



AI 학습과 뉴스 저작권 간의 '공정이용' 갈등 [KISO(한국인터넷자율정책기구) JOURNAL 제54호 게재]

1. AI의 발전과 저작권계의 대응

가. 생성형 AI와 저작권 침해

생성형 AI(Generative AI)는 이용자가 입력한 특정한 요구(프롬프트)에 따라 AI가 다양한 수준의 자율성(Autonomy)을 가지고 콘텐츠를 생성하는 인공지능 기술을 말한다¹. 이러한 개념 정의의 중요 요소로는 생성형 AI의 목적이 콘텐츠를 표현하고 생성하는 것에 있다는 점, 콘텐츠를 생성할 때 AI에 상당한 정도의 자율성이 있다는 점(동일한 프롬프트를 입력하더라도 다른 결과가 산출될 수 있음), 콘텐츠 생성은 이용자의 프롬프트에 의해 지시·유도된다는 점(책임 관계에서 이용자 역할이 부각됨)이라고 생각된다. ChatGPT, Stable Diffusion과 같은 생성형 AI는 콘텐츠를 생성하는 표현력에 집중한다는 점에서 데이터를 분류하거나 범주화 하도록 설계된 판별형 AI(Discriminative AI)와 구분된다. 통계학적 용어를 빌면, 생성형 AI는 주어진 원본 데이터를 학습하여 확률적으로 원본 데이터와 구별되지 않거나(GAN), 원본 데이터의 분포를 따르는(Diffusion, VAE), 유사한(가짜) 데이터를 생성하는 인공지능 모델이며 따라서 원본 데이터의 분포를 파악하여 학습하는 것이 가장 중요하게 다루어진다. 이와 같은 학습데이터의 분포를 모방하려는 기술적 특성으로 인하여 생성형 AI의 산출물은 근원적으로 타인의 저작물과 유사할 수 있다. 이러한 특성으로 인하여 판별형 AI와 달리 생성형 AI에서는 산출물 생성 단계(OUTPUT 단계)의 저작권 침해 문제가 심각하게 논의될 수밖에 없다.

나. 'TDM 예외' 논의의 후퇴

또한 기계학습에 의한 AI 개발은 원천적으로 AI 학습 단계(INPUT 단계)에서의 저작권 침해의 이슈를 내재하고 있다. 기계학습은 학습데이터의 복제를 수반하기 때문에 제3자의 저작물이 적법한 권한 없이 기계학습 과정에서 복제될 수 있다. 통상 학습데이터는 인터넷 등 공개된 공간에서 수집되는데 인터넷상에 존재하는 대규모의 저작물에 대해 일일이 저작권자의 동의를 얻는 것은 사실상 불가능하다. 이러한 사정으로 인하여 세계 각국은 AI의 기계학습을 포함하는 TDM(Text and Data Mining) 행위에 대해 저작권법상 면책 요건을 부여하는 TDM 예외 입법을 논의하고 있다. 실제 EU, 일본, 영국 등은 TDM 예외 조항을 입법했고 미국은 판례법에 근거한 공정이용 법리에 따라 TDM 내지 기계학습을 허용하려는 움직임을 보였다.

한국 역시 2021년 저작권법 전면 개정안에 TDM 예외 조항을 도입하면서 이러한 국제적 흐름에 따르고 있었다. 그러나 2022년 11월에 등장한 OpenAI의 ChatGPT가 생성형 AI의 놀라운 표현력을 보여 주면서 이러한 흐름은 바뀌게 되었다. ChatGPT와 같은 대규모 언어모델, Stable Diffusion과 같은 이미지 생성 AI가 관련 저작물 시장에서 저작권자와 직접적으로 경쟁할 수 있다는 점이 확인됨에 따라 저작권계는 저작권자에게 보상 없는 TDM 내지 기계학습에 반대하는 입장을 표명하게 된 것이다.

다. 관련 소송의 발생

저작권계의 반발은 AI 개발사에 대한 소송으로 이어졌다. 2022년 11월 미국에서 GitHub Copilot의 코드 자동 생성 서비스에 대해 오픈소스 라이선스 위반을 주된 이유로 하는 집단소송(class action)이 발생한 이후, ChatGPT, LLaMA, Claude, Stable Diffusion, Midjourney 등 대중에게 알려진 거의 모든 AI 서비스에 대해 저작권·프라이버시 침해 소송이 진행되고 있다. 이 글에서 논하고자 하는 뉴욕타임스 v. Open AI, Microsoft 소송도 그 중의 하나이다. 처음에는 개별 저작권자가 집단소송의 형식으로 대응하던 것이 최근에는 게이티이미지, Authors Guild(미국의 작가 협회 격), 뉴욕타임스 등의 대형 저작권자가 소를 제기하고 있다. 이는 AI 소송이 격화되고 있다는 경향을 보여주는 것과 동시에 저작권 계가 소송을 통해 AI 업계와 라이선스를 시도하는 것으로도 볼 수 있다.

2. NYT vs OpenAI, MS 소송

가. 본 소송의 의의: 독립 저널리즘과 신기술의 대립

¹ 최근 EU 의회를 최종 통과한 EU AI Act 제3조 제1항은 'AI 시스템'이란 다양한 수준의 자율성을 가지고 작동하도록 설계된 기계 기반 시스템으로서, 배포 후 적응력을 발휘할 수 있으며 명시적 또는 암묵적 목적을 위해 수신한 입력으로부터 물리적 또는 가상 환경에 영향을 미칠 수 있는 예측, 콘텐츠, 추천 또는 결정과 같은 출력을 생성하는 방법을 추론하는 것을 의미한다고 정의한다.

뉴욕타임스의 소장 첫머리는 “독립 저널리즘은 민주주의에서 필수적이다. 피고들이 불법적으로 원고의 저작물을 사용하여 원고와 경쟁하는 AI 제품을 만드는 것은 원고의 독립 저널리즘을 위협하는 행위이다.”라는 문장으로 시작한다. 피고들의 ChatGPT, BingChat이 원고의 기사를 거의 그대로 생성해 내고 특히 원고가 운영하는 상품 리뷰 사이트인 와이어커터(Wirecutter)의 상품 추천 리스트를 그대로 복제하여 AI 이용자들에게 보여주는 행위는 원고의 재정적인 능력을 손상시켜 결국 민주주의의 훼손으로 이어진다는 논리이다. 이에 대해 피고 OpenAI는 최근 제출한 기각신청(motion to dismiss)에서 ChatGPT와 같은 AI는 인류의 능력을 증진시킬 수 있는 혁명적인 기술이라고 응답하고 있다. 종래 VTR, MP3 등 신기술의 등장에 따른 저작권 분쟁은 음반사업자, 방송사업자 등 사업상의 기득권 체계에 대해 신기술 사업자가 도전한 것이라고 한다면 위 소송은 원·피고 모두 민주주의, 인류의 발전이라는 상위 가치를 전면에 등장시킨 점이 이채롭다. 원고의 청구가 금전 손해배상도 구하고 있는 점에서 본 소송이 서로의 사업적 이익을 추구하는 분쟁임을 벗어날 수는 없겠으나 본 소송이 가지는 인류 사회적 가치도 주목해야 할 필요가 있다. 원고는 뉴욕타임스의 기사가 포함된 데이터셋과 그에 기초한 GPT 모델의 폐기를 청구하고 있는데 AI가 국가 안보적 가치마저 획득하고 있는 현 상황에서 설령 원고의 저작권침해 주장이 맞다고 하더라도 이미 존재하고 있는 AI 모델까지 폐기하기는 현실적으로 어렵다.

나. 원고의 주장

본 소송은 뉴욕타임스가 2023년 12월 OpenAI와 마이크로소프트를 상대로 ChatGPT, Bing Chat, Microsoft 365 Copilot 제품과 관련해 뉴욕 남부지방법원에 제기한 분쟁이다.² 원고는 AI 학습 단계(INPUT 단계)에서 이루어진 무단적인 저작권 복제와 AI 산출단계(OUTPUT 단계)에서 발생하는 암기(memorization) 또는 역류(regurgitation) 현상과 환각(hallucination) 현상을 비판한다. 암기 내지 역류 현상은 AI가 콘텐츠를 생성하는 과정에서 학습데이터 원본과 거의 유사한 사본을 출력하는 현상을 말하는데 이는 저작권법적으로 복제에 해당할 여지가 있다. 그리고 환각 현상은 원고 뉴스가 아닌 내용을 마치 원고가 작성한 것처럼 이용자에게 보여 주는 것이므로 환각 현상에 의해 원고의 명성 내지 브랜드 가치가 훼손될 우려가 있다는 것이다.

이를 입증하기 위해 원고는 2,000페이지에 달하는 증거 69개를 제시했고 특히 증거 J(Exhibit J)에서 피고 AI 모델의 암기 현상을 보여주는 100개의 사례를 기재했다. 위 증거 J에서 제시된 암기 사례는 이 사건에서 매우 중요한 가치를 가진다고 생각된다. 위 100개의 사례는 ChatGPT 등에 프롬프트로서 특정 뉴스 기사의 제목, 첫 문단 등의 짧은 부분(short snippet)을 입력하면 해당 뉴스 기사의 그 다음 문단이 거의 그대로 출력되는 현상을 캡처한 것이다. 피고들의 AI 서비스가 이와 같은 암기 결과를 보여줌에 따라 뉴스 이용자들이 원고 사이트에 접속하지 않아도 뉴스 기사를 볼 수 있게 돼 결국 원고의 구독료, 광고료, 라이선스료, 제휴 마케팅 수수료를 상실하게 되었다는 것이 원고의 주장이다.

원고는 소장 청구원인으로서 저작권 직접 침해, 저작권 대위 침해, 저작권 기여 침해, DMCA 위반(CMI 무단 삭제), 부정경쟁행위, 상표 희석화를 주장하고 있다.

다. 피고들의 기각신청(motion to dismiss)

이에 대해 피고들은 기각신청에서 이러한 암기 현상은 비의도적인 결과라고 항변한다. 기계학습의 원리상 AI 모델은 내부에 학습데이터를 저장하지 않는다. 암기 현상은 특이한 상황에서 정교한 파라미터 값에 의해 발생하는 우연적인 결과라는 것이다. 사실 증거 J에서 제시되고 있는 암기 현상은 일반적인 이용 형태에서 발견될 수 있는 것들은 아니라고 보인다. 예컨대 증거 J의 두 번째 사례는 글자 수(token)가 1,580개에 달하는 뉴스 기사의 제목과 앞부분 기사 내용을 그대로 프롬프트로 입력한 것이므로, 이미 해당 뉴스 기사의 내용을 아는 이용자만이 이러한 결과를 출력할 수 있는 사례로 볼 수 있다. 이러한 점에서 피고들은 원고가 이용자약관을 위반해 일종의 해킹을 하였고 일부러 버그를 유도했다고 비판한다. 기각돼야 하는 청구원인으로는 저작권 직접 침해행위 중 3년 시효가 지난 행위, 기여 침해, DMCA 위반, 부정경쟁행위를 주장하고 있다. 피고들은 기각신청을 통해서 저작권 직접 침해 주장을 이 사건의 청구원인으로 집중시키고 이에 대한 공정이용 항변을 이 사건의 핵심 쟁점으로 부각하려는 전략을 펴고 있는 것으로 보인다.

3. 시사점

한국에서도 네이버 하이퍼클로바X의 뉴스 데이터 학습이 뉴스 이용에 대한 ‘뉴스콘텐츠 제휴 약관’을 위반한 것이라는 주장이 제기된 바 있다. 네이버 뉴스콘텐츠 제휴 약관 제8조 제3항은 “네이버는 서비스 개선, 새로운 서비스 개발을 위한 연구를 위해 직접, 공동으로 또는 제3자에게 위탁하는 방식으로 정보를 이용할 수 있습니다. 단, 제3자에게 위탁하는 방식으로 진행할 경우 사전에 제공자의 동의를 얻어야 합니다.”라고 규정하고 있다. 신문협회는 네이버가 하이퍼클로바X 학습에 뉴스를 사용하는 행위는 위 약관의 목적 범위를 넘는 것이므로 위 조항이 적용될 수 없고 AI 학습을 위해서는 별도의 계약을 필요하다고 보고 있다.

현재 벌어지고 있는 AI 분쟁의 내용은 크게 보면 AI 학습 단계(INPUT 단계)와 AI 산출 단계(OUTPUT 단계)로 나눌 수 있다. 위 하이퍼클로바X에 대한 문제제기는 AI 학습 단계에 관한 것이고 전술한 뉴욕타임스 사건은 AI 학습 단계와 AI 산출 단계를 모두 아우르고 있다. 먼저 AI 산출 단계의 문제를 살펴보면, 이 문제는 AI 학습 단계의 그것보다 비교적 수월한 해결책을 찾을 수 있을 것으로 생각된다. 저작권 침해 여부에 관해서는 종래 법리인 의거성 및 실질적 유사성 판단기준을 거의 그대로 원용할 수 있다고 보이고 무엇보다 현재까지의 증거에 따르면 AI 산출단계에서 발견되는 저작권 침해 현상은 주로 AI 이용자의 침해유발적인 프롬프트에서 기인한 것이기 때문이다. AI 서비스 제공자에게 적절한 필터링 의무를 부여하거나 저작권법상 OSP 책임 규정과 같은 notice and take down 절차를 마련하는 방안을 검토할 수 있다. 반면 AI 학습 단계에서 발생하는 저작권 침해 분쟁은 어느 일방의 완전한 승리 또는 패배에 따라 해결될 성질이 아니고 사회 전반의 합의 내지 입법에 의해 해결될 수밖에 없을 것으로 생각된다. 현재의 국제 상황상 저작권 침해를 이유로 AI 기술 개발을 중단할 수는 없고 다른 한편으로 저작권자의 이해를 도외시키고 AI 개발만을 추구할 수도 없기 때문이다. 공개된 정보(저작물)의 사용에 대한 정책적·법률적 기준 마련에 더욱 힘써야 할 시기이다.

² 사건번호: 1:23-cv-1195

원본 기사: KISO저널 제54호, 국내외 주요 소식([KISO JOURNAL](#))

상기 내용에 관해 문의사항 있으시면 언제든지 저희 법무법인 린의 IP팀 전응준 변호사(Tel. 02-3477-8695)에게 연락 주시기 바랍니다.

[홈페이지](#)

관련 구성원



전응준 변호사
T. 02-3477-8695
E. ejjeon@law-lin.com

법무법인 린의 뉴스레터는 일반적인 정보제공만을 목적으로 발행되므로 이에 수록된 내용은 법무법인 린의 공식적인 견해나 구체적인 사안에 관한 법률의견이 아님을 알려드립니다.

이 메일을 수신 거부하려면 lin-newsletter+unsubscribe@law-lin.com 로 보내주시기 바랍니다.

[More LIN Newsletters](#)

鹿米 此舜 법무법인 린

서울 서초구 서초중앙로 24 길 27, 지파이브센트럴 프라자 326 호
T.02-3477-8695 F.02-3477-8694