

# 085 AI 레드티밍

AI Red Teaming

## AI의 취약점을 공격자 관점에서 시험해 위험을 식별하는 검증 절차

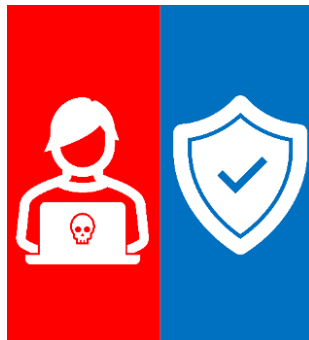
- 비정상 입력·악의적 프롬프트·사회적 편향 등 실제 위협 상황을 모의해, AI 시스템의 안전성과 신뢰성을 사전에 평가하는 점검 활동
- 단순 성능 시험이 아닌, AI의 예상 밖 행동과 사회적 영향을 탐지·완화하기 위한 선제적 검증 체계

### AI 레드티밍의 개요

AI 레드티밍은 AI 시스템의 위험 요인과 취약점을 사전에 식별·검증해 안전성과 신뢰성을 확보하는 과정입니다. 원래 군사·보안 분야에서 '적의 시각으로 방어 체계를 점검한다'는 의미로 쓰이던 레드팀(red team) 개념을 AI 안전 관리에 적용한 것입니다. 이 과정에서는 공격자 관점에서 AI 모델을 시험해 비정상 입력이나 악의적 명령에 대한 반응을 분석하고, 편향된 응답, 유해 콘텐츠, 정보 왜곡, 보안 우회 등의 문제를 찾아냅니다. 즉, AI 레드티밍은 단순한 오류 탐색이 아니라 AI의 윤리·보안·신뢰성을 검증하는 핵심 절차입니다.

### AI 레드티밍 수행 방식

AI 레드티밍은 보통 공격적 시험, 시나리오 검증, 위험 분석의 세 단계로 진행됩니다. 공격적 시험은 비정상 입력이나 우회 프롬프트를 주어 오작동 여부를 점검하고, 시나리오 검증은 실제 사용 환경을 모의해 편향된 출력이나 부적절한 응답을 평가합니다. 이후 위험 분석 단계에서는 탐지된 문제를 분류하고 대응 정책, 데이터 필터링 등 개선책을 마련합니다. 이러한 검증은 내부 보안팀 외에도 외부 전문가나 독립 기관이 참여해 객관성을 확보하며, 최근에는 정부나 민간의 AI 안전성 평가 제도와 연계해 정기적 검증 체계로 확산되고 있습니다.



### AI 레드티밍의 중요성

AI 레드티밍은 생성형 AI 확산에 따라 필수적인 신뢰성 검증 절차로 자리 잡고 있습니다. 생성형 AI 모델 내부의 작동 원리를 완전히 통제하기 어렵고, 예측 불가능한 출력이 사회적 문제로 이어질 수 있기 때문에, 기술적 안전성과 사회적·윤리적 영향을 함께 평가해야 합니다. 결국 AI 레드티밍은 안전하고 책임 있는 AI 구현을 위한 선제적 관리 체계이자, 기술 혁신과 사회적 신뢰를 함께 확보하는 핵심 기반입니다.