

IoT Device Labeling Using Large Language Models

Bar Meyuhas
Reichman University

Anat Bremler-Barr
Tel-Aviv University

Tal Shapira
Tel-Aviv University

Abstract

The IoT market is diverse and characterized by a multitude of vendors that support different device functions (e.g., speaker, camera, vacuum cleaner, etc.). Within this market, IoT security and observability systems use real-time identification techniques to manage these devices effectively. Most existing IoT identification solutions employ machine learning techniques that assume the IoT device, labeled by both its vendor and function, was observed during their training phase. We tackle a key challenge in IoT labeling: how can an AI solution label an IoT device that has never been seen before and whose label is unknown?

Our solution extracts textual features such as domain names and hostnames from network traffic, and then enriches these features using Google search data alongside catalog of vendors and device functions. The solution also integrates an auto-update mechanism that uses Large Language Models (LLMs) to update these catalogs with emerging device types. Based on the information gathered, the device’s vendor is identified through string matching with the enriched features. The function is then deduced by LLMs and zero-shot classification from a predefined catalog of IoT functions.

In an evaluation of our solution on 97 unique IoT devices, our function labeling approach achieved HIT1 and HIT2 scores of 0.7 and 0.77, respectively. As far as we know, this is the first research to tackle AI-automated IoT labeling.

1 Introduction

The rapid proliferation of Internet of Things (IoT) technology in various sectors has created new security and management challenges. Correctly identifying an IoT device plays a crucial role in IoT security and observability solutions [16, 25, 31]. Extensive research in both industry and academia has been dedicated to IoT identification algorithms [2, 18, 19, 23, 27], which are crucial for many IoT security and observability solutions. These algorithms predominantly rely on labeled IoT datasets for the real-time identification of already **seen**

devices. However, given the sheer diversity of IoT devices, collecting labeled data is a formidable challenge.

This paper introduces a method that passively monitors network traffic to automatically label unseen IoT devices by identifying their vendor and function. To the best of our knowledge, this is the first paper to address the IoT labeling problem. The motivation to address this issue was born out of discussions with IoT security companies who are having difficulty identifying IoT devices in the wild. In the IoT realm, maintaining a lab with all known IoT devices simply isn’t feasible. Moreover, the increasing number of cybersecurity attacks and new regulations make IoT vulnerability management important to various organizations nowadays. IoT security considers a vulnerability to be potential when it applies to a specific device type and vendor/model. Therefore, the first stage in any vulnerability management program is to identify and construct an inventory of assets, particularly IoT devices owned by the organization [34].

Labeling IoT devices with their vendor and function provides crucial visibility, ensuring administrators can effectively oversee their network. Our solution also introduces the potential for tailored security policies that can be formulated based on precise device functions, such as ones that give a device labeled ‘smart doorbell’ limited external communication but grant one labeled ‘smart TV’ broader access.

To the best of our knowledge, the high-tech industry currently relies on a single solution for IoT device labeling, with Fing [15] standing out as the leading product used by IT teams to label devices in their networks. Our solution has surpassed Fing’s performance in achieving accurate results. Their approach to the labeling challenge is through crowd-sourcing. Their solution allows users, such as network owners or administrators, to label devices that the monitoring app couldn’t accurately identify. However, this solution needs extensive network coverage to work optimally, and a recent study indicated its limitations in accuracy [26].

In this study, we demonstrate how recent advancements in Large Language Models (LLMs) can address the challenging task of IoT device labeling. Our algorithm uses a catalog of

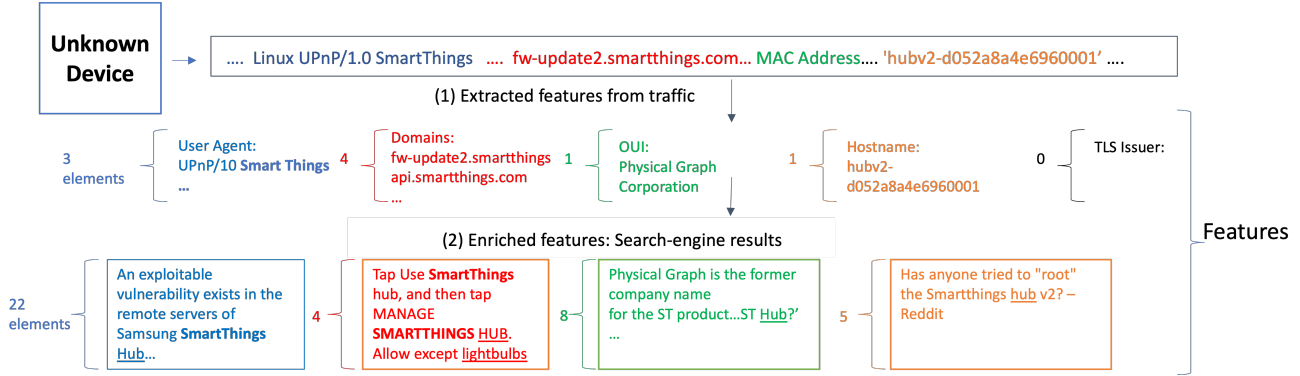


Figure 1: Example of Features for the SmartThing Hub: First, we present the features derived from the traffic, followed by a sample of the enriched features (the color correlates between the feature and the enriched feature). Words relevant to the vendor label decision are highlighted in bold, and those relevant to function decisions are underlined.

known vendors and the various IoT functions they produce to identify known devices. However, by harnessing the power of LLMs, it can automatically update this catalog when a new device type appears.

The algorithm starts by extracting textual strings from the network’s traffic log. Specifically, it extracts the following features: domains to which the devices connect, hostnames, TLS issuers, OUI (after the MAC lookup [22]), and user-agents. Note that different amounts of textual values can be associated with each feature type. For instance, a single IoT device can connect to multiple domains and may have different TLS issuers for each domain.

The algorithm then takes the text features and feeds them into a search engine so the search results’ descriptions can be used to enrich the data. Because these results typically include pages from IoT manufacturers, online shopping websites selling the device, and security blogs discussing the device’s features, such searches often yield references to the vendors and functions of the IoT devices. Figure 1 shows an example of the enriched features, and how they reveal the vendor and function. The enriched features give us a large text dataset in which we use string matching to find the vendor of the device and LLMs to determine the function of the device. Note, however, that different enriched features can result in different labels, as seen in our example. Choosing the most accurate label among several possible ones is one of the challenges addressed by our solution. It was vital that our algorithm determine how to weigh the answers from different enriched feature values, according to the feature type. To address this challenge, we employed a machine-learning approach.

Our labeling algorithm consists of two steps: using string matching to identify the vendor and using LLMs with zero-shot classification to label the function. In the first step, we identify the vendor by performing string matching of the text results in our dataset against a catalog of vendors and create labels based on the most common ones. In our dataset, which comprises 97 unique IoT devices from 55 distinct vendors,

the algorithm achieves a high accuracy rate: 86% for HIT1 and 89% for HIT2 ¹. Although there is a misconception that the OUI information can reveal the vendor, we observed that it achieves only a 64% accuracy. This is primarily because in IoT, the OUI often identifies the NIC vendor rather than the actual device vendor.

In the second step, the algorithm receives the catalog of potential functions for the vendor identified in the first step. Here, we leverage the fact that most IoT vendors specialize in producing a limited range of device functions. We employ large language models (LLMs) and use zero-shot classification [12] to map each feature value to its relevant functions. Specifically, we used Roberta [13] to classify each feature into possible functions and gathered the confidence scores of these classifications. Our method then aggregates these scores to establish an overall confidence level for each label. The label with the highest confidence level is then assigned as the device’s function. In our dataset, which includes 21 functions, the algorithm achieves an accuracy of 70% for HIT1 and 77% for HIT2. Since our algorithm works on textual features, it can provide explanations to humans who wish to verify its result, by outputting the set of features upon which its decision was based.

2 IoT Labeling Requirements

The requirements we outlined for an efficient and automated IoT labeling solution are as follows:

- **Universality:** The algorithm is aimed at processing data from IoT devices it has never before encountered. The algorithm leverages the use of a catalog that contains possible vendors and functions they support. The algorithm is able to update this catalog when the device type does not already appear in the catalog.

¹HIT2 measures whether the correct vendor was among the top two suggested labels.

- **Accuracy:** High accuracy is paramount in IoT labeling due to its critical role in enhancing network observability and security.
- **Explainability:** The solution incorporates a confidence level and justification for each label, allowing for potential human verification. This aspect caters to scenarios where it is better to have an inaccurate label than one that is incorrect.
- **Passive solution:** The IoT labeling operates without actively probing or querying the device for information, as is done in Shodan [40] and Censys [9]. We note that in many cases, this is not a possible or informative solution.
- **Offline process:** The IoT labeling operates offline by collecting the features over long periods (hours or days). The algorithm labeling running time is not a critical parameter since it is run only when it encounters a new device that was not recognized by the system.

3 Background: Zero-Shot Classification

In traditional ML classification tasks, a model is trained on a specific set of classes using samples from each type. In our work, each of these classes is equivalent to a new IoT type. If a new class emerges after the model has been trained, the entire model either needs to be retrained from scratch or fine-tuned, necessitating new labeled data for that class. This iterative process is both time consuming and labor intensive, which makes it not sustainable, especially in dynamic environments like IoT. In the IoT realm, it is generally not feasible to obtain samples from every type, which means the set of classes will evolve and expand over time.

Zero-shot classification is unique in that it classifies text without requiring explicit training on the target classes. Instead of learning representations only for classes that were seen during the training phase, zero-shot classification leverages the semantic relationships between classes. It often uses auxiliary information such as class attributes or textual descriptions to perform the classification. This enables the model to infer classes it has never seen during training, create new labels, and efficiently update the model with new classes. In our case, the IoT catalogs that serve as the input to our algorithm need to be updated based on the zero-shot classification; however, the model itself and its optimized parameters do not require re-training once the catalog has been updated.

In essence, zero-shot classification, with its ability to dynamically adapt to new classes and its minimal reliance on exhaustive labeled data, provides a solution to the challenges of IoT labeling.

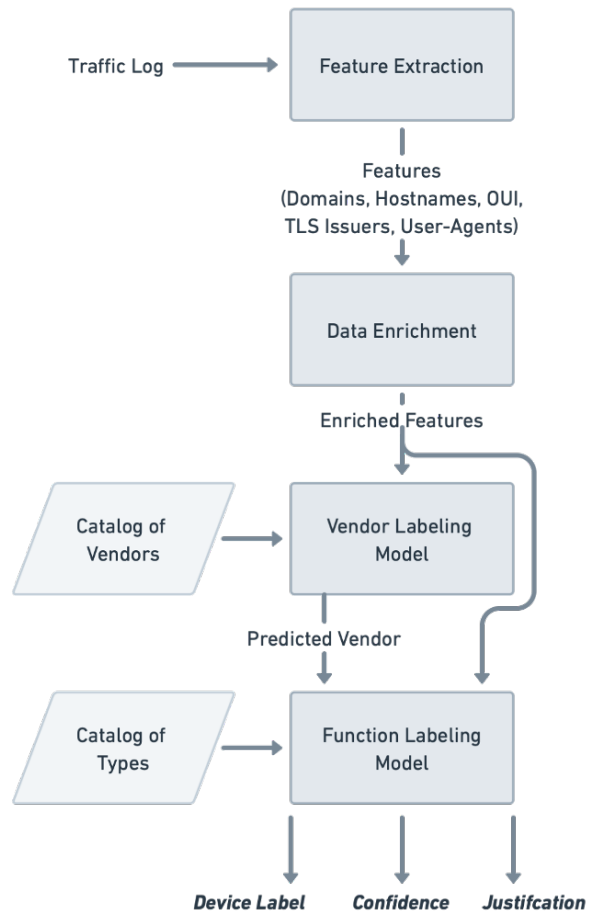


Figure 2: A schematic illustration of our IoT labeling solution. First, features are being extracted and then enriched. Second, we perform our vendor and function models labeling. The system’s output is label, confidence and justification for each device.

4 Dataset

We trained and tested our IoT labeling algorithm on a dataset that amalgamates five distinct open-source datasets ([28], [36], [42], [35], and [8]). These datasets were selected by previous research to provide a representation for identification-oriented tasks. The dataset includes traffic logs of 161 IoT devices², from 55 distinct vendors, categorized across 21 different functions. We note that the traffic logs are of different lengths, varying from hours to days.

For our testing, we randomly selected a single device from each group of devices that shared the same type (i.e., vendor and function), in order to avoid bias arising from multiple repeated instances. All the subsequent experiments discussed in this paper were based on analyses conducted on this collection of 97 unique devices that appear in the dataset, unless explic-

²Each device can be identified uniquely by its MAC address

itly stated otherwise. Nonetheless, our experiments showed that devices of the same type often have lower similarity when it comes to feature values, as detailed in Section 6. This outcome likely occurred because a single device type can encompass a range of different IoT models. For example, you can have various camera models from the same IoT vendor. Because we lacked the information to determine whether the IoT devices shared the same model, we decided to take the restrictive approach and include only one device per type.

The devices in our datasets are labeled with vendor and function since they were recorded in the lab. An example of a label with vendor and function would be: Belkin Plug. By manually examining the device’s name on the internet, we were able to translate it into a function. This labeled dataset provided us with a ground truth basis for our subsequent experiments and models.

4.1 Enriched Features

Our labeling algorithm treats and processes the traffic from each IoT device separately. Based on the traffic, it derives all the feature values from all the feature types (hostname, domains, TLS Issuers, OUI, User-Agents). The feature values then undergo an enrichment process in which a search query is performed for each feature value; our solution extracts at most the top 10 results from this query. We executed our data enrichment phase by using SerpAPI [39] to access Google search results across the web. During this process, most of the feature values returned less than 10 results, but 60% returned more than 8 results, as shown in Figure 3.

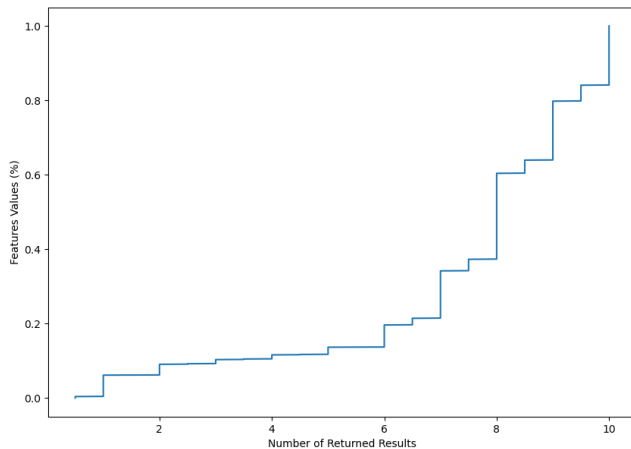


Figure 3: CDF distribution showing the number of results returned per feature value

The structure of each search result from SerpAPI [39] is composed of a title and snippet of the web page; both of these were used as enriched data. Google uses a number of different sources to automatically determine the title of the search result. The snippets are automatically created from the page

content [17], which usually contains the search query itself. In our case, the search result data often encompassed sites owned by the device vendor and occasionally linked to public forums where device users discuss their device’s network behavior. In some cases, we even found links to eCommerce platforms that sell devices, which enabled the identification of specific device models.

We denote the enriched version of feature type f , by f^+ . Correspondingly, we define the vector of the enriched feature as follows:

Definition 1 Let $f_{t,p}^+$ be the vector $= (f_{t,p,1}^+ \dots f_{t,p,k}^+)$ of the enriched feature values vector with k top search results on feature value $f_{t,p}$ of feature type t . $f_{t,m,j}^+$ is the j -th enriched feature returned from a search on $f_{t,m}$ feature t value.

Table 1 summarizes the terminology used throughout the paper. All the devices in the dataset, including all the features with their enrichment data, are publicly available in JSON format at [3] for further research into IoT device labeling. The format of the JSON is given in Appendix, Table 1.

5 Labeling Algorithm

In this section, we present the labeling algorithm. To determine the device type, we first determine the vendor label and then find the function by examining which IoT device functions the identified vendor manufactures. Vendor and function labeling essentially operate in similar ways. Each enriched feature for a single IoT device is matched against the corresponding vendor or function catalog³. We checked several matching algorithms in our experiments, including string matching, Roberta model [11], and ChatGPT [32]. When it comes to vendor labeling, the naive string matching provided the best performance, while the Roberta model outperformed the other options for function labeling. Although Roberta is adept at handling context-driven classification, it’s not ideal for vendor labeling. This is primarily due to two reasons. First, zero-shot classification is not designed to effectively handle hundreds of distinct labels, as in the case of vendors. Second, the vendor classification task depends less on context and more on specific strings, making string-matching approaches more suitable. The Roberta model is suitable to use context and semantic similarity for the function labels because (a) it is based on Transformer architecture [13] and (b) the function catalog size is small enough for it to handle.

Next, for each enriched feature value (i.e., search result) and potential label, the algorithm produces a confidence score indicating the likelihood that the enriched feature value correctly classified the potential label.

Formally:

³We focus only on enriched features since our experiments show that adding information from the basic features does not improve accuracy, mainly because the search value typically appears in the enriched data.

Term	Description
$V = \{V_1, V_2, \dots, V_{n_V}\}$	Catalog of Vendors of length n_V
$F = \{F_1, V_2, \dots, F_{n_F}\}$	Catalog of functions of length n_F
$T = \{T_1, T_2, \dots, T_{n_T}\}$	Catalog of types of length n_F where $T \subseteq V \times F$
$FT = \{hostname, domains, TLS, OUI, UserAgent\}$	Feature type set (size of 5)
$f_t = (f_{t,1} \dots f_{t,m})$	Feature values vector of feature type $t \in FT$
$f_{t,p}^+ = (f_{t,p,1}^+ \dots f_{t,p,k}^+)$	Enriched feature values vector with k top search result on $f_{t,p}$
$S_{F_m, f_{t,p,l}^+}$	Confidence scores labeling $f_{t,p,l}^+$ with F_m
w_t	Weight per feature type $t \in FT$
θ_t	Confidence threshold per feature type $t \in \theta$

Table 1: Terms and Terminology for Labeling Algorithms and Scoring Techniques

Definition 2 Let $S_{label, f_{t,m,j}^+}$ denote the confidence score that the enriched feature value $f_{t,i,j}^+$ corresponds to the label. The label can be a vendor label, i.e., $V_m \in V$ in vendor labeling or can be a function label, i.e., $F_m \in F$ in function labeling.

In vendor labeling, the algorithm produces a confidence score based on the number of times a label correctly matches the enriched feature.

For the function labeling, we used the Roberta confidence score [11]. To enhance accuracy, the algorithm performs the labels of only functions associated with the identified vendor, when available; otherwise, it performs full check with all the functions in the catalog.

Next, our algorithm calculates an aggregated score for each potential label. This is done by weighting the confidence score based on the original feature type of the enriched feature and the number of results returned for that feature type. The optimization of this scoring is detailed in the next section. The label with the highest score is then selected.

Algorithm 1 Vendor Labeling Algorithm

- 1: **Input:** A traffic log of an IoT device with enriched feature values $f_{t,i,j}^+$ for all possible value types $t \in FT$.
- 2: **Perform:**
- 3: **for** each label vendor $V_m \in V$ **do**
- 4: **for** every enriched feature value $f_{t,i,j}^+$ **do**
- 5: Let $S_{V_m, f_{t,m,j}^+}$ be the confidence score,
- 6: the number of times the string appears
- 7: V_m at the enriched feature $f_{t,i,j}^+$.
- 8: **end for**
- 9: Calculate the aggregate score of V_m based on the confidence score across all enriched feature values and types.
- 10: **end for**
- 11: **return** the vendor label with the highest score.

To assess the contribution of each feature type to the overall accuracy, we calculated the accuracy, which is the ratio of correct device labels obtained when running the labeling

Algorithm 2 Algorithm for Function labeling

- 1: **Input:**
- 2: A traffic log of an IoT device with enriched feature values $f_{t,i,j}^+$ for all possible value types $t \in FT$.
- 3: Let *Vendor* be the identified vendor after the vendor labeling.
- 4: **Perform:**
- 5: Let *FL* be a set of candidate relevant function labels in
- 6: type catalog $FL = \{F_m \mid (vendor, F_m) \in T\}$.
- 7: If $FL = \emptyset$ then $FL = F$, i.e., all possible functions.
- 8: Calculate the confidence values:
- 9: **for** each label function $F_m \in FL$ **do**
- 10: **for** each enriched feature value $f_{t,i,j}^+$ **do**
- 11: Given the Roberta model to labeling $f_{t,m,j}^+$
- 12: against possible labels *FL*.
- 13: Let $S_{F_m, f_{t,m,j}^+}$ be the confidence score of $F_m \in FL$.
- 14: **end for**
- 15: Aggregate the score of F_m based on the confidence score across all enriched feature values and types.
- 16: **end for**
- 17: **return** *Function*, the function label with the highest score.
- 18: **return** the type of the IoT, (Vendor, Function).

algorithm on each specific feature type. Figure 4a shows the accuracy of string matching in predicting the vendor and Figure 4b shows the accuracy of Roberta in predicting the device function. As can be observed, different enriched features exhibit varying levels of accuracy, where Domains⁺ and Hostname⁺ have the highest accuracy in both vendor and function labeling algorithms. Moreover, as we noted before, the enriched feature is more informative than the basic feature. We also found that it is harder to predict the function than the vendor. One clear outcome is that we must weight the label differently according to the feature type and concentrate on the enriched features.

Our vendor labeling algorithm is outlined in Algorithm 1,

and our function labeling algorithm is detailed in Algorithm 2.

6 Experiments Results

In this section, we present the experimental results of our IoT labeling. We compare our solution to a few other common methods that are common in industry and academia. First is Fing product [15], a labeling solution that is common in the industry, and is used by leading cyber-security firms. Second is the OUI method, based on the MAC address. In addition, we present a few versions of our method with more naive approaches. To evaluate the effectiveness of our methods, we present HIT1 which is the ratio of devices for which the method correctly predicted its label to the total number of devices in the dataset. Similarly, we present HIT2, the ratio of devices for which the correct label is in the top two labels predicted. We also present the ratio of devices that could not be labeled, where the algorithm did not return any results, for example, if there were no features or no enriched features above the threshold. In all the techniques that required optimization learning, we used 5-fold cross-validation to ensure the training and testing were done on different devices.

We found a mostly naive approach to vendor labeling that uses OUI, the first 24 bits of a MAC address that identify the vendor of a network interface card (NIC) or network adapter [22]. This technique does not require the vendor catalog but has a lower accuracy of 0.64 since it often identifies the NIC vendor, which is usually different from IoT vendor. We also checked GPT-4’s ability to predict the vendor based on device features. Given the limited number of tokens that can be used with this large language model, we restricted our input to the original dataset features, without enrichment. Because GPT-4 learns the Internet, it has enrichment information at some level.

Our vendor labeling uses string matching based on all the enriched features. This method yielded the highest results with an HIT1 accuracy of 0.86 and an HIT2 accuracy of 0.89. The GPT-4 model, without the assistance of a vendor catalog, achieved a respectable accuracy of 0.83 (HIT1) and 0.86 (HIT2). This suggests that while GPT-4’s contextual understanding is excellent, more straightforward methods like string matching outperform it for vendor labeling.

Figure 5 shows the confidence score in our vendor labeling. It is evident that these are definite results, with a high occurrence of the string in the data; in most cases, the highest score is much higher than the second score. Note that the y-axis is in log scale.

We conducted an experiment to measure the accuracy of our labeling algorithm across each feature and gauge the impact of different features on the method’s success. In Figure 6 we show the number of correct results as a function of the number of returned results, where zero indicates no enrichment. It is evident that while enrichment enhances the labeling

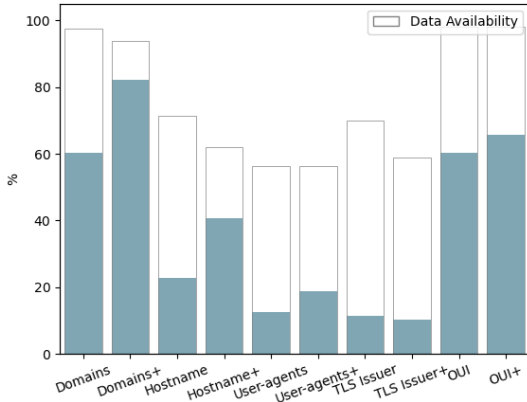
of vendors using the string matching method, the improvement is not as substantial compared to how much enrichment impacted the function labeling method. Furthermore, some features, such as TLS Issuer, appear to be less affected, and enrichment does not elevate their performance at all. We did not see any notable improvement for vendor labeling, when we performed the same optimization (weighting features), described earlier this section.

Table 3 shows the results for function labeling. We compared our method’s results with several variations of labeling method using more simple string matching or advanced (GPT-4) methods. When analyzing the performance of the methods using the full catalog of functions, Roberta performs much better than the string matching method. We also observed that without enriched features, Roberta performance is very low (HIT1 accuracy of 0.16). However, with enriched features, Roberta exhibits a commendable improvement, achieving an HIT1 accuracy of 0.56 and an HIT2 accuracy of 0.65. When functions are restricted to those associated with the vendor (represented by a ‘V’ in Table 3), the algorithm receives the catalog of potential functions for the identified vendor, and leverages the fact that most IoT vendors specialize in producing a limited range of functions. The GPT-4 model, without the assistance of a type catalog, achieved an accuracy of 0.65 (HIT1). When we used the GPT model together with the type catalog, there was a slight improvement. We assume that the type catalog had only a slight impact on the GPT model because GPT had this knowledge when labeling, even without the type catalog. Roberta, based on the types catalog using all enriched features and weighted using our optimization process, surfaces as the top performer with an HIT1 accuracy of 0.7 and an HIT2 accuracy of 0.77; this significantly outperformed the string matching method.

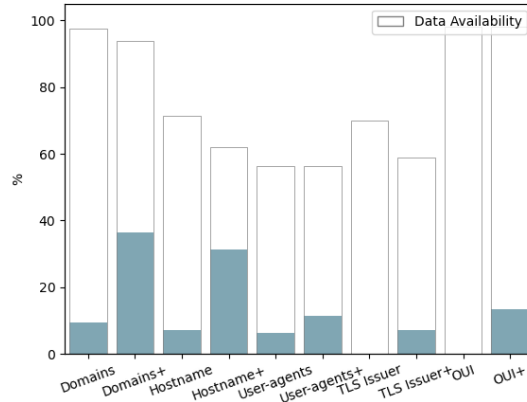
Accuracy of Type Labeling The accuracy of type labeling refers to the successful labeling of both a device’s vendor and function. Utilizing a combination of string matching for vendor labeling and Roberta for function labeling, the best accurate method for each task, yield a HIT1 accuracy of 0.6, a HIT2 accuracy of 0.7.

Results for All vs. Unique Devices As mentioned in Section 4, we randomly selected devices for type in order to avoid repetitive devices that may bias our results. In this subsection, we show the results of our method on all the devices alongside a description of the diversity of the data within devices of the same vendor and function. In Tables 2 and 3 we present the results of our labeling method on *all* the devices in parentheses. Our function labeling method achieved 0.74 for HIT1 and 0.85 for HIT2. While our vendor labeling resulted in 0.89 and 0.92 for HIT1 and HIT2, respectively.

To verify that no recurring devices were included in our comprehensive collection, we employed the Jaccard similarity coefficient, which quantifies the degree of overlap between two sets. Each device was scrutinized for similarities across all features. This process allowed us to calculate an average



(a) Vendor Labeling with String-matching Algorithm.



(b) Function Labeling with Roberta Algorithm.

Figure 4: Comparative analyses of labeling accuracy per feature, indicated by filled bars and availability, indicated by hollow bars (represents the percentage that the feature exists across the dataset). Figure 4b presents the accuracy of the device function while Figure 4a presents the vendor.

Method	Catalog of Vendors	Features	Accuracy HIT1	Accuracy HIT2
OUI	N	MAC	0.64	-
Our IoT labeling	Y	Domains⁺, Hostname⁺, TLS⁺, User-Agents⁺, OUI⁺	0.86 (0.89)	0.89 (0.92)
Our IoT labeling	Acquired	Domains ⁺ , Hostname ⁺ , TLS ⁺ , User-Agents ⁺ , OUI ⁺	0.76	0.8
GPT-4	N	Domains, Hostname, TLS, User-Agents, OUI	0.83	0.86
Fing	N	Hostname, User-Agents, MAC	0.76	-

Table 2: Accuracy of vendor labeling, where the catalog of vendors is given (Y=yes) or (N-not) or Acquired on-the-fly. The accuracy results for the whole dataset appear in parentheses

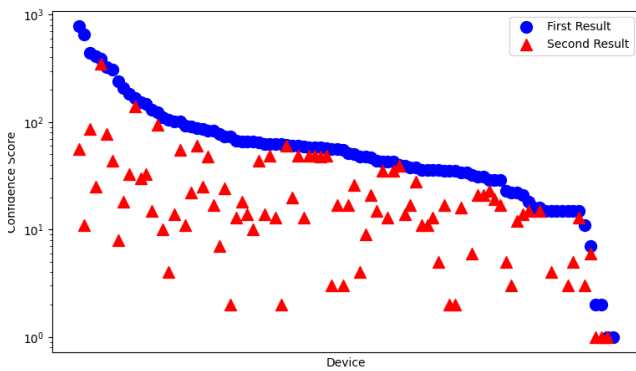


Figure 5: Confidence score per device of our vendor labeling algorithm (equals to the number of label matching in the enriched data), ordered by the first score matching. The highest score is marked in a circle and the second highest is marked in a triangle.

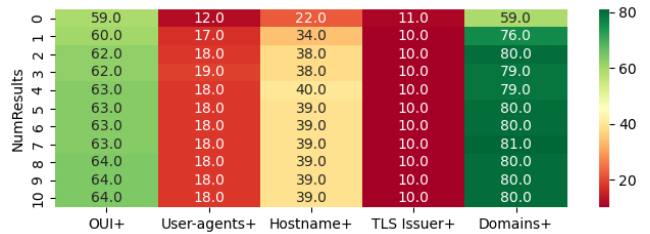


Figure 6: The distribution of correct results of vendor labeling for various features using the string matching method. Darker shades indicate higher numbers of correct results.

Method	Catalog of Functions	Features	Accuracy HIT1	Accuracy HIT2
String Matching	A	Domains ⁺ , Hostname ⁺ , TLS ⁺ , User-Agents ⁺ , OUI ⁺	0.49	0.63
Roberta	A	Domains ⁺ , Hostname ⁺ , TLS ⁺ , User-Agents ⁺ , OUI ⁺	0.56	0.65
Roberta	A	Domains, Hostname, TLS, User-Agents, OUI	0.16	0.2
GPT-4	A	Domains, Hostname, TLS, User-Agents, OUI	0.65	0.75
GPT-4	V	Domains, Hostname, TLS, User-Agents, OUI	0.67	0.81
Fing	N	Hostname, User-Agents, MAC	0.57	-
Our Function Labeling	V	Domains⁺, Hostname⁺, TLS⁺, User-Agents⁺, OUI⁺	0.7 (0.74)	0.77 (0.84)

Table 3: Function Labeling Accuracy. 'A' denotes accuracy when all items in the function catalog are given, while 'V' represents accuracy when only the functions associated with the device vendor are provided.

Table 4: Similarity Measures for Features. This table summarizes the average similarity and standard variation for features within groups of devices (same vendor and function).

Key	Average Similarity \pm SV
Domains	0.266300 \pm 0.250173
User-agents	0.075917 \pm 0.240749
Hostname	0.175000 \pm 0.352959
TLS Issuers	0.531568 \pm 0.382723

similarity for each pair of devices based on these keys. This process facilitated the calculation of a similarity measure within groups of devices sharing the same vendor and function. Devices from the same vendor and function are naturally expected to share certain similarities, yet they can still exhibit distinguishing patterns due to different factors such as device model, hardware configuration, and software version.

Our findings, as depicted in Table 4, confirmed that the average similarity between devices was reasonably low, thereby suggesting a high degree of uniqueness across the samples in the same group.

7 Related Work

The field of device identification is vast and diverse, and various strategies and techniques have been developed to address different aspects of this problem. The majority of work focuses on identifying known or previously seen devices. Also, some previous researches focus on the identification of specific actions (turn on/off, etc.) [41] or a specific unique model of a device, all presented in Table 5.

In contrast, this paper focuses on the task of labeling the vendor and the function of unknown or previously unseen IoT devices. Existing methodologies for device identification rely on labeled data, meaning that prior knowledge about the devices under scrutiny is a prerequisite. These details are predominantly gathered through passive or active [7, 14, 45] methods, where network samples produced by a device are monitored and subsequently used to generate an identification pattern [6]. This approach has been adopted in numerous

studies, as highlighted by [2, 5, 18, 19, 27, 33]. All mentioned works identify only seen devices. Our research diverges from these conventional methodologies and addresses a more intricate and advanced issue, the labeling of previously unseen devices.

There are a limited number of studies that have ventured into the realm of identifying unseen devices. Wang et al. [44] present a method to identify unseen devices using active probing methods. Le et al. [23] proposed two algorithms for vendor identification of unseen devices; they used a blend of domain owners, certificate owners, and the OUI among other parameters, to reach a correct vendor classification rate of 72 devices out of 94 devices (76%).

Type classification was done on seen data only with an extensive training phase, which is both costly and labor intensive. Bai et al. [4] uses ML, specifically time-series, to classify device types. They segmented each device traffic into fixed-time sub-flows. The final dataset consisted of 15 devices belonging to 4 categories and achieved 75% accuracy. Table 5, presents a summary of several papers, organized into the 3 categories (data, output, and methodology).

We note that in the IETF's Manufacturer Usage Description (MUD) framework [24], the IoT device self-identifies, and an allow-list (access control list) is fetched. However, there has been no real-world adoption of MUD to date.

8 Conclusion & Future Work

This study demonstrates that, by combining recent advances in LLM and NLP, we can effectively label unseen IoT devices in the wild. Over a set of 97 unique IoT devices, our function labeling approach achieved HIT1 and HIT2 scores of 0.7 and 0.81, respectively.

In the IoT realm, it is generally not feasible to obtain samples from every type. We used zero-shot classification models, leveraging the fact that it allows the update of the labels with no need of re-training the model itself. As far as we know, this work is the first to address the IoT device labeling problem.

For future work, there is room for exploration in the realm of LLM explainability. We propose additional refinements to the LLM to deliver improved explainability by fine-tuning

Paper	Data	Output	Methodology
	Seen(S)\ Unseen(U)	Type(T) \ Vendor(V) \ Unique (U) \ Actions (A)	Active (A) \ Passive(P)
[23]	S U	T V	P P
[4]	U	T	P
[1]	S	T	P+A
[5]	S	T	P
[29]	S	U	P
[30]	S	U	P
[44]	U	T+V	A
[27]	S	U	P
[33]	S	U	P
[43]	S	U	P
[41]	S	U+A	P
[37]	S	U	P
[21]	S	T+V	P
[20]	S	U	P
[45]	S	T+V+U	A

Table 5: Existing IoT identification mechanisms based on 3 key categories data (Seen vs Unseen), output (Type, Vendor, Unique, specific Actions) and methodology (Active vs Passive).

the model over multiple examples with human-written explanations. The expected outcome is that these advancements will not only facilitate a deeper understanding of the labeling procedure but will also foster greater confidence in the output of the model.

We also plan to ascertain the validity of the explanations provided by our model. By manually tagging these explanations, we can measure their semantic similarity to human-provided explanations by employing a regression model such as DistilBERT [38] as the scoring function [10]. Quantifying this similarity will allow us to assign a score to each explanation, providing a reference for future improvements to our model.

We also intend to explore refinements in the features that are collected and enriched. For example, some of the domain names are AWS services that are used by many vendors, and they distract the labeling algorithm. In our ongoing work, we identify these domains and omit them from the catalog of features.

References

[1] Nesrine Ammar, Ludovic Noirie, and Sébastien Tixeuil. Network-protocol-based iot device identification. In *2019 Fourth International Conference on Fog and Mo-*

bile Edge Computing (FMEC), pages 204–209. IEEE, 2019.

[2] Sandhya Aneja, Nagender Aneja, and Md Shohidul Islam. Iot device fingerprint using deep learning. In *2018 IEEE international conference on internet of things and intelligence system (IOTAIS)*, pages 174–179. IEEE, 2018.

[3] Anonymous. Dataset devices json: Enriched data., 2023.

[4] Lei Bai, Lina Yao, Salil S Kanhere, Xianzhi Wang, and Zheng Yang. Automatic device classification from network traffic streams of internet of things. In *2018 IEEE 43rd conference on local computer networks (LCN)*, pages 1–9. IEEE, 2018.

[5] Jiaqi Bao, Bechir Hamdaoui, and Weng-Keen Wong. Iot device type identification using hybrid deep learning approach for increased iot security. In *2020 International Wireless Communications and Mobile Computing (IWCMC)*, pages 565–570. IEEE, 2020.

[6] Genevieve Bartlett, John Heidemann, and Christos Papadopoulos. Understanding passive and active service discovery. In *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*, pages 57–70, 2007.

[7] Sergey Bratus, Cory Cornelius, David Kotz, and Daniel Peebles. Active behavioral fingerprinting of wireless devices. In *Proceedings of the first ACM conference on Wireless network security*, pages 56–61, 2008.

[8] Anat Bremler-Barr, Bar Meyuhas, and Ran Shister. One mud to rule them all: Iot location impact. In *NOMS 2022-2022 IEEE/IFIP Network Operations and Management Symposium*, pages 1–5, 2022.

[9] Censys. Censys - internet security data for everyone. <https://censys.com/>, 2023. Accessed: [Insert date you accessed the website].

[10] Lingjiao Chen, Matei Zaharia, and James Zou. Frugal-gpt: How to use large language models while reducing cost and improving performance, 2023.

[11] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*, 2019.

[12] Hugging Face. What is Zero-Shot Classification? <https://huggingface.co/tasks/zero-shot-classification>, Aug 2020. Accessed: 2023-Aug-02.

- [13] Hugging Face. xlm-roberta-large-xnli model and fine-tunes it on a combination. <https://huggingface.co/joeddav/xlm-roberta-large-xnli>, June 2020. Accessed: 2022-Feb-02.
- [14] Xuan Feng, Qiang Li, Haining Wang, and Limin Sun. Acquisitional rule-based engine for discovering iot devices. In *27th USENIX Security Symposium (USENIX Security 18)*, pages 327–341, 2018.
- [15] Fing. Device recognition user guide. https://www.fing.com/images/uploads/general/Device_Recognition_User_Guide.pdf, 2018.
- [16] 2023 FIREDOME. Iot endpoint protection platform, 2023. Accessed: 2023-07-28.
- [17] Google. Google search central, 2023. [Online; accessed 13-Oct-2023].
- [18] Hang Guo and John Heidemann. Ip-based iot device detection. In *Proceedings of the 2018 workshop on IoT security and privacy*, pages 36–42, 2018.
- [19] Ibbad Hafeez, Markku Antikainen, Aaron Yi Ding, and Sasu Tarkoma. Iot-keeper: Detecting malicious iot network activity using online traffic analysis at the edge. *IEEE Transactions on Network and Service Management*, 17(1):45–59, 2020.
- [20] Muhammad Hamza, Syed Mashhad M Geelani, Qamar Nawaz, Asif Kabir, and Isma Hamid. Clustering of iot devices using device profiling and behavioral analysis to build efficient network policies. *Mehran University Research Journal Of Engineering & Technology*, 40(2):335–345, 2021.
- [21] Pratibha Khandait, Neminath Hubballi, and Bodhisatwa Mazumdar. Iothunter: Iot network traffic classification using device specific keywords. *IET Networks*, 10(2):59–75, 2021.
- [22] Deniel Laurent. Wireshark manufacturer database.
- [23] Franck Le, Jorge Ortiz, Dinesh Verma, and Dilip Kandlur. Policy-based identification of iot devices’ vendor and type by dns traffic analysis. In *Policy-Based Autonomous Data Governance*, pages 180–201. Springer, 2019.
- [24] Eliot Lear, Ralph Droms, and Dan Romascanu. Manufacturer Usage Description Specification. RFC 8520, March 2019.
- [25] Check Point Software Technologies Ltd. Quantum iot protect, 2023. Accessed: 2023-05-22.
- [26] A. Mandalari, H. Haddadi, D. J. Dubois, and D. Choffnes. Protected or porous: A comparative analysis of threat detection capability of iot safeguards. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 3061–3078, may 2023.
- [27] Samuel Marchal, Markus Miettinen, Thien Duc Nguyen, Ahmad-Reza Sadeghi, and N Asokan. Audi: Toward autonomous iot device-type identification using periodic communication. *IEEE Journal on Selected Areas in Communications*, 37(6):1402–1412, 2019.
- [28] Markus Miettinen, Samuel Marchal, Ibbad Hafeez, N. Asokan, Ahmad-Reza Sadeghi, and Sasu Tarkoma. Iot sentinel: Automated device-type identification for security enforcement in iot. In *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*, pages 2177–2184, 2017.
- [29] Markus Miettinen, Samuel Marchal, Ibbad Hafeez, N Asokan, Ahmad-Reza Sadeghi, and Sasu Tarkoma. Iot sentinel: Automated device-type identification for security enforcement in iot. In *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*, pages 2177–2184. IEEE, 2017.
- [30] Nizar Msadek, Ridha Soua, and Thomas Engel. Iot device fingerprinting: Machine learning based encrypted traffic analysis. In *2019 IEEE wireless communications and networking conference (WCNC)*, pages 1–8. IEEE, 2019.
- [31] Palo Alto Networks. Introduction to iot security. <https://docs.paloaltonetworks.com/iot/iot-security-admin/iot-security-overview/introduction-to-iot-security#id235c0535-8183-48ce-b591-9f4f512fa914>. Accessed: 2023-05-19.
- [32] OpenAI. Openai gpt. <https://openai.com/research/gpt-4>, 2023. Accessed: 2023-Oct-01.
- [33] Jorge Ortiz, Catherine Crawford, and Franck Le. Devicemien: network device behavior modeling for identifying unknown iot devices. In *Proceedings of the International Conference on Internet of Things Design and Implementation*, pages 106–117, 2019.
- [34] Inc. Palo Alto Networks. Palalto: Iot security administrator’s guide.
- [35] Roberto Perdisci, Thomas Papastergiou, Omar Alrawi, and Manos Antonakakis. Iotfinder: Efficient large-scale identification of iot devices via passive dns traffic analysis. In *2020 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 474–489, 2020.

- [36] Said Jawad Saidi, Anna Maria Mandalari, Roman Kolcun, Hamed Haddadi, Daniel J Dubois, David Choffnes, Georgios Smaragdakis, and Anja Feldmann. A haystack full of needles: Scalable detection of iot devices in the wild. In *Proceedings of the ACM Internet Measurement Conference*, pages 87–100, 2020.
- [37] Said Jawad Saidi, Anna Maria Mandalari, Roman Kolcun, Hamed Haddadi, Daniel J Dubois, David Choffnes, Georgios Smaragdakis, and Anja Feldmann. A haystack full of needles: Scalable detection of iot devices in the wild. In *Proceedings of the ACM Internet Measurement Conference*, pages 87–100, 2020.
- [38] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020.
- [39] SerpAPI. Google search api. <https://serpapi.com/>.
- [40] Shodan. Shodan - the search engine for the internet of things. <https://www.shodan.io/>, 2023.
- [41] Arunan Sivanathan, Habibi Gharakheili, Hassan, and Vijay Sivaraman. Managing iot cyber security using programmable telemetry and machine learning. *IEEE Transactions on Network and Service Management*, 17(1):60–74, 2020.
- [42] Arunan Sivanathan, Hassan Habibi Gharakheili, Franco Loi, Adam Radford, Chamith Wijenayake, Arun Vishwanath, and Vijay Sivaraman. Classifying iot devices in smart environments using network traffic characteristics. *IEEE Transactions on Mobile Computing*, 18(8):1745–1759, 2019.
- [43] Arunan Sivanathan, Hassan Habibi Gharakheili, and Vijay Sivaraman. Inferring iot device types from network behavior using unsupervised clustering. In *2019 IEEE 44th Conference on Local Computer Networks (LCN)*, pages 230–233. IEEE, 2019.
- [44] Ruimin Wang, Haitao Li, Jing Jing, Liehui Jiang, and Weiyu Dong. Wysiwyg: Iot device identification based on webui login pages. *Sensors*, 22(13):4892, 2022.
- [45] Kai Yang, Qiang Li, and Limin Sun. Towards automatic fingerprinting of iot devices in the cyberspace. *Computer Networks*, 148:318–327, 2019.