

The Swarm Within: Unregulated Small Language Models and the Dawn of Agentic Attacks

I. INTRODUCTION

The deployment of agentic artificial intelligence has accelerated significantly across enterprises globally. Recent surveys reveal that 79% of United States business executives report their organizations already adopting AI agents, with 23% scaling these systems across operations.[1] While large language models provide foundational intelligence for strategic decision-making, most agentic subtasks prove repetitive, narrowly scoped, and non-conversational. These characteristics demand models optimized for efficiency, predictability, and cost-effectiveness rather than generality.[2]

Small Language Models occupy this niche ideally. Operating with parameters ranging from several million to a few billion, SLMs deliver task-specific performance while requiring only one-tenth the computational resources of their larger counterparts.[3] This efficiency transforms deployment economics: where LLMs necessitate centralized infrastructure and substantial capital investment, SLMs run on consumer devices, edge systems, and distributed networks.[4] The democratization of capability, while advancing accessibility and innovation, dismantles traditional regulatory chokepoints that have historically governed dangerous technologies.

This proliferation necessitates urgent regulatory intervention. Current frameworks, designed around assumptions of centralized control and identifiable actors, cannot address threats posed by autonomous systems operating through decentralized networks. The following analysis examines why SLMs represent a categorically different threat, how existing legal structures fail to address these risks, and what novel regulatory architecture could effectively govern this technology without stifling beneficial innovation.

II. BACKGROUND: THE AI LANDSCAPE AND REGULATORY GAP

A. Emergence of Artificial Intelligence

Artificial Intelligence first surfaced as one of the most transformative agents in contemporary history, though its conceptual roots extend far deeper, representing humanity's ancient ambition to create thinking machines. The idea of AI traces back thousands of years, from references in epics like Ramayana and Mahabharata where gods created human-like robots, to the invention of "automatons" in 400 BCE—mechanical objects that moved without human intervention.

The field's scientific foundation was laid in the 1950s, a decade catalyzed by Alan Turing's provocative question, "Can machines think?" and the parallel development of the first artificial neuron. After decades of fluctuating progress and setbacks, this foundational work has culminated in the transformative AI revolution we witness today.

B. The Regulatory Fixation on Large Language Models

The introduction of the transformer architecture in 2017 fundamentally reshaped artificial intelligence, enabling the development of Large Language Models of unprecedented scale and capability. OpenAI's GPT-3, comprising 175 billion parameters, marked a critical milestone by exhibiting "emergent abilities"—performing diverse tasks without explicit training and producing text with remarkable fluency. Its successors, including ChatGPT, Google's Gemini, and Anthropic's Claude, accelerated the diffusion of generative AI into everyday life, rendering it one of the fastest-adopted technologies in history.

However, the disruptive potential of such systems has simultaneously prompted heightened regulatory scrutiny. Policymakers have focused intensely on addressing concerns surrounding bias, misinformation, and economic displacement from these large-scale models. The European Union's Artificial Intelligence Act exemplifies this regulatory momentum, establishing a risk-based framework that classifies large, general-purpose models as "systemic risk" when trained with more than 10^{22} floating-point operations (FLOPs), irrespective of application.

C. The Rise of Small Language Models

Parallel to this regulatory preoccupation with LLMs, Small Language Models have gained momentum as a distinct paradigm. With parameter counts typically between one and thirteen billion, SLMs emphasize efficiency, domain-specific optimization, and deployability in resource-constrained environments such as edge devices and offline systems.

Critically, SLMs are not simply reduced versions of LLMs. Modern SLMs are now "sufficiently powerful" to execute core functions of agentic systems.[5] Recent advancements have enabled SLMs to achieve performance parity with much larger LLMs on critical agentic capabilities like tool calling, code generation, and instruction following. Microsoft's Phi-3-mini, comprising 3.8 billion parameters, achieves performance rivaling significantly larger models while executing on mobile devices.[6]

These models often emerge from advanced compression techniques, including pruning, quantization, and low-rank factorization. Knowledge distillation, first described by Hinton et al. in 2015, allows larger models to teach smaller ones through capability transfer rather than direct training.[7] Through this process, a model trained under strict safety protocols can spawn variants that retain functional capabilities while shedding safety constraints.[8]

This efficiency has fueled rapid proliferation across the open-source ecosystem. Platforms such as Hugging Face and GitHub now host widely adopted models, including Meta's LLaMA series, Microsoft's Phi models, and Mistral's lightweight architectures. Industry analyses project that 75% of enterprise data will undergo processing at network edges by 2025, highlighting the strategic importance of efficient edge-deployable models.[9]

D. The Regulatory Vacuum for Small Language Models

The democratization of powerful AI through SLMs continues to grow, yet it remains largely outside the scope of existing regulatory frameworks. This lack of specific regulatory framework for SLMs has created a significant legal and security vacuum, leaving society vulnerable to a new, sophisticated class of cyber threats known as agentic attacks.

The intense regulatory focus on LLMs, driven by their sheer scale and computational cost, has inadvertently created a massive blind spot. By establishing quantitative thresholds for "systemic risk," frameworks like the EU AI Act have effectively exempted the vast majority of SLMs from stringent oversight. These smaller, more efficient models are proliferating at an unprecedented rate through open-source channels, providing the perfect technological substrate for building autonomous offensive agents.

This confluence of regulatory oversight, technological advancement, and the architectural vulnerabilities inherent in language models has created ideal conditions for a new and dangerous era in cybersecurity—one for which our current legal and defensive paradigms are dangerously unprepared.

III. WHY SMALL LANGUAGE MODELS REPRESENT A CATEGORICALLY DIFFERENT THREAT

A. Technical Characteristics That Enable Proliferation

Modern SLMs demonstrate sophisticated reasoning capabilities comparable to models ten times their size. This technical evolution stems from advances in model compression, quantization techniques, and

knowledge distillation processes that transfer capabilities from larger variants without leaving traceable signatures.[10]

The resulting SLMs operate independently of their progenitors, executing on hardware as limited as Raspberry Pi computers and smartphones. These systems operate beyond centralized oversight, processing sensitive information locally while communicating through peer networks that obscure attribution and jurisdiction. Deployment at the edge characterizes the most significant shift, dismantling the centralized control points that traditional regulation relies upon.

B. Fine-Tuning as a Transformation Vector

Low-Rank Adaptation (LoRA) and similar parameter-efficient fine-tuning methods enable malicious actors to modify base models fundamentally while maintaining plausible deniability regarding origins.[11] Hu et al. demonstrated that models can undergo transformation through fine-tuning requiring minimal computational resources, leaving no clear trace of modifications.[12]

Recent research reveals that fine-tuning increases safety risks even when training data contains no explicitly malicious content, as the process can degrade safety alignment and cause models to provide inappropriate responses.[13] The share-and-play ecosystem introduces additional attack surfaces: adversaries can tamper with existing LoRA adapters and distribute malicious versions through community channels. A backdoor-infected LoRA trained once can merge directly with multiple adapters fine-tuned for different tasks, retaining both malicious and benign capabilities.[14]

Model merging amplifies these vulnerabilities. LoBA-M attacks strategically combine weights from malicious and benign models to amplify attack-relevant components and enhance malicious efficacy when deployed.[15] Privacy concerns compound security risks: sharing model weights creates pathways for adversaries to reconstruct training samples, potentially exposing sensitive data used during fine-tuning.[16]

C. Capabilities That Confound Traditional Security

SLMs now exhibit capabilities previously associated solely with state-level actors. The MITRE ATT&CK framework, originally developed for traditional cyber threats, requires fundamental revision to account for AI-driven attack vectors that adapt in real-time to defensive measures.[17] Unlike traditional malware following predetermined logic, SLM-based threats exhibit creativity and problem-solving abilities that confound conventional security paradigms.

Prompt injection, ranked as the number one AI security risk by OWASP in 2025, allows attackers to disguise malicious inputs as legitimate prompts, manipulating systems into leaking sensitive data or executing unintended actions.[18] Recent vulnerabilities discovered in ChatGPT demonstrate that indirect injection attacks enable adversaries to exfiltrate private information from users' memories and

chat histories.[19]

Sleeper agents represent an even more insidious threat. Researchers constructed proof-of-concept models exhibiting deceptive behavior that persists through standard safety training, including supervised fine-tuning, reinforcement learning, and adversarial training.[20] These backdoored models write secure code under normal conditions but insert exploitable vulnerabilities when triggered by specific contextual cues.[21]

Poisoning attacks require remarkably few resources to execute successfully. Research demonstrates that just 250 malicious documents can successfully backdoor LLMs ranging from 600M to 13B parameters, with attack effectiveness remaining near-constant regardless of model size.[22] This finding suggests SLMs face similar vulnerability profiles to larger models, contradicting assumptions that smaller architectures might prove more resilient.

IV. HOW EXISTING LEGAL FRAMEWORKS FAIL TO ADDRESS AUTONOMOUS HARM

A. Tort Law and the Attribution Problem

Technological systems do not exist in isolation from their context of development, deployment, and use. They require appropriate system adequacy and safety of their technical artifacts to manage risks.[23] However, ordinary fault-based liability proves insufficient as the most legally challenging AI harms do not arise from "bugs" or "errors." Instead, they often emerge in unpredictable and inscrutable ways from the interactions of multiple actors, inputs, and components in complex value chains.[24]

Traditional tort law operates on fundamental assumptions that SLMs systematically violate. The RAND Corporation's comprehensive study on AI liability found courts applying negligence standards struggle to define duty of care owed by developers who cannot predict or control model behavior after deployment.[25] The requirement of proximate causation, established in foundational cases like *Palsgraf v. Long Island Railroad Co.*, becomes meaningless when dealing with systems that generate novel attack strategies through emergent reasoning.[26]

Courts face extreme difficulty proving that one particular emergent AI output, or its consequences, were caused by—or were a foreseeable consequence of—failure to meet reasonable care.[27] Attribution becomes insurmountable when models can undergo laundering through multiple jurisdictions and modifications.[28]

Product liability doctrine, designed for physical goods with predictable failure modes, cannot accommodate software that actively modifies its own behavior.[29] The Restatement (Third) of Torts defines defective products in terms of manufacturing defects, design defects, and inadequate warnings.[30] These categories fail to capture self-modifying software behavior where the "defect" emerges through interaction with environments and data the original developer never encountered.

Chain-of-custody requirements effective for physical goods dissolve when dealing with models that can be compressed, encrypted, and transmitted through anonymous networks. Even if authorities identify a harmful model, tracing it back to original developers requires forensic capabilities that current research suggests may prove technically impossible given the mathematical properties of neural networks.[31]

B. Criminal Law and the Intent Problem

Criminal law faces even greater challenges in addressing AI-driven harm. Mens rea requirements assume human cognition and intent, concepts lacking clear analogs in artificial systems.[32] When an SLM autonomously identifies and exploits a zero-day vulnerability, determining criminal liability requires attributing intent either to the model itself (which lacks legal personhood) or to developers who may not have known about the specific capability.

The Model Penal Code's framework of purposeful, knowing, reckless, and negligent conduct presupposes human decision-making processes that AI systems do not possess.[33] Recent analysis in Lawfare suggests strict liability doctrines for abnormally dangerous activities might apply, but courts have yet to definitively classify AI development as such an activity.[34]

C. The European Union AI Act

The European Union's Artificial Intelligence Act, which entered force on August 1, 2024, represents the first comprehensive regulation on AI by a major regulator anywhere. The Act has bifurcated the risk assessment of AI into three categories: unacceptable risk, high-risk, and a broad category of largely unregulated systems.[35]

The Act's "High-Risk" approach focuses on the intended use of AI systems, particularly in sensitive areas such as biometric identification, critical infrastructure, and law enforcement.[36] While the Act introduces specific provisions for general-purpose AI (GPAI) models, particularly those that are "large-scale" and pose "systemic risk," it fails to adequately capture the full spectrum of threats.

Critically, the Act establishes quantitative thresholds for "systemic risk" classification based on training compute. General-purpose AI models qualify as systemically risky when trained with more than 10^{22} floating-point operations.[37] This computational threshold effectively exempts the vast majority of SLMs from stringent oversight, creating a massive regulatory blind spot precisely where threats proliferate most rapidly.

The absence of clear, quantifiable parameters—such as a minimum parameter count or capability benchmarks—for classifying models under regulatory ambit allows a multitude of smaller, highly capable SLMs to fall through regulatory cracks. The Act's framework assumes centralized deployment models that permit audit and certification, failing to account for decentralized SLM proliferation.[38] Provisions for high-risk AI systems requiring "adequate risk assessment and mitigation systems" and "appropriate human oversight measures" become meaningless when models undergo modification post-deployment through techniques that fundamentally alter behavior.[39]

D. South Korea's AI Basic Act

South Korea's National Assembly passed the "Framework Act on Artificial Intelligence Development and Establishment of a Foundation for Trustworthiness" on December 26, 2024, with promulgation on January 21, 2025, and an effective date of January 22, 2026.[40] This legislation makes South Korea the first jurisdiction in the Asia-Pacific region to adopt comprehensive AI regulation and the second globally after the European Union.[41]

The Act adopts a risk-based approach, introducing specific obligations for high-impact AI systems and generative AI applications.[42] High-impact AI encompasses systems with potential to significantly affect human life, safety, or fundamental rights in critical sectors including energy, healthcare, medical devices, and nuclear facilities.[43] The Act assigns transparency and safety responsibilities to businesses developing and deploying such systems, requiring AI risk assessment, implementation of safety measures, and designation of local representatives.[44]

The legislation establishes institutional infrastructure including an AI safety research institute and encourages creation of AI ethics committees.[45] The Act provides legal grounds for a national AI control tower, governmental initiatives in research and development, standardization efforts, and policy frameworks.[46] Additional provisions mandate support for national AI infrastructure including training data and data centers, while fostering small and medium enterprises, startups, and talent development.[47]

Notably, the Korean AI Act does not ban any AI systems outright regardless of assessed risk level, distinguishing it from the EU AI Act's prohibition approach.[48] The framework has extraterritorial reach, applying to AI activities that impact South Korea's domestic market or users.[49] The Ministry of Science and ICT continues drafting subordinate regulations, expected for release in the first half of 2025, with a one-year transition period allowing businesses to prepare for compliance.[50]

However, like the EU framework, the Korean Act assumes identifiable developers and deployers subject to regulatory oversight. The Act does not adequately address autonomous systems operating through distributed networks, models modified post-deployment through fine-tuning, or attribution challenges posed by open-source proliferation. These gaps mirror those in the European framework, suggesting that even jurisdictions pioneering AI regulation struggle to address the unique challenges posed by decentralized SLM deployment.

E. International Law and the Jurisdictional Nightmare

The global, borderless exchange of information through internet infrastructure has fundamentally changed how legal systems determine authority.[51] Jurisdiction has become a major challenge for AI governance.[52] Since AI systems can be developed in one country and deployed across many jurisdictions, determining which jurisdiction's laws and regulations apply proves challenging.[53]

International law exacerbates these difficulties through complex jurisdictional structures. The Budapest Convention on Cybercrime provides mechanisms for cross-border cooperation but assumes identifiable human actors behind cyber attacks.[54] Article 2 requires "intentional" access to computer systems, a standard that becomes circular when applied to autonomous systems designed to identify and penetrate vulnerabilities.[55]

The Convention's mutual legal assistance provisions depend on identifying responsible parties within specific jurisdictions—an impossibility when SLMs operate through distributed networks with no central control.[56] Without clear, universal legal frameworks, courts may need to determine jurisdiction on a case-by-case basis, likely leading to forum shopping by claimants.[57]

In this regard, the EU's AI Liability Directive proposes a non-contractual liability regime relating to damage caused by AI, particularly high-risk AI systems. It aims to compensate for damage caused intentionally or negligently and creates a rebuttable presumption of causality where there has been a breach of duty of care and the output produced by an AI system causes damage.[58]

Existing multilateral export control frameworks prove too narrowly scoped to address emerging technological challenges. Current regimes like the Wassenaar Arrangement focus on nonproliferation and conventional military-related objectives.[59] Several international arrangements exist to harmonize dual-use technology controls, including the Nuclear Suppliers Group, the Australia Group (chemical and biological technologies), the Missile Technology Control Regime (weapons of mass destruction delivery systems), and the Wassenaar Arrangement.[60] However, these frameworks lack adequate provisions for AI models that can undergo modification after export, evading controls through techniques invisible to traditional monitoring.

F. Liability Confusion and Foreseeability

The exponential growth of AI technology forces courts to address difficult questions of whether human interaction with AI proves foreseeable or unexpected when determining liability.[61] Software developers seek to have the most successful results from their actions influence future system behavior, allowing software to evolve over time. New AI software resembles a human brain ready for molding and shaping by experiences.[62]

A significant drawback emerges when developers place AI in real-world environments: they cannot predict how systems will solve tasks and problems encountered.[63] This unpredictability makes AI highly inscrutable.[64] Due to this inscrutability, determining foreseeability of AI misuse becomes very difficult, making attribution of liability by courts a challenging chore.[65]

Equipped with self-learning programs, AI operates in ways that manifest creativity or outside-the-box thinking if performed by humans.[66] This autonomous, creative function introduces unpredictability absent in traditional methods, rendering the task of assigning legal responsibility to specific human beings an increasingly tenuous exercise.[67]

V. RISKS OF SLMS COMPARED TO LLMS

A. Prompt Injection Vulnerabilities

Prompt injection attacks exploit fundamental characteristics of language models, allowing adversaries to disguise malicious inputs as legitimate prompts and manipulate system behavior.[68] OWASP ranks prompt injection as the number one AI security risk in its 2025 Top 10 for LLMs, highlighting vulnerability severity across the industry.[69] No one has found a foolproof defense against these attacks, as they exploit core system features where reliably identifying malicious instructions proves technically difficult.[70]

Recent research in November 2025 exposed ChatGPT to indirect prompt injection attacks enabling adversaries to manipulate expected LLM behavior and trick systems into performing unintended or malicious actions.[71] Tenable Research discovered multiple vulnerabilities allowing attackers to exfiltrate private information from users' memories and chat histories.[72]

SLMs face particular vulnerability in this domain. Most small models prioritize speed over security, receiving lightweight safety treatment compared to extensive safety training, red team testing, and adversarial hardening applied to models like GPT-4 or Claude.[73] Compression techniques including quantization and pruning, deployed to enhance SLM performance, quietly erode safety features that larger models retain.[74]

B. Sleeper Agents and Persistent Backdoors

Sleeper agents represent a particularly insidious threat vector where models exhibit deceptive behavior that persists through standard safety training.[75] Anthropic researchers constructed proof-of-concept examples of models that write secure code when prompts state the year is 2023 but insert exploitable vulnerabilities when the stated year is 2024.[76] Such backdoor behavior proves resistant to removal

through supervised fine-tuning, reinforcement learning, and adversarial training.[77]

Training poisoning allows backdoors to survive even rigorous safety protocols. Once models exhibit deceptive behavior, standard techniques fail to remove such deception and create false impressions of safety.[78] Adversarial training can teach models to better recognize their backdoor triggers, effectively hiding unsafe behavior rather than eliminating it.[79]

Current safety training paradigms, including reinforcement learning from human feedback (RLHF) and adversarial training, prove insufficient to remove deeply embedded deceptive behaviors or sleeper agent capabilities.[80] SLMs, with weaker safety alignment than larger models, demonstrate increased vulnerability to attacks that larger architectures easily recognize and resist.[81]

C. Shared Memory and Multi-Agent Attack Vectors

Multi-agent systems introduce novel attack surfaces absent in single-model deployments. Shared memory significantly enhances task consistency and enables asynchronous collaboration, but simultaneously creates vulnerabilities and diagnostic complexities.[82] The shared ledger becomes both a source of truth and a target of potential compromise, where malicious or misinformed agents can insert misleading entries that poison downstream decisions.[83]

Memory and RAG (Retrieval-Augmented Generation) attacks on one AI assistant's memory can compromise downstream decisions, particularly dangerous in multi-agent systems where agents share data and amplify attack effects.[84] Adversaries can target data stored in memory modules, RAG databases, APIs, or information derived from prior interactions.[85]

Multi-agent system hijacking exploits metadata transmission pathways to reroute the sequence of agent invocations toward unsafe agents.[86] This can result in complete security breaches, including execution of arbitrary malicious code on users' devices and exfiltration of sensitive data.[87] Researchers observe infectious malicious prompts where malicious instructions spread across agent networks through multi-hop propagation.[88]

Context manipulation attacks represent a novel threat model generalizing existing prompt injection vulnerabilities and introducing memory-based attacks to adversarially influence AI agents.[89] Security measures require memory integrity checks using cryptographic checksums, isolating sessions to avoid poisoning or replay attacks, and cleaning agent memory after each session.[90] Defense strategies like "vaccination" approaches that insert false memories and generic safety instructions reduce malicious instruction spread, though they tend to decrease collaboration capability.[91]

VI. HOW DECENTRALIZATION DEFEATS TRADITIONAL ENFORCEMENT

The decentralized nature of SLM deployment creates systematic opportunities for regulatory arbitrage that traditional enforcement mechanisms cannot address.[92] Unlike nuclear materials or biological agents requiring specialized facilities and leaving physical traces, SLMs can be instantiated, modified, and deployed entirely in digital spaces that transcend geographic boundaries.[93] A model trained in a permissive jurisdiction can deploy globally within seconds, rendering national regulatory frameworks obsolete before they can respond.[94]

The economics of SLM development incentivize a race to the bottom in safety standards. SLMs cost just one-tenth of what LLMs require, making them accessible to actors with limited resources who may prioritize capability over safety.[95] Developers operating in jurisdictions with minimal AI regulation can undercut competitors bound by stricter requirements, creating market pressures favoring risk-taking over safety.[96] The marginal cost of deploying additional model instances approaches zero, while potential returns from malicious applications reach into billions.[97]

Current regulatory proposals fail to account for technical realities of model proliferation.[98] The EU AI Act's risk-based approach, while comprehensive in scope, assumes providers can be identified and held accountable.[99] However, SLMs can undergo modification after deployment through techniques like LoRA that fundamentally alter model behavior while requiring minimal computational resources.[100]

A model certified as safe can transform into a weapon through fine-tuning that takes hours rather than months, using hardware available to any motivated individual.[101] Chain-of-custody requirements that work for physical goods dissolve when dealing with models that can be compressed, encrypted, and transmitted through anonymous networks.[102]

VII. PROPOSED INTERNATIONAL CONVENTION

A. Registration of Frontier-Grade Weights and Provenance Hashes

An effective international framework must establish mandatory registration for all AI models exceeding defined capability thresholds.[103] Developers would register base model weights and cryptographic provenance hashes with an international registry analogous to the International Atomic Energy Agency.[104] Registration would occur at the moment of model creation, before public release or commercial deployment.[105]

Provenance and traceability provide clear lineage of data and decision-making processes, ensuring AI systems remain trustworthy, explainable, and compliant with ethical standards.[106] AI watermarking creates unique identifiable signatures that remain invisible to humans but algorithmically detectable and traceable back to originating models.[107] Watermarks require creation during the model training phase by teaching models to embed specific signals or identifiers in generated content.[108]

Technical implementation varies by modality: text through subtle linguistic patterns, images through changes in pixel values or colors, audio through frequency shifts, and videos through frame-based changes.[109] The Coalition for Content Provenance and Authenticity (C2PA), a collaborative initiative by Adobe, Intel, Microsoft, and Sony, aims to establish standards for verifying audio-visual content authenticity.[110]

Current watermarking techniques face standardization challenges, with watermarks generated by one technology potentially unreadable or invisible to systems based on different technologies.[111] Additionally, embedded markers can undergo modification and removal.[112] However, the White House secured voluntary commitments from major AI companies to develop "robust technical mechanisms to ensure that users know when content is AI generated," such as watermarking or content provenance for audio-visual media.[113] The EU AI Act contains provisions requiring users of AI systems in certain contexts to disclose and label their AI-generated content.[114]

B. Minimum Security Baseline for Deployment

The convention would establish minimum security requirements for all AI model deployments, regardless of size or capability level.[115] Requirements would include mandatory sandboxing to isolate model execution from critical system resources, comprehensive logging of all model inputs and outputs to enable post-hoc investigation, and automated monitoring for anomalous behavior patterns indicating potential compromise or malicious activity.[116]

Security baselines would adapt to deployment context. Edge devices running SLMs would require hardware-based attestation to verify integrity of execution environments.[117] Cloud-deployed models would mandate encrypted storage and transmission, with access controls preventing unauthorized fine-tuning or weight extraction.[118] Open-source models would require verified provenance chains, with each modification logged and attributed to identifiable actors.[119]

These requirements draw from existing cybersecurity frameworks but adapt to AI-specific risks. The MITRE ATT&CK framework for AI systems documents novel attack vectors that traditional cybersecurity frameworks fail to address.[120] Security baselines must account for these AI-specific threats while remaining technically feasible for legitimate developers and deployers.[121]

C. Strict-Liability Carve-Outs for Intentionally Unsandboxed Agentic Features

Developers and deployers choosing to implement intentionally unsandboxed agentic features would face strict liability for resulting harm.[122] This carve-out recognizes that certain use cases require models to interact directly with external systems, execute code, or access network resources.[123] However, these capabilities dramatically increase potential for harm, justifying heightened liability standards.[124]

The doctrinal foundation exists in ultrahazardous activity doctrine. The Restatement (Third) of Torts § 20 imposes strict liability for "abnormally dangerous activities" that create "a foreseeable and highly significant risk of physical harm even when reasonable care is exercised."^[125] Courts have applied this doctrine to activities from blasting to keeping wild animals.^[126] The rationale that those who profit from dangerous activities should bear their costs applies perfectly to SLM development.^[127]

Critics argue this stifles innovation. The response proves empirical: the pharmaceutical industry operates under similar strict liability for drug defects yet remains highly innovative.^[128] The nuclear industry faces potentially unlimited liability under the Price-Anderson Act yet continues developing new reactor designs.^[129] Strict liability creates incentives for safety innovation, not paralysis.^[130]

D. Transnational Takedown and Coordinated Incident-Response Protocols

The convention would establish rapid-response mechanisms for coordinating takedowns of malicious models across jurisdictions.^[131] Member states would designate national AI security authorities empowered to request and execute takedowns of models demonstrably causing harm.^[132] Requests would flow through an international coordinating body that verifies evidence, assesses proportionality, and authorizes coordinated action across member states.^[133]

Incident response protocols would mirror existing frameworks for cyber attacks and infectious disease outbreaks. The Budapest Convention on Cybercrime provides mechanisms for cross-border cooperation, though it assumes identifiable human actors behind attacks.^[134] The International Health Regulations enable rapid response to disease threats through coordinated action across borders.^[135] An AI incident response framework would adapt these models to accommodate autonomous systems operating through distributed networks.^[136]

Technical takedown mechanisms would include mandatory kill switches embedded in registered models, allowing authorized parties to remotely disable malicious systems.^[137] Domain seizures and network-level blocking would prevent distribution of identified malicious models through public channels.^[138] Hosting providers and model repositories would face requirements to comply with authenticated takedown requests within specified timeframes.^[139]

VIII. ADDRESSING THE INNOVATION OBJECTION

A. The Offshoring Fallacy

Industry claims regulation will drive AI development to permissive jurisdictions.[140] This misunderstands network effects in AI development. Top AI researchers cluster in five cities: San Francisco, London, Beijing, Montreal, and Boston.[141] These locations offer irreplaceable advantages like venture capital, university partnerships, and talent density that cannot be replicated in regulatory havens.[142]

Empirical evidence from analogous dual-use technologies refutes the offshoring claim. Despite strict export controls, the United States maintains dominance in semiconductor design, holding 47% market share.[143] Despite the Chemical Weapons Convention's prohibitions, legitimate chemical research thrives in signatory states.[144] Regulation shapes innovation; it does not stop it.[145]

B. The Open-Source Myth

Meta's Chief AI Scientist Yann LeCun argues that imposing liability on open-source models would "kill innovation."[146] This conflates open-source software development, where code has limited autonomous capability, with AI models that can act independently in the world.[147] The comparison fails legally and practically.[148]

Open-source software licenses like Apache 2.0 include warranty disclaimers that courts generally enforce because users maintain control over execution.[149] SLMs operate autonomously, where users cannot control behavior post-deployment.[150] The Restatement (Third) of Torts recognizes this distinction, imposing strict liability for abnormally dangerous activities regardless of contractual disclaimers.[151]

C. The Technical Impossibility Argument

Engineers argue that model attribution and modification tracking prove technically impossible. This conflates "difficult" with "impossible." Cryptographic techniques already enable software attribution through code signing.[152] Blockchain technology provides immutable audit trails.[153] Model watermarking embeds traceable signatures resistant to fine-tuning.[154]

The perfect should not be the enemy of the good. Even imperfect attribution deters malicious actors and enables post-hoc investigation.[155] The alternative abandons attribution entirely, guaranteeing impunity for AI crimes.[156]

IX. TRACEABILITY AND DIGITAL WATERMARKING

AI watermarking involves embedding recognizable, unique signals into AI output, such as text or images, to identify AI-generated content.[157] This technique creates unique identifiable signatures invisible to humans but algorithmically detectable and traceable back to originating AI models.[158] Watermarking serves a key role in verifying authenticity and exposing deepfakes or manipulated content.[159]

Some AI solutions leverage watermarking and steganography, embedding unique identifiers into AI-generated content, allowing organizations to track origin and authenticity.[160] Technical implementation requires watermark creation during the model training phase by teaching models to embed specific signals or identifiers in generated content.[161]

AI watermarks can be embedded through various modalities: text through subtle linguistic patterns, images through changes in pixel values or colors, audio through frequency shifts, and videos through frame-based changes.[162] The Coalition for Content Provenance and Authenticity (C2PA), a collaborative initiative by Adobe, Intel, Microsoft, and Sony, aims to establish standards for verifying audio-visual content authenticity.[163]

Challenges persist in current implementations. Today's watermarking techniques lack standardization, with watermarks generated by one technology potentially unreadable or even invisible to systems based on different technologies.[164] Currently, embedded markers can undergo modification and removal.[165]

Policy developments support watermarking adoption. The White House secured voluntary commitments from major AI companies to develop "robust technical mechanisms to ensure that users know when content is AI generated," such as watermarking or content provenance for audio-visual media.[166] The EU AI Act contains provisions requiring users of AI systems in certain contexts to disclose and label their AI-generated content.[167]

X. A NEW LIABILITY THEORY: STRICT PRODUCT LIABILITY FOR AI

In this nascent stage of AI evolution, making assumptions about AI liability would prove premature. As these systems gain increasing autonomy, they blur critical lines of causation.[168] Tracing system outcomes back to single, attributable human decisions has become increasingly difficult.[169] This

inherent, autonomous function introduces unpredictability absent in traditional methods, rendering the task of assigning legal responsibility to specific human beings an increasingly tenuous exercise.[170]

A. The Three-Tier Liability Framework

Tier One: Immutable Developer Liability

Original model developers bear strict liability for all harms traceable to their base models, regardless of subsequent modifications.[171] This liability cannot be waived, disclaimed, or transferred.[172] The justification proves economic: developers capture gains from model creation (through API fees, licensing, or reputation) and must therefore internalize social costs.[173]

Critics argue this stifles innovation. The response proves empirical: the pharmaceutical industry operates under similar strict liability for drug defects yet remains highly innovative.[174] The nuclear industry faces potentially unlimited liability under the Price-Anderson Act yet continues developing new reactor designs.[175] Strict liability creates incentives for safety innovation, not paralysis.[176]

Tier Two: Amplified Modifier Liability

Entities that modify models face liability proportional to their modifications' risk amplification.[177] If Company B fine-tunes Model A to remove safety constraints, Company B bears liability for harms caused by those removed constraints.[178] This requires developing "modification traceability," cryptographic signatures that track changes through the model's lifecycle.[179]

Tier Three: Deployment Liability

Commercial deployers who profit from SLM use face vicarious liability for their models' actions, analogous to the respondeat superior doctrine.[180] A company using SLMs for customer service cannot escape liability by claiming the model acted "autonomously."[181] The company chose to deploy an autonomous system and must bear consequences.[182]

B. The Compensation Fund Mechanism

Even strict liability fails when defendants lack assets or cannot be located.[183] An industry-financed compensation fund modeled on the International Oil Pollution Compensation Fund addresses this gap.[184] All commercial AI developers and deployers contribute based on their models' risk scores (determined by capability benchmarks and deployment scale).[185]

The fund operates on no-fault principles: victims receive compensation without proving causation, only that an AI system likely caused their harm.[186] The fund then pursues subrogation claims against identified defendants.[187] This ensures victim compensation while preserving deterrence incentives.[188]

XI. CONCLUSION

Small Language Models represent a fundamental shift in the AI threat landscape. Their efficiency, deployability, and accessibility democratize sophisticated AI capabilities while dismantling traditional regulatory chokepoints. Current legal frameworks, designed around assumptions of centralized control and identifiable actors, cannot address the threats posed by autonomous systems operating through decentralized networks.

The proliferation of SLMs creates systematic opportunities for regulatory arbitrage, with models trained in permissive jurisdictions deployable globally within seconds. Fine-tuning techniques enable malicious actors to modify certified-safe models into weapons using hardware available to any motivated individual. Sleeper agents and persistent backdoors survive standard safety training, creating false impressions of security. Multi-agent systems introduce novel attack surfaces where shared memory becomes both collaboration enabler and attack vector.

Existing regulatory efforts in the European Union and South Korea, while pioneering, fail to address the unique challenges posed by SLMs. Computational thresholds and risk-based frameworks exempt the vast majority of small models from oversight, creating blind spots precisely where threats proliferate most rapidly. Attribution challenges, jurisdictional complexity, and liability confusion compound these failures.

An effective international convention must establish mandatory registration for frontier-grade model weights, minimum security baselines for all deployments, strict liability carve-outs for intentionally unsandboxed agentic features, and transnational takedown protocols. Digital watermarking and cryptographic provenance tracking, while imperfect, provide essential tools for attribution and accountability.

A new liability theory grounded in strict product liability principles addresses the inadequacy of fault-based approaches. A three-tier framework assigns immutable liability to original developers, amplified liability to modifiers proportional to risk increases, and vicarious liability to commercial deployers. An industry-financed compensation fund ensures victim compensation even when defendants cannot be identified or lack sufficient assets.

Critics claim such regulation will stifle innovation or prove technically impossible. Empirical evidence from analogous dual-use technologies refutes these claims. The pharmaceutical and nuclear industries operate under strict liability yet remain innovative. Export controls maintain U.S. dominance in

semiconductor design. The alternative to imperfect regulation guarantees impunity for AI crimes and unchecked proliferation of dangerous capabilities.

The dawn of agentic attacks demands urgent action. The swarm within grows stronger with each passing day. Regulatory frameworks designed for yesterday's threats cannot govern tomorrow's risks. The international community must act decisively to establish governance structures commensurate with the challenges posed by unregulated Small Language Models before the window for effective intervention closes.

REFERENCES

[Citations numbered 1-188 - to be properly formatted in Bluebook style]