**The History and Future of CAPTCHA**

Adam M. Klions

Programming Languages SCS 250 01

Professor Weakland

December 1, 2021

As companies become more reliant on technology, they face increasing challenges to secure their proprietary information and customer data. Studies have shown that 94.2% of websites have been attacked by an automated system in attempts to steal personal information, passwords, or to just send spam and advertisements. (Zeifman 2017) These automated systems, known as bots, make up about 25% of all internet traffic in 2020 and target websites of every size. (Imperva 2021) CAPTCHAs have been used for the last twenty years as a way to prevent these bots from accessing websites with requests. The name CAPTCHA stands for ***Completely Automated Public Turing test to tell Computers and Humans Apart.*** CAPTCHA and its successors are used by over six million websites today. Over the years CAPTCHAs have evolved to both allow for more ease of use for legitimate users while also better preventing computers from circumventing them.

The first system like CAPTCHA was first created in 1997 by AltaVista in an effort to stop bots from skewing their search results through bulk URL submissions. Three years later Yahoo asked researchers at Carnegie Mellon for a solution to a similar problem, bots creating accounts to spam advertisements to actual users. Those researchers developed CAPTCHA. At the time, CAPTCHA was just a randomly selected dictionary word or just a string of random characters with distortions and high noise backgrounds that humans could easily identify but computer text recognition could not, like the one shown in figure 1(Hasan). As CAPTCHA's use grew, its creator, Luis von Ahn, (2011)  believed that CAPTCHA was wasting people's time and energy and

that CAPTCHA had the potential to both protect against spam and also utilize that spent effort. In 2007 his team began to work on reCAPTCHA. This project would be acquired by Google two years later. reCAPTCHA version 1, shown in figure 2, (Google) worked with the same basic design as its predecessor, with one major change, instead of a randomly selected string it would use two words, one that the computer already knew and one that the computer took from a scanned text that it did not recognize. This allowed reCAPTCHA to identify unknown words in texts through mass collaboration. (Ahn 2011)

reCAPTCHA is still in use today but looks very different than it did before. Instead of using words from scanned texts like in reCAPTCHA v1, reCAPTCHA v2, figure 3, (Surana 2021) takes images from Google Street View and has the user identify objects within these images. The data gathered from this is, according to Google, used to help improve machine learning and image recognition. It was later modified to only require the user to click on a checkbox, shown in figure 4. (Google) It uses data such as mouse movements, cookies, browsing history, and other metrics that Google hasn't disclosed for security reasons to determine a score ranging between 0 and 1. A score of 0 means the user is most likely a bot and a score of 1 means the user is most likely a human. Based on the score there are three actions that can be taken: 1) If the score is too low the website can deny access to the user. 2) If the score is high enough, it can allow access to the user. 3) If the score is between the two thresholds, it can have the user complete an image recognition test to further verify whether or not they are a human. (Datadome 2021) A third version was created in 2018 that runs entirely in the background using similar metrics as v2 and only requires user input if the user is

Figure 1: Early Text-Based CAPTCHA
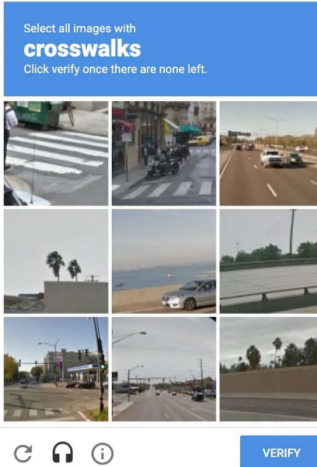

Figure 2: reCAPTCHA v1
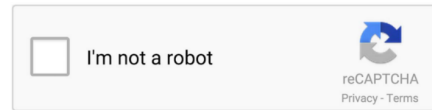

Figure 3: reCAPTCHA v2


Figure 4: reCAPTCHA v2 Checkbox

deemed to have a high probability of being a computer. This design alleviates one of the biggest complaints about reCAPTCHA in that it wastes legitimate user's time.

Companies are in a constant arms race against the evolution of the computer's abilities to solve CAPTCHAs. There are two means which people use today to circumvent CAPTCHAs: CAPTCHA farms, and sufficiently advanced AI designed to defeat CAPTCHAs. CAPTCHA farms are companies that pay thousands of workers to manually solve CAPTCHAs for the bots at pennies per solve. (Netacea 2021) There are some solutions to this issue, such as validating that the client's ip is the same throughout the process, but the implementation for these solutions aren't used everywhere. The other problem CAPTCHA designers face is the fact that machines are getting more effective at solving CAPTCHAs on their own. Researchers at Columbia University have already created a program that can solve reCAPTCHA v2 with around a 70.78% success rate. (Sivakorn et al. 2016) Technology is quickly catching up to the capabilities of humans. According to Si Hong Loo (2012), "A good CAPTCHA has to satisfy three basic properties. The tests must be: Easy for humans to pass, easy for a

server to generate and determine, and difficult for a machine to pass." The difference between the first and third property is growing narrower which makes designing new tests difficult.

3D CAPTCHAs are an idea of a relatively new test that would be very difficult for computers to solve but at the cost of being more difficult to generate and taking longer for humans to solve. These CAPTCHAs involve taking a 3D model of some object or text and having the user rotate it in 3D space in order to identify it. Most companies seem to be going in the opposite direction though, favoring ease of access over catching every bot. (Si Hong Loo 2012)

SenCAPTCHA is a novel idea utilizing the sensors in mobile devices in order to create tests that can be used on devices where other kinds of tests are too difficult to implement. SenCAPTCHA uses the orientation sensors in smartphones and smart watches to have the user roll a ball onto the eyes of a randomly selected animal. Then it will check the path that the ball took and determine whether a human or machine completed the task. It also features a series of image mutations as a way to create more tests from a small selection of pictures. The use of animals will make facial recognition and eye detection difficult for machines due to the high variance of facial structures



**Figure 5: SenCAPTCHA**

between different species. There is also an option to complete the test with just a mouse to move the ball which allows for any device to be able to be used. (Feng et al. 2020)

The future of spam prevention is uncertain. As technology improves it will become harder and harder to differentiate between humans and machines. More advanced tests and techniques will be required in order to stop spam, but this comes at the risk of creating too many false positives where actual humans may be prevented from accessing a website due to them being unable or unwilling to complete the CAPTCHA. SenCAPTCHA shows that we still have ways to create tests that machines will struggle with, but in the future more creative solutions will be required. Most developers seem to be moving away from the obstructive tests and instead favoring background monitoring and only interrupting workflow if the website believes the user may be a bot. Overall this is a good direction to move in, due to a fair number of users disliking and misunderstanding CAPTCHAs.

References

Ahn, L. b., (2011, April). *Massive-scale Online Collaboration* [Video]. TED Conferences.

    https://www.ted.com/talks/luis_von_ahn_massive_scale_online_collaboration

Chew, M., & Tygar, J. D. (2004, September). Image recognition captchas. *International*

    *Conference on Information Security* 268-279. Springer, Berlin, Heidelberg.

Feng, Y., Cao, Q., Qi, H., & Ruoti, S. (2020). SenCAPTCHA: A Mobile-First CAPTCHA

    Using Orientation Sensors. *Proceedings of the ACM on Interactive, Mobile,*

    *Wearable and Ubiquitous Technologies*, *4*(2), 1–26.

    https://doi.org/10.1145/3397312

Hasan, W., (2016). A Survey of Current Research on CAPTCHA. International Journal of

    Computer Science & Engineering Survey. 7. 1-21. 10.5121/ijcses.2016.7301.

Imperva. (2021, April 13). *Bad Bot Report 2021 The Pandemic of the Internet*. Imperva.

    Retrieved from

    https://www.imperva.com/resources/resource-library/reports/bad-bot-report/.

Loo, S. H. (2012). An Insight into CAPTCHA. *Introduction to Information and System*

    *Security*, 31–36.

Netacea. (2021, October 15). *Part Two: What are Captcha Farms?* Netacea. Retrieved

    November 30, 2021, from https://www.netacea.com/blog/what-are-captcha-farms/.

*ReCAPTCHA v2 vs V3: Efficient bot protection? [2021 update]*. Datadome. (2021).

    Retrieved November 30, 2021, from

    https://datadome.co/bot-detection/recaptchav2-recaptchav3-efficient-bot-protection

    /.

Sivakorn, S., Polakis, J., & Keromytis, A. D. (2016). *I'm not a human: Breaking the Google reCAPTCHA*. Blackhat. Retrieved November 30, 2021, from https://www.blackhat.com/docs/asia-16/materials/asia-16-Sivakorn-Im-Not-a-Human-Breaking-the-Google-reCAPTCHA-wp.pdf.

Surana, P. (2021, January 17). Integrating reCAPTCHA with Next.js. Retrieved November 30, 2021, from https://prateeksurana.me/blog/integrating-recaptcha-with-next/.

Zeifman, I. (2017, January 27). *Bot traffic report 2016*. Imperva. Retrieved November 30, 2021, from https://www.imperva.com/blog/bot-traffic-report-2016/?redirect=Incapsula.