

Christopher Choy*, David Zagardo, Olivia Witanowska, and Tresna Mulatua

What the FLoC? Literature Review and Design Proposal for Federated Learning of Cohorts

Abstract: We investigate the Federated Learning of Cohorts (FLoC), Google’s current keystone for its plan to replace third-party cookies. Our investigations included in-depth literature reviews as well as conversations with experts on privacy economics, cryptography, and web tracking. We deliver a discussion on the ongoing debate surrounding FLoC from its announcement to its current trial phases. We ultimately find several faults with FLoC and view it as unsustainable in its current form. Finally, we propose a design for improving FLoC’s privacy guarantees while still providing good advertiser utility.

Keywords: federated learning, cohorts, targeted advertising, online tracking, differential privacy

DOI Editor to enter DOI

Received ..; revised ..; accepted ...

1 Introduction

As user demand for online privacy has increased, major platforms and service providers have begun to take a stand against indiscriminate web tracking. Google is one of the latest to join this cause by announcing their intent to phase out support for the third-party cookie. Still, Google is a major ad-tech company and remains reticent to remove online tracking and user profiling entirely. They proposed the Privacy Sandbox, an initiative to design, test, and adopt a series of tools that will preserve both privacy and the open web. Central to their new model of delivering targeted content and ads while preserving individual privacy are cohorts, k-anonymous identifiers in a user’s browser that represented said user’s interests. The Federated Learning of Cohorts

(FLoC) is their proposal for a privacy-preserving mechanism that calculates these user cohorts.

As a major player in the browser market, Google adopting FLoC will redefine online privacy for billions of users around the world. However, many critiques have cautioned against FLoC while some have outright denounced it as actually worsening privacy. Understanding both the potential benefits and harms of FLoC then becomes a necessary and important endeavor for anyone seeking to make a decision on its adoption and usage. This includes Google engineers and product managers, policymakers examining the state of online tracking, privacy advocates seeking to deliver nuanced critiques or recommendations, and users debating whether they want to opt-out. Notably, this endeavor goes beyond FLoC’s whitepaper. A robust understanding considers FLoC in the context of the greater Privacy Sandbox, the economic market forces at play, as well as how it compares to existing third-party cookie tracking and other alternatives.

The goal of our study is to provide an in-depth coverage of FLoC in context as well as some novel considerations. To that end, we conduct a literature review of existing critiques and specifications for FLoC. In addition, we synthesize this review with opinions and insights from 6 experts on industry-level privacy engineering, privacy and web economics, and cryptography. In this paper we discuss our answers and investigation methodology to several key questions that drove our research: Why did FLoC initially fail and will it fail again? Is behavioral advertising really a necessary condition to the open web as Google implies? And what alternatives are there for privacy-preserving targeted advertising? Finally, we offer several design proposals for improving FLoC by mitigating its privacy risks while continuing to provide good advertising accuracy.

Beyond here, our paper is divided into 4 sections. First, our background section which provides a robust review of necessary concepts like the current state of online targeted advertising. Second, our discussion section which details our findings from pursuing the answers to our key research questions. Third, our design proposals for improving FLoC. And finally, our conclusion with

*Corresponding Author: Christopher Choy: Affil, E-mail: cchoy2@andrew.cmu.edu

David Zagardo: Affil, E-mail: dzagardo@andrew.cmu.edu

Olivia Witanowska: Affil, E-mail: owitanow@andrew.cmu.edu

Tresna Mulatua: Affil, E-mail: tmulatua@andrew.cmu.edu

closing arguments, study limitations, and acknowledgements.

2 Background

Our work assumes a general understanding of three main subject bodies for which we provide a brief overview here. First, we detail the current landscape of online targeted advertising. This includes knowledge on third-party cookies, behavioral vs. contextual advertising, advertiser controls, and past studies on user preferences. Second, we summarize Google's Privacy Sandbox initiative—the parent project for which FLoC is a key portion. This includes the motivations for the initiative, other relevant technologies aside from FLoC, and the current project timeline. Finally, we explain the Federated Learning of Cohorts (FLoC). This includes summaries of relevant portions from its white paper and its troubled past.

2.1 Online Targeted Advertising

There exists a great deal of prior literature on web tracking and targeted (a.k.a. tailored) advertising. For our paper, we provide a brief background summary on some of the most relevant portions of that research space. Specifically, we explain how third-party cookies are used in tracking, how users feel about web tracking and targeted advertising, and how advertisers reach users today. In section 2.2, we detail how the Privacy Sandbox is impacting or could impact this space.

2.1.1 Third-Party Cookies

An HTTP cookie is a "globally unique pseudonymous device identifier" encoded into any stateful web technology [24]. HTTP is a stateless protocol, so cookies are typically embedded into browsers for the purposes of retaining that browser's state as it navigates through pages in a domain and switches between websites. This allows for functionality like staying logged in and maintaining a shopping cart while navigating through a website like Facebook or Amazon. Cookies for these purposes are typically embedded into a user's browser by the publisher of the website they are visiting. These cookies are called first-party cookies.

Third-party cookies are loaded into an end user's browser by a separate party than the website publisher. These third parties can take the form of ad-tech platforms seeking to learn your behavior across multiple pages rather than just one. Any website a user visits that has the third party server's code in it can query for these cookies from the user's browser and learn details of their visit. The most insidious part of third party cookies is the fact that this loading and reading of cookies into a user's browser is often done without said user's knowledge or consent.

The key trait about cookies for this paper is that they are intended to be uniquely identifying. Any website that knows what cookies to ask for can track and immediately know who an individual is. This is double edged sword. Being able to uniquely identify individuals across the web is good for authentication and web tailoring, but it is horrible for privacy.

2.1.2 User Preferences

Previous studies have shown that a majority of users are against online tracking advertising. In a 2009 US phone survey, 87% of respondents would not want advertising based on tracking [33]. In a 2010 study by McDonald and Cranor, 64% of survey participants noted that targeted advertising was invasive [25]. However, more recent studies may reveal changing trends or conflicting information. A 2015 study revealed that half of the participants agreed that website advertising is necessary to enjoy free services on the internet [17]. In addition, that study showed 37% of participants disliked receiving targeted advertisements, 23% liked it, and 40% were neutral about it.

Such changes in statistics suggest that either the methodologies in between these studies were sufficiently different or that changes in privacy preferences has changed over time. As entities like ad-tech companies aim to further normalize privacy violations, it becomes more important for works like this one to fully review technologies before they are rolled out .

2.1.3 Targeting ads to users

One of the most effective ways for a company to get their message across is through online advertising. Companies need to reach the right audience in order to make the most out of their advertising. Currently, companies can target specific users based on the traits, interests, and

preferences that they discover through tracking third-party cookies. Through the profiles that they build, advertisers can decide whether it would be worth showing a user a specific ad or not. Advertisers can use behavioral and/or contextual advertising to show relevant ads to users. Behavior advertising relies on monitoring user interests through their viewing behaviors, such as displaying a car ad to a user who had previously searched for cars online. In contextual advertising, advertisers will place their ads in the most relevant areas, like displaying their car ad on an automotive forum. By balancing both forms of advertising, marketers can effectively target users before and after they leave particular websites [5].

Re-targeting, also known as remarketing, is an online advertising method that allows companies to follow users all over the web using cookies. A small amount of code called a pixel is embedded into the website that puts an anonymous cookie onto the site's new visitors. When the visitor leaves, the cookie continues to follow them across the different sites that they visit. Cookies will let retargeting providers know what ad to serve a user on a website, ensuring that the company's particular ad is being shown to a user that previously visited their website [7]. This method of advertising allows marketers to effectively spend their advertising budget by focusing on people who have already shown interest in their particular brand.

When a user visits a website, an auction is held to decide what goes in the ad space that is available. Real-Time Bidding involves supply-side platforms (SSP) who offer advertising spaces and demand-side platforms (DSP) who are looking to buy space. When a user enters a website containing ad space, the ad exchange is notified of the available space. The DSPs then bid on the space for their advertisers and the highest offer is able to retrieve and display their winning impression. The entire auction and loading of the impression takes less than a second so users aren't even aware it's happening until the ad is displayed. Fig. 1 shows the process that happens behind the scenes with real-time bidding. While an impression does not guarantee the company that someone has clicked on their ad, real-time bidding does allow them to automate the strategic planning process to buy advertising space through target filtering [21] [7]. Advertisers not only get to benefit through automation, but statistically targeting users this way has resulted in high return on investment [8].

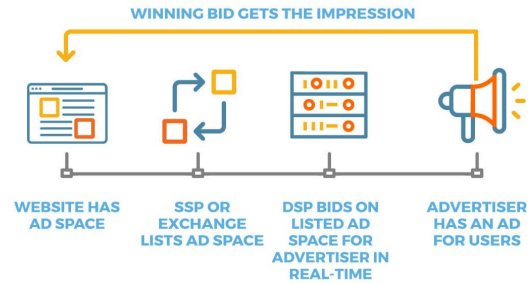


Fig. 1. Sketch of the ad bidding process, from CookiePro [9]. Websites offer ad space and advertisers bid to occupy that space.

2.2 The Privacy Sandbox

The Privacy Sandbox is Google's current initiative to build a more private web while still maintaining a healthy ad-ecosystem. In particular, they aim to phase out usage of the third-party cookie and replace it with their own fleet of technologies. Since FLoC does not and will not ever exist in a vacuum, it is imperative to understand how its cohort IDs will be used by the surrounding infrastructure to change the way online advertising is done today. In this section, we detail out why Google likely began this initiative in the first place, how the Privacy Sandbox will change the advertising methodology, and what mitigations have been proposed for improving individual user privacy protections.

2.2.1 Initial idea and motivations

The third-party cookie is great for the ad-ecosystem because it enables data brokers to track user activity across websites and build a comprehensive user profile that can then be sold to advertisers for targeted advertising. For Google in particular, there is a vested interest in maintaining a healthy ad-supported ecosystem. In 2020, Google Ads generated \$147 billion, 80% of Alphabet's revenue [22]. If Google fails to support advertisers, they can expect to lose a good portion of them to other "Walled Garden" platforms like Facebook. Advertisers function purely on economics; they want the best return-on-investment for their ads and targeted advertising is huge optimization for this.

However, Google is under pressure from competitors, legislature, and users to deliver on their promise to phase out support for the third-party cookie. "Users are demanding greater privacy—including transparency, choice, and control over how their data is used—and it's clear the web ecosystem needs to evolve to meet

these increasing demands," they wrote in a 2020 blog post [31]. As of June 2021, Google Chrome held 60% of the browser market share over their major competitors Safari (20%) and Firefox (3%) [22]. If they fail to meet user expectations, they risk losing their dominant position in the browser market.

2.2.2 Targeting and measuring ad performance

There are a couple of proposals working together alongside FLoC, such as FLEDGE. An extension of the proposal TURTLEDOVE, FLEDGE will allow for retargeting. Even after FLoC algorithms put a user into an interest cohort, advertisers may request browsers to have the user identify with a unique cohort which will allow more targeted advertisement when they visit a new site. Most importantly, these functions will be completed on the browser side unlike previously where it was done on the server side and involved trackers [14]. When requesting a user to be put into a particular cohort, the advertiser is not able to see the user, but rather is suggesting to the browser that the person should be associated with a different cohort based on what the advertiser was showing.

Another proposal is Attribution Reporting API (AR). AR is an API that is embedded inside the user's browser to measure two types of events that are linked to the user when visiting a website [27]. It works by automatically generating two types of reports that contain information about every user activity in the website. The first report is a report that contains information about the user activity when they click on some product; click-through report. The other report contains information about what information is viewed, it's called the view-through report.

Those two reports will be processed by the browser to create an event level report that consists of the summary of an individual user activity, such as click and view on a content. On another case, the browser will create another report called the aggregate report that works with the same principle as the event level report, but instead of specific particular user, the report gathers information from a group of users. All of these processes are happening inside the user's browser. Google claims that there is no PII included inside the report, and all of the browsing interest is kept by the browser locally on the browser. The browser then will send the encrypted report data to the advertiser.

2.2.3 Privacy hardening

Another proposal called User-Agent Client Hints enables servers to request information about the user's device without the need of parsing unstructured User-Agent Strings. By requiring the server to request specific information, UA-CH aims to reduce fingerprinting as the result of overly detailed UA strings [26]. The sandbox also contains a new API for Trust Tokens, a way for users to be authenticated across multiple sites without the use of passive tracking. Normally, users will receive tokens from various sites that they visit. However, in the new API a single token can be used to authenticate a user on multiple websites. Issuer websites can issue a trust token to a user's browser if they determine they aren't a bot through actions such as continuous account usage or completing transactions. The tokens can later be redeemed in other contexts to confirm a user's authenticity. Trust tokens are encrypted so websites are unable to identify individuals using them [13]. The API is beneficial for advertisers as it protects them from bots and fraud while also protecting the privacy of the users visiting.

2.3 Federated Learning of Cohorts

Google's proposed design for Federated Learning of Cohorts is of course crucial to understand for deriving and understanding both the critiques and improvement proposals we deliver later in the paper. Here we summarize relevant portions from the FLoC white paper [1] as well as the existing history behind its adoption timeline, in addition to providing simplified definitions of relevant terms and technologies.

Federated learning is a method with which one can train machine learning models without requiring a user's raw data to be sent to a centralized server. [2] With federated learning, developers are able to train their models on user devices. This approach allows the user's raw data to remain private, while the device (typically a cell phone) releases the training results to the centralized server. The models are then tested on the user's device using their private data, which allows the developer to gain an understanding of how accurate the model is. Through an iterative process of improvements, these models are updated and then made to be static.

Cohorts are "groups of users with similar interests." In the context of FLoC, Google generates cohorts of users to determine how it will distribute advertisements.

Locality sensitive hash functions are algorithms that have high probability to distribute similar elements to the same bins. This is contrary to the conventional understanding or usage of hash algorithms, in which one typically wishes for a hash algorithm to be indeterminate of an element's bits. Locality sensitive hashing is designed for collisions.[23]

Formally, a locally sensitive hashing algorithm takes an input v in the space of inputs \mathbb{V} and outputs a hash h within the space of all hashes \mathbb{H} . The space of \mathbb{H} is much smaller than the space of \mathbb{V} . Locally sensitive hashing algorithms have been applied to many common computer science problems like the Nearest Neighbor Search, Audio and Image Similarity, and Hierarchical Clustering. It follows naturally that the Chromium team has considered SimHash as their clustering algorithm of choice, in addition to the fact that it was developed and patented by Google.

2.3.1 How it works

The intuition is to “hide” users from individual tracking and targeting by assigning them into an interest-based cohort that is shared with some k other users at minimum; thus granting users k -anonymity. Ideally, the cohort identifier value is the only information shared with advertisers for the purposes of ad targeting. Additionally, the derivation of a user's cohort ID from their browsing history should occur locally on the user's device. Ideally, no raw browsing data will be sent to a centralized server.

The precise principles behind FLoC are as follows. Cohorts should prevent individual cross-site tracking and consist of users with similar browsing behavior. Computation of cohort assignments should be unsupervised with parameter values like minimum cohort size that can easily be explained. Computation should also be simple enough to accommodate low system requirements since it aims to leverage federated learning.

The white paper provides the following criteria for evaluating FLoC implementations: privacy, utility, and centralization. Privacy is defined as the fraction of users in large cohorts, utility is the measure of similarity between users within a cohort, and centralization assesses how much information needs to be sent to a central server to calculate cohort ID.

The paper discusses two possible implementations for cohort calculation. The first method is to feature extract a user's activity into a vector that is then hashed into p -bit vector using the SimHash algorithm. A useful

feature of this locally-sensitive hashing algorithm is that similar activity vectors will be given hash values close to each other. Cohorts in this case can be defined in two ways. The first definition is to declare all users with matching hashes (collisions) are in the same cohort. This is problematic because there's no way to guarantee a minimum cohort size, so any privacy guarantees would be limited to a probabilistic guarantee of collision. The second definition uses SortingLSH to sort the hashes and define cohorts by grouping the hashes into cohorts with minimum sizes. This allows users with differing, but similar, hashes to end up in the same cohort. The advantage here is that it grants the ability to fine-tune cohort sizings. The disadvantage to this approach is that it requires some centralization of the hashes to compute the cohorts.

The second method for cohort calculation is to use affinity hierarchical clustering with centroids. This method leverages user similarity graphs to inform cluster creation. Given the same vectors as the SimHash approach, this approach constructs a nearest neighbor graph with cosine similarity-weighted edges. Cohorts are then formed by grouping together nearby nodes until a desired cluster size is reached. The advantage to this approach is that cohorts are formed much more accurately with fine control over the cohort sizes. The disadvantage is that this approach currently requires centralized processing of user data (i.e. FL requirements are not met).

The paper found that compared to forming cohorts by random grouping, generating cohorts based on common user interest with FLoC could produce a 70% increase in accuracy while maintaining high levels of anonymity.

2.3.2 Deployment History

Federated Learning of Cohorts (FLoC) was initially unveiled by Google in January 2021 [11]. The initial timeline was for FLoC to be rolled out in origin trials in March [6]. However, after mass backlash from users, advertisers, EFF, and W3C, Google noted a “need to move at a responsible pace” and updated their timeline in a June 2021 announcement [10]. Now, Google plans to begin phasing out third-party cookie support in late-2022, taking the time to test and iterate on FLoC in the meantime [4].

Today, trial versions of FLoC are live in some Google Chrome browsers. The latest origin trials for it closed in mid-2021. Starting in Q1 2022, they will transition from the discussion phase of FLoC to the test-

ing phase. If nothing changes, we expect that the testing phase will enable FLoC by default in the Chrome browsers of select regions of users. During this time, third-party cookies will likely still be enabled and the privacy of individuals will be impacted.

3 Discussion

In the course of our investigations, we sought to answer several major questions. Why did FLoC initially fail and do we expect it to fail again? Is FLoC better than third-party cookies for privacy? Is the root idea behind the Privacy Sandbox, that 'behavioral advertising is necessary' flawed? What alternatives for privacy-preserving targeted advertising could Google pursue instead? The answers to these questions are non-obvious and required significant research, understanding, and discussion. Even then, the answers were not necessarily satisfying. To answer these questions, we looked to prior literature, open discussions with experts, and our own expertise. For our experts, we spoke with both professors and engineers in the field.

3.1 Criticisms of FLoC

FLoC has been widely criticized by advertisers, privacy advocates, and Google's competitors. We enumerate the best arguments from each in this section and combine them with our own explanations for why the flaws these criticisms point out are bad for FLoC and privacy.

3.1.1 Privacy advocate responses

The goal of FLoC is to avoid letting trackers access specific pieces of information that they can tie to specific people. However, FLoC does not exist in a vacuum and cohort IDs can be cross-correlated with other data and metadata collected by websites, advertisers, and Google to expose individuals to privacy harms.

One weakness of FLoC is that it makes browser fingerprinting of users easier. When users query websites, the website is not only collecting user's information within the cookie (first-party cookie) but also user's cohort ID which contains the user's browsing interest. While fingerprinting was previously possible with third-party cookies, cohort IDs function as an additional data

point to the user so that it makes user more unique even in between couple thousand group members.

Another FLoC weakness is that it can be reverse-engineered over time as websites and advertisers collect information through correlating cohort IDs with specific services. SSO services like Facebook, Twitter, and Google can allow these providers to learn what websites a user is accessing with what cohort ID. Websites that receive a large number of requests with the same cohort ID may also form conclusions about the particular interests for a given cohort ID. Over time, this ability to link cohorts with explicit interests will only become easier as data sets become more longitudinal. This means that websites and advertisers will actually be able to know more about a user on their first contact than they would otherwise. This runs contrary to the idea that users should have a right to present different aspects of their identity in different contexts. For example, if a user visits for medical information the health institution doesn't have a reason to know about user political affiliation.

Interest-based cohort IDs are also dangerous as they can leak sensitive information. FLoC uses an unsupervised algorithm to create its clusters. This means no one will have direct control over how people are grouped together. Since cohorts are interest-based it is likely that users will be grouped into cohorts that can reveal sensitive characteristics such as ethnicity, age, political orientation, sexual preferences, and income. This information may also be dangerously correlated with websites related to substance abuse, financial hardship or mental health support.

The research has found that FLoC is not necessarily a better solution to replace 3P cookie. Compared to 3P cookie FLoC adds an additional information to the advertiser, it adds more data point such as the information that contain in the Cohort ID. FLoC also broadcast the information to any website that user visits that previously not available in 3P cookie. Furthermore, FLoC could potentially give an obscure way on how advertiser will see the user and that could enable them to do targeted-advertising to people with sensitive characteristic.

3.1.2 Advertiser responses

FLoC gives more power to Google to control online advertising market. User information and activity on the websites will no longer can be accessible by the ad tech company directly using their 3P cookie. User's informa-

tion can only be process through the Cohort ID which only Google can decrypt the information inside of it.

Furthermore, for advertisers FLoC has shown clear limitations compared to third-party cookies. First, it allows for only interest-based ad targeting. This makes other targeting techniques, like ad sequencing and frequency capping, difficult once the third-party cookie is removed. Ad sequencing is a method for showing a specific user a sequence of ads in a particular order. This allows advertisers to do more cost-effective things like introducing their brand to a user with a long video ad and then reinforcing that message with short video ads. Frequency capping can also improve advertiser's budgets by limiting the number of times an ad is shown to a particular person within a timeframe[18].

Advertisers and regulators are also concerned about how useful FLoC is for reinforcing Google's dominant position in the ad market. As Google controls how cohort IDs are calculated, they may tweak the algorithm to their liking. They are also the only ones given a longitudinal view of users' cohorts over time as they continue to have access to raw browser data. FLoC and much of the privacy sandbox allows for Google to act as both referee and player in the advertising market.

3.1.3 Competitor responses

One of the main ideas why Google plan to eliminate 3P cookie by initiating privacy sandbox is because it can track user on multiple websites and read users' behavioral activity across websites, known as cross-site cookies. Behavioral advertising on the Web is primarily enabled by "cross-site cookies". By looking into the structure of FLoC, many browser technology companies like Firefox, Edge, and Safari are giving a sign that they will not implement Google's solution into their system.

For example, Firefox in their research report for FLoC mentioned, by combining a measurement of longitudinal evolution of Cohort IDs and pre-existed user's identifier such as IP address, geolocation and browser fingerprinting, FLoC is still significantly define user behavior on the internet[30]. Moreover, FLoC is magnifying user's ability to be tracked if someone with state-based tracking system track them because it enables trackers to gather up a profile of each user's activity on the internet and real life, thus having severely detrimental effects on user privacy. Firefox found that FLoC is adding more data point in user when they are browsing on the internet, and that is potentially could worsen

fingerprinting[30]. As a result, Firefox is blocking FLoC by default in their browser.

As for Microsoft, they believe the important thing in the process of replacing 3P cookie is to give user more choice on the transparent, control and privacy which they do not see on FLoC, thus they do not have any plan in the moment to implement FLoC. Also, FLoC is potentially leaked user information by giving signal in their ID-based group, so Microsoft blocked by default in the Edge browser. As a solution, Microsoft plan to implement their own proposal to replacing 3P cookie by implementing PARAKEET.

While on Safari, Apple believe that they will not implement FLoC in the near future but they still in the considerations for every proposal that available right now. At the moment, the browser has implemented its own privacy preserving feature called the Intelligent Tracking Prevention (ITP)[16]. This feature enables browser to block every type of 3P cookie without any exception. That means, there will be no cookie or website that will follow user while they are browsing on the internet.

DuckDuckGo took a more aggressive denouncement of FLoC by adding a Block FLoC feature to their browser extensions. The extension sends opt-out requests, setting page header, and overriding the Javascript API that calls for the cohort IDs. From our DuckDuckGo expert, "We're essentially deleting `Document.prototype.interestCohort`. As you hypothesized, in case the opt-out isn't honored, it'll still stop the data collection of the cohort ID."

3.2 The Necessity of Tracking

The mere fact that we are even discussing Google FLoC and offering design improvements has the dangerous power to perpetuate the idea that online behavioral advertising (OBA) is a necessary part of the internet experience. We have found through our review of the literature that there is no solid evidence to support this idea. In talking with numerous experts, none seem to support the idea that OBA is necessary for a thriving economy. In fact, there are new marketing firms, like Kobler in Norway for example, that purport to be able to provide similar (if not better) conversion metrics based solely on contextual advertising. We find ourselves in a prisoner's dilemma, where publishers and ad tech firms feel they must all abandon OBA at the same time, for the fear of being unable to keep up with conversion metrics provided by other advertisers.

A plethora of agencies, both private and public, have concluded that consumers are harmed by the competition in the online advertising space. A study published by the Competition and Markets Authority in the United Kingdom found that “competition is not working well in these markets, leading to substantial harm for consumers and society as a whole.” [3]

Major social media and technology companies have made the divisive decision to block research on political ads. [32] This dangerously Orwellian narrative echoes the CDC’s inability to conduct research on gun violence for 25 years as lobbied by the NRA.

It is pedantic to say that data collected from users is not always done so for the benefit of the user. In a landmark revelation earlier this year, it was brought to light that Facebook has chosen to act in opposition to the results of internal research. [28] It comes as no surprise that these studies confirmed the suspicions that many have held for the past decade: Facebook’s platform actively prioritizes content that keeps users engaged, as opposed to content that is naturally disseminated, beneficial, or healthy for society at large.

By pushing the notion that online behavioral advertising is a necessary part of the internet, our society will struggle to shift away from the framework of decisional interference and subversion of autonomy that has become the de facto standard in social media. Tech companies, by societal necessity, must strive to pivot from their existing business models, and recognize that they have a humanitarian duty to protect those that use their platforms.

3.3 Alternatives to FLoC

Google has not been the only company interested in addressing privacy issues related to online advertising. With the Privacy Sandbox, Google will be changing the environment by retaining user information on the browser rather than sharing with ad networks. However, others have suggested keeping it the same [29]. Other proposals such as PARAKEET have juggled with the idea of adding API to anonymize data coming from both the browser and the ad network, while others suggested regulation of data by a neutral institution like seen in GARUDA. While not necessarily solutions, these alternative proposals show the different perspectives on how we could handle the issue of user privacy.

3.3.1 Private Anonymized Requests for Ads that Keep Efficacy and Enhanced Transparency (PARAKEET)

Private Anonymized Requests for Ads that Keep Efficacy and Enhanced Transparency (PARAKEET) is Microsoft’s proposed alternative to third-party cookies. Working alongside Harvard professor Gary King, they proposed PARAKEET, a set of APIs that would ensure data is kept private using differential privacy while still allowing it to be used for research purposes by government and private sectors to gain novel insights [19]. They will also remove the need for cross-site identity tracking while still allowing for sufficient understanding of user interests for advertisers [35].

The browser, along with the browser-provided service, will be responsible for anonymization of ad requests containing user’s interests going to the ad network. It will anonymize information such as the context provided by the publisher, user’s geographic information and other client-specific details. By using differential privacy, they will enforce limits on the identifiability of user information contained in ad requests [29] [35]. The anonymized information within the ad request will then be used by the ad network to perform ad matching.

PARAKEET will support common methods used for ad matching such as: retargeting, lookalike targeting, in-market audience, and contextual targeting. For interest-based ad targeting, Microsoft created an API similar to TURTLEDOVE which will allow advertisers to add users to an ad interest group for a requested amount of time [29]. This information is then stored on the browser to be later used to compute the user’s ad interests with the help of the PARAKEET service. PARAKEET API will create anonymized differentially private ad interest vectors that will be stored on the browser for use in ad requests to the ad network. Users will also have control for both actual and anonymized ad interests and features [35].

Unlike in Google’s Privacy Sandbox, when an ad request is initiated the decision process is done on the ad network like it currently is. The user’s interest vector is applied to the request before further anonymization in the proxy server created by PARAKEET. Anonymized ad requests are sent to the ad network where a chosen ad is sent back through the proxy server to the publisher on the user’s browser, therefore not requiring any on-device processing [29]. This is unlike Google’s proposal, where the decision making is all located locally on the user’s browser. Additionally, ad networks respond with an ad bundle containing an encoded ad click URL. The URL

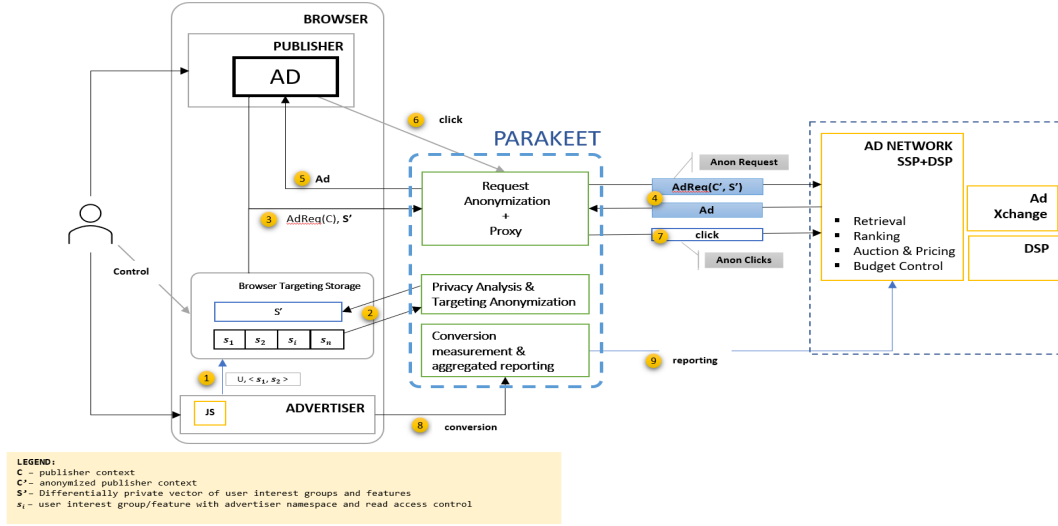


Fig. 2. Environment containing PARAKEYT API, as described by Microsoft [20].

will be rendered on the fenced frame on the publisher's site but won't allow any exchange of information with the page containing the frame to protect user privacy [35].

Microsoft hopes to challenge Google's Privacy Sandbox by keeping the ecosystem functioning the same way it currently does, with marketers still receiving user data, however the information is more anonymous. They believe that with their proposal they can improve end-user privacy while still allowing sites the ability to gain business through ad funding [35].

3.3.2 Governance of Ad Request by a Union of Diverse Actors (GARUDA)

The idea of Governance of Ad Request by a Union of Diverse Actor (GARUDA) is to have a centralized ad server that will integrate every stakeholders in internet-based advertisement industry. The centralized server works as gateway for each stakeholders so every data that goes through the server can be easier to maintain and audit. The model is suggesting for using a centralized server because it has 2 distinctive advantages compares to other such as, it is easier to build trust between stakeholders and it doesn't to change much to work properly. That means, the stakeholders like publishers and ad tech company will be easier to trust each other through the system[15].

For the proposal can be work properly the centralized server must be trusted by all stakeholders. There should not be just one party to operate the server, oth-

erwise it will be no trust between entities. The trust can be achieved by creating an independent board of governance from each party of the entities that bounded with a regulation or contract. This board then will supervise the institution which runs on daily basis by the legal entity. Legal entity is in charge for maintaining and deploying the trusted server and the network environment which GARUDA will run on, such as PARAKEYT.

3.4 Closing Discussion Remarks

Given all of this information, it becomes clear the lack of a clear answer to each of our research questions. What we formulate here then are our own opinion answers based on the collaborative research and discussions we had throughout this semester.

In it's current state, we do not believe that FLoC will ever fully replace third-party cookies despite Google's current timeline insisting that it will. The landscape has shown that advertisers, privacy advocates, and competitors have all found fault and refused to adopt it. If advertisers flee to other walled gardens like Facebook or a majority of their users opt-out and deploy DuckDuckGo's FLoC blocker, Google stands to lose an incredible amount of ad revenue. This being the case, it is likely that Google may only fully adopt FLoC while third-party cookies are still enabled in their browser. In this case, the presence of cohort IDs decidedly worsens user privacy as it enables easier browser fingerprinting. Another possible scenario is that FLoC may be adopted

only to be quickly replaced with an alternative solution similar to PARAKEET or GARUDA.

4 Design Proposals

In addition to the above research questions, we also had one last major research question, "if FLoC is truly terrible, can it be salvaged or improved enough to be viable?" To that end, we explored applying concepts from our lectures this semester on information security and differential privacy to try and improve FLoC.

FLoC is imperfect in its current state. Seeking perfection is dangerous, but revealing sensitive information about at-risk groups is an unacceptable flaw that must be amended if FLoC is to be safe, effective, and an acceptable contribution to the digital marketing environment.

We detail here several different design proposals for improving the current model of FLoC. Our first proposal is to apply differential privacy to the users' feature vectors. Our second proposal is to utilize hamming distance calculations to improve the accuracy of user cohorts. Our final proposal is to deceive and confound trackers through uncertainty reporting of cohort IDs.

For each of these proposals, we detail the weakness it seeks to overcome, the intuition and theory behind the proposal, and our own theoretical evaluation of the proposals strengths and weaknesses.

4.1 Differential Privacy

The idea of interest-based targeted advertising is overall unconcerned with users occasionally seeing ads that aren't actually relevant to them. In other words, occasional false positives and false negatives during ad auctions are acceptable. From a privacy enhancing perspective, this type of noise can actually be preferred. Thus, FLoC is a good candidate for applying differential privacy to because we seek good statistical utility at an aggregate-level and good privacy at an individual-level.

Differential privacy can be applied to anything. Whether or not it will be effective hinges on appropriately identifying what it is that one wishes to protect. FLoC is a difficult candidate for differential privacy. At the end of the day, in order for online behavioral advertising to be accurate or useful, a user's interests must be disclosed in some way, shape, or form. This points

to an inherent design flaw, and one that, on its own, differential privacy is unable to solve.

Ignoring the fact that research shows users of the internet do not want to be tracked, we pose a handful of questions below, and consider design proposals incorporating differential privacy that aim to provide appropriate solutions in subsequent sections. [12]

4.1.1 Differential Privacy

We must first frame the application of differential privacy through the lens of the following question:

What is it that we wish to protect?

Without identifying what we consider to be private, we are unable to successfully apply DP. Below are four specific scenarios where we detail what one might consider to be private.

1. *User Browsing History*: Here, we consider a user's browsing history feature *Vector* to be the information that we wish to keep private.

2. *User Interests*: Here, we consider user interests to be the aggregated browsing history subsets that map to specific interests.

3. *User Cohort*: Here, we consider a user's cohort (determined by Google) to be the private information we aim to protect.

4. *User Behavior*: Here, we consider user behavior to be reflected in the advertisements a user is shown.

4.1.2 Differentially Private Feature Vectors

We found a suitable way to enforce a privacy guarantee on the feature vector representing an individual's browsing history would be by adding Laplace noise to it. With this approach, the exact websites a user visits will be kept secret from the cohort calculation.

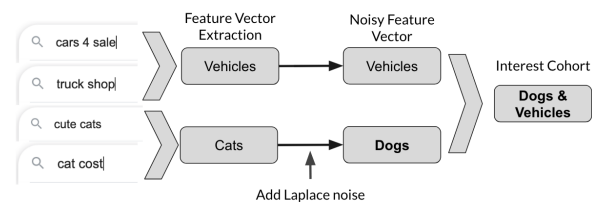


Fig. 3. A simplified diagram of the effect of noising a user's feature vector. Add Laplace noise to the feature vector to noise a user's interests

Notably, this model assumes a few things about our feature vector and there are a few different cases we considered. The first case is that our feature vector is a bit vector where each bit represents whether a user visited a particular web page or searched for a particular keyword. The length of this bit vector does not need to be the number of all websites on the internet as one of our experts was concerned about in this model. It can be restricted to the length of a subset of websites Google is interested in using for profiling (e.g. cnn.com, foxnews.com, nytimes.com, amazon.com). If our feature vector is a bit vector, we could actually apply Randomized Response instead of the Laplace mechanism.

The second case is that the feature vector is a sequence of score functions where each score function represents the number of times a user encountered or searched for a particular keyword (e.g. kitchenware) or keywords directly related to a keyword (e.g. spoon, pan). In this case, the problem begins to represent heavy hitter problems. If we aim to define our cohort interests by selecting the highest scoring keyword or set of high scoring keywords, then the Exponential Mechanism and Report Noisy Max algorithm can provide suitable means of obtaining a noisy answer for this.

The drawbacks to this approach is that, despite it keeping the precise websites a user visits private, it does not greatly perturb the resulting cohort ID. Since cohorts are formed using locally-sensitive hashing, similar interests result in "nearby" cohort IDs. We may obfuscate a user's true cohort ID and interests, but they will still land in the neighborhood of it. Therefore, adversaries will still learn a person's actual general interests even if that's just, for example, "cats" and not "cats riding bikes".

4.1.3 Privacy and Utility Guarantees

Applying noise at different stages of the aggregation process will incur different levels of privacy and utility trade off. Considering DP in production with Laplace or Gaussian Mechanisms, utility is maximized when the noise is applied over the entire data set. Below, we show the trade-off in utility that is incurred through a histogram function when applying noise locally and centrally. We see that we have an additive factor of m in the numerator. Heuristically this makes sense, given that we apply noise over every single bit within the user's vector within the local model.

Central Model:

$$b = \frac{GS_f}{\epsilon} = \frac{\Delta}{\epsilon} = \frac{2}{\epsilon}$$

Local Model:

$$b' = \frac{GS_f}{\epsilon'} = \frac{\Delta}{\epsilon'} = \frac{1}{(\epsilon/m)} = \frac{m}{\epsilon}$$

4.2 Clustering On a Distantly Federated Island's Secret Hamming Distance (CODFISH_d)

A principle design consideration within FLoC is Federated Learning. At every turn, Google proposes offloading more and more data to a centralized server, which fails at keeping true to the proposition's namesake. The principle considerations they note for doing so lie in two different factors:

1. On-device graph centroid clustering is computationally intensive for mobile devices [1]
2. The proposed mechanism for grouping required knowledge of every other user's feature vector [1]

In this section, we propose a method for determining a user's "nearest-neighbor" based on SimHash and Hamming Distance. The "nearest-neighbor" problem was described as the "post office" problem by Don Knuth in "The Art of Computer Programming." Within a city, a post office must determine within which ZipCode to place a residence. By computing the distance from the center of every ZipCode and returning a minimum value, a residence could be assigned their appropriate ZipCode. We incorporate this idea into our design proposal, viewing a user's ZipCode as an Island. A user's hash places them within the space of hashes, adrift at sea. Our algorithm assigns a user's hash to an Island within the sea of hashes. In future work, our algorithm could be refined by allowing each Island to include references to adjacent Islands, which would allow for a shorter running time complexity granted by the need to calculate fewer comparisons necessary to reach a user's Island.

Hamming Distance is defined as:

$$d(v, v') = \sum_{k=1}^p v_k \oplus v'_k$$

In our case, the function operates on two p -dimensional bit vectors of the same size. At each index k , v_k is checked against v'_k . If the bits are identical (e.g. 1 and 1), no change to the distance is made. If the bits

are different (e.g. 1 and 0), the sum is incremented by 1. Smaller Hamming Distances indicate greater levels of similarity. Larger Hamming Distances indicate smaller levels of similarity.

Our proposal was inspired by an online article showing how Hamming Distance can be used to increase the accuracy of SimHash in determining file similarity. [34] Rather than files, we will be considering bit vectors, which allows for a much cleaner implementation of the Hamming Distance and SimHash relationship.

Our proposal makes the following assumptions:

1. *The Space of Islands*: The underlying, critical assumption behind our design proposal lies in the idea that Google will create a Rainbow Table of equally spaced *SimHashes* from equally spaced *Vectors*. It should be noted that these *SimHashes* will not be a user's final cohort ID in our design proposal. Moving forward, we will refer to these "hash nodes" as a user's *Island* within the space of all hashes. To wit, a user's *Island* represents their approximate location within the sea of all hashes \mathbb{H} . Mathematically, we represent this in the manner shown below:

$$\mathbb{I} = (I_1, I_2, I_3, \dots, I_{9,998}, I_{9,999}, I_{10,000})$$

With $H(i_1) = I_1$, $H(i_2) = I_2$, $H(i_3) = I_3$, where $H(i)$ is the hash of an *Island* vector i , where $d(i_1, i_2) = d(i_2, i_3) = d(i_{p-1}, i_p)$.

There are approximately 2,000,000,000 users of Google Chrome. Assuming 10,000 equally spaced *Islands*, this will result in the average *Island* having a size of 200,000 users. These 10,000 bit vectors and their corresponding hashes will be stored on a user's device. Considering storage implications, we recommend the size of each bit vector to be approximately the size of a cookie; 4,093 bytes, or roughly 4,000,000 bits. Storing 10,000 bit vectors of such length would incur storage costs of approximately 40mb. We make the assumption that this storage overhead is acceptable and will incur negligible performance loss in user experience.

2. *The Vector*: A user's browsing history will be a sparsely populated bit vector v of size p , where p is the space of the 4,000,000 most popular websites. A bit of 0 corresponds to no visit, and a bit of 1 corresponds to any number of visits greater than or equal to 1. The current estimate of the size of the internet is approximately 1,700 million websites. We move forward on the assumption that a bit vector of dimension 4 million is enough to capture novel differences in browsing history

across users. This bit vector is generated and kept up to date on the user's device. It should be noted that there is no way to determine how Google is presently generating a user's bit vector, and application of our design proposal may or may not be possible based on Google's final implementation of FLoC.

3. *The Hash*: The user's hash is locally calculated from the *Vector* using the *SimHash* algorithm.

4. *Federated Learning*: All private data (including the user's *Hash*, *Vector*, and pure *Hamming Distance*) will never leave the user's device. We will consider a user's *Island* to be non-private information for the sake of our design proposal.

5. *Sensitive Web Activity*: We make the assumption that Google will be black-listing specific search queries or browsing activity associated with sensitive information. As a result, this data will not be used to calculate a user's feature vector, nor will it be used to determine the space of all hashes. We know that this is not a fool-proof or even scientifically-grounded idea. Fairness through unawareness fails, of this we are certain. We may choose to remove from our feature vector calculation websites containing keywords relating to health, race, gender, or other sensitive attributes, but this does not preclude inherent biases within user behavior that have been cemented in cultural norms through centuries of social inequity. To put it plainly; even without the use of sensitive browsing history, it may still be possible to determine a user's sensitive data.

6. *Future Work*: We have considered it to be outside the scope of this semester-long project to apply these methods, and move forward under the assumption that run-time analysis of our proposal would be conducted in future work.

The following steps outline our algorithmic implementation of the above assumptions, tools, and proposed methodology.

1. *Calculate Hamming Distance*: The first step is to calculate the Hamming Distance of a user's *Vector* from each and every locally stored *Island*.

2. *Apply Gaussian Mechanism to Hamming Distance*: Next, we apply Gaussian (or Laplacian) Noise locally to the set of all Hamming Distances, and clip each distance to a minimum value (discussed in detail below).

3. *Return Nearest Neighbors*: The final step in our implementation of the algorithm is to return to Google a differentially private approximation of the user’s Island within the set of all Islands. Specifically, we send the following:

$$\tilde{d}_k, \mathbb{I}_1, \text{ and } \mathbb{I}_2$$

Where \tilde{d}_k is a user’s differentially private Hamming Distance, and \mathbb{I}_1 and \mathbb{I}_2 are the user’s 2 closest *Islands*. By returning a differentially private approximation of a user’s space within the set of all Islands, we are able to guarantee two things:

1. *Utility*: Google will be able to determine a user’s approximate location within the set of all hashes. Knowing the user’s differentially private Hamming Distance from two specific nodes, Google can reconstruct a distribution of the user’s location throughout the space. From this information, they will be able to determine which users have which interests, and carve out space for each cohort in a manner that protects at-risk populations.

2. *Measurable Privacy*: While this does not prevent Google from determining a user’s interests, we are able to keep the entirety of the process in alignment with federated learning. The browsing history, hash, and bit vector never leave the user’s device.

Clipping the Hamming Distance: Below, we explain the privacy advantages of clipping to a minimum value.

To account for the scenario where a user’s *Vector* is nearly identical to one of the *ZipCode* vectors, we clip the Hamming Distance to a minimum value of:

$$d(i_k, i_{k-1}) = \sum_{k=1}^p i_k \oplus i_{k-1} = I_d$$

$$\tilde{d}_k = \max \left\{ d_k, \frac{I_d}{2} \right\}$$

Where I_d is the equal distance between each *Island*, d_k is a user’s differentially private Hamming Distance, and the value \tilde{d}_k is a clipped value of Hamming Distance.

By construct, \tilde{d}_k can never be less than $\frac{I_d}{2}$, as $d_k \in [0, \frac{I_d}{2})$ will always be clipped to $I_d/2$. When the noisy value of d_k is greater than $I_d/2$, the value of d_k will be returned as the user’s Hamming Distance.

The user is then able to calculate their nearest neighbors using their own *Vector* in conjunction with their Hamming Distance. This will guarantee that any

one user will *always* report at least 2 individual *Islands* as their nearest neighbors.

If we are able to translate the one-dimensional representation of a user’s *Vector* into a two-dimensional representation, we could guarantee that \tilde{d}_i would *always* return at least 4 individual *Islands* when $I_d/2$ (our clip value) is equal to half the length of the hypotenuse of the smallest coordinate square while providing similar levels of accuracy. We leave the implementation of this idea for future work.

4.3 Cohort Uncertainty Reporting Entropy

The aforementioned approaches grant some privacy to users while improving cohort accuracy. However, they do little to really mitigate FLoC’s fatal flaw: cohort IDs share a person’s interests publicly. So long as a user’s cohort ID is close to their true cohort ID, any web provider can learn their true general interests. Inspired by moving target defense and deception security protocols, we explored the idea of completely and temporally replacing a user’s true cohort ID with a fake cohort ID.

We propose Cohort Uncertainty Reporting Entropy (CURE) as a mechanism for managing and changing out a user’s cohort ID that operates entirely in the user’s local browser. CURE does not modify the FLoC process and allows for the calculation of a user’s cohort to proceed normally. Once the true cohort ID is calculated, CURE constructs an n length list of cohort IDs that may or may not contain the user’s true cohort ID with some probability $p < 1$. The selection of these false cohort IDs from the total space of cohort IDs should be independent, uniform, and random. CURE will then cycle through which cohort ID it reports, choosing to swap either on a temporal basis after time t or on an event basis after some event like closing the browser. When cohorts are recalculated, a new true cohort ID is calculated and a new set of cohort IDs is drawn.

An advantage to CURE is that it provides easily configurable privacy and utility guarantees for users and advertisers respectively. It can also be customized to individual browsers and websites in the event a user wants to select a certain privacy level or white-list targeted ads for a particular website. We can calculate the utility-privacy ratio, u , given by CURE with the formula,

$$u = \frac{p}{n}$$

The lower the value of u , the more privacy is granted to a user as u represents the expected amount of time a user’s true cohort ID will be broadcast for ad auctions.

A more advertiser utility focused configuration might set $p = .8$ and $n = 2$ to have a user see accurately targeted ads 40% of the time. A more privacy focused configuration might set $p = .5$ and $n = 5$ to have their true cohort shown only 10% of the time.

CURE greatly increases individual user privacy even at high utility-privacy ratios. This is because even at $n = 2, p = 1$ it grants a level of uncertainty and plausible deniability about what a user's true cohort is. An adversary will not be certain of a user's actual interests and thus cannot tie sensitive characteristics to them. In addition, CURE also effectively increases the population size of a user's cohort. At decently high levels of n , cohort IDs may fail to be a usable feature in browser fingerprinting completely.

CURE decreases advertiser utility, but unlike other privacy advocate approaches, it still allows for *some* behavioral advertising. In addition, the easily calculable and adjusted utility-privacy ratio means ad tech providers, publishers, and merchants can better estimate the increases in cost to them.

It is important that p , the probability of the true cohort being in the cohort set, never be 1. In a worse-case scenario, CURE will randomly select a set of cohorts all in the same neighborhood. It is important then that a user still have plausible deniability about those interests being truly associated with them. In another version of CURE, the set of cohorts could be constrained to be some sum total distance from each other. Figuring these particular implementation details out is best left to future work.

Another downside of CURE is that, like FLoC, it does not exist in a vacuum. When other privacy sandbox tools depend on a more consistent cohort ID being broadcast for advertiser techniques like re-targeting, CURE will greatly disrupt their capabilities.

5 Conclusion

This paper detailed an in-depth overview of Google's Federated Learning of Cohorts including its surrounding Privacy Sandbox and the critiques that caused it to initially fail. From there, we spoke with several experts in the form of privacy engineers and professors for their expertise on privacy economics, industry perspectives, and cryptography. Using those conversations as a spring board, we synthesized our own answers and discussions to our three favorite key research questions on FLoC's overall feasibility. Finally, we proposed three different

improvements for FLoC that could mitigate its privacy harms while still retaining good advertiser utility. We hope that this paper may serve as a comprehensive reference and be used to inspire future research and policy decisions.

5.1 Limitations and Future Work

One of our biggest challenges from the onset was the sheer scope of the space we wanted to dive into. We aimed to create a decently comprehensive review of all the critiques of Google FLoC, but we vastly underestimated the amount of relevant material. In addition to increasing our research load, it also made it difficult to conceive of important or novel things to do that hadn't already been done before. All of this was further bound by the simple limitation of lacking time, manpower, and experience.

We found throughout our proposal design process that not knowing the specifics of how FLoC is implemented by Google greatly limited our own designs. We had to make several generalizations and assumptions about the total space of cohort identifiers, how feature vectors were extracted from browser behavior, and the compute time for certain processes.

Early on in the project, we played around with the idea of mocking up a privacy sandbox ecosystem and modeling reconstruction attacks. However, we found we needed to make too many assumptions and simplifications about the implementation of these parts. Anything we would have produced would have likely be insignificant or a poor representation. If we had more details here, we would have loved to better explore this space and possibly do things like testing how robust DuckDuckGo's Javascript blocker is against changes to FLoC's signature.

Another limitation was that none of us were cryptographic experts. While we could apply portions of what we knew about differential privacy, the process was still complex enough that we were unable to validate our own answers or know of potentially better ways of applying it. That said, our proposal was still a good educational exercise in differential privacy considerations. Given more time, we would have liked to work through more formal write-ups and proofs that we could present to our cryptographic experts for feedback.

Given more time, we would also have continued interviewing more experts for their perspectives. We would have loved to speak with engineers on the actual

Google Privacy Sandbox team as well as workers at the EFF who loudly renounced FLoC from the beginning.

6 Acknowledgements

We would like to thank Professors Norman Sadeh and Ehab Al-Shaer for their project guidance, relevant lecture materials, and extraordinarily helpful office hours with us. We also thank our amazing experts Peter Doljanski, Tim Libert, Alessandro Acquisti, Steven Wu, and Aloni Cohen for their time and conversations with us. Without them, this project would have been much less interesting to conduct. Finally, we would like to thank Carnegie-Mellon University for providing us with project resources including Zoom subscriptions, group meeting rooms, and our professors.

References

- [1] Evaluation of cohort algorithms for the flocc api. URL <https://github.com/google/ads-privacy/blob/master/proposals/FLoC/FLoC-Whitepaper-Google.pdf>.
- [2] Federated learning. URL <https://federated.withgoogle.com>.
- [3] Online platforms and digital advertising market study. URL <https://www.gov.uk/cma-cases/online-platforms-and-digital-advertising-market-study>.
- [4] URL <https://privacysandbox.com/>.
- [5] Behavioral targeting vs. contextual targeting: Which is better? URL <https://blog.edificeautomotive.com/behavioral-targeting-vs.-contextual-targeting-which-is-better>.
- [6] URL <https://developer.chrome.com/blog/privacy-sandbox-update-2021-jan/>.
- [7] What is retargeting and how does it work? URL <https://retargeter.com/what-is-retargeting-and-how-does-it-work/>.
- [8] What is real-time bidding? how do i set up rtb?: Adjust. URL <https://www.adjust.com/glossary/real-time-bidding/>.
- [9] Real time bidding, Feb 2020. URL <https://www.cookiepro.com/wp-content/uploads/2020/02/real-time-bidding.jpg>.
- [10] Jun 2021. URL <https://blog.google/products/chrome/updated-timeline-privacy-sandbox-milestones/>.
- [11] Building a privacy-first future for web advertising, Jan 2021. URL <https://blog.google/products/ads-commerce/2021-01-privacy-sandbox/>.
- [12] A. Acquisti, L. Brandimarte, and G. Loewenstein. *Secrets and Likes: The Drive for Privacy and the Difficulty of Achieving It in the Digital Age*. Number ID 3688497. Sep 2020. URL <https://papers.ssrn.com/abstract=3688497>.
- [13] S. D. Advocate. Trust tokens, . URL <https://developer.chrome.com/docs/privacy-sandbox/trust-tokens/>.
- [14] S. D. Advocate. Fledge, . URL <https://developer.chrome.com/docs/privacy-sandbox/fledge/>.
- [15] R. Berjon. Governance of ad requests by a union of diverse actors (garuda), 01 2021. URL <https://darobin.github.io/garuda/>.
- [16] B. Budington. Apple's new webkit policy takes a hard line for user privacy, 08 2019. URL <https://www.eff.org/deeplinks/2019/08/apples-new-webkit-policy-takes-hard-line-user-privacy>.
- [17] F. Chanchary and S. Chiasson. User perceptions of sharing, advertising, and tracking. In *Eleventh Symposium On Usable Privacy and Security (SOUPS 2015)*, pages 53–67, Ottawa, July 2015. USENIX Association. ISBN 978-1-931971-249. URL <https://www.usenix.org/conference/soups2015/proceedings/presentation/chanchary>.
- [18] Google. About video ad sequencing, 2021. URL <https://support.google.com/google-ads/answer/9161595?hl=en>.
- [19] J. Kahan. Microsoft and harvard's institute for quantitative social science collaboration develops open data differential privacy platform, opens new research, Feb 2021. URL <https://www.linkedin.com/pulse/microsoft-harvards-institute-quantitative-social-science-john-kahan/>.
- [20] KeldaAnders. Parakeet, Feb 2021. URL <https://github.com/WICG/privacy-preserving-ads/blob/main/diagrams/Overallflow.png>.
- [21] W. Kenton. What is an impression?, May 2021. URL <https://www.investopedia.com/terms/i/impression.asp>.
- [22] S. Liu. Global market share held by the leading web browser versions as of june 2021, Jun 2020. URL <https://www.statista.com/statistics/268299/most-popular-internet-browsers/>.
- [23] G. S. Manku, A. Jain, and A. Das Sarma. Detecting near-duplicates for web crawling. In *Proceedings of the 16th international conference on World Wide Web - WWW '07*, page 141. ACM Press, 2007. ISBN 9781595936547. 10.1145/1242572.1242592. URL <http://portal.acm.org/citation.cfm?doid=1242572.1242592>.
- [24] J. R. Mayer and J. C. Mitchell. Third-party web tracking: Policy and technology. In *2012 IEEE Symposium on Security and Privacy*, pages 413–427, 2012. 10.1109/SP.2012.47. URL <https://ieeexplore.ieee.org/abstract/document/6234427>.
- [25] A. McDonald and L. Cranor. Beliefs and behaviors: Internet users' understanding of behavioral advertising. 08 2010.
- [26] R. Merewood. Improving user privacy and developer experience with user-agent client hints. URL <https://web.dev/user-agent-client-hints/>.

- [27] M. Nalpas, 08 2021. URL <https://developer.chrome.com/docs/privacy-sandbox/attribution-reporting-introduction/#use-cases-and-features>.
- [28] W. Oremus. Msn, 2020. URL <https://www.msn.com/en-us/news/technology/facebook-keeps-researching-its-own-harms-and-burying-the-findings/ar-AAOvQpK>.
- [29] Rahul and Surur. Microsoft prefers their parakeet to google's floc, Apr 2021. URL <https://mspoweruser.com/microsoft-prefers-their-parakeet-to-googles-floc/>.
- [30] E. Rescorla and M. Thomson. Technical comments on floc privacy, 06 2021. URL https://mozilla.github.io/ppa-docs/floc_report.pdf.
- [31] J. Schuh. Building a more private web: A path towards making third party cookies obsolete, Jan 2020. URL <https://blog.chromium.org/2020/01/building-more-private-web-path-towards.html>.
- [32] M. Swant. Top lawmakers and consumer advocates condemn facebook's decision to block academic research on political ads. URL <https://www.forbes.com/sites/martyswant/2021/08/04/top-lawmakers-and-consumer-advocates-condemn-facebooks-decision-to-block-academic-research-on-political-ads/>.
- [33] J. Turow, J. King, C. J. Hoofnagle, A. Bleakley, and M. Hennessy. Americans reject tailored advertising and three activities that enable it. Sep 2009. 10.2139/ssrn.1478214. URL <https://doi.org/10.2139/ssrn.1478214>.
- [34] B. Whitmore. Simhash and solving the hamming distance problem: explained, Aug 2019. URL <http://benwhitmore.altervista.org/simhash-and-solving-the-hamming-distance-problem-explained/>.
- [35] Wicg. `privacy-preserving-ads/parakeet.md` at main · wicg/privacy-preserving-ads, Oct 2021. URL <https://github.com/WICG/privacy-preserving-ads/blob/main/Parakeet.md#introduction>.