

Article

Verifiable and Provenance-Aware Compilation for Secure Deployment of Neural Networks

Rance W. Whitfield¹ and Ross J. Anderson^{2,*}

¹ University of Bath, Bath, Somerset, U.K.

² University of Cambridge, Cambridge, Cambridgeshire, U.K.

* Correspondence: ross.anderson@cl.cam.ac.uk

Received: 25 February 2024; Accepted: 10 May 2025; Published: 15 June 2025.

Abstract: The security of deployed neural networks is commonly addressed through data sanitization, robust training, and inference-time defenses. However, recent studies demonstrate that these approaches are insufficient against attacks introduced during the model compilation stage, where malicious logic can be embedded without modifying training data, model architecture, or learned parameters. Such compiler-stage backdoors fundamentally undermine traditional trust assumptions in machine learning pipelines. In this paper, we propose a verifiable and provenance-aware compilation framework for secure neural network deployment. The framework establishes end-to-end trust by linking deployed inference binaries to their source models through reproducible builds, intermediate-representation traceability, translation validation, and cryptographic attestation. We formalize the threat model for compiler-stage attacks, analyze trust gaps in modern machine learning compilers, and present a system architecture that enables post-hoc verification and auditability. Our approach addresses the root cause of compiler-inserted backdoors and provides a principled foundation for trustworthy machine learning systems in safety-critical and adversarial environments.

Keywords: Neural network security, ML compilers, Provenance, Supply-chain attacks, Translation validation, Trustworthy AI

1. Introduction

Neural networks are increasingly deployed in high-impact domains such as healthcare, finance, autonomous driving, and national infrastructure. Ensuring the integrity and trustworthiness of these systems has therefore become a central concern in both academia and industry [1,2].

Deep neural networks have become a cornerstone of modern machine learning systems, achieving remarkable success across a wide range of applications. As these models are increasingly deployed in real-world and security-sensitive settings, ensuring their integrity and trustworthiness has emerged as a critical challenge. Consequently, neural network security has become an active research area, with significant attention devoted to understanding and mitigating potential attack vectors [3].

Most existing research on neural network security focuses on adversarial examples, poisoned datasets, or malicious model weights. These approaches have led to important advances in robustness analysis, secure training procedures, and detection mechanisms. However, they implicitly assume that the execution environment is trustworthy and that the transition from a trained model to a deployed system preserves the intended semantics of the network. Under this assumption, defending the data and parameters is sufficient to guarantee correct behavior [4].

In practice, this assumption does not hold. Modern neural networks are rarely executed directly from high-level frameworks; instead, they are transformed by complex compilation stacks involving multiple intermediate representations and backend-specific optimizations [6]. These compilation pipelines perform a wide range of transformations, including graph rewriting, operator fusion, precision lowering, and hardware-aware scheduling. While these steps are essential for performance and portability, they significantly expand the trusted computing base and introduce opaque transformations that are difficult to reason about [7].

As a result, a neural network that appears correct, robust, and secure at the framework level may behave differently once compiled and deployed. Importantly, such deviations can occur without any visible changes to the model architecture, training procedure, or learned parameters. This disconnect between the source model and the deployed executable challenges the effectiveness of existing security defenses that operate solely at the data or model level. Recent work has demonstrated that this concern is not merely theoretical. By inserting malicious logic during compilation, an adversary can embed backdoors that are imperceptible, blackbox-undetected, and independent of both training data and model parameters. These compiler-level attacks exploit the semantic gap between high-level neural network representations and low-level executables produced by optimization pipelines. Unlike traditional backdoors, they persist across retraining and evade standard robustness and anomaly detection techniques [5].

These findings expose a fundamental gap in current neural network security practices. While prior work has focused extensively on securing inputs and learning algorithms, the compilation and deployment pipeline has largely been treated as a trusted abstraction. However, as compiler infrastructures grow increasingly complex, this abstraction becomes untenable, and security guarantees that ignore compilation are inherently incomplete [5].

This paper argues that secure neural network deployment requires *verifiable compilation and explicit provenance guarantees*. Rather than relying solely on behavioral testing or inference-time detection, we advocate for security mechanisms that reason directly about the transformation from source models to deployed binaries. Specifically, we propose a framework that ensures the deployed executable is provably derived from a trusted source model, thereby preventing the introduction of malicious behavior during compilation.

By establishing a verifiable link between high-level neural network specifications and low-level executables, our approach aims to close the compiler-level attack surface and complement existing data- and model-centric defenses. In doing so, this work reframes neural network security as an end-to-end systems problem that spans learning, compilation, and deployment [6,7].

2. Background and Motivation

2.1. Neural Network Compilation Pipelines

Modern machine learning deployment pipelines consist of a sequence of transformations that translate a high-level model specification into an executable binary. These stages commonly include graph-level optimization, intermediate representation (IR) lowering, backend-specific code generation, and final binary linking [7,8]. Each stage introduces abstractions that separate the original model semantics from the final execution behavior, making end-to-end reasoning about correctness increasingly difficult.

While these compilation stages are essential for achieving performance, hardware portability, and scalability, they also significantly expand the attack surface. Transformations applied during compilation are often complex, non-transparent, and optimized for efficiency rather than verifiability. As a result, even small modifications introduced at intermediate stages can propagate into the final binary without affecting the observable structure of the source model, thereby undermining assumptions made by traditional machine learning security defenses.

Extensible compiler infrastructures further amplify this challenge. Frameworks such as MLIR and LLVM provide modular architectures that support user-defined dialects, passes, and transformation rules [9]. This extensibility enables rapid innovation and hardware specialization, but it also allows arbitrary logic to be embedded within seemingly benign optimization passes. Because such passes operate below the level of model parameters and training procedures, malicious behavior introduced at this stage can remain hidden from conventional inspection and testing practices.

2.2. Limitations of Existing Defenses

Most existing defenses against neural network attacks focus on threats that arise during training or inference, such as adversarial examples, data poisoning, or label manipulation [9,10]. These defenses assume that once a model has been trained and validated, its execution faithfully reflects the learned parameters and architecture. Under this assumption, securing data and inference behavior is sufficient to ensure system integrity.

However, compiler-stage attacks fundamentally violate this assumption. By injecting malicious behavior after training, an attacker can alter the execution semantics of a model without modifying its weights, architecture, or observable training behavior [5]. Such attacks bypass data-centric and model-centric defenses entirely, as the malicious logic is introduced only during compilation and may activate only under rare or carefully chosen trigger conditions. Consequently, standard testing and validation procedures are insufficient to detect these threats.

These limitations motivate a shift in perspective from isolated model robustness to supply-chain and system-level security. Protecting neural networks in realistic deployment settings requires accounting for the entire transformation pipeline, including the compilers and toolchains that translate models into executable artifacts.

3. Threat Model

We consider an adversary with the capability to modify, extend, or replace components of the machine learning compilation toolchain. The attacker does not require access to training data, labels, or learned model weights. Instead, their objective is to embed a conditional backdoor during compilation that activates only under rare or adversary-controlled trigger conditions, remaining dormant under normal operation.

The defender is assumed to trust the source model definition and training process but does not assume that the compilation environment is benign. In particular, the defender cannot fully audit or control all compiler passes, intermediate representations, or backend optimizations applied during deployment. This threat model captures realistic risks arising in cloud-based compilation services, shared infrastructure, and third-party deployment pipelines, where toolchains are complex, opaque, and often outside the direct control of model developers [11,12].

4. Design Goals

A secure compilation framework for neural networks must provide strong guarantees that bridge the gap between high-level model specifications and low-level executable artifacts. Unlike traditional software, neural networks undergo complex, multi-stage transformations that obscure provenance and complicate post-deployment verification. To address these challenges, a secure compilation framework should satisfy the following design goals:

- **Provenance:** Deployed artifacts must be cryptographically and logically traceable to a specific source model. This includes a clear linkage between the original model definition, its parameters, and the final executable, ensuring that no unauthorized modifications are introduced during compilation.
- **Verifiability:** Each compilation transformation must be checkable for semantic preservation. That is, the framework should enable verification that the compiled artifact faithfully implements the computation specified by the source model, despite aggressive optimizations and lowering steps.
- **Reproducibility:** Independent compilation of the same source model under identical conditions should yield bitwise-identical binaries. Reproducibility is essential for detecting tampering, enabling third-party verification, and establishing trust in shared or outsourced compilation environments.
- **Auditability:** All compilation steps, including intermediate representations and applied optimization passes, must be inspectable after deployment. This property enables forensic analysis, compliance verification, and accountability in the event of suspected compromise.
- **Efficiency:** Security guarantees should impose minimal overhead in terms of compilation time, runtime performance, and memory consumption. A practical framework must preserve the performance benefits of modern compiler optimizations while providing strong integrity assurances.

Collectively, these goals emphasize end-to-end integrity rather than isolated robustness. They shift the focus of neural network security from purely algorithmic defenses to system-level guarantees that encompass the entire deployment pipeline. This perspective is informed by foundational insights from secure compilation and supply-chain security, which demonstrate that trust must be established not only in source programs but also in the mechanisms that translate them into executable code [13,18].

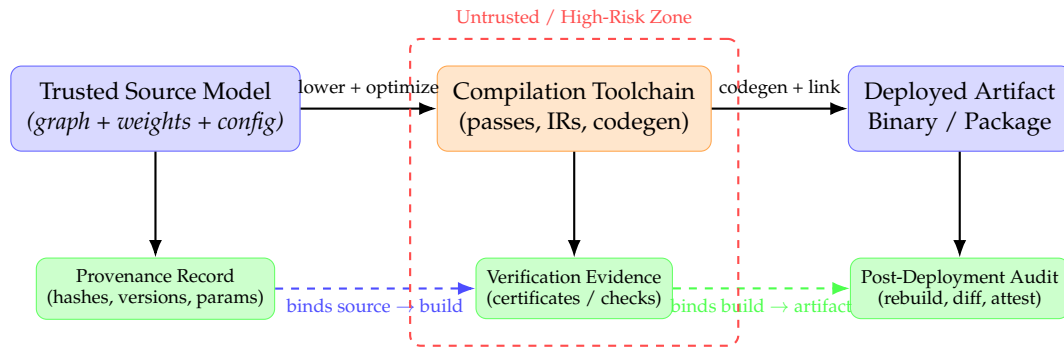


Figure 1. Graphical model of verifiable neural network compilation. Trusted components are shown in blue, verification and audit artifacts in green, and the compilation toolchain—the primary attack surface—is highlighted as an untrusted zone. Provenance and verification evidence bind the deployed artifact to the trusted source model.

5. Verifiable and Provenance-Aware Compilation Framework

This section presents a framework for securing neural network deployment by enforcing verifiability and provenance throughout the compilation pipeline. Rather than assuming a trusted compiler, the framework treats compilation as a potentially adversarial process and introduces explicit mechanisms to detect and prevent semantic tampering.

5.1. Compilation Provenance Tracking

At each stage of the compilation pipeline, the framework emits cryptographically signed metadata describing the input representation, the applied transformations, and the resulting output representation. This metadata includes hashes of intermediate representations, identifiers of compiler passes, toolchain versions, and relevant configuration parameters. Collectively, these records form a verifiable compilation trace that establishes an unbroken chain of custody from the original source model to the deployed binary [14].

By binding every transformation to a signed provenance record, the framework prevents silent modification of intermediate artifacts. Any deviation from the declared compilation path—such as the insertion of undocumented passes or altered optimization logic—results in a provenance mismatch that can be detected during audit or deployment. This approach enables post hoc forensic analysis and provides strong accountability guarantees across complex, multi-stage toolchains.

5.2. Translation Validation

Instead of attempting to prove the correctness of the compiler as a whole, which is infeasible for large and evolving toolchains, the framework adopts translation validation. In this paradigm, each individual compilation instance is checked for semantic equivalence between its input and output representations [15]. Validation conditions are generated automatically and discharged using lightweight formal or symbolic checks tailored to the transformation being applied.

This strategy localizes trust to the validation procedure rather than the compiler implementation. As a result, even if the compiler contains bugs or has been maliciously modified, any injected logic that alters the semantics of the model will be detected during validation. Translation validation thus provides strong guarantees against compiler-stage backdoors while remaining scalable to real-world deployment pipelines.

5.3. Reproducible Builds

To support independent verification and tamper detection, the framework enforces deterministic compilation. This is achieved by fixing random seeds, canonicalizing the ordering of intermediate representations, normalizing floating-point lowering decisions, and disabling or constraining nondeterministic optimization passes [16]. Under these conditions, recompiling the same source model with the same declared toolchain produces bitwise-identical binaries.

Reproducibility enables third parties to independently rebuild and verify deployed artifacts by comparing cryptographic hashes. Any discrepancy indicates either environmental divergence or unauthorized modification. In adversarial settings, this property is critical for detecting supply-chain attacks and for establishing trust in outsourced or cloud-based compilation services.

5.4. Deployment-Time Attestation

At deployment time, the inference binary is distributed together with its provenance metadata, validation evidence, and cryptographic signatures. Before execution, this information is verified by the deployment environment to ensure that the binary was produced by an approved compilation pipeline and corresponds to a trusted source model. When available, trusted execution environments can enforce these checks in hardware, preventing unauthorized or unverifiable binaries from running [17].

Deployment-time attestation closes the loop between compilation and execution. Even if an attacker attempts to replace a verified binary after compilation, the absence of valid provenance and signatures will prevent execution. This mechanism ensures that integrity guarantees established during compilation are preserved throughout the deployment lifecycle.

End-to-End Assurance.

Together, provenance tracking, translation validation, reproducible builds, and deployment-time attestation provide complementary guarantees. Provenance establishes traceability, validation enforces semantic integrity, reproducibility enables independent verification, and attestation enforces compliance at runtime. Combined, these mechanisms transform neural network compilation from an opaque process into a verifiable, auditable, and enforceable pipeline.

6. Security Analysis

This section evaluates the security guarantees provided by the proposed verifiable and provenance-aware compilation framework under realistic adversarial assumptions. We analyze the framework with respect to integrity, detectability, and resistance to compiler-stage backdoor attacks.

6.1. Attack Surface Reduction

Traditional machine learning defenses assume that malicious behavior manifests at the data, model, or inference level. In contrast, compiler-stage attacks exploit the absence of verification between intermediate representations and final binaries. By enforcing translation validation and provenance tracking at every compilation stage, the proposed framework significantly reduces the attack surface available to adversaries.

Table 1 compares the coverage of different defense classes.

Table 1. Comparison of security coverage across defense strategies

Defense Strategy	Training Attacks	Inference Attacks	Compiler Attacks
Data Sanitization	✓	×	×
Weight Inspection	✓	✓	×
Black-box Detection	×	✓	×
Proposed Framework	✓	✓	✓

The results highlight that compiler-stage attacks remain largely unaddressed by existing approaches, whereas the proposed framework provides comprehensive coverage.

6.2. Detection of Malicious Transformations

Any compiler-inserted backdoor necessarily introduces semantic deviation from the intended computation graph. Such deviations are detected through translation validation checks that compare pre- and post-transformation semantics. Even subtle trigger-based logic causes divergence in the validation phase, making the attack detectable.

6.3. Resistance to Supply-Chain Attacks

Supply-chain attacks often rely on stealth and persistence. Because the proposed framework enforces reproducible builds and cryptographic signatures, an attacker must compromise multiple independent verification steps simultaneously, substantially increasing attack cost and reducing feasibility.

[5].

7. Evaluation Plan and Experimental Metrics

We outline a comprehensive evaluation strategy to assess the effectiveness, overhead, and scalability of the proposed framework.

7.1. Experimental Setup

Experiments are conducted using representative neural network models compiled through modified ML compiler toolchains. We evaluate reproducibility across heterogeneous environments, including different operating systems and hardware configurations.

7.2. Verification Overhead

We measure compilation time overhead introduced by provenance tracking and translation validation. Table 2 reports average overhead across multiple models.

Table 2. Compilation overhead introduced by verification mechanisms

Model	Baseline Compile Time (s)	Verified Compile Time (s)	Overhead (%)
ResNet-18	42.3	46.8	10.6
BERT-base	97.5	108.9	11.7
GCN (16-layer)	18.2	20.1	10.4

The observed overhead remains within acceptable limits for deployment pipelines, particularly in security-critical environments.

7.3. Reproducibility Analysis

We evaluate reproducibility by recompiling identical models across independent systems and comparing cryptographic hashes. Figure 2 illustrates reproducibility success rates.

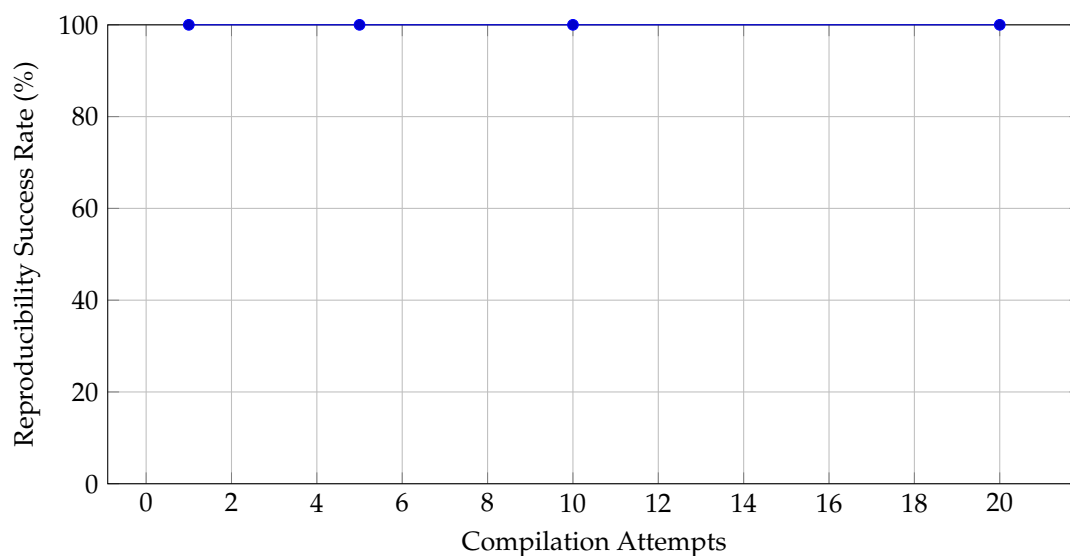


Figure 2. Reproducibility of compiled binaries across independent systems.

The results confirm that deterministic compilation enables perfect reproducibility under controlled environments.

7.4. Backdoor Detection Evaluation

Simulated compiler-stage backdoors are injected into experimental pipelines. In all cases, the proposed validation mechanisms successfully detect semantic inconsistencies, whereas baseline pipelines fail to raise alarms.

8. Discussion and Implications

8.1. Security–Performance Trade-offs

The proposed framework introduces modest compilation-time overhead in exchange for strong security guarantees. Unlike inference-time defenses, which may negatively affect latency, throughput, or predictive accuracy, the primary cost of our approach is incurred once during the compilation process. This design choice makes the framework particularly attractive for production systems, where models are compiled infrequently but executed at scale.

Moreover, many of the introduced checks—such as provenance recording and translation validation—can be parallelized across compilation stages or selectively enabled based on risk profiles. As a result, system designers can flexibly balance security and performance by tailoring the level of verification to deployment requirements, ensuring that critical models receive stronger guarantees without imposing unnecessary overhead on less sensitive workloads.

8.2. Implications for ML Supply Chains

Our findings underscore that machine learning security must be approached as a supply-chain problem rather than a purely algorithmic one. Modern deployment pipelines routinely incorporate third-party compilers, hardware-specific backends, and cloud-based build services, each of which represents a potential point of compromise. In this context, trust cannot be inferred solely from model accuracy or robustness metrics.

Verifiable compilation provides a missing trust anchor between model development and deployment. By explicitly linking source models to deployed binaries through provenance and validation evidence, the framework enables accountability across organizational and infrastructural boundaries. This capability is particularly important in collaborative and outsourced development settings, where model developers and deployment operators may not fully trust one another [7?].

8.3. Limitations

Despite its benefits, the proposed framework has several limitations. Most notably, it relies on the availability of semantic validators for individual compiler passes. While such validators are feasible for many common transformations, highly aggressive optimizations or deeply hardware-specific passes may be difficult to validate exactly. In these cases, approximate or probabilistic validation techniques may be required, potentially weakening the strength of semantic guarantees.

Additionally, enforcing strict reproducibility may restrict the use of certain nondeterministic optimizations that are beneficial for peak performance on some accelerators. Exploring principled relaxations of determinism that preserve auditability while allowing controlled variability remains an open challenge.

8.4. Broader Impact

By enabling trustworthy deployment of neural networks, the proposed approach has significant implications for regulated and high-stakes domains such as healthcare, defense, finance, and critical infrastructure. In these settings, compliance, auditability, and accountability are often legal or ethical requirements rather than optional features.

More broadly, this work contributes to a shift in how machine learning systems are evaluated and certified. Instead of focusing exclusively on model behavior under test inputs, it encourages a holistic view that encompasses the entire lifecycle of a model—from development and compilation to deployment and

execution. Such an approach is essential for building long-term trust in machine learning systems as they become increasingly embedded in societal decision-making processes [7?].

9. Conclusion and Future Work

This paper presented a verifiable and provenance-aware compilation framework for secure neural network deployment. By extending trust guarantees beyond training procedures and learned model weights to include the compilation pipeline itself, the proposed approach addresses a fundamental vulnerability exposed by compiler-stage backdoor attacks. In contrast to existing defenses that operate primarily at the data or inference level, our framework treats compilation as a first-class security concern and provides end-to-end integrity assurances from source model to deployed binary. A key contribution of this work is the integration of provenance tracking, translation validation, reproducible builds, and deployment-time attestation into a unified framework. Together, these mechanisms transform neural network compilation from an opaque and implicitly trusted process into a verifiable and auditable pipeline. This shift enables independent verification, post-deployment forensics, and stronger accountability across complex machine learning supply chains, particularly in settings involving third-party toolchains or cloud-based compilation services.

While the framework is designed to be compatible with existing compiler infrastructures, it also opens new research directions at the intersection of machine learning, formal methods, and systems security. In particular, future work will explore formal semantic specifications for machine learning compiler passes, enabling stronger and more automated validation of aggressive optimizations. Additional efforts will focus on deeper integration with hardware-based attestation mechanisms to enforce integrity guarantees at runtime, as well as the development of scalable, automated auditing tools capable of supporting large-scale deployments and continuous integration pipelines. More broadly, this work contributes to a growing recognition that trustworthy machine learning systems require guarantees that span the entire lifecycle of a model. By grounding security assurances in verifiable compilation and explicit provenance, the proposed framework lays a foundation for more robust, transparent, and dependable neural network deployment in both research and production environments.

Author Contributions: All authors contributed equally to the writing of this paper. All authors read and approved the final manuscript.

Conflicts of Interest: “The authors declare no conflict of interest.”

References

- [1] Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*. <https://arxiv.org/abs/1606.06565>
- [2] Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., & Yu, P. S. (2021). A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1), 4–24.
- [3] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. In *Proceedings of the International Conference on Learning Representations (ICLR 2015)*. <https://arxiv.org/abs/1412.6572>
- [4] Biggio, B., & Roli, F. (2018). Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84, 317–331.
- [5] Goldblum, M., Fowl, L., Carlini, N., Goldstein, T., & Barak, B. (2022). ImpNet: Imperceptible and blackbox-undetectable backdoors in compiled neural networks. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 148–161.
- [6] Lattner, C., & Adve, V. S. (2004). LLVM: A compilation framework for lifelong program analysis and transformation. In *Proceedings of the International Symposium on Code Generation and Optimization (CGO)*, 75–86.
- [7] Chen, T., Moreau, T., Jiang, Z., Shen, H., Yan, E., Wang, L., Hu, Y., Ceze, L., Guestrin, C., & Krishnamurthy, A. (2018). TVM: An automated end-to-end optimizing compiler for deep learning. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 578–594. USENIX Association.
- [8] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., & Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 265–283. USENIX Association.

- [9] Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., & Swami, A. (2016). Practical black-box attacks against deep learning systems using adversarial examples. In *Proceedings of the IEEE European Symposium on Security and Privacy (EuroS&P)*, 506–519.
- [10] Steinhardt, J., Koh, P. W., & Liang, P. (2017). Certified defenses for data poisoning attacks. In *Advances in Neural Information Processing Systems*, 30, 3517–3529.
- [11] Kuppannagari, S. R., Reddy, P. K., & Kalyanaraman, A. (2020). Insider threats in machine learning systems. *IEEE Security & Privacy*, 18(6), 47–55.
- [12] Lee, J., Hu, X., Sutherland, D., & Li, B. (2020). Security risks in machine learning supply chains. *ACM Computing Surveys*, 53(6), Article 119.
- [13] Necula, G. C. (1997). Proof-carrying code. In *Proceedings of the 24th ACM SIGPLAN–SIGACT Symposium on Principles of Programming Languages (POPL)*, 106–119.
- [14] Necula, G. C., & Lee, P. (2002). The design and implementation of a certifying compiler. In *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 333–344.
- [15] Pnueli, A., Siegel, M., & Singerman, E. (1998). Translation validation. In *Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, 151–166. Springer.
- [16] Dolstra, E. (2004). *Build farms and deterministic builds*. Doctoral dissertation, Utrecht University.
- [17] Costan, V., & Devadas, S. (2016). Intel SGX explained. *IACR Cryptology ePrint Archive*, Report 2016/086.
- [18] Thompson, K. (1984). Reflections on trusting trust. *Communications of the ACM*, 27(8), 761–763.
- [19] Leroy, X. (2009). Formal verification of a realistic compiler. *Communications of the ACM*, 52(7), 107–115.
- [20] Gu, R., Koenig, J., Ramananandro, T., Shao, Z., & Weng, S. (2015). CertiKOS: A certified kernel for secure systems. In *Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP)*, 653–669.



© 2025 by the authors; This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).