



2 메타, MTIA 4세대 공개, 빅테크 AI 칩 자립 경쟁 가속

→ 메타 MTIA 순차 배포 확정, 추론 인프라 자립 구조 본격화

- 메타, 자체 개발 AI 반도체 'MTIA' 시리즈 배포 로드맵 공식화
 - 메타가 2026년 3월 자체 AI 칩 'MTIA(Meta Training and Inference Accelerator)' 4종 (MTIA 300·400·450·500)을 공개하며, 2026~2027년까지 데이터센터 순차 배포 계획 발표
 - MTIA 300은 페이스북·인스타그램 R&R(추천·광고랭킹) 시스템에 이미 양산 투입되었으며, 2027년 말까지 MTIA 400(Iris), 450(Arke), 500(Astrid)을 6개월 주기로 순차 출시할 예정
 - 메타는 AI 모델 학습에는 엔비디아·AMD 등 외부 칩을 활용하고 서비스 운영 단계의 추론 작업은 MTIA로 처리하는 '이원화 전략'으로 비용·성능·공급 안정성 3축을 동시 관리하는 구조를 구축
- AI 추론 비용폭증과 공급망 리스크가 자체 칩 설계 촉진
 - AI 서비스 전체 운영 비용의 약 70~80%는 '추론' 단계에서 발생하는 가운데, 범용 GPU는 추론에 비효율적 구조로 설계되어 있어 대규모 서비스 운영 시 비용 부담이 급격히 증대
 - 엔비디아 GPU 수요 폭증으로 공급 부족과 가격 상승이 지속, 단일 공급사 의존에 따른 수급 리스크가 전략적 위협 수준으로 증대됨에 따라 공급망 다양화 필요성이 임계점에 도달
 - 메타의 2025년 설비투자(CAPEX) 규모가 역대 최대치에 달하면서, 자체 칩 설계를 통한 추론 비용 절감·공급망 다양화·가격 협상 레버리지 확보가 복합적 전략 목적으로 부상

→ 메타의 탈(脫) 엔비디아 승부수, 차세대 MTIA 진화 로드맵과 3대 설계 원칙

- (추론 우선) 메타 서비스에 특화된 '추론 우선' 설계로 비용 효율 극대화
 - 범용 GPU가 대규모 사전 학습을 1순위로 설계된 것과 달리, MTIA 450·500은 GenAI 추론을 최우선으로 설계한 후 R&R 학습·추론 등 다른 워크로드에 확장 적용하는 '추론 우선' 철학 채택
 - 추론 성능의 핵심 병목인 메모리 대역폭과 저장밀도 연산에 집중 투자하여 MTIA 500까지 기존 MTIA 300 대비 HBM 대역폭 4.5배 확대, MX4 연산 성능 25배 이상 향상 전망

- MX4·MX8 등 추론에 최적화된 저정밀 데이터 타입을 각 세대에 맞게 공동 설계하여, 연산 정밀도 손실을 최소화하면서 처리 속도와 에너지 효율을 동시에 극대화
- (고속 대응) 모듈형 칩렛 설계로 AI 모델 진화 속도에 실시간 대응
 - 메타는 컴퓨터·네트워크·HBM 칩렛을 기능별로 분리 설계해 각 모듈을 개별 교체·업그레이드할 수 있는 모듈형 칩렛 구조를 도입, 이를 통해 업계 평균 대비 절반 수준인 6개월 이하 출시 주기 실현
 - MTIA 400·450·500은 동일 세시·랙·네트워크 인프라를 공유하는 드롭인 (Drop-in) 호환 구조를 채택, 신세대 칩을 기존 데이터센터에 즉시 투입 가능하도록 설계함으로써 배포 속도 극대화
 - AI 모델 발전 속도가 기존 반도체 개발 주기를 초과하는 상황에서, 메타는 단일 설계에 장기간 의존하기보다 반복적 개선 방식을 통해 최신 AI 기술 변화에 신속 대응하는 전략을 채택
- (도입 편의) 업계 표준 소프트웨어 기반으로 마찰 없는 도입 환경 구현
 - PyTorch, vLLM, Triton 등 업계 표준 소프트웨어와 호환 설계로 기존 GPU 기반 모델을 MTIA 전용 코드 수정 없이 이식 가능하도록 구현함으로써 개발자 진입 장벽 및 전환 비용 최소화
 - OCP(Open Compute Project) 표준을 준수하여 기존 데이터센터에 즉시 배포 가능한 환경을 구축하였으며, 에이전트 기반 AI 시스템을 활용한 커널 자동 생성·최적화 체계도 함께 구현
 - 메타가 주도해 설립한 PyTorch 생태계를 MTIA에 직접 활용하는 구조로, 자사 소프트웨어 자산과 하드웨어 전략을 긴밀히 연동한 수직 통합형 AI 인프라 체계 구축
- ➔ 빅테크 전반으로 확산되는 커스텀 AI 칩 경쟁
 - 메타·MS까지 합류하며 AI 반도체 자립화 경쟁이 전면 확산
 - 구글은 2015년 AI 전용 TPU를 처음 도입하며 ASIC 선구자로 자리잡았으며, 이후 자사 클라우드(GCP) 고객에게도 개방해 외부 수익화까지 병행하는 전략으로 ASIC 투자의 경제적 타당성을 업계에 입증
 - 아마존은 2018년 자체 칩 발표 이후 AWS 데이터센터용 트레이니엄(학습)·인퍼렌시아(추론)를 분리 개발하며 학습·추론 분리 전략의 선례를 제시, 이후 빅테크들의 이원화 전략 확산에 영향



- 메타·MS까지 ASIC 개발에 합류하면서 AI 반도체 시장은 엔비디아 GPU 중심의 단일 공급 구조에서 빅테크가 직접 설계에 참여하는 다극 경쟁 구도로 빠르게 재편되는 양상
- 구글·MS·메타, 칩 활용 범위·설계 철학·생태계 전략에서 뚜렷한 차별화 보유
 - 세 기업 모두 TSMC 위탁생산과 PyTorch·Triton 등 오픈 소프트웨어 생태계 활용이라는 공통 전략을 공유하나, 칩 개발 주기와 외부 고객 제공 여부에서 차별점이 존재
 - (구글 TPU v7 Ironwood) TPU v5p 대비 최대 10배 성능을 구현하며 최대 9,216개 칩 슈퍼포드로 확장 가능한 구조를 채택하였으며, 엔트로픽 등 외부 AI 기업도 활용 가능
 - (MS Maia 200) TSMC 3nm 공정으로 제작되어 FP4 성능 10+ 페타플롭스를 구현하고 AWS Trainium 3 대비 FP4 성능 3배 우위 달성, GPT-5.2·Microsoft Copilot 등 핵심 서비스에 직접 투입되고 Azure 클라우드 고객 개방도 예정

➔ 빅테크 ASIC 확산이 AI 반도체 공급망에 미치는 영향

- 단순 ‘칩 수요자’에서 ‘설계자’로 진화하는 빅테크
 - 메타는 미국 루이지애나·오하이오·인디애나 등에 초대형 데이터센터를 건설 중이며, 30개 운영·계획 데이터센터 중 26개를 미국 내에 배치하며 AI 인프라 자립 기반을 전면 구축
 - 엔비디아·AMD 등 외부 공급사에 전적으로 의존하던 빅테크가 직접 반도체를 설계·배포하면서, AI 인프라의 핵심 권력이 칩 공급사에서 서비스 기업으로 이동하는 구조적 전환이 가속화
 - TSMC(생산)와 브로드컴(설계)이 빅테크의 맞춤형 AI 칩 개발을 전폭 지원 하면서, 기존 엔비디아 독점 체제를 허물고 ‘빅테크-파운드리’ 연합 중심의 생태계 재편이 가속화
- HBM 고객 다변화로 메모리 기업의 협상력도 동반 강화
 - 빅테크가 칩 설계·공급망 통제권을 직접 확보하면서, 자체 칩 설계 역량 자체가 AI 서비스 경쟁력을 좌우하는 핵심 변수로 부상하고 AI 인프라의 주도권이 빅테크로 본격 이동
 - 메타는 MTIA 시리즈부터 HBM을 신규 도입해 MTIA 500 기준 HBM 용량을 최대 512GB까지 확대할 예정이며, 글로벌 메모리 기업들의 새로운 주요 고객으로 부상

- 엔비디아 독점 공급 구조에서 구글·MS·메타·아마존 등 복수 빅테크로 HBM 수요처가 분산되며, 메모리 기업들의 가격 협상력 및 수주 안정성 동반 강화
- 시장조사 업체 가트너는 AI 가속기 시장의 분절화가 메모리 기업들에게 고객 사별 맞춤형 HBM이라는 새로운 고부가가치 시장을 열어줄 것으로 분석하며, 한국 메모리 산업의 영향력 확대를 전망
- AI 인프라 주도권 경쟁이 모델·칩·공급망 전면전으로 확대
 - 엔비디아는 ASIC 확산에 맞서 오픈 AI 모델 ‘네모트론3 슈퍼’를 공개하고 자사 블랙웰 GPU에 최적화하는 전략으로 맞대응, AI 모델 회사가 칩을 만들고 칩 회사가 AI 모델을 만드는 ‘교차 경쟁’ 구도가 본격화
 - 엔비디아는 블랙웰 GPU 기반 추론 속도를 호퍼 대비 최대 4배 향상하고 AI 클라우드 기업 네비우스에 20억 달러를 투자하며, 칩 판매를 넘어 인프라 플랫폼 기업으로의 영역 확장을 본격화
 - 빅테크의 이원화 전략(자체 칩+외부 GPU 병행)은 추론 비용 절감과 공급망 협상 레버리지를 동시에 부여하는 구조로, AI 인프라 주도권 경쟁이 모델·칩·공급망을 아우르는 전면전으로 확대되는 양상

출처: 한국일보 외(2026.32.)

<https://www.hankookilbo.com/news/article/A2026031209030004604?did=NA>
https://www.choicestock.co.kr/stock/news_view/115889?bu=
<https://www.aitimes.kr/news/articleView.html?idxno=39047>
<http://www.dailysisa.com/news/articleView.html?idxno=51582>
<https://blogs.microsoft.com/blog/2026/01/26/maia-200-the-ai-accelerator-built-for-inference/>
<https://www.aitimes.kr/news/articleView.html?idxno=37113>
<https://about.fb.com/news/2026/03/expanding-metas-custom-silicon-to-power-our-ai-workloads/>
<https://ai.meta.com/blog/meta-mtia-scale-ai-chips-for-billions/>
<https://wccftech.com/meta-sprays-out-four-mtia-ai-chips-in-two-years/>
<https://www.techinasia.com/news/meta-expands-in-house-ai-chips>
<https://www.wired.com/story/meta-unveils-four-new-chips-to-power-its-ai-and-recommendation-systems/>
<https://www.cnbc.com/2026/03/11/meta-ai-mtia-chip-data-center.html>