#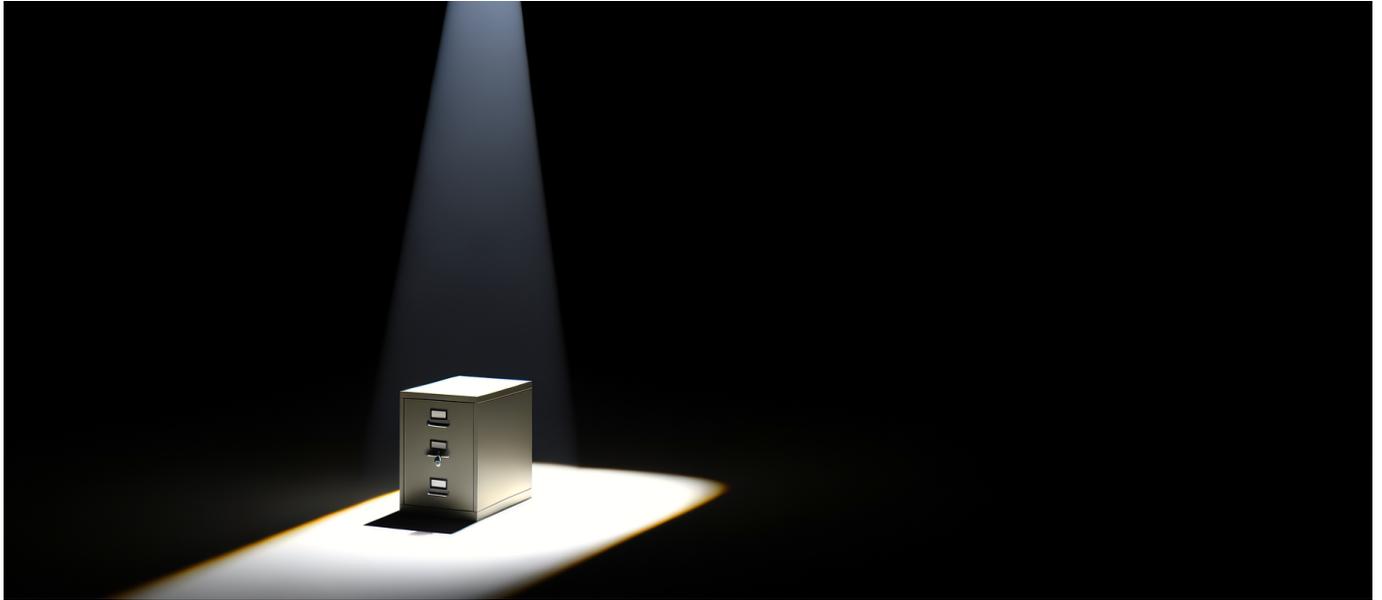 Congress Introduces H.R. 7209 TRAIN Act: Copyright Holders Can Now Subpoena AI Training Data with Court-Issued Warrants

For the first time in U.S. history, individual copyright holders can legally compel OpenAI, Anthropic, and Google to reveal whether their specific works were used to train AI models. Refuse to answer, and courts will presume you're guilty of copying.

## The News: A Legal Weapon Takes Shape

On January 22, 2026, Representatives Madeleine Dean (D-PA) and Nathaniel Moran (R-TX) introduced H.R. 7209, the Transparency, Responsibility, and Accountability for AI Training Act—better known as the TRAIN Act. This bipartisan bill creates an administrative subpoena process that fundamentally alters the legal calculus for every company building generative AI systems.

The mechanism is brutally simple. Any copyright owner can file a sworn declaration

with a U.S. district court clerk stating a good-faith belief that their copyrighted work was used to train an AI model. The clerk issues the subpoena. The AI developer must respond.

Here's where it gets teeth: [per the bill text](#), developers who fail to comply face a "rebuttable presumption of copying." In plain English, if you don't answer the subpoena, courts will legally assume you stole the work until you prove otherwise. That's not a procedural nuisance—it's a complete inversion of how copyright litigation has functioned for decades.

The bill applies to "generative AI models including original, retrained, or fine-tuned versions." This language matters enormously. It covers not just foundation models but every derivative work built on top of them. Fine-tuned your medical diagnosis assistant on GPT-4? You're potentially in scope. Retrained an open-source model on proprietary data? Same rules apply.

One critical constraint prevents fishing expeditions: subpoenas are limited to the requester's own copyrighted works. A photographer can ask if their photos were used. They cannot demand a full inventory of every Getty Images photograph in the training set. Bad-faith requests trigger sanctions under Federal Rule of Civil Procedure 11, which carries real financial penalties for attorneys and clients who abuse the process.

The effective date deserves attention. Unlike most federal legislation, H.R. 7209 takes effect immediately upon enactment with no phase-in period. If this passes, the subpoenas can start the next day.

## Why This Changes Everything

The TRAIN Act represents the first federal legislative mechanism giving individual creators legal discovery power over AI training data. Before this bill, creators had one path to learn if their work trained AI: file a lawsuit and hope for discovery. That required substantial legal resources, a viable infringement theory, and the willingness to bet years of litigation costs on an uncertain outcome.

Now? A freelance illustrator in Tulsa can file a sworn declaration with the district court clerk for the cost of a filing fee. Major AI labs will face subpoenas from thousands of individual creators, each demanding confirmation about their specific works.

[The National Law Review analysis](#) calls this "pulling back the curtain" on AI training data. That's underselling it. This bill hands copyright holders a crowbar.

## The Winners

**Individual creators** gain asymmetric power they've never possessed. A novelist suspecting their books were scraped no longer needs to speculate. They can demand answers with legal backing.

**Copyright collectives and guilds** will become coordinating mechanisms for mass subpoena campaigns. Expect the Authors Guild, SAG-AFTRA, and professional photography associations to create streamlined filing systems for their members.

**Litigation finance firms** will smell blood. The rebuttable presumption provision transforms marginal infringement cases into viable investments. If AI companies refuse to comply with subpoenas, plaintiffs enter court with the legal equivalent of a head start.

**Enterprise AI vendors with clean data** suddenly have a competitive moat. If you can prove your training data was properly licensed, you can market compliance as a feature while competitors scramble.

## The Losers

**Foundation model providers** face an existential documentation challenge. Can OpenAI demonstrate definitively whether any specific work from the pre-2022 internet was or wasn't in their training corpus? For models trained on Common Crawl data at scale, that question may not have a clean answer.

**Open-source model ecosystems** inherit liability risk they never anticipated. If Mistral or Meta release weights for a model trained on questionable data, every downstream user of those weights becomes a potential subpoena target.

**Startups using third-party APIs** face a legal supply chain problem. Your fine-tuned model sits on top of Claude or GPT-4. The foundation model gets subpoenaed. Does your compliance obligation extend to data you never controlled and can't audit?

# Technical Implications: Why Compliance Is Harder Than It Sounds

Let's get concrete about what this bill actually demands from an engineering perspective.

Modern large language models are trained on datasets measured in terabytes of text. GPT-4's training data is estimated at 1-2 trillion tokens. Claude's training corpus has never been publicly detailed. These datasets were assembled through web scraping, licensed data partnerships, books3-style torrent aggregations, and dozens of other sources that predate any serious provenance tracking.

**The core problem: most foundation models cannot generate an authoritative list of every document in their training set.**

This isn't negligence—it's an artifact of how the field developed. When researchers at Google published the original Transformer paper in 2017, nobody was thinking about per-work attribution. When OpenAI trained GPT-2 in 2019, the legal consensus held that training on copyrighted data likely constituted fair use. Data governance was an afterthought.

Consider what compliance actually requires:

**1. Deduplication and identification.** Training datasets contain duplicates, near-duplicates, and fragments. If someone subpoenas asking whether their 50,000-word novel was used, you need systems that can identify partial matches, reformatted versions, and excerpts. This is a non-trivial semantic search problem across petabyte-scale data.

**2. Provenance chains for derived data.** Many training sets are aggregations of other datasets. Common Crawl contains data from web pages that themselves contain quoted or copied content. If a blog post republished an entire New York Times article, and that blog post was in Common Crawl, and Common Crawl was in your training mix—was the New York Times article in your training data? Technically yes. Can you detect that? Probably not without substantial new infrastructure.

**3. Model-level attribution.** Even if you know what was in the training data, can you prove a specific work influenced model weights? For some embedding-based

retrieval approaches, this is tractable. For models trained end-to-end with gradient descent across trillions of tokens? The influence of any single document on final model behavior is genuinely difficult to quantify.

**4. Fine-tuning lineage.** The bill explicitly covers "retrained or fine-tuned versions." If you fine-tuned a model, you need complete documentation of your fine-tuning data. If you fine-tuned a model that someone else fine-tuned, you need their documentation too. Supply chains get long fast.

## What Good Compliance Architecture Looks Like

Organizations that survive this regulatory environment will have built systems with several key properties:

**Immutable data manifests.** Every training run gets a cryptographically signed manifest of every document in the dataset, with content hashes, source URLs, access timestamps, and license metadata. This needs to be baked into your training pipeline infrastructure, not bolted on afterward.

**Content fingerprinting at ingestion.** When data enters your pipeline, generate robust perceptual hashes (for images) or semantic fingerprints (for text) that can match against query works even when formatting differs.

**Audit-ready query interfaces.** Legal teams will need to run queries like "was any work by author X published between dates Y and Z in training set W?" in reasonable time. This is an internal search engine problem that most organizations haven't invested in.

**Negative certificates.** Equally important: the ability to prove something was NOT in your training data. This requires your data manifest to be complete and auditable, which means tracking exclusion decisions as carefully as inclusion decisions.

The engineering lift here is substantial. For organizations that haven't been tracking this data, there may be no reconstruction path. You cannot retroactively build a manifest for a training set you assembled three years ago and didn't document.

# How This Differs From California's AB 2013

Context matters here. California's AB 2013 went into effect on January 1, 2026—just three weeks before the TRAIN Act introduction. Both address AI training transparency. They are not the same.

[AB 2013 requires public summaries](#) of training data categories. AI providers must publish high-level descriptions of what types of data trained their models. "We used publicly available web data, licensed book corpora, and user-submitted prompts" satisfies AB 2013.

The TRAIN Act requires work-specific disclosure in response to individual subpoenas. "Did you train on THIS SPECIFIC PHOTOGRAPH, image ID 847293847, taken by John Smith on April 15, 2019?" is the level of granularity in play.

AB 2013 is a disclosure regime. The TRAIN Act is a discovery mechanism.

The compliance burden differs by orders of magnitude. Meeting AB 2013 requirements is a documentation and PR exercise. Answering TRAIN Act subpoenas requires technical infrastructure that most organizations don't have.

There's also a jurisdiction difference worth noting. AB 2013 applies to AI systems "deployed in California" or serving California residents—essentially any commercial AI product. The TRAIN Act, if passed, would apply federally but only activates when a copyright holder files. Companies could theoretically ignore AB 2013's broader transparency requirements while still facing targeted TRAIN Act subpoenas.

# The Contrarian Take: What Most Coverage Gets Wrong

The legal commentary I've read focuses on this bill as a win for copyright holders and a loss for AI companies. That framing misses the more interesting dynamics at play.

**This bill may actually help well-resourced AI labs and hurt their smaller competitors.**

OpenAI, Anthropic, and Google have the engineering talent, legal resources, and

operational scale to build compliance infrastructure. They've also been negotiating data licensing deals for years, giving them at least partial cover for portions of their training data. The TRAIN Act creates massive fixed costs that they can absorb and their competitors cannot.

Consider a startup that fine-tunes open-source models for vertical applications. They didn't curate the base model's training data. They may not even have access to documentation about what was in it. Under the TRAIN Act, they face identical subpoena obligations to companies that trained their own models from scratch. But they have no way to answer.

If you're Anthropic, this might be the best competitive moat you never had to build. Every model trained on sketchy data becomes legally radioactive. The "train on everything and ask forgiveness later" era ends. The companies who were already doing compliance work gain advantage.

**The rebuttable presumption provision may not survive constitutional scrutiny.**

Creating a legal presumption of copying based on failure to respond to a subpoena raises due process questions that legal scholars will litigate for years. The Fifth Amendment protects against compelled self-incrimination. While this applies more cleanly to criminal contexts, forcing disclosure that could establish civil liability, or else face presumptive guilt, sits in genuinely uncertain constitutional territory.

I expect this provision specifically to face legal challenges. It may be severed from the rest of the bill even if the subpoena mechanism survives.

**The "sworn declaration" requirement is doing less work than it appears.**

The good-faith belief standard is subjective. A creator who publishes their work online can reasonably believe it was scraped, because most text and images online were scraped into someone's training set at some point. This isn't bad faith—it's statistical probability.

The limitation to "only their own works" prevents true fishing expeditions, but it doesn't prevent coordinated campaigns. If 50,000 members of a photographers' guild each file declarations about their own works, you have 50,000 subpoenas. The compliance burden scales linearly with the number of motivated creators. There are

a lot of motivated creators.

# Practical Implications: What Technical Leaders Should Do Now

If you're a CTO, VP of Engineering, or technical founder at a company building or deploying generative AI, here's your action list:

## Immediate (Next 30 Days)

**Audit your training data documentation.** Do you know what's in your training sets? Can you produce a manifest? If you fine-tuned third-party models, do you have contractual provisions giving you access to base model documentation?

**Identify your exposure surface.** How many copyrighted works are plausibly in your training data? What categories of creators are most represented? Photographers, writers, musicians, and visual artists have different organizing structures and different likelihoods of filing subpoenas.

**Review your legal counsel's AI expertise.** Copyright law around AI training is genuinely novel. Generalist IP attorneys may not have current expertise. If your counsel's guidance on this bill is "wait and see what happens," find additional counsel.

## Medium-Term (Next 90 Days)

**Build or buy compliance infrastructure.** Content fingerprinting, manifest generation, and audit interfaces are table stakes if this passes. The build vs. buy decision depends on your scale. Vendors in the data governance space will launch TRAIN Act compliance products within weeks of passage.

**Evaluate your model supply chain.** If you depend on third-party models, understand their training data practices and contractual obligations. Request documentation proactively. If they can't provide it, assess whether that dependency is sustainable.

**Model the cost of non-compliance.** What does a rebuttable presumption of copying cost you in actual litigation? For many companies, the answer is

"potentially existential." For others, it's a manageable settlement cost. Know which category you're in.

### Strategic (Next 12 Months)

**Consider synthetic and licensed-only training approaches.** Models trained exclusively on synthetic data, properly licensed datasets, or public domain materials face zero exposure under the TRAIN Act. This has always been the cleanest legal position. It's becoming the only defensible one.

**Watch for legislative amendments.** This bill is two days old. It will change as it moves through committee. The rebuttable presumption provision is the most aggressive element and the most likely to be modified. The subpoena mechanism itself has broader support.

**Invest in provenance technology.** Content authenticity standards (C2PA, etc.), blockchain-based attribution systems, and other provenance technologies become more valuable in a world where you need to prove what you trained on. The infrastructure investments you make now may become compliance requirements later.

# Where This Goes in 6-12 Months

Prediction is dangerous, but informed speculation is useful. Here's my read on likely developments:

**The TRAIN Act passes in some form by Q3 2026.** Bipartisan AI legislation is rare enough that sponsors will work hard to keep consensus. The mechanism is simple enough to explain to non-technical legislators. Copyright protection polls well. Opposing this bill means voting against creators and for Big Tech. Few representatives will take that trade.

**The rebuttable presumption provision gets modified.** Either in committee or through early legal challenges, the automatic presumption of copying will soften. Likely compromise: a presumption of "good cause" for discovery rather than presumption of infringement itself. Still powerful, less constitutionally fraught.

**We see the first subpoena campaigns within 60 days of passage.** Organized groups of creators—probably visual artists or journalists—will file coordinated

subpoenas against major labs as both practical discovery and symbolic action. The responses (or non-responses) will set early precedent.

**At least one major litigation settlement references TRAIN Act subpoena responses.** The Sarah Silverman case, the New York Times case, or one of the other pending lawsuits will incorporate evidence obtained through this mechanism. This establishes the practical value of the law beyond its deterrent effect.

**Open-source model release practices change significantly.** Meta, Mistral, and others releasing model weights will face pressure to release training manifests alongside them. "Open weights without open data documentation" becomes legally hazardous for anyone deploying those weights commercially.

**A secondary market emerges for "clean" training data.** Data brokers who can certify licensing provenance for training-ready corpora will command premium prices. The existing market for licensed data accelerates. Ventures in this space become significantly more fundable.

**International coordination questions arise.** Models trained overseas but deployed in the U.S. face TRAIN Act subpoenas. How do you compel compliance from a company headquartered in China or the UAE? The enforcement mechanics get messy quickly.

# The Deeper Structural Shift

Zoom out from the legal mechanics and consider what the TRAIN Act represents structurally.

For the past decade, AI development has operated on an implicit social contract: researchers could train on the internet's accumulated knowledge, and society would benefit from the resulting capabilities. That contract was never explicitly negotiated. Creators never consented. But enforcement mechanisms didn't exist, so the training continued.

The TRAIN Act is one of several simultaneous signals that this implicit contract is being renegotiated. The European AI Act imposes training data requirements. Multiple ongoing lawsuits are testing fair use boundaries. Licensing deals between AI companies and publishers suggest that even model developers see the legal landscape shifting.

We're moving from a regime of "train first, ask questions never" to one of "document everything, prove your licensing." That transition will be expensive, slow, and incomplete. Models trained before documentation requirements existed cannot retroactively become compliant.

This creates a temporal discontinuity in the AI industry. Pre-2026 models carry legal risk that cannot be fully mitigated. Post-2026 models can be built with compliance in mind from day one. The transition period will be messy for everyone.

For technical leaders, the strategic question isn't whether to comply—it's how fast you can build systems that make compliance tractable. The companies that solve this problem first gain a moat that no amount of GPU spend can replicate.

**The era of AI training as a purely technical challenge is over—it's now a legal and operational one, and the organizations that adapt fastest will define the industry's next decade.**