

VisPAD : Visualization and Pattern Discovery tool for fighting Human Trafficking

Pratheeksha Nair
pratheeksha.nair@mail.mcgill.ca
SCS, McGill University & Mila
Montreal, Canada

Reihaneh Rabbany
rrabba@cs.mcgill.ca
SCS, McGill University & Mila
Montreal, Canada

Lars Thørvæld
The Thørvæld Group
Hekla, Iceland
larst@affiliation.org

Valerie Béranger
Inria Paris-Rocquencourt
Rocquencourt, France

Aparna Patel
Rajiv Gandhi University
Doimukh, Arunachal Pradesh, India

Huifen Chan
Tsinghua University
Haidian Qu, Beijing Shi, China

Charles Palmer
Palmer Research Laboratories
San Antonio, Texas, USA
cpalmer@prl.com

John Smith
The Thørvæld Group
Hekla, Iceland
jsmith@affiliation.org

Julius P. Kumquat
The Kumquat Consortium
New York, USA
jpkumquat@consortium.net

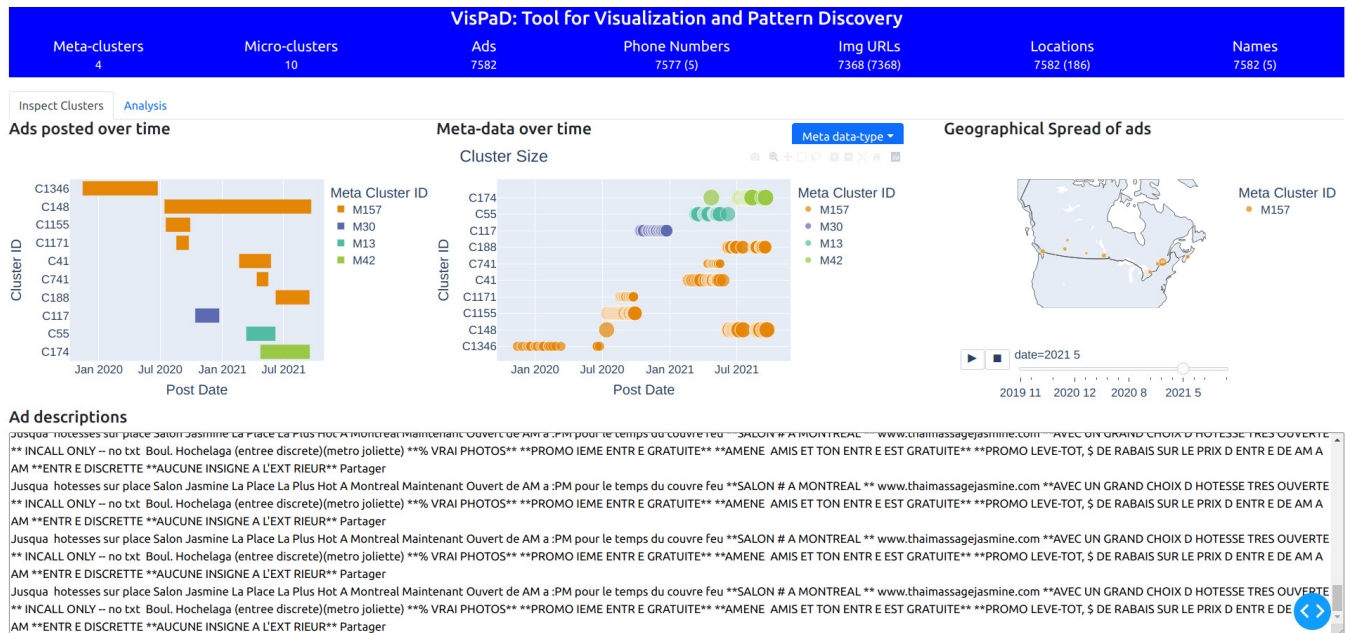


Figure 1: VisPAD showing a group of selected micro-clusters of escort ads. The top blue banner summarizes the information from selected micro-clusters. In clock-wise order starting from top left is the frequency of ads posted between the earliest and latest post, the meta-data information of micro-clusters over time, the geographical spread of ads and the ad descriptions.

ABSTRACT

A clear and well-documented \LaTeX document is presented as an article formatted for publication by ACM in a conference proceedings or journal publication. Based on the “acmart” document class, this article presents and explains many of the common variations, as well as many of the formatting elements an author may use in the preparation of the documentation of their work.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

KEYWORDS

datasets, neural networks, gaze detection, text tagging

ACM Reference Format:

Pratheeksha Nair, Reihaneh Rabbany, Lars Thørväld, Valerie Béranger, Aparna Patel, Huifen Chan, Charles Palmer, John Smith, and Julius P. Kumquat. 2018. VisPAD : Visualization and Pattern Discovery tool for fighting Human Trafficking. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 4 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Human trafficking (HT) is a nefarious crime affecting over 4.8 million people in the world [1] that often goes undetected and is difficult to tackle. A report on the involvement of technology in recruiting human trafficking victims[4] shows that a large number of the victims are advertised online on escort websites[7] with over fifty percent of the victims having no say in the content of their advertisements (ads).

These escort websites also contain ads from legitimate sex workers (posted at their own will), massage parlors (which may or may not include HT victims) and a large number of fake ads. These fake ads could be posted by scammers with the purpose of taking money from customers seeking escorts or spammers whose motive is unclear. These fake ads sometimes show characteristics similar to HT and make it more difficult to track down the suspicious ones. Separating each of these types (called *modus operandi* or M.O), especially removing spam and scam ads, can help in producing a more useful sample for identifying suspicious groups of ads.

Since traffickers entirely control the ad content for most of their victims [7], the ads posted by the same trafficker tend to have similarities, especially in the text template. This idea has been the cornerstone of several works[2, 3, 6] that are focused on detecting HT from escort ads. However, the problem becomes more complex as other M.Os (such as *spam* by the same spammer for example) may also contain similarities in the text. Moreover, these M.Os may also keep evolving as the traffickers adapt their strategies over time to avoid detection.

Previous work such as InfoShield [2] identifies groups of ads (called ‘micro-clusters’) that are related based on shared text templates. Apart from descriptions, the ads are also associated with information such as phone numbers, email addresses, social media tags, location of posting, etc which are known as ‘meta-data’. Different micro-clusters that share the same meta-data can further be grouped together into ‘meta-clusters’[8]. Having obtained these micro-clusters, domain experts look for certain characteristics such as the geographical spread, timeline of posting, presence of invalid URLs (potentially spam), occurrence of certain trigger words such as ‘gfe’, etc to identify suspicious ads for further investigation[5].

In this paper we introduce, for the first time to the best of our knowledge, micro-cluster characterization and representation as n -dimensional vectors in a feature space. We also introduce a novel tool, VisPAD , for visualizing these micro-cluster embeddings and

Table 1: Features used for micro-cluster characterization

Feature	Description
Cluster Size	Number (#) of ads
Phone Count	# of unique phone numbers
Location Count	# of distinct locations
Location Radius	Maximum radius of ad locations
Phone Entropy	Entropy of phone number occurrences
Person Name Count	# of names in the ad description
Valid URL count	# of valid URLs in the ad description
Invalid URL count	# of invalid URLs in the ad description
Ads per week	# of ads posted per week

helps discover patterns of specific behaviours. More specifically, this work makes the following main contributions:

- (1) *Cluster visualization platform*: VisPAD is an interactive tool for visualizing cluster embeddings and related information with respect to escort ads in the context of HT detection.
- (2) *Pattern discovery tool*: VisPAD allows a user to inspect specific, multi-faceted information of groups of clusters with suspicious/interesting patterns.
- (3) *Cluster characterization*: We introduce micro-cluster characterization as multi-dimensional vectors (meaning representations) for HT detection. These vectors can also act as features for developing learning algorithms

2 OVERVIEW OF VISPAD

In this demo paper, we introduce VisPAD , an open-source web application¹ for visualization of micro-cluster embeddings and inspection of suspicious patterns.

2.1 Cluster characterization

Each micro-cluster is characterized by nine feature values determined based on discussions with domain experts [5]. These features are mentioned in Table 1. The cluster size is a good indicator of which micro-clusters to pay attention to. Smaller micro-clusters (≤ 10) tend to usually be noise or independent escort ads. A larger location count and radius within a short duration may be indicative of spam. A high person name count is considered suspicious as ads indicating multiple escorts may potentially be HT [2]. The presence of several invalid URLs in the text can also be indicative of spam/fake ads. However, these are not definite indicators of these M.Os and need to be analyzed in conjunction with other information such as ad descriptions, posting frequency and meta-data information over time. VisPAD compiles all this information and presents them through a single interface.

2.2 Design and Usage Scenario

Figure 1 shows the landing page a user (e.g, a law enforcement agent, HT domain expert) first sees when they open the tool.

¹<https://github.com/ComplexData-MILA/MicroViz>

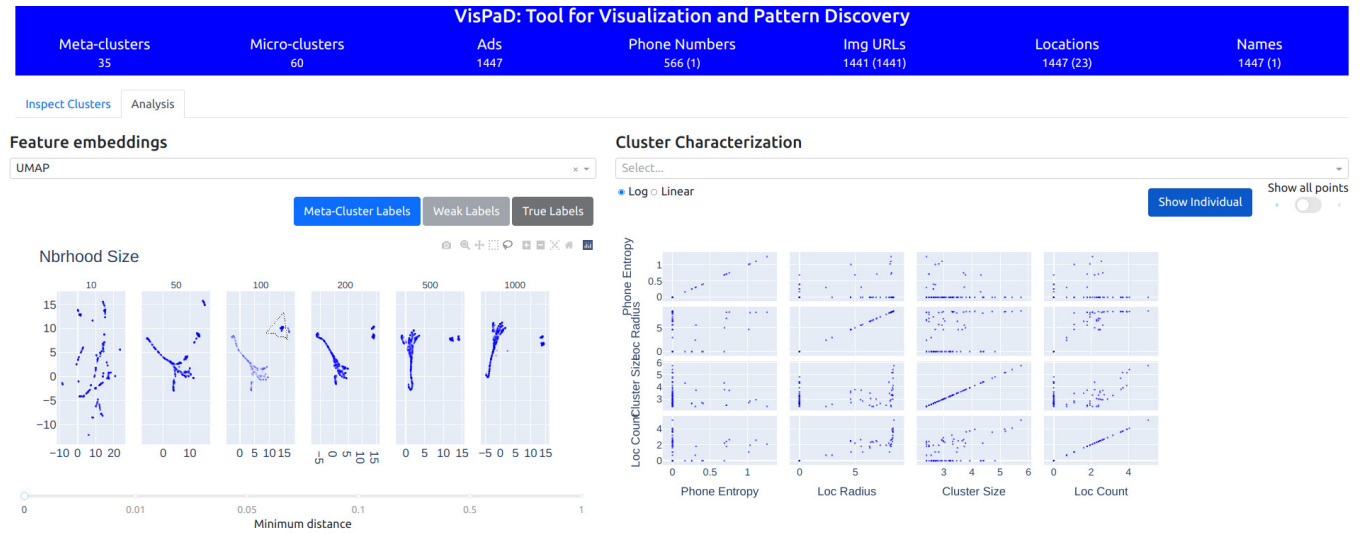


Figure 2: VisPAD showing cluster embeddings in 2D and pair-plots of cluster characterization. The most relevant 4 features out of 9 are displayed (on the right) based on the selection made by a user (indicated by dotted lines on the left side figure). Each of the subfigures and the slider at the bottom allow a user to fine-tune the parameter settings for UMAP dimensionality reduction.

Summary Panel. On the top blue banner, they see a summary of the ten largest micro-clusters of ads in that dataset², observing that there are 4 groups of meta-clusters (micro-clusters connected by meta-data) with only 5 unique phone numbers and names and 186 unique locations³.

Inspect Clusters. They then move on to the *Inspect Clusters* tab which is open on the landing page to first peruse the chart on the left titled ‘*Ads posted over time*’. They may hover over the colored block of micro-cluster C148 (which has the longest posting period) to see a tooltip displaying the number of ads posted between July 2020 and September 2021. They can further track how the cluster size of C148 changes in the same time range in the middle chart titled ‘*Meta-data over time*’. They may also observe other meta-data information by selecting from the dropdown menu ‘*Meta data-type*’. The analyst may also inspect the locations where the ads were posted over time by studying the rightmost bubble chart titled ‘*Geographical spread of ads*’. They may notice that the size of a bubble indicates how many ads were posted in that location and hovering over a it, gives them the name of the location. They can drag the slider to choose the time-period of their interest. Next they may read the scrollable ‘*Ad description*’ panel at the bottom of the page which displays the text present in the ads. They may notice the repeated mention of a Thai massage salon and together with the nation-wide spread of ads and very few number of associated phone numbers, they conclude that the largest groups of micro-clusters cover massage spas.

Analysis.

- (1) *Feature embedding visualization* in 2-dimensions (though ICA, UMAP or TSNE) allows identification of latent patterns.

²Due to the sensitive nature of the data, it cannot be publicly shared. For this paper, the data used was crawled by us from a popular escort website based in Quebec, Canada.

³The number within brackets indicates unique count

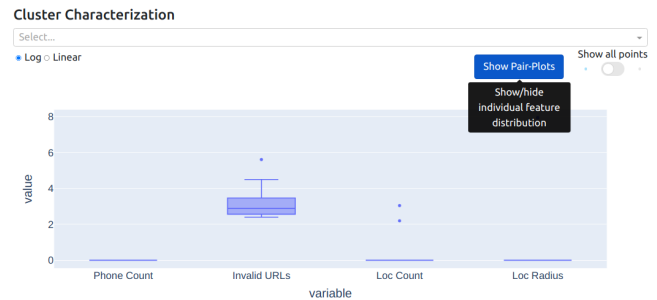


Figure 3: Caption

Also allows to choose between different parameter settings for UMAP and TSNE to control the visualization style. Give example of interesting pattern from UMAP/TSNE

- (2) *Compare across different groupings of micro-clusters* based on shared meta-data or weak labels. This allows a user to consider multiple kinds of labeling strategies before making conclusions.
- (3) *Cluster characterization pair-plots* in log/linear scale allows the user to visualize/identify patterns between pairs of cluster features and inspect interesting ones further.
- (4) *Displaying the most relevant features to a selected set of micro-clusters.* If a certain set of micro-clusters form a pattern in the low-dimensional space, which high-dimensions features contribute most to it? A classic Linear Regression model is applied to identify the top 4 most important features that separate the selected micro-clusters from the rest.

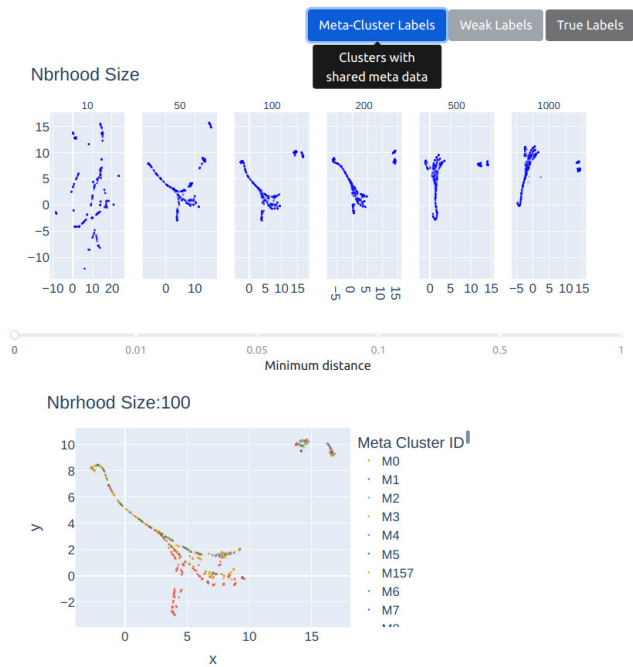


Figure 4: Caption

3 CONCLUSION

REFERENCES

- [1] International Labour Office. 2012. ILO Global Estimate of Forced Labour. http://www.ilo.org/wcmsp5/groups/public/---ed_norm/---declaration/documents/publication/wcms_182004.pdf.
- [2] Meng-Chieh Lee, Catalina Vajiac, Aayushi Kulshrestha, Sacha Levy, Namyong Park, Cara Jones, Reihaneh Rabbany, and Christos Faloutsos. 2021. InfoShield: Generalizable Information-Theoretic Human-Trafficking Detection. In *IEEE ICDE*. IEEE.
- [3] Lin Li, Olga Simek, Angela Lai, Matthew P. Daggett, Charlie K. Dagli, and Cara Jones. 2018. Detection and Characterization of Human Trafficking Networks Using Unsupervised Scalable Text Template Matching. In *IEEE BigData*. IEEE, 3111–3120.
- [4] SELL DOMESTIC MINOR. 2015. A REPORT ON THE USE OF TECHNOLOGY TO RECRUIT, GROOM AND SELL DOMESTIC MINOR SEX TRAFFICKING VICTIMS. (2015).
- [5] Andreas Olligschlager and Cara Jones. 2021. Personal communication. (2021).
- [6] Reihaneh Rabbany, David Bayani, and Artur Dubrawski. 2018. Active Search of Connections for Case Building and Combating Human Trafficking. In *KDD*. ACM, 2120–2129.
- [7] Thorn. 2015. Report on the Use of Technology to Recruit, Groom and Sell Domestic Minor Sex Trafficking Victims. https://2715111qnwey246mkc1vzqg0-wpengine.netdna-ssl.com/wp-content/uploads/2015/02/Survivor_Survey_r5.pdf.
- [8] Catalina Vajiac, Andreas Olligschlager, Yifei Li, Pratheeksha Nair, Meng-Chieh Lee, Namyong Park, Reihaneh Rabbany, Duen Horng Chau, and Christos Faloutsos. 2021. TRAFFICVIS: Fighting Human Trafficking through Visualization. (2021).