

# Web3.0을 기반으로 시를 활용하여 저작권을 보호하는 방안

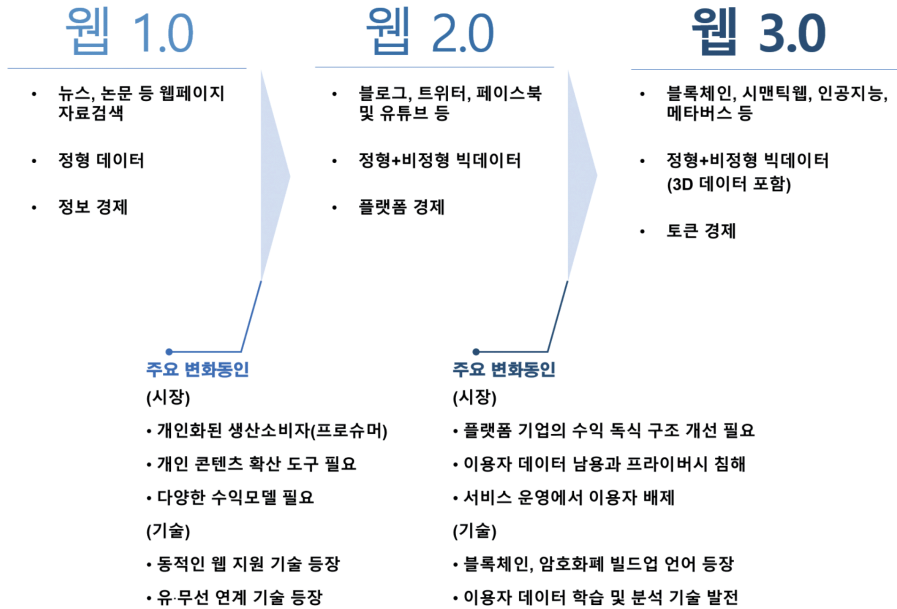
전시형 (서울과학기술대학교)

<p>목 차</p> <p>1. 서 론</p> <p>2. Web3.0 기반 문서 탐색의 이면</p>	<p>3. Web3.0 기반의 표절 적발</p> <p>4. 결 론</p>
--	--

## 1. 서 론

웹3.0은 웹2.0의 한계를 극복하고 사용자의 데이터 주권과 개인화된 정보 제공을 위한 지능형 웹 기술을 말한다. (그림 1)과 같이 웹3.0은 블록체

인, 암호화폐, NFT, 인공지능, 메타버스 등 다양한 기술들을 융합했다. 개방적이고, 안전하고, 탈중앙화되어 있으며, 공유경제를 가능하게 하는 웹을 추구한다[1]. 웹3.0은 컴퓨터가 시맨틱 웹 기술을 이용하여 웹페이지에 담긴 내용을 이해하고 상황



(그림 1) 웹 패러다임 변화[1]

과 맥락에 맞는 개인 맞춤형 정보를 제공하는 지능형 웹 기술을 말하기도 한다[2].

시맨틱 웹의 구성기술인 Web Ontology Language(OWL)는 웹 문서(고유 링크를 가짐)의 특성과 문서 간의 관계를 나타내는 규격으로 기능하기에는 충분하지만, 이전의 기술로는 문서마다 OWL로 문서를 지식화하는 작업이 수고로워서 고비용을 소요할 수밖에 없었다. 이제 AI 기술을 통해 시맨틱웹을 구현하는 데에 돌파구가 열리며 필요한 지식을 찾기는 더욱 쉬워지게 되었으나, 반대급부로 잘 팔리는 분야의 지식을 표절하여 재생산함으로써 원저작자에게 금전적 손실을 입히는 사례가 늘어나고 있다. 블록체인 기술을 통해 저작물의 최초 생산 여부를 증명하는 시스템을 만들 수는 있지만, 단순 복제가 아니라 약간의 변조를 가하는 본격적인 범죄에는 대처하기가 어려운 형국이다. 벌써부터 ‘노아AI 표절 사건’이나 ‘챗지피티와 유튜브 쇼츠로 떼돈 버는 법’같은 기술 악용 사례가 나오고 있다. 본 기고에서는 웹 기반에서 가능한 기술로 저작권을 보호하는 방안을 제언하고자 한다.

## 2. Web3.0 기반 문서 탐색의 이면

웹2.0 플랫폼 사업자들은 사용자를 대규모로 유치하는 웹 환경을 제공하는 대가로 사용자 데이터를 소유했고, 이를 통해 지속적으로 수익을 창출하는 중이다. 오랜 시간 축적해 온 데이터는 거대한 자산이 되었다. 웹2.0의 중앙집중적이고 폐쇄적인 플랫폼은 정작 정보를 생산한 사용자와는 무관하게 플랫폼을 소유한 기업이 무소불위의 권력을 갖게 되었다[2].

웹3.0은 1998년 팀 버너스리가 제안한 개념인 ‘시맨틱 웹(Semantic Web)’이 시초이다. 웹2.0이 해결하지 못하는 문제를 해결하려는 논의를 지속해 오면서 웹3.0은 ‘시맨틱 웹’과 ‘탈중앙화’, 프로토콜 경제를 근간으로 하는 메타버스 개념까지 포함한다[2].

그 중 시맨틱 웹은 불필요한 정보를 제거하며 사용자의 성향과 검색 목적에 따른 맞춤형 정보를 신속하게 도출할 수 있다. 웹2.0의 키워드 검색 기능은 원하는 정보를 찾기 위해 수십 페이지까지 넘겨야 하는 수고로움이 있었지만, 웹3.0은 메타데이터의 집합인 온톨로지(Ontology)를 활용하여 데이터를 개념적으로 연결하고 분석하고자 한다.

〈표 1〉 웹3.0 특징(3+2)[2]

특성		내용	제공가치와 의미
기본 특성	탈중앙화	<ul style="list-style-type: none"> <li>중앙 통제기관(중개자) 없는 거래환경 제공</li> <li>자율적·민주적 운영규칙 결정</li> <li>(주요 기술) 블록체인, DAO, 암호화폐 등</li> </ul>	<ul style="list-style-type: none"> <li>데이터 독과점 극복</li> <li>운영 투명성 및 사용자 권한 강화</li> </ul>
	데이터 소유권	<ul style="list-style-type: none"> <li>중앙 통제기관이 소유하던 데이터를 사용자가 소유</li> <li>(주요 기술) NFT, 암호화폐, Dapp 등</li> </ul>	<ul style="list-style-type: none"> <li>데이터 소유권 증명</li> <li>수익 실현</li> </ul>
	높은 보안성	<ul style="list-style-type: none"> <li>중앙 서버가 필요 없는 데이터 분산저장</li> <li>프로토콜 기반 무보증·무허가 참여</li> <li>(주요 기술) Dapp, 스마트 컨트랙트, 엣지컴퓨팅 등</li> </ul>	<ul style="list-style-type: none"> <li>보안성, 신뢰성 향상</li> <li>참여 가능 대상 확대</li> </ul>
구현 특성	지능화 서비스	<ul style="list-style-type: none"> <li>사용자에게 맞춤형 지능형 서비스 제공</li> <li>(주요 기술) 초거대 AI, 시맨틱웹 등</li> </ul>	<ul style="list-style-type: none"> <li>편의성, 생산성, 효율성 증대</li> </ul>
	확장된 미디어 인터페이스	<ul style="list-style-type: none"> <li>현실세계와 가상세계가 융합된 공간 제공</li> <li>(주요 기술) 메타버스(AR/VR/XR), 라이프 로깅 기술 등</li> </ul>	<ul style="list-style-type: none"> <li>높은 몰입감</li> <li>새로운 사용자 경험</li> </ul>

온톨로지는 시맨틱 웹 언어로, 사물, 사물의 그룹 및 사물 간의 관계에 대한 풍부하고 복잡한 지식을 표현하기 위해 만들었다. 이를 지식의 일관성을 검증하거나 추론을 수행하는 데 사용함으로써 더욱 효과적으로 정보를 검색한다[3].

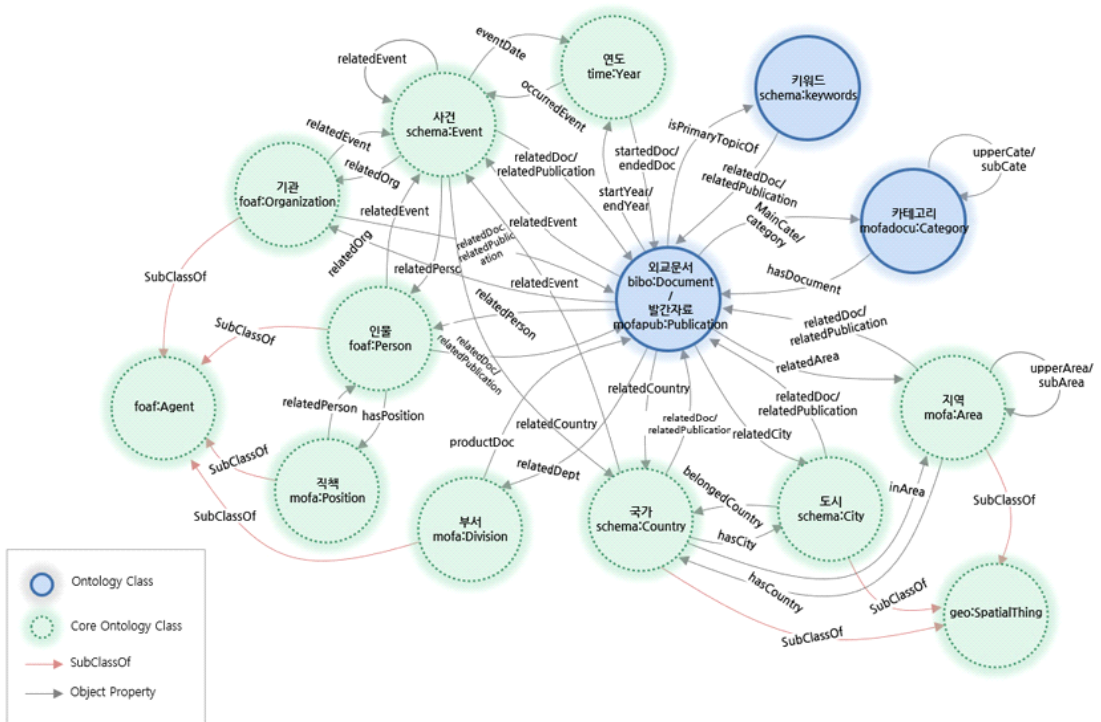
OWL(Web Ontology Language)은 온톨로지를 만들기 위한 지식 표현 언어의 한 계열이다. 2004년도에 주창했음에도 작성에 수고가 많이 들고, HTML의 META tag처럼 오용 및 악용 가능성이 높아 유명세만큼 널리 활용했다고 보기는 힘들었다. 최근 들어 AI 기술이 발전하면서 이를 활용한 OWL 작성에 대한 연구가 이어지고 있다.

눈부시게 뛰어난 성능을 자랑하는 최근의 LLM(Large Language Model) 개발 이전에도 Machine Learning을 활용하여 OWL을 자동으로 생성하는 방법[4], 자연어처리 기법이 발전하면서 작성한 문서를 파싱하여 개념, 속성, 관계 등을 추

출하고 이를 OWL로 변환하는 연구[5]가 있었으며, 딥러닝(Deep learning) 기반의 기계번역 개념을 적용한 사례를 활용하여, 수작업의 의존성이 감소한 방법으로 텍스트로부터 온톨로지를 생성하는 방법을 구현한 연구[6]가 나오기도 했다.

1998년 팀 버너스리가 제안한 개념인 ‘시맨틱 웹(Semantic Web)’이 이제서야 빛을 보기 시작했다고 해도 과언이 아닐 것이다. 다만, 사람뿐만 아니라 기계에게도 쉬워진 검색은 단순 복제와 같은 직접적인 표절을 넘어 아비한 변조 작업을 통해 원저작자에게 극심한 금전적 손실을 입히는 악용의 도구로 활용이 되는 형국이다.

2023년 2월, 온라인 공간을 뜨겁게 달군 이른바 ‘노아AI 표절 논란’의 과정의 다음과 같다. 섬네일 이미지뿐 아니라 유명 유튜브 채널의 영상을 스크립트(대본)까지 그대로 베껴 만든 불법행위로 수익을 창출한 ‘섬네일과 주제 베끼기’를 노하우라



(그림 2) (예시) 온톨로지 Class 연관도(출처: 외교부 OPEN DATA)[2]

고 퍼뜨린 유명인사는 자신의 꿈을 유료 강의로 판매했고, 베낄 만한 썸네일과 영상을 찾아내는 유료 ‘유튜브 데이터 분석 서비스’를 만들어 돈을 벌었다. 이것이 노아AI로, 유튜브 데이터를 기반으로 인기 있는 채널과 영상을 쉽게 검색하고 추출하게 해주는 서비스다. 노아AI를 활용한 이들은 ‘표절 유튜브 영상’을 생산해냈다. 구독자 142만 명을 보유한 과학 유튜버 ‘리뷰잉이’가 피해자로서 유튜브에 게재한 고발 영상으로 인해, 다행히 노아AI 서비스는 폐쇄하고 이를 주도한 유튜버는 잠적했으나[7] 여전히 문제의 불씨는 남았다.

실제로 한 유튜버는 2월 17일 ‘챗지피티와 유튜브 쇼츠로 떼돈 버는 법’이라는 영상을 올리며 AI로 대본뿐 아니라 이미지까지 모두 해결하는 방법을 공개했다. 한국어로 질문을 작성한 뒤 구글 번역기를 이용해 영어로 번역한다. 챗지피티로 그럴듯한 답변을 추출한 뒤 다시 구글 번역기로 한국어 대본을 만든다. DeepAI 사이트를 이용해 영상에 들어갈 이미지를 생성한다. 글과 이미지를 합쳐 쇼츠용 영상을 제작해 배포한다[7]. 최근에는 생성만이 아니라 가공과 변조에 해당하는 요약 목적으로도 ChatGPT 같은 Generative AI를 활용

하기에 표절 콘텐츠를 대량으로 양산하는 게 가능하다. LLM에 논문 전문을 입력하여 생성한 초록의 품질이 꽤 훌륭했다는 연구도 벌써 나왔다.[8] 긴 시간 공을 들여서 콘텐츠를 제작하는 개인 사업자는 대처하기가 너무나도 까다로워졌다.

콘텐츠 산업을 보호하려면 콘텐츠 플랫폼 기업만의 자정작용을 기대해서는 곤란하다. 음원과 영상을 부당하게 사용한 콘텐츠를 적발하는 기능과는 달리 ‘노아AI 표절 사건’과 같이 변조하는 방식의 표절에 대해서는 YouTube정도 되는 빅테크 기업도 여론이 들끓기 전까지는 손을 놓고 있었기 때문이다. 야비한 표절을 적발하는 기술은 물론 플랫폼 기업에게 협조를 끌어내는 정부 정책도 필요하다.

### 3. Web3.0 기반의 표절 적발

블록체인이나 웹3.0을 통해 저작권을 지키는 방법에 대한 연구는 많다. 블록체인 기술을 활용하여 저작물의 소유권을 증명하고, 저작권 침해 사례를 추적하고, 저작권 분쟁을 해결하는 등의 방식으로 저작권을 보호할 수 있다. 사용자의 멀티미디어 콘텐츠 소유권 및 이용 권한을 증명하여



(그림 3) ©리뷰잉이 유튜브 영상 갈무리[7]

저작권 침해를 방지하는 식으로 발전을 피하기도 했다[9]. 스마트 계약을 활용하여 데이터 이용에 대한 제어 권한을 부여하는 등의 방식으로 개인정보를 보호할 수도 있다.

위와 같이 원본을 무단으로 사용하는 침해에 대해서는 블록체인 기술이 도움이 되겠으나, 변조를 가하는 보다 악의적인 표절을 막기에는 부족하다. 콘텐츠 제작 및 발표 시점에 대한 투명성을 제공하는 것까지가 블록체인의 역할이라 할 수 있겠다. 문서의 소유권과 이용 권한을 확실하게 증명할 수 있어, 저작권 침해와 관련된 분쟁을 예방할 수 있다.

다행히 딥러닝 기술이 발전하면서 여러모로 각광받는 LLM(Large Language Model) 혹은 foundation model인 BERT 모델을 활용한 연구[10]도 나왔다. 다만 LLM 모델은 문장 벡터 간의 유사성을 측정하는 방법과 잠재적 정렬을 수행하는 방법 등의 세부적인 구현 방법에 따라 성능이 달라질 수 있다. BERT 수준의 LLM으로는 표절 여부를 판단하는데 있어서 다른 방법과 함께 활용해야 효과적일 것이다.

전세계를 떠들썩하게 하고 있는 ChatGPT는 GPT-3.5 (Generative Pre-trained Transformer) 모델을 기반으로 하며 이전의 LLM보다 훨씬 더 많은 데이터와 1,750억 개의 파라미터로 사전 학습을 수행했다. 대규모 텍스트 데이터를 활용하여 사전 학습된 언어 모델을 구축하여, 다양한 자연어 처리 작업에 활용하고 있다. 앞서 말했다시피 표절 도구로서의 가능성을 먼저 보였으며 표절 적발 같은 주제보다는 다른 방향으로 활용하고 있는 추세이다.

그럼에도 OpenAI의 GPT-4, Google의 LaMDA 등 Hyper scale LLM이 보여주는 무서운 발전속도를 감안하면 조만간 더욱 정교한 표절 적발 방법을 개발할 것으로 기대가 된다. 인터넷으로 접근 가능한 모든 문서에 대한 표절 적발은 힘들지라도, 명시적으로 저작권 보호를 요청한 저작

물 DB를 기준으로 한 표절 적발은 비교적 용이할 거라 추측한다. 더불어 텍스트와 이미지와 음성을 모두 이해하는 멀티모달(Multi Modal) AI가 속속 등장하는 만큼 멀티미디어의 표절을 적발하는 기능의 개발 역시 얼마 남지 않았기를 바란다. 참고로, 생성 AI의 산출물을 식별하는 일 역시 중요하지만, 표절 적발과는 다른 맥락이므로 이 글에서는 언급하지 않겠다.

## 4. 결 론

GPT-4, LaMDA 같은 초거대 AI 모델은 학습한 데이터와 유사한 데이터에 대해서는 높은 성능을 보인다. 한국어 데이터에 대해서는 다소 성능이 떨어질 가능성이 있다. 새로운 유형의 표절 방법이 등장하거나, 다양한 언어나 문체의 텍스트를 다룰 때 역시 성능저하가 발생할 수 있다. 글로벌 빅테크 기업은 아무래도 한국어에 대한 지원이 상대적으로 미흡할 수밖에 없기도 하다.

기술 외적인 문제 역시 작지 않다. 저작권 침해의 판단 근거가 되는 블록체인 기반 저작권 관리 역시 사기업이 주도하기는 힘들며 기존의 저작권 관련 기관에서 다루던 수준을 상회하는 업무이다. 영상, 음성을 포함하여 사용자들이 올린 신규 콘텐츠가 저작권 정보 블록체인에 등록된 기존 콘텐츠를 표절했는지 여부를 검증하기 위해 AI 기능을 구동하는 IT 인프라 비용은 콘텐츠 플랫폼 기업에게 부담이 될 수밖에 없다.

때문에 콘텐츠 산업의 지속적인 발전을 위해 정부가 나서야만 한다. 한국형 초거대 AI를 만들어야 할지는 모르겠으나, 글로벌 빅테크 기업의 LLM 차기 버전에 한국어 데이터를 선별하여 제공하며 학습에 반영하게 하는 등의 협조를 요청하는 것도 방법 중 하나이다. 국내외 콘텐츠 플랫폼 기업에게 저작권 침해여부 모니터링을 강권할 필요도 있다. 공공기관을 통하거나, 콘텐츠 플랫폼



기업 간에 저작권 관리 블록체인을 운영하게 하는 등의 방안도 강구하길 바란다. 강요만 할 게 아니라 상응하는 지원책을 마련해야 바람직하겠다. 결코 쉬운 일은 아니다.

앞으로 인간이 얼마나 더 창작을 할 수 있을지 장담할 수 있는 사람은 이제 없을 듯하다. 기술 발전에 따라 창작을 한다는 개념이 앞으로도 계속 바뀌어 나가긴 하겠다. 그렇다 하더라도 야비한 콘텐츠 도둑때문에 인간의 창작의욕이 꺾여서는 안 된다. 글자 그대로 하루가 멀다 하고 변화하는 AI 기술 트렌드를 감안하여 콘텐츠 진흥책을 보다 현실적으로 마련해야 한다. 늦어져서는 안 될 것이다.

### 참 고 문 헌

- [1] 박정렬, 최새술. "Web3.0 Reboot: Issues and Prospects". *Electronics and Telecommunications Trends*, 37(2), 73-82. 2022.
- [2] 함대훈. "인터넷, 웹3.0으로의 진화". 삼성 SDS 인사이트 리포트, 2022년 5월.
- [3] OWL Working Group. "Web Ontology Language (OWL)". <https://www.w3.org/OWL/>. 2012.
- [4] Lubani, M., Noah, S. A. M., & Mahmud, R. "Ontology population: Approaches and design aspects". *Journal of Information Science*, 45(4), 502-515. 2019.
- [5] Gurbuz, O., Rabhi, F., & Demirors, O. "Process ontology development using natural language processing: a multiple case study". *Business Process Management Journal*, 25(6), 1208-1227. 2019.
- [6] 신유진, 이지향. "딥러닝 기반 기계번역 개념을 활용한 Text-to-Ontology 변환 사례". 한국정보처리학회 학술대회논문집, 28(2), 891-894. 2021.
- [7] 김동인. "노아시 표절 사건'이 유튜브 생태계

에 던지는 질문". 시사IN, 807호. 2023.

- [8] Gao, C. A., Howard, F. M., Markov, N. S., Dyer, E. C., Ramesh, S., Luo, Y., & Pearson, A. T. "Comparing scientific abstracts generated by ChatGPT to original abstracts using an artificial intelligence output detector, plagiarism detector, and blinded human reviewers". *bioRxiv*, 2022-12. 2022.
- [9] A. Vishwa, F. K. Hussain, "A Blockchain based approach for multimedia privacy protection and provenance". 2018 IEEE Symposium Series on Computational Intelligence (SSCI), Bangalore, India, 1941-1945. 2018.
- [10] Bohra, A., & Barwar, N. C. "A Deep Learning Approach for Plagiarism Detection System Using BERT". In *Congress on Intelligent Systems: Proceedings of CIS 2021, Volume 2* (pp. 163-174). Singapore: Springer Nature Singapore, 2022.

### 저 자 약 력



전 시 형

이메일: sihyoung.jum@gmail.com

- 2003년 동국대학교 컴퓨터멀티미디어공학 (학사)
- 2022년 서울과학기술대학교 산업정보시스템공학 (석사)
- 2003년~2022년 롯데정보통신 / AI기술팀장
- 2022년~현재 메조미디어 / 데이터테크팀장
- 관심분야: 인공지능, 빅데이터, 예측, 최적화, 추천, 광고, 마케팅, 스마트 팩토리