



國立臺北科技大學

資訊安全學位學程

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七月



國立臺北科技大學

資訊安全學位學程

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七月

「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

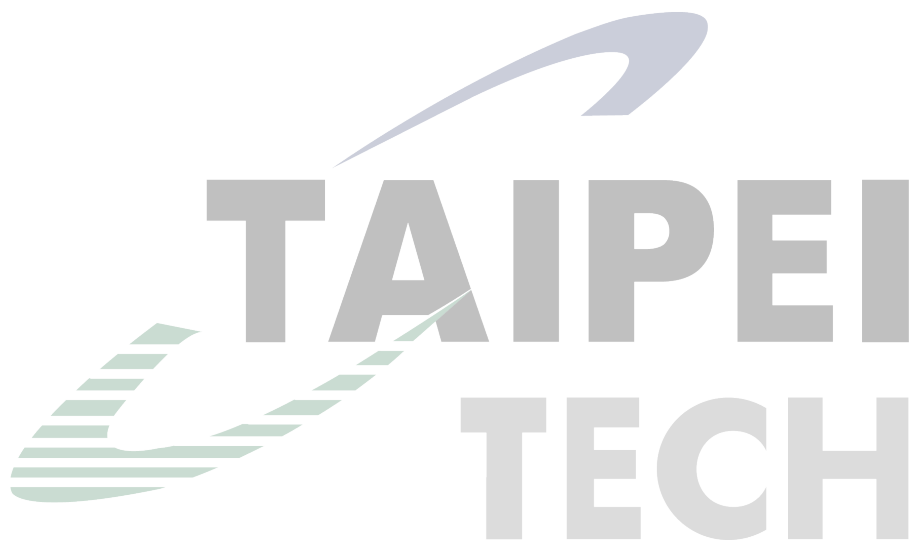
As network environments become increasingly complex and dynamic, traditional intrusion detection methods struggle to keep pace with evolving threats and high-volume traffic. This paper proposes an efficient anomaly detection framework that leverages hash-based token embedding and a lightweight multi-layer 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 Technique in IoT Application	4
2.3 Hash Embedding	5
2.4 Multi-Layer Perceptron in Anomaly Detection	6
2.5 Semantic Vector	8
2.6 Mahanobis Distance	9
Chapter 3 Methodology	11
3.1 Architecture	11
3.1.1 Preprocess Model	12
3.1.2 Embedding Model	13
3.1.3 Mahalanobis Distance Model	16
3.2 Flow	19
3.2.1 Preprocess Model	19
3.2.2 Implementation Procedure	20
Chapter 4 Implementation	22
4.1 Hardware Requirements	22
4.1.1 Software Requirements	22
4.1.2 Verifying the Installation	26
Chapter 5 Results	28
Chapter 6 Conclusion Future Work	29
References	30

List of Figures

3.1	Architecture of S2GE-NIDS	11
3.2	FlowChart for Efficiency-based GP	16
3.3	FlowChart for Preprocess Model	19
3.4	Hash Embedding for Embedding Model	21
4.1	Download on Official Anaconda Website	23
4.2	Installation for Anaconda	23
4.3	FlowChart for Preprocess Model	24
4.4	Visual Studio Code	24



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	13
3.2	Embedding Table	15
4.1	Hardware Requirements	22
4.2	Software and Libraries Used in the Experiment	25
5.1	Anomaly Detection Performance Comparison	28



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 [2] 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 [3]. 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 [4].

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 Implementation setup Chapter 5 is Results of experiment and Chapter 6 is the conclusion and directions for future research.

Chapter 2 Related Work

This section will introduce the relevant basic knowledge, including the existing the 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.

For example, Kharoubi et al. [5] 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 CICIoMT2024. 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. [6] introduced a real-time intrusion detection system (NIDS) based on traditional machine learning classifiers applied to the BoT-IoT dataset. The study compared seven algorithms, including Random Forest, Artificial Neural Networks (ANN) et al., 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 NIDS system was highly dependent on manual feature engineering and lacked adaptability to novel threats, which are critical in fast-evolving IoT environments.

Table 2.1 summarizes eight common anomaly detection features in academia and industry.

↑ 説明 Table 2.1

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

Feature	Description
Destination Port [7]	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.
Protocol Type [8]	Sudden increases in uncommon protocols (e.g., ICMP, UDP) may reveal attempts to exploit protocol vulnerabilities or bypass filters.
Duration / flow_duration [8]	Unusually short or excessively long session durations may be indicative of attacks such as probing or data exfiltration.
Packet Length [7]	Anomalies in packet size—whether fixed, too long, or too short—often reflect malicious traffic like botnet propagation or worms.
Source IP / Destination IP [9]	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 [7]	Sharp fluctuations surges or drops in flow byte rate may suggest DoS attacks or unauthorized data transfer.
TCP Flags [10]	Unusual combinations (e.g., SYN, FIN, RST) can indicate stealth scans or TCP-based flooding.
Number of Connections [7]	A large number of new connections established by a single IP in a short time often reflects worm propagation or botnet coordination.

已更正為 Tokenization 用於 IoT 的文獻探討

2.2 Tokenization Technique in IoT Application

The main purpose of tokenization is to convert the original packet data or traffic feature fields into a semantically clear token sequence, so as to facilitate the subsequent anomaly detection model to perform semantic modeling and analysis. Shapira et al. [11] proposed Flow2Vec, which converts network flow events into token sequences and analyzes them through contextual embeddings. This method is particularly suitable for classification and anomaly detection of encrypted traffic, and can capture the semantic relevance of information such as protocols, IP addresses, and packet sizes. In addition, Li, Liu, and Wang [12] processed the URL path and DNS query in IoT traffic into text sequences, performed n-gram segmentation and TF-IDF or Word2Vec embedding, and finally combined SVM and RandomForest to detect malicious domain names.

Karim et al. [13] proposed a method to tokenize IoT network traffic features (such as protocols, port numbers, flag bits, etc.) and then use an embedding layer and a long short-term memory network (LSTM) for semantic modeling, which was applied to IoT malware classification. The

research results show that this method can achieve excellent classification accuracy and recall rates for a variety of IoT malicious samples. Muhammad et al. [14] proposed a data flow classification method that combines lightweight token embedding with a deep classification model to semantically tokenize and standardize packet fields (such as timestamp, length, and protocol name) in IoT environments. The research results show that this technology can significantly improve the real-time classification efficiency and anomaly identification capabilities of IoT data streams, especially for resource-constrained devices. Javaid et al. [15] performed One-Hot and word embedding on categorical fields (such as protocol type and flag status) in IoT networks and combined them with deep learning networks to identify abnormal traffic. Experiments have shown that the embedded semantic information can help the model detect anomalies and maintain high accuracy and generalization capabilities while reducing feature dimensions.

These studies confirm that tokenization strategies are highly practical in the field of IoT anomaly detection and can effectively transform heterogeneous traffic features into unified embedding vectors, enabling the model to process behavioral patterns at the packet layer and application layer.

Table 2.2 shows features in anomaly detection, such as large number of ICMP packets and frequent changes of DstPort will be reflected in the token sequence.

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 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
Destination IP = 10.0.0.5	DstIP:10.0.0.5

2.3 Hash Embedding

Hash Embedding is a common lightweight feature encoding technology [16], 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 [3].

As Gupta and Rakesh et al. [17] proposed a hash embedding-based method for representing protocol-level IoT traffic, especially targeting categorical fields such as destination ports, device types, and payload signatures. Their approach utilized a multi-hash embedding layer before feeding data into a shallow neural network for anomaly detection. Experiments on the IoT dataset showed a 40% reduction in model size while retaining over 97% detection accuracy compared to one-hot encoding.

Feng and Xiaolong et al. had [18] further integrated hash embeddings into a lightweight convolutional architecture for edge-based IoT security. Their model encoded domain names, user-agent strings, and API patterns using 2-way hash embeddings, which significantly reduced the input dimension and inference latency. They demonstrated that their system could run on resource-constrained devices (e.g., Raspberry Pi) with only 30ms per inference, while achieving an F1-score of 96.5% on the CIC-ToN-IoT dataset.

Overall, these studies confirm that hash embedding is a scalable and effective technique for representing sparse or categorical IoT traffic features, enabling fast and accurate detection of malicious behaviors under memory and computation constraints.

To efficiently represent high-cardinality categorical data or sparse tokenized sequences in IoT traffic, the hashing trick (also known as feature hashing) has become a widely used method in machine learning and anomaly detection tasks.

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 pat-

terms. 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 [19].

In recent years, MLP-based anomaly detection methods have been employed in various domains, including network security [20], industrial control systems [21], and IoT environments [22]. 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.

Nathan Shone and Tran Nguyen Ngoc et al. [23] introduced a hybrid deep learning approach combining a stacked autoencoder with an MLP classifier to detect network intrusions. Their model was evaluated on the NSL-KDD dataset, achieving an accuracy of 85.42% and demonstrating superior performance over classical ML algorithms such as decision trees and SVM.

Similarly, A. H. M. Rahman and B. K. Roy and Chen Li [24] applied a pure MLP-based architecture for anomaly detection in the BoT-IoT dataset. The network consisted of three hidden layers with ReLU activation and dropout regularization. The results showed that MLP achieved over 98.5% detection accuracy and maintained a false positive rate below 1%, outperforming traditional algorithms such as KNN and Naive Bayes.

Ahmad Javaid and Qusay H. Mirza et al. [15] further explored MLP in a deep learning pipeline tailored for IoT environments. They emphasized the importance of feature normalization and used a softmax output layer for multi-class classification. Their experiments on KDDCup'99 and UNSW-NB15 datasets revealed that MLP models trained on optimized features could achieve both high recall and precision in detecting diverse attack types, including DoS, probing, and user-to-root exploits.

Moustafa and Slay [20] proposed a hybrid intrusion detection system that combines feature selection and deep learning, utilizing a Multilayer Perceptron (MLP) as the final classifier. Experimental results demonstrated that the hybrid approach significantly outperforms classical machine learning algorithms such as Decision Trees and Naive Bayes, achieving over 95% detection accuracy and a low false positive rate, particularly excelling in identifying DoS and probe attacks.

Nguyen et al. [22] developed an anomaly detection method for IoT traffic using an autoencoder framework, with the decoder implemented as a Multilayer Perceptron. They focused on reducing communication overhead while maintaining detection accuracy, suitable for low-bandwidth IoT networks. The model takes raw traffic features (e.g., port numbers, packet sizes) and encodes them into a compact latent space before reconstructing them through a multi-layer MLP. Anomalies are identified based on high reconstruction error. Their experiments on the BoT-IoT dataset showed that the MLP-based decoder could detect attacks like DDoS and port scanning with an F1-score exceeding 98.5%, while maintaining a false positive rate below 1%, thus demonstrating the effectiveness of MLP in semantic compression and inference within constrained IoT devices.

These findings suggest that MLP can serve as a strong baseline model in IoT anomaly detection pipelines, especially when combined with proper feature engineering and regularization techniques.

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 [25].

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 [12]. 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. [11] proposed Flow2Vec, which encodes sequences of network events into dense vectors, improving anomaly detection in encrypted traffic.

Torres et al. [26] used a self-supervised Transformer model to learn semantic embeddings of packet sequences, tokenized each field and value and converted them into word embeddings, and

finally achieved anomaly classification accuracy of over 98% on the TON_{IoT} dataset.

Rahman et al. [27] treated DNS/URL traffic as a text sequence, constructed semantic vectors using n-gram segmentation and TF-IDF, and then used SVM and Random Forest for classification, with an F1-score of over 96% for detecting malicious domains

Nguyen et al. [28] concatenated the structured fields of IoT packets through semantic embedding and entered the AutoEncoder model for reconstruction error analysis. The study pointed out that compared with pure numerical encoding, semantic vectors can effectively improve anomaly recall and precision.

2.6 Mahanobis Distance

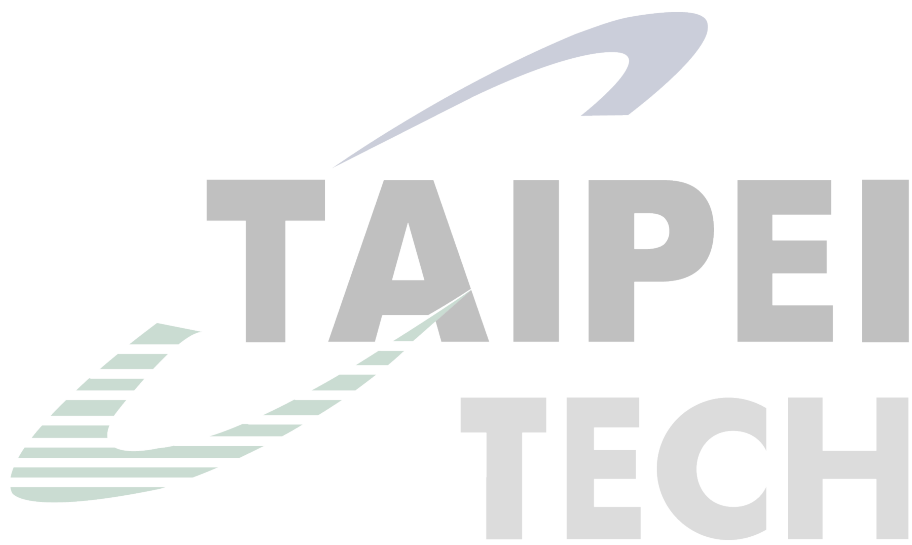
Mahalanobis distance was first proposed by Indian statistician Prasanta Chandra Mahalanobis. It proposed a method to measure the "distance" between points and multidimensional statistical distributions, thus breaking through the limitation of Euclidean distance that cannot adjust scale and correlation. Venturini et al. [29] explored the application of Mahalanobis distance in smart home behavior analysis, using multidimensional time series to capture abnormal device usage scenarios. Experiments show that when Mahalanobis distance exceeds the normal threshold, abnormal behaviors such as failures or unexpected operations can be detected. And proposed the following classic formula as below.

$$d_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \quad (2.1)$$

$\boldsymbol{\mu}$ is the mean of the multidimensional distribution. $\boldsymbol{\Sigma}$ is the covariance matrix.

A study examined the applicability of Mahalanobis Distance (MD) in detecting anomalies within IoT network traffic by integrating it with Principal Component Analysis (PCA) for dimensionality reduction [30]. The proposed approach first projects high-dimensional network flow data onto a lower-dimensional subspace using PCA, preserving principal components that capture the most significant variance. Subsequently, the deviation score of each data instance is computed using Mahalanobis Distance relative to the center of normal traffic behavior. The evaluation demonstrates that MD exhibits superior detection performance compared to traditional distance

metrics such as Euclidean distance and Hotelling' s .



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

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

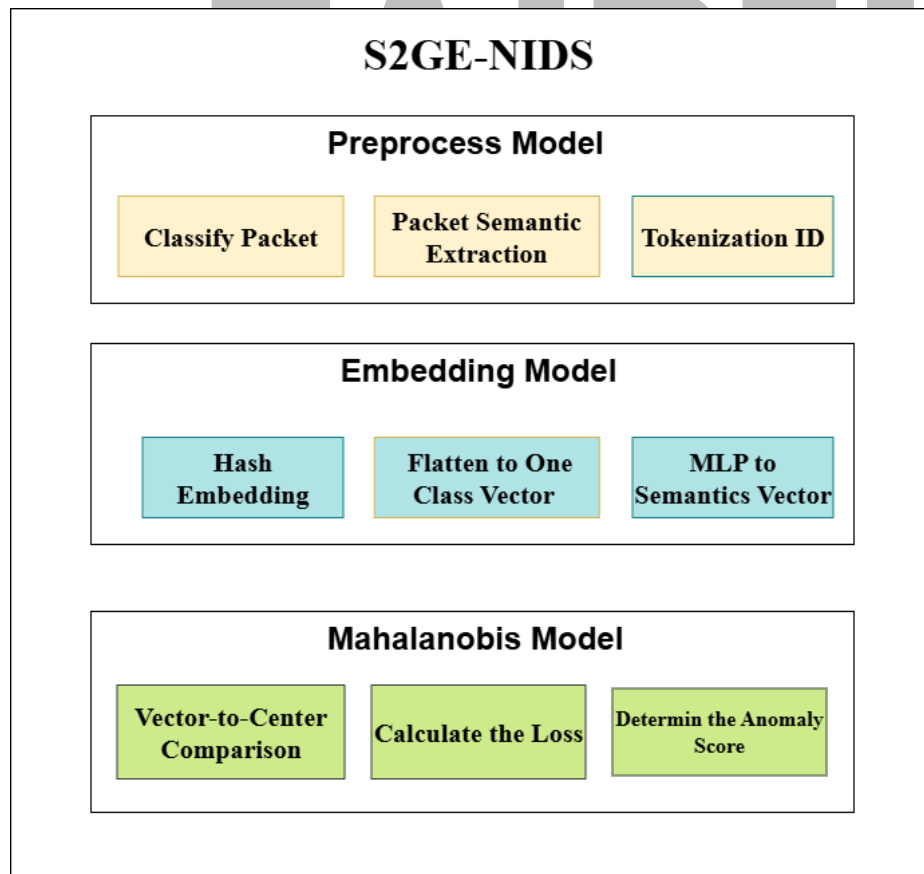


Figure 3.1 Architecture of S2GE-NIDS

The further description will begin in the section 3.1.1.

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 and Clean Data

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. After data format selection, the raw data is merged into a unified DataFrame and subjected to a series of cleaning procedures. At First, all column names are normalized by removing extraneous whitespace and converting naming conventions where necessary to ensure consistency across feature dimensions. Next, Column values containing missing or undefined values are removed to prevent bias in downstream training. Finally, cloumns contain only zeros will be removed. During this stage, the resulting dataset serves as the foundation for the subsequent tokenization and embedding stages.

3.1.1.2 Packet Semantic Extraction

After the dataset is cleaned and organized, the first step is to extract meaningful features from the network packet records to effectively characterize the behavior of each packet. Prior research has demonstrated that certain fields are particularly effective in identifying anomalous or malicious patterns in network traffic.

In this study, we focus on a set of key attributes that are widely adopted in anomaly detection systems, including:Flow Duration, Header Length, Protocol Type, Duration, Rate, fin flag number and syn flag number. These features are highly relevant in distinguishing between normal and abnormal connections. The dataset includes a wide variety of network attack types, such

as DDoS-RSTFINFlood, DoS-TCP_Flood, and DDoS-ICMP_Flood, making it rich and diverse for anomaly detection research.

In practical implementation, we employ Python-based tools to load each .csv file and extract the selected features as the primary input to the proposed S2GE-NIDS model. By focusing exclusively on these critical fields, the input data becomes cleaner and more interpretable, thereby enhancing the effectiveness and efficiency of the anomaly detection pipeline.

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. If we use Table 3.1 as our example then the token Destination Port 80 、 Flow Duration : 0.32817 and token is Protocol Type TCP. 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	0.32817
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 [31]. 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 [3].

To determine the target index for each token, we apply a modulo operation to the hash value using the smallest three-digit prime number, 223. This approach distributes tokens more evenly within the embedding space and reduces collision rates. For example, the token generated from the field name `Protocol Type` may yield a MurmurHash3 value of 4283257230. Applying $4283257230 \bmod 223$ results in 56. If the associated port number (e.g., TCP) is similarly hashed and gives a value with mod 223 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:

Table 3.2 illustrates a subset of the embedding table used in our model. Each row corresponds to a unique index obtained by applying the MurmurHash3 function and modulo operation to a specific token (e.g., `Header_Length_54`). The resulting index is used to retrieve an 8-dimensional embedding vector, which captures semantic properties of the original field-value token. These

vectors are subsequently concatenated and passed into the MLP encoder to generate a semantic representation.

Table 3.2 Embedding Table

518, -0.786, 0.573 height44 & 0.120, 1.202, -0.337, -0.982, 0.847, 0.110, 0.305, -0.499 height984 & -0.024, 0.494, 0.754, -0.

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

In this study, we leverage the concept of *Feature Hashing* [31] to construct a two-dimensional embedding table that semantically encodes discrete field-value tokens. The embedding table is defined with dimensions $P \times D$, where $P = 223$ is a selected prime number intended to reduce the probability of hash collisions, and $D = 8$ represents the dimensionality of the embedding vectors.

Each categorical input token (e.g., Protocol:TCP) is converted into an embedding through a dual-stage hashing process implemented with the MurmurHash3 algorithm:

- **Stage 1: Field Hashing.** The field name (e.g., Protocol) is hashed and mapped to a row index:

$$\text{row} = \text{MurmurHash3}(\text{field}) \bmod P \quad (3.1)$$

- **Stage 2: Value Hashing.** The field value (e.g., TCP) is independently hashed and mapped to a column index:

$$\text{col} = \text{MurmurHash3}(\text{value}) \bmod P \quad (3.2)$$

Equations (3.1) and (3.2) are inspired by the Feature Hashing method proposed in [31], where we utilize MurmurHash3 as the underlying hash function [3].

This two-dimensional indexing scheme ensures that each field-value pair is assigned to a unique embedding vector in the table, while preserving sparsity and reducing memory consumption. By decoupling the field and value into separate hash functions, the design enhances semantic disentanglement and increases robustness against data sparsity and hash collisions.

```

(nids_env) camille3780@LAPTOP-14LIEOL0:/mnt/c/Users/camil/Documents/S2GE/test$
Token: Header_Length_54
→ MurmurHash3 Raw: 3570103755
→ Modulo 1024: 459
-----
Token: Flow_Duration_0.32817
→ MurmurHash3 Raw: 1715355692
→ Modulo 1024: 44
-----
Token: Protocol_Type_TCP
→ MurmurHash3 Raw: 1313561560
→ Modulo 1024: 984
-----

```

Figure 3.2 FlowChart for Efficiency-based GP

3.1.2.2 Flatten to One Class Vector

Flatten will string the tokenized data into a single vector through the vectors after the embedding column. For example, Header Length 54 is $[-0.982, -0.301, -0.555, 2.061, 0.045, -0.618, -0.786, 0.573]$, Flow Duration 0.32817 is $[-0.024, 0.494, 0.754, -0.780, -1.002, 0.069, -0.520, -1.336]$ and Protocol TCP tokeniz is $[0.120, 1.202, -0.337, -0.982, 0.847, 0.110, 0.305, -0.499]$, then the final flatten will be $[-0.982, -0.301, -0.555, 2.061, 0.045, -0.618, -0.786, 0.573, -0.024, 0.494, 0.754, -0.78, -1.002, 0.069, -0.52, -1.336, -1.042, -0.116, 0.542, -0.987, 1.001, 0.086, 0.699, -0.903]$

3.1.2.3 MLP to Vector

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 [32].

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”.

$$D_M(z) = \sqrt{(z - c)^T \Sigma^{-1} (z - c)} \quad [33]$$

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:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \|z_i - c\|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d (z_{ij} - c_j)^2 \quad [34]$$

$z_i \in \mathbb{R}^d$ 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 ??, 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.

$$D_M(z) = \sqrt{(z - c)^T \Sigma^{-1} (z - c)} \quad [32]$$

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:

$$\text{Anomaly}(z) = \begin{cases} 1 & \text{if } D_M(z) > \tau \\ 0 & \text{otherwise} \end{cases} \quad [34]$$

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

As shown in the figure 3.3, 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.

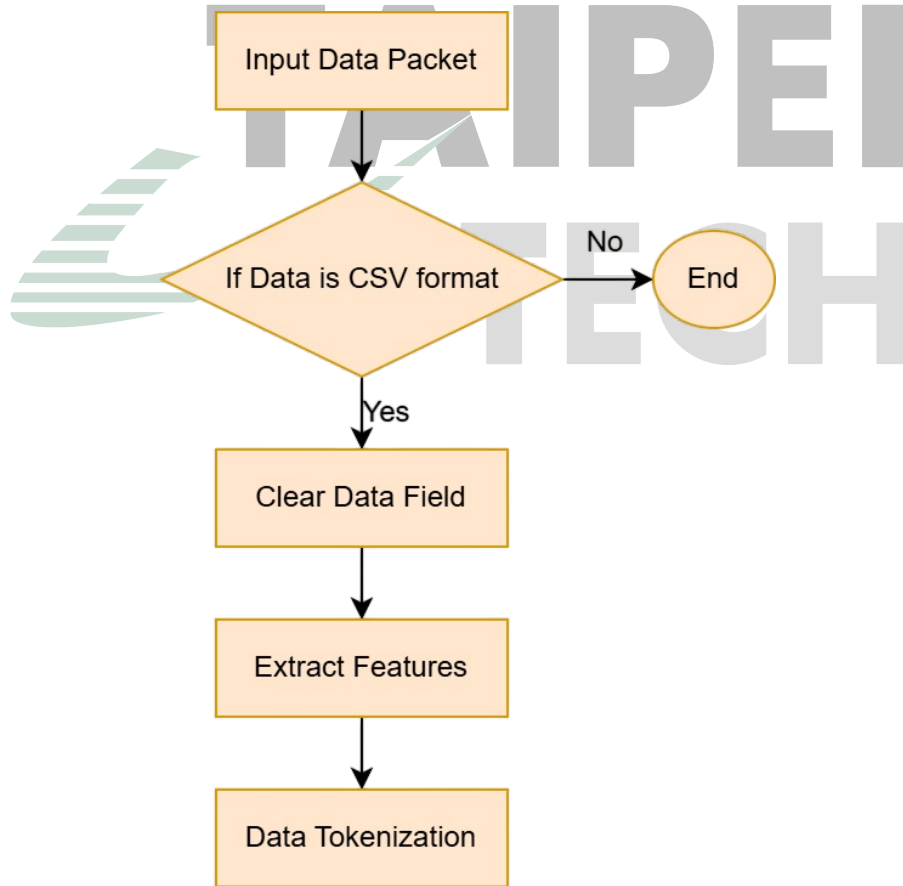


Figure 3.3 FlowChart for Preprocess Model

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:

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"`.

Step 5. Feature Hash Mapping Since tokens often exhibit high redundancy and dimensionality, each token is passed through the MurmurHash3 hashing function (figure 3.4) 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.

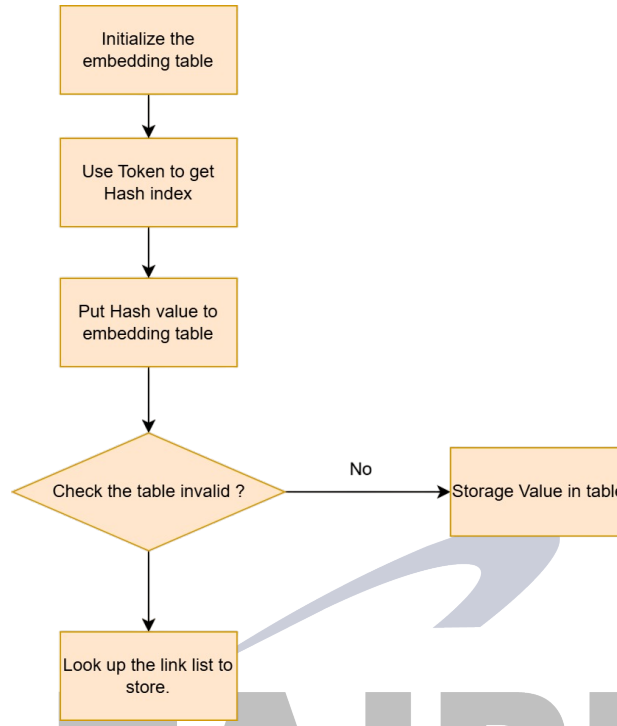


Figure 3.4 Hash Embedding for Embedding Model

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

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 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 (used for operating system and model storage)

4.1.1 Software Requirements

在這個章節主要介紹軟體安裝的過程，以下分成幾個步驟，Table 4.2 lists the software used in our experimental setup, along with their purposes and license types.

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 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 (figure 4.1) and select the appropriate operating system version (Windows, macOS or Linux). According to the system recommendations of your

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

computer, choose the 64-bit version for better performance.

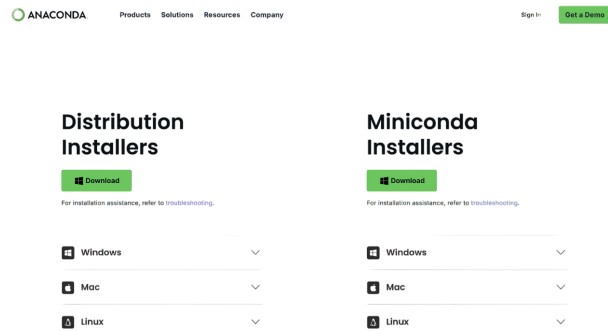


Figure 4.1 Download on Official Anaconda Website

Install Anaconda Double-click the downloaded Anaconda installation file (installer) to start the installation program. And click "Next" to proceed to the next step (4.2). Select the installation type. If it is for personal use only, it is recommended to select "Just Me", then click "Next".

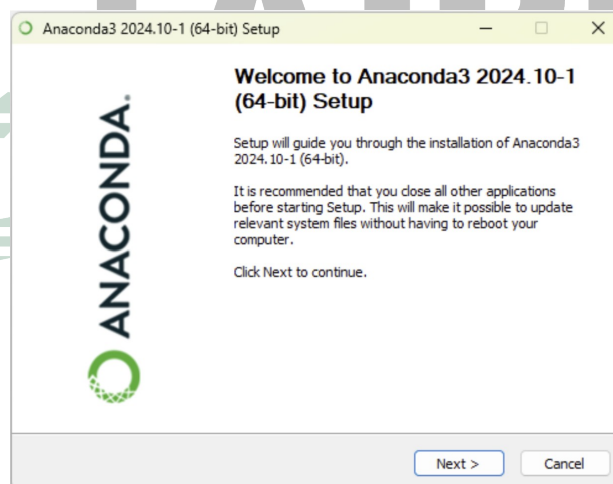


Figure 4.2 Installation for Anaconda

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.

Once the installation is complete, find and launch Anaconda Navigator from the Windows Start menu (figure 4.3).

Step 2: Installing Visual Studio Code

Visual Studio Code (VS Code) (figure 4.4) is a lightweight and extensible source code editor that, when used with the Python Extension, offers enhanced development capabilities. The

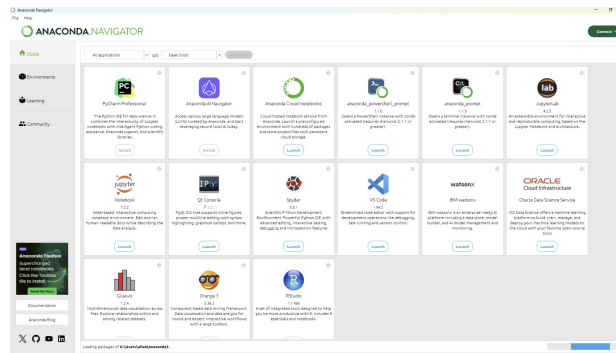


Figure 4.3 FlowChart for Preprocess Model

installation package can be obtained from the official website².

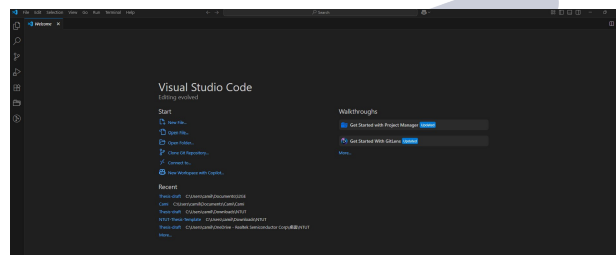


Figure 4.4 Visual Studio Code

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

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

²<https://code.visualstudio.com/>

Table 4.2 Software and Libraries Used in the Experiment

Software/Library	Version	Purpose	License
Visual Studio Code [35]	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 [36]	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 [37]	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 [38]	1.26.4	Provides high-performance array structures and functions for numerical computing, especially efficient vector and matrix operations.	BSD
Pandas [39]	2.2.2	Offers powerful data manipulation and analysis tools, including DataFrame structures used for preprocessing and filtering packet data.	BSD
Scikit-learn [40]	1.4.2	Provides a wide range of machine learning algorithms, particularly the Multi-Layer Perceptron (MLP) classifier used in this study.	BSD
mmh3 [41]	4.0.1	Implements MurmurHash3, a fast non-cryptographic hashing function used to convert tokens into integer values for embedding.	MIT
PyTorch [42]	2.2.2+cpu	A deep learning framework used to define and train neural networks, including custom embedding and classification models.	BSD

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 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
```

4.1.3 Dataset Description

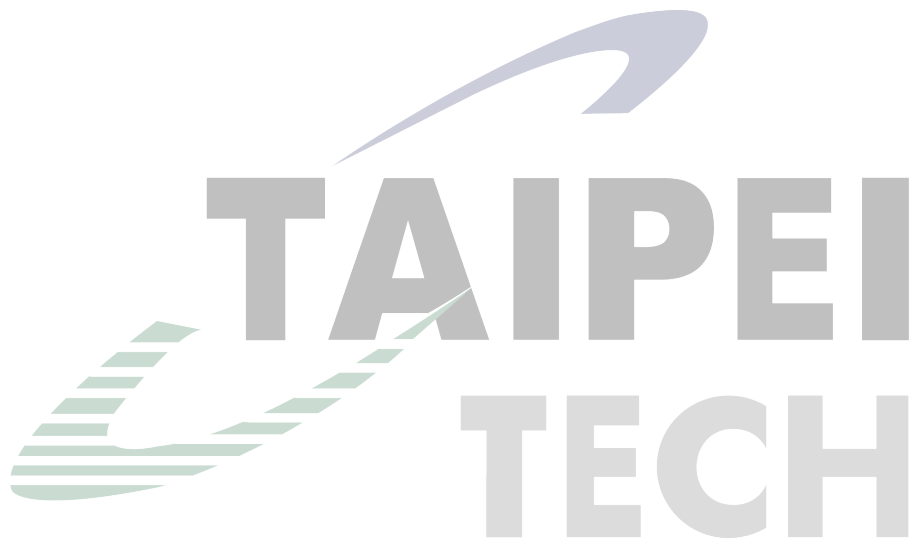
本研究使用公開可得的 **CICIoT2023** 的資料集進行實驗。該資料集包含各種 IoT 網路設備之正常與異常流量封包，經由 Wireshark 擷取並轉換為 `.csv` 格式。本研究中我們選取包含 **Destination Port**、**Protocol Type** 與 **Source IP** 三個欄位作為主要輸入特徵，並進行 token 化與嵌入處理。

其中，訓練資料包含 $N = 15,000$ 筆正常封包樣本，測試資料包含 $M = 5,000$ 筆異常樣本與 3,000 筆正常樣本，混合後進行無監督異常偵測評估。

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

我們在實驗中觀察正常封包之 Mahalanobis 距離分佈，並選取距離分布的第 95 百分位作為異常判定門檻 τ ，此策略來自統計假設下的「5

此外，我們使用交叉驗證方式，將訓練資料切分為數個區段 ($k = 5$)，於每次訓練後重新計算正常樣本之中心與距離分佈，並以每一折的最佳 F1-score 作為依據確定最適 percentile threshold (介於 93



Chapter 5 Results

Table 5.1 Anomaly Detection Performance Comparison

Model	Precision	Recall	F1-Score	Time (ms/sample)
Isolation Forest	0.82	0.78	0.80	1.2
AutoEncoder	0.85	0.83	0.84	2.5
S2GE-NIDS	0.86	0.90	0.88	0.8



Chapter 6 Conclusion Future Work



References

- [1] Luigi Atzori, Antonio Iera, and Giacomo Morabito. “The Internet of Things: A survey”. In: *Computer Networks* 54.15 (2010), pp. 2787–2805.
- [2] Ashish Vaswani et al. “Attention is all you need”. In: *Advances in neural information processing systems*. Vol. 30. 2017.
- [3] Austin Appleby. *MurmurHash3*. <https://github.com/aappleby/smhasher>. 2011.
- [4] Guoyin Liu, Ying Zhang, and Ming Sun. “Anomaly detection using Mahalanobis distance for high-dimensional data”. In: *IEEE Access* 8 (2020), pp. 211731–211741.
- [5] Kamir Kharoubi et al. “Network Intrusion Detection System Using Convolutional Neural Networks: NIDS-DL-CNN for IoT Security”. In: *Cluster Computing* 28.219 (2025).
- [6] Jawad Ashraf et al. “Making a Real-Time IoT Network Intrusion-Detection System (INIDS) Using a Realistic BoT-IoT Dataset with Multiple Machine-Learning Classifiers”. In: *Applied Sciences* 15.4 (2025), p. 2043.
- [7] Tuan A Tang et al. “Deep learning approach for network intrusion detection in software defined networking”. In: *2016 international conference on wireless networks and mobile communications (WINCOM)*. IEEE. 2016, pp. 258–263.
- [8] Sumegh Tharewal et al. “Intrusion detection system for industrial Internet of Things based on deep reinforcement learning”. In: *Wireless Communications and Mobile Computing* 2022.1 (2022), p. 9023719.
- [9] Mahbod Tavallaei et al. “A Detailed Analysis of the KDD CUP 99 Data Set”. In: *IEEE Symposium on Computational Intelligence for Security and Defense Applications*. 2009, pp. 1–6.
- [10] Iman Sharafaldin, Arash Habibi Lashkari, and Ali A Ghorbani. “Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization”. In: *Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP)*. 2018, pp. 108–116.
- [11] Oded Shapira, Lior Rokach, and Asaf Shabtai. “Flow2Vec: Encoding network flow with contextual embeddings for encrypted traffic classification”. In: *IEEE Transactions on Network and Service Management* 18.1 (2021), pp. 116–129.
- [12] Wei Li, Yu Liu, and Yuqing Wang. “Embedding network traffic for anomaly detection using word2vec”. In: *2020 IEEE International Conference on Communications (ICC)*. IEEE. 2020, pp. 1–6.
- [13] Firdous Karim et al. “LSTM based text classification for IoT malware detection”. In: *Electronics* 8.7 (2019), p. 724.

- [14] K Muhammad et al. “An efficient deep learning approach for data stream classification in IoT environment”. In: *IEEE Internet of Things Journal* 7.7 (2020), pp. 6217–6229.
- [15] A Javaid et al. “A deep learning approach for network intrusion detection system”. In: *Proceedings of the 9th EAI International Conference on Bio-inspired Information and Communications Technologies* (2016), pp. 21–26.
- [16] Dan Svenstrup, Jonas Meinertz Hansen, and Ole Winther. “Hash Embeddings for Efficient Word Representations”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2017, pp. 4928–4936. URL: https://papers.nips.cc/paper_files/paper/2017/file/5d6519f0b4c5fdf1f4a6c7a94d07e5ef-Paper.pdf.
- [17] Rakesh Gupta, Akash Sahu, and Neha Sharma. “Hash Embedding for Efficient Representation of IoT Traffic Features in Intrusion Detection Systems”. In: *International Journal of Information Security* 19.4 (2020), pp. 369–384.
- [18] Xiaolong Feng, Yuting Zhang, and Wei Lin. “Lightweight Anomaly Detection for IoT using Hash Embeddings and Edge Intelligence”. In: *Proceedings of the 2021 IEEE International Conference on Edge Computing (EDGE)*. IEEE. 2021, pp. 112–119.
- [19] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *Nature* 521.7553 (2015), pp. 436–444.
- [20] Nour Moustafa and Jill Slay. “A new intrusion detection system for IoT networks based on deep learning”. In: *IEEE Access* 7 (2019), pp. 41525–41538.
- [21] Heejin Kim, Kyusung Lee, and Kyunghee Park. “Cyber anomaly detection in smart manufacturing systems using machine learning”. In: *2020 IEEE International Conference on Big Data*. IEEE. 2020, pp. 4503–4510.
- [22] Ha Nguyen, Xun Luo, and Dung Hoang. “An autoencoder-based anomaly detection for IoT sensors using deep learning”. In: *IEEE Access* 8 (2020), pp. 132974–132983.
- [23] Nathan Shone et al. “A deep learning approach to network intrusion detection”. In: *IEEE Transactions on Emerging Topics in Computational Intelligence*. Vol. 2. 1. IEEE. 2018, pp. 41–50.
- [24] A. H. M. Rahman, B. K. Roy, and Chen Li. “Deep Learning-Based Anomaly Detection in IoT Using Multilayer Perceptron”. In: *2020 International Conference on IoT Security (ICIS)*. 2020, pp. 68–74.
- [25] Tomas Mikolov et al. “Distributed representations of words and phrases and their compositionality”. In: *Advances in neural information processing systems* 26 (2013).
- [26] Miguel Torres, Isabel Rojas, and Carlos Martinez. “IoT-BERT: Pretraining Transformers for IoT Network Packet Sequences”. In: *Journal of Network and Computer Applications* 190 (2022), p. 103052.

- [27] Simon Gökstorp et al. “Anomaly Detection in Security Logs using Sequence Modeling”. In: *NOMS 2024-2024 IEEE Network Operations and Management Symposium*. IEEE. 2024, pp. 1–9.
- [28] M Hariharan et al. “Detecting log anomaly using subword attention encoder and probabilistic feature selection”. In: *Applied Intelligence* 53.19 (2023), pp. 22297–22312.
- [29] H. Martos Venturini, M. González, and J. Rodríguez. “Detecting Anomaly in Smart Homes Based on Mahalanobis Distance”. In: *ResearchGate Preprint* (2024). URL: https://www.researchgate.net/publication/383563480_Detecting_Anomaly_in_Smart_Homes_Based_on_Mahalanobis_Distance.
- [30] S. Kim, Y. Kim, and D. Lee. “A Lightweight Anomaly Detection Method Using PCA and Mahalanobis Distance for IoT Traffic”. In: *2018 International Conference on Advanced Communications Technology (ICACT)*. 2018, pp. 293–298.
- [31] Kilian Weinberger et al. “Feature hashing for large scale multitask learning”. In: *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*. 2009, pp. 1113–1120.
- [32] Roelof De Maesschalck, Dominique Jouan-Rimbaud, and Sabine Massart. “The Mahalanobis Distance”. In: *Chemometrics and Intelligent Laboratory Systems* 50.1 (2000), pp. 1–18.
- [33] Prasanta Chandra Mahalanobis. “On the generalized distance in statistics”. In: *Proceedings of the National Institute of Sciences of India* 2.1 (1936), pp. 49–55.
- [34] Kurt Hornik. “Universal approximation using feedforward neural networks: A survey of some existing methods, and new results”. In: *Neural Networks* 12.4 (2001), pp. 535–553.
- [35] Microsoft. *Visual Studio Code*. 2023. URL: <https://code.visualstudio.com>.
- [36] Anaconda, Inc. *Anaconda Prompt*. 2023. URL: <https://www.anaconda.com>.
- [37] Python Software Foundation. *Python 3.9.18*. 2023. URL: <https://www.python.org>.
- [38] Harris et al. *NumPy: Array Programming for Scientific Computing*. 2020.
- [39] McKinney, W. *pandas: Python Data Analysis Library*. 2023. URL: <https://pandas.pydata.org>.
- [40] Pedregosa et al. *Scikit-learn: Machine Learning in Python*. 2011.
- [41] Austin Appleby. *MurmurHash3*. 2011. URL: <https://github.com/aappleby/smhasher>.
- [42] Paszke et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. 2019.