

# ConfuGuard: Using Metadata to Detect Active and Stealthy Package Confusion Attacks Accurately and at Scale

Anonymous Author(s)

## ABSTRACT

Package confusion attacks such as typosquatting threaten software supply chains. Attackers make packages with names that syntactically or semantically resemble legitimate ones, tricking engineers into installing malware. While prior work has developed defenses against package confusions in some software package registries, notably NPM, PyPI, and RubyGems, gaps remain: high false-positive rates; generalization to more software package ecosystems; and insights from real-world deployment.

In this work, we introduce ConfuGuard, a solution designed to address the challenges posed by package confusion threats. We begin by presenting the first empirical analysis of benign signals derived from prior package confusion data, uncovering their threat patterns, engineering practices, and measurable attributes. We observed that 13.3% of real package confusion attacks are initially stealthy, so we take that into consideration and refined the definitions. Building on state-of-the-art approaches, we extend support from three to six software package registries, and leverage package metadata to distinguish benign packages. Our approach significantly reduces 64% false-positive (from 77% to 13%), with acceptable additional overhead to filter out benign packages by analyzing the package metadata. ConfuGuard is in production at our industry partner, whose analysts have already confirmed 301 packages detected by ConfuGuard as real attacks. We share lessons learned from production and provide insights to researchers.

## 1 INTRODUCTION

Software package registries (SPRs) are indispensable in modern development. They help engineers develop applications by integrating open-source packages, reducing costs and accelerating product cycles [47, 52, 67]. SPRs are used by industry and governments [12, 57], including in AI and safety-critical applications [27, 55]. The resulting software supply chain is an attractive target — if an adversary can cause an engineer to install a malicious package, the attacker may gain access to downstream users’ resources. One of their tactics is *package confusion attacks* [28, 41]: creating malicious packages with mis-leading names [23, 30, 44]. Figure 1 shows the attack model, covering naming tactics (lexical, syntactic, and semantic confusion) and attack timing (active vs. stealthy).

Researchers have proposed various methods to mitigate package confusion attacks. Early solutions detected packages with lexically similar names [50, 58, 65]. State-of-the-art approaches have begun exploring semantic similarity [41]. However, as criticized by Ohm *et al.* [43], prior work did not consider the false positive rate which is important for industry use [43]. We observe two other shortcomings: prior work focuses on limited SPRs, notably NPM and PyPI; and the literature lacks any experience reports from production deployment.

With our industry partner, we developed a novel system to detect package confusion attacks. To address the problem of high

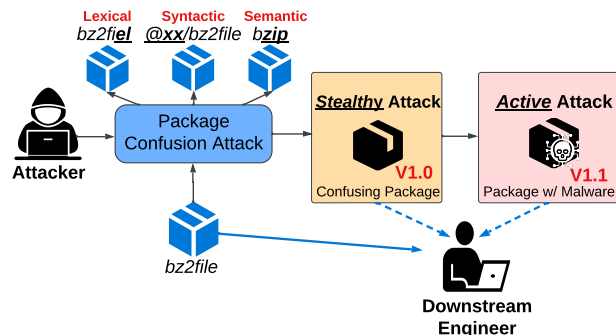


Figure 1: Threat model illustrating common active and stealth package confusion attacks. These malicious packages can mislead developers into installing them or being included as dependencies in downstream packages.

false positive rates, we measured the engineering practices of package confusion attackers, and propose signals, based on package metadata, that distinguish benign packages from genuine (though perhaps stealthy) threats. To generalize prior work to more SPRs, we analyzed package naming practices and confusion attacks in six SPRs, and identified novel aspects of these SPRs that expand the attack surface for confusion attacks. To make our approach production-grade, we develop embedding-based name comparisons and a novel nearest-neighbor search algorithm that provides greater efficiency and accuracy.

We integrated these techniques in the ConfuGuard system. We evaluated ConfuGuard on previous attack datasets as well as several months’ of production data. We report that ConfuGuard achieves a false positive rate of 13%, an improvement of 64% over the state of the art. On a modern server, ConfuGuard has average latency of 6.8 seconds per package. In production, ConfuGuard identified 301 confirmed package confusion attacks in its first three months.

To summarize our contributions:

- We refine the definition of package confusion attacks to encompass both active and stealthy threats.
- We generalize prior analyses to several new SPRs, noting the effect of package naming hierarchies on naming attack tactics.
- We propose novel algorithms to improve attack detection accuracy (embedding) and efficiency (nearest neighbor search).
- We open-source our production dataset, ConfuDB, which includes 1,561 production data triaged by security analysts.
- We share lessons learned from using ConfuGuard in production.

**Significance:** Package confusion attacks are common yet challenging to detect. This paper extends current knowledge by introducing a novel detection system, designed through empirical analysis of previous data and insights gained from ongoing deployment. ConfuGuard is used in production by our industry partner.

## 2 BACKGROUND AND RELATED WORK

Here we discuss software package registries and naming conventions (§2.1), package confusion attacks (§2.2), and defenses (§2.3).

### 2.1 SPRs and Package Naming

Much of modern software development depends on community ecosystems of libraries and frameworks that facilitate software reuse [67, 69]. We refer to these reusable software artifacts as *software packages*. These packages are distributed by *software package registries (SPRs)* [49, 52]. Table 1 shows the SPRs for six popular ecosystems comprising over 7 million software packages. Some SPRs restrict package publication to individual accounts, so the publisher’s identifier does not appear in the package name. Others allow grouping packages into namespaces or organizations, which are incorporated into the package name. We call the former “1-level” and the latter “hierarchical”.

SPRs typically act as infrastructure and leave it to the engineering community to establish norms for interaction. One aspect of their “hands-off” approach is that SPRs impose almost no constraints on package names — they define an allowed character set and forbid duplicate names [2, 3, 6–9]. This flexibility introduces two related challenges for engineers. First, based on application needs (e.g., traditional software vs. machine learning), engineers in different SPRs have opted for different naming conventions, which can enhance readability, maintainability, and collaboration [18, 26], but which may also confuse newcomers. Second, as we detail next, SPRs’ permissive naming practices enables adversaries to choose package names to deliberately confuse engineers.

**Table 1: Ecosystems examined in this study, highlighting their primary domains, naming conventions, and SOTA works on package confusion detections. Registry sizes are as of Jan. 2025. An example 1-level naming pattern is PyPI’s requests module. An example hierarchical name is Hugging Face’s google-bert/bert-base-uncased.**

Registry (# pkgs)	Domain	Name Structure	Prior Work
PyPI (619K)	Python	1-level	[41, 65]
RubyGems (212K)	Ruby	1-level	[41]
Maven (503K)	Java	Hierarchical	N/A
Golang (175K)	Go	Hierarchical	N/A
Hugging Face (1.1M)	AI Models	Hierarchical	[26]
NPM (5.1M)	JavaScript	Both	[41, 58]

### 2.2 Package Confusion Attacks and Taxonomy

*Package confusion* is an attack in which adversaries seek to persuade engineers to install malicious packages [41].<sup>1</sup> Attackers typically choose a package name similar to a legitimate one that meets the engineer’s needs [41]. Attacks may be multi-pronged, e.g., using advertising campaigns to increase their chances of fooling engineers [25, 51]. Package confusion attacks are an ongoing concern, with hundreds of packages detected by researchers in Socket [4], ReversingLabs [48], Orca [1], and others [15].

<sup>1</sup>Package confusion is the SPR-specific manifestation of the general naming confusion attack, which affects all IT endeavors including DNS (domain names) [20, 22, 62], BNS (Blockchain names) [39], and container images (image names) [10, 33].

Table 2 summarizes the state-of-the-art taxonomy of package confusion techniques, based on Neupane *et al.* [41]. The variations in naming structure in different SPRs also affect the nature of package confusion by SPR. In SPRs with 1-level names, as shown in prior work, confusion attacks look like the top portion of the table. In hierarchical SPRs (this work), the attacker has a larger naming surface to exploit (bottom of table).

**Table 2: Common package confusion techniques with examples. The top section presents some of the taxonomy proposed by Neupane *et al.* [41]. We added four patterns (bottom, cf. §4.1), generalizing the prior taxonomy to attacks in other SPRs. The examples for those four patterns are confusion threats flagged by ConfuGuard.**

Technique	Example (basis → attack)
1-step Levenshtein Distance	crypto → crypt
Sequence reordering	python-nmap → nmap-python
Scope confusion	@cicada/render → cicada-render
Semantic substitution	bz2file → bzip
Alternate spelling	colorama → colourama
<b>Impersonation Squatting</b>	meta-llama/Llama-2-7b-chat-hf → facebook-llama/Llama-2-7b-chat-hf
<b>Compound Squatting</b>	@typescript-eslint/eslint-plugin → @typescript-eslint/eslint
<b>Domain Confusion (Golang)</b>	github.com/prometheus/prometheus → git.luolix.top/prometheus/prometheus
<b>Command Squatting</b>	NPM i --help → NPM i help

### 2.3 Defenses Against Package Confusion

Package confusion attacks can be mitigated at many points in the reuse process, e.g., by scanning for malware in registries [64] or during dependency installation [31] or by sandboxing dependencies at runtime [14, 21, 63]. Such techniques provide good security properties, but are computationally expensive. Our approach is part of a class of cheaper techniques that leverage the most pertinent attribute of a package confusion attack: the similarity of the name to that of another package. We describe the literature in this vein.

**2.3.1 Approaches.** In the academic literature, the earliest name-oriented defenses were purely lexical: Taylor *et al.* [58] and Vu *et al.* [65] measured similarity using Levenshtein distance [32] to identify PyPI packages with minor textual alterations from one another. To detect confusion attacks based on name semantics (e.g., bz2file vs. bzip), Neupane *et al.* [41] applied FastText embeddings [16] to enable an abstract comparison of package names.

In acknowledgment of package confusion attacks, industry has also offered some defenses. NPM and PyPI both use lexical distances to flag packages or otherwise reduce user errors [42, 46, 70]. Some other companies also offer relevant security tools. For example, Stacklok combines Levenshtein distance with repository and author activity metrics to assign a risk score to suspicious packages [54]. Microsoft’s OSSGadget generates lexical string permutations and verifies the existence of mutated package names on SPRs to detect potential package confusion attacks [38].

**2.3.2 Knowledge gaps.** Despite these advances, key challenges persist that hinder the accuracy and reliability of current package confusion detection approaches. The primary concern was

articulated by Ohm *et al.* [43]: these techniques have high false-positive rates. Second, detailed empirical treatments have so far focused exclusively on package confusion attacks in NPM, PyPI, and RubyGems [41, 59, 65], leaving unknown the confusion attacks in SPRs with hierarchical naming structures. Finally, although summaries of package confusion defenses in production are reassuring, the lack of detailed production data makes their utility unclear.

### 3 PROBLEM STATEMENT

In this section, we give our definition of package confusion threats (§3.1), our threat model (§3.2), and our system requirements (§3.3).

#### 3.1 Refined Definition of Package Confusion

We refine the definition of package confusion threat by emphasizing the manifestation of malicious intent. Our refined definition is supported by our empirical analysis of real attacks in §4.1.

##### Definition: Package Confusion Threat

A **package confusion threat** is a software package whose name mimics a trusted resource, with the intent of deceiving developers into installing code that is actively malicious or may become so (stealthy).<sup>2</sup>

The mimicry may be *lexical* (i.e., typosquatting), *syntactic*, or *semantic*, but the intent is to confuse an engineer.

Because our definition considers threats, not just attacks, the current absence of malware does not imply future innocence. In contrast, prior work relies on malware detectors and only finds active attacks [41, 59]. Incorporating stealthy threats into this definition raises the possibility of high false positive rates. To address this, we analyzed false positives from prior package confusion data (§4.2) to identify signals that distinguish benign packages from stealthy threats.

#### 3.2 Threat Model

We describe the context and threats considered in this work.

**Context Model:** We focus on software package registries, including SPRs hosting traditional software packages and pre-trained AI models. The relevant properties are given in §2.1.

**Threat Model:** We include some threats and exclude others.

- *In-scope:* We consider attacks where packages with deceptively similar names to legitimate packages are published in the software package registry. As our definition permits, the attacker may initially publish non-malicious content (stealthy threat) and later introduce malware (active threat).
- *Out-of-scope:* We exclude attacks using non-confusing package names, or that compromise existing legitimate packages. These attacks must be mitigated by other measures [44].

This threat model is stronger than those of existing package confusion detection techniques [41, 58, 65], because we permit attacks to start with stealthy (non-malicious) package content.

<sup>2</sup>For comparison, Neupane *et al.* wrote: “Package confusion attack: a malicious package is created that is designed to be confused with a legitimate target package and downloaded by mistake” [41]. This requires the package to be actively malicious.

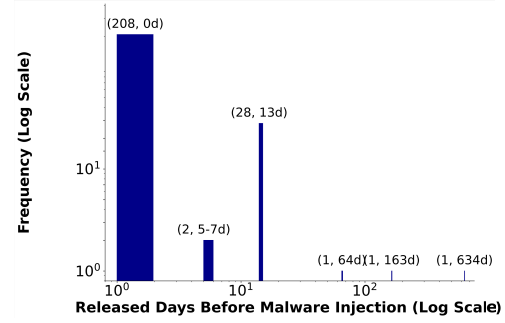


Figure 2: Distribution of time until malware, in packages with confusing names. 13.3% (32/240) of attacks occur  $\geq 5$  days after release.

#### 3.3 Requirements

A detection system must meet both security and operational goals. Working with our industry partner, we identified four requirements:

**Req<sub>1</sub>: High precision and recall.** A threat detection system must balance precision (low false positive rate) against recall (low false negative rate). False negatives enable attacks, but false positives lead to alert fatigue [13].

**Req<sub>2</sub>: Efficient and Timely Detection.** The system must scale to large registries — low-latency checks enabling real-time feedback.

**Req<sub>3</sub>: Compatibility Across Ecosystems.** The system must support many SPRs with their diverse naming schemes (§2.1).

**Req<sub>4</sub>: High Recall of Stealth Package Confusion Attacks.** The system must identify both active and stealth package confusions.

### 4 ANALYSIS OF CONFUSION ATTACKS

To address the high false positive rates of previous detection systems for package confusion, Ohm *et al.* recently called for a deeper analysis of the available evidence [43]. We respond to their call in two ways: describing attackers’ practices (§4.1), and reporting metadata features that distinguish true from false positives (§4.2).

#### 4.1 Attackers’ Practices

We examined first the role of stealth in confusion attacks, and then the attack variations in SPRs beyond NPM, PyPI, and RubyGems. Our findings informed requirements 1, 3, and 4 (§3.3).

**4.1.1 Stealthy Confusion Attacks.** We analyzed the versions of the 240 confirmed NPM package confusion “true-positives”, defined as confusing packages that include malware. These data were collected and reported by Neupane *et al.* [41]. We analyzed the number of days these malicious packages were available before malware was injected. We estimated the injection time based on the last version updated before removal by the SPR.

Figure 2 illustrates our findings on the release time before malware injection. Most package confusion attacks injected malware on the first day of release, e.g., in the initial version. However, in 13.3% of the packages, attackers undertook stealth attacks, deploying malware after days or weeks. Attackers may use this strategy to grow a userbase before exploitation. Prior work on package confusion did not account for this behavior.



4.1.2 *Extending the Taxonomy to New SPRs.* Prior work has examined malicious naming strategies in only a few SPRs. Our industry partner’s analysts have studied package confusion attacks in many other SPRs, including all those in Table 1. We discussed Neupane *et al.*’s taxonomy with them, and compared it to the real-world attack patterns they have observed. Through this dialogue, we identified four additional mechanisms for package confusion attacks. Definitions follow; see §2.2 for examples.

- **Impersonation Squatting.** In hierarchical SPRs, attackers impersonate legitimate maintainers or organizations by registering a similar author or groupId.
- **Compound Squatting.** Making multiple coordinated edits to a hierarchical name, e.g., altering both scope and delimiters at once. We consider this distinct because the compounded changes complicate conventional lexical comparisons.
- **Domain Confusion.** In SPRs where package names include URLs (golang), attackers may register domains resembling those of trusted mirrors or proxies.
- **Command Squatting.** Registering package names that mimic command-line options in other packages or utilities.

## 4.2 Metadata Features of Package Confusion

As discussed in §2.3, almost all package confusion detectors make use only of the names of packages. While this allows those systems to avoid costly content analysis, it leads to high false positive rates. We propose to integrate metadata to distinguish true from false positives. Although metadata signals can be circumvented, our approach aims to raise the attacker’s cost, reflecting the typical cat-and-mouse dynamics in open package ecosystems [30].

4.2.1 *Method.* We randomly sampled 626 packages from the original 640,482 false positives reported by Neupane *et al.* [41], which have suspicious package names but are not flagged as malware in their respective security analyses. This sample size gives a confidence level of 99% with a margin of error of 5% on the resulting distribution of causes.

To identify possible metadata signals in the existing package confusion data, we began by having two researchers — experts in software supply chains — analyze 200 of these packages. They analyzed each package’s metadata (READMEs, maintainers, versions, etc.). Each analyst independently proposed possible features based on this analysis (codebook) and then they discussed this codebook together to reach agreement on terms and definitions. To test validity, they then independently applied this codebook to the 200 packages and measured agreement using Cohen’s Kappa [19]. The initial inter-rater reliability was 0.6 (“substantial” [37]). The researchers subsequently discussed to resolve discrepancies and refine their analysis. Through discussion, they reached consensus on measurable attributes that could indicate malicious intent or the possibility of a stealthy attack.

Based on the high agreement in this process, one of these researchers analyzed the remaining 426 packages.

4.2.2 *Results.* Of the 626 packages sampled, 601 were still accessible. Among these, we identified 464 benign packages (77.2%) and what we believe are 137 stealthy attacks (22.8%). The progression of our results and the manually labeled dataset are available in

our artifact (§8). We identified 11 metadata features of the packages used in confusion attacks, as well as benign signals. These attributes include factors such as a distinct purpose, adversarial package naming, and the comprehensiveness of available metadata. See Table 3 for the features leveraged in ConfuGuard.

## 5 CONFUGUARD DESIGN AND IMPLEMENTATION

We introduce ConfuGuard, our novel detector for package confusion attacks. ConfuGuard is designed to meet the system requirements (§3.3). Addressing the insights from our empirical study (§4), ConfuGuard integrates both syntactic and semantic name analysis, hierarchical naming checks, and metadata-based verification to enhance threat detection and mitigation. ConfuGuard is intended for *backend* use, in package registries or similar platforms, and so it prioritizes accuracy (Req<sub>1</sub>) over latency (Req<sub>2</sub>).

Figure 3 shows the five main components of ConfuGuard:

- (1) A metadata database, ensuring real-time awareness of new and evolving packages (Req<sub>4</sub>).
- (2) A fine-tuned embedding database to capture domain-specific semantic name similarities (Req<sub>1,3</sub>).
- (3) Using popularity metrics to protect high-value targets (Req<sub>2,3</sub>).
- (4) Checking for syntactic and semantic confusion strategies (Req<sub>2</sub>).
- (5) Filtering out benign packages using package metadata (Req<sub>1,4</sub>).

The first three components provide infrastructure to support our analysis (components 4 and 5). Next, we detail the rationale and implementation of each part of ConfuGuard.

### 5.1 Part ①: Package Metadata Database

5.1.1 *Rationale.* SPRs grow rapidly, and package confusion attacks are initiated regularly. Comprehensive and regularly updated metadata ingestion enable early threat detection.

5.1.2 *Approach.* ConfuGuard relies on a database with package names, version histories, commit logs, license info, maintainer records, and other publicly-accessible metadata, all updated on a weekly basis. Regular updates reduce concerns about stale data and ensure our results are up to date. We specifically use our industry partner’s private database, which consolidates metadata from NPM, PyPI, RubyGems, Maven, Golang, and Hugging Face.

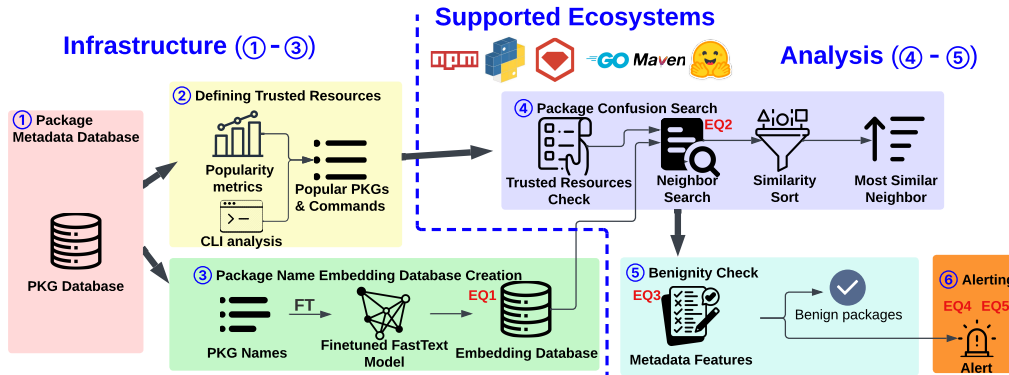
### 5.2 Part ②: Defining Trusted Resources

5.2.1 *Rationale.* Confusion attacks mimic trusted resources. If an attacker chooses a name that mimics an untrusted package, little harm can result. Untrusted packages can be compared against trusted ones to focus on the threats of highest potential impact.

5.2.2 *Approach.* Like prior works [41, 59], we operationalize trust in terms of popularity. Software engineers commonly make use of popularity signals as a proxy for trustworthiness [17], and thus attackers choose names similar to popular packages. Our specific popularity measure depends on the SPR. Many SPRs — in our case, NPM, PyPI, RubyGems, and Hugging Face — offer weekly or monthly download counts. Our other two supported SPRs, Maven and Golang, do not publish download metrics. For these, we make use of the `ecosyste.ms` database, which estimates popularity using indicators like stargazers, forks, and dependent repositories [40]. The thresholds were set based on production experience. We chose

**Table 3: Overview of the 13 metadata-based verification rules. Each rule includes a description of its purpose and the specific implementation steps taken to verify flagged packages. The final three rules (R12–R14) were added as part of our refinement process based on further observed false-positive patterns after we deployed our system in production.**

Rule	Description	Implementation
<b>R1: Obviously unconfusing</b>	Determine if a package name conveys an intentional, brand-specific identity (e.g., <code>catplotlib</code> ), clearly differentiating it from deceptive imitations.	Cross-reference flagged names with known legitimate projects and use LLM analysis to assess whether the naming convention is deliberate and unambiguous for developers.
<b>R2: Distinct Purpose</b>	Distinguish packages with different functionalities, even if names are superficially similar (e.g., <code>lodash-utils</code> vs. <code>lodash</code> ).	Extract package descriptions and calculate semantic similarity using TF-IDF cosine scores. A score below 0.5 indicates distinct purposes, reducing suspicion of deception.
<b>R3: Fork Package</b>	Detect benign forks sharing near-identical code or metadata with a popular package.	Compare README files, version histories, and file structures for high overlap. Similarities without malicious edits suggest harmless forks.
<b>R4: Active Development/Maintained</b>	Determine if packages are frequently updated or actively maintained by multiple contributors, which are less likely to have malicious intent.	Retrieve metadata for the last update, commit history, and version count. Classify packages with recent updates (e.g., within 30 days) or more than five versions as legitimate.
<b>R5: Comprehensive Metadata</b>	Identify packages missing critical metadata elements, such as licenses, maintainers, or homepages, which are typical of legitimate projects.	Check for the presence of licenses, contact details, and repository links.
<b>R6: Overlapped Maintainers</b>	Distinguish legitimate extensions or rebrands by verifying if the flagged package shares maintainers with the legitimate one.	Match maintainer identifiers (e.g., email, GitHub handle) between flagged and legitimate packages. Overlapping maintainers suggest legitimate intent.
<b>R7: Adversarial Package Name</b>	Filter out name pairs with significant length differences, as these often indicate unrelated projects rather than covert mimicry.	Compare string lengths of flagged and legitimate package names. A difference exceeding 30% indicates likely unrelated naming.
<b>R8: Well-known Maintainers</b>	Trust packages maintained by reputable and recognized authors/organizations.	Leverage knowledge in LLM training data to identify if a maintainer is trustworthy in the community.
<b>R9: Clear Description</b>	Verify that each package includes a README or equivalent documentation detailing its purpose, usage, and key features.	Analyze repository metadata to ensure there is a clear, summarized description of the package; packages lacking such documentation should be flagged for further review.
<b>R10: Has Malicious Intent</b>	Identify packages that deliberately mimic legitimate ones through near-identical descriptions, or have obviously suspicious content.	Use LLM to analyze the package name and description to triage whether there is obvious malicious intent in the package.
<b>R11: Experiment/Test Package</b>	Identify packages that are used for test or experiment purposes only.	Use LLM to analyze package descriptions and determine whether it is an experiment/test.
<b>R12: Package Relocation</b>	Account for legitimate package relocations, common in hierarchical registries like Maven.	Parse metadata files ( <code>pom.xml</code> ) for <code>&lt;relocation&gt;</code> tags or analogous fields. Identify and ignore renamed or migrated projects.
<b>R13: Organization Allowed List</b>	Prevent false positives by excluding packages published by trusted or verified organizations.	Maintain an allowedlist of approved organizations. If a flagged package is published under an allowed organization (e.g., <code>@oxc-parser/binding-darwin-arm64</code> ), it should be considered legitimate, comparing to <code>binding-darwin-arm64</code> .
<b>R14: Domain Proxy/Mirror</b>	Account for legitimate proxies or mirrors in the ecosystem, common in hierarchical registries like Golang.	Maintain an allowed list of recognized domains that serve as proxies or mirrors. If a given package is published under a valid domain (e.g., <code>gopkg.in/go-git/go-git</code> ), consider it legitimate when compared to its primary source ( <code>github.com/go-git/go-git</code> ).



**Figure 3: Overview of the ConFuGuard design with five primary components. The main novelty of ConFuGuard is the techniques of Part 4 and 5. We also improve on prior work in Part 2. and Part 1 (DB) can be constructed using public registry APIs. Part 3’s approach for popular packages is from prior work, adapted for the broader set of SPRs we handle [41, 58]. The red texts indicate the evaluation questions (EQs).**

5,000 weekly downloads for NPM, PyPI, and RubyGems; 1,000 for Hugging Face due to its costly usage; and an avg\_ranking score of 10 for Maven and Golang. We then adjusted the Golang thresholds to 4 to balance production latency.

Adversaries may inflate the popularity statistics of their packages [25], which could cause a malicious (or decoy) package to become trusted. If our partner’s analyst team flags any package as suspicious (§5.6), it is removed from the list of trusted packages.

As a secondary definition of trust, we analyzed the command-line interfaces of tooling for each SPR, e.g., the “npm” and “mvn” commands. These interfaces contain built-in keywords and commands which a careless engineer might confuse or mis-type for package names, e.g., mistaking `npm i help` (install a package named

`help`) for `npm i --help` (help command). All such keywords are included as trusted resources.

### 5.3 Part ③: Package Name Embeddings

**5.3.1 Rationale.** Detecting maliciously similar names requires accurately capturing subtle lexical, syntactic, and semantic variations. Traditional Levenshtein edit distance methods often fail to account for domain-specific semantic nuances (e.g., `meta-llama` vs. `facebook-llama`), while generic embedding models can introduce inaccuracies, resulting in higher false-positive or false-negative rates. Prior work apply embedding for word-wise semantic similarity check [41], while we would also like to use embedding for the whole name to support efficient neighbor search (§5.4). Robust,

fine-tuned embeddings address these shortcomings by providing enhanced semantic sensitivity, reducing erroneous alerts. Moreover, this embedding approach is also generic and unified to support various naming convention of SPRs.

**5.3.2 Approach.** An embedding fine-tuned on real package names enhances semantic sensitivity, enabling the detection of adversarial or suspicious names. This capability is integral to assessing the risk associated with a package and serves as a cornerstone of our *intent-centric* approach. Following [41], we build upon FastText [16], starting with the pre-trained model pre-trained on cc.en.300.bin, and *fine-tune* it using all (totally ~9.1 million) package names extracted from the metadata database in Nov. 2024 (§5.1). By fine-tuning, our version is expected to better capture domain-specific subwords.

We then use the fine-tuned embedding model to create an embedding vector database utilizing the vector format provided by pgvector due to its efficient vector operations for databases [11]. Given a package name, we remove its delimiters, pass the concatenated string to the embedding model, and store the extracted embedding in our database. We note that for hierarchical names from Golang (domain/author/repository), Maven (groupId:artifactId), and Hugging Face (author/model), we split names to create both author and package identifier embeddings per package. This approach mitigates high false-positive rates that occur when a single embedding causes similar packages from the same organization or author to be incorrectly classified as confusing in §5.4.

The complete embedding database occupies 24 GB, with each embedding vector corresponding to a single package name (or its author name and package identifier). This setup not only facilitates rapid query-based lookups but also supports subsequent steps in package neighbor searching. Visualization of the extracted embedding vectors is available in our artifact (§8).

## 5.4 Part ④: Package Confusion Search

Having obtained the appropriate infrastructure (Parts 1–3), we must now detect package confusion threats efficiently yet accurately. In Part 4, we efficiently flag possible threats, and then in Part 5 we filter out likely-benign ones.

**5.4.1 Rationale.** Modern SPRs are enormous, so any detection scheme must be designed with scalability in mind. Given the set of all packages in an ecosystem, we follow prior work [41, 58] by comparing the (relatively large) “long tail” of untrusted resources to the (relatively small) set of trusted resources.

### 5.4.2 Approach.

**Trusted Resources Selection.** We compare untrusted resources (*i.e.*, unpopular packages) against all names in the popular list. For packages in the trusted resources list (*i.e.*, popular packages), we compare it only with “more” trusted packages — for popularity, we use download rates  $\geq 10\times$  higher or ecosyste.ms scores  $\geq 2\times$  higher.

**Neighbor Search (Package Name Similarity Search).** We perform a distance search between each resource and the set of trusted resources. Syntactic similarity is measured using Levenshtein distance with a threshold of 2, while semantic similarity is determined using embedding distance. For semantic similarity, we apply a cosine similarity threshold of 0.93, as suggested in §6.2. For hierarchical names,

we adjust the similarity thresholds to 0.99 for package identifiers and 0.9 for author names to better detect compound squatting attacks. To efficiently measure semantic similarity, we leverage the *Approximate Nearest Neighbors* method with the *HNSW* index [36] in PostgreSQL. We opted for *HNSW* over *IVFFlat* based on benchmarking results that show faster search speeds and minimal vector perturbation [60]. By partitioning the embedding space into multiple clusters, the *HNSW* index limits distance computations to a smaller subset of candidate packages, rather than exhaustively comparing all pairs (**Req3**).

**Similarity Sort and Most Similar Neighbor(s).** Inspired by Neupane *et al.* [41], we define a similarity function that ranks neighbors based on Levenshtein distance, n-gram similarity, phonetics, substring matching, and fuzzy ratio. This sorting ensures the system accurately identifies the most likely attack target. If the previous steps produce any nearby neighbors, we consider these as possible attack targets. We treat the nearest neighbors as the most likely targets. The suspicious package and these nearest neighbors are propagated to the benignity check (§5.5). Currently we consider only the two nearest neighbors, balancing accuracy against speed.

## 5.5 Part ⑤: Benignity Filter

**5.5.1 Rationale.** Prior work used purely string-oriented methods, which often misclassify harmless or beneficial packages as confusion attacks. To mitigate false positives, we developed a metadata-driven benignity check inspired by a recent metadata-aware malware detector [56]. We seek to filter out legitimate packages based on explainable heuristics, avoiding unnecessary alerts.

**5.5.2 Approach.** Our analysis of false-positive data (§4.2) showed many reasons for which legitimate packages may have names resembling other trusted resources. We iteratively developed 10 rules that utilize metadata features (*e.g.*, version history, maintainers) to distinguish package confusion attacks from benign packages. We apply these rules to assess the benignity of fetched neighbors from §5.4 and reduce the false positive rate. We added several additional rules during deployment. During implementation, we identify rules that require description analysis (*e.g.*, R1, R2, R3) or additional knowledge (*e.g.*, R8). To handle these cases, we use an LLM to analyze each package and generate the corresponding rules. We provide the package description and relevant metadata as input, refining the prompt iteratively based on feedback from security analysts. The final prompts are available in our artifact (§8).

Table 3 summarizes the goals and implementation details for each rule. Once our name-based detector flags a package as suspicious, we retrieve its metadata and apply benignity check rules, using a simple weighted sum to calculate a risk score. We have explored both manually-assigned weights and learning them from regression. The simple scoring function facilitates feedback from analysts.

## 5.6 Part ⑥: Alerting in Production

**Use in Production.** ConfuGuard is integrated into our industry partner’s security scanning system for software packages. They had previously relied primarily on content-based malware scanners, but these are so computationally costly that they cannot be applied at scale. Their former approach for scalable detection of



confusion attacks was a lexical check (Levenshtein distance) which had an unacceptably high false positive rate. ConfuGuard replaced that approach and complements the malware scanners. A team of analysts triages alerts from ConfuGuard and the other scanners to make a threat assessment. If they deem a package to be malicious, that information is propagated to our partner’s customers through their security feeds.

*Experience-Driven Optimizations.* We trialed several versions and parameterizations of ConfuGuard over several months to enhance its performance. We share four examples. First, to improve accuracy across ecosystems, we tailored the neighbor search query (Part 4). For instance, lengthy identifiers, *e.g.*, Maven’s `artifact_id` and Golang’s `domain name`, diminish the effectiveness of embedding similarity. To address this, we opted to compute separate similarity measurements for each component of the package names and apply an edit distance threshold. If the difference exceeds two-thirds of the original name’s length, it is unlikely to confuse users. Second, within our threat model, we determined that changes to both identifiers in Maven, Golang, and Hugging Face were unlikely to confuse users and no such attacks were observed. Consequently, we excluded these from our algorithm. Third, we refined our metadata checks for benignity. We added rules R11, R12, and R13 based on observed patterns. Fourth, for performance, we experimented with quantized versions of the embedding model. Ultimately we found the latency gain a poor bargain for the resulting accuracy degradation.

## 6 EVALUATION

To evaluate ConfuGuard, we pose five Evaluation Questions (EQs) to assess its performance at the component level (EQ1–EQ3) and the system level (EQ4–EQ5). At the component level, we measure the effectiveness of novel mechanisms introduced in ConfuGuard. At the system level, we evaluate its integrated functionality and scalability compared to baseline tools (§6.1), and ability to detect real-world package confusion threats. Additionally, we compare our approach to SOTA methods to benchmark its effectiveness. An overview of the evaluation process is illustrated in Figure 3.

*Component-level.* We assess how individual components contribute effectively to the overall system.

- **EQ1: Performance of Embedding Model.** What is the accuracy and efficiency of our embedding model? (§5.2)
- **EQ2: Neighbor Search Accuracy.** How effective is the neighbor search approach? (§5.4)
- **EQ3: Metadata Verification Accuracy.** How much does the benignity filter reduce false positive rates? (§5.5)

*System-level.* We examine the performance of the full ConfuGuard system and compare to other approaches.

- **EQ4: Discovery.** Can ConfuGuard identify previously unknown package confusion threats?
- **EQ5: Baseline Comparison.** How does ConfuGuard perform in terms of accuracy and latency compared to SOTA tools?

All experiments run on a server with 32 CPU cores (Intel Xeon CPU @ 2.80GHz) and 256 GB of RAM. Notably, the training and fine-tuning of FastText models do not require GPUs.

## 6.1 Baseline and Evaluation Datasets

*State-of-the-Art Baselines.* We compare our system to the Levenshtein distance approach [65], and Typomind [41]. We consider OSSGadget [38] out of scope because it only handles lexical confusions and its latency for long package names makes it unsuitable for production. These are all state-of-the-art open-source tools.

*Evaluation Datasets.* We evaluate using two datasets:

- **NeupaneDB:** 1,840 packages from [41], including 1,239 confirmed real-world package confusion attacks, and 601 manually labeled data we analyzed in §4.2.
- **ConfuDB:** 1,561 packages triaged by security analysts, collected during the development and refinement of ConfuGuard (§5.6).

## 6.2 EQ1: Perf. of Embedding Model

We measured effectiveness and efficiency of our embedding model.

### 6.2.1 Effectiveness.

*Method.* We evaluate the effectiveness of embedding-based similarity detection by comparing three approaches:

- (1) *Levenshtein-Distance*, calculates the number of single-character edits required to change one package name to another.
- (2) *Pre-trained FastText* [16] (cc.en.300.bin), used in the SOTA work Typomind [41], employs the general-purpose embedding model cc.en.300.bin [16] to capture semantic relationships.
- (3) *Fine-tuned FastText (Ours)*, which we have adapted using a corpus specifically composed of package names to better capture domain-specific similarities.

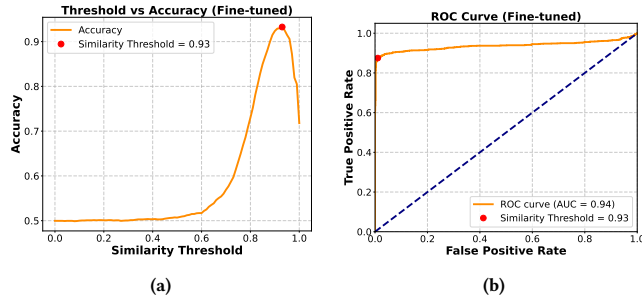
To systematically compare these methods, we construct a balanced test set consisting of both positive and negative pairs, with each category containing 1,239 data points.

- **Positive Pairs:** 1,239 real package confusions from NeupaneDB [41].
- **Negative Pairs:** Created by randomly pairing unrelated package names from registries, ensuring low similarity scores across all tested methods. These pairs are guaranteed not to represent package confusion relationships.

For each approach, we applied a predefined similarity threshold to classify package pairs as potential package confusion attacks. The threshold was selected via a grid search to optimize Precision and Recall, ensuring effective identification of true confusions. Pairs with similarity scores above the threshold were classified as positive, while those below were classified as negative. False positives were subsequently filtered using our false-positive verification process in Step ⑤. We compute Precision, Recall, and F1 scores for each approach to assess their performance.

**Table 4: Similarity effectiveness of using distance and embedding. We apply a grid search which automatically determines the optimal threshold to maximize F1 scores for positive pairs while maintaining relatively high F1 scores for negative pairs.**

Model	Positive Pairs			Negative Pairs			Overall Score
	P	R	F1	P	R	F1 Score	
Levenshtein Distance	1.00	0.80	0.89	1.00	1.00	1.00	0.94
Pretrained	1.00	0.85	0.92	1.00	0.98	0.99	0.96
Fine-Tuned (Ours)	1.00	<b>0.90</b>	<b>0.95</b>	1.00	0.96	0.98	<b>0.96</b>



**Figure 4: Performance Metrics of the Fine-Tuned Model: (a) Threshold Accuracy and (b) ROC Curve.**

*Result.* Table 4 presents the comparison between our fine-tuned model and baseline approaches. The results show that our approach outperforms both the Levenshtein Distance and Pretrained models, achieving the highest F1 score for positive pairs (0.95) while maintaining strong performance on negative pairs (0.98).

Figure 4 illustrates the performance metrics of our approach through a grid search, highlighting the relationship between similarity threshold and accuracy, as well as the ROC curve for the fine-tuned model. The optimal similarity threshold of 0.93, identified through this search, is used in the production system. With an AUC of 0.94, our model demonstrates a balanced trade-off between precision and recall for both positive and negative pairs.

#### 6.2.2 Efficiency.

*Method.* To assess the efficiency of embedding database creation, we evaluated the overhead of embedding database creation step (§5.3). We evaluated the embedding database creation efficiency by measuring the throughput, latency, and overall overhead associated.

**Table 5: Evaluation of embedding model efficiency, HNSW indexing overhead, and total embedding database creation overhead.**

Ecosystem	Throughput (pkgs/s)	Latency (ms)	Indexing Time (s)	Total Overhead (s)
NPM	7705.99	0.13	7.42	650.19
PyPI	13581.15	0.07	4.39	46.85
RubyGems	8694.68	0.12	4.07	28.18
Maven	4887.19	0.20	4.91	134.47
Golang	8490.36	0.12	7.80	232.93
Hugging Face	5027.95	0.20	4.46	166.67

*Result.* Table 5 presents the efficiency measurement. Overall, the quantized models demonstrate strong performance with minimal overhead. The HNSW indexing adds negligible processing time, requiring less than 10 seconds additional overhead per table.

### 6.3 EQ2: Neighbor Search Accuracy

*Methods.* We use a threshold of 0.93 obtained from EQ1. To evaluate neighbor search performance, we analyzed suspicious packages identified by Typomind.

*Results.* Our neighbor search algorithm accurately detected 99% of the 1,239 real-world package confusion attacks from NeupaneDB [41]. Moreover, our method effectively flagged impersonation squatting attacks targeting hierarchical package names, notably identifying

reported cases on Hugging Face [45] and Golang [5] that earlier methods overlooked. These findings confirm that our neighbor search algorithm achieves state-of-the-art performance.

### 6.4 EQ3: Benignity Check Accuracy

*Method.* Starting with the raw embedding output (i.e., packages flagged for name-based suspicion), we apply the benignity check described in §5.5.

We first use cross-validation to learn the weights and measure false positive rates (FPR) using NeupaneDB. We then computed SHAP values by learning the weighed sum through regression to quantify the importance of each metric defined in Table 2 [35]. The resulting SHAP plot provides clear insights into how these metrics contribute to our system’s decision-making process. Additionally, we use regression to learn the weights

*Result.* Figure 5 shows the SHAP value plot which indicates that our rules used in the benignity filter are useful. The plot highlights that distinct purpose (R2) and suspicious intent (R10) are the most influential features driving the model’s predictions, suggesting that packages showing suspicious or unusual behavior strongly push the model toward a malicious classification. Conversely, features such as known maintainers (R8) and fork packages (R3) generally pull the model’s output toward benign, indicating they lower the likelihood of a threat. Meanwhile, no clear description also shows a notable positive effect, aligning with the idea that missing documentation can be an indicator of malicious intent (R9).

We use cross-validation to learn parameters and estimate the trade-off between reducing false positives and increasing false negatives. Our benignity filter correctly flagged 399 out of 464 (86%) false positives in the human labeled set. In the real-world attack set, it failed to classify 53 true positives out of 443.<sup>3</sup> A detailed investigation revealed two main causes: (1) missing metadata and (2) attacks that had been removed, with the package now maintained by npm, which the filter considers a valid maintainer.

These results confirm that supplementary heuristics beyond raw name similarity reduce false positives while retaining high recall for genuinely malicious confusions.

### 6.5 EQ4: Discovery of New Package Confusions

*Method.* To evaluate the effectiveness of our package confusion detection system, we deployed ConFuGuard in a production environment for three months. During this period, flagged packages were analyzed using a commercial malware scanner and reviewed by threat analysts for detailed insights.

*Result.* Table 6 presents the package confusion attack that ConFuGuard identified during three months deployment period (i.e., ConFuDB). Within this database, 1,260 packages were labeled as false positives.<sup>4</sup> Among the remaining 301 suspicious packages, 230 (76%) exhibited stealthy confusion attack behavior, 63 (21%) were identified as malware, 6 (2%) were categorized as vulnerabilities, and 2 (0.7%) fell into other anomaly categories.

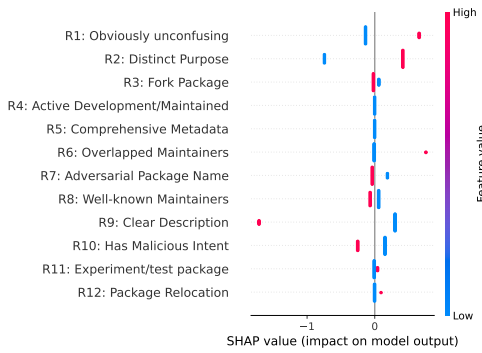
<sup>3</sup>Only 443 out of 1,239 real attacks still have available metadata for our analysis.

<sup>4</sup>The high false positive rate in ConFuDB is primarily due to the difficulty of confirming stealthy confusion attacks, which were left unreviewed by analysts.



**Table 6: Accuracy metrics for ConfuGuard and Typomind, detailing true/false positives and negatives, as well as recall, precision, F1 score, and accuracy for active, stealthy, and benign threats. Although Typomind flags all packages as suspicious, our benignity filter substantially reduces false positives in production. Due to availability of metadata, we evaluated ConfuGuard on 1,043 packages from NeupaneDB and 1,202 ConfuDB. We note that ConfuGuard exhibits lower benign classification accuracy in ConfuDB because it flags more complex cases (e.g., compound squatting) as potential threats, a nuance that Typomind does not capture.**

Dataset	Total	Threat Type	Sub-total	ConfuGuard								Typomind							
				TP	FP	TN	FN	Recall	Prec.	F1	Acc.	TP	FP	TN	FN	Recall	Prec.	F1	Acc.
NeupaneDB	1,840	Active	1,239	386	0	0	57	0.87	1.00	0.93	0.87	1000	0	0	239	0.80	1.00	0.89	0.81
		Stealthy	137	80	0	0	56	0.59	1.00	0.74	0.59	137	0	0	0	1.00	1.00	1.00	1.00
		Benign	464	0	72	391	0	0.00	0.00	0.00	0.84	0	464	0	0	0.00	0.00	0.00	0.00
ConfuDB	1,561	Active	71	41	0	0	4	0.91	1.00	0.95	0.91	40	0	0	31	0.56	1.00	0.72	0.56
		Stealthy	230	110	0	0	24	0.82	1.00	0.90	0.82	111	0	0	113	0.50	1.00	0.66	0.50
		Benign	1,260	0	532	491	0	0.00	0.00	0.00	0.48	0	322	612	0	0.00	0.00	0.00	0.65



**Figure 5: SHAP value plot for metadata verification using learnable features (i.e., excluding R13, R14) from Table 3. The x-axis shows each feature’s impact on model output (positive values are towards benignity). Feature values are color-coded: red indicates high values with a strong, direct impact on the model output.**

Analysis of production data shows that ConfuGuard effectively detects advanced threats (§4.1) – including compound, impersonation, and command squatting – across platforms such as Maven, Golang, and Hugging Face. For example, it identified an *impersonation squatting* attack on Maven: the malicious package `io.github.leetcnunch:scribejava-core` mimicked the legitimate `com.github.scribejava`. Although the package content appears identical at first glance, it contains code that injects a network call to steal user credentials. Similarly, our system detected malware on a Golang package, where the username `boltdb-go` attempted to impersonate `boltdb`. ConfuGuard also flagged suspicious Hugging Face packages, such as `TheBlock/Mistral-7B-Instruct-v0.2-GGUF`, mimicking benign package, `TheBloke/Mistral-7B-Instruct-v0.2-GGUF`, which has 92K monthly downloads. We were also able to capture *domain confusions*, such as Go packages from `github.phpd.cn` and `github.hscsec.cn` which look like proxies of official Golang packages but contain malware that either sets up a rogue MySQL server to steal sensitive files or intercepts and downgrades secure traffic to redirect data to attacker-controlled endpoints. Additionally, we uncovered a compound squatting attack on `@typescript-eslint/eslint-plugin`,

where an attacker used a similar namespace and package identifier (`@typescript-eslint/eslint`) to mislead users.

## 6.6 EQ5: System Accuracy and Efficiency

**Method.** We evaluated the end-to-end system accuracy and efficiency (i.e., latency and throughput) by running ConfuGuard and Typomind on NeupaneDB and ConfuDB.

**Result.** Table 6 presents the accuracy measurements. Across the datasets, ConfuGuard generally exhibits stronger performance in detecting active threats, achieving higher recall than Typomind while both maintain perfect precision. In the case of stealthy threats, the results are mixed: Typomind excels in the NeupaneDB dataset with a perfect recall, yet its performance declines in ConfuDB where ConfuGuard outperforms it. For benign cases, the accuracy differs notably, as Typomind misclassifies benign packages in NeupaneDB, whereas it performs better on ConfuDB, though still lower than ConfuGuard’s performance in that category. Overall, ConfuGuard demonstrates strengths in detecting active and stealthy threats across datasets, reducing the FPR from 77% (§4.2) to 13% (72/543).

**Table 7: System latency and throughput measurements. ConfuGuard has higher latency due to the additional benignity filter. Our industry partner considers this performance overhead acceptable.**

	ConfuGuard	Typomind
Latency (ms/pkg)	6816	120
Throughput (#pkgs/s)	0.08	0.26

Table 7 indicates that while ConfuGuard incurs significantly higher latency (6816 ms per package vs. 120 ms for Typomind) and lower throughput (0.08 vs. 0.26 packages/s), these trade-offs are acceptable in production because the additional latency—stemming from benignity filtering and reliance on LLM-based metadata verification—substantially enhances detection accuracy by reducing false positives. Despite slower processing due to dependencies on the OpenAI API and o3-mini, the system scales effectively across large registries, and future work may optimize inference times without compromising accuracy. Notably, ConfuGuard significantly reduces the workload on analysts, resulting in notable cost savings.

## 7 DISCUSSION

### 7.1 Three Lessons Learned from Production

*Lesson 1: ConfuGuard prevents real confusion attacks.* Over the past two months, we deployed the system in our industry partner’s production environment, during which we identified and confirmed 301 threats, with 37,266 additional threats under review.

*Lesson 2: False Positives Harm Our Customers.* Table 3 represents the most recent metadata verification rules. In deployment, we found additional cases where the package has a very similar name and the README of their package was missing which made our system classify it as suspicious stealth attack. That package was owned by one of our partner’s customers and served as a “transitive package”. The customer felt that being flagged harmed their reputation. We added an allow-list tied to our partner’s customers.

*Lesson 3: Ontology Matters.* ConfuGuard is deployed as part of a security analysis pipeline, with human analysts to triage the results (§5.6). Analysts need to be able to interpret the results of ConfuGuard. The analysts feel that the package confusion categories outlined in Table 2 do not provide sufficient explanatory power. We propose to refine the confusion taxonomy to consider the malicious content and intent behind each package, including a risk-level classification.

### 7.2 Limitations and Security Analysis

This section discusses our system limitations and how attackers might bypass ConfuGuard.

*Gaming Metrics.* Our system relies on software metrics to gauge the likelihood that a package is a confusion attack. These metrics might be gamed. There has been little formal study of the feasibility of gaming these metrics, but recent work suggests both the possibility and some real-world examples [25].

*Limitations in Neighbor Search.* One significant limitation of ConfuGuard is its inability to handle short names or acronyms. *FastText* struggles with short words (e.g., `xml` vs. `yml` have a similarity score of 0.7). The model’s reliance on character  $n$ -grams often fails to capture subtle similarities effectively in such cases, providing an avenue for attackers to exploit short package names. To mitigate this, we implemented a list of potential substitutions to identify cases of visual or phonetic ambiguity, and we combine embedding similarity with Levenshtein distance. This hybrid approach improves neighbor search for short names, but increases the computational cost does not fully resolve the concern.

*Bypassing Metadata Verification.* Using an LLM for benign filtering introduces correctness risks (hallucination, jailbreaking [68]). Adversaries might exploit this by tailoring their package name or metadata to persuade the LLM that the package is trustworthy using techniques like prompt injection or model hijacking [34, 71].

### 7.3 Future Directions

*Enhancing Representations for Package Confusion Detection:* Improving the representation of package names is crucial for more robust detection. While *FastText* captures semantic similarity through subword embeddings, it struggles with typographical variations, particularly for short words or acronyms. Fine-tuning *FastText* or

training a more efficient model on domain-specific corpora including both correct and misspelled terms can address these limitations. Augmenting training data with synthetic typos and incorporating typo normalization or correction techniques before embedding generation can significantly reduce errors. Advanced models, such as transformer-based architectures fine-tuned with contrastive learning on typo-specific datasets, present a promising alternative for enhancing detection accuracy and reliability. This approach has proven effective in combating domain typosquatting, but no research has been conducted targeting package typosquatting [29]. One challenge is the limited availability of data for verifying package squatting cases. LLMs might help here [61].

*Mitigating LLM Hallucination in Code Generation.* The increasing use of LLMs for code generation has introduced new challenges, as these models often hallucinate package names or generate commands for nonexistent or maliciously similar packages [53]. These hallucinations pose serious threats to the security of the software supply chain [61, 66]. Addressing this issue requires implementing typo or hallucination correction mechanisms in LLM-based package recommenders. Verifying package legitimacy, detecting typos, and integrating contextual checks can prevent the propagation of incorrect package names, reducing the risks associated with hallucinations.

*Meta-Learning for Malicious Package Detection.* Meta-learning approaches offer significant potential for improving malicious package detection. By leveraging anomaly detection techniques and metadata analysis [24], systems can dynamically adapt to evolving attack strategies. Meta-learning frameworks could analyze patterns across registries and rapidly identify emerging threats, enhancing the scalability and robustness of detection systems. Integrating such frameworks will be key to staying ahead of increasingly sophisticated attackers.

## 8 CONCLUSION

We present ConfuGuard, an embedding-based package confusion detection system. Based on real-world attack patterns, we refined the package confusion definition and developed a taxonomy based on engineering practices. ConfuGuard is being used in production at our industry partner and contributed to 301 confirmed package confusion threats in three months. Compared to SOTA methods, our system is good at capturing additional confusion categories, achieves a substantially lower false-positive rate, and maintains acceptable latency, making it well-suited for deployment on SPR backends while remaining effective for frontend on-demand requests. We shared our insights from production experience, customer feedback, the need for an improved ontology, and outlined future directions.

## DATA AVAILABILITY AND RESEARCH ETHICS

Our artifact is available at: <https://github.com/confuguard/confuguard>. For replicability, we include all results and the code of the system, except the commercial metadata database (§5.1).

Our work poses limited ethical concerns. We analyzed public packages and disclosed suspicious ones to our industry partner.

# REFERENCES

- [1] Dependency confusion: Supply chain attacks. <https://orca.security/resources/blog/dependency-confusion-supply-chain-attacks/>. Accessed: 2025-03-12.
- [2] Guide to maven naming conventions. <https://maven.apache.org/guides/mini/guide-naming-conventions.html>.
- [3] Hugging face naming limitations. <https://hf.co/help/r/jfrog-artifactory-documentation/hugging-face-naming-limitations>.
- [4] Malicious npm package typosquats popular typescript eslint plugin. <https://socket.dev/blog/malicious-npm-package-typosquats-popular-typescript-eslint-plugin>. Accessed: 2025-03-12.
- [5] Malicious package exploits go module proxy caching for persistence. <https://socket.dev/blog/malicious-package-exploits-go-module-proxy-caching-for-persistence>. Accessed: 2025-03-12.
- [6] npm package json: name. <https://docs.npmjs.com/cli/v9/configuring-npm/package-json#name>.
- [7] Package names. <https://go.dev/blog/package-names>.
- [8] Pep 423 – metadata for python software packages 2.0. <https://peps.python.org/pep-0423/>.
- [9] Rubygems.org: Gem naming conventions. <https://guides.rubygems.org/rubygems-org-gem-naming/>.
- [10] Whoami: A cloud image name confusion attack. <https://securitylabs.datadoghq.com/articles/whoami-a-cloud-image-name-confusion-attack/>.
- [11] pgvector: A vector extension for postgresql. <https://github.com/pgvector/pgvector>, 2023.
- [12] General Services Administration (GSA) 18F. 18f open source policy. <https://18f.gsa.gov/open-source-policy/>, 2025.
- [13] Agency for Healthcare Research and Quality. Alert fatigue. <https://psnet.ahrq.gov/primer/alert-fatigue>, 2019.
- [14] Paschal C Amusuo, Kyle A Robinson, Tanmay Singla, Huiyun Peng, Aravind Machiry, Santiago Torres-Arias, Laurent Simon, and James C Davis. ZTDJVA: Mitigating software supply chain vulnerabilities via zero-trust dependencies. In *ICSE*, 2025.
- [15] Alex Birsan. Dependency confusion: How i hacked into apple, microsoft and dozens of other companies, 2021.
- [16] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- [17] Hudson Borges and Marco Tulio Valente. What’s in a GitHub Star? Understanding Repository Starring Practices in a Social Coding Platform. In *Journal of Systems and Software (JSS)*, 2018.
- [18] Simon Butler, Michel Wermelinger, Yijun Yu, and Helen Sharp. Exploring the influence of identifier names on code quality: An empirical study. In *2010 14th European Conference on Software Maintenance and Reengineering*, pages 156–165. IEEE, 2010.
- [19] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.
- [20] Conor Coyle. How google chrome can identify typos in urls to keep you safe online, 2023.
- [21] Gabriel Ferreira, Limin Jia, Joshua Sunshine, and Christian Kastner. Containing Malicious Package Updates in npm with a Lightweight Permission System. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, pages 1334–1346. IEEE, May 2021. Place: Madrid, ES.
- [22] Evgeniy Gabrilovich and Alex Gontmakher. The homograph attack. *Communications of the ACM*, 2002.
- [23] Yacong Gu, Lingyun Ying, Yingyuan Pu, Xiao Hu, Huajun Chai, Ruimin Wang, Xing Gao, and Haixin Duan. Investigating package related security threats in software registries. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 1578–1595. IEEE, 2023.
- [24] Sajal Halder, Michael Bewong, Arash Mahboubi, Yinhao Jiang, Md Rafiqul Islam, Md Zahid Islam, Ryan HL Ip, Muhammad Ejaz Ahmed, Gowri Sankar Ramachandran, and Muhammad Ali Babar. Malicious package detection using metadata information. In *Proceedings of the ACM on Web Conference (WWW’24)*, pages 1779–1789, 2024.
- [25] Hao He, Haoqin Yang, Philipp Burckhardt, Alexandros Kapravelos, Bogdan Vasilescu, and Christian Kästner. 4.5 million (suspected) fake stars in github: A growing spiral of popularity contests, scams, and malware. *arXiv:2412.13459*, 2024.
- [26] Wenxin Jiang, Chingwo Cheung, George K Thiruvathukal, and James C Davis. Exploring naming conventions (and defects) of pre-trained deep learning models in hugging face and other model hubs. *arXiv:2310.01642*, 2023.
- [27] Wenxin Jiang, Nicholas Synovic, Matt Hyatt, Taylor R. Schorlemmer, Rohan Sethi, Yung-Hsiang Lu, George K. Thiruvathukal, and James C. Davis. An empirical study of pre-trained model reuse in the hugging face deep learning model registry. In *IEEE/ACM 45th International Conference on Software Engineering (ICSE’23)*, 2023.
- [28] Berkay Kaplan and Jingyu Qian. A Survey on Common Threats in npm and PyPi Registries, August 2021. *arXiv:2108.09576 [cs]*.
- [29] Takashi Koide, Naoki Fukushima, Hiroki Nakano, and Daiki Chiba. Phishreplicant: A language model-based approach to detect generated squatting domain names. In *Proceedings of the 39th Annual Computer Security Applications Conference*, pages 1–13, 2023.
- [30] Piergiorgio Ladisa, Henrik Plate, Matias Martinez, and Olivier Barais. SoK: Taxonomy of Attacks on Open-Source Software Supply Chains. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 1509–1526, May 2023.
- [31] Jasmine Latendresse, Suhaib Mujahid, Diego Elias Costa, and Emad Shihab. Not all dependencies are equal: An empirical study on production dependencies in npm. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2022.
- [32] Vladimir I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966.
- [33] Guannan Liu, Xing Gao, Haining Wang, and Kun Sun. Exploring the uncharted space of container registry typosquatting. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 35–51, 2022.
- [34] Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Zihao Wang, Xiaofeng Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, et al. Prompt injection attack against llm-integrated applications. *arXiv:2306.05499*, 2023.
- [35] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 2017.
- [36] Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836, 2018.
- [37] Mary L McHugh. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3):276–282, 2012.
- [38] Microsoft. Microsoft ossgadget. <https://github.com/microsoft/OSSGadget>.
- [39] Muhammad Muzammil, Zhengyu Wu, Lalith Harisha, Brian Kondracki, and Nick Nikiforakis. Typosquatting 3.0: Characterizing squatting in blockchain naming systems. *arXiv:2411.00352*, 2024.
- [40] Andrew Nesbitt. Ecosyste.ms database: A comprehensive dataset for software ecosystem analysis, 2025.
- [41] Shradha Neupane, Grant Holmes, Elizabeth Wyss, Drew Davidson, and Lorenzo De Carli. Beyond typosquatting: an in-depth look at package confusion. In *USENIX Security ’23*, pages 3439–3456, 2023.
- [42] npm, Inc. Threats and mitigations. <https://docs.npmjs.com/threats-and-mitigations#by-typosquatting--dependency-confusion>, 2024.
- [43] Marc Ohm and Charlene Stuke. Sok: Practical detection of software supply chain attacks. In *Proceedings of the 18th International Conference on Availability, Reliability and Security*, pages 1–11, 2023.
- [44] Chinenye Okafor, Taylor R. Schorlemmer, Santiago Torres-Arias, and James C. Davis. SoK: Analysis of Software Supply Chain Security by Establishing Secure Design Properties. In *Proceedings of the 2022 ACM Workshop on Software Supply Chain Offensive Research and Ecosystem Defenses*, pages 15–24, Los Angeles CA USA, November 2022. ACM.
- [45] ProtectAI. Unveiling ai supply chain attacks on hugging face, n.d. Accessed: 2025-01-15.
- [46] Python Software Foundation. Acceptable use policy. <https://policies.python.org/pypi.org/Acceptable-Use-Policy/>, 2024.
- [47] Eric S. Raymond. *The cathedral & the bazaar: musings on Linux and open source by an accidental revolutionary*. O’Reilly, Beijing ; Cambridge, Mass, 1st ed edition, 1999.
- [48] ReversingLabs. R77 toolkit: Typosquatting attack in npm ecosystem, 2023.
- [49] Taylor R Schorlemmer, Kelechi G Kalu, Luke Chigges, Kyung Myung Ko, et al. Signing in four public software package registries: Quantity, quality, and influencing factors, 2024.
- [50] Lee Joon Sern and Yam Gui Peng David. Typoswype: An imaging approach to detect typo-squatting. In *2021 11th IFIP International Conference on New Technologies, Mobility and Security (NTMS)*, pages 1–5. IEEE, 2021.
- [51] Ax Sharma. PyPI crypto-stealer targets windows users, revives malware campaign. *Sonatype Blog*, May 2024.
- [52] César Soto-Valero, Amine Benelallam, Nicolas Harrand, Olivier Barais, and Benoit Baudry. The Emergence of Software Diversity in Maven Central. In *International Conference on Mining Software Repositories (MSR)*, 2019.
- [53] Joseph Spracklen, Raveen Wijewickrama, AHM Sakib, Anindya Maiti, and Mur-tuza Jadhwal. We have a package for you! a comprehensive analysis of package hallucinations by code generating llms. *arXiv:2406.10279*, 2024.
- [54] Stacklok. Detecting typosquatting attacks on open source packages using levenshtein distance and activity data. <https://www.blackduck.com/resources/analyst-reports/open-source-security-risk-analysis.html>, n.d.
- [55] John Stenbit. Open source software (oss) in the department of defense (dod). *Department of Defense, Memorandum*, 2003.
- [56] Xiaobing Sun, Xingan Gao, Sicong Cao, Lili Bo, Xiaoxue Wu, and Kaifeng Huang. 1+ > 2: Integrating deep code behaviors with metadata features for malicious pypi package detection. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pages 1159–1170, 2024.
- [57] Synopsys, Inc. 2024 open source security and risk analysis report. <https://www.blackduck.com/resources/analyst-reports/open-source-security->



- [risk-analysis.html](#), 2024.
- [58] Matthew Taylor, Ruturaj Vaidya, Drew Davidson, Lorenzo De Carli, and Vaibhav Rastogi. Defending against package typosquatting. In *Network and System Security: 14th International Conference (NSS)*.
- [59] Matthew Taylor, Ruturaj K Vaidya, Drew Davidson, Lorenzo De Carli, and Vaibhav Rastogi. Spellbound: Defending against package typosquatting. *arXiv:2003.03471*, 2020.
- [60] Tembo. Vector indexes in pgvector. <https://tembo.io/blog/vector-indexes-in-pgvector>, 2024.
- [61] Christopher Tozzi. Package hallucination: The latest, greatest software supply chain security threat? *IDC Blog*, April 22 2024.
- [62] Marcin Ulikowski. dnstwist: Domain name permutation engine for detecting homograph phishing attacks, typo squatting, and brand impersonation. <https://github.com/elceef/dnstwist>, 2025.
- [63] Nikos Vasilakis, Ben Karel, Nick Roessler, Nathan Dautenhahn, Andre DeHon, and Jonathan M. Smith. BreakApp: Automated, Flexible Application Compartmentalization. In *Proceedings 2018 Network and Distributed System Security Symposium*, San Diego, CA, 2018. Internet Society.
- [64] Duc-Ly Vu, Fabio Massacci, Ivan Pashchenko, Henrik Plate, and Antonino Sabetta. LastPyMile: identifying the discrepancy between sources and packages. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 780–792, Athens Greece, August 2021. ACM.
- [65] Duc-Ly Vu, Ivan Pashchenko, Fabio Massacci, Henrik Plate, and Antonino Sabetta. Typosquatting and combosquatting attacks on the python ecosystem. In *EuroS&P Workshops*, pages 509–514. IEEE, 2020.
- [66] Vulcan Cyber. Can you trust chatgpt’s package recommendations? <https://vulcan.io/blog/ai-hallucinations-package-risk>, April 17 2023.
- [67] Erik Wittern, Philippe Suter, and Shriram Rajagopalan. A look at the dynamics of the JavaScript package ecosystem. In *International Conference on Mining Software Repositories (MSR)*, 2016.
- [68] Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaxing Song, Ke Xu, and Qi Li. Jailbreak attacks and defenses against large language models: A survey. *arXiv:2407.04295*, 2024.
- [69] Markus Zimmermann, Cristian-Alexandru Staicu, and Michael Pradel. Small World with High Risks: A Study of Security Threats in the npm Ecosystem. In *USENIX Security Symposium*, 2019.
- [70] Tzachi Zornstein and Yehuda Gelb. A new, stealthier type of typosquatting attack spotted targeting NPM. *Checkmarx Blog*, May 2023.
- [71] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv:2307.15043*, 2023.