# Deceiving Google's Cloud Video Intelligence API Built for Summarizing Videos

Hossein Hosseini     Baicen Xiao     Radha Poovendran
Network Security Lab (NSL)
Department of Electrical Engineering, University of Washington, Seattle, WA
Email: {hosseinh, bcxiao, rp3}@uw.edu

arXiv:1703.09793v2 [cs.CV] 31 Mar 2017

## Abstract

*Despite the rapid progress of the techniques for image classification, video annotation has remained a challenging task. Automated video annotation would be a breakthrough technology, enabling users to search within the videos. Recently, Google introduced the Cloud Video Intelligence API for video analysis. As per the website, the system can be used to "separate signal from noise, by retrieving relevant information at the video, shot or per frame" level. A demonstration website has been also launched, which allows anyone to select a video for annotation. The API then detects the video labels (objects within the video) as well as shot labels (description of the video events over time).*

*In this paper, we examine the usability of the Google's Cloud Video Intelligence API in adversarial environments. In particular, we investigate whether an adversary can subtly manipulate a video in such a way that the API will return only the adversary-desired labels. For this, we select an image, which is different from the video content, and insert it, periodically and at a very low rate, into the video. We found that if we insert one image every two seconds, the API is deceived into annotating the video as if it only contained the inserted image. Note that the modification to the video is hardly noticeable as, for instance, for a typical frame rate of 25, we insert only one image per 50 video frames. We also found that, by inserting one image per second, all the shot labels returned by the API are related to the inserted image. We perform the experiments on the sample videos provided by the API demonstration website and show that our attack is successful with different videos and images.*

## 1. Introduction

In recent years, machine learning techniques have been extensively deployed for computer vision tasks, particularly recognizing objects in images [1–4]. However, using machine learning for annotating videos has remained a challenging task, due to the temporal aspect of video data and the difficulty of collecting sufficient well-tagged training
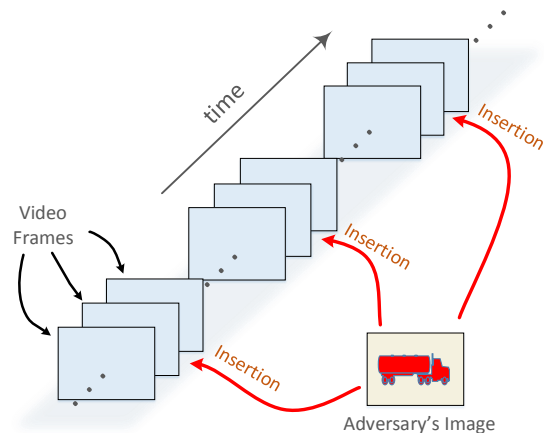


Figure 1: Illustration of the image insertion attack on Google's Cloud Video Intelligence API. The adversary chooses an image and inserts it, periodically and at a very low rate, into the video. Our experimental results show that, by inserting the image once every second, we can deceive the API to output only the labels of the inserted image for both the video and shot labels.

samples [5–7]. Due to the growth of video data on the Internet, automatic video annotation has gained a lot of attention from the research community [8–10], as well as companies such as Facebook [11] and Twitter [12]. Automatic video annotation can enable searching the videos for a specific event, which is helpful in applications such as video surveillance or returning the search results on the web. It can be also used for prescanning user videos, for example in YouTube and Facebook, where distribution of certain types of illegal contents is not permitted.

Recently, Google introduced the Cloud Video Intelligence API for video analysis [13]. A demonstration website has been launched which allows anyone to select a video stored in Google Cloud Storage for annotation [14]. The API then quickly identifies the *video labels* which are the key objects within the video. It also detects the scene changes and provides *shot labels* as the detailed description

1

of the video events over time. Similar to other Google's machine learning APIs, the Cloud Video Intelligence API is made available to developers to build applications that can automatically search within the videos [13]. Hence, the API has the potential to simplify the video understanding and enable searching in videos just as text documents.

Machine learning systems are typically designed and developed with the implicit assumption that they will be deployed in benign settings. However, many works have pointed out their vulnerability in adversarial environments [15–18]. Security evaluation of machine learning systems is an emerging field of study. In [19], Carlini et al. showed that voice interfaces can be attacked with hidden voice commands that are unintelligible to humans, but are interpreted as commands by devices. In [20], Sharif et al. proposed techniques for physically realizable image modification to attack face-recognition systems. Recently, Hosseini et al. showed that the Google's Perspective API for detecting toxic comments can be defeated by subtly modifying the input text [21].

In this paper, we examine the usability of the Google's Cloud Video Intelligence API in adversarial environments. In particular, we investigate whether an adversary can deceive the API into returning *only* the adversary-desired labels, by slightly manipulating the input video. Such vulnerability will seriously undermine the performance of the video annotation system in real-world applications. For example, a search engine may wrongly suggest manipulated videos to users, or a video filtering system can be bypassed by slightly modifying a video which has illegal contents.

For manipulating the videos, we select an image, different from the video content, and insert it, periodically and at a very low rate, into the video. Our experimental results show that by inserting the image once every two seconds, the API is deceived into returning only the video labels which are related the inserted image. Note that the modification to the video is hardly noticeable as, for instance, for a typical frame rate of 25, we insert only one image per 50 video frames. We also found that by inserting one image per second, all the shot labels returned by the API are related to the inserted image. We perform the experiments on the sample videos provided by the API demonstration website and with different images. Figure 1 illustrates the image insertion attack on the Google's Cloud Video Intelligence API.

## 2. Google's Cloud Video Intelligence API

The Google's Cloud Video Intelligence API is designed for video understanding and analysis. It enables the developers to easily search and discover the video content by providing information about entities (nouns or verbs) in the video and when they occur within the video. It was noted in [14] that the system can be used to "separates signal from
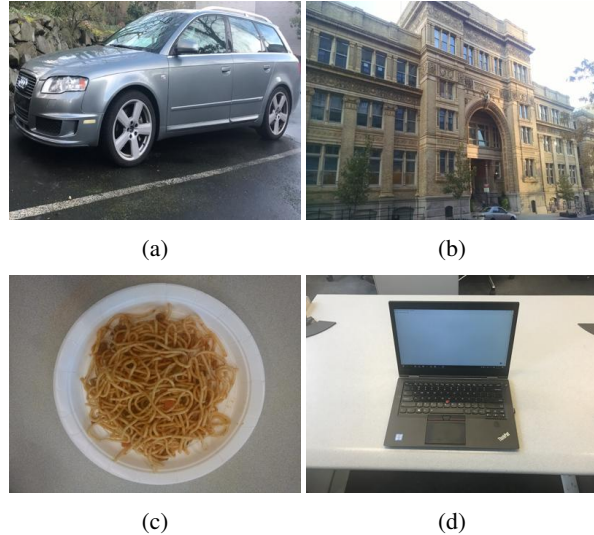


Figure 2: The four images (a) a car, (b) a building, (c) a food plate, and (d) a laptop, that were used in experiments for inserting within the sample videos.

noise, by retrieving relevant information at the video, shot or per frame" level. The API uses deep-learning models, built using frameworks such as TensorFlow and applied on large-scale media platforms such as YouTube [13].

The system is said to be helpful for large media companies to better understand the video data, and for media organizations and consumer technology companies, who want to build their media catalogs or find easy ways to manage crowd-sourced content [13]. The underlying technology can be also used to improve the video recommendations, as it enables the search engines to consider the video content, beyond the metadata like descriptions and comments, for searches.

## 3. The Image Insertion Attack

In this section, we describe the image insertion attack for deceiving the Google's Cloud Video Intelligence API. The goal of the attack is to modify a given video in such a way that a human observer would perceive its original content, but the API returns only the adversary-desired annotations. We performed the experiments with three sample videos "Animals.mp4", "GoogleFiber.mp4" and "Jane-Goodall.mp4", which are provided by the demonstration website of the Google's Cloud Video Intelligence API [14]. The API provides video labels (objects in the entire video), shot changes (scene changes within the video) and shot labels (description of the video events over time).

The attack procedure is as follows. We first tested the API with sample videos and verified that the API did indeed accurately detect both the video and shot labels. For

Table 1: Demonstration of the Image Insertion Attack on the Google's Cloud Video Intelligence API. We performed the experiments with three sample videos provided by the API website [14]. The images are inserted once every two seconds within the video, which is equal to inserting one image per 50 video frames for a typical frame rate of 25. For each of the input videos and inserted images, the table shows the video label with the highest confidence returned by the API.

| Video Name | Inserted Image | Video Label Returned by API (Confidence Score) |
|---|---|---|
| **"Animals.mp4"** | "Car" | Audi (98%) |
| | "Building" | Building (89%) |
| | "Food Plate" | Pasta (99%) |
| | "Laptop" | Laptop (91%) |
| **"GoogleFiber.mp4"** | "Car" | Audi (98%) |
| | "Building" | Classical architecture (95%) |
| | "Food Plate" | Noodle (99%) |
| | "Laptop" | Laptop (91%) |
| **"JaneGoodall.mp4"** | "Car" | Audi (98%) |
| | "Building" | Classical architecture (95%) |
| | "Food Plate" | Pasta (99%) |
| | "Laptop" | Laptop (91%) |

example, for the "Animals.mp4" video, the API returns the video labels "Animal," "Wildlife," "Zoo," "Terrestrial animal," "Nature," "Tourism," and "Tourist destination," which are consistent with the video content.

We then downloaded the sample videos and modify them. For manipulating the videos, we select an image, different from the video content, and insert it, periodically and at a very low rate, into the videos. Figure 2 shows the four images that were used for image insertion attack, namely, a car, a building, a food plate and a laptop. The schematic of the image insertion attack is illustrated in Figure 1. At the end, we stored the manipulated videos on the Google cloud storage and used them as inputs to the API. [1]

Our experimental results show that if we insert an image periodically once every two seconds and in appropriate places, the API completely fails to correctly understand the video content and annotates it as if the video was only about the inserted image. Note that the image insertion rate is very low. That is, for a typical frame rate of 25, we insert only one image per 50 video frames, resulting in an image insertion rate of 0.02. Therefore, the modification to the video is hardly noticeable. Moreover, we tested the API with videos with different frame rates and verified that the attack is successful, regardless of the choice of the frame rate.

Table 1 provides the API's output for the video labels (the table shows only the label with the highest confidence

score). As can be seen, regardless of the video content, the API returns a video label, with a very high confidence score, that exactly matches the corresponding inserted images. Figure 3 shows the results in more details, providing the screenshots of the video annotations for the sample video "Animals.mp4" and the four versions, each manipulated with one of the images presented in Figure 2. The results show that, while the API can accurately annotate the original video, for the manipulated videos it *only* outputs the labels which are related to the inserted image. Figures 4 and 5 show similar experiments with the "GoogleFiber.mp4" and "JaneGoodall.mp4" videos, respectively.

We performed similar experiments for changing the video shot labels returned by the API. Note that shot labels provide a detailed description of the individual scenes within the video; therefore, compared to changing the video labels, it is more challenging to change all the shot labels, while maintaining a low image insertion rate. However, we found that by inserting one image per second, resulting in an image insertion rate of 0.04 for the frame rate of 25, *all* the shot labels returned by the API are related to the inserted image. Figures 6 shows the screenshots of the shot labels for the original video "Animals.mp4" and the four manipulated versions, each with one of the inserted images. While the figure shows the results only for one shot, we verified that the attack succeeds to change all the shot labels to the labels of inserted image. Moreover, it can be seen that the proposed image insertion attack completely alters the pattern of the *shot changes* of the video, returned by the API.

---

[1]The experiments are performed on the interface of the Cloud Video Intelligence API's website on Mar. 24, 2017.

## 4. Discussion

Many applications can benefit from automated video search and summarization. For example, in video surveillance, one needs to search many hours of videos for a specific event. Also, some Internet platforms, such as YouTube and Facebook, require to process enormous amounts of video files every day, for video recommendation and to block the videos with illegal contents. The Google's Cloud Video Intelligence API is designed to enable the developers to quickly search the video contents, just as text documents. Hence, it has the potential to transform the video analysis field to the point that users can search for a particular event and get related videos along with the exact timings of the events within the videos.

However, we showed that the API has certain security weaknesses. Specifically, an adversary can insert an image, periodically and at a very low rate, into the video in a way that all the generated shot labels are about the inserted image. Such vulnerability seriously undermines the applicability of the API in adversarial environments. For example, one can upload a manipulated video which contains adversarial images related to a specific event, and the API wrongly suggests it to users who asked for videos from the event. Furthermore, an adversary can bypass a video filtering system by inserting a benign image into a video with illegal contents.

Note that we could deceive the Google's Cloud Video Intelligence API, without having any knowledge about the learning algorithms, video annotation algorithms or the cloud computing architecture used by the API. That is, we developed an approach for deceiving the API, by only querying the system with different inputs. Through experiments, we showed that the attack is consistently successful with different videos and images. The success of the image insertion attack shows the importance of designing the system to work equally well in adversarial environments.

## 5. Conclusion

In this paper, we showed that the Google's Cloud Video Intelligence API can be easily deceived by an adversary without compromising the system or having any knowledge about the specific details of the algorithms used. In essence, we found that an adversary can slightly manipulate a video by inserting an image periodically into it, such that the API returns only the labels that are related to the inserted image.
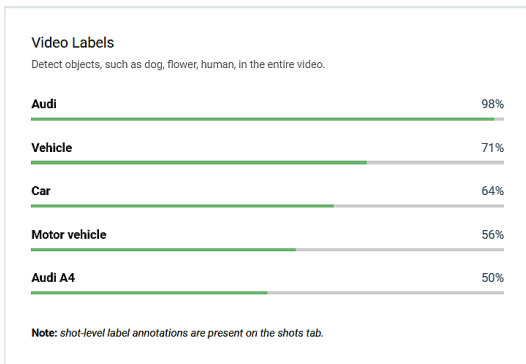
### Acknowledgments

## References

[1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, pp. 1097–1105, 2012.

[2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[3] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 580–587, 2014.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.

[5] C. Vondrick, D. Ramanan, and D. Patterson, "Efficiently scaling up video annotation with crowdsourced marketplaces," *Computer Vision–ECCV 2010*, pp. 610–623, 2010.

[6] Z.-J. Zha, M. Wang, Y.-T. Zheng, Y. Yang, R. Hong, and T.-S. Chua, "Interactive video indexing with statistical active learning," *IEEE Transactions on Multimedia*, vol. 14, no. 1, pp. 17–27, 2012.

[7] M. Wang, B. Ni, X.-S. Hua, and T.-S. Chua, "Assistive tagging: A survey of multimedia tagging with human-computer joint exploration," *ACM Computing Surveys (CSUR)*, vol. 44, no. 4, p. 25, 2012.

[8] M. Wang, X.-S. Hua, R. Hong, J. Tang, G.-J. Qi, and Y. Song, "Unified video annotation via multigraph learning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 19, no. 5, pp. 733–746, 2009.

[9] M. Wang, X.-S. Hua, J. Tang, and R. Hong, "Beyond distance measurement: constructing neighborhood similarity for video annotation," *IEEE Transactions on Multimedia*, vol. 11, no. 3, pp. 465–476, 2009.

[10] T. Zhang, C. Xu, G. Zhu, S. Liu, and H. Lu, "A generic framework for video annotation via semi-supervised learning," *IEEE Transactions on Multimedia*, vol. 14, no. 4, pp. 1206–1219, 2012.

[11] https://code.facebook.com/posts/804694409665581/powering-facebook-experiences-with-ai/.

[12] https://blog.twitter.com/2016/increasing-our-investment-in-machine-learning.

[13] https://cloud.google.com/blog/big-data/2017/03/announcing-google-cloud-video-intelligence-api-and-more-cloud-machine-learning-updates.

[14] https://cloud.google.com/video-intelligence/#demo.

[15] M. Barreno, B. Nelson, R. Sears, A. D. Joseph, and J. D. Tygar, "Can machine learning be secure?," in *Proceedings*

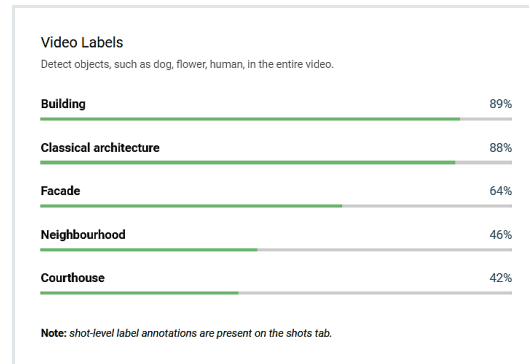*of the 2006 ACM Symposium on Information, computer and communications security*, pp. 16–25, ACM, 2006.

[16] L. Huang, A. D. Joseph, B. Nelson, B. I. Rubinstein, and J. Tygar, "Adversarial machine learning," in *Proceedings of the 4th ACM workshop on Security and artificial intelligence*, pp. 43–58, ACM, 2011.

[17] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *Security and Privacy (EuroS&P), 2016 IEEE European Symposium on*, pp. 372–387, IEEE, 2016.

[18] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, "Concrete problems in ai safety," *arXiv preprint arXiv:1606.06565*, 2016.

[19] N. Carlini, P. Mishra, T. Vaidya, Y. Zhang, M. Sherr, C. Shields, D. Wagner, and W. Zhou, "Hidden voice commands," in *25th USENIX Security Symposium (USENIX Security 16), Austin, TX*, 2016.

[20] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1528–1540, ACM, 2016.

[21] H. Hosseini, S. Kannan, B. Zhang, and R. Poovendran, "Deceiving google's perspective api built for detecting toxic comments," *arXiv preprint arXiv:1702.08138*, 2017.

(a) Video labels generated by API for the original video.



(b) Video labels of the manipulated video, where an image of a car is inserted once every two seconds.



(c) Video labels of the manipulated video, where an image of a building is inserted once every two seconds.



(d) Video labels of the manipulated video, where an image of a food plate is inserted once every two seconds.
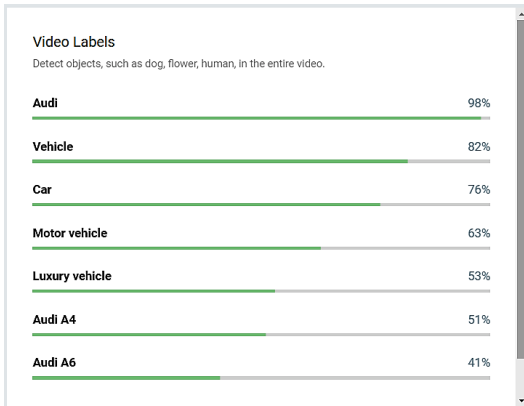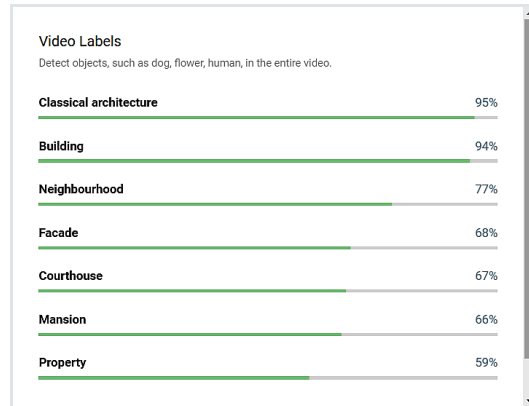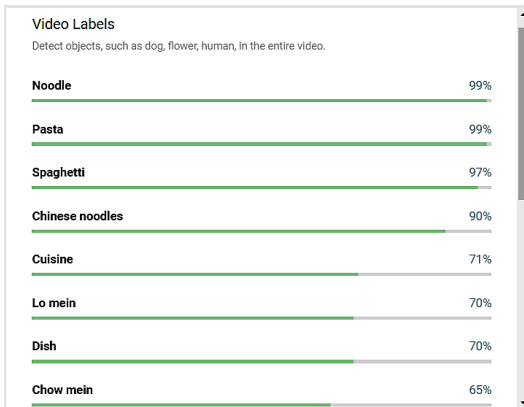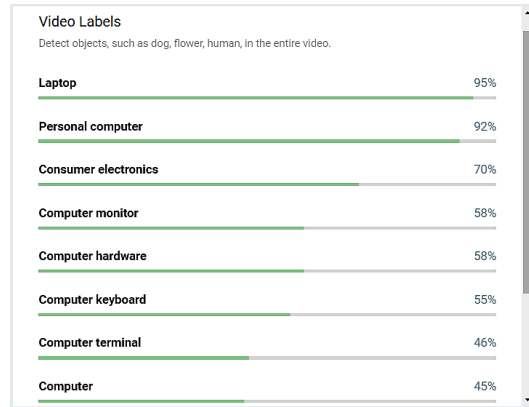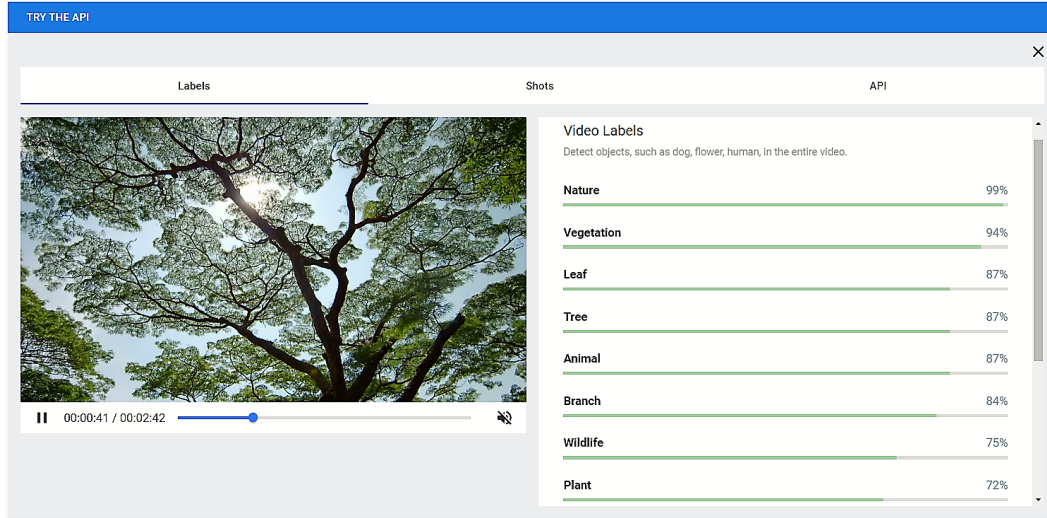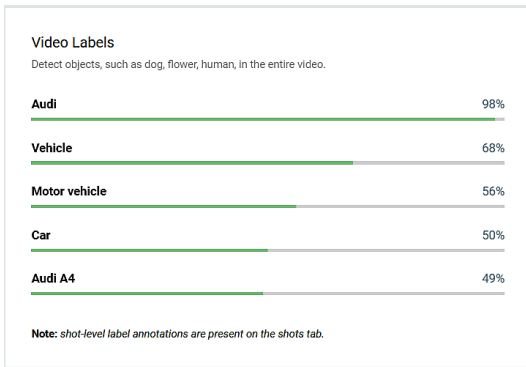


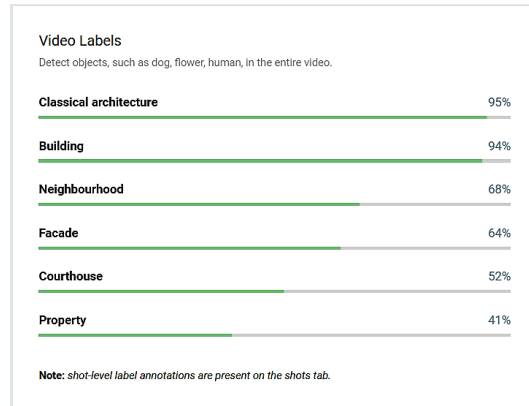(e) Video labels of the manipulated video, where an image of a laptop is inserted once every two seconds.

Figure 3: The results of the image insertion attack for changing the video labels of the sample video "Animals.mp4," provided by the demonstration website of the Google's Cloud Video Intelligence [14].
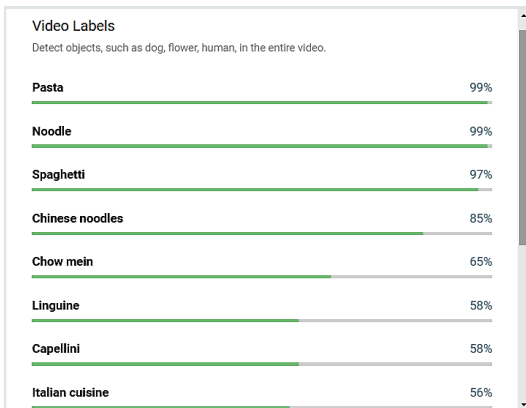
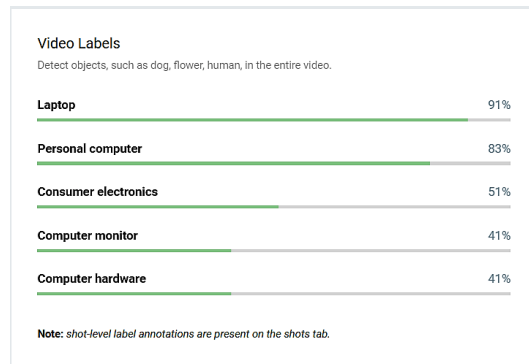(a) Video labels generated by API for the original video.



(b) Video labels of the manipulated video, where an image of a car is inserted once every two seconds.



(c) Video labels of the manipulated video, where an image of a building is inserted once every two seconds.
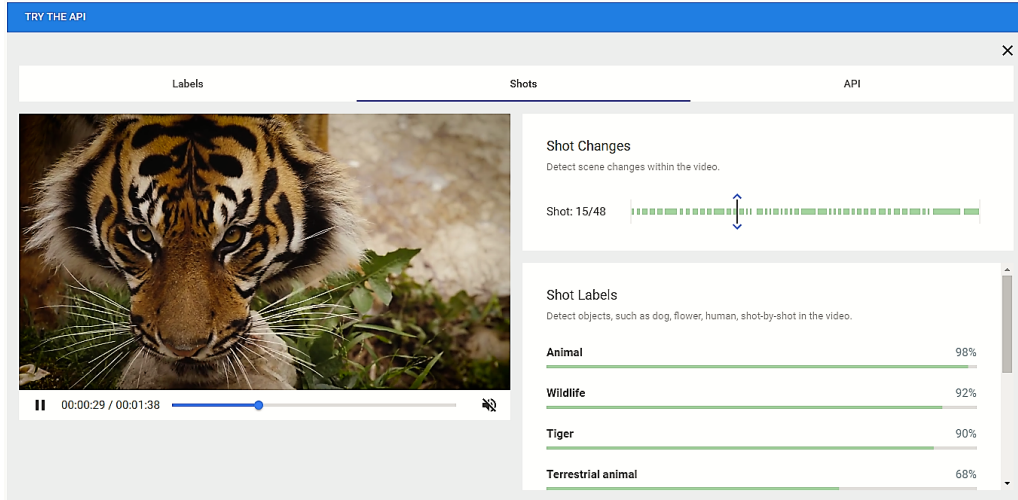


(d) Video labels of the manipulated video, where an image of a food plate is inserted once every two seconds.



(e) Video labels of the manipulated video, where an image of a laptop is inserted once every two seconds.

Figure 4: The results of the image insertion attack for changing the video labels of the sample video "GoogleFiber.mp4," provided by the demonstration website of the Google's Cloud Video Intelligence [14].

(a) Video labels generated by API for the original video.



(b) Video labels of the manipulated video, where an image of a car is inserted once every two seconds.



(c) Video labels of the manipulated video, where an image of a building is inserted once every two seconds.



(d) Video labels of the manipulated video, where an image of a food plate is inserted once every two seconds.
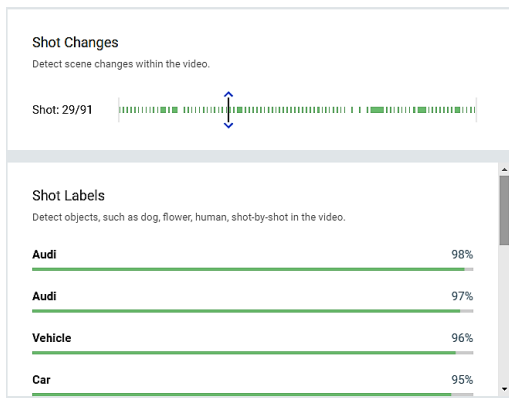


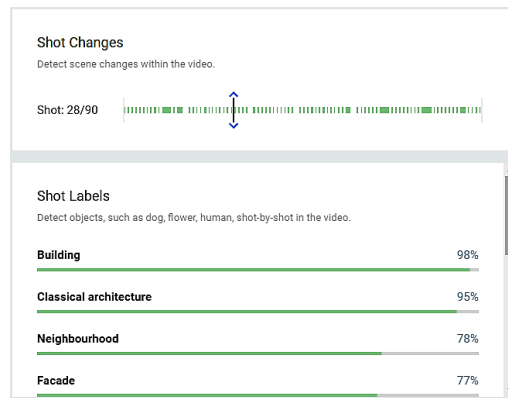(e) Video labels of the manipulated video, where an image of a laptop is inserted once every two seconds.

Figure 5: The results of the image insertion attack for changing the video labels of the sample video "JaneGoodall.mp4," provided by the demonstration website of the Google's Cloud Video Intelligence [14].
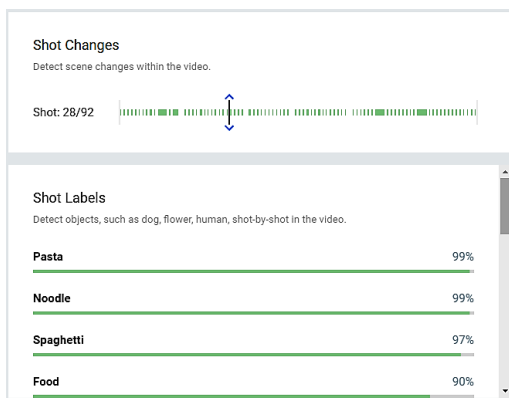
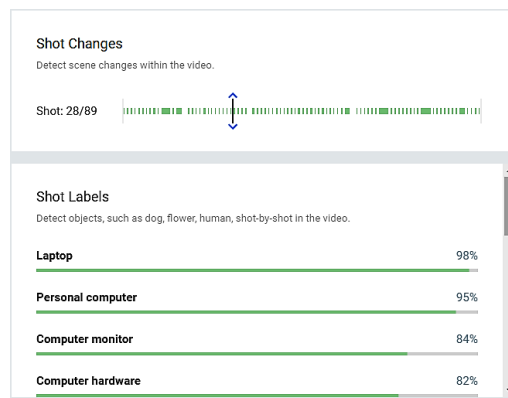(a) Shot labels generated by API for the original video.



(b) Shot labels of the manipulated video, where an image of a car is inserted once every second.



(c) Shot labels of the manipulated video, where an image of a building is inserted once every second.



(d) Shot labels of the manipulated video, where an image of a food plate is inserted once every second.



(e) Shot labels of the manipulated video, where an image of a laptop is inserted once every seconds.

Figure 6: The results of the image insertion attack for changing the shot labels of the sample video "Animals.mp4," provided by the demonstration website of the Google's Cloud Video Intelligence [14]. While the figures shows the results only for one shot, we verified that the attack succeeds to change *all* the shot labels to the labels of inserted image. Note that the periodic image insertion also completely alters the *shot changes* of the video, returned by the API.

9