

To show or not to show: Redacting sensitive text from videos of electronic displays

ABHISHEK MUKHOPADHYAY*, CPDM, Indian Institute of Science, Bangalore, Karnataka, India

SHUBHAM AGARWAL, PulseLabs AI, Salt Lake City, Utah, United States

PATRICK DYLAN ZWICK, PulseLabs AI, Salt Lake City, Utah, United States

PRADIPTA BISWAS, CPDM, Indian Institute of Science, Bangalore, Karnataka, India



Fig. 1. An example of redacting names, and phone numbers using Google Cloud Vision. (a) The image on the left is without redaction, while (b) the image on the right has names and phone numbers automatically redacted.

With the increasing prevalence of video recordings there is a growing need for tools that can maintain the privacy of those recorded. In this paper, we define an approach for redacting personally identifiable text from videos using a combination of optical character recognition (OCR) and natural language processing (NLP) techniques. We examine the relative performance of this approach when used with different OCR models, specifically Tesseract and the OCR system from Google Cloud Vision (GCV). For the proposed approach the performance of GCV, in both accuracy and speed, is significantly higher than Tesseract. Finally, we explore the advantages and disadvantages of both models in real-world applications.

CCS Concepts: • **Computing methodologies** → **Natural language processing**; *Information extraction*; *Object recognition*.

Additional Key Words and Phrases: Redaction, Optical Character Recognition, Named Entity Recognition, Google Cloud Vision, Tesseract

ACM Reference Format:

Abhishek Mukhopadhyay, Shubham Agarwal, Patrick Dylan Zwick, and Pradipta Biswas. 2022. To show or not to show: Redacting sensitive text from videos of electronic displays. In *AutoUI'22*, Seoul, South Korea, 6 pages. <https://doi.org/3544999.3552529>

*This work was carried out during the internship at Pulse Labs AI.

Correspondence: abhishekmkh@iisc.ac.in, shubham.agarwal@pulselabs.ai, dylan.zwick@pulselabs.ai, pradipta@iisc.ac.in

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.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

1 INTRODUCTION

The massive improvements over the past few decades in digital camera technology, ranging from reductions in camera costs to enhancements in image quality to greater efficiency in video storage and transmission, have led to an enormous increase in the amount of available video data. This, combined with major developments in computer vision and machine learning technology, has created enormous opportunities to make life better through the collection and utilization of this video data. Potential applications here range from improved security to interactive entertainment. However, the collection and utilization of this data also entails ethical privacy concerns and the potential for unwanted intrusion into people's lives without their permission. One way to attempt to achieve the benefits of more omnipresent video collection while mitigating the intrusion on privacy is through the automatic redaction of personally identifiable information (PII). This means automatically removing or obscuring content from video data that can be used to identify an individual while maintaining as much other video data as possible.

A relatively new context generating a significant amount of video data is the cabins of automobiles. With the emergence of automatically responsive infotainment systems, there is the potential to improve the safety and overall driver experience by quickly integrating data from driver movement into what the infotainment system displays. However, building and improving algorithms to achieve this requires the recording and analysis of a significant amount of in-cabin video data, which drivers may be uncomfortable sharing. One way to decrease this lack of comfort is through the removal of PII from these in-cabin videos.

Within in-cabin videos, textual PII mostly appears on the infotainment system display. For example, if the infotainment system is used to call a friend or navigate to a family member's house, the system could display PII as part of this task. The goal of this paper is to investigate how this type of PII can be automatically removed without removing other parts of the captured data.

The data used for this investigation comes from a set of usability studies conducted by Pulse Labs AI¹. Pulse Labs is a company that specializes in using AI techniques to improve driver safety and experience. For these usability studies, drivers pointed cameras at their infotainment systems and recorded themselves performing common tasks like calling a friend or navigating to a business. Image (a) from Figure 1 above is an example frame from one of these videos.

Earlier research [1, 2, 5, 6] on removing PII from videos has focused on anonymizing faces. In contrast, this paper focuses on redacting PII (like names, email addresses, or phone numbers) from text appearing in videos. Figure 1 shows an example of this process in action.

Our contributions in this work can be summarized as:

- We extend image redaction techniques by applying them to electronic display system interaction videos using a combination of optical character recognition (OCR) and natural language processing (NLP).
- We use real-world datasets to investigate the redaction of sensitive textual PII from electronic display system videos.
- We provide a systematic comparison of the performance of two popular OCR models, Tesseract and GCV, in the context of in-cabin video interactions.

2 PROPOSED APPROACH

Our approach aims to remove sensitive information from the electronic display system present within videos on a frame-by-frame basis. In this procedure, we first use an OCR model to identify text, and the box-coordinates of that

¹<https://pulselabs.ai>

text, within a frame. Next, we use NLP techniques to identify the presence of entities within the text, in our case names, phone numbers, and email addresses. Specifically, we use regular expressions² (RegEx) to identify email addresses and phone numbers based on pattern matching, and deep learning based named-entity recognition (NER) to detect names. Finally, we apply redaction by adding green boxes at the identified coordinates of different entities. All the individual redacted frames are converted back to video by replacing their corresponding frames in the original unredacted video.

In our research, we experimented with two OCR models - Tesseract³ [4], and Google Cloud Vision⁴. The NER tool used in our approach was spaCy⁵ [3] in Python. Figure 2 provides an overview of our complete pipeline.

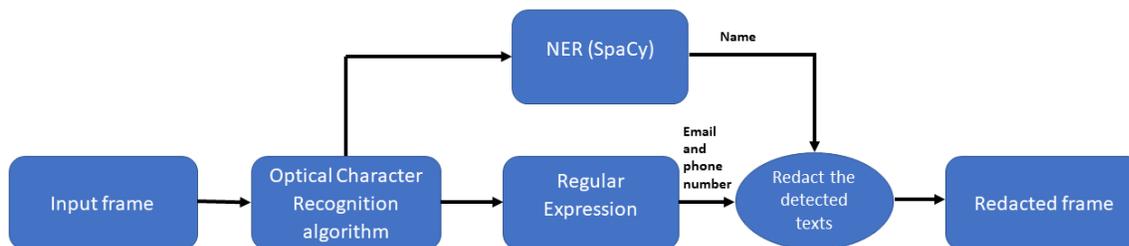


Fig. 2. Pipeline for proposed redaction method.

An example of a frame with redaction performed using Tesseract for OCR, using GCV for OCR, and using manual identification, is provided in Figure 3. Here, we find that while the OCR from GCV (See Figure 3c) is able to correctly identify and anonymize both the name and phone number, the Tesseract OCR (Figure 3d) fails to identify the name, leading to partial redaction and insufficient anonymization.

3 RESULTS AND DISCUSSION

The data for this analysis was taken from different sets of videos. Each video set featured user interactions that tended to surface textual PII of a specific type. These three types were names, phone numbers, and email addresses. From each of these sets of videos, a test dataset was created by randomly sampling 100 frames from the videos within the set. Each frame was then examined manually, and the textual entity of interest was identified and marked. The resulting datasets contained 209, 183, and 214 instances of names, phone numbers, and email addresses respectively.

We then applied our PII redaction approach to each frame twice, once using Tesseract and once using GCV. We compared the PII identifications created with the proposed approach with those created manually, and evaluated the performance of the proposed approach using metrics such as accuracy, precision, and recall. Note that for any frame it was possible to get a count of true positives, false positives, and false negatives, but not true negatives, as the number of times our approach correctly did not detect PII is not well defined. Consequently, for this analysis accuracy is defined as $TP/(TP + FP + FN)$, where TP , FP , and FN indicate true positive, false positive, and false negative, respectively. Precision, $TP/(TP + FP)$, and recall, $TP/(TP + FN)$, are defined as usual.

Across the three datasets, we observed that GCV (mean accuracy: 91.21%) obtained a significantly higher accuracy than Tesseract (mean accuracy: 54.21%) (Figure 4a). For the three types of textual PII we focused on, the accuracy of

²Regular expressions are a way to describe a pattern of text for which we're looking. For example, to detect the pattern $x@y.z$, where x , y , and z are words, we use the regular expression $[\wedge\wedge\d] + \@[\wedge\wedge\d] + \.[\wedge\wedge\d]+$.

³<https://github.com/tesseract-ocr/tesseract>

⁴<https://cloud.google.com/vision>

⁵<https://spacy.io/>

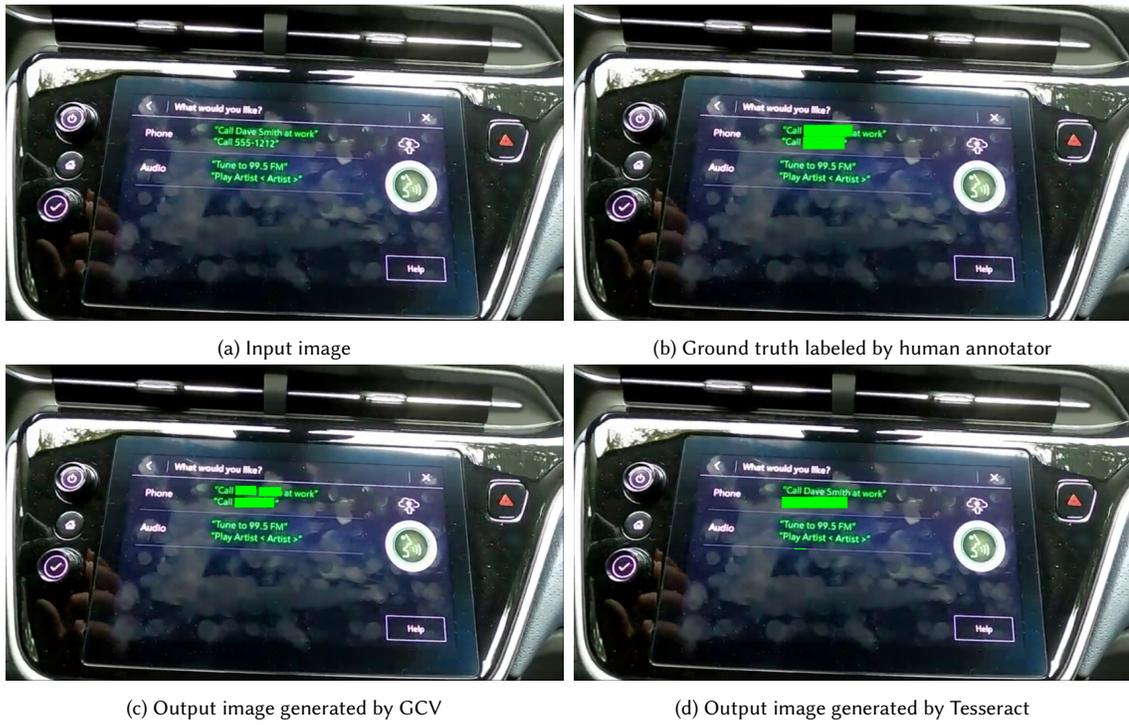


Fig. 3. Overview of image samples and corresponding outputs.

Tesseract was much lower than GCV (Figure 4b) except on email addresses, where the two models were much closer. The precision and recall of the two OCR models for every type of textual PII are given in Table 1.

Finally, we compared the average frame processing speed for each OCR model, and found that GCV took on average 1.09 seconds to process a frame, while Tesseract on average took 1.25 seconds. All processing was done using NVIDIA RTX 2070 processor. Note the proposed model is not limited to detecting textual PII within automotive infotainment systems, and could be applied without changes to other types of video like mobile screen recordings (see Figure 5).

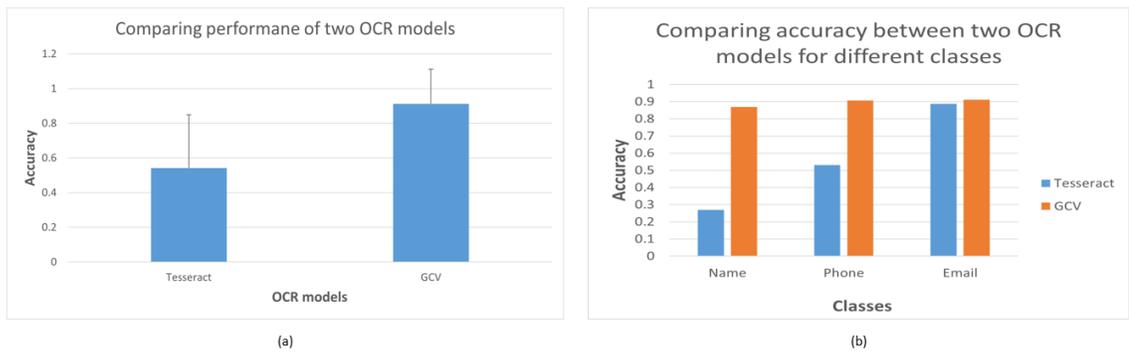


Fig. 4. Accuracy analysis graphs. (a) Overall accuracy between two OCR models; (b) Class-wise performance.

Table 1. Precision and Recall table for two OCR models

Classes	Tesseract		GCV	
	Precision	Recall	Precision	Recall
Email address	1	0.89	1	0.91
Name	0.42	0.43	0.88	0.98
Phone	0.90	0.56	1	0.91

Tesseract performed significantly worse than GCV, and also required more processing time per frame. We believe the difference in accuracy is because GCV was much better at separating punctuation within names and phone numbers (Figure 1). This deficiency in Tesseract significantly impeded the ability of RegEx-based pattern matching or spaCy-based NER to detect the textual PII of interest. However, this issue was not as prevalent within email addresses, which can explain why the two models were much closer for that type of textual PII. In addition, Tesseract struggled to accurately detect textual information within blurred images.

However, Tesseract is not entirely without virtues. It is open-sourced, free to use software that can be downloaded and deployed offline, while Google charges for the use of the GCV API, which can only be accessed online. If cost or offline access are priorities, Tesseract may still have some appeal.

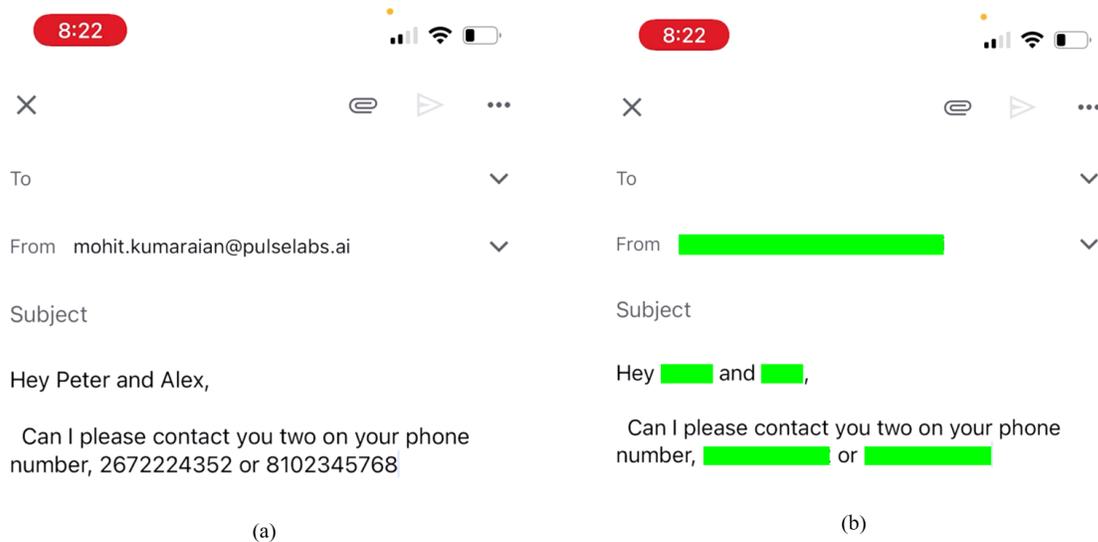


Fig. 5. An example of redacting names, email addresses and phone numbers using Google Cloud Vision on mobile screen recording.

4 CONCLUSION

In this paper, we proposed an approach for the automated detection of textual PII within recorded videos, based upon the use of OCR and NLP techniques. Our initial results show that we can successfully redact most of the textual entities in the recorded videos with a high accuracy. For the two OCR models, Tesseract and GCV, compared in this paper the performance of GCV with the proposed approach was significantly better. Our results are an initial step towards the

automatic removal of PII from in-cabin automotive video data, which will be valuable as the need to respect driver privacy while improving the computer vision models used by automotive infotainment systems become increasingly important.

ACKNOWLEDGMENTS

We would like to thank Abhishek Suthan, Mohit Kumaraian, Mohan Karthik and the Pulse Labs team for the fruitful discussions during the initial part of the project.

REFERENCES

- [1] Oran Gafni, Lior Wolf, and Yaniv Taigman. 2019. Live face de-identification in video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9378–9387.
- [2] Anjith George and Sébastien Marcel. 2021. On the effectiveness of vision transformers for zero-shot face anti-spoofing. In *2021 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–8.
- [3] Matthew Honnibal and Ines Montani. 2018. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. (2018). To appear.
- [4] Anthony Kay. 2007. Tesseract: an open-source optical character recognition engine. *Linux Journal* 2007, 159 (2007), 2.
- [5] Zhongzheng Ren, Yong Jae Lee, and Michael S Ryoo. 2018. Learning to anonymize faces for privacy preserving action detection. In *Proceedings of the european conference on computer vision (ECCV)*. 620–636.
- [6] Yuanyuan Xu, Wan Yan, Genke Yang, Jiliang Luo, Tao Li, and Jianan He. 2020. CenterFace: joint face detection and alignment using face as point. *Scientific Programming* 2020 (2020).