

Dazed & Confused: A Large-Scale Real-World User Study of reCAPTCHA_{v2}

Andrew Searles
searlesa@uci.edu
UC Irvine

Renascence Tarafder Prapty
rprapty@uci.edu
UC Irvine

Gene Tsudik
gene.tsudik@uci.edu
UC Irvine

ABSTRACT

Since about 2003, CAPTCHAs have been widely used as a barrier against bots, while simultaneously annoying great multitudes of users worldwide. As their use grew, techniques to defeat or bypass CAPTCHAs kept improving, while CAPTCHAs themselves evolved in terms of sophistication and diversity, becoming increasingly difficult to solve for both bots and humans. Given this long-standing and still-ongoing arms race, it is important to investigate usability, solving performance, and user perceptions of modern CAPTCHAs. In this work, we do so via a large-scale (over 3,600 distinct users) 13-month real-world user study and post-study survey. The study, conducted at a large public university, was based on a live account creation and password recovery service with currently prevalent CAPTCHA type: reCAPTCHA_{v2}.

Results show that, with more attempts, users improve in solving checkbox challenges. For website developers and user study designers, results indicate that the website context directly influences (with statistically significant differences) solving time between password recovery and account creation. We consider the impact of participants' major and education level, showing that certain majors exhibit better performance, while, in general, education level has a direct impact on solving time. Unsurprisingly, we discover that participants find image challenges to be annoying, while checkbox challenges are perceived as easy. We also show that, rated via System Usability Scale (SUS), image tasks are viewed as "OK", while checkbox tasks are viewed as "good".

We explore the cost and security of reCAPTCHA_{v2} and conclude that it has an immense cost and no security. Overall, we believe that this study's results prompt a natural conclusion: *reCAPTCHA_{v2} and similar reCAPTCHA technology should be deprecated.*

1 INTRODUCTION

Many types of Internet-based activities and services require verification of human presence, e.g., ticket sales, reservations, and account creation. Left unchecked, bots will gobble up most resources available through such activities: they are much faster and way more agile than any human or a group thereof. This problem is not new: the first seminal step to combat it took place in 2003 when von Ahn et al. [71] proposed CAPTCHA as an automated test that is supposed to be easy for humans to pass, yet difficult or impossible for computer programs (aka bots) at the time. The key conjecture underlying the CAPTCHA concept is that, if a computer program successfully solved CAPTCHAs, then the same program could be repurposed to solve some computationally hard AI problem.

This seemed to be a win-win situation: either CAPTCHAs attest to genuine human presence or they spur a significant advance in AI technology. Furthermore, CAPTCHAs were touted as a tool for the common good, since human-based solutions helped with difficult

(for computers) and useful tasks, such as recognizing blurred text that confounded OCR algorithms, or labeling photos with names of objects appearing in them in order to aid image classification.

Another major advance occurred in 2007 when von Ahn et al. introduced reCAPTCHA [72]. reCAPTCHA was designed to reuse challenge results as a form of human-based data labeling for advancing machine learning. Google acquired reCAPTCHA in late 2009 [4] and, by June 2010, it was reported that reCAPTCHA had over 100 million distinct daily users [5]. Assuming that this number stayed constant since 2010 (though it most likely grew significantly), over half a trillion reCAPTCHAs have been solved in the meantime. This collectively amounts to an immense human cost.

However, almost from the start, an "arms race" began between bot and CAPTCHA developers. Most early CAPTCHA types were based on recognition of distorted text. Unfortunately, as a consequence of rapid advances in machine learning and computer vision, bots evolved to quickly and accurately recognize and classify distorted text [39, 45, 75], reaching over 99% accuracy by 2014 [43, 64]. To this end, in 2012 Google switched from distorted text to image classification, using images from the Google Street View project [61]. This transition ended in 2014 with the introduction of reCAPTCHA_{v2} [6], which uses a two-step process: (1) a combination of behavioral analysis and a simple checkbox, and (2) image classification tasks as a fallback for users who fail the checkbox challenge [22]. By 2016, both (1) and (2) were defeated with a high degree of accuracy by bots [65].

Regardless of its diminished efficacy, reCAPTCHA remains to be the prevalent CAPTCHA type on the Internet [16], deployed on over 13 million websites in 2023. It is therefore important to periodically evaluate and quantify its impact in terms of usability, solving performance, and user perceptions.

Several prior CAPTCHA user studies explored solving performance, e.g., [29, 32, 35, 37, 41, 42, 46, 50, 54, 59, 62, 63, 68, 69]. Also, [35, 42, 59] looked into usability of CAPTCHAs via the well-known SUS scale. [29, 37, 41, 50, 54, 63, 68] studied user preferences related to CAPTCHA types. However, only two recent (2019/2023) user studies [63, 68] involved reCAPTCHA_{v2}. However, [68] had relatively few participants (40), used unclear methodology, and did not consider usability. [63] presents interesting comparison points discussed in Section 5. Most other user studies [35, 37, 41, 42, 50, 59] were conducted on newly proposed (and therefore, mocked-up) CAPTCHA types.

Furthermore, many previous CAPTCHA studies [32, 35, 42, 46, 50, 59, 63] were conducted on Amazon Mechanical Turk (MTurk) [15], which exhibits data quality issues [73]. Also, all these studies involved some bias, since participants were informed about study goals, i.e., they were selected based on their willingness to solve CAPTCHAs, for a certain monetary reward.

The above discussion motivates the work presented in this paper, the centerpiece of which is a large-scale (> 3, 600 participants) 13-month IRB-approved user study of reCAPTCHA v2. The study was conducted using a live account creation (and password recovery) service with unaware participants who, for the most part, have never before used this service. Results of the study yield some interesting observations that might be of interests to CAPTCHA designers as well as websites using (or considering the use of) CAPTCHAS.

Main contributions of this work are:

- A comprehensive quantitative analysis of solving time and how it relates to certain dimensions. In particular, this is the first study to obtain multiple solving attempts per person. It shows that form-specific checkbox solving time improves with more attempts, with the first attempt being 35% slower than the 10th, shown in Tables 7 and 8. We also show statistically significant changes in checkbox solving time based on the type of service, with password recovery being faster, as shown in Tables 4, 5 and 6. With respect to educational level¹, there is a direct trend from freshmen (slowest) to seniors (fastest) at solving reCAPTCHA v2 as shown in Tables 9, 10 and 11. In terms of participants' major (field of study), there were minor trends with statistical significance of technical (aka STEM) majors solving time being faster than that of non-technical majors, as shown in Tables 12, 13 and 14.
- An in-depth qualitative analysis of reCAPTCHA v2 usability for both checkbox-only and checkbox-and-image combination. Results demonstrate that 40% of participants found the image version to be annoying (or very annoying), while <10% found the checkbox version annoying. SUS data shows that image results have a mid-score of 58, while checkbox has a score of 78, with 90 being the highest score observed. Based on the open-ended feedback represented in a *word cloud*, participants' most frequent term for the checkbox version was “easy” and, for the image version – “annoying”.
- A detailed discussion of the cost and security of reCAPTCHA v2 (Section 6). Our security analysis shows a blatant vulnerability [47], the ease of implementing large-scale automation [66], usage of privacy invasive tracking cookies [66], and weakness of security premise of fallback (image challenge) [38]. Our cost analysis investigates total human time spent solving reCAPTCHA v2, human labor, network traffic, electricity usage, potential profits and the corresponding environmental impact. There have been at least 512 billion reCAPTCHA v2 sessions, taking 819 million hours, which translates into at least \$6.1 billion USD in free wages. Traffic resulting from reCAPTCHA v2 consumed 134 Petabytes of bandwidth, which translates into about 7.5 million kWhs of energy, corresponding to 7.5 million pounds of CO₂ pollution.

Organization: Section 2 provides some background on current CAPTCHA types and System Usability Scale (SUS). Then, Section 3 describes the methodology, design, ethics, and implementation of the user study. Next, Section 4 presents the results and their analysis.

¹In the American undergraduate system, “freshmen” are 1st-year students, “sophomore” – 2nd, “junior” – 3rd, and “senior” – 4th.

Then, Section 5 contextualizes our results against previous user studies. Next, Section 6 presents the cost and security analysis. Section 7 concludes the paper.

2 BACKGROUND

This section overviews the CAPTCHA landscape and System Usability Scale (SUS). Given familiarity with these topics, it can be skipped with no loss of continuity.

2.1 CAPTCHAS

A recent survey by Guerar et al. [44] is a comprehensive overview of the current CAPTCHA landscape. It proposes a ten-group classification to encompass all current and emerging schemes: Text-based, Image-based, Audio-based, Video-based, Game-based, Slider-based, Math-based, Behavior-based, Sensor-based, and Liveness-detection. It also discusses usability, attack resilience, privacy, and open challenges for each class. Since this paper focuses on behavior and image-based CAPTCHAS (Which are used in reCAPTCHA v2), we summarize them below. For the rest, we refer to [44].

Text-based CAPTCHAS are the earliest type, originally proposed by Von Ahn et al. [71]. They present the user with an image containing a random sequence of visually-distorted alphanumeric characters, possibly combined with other visual elements, e.g., lines and dots. The user is required to correctly identify the characters and type them into a text field. An example of text-based reCAPTCHA is shown in Figure 1. The idea is that humans should be significantly better than bots at recognizing distorted characters. However attacks on text-based CAPTCHAS have been quite successful and widely studied [33, 40, 43, 55, 67, 76]. For example, [76] and [55] achieve over 97% accuracy on certain text-based CAPTCHAS within fractions of a second, using machine learning. Furthermore, [43] achieves 99.8% accuracy on reCAPTCHA schemes of the time. Although research has shown that basic text-based CAPTCHAS are no longer effective in distinguishing humans from machines, they are still widely used.

Image-based CAPTCHAS typically require users to perform an image classification task, such as selecting images that match the accompanying written description. Most popular instances are hCAPTCHA [19] and reCAPTCHA [21] version 2 onward. Examples are shown in Figure 2, Figure 3 and Figure 5. The difficulty of these CAPTCHAS is associated with that of computer vision-based image classification. At the time of the introduction of these CAPTCHAS types, corresponding problems were not easily solvable by machines. However as computer vision research advanced, attacks on image-based CAPTCHAS became more successful. Concrete attacks include [26, 48, 49, 57, 66, 74], some of which report success rates of 85% for reCAPTCHA and 96% for hCAPTCHA.

Behavior-based (or invisible) CAPTCHAS are newer: they either require users to click a box (e.g., “I am not a robot”), or are completely invisible/transparent to the user. Instead of a visual challenge, they rely on client-side scripts and other opaque techniques to collect, in the background, historical behavioral information about the user. This information is sent to the CAPTCHA provider, which uses various heuristic-based techniques to identify bot-like behavior. For instance, Google’s popular No-CAPTCHA reCAPTCHA: “actively considers a user’s entire engagement with



Figure 1: reCAPTCHA v1 distorted text CAPTCHA [21]

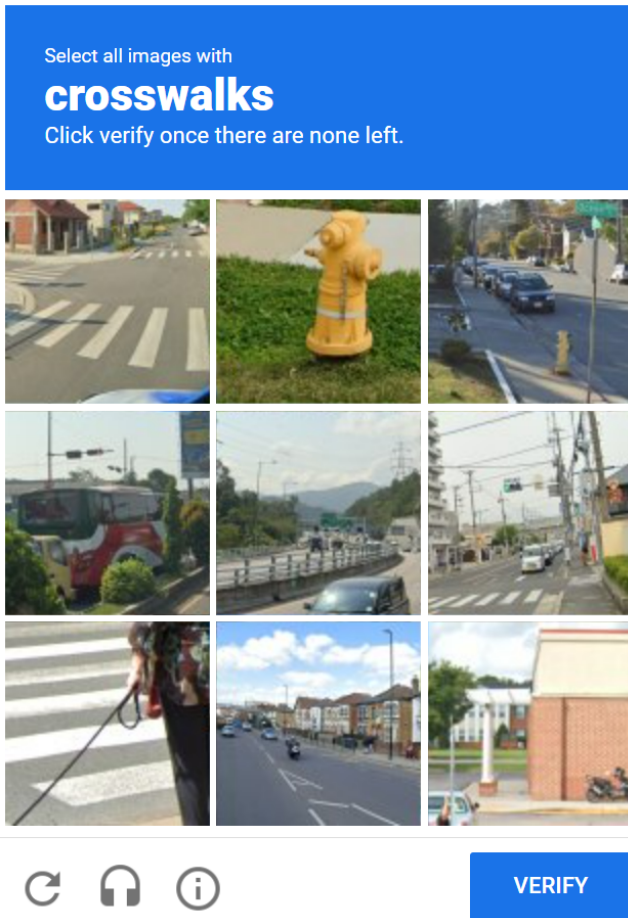


Figure 2: Image Labeling Task CAPTCHA [21]

the CAPTCHA – before, during, and after – to determine whether that user is a human” [22]. Sivakorn, et al. [65], evaluate reCAPTCHA risk analysis system and determine that Google tracks cookies, browsing history, and browser environment, e.g., canvas rendering, user-agent, screen resolution and mouse. [65] also showed that legitimate cookies can be automatically farmed to attack reCAPTCHA v2 with 100% success on a large scale.

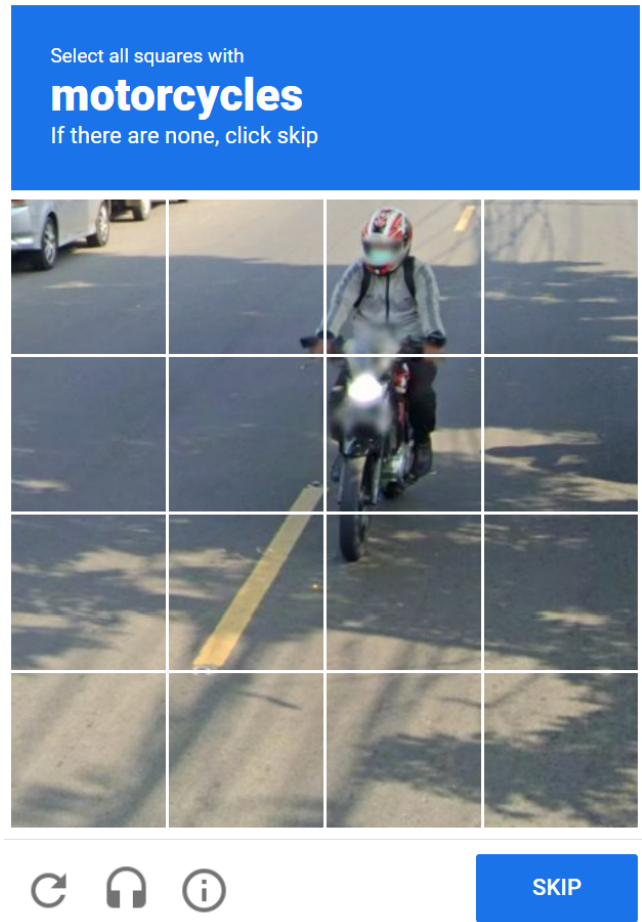


Figure 3: Image Bounding Box Task CAPTCHA [21]

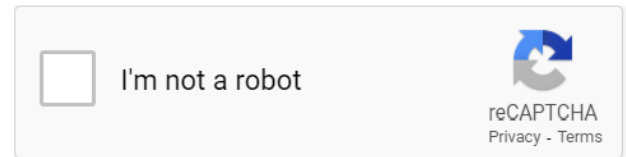


Figure 4: v2 checkbox CAPTCHA [22]

2.2 System Usability Scale (SUS)

System Usability Scale (SUS), shown in Figure 6, is a classical and popular survey method designed to assess usability of various systems or products. Proposed by Brooke, et al. [31] in 1996, it consists of ten statements: five positive and five negative. Each statement is on a 5-point Likert scale ranging from *Strongly Disagree* (1) to *Strongly Agree* (5).

SUS is widely used to measure usability of a wide range of products and systems, from everyday products (such as phones, fitness bands, and appliances [53, 56]) to websites, software, mobile apps and even CAPTCHAs [36, 51–53, 60, 70]. SUS is very popular because

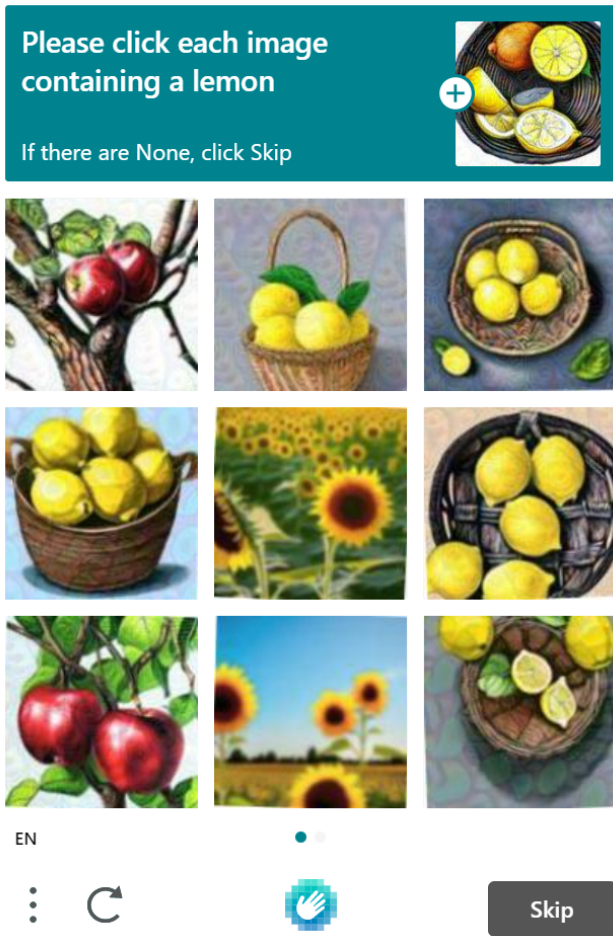


Figure 5: hCAPTCHA [19]

- (1) I think that I would like to use this system frequently.
- (2) I found the system unnecessarily complex.
- (3) I thought the system was easy to use.
- (4) I think that I would need the support of a technical person to be able to use this system.
- (5) I found the various functions in this system were well integrated.
- (6) I thought there was too much inconsistency in this system.
- (7) I would imagine that most people would learn to use this system very quickly.
- (8) I found the system very cumbersome to use.
- (9) I felt very confident using the system.
- (10) I needed to learn a lot of things before I could get going with this system.

Figure 6: System Usability Scale (SUS) [19]

of its simplicity and conciseness. Participants tend to easily understand and quickly complete the SUS questionnaire. The process of calculating scores is also very straightforward:

- For odd-numbered statements, subtract 1 from each response value
- For even-numbered statements, subtract each response value from 5

- Sum up all response values and multiply the result by 2.5 This yields a SUS score between 0 and 100 for each participant.

To associate a given usability level with individual scores, [28] provides adjective scaling, shown in Table 1. This scale consists of seven usability levels starting from the worst imaginable usability and going up to the best imaginable usability.

Table 1: Adjective Ratings of SUS Scores

Adjective	Mean SUS Score
Worst Imaginable	12.5
Awful	20.3
Poor	35.7
OK	50.9
Good	71.4
Excellent	85.5
Best Imaginable	90.9

3 THE USER STUDY

Recall that the goals of the user study are to measure solving times, error rates, and user perceptions of reCAPTCHA v2, the currently prevalent CAPTCHA type.

3.1 The Setting

This study was conducted continuously over the period of roughly 13 months in the 2022-2023 time-frame. It took place on a campus of the University of California Irvine, though the scope was limited to one specific school. The term *school* denotes an organizational entity that includes two or more academic departments. The university contains a number of such schools, e.g., School of Engineering, School of Law, and School of Humanities.

The specific school hosting our study is called SICS: *School of Information & Computer Sciences*. SICS includes several departments, all somehow related to Computer Science. SICS offers a number of fairly typical undergraduate (BS) and graduate (MS and PhD) programs.

For many years, SICS requires for every person, who for the first time, enrolls in any SICS course, to create a SICS-specific user account via the school's web interface. A typical scenario is that a student who enrolls in at least one SICS course in their entire university career, would create a SICS account **only once**. Consequently, a student who wants to create a SICS account has not previously engaged in SICS account creation, meaning that they have no knowledge of the workflow involved, and no expectations of either seeing or not seeing CAPTCHAs as part of the process.

This motivates the key feature of our user study: introduction (insertion) of reCAPTCHA v2 into the SICS account management workflow. This actually involves two separate services: (1) account creation for new users, and (2) password recovery for users with existing accounts. This was accomplished with the much-appreciated help and cooperation of the SICS IT Department.

As mentioned earlier, the study ran for about 13 months. This is because we wanted to include as many distinct users as possible. Since the yearly academic calendar has multiple terms, we aimed

to catch the beginning of each term (and a week or so prior to it), since this is the time when the bulk of new account creation and password recovery activity typically takes place.

3.2 Justification

We now discuss the rationale for the user study setting. Clearly, an ideal and comprehensive CAPTCHA user study would be as inclusive as possible, comprising a true cross-section of the world population. Whereas, our study targeted participants are (mostly) university students, including undergraduates who range from incoming (freshmen) to graduating (seniors), as well as graduate students enrolled in a variety of programs (MS, MA, MBA, MFA, JD, MD, PhD). The latter are split among so-called *professional* degree programs, e.g., MBA, JD, MD, and some MS/MA, while others are in regular degree programs, e.g., PhD, MFA, and some MS/MA. Such participants are surely not representative of the world, or even national, user population. Nonetheless, we conjecture that data stemming from this admittedly narrow population segment is useful, since it reflects an “optimistic” perception of CAPTCHAs. This is because young and tech-savvy users represent the most agile populations segment and the one most accustomed to dealing with CAPTCHAs, due to their heavy Internet use. Thus, by studying various (not generally positive) impact factors of CAPTCHAs, we prefer to err on the side of the population that is intuitively the least allergic to CAPTCHA use.

Some reasons for our study setting are fairly obvious. In particular, it would have been very challenging, if not impossible, to convince any other organization to introduce CAPTCHAs into its service workflow, or to allow us to collect data about their current CAPTCHA use. Alternatively, one could imagine approaching Google and requesting access to the centralized reCAPTCHA v2 service. This would have been ideal since it would give us access to a huge number of diverse reCAPTCHA v2 users worldwide. Indeed, we attempted to do this. However Google’s legal team denied our request to gain access to large-scale data from reCAPTCHA v2. There is very likely a natural counter-incentive for Google (or any other CAPTCHA provider) to cooperate with outside researchers in a user study, since doing so might reveal certain negative aspects of the service. Another possibility would have been to create our own brand new service and use CAPTCHA to guard access to it, thus hoping to attract prospective users of broad demographics. While theoretically plausible, doing so would be prohibitively time and effort-consuming for academic researchers.

Finally, even with our somewhat narrow target demographic of university students, the user study could have been more latitudinal, i.e., it could span multiple universities in various parts of the world. This would have yielded more valuable results across political, cultural and linguistic boundaries. However, this would have been a massive effort requiring careful coordination with, and participation of, both researchers and IT departments in each university.

3.3 The Website

The SICS website used in the study is hosted within the university network. In order to create a SICS account, a user must first login to the campus VPN with their university account. This allows us to claim, with high confidence, that all collected data stemmed from

real human users, who are, for the most part, students (see Section 3.5 below).

The back-end is a basic PHP server that serves HTML and JavaScript. It is maintained by SICS IT department. The account creation service includes a form requesting basic student information, e.g., name and student ID. The password recovery service includes a form requesting existing account information. In both cases, reCAPTCHA v2 was initially hidden and rendered after clicking the submit button. Basic website workflows for account creation and password recovery are described in Appendix A.

All timing events were measured using JavaScript native Date library, which has millisecond precision. JavaScript was used to block form submission, such that an initial timing event is recorded and a reCAPTCHA v2 is rendered simultaneously. Initially, a behavior-based click box CAPTCHA (Figure 4) is presented. In order to solve it, a user clicks the checkbox sending data to Google reCAPTCHA v2 site. It either approves the request or presents an image-based challenge. Upon reCAPTCHA v2 validation, a second timing event is captured and the form is submitted.

Solving time is thus comprised of the time interval starting from CAPTCHA rendering until the client browser receives a successful validation response from Google reCAPTCHA v2 service. (This includes image challenges and failed solution attempts.) Upon successful form submission, the IT database stores these two timestamps along with the form information.

3.4 Directory Crawler

Recall that the study involved unwitting participants, i.e., unaware of both existence and purpose of the study. In order to subsequently obtain demographic information about each participant, we created a JavaScript crawler that automatically searches the university directory using email addresses. This directory is publicly available from both inside and outside the university network. Information gathered by the crawler includes major and college education level (freshman, sophomore, junior, senior, or graduate) of each participant.

3.5 Logistics & Data Cleaning

In total, the SICS IT department supplied 9,169 instances of account creation and password recovery with reCAPTCHA v2 solving time data. The original form data was larger, since it included errors, such as incomplete forms and incorrect values. Each record (form) has the following fields: database ID, date and time, student ID, email address, service, and timing. Starting with 9,169 instances, we filtered results using the directory crawler, labeling entries with student IDs that were not found and correcting student IDs with minor typos. A total of 229 entries were labeled as none for student ID and 295 student ID typos were corrected.

Successful form submissions have certain constraints, e.g., field formatting. If a person enters erroneous data that does not fit the constraints, they still have to solve a reCAPTCHA v2 before the form is submitted. Cases of multiple submissions occurred because of unsuccessful attempts to enter form data. For some entries, there were small typos, though mixed with temporal evidence they were correctable.

28 records were removed, since each had solving time of > 60 seconds which adds a high degree of variance. We ended up with 9,141 valid records of which 8,915 correspond to 3,625 unique participants. 226 entries, labeled as none for student ID, are not included among the unique participants, attempts, educational level, and major analysis. Of the 8,915, 231 form submissions correspond to 52 unique non-students (i.e., faculty or staff) and are not included in the educational level and major analysis. For the purposes of the educational level and major analysis, 3,573 unique students completed 8,631 reCAPTCHA2 challenges.

3.6 Post-Study Survey

After the completion of the study, we randomly selected and contacted, by email, 800 participants in order to solicit feedback on their reCAPTCHA2 experience via a survey (a Google form). In the end, a total of 108 completed the survey. The incentive was an \$5 Amazon gift card. The survey collected answers to SUS questions regarding both checkbox and image CAPTCHAS. It also collected information about (more detailed) demographics, frequency and nature of internet usage, as well as preferences and opinions about checkbox and image CAPTCHAS.

3.7 Ethical Considerations

The user study was duly approved by the university’s Institutional Review Board (IRB). Collection of student email addresses for recruitment and demographic analysis purposes was also explicitly approved. Since prospective participants were not pre-informed of their participation in the study, two additional documents were filed and approved by the IRB: (1) “Use of deception/incomplete disclosure” and (2) “Waiver or Alteration of the Consent”. Study participants who completed the post-study survey were compensated US\$5 for about 5 minutes of their time. This was also IRB-approved.

No personally identifiable information (PII) was used in the demographics analysis.

After the completion of the study, all participants were informed, by email, of their participation and the purpose of the study. They were also informed that some basic demographic information about them that was collected via campus directory lookup.

4 RESULTS & ANALYSIS

This section presents the results of the user study based on the live service experiment. We consider both quantitative (solving time) and qualitative (SUS, rating, feedback) data to provide a comprehensive analysis of reCAPTCHA2 usability.

4.1 University Demographics

Student population of the university is large and diverse. Figure 7 summarizes its demographics. We use university demographics, because students from multiple departments who take any SICS course create accounts. Thus, demographics about SICS students would not be enough. Moreover, the university does not maintain or provide SICS-specific demographics.

According to recent statistics, the total number of students is ≈ 36,000 of whom 54% are female, 44.6% are male, plus 1.4% are non-binary or unstated. In terms of ethnicity, the rough breakdown is: 34% Asian, 24% Hispanic, 17% international, 15.44% White, 2.23%

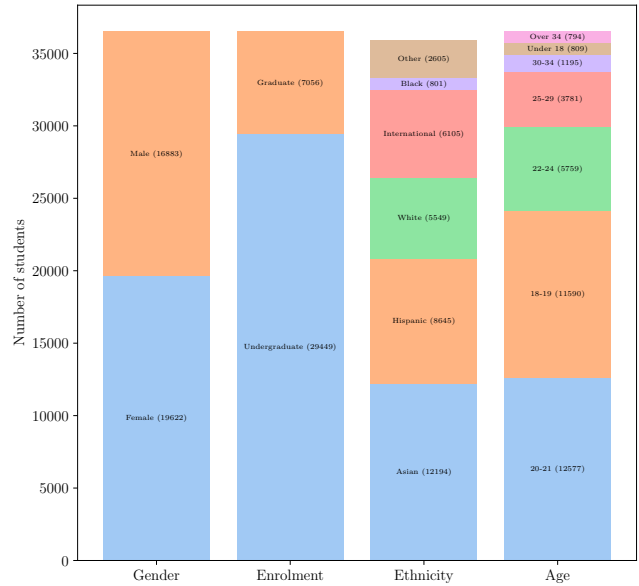


Figure 7: University Demographics

Black, and 7.25% other ethnic groups. The split between undergraduate and graduate students is 78.10% to 21.9%.

As far as the educational level, freshmen constitute 14% of the student body, sophomores – 15%, juniors – 21%, and seniors – 28%. The rest (≈ 22%) are graduate students. Interestingly, the age range of the student population is very wide, ranging from under 18 to over 64. Nonetheless, the majority (82%) fall into the 18 – 24 age range.

We also consider demographics of the 108 participants who engaged in the post-study survey (see Figure 8). The gender split is 62 (57.4%) male, 44 (40.7%) female, and 2 (1.9%) non-binary. The age of participants ranges from 18 to 30 with the majority (87.04%) under 25. Participants were also asked about their highest level of education. All participants have at least a high school degree. 58 participants (53.7%) are undergraduates and 50 (46.3%) are graduate students. All participants use the Internet daily and the main purpose of Internet usage for the majority (57.4%) is education. Finally, the country of residence for most (82.4%) participants is the United States, which is directly in line with the 17% international students from the overall university demographics.

Unfortunately, similarly detailed demographics for participants who solved reCAPTCHA2 as part of the main live experiment are unknown. However, the demographics of the 108 who participated in the post-study survey, closely resemble those of the overall campus total population in terms of gender, age, and educational level. Therefore, it is reasonable to assume that the demographics of all participants are the same, or very similar.

4.2 reCAPTCHA2 Dashboard Data

Google provides reCAPTCHA2 analytic data for website operators via a dashboard [18]. With it, website operators can generate a key-pair necessary for implementing reCAPTCHA2 on a web page. Difficulty setting can also be chosen on the dashboard. We used the

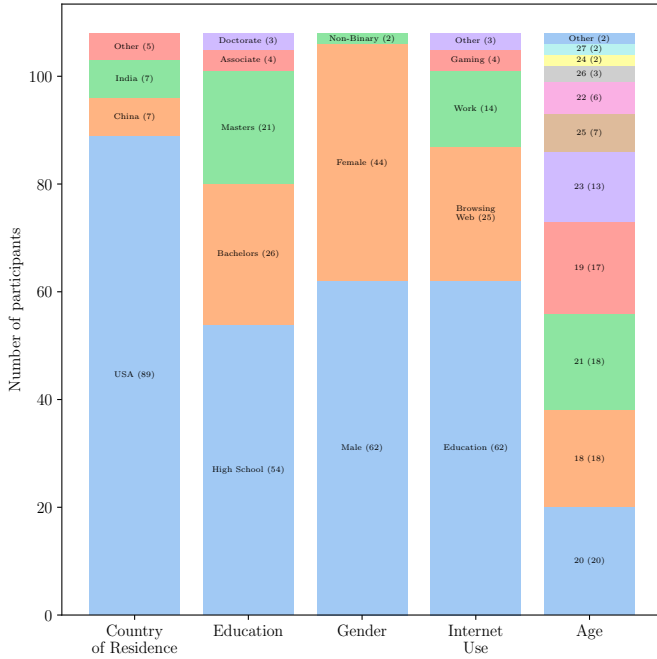


Figure 8: Demographics of post-study survey participants

“easy” setting in all experiments. The admin console allows for data to be downloaded in CSV format with the following fields per day:

no CAPTCHAs, Passed CAPTCHAs, Failed CAPTCHAs, Total Sessions, Failed Sessions, Average Score, and Average Response Time.

Table 2: Google’s reCAPTCHA dashboard data

no CAPTCHAs (checkbox)	7629
Passed CAPTCHAs (Image)	1890
Failed CAPTCHAs (Image)	143
Total Sessions	9538
Failed Sessions	19
Image accuracy	92.96%
Behavior accuracy	79.98%

Average score and response time are highly sparse and only appear on days with over 400 total sessions. Table 2 shows a sum for all days when data was collected over the entire study period. The image accuracy of 93% is computed as:

$$\frac{(\#passed\ CAPTCHAS)}{(\#passed\ CAPTCHAS + \#failed\ CAPTCHAS)}$$

The behavioral accuracy of 80% is computed as:

$$\frac{(\#ofCAPTCHAS)}{(total\#sessions)}$$

Notably, there are 9,538 CAPTCHA sessions reported by the admin console data, while we were supplied with 9,169 sessions, meaning that 369 form submissions has incomplete data or resulted in an

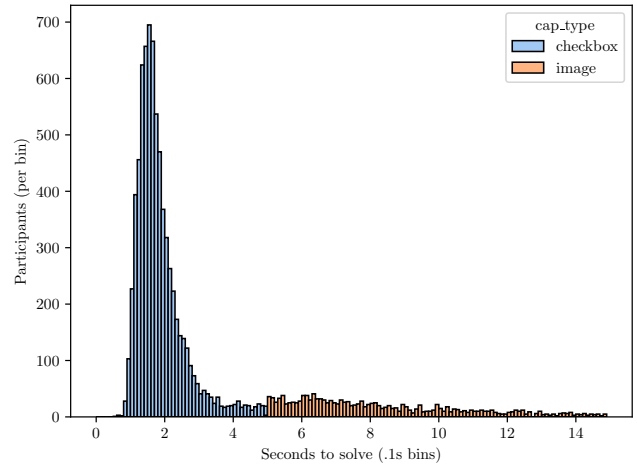


Figure 9: Timing results in bins of .1 seconds

error. This is likely due to incomplete sessions, e.g., refreshing before validation, or other form submission errors.

Table 3: Agglomerated solving time for reCAPTCHA Mode

mode	Count	Mean	Median	Std	Var	Max	Min
behavior	7334	1.85	1.67	0.71	0.50	4.99	0.51
image	1807	10.3	8.20	6.54	42.8	59.8	4.99
total	9141	3.53	1.83	4.50	20.3	59.8	0.51

4.3 Solving Time

Solving time for reCAPTCHA v2 is measured from the initial display to the successful verification. Data for solving time is split based on behavioral accuracy of 80% in Table 2. Since all tasks require a checkbox and some also require an image task, we assume that the 80% fastest solving times correspond to checkbox interactions. This split is also noted in the recent work by Searles, et. al [63]. All timing for image-based results is therefore a combination of check-box and image tasks.

Table 3 shows the results of 7,334 behavior and 1,807 images based on this split. The mean solving time for behavioral CAPTCHAs is 1.85 seconds, while the image mean solving time is 10.3 seconds. The latter corresponds to a notable 557% increase.

Looking at Figure 9, there is a sharp drop-off in solving time starting around 2, and ending at 5, seconds: it hits a low and then goes back up slightly. The split point for image and behavior is about 5 seconds, which matches the drop-off point, thus strengthening the accuracy of the split. Figure 10 shows timing results after the image split. Notably, image and checkbox data follow similar patterns of distribution.

Solving time can also be partitioned into the following dimensions, based on collected data: Service, Attempts, Educational Level, and Major. This is done separately for image and behavior, across those listed dimensions in Tables 4, 5, 7, 8, 9, 10, 12 and 13.

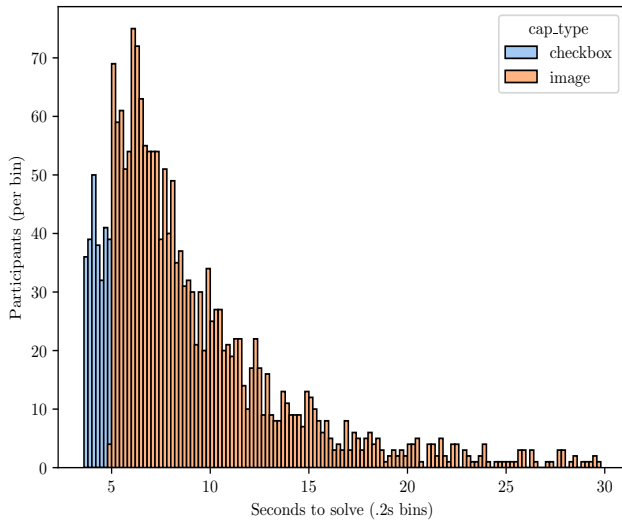


Figure 10: Image timing results in bins of .2 seconds

4.3.1 *Statistical Testing.* For the sake of statistical validity, we apply a series of standard statistical tests to solving times. We run all of the following statistical tests on both image and checkbox data separately. Statistical methods were applied using python’s `scipy` [23] library. With a null hypothesis that solving times adhere to a normal distribution, we performed the *Shapiro-Wilk normality test*. For both image and checkbox cases, results showed that we can reject the null hypothesis ($p < 0.001$). With a null hypothesis that the skewness is the same as that of a corresponding normal distribution, we ran the timing data with *skewtest*. For both image and checkbox, results reject the null hypothesis in favor of the alternative: the distribution of solving times is skewed ($p < 0.001$) to the right. With a null hypothesis that the kurtosis is the same as that of a normal distribution, we used the *tailedness test*. For both image and checkbox, results show the samples were drawn from a population that has a heavy-tailed distribution ($p < 0.001$). We used the *Brown Forsythe test* to compare equality of variance between image and checkbox, which shows that they do not exhibit equal variance. We used the *Kruskal-Wallis test* with the Holm-Bonferroni method to adjust for family-wise error in order to test the equality of mean between modes, services, attempts, majors, and educational level. Significant results are included in Figures 11, 12, and 13.

Table 4: Checkbox solving time in seconds for each service

Service	Count	Mean	Median	Std	Var	Max	Min
Password Reset	2654	1.67	1.51	0.65	0.42	4.99	0.51
Account Creation	4680	1.96	1.76	0.71	0.51	4.97	0.86

4.3.2 *Services.* As mentioned earlier, the website had two services that invoked reCAPTCHA: password recovery and account creation. Tables 4, 5, and 6 show results from these two CAPTCHA interactions. There were 6, 155 account creation, and 2, 986 password

Table 5: Image solving time in seconds for each service

Service	Count	Mean	Median	Std	Var	Max	Min
Password Reset	332	10.4	8.01	6.59	43.5	43.5	5.01
Account Creation	1475	10.3	8.23	6.53	42.7	59.8	4.99

Table 6: Total solving time in seconds for each service

Service	Count	Mean	Median	Std	Var	Max	Min
Password Reset	2986	2.63	1.58	3.56	12.7	43.5	0.51
Account Creation	6155	3.97	2.00	4.84	23.4	59.8	0.86

recovery, form submissions. Notably, for behavioral results, the Kruskal-Wallis test shows statistically significant differences between account creation and password recovery with a $p = 1.1e^{-115}$. Students who interacted with the account creation service solved behavioral CAPTCHAS 17% slower than those who interacted with the password recovery service. Additionally, 50% more time was spent solving reCAPTCHA during account creation than during password recovery. Total results are also statistically significant with $p = 6.7e^{-162}$. However, since 90% of students who interacted with the latter have already interacted with the account creation service, these results may be conflated by multiple prior attempts. For the image case, the Kruskal-Wallis Test yielded no statistically significant results.

Table 7: Solving time for number of checkbox attempts

Attempt	Count	Mean	Median	Std	Var	Max	Min
1	2888	2.02	1.80	0.73	0.54	4.97	0.94
2	1293	1.84	1.67	0.65	0.42	4.97	0.62
3	751	1.80	1.63	0.66	0.44	4.95	0.80
4	513	1.73	1.55	0.63	0.40	4.89	0.78
5	371	1.73	1.57	0.70	0.49	4.92	0.89
6	272	1.61	1.47	0.58	0.34	4.57	0.84
7	212	1.67	1.52	0.65	0.43	4.90	0.64
8	167	1.66	1.52	0.65	0.43	4.65	0.84
9	127	1.60	1.48	0.57	0.33	4.09	0.88
10	112	1.56	1.44	0.63	0.39	4.97	0.85
11	94	1.63	1.41	0.76	0.57	4.90	0.88
12	67	1.61	1.46	0.68	0.46	4.47	0.51
13	52	1.58	1.37	0.70	0.49	4.49	0.96
14	37	1.53	1.45	0.63	0.40	4.62	0.92
15	28	1.51	1.41	0.56	0.31	3.88	0.88

Table 8: Solving time for number of image attempts

Attempt	Count	Mean	Median	Std	Var	Max	Min
1	1264	10.5	8.36	6.60	43.5	58.9	4.99
2	260	10.9	8.16	7.47	55.8	55.5	5.00
3	93	9.30	8.16	4.09	16.7	29.2	5.00
4	45	10.0	7.77	8.41	70.7	59.8	5.21
5	25	8.76	7.48	4.56	20.8	26.4	5.12
6	15	7.26	6.06	2.33	5.44	12.3	5.18

4.3.3 *Attempts.* Interestingly, some participants submitted forms multiple times. For behavior-based challenges, the average number of attempts was 3.52, and 1.73 for image-based ones. Tables 7 and

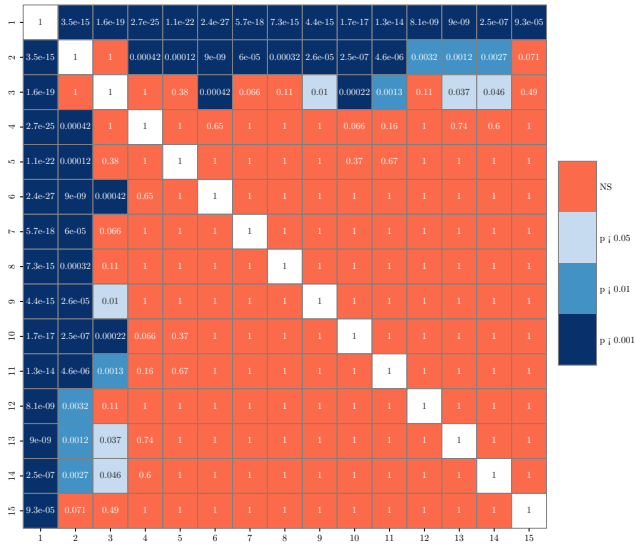


Figure 11: Kruskal-Wallis results for checkbox attempts

8 show timing results over multiple attempts. The highest number of attempts was 37 for behavior-based challenges and 20 for image-based ones. Behavioral results from the Kruskal-Wallis test in Figure 11 show that there is a statistically significant difference between the first and subsequent attempts ($p < .001$). While for the second attempt there is a statistically significant difference ($p < .001$) between all other attempts except the third. In general, this data shows that checkbox solving time decreases with more attempts, meaning that humans improve at solving checkbox challenges.

We observe an interesting behavioral phenomena whereby participants react faster when they know what to expect. However, average image results show a slight increase on the second attempt, while subsequent attempts decrease. This may be attributed to reCAPTCHA2 presenting a more difficult challenge on the second attempt. Image results from the Kruskal-Wallis test show no statistically significant differences between image attempts. This is likely due to the drop-off in the number of participants who solved multiple image challenges.

Table 9: Checkbox solving time for educational levels

Level	Count	Mean	Median	Std	Var	Max	Min	Total %
Freshmen	479	1.89	1.68	0.77	0.59	4.96	0.95	62.1%
Sophomore	1198	1.93	1.73	0.73	0.54	4.99	0.91	71.3%
Junior	1912	1.84	1.66	0.69	0.48	4.96	0.51	85.2%
Senior	2399	1.79	1.63	0.68	0.46	4.95	0.64	87.4%
Graduate	932	1.95	1.76	0.70	0.49	4.96	0.91	78.6%

4.3.4 Educational Level. Educational level was obtained via the website crawler, as described in Section 3.4. Tables 9, 10, and 11 present data for different educational levels. For the split checkbox and image data there are only very minor differences between solving times based on the educational level. In terms of statistical

Table 10: Image solving time for different educational levels

Level	Count	Mean	Median	Std	Var	Max	Min	Total %
Freshmen	294	10.5	8.42	6.24	38.9	56.2	5.01	37.9%
Sophomore	483	10.3	8.08	7.28	52.9	59.8	5.00	28.7%
Junior	334	10.3	8.18	6.00	36.0	45.4	5.00	14.8%
Senior	346	10.2	8.05	6.20	38.5	43.9	4.99	12.6%
Graduate	254	10.7	8.49	6.85	46.9	50.0	5.01	21.4%

Table 11: Total solving time for different educational levels

Level	Count	Mean	Median	Std	Var	Max	Min
Freshmen	773	5.15	2.33	5.69	32.4	56.2	0.95
Sophomore	1681	4.33	2.05	5.47	29.9	59.8	0.91
Junior	2246	3.09	1.77	3.84	14.7	45.4	0.51
Senior	2745	2.85	1.71	3.62	13.1	43.9	0.64
Graduate	1186	3.82	1.97	4.83	23.3	50.0	0.91

Table 12: Checkbox solving time for various majors

Major	Count	Mean	Median	Std	Var	Max	Min
CmptSci	2658	1.80	1.63	0.66	0.44	4.99	0.62
CSE	754	1.83	1.63	0.77	0.59	4.97	0.64
SW Engr	635	1.79	1.59	0.70	0.50	4.92	0.51
Undclrd	613	1.91	1.70	0.75	0.56	4.97	0.95
MCS	326	2.01	1.85	0.68	0.47	4.96	0.91
DataSci	273	1.98	1.74	0.80	0.65	4.96	1.03
IN4MATX	225	1.88	1.70	0.66	0.44	4.96	1.01
BIM	171	1.96	1.73	0.81	0.66	4.94	0.89
GameDes	139	1.80	1.60	0.74	0.55	4.79	0.77
Math	121	1.84	1.72	0.59	0.35	3.93	1.00
MofData	102	1.90	1.70	0.72	0.51	4.60	1.03
EngrCpE	82	1.90	1.74	0.66	0.44	4.36	0.98
PSW ENG	79	2.03	1.81	0.72	0.52	4.69	0.91
Bus Adm	77	1.92	1.77	0.78	0.61	4.78	0.88
CSGames	77	1.73	1.55	0.66	0.44	4.76	0.88
BusEcon	62	1.98	1.77	0.66	0.44	4.95	0.95
Bio Sci	60	1.89	1.69	0.74	0.55	4.89	0.99
Stats	47	1.73	1.56	0.61	0.37	4.59	1.11
Cog Sci	37	1.97	1.94	0.56	0.31	3.39	0.97
Net Sys	34	2.13	1.96	0.69	0.47	3.92	1.33
Psych	29	1.70	1.59	0.44	0.19	2.62	1.05
Engr ME	26	2.06	1.85	0.54	0.29	3.44	1.39

significance, Figure 12 shows statistically significant differences in total solving time for all educational levels. In terms of total time, freshmen are the slowest – 80% slower than seniors. There is a direct trend from freshman to seniors showing a reduction in solving time. Similarly, there is a trend of the total ratio of image to checkbox challenges.

4.3.5 Majors. Majors of the study participants (i.e., disciplines they study) were obtained through the website crawler, as described in Section 3.4. Tables 12, 14, and 13 present solving times for participants with various majors. Although there are 62 majors in total, Tables 12, 14, and 13 only show 22 majors. This is because each of the remaining 40 majors had < 20 reCAPTCHA2 sessions. As the Kruskal-Wallis test in Figure 13 shows, only 8 majors had statistically significant differences in terms of checkbox solving behavior. Among these, Computer Science had the lowest, and Informatics – the highest, total average solving time.



Figure 12: Kruskal-Wallis results for total and educational level

Table 13: Image solving time for various majors

Major	Count	Mean	Median	Std	Var	Max	Min
CmptSci	527	10.2	8.13	6.12	37.4	44.5	4.99
Undclrd	239	11.0	8.15	8.25	68.0	59.8	5.02
CSE	196	9.94	8.26	5.70	32.5	42.0	5.03
SW Engr	161	9.64	7.63	5.51	30.3	45.4	5.04
DataSci	90	10.2	8.51	5.84	34.1	41.2	5.01
MCS	78	10.5	8.65	6.01	36.2	38.1	5.05
IN4MATX	62	12.4	8.37	9.82	96.4	50.9	5.14
BIM	55	9.47	8.35	4.25	18.0	25.5	5.02
GameDes	47	12.0	8.58	10.8	117	56.2	5.16
MofData	29	9.71	8.99	4.48	20.1	25.6	5.39
Math	28	10.6	8.98	5.08	25.8	28.9	5.34
EngrCpE	24	10.9	9.11	4.27	18.2	20.1	5.42
PSW ENG	18	8.98	7.58	4.30	18.5	21.3	5.01
Stats	18	8.84	6.89	4.60	21.1	23.9	5.73
BusEcon	16	9.74	9.39	4.45	19.8	21.3	5.11
Bio Sci	15	12.0	11.2	5.70	32.5	23.0	5.72
Bus Adm	12	8.78	8.87	1.87	3.51	13.1	5.56
CSGames	8	10.7	10.1	4.04	16.4	16.1	5.90
Engr ME	8	9.26	7.85	4.20	17.6	18.4	5.38
Cog Sci	7	8.15	6.66	4.17	17.4	16.8	5.00
Net Sys	5	21.0	9.47	18.5	343	50.0	7.16
Pysch	3	6.41	6.29	1.18	1.40	7.65	5.30

4.4 Survey Results

We now discuss the study results pertaining to usability, preferences, and opinions about reCAPTCHA2. An interactive version of the google form we used is available at [20]. 800 randomly selected study participants were contacted by email, with the goal of obtaining at least 100 respondents. In the end, a total of 108 completed the survey. Two solving scenarios are considered:

Checkbox only Only the checkbox challenge: after clicking the checkbox, no image challenge was served. This applies to 42 participants.

Table 14: Total solving time for various majors

Major	Count	Mean	Median	Std	Var	Max	Min
CmptSci	3185	3.19	1.75	4.05	16.4	44.5	0.62
CSE	950	3.51	1.81	4.23	17.9	42.0	0.64
Undclrd	850	4.47	2.03	6.02	36.2	59.8	0.95
SW Engr	796	3.38	1.75	4.06	16.5	45.4	0.51
MCS	404	3.65	2.08	4.32	18.7	38.1	0.91
DataSci	362	3.98	2.02	4.55	20.7	41.2	1.03
IN4MATX	287	4.14	1.89	6.29	39.5	50.9	1.01
BIM	226	3.79	1.97	3.91	15.3	25.5	0.89
GameDes	186	4.38	1.86	7.03	49.4	56.2	0.77
Math	147	3.50	1.89	4.11	16.9	28.9	1.00
MofData	131	3.63	1.94	3.92	15.3	25.6	1.03
EngrCpE	106	3.93	1.98	4.31	18.6	20.1	0.98
PSW ENG	97	3.32	1.93	3.33	11.1	21.3	0.91
Bus Adm	89	2.85	1.83	2.55	6.52	13.1	0.88
CSGames	85	2.43	1.61	2.60	6.77	18.4	0.88
BusEcon	78	3.57	2.06	3.76	14.1	21.3	0.95
Bio Sci	75	3.90	1.87	4.80	23.0	23.0	0.99
Stats	65	3.70	1.68	4.02	16.2	23.9	1.11
Cog Sci	44	2.96	2.02	2.81	7.90	16.8	0.97
Net Sys	39	4.55	2.07	8.81	77.6	50.0	1.33
Engr ME	34	4.08	2.16	4.18	17.5	16.1	1.39
Psych	32	2.14	1.65	1.49	2.21	7.65	1.05

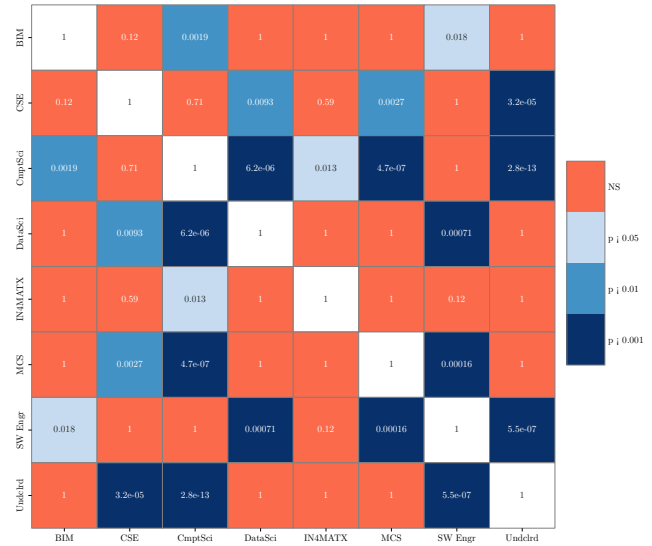


Figure 13: Kruskal-Wallis results for total and major

Checkbox+image Both checkbox and image challenges: after clicking the checkbox, an image challenge was served. This applies to 66 participants.

4.4.1 System Usability Scale (SUS) Score Analysis. Table 15 reports the SUS score for both scenarios. Results from individual SUS statements are not analyzed, since they do not provide meaningful information [28, 31].

SUS checkbox scores are: 78.51 for checkbox only, and 76.21 for checkbox+image. Referring to Table 1, the usability level for checkbox in both scenarios is “Good”. We thus conclude that for checkbox, the SUS score and the usability level do not vary depending on the solving scenario, i.e., whether or not an image challenge is served afterwards. On the other hand, the SUS score of image

Table 15: SUS Scores for reCAPTCHA v2

Solving Scenario	reCAPTCHA Type	SUS Score
Checkbox only	Checkbox	78.51
Checkbox+image	Checkbox	76.21
Checkbox+image	Image	58.90

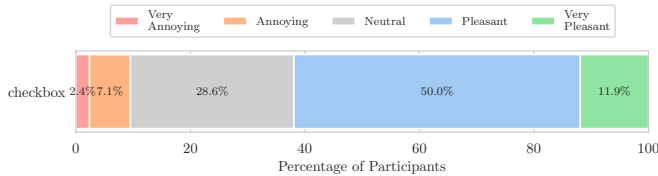


Figure 14: Preference score for checkbox only scenario

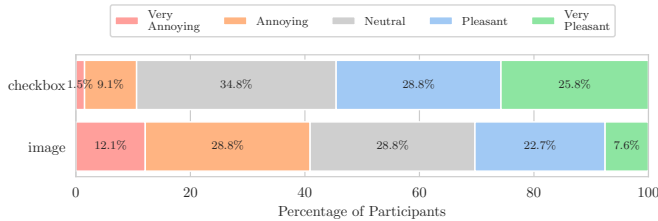


Figure 15: Preference score for checkbox+image scenario

is 58.90 and the usability level is “OK”. This difference is likely influenced by the difficulty of the task, since clicking a checkbox is surely much simpler than classifying an image. We observed that solving image challenges takes 557% longer than checkbox.

4.4.2 Preference Analysis. Besides SUS questions, the post-study questionnaire asked about participants’ preferences regarding checkbox and image versions. Specifically, they were asked to provide opinions using a custom scale. Figure 14 and Figure 15 show the preferences in both scenarios.

Majority of participants in both scenarios (61.9% and 54.6%, respectively) find checkbox either “pleasant” or “very pleasant”. Whereas, only a minority (30.3%) find image “pleasant” or “very pleasant”. A significantly larger percentage of participants think that image is “annoying” or “very annoying”. Similar to SUS scores, the preference for checkbox does not change when it is followed by an image.

We also compute quantitative scores in order to rate reCAPTCHA v2 based on participants’ preferences. For that, the preference scale is converted into a five-point Likert scale with “very annoying” corresponding to 1 and “very pleasant” corresponding to 5. The rating of checkbox is: 3.62 for checkbox only, and 3.68 for checkbox+image, scenario. Similar to the SUS score, the rating of checkbox is independent of the solving scenario. The rating of image is appreciably lower, at 2.84.

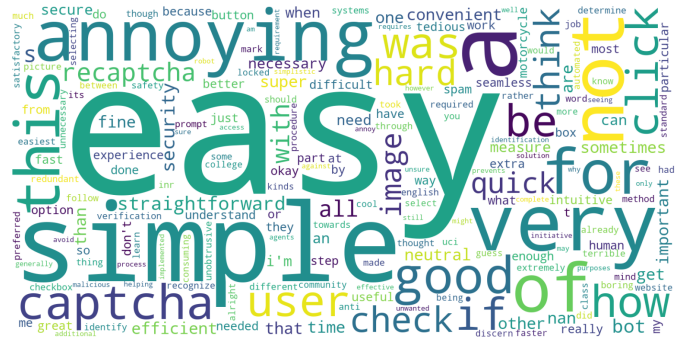


Figure 16: Word cloud from feedback on checkbox



Figure 17: Word cloud from feedback on image

Moreover, comparing preference scores from Figure 14 and Figure 15 with SUS scores from Table 15 we observe a trend for both checkbox and image: checkbox is more usable and rated positively, while image is less usable and rated negatively. This leads to an unsurprising conclusion that participants’ preference for a given reCAPTCHA v2 type is correlated with its usability level.

4.4.3 Qualitative Feedback. In the final part of the survey, participants were asked to provide open-ended feedback about checkbox and image using at least one word. Using collected feedback, we generate word clouds for both.

The most prominent words for checkbox challenges in Figure 16 are “easy” and “simple”. Other significant positive words are “good” and “quick”. Nevertheless, checkbox is still labeled “hard” and “annoying” by some participants.

Figure 17 shows that the most prominent word describing image is “annoying”, while a small fraction of participants label it as “good” and “easy”. We acknowledge that the custom scale used in the survey might possibly introduce bias toward the word “annoying”. However, the scale also includes the word “pleasant” and that word is not present in the word cloud as a positive opinion. Instead, participants used other positive-sounding words, such as “easy” or “simple”.

Negative words from the checkbox cloud and positive ones from the image cloud indicate that neither reCAPTCHA v2 type is universally liked or disliked.

5 COMPARISON WITH RELATED WORK

Prior results include [29, 32, 35, 37, 41, 42, 46, 50, 54, 59, 62, 63, 68, 69]. They present average solving times for various CAPTCHAs, ranging from 3.1 to 47 seconds. Compared with our observed mean solving time of 1.8 seconds for checkbox reCAPTCHA_{v2}, previous results are 1.7 to 26 times slower. For image reCAPTCHA_{v2}, our mean solving time is 10 seconds, which is 3.3 times slower than the fastest, and 5 times faster than the slowest, previously reported results. The fast solving time may be related to the trend noted in [32, 63] of age influencing solving time: younger participants seem to solve faster than older ones. Since our population is mostly university students (aged 18 – 25), our results re-confirm this trend.

Table 16: Comparison with results from prior user studies evaluating reCAPTCHA_{v2}: checkbox (C), image (I), total (T). Mean in seconds

Study	Unique users	reCAPTCHA _{v2} s solved	Mean	Accuracy
Ours	3,625	9,141	10.4 (I), 1.85 (C), 3.53 (T)	93% (I), 80% (C)
[63]	1,400	2,800	14-26 (I), 3.1-4.9 (C)	71-81% (I), 71-85% (C)
[68]	40	40	3.1 (T)	None

To the best of our knowledge, only two prior efforts studied reCAPTCHA_{v2}: [63] and [68]. Table 16 shows a direct comparison of the results. However, [68] provides a very limited data set of ($n = 40$) reCAPTCHA_{v2}, containing only total times. Whereas, [63], provides the following points of comparison:

- (1) Amazon Mturk vs “real world” participants
- (2) Participant awareness vs unawareness of the study existence and purpose
- (3) Mock vs real account creation;
- (4) Preferences/Rating

Webb et al. [73] reported several points of concern about the quality of data collected from MTurk [15]. Our data and results are derived from a real-world scenario of actual users creating real accounts for a real service. However, since both this work and [63] implement reCAPTCHA_{v2} in a similar way, some interesting conclusions can be drawn regarding the efficacy of Mturk data. Mturk users in [63] solved easy checkbox challenges 1.7 – 2.7 times slower than our participants. They also solved easy image challenges 1.6 – 2.6 times slower than our participants. Another consideration is network speed, since MTurkers were participating in the [63] study over the Internet. In contrast, our study was conducted with most² participants being in close network proximity. Therefore, it would explain why Mturk results are slower since they can originate anywhere in the world, according to demographics reported in [63]. This may also skew our results to be faster than the actual total reCAPTCHA_{v2} solving time.

[63] showed that participants’ awareness of the true purpose of the study could alter solving times. The solving time of participants who thought that they were participating in an account

²Recall that VPN use was required to create an account or recover a password, thus taking part in our study. Although most participants were on campus, some were remote. The exact number of the latter is unknown.

creation study was up to 57.5% slower than those who knew that they were participating in a study about solving CAPTCHAs. Account creation solving times (in seconds) for easy reCAPTCHA_{v2} of [63] are 4.9 for checkbox, and 26.3 for image. In contrast, our results are 2.02 for checkbox and 10.5 for image (on the first attempt). This translates into an average of 2.5 times slowdown for both challenge types. On the first attempt, our participants show the slowest mean and the least awareness (no study information presented) for checkbox challenges across significant groups ($n = 2, 888$), and our results show that solving time improves with subsequent attempts. Whereas, in [63] participants solved 10 CAPTCHAs (among them 2 reCAPTCHA_{v2}s) in sequence, which could lower the timing due to the repeated attempts bias.

Also, [63] observed a lot of task abandonment, which might be due to the mocked-up (fake) account creation in that study. This is unlike our case where participants must create accounts due to the SICS school-wide policy. Thus, they must complete the form with successful post-validation by the back-end server. (In other words, abandonment is not an option).

[63] did not validate form information during account creation form submission beyond checking form field constraints, which could significantly alter the user study experience. Since our high average multiple attempts per participant of 3.52 for checkbox and 1.73 for image was likely due to failed post-validation by the back-end server.

Study participants in [63] rated reCAPTCHA_{v2} on a Likert scale, from “least enjoyable” to “most enjoyable”. Results showed that checkbox was the most enjoyable out of all CAPTCHAs, while image challenges were the least so. The term “enjoyable” is synonymous with pleasant (the opposite of “annoying”), which presents a point of comparison. Our results in Figures 14 and 15 are very similar in terms of positive and negative responses, thereby confirming results of [63]

6 DISCUSSION

6.1 Cost Analysis

We now attempt to quantify various costs incurred by global use of reCAPTCHA on the internet. In particular we want to estimate the total time spent solving reCAPTCHAs, the overall amount of human labor, network traffic (bandwidth), power consumption and the consequent environmental impact. Note that, in the informal analysis below, we consider all estimates to be generous lower bounds.

Given that historic average solving time for distorted-text CAPTCHAs (same type used by reCAPTCHA v1) was 9.8 seconds and the conservative rate of 100 million reCAPTCHAs per day [5], 980 million seconds per day were spent solving reCAPTCHA v1s. For reCAPTCHA v1, it lived from 2009-2014 for 5 years amounting to 183 billion reCAPTCHA v1 sessions, taking 1.79 trillion seconds, or 497 million hours of human time spent solving reCAPTCHA v1. Given that the US federal minimum wage is \$7.5, this roughly yields \$3.7 billion in free wages.

Given that average solving time for all reCAPTCHA_{v2} sessions is 3.53 seconds and the conservative rate of 100 million reCAPTCHAs per day [5], 353 million seconds per day are spent solving reCAPTCHA_{v2}s. For reCAPTCHA_{v2}, it has been 9 years amounting

to 329 billion reCAPTCHA sessions taking 1.16 trillion seconds, or 322 million hours of human time spent solving reCAPTCHA. Given that the US federal minimum wage is \$7.5, this roughly yields \$2.4 billion in free wages.

Assuming un-cached scenarios from our technical analysis (see Appendix B), network bandwidth overhead is 408 KB per session. This translates into 134 trillion KB or 134 Petabytes (194 x 1024 Terrabytes) of bandwidth. A recent (2017) survey [27] estimated that the cost of energy for network data transmission was 0.06 kWh/GB (Kilowatt hours per Gigabyte). Based on this rate, we estimate that 7.5 million kWh of energy was used on just the network transmission of reCAPTCHA data. This does not include client or server related energy costs. Based on the rates provided by the US Environmental Protection Agency (EPA) [24] and US Energy Information Administration (EIA) [17], 1 kWh roughly equals 1-2.4 pounds of CO2 pollution. This implies that reCAPTCHA bandwidth consumption alone produced in the range of 7.5-18 million pounds of CO2 pollution over 9 years.

In total from reCAPTCHA v1 and reCAPTCHA: There have been at least 512 billion reCAPTCHA sessions taking 2.95 trillion seconds, or 819 million hours, which is at least \$6.1 billion USD in free wages. Out of the 329 billion reCAPTCHA sessions, (Our rate of 20%) at least 65.8 billion would have been image challenges, while 263.2 million would have been checkbox challenges. Thus 250 billion challenges would have resulted in labeled data. According to Google, the value of 1,000 items of labeled data is in the \$35-129 USD range [14], which would be worth at least \$8.75-32.3 billion USD per each sale.

Lastly, we look into the economics of tracking cookies, another main by-product of reCAPTCHA. Tracking cookies play an ever-increasing role in the rapidly growing online advertisement market. According to Forbes [3], digital ad spending reached over \$491 billion globally in 2021, and more than half of the market (51%) heavily relied on third-party cookies for advertisement strategies [1]. The expenditure on third-party audience data (collected using tracking cookies) in the United States reached from \$15.9 billion in 2017 to \$22 billion in 2021 [2]. More concretely, the current average value life-time of a cookie is €2.52 or \$2.7 [58]. Given that there have been at least 329 billion reCAPTCHA sessions, which created tracking cookies, that would put the estimated value of those cookies at \$888 billion dollars.

6.2 Security Analysis

In the following subsections 6.3 and 6.4, we discuss different attacks that have been performed successfully against reCAPTCHA. We consider behavior-based, image, and audio challenges in reCAPTCHA and reCAPTCHA3. Table 17 shows a direct comparison of the time and accuracy for humans and bots.

Table 17: Humans vs. bot solving time (seconds) and accuracy (percentage) for reCAPTCHA.

Type	Human				Bot	
	Time	Acc	Time	Acc	Time	Acc
checkbox	1.85	80%	3.1-4.9 [63]	85% [63]	1.4 [66]	100% [66]
image	10.4	93%	16-26 [63]	81% [63]	17.5 [49]	85% [49]

6.3 reCAPTCHA

reCAPTCHA presents three different types of captcha challenges to the users: behavior-based (checkbox) challenge, image challenge, and audio challenge. Unfortunately, each of these captcha types has been proven vulnerable to attacks.

6.3.1 Checkbox Challenge. With the introduction of reCAPTCHA, came a new serious vulnerability in the form of click-jacking [47]. Adversaries can make "trustworthy" users generate g-recaptcha-response-s, which can be automatically used to pass challenges, ultimately making a Bot's job infinitely easier!

Sivakorn, et al. [66] perform an in-depth analysis of the risk analysis system of reCAPTCHA and implement an attack to manipulate it. Based on this analysis and implementation:

- (1) Google primarily uses tracking cookies in the risk analysis system.
- (2) At least 63,000 valid cookies can be automatically created per day per IP address.
- (3) 9 days after a cookie creation, checkbox attempts using the cookie will succeed.
- (4) 52,000-59,000 checkbox challenges can be solved with 100% accuracy per day per IP address.
- (5) The average solution time is 1.4 seconds with 100% accuracy, shown in Table 17.

Given the blatant vulnerability [47], ease of implementing large-scale automation [66], and usage of privacy invasive tracking cookies reCAPTCHA checkbox presents itself as a complete vulnerability disguised as a security tool. Google was previously sued 22.5 million for **secretly** adding tracking cookies to apple users devices [8]. It can be concluded that the true purpose of reCAPTCHA is as a tracking cookie farm for advertising profit masquerading as a security service.

6.3.2 Image Challenge. Image-labeling challenges have been around since 2004 with the introduction of Image Recognition CAPTCHAs by Chew et al. [34]. 6 years later, in 2010 Fritsch et al. [38] published an attack that beat the prevalent image CAPTCHAs of the time with 100% accuracy. At this point, it could be concluded that image recognition was no longer difficult to solve automatically with a computer. However in 2014 with the introduction of reCAPTCHA, the fall-back security method was an image challenge, which had been proven insecure 4 years prior. The idea is that if your cookies aren't valuable enough then reCAPTCHA would present an image labeling task. This wouldn't make sense as a security service, yet it would make sense given that obtaining labeled image data is highly valuable and is even sold by Google [14]. The conclusion can be extended that the true purpose of reCAPTCHA is a free image-labeling labor and tracking cookie farm for advertising and data profit masquerading as a security service.

Consequently, [66] and [49] investigate and successfully implement automated solutions to reCAPTCHA's image labeling task. In 2016, [66] showed that a plethora of automated services, including Google's own Google Reverse Image Search (GRIS), could be used to automatically complete reCAPTCHA's image labeling tasks. [66] also implemented its own easy solver with 70.8% accuracy at 19.2 seconds per reCAPTCHA image labeling task. In 2020, [49] also showed that many automated services, including Google's

own Google Cloud Vision, could be used to automatically complete reCAPTCHA v2s image labeling tasks with reasonable speed and accuracy. [49] similarly implemented an attack, achieving a high level of speed (17.5 seconds) and accuracy (85%), shown in Table 17.

6.3.3 Audio Challenge. As part of reCAPTCHA v2, Google introduced accessibility options allowing users to use audio CAPTCHAs, instead of image-based ones. Unsurprisingly, these audio CAPTCHAs introduce an accessibility side-channel, especially apparent due to advances in speech-to-text technology.

In 2017, Bock, et al. [30] introduced an automated system called *unCaptcha* which can solve audio challenges with 85.15% accuracy and 5.42 seconds average solving time. Similar to other attacks [30] uses Google’s own voice recognition technology as a means to break audio challenges.

6.4 reCAPTCHA v3

reCAPTCHA v3 was introduced in 2018 [7] proposing the returning of a score, which website developers could use to decide whether to prompt with a challenge or perform some other action. Challenges types served by reCAPTCHA v3 are the same as reCAPTCHA v2. Also, there is no discernible difference between reCAPTCHA v2 and reCAPTCHA v3 in terms of appearance or perception of image challenges and audio challenges. Hence, attacks targeting reCAPTCHA v2 image/audio challenges are also applicable for those of reCAPTCHA v3. However, assuming that the risk analysis system was updated from reCAPTCHA v2 to reCAPTCHA v3, breaking behavior-based challenges of reCAPTCHA v3 might require new techniques. In 2019, Akrou, et al. [25] presented a reinforcement learning (RL) based attack breaking reCAPTCHA v3’s behavior-based challenges, obtaining high scores (.9+), with 97% accuracy and only requiring 2,000 data points as a training set.

7 SUMMARY

Over 13 years passed since reCAPTCHA’s initial appearance and its current prevalence is undeniable. It is thus both timely and important to investigate its usability. This paper presents a real-world user study with over 3,600 unbiased (unwitting) participants solving over 9,000 reCAPTCHA v2 challenges. We explore four new dimensions of reCAPTCHA v2 solving time: # of attempts, service type, as well as educational level and major. Results show that:

- Participants improve in terms of solving time with more attempts, for checkbox challenges.
- The service/website setting is an important consideration for researchers and web developers, since it has a statistically significant effect on solving time.
- Educational level directly impacts solving time.
- There were minor trends with statistical significance of participants with technical (STEM) majors solving time being faster than that of others.

In terms of usability, the post-study survey results show that the checkbox challenge gets an average SUS score of 77. This is considered to be acceptable and preferred by many participants over the image challenge, which has an average SUS score of 59. Notably, participants found the image challenge to be annoying.

In terms of cost, we estimate that – during over 13 years of its deployment – 819 million hours of human time has been spent on

reCAPTCHA, which corresponds to at least \$6.1 billion USD in wages. Traffic resulting from reCAPTCHA consumed 134 Petabytes of bandwidth, which translates into about 7.5 million kWhs of energy, corresponding to 7.5 million pounds of CO₂. In addition, Google has potentially profited \$888 billion USD from cookies and \$8.75-32.3 billion USD per each sale of their total labeled data set.

In terms of security reCAPTCHA v2 presents:

- Click-jacking (a blatant vulnerability) [47]
- Trivial implementation of large-scale automation attacks [66]
- Weakness of security premise of fallback (image challenge) [38, 49, 66]
- Usage of privacy invasive tracking cookies (for security) [66]

Ultimately, given these points it can be concluded that reCAPTCHA v2 presents no real security.

Given that: (1) reCAPTCHA v2 is negatively perceived by most users, (2) its immense cost, and (3) its susceptibility to bots, our results prompt a natural conclusion:

reCAPTCHA v2 and similar reCAPTCHA technology should be deprecated.

REFERENCES

- [1] [n. d.]. Degree of reliance on third-party cookies in digital advertising in the United States as of July 2021. <https://www.statista.com/statistics/1222230/reliance-cookie-advertising-usa/>.
- [2] [n. d.]. Spending on third-party audience data supporting marketing related efforts in the United States from 2017 to 2021, by type. <https://www.statista.com/statistics/1202754/third-party-audience-data-spending-usa/>.
- [3] [n. d.]. The Truth In User Privacy And Targeted Ads. <https://www.forbes.com/sites/forbestechcouncil/2022/02/24/the-truth-in-user-privacy-and-targeted-ads/>.
- [4] 2009. Teaching computers to read: Google acquires reCAPTCHA. <https://web.archive.org/web/2012051113750/http://googleblog.blogspot.com/2009/09/teaching-computers-to-read-google.html>
- [5] 2010. recaptcha FAQ from 2010 archived. <https://web.archive.org/web/20100629174402/http://www.google.com:80/recaptcha/faq>
- [6] 2014. Are you a robot? Introducing “No CAPTCHA reCAPTCHA”. <https://security.googleblog.com/2014/12/are-you-robot-introducing-no-captcha.html>.
- [7] 2018. reCAPTCHA v3. <https://developers.google.com/search/blog/2018/10/introducing-recaptcha-v3-new-way-to>.
- [8] 2019. Google will pay 22.5 million to settle FTC charges it misrepresented privacy assurances to users of Apple’s Safari Internet Browser. <https://www.ftc.gov/news-events/news/press-releases/2012/08/google-will-pay-225-million-settle-ftc-charges-it-misrepresented-privacy-assurances-users-apples>
- [9] 2023. <https://www.google.com/chrome/>
- [10] 2023. <https://www.pingdom.com/>
- [11] 2023. <https://www.webpagetest.org/>
- [12] 2023. <https://www.jitbit.com/macro-recorder/mouse-recorder/>
- [13] 2023. <https://playwright.dev/>
- [14] 2023. AI Platform Data Labeling Service pricing. <https://cloud.google.com/ai-platform/data-labeling/pricing>
- [15] 2023. Amazon Mechanical Turk. <https://www.mturk.com/>.
- [16] 2023. CAPTCHA Usage Distribution on the Entire Internet. <https://trends.builtwith.com/widgets/captcha/traffic/Entire-Internet>.
- [17] 2023. Energy Information Administration FAQ. <https://www.eia.gov/tools/faqs/faq.php?id=74&t=11>
- [18] 2023. Google reCAPTCHA admin dashboard. <https://www.google.com/u/2/recaptcha/admin/site>
- [19] 2023. hCaptcha. <https://www.hcaptcha.com/>.
- [20] 2023. The post study survey via google forms (interactive version) DO NOT SUBMIT PERSONAL INFO. https://docs.google.com/forms/d/e/1FAIpQLSejdiHyw2z3YpjxTYMXeOTrn6ZC8Az6ockPm4b9lbhsvQ77gg/viewform?usp=sf_link
- [21] 2023. reCAPTCHA. <https://www.google.com/recaptcha/about/>.
- [22] 2023. reCAPTCHA v2. <https://developers.google.com/recaptcha/docs/display>.
- [23] 2023. SciPy is an open-source software for mathematics, science, and engineering. <https://scipy.org/>
- [24] 2023. USA Environmental Protection Agency Greenhouse Gas Calculator. <https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator>
- [25] Ismail Akrou, Amal Feriani, and Mohamed Akrou. 2019. Hacking google recaptcha v3 using reinforcement learning. *arXiv preprint arXiv:1903.01003* (2019).

- [26] Fatmah H. Alqahtani and Fawaz A. Alsulaiman. 2020. Is image-based CAPTCHA secure against attacks based on machine learning? An experimental study. *Computers & Security* 88 (2020), 101635. <https://doi.org/10.1016/j.cose.2019.101635>
- [27] Joshua Aslan, Kieren Mayers, Jonathan G. Koomey, and Chris France. 2018. Electricity Intensity of Internet Data Transmission: Untangling the Estimates. *Journal of Industrial Ecology* 22, 4 (2018), 785–798. <https://doi.org/10.1111/jiec.12630> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1111/jiec.12630>
- [28] Aaron Bangor, Philip Kortum, and James Miller. 2009. Determining what individual SUS scores mean: Adding an adjective rating scale. *Journal of usability studies* 4, 3 (2009), 114–123.
- [29] J. P. Bigham and A.C. Cavender. 2009. Evaluating Existing Audio CAPTCHAs and an Interface Optimized for Non-Visual Use. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Boston, MA, USA) (CHI '09). ACM, New York, NY, USA, 1829–1838. <https://doi.org/10.1145/1518701.1518983>
- [30] Kevin Bock, Daven Patel, George Hughey, and Dave Levin. 2017. unCaptcha: A Low-Resource Defeat of reCaptcha's Audio Challenge. In *11th USENIX Workshop on Offensive Technologies (WOOT 17)*. USENIX Association, Vancouver, BC.
- [31] John Brooke et al. 1996. SUS-A quick and dirty usability scale. *Usability evaluation in industry* 189, 194 (1996), 4–7.
- [32] Elie Bursztein, Steven Bethard, Celine Fabry, J. C. Mitchell, and Dan Jurafsky. 2010. How Good Are Humans at Solving CAPTCHAs? A Large Scale Evaluation. In *2010 IEEE Symposium on Security and Privacy*. 399–413. <https://doi.org/10.1109/SP.2010.31>
- [33] Jun Chen, Xiangyang Luo, Yanqing Guo, Yi Zhang, and Daofu Gong. 2017. A Survey on Breaking Technique of Text-Based CAPTCHA. *Security and Communication Networks* (12 2017). <https://doi.org/10.1155/2017/6898617>
- [34] Monica Chew and J. D. Tygar. 2004. Image Recognition CAPTCHAs. In *Information Security*, Kan Zhang and Yuliang Zheng (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 268–279.
- [35] Yunhe Feng, Qing Cao, Hairong Qi, and Scott Ruoti. 2020. SenCAPTCHA: A Mobile-First CAPTCHA Using Orientation Sensors. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 2, Article 43 (jun 2020), 26 pages. <https://doi.org/10.1145/3397312>
- [36] Yunhe Feng, Qing Cao, Hairong Qi, and Scott Ruoti. 2020. SenCAPTCHA: A mobile-first CAPTCHA using orientation sensors. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 2 (2020), 1–26.
- [37] Christos A. Fidas, Artemios G. Voyiatzis, and Nikolaos M. Avouris. 2011. On the Necessity of User-Friendly CAPTCHA. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Vancouver, BC, Canada) (CHI '11). ACM, New York, NY, USA, 2623–2626. <https://doi.org/10.1145/1978942.1979325>
- [38] Christoph Fritsch, Michael Netter, Andreas Reisser, and Günther Pernul. 2010. Attacking Image Recognition Captchas. In *Trust, Privacy and Security in Digital Business*, Sokratis Katsikas, Javier Lopez, and Miguel Soriano (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 13–25.
- [39] Haichang Gao, Wei Wang, and Ye Fan. 2012. Divide and conquer: an efficient attack on Yahoo! CAPTCHA. In *2012 IEEE 11th International Conference on Trust, Security and Privacy in Computing and Communications*. IEEE, 9–16.
- [40] Haichang Gao, Jeff Yan, Fang Cao, Zhengyua Zhang, Lei Lei, Mengyun Tang, Ping Zhang, Xin Zhou, Xuqin Wang, and Jiawei Li. 2016. A Simple Generic Attack on Text Captchas. In *Network and Distributed System Security Symposium (NDSS)*. San Diego, California, United States.
- [41] H. Gao, D. Yao, H. Liu, X. Liu, and L. Wang. 2010. A Novel Image Based CAPTCHA Using Jigsaw Puzzle. In *2010 13th IEEE International Conference on Computational Science and Engineering*. 351–356. <https://doi.org/10.1109/CSE.2010.53>
- [42] Song Gao, Manar Mohamed, Nitesh Saxena, and Chengcui Zhang. 2019. Emerging-Image Motion CAPTCHAs: Vulnerabilities of Existing Designs, and Countermeasures. *IEEE Transactions on Dependable and Secure Computing* 16, 6 (2019), 1040–1053. <https://doi.org/10.1109/TDSC.2017.2719031>
- [43] Ian J Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. 2014. Multi-digit number recognition from street view imagery using deep convolutional neural networks. *arXiv preprint arXiv:1312.6082* (2014).
- [44] Meriem Guerar, Luca Verderame, Mauro Migliardi, Francesco Palmieri, and Alessio Merlo. 2021. Gotta CAPTCHA 'Em All: A Survey of Twenty years of the Human-or-Computer Dilemma. *CoRR* abs/2103.01748 (2021). arXiv:2103.01748 <https://arxiv.org/abs/2103.01748>
- [45] Carlos Javier Hernandez-Castro and Arturo Ribagorda. 2010. Pitfalls in CAPTCHA design and implementation: The Math CAPTCHA, a case study. *Computers & Security* 29, 1 (2010), 141–157.
- [46] Chien-Ju Ho, Chen-Chi Wu, Kuan-Ta Chen, and Chin-Laung Lei. 2011. DevilTyper: A Game for CAPTCHA Usability Evaluation. *Comput. Entertain.* 9, 1, Article 3 (apr 2011), 14 pages. <https://doi.org/10.1145/1953005.1953008>
- [47] Egor Homakov. 2014. The No CAPTCHA problem. <https://homakov.blogspot.com/2014/12/the-no-captcha-problem.html>
- [48] Md Imran Hossen and Xiali Hei. 2021. A Low-Cost Attack against the hCaptcha System. *CoRR* abs/2104.04683 (2021). arXiv:2104.04683 <https://arxiv.org/abs/2104.04683>
- [49] Md Imran Hossen, Yazhou Tu, Md Fazle Rabby, Md Nazmul Islam, Hui Cao, and Xiali Hei. 2020. An Object Detection based Solver for Google's Image reCAPTCHA v2. In *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*. USENIX Association, San Sebastian, 269–284.
- [50] Mohit Jain, Rohun Tripathi, Ishita Bhansali, and Pratyush Kumar. 2019. Automatic Generation and Evaluation of Usable and Secure Audio reCAPTCHA. In *The 21st International ACM SIGACCESS Conference on Computers and Accessibility* (Pittsburgh, PA, USA) (ASSETS '19). Association for Computing Machinery, New York, NY, USA, 355–366. <https://doi.org/10.1145/3308561.3353777>
- [51] Ayca Kaya, Reha Ozturk, and Cigdem Altin Gumussoy. 2019. Usability measurement of mobile applications with system usability scale (SUS). In *Industrial Engineering in the Big Data Era: Selected Papers from the Global Joint Conference on Industrial Engineering and Its Application Areas, GJCIE 2018, June 21–22, 2018, Nevsehir, Turkey*. Springer, 389–400.
- [52] Brandy Klug. 2017. An overview of the system usability scale in library website and system usability testing. *Weave: Journal of Library User Experience* 1, 6 (2017).
- [53] Philip T Kortum and Aaron Bangor. 2013. Usability ratings for everyday products measured with the system usability scale. *International Journal of Human-Computer Interaction* 29, 2 (2013), 67–76.
- [54] Kat Krol, Simon Parkin, and M. Angela Sasse. 2016. Better the Devil You Know: A User Study of Two CAPTCHAs and a Possible Replacement Technology. In *2016 NDSS Workshop on Usable Security*. 1–10. <https://doi.org/10.14722/usec.2016.230013>
- [55] Chunhui Li, Xingshu Chen, Haizhou Wang, Peiming Wang, Yu Zhang, and Wenxian Wang. 2021. End-to-end attack on text-based CAPTCHAs based on cycle-consistent generative adversarial network. *Neurocomputing* 433 (2021), 223–236. <https://doi.org/10.1016/j.neucom.2020.11.057>
- [56] Jun Liang, Deqiang Xian, Xingyu Liu, Jing Fu, Xingting Zhang, Bizhou Tang, Jianbo Lei, et al. 2018. Usability study of mainstream wearable fitness devices: feature analysis and system usability scale evaluation. *JMIR mHealth and uHealth* 6, 11 (2018), e11066.
- [57] David Lorenzi, Jaideep Vaidya, Emre Uzun, Shamik Sural, and Vijayalakshmi Atluri. 2012. Attacking Image Based CAPTCHAs Using Image Recognition Techniques. In *Information Systems Security*, Venkat Venkatakrishnan and Diganta Goswami (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 327–342.
- [58] Klaus M Miller and Bernd Skiera. 2023. Economic Consequences of Online Tracking Restrictions: Evidence from Cookies. *International Journal of Research in Marketing* (2023).
- [59] Manar Mohamed, Song Gao, Nitesh Saxena, and Chengcui Zhang. 2014. Dynamic Cognitive Game CAPTCHA Usability and Detection of Streaming-Based Farming. In *2014 NDSS Workshop on Usable Security*. 1–10. <https://doi.org/10.14722/usec.2014.230021>
- [60] Debajyoti Pal and Vajirasak Vanijja. 2020. Perceived usability evaluation of Microsoft Teams as an online learning platform during COVID-19 using system usability scale and technology acceptance model in India. *Children and youth services review* 119 (2020), 105535.
- [61] Sarah Perez. 2012. Google now using recaptcha to decode street view addresses. <https://techcrunch.com/2012/03/29/google-now-using-recaptcha-to-decode-street-view-addresses/>
- [62] Steven A. Ross, J. Alex Halderman, and Adam Finkelstein. 2010. Sketcha: A Captcha Based on Line Drawings of 3D Models. In *Proceedings of the 19th International Conference on World Wide Web*. ACM, New York, NY, USA, 821–830. <https://doi.org/10.1145/1772690.1772774>
- [63] Andrew Searles, Yoshimichi Nakatsuka, Ercan Ozturk, Andrew Paverd, Gene Tsudik, and Ai Enkoji. 2023. An Empirical Study & Evaluation of Modern CAPTCHAs. In *32nd USENIX Security Symposium (USENIX Security 23)*. USENIX Association, Anaheim, CA, 3081–3097. <https://www.usenix.org/conference/usenixsecurity23/presentation/searles>
- [64] Vinay Shet. 2014. Street View and reCAPTCHA technology just got smarter. <https://security.googleblog.com/2014/04/street-view-and-recaptcha-technology.html>
- [65] Suphannee Sivakorn. 2016. I'm Not a Human: Breaking the Google reCAPTCHA.
- [66] Suphannee Sivakorn, Iasonas Polakis, and Angelos D. Keromytis. 2016. I am Robot: (Deep) Learning to Break Semantic Image CAPTCHAs. In *2016 IEEE European Symposium on Security and Privacy (EuroSP)*. 388–403. <https://doi.org/10.1109/EuroSP.2016.37>
- [67] Mengyun Tang, Haichang Gao, Yang Zhang, Yi Liu, Ping Zhang, and Ping Wang. 2018. Research on Deep Learning Techniques in Breaking Text-Based Captchas and Designing Image-Based Captcha. *IEEE Transactions on Information Forensics and Security* 13, 10 (2018), 2522–2537. <https://doi.org/10.1109/TIFS.2018.2821096>
- [68] Nitirat Tanthavech and Apichaya Nimkoumpai. 2019. CAPTCHA: Impact of Website Security on User Experience. *ICIIT '19: Proceedings of the 2019 4th International Conference on Intelligent Information Technology* (02 2019), 37–41. <https://doi.org/10.1145/3321454.3321459>
- [69] Erkam Uzun, Simon Chung, Irfan Essa, and Wenke Lee. 2018. rtCaptcha: A Real-Time Captcha Based Liveness Detection System. In *Network and Distributed System Security Symposium (NDSS)*. San Diego, California, United States. <https://doi.org/10.14722/ndss.2018.23253>
- [70] Prokopia Vlachogianni and Nikolaos Tselios. 2022. Perceived usability evaluation of educational technology using the System Usability Scale (SUS): A systematic

- review. *Journal of Research on Technology in Education* 54, 3 (2022), 392–409.
- [71] Luis von Ahn, Manuel Blum, Nicholas J. Hopper, and John Langford. 2003. CAPTCHA: Using Hard AI Problems for Security. In *Advances in Cryptology – EUROCRYPT 2003*, Eli Biham (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 294–311.
- [72] Luis von Ahn, Benjamin Maurer, Colin McMillen, David Abraham, and Manuel Blum. 2008. reCAPTCHA: Human-Based Character Recognition via Web Security Measures. *Science* 321, 5895 (2008), 1465–1468. <https://doi.org/10.1126/science.1160379> arXiv:<https://www.science.org/doi/pdf/10.1126/science.1160379>
- [73] Margaret A Webb and June P Tangney. 2022. Too Good to Be True: Bots and Bad Data From Mechanical Turk. *Perspectives on Psychological Science* (2022).
- [74] Haiqin Weng, Binbin Zhao, Shouling Ji, Jianhai Chen, Ting Wang, Qinning He, and Raheem Beyah. 2019. Towards understanding the security of modern image captchas and underground captcha-solving services. *Big Data Mining and Analytics* 2, 2 (2019), 118–144. <https://doi.org/10.26599/BDMA.2019.9020001>
- [75] Jeff Yan and Ahmad Salah El Ahmad. 2008. A Low-cost Attack on a Microsoft CAPTCHA. In *Proceedings of the 15th ACM conference on Computer and communications security*. 543–554.
- [76] Yang Zi, Haichang Gao, Zhouhang Cheng, and Yi Liu. 2020. An End-to-End Attack on Text CAPTCHAs. *IEEE Transactions on Information Forensics and Security* 15 (2020), 753–766. <https://doi.org/10.1109/TIFS.2019.2928622>

A WORKFLOW

In this appendix we show basic workflows for account creation and password recovery processes that participants followed in the user study.

A.1 Account Creation

Figures 18, 19, 20, 21 constitute the workflow of the account creation process.

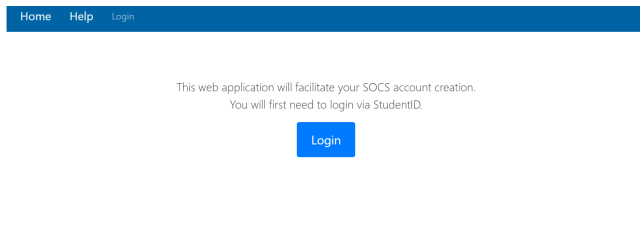


Figure 18: Initial login page

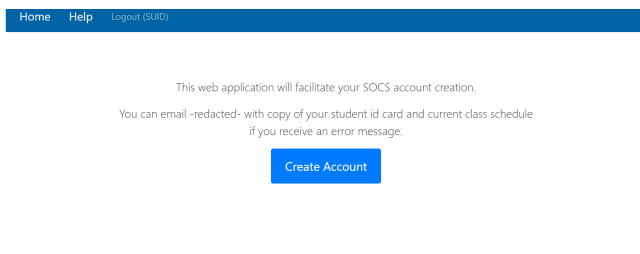


Figure 19: Initial Account Creation Page

A.2 Password Recovery

Figures 22 and 23 present the workflow of the password recovery process.

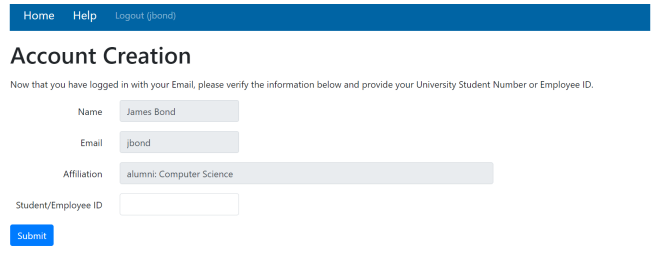


Figure 20: Account creation form

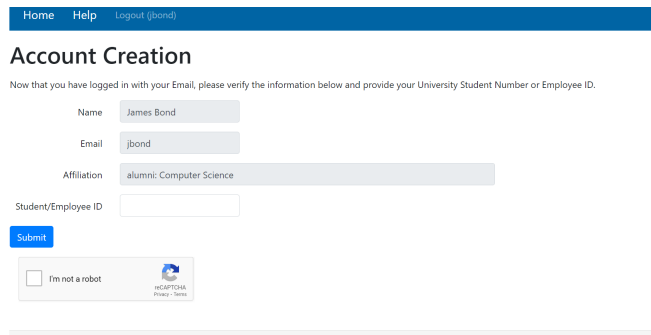


Figure 21: AC form after clicking submit

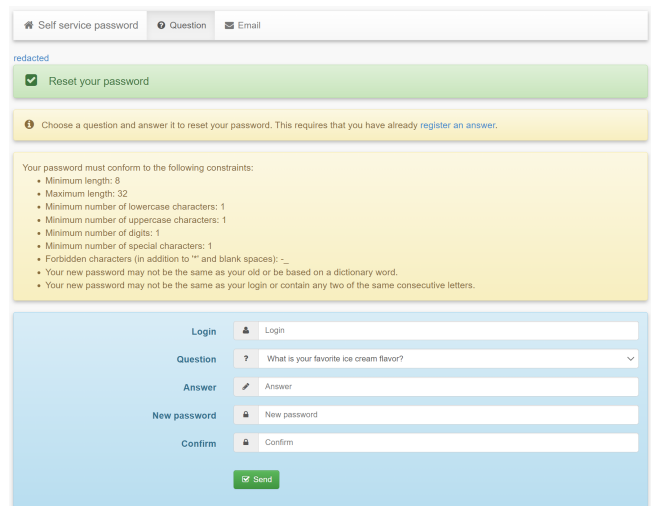


Figure 22: Password Recovery form

Figure 23: Password Recovery form after clicking submit

B NETWORK ANALYSIS OF RECAPTCHA

This Appendix contains a high-level technical analysis of reCAPTCHA. It has been considered in [66], which described the display method and workflow of reCAPTCHA with the emphasis on security aspects. Whereas, our goal is to: (1) determine various overhead factors incurred whenever a web-page uses reCAPTCHA, and (2) investigate reCAPTCHA’s automation detection capability. To this end, we performed black box program and network traffic analyses for common usage scenarios. We used two simple web pages for this purpose:

- Baseline page without any CAPTCHAS. This page is called *simple.html* and its source code is shown in Figure 24.

```
<html>
<head> <title> Simple Web Page </title> </head>
<body> <h4> A minimal web page </h4> <br /> </body>
</html>
```

Figure 24: Source code of simple.html

```
<html>
<head> <title> reCAPTCHA Difficult </title>
<script src="https://www.google.com/recaptcha/api.js"
  async defer></script>
</head>
<body>
<h4>A minimal web page</h4>
<div class="g-recaptcha"
  data-sitekey="obtained-site-key"></div> <br />
</body>
</html>
```

Figure 25: Source code of recaptcha.html

- A page similar to the baseline page, except with an additional reCAPTCHA. This page is called *recaptcha.html* and its source code is shown in Figure 25. As evident from the figure, integrating reCAPTCHA into a web page is very easy and straightforward.

These pages were visited using Google Chrome browser [9] and each usage scenario was repeated at least ten times. Browsing was performed in both guest and normal (profile logged-in) modes. Relevant information in the format of a *.har* file was collected for each scenario using Chrome DevTools.

The rest of this section describes the findings. Notations are summarized in Table 18.

Table 18: Notation Summary

Notation	Description
g1	https://www.google.com/recaptcha
g2	https://www.gstatic.com/recaptcha/releases/vkGiR-M4noX1963Xi_DB0JeI
g3	https://www.gstatic.com/recaptcha/api2/
g4	https://www.google.com/recaptcha/api2/
g5	https://fonts.gstatic.com/s/roboto/v18/
dv	different values

B.1 Page load Latency

Table 19 shows additional API calls made while loading *recaptcha.html* webpage.

Table 19: reCAPTCHA API Calls during page load

Request URL	Content-Length (B)
g1/api.js	554
g2/recaptcha_en.js	166822
g4/anchor?ar=[dv]	27864 (average)
g2/styles_ltr.css	24605
g2/recaptcha_en.js	166822
g3/logo_48.png	2228
g4/webworker.js?hl=[dv]	112
g2/recaptcha_en.js	166822
g4/bframe?hl=[dv] &v=[dv]&k=[dv]	1141-1145
g2/styles_ltr.css	24605
g2/recaptcha_en.js	166822
Network Overhead	254.01 KB-316.64KB

There are also 2-to-6 calls to g5 for downloading various web fonts. Content length for each of these calls is 15340, 15344, and 15552 bytes. Even though multiple calls are made to download *recaptcha_en.js* and *styles_ltr.css*, only the first call downloads the file, if necessary. These observations are taken into account when computing network overhead in Table 19.

Moreover, *api.js*, *recaptcha_en.js*, *styles_ltr.css*, *logo_48.png*, and web fonts are often served from the cache. Table 19 provides an upper bound on network overhead for page load. Average network overhead is computed by extracting actual network transmission during page load from collected *.har* files. Table 20 shows the results.

Table 20: recaptcha.html load network overhead

Scenario	Page Name	Page Size(KB)
First load	simple.html	0.631KB
First load	recaptcha.html	408.5KB
Network Overhead		407.869KB
Subsequent loads	simple.html	0.241 KB
Subsequent loads	recaptcha.html	29.56 KB
Network overhead		29.319 KB

We investigated load latency using Chrome DevTools, pingdom.com [10], and webpagetest.com [11]. Table 21 presents the results. Latency computed using Chrome DevTools is the highest since Chrome DevTools determines the load time of simple.html and recaptcha.html in the same network where the concerned web pages are hosted. Observation shows that load latency increases as the distance between the user and the hosted webpage decreases (in terms of hops).

Table 21: recaptcha.html load latency

Measurement Tool	Page Name	load Time
Chrome DevTools	simple.html	51.16ms
Chrome DevTools	recaptcha.html	425.81ms
Time Overhead		374.65ms, 732.31%
pingdom.com	simple.html	375ms
pingdom.com	recaptcha.html	796ms
Time Overhead		471ms, 125.6%
webpagetest.org	simple.html	814.22ms
Subsequent Loads	recaptcha.html	2074.78ms
Latency		1260.56ms, 154.82%

B.2 Checkbox Click Overhead

Table 22 shows additional API calls made after checkbox is clicked. In this scenario, image CAPTCHA is not served to the user.

Table 22: reCAPTCHA API Calls after checkbox click

Request URL	Content-Length (B)	
g4/reload?k=[dv]	23844.67 (average)	
g4/userverify?k=[dv]	580.56 (average)	
g3/refresh_2x.png	600	
g3/audio_2x.png	530	
g3/info_2x.png	665	
g5/[font].woff2	15552	
Network Overhead		24.43 KB-41.77KB

In some cases, only the first two calls are made. Even when other calls are made, files are normally served from the cache, so there is no network traffic. Files are downloaded only in the first-ever attempt to solve reCAPTCHA in a given client browser. Table 22 depicts upper and lower bounds for the network overhead.

B.3 reCAPTCHA Image load Overhead

Table 23 shows additional API calls made when checkbox is clicked and an image CAPTCHA is loaded. It also provides the upper bound and the lower bound of the network overhead due to these calls.

Table 23: reCAPTCHA API Calls for image load

Request URL	Content-Length (B)	
g4/reload?k=[dv]	24439.16667 (average)	
g3/refresh_2x.png	600	
g3/audio_2x.png	530	
g3/info_2x.png	665	
g4/payload?p=[dv]	39589.45455 (average)	
Network Overhead		64.03 KB-96.72KB

In some cases, two calls are made to g5 to download web fonts; content length is 15340 and 15552 bytes, respectively. Also, refresh_2x.png, audio_2x.png, info_2x.png, and web fonts are often served from the cache instead of being downloaded.

B.4 Image Solution Verification Overhead

Table 24 shows additional API calls made when an image CAPTCHA solution is verified. In case of a correct solution, only the third call from Table 24 requires network transmission and thus incurs network overhead. In case of a wrong solution, the last call from Table 24 is made, which requires network transmission and adds to network overhead. In both cases, other calls are usually served from the cache. In some instances, when a wrong solution occurs, only the third and fifth calls from Table 24 are made.

Table 24: reCAPTCHA API Calls for correct image solution

Case	Request URL	Content-Length (B)
Both	g3/refresh_2x.png	600
Both	g3/audio_2x.png	530
Both	g4/userverify?k=[dv]	595.88
Both	g3/info_2x.png	665
Wrong Solution	g4/payload?p=[dv]	40922.167 (average)
Correct Solution Network Overhead		0.6KB
Wrong Solution Network Overhead		41.58KB

B.5 reCAPTCHA Expiration Overhead

Table 25 shows additional API calls made after a reCAPTCHA solution expires. Only the first and seventh calls (g4/anchor and g4/bframe) require network transmission and are considered in network overhead. Other calls are served from the cache.

Summary: Results of evaluating network overhead for various reCAPTCHA usage scenarios are summarized in Table 26. As evident from these results, using reCAPTCHA incurs considerable network and timing overhead.

Table 25: reCAPTCHA API Calls for reCAPTCHA expiration

Request URL	Content-Length (B)
g4/anchor?ar=[dv]	27864 (average)
g2/styles__ltr.css	24605
g2/recaptcha__en.js	166822
g3/logo_48.png	2228
g4/webworker.js?hl=[dv]	112
g2/recaptcha__en.js	166822
g4/bframe?hl=[dv] &v=[dv]&k=[dv]	1141-1145
g2/styles__ltr.css	24605
g2/recaptcha__en.js	166822
Network Overhead	29KB

Table 26: Summary of reCAPTCHA Network Overhead

Scenario	Network Overhead(KB)
First time Page Load	408.5
Subsequent Page Loads	29.319
Checkbox Click	24.43-41.77
Image Load	64.03-96.72
Image Correct Solution Verification	0.6
Image Wrong Solution Verification & New Image load	41.58
Solution Expiration	29

B.6 Automation Detection

Finally, we briefly looked into automation detection capability of reCAPTCHA. Specifically, checkbox click is performed through Jit-bit mouse macro recorder [12] and playwright automated headless Chrome browser [13]. Interestingly, the use of the mouse macro is not considered as suspicious bot activity by reCAPTCHA. When checkbox is clicked and the page is reloaded in quick succession, an image CAPTCHA is served on around 14 tries, regardless of whether the tasks were performed manually or via the mouse macro. However, performing the same tasks via Playwright Chrome browser is considered suspicious – an Image CAPTCHA is served upon the first request.