

Using Deep Learning to Solve Google reCAPTCHA v2's Image Challenges

Dylan Wang
Department of Computer Science
San José State University
San José, CA, USA
dylan.wang@sjsu.edu

Melody Moh
Department of Computer Science
San José State University
San José, CA, USA
melody.moh@sjsu.edu
<https://orcid.org/0000-0002-8313-6645>

Teng-Sheng Moh
Department of Computer Science
San José State University
San José, CA, USA
teng.moh@sjsu.edu
<https://orcid.org/0000-0002-2726-102X>

Abstract—The most popular CAPTCHA service in use today is Google reCAPTCHA v2, whose main offering is an image-based CAPTCHA challenge. This paper looks into the security measures used in reCAPTCHA v2's image challenges and proposes a deep learning-based solution that can be used to automatically solve them. The proposed method is tested with both a custom object-detection deep learning model as well as Google's own Cloud Vision API, in conjunction with human mimicking mouse movements to bypass the challenges. The paper also suggests some potential defense measures to increase overall security and other additional attack directions for reCAPTCHA v2.

Keywords—captcha, reCAPTCHA, deep learning, image processing, object-detection, internet security

I. INTRODUCTION

Completely Automated Public Turing Tests to tell Computers and Humans Apart (CAPTCHA) are a common occurrence when browsing the internet these days, and they are a common security measure employed by websites to defend against bots and other forms of autonomous traffic. They are tests that are designed to identify potential autonomous traffic from real human users. The general goal when designing a captcha is that it should be easy for humans and hard for computers. Originally captchas were text-based and had users identify alphanumeric characters, but these have been surpassed in popularity by image-based ones. With image-based captchas, the user usually must perform some operation on an image or set of images, such as identifying an object. The biggest reason for the switch is the difference in complexity. For human users there is not a significant difference in difficulty between identifying images and identifying text, but for computers there is a significant uptick in complexity. This resulted in many captcha services to start migrating from using text to images, which calls into question the overall security of image-based captchas.

The most popular service in use today is Google's captcha service reCAPTCHA v2 which mainly uses image-based challenges, and the service reportedly serves up millions of captcha challenges every day [1]. reCAPTCHA v2 has two main steps. Firstly, the user must click on a checkbox to initiate the service. At this time, Google performs a risk analysis on the user and returns a score based on several factors. These factors

include behavior of user on current webpage, cookies and browsing history on Google services, etc. A higher score represents the likelihood that the user is a human and if the returned score meets a certain threshold, the user receives no captcha challenge and is automatically passed through, called noCAPTCHA. If the score didn't meet the threshold, the user is presented with an image-based captcha challenge, with an audio-based captcha challenge available as an alternative. Once the user successfully completes the challenge they are verified.

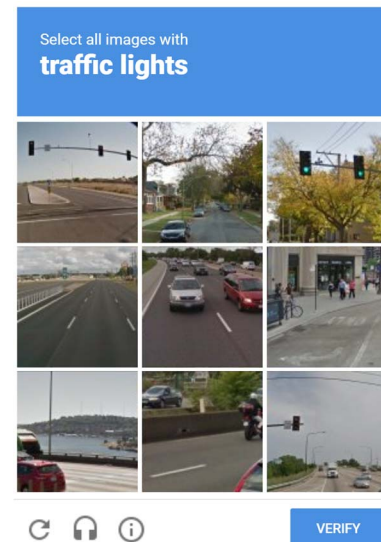


Fig. 1. Example of a captcha image challenge from reCAPTCHA v2

This paper will investigate the security of reCAPTCHA v2, focusing on its image-based challenges. It will identify security measures used on the given images as well as the interactions with the challenge elements. A method is proposed to automatically solve reCAPTCHA v2's image-based challenges using deep learning. Two deep learning methods will be implemented with this solution, a custom trained object-detector and an off-the-shelf object-detection API. This method is then evaluated by testing it on the live reCAPTCHA v2 service. Afterwards, the paper will identify possible additional measures that could be added to increase the security of the image

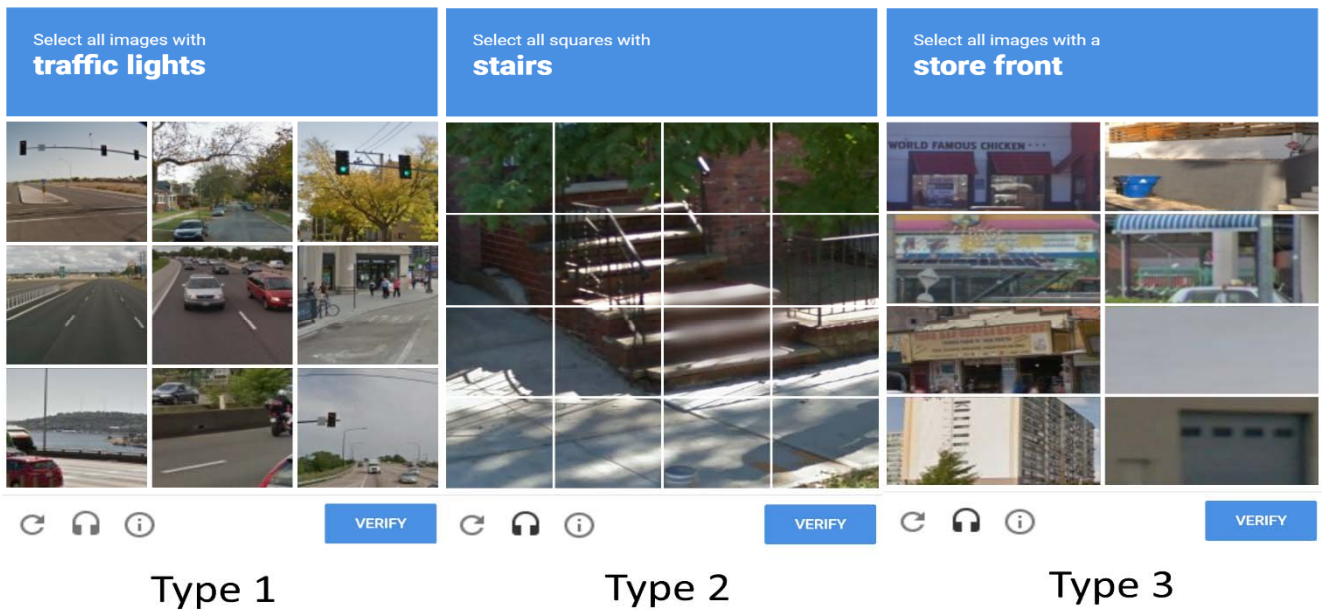


Fig. 2. The 3 different types of reCAPTCHA image challenges

challenges, and it will also suggest additional attack directions for bypassing reCAPTCHA v2.

II. BACKGROUND

For the proposed reCAPTCHA solver, the custom deep learning solution utilizes You-only-look-once version 3 (YOLOv3) object-detection, and the API used is Google's own Cloud Vision API. YOLOv3 is a convolutional neural network (CNN) based object-detection method that maintains accuracy while significantly increasing speed over other CNN based architectures like single-shot multibox detection or region-based CNNs [2]. Google's Cloud Vision offers a high performing pre-trained machine learning model for object-detection through an easy to use REST API [3].

For reCAPTCHA v2's image challenges, the user is given a text and sometimes picture clue and must identify all images containing said clue from a set of given candidates. It is important to note that the user doesn't have to necessarily achieve 100% accuracy, as it is possible to miss some images and still pass the captcha challenge. The image challenges can be further broken down into 3 distinct types which should be taken into consideration.

A. Type 1 – 3 x 3

The first and most common type will be called "3 x 3" or "3 by 3". In this type, the challenge gives the user a 3 by 3 grid of square and distinct images. Additionally, in this subtype after selecting a correct image it may be replaced by a new image that must also be considered. The user will be asked to continue selecting all images with the clue until none remain. This type uses a set of about twelve possible clues.

B. Type 2 – 4 x 4

The second type will be called "4 x 4" or "4 by 4". In this arrangement, the user is given 1 single square image which is divided into a 4 by 4 grid of smaller sections. From testing it seems to share the same clue set as type 1. Unlike type 1

however, type 2 will not have new images be phased in after a successful selection.

C. Type 3 – 4 x 2

The last and least common type is "4 x 2" or "4 by 2". With this version, the user is given a 4 x 2 grid of rectangular and distinct images. This type generally has its own set of clues that isn't used in the other two types. Three potential clues were identified for type 3. Also, like type 1, new images might be phased in to replace correctly clicked images.

Types 1 and 3 can be further broken down into two additional configurations, one configuration being a static set of images, and the other being the aforementioned new images replacing old images.

Once the user has successfully completed a challenge and submitted it, they may be given a totally new additional challenge which they must also complete. This new challenge can be of different type from the previous challenge. Whether or not the user is given additional challenges seems to be based upon reCAPTCHA's risk detection, giving higher risk users more challenges, though most users are verified after only 1 or 2 challenges. This risk analysis score seems to also affect the images themselves, with more noise being added for users who are considered riskier. No other changes to the image besides noise was observed.



Fig. 3. Example of a reCAPTCHA image with and without noise

III. RELATED WORKS

While there are numerous works researching text-based captchas, works focusing image-based captchas are much more limited, and there are even fewer focused on creating an attack to solve real image captcha services. For older revisions of reCAPTCHA, several papers exist such as [4] but the service has changed significantly since then. For the most recent work on reCAPTCHA v2, there is a 2018 paper that created an attack against its image challenges, also utilizing a form of deep learning [5]. Our proposed solution's testing results cannot be directly compared to this other work however, due the differing testing conditions and large updates made the reCAPTCHA v2 since then. This is because reCAPTCHA v2 is a live service and has changed quite a bit in even in the one-year gap, so the performance of [5] will have changed if run today, and its source code isn't available for updated testing. In terms of results, [5] concluded that image-based captchas were quite vulnerable to potential attack and required additional security measures or alternatives.

IV. PROPOSED SOLUTION

The proposed solution is to utilize human mimicking mouse movement in combination with deep learning object-detection to complete reCAPTCHA's image challenges. The mouse movements are used to try and initially receive a better risk analysis score before using the object-detection to detect the given clue in the image challenges.

A. Human Mouse Movements

The humanlike mouse movement is achieved by moving the mouse along a Bezier curve calculated with randomly selected points on the container webpage. reCAPTCHA can see this movement data when doing its risk analysis, and the intent with this stage is to receive a better risk score. While this will not be enough to trigger the automatic noCAPTCHA verification without additional steps, it should be enough to receive a slightly better score than without any initial movement at all. This better score will result in the following image challenge having less noise, which should make it easier for our object-detection method to classify images.

B. Image Challenge

Once the image challenge is received by the client after the initial risk analysis, the first step is to extract the text clue. This is stored in plaintext, so it is trivially extractable. Secondly, the solution extracts the given image set. The images are located and screenshotted by the solution on the webpage using the initial location of the reCAPTCHA widget as a starting point. Google dynamically generates the given set of images as one unique image file for each challenge request, so it needs to be cropped and spliced before being fed into the object-detector.

a) Type 1 and Type 3: For type 1 and type 3 challenges, the 9 and 8 distinct images respectively are cropped and stored individually. These images are then processed by the object-detection method according to the clue, and any images that contain the clue label is selected in reCAPTCHA. After a click, the HTML code is scanned to detect if a new image is going to be phased in to replace the selected image. If the HTML element containing the image doesn't change after a click, then

the set of images is static so the solution simply selects all labeled images and submits. If the HTML element does change, then that means new images will be added in that must also be labeled through the object detector. In this case the solution will continue until no newly added images contain the clue label, then submit afterwards.

b) Type 2: For type 2 challenges, the original image is spliced together from the given generated image, effectively removing the gridlines. This complete image is then sent through the object-detector to receive bounding boxes with labels. The solution then breaks the image back into the 4 by 4 sections, and any sections containing part of a bounding box for the clue label is selected. This type of challenge always features a static set of images, so once the indicated sections are clicked it will be submitted.



Fig. 4. An image pre- and post- YOLO object-detection

After submission of a challenge, the proposed solution will wait and see if reCAPTCHA verifies or instead sends another challenge. If it verifies, then the solution is done, otherwise the solution will repeat the solving process with the new challenge until it does receive verification.

C. Object-Detection

Two different custom YOLOv3 models will be created for testing, created using two different training sets. One data set was obtained from ImageNet [6], supplemented by Google image searches for any clue categories not available on ImageNet. The second training set will be composed of images taken directly from Google reCAPTCHA v2 challenges. For each clue category and image source, 300 images were used for the training set, and an additional 200 for a testing set. Transfer learning was used, with the initial pre-trained weights coming from a YOLOv3 model trained on the Open Images data set by Google.

This solution's YOLO-ImageNet model achieved an average training accuracy of ~96.8% with 3-fold cross validation and took about 12 hours to train on the entire training set. The YOLO-Google model achieved ~96.3% and similarly took around 12 hours to train. The training was done on a machine with an AMD Ryzen 7 1700x CPU, a Nvidia GTX 1080 Ti GPU, and 16 GB of RAM. On their respective testing sets, the YOLO-ImageNet model achieved an accuracy of ~93.8% and the YOLO-Google model achieved ~95.6%. For reference, the Cloud Vision API achieved an accuracy of ~96.7% and ~97.3% on those same respective testing sets.

V. EVALUATION

To evaluate the proposed solution, it will be run on Google's reCAPTCHA v2 demo site 200 times each for both YOLO setups and Cloud Vision. Additionally, it will also be tested with and without the initial human mimicking mouse movement to evaluate its effectiveness. A run will be considered a success if it receives verification from Google, and a failure anytime it does not. A VPN was used between each run to switch IPs in order to avoid Google flagging the test site's actual IP as suspicious, which would result in throttled service requests and increased risk scores.

TABLE I. OVERALL SUCCESS RATE OF PROPOSED SOLUTION

Object-Detection Solution			Success Rate
YOLOv3	ImageNet	w/ movement	61%
		w/o movement	38%
	Google	w/ movement	68%
		w/o movement	41%
Cloud Vision		w/ movement	73%
		w/o movement	42%

As the results show in Table I, adding the initial mouse movements greatly increases the success rate, a 23% and 27% increase for YOLO and 31% increase for Cloud Vision, most likely due to the decreased chance to encounter noise in the given images. Noise is highly detrimental to the performance of object-detection methods, so minimizing any added noise was a key addition to the solution. As for the performance with the mouse movements, both YOLO models and Cloud Vision achieved greater than 60% success rate, with YOLO at 61% and 68%, and Cloud Vision at 73%.

Additionally, the YOLO model trained on actual reCAPTCHA images performed ~7% better than the one trained on the ImageNet set. This is most likely explained by the composition of the images of each set. The ImageNet images are composed of all different sorts of images, including both staged and candid pictures that are not necessarily representative of the types of images the model might encounter. The training set taken from actual reCAPTCHA v2 challenges however is obviously much more indicative of the types of images the model can expect to see in practice.

The results shown in Table II and III reveal some additional aspects about the strengths and weaknesses of the proposed solution and reCAPTCHA v2. First is the addition of the mouse movement, which seemed to greatly decrease the chance of encountering additional challenges after the first, as well as affect the chance of encountering type 1 challenges. There was a more than 20% increase in likelihood of encountering a type 1 challenge when there was no mouse movement. It also seemed to cause great change in performance for type 1 challenges but not types 2 and 3, because noise seems to only be added to type 1 images.

TABLE II. SUCCESS RATE FOR EACH TYPE OF CHALLENGE

Object-Detection Solution			Success Rate		
			Type 1	Type 2	Type 3
YOLOv3	ImageNet	w/ movement	47%	89%	24%
		w/o movement	16%	90%	26%
	Google	w/ movement	61%	90%	33%
		w/o movement	21%	91%	34%
Cloud Vision		w/ movement	67%	91%	31%
		w/o movement	20%	88%	30%

TABLE III. BREAKDOWN OF CHALLENGE TYPES ENCOUNTERED FOR EACH SOLUTION'S 200 RUNS

Object-Detection Solution			Number of Challenges			
			Type 1	Type 2	Type 3	Total
YOLO v3	Image Net	w/ move-ment	152 54%	92 33%	38 13%	282 100%
		w/o move-ment	280 74%	61 16%	40 10%	381 100%
	Google	w/ move-ment	150 52%	102 35%	38 13%	290 100%
		w/o move-ment	284 75%	64 17%	30 8%	378 100%
Cloud Vision		w/ move-ment	141 51%	97 35%	39 14%	277 100%
		w/o move-ment	298 75%	71 18%	28 7%	397 100%

For type 2 challenges, the high success rate could be attributed to the nature of the image compared to the other types. Type 2 is a singular larger image, as opposed to the many smaller images of the other types. Larger images are easier for object-detectors to work on, and less overall images means less complexity. For type 3, the low success rate is most likely attributed to the composition of the image. The images given for type 3 are generally more obscured and more tightly cropped, meaning less contextual information for models to use.

As for the different object-detection methods evaluated, they all performed similarly on type 2 but Cloud Vision has a small advantage on type 1 and YOLO-Google has a small advantage on type 3. The difference in type 1 performance seems to be the biggest factor in the difference in overall performance, as they are the most common type. However, both YOLO models and

Cloud Vision performed significantly worse on type 1 images with the noise involved. When the solutions include mouse movement, all object-detection methods found type 2 the easiest and type 3 the hardest to correctly identify.

VI. SUGGESTIONS FOR IMPROVING reCAPTCHA v2

From the results of the experiments, as well as personal use, a couple of possible improvements to reCAPTCHA's overall performance as a security tool have presented themselves.

A. Store Clue Differently

Currently, reCAPTCHA v2 stores its clues as plaintext, which is perhaps the least secure way for it to be stored. We recommend storing and presenting the clue in another fashion, such as an image of the text, or perhaps by displaying the clue itself using only a representative image. This change would not greatly impact human usability, but it would add an extra layer of complexity that automated systems would have to handle.

B. Remove Type 2 Challenges

It seems that the difficulty of type 2 challenges is lacking in comparison to type 1 and 3, with the proposed solution achieving about ~90% success rates. It might be better to remove them outright and focus on type 1 and type 3 challenges, as the solution struggled much more with them.

C. Image Transformations

As shown by the decrease in performance when noise was involved, such transformations upon an image causes issues for deep learning object-detectors. Perhaps additional image transformations could also be implemented to help throw off these machine learners, such as slight rotations or color shifts. Again, these types of modifications increase the difficulty for machines considerably more than it does for humans.

D. Context-Sensitive Challenges

Instead of having the user just perform identification on an image, maybe have the user do some task that also involves deciphering the context of the image. For example, instead of asking the user to determine whether or not a picture contains a fruit, instead ask the user to count the number of fruits in the picture. Another example is instead of asking the user to identify whether a photo has a person in it, ask the user to determine whether the person in the image is smiling or frowning. Context-sensitive challenges like these would require redesigning most current automated systems to adapt for each new variation, as they would need to do more than just image classification/detection. However, for humans these come naturally and should not be a large step up in complexity.

VII. SUGGESTIONS FOR BYPASSING reCAPTCHA v2

In addition to suggestions for improvement, some suggestions for alternative methods for bypassing and attacking reCAPTCHA have presented themselves as well.

A. Focus on Type 1 and Type 3 Performance

Unlike type 2 challenges, types 1 and especially 3 are still relatively difficult for object-detection models to work out, so

focusing on improving performance in these areas will increase overall performance greatly.

B. Focus on Google Cloud Vision

As the experiments indicated, Cloud Vision performed very well in the solution. One possible way forward is to focus on using the API instead of creating a custom object-detection method. Cloud Vision is trained and maintained by Google's vast resources, much more than any one person could reasonably obtain, and could possibly increasingly outperform most custom solutions in the future. Additionally, it is possible that Google uses reCAPTCHA v2 in part to help train Cloud Vision, so it should adjust to changes more easily.

C. Attack the Audio Alternative

In its mission to stay accessible to those with disabilities, reCAPTCHA as well as many other similar services offer an audio alternative to their main visual captcha challenges. So, despite the advancements in computer vision, it might be easier to attempt to bypass reCAPTCHA through their audio challenge as opposed to their image challenge. Part of this is due to the fact that the overall complexity of their audio challenges is less than that of their image ones.

D. Initial Risk Analysis

Being able to solve image challenges would be unnecessary if a reliable way to pass the initial risk analysis is found. Finding a way to consistently improve your risk analysis score beyond mimicking human mouse movement to receive noCAPTCHA could make the ability to solve these challenges meaningless.

VIII. CONCLUSION

To conclude, Google reCAPTCHA v2's image challenges can be reliably solved by a computer, with the proposed solution achieving a high success rate in bypassing them. The solution can solve them consistently with not only a custom trained object-detection deep learning model achieving ~68% bypass rate, but also with Google's own Cloud Vision API achieving a ~73% bypass rate. Additionally, several suggestions were discussed to potentially improve reCAPTCHA's security, as well as several suggestions for future work in bypassing methods.

REFERENCES

- [1] reCAPTCHA, <https://developers.google.com/recaptcha/>
- [2] Redmon, J., & Farhadi, A. (2018). "YOLOv3: An Incremental Improvement", <https://pjreddie.com/media/files/papers/YOLOv3.pdf>
- [3] Google Cloud Vision, <https://cloud.google.com/vision/>
- [4] S. Sivakorn, J. Polakis, and A. D. Keromytis, "I'm not a human: Breaking the google recaptcha," Black Hat, (i), pp. 1--12, 2016.
- [5] BinbinZhao,HaiqinWeng,ShoulingJi,JianhaiChen,TingWang, Qinning He, Raheem Beyah. 2018. Towards Evaluating the Security of Real-World Deployed Image CAPTCHAs. In AISec '18: 11th ACM Workshop on Artificial Intelligence and Security Oct. 19, 2018, Toronto, ON, Canada. ACM, NY, NY, USA, 12 pages. <https://doi.org/10.1145/3270101.3270104>
- [6] ImageNet, <http://www.image-net.org/>