

# 098 AI 추론(Inference)

AI Inference

## 학습된 AI 모델이 입력 데이터를 기반으로 결과를 계산·출력하는 과정

- 훈련을 통해 학습된 모델을 기반으로 새로운 데이터를 분석하고 예측·생성하는 AI의 실행 단계
- AI 서비스에서 응답 속도와 정확성을 결정짓는 핵심 기술로, 실제 작동 성능을 좌우하는 과정

### ● AI Inference 개요

AI Inference는 학습이 완료된 AI 모델이 실제로 데이터를 받아 결과를 도출하는 실행 단계를 말합니다. 예를 들어 이미지 인식 모델이 사진을 보고 사물을 식별하거나, 언어모델이 사용자의 질문에 답변을 생성하는 모든 과정이 Inference에 해당합니다. 이 단계에서 AI는 학습 중에 형성한 가중치(weight)와 패턴을 활용해 입력 데이터를 분석하고, 가장 가능성 높은 출력을 계산합니다. 따라서 AI Inference는 AI가 실제 서비스를 제공하는 순간이자, 모델의 성능·응답 속도·비용 효율을 직접 결정하는 기술적 핵심이라 할 수 있습니다. 이 때문에 생성형 AI 확산 이후, 텍스트·이미지·음성 등 대규모 데이터를 실시간으로 처리해야 하는 환경에서 저지연·고효율·고신뢰성 추론이 중요해지고 있습니다.

### ● AI Inference의 작동 방식

AI Inference는 입력 데이터를 받아 순전파를 포함한 토큰 생성, attention 계산 등 모델 구조에 따라 다양한 계산 과정을 수행하여 결과를 출력합니다. 이는 학습 중 구축된 신경망 구조를 그대로 사용하되, 가중치를 수정하지 않고 예측만 수행하는 방식입니다. 대형 모델일수록 연산량이 많기 때문에, 빠르고 효율적인 추론을 위해 GPU, TPU, NPU 같은 고성능 연산 장치가 필수적으로 사용됩니다. 또한 추론 효율을 높이기 위해 모델 경량화, 양자화, 배치 처리 등 다양한 최적화 기법이 적용됩니다. 예를 들어 모바일 기기나 에지 환경에서는 경량화된 모델을 통해 연산 속도와 전력 효율을 높이는 것이 중요합니다. 클라우드 환경에서는 여러 요청을 동시에 처리하는 병렬 추론 구조가 사용됩니다.

### ● AI Reasoning & AI Inference

AI Reasoning과 AI Inference는 국어로 모두 'AI 추론'으로 번역되지만, AI가 정보를 처리해 결과를 산출하는 과정에서 서로 다른 역할을 수행합니다. Reasoning은 주어진 정보와 규칙을 기반으로 추론 경로와 단계적 근거를 구성하는 과정으로, 결과가 어떻게 도출되었는지 설명할 수 있는 구조를 만들어냅니다. 반면 Inference는 학습된 모델이 입력 데이터를 받아 최종 출력을 계산하는 실행 단계로, 학습된 지식을 실제 문제 해결에 적용합니다. Reasoning이 추론 과정의 정교함과 근거 생성에 중점을 둔다면, Inference는 이러한 추론 결과를 빠르고 효율적으로 산출하는 계산 효율성에 초점을 둡니다. 두 과정은 역할은 다르지만 상호 보완적으로 기능하며, Reasoning이 고도화될수록 결과의 일관성과 설명 가능성이 높아지고, Inference가 최적화될수록 실제 서비스 환경에서 그 결과를 더 신속하고 안정적으로 제공할 수 있습니다.