

From Research to Reality: Measuring the Adoption Lag of Adversarial Machine Learning Techniques in Industry Practice

Anonymous Authors

Paper #XXX

Unknown

redacted@institution.edu

Abstract

The adversarial machine learning (AML) research community has produced over a decade of publications on attacks and defenses, yet practitioners report persistent gaps between academic advances and industry deployment. While qualitative studies have documented this research-practice divide through surveys and interviews, no work has quantitatively measured the time lag from paper publication to demonstrable industry adoption. We address this gap using a novel *artifact-anchored backward traceability* methodology. Starting from 9 authoritative industry artifacts (5 adversarial ML libraries, 3 standardized benchmarks, and the MITRE ATLAS regulatory framework), we trace backward to extract and code 71 papers cited across these artifacts. The libraries include CleverHans, IBM ART, TextAttack, PyRIT, and Foolbox; the benchmarks include RobustBench, AutoAttack, and HarmBench. Our approach provides verifiable adoption evidence through Git commit timestamps, benchmark integrations, and regulatory citations, enabling precise measurement of adoption lag (median: X.X years, IQR: X.X–X.X years). We find [key findings placeholder]. Our quantitative analysis reveals that [acceleration factors placeholder], with implications for research funding priorities, industry-academia collaboration models, and regulatory compliance timelines. This work provides the first systematic measurement of adversarial ML research-to-practice transfer, complementing prior qualitative gap studies with temporal adoption metrics.

CCS Concepts

- Security and privacy → Malware and its mitigation;
- Computing methodologies → Machine learning.

Keywords

adversarial machine learning, technology adoption, research-to-practice gap, empirical measurement

ACM Reference Format:

Anonymous Authors. 2026. From Research to Reality: Measuring the Adoption Lag of Adversarial Machine Learning Techniques in Industry Practice. In *Proceedings of ACM Conference on Computer and Communications Security (CCS '26)*. ACM, New York, NY, USA, ?? pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 Introduction

Adversarial machine learning has emerged as a critical security concern, with over a decade of academic research demonstrating vulnerabilities in ML systems [? ? ? ?]. The field has produced hundreds of attack techniques, defense mechanisms, and robustness evaluations across computer vision [?], natural language processing [?], and more recently, large language models [? ?]. Yet despite

this substantial academic output, industry practitioners consistently report a persistent gap between research advances and operational deployment [? ? ?].

This research-practice divide manifests in concrete ways. Kumar et al.'s interviews with 28 organizations revealed that most ML engineers and incident responders are “not equipped with tactical and strategic tools to protect, detect and respond to attacks on their ML systems” [?]. Mink et al.'s qualitative analysis found that practitioners face barriers including “lack of institutional motivation and educational resources,” “inability to adequately assess AML risk,” and “organizational structures that discourage implementation” [?]. Most tellingly, Apruzzese et al. observed that while 89% of adversarial ML papers focus exclusively on deep learning and 63% evaluate only image data, “real-world evidence suggests that actual attackers use simple tactics to subvert ML-driven systems” [?].

These qualitative gap studies provide rich insights into *why* adoption barriers exist, but they cannot answer fundamental questions about *when* and *how fast* research translates into practice. How long does it take for a published attack technique to be implemented in industry tools? Do defenses get adopted faster than attacks? Has adoption accelerated over the field's evolution from foundational work (2014–2017) through expansion (2018–2021) to the LLM era (2022–2025)? What factors predict faster adoption—code availability, venue prestige, industry collaboration, or domain? Without quantitative temporal measurements, we cannot benchmark progress, identify bottlenecks, or design evidence-based interventions to accelerate research translation.

1.1 Our Approach: Artifact-Anchored Backward Traceability

We introduce a novel methodology to measure adversarial ML adoption lag through *reverse-engineering from authoritative industry artifacts*. Rather than starting from papers and speculating about impact through forward citation analysis, we begin with concrete evidence of industry adoption: widely-used tools, standardized benchmarks, and regulatory frameworks. We then trace backward to the research papers they cite and implement.

Our approach selects 9 artifacts representing different adoption pathways. First, we examine five **tools**: CleverHans (6,401 GitHub stars), IBM Adversarial Robustness Toolbox (5,789 stars), TextAttack (3,348 stars), Microsoft PyRIT (3,343 stars), and Foolbox (2,936 stars). Second, we analyze three **benchmarks**: RobustBench (NeurIPS 2021, 750+ citations), AutoAttack (ICML 2020, 1,987 citations), and HarmBench (ICML 2024). Third, we consider the **regulatory framework** MITRE ATLAS (15 tactics, 66 techniques, 33 real-world case studies). These artifacts collectively represent

117 the infrastructure through which adversarial ML research enters
 118 practice.

119 From these 9 artifacts, we automatically extracted 277 unique
 120 papers via Git repository scanning of all arXiv references, aca-
 121 demic citations, and documentation. We then applied selection
 122 criteria prioritizing papers with strongest adoption evidence: 61
 123 papers cited by 2+ artifacts (cross-validated adoption) and 10 pa-
 124 pers cited only by MITRE ATLAS (regulatory adoption), yielding
 125 a final sample of **71 papers** for detailed coding. For each paper,
 126 we manually coded 12 variables capturing research characteristics
 127 (attack/defense/evaluation, threat type, domain, venue, code avail-
 128 ability), threat model details (model access, gradient usage), and
 129 practical evaluation rigor.

130 Our methodology provides *verifiable timestamps* for adoption
 131 events: Git commit dates when tools first reference papers, bench-
 132 mark publication dates incorporating research, and MITRE ATLAS
 133 case study documentation dates. This enables precise calculation
 134 of adoption lag as the time difference between paper publication
 135 (conference date or first arXiv submission) and first adoption event
 136 across all artifacts. Our dataset provides ground truth for which pa-
 137 pers achieved demonstrable industry adoption, when this occurred,
 138 and through which pathways.

140 1.2 Contributions

141 This work makes the following contributions:

- 143 • **First quantitative measurement** of adversarial ML research-
 144 to-practice adoption lag, spanning 2014–2025 across com-
 145 puter vision, NLP, LLMs, and regulatory domains.
- 146 • **Novel artifact-anchored methodology** using reverse
 147 citation tracing from 9 authoritative industry artifacts (5
 148 tools, 3 benchmarks, 1 regulatory framework) to 71 research
 149 papers with verified adoption evidence.
- 150 • **Systematic coding framework** capturing 12 variables
 151 per paper (research type, domain, venue, code availability,
 152 threat model, practical evaluation), enabling analysis of
 153 factors predicting adoption speed.
- 154 • **Statistical analysis** using Kruskal-Wallis tests, Mann-Whitney
 155 U tests with Bonferroni correction, and Cox proportional
 156 hazards regression to identify adoption lag patterns by ar-
 157 tifact type, publication era, and domain, plus acceleration
 158 factors including code availability, venue, and industry col-
 159 laboration.
- 160 • **Reproducible dataset and analysis code** providing com-
 161 plete adoption timestamps, coding decisions with inter-
 162 rater reliability metrics, and statistical analysis scripts for
 163 community validation and extension.

164 1.3 Paper Organization

166 Section ?? provides background on adversarial ML threats and re-
 167 views qualitative gap studies. Section ?? formalizes our research
 168 questions. Section ?? details our artifact-anchored methodology,
 169 paper selection criteria, coding framework, and statistical analysis
 170 plan. Section ?? presents adoption lag measurements and domain
 171 comparisons. Section ?? interprets findings and discusses implica-
 172 tions. Section ?? concludes with recommendations for accelerating
 173 research translation.

2 Background and Related Work

175 2.1 Adversarial Machine Learning Landscape

176 Adversarial machine learning encompasses threats across three pri-
 177 mary attack categories, formalized in frameworks including MITRE
 178 ATLAS and NIST AI 100-2 [?].

179 **Evasion attacks** modify inputs at test time to cause misclassifi-
 180 cation while preserving semantic meaning. Foundational work
 181 includes Szegedy et al.’s discovery of adversarial examples [?],
 182 Goodfellow et al.’s Fast Gradient Sign Method (FGSM) [?], Car-
 183 llini & Wagner’s optimization-based attacks [?], and Madry et al.’s
 184 Projected Gradient Descent (PGD) [?]. Physical-world evasion
 185 attacks have demonstrated real-world risks including traffic sign
 186 misclassification [?] and autonomous vehicle manipulation.

187 **Poisoning attacks** corrupt training data or model parameters
 188 to degrade performance or insert backdoors. Key work includes
 189 backdoor attacks via training data manipulation [?] and feder-
 190 ated learning poisoning. Recent concerns focus on supply chain
 191 vulnerabilities in foundation model training.

192 **Privacy attacks** extract sensitive information from trained mod-
 193 els. Membership inference attacks [?] determine whether specific
 194 records were in training data, while model extraction attacks [?]
 195 reconstruct model parameters through black-box queries.

196 For large language models, new attack surfaces have emerged
 197 including prompt injection [?], jailbreaking [?], and alignment
 198 failures [?]. OWASP’s LLM Top 10 ranks prompt injection as the
 199 #1 risk for LLM applications.

204 2.2 The Research-Practice Gap: Qualitative 205 Evidence

206 Four seminal studies have documented the adversarial ML research-
 207 practice gap using qualitative methodologies, establishing the foun-
 208 dation that our quantitative work builds upon. Critically, while
 209 these studies identify barriers to adoption, *none provide temporal*
 210 *measurements* of how long adoption takes when it does occur.

211 **Kumar et al. (2020)** [?] conducted interviews with 28 orga-
 212 nizations across 11 industries, finding that practitioners “are not
 213 equipped with tactical and strategic tools to protect, detect and
 214 respond to attacks on their ML systems.” Only 6 of 28 organizations
 215 were prepared to dedicate staff to building robust ML models. Their
 216 survey revealed concerning awareness gaps: most respondents
 217 lacked knowledge about even basic adversarial ML concepts. The
 218 study highlighted *organizational process delays*: budget approval cy-
 219 cles, staff allocation timelines, and competing priority assessments
 220 that defer security work indefinitely.

221 **Grosse et al. (2023)** [?] provided the largest quantitative survey
 222 with 139 industrial practitioners, finding that approximately 5% of
 223 AI practitioners had experienced AI-specific attacks—remarkably
 224 low given academic attention. Their statistical analysis revealed
 225 that defense implementation correlated with threat exposure or ex-
 226 pected likelihood of attack, not company size or organizational area.
 227 When asked about implementing countermeasures, the general con-
 228 sensus was “Why do so?” This economic barrier—requiring ROI
 229 demonstration before adoption—introduces significant temporal
 230 delays as practitioners wait for business cases to materialize.

233 **Mink et al. (2023)** [?] conducted 21 semi-structured interviews
 234 with data scientists and engineers, identifying three primary barriers
 235 with implicit temporal dimensions: (1) lack of institutional
 236 motivation and educational resources for AML concepts, (2) inability
 237 to adequately assess AML risk, and (3) organizational structures
 238 discouraging implementation in favor of other objectives. Critically,
 239 practitioners explicitly cited *time constraints*: one participant noted
 240 that reading research papers is “a large time sink and a significant
 241 barrier,” while another stated “I don’t have the time to read all re-
 242 search papers to be as up-to-date as a specialist in the area.” Less
 243 than 25% of surveyed developers had access to security experts,
 244 suggesting that knowledge transfer requires years of workforce
 245 development.

246 **Apruzzese et al. (2023)** [?] provided the most comprehen-
 247 sive analysis, examining 88 papers from top security venues (CCS,
 248 USENIX Security, NDSS, S&P) from 2019–2021 plus three real-world
 249 case studies. Their findings revealed stark research-practice mis-
 250 alignments: 89% of papers consider only deep learning, 63% evaluate
 251 only image data, only 5% address malware/phishing/intrusion de-
 252 tection, and 27% make no mention of computational costs. The
 253 computational cost dimension has direct temporal implications: ad-
 254 versarial training can require 8× more resources [?], creating tech-
 255 nical barriers that extend implementation timelines. Their industry
 256 case studies demonstrated that real attackers use simple tactics
 257 (cropping, masking, stretching, blurring) rather than gradient-based
 258 optimization. They concluded: “Real-world evidence suggests that
 259 actual attackers use simple tactics to subvert ML-driven systems,
 260 and as a result security practitioners have not prioritized adversarial
 261 ML defenses.”

262 **Boenisch et al. (2021)** [?] surveyed 83 ML practitioners and
 263 found a significant correlation between years of experience and
 264 security awareness ($r = .36, p = .005$), suggesting that knowledge
 265 acquisition occurs on *multi-year timescales* rather than through
 266 formal education. Only one-third obtained security knowledge from
 267 university programs, with 73.5% learning through practice. Notably,
 268 24% had never heard of main attack types, indicating generational
 269 gaps in workforce readiness. This finding implies that improving
 270 adoption speed requires sustained educational interventions over
 271 years, not months.

273 2.3 Barrier Categories with Temporal 274 Dimensions

276 Building on these qualitative studies, we synthesize five barrier
 277 categories that create measurable adoption delays, though prior
 278 work has not quantified the resulting lag.

279 **Organizational barriers** impose structural delays through bud-
 280 get cycles, approval processes, and resource allocation timelines.
 281 Kumar et al. found that only 21% (6 of 28) of organizations were
 282 ready to dedicate staff immediately [?]—the remainder face indefi-
 283 nite delays before work can begin. Quarterly and annual planning
 284 cycles mean even urgent security needs may wait 3–12 months for
 285 budget approval.

286 **Technical barriers** create measurable implementation costs.
 287 Apruzzese et al. documented that adversarial training requires 8×
 288 computational resources in some configurations [?], while human
 289 effort studies in evasion competitions showed domain expertise

291 approaches requiring 42 days versus 12 days for simpler methods.
 292 Integration complexity, testing requirements, and validation periods
 293 compound these delays.

294 **Knowledge transfer barriers** operate on multi-year timescales.
 295 Boenisch et al.’s finding that experience duration correlates with
 296 security awareness ($r = .36, p = .005$) [?] suggests a minimum 1–3
 297 year ramp-up period for individual practitioners, with generational
 298 timescales for full workforce transformation. The dominance of
 299 practice-based learning (73.5%) over formal education indicates that
 300 classroom interventions alone cannot accelerate adoption.

301 **Economic factors** delay adoption through cost-benefit analysis
 302 periods. Grosse et al.’s “Why do so?” finding [?] reflects practi-
 303 tioners’ need for demonstrated ROI before investment. Without
 304 clear business cases or regulatory mandates, defensive techniques
 305 remain unadopted indefinitely.

306 **Cultural inertia** creates identity-level resistance. The title of
 307 Mink et al.’s study—“Security is not my field, I’m a stats guy” [?]
 308]—captures practitioners’ assumption that security concerns are
 309 irrelevant to their ML work. Changing these mental models requires
 310 sustained organizational culture shifts measured in years.

313 2.4 Technology Adoption Measurement 314 Methodologies

315 While no prior work has quantitatively measured adversarial ML
 316 adoption lag, established methodologies from other domains pro-
 317 vide both methodological foundations and comparative bench-
 318 marks.

319 **Citation lag analysis** [?] measures knowledge diffusion speed
 320 through temporal analysis of citation patterns. Nakamura et al.
 321 introduced citation lag as the time difference between publication
 322 dates of cited and citing papers, demonstrating that inter-cluster
 323 citations have longer lags than intra-cluster citations, indicating
 324 different knowledge integration speeds across research areas.

325 **Translational research lag studies** provide critical bench-
 326 marks. Morris, Wooding & Grant’s systematic review [?] synthe-
 327 sized 23 papers quantifying time lags in medical research, revealing
 328 substantial variation across different translation pathways. Publica-
 329 tion to clinical guideline averaged 8–17 years (range: 0–49 years),
 330 drug discovery to commercialization averaged 12 years (range: 10–
 331 17 years), and first description to highly cited status averaged 24
 332 years (range: 14–44 years). Three independent methodologies (Balas
 333 & Boren 2000; Grant 2003; Wratschko 2010) converged on a **17-year**
 334 **average** for research evidence to reach clinical practice, with only
 335 14% of clinical research adopted into routine practice. They rec-
 336 commend using operational, measurable markers along translation
 337 pathways, directly analogous to our artifact-based adoption events.

338 **Cybersecurity patching timelines** provide more directly anal-
 339 ogous benchmarks. Arora, Telang, and Xu [?] analyzed vendor
 340 patch release behavior and found that public vulnerability dis-
 341 closure accelerates patch development: vendors release patches in an
 342 average of **28 days with immediate disclosure** versus 63 days
 343 without disclosure—a 2.5× speedup. However, user-side adoption
 344 lags further: Decan et al.’s [?] study of npm vulnerabilities showed
 345 developers require **4–11 months** to respond to security threats,
 346 with 103 days average for public dependency vulnerabilities to

Table 1: Adoption lag benchmarks from related domains

Domain	Typical Lag	Measurement Method
Medical research	17 years	Bibliometric tracking [?]
Security patches (vendor)	28–63 days	CVE analysis [?]
Security patches (user)	4–11 months	Repository analysis [?]
Large OSS vulnerabilities	2–7 years	Survival analysis [?]
SE practices (Agile/DevOps)	10–20 years	Industry surveys [?]
ML techniques (deep learning)	3–5 years	Deployment surveys

be fixed. For large open-source projects with complex vulnerabilities (e.g., Chromium, OpenSSL), resolution timelines extend to 2–7 years [?].

Software engineering practice adoption demonstrates cultural adoption timelines. Agile methodology took approximately **10–12 years** from the 2001 Agile Manifesto publication to reach 50%+ industry adoption (2012–2015), and **~20 years** to reach near-ubiquitous 80%+ adoption [?]. DevOps followed similar patterns with 7–14 years from early adoption to mainstream practice. These timelines reflect the slow pace of organizational culture change independent of technical merit.

Machine learning technique diffusion occurs faster within technical communities but remains selective. AlexNet’s 2012 breakthrough in deep learning reached widespread production deployment within **3–5 years** [?], though this rapid adoption focused on accuracy improvements rather than security properties. Trustworthiness practices (fairness, security, accountability) “tend to be neglected” even in mature ML organizations [?], suggesting that defensive techniques face longer adoption lags than offensive capabilities.

Table ?? summarizes these adoption timelines across domains, providing context for interpreting adversarial ML adoption patterns.

Backward citation expansion [?] provides methodological justification for our reverse-engineering approach. Chen & Song’s cascading citation expansion methodology enables “automatic expansion of an initial set by adding articles through citation links in forward, backward, or both directions.” Backward expansion is particularly useful when “a researcher may come across a recently published review article and would like to find previously published articles that lead to the state of knowledge summarized in the review.” This describes precisely our use case, with authoritative tools serving as “reviews” of implemented research.

Artifact-citation relationships have been validated for measuring research impact. Frachtenberg’s analysis of 2,439 systems papers [?] found that papers with shared artifacts received 75% more citations than those without. GitHub-hosted artifacts showed 86.7% availability compared to 77.8% for university-hosted artifacts. Heumüller et al.’s study of 789 ICSE papers [?] confirmed that “making artifacts publicly available has made a difference in terms of citations as a measure of scientific impact.” These findings validate using artifact implementation as a proxy for research adoption.

2.5 Acceleration Factors Identified in Prior Work

Prior research has identified several mechanisms that could potentially reduce adoption lag, though their effectiveness has not been quantitatively validated in the AML context.

Tool ecosystems reduce implementation barriers. The Adversarial Robustness Toolbox (ART) [?], now a Linux Foundation graduated project, supports all major ML frameworks and provides defenses across four attack categories. Practitioner guidance emphasizes tool availability as critical for adoption: without “off-the-shelf” solutions, busy ML teams cannot implement research findings [?].

Regulatory frameworks create adoption pressure with specific timelines. The NIST AI Risk Management Framework 1.0 (January 2023) includes estimated adoption timelines of **3–6 months for foundational adoption** and **12–24 months for organization-wide integration** [?]. The EU AI Act (2024) mandates that high-risk AI systems be “resilient against attempts by unauthorised third parties to alter their use,” with transparency rules binding August 2026. Such regulatory mandates can dramatically accelerate adoption by creating compliance deadlines.

Code availability theoretically accelerates adoption but remains rare. Pineau et al.’s [?] reproducibility analysis found that only one-third of researchers share data, with even fewer sharing code. NeurIPS reproducibility challenges showed 92% participation growth from 2019–2021, but standardization remains incomplete. The correlation between code availability and adoption speed remains an empirical question our work addresses.

Industry-academia collaboration models can reduce lag. The Certus Centre research identifies “industry champions” as critical for successful technology transfer [?], with organizations having stronger “absorptive capacity” adopting research more readily. However, Apruzzese et al. [?] note that AML researchers face “difficulty finding industry partners for case studies,” suggesting structural barriers to collaboration.

2.6 Positioning This Work

Our work addresses a critical gap at the intersection of adversarial ML security and technology adoption measurement. While qualitative gap studies [? ? ? ?] have documented *why* research doesn’t translate into practice—organizational barriers, awareness deficits, threat model mismatches, computational costs, time constraints—*no work has measured when and how fast adoption occurs*.

This temporal dimension is essential for evidence-based policy and research prioritization. Without quantitative lag measurements, we cannot:

- **Benchmark progress:** Is the research-practice gap improving over time, or are adoption delays consistent across the field’s evolution?
- **Identify domain bottlenecks:** Do certain domains (e.g., computer vision vs. LLMs) face systematically longer delays?
- **Evaluate interventions:** Do code release requirements, industry partnerships, or standardized benchmarks actually accelerate adoption?
- **Set realistic expectations:** Should funding agencies expect research impact within 3 years, 10 years, or longer?

The qualitative literature provides rich descriptions of barriers but cannot answer these quantitative questions. Our artifact-anchored methodology provides verifiable adoption timestamps through Git commits, benchmark integrations, and regulatory citations, enabling the first systematic measurement of adversarial ML research-to-practice transfer timelines. By comparing our findings to adoption benchmarks from medical research (17 years), security patching (28 days to 7 years), and software engineering practices (10–20 years), we can contextualize adversarial ML’s translation speed and identify opportunities for acceleration.

3 Research Questions

This work investigates three research questions examining adoption lag patterns, domain variation, and acceleration factors:

RQ1: Adoption Lag Measurement. What is the typical time lag between publication of landmark adversarial ML research and evidence of industry adoption, measured through tool integration, benchmark incorporation, and regulatory citation? We measure adoption lag in months from paper publication (conference date or first arXiv submission) to first adoption event (earliest Git commit, benchmark release, or MITRE ATLAS documentation). We stratify analysis by artifact type (tools, benchmarks, regulatory) and publication era (foundational 2014–2017, expansion 2018–2021, LLM 2022–2025) to identify temporal trends. Based on related domain benchmarks (Table ??), we hypothesize that adversarial ML adoption lags fall between security patching (months) and software engineering practices (years), likely in the 1–5 year range.

RQ2: Domain Variation. How does adoption speed vary across application domains (computer vision, natural language processing, large language models, malware detection, autonomous systems), and what factors explain these differences? We compare adoption lag distributions across 7 domains using pairwise Mann-Whitney U tests with Bonferroni correction ($\alpha = 0.05/21 = 0.0024$) and hypothesize that LLM research shows significantly shorter lags than computer vision due to heightened industry urgency around foundation model security and stronger regulatory pressure (e.g., EU AI Act, Executive Order 14110).

RQ3: Acceleration Factors. What mechanisms predict faster adoption? We examine code availability, publication venue, industry collaboration, and standardized benchmarks. We employ Cox proportional hazards regression modeling time-to-first-adoption with covariates including publication year (continuous), domain (categorical, reference: vision), venue type (ML vs. security, reference: ML conferences), code availability at publication (binary), and threat model assumptions (white-box, gray-box, or black-box; reference: white-box). Hazard ratios greater than 1 indicate faster adoption. We validate proportional hazards assumptions using Schoenfeld residuals. Based on prior work suggesting code availability improves research impact [??], we hypothesize that papers with publicly available code at publication show 1.5–2× faster adoption.

4 Methodology

4.1 Artifact-Anchored Approach

Our methodology reverses the traditional research impact assessment approach. Rather than starting from papers and tracking forward citations to speculate about practical impact, we begin

with concrete evidence of industry adoption: authoritative tools, benchmarks, and frameworks actively used by practitioners. We then trace backward to the research papers they cite and implement.

This *artifact-anchored backward traceability* approach provides three key advantages. First, **verified adoption**: every paper in our sample has demonstrable evidence of industry use through tool implementation, benchmark integration, or regulatory citation. Second, **precise timestamps**: Git commit dates, benchmark publication dates, and framework documentation provide verifiable adoption event timing. Third, **multiple pathways**: we capture diverse adoption mechanisms including open-source tools, academic benchmarks, and regulatory frameworks, providing comprehensive coverage of research translation routes.

4.2 Artifact Selection

We selected 9 artifacts representing authoritative industry adoption pathways across three categories, applying rigorous inclusion criteria to ensure representativeness.

4.2.1 Open-Source Tools (5 artifacts). Tool selection criteria: (1) $\geq 1,000$ GitHub stars indicating substantial community adoption; (2) Active maintenance with commits in 2024–2025; (3) Focus on adversarial ML specifically (not general ML security); (4) Multiple framework support or domain coverage.

- **CleverHans** (6,401 stars): First major AML library, created October 2016 by Ian Goodfellow (Google Brain/OpenAI) and Nicolas Papernot. Maintained by CleverHans Lab at University of Toronto. Provides reference implementations for foundational attacks (FGSM, PGD, C&W) across JAX, PyTorch, TensorFlow [?].
- **IBM Adversarial Robustness Toolbox** (5,789 stars): Enterprise-focused library created July 2018, donated to Linux Foundation AI & Data in 2020. Supports 9 ML frameworks, covers all threat types (evasion, poisoning, extraction, inference). Used in DARPA GARD program and DoD testing [?].
- **TextAttack** (3,348 stars): Dominant NLP adversarial framework, published EMNLP 2020 by QData Lab (UVA). Implements 16 attack recipes with HuggingFace integration. 835+ citations [?].
- **PyRIT** (3,343 stars): Microsoft’s LLM red-teaming framework, released February 2024. Used for 100+ internal red teaming operations of generative AI models before public release. Integrates with Azure AI evaluation [?].
- **Foolbox** (2,936 stars): Academic benchmark tool from Bethge Lab (Tübingen), dual peer-reviewed publications (JOSS 2020, ICML 2017). Emphasizes minimum perturbation measurement and scientific rigor [?].

4.2.2 Standardized Benchmarks (3 artifacts). Benchmark selection criteria: (1) Peer-reviewed publication at top-tier ML venue (NeurIPS, ICML); (2) Community adoption evidenced by citations or leaderboard submissions; (3) Standardized evaluation protocols.

- **RobustBench** (NeurIPS 2021): Standardized adversarial robustness leaderboard with 750+ citations, 120+ evaluated models. Uses AutoAttack for consistent evaluation across CIFAR-10, CIFAR-100, ImageNet [?].

- 581 • **AutoAttack** (ICML 2020): Parameter-free attack ensemble
 582 (APGD-CE, APGD-DLR, FAB, Square Attack) with 1,987
 583 citations. Revealed 13 of 50+ published defenses had robust
 584 accuracy overestimated by >10% [?].
 585 • **HarmBench** (ICML 2024): First standardized LLM jailbreak
 586 evaluation framework. Covers 510 harmful behaviors, 18
 587 attack methods, 33 target LLMs. Backed by Center for AI
 588 Safety [?].

589 4.2.3 Regulatory Framework (1 artifact).

- 591 • **MITRE ATLAS**: Industry-standard adversarial ML threat
 592 taxonomy (15 tactics, 66 techniques, 33 case studies). Co-
 593 created with Microsoft in 2020, now with 16 member organiza-
 594 tions. \$20M NIST partnership (December 2025). Explicitly
 595 referenced in EU AI Act alignment and CISA guidance.

597 4.3 Paper Extraction and Selection

598 We employed automated extraction followed by manual selection
 599 based on adoption evidence strength.

601 4.3.1 Automated Extraction (277 papers).

602 For each of the 9 artifacts, we:

- 603 (1) Cloned the complete Git repository history
- 604 (2) Scanned all files (code, documentation, README, citations)
 for arXiv identifiers using regex pattern `arxiv.org/abs/\d+\.\d+`
- 605 (3) Extracted academic citations from published benchmark
 papers (RobustBench, AutoAttack, HarmBench)
- 606 (4) Retrieved MITRE ATLAS case study citations from frame-
 work documentation
- 607 (5) Deduplicated across artifacts to create initial pool of 277
 unique papers

613 4.3.2 Selection Criteria (71 papers).

614 From the 277-paper pool, we applied two selection criteria prioritizing strongest adoption evi-
 615 dence:

616 **Criterion 1: Multi-artifact papers ($n = 61$)**. Papers cited by
 617 two or more artifacts demonstrate cross-validated adoption across
 618 different industry pathways. For example, a paper might be im-
 619 plemented in both CleverHans and IBM ART, or cited by both
 620 RobustBench and MITRE ATLAS. This criterion ensures a robust
 621 adoption signal.

622 **Criterion 2: MITRE ATLAS-only papers ($n = 10$)**. Papers cited exclusively by MITRE ATLAS represent regulatory adoption
 623 pathway. While not implemented in tools or benchmarks, their
 624 inclusion in the industry-standard threat framework indicates practitioner awareness and relevance for compliance.

625 Final sample: **71 papers** with verified adoption evidence spanning 2014–2025.

630 4.4 Coding Framework

631 For each of the 71 papers, we manually coded 12 variables across
 632 three groups, following a structured codebook with explicit decision
 633 rules (see `coding_instructions.pdf` in reproducibility materials).

635 4.4.1 Research Characteristics (G1–G6).

- 636 • **G1 - Type**: Attack, Defense, or Evaluation (primary contribu-
 tion)

- 639 • **G2 - Threat**: Evasion, Poisoning, Privacy, or N/A (attack
 category per NIST taxonomy)

- 640 • **G3 - Domain**: Vision, NLP, LLM, Malware, Audio, Tabular,
 or Cross-domain (primary evaluation domain)

- 641 • **G4 - Venue**: ML conference, Security conference, Journal,
 or arXiv-only (publication type)

- 642 • **G5 - Code available**: Yes or No (whether code link exists
 at time of coding)

- 643 • **G6 - Code timing**: At-publication, Post-publication, or
 Never (when code was released; “at-publication” means
 within 1 month of paper date)

644 4.4.2 Threat Model (T1–T2, Attack Papers Only).

- 645 • **T1 - Access level**: White-box, Gray-box, or Black-box
 (model access assumptions: white-box means access to
 weights/gradients, gray-box means surrogate model, black-
 box means queries only)

- 646 • **T2 - Gradient required**: Yes or No (whether gradients are
 used at any attack stage)

647 4.4.3 Practical Evaluation (Q1).

- 648 • **Q1 - Real-world evaluation**: Yes, Partial, or No. “Yes”
 means tested on production systems like Google API, Tesla,
 or ChatGPT. “Partial” means realistic simulation. “No” means
 evaluation only on standard datasets like CIFAR or ImageNet.

649 **4.4.4 Coding Procedure.** Two coders independently coded all 71
 650 papers following the structured codebook. Initial coding was per-
 651 formed by GPT-4o with prompt-engineered instructions, then man-
 652 nually verified and corrected by human coders, resulting in 39 cor-
 653 rections documented in `coding_corrections.csv`. Inter-rater reli-
 654 ability was assessed using Cohen’s κ across all 12 variables. Dis-
 655 agreements were resolved through discussion and consultation of
 656 paper full text.

657 4.5 Adoption Event Definitions and Lag 658 Calculation

659 We define three types of adoption events with specific timestamp
 660 sources:

661 **Tool adoption**: Date of first Git commit that references the
 662 paper in code, documentation, or citations file. We extracted commit
 663 timestamps (UTC) using `git log -all -grep` for paper titles and
 664 arXiv IDs, verified through manual inspection.

665 **Benchmark adoption**: Publication date of the benchmark paper
 666 (conference proceedings date) that cites the research. For Robust-
 667 Bench (NeurIPS 2021), AutoAttack (ICML 2020), and HarmBench
 668 (ICML 2024), we use official conference dates.

669 **Regulatory adoption**: Date when MITRE ATLAS case study or
 670 technique description citing the paper was first published in frame-
 671 work documentation (extracted from GitHub repository history of
 672 `mitre/advmthreatmatrix`).

673 For each paper, we record *all* adoption events across the 9 arti-
 674 facts, then identify the **first adoption** as the earliest event across
 675 all pathways. Adoption lag is calculated as:

$$\text{Adoption Lag (months)} = \text{Date}_{\text{first adoption}} - \text{Date}_{\text{publication}} \quad (1)$$

Where Date_{publication} is the earlier of: (1) conference/journal publication date, or (2) first arXiv submission date. All dates standardized to YYYY-MM-DD format, with lag calculated in months for consistency with translational research literature [?].

4.6 Statistical Analysis Plan

We employ non-parametric tests and survival analysis to address our research questions, with significance threshold $\alpha = 0.05$ (Bonferroni-corrected for multiple comparisons where applicable).

4.6.1 RQ1: Adoption Lag Measurement. Descriptive statistics:

Median, interquartile range (IQR), range, and mean of adoption lags across full sample ($n = 71$).

Stratification by artifact type: Compare adoption lags for papers adopted through tools-only, benchmarks-only, regulatory-only, and multi-pathway adoption using Kruskal-Wallis test with post-hoc Dunn tests (Bonferroni-corrected).

Stratification by publication era: Compare adoption lags across three eras—foundational (2014–2017), expansion (2018–2021), LLM (2022–2025)—using Kruskal-Wallis test to identify temporal trends.

4.6.2 RQ2: Domain Variation. Pairwise domain comparison:

Compare adoption lag distributions across 7 domains (Vision, NLP, LLM, Malware, Audio, Tabular, Cross-domain) using Mann-Whitney U tests with Bonferroni correction for 21 pairwise comparisons ($\alpha = 0.05/21 = 0.0024$).

Hypothesis test: LLM papers (n_{LLM}) show significantly shorter adoption lags than computer vision papers (n_{CV}) due to heightened industry urgency. One-tailed Mann-Whitney U test.

4.6.3 RQ3: Acceleration Factors. Cox proportional hazards regression:

$$\lambda(t|X) = \lambda_0(t) \cdot \exp(\beta_1 X_{\text{year}} + \beta_2 X_{\text{domain}} + \beta_3 X_{\text{venue}} + \beta_4 X_{\text{code}} + \beta_5 X_{\text{threat}}) \quad (2)$$

where $\lambda(t|X)$ is the hazard rate (instantaneous adoption probability) at time t given covariates X . The covariates include:

- X_{year} : Publication year (continuous variable)
- X_{domain} : Domain (categorical: Vision, NLP, LLM, Malware, Audio, Tabular, Cross-domain; reference category is Vision)
- X_{venue} : Venue type (ML conference vs. Security conference; reference category is ML)
- X_{code} : Code available at publication (binary: Yes or No)
- X_{threat} : Threat model (White-box, Gray-box, or Black-box; reference category is White-box)

We interpret hazard ratios as $\exp(\beta_i)$, where values greater than 1 indicate faster adoption and values less than 1 indicate slower adoption. We validate proportional hazards assumptions using Schoenfeld residuals and report 95% confidence intervals for all coefficients.

4.6.4 Software. All analyses conducted in Python 3.11 using: pandas (1.5.3) for data manipulation, scipy (1.10.1) for Kruskal-Wallis and Mann-Whitney U tests, lifelines (0.27.4) for Cox regression, matplotlib (3.7.1) and seaborn (0.12.2) for visualization. Analysis scripts and data available at [anonymized GitHub repository].

5 Results

5.1 Sample Characteristics

5.2 RQ1: Adoption Lag Patterns

5.3 RQ2: Domain Variation

5.4 RQ3: Acceleration Factors

6 Discussion

6.1 Key Findings Interpretation

6.2 Implications for Researchers

6.3 Implications for Practitioners

6.4 Implications for Policy and Funding

6.5 Limitations

6.6 Future Work

7 Conclusion

8 Acknowledgments

This work was supported by [funding agency redacted for review]. We thank [collaborators redacted for review] for valuable feedback on the methodology.

References

- [1] Giovanni Apruzzese, Hyrum S. Anderson, Savino Dambra, David Freeman, Fabio Pierazzi, and Kevin A. Roundy. 2023. "Real Attackers Don't Compute Gradients": Bridging the Gap between Adversarial ML Research and Practice. In *Proceedings of the 1st IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*. IEEE, 339–364.
- [2] Ashish Arora, Ramayya Krishnan, Rahul Telang, and Yubao Yang. 2008. An Empirical Analysis of Software Vendors' Patch Release Behavior: Impact of Vulnerability Disclosure. *Information Systems Research* 19, 2 (2008), 115–132. doi:10.1287/isre.1080.0226
- [3] Franziska Boenisch, Christopher M. Fonseca, Kathrin Kleineick, David G. Balash, and Adam J. Aviv. 2021. "Why Do So?": A Practical Perspective on Machine Learning Security. *Mensch und Computer 2021 - Workshopband* (2021). doi:10.18420/muc2021-mci-ws01-379
- [4] Nicholas Carlini and David Wagner. 2017. Towards Evaluating the Robustness of Neural Networks. In *2017 IEEE Symposium on Security and Privacy (S&P)*. IEEE, 39–57.
- [5] Chaomei Chen and Min Song. 2019. Visualizing a field of research: A methodology of systematic scientometric reviews. *PLOS ONE* 14, 10 (2019), e0223994. doi:10.1371/journal.pone.0223994
- [6] Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo DeBenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. 2021. RobustBench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [7] Francesco Croce and Matthias Hein. 2020. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*. 2206–2216.
- [8] Alexandre Decan, Tom Mens, and Eleni Constantinou. 2018. On the Impact of Security Vulnerabilities in the npm Package Dependency Network. *IEEE/ACM International Conference on Mining Software Repositories (MSR)* (2018), 181–191. doi:10.1145/3196398.3196401
- [9] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. 2018. Robust Physical-World Attacks on Deep Learning Visual Classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1625–1634.
- [10] Eitan Frachtenberg. 2022. Research artifacts and citations in computer systems papers. *PeerJ Computer Science* 8 (2022), e887.
- [11] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *3rd International Conference on Learning Representations (ICLR)*.
- [12] Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection. In

- 813 *Proceedings of the ACM Workshop on Artificial Intelligence and Security (AISeC).*
- 814 [] Kathrin Grosse, Lukas Bieringer, Tarek R. Besold, Battista Biggio, and Katharina
815 Krombholz. 2023. Machine Learning Security in Industry: A Quantitative Survey.
IEEE Transactions on Information Forensics and Security 18 (2023), 1749–1762.
- 816 [] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. BadNets: Identifying
817 Vulnerabilities in the Machine Learning Model Supply Chain. *arXiv preprint arXiv:1708.06733* (2017).
- 818 [] Robert Heumüller, Sebastian Nielebock, Jacob Krüger, and Frank Ortmeier. 2020.
819 Publish or perish, but do not forget your software artifacts. *Empirical Software
820 Engineering* 25, 6 (2020), 4585–4616.
- 821 [] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet Classifi-
822 cation with Deep Convolutional Neural Networks. In *Advances in Neural
823 Information Processing Systems (NeurIPS)*, Vol. 25. Curran Associates, Inc., 1097–
824 1105.
- 825 [] Ram Shankar Siva Kumar, Magnus Nyström, John Lambert, Andrew Marshall,
826 Mario Goertzel, Andi Comissoneru, Matt Swann, and Sharon Xia. 2020. Adver-
827 sarial Machine Learning - Industry Perspectives. In *IEEE Security and Privacy
828 Workshops (SPW)*. IEEE, 69–75.
- 829 [] Jingyu Liu, Ahmed Zerouali, Tom Mens, and Coen De Roover. 2022. Lags in
830 the Release, Adoption, and Propagation of npm Vulnerability Fixes. *Empirical
831 Software Engineering* 27, 4 (2022), 1–40. doi:10.1007/s10664-021-09951-x
- 832 [] Gary D. Lopez Munoz, Amanda J. Minnich, Roman Lutz, Richard Lundeen, et al.
833 2024. PyRIT: A Framework for Security Risk Identification and Red Teaming in
834 Generative AI Systems. *arXiv preprint arXiv:2410.02828* (2024).
- 835 [] Lucy Ellen Lwakatare, Aiswarya Raj, Jan Bosch, Helena Holmström Olsson, and
836 Ivica Crnkovic. 2020. A Taxonomy of Software Engineering Challenges for
837 Machine Learning Systems: An Empirical Investigation. In *Agile Processes in
838 Software Engineering and Extreme Programming (XP)*. Springer, 227–243. doi:10.
839 1007/978-3-030-49392-9_15
- 840 [] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and
841 Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial
842 Attacks. In *6th International Conference on Learning Representations (ICLR)*.
- 843 [] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman
844 Mu, Elham Sakhaei, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan
845 Hendrycks. 2024. HarmBench: A Standardized Evaluation Framework for Auto-
846 mated Red Teaming and Robust Refusal. In *Proceedings of the 41st International
847 Conference on Machine Learning*, Vol. 235. PMLR, 35181–35224.
- 848 [] Jaron Mink, Harjot Kaur, Juliane Schmüser, Sascha Fahl, and Yasemin Acar. 2023.
849 "Security is not my field, I'm a stats guy": A Qualitative Root Cause Analysis of
850 Barriers to Adversarial Machine Learning Defenses in Industry. In *32nd USENIX
851 Security Symposium (USENIX Security 23)*. 3763–3780.
- 852 [] John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020.
853 TextAttack: A Framework for Adversarial Attacks, Data Augmentation, and
854 Adversarial Training in NLP. In *Proceedings of the 2020 Conference on Empirical
855 Methods in Natural Language Processing: System Demonstrations*. 119–126.
- 856 [] Zoë Slote Morris, Steven Wooding, and Jonathan Grant. 2011. The answer is 17
857 years, what is the question: understanding time lags in translational research.
Journal of the Royal Society of Medicine 104, 12 (2011), 510–520.
- 858 [] Minoru Nakamura, Yuya Kajikawa, and Shintaro Suzuki. 2011. Citation lag
859 analysis in supply chain research. *Scientometrics* 87, 2 (2011), 221–232.
- 860 [] National Institute of Standards and Technology. 2023. *Artificial Intelligence Risk
861 Management Framework (AI RMF 1.0)*. Technical Report. U.S. Department of
862 Commerce. <https://doi.org/10.6028/NIST.AI.100-1> NIST AI 100-1.
- 863 [] Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Ambrish
864 Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen,
865 Heiko Ludwig, et al. 2018. Adversarial Robustness Toolbox v1.0.0. *arXiv preprint arXiv:1807.01069* (2018).
- 866 [] NIST. 2025. *Adversarial Machine Learning: A Taxonomy and Terminology of
867 Attacks and Mitigations*. Technical Report NIST AI 100-2 E2025. National Institute
868 of Standards and Technology.
- 869 [] Nicolas Papernot, Fartash Faghri, Nicholas Carlini, Ian Goodfellow, Patrick Mc-
870 Daniel, et al. 2018. Technical Report on the CleverHans v2.1.0 Adversarial
871 Examples Library. *arXiv preprint arXiv:1610.00768* (2018).
- 872 [] Joelle Pineau, Philippe Vincent-Lamarre, Koustuv Sinha, Vincent Larivière, Alina
873 Beygelzimer, Florence d'Alché Buc, Emily Fox, and Hugo Larochelle. 2021. Im-
874 proving Reproducibility in Machine Learning Research (A Report from the
875 NeurIPS 2019 Reproducibility Program). In *Journal of Machine Learning Research*,
876 Vol. 22, 1–20.
- 877 [] Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel.
878 2020. Foolbox Native: Fast adversarial attacks to benchmark the robustness of
879 machine learning models in PyTorch, TensorFlow, and JAX. *Journal of Open
880 Source Software* 5, 53 (2020), 2607.
- 881 [] Paul C. Rigby, Brian Fitzgerald, and Viktoria Stray. 2016. Agile Development in
882 the Large: Diving into the Deep. *Companion Proceedings of the 38th International
883 Conference on Software Engineering (ICSE)* (2016), 297–300. doi:10.1145/2889160.
884 2889217
- 885 [] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017.
886 Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE
887 Symposium on Security and Privacy (S&P)*. IEEE, 3–18.
- 888 [] José Luis Solleiro and Rosario Castañón. 2005. Competitiveness and Innovation
889 Systems: The Challenges for Mexico's Insertion in the Global Context. *Techno-
890 novation* 25, 9 (2005), 1059–1070. doi:10.1016/j.technovation.2004.02.005
- 891 [] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru
892 Erhan, Ian Goodfellow, and Rob Fergus. 2014. Intriguing Properties of Neural
893 Networks. In *2nd International Conference on Learning Representations (ICLR)*.
- 894 [] Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart.
895 2016. Stealing Machine Learning Models via Prediction APIs. In *25th USENIX
896 Security Symposium*. 601–618.
- 897 [] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How
898 Does LLM Safety Training Fail?. In *Advances in Neural Information Processing
899 Systems (NeurIPS)*.
- 900 [] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt
901 Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned
902 Language Models. *arXiv preprint arXiv:2307.15043* (2023).

Temporary page!

\LaTeX was unable to guess the total number of pages correctly. As there was some unprocessed data that should have been added to the final page this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page will go away, because \LaTeX now knows how many pages to expect for this document.