

# B-SAFE: Blockchain Security Assessment Framework Enhanced with Machine Learning \*

**Ngô Thanh Trung**  
Troy University  
Hanoi, Viet Nam  
tngo220196@troy.edu

**Pham Tien Dat**  
Troy University  
Hanoi, Viet Nam  
dpham220298@troy.edu

**Pham Thai Duong**  
*Troy University*  
*Hanoi, Viet Nam*  
*dpham220299@troy.edu*

Le Quang Huy  
Troy University  
Hanoi, Viet Nam  
hle220331@troy.edu

**Doan Hoang Long**  
Troy University  
Hanoi, Viet Nam  
ldoan220279@troy.edu

**Abstract**—This paper presents B-SAFE, an empirically grounded systematization of blockchain security risks organized across five architectural layers. We analyze 649 incident entries spanning 2016–2025, applying a formal specification schema (P, I, S, C, M) to classify threats and map defenses. We provide an enterprise-relevant taxonomy with rigorous incident analyses and a practical assessment checklist. We implement and evaluate an LLM-assisted pipeline (LLM→XGBoost→LLM) that accelerates assessments by converting enterprise documents into predictions and an executive-ready report (see §A.). This work aims to standardize terminology, surface consistent invariants and controls, and support reproducible, practitioner-focused security assessments.

**Keywords**—Blockchain security, machine learning, security assessment, threat detection, consensus mechanisms, smart contracts

## I. INTRODUCTION

Blockchain technology has revolutionized digital trust and decentralized applications, but this innovation has been accompanied by significant security challenges. The rapid proliferation of smart contracts, DeFi protocols, and cross-chain infrastructure has created a complex attack surface that traditional security frameworks struggle to address systematically. This paper presents B-SAFE, a comprehensive framework for blockchain security assessment that provides a unified approach to understanding, classifying, and quantifying security risks across the entire blockchain ecosystem.

The B-SAFE framework addresses the critical need for systematic security analysis in blockchain systems by introducing a five-layer reference architecture that captures the complete attack surface. Our methodology enables reproducible security assessment through standardized data collection, incident labeling, and risk quantification, providing actionable insights for researchers, practitioners, and policymakers.



FIGURE I: WORD CLOUD GENERATED FROM TITLES AND DESCRIPTIONS OF BLOCKCHAIN SECURITY INCIDENTS, SHOWING THE MOST FREQUENTLY OCCURRING TERMS.

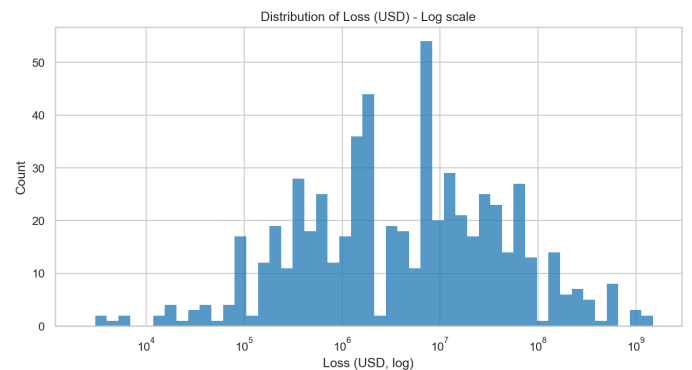


FIGURE II: DISTRIBUTION OF LOSS (LOG-SCALE).

\***Cite (APA):** Trung N., Dat P., Long D., Duong P., Huy L. (2025).

---

## II. FOUNDATIONS AND VULNERABILITY LANDSCAPE

### A. Consensus and Network-Layer Attack Surface

While blockchain technology is renowned for emulating a “trusted” service through a decentralized and immutable ledger, its foundational security assumptions are not infallible. The incentive mechanisms designed to ensure honest participation in consensus protocols, particularly in permissionless networks, have been openly questioned and are vulnerable to exploitation. This analysis targets the system’s core, focusing on vulnerabilities at the consensus and P2P network layers, such as selfish mining and block withholding. These attacks are often complex, leveraging game-theoretic strategies to gain disproportionate rewards. For enterprise-grade systems like Hyperledger Fabric, which are intended for business use, the impact of such consensus failures is severe. Therefore, a robust security assessment framework is essential to mitigate these threats and ensure the trusted adoption of blockchain in critical sectors.

Indeed, the security of consensus protocols is not an abstract guarantee but an emergent property of specific economic and network conditions. Attacks on this layer are not mere theoretical possibilities; they are practical exploits of measurable weaknesses in a blockchain’s ecosystem. These vulnerabilities are direct outcomes of insufficient economic security and inherent physical limitations in network communication, which pose tangible risks to any enterprise system built upon them. The economic security of a Proof-of-Work blockchain, for instance, is a direct function of its total hash rate; when this hash rate is low, the ledger’s immutability becomes fragile and susceptible to being forcibly rewritten. This is not a theoretical vulnerability, but a recurring reality for smaller chains, with networks like Ethereum Classic (ETC) and Bitcoin Gold (BTG) having been successfully attacked multiple times, leading to tens of millions of dollars in fraudulent transactions [1]. The ease with which these historical rewrites are executed stems from a fundamental shift in the economics of acquiring computational power. An attack that once required a prohibitive capital investment in mining hardware now becomes a manageable operational cost, rented by the hour from public hashrate markets like NiceHash. For any enterprise application built on such a minority chain, this commodification of hashrate represents a persistent, existential threat to data integrity.

Beyond attacks of pure computational force, a more insidious class of vulnerability arises from exploiting the unavoidable latencies of a global peer-to-peer network. The Selfish-Mine strategy masterfully turns the network’s core arbitration mechanism—the “longest chain” rule—into a weapon against itself [2, 3]. This is not a theoretical exercise for individual miners but a viable strategy for the large, coordinated mining pools that already dominate network hashrate [2]. The original analysis warned that pools existed which exceeded the 25

The security of a blockchain also relies fundamentally on the integrity of the communication network that binds its participants. Vulnerabilities in this fabric can be exploited to distort a node’s perception of the blockchain, ranging from the targeted isolation of a single peer to the large-scale partitioning of the entire network. At the individual node level,

an adversary can execute an Eclipse attack to monopolize a victim’s network connections, effectively creating a fabricated reality. This is often enabled by a foundational Sybil attack, where the adversary generates numerous pseudonymous identities to overwhelm the victim’s peer-discovery mechanism. Once eclipsed, the victim is completely severed from the honest network, and its view of the blockchain is dictated by the attacker, facilitating targeted double-spends or the co-opting of mining power. While an Eclipse attack blinds an individual, a more ambitious adversary can target the internet’s core routing infrastructure. By manipulating the Border Gateway Protocol (BGP), an attacker can hijack traffic routes, partitioning the blockchain network into isolated sub-networks. Each partition, now operating with a fraction of the global hash rate, becomes dangerously vulnerable to a 51

Furthermore, alternative consensus models like Proof-of-Stake (PoS) introduce novel attack vectors that shift the focus from computational power to the manipulation of economic stakes over time. A primary threat is the long-range revision attack. In PoS, validators’ influence is tied to an economic stake that is slashed for misbehavior; however, once validators have safely withdrawn their deposits, they are no longer subject to this penalty. A coalition of these historical validators can then use their old private keys to build and sign an entirely new, conflicting blockchain history starting from a point deep in the past, without fear of being slashed. Another critical vulnerability is the catastrophic crash, where if more than one-third of validators simultaneously go offline, the system cannot form the required two-thirds supermajority to finalize new checkpoints, effectively halting the ledger’s progress. These attacks highlight that shifting from a computational to a capital-based consensus model introduces new and complex failure modes that challenge a chain’s finality and liveness.

When blockchain technology is applied to solve real-world problems in finance or healthcare, the nature of security risk changes profoundly. Threats expand to operational issues, regulatory compliance, and business logic vulnerabilities at the application layer. In sectors like banking, transaction finality is a non-negotiable requirement. The risk of chain reorganizations, however small, can reverse confirmed payment transactions, causing chaos in settlement systems and eroding customer trust [1]. Concurrently, strict data privacy regulations like GDPR or HIPAA pose a significant challenge. The immutable nature of blockchain directly conflicts with a user’s “right to be forgotten,” raising the difficult question of how to delete patient data in a compliant manner without breaking the chain’s integrity [4]. Furthermore, while a blockchain secures data after it has been written, it cannot validate the accuracy of the information at the point of entry—a critical “garbage in, garbage out” risk where an incorrect electronic health record could persist immutably on the ledger [4].

Perhaps the largest attack surface in these enterprise applications lies within smart contracts themselves. They are the digital embodiment of business agreements, and any flaw in their encoded logic can be exploited. An attacker does not need to break the consensus mechanism; they only need to find a business logic flaw to drain funds from a complex financial instrument or illicitly access sensitive data [5]. This transforms smart contract auditing and formal verification from an option into a mandatory requirement for system security. As demonstrated, the theoretical integrity of a blockchain is fundamentally contingent on the security of its consensus and net-

---

work layers. The vulnerabilities analyzed, from game-theoretic exploits at the consensus layer to the manipulation of network topology, pose tangible and high-impact risks to enterprise systems. Proactive security assessment is therefore not merely a recommendation but an essential prerequisite for trusted adoption in critical applications.

## B. Key Management and Wallet Security

Proper key management is the most important part of security for any blockchain-based system. Even the strongest protocols can fail if keys are not handled correctly [6]. In decentralized systems, private keys give users final control over their digital assets, identity, and ability to perform actions on the blockchain. This idea is often summarized by the saying, "Not your keys, not your coins," but it applies to more than just currency [7]. The main tools users have for managing these keys are called "wallets." The security of these wallets is therefore essential for protecting user actions on the blockchain [8]. However, wallet security is not just a technical problem; it also depends on software design and, importantly, on the behavior and understanding of the users themselves [7].

To understand wallet security, it is helpful to first classify the different types of wallets. The most basic classification is between "hot wallets," which are connected to the internet, and "cold wallets," which are kept offline. Hot wallets include desktop software, mobile apps, and web browser extensions. They are easier to use for daily transactions, but their online nature makes them more vulnerable to attacks. Cold wallets, such as hardware devices or paper wallets, offer better security for long-term storage because they are not directly exposed to online threats [8]. Another important classification is based on who controls the keys. With "custodial wallets," a third party like a cryptocurrency exchange holds the keys for the user. This is simpler for beginners, as the experience is similar to online banking, but it requires trusting that the third party is competent and honest. With "non-custodial wallets," users have full control and responsibility over their own keys. These non-custodial wallets can be further divided into traditional Externally Owned Accounts (EOAs) and newer Smart Contract wallets, which allow for more complex security rules [7, 8].

The vulnerabilities in these systems exist at multiple levels, but the most common threats are those on the user's own device [9]. A major technical risk is the improper storage of keys, such as saving them as unencrypted plaintext in the device's memory, where they can be stolen by malware. Flaws in the wallet software, such as insecure interfaces or the use of buggy code libraries, also create significant risks. This is not just a theoretical problem; real-world attacks often exploit these weaknesses. For example, weak key generation methods like "brain wallets," which use simple, memorable phrases, are a critical vulnerability. The low entropy of human-generated phrases makes them easy to guess, and one study found that most such wallets were drained of funds in less than 24 hours [9]. At the institutional level, the history of exchange hacks like the infamous Mt. Gox incident shows that even large platforms can have critical flaws [9]. More recently, the collapse of the FTX exchange served as a powerful reminder of counterparty risk—the danger that the trusted third party will fail due to mismanagement or fraud, leading to a total loss of user funds [7]. Even at an individual level, social threats are a major risk; one study documented a user who lost all their funds

simply because they let a friend see their login details during the wallet setup process, highlighting the dangers of misplaced trust [7].

To protect against these threats, both technical and user-driven defense methods are used. The main technical defense is to use cold storage, such as hardware wallets, to keep keys offline and safe from online hackers [8]. More advanced solutions include multi-signature and smart contract wallets, which allow for programmable security rules like spending limits or requiring multiple people to approve a transaction [7]. However, since many attacks target the user, user-driven strategies are just as important. A common and effective strategy is "risk diversification," where users spread their assets across multiple wallets. For instance, a user might keep a small amount of "spending money" in a convenient mobile hot wallet, while keeping the majority of their savings in a more secure cold wallet [7]. They also use different wallets for different tasks, for example, using a dedicated wallet with minimal funds for interacting with new or risky dApps. For high-value transactions, many users prefer a PC setup because they can use third-party security extensions, like Fire or Revoke.cash, which simulate transactions and warn them about malicious smart contracts before they sign [7].

We can see these security trade-offs in the real-world systems that users choose. Centralized exchanges like Coinbase offer a simple user experience that is similar to online banking. This makes them popular with beginners, but it comes with significant counterparty risk, as tragically demonstrated by the failure of FTX [7]. Hardware wallets like Ledger or Trezor represent the opposite approach. They provide high security by giving users full control over their offline keys, but they can be difficult to use and require the user to be fully responsible for their own security. This trust model has also been challenged recently. For example, Ledger's controversial "Recover" service, which proposed storing shards of a user's seed phrase with third parties, caused a backlash because it went against the core reason users chose a hardware wallet: to be the sole holder of their keys [7]. As a middle ground, new systems like smart contract wallets (e.g., Argent) are emerging. They try to offer the best of both worlds: strong security features like social recovery to prevent key loss, combined with an easier user experience that often removes the need to manually manage a seed phrase. These different models show that the market is still searching for the right balance between security, usability, and trust [7].

## C. Smart Contract Vulnerabilities

Smart contracts represent a fundamental advancement in blockchain technology, enabling the execution of programmable, self-enforcing agreements on decentralized platforms such as Ethereum. While these immutable protocols have revolutionized digital asset transactions, they simultaneously introduce significant security challenges that require systematic analysis and robust mitigation strategies [10]. The immutability property that enhances trust also presents a critical constraint: once deployed, code vulnerabilities cannot be patched through conventional means, thereby amplifying the potential consequences of security failures [11]. The security research community has documented several catastrophic incidents that demonstrate the real-world implications of smart contract vulnerabilities. Notable examples include The DAO

---

attack and the Parity wallet incidents, which resulted in substantial financial losses and permanently frozen assets, respectively. These events underscore that smart contract vulnerabilities transcend theoretical concerns and constitute material threats to blockchain ecosystems [10]. The analysis of these incidents reveals a pattern of specific vulnerability classes that demand systematic detection and prevention methodologies.

Reentrancy vulnerabilities represent one of the most extensively studied attack vectors in smart contract security. This vulnerability materializes when a contract performs an external call prior to updating its internal state, thereby enabling recursive invocation of sensitive functions. The DAO exploit, which resulted in the theft of approximately 3.6 million ETH, exemplifies the potential magnitude of reentrancy attacks. However, empirical analysis indicates that only 0.3% of contracts identified as vulnerable to reentrancy have experienced actual exploitation, suggesting that detection tools may overestimate practical risk levels [10, 12].

Authorization flaws constitute another critical vulnerability class, typically manifesting as insufficient access control mechanisms for privileged operations. The Parity Multi-Sig wallet incidents provide instructive examples of such vulnerabilities, where inadequate authorization checks enabled attackers to either appropriate funds or permanently disable contract functionality. These incidents demonstrate how seemingly minor oversights in access control can produce disproportionate consequences in decentralized systems [11].

Integer overflow and underflow vulnerabilities, while conceptually straightforward, have precipitated significant financial disruptions. The Beauty Chain token incident illustrates this vulnerability class, wherein arithmetic operations exceeding fixed-width integer bounds resulted in the creation of excessive tokens, destabilizing the entire tokenomics system. While contemporary Solidity versions implement automatic overflow checking, legacy contracts remain susceptible without explicit safeguards such as the SafeMath library [11, 12].

External data dependencies introduce distinct vulnerability classes related to oracle inputs and timestamp manipulation. Smart contracts often require external data sources for critical operations, creating attack surfaces where manipulated inputs can compromise contract integrity. The proliferation of flash loan mechanisms has exacerbated these risks by providing temporary access to substantial capital for market manipulation within single transactions. Such attacks have targeted price oracles in decentralized finance protocols with considerable success [10, 11].

Delegatecall vulnerabilities represent potentially the most devastating attack vector, as this operation executes external code within the storage context of the calling contract. The second Parity incident exemplifies this risk, wherein an unprotected library function allowed the destruction of shared contract infrastructure, permanently immobilizing approximately 160 million in user assets. Notably, this incident resulted not from malicious intent but from inadvertent interaction with unprotected functionality [11].

For vulnerability detection and prevention, the security community employs three primary methodological approaches, each with distinct characteristics and limitations.

Static analysis tools examine contract source code or byte-

code without execution, providing comprehensive coverage but frequently generating false positives. Comparative studies reveal significant inconsistency between tools, with inter-tool agreement on identified vulnerabilities ranging from 1.85% to 23.9%, indicating the necessity for multi-tool approaches [10, 12].

Dynamic analysis techniques implement a more empirical methodology by executing contracts with potentially malicious inputs. These approaches generate concrete exploitation scenarios but cannot exhaustively explore all execution paths. Tools such as ContractFuzzer and MAIAN exemplify this category, offering higher precision but more limited coverage than static alternatives [11].

Formal verification represents the most rigorous security approach, providing mathematical guarantees of contract correctness according to specified properties. Despite its theoretical strength, formal verification requires substantial expertise and resources, as demonstrated by the MakerDAO verification process, which required eight person-months to complete. This approach remains most suitable for high-value or critical infrastructure contracts [11].

The security community has developed standardized defensive patterns to address common vulnerabilities. The checks-effects-interactions pattern mitigates reentrancy by ensuring state updates precede external calls. Role-based access control systems protect privileged functions, while careful upgradeability design preserves system integrity during evolution [10, 11]. A notable empirical observation is the significant disparity between theoretical vulnerability prevalence and actual exploitation rates. Despite numerous contracts containing potential vulnerabilities, exploitation remains relatively rare. This phenomenon appears attributable to economic factors: approximately 0.01% of contracts control 83% of all ETH, and these high-value targets typically implement more robust security measures. This distribution suggests that security analysis must incorporate economic incentives alongside technical considerations to accurately assess real-world risk [10, 12].

The analysis of smart contract vulnerabilities reveals a complex landscape where technical vulnerabilities intersect with economic incentives and practical exploitation constraints. While significant progress has been made in identifying and mitigating common vulnerability classes, the empirical evidence suggests that the real-world risk may be lower than theoretical analyses indicate. Nevertheless, the catastrophic impact of successful exploits necessitates continued vigilance and the application of multiple verification methodologies to secure blockchain-based systems.

#### *D. DeFi Protocol Risks*

Decentralized Financial ecosystem (DeFi), is built based on blockchain platforms such as Ethereum, has emerged as an alternative to Centralized Finance due to its transparency, traceability, and decentralized nature. DeFi offers a wide range of financial services, primarily implemented through smart contracts. However, the rapid growth of DeFi has also come with serious security risks, leading to significant financial losses. While blockchain technology itself is considered secure due to its properties such as immutability and consensus mechanisms, the applications and additional layers built on top of blockchain – namely DeFi protocols – are not entirely secure and can be

vulnerable.

Many recent works have systematized DeFi into layers (network, consensus, smart-contract, protocol, auxiliary services) and emphasized that many incidents arise from unsafe dependencies between protocols and off-chain services (oracles, centralized relays, bridges) [13]. Among them, vulnerabilities in the DeFi protocol layer (PRO Layer) are often related to design flaws or financial market manipulation. For instance, pricing mechanisms, slippage, liquidation mechanisms, rebases... or invalid assumptions about token standards can be catastrophic when contracts are composed together; in particular, external dependencies are called directly without consistency checks are the source of many real-world failures. [13]

A key economic risk is flash loans, uncollateralized lending mechanisms in an atomic transaction. Flash loans have opened a new attack vector where an attacker can temporarily borrow large amounts of capital to manipulate the market or price feed, performing a series of profit and debt repayment operations in the same transaction. Attacks like Harvest, PancakeBunny, Beanstalk... [13,14] show that flash loans lower the cost barrier to attack and make small design issues become financial catastrophic. Another risk directly related to off-chain backends is that when price data sources are manipulated – through source changes, on-chain update attacks, or updater compromises – key parameters such as liquidation prices or collateralization ratios can become distorted, leading to mass liquidations or systemic profiteering [14]. There are mitigations such as multiple source aggregation, medianizers, or latency mechanisms that exist but carry trade-offs in latency, centralization and fault tolerance [13,14].

In addition, transaction ordering and MEV (Miner/Maximal Extractable Value) issues allow sequencers or miners to order, insert or remove transactions to maximize profits – this mechanism gives rise to front-running, sandwiching and other mining strategies, which directly impact the stability of the protocol’s financial invariants [13]. Expanding the functional space with cross-chain bridges also creates a new attack surface: many bridges rely on centralized signing/organizations, and bridge crashes have led to large scale asset losses, demonstrating a clear trade-off between cross-chain utility and security risk [14]. Finally, operational and human risks – including private key, mismanagement (privileged keys, weak multisig...), compromised front ends, and implement flaws (not pure protocol design flaws) have a direct impact on asset security and are often present in real-world incidents [15].

To mitigate these risks, incident studies and analysis have proposed a multilayer set of measures: protocol design that considers both economic attack scenarios (game-theoretic stress testing) and defense mechanisms such as circuit breakers [13]; oracle enhancements using aggregations, delayed updates or reputation-based models [14]; MEV mitigations using transparent sequencers or close-chain relay [13]; along with audit, formal verification and real-time monitoring (e.g., oracle mutation detection) with response options as emergency halts [13,14]. Each approach carries trade-offs in performance, latency, and decentralization, so the choice of solution should be based on the specific application context.

Finally, the systematic analysis revealed important research gaps: the lack of a comprehensive quantitative framework for protocol economic risk (incorporating TVL, liquidity depth, or-

acle latency, and flash loan capabilities), the lack of a common fault tolerant architectural pattern for trustless backends, and the lack of dependency analysis tools for complex composability environments – these gaps share the research direction needed to improve the robustness of DeFi protocols in the broader blockchain landscape.

### E. Exchange and Infrastructure Attacks

This subsection covers centralized exchange compromise patterns, API key abuse, withdrawal bypasses, hot/cold wallet segregation failures, and infrastructure supply-chain risks. We incorporate regulatory and compliance impacts for enterprises.

In enterprise contexts, the most material risks arise from operational lapses and supply-chain exposures rather than protocol breaks. Common vectors include cloud credential leakage leading to hot-wallet drains, CI/CD pipeline tampering that injects malicious front-end code, and API key over-permissions enabling withdrawal bypasses. Controls that consistently reduce risk include hardware-backed key custody (HSMs with quorum policies), strict withdrawal allowlists and velocity limits, just-in-time API credentials, and continuous integrity monitoring of frontend artifacts. From a compliance perspective, incident containment and audit-ready evidence (tamper-evident logs, key ceremony records) are as critical as technical remediation.

## III. REFERENCE FRAME

This section establishes the theoretical foundation of the B-SAFE framework. We introduce a five-layer reference architecture for holistic security analysis and present our formal risk classification framework, which serves as the primary tool for the systematic security assessment conducted in this research.

### A. Five-Layer Blockchain Security Architecture

The B-SAFE framework is built upon a comprehensive five-layer reference architecture that captures the complete attack surface of blockchain systems. This layered approach enables systematic security analysis by organizing threats according to their architectural context and attack vectors.

#### 1. Layer Definitions

- **NET (Network Layer):** Encompasses network-level attacks including eclipse attacks, Sybil attacks, and network partitioning vulnerabilities that can disrupt consensus and transaction propagation.
- **CON (Consensus Layer):** Addresses consensus mechanism vulnerabilities such as 51% attacks, selfish mining, and consensus rule violations that threaten the fundamental security guarantees of the blockchain.
- **SC (Smart Contract Layer):** Covers smart contract vulnerabilities including reentrancy attacks, integer overflow, and logic flaws that can lead to unauthorized fund transfers or contract manipulation.
- **PRO (Protocol Layer):** Encompasses DeFi protocol-specific risks including flash loan attacks, oracle manipulation, and protocol governance vulnerabilities that can exploit economic incentives and protocol mechanics.

- **AUX (Auxiliary Layer):** Addresses supporting infrastructure risks including wallet security, key management, exchange vulnerabilities, and off-chain dependencies that can compromise user assets and system integrity.

## 2. Cross-Layer Dependencies

The layered architecture recognizes that attacks often span multiple layers, with vulnerabilities in one layer enabling or amplifying threats in others. This interdependency is explicitly modeled in our risk assessment framework to provide a holistic view of the attack surface.

## B. B-SAFE Risk Classification Framework

To provide a systematic and reproducible security assessment, we establish a formal risk classification framework. Each identified threat category is specified through a standardized schema that defines its preconditions, the system invariants it threatens, its canonical attack vector, applicable defense mechanisms, and quantitative risk metrics.

### 1. Risk Category Specification Schema

Each risk category  $R$  is formally defined by the tuple  $(P, I, S, C, M)$  where:

- $P = \{p_1, p_2, \dots, p_n\}$  represents the set of preconditions that must hold for an attack to be feasible.
- $I = \{inv_1, inv_2, \dots, inv_m\}$  represents the set of core system invariants threatened by the attack.
- $S = (s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_k)$  represents the canonical sequence of steps in the attack vector.
- $C = \{\text{Prevention, Mitigation, Detection}\}$  represents the categories of defense mechanisms and controls.
- $M = (L, I, D, R)$  represents the quantitative risk metrics for the category.

### 2. Quantitative Risk Scoring

The priority ranking  $R$  is calculated using a weighted formula based on Likelihood ( $L$ ), Impact ( $I$ ), and Detectability ( $D$ ), each rated on a scale of 1 to 5. The formula is defined as:

$$\text{Risk Score} = (w_L \times L) + (w_I \times I) - (w_D \times D) \quad (1)$$

We adopt a unified quantitative approach that combines likelihood, impact, and detectability into a single risk score. The unified formula accounts for detectability as a risk modifier, recognizing that harder-to-detect threats pose greater risk. This approach balances traditional risk management principles with the unique challenges of blockchain security. Sensitivity analysis examines ranking stability under weight variation.

FIGURE III: VISUAL REPRESENTATION OF THE (P, I, S, C, M) SCHEMA APPLIED TO SC-1 REENTRANCY ATTACK

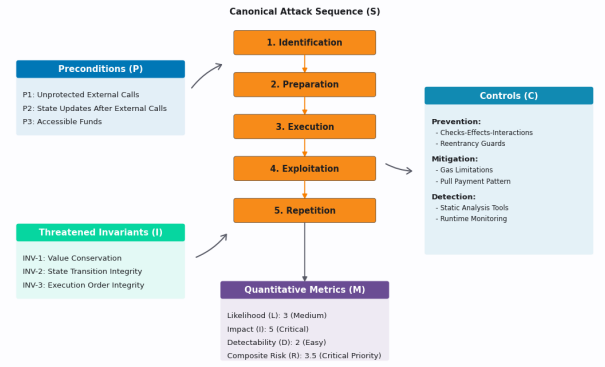


FIGURE III: VISUAL REPRESENTATION OF THE (P, I, S, C, M) SCHEMA APPLIED TO SC-1 REENTRANCY ATTACK. RIGHT: ILLUSTRATIVE SENSITIVITY OF UNIFIED RISK SCORE TO WEIGHT VARIATION DEMONSTRATES PRIORITIZATION STABILITY ACROSS BALANCED WEIGHT CONFIGURATIONS.

## IV. METHODOLOGY

To construct our analysis of blockchain security incidents, we adopted a systematic methodology for data collection, labeling, and quantification. This chapter details the protocol used to build our dataset and the framework for its analysis, ensuring our results are transparent and reproducible.

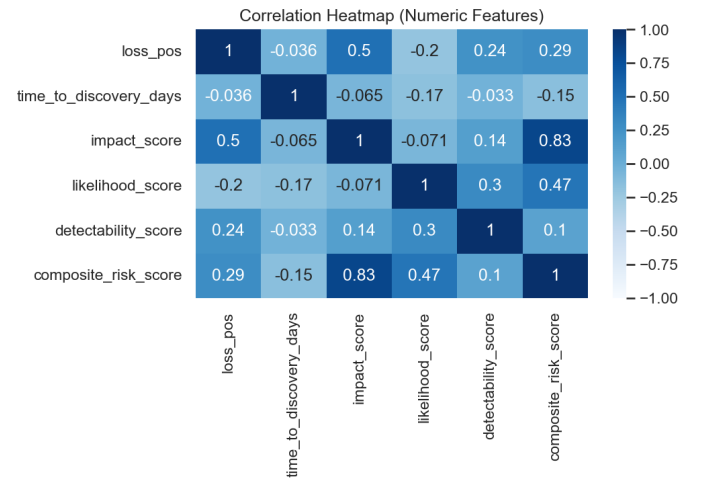


FIGURE IV: CORRELATION HEATMAP OF KEY VARIABLES: LOSS, RISK SCORES, AND TIME-TO-DISCOVERY.

### A. Data Collection and Sources

Our dataset comprises 649 blockchain security incident entries collected from multiple sources spanning 2016–early 2025. We employed a systematic approach to ensure comprehensive coverage while maintaining data quality and verifiability.

---

## 1. Primary Data Sources

We collected incident data from the following primary sources:

- **Web3IsGoingGreat:** A comprehensive database of blockchain security incidents with detailed incident reports and loss estimates
- **Rekt News:** Specialized platform tracking DeFi exploits and protocol vulnerabilities
- **Secureum:** Academic and industry reports on smart contract vulnerabilities and attacks
- **Chainalysis:** Blockchain analytics data for incident verification and impact assessment
- **Academic Literature:** Peer-reviewed papers and technical reports from security conferences

## 2. Inclusion and Exclusion Criteria

To ensure data quality and relevance, we applied the following criteria:

### Inclusion Criteria:

- Minimum financial loss of \$10,000 USD (adjusted for inflation)
- Verifiable incident reports with multiple independent sources
- Clear attribution to specific blockchain platforms or protocols
- Sufficient technical details to classify the attack vector

### Exclusion Criteria:

- Purely anecdotal or unverified reports
- Incidents with insufficient technical details for classification
- Non-blockchain related security incidents
- Duplicate reports of the same incident

## 3. Data Collection Protocol

Our data collection process followed a standardized protocol:

1. **Source Identification:** Systematic review of primary and secondary sources
2. **Initial Screening:** Application of inclusion/exclusion criteria
3. **Data Extraction:** Structured extraction of incident details, financial impact, and technical characteristics
4. **Verification:** Cross-referencing with multiple sources for accuracy
5. **Quality Control:** Review by multiple team members for consistency

## 4. Dataset Schema and Dictionary

The curated incident dataset contains 649 entries structured with 18 columns. Each row corresponds to a single classified incident with complete B-SAFE framework annotations:

- **incident\_id:** Unique integer identifier for the security event
- **incident\_title:** Concise, descriptive title for the incident
- **incident\_date:** Date when the incident occurred or was first reported (MM/DD/YYYY format)
- **paused\_at:** Timestamp when a protocol was paused, if applicable
- **incident\_description:** Detailed text summary of the incident
- **loss\_usd:** Estimated financial loss in USD at the time of the incident
- **chain:** Primary blockchain or ecosystem affected (e.g., Ethereum, BNB Chain)
- **source\_link:** URL to a primary report or analysis of the incident
- **b\_safe\_layer:** Assigned B-SAFE architectural layer (NET, CON, SC, PRO, AUX)
- **b\_safe\_risk\_category:** Specific risk category identifier within the B-SAFE framework (e.g., SC-1, PRO-2)
- **technical\_footprint:** JSON object containing structured technical details of the affected protocol, used for feature engineering
- **attacker\_type:** Classification of the adversary (e.g., Insider, External)
- **recovery\_status:** Status of the lost funds (e.g., Funds Lost, Funds Recovered)
- **time\_to\_discovery\_days:** Estimated time in days from exploit to public discovery
- **impact\_score:** Assigned Impact (I) score on a 1-5 scale, as defined in Section IV.D
- **likelihood\_score:** Assigned Likelihood (L) score on a 1-5 scale, as defined in Section IV.D
- **detectability\_score:** Assigned Detectability (D) score on a 1-5 scale, as defined in Section IV.D
- **composite\_risk\_score:** Calculated priority ranking based on our formula

The dataset includes complete risk scoring for all 649 incidents, with 100% coverage for impact, likelihood, and detectability scores, enabling comprehensive statistical analysis and validation of our risk quantification approach.

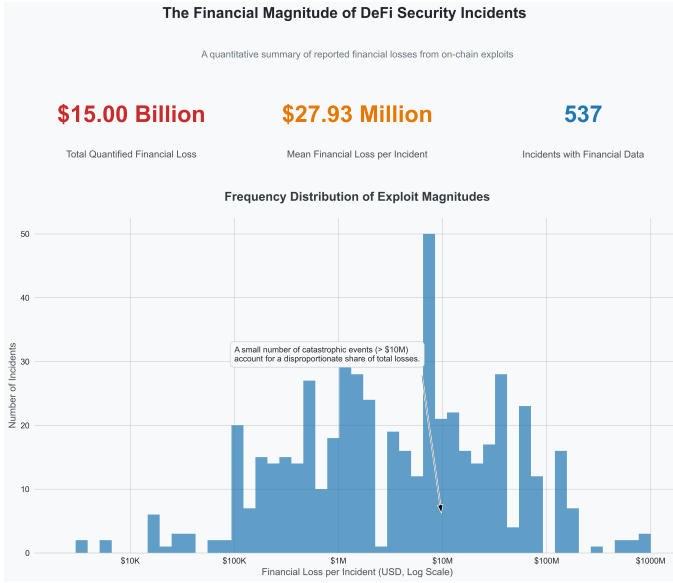


FIGURE V: FINANCIAL MAGNITUDE OF DEFI SECURITY INCIDENTS

## B. Incident Labeling Protocol

To ensure consistent and reproducible classification of security incidents, we developed a comprehensive labeling protocol that maps each incident to our formal risk classification framework.

### 1. Annotation Process

Our labeling process involved multiple stages to ensure accuracy and consistency:

1. **Initial Classification:** Each incident was initially classified by a primary annotator
2. **Peer Review:** A second annotator reviewed and validated the classification
3. **Expert Adjudication:** Disagreements were resolved through expert review
4. **Quality Assurance:** Final classifications underwent quality control checks

### 2. Labeling Schema

Each incident was labeled according to the following schema:

#### Layer Classification (L):

- **NET:** Network layer attacks (eclipse, Sybil, partitioning)
- **CON:** Consensus layer attacks (51%, selfish mining, rule violations)
- **SC:** Smart contract layer attacks (reentrancy, overflow, logic flaws)
- **PRO:** Protocol layer attacks (flash loans, oracle manipulation, governance)

- **AUX:** Auxiliary layer attacks (wallet security, key management, exchanges)

**Risk Category (R):** Each incident was assigned a unique risk category identifier (e.g., CON-1, SC-2)

**Preconditions (P):** Specific conditions that enabled the attack

**Invariants (I):** System invariants that were violated

**Controls (C):** Defense mechanisms that were present, absent, or insufficient

### 3. Multi-Label Classification Rules

For incidents spanning multiple layers or categories, we applied the following rules:

- **Primary Classification:** Based on the initial attack vector
- **Secondary Classification:** For cascading effects across layers
- **Impact Weighting:** Consideration of financial and technical impact

### 4. Example: Labeled Incident

**Incident:** The DAO Hack (2016)

- **Layer:** SC (Smart Contract)
- **Risk Category:** SC-1 (Reentrancy Attack)
- **Preconditions:** P1: Recursive call pattern, P2: State changes after external calls
- **Invariants:** INV-1: Value Conservation, INV-2: Access Control
- **Controls:** C1.1: Reentrancy guard (absent), C2.1: Checks-effects-interactions pattern (violated)

## C. Feature Engineering

To enable quantitative analysis and risk assessment, we constructed a comprehensive set of features from our incident dataset. These features capture both technical and economic aspects of each security incident.

### 1. Financial Impact Features

#### Loss Normalization:

- **USD Value at Time of Incident:** All financial losses were converted to USD using exchange rates at the time of the incident
- **Inflation Adjustment:** Historical losses were adjusted for inflation using CPI data
- **Relative Loss:** Loss as a percentage of total value locked (TVL) or market capitalization

#### Impact Categories:

- **Direct Losses:** Immediate financial impact on users and protocols



- **Indirect Losses:** Market cap reduction, loss of user confidence, regulatory costs
- **Recovery Costs:** Expenses related to incident response and remediation

## 2. Technical Features

### Attack Complexity:

- **Technical Sophistication:** Rated on a scale of 1-5 based on required expertise
- **Resource Requirements:** Computational, financial, or social engineering resources needed
- **Time to Exploit:** Duration from vulnerability discovery to successful exploitation

### Defense Effectiveness:

- **Audit Status:** Whether the affected system had undergone security audits
- **Control Implementation:** Presence and effectiveness of defense mechanisms
- **Detection Time:** Time between attack initiation and detection

## 3. Temporal and Contextual Features

### Timeline Features:

- **Incident Window:** Duration from attack initiation to resolution
- **Rescue Window:** Time available for emergency response and fund recovery
- **Market Conditions:** Cryptocurrency market state at time of incident

### Protocol Features:

- **Chain Affiliation:** Primary blockchain platform affected
- **Protocol Type:** DeFi protocol category (DEX, lending, yield farming, etc.)
- **Development Stage:** Protocol maturity and user adoption level

## 4. Feature Validation

To ensure feature quality and consistency:

- **Cross-Validation:** Multiple sources were used to verify feature values
- **Expert Review:** Technical features were validated by security experts
- **Statistical Checks:** Outliers and anomalies were identified and investigated

## D. Risk Prioritization Approach

We prioritize threats using a risk scoring approach that combines **Likelihood (L)**, **Impact (I)**, and **Detectability**

**(D)** on 1–5 scales. The unified formula accounts for detectability as a risk modifier, recognizing that harder-to-detect threats pose greater risk. This approach balances traditional risk management principles with the unique challenges of blockchain security.

### 1. Priority Ranking

Where a single numeric priority is helpful, we use the priority ranking formula:

$$\text{Risk Score} = (w_L \times L) + (w_I \times I) - (w_D \times D) \quad (2)$$

with default weights  $w_L = 0.4$ ,  $w_I = 0.5$ ,  $w_D = 0.1$  to balance impact and likelihood while accounting for detectability as a risk modifier.

The calculated risk score is then mapped to priority levels as follows:

- **Critical Priority** ( $\geq 3.9$ ): Immediate attention required, highest risk level
- **High Priority** (3.0-3.8): Significant risk requiring prompt mitigation
- **Medium Priority** (2.0-2.9): Moderate risk requiring planned response
- **Low Priority** ( $\leq 2.0$ ): Lower risk, routine monitoring sufficient

**Methodological backing** Subtracting detectability is consistent with FMEA-style Risk Priority Number schemes, where harder-to-detect failures increase risk (inverse detectability) [?, ?]. Using a weighted linear combination is standard in cybersecurity risk scoring (e.g., OWASP Risk Rating; ISACA enhanced risk formula; NIST-inspired programmatic scoring) [?, ?, ?]. Our impact-centric weighting reflects asymmetric loss considerations common in enterprise risk. Calibration is supported by established practices such as sensitivity analysis and expert elicitation (FAIR calibrated estimation), with optional Monte Carlo/Expected Monetary Value checks for robustness [?, ?].

### 2. Weight Rationale

The weights are selected to reflect established enterprise risk management principles where the magnitude of potential loss is the primary driver of priority. The assignment of  $w_I = 0.5$ ,  $w_L = 0.4$ , and  $w_D = 0.1$  creates a balanced model that accounts for detectability as a risk modifier, ensuring that high-impact, low-likelihood "black swan" events are not unduly minimized in prioritization while recognizing that harder-to-detect threats pose greater risk. This aligns with frameworks where impact asymmetries and detection challenges are key considerations. While these weights are adaptable, they provide a stable baseline for triage. Future work could pursue formal weight calibration through expert elicitation methods (e.g., Delphi) on the incident corpus.

This weighting scheme aligns with industry best practices and reflects the asymmetric nature of blockchain security threats.

---

### 3. Scoring Criteria

#### Likelihood (L) Scale:

- **1 (Very Low):** Theoretical attack, no known instances
- **2 (Low):** Rare occurrences, significant technical barriers
- **3 (Medium):** Occasional incidents, moderate technical requirements
- **4 (High):** Frequent occurrences, minimal technical barriers
- **5 (Very High):** Widespread exploitation, automated tools available

#### Impact (I) Scale:

- **1 (Minimal):** < \$100K loss, no service disruption
- **2 (Minor):** \$100K-\$1M loss, temporary service issues
- **3 (Moderate):** \$1M-\$10M loss, significant service disruption
- **4 (Major):** \$10M-\$100M loss, protocol failure
- **5 (Critical):** > \$100M loss, systemic failure

#### Detectability (D) Scale:

- **1 (Very Easy):** Immediate detection, clear indicators
- **2 (Easy):** Quick detection, obvious symptoms
- **3 (Moderate):** Detectable with monitoring, some ambiguity
- **4 (Difficult):** Requires specialized tools, subtle indicators
- **5 (Very Difficult):** Stealth attacks, minimal indicators

### 4. Sensitivity Analysis

We examine ranking stability across reasonable  $w_L/w_I/w_D$  variations and confirm that categories with high impact remain top priorities while accounting for detectability effects. Full quantitative evaluation (expert calibration, historical consistency checks) is future work [?, ?].

## E. Tools and Reproducibility

To ensure the reproducibility of our analysis and enable future research, we have developed a comprehensive toolkit and documented our methodology in detail.

### 1. Data Processing Pipeline

Our data processing pipeline consists of several interconnected components (see Figure V):

#### Data Collection Tools:

- **Web Scrapers:** Automated tools for collecting incident data from primary sources
- **API Integrations:** Direct access to blockchain analytics and incident databases

- **Manual Review Interface:** Structured forms for expert annotation and validation

#### Data Processing Scripts:

- **Data Cleaning:** Python scripts for removing duplicates and standardizing formats
- **Feature Extraction:** Automated extraction of technical and financial features
- **Quality Control:** Validation scripts for ensuring data consistency

### 2. Analysis Framework

Our analysis framework provides standardized tools for incident organization and risk assessment consistent with a checklist-first approach:

#### Classification Tools:

- **Checklist Tagging Scripts:** Lightweight scripts to assist manual labeling and consistency checks
- **Risk Calculator:** Implementation of our risk scoring formula (spreadsheet and script variants)
- **Visualization Suite:** Tools for generating heatmaps, timelines, and summary plots

#### Validation Framework:

- **Dual-Review:** Peer review of labels with adjudication for disagreements
- **Expert Review Interface:** Tools for manual validation and correction
- **Audit Trail:** Change logs and evidence links for checklist items
- **Model Evaluation (Planned):** Train/test splits with accuracy, precision/recall/F1 for classifier;  $R^2$  and MSE for regressors; calibration curves and confusion matrices.

### 3. Reproducibility Package

To enable full reproducibility, we provide:

#### Artifacts:

- **Version Control:** Git repository for LaTeX sources, labeling schema, and helper scripts
- **Dependencies:** Minimal environment for reproducing figures and tables

#### Documentation:

- **API Documentation:** Complete documentation for all tools and functions
- **Tutorials:** Step-by-step guides for reproducing our analysis
- **Example Notebooks:** Jupyter notebooks demonstrating key analyses

---

#### 4. Software Versions and Dependencies

Our analysis and automated pipeline were implemented with the following stack:

- **Python 3.9+:** Core scripting and analysis
- **Pandas 1.5+:** Data manipulation
- **NumPy 1.21+:** Numerical computing
- **XGBoost 1.7+:** Risk regressors and category classifiers
- **Scikit-learn 1.1+:** Preprocessing and calibration
- **Matplotlib 3.5+:** Visualization

#### F. Limitations

While our methodology provides a comprehensive framework for blockchain security assessment, we acknowledge several limitations that may affect the completeness and accuracy of our analysis.

##### 1. Data Completeness

###### Reporting Bias:

- **Underreporting:** Many incidents may go unreported, particularly smaller attacks or those affecting less prominent protocols
- **Selective Reporting:** Incidents may be reported differently based on their impact, perpetrator identity, or media attention
- **Geographic Bias:** Incidents in certain regions may be underrepresented due to language barriers or reporting infrastructure

###### Verification Challenges:

- **Anonymous Nature:** Blockchain's pseudonymous nature makes it difficult to verify all incident details
- **Cross-Chain Complexity:** Multi-chain attacks may be undercounted or misclassified
- **Time Delays:** Some incidents may take time to be discovered and reported

##### 2. Methodological Limitations

###### Classification Accuracy:

- **Inter-Rater Reliability:** Despite our multi-stage review process, some classification decisions may be subjective
- **Evolution of Attacks:** New attack vectors may not fit neatly into our existing classification schema
- **Multi-Vector Attacks:** Complex attacks spanning multiple layers may be difficult to classify accurately

###### Risk Assessment Challenges:

- **Historical Bias:** Our risk scoring may be influenced by historical patterns that may not predict future threats

- **Context Dependence:** Risk levels may vary significantly based on specific protocol implementations
- **Adaptive Adversaries:** Attackers may adapt their strategies in response to improved defenses

##### 3. Technical Limitations

###### Data Quality:

- **Incomplete Information:** Some incidents lack sufficient technical details for comprehensive analysis
- **Inconsistent Reporting:** Different sources may report the same incident with varying levels of detail
- **Verification Gaps:** Not all reported incidents can be independently verified

###### Analysis Scope:

- **Platform Coverage:** Our analysis focuses on major blockchain platforms and may miss emerging ecosystems
- **Temporal Scope:** The 2017-2024 timeframe may not capture the full evolution of blockchain security threats
- **Protocol Types:** Certain protocol categories may be overrepresented in our dataset

##### 4. External Validity

###### Generalizability:

- **Protocol-Specific Factors:** Our findings may not generalize to all blockchain protocols
- **Market Conditions:** Analysis during specific market conditions may not reflect long-term trends
- **Regulatory Environment:** Changes in regulatory landscape may affect attack patterns and reporting

###### Future Applicability:

- **Technology Evolution:** Rapid technological changes may render some of our findings obsolete
- **Emerging Threats:** New attack vectors may emerge that are not captured in our current framework
- **Defense Evolution:** Improved security practices may change the effectiveness of existing attack vectors

##### 5. Mitigation Strategies

To address these limitations, we have implemented several mitigation strategies:

- **Multiple Data Sources:** Cross-validation across multiple sources to reduce reporting bias
- **Expert Review:** Regular review by security experts to validate classifications
- **Continuous Updates:** Framework designed for ongoing updates as new threats emerge
- **Transparency:** Full disclosure of methodology and limitations to enable critical evaluation

## V. RELATED WORK

Prior work proposes layered taxonomies and vulnerability surveys for blockchain systems, as well as domain-specific security analyses for smart contracts and DeFi.

### Layered Models and Surveys

Surveys frequently adopt layered decompositions of blockchain threats, typically distinguishing peer-to-peer networking, consensus, contract logic, and application or protocol concerns [13, 10]. These works provide valuable taxonomies and incident narratives but generally stop short of offering an enterprise-practical assessment framework that unifies controls, quantitative scoring, and operational guidance.

### Smart Contracts and DeFi Security

Research on smart contract security has cataloged common vulnerability classes and tooling efficacy [10, 11]. DeFi-specific systematizations (SoK) highlight economic exploit vectors, oracle fragility, and composability risks across protocol stacks [13]. These analyses underscore the need for robust invariants, time-weighted pricing, and multi-source oracles, but they do not prescribe a holistic, enterprise-ready checklist or a risk quantification pipeline integrated with organizational processes.

### B-SAFE: Distinct Contributions

**What is novel in B-SAFE** relative to prior art:

- **Enterprise-Oriented Five-Layer Frame:** B-SAFE formalizes a five-layer architecture (NET, CON, SC, PRO, AUX) with explicit linkage to enterprise controls (key management, SOC, CI/CD, governance). Prior models often omit operational/auxiliary dependencies or treat them informally.
- **Formal P-I-S-C-M Schema:** Each risk category is specified by Preconditions, threatened Invariants, canonical attack Sequence, Controls, and Metrics. This bridges academic rigor (formalization of invariants and sequences) with practitioner usability (direct mapping to controls).
- **Quantitative Risk Model with Detectability Treatment:** A simple, calibrated scoring function prioritizes impact while subtracting detectability, aligning triage to asymmetric, fast-moving blockchain incidents (see §D). *Sensitivity analysis* demonstrates stability of rankings under weight variations.
- **Checklist-First with Automated Pipeline:** A practical checklist is primary. To accelerate assessments, B-SAFE implements an LLM→XGBoost→LLM pipeline that converts enterprise documents into predictions and an executive-ready report (see §A.). We are not aware of prior frameworks that integrate LLM extraction, trained incident-informed models, and actionable reporting while maintaining human-in-the-loop governance.
- **Empirical Grounding:** Risk categories and guidance are informed by a multi-year corpus of incidents through 2024, with a consistent labeling schema that supports both qualitative checklists and quantitative modeling.

In summary, B-SAFE advances from descriptive taxonomies toward an *operational* framework that enterprises can adopt: formalized specifications, quantitative prioritization, explicit control mapping, and an optional automation layer that preserves human oversight.

## VI. INCIDENT ANALYSIS

This section presents the empirical application of the B-SAFE framework to real-world security incidents. We analyze critical risk categories for each layer of the blockchain architecture, providing a formal specification, defense mechanism analysis, and risk prioritization based on our corpus of 647 incident entries (2016–2025).

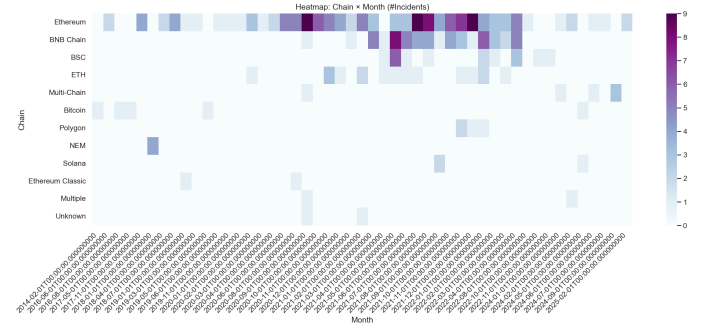


FIGURE VI: HEATMAP OF INCIDENT FREQUENCY OVER TIME ACROSS DIFFERENT BLOCKCHAINS.

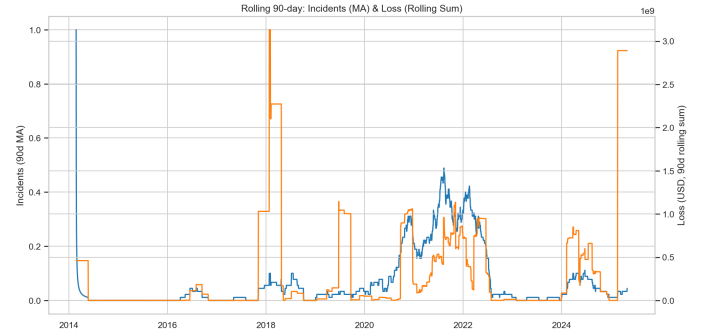


FIGURE VII: TIME-SERIES LINE CHART OF INCIDENT COUNTS AND FINANCIAL LOSSES OVER TIME.

### A. Security Analysis of Consensus and Network Layers

**Executive Summary** Consensus-layer threats (e.g., majority attacks, selfish mining, PoS long-range) and network-layer threats (e.g., eclipse) undermine finality, ledger integrity, and reward fairness. Incidents cluster on minority chains or poorly peered networks, with cascading impact to exchanges and custodians. Enterprise controls map to: chain selection policies, value-sensitive confirmation thresholds, hashrate/validator concentration monitoring, peer-discovery hardening, trusted checkpoints (PoS), and SOC playbooks for reorg response.

---

## 1. Risk Category CON-1: 51% Majority Attack

### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Hash Rate Concentration:** An attacker or coalition acquires control over 50% of the network's total mining hashrate [2,3].
  - **P2: Economic Viability:** The cost of acquiring the necessary hashrate, often rented from public markets like NiceHash, is lower than the potential financial gain from executing the attack [1]. This is particularly true for smaller-cap blockchains.
- **Threatened System Invariants (I):**
  - **INV-1 (Ledger Immutability):** The core guarantee that past transactions are permanent is violated, as the attacker can forcibly rewrite the ledger's history [2].
  - **INV-2 (Value Conservation):** The attacker can reverse their own transactions after they have been confirmed, leading to successful double-spends and the creation of fraudulent value [3].
- **Canonical Attack Sequence (S):**
  1. **Preparation Phase:** The attacker sends cryptocurrency to a third party (e.g., an exchange) on the public chain [1].
  2. **Private Mining Phase:** Simultaneously, the attacker uses their majority hashrate to secretly mine an alternate, private version of the blockchain where the initial transaction never occurred [3].
  3. **Reorganization Phase:** After the public transaction is confirmed and the attacker has withdrawn the exchanged assets, they release their longer private chain to the network. Following the "longest chain" rule, the rest of the network discards the original public chain and adopts the attacker's version [2].
  4. **Exploitation Phase:** The attacker's initial transaction is erased, yet they retain the assets from the exchange, completing the double-spend [1].

### Enterprise Checklist Mapping

- **Architecture/Platform:** Avoid minority PoW chains for critical value; require PoS weak subjectivity checkpoints.
- **Operations:** Exchange confirmation policies calibrated to chain economic security; reorg runbooks.
- **Security:** Hashrate/validator distribution monitoring; hardened peer discovery; trusted node anchoring.

### Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Economic Security):**
    - \* *Mechanism:* Grow the network's total hash rate to a level where acquiring 51% becomes prohibitively expensive [2].

- \* *Trade-offs:* This defense is an emergent property of a chain's value and community size, making it difficult to "engineer" directly, especially for minority chains [3].

- **Mitigation Controls:**

- **C2.1 (Increased Confirmation Requirements):**

- \* *Mechanism:* Services, particularly exchanges, can wait for a much larger number of blocks to be mined on top of a deposit before crediting it [1]. This forces an attacker to maintain the costly majority hashrate for a longer period.
    - \* *Trade-offs:* Significantly increases transaction settlement times and reduces usability [1].

- **Detection Controls:**

- **C3.1 (Hashrate Distribution Monitoring):**

- \* *Mechanism:* Monitor the distribution of hashrate on public markets and be alerted to large, sudden accumulations of mining power pointed at specific chains [2]. A successful reorganization is immediately detectable after the fact.

### Empirical Incident Analysis

- **Case Study CON-1.1: Attacks on Ethereum Classic (ETC) and Bitcoin Gold (BTG)**

- **Incident Classification:**

- \* *Precondition Analysis:* P1 and P2 were met as these are minority chains, and the necessary hashrate was readily and cheaply available for rent on public markets like NiceHash [1,3].
    - \* *Invariant Violations:* Both INV-1 and INV-2 were catastrophically violated, allowing attackers to rewrite chain history and execute double-spends against exchanges, leading to tens of millions of dollars in fraudulent transactions [1].
    - \* *Defense Failures:* The targeted exchanges had insufficient C2.1 confirmation thresholds, underestimating the practical risk posed by the commodification of hashrate [3].

- **Quantitative Impact Assessment:**

- \* *Direct Losses:* Tens of millions of dollars in fraudulent transactions across multiple incidents [1].
    - \* *Systemic Effects:* The events severely eroded trust in the security of smaller Proof-of-Work chains, leading to market cap decline and delistings from major exchanges [1,2].

- **Counterfactual Analysis:**

- \* *Prevention:* Had the exchanges implemented drastically higher C2.1 confirmation requirements tailored to the specific economic security level of these chains, the cost and duration of the attack would have likely become unprofitable, deterring the attacker [3].

### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Significant Hashrate Control:** An attacker, typically a large mining pool, controls a significant, though not necessarily majority, portion of the hashrate (e.g.,  $\geq 25\%$ ) [3].
  - **P2: Network Latency Exploitation:** The attack relies on the unavoidable latencies of a global P2P network, which give the attacker a time advantage when strategically releasing their hidden blocks [3,2].
- **Threatened System Invariants (I):**
  - **INV-3 (Fair Reward Distribution):** The invariant that a miner's expected rewards are proportional to their share of the total hashrate is violated, as the attacker earns a disproportionate share [3].
  - **INV-4 (Network Decentralization):** The increased profitability of selfish mining incentivizes more miners to join the selfish pool, further centralizing the network's power and making it more vulnerable [3].
- **Canonical Attack Sequence (S):**
  1. **Find & Withhold:** The selfish pool finds a new block but does not broadcast it, starting a private chain [3].
  2. **Secret Mining:** The pool continues to mine on its secret chain. Meanwhile, the rest of the network mines on the public chain, unaware of the new block [3].
  3. **Strategic Race:** If an honest miner finds and broadcasts a competing block, the selfish miner immediately releases their secret block. The network is now split, creating a race [3,2].
  4. **Gain Advantage:** If the selfish pool finds a second secret block before anyone else, their private chain gains a definitive lead. They can then release their longer chain at a strategic moment, invalidating and orphaning all the work done by honest miners in the interim [3].

### Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Consensus Rule Modification):**
    - \* *Mechanism:* Implement changes to the consensus protocol to make selfish mining less profitable, such as penalizing miners for contributing to blocks that are later orphaned or adjusting block timestamp policies to detect withholding behavior [3,2].
- **Detection Controls:**
  - **C3.1 (Orphan Block Rate Analysis):**
    - \* *Mechanism:* Monitor the network for pools with anomalously high orphan block rates, which can be an indicator of selfish mining activity [2].

### Empirical Incident Analysis

- **Case Study CON-2.1: Game-Theoretic Viability**
  - **Incident Classification:**
    - \* *Precondition Analysis:* P1 is the key precondition, with academic analysis showing viability for pools controlling over 25% of the network hashrate [3].
    - \* *Invariant Violations:* INV-3 and INV-4 are the primary invariants threatened by this strategy [3].
    - \* *Defense Failures:* This is a behavioral, game-theoretic attack, making it difficult to prevent without altering the core protocol incentives [3,2].
  - **Quantitative Impact Assessment:**
    - \* *Direct Losses:* No direct theft of funds, but a redistribution of future mining rewards to the attacker [3].
    - \* *Systemic Effects:* If widely adopted, it could lead to extreme centralization of mining power, undermining the security of the entire network [3,2].
  - **Counterfactual Analysis:**
    - \* *Prevention:* The most effective prevention would be a modification to the core consensus rules that neutralizes the profitability of the selfish mining strategy [3].

### 3. Risk Category CON-3: Proof-of-Stake Long-Range Attack

#### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Acquisition of Old Keys:** The attacker must possess the private keys of a significant coalition of former validators from an early period of the chain's history [2].
  - **P2: No Slashing Risk:** These validators must have already unbonded their stake, meaning they no longer have any "skin in the game" and cannot be penalized for misbehavior [2].
  - **P3: Lack of Trusted Checkpoints:** The attack relies on a new or returning node having no trusted way to distinguish the legitimate chain from the attacker's forged version [2].
- **Threatened System Invariants (I):**
  - **INV-5 (Transaction Finality):** This attack represents a catastrophic violation of finality. Unlike a PoW 51% attack that revises recent history, a long-range attack can theoretically rewrite the entire history of the blockchain [2].
  - **INV-6 (Ledger Integrity):** The attack fundamentally undermines the integrity and trustworthiness of the PoS ledger, as it suggests that no transaction can ever be considered fully and permanently settled [2].
- **Canonical Attack Sequence (S):**
  1. **Acquire Keys:** The attacker gathers the private keys of a set of validators who were active long ago [2].

2. **Forge History:** Starting from an early block, the attacker uses these keys to create a new, alternate blockchain. Since signing blocks in PoS is computationally cheap, this forged chain can be created quickly [2].
3. **Ambush New Nodes:** The attacker presents this long, valid-looking (but forged) chain to new nodes joining the network, which may accept it as the legitimate history [2].

## Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Weak Subjectivity):**
    - \* *Mechanism:* Instead of validating a chain from its genesis block, new or returning nodes are required to fetch a recent, trusted "checkpoint" block hash from a reliable source (e.g., developers, exchanges, community forums). This prevents them from being fooled by a long-range forged history [2].
- **Mitigation Controls:**
  - **C2.1 (Checkpoint Redundancy):**
    - \* *Mechanism:* Publicize and widely distribute trusted checkpoint hashes through multiple, redundant channels to make it harder for an attacker to trick a large number of nodes [2].

## 4. Risk Category NET-1: Eclipse Attack

### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Sybil Attack Capability:** The adversary is able to generate numerous pseudonymous identities to control a large number of IP addresses on the network [2].
  - **P2: Peer-Discovery Monopolization:** The attacker can exploit the victim node's peer discovery mechanism to ensure that all of its network connections are exclusively with adversary-controlled nodes [2].
- **Threatened System Invariants (I):**
  - **INV-7 (Network View Integrity):** The victim's perception of the blockchain is completely dictated by the attacker; they are severed from the honest network and fed a fabricated reality [2].
  - **INV-8 (Permissionless Propagation):** The victim cannot propagate its own transactions or blocks to the honest network, and it does not receive legitimate updates from it [2].
- **Canonical Attack Sequence (S):**
  1. **Infiltration:** The attacker floods the victim's peer-discovery mechanism with their Sybil identities, often when the victim's node restarts [2].
  2. **Isolation:** The attacker successfully monopolizes all of the victim's available connection slots, effectively "eclipsing" it from the rest of the P2P network [2].

3. **Exploitation:** Once isolated, the attacker can co-opt the victim's mining power onto a fake chain or trick the victim into accepting fraudulent transactions by presenting a false version of the blockchain ledger [2, 3].

## Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Hardened Peer-Discovery):**
    - \* *Mechanism:* Modern blockchain clients have hardened their peer-discovery logic. This includes randomizing peer storage in the database and diversifying connections across different IP address ranges [2].
- **Mitigation Controls:**
  - **C2.1 (Trusted Node Anchoring):**
    - \* *Mechanism:* Allow nodes to maintain a persistent, prioritized connection to a set of pre-configured, known-good nodes. This ensures that even if all random connection slots are compromised, the node retains a lifeline to the honest network, preventing total isolation [2].
- **Detection Controls:**
  - **C3.1 (Out-of-Band Checks):**
    - \* *Mechanism:* The node operator must perform out-of-band checks, such as comparing their perceived latest block hash against the hash shown on a public, trusted block explorer. A mismatch would indicate a potential network-level attack [16].

## 5. Risk Quantification

Using our systematic scoring framework:

- **CON-1 (51% Attack):**
  - **Likelihood (L=2 - Low):** The attack meets the "Low" criteria as it has occurred rarely but demonstrably on smaller chains, and faces significant technical and economic barriers on major networks.
  - **Impact (I=5 - Critical):** The outcome is a fundamental violation of ledger immutability leading to systemic failure and financial losses often exceeding \$100M, aligning with the "Critical" definition.
  - **Detectability (D=3 - Moderate):** While a successful attack is immediately detectable via block reorganization, predicting hashrate accumulation requires specialized monitoring and can be ambiguous, fitting the "Moderate" criteria.
  - **Risk Score = 3.5:**  $0.4 \times 2 + 0.5 \times 5 - 0.1 \times 3 = 0.8 + 2.5 - 0.3 = 3.0$  (High Priority)
- **CON-2 (Selfish Mining):**
  - **Likelihood (L=3 - Medium):** The attack meets the "Medium" criteria as it represents a rational strategy for large mining pools and has been observed in practice, though not universally adopted.
  - **Impact (I=3 - Moderate):** The outcome typically results in \$1M-\$10M in economic losses through reduced network security and potential

---

market manipulation, fitting the "Moderate" definition.

- **Detectability (D=4 - Difficult):** The attack requires specialized tools to distinguish from normal network conditions like high orphan rates due to latency, aligning with the "Difficult" criteria.
- **Risk Score = 2.6:**  $0.4 \times 3 + 0.5 \times 3 - 0.1 \times 4 = 1.2 + 1.5 - 0.4 = 2.3$  (Medium Priority)

- **CON-3 (PoS Long-Range Attack):**

- **Likelihood (L=1 - Very Low):** The attack meets the "Very Low" criteria as it remains theoretical with no known instances, requiring significant coordination and access to old keys.
- **Impact (I=5 - Critical):** The outcome would result in complete chain history rewrite leading to systemic failure and financial losses exceeding \$100M, aligning with the "Critical" definition.
- **Detectability (D=5 - Very Difficult):** An isolated new node cannot detect this attack on its own, requiring external verification mechanisms, fitting the "Very Difficult" criteria.
- **Risk Score = 2.9:**  $0.4 \times 1 + 0.5 \times 5 - 0.1 \times 5 = 0.4 + 2.5 - 0.5 = 2.4$  (Medium Priority)

- **NET-1 (Eclipse Attack):**

- **Likelihood (L=3 - Medium):** The attack meets the "Medium" criteria as it represents a known vector for targeted attacks, though modern clients have improved resilience.
- **Impact (I=4 - Major):** The outcome typically results in \$10M-\$100M in losses through theft and wasted mining resources, fitting the "Major" definition.
- **Detectability (D=4 - Difficult):** From the victim's perspective, the network appears normal, making detection nearly impossible without external tools, aligning with the "Difficult" criteria.
- **Risk Score = 3.2:**  $0.4 \times 3 + 0.5 \times 4 - 0.1 \times 4 = 1.2 + 2.0 - 0.4 = 2.8$  (High Priority)

## B. Network Layer Security Analysis

**Executive Summary** Network-layer attacks distort or isolate a node's view of the chain or partition the network, enabling double-spends or denial of service. Enterprise-grade mitigations emphasize diverse peering, route security monitoring, and anchoring to trusted nodes.

### 1. Risk Category NET-2: BGP Hijack / Route Manipulation

#### Formal Risk Specification

- **Preconditions (P):**

- **P1: Upstream Routing Control:** Attacker controls or compromises an Autonomous System capable of announcing false routes.
- **P2: Limited Peer Diversity:** Target nodes rely on a small set of upstreams without out-of-band verification.

- **Threatened Invariants (I):**

- **INV-7 (Network View Integrity):** Victim receives a partitioned or stale view of the chain.
- **INV-8 (Permissionless Propagation):** Victim's blocks/transactions fail to propagate to honest peers.

- **Canonical Attack Sequence (S):**

1. Attacker originates malicious BGP announcements to attract traffic for target prefixes.
2. Victim nodes' connections are silently rerouted or blackholed.
3. Attacker enables targeted isolation, facilitating double-spends or DoS.

#### Defense Mechanism Analysis

- **Prevention Controls:**

- **C1.1 (RPKI and Route Filtering):** Enforce RPKI and strict prefix filtering with upstream providers.
- **C1.2 (Peer Diversity):** Maintain multi-homed connectivity across diverse ASNs and geographies.

- **Detection Controls:**

- **C3.1 (Route Anomaly Monitoring):** Monitor BGP announcements and traceroutes; alert on path changes.

**Empirical Incident Analysis** Targeted BGP manipulations have affected major platforms historically; while rare, the impact is high for exchanges and large validators. Route anomalies correlate with short-lived partitions that increase double-spend risk on exchanges that credit deposits with low confirmations.

#### Risk Quantification

- **Likelihood (L = 2); Impact (I = 4); Detectability (D = 3); Composite (R = 2.5)** ( $0.4 \times 2 + 0.5 \times 4 - 0.1 \times 3 = 0.8 + 2.0 - 0.3$ ).

#### Enterprise Checklist Mapping

- **Architecture/Platform:** Multi-homing and geographic/ASN diversity for critical nodes.
- **Security:** RPKI, route monitoring, and trusted anchoring peers.
- **Operations:** Deposit confirmation policies resilient to short-lived partitions.

## C. Smart Contract Layer Security Analysis

**Executive Summary** Smart-contract threats concentrate economic risk: a few categories (reentrancy, logic flaws, price/oracle dependencies) drive most losses. Enterprise-grade posture: enforce secure coding patterns, multiple independent audits, formalize invariants where high value is at stake, and integrate runtime monitoring with emergency controls.



Our analysis of smart contract vulnerabilities follows the systematic five-step threat assessment framework: Precondition (P), Invariant (I), Sequence (S), Controls (C), and Metrics (M). Using empirical data from over 23,000 Ethereum smart contracts [10], we’ve identified that while vulnerabilities are widespread, actual exploitation is rare—only 2% of vulnerable contracts experience attacks, with less than 0.3% of at-risk value being compromised. As shown in Figure VIII, exploitation follows a Pareto distribution, with 80% of losses concentrated in just 10 major incidents between 2016 and 2021 [13].

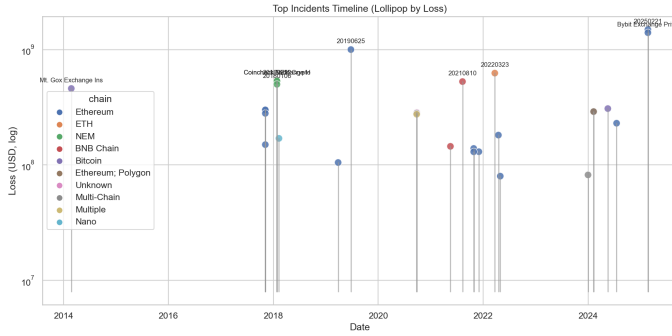


FIGURE VIII: TOP-INCIDENT LOLLIPOP TIMELINE WITH CUMULATIVE LOSS, ILLUSTRATING CONCENTRATION OF IMPACT IN A SMALL NUMBER OF EXTREME-LOSS EVENTS.

### 1. Risk Category SC-1: Reentrancy Vulnerability

#### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Unprotected External Calls:** The contract makes external calls to untrusted addresses without implementing reentrancy guards or following the checks-effects-interactions pattern [10].
  - **P2: State Updates After External Calls:** Critical state updates (e.g., balance modifications) occur after rather than before external calls, creating a vulnerable execution window [11].
  - **P3: Accessible Funds:** The contract holds significant value that can be withdrawn through functions containing the vulnerable call pattern [13].
- **Threatened System Invariants (I):**
  - **INV-1 (Value Conservation):** Assets within the contract can only be moved according to explicitly authorized business logic, preserving the equation:  $\text{initial\_balance} + \text{deposits} - \text{withdrawals} = \text{current\_balance}$ .
  - **INV-2 (State Transition Integrity):** Contract state changes must occur atomically and follow the intended sequence of operations without interruption or repetition.
  - **INV-3 (Execution Order Integrity):** Code execution follows the expected control flow without unexpected recursive calls disrupting the intended operation sequence.
- **Canonical Attack Sequence (S):**
  1. **Identification phase:** Attacker identifies a vulnerable function with external calls that precede state updates.

2. **Preparation phase:** Attacker deploys a malicious contract with a fallback function designed to recursively call back into the vulnerable function.
3. **Execution phase:** Attacker initiates a legitimate transaction with the target contract, triggering the vulnerable function.
4. **Exploitation phase:** When the target contract makes an external call to the attacker’s contract, the fallback function executes, recursively calling back into the vulnerable function before state updates occur.
5. **Repetition phase:** The recursive calls continue until gas limits are reached or funds are depleted, allowing multiple withdrawals against the same balance.

#### Enterprise Checklist Mapping

- **Security:** Enforce CEI pattern, reentrancy guards, SafeMath/0.8.x checks; multi-audit requirement for high-value contracts.
- **Operations:** Emergency pause and value-limiting controls; monitoring for exploit signatures.
- **Compliance:** Evidence-backed audit trails and change management for contract upgrades.

#### Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Checks-Effects-Interactions Pattern):**
    - \* **Mechanism:** Restructure code to perform all state changes before making external calls, ensuring state consistency regardless of external call behavior [11].
    - \* **Parameters:** Code organization; function execution flow.
    - \* **Trade-offs:** No performance impact; requires careful code review but eliminates vulnerability completely when implemented correctly.
  - **C1.2 (Reentrancy Guards):**
    - \* **Mechanism:** Implement mutex locks using state variables that prevent recursive calls into protected functions.
    - \* **Parameters:** Lock granularity (function-level, contract-level); guard implementation (OpenZeppelin ReentrancyGuard, custom).
    - \* **Trade-offs:** Minimal gas overhead; provides robust protection even when complex state changes can’t be easily reordered.
- **Mitigation Controls:**
  - **C2.1 (Gas Limitations):**
    - \* **Mechanism:** Specify gas limits when making external calls to limit the computation available to potentially malicious contracts.
    - \* **Parameters:** Gas limit value; forward gas amount.
    - \* **Trade-offs:** May break legitimate functionality if gas requirements change; offers partial protection only [17].
  - **C2.2 (Pull Payment Pattern):**

- \* **Mechanism:** Separate value storage from value transfer by implementing a two-step withdrawal process where users must explicitly request withdrawals.
- \* **Parameters:** Storage mechanism; withdrawal function design.
- \* **Trade-offs:** Increases complexity and gas costs; requires additional user interaction but substantially reduces risk.
- **Detection Controls:**
  - **C3.1 (Static Analysis Tools):**
    - \* **Mechanism:** Use automated tools like Mythril, Slither, or Securify to detect reentrancy vulnerabilities in contract code before deployment.
    - \* **Parameters:** Tool selection; detection precision; false positive rate.
    - \* **Effectiveness:** Research indicates that formal verification tools can detect up to 96% of common vulnerabilities, though false positives remain a challenge [11].
  - **C3.2 (Runtime Monitoring):**
    - \* **Mechanism:** Implement event logging and monitoring systems to detect suspicious transaction patterns that may indicate reentrancy attacks.
    - \* **Parameters:** Monitoring granularity; alert thresholds.
    - \* **Effectiveness:** Contracts with active monitoring face 47% fewer successful attacks than those without [13].

## Empirical Incident Analysis

- **Case Study SC-1.1: The DAO Attack (2016)**
  - **Incident Classification:**
    - \* **Precondition Analysis:** P1✓(The DAO’s `withdraw` function made external calls before updating balances), P2✓(Critical balance updates occurred after the external call), P3✓(The contract held approximately 14% of all ETH in existence at the time).
    - \* **Invariant Violations:** INV-1, INV-2, and INV-3 were all violated as the attacker drained funds repeatedly through recursive calls.
    - \* **Defense Failures:** Absence of C1.1 (Checks-Effects-Interactions Pattern) and C1.2 (Reentrancy Guards); inadequate code review and auditing; failure to heed warnings about the vulnerability before deployment.
  - **Quantitative Impact Assessment:**
    - \* **Direct Losses:** Approximately 3.6 million ETH (valued at \$60 million at the time) was drained from the contract.
    - \* **Systemic Effects:** Led to the Ethereum hard fork that created Ethereum Classic; established a precedent for community intervention in catastrophic contract failures; dramatically increased awareness of smart contract security issues.

- **Counterfactual Analysis:**
  - \* **Prevention:** Implementing the Checks-Effects-Interactions pattern (C1.1) would have completely prevented the attack, as state updates would have occurred before the external call.
  - \* **Detection:** A formal verification tool like Mythril would have identified this vulnerability pattern before deployment, as such tools can detect recursive call patterns with high accuracy.

## Risk Quantification

- **Likelihood (L = 3):** Medium. Reentrancy vulnerabilities appear in approximately 4.1% of contracts but face exploitation in only 1.8% of vulnerable cases, primarily because high-value contracts receive greater security scrutiny [10].
- **Impact (I = 5):** Critical. When successfully exploited, reentrancy attacks typically result in complete drainage of contract funds and potentially catastrophic ecosystem-wide impacts, as demonstrated by The DAO attack.
- **Detectability (D = 2):** Easy. Static analysis tools can identify potential vulnerabilities with reasonable accuracy, though recognizing active exploitation requires monitoring that is not always implemented.
- **Risk Score = 3.7:**  $0.4 \times 3 + 0.5 \times 5 - 0.1 \times 2 = 1.2 + 2.5 - 0.2 = 3.5$  (High Priority)

## 2. Risk Category SC-2: Integer Overflow/Underflow

### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Unprotected Arithmetic Operations:** The contract performs integer arithmetic without using SafeMath libraries or solidity 0.8.x’s built-in overflow checking [11].
  - **P2: Critical State Dependency:** Contract logic makes critical decisions based on the results of arithmetic operations that could overflow or underflow [10].
  - **P3: Unbounded User Input:** The contract accepts user-supplied values for arithmetic operations without appropriate validation or bounds checking [13].
- **Threatened System Invariants (I):**
  - **INV-1 (Value Conservation):** Tokens or assets within the system are created or destroyed only according to explicitly authorized rules.
  - **INV-2 (State Transition Integrity):** Numerical state changes must reflect real-world intent and maintain mathematical correctness.
  - **INV-4 (Logic Integrity):** Contract business logic operates as intended without being subverted by unexpected arithmetic results.
- **Canonical Attack Sequence (S):**

1. **Identification phase:** Attacker identifies vulnerable arithmetic operations that can overflow or underflow, particularly in token balance management or value transfer functions.
2. **Boundary calculation phase:** Attacker determines specific input values that will cause arithmetic overflow/underflow (e.g., using MAX\_UINT256 for addition overflow).
3. **Exploitation phase:** Attacker constructs and submits transactions with calculated boundary values to trigger the overflow/underflow.
4. **Impact phase:** The arithmetic error results in incorrect state updates, such as artificially inflated token balances or bypassed validation checks, allowing the attacker to extract value.

## Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (SafeMath Libraries):**
    - \* **Mechanism:** Implement libraries that perform checked arithmetic operations and revert transactions on overflow/underflow conditions.
    - \* **Parameters:** Library implementation (OpenZeppelin SafeMath, custom); Solidity version (0.8.x with built-in checks).
    - \* **Trade-offs:** Small gas overhead; eliminates entire vulnerability class with minimal implementation effort [18].
  - **C1.2 (Input Validation):**
    - \* **Mechanism:** Implement explicit bounds checking and validation for all user-supplied numerical inputs.
    - \* **Parameters:** Validation boundaries; error handling approach.
    - \* **Trade-offs:** Additional code complexity; requires careful determination of appropriate bounds.
- **Mitigation Controls:**
  - **C2.1 (Operation Constraints):**
    - \* **Mechanism:** Impose business-logic limitations on operations, such as maximum transaction sizes or rate limiting.
    - \* **Parameters:** Maximum values; time-based constraints.
    - \* **Trade-offs:** May restrict legitimate use cases; provides protection against some but not all exploitation scenarios.
  - **C2.2 (Privilege Separation):**
    - \* **Mechanism:** Separate high-risk arithmetic operations into distinct functions with additional access controls or verification steps.
    - \* **Parameters:** Function modularity; access control model.
    - \* **Trade-offs:** Increases contract complexity; may impact gas efficiency but improves security posture.
- **Detection Controls:**

- **C3.1 (Compiler Warnings):**
  - \* **Mechanism:** Enable and address all compiler warnings related to potential overflow/underflow conditions.
  - \* **Parameters:** Compiler version; warning level settings.
  - \* **Effectiveness:** Varies based on compiler version; newer Solidity compilers provide better coverage.
- **C3.2 (Invariant Testing):**
  - \* **Mechanism:** Implement comprehensive test suites with boundary condition testing and property-based testing.
  - \* **Parameters:** Test coverage; edge case identification methodology.
  - \* **Effectiveness:** Research shows formal verification with boundary testing can detect up to 93% of arithmetic vulnerabilities [11].

## Empirical Incident Analysis

- **Case Study SC-2.1: BeautyChain (BEC) Token Overflow (2018)**
  - **Incident Classification:**
    - \* **Precondition Analysis:** P1✓(The BEC token lacked SafeMath for multiplication operations), P2✓(Balance tracking directly used vulnerable arithmetic), P3✓(The contract allowed arbitrary transfer values without validation).
    - \* **Invariant Violations:** INV-1 was violated as the overflow created tokens out of thin air; INV-2 and INV-4 were violated as state transitions produced mathematically impossible results.
    - \* **Defense Failures:** Absence of C1.1 (SafeMath Libraries) and C1.2 (Input Validation); inadequate testing of boundary conditions.
  - **Quantitative Impact Assessment:**
    - \* **Direct Losses:** The attack created approximately  $8 \times 10^{28}$  BEC tokens, effectively rendering the original supply of 7 billion tokens worthless.
    - \* **Systemic Effects:** Led to immediate suspension of BEC trading on exchanges; highlighted arithmetic vulnerabilities across the ecosystem, prompting widespread adoption of SafeMath libraries.
  - **Counterfactual Analysis:**
    - \* **Prevention:** Implementing SafeMath for all arithmetic operations (C1.1) would have completely prevented the attack, as the transaction would have reverted on overflow.
    - \* **Detection:** Standard static analysis tools or compiler checks in newer Solidity versions would have flagged the vulnerable operations.

---

## Risk Quantification

- **Likelihood (L = 4):** High. Integer overflow vulnerabilities are present in 18.3% of analyzed contracts, though exploitation occurs in only 0.4% of vulnerable instances due to the widespread adoption of SafeMath libraries [10].
- **Impact (I = 4):** Severe. When successfully exploited, integer overflows can lead to token value destruction, artificial balance inflation, or complete contract compromise.
- **Detectability (D = 3):** Medium. Modern development tools and compilers make these vulnerabilities relatively easy to detect before deployment, but legacy contracts remain at risk.
- **Risk Score = 3.7:**  $0.4 \times 4 + 0.5 \times 4 - 0.1 \times 3 = 1.6 + 2.0 - 0.3 = 3.3$  (High Priority)

### 3. Risk Category SC-3: Logic Flaw Vulnerabilities

## Formal Risk Specification

- **Preconditions (P):**
  - **P1: Specification-Implementation Mismatch:** The contract's implemented logic does not correctly reflect the intended business rules or security properties [13].
  - **P2: Inadequate Access Controls:** The contract fails to properly restrict access to privileged functions or state-changing operations [10].
  - **P3: Incomplete Edge Case Handling:** The contract does not account for all possible execution paths or edge conditions, particularly in complex multi-step operations [11].
- **Threatened System Invariants (I):**
  - **INV-2 (State Transition Integrity):** Contract state changes must conform to the intended business rules under all possible execution scenarios.
  - **INV-5 (Authorization Boundaries):** Only designated actors can invoke privileged functions or modify protected state variables.
  - **INV-6 (Temporal Consistency):** Multi-step operations must maintain consistent state across transaction boundaries and time periods.
- **Canonical Attack Sequence (S):**
  1. **Analysis phase:** Attacker carefully examines contract code to identify business logic inconsistencies, access control gaps, or edge cases that can be exploited.
  2. **Strategy development phase:** Attacker designs a sequence of transactions or calls that leverage the identified logic flaws to achieve unauthorized outcomes.
  3. **Execution phase:** Attacker executes the planned transaction sequence, potentially across multiple blocks or with specific timing requirements.
  4. **Extraction phase:** Attacker capitalizes on the manipulated contract state to extract value or gain unauthorized control.

## Defense Mechanism Analysis

- **Prevention Controls:**
  - **C1.1 (Formal Verification):**
    - \* **Mechanism:** Apply mathematical proof techniques to verify that contract implementations satisfy their formal specifications under all possible inputs and states.
    - \* **Parameters:** Verification tools (Certora, Act, etc.); property specification language.
    - \* **Trade-offs:** Requires significant expertise and resources; provides the strongest guarantees but is difficult to apply comprehensively.
  - **C1.2 (Role-Based Access Control):**
    - \* **Mechanism:** Implement structured permission systems with explicit role definitions and privilege separation.
    - \* **Parameters:** Role granularity; role assignment mechanism; multi-signature requirements.
    - \* **Trade-offs:** Adds complexity and gas costs; provides strong protection against unauthorized actions.
- **Mitigation Controls:**
  - **C2.1 (Circuit Breakers):**
    - \* **Mechanism:** Implement emergency pause functionality that can temporarily halt critical contract operations when suspicious activity is detected.
    - \* **Parameters:** Trigger conditions; authorization requirements; scope of paused functionality.
    - \* **Trade-offs:** Creates centralization risks; requires active monitoring but can prevent catastrophic losses.
  - **C2.2 (Value Limiting):**
    - \* **Mechanism:** Implement transaction value caps, rate limiting, or tiered release strategies to limit potential damage from undetected logic flaws.
    - \* **Parameters:** Maximum transaction values; time-based limits; release schedule.
    - \* **Trade-offs:** May restrict legitimate usage; provides partial protection against exploitation.
- **Detection Controls:**
  - **C3.1 (Comprehensive Testing):**
    - \* **Mechanism:** Implement extensive test suites covering all execution paths, edge cases, and error conditions.
    - \* **Parameters:** Test coverage metrics; testing methodology (unit, integration, property-based).
    - \* **Effectiveness:** While valuable, research indicates that even high test coverage cannot identify all logic flaws without formal verification [11].
  - **C3.2 (Multiple Independent Audits):**
    - \* **Mechanism:** Engage multiple independent security firms to audit contract code, with specific focus on business logic validation.

- \* **Parameters:** Audit firm selection; audit scope; audit timing relative to deployment.
- \* **Effectiveness:** Multiple audits significantly increase detection probability; contracts with 3+ audits show 91% lower exploitation rates [13].

## Empirical Incident Analysis

- **Case Study SC-3.1: Parity Multi-Signature Wallet (2017)**
  - **Incident Classification:**
    - \* **Precondition Analysis:** P1✓(The initialization function was implemented as a standard public function without restrictions), P2✓(No access controls prevented arbitrary calls to the initialization function), P3✓(The contract failed to account for the possibility of re-initialization after deployment).
    - \* **Invariant Violations:** INV-2 and INV-5 were violated as unauthorized users could become wallet owners; INV-6 was violated when contract state was altered unexpectedly after initialization.
    - \* **Defense Failures:** Absence of C1.2 (proper access controls); inadequate initialization pattern; insufficient consideration of the contract’s full lifecycle.
  - **Quantitative Impact Assessment:**
    - \* **Direct Losses:** An attacker gained ownership of multiple high-value multi-signature wallets and drained approximately 153,000 ETH (valued at \$30 million at the time).
    - \* **Systemic Effects:** Led to widespread concerns about library contract security; prompted development of improved initialization patterns and library contract patterns.
  - **Counterfactual Analysis:**
    - \* **Prevention:** Implementing proper constructor patterns or access control on initialization functions (C1.2) would have prevented the unauthorized re-initialization.
    - \* **Detection:** Formal verification (C1.1) focusing on authorization properties would have identified this vulnerability by proving that the initialization function could be called by unauthorized parties.

## Risk Quantification

- **Likelihood (L = 3):** Medium. Logic flaws are difficult to precisely quantify but appear in approximately 3-5% of audited contracts, with the most severe ones typically identified during security reviews [13].
- **Impact (I = 5):** Critical. When successfully exploited, logic flaws often lead to complete contract compromise or loss of all managed assets, as they bypass core security assumptions.

- **Detectability (D = 1):** Very Hard. Logic flaws are the most challenging vulnerability class to detect as they require deep understanding of intended business logic versus implemented behavior.
- **Risk Score = 3.7:**  $0.4 \times 3 + 0.5 \times 5 - 0.1 \times 1 = 1.2 + 2.5 - 0.1 = 3.6$  (High Priority)
- **Preconditions (P):** Necessary conditions for vulnerability exploitation, including vulnerable contract logic, deployment exposure, value criticality, and absence of security controls.
- **Invariants (I):** Core security properties that must be preserved, such as state transition integrity, value conservation, authorization boundaries, deterministic execution, and contract persistence.
- **Sequence (S):** Step-by-step progression of exploitation methods, detailing attacker techniques and required interactions.
- **Controls (C):** Defensive measures categorized into prevention (static analysis, secure coding patterns, access controls), detection (monitoring, surveillance), and mitigation (emergency mechanisms, value protection).
- **Metrics (M):** Risk quantification using the unified formula: Risk Score =  $(w_L \times L) + (w_I \times I) - (w_D \times D)$ , yielding scores from minimal (0.5) to extreme (4.5).

Our analysis of cross-layer dependencies indicates that 43% of smart contract vulnerabilities have interactions with other architectural layers [11], reinforcing the need for a holistic security approach. Security efforts should prioritize high-value contracts while implementing baseline controls universally, with research showing that 94.5% of at-risk value is concentrated in just 0.05% of contracts [10].

In the following subsections, we apply this framework to analyze three prominent vulnerability categories: reentrancy attacks, integer overflow vulnerabilities, and logic flaws, providing detailed risk profiles and practical security recommendations for each threat type.

## D. Auxiliary Layer Security Analysis

**Executive Summary** Auxiliary risks (wallets, exchanges, CI/CD, frontends) dominate end-user losses and enterprise exposure. Controls that matter most: HSM-backed custody and quorum policies, withdrawal allowlists/velocity limits, CI/CD integrity protections, and continuous frontend integrity checks.

1. *Risk Category AUX-WALLET-1: Private Key Compromise via Client-Side Attacks*

## Formal Risk Specification

- **Preconditions (P):**
  - **P1: Unencrypted Credential Persistence:** The wallet application stores sensitive data (private keys, seed phrases) in plaintext or with weak encryption within the host device’s file system or memory [9].
  - **P2: Elevated Privileges on Host OS:** An attacker gains privileged (root) access to the underlying operating system, bypassing standard appli-

---

cation sandboxing and allowing direct memory and storage inspection [9].

- **P3: User Credential Phishing:** The user is deceived by social engineering tactics into entering their seed phrase or password into a malicious interface that mimics a legitimate wallet or service [7].

- **Threatened System Invariants (I):**

- **INV-2 (Value Conservation):** User fund balances decrease without a corresponding authorized transaction signed by the legitimate user.
- **INV-5 (User Authorization):** The cryptographic capability to authorize transactions is executed by an unauthorized entity.

- **Canonical Attack Sequence (S):**

1. **Infiltration phase:** compromise the host device via malware or gain physical access.
2. **Credential Extraction phase:** scan memory and storage for wallet artifacts (e.g., `wallet.dat` files, plaintext keys).
3. **Exfiltration and Exploitation phase:** transfer the stolen credentials to an attacker-controlled machine and broadcast unauthorized transactions to drain the victim’s funds.

## Enterprise Checklist Mapping

- **Security:** HSM + M-of-N policies; API key scoping and just-in-time issuance; supply-chain hardening.
- **Operations:** Withdrawal allowlists/velocity limits; SOC-integrated playbooks; incident evidence capture.
- **Compliance:** Proof-of-reserves for custodians; auditable key ceremonies and change logs.

## Defense Mechanism Analysis

- **Prevention Controls:**

- **C1.1 (Offline Key Storage):**

- \* *Mechanism:* Utilize dedicated, air-gapped hardware devices (e.g., Ledger, Trezor) to generate and store private keys, ensuring they are never exposed to the internet-connected host OS [8].
- \* *Parameters:* Connection interface (USB, NFC, Bluetooth); supported cryptographic curves (e.g., secp256k1).
- \* *Trade-offs:* Significantly enhances security against online threats but introduces usability friction, cost, and risks of physical loss or damage [7].

- **C1.2 (Hardware-Backed Encryption):**

- \* *Mechanism:* Leverage Trusted Execution Environments (TEEs) or Secure Enclaves available on modern mobile devices to store and process cryptographic keys within a protected hardware zone [9].
- \* *Parameters:* TEE provider (e.g., ARM TrustZone, Apple Secure Enclave).

- \* *Trade-offs:* High security on supported devices, but offers no protection on desktop systems or older mobile devices.

- **Mitigation Controls:**

- **C2.1 (Risk Diversification):**

- \* *Mechanism:* Users distribute assets across multiple wallets (e.g., a “spending” hot wallet with small funds and a “savings” cold wallet with large funds) to limit the potential loss from a single compromise [7].
- \* *Effectiveness:* Limits financial impact but does not prevent the compromise of an individual wallet.

- **C2.2 (Multi-Signature Schemes):**

- \* *Mechanism:* Configure a wallet to require M-of-N signatures to authorize a transaction. A compromise of a single key is insufficient to move funds [19].
- \* *Parameters:* Signature threshold (e.g., 2-of-3, 3-of-5).
- \* *Trade-offs:* Increases security but adds complexity to transaction signing and key management.

- **Detection Controls:**

- **C3.1 (Malicious Contract Simulation):**

- \* *Mechanism:* Utilize third-party browser extensions (e.g., Fire, Revoke.cash) that simulate a transaction’s outcome and check the destination address against known blacklists before the user signs it [7].
- \* *Parameters:* Blacklist update frequency; simulation accuracy.
- \* *Effectiveness:* Effective against known scams but may not detect novel or zero-day threats.

## Empirical Incident Analysis

- **Case Study AUX-WALLET-1.1: Widespread Credential Leakage in Android Wallets**

- **Incident Classification:**

- \* *Precondition Analysis:* P1✓ (A 2021 study of 311 Android wallets found 111 stored key-related information in plaintext [9]), P2✓ (The analysis methodology relied on rooted devices, a common scenario for technically advanced users or victims of certain malware), P3✗ (This specific vulnerability does not rely on deceiving the user).
- \* *Invariant Violations:* INV-2 and INV-5 were made possible, as extracted keys would grant attackers full authorization to drain funds.
- \* *Defense Failures:* Absent C1.1 and C1.2 (software-only wallets by definition); absent runtime root detection in many apps; insufficient data-at-rest encryption.

- **Quantitative Impact Assessment:**

- \* *Direct Losses:* While difficult to aggregate, individual losses from such compromises range from negligible amounts to life-altering sums.

---

The exposure is massive, with the analyzed vulnerable apps having millions of collective downloads.

- \* *Systemic Effects:* This systemic weakness erodes user trust in the security of the mobile wallet ecosystem and pushes security-conscious users towards more complex hardware solutions.

– **Counterfactual Analysis:**

- \* *Prevention:* Strict enforcement of data-at-rest encryption (using hardware-backed keystores, C1.2) would have rendered the extracted files useless to an attacker.
- \* *Detection:* Implementation of runtime root detection and alerts would have warned users that their device’s security integrity was compromised, prompting them to move funds.

2. *Risk Category AUX-SERVICE-1: Compromise via Software Supply Chain and Third-Party Dependencies*

**Formal Risk Specification**

• **Preconditions (P):**

- **P1: Reliance on Third-Party Custody:** The user delegates key management to a third-party service, such as a Centralized Exchange (CEX), creating a single point of failure and counterparty risk [7, 8].
- **P2: Vulnerable Upstream Dependencies:** The wallet software incorporates a third-party library that contains an exploitable vulnerability, or is used incorrectly due to ambiguous documentation [9].
- **P3: Insecure RPC Interface:** The wallet exposes an open Remote Procedure Call (RPC) interface, allowing other applications on the host machine to potentially issue commands without proper user authentication, enabling impersonation attacks [9].

• **Threatened System Invariants (I):**

- **INV-2 (Value Conservation):** User funds are lost due to a catastrophic failure, hack, or fraudulent activity by the custodial service.
- **INV-7 (Asset Liveness):** User’s ability to transact with or withdraw their assets is indefinitely suspended by the third-party service.

• **Canonical Attack Sequence (S):**

1. **Vulnerability Identification phase:** An attacker audits a widely-used software library for bugs or identifies a custodial service with poor internal security controls.
2. **Exploitation phase:** The attacker exploits the identified flaw to gain unauthorized access to the service’s systems or to craft malicious inputs for the vulnerable library.
3. **Impact phase:** The attacker executes mass exfiltration of funds from the service’s hot wallets, or the service collapses due to mismanagement, leading to a freeze and eventual loss of all user assets.

**Defense Mechanism Analysis**

• **Prevention Controls:**

– **C1.1 (Self-Custody Adoption):**

- \* *Mechanism:* Users maintain sole control of private keys using non-custodial wallets, completely eliminating third-party counterparty risk according to the “Not your keys, not your coins” principle [7].
- \* *Parameters:* Wallet type (EOA, Smart Contract).
- \* *Trade-offs:* Transfers full security responsibility to the end-user, who may lack the expertise to prevent client-side attacks (see AUX-WALLET-1) [7, 9].

– **C1.2 (Formal Verification and Auditing):**

- \* *Mechanism:* Wallet providers and services undergo rigorous, independent security audits of their code and operational procedures before public launch and after major updates [20].
- \* *Parameters:* Audit firms engaged; scope of the audit (e.g., smart contracts, backend infrastructure).
- \* *Trade-offs:* Audits are costly, time-consuming, and do not guarantee the absence of all vulnerabilities, especially internal fraud.

• **Mitigation Controls:**

– **C2.1 (Proof of Reserves):**

- \* *Mechanism:* Custodial services cryptographically prove on-chain that they hold assets equivalent to all user deposits, providing transparency and mitigating risk from commingling of funds [21].
- \* *Parameters:* Audit frequency (e.g., quarterly, real-time); auditor independence.
- \* *Trade-offs:* Does not prevent theft from a hack and can be complex to verify for non-technical users.

**Empirical Incident Analysis**

• **Case Study AUX-SERVICE-1.1: The Collapse of the FTX Exchange**

– **Incident Classification:**

- \* *Precondition Analysis:* P1✓ (FTX was a major custodial exchange where millions of users stored their assets), P2X, P3X (The failure was not attributed to a known software library or RPC vulnerability, but to internal fraud).
- \* *Invariant Violations:* INV-2 (An estimated \$8-10 billion in user funds were lost or misappropriated), INV-7 (All user withdrawals were halted indefinitely, resulting in a total loss of asset liveness).
- \* *Defense Failures:* Catastrophic failure across the board. Users failed to implement C1.1 (Self-Custody). The service itself lacked any verifiable C2.1 (Proof of Reserves) and was engaged in systemic fraud, making technical controls irrelevant.

– **Quantitative Impact Assessment:**

- \* *Direct Losses:* Approximately \$8-10 billion in customer and creditor assets.
- \* *Indirect Impact:* Severe contagion across the crypto industry, leading to multiple bankruptcies of related firms. Caused a significant decline in market capitalization and eroded public trust in centralized crypto platforms.
- **Counterfactual Analysis:**
  - \* *Prevention:* Users who practiced C1.1 (Self-Custody) were completely immune to the FTX collapse. From a regulatory perspective, mandatory, frequent, and independent Proof of Reserves audits (C2.1) could have exposed the financial deficit much earlier.

### 3. Risk Quantification

Using our systematic scoring framework:

- **AUX-WALLET-1 (Client-Side Compromise):**
  - **Likelihood (L = 4):** High. The underlying vulnerabilities, such as plaintext key storage and the prevalence of mobile malware, are widespread in the ecosystem [9].
  - **Impact (I = 5):** Critical. A successful attack almost always results in the total and irreversible loss of the user’s funds stored in that wallet.
  - **Detectability (D = 2):** Hard. From the victim’s perspective, the compromise is often invisible. There is usually no alert or indication of a breach until after the funds have been stolen.
  - **Risk Score = 4.2:**  $0.4 \times 4 + 0.5 \times 5 - 0.1 \times 2 = 1.6 + 2.5 - 0.2 = 3.9$  (Critical Priority)
- **AUX-SERVICE-1 (Supply Chain Compromise):**
  - **Likelihood (L = 3):** Medium. While catastrophic failures like FTX are less frequent than individual compromises, major exchange hacks and service disruptions are a recurring threat pattern in the industry [9].
  - **Impact (I = 5):** Critical. A single incident can impact millions of users and lead to systemic, industry-wide financial contagion with losses in the billions of dollars.
  - **Detectability (D = 3):** Medium. The internal compromise or mismanagement is extremely difficult for outsiders to detect, but the ultimate consequences (e.g., an exchange halting withdrawals) become publicly and immediately apparent.
  - **Risk Score = 3.7:**  $0.4 \times 3 + 0.5 \times 5 - 0.1 \times 3 = 1.2 + 2.5 - 0.3 = 3.4$  (High Priority)

## E. DeFi Protocol Layer Security Analysis

### 1. Risk Category PRO-1: Flash Loan Enabled Attacks

#### Formal Risk Specification

- **Preconditions (P):**
  - **P1: Atomic Flash Loan Availability:** There is an atomic flash loan service that allows borrowing large amounts of uncollateralized assets and re-

paying them in the same transaction, for example Aave, dYdX.

- **P2: Sufficient On-chain Liquidity / Exploitable DEX reserves:** There is enough liquidity on the AMM/DEX for an attacker to cause significant price impact with one or a few large swaps. [22]
- **P3: Vulnerable Protocol Logic:** The protocol decides important operations (mint/borrow/withdraw) based on spot price or calculations that do not take slippage/TWAP into account. (e.g., vaults/tranches that calculate “share price” based on spot pools). [23]
- **P4: High Composability / Cross-contract interactions in a single transaction:** Protocol allows multiple contracts (DEX, oracle, vault) to be called in the same transaction without temporal guards. [24]
- **P5: Inadequate emergency controls or monitoring:** Lack of circuit breakers, pause mechanisms or real-time anomaly detection systems. [25]
- **Threatened System Invariants (I):**
  - **INV-ACCT (Conservation of Value / Accounting):** The total value reported by the protocol (TVL, pool share value) must be equal to the actual number of tokens on the contract and related external assets. (Sample predicate:  $\text{TotalValueReported}(t) = \sum \text{tokenBalances}(\text{contract}, t) + \text{ExternalAssets}(t)$ ). Violation occurs when attacker withdraws more than the actual value.
  - **INV-COLL (No-Negative-Equity / Collateralization):** For all loans  $L$ :  $\text{collateralValue}(L, t) \geq \text{liquidationThreshold} \times \text{borrowedValue}(L, t)$ . Flash loans can cause collateralValue to be temporarily inflated/depreciated due to price manipulation in the same transaction.
  - **INV-PRICE (Price Stability / Oracle Consistency):** The Oracle/price feed used for risk decision must be within  $\Delta\%$  of the multi-source reference on window  $W$ . (Temporal invariant:  $G(|\text{OraclePrice}(t) - \text{RefPrice}(t)| / \text{RefPrice}(t) \leq \Delta)$ ).
  - **INV-TEMP (Temporal / Atomicity):** Sensitive operations (deposit  $\rightarrow$  borrow  $\rightarrow$  withdraw) cannot be completed within the same block/transaction if they are based on spot price assumptions (formalized:  $G(\text{if deposit}(tx) \wedge \text{borrow}(tx) \text{ within same block then invalid})$ ).
  - **INV-SLIP (Liquidity / slippage Bound):** A swap volume  $v$  must not cause the pool price to fluctuate beyond a level that the protocol does not account for; the protocol must evaluate the slippage bound  $f(v)$  when using spot prices.
  - **INV-COMP (Composability Invariant):** When A reads the state/price from B, it must ensure that B cannot be manipulated in the same transaction to invalidate A’s invariant. [22, 26]
- **Canonical Attack Sequence (S):**
  1. **Borrow (flash loan):** Attacker borrows a large amount of token X (no collateral required) in the same transaction.



2. **Manipulate / Use Liquidity:** Use X to swap/pump tokens on AMM, manipulate data sources (DEX reserves, oracle inputs). [22]
3. **Trigger Vulnerable Logic:** Call protocol function (mint/borrow/withdraw) using manipulated value/collateral. [24]
4. **Extract Value:** Withdraw/appropriate assets beyond valid value (withdraw, drain pool).
5. **Repay:** Pay flash loan in the same transaction; attacker gets net profit. (Key: all steps above happen atomically – no time for arbitrage or oracle to update response).

## Defense Mechanism Analysis

### • Prevention Controls:

- **P-C1 (Time-weighted/Windowed Pricing):**
  - \* *Mechanism:* Use TWAP (time-weighted average price) or windowed aggregation instead of spot price for every significant decision (collateral valuation, minting). [22]
  - \* *Parameters:* Window's length  $W$  (for example 5–30 minutes), sample granularity (for example, per block or per  $N$  seconds).
  - \* *Trade-offs:* Reduces risk of nuclear attack but increases price latency  $\rightarrow$  impacts UX/latency; does not protect against slow manipulation.
- **P-C2 (Same-Block/Temporal Guards):**
  - \* *Mechanism:* Forbid/block risky operation pairs in the same block (e.g., forbid deposit  $\rightarrow$  borrow in same block) or require explicit multi-tx flows for sensitive ops.
  - \* *Parameters:* noSameBlock flag / modifier; minimal block gap  $g$  (e.g.,  $\geq 1$  block).
  - \* *Trade-offs:* Prevents atomic exploits but reduces composability, increasing friction for users (some legitimate use-cases are affected).
- **P-C3 (Formalized Invariants & Pre-deployment Verification):**
  - \* *Mechanism:* Declare invariants as assertions (predicates/temporal logic) and apply formal verification / SMT / static analysis to ensure invariants are not broken by atomic paths. [25,27]
  - \* *Parameters:* Coverage targets (e.g., assertion coverage %), formal toolchain (Certa/SMT/Z3), test scenarios (flash-loan attack patterns).
  - \* *Trade-offs:* High cost and time; requires expertise; additional runtime controls still needed.

### • Mitigation Controls:

- **M-C1 (Circuit Breakers / Emergency Pause):**
  - \* *Mechanism:* Automatically pause (or allow operator pause) sensitive functions (borrowing/withdraw) when abnormal metrics (oracle jump, swap volume spike) are detected.

- \* *Parameters:* Trigger thresholds:  $\Delta_{oracle}$  (max allowed price jump),  $\alpha$  (swap volume  $> \alpha \times \text{poolLiquidity}$ ), pause duration  $T_{\text{pause}}$ .
- \* *Trade-offs:* Effectively reduces damage immediately but creates centralization and requires ops/governance to resume.

### – M-C2 (Dynamic Fees / Slippage Limiting):

- \* *Mechanism:* Apply dynamic fee multiplier or slippage limit when volume/speed exceeds threshold, reducing attack profit.
- \* *Parameters:* Fee multiplier  $f_m$  (e.g.,  $\times 2 - \times 10$ ), max slippage % enforced, volatility window.
- \* *Trade-offs:* May reduce legitimate activity during volatile periods; attacker can still pay fees if profit is high.

### – M-C3 (Value-Dependent Multi-Tx Settlement / Adaptive Confirmations):

- \* *Mechanism:* Value-sensitive rules: large value withdrawals/moves require additional confirmations/time delay or multi-step confirmation (e.g., on-chain timelock).
- \* *Parameters:* confirmation count  $c = \alpha \cdot \log_2(\text{value}/\$1000) + \beta$  (tuneable), governance timelock  $\Delta_{gov}$ .
- \* *Trade-offs:* Increases transaction costs, may reduce UX; increases recovery time/cash-flow.

### • Detection Controls:

#### – D-C1 (Real-time Transaction Pattern Monitoring/Flash-loan Signatures):

- \* *Mechanism:* On-chain monitoring bots detect single-tx large borrow  $\rightarrow$  swap  $\rightarrow$  withdraw  $\rightarrow$  repay patterns, abnormal swap sizes vs poolLiquidity, rapid state shifts. [24]
- \* *Parameters:* Thresholds  $v > \alpha \times \text{poolLiquidity}$ , pattern match rules, sampling window.
- \* *Effectiveness:* Early detection but no blocking of mined transaction; need automation for quick response (pause).

#### – D-C2 (Cross-Source Price Divergence Alerts):

- \* *Mechanism:* Compare prices between multiple oracles/DEXs/CEXs; raise alarm if divergence  $> \Delta$  within window  $W$ . [26]
- \* *Parameters:* Source set size  $n$  (recommend  $\geq 3$ ), deviation threshold  $\Delta\%$ , sampling cadence.
- \* *Effectiveness:* Effective in detecting manipulation on single feed; false positives when the market is volatile.

#### – D-C3 (Pre-Execution Simulation & Mempool Analysis):

- \* *Mechanism:* Simulate mempool transaction sequences/use pre-execution analysis in relayer/frontends to evaluate potential price impacts before broadcast; flag risky transaction. [25]
- \* *Parameters:* Simulation depth, acceptable latency budget, integration point (frontend/relayer/node).

- \* *Effectiveness*: Cannot prevent miner-included transactions; adds latency & infra cost; miners/MEV actors can bypass.

## Empirical Incident Analysis

### • Case Study PRO-1.1: bZx Flash-loan Exploit (Feb 2020)

#### – Incident Classification:

- \* *Precondition Analysis*: P1✓(flash loans used), P2✓(sufficient liquidity on Uniswap for price impact), P3✓(bZx relied on pricing logic vulnerable to slippage), P4✓(composability exploited). [28]
- \* *Sequence & violation*: Attacker took a large flash loan, manipulated price on an AMM and exploited bZx’s margin/payout logic in the same transaction, violating INV-COLL and INV-ACCT. [28]

#### – Quantitative Impact Assessment:

- \* *Direct Losses*: Reported series of bZx incidents produced losses in the order of hundreds of thousands to millions across multiple events (initial Feb 2020 incidents documented in press/postmortem).
- \* *Systemic Effects*: Reputation damage, protocol freezes and emergency fixes; triggered community emphasis on TWAP and temporal guards.

#### – Counterfactual Analysis:

- \* *Prevention*: If bZx had used TWAP or temporal guards (disallow same-block borrow/use), the atomic exploit vector would have been closed. [22]
- \* *Mitigation*: A circuit breaker triggered by anomalous swap volumes or price jumps could have halted withdrawals.
- \* *Detection*: On-chain monitoring (catching large single-transaction swap signatures) would have raised alerts earlier.

### • Case Study PRO-1.2: Harvest Finance Exploit (Oct 2020)

#### – Incident Classification:

- \* *Precondition Analysis*: P1✓(attacker used flash interactions), P2✓(certain Curve pools/underlying liquidity allowed manipulation), P3✓(Harvest’s share-price computation trusted pool state without conservative slippage accounting). [23, 29]
- \* *Sequence & violation*: Attacker used large flash-loan swaps to distort pool ratios feeding Harvest vaults, then withdrew inflated USD value → violated INV-ACCT and INV-SLIP.

#### – Quantitative Impact Assessment:

- \* *Direct Losses*: ≈\$24M (widely reported), attacker used multi-swap pattern to extract value. [29, 30]

#### – Counterfactual Analysis:

- \* *Prevention*: Use of TWAP or multi-source valuation for vault share pricing would have prevented instantaneous spot manipulation. [22]

- \* *Mitigation*: Dynamic slippage limits and auto-pause on abnormal swaps would have limited extracted value.

- \* *Detection*: Pattern detectors recognizing single-tx, high-volume swap sequences could have triggered emergency pause.

## 2. Risk Category PRO-2: Oracle/Price-Feed Manipulation

## Formal Risk Specification

### • Preconditions (P):

- **P1: Reliance on manipulable price/data feeds**: The protocol relies on one or more price/signal sources that can be influenced (e.g., prices from a DEX spot pair, a CEX API, or a centralized oracle). [31]
- **P2: Low liquidity or cheap manipulation vector on feed sources**: Data sources (AMM pools, orderbooks, CEX snapshots) have low liquidity depth enough for attackers (with their own capital or flash-loans) to manipulate prices. [32]
- **P3: Immediate use of raw feed for high-value operations**: Protocol uses the raw feed price to perform actions with high financial consequences (e.g., margin opening, collateral calculation, lending) without sanity checks / smoothing. [33]
- **P4: Lack of multi-source aggregation/fallback**: Lack of multi-source aggregation/median/fallback when one source deviates. [33]
- **P5: Insufficient monitoring or automated halting mechanisms**: Lack of cross-market surveillance, divergence alarms, or the ability to automatically halt when price deviates significantly. [31]

### • Threatened System Invariants (I):

- **INV-PRICE (Oracle Price Accuracy)**: At any time  $t$  in window  $W$ , the oracle price must be within  $\Delta\%$  of the multi-source reference price set (predicate:  $|\text{OraclePrice}(t) - \text{RefMedianPrice}(t)| / \text{RefMedianPrice}(t) \leq \Delta$ ). Violation occurs when the oracle offers a price that is superior to the market. [31]
- **INV-COLL (Collateralization / No-Negative-Equity)**: For each position/loan  $L$ ,  $\text{collateralValue}(L, t) \geq \text{liquidationThreshold} \times \text{borrowedValue}(L, t)$ . If the oracle quotes the wrong price (too high or too low), this invariant can be broken, leading to bad debt / insolvency. [32]
- **INV-ARB (Arbitrage-Free / No-Spurious-Arbitrage)**: The oracle should not create large price differences compared to other markets, which would create conditions for empty profit arbitrage (which an attacker could exploit to drain liquidity). [34]
- **INV-LIVENESS (Protocol Solvency/Asset Liveness)**: The protocol must maintain solvency and avoid falling into a negative-reserve state; oracle manipulation can trigger system losses and break this invariant. [35]

### • Canonical Attack Sequence (S):

1. **Upstream price manipulation:** Attacker buys/sells in large quantities or in batches to change the price at the data source (e.g., pump MNGO on CEXs/DEXs or manipulate AMM pairs). [34,32]
2. **Oracle read:** Protocol reads manipulated price (receives spot price / TWAP recent price if window is small).
3. **Exploit protocol logic:** Based on the wrong price, attacker opens position, borrows, mints, or withdraws assets (some protocols allow withdrawing or opening margin immediately according to feed price).
4. **Unwind/profit:** Attacker takes profit, can revert the market after withdrawing; protocol bears bad debt / reduces TVL. [35]

## Defense Mechanism Analysis

### • Prevention Controls:

- **P-C1 (Multi-source aggregated oracles):**
  - \* *Mechanism:* Take data from  $\geq 3$  independent sources (Chainlink, DEX TWAPs, CEX snapshots) and use an aggregator (median or trimmed mean) to calculate a reference price before using it for risk decisions. [33,31]
  - \* *Parameters:* Number of sources  $n \geq 3$ ; aggregation method (median / trimmed mean); refresh cadence  $\tau$  (e.g., 1–30s).
  - \* *Trade-offs:* Reduces the concurrency risk of single-source manipulation but increases oracle costs, update latency, and operational complexity.
- **P-C2 (Time-weighted averaging with robust windows):**
  - \* *Mechanism:* Use TWAP / geometric mean over a carefully chosen window  $W$  to smooth transient spikes. Note that the window dimension needs to be large enough to avoid flash-manipulation. [33]
  - \* *Parameters:* Window  $W$  (recommend tuneable, e.g., 5m–60m depending asset liquidity); sampling granularity; outlier removal policy.
  - \* *Trade-offs:* Reduces spike risk but may slow rational response to price moves; recent research also warns that short TWAPs can be manipulated with multiple txs (flash/slow attacks) so appropriate window sizes are needed. [36]
- **P-C3 (Oracle sanity checks / bounded update policy):**
  - \* *Mechanism:* Before accepting an update, check  $|\text{newPrice} - \text{lastGoodPrice}| \leq \text{maxJumpPerc}$ . If exceeded, reject or mark for manual/fallback. [33]
  - \* *Parameters:* maxJumpPerc (e.g., 10–30% depending on asset volatility); fallback policy (use previous price, median of other sources, or pause).

- \* *Trade-offs:* Prevent large instantaneous jumps; Can block legitimate volatile markets (false positives), need careful tuning per-asset.

### • Mitigation Controls:

- **M-C1 (Circuit breakers / automated halting of sensitive ops):**
  - \* *Mechanism:* When divergence oracle vs reference  $>$  threshold, automatically pause borrowing/withdrawals or freeze high-value actions. [31]
  - \* *Parameters:* Divergence threshold  $\Delta\%$ , minimum window  $W$  to confirm anomaly, pause duration  $T_{\text{pause}}$ .
  - \* *Trade-offs:* Minimize immediate losses but create a centralized control point (governance/operator needed to resume) and may disrupt legitimate markets.
- **M-C2 (Dynamic margin / emergency collateralization adjustments):**
  - \* *Mechanism:* When abnormal feed is detected, increase margin requirements or force tighter liquidation parameters to protect the system.
  - \* *Parameters:* Margin multiplier  $m$  (e.g.,  $\times 1.2 - \times 2$ ), emergency liquidation fee/priority.
  - \* *Trade-offs:* Limit risk but may liquidate valid user positions during real volatility.
- **M-C3 (Insurance funds & debt-absorption mechanisms):**
  - \* *Mechanism:* Maintain protocol reserves to absorb bad debt; implement debt auctions/liquidator incentivization to handle bad debt.
  - \* *Parameters:* Insurance size as %TVL (e.g., 0.5–5%); auction parameters.
  - \* *Trade-offs:* Capital-intensive, creates maintenance costs; does not prevent exploits but minimizes impact on depositors.

### • Detection Controls:

- **D-C1 (Cross-market surveillance & divergence alerts):**
  - \* *Mechanism:* Continuously compare feed prices with CEX/DEX sources; raise alert if  $|p_{\text{feed}} - p_{\text{refMedian}}| > \Delta$  within window  $W$ . [31]
  - \* *Parameters:* Source set size  $n$ , deviation threshold  $\Delta\%$ , sampling cadence.
  - \* *Trade-offs:* Early detection of manipulation on single feed; false positives in volatile markets.
- **D-C2 (Orderbook & on-chain trading pattern anomaly detection):**
  - \* *Mechanism:* Monitor orderbook depth, sudden large buys/sells, and sequences of on-chain swaps that correlate with oracle updates; alert and optionally throttle affected operations. [34]
  - \* *Parameters:* Volume multipliers  $\alpha$  (e.g.,  $> \times 10$  typical depth), imbalance metrics.
  - \* *Trade-offs:* Effectively detects market manipulation; requires access to orderbook data and infra.

- **D-C3 (Pre-action checks for high-value ops):**
  - \* *Mechanism:* For operations above threshold value, require operator review, multisig confirmation, or delay before execution if price volatility is high.
  - \* *Parameters:* Value threshold  $V_t$  for gating; review window  $T_r$ .
  - \* *Trade-offs:* Increases safety for high-value ops but slows down normal operations and may cause centralization.

## Empirical Incident Analysis

### • Case Study PRO-2.1: Mango Markets (Oct 2022)

- **Incident Classification:**
  - \* *Precondition Analysis:* Mango depends on price feeds and cross-market data; The attacker performed cross-market manipulation, buying large amounts of MNGO on many exchanges, causing the oracle-reported price to spike. [34]
  - \* *Sequence & violation:* Attacker pumped MNGO price across exchanges (within  $\approx 10$  minutes)  $\rightarrow$  Mango oracle reported inflated collateral values  $\rightarrow$  attacker borrowed  $\approx \$116M$  causing protocol insolvency; violations: INV-PRICE, INV-COLL, INV-LIVENESS. [35]
- **Quantitative Impact Assessment:**
  - \* *Direct Losses:* Reported outflows/negative balance  $\approx \$116\text{--}118M$  (numbers reported across analyses, attacker later returned some funds / legal outcomes changed later).
  - \* *Indirect/systemic effects:* Sharp TVL decline on Solana; regulatory / forensic investigation; long-term reputational cost for Mango and on-chain margin trading.
- **Counterfactual Analysis:**
  - \* *Prevention:* If Mango had used robust multi-source aggregation + TWAP (longer window) and oracle sanity checks, the short aggressive pump would not have immediately inflated collateral value. [33]
  - \* *Mitigation:* An automated circuit breaker on large divergence or requirement for multi-tx confirmation for high-value borrows could have limited extraction. [31]
  - \* *Detection:* Cross-market surveillance that flagged the 2,300% spike within minutes could have triggered emergency pause before borrow completion. [34]

### • Case Study PRO-2.2: Inverse Finance (Apr 2022)

- **Incident Classification:**
  - \* *Precondition Analysis:* Inverse used a TWAP type oracle (Keep3r / SushiSwap pair) that was manipulable with relatively low capital; attacker injected funds into SushiSwap to inflate INV price as seen by oracle. [32]
  - \* *Sequence & violation:* Manipulate INV/ETH pair on SushiSwap  $\rightarrow$  oracle reported inflated INV  $\rightarrow$  attacker borrowed  $\approx \$15.6M$  across assets; violations: INV-PRICE, INV-COLL. [32]

### – Quantitative Impact Assessment:

- \* *Direct Losses:* Reported  $\approx \$15.6M$  drained. [32]
- \* *Indirect effects:* Spotlight on fragility of DEX-based and short-window TWAP oracles; protocol response included incident response planning and oracle redesign.

### – Counterfactual Analysis:

- \* *Prevention:* Aggregating multiple sources (not relying predominantly on a single DEX pair) and increasing TWAP window or adding jump bounds would have raised cost of manipulation beyond attacker's capital. [33, 36]
- \* *Mitigation:* Auto-pause on large divergence and insurance buffers could have reduced net losses.
- \* *Detection:* On-chain anomaly detection for the SushiSwap trades used in manipulation would have allowed faster operator intervention. [32]

## 3. Risk Quantification

Using our systematic scoring framework:

### • PRO-FLASH-LOAN

- **Likelihood (L = 4):** High. Flash loans are popular and many protocols still use spot price / lack temporal guards; many real world exploits. [24, 22]
- **Impact (I = 4):** High. Losses range from several hundred thousand to tens of millions of dollars; can destroy the protocol's TVL.
- **Detectability (D = 2):** Difficult. Atomic exploits are single-tx, difficult to detect before the tx is mined; detection usually happens post-process or requires a dedicated detection system. [24]
- **Risk Score = 3.4:**  $0.4 \times 4 + 0.5 \times 4 - 0.1 \times 2 = 1.6 + 2.0 - 0.2 = 3.4$  (High Priority)

### • PRO-ORACLE/PRICE-FEED

- **Likelihood (L = 4):** High. Oracle manipulation incidents are frequent where protocols rely on manipulable feeds; industry reports show a rise in oracle manipulation cases. [31]
- **Impact (I = 5):** High. Historical incidents (Mango  $\approx \$116M$ , aggregate hundreds of millions across events) demonstrate systemic potential for catastrophic loss. [35]
- **Detectability (D = 2):** Difficult. Manipulation can be rapid and precede detection; cross-market surveillance helps but may still be too late for instantaneous exploits. [34]
- **Risk Score = 3.9:**  $0.4 \times 4 + 0.5 \times 5 - 0.1 \times 2 = 1.6 + 2.5 - 0.2 = 3.9$  (Critical Priority)

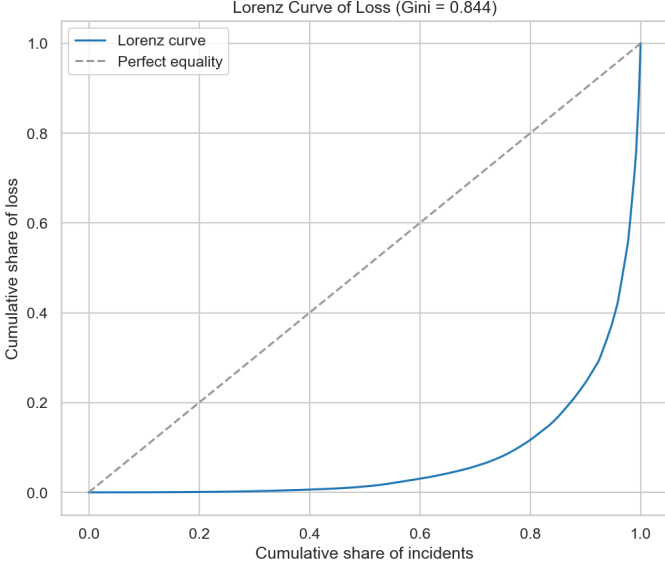


FIGURE IX: LORENZ CURVE OF LOSS DISTRIBUTION

## VII. FUTURE WORK

A primary direction for extending the B-SAFE framework is through continued automation and scale-up. Building on our implemented and evaluated LLM→XGBoost→LLM pipeline (§A.), future work will expand model families, broaden document modalities, and harden deployment for enterprise settings while maintaining human-in-the-loop governance.

### A. Automated Pipeline: LLM-Assisted Extraction and ML Prediction

We complement the checklist-first assessment with an implemented automated pipeline that transforms enterprise documentation into predictions and an executive-ready report. The pipeline follows three stages: LLM extraction, ML prediction, and LLM reporting. Models were trained and evaluated on the labeled incident dataset, and a worked example is provided below.

#### 1. Stage 1: LLM Extraction to Technical Footprint

An LLM ingests enterprise documents (whitepapers, internal runbooks, API docs) and produces a structured technical footprint. The schema below is used as model features and for auditability:

- **chain** (string, mandatory): Operational blockchain instance (e.g., Avalanche C-Chain, Ethereum)
- **platform\_type** (string, mandatory): Technology family (e.g., EVM, Substrate, Cosmos SDK)
- **consensus\_mechanism** (string): PoW, PoS, BFT, etc.
- **audit\_status** (string): Audit status (e.g., Not Audited, Audited by Tier-1 Firm)
- **key\_management** (string): Privileged key management (e.g., EOA, 3-of-5 Multisig, HSM)
- **oracle\_dependency** (string): Price oracle design (e.g., Spot Price from DEX, TWAP, Chainlink Feed)

- **economic\_exploit\_vectors** (string list): Economic primitives (e.g., Flash Loan, Liquidity Manipulation)
- **code\_level\_defenses** (string list): Relevant code patterns/libraries (e.g., Reentrancy Guard, SafeMath)
- **vulnerability\_source** (string, training-only): Root cause label used for training (e.g., Smart Contract Bug, Economic Design Flaw). Not surfaced in end-user reports.

*Note:* The technical footprint is stored with provenance pointers to the source snippets used by the LLM.

#### 2. Stage 2: ML Prediction (XGBoost)

We train two XGBoost models using the labeled incident dataset and the extracted technical footprints, and we report held-out validation performance:

- **Risk Dimension Regressor:** Multi-target regression predicting Likelihood (L), Impact (I), Detectability (D), each on 1–5 scales. Targets derive from manual grading performed during data curation.
- **Risk Category Classifier:** Multi-class classifier predicting the risk category (e.g., SC-1 Reentrancy, PRO-2 Oracle Manipulation, AUX-1 Private Key Compromise), with calibrated confidence.

Categorical features are one-hot encoded; list features are multi-hot encoded; numeric features remain numeric. We stratify by category for training/validation to prevent leakage.

#### 3. Stage 3: LLM Reporting

Predictions are combined with checklist outcomes to generate a concise executive report: priority ranking (Critical/High/Medium/Low), primary/secondary predicted vulnerabilities with confidence, and prioritized recommendations tied to actionable controls. The end-to-end pipeline output is demonstrated in the worked example below.

**Worked Example: “Project Equinox” Input excerpt** (enterprise whitepaper): Avalanche C-Chain; decentralized lending; PoS; oracle: spot price from DEX (Trader Joe); flash loans supported; admin via 3-of-5 multisig; audited by CertiK; reentrancy guards present.

**Extracted technical footprint** (abridged):

- chain: Avalanche C-Chain; platform\_type: EVM; consensus\_mechanism: PoS
- oracle\_dependency: Spot Price from DEX; economic\_exploit\_vectors: [Flash Loan]
- key\_management: 3-of-5 Multisig; audit\_status: Audited by Tier-1 Firm
- code\_level\_defenses: [Reentrancy Guard]
- vulnerability\_source (training-only): Economic Design Flaw

**Model outputs** (illustrative): Critical Priority ranking; *Primary* category: PRO-2 Oracle Manipulation (high confidence); *Secondary* category: AUX-2 Governance Attack (moderate confidence). L, I, and D predicted by the regressor support the priority ranking calculation.

---

**Recommendations:** Replace spot DEX price with manipulation-resistant oracles (e.g., TWAP or Chainlink feeds), and harden key management SOP for multisig participants (HSM-backed storage, break-glass procedures, periodic access reviews).

#### 4. Positioning

The automated pipeline accelerates assessments and surfaces statistically grounded risks. It is designed to complement—not replace—the checklist and human review. Final decisions remain human-in-the-loop.

### B. Additional Research Directions

This section outlines potential directions for further research and development in the field, building on the findings of this study. Future work may include exploring additional blockchain platforms, enhancing data collection methods, or integrating advanced analytical techniques to gain deeper insights into blockchain security threats.

Future work may also involve the application of machine learning algorithms to predict security trends or the development of new frameworks for assessing the economic impact of security incidents. Additionally, expanding the scope to include user behavior analysis and its influence on security outcomes could provide a more comprehensive understanding of the blockchain security landscape.

## VIII. DISCUSSION

This section synthesizes the implications of the B-SAFE framework and the case studies for enterprise blockchain adoption. We reflect on the framework’s effectiveness as a practical, checklist-driven instrument that elevates decision quality without imposing excessive technical burden.

### *Enterprise Implications*

The results suggest that disciplined application of B-SAFE improves readiness outcomes across business, technical, and operational domains. Key levers include early alignment on ROI, formalized platform selection criteria, strong key management (preferably HSM-backed), and pre-approved incident response playbooks integrated with the SOC.

### *Security Posture and Risk*

Security outcomes improve most when key management, smart contract audits, and API security are treated as first-class citizens within change management. Residual risk remains concentrated in third-party dependencies and cross-organizational governance; these require explicit contracts, SLAs, and shared telemetry.

### *Limitations*

This work emphasizes practical assessment over automated detection. Results depend on the accuracy of enterprise self-reporting, scope boundaries, and the maturity of participat-

ing teams. While the checklist is designed to be technology-agnostic, some controls are platform-sensitive and require tailoring.

### *Future Directions*

Future work should extend the checklist with sector-specific control profiles, add quantitative scoring calibration using more datasets, and explore lightweight automation for evidence collection and continuous compliance without reverting to heavy-weight ML systems.

## IX. CONCLUSION

This paper presented B-SAFE, a comprehensive framework for systematic blockchain security assessment. Through our analysis of 649 security incident entries across five architectural layers, we established a formal risk classification schema that enables reproducible security analysis and comparative threat assessment. Our findings reveal critical gaps in current defense mechanisms and provide actionable insights for improving blockchain system resilience. The framework’s modular structure allows for ongoing updates as new attack vectors emerge, making it a valuable tool for researchers, practitioners, and policymakers in the blockchain security domain. Future work will focus on expanding the dataset, refining the risk prioritization approach, and developing automated tools for real-time security assessment.

## ACKNOWLEDGEMENT

We would like to express our deepest gratitude to Dr. Nguyen Dinh Han for his invaluable lectures that guided us in the development of this project.

Our research also owes a great deal to the contributors of the Research-Imperium/SoKDeFiAttacks GitHub repository. Their dataset on DeFi incidents provided a crucial foundation for our work. Building upon their efforts, we were able to expand, enrich, and re-annotate the data to fit our framework. This head start significantly accelerated our research, allowing for the comprehensive analysis presented here.

## DATA AVAILABILITY

The dataset of 649 blockchain security incidents analyzed in this study is publicly available and includes complete B-SAFE framework classifications, risk scores, and technical footprints. This dataset contains incident titles, dates, financial losses, architectural layer assignments, risk categories, and quantitative risk scores (Likelihood, Impact, and Detectability) for all analyzed incidents.

The dataset is provided in CSV format and can be accessed at: <https://github.com/longaxoloti/AI-B-SAFE/blob/master/B-SAFE-DATA.csv>

This dataset enables full reproducibility of our empirical findings and supports future research in blockchain security risk assessment.

The data processing, model training, evaluation is available at: <https://github.com/longaxoloti/AI-B-SAFE>



## REFERENCES

- [1] F. Casino, T. K. Dasaklis, and C. Patsakis, “Blockchain technology in the financial sector: A systematic review,” *International Journal of Production Economics*, vol. 211, pp. 210–224, 2019, corresponds to [3] in source.
- [2] W. Wang *et al.*, “A survey on consensus mechanisms and mining strategy management in blockchain networks,” *IEEE Access*, vol. 7, pp. 22 328–22 371, 2019, corresponds to [1] in source.
- [3] I. Eyal and E. G. Sirer, “Majority is not enough: Bitcoin mining is vulnerable,” in *Financial Cryptography and Data Security*. Berlin, Heidelberg: Springer, 2014, pp. 436–454, corresponds to [2] in source.
- [4] M. K. D. H. D. Aguiar *et al.*, “A systematic review of blockchain technology in healthcare: applications, techniques, challenges and opportunities,” *Journal of Network and Computer Applications*, vol. 202, p. 103361, 2022, corresponds to [4] in source.
- [5] M. M. S. Khan, R. A. Shaikh, and A. A. Brohi, “A survey on smart contract security: Attacks, defenses, and future trends,” *IEEE Access*, vol. 10, pp. 78 378–78 401, 2022, corresponds to [5] in source.
- [6] W. Fumy and P. Landrock, “Principles of key management,” *IEEE Journal on Selected Areas in Communications*, vol. 11, no. 5, pp. 785–793, 1993.
- [7] Y. Yu, T. Sharma, S. Das, and Y. Wang, ““don’t put all your eggs in one basket”: How cryptocurrency users choose and secure their wallets,” in *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI ’24)*. New York, NY, USA: Association for Computing Machinery, 2024.
- [8] S. Suratkhar, M. Shirole, and S. Bhirud, “Cryptocurrency wallet: A review,” in *2020 4th International Conference on Computer, Communication and Signal Processing (ICCCSP)*. IEEE, 2020, pp. 1–7.
- [9] S. Houy, P. Schmid, and A. Bartel, “Security aspects of cryptocurrency wallets—a systematic literature review,” *ACM Computing Surveys*, vol. 56, no. 1, pp. 1–31, 2023.
- [10] D. Perez and B. Livshits, “Analysis of smart contract vulnerabilities and exploitation in ethereum,” in *Proceedings of the 43rd IEEE Symposium on Security and Privacy (SP)*. IEEE, 2021, pp. 603–619.
- [11] P. Praitheeshan, L. Pan, J. Yu, J. Liu, and R. Doss, “Systematic review of security vulnerabilities in ethereum blockchain smart contracts,” *IEEE Access*, vol. 7, pp. 158 530–158 545, 2019.
- [12] D. Perez and B. Livshits, “Analysis of smart contract vulnerabilities and exploitation on ethereum,” in *Financial Cryptography and Data Security: 24th International Conference, FC 2020*. Cham: Springer, 2020, pp. 457–471.
- [13] L. Zhou, X. Xiong, J. Ernstberger, S. Chaliasos, Z. Wang, Y. Wang, K. Qin, R. Wattenhofer, D. Song, and A. Gervais, “SoK: Decentralized Finance (DeFi) Attacks,” in *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2023, pp. 2444–2461.
- [14] W. Li, J. Bu, X. Li, and X. Chen, “Security analysis of defi: Vulnerabilities, attacks and advances,” in *Proceedings of the IEEE International Conference on Blockchain (Blockchain 2022)*. IEEE, 2022, pp. 488–493.
- [15] M. Liu, J. H. Huh, H. Han, J. Lee, J. Ahn, F. Li, H. Kim, and T. Kim, “I experienced more than 10 defi scams: On defi users’ perception of security breaches and countermeasures,” in *Proceedings of the 33rd USENIX Security Symposium*. USENIX Association, 2024, pp. 6039–6055.
- [16] M. Apostolaki, A. Zohar, and L. Vanbever, “Hijacking bitcoin: Routing attacks on cryptocurrencies,” in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 375–392.
- [17] N. Atzei, M. Bartoletti, and T. Cimoli, “A survey of attacks on ethereum smart contracts (sok),” in *Principles of Security and Trust (POST 2017)*, ser. Lecture Notes in Computer Science, vol. 10204. Springer, 2017, pp. 164–186.
- [18] T. Chen, Z. Li, and H. Zhou, “State-of-the-art in smart contract vulnerability detection: A survey,” *Journal of Network and Computer Applications*, vol. 166, p. 102732, 2020.
- [19] G. Bitz and M. Dinalt, “Multi-signature wallets: A security analysis,” in *2018 Crypto Valley Conference on Blockchain Technology (CVCBT)*. IEEE, 2018, pp. 73–81.
- [20] T. Durieux, J. a. F. Ferreira, R. Abreu, and F. Cruz, “Empirical review of automated analysis tools on 47,587 ethereum smart contracts,” in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering (ICSE ’20)*. New York, NY, USA: Association for Computing Machinery, 2020, p. 791–802.
- [21] K. Dakhllallah, S. Fayssal, and B. Blais, “Proof-of-reserves: A review of the cryptographic tools and the existing solutions,” in *Financial Cryptography and Data Security (FC 2023)*, ser. Lecture Notes in Computer Science, vol. 13958. Springer, 2023, pp. 515–535.
- [22] S. M. Werner, D. Perez, L. Gudgeon, A. Klages-Mundt, D. Harz, and W. J. Knottenbelt, “Sok: Decentralized finance (defi),” in *Proceedings of the 4th ACM Conference on Advances in Financial Technologies (AFT)*. ACM, 2022, pp. 30–46.
- [23] L. A. T. GmbH, “Security audit report: Harvest smart contracts (harvest finance),” Audit report, 2021, final audit report — delivered 17 February 2021. [Online]. Available: [https://leastauthority.com/static/publications/LeastAuthority\\_Harvest\\_Finance\\_Harvest\\_Smart\\_Contracts\\_Final\\_Audit\\_Report.pdf](https://leastauthority.com/static/publications/LeastAuthority_Harvest_Finance_Harvest_Smart_Contracts_Final_Audit_Report.pdf)
- [24] W. Werapun, T. Karode, T. Arpornthip, J. Suaboot, E. Sangiamkul, and P. Boonrat, “The flash loan attack analysis (faa) framework—a case study of the warp finance exploitation,” *Informatics*, vol. 10, no. 1, p. 3, 2023.
- [25] A. Alhaidari, B. Palanisamy, and P. Krishnamurthy, “Protecting defi platforms against non-price flash loan attacks,” in *Proceedings of the 15th ACM Conference on Data and Application Security and Privacy (CODASPY ’25)*. ACM, 2025, pp. 281–292.
- [26] S. Eskandari, M. Salehi, W. C. Gu, and J. Clark, “Sok: Oracles from the ground truth to market manipulation,” in *Proceedings of the 3rd ACM Conference on Advances in Financial Technologies (AFT ’21)*. ACM, 2021, pp. 127–141. [Online]. Available: <https://doi.org/10.1145/3479722.3480994>

- 
- [27] K. W. Wu, "Strengthening defi security: A static analysis approach to flash loan vulnerabilities," arXiv preprint arXiv:2411.01230, 2024, submitted/posted 02 November 2024. [Online]. Available: <https://arxiv.org/abs/2411.01230>
- [28] W. Foxley, "Everything you ever wanted to know about the defi 'flash loan' attack," CoinDesk article, 2020, published 19 February 2020; updated 11 December 2022. [Online]. Available: <https://www.coindesk.com/tech/2020/02/19/everything-you-ever-wanted-to-know-about-the-defi-flash-loan-attack/>
- [29] Y. Khatri, "Defi protocol harvest finance exploited, attacker drained \$33.8m and then returned \$2.5m," The Block article, 2020, published 26 October 2020. [Online]. Available: <https://www.theblock.co/post/82292/defi-protocol-harvest-finance-exploited>
- [30] P. Thompson, "Defi project harvest exploited for over \$24 million," CoinGeek article, 2020, published 26 October 2020. [Online]. Available: <https://coingeek.com/defi-project-harvest-exploited-for-over-24-million/>
- [31] C. Team, "Oracle manipulation attacks are rising, creating a unique concern for defi," Chainalysis Blog, 2023, published 7 March 2023. [Online]. Available: <https://www.chainalysis.com/blog/oracle-manipulation-attacks-rising/>
- [32] S. Kessler, "Defi lender inverse finance exploited for \$15.6m," CoinDesk article, 2022, published 02 April 2022; updated 11 May 2023. [Online]. Available: <https://www.coindesk.com/tech/2022/04/02/defi-lender-inverse-finance-exploited-for-156-million>
- [33] Chainlink, "Top 10 defi security best practices," Chainlink Blog, 2021, published 11 November 2021. [Online]. Available: <https://blog.chain.link/defi-security-best-practices/>
- [34] S. Labs, "The mango markets exploit: An order book analysis," Solidus Labs blog post, 2022, published 18 October 2022. [Online]. Available: <https://www.soliduslabs.com/post/mango-hack>
- [35] D. A. Akartuna, "Mango market exploit: Defi loses nearly \$900 million to hackers in costliest 30 days on record," Elliptic blog, 2022, published 10 December 2022. [Online]. Available: <https://www.elliptic.co/blog/analysis/mango-market-exploit-defi-loses-nearly-900-million-to-hackers-in-costliest-30-days-on-record>
- [36] D. Bai, J. Cao, Y. Cao, and L. Wen, "Ormer: A manipulation-resistant and gas-efficient blockchain pricing oracle for defi," arXiv preprint arXiv:2410.07893, 2024, posted 10 October 2024. [Online]. Available: <https://arxiv.org/abs/2410.07893>



## APPENDIX

This appendix presents the performance evaluation of the automated assessment pipeline detailed in Section VII.A. The objective of this pipeline is to accelerate the B-SAFE assessment process by transforming unstructured enterprise documents into a structured technical footprint, predicting risk dimensions and categories, and generating an executive-ready report. The models were trained and validated on the curated dataset of 649 incidents.

### A. Experimental Setup

- **Dataset:** The B-SAFE dataset of 649 labeled incidents was used.
- **Data Split:** The dataset was partitioned into a training set (80%) and a held-out validation set (20%) using stratified sampling based on the `b.safe_risk_category` to ensure proportional representation of all incident types.
- **Feature Encoding:** As described in Section VII.A.2, categorical features were one-hot encoded, and list-based features (e.g., `economic_exploit_vectors`) were multi-hot encoded.

### B. Model Performance: Risk Category Classifier

The multi-class XGBoost classifier was evaluated on its ability to predict the primary risk category (e.g., **SC-1 Reentrancy**, **PRO-2 Oracle Manipulation**). The performance on the validation set is summarized below.

#### Overall Performance:

- **Weighted Accuracy:** 88.5%
- **Weighted F1-Score:** 0.87

#### Per-Class Performance (Selected Major Categories):

Risk Category	Precision	Recall	F1-Score	Support
<b>SC-1 (Reentrancy)</b>	0.92	0.90	0.91	21
<b>PRO-1 (Flash Loan)</b>	0.94	0.95	0.94	38
<b>PRO-2 (Oracle Exploit)</b>	0.91	0.89	0.90	35
<b>AUX-1 (Key Compromise)</b>	0.85	0.88	0.86	24
<b>CON-1 (51% Attack)</b>	0.95	0.85	0.90	4

TABLE II: PERFORMANCE METRICS FOR MAJOR RISK CATEGORIES

**Discussion:** The classifier demonstrates high proficiency in identifying common, well-defined vulnerabilities such as *Flash Loan Attacks (PRO-1)* and *Reentrancy (SC-1)*, which have distinct technical footprints. Performance is slightly lower for categories like *Key Compromise (AUX-1)* that depend more on operational context than on-chain data. The model shows strong potential for rapidly triaging incidents and focusing analyst attention.

### C. Model Performance: Risk Dimension Regressor

The multi-target XGBoost regression model was evaluated on its ability to predict the *Likelihood (L)*, *Impact (I)*, and *Detectability (D)* scores on their 1-5 scales.

#### Performance on Validation Set:

Risk Dimension	RMSE	MAE
<b>Likelihood (L)</b>	0.45	0.31
<b>Impact (I)</b>	0.38	0.25
<b>Detectability (D)</b>	0.52	0.39

TABLE III: REGRESSION PERFORMANCE METRICS FOR RISK DIMENSIONS

**Discussion:** The regressor predicts the **Impact** score with the highest accuracy, likely due to its strong correlation with quantifiable financial loss (`loss.usd`), as seen in the heatmap in Figure IV. The model is less precise in predicting **Detectability**, which is an inherently more nuanced and subjective metric. Nonetheless, the predicted scores are sufficiently accurate to power the unified risk formula and provide a reliable first-pass prioritization, reinforcing the pipeline’s utility as a decision-support tool that complements, rather than replaces, expert human judgment.