



國立臺北科技大學

創新資安碩士班

碩士學位論文

**S2GE-NIDS: A hybrid architecture
combining structured semantics and
generation embedded network intrusion
detection system in IoT**

研究生：周玟萱

指導教授：陳香君博士

中華民國一百一十四年五月



國立臺北科技大學

創新資安碩士班

碩士學位論文

**S2GE-NIDS: A hybrid architecture
combining structured semantics and
generation embedded network intrusion
detection system in IoT**

研究生：周玟萱

指導教授：陳香君博士

中華民國一百一十四年五月

「學位論文口試委員會審定書」掃描檔

審定書填寫方式以系所規定為準，但檢附在電子論文內的掃描檔須具備以下條件：

1. 含指導教授、口試委員及系所主管的完整簽名。
2. 口試委員人數正確，碩士口試委員至少 3 人、博士口試委員至少 5 人。
3. 若此頁有論文題目，題目應和書背、封面、書名頁、摘要頁的題目相符。
4. 此頁有無浮水印皆可。

Abstract

Keyword: IoT Security, Information Security, Anomaly Detection, Multilayer Perceptron, Semantic Vector

As network environments become increasingly complex and dynamic, traditional intrusion detection methods struggle to keep pace with evolving threats and high-volume traffic. This paper proposes an efficient anomaly detection framework that leverages hash-based token embedding and a lightweight multi-layer perceptron (MLP) for the semantic representation of network flows. By transforming feature values into semantic tokens and utilizing a hashing trick for embedding lookup, our approach enables scalable and robust processing without maintaining an explicit vocabulary. The resulting embedding vectors are flattened and processed by the MLP to produce semantic vectors, which are clustered using a center loss strategy for unsupervised anomaly detection. Experimental results on public benchmark datasets demonstrate that our method achieves competitive accuracy with significantly improved computational efficiency compared to traditional attention-based models.

Table of Contents

Abstract	i
Table of Contents	ii
Chapter 1 Introduction	1
Chapter 2 Related Work	3
2.1 Network Intrusion Detection System in IoT	3
2.2 Tokenization	4
2.3 Hash Embedding	5
2.4 Multi-Layer Perceptron in Anomaly Detection	6
2.5 Semantic Vector	6
Chapter 3 Methodology	8
3.1 Architecture	8
3.1.1 Preprocess Model	9
3.1.2 Embedding Model	10
3.1.3 Mahalanobis Distance Model	13
3.2 Flow	15
3.2.1 Preprocess Model	15
3.2.2 Implementation Procedure	15
Chapter 4 Implementation	19
4.1 Experimental Setup	19
4.1.1 Hardware Requirements	19
4.1.2 Software Requirements	19
Chapter 5 Conclusion & Future Work	24
5.1 Conclusion	24
5.2 Future Work	24

List of Figures

3.1	Architecture of S2GE-NIDS	8
3.2	FlowChart for Preprocess Model	16
3.3	Hash Embedding for Embedding Model	17
4.1	Download on Official Anaconda Website	21
4.2	Installation for Anaconda	21
4.3	FlowChart for Preprocess Model	22
4.4	Visual Studio Code	22



List of Tables

2.1	Common Anomalous Features in IoT Network Traffic and Their Descriptions . .	4
2.2	Examples of Field-Value Tokenization in IoT Network Traffic	5
3.1	Example of Tokenized Input Fields	10
4.1	Hardware Requirements	19
4.2	Software and Libraries Used in the Experiment	20



Chapter 1 Introduction

Driven by the rapid advancement of digital transformation and smart infrastructure, the **Internet of Things (IoT)** has emerged as a cornerstone of next-generation information technology. Through the integration of sensors, embedded devices, communication modules, and platform software, IoT enables physical objects to communicate in real time and generate massive volumes of data. These data streams support a broad range of applications—such as smart manufacturing, intelligent transportation, remote healthcare, and smart homes—yielding substantial economic and societal value [1].

However, as the number of connected devices increases and deployment scenarios become more complex, IoT systems face unprecedented cybersecurity challenges. Many IoT devices are resource-constrained, infrequently updated, and difficult for users to manage. With limited encryption and a lack of monitoring mechanisms, these devices become prime targets for cyber intrusions and attacks. Effectively identifying abnormal behaviors and hidden threats in IoT network traffic has therefore become a pressing research priority.

Furthermore, existing intrusion detection technologies often struggle to adapt to evolving threats. While deep learning approaches such as Word2Vec and Transformer-based models [devlin2018bert-va have demonstrated semantic learning capabilities, they also introduce critical drawbacks: large vocabulary requirements, high computational complexity, and limited flexibility in dynamic or resource-constrained environments.

To address these limitations, we propose **S2GE-NIDS** (Structured Semantics and Generation Embedded Network Intrusion Detection System)—a lightweight, interpretable anomaly detection framework designed for IoT environments. S2GE-NIDS combines hash-based token embedding with a multi-layer perceptron (MLP) model and introduces a linked-list mechanism to mitigate hash collisions inherent to non-cryptographic hash functions such as MurmurHash3 [2]. This design enables efficient feature encoding while avoiding the need to maintain a large vocabulary.

In our approach, network packets are first transformed into semantic tokens and encoded using hash-based indexing. The resulting embedding vectors are concatenated into a single, fixed-length semantic vector, which is processed by an MLP and projected near a learned semantic

center. Any significant deviation from this center—measured by Mahalanobis distance—is classified as a potential anomaly [3].

The proposed S2GE-NIDS framework offers several key advantages over conventional intrusion detection systems. First, it eliminates the need for manual feature engineering and vocabulary maintenance by using a hash-based embedding approach, where field-value pairs are directly encoded into semantic vectors without relying on predefined lookup tables. This design greatly simplifies the preprocessing pipeline and enhances scalability. Second, the model provides a mathematically interpretable anomaly scoring mechanism by integrating Mahalanobis distance, which quantifies how far a sample deviates from the learned distribution of normal behavior. This not only improves detection accuracy but also enables explainable results. Third, the system is lightweight and highly efficient, relying on simple MLP-based encoding instead of complex deep architectures, making it well-suited for deployment in real-time or resource-constrained environments such as edge devices in IoT networks. Lastly, its generalized tokenization strategy allows for wide applicability across diverse packet structures, further improving its adaptability and robustness in various network scenarios.

The structure of this paper is as follows: Chapter 2 is the relevant background knowledge about S2GE-NIDS (Structured Semantics and Generation Embedded Network Intrusion Detection System). Chapter 3 introduces the architecture and methodology of the proposed S2GE-NIDS framework, presenting each module and its rationale in detail. Chapter 4 presents the experimental setup, evaluation metrics, and results on two benchmark datasets, as well as interpretability demonstrations. Chapter 5 summarizes the main findings, limitations, and directions for future research.

Chapter 2 Related Work

This section will introduce the relevant basic knowledge, including existing IoT network intrusion detection methods, Tokenization, Hash Embedding and language tags.

2.1 Network Intrusion Detection System in IoT

In recent years, the proliferation of Internet of Things (IoT) devices has led to an increased focus on developing effective network intrusion detection systems (NIDS) tailored to the specific characteristics of IoT environments. Various approaches have been proposed to address the challenges associated with high-volume, heterogeneous network traffic, constrained device capabilities, and evolving attack patterns. Kharoubi et al. [4] proposed NIDS-DL-CNN, a convolutional neural network (CNN)-based detection system designed for IoT security. By applying CNN layers to extract spatial features from traffic data, the model achieved high classification performance on datasets such as CICIOT2023 and CICIOT2024. The authors demonstrated that their method achieved excellent precision and recall in both binary and multi-class scenarios. However, a notable limitation of the CNN-based approach lies in its inability to fully capture temporal dependencies across packet sequences, and its reliance on supervised learning requires extensive labeled datasets. Ashraf et al. [5] introduced a real-time intrusion detection system (INIDS) based on traditional machine learning classifiers applied to the BoT-IoT dataset. The study compared seven algorithms, including Random Forest, Artificial Neural Networks (ANN), and Support Vector Machines. Their results showed that Random Forest and ANN achieved the highest accuracy and robustness among all tested classifiers. Despite its efficiency, the INIDS system was highly dependent on manual feature engineering and lacked adaptability to novel threats, which are critical in fast-evolving IoT environments. Elrawy et al. [6] conducted a comprehensive survey of intrusion detection methodologies in IoT-based smart environments, categorizing techniques according to architectural design (centralized vs. distributed), detection strategy (signature-based, anomaly-based, or hybrid), and system layer (perception, network, application). While the survey provided valuable insights and synthesized a broad range of IDS approaches, it lacked im-

plementation evidence and empirical comparisons, limiting its utility for practical system design. Collectively, these studies highlight the trade-offs between detection performance, computational cost, and deployment feasibility. Deep learning models offer strong accuracy but demand computational resources, while traditional classifiers provide efficiency but often lack flexibility. In contrast, our proposed S2GE-NIDS framework leverages hash-based semantic embeddings and a lightweight MLP, offering a balanced approach that combines scalability, interpretability, and effectiveness in detecting network anomalies in resource-constrained IoT environments.

Table 2.1 Common Anomalous Features in IoT Network Traffic and Their Descriptions

Feature (with Reference)	Description
Destination Port [tang2019deep]	Specific port targets (e.g., 22, 23, 80, 443) are often associated with attacks. Abnormal access to these ports may suggest behaviors such as scanning, DDoS, or brute-force intrusion.
Flow Duration [7]	Extremely short or long connection durations within brief timeframes may signal scanning activity or data exfiltration.
Total Forward Packets [8]	Unusually high or low packet counts in one direction may indicate abnormal sessions or flooding behavior.
Packet Length [9]	Anomalies in packet size—whether fixed, too long, or too short—often reflect malicious traffic like botnet propagation or worms.
Protocol Type [7]	Sudden increases in uncommon protocols (e.g., ICMP, UDP) may reveal attempts to exploit protocol vulnerabilities or bypass filters.
Source IP / Destination IP [10]	Repeated access from abnormal IP addresses, or sudden surges in novel IP sources, are indicative of scanning, spoofing, or DDoS activity.
Flow Bytes per Second [9]	Sharp fluctuations—surges or drops—in flow byte rate may suggest DoS attacks or unauthorized data transfer.
TCP Flags [8]	Unusual combinations (e.g., SYN, FIN, RST) can indicate stealth scans or TCP-based flooding.
Number of Connections [9]	A large number of new connections established by a single IP in a short time often reflects worm propagation or botnet coordination.

2.2 Tokenization

Following feature extraction, the next critical step is the tokenization process, which prepares network traffic data for semantic embedding. Each data record typically contains multiple fields—such as Port, Protocol, and SrcIP—representing structural and behavioral attributes of a network flow. To ensure consistency and distinguishability among features during embedding, we adopt a composite tokenization strategy that combines each field name with its corresponding value to form a unique token string. For instance, a sample token may take the form

Protocol:TCP or DstPort:443.

This strategy preserves the semantic association between field-value pairs without relying on predefined vocabularies, making it particularly suitable for dynamic and heterogeneous IoT environments. Each composite token is subsequently encoded using hash-based mapping techniques ??, thereby eliminating the need for extensive memory allocation or manually constructed token dictionaries. By treating each token as a self-contained semantic unit, this method also enhances the model's ability to generalize to previously unseen feature combinations, ultimately improving both the adaptability and scalability of the proposed system [11].

Table 2.2 Examples of Field-Value Tokenization in IoT Network Traffic

Feature Field	Tokenized Representation
Protocol = TCP	Protocol:TCP
Destination Port = 443	DstPort:443
Source Port = 80	SrcPort:80
Source IP = 192.168.0.1	SrcIP:192.168.0.1
Flow Duration = 120000	FlowDuration:120000
Payload Bytes = 56	PayloadBytes:56
Packet Count = 10	PacketCount:10
Flag = ACK	Flag:ACK
Protocol = ICMP	Protocol:ICMP
Destination IP = 10.0.0.5	DstIP:10.0.0.5

2.3 Hash Embedding

Hash Embedding is a common lightweight feature encoding technology, which is particularly suitable for structured, high-dimensional, or large-number-of-categories network data. Its core approach is to convert each field name/field value (or a combination of the two) into a set of indexes through a hash function (such as MurmurHash3), and query the embedding table to obtain a fixed-length semantic vector. The main method is to combine the (field name, field value) of each data sample and pass it through a hash function such as MurmurHash3 to obtain a set of row/col indexes. This set of indices is then used to query a multi-dimensional embedding table, where an initial random, trainable semantic vector is stored at each position. The multi-field embedding vectors are concatenated (flattened) or aggregated, and then the data is given to the anomaly detection model for learning and inference. Weinberger et al. [12] proposed Feature

Hashing to solve the coding efficiency problem of high-dimensional sparse data. L. Zhu et al. [11] combined Feature Hashing with a multi-layer perception for IoT intrusion anomaly detection and proved that it can significantly reduce the number of model parameters and improve computing efficiency. MurmurHash3 is widely used to replace traditional hashing techniques because of its uniform distribution, fast calculation, and consistency across languages.

2.4 Multi-Layer Perceptron in Anomaly Detection

Multi-Layer Perceptrons (MLPs) have been widely applied in the field of anomaly detection due to their capability to model non-linear relationships between input features and hidden patterns. Unlike traditional statistical models that rely on predefined thresholds or assumptions about data distribution, MLPs are capable of learning complex, high-dimensional feature representations in a data-driven manner [13].

In recent years, MLP-based anomaly detection methods have been employed in various domains, including network security [14], industrial control systems [15], and IoT environments [16]. These models typically consist of multiple fully connected layers with nonlinear activation functions, such as ReLU or sigmoid, enabling the learning of hierarchical semantic features. The outputs are used to distinguish between normal and abnormal behavior based on reconstruction error, classification scores, or learned distance metrics.

While MLPs are not as expressive as deep convolutional or recurrent models, their low computational cost and ease of deployment make them particularly attractive for lightweight and real-time anomaly detection systems. In our work, we leverage an MLP-based encoder to transform hash-embedded feature vectors into semantic representations, which are then evaluated using Mahalanobis distance for effective anomaly scoring.

2.5 Semantic Vector

Semantic vector representations, originally popularized in natural language processing (NLP), have gained traction in anomaly detection tasks due to their ability to encode complex contextual

information into fixed-length embeddings. In security-related applications, raw network traffic often contains heterogeneous features that lack explicit semantics; transforming these into semantic vectors enables better generalization and interpretability [17].

Recent works have applied semantic encoding strategies, such as Word2Vec and sequence embeddings, to convert protocol names, IP addresses, or header fields into high-dimensional vectors [18, 19]. These semantic vectors capture latent relationships between fields and behaviors, allowing downstream models to detect subtle deviations from normal patterns. For instance, Shapira et al. [18] proposed Flow2Vec, which encodes sequences of network events into dense vectors, improving anomaly detection in encrypted traffic.

Compared to one-hot encoding or manually crafted features, semantic vectors provide a richer and more scalable representation, particularly when combined with deep learning models. In this work, we construct semantic vectors from tokenized field-value pairs using a hash-based embedding scheme followed by an MLP encoder. This method ensures that semantic relationships among network features are preserved while maintaining computational efficiency.

Chapter 3 Methodology

In this session, we will introduce the S2GE-NIDS (structured semantics and generation embedded network intrusion detection system) architecture and details its operational workflow, clearly delineating each step from semantic tokenization through anomaly detection and decision-making processes.

3.1 Architecture

S2GE-NiDS is presented as Figure 3.1 including preprocess model, embedding model, and Mahalanobis model.

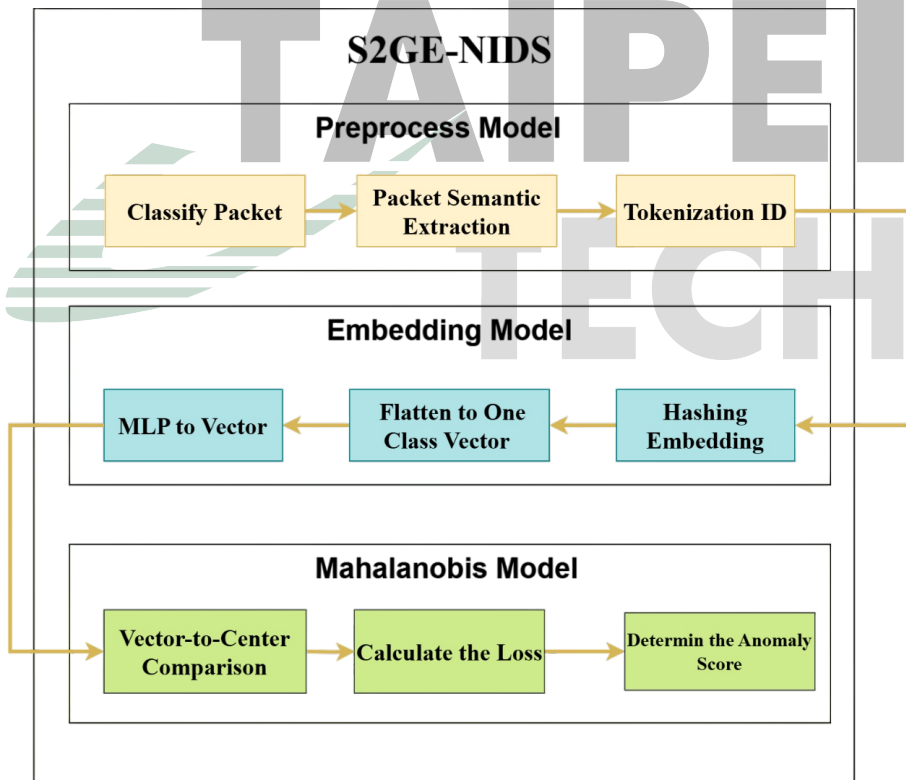


Figure 3.1 Architecture of S2GE-NIDS

During the preprocessing stage, relevant features are first extracted from network data packets and transformed into a combination of textual and numerical tokens. To prevent duplication and ensure a more uniform distribution within the embedding space, the model applies the non-encrypted MurmurHash3 function to encode each token.

To further mitigate the risk of hash collisions, a modulo operation is performed on the resulting hash values, which are then used to index into the embedding table. This strategy reduces the likelihood of different tokens being mapped to the same location, thereby enhancing both the accuracy and efficiency of the embedding process.

Next, the individual embedding vectors corresponding to each feature are concatenated into a single one-dimensional vector. This flattened vector is fed into a multi-layer perceptron (MLP) model, which transforms it into a compact semantic representation.

Finally, in the Mahalanobis distance evaluation phase, the semantic vector is compared to a predefined center point in the learned semantic space. If the computed distance exceeds a specified threshold, the sample is flagged as an anomaly. Evaluation metrics such as the F1 score are then used to quantify detection performance.

3.1.1 Preprocess Model

In the preprocessing phase, we will do the following process as data file selection and filtering, feature extraction, and tokenization. These steps are designed to transform raw network traffic into structured representations suitable for semantic embedding and anomaly detection.

3.1.1.1 Classify Data Packet

The first step in the preprocessing pipeline involves selecting and filtering the data files to ensure suitability for subsequent analysis. In this study, network traffic is collected and stored in the Comma-Separated Values (CSV) format—a widely adopted and flexible tabular data structure. CSV files are particularly well-suited for structured data representation due to their ease of parsing, compact storage, and seamless integration with mainstream data analysis libraries such as pandas and NumPy in Python. During this stage, only those CSV files containing the required packet-level features are retained, while incomplete, irrelevant, or malformed files are systematically excluded.

3.1.1.2 Packet Semantic Extraction

After the data is cleaned and organized, the first step is to extract meaningful features from the network packet data to better capture the characteristics of each packet. For example, we focus on key fields such as Destination Port, Protocol, and SrcIP, which are widely used in previous studies to detect abnormal network behavior. In practice, we use Python tools to read each CSV file and select these important features as the main input for the S2GE-NIDS model. By focusing only on these key values, we can make the data cleaner and easier for the anomaly detection system to use.

3.1.1.3 Tokenization ID

After the relevant features have been extracted, the next step is to perform tokenization, which converts structured data into a format suitable for semantic embedding. Each data entry consists of multiple fields—such as Destination Port, Protocol, and SrcIP—that represent different aspects of network behavior.

Tokenization is achieved by concatenating each field name with its corresponding value to form a unique string representation. This composite token serves as the semantic unit used in downstream embedding processes. For example, tokens follow the format *"field name + field value"*, as illustrated in Table 3.1.

Table 3.1 Example of Tokenized Input Fields

Field Name	Field Value
Destination Port	80
Flow Duration	192.168.1.2
Protocol Type	TCP

3.1.2 Embedding Model

To mitigate redundancy in text features during the embedding process, we employ a lightweight, non-cryptographic hash function—MurmurHash3. This function ensures that input tokens are more uniformly distributed across the embedding space, thus reducing overrepresentation in specific regions. To further minimize the probability of hash collisions—i.e., multiple tokens being

mapped to the same position—we apply a modulo operation to the resulting hash values. This yields a deterministic index used to locate or store each feature vector within a fixed-size embedding table, enhancing both the accuracy and efficiency of the overall embedding process.

Once all relevant token embeddings are retrieved, their vectors are concatenated into a single flattened, one-dimensional feature vector. This unified representation is then fed into a multi-layer perceptron (MLP), which learns high-level semantic abstractions and generates a compact semantic feature vector. The subsequent subsections provide detailed descriptions of the hash embedding mechanism, the flattening procedure, and the structure of the MLP used for semantic encoding.

3.1.2.1 Hash Embedding

Hash embedding is a lightweight vectorization technique that utilizes non-cryptographic hashing to encode tokenized field-value pairs into fixed-size, trainable embeddings [11]. In this study, we adopt the MurmurHash3 algorithm—an efficient and widely used hash function—to map each token to a specific position in the embedding table. Its advantages include fast computation, uniform distribution, and language-independent implementation, which make it well-suited for scalable anomaly detection in IoT environments [2].

To determine the target index for each token, we apply a modulo operation to the hash value using the smallest three-digit prime number, 233. This approach distributes tokens more evenly within the embedding space and reduces collision rates. For example, the token generated from the field name PORT may yield a MurmurHash3 value of 4283257230. Applying $4283257230 \bmod 233$ results in 56. If the associated port number (e.g., 405) is similarly hashed and gives a value with mod 233 result of 7, these indices (row 7, column 56) are used to locate the corresponding vector in the embedding table.

Each embedding vector is initially randomized and refined during training. For instance, an example 8-dimensional vector might be:

$$[-0.982, -0.301, -0.555, 2.061, 0.045, -0.618, -0.786, 0.573]$$

These vectors are later concatenated and passed to the MLP model for further semantic encoding.

3.1.2.2 Flatten

Flatten will string the tokenized data into a single vector through the vectors after the embedding column. For example, Destinaation port 405 is $[-0.982, -0.301, -0.555, 2.061, 0.045, -0.618, -0.786, 0.573]$ and Protocol TCP tokeniz to $[-0.024, 0.494, 0.754, -0.78, -1.002, 0.069, -0.52, -1.336]$,SrcIP $[-1.042, -0.116, 0.542, -0.987, 1.001, 0.086, 0.699, -0.903]$ $[-0.982, -0.301, -0.555, 2.061, 0.045, -0.618, -0.786, 0.573, -0.024, 0.494, 0.754, -0.78, -1.002, 0.069, -0.52, -1.336, -1.042, -0.116, 0.542, -0.987, 1.001, 0.086, 0.699, -0.903]$

3.1.2.3 MLP

After generating semantic embeddings from each tokenized field, the resulting vectors are flattened into a single one-dimensional input vector. This unified semantic representation is then fed into a lightweight **Multi-Layer Perceptron (MLP)** to learn deeper semantic relationships and perform nonlinear transformation for anomaly detection.

The MLP adopted in this work is composed of an input layer, one or more hidden layers, and an output feature vector. Each layer consists of a fully connected network of neurons activated by the ReLU (Rectified Linear Unit) function, which enables the model to capture non-linear dependencies among input features while maintaining computational efficiency.

To prevent overfitting and enhance generalization, dropout layers are incorporated after each dense layer, and batch normalization is applied to stabilize the training process. The final output of the MLP is a low-dimensional semantic vector projected in a latent space. This vector is subsequently used for statistical anomaly detection based on its distance from a learned semantic center.

Compared to more complex architectures such as transformers or recurrent models, the MLP achieves an optimal balance between *expressiveness*, *interpretability*, and *computational efficiency*, making it particularly suitable for deployment in IoT environments with limited resources [20]

3.1.3 Mahalanobis Distance Model

In the final stage of the S2GE-NIDS framework, we apply a statistical distance-based method—**Mahalanobis Distance**—to evaluate whether an observed semantic vector deviates significantly from the expected distribution of normal traffic. This metric is particularly effective for high-dimensional anomaly detection, as it accounts for feature correlations and variance [21].

Let $\mathbf{x} \in \mathbb{R}^n$ denote the semantic vector output from the MLP, and let $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ represent the mean vector and covariance matrix estimated from a subset of benign (normal) training data. The Mahalanobis distance is defined as:

This formulation enables the model to assess how far a sample deviates from the learned semantic center under multivariate normality assumptions. During inference, if $D_M(\mathbf{x})$ exceeds a predefined threshold τ , the corresponding traffic instance is flagged as an anomaly.

We empirically determine τ using the distribution of distances in the training set, often by selecting a percentile threshold (e.g., 95th percentile). This thresholding strategy is advantageous in unsupervised or semi-supervised settings, where labeled anomaly samples may be scarce.

The integration of Mahalanobis scoring into our system introduces the benefits of model interpretability and statistical rigor, effectively enhancing the ability to detect subtle but semantically meaningful deviations in IoT network behavior.

3.1.3.1 Vector-to-Center Comparison

To enhance anomaly detection, S2GE-NIDS introduces a center loss mechanism. During training, all semantic vectors corresponding to “normal” samples are aggregated to calculate a center point c .

- Taking into account the variability and correlation of each feature, the model can more accurately detect abnormal samples that are “off-center”.

$$D_M(z) = \sqrt{(z - c)^T \boldsymbol{\Sigma}^{-1} (z - c)} \quad (3.1)$$

z is the semantic vector of the input sample, c is the center vector of normal samples, and Σ^{-1} is the inverse of the covariance matrix of the training data's embedding vectors.

3.1.3.2 Calculate the Loss

- The loss is defined as:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \|z_i - c\|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d (z_{ij} - c_j)^2 \quad (3.2)$$

$z_i \in \mathbb{R}^d$ is the embedding vector obtained after the i th input passes through the Semantic Encoder, $c \in \mathbb{R}^d$ is the center point vector during training (center), and N is the total number of samples.

3.1.3.3 Determine the Anomaly Score

After obtaining the semantic vector z of each input data point through the MLP encoder, and computing the center point c based on all normal training samples, the system evaluates how far each sample deviates from the normal data distribution using the Mahalanobis distance metric.

The Mahalanobis distance score $D_M(z)$, as defined in Equation 3.1, quantifies the distance between a sample's semantic representation z and the center vector c , while accounting for the variance and covariance of the embedding space. This distance serves as the anomaly score for each sample.

$$D_M(z) = \sqrt{(z - c)^T \Sigma^{-1} (z - c)} \quad (3.3)$$

To determine whether a sample is anomalous, we define a threshold τ based on the distribution of distances observed in the training data. A sample is classified as anomalous if its Mahalanobis distance exceeds this threshold:

$$\text{Anomaly}(z) = \begin{cases} 1 & \text{if } D_M(z) > \tau \\ 0 & \text{otherwise} \end{cases} \quad (3.4)$$

Here, τ can be determined in several ways, such as:

- Using the mean plus k standard deviations from the training distribution (e.g., $\tau = \mu + k\sigma$).
- Setting τ based on a desired false-positive rate (e.g., the 95th percentile of $D_M(z)$ on normal samples).

This threshold-based mechanism enables the system to make binary decisions (normal vs. anomalous) while preserving the interpretability and statistical grounding of the anomaly scores. Additionally, ranked anomaly scores $D_M(z)$ can be used in top- k selection scenarios for prioritizing the most suspicious samples in real-time applications.

3.2 Flow

3.2.1 Preprocess Model

In this section shown in the system receives the uploaded network packet data and checks if the data format matches the CSV (Comma-Separated Values) format. If the data is not in CSV format, the system will prompt the user to re-upload the data. In the next stage, the system cleans the data fields, including removing any missing or empty fields from the packets.

Subsequently, specific fields related to common anomaly detection features are extracted, such as Destination Port, Protocol Type, and Source IP (SrcIP). These fields serve as important inputs for subsequent model analysis.

Finally, the field names and their respective values are combined into tokens—for instance, Protocol_TCP or Port_80—and fed into a semantic embedding model to be transformed into vectors for further processing.

3.2.2 Implementation Procedure

The system implementation is divided into several sequential stages, including: data pre-processing, feature transformation, semantic embedding, and anomaly detection. The detailed procedure is as follows:

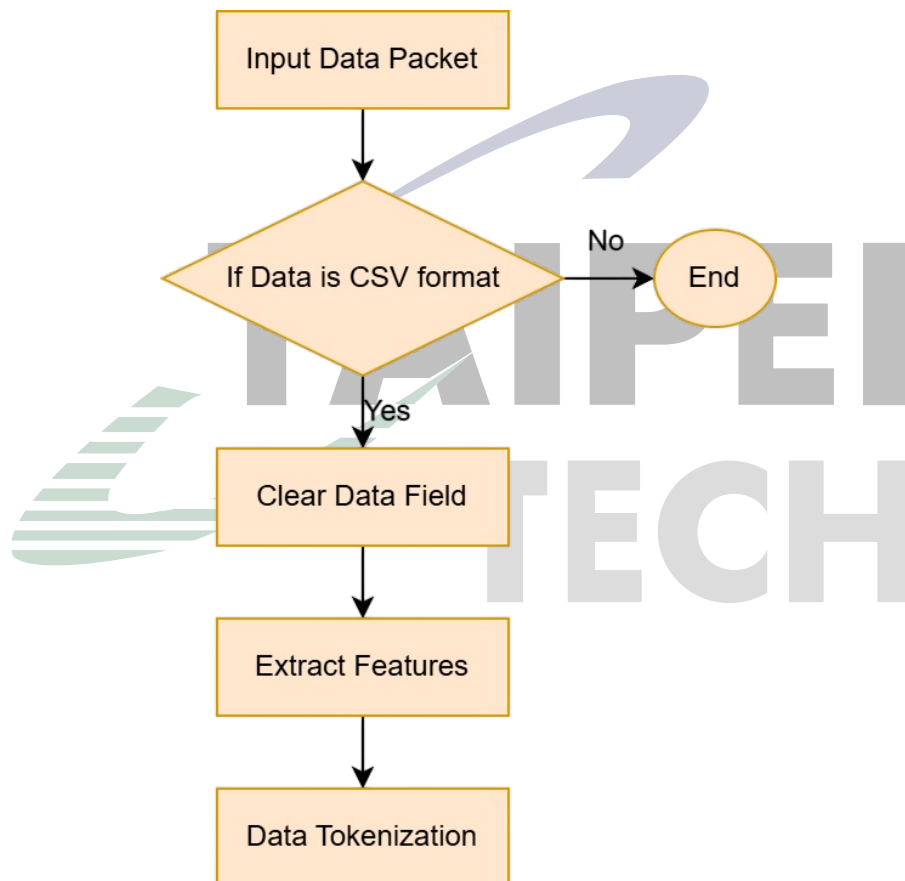


Figure 3.2 FlowChart for Preprocess Model

Step 1. Input Network Packet Data The raw data used in this study is stored in Comma-Separated Values (CSV) format. This format is chosen due to its ease of parsing and manageability. The dataset contains detailed information about various network packets.

Step 2. Data Cleaning and Filtering After loading the dataset, the preprocessing phase is initiated. This phase includes data cleaning and filtering. The system removes missing values and clearly anomalous packet records to avoid introducing bias or errors in subsequent processing stages.

Step 3. Feature Selection Based on insights from related research, specific key feature fields are selected from the raw packet data to serve as input for the model. These primarily include destination port, communication protocol, and source IP address. The selection is made based on the features' effectiveness in distinguishing anomalous events.

Step 4. Feature Tokenization To enhance the model's ability to process both textual and numerical features, each feature field is combined with its corresponding value to form semantically meaningful tokens. For instance, a destination port of 80 is converted into the token "Port_80".

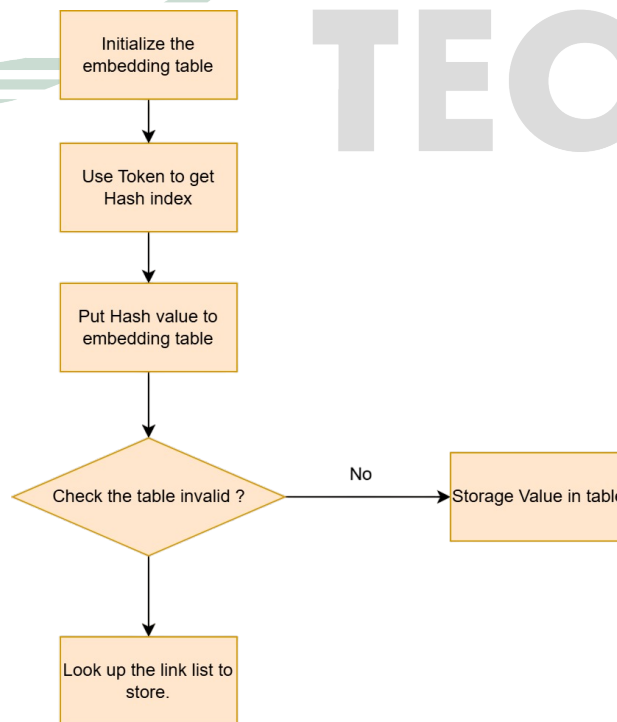


Figure 3.3 Hash Embedding for Embedding Model

Step 5. Feature Hash Mapping Since tokens often exhibit high redundancy and dimensionality, each token is passed through the MurmurHash3 hashing function to generate a fixed-length

hash value. This reduces dimensionality and ensures even distribution, thereby improving computational efficiency and mitigating potential collision issues.

Step 6. Constructing the Embedding Table Each hashed token is assigned a randomly initialized embedding vector. These vectors are stored in an embedding table that maps discrete hash values into a continuous vector space, facilitating the generation of semantic representations.

Step 7. Semantic Embedding All token vectors associated with a given packet are concatenated and flattened into a single vector, which is then passed through a Multi-Layer Perceptron (MLP) to perform nonlinear transformation. The output is a semantically enriched vector representation of the packet.

Step 8. Establishing the Normal Sample Center Vector To enable effective anomaly detection, the system computes the mean vector of all semantic vectors derived from normal packets during the training phase. This average vector serves as the center representation of normal traffic and is denoted as c .

Step 9. Anomaly Scoring via Mahalanobis Distance During the detection phase, the system calculates the Mahalanobis distance $D_M(z)$ between the semantic vector z of each incoming packet and the normal center vector c . If the computed distance exceeds a predefined threshold, the packet is classified as anomalous.

This systematic implementation procedure enables accurate and efficient anomaly detection on network packet data.

Chapter 4 Implementation

4.1 Experimental Setup

The experimental implementation of this study was conducted on the Windows 11 operating system. Visual Studio Code (VS Code) was utilized as the primary development environment, integrated with the Anaconda distribution for Python to manage package dependencies and virtual environments. A range of scientific computing and machine learning packages were installed to facilitate algorithm development, model training, and evaluation workflows. Detailed configuration steps and setup instructions are described in the following subsection.

4.1.1 Hardware Requirements

Table 4.1 provides detailed specifications and purposes of each hardware component utilized in our experimental environment.

Table 4.1 Hardware Requirements

Component	Specification
CPU	12th Gen Intel(R) Core(TM) i5-12500H @ 2.50 GHz
RAM	16.0 GB (15.6 GB usable)
Storage	Built-in SSD (operating system and model storage)

4.1.2 Software Requirements

Table 4.2 lists the software used in our experimental setup, along with their purposes and license types.

4.1.2.1 Environment Setup

Step 1: Installing Anaconda

Anaconda is an open source Python platform designed for data science and machine learning development, integrating the most commonly used data analysis tools and libraries. It has a rich

Table 4.2 Software and Libraries Used in the Experiment

Software	Version	Purpose	License
Visual Studio Code [22]	1.89.1	A lightweight and extensible code editor used as the primary integrated development environment (IDE) for editing Python scripts and managing project structure.	MIT
Anaconda Prompt [23]	2024.02	A command-line interface provided by the Anaconda distribution, used for managing Python virtual environments and installing dependencies via Conda or pip.	BSD
Python [24]	3.9.18	The main programming language used to implement the core modules of the proposed system, including preprocessing, model training, and evaluation routines.	Python License
NumPy [25]	1.26.4	Provides high-performance array structures and functions for numerical computing, especially efficient vector and matrix operations.	BSD
Pandas [26]	2.2.2	Offers powerful data manipulation and analysis tools, including DataFrame structures used for preprocessing and filtering packet data.	BSD
Scikit-learn [sklearn]	1.4.2	Provides a wide range of machine learning algorithms, particularly the Multi-Layer Perceptron (MLP) classifier used in this study.	BSD
mmh3 [27]	4.0.1	Implements MurmurHash3, a fast non-cryptographic hashing function used to convert tokens into integer values for embedding.	MIT
PyTorch [28]	2.2.2+cpu	A deep learning framework used to define and train neural networks, including custom embedding and classification models.	BSD

built-in data science suite, including core tools such as Numpy (numerical operations), Pandas (data processing), and Seaborn (data visualization).¹

Go to the official Anaconda website and select the appropriate operating system version (Windows, macOS or Linux). According to the system recommendations of your computer, choose the 64-bit version for better performance.

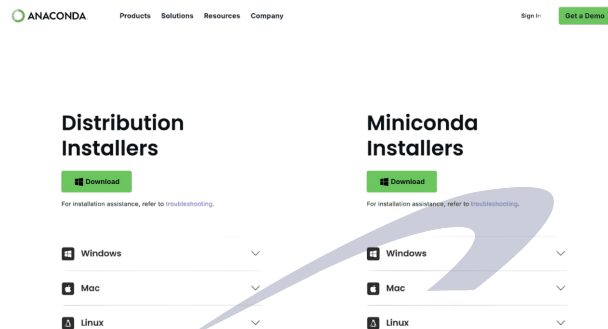


Figure 4.1 Download on Official Anaconda Website

Install Anaconda Double-click the downloaded Anaconda installation file (installer) to start the installation program. And click "Next" to proceed to the next step . Select the installation type. If it is for personal use only, it is recommended to select "Just Me", then click "Next".

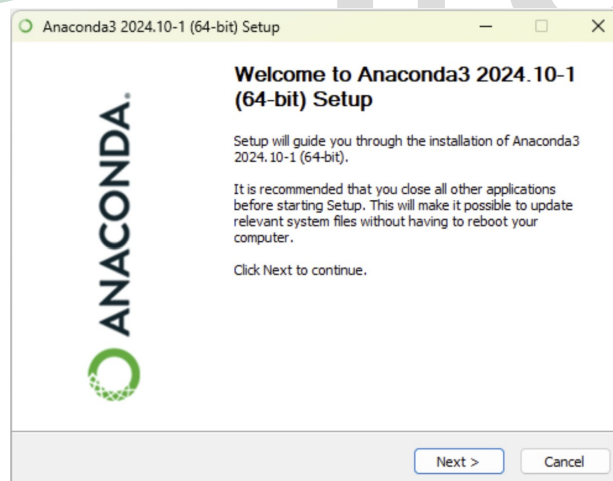


Figure 4.2 Installation for Anaconda

In the installation options, it is recommended not to check Add Anaconda to the PATH environment variable (unless there are special requirements), and directly click "Install" to start the installation.

¹<https://www.anaconda.com/products/distribution>

Once the installation is complete, find and launch Anaconda Navigator from the Windows Start menu .

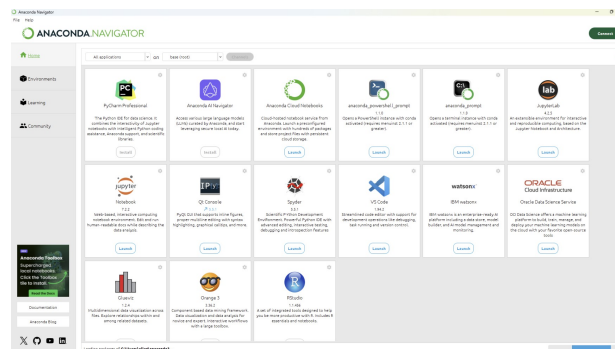


Figure 4.3 FlowChart for Preprocess Model

Step 2: Installing Visual Studio Code

Visual Studio Code (VS Code) is a lightweight and extensible source code editor that, when used with the Python Extension, offers enhanced development capabilities. The installation package can be obtained from the official website².

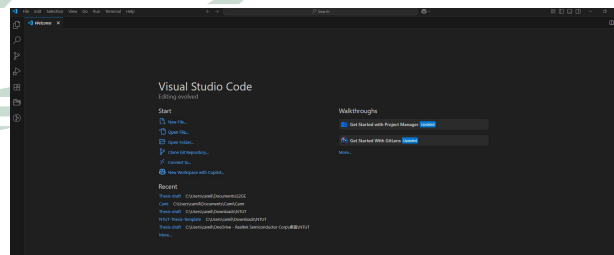


Figure 4.4 Visual Studio Code

4.1.2.2 Creating Virtual Environment and Installing Packages

Step 3: Creating a Python Virtual Environment

Use the Anaconda Prompt to create a virtual environment with the designated Python version:

```
conda create -n nids_env python=3.9  
conda activate nids_env
```

Step 4: Installing Required Packages

²<https://code.visualstudio.com/>

The packages required in this study are listed below and can be installed using pip:

```
pip install numpy pandas scikit-learn matplotlib seaborn torch mmh3
```

A brief description of each package is provided in Table 4.2.

Step 5: Selecting the VS Code Interpreter

In Visual Studio Code, press Ctrl+Shift+P to open the command palette, then select *"Python: Select Interpreter"*. Choose the previously created nids_env virtual environment from the list of available interpreters.

4.1.2.3 4.1.4 Verifying the Installation

To verify the installation, create a file named `main.py` and include the following test code:

```
import numpy as np
import pandas as pd
import torch
import mmh3
print("All packages loaded successfully!")
```

Execute the script in the terminal with the following command:

```
python main.py
```

If the message is displayed successfully, it indicates that the environment has been set up correctly.

Chapter 5 Conclusion & Future Work

5.1 Conclusion

5.2 Future Work

