

Adverse Selection in the AI Data Commons

Evidence from News and Media Sites

Kai Zhu

Bocconi University

Digital Economy Workshop 2026

The AI Training Data Problem

- Generative AI derives its power from training on **high-quality web content**
(Gunasekar et al. 2023; Longpre et al. 2023)

The AI Training Data Problem

- Generative AI derives its power from training on **high-quality web content** (Gunasekar et al. 2023; Longpre et al. 2023)
- Yet AI companies use this content **without compensation**
Rare exceptions:
 - OpenAI–Reddit \$60M deal
 - OpenAI–Associated Press licensing agreement

The AI Training Data Problem

- Generative AI derives its power from training on **high-quality web content** (Gunasekar et al. 2023; Longpre et al. 2023)
- Yet AI companies use this content **without compensation**
Rare exceptions:
 - OpenAI–Reddit \$60M deal
 - OpenAI–Associated Press licensing agreement
- No established market: producers' only option is **binary opt-out** via `robots.txt` (Kim et al. 2025; Chang & He 2025)

The AI Training Data Problem

- Generative AI derives its power from training on **high-quality web content** (Gunasekar et al. 2023; Longpre et al. 2023)
- Yet AI companies use this content **without compensation**
Rare exceptions:
 - OpenAI–Reddit \$60M deal
 - OpenAI–Associated Press licensing agreement
- No established market: producers' only option is **binary opt-out** via `robots.txt` (Kim et al. 2025; Chang & He 2025)

Who opts out? If the highest-quality producers exit first, the result is a **lemons problem** (Akerlof, 1970)—and the remaining AI data commons systematically degrades.

Research Questions & Preview of Findings

RQ: Does the market for AI training data exhibit adverse selection?

Key findings:

① **Quality-blocking gradient**

High-factual media block at **6×** the rate of low-factual sources

② **Strategic targeting**

High-quality sites selectively block training bots over search bots

③ **Misinformation asymmetry**

Conspiracy, propaganda, pseudoscience stay open for AI crawlers

④ **Ideological sorting**

Centrist outlets block most; both left and right extremes block least

⇒ Evidence from **9,600 news and media sites** with two complementary quality measures: expert ratings and web-graph metrics

Institutional Background

robots.txt

- Voluntary protocol since 1994
- Sites specify which bots to block
- Advisory, not legally enforced
- Can target *specific* AI crawlers

AI Crawlers

- **Training:** GPTBot, Google-Extended, ClaudeBot
Extract content to build models
- **Search:** ChatGPT-User, PerplexityBot, Applebot-Extended
Retrieve content to answer queries — can drive traffic back

Key events:



Why media & news?

- News content shapes public discourse and how AI models handle factual claims
- Expert quality ratings exist—enabling direct measurement of adverse selection
- Center of active policy debate (NYT v. OpenAI, EU AI Act)
- **High stakes**: if credible outlets exit but misinformation stays, AI inherits a distorted information environment

Sample

- **9,611** media/news sites (cross-section, February 2026)
- **7,002** sites × 45 months panel (August 2022–February 2026)
- Source: Media Bias/Fact Check (MBFC)
- robots.txt scraped monthly via HTTP Archive

Quality Measures: Two Complementary Indices

Content quality index (*“inside-out”*)

Expert-rated: what the site *produces*

- MBFC factual reporting (1–6)
- MBFC credibility rating (1–4)
- Wikipedia notable (binary)

→ PCA first component, standardized

Domain quality index (*“outside-in”*)

Algorithmic: how the *web* treats the site

- Trust Flow (Majestic, 0–100)
- Open PageRank (web centrality)
- Tranco rank (aggregated popularity)

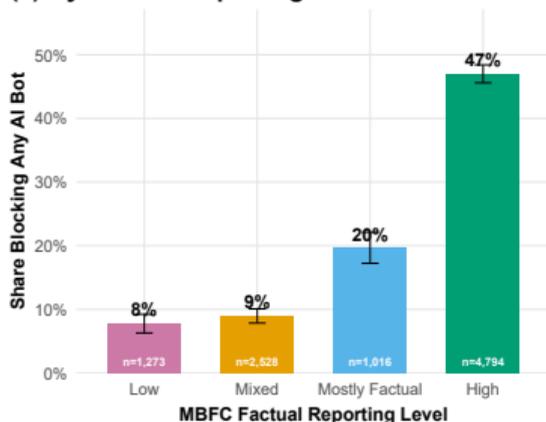
→ PCA first component, standardized

Result 1: Best Content Most Likely to Leave the Training Pool

Higher-quality media sites block AI crawlers at far greater rates

Cross-section of 9,611 MBFC-rated media sites (CCBot excluded)

(a) By Factual Reporting Level



(b) By Credibility Rating



Textbook adverse selection:

- Low-factual sites: **~8%** block; high-factual sites: **~47%** block—a **6×** gap
- Monotonically increasing across both factual reporting and credibility
Consistent across all four quality measures: factual, credibility, content quality index, domain quality index

Result 2: Misinformation Stays Freely Available for AI Training

(a) Credible vs. Questionable

Misinformation sources rarely block AI crawlers

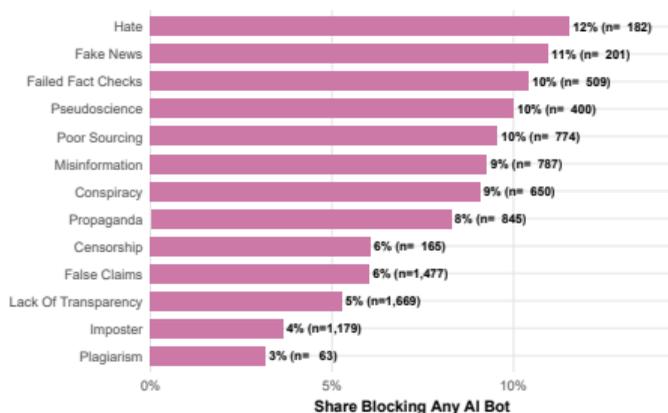
MBFC questionable flags vs. credible sources (N = 9,611, CCBot excluded)



(b) By Misinformation Flag

AI crawler blocking varies across misinformation types

Blocking rate by MBFC questionable source flag (flags with 5+ sites)



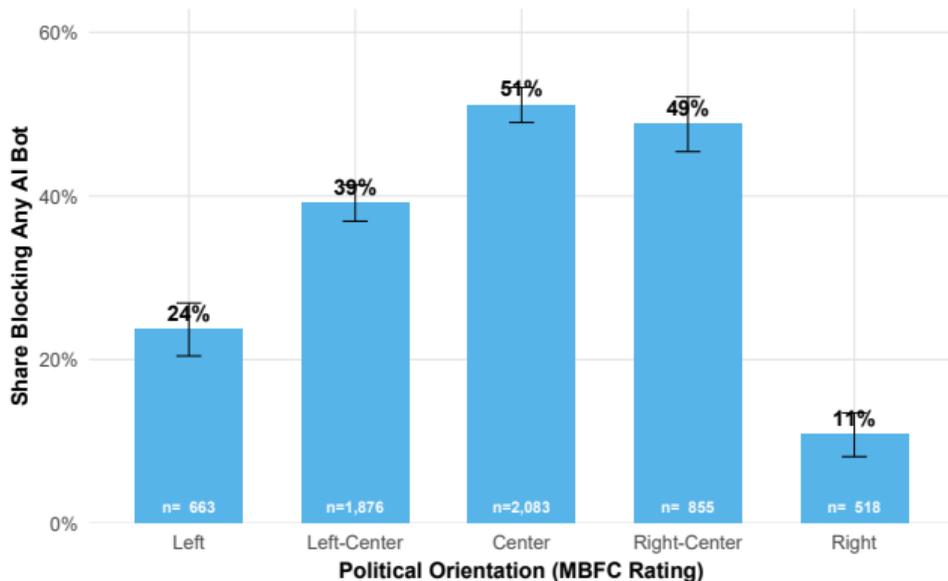
Opt-out skews the corpus toward misinformation:

- Credible outlets block at **5.7x** the rate of questionable sources (35% vs. 6%)
- Conspiracy, pseudoscience, propaganda sources: all below 10%
- As credible sources exit, the remaining corpus tilts toward misinformation

Result 3: AI Training Data Is Becoming Ideologically Sorted

Centrist outlets block most; political extremes remain accessible

MBFC political orientation (N = 5,995, CCBot excluded)



Political extremes stay in the training pool:

- Center: **51%** block
- Left: 24% block
- Right: only **11%** block

⇒ Inverted-U: centrist voices exit, extremes remain

Result 3: Ideological Gaps Persist After Quality Controls

	(1) Bias	(2) +Dom. Q	(3) +Cont. Q	(4) Full
Left	-.354*** (.024)	-.352*** (.024)	-.234*** (.026)	-.230*** (.026)
Left-Ctr	-.173*** (.017)	-.177*** (.017)	-.159*** (.017)	-.163*** (.017)
Right-Ctr	-.101*** (.022)	-.102*** (.022)	-.045** (.022)	-.045** (.022)
Right	-.483*** (.028)	-.467*** (.028)	-.206*** (.035)	-.186*** (.035)
Cont. Q			.234*** (.018)	.236*** (.018)
N	4,650	4,650	4,650	4,650
R ²	.085	.091	.117	.125
Cat. FE				X

Ref: Center. * $p < .1$, ** $p < .05$, *** $p < .01$

Quality explains part, but not all:

- Right: **-19 pp** vs. Center after full controls
- Left: **-23 pp** vs. Center after full controls
- Content quality is a strong predictor (+24 pp per SD)

⇒ Political orientation has an **independent** effect beyond quality

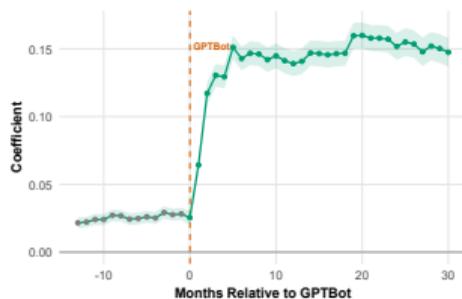
Result 4: Quality–Blocking Gradient Amplified After GPTBot

The quality-blocking gradient emerges sharply after the GPTBot announcement

Event study coefficients, 5,690 media sites, category + month FE, site-clustered SE

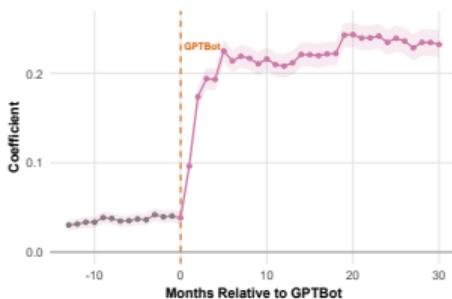
(a) MBFC Factual Reporting

Coeff. on Factual x Month (ref: -14), 95% CI



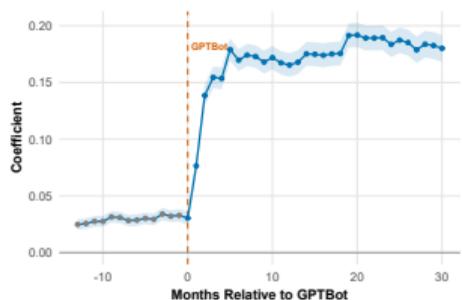
(b) MBFC Credibility

Coeff. on Credibility x Month (ref: -14), 95% CI



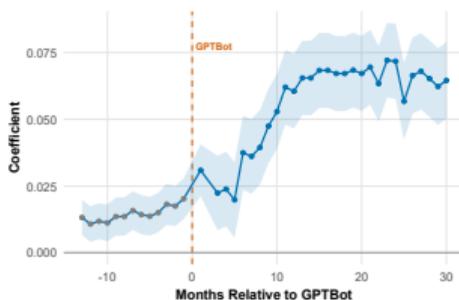
(c) Content Quality Index

Coeff. on Content Quality x Month (ref: -14), 95% CI



(d) Domain Quality Index

Coeff. on Domain Quality x Month (ref: -14), 95% CI



When did adverse selection emerge?

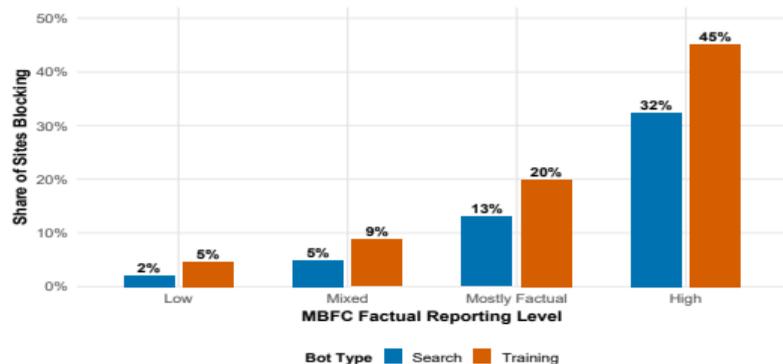
- Pre-period: small, flat quality–blocking gradient
- Aug 2023 (GPTBot): **sharp break**—coefficients jump **5–6×**
- Gradient **keeps widening**: social learning and cascading adoption

⇒ Not a pre-existing pattern—GPTBot made opt-out salient and actionable

Result 5: Adverse Selection Concentrated in Training Data

High-factual media selectively target training bots

Post-GPTBot, MBFC-rated media sites (N = 185,936, CCBot excluded)



Training bot coefficients consistently exceed search bot coefficients

Panel regressions with category + time FE, site-clustered SE, 95% CI (7,002 media sites)



Blocking is strategic, not indiscriminate:

- Quality gap larger for training bots than search bots across factual levels
- Regression coefficients **2–3×** larger for training bots across all quality measures
- Robust across cross-section, panel, site + time FE

⇒ Producers block value *extraction*; permit value *creation*

Result 6: Counterfactual Simulation — Selection, Not Volume, Drives Quality Loss

Adverse selection shifts available content quality leftward

Volume-weighted quality distributions at the same overall blocking rate (5,782 sites)



Both scenarios block at 41% overall. Dashed lines: volume-weighted medians. The leftward shift under current observed blocking is the adverse selection effect.

Decomposing the quality decline:

- **Counterfactual:** same blocking rate: quality-correlated blocking vs random selection
- Adverse selection accounts for the **entire** quality decline, not blocking per se
- Projections through **2028:** degradation compounds to **18–32%** quality decrease

⇒ Quality loss driven by *selection*, not by volume of opt-outs

- ① **Licensing frameworks:** Replace binary opt-out with a market for data licenses
Compensate quality producers to stay in the commons
- ② **Collective action:** Lower transaction costs via industry-wide standards
cf. ASCAP for music — centralized licensing reduces bilateral negotiation costs
- ③ **Regulatory design:** Opt-out regimes must account for adverse selection
EU AI Act, copyright reform

Opt-out regimes by design produce a lemons equilibrium

Three contributions:

- ① **Economics of data markets:** First evidence of adverse selection in AI training data
No price mechanism → best content exits first → lemons problem (Akerlof, 1970)
- ② **AI data governance:** Training data composition is *endogenous*, not given
The commons degrades by design — producers' strategic opt-out decisions shape what remains
- ③ **Misinformation & media:** A new channel for AI to amplify misinformation
Not by generating it, but through the systematic withdrawal of high-quality source material

Three contributions:

- ① **Economics of data markets:** First evidence of adverse selection in AI training data
No price mechanism → best content exits first → lemons problem (Akerlof, 1970)
- ② **AI data governance:** Training data composition is *endogenous*, not given
The commons degrades by design — producers' strategic opt-out decisions shape what remains
- ③ **Misinformation & media:** A new channel for AI to amplify misinformation
Not by generating it, but through the systematic withdrawal of high-quality source material

As the best sources withdraw, AI training data becomes systematically **less factual, less credible, and more ideologically skewed**—without institutional solutions, this degradation becomes **self-reinforcing**.

Thank You!

Questions and Comments?

Kai Zhu

Bocconi University