

# 081 AI 가속기

AI Accelerator

## AI 연산을 고속·고효율로 처리하는 전용 장치

- GPU·NPU·TPU 등을 묶어 부르는 말로, 많은 계산을 동시에 처리해 대규모 학습과 실시간 추론을 빠르게 수행하도록 설계된 핵심 하드웨어 구성
- 데이터 센터와 에지에서 전력·냉각 체계와 함께 운용되는 AI 서비스 기반 요소

### ● AI 가속기의 개념

AI 가속기는 여러 계산을 동시에 수행하도록 설계된 장치로, 대규모 연산을 효율적으로 처리해 학습 속도와 응답 성능을 높입니다. 중앙처리장치(CPU)가 순차적으로 작업을 수행한다면, 가속기는 수많은 연산을 병렬로 수행해 AI 학습과 추론을 가속합니다. 대표적인 형태로는 GPU, NPU, TPU가 있으며, 성능은 연산 칩뿐 아니라 메모리, 전력, 냉각, 네트워크가 얼마나 효율적으로 결합되는지에 따라 달라집니다. 즉, AI 가속기는 연산 코어와 메모리, 연결망, 전력 공급, 냉각 체계가 하나로 통합된 시스템 단위 장치로 이해할 수 있습니다.

### ● AI 가속기와 AI 반도체의 차이

AI 반도체는 실제 연산을 수행하는 칩 수준의 부품이고, AI 가속기는 그 반도체를 탑재해 작동하도록 만든 장치 수준의 구성체입니다. 가속기 내부에는 반도체 칩 외에도 고속 메모리(HBM·D램), 전력 공급 장치, 냉각 시스템, 네트워크 연결부, 제어용 소프트웨어가 함께 포함됩니다. 반도체가 성능의 '엔진'이라면, 가속기는 그 엔진이 안정적으로 작동하도록 돋는 '전체 장치'입니다. 따라서 반도체는 회로 구조와 처리 효율 같은 기술 사양이 중심이지만, 가속기는 실제 환경에서의 운용 효율과 안정성이 핵심입니다.

### ● AI 가속기의 중요성

AI 가속기는 대규모 모델 학습과 실시간 서비스 운영을 가능하게 하는 핵심 장치로, 학습 단계에서는 방대한 데이터를 병렬로 처리해 시간을 단축하고, 서비스 단계에서는 빠른 응답 속도로 사용자 경험을 향상시킵니다. 또한 전력 효율과 냉각 성능은 운영비와 환경 영향을 결정하는 중요한 요소입니다. 효율적인 인프라를 구축한 기관일수록 안정성과 지속 가능성이 높으며, 에지 단말에서는 제한된 자원에서도 성능을 유지하는 경량형 NPU가 활용되어 클라우드와 현장이 유기적으로 연결된 AI 생태계를 형성하고 있습니다.



출처: 조선일보