

Detecting Slang on the Dark Web Based on Word Co-occurrence Relationships in Anchor Texts

Authors:

Ken Shiozawa, Hiroo Hayashi, Soramichi Akiyama, and Marie Katsurai

Published in:

The 40th International Conference on Information Networking (ICOIN)

Link to IEEE Xplore:

To Be Added

Note:

The copyright of this work is held by IEEE

Detecting Slang on the Dark Web Based on Word Co-occurrence Relationships in Anchor Texts

Ken Shiozawa

*Graduate School of Science and Engineering
Doshisha University
Kyoto, Japan
shiozawa24@mm.doshisha.ac.jp*

Soramichi Akiyama

*College of Information Science and Engineering
Ritsumeikan University
Osaka, Japan
s-akym@fc.ritsumei.ac.jp*

Hiroo Hayashi

*Graduate School of Science and Engineering
Doshisha University
Kyoto, Japan
hayashi22@mm.doshisha.ac.jp*

Marie Katsurai

*Faculty of Science and Engineering
Doshisha University
Kyoto, Japan
katsurai@mm.doshisha.ac.jp*

Abstract—The dark web, accessed through anonymizing networks such as Tor, enables various illicit activities including drug and firearm trafficking, cyberattack services, and the exchange of personal confidential information. While relevant stakeholders worldwide are actively addressing these challenges, establishing effective investigation methodologies remains difficult due to the widespread use of specialized slang terminology designed to evade surveillance systems. This paper introduces a novel approach for detecting criminal slang by leveraging hyperlink structures within dark web content. Our method collects anchor tags from HTML files and analyzes their hyperlinks to identify words that have strong relationships with known crime-related terms. Specifically, we assume that words in anchor texts that share common link destinations have “co-occurrence relationships via hyperlinks” and employ pointwise mutual information to quantify their strength. We conducted experiments using 251,892 HTML files crawled from the dark web to compare the proposed method with both a simple co-occurrence method applied to anchor texts only and a large language model. Our method achieved higher accuracy in identifying verified slang terms by effectively filtering noisy words while maintaining competitive detection rates.

Index Terms—dark web, slang detection, anchor text, hyperlink analysis, word co-occurrence analysis

I. INTRODUCTION

The Web is commonly divided into the surface web and the deep web. The surface web consists of pages that can be accessed via search engines, such as homepages, social networking services, and blogs. On the other hand, the deep web consists of pages that cannot be accessed via search engines, such as private database contents, corporate intranets, encrypted file repositories, and includes content that website owners have blocked from search engines or secured with login requirements. The dark web is a small part of the deep web that is accessible only through anonymizing networks like Tor (The Onion Router) or I2P (The Invisible Internet Project) and specialized applications.

Dark web users can trade services and products without revealing their identities and without coming into contact with

sellers. As a result, the sale of drugs, child pornography, weapons, credit card information, and other illegal items is rampant. Law enforcement agencies, security companies, and specialized cybercrime units around the world are working to identify websites and their operators in order to detect and crack down on such activities [1]. However, those who engage in illegal transactions on the dark web are wary of such arrests and have developed multiple strategies to avoid detection. Among these methods, they commonly attempt to evade surveillance by using slang terminology. Slang refers to special terms used by limited groups of people to express specific concepts. For example, common drug slang terms include “coke” and “snow” for cocaine, and “blueberry” and “grass” for marijuana. By using these slang terms, they conceal the nature of their transactions from monitors and hinder detection through keyword searches using general vocabulary. When the meaning of slang terms becomes widely known, new slang terms often emerge to replace them. Therefore, there is a need for technology that can automatically detect slang terms from text on the dark web. Conventional studies used neural network-based language models to calculate textual features of words on the dark web and then detected words whose usage differs from standard texts [2], [3]. However, these methods require language model training based on large document collections, making it difficult to efficiently learn the meanings of newly emerged slang terms. The computational cost and time required for retraining models to adapt to evolving criminal terminology pose significant practical challenges.

This paper proposes a simple method for detecting slang on the dark web by utilizing hyperlinks and anchor text within pages. The proposed method assumes that there is a strong semantic relationship between certain crime-related terms and words that share common link destinations. Our method assumes that when some page authors directly express crime-related terms using anchor texts without employing slang, other authors may use slang alternatives when linking

to the same content. Based on this assumption, we collected anchor tags from HTML files on the dark web and calculate the co-occurrence frequency of words within anchor texts that point to the same page using the href attribute. This approach enables us to discover alternative expressions, i.e., slang terms, for crime-related terms by analyzing these co-occurrence patterns.

To verify the effectiveness of the proposed method, we conducted slang detection experiments on 251,892 HTML files. The results demonstrated that the proposed method effectively discovers slang terms for specified crime-related terms and successfully detects numerous alternative terms for crime-related substances, including slang, chemical names, and product names, while significantly outperforming baseline approaches in terms of accuracy and noise reduction.

II. RELATED WORK

Several related works have vectorized the meanings of words on the dark web and discovered words that are used in ways that differ from their original meanings by calculating the similarity between word vectors. In particular, researchers have developed methods to detect words with different meanings in the dark web and surface web by constructing separate corpora for each domain, training word vectors on each corpus, and identifying words whose vector representations differ significantly between the two domains as candidate slang words. Yuan et al. [2] extended the word2vec architecture [4] to accept inputs from two corpora, enabling semantic comparison of words between different corpora in the output layer. Seyler et al. [5] first obtained semantic vectors for words in one corpus, then searched for general words similar to each word in the other corpus, performing this process bidirectionally. They considered words to be slang when their respective nearest neighbors differed. Ke et al. [6] trained BERT [7] using a self-constructed Chinese dark web corpus and calculated semantic similarity between slang and general expressions. Specifically, using 10 types of slang terms representing drugs and weapons as clues, words with a high degree of similarity to them were selected as slang candidates, and those with a low degree of similarity between the dark web corpus and general usage were considered to be slang. While these neural network-based approaches have shown promising results, they require substantial computational resources and large-scale corpora, making rapid adaptation to newly emerged slang terminology challenging.

In research on information retrieval not limited to the dark web, methods have been proposed for representing the meaning of words using web page tags and text to analyze linking patterns and extract semantic relationships. Anchor tags on web pages have attracted attention because they contain rich semantic information, with anchor text frequently serving as human-generated labels that capture alternative names, synonyms, and descriptive expressions for the target content. Examples include applications to synonym extraction [8], translation extraction [9], and personal name alias

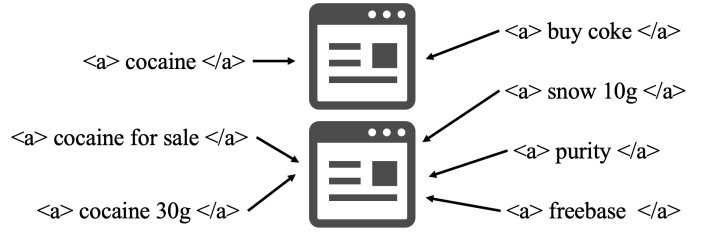


Fig. 1. Overview of word co-occurrence relationships based on hyperlinks. When anchor texts from different web pages point to the same destination, words appearing in these texts are collected as link-sharing words. PMI calculation then quantifies the semantic relationship between seed keywords (e.g., ‘cocaine’) and potential slang terms (e.g., ‘coke’, ‘snow’, ‘freebase’) based on their co-occurrence patterns in shared link destinations.

extraction [10]–[12]. In particular, Bollegala et al. [11] automatically extracted aliases or nicknames of specific individuals by utilizing word co-occurrence, lexical pattern frequencies, and page counts in anchor texts when links from different web pages pointed to the same web page. In this study, we also collect anchor texts from the dark web and apply a similar approach to the detection of slang terms related to crime.

Several recent works have applied large language models (LLMs) to slang detection. Sun et al. [13] constructed a dataset from movie subtitles to evaluate LLMs’ ability to understand slang, showing that Generative Pre-trained Transformer (GPT) models performed better than BERT but required fine-tuning for optimal performance. Fillies et al. [14] proposed a simple LLM-based solution for detecting “Algospeak” (a variant language for content moderation evasion), demonstrating that it can be effectively decoded using GPT-4 and prompt engineering alone. Sun et al. [15] developed a framework that combines contextual information and knowledge of how slang meanings extend from conventional word meanings to improve performance in both slang interpretation and machine translation. In our experiments, we use the LLM-based slang interpretation method from Sun et al. [13] as a baseline to evaluate the effectiveness of our proposed method.

III. PROPOSED METHOD

In this section, we propose a method for detecting slang based on hyperlinks on the dark web. Anchor text with the same link destination on a web page can be considered to contain related words. For example, as shown in Figure 1, by collecting anchor tags that point to the same web page as links with “cocaine” as anchor text, it is expected that alternative terms for cocaine (such as “coke”) can be collected. For convenience, we refer to the set of words collected from anchor texts that share common link destinations as **link-sharing words**. Specifically, given a seed keyword, link-sharing words are all the words (excluding the seed keyword itself) that appear in anchor texts pointing to the same destination URLs as anchor texts containing the seed keyword. This concept is fundamental to our method because it captures words that are used in similar contexts to refer to the same web content, thereby identifying potential alternative expressions including

slang terms, product variants, and related terminology. The proposed method sets known crime-related terms as seed keywords and collects related words that share the same link destinations (see Section III-A). Next, it quantifies the strength of relationships between seed keywords and related words, extracting words with higher values as those with high potential to be slang (see Section III-B).

A. Collecting shared words linked to seed keywords

First, we collect anchor tags containing seed keywords from HTML data obtained by crawling the dark web, and extract all other anchor tags that point to the same destinations. Next, we eliminate unnecessary words that could cause noise in the proposed method. Specifically, we perform the following preprocessing on the set of anchor texts obtained from the anchor tags:

- 1) Separate each anchor text with spaces and convert it into a list of tokens (i.e., individual words).
- 2) Remove tokens matching NLTK¹ stop words, numeric-only tokens, special symbol-only tokens, and URL tokens.
- 3) Remove tokens longer than 30 characters.
- 4) Remove anchor text with more than 30 tokens.

The set of tokens obtained after applying the above preprocessing is considered to be the seed keyword’s link-sharing words.

B. Slang extraction

Next, we extract slang from the collected link-sharing words. Simply counting co-occurrence frequencies between seed keywords and link-sharing words would result in commonly used words being ranked at the top. Therefore, the proposed method uses pointwise mutual information (PMI), which considers the occurrence frequency of each word itself, to quantify the strength of relationships between seed keywords and each link-sharing word. Let x represent a seed keyword and y represent any of its link-sharing words. Given $P(x)$ as the occurrence probability of word x and $P(x, y)$ as the joint probability of words x and y , the pointwise mutual information PMI of words x and y is calculated by the following equation:

$$\begin{aligned} \text{PMI}(x, y) &= \log \left(\frac{P(x, y)}{P(x)P(y)} \right) \\ &= \log \left(\frac{C(x, y) \cdot N}{C(x) \cdot C(y)} \right), \end{aligned} \quad (1)$$

where $C(x)$ represents the number of pages linked from word x , $C(x, y)$ represents the number of pages linked from both words x and y , and N represents the total number of pages in the dataset. The larger the value of $\text{PMI}(x, y)$, the stronger the relationship between link-sharing word y and seed keyword x . Therefore, by fixing x and sorting $\text{PMI}(x, y)$ in descending order, examining the top-ranked y values is expected to reveal slang terms that are alternative expressions of x . When anchor texts contain multiple words, we extract all individual tokens

from each anchor text during preprocessing and calculate PMI values for each token separately. As demonstrated in Figure 1, when anchor texts such as ‘buy coke’, ‘snow 10g’, and ‘purity freebase’ all point to the same destination page as the seed keyword ‘cocaine’, we extract individual tokens: ‘buy’, ‘coke’, ‘snow’, ‘10g’, ‘purity’, and ‘freebase’. Each extracted token is then evaluated independently against the seed keyword ‘cocaine’ using the PMI calculation. This token-level analysis enables the detection of slang terms that may be embedded within longer commercial phrases commonly found in dark web marketplaces.

IV. EXPERIMENTS

We conducted experiments to verify the effectiveness of the proposed method for detecting slang terms on the dark web. Six words representing illegal drugs were selected as seed keywords, as shown in Table I. These drug names were extracted from the official list published by the U.S. Drug Enforcement Administration². The following subsections describe the dataset construction and preprocessing (Section IV-A), baseline methods for comparison (Section IV-B), and experimental results with discussion (Section IV-C).

A. Dataset construction and anchor text preprocessing

The experimental dataset consists of 251,892 HTML files obtained from 14,363,466 URLs using the large-scale dark web crawling method proposed by Kimura et al. [16]. The number of HTML files is smaller than that of URLs because many URLs returned 404 errors or were unreachable due to invalid domains. To investigate the languages used within this dataset, we analyzed the lang attribute in HTML tags. For pages without a specified lang attribute, language detection was performed using Python’s langdetect library³. We then determined the primary language of each domain based on the language distribution of its successfully retrieved HTML pages. Consistent with previous dark web studies [17]–[19], English domains were predominant, accounting for 95.6% of the 4,655 domains from which HTML was successfully retrieved. Therefore, we limited our analysis to English-language domains.

We extracted 18,139,380 anchor tags from the HTML data using the Python library BeautifulSoup⁴. After applying the preprocessing steps described in Section III-A, the number of anchor tags was reduced to 16,675,088.

B. Baseline methods

To validate the effectiveness of our proposed method, we compared it with the following two baseline methods that represent different paradigms for slang detection.

Baseline 1 (w/o shared links): This method simply calculates the PMI of words with seed keywords in anchor texts without considering shared link relationships. For example, given anchor texts “buy cocaine here,” “cocaine for sale,” and

²<https://www.dea.gov/sites/default/files/2018-07/DIR-022-18.pdf>

³<https://pypi.org/project/langdetect/>

⁴<https://pypi.org/project/beautifulsoup4/>

¹<https://www.nltk.org/>

TABLE I
STATISTICS OF SEED KEYWORDS AND ANCHOR TAG OCCURRENCES.

| Seed Keyword | Total Seed Tags |
|-----------------|-----------------|
| cocaine | 1484 |
| MDMA | 378 |
| heroin | 122 |
| LSD | 699 |
| marijuana | 253 |
| methamphetamine | 14 |

TABLE II
PROMPT USED FOR LLM BASELINE METHOD.

| |
|--|
| Is the following word slang for {Seed_keyword}? Answer only ‘Yes’ or ‘No’. |
| Word: {word} |
| Answer: |

“get cocaine now” that may link to different websites, this approach counts the co-occurrence frequency of “buy,” “sale,” and “get” with “cocaine” regardless of their link destinations. Words are ranked by PMI, and the top-ranked words are extracted as potential slang terms.

Baseline 2 (LLM): This method uses Meta-Llama-3-8B-Instruct [20] to determine whether tokens in anchor texts are slang or not. For each token, the LLM is prompted to make a binary decision about whether the word represents slang for a given seed keyword. The prompt used is shown in Table II, where {Seed_keyword} is replaced with each drug name and {Word} is replaced with tokens from the anchor text.

These baselines allow us to evaluate our approach against both traditional co-occurrence methods and a state-of-the-art language model.

C. Slang detection results

We calculated the strength of word relationships to seed keywords using our proposed method and the two baseline approaches. Tables III and IV show the top 20 words extracted by the proposed method and PMI-based baseline, respectively. When there are insufficient link-sharing words or anchor text occurrences, fewer than 20 words are listed. Words verified as actual slang terms by the first and second authors through web search are shown in bold.

Our proposed method successfully identified many genuine slang terms across all drug categories. As shown in Table III, our method found well-known slang terms such as “fish-scale” and “freebase” for cocaine, “molly” and “ecstasy” for MDMA, “white” and “tar” for heroin, “blotters” and “tabs” for LSD, and “shatter” and “delta” for marijuana. Even for methamphetamine, which had very few anchor tag occurrences (only 14 total as shown in Table I), our method still found relevant terms such as “uncut,” “ice,” and “meth.” In contrast, Baseline 1 (w/o shared links) produced poor results, as shown in Table IV. The top-ranked words were dominated by weight and quantity indicators like “3gr,” “10g,” “1kg,” “2g,” “25g,” “250mg,” vendor-related terms like “vendershop,” and location-related terms like “USA,” “JAPAN,” and “INDIA.”

This happened because Baseline 1 simply counted how often words appeared together in anchor texts without considering whether they linked to the same destination pages. Even when PMI values were calculated for these co-occurrences, the method could not distinguish between meaningful drug slang and these marketplace-related terms. Our method performed much better at finding real slang terms. For example, when looking for cocaine slang, Baseline 1 found only two genuine slang terms in its top 20 results, while our method found six genuine slang terms. Baseline 2 (LLM) showed different problems, as shown in Table V. Unlike Tables III and IV, which present the top 20 ranked terms with their scores, Table V focuses on detection accuracy, showing the large discrepancy between the total number of detected terms and correctly identified ones (e.g., only 42 out of 320 detected terms for cocaine were correctly identified). Meta-Llama-3-8B-Instruct found many more potential slang candidates (300–663 words per drug), but most of them were wrong. For cocaine, the model identified 320 possible slang terms, but only 42 of them were actually correct slang (13% accuracy). The accuracy was similarly low for other drugs: MDMA (5.3%), heroin (4.3%), LSD (6.7%), marijuana (7.5%), and methamphetamine (5.5%). This poor performance occurred because the model made decisions based only on individual words using the simple prompt shown in Table II, without enough context to make accurate judgments.

Our proposed method achieved much better accuracy than both baseline methods. The PMI calculation helped filter out noisy words while keeping genuine slang terms. By analyzing which words appear in links pointing to the same web pages, our method could identify words that are truly related to the target drug rather than just frequently mentioned together. This approach successfully reduced the marketplace-related noise that dominated the PMI-based method and avoided the high error rate of the language model approach. However, our method had some limitations: it sometimes found slang terms for multiple drugs instead of just the target drug. For example, when searching for cocaine slang, our method also found some terms related to other drugs like MDMA and heroin. To solve this problem and find slang terms only for the specific target drug, future work would need to analyze the actual content of the linked web pages, not just the anchor text.

V. CONCLUSION AND FUTURE WORK

In this paper, we proposed a simple method based on hyperlinks for detecting slang terms on the dark web. Our approach used the assumption that anchor texts pointing to the same destination page share semantic relationships, enabling the identification of alternative expressions for crime-related terms without requiring extensive corpus-based training. The proposed method collected anchor tags from dark web HTML files and used PMI to quantify relationships between seed keywords and link-sharing words. Through experiments on 251,892 HTML files, we confirmed that words with strong PMI relationships to seed keywords were actually verified slang terms. Our method demonstrated much better accuracy in

identifying verified slang terms compared to both frequency-based and LLM baselines while effectively reducing weight-related noise words that dominated conventional approaches.

However, our method had some limitation that need to be addressed in future work. First, the method extracted slang terms for other drugs and drug-related terms beyond the target substance. For example, when searching for cocaine-related slang, the method also identified terms related to other drugs such as MDMA and heroin. To exclude these and extract slang terms only for the target drug, it would be necessary to identify the content of linked web pages. Second, in our preprocessing steps, we removed numeric-only tokens to reduce noise, but some slang terms consist entirely of numbers (such as drug code names or street numbers). Future work should consider methods to distinguish between meaningful numeric slang and noise numbers. Additionally, future research will focus on improving slang detection performance by exploring the integration of link destination content analysis to enhance specificity and developing more sophisticated filtering techniques to better identify legitimate slang terms.

ACKNOWLEDGMENT

We would like to express our gratitude to Mr. Yuuki Kimura of the Graduate School of Information Science and Engineering, Ritsumeikan University, for providing us with data collected from the dark web. This study was carried out using the TSUBAME4.0 supercomputer at Institute of Science Tokyo.

REFERENCES

- [1] R. Basheer and B. Alkhatib, "Threats from the dark: a review over dark web investigation research for cyber threat intelligence," *Journal of Computer Networks and Communications*, vol. 2021, pp. 1–21, 2021.
- [2] K. Yuan, H. Lu, X. Liao, and X. Wang, "Reading thieves' cant: automatically identifying and understanding dark jargons from cybercrime marketplaces," in *27th USENIX Security Symposium (USENIX Security 18)*, 2018, pp. 1027–1041.
- [3] Y. Jin, E. Jang, J. Cui, J.-W. Chung, Y. Lee, and S. Shin, "DarkBERT: A language model for the dark side of the internet," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2023, pp. 7515–7533.
- [4] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems*, vol. 26. Curran Associates, Inc., 2013, pp. 3111–3119.
- [5] D. Seyler, W. Liu, X. Wang, and C. Zhai, "Towards dark jargon interpretation in underground forums," in *Advances in Information Retrieval*, ser. Lecture Notes in Computer Science, vol. 12657. Springer, 2021, pp. 393–400.
- [6] L. Ke, P. Xiao, X. Chen, S. Yu, X. Chen, and H. Wang, "A novel cross-domain adaptation framework for unsupervised criminal jargon detection via pre-trained contextual embedding of darknet corpus," *Expert Systems with Applications*, vol. 242, p. 122715, 2024.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, vol. 1. Association for Computational Linguistics, 2019, pp. 4171–4186.
- [8] Z. Chen, S. Liu, L. Wenxin, G. Pu, and W.-Y. Ma, "Building a web thesaurus from web link structure," in *Proceedings of the 26th annual international ACM SIGIR Conference on Research and Development in Informaion Retrieval*, 2003, pp. 48–55.
- [9] W.-H. Lu, L.-F. Chien, and H.-J. Lee, "Anchor text mining for translation of web queries: A transitive translation approach," *ACM Transactions on Information Systems (TOIS)*, vol. 22, no. 2, pp. 242–269, 2004.
- [10] D. Bollegala, Y. Matsuo, and M. Ishizuka, "A co-occurrence graph-based approach for personal name alias extraction from anchor texts," in *Proceedings of the Third International Joint Conference on Natural Language Processing: Volume-II*, 2008.
- [11] Y. Matsuo, M. Ishizuka, and D. Bollegala, "Automatic discovery of personal name aliases from the web," *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 6, pp. 831–844, 2010.
- [12] R. Jayabhaduri *et al.*, "Automatic discovery of association orders between name and aliases from the web using anchor texts-based co-occurrences," *International Journal of Computer Applications*, vol. 41, no. 19, 2012.
- [13] Z. Sun, Q. Hu, R. Gupta, R. Zemel, and Y. Xu, "Toward informal language processing: Knowledge of slang in large language models," in *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Mexico City, Mexico: Association for Computational Linguistics, 2024, pp. 1683–1701.
- [14] J. Fillies and A. Paschke, "Simple LLM based approach to counter algospeak," in *Proceedings of the 8th Workshop on Online Abuse and Harms (WOAH 2024)*. Mexico City, Mexico: Association for Computational Linguistics, 2024, pp. 136–145.
- [15] Z. Sun, R. Zemel, and Y. Xu, "Semantically informed slang interpretation," in *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Seattle, United States: Association for Computational Linguistics, 2022, pp. 5213–5231.
- [16] Y. Kimura, S. Akiyama, A. Inomata, and T. Uehara, "On collecting onion server fingerprints and identification of their operators," in *International Conference on Software Quality, Reliability, and Security Companion (QRS-C)*, 2023, pp. 540–548.
- [17] M. W. Al Nabki, E. Fidalgo, E. Alegre, and I. De Paz, "Classifying illegal activities on tor network based on web textual contents," in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, 2017, pp. 35–43.
- [18] M. W. Al-Nabki, E. Fidalgo, E. Alegre, and L. Fernández-Robles, "Torank: Identifying the most influential suspicious domains in the tor network," *Expert Systems with Applications*, vol. 123, pp. 212–226, 2019.
- [19] Y. Jin, E. Jang, Y. Lee, S. Shin, and J.-W. Chung, "Shedding new light on the language of the dark web," *arXiv preprint arXiv:2204.06885*, 2022.
- [20] M. AI, "Introducing meta llama 3: The most capable openly available llm to date," 2024. [Online]. Available: <https://ai.meta.com/blog/meta-llama-3/>

TABLE III
TOP 20 WORDS RELATED TO SEED KEYWORDS USING THE PROPOSED METHOD (PMI). BOLD WORDS INDICATE VERIFIED SLANG TERMS.

| No. | cocaine | | MDMA | | heroin | | LSD | | marijuana | | methamphetamine | |
|-----|------------------|------|------------------|------|-----------------|------|-----------------|------|-------------------|-------|-----------------|------|
| | Word | PMI | Word | PMI | Word | PMI | Word | PMI | Word | PMI | Word | PMI |
| 1 | meths | 8.37 | champagne | 8.93 | freebase | 9.14 | gummies | 9.33 | dispensary | 10.00 | uncut | 8.16 |
| 2 | freebase | 8.37 | cola | 8.82 | undermarket | 8.19 | meths | 9.29 | pens | 9.72 | select | 6.44 |
| 3 | stimulant | 8.37 | freebase | 8.38 | afgan | 8.04 | asteroid | 9.22 | shatter | 9.72 | ice | 5.92 |
| 4 | tronger | 8.37 | afgan | 7.28 | afghan | 7.98 | mcg | 9.08 | glow | 9.38 | usa | 5.74 |
| 5 | blosters | 8.33 | fishscale | 7.13 | fishscale | 7.89 | blotters | 9.08 | edibles | 7.64 | options | 5.66 |
| 6 | asteroid | 8.29 | xtc | 7.03 | champagne | 7.12 | potassium | 8.77 | vape | 6.31 | meth | 5.42 |
| 7 | fishscale | 8.03 | ritalin | 6.77 | tar | 6.46 | cocs | 8.64 | thc | 6.12 | product | 5.26 |
| 8 | potassium | 7.85 | ecstasy | 5.71 | dmt | 5.64 | cyanide | 8.64 | grade | 5.83 | | |
| 9 | cyanide | 7.72 | molly | 5.16 | ketamine | 5.49 | cali | 8.64 | medical | 5.80 | | |
| 10 | merch | 7.68 | dmt | 4.88 | sw | 5.48 | blotter | 8.64 | cannabis | 5.74 | | |
| 11 | thcservice | 7.68 | ketamine | 4.73 | uncut | 5.24 | freebase | 8.64 | sell | 5.25 | | |
| 12 | bolivian | 7.56 | dutch | 4.71 | pure | 4.41 | afghan | 8.35 | cartridges | 5.17 | | |
| 13 | afghan | 7.43 | gram | 4.65 | powder | 4.16 | tabs | 8.16 | rolls | 4.86 | | |
| 14 | serotonin | 7.27 | hq | 4.53 | select | 3.93 | coke | 7.95 | wax | 4.80 | | |
| 15 | columbia | 7.25 | lsd | 4.40 | diamond | 3.90 | ald | 7.95 | delta | 4.79 | | |
| 16 | purity | 7.10 | terms | 4.25 | white | 3.77 | goblin | 7.78 | japan | 4.38 | | |
| 17 | coke | 7.06 | mg | 4.00 | mdma | 3.51 | erowid | 7.72 | seeds | 4.35 | | |
| 18 | charlieuk | 6.76 | speed | 3.86 | black | 3.46 | afgan | 7.54 | flowers | 4.27 | | |
| 19 | spain | 6.63 | pills | 3.41 | real | 3.33 | capsules | 7.54 | vendor | 4.26 | | |
| 20 | peruvian | 6.47 | update | 3.29 | speed | 3.29 | fishscale | 7.39 | cocaine | 4.08 | | |

TABLE IV
TOP 20 WORDS RELATED TO SEED KEYWORDS USING BASELINE 1 (W/O SHARED LINKS). BOLD WORDS INDICATE VERIFIED SLANG TERMS.

| No. | cocaine | | MDMA | | heroin | | LSD | | marijuana | | methamphetamine | |
|-----|------------------|------|------------------|------|--------------|-------|-----------------|------|----------------|-------|-----------------|------|
| | Word | PMI | Word | PMI | Word | PMI | Word | PMI | Word | PMI | Word | PMI |
| 1 | 3gr | 8.87 | maserati | 9.98 | afgan | 10.37 | aurum | 9.70 | 1000g | 10.70 | uncut | 8.10 |
| 2 | vendershop | 8.87 | champagne | 9.62 | afghan | 8.61 | blosters | 9.70 | napoli | 10.70 | ice | 6.51 |
| 3 | colu | 8.87 | cola | 9.36 | 10g | 8.47 | 125ug | 9.70 | shatter | 10.29 | usa | 6.02 |
| 4 | kilo | 8.87 | 2g | 9.06 | undermarket | 8.39 | 200mcg | 9.70 | dispensary | 10.14 | meth | 5.73 |
| 5 | 50grams | 8.87 | 10g | 8.77 | 1g | 7.74 | offeremoji | 9.70 | glow | 9.67 | | |
| 6 | charlieuk | 8.82 | xtc | 8.09 | tar | 7.43 | cocs | 9.68 | pens | 9.60 | | |
| 7 | fishscale | 8.74 | ritalin | 7.63 | 5g | 7.29 | 200ug | 9.61 | edibles | 8.85 | | |
| 8 | 1g | 8.53 | 25g | 7.58 | drug | 6.61 | mcg | 9.48 | cocaine | 7.48 | | |
| 9 | 2g | 8.36 | 250mg | 7.21 | sw | 6.39 | blotters | 9.45 | vape | 6.65 | | |
| 10 | bolivian | 8.18 | 5g | 7.06 | uncut | 5.29 | 250ug | 9.36 | rolls | 6.21 | | |
| 11 | thcservice | 8.17 | molly | 5.83 | powder | 4.98 | blotter | 9.01 | wax | 6.06 | | |
| 12 | purity | 7.96 | ecstasy | 5.66 | market | 4.61 | 300ug | 9.01 | medical | 6.01 | | |
| 13 | columbia | 7.76 | gr | 5.14 | pure | 4.35 | cali | 9.01 | india | 6.00 | | |
| 14 | 1kg | 7.48 | dutch | 5.09 | diamond | 4.34 | ld | 8.79 | grade | 5.83 | | |
| 15 | marijuana | 7.48 | gram | 4.87 | grade | 4.11 | tabs | 8.72 | thc | 5.70 | | |
| 16 | sapin | 7.46 | mg | 4.85 | real | 4.05 | ug | 8.38 | japan | 5.11 | | |
| 17 | nextpress | 7.08 | terms | 4.78 | online | 3.74 | erowid | 8.10 | seeds | 5.00 | | |
| 18 | poland | 6.95 | hq | 4.72 | photo | 3.57 | ritalin | 7.35 | store | 5.00 | | |
| 19 | scale | 6.40 | lsd | 3.75 | buy | 3.22 | tab | 6.90 | seed | 4.95 | | |
| 20 | cook | 6.11 | sale | 3.21 | asian | 1.87 | ald | 6.76 | cannabis | 4.95 | | |

TABLE V
SLANG DETECTION PERFORMANCE OF BASELINE 2 (LLM).

| Drug Name | LLM-Detected | Correctly Identified | Examples |
|-----------------|--------------|----------------------|----------------------------------|
| cocaine | 320 | 42 | coke, snow, fishscale, charlieuk |
| MDMA | 300 | 16 | ecstasy, molly, xtc, speed |
| heroin | 323 | 14 | black, white, horse, junk |
| LSD | 356 | 24 | acid, blotter, tabs, lucy |
| marijuana | 663 | 50 | weed, pot, ganja, hash |
| methamphetamine | 488 | 27 | meth, ice, speed, crystal |