

049 온디바이스 AI

On-Device AI

클라우드 대신 단말기 내부에서 직접 AI 연산을 수행하는 기술

- 스마트폰, IoT 기기 등 사용자 단말에서 AI 모델을 실행해 실시간 판단·예측·생성을 수행하는 기술
- 데이터를 외부로 전송하지 않아 속도와 보안성을 동시에 확보

● 온디바이스 AI의 개념

온디바이스 AI는 인공지능 모델이 클라우드 서버가 아닌 스마트폰·IoT 기기·웨어러블·차량 등 단말기 내부에서 직접 연산을 수행하는 기술을 의미합니다. 기존 AI 서비스가 데이터를 서버로 전송해 분석했다면, 온디바이스 AI는 이 과정을 기기 안에서 처리해 통신 지연을 최소화하고 개인정보 유출 위험을 줄입니다. 스마트폰의 음성 인식, 카메라 피사체 인식, 자율주행 보조 시스템, 웨어러블의 건강 분석 등에서 활용됩니다. 즉, 클라우드 중심의 '중앙집중형 AI'를 넘어, 기기 자체를 지능화된 연산 주체로 전환하는 기술입니다.

● 온디바이스 AI가 주목받는 이유

온디바이스 AI가 주목받는 이유는 실시간성·보안성·개인화를 모두 강화했기 때문입니다. 데이터가 기기 내부에서 처리되어 응답 속도가 빠르고, 네트워크가 불안정한 환경에서도 안정적으로 작동합니다. 외부 서버 전송이 필요 없어 개인정보 보호에도 유리하며, 사용자의 로컬 데이터를 기반으로 맞춤형 서비스를 구현할 수 있습니다. 이러한 특성 덕분에 제조사와 반도체 기업들은 NPU를 내장하고 모델을 기기에서 직접 실행하는 방향으로 전환하고 있습니다. 또한 경량화 모델과 저전력 연산 기술의 발전으로 작은 기기에서도 고성능 AI를 구현할 수 있게 되었습니다.

● 온디바이스 AI의 한계

온디바이스 AI는 단말기의 저장 용량과 연산 능력의 한계로 제약이 존재합니다. LLM이나 복잡한 딥러닝 모델은 기기에 직접 탑재하기 어렵고, 정기적인 업데이트가 필요합니다. 또한 기기 간 사양 차이로 동일한 모델이라도 성능 편차가 생기며, 최적화 수준에 따라 속도나 정확도가 달라질 수 있습니다. 제조사와 운영체제 간의 호환성 문제 역시 해결해야 할 과제입니다.

● 온디바이스 AI의 전망

온디바이스 AI는 AI의 탈중앙화와 개인화 시대를 여는 핵심 기술로 평가됩니다. 앞으로는 클라우드와 단말이 협력하는 하이브리드 AI 구조가 확산되며, 중앙 서버의 계산력과 기기 내부의 실시간 처리가 결합된 생태계가 구축될 것입니다. 또한 초저전력 반도체, 경량화언어모델, 연합학습 기술이 발전하면서 개인 데이터를 로컬에서 안전하게 학습·활용하는 프라이버시 중심의 AI 환경이 보편화될 것으로 전망됩니다.

관련 용어

에지 AI (Edge AI)

에지 AI는 데이터가 생성되는 지점(에지)에서 AI 연산을 수행하는 기술을 말합니다. 클라우드로 데이터를 전송하지 않고, 가까운 네트워크 단이나 기기에서 직접 분석해 응답 속도를 높이고 대역폭 부담을 줄입니다. 예를 들어 공장 센서, CCTV, 자율주행차 카메라 등에서 수집된 정보를 현장에서 즉시 판단하는 방식입니다. 이는 온디바이스 AI보다 범위가 넓은 개념으로, 개별 기기뿐 아니라 게이트웨이·로컬 서버 등 인접 장비까지 포함합니다. 에지 AI는 실시간성·보안성·네트워크 효율을 동시에 확보할 수 있어, 산업 자동화와 IoT 시대의 핵심 인프라로 주목받고 있습니다.

관련 용어

임베디드 AI (Embedded AI)

임베디드 AI는 AI 알고리즘을 기기 내부의 전자회로나 하드웨어에 직접 탑재해 작동시키는 형태를 의미합니다. 주로 마이크로컨트롤러(MCU)나 전용 칩셋에 AI 모델을 내장하여, 별도의 네트워크 연결 없이도 데이터 인식·분석이 가능합니다. 예를 들어 카메라가 자동으로 얼굴을 인식하거나 가전제품이 사용 패턴을 스스로 학습하는 기능이 이에 해당합니다. 임베디드 AI는 하드웨어에 최적화된 초경량 모델을 사용하기 때문에 연산 속도가 빠르고 전력 소비가 적습니다. 온디바이스 AI보다 하드웨어 종속성이 강하며, 제한된 환경에서도 작동하는 초소형·고효율 AI 기술로 평가됩니다.

관련 용어

클라우드 AI (Cloud AI)

클라우드 AI는 대규모 서버나 데이터센터에서 AI 모델을 구동하고, 네트워크를 통해 서비스를 제공하는 구조입니다. 기기에서 데이터를 수집해 중앙 서버로 전송하고, 거기서 연산·분석·추론을 수행한 뒤 결과를 다시 전달합니다. 고성능 GPU, LLM, 방대한 데이터가 필요한 AI 서비스는 대부분 이 구조를 기반으로 합니다. 클라우드 AI의 강점은 연산 능력과 확장성이 뛰어나다는 점이지만, 네트워크 지연과 개인정보 유출 가능성의 단점으로 꼽힙니다. 온디바이스·에지 AI가 이러한 한계를 보완하는 형태로 발전하고 있습니다.

관련 용어

하이브리드 AI (Hybrid AI)

하이브리드 AI는 클라우드와 온디바이스·에지 AI의 장점을 결합한 협력형 구조입니다. 중앙 서버가 복잡한 연산과 대규모 학습을 담당하고, 단말기나 에지 기기가 실시간 분석과 즉각적인 응답을 처리합니다. 예를 들어 스마트폰 음성 인식 같은 간단한 명령은 기기 내부에서 처리하고, 복잡한 질의는 클라우드로 전송해 고급 연산을 수행하는 식입니다. 이를 통해 속도·보안·정확성을 모두 확보할 수 있습니다. 하이브리드 AI는 AI의 분산 처리 구조를 완성하는 모델로, 향후 지능형 네트워크와 협력형 학습의 핵심 기술로 주목받고 있습니다.