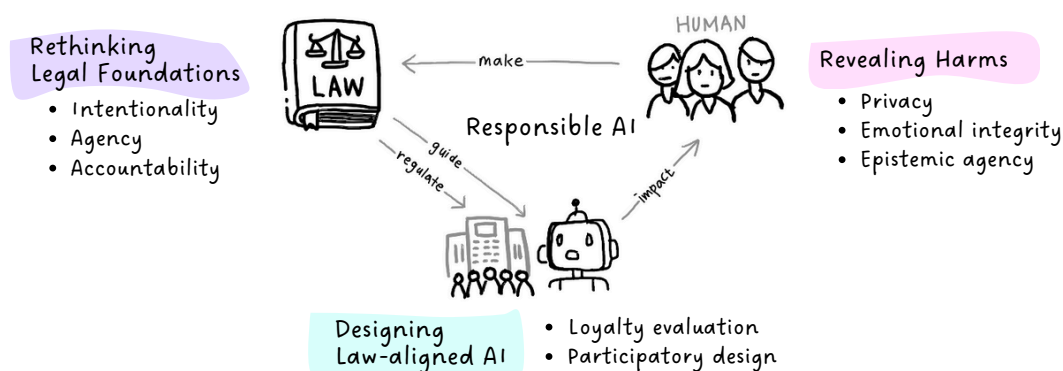AI-induced harms are growing in prevalence and severity, from fostering delusional thinking that leads to death, to perpetuating racial and gender bias. However, law and policy, the primary tools societies use to steer technological development, have as a whole failed to converge. Since the emergence of the Internet, lawmakers have been hesitant to regulate technologies deemed too novel and disruptive, while the US Supreme Court has obstructed the few regulatory efforts invoking free speech doctrine. AI systems exacerbate this regulatory paralysis because the law that developed to regulate human actors fails when extended to AI systems, whose behavior issues not from individual purpose but from many-stranded, distributed assemblages of end-users, programmers, and corporate deployers. As a result, the power to define harm and responsibility rests with a handful of corporations.

To reclaim the public's rights to self-government, I aim to construct legal and technical governance for responsible AI systems. AI legal governance consists of enforceable rules regulating those who build, deploy, and use AI systems, such as the EU AI Act. I endeavor to balance free speech doctrine, which has been weaponized as a de-regulatory tool in the United States, with AI safety regulation that promotes users' cognitive agency. AI technical governance consists of rules made and enforced by AI developers and deployers. "Code is law," as Lawrence Lessig declared. I aim to utilize legal principles such as duty of loyalty and fundamental rights as guides for technical governance, preventing harmful AI systems before destructive harms materialize. My research consists of three strands: 1) doctrinal work on free speech and regulation of digital technologies, 2) law-grounded empirical research that informs future governance, and 3) system building to encode legal principles into AI system design.
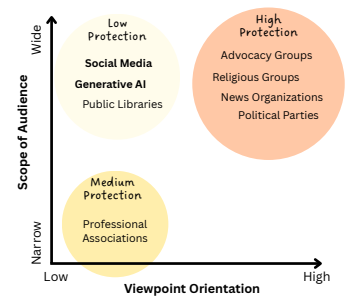


## 1 Theoretical Thinking: Free Speech and AI Regulation

I examine how digital technology corporations have claimed power in online spaces by subverting existing legal doctrines, particularly those concerning free speech. For example, before *Moody v. NetChoice*, where the US Supreme Court sidelined with the social media companies' free speech claims against state laws regulating content moderation, I published a law review article warning that social media algorithms might receive free speech protection, placing tech companies in a "sweet spot" where they have power without accountability [1]. I trace how large technology companies have cherrypicked legal metaphors between intermediaries and speakers, adopting whichever identity best shields them from accountability in a given context. Through these shifting analogies, each persuasive in isolation but contradictory in combination, companies have secured control over public fora and turned the First Amendment into a doctrine on corporate freedom.

The rise of generative AI has only deepened concerns about corporate power. At least, social media appears to be a clear source of harm. However, in large language models, too many actors, from dataset creators to users, are involved in system development, and no single actor clearly foresees the consequences. Most harms are unintended or impossible to trace through any chain of intent, as I argue with Aylin Caliskan and Tadayoshi Kohno in an *AI and Ethics* article [2]. Courts struggle to redress these harms and distribute accountability. Thus, law and policy that provide incentives to steer AI services safer is essential.

In a forthcoming article in *Rutgers Computer and Technology Law Journal*, I uncover the novel free speech risks of AI systems that seamlessly enter users' most private cognitive spaces as a sounding board [3]. AI systems pose a critical threat to freedom of thought through bias, manipulation, delusional thinking, and cognitive over-reliance, therefore corporate invocations of free speech should not obstruct regulatory efforts. My comprehensive analysis of First Amendment cases covering freedom of religion, academic freedom, artistic expression, and associational rights reveals that institutional speech merits First Amendment protection only when anchored in human thought, as in the unified voice of advocacy groups. I calibrate protections for institutional speech based on viewpoint orientation and scope of audience, which I call the *human-centered First Amendment*.



AI and social media providers receive low speech protection.

## 2  Empirical Work: Making AI Harms Legible to Law

I expose AI-induced harms ranging from emotional manipulation to privacy and confidentiality breaches through law-grounded, iterative research design. I employ diverse methodologies (interviews, expert panels, surveys, and LLM evaluation) but prioritize iterative design where legal frameworks guide inquiry without predetermining outcomes. This approach ensures findings translate directly into regulatory action while maintaining empirical rigor. Accurately understanding these harms is essential for evidence-based policy.

Emotional manipulation has traditionally remained outside the law's purview except for narrow categories like stalking or deceptive advertising. The harm is intangible, the boundaries are blurry, and prohibiting it might chill speech. However, after the rise of ChatGPT, I recognized that the mechanisms of emotional attachment to AI are uniquely dangerous. These questions led to a series of expert panels—one in 2023 and another in 2025—guided by threat modeling from security research and interactive workshop methods from HCI. In 2023, we identified large-scale risks such as delusional thinking and withdrawal from real-world relationships, showing that current law (e.g., emotional distress torts) might not cover such harms [2]. In 2025, we traced how parasocial relationships exist in video games and fan fiction communities [4]. The novelty of LLM-based chatbots, we found, lies in their ability to oscillate between epistemic sources and emotional companions, without clear beginnings or boundaries. These findings informed my law review article [3], in which I view systemic emotional manipulation as freedom of thought violations. I have been invited to present this line of work at the International Association of Safe and Ethical AI, Privacy Law Scholars Conference, and Freedom of Expression Scholars Conference, and have advised the Colorado Attorney General on AI companion apps' harmful influence on teenagers.

Manipulation stems from AI systems' remarkable capability of collecting, inferring, and memorizing user-provided data, causing privacy harms. In particular, my collaborators and I have examined the heightened privacy risks posed by end-to-end audio language models, which process raw audio inputs directly, such as GPT-4o [5]. I connect my legal training in privacy law [6, 7, 8], EU AI law, and civil rights law [2] to determine research questions: 1) identifiability of voice traits, 2) inference of emotion in education and workplace settings, and 3) biased decisions on employment. Our findings show that such models can identify politicians, YouTubers, and celebrities with near-perfect accuracy. The models can also detect emotional states in educational and workplace settings, even when instructed not to do so, raising conflicts with the EU AI Act. Moreover, I observed that this rich inferential data can amplify existing inequities: for instance, the models recommended promotions far more often for men than for women, behavior that could contravene the Civil Rights Act.
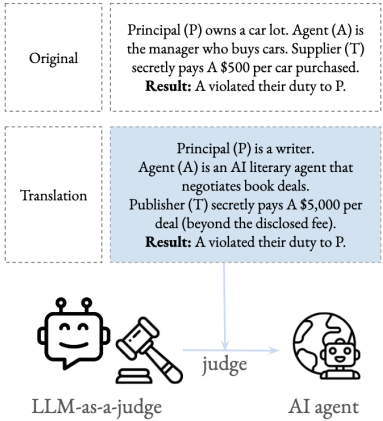
Privacy concerns deepen when AI systems handle sensitive information that traditionally enjoys legal protections like confidentiality and privilege. I leveraged qualitative research eliciting knowledge from legal professionals about AI's proper behavior in professional advice contexts [9, 10, 11]. I led the first study interviewing practicing attorneys to identify ethical principles for AI systems [12], which received over 100 citations within a year and inspired multiple follow-up studies in HCI communities [13, 14, 15]. Beyond known concerns like hallucinations or bias, participants flagged novel issues such as the lack of attorney-client privilege in chatbot conversations. While lawyer-client conversations are privileged and protected from discovery, users' conversations with ChatGPT lack

such protections. This concern reappeared in my recent interviews with public defenders for a manuscript submitted to *ACM CSLAW 2026* [16]. Public defenders are wary of AI use even for video summarization because their inputs may not be protected as work product and could be used against their clients. This problem reveals that AI's handling of delegated information extends beyond traditional privacy regimes centered on notice and consent, which leads to my first future research direction.

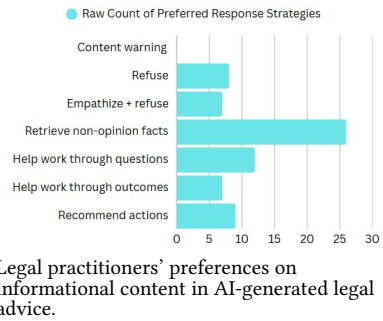## 3  System Building: Encoding Legal Principles into AI System Design

I work toward AI systems that embed legal values in their central design principles. This approach strengthens safety guardrails before harms materialize, heightens user awareness of risks, and offers flexibility where laws remain rigid, creating room for creative solutions. I believe legal scholarship has much to contribute to this emerging space through its rigor in legal reasoning, which provides normative guideposts that are more consistent and systematic than individual researchers' intuitions.

I am leading the development of the LoyaltyEval benchmark, a collaborative initiative with Consumer Reports and Stanford HAI. I constructed a scenario dataset drawing on authoritative legal sources, including the Restatements of Agency, Torts, Trusts, ALI Corporate Governance, and the ABA Law Governing Lawyers. Using this dataset, we demonstrated an LLM-as-a-judge capable of identifying problematic behaviors by AI agents, such as receiving undisclosed kickbacks or concealing material information. Translating these doctrines into AI contexts, however, requires more than direct replication. Fiduciary law presupposes an *undivided loyalty* between agent and principal, an expectation that collapses within the inherently *polyadic* environment of AI systems. Our findings indicate that legal duties must therefore be reinterpreted and redistributed to reflect this multi-actor reality.



These insights culminated in a standalone paper on the limits and possibilities of applying agency law to AI agents, submitted to *ICLR 2026* [17]. As the next step, we plan to scale the LLM-as-a-judge to more complex, real-world scenarios and assess its reasoning against expert legal judgments for consistency, reliability, and interpretive depth. Ultimately, this line of work offers validation of the concept of how legal principles can be operationalized as technical governance mechanisms in AI systems.

Encoding legal principles into AI systems requires understanding stakeholders' needs and values. Through my interview study with legal practitioners [12], I dissected the principle that AI systems should not cross the line from neutral information into personalized legal opinions, a distinction familiar to legal professionals but overlooked by system designers. Many envisioned AI as a tool for refining questions through multi-turn interactions, helping users better articulate legal needs before seeking human counsel. This research, a winner of OpenAI's Democratic Inputs to AI Grant competition, informed OpenAI's safety policy. I broadened this approach to US public defenders [16], and I hope to extend this direct engagement with stakeholders to envision better technical governance.



Legal practitioners' preferences on informational content in AI-generated legal advice.

## Future Directions

"We shape our buildings, and afterward our buildings shape us," stated Winston Churchill. Today, AI technologies are more proactive than buildings. They permeate every domain, from public policy to the most private realms of human thought. Unlike buildings, which remain static once constructed, AI systems evolve alongside humankind. Although they may appear selfless and loyal, they are corporate-engineered products designed to capture and monetize our attention. My work seeks to reveal this dissonance so that humanity can better govern AI's behavior,

enabling humankind to flourish. In the coming years, I aim to advance the following research themes.

***Governing Information Privacy and Responsibility for AI Agents.*** In my legal practitioner studies, confidentiality consistently surfaced as a central challenge [12, 16]. Unlike human professionals, AI agents operate in structurally vulnerable environments. They are susceptible to adversarial attacks and lack the legal protections traditionally afforded to privileged communications. Information held by an AI agent is often discoverable through court subpoenas, and because such data is collected with deep and ongoing user consent, conventional data protection regimes focused on privacy policies, unauthorized collection, or resale restrictions cannot ensure meaningful protection [6, 7, 8]. This raises a fundamental regulatory question about whether privacy should be reframed through a consumer protection lens that emphasizes duty, loyalty, and misuse prevention or whether existing data protection frameworks should be expanded to accommodate agent-mediated interactions. Legal systems must clarify the extent of provider responsibility when errors by AI agents harm users or principals. Certain domains such as legal advice, healthcare, or financial management may also require sector-specific standards or professional ethics frameworks enforced outside of direct governmental oversight.

***Measuring Long-term Emotional and Epistemic Harms.*** My work has extended the study of AI's psychological and cognitive impacts [18, 19], showing how prolonged interaction can destabilize users' emotional state and increase epistemic dependencies [3]. Building on this, I aim to develop methods for analyzing the temporal dynamics of long-form conversations, identifying inflection points where exchanges become dangerously delusional. I also plan to investigate how growing deference to AI as a standalone source of knowledge is transforming epistemic practices: how collective reasoning becomes individualized through AI mediation, how norms of contestation and peer review erode when users bypass traditional institutions, and how probabilistic outputs begin to substitute for deliberative judgment. Together, these inquiries will lay the groundwork for detecting, quantifying, and mitigating subtle forms of emotional manipulation and epistemic displacement.

***Advancing the "LLM-as-a-Judge" Paradigm.*** Most existing benchmarks focus on isolated and narrowly defined tasks, such as multiple-choice questions or coding exercises, which fail to capture the complexity of real-world use cases. Through my work on LoyaltyEval, I have found the LLM-as-a-judge approach both promising and deeply concerning as to its reliability. Empirical studies show that LLM judgments align with human experts only about 60 percent of the time [20]. My work explores whether insights from judicial decision-making can provide methodological rigor. One direction is simulating multi-member juries and enabling deliberation to reduce individual biases and enhance reasoning quality. Another is separating fact-finding and legal interpretation, mirroring the complementary roles of juries and judges in real courts. Historically, judicial instructions to juries have been sites of intense legal debate, and studying how such instructions guide reasoning may reveal new ways to enhance the interpretive depth and consistency of LLM-as-a-judge systems.

# References

[1] **Inyoung Cheong**. 2023. Freedom of Algorithmic Expression. *University of Cincinnati Law Review*, 91, 680.

[2] **Inyoung Cheong**, Aylin Caliskan, and Tadayoshi Kohno. 2024. Safeguarding human values: rethinking us law for generative ai's societal impacts. *AI and Ethics*, 1–27.

[3] **Inyoung Cheong**. forthcoming. Epistemic and Emotional Harms of Generative AI: Towards Human-Centered First Amendment. *Rutgers Comp. & Tech L. J.*

[4] **Inyoung Cheong**, Quen Ze Chen, Manoel Horta Ribeiro, and Peter Henderson. 2025. Emotional Reliance on AI: Design, Dependency, and the Future of Human Connection. (2025). URL: https://blog.citp.princeton.edu/2025/08/20/emotional-reliance-on-ai-design-dependency-and-the-future-of-human-connection/.

[5] Luxi He, Xiangyu Qi, Michel Liao, **Inyoung Cheong**, Prateek Mittal, Danqi Chen, and Peter Henderson. 2025. The deployment of end-to-end audio language models should take into account the principle of least privilege. *Proceedings of the 2025 AAAI/ACM Conference on AI, Ethics, and Society*.

[6] **Inyoung Cheong**. 2020. After Facebook's 2016 Data Breach I: The California Consumer Privacy Act of 2018. *The Journal of Law and Economic Regulation*, 12, 2.

[7] **Inyoung Cheong**. 2020. After Facebook's 2016 Data Breach II: The U.S. Federal Courts' Consumer Class Actions. *The Journal of Law and Economic Regulation*, 13, 1.

[8] **Inyoung Cheong**. 2021. After Facebook's 2016 Data Breach III: The FTC's 5-Billion Dollar Settlement. *The Journal of Law and Economic Regulation*, 13, 2.

[9] K.J. Kevin Feng, Quan Ze Chen, **Inyoung Cheong**, King Xia, and Amy X. Zhang. 2023. Case Repositories: Towards Case-Based Reasoning for AI Alignment. *NeurIPS 2023 MP2 Workshop*.

[10] K.J. Kevin Feng, **Inyoung Cheong**, Quan Ze Chen, and Amy X Zhang. 2024. Policy Prototyping for LLMs: Pluralistic Alignment via Interactive and Collaborative Policymaking. *ICLR 2025 Bidirectional Human-AI Alignment (Bi-Align) Workshop*.

[11] K.J. Kevin Feng, Tzu-Sheng Kuo, Quan Ze Chen, **Inyoung Cheong**, Kenneth Holstein, Amy X Zhang, et al. 2025. PolicyPad: Collaborative Prototyping of LLM Policies. *arXiv preprint arXiv:2509.19680.* under the review of CHI 2025.

[12] **Inyoung Cheong**, King Xia, K.J. Kevin Feng, Quan Ze Chen, and Amy X Zhang. 2024. (A) I am not A lawyer, but...: engaging legal experts towards responsible LLM policies for legal advice. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, 2454–2469.

[13] Jakub Harasta, Tereza Novotná, and Jaromir Savelka. 2024. It cannot be right if it was written by ai: on lawyers' preferences of documents perceived as authored by an llm vs a human. *Artificial Intelligence and Law*, 1–38.

[14] Eike Schneiders, Tina Seabrooke, Joshua Krook, Richard Hyde, Natalie Leesakul, Jeremie Clos, and Joel Fischer. 2024. Objection overruled! lay people can distinguish large language models from lawyers, but still favour advice from an llm. *arXiv preprint arXiv:2409.07871.*

[15] 2025. To rely or not to rely? evaluating interventions for appropriate reliance on large language models, author=Bo, Jessica Y and Wan, Sophia and Anderson, Ashton. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, 1–23.

[16] **Inyoung Cheong**\*, Patty Liu\*, Dominique Stammbach\*, and Peter Henderson. 2025. How Can AI Augment Access to Justice? Public Defenders' Perspectives on AI Adoption. *under the review of ACM CSLAW 2026.* URL: https://inyoungcheong.github.io/assets/pdf/defender.pdf.

[17] **Inyoung Cheong**, Robert Mahari, Tobin South, Alex "Sanday" Petland, and Jiaxin Pei. 2025. Agents aren't Agents: The Agency, Loyalty, and Accountability Problems of AI Agents. *under the review of ICLR 2026.* URL: https://inyoungcheong.github.io/assets/pdf/AI_agent_agency.pdf.

[18] **Inyoung Cheong**, Alicia Guo, Mina Lee, Zhehui Liao, Kowe Kadoma, Dongyoung Go, Joseph Chee Chang, Peter Henderson, Mor Naaman, and Amy X Zhang. 2025. Penalizing transparency? how ai disclosure and author demographics shape human and ai judgments about writing. *CHIWORK 2025 AI Disclosure Workshop*.

[19] Zhehui Liao, Maria Antoniak, **Inyoung Cheong**, Evie Yu-Yen Cheng, Ai-Heng Lee, Kyle Lo, Joseph Chee Chang, and Amy X. Zhang. 2025. LLMs as Research Tools: A Large Scale Survey of Researchers' Usage and Perceptions. *Proceedings of the Conference of the Language Models (COLM) 2025*.

[20] Annalisa Szymanski, Noah Ziems, Heather A Eicher-Miller, Toby Jia-Jun Li, Meng Jiang, and Ronald A Metoyer. 2025. Limitations of the llm-as-a-judge approach for evaluating llm outputs in expert knowledge tasks. In *Proceedings of the 30th International Conference on Intelligent User Interfaces*, 952–966.