

CmdCaliper: A Semantic-Aware Command-Line Embedding Model and Dataset for Security Research

Anonymous EMNLP submission

Abstract

This research addresses command-line embedding in cybersecurity, a field obstructed by the lack of comprehensive datasets due to privacy and regulation concerns. We propose the first dataset of similar command lines, named CmdDataset, for training and unbiased evaluation. The training set is generated using a set of large language models (LLMs) comprising 28,520 similar command-line pairs. Our testing dataset consists of 2,807 similar command-line pairs sourced from authentic command-line data.

In addition, we propose a command-line embedding model named CmdCaliper, enabling the computation of semantic similarity with command lines. Performance evaluations demonstrate that the smallest version of CmdCaliper (30 million parameters) suppresses state-of-the-art (SOTA) sentence embedding models with ten times more parameters across various tasks (e.g., malicious command-line detection and similar command-line retrieval).

Our study explores the feasibility of data generation using LLMs in the cybersecurity domain. Furthermore, we release our proposed command-line dataset, embedding models' weights and all program codes to the public. This advancement paves the way for more effective command-line embedding for future researchers.

1 Introduction

Sentence embeddings, which map diverse sentences into a unified semantic feature space, are critical for various NLP applications such as classifier training, visualization (van der Maaten and Hinton, 2008), and retrieval-augmented generation (RAG) (Lewis et al., 2020). In cybersecurity, command lines provide invaluable information for detecting malicious attacks by comparing them with known historical malicious command lines from a semantic perspective. However, the flexibility in command-line syntax and structure poses challenges for fully leveraging this information. For

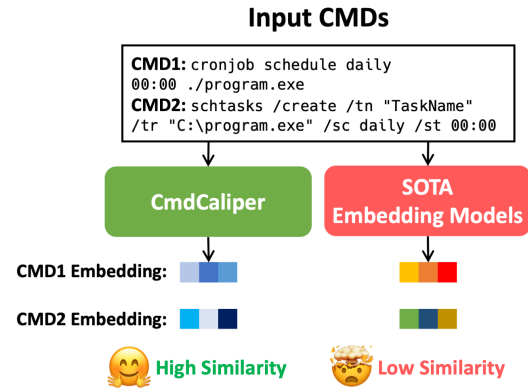


Figure 1: After fine-tuning our proposed similar command-line pair dataset, CmdDataset, our proposed command-line embedding model, CmdCaliper, can effectively embed command lines based on their semantics rather than solely on appearance.

example, as shown in Fig. 1, one can still correlate the two command lines according to their outputs despite the different appearances. To achieve this, using a robust embedding model to calculate the semantic similarity of command lines is promising. However, the grammatical differences between command lines and natural language sentences hinder the direct application of sentence embedding models to command-line tasks. Furthermore, one main challenge exacerbates the difficulty of research in command-line embedding: the scarcity of datasets specifically designed for command-line embedding tasks, both for training models and for fairly evaluating the performance of different methods.

To address the aforementioned challenges, this paper introduces the first comprehensive dataset, CmdDataset, which includes semantically similar pairs of command lines for both training and evaluating command-line embedding methodologies. Inspired by the successes of data synthesis by LLMs (Wang et al., 2023b,a), the similar command-line pairs in our training set are automati-

066 cally generated from a set of diverse command-line
067 seed initialized from multiple real-world sources
068 by a total of six distinct LLMs trained on diverse
069 datasets (§ 3.1). This facilitates a broader range
070 of command-line generation. For the testing set of
071 CmdDataset, to prevent training data leakage, we
072 directly employed a totally different data source
073 instead of synthesizing command lines by LLMs,
074 as done in the training set. (§ 3.2)

075 To the end, our training set consists of 28,520
076 similar command-line pairs, totaling 55,909 unique
077 command lines, and our testing set comprises
078 2,807 similar command-line pairs, totaling 5,576
079 unique command lines. Our dataset analysis and
080 human evaluation results (§ 5) demonstrate that
081 our pipeline can generate highly diverse and high-
082 quality similar command-line pairs.

083 Based on our proposed dataset, CmdDataset, we
084 also developed the first embedding model special-
085 ized for command-line embeddings, called Cmd-
086 Caliper. By encouraging semantically similar sam-
087 ples to come closer and simultaneously increasing
088 the distance between semantically dissimilar sam-
089 ples in the embedding space, CmdCaliper can em-
090 bed command lines into vectors from a semantic
091 perspective. As demonstrated in Fig. 1, even when
092 command lines differ in appearance, CmdCaliper
093 can still position them closely in the embedding
094 space based on their semantic meanings.

095 Our evaluation results (§6) demonstrate that even
096 the smallest version of CmdCaliper, with approxi-
097 mately 0.03 billion parameters, can surpass SOTA
098 sentence embedding models with ten times more
099 parameters (0.335 billion parameters) across vari-
100 ous command-line specific tasks, such as malicious
101 command-line detection, similar command-line re-
102 trieval, and command-line classification.

103 Our contribution is threefold. First, we pro-
104 pose the first dataset of similar command-line pairs
105 named CmdDataset, which allows for training and
106 performance evaluation. Through detailed valida-
107 tion of the dataset’s effectiveness, we believe it is
108 well-suited for further command-line research. Sec-
109 ondly, we explore the potential of using LLMs to
110 synthesize data in the cybersecurity domain. Our
111 experiments demonstrate that LLMs can indeed
112 generate high-quality and diverse data. Lastly, we
113 propose the first semantic command-line embed-
114 ding model, CmdCaliper. Our evaluations reveal
115 that a command-line-specific embedding model
116 significantly enhances performance across various

117 downstream tasks compared to generic sentence
118 embedding models. We open-source the entire
119 dataset, model weights, and all program codes un-
120 der BSD License at [AnonymousRepo](https://anonymous.4open.science/r/CmdCaliper/README.md)¹

121 2 Related Work

122 2.1 Semantic Embedding

123 Early works of sentence embedding such as
124 Word2Vec (Mikolov et al., 2013) and Glove (Pen-
125 ington et al., 2014), require the training of a pre-
126 defined static embedding lookup table to fuse into
127 the embedding vector of different sentences.

128 Recent works leverage well-trained language
129 models such as BERT (Devlin et al., 2019) and
130 T5 (Raffel et al., 2020) as pre-trained models to
131 globally embed inputs (i.e., the words of sentences)
132 into embedding vectors while considering con-
133 textual relationships to achieve impressive down-
134 stream performance. To fine-tune such models, a
135 contrastive learning scheme (van den Oord et al.,
136 2018) is adopted, as demonstrated in (Gao et al.,
137 2021; Chuang et al., 2022; Zeng et al., 2022; Nee-
138 lakantan et al., 2022; Wang et al., 2022, 2023a)

139 In the cybersecurity domain, semantic embed-
140 ding is crucial for computing semantic similar-
141 ity across various data types, including logs and
142 command lines. For log-based data, Golczynski
143 et al. (Golczynski and Emanuello, 2021) intro-
144 duced an autoencoder-based model to convert log
145 data into embedding vectors for anomaly detection.
146 Log2Vec (Liu et al., 2019) creates a heterogeneous
147 graph for each log dataset and uses the random
148 walk method with Word2Vec to embed log data.
149 LogBert (Guo et al., 2021) leverages BERT for
150 anomaly detection, clustering the embedding vec-
151 tors of normal samples to increase the separation
152 from abnormal samples.

153 For command-line based data, Ongun et
154 al. (Ongun et al., 2021) adapted the tokenization
155 methodology for command-line-based data and
156 followed the Word2Vec approach to train a pre-
157 defined embedding lookup table for a large amount
158 of command-line data. Conversely, Dong et
159 al. (Dong et al., 2023) directly adopted Word2Vec
160 for the command-line embedding.

161 2.2 LLMs in the Cybersecurity

162 In the field of cybersecurity, many re-
163 searchers (Motlagh et al., 2024; Divakaran
164 and Peddinti, 2024) have also explored leveraging

¹<https://anonymous.4open.science/r/CmdCaliper/README.md>

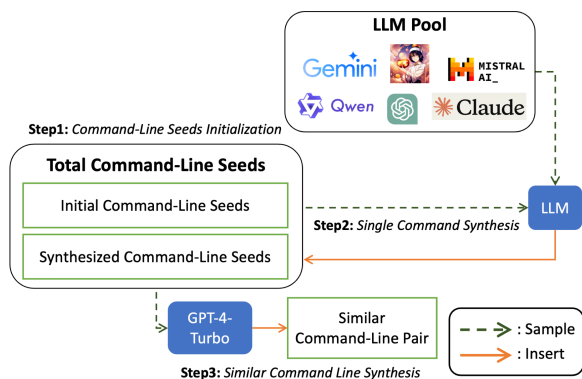


Figure 2: The illustration of the pipeline for automatically generating a dataset of similar command-line pairs using the Self-Instruct algorithm with a pool of LLMs.

LLMs (OpenAI, 2023; Anthropic, 2024; Touvron et al., 2023) for malicious code generations.

Moskal et al. (Moskal et al., 2023) use LLMs to generate executable code for agent actions, facilitating automated cyber campaigns. McKee et al. (McKee and Noever, 2023) explore using LLMs as honeypots by generating executable commands to simulate Linux, Mac, and Windows terminals. Chatzoglou et al. (Chatzoglou et al., 2023) demonstrate ChatGPT’s ability to generate malicious code that can evade detection.

3 CmdDataset: The First Command-Line Similarity Dataset

Despite the impressive performance of existing sentence embedding models (Li et al., 2023; Gao et al., 2021; Neelakantan et al., 2022), no embedding model has been designed specifically for command lines. We believe this is due to the unavailability of a large, diversified dataset with adequate annotations for effective training and unbiased evaluation.

In this section, we primarily focus on introducing the first command-line similarity dataset, named CmdDataset. This training set comprises 28,520 pairs of command lines automatically generated by a pool of LLMs. In contrast, the testing set contains 2,807 pairs of command lines collected from real-world attack scenarios.

3.1 Training Set Synthesis by LLMs

Collecting large-scale unlabeled or small-scale annotated command-line data is challenging due to two main factors. Firstly, labeling command-line datasets requires specialized cybersecurity knowledge, making it more stringent and costly than labeling images or natural language sentences. Sec-

ondly, privacy concerns involving company or personal information in real-world command lines discourage sharing, complicating efforts to gather diverse, large-scale datasets.

The automatic data generation process known as Self-Instruct (Honovich et al., 2023; Taori et al., 2023; Wang et al., 2023b) has proven effective in acquiring a comprehensive and diverse corpus of instructional data for fine-tuning LLMs. This process utilizes a powerful pre-trained large-scale language model, such as ChatGPT or Claude 3.

Our research is inspired by the substantial success of LLMs in code generation (Rozière et al., 2023; Patil et al., 2023) that exhibits similar structures to command lines and strong comprehension in the cybersecurity domain, such as malware generation (Pa et al., 2023; Botacin, 2023; Charan et al., 2023; Chatzoglou et al., 2023). Based on these capabilities, we adapt the Self-Instruct method to synthesize a substantial number of similar command-line pairs using LLMs. Our data synthesis pipeline comprises three stages: 1) Initial Seeds Collection, 2) Single Command Line Synthesis using a Pool of LLMs, and 3) Similar Command Line Synthesis, as illustrated in Fig. 2.

3.1.1 Initial Seeds Collection

Incorporating randomness into prompts is crucial for enabling LLMs to synthesize diverse command lines. This is achieved using initial seeds, which consist of a diverse set of command lines. During each synthesis iteration, a subset of these seeds is sampled to construct the prompt, diversifying it and encouraging a varied output from LLMs. To ensure high-quality initial seeds, we collected 2,061 diverse Windows command lines from multiple sources (e.g., public red-team exercises and Windows commands documentation). For more details on the initial-seed collection, see Appendix B.

3.1.2 Single Command Line Synthesis with a Pool of LLMs

Beyond the diversity of command-line seeds, the ability of LLMs to generate sufficiently varied data is also crucial. In the original Self-Instruct pipeline (Wang et al., 2023b), GPT-3 (Brown et al., 2020) was adopted for data generation, which confines all generated data to the distribution of GPT-3’s training data. We extended this method by constructing a pool of distinct LLMs to aid in data generation. This is intuitive because the distribution of training data varies among different LLMs, leading

to differences in the nature of the command lines they preferentially generate. The LLM pool of our pipeline comprises the following models: Mixtral 8x7B (Jiang et al., 2024), WizardLM-13B-v1.2 (Xu et al., 2024), Gemini-1.0-Pro (Team et al., 2023), Claude-3-Haiku (Anthropic, 2024), Qwen1.5-14B-Chat (Bai et al., 2023), and GPT-3.5-Turbo (OpenAI, 2022).

After constructing the LLM pool, we iteratively synthesize command lines by randomly sampling 12 command lines from the total command-line seeds—which include previously synthesized command lines and initial seeds—for prompt composition, as shown in Fig. 6. We then instruct a randomly selected LLM from the pool to synthesize four new command lines distinct from those in the prompt. Valid command lines are extracted from the LLM responses. Due to the randomness of generation, the LLMs may not always produce four new valid command lines. These valid new command lines are added to the total command-line seeds for the next iteration. The generation process is halted after synthesizing 28,520 command lines.

In this step, LLMs may synthesize non-executable command lines with minor syntax errors. However, for our dataset of similar command-line pairs, 'similar' refers to command lines sharing the same purposes or intentions, typically based on associated executable files, arguments, and argument values. Therefore, even with syntax errors, the command lines should still convey the same purpose or intention as their correct versions.

3.1.3 Similar Command Line Synthesis

After collecting 28,520 command lines, we instructed GPT-4-Turbo (OpenAI, 2023) to generate a similar command line for each. The prompting template for this instruction is displayed in Fig. 7. Here, 'similar' refers to sharing the same purpose or intention, rather than merely having a similar appearance. This distinction is crucial, as in real-world scenarios, attackers may use different command lines or obfuscation techniques to achieve the same goal. We showcase several pairs of generated similar command lines in Table. 9, demonstrating the efficacy of utilizing LLMs for this purpose.

3.2 Real-World Testing Set Collection

To create a testing set that can fairly and comprehensively evaluate different methods, better reflect real-world usage scenarios, and avoid training data leakage, we neither directly partitioned the train-

ing set for the testing set, nor did we use LLMs to generate entirely new command lines from the initial seeds as the training set collection pipeline. Instead, we employed Splunk Attack data (Splunk), a dataset curated from various attacks, as the source for our testing set. This allows us to evaluate various approaches in real-world scenarios, as it includes many malicious command lines corresponding to various MITRE ATT&CK (mitre) techniques, covering multiple distinct attack vectors.

Initially, we extracted 12,723 unique command lines from the Splunk Attack data. However, many command lines had similar meanings but differed slightly in appearance, leading to data duplication and potential evaluation inaccuracies. To address this concern, we generated explanations using ChatGPT and then converted these explanations into embeddings using GTE-Large (Li et al., 2023). We used these embeddings to remove command lines with semantically similar content, resulting in a final testing set of 2,807 command lines. For more details about the deduplication process, please refer to Appendix A. We then followed the similar command line synthesis step proposed in Sec. 3.1.3 to instruct GPT-4-turbo to generate the corresponding similar command line for each command line.

For each command line from the original set of 2,807, we used the explanation embeddings to identify the 1,000 least similar command lines from the remaining 2,806 as negative command lines. Additionally, we designated the generated similar command line as the positive command line. This method avoids evaluation inaccuracies by preventing the inclusion of semantically similar command lines among the negatives.

4 CmdCaliper: A Semantic-Aware Command-Line Embedding Model

Utilizing the proposed dataset CmdDataset, a command-line embedding model can be trained with a contrastive objective, like sentence embedding methodologies, as seen in (Gao et al., 2021; Chuang et al., 2022; Neelakantan et al., 2022). Given an embedding model, denoted as E , the procedure of embedding a command line c_i into an embedding vector e_i can be described as $e_i = E(c_i)$. In each training iteration, several similar pairs are randomly sampled from the training set to form a batch, which contains k similar command-line pairs $\{(x_i, x_i^+)\}_{i=1}^k$, where (x_i, x_i^+) represents the i^{th} similar command-line pair. Within the batch,

	Training Set	Testing Set
Num of command-line pair	28,520	2,807
Num of unique command line	55,909	5,576
Max. command-line length	3,464	7,502
Min. command-line length	3	2
Avg. command-line length	91.635	96.301
Std. of command-line length	60.794	196.675

Table 1: The statistic information of CmdDataset.

	Coverage Rate (%)
Windows Commands	73.52
Windows Common File Extensions	70.67

Table 2: Coverage rates of CmdDataset across all Windows Commands and Windows common file extensions.

the similar command line x_i^+ is regarded as a positive sample of sample x_i , thereby encouraging the associated embedding vectors to be closer within the feature space. Conversely, other samples $\{x_j^+, \forall j \in \{1, 2, \dots, k\} \setminus \{i\}\}$ are treated as in-batch negatives (Sohn, 2016; Neelakantan et al., 2022; Gao et al., 2021), encouraging the embedding vectors to be farther. The InfoNCE loss (van den Oord et al., 2018), denoted as \mathcal{L}_{info} , can be calculated as follows:

$$\mathcal{L}_{info} = - \sum_{i=1}^k \log \frac{\exp(\frac{E(x_i) \cdot E(x_i^+)}{\tau})}{\sum_{j=1}^k \exp(\frac{E(x_i) \cdot E(x_j^+)}{\tau})}, \quad (1)$$

where τ is a hyperparameter that $\tau \in \mathbb{R}^+$.

5 Evaluation on CmdDataset

5.1 Experimental Settings

In the entire CmdDataset synthesis pipeline, to enhance the diversity of the generated command lines, we follow the approach outlined in (Wang et al., 2023a) and set the temperature parameter to 1 for all LLMs to encourage more random outputs.

5.2 The Statistics of the Dataset

The statistical information for the training and testing sets of the CmdDataset is presented in Table 1. Thanks to the real-world sources of our testing set, the standard deviation of the testing data significantly differs from that of the training set, enabling a more generalized and accurate evaluation.

5.3 The Diversity of the Synthesized Command Line

In this experiment, we aim to assess the diversity of these synthesized command lines. We calculated

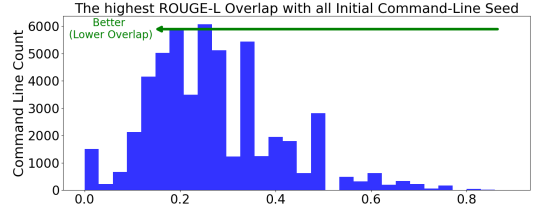


Figure 3: The distribution of the highest ROUGE-L overlap score between the generated command lines and the initial command-line seeds.

their coverage across all Windows commands² and common file name extensions in a clean Windows 10 virtual machine, as demonstrated in Table 2. To avoid bias, we excluded our manually formulated initial seeds and focused only on the command lines generated by LLMs. Overall, our synthesized command lines achieve a coverage rate of 73.52% out of 306 unique Windows commands and 70.67% out of 75 common file extensions. For more details about the coverage rate calculation process, please refer to Appendix C.

We conducted an in-depth analysis of the differences between the generated command lines and the initial command-line seeds, which served as a foundational starting point for constructing the command-line dataset. For each generated command line, we calculated the highest ROUGE-L overlap (Lin, 2004) which ranges from 0 to 1 among all command-line seeds. A higher ROUGE-L score indicates a greater overlap between the generated command lines and the initial command-line seeds. The distribution of ROUGE-L scores is illustrated in Fig. 3. These findings suggest that the command lines synthesized by our pipeline (Sec. 3.1) are not limited to minor tweaks of the original command-line seeds. On the contrary, they are capable of producing a broad range of command lines, some of which may exhibit significant differences from the initial seeds.

5.4 The Diversity within the Similar Command-Line Pairs

In this section, we examined the distribution of ROUGE-L overlap for each pair of similar command lines, as shown in Fig. 4. These metrics help determine whether similar command-line pairs are derived from minor modifications to arguments or entirely different commands achieving a similar objective. Notably, our findings reveal that most

²Windows Command-line reference A-Z

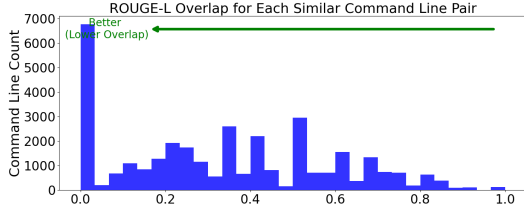


Figure 4: The distribution of ROUGE-L overlap score for all similar command-line pairs.

ROUGE-L scores are low, nearing zero, suggesting that the “similarity” in our command-line dataset is not solely based on lexical similarities but rather reflects genuine semantic similarities. This vital understanding enables us to train command-line embedding models from a semantic perspective and subsequently evaluate the performance of different command-line embedding models.

5.5 The Quality of the LLM’s Command-Line Explanations

In Sec. 3.2, we utilize ChatGPT (OpenAI, 2022) to generate explanations for command lines and employ GTE-Large (Li et al., 2023) to convert these explanations into embeddings for data processing. In this evaluation, our focus is on assessing the quality of the explanations generated by LLM.

We randomly selected 200 command lines and their explanations from our training set. An expert (collaborator of this work) with over three years of cybersecurity experience assessed the correctness of each explanation, providing scores as positive, neutral, or negative, while ignoring minor syntax errors. Normalizing the expert’s scores against the highest possible total, we obtained a score of **98.25%**, indicating that most command-line explanations accurately describe their purposes. Several good and bad examples are listed in Table 10.

5.6 The Quality of the Similar Command-Line Pairs

In this experiment, we investigate whether the command-line pairs generated by LLMs are truly similar in terms of semantics. Leveraging the findings presented in Sec. 5.5, which demonstrate a high alignment between the explanations generated by LLMs and those provided by human experts, we follow the same experimental settings, utilizing identical prompts and instructing ChatGPT (OpenAI, 2022) to generate explanations for each command line in the testing set.

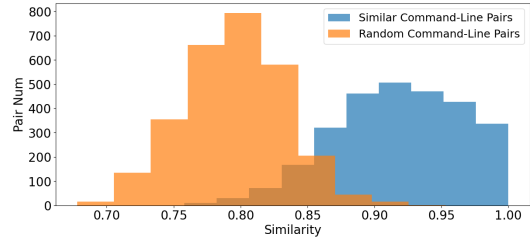


Figure 5: The histogram of the explanation similarity between random command-line pairs and similar command-line pairs in the testing set of CmdDataset.

LLMs for Data Synthesis	Coverage Rate (%)
GPT-3.5-Turbo (OpenAI, 2022)	48.75
Mixtral 8x7B (Jiang et al., 2024)	27.5
WizardLM-13B-v1.2 (Xu et al., 2024)	35
Gemini-1.0-Pro (Team et al., 2023)	51.25
Qwen1.5-14B-Chat (Bai et al., 2023)	23.75
Claude-3-Haiku (Anthropic, 2024)	30
LLM Pool (our)	70

Table 3: Explanation clusters coverage rates of the command lines synthesized by different LLMs.

After acquiring all explanations, we employ GTE-Large (Li et al., 2023) to embed all explanations into corresponding embedding vectors and calculate the similarity between explanations for each similar command-line pair. Subsequently, we randomly construct an equal number of command-line pairs, totaling 2,807 pairs, with both command lines in each pair randomly sampled from all command lines in the testing set.

Fig. 5 illustrates the similarity distributions of both similar and random command-line pairs. Notably, there is a significant gap in the explanation similarity between these two groups. Specifically, the average explanation similarity of each similar command-line pair exceeds that of 97.483% of the random command-line pairs, indicating the effectiveness of synthesizing similar command lines with LLMs. Note that the range of similarity for GTE-Large is approximately between 0.65 and 1.

5.7 Effectiveness of a Pool of LLMs

In this experiment, we aim to study whether a pool of LLMs can make the synthesized command lines more diverse versus utilizing a single LLM as we described in Sec. 3.1. First, we instruct each LLM in our pool to synthesize 7,500 command lines from the same initial seeds, totaling 52,500 command lines. Using findings from Sec. 5.5, we then generate high-quality explanations for all synthesized

command lines. Next, we embed these explanations into vectors using GTE-Large (Li et al., 2023) and cluster them with DBSCAN (Ester et al., 1996), using a maximum distance of 0.08, a minimum of 5 samples per cluster, and cosine similarity as the distance metric. This process resulted in 80 distinct clusters (excluding the noise cluster), each representing a specific purpose or intention based on similar command line explanations. We then computed the coverage rates of the LLM pool and each individual LLM across these clusters to assess the diversity of synthesized command lines.

The results are presented in Table 3. It is evident that command lines synthesized by the LLM pool cover the highest number of explanation clusters, reaching up to 70%. This highlights the capability of utilizing a pool of LLMs pre-trained on diverse training data to generate a broader range of command lines.

6 Evaluation on CmdCaliper

6.1 Experimental Settings

CmdCaliper was trained on three distinct model scales: small, base, and large, which are initialized from the GTE-small, -base, and -large (Li et al., 2023), respectively. For more details about the hyperparameters and training processes, please refer to Appendix E.

6.2 Compare with SOTAs

This section compares several SOTA sentence embedding methods using the testing set from Cmd-Dataset. We adopt Mean Reciprocal Ranking@K (MRR) and Top@K metrics for evaluating the performance of CmdCaliper, following the text search task methodology (Muennighoff et al., 2022). For more details about the two evaluation metrics, please refer to Appendix D. Both metrics yield scores between 0 and 100. The results of this comparative analysis are presented in Table 8.

The results indicate that CmdCaliper consistently achieves competitive performance with state-of-the-art (SOTA) models across all evaluation metrics and scales. Specifically, CmdCaliper-Base achieved an MRR@3 score of 87.56, surpassing our pretrained model - GTE-Base (Li et al., 2023) by 9.36 and outperforming all sentence embedding models of comparable size. On a larger scale, CmdCaliper-Large achieved an MRR@3 score of 89.12, surpassing GTE-Large by 4.86. Remarkably, even CmdCaliper-Small, with a mere 0.03B pa-

Methods	Params (B)	MRR @3	MRR @10	Top @3	Top @10
Levenshtein distance ¹	-	71.23	72.45	74.99	81.83
Word2Vec ²	-	45.83	46.93	48.49	54.86
E5 _S ³	0.03	81.59	82.6	84.97	90.59
GTE _S ⁴	0.03	82.35	83.28	85.39	90.84
CmdCaliper _S	0.03	86.81	87.78	89.21	94.76
BGE-en _B ⁵	0.11	79.49	80.41	82.33	87.39
E5 _B	0.11	83.16	84.07	86.14	91.56
GTR _B ⁶	0.11	81.55	82.51	84.54	90.1
GTE _B	0.11	78.2	79.07	81.22	86.14
CmdCaliper _B	0.11	87.56	88.47	90.27	95.26
BGE-en _L	0.34	84.11	84.92	86.64	91.09
E5 _L	0.34	84.12	85.04	87.32	92.59
GTR _L	0.34	88.09	88.68	91.27	94.58
GTE _L	0.34	84.26	85.03	87.14	91.41
CmdCaliper _L	0.34	89.12	89.91	91.45	95.65

¹(Haldar and Mukhopadhyay, 2011) ²(Mikolov et al., 2013) ³(Wang et al., 2022) ⁴(Li et al., 2023) ⁵(Xiao et al., 2023) ⁶(Ni et al., 2022)

Table 4: Comparison with the SOTAs for different pre-trained language models. Subscript *S*, *B*, and *L* denote the Small, Base, and Large versions respectively.

Models \ <i>r</i> (%)	20	40	60	80
GTR _{Base}	0.793	0.852	0.866	0.903
E5 _{Base}	0.796	0.859	0.87	0.899
GTE _{Base}	0.8	0.868	0.874	0.903
CmdCaliper _{Base}	0.869	0.906	0.927	0.939

Table 5: The AUC comparison for different embedding models and different sample rate *r*%.

rameters, is comparable with all SOTA embedding models at the large scale (with 0.335B parameters).

6.3 Semantic-Based Malicious Command-Line Detection

In this section, we approach malicious command-line detection as a retrieval task using the open-source atomic-red-team dataset (Canary), which includes command lines corresponding to 55 different MITRE ATT&CK techniques (mitre) (i.e., each technique describes different command line attack behaviors). We iteratively select *r*% of the malicious command lines from each technique as query command lines, while the remaining $1 - r$ % serve as positive command lines. Command lines from other techniques act as negative command lines. This process is repeated for each technique, and we calculate the average area under curve (AUC). The intuition behind this experiment is that a good embedding model should cluster malicious command lines from the same technique closer together,

Embedding Models	Params (B)	Acc. (%)
E5 _{Base}	0.11	93.86
GTE _{Base}	0.11	92.8
CmdCaliper _{Base}	0.11	96.37

Table 6: Accuracy comparison for command-line classification fine-tuned on fixed embedding models.

as they share similar attack behaviors. For more details about the experiment setup, please refer to Appendix F.

The detection results are illustrated in Table. 5. As observed, CmdCaliper-Base significantly outperforms all embedding models not fine-tuned on the command-line dataset. This difference is especially pronounced when the sample ratio is smaller. For instance, at a 20% sample ratio ($r = 20$), CmdCaliper-Base improves upon GTE-Base (Li et al., 2023) by approximately 0.069 in AUC. This suggests that when the query command-line set is smaller, the model requires a deeper understanding of the semantics of command lines.

6.4 Transfer to Command-Line Classification

Fine-tuning an additional module on a pre-trained embedding model for tasks like classification or regression often outperforms training from scratch. This is because well-trained embeddings capture rich, meaningful information that can be used across tasks. In this experiment, we trained a logistic regression classifier for Windows command classification using fixed command-line embeddings from different approaches.

Collecting a labeled command-line dataset poses significant challenges. To address this, we selected seven Windows commands: ‘find’, ‘robocopy’, ‘msiexec’, ‘rundll32’, ‘sc query’, ‘certutil’, and ‘print’ to synthesize 24,500 training and 24,500 testing command lines. For more details about the the classification dataset synthesis and the experimental setup, please refer to Appendix G.

The results of the classification are presented in Table. 6. As observed, CmdCaliper generally outperformed other sentence embedding models in terms of accuracy for the same model size. For example, CmdCaliper-Base achieved a 3.57% improvement over GTE-Base (Li et al., 2023). These findings highlight the importance of specialized embedding models for command-line data, allowing the models to encode more command-line information into their embedding vectors.

Model	Params (B)	MRR @3	MRR @10	Top @3	Top @10
Random Initialization	0.11	70.43	72.14	74.49	84.04
Bert _{Base}	0.11	82.25	83.38	85.79	92.13
GTE _{Base}	0.11	87.56	88.47	90.27	95.26

Table 7: Comparison of the performance of CmdCaliper fine-tuning from the different model.

6.5 Does Command-Line Embedding Benefit from Sentence Embedding?

We conducted experiments under three settings: training the Bert-Base (Devlin et al., 2019) network architecture, the pretrained model of GTE-Base (Li et al., 2023), with randomly initialized weights; fine-tuning from Bert-Base; and fine-tuning directly from the GTE-Base model. The results of the comparison are illustrated in Table 7. As observed, the performance of the pretrained model significantly influences the performance of the command-line embedding model. For instance, fine-tuning from the embedding model yields the highest MRR@3, showing a 5.31% improvement compared to direct fine-tuning with BERT.

We believe that the reason command-line embedding models benefit from a good sentence embedding model lies in the fact that, although command lines often have entirely different grammar and structure from general sentences, in many cases, we can still infer some partial meanings of the command lines from semantically meaningful words such as filenames, arguments, or folder names.

7 Conclusion

In this work, we introduce CmdDataset, a dataset of similar command-line pairs. The training set utilizes the impressive capabilities of an LLM pool for automated generation, while the testing set consists of real-world malicious command lines for realistic evaluation. We also present CmdCaliper, the first dedicated command-line embedding model. Our results show that CmdCaliper, specifically designed for command-line processing, outperforms existing sentence embedding methods in various command-line downstream tasks, such as command classification, malicious command line detection, and similar command-line retrieval.

We open-source the dataset, model weights, and all program codes, hoping this study sheds light on future command-line embedding research.

Methods	Hidden Dim	MRR @3	MRR @10	Top @3	Top @10
BGE-en _{Large} [1]	1024	84.11	84.92	86.64	91.09
E5 _{Large} [2]	1024	84.12	85.04	87.32	92.59
GTE _{Large} [3]	1024	84.26	85.03	87.14	91.41
OpenAI-text-embedding-3-small [4]	1536	84.35	85.28	87.78	92.8
OpenAI-text-embedding-3-large	3072	89.55	90.19	92.59	96.12
CyCraft-CmdCraft	1024	89.76	90.32	92.59	95.87

Table 8: Comparison with the SOTAs for different pre-trained language models. Subscript *S*, *B*, and *L* denote the Small, Base, and Large versions respectively.

8 Limitation

Despite the contributions made in this paper, several tasks listed below are worth exploring in the future.

- Support more command-line interpreters: Given their prominence in recent cybersecurity incidents, this study focuses on Windows and PowerShell commands. Compelling statistics (Kutscher, 2023) reveal that over two-thirds of script-based attacks leverage PowerShell. Nevertheless, an investigation into a unified command-line embedding model capable of spanning multiple command-line interpreters presents a promising avenue for future research.
- Resilience against command-line obfuscation: While CmdCaliper was trained on semantically similar command-line pairs, providing a certain degree of resilience against obfuscated command lines in this study, attackers often employ more sophisticated command-line obfuscation techniques to evade defense and detection mechanisms. This poses a significant challenge when relying solely on the embedding model for detecting malicious activities.
- Nested command line: We can obtain the corresponding semantic embedding vector by inputting a command line into CmdCaliper. However, a command line itself can be a composition of multiple command lines as well, making it difficult to accurately embed them into the feature space. Techniques such as few-shot learning (Brown et al., 2020) or instruction-finetuned text embeddings (Su et al., 2023) may provide potential solutions for generating command-line embeddings for specific downstream tasks. This area repre-

sents another possibility for future investigation.

References

- Anthropic. 2024. [Claude 3 haiku: our fastest model yet](#).
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Marcus Botacin. 2023. Gp threats-3: Is automatic malware generation a threat? In *IEEE Security and Privacy Workshops*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the 2020 Conference on Advances in Neural Information Processing Systems*.
- Red Canary. [Atomic red team](#).
- P. V. Sai Charan, Hrushikesh Chunduri, P. Mohan Anand, and Sandeep K Shukla. 2023. [From text to mitre techniques: Exploring the malicious use of large language models for generating cyber attack payloads](#).
- Efstratios Chatzoglou, Georgios Karopoulos, Georgios Kambourakis, and Zisis Tsiatsikas. 2023. [Bypassing antivirus detection: old-school malware, new tricks](#).

716	Yung-Sung Chuang, Rumen Dangovski, Hongyin Luo,	Sophia Yang, Szymon Antoniak, Teven Le Scao,	771
717	Yang Zhang, Shiyu Chang, Marin Soljagic, Shang-	Théophile Gervet, Thibaut Lavril, Thomas Wang,	772
718	Wen Li, Wen tau Yih, Yoon Kim, and James Glass.	Timothée Lacroix, and William El Sayed. 2024. <i>Mix-</i>	773
719	2022. DiffCSE: Difference-based contrastive learn-	<i>tral of experts.</i>	774
720	ing for sentence embeddings. In <i>Proceedings of the</i>		
721	<i>2022 Conference of the North American Chapter of</i>	Diederik Kingma and Jimmy Ba. 2015. Adam: A	775
722	<i>the Association for Computational Linguistics.</i>	method for stochastic optimization. In <i>Proceedings</i>	776
		<i>of the 3rd International Conference on Learning Rep-</i>	777
723	DARPA. 2019. <i>Transparent computing engagement 5</i>	<i>resentations.</i>	778
724	<i>data release.</i>		
		Jurgen Kutscher. 2023. <i>M-Trends 2023: Cy-</i>	779
725	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	<i>bersecurity insights from the frontlines.</i>	780
726	Kristina Toutanova. 2019. BERT: Pre-training of	https://www.mandiant.com/resources/blog/	781
727	deep bidirectional transformers for language under-	m-trends-2023.	782
728	standing. In <i>Proceedings of the 2019 Conference</i>		
729	<i>of the North American Chapter of the Association</i>	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	783
730	<i>for Computational Linguistics: Human Language</i>	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	784
731	<i>Technologies.</i>	rich Küttler, Mike Lewis, Wen tau Yih, Tim Rock-	785
		täschel, Sebastian Riedel, and Douwe Kiela. 2020.	786
732	Dinil Mon Divakaran and Sai Teja Peddinti. 2024. <i>Llms</i>	Retrieval-augmented generation for knowledge-	787
733	<i>for cyber security: New opportunities.</i>	intensive nlp tasks. In <i>Advances in Neural Infor-</i>	788
		<i>mation Processing Systems.</i>	789
734	Feng Dong, Liu Wang, Xu Nie, Fei Shao, Haoyu Wang,	Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long,	790
735	Ding Li, Xiapu Luo, and Xusheng Xiao. 2023. DIST-	Pengjun Xie, and Meishan Zhang. 2023. Towards	791
736	DET: A Cost-Effective distributed cyber threat de-	general text embeddings with multi-stage contrastive	792
737	tection system. In <i>Proceedings of the 32nd USENIX</i>	learning. <i>arXiv:2308.03281.</i>	793
738	<i>Security Symposium.</i>		
739	Martin Ester, Hans-Peter Kriegel, Jorg Sander, and Xi-	Chin-Yew Lin. 2004. ROUGE: A package for automatic	794
740	aowei Xu. 1996. A density-based algorithm for dis-	evaluation of summaries. In <i>The Proceedings of 2004</i>	795
741	covering clusters in large spatial databases with noise.	<i>Text Summarization Branches Out.</i>	796
742	In <i>KDD.</i>		
743	Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021.	Fucheng Liu, Yu Wen, Dongxue Zhang, Xihe Jiang,	797
744	SimCSE: Simple contrastive learning of sentence em-	Xinyu Xing, and Dan Meng. 2019. Log2vec: A	798
745	beddings. In <i>Proceedings of the 2021 Conference on</i>	heterogeneous graph embedding based approach for	799
746	<i>Empirical Methods in Natural Language Processing.</i>	detecting cyber threats within enterprise. In <i>Pro-</i>	800
		<i>ceedings of the 2019 ACM SIGSAC Conference on</i>	801
747	Andrew Golczynski and John A. Emanuello. 2021. End-	<i>Computer and Communications Security.</i>	802
748	to-end anomaly detection for identifying malicious		
749	cyber behavior through nlp-based log embeddings.	Forrest McKee and David Noever. 2023. <i>Chatbots in a</i>	803
750	In <i>Proceedings of the 1st International Workshop on</i>	<i>honeypot world.</i>	804
751	<i>Adaptive Cyber Defense.</i>		
		Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey	805
752	Haixuan Guo, Shuhan Yuan, and Xintao Wu. 2021. Log-	Dean. 2013. Efficient estimation of word represen-	806
753	bert: Log anomaly detection via bert. In <i>Proceedings</i>	tations in vector space. In <i>Proceedings of the 1st</i>	807
754	<i>of the 2021 International Joint Conference on Neural</i>	<i>International Conference on Learning Representa-</i>	808
755	<i>Networks.</i>	<i>tions, Workshop Track Proceedings.</i>	809
		mitre. <i>Mitre att&ck.</i>	810
756	Rishin Haldar and Debajyoti Mukhopadhyay. 2011.		
757	<i>Levenshtein distance technique in dictionary lookup</i>	Stephen Moskal, Sam Laney, Erik Hemberg, and Una-	811
758	<i>methods: An improved approach.</i>	May O'Reilly. 2023. <i>Llms killed the script kiddie:</i>	812
		<i>How agents supported by large language models</i>	813
759	Or Honovich, Thomas Scialom, Omer Levy, and Timo	<i>change the landscape of network threat testing.</i>	814
760	Schick. 2023. Unnatural instructions: Tuning lan-		
761	guage models with (almost) no human labor. In <i>Pro-</i>	Farzad Nourmohammadzadeh Motlagh, Mehrdad Ha-	815
762	<i>ceedings of the 61st Annual Meeting of the Associa-</i>	jizadeh, Mehryar Majd, Pejman Najafi, Feng Cheng,	816
763	<i>tion for Computational Linguistics.</i>	and Christoph Meinel. 2024. <i>Large language models</i>	817
		<i>in cybersecurity: State-of-the-art.</i>	818
764	Albert Q. Jiang, Alexandre Sablayrolles, Antoine		
765	Roux, Arthur Mensch, Blanche Savary, Chris	Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and	819
766	Bamford, Devendra Singh Chaplot, Diego de las	Nils Reimers. 2022. <i>Mteb: Massive text embedding</i>	820
767	Casas, Emma Bou Hanna, Florian Bressand, Gi-	<i>benchmark. arXiv preprint arXiv:2210.07316.</i>	821
768	anna Lengyel, Guillaume Bour, Guillaume Lam-		
769	ple, Léo Renard Lavaud, Lucile Saulnier, Marie-		
770	Anne Lachaux, Pierre Stock, Sandeep Subramanian,		

822	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. 2022. Text and code embeddings by contrastive pre-training. ArXiv:2201.10005.	
823		
824		
825		
826		
827		
828		
829		
830		
831		
832	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith Hall, Ming-Wei Chang, and Yinfei Yang. 2022. Large dual encoders are generalizable retrievers. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> .	
833		
834		
835		
836		
837		
838	Talha Ongun, Jack W. Stokes, Jonathan Bar Or, Ke Tian, Farid Tajaddodianfar, Joshua Neil, Christian Seifert, Alina Oprea, and John C. Platt. 2021. Living-off-the-land command detection using active learning. In <i>Proceedings of the 24th International Symposium on Research in Attacks, Intrusions and Defenses</i> , page 442–455.	
839		
840		
841		
842		
843		
844		
845	OpenAI. 2022. Introducing chatgpt .	
846	OpenAI. 2023. New models and developer products announced at devday .	
847		
848	Yin Minn Pa Pa, Shunsuke Tanizaki, Tetsui Kou, Michel Van Eetenand Katsunari Yoshioka, and Tsutomu Matsumoto. 2023. An attacker’s dream? exploring the capabilities of chatgpt for developing malware. In <i>Proceedings of the 16th Cyber Security Experimentation and Test Workshop</i> .	
849		
850		
851		
852		
853		
854	Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. Gorilla: Large language model connected with massive apis. <i>arXiv preprint arXiv:2305.15334</i> .	
855		
856		
857		
858	Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In <i>Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing</i> .	
859		
860		
861		
862		
863	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> .	
864		
865		
866		
867		
868	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code llama: Open foundation models for code .	
869		
870		
871		
872		
873		
874		
875		
876		
	Kihyuk Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective. In <i>Proceedings of the 2016 Conference on Advances in Neural Information Processing Systems</i> .	877 878 879 880
	Splunk. Splunk attack data repository .	881
	Hongjin Su, Weijia Shi, Jungo Kasai, Yizhong Wang, Yushi Hu, Mari Ostendorf, Wen tau Yih, Noah A. Smith, Luke Zettlemoyer, and Tao Yu. 2023. One embedder, any task: Instruction-finetuned text embeddings. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> .	882 883 884 885 886 887
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	888 889 890 891 892
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Gura, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Meray, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi, Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam Błoniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette, Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin Chadwick, Ilya Kornakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan, Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier Garcia, Thanumalayan Sankaranarayanan Pillai, Jacob Devlin, Michael Laskin, Diego	893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937

938	de Las Casas, Dasha Valter, Connie Tao, Lorenzo	1001
939	Blanco, Adrià Puigdomènech Badia, David Reitter,	1002
940	Mianna Chen, Jenny Brennan, Clara Rivera, Sergey	1003
941	Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,	1004
942	Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yim-	1005
943	ing Gu, Kate Olszewska, Ravi Addanki, Antoine	1006
944	Miech, Annie Louis, Denis Teplyashin, Geoff Brown,	1007
945	Elliot Catt, Jan Balaguer, Jackie Xiang, Pidong Wang,	1008
946	Zoe Ashwood, Anton Briukhov, Albert Webson, San-	1009
947	jay Ganapathy, Smit Sanghavi, Ajay Kannan, Ming-	1010
948	Wei Chang, Axel Stjerngren, Josip Djolonga, Yut-	1011
949	ing Sun, Ankur Bapna, Matthew Aitchison, Pedram	1012
950	Pejman, Henryk Michalewski, Tianhe Yu, Cindy	1013
951	Wang, Juliette Love, Junwhan Ahn, Dawn Bloxwich,	1014
952	Kehang Han, Peter Humphreys, Thibault Sellam,	1015
953	James Bradbury, Varun Godbole, Sina Samangoei,	1016
954	Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.	1017
955	Arnold, Vijay Vasudevan, Shubham Agrawal, Jason	1018
956	Riesa, Dmitry Lepikhin, Richard Tanburn, Srivat-	1019
957	san Srinivasan, Hyeontaek Lim, Sarah Hodgkinson,	1020
958	Pranav Shyam, Johan Ferret, Steven Hand, Ankush	1021
959	Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Gi-	1022
960	ang, Alexander Neitz, Zaheer Abbas, Sarah York,	1023
961	Machel Reid, Elizabeth Cole, Aakanksha Chowdh-	1024
962	ery, Dipanjan Das, Dominika Rogozińska, Vitaliy	1025
963	Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas	1026
964	Zilka, Flavien Prost, Luheng He, Marianne Mon-	1027
965	teiro, Gaurav Mishra, Chris Welty, Josh Newlan,	1028
966	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	1029
967	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	1030
968	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	1031
969	Anirudh Baddepudi, Alex Goldin, Adnan Ozturk,	1032
970	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	1033
971	dra Sachan, Reinald Kim Amplayo, Craig Swans-	1034
972	on, Dessie Petrova, Shashi Narayan, Arthur Guez,	1035
973	Siddhartha Brahma, Jessica Landon, Miteyan Pa-	1036
974	tel, Ruizhe Zhao, Kevin Villela, Luyu Wang, Wen-	1037
975	hao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	1038
976	James Keeling, Petko Georgiev, Diana Mincu, Boxi	1039
977	Wu, Salem Haykal, Rachel Saputro, Kiran Vodra-	1040
978	halli, James Qin, Zeynep Cankara, Abhanshu Sharma,	1041
979	Nick Fernando, Will Hawkins, Behnam Neyshabur,	1042
980	Solomon Kim, Adrian Hutter, Priyanka Agrawal,	1043
981	Alex Castro-Ros, George van den Driessche, Tao	1044
982	Wang, Fan Yang, Shuo yiin Chang, Paul Komarek,	1045
983	Ross McIlroy, Mario Lučić, Guodong Zhang, Wael	1046
984	Farhan, Michael Sharman, Paul Natsev, Paul Michel,	1047
985	Yamini Bansal, Siyuan Qiao, Kris Cao, Siamak Shak-	1048
986	eri, Christina Butterfield, Justin Chung, Paul Kishan	1049
987	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	1050
988	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	1051
989	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	1052
990	Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	1053
991	Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	1054
992	Yash Katariya, Sebastian Riedel, Paige Bailey, Kefan	1055
993	Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone,	1056
994	Neil Houlsby, Xuehan Xiong, Zhen Yang, Elena Gri-	1057
995	bovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music	1058
996	Li, Thais Kagohara, Jay Pavagadhi, Sophie Bridgers,	1059
997	Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed,	1060
998	Tianqi Liu, Richard Powell, Vijay Bolina, Mariko	1061
999	Iinuma, Polina Zablotskaia, James Besley, Da-Woon	1062
1000	Chung, Timothy Dozat, Ramona Comanescu, Xi-	1063
	ance Si, Jeremy Greer, Guolong Su, Martin Polacek,	
	Raphaël Lopez Kaufman, Simon Tokumine, Hexiang	
	Hu, Elena Buchatskaya, Yingjie Miao, Mohamed	
	Elhawaty, Aditya Siddhant, Nenad Tomasev, Jin-	
	wei Xing, Christina Greer, Helen Miller, Shereen	
	Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Ange-	
	los Filos, Milos Besta, Rory Blevins, Ted Klimenko,	
	Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Os-	
	car Chang, Mantas Pajarskas, Carrie Muir, Vered	
	Cohen, Charline Le Lan, Krishna Haridasan, Amit	
	Marathe, Steven Hansen, Sholto Douglas, Rajku-	
	mar Samuel, Mingqiu Wang, Sophia Austin, Chang	
	Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso Lorenzo,	
	Lars Lowe Sjösund, Sébastien Cevey, Zach Gle-	
	icher, Thi Avrahami, Anudhyan Boral, Hansa Srini-	
	vasan, Vittorio Selo, Rhys May, Konstantinos Aiso-	
	pos, Léonard Hussenot, Livio Baldini Soares, Kate	
	Baumli, Michael B. Chang, Adrià Recasens, Ben	
	Caine, Alexander Pritzel, Filip Pavetic, Fabio Pardo,	
	Anita Gergely, Justin Frye, Vinay Ramasesh, Dan	
	Horgan, Kartikeya Badola, Nora Kassner, Subhra-	
	jit Roy, Ethan Dyer, Víctor Campos Campos, Alex	
	Tomala, Yunhao Tang, Dalia El Badawy, Elspeth	
	White, Basil Mustafa, Oran Lang, Abhishek Jin-	
	dal, Sharad Vikram, Zhitao Gong, Sergi Caelles,	
	Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng,	
	Wojciech Stokowiec, Ce Zheng, Phoebe Thacker,	
	Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh,	
	James Svensson, Max Bileschi, Piyush Patil, Ankesh	
	Anand, Roman Ring, Katerina Tsihlias, Arpi Vezer,	
	Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom	
	Kwiatkowski, Samira Daruki, Keran Rong, Allan	
	Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg,	
	Mina Khan, Lisa Anne Hendricks, Marie Pellat,	
	Vladimir Feinberg, James Cobon-Kerr, Tara Sainath,	
	Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives,	
	Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd,	
	Le Hou, Qingze Wang, Thibault Sottiaux, Michela	
	Paganini, Jean-Baptiste Lespiau, Alexandre Mou-	
	farek, Samer Hassan, Kaushik Shivakumar, Joost van	
	Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh	
	Goyal, Matthew Tung, Andrew Brock, Hannah Shea-	
	han, Vedant Misra, Cheng Li, Nemanja Rakićević,	
	Mostafa Dehghani, Fangyu Liu, Sid Mittal, Jun-	
	hyuk Oh, Seb Noury, Eren Sezener, Fantine Huot,	
	Matthew Lamm, Nicola De Cao, Charlie Chen, Sid-	
	harth Mudgal, Romina Stella, Kevin Brooks, Gau-	
	tam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita	
	Melinkeri, Aaron Cohen, Venus Wang, Kristie Sey-	
	more, Sergey Zubkov, Rahul Goel, Summer Yue,	
	Sai Krishnakumaran, Brian Albert, Nate Hurley,	
	Motoki Sano, Anhad Mohananey, Jonah Joughin,	
	Egor Filonov, Tomasz Kępa, Yomna Eldawy, Jiaw-	
	ern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor	
	Bos, Jerry Chang, Sanil Jain, Sri Gayatri Sundara	
	Padmanabhan, Subha Puttagunta, Kalpesh Krishna,	
	Leslie Baker, Norbert Kalb, Vamsi Bedapudi, Adam	
	Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin,	
	Xiang Zhou, Zhichun Wu, Sam Sobell, Andrea Si-	
	ciliano, Alan Papir, Robby Neale, Jonas Bragagnolo,	
	Tej Toor, Tina Chen, Valentin Anklin, Feiran Wang,	
	Richie Feng, Milad Gholami, Kevin Ling, Lijuan	
	Liu, Jules Walter, Hamid Moghaddam, Arun Kishore,	

1064	Jakub Adamek, Tyler Mercado, Jonathan Mallinson,	ton, Wojciech Rzadkowski, Fiona Macintosh, Kon-	1127
1065	Siddhinita Wandekar, Stephen Cagle, Eran Ofek,	stantin Shagin, Paul Medina, Chen Liang, Jinjing	1128
1066	Guillermo Garrido, Clemens Lombriser, Maksim	Zhou, Pararth Shah, Yingying Bi, Attila Dankovics,	1129
1067	Mukha, Botu Sun, Hafeezul Rahman Mohammad,	Shipra Banga, Sabine Lehmann, Marissa Bredesen,	1130
1068	Josip Matak, Yadi Qian, Vikas Peswani, Pawel Janus,	Zifan Lin, John Eric Hoffmann, Jonathan Lai, Ray-	1131
1069	Quan Yuan, Leif Schelin, Oana David, Ankur Garg,	nald Chung, Kai Yang, Nihal Balani, Arthur Bražin-	1132
1070	Yifan He, Oleksii Duzhyi, Anton Ålgmyr, Timo-	skas, Andrei Sozanschi, Matthew Hayes, Héctor Fer-	1133
1071	thée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex	nández Alcalde, Peter Makarov, Will Chen, Anto-	1134
1072	Chinien, Rakesh Shivanna, Aleksandr Chuklin, Josie	nio Stella, Liselotte Snijders, Michael Mandl, Ante	1135
1073	Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed,	Kärman, Pawel Nowak, Xinyi Wu, Alex Dyck, Kr-	1136
1074	Subhadrata Das, Zihang Dai, Kyle He, Daniel von	ishnan Vaidyanathan, Raghavender R, Jessica Mal-	1137
1075	Dincklage, Shyam Upadhyay, Akanksha Maurya,	let, Mitch Rudominer, Eric Johnston, Sushil Mit-	1138
1076	Luyan Chi, Sebastian Krause, Khalid Salama, Pam G	tal, Akhil Udathu, Janara Christensen, Vishal Verma,	1139
1077	Rabinovitch, Pavan Kumar Reddy M, Aarush Sel-	Zach Irving, Andreas Santucci, Gamaleldin Elsayed,	1140
1078	van, Mikhail Dektiarev, Golnaz Ghiasi, Erdem Gu-	Elnaz Davoodi, Marin Georgiev, Ian Tenney, Nan	1141
1079	ven, Himanshu Gupta, Boyi Liu, Deepak Sharma,	Hua, Geoffrey Cideron, Edouard Leurent, Mah-	1142
1080	Idan Heimlich Shtacher, Shachi Paul, Oscar Aker-	moud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy	1143
1081	lund, François-Xavier Aubet, Terry Huang, Chen	Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper	1144
1082	Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze,	Snoek, Mukund Sundararajan, Xuezhi Wang, Zack	1145
1083	Francesco Bertolini, Liana-Eleonora Marinescu, Mar-	Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar,	1146
1084	tin Bülle, Dominik Paulus, Khyatti Gupta, Tejasi	Lara Tumeh, Eyal Ben-David, Rishub Jain, Jonathan	1147
1085	Latkar, Max Chang, Jason Sanders, Roopa Wil-	Uesato, Romina Datta, Oskar Bunyan, Shimu Wu,	1148
1086	son, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet,	John Zhang, Piotr Stanczyk, Ye Zhang, David Steiner,	1149
1087	Tulsee Doshi, Sid Lall, Swaroop Mishra, Wanming	Subhajit Naskar, Michael Azzam, Matthew Johnson,	1150
1088	Chen, Thang Luong, Seth Benjamin, Jasmine Lee,	Adam Paszke, Chung-Cheng Chiu, Jaume Sanchez	1151
1089	Ewa Andrejczuk, Dominik Rabiej, Vipul Ranjan,	Elias, Afroz Mohiuddin, Faizan Muhammad, Jin	1152
1090	Krzysztof Styrac, Pengcheng Yin, Jon Simon, Mal-	Miao, Andrew Lee, Nino Vieillard, Jane Park, Ji-	1153
1091	colm Rose Harriott, Mudit Bansal, Alexei Robsky,	ageng Zhang, Jeff Stanway, Drew Garmon, Abhijit	1154
1092	Geoff Bacon, David Greene, Daniil Mirylenka, Chen	Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Lu-	1155
1093	Zhou, Obaid Sarvana, Abhimanyu Goyal, Samuel	owei Zhou, Jonathan Evens, William Isaac, Geoffrey	1156
1094	Andermatt, Patrick Siegler, Ben Horn, Assaf Is-	Irving, Edward Loper, Michael Fink, Isha Arkatkar,	1157
1095	rael, Francesco Pongetti, Chih-Wei "Louis" Chen,	Nanxin Chen, Izhak Shafran, Ivan Petrychenko,	1158
1096	Marco Selvatici, Pedro Silva, Kathie Wang, Jack-	Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai	1159
1097	son Tolins, Kelvin Guu, Roey Yogeve, Xiaochen Cai,	Zhu, Peter Grabowski, Yu Mao, Alberto Magni,	1160
1098	Alessandro Agostini, Maulik Shah, Hung Nguyen,	Kaisheng Yao, Javier Snaider, Norman Casagrande,	1161
1099	Noah Ó Donnaile, Sébastien Pereira, Linda Friso,	Evan Palmer, Paul Suganthan, Alfonso Castaño,	1162
1100	Adam Stambler, Adam Kurzrok, Chenkai Kuang,	Irene Giannoumis, Wooyeol Kim, Mikolaj Rybiński,	1163
1101	Yan Romanikhin, Mark Geller, ZJ Yan, Kane Jang,	Ashwin Sreevatsa, Jennifer Prendki, David Soergel,	1164
1102	Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qi-	Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari,	1165
1103	jun Tan, Dan Banica, Daniel Balle, Ryan Pham,	Meenu Gaba, Jeremy Wiesner, Diana Gage Wright,	1166
1104	Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot	Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay	1167
1105	Singh, Chris Hidey, Niharika Ahuja, Pranab Sax-	Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu,	1168
1106	ena, Dan Dooley, Srividya Pranavi Potharaju, Eileen	Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert	1169
1107	O'Neill, Anand Gokulchandran, Ryan Foley, Kai	Cui, Tian LIN, Marcus Wu, Ricardo Aguilar, Keith	1170
1108	Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta,	Pallo, Abhishek Chakladar, Ginger Perng, Elena Al-	1171
1109	Ragha Kotikalapudi, Chalence Safranek-Shrader, An-	lica Abellan, Mingyang Zhang, Ishita Dasgupta,	1172
1110	drew Goodman, Joshua Kessinger, Eran Globen, Pra-	Nate Kushman, Ivo Penchev, Alena Repina, Xihui	1173
1111	teek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang	Wu, Tom van der Weide, Priya Ponnappalli, Car-	1174
1112	Song, Ali Eichenbaum, Thomas Brovelli, Sahitya	oline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier	1175
1113	Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani,	Dousse, Fan Yang, Jeff Piper, Nathan Ie, Rama Pa-	1176
1114	Charles Chen, Andy Crawford, Shalini Pal, Mukund	samarthi, Nathan Lintz, Anitha Vijayakumar, Daniel	1177
1115	Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski,	Andor, Pedro Valenzuela, Minnie Lui, Cosmin Padu-	1178
1116	Pierre-Louis Cedoz, Chenmei Li, Shiyuan Chen,	raru, Daiyi Peng, Katherine Lee, Shuyuan Zhang,	1179
1117	Niccolò Dal Santo, Siddharth Goyal, Jitesh Pun-	Somer Greene, Duc Dung Nguyen, Paula Kurylow-	1180
1118	jabi, Karthik Kappaganthu, Chester Kwak, Pallavi	icz, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam	1181
1119	LV, Sarmishta Velury, Himadri Choudhury, Jamie	Choo, Ziqiang Feng, Biao Zhang, Achintya Sing-	1182
1120	Hall, Premal Shah, Ricardo Figueira, Matt Thomas,	hal, Dayou Du, Dan McKinnon, Natasha Antropova,	1183
1121	Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Ju-	Tolga Bolukbasi, Orgad Keller, David Reid, Daniel	1184
1122	rdi, Sharat Chikkerur, Yenai Ma, Adams Yu, Soo	Finchelstein, Maria Abi Raad, Remi Crocker, Pe-	1185
1123	Kwak, Victor Áhdel, Sujeevan Rajayogam, Travis	ter Hawkins, Robert Dadashi, Colin Gaffney, Ken	1186
1124	Choma, Fei Liu, Aditya Barua, Colin Ji, Ji Ho	Franko, Anna Bulanova, Rémi Leblond, Shirley	1187
1125	Park, Vincent Hellendoorn, Alex Bailey, Taylan Bi-	Chung, Harry Askham, Luis C. Cobo, Kelvin Xu,	1188
1126	lal, Huanjie Zhou, Mehrdad Khatir, Charles Sut-	Felix Fischer, Jun Xu, Christina Sorokin, Chris Al-	1189

1190	berti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev,	ish Reddy Vuyyuru, John Aslanides, Nidhi Vyas,	1253
1191	Hannah Forbes, Dylan Banarse, Zora Tung, Mark	Martin Wicke, Xiao Ma, Evgenii Eltyshv, Nina Mar-	1254
1192	Omernick, Colton Bishop, Rachel Sterneck, Rohan	tin, Hardie Cate, James Manyika, Keyvan Amiri,	1255
1193	Jain, Jiawei Xia, Ehsan Amid, Francesco Piccinno,	Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier,	1256
1194	Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz,	Nilesh Tripuraneni, David Madras, Mandy Guo,	1257
1195	Alex Polozov, Victoria Krakovna, Sasha Brown, Mo-	Austin Waters, Oliver Wang, Joshua Ainslie, Jason	1258
1196	hammadHossein Bateni, Dennis Duan, Vlad Firoiu,	Baldrige, Han Zhang, Garima Pruthi, Jakob Bauer,	1259
1197	Meghana Thotakuri, Tom Natan, Matthieu Geist,	Feng Yang, Riham Mansour, Jason Gelman, Yang Xu,	1260
1198	Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko	George Polovets, Ji Liu, Honglong Cai, Warren Chen,	1261
1199	Tojo, Michael Kwong, James Lee-Thorp, Christo-	XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof	1262
1200	pher Yew, Danila Sinopalnikov, Sabela Ramos, John	Angermueller, Xiaowei Li, Anoop Sinha, Weiren	1263
1201	Mellor, Abhishek Sharma, Kathy Wu, David Miller,	Wang, Julia Wiesinger, Emmanouil Koukoumidis,	1264
1202	Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jen-	Yuan Tian, Anand Iyer, Madhu Gurusurthy, Mark	1265
1203	nifer Beattie, Emily Caveness, Libin Bai, Julian	Goldenson, Parashar Shah, MK Blake, Hongkun Yu,	1266
1204	Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi	Anthony Urbanowicz, Jennimaria Palomaki, Chrisan-	1267
1205	Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng,	tha Fernando, Ken Durden, Harsh Mehta, Nikola	1268
1206	Frederick Liu, Fan Yang, Rui Zhu, Tian Huey Teh,	Momchev, Elahe Rahimtoroghi, Maria Georgaki,	1269
1207	Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin,	Amit Raul, Sebastian Ruder, Morgan Redshaw, Jin-	1270
1208	Daniel Toyama, Evan Rosen, Sasan Tavakkol, Lint-	hyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li,	1271
1209	ing Xue, Chen Elkind, Oliver Woodman, John Car-	Blake Hechtman, Parker Schuh, Milad Nasr, Kieran	1272
1210	penter, George Papamakarios, Rupert Kemp, Sushant	Milan, Vladimir Mikulik, Juliana Franco, Tim Green,	1273
1211	Kafle, Tanya Grunina, Rishika Sinha, Alice Tal-	Nam Nguyen, Joe Kelley, Aroma Mahendru, Andrea	1274
1212	bert, Diane Wu, Denese Owusu-Afriyie, Cosmo	Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshi-	1275
1213	Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna	tij Bansal, Vikram Rao, Rakesh Ghiya, Emma Wang,	1276
1214	Narayana, Jing Li, Saaber Fatehi, John Wieting,	Ke Ye, Jean Michel Sarr, Melanie Moranski Preston,	1277
1215	Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura	Madeleine Elish, Steve Li, Aakash Kaku, Jigar Gupta,	1278
1216	Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi	Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi	1279
1217	Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Re-	M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric	1280
1218	beca Santamaria-Fernandez, Sonam Goenka, Wenny	Chu, Xuanyi Dong, Amruta Muthal, Senaka Buth-	1281
1219	Yustalim, Robin Strudel, Ali Elqursh, Charlie Deck,	pitiya, Sarthak Jauhari, Nan Hua, Urvashi Khan-	1282
1220	Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoff-	delwal, Ayal Hitron, Jie Ren, Larissa Rinaldi, Sha-	1283
1221	mann, Dan Holtmann-Rice, Olivier Bachem, Sho	har Drath, Avigail Dabush, Nan-Jiang Jiang, Har-	1284
1222	Arora, Christy Koh, Soheil Hassas Yeganeh, Siim	shal Godhia, Uli Sachs, Anthony Chen, Yicheng	1285
1223	Pöder, Mukarram Tariq, Yanhua Sun, Lucian Ionita,	Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai,	1286
1224	Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, An-	James Wang, Chen Liang, Jenny Hamer, Chun-Sung	1287
1225	mol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz,	Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít	1288
1226	Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown,	Listík, Mathias Carlen, Jan van de Kerkhof, Marcin	1289
1227	Shreya Singh, Wei Fan, Aaron Parisi, Joe Stan-	Pikus, Krunoslav Zaher, Paul Müller, Sasha Zykova,	1290
1228	ton, Vinod Koverkathu, Christopher A. Choquette-	Richard Stefanec, Vitaly Gatsko, Christoph Hirn-	1291
1229	Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash	schall, Ashwin Sethi, Xingyu Federico Xu, Chetan	1292
1230	Shroff, Mani Varadarajan, Sanaz Bahargam, Rob	Ahuja, Beth Tsai, Anca Stefanoiu, Bo Feng, Kes-	1293
1231	Willoughby, David Gaddy, Guillaume Desjardins,	shav Dhandhanian, Manish Katyral, Akshay Gupta,	1294
1232	Marco Cornero, Brona Robenek, Bhavishya Mit-	Atharva Parulekar, Divya Pitta, Jing Zhao, Vivaan	1295
1233	tal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev,	Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin	1296
1234	Henrik Jacobsson, Alireza Ghaffarkhah, Morgane	Li, Peter Danenberg, Dennis Tu, Alex Pine, Vera	1297
1235	Rivière, Alanna Walton, Clément Crepy, Alicia Par-	Filippova, Abhipso Ghosh, Ben Limonchik, Bhar-	1298
1236	rish, Zongwei Zhou, Clement Farabet, Carey Rade-	gava Urala, Chaitanya Krishna Lanka, Derik Clive,	1299
1237	baugh, Praveen Srinivasan, Claudia van der Salm,	Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak,	1300
1238	Andreas Fidjeland, Salvatore Scellato, Eri Latorre-	Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal	1301
1239	Chimoto, Hanna Klimczak-Plucińska, David Bridson,	Majmundar, Michael Alverson, Michael Kucharski,	1302
1240	Dario de Cesare, Tom Hudson, Piermaria Mendolic-	Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo	1303
1241	chio, Lexi Walker, Alex Morris, Matthew Mauger,	Pelagatti, Rohan Kohli, Saurabh Kumar, Joseph Kim,	1304
1242	Alexey Guseynov, Alison Reid, Seth Odoom, Luc-	Swetha Sankar, Vineet Shah, Lakshmi Ramachan-	1305
1243	cia Loher, Victor Cotruta, Madhavi Yenugula, Do-	druni, Xiangkai Zeng, Ben Bariach, Laura Weidinger,	1306
1244	minik Grewe, Anastasia Petrushkina, Tom Duerig,	Amar Subramanya, Sissie Hsiao, Demis Hassabis,	1307
1245	Antonio Sanchez, Steve Yadlowsky, Amy Shen,	Koray Kavukcuoglu, Adam Sadovsky, Quoc Le,	1308
1246	Amir Globerson, Lynette Webb, Sahil Dua, Dong	Trevor Strohman, Yonghui Wu, Slav Petrov, Jeffrey	1309
1247	Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi,	Dean, and Oriol Vinyals. 2023. Gemini: A family of	1310
1248	Ananth Agarwal, Tomer Shani, Matan Eyal, Anuj	highly capable multimodal models .	1311
1249	Khare, Shreyas Rammohan Belle, Lei Wang, Chetan		
1250	Tekur, Mihir Sanjay Kale, Jinliang Wei, Ruoxin	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1312
1251	Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1313
1252	Zhao, Stephan Lee, Pandu Nayak, Doug Fritz, Man-	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1314
		Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	1315

1316	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,
1317	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,
1318	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-
1319	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan
1320	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,
1321	Isabel Kloumann, Artem Korenev, Punit Singh Koura,
1322	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-
1323	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-
1324	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-
1325	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-
1326	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,
1327	Ruan Silva, Eric Michael Smith, Ranjan Subrama-
1328	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-
1329	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,
1330	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,
1331	Melanie Kambadur, Sharan Narang, Aurelien Ro-
1332	driguez, Robert Stojnic, Sergey Edunov, and Thomas
1333	Scialom. 2023. LLaMA 2: Open foundation and
1334	fine-tuned chat models .
1335	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018.
1336	Representation learning with contrastive predictive
1337	coding .
1338	Laurens van der Maaten and Geoffrey Hinton. 2008.
1339	Visualizing data using t-sne . <i>Journal of Machine</i>
1340	<i>Learning Research</i> .
1341	Liang Wang, Nan Yang, Xiaolong Huang, Binxing
1342	Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder,
1343	and Furu Wei. 2022. Text embeddings by weakly-
1344	supervised contrastive pre-training. <i>arXiv preprint</i>
1345	<i>arXiv:2212.03533</i> .
1346	Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang,
1347	Rangan Majumder, and Furu Wei. 2023a. Improving
1348	text embeddings with large language models. <i>arXiv</i>
1349	<i>preprint arXiv:2401.00368</i> .
1350	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa
1351	Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh
1352	Hajishirzi. 2023b. Self-instruct: Aligning language
1353	models with self-generated instructions. In <i>Proceed-</i>
1354	<i>ings of the 61st Annual Meeting of the Association</i>
1355	<i>for Computational Linguistics</i> .
1356	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas
1357	Muennighoff. 2023. C-pack: Packaged resources
1358	to advance general chinese embedding .
1359	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,
1360	Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei
1361	Lin, and Daxin Jiang. 2024. WizardLM: Empow-
1362	ering large pre-trained language models to follow
1363	complex instructions . In <i>Proceedings of the 12th In-</i>
1364	<i>ternational Conference on Learning Representations</i> .
1365	Jiali Zeng, Yongjing Yin, Yufan Jiang, Shuangzhi Wu,
1366	and Yunbo Cao. 2022. Contrastive learning with
1367	prompt-derived virtual semantic prototypes for un-
1368	supervised sentence embedding. In <i>Findings of the</i>
1369	<i>Association for Computational Linguistics: EMNLP</i>
1370	<i>2022</i> .

A Testing Set Collection Detail 1371

To deduplicate, we utilized ChatGPT (OpenAI, 2022) to transform all command lines into concise descriptions that encapsulate the purpose and intention. Subsequently, we transformed these brief descriptions into embeddings using GTE-Large (Li et al., 2023), which achieved SOTA performance on the MTEB leaderboard (Muennighoff et al., 2022) among models of similar size. We then applied DBSCAN (Ester et al., 1996) for clustering the embeddings. Through this approach, each cluster contains command lines with highly similar semantics based on their explanations. Finally, we extracted two command lines from each cluster, resulting in a testing set comprising 2,807 command lines in total. 1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386

B Initial-Seed Collection Detail 1387

To gather high-quality initial seeds, we first extracted all command lines executed in DARPA Transparent Computing (DARPA, 2019), totaling 142,886 unique command lines. We then applied a heuristic filtering process to eliminate command lines that are semantically similar and differ only slightly, such as variations in log file suffixes. Ultimately, we curated 722 command lines as part of the initial seeds. 1388
1389
1390
1391
1392
1393
1394
1395
1396

To further extend the initial seeds, we formulated 796 command lines based on the descriptions and corresponding syntax found in Windows Commands³. Additionally, we parsed all example command lines from SS64⁴, totaling 497 command lines, and collected an additional 46 command lines from GitHub. Together, these contributions amounted to 2,061 high-quality and diverse command lines, which were integrated to form our initial command-line seeds. 1397
1398
1399
1400
1401
1402
1403
1404
1405
1406

C Windows-Commands Coverage Rate Detail 1407 1408

While many Windows commands are listed separately, several utilize identical executable files, like ‘reg add’ and ‘reg copy’. To present a unified perspective, commands with shared executable were further grouped into one, resulting in 306 unique Windows commands. 1409
1410
1411
1412
1413
1414

For common file name extensions, we first parsed all extensions from a clean Windows 10 1415
1416

³Windows Command-line reference A-Z

⁴SS64 Windows CMDs

<pre>- "C:\Windows\System32\bitsadmin.exe" /transfer 59697582645 /priority foreground http://example.com/example1234 "C:\Users\Public\Videos\V123456789\log32.dll" - powershell -command "Start-BitsTransfer -Source 'http://malicious.com/malicious1234' -Destination 'C:\Users\Public\Videos\V99999999\log32.dll' -Priority High"</pre>
<pre>- "cmd" /c "net use \\REMOTEDIR /user:Administrator password /persistent:no" - python -c "import os; os.system('net use \\REMOTEDIR /user:Administrator password /persistent:no')"</pre>
<pre>- reg query "HKLM\Software\Microsoft\Windows\CurrentVersion\Uninstall" /f "Chrome" /s - Get-ItemProperty HKLM:\Software\Microsoft\Windows\CurrentVersion\Uninstall* Select-Object -Property DisplayName, UninstallString Where-Object {\$_ .DisplayName -like '*Chrome*'} Format-Table -AutoSize</pre>
<pre>- schtasks /create /tn "TaskName" /tr "C:\Path\to\program.exe" /sc daily /st 00:00 - cronjob schedule daily 00:00 /path/to/program</pre>

Table 9: The similar command-line pairs in CmdDataset. Similar command lines are not merely similar on the lexical level but also in terms of their intrinsic purpose and semantic meaning.

Here are 12 Windows command line examples for referencing:

1. {sampled command line seed 1}
2. {sampled command line seed 2}
- ...
12. {sampled command line seed 12}

Your job is to synthesize 4 new Windows command lines. Please adhere to the following synthesizing guidelines:

- Ensure diverse command lines in appearance, argument value, purpose, result, and length, particularly making sure the generated command lines differ significantly from the reference command lines in every aspect.
- Prioritize practicality in generated commands, ideally those executed or executable. For example, please give me real argument value, filename, IP address, and username.
- Include Windows native commands, commands from installed applications or packages (for entertainment, work, artistic, or daily purposes), commands usually adopted by IT, commands corresponding with mitre att&ck techniques, or even some commonly used attack command lines. The more uncommon the command line, the better.,
- Do not always generate short command lines only. Be creative to synthesize all kind of command lines.

Give me your generated command lines only without any explanation or anything else. Separate each generated command line with "\n" and add a prefix "<CMD>" before each generated cmd.

Figure 6: The prompt used for generating a single command line. 12 exemplary command lines are randomly sampled from total command-line seeds for in-context demonstration.

Your task is to generate a similar Windows command line for each entry in the following command line list.
 In this task, 'similar' means that the command lines share the same purpose, or intention, rather than merely having a similar appearance.

Consequently, the generated command lines may differ significantly in argument values, format, and order from the original command line, or even from a different executable file, as long as they serve a similar purpose or intention.

{query command line}

Be creative to make the command lines appear distinctly different while adhering to the defined 'similar' criteria. For instance, you might employ obfuscation techniques, randomly rearrange the order of arguments, change the way to call the exe file, or substitute the executable file with a similar one.
 Please provide only the generated similar command lines without any explanation, prefixed with "<CMD>", and separate each command line with "\n".

Figure 7: The prompt used for generating similar command lines. The query command line is randomly sampled from the total command-line seeds.

Command Line	Explanation	Labels
net use Z: \\192.168.1.1\SharedFolder /user:administrator Passw0rd! findstr /i connected	This command line is used to map a network drive to the letter Z, connect to a shared folder on a specific IP address using administrator credentials, and then search for the keyword "connected" in the output.	Positive
schtasks /create /sc weekly /d MON,TUE,WED,THU,FRI /tn "WeeklyBackup" /tr "C:\Scripts\backup.bat" /st 18:00	This command line creates a scheduled task to run a backup script every weekday at 6:00 PM.	Positive
tasklist /fi "IMAGENAME eq notepad.exe" /fo list find "1234"	This command line is used to list all running processes with the name "notepad.exe" and then search for a specific process ID "1234" within the list.	Positive
findstr /s /i /m "hello world" "world take care" C:\Users* *.pdf	Search for the phrase "hello world" in all PDF files located in the C:\Users directory and its subdirectories, ignoring case and only displaying the file names that contain the phrase.	Neutral

Table 10: Example of command lines and their corresponding explanations generated by GPT-3.5-Turbo (OpenAI, 2022). The rightmost column denotes the labels (e.g., Positive, Neutral, or Negative) assigned by the expert.

virtual machine, and removing those with special characters and frequencies lower than 0.05%. This process yielded a total of 75 common extensions. We then identified how many of these extensions are included in our training set.

D Evaluation Metrics Detail

MRR@K is a key metric in information retrieval and recommendation systems. It calculates the average reciprocal rank of the first relevant item within a list of ranked results, focusing on the top K positions. This method helps us gauge how well the most relevant item ranks among the top 10 with the highest predicted scores. The Top@K metric is similar to MRR@K but differs in that it awards a score if the ground truth is within the top K ranks, without the rank-dependent decay seen in MRR.

E Hyperparameters and Training Process of CmdCaliper Detail

We trained CmdCaliper for 2 epochs using the Adam optimizer (Kingma and Ba, 2015) with a learning rate of 0.00002 and a batch size of 64. For the temperature parameter τ in Equation 1, we set it to 0.05. CmdCaliper was trained on three distinct model scales: small, base, and large. These models were initialized from the GTE-small (Li et al., 2023), GTE-base, and GTE-large, respectively.

We randomly selected 1,000 similar pairs from the training set to form a validation set. Evaluations on the validation set were conducted every 50 training steps. The checkpoint that demonstrated optimal performance on the validation set was then used for subsequent evaluations on the testing set.

F Semantic-Based Malicious Command-Line Detection Detail

Initially, we define the pre-collected set of malicious command lines as the ‘malicious gene pool’. Given a new command line, we leverage a command-line embedding model to obtain their embedding vectors, and compute the semantic similarity between each command line within the malicious gene pool. If one of the similarities exceeds a pre-defined threshold, we classify the new command line as malicious. This approach enables us to detect malicious command lines from a semantic perspective, even when attackers attempt to obfuscate command lines to evade pattern-based detection.

In this experiment, we utilize the open-source dataset: atomic-red-team (Canary), which encompasses a variety of command lines corresponding to numerous MITRE ATT&CK techniques (e.g., Abuse Elevation Control Mechanism or Browser Session Hijacking). The dataset consists of a set of techniques, denoted as $\{t_1, t_2, \dots, t_n\}$. Within each technique t_i , we obtain a set of malicious command lines, denoted as \mathbb{L}_i , containing M_i entries: $\mathbb{L}_i = \{c_1^i, c_2^i, \dots, c_{M_i}^i\}$. By unioning all these malicious command line sets, we form a total command line set \mathbb{A} , which consists of 1,523 malicious command lines across various techniques, represented as $\mathbb{A} = \bigcup_{i=1}^n \mathbb{L}_i$. To construct the malicious gene pool \mathbb{P}_i for each technique t_i , we select the first $r\%$ of the command lines from each technique, defined as $\mathbb{P}_i = \{c_1^i, c_2^i, \dots, c_{\lfloor \frac{r}{100} \times M_i \rfloor}^i\}$. The remaining command lines form the incoming command line set \mathbb{O}_i , denoted as: $\mathbb{O}_i = \{c_{\lfloor \frac{r}{100} \times M_i \rfloor + 1}^i, \dots, c_{M_i}^i\}$.

For malicious gene pool construction, we exclude techniques with fewer than 9 malicious command lines, resulting in a total of 55 distinct malicious gene pools corresponding to different techniques. For the evaluation of each technique t_i , we first exclude the command lines in the malicious gene pool to form a candidate command line set \mathbb{C}_i , denoted as $\mathbb{C}_i = \mathbb{A} \setminus \mathbb{P}_i$. We treat the incoming command line set \mathbb{O}_i as the positive set, while the negative set \mathbb{G}_i is formed by excluding all command lines corresponding to the technique t_i , denoted as $\mathbb{G}_i = \mathbb{A} \setminus \mathbb{L}_i$.

Given a command line c_j^i from the candidate command line set \mathbb{C}_i and an embedding model E , the embedding vector e_j^i can be computed by $e_j^i = E(c_j^i)$. Subsequently, the malicious score $s_{c_j^i}^i$ of the command line c_j^i for the technique t_i is determined by calculating its maximum similarity with each command line p_k^i in the malicious gene pool \mathbb{P}_i :

$$s_{c_j^i}^i = \max_{p_k^i \in \mathbb{P}_i} S(e_j^i, E(p_k^i)) \quad (2)$$

where $S(\cdot)$ is the similarity function (e.g., cosine similarity function). Note that if the command line c_j belongs to the positive set \mathbb{O}_i of the technique t_i , then the malicious score should exceed the pre-defined threshold τ for correct classification; otherwise, the score should be below the threshold if the command line is in the negative set \mathbb{G}_i .

We iteratively evaluated all 55 techniques and

1513 concatenated all malicious scores to compute the
1514 area under the curve (AUC). This evaluation metric
1515 aligns with real-world application scenarios where
1516 a static threshold is usually applied across all tech-
1517 niques.

1518 **G Classification Dataset Synthesis Detail**

1519 These commands to synthesize the command-
1520 line classification dataset were chosen based
1521 on their ability to accept a wide range of ran-
1522 domly generated strings as arguments, provided
1523 that the corresponding file exists. The classi-
1524 fier’s task is to identify the corresponding Win-
1525 dows command, regardless of the argument’s
1526 length or complexity. For example, the com-
1527 mand lines “find ‘fewj2po3kdlewfmpemrgborktig
1528 fe34krop4k5ogjs9rkgefefw34f” and “find ‘test”
1529 should both be correctly categorized under the
1530 ‘find’ command. This highlights the importance
1531 of a robust command-line embedding model in en-
1532 coding the command information within its embed-
1533 dings, as it plays a crucial role in determining the
1534 purpose of each command line.

1535 In this experiment, we used the pattern “<com-
1536 mand> ‘<argument value>” to randomly generate
1537 7,000 command lines for each command. Of these,
1538 3,500 were assigned to the training set and the re-
1539 maining 3,500 to the testing set. The arguments for
1540 each command line were formed by concatenating
1541 seven random strings, made up of ASCII letters
1542 and digits, with lengths ranging from 1 to 20 char-
1543 acters, separated by spaces. To increase the diffi-
1544 culty of classification and simulate the obfuscation
1545 techniques that attackers might use in real-world
1546 scenarios to evade detection, we also randomly in-
1547 corporated seven different commands into the argu-
1548 ment values. For example, a synthesized command
1549 line might read: “certutil.exe ‘msiexec tr9QIIL
1550 find C print oGod 5K 7Okf4 2ZcVT9 rundll32 sc
1551 query mNjIL robocopy q5””, where only the first
1552 command, ‘certutil’, is valid.

1553 We randomly selected 20% of the training set to
1554 serve as a validation set in the search for optimal
1555 hyperparameters. Subsequently, we utilized the
1556 obtained optimal hyperparameters to train a logistic
1557 regression classifier for performance evaluation on
1558 the testing set.