# Children as Canaries: The Pediatric Sentinel Effect in Algorithmic Harm

**Travis Gilly**
Real Safety AI Foundation
t.gilly@ai-literacy-labs.org

December 2025

Working Paper

## Abstract

Children function as sentinel populations for algorithmic harm in the same way they serve as early indicators of environmental toxicity in public health research. This paper introduces the Pediatric Sentinel Effect in Algorithmic Harm, a framework proposing that harms experienced by children from AI and algorithmic systems reliably precede recognition of those same harms in adult populations. Drawing on analysis of 40 source documents (including UNICEF regional assessments, the UN Convention on the Rights of the Child, and the AI Incidents Database) comprising over 7,300 extracted harm terms, we identify three domains where the sentinel effect is empirically observable: biometric surveillance (school facial recognition deployment preceded adult wrongful arrest cases by eight months), predictive risk systems (child welfare algorithms established precedents for administrative data misuse now spreading to adult contexts), and recommender systems (documented mental health impacts on adolescents preceded clinical recognition of algorithmic addiction in adults). We propose three mechanisms driving this effect: mandatory institutional exposure reducing children's capacity to opt out, cognitive developmental factors increasing susceptibility to algorithmic manipulation, and paradoxical data vulnerability wherein the most legally protected population generates the most behavioral data. The framework has limitations in economic and labor domains where adults experience harms first. We argue that child-focused AI safety testing should be treated as a leading indicator methodology for population-level risk assessment, and that existing pediatric sentinel frameworks from toxicology and environmental health offer immediate methodological templates for AI governance.

## 1. Introduction

Public health researchers have long recognized children as sentinel populations for environmental hazards. The developmental sensitivity of pediatric physiology means children manifest symptoms of toxic exposure at lower thresholds and shorter latencies than adults. Lead poisoning, air pollution

effects, and endocrine disruption were each identified in children before their population-wide impacts were understood. This paper argues that an analogous phenomenon occurs with algorithmic systems: children experience AI-related harms earlier, more severely, and more visibly than adult populations, making them predictive indicators for broader societal risks.

The "canary in the coal mine" metaphor is not merely rhetorical. Safety advocates have begun explicitly applying this framing to youth experiences with AI systems, particularly regarding image-based abuse and mental health deterioration from social platforms. However, no unified academic framework has synthesized the evidence across domains or connected this pattern to established sentinel population theory from public health.

This working paper presents preliminary findings from a systematic analysis of harm vocabularies extracted from 40 child-focused policy documents, cross-referenced with incident-level data from the AI Incidents Database. We propose the Pediatric Sentinel Effect in Algorithmic Harm as a framework for understanding why children's experiences should inform AI governance for all populations.

## 2. The Pediatric Sentinel Framework

Sentinel population methodology in toxicology identifies subgroups whose characteristics make them early indicators of environmental risk. Children qualify due to higher exposure rates relative to body mass, developing organ systems with greater sensitivity, and longer time horizons for chronic effects to manifest.

We propose that analogous characteristics position children as algorithmic sentinels:

**First**, children experience mandatory institutional exposure to algorithmic systems. Unlike adults who retain some capacity to avoid consumer technologies, children are subject to compulsory education systems that increasingly deploy EdTech platforms, surveillance tools, and predictive analytics without meaningful consent mechanisms. Schools function as deployment environments with lower procurement oversight than police departments or healthcare systems.

**Second**, developmental cognitive factors increase susceptibility. The prefrontal cortex, responsible for critical evaluation and impulse regulation, continues developing through adolescence and into early adulthood. AI systems designed for adult cognition (including outputs prone to hallucination, validation loops optimized for engagement, and persuasive interfaces) interact with developing minds in ways their designers did not model.

**Third**, a paradox of data protection creates heightened vulnerability. Frameworks like GDPR classify children's data as requiring special protection, yet gamified educational platforms and social media generate extensive behavioral data from minors. Children are simultaneously the most legally protected and most extensively tracked demographic.

## 3. Evidence Across Three Domains

Our analysis identifies three domains where timeline evidence supports the sentinel effect.

**Biometric surveillance:** The Lockport City School District in New York activated facial recognition systems in May 2019, immediately revealing accuracy failures for children's developing facial structures

and triggering privacy litigation from the New York Civil Liberties Union. The controversy prompted New York State to pass Senate Bill S5140B in December 2020, establishing a moratorium on biometric identifying technology in schools; the lawsuit was subsequently dismissed as moot after the legislation achieved the remedy sought. The first high-profile adult wrongful arrest case attributable to facial recognition (Robert Williams, Detroit) occurred in January 2020, approximately eight months after children in Lockport were already being scanned. Schools' lower oversight thresholds enabled deployment that normalized the technology before police applications faced equivalent scrutiny.

**Predictive risk assessment:** The Allegheny Family Screening Tool began predicting child maltreatment risk in 2016 using administrative data including food stamp utilization and mental health service history. While the COMPAS criminal recidivism algorithm received attention the same year following ProPublica's "Machine Bias" investigation, the child welfare tool established broader precedents for civil-context prediction using poverty indicators rather than behavioral evidence. These methodologies are now migrating to adult domains including hiring, insurance, and housing.

**Recommender systems:** Internal Facebook documents revealed that by 2019 the company possessed research showing Instagram's algorithms exacerbated body image disorders. A March 2020 internal presentation found that 32% of teen girls who felt bad about their bodies reported Instagram made them feel worse; among teens who reported suicidal thoughts, 13% of British users and 6% of American users traced the issue to Instagram. Adult experiences of political polarization and algorithmic radicalization were contemporaneous, but clinical proof of the engagement optimization model's psychological toxicity emerged from adolescent data. The "teen mental health crisis" framing provided the evidentiary foundation later applied to understand adult platform addiction.

## 4. Mechanisms

Three mechanisms explain why harms manifest in children before adults achieve recognition of equivalent harms.

**Captive population dynamics** mean children cannot exit harmful systems. An adult dissatisfied with an employer's AI-driven management can seek other employment. A child assigned to a school deploying surveillance technology has no equivalent option. Mandatory attendance laws create captive user bases for algorithmic experimentation.

**Cognitive mismatch** between AI system design assumptions and child development creates amplified impact. Large language model outputs calibrated against adult reasoning patterns may be dismissed by adult users as unreliable but accepted as authoritative by children lacking developed critical evaluation capacity. The same system produces different harm levels depending on the developmental stage of the user.

**Institutional data hunger** in educational and child welfare contexts drives aggressive data collection. Platforms marketed as educational tools often collect behavioral data at granularities that would trigger regulatory concern in employment or healthcare contexts. The framing of data collection as serving child welfare reduces scrutiny while increasing exposure.

## 5. Scope and Limitations

The pediatric sentinel effect has clear boundaries. In domains driven by economic participation, adults serve as leading indicators because children are not yet participants.

Employment discrimination through algorithmic hiring systems impacts adults first by definition; children do not apply for jobs. The EEOC's settlement with iTutorGroup in August 2023 over automated age-based rejection of applicants, and the ongoing Mobley v. Workday litigation (proceeding as a collective action as of 2024), illustrate harms impossible to observe in pediatric populations.

Labor displacement from generative AI systems first affected adult professional illustrators and writers whose work was scraped for training data. While children's creative output is also appropriated, the economic harm registered immediately in adult professional communities.

These limitations do not invalidate the sentinel framework but define its domain of applicability: surveillance, predictive systems, and engagement-optimized platforms deployed in institutional contexts where children lack exit options.

## 6. Implications for AI Governance

If children function as sentinel populations for algorithmic harm, several governance implications follow.

**Child-focused AI safety testing should be treated as leading indicator methodology.** Harms detected in school deployments, child welfare systems, and youth-oriented platforms predict harms that will later appear in adult contexts. Regulatory attention should flow toward, not away from, pediatric applications.

**Existing pediatric sentinel methodology from environmental health offers immediate templates.** The infrastructure for monitoring children as indicators of toxic exposure could be adapted for algorithmic harm surveillance with appropriate translation of metrics.

**Mandatory institutional exposure creates regulatory leverage.** Schools and child welfare agencies operate under public oversight in ways that private adult contexts do not. Procurement requirements for algorithmic systems in child-serving institutions could establish standards that propagate to other domains.

## 7. Conclusion

The Pediatric Sentinel Effect in Algorithmic Harm is not a metaphor but a methodological framework grounded in established public health practice. Children's experiences with AI systems are not merely concerning in their own right but predictive of population-level risks. The evidence across biometric surveillance, predictive risk assessment, and recommender systems demonstrates consistent patterns of earlier harm manifestation in pediatric populations. AI governance frameworks that treat child safety as a specialized concern rather than a leading indicator are methodologically incomplete.

## References

Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016, May 23). Machine Bias. *ProPublica.* https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

Equal Employment Opportunity Commission. (2023, August 9). iTutorGroup to Pay $365,000 to Settle EEOC Discriminatory Hiring Suit [Press release]. https://www.eeoc.gov/newsroom/itutorgroup-pay-365000-settle-eeoc-discriminatory-hiring-suit

Eubanks, V. (2018). *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor.* St. Martin's Press.

5Rights Foundation. (2021). *But How Do They Know It Is a Child? Age Assurance in the Digital World.* https://5rightsfoundation.com/

Hill, K. (2020, June 24). Wrongfully Accused by an Algorithm. *New York Times.* https://www.nytimes.com/2020/06/24/technology/facial-recognition-arrest.html

Mobley v. Workday, Inc., No. 23-cv-00770-RFL (N.D. Cal. July 12, 2024) (order denying motion to dismiss).

New York State Senate. (2020). Senate Bill S5140B: An act to amend the education law, in relation to the use of biometric identifying technology. Signed December 22, 2020, Chapter 349. https://www.nysenate.gov/legislation/bills/2019/s5140

UNICEF Office of Research - Innocenti. (2020). *Policy Guidance on AI for Children.* https://www.unicef.org/globalinsight/reports/policy-guidance-ai-children

Weiss, E. (2022, July 6). NYCLU Finds New York Schools Not Sticking to Facial Recognition Ban. *ID Tech Wire.* https://idtechwire.com/nyclu-finds-new-york-schools-not-sticking-facial-recognition-ban-070602/

Wells, G., Horwitz, J., & Seetharaman, D. (2021, September 14). Facebook Knows Instagram Is Toxic for Teen Girls, Company Documents Show. *Wall Street Journal.* https://www.wsj.com/articles/facebook-knows-instagram-is-toxic-for-teen-girls-company-documents-show-11631620739