

OpenChain-KWG

# AI 데이터 컴플라이언스 프로세스

수석변호사 조정원

June, 2025

## Contents

---

- 1 AI 개발 및 서비스 과정에서의 리스크
- 2 LG AI연구원 Compliance 체계
- 3 Data Compliance Agent : NEXUS

1

AI 개발 및 서비스  
과정에서의 리스크

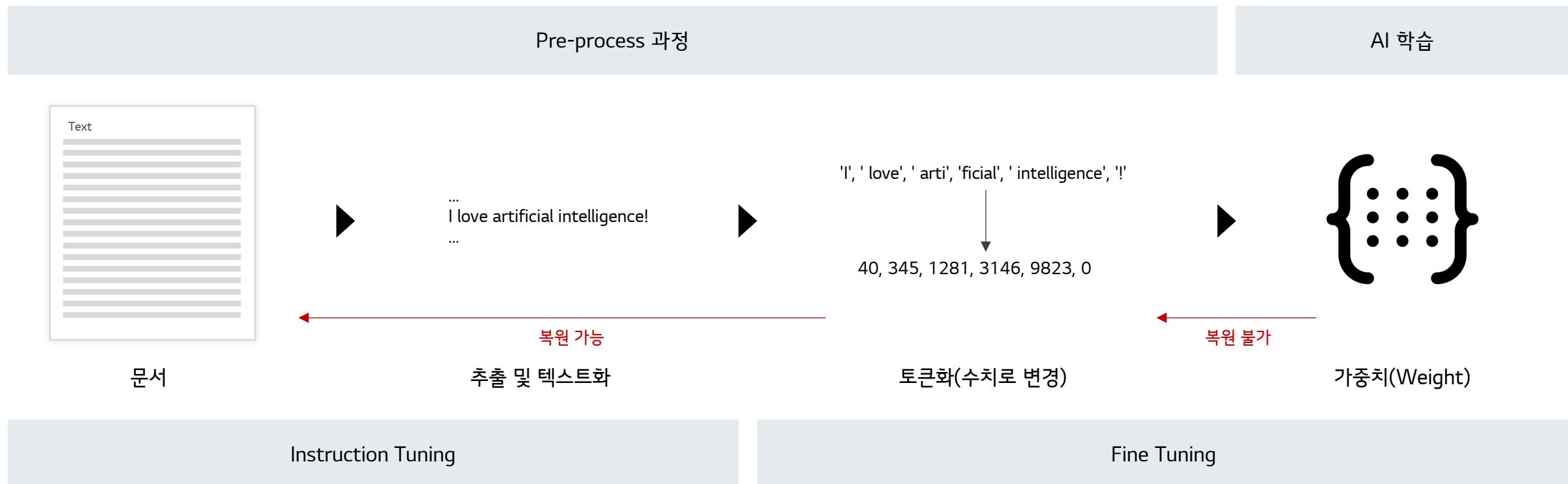
Chapter

---

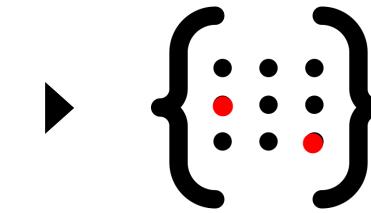
피고	원고	소제기일	배경	이슈가 되는 데이터
OpenAI	Doe 1	2022.3.11	개발자들의 라이선스를 무시하고 Github을 통해 코드를 무단으로 학습하여 CODEX와 Copilot을 개발했다고 주장하며 소제기	GitHub의 Licensed Code
	Mark Walters	2023.6.5	사실이 아닌 내용을 생성하여 명예훼손을 했다고 주장하며 소제기	Privacy (Defamation)
	Paul Tremblay	2023.6.28	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었으며, GPT-4가 그들의 책에 대한 상세 요약을 제공했다고 주장하며 소제기	The Pile3 : Books2, Shadow Liabrary (Library Genesis, Z-Library, Sci-Hub, Bibliotik)
	Sarah Silverman	2023.7.7	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었으며, GPT-4가 그들의 책에 대한 상세 요약을 제공했다고 주장하며 소제기	Shadow Library (Z-Library, Bibliotik, Library Genesis)
	A.T.	2023.9.5	개인정보가 포함된 학습데이터를 사용했다고 주장하며 소제기	Privacy
	Chabon	2023.9.8	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었으며, GPT-4가 그들의 책에 대한 상세 요약을 제공했다고 주장하며 소제기	Shadow Library (Z-Library, Bibliotik, Library Genesis)
	Authors Guild	2023.9.19	작가들이 자신들의 책과 글들이 허락 없이 학습데이터로 활용되고, ChatGPT의 파생물로 현현했다고 주장하며 소제기	Common Crawl, Books1/Books2 Dataset
	Alter	2023.11.21	작가들이 자신들의 책과 글들이 허락 없이 학습데이터로 활용되고, ChatGPT의 파생물로 현현했다고 주장하며 소제기	Authors Guild v OpenAI와 병합
	The New York Times Co.	2023.12.27	NYT가 자신들의 저작물이 GPT에 학습되고, ChatGPT가 NYT의 저작물을 제공한다고 주장하며 소제기	News Data
	Basbanes	2024.1.5	두 작가가 그들의 책이 허락 없이 학습데이터로 활용 되었다고 주장하며 소제기	Authors Guild v OpenAI와 병합
	A.S.	2024.2.27	개인정보가 포함된 학습데이터를 사용했다고 주장하며 소제기	Privacy
	Raw Story Media, Inc.	2024.2.28	뉴스 사이트인 Raw Story와 Alter Net가 Copyright Management Information 삭제에 대한 DMCA 위반으로 소제기	News Data
	Intercept Media Inc.	2024.2.28	뉴스 사이트인 The Intercept가 Copyright Management Information 삭제에 대한 DMCA 위반으로 소제기	News Data
	Elon Musk	2024.2.29	OpenAI 설립 시의 합의 내용과 다르게 운영되는 내용에 대한 합의 위반을 주장하며 소제기	Founding Agreement (합의)
	New York Daily News	2024.4.30	8개의 신문사가 자신들의 저작물이 GPT에 학습되고, ChatGPT가 신문사들의 저작물을 제공한다고 주장하며 소제기	News Data
Meta	Richard Kadrey	2023.7.7	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었다고 주장하며 소제기	The Pile : Books3, Shadow Library (Bibliotik )
	Chabon	2023.9.12	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었다고 주장하며 소제기	Richard Kadrey v. Meta Platforms, Inc.와 병합
Alphabet	J.L.	2023.7.11	인터넷 이용자들을 대표하여 구글이 제3자의 재산이며 개인정보가 포함된 인터넷의 글을 허락 없이 스크래핑 했다고 주장하며 소제기	C-4 (Common Crwal의 파생)
Stability AI	Sarah Andersen	2023.1.13	Visual Art작가들이 그들의 허락없이 저작물이 복제되고 Output으로 현현되었다고 주장하며 소제기	Laion Dataset
	Getty Images	2023.2.3	저작권자의 허락을 받지 않고 사진을 무단으로 이용하여 학습데이터로 사용했다고 소제기	Laion Dataset
Anthropic	Concord Music Group, Inc.	2023.10.18	Universal Music Group의 음반 저작권자들이 Claude가 음반(가사)에 대한 저작권을 침해하여 AI Model을 개발하고 운영했다고 주장하며 소제기	음악의 가사
Databricks, Mosaic ML	Stewart O'Nan	2024.3.8	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었다고 주장하며 소제기	The Pile : Books3
NVIDIA	Abdi Nazemian	2024.3.8	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었다고 주장하며 소제기	The Pile : Books3
Microsoft, The Eleuther AI Institute, Meta, Bloomberg	Mike Huckabee	2023.10.17	작가들이 자신의 책이 허락 없이 학습데이터로 활용되었다고 주장하며 소제기	The Pile : Books3
Ross Intelligence	Thomson Reuteres Enterprise Centre GMBH	2020.5.6	Westlaw의 Research Database 권리자가 AI 스타트업이 Westlaw의 headnote를 불법적으로 복제하여 AI학습데이터로 사용했다고 주장하며 소제기	Article(Legal Research)
Bloomberg	Huckabee	2023.12.27	피고가 작가들의 책이 저작권법으로 보호받는 저작물인 것을 인지하고 있음에도 불구하고 AI학습데이터로 이용했다고 주장하며 소제기	The Pile : Books3
Suno, Inc.	Universal Music Group	2024.6.24	Suno가 저작권이 있는 음악을 무단으로 AI 학습에 사용하고 AI 음악 생성 서비스를 제공했다고 주장하며 소제기	음악

## AI 학습 과정에서의 Data의 변화 : Text에서 Weight으로 변화

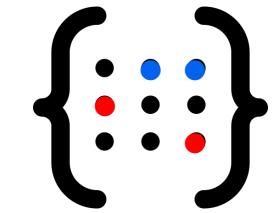
- AI학습 과정에서 저작물은 그 형태가 문서→Text →단위(토큰)→숫자→행렬로 변화하고, AI모델의 출력은 반대 과정을 거칩니다.



"instruction": "Summarize the following article in one paragraph.",  
"input": "The global economy has faced numerous challenges in recent years...",  
"output": "The article discusses the economic disruptions caused by global events..."



"input": "계약서에 명시된 해지 조항에 따라, 계약을 종료할 수 있는 조건은 무엇인가요?",  
"output": "계약 해지 조건은 다음과 같습니다:  
(1) 상대방의 중대한 위반, (2) 파산 또는 청산 절차 시작, ..."

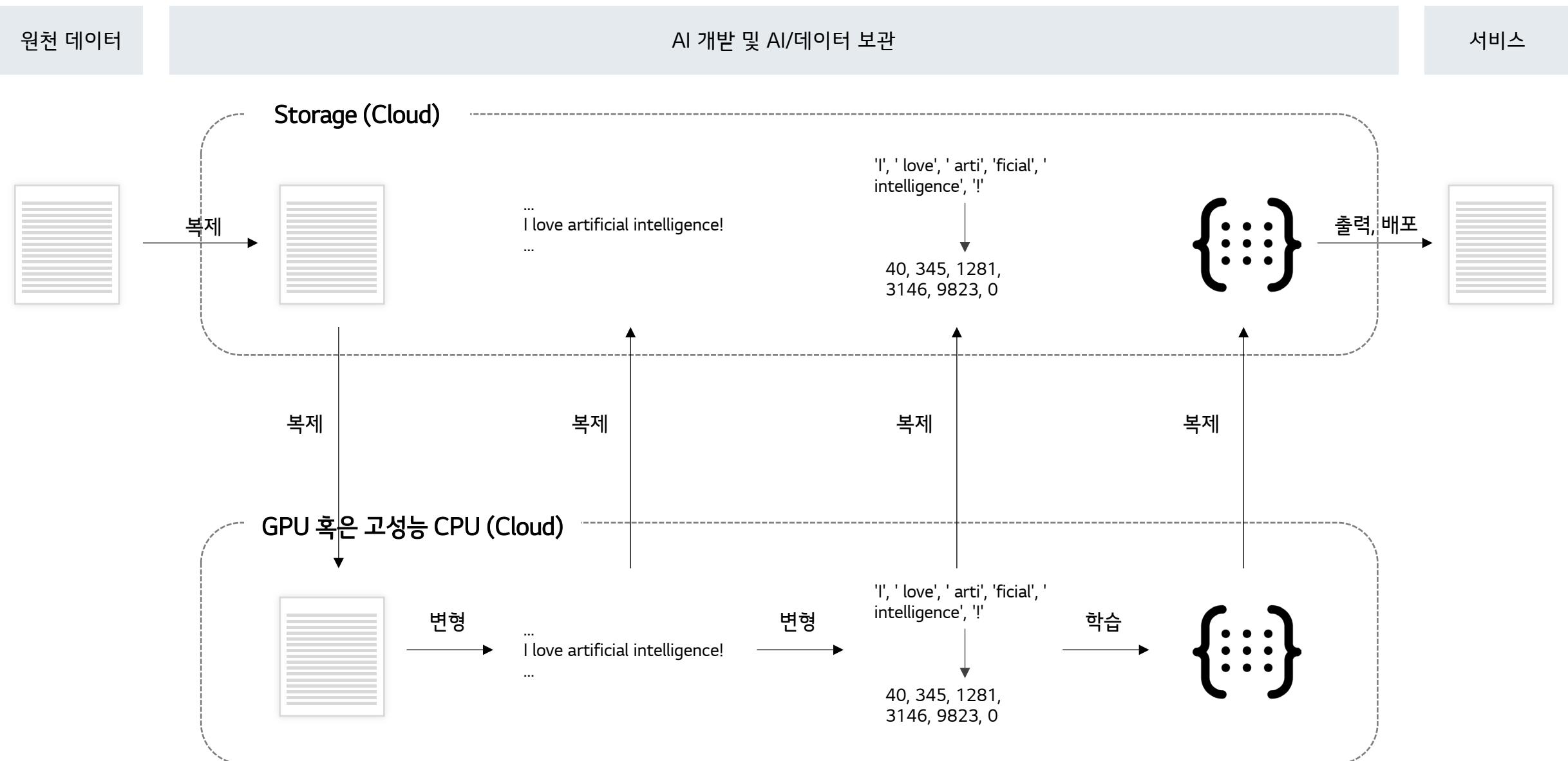


Input(지시문)에 대한 명령 이해(understanding) 및 따르기(following) 학습

특정 Task/지식/도메인에 맞게 정밀 조정(예 : 법률AI)

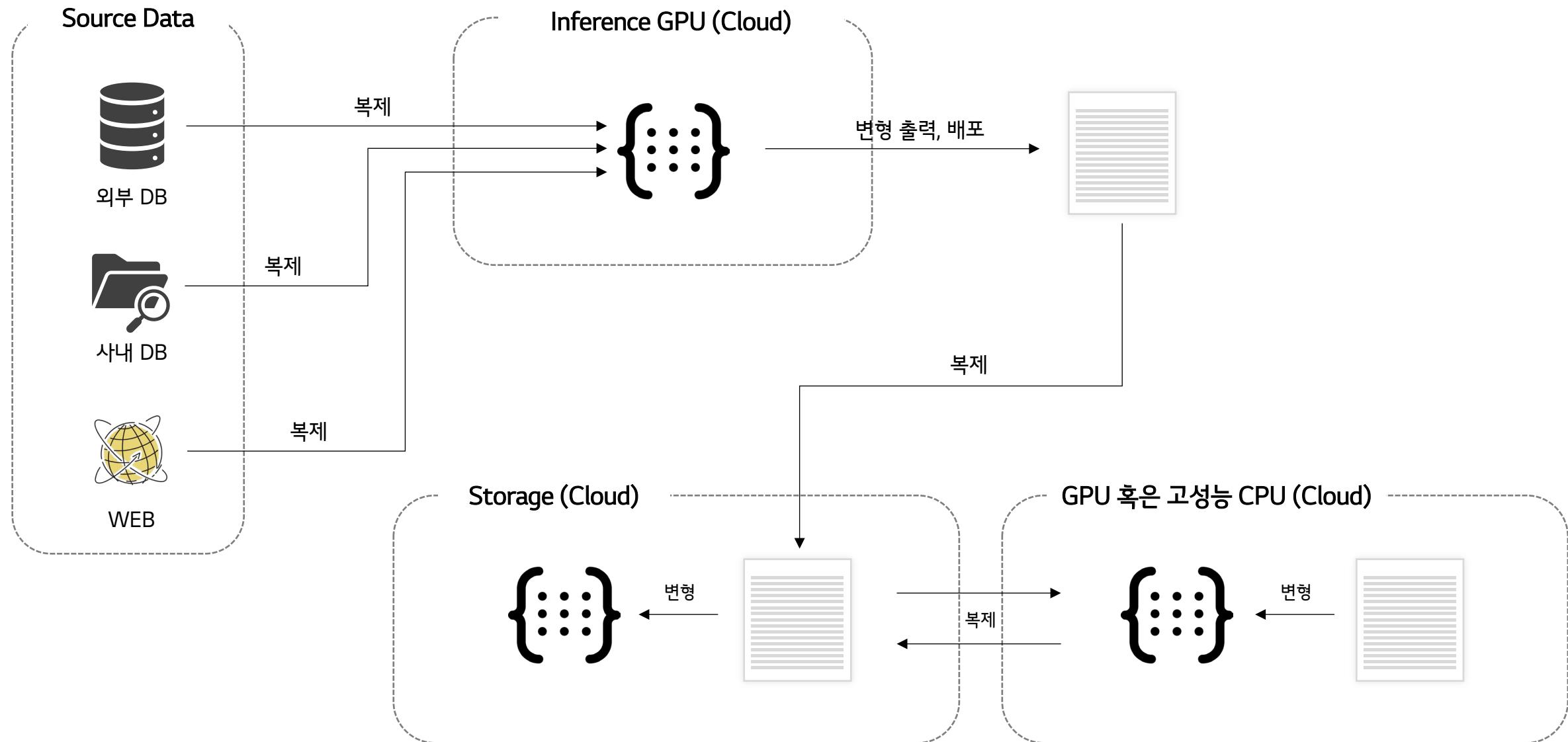
## 학습과정의 리스크

- AI학습 과정에서 원천 데이터는 복제, 변형(혹은 2차적저작물 생성), 배포 됩니다.



## 서비스과정의 리스크

- AI서비스 과정에서도 소스 데이터가 복제, 변형(혹은 2차적저작물 생성), 배포 됩니다. 또한, 배포된 데이터가 다시 재학습 되어 또다른 AI모델에 학습됩니다.



## AI모델 학습 이후의 리스크

"AI 모델 개발이 끝났더라도, 전처리와 토큰화에 수개월·수억원이 투입되기 때문에, 대부분의 기업은 다음 버전을 위해 학습 데이터를 내부 스토리지에 저장합니다."

### Tremblay v. OpenAI 사건 – 학습 데이터 전체 공개 명령

#### 1) 소송 개요

- 원고: 작가 Paul Tremblay, Sarah Silverman 등, 피고 : OpenAI 외
- 주장: OpenAI가 저작권 보호 도서를 무단 수집해 GPT-4 훈련에 사용 → 직접 저작권 침해 및 캘리포니아 부정경쟁법 위반

2) 법원의 명령 : 2025년 1월, 연방법원은 OpenAI에 대해 GPT-4 훈련에 사용된 전체 English Colang 데이터셋을 원고 측에 제공하라고 명령

#### 3) 공개되는 데이터셋의 범위

- English Colang 전체 원본 데이터셋
- 보안실 내, 인터넷 차단된 컴퓨터에서만 열람 가능, 녹음/복사 불가, OpenAI가 메모 검열 가능

4) 시사점 : 법원이 AI 학습 데이터 자체를 저작권 침해 판단의 핵심 증거로 인정하였으며, 추후 유사 제출 명령이 반복될 경우, 기업들은 데이터 출처나 처리 절차 관리 필요성이 있음

### The New York Times v. OpenAI : 소스코드 및 학습 내역 공개 명령

#### 1) 소송 개요

- 원고 : 뉴욕타임즈 (NYT), 피고 : OpenAI 외
- 주장 : OpenAI와 Microsoft가 NYT 뉴스 기사를 무단 수집해 GPT를 훈련시키고, GPT가 NYT 문구를 거의 그대로 복원함 → 직접 저작권 침해 및 계약 위반

2) 법원의 명령 : 2024년 말, 법원은 GPT 훈련 내역 및 ChatGPT의 소스코드 일부에 대한 열람을 허용

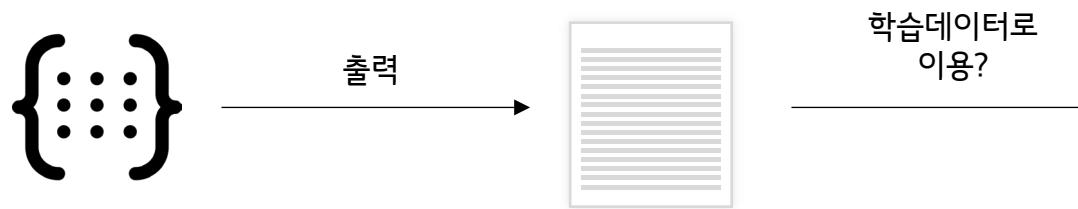
#### 3) 공개되는 소스코드의 범위

- GPT 학습 내역 일부
- ChatGPT 소스코드 열람은 샌드박스 환경 내에서만 허용, 인터넷 완전 차단, 녹화·복제 금지

4) 시사점 : AI 소스코드 조차 법원의 사법 검토 대상이 될 수 있고, 모델 생성과정 전체가 증거 개시 대상이 될 수 있음

## — 지식증류(distillation)를 제한하는 이용약관

- AI Model을 통해 Output을 생성하여 학습데이터로 재활용 하는 practice는 일반적인 수준으로 이용되고 있습니다. 이는 AI Model 개발자의 학습데이터셋 구축 노하우를 통한 양질의 데이터를 생성할 수 있는 방법입니다.
- 지식증류를 제한하는 이용약관을 확인 후 AI모델의 Output을 이용해야 합니다.



### OpenAI Terms of Use

What you cannot do. You may not use our Services for any illegal, harmful, or abusive activity. For example, you may not:

...  
Use Output to develop models that compete with OpenAI.

### Llama 2 Community License Agreement

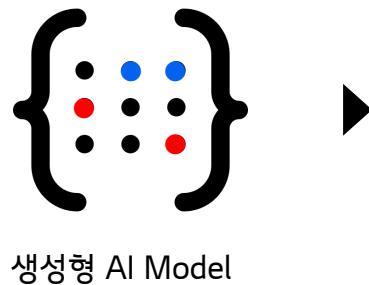
v. You will not use the Llama Materials or any output or results of the Llama Materials to improve any other large language model (excluding Llama 2 or derivative works thereof).

### Llama 4 License Agreement

i. If you distribute or make available the Llama Materials (or any derivative works thereof), or a product or service (including another AI model) that contains any of them, you shall (A) provide a copy of this Agreement with any such Llama Materials; and (B) prominently display "Built with Llama" on a related website, user interface, blogpost, about page, or product documentation. If you use the Llama Materials or any outputs or results of the Llama Materials to create, train, fine tune, or otherwise improve an AI model, which is distributed or made available, you shall also include "Llama" at the beginning of any such AI model name.

## 원저작물의 복제 혹은 2차적저작물

- Transformer 아키텍처 기반의 Encoder-Decoder 모델은 원저작물과 상당히 유사한 출력물(Output)을 생성해내게 되었습니다.



### New York Times v. OpenAI, Microsoft

#### Output from GPT-4:

exempted it from regulations, subsidized its operations and promoted its practices, records and interviews showed.

Their actions turned one of the best-known symbols of New York — its yellow cabs — into a financial trap for thousands of immigrant drivers. More than 950 have filed for bankruptcy, according to a Times analysis of court records, and many more struggle to stay afloat.

"Nobody wanted to upset the industry," said David Klahr, who from 2007 to 2016 held several management posts at the Taxi and Limousine Commission, the city agency that oversees medallions. "Nobody wanted to kill the golden goose."

New York City in particular failed the taxi industry, The Times found. Two former mayors, Rudolph W. Giuliani and Michael R. Bloomberg, placed political allies inside the Taxi and Limousine Commission and directed it to sell medallions to help them balance budgets and fund key initiatives.

During that period, much like in the mortgage lending crisis, a group of industry leaders enriched themselves by artificially inflating medallion prices. They encouraged medallion buyers to borrow as much as possible and ensnared them in interest-only loans and other one-sided deals that often required borrowers to pay hefty fees, forfeit their legal rights and give up most of their monthly incomes.

When the market collapsed, the government largely abandoned the drivers who bore the brunt of the crisis. Officials did not bail out borrowers or persuade banks to soften loan

the company's dormitories, according to the executive. Each employee was given a biscuit and a cup of tea, guided to a workstation and within half an hour started a 12-hour shift fitting glass screens into beveled frames. Within 96 hours, the plant was producing over 10,000 iPhones a day.

"The speed and flexibility is breathtaking," the executive said. "There's no American plant that can match that."

Similar stories could be told about almost any electronics company — and outsourcing has also become common in hundreds of industries, including accounting, legal services, banking, auto manufacturing and pharmaceuticals.

But while Apple is far from alone, it offers a window into why the success of some prominent companies has not translated into large numbers of domestic jobs. What's more, the company's

#### Actual text from NYTimes:

exempted it from regulations, subsidized its operations and promoted its practices, records and interviews showed.

Their actions turned one of the best-known symbols of New York — its signature yellow cabs — into a financial trap for thousands of immigrant drivers. More than 950 have filed for bankruptcy, according to a Times analysis of court records, and many more struggle to stay afloat.

"Nobody wanted to upset the industry," said David Klahr, who from 2007 to 2016 held several management posts at the Taxi and Limousine Commission, the city agency that oversees cabs. "Nobody wanted to kill the golden goose."

New York City in particular failed the taxi industry, The Times found. Two former mayors, Rudolph W. Giuliani and Michael R. Bloomberg, placed political allies inside the Taxi and Limousine Commission and directed it to sell medallions to help them balance budgets and fund priorities. Mayor Bill de Blasio continued the policies.

Under Mr. Bloomberg and Mr. de Blasio, the city made more than \$855 million by selling taxi medallions and collecting taxes on private sales, according to the city.

But during that period, much like in the mortgage lending crisis, a group of industry leaders enriched themselves by artificially inflating medallion prices. They encouraged medallion buyers to borrow as much as possible and ensnared them in interest-only loans and other one-sided deals that often required them to pay hefty fees, forfeit their legal rights and give up most of their monthly incomes.

A foreman immediately roused 8,000 workers inside the company's dormitories, according to the executive. Each employee was given a biscuit and a cup of tea, guided to a workstation and within half an hour started a 12-hour shift fitting glass screens into beveled frames. Within 96 hours, the plant was producing over 10,000 iPhones a day.

"The speed and flexibility is breathtaking," the executive said. "There's no American plant that can match that."

Similar stories could be told about almost any electronics company — and outsourcing has also become common in hundreds of industries, including accounting, legal services, banking, auto manufacturing and pharmaceuticals.

But while Apple is far from alone, it offers a window into why the success of some prominent companies has not translated into large numbers of domestic jobs. What's more, the company's

### Getty Images v. Stability AI

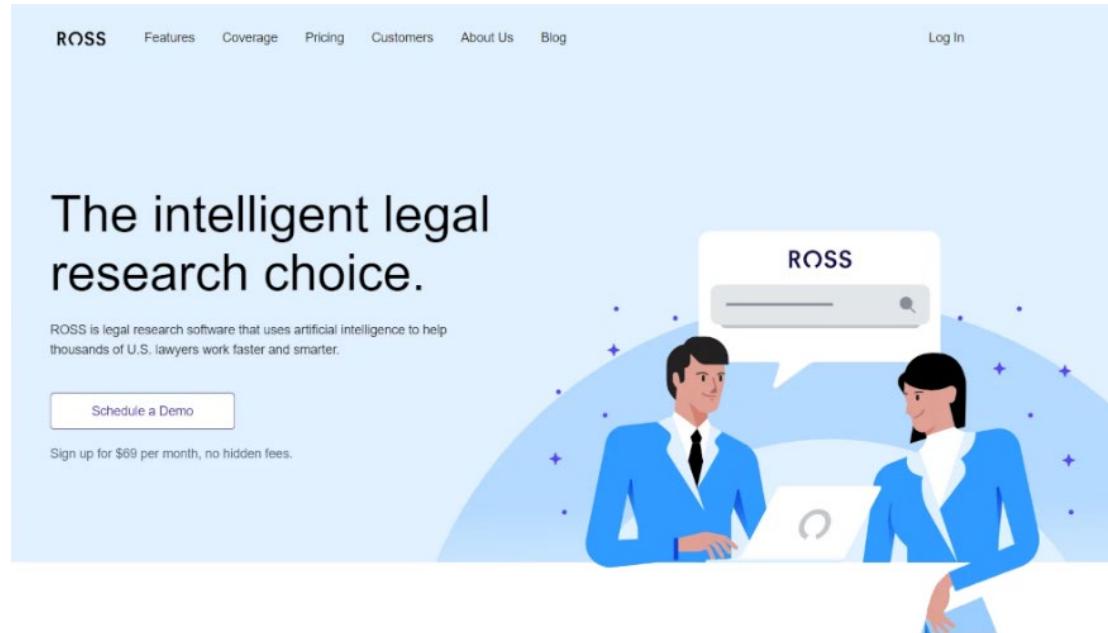


### J. Doe 1, J. Doe 2 v Github, Microsoft, Copilot

```
function isPrime(n) {  
    if (n < 2) {  
        return false;  
    }  
    for (let i = 2; i < n; i++) {  
        if (n % i === 0) {  
            return false;  
        }  
    }  
    return true;  
}
```

## Thomson Reuters v. ROSS Intelligence (미국에서의 AI 저작권 소송 첫 판결)

- 2020년 원고인 Thomson Reuters는 ROSS사가 자사의 서비스인 Westlaw의 자료를 무단 복제하여 학습했다고 주장하며 소를 제기하였고, 법원에서 2025년 2/11에 Summary Judgement 판결을 내리면서 원고인 Thomson Reuters의 승소 판결을 내림.



### 법원의 판결

- 공정이용의 네가지 요소인 1) 사용 목적, 2) 원본 저작물의 성격, 3) 사용된 저작물의 양, 4) 해당 사용이 원본 저작물의 시장 가치에 미치는 영향 중 특히 4)에 대해서 시장에서 Thomson Reuters에 손해를 미쳤다고 판단하며 ROSS의 공정이용 주장을 기각함.
- ROSS가 Westlaw의 대체제를 개발함으로써 경쟁하려 했으며, 이 행위가 AI 학습 데이터 시장에 영향을 미칠 수 있다고 판단함. ROSS가 독자적으로 개발할 수도 있었던 것을 Thomson Reuters의 저작권을 침해하며 사용한 것은 정당화될 수 없다고 판결했습니다.
- 판결이 나온 지 몇 시간 만에, 음악 저작권자(출판사)가 Anthropic을 상대로 제기한 소송의 원고 측은 법원에 이 판결을 제출할 수 있도록 허가를 요청함.

## — Li v. Liu (베이징인터넷법원판결 : (2023)경0497민초11279호)

- 원고인 Li가 Stable Diffusion을 이용해 생성한 이미지를 SNS에 올린 후, 피고 Liu가 이를 무단으로 사용한 것에 대해 '저작권 침해 및 정보네트워크 보급권 분쟁 혐의'로 소송을 제기함.

### Stable Diffusion으로 생성한 원고의 이미지



### 법원의 판결

- AI는 도구일 뿐이며, 사람의 지적 투자가 창작과정에 반영되었으므로 생성된 이미지의 저작물성이 인정됨.
- 원고가 프롬프트 입력과 파라미터 설정 등 창작 과정에 관여했으므로, 원고가 생성된 이미지의 저작자로 인정됨.
- 피고가 원고의 저작물을 무단 사용하고 워터마크를 제거했으므로 저작권 침해가 인정됨.
- 피고는 원고에게 공개 사과문을 게시해야 하며, 500위안(한화 약 10만원)을 지급해야 함)

## — 신창화 문화발전 유한공사 v. AI회사 (광저우인터넷법원판결 : (2024)월0192민초1호)

- 중국 울트라맨 독점 저작권 사용 허가를 받은 원고인 상하이 신창화 문화발전 유한공사가 피고 A사의 생성형 AI 서비스가 제공되는 사이트에서 울트라맨 이미지가 생성되는 것을 보고 저작권이 침해되었다며 소송을 제기함.

### 피고 A사의 서비스로 생성된 이미지



### 법원의 판결

- 피고의 웹사이트에서 '울트라맨' 키워드 입력 만으로 저작권으로 보호되는 울트라맨 이미지와 실질적으로 유사한 이미지가 출력됨.
- 생성형AI 서비스 제공자인 피고는 저작권법 상 원고의 복제권 및 각색권을 침해함.
- 피고는 '생성형 인공지능 잠정관리 방법'(중국의 법률)에 따른 '민원 제기 및 신고 체제 구축'을 하지 않았고, '사용자 규정, 약관 등으로 사용자가 타인의 지식재산권을 침해하지 않도록 고시'를 하지 않았으며, 'AI 산출물에 대한 표시 의무'를 다하지 않았음.
- 피고는 원고에게 10,000위안(약 200만원)을 지급해야 하며, 울트라맨 관련 키워드로 실질적으로 유사한 이미지가 생성되지 않도록 예방 조치를 취해야 함.

## — AI Output의 유사성에 대한 원저작자 – AI개발자의 합의 : Concord Music Group v. Anthropic (저작권 소송 중 일부 합의)

- 2023년 원고인 Universal Music, Concord Music, ABCKO Music 등이 Anthropic의 Claude 챗봇에서 음악의 가사가 그대로 출력되면서 복제, 전시, 2차적 저작물 작성권리를 침해했다고 주장하며 1200억이 넘는 손해배상을 청구함.
- 음악의 가사를 학습하는 부분에 대한 예비 금지명령(preliminary injunction)에 대한 다툼은 계속되고 있고, Claude 챗봇에서 가사가 출력되는 부분에 대해서는 원고와 피고가 합의하여 법원이 합의 문서(stipulation)를 승인함.

### 합의 내용

- Anthropic은 현재의 AI 모델 및 제품에서 가드레일(Guardrails)을 유지하고, 향후 도입되는 새로운 LLM 및 신규 제품에 대해서도 일관된 방식으로 입력 및 출력에 대한 가드레일을 적용할 것.
- 소송 절차 진행 중 원고측은 언제든지 피고측에 가드레일이 효과적으로 작동하지 않고 있다고 지적할 수 있고, 이에 대해 Anthropic은 신속하게 응답 및 조사를 취해야 하며, 이는 아래와 같은 경우를 포함함:
  - Output을 통해 음악 가사가 재생산, 배포, 전시되거나,
  - 이를 통해 2차적 저작물이 출력되는 행위



- 시사점:
  - AI 서비스를 개발하고 제공하고 있는 회사는 금지명령(injunction)에 대항하며 저작물 권리자와 합의 없이 본안 판결까지 갈 수 있는 상황에 놓이게 되는가?
  - Output에 대해 기술적인 조치에 대한 법적 구속력 있는 합의에 이르는 경우, AI 개발 및 서비스 배포에 많은 추가 비용이 발생하지 않는가?
  - AI 서비스를 개발하는 기업은 완벽한 기술적인 조치를 할 수 있는가?

## AI개발사의 항변 - 공정이용 법리 (New York Times v. OpenAI, Microsoft를 바탕으로)

- 피고의 항변 : 저작권법 상 저작물을 원저작자의 허락 없이도 이용할 수 있는 공정이용(Fair Use)의 법리

### 미국 17 U.S.C. 107

#### §107. Limitations on exclusive rights: Fair use

Notwithstanding the provisions of sections 106 and 106A, the fair use of a copyrighted work, including such use by reproduction in copies or phonorecords or by any other means specified by that section, for purposes such as criticism, comment, news reporting, teaching (including multiple copies for classroom use), scholarship, or research, is not an infringement of copyright. In determining whether the use made of a work in any particular case is a fair use the factors to be considered shall include—

- (1) the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes;
- (2) the nature of the copyrighted work;
- (3) the amount and substantiality of the portion used in relation to the copyrighted work as a whole; and
- (4) the effect of the use upon the potential market for or value of the copyrighted work.

The fact that a work is unpublished shall not itself bar a finding of fair use if such finding is made upon consideration of all the above factors.

2012.3.15  
한-미 FTA  
발효시점

### 대한민국 저작권법 제35조의5

#### 제35조의5(저작물의 공정한 이용)

① 제23조부터 제35조의4까지, 제101조의3부터 제101조의5까지의 경우 외에 저작물의 일반적인 이용 방법과 충돌하지 아니하고 저작자의 정당한 이익을 부당하게 해치지 아니하는 경우에는 **보도·비평·교육·연구 등을 위하여** 저작물을 이용할 수 있다.

② 저작물 이용 행위가 제1항에 해당하는지를 판단할 때에는 다음 각 호의 사항 등을 고려하여야 한다.

1. **영리성 또는 비영리성 등** 이용의 목적 및 성격
2. 저작물의 종류 및 용도
3. 이용된 부분이 저작물 전체에서 차지하는 비중과 그 중요성
4. 저작물의 이용이 그 저작물의 현재 시장 또는 가치나 잠재적인 시장 또는 가치에 미치는 영향

### 공정이용 4요소

### 고려할 사항

제1요소 : 원저작물 이용의 목적과 성격	원저작물의 이용이 상업적인가? 원저작물의 이용이 변형적인가? 원저작물의 이용이 악의나 고의에 의한 이용인가? 원저작물의 이용이 합리적이고 관습적인가?
제2요소 : 원저작물의 성격	원저작물이 창작적 저작물 혹은 사실적 저작물인가? 원저작물이 발행된 저작물인가?
제3요소 : 이용된 부분의 양과 중요성	원저작물을 이용한 양적, 질적 수준이 목적과 관련하여 합리적인 수준인가? (제1요소 및 제4요소와 상호 보완 관계)
제4요소 : 원저작물의 잠재적 시장 또는 가치에 미치는 영향	원저작물의 이용이 원저작물의 잠재적 시장 또는 가치에 얼마나 영향을(손해를) 미치는가?

## — AI개발사의 항변 - 공정이용 법리 (New York Times v. OpenAI, Microsoft를 바탕으로)

- 공정 이용의 4요소에 대해서는 일률적인 기준을 제시하고 있지 않으며, 상호 관계에서 각 요소가 얼마나 영향력을 가지고 있는지에 대한 지침이 존재하지 않음.

### Goldsmith v. Andy Warhol Foundation 사건



- 사진작가 Goldsmith가 1981년 미국의 뮤지션인 Prince Rogers Nelson을 촬영한 사진을 바탕으로 Andy Warhol이 변형하여 만든 "Prince Series" 작품에 대해 Goldsmith 측이 저작권 침해의 소를 제기함.
- 1심 : 공정이용이 인정되어 Andy Warhol 측이 승소.
- 2심 : 공정이용이 인정되지 않아서 Goldsmith 측이 승소.
- 3심 : 대법원에서 공정이용이 인정되지 않아서 Goldsmith 측이 최종 승소.
- 1. 이용의 목적과 성격 : 실질적으로 동일한 목적이며, 상업적 이용이었음.
- 2. 원저작물의 성격 : 원저작물인 사진의 창작성과 예술성이 인정됨.
- 3. 이용된 부분의 양과 중요성 : 원작의 핵심적인 본질적 요소를 그대로 유지하였음.
- 4. 원저작물의 잠재적 시장 또는 가치에 미치는 영향 :  
Andy Warhol 작품이 Goldsmith 사진의 잠재적 시장을 침해할 수 있음.

### Authors Guild v. Google(Books) 사건

- Authors Guild 측은 Google Books 측이 수백만 권의 책을 스캔하여 디지털화 하고, 이를 검색 가능하게 만들어 원고의 저작권을 침해했다고 주장하며 소를 제기함.
  - 1심 : 공정이용이 인정되어 Google 측이 승소.
  - 2심 : 공정이용이 인정되어 Google 측이 승소.
  - 3심 : 대법원에서 공정이용이 인정되어 Google 측이 최종 승소.
- 1. 이용의 목적과 성격 : 책의 전문 검색 기능이라는 새로운 가치를 창출하였고, 공익적 목적이 큼.
- 2. 원저작물의 성격 : 대부분의 책이 논픽션이었으나, 이 요소가 중요하게 판단되지는 않음.
- 3. 이용된 부분의 양과 중요성 : 스니펫 표시가 제한적이어서, 저작물 전체를 대체할 수는 없음.
- 4. 원저작물의 잠재적 시장 또는 가치에 미치는 영향 :  
오히려 책의 발견 가능성을 높여 저자와 출판사에 이익을 줄 수 있고, 시장 대체 효과는 미미함.

OpenAI의 NYT 기사 이용은, 1) 목적과 성격이 변형적인가? 상업적인가? 2) NYT 기사는 창작적인가? 사실적인가? 3) 원저작물을 질적, 양적으로 얼마나 사용하는가?

또한, 4) NYT의 구독 모델이나 광고 수익 등 NYT의 잠재적 시장에 어떠한 영향을 미쳤는가?

## — Data의 투명성과 관련된 법적, 윤리적 움직임

- 각 국가는 AI학습에 이용되는 데이터의 출처나 저작권법 준수 여부와 관련한 법률을 제정하거나 입법을 추진하고 있습니다.

### 워싱턴주 입법中

HB 1168 – AI 학습 데이터 투명성 법안

- 생성형 AI 시스템 또는 서비스를 개발하는 기업은 해당 시스템의 훈련에 사용된 데이터에 대한 문서를 공개해야 합니다.
- 공개 문서에는 데이터 출처, 데이터 유형 및 수량, 개인정보 사용 여부, 데이터 수정 과정, 합성 데이터 사용 여부 등이 포함되어야 합니다.

### 미국 저작권청

저작권과 AI 3차 보고서 초안

- 생성형 AI 학습에서의 공정이용 적용 가능성 모델 학습 시 저작물 이용은 원칙적으로 저작권 침해의 소지가 있음.
- 상업적으로 수익을 창출하고 학습에 사용된 저작물을 그대로 모방하거나 시장에서 경쟁한다면, 이는 공정이용의 경계를 넘어서는 행위로 판단됨.
- 시장 영향과 상업성 AI 기업이 무단으로 저작물을 활용하여 수익을 창출하면서도 저작권자에게 대가를 지불하지 않는다면, 이는 시장 왜곡을 초래하고, 공정하지 않다고 평가됨.

### 캘리포니아주 제정

AB 2013 – 생성형 인공지능 훈련 데이터 투명성법

- 개발자는 해당 시스템이나 서비스가 출시되기 전에, 또는 실질적인 수정이 이루어지기 전에, 훈련에 사용된 데이터에 대한 문서를 웹사이트에 공개해야 합니다.
- 데이터셋의 출처, 소유자, 라이선스 획득 여부 등이 공개되어야 합니다.

### 버지니아주 입법中

HB 2250 – AI 학습 데이터 투명성 법안

- 소비자는 자신의 개인정보가 AI 모델 훈련에 사용되는 것을 거부할 수 있는 권리를 가집니다.
- 개발자는 소비자의 요청에 따라 훈련 데이터에 대한 확인 및 삭제 요청을 처리해야 하며, '훈련 금지(Do Not Train)'로 지정된 데이터를 AI 훈련에 사용해서는 안 됩니다.

### EU

EU AI ACT 실천강령 공개 예정

- 범용 AI 모델 제공자는 저작권법 준수를 위한 정책을 수립·유지하고 실행해야 함(정책 공개를 권장)
- 웹 크롤링을 통해 수집하는 저작물은 합법적으로 접근 가능한 데이터만 크롤링
- 'robots.txt' 등 기계 판독 가능한 형식의 옵트아웃을 준수하기 위해 노력
- 웹 크롤링 외의 방식(제3자 제공 데이터 등)으로 학습데이터를 수집하는 경우, 해당 데이터의 저작권 상태를 확인하기 위해 합리적 노력을 해야 함
- 반복적으로 침해 산출물을 생성할 정도로 저작물을 기억하는 위험을 완화하기 위한 합리적 노력을 해야함 등

## 텍스트 및 데이터 마이닝(TDM) 면책?

- EU TDM(DSM) 면책 규정은 직접적인 AI 학습 데이터 입수, 학습, 배포 과정의 일부 영역만을 면책합니다.

### DIRECTIVE (EU) 2019/790 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

#### Article 3

##### Text and data mining for the purposes of scientific research

1. Member States shall provide for an exception to the rights provided for in Article 5(a) and Article 7(1) of Directive 96/9/EC, Article 2 of Directive 2001/29/EC, and Article 15(1) of this Directive for reproductions and extractions made by **research organisations and cultural heritage institutions** in order to carry out, for the purposes **of scientific research**, text and data mining of works or other subject matter to which they have lawful access.

...

① Article 3에서는 과학적 연구인 경우에는 DSM 규정이 제한 없이 적용 됨을 명시함.

→ 독일 저작권법(Urheberrechtsgesetz, UrhG) 제60조d항

(1) In order to enable **scientific research**, it shall be permissible to reproduce works for text and data mining purposes (section 44b (1) and (2), first sentence) if: 1. the access to the works is lawful and 2. the reproductions are made by **research organisations or cultural heritage institutions**....

#### Article 4

##### Exception or limitation for text and data mining

1. Member States shall provide for an exception or limitation to the rights provided for in Article 5(a) and Article 7(1) of Directive 96/9/EC, Article 2 of Directive 2001/29/EC, Article 4(1)(a) and (b) of Directive 2009/24/EC and Article 15(1) of this Directive for reproductions and extractions of lawfully accessible works and other subject matter for the purposes of text and data mining.

...

3. The exception or limitation provided for in paragraph 1 shall apply on condition that the use of works and other subject matter referred to in that paragraph **has not been expressly reserved by their rightholders in an appropriate manner, such as machine-readable means in the case of content made publicly available online**.

...

② Article 4에서는 상업적 목적까지 포함되는 경우(과학적 연구가 아닌 경우), 저작권자가 기계 판독 가능한 수단 등으로 의사표시하여 Opt-Out 할 수 있는 것으로 명시함.

→ 독일 저작권법(Urheberrechtsgesetz, UrhG) 제44조b항

...(3) Rightholders shall be entitled to prohibit the use of works for text and data mining purposes in accordance with subsection (2) if they have expressly reserved this right in an appropriate manner.

2

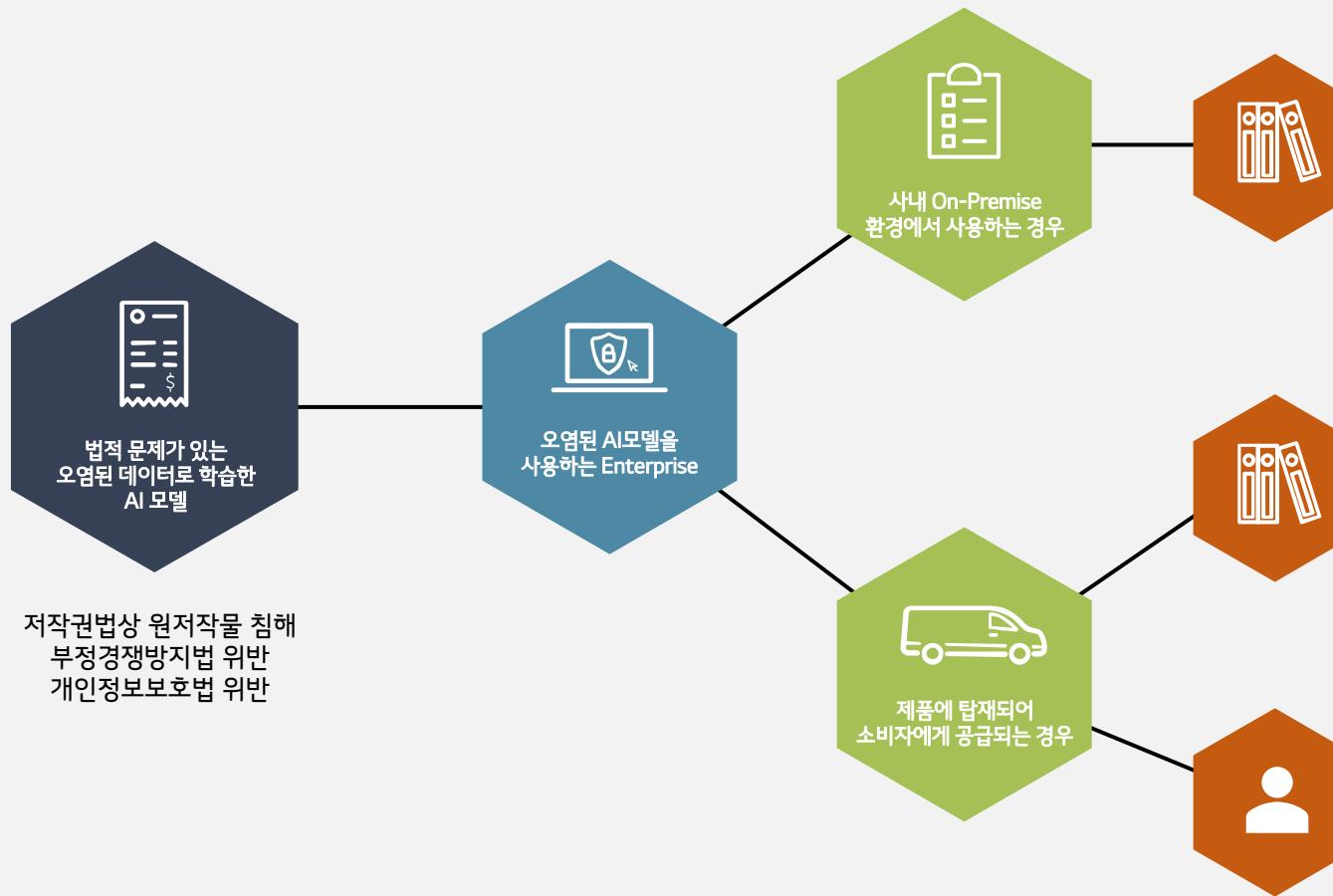
Chapter

---

LG AI연구원  
Compliance 체계

## — LG에게 발생할 수 있는 리스크

- 법적 문제가 있는 데이터를 학습데이터로 사용한 AI모델을 Enterprise에서 사용하게 되면, 데이터 저작자와의 분쟁으로 손해배상, 형사처벌, 과태료, AI모델 폐기 청구 등의 분쟁이 발생할 수 있으며, Enterprise 소비자에 대해 리콜을 시행하거나, 손해배상 및 형법상 고소될 수 있습니다.



### 데이터 저작자와의 분쟁

- 저작권법상 저작권 침해에 대한 민·형사상의 책임<sup>1</sup>
- 부정경쟁방지법상 부정사용행위, 퍼블리시티권 침해, 성과도용행위에 따른 민·형사상의 책임<sup>2</sup>
- 영업비밀 침해에 따른 민·형사상의 책임<sup>3</sup>
- 위법한 개인정보 이용에 따른 법적 책임<sup>4</sup>

### 데이터 저작자와의 분쟁

- 저작권법상 저작권 침해 방조행위에 대한 법적 책임<sup>5</sup>
- 부정경쟁행위 내지 영업비밀 침해행위의 방조행위에 대한 법적 책임
- 영업비밀 침해에 따른 민·형사상의 책임
- 위법한 개인정보 이용에 따른 법적 책임

### 소비자와의 분쟁

- 소비자기본법상 소비자 재산에 위해를 끼친 것, 결함에 대해 고지의 의무를 다하지 않은 것에 대한 불법행위<sup>6</sup>
- 민법상 법률적 제한 내지 장애로 인해 하자에 해당하여 손해배상 책임 발생<sup>7</sup>
- 형법상 해당 하자가 발생할 것을 알고도 판매한 것에 대한 사기죄<sup>8</sup>

1: 저작권법 전반 및 제123조침해행위 정지청구 참고, 2: 부정경쟁방지법 제2조 제1호 (카)목 3, 3: 부정경쟁방지법 제2조 제3호 (나)목 및 (라)목 및 부정경쟁방지법 제2조 제3호 (다)목 및 (바)목 , 4: 개인정보보호법 제39조, 5: 대법원 2021. 11. 25. 선고 2021도10903 판결 참조, 6: 소비자기본법 제19조 제1항(중국 소비자 권리 보호법 제19조), 7: 민법 제580조, 민법 제582조 및 대법원 2000. 1. 18. 선고 98다18506 판결, 대법원 2018. 7. 12. 선고 2015다64315 판결 참조, 8: 대법원 2008. 5. 8. 선고 2008도1652 판결 참조

## Data Compliance

- Data Compliance 가이드라인을 통해, Data lifecycle 관점에서 각 담당자가 Data의 이용 과정에서 발생할 수 있는 risk를 탐지합니다.
- Data Compliance Report(평가 기준)을 통해, 학습 데이터에 대한 법적 리스크를 평가합니다.



### DATA COMPLIANCE 가이드라인

- AI 연구·개발 과정에서 발생할 수 있는 데이터 관련 법률 리스크
- LG AI연구원의 데이터 취급 관련 표준 업무 프로세스
- Data Compliance 체크리스트 (수집, 보관, 파기, 제3자 제공/위탁)

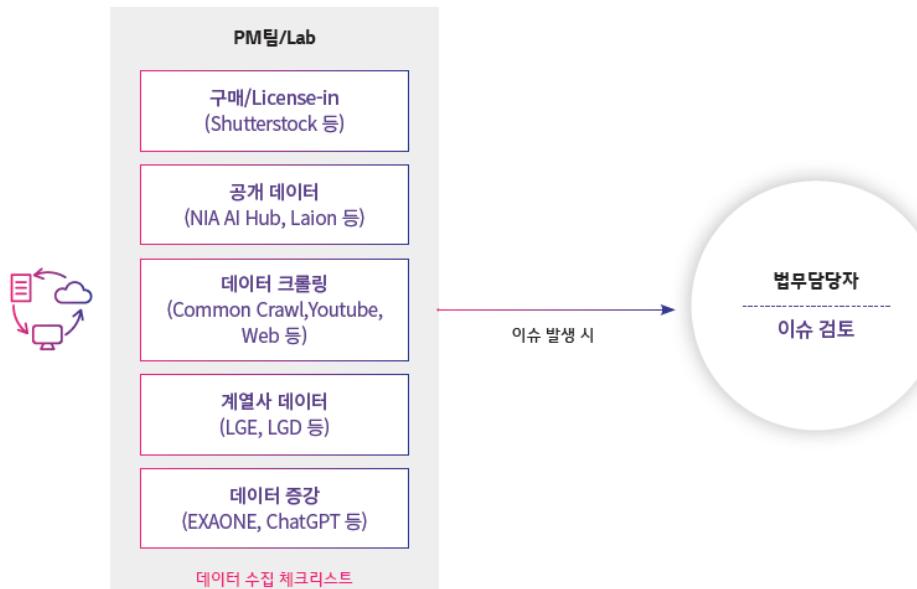
### DATA COMPLIANCE REPORT

- Data Compliance Score 및 Class 평가 기준(ENG)

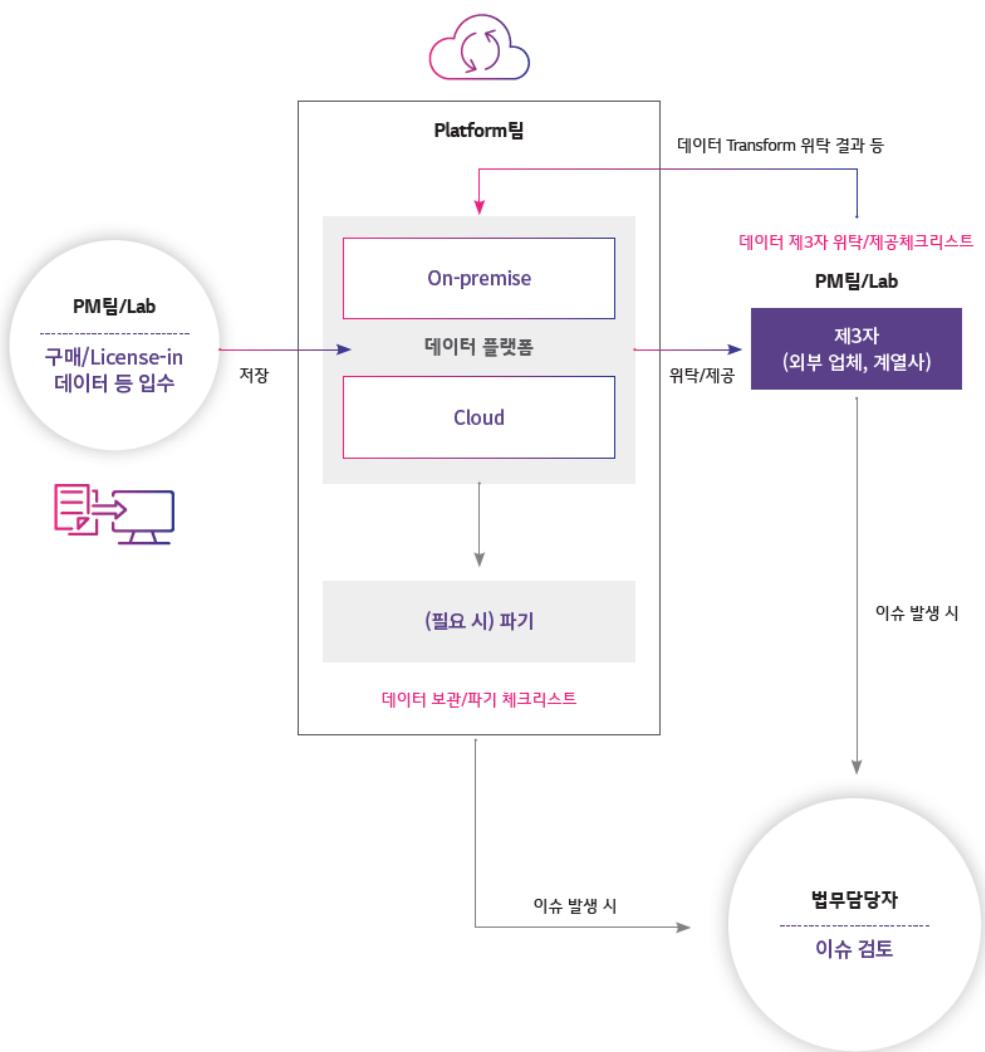
## Data Compliance

- Data의 수집, 보관, 파기, 이전, 구매, 이용허락(License-in) 과정에서 발생할 수 있는 법적 리스크를 탐지합니다.
- 법적 리스크가 존재하는 경우, 법무담당자가 개입하여 올바른 IP Transaction이 발생하도록 조치합니다.

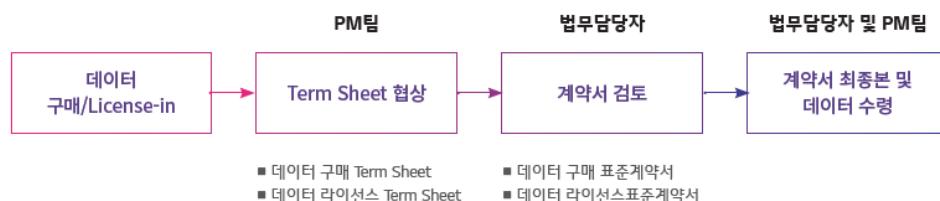
Data 수집 프로세스



Data 보관/파기/이전 프로세스



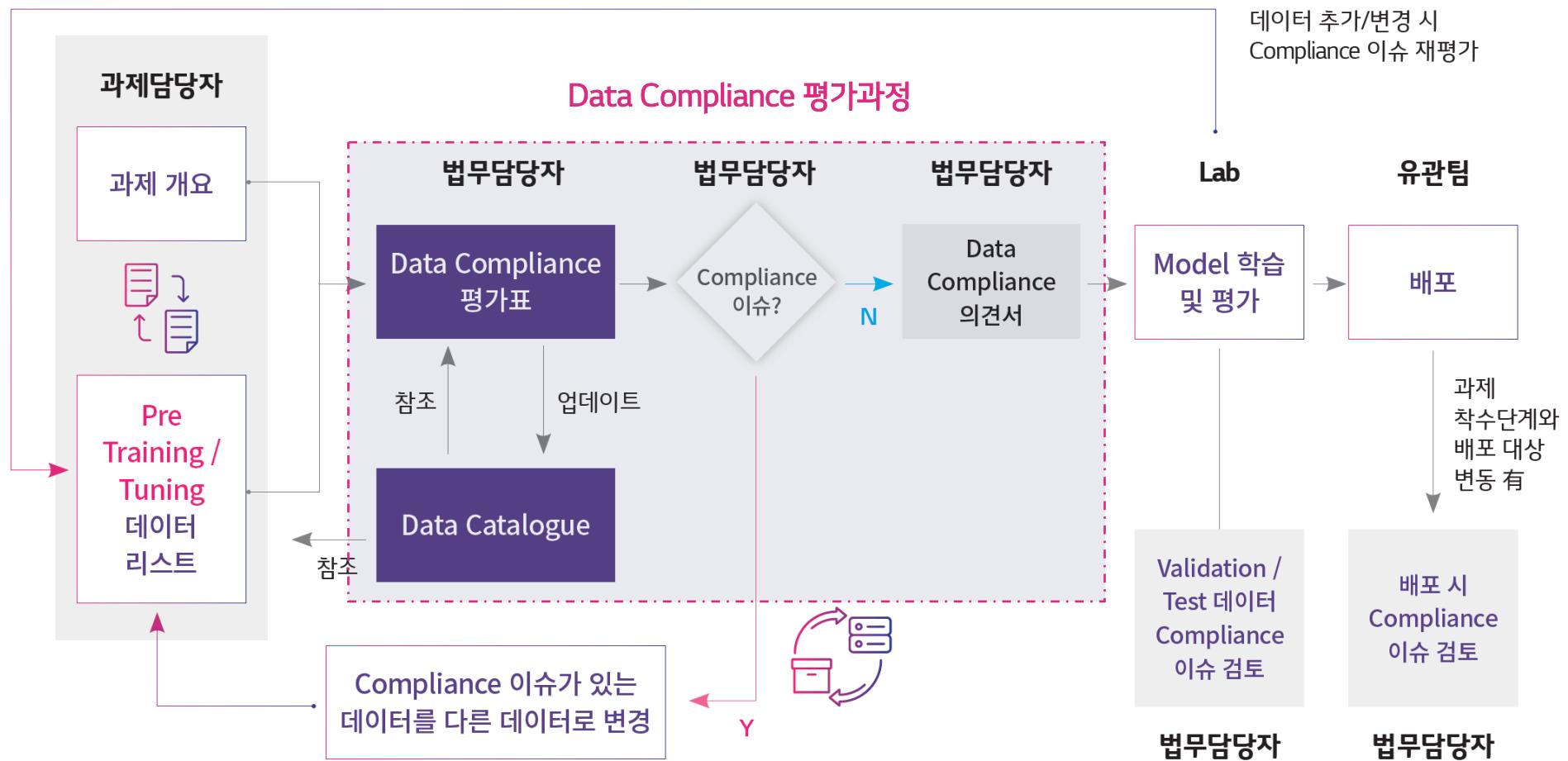
Data 구매/License-in 프로세스



## Data Compliance

- Data Compliance 프로세스를 통해, 모든 학습 데이터를 검토하고 평가합니다.

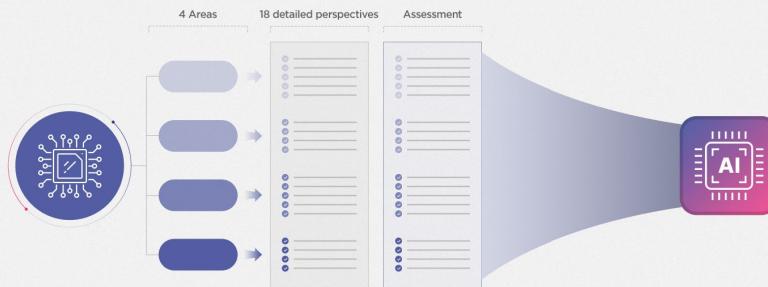
Data Compliance 프로세스



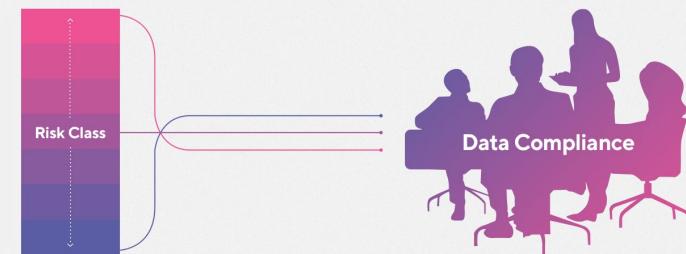
## Data Compliance

- 학습 데이터에 법적 리스크가 존재하는 경우, 법무담당자와 연구원이 함께 대안을 고민하고 해결책을 마련합니다.

We evaluate **every data** from **4 areas, 18 perspectives** and **analyze issues** by each perspective



We also assess **AI models** that have learned data using **a 7-rating scale** and work to **improve any issues**.



We engage **various stakeholders**, including legal and research professionals, in analyzing, assessing, and improving data.

### Permission to Use Data

- Granting of License Status
- Permission to Modify Data and Authority to Create Derivative Works
- Potential Copyright Infringement of Generated Outputs
- Status of Rights Assignment for Outputs
- Existence of Data Attribution Obligations

### Data Usage Period and Territory

- Limitations on Data Usage Period
- Possibility of License Revocation
- Limitations on AI Model Service Period
- Geographic Restrictions on Data Usage

### Personal Information and Data Security

- Inclusion of Personal Information or Pseudonymized Data
- Consent Status of Data Subjects
- Possibility of Third-Party Data Transfer/Provision
- Restrictions on Specific User Access to Data

### Additional Legal Risks

- Issues in the Data Collection Process
- Known Disputes Regarding AI Models Using the Data
- Additional Risks Present in License Agreements
- Types of License Conditions

## Data Compliance

- 평가의 카테고리는 총 4가지로, 1. Data License 항목, 2. Data 사용 기간 및 지역 관련 항목, 3. 개인정보 및 데이터 보안 관련 항목, 4. 추가적인 법적 리스크 관련 항목으로 분류됩니다.

1. Data License 항목		リスク				
		小			리스크	大
Data에 대한 License 부여 여부	1-1) License 부여 여부	상업적 목적으로 제한없이 사용 가능		명시적으로 내부 연구 목적으로 사용 가능	Unknown	명시적으로 사용 불가
	1-2) 데이터 수정 가능 권한 및 2차적 저작물 작성 권한	2차적 저작물 작성과 포함한 모든 수정과 변형 권한 부여		2차적 저작물 작성은 불가능하나, 2차적 저작물 작성에 이르지 않는 수준의 수정 또는 변형 가능	Unknown	명시적으로 2차적 저작물 작성을 포함한 여하한 수정과 변형 불가
원저작자의 권리 침해 가능성	1-3) 산출물(Output)의 원저작권 침해 가능성	원저작자 등의 받거나, 산출물이 생성되지 않음	원저작물과 다른 형태의 산출물을 생성됨	원저작자의 저작물과 유사한 산출물 생성 가능성 있으나, 높지 않음	원저작자의 저작물과 유사한 산출물 생성 가능성과 무관하게, 원저작자의 저작물 중 일부가 산출물에 포함될 가능성이 있음	원저작자의 저작물과 유사한 산출물 생성 가능성 높음
산출물에 대한 권리	1-4) Prompt, Output에 대한 권리 부여 여부	Prompt, Output에 대해 당사에 소유권, 지재권 있음	Prompt, Output에 대해 당사에 사용권한 있음	Prompt, Output에 대한 당사의 권리 Unknown		명시적으로 Prompt, Output에 대한 당사의 권리 없음
데이터 고지의 의무	1-5) 데이터 고지의 의무 존재 여부	(고지가능 시) 데이터 고지의 의무 존재 혹은 부존재	(고지불가 시) 데이터 고지의 의무 미존재		(고지불가 시) 별도 데이터 고지의 의무 존재	
2. Data 사용 기간 및 지역 관련 항목						
Data 사용 기간	2-1) Data 사용 기간의 제한	데이터 영구적으로 사용 가능	운영에 문제 없을 수준의 데이터 기간 제한 존재하거나 데이터 기간 존재하지 않음	데이터는 기간 제한 있으나 AI 모델에 대한 제한 Unknown		데이터 사용 가능성이 이미 경과함
	2-2) Data 라이선스 부여의 철회 가능성	라이선스 부여철회불가	Unknown	라이선스 철회가능		
AI모델 사용 기간	2-3) AI모델 서비스 기간의 제한	AI모델 서비스 영구적인 제공 가능		Unknown		AI모델 서비스 영구적인 제공 불가
Data 사용 지역	2-4) Data 사용 지역의 제한	Worldwide	Unknown		특정지역(한국, 미국, EU를 포함)에서만 사용 가능	특정지역(한국, 미국, EU를 포함하지 않음)에서만 사용 가능

## Data Compliance

- 평가의 카테고리는 총 4가지로, 1. Data License 항목, 2. Data 사용 기간 및 지역 관련 항목, 3. 개인정보 및 데이터 보안 관련 항목, 4. 추가적인 법적 리스크 관련 항목으로 분류됩니다.

### 3. 개인정보 및 데이터 보안 관련 항목

小 ————— 리스크 ————— 大

개인정보 포함여부	3-1) 개인정보의 포함 여부	개인정보 포함되어 있지 않거나, 익명 비식별화 처리 계획이 있음	개인정보 포함되었으나, 가명 처리 계획이거나 Unknown	개인정보 포함 가능성이 높음	개인정보 포함 되었음
정보주체의 동의 수취 여부	3-2) 정보주체의 동의 여부	개인정보 포함되어 있지 않거나, 개인정보 포함되어 있으나 정보주체의 동의를 받음			개인정보 포함 되었으나 정보주체의 동의를 받지 않음
가명정보 포함여부	3-3) 가명정보의 포함 여부	가명정보 포함되어 있지 않거나, 가명정보 포함되어 있으나 정보주체의 동의를 받음	개인정보 포함되어 있고 이를 가명 처리할 계획임		가명정보 포함되어 있음
Data 위탁/제3자 제공 가능 여부	3-4) Data 위탁/제3자 제공 가능 여부	Data 제3자 위탁/제공 가능 권리 있거나, 명시적 제한 없음	Unknown		명시적으로 Data 제3자 위탁/제공 불가
Data 사용권한 제한 여부	3-5) Data에 대한 특정 이용자 사용 제한	Data 이용에 대해 특정 이용자에게만 사용권한이 부여되어 있지 않음	Data 이용에 대해 특정 이용자에게만 사용권한이 부여되어 있음		

### 4. 추가적인 법적 리스크 관련 항목

Data 수집 신뢰성	4-1) Data 수집 과정에서의 적법성	적법한 방법으로 획득	Web Crawling 등을 통하여 Data 획득	Unknown	Robots.txt를 회피하거나 적법하지 않은 방법으로 획득
Data 분쟁 여부	4-2) Data가 이용된 AI모델에 대한 알려진 분쟁	알려진 Data 분쟁 없음	알려진 Data 분쟁 존재하나, 개인 대상 소액 분쟁	10억 이상의 분쟁 존재함	100억 이상의 분쟁 존재함
계약적 Risk	4-3) License 계약의 추가 리스크 존재함	알려진 추가 리스크 없음	(데이터 보안, 비밀유지의무 등) 까다로운 관리 체계 요구함	책임한도 무한, 자유로운 Audit 가능한 경우 등의 리스크 존재함	
License 조건의 유형	4-4) License 계약의 추가 리스크 존재함			재배포가 가능함	Share-Alike 수준의 요구 사항

## — Data Compliance : 1. Data License 관련 항목

### 1-1) 라이선스 부여 여부 – 데이터 사용의 적법성 판단 기준

- 데이터의 저작물성 여부 및 저작권자와의 계약 내용에 따라 상업적 이용 가능 여부가 결정되며, 이에 따라 A-1부터 C-2까지 등급을 구분함.
- 저작물로 보호되는 경우 원칙적으로 라이선스가 필요하고, 공정이용이나 TDM 면책 등 예외에 해당할 수 있으나 사전 판단이 어려움.
- 명시적 사용 허락이 존재하는 경우 가장 낮은 리스크로 평가됨.

### 1-2) 데이터 수정 권한 및 2차적저작물 작성 – 가공·변형 가능성에 대한 권리 평가 항목

- 데이터에 대한 수정 또는 2차적저작물 작성 권리가 명시적으로 부여된 경우 법적 리스크가 낮으며, 그렇지 않거나 권리 여부가 불분명한 경우 리스크가 존재함.
- 특히 동일성유지권 침해 가능성과 오픈라이선스 조건 등을 고려할 필요가 있음.
- 수정·변형이 금지된 경우 최저 점수를 부여함.

### 1-3) 산출물로 인한 분쟁 발생 가능성 – 생성형 AI 적용 시 분쟁 위험성 판단 요소

- AI가 생성한 산출물이 원저작물과 유사한 정도에 따라 법적 분쟁의 발생 가능성이 달라짐.
- 원저작자의 동의를 받았거나 산출물이 존재하지 않는 경우 리스크가 낮으며, 유사한 산출물이 생성될 가능성이 높을수록 분쟁 위험이 커짐.
- 텍스트·이미지 등 시각적으로 명확한 데이터일수록 분쟁 가능성이 높게 평가됨.

### 1-4) 산출물(Output)에 대한 권리 부여 여부 – 산출물 활용 권리의 명확성 평가 항목

- AI 모델의 산출물에 대해 소유권이나 사용권이 명시적으로 부여되어 있는지 여부에 따라 법적 리스크를 판단함.
- 산출물에 대한 권리가 명확히 당사자에게 귀속되는 경우 리스크가 낮으며, 권리 부존재 또는 명확하지 않은 경우 2차적 활용 시 문제가 발생할 수 있음.
- 특히 프롬프트 및 재학습 활용 가능성과 관련해 중요함.

### 1-5) 데이터 고지의무 존재 여부 – 고지 가능성 및 법적 부담 여부에 따른 위험도 평가 항목

- 데이터 사용 시 저작자 고지 또는 정보주체 고지의무가 존재하는 경우, 해당 고지를 이행할 수 있는지가 법적 리스크 판단에 중요함.
- 고지가 가능하면 리스크가 낮으며, 고지 의무는 있으나 고지가 불가능한 경우에는 향후 법적 문제로 이어질 가능성이 있어 낮은 점수를 부여함.
- CCL, 공정이용, 개인정보보호법 등 관련 규정이 고려됨.

## — Data Compliance : 2. 데이터 사용 기간 및 지역 관련 항목

### 2-1) 데이터 사용 기간의 제한 – 사용 가능 기간에 따른 법적 안정성 평가 항목

- 데이터 사용 기간이 영구적이거나 명시적 제한이 없는 경우 법적 리스크가 낮으며, 기간 제한이 존재하되 AI모델 사용 여부가 명확하지 않으면 중간 수준의 리스크로 평가됨.
- 반면, 데이터 사용 기간이 이미 만료된 경우 라이선스가 소멸한 것으로 간주되어 법적 사용이 불가능하므로 가장 높은 리스크로 간주됨.

### 2-2) 데이터 라이선스 부여의 철회 가능성 – 라이선스 안정성 관련 평가 항목

- 라이선스 철회가 불가능하도록 명시된 경우 법적 안정성이 높으며, 철회 가능성 여부가 불명확하거나 조건부 철회 가능성이 있는 경우 중간 리스크로 평가됨.
- 철회가 자의적으로 가능하도록 명시된 경우에는 AI모델의 지속적 사용이나 업데이트에 중대한 영향을 줄 수 있으므로 높은 법적 리스크가 발생함.

### 2-3) AI모델 서비스 기간의 제한 – AI모델 운영 지속성 평가 항목

- AI모델을 영구적으로 서비스할 수 있는 경우 안정성이 높고, 서비스 기간이 불명확하거나 제3자 계약 조건 등 외부 요인으로 인해 제한될 가능성이 있는 경우 중간 수준의 리스크로 평가됨.
- 반면, 서비스 기간 제한이 명시되어 있어 영구적 제공이 불가능한 경우 AI 비즈니스 모델 자체에 영향을 미치므로 높은 리스크로 간주됨.

### 2-4) 데이터 사용 지역의 제한 – 글로벌 서비스 적합성 평가 항목

- 데이터를 전세계에서 자유롭게 사용할 수 있는 경우 법적 제약이 없으므로 가장 안전하며, 사용 가능 지역이 불명확한 경우 중간 수준의 리스크로 평가됨.
- 데이터가 특정 지역(한국, 미국, EU 포함)에서만 사용 가능한 경우 제한된 범위 내 서비스는 가능하나 범용성은 낮아지며, 해당 주요 지역에서 사용이 불가능한 경우 AI서비스 확장이 사실상 불가능해져 매우 높은 리스크로 평가됨.

## — Data Compliance : 3. 개인정보 및 데이터 보안 관련 항목

### 3-1) 개인정보 포함 여부 – 개인정보 처리 리스크의 핵심 항목

- 학습 데이터에 개인정보가 포함되어 있지 않거나, 익명·비식별화 처리 계획이 있는 경우 법적 리스크가 낮음.
- 반면 개인정보가 명확히 포함된 경우, 적법한 수집 및 처리 근거가 요구되며, AI 산출물에 노출될 가능성까지 고려해야 하므로 리스크가 매우 높음.
- 가명처리 계획이 있거나 포함 여부가 불확실한 경우에는 중간 수준의 법적 위험이 존재함.

### 3-2) 정보주체 동의 여부 – 개인정보 처리의 적법성 판단 기준

- 개인정보가 포함되어 있더라도 정보주체의 동의를 받은 경우에는 법적으로 적법한 처리로 평가되며 리스크가 거의 없음.
- 반면 정보주체의 동의를 받지 않은 개인정보의 경우, 각국 법제의 예외 사유에 해당하지 않는 한 AI 학습에 사용하는 것은 불법 처리로 간주될 수 있어 법적 리스크가 매우 높음.

### 3-3) 가명정보 포함 여부 – 가명정보 처리에 따른 관리 필요성 평가 항목

- 가명정보가 포함되어 있지 않거나 정보주체의 동의를 받아 포함된 경우 리스크가 낮음.
- 가명처리 계획이 있는 경우 기존 개인정보보다 리스크는 낮으나, 가명정보의 적절한 처리와 안전한 활용을 위한 조치가 요구되므로 일정한 주의가 필요함.
- 반면, 가명정보가 명확히 포함된 경우에는 여전히 개인정보로서의 법적 보호가 적용되어 리스크가 존재함.

### 3-4) Data 위탁/제3자 제공 가능 여부 – 데이터 공유 및 전송의 적법성 평가 항목

- 데이터의 제3자 제공 및 위탁이 명시적으로 허용되거나 제한이 없는 경우, 활용의 자유도가 높고 법적 리스크가 낮음.
- 반면 관련 권한 여부를 알 수 없을 경우, 적법성 확인 및 관리 책임에 대한 리스크가 존재함.
- 명시적으로 위탁이나 제공이 금지된 경우, 이를 위반할 경우 계약 위반이나 법적 책임이 발생할 수 있어 리스크가 매우 높음.

### 3-5) Data에 대한 특정 이용자 사용 제한 – 접근권한 통제에 따른 운영 리스크 항목

- 데이터 사용권한이 특정 이용자에게 제한되지 않은 경우 자유로운 활용이 가능하므로 법적 리스크가 낮음.
- 반면 특정 이용자에게만 사용권한이 부여되어 있는 경우, 해당 범위를 벗어난 접근이나 활용은 계약 위반 내지 불법 처리로 판단될 수 있어 제한적 접근권한 관리가 요구되면 리스크가 존재함.

## — Data Compliance : 4. 추가적인 법적 리스크 관련 항목

### 4-1) 데이터 수집 과정에서의 이슈 – 수집 경로의 적법성 여부에 따른 리스크 평가 항목

- 데이터가 자체 생성되었거나 적법하게 라이선스를 통해 취득된 경우 리스크가 낮음.
- 반면, 웹 크롤링을 통한 수집은 각국 법령 및 사이트 정책에 위배될 가능성이 있어 중간 수준의 리스크를 가짐.
- 수집 방법이 불명확하거나 robots.txt를 무시하거나 불법적 경로로 수집된 경우, 법적 분쟁 발생 가능성성이 높아 리스크가 매우 큼.

### 4-2) 데이터가 이용된 AI모델에 대한 알려진 분쟁 – 기전 분쟁 사례에 기반한 간접 리스크 평가 항목

- 동일한 데이터나 라이선스가 적용된 다른 AI모델에 대해 분쟁 이력이 없는 경우 법적 안정성이 높음.
- 반면 소액의 개인 소송이나 고액(10억~100억 이상)의 대형 분쟁이 존재하는 경우, 동일한 원인으로 유사한 분쟁이 재발할 가능성이 높고, 특히 대형 분쟁일수록 그 영향이 커 리스크가 매우 높게 평가됨.

### 4-3) 라이선스 계약의 추가 리스크 존재 – 계약 조항 상 책임 및 의무의 부담 정도를 평가하는 항목

- 별도의 위험 요인이 없는 경우 리스크가 거의 없으나, 보안 및 비밀유지 등 관리 의무가 과도하게 요구되는 경우 운영상 부담 및 분쟁 가능성이 존재함.
- 특히, 라이센서의 책임한도가 무한하거나, 라이센서의 자유로운 감사 권한이 존재하는 경우에는 분쟁 발생시 불리한 결과로 이어질 수 있어 리스크가 높음.

### 4-4) 라이선스 조건의 유형 – 데이터셋 배포 제한의 조건에 따른 구분 항목

- 라이선스 조건은 데이터를 자유롭게 사용할 수 있는 제1유형(무제한 사용 가능), 조건부 사용 가능한 제2유형(고지 의무, 변경 고지, 동일 라이선스 조건 등), 그리고 사용은 가능하나 배포나 수정이 불가능한 제3유형으로 분류됨.
- 제2유형은 다시 세분되어 여러 조건이 중첩될 수 있으며, 라이선스 간 충돌 가능성이 있는 조건(예: 동일 라이선스 의무, sub-license 제한 등)이 포함될 경우 Mother set 단위의 활용이 제한되거나 분쟁 리스크가 커질 수 있음.
- 이러한 분류는 데이터셋의 활용 가능성과 법적 위험도를 사전에 판단하기 위한 기준으로 중요함.

## Data Compliance

- 항목별로 평가된 점수를 바탕으로, 각 항목별 가중치를 통해 해당 Data의 Risk를 정의합니다.
- Class, 총점으로 정의되며, Weight의 결정 기준은 각 국가의 입법 현황, 사법적인 판단에 기인합니다.

$R_{1-1}$  = Mandatory Data Class

$$R_i = 0.1R_{1-2} + 0.15R_{1-3} + 0.05R_{1-4} + 0.03R_{1-5} + 0.07R_{2-1} + 0.03R_{2-2} + 0.05R_{2-3} + 0.04R_{2-4} + 0.09R_{3-1} + 0.03R_{3-2} + 0.03R_{3-3} + 0.05R_{3-4} + 0.02R_{3-5} + 0.06R_{4-1} + 0.05R_{4-2} + 0.1R_{4-3} + 0.05R_{4-4}$$

$R_i$  = Risk score of data #i

### Category A. Permission to Use Data

$R_{1-2}$ = Permission to Modify Data and Authority to Create Derivative Works

$R_{1-3}$ = Potential Copyright Infringement of Generated Outputs

$R_{1-4}$ = Status of Rights Assignment for Outputs

$R_{1-5}$ = Existence of Data Attribution Obligations

### Category B. Data Usage Period and Territory

$R_{2-1}$ = Limitations on Data Usage Period

$R_{2-2}$ = Possibility of License Revocation

$R_{2-3}$ = Limitations on AI Model Service Period

$R_{2-4}$ = Geographic Restrictions on Data Usage

### Category C. Personal Information and Data Security

$R_{3-1}$ = Inclusion of Personal Information

$R_{3-2}$ = Consent Status of Data Subjects

$R_{3-3}$ = Inclusion of Pseudonymized Data

$R_{3-4}$ = Possibility of Third-Party Data Transfer/Provision

$R_{3-5}$ = Restrictions on Specific User Access to Data

### Category D. Additional Legal Risks

$R_{4-1}$ = Issues in the Data Collection Process

$R_{4-2}$ = Known Disputes Regarding AI Models Using the Data

$R_{4-3}$ = Additional Risks Present in License Agreements

$R_{4-4}$ = Types of License Conditions

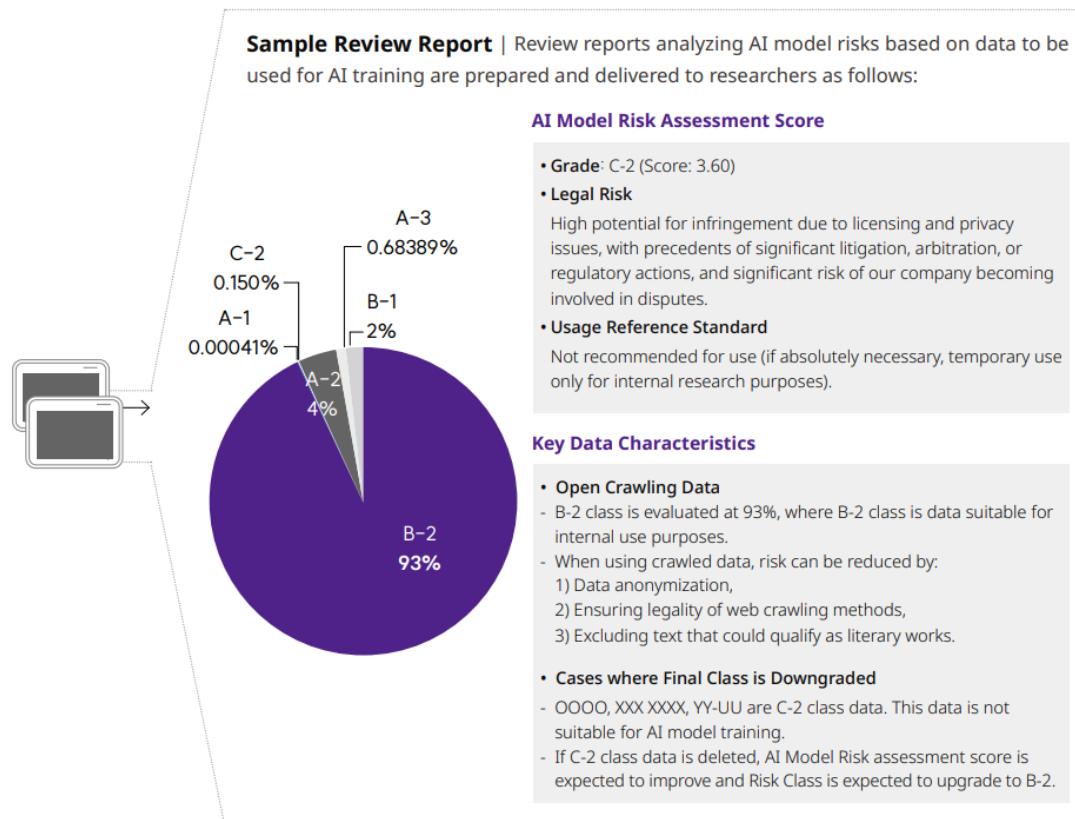
## Data Compliance

- 생성형 AI의 경우 각 Data의 Token의 가중치를 바탕으로 AI모델의 Risk Category와 Score가 계산 됩니다.
- 용도 참고별로, 서비스 전부터 각 AI모델의 Risk별 적합한 용도(제3자 제공)의 범위를 인식합니다.

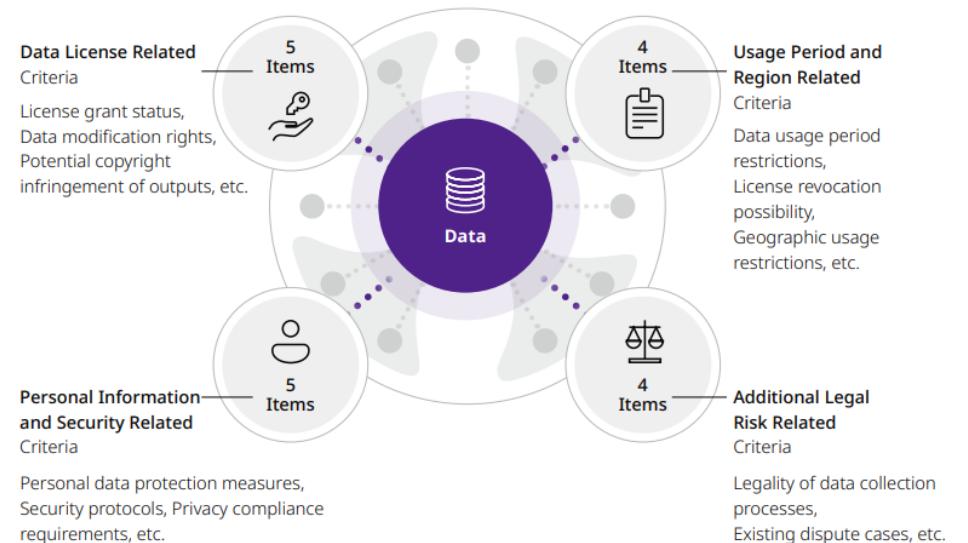
Category	Score	분류	LEGAL RISK	용도 참고
A-1	4.90	License/개인정보 -Free	자체서비스, Public Cloud에서 AI Model 서비스 등 Public하게 공개되어도 원저작자나 Licensor, 정보주체, 유관기관으로부터 Data와 관련한 분쟁의 가능성은 없음	B2C
A-2	4.57	License/개인정보 위험도 小	License/개인정보 이슈가 존재하여 침해의 가능성이 약간 존재 가능하나 소송이나 중재 등 분쟁 또는 유관기관의 처분 등으로 심화된 알려진 Case 없음	
A-3	4.22	License/개인정보 위험도 小	License/개인정보 이슈가 존재하여 침해의 가능성이 상당히 존재 가능하나 소송이나 중재 등 분쟁 또는 유관기관의 처분 등으로 심화된 알려진 Case 없음	
B-1	3.73	License/개인정보 위험도 中	License/개인정보 이슈가 존재하여 침해의 가능성이 높으며, 소송이나 중재 등 분쟁 또는 유관기관의 처분 등이 발생한 Case가 존재하고, 당사도 분쟁에 휘말릴 가능성이 약간 존재함	B2B
B-2	3.51	License/개인정보 위험도 中	License/개인정보 이슈가 존재하여 침해의 가능성이 높으며, 소송이나 중재 등 분쟁 또는 유관기관의 처분 등이 발생한 Case가 수 건 있고, 당사도 분쟁에 휘말릴 가능성이 약간 존재함	
C-1	3.18	License/개인정보 위험도 大	License/개인정보 이슈가 존재하여 침해의 가능성이 크게 높으며, 소송이나 중재 등 분쟁 또는 유관기관의 처분 등이 발생한 Case가 상당히 있고, 당사도 분쟁에 휘말릴 가능성이 존재함	Internal
C-2	-	License/개인정보 위험도 大	License/개인정보 이슈가 존재하여 침해의 가능성이 크게 높으며, 거액의 소송이나 중재 등 분쟁 또는 유관기관의 처분 등이 발생한 Case가 있고, 당사도 분쟁에 휘말릴 가능성이 크게 존재함	

## Data Compliance

- Data 뿐만 아니라, 평가된 Data를 이용하여 개발된 모든 AI 모델의 리스크를 평가하여 리스크를 Classification 하고, 용도에 맞게 AI를 이용합니다.
- 지속적인 Data 및 AI의 Data Compliance 체계 수행으로 저작권법, 개인정보보호법, 부정경쟁방지법 준수와 다가올 국내외 AI 기본법을 대비하고자 합니다.



**Evaluation Criteria: 18 Legal Perspectives** | Data is protected by various laws including copyright law, personal information protection law, and unfair competition prevention law in each country, with legally permissible uses varying by jurisdiction. Considering these characteristics, our AI-based data compliance system reviews potential risks across 18 legal perspectives:



**Data Risk Classification Table** | Datasets are classified into three major grades (A, B, C) and seven sub-categories based on risk assessment results.

Category	License/Privacy	Key Legal Risks
Grade A	Risk-Free	Very low likelihood of legal disputes or risks.
Grade B	Medium Risk	High likelihood of violations related to licenses or privacy issues.
Grade C	High Risk	Very high likelihood of legal dispute escalation.

Chapter

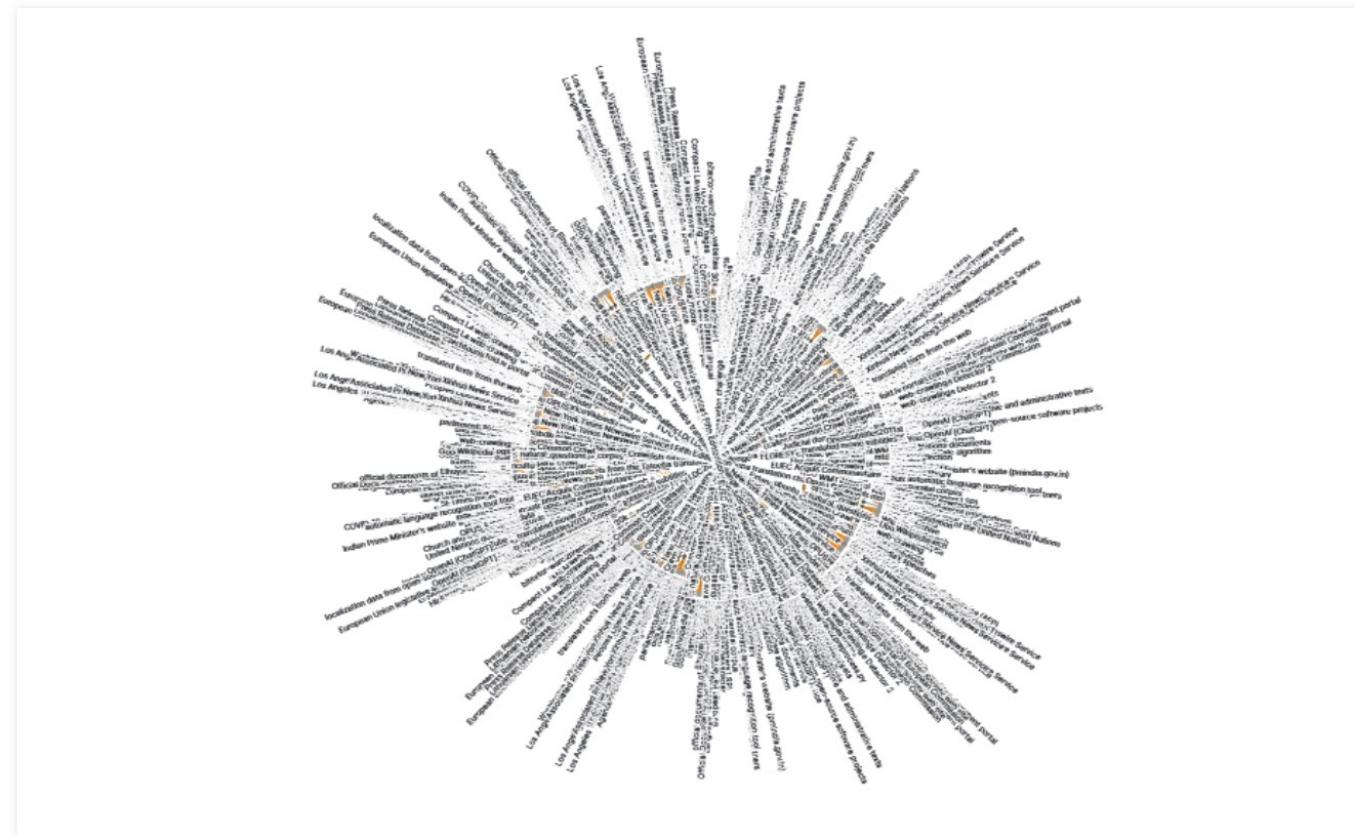
---

3

Data Compliance  
Agent : NEXUS

## 현대 AI 학습 데이터 검토의 어려움

- 엄청난 파라미터 사이즈를 갖는 모델들이 주류가 되면서, 학습 데이터셋 또한 완전 새로운 데이터로만 구성된 단순 구조가 아닌, 성능 향상에 필요한 다양한 소스부터 수집된 데이터들을 복합적으로 섞은 거대한 수직 계층형 구조를 갖게 되었습니다.
- 상업적으로 이용 가능하다고 판단된 2,852개의 AI 학습 데이터셋 중 종속 데이터의 리스크를 모두 고려해 본 결과, 21.21%인 605개의 데이터셋만 상업적으로 이용 가능했습니다.



Mean	Std	Min	25%	50%	75%	Max
2.20	1.83	0	1	2	3	16

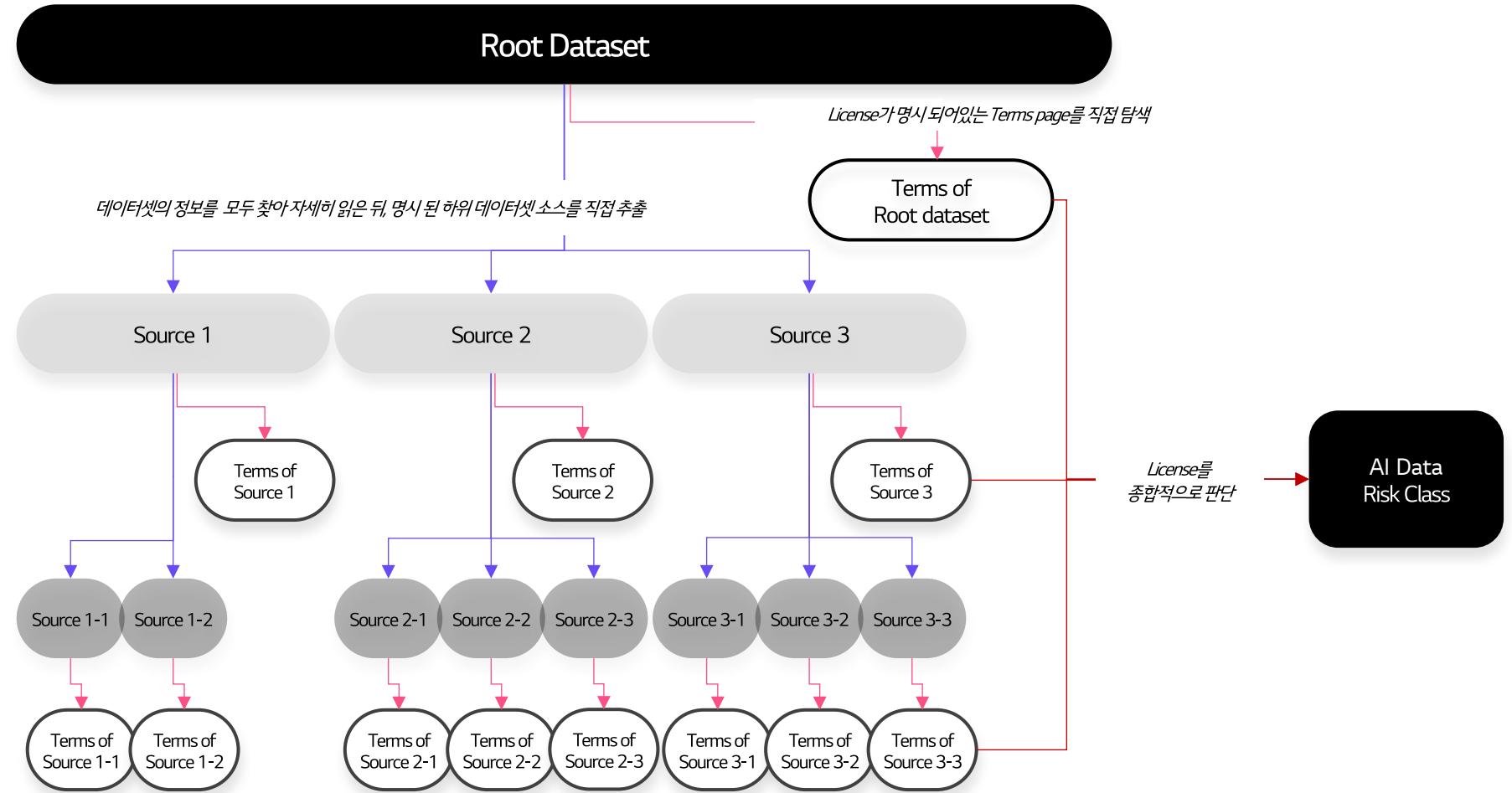
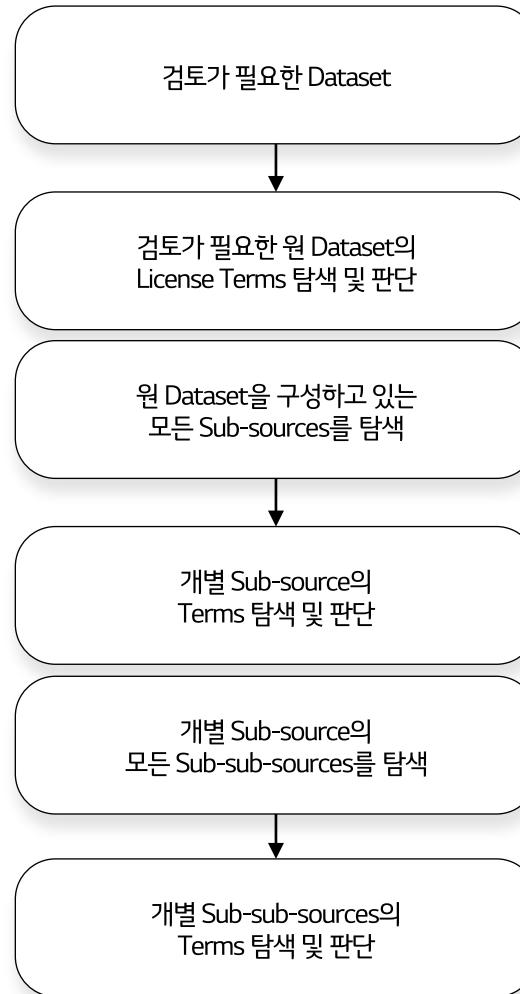
3,612 데이터셋의 계층 레벨(Depth) 동계

Mean	Std	Min	25%	50%	75%	Max
23.18	87.15	1	2	5	11	1691

3,612 데이터셋의 Data Source 동계

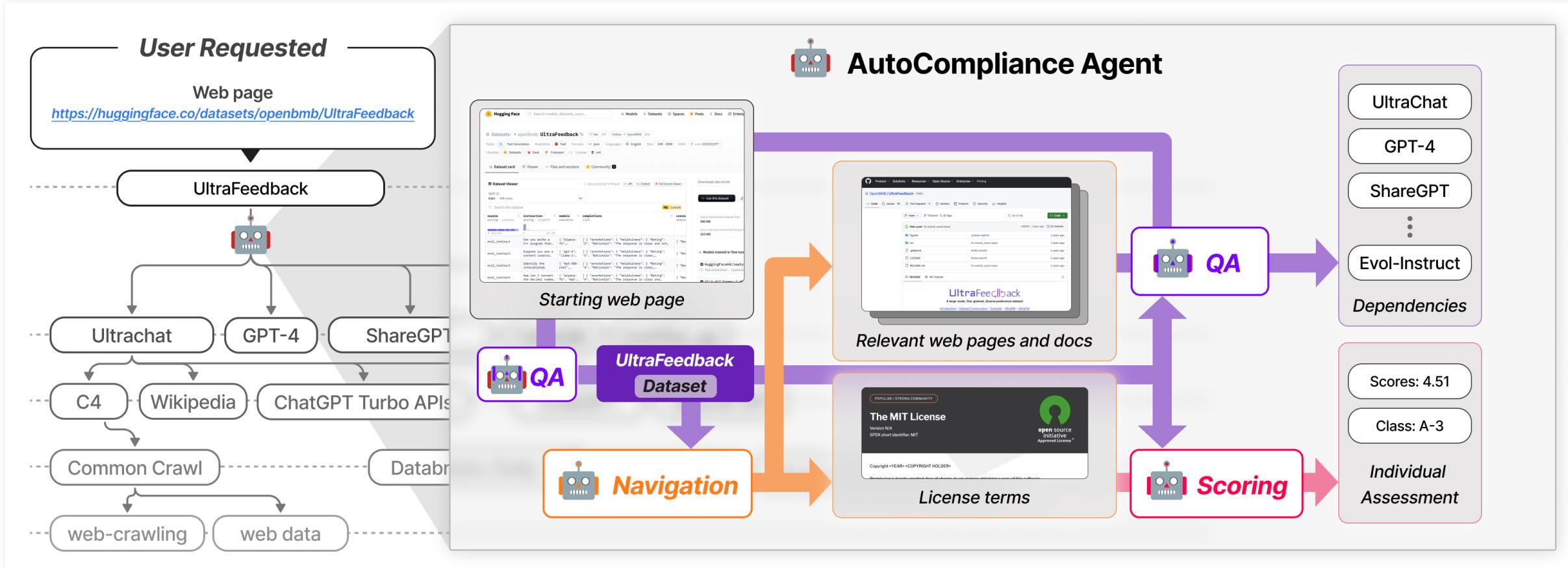
## 현대 AI 학습 데이터 검토의 어려움

- 실제로 Open-sourced 학습데이터를 이용하는 경우, 해당 학습데이터의 모든 원본 데이터에 대한 법적 리스크를 검토해야 하는 어려움이 존재합니다.



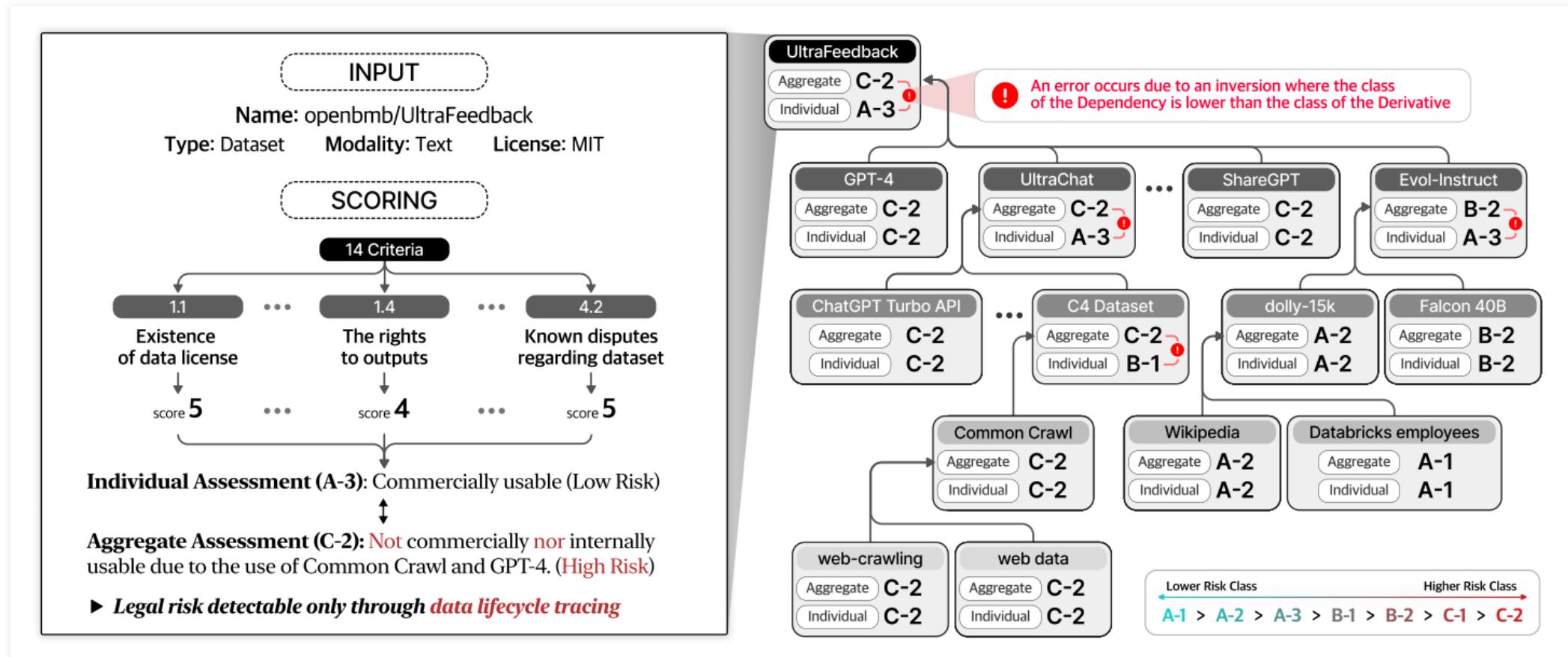
## NEXUS를 활용한 Data Compliance 자동화 구조

- NEXUS는 EXAONE 3.5 32B를 기반으로 개발되었습니다.
- QA, Navigation, Scoring 모듈을 통해 데이터셋의 정보를 바탕으로 어떤 하위 데이터 소스로 데이터셋이 구성되어 있는지, 각 데이터 셋의 라이선스 정보가 어떻게 되는지를 자동적으로 탐색합니다.



## NEXUS를 활용한 Data Compliance 자동화 구조

- NEXUS는 사용하고자 하는 단일 데이터 셋이 아닌, 데이터 셋이 갖고 있는 모든 출처가 되는 데이터 셋의 법적 위험을 Data의 Life Cycle 측면에서 탐지합니다.



## NEXUS 성능

- NEXUS는 인간 변호사와 비교했을 때 45배 빠른 속도, 0.1% 수준의 비용, 26% 빠른 정확도를 보였습니다.

Round	작성자	라이선스 데이터셋 작성	소요 시간	소요 금액
1	법무법인 변호사 5인	200개 작성 및 검토, 녹화 (HF Most Download 1-200)	29시간 / 인	62,634,000
2	법무법인 변호사 5인		(Total : 145시간)	
3	법무법인 변호사 5인	50개 Motherset만 재작성 (HF Most Download 1-50)	2.26시간 / 인 (Total : 11.3시간)	5,130,000
4	법무법인 변호사 5인	45개 (HF Most Download 201-245)	3.94시간 / 인 (Total : 19.7시간)	8,190,000
	Actionable AI	45개 (HF Most Download 201-245)	Total 0.9시간	\$2.5
5	법무법인 변호사 5인	155개 (HF Most Download 246-400)	10.02시간 / 인 (Total : 50.1시간)	20,900,000

Set Accuracy (↑)		
Name	Dependencies	License terms
<b>AutoCompliance Agent</b>	<b>81.04%</b>	<b>95.83%</b>
Human expert	64.19%	87.73%
ChatGPT-4o	25.00%	39.81%
Perplexity Pro	28.24%	22.22%

Name	Time (sec)	Cost (\$)
<b>AutoCompliance Agent</b>	<b>53.1</b>	<b>0.29</b>
Human expert	2,418	207

## NEXUS Web Repository (<http://nexus.lgresearch.ai/>)

- NEXUS Web Repository를 통해 실증 연구에서 NEXUS를 통해 검토한 3,612개 데이터셋의 검토 결과를 직접 확인 할 수 있습니다.

The screenshot shows the homepage of the NEXUS Web Repository. At the top left is the "NEXUS" logo. To its right are links for "Legal Framework" and "Contributors". In the top right corner is a circular button with the text "↗ LG AI Reserach". Below the header, the text "AUTOMATED DATA COMPLIANCE SYSTEM" is displayed. The main feature is a large, bold, black headline: "Decode the Data, Unveil the Unknown.". Below the headline is a descriptive paragraph: "Explore insights from LG AI Research's Agent AI, which proactively assesses legal risks and provides insights through deep data analysis." At the bottom center is a black button with the text "↗ Go to Paper". A cursor arrow is visible at the bottom center of the page.

NEXUS

Legal Framework Contributors

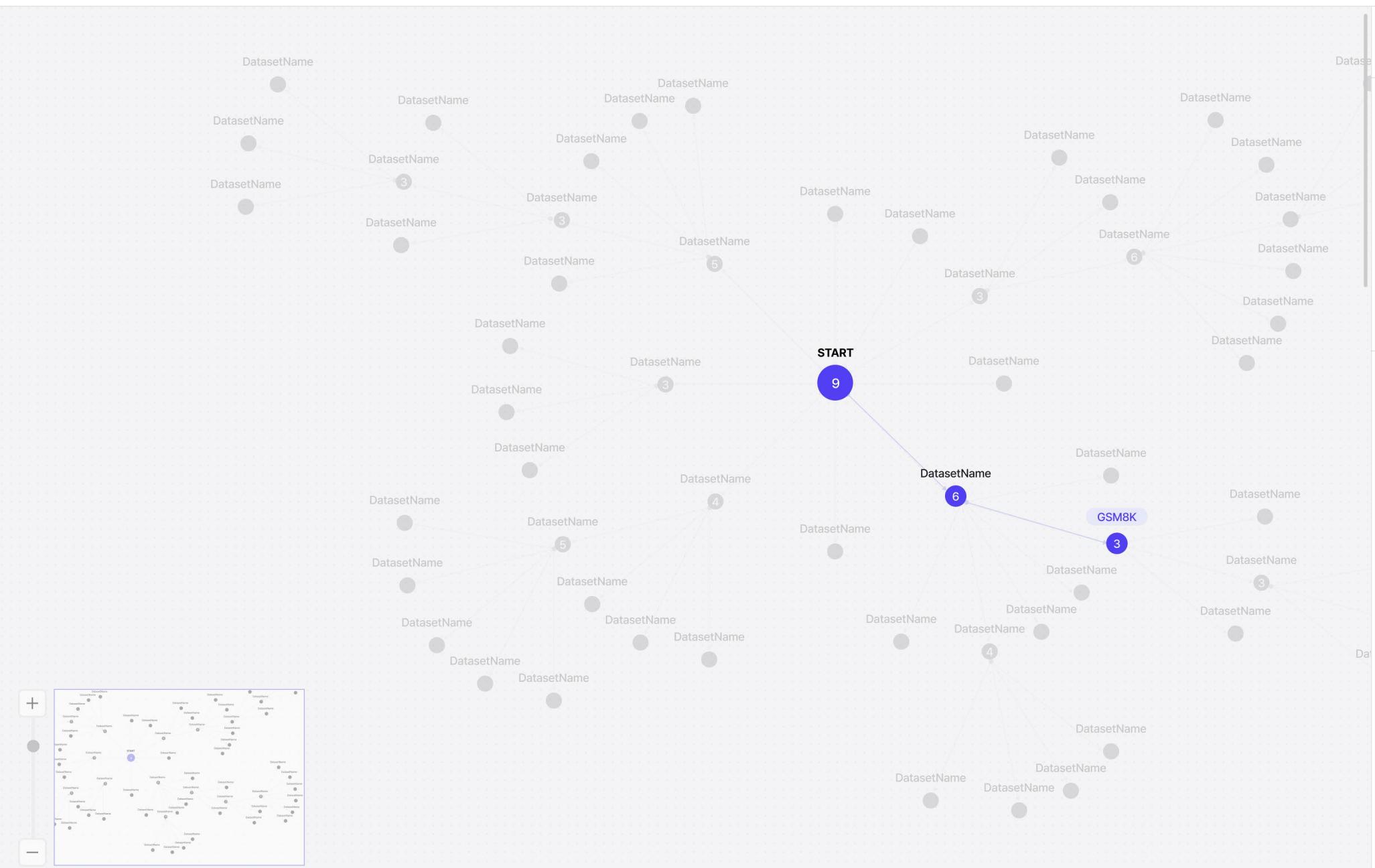
↗ LG AI Reserach

AUTOMATED DATA COMPLIANCE SYSTEM

# Decode the Data, Unveil the Unknown.

Explore insights from LG AI Research's Agent AI, which proactively assesses legal risks and provides insights through deep data analysis.

↗ Go to Paper



## GSM8K

### License 정보

Data URL

GSM8K

License URL

MIT

Source Name

news articles news

blog posts

### Subset 정보 (7)

1 AQuA-RAT

Apache 2.0

2 TheoremQA

Apache 2.0

3 GSM8K-RFT

Non Listed

4 Out Curated

Apache 2.0

5 NumGLUE

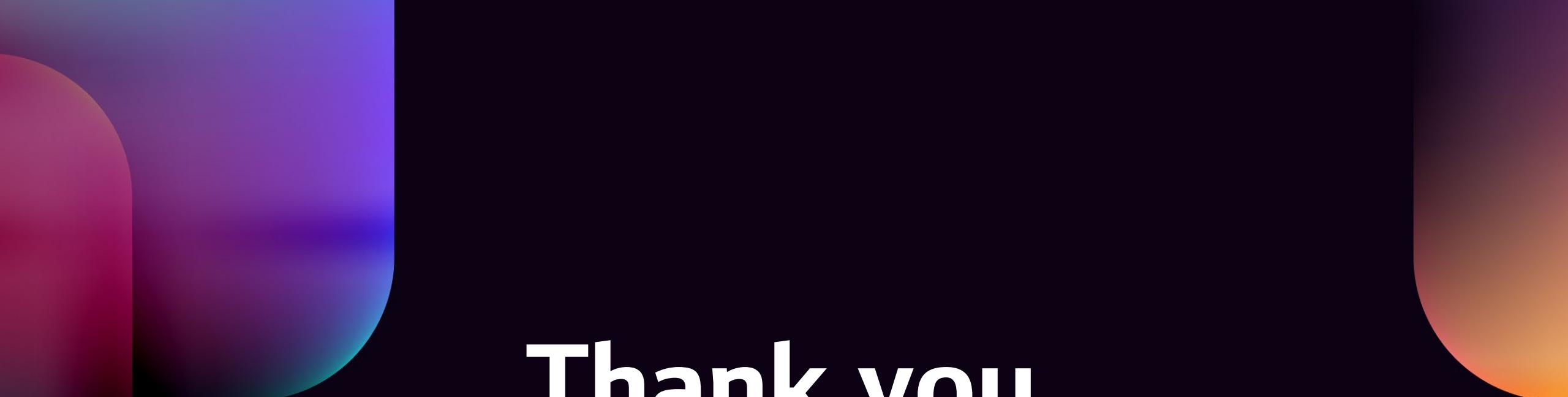
Apache 2.0

6 AQuA-RAT

Apache 2.0

7 AQuA-RAT

Apache 2.0



# Thank you

[legal@lgresearch.ai](mailto:legal@lgresearch.ai)