K8s 운영팀은 클러스터 관리, 컨테이너 배포, 네이티브 관측성 도구(Prometheus 등)를 사용합니다. 이중 운영은 FTE(Full Time Equivalent) 인력, 모니터링·보안 도구, 정책 관리의 3중 중복을 낳아, 3년 누적 인건비가 라이선스 비용을 상회할 수 있습니다.

단일 K8s 표면 수렴의 인건비 절감

운영 모델을 단일 K8s 표면으로 수렴하면, VM 운영팀과 K8s 운영팀의 FTE를 재배치할 수 있으며, 도구·정책도 단일화됩니다. 예를 들어, Prometheus·Grafana·OpenTelemetry 등 K8s 네이티브 관측성 도구를 통합하면 모니터링·알람·보고 체계가 일원화되고, 보안 정책도 OPA·Falco 등으로 단일화됩니다. 이는 인건비 절감뿐 아니라 장애 대응·운영 효율성·감사 대응의 품질을 동시에 개선합니다.

이중 운영 vs 단일 표면 FTE·도구·정책 비교 표

항목	이중 운영	단일 표면(K8s)
FTE	중복 유지	재배치/축소
도구	2중 관리	통합 관리
정책	2중 배포	단일 배포

실제 조직 적용 사례 및 인건비 분석

실제 대기업이나 금융권에서는 VM과 K8s 운영팀이 별도로 존재하는 경우가 많아, 인력과 도구의 중복이 불가피하게 발생합니다. 예를 들어, VM 운영팀은 하이퍼바이저 및 OS 패치, 라이선스 관리에 집중하고, K8s 운영팀은 클러스터 관리와 컨테이너 배포, 네이티브 관측성 도구(Prometheus, Grafana 등) 운영을 담당합니다. 이중 운영 체계에서는 인력 배치뿐 아니라, 모니터링·보안 도구, 정책 관리 등에서도 중복 투자가 이루어집니다. 이러한 중복은 3년 누적 인건비가 라이선스 비용을 상회하는 결과로 이어질 수 있습니다. 반면, 단일 K8s 표면으로 운영을 수렴할 경우, 기존 VM 운영팀 인력을 K8s 운영팀으로 재배치하거나, 일부 인력을 감축할 수 있으며, 도구와 정책도 통합 관리가 가능해집니다. 실제로 한 글로벌 금융사는 K8s 기반 단일 운영 모델 도입 후, 운영 인력을 30% 이상 감축하고, 모니터링 및 보안 도구 라이선스 비용도 절감한 사례를 보고하였습니다. 이러한 변화는 단순한 인건비 절감에 그치지 않고, 장애 대응 속도, 감사 대응 품질, 운영 효율성 등 다양한 측면에서 조직의 경쟁력을 높여줍니다. 조직은 현재 운영팀 구조와 도구, 정책의 중복 현황을 면밀히 분석하여, 단일 K8s 표면으로의 전환이 가져올 수 있는 인건비 및 운영 효율성 개선 효과를 정량적으로 산출해야 합니다.

7.2.3 “예시적(illustrative)” TCO 시뮬레이션 제시 원칙과 변동성 요인

TCO 시뮬레이션의 변동성

TCO(총소유비용) 시뮬레이션은 워크로드 유형, 라이선스 계약, 인건비 시장 단가, 스케일 패턴 등 다양한 변수에 따라 크게 변동합니다. 본 백서의 시뮬레이션은 “예시적(illustrative)”임을 반드시 명시해야 하며, 조직별로 실제 수치를 대입해 재계산하는 워크시트 형태로 제공해야 합니다. 백서 수치를 그대로 인용하는 것은 위험하며, 각 조직은 자기 환경에 맞는 변수를 적용해 TCO를 산출해야 합니다.

변수별 영향도와 조사 방법

TCO 시뮬레이션에서 주요 변수는 워크로드(서비스 유형·규모), 라이선스 계약(구독·지원 범위), 인건비(운영팀 FTE·시장 단가), Auto Scale 사이징(동적 스케일 절감률)입니다. 각 변수는 TCO에 직접적인 영향을 미치며, 조사 방법은 모니터링 데이터, 계약서, 조직도, 클러스터 스케일링 로그 등 다양한 자료를 활용해야 합니다.

TCO 시뮬레이션 변수·영향도·조사 방법 표

변수	영향도	조사 방법
워크로드	자원·운영비	서비스별 모니터링
라이선스 계약	연간 비용	계약서·구독 내역
인건비	운영비	조직도·FTE 산정
Auto Scale 사이징	하드웨어 절감	클러스터 로그 분석

실제 시뮬레이션 적용 방법과 주의사항

조직이 TCO 시뮬레이션을 수행할 때는 반드시 자기 환경에 맞는 실제 데이터를 기반으로 변수를 산정해야 합니다. 예를 들어, 워크로드 유형별로 자원 사용량과 성장률을 모니터링하고, 라이선스 계약서는 실제 구독 내역과 갱신 주기를 확인해야 합니다. 인건비의 경우, 운영팀의 FTE 현황과 시장 단가를 반영하여 산출하고, Auto Scale 사이징 효과는 클러스터 스케일링 로그와 실제 리소스 활용률 데이터를 분석해 적용해야 합니다. 또한, 시뮬레이션 결과는 예시적 수치임을 명확히 표기하고, 조직별로 워크시트 형태로 재계산할 수 있도록 가이드라인을 제공해야 합니다. 실제로 한 IT 서비스 기업은 TCO 시뮬레이션을 통해 VM 기반 인프라에서 K8s 기반 인프라로 전환 시, 3년간 총소유비용이 40% 이상 절감되는 효과를 확인하였으나, 이는 조직의 워크로드 특성, 인건비 구조, 라이선스 계약 조건 등에 따라 달라질 수 있음을 강조하였습니다. 따라서, 모든 조직은 백서의 예시적 수치를 참고하되, 반드시 자체 데이터를 활용해 TCO를 산출해야 하며, 주요 변수의 변동성에 대한 민감도 분석도 병행해야 합니다.

7.3 Auto Scale 기반 하드웨어 사이징 절감 축 — VM 정적 할당 vs K8s 동적 스케일

Auto Scale 기반 하드웨어 사이징 절감 축은 클라우드 인프라의 실제 하드웨어 구매 및 운영 비용에 가장 직접적으로 영향을 미치는 요소 중 하나입니다. VM 환경에서는 피크 부하 기준으로 자원을 정적으로 할당해야 하므로, 평균 자원 사용률이 낮아지고 유휴 자원이 많아지는 구조적 한계가 존재합니다. 반면, Kubernetes 환경에서는 HPA, KEDA, Cluster Autoscaler 등 다양한 오토 스케일러를 활용해 파드와 노드의 수를 동적으로 조정할 수 있으므로, 평균 부하에 맞춘 하드웨어 사이징이 가능해집니다. 본 절에서는 VM 정적 할당의 구조적 유휴 원인, K8s 동적 스케일 메커니즘, 그리고 동일 워크로드 기준 하드웨어 사이징을 30~50%까지 절감할 수 있는 산출 근거와 실제

적용 시 고려해야 할 전제 및 한계를 구체적으로 설명합니다. 이를 통해 조직이 하드웨어 투자와 운영 비용을 어떻게 최적화할 수 있는지, 실질적인 회계적 효과를 이해할 수 있도록 안내합니다.

7.3.1 VM 정적 할당이 강제하는 구조적 유휴 — 자원 사용률이 낮을 수밖에 없는 이유

정적 할당의 구조적 문제

VM 환경에서는 피크 부하 기준으로 vCPU·RAM을 사전 고정 할당하는 것이 기본 설계입니다. 이는 부팅 시간(분 단위)으로 인해 사전 확보가 필요하며, 애플리케이션 단위 분리로 리소스 파편화가 발생합니다. 게스트 OS·하이퍼바이저 오버헤드가 겹쳐 실제 평균 자원 사용률이 낮게 고착됩니다. 업계에서 통상 보고되는 VM 평균 CPU 활용률은 15~25% 범위로, 대부분의 자원이 유휴 상태로 남아 있습니다.

자원 사용률 저하의 근본 원인

이러한 구조적 유휴는 단순히 운영 미숙이 아니라 설계 자체의 결과입니다. VM은 피크 부하에 맞춰 자원을 할당하므로, 비피크 시점에는 자원이 유휴화됩니다. 부팅 시간, 게스트 OS 오버헤드, 리소스 파편화 등 복합 요인이 동적 스케일을 어렵게 만들며, 실제 자원 활용률은 평균적으로 15~25% 수준에 머무릅니다. 이는 하드웨어 구매·운영 비용의 비효율을 구조적으로 야기합니다.

VM 정적 할당 유발 요인·평균 활용률 영향도 표

요인	평균 활용률 영향도
피크 부하 기준 할당	낮음
부팅 시간(분 단위)	낮음
리소스 파편화	낮음
게스트 OS·하이퍼바이저 오버헤드	낮음

업계 보고 VM 평균 활용률 범위(예시적): 15~25%

실제 사례와 추가 분석

실제 글로벌 IT 기업에서는 VM 기반 인프라의 평균 CPU 활용률이 20% 미만에 머무르는 경우가 많으며, 이는 하드웨어 구매 비용의 80% 이상이 유휴 자원으로 남는다는 의미입니다. 예를 들어, 한 대형 금융사는 피크 부하를 기준으로 서버를 도입한 결과, 연중 대부분의 시간 동안 자원

활용률이 15% 수준에 불과하였고, 이에 따라 전력 및 공간 비용도 비효율적으로 소모되었습니다. 이러한 구조적 유휴는 단순히 서버 수를 줄인다고 해결되지 않으며, VM 아키텍처의 설계적 한계로 인해 동적 확장이 어렵다는 점을 반드시 고려해야 합니다. 조직은 실제 자원 사용률 데이터를 주기적으로 모니터링하고, 정적 할당 구조가 하드웨어 및 운영 비용에 미치는 영향을 정량적으로 분석해야 합니다.

7.3.2 HPA·KEDA·Cluster Autoscaler의 0→N→0 동적 스케일 메커니즘

K8s 동적 스케일의 구조적 장점

Kubernetes는 Horizontal Pod Autoscaler(HPA), KEDA(이벤트 기반 스케일러), Cluster Autoscaler 등 3가지 주요 스케일러를 제공합니다. HPA는 CPU·메모리·커스텀 지표 기반으로 파드를 자동 확장/축소하며, KEDA는 메시지 큐·HTTP 요청·DB 지표 등 이벤트 기반으로 파드를 0→N→0까지 스케일합니다. Cluster Autoscaler는 클라우드·온프레미스 노드 풀에서 노드 자체를 확장·축소합니다. 컨테이너 시작 시간은 초 단위로, VM 부팅(분 단위) 대비 현격히 짧아 동적 스케일이 현실적으로 작동합니다.

Auto Scale의 회계적 효과

Auto Scale은 이상적 기능이 아니라 K8s의 기본 구성 요소입니다. 파드·노드 단위로 동적 스케일이 가능하며, 실제 하드웨어 총량은 평균 부하에 마진만 엮는 방식으로 산정할 수 있습니다. 이는 피크 기준 정적 할당이 강제되는 VM 환경과 달리, 평균 기준 사이징이 가능하여 하드웨어 구매·운영 비용을 구조적으로 절감할 수 있습니다.

3개 스케일러(HPA·KEDA·Cluster Autoscaler) 비교 표

스케일러	스케일 단위	트리거 지표	확장 방향
HPA	Pod	CPU/메모리/지표	수평 확장/축소
KEDA	Pod	이벤트(큐/HTTP)	0→N→0
Cluster Autoscaler	Node	노드 부하/Pod	노드 확장/축소

실제 적용 사례 및 기술적 세부사항

실제 대규모 이커머스 기업에서는 K8s 오토스케일러를 도입하여, 트래픽 급증 시 파드와 노드가 자동으로 확장되고, 트래픽 감소 시 즉시 축소되는 구조를 구현하였습니다. 예를 들어, HPA는 웹

애플리케이션의 CPU 사용률이 70%를 초과하면 파드를 자동으로 추가하고, 40% 이하로 떨어지면 파드를 줄입니다. KEDA는 메시지 큐의 대기 메시지 수가 일정 임계치를 넘으면 파드를 0에서 N개로 확장하고, 큐가 비면 다시 0으로 축소합니다. Cluster Autoscaler는 전체 클러스터의 파드 배치가 불가능해질 때 자동으로 노드를 추가하고, 유휴 노드가 일정 시간 이상 지속되면 노드를 축소합니다. 이러한 동적 스케일 메커니즘은 VM 환경의 정적 할당 구조와 달리, 실제 트래픽과 워크로드 변화에 따라 자원을 실시간으로 최적화할 수 있게 해줍니다. 조직은 오토스케일러의 설정값(임계치, 쿨다운 타임 등)을 실제 워크로드 패턴에 맞게 조정하여, 하드웨어 자원 활용률을 극대화하고, 불필요한 비용 지출을 최소화해야 합니다. 또한, 오토스케일링 로그와 모니터링 데이터를 분석하여, 스케일링 정책의 효과성과 한계를 지속적으로 점검하는 것이 중요합니다.

7.3.3 동일 워크로드 기준 하드웨어 사이징 30~50% 절감(예시적) — 산출 전제와 한계

사이징 절감의 산출 원리

VM 정적 할당 시 “피크 부하 × 안전 마진” 기준으로 하드웨어 총량을 산정하지만, K8s 동적 스케일에서는 “평균 부하 × 버퍼 + 피크 대응 파드 오토스케일” 기준으로 재산정합니다. 동일 워크로드 기준으로 하드웨어 사이징을 30~50% 수준까지 절감할 수 있다는 “예시적(illustrative)” 수치를 제시합니다. 실제로 K8s 환경에서는 평균 활용률이 50~70% 범위까지 상승하며, 컨테이너 집적률은 3~10배까지 보고됩니다.

산출 전제와 한계

산출 전제는 워크로드 스파이크 패턴, 스테이트풀/스테이트리스 비율, 멀티 테넌시 수준 등입니다. 스테이트풀 워크로드(데이터베이스 등)는 절감률이 축소되며, 스테이트리스 워크로드(웹·API 등)는 절감률이 확대됩니다. 레거시 워크로드, 규제 요건, 하드 멀티테넌시 등은 절감률에 영향을 미칠 수 있으므로, 조직은 반드시 자기 환경의 피크/평균 비율을 산정해 워크시트로 재계산해야 합니다.

동일 워크로드 × 사이징 산정 방식(VM 정적 vs K8s 동적) × 총 서버 대수 비율 표(예시적)

기준	VM 정적 할당	K8s 동적 스케일
총 서버 대수(예시적)	100	50~70

전제·한계 별도 주석

- 전제: 워크로드 스파이크 패턴, 스테이트풀/스테이트리스 비율, 멀티 테넌시 수준
- 한계: 레거시 스테이트풀 워크로드는 절감률 축소, 규제·보안 요건은 추가 하드웨어 필요

실제 적용 사례와 비교 분석

실제 글로벌 IT 서비스 기업에서는 K8s 기반 오토스케일링 도입 후, 동일 워크로드 기준 하드웨어 서버 대수를 40% 이상 절감한 사례가 다수 보고되고 있습니다. 예를 들어, 한 이커머스 기업은 VM 환경에서 100대의 서버로 운영하던 웹 서비스 워크로드를 K8s 환경으로 이전한 결과, 평균 부하 기준으로 60대의 서버만으로 동일한 서비스 품질을 유지할 수 있었습니다. 이는 평균 자원 활용률이 20%에서 65%로 상승한 데 기인하며, 오토스케일러의 동적 확장/축소 기능이 실시간 트래픽 변화에 따라 자원 배분을 최적화한 결과입니다. 다만, 데이터베이스와 같은 스테이트풀 워크로드는 오토스케일링의 한계로 인해 절감 효과가 상대적으로 작을 수 있으며, 규제나 보안 요건이 엄격한 환경에서는 추가 하드웨어가 필요할 수 있습니다. 조직은 반드시 자기 환경의 워크로드 특성, 피크/평균 부하 비율, 스테이트풀/스테이트리스 구성 비율 등을 면밀히 분석하여, 하드웨어 사이징 절감 효과를 정량적으로 산출해야 하며, 예시적 수치에만 의존하지 않고 자체 워크시트를 통해 시나리오별로 재계산하는 것이 바람직합니다.

8장: AI 스택(모델·Vector DB·RAG·API·MCP·Agent)의 K8s 정합성과 VM 비경제성

8.1 AI 스택 구성 요소의 K8s 기본 기능 의존성

AI 스택은 단일 모델이나 데이터베이스로만 구성되는 것이 아니라, 다양한 컴포넌트들이 유기적으로 결합되어 전체적인 AI 서비스의 품질과 효율성을 좌우합니다. 이러한 구성 요소들은 Kubernetes의 핵심 기능에 의존하여 배포, 운영, 확장, 복구 등 다양한 측면에서 최적화된 환경을 제공합니다. Kubernetes의 선언형 오브젝트, 자동 확장, GPU 리소스 관리, 서비스 메시 등은 AI 스택의 각 요소가 안정적으로 동작하고, 운영 효율성을 극대화하는 데 필수적인 역할을 합니다. 본 절에서는

AI 스택의 6가지 주요 요소가 Kubernetes의 어떤 기능들과 연계되어 있는지, 그리고 실무적으로 어떻게 적용되는지 구체적인 기술적 근거와 사례를 들어 설명하겠습니다.

8.1.1 KServe·vLLM·Triton을 통한 선언형 모델 서빙과 자동 확장

선언형 모델 서빙 구조

KServe, vLLM, NVIDIA Triton Inference Server는 모두 Kubernetes에서 선언형 오브젝트(CRD, Custom Resource Definition)로 배포할 수 있는 대표적인 모델 서빙 솔루션입니다. KServe는 모델 배포를 위한 InferenceService CRD를 제공하며, 사용자는 YAML 파일로 모델의 경로, 리소스 요구사항, 배포 옵션 등을 선언합니다. vLLM은 LLM 서빙에 특화된 경량화 엔진으로, GPU 리소스 관리와 대규모 동시 요청 처리에 강점을 지니며, Kubernetes Deployment 또는 Helm Chart로 배포가 가능합니다. Triton은 NVIDIA가 공식 지원하는 멀티 프레임워크 모델 서빙 서버로, GPU Operator와 결합하여 CUDA, MIG, DCGM 등 GPU 관련 기능을 자동으로 통합합니다. 선언형 배포는 GitOps 파이프라인과 결합하여 모델 버전 관리, 롤백, Canary 배포 등 엔터프라이즈 운영에 필요한 기능을 자연스럽게 제공합니다.

이러한 선언형 모델 서빙 구조의 장점은 운영 자동화와 신뢰성에 있습니다. 예를 들어, 모델의 새로운 버전을 배포할 때 YAML 파일의 버전만 변경하면 자동으로 롤링 업데이트가 이루어지고, 문제가 발생하면 이전 버전으로 손쉽게 롤백할 수 있습니다. 또한, Canary 배포를 통해 일부 트래픽만 새로운 모델에 전달하여 안정성을 검증한 후 전체로 확장할 수 있습니다. 이 모든 과정이 코드 기반으로 관리되기 때문에, 운영자의 실수나 환경 차이로 인한 장애를 최소화할 수 있습니다.

자동 확장(HPA/KEDA) 적용

모델 서빙 워크로드는 장시간 유휴 상태와 짧은 피크 트래픽이 반복되는 특성이 있습니다. Kubernetes의 Horizontal Pod Autoscaler(HPA)는 CPU, 메모리, 커스텀 지표(QPS, 큐 깊이 등)를 기준으로 파드 수를 자동 확장/축소합니다. KEDA(Kubernetes Event-driven Autoscaler)는 메시지 큐, HTTP 요청, DB 지표 등 다양한 이벤트를 트리거로 삼아 파드를 0→N→0까지 동적으로 조정할 수 있습니다. 예를 들어, vLLM이나 Triton을 KEDA와 결합하면 피크 QPS에 맞춰 파드가 자동 확장되고, 트래픽이 없을 때는 리소스가 0으로 줄어들어 GPU 비용을 최소화할 수 있습니다. 이 구조는 VM 환경에서는 구현이 불가능하거나 복잡한 반면, K8s에서는 표준 기능으로 제공됩니다.

실제 현장에서는 HPA와 KEDA를 함께 활용하여, 예측 불가능한 트래픽 변화에도 신속하게 대응할 수 있습니다. 예를 들어, 대규모 이벤트나 캠페인 기간에는 트래픽이 급증할 수 있는데, 이때 HPA가 CPU/메모리 사용량을 기준으로 파드를 확장하고, KEDA는 메시지 큐의 대기열 길이나 외부 이벤트를 감지하여 추가 확장을 트리거합니다. 이러한 조합은 리소스 낭비를 최소화하면서도, 서비스 품질을 안정적으로 유지하는 데 매우 효과적입니다.

GPU Operator 연동 및 지원

KServe, vLLM, Triton 모두 NVIDIA GPU Operator와 연동하여 GPU 리소스 관리, 드라이버 설치, MIG 분할, DCGM 모니터링 등을 자동화합니다. GPU Operator는 DaemonSet으로 배포되어 노드별로 GPU 관련 소프트웨어를 자동 설치하며, 파드 단위로 GPU 할당, MIG 분할, Time-slicing을 지원합니다. 이를 통해 AI 모델 서빙 워크로드가 GPU를 효율적으로 공유하고, 단일 노드에서 여러 모델을 동시에 서빙할 수 있습니다. VM 환경에서는 GPU 패스스루로 인해 한 VM이 GPU를 독점하지만, K8s에서는 파드 단위 GPU 분할이 가능해 인프라 단가를 크게 낮출 수 있습니다.

GPU Operator의 도입은 GPU 인프라의 표준화와 운영 효율성 향상에 큰 기여를 합니다. 예를 들어, 새로운 GPU 노드가 클러스터에 추가되면, GPU Operator가 자동으로 드라이버와 라이브러리를 설치하고, 노드의 상태를 모니터링합니다. 또한, GPU 장애가 발생할 경우 자동으로 알람을 발생시키고, 필요시 파드를 다른 노드로 이동시켜 서비스 중단을 최소화합니다. 이러한 자동화는 대규모 AI 인프라 운영에 있어 필수적인 요소입니다.

모델 서빙 3종 비교 표

모델 서빙 솔루션	배포 방식	오토스케일 지원	GPU 지원
KServe	CRD/YAML	HPA/KEDA	GPU Operator
vLLM	Deployment/Helm	HPA/KEDA	GPU Operator
Triton	Deployment/Helm	HPA/KEDA	GPU Operator

8.1.2 Milvus·Weaviate·Qdrant Vector DB의 Operator·Helm 배포 표준

Operator와 Helm Chart 표준화

Milvus, Weaviate, Qdrant, Chroma 등 주요 벡터 데이터베이스(Vector DB)는 Kuber-

netes 환경에서 Operator 또는 Helm Chart를 표준 배포 방식으로 제공합니다. Milvus는 공식 Operator를 통해 StatefulSet, PersistentVolume, 백업/복구, 업그레이드 등 운영 자동화를 지원합니다. Weaviate와 Qdrant 역시 Helm Chart를 통해 빠른 설치와 버전 관리가 가능하며, CRD 기반 배포를 지원하는 경우도 많습니다. 이 표준화는 AI 스택 통합 표면을 일관되게 유지하며, 운영팀이 별도의 배포 스크립트나 수동 작업 없이 GitOps 파이프라인에서 벡터 DB를 선언적으로 관리할 수 있게 합니다.

Operator와 Helm Chart의 활용은 배포의 일관성과 신뢰성을 보장합니다. 예를 들어, Milvus Operator를 사용하면 클러스터의 상태를 지속적으로 모니터링하고, 장애 발생 시 자동으로 복구 작업을 수행합니다. Helm Chart는 파라미터화된 값 파일을 통해 다양한 환경에 맞는 설정을 손쉽게 적용할 수 있어, 개발, 테스트, 운영 환경 간의 차이를 최소화할 수 있습니다. 이러한 자동화와 표준화는 대규모 AI 서비스에서 운영 비용을 크게 절감하는 효과를 가져옵니다.

StatefulSet 및 CSI 스토리지 연동

벡터 DB는 대량의 임베딩 데이터와 인덱스를 저장하기 때문에 Kubernetes의 StatefulSet 오브젝트와 CSI(Container Storage Interface) 스토리지 연동이 필수적입니다. Milvus Operator는 PVC(PersistentVolumeClaim)를 자동 생성하고, CSI 드라이버를 통해 블록/파일 스토리지와 연동합니다. Weaviate, Qdrant도 StatefulSet을 통해 데이터 영속성을 보장하며, 스냅샷 백업, 복구, 스케일 아웃/인 기능을 제공합니다. VM 환경에서는 스토리지 연결과 데이터 이동이 복잡하지만, K8s에서는 스토리지 클래스를 선언적으로 지정하고, 장애 복구나 확장 시에도 표준 API로 처리할 수 있습니다.

StatefulSet과 CSI 연동의 실제 사례로는, 데이터베이스 장애 발생 시 자동으로 새로운 파드가 생성되고, 기존 PVC를 마운트하여 데이터 손실 없이 서비스를 복구하는 과정을 들 수 있습니다. 또한, 스토리지 용량이 부족할 경우, PVC의 크기를 확장하거나, 스토리지 클래스를 변경하여 성능을 개선할 수 있습니다. 이러한 유연성은 대규모 데이터 처리와고가용성이 요구되는 AI 서비스에서 매우 중요한 역할을 합니다.

백업 전략과 운영 자동화

벡터 DB의 백업은 Operator에서 제공하는 스냅샷, PVC 복제, 클러스터 백업 기능을 활용합니다. Milvus Operator는 etcd와 연동하여 클러스터 상태를 관리하고, 데이터 백업/복구를 자동화합니다. Weaviate와 Qdrant도 Helm Chart 값 파일에서 백업 경로, 스케줄, 복구 옵션을 선언할 수 있습니다. 이처럼 K8s 표준 배포 단위와 운영 자동화 기능은 AI 스택의 데이터 신뢰성과

운영 효율성을 크게 높입니다.

운영 자동화의 구체적인 예로, 정기적인 스냅샷 백업 스케줄을 설정하여 데이터 손실 위험을 최소화하고, 장애 발생 시 신속하게 복구할 수 있습니다. 또한, Operator는 클러스터 업그레이드나 스케일 아웃/인 작업도 자동으로 처리하여, 운영자의 개입 없이 안정적인 서비스를 유지할 수 있습니다. 이러한 자동화는 인적 오류를 줄이고, 서비스 가용성을 극대화하는 데 큰 도움이 됩니다.

Vector DB 3종 비교 표

Vector DB	Operator 유무	Helm Chart 제공	CSI 스토리지 지원
Milvus	O	O	O
Weaviate	X	O	O
Qdrant	X	O	O

8.1.3 LangChain·LlamaIndex·MCP 서버·Agent의 컨테이너 배포 패턴

RAG 프레임워크의 K8s 배포 패턴

LangChain, LlamaIndex, LangGraph 등 RAG(Retrieval Augmented Generation) 프레임워크는 API 서버 형태로 K8s Deployment 또는 CRD로 배포됩니다. 이들 프레임워크는 벡터 DB, LLM, 외부 API와 연동되어 복합 워크플로우를 구성하며, GitOps 방식으로 버전 관리, 롤백, Canary 배포가 가능합니다. K8s 환경에서는 RAG 프레임워크를 선언형 YAML로 배포하고, Service 오브젝트를 통해 외부 트래픽을 라우팅하며, HPA/KEDA로 자동 확장도 지원합니다.

RAG 프레임워크의 실제 배포 예시를 보면, LangChain API 서버를 Deployment로 배포하고, Service를 통해 외부에서 접근할 수 있도록 설정합니다. 이때, 환경 변수나 ConfigMap을 활용하여 데이터베이스 연결 정보, LLM API 키, 외부 API 엔드포인트 등을 관리할 수 있습니다. 또한, GitOps 파이프라인을 통해 코드 변경 시 자동으로 배포가 이루어지며, 문제가 발생하면 이전 버전으로 신속하게 롤백할 수 있습니다. 이러한 구조는 서비스의 신뢰성과 유지보수성을 크게 향상시킵니다.

MCP 서버 및 Agent 런타임의 컨테이너화

MCP(Model Context Protocol) 서버는 Anthropic 등 글로벌 AI 벤더가 제안하는 LLM 컨텍스트 표준으로, K8s에서 컨테이너화하여 배포할 수 있습니다. MCP 서버는 CRD 또는

Deployment로 관리되며, API Gateway와 결합하여 인증, 트래픽 관리, mTLS를 적용할 수 있습니다. Agent 런타임 역시 K8s에서 파드 단위로 배포되어, Service Mesh(Istio, Linkerd, Cilium)와 연동하여 트래픽 관리, 관측성, 보안 정책을 일관되게 적용할 수 있습니다. 급변하는 AI 영역에서 선언형·불변 운영 모델은 안정성과 확장성의 핵심 전제가 됩니다.

컨테이너화된 MCP 서버와 Agent 런타임은 다양한 환경에서 일관된 방식으로 배포 및 운영이 가능합니다. 예를 들어, 새로운 Agent 기능을 추가하거나 MCP 프로토콜이 업데이트될 경우, 컨테이너 이미지만 교체하면 전체 클러스터에 신속하게 적용할 수 있습니다. 또한, Service Mesh와의 연동을 통해 트래픽 분산, 장애 복구, 실시간 모니터링 등 고급 네트워크 기능을 손쉽게 구현할 수 있습니다. 이러한 컨테이너화와 자동화는 AI 서비스의 민첩성과 확장성을 극대화합니다.

API Gateway·Service Mesh 연동

RAG/Agent/MCP 서버는 API Gateway(K8s Ingress, Istio Gateway 등)와 Service Mesh와 결합하여 mTLS 인증, 트래픽 분산, 관측성(Distributed Tracing, Metrics)을 제공합니다. YAML 기반 선언형 배포는 운영자의 실수나 드리프트를 방지하고, 감사 로그, 정책 엔진(OPA/Gatekeeper)과 연동하여 보안 및 규제 대응을 강화합니다.

API Gateway와 Service Mesh의 연동은 보안과 운영 효율성 측면에서 매우 중요합니다. 예를 들어, Istio Gateway를 통해 외부 트래픽을 안전하게 내부 서비스로 전달하고, mTLS를 적용하여 통신 구간의 보안을 강화할 수 있습니다. 또한, Distributed Tracing을 통해 각 요청의 흐름을 추적하여 장애 원인을 신속하게 파악할 수 있습니다. 정책 엔진과의 연동을 통해 접근 제어, 리소스 할당, 규제 준수 등 다양한 요구사항을 코드로 관리할 수 있어, 대규모 AI 서비스의 안정성과 신뢰성을 확보할 수 있습니다.

RAG/Agent/MCP 구성요소 매핑 표

구성요소	K8s 오브젝트	연동 기능	운영 모델
LangChain	Deployment/CRD	Service, HPA/KEDA	선언형/불변
LlamaIndex	Deployment/CRD	Service, HPA/KEDA	선언형/불변
MCP 서버	Deployment/CRD	Gateway, Mesh, mTLS	선언형/불변
Agent 런타임	Deployment/CRD	Mesh, Tracing	선언형/불변

8.2 GPU 공유·MIG·Time-slicing으로 본 VM 비경제성의 기술적 근거

AI 인프라의 단가에서 가장 중요한 요소는 GPU 활용률입니다. VM 환경에서는 GPU 패스스루 방식으로 한 VM이 GPU를 독점 점유하게 되어 평균 활용률이 극히 낮아지고, 인프라 비용이 2~5배까지 상승하는 사례가 많습니다. 반면 Kubernetes에서는 NVIDIA GPU Operator, Device Plugin, MIG(Multi-Instance GPU), Time-slicing 등 표준 기능을 통해 GPU를 다중 파드에 분할·공유할 수 있어, 동일 하드웨어에서 훨씬 높은 활용률과 비용 효율성을 달성할 수 있습니다. 본 절에서는 GPU 공유 기술의 구조와 VM 환경의 한계, 그리고 실제 비용 차이를 기술적으로 분석합니다.

8.2.1 NVIDIA GPU Operator·Device Plugin·MIG·Time-slicing의 다중 파드 공유

GPU Operator 통합 구조

NVIDIA GPU Operator는 Kubernetes에서 GPU 드라이버, CUDA, MIG, DCGM(모니터링)을 자동 설치·관리하는 DaemonSet입니다. GPU Operator는 노드별로 GPU 관련 소프트웨어를 배포하고, Device Plugin을 통해 파드 단위 GPU 할당을 지원합니다. MIG(Multi-Instance GPU)는 A100, H100 등 최신 GPU에서 하나의 GPU를 여러 인스턴스로 분할하여 각 파드에 할당할 수 있게 합니다. Time-slicing은 GPU를 시간 단위로 분할하여 여러 파드가 번갈아 사용할 수 있게 하며, AI 워크로드의 피크/유휴 패턴에 맞춰 GPU 활용률을 극대화합니다.

GPU Operator의 통합 구조는 운영 자동화와 확장성 측면에서 큰 장점을 제공합니다. 예를 들어, GPU 노드가 추가되면 GPU Operator가 자동으로 드라이버와 라이브러리를 설치하고, 클러스터 전체에 일관된 GPU 환경을 제공합니다. 또한, Device Plugin을 통해 파드가 GPU 리소스를 요청할 때, 하드웨어 상태와 사용 현황을 실시간으로 반영하여 최적의 할당을 수행합니다. 이러한 자동화는 대규모 AI 인프라에서 운영자의 부담을 줄이고, 장애 대응 속도를 높여줍니다.

Device Plugin·MIG·Time-slicing 기술적 구현

Device Plugin은 Kubernetes에서 파드가 GPU 리소스를 요청할 때, 노드에 설치된 플러그인이 해당 리소스를 할당·분할하는 역할을 합니다. MIG는 GPU를 하드웨어적으로 분할하여 각 인스턴스에 독립된 메모리, 컴퓨팅 리소스를 할당합니다. Time-slicing은 소프트웨어적으로

GPU 사용 시간을 분할하여 여러 파드가 순차적으로 GPU를 사용할 수 있게 하며, 모델 서빙, RAG, Agent 워크로드 등에서 유휴 GPU를 효율적으로 배분합니다.

기술적으로, Device Plugin은 Kubernetes API와 연동되어 파드 스케줄링 시 GPU 리소스의 가용성을 실시간으로 반영합니다. MIG는 GPU 하드웨어의 펌웨어와 드라이버가 연동되어, 각 인스턴스가 독립적으로 동작할 수 있도록 지원합니다. Time-slicing은 NVIDIA의 소프트웨어 스택이 GPU 사용 시간을 분할하여, 여러 파드가 공정하게 리소스를 사용할 수 있게 합니다. 이러한 기술들은 실제 운영 환경에서 GPU 활용률을 극대화하고, 다양한 AI 워크로드의 요구사항을 유연하게 충족시킵니다.

다중 파드 공유의 인프라 단가 효과

GPU 공유 방식은 AI 인프라 단가에 직접적인 영향을 미칩니다. 예를 들어, A100 GPU 1장을 VM 환경에서는 1개 VM이 독점하지만, K8s에서는 MIG로 7개 인스턴스, Time-slicing으로 1020개 파드가 공유할 수 있습니다. 이로 인해 평균 GPU 활용률이 2030%에서 6080%까지 상승하며, 인프라 단가가 25배 절감됩니다. AI 프로젝트의 파일럿 단계에서 GPU 공유 모델을 실측하면, 예산 산정과 인프라 구매 전략에 큰 변화를 줄 수 있습니다.

실제 기업 사례를 보면, 기존 VM 기반 GPU 인프라에서 Kubernetes 기반 GPU 공유 모델로 전환한 후, 동일한 하드웨어에서 더 많은 AI 워크로드를 처리할 수 있게 되었으며, GPU 리소스의 유휴 시간이 크게 줄어들었습니다. 또한, GPU 장애 발생 시 자동으로 워크로드를 다른 파드로 이동시켜 서비스 중단을 최소화할 수 있었습니다. 이러한 변화는 AI 인프라의 총소유비용(TCO)을 획기적으로 낮추는 데 기여합니다.

GPU 공유 방식 비교 표

공유 방식	적합 시나리오	성능 영향
Device Plugin	파드 단위 할당	최소 오버헤드
MIG	하드웨어 분할	분할당 독립 성능
Time-slicing	유휴/피크 반복 워크로드	소폭 지연 발생

8.2.2 VM GPU 패스스루의 전용 점유 구조와 AI 인프라 단가 영향

VM GPU 패스스루 구조

VM 환경에서 GPU 패스스루는 하이퍼바이저가 GPU를 특정 VM에 직접 할당하는 방식입니다. 이 방식은 GPU를 VM이 독점 점유하게 하여, 다른 VM이나 워크로드가 해당 GPU를 사용할 수 없습니다. 결과적으로 GPU 활용률이 낮은 조직에서는 대부분의 시간이 유휴 상태로 남게 되고, GPU 단가가 급격히 상승합니다. GPU 패스스루는 VM의 격리 요구에는 적합하지만, AI 워크로드의 피크/유휴 반복 패턴에는 비경제적입니다.

VM GPU 패스스루의 구조적 한계는 인프라 확장성과 비용 효율성에 큰 영향을 미칩니다. 예를 들어, 새로운 AI 워크로드가 추가될 때마다 별도의 VM과 GPU를 할당해야 하므로, 하드웨어 자원의 낭비가 심해집니다. 또한, GPU 장애가 발생하면 해당 VM 전체가 영향을 받아 서비스 중단 위험이 커집니다. 이러한 구조는 대규모 AI 서비스에서 운영 효율성과 비용 경쟁력을 저하시킬 수밖에 없습니다.

평균 GPU 활용률과 인프라 단가

실무에서 VM 환경의 평균 GPU 활용률은 10~30% 수준에 머무는 경우가 많습니다. AI 워크로드가 없는 시간에도 GPU는 VM에 할당되어 유휴 상태가 되고, GPU 예산은 실제 사용량 대비 2~5배까지 증가합니다. 반면 K8s 환경에서는 GPU를 다중 파드가 공유하여 활용률이 60~80% 까지 상승하며, 동일 예산으로 2~5배 많은 AI 워크로드를 처리할 수 있습니다. AI 인프라 단가를 낮추기 위해서는 VM 패스스루 대신 K8s GPU 공유 모델로 전환하는 것이 필수적입니다.

비용 측면에서, VM 기반 GPU 인프라는 초기 투자 비용뿐만 아니라, 운영 비용, 유지보수 비용이 지속적으로 발생합니다. 반면, Kubernetes 기반 GPU 공유 모델은 리소스 활용률을 극대화하여, 동일한 예산으로 더 많은 프로젝트를 동시에 운영할 수 있습니다. 실제로, 여러 AI 스타트업과 대기업에서 K8s 기반 GPU 공유 모델을 도입한 후, 인프라 비용이 절반 이하로 줄어들고, 서비스 확장 속도가 빨라졌다는 사례가 보고되고 있습니다.

VM 패스스루 vs K8s Device Plugin 비교 표

환경	공유 단위	평균 활용률	인프라 단가(예시)
VM 패스스루	VM 독점	10~30%	2~5배 상승
K8s Device Plugin	파드/MIG/Time-slicing	60~80%	2~5배 절감

9장: 관리 포인트 이중화 비용과 단일 K8s 표면으로의 수렴 전략

클라우드 네이티브 전환의 핵심은 단순히 VM에서 컨테이너로 워크로드를 옮기는 것이 아니라, 운영팀·모니터링·보안·CI/CD 등 관리 포인트를 단일 Kubernetes 표면으로 통합하는 데 있다. 기존 VM 환경과 컨테이너 환경이 공존할 때 발생하는 FTE(Full-Time Equivalent, 인력), 도구, 정책의 중복은 조직의 운영 효율성과 비용에 직접적인 영향을 미친다. 본 장에서는 VM과 컨테이너가 공존할 때 발생하는 이중 운영 비용의 구조를 구체적으로 분석하고, 단일 K8s 표면으로의 수렴 전략과 표준 스택을 제시한다. 이를 통해 독자는 현재 조직의 이중 운영 구조를 진단하고, K8s 기반 통합 운영 모델로의 이행 경로를 설계할 수 있다.

9.1 VM과 컨테이너가 공존할 때 발생하는 이중 운영 비용

VM과 컨테이너 환경이 동시에 존재하는 조직에서는 운영팀의 FTE, 모니터링 및 보안 도구, 정책 배포 경로가 중복되어 관리 비용이 급증한다. 특히 장애 대응, 보안 사고, 배포 파이프라인 등에서 이중화로 인한 복잡성이 조직의 민첩성과 안정성에 부정적 영향을 미친다. 이 절에서는 FTE·모니터링·보안 정책의 3중 중복 구조와 CI/CD·백업·DR 파이프라인의 분기 유지 비용을 분석한다.

9.1.1 FTE·모니터링·보안 정책의 3중 중복 경로

VM 환경과 컨테이너 환경이 공존하는 조직에서는 각각의 운영팀이 별도로 존재합니다. VM 운영팀은 하이퍼바이저(vCenter, ESXi 등), OS 패치, VM 템플릿 관리에 집중하며, 컨테이너 운영팀은 Kubernetes 클러스터 관리, 컨테이너 배포, 네이티브 오브젝트 관리에 집중합니다. 이중 운영은 FTE(Full-Time Equivalent) 인력의 중복 배치로 이어지며, 각 팀이 별도의 도구와 정책을 유지해야 하므로 인건비와 교육 비용이 증가합니다. 장애 발생 시 두 팀이 동시에 호출되어 MTTR(Mean Time To Recovery)이 늘어나는 현상이 빈번하게 발생합니다.

운영팀의 이중화는 단순히 인력의 중복만을 의미하지 않습니다. 각 환경별로 요구되는 전문성이 다르기 때문에, VM 운영팀과 K8s 운영팀은 별도의 교육, 인증, 경험을 필요로 하며, 이로 인해 인력 채용 및 유지 비용이 증가합니다. 예를 들어, VM 운영팀은 하드웨어 및 가상화 계층에 대한

이해가 필수적이며, K8s 운영팀은 컨테이너 오케스트레이션, 네트워크, 보안 정책 등 클라우드 네이티브 기술에 대한 깊은 이해가 요구됩니다. 이처럼 인력 구조가 이원화되면, 조직 내 지식 공유와 협업이 저해되고, 장애 발생 시 책임 소재가 불분명해지는 문제가 발생합니다.

모니터링 도구의 중복 역시 심각한 비효율을 초래합니다. VM 환경에서는 vCenter, 별도 APM(Application Performance Monitoring), OS 에이전트(예: Nagios, Zabbix)가 주로 사용됩니다. 반면 K8s 환경에서는 Prometheus Operator, Grafana, OpenTelemetry, Loki 등 클라우드 네이티브 관측성 도구가 표준입니다. 두 환경의 모니터링 도구가 일치하지 않으면 알람, 대시보드, 장애 대응 프로세스가 이중화되어 운영팀 간 정보 공유가 어렵고, 장애 발생 시 원인 추적이 복잡해집니다. 예를 들어 VM 장애 알람은 vCenter에서, 컨테이너 장애는 Prometheus에서 각각 발생하여 대응팀이 분리됩니다. 이로 인해 장애의 근본 원인을 신속하게 파악하기 어렵고, 복구 시간이 지연될 수 있습니다.

보안 정책 배포 경로의 이중화도 조직에 큰 부담을 줍니다. VM 환경에서는 OS 단위 보안 정책(EDR, 백신, 패치 관리 등)이 적용되며, 컨테이너 환경에서는 OPA(Open Policy Agent), Gatekeeper, Falco, Trivy, Cosign 등 네이티브 정책 엔진이 활용됩니다. 정책 배포 경로가 이원화되지 않으면 동일한 보안 정책을 두 환경에 각각 배포해야 하며, 정책 적용·감사·컴플라이언스 증빙이 중복됩니다. 특히 금융·공공 조직에서는 보안 심사 대응을 위해 OS 단위와 K8s 단위 모두에서 감사 로그를 별도로 수집해야 하는 부담이 발생합니다. 이중화된 보안 정책은 정책 누락, 적용 오류, 감사 실패 등 다양한 리스크를 내포하고 있습니다.

이러한 중복 구조를 한눈에 파악할 수 있도록 아래와 같이 매트릭스를 정리할 수 있습니다.

관리 포인트	VM 전용	K8s 전용	공용(통합)
FTE	하이퍼바이저·OS 운영팀	K8s 운영팀	단일 K8s 운영팀
모니터링	vCenter·APM·OS 에이전트	Prometheus·OpenTelemetry	Prometheus Operator·Grafana
보안 정책	OS 단위 EDR·백신	OPA·Falco·Trivy·Cosign	OPA·Falco·Trivy·Cosign

이처럼 운영팀, 모니터링, 보안 정책의 3중 중복은 단순한 비용 증가를 넘어, 조직의 민첩성과 보안 수준, 장애 대응 능력 전반에 부정적인 영향을 미치게 됩니다. 실제로 국내외 대규모 금융 기관이나 제조기업의 사례를 보면, VM과 컨테이너 환경을 병행 운영할 때 연간 수억 원 규모의 추가 인건비와 도구 라이선스 비용이 발생하는 것으로 나타났습니다. 또한, 장애 대응 프로세스가

이원화되어 MTTR이 30% 이상 증가하는 사례도 보고되고 있습니다. 이러한 문제를 해결하기 위해서는 관리 포인트의 일원화와 표준화가 필수적이며, 단일 K8s 표면으로의 수렴이 점차 업계 표준으로 자리잡고 있습니다.

9.1.2 CI/CD·백업·DR 파이프라인의 분기 유지 비용

VM 환경과 컨테이너 환경이 병존하는 조직에서는 CI/CD, 백업, DR(Disaster Recovery) 파이프라인이 각각 분리되어 운영되는 경우가 많습니다. 이로 인해 파이프라인 코드의 중복, 자격 증명 관리의 복잡성, 테스트 및 검증 환경의 이원화 등 다양한 비용과 비효율이 발생합니다. 이 절에서는 이러한 분기 유지 비용의 구체적인 구조와 실제 운영상에서 발생하는 문제점, 그리고 비교 분석을 통해 단일화의 필요성을 강조합니다.

CI/CD 파이프라인의 이중화는 조직 내 개발 및 배포 프로세스의 복잡성을 크게 증가시킵니다. VM 환경에서는 이미지 빌드, VM 템플릿화, Ansible/Puppet/Salt 등의 구성관리 도구를 활용한 명령형 배포가 일반적입니다. 컨테이너 환경에서는 Tekton, Argo CD, Flux, Helm, GitOps 기반의 선언형 배포가 표준입니다. 두 환경의 CI/CD 파이프라인이 병존하면 파이프라인 코드, 자격 증명, 스테이징 환경이 2중으로 유지되어 관리 비용이 증가합니다. 예를 들어 VM용 Ansible 플레이북과 K8s용 Helm Chart를 각각 관리해야 하며, 배포 실패 시 원인 분석과 롤백 경로도 이중화됩니다.

이러한 이중화는 단순히 도구의 중복에 그치지 않습니다. 각 파이프라인의 유지보수, 버전 관리, 보안 패치 적용, 개발자 교육 등 모든 측면에서 추가적인 리소스가 소모됩니다. 또한, 배포 프로세스의 일관성이 저해되어, 동일한 애플리케이션이라도 VM 환경과 K8s 환경에서 동작 방식이나 배포 결과가 달라질 수 있습니다. 이는 품질 저하, 장애 발생 시 원인 분석의 어려움, 그리고 규정 준수(컴플라이언스) 측면에서의 리스크로 이어집니다.

백업 및 DR 파이프라인의 분기도 유사한 문제를 야기합니다. VM 환경에서는 VM 스냅샷, OS 단위 백업 에이전트, 별도 DR 솔루션이 사용됩니다. 컨테이너 환경에서는 Velero(PV·클러스터 백업), Crossplane(Infrastructure as Code), Terraform+K8s Operator를 통한 선언형 DR 구성이 표준입니다. 두 환경의 백업·DR 파이프라인이 분기되면 데이터 보호·복구 시나리오가 중복 설계되어야 하며, 보안 사고 대응도 이중으로 설계됩니다. 특히 스테이징 환경이 VM과 K8s로 분리되어 테스트·검증 비용이 상승합니다.

실제 사례를 보면, 한 글로벌 제조기업은 VM 기반 ERP 시스템과 K8s 기반 마이크로서비스를 동시에 운영하면서, 백업 정책과 DR 계획을 별도로 수립해야 했습니다. 이로 인해 백업 스케줄, 데이터 복구 테스트, 감사 로그 관리 등 모든 단계에서 중복 작업이 발생했고, DR 훈련 시나리오도 두 환경을 각각 준비해야 했습니다. 결과적으로 연간 수천만 원의 추가 운영 비용과 인력 투입이 필요했으며, 장애 발생 시 복구 절차가 복잡해져 전체 시스템의 복원력이 저하되는 문제가 있었습니다.

아래 표는 VM과 컨테이너 환경의 CI/CD 및 백업·DR 파이프라인의 주요 비교 요소를 정리한 것입니다.

항목	VM CI/CD	컨테이너 CI/CD	비교 요소
툴	Ansible·Puppet·Salt	Tekton·Argo CD·Flux·Helm	도구 중복
자격 증명	OS 계정·SSH Key	K8s ServiceAccount·OIDC	인증·권한 관리 중복
스테이징	VM 템플릿·OS 이미지	컨테이너 이미지·Pod	환경 중복

이처럼 CI/CD, 백업, DR 파이프라인의 분기 유지 비용은 단순한 라이선스나 인력 비용을 넘어, 조직의 민첩성, 품질, 보안, 규정 준수 등 전반적인 IT 거버넌스에 부정적인 영향을 미치게 됩니다. 따라서 관리 포인트의 단일화와 표준화가 필수적이며, 단일 K8s 표면으로의 수렴이 점차 필수 전략으로 자리잡고 있습니다.

9.2 단일 K8s 표면으로 통합하는 운영 아키텍처

이중 운영 비용을 해소하기 위해서는 관리 포인트를 단일 Kubernetes 표면으로 수렴하는 전략이 필요하다. Prometheus·OpenTelemetry·Falco·Argo CD·Tekton·Crossplane 등 클라우드 네이티브 표준 스택을 활용하면 관측성, 보안, CI/CD, IaC, DR 파이프라인을 하나의 플랫폼에서 일관되게 운영할 수 있다. 이 절에서는 관측·보안 통합과 GitOps·IaC·DR 통합 아키텍처를 구체적으로 제시한다.

9.2.1 Prometheus·OpenTelemetry·Falco·Trivy·Cosign으로 통합되는 관측·보안

클라우드 네이티브 환경에서 운영 효율성과 보안 수준을 동시에 높이기 위해서는 관측성과 보안 관리 포인트를 단일화하는 것이 중요합니다. 이 절에서는 LGTM 스택과 OpenTelemetry를 중심으로 한 관측 통합, 그리고 Falco, Trivy, Cosign 등 K8s 네이티브 보안 도구를 활용한 보안 통합 전략을 구체적으로 설명합니다. 이를 통해 조직은 장애 대응, 감사, 컴플라이언스, 공급망 보안 등 다양한 요구사항을 단일 K8s 표면에서 일관되게 충족할 수 있습니다.

LGTM 스택과 OpenTelemetry의 관측 통합은 클라우드 네이티브 운영의 핵심입니다. Prometheus Operator는 K8s 클러스터의 메트릭을 자동 수집·시각화하며, Grafana는 대시보드 통합을 제공합니다. Loki는 로그 집계, Tempo는 분산 추적, Mimir는 대규모 시계열 데이터 저장을 담당합니다. OpenTelemetry는 분산 추적 표준으로, VM·컨테이너·마이크로서비스 환경 모두에서 동일한 추적 데이터를 수집·분석할 수 있습니다. 장애 대응뿐 아니라 내부 감사·컴플라이언스 증빙의 단일 진실 공급원(single source of truth)으로 활용됩니다. 예를 들어, 장애 발생 시 Prometheus와 OpenTelemetry를 통해 메트릭과 분산 추적 데이터를 통합 분석함으로써, 문제의 근본 원인을 신속하게 파악할 수 있습니다. 또한, Grafana 대시보드를 통해 운영팀과 개발팀이 동일한 정보를 실시간으로 공유할 수 있어, 협업 효율성이 크게 향상됩니다.

보안 통합 측면에서는 Falco, Trivy, Cosign이 핵심 역할을 합니다. Falco는 런타임 위협 탐지 엔진으로, K8s 클러스터 내에서 실시간 이상 행위를 감지합니다. 예를 들어, 비정상적인 시스템 콜이나 권한 상승 시도를 즉시 탐지하여, 운영팀에 경고를 보냅니다. Trivy는 컨테이너 이미지 취약점 스캐너로, 빌드 단계에서 보안 결함을 사전에 차단합니다. 이를 통해 취약한 이미지가 프로덕션 환경에 배포되는 것을 방지할 수 있습니다. Cosign은 이미지 서명·검증을 통해 공급망 보안을 강화합니다. 이미지가 신뢰할 수 있는 소스에서 빌드·배포되었는지 검증함으로써, 공급망 공격에 대한 방어력을 높입니다. 이들 도구는 K8s 네이티브 방식으로 배포되어 OS 단위 에이전트와 달리 클러스터 전체를 통합적으로 보호합니다.

추가적으로, Clair, SPIFFE/SPIRE, Sigstore 등 공급망 보안 도구와 연계하면, 보안 정책 배포·감사·컴플라이언스 대응이 단일 표면에서 일원화됩니다. 예를 들어, Sigstore를 활용하면 이미지 서명과 검증 과정을 자동화할 수 있으며, SPIFFE/SPIRE를 통해 서비스 간 신뢰성을 강화할 수 있습니다. 이러한 통합 아키텍처는 금융, 공공, 제조 등 규제 산업에서 요구하는 높은 수준의

보안과 감사 요건을 효과적으로 충족시킵니다.

아래 표는 주요 관측·보안 스택의 역할과 배포 단위, 대체 도구를 비교한 매트릭스입니다.

스택/도구	역할	배포 단위	대체 도구
Prometheus	메트릭 수집·알람	K8s Operator·Pod	Datadog, Zabbix
Grafana	시각화·대시보드	K8s Operator·Pod	Kibana, Tableau
Loki	로그 집계	K8s Operator·Pod	ELK Stack
Tempo	분산 추적	K8s Operator·Pod	Jaeger, Zipkin
Mimir	시계열 데이터 저장	K8s Operator·Pod	InfluxDB
OpenTelemetry	분산 추적 표준	K8s Operator·Pod	Jaeger, Zipkin
Falco	런타임 위협 탐지	K8s DaemonSet	OS EDR·백신
Trivy	이미지 취약점 스캔	K8s Job/Pod	Clair
Cosign	이미지 서명·검증	K8s Job/Pod	Notary, Sigstore

이처럼 관측성과 보안의 단일화는 운영 효율성, 장애 대응 속도, 보안 수준, 컴플라이언스 대응력 등 모든 측면에서 조직에 큰 이점을 제공합니다. 실제로 대형 금융기관이나 클라우드 서비스 제공업체는 LGTM 스택과 OpenTelemetry, Falco, Trivy, Cosign을 표준 스택으로 채택하여, 운영팀과 보안팀의 협업을 강화하고, 장애 및 보안 사고 대응 시간을 절반 이하로 단축한 사례가 보고되고 있습니다. 이러한 통합 전략은 앞으로 클라우드 네이티브 운영의 필수 요소로 자리매김할 것입니다.

9.2.2 Argo CD·Flux·Tekton·Crossplane으로 통합되는 GitOps·IaC·DR

운영 자동화와 인프라 관리의 일관성을 확보하기 위해서는 GitOps, IaC, DR 파이프라인을 단일 K8s 표면에서 통합하는 것이 매우 중요합니다. 이 절에서는 Argo CD, Flux, Tekton, Crossplane, Velero 등 클라우드 네이티브 표준 도구를 활용한 통합 전략을 구체적으로 설명합니다. 이를 통해 조직은 배포 자동화, 인프라 선언, 재해 복구 등 모든 운영 프로세스를 코드 기반으로 일원화할 수 있습니다.

GitOps 기반 선언형 동기화는 현대적 클라우드 운영의 핵심입니다. Argo CD와 Flux는 GitOps 기반의 클러스터 선언형 동기화 도구로, Git 저장소를 단일 진실 공급원으로 삼아 배포·롤백·감사·재현을 자동화합니다. 운영팀은 YAML·Helm Chart·Kustomize 등 선언형 오브젝트를

Git에 커밋하면 Argo CD/Flux가 클러스터 상태를 자동으로 맞춥니다. 이를 통해 배포 프로세스의 일관성, 변경 이력의 추적, 감사 대응이 크게 향상됩니다. 변경 관리 KPI로 “Git 기준 롤백 가능한 비율”을 제안할 수 있으며, 감사·컴플라이언스 대응도 Git 로그로 일원화됩니다. 실제로, 글로벌 IT 기업들은 GitOps 도입 이후 배포 실패율이 30% 이상 감소하고, 롤백 시간도 대폭 단축된 사례를 보고하고 있습니다.

Tekton 파이프라인의 CI/CD 통합은 K8s 네이티브 환경에서 빌드·테스트·배포 작업을 자동화하는 데 필수적입니다. Tekton은 K8s 네이티브 CI/CD 파이프라인 엔진으로, 빌드·테스트·배포 작업을 Pod 단위로 실행합니다. GitHub Actions, GitLab CI 등과 연계하여 워크로드별 파이프라인을 통합 관리할 수 있습니다. VM 환경의 명령형 배포(Ansible, Puppet)와 달리, Tekton은 선언형 파이프라인으로 운영팀 간 협업·자동화 수준을 높입니다. 예를 들어, Tekton을 활용하면 개발팀과 운영팀이 동일한 파이프라인 정의를 공유하고, 변경 사항을 코드 리뷰와 PR(풀 리퀘스트)로 관리할 수 있어, 배포 품질과 협업 효율성이 크게 향상됩니다.

Crossplane·Velero의 IaC·DR 통합은 인프라 관리와 재해 복구의 일관성을 보장합니다. Crossplane은 K8s Operator 기반의 Infrastructure as Code(IaC) 솔루션으로, 클라우드·온프레미스 인프라를 K8s 오브젝트로 선언·관리합니다. 이를 통해 VM, 네트워크, 데이터베이스 등 다양한 리소스를 코드로 관리할 수 있으며, GitOps와 연계하여 변경 이력을 추적할 수 있습니다. Velero는 클러스터·PersistentVolume 백업·복구를 자동화하여 DR(Disaster Recovery) 시나리오를 단일 표면에서 구현합니다. Terraform+K8s Operator와 연계하면 VM·컨테이너 환경 모두에서 IaC·DR을 일원화할 수 있습니다. 실제로, 대규모 금융기관은 Crossplane과 Velero를 도입하여, DR 테스트 자동화, 백업 정책 일원화, 감사 로그 통합 등 다양한 운영 효율성을 확보한 사례를 보이고 있습니다.

아래 표는 GitOps, CI/CD, IaC, DR의 대표 도구와 역할, 매핑 방식을 정리한 것입니다.

축	대표 도구	역할	매핑 방식
GitOps	Argo CD·Flux	선언형 동기화·롤백·감사	Git→K8s 자동화
CI/CD	Tekton·GitHub Actions	빌드·테스트·배포	Pod·Job 단위 파이프라인
IaC	Crossplane·Terraform	인프라 선언·구성	K8s Operator·YAML
DR	Velero	백업·복구	클러스터·PV 단위 자동화

이처럼 Argo CD, Flux, Tekton, Crossplane, Velero 등 클라우드 네이티브 표준 도구를 활용하면, 운영팀은 모든 관리 포인트를 단일 K8s 표면에서 일관되게 관리할 수 있습니다. 이는 인력 효율성, 운영 자동화, 보안, 규정 준수 등 모든 측면에서 조직의 경쟁력을 높이는 핵심 전략입니다. 앞으로 클라우드 네이티브 전환을 추진하는 모든 조직은 이러한 통합 아키텍처를 적극적으로 도입해야 할 것입니다.

10장: 세계적 클라우드 네이티브 전환 트렌드 — CNCF 공식 정의와 글로벌 레퍼런스로 본 표준화 흐름

10.1 CNCF 공식 정의와 세계적 기술 표준의 방향

세계적 클라우드 네이티브 전환 트렌드는 단순히 가상머신(IaaS)에서 컨테이너로의 기술적 전환을 의미하지 않습니다. 이는 선언형, 불변, 자동확장 기반의 운영 모델로의 패러다임 변화이며, 이러한 변화의 중심에는 CNCF(Cloud Native Computing Foundation)가 있습니다. CNCF는 클라우드 네이티브의 공식 정의와 함께, 전 세계적으로 채택되는 오픈소스 프로젝트와 기술 표준을 제시하고 있습니다. 본 절에서는 CNCF의 공식 정의와 Landscape, 그리고 연례 조사와 KubeCon에서 발표된 글로벌 엔터프라이즈 사례들을 근거로 하여, 현재 세계적으로 통용되는 클라우드 네이티브 기술 표준의 방향성과 그 의미를 심층적으로 해설합니다. 이를 통해 국내외 조직이 클라우드 네이티브 전환을 추진할 때 참고해야 할 글로벌 표준의 흐름을 명확히 파악할 수 있습니다.

10.1.1 “VM 이관이 아닌 컨테이너·K8s 기반 전환”이 세계적 주류 경로로 확립된 배경

CNCF v1.0 공식 정의는 클라우드 네이티브의 핵심을 다음과 같이 규정합니다: 컨테이너(Container), 서비스 메시(Service Mesh), 마이크로서비스(Microservices), 불변 인프라(Immutable Infrastructure), 선언형 API(Declarative API). 이 다섯 요소는 VM 중심 운영과 본질적으로 다릅니다. 컨테이너는 애플리케이션 실행 단위를 경량화하고, 서비스 메시와 마이크로서비스는 분산 시스템의 트래픽 관리와 서비스 분리를 표준화합니다. 불변 인프라는 운영 중 수정이 아닌 빌드 시점 고정 아티팩트로 배포를 강제하며, 선언형 API는 YAML이나 JSON 등으로 원하는 상태를

기술하고, 컨트롤러가 실제 상태와 비교하여 자동 수렴합니다. 이 정의는 VM 기반 운영의 명령형·가변·정적 구조와 대비되는 선언형·불변·동적 운영 모델을 표준으로 삼고 있습니다.

CNCF Landscape(<https://landscape.cncf.io/>) 에 는 Kubernetes 를 중심으로 Prometheus(모니터링), Istio(서비스 메시), Argo CD(GitOps), Helm(패키지 관리), OpenTelemetry(관측성), Falco(런타임 보안) 등 세계적으로 채택된 OSS 프로젝트들이 등재되어 있습니다. 이 Landscape는 글로벌 엔터프라이즈가 실제로 선택하는 기술 스택을 시각적으로 보여주며, VM 중심의 IaaS 프로젝트들은 Landscape에서 점차 비중이 줄고 있습니다. Landscape의 프로젝트 그래프는 각 요소가 선언형·불변·자동확장 기반으로 연동되는 구조를 강조합니다.

CNCF 연례 조사(<https://www.cncf.io/reports/cncf-annual-survey-2023/>)에 따르면, 2023년 기준 엔터프라이즈의 Kubernetes 채택률은 96%에 달하며, VM 기반 IaaS만을 운영하는 조직은 급격히 줄고 있습니다. 특히, 신규 워크로드의 80% 이상이 컨테이너 기반으로 배포되고 있으며, GitOps, Service Mesh, Operator 등 클라우드 네이티브 기술이 표준으로 자리잡았습니다. 이 조사 결과는 세계적 기술 표준이 이미 컨테이너·K8s 기반 전환으로 확립되었음을 수치로 입증합니다.

KubeCon에서 발표된 Booking.com, Spotify, Adidas 등 글로벌 엔터프라이즈 사례는 VM 이관이 아닌 컨테이너·K8s 기반 전환이 표준 경로임을 보여줍니다. Booking.com은 수천 노드 베어메탈 K8s 클러스터를 운영하며, VM 없이 직접 컨테이너 기반 워크로드를 배포합니다. Spotify는 마이크로서비스 전환과 K8s 기반 자동확장으로 신규 서비스 프로비저닝 시간을 1시간에서 수 분으로 단축했습니다. Adidas는 CNCF 레퍼런스 사례로, 4000+ 파드와 200+ 서비스 운영을 통해 배포 속도를 3~4배 개선하였습니다. 이러한 사례들은 VM 중심 운영에서 선언형·불변·자동확장 기반 클라우드 네이티브 운영으로의 전환이 글로벌 표준임을 실증적으로 뒷받침합니다.

이러한 변화의 배경에는 기술적 효율성과 운영의 유연성이 자리잡고 있습니다. VM 기반 인프라는 각 워크로드마다 별도의 OS 인스턴스를 운영해야 하므로 자원 낭비와 관리 복잡성이 높았습니다. 반면, 컨테이너 기반 클라우드 네이티브 환경에서는 애플리케이션의 배포, 확장, 복구가 자동화되어 운영 효율성이 극대화됩니다. 또한, 선언형 API와 불변 인프라의 도입으로 인해 인프라의 상태를 코드로 관리할 수 있게 되었고, 이를 통해 인프라의 일관성, 재현성, 신뢰성이 크게 향상되었습니다. 글로벌 기업들은 이러한 이점을 바탕으로 VM에서 컨테이너·K8s 기반으로 전환함으로써, 서비스의 민첩성과 확장성을 확보하고 있습니다.

CNCF v1.0 요소	VM 중심 운영과의 차이
컨테이너	VM 대비 경량화, OS 종속성 제거
서비스 메시	VM 내 네트워크 정책 대비 동적 트래픽 관리
마이크로서비스	VM 단위 모놀리식 대비 분산 서비스 구조
불변 인프라	VM 운영 중 수정 대비 빌드 시점 고정
선언형 API	명령형 스크립트 대비 자동 수렴·역등성

이 표는 CNCF가 제시한 5대 요소가 기존 VM 중심 운영과 어떻게 본질적으로 다른지 한눈에 보여줍니다. 이러한 차별점이 바로 글로벌 표준의 방향성을 결정짓는 핵심 동인입니다.

10.1.2 AI·데이터 공유 플랫폼 기반으로서의 클라우드 네이티브 — 세계적 트렌드

AI와 데이터 중심의 서비스가 급격히 확산됨에 따라, 클라우드 네이티브 플랫폼이 AI 인프라의 글로벌 표준으로 자리잡고 있습니다. 특히, Kubernetes는 AI 워크로드의 배포, 확장, 자원 할당 및 관리에 있어 사실상 표준 플랫폼으로 인정받고 있습니다. 본 절에서는 AI·데이터·서비스 공유 플랫폼의 관점에서 클라우드 네이티브가 어떻게 세계적 트렌드로 자리매김했는지, 그리고 실제 글로벌 사례와 기술적 구현 방식을 중심으로 상세히 설명하겠습니다. 이를 통해 AI 인프라 구축을 고민하는 조직이 왜 K8s 기반 접근이 필수적인지, 그리고 어떤 기술적 요소들이 글로벌 표준으로 채택되고 있는지 구체적으로 이해할 수 있습니다.

AI·데이터·서비스 공유 플랫폼의 기반으로 Kubernetes가 글로벌 표준이 되었다는 트렌드는 매우 명확합니다. NVIDIA GPU Operator(<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>)는 K8s 환경에서 GPU 드라이버·CUDA·MIG·DCGM을 자동 설치·관리 하며, Kubeflow, KServe, Ray 등 AI 워크로드 프레임워크가 K8s Operator·Helm Chart 형태로 배포됩니다. Vector DB(Milvus, Weaviate, Qdrant 등) 역시 Helm Chart와 Operator로 K8s에 배포되어, StatefulSet·CSI 스토리지·백업 전략과 결합됩니다. MCP(Model Context Protocol) 서버와 AI Agent 런타임도 API Gateway·Service Mesh와 연동되어 K8s 네이티브로 배포되는 것이 표준입니다.

세계적 AI 인프라 로드맵에서는 다기관 API, 공유 Vector DB, 공통 MCP 서버가 K8s 기반으로 통합되는 패턴이 확산되고 있습니다. 예를 들어, 여러 조직이 동일한 Vector DB(예: Qdrant)를 K8s 클러스터에서 운영하며, MCP 서버와 Agent 런타임을 컨테이너로 배포해 자동

확장(HPA·KEDA)과 GPU 공유(MIG·Time-slicing)를 활용합니다. 이 구조는 AI 워크로드의 선언형·불변·자동확장 특성을 극대화하며, VM 기반 운영에서는 불가능한 수준의 탄력성과 경제성을 제공합니다. 국내 공공·금융의 AI 전환 역시 이 글로벌 표준을 따르는 것이 상호운용성·확장성 측면에서 유리합니다.

AI 인프라의 요구사항은 크게 모델 서빙, 대용량 데이터 저장 및 검색(Vector DB), RAG(검색 증강 생성) 프레임워크, API Gateway, MCP 서버 및 에이전트 관리 등으로 나눌 수 있습니다. 각 항목 별로 K8s 기반의 오픈소스 솔루션과 글로벌 레퍼런스가 존재하며, 실제로 Booking.com, Spotify, Adidas, CERN 등은 이러한 구조를 적극적으로 도입하고 있습니다. 예를 들어, Booking.com은 KServe와 vLLM을 활용해 대규모 모델 서빙을 자동화하였고, Adidas는 Milvus Operator 기반의 Vector DB를 통해 AI 서비스의 확장성과 신뢰성을 확보하였습니다. 또한, 글로벌 금융기관들은 LangChain, LlamaIndex와 같은 RAG 프레임워크를 K8s CRD로 배포하여, 데이터 보안과 확장성을 동시에 달성하고 있습니다. 이러한 사례들은 AI 인프라 요구사항이 K8s 기반에서 어떻게 효과적으로 충족되는지 보여주는 대표적인 예시입니다.

AI 인프라 요구 항목	K8s 대응 기능	글로벌 레퍼런스 사례
모델 서빙	KServe·vLLM·Triton Operator	Booking.com·Spotify
Vector DB	Milvus·Weaviate·Qdrant Operator	Adidas·CERN
RAG 프레임워크	LangChain·LlamaIndex CRD	글로벌 금융기관
API Gateway	Istio·NGINX Ingress	KubeCon 발표 사례
MCP 서버·Agent	Deployment·Service Mesh	글로벌 AI 연구기관

이 표는 AI 인프라의 주요 요구사항과 이에 대응하는 K8s 기능, 그리고 실제 글로벌 적용 사례를 정리한 것입니다. 이를 통해 AI·데이터 플랫폼의 표준화 흐름이 왜 Kubernetes 중심으로 수렴하고 있는지, 그리고 각 기술 요소가 어떻게 상호 연계되는지 명확히 이해할 수 있습니다.

10.2 글로벌 공공·엔터프라이즈 사례와 국내 확산 동향

클라우드 네이티브 전환이 세계적 기술 표준으로 자리잡은 현실은 글로벌 공공·엔터프라이즈 사례와 국내 금융·공공 확산 동향에서 명확히 드러납니다. 최근 몇 년간 글로벌 대형 기관과 기업들은 기존 VM 기반 인프라에서 벗어나, K8s 기반의 선언형·불변·자동확장 운영 모델로 대규모 전환을 이루

고 있습니다. 이러한 변화는 단순한 기술적 유행이 아니라, 실제 운영 효율성과 비즈니스 민첩성, 그리고 확장성 측면에서 명확한 이점을 제공하기 때문입니다. 본 절에서는 CERN, Booking.com, Spotify, Adidas 등 세계적 사례와 더불어, 국내 금융·공공 부문에서의 클라우드 네이티브 확산 동향과 MSAP.ai의 채택 흐름을 비교 분석하여, 클라우드 네이티브 전환의 현실성과 그 의미를 구체적으로 확인합니다.

10.2.1 CERN·Spotify·Booking.com·Adidas 등 글로벌 공공·엔터프라이즈 클라우드 네이티브 전환 사례

CERN은 세계 최대의 과학 연구기관으로, LHC(대형 하드론 충돌기) 데이터 처리에 K8s 기반 대규모 인프라를 활용합니다. OpenStack과 K8s 혼합 운영을 통해 VM과 컨테이너 워크로드를 병행하며, 수천 노드 규모의 클러스터를 운영합니다. 이 사례는 과학 연산 분야에서도 VM 기반에서 K8s 기반으로의 전환이 현실적임을 보여줍니다.

Booking.com은 KubeCon EU 2022에서 베어메탈 K8s 클러스터를 수천 노드 규모로 운영하는 사례를 발표했습니다. VM 없이 직접 컨테이너 기반 워크로드를 배포하며, 자동확장·불변 인프라·관측성 스택을 완비하여 운영 효율과 배포 속도를 크게 향상시켰습니다.

Spotify는 VM 기반 모놀리식 운영에서 마이크로서비스·K8s 기반 자동확장으로 전환하여 신규 서비스 프로비저닝 시간을 1시간에서 수 분으로 단축하였습니다. 이 과정에서 Service Mesh, GitOps, Operator 패턴을 적극 도입하여 운영 모델을 완전히 현대화했습니다.

Adidas는 CNCF 공식 레퍼런스 사례로, 2019년부터 4000+ 파드와 200+ 서비스 운영을 통해 배포 속도를 3~4배 개선하였습니다. K8s 기반의 선언형·불변·자동확장 운영 모델을 통해 장애 복구 시간(MTTR)도 50% 이상 단축하였습니다.

미국 대형 통신사들은 OpenStack Helm과 K8s 조합으로 핵심 네트워크 워크로드를 이전하였으며, OpenStack 워크로드조차 K8s 운영 모델로 수렴하고 있습니다. 이는 VM 기반 IaaS만을 고집하는 조직이 점차 줄고 있음을 보여줍니다.

이러한 글로벌 사례들은 단순히 기술 도입을 넘어, 조직의 비즈니스 민첩성과 운영 효율성을 극대화하는 방향으로 클라우드 네이티브 전환이 이루어지고 있음을 시사합니다. CERN의 경우, 대규모 과학 데이터 처리에서 K8s의 자동확장성과 자원 효율성을 극대화하였고, Booking.com은 베어메탈 환경에서 K8s를 활용해 VM 없이도 대규모 서비스를 안정적으로 운영하고 있습니다.

Spotify는 마이크로서비스 구조와 K8s의 결합을 통해 신속한 서비스 출시와 운영 자동화를 실현하였으며, Adidas는 GitOps와 K8s를 결합해 배포 속도와 장애 대응력을 크게 높였습니다. 미국 대형 통신사들은 네트워크 워크로드까지도 K8s 기반으로 이전함으로써, 네트워크 인프라의 유연성과 확장성을 확보하고 있습니다.

이처럼 다양한 산업 분야에서 K8s 기반 클라우드 네이티브 전환이 이루어지고 있으며, 이는 글로벌 표준으로 자리잡고 있습니다. 각 조직의 전환 범위와 플랫폼 선택은 다르지만, 공통적으로 선언형·불변·자동확장 운영 모델을 도입함으로써 운영 효율성과 비즈니스 경쟁력을 크게 강화하고 있습니다.

조직	성격	전환 범위	플랫폼
CERN	과학 연구기관	LHC 데이터 처리	OpenStack + K8s
Booking.com	글로벌 서비스	전체 워크로드	베어메탈 K8s
Spotify	음악 서비스	마이크로서비스 전체	K8s + Service Mesh
Adidas	제조/유통	4000+ 파드, 200+ 서비스	K8s + GitOps
미국 대형 통신사	네트워크	핵심 네트워크 워크로드	OpenStack Helm + K8s

이 표는 각 조직의 성격, 전환 범위, 그리고 도입한 플랫폼을 정리한 것으로, 클라우드 네이티브 전환이 특정 산업에 국한되지 않고 다양한 분야에서 폭넓게 확산되고 있음을 보여줍니다.

10.2.2 국내 금융·공공의 클라우드 네이티브 확산 동향과 MSAP.ai 채택 흐름

국내에서도 금융·공공 부문을 중심으로 클라우드 네이티브 전환이 빠르게 확산되고 있습니다. 특히, 금융권은 보안과 규제 준수라는 엄격한 요건에도 불구하고, K8s 기반 운영의 효율성과 확장성, 그리고 자동화의 이점을 인식하여 점진적으로 전환을 추진하고 있습니다. 본 절에서는 국내 금융권의 K8s 기반 전환 사례, MSAP.ai의 국산 PaaS 기준 채택 확대, 그리고 금융 규제 환경에서의 K8s 기반 운영 가능성 입증 사례를 구체적으로 살펴보겠습니다.

국내 대형 은행들은 해외 상용 엔터프라이즈 K8s 디스트리뷰션 A를 채택하여 디지털 전환을 추진하고 있습니다. 기존 금융기관도 VM 기반 운영에서 컨테이너·K8s 기반으로 점진적 전환을 진행 중이며, 신규 디지털 금융사는 출범 초기부터 K8s 기반 운영을 선택하고 있습니다. 이 과정에서 컴플라이언스·규제 대응이 필수적으로 요구되지만, K8s 기반 운영이 충분히 가능성이 입증되고

있습니다.

MSAP.ai(투라인클라우드)는 국내 공공·금융에서 채택을 확대하고 있는 국산 PaaS 기준 제품입니다. CNCF Certified Kubernetes Conformance Program 등재를 통해 표준 K8s API·기능 호환성을 공식 검증받았으며, 한글화 UI/문서, 국산 인증, 국내 지원 체계를 갖추고 있습니다. 국내 규제 환경에 맞춘 운영 모델과 AI 워크로드 정합성, 총소유비용(TCO) 측면에서 해외 상용 디스트리뷰션 대비 우위를 갖고 있습니다.

금융 규제 준수가 K8s 도입의 장애가 아니라는 사실은 국내·글로벌 사례로 이중 입증됩니다. 국내 금융권의 K8s 기반 운영 사례는 규제 환경에서도 선언형·불변·자동확장 기반 클라우드 네이티브 전환이 충분히 가능함을 보여주며, MSAP.ai는 국산 PaaS 기준 옵션으로 자리잡고 있습니다.

실제로, 국내 대형 은행은 2021년부터 전체 워크로드를 해외 상용 K8s 디스트리뷰션 A로 전환하기 시작했으며, 기존 금융기관들은 2022년부터 MSAP.ai를 활용해 점진적으로 컨테이너 전환을 추진하고 있습니다. 신규 디지털 금융사는 2023년 출범과 동시에 K8s 업스트림을 기반으로 IT 인프라를 설계하여, 초기부터 클라우드 네이티브 운영 모델을 도입하였습니다. 이러한 흐름은 국내 금융권이 글로벌 표준을 적극적으로 수용하면서도, 국내 규제와 환경에 맞는 맞춤형 솔루션을 선택하고 있음을 보여줍니다.

특히, MSAP.ai는 한글화된 UI와 문서, 국내 인증 및 지원 체계, 그리고 AI 워크로드와의 정합성 등에서 차별화된 경쟁력을 갖추고 있습니다. 이를 통해 국내 금융·공공기관은 글로벌 수준의 클라우드 네이티브 인프라를 구축함과 동시에, 국내 환경에 최적화된 운영과 지원을 받을 수 있습니다. 또한, K8s 기반 운영의 자동화와 확장성, 그리고 불변 인프라의 도입으로 인해, 금융 규제 환경에서도 안정적이고 효율적인 서비스 제공이 가능함이 입증되고 있습니다.

조직	플랫폼	전환 범위	시작 시점
대형 은행	해외 디스트리뷰션 A	전체 워크로드	2021년
기존 금융기관	MSAP.ai	점진적 컨테이너 전환	2022년
신규 금융사	K8s 업스트림	초기부터 K8s 기반	2023년

이 표는 국내 금융권의 주요 조직별로 도입한 플랫폼, 전환 범위, 그리고 시작 시점을 정리한 것으로, 국내에서도 클라우드 네이티브 전환이 본격적으로 확산되고 있음을 보여줍니다. 앞으로도 이러한 흐름은 더욱 가속화될 것으로 예상되며, MSAP.ai와 같은 국산 PaaS의 역할이 더욱 중요

해질 것입니다.

11장: VM이 여전히 유효한 예외 케이스와 KubeVirt 흡수 전략

11.1 VM을 유지해야 하는 기술·규제 조건

VM 환경에서 컨테이너로의 전환이 조직의 목표 아키텍처로 자리잡고 있지만, 모든 워크로드가 컨테이너화에 적합한 것은 아니다. 실제로 기술적 제약이나 규제 요건으로 인해 VM을 유지해야 하는 예외 케이스가 존재하며, 이 비율은 전체 워크로드의 10~30%에 달할 수 있다. 이러한 예외 케이스는 단순히 기술적 한계에서 비롯되는 것이 아니라, 금융·공공 등 규제 산업에서 요구하는 하드 멀티테넌시, 특수 커널 및 드라이버 요구, 실시간 OS, 그리고 Windows 워크로드와 같은 운영 환경의 제약에서 발생한다. 이 절에서는 VM을 유지해야 하는 주요 조건과, 컨테이너화가 불가능하거나 비효율적인 워크로드의 판단 기준을 구체적으로 제시한다. 또한, 하드 멀티테넌시 규제에 대응하기 위한 Kata Containers와 gVisor 같은 대체 격리 옵션을 소개함으로써, K8s 기반 환경에서 관리 포인트 이중화를 최소화하는 전략을 설명한다.

11.1.1 상이 커널 요구·레거시 드라이버·실시간 OS·Windows 워크로드

특정 워크로드가 컨테이너 환경에 적합하지 않은 대표적인 이유는 커널 버전, 드라이버, 운영체제(OS) 종속성, 그리고 실시간 처리나 Windows 환경과 같은 특수 요구 때문이다. 이 섹션에서는 각각의 기술적 제약이 실제로 어떻게 VM 환경을 필요로 하게 만드는지, 그리고 이를 판단할 수 있는 구체적인 기준을 제시합니다. 또한, 실제 현업에서 자주 마주치는 사례와 함께, 컨테이너로의 전환이 어려운 워크로드의 특성을 분석하여 조직이 합리적으로 VM을 유지할 필요가 있는 상황을 명확히 구분할 수 있도록 안내합니다.

커널 버전 상이 요구

특정 워크로드는 커널 버전이나 기능에 대한 강한 종속성을 가진다. 예를 들어, 실시간 처리 시스템이나 특수 하드웨어 연동이 필요한 환경에서는 커널 모듈이나 특정 버전의 커널이 필수적이다. 컨테이너는 호스트 커널을 공유하기 때문에, 이러한 요구를 충족시키기 어렵다. VM은

하이퍼바이저 위에 독립된 Guest OS를 올려 커널 버전과 설정을 자유롭게 맞출 수 있으므로, 커널 상이 요구가 있는 워크로드에 적합하다. 판단 기준은 워크로드가 커널 모듈 설치, 특정 커널 패치, 또는 실시간 커널 기능을 반드시 요구하는지 여부다.

실제로, 제조업의 공장 자동화 시스템이나 금융권의 초저지연 트랜잭션 처리 시스템 등에서는 커널의 특정 기능이나 패치가 필수적입니다. 컨테이너 환경에서는 이러한 커널 패치나 모듈을 개별적으로 적용할 수 없으므로, VM을 통해 별도의 커널 환경을 구성해야만 합니다. 또한, 일부 보안 솔루션이나 네트워크 장비 연동 소프트웨어는 커널의 특정 버전에만 호환되는 경우가 많아, 이러한 요구가 있는 경우 VM이 사실상 유일한 선택지가 됩니다.

레거시 드라이버 종속

ERP, 제조, 금융 등에서 사용되는 레거시 소프트웨어는 특정 드라이버, 커널 모듈, 또는 OS 버전에 강하게 의존한다. 컨테이너 환경에서는 호스트 커널의 드라이버만 사용 가능하므로, Guest OS에서만 지원되는 드라이버가 필요한 경우 VM이 필수적이다. 특히 하드웨어 연동(PCI 장치, 특수 네트워크 카드 등)이나 폐쇄형 소프트웨어가 드라이버를 직접 설치하는 경우, VM 환경이 유일한 선택지가 된다. 판단 기준은 드라이버가 호스트 커널에 직접 설치될 수 없는지, 또는 컨테이너에서 접근이 불가능한지 여부다.

예를 들어, 오래된 금융 코어 시스템이나 제조 현장의 장비 제어 소프트웨어는 특정 하드웨어와 연동되는 드라이버가 반드시 필요합니다. 이러한 드라이버는 최신 커널이나 컨테이너 환경에서 지원되지 않거나, 커널 모듈을 직접 설치해야만 동작하는 경우가 많습니다. 실제로, 일부 특수 네트워크 카드나 스토리지 컨트롤러는 공급사가 제공하는 드라이버가 특정 리눅스 배포판과 커널 버전에만 호환되어, 컨테이너 환경에서는 사용할 수 없습니다. 이처럼 레거시 드라이버 종속성이 강한 워크로드는 VM 환경에서만 안정적으로 운영이 가능합니다.

실시간 OS 및 특수 워크로드

실시간 OS(Real-Time Operating System)는 엄격한 시간 제약과 커널 튜닝이 요구된다. 컨테이너는 일반적인 Linux 커널을 공유하므로, 실시간 OS의 요구를 충족시키기 어렵다. VM은 실시간 OS 이미지를 Guest로 올려 커널 레벨의 제어와 튜닝을 가능하게 한다. 판단 기준은 워크로드가 실시간 스케줄링, 커널 튜닝, 또는 특수 OS 기능을 요구하는지 확인하는 것이다.

실시간 제어가 필요한 산업용 로봇, 의료기기, 통신 장비 등에서는 마이크로초 단위의 응답성과 커널 레벨의 스케줄링 제어가 필수적입니다. 이러한 요구는 컨테이너 환경에서 제공하기 어렵기 때문에, 실시간 커널이 적용된 OS를 VM에 올려 운영하는 것이 일반적입니다. 또한, 일부 특수

목적의 OS(예: QNX, RTEMS 등)는 컨테이너 환경에서 지원되지 않으므로, VM을 통한 격리와 커널 제어가 반드시 필요합니다.

Windows 워크로드

Windows 컨테이너는 호스트와 컨테이너의 커널 버전이 일치해야 하며, Linux 컨테이너만큼 성숙하지 못했다. Windows 기반 ERP, Active Directory, 레거시 업무 시스템은 VM 환경에서 독립적으로 운영하는 것이 현실적이다. 판단 기준은 해당 워크로드가 Windows OS에 종속적이고, 컨테이너화 시 커널 버전 일치·호환성 문제가 발생하는지 여부다.

실제로, 국내외 많은 기업에서는 여전히 Windows 기반의 ERP, 그룹웨어, 인증 시스템(Active Directory) 등을 운영하고 있습니다. Windows 컨테이너는 기능적으로 제한적이고, 커널 버전 호환성 문제로 인해 다양한 Windows 애플리케이션을 안정적으로 운영하기 어렵습니다. 또한, 일부 레거시 소프트웨어는 Windows Server의 특정 버전에서만 동작하며, 컨테이너화가 기술적으로 불가능한 경우가 많습니다. 이러한 환경에서는 VM을 통해 Windows OS를 독립적으로 운영하는 것이 가장 현실적인 선택입니다.

ERP 등 OS 종속 워크로드

ERP, 금융 코어 시스템 등은 OS 버전, 패치, 보안 정책에 강하게 의존한다. 컨테이너로의 전환이 기술적으로 불가능하거나, 내부통제·감사 요건 등 규제적 이유로 VM 환경이 요구될 수 있다. 판단 기준은 OS 종속성이 기술·규제적으로 해소될 수 있는지, 아니면 VM 환경이 불가피한지에 대한 정책적 검토다.

ERP 시스템이나 금융권의 핵심 업무 시스템은 특정 OS 버전, 패치 레벨, 보안 정책에 따라 엄격하게 관리됩니다. 예를 들어, SAP ERP는 특정 리눅스 배포판과 커널 버전에서만 공식 지원을 제공하며, 보안·감사 요건으로 인해 임의의 커널 변경이나 컨테이너화가 허용되지 않는 경우가 많습니다. 또한, 내부통제와 외부감사 기준에 따라 OS 환경의 변경이 제한되거나, 시스템 변경 이력이 엄격히 관리되어야 하는 경우 VM 환경이 필수적입니다. 이처럼 OS 종속성이 강한 워크로드는 정책적·기술적 검토를 거쳐 VM 환경을 유지하는 것이 바람직합니다.

11.1.2 하드 멀티테넌시 규제와 Kata Containers·gVisor 대체 조건

컨테이너 환경의 확산에도 불구하고, 금융·공공 등 규제 산업에서는 여전히 하드 멀티테넌시와 같은 강력한 격리 요건이 요구되는 경우가 많습니다. 이 섹션에서는 하드 멀티테넌시 규제의 구

체적 요구사항과, 이를 충족시키기 위한 대체 격리 기술인 Kata Containers, gVisor의 구조와 적용 시나리오를 상세히 설명합니다. 또한, 이러한 기술들이 K8s 환경에서 어떻게 통합 관리될 수 있는지, 그리고 오버헤드와 적합 시나리오에 대한 실질적인 판단 기준을 제시하여, 조직이 규제 대응과 운영 효율성 사이에서 최적의 선택을 할 수 있도록 안내합니다.

하드 멀티테넌시 규제 요건

금융·공공 등에서는 하이퍼바이저 수준의 격리(하드 멀티테넌시)를 요구하는 경우가 많다. 컨테이너는 기본적으로 커널을 공유하므로, 커널 취약점이 발생하면 전체 테넌트에 영향을 줄 수 있다. VM은 하이퍼바이저를 통해 OS 단위 격리를 제공하므로, 규제 심사에서 요구하는 격리 수준을 만족한다. 판단 기준은 내부통제·감사·보안 심사에서 하이퍼바이저 격리 요구가 명시되어 있는지 여부다.

실제로, 금융감독원이나 정부기관의 보안 심사에서는 하이퍼바이저 기반의 격리를 명확히 요구하는 경우가 있습니다. 예를 들어, 여러 고객사의 데이터를 동시에 처리하는 SaaS 환경이나, 민감한 개인정보를 다루는 공공 시스템에서는 커널 취약점으로 인한 테넌트 간 침해 가능성을 최소화해야 하므로, VM 기반의 하드 멀티테넌시가 필수적입니다. 이러한 규제 요건은 컨테이너 환경만으로는 충족하기 어려우며, VM 또는 하이퍼바이저 기반의 대체 기술이 필요합니다.

Kata Containers의 격리 구조

Kata Containers는 컨테이너를 경량 VM으로 실행하여 하이퍼바이저 수준의 격리를 제공한다. 이는 K8s 환경에서 VM과 컨테이너의 장점을 결합한 방식으로, 보안·규제 요구를 충족시키면서도 관리 포인트를 K8s 표면으로 통합할 수 있다. 다만, 베어메탈 컨테이너 대비 5~15%의 CPU·메모리 오버헤드가 발생한다는 공식 벤치마크가 있다. 적합 시나리오는 하드 멀티테넌시가 필수적이고, 약간의 성능 저하를 감내할 수 있는 환경이다.

Kata Containers는 각 컨테이너를 독립적인 경량 VM으로 실행함으로써, 컨테이너 간 커널 공유로 인한 보안 취약점을 원천적으로 차단합니다. 예를 들어, 금융권의 멀티테넌트 SaaS 서비스나, 민감 데이터 처리가 필요한 환경에서는 Kata Containers를 통해 하이퍼바이저 수준의 격리를 제공하면서도, K8s의 자동화·확장성·관찰 가능성 등 클라우드 네이티브의 이점을 함께 누릴 수 있습니다. 단, 오버헤드가 존재하므로, 성능 요구가 극단적으로 높은 워크로드에는 신중한 검토가 필요합니다.

gVisor의 커널 격리

gVisor는 사용자 공간에서 커널 API를 재구현하여, 컨테이너마다 별도의 커널 경계를 제공

한다. 일반 컨테이너 대비 보안 격리가 강화되지만, I/O·시스템 콜 성능이 다소 저하된다. 적합 시나리오는 보안 격리와 성능 사이의 트레이드오프를 허용할 수 있는 환경이다.

gVisor는 Google에서 개발한 샌드박스 런타임으로, 각 컨테이너가 호스트 커널과 직접 상호작용하지 않고, 사용자 공간에서 재구현된 커널 API를 통해 간접적으로 자원에 접근하도록 설계되어 있습니다. 이를 통해 커널 취약점으로 인한 보안 사고를 예방할 수 있으며, 다중 테넌트 환경에서 보안 격리를 강화할 수 있습니다. 다만, 시스템 콜이나 I/O 작업의 성능이 일부 저하될 수 있으므로, 보안이 최우선인 환경이나, 약간의 성능 저하를 감내할 수 있는 워크로드에 적합합니다.

K8s 표면에서의 규제 대응

Kata Containers와 gVisor는 K8s CRD 또는 RuntimeClass를 통해 통합 관리가 가능하다. 이를 통해 VM과 컨테이너의 관리 포인트 이중화를 피하고, 규제 대응을 K8s 표면에서 일원화할 수 있다. 판단 기준은 해당 워크로드가 하드 멀티테넌시를 요구하면서도 K8s 기반 관리가 가능한지 여부다.

Kubernetes에서는 RuntimeClass를 활용해 각 워크로드별로 런타임을 지정할 수 있습니다. 예를 들어, 보안이 중요한 워크로드에는 Kata Containers나 gVisor 런타임을 적용하고, 일반 워크로드에는 표준 컨테이너 런타임을 사용하는 식으로 유연하게 운영이 가능합니다. 이를 통해 운영팀은 단일 K8s API와 도구 체계에서 모든 워크로드를 관리할 수 있으며, 규제 대응과 운영 효율성을 동시에 확보할 수 있습니다.

오버헤드와 적합 시나리오

Kata Containers는 베어메탈 컨테이너 대비 5~15% 오버헤드가 있지만, 하이퍼바이저 수준 격리가 필수적인 경우에는 감수할 수 있는 범위다. gVisor는 보안 격리와 성능 저하 사이의 균형이 필요한 시나리오에 적합하다. VM 환경이 불가피한 경우 외에는, K8s 표면에서의 격리 대체 옵션을 적극 검토할 필요가 있다.

실제로, 대형 금융기관이나 클라우드 서비스 제공업체에서는 하드 멀티테넌시가 요구되는 워크로드에 Kata Containers를 적용하여, 보안과 규제 요건을 충족시키면서도 K8s 기반의 자동화와 확장성을 활용하고 있습니다. gVisor는 상대적으로 경량화된 보안 격리가 필요한 환경, 예를 들어 개발·테스트 환경이나, 보안이 중요한 일부 서비스에 적용할 수 있습니다. 이처럼 오버헤드와 보안 요구 수준을 종합적으로 고려하여, 조직의 정책과 워크로드 특성에 맞는 격리 옵션을 선택하는 것이 중요합니다.

11.2 KubeVirt로 잔존 VM을 K8s 오브젝트로 흡수하는 아키텍처

컨테이너화가 불가능한 워크로드 또는 규제·기술적 예외 케이스에 대해서는, VM을 완전히 폐지하는 대신 KubeVirt를 활용하여 K8s 오브젝트로 흡수하는 전략이 현실적이다. KubeVirt는 VM을 Kubernetes Custom Resource Definition(CRD)로 관리하며, Pod와 VM을 동일 API·파이프라인에서 선언형으로 운영할 수 있게 한다. 이를 통해 VM의 관리 포인트를 K8s 표면으로 통합하고, GitOps·관찰 가능성·보안 정책 등 클라우드 네이티브의 장점을 VM에도 적용할 수 있다. 이 절에서는 KubeVirt의 성능 오버헤드 벤치마크와 적합 워크로드 판단 기준, 그리고 Pod와 VM을 동일 API·GitOps 파이프라인에서 관리하는 패턴을 기술적으로 설명한다.

11.2.1 KubeVirt IOPS·네트워크 오버헤드 벤치마크와 적합 워크로드 판단

KubeVirt를 도입할 때 가장 중요한 고려사항 중 하나는 VM 성능 오버헤드와 실제로 어떤 워크로드에 적합한지에 대한 명확한 판단입니다. 이 섹션에서는 KubeVirt의 VM 관리 구조와 공식 벤치마크 결과를 바탕으로, IOPS(초당 입출력 작업 수) 및 네트워크 지연에 대한 실질적인 수치를 제시합니다. 또한, KubeVirt가 적합한 워크로드와 부적합한 워크로드의 구체적인 예시를 들어, 조직이 전환 시 예상할 수 있는 성능 및 운영상의 변화를 현실적으로 평가할 수 있도록 안내합니다. 마지막으로, KVM과 KubeVirt의 비교 표를 통해 관리 방식과 성능 차이를 한눈에 파악할 수 있도록 정리합니다.

KubeVirt의 VM 관리 구조

KubeVirt는 VM을 Kubernetes CRD(VirtualMachine)로 정의하여, Pod와 동일한 오브젝트 관리 체계에 통합한다. VM의 생성, 삭제, 스케일, 네트워크, 스토리지 등 모든 운영을 K8s API와 YAML 선언형 방식으로 처리할 수 있다. 이는 VM 관리의 자동화와 GitOps 파이프라인 통합을 가능하게 한다.

KubeVirt의 가장 큰 특징은 VM을 K8s의 리소스 오브젝트로 취급함으로써, 기존 VM 관리의 복잡성을 크게 줄이고, 선언형 인프라 관리의 이점을 VM에도 적용할 수 있다는 점입니다. 예를 들어, VM의 라이프사이클(생성, 삭제, 스케일링 등)을 K8s의 컨트롤러와 오토스케일러로 자동화할 수 있으며, 네트워크와 스토리지 역시 K8s의 CNI, CSI 표준을 그대로 활용할 수 있습니다. 이를 통해 운영팀은 VM과 컨테이너를 동일한 도구와 프로세스로 관리할 수 있게 됩니다.

IOPS·네트워크 오버헤드 공식 벤치마크

Red Hat 공식 벤치마크에 따르면, KubeVirt에서 실행되는 VM의 IOPS는 네이티브 KVM 대비 약 3% 내의 오버헤드가 발생하며, 네트워크 지연은 약 5% 내에서 추가된다. 이는 대부분의 엔터프라이즈 워크로드에서 허용 가능한 수준이며, 극한의 성능 요구가 없는 한 KubeVirt를 통한 VM 흡수가 현실적이다. 예를 들어, ERP, 금융 코어, 특수 커널 요구 워크로드 등은 KubeVirt 환경에서 안정적으로 운영 가능하다.

실제 사례로, 한 대형 금융기관에서는 기존 KVM 기반 VM을 KubeVirt 환경으로 이전하면서, ERP 및 코어 बैं킹 시스템의 트랜잭션 처리량이 2~3% 내에서만 감소하는 것을 확인하였으며, 네트워크 지연 역시 5% 이내로 관리되었습니다. 이는 실시간 트레이딩이나 HPC(고성능 컴퓨팅)와 같이 극한의 성능이 요구되는 워크로드를 제외하면, 대부분의 업무 시스템에서 충분히 수용 가능한 수준임을 의미합니다. 또한, KubeVirt 환경에서는 VM의 자동화, 롤링 업데이트, 선언형 배포 등 클라우드 네이티브의 이점을 함께 누릴 수 있습니다.

적합 워크로드와 부적합 워크로드

적합 워크로드는 다음과 같다:

- 커널·드라이버 요구, 실시간 OS, Windows 워크로드 등 컨테이너화가 불가능한 워크로드
- 하드 멀티테넌시 규제 대응이 필요한 워크로드
- ERP 등 OS 종속성이 강한 시스템

부적합 워크로드는 극한의 지연·IOPS 요구(예: 초고속 트레이딩, HPC 등)나, KubeVirt 오버헤드가 허용되지 않는 환경이다.

적합한 워크로드의 실제 예로는, 제조업의 장비 제어 시스템, 금융권의 ERP 및 코어뱅킹, 공공기관의 민감 데이터 처리 시스템 등이 있습니다. 이들은 OS, 커널, 드라이버 종속성이 강하고, 규제 요건으로 인해 컨테이너화가 불가능하지만, KubeVirt를 통해 K8s 기반의 자동화와 통합 관리를 실현할 수 있습니다. 반면, 초저지연이 필수적인 고빈도 트레이딩 시스템이나, 대규모 병렬 연산이 필요한 HPC 워크로드는 KubeVirt의 오버헤드가 성능에 치명적일 수 있으므로, 별도의 KVM 환경을 유지하는 것이 바람직합니다.

실행 전략과 현실적 전환 비용

잔존 VM 수에 KubeVirt 벤치마크를 적용하여, 전환 후 성능·운영 비용·관리 포인트 변화 등을 산정할 수 있다. KubeVirt Early Majority 단계에 진입한 만큼, 대규모 엔터프라이즈에서도 실무

적용이 가능하며, VM을 버리지 않고 K8s로 흡수하는 현실적 실행 전략이 된다.

전환 전략 수립 시에는, 각 VM 워크로드의 성능 요구사항을 KubeVirt 공식 벤치마크와 비교하여, 예상되는 오버헤드와 운영 효율성 개선 효과를 함께 평가해야 합니다. 예를 들어, 100대의 VM 중 80대가 KubeVirt로 무리 없이 이전 가능하다면, 관리 포인트 통합과 자동화로 인한 운영 비용 절감 효과가 오버헤드로 인한 성능 저하보다 훨씬 클 수 있습니다. 또한, KubeVirt는 이미 글로벌 대기업과 금융기관에서 실무 적용이 활발히 이루어지고 있어, 기술적 안정성과 지원 생태계도 충분히 확보되어 있습니다.

KVM vs KubeVirt 비교 표

관리 단위	IOPS 오버헤드	네트워크 지연	운영 방식
KVM	0%	0%	하이퍼바이저 API
KubeVirt	~3% 내	~5% 내	K8s CRD/YAML

이 표를 통해 KubeVirt와 KVM의 성능 및 관리 방식 차이를 한눈에 비교할 수 있습니다. KubeVirt는 약간의 오버헤드가 있지만, K8s 기반의 선언형 관리와 자동화, 통합 관찰 가능성 등에서 큰 이점을 제공합니다.

11.2.2 Pod와 VM을 동일 API·동일 GitOps 파이프라인에서 관리하는 패턴

KubeVirt를 활용하면, Pod와 VM을 동일한 Kubernetes API와 GitOps 파이프라인에서 일관되게 관리할 수 있습니다. 이 섹션에서는 KubeVirt의 CRD 구조와 선언형 YAML 관리 방식, 그리고 GitOps 도구(Argo CD, Flux 등)를 통한 통합 배포·운영 패턴을 구체적으로 설명합니다. 또한, 관찰 가능성 및 보안 스택의 통합 적용 사례와, 단일 표면 수렴이 조직의 운영 효율성에 미치는 긍정적 효과를 실제 현업 사례와 함께 분석합니다. 마지막으로, Pod와 VirtualMachine CRD의 기능적 차이를 표로 정리하여, 운영팀이 실질적으로 어떤 이점을 누릴 수 있는지 명확히 제시합니다.

KubeVirt CRD와 선언형 YAML 관리

KubeVirt는 VM을 VirtualMachine CRD로 정의하고, YAML 파일을 통해 선언형으로 관리한다. Pod와 동일하게 `kubectl apply`로 배포·수정·삭제가 가능하며, GitOps(Argo CD, Flux) 파이프라인에서 VM과 Pod를 동일하게 관리할 수 있다. 이를 통해 운영팀은 단일 표면에서 모든 워크로드를 관리하며, 배포·롤백·업그레이드·감사 추적이 일관된 방식으로 이루어진다.

KubeVirt의 VirtualMachine CRD는 Pod와 거의 동일한 방식으로 정의되며, VM의 스펙(이미지, CPU, 메모리, 네트워크, 스토리지 등)을 YAML로 선언할 수 있습니다. 운영팀은 Git 저장소에 VM 정의 파일을 관리하고, 변경 사항을 PR(풀 리퀘스트)로 검토·승인한 뒤, GitOps 도구를 통해 자동 배포할 수 있습니다. 이 과정에서 배포 이력, 롤백, 변경 감사 등 모든 관리가 코드 기반으로 일관되게 이루어집니다.

동일 GitOps 파이프라인 적용

Argo CD, Flux 등 GitOps 도구를 활용하면, VM과 컨테이너 워크로드를 동일한 Git 저장소와 파이프라인에서 관리할 수 있다. 변경 관리, 롤백, 감사 로그, 배포 리드타임 단축 등 클라우드 네이티브의 이점을 VM에도 적용할 수 있다. VM의 선언형 배포는 기존의 명령형 VM 템플릿·스냅샷 방식보다 운영 품질과 자동화 수준이 높다.

실제 현업에서는, 기존 VM 환경에서 수동으로 템플릿을 복제하거나, 스냅샷을 관리하던 방식에서 벗어나, GitOps 기반의 선언형 배포로 전환함으로써 배포 오류와 인적 실수를 크게 줄이고 있습니다. 예를 들어, 신규 VM 배포 시 YAML 파일만 수정하면 자동으로 K8s 클러스터에 적용되며, 롤백도 Git 커밋 단위로 손쉽게 수행할 수 있습니다. 이로 인해 배포 리드타임이 단축되고, 운영 효율성이 크게 향상됩니다.

관찰·보안 스택의 통합 적용

Prometheus, OpenTelemetry, Falco, Trivy, Cosign 등 관찰 가능성 및 보안 스택을 Pod와 VM 모두에 동일하게 적용할 수 있다. VM의 상태, 로그, 메트릭, 보안 이벤트를 K8s 표면에서 통합 모니터링하며, 정책 엔진(OPA/Gatekeeper)도 동일하게 적용 가능하다.

예를 들어, Prometheus를 활용해 VM과 Pod의 메트릭을 통합 수집·분석할 수 있으며, OpenTelemetry로 분산 추적을 구현할 수 있습니다. Falco, Trivy, Cosign 등 보안 도구 역시 VM과 컨테이너 모두에 적용되어, 취약점 진단, 런타임 보안 이벤트 탐지, 이미지 서명 검증 등을 일관되게 수행할 수 있습니다. 또한, OPA/Gatekeeper 정책 엔진을 통해 VM과 Pod의 배포 정책, 보안 기준을 통합적으로 관리할 수 있습니다.

단일 표면 수렴의 기술적 완결

Pod와 VM을 동일 API·GitOps 파이프라인에서 관리하는 구조는 관리 포인트 이중화(9장)를 기술적으로 완결한다. 운영팀의 FTE·도구·정책 중복이 해소되고, 배포·감사·보안·관찰 등 모든 운영이 단일 표면에서 일관되게 이루어진다. 이는 K8s 기반 아키텍처의 궁극적 목표 상태다.

실제로, 대형 엔터프라이즈에서는 VM과 컨테이너를 각각 별도의 관리 체계로 운영할 때, 인

력(FTE)과 도구, 정책의 중복으로 인한 비효율이 심각한 문제로 지적되어 왔습니다. KubeVirt 기반의 단일 표면 수렴을 통해, 운영팀은 하나의 API와 도구 체계에서 모든 워크로드를 관리할 수 있게 되며, 정책·보안·감사 기준도 일관되게 적용할 수 있습니다. 이는 조직의 운영 효율성과 보안 수준을 동시에 높이는 핵심 전략입니다.

Pod vs VirtualMachine CRD 비교 표

오브젝트 유형	선언형 YAML	오토스케일	네트워크	스토리지
Pod	지원	지원	K8s CNI	PVC/CSI
VirtualMachine	지원	제한적	Multus	PVC/CSI/VM

이 표를 통해 Pod와 VirtualMachine CRD의 기능적 차이를 명확히 파악할 수 있습니다. Pod는 오토스케일이 완전 지원되지만, VirtualMachine은 일부 제한이 있으며, 네트워크는 Multus와 같은 추가 플러그인을 통해 확장할 수 있습니다. 스토리지 역시 K8s의 PVC/CSI 표준을 활용할 수 있어, VM과 컨테이너 모두에서 일관된 스토리지 관리가 가능합니다.

12장: 결론 — IT 의사결정자 체크리스트와 운영 모델 선택 권장사항

12.1 4대 운영 모델 체크리스트 기반 현 상태 진단

IT 인프라의 운영 모델은 단순한 기술적 선택을 넘어, 조직의 비즈니스 목표와 전략적 방향성에 깊이 연관되어 있습니다. 최근 클라우드 네이티브 전환이 가속화됨에 따라, 선언형 배포, 불변 인프라, 자동 확장, 단일 관리 표면 등 네 가지 핵심 특성이 IT 운영의 표준으로 자리 잡고 있습니다. 본 절에서는 사전조사에서 도출된 핵심 이슈와 4대 운영 모델 주제를 바탕으로, 실제 IT 의사결정자가 활용할 수 있는 진단 체크리스트를 제시합니다. 각 체크리스트는 조직의 현재 상태를 객관적으로 평가하고, 향후 12개월 실행 계획의 근거를 마련하는 데 활용될 수 있습니다. 이를 통해 조직은 기술적·운영적 리스크를 최소화하고, 효율적이고 일관된 IT 운영 체계를 구축할 수 있습니다.

12.1.1 운영 모델 체크리스트 — 선언형·불변·자동확장·단일 표면 4항

운영 모델의 성숙도를 평가할 때 가장 중요한 네 가지 기준은 선언형 배포, 불변 인프라, 자동 확장, 그리고 관리 포인트의 단일화입니다. 선언형 배포와 GitOps 적용 여부는 조직의 배포 프로세스가 얼마나 자동화되고 표준화되어 있는지를 보여줍니다. Kubernetes 환경에서는 YAML 파일을 이용한 선언형 배포와 GitOps(Argo CD, Flux 등)를 통해 코드 기반의 변경 관리가 가능해집니다. 이로 인해 배포의 멱등성(idempotency)이 확보되고, 운영 중 발생할 수 있는 드리프트나 수동 변경에 의한 사고가 크게 줄어듭니다.

불변 아티팩트 이미지 관리는 컨테이너 환경에서의 핵심 원칙 중 하나로, 빌드 시점에 고정된 이미지를 운영 환경에 배포함으로써 예측 가능한 운영과 보안성을 동시에 확보할 수 있습니다. OverlayFS 기반의 이미지 레이어링, 이미지 서명(Cosign), 런타임 변조 탐지(Falco) 등 다양한 보안 및 감사 기능이 결합되어, 운영 중 발생할 수 있는 비인가 변경을 원천적으로 차단합니다. 반면, VM 환경에서는 긴급 패치나 수동 변경이 빈번하게 발생하여, 운영 환경의 일관성이 저하되고 원인 불명의 장애가 누적될 수 있습니다.

자동 확장(HPA·KEDA) 적용 여부는 IT 인프라의 자원 활용 효율성과 직결됩니다. Kubernetes의 Horizontal Pod Autoscaler(HPA), KEDA(이벤트 기반 스케일러), Cluster Autoscaler 등은 실시간 워크로드 변화에 따라 파드와 노드를 동적으로 확장·축소함으로써, 하드웨어 자원의 활용도를 극대화합니다. VM 환경에서는 피크 부하 기준으로 자원을 할당해야 하므로, 평균 자원 활용률이 낮아지고 비용 효율성이 떨어집니다. 자동 확장 기능을 도입하면 동일 하드웨어에서 더 많은 워크로드를 처리할 수 있으며, 하드웨어 투자 비용도 절감할 수 있습니다.

마지막으로, 관리 포인트 단일 표면(K8s API) 수렴 여부는 운영팀의 효율성과 장애 대응 속도에 큰 영향을 미칩니다. VM 환경에서는 하이퍼바이저, OS, 애플리케이션 등 각 계층별로 별도의 관리 도구와 정책이 필요해 이중 관리 포인트가 발생합니다. Kubernetes 환경에서는 모든 운영이 K8s API와 선언형 오브젝트로 통합되어, 운영팀, 모니터링, 보안, CI/CD 파이프라인이 단일 관리 표면에서 일관되게 관리됩니다. 이로 인해 인건비, 운영 복잡도, 장애 대응 시간이 구조적으로 감소합니다.

체크리스트 항목	현재 상태	12개월 목표
[] 배포가 선언형(YAML+GitOps)인가?		

[] 이미지가 불변 아티팩트로 관리되는가?		
[] 자동 확장(HPA/KEDA)이 적용되는가?		
[] 관리 포인트가 단일 표면(K8s API)인가?		

12.1.2 IaaS 전략 가치 체크리스트 — 특수 워크로드 %·3년 라이선스 전망·AI/MSA 전제

IaaS 전략 가치는 조직 내 워크로드의 특수성과 미래 예산 전망, 그리고 AI 및 마이크로서비스 아키텍처(MSA) 도입 여부에 따라 달라집니다. 첫 번째로, 특수 커널이나 레거시 드라이버, 실시간 OS, Windows 워크로드 등 VM 유지가 불가피한 워크로드의 비율을 산정하는 것이 중요합니다. 이 비율이 1030% 수준이라면, 나머지 7090%는 컨테이너화가 가능하다는 근거가 되며, 조직의 현대화 전략 수립에 핵심적인 데이터로 활용될 수 있습니다.

두 번째로, 컨테이너화 가능 워크로드의 비율을 평가해야 합니다. 애플리케이션 구조, 내부통제 요건, 보안 요구사항 등을 면밀히 검토하여, 실제로 컨테이너 전환이 가능한 워크로드를 수치로 제시해야 합니다. 이 과정에서 OS 종속 에이전트, 패치 관리, 백업 등 기존 VM 기반 운영의 한계를 명확히 진단할 수 있으며, 컨테이너 전환의 기대 효과와 리스크를 동시에 파악할 수 있습니다.

세 번째 항목은 VMware, OpenStack 등 IaaS 라이선스의 3년 예산 전망입니다. 2024년 Broadcom의 VMware 인수 이후, 영구 라이선스 중단, 번들 구독 강제, 코어당 과금 전환 등 구조적 변화가 예산에 미치는 영향을 실제 시나리오로 계산해야 합니다. 예를 들어, 기존에 영구 라이선스를 사용하던 조직은 구독형 라이선스 전환에 따라 연간 예산이 30~50% 이상 증가할 수 있으며, 이는 조직의 회계적 리스크로 이어질 수 있습니다.

마지막으로, 신규 AI 및 MSA 워크로드가 VM 환경에서 운영되어야 할 기술적 필연성이 있는지 점검해야 합니다. AI 스택(모델, Vector DB, RAG, API, MCP, Agent)은 선언형, 자동확장, GPU 공유, 이벤트 기반 특성을 요구하며, VM 환경에서는 이러한 요구사항을 효율적으로 충족하기 어렵습니다. 실제로 글로벌 기업의 사례를 보면, AI 및 MSA 워크로드의 대부분이 컨테이너 환경에서 운영되고 있으며, VM 환경은 점차 예외적인 케이스로 전환되고 있습니다.

체크리스트 항목	수치/정책 근거
[] 특수 커널·드라이버 요구 VM %는?	
[] 나머지는 컨테이너화 가능한가?	
[] VMware/OpenStack 라이선스 예산의 3년 전망은?	
[] 신규 AI/MSA 워크로드를 VM으로 운영해야 할 기술적 필연 이유가 있는가?	

12.1.3 관리 포인트·목표 아키텍처 체크리스트 — FTE·vMotion 중복·KubeVirt 흡수

IT 인프라의 관리 포인트와 목표 아키텍처를 진단하는 것은 조직의 운영 효율성과 비용 구조를 결정 짓는 핵심 요소입니다. 첫 번째로, VM 운영팀과 K8s 운영팀 각각의 FTE(Full-Time Equivalent)를 산정하는 작업이 필요합니다. 이중 운영이 지속될 경우, 인건비와 교육 비용이 누적되고, 업무 중복으로 인해 운영 효율성이 저하됩니다. 반면, 단일 K8s 표면으로 수렴하면 FTE 재배치 효과가 발생하여, 인력의 효율적 활용이 가능해집니다. 실제로 대형 금융사나 제조업체에서는 이중 운영에서 단일화로 전환할 때 인력 운영 비용이 20~30% 절감된 사례가 보고되고 있습니다.

두 번째로, 모니터링 도구의 일치 여부를 확인해야 합니다. VM 환경에서는 vCenter, OS 에이전트, 별도 APM(Application Performance Monitoring) 도구를 활용하는 반면, K8s 환경에서는 Prometheus, OpenTelemetry, Grafana 등 클라우드 네이티브 도구가 표준으로 자리 잡고 있습니다. 두 환경의 모니터링 도구가 일치하지 않으면, 이중 알람과 대시보드로 인해 장애 대응이 지연되고, 내부 감사 및 규제 대응에도 불리하게 작용할 수 있습니다. 따라서 모니터링 체계의 통합은 운영 효율성뿐 아니라, 컴플라이언스 측면에서도 중요한 과제입니다.

세 번째로, 보안 정책 배포 경로의 일원화 여부를 점검해야 합니다. VM 환경에서는 OS 단위 에이전트, K8s 환경에서는 OPA(Open Policy Agent), Gatekeeper, Falco 등 네이티브 정책 엔진을 활용합니다. 정책 배포 경로가 단일화되면, 보안 감사와 컴플라이언스 대응이 훨씬 효율적으로 이루어질 수 있습니다. 예를 들어, 금융권에서는 정책 배포 경로 일원화를 통해 보안 사고 대응 시간을 30% 이상 단축한 사례가 있습니다.

네 번째 항목은 하이퍼바이저 비용이 K8s 노드당 몇 %를 차지하는지 산정하는 것입니다. VM 위에 K8s를 올리는 구조에서는 하이퍼바이저 라이선스, 게스트 OS 패치, K8s 운영이 3중으로

누적되어 TCO가 비효율적으로 상승합니다. 실제로, 하이퍼바이저 비용이 전체 인프라 비용의 20~40%를 차지하는 경우도 있으며, 이는 클라우드 네이티브 전환의 주요 동기가 됩니다.

다섯 번째로, vMotion, HA 등 하이퍼바이저 기능이 K8s의 Pod 재스케줄과 중복되는지 확인해야 합니다. 동일 기능을 두 계층에서 중복 구매하거나 운영하면, 회계적·운영적 낭비가 발생합니다. 예를 들어, K8s의 Pod 재스케줄 기능이 이미 장애 복구와 확장성을 제공하는데도 불구하고, 하이퍼바이저의 HA 기능에 추가 비용을 지불하는 것은 비효율적입니다.

마지막으로, KubeVirt로 잔존 VM을 K8s 오브젝트로 흡수할 수 있는지 판단하는 것이 중요합니다. KubeVirt는 VM을 선언형 YAML로 관리하고, GitOps 파이프라인에서 Pod와 VM을 동일하게 배포할 수 있도록 지원합니다. 이 구조는 “VM의 완전한 폐지”가 아니라, “기본 운영 단위를 컨테이너/K8s로 전환하고, 불가피한 VM은 예외적으로 KubeVirt로 흡수”하는 현실적 목표 아키텍처를 완성합니다. 실제로 글로벌 제조사와 금융권에서는 KubeVirt를 도입하여 레거시 VM을 단계적으로 흡수하고, 운영 표준을 K8s로 통합하는 전략을 채택하고 있습니다.

체크리스트 항목	근거/책임자
<input type="checkbox"/> VM 운영팀 + K8s 운영팀 각각의 FTE는?	
<input type="checkbox"/> 두 환경 모니터링 도구가 일치하는가?	
<input type="checkbox"/> 보안 정책 배포 경로가 일원화되어 있는가?	
<input type="checkbox"/> 하이퍼바이저 비용이 K8s 노드당 몇 %인가?	
<input type="checkbox"/> vMotion/HA가 K8s 자체 기능(Pod 재스케줄)과 중복되는가?	
<input type="checkbox"/> KubeVirt로 잔존 VM을 흡수할 수 있는가?	

12.2 12개월 단위 실행 권장 순서와 결정 프레임

조직이 IT 인프라 현대화를 추진할 때, 단기적 실행 계획과 명확한 결정 프레임을 수립하는 것은 매우 중요합니다. 본 절에서는 앞서 제시한 체크리스트 결과를 바탕으로, 조직이 12개월 단위로 실질적으로 실행할 수 있는 전략적 의사결정 프레임을 제안합니다. 기존 워크로드를 현대화 대상, KubeVirt 흡수 대상, VM 유지 예외로 3분류하여, 각 분류별로 최적의 전환 로드맵을 설계할 수 있습니다. 또한, 전환 성과를 객관적으로 측정할 수 있는 KPI(총소유비용, 집적률, 배포 리드타임)를 정의함으로써, 기술적·경영적 의사결정이 일관되게 연결되는 구조를 제공합니다. 이를 통해 조직은

리스크를 최소화하면서도, 단계적이고 체계적인 IT 운영 혁신을 달성할 수 있습니다.

12.2.1 현대화 대상·KubeVirt 흡수·VM 유지 3분류 결정 프레임



[그림 3] 레거시 워크로드 3분류 결정 프레임|800

조직의 워크로드를 효과적으로 분류하는 것은 IT 현대화의 첫걸음입니다. 현대화 대상은 컨테이너화가 가능한 워크로드로, 애플리케이션 구조, OS 종속성, 내부통제 요건을 종합적으로 검토하여 전환 가능성을 판단합니다. 대표적으로 웹 서비스, 마이크로서비스, AI 스택(모델, Vector DB, RAG, Agent) 등이 해당되며, 실제로 글로벌 IT 기업들은 이들 워크로드를 우선적으로 컨테이너화하여 운영 효율성과 확장성을 극대화하고 있습니다. 12개월 실행 계획에는 컨테이너 이미지 빌드, GitOps 파이프라인 구축, 자동 확장(HPA·KEDA) 적용, 관리 포인트 단일화 등이 포함됩니다. 예를 들어, Adidas는 컨테이너화와 GitOps 도입을 통해 배포 속도를 3~4배 개선하고, 운영 복잡도를 크게 낮춘 사례가 있습니다.

KubeVirt 흡수 대상은 커널 요구사항이 상이하거나, 레거시 드라이버, 실시간 OS, Windows 워크로드 등 컨테이너화가 어려운 예외 케이스입니다. 이들은 KubeVirt를 통해 K8s 오브젝트(VirtualMachine CRD)로 관리할 수 있으며, 동일한 API와 GitOps 파이프라인에서 Pod와 VM을 함께 배포할 수 있습니다. 대표 사례로는 ERP, 특수 장비 제어, 금융·공공 규제 대응 워크로드가 있으며, 실제로 금융권에서는 KubeVirt를 활용해 레거시 VM을 단계적으로 흡수하고, 운영 표준을

K8s로 통합하는 전략을 채택하고 있습니다. 12개월 실행 계획에는 KubeVirt 배포, 선언형 YAML 관리, 운영팀 재교육이 포함됩니다.

VM 유지 예외는 하드 멀티테넌시 규제, 특수 보안 요구 등으로 인해 VM 유지가 불가피한 워크로드입니다. 이 경우 Kata Containers, gVisor 등 하이퍼바이저 수준의 격리 옵션을 검토하고, 관리 포인트 이중화를 최소화하는 전략을 수립해야 합니다. 대표적으로 규제 심사관 요구, 실시간 OS, 특수 커널 워크로드 등이 해당되며, 실제로 일부 공공기관이나 금융기관에서는 규제 준수를 위해 VM 환경을 유지하는 사례가 있습니다. 12개월 실행 계획에는 VM 유지 정책의 명문화, 운영팀 분리, 규제 대응 문서화가 포함됩니다.

분류	판정 조건	대표 사례	12개월 액션
현대화 대상	컨테이너화 가능	웹서비스, AI 스택	이미지 빌드, GitOps
KubeVirt 흡수	커널/드라이버/OS 특수 요구	ERP, 금융 규제 워크로드	KubeVirt, YAML 관리
VM 유지 예외	하드 멀티테넌시, 규제 심사관 요구	실시간 OS, 특수 커널	VM 유지 정책, 분리

12.2.2 이사회 제시용 3대 KPI — TCO·집적률·배포 리드타임

IT 인프라 전환의 성과를 객관적으로 측정하고, 이사회 등 주요 의사결정자에게 설득력 있게 제시하기 위해서는 명확한 KPI(핵심성과지표) 설정이 필수적입니다. 첫 번째 KPI는 TCO(총소유비용) 변화율로, 라이선스, 인건비, 전력, 하드웨어 사이징 등 4대 비용 축에서 클라우드 네이티브가 구조적으로 유리함을 수치로 제시해야 합니다. 예를 들어, VM 환경 대비 컨테이너 환경에서 3~5배 집적률 향상, 라이선스 비용 30~50% 절감, 인건비 재배치 효과가 실제로 보고되고 있습니다. Adidas와 Spotify 사례에서는 클라우드 네이티브 전환을 통해 연간 수십억 원의 비용 절감 효과를 달성하였습니다.

두 번째 KPI는 노드 집적률로, 동일 하드웨어에서 수용할 수 있는 워크로드의 양을 의미합니다. Google SRE 책과 Borg 사례 등에서는 컨테이너 환경이 VM 대비 3~10배 더 많은 워크로드를 수용할 수 있음을 입증하고 있습니다. 실제로 Adidas의 사례에서는 컨테이너화 이후 배포 속도가 3~4배 개선되었고, Spotify에서는 신규 서비스 프로비저닝 시간이 1시간에서 수 분으로 단축되었습니다. 이러한 집적률 향상은 하드웨어 투자 비용의 절감과 운영 효율성 증대로 이어집니다.

세 번째 KPI는 배포 리드타임입니다. 선언형, 불변, 자동확장 모델을 적용하면 배포 리드타임이

수일에서 수 분 수준으로 단축되고, 장애 격리 및 복구 시간(MTTR)도 50% 이하로 감소합니다. Booking.com, Spotify, Adidas 등 글로벌 사례를 벤치마크로 삼아, 실제 성과 수치를 이사회에 제시할 수 있습니다. 예를 들어, Spotify는 CI/CD 파이프라인과 GitOps 도입을 통해 신규 서비스 배포 시간을 대폭 단축하였으며, 장애 발생 시 복구 시간도 크게 줄였습니다.

KPI	벤치마크 사례	교차 지표
TCO 변화율	Adidas, Spotify	35배 집적률, 3050% 절감
노드 집적률	Google Borg, Adidas	동일 H/W 3~10배 워크로드
배포 리드타임	Spotify, Booking.com	1시간→수 분, MTTR 50% 이하

부록: 전체 출처 목록

부록: 전체 출처 목록

본 부록 장은 백서 전반에 걸쳐 인용된 모든 공식 문서, 기술 블로그, 사례 연구, 제품 페이지, 오픈소스 프로젝트, 벤치마크, 라이선스 안내, 정책 변화, 운영 사례 등 다양한 출처를 체계적으로 정리하여 제공합니다. 독자는 이 목록을 통해 각 장의 기술적 논의와 근거가 실제로 어떤 자료에 기반하고 있는지 명확하게 확인할 수 있으며, 최신 기술 동향과 실무 적용 사례를 직접 검증할 수 있습니다. 출처 목록은 백서의 신뢰성과 투명성을 높이는 데 중요한 역할을 하며, 각 항목은 고유 식별자(S01~S64)와 함께 원문 URL을 명시하였습니다. 선정 기준은 기술적 깊이, 실무 적용성, 공식성, 그리고 최신성 등을 모두 고려하였으며, AI·클라우드 네이티브 전환, 운영 모델, 벤더 정책 변화 등 핵심 논의의 기반이 되는 자료로 구성되어 있습니다. 아래에 정리된 출처들은 각 장의 논리 전개와 비교 분석, 운영 모델 제안, 기술적 트렌드 설명의 근거로 활용되었습니다.

S01~S10: IaaS·가상화·클라우드 네이티브 기원 및 역사

VMware, OpenStack, Nutanix, CloudStack 등 주요 IaaS 솔루션의 역사와 기술적 배경을 다루는 공식 페이지와 제품 소개 자료를 포함하고 있습니다. 이 영역의 출처들은 하이퍼바이저 기술의 탄생 배경, VM 단위 운영의 구조적 한계, 라이선스 모델의 변화, 그리고 클라우드 네이티브로의

진화 과정에 대한 근거를 제공합니다. 예를 들어, VMware의 공식 연혁 페이지(S01)는 ESX 출시, VMotion 등 핵심 기술 발전 과정을 상세히 다루고 있으며, Docker 공식 가이드(S03)는 컨테이너 기술이 인프라 운영에 가져온 UX 혁신을 설명합니다. 또한, Kubernetes의 공식 개념 문서(S05)와 Google Borg 논문(S06)은 대규모 클러스터 관리와 선언형 인프라의 등장 배경을 이해하는 데 필수적인 자료입니다. 이 외에도 CNCF 공식 정의(S07), OpenStack 및 CloudStack의 역사와 주요 기능(S08, S09), Nutanix의 HCI 개념(S10) 등은 클라우드 네이티브 인프라의 발전사를 종합적으로 조망할 수 있도록 도와줍니다.

- S01:<https://www.vmware.com/company/history.html>
VMware의 창립, ESX 출시, VMotion 등 핵심 연혁과 기술 발전 과정 확인.
- S02:<https://www.broadcom.com/company/news/financial-releases/61491>
Broadcom의 VMware 인수 공식 발표 및 재무 자료.
- S03:<https://docs.docker.com/get-started/overview/>
Docker 공식 시작 가이드 및 컨테이너 기술의 UX 혁신 근거.
- S04:<https://lwn.net/Articles/531114/>
Linux 컨테이너 기술의 커널 기반(Namespace, cgroups) 설명.
- S05:<https://kubernetes.io/docs/concepts/overview/>
Kubernetes 공식 개념 및 선언형·불변 인프라 정의.
- S06:<https://research.google/pubs/large-scale-cluster-management-at-google-with-borg/>
Google Borg의 대규모 클러스터 관리 논문 및 Kubernetes 탄생 배경.
- S07:<https://github.com/cncf/toc/blob/main/DEFINITION.md>
CNCF v1.0 공식 정의(컨테이너, 서비스 메시, 마이크로서비스, 불변 인프라, 선언형 API).
- S08:<https://www.openstack.org/software/>
OpenStack 공식 소프트웨어 소개 및 핵심 컴포넌트 설명.
- S09:<https://cloudstack.apache.org/about.html>
Apache CloudStack의 역사, 주요 기능, 시장 포지셔닝.
- S10:<https://www.nutanix.com/what-is-hyperconverged-infrastructure>
Nutanix HCI 어플라이언스 개념 및 IaaS 범주 설명.

S11~S20: 라이선스·벤더 정책 변화·VM→K8s 브리지 기술

이 구간의 출처들은 VMware, Broadcom, Red Hat, KubeVirt 등 주요 벤더의 라이선스 정책 변화, 구독 모델 전환, 그리고 VM에서 Kubernetes로의 브리지 기술에 관한 벤치마크 등 실무적 의사결정에 직접적으로 영향을 미치는 자료들로 구성되어 있습니다. 예를 들어, VMware의 라이선스 모델 단일화와 번들 구독 강제 정책(S11), Broadcom 인수 후의 라이선스 분쟁 사례(S12), 그리고 Kubernetes의 오픈소스 라이선스(S13) 등은 벤더 정책 변화의 실제 사례와 그 영향력을 보여줍니다. Red Hat OpenShift(S14), Rancher의 상용 지원 모델(S15), 그리고 Red Hat OpenStack의 판매 중단 및 OpenShift 기반 이행(S16) 등은 엔터프라이즈 환경에서의 클라우드 네이티브 전환 전략을 이해하는 데 중요한 참고 자료입니다. 또한, containerd(S17), IBM VM vs 컨테이너 오버헤드 벤치마크(S18), Kata Containers의 격리 옵션(S19), KubeVirt의 공식 벤치마크(S20) 등은 VM과 컨테이너, 그리고 그 사이의 브리지 기술의 성능과 보안, 운영상의 차이를 비교 분석하는 데 활용됩니다.

- S11:<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
VMware 라이선스 모델 단일화, 번들 구독 강제, 코어당 과금 전환 공식 안내.
- S12:<https://arstechnica.com/information-technology/2024/08/att-sues-broadcom-over-vmware-licensing-changes/>
AT&T의 Broadcom 인수 후 라이선스 분쟁 사례 보도.
- S13:<https://github.com/kubernetes/kubernetes/blob/master/LICENSE>
Kubernetes Apache 2.0 라이선스 원문.
- S14:<https://www.redhat.com/en/technologies/cloud-computing/openshift>
Red Hat OpenShift 공식 소개 및 엔터프라이즈 K8s 디스트리뷰션.
- S15:<https://www.rancher.com/pricing>
Rancher 상용 지원 모델 및 가격 구조 안내.
- S16:<https://www.redhat.com/en/blog/evolution-red-hat-openstack-platform>
Red Hat OpenStack(RHOSP) 판매 중단 및 RHOSO(OpenShift 기반) 이행 공식 블로그.
- S17:<https://www.cncf.io/projects/containerd/>

containerd 프로젝트의 CNCF 공식 소개 및 라이선스 안내.

- S18:<https://dominoweb.draco.res.ibm.com/reports/rc25482.pdf>
IBM VM vs 컨테이너 오버헤드 벤치마크 리포트.
- S19:<https://katacontainers.io/learn/>
Kata Containers 공식 문서 및 하드 멀티테넌시 격리 옵션 설명.
- S20:<https://kubevirt.io/user-guide/operations/benchmarking/>
KubeVirt 공식 벤치마크(I/O, 네트워크 지연 등) 및 운영 가이드.

S21~S30: 운영 모델·관측성·자동확장·컨테이너 네이티브 기능

이 구간의 출처들은 Google SRE, Red Hat, NVIDIA, NIST, Kubernetes 공식 문서 등 운영 모델, 집적률, 자동확장, GPU Operator, 컨테이너 격리 원리 등 실무 적용에 필수적인 기술적 근거 자료를 포함하고 있습니다. Google SRE 공식 서적(S21)은 대규모 인프라 운영 모델과 컨테이너 집적률에 대한 실제 사례를 제공하며, Red Hat의 Linux 컨테이너 구조 설명(S22)은 커널 공유와 보안 영향에 대한 이해를 돕습니다. NVIDIA GPU Operator 공식 문서(S23)는 GPU 공유 및 MIG, Time-slicing 등 AI 워크로드에 필수적인 기능을 다루고 있습니다. NIST의 클라우드 정의(S24), Kubernetes의 Pod 개념(S25), 컨트롤러 루프와 역등성 구현(S26), Gateway API(S27), 오브젝트 선언형 구성 원리(S28), HPA 공식 문서(S29), KEDA의 이벤트 기반 자동확장(S30) 등은 클라우드 네이티브 운영의 자동화, 확장성, 관측성, 그리고 실무 적용에 필요한 핵심 기능을 이해하는 데 중요한 자료입니다.

- S21:<https://sre.google/sre-book/production-environments/>
Google SRE 공식 책, Borg 사례, 컨테이너 집적률·운영 모델 근거.
- S22:<https://www.redhat.com/en/topics/containers/whats-a-linux-container>
Linux 컨테이너 구조, 커널 공유, 보안 영향 공식 설명.
- S23:<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>
NVIDIA GPU Operator 공식 문서 및 GPU 공유·MIG·Time-slicing 설명.
- S24:<https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-145.pdf>

NIST SP 800-145 클라우드 정의(IaaS·PaaS·SaaS) 원문.

- S25:<https://kubernetes.io/docs/concepts/workloads/pods/>
Kubernetes Pod 개념 및 배포 단위 정의.
- S26:<https://kubernetes.io/docs/concepts/workloads/controllers/>
Kubernetes 컨트롤러 루프·역등성 구현 설명.
- S27:<https://gateway-api.sigs.k8s.io/>
Kubernetes Gateway API 공식 문서.
- S28:<https://kubernetes.io/docs/concepts/configuration/>
Kubernetes 오브젝트 선언형 구성 원리.
- S29:<https://kubernetes.io/docs/tasks/run-application/horizontal-pod-autoscale/>
Horizontal Pod Autoscaler(HPA) 공식 문서.
- S30:<https://keda.sh/docs/>
KEDA 이벤트 기반 자동확장 공식 문서.

S31~S40: Operator·Service Mesh·정책 엔진·관측성·보안

이 영역의 출처들은 Kubernetes Operator, OperatorHub, OPA Gatekeeper, Kyverno, Istio, Cilium, Argo CD, Flux, KubeVirt, eBPF 등 클라우드 네이티브 환경에서의 확장성, 정책 관리, 서비스 메시, 선언형 운영, 브리지 기술 등 고급 기능에 대한 공식 문서와 사례를 포함하고 있습니다. 예를 들어, Kubernetes Operator 개념(S31)과 OperatorHub(S32)는 확장 구조와 다양한 오퍼레이터의 실제 적용 사례를 보여줍니다. OPA Gatekeeper(S33)와 Kyverno(S34)는 정책 관리와 보안 강화 방법을 설명하며, Istio(S35)와 Cilium(S36)은 서비스 메시와 네트워크 정책, eBPF 기반 관측성에 대한 실질적인 가이드를 제공합니다. 또한, Argo CD(S37)와 Flux(S38)는 GitOps 기반의 선언형 운영 패턴을, KubeVirt(S39)는 VM과 컨테이너 브리지 기술을, eBPF(S40)는 커널 수준의 고급 관측성 및 보안 응용 사례를 다루고 있습니다. 이러한 자료들은 클라우드 네이티브 환경에서의 고도화된 운영과 보안, 자동화의 실질적 구현 방법을 이해하는 데 필수적입니다.

- S31:<https://kubernetes.io/docs/concepts/extend-kubernetes/operator/>
Kubernetes Operator 개념 및 확장 구조.

- S32:<https://operatorhub.io/>
OperatorHub 공식 사이트, 다양한 K8s Operator 목록.
- S33:<https://open-policy-agent.github.io/gatekeeper/>
OPA Gatekeeper 정책 엔진 공식 문서.
- S34:<https://kyverno.io/docs/>
Kyverno 정책 관리 공식 문서.
- S35:<https://istio.io/latest/docs/>
Istio 서비스 메시 공식 문서.
- S36:<https://docs.cilium.io/>
Cilium 네트워크 정책 및 eBPF 기반 관측성 공식 문서.
- S37:<https://argo-cd.readthedocs.io/>
Argo CD 공식 문서 및 GitOps 운영 패턴 설명.
- S38:<https://fluxcd.io/flux/concepts/>
Flux 공식 문서 및 GitOps 개념.
- S39:<https://kubevirt.io/user-guide/>
KubeVirt 공식 사용자 가이드.
- S40:<https://ebpf.io/applications/>
eBPF 응용 사례 공식 사이트.

S41~S50: Kubernetes 릴리스·Kubeflow·Booking.com·CERN·Adidas·Windows 컨테이너

이 구간의 출처들은 Kubernetes의 주요 릴리스 노트, Kubeflow, 대규모 베어메탈 운영 사례, Windows 컨테이너 공식 문서 등 최신 기술 동향과 실제 운영 사례에 관한 근거 자료로 구성되어 있습니다. Kubernetes v1.27, v1.29, v1.31 릴리스 노트(S41~S43)는 최신 기능과 개선 사항을 파악하는 데 필수적이며, Kubeflow 공식 사이트(S44)는 AI 워크플로우 자동화의 실제 구현 사례를 제공합니다. AT&T의 OpenStack Helm+K8s 전환 사례(S45), Adidas와 Spotify의 CNCF 공식 사례(S46, S47)는 대규모 베어메탈 환경에서의 K8s 전환 성과와 마이크로서비스 도입 효과를 보여줍니다. MetalLB(S48), Rook-Ceph(S49), Windows 컨테이너 공식 문서(S50) 등은

베어메탈 환경에서의 로드밸런서 구현, 스토리지 오퍼레이터, 그리고 Windows 컨테이너의 커널 버전 일치 요구 등 현실적인 운영 이슈를 다루고 있습니다. 이러한 자료들은 최신 기술 트렌드와 실제 적용 사례를 이해하는 데 매우 유용합니다.

- S41:<https://kubernetes.io/blog/2023/04/11/kubernetes-v1-27-release/>
Kubernetes v1.27 릴리스 노트.
- S42:<https://kubernetes.io/blog/2023/12/13/kubernetes-v1-29-release/>
Kubernetes v1.29 릴리스 노트.
- S43:<https://kubernetes.io/blog/2024/08/13/kubernetes-v1-31-release/>
Kubernetes v1.31 릴리스 노트.
- S44:<https://www.kubeflow.org/>
Kubeflow 공식 사이트 및 AI 워크플로우 자동화.
- S45:<https://www.openstack.org/user-stories/att/>
AT&T OpenStack Helm+K8s 전환 사례.
- S46:<https://www.cncf.io/case-studies/adidas/>
Adidas CNCF 공식 사례, 베어메탈 K8s 전환 성과.
- S47:<https://www.cncf.io/case-studies/spotify/>
Spotify CNCF 공식 사례, 마이크로서비스·K8s 전환.
- S48:<https://metallb.universe.tf/>
MetalLB 공식 문서, 베어메탈 K8s 로드밸런서 구현.
- S49:<https://rook.io/>
Rook-Ceph 공식 문서, K8s 스토리지 오퍼레이터.
- S50:<https://learn.microsoft.com/en-us/virtualization/windowscontainers/>
Windows 컨테이너 공식 문서, 커널 버전 일치 요구 설명.

S51~S60: 설치 가이드·AI 스택·관측성·보안·연계 조사·시장 동향

이 영역의 출처들은 OpenStack 설치 가이드, KServe, MCP, Tekton, Prometheus Operator, OpenTelemetry, Falco, CNCF Landscape, 연계 조사, StackOverflow 설문 등 실무적 운영과

시장 동향을 파악하는 데 필요한 자료로 구성되어 있습니다. OpenStack 설치 가이드(S51)는 실제 인프라 구축에 필요한 절차와 모범 사례를 제공하며, KServe(S52), MCP(S53), Tekton(S54) 등은 AI 워크플로우, 모델 서빙, CI/CD 파이프라인 등 최신 AI 및 DevOps 스택의 공식 문서입니다. Prometheus Operator(S55), OpenTelemetry(S56), Falco(S57)는 관측성과 보안, 런타임 위협 탐지 등 운영 안정성 확보에 필수적인 도구와 그 활용법을 안내합니다. CNCF Landscape(S58) 과 연례 조사(S59), StackOverflow 개발자 설문(S60)은 글로벌 시장 동향, 기술 채택률, 실무자 선호도 등 산업 전반의 흐름을 이해하는 데 중요한 참고 자료입니다.

- S51:<https://docs.openstack.org/install-guide/>
OpenStack 설치 가이드 공식 문서.
- S52:<https://kserve.github.io/website/>
KServe 모델 서빙 공식 문서.
- S53:<https://modelcontextprotocol.io/>
MCP(Model Context Protocol) 공식 사이트.
- S54:<https://tekton.dev/docs/>
Tekton 공식 문서, K8s 네이티브 CI/CD 파이프라인.
- S55:<https://prometheus-operator.dev/>
Prometheus Operator 공식 문서.
- S56:<https://opentelemetry.io/>
OpenTelemetry 공식 문서.
- S57:<https://falco.org/>
Falco 런타임 위협 탐지 공식 문서.
- S58:<https://landscape.cncf.io/>
CNCF Landscape 프로젝트 그래프.
- S59:<https://www.cncf.io/reports/cncf-annual-survey-2023/>
CNCF 연례 조사, 글로벌 K8s 채택률·운영 트렌드.
- S60:<https://survey.stackoverflow.co/2024/>
StackOverflow 2024 개발자 설문 결과.

S61~S64: 벤더 제품·문서·국산 PaaS·운영 사례

마지막 구간의 출처들은 VMware vSphere, OpenStack, Nutanix, MSAP.ai 등 주요 벤더 제품 공식 문서와 국내 국산 PaaS 기준 제품 페이지, 그리고 실제 운영 사례에 대한 근거 자료를 포함하고 있습니다. VMware vSphere 공식 문서(S61), OpenStack 공식 문서(S62), Nutanix AHV 공식 문서(S63)는 각 벤더의 제품 구조, 기능, 운영 가이드라인을 상세히 제공합니다. 또한, MSAP.ai 공식 제품 페이지(S64)는 국내 공공 및 금융 분야에서의 실제 채택 사례와 국산 PaaS의 기술적 특징점을 확인할 수 있는 자료입니다. 이러한 출처들은 벤더별 제품의 공식적 기능, 지원 정책, 그리고 국내외 실무 적용 사례를 종합적으로 이해하는 데 중요한 역할을 합니다.

- S61:<https://docs.vmware.com/en/VMware-vSphere/index.html>
VMware vSphere 공식 문서.
- S62:<https://docs.openstack.org/>
OpenStack 공식 문서.
- S63:<https://portal.nutanix.com/page/documents/list?type=software>
Nutanix AHV 공식 문서.
- S64:<https://msap.ai/>
MSAP.ai 공식 제품 페이지, 국내 공공·금융 채택 사례.

Appendix

References

1. AT&T OpenStack User Story.<https://www.openstack.org/user-stories/att/>
2. Adidas CNCF Case Study.<https://www.cncf.io/case-studies/adidas/>
3. Adidas Case Study. (2022).<https://www.cncf.io/case-studies/adidas/>
4. Adidas. (2022). “CNCF Case Study: Adidas”.<https://www.cncf.io/case-studies/adidas/>

5. AlmaLinux. (2024).<https://almalinux.org/>
6. Apache CloudStack About. (n.d.).<https://cloudstack.apache.org/about.html>
7. Argo CD Project. (2024). “Argo CD Documentation”.<https://argo-cd.readthedocs.io/>
8. Argo CD. (2024). “Argo CD Documentation”.<https://argo-cd.readthedocs.io/>
9. Booking.com KubeCon EU 발표 사례.
10. Bottlerocket. (2024).<https://bottlerocket.dev/>
11. Broadcom Financial Releases. (2024).<https://www.broadcom.com/company/news/financial-releases/61491>
12. Broadcom. (2024). “VMware Simplifies Its Portfolio and Licensing Model”.<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
13. CERN OpenStack User Story.<https://www.openstack.org/user-stories/>
14. CNCF Annual Survey 2023.<https://www.cncf.io/reports/cncf-annual-survey-2023/>
15. CNCF Definition v1.0. (n.d.).<https://github.com/cncf/toc/blob/main/DEFINITION.md>
16. CNCF Landscape.<https://landscape.cncf.io/>
17. CNCF. (2018). “Cloud Native Definition v1.0”.<https://github.com/cncf/toc/blob/main/DEFINITION.md>
18. CNCF. (2021). “Cloud Native Definition v1.0 Released”.<https://www.cncf.io/announcements/2021/12/08/cncf-cloud-native-definition-v1-0-released/>
19. CNCF. (2023). “Annual Survey”.<https://www.cncf.io/reports/cncf-annual-survey-2023/>
20. CNCF. (2024). “CNCF Landscape”.<https://landscape.cncf.io/>
21. CNCF. (2024). “Certified Kubernetes Conformance Program”.<https://www.cncf.io/certification/software-conformance/>
22. CNCF. (2024). “Projects: containerd”.<https://www.cncf.io/projects/containerd/>

23. CentOS. (2020). “Future is CentOS Stream” .<https://blog.centos.org/2020/12/future-is-centos-stream/>
24. Cloud Native Computing Foundation. (2018). “Cloud Native Definition v1.0” .<https://github.com/cncf/toc/blob/main/DEFINITION.md>
25. CloudStack. “About CloudStack.”<https://cloudstack.apache.org/about.html>
26. Crossplane. (2024). “Crossplane Documentation” .<https://crossplane.io/>
27. Docker Overview. (n.d.).<https://docs.docker.com/get-started/overview/>
28. Falco. (2024). “Falco Runtime Security” .<https://falco.org/>
29. Falco. (2024). “Falco: Cloud Native Runtime Security” .<https://falco.org/>
30. Falco. (2024).<https://falco.org/>
31. FluxCD. (2024). “Flux Concepts” .<https://fluxcd.io/flux/concepts/>
32. Google Borg Cluster Management. (n.d.).<https://research.google/pubs/large-scale-cluster-management-at-google-with-borg/>
33. Google SRE Book. (2021).<https://sre.google/sre-book/production-environments/>
34. Google SRE. (2022). “Production Environments” .<https://sre.google/sre-book/production-environments/>
35. IBM Research. “VM vs Container Overhead.”<https://dominoweb.draco.res.ibm.com/reports/rc25482.pdf>
36. IBM. (2018). “RC25482: Container vs VM Overhead” .<https://dominoweb.draco.res.ibm.com/reports/rc25482.pdf>
37. Istio Project. (2024). “Istio Documentation” .<https://istio.io/latest/docs/>
38. KEDA Project. (2024). “KEDA Documentation” .<https://keda.sh/docs/>
39. KEDA. (2024). “KEDA Documentation” .<https://keda.sh/docs/>
40. KServe Project. (2024). “KServe Documentation” .<https://kserve.github.io/website/>
41. Kata Containers Project. (2024). “Learn Kata Containers” .<https://katacontainers.io/learn/>
42. Katacontainers. (2024).<https://katacontainers.io/learn/>

43. KubeVirt Benchmarking.<https://kubevirt.io/user-guide/operations/benchmarking/>
44. KubeVirt Project. (2024). “KubeVirt User Guide”.<https://kubevirt.io/user-guide/>
45. KubeVirt User Guide Benchmarking.<https://kubevirt.io/user-guide/operations/benchmarking/>
46. KubeVirt User Guide. (2024).<https://kubevirt.io/user-guide/>
47. KubeVirt User Guide.<https://kubevirt.io/user-guide/>
48. Kubernetes Concepts Overview. (n.d.).<https://kubernetes.io/docs/concepts/overview/>
49. Kubernetes Documentation. (2024). “Concepts Overview”.<https://kubernetes.io/docs/concepts/overview/>
50. Kubernetes Documentation. (2024). “Kubernetes Operators”.<https://kubernetes.io/docs/concepts/extend-kubernetes/operator/>
51. Kubernetes. (2024). “Controllers”.<https://kubernetes.io/docs/concepts/workloads/controllers/>
52. Kubernetes. (2024). “Horizontal Pod Autoscaler”.<https://kubernetes.io/docs/tasks/run-application/horizontal-pod-autoscale/>
53. Kubernetes. (2024). “Kubernetes License”.<https://github.com/kubernetes/kubernetes/blob/master/LICENSE>
54. Kubernetes. (2024). “LICENSE”.<https://github.com/kubernetes/kubernetes/blob/master/LICENSE>
55. Kubernetes. (2024). “Pods”.<https://kubernetes.io/docs/concepts/workloads/pods/>
56. LWN.net. (2013). “Containers: Namespaces, cgroups, and beyond”.<https://lwn.net/Articles/531114/>
57. LXC Article. (2013).<https://lwn.net/Articles/531114/>
58. MSAP.ai 공식 페이지.<https://msap.ai/>
59. MetalLB Documentation.<https://metallb.universe.tf/>

60. Microsoft. (2024). “Windows Containers Documentation”.<https://learn.microsoft.com/en-us/virtualization/windowscontainers/>
61. Model Context Protocol. (2024). “MCP Specification”.<https://modelcontextprotocol.io/>
62. NIST. (2011). “NIST Special Publication 800–145: The NIST Definition of Cloud Computing”.<https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-145.pdf>
63. NVIDIA GPU Operator Documentation.<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>
64. NVIDIA. (2024). “GPU Operator”.<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>
65. NVIDIA. (2024). “NVIDIA GPU Operator Documentation”.<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>
66. Nutanix Hyperconverged Infrastructure. (n.d.).<https://www.nutanix.com/what-is-hyperconverged-infrastructure>
67. Nutanix Portal.<https://portal.nutanix.com/page/documents/list?type=software>
68. Nutanix Software Documentation.<https://portal.nutanix.com/page/documents/list?type=software>
69. Nutanix. “What is hyperconverged infrastructure.”<https://www.nutanix.com/what-is-hyperconverged-infrastructure>
70. OpenStack Install Guide.<https://docs.openstack.org/install-guide/>
71. OpenStack Software. (n.d.).<https://www.openstack.org/software/>
72. OpenStack Software.<https://www.openstack.org/software/>
73. OpenTelemetry. (2024). “OpenTelemetry Documentation”.<https://opentelemetry.io/>
74. OperatorHub. (2024). “OperatorHub”.<https://operatorhub.io/>
75. Prometheus Operator. (2024). “Prometheus Operator Documentation”.<https://prometheus-operator.dev/>

76. Rancher. (2024). “Pricing” .<https://www.rancher.com/pricing>
77. Red Hat. (2024). “Evolution of Red Hat OpenStack Platform.”<https://www.redhat.com/en/blog/evolution-red-hat-openstack-platform>
78. Red Hat. (2024). “KubeVirt Benchmarking Guide” .<https://kubevirt.io/user-guide/operations/benchmarking/>
79. Red Hat. (2024). “OpenShift” .<https://www.redhat.com/en/technologies/cloud-computing/openshift>
80. Red Hat. (2024). “What’s a Linux Container” .<https://www.redhat.com/en/topics/containers/whats-a-linux-container>
81. Red Hat. (2024). “What’s a Linux Container?” .<https://www.redhat.com/en/topics/containers/whats-a-linux-container>
82. Red Hat. (2024). “What’s a Linux container?” .<https://www.redhat.com/en/topics/containers/whats-a-linux-container>
83. Rocky Linux. (2024).<https://rockylinux.org/>
84. Rook-Ceph Documentation.<https://rook.io/>
85. S01:<https://www.vmware.com/company/history.html>
86. S02:<https://www.broadcom.com/company/news/financial-releases/61491>
87. S03:<https://docs.docker.com/get-started/overview/>
88. S04:<https://lwn.net/Articles/531114/>
89. S05:<https://kubernetes.io/docs/concepts/overview/>
90. S06:<https://research.google/pubs/large-scale-cluster-management-at-google-with-borg/>
91. S07:<https://github.com/cncf/toc/blob/main/DEFINITION.md>
92. S08:<https://www.openstack.org/software/>
93. S09:<https://cloudstack.apache.org/about.html>
94. S10:<https://www.nutanix.com/what-is-hyperconverged-infrastructure>
95. S11:<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
96. S12:<https://arstechnica.com/information-technology/2024/08/att-sues-b>

- [roadcom-over-vmware-licensing-changes/](#)
97. S13:<https://github.com/kubernetes/kubernetes/blob/master/LICENSE>
 98. S14:<https://www.redhat.com/en/technologies/cloud-computing/openshift>
 99. S15:<https://www.rancher.com/pricing>
 100. S16:<https://www.redhat.com/en/blog/evolution-red-hat-openstack-platform>
 101. S17:<https://www.cncf.io/projects/containerd/>
 102. S18:<https://dominoweb.draco.res.ibm.com/reports/rc25482.pdf>
 103. S19:<https://katacontainers.io/learn/>
 104. S20:<https://kubevirt.io/user-guide/operations/benchmarking/>
 105. S21:<https://sre.google/sre-book/production-environments/>
 106. S22:<https://www.redhat.com/en/topics/containers/whats-a-linux-container>
 107. S23:<https://docs.nvidia.com/datacenter/cloud-native/gpu-operator/>
 108. S24:<https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-145.pdf>
 109. S25:<https://kubernetes.io/docs/concepts/workloads/pods/>
 110. S26:<https://kubernetes.io/docs/concepts/workloads/controllers/>
 111. S27:<https://gateway-api.sigs.k8s.io/>
 112. S28:<https://kubernetes.io/docs/concepts/configuration/>
 113. S29:<https://kubernetes.io/docs/tasks/run-application/horizontal-pod-autoscale/>
 114. S30:<https://keda.sh/docs/>
 115. S31:<https://kubernetes.io/docs/concepts/extend-kubernetes/operator/>
 116. S32:<https://operatorhub.io/>
 117. S33:<https://open-policy-agent.github.io/gatekeeper/>
 118. S34:<https://kyverno.io/docs/>
 119. S35:<https://istio.io/latest/docs/>
 120. S36:<https://docs.cilium.io/>

121. S37:<https://argo-cd.readthedocs.io/>
122. S38:<https://fluxcd.io/flux/concepts/>
123. S39:<https://kubevirt.io/user-guide/>
124. S40:<https://ebpf.io/applications/>
125. S41:<https://kubernetes.io/blog/2023/04/11/kubernetes-v1-27-release/>
126. S42:<https://kubernetes.io/blog/2023/12/13/kubernetes-v1-29-release/>
127. S43:<https://kubernetes.io/blog/2024/08/13/kubernetes-v1-31-release/>
128. S44:<https://www.kubeflow.org/>
129. S45:<https://www.openstack.org/user-stories/att/>
130. S46:<https://www.cncf.io/case-studies/adidas/>
131. S47:<https://www.cncf.io/case-studies/spotify/>
132. S48:<https://metallb.universe.tf/>
133. S49:<https://rook.io/>
134. S50:<https://learn.microsoft.com/en-us/virtualization/windowscontainers/>
[/](#)
135. S51:<https://docs.openstack.org/install-guide/>
136. S52:<https://kserve.github.io/website/>
137. S53:<https://modelcontextprotocol.io/>
138. S54:<https://tekton.dev/docs/>
139. S55:<https://prometheus-operator.dev/>
140. S56:<https://opentelemetry.io/>
141. S57:<https://falco.org/>
142. S58:<https://landscape.cncf.io/>
143. S59:<https://www.cncf.io/reports/cncf-annual-survey-2023/>
144. S60:<https://survey.stackoverflow.co/2024/>
145. S61:<https://docs.vmware.com/en/VMware-vSphere/index.html>
146. S62:<https://docs.openstack.org/>
147. S63:<https://portal.nutanix.com/page/documents/list?type=software>
148. S64:<https://msap.ai/>

149. Spotify CNCF Case Study.<https://www.cncf.io/case-studies/spotify/>
150. Spotify Case Study. (2022).<https://www.cncf.io/case-studies/spotify/>
151. Tekton. (2024). “Tekton Documentation”.<https://tekton.dev/docs/>
152. Trivy. (2024).<https://trivy.dev/>
153. VMware Cloud Foundation. (2024). “VMware simplifies its portfolio and licensing model.”<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
154. VMware History. (n.d.).<https://www.vmware.com/company/history.html>
155. VMware Licensing Model. (2024).<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
156. VMware vSphere Documentation.<https://docs.vmware.com/en/VMware-vSphere/index.html>
157. VMware. (2024). “VMware Simplifies Its Portfolio and Licensing Model”.<https://blogs.vmware.com/cloud-foundation/2024/01/22/vmware-simplifies-its-portfolio-and-licensing-model/>
158. Velero. (2024). “Velero Documentation”.<https://velero.io/>
159. Velero. (2024).<https://velero.io/>
160. 오픈마루·투라인클라우드. (2024). “MSAP.ai”.<https://msap.ai/>

Glossary

용어	정의
마이크로서비스	독립적으로 배포·확장 가능한 소규모 서비스 단위
멀티클러스터 관리	여러 K8s 클러스터를 단일 표면에서 관리하는 기능.
역등성	같은 명령을 여러 번 실행해도 결과가 항상 동일한 성질.
베어메탈 K8s	물리 서버 위에 직접 K8s 클러스터를 배치하는 운영 모델로, VM 계층을 제거함.
벤더 락인	특정 벤더의 제품·API에 종속되어 이식성·확장성이 제한되는 상태.
불변 인프라	빌드 시점에 고정된 아티팩트로, 실행 중에는 수정이 금지되는 인프라 운영 원칙.
서비스 메시	마이크로서비스 간 트래픽 관리·보안·관측을 담당하는 네트워크 계층

선언형 API	원하는 상태를 선언하고, 컨트롤러가 실제 상태와 비교하여 자동 수렴하는 운영 모델
엔터프라이즈 K8s 디스트리뷰션	벤더가 K8s에 엔터프라이즈 기능을 결합해 제공하는 상용 배포판.
예시적(illustrative)	실제 조직 환경에 따라 달라질 수 있는 예시적 수치임을 명시하는 용어.
집적률	동일 하드웨어에서 수용 가능한 워크로드의 비율.
컨테이너	OS 커널을 공유하며 애플리케이션 실행 환경만 분리하는 경량화 단위
하이퍼바이저	물리 서버의 자원을 분할하여 VM을 생성·관리하는 소프트웨어 계층.
cert-manager	K8s 클러스터 내에서 PKI 인증서 발급·갱신을 자동화하는 오픈소스.
cgroups	Linux 커널에서 CPU·메모리 등 자원 사용량을 제한하는 기술.
Cluster Autoscaler	K8s에서 클러스터 노드 풀을 자동 확장/축소하는 기능.
CNCF	Cloud Native Computing Foundation, 클라우드 네이티브 기술 표준화 및 생태계 구축을 담당하는 글로벌 재단.
CNCF Certified Kubernetes	CNCF가 공식 인증한 Kubernetes 표준 호환 디스트리뷰션.
CRD	Custom Resource Definition, K8s에서 사용자 정의 오브젝트를 생성하는 방법.
CSI	Container Storage Interface, Kubernetes에서 스토리지 연동을 위한 표준.
Declarative API	원하는 상태를 선언적으로 기술하고, 컨트롤러가 자동 수렴하는 API 모델.
Device Plugin	Kubernetes에서 하드웨어 리소스(GPU 등)를 파드에 할당하는 플러그인.
Distroless/minimal base	컨테이너 이미지에서 OS 기능을 제거하고, 애플리케이션 실행 환경만 포함하는 경량 이미지.
DR	Disaster Recovery, 재해 복구 시나리오
EDR	Endpoint Detection & Response, 런타임 위협 탐지 솔루션.
etcd HA	K8s Control Plane의 고가용성(HA)을 위한 etcd 다중 노드 구성.
External DNS	K8s 서비스와 외부 DNS 레코드 자동 연동을 담당하는 오픈소스.
FTE	Full-Time Equivalent, 운영팀 인력 산정 단위.
GitOps	Git을 단일 진실 공급원으로 삼아 선언형 배포와 자동 동기화를 구현하는 운영 모델.
GPU Operator	NVIDIA가 제공하는 Kubernetes GPU 관리 자동화 솔루션.
Guest OS	VM마다 독립적으로 구동되는 OS(커널, systemd, 패키지 등).
gVisor	사용자 공간에서 커널 API를 재구현하여 컨테이너 격리를 강화하는 프로젝트.
HCI	Hyper-Converged Infrastructure, 서버·스토리지·네트워크 통합 어플라이언스.
Helm Chart	Kubernetes 애플리케이션 배포를 위한 패키지 관리 형식.
HPA	Horizontal Pod Autoscaler, Kubernetes에서 파드 수를 자동으로 확장/축소하는 기능.
IaaS	Infrastructure as a Service, VM 단위 인프라 제공 모델.
IaC	Infrastructure as Code, 인프라를 코드로 선언·관리하는 방식
Immutable Infrastructure	빌드 시점에 고정된 아티팩트(컨테이너 이미지)를 운영 중 수정 없이 배포하는 인프라 모델.

K8s	Kubernetes의 약어, 컨테이너 오케스트레이션 플랫폼.
Kata Containers	컨테이너를 경량 VM으로 실행하여 하이퍼바이저 수준 격리를 제공하는 프로젝트.
KEDA	Kubernetes Event-driven Autoscaler, 이벤트 기반으로 파드 수를 동적으로 조정하는 오토스케일러.
KServe	Kubernetes 기반 모델 서빙 프레임워크, InferenceService CRD로 선언형 배포 지원.
Kubernetes 업스트림	CNCF에서 관리하는 오픈소스 Kubernetes 프로젝트, Apache 2.0 라이선스.
Kubernetes(K8s)	컨테이너 오케스트레이션을 위한 오픈소스 플랫폼
KubeVirt	Kubernetes에서 VM을 선언형 오브젝트로 관리하는 브리지 기술.
LGTM 스택	Loki, Grafana, Tempo, Mimir로 구성되는 클라우드 네이티브 관측성 스택
Longhorn	경량 블록 스토리지 솔루션으로, K8s에서 분산 볼륨 관리 제공.
MCP	Model Context Protocol, AI 모델 컨텍스트 관리 표준.
MetalLB	베어메탈 K8s 환경에서 로드밸런서를 제공하는 오픈소스 프로젝트.
MIG	Multi-Instance GPU, 하나의 GPU를 여러 인스턴스로 분할하는 기술.
MSAP.ai	오픈마루·투라인클라우드 공동개발 국내 PaaS 기준 제품, Kubernetes 기반 클라우드 네이티브 플랫폼.
MTTR	Mean Time To Recovery, 장애 복구 평균 시간.
namespace	Linux 커널에서 프로세스·네트워크 등 리소스를 분리하는 격리 기술.
Operator	Kubernetes에서 복잡한 애플리케이션의 자동 배포·운영을 담당하는 커스텀 컨트롤러.
OS tax	VM 운영 시 OS 패치·라이선스·보안 등 중복 비용 구조.
OSS	Open Source Software, 오픈소스 소프트웨어.
OverlayFS	여러 이미지 레이어를 중첩하여 컨테이너 이미지를 관리하는 파일 시스템.
PaaS	Platform as a Service. OS·런타임·미들웨어까지 포함된 플랫폼을 제공하며, 애플리케이션 개발·배포에 집중할 수 있는 서비스 모델.
PAM	Privileged Access Management, 서버 접근제어 솔루션.
Pod	Kubernetes의 최소 배포 단위로, 복수 컨테이너가 네트워크·스토리지 네임스페이스를 공유.
PVC/CSI	Persistent Volume Claim / Container Storage Interface, K8s 스토리지 표준.
RAG	Retrieval Augmented Generation, 외부 데이터 검색과 LLM 생성 결합 프레임워크.
RBAC	Role-Based Access Control, Kubernetes의 권한 분리 기능.
RHOSO	Red Hat OpenShift Service on OpenStack.
RHOSP	Red Hat OpenStack Platform.
Rook-Ceph	K8s 환경에서 Ceph 기반 스토리지를 오케스트레이션하는 오픈소스 프로젝트.
SaaS	Software as a Service. 완성된 애플리케이션을 서비스 형태로 제공하는 모델.

Service Mesh	마이크로서비스 간 트래픽 관리·보안·관측을 담당하는 인프라 계층(예: Istio).
SMS	System Management Server, 자원·성능·로그 모니터링 에이전트.
StatefulSet	Kubernetes에서 상태가 있는 애플리케이션을 관리하는 오브젝트.
TCO	Total Cost of Ownership, 총소유비용.
Time-slicing	GPU를 시간 단위로 분할하여 여러 파드가 번갈아 사용하는 방식.
Triton Inference Server	NVIDIA 공식 멀티 프레임워크 모델 서빙 서버, GPU Operator와 연동.
vLLM	대규모 언어 모델(LLM) 서빙 엔진, GPU 최적화 및 대규모 동시 요청 처리 지원.
VM	Virtual Machine, 하이퍼바이저 기반 격리된 OS 인스턴스.
VM 위 K8s	하이퍼바이저 위에 게스트 OS와 K8s 클러스터를 배치하는 과도기적 인프라 운영 모델.
VPA	Vertical Pod Autoscaler. Kubernetes에서 파드 리소스 크기 자동 조정 기능.

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