# TraffickCam: Explainable Image Matching For Sex Trafficking Investigations

**Abby Stylianou[1], Richard Souvenir[2] and Robert Pless[3]**

[1]Saint Louis University, [2]Temple University, [3]George Washington University

abby.stylianou@slu.edu,souvenir@temple.edu,pless@gwu.edu

## Abstract

Investigations of sex trafficking sometimes have access to photographs of victims in hotel rooms. These images directly link victims to places, which can help verify where victims have been trafficked or where traffickers might operate in the future. Current machine learning approaches give promising results in image search to find the matching hotel. This paper explores approaches to make this end-to-end system better support government and law enforcement requirements, including improved performance, visualization approaches that explain what parts of the image led to a match, and infrastructure to support exporting the results of a query.

## Introduction

Modern large-scale image matching approaches offer opportunities to search visual information at scales unimaginable just a few years ago. It now is possible to automatically search through databases containing billions of images, not just to retrieve generically similar images (e.g. Google Image Search), but also to find and recognize specific people and places. This technology provides some incredible opportunities (Zheng et al. June 2009) and also challenges societal expectations of anonymity and privacy (Oh et al. 2016).

When these technologies are used for official business in the public sector, there is an increasing need to provide explanations of *why* machine learning algorithms give the results they do. In the case of image analysis, substantial work has been done to visualize the image regions responsible for classification results, but much less work has explored approaches for image matching.

This paper details work building tools for investigators to determine the hotel where pictures of sex trafficking victims were taken. Our system, TraffickCam, has found that carefully implemented baseline approaches can be effective at this challenging problem (Stylianou et al. 2015; Stylianou, Souvenir, and Pless 2017; Stylianou et al. 2019). This paper shares work in progress to improve the performance of this image search and visualize what parts of the image are most important for an image match, in order to give investigators greater insight into why the system is suggesting a particular result, and therefore, perhaps, if the system should be trusted. This supports advice from a recent formal study of how human trafficking investigators interact with modern software support tools (Deeb-Swihart, Endert,
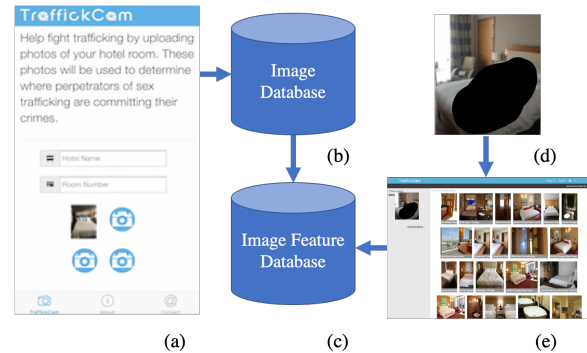


Figure 1: The TraffickCam system consists of a smartphone application (a) to collect relevant imagery of hotel rooms to support investigations of human trafficking. To date, the TraffickCam image database (b) consists of over 3 million hotel room images collected from this application and other publicly available sources of hotel room photos. Deep metric learning is used to convert these images into a searchable index of image features (c). This index supports a law enforcement platform where investigators can upload masked off images from trafficking investigations (d) to retrieve the most similar images in the TraffickCam database (e).

and Bruckman 2019) that concludes: "When designing tools for law enforcement, it is important to choose algorithms that are human interpretable and design visualizations that help officers get an intuition for how the process works."

Informal interactions with users of our system and careful examination of the problem domain highlights that learning approaches focused on whole image matching may have limited success in realistic conditions because limited parts of the room may be visible, or, where the database of hotel room images does have examples of similar rooms, they may share some but not all of the same design elements (carpet, lights, headboard, etc.). Therefore, we also explore initial efforts at expanding a visualization tool that highlights the specific elements of the hotel room images match.

We explore the development of search and visualization tools that allow an end-to-end solution to respect the needs of law enforcement and government end users, including the ability to explain and report on why the system is giving the
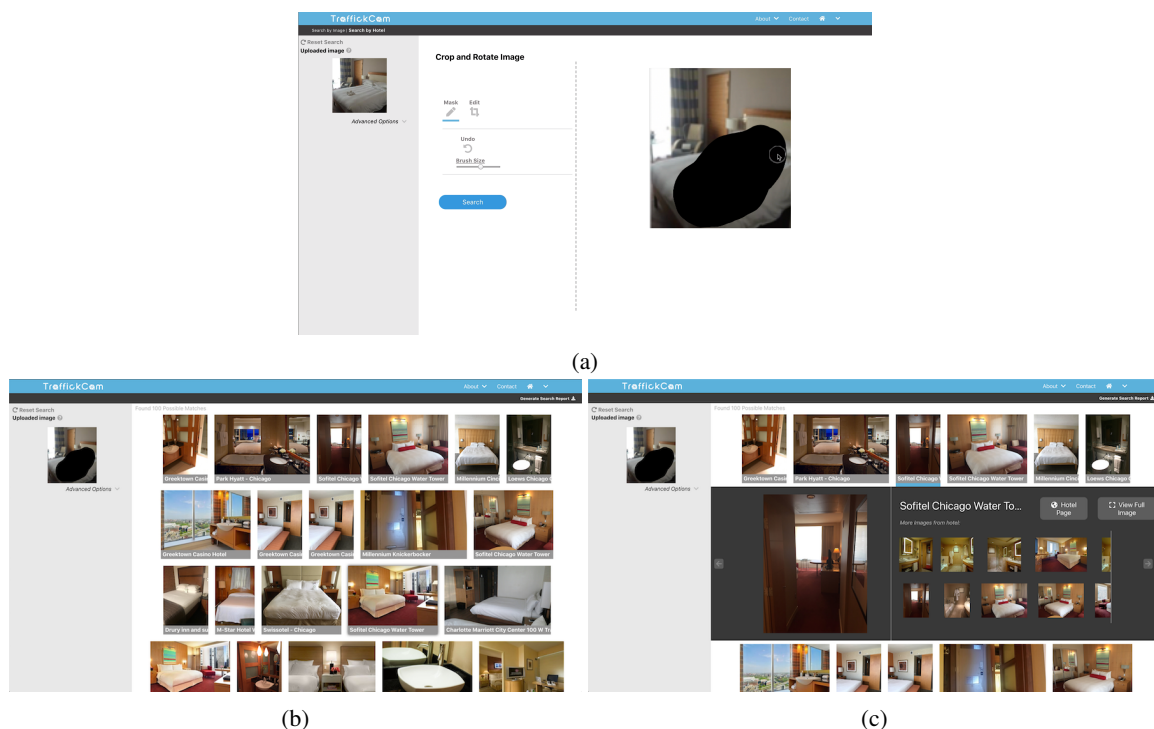
Figure 2: (a) The TraffickCam investigator search platform first allows an investigator to mask off any sensitive content in the search image. (b) The search returns the most similar images in the TraffickCam database (searches can be limited by geographic extents or by search terms). (c) Investigators can click on a relevant image to see other images at that hotel.

results that it does. Our specific contributions include:

- updating deep learning training strategies to give better, more interpretable results,

- improvements to visualization of image matching that highlight important image regions and specifically what regions of images are considered to match, and

- a search interface infrastructure that supports reporting results for investigations.

## Related Work

A limited amount of research has been published on ways to integrate machine learning (ML) to support public sector efforts to fight sex trafficking. Outside of image analysis, there is text data analysis from the advertising site Backpage (Alvari, Shakarian, and Snyder 2017) with ML to highlight online advertisements that might be related to sex trafficking, and multi-modal classification of whether online ads are offering paid sex services (Tong et al. 2017). This analysis of individual ads supports work to build knowledge graphs from online escort ads to discover connections that suggest human trafficking based on network properties (Szekely et al. 2015).

Specific work on image analysis to recognize hotels depicted in advertising images started with local image features (Stylianou et al. 2015) and progressed to deep learning features (Stylianou, Souvenir, and Pless 2017). Recently,

Hotels50k, a public dataset focused on the hotel identification problem was released (Stylianou et al. 2019). In this work, we use this dataset to explore visualization that assist law enforcement in the verification and validation of the image search results.

Visualizations of deep networks for image analysis have focused almost exclusively on classification networks (Zhang and Zhu 2018), highlighting the region of the image most responsible for the classification. In this paper, we build on one of the few approaches to visualizing similarity networks (Stylianou, Souvenir, and Pless 2019), which are the types of networks most commonly used for large-scale image matching.

## TraffickCam System Design

The overall system includes a mobile app for crowd-sourcing data (described in (Stylianou et al. 2015)), a search engine, a web-front end to support that search, and a report generation system. We also describe a collection of visualization tools that provide intuition about system behavior. These components are modular, which allows for flexible system design and component-wise updates that do not affect the rest of the system. In this section, we review improvements to the underlying search implementation, visualization tools and a system that help investigators generate reports about the search results.

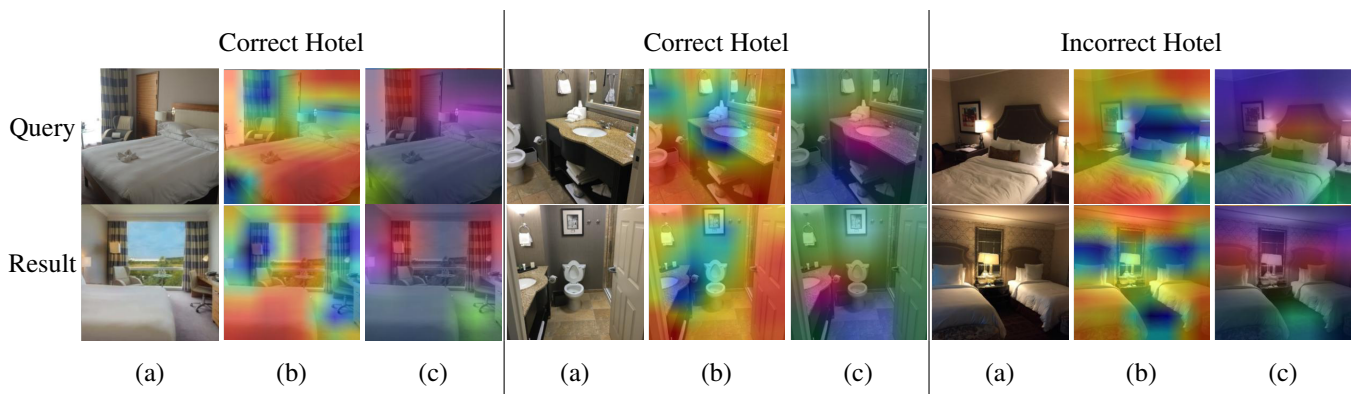|  | Correct Hotel | | | Correct Hotel | | | Incorrect Hotel | | |
|---|---|---|---|---|---|---|---|---|---|
| Query | | | | | | | | | |
| Result | | | | | | | | | |
| | (a) | (b) | (c) | (a) | (b) | (c) | (a) | (b) | (c) |

Figure 3: Example visualization results. Three examples of a query image (top), and the closest matching result image (bottom). For each, we show (b) figures that highlight which parts of image are most important for the classification (scaled from red to blue where blue is most important), and (c) a PCA based color-coding scheme that colors matching parts of the two images to be the same.

## Improved Metric Learning for Hotel Matching

State of the art approaches to large-scale image matching learn deep networks that map images to a feature space. The training optimizes the network so images from the same hotel are mapped to nearby locations in the feature space and images from different hotels are farther apart. The hotel matching problem is particularly hard, because multiple rooms from the same hotel may look very different, and even multiple images from the same room may look different (for example if one image includes the bed and the other is of the bathroom). State of the art results on the Hotels-50K dataset use Resnet-50 (He et al. 2015) with careful data augmentation and pre-processing to give about 8% accuracy in finding the correct hotel as the first result returned by the system (a factor of 4000 improvement over the chance result of 1/50000) (Stylianou et al. 2019).

The loss function with which the network was trained was a "batch-all" formulation that forces all images from a single hotel to be mapped to the same location. Because images from one hotel may look quite different, this requires the network to memorize different viewpoints of the training data and therefore does not generalize as well. Our best result have come from a refined process for training the network that focuses on "Easy Positive" triplet loss (Xuan, Stylianou, and Pless 2019). This forces the network to match only the most similar images from a hotel. This improves the generalization ability and gives about 16% accuracy in finding the correct hotel as the first response on training data.

## Improved Visualization Tools

Resnet-50 is a common deep learning architecture that includes a collection of convolutional layers following by a Global Average Pooling (GAP) layer that converts a spatial feature layout into a single vector to describe the image. In the current Traffickcam system, this final convolutional layer is $7 \times 7 \times 2048$ and the global pooling converts this to a single 2048 element vector. This final vector is used to compute the similarity between images. Previous work (Stylianou, Sou-

venir, and Pless 2019) decomposed the final similarity calculation between two 2048-element vectors to give back $7 \times 7$ resolution maps of the parts of the image that were most important for the similarity assessment.

Figure 3 shows example queries (top) and results that were shown to be similar (bottom). Visualization of which parts of an image results led them to be judged similar are shown labelled as (b), with the most important regions shown in blue. This gives a general assessment of which parts of the image are important, but does not specify specifically which components of the query image correspond to particular components of the result image.

In Figure 3, the color coded images labeled (c) are early results from a visualization approach that attempts to gain more insight about the specific relationship between images. To construct these visualizations for a pair of images that are similar, we consider the $7 \times 7 \times 2048$ convolutional layer as 49 vectors (each of length 2048) describing the local features of each image. These 49 vectors from each image are concatenated to get 98 total vectors. We can then run PCA on these vectors to get principle components (each of size 2048) that describe patterns of how the the local image description vary within and between the images. We then display the top 3 components as an RGB color mapped back to the locations on the image pair. This visualization has the property that similarly colored locations in the visualization actually have similar representations in the feature space, and therefore correspond to each other.

On the left set of results from Figure 3, this visualization approach shows that the network is matching the curtains (because both images have curtains highlighted in red), the floor (in green), and the headboard/wall light (in purple). The center image matches the floors, the sink, but not, for example, the toilet. The third set of results, on the right, are from a pair of images with a high similarity score but from different hotels. The PCA-based visualization approach highlights the similar headboard in red, but incorrectly shows that the network is mapping the linens and pillows on the bed in the query image to the carpet in the result
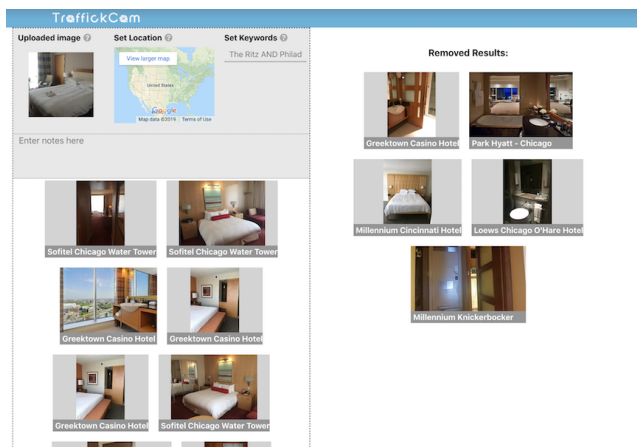
Figure 4: Investigators can select particular results from a query to include in a summary report. This report contains the masked query image, search criteria, any notes that the investigator adds, and the selected results. This report can be saved or printed directly from the website.

image (highlighted in green).

## Automatic Report Generation

In the context of an investigative process, the result of using such a system is a report to be shared with law enforcement agents or other stakeholders who may not be knowledgeable in machine learning or image retrieval. Thus, it is important that the final report is self-contained and convincing. Traffickcam provides functionality for summarizing and curating the results of a query. Figure 4 shows the automated report generation page, which includes the top matching images, information on the most likely hotel matches, and a space for the analyst to include notes. Additionally, the analyst has the ability to add, remove, or re-arrange the matching images returned by the automated system in order to frame the results in a way that best supports the conclusions and/or conveys the appropriate level of certainty.

## Conclusion

TraffickCam is an existing system to support sex trafficking investigations by identifying the hotel where an image was taken. It is currently in use at the National Center for Missing and Exploited Children. This report details work in progress to improve this system both in terms of system accuracy and to make it better support the public sector use cases through improved results, and increased explainability of those results.

One interesting issue is that evaluating search results in this domain is not necessarily best done by comparing results to ground truth (i.e., the evaluation metric proposed in (Stylianou et al. 2019)). The primary use of the system is to suggest possible hotel matches, which are then evaluated by an investigator. If the system matches the query to the correct hotel, but shows views of that hotel where the matching features are not obvious to the investigator, they

are likely to miss this match, even though the automated part of the search was effective. Finding ways to integrate visualization tools within the search results and evaluating the end-to-end system performance are important next steps.

## References

[Alvari, Shakarian, and Snyder 2017] Alvari, H.; Shakarian, P.; and Snyder, J. K. 2017. Semi-supervised learning for detecting human trafficking. *Security Informatics* 6(1):1.

[Deeb-Swihart, Endert, and Bruckman 2019] Deeb-Swihart, J.; Endert, A.; and Bruckman, A. 2019. Understanding law enforcement strategies and needs for combating human trafficking. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 331. ACM.

[He et al. 2015] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Deep residual learning for image recognition. *CoRR* abs/1512.03385.

[Oh et al. 2016] Oh, S. J.; Benenson, R.; Fritz, M.; and Schiele, B. 2016. Faceless person recognition: Privacy implications in social media. In *European Conference on Computer Vision*, 19–35. Springer.

[Stylianou et al. 2015] Stylianou, A.; Norling-Ruggles, A.; Souvenir, R.; and Pless, R. 2015. Indexing open imagery to create tools to fight sex trafficking. *2015 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)* 00:1–6.

[Stylianou et al. 2019] Stylianou, A.; Xuan, H.; Shende, M.; Souvenir, R.; and Pless, R. 2019. Hotels-50k: A global hotel recognition dataset. In *AAAI*.

[Stylianou, Souvenir, and Pless 2017] Stylianou, A.; Souvenir, R.; and Pless, R. 2017. Traffickcam: Crowdsourced and computer vision-based approaches to fighting sex trafficking. *2017 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*.

[Stylianou, Souvenir, and Pless 2019] Stylianou, A.; Souvenir, R.; and Pless, R. 2019. Visualizing deep similarity networks. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2029–2037. IEEE.

[Szekely et al. 2015] Szekely, P.; Knoblock, C. A.; Slepicka, J.; Philpot, A.; Singh, A.; Yin, C.; Kapoor, D.; Natarajan, P.; Marcu, D.; Knight, K.; et al. 2015. Building and using a knowledge graph to combat human trafficking. In *International Semantic Web Conference*, 205–221. Springer.

[Tong et al. 2017] Tong, E.; Zadeh, A.; Jones, C.; and Morency, L.-P. 2017. Combating human trafficking with deep multimodal models. *arXiv preprint arXiv:1705.02735*.

[Xuan, Stylianou, and Pless 2019] Xuan, H.; Stylianou, A.; and Pless, R. 2019. Improved embeddings with easy positive triplet mining. *arXiv preprint arXiv:1904.04370*.

[Zhang and Zhu 2018] Zhang, Q.-s., and Zhu, S.-C. 2018. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology & Electronic Engineering* 19(1):27–39.

[Zheng et al. June 2009] Zheng, Y.-T.; Zhao, M.; Song, Y.; Adam, H.; Buddemeier, U.; Bissacco, A.; Brucher, F.; Chua, T.-S.; and Neven, H. June, 2009. Tour the world: building a web-scale landmark recognition engine. In *IEEE Conference on Computer Vision and Pattern Recognition*.