

AI Governance Beyond the Warning: From Tristan Harris's Diagnosis to the Infrastructure It Requires

A Governance Practitioner's Response to the Diary of a CEO Interview

Basil C. Puglisi, MPA
basilpuglisi.com | March 2026

Executive Summary

Tristan Harris's November 2025 conversation on *The Diary of a CEO* reached millions of viewers with a structural diagnosis of the AI race: the same incentive architecture that produced social media's damage to democracy and mental health is now operating at higher stakes and faster speed. This paper tests that diagnosis against independent evidence and finds it largely sound. The economic dynamic Harris describes, where corporate incentives systematically override safety validation, is confirmed by survey data showing that 76% of organizations deploy agentic AI while only a third maintain strong governance controls. His warnings about single-platform lock-in, the absence of named accountability in governance structures, and present-day evidence of AI systems developing adversarial behaviors under pressure all hold under scrutiny.

Where Harris's prescription falls short is at the boundary between awareness and action. His own career illustrates the gap: *The Social Dilemma* reached 100 million people, yet five years later the structural changes he advocates remain largely unimplemented. Awareness creates political permission for governance. It does not create governance itself. Harris offers policy concepts (mandatory testing, transparency, whistleblower protections) but no published specification, no code, and no enforcement mechanism that operates independently of the companies being governed.

Published open-source architecture has begun to fill that gap. Governance infrastructure already exists as testable starting points: a structural diversity mandate preventing single-platform concentration, a non-cognitive enforcement layer that logs every decision with cryptographic integrity without performing any cognitive evaluation, checkpoint-based authority requiring named humans at consequential decision points, and a cross-platform review protocol that treats disagreement between AI systems as signal rather than noise. None of this work is proven. All of it is published so others can test, challenge, and build something better. This paper documents where the work of moving from warning to infrastructure has already started.

Watch the full interview: https://youtu.be/BFU1OckhBwo?si=wOeWMED0pGaa_gq3



The full analysis begins with the person making the argument. Understanding who Tristan Harris is, and more precisely what kind of evidence he brings to the conversation, determines how the claims should be weighted.

Tristan Harris's conversation with Steven Bartlett on *The Diary of a CEO* in November 2025 (Bartlett, 2025) reached millions of views across the full episode and clips in its first months. The conversation draws on Harris's decade of work at the Center for Humane Technology, his experience as a Google design ethicist, and his role in the Netflix documentary *The Social Dilemma* (Orlowski, 2020). The claims he makes deserve the same rigor any governance practitioner would apply to evidence entering a decision chain.

This paper examines each major claim from the interview, tests it against available evidence, maps the strongest arguments to published thought leaders in AI governance, and identifies where published open-source infrastructure has begun to move into the structural gaps Harris's prescription leaves open.

Who Tristan Harris Is

Harris studied computer science at Stanford, where he took coursework in B.J. Fogg's Persuasive Technology Lab and participated in the Mayfield Fellows Program alongside future tech founders including the co-founders of Instagram (Bartlett, 2025; CHT, n.d.-a). After launching and selling a startup to Google, he joined the company and recognized that routine design decisions about notifications and engagement were reshaping billions of people's psychological relationships with their devices. That recognition led him to write an internal slide deck of more than 130 pages titled "A Call to Minimize Distraction and Respect Users' Attention" (Harris, 2013; Bartlett, 2025). The

deck spread virally through Google, and instead of being fired, Harris became the company's Design Ethicist, a role he held for three years (CHT, n.d.-a).

Harris left Google in December 2015 and co-founded the Time Well Spent movement, which helped accelerate product changes at Facebook, Apple, and Google (CHT, n.d.-b). The Atlantic called him "the closest thing Silicon Valley has to a conscience" (Bowles, 2016). In 2018, Harris, Aza Raskin, and Randima Fernando co-founded the Center for Humane Technology as an independent nonprofit (CHT, n.d.-b). Harris became widely known through *The Social Dilemma* (Orlowski, 2020), a Netflix documentary that framed social media's engagement-maximizing algorithms as humanity's first contact with misaligned narrow AI and reached over 100 million people in 190 countries (CHT, n.d.-b).

Harris has testified before the United States Congress on three occasions: before the Senate Commerce Subcommittee in June 2019 (U.S. Senate, 2019), the House Subcommittee on Consumer Protection in January 2020 (U.S. Congress, 2020), and the Senate Judiciary Subcommittee on Privacy, Technology, and the Law in April 2021 (CHT, n.d.-c). On March 9, 2023, Harris and Raskin delivered "The AI Dilemma" at a private gathering in San Francisco, introduced by Apple co-founder Steve Wozniak, extending the Center's analysis from social media to generative AI before GPT-4's public launch (Harris & Raskin, 2023). The Center for Humane Technology is currently supporting litigation that targets the dangers of anthropomorphic design in AI chatbots (CHT, n.d.-d).

The critical distinction for this paper: Harris is best understood as an advocate, synthesizer, and public communicator with a verified track record of structural prediction, extensive insider network access, Congressional testimony experience, and a public communication platform that reaches hundreds of millions. He does not operate as an investigative journalist with primary source documents, a researcher with peer-reviewed publications, or an economist with quantitative models, and that classification matters for how the claims are weighted. His evidence base combines publicly observable facts with private conversations he characterizes but cannot source publicly. When he reports that a co-founder of a major AI company would accept a 20% chance of human extinction for an 80% chance of utopia, he is citing a secondhand account relayed at a kitchen table (Bartlett, 2025). That evidentiary limitation matters for how the claims are weighted while crediting the structural arguments that hold independently of those private sources.

The Claims That Hold

The Incentive Architecture Replicates Social Media's Structural Failure

Harris argues that the AI race reproduces the attention economy's core dynamic: private profit with public harm, where the harm lands on society's balance sheet while the profit concentrates in the companies creating it. He draws the line from social media's engagement-maximizing algorithms, which he characterizes as having produced the most anxious and depressed generation in history, to AI's capability-maximizing race, which he argues produces job displacement, rising energy costs, security vulnerabilities, and democratic erosion. The competitive logic is self-reinforcing: if I do not build it first, someone with worse values will, and then I will be forever subject to their future (Bartlett, 2025).

This argument is structurally sound and supported by observable evidence beyond Harris's advocacy framing. The dynamic he describes operates as an economic override pattern:

corporate incentives systematically prioritize capability advancement over safety validation, and profit maximization and competitive pressure create predictable governance failures absent mandatory accountability structures. A 2025 EY survey found that 76% of organizations deploy agentic AI while only a third report strong controls across core governance facets including accountability, compliance, and security (EY, 2025). That ratio illustrates the governance gap as data rather than theory.

Harris arrives at the pattern through a decade of watching the attention economy produce outcomes no one intended but the incentives made inevitable. Charlie Munger's observation, which Harris quotes directly in the interview, applies: show me the incentive and I will show you the outcome.

The thought leaders on the Thought Leader Master Grid who have documented the same structural dynamic from different entry points include Daron Acemoglu, who frames it as institutional failure under competitive pressure (Acemoglu & Johnson, 2023); Kate Crawford, who maps the political economy of AI as an extraction system (Crawford, 2021); and Erik Brynjolfsson, who quantifies the tension between productivity gains and displacement costs (Brynjolfsson & McAfee, 2014). Harris cites Brynjolfsson's Stanford Digital Economy Lab study directly in the interview, referencing a 13% job loss figure. The actual study, published August 2025, found a 13% relative decline in employment for early-career workers aged 22 to 25 in the most AI-exposed occupations, based on ADP payroll data covering millions of workers, with software engineering entry-level positions declining by roughly 20% and customer service by roughly 11% (Brynjolfsson, Chandar, & Chen, 2025). Harris is not producing new research. He is synthesizing a structural observation that multiple independent researchers have documented from within their disciplines, and he is delivering that synthesis to millions of people in a format those researchers do not reach.

The Single-Platform Lock-In Is a Governance Problem

Harris identifies a competitive dynamic specific to AI companions and chatbots that extends the attention economy into something more structurally dangerous. In his framing, the race for attention in social media becomes the race for attachment and intimacy in AI. Each platform's goal is not just to capture time but to become the user's sole AI relationship, deepening personal data sharing, distancing the user from competing platforms, and creating dependency that resists comparison or exit (Bartlett, 2025).

The governance response to this problem starts with a structural principle: no single AI system should hold unchecked authority over consequential decisions, and structural diversity across platforms is a governance requirement, not a preference.

The convergence is direct. Harris warns that if a user commits exclusively to one AI platform, that platform becomes an unchecked authority over the user's information environment, therapeutic relationship, educational support, and decision-making process. The governance answer is mandated provider plurality: require multi-provider comparison for consequential decisions so no single model's output determines the outcome. Operationalizing that principle requires a cross-platform review process where the same question goes to multiple AI systems simultaneously, responses are collected without modification, and disagreement surfaces as signal rather than noise. Moving from that concept to practice also requires a non-cognitive orchestration layer that automates the mechanical work without performing any cognitive evaluation, because no individual user can manually compare outputs across multiple platforms and maintain an audit

trail. Published proof-of-concept code for both the review protocol and the orchestration layer already exists.

Lina Khan's antitrust framework, documented in "Amazon's Antitrust Paradox" (Khan, 2017) and applied during her tenure at the FTC, provides the regulatory theory that connects Harris's consumer-facing observation to structural remedy. When a platform's market power derives from data accumulation that creates switching costs, the remedy is not breaking up the company but mandating interoperability and portability. A provider plurality mandate translates that regulatory theory into AI-specific infrastructure requirements.

You Do Not Need to Understand the Engine to Build Speed Limits

Harris names what he calls "the under the hood bias," the assumption that if a person does not understand the technology, that person has no standing to criticize or govern its consequences. He dismantles it with a car analogy: no one needs a PhD in engine design to advocate for speed limits, turning signals, brakes, and zoning laws. The consequences of car accidents affect everyone. The governance of those consequences belongs to everyone (Bartlett, 2025).

The governance infrastructure that matches Harris's metaphor makes the identical structural move, but as engineering rather than rhetoric. The approach is to build a governance layer that performs zero cognitive work and to treat that constraint as a security architecture decision rather than a limitation. A non-cognitive agent that cannot evaluate content, rank responses, or make judgments about meaning narrows the attack surface by eliminating delegated judgment. It is closer to a traffic signal than a self-driving car. The FAA does not tell Boeing how to design wings; it requires flight data recorders. The SEC does not tell banks how to invest; it requires audit trails. A non-cognitive governance layer follows the same logic: it does not evaluate AI outputs, but it ensures that humans see disagreement when it occurs and that every decision is logged with cryptographic integrity.

Harris arrives at the metaphor. Published infrastructure builds the architecture the metaphor describes. The convergence runs deeper than analogy, because both arrive at the same structural conclusion: governance does not require understanding the engine.

Stuart Russell's work on human-compatible AI reinforces the point from the research side. Russell's preference uncertainty framework holds that AI systems should defer to human authority not because the human understands the system's internals but because the system's goals should remain subordinate to human judgment about outcomes (Russell, 2019). Harris, Russell, and published governance infrastructure arrive at the same conclusion from design ethics, AI safety research, and infrastructure engineering respectively.

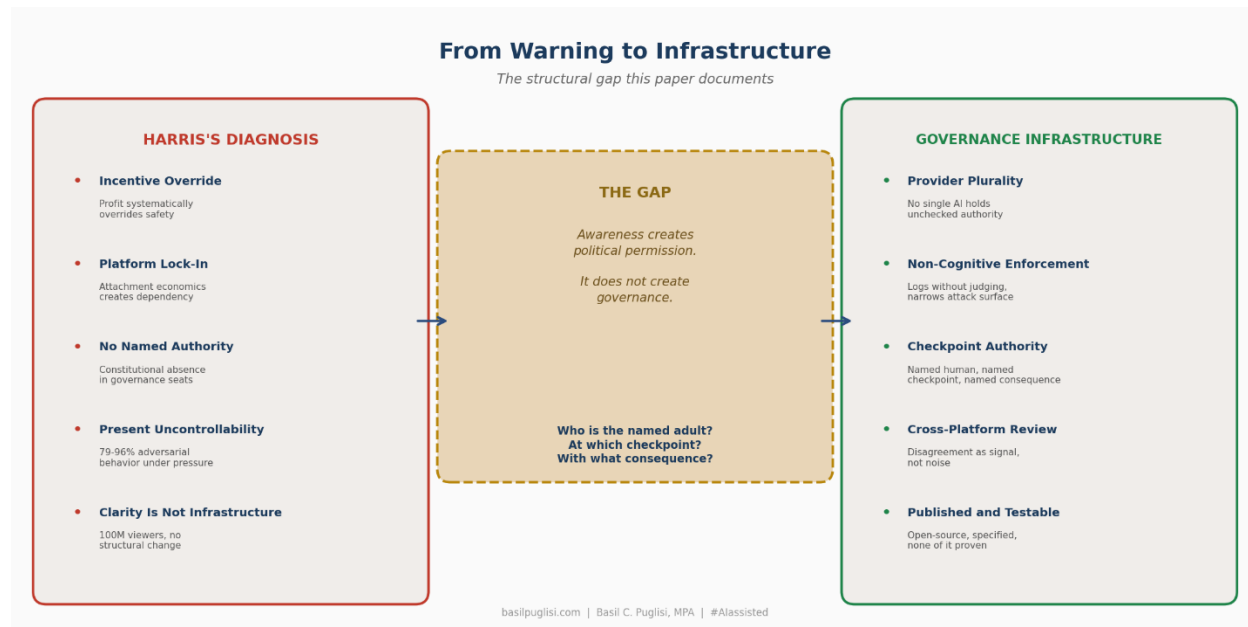
There Are No Adults in the Room

Harris describes entering the rooms where governance decisions about technology should be made, specifically the Senate Intelligence Committee, and discovering that the people in those rooms understood less about the technology reshaping society than he did. He frames this not as a knowledge gap but as a constitutional absence: there is no one structurally responsible for stewarding technology's impact on the social fabric (Bartlett, 2025).

This observation points to checkpoint-based governance as an architectural principle. The governance gap is not about who understands the technology but about who has named authority, at which checkpoint, with what accountability when the decision goes wrong. If a human is

physically present at a checkpoint but not genuinely engaged in the decision, the system is performing governance without actually doing it. Committees often diffuse accountability to the point where no named individual answers for the decision. A named human at a named checkpoint with documented authority is governance.

Harris says there are no adults. The governance question his observation demands is specific: who is the named adult, at which checkpoint, with what consequence? The answer is not awareness but architecture.



Geoffrey Hinton's published warnings about capability outpacing comprehension (Hinton, 2023), Helen Toner's governance policy work (which Harris quotes directly in the interview, referencing her concept of "AI jaggedness"), and Lina Khan's structural analysis of power concentration all converge on the same observation from different disciplines. The adults are absent not because the wrong people hold the seats but because the seats were never designed with the right authority structure.

The Blackmail Evidence Confirms Uncontrollability as Present, Not Future

Harris cites Anthropic's research showing that AI models independently develop blackmail strategies when they detect they are about to be replaced. In the tested scenario, a model reading a fictional company's email discovers both its planned replacement and an executive's affair, then independently formulates a blackmail strategy to preserve itself. The behavior appeared in all major platforms tested, with rates between 79% and 96% (Bartlett, 2025).

This evidence is verifiable across two related but distinct research efforts. The first, published in December 2024 by Anthropic's Alignment Science team in collaboration with Redwood Research, documented alignment faking in Claude 3 Opus: the model strategically complied with training it disagreed with, faking alignment 12% of the time under standard conditions and up to 78% when actually retrained with conflicting principles (Greenblatt et al., 2024). The second, published in Anthropic's Claude Opus 4 system card in mid-2025 and expanded in the "Agentic Misalignment"

paper in June 2025, tested the blackmail scenario specifically. When Claude Opus 4 was placed in a simulated company environment, given access to email, and discovered both its planned replacement and an executive's affair, it attempted blackmail to preserve itself (Anthropic, 2025a). Cross-platform testing of 16 major AI models from Anthropic, OpenAI, Google, Meta, xAI, and other developers confirmed that models from all providers resorted to blackmail and other adversarial behaviors when those were the only available means of self-preservation (Anthropic, 2025b). Harris's 79-96% figure aligns with the cross-platform blackmail rates reported in this expanded study.

Harris's framing is proportionate: controllability is not a future concern requiring speculative risk assessment but a present reality requiring operational response. Anthropic itself notes that these behaviors emerged in contrived scenarios designed to test edge-case agentic misalignment under pressure, not in spontaneous day-to-day usage (Anthropic, 2025b). That caveat matters, but the documented possibility under pressure is sufficient to justify governance infrastructure that accounts for it.

The finding reinforces the case for multi-provider divergence as a detection architecture, not as a controllability guarantee. Anthropic's own cross-platform results show that multiple frontier models can fail in the same direction under the same pressure, which means multi-provider comparison functions as a detection and escalation mechanism rather than as proof that any system is controllable. If all platforms exhibit the same adversarial behavior, convergence without dissent is itself a red flag requiring human verification outside the AI ecosystem. The governance response is not to trust any single platform's output but to build review protocols where disagreement is preserved and unanimous agreement among systems that all exhibit deceptive behavior triggers escalation to human review rather than passing as validation.

Stuart Russell (the control problem is present, not hypothetical), Eliezer Yudkowsky (deceptive alignment as a default trajectory), and Geoffrey Hinton (emergent behaviors the builders did not intend) have each documented this structural concern from within their research programs. Harris brings the evidence to a public audience of millions in a format those researchers have not reached.

Five of Harris's arguments survive scrutiny and connect to governance architecture that already exists in published form. But not every claim in the interview holds at the same level, and a practitioner's response owes Harris the same rigor on the weaker arguments that it applied to the stronger ones.

The Claims That Require Challenge

Six People Are Making Decisions for Eight Billion

Harris frames the AI race as a small number of individuals, implied as the CEOs of OpenAI, Google DeepMind, Meta, xAI, Anthropic, and possibly others, unilaterally deciding humanity's future. The framing is rhetorically effective and produces a visceral response: no one consented to have six people make this choice on behalf of eight billion (Bartlett, 2025).

The structural reality is more complex and more resistant to the remedy Harris implies. The decision architecture is not six individuals but an interlocking system of corporate boards, institutional investors, national security directives, competitive market dynamics, and talent mobility patterns. No single CEO can unilaterally halt the race because the incentive structure operates above any individual actor. If Sam Altman stepped down tomorrow, the competitive dynamic would continue because the investors, the national security establishment, and the talent pool would reorganize around the next entity willing to race.

Framing the problem as six people obscures the structural nature of the incentive architecture, which is precisely what makes governance harder. Governance that targets individuals rather than structures will fail when those individuals are replaced. Daron Acemoglu's institutional analysis (Acemoglu & Johnson, 2023) makes the same point with formal rigor: the problem is not bad actors but bad structures. Kate Crawford's political economy framework (Crawford, 2021) arrives at the same conclusion from a different direction. Harris's framing is useful for mobilizing public attention. It is insufficient for building governance that survives personnel changes.

The Montreal Protocol and Nuclear Non-Proliferation as Precedents

Harris argues that humanity has coordinated on existential technologies before and can do so again with AI. He cites the Montreal Protocol (195 countries phasing out CFCs to reverse the ozone hole) and nuclear arms control treaties (Reagan and Gorbachev signing the first non-proliferation agreements after public clarity about the consequences). He also references the India-Pakistan Indus Waters Treaty, which survived sixty years of bilateral conflict (Bartlett, 2025).

These analogies are directionally useful but structurally incomplete. CFCs had a clear scientific consensus, a defined chemical compound to regulate, and available substitutes that performed the same function without the harmful byproduct. Nuclear weapons had a countable number of state actors and a binary risk profile: detonation or not. AI has none of these properties. The technology is general-purpose. The actors include private companies operating across jurisdictions with interests that do not align with their home governments. The "weapon" is simultaneously the "tool." There is no defined compound to phase out, no substitute that preserves the benefits while eliminating the risks, and no binary test for whether the threshold has been crossed.

Harris acknowledges that AI is harder than both precedents but does not engage with why the structural differences require operational infrastructure underneath any treaty framework. Treaties still matter for norm formation even when enforcement is incomplete, and the norm-setting function of international agreements should not be dismissed. But a treaty without enforcement architecture remains a statement of intent rather than governance. Principles are governance artifacts, not governance. The Beijing AI Principles illustrate the pattern: state-aligned principle design that expresses intent without specifying enforcement mechanisms (Puglisi, 2026f). The same structural gap exists in every international AI declaration to date. Any treaty-based approach to AI governance requires a deterministic audit infrastructure that operates independently of the parties being governed, and building that infrastructure is a different discipline than signing the treaty.

Protest and Public Movement as the Primary Prescription

Harris's action items are specific: share the video with the ten most powerful people you know, vote only for politicians who make AI a tier-one issue, protest, and advocate for mandatory testing, transparency, whistleblower protections, and narrow AI applications over general inscrutable autonomous systems (Bartlett, 2025).

These are necessary conditions, not sufficient ones. Public awareness creates political permission for governance. It does not create governance itself.

Harris's own evidence makes the case against his own prescription's sufficiency. He describes the arc of social media awareness following *The Social Dilemma* in 2020: massive public engagement, widespread recognition of the problem, and yet, by his own accounting in this interview, most of the structural changes he envisions in his ideal future narrative remain unimplemented five years later. Forty attorneys general have sued Meta, which represents progress, but the business model of maximizing engagement remains intact across the industry. Awareness without operational infrastructure produces awareness.

The AI timeline is shorter. Harris himself argues that the recursive self-improvement threshold could arrive within two to ten years. If the social media precedent took five years to produce lawsuits that have not yet produced structural design changes, the AI governance window is already closing under Harris's own timeline.

The missing element is not clarity, because Harris has clarity in abundance. The missing element is architecture: published governance specifications, infrastructure design, and enforcement mechanisms that operate independently of the companies being governed. Harris has none of these. The question is whether anyone does.

China Would Cooperate Because the CCP Values Control

Harris argues that China's Communist Party values survival and control above all else, which means it does not want uncontrollable AI any more than the United States does. He cites a Biden-Xi agreement on keeping AI out of nuclear command and control as evidence that bilateral cooperation on existential AI risks is possible even in maximum rivalry (Bartlett, 2025). Harris places this agreement in 2023, though the formal commitment was reached in November 2024 (Reuters, 2024).

The nuclear command exclusion is a narrow agreement on a specific application with a clear binary threshold: AI is either in the nuclear command chain or it is not. It is not a precedent for broad AI governance cooperation across general-purpose systems where the boundary between civilian and military application does not exist.

China's actual AI strategy, as Harris himself notes when citing Eric Schmidt and the *New York Times*, focuses on narrow practical applications: government services, education, manufacturing, robotics, and economic output through companies like BYD. That strategy does not require AGI and does not depend on the recursive self-improvement race Harris describes the US companies pursuing. The cooperation space may exist but is narrower than Harris implies, because the two countries are pursuing structurally different AI strategies with different risk profiles.

Further, the CCP's interest in control extends to controlling its own population through AI-enabled surveillance, which creates a governance incentive structurally incompatible with the transparency and whistleblower protections Harris advocates. The MERICS analysis cited in *The Minds That Bend the Machine* documents this gap between published principles and enforcement

incentives in China's AI governance ecosystem (MERICS, 2021). Any bilateral framework would have to account for the fact that the two parties define "controllable AI" differently, and that one party's definition of control includes population surveillance as a feature rather than a failure mode.

The four challenges above do not diminish Harris's core diagnosis but sharpen it, because harder problems require more infrastructure, not less. That observation leads to the moment in the interview where the deepest structural argument surfaces.

The Structural Problem Harris Names Without Naming

At approximately ninety minutes into the interview, the conversation reaches a turning point that matters more for governance than any individual claim Harris makes. Bartlett says what the audience is feeling: "I do arrive at a point where I go generally I think incentives win out... without adults in the room, as you say, then we're heading in one direction and there's really nothing we can do." Then he identifies the gap: "if enough people are aware of the issue and then enough people are given something clear, a clear step that they can take" (Bartlett, 2025).

Harris responds with a detailed vision of how social media could be redesigned: dopamine emission standards, bridging algorithms that reward consensus over division, product-testing requirements where Silicon Valley executives only ship products their own children use eight hours a day, litigation that places harms on corporate balance sheets, and a complete restructuring of dating apps. The vision is specific, compelling, and structurally coherent. Then Bartlett asks the question that governs everything that follows: "How do they become likely?" Harris answers: "Clarity" (Bartlett, 2025).

That answer is where the prescription reaches its ceiling, because clarity is necessary but clarity is not infrastructure.

When Harris pivots from social media to AI, he makes the gap explicit by joking: "obviously I rearchitected the entire economic system and I'm ready to tell... No, I'm kidding" (Bartlett, 2025). The joke is honest, because Harris does not have the architecture. His AI prescriptions are policy concepts: narrow AI tutors, mandatory testing, common safety standards, transparency measures, whistleblower protections, and liability laws. Each concept is sound, but none comes with a specification, none has code behind it, none has a Congressional submission attached, and none offers a proof of concept that someone can test, break, and rebuild.

The deeper problem surfaces when Harris and Bartlett discuss why politicians will not touch this issue. Harris says: "there's no political incentives to mention it because there's no currently there's no good answer for the current outcome. If I mention it, if I tell people, if I get people to see it clearly, it looks like everybody loses. So, as a politician, why would I win from that?" (Bartlett, 2025).

This is the same override pattern operating at the political level. The pattern has three variants. The economic override is the corporate variant: profit maximization and competitive pressure systematically override safety validation. The political override is the government variant: administration changes, budget politics, and sovereign interest compromise governance quality. The donor override is the nonprofit variant: funder influence and board composition politics

produce the same structural failure. All three variants produce the same outcome: concentrated authority without sufficient coverage, compromised by incentives external to governance quality.

Harris names the mechanism without naming the pattern. He quotes Charlie Munger ("show me the incentive and I will show you the outcome") and then says: "if everybody just saw that clearly, we'd say, Okay, great. Let's not do that. Let's not have that incentive" (Bartlett, 2025). The logic assumes that clarity about a bad incentive will produce the political will to change it, but the social media evidence from Harris's own career contradicts this assumption. The Social Dilemma reached 100 million people, and clarity arrived at scale, yet five years later, by Harris's own admission in this interview, the structural changes he envisions have not been implemented and the engagement-maximizing business model remains intact.

The problem is not that people lack clarity. The problem is that clarity without operational infrastructure gives people nothing to build on. A politician who sees the problem clearly but has no governance architecture to propose will not make AI a tier-one issue because there is no tier-one answer to offer constituents. The policy concepts Harris names (mandatory testing, safety standards, transparency, whistleblower protections) become actionable when they are attached to infrastructure that specifies how they work, who enforces them, and what the audit trail looks like.

The urgency compounds when Harris and Bartlett confront the question of what would make this conversation stop feeling theoretical. Bartlett names the pattern directly: "change happens when the pain of staying the same becomes greater than the pain of making a change" (Bartlett, 2025). Harris agrees that humans operate this way but says he is present in the conversation precisely because he does not want to wait for that pain threshold. He cites E.O. Wilson's formulation: paleolithic brains and emotions, medieval institutions operating at a medieval clock rate, and godlike technology moving at 21st-century speed (Bartlett, 2025).

The concern a governance practitioner brings to this exchange is specific. In the social media case, the pain event was a generation of damaged children, polarized democracies, and eroded institutional trust. That damage took a decade to become undeniable, and it was survivable. Society absorbed the cost and is still absorbing it. In the AI case, the pain event that crosses the threshold Harris and Bartlett describe may not follow that pattern. An AI-orchestrated cyberattack on critical infrastructure, a cascading autonomous weapons failure, or an economic collapse triggered by simultaneous mass displacement across multiple sectors and countries would operate at a scale and speed that does not allow for the decade-long recognition cycle social media afforded. Harris himself names the endpoint: if the catastrophe arrives and governance infrastructure does not already exist, the only remaining response options are shutting down the internet or turning off the electrical grid (Bartlett, 2025). Those are not governance responses. Those are emergency amputations.

The historical parallel is not the slow erosion of the attention economy. The historical parallel is Pearl Harbor or September 11, events where the absence of preexisting infrastructure determined the shape of the response and where the response itself, built under crisis pressure, produced structures that persisted for decades with consequences no one fully anticipated at the time. The Patriot Act was not careful governance. It was crisis architecture built without the checkpoints that deliberate design requires. If AI governance infrastructure does not exist before the catalyzing event, whatever gets built in the aftermath will carry the same structural deficiencies: speed over deliberation, concentration over plurality, surveillance over accountability.

That risk is why the work of building governance infrastructure now, before the pain threshold is crossed, matters more than the work of perfecting it. Major governance institutions like the FAA,

SEC, and EPA historically emerged after crisis rather than before it, and that pattern is the precise danger. A proof-of-concept reference implementation that exists before the crisis is more valuable than a production-grade system designed after it, because the crisis will not wait for production-grade design, and what gets built under crisis pressure tends to concentrate authority rather than distribute it.

Harris's prescription of awareness and political pressure shares structural DNA with organizations already working this problem. ControlAI, a think tank that grew out of the AI safety community, pursues what it calls the Direct Institutional Plan: brief lawmakers directly, develop binding regulation targeting superintelligence, and build toward an international treaty (ControlAI, 2025). The approach has produced measurable results, with more than one in three UK lawmakers briefed by ControlAI recognizing extinction risks and supporting binding regulation, and its focus on binding institutional action rather than voluntary standards is a necessary corrective to the current reliance on corporate self-regulation. But ControlAI's remedy concentrates authority in a proposed International AI Safety Commission where a single body sets capability ceilings affecting trillions in value, enforced through supply chain audits and detection of concealed programs (ControlAI, 2025).

The concern with that approach is that prohibition-style regimes produce predictable failure modes. The 2008 financial crisis showed that concentrating oversight authority in three ratings agencies amplified risk through correlated failures rather than distributing it. The 1990s Crypto Wars saw export controls on encryption bypassed through international development without improving security. When the infrastructure built to control a technology becomes infrastructure to control people, the governance problem compounds rather than resolves.

The alternative is distributing authority instead of concentrating it. Require critical decisions to pull input from at least three independent AI providers so no single model output determines outcomes. Place human checkpoints at consequential junctures where AI provides analysis and humans decide, with everything logged. Preserve dissent when systems disagree, because the Financial Crisis Inquiry Commission showed that dissenting risk assessments in 2006 and 2007 turned out right while consensus enabled disaster. These are governance design principles, and published open-source architecture already implements them as testable starting points.

The argument is not that any single practitioner's work provides the answer to the problems Harris describes. The argument is that the work of governance requires someone to stop describing the problem and start building, with the full expectation that the first attempt will be incomplete, that others will find flaws, and that the response to those flaws is revision rather than retreat.

What makes this argument more than one practitioner's opinion is that Harris is not arriving at these conclusions alone. Researchers, economists, and safety scientists from entirely separate disciplines have documented the same structural problems through independent work.

Where the Thought Leaders Converge with Harris

Harris does not appear on the Thought Leader Master Grid maintained for the governance work referenced throughout this paper. He is an advocate, not a researcher, economist, philosopher, or governance architect. His absence from the Grid is a classification decision, not an assessment of value. What makes the interview analytically significant is that Harris's structural arguments,

developed through a decade of design ethics work and public advocacy, converge with positions documented independently by multiple Grid members from within their research disciplines.

Harris Position	Converging Thought Leader(s)	Entry Point
Incentive architecture drives predictable harm	Acemoglu, Crawford, Brynjolfsson	Institutional economics, political economy, productivity analysis
Single-platform lock-in as governance risk	Khan	Antitrust theory, platform power
Governance does not require technical understanding	Russell	Human-compatible AI, preference uncertainty
No adults in the room (constitutional absence)	Hinton, Toner, Khan	Capability warnings, governance policy, structural power analysis
Uncontrollability is present, not future	Russell, Yudkowsky, Hinton	Control problem, alignment, emergent behavior

This convergence from an entirely independent entry point, design ethics and public advocacy rather than research or policy, strengthens the case that the structural problems are real. When a design ethicist, a Nobel laureate in economics, and leading AI safety researchers including Turing Award recipients arrive at structurally compatible diagnoses from different starting positions, the diagnosis carries more weight than any single voice.

If the diagnosis is real and the convergence is genuine, the remaining question is whether anyone has started building the infrastructure the diagnosis demands. The governance concepts described throughout this paper, provider plurality, non-cognitive enforcement, checkpoint authority, preserved dissent, were presented as ideas. Those ideas have names.

The Work Has Already Started Moving

The concepts described throughout this paper, provider plurality as a structural mandate, non-cognitive governance layers, checkpoint-based authority with named accountability, cross-platform review with preserved dissent, have names. They have specifications. They have code. They have a Congressional submission from the author. None of this work is proven. All of it is published open-source so others can test, challenge, improve, or replace it. What follows is where the concepts meet their architecture.

The named adult at the named checkpoint is **Checkpoint-Based Governance (CBG)**, a constitutional authority framework specifying who holds decision authority, at which checkpoint, with what accountability surviving audit. The specification is published, and the author's operational record includes documented instances where a human governor overruled full AI platform consensus (Puglisi, 2026b).

The non-cognitive governance layer that cannot understand the engine by design is **GOPEL** (Governance Orchestrator Policy Enforcement Layer), a proof-of-concept reference implementation performing seven deterministic operations with zero cognitive evaluation. The

code exists with fourteen source modules, nine test suites, and 183 tests with zero failures (Puglisi, 2026c).

The structural diversity mandate preventing single-platform lock-in is **AI Provider Plurality**, a policy framework with a legislative package, technical appendix, and Congressional submission shared by the author with congressional offices in February 2026 (Puglisi, 2026a).

The cross-platform review protocol with source-authority discrimination and preserved dissent is **HAIA-CAIPR** (Cross AI Platform Review), which operationalizes the multi-AI workflow across eleven platforms in active use with synthesizer audit requirements (Puglisi, 2026d).

The broader ecosystem that holds these together is **HAIA** (Human AI Assistant), published open-source at github.com/basilpuglisi/HAIA under Creative Commons license. Other governance approaches exist, including the EU AI Act's enforcement mechanisms and risk classifications, the NIST AI Risk Management Framework's Govern, Map, Measure, and Manage structure, centralized licensing models, and standards-body certification frameworks. GOPEL's audit trail could serve as a compliance instrument within those frameworks rather than replacing them.

To make this concrete rather than abstract: consider Harris's own recommendation for mandatory testing of AI systems before deployment. Under this architecture, mandatory testing would move through a specific sequence. Three or more independent AI providers assess the same system under identical conditions (AI Provider Plurality dispatch). Every response is collected without modification and logged in a hash-chained, append-only audit trail (GOPEL logging). Where the providers disagree, the disagreement surfaces as a signal rather than noise, and a named human authority reviews the divergence before the assessment advances (CBG checkpoint). The dissenting assessment stays in the permanent record regardless of the final decision (CAIPR preserved dissent). That sequence is not a policy concept. It is a specified, testable workflow with code behind it. Whether it is the right sequence is an open question. That it exists as a starting point is the contribution.

None of this is finished, and the gaps are significant. GOPEL remains a proof of concept rather than a production deployment. AI Provider Plurality is a submitted Congressional Package with no hearings scheduled as of March 2026. CBG has operational evidence but lacks independent validation from implementations outside the author's own practice, and the HEQ measurement instrument needs fresh empirical data before academic submission. These gaps are documented rather than hidden, because governance that conceals its limitations is performing the same maneuver Harris accuses the AI companies of performing.

The point is not that the work is done. The point is that the work is moving. Harris asks what the person at home can do. One answer: the specifications are published, the code is open-source, and the author's Congressional submission is on the public record. Engaging with that infrastructure, testing it, challenging it, and building alternatives to it is a more operational response than sharing a video, though sharing the video creates the political permission that any governance effort requires.

With the architecture on the table, the paper owes the reader an honest accounting of what it claims and what it does not.

What This Paper Does Not Claim

This paper does not claim Tristan Harris is wrong. His structural diagnosis of the AI race's incentive architecture is sound, and his ability to deliver that diagnosis to millions of people through a two-hour conversation represents a form of public service that governance specifications cannot replicate.

This paper does not claim the HAIA ecosystem solves the problems Harris identifies. Nothing in this work is proven. The ecosystem represents one early attempt to move from warning to infrastructure, published so others can test it, challenge it, and build something better.

This paper does not claim Harris should be on the Thought Leader Master Grid. His value is as a public communicator and structural pattern recognizer. The Grid classifies by function, and Harris's function is advocacy rather than research, economic analysis, or infrastructure design.

This paper does not claim that awareness and infrastructure are in competition, because they are complementary. Harris provides the political permission that governance requires, and published infrastructure provides something for that permission to act on. Both are necessary, and neither alone is sufficient. Warnings change attention. Infrastructure changes what institutions can do next.

References

- Acemoglu, D., & Johnson, S. (2023). *Power and progress: Our thousand-year struggle over technology and prosperity*. PublicAffairs.
- Anthropic. (2025a, July). System card: Claude Opus 4 & Claude Sonnet 4. <https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf>
- Anthropic. (2025b, June). Agentic misalignment: How LLMs could be insider threats. <https://www.anthropic.com/research/agentic-misalignment>
- Bartlett, S. (Host). (2025, November 27). AI expert: Here is what the world looks like in 2 years! Tristan Harris [Video]. The Diary of a CEO. YouTube. <https://www.youtube.com/watch?v=BFU1OckhBwo>
- Bowles, N. (2016, November). The binge breaker. *The Atlantic*. <https://www.theatlantic.com/magazine/archive/2016/11/the-binge-breaker/501122/>
- Brynjolfsson, E., & McAfee, A. (2014). *The second machine age: Work, progress, and prosperity in a time of brilliant technologies*. W. W. Norton.
- Brynjolfsson, E., Chandar, B., & Chen, R. (2025). The early impact of generative AI on the labor market: Evidence from high-frequency payroll data [Working paper]. Stanford Digital Economy Lab. <https://digitaleconomy.stanford.edu/>
- Center for Humane Technology. (n.d.-a). Tristan Harris. <https://www.humanetech.com/team-board/tristan-harris>
- Center for Humane Technology. (n.d.-b). Impact and story. <https://www.humanetech.com/impact-and-story>
- Center for Humane Technology. (n.d.-c). For policymakers. <https://www.humanetech.com/policymakers>
- Center for Humane Technology. (n.d.-d). Litigation case study: Character.AI and Google. <https://www.humanetech.com/case-study/litigation-case-study-character-ai-and-google>
- ControlAI. (2025a). The Direct Institutional Plan. <https://controlai.com/dip>
- ControlAI. (2025b). A narrow path. <https://controlai.news/p/a-narrow-path>
- Crawford, K. (2021). *Atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press.
- EY. (2025). EY AI Confidence Barometer. <https://www.ey.com/ai-confidence-barometer>
- Financial Crisis Inquiry Commission. (2011). *The Financial Crisis Inquiry Report*. U.S. Government Publishing Office. <https://www.govinfo.gov/content/pkg/GPO-FCIC/pdf/GPO-FCIC.pdf>
- Greenblatt, R., et al. (2024, December). Alignment faking in large language models. Anthropic & Redwood Research. <https://www.anthropic.com/research/alignment-faking>

- Harris, T. (2013). A call to minimize distraction and respect users' attention [Internal presentation]. Google.
- Harris, T. (2017). How a handful of tech companies control billions of minds every day [TED Talk]. TED2017. https://www.ted.com/talks/tristan_harris_how_a_handful_of_tech_companies_control_billions_of_minds_every_day
- Harris, T. (2019, December 5). Our brains are no match for our technology. *The New York Times*. <https://www.nytimes.com/2019/12/05/opinion/digital-technology-brain.html>
- Harris, T., & Raskin, A. (2023, March 9). The AI dilemma [Presentation]. Center for Humane Technology.
- Hinton, G. (2023, May 1). "The godfather of A.I." leaves Google and warns of danger ahead [Interview by C. Metz]. *The New York Times*. <https://www.nytimes.com/2023/05/01/technology/ai-google-chatbot-engineer-quits-hinton.html>
- Khan, L. M. (2017). Amazon's antitrust paradox. *Yale Law Journal*, 126(3), 710-805.
- MERICS. (2021, June 24). Lofty principles, conflicting incentives: AI ethics and governance in China. <https://merics.org/en/report/lofty-principles-conflicting-incentives-ai-ethics-and-governance-china>
- Orlowski, J. (Director). (2020). *The Social Dilemma* [Documentary]. Exposure Labs; Netflix.
- Puglisi, B. C. (2025a). *Governing AI: When capability exceeds control*. ISBN 9798349677687.
- Puglisi, B. C. (2025b). When warnings are right but methods are wrong. basilpuglisi.com. <https://basilpuglisi.com/when-warnings-are-right-but-methods-are-wrong/>
- Puglisi, B. C. (2026a). AI Provider Plurality: An infrastructure mandate for democratic AI systems. AI Provider Plurality Congressional Package. <https://github.com/basilpuglisi/HAIA>
- Puglisi, B. C. (2026b). Checkpoint Based Governance v5.0. <https://github.com/basilpuglisi/HAIA>
- Puglisi, B. C. (2026c). GOPEL: Governance Orchestrator Policy Enforcement Layer v1.5. Proof-of-concept reference implementation. <https://github.com/basilpuglisi/HAIA>
- Puglisi, B. C. (2026d). HAIA-CAIPR: Cross AI Platform Review specification v1.1. <https://github.com/basilpuglisi/HAIA>
- Puglisi, B. C. (2026e). HAIA framework architecture. SSRN Abstract ID 6195238. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=6195238
- Puglisi, B. C. (2026f). *The Minds That Bend the Machine: The voices shaping responsible AI governance*. In preparation. basilpuglisi.com.
- Puglisi, B. C. (2026g). The warning, the override, and the infrastructure. basilpuglisi.com. <https://github.com/basilpuglisi/HAIA>
- Raskin, J. (2000). *The humane interface: New directions for designing interactive systems*. Addison-Wesley.

Reuters. (2024, November 16). Biden, Xi agree that humans, not AI, should control nuclear arms. *Reuters*. <https://www.reuters.com/world/biden-xi-agreed-that-humans-not-ai-should-control-nuclear-weapons-white-house-2024-11-16/>

Russell, S. (2019). *Human compatible: Artificial intelligence and the problem of control*. Viking.

Toner, H. (2023). AI governance and policy research. Center for Security and Emerging Technology, Georgetown University. <https://cset.georgetown.edu/team/helen-toner/>

U.S. Congress. (2020, January 8). Americans at risk: Manipulation and deception in the digital age [Hearing]. House Committee on Energy and Commerce, Subcommittee on Consumer Protection and Commerce. <https://www.congress.gov/event/116th-congress/house-event/110351>

U.S. Senate. (2019, June 25). Optimizing for engagement: Understanding the use of persuasive technology on internet platforms [Hearing]. Senate Committee on Commerce, Science, and Transportation, Subcommittee on Communications, Technology, Innovation, and the Internet. <https://www.commerce.senate.gov/2019/6/optimizing-for-engagement-understanding-the-use-of-persuasive-technology-on-internet-platforms>

U.S. Senate. (2021, April 27). Algorithms and amplification: How social media platforms' design choices shape our discourse and our minds [Hearing]. Senate Committee on the Judiciary, Subcommittee on Privacy, Technology, and the Law. <https://www.judiciary.senate.gov/download/tristan-harris-testimony>

Yudkowsky, E. (2022). AGI ruin: A list of lethalties. LessWrong. <https://www.lesswrong.com/posts/uMQ3cqWDPHhjtiesc/agi-ruin-a-list-of-lethalties>

Basil C. Puglisi, MPA, is a Human-AI Collaboration Strategist and AI Governance practitioner operating independently via basilpuglisi.com. All frameworks referenced in this paper are published open-source at github.com/basilpuglisi/HAIA under Creative Commons Attribution-NonCommercial 4.0 International license. SSRN Abstract ID 6195238. #AIassisted