

# Gordon: Privacy Budget Management for W3C’s Privacy-Preserving Attribution API

## Abstract

Privacy-preserving advertising APIs like Privacy-Preserving Attribution (PPA) promise to enhance web privacy while enabling effective ad measurement. PPA replaces cross-site tracking with encrypted reports governed by differential privacy (DP), but current designs lack a principled approach to privacy budget management—creating uncertainty around critical design decisions. We present *Gordon*, a privacy budget manager for PPA that clarifies per-site budget semantics and introduces a global budgeting system grounded in resource isolation principles. *Gordon* enforces utility-preserving limits via quota budgets and improves global filter utilization through a novel batched scheduling algorithm. Together, these mechanisms establish a robust foundation for enforcing privacy protections in adversarial environments. We implement *Gordon* in Rust and Firefox and evaluate it on real-world ad data, demonstrating its resilience and effectiveness.

## 1 Introduction

Privacy-preserving advertising APIs, now under development and standardization in major browsers via the W3C, offer a rare opportunity to enhance online privacy while sustaining the web’s primary funding model. Historically, browsers have lacked structured support for ad-related tasks like *conversion attribution measurement*, which requires linking ads viewed on content sites to purchases made on seller sites—a cross-origin function fundamentally at odds with the same-origin principle that underpins browser design. This lack of support for the advertising workload has fueled widespread cross-site tracking through third-party cookies, fingerprinting, and other workarounds. The goal of the new APIs is to provide a structured, privacy-preserving alternative that aligns with browser principles while meeting advertising needs. However, these APIs remain in early stages, with technical challenges still unresolved—creating an opportunity for academic contribution.

Such collaborations have already had impact, underscoring that the space is ripe for foundational work. The *Cookie Monster* paper, presented at SOSP last year [8], introduced the first formal framework based on individual differential privacy (individual DP) [2] to systematically analyze and optimize these APIs—a framework later adopted by Google in privacy analysis of its ARA API [*ara-paper*]. That same Cookie Monster framework now underpins Privacy-Preserving Attribution (PPA) [*ppa*], the API standard being drafted by Private Advertising Technology Working Group (PATWG), a W3C working group that includes representatives from all browsers [*patwg*]. We are active participants in PATWG, tackling technical challenges from a scientific perspective to help advance the APIs’

practicality under strong privacy guarantees.

In this paper, we address a key open challenge: *privacy budget management in PPA*. PPA replaces cross-site tracking with a system where content sites register ads with the browser, seller sites request encrypted reports, and reports are only accessible via DP aggregation using secure multi-party computation or a trusted execution environment. Before sending an encrypted report, the browser deducts privacy loss from a *per-site privacy budget*, limiting how much new information a site can infer about a user. While PPA, through Cookie Monster’s algorithm, optimizes privacy loss accounting within each per-site budget using individual DP, it does not address how to manage these granular budgets to balance privacy with utility in an adversarial advertising ecosystem.

The absence of a principled approach to privacy budget management has led to unresolved questions within PATWG, creating uncertainty in key design decisions. For instance, should some sites get budget while others do not—and if so, based on what criteria?<sup>1</sup> Should there be a cap on how many sites are allocated budget, and if so, how can we prevent a denial-of-service attack where one entity exhausts it?<sup>2</sup> Should API invocations be rate-limited to prevent privacy or DoS attacks? To date, there is no consensus, largely due to the lack of a foundation to drive solutions.

We describe *Gordon*, a *privacy budget manager for PPA* that addresses semantic gaps in per-site privacy loss accounting and challenges introduced by the coarse-grained global budget PPA incorporates to protect users against adversaries controlling many sites. To clarify the semantics of PPA’s ambiguous per-site budgeting—often affected by shifting roles of third parties in the advertising ecosystem—we propose changes to the PPA interface, protocol, and terms of use, some of which have already been accepted by PPA.

For the global budget, the challenge is configuring and managing it to support benign workloads while resisting depletion by malicious actors. Our insight is to treat the global privacy budget as a *shared resource*—analogous to traditional computing resources but governed by privacy constraints—and to apply classic resource isolation techniques, such as quotas and batched scheduling [drf, 5], in this new domain. Beyond per-site and global budgets, *Gordon* introduces *quota budgets* that regulate global-budget consumption, ensuring graceful utility degradation for benign sites under attack. It does so by forcing adversaries to operate within expected workload bounds—which they can currently evade to wreak havoc on PPA’s global budget. Further, to address underutilization of

<sup>1</sup> Live discussion in W3C’s PAT community group, April 2024.

<sup>2</sup> <https://github.com/w3c/ppa/issues/69>, January 2025.

the global privacy budget caused by static quota partitioning, we propose a batched scheduling algorithm that reallocates unused capacity to otherwise-blocked requests. Together, these mechanisms establish a principled, practical foundation for PPA and give browsers a basis for enforceable defenses, along with guidance on where to focus.

We implement Gordon in two components: (1) `pdslib`, a generic on-device individual DP library that subsumes Cookie Monster and extends it with Gordon's budget management, and (2) integration into Mozilla Firefox's Private Attribution, a minimal PPA implementation. Upon release, these prototypes will serve as reference implementations for PPA, a service the PATWG has acknowledged as valuable.

We evaluate Gordon on a dataset from the Criteo ad-tech company. Results show that: **TODO(Pierre): (1) to occupy the paragraph for space budgeting, I add some lorem ipsum; (2) Lorem ipsum dolor sit amet, consectetur adipiscing elit; (3) Sed do eiusmod tempor incididunt ut labore et dolore magna aliqua; and (4) Ut enim ad minim veniam.**

## 2 PPA Overview and Gaps

### 2.1 PPA architecture

Fig. 1(a) illustrates the architecture of *Privacy-Preserving Attribution (PPA)*, W3C's browser-based API that enables *conversion attribution measurement* while preserving user privacy. Traditionally, browsers enforce a *same-origin policy*, while conversion attribution—the process of determining whether users who see an ad later make a purchase—is inherently *cross-origin*. It requires linking ad impressions shown on content sites (e.g., *news.ex*, *blog.ex*) to conversions occurring on advertiser sites (e.g., *shoes.ex*). In the absence of a structured API for this, advertisers rely on workarounds like third-party cookies, fingerprinting, and backend data exchanges—bypassing browser policies to accommodate workloads misaligned with current API structures.

PPA addresses this gap by enabling *cross-origin ad measurement* while preserving *single-origin privacy*, using differential privacy (DP) and secure aggregation via secure multi-party computation (MPC) or a trusted execution environment (TEE). This design bounds cross-origin information leakage, allowing the API to support effective ad measurement while upholding the intention of the browser's same-origin policy.

PPA defines four principals. **Impression sites** (*news.ex*, *blog.ex*) are content sites where ads are displayed. These sites register ad impressions with the browser using the function `saveImpression()`. **Conversion sites**, a.k.a. **advertiser sites** (*shoes.ex*), are sites where purchases or other conversions occur. When a user does a conversion, these sites invoke `measureConversion()` to link the event to any relevant prior ad impressions. **Intermediary sites** (*r1.ex*, *r2.ex*) are adtechs, typically embedded as frames in impression and conversion sites, that facilitate ad delivery and measurement. Unlike traditional tracking-based adtechs, they don't collect cross-site data directly but receive encrypted reports via a third

function, `getReport()`, which they then submit for secure aggregation. **Aggregation services** (e.g., *divviup.org*) are trusted MPC/TEE services that aggregate encrypted reports, applying DP to produce aggregated conversion metrics while ensuring no single entity can reconstruct individual user data.

### 2.2 Example workflow

Fig. 1(b) shows an example workflow for PPA, consisting of six steps (the same steps are also marked in the Fig. 1(a) architecture). The example entails an advertiser, *shoes.ex*, that launches an ad campaign to promote a new product. To compare the effectiveness of two ad creatives—a colorful ad highlighting the shoe's design and a black-and-white ad emphasizing materials and comfort—*shoes.ex* partners with two placement adtechs, *r1.ex* and *r2.ex*. Each adtech places the ads on content sites, e.g., *r1.ex* on *blog.ex* and *r2.ex* on *news.ex*. In addition to placing ads, these adtechs provide a *measurement service* that allows *shoes.ex* to compare the performance of its creatives within their respective networks.

① When a user visits *blog.ex*, *r1.ex* displays the colorful ad and registers the impression by calling `saveImpression()` with the parameters shown in the figure. ② Later, the user visits *news.ex*, where *r2.ex* displays the black-and-white ad and registers it by also calling `saveImpression()`. These impressions are stored *locally in the browser* within an *Impression Store*, along with important metadata, shown in Fig. 1(c).

③ Subsequently, if the user visits *shoes.ex* and purchases the shoes for \$60, the site invokes `measureConversion()` with the parameters shown in Fig. 1(b). This function searches the *Impression Store* in the browser for *relevant impressions*, matching the `impressionSite` and `conversionSite` metadata of the impressions to the parameters of `measureConversion()`. It then generates an *attributionObject*, which encapsulates the attribution histogram and manages privacy loss accounting. Assuming that PPA applies *uniform attribution*, it will assign the \$60 conversion value equally between the two registered impressions, assigning \$30 to each and resulting in the following attribution histogram: {1:30, 2:30}.

④ The *attributionObject* is *lazy*, i.e., no privacy loss occurs until it is used to request a report. To support DP queries, *shoes.ex* hands over the *attributionObject* to the *r1.ex* and *r2.ex* contexts within the browser, which invoke `attributionObject.getReport()`, specifying the aggregation service they intend to use (from a list of such services trusted by the browser). The browser processes these invocations by: (1) filtering the attribution histogram so that each intermediary only sees its own contributions (*r1.ex* gets {1:30}, *r2.ex* gets {2:30}); (2) encrypting the report and secret-sharing it (if MPC is used), while attaching some critical parameters as authenticated data, such as `epsilon` and `maxValue`; and (3) performing privacy loss accounting before sending the encrypted reports over to the intermediaries.

⑤ As more users purchase *shoes.ex*'s advertised product, additional encrypted reports are generated, each contain-

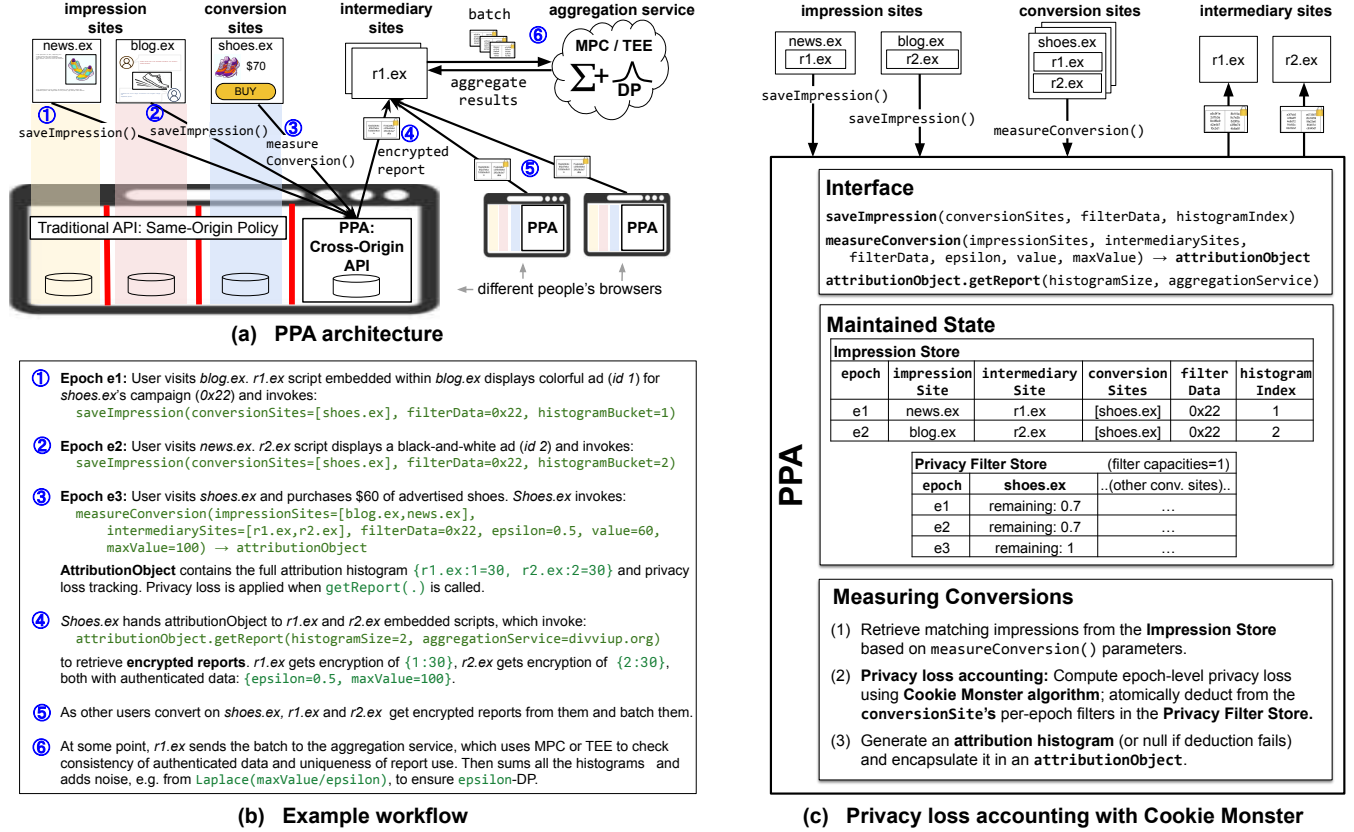


Fig. 1. PPA overview.

ing zero, one, or two attributed ads. ⑥ The intermediaries batch these reports and submit them to an aggregation service, which performs the final step: (1) validating the reports, ensuring all parameters in authenticated data match and that no report is reused; (2) summing the attribution values; and (3) applying DP, adding noise (such as from a Laplace distribution with scale  $\text{maxValue}/\epsilon$ ) to protect individual users. The resulting *noised, aggregated conversion metrics* are then provided to *r1.ex* and *r2.ex*, which relay the ad-effectiveness comparison back to *shoes.ex*, helping it discern which of the colorful vs. black-and-white ads leads to higher revenue.

### 2.3 Privacy loss accounting with Cookie Monster

PPA enforces privacy using the individual differential privacy (individual DP) framework from the Cookie Monster paper [cm-paper], which tracks each user's privacy loss separately and optimizes for on-device attribution. This is a key departure from traditional DP, which maintains a single global guarantee across users. individual DP allows PPA to bound privacy loss more efficiently—based only on the actual contribution of a device to a query.

Within each browser, PPA enforces individual DP at the *epoch* level, dividing the impression stream into time intervals (e.g., a week), each with its own privacy budgets. Each device maintains an *Impression Store* to log impressions per epoch and a *Privacy Filter Store* to track per-epoch budgets. A *privacy filter* acts as the epoch's budget manager: it deducts

privacy loss only if sufficient budget remains and only when data from that epoch contributes to a query; if depleted, it blocks further use of that epoch's data. Importantly, PPA maintains separate epoch-level privacy filters *per site*, a design choice that we show raises budget management questions.

Fig. 1(c) shows these internal components and illustrates how privacy loss is computed and enforced. When a conversion occurs (`measureConversion()`), the browser uses the Cookie Monster algorithm to: (1) retrieve all relevant impressions from the *Impression Store*, grouped by epoch; (2) compute *individual privacy loss per epoch*, using  $\text{value} / \text{maxValue} * \epsilon$  if an epoch has at least one relevant impression, or zero otherwise; and (3) attempt atomic deduction of this loss across all contributing epochs from the *conversion site's filters*, returning a null attribution if deduction fails, and the real one otherwise.

For example, in Fig. 1(b), epochs *e1* and *e2* each incur an individual privacy loss of  $\text{value} / \text{maxValue} * \epsilon = 0.3$ , while traditional DP would charge the full  $\epsilon = 0.5$  loss. Even better, epoch *e3*, which contains no relevant impressions, incurs *zero* individual privacy loss. Although this process is nominally part of `measureConversion()`, in practice it is deferred until `getReport()`: if no report is requested, no privacy loss is incurred. Fig. 1(c) shows the resulting filter state after *r1.ex* and *r2.ex* request their reports: assuming an initial filter capacity of 1, the conversion site *shoes.ex* re-

tains 0.7 budget in  $e_1$  and  $e_2$ , and the full 1 in  $e_3$ . In contrast, standard DP would leave only 0.5 in each epoch.

This example highlights how individual DP limits privacy loss based on actual contributions. To further understand this dynamic—which is significant for our own system’s design—we introduce a stock-and-flow analogy that captures the behavioral pattern that individual DP induces in PPA.

## 2.4 Stock-and-flow pattern

A key informal argument for PPA’s practicality, voiced in PATWG discussions, is that individual DP accounting naturally limits privacy consumption by tying it to *user actions on both impression and conversion sites*. Non-zero privacy loss arises only when both an impression (signifying a user visit to an impression site) and a conversion (a visit to a conversion site) are present. This induces a *stock-and-flow pattern*: *privacy stock* is created on impression sites as impressions are saved, and *privacy flow* is triggered on conversion sites when reports are requested over those impressions—*both gated by user actions*. PATWG discussions generally acknowledge that users who engage with more impression and conversion sites should incur more privacy loss—up to a limit, discussed next.

## 2.5 Global privacy filter

PPA acknowledges that relying solely on per-site filters risks exposing users to adversaries capable of coordinating API activity across multiple sites. Such behavior amplifies information gain from attribution, proportional to the number of sites involved. Since per-site filters impose no bound on this, PPA proposes “safety limits”—per-epoch global filters that span site boundaries—originally suggested by [6]. While the spec gives no detail on how to manage these filters—a gap this paper addresses (see next section)—our input has shaped the spec’s guiding principles: (1) global filter capacities must be much larger than per-site budgets, by necessity; and (2) these filters should “remain inactive during normal browsing and [trigger] only under high-intensity use or attack” [ppa].

## 2.6 Foundational gaps

We identify two key gaps in PPA related to managing its two filter types—per-site and global—which we address in Gordon.

**Gap 1: Unclear semantics of per-site filters.** PPA adopts Cookie Monster’s accounting model, which tracks privacy loss *per querier*, but is ambiguous about who counts as a querier in real-world deployments. For **single-advertiser queries**, PPA maps the querier to the conversion site—e.g., an intermediary requests a report on behalf of a specific advertiser like *shoes.ex*, and privacy loss is charged to that advertiser’s budget. Yet intermediaries also receive these reports and may reuse them for their own purposes, raising the question of whether they too should be considered queriers. The ambiguity grows with PPA’s planned support for **cross-advertiser queries**, where intermediaries aim to optimize across multiple advertisers (e.g., training models to choose the best ad for a given context). Since intermediaries directly

benefit from such queries, PATWG plans to charge privacy loss against their own budgets. This blurs the boundary between client-serving and self-serving queries, complicating the semantics of per-site accounting and increasing the risk of report misuse. In PATWG discussions, the global filter is often cited as a fallback, offering clearer semantics even when per-site guarantees break down—but this introduces its own configuration and management challenges, which we explore next. This paper proposes changes to the PPA API, protocol, and terms of use to clarify per-site semantics and the assumptions under which they remain sound (§4.1).

**Gap 2: Lack of mechanisms to manage the global filter.** A critical yet under-specified part of PPA is the configuration and management of the global filter, a shared resource across all parties requesting reports from a browser. This raises two key challenges: (1) how to set its capacity to support benign workloads and (2) how to prevent malicious actors from depleting it—either to boost their own utility or to deny service to others (e.g., competitors). While per-site budgets cap consumption per domain, they offer weak protection, as domain names are cheap and easily acquired. PATWG-discussed mitigations range from requiring sites to register with a trusted authority to browser-side heuristics for identifying illegitimate use of the API. But site registration faces resistance from some industry participants for undermining the API’s open nature while heuristics rely on notions of “legitimacy” that are hard to define, especially for a nascent API with no deployment history and potentially valuable, unforeseen use cases. For instance, should the number of invocations be limited? Over what period and to what value? Should access to device-side budgets be restricted? On what grounds? While discussion in PATWG continues, we argue that the group lacks a foundation—a minimal set of principled mechanisms with well-defined properties under clear assumptions—to guide browsers toward targeted, defense-in-depth strategies that are both protective and not over-constraining for the API. This paper contributes such a foundation, from the vantage point of PPA’s internal privacy budget management (§4.2).

## 3 Gordon Overview

We address PPA’s gaps by (1) clarifying the two distinct threat models that per-site and global guarantees address (§3.1) and (2) introducing Gordon to both restore the semantics of per-site filters and manage the global filter to support legitimate use while limiting abuse (§3.3). §3.2 introduces an example.

### 3.1 Threat model

PPA and Gordon share similar threat models. Users trust their OS, browser, and browser-supported aggregation services. They extend limited trust to first-party sites they visit intentionally—i.e., through *explicit actions* like direct navigations or clicks—granting them access to first-party data and cookies. Embedded intermediaries are not trusted at all, and no site—first-party or otherwise—is trusted with cross-site data.



As API designers, we must address two threat levels. The first is **intended use**, which assumes well-intentioned actors. Our goal here is to make compliance easy through careful API design and well-defined semantics. In the security literature, such actors are termed *honest-but-curious*: they follow the protocol but aim to extract as much information as permitted. Because some rules cannot be enforced by protocol alone, the API must include terms of use to close this gap. We define honest-but-curious adversaries as those who respect both the protocol and its terms of use.

PPA's *per-site filters* are meant to provide strong privacy guarantees against individual honest-but-curious sites. However, ambiguities in the current API and the lack of formal terms of use leave these guarantees semantically underspecified. Gordon addresses these gaps directly.

The second level involves **adversarial use**, where actors subvert the protocol and terms of use to extract excessive user information or maximize query utility through unauthorized budget consumption. Per-site budgets offer some protection when queriers operate independently or with limited coordination, leveraging DP's compositionality. But they fail under *large-scale Sybil attacks*, where an adversary registers many fake domains to bypass per-site caps. For example, a malicious conversion site X may use automatic redirection to cycle through Sybils, each triggering a single-advertiser report that maxes out its respective filter—multiplying the user's privacy loss by the number of Sybils.

PPA's *global filter* is designed to mitigate large-scale Sybil attacks by enforcing a coarser-grained budget. However, it introduces a new vulnerability: *denial-of-service (DoS) depletion attacks*. A malicious actor can deliberately exhaust the global budget, blocking legitimate queries—either to boost their own utility or harm competitors. These attacks can mirror the Sybil strategies used against per-site filters. §4.2 gives example attacks to which PPA is currently vulnerable.

Gordon embeds resilience directly into privacy budget management to defend against DoS depletion. While this layer alone does not provide complete end-to-end protection, it establishes a strong foundation and clarifies what browsers must enforce to achieve it. Building on the stock-and-flow model from §2.4, Gordon assumes that under intended use, privacy consumption is driven by explicit user actions—such as navigations or clicks—on distinct content and conversion first-party sites. As long as benign usage adheres to this pattern, Gordon fulfills PPA's guiding principles for the global filter (§2.5): supporting normal workloads under benign conditions and degrading gracefully under attack.

For this graceful degradation to hold in practice, two assumptions must be enforced: (1) browsers can reliably distinguish intentional user actions from automatic navigations, and (2) malicious actors cannot easily induce large numbers of users to intentionally visit many distinct attacker-controlled domains. If these assumptions fail, Gordon still upholds its privacy guarantees, but its DoS resilience will diminish.

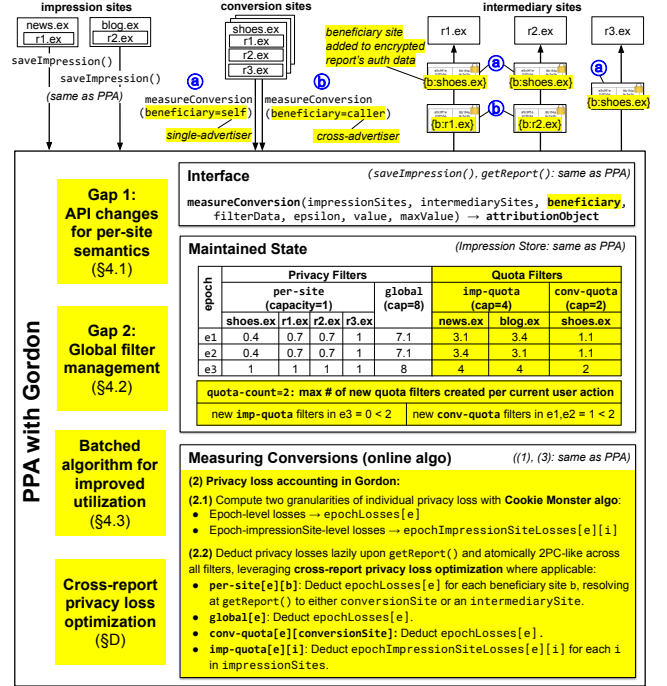


Fig. 2. Gordon architecture. Changes vs. PPA (Fig. 1(c)) in yellow.

### 3.2 Running example

We update the *shoes.ex* example to support *cross-advertiser queries*, a feature PPA plans to add soon. Our Gordon design anticipates this shift, which significantly impacts privacy budget management. To reflect this, we modify the example: *shoes.ex* contracts with *r1.ex* and *r2.ex* for ad placement and evaluation as before, but now *r1.ex* and *r2.ex* also optimize placements across advertisers and content sites. They will each therefore be interested in obtaining two encrypted reports for each conversion: one for single-advertiser measurement on behalf of *shoes.ex* and one for cross-advertiser optimization on their own behalf. Additionally, we introduce *r3.ex*, which focuses solely on single-advertiser measurements and specializes in cross-intermediary reporting, providing a complete view of *shoes.ex*'s ad performance across the two placement intermediaries *r1.ex* and *r2.ex*. *r3.ex* will require only one encrypted report for the single-advertiser measurement on *shoes.ex*'s behalf.

### 3.3 Gordon architecture

Fig. 2 shows Gordon's architecture, with proposed changes to PPA highlighted in yellow (relative to Fig. 1(c)). Gordon modifies all three layers of PPA: the interface, the privacy filter architecture, and how privacy loss is accounted for during conversion measurement and report requests. These changes span four major conceptual shifts (yellow boxes on the left).

**Gap 1: API changes for per-site semantics (§4.1).** We modify the API, protocol, and terms of use to eliminate ambiguity in budget attribution. Specifically, we introduce a *beneficiary site* parameter, authenticate it to the aggregation service, and enforce its use—both technically and contractually—so that

reports can only support the intended site's DP queries. These changes prevent intermediaries from misusing reports funded by conversion sites, restoring clear per-site semantics for parties that comply with the protocol and its terms. This addresses PPA's Gap 1 from §2.6.

**Gap 2: Global filter management (§4.2).** Without changing the API or protocol, we rework PPA's internal state to restore its intended stock-and-flow model of privacy loss, where user actions create "stock" at impression sites and trigger "flow" at conversion sites. Depletion attacks break this structure by automating flows or collapsing domain roles, draining the global filter without real user input. To defend against this, we introduce three quotas: one limits how much stock an impression site can create, another caps how much flow a conversion site can trigger, and a third bounds how many new sites can participate per user action. These quotas don't just limit indirect proxies (like API calls or intermediaries); they act directly on the core protected resource—the global filter—enforcing a privacy budget flow tightly coupled to actual user behavior and curbing adversarial misuse. We show that this leads to graceful degradation of utility for benign workloads under attack. This addresses PPA's Gap 2.

**Batched algorithm to improve utilization (§4.3).** Static quotas can underutilize the global privacy budget, especially when some impression sites see little demand while others face heavy load. To address this inefficiency, Gordon introduces a batched scheduling algorithm that collects unserved requests over a time interval and reallocates unused impression-site quota toward them at the end of each interval. This approach still provides DoS resilience and is provably Pareto efficient within each device epoch.

**Cross-report privacy loss optimization (§D).** Gordon introduces an optimization that reduces overcounting when multiple reports for the same conversion involve disjoint impression sets (e.g., across intermediaries). Since such reports reveal no more than the original attribution histogram to shared filters, Gordon accounts for them jointly across privacy and quota filters. We defer details to Appendix D, as they are not vital to understanding our main contributions. §D.1 illustrates how Fig. 2's filter states are calculated from it.

## 4 Detailed Design

### 4.1 API changes for per-site semantic (Gap 1)

We begin by addressing ambiguities in PPA's per-site filters, which aim to ensure privacy against honest-but-curious actors but currently fall short. Intermediaries like *r1.ex* and *r2.ex* can request reports on behalf of *shoes.ex*, causing PPA to deduct privacy loss from *shoes.ex*'s budget—even though *r1.ex* and *r2.ex* receive the reports and may reuse them for their own analytics. If those same intermediaries later run cross-advertiser queries (e.g., to train a model to choose between ads for *shoes.ex*, *toys.ex*, or *tvx.ex* based on content-site context), PPA charges their budgets directly. But when a single entity serves both roles, the line between client-serving

and self-serving blurs. Even honest actors may be tempted to misuse reports charged to others. Conversion sites may also shard themselves into subdomains (e.g., *shoes-cart.ex*, *shoes-purchase.ex*) to extend their budget. Without clear constraints on report use, per-site accounting loses semantic integrity.

Gordon changes PPA's API, protocol, and terms of use to clarify the *beneficiary* for each DP query. In the **API**, we add a `beneficiary` parameter to `measureConversion()`. During `getReport()`, browsers resolve the beneficiary: to the conversion site for single-advertiser measurement, or to the requesting intermediary for cross-advertiser optimization. Privacy loss is then charged to the beneficiary's per-epoch filters, which are created as needed. Under the honest-but-curious model, we permit unrestricted filter creation. In the **protocol**, the beneficiary is included in the report's authenticated data, and aggregators are required to reject any batch with inconsistent beneficiaries. This blocks intermediaries from reusing reports charged to other clients' budgets. In the **terms of use**, we prohibit using DP results tied to one `beneficiarySite` to benefit another. Reports and results must remain siloed by beneficiary, even across shared infrastructure. This prohibits report-sharing among sharded identities (e.g., *shoes-cart.ex*, *shoes-purchase.ex*), cross-company collusion, and Sybil behavior (§3.1). Honest-but-curious sites will avoid these.

Together, these changes make per-site privacy loss accounting meaningful and enforceable under clear assumptions. PPA has already adopted the protocol update, and we plan to propose the API and policy changes shortly. Appendix A models beneficiary-based behavior, as required for subsequent proofs.

**Example.** In Fig. 2, *shoes.ex* issues two `measureConversion()` calls for a \$60 purchase: (a) one for its own use (`beneficiary = self`), and (b) one for intermediaries (`beneficiary = caller`). In the first case, intermediaries like *r1.ex*, *r2.ex*, and *r3.ex* request reports on behalf of *shoes.ex*, which deduct from its budget. In the second, *r1.ex* and *r2.ex* request reports on their own behalf, triggering deductions from their own budgets. *r3.ex* does not participate and preserves its budget. Each encrypted report includes the beneficiary in authenticated data: `b:shoes.ex` for single-advertiser use (a); `b:r1.ex` and `b:r2.ex` for cross-advertiser reports (b).

### 4.2 Online global filter management (Gap 2)

With clarified semantics, per-site filters offer *strong privacy protection against honest-but-curious sites*, assuming tight configuration (e.g., capacity  $\epsilon_{\text{per-site}} = 1$ ). But non-compliant behavior remains possible, making the global filter essential to safeguard against worst-case privacy loss—i.e., an adversary capable of accessing and combining results from *all sites*. To enforce both DP guarantees, Gordon implements a two-phase commit-like algorithm: data-driven (non-null) reports are returned only if they can be atomically funded by both per-site and global filters; otherwise, null reports are returned. Appendix B.2 formalizes this algorithm and its dual cross-granularity (individual) DP guarantee—a relevant property in

practice that, to our knowledge, has never been formalized.

A key challenge in managing the global filter is balancing competing goals: supporting benign workloads, resisting depletion attacks on this shared resource, and minimizing its guarantee to offer the strongest privacy protections that practical deployment can afford. We exemplify anticipated depletion attacks, then discuss limitations of existing defenses.

**DoS depletion attacks.** An adversary  $X$  may attempt to exhaust the global budget—either to boost their own utility or to disrupt others’. This threat already exists in PPA through single-advertiser queries and will grow with cross-advertiser support. To carry it out,  $X$  registers  $s = \epsilon_{\text{global}} / \epsilon_{\text{per-site}}$  Sybil domains and distributes queries across them.

*Attack 1: Cross-advertiser reports.*  $X$  builds a site embedding the  $s$  Sybil domains as intermediaries. When user  $u$  visits, the site: (1) registers  $s$  impressions with  $X$  as the conversion site and a different Sybil as intermediary; and (2) has each intermediary request a cross-advertiser report, exhausting its per-site budget. This drains the global budget for  $u$ . If many users visit  $X$  once per epoch,  $X$  can disrupt others’ measurements for that epoch. If users continue arriving across epochs,  $X$  can sustain disruption, mounting a persistent attack with one popular site and just one visit per user per epoch. PPA isn’t currently vulnerable, lacking cross-advertiser support. But a similar attack works using single-advertiser reports:

*Attack 2: Single-advertiser reports.* Here, the Sybils serve as both impression and conversion domains. When  $u$  visits  $X$ ,  $X$  auto-redirects  $s$  times, switching domains to register impressions and trigger single-advertiser reports in lock step; each Sybil can register impressions for a different Sybil as the conversion site. As before, this depletes the global budget. While aggressive redirection may be heuristically flagged, redirection is too common for browsers to block outright.

*Attack 3: Single-advertiser reports, subtler version.* Upon user  $u$ ’s visit,  $X$  (1) registers  $s$  impressions with Sybil conversion sites, and (2) redirects once to load a new Sybil domain that requests a report. Reports use maximum attribution windows to draw budget across past epochs via impressions previously registered by  $X$ . If many users visit  $X$ ’s site roughly  $s$  times during each epoch’s data lifetime (typically months),  $X$  can sustain global budget depletion—again with one site but requiring multiple user visits.

**Limitations of existing defenses.** Some behaviors in these attacks clearly exceed reasonable use and should be disabled. (1) PPA is designed for cross-site measurement, so queries with identical `impressionSite` and `conversionSite` (Attack 1) should be disallowed. (2) A single user action shouldn’t simultaneously register an impression and trigger conversion measurement of it—even across domains (Attack 1). (3) Excessive redirection (Attack 2) should be detectable. (4) Allowing an epoch’s entire global budget to be exhausted in seconds is a fundamental flaw (Attacks 1 and 2). While these heuristics offer minimal protection, more principled defense

is needed for subtler abuse like Attack 3.

In PATWG discussions, several mitigations have been proposed: restricting which sites receive per-site filters (e.g., via mandatory registration), rate-limiting API calls, or capping the number of sites granted filters per epoch. While potentially useful, these measures risk over-constraining a nascent, evolving workload. Mandatory registration could limit access to the API, undermining the web’s openness. Hard limits on per-site impression counts are tricky: some sites show many ads, others few. The same goes for conversions, which may range from rare purchases to frequent landing-page visits. Capping intermediaries per conversion could constrain advertisers’ ability to work with diverse partners. And if the API is repurposed beyond advertising—e.g., to measure engagement or reach—workload patterns may evolve further. Fixed constraints that seem reasonable today could stifle innovation or penalize legitimate new uses.

**Our approach: Enforce stock-and-flow.** We aim for defenses that make *minimal assumptions about workloads*, enabling browsers to provide strong protection without overly restricting the API. In §2.3, we introduce a stock-and-flow pattern for PPA’s intended use: privacy loss is driven by explicit user actions—like navigations or clicks—across distinct impression and conversion domains. The above attacks break this pattern by automating flows and collapsing domain roles.

We restore the pattern via *quotas*: impression-site quotas cap stock creation, conversion-site quotas cap triggered flow, and a count-based limit bounds the number of new sites that can create the preceding quotas from a single user action. Unlike indirect metrics (e.g., API call frequency, number of intermediaries, or domains with per-site filters), our first two quotas operate directly on the protected resource: the global filter. Each represents a *share* of the global budget calibrated to a browser-defined “normal” workload. The third quota anchors the stock-and-flow pattern to explicit user action. Together, these quotas constrain adversaries to operate within the contours of normal workloads, preventing global budget drainage by a single site with limited user interaction.

**Gordon quota system.** Fig. 2 (yellow background) highlights the internal state maintained by Gordon to manage global privacy filters in PPA. Appendix B formalizes the system’s behavior and proves its privacy and resilience properties. We use two types of quotas: (1) *quota filters*, `imp-quota` and `conv-quota`, which are implemented as DP filters—not for privacy accounting, but to regulate global filter consumption, a novel use in DP literature; (2) a standard *count-based quota* limiting the number of new quota filters created per user action.

The impression-site quota filter, `imp-quota`, is scoped per impression site and per epoch. It bounds the global privacy loss from flows that use stock created by impressions from that site. When site  $i$  first calls `saveImpression()` in epoch  $e$ , Gordon creates `imp-quota[e][i]`, with capacity set to a share of the global filter. This quota is consumed only if a

<b>“Normal” workload parameters:</b>	
<b>M:</b>	max # of impression sites in an epoch contributing to non-zero loss in epoch.
<b>N:</b>	max # of conversion sites that request non-zero loss from an epoch.
<b>n:</b>	max # of conversion sites that request non-zero loss from a single (epoch, impression site) pair.
<b>r:</b>	max budget consumed by an intermediary’s cross-advertiser queries on a single conversion site, as a fraction of the intermediary’s $\epsilon_{\text{per-site}}$ .
<b>Filter</b>	<b>Capacity configuration</b>
Per-site filter	$\epsilon_{\text{per-site}}$ : configuration parameter
Global filter	$\epsilon_{\text{global}} = \max(N, n \cdot M)(1 + r)\epsilon_{\text{per-site}}$
Impression-site quota	$\epsilon_{\text{imp-quota}} = n(1 + r)\epsilon_{\text{per-site}}$
Conversion-site quota	$\epsilon_{\text{conv-quota}} = (1 + r)\epsilon_{\text{per-site}}$

Tab. 1. Gordon filter configurations.

later report matches an impression from  $i$  in that epoch—that is, if  $i$ ’s stock is used.

The conversion-site quota filter,  $\text{conv-quota}$ , is scoped per conversion site and per epoch. It bounds global privacy loss from flows initiated by conversions on that site.  $\text{conv-quota}[e][c]$  is created when site  $c$ , in or after epoch  $e$ , first calls `measureConversion()` in a way that could incur non-zero loss in  $e$ . It is consumed on `getResult()`—i.e., a flow occurs.

Fig. 2 sketches Gordon’s privacy loss accounting algorithm (box “Measuring Conversions”; full version in Appendix B.1). Per-epoch individual privacy losses are first computed using the Cookie Monster algorithm. Then, for each `getReport()` call, Gordon attempts to deduct losses across all relevant filters in an atomic transaction per epoch: success only alters state if all checks pass. A non-null report is returned only if all checks succeed across all epochs. For each epoch  $e$ , relevant filters include the beneficiary’s  $\text{per-site}$  filter, the  $\text{global}$  filter, the conversion site’s  $\text{conv-quota}$ , and an  $\text{imp-quota}$  for each impression site with non-zero loss. To efficiently enforce impression-site quotas, we compute loss at the (epoch, impression site) level and charge it to the corresponding  $\text{imp-quota}$ . Although total quota capacities may exceed the global filter at any moment, Gordon’s atomic checks ensure global budget is never breached.

Quota filters cap how much each first-party site contributes to global privacy consumption. In a world without automatic redirects—where every domain change reflects a user action—this would suffice to reestablish user-driven stock-and-flow. But redirects are pervasive, so we allow a bounded number of first-party domains to trigger new quota creation after a single explicit user action. This bound,  $\text{quota-count}$ , is configurable and expected to be small (e.g., 2 or 3). We also recommend disallowing a single domain from registering both an impression and a conversion on the same user action.

**Configuration to “normal” workload.** How should filters be configured to avoid disrupting benign workloads? We take three steps. First, we define four browser-adjustable parameters describing expected workload scale ( $N, M, n, r$ ), defined in Table 1. Second, given these parameters and  $\epsilon_{\text{per-site}}$ , we express constraints the other capacities must meet to support this workload:  $\epsilon_{\text{conv-quota}} \geq (1+r)\epsilon_{\text{per-site}}$ ;  $\epsilon_{\text{imp-quota}} \geq n \cdot \epsilon_{\text{conv-quota}}$ ;

$\epsilon_{\text{global}} \geq \max(N \cdot \epsilon_{\text{conv-quota}}, M \cdot \epsilon_{\text{imp-quota}})$ . Third, we derive capacity formulas from these constraints, as shown in Table 1.

**Resilience to DoS depletion.** We prove the following:

**Theorem 1** (Resilience to DoS depletion (proof in B.3)). *Consider an adversary who manages to create  $M^{\text{adv}}$  and  $N^{\text{adv}}$   $\text{imp-quota}$  and  $\text{conv-quota}$  filters, respectively. The maximum budget  $\epsilon_{\text{global}}^{\text{adv}}$  that the adversary can consume from the global filter on a device  $d$  is such that:*

$$\epsilon_{\text{global}}^{\text{adv}} \leq \min(M^{\text{adv}} \epsilon_{\text{imp-quota}}, N^{\text{adv}} \epsilon_{\text{conv-quota}}).$$

This blocks *Attack 1*, where all impressions and conversions occur under one domain, yielding  $M^{\text{adv}} = N^{\text{adv}} = 1$  and capping consumption at  $\min(\epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}})$ , far from depletion. The  $\text{quota-count}$  bound blocks *Attack 2*, where a single user visit triggers automatic redirection. This bound—small by design—limits how many quota filters can be created per user action, allowing only modestly more budget use than in *Attack 1*.

In general cases like *Attack 3*, an adversary who receives  $U_{\text{adv}}$  interactions from user  $u$  during epoch  $e$  can create at most  $M^{\text{adv}} + N^{\text{adv}} \leq \text{quota-count} \cdot U_{\text{adv}}$ . Our quotas ensure *graceful degradation* for benign workloads as a function of  $U_{\text{adv}}$ . We prove the following:

**Theorem 2** (Graceful degradation (proof here)). *Consider a adversary collecting  $U^{\text{adv}}$  user actions on sites under their control for device  $d$ . Under the configuration of Table 1, the budget  $\epsilon_{\text{global}}^{\text{adv}}$  that this adversary can consume from the global filter is upper-bounded by:*

$$\epsilon_{\text{global}}^{\text{adv}} \leq (1 + r)\epsilon_{\text{per-site}} \times \frac{n}{1 + n} (\text{quota-count} \times U^{\text{adv}}).$$

*Proof.* Thm. 1 implies the most efficient way to allocate the  $\text{quota-count} \times U^{\text{adv}} = M^{\text{adv}} + N^{\text{adv}}$  filter creations available to the attacker is such that  $M^{\text{adv}} \epsilon_{\text{imp-quota}} = N^{\text{adv}} \epsilon_{\text{conv-quota}}$ , or  $M^{\text{adv}} n(1 + r)\epsilon_{\text{per-site}} = N^{\text{adv}} (1 + r)\epsilon_{\text{per-site}}$ . This yields  $M^{\text{adv}} = \frac{1}{n+1} \text{quota-count} \times U^{\text{adv}}$  and  $N^{\text{adv}} = \frac{n}{n+1} \text{quota-count} \times U^{\text{adv}}$ . Applying Thm. 1 concludes the proof.  $\square$

### 4.3 Batched scheduling to improve utilization

Static quota partitioning can underutilize the global filter—even in benign scenarios. Suppose a device visits only two impression sites: *news.ex*, with many conversions, and *blog.ex*, with just one. The lone advertiser of *blog.ex* is limited by  $\epsilon_{\text{per-site}}$ , leaving much of its  $\text{imp-quota}$  unused, while the many advertisers of *news.ex* are bottlenecked by the  $\text{imp-quota}$  of *news.ex*. As a result, significant global budget remains idle, despite no added privacy risk from reallocating it to blocked *news.ex* advertisers. This limitation, flagged by PATWG participants, motivates our algorithmic solution that improves utilization while maintaining resilience to depletion.

Dynamically adjusting quotas based on observed demand would invite attacks, but we observe that if PPA supports *batched mode*—collecting requests over a period of time and servicing them gradually—we can make smarter schedul-



ing decisions. In particular, we can gradually release unused impression-site quota to support otherwise-blocked requests. The challenge is to (1) preserve some formal resilience guarantees, and (2) avoid scheduling decisions that depend on cross-epoch filter state, which would violate individual DP semantics. We present an algorithm that satisfies both constraints and shows significant utilization gains in evaluation.

**Algorithm.** Algorithm 1 outlines the approach (full version in Appendix C). We divide each epoch data’s lifetime into *T scheduling intervals* (e.g., one week). We extend the PPA API to support a *response time*—the interval after which a report is returned. For privacy, reports are only delivered at their response time; unscheduled requests yield encrypted null reports. Each interval has three phases: (1) *Initialization*: We release a portion  $\epsilon_{\text{global}}/T$  of the

global budget, adding it to any leftover from prior intervals. With both `imp-quota` and `conv-quota` filters active, we try allocating queued requests using this budget. `TryAllocate()` decides whether to attempt allocation for a request  $r$  based solely on public metadata (as required for individual DP); if yes, it removes  $r$  from the queue and applies all active filters. (2) *Online*: As requests arrive, and with both quotas on, we decide immediately based on the same process, whether to allocate or queue them. (3) *Batch*: At the interval’s end, we disable `imp-quota` (keeping `conv-quota`), sort the queue via a max-min-fairness heuristic, and allocate requests one-by-one until no more succeed. `TryAllocate()` always attempts allocation if a request’s response is due.

**Sorting the queue.** Inspired by max-min fairness [5], we sort requests by the impression site with the least estimated budget consumption so far—based only on public metadata, per individual DP constraints. Within each site, requests are ordered by ascending requested privacy budget. For multi-site requests, we sort by site with lowest estimated budget usage.

**Resilience to DoS depletion.** During the online phase, both the `imp-quota` and `conv-quota` quotas are active, preserving the same resilience properties as in Thm. 2. In the batch phase, we lift the `imp-quota`, allowing any unused global filter capacity released so far to be reallocated. This helps support constrained benign workloads—such as some of *news.ex*’s advertisers—but also opens the door to adversarial exploitation. Nonetheless, consumption remains bounded by `conv-quota`, yielding the following bound on adversarial consumption:

$$\epsilon_{\text{global}}^{\text{adv}} \leq (1 + r) \epsilon_{\text{per-site}} \times \text{quota-count} \times (U^{\text{adv}} - 1).$$

This bound is likely pessimistic: successful attacks would need to (1) time online queries to deplete budget ahead of legitimate requests, (2) coordinate hybrid strategies across the online and batch phases, and (3) overcome the scheduler’s sorting mechanism, which favors low-budget and underrepresented impression sites. We leave tighter formal analysis of the batched algorithm’s resilience for future work.

**Pareto efficiency.** Our claim that batching improves global-filter utilization is substantiated both empirically in §6.5 and through a Pareto efficiency property we prove for the batch phase of our preceding algorithm:

**Theorem 3** (Pareto efficiency (proof in C.3)). *The batched global filter scheduling algorithm is Pareto efficient with regards to the global filter budget within a device epoch.*

#### 4.4 Recommendations for PATWG

Gordon provides browsers with foundational building blocks for defending against DoS depletion attacks on PPA’s global filter—though not an end-to-end solution. Operating within the budget management layer, our techniques offer built-in resilience independent of specific web attack vectors. However, they rest on assumptions—namely, that attackers cannot easily induce many users to visit many attacker-controlled domains—which browsers must enforce to achieve full protection. Our threat model (§3.1) leaves enforcement out of scope, but Gordon establishes a foundation to drive end-to-end solutions, which has so far been lacking in PATWG, hindering its progress. We conclude this section with a set of **Don’ts** and **Do’s**, some addressing directions raised in PATWG.

**Don’ts:** (1) *Don’t rate-limit API invocations*: This is not directly useful and risks stifling benign use cases. In Gordon, sites may register arbitrary impressions and conversions, and intermediaries may request any number of reports. The true limit is on how many distinct *domains* can act after a single user action. (2) *Don’t limit the number of per-site filters*: These are meant to track privacy loss from honest-but-curious sites. They are not suitable levers for defending the global filter from malicious actors trying to deplete it. (3) *Don’t require intermediary registration*: With proper budget management—by Gordon and by first parties managing their own quotas—intermediaries do not impact privacy or resilience guarantees. While Gordon leaves to future work the ability for first parties to control how intermediaries consume their quota, we believe this can be done rigorously, further reducing the need for intermediary registration with a PPA authority.

**Do’s:** Focus on *detecting and disabling patterns of site sharding across domains*. For example, sites may attempt to shard themselves—e.g., routing each user interaction through a distinct domain—to inflate their quota access and deplete the global filter. While some level of sharding is inevitable (e.g., legitimate third-party integrations like shopping carts), aggressive self-sharding for DoS purposes should be explicitly prohibited. First, PPA should ban such behavior in its terms of use, which large, legitimate sites are likely to re-

spect. Second, browsers should develop heuristics to detect noncompliant patterns and block offending sites from using the API. One possible signal is when a landing site frequently links to dynamically changing domains that invoke the API before returning users to the same main site. Third, Gordon's sorting algorithm could be extended to penalize suspicious-but-not-yet-blocked behavior. These are examples of concrete, actionable directions that PATWG can now pursue based on the resilience foundation provided by Gordon.

## 5 Prototype

We implement Gordon in two components: (1) `pdslib`, a general-purpose on-device individual DP library in Rust, and (2) its integration into Firefox's Private Attribution, a basic PPA prototype. **pdslib**: Gordon's core logic lives in `pdslib`, a Rust library for privacy budget management designed for broader individual DP use cases beyond advertising—e.g., location services in mobile apps. `pdslib` provides a generic interface: clients (sites or apps) register events (e.g., ad views, location visits) and request reports (e.g., attributions, model updates), receiving encrypted responses under strict privacy and isolation constraints. It implements all filters, quotas, privacy accounting, batching, and cross-report optimizations. Gordon is a PPA-specific instantiation—a 350 LoC shim atop `pdslib`'s 2k LoC, specializing its generics to the PPA spec.

**Firefox integration**: We integrate `pdslib` and the Gordon shim into Firefox's Private Attribution, replacing its primitive report-count-based accounting with full privacy loss tracking (Firefox's PA lacks even Cookie Monster logic). Integration mainly required modifying a single file—XXX. Appendix E shows a Firefox extension dashboard we built to visualize filter and quota usage. We plan to open-source `pdslib`, the shim, and the integration to support PATWG and broader private aggregation cases.

## 6 Evaluation

We seek to answer the following questions: **(Q1)** What parameters define “normal” operation in the Criteo workload? **(Q2)** Do query error rates vary with different quota capacities? **(Q3)** Can quotas preserve low error rates for benign queries under DoS attacks? **(Q4)** Do quotas lead to under-utilization, and can batching mitigate this?

### 6.1 Methodology

**Pierre**: Will shrink to 2.5 pages and move details to appendix.

**Dataset**. We evaluate Gordon on CriteoPrivateAd [7], a dataset released by Criteo specifically to evaluate private advertising systems. The dataset is a sample of 30 days production data collected by Criteo using third-party cookies. It contains 100M impressions displayed on various impression sites (`publisher_id` feature) for different conversion sites (`campaign_id`, which is a good proxy for advertiser domain according to [7]). It only involves a single intermediary (Criteo itself). Each impression contains rich features, including a device identi-

fier that is reset every day, contextual and user features, and attribution information such as whether an impression led to a click, a visit or a sale on a conversion site. This attribution information allows us to generate a list of conversions for each device. Since the Criteo dataset is subsampled at the impression level and not device level, most devices contain only a single impression. Criteo provides a device-level distribution of their real traffic, along with a resampling methodology to extract a smaller dataset from CriteoPrivateAd that follows the device-level distribution. We follow their methodology to obtain a dataset with XXX impressions and XXX conversions. We use the first 10 days of the dataset to design our algorithms and calibrate parameters, and we report results on the previously-unseen last 20 days.

**Evaluation scenario**. We define single-advertiser measurement queries. For each conversion, the conversion site requests a report attributing the conversion to the last touch impression. Impressions are grouped into 5 buckets according to the `features_ctx_not_constrained_0` feature present in Criteo. After aggregation, each advertiser receives a histogram counting how many conversions were attributable to each bucket. The feature could correspond to geographical region, device type or user category.

**Metric**. We adopt  $RMSRE_\tau$ —relative root mean square error truncated at  $\tau$ —to measure DP histogram error, following [1].

**Benign workload process**. Attribution requests need to specify a requested budget  $\epsilon$ . First, we decide on an aggregation batch size  $b$  for each conversion site, based on the average number  $b_0$  of conversions per day for each conversion site. Since it is unfeasible to get high-accuracy results with reasonable epsilon on small batches, we only keep the 73 conversion sites that have  $b_0 \geq 100$ . Then, we take  $b = \min(5000, 10b_0)$  to expect at least one batch over 10 days. Finally, we derive  $\epsilon$  from  $b$  using a target error  $RMSRE_{\tau=0.05 \times b} = 5\%$ .

**Attack workload process**. We augment the benign scenario with a DoS depletion attack akin to Attacks 2 and 3 (since our dataset does not permit cross-advertiser queries for more than one intermediary site). First, we create a malicious impression site  $X$ , by duplicating the impression site  $X^*$  that had the most impressions in Criteo. This ensures that a non-negligible fraction of the devices visit  $X$  at some point, thus enabling the attack to have visible effect on query accuracy. Next, we create new malicious conversion sites  $Y_1^1, \dots, Y_1^s, \dots, Y_k^1, \dots, Y_k^s$  by duplicating  $s$  times each of the  $k = 10$  conversion sites  $Y_1^*, \dots, Y_k^*$  with most conversions attributed to  $X^*$  in Criteo. Finally, we create impressions for  $X$  by copying all the impressions  $X^*$  has for  $Y_1^*, \dots, Y_s^*$ , and we create conversions for  $Y_i^j$  by copying all the conversions that  $Y_i^*$  has for  $X$ .  $s = 1$  corresponds to Attack 3, since conversions happen over time, following the same pattern as the benign conversions  $Y_1^*, \dots, Y_k^*$  have for  $X^*$ . A larger value for  $s$  corresponds to Attack 2, since we have a chain of  $s$  conversion sites  $Y_i^1, \dots, Y_i^s$  (using automatic redirection or tricking the user into clicking)

instead of a single conversion site  $Y_i^*$ .

**Baselines.** We compare Gordon against two baselines. The first baseline is **Cookie Monster**, which only maintains per-site filters. The second baseline, which we call **PPA**, is Cookie Monster augmented by a global filter, which corresponds to the current state of the PPA draft specification with the global filter as a safety limit.

**Defaults.** Unless specified otherwise, we take  $\epsilon_{\text{per-site}} = 1$ ,  $\epsilon_{\text{conv-quota}} = 1$ ,  $\epsilon_{\text{imp-quota}} = 4$ ,  $\epsilon_{\text{global}} = 8$ , corresponding to the 95th percentile in Tab. 2.

## 6.2 “Normal” workload parameters in Criteo (Q1)

In §4.2 we established sufficient conditions under which filters do not harm normal operation, thanks to parameters  $r, N, M$  and  $n$ . In our Criteo workload we have  $r = 0$  since we only consider measurement queries, and can thus take  $\epsilon_{\text{conv-quota}} = \epsilon_{\text{per-site}}$ . However, determining the values of  $N, M$  and  $n$  is not immediate since it would require executing a benign workload to count how many non-zero reports each device-epoch sent to or from various parties. **Pierre: TODO: do that, it’s tighter and easier to explain in the text. Keeping rough statistics as a placeholder, but they are overestimations (e.g., count of “authorized” conversion sites for each impression site, instead of actual conversions that happen).** Instead, we can derive upper bounds  $\tilde{N} \geq N, \tilde{M} \geq M, \tilde{n} \geq n$  by computing simple statistics from a dataset of impressions and conversions, as detailed next. For  $\tilde{N}$ , we compute a percentile of the distribution of number of unique conversion sites in each device-epoch, and multiply by the maximum attribution window length (2 in Criteo). This gives a crude upper bound on the maximum number of times most devices are queried, which is itself an upper bound on the maximum number of times most devices are queried *with non-zero loss*. For  $\tilde{M}$ , we take a percentile of the number of unique impressions sites across device-epochs. For  $\tilde{n}$ , we take a percentile of the number of unique conversion sites across impression sites and device-epochs, and multiply by the maximum attribution window length.

%ile	$\tilde{N}$	$\tilde{M}$	$\tilde{n}$	$\epsilon_{\text{global}}$	$\epsilon_{\text{imp-quota}}$
50	2	1	2	2	2
90	4	2	4	8	4
95	4	2	4	8	4
99	6	3	6	18	6
100	12	7	14	98	14

Tab. 2. “Normal” workload configs for Criteo. **TODO(Pierre):** Pls wrap text around this table.

## 6.3 Query errors under normal workload (Q2)

We vary  $\epsilon_{\text{imp-quota}}$  and measure its impact on query error. Fig. 3a reports both median and 95th percentile RMSRE. Since Cookie Monster and PPA lack a impression-site quota filter, their errors remain constant across  $\epsilon_{\text{imp-quota}}$  values. Moreover, Cookie Monster and PPA exhibit identical error, as we set  $\epsilon_{\text{global}}$  high enough to deactivate PPA’s global filter

—effectively reducing it to Cookie Monster. In contrast, Gordon’s error increases at low  $\epsilon_{\text{imp-quota}}$ , as the impression-site quota blocks some reports. For  $\epsilon_{\text{imp-quota}} \geq 2$ , the filter no longer affects benign-query error, suggesting that reasonably-configured quotas can preserve utility.

Fig. 3b attributes this error to each filter type, for Gordon. RMSRE in queries comes from two sources: DP noise (variance), and missing values due to filters and quotas (bias). We focus on the error introduced by filters and quotas. For each report, we compute how many and which filters were out-of-budget while computing the report. If any per-site filter filter is out-of-budget, we label the whole report as potentially biased and attribute the error to the per-site filter filter. We look at global filter filters next, then conversion-site quota filters and finally impression-site quota filters. Next, for each query we compute the number of reports in each category, and divide it by the number of reports in the query. This gives us the fraction of reports in the query affected by the per-site filter (resp. global filter, impression-site quota, conversion-site quota) filter. Finally, we average the fractions over all the queries in the workload.

For  $\epsilon_{\text{imp-quota}} = 1$ , almost a third of reports in each query are filtered out by the impression-site quota filter. This explains the high error seen at low values of  $\epsilon_{\text{imp-quota}}$  in Fig. 3a. For  $\epsilon_{\text{imp-quota}} = 2$  and higher, almost all the biased reports come from per-site filter filters. These filters also exist in Cookie Monster and PPA, which explains why Gordon’s error converges to Cookie Monster’s error in Fig. 3a.

## 6.4 Query errors under DoS attack (Q3)

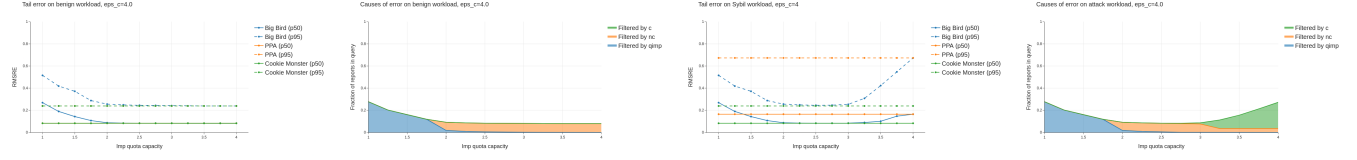
Fig. 3c shows that impression-site quota limits the impact X has on the c-filter, for an attack with  $s = 10$ . We do not include X’s queries in the error. For a well-configured impression-site quota filter, honest queriers are perfectly isolated from X, and are still able to run their own queries without being perturbed by the impression-site quota.

Fig. 3d further drills down into the causes for high error at both ends of the graph in Fig. 3c. For low  $\epsilon_{\text{imp-quota}}$ , the error comes from the impression-site quota filter, as in Fig. 3b. The attack has no direct impact on error, but impression-site quota is hurting even honest queries. For high  $\epsilon_{\text{imp-quota}}$ , the error comes from the global filter filter. This is because the impression-site quota filter is too loose and lets X deplete the global filter filter, thereby denying reports for honest queriers.

## 6.5 Batched system evaluation (Q4)

## 7 Related Work

Our primary contribution is advancing the **PPA API**—an in-progress W3C standard likely to become the default for browser-based advertising measurement, and thus a critical part of the web’s infrastructure. Gordon addresses two foundational gaps in PPA: it clarifies the ambiguous semantics of per-site filters and introduces a system for configuring and managing both these filters and the global filter. This



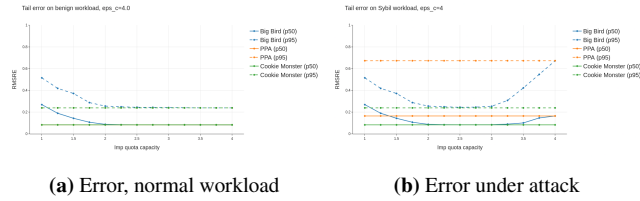
(a) Error, normal workload

(b) Error causes, normal workload

(c) Error under attack

(d) Error causes under attack

**Fig. 3. Quota system evaluation.** (a), (b): Query error and its root causes in a benign case. (c), (d): Benign-query error and root causes under attack. **Roxana:** Changes to graphs: (1) Move the legends inside the graphs (there's space on top given the  $[0,1]$  y axis (which you should keep!). You can order the legend entries 2x2 to fit atop each graph. (2) Enlarge all fonts and line widths and point sizes significantly so they become visible when areas are squeezed. (3) Remove graph titles. (4) Turn everything into black-and-white – no color differentiations pls! You can differentiate based on point types (consistent across the colors hence system types). For area graphs, you can use different textures in place of colors. (5) Make the PPA points much larger than the Cookie Monster points so you can see them from behind CM. (6) Rename Big Bird to Gordon. (7) Label Cookie Monster as “Cookie Monster  $\epsilon_{\text{global}} = \infty$ ”, PPA as “PPA  $\epsilon_{\text{global}} = 4$ ”, Gordon as “Gordon  $\epsilon_{\text{global}} = 4$ ”. (8) x axis label: “Impression site quota capacity ( $\epsilon_{\text{imp-quota}}$ ). x axis label: RMSRE becomes “RMSRE (error, lower is better)”; “Fraction of reports in query” becomes “% of reports in query impacted by a filter.” For (b) and (d), legend should say: “per-site filter,” “global filter,” “quota filter”. (9) Please transform y axis for (b) and (d) to 0-100%, to not confuse it with the other 0-1 RMSRE axes. **Pierre:** Also probably combine all the 4 graphs as a single pdf on the Python side, to share captions etc.



(a) Error, normal workload

(b) Error under attack

**Fig. 4. Batched system evaluation.** **Roxana:** BOGUS FIGURES, DO NOT LOOK.

system maintains strong privacy guarantees, supports benign workloads, and resists global filter depletion attacks.

From the standpoint of prior foundational contributions to PPA and related private advertising APIs, the most closely related work is Cookie Monster [8], which efficiently accounts for privacy loss within per-epoch filters tied to individual “queriers.” Gordon goes further, showing how to *manage* these fine-grained filters alongside the global filter. It also introduces a new cross-report privacy loss optimization, whereas Cookie Monster’s techniques apply only to single reports. Other foundational contributions include Google’s ARA papers [ara-paper, other-ara-papers-on-dp], work on the MPC components of these APIs [prio, mpc-ara-papers, ask-ben-about-ipa-mpc-paper], and informal proposals for now-obsolete API designs like IPA, PAM, and Hybrid [hybrid, 3, 6]. More broadly, there is significant research on privacy-preserving ad targeting [sebastian-angels-works].

Beyond ad measurement, this paper contributes to the broader challenge of **privacy budget management**, a fundamental yet underexplored aspect of DP systems. While budget allocation *within* a single DP query has been well studied [multiple-papers], managing budgets *across* queries from mutually distrustful parties remains less understood. Related work includes scheduling of shared global budgets [pack, 5, 4], from which our batched scheduling design draws inspiration but differs in key ways. Gordon operates under epoch-level individual DP, introducing constraints that existing privacy schedulers cannot handle—such as the need to not rely

on cross-epoch available-budget information to make scheduling decisions. Further, these systems lack support for queriers issuing multiple tasks and do not address DoS depletion. Further afield, **TODO(Pierre): Can you fill in? I didn’t read these papers.** <https://dl.acm.org/doi/10.14778/3467861.3467870> and <https://dl.acm.org/doi/10.1145/3626761>.

Gordon builds upon and contributes to literature on **privacy filters** [empty citation], especially individual filters [empty citation]. These are basic budget management primitives covering composition accounting and halting to enforce DP. We use them both conventionally (to enforce privacy, as others do [multiple-papers]) and in a new role: as quotas that limit consumption to preserve isolation across queriers. A core contribution is our formalization of managing filters at multiple granularities—per-site and global—to support simultaneous DP guarantees. While essential in practice, such guarantees had not been formalized, nor had prior work shown how to configure and manage them. This paper does both.

## 8 Conclusions

## References

- [1] Hidayet Aksu et al. “Summary Reports Optimization in the Privacy Sandbox Attribution Reporting API”. In: *Proc. Priv. Enhancing Technol.* 2024.4 (2024), pp. 605–621. DOI: 10.56553/POPETS-2024-0132. URL: <https://doi.org/10.56553/popets-2024-0132>.
- [2] Hamid Ebadi, David Sands, and Gerardo Schneider. “Differential Privacy: Now It’s Getting Personal”. In: *Proceedings of the 42nd Annual ACM SIGPLAN SIGACT Symposium on Principles of Programming Languages*. POPL ’15: The 42nd Annual ACM SIGPLAN SIGACT Symposium on Principles of Programming Languages. Mumbai India: ACM, Jan. 14, 2015, pp. 69–81. ISBN: 978-1-4503-3300-9. DOI: 10.1145/2676726.2677005.
- [3] *Interoperable Private Attribution (IPA)*. <https://github.com/patcg-individual-drafts/ipa>. 2022.
- [4] Nicolas Küchler et al. “Cohere: Privacy Management in Large Scale Systems”. In: *CoRR* abs/2301.08517 (2023). DOI: 10.48550/ARXIV.2301.08517. arXiv: 2301.08517. URL: <https://doi.org/10.48550/arXiv.2301.08517>.
- [5] Tao Luo et al. “Privacy Budget Scheduling”. In: *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*. USENIX Association, July 2021, pp. 55–74. ISBN: 978-1-939133-22-9. URL: <https://www.usenix.org/conference/osdi21/presentation/luo>.
- [6] *Private Ad Measurement (PAM)*. <https://github.com/patcg-individual-drafts/private-ad-measurement>. 2023.
- [7] Mehdi Sebbar et al. *CriteoPrivateAd: A Real-World Bidding Dataset to Design Private Advertising Systems*. 2025. arXiv: 2502.12103 [cs.CR]. URL: <https://arxiv.org/abs/2502.12103>.
- [8] Pierre Tholoniati et al. “Cookie Monster: Efficient On-Device Budgeting for Differentially-Private Ad-Measurement Systems”. In: *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles*. SOSP ’24. New York, NY, USA: Association for Computing Machinery, Nov. 15, 2024, pp. 693–708. ISBN: 9798400712517. DOI: 10.1145/3694715.3695965.



## A API changes for per-site semantic (Gap 1)

This section formalizes API changes to clarify the per-site semantics. Starting from Cookie Monster's formalism, we adapt it to capture Gordon's notion of beneficiaries. While this section does not present a standalone result, its formalism underpins the main theorems proved later and evoked in the body of the paper. We begin by mapping terminology between Cookie Monster and Gordon (§A.1), then outline modeling restrictions (§A.2), after which we define Gordon's data (§A.3) and query models (§A.4), and conclude with sensitivity analyses needed in subsequent sections (§A.5).

### A.1 Sites and roles

In this paper, we align with the terminology of PPA, which differs from that of Cookie Monster, to better align with PPA. Moreover, we make sites appear explicitly in the data and query model. We take a set of sites  $\mathcal{S}$  (e.g., domain name). The same site can appear under different roles:

- impression site: site where an impression occurs (publisher in Cookie Monster)
- conversion site: site where a conversion occurs (advertiser in Cookie Monster)
- beneficiary site: site that receives the results of a DP query (querier in Cookie Monster)

### A.2 Modeling restrictions

We focus here on single-beneficiary settings, captured by the following assumptions, which we then relax in our study of the cross-report optimization in §D:

- The set of public events for a beneficiary site  $b$  is the set of all conversions where  $b$  appears as beneficiary site. In other words, we hardcode  $P = C_b$ . This prevents publisher reports or certain adtech optimization queries, that were technically supported by Cookie Monster.<sup>1</sup>
- Each report has a single beneficiary site: either the conversion site (measurement report) or another site (optimization report). In this formalism, a conversion site can register the same conversion data multiple times with different beneficiary sites if it wants to execute both measurement and optimizations.

### A.3 Data model

A database  $D$  is a set of device-epoch records where each record  $x = (d, e, F) \in \mathcal{X} = \mathcal{D} \times \mathcal{E} \times \mathcal{P}(\mathcal{S} \times \mathcal{I} \cup \mathcal{S} \times \mathcal{S} \times \mathcal{C})$  contains a device  $d$ , an epoch  $e$  and a set of impression and conversion events  $F$ . Each event  $f \in F$  contains the site (impression site  $i$  or conversion site  $c$ ) where the event occurred:  $f = (i, \text{imp}) \in \mathcal{S} \times \mathcal{I}$  or  $f = (c, b, \text{conv}) \in \mathcal{S} \times \mathcal{S} \times \mathcal{C}$ . Additionally, conversions contain the beneficiary site  $b$  that will receive the conversion report.<sup>2</sup>

**Definition 1** (Filter ( $\mathcal{F}_x$ )). For each device-epoch record  $x = (d, e, F)$ , we maintain several types of privacy filters:

- $\mathcal{F}_x^{\text{nc}[b]}$ : Non-collusion filter for beneficiary site  $b$ , with capacity  $\epsilon_{\text{nc}[b]}$

- $\mathcal{F}_x^c$ : Collusion filter with capacity  $\epsilon_c$ .
- $\mathcal{F}_x^{q\text{-conv}[c]}$ : Conversion site quota filter for site  $c$ , with capacity  $\epsilon_{\text{conv-quota}}$
- $\mathcal{F}_x^{q\text{-imp}[i]}$ : Impression site quota filter for site  $i$ , with capacity  $\epsilon_{\text{imp-quota}}$

Each filter maintains a state of consumed budget and provides operations:

- $\text{canConsume}(\epsilon)$ : Returns *TRUE* if the filter can accommodate additional privacy loss  $\epsilon$ . In particular, letting  $\mathbb{I}_{\text{pass}}$  be how it's defined in lemma 1,
  - For a per-site filter, global filter, or *conv-quota* filter, return *TRUE* iff it satisfies,

$$\epsilon_x^t \leq \epsilon_{\text{initial}} - \sum_{k \in [t-1]} \epsilon_x^k \cdot \mathbb{I}_{\text{pass}}[k], \quad (1)$$

where  $\epsilon_{\text{initial}}$  is the respective initialized privacy budget of the global filter and conversion-site quota filters in Line 8 of Alg. 3, and  $\epsilon_x^t$  is the privacy loss at step  $k$  from the corresponding filter (epoch-level budget consumption is the same as "ComputeIndividualBudget" defined in appendix D of [8]).

- For an impression-site quota filter, return *TRUE* iff it satisfies

$$\epsilon_x^{i,t}[i] \leq \epsilon_{\text{imp-quota}} - \sum_{k \in [t-1]} \epsilon_x^{i,k}[i] \cdot \mathbb{I}_{\text{pass}}[k], \quad (2)$$

where  $\epsilon_x^{i,t}[i]$  is the site-level privacy loss at step  $k$  from the *imp-quota* filter (site-level budget consumption in Alg. 5).

- $\text{tryConsume}(\epsilon)$ : Deducts privacy loss  $\epsilon$  from the filter's remaining capacity iff  $\text{canConsume}$  on all relevant filters return *TRUE*, and the amount consumed is by basic compositions of pure DP filters.

### A.4 Query model

**Definition 2** (Attribution function, adapted from Cookie Monster). Fix a set of relevant impression sites  $\mathbf{i}_A \subset \mathcal{S}$  and a set of impressions<sup>3</sup> relevant to the query  $F_A \subset \mathbf{i}_A \times \mathcal{I}$ . Fix  $k, m \in \mathbb{N}^*$  where  $k$  is a number of epochs. An attribution function is a function  $A : \mathcal{P}(\mathcal{I})^k \rightarrow \mathbb{R}^m$  that takes  $k$  event sets  $F_1, \dots, F_k$  from  $k$  epochs and outputs an  $m$ -dimensional vector  $A(F_1, \dots, F_k)$ , such that only relevant events contribute to  $A$ . That is, for all  $(F_1, \dots, F_k) \in \mathcal{P}(\mathcal{I})^k$ , we have:

$$A(F_1, \dots, F_k) = A(F_1 \cap F_A, \dots, F_k \cap F_A). \quad (3)$$

**Definition 3** (Report identifier and attribution report, same as Cookie Monster). Fix a domain of report identifiers  $\mathbb{Z}$ . Consider a mapping  $d(\cdot)$  from report identifiers  $R$  to devices  $\mathcal{D}$  that gives the device  $d_r$  that generated a report  $r$ .<sup>4</sup>

Given an attribution function  $A$ , a set of epochs  $E$  and a report identifier  $r \in \mathbb{Z}$ , the attribution report  $\rho_{r,A,E}$ , or  $\rho_r$  for short, is a function over the whole database  $D$  defined by:

$$\rho_r : D \in \mathbb{D} \mapsto A(D_{d_r}^E). \quad (4)$$

**Definition 4** (Query, same as Cookie Monster). *Consider a set of report identifiers  $R \subset \mathbb{Z}$ , and a set of attribution reports  $(\rho_r)_{r \in R}$  each with output in  $\mathbb{R}^m$ .<sup>5</sup> The query for  $(\rho_r)_{r \in R}$  is the function  $Q : \mathbb{D} \rightarrow \mathbb{R}^m$  is defined as  $Q(D) := \sum_{r \in R} \rho_r(D)$  for  $D \in \mathbb{D}$ .*

### A.5 Sensitivity analyses

The Cookie Monster paper analyzes global and individual sensitivities of queries at device-epoch level (§C of [8]). In Gordon, we additionally need such analyses at device-epoch-site-level, as some of our quota systems operate at device-epoch-site granularity (notably the impression-site quotas)—see Alg. 5. This section provides analysis for these sensitivities: Thm. 4 and give the global sensitivities of generating reports and answering queries, respectively. Thm. 8 and Thm. 9 give the individual sensitivity versions.

We first start global sensitivity analysis of generating reports.

**Theorem 4** (Global sensitivity of reports per epoch-site). *Fix a report identifier  $r$ , a device  $d$ , a set of epochs  $E_r = \{e_1^{(r)}, \dots, e_k^{(r)}\}$ , and a set of sites  $I_r^{(e)} = \{i_1^{(e)}, \dots, i_{m_e}^{(e)}\}$  for each epoch  $e \in E_r$ . Let  $A(\cdot)$  be the attribution function of interest, so that the corresponding report is  $\rho : D \mapsto A(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ . We have  $\Delta(\rho)$*

$$= \max_{\substack{i \in [k], e := e_i^{(r)}, j \in [m_e], e' = e_k^{(r)} \\ \mathcal{F} := (F_{1,1}, \dots, F_{k,m_{e'}}) \subseteq (\mathcal{S} \times \mathcal{S} \times \mathcal{C} \cup \mathcal{S} \times \mathcal{I})^k}} \|A(\mathcal{F} : F_{i,j} := \emptyset) - A(\mathcal{F})\|_1. \quad (5)$$

*Proof.* Fix a report  $\rho$ . Let  $k$  be  $|E_r|$ . Then, for  $i \in [k]$ , we enumerate through the epochs in  $E_r$ , and call the  $i$ -th epoch  $e_i$ . By definition of global sensitivity:

$$\Delta(\rho) = \max_{D, D' \in \mathbb{D} : \exists x \in \mathcal{X}, D' = D + x} \|\rho(D) - \rho(D')\|_1, \quad (6)$$

from which we expand what  $\rho$  function does:

$$\begin{aligned} \Delta(\rho) &= \max_{\substack{D, D' \in \mathbb{D} : \\ \exists x \in \mathcal{X}, D' = D + x}} \|A(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) - A((D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}})\|_1 \\ &= \max_{\substack{D, D' \in \mathbb{D} : \\ \exists x = (d, e, F) \in \mathcal{X} : \\ e \in E_r, D' = D + x, \\ \forall (i, \text{imp}) \in F : i \in I_r^{(e)}}} \|A(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) - A((D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}})\|_1, \end{aligned} \quad (7)$$

because for  $(d', e', F) \in \mathcal{X}$  with  $d' \neq d$  or  $e' \notin E_r$ , or any  $(i, \text{imp}) \in F$  where  $i \notin I_r^{(e)}$ , they would not be considered as in the computation of  $A(\cdot)$  and  $A(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) = A((D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ .

Next, we show that the following sets are equal:

- $\{(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}, (D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) : D, D' \in \mathbb{D} : \exists x = (d, e, F) \in \mathcal{X} : e \in E_r, D' = D + x, \forall (i, \text{imp}) \in F : i \in I_r^{(e)}\}$ .
- $\{(\mathcal{F} : F_{i,j} = \emptyset, \mathcal{F}) : i \in [k], e := e_i^{(r)}, j \in [m_e], e' = e_k^{(r)}, \mathcal{F} := (F_{1,1}, \dots, F_{k,m_{e'}}) \subseteq (\mathcal{S} \times \mathcal{S} \times \mathcal{C} \cup \mathcal{S} \times \mathcal{I})^k\}$

We show that for all instances on the one side, there exists a corresponding instance on the other. In one direction, take any tuple from the first set, where  $(D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}} = D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}} + x$ , for some  $x = (d, e, F)$ . Firstly, the  $x$  satisfies  $d = d_r$ ,  $e \in E_r$ , and  $\forall f = (i, \text{imp}) \in F, i \in I_r^{(e)}$ . Therefore, we construct  $\mathcal{F}$  as follows:

$$F_{i,j} = \{f \in F : f = (j, \text{imp}) \text{ or } (c, q, \text{conv}) \text{ is relevant to site } j\}.$$

Note that every  $F_{i,j}$  is set independently from  $x$ , except for the one that contains the events for site  $j_0$  such that  $F$  contains either  $(j_0, \text{imp})$  or  $(c, q, \text{conv})$  that is associated with  $j_0$ , in epoch  $i_0$  such that  $e_{i_0}^{(r)} = e \in E_r$ . Since  $x \notin D$ , we know that  $F_{i_0, j_0}$  is set to  $\emptyset$  on the side corresponding to  $D$ , and it is set to contain all events related to  $j_0$  in epoch  $i_0$  on the side corresponding to  $D'$ . That is, the corresponding instance in the second set would be  $(\mathcal{F} : F_{i_0, j_0} = \emptyset, \mathcal{F})$ .

Conversely, let  $(\mathcal{F} : F_{i,j} = \emptyset, \mathcal{F})$  be from the second set. Then, we know the set of events corresponding to site  $j$  in epoch  $i$  is empty for the first element. So, we let  $x := (d, e, F) \in \mathcal{X}$ , where  $d$  is the same devices as was fixed,  $e = e_i^{(r)}$ , and  $F$  contains events related to site  $j$  in epoch  $i$ . Then, we let  $(D, D') \in \mathbb{D} \times \mathbb{D}$ , where  $D' = D + x$  and everything in  $D$  and  $D'$  other than  $x$  corresponds to the rest of the  $F_{i,j}$  that are the same in both  $\mathcal{F}$  and  $\mathcal{F} : F_{i,j} = \emptyset$ . As such,  $(D_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}}, (D')_d^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$  is the corresponding element in the first set.

Finally, we substitute the first set in equation line (8) by the second set accordingly, and the resulting max should be equal as the two sets are equivalent, and the resulting form will be the one claimed.  $\square$

Immediately, we can acquire the following corollary about global sensitivity of generating reports, with attribution functions that have  $m$ -dimensional outputs as they are in our setting.

**Corollary 5.** *Suppose the attribution function  $A(\cdot)$  has an  $m$ -dimensional outputs, and each dimension of any  $\mathcal{F} \subseteq (\mathcal{S} \times \mathcal{S} \times \mathcal{I} \cup \mathcal{S} \times \mathcal{C})^k$  satisfy  $\forall i \in [m], A(\mathcal{F})_i \in [0, A^{\max}]$ , then we have  $\Delta(\rho) \leq mA^{\max}$ .*

*Proof.* Given the form in Thm. 4, we know  $\Delta(\rho)$  equals the max over  $\|A(\mathcal{F} : F_{i,j} = \emptyset) - A(\mathcal{F})\|_1 = \sum_{i=1}^m |A(\mathcal{F} : F_{i,j} = \emptyset)_i - A(\mathcal{F})_i|$ . But, we know for each dimension,  $A(\mathcal{F})_i \in [0, A^{\max}]$ , so the absolute value of the difference in each dimension is  $\in [0, A^{\max}]$ . Thus, the sum over  $m$  dimensions is  $\|A(\mathcal{F} : F_{i,j} = \emptyset) - A(\mathcal{F})\|_1 \leq mA^{\max}$ .  $\square$

Next, we analyze the global sensitivity of answering queries.

**Theorem 6** (Global sensitivity of queries per epoch-site). *Let  $R$  be the set of report identifiers relevant to a query  $Q$ . Then, we write  $(\rho_r, d_r, E_r)_{r \in R}$  as the reports, devices, and epoch windows corresponding to each  $r \in R$ . Then, we have*

$$\Delta(Q) \leq \max_{(d,e,F) \in \mathcal{X}} \sum_{\substack{r \in R: d=d_r, e \in E_r, \\ \exists i \in I_r^{(e)}: F \text{ contains events related to } i}} \Delta(\rho_r).$$

*Proof.* By definition of global sensitivity,

$$\Delta(Q) = \max_{D, D' \in \mathbb{D}: \exists x \in \mathcal{X}, D' = D + x} \|Q(D) - Q(D')\|_1 \quad (9)$$

$$= \max_{x \in \mathcal{X}} \max_{D, D' \in \mathbb{D}: D' = D + x} \|Q(D) - Q(D')\|_1. \quad (10)$$

Let  $x = (d, e, F) \in \mathcal{X}$ , where  $F$  contains events associated with site  $j$ . For  $r \in R$  such that  $d \neq d_r$  or  $e \notin E_r$ , or  $F$  contains events related to site  $i \notin I_r^{(e)}$ , we have  $\rho(D) = \rho(D')$ . So, by applying triangle-inequality in the mean time:

$$\|Q(D) - Q(D')\|_1 \leq \sum_{r \in R: d=d_r, e \in E_r, i \in I_r^{(e)}} \|\rho_r(D) - \rho_r(D')\|_1 \quad (11)$$

$$\leq \sum_{r \in R: d=d_r, e \in E_r, i \in I_r^{(e)}} \Delta(\rho_r), \quad (12)$$

where the second inequality is by definition of  $\Delta(\rho_r)$ . Note that this bound is independent from the choice of  $D$  and  $D'$ , so we take the max over  $(d, e, i)$  and get

$$\Delta(Q) \leq \max_{d,e,i} \sum_{r \in R: d=d_r, e \in E_r, i \in I_r^{(e)}} \Delta(\rho_r).$$

But, the choice of  $i$  can be used to set  $F$  accordingly, so, equivalently, this suffices to prove the theorem.  $\square$

We can tighten the upper bound in the preceding theorem if we additionally require each device-epoch-site to participate in at most one report.

**Corollary 7.** *If each device-epoch-site participates in at most one report, then  $\Delta(Q) = \max_{r \in R} \Delta(\rho_r)$ .*

*Proof.* Since each device-epoch-site participates in at most one report,  $\sum_{r \in R: d=d_r, e \in E_r, i \in I_r^{(e)}} \Delta(\rho_r) = \Delta(\rho_r)$ . So,

$$\Delta(Q) \leq \max_{(d,e,F) \in \mathcal{X}} \Delta(\rho_r).$$

But then,  $\Delta(Q)$  must match the sensitivity of one of the  $r$ 's in this case, so the inequality is tight:  $\Delta(Q) = \max_{r \in R} \Delta(\rho_r)$ . Particularly,  $\forall r, \exists D \sim D'$ , such that  $\|\rho_r(D') - \rho_r(D)\|_1 = \Delta(\rho_r)$ .  $\square$

Next, we shift our gears to individual sensitivity analyses. We first analyze the individual sensitivity of generating reports.

**Theorem 8** (Individual sensitivity of reports per epoch-site). *Fix a report identifier  $r$ , a device  $d_r$ , a set of epochs  $E_r = \{e_1^{(r)}, \dots, e_k^{(r)}\}$ , a set of sites  $I_r^{(e)} = \{i_1^{(e)}, \dots, i_{m_e}^{(e)}\}$  for each epoch  $e \in E_r$ , an attribution function  $A$  with relevant events*

$F_A$ , and the corresponding report  $\rho : D \mapsto A(D_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ . Fix a device-epoch record  $x = (d, e, F) \in \mathcal{X}$ , where  $F \subseteq S \times S \times C \cup S \times I$ , so that  $x_i = (d, e, F_i)$  is the projection where  $F_i$  contains only events related to site  $i$ .

If the report requests a single epoch  $E_r = \{e_r\}$  as well as a single site in the one epoch,  $I_r^{(e_r)} = \{i_r\}$ , then we have:

$$\Delta_{x_i}(\rho) = \begin{cases} \|A(F_i) - A(\emptyset)\|_1 & , \text{ if } d = d_r, e = e_r \text{ and } i = i_r \\ 0 & , \text{ otherwise.} \end{cases} \quad (13)$$

In particular, it should be intuitive to see why  $\mathcal{F}$  is a single  $F_i$  in the first case, because we only have one epoch which contains one site, so there should be only one set of events corresponding to the one site in the one epoch.

Otherwise, either  $|E_r| \geq 2$  or  $|I_r^{(e)}| \geq 2$  for some  $e \in E_r$ , or both, and so we have:

$$\Delta_{x_i}(\rho) \leq \begin{cases} \Delta(\rho) & , \text{ if } d = d_r, e \in E_r, i \in I_r^{(e)} \text{ and } F_i \cap F_A \neq \emptyset \\ 0 & , \text{ otherwise.} \end{cases} \quad (14)$$

*Proof.* Fix a report  $\rho$  and  $x_i = (d, e, F_i) \in \mathcal{X}$ . Consider any  $D, D' \in \mathbb{D}$  such that  $D' = D + x_i$ . We have  $\rho(D) = A(D_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$  and  $\rho(D') = A((D')_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ .

- First, if  $d \neq d_r, e \notin E_r$ , or  $i \notin I_r^{(e)}$ , then  $(D')_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}} = D_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}}$ . Hence,  $\|\rho(D) - \rho(D')\|_1 = 0$  for all such  $D, D'$ , which implies  $\Delta_{x_i}(\rho) = 0$ .
- Next, suppose that the report requests a single epoch  $E_r = \{e_r\}$  with a single site  $I_r^{(e_r)} = \{i_r\}$ :
  - If  $d = d_r, e = e_r$ , and  $i = i_r$ , then since  $D + x_i = D'$ , we must have  $(d_r, e_r, F_i) \notin D$ , and thus  $D_{d_r}^{e_r, i_r} = \emptyset$ . On the other hand,  $(D')_{d_r}^{e_r, i_r} = F_i$  (restricted to events relevant to site  $i_r$ ). Thus,  $\|\rho(D) - \rho(D')\|_1 = \|A(F_i) - A(\emptyset)\|_1$ .
  - If  $d \neq d_r, e \neq e_r$ , or  $i \neq i_r$ , then  $(d, e, F_i)$  doesn't change the outcome and  $(D')_{d_r}^{i_r} = D_{d_r}^{i_r}$ . Hence,  $\|\rho(D) - \rho(D')\|_1 = 0$ .
- Now, suppose that the report requests either an arbitrary range of epochs  $E_r$  each of whom has at least one site, or a single epoch that has multiple sites  $I_r^{(e_r)}$ :
  - If  $d \neq d_r, e \notin E_r$ , or  $i \notin I_r^{(e)}$ , then  $A((D')_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) = A(D_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ , i.e.,  $\|\rho(D') - \rho(D)\|_1 = 0$ .
  - If we have  $d = d_r, e = e_j^{(r)} \in E_r$ , and  $i \in I_r^{(e)}$ , but  $F_i$  is simply not related to the attribution request, i.e.  $F_i \cap F_A = \emptyset$ . Then, by definition of  $F_A$ , we have  $A((D')_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}}) = A(D_{d_r}^{E_r, \{I_r^{(e)}\}_{e \in E_r}})$ , i.e.,  $\|\rho(D) - \rho(D')\|_1 = 0$ .
  - Otherwise, it must be the case that  $d = d_r, e = e_j^{(r)} \in E_r, i \in I_r^{(e)}$  and  $F_i \cap F_A \neq \emptyset$  and there are events in

the intersection that is related to some site  $i$  in epoch  $e$ , so we have:

$$\|\rho(D) - \rho(D')\|_1 = \|A(\mathcal{F} : F_{j,i} = \emptyset) - A(\mathcal{F})\|_1, \quad (15)$$

where  $j$  is the index of epoch  $e$  in  $E_r$ , and  $F_{j,i}$  represents the relevant events for site  $i$  in epoch  $e_j^{(r)}$ .

The first two cases are independent over choices of  $D \sim D'$ , so taking the max over such choices still gives  $\Delta_{x_i}(\rho) = 0$ . Unfortunately, the third identity does depend on the choice of  $D \sim D'$ , and taking the max only gives the general definition of global sensitivity, in the worst case. Particularly,

$$\Delta_{x_i}(\rho) = \max_{\mathcal{F}=\{F_{j,i}\}: F_{j,i} \subseteq S \times S \times C \cup S \times I} \|\rho(D) - \rho(D')\|_1 \quad (16)$$

$$= \max_{\mathcal{F}=\{F_{j,i}\}: F_{j,i} \subseteq S \times S \times C \cup S \times I} \|A(\mathcal{F} : F_{j,i} = \emptyset) - A(\mathcal{F})\|_1 \quad (17)$$

$$\leq \Delta(\rho), \quad (18)$$

where the last inequality is tight due to the definition of global sensitivity and Thm. 4, for the worst case choice of  $x_i$  in general.  $\square$

Finally, we analyze the individual sensitivity of answering queries.

**Theorem 9** (Individual sensitivity of queries per device-epoch-site). *Fix a query  $Q$  with corresponding report identifiers  $R$  and reports  $(\rho_r)_{r \in R}$ . Fix a device-epoch-site record  $x_i = (d, e, F_i) \in \mathcal{X}$ , where  $F_i$  contains only events related to site  $i$ . We have:*

$$\Delta_{x_i}(Q) \leq \sum_{r \in R} \Delta_{x_i}(\rho_r) \quad (19)$$

*In particular, if  $x_i$  participates in at most one report  $\rho_r$ , then  $\Delta_{x_i}(Q) = \Delta_{x_i}(\rho_r)$ .*

*Proof.* Take  $D, D' \in \mathcal{D}$  such that  $D' = D + x_i$ . By the triangle inequality:

$$\Delta_{x_i}(Q) = \max_{D'=D+x_i \in \mathcal{D}} \|Q(D) - Q(D')\|_1 \quad (20)$$

$$= \max_{D'=D+x_i \in \mathcal{D}} \left\| \sum_{r \in R} \rho_r(D) - \rho_r(D') \right\|_1 \quad (21)$$

$$\leq \sum_{r \in R} \max_{D'=D+x_i \in \mathcal{D}} \|\rho_r(D) - \rho_r(D')\|_1 \quad (22)$$

$$\leq \sum_{r \in R} \Delta_{x_i}(\rho_r), \quad (23)$$

where the last inequality is by the definition of individual sensitivity.

When  $x_i$  participates in at most one report  $\rho_{r_0}$ , we have  $\Delta_{x_i}(\rho_r) = 0$  for all  $r \neq r_0$ . Therefore,  $\Delta_{x_i}(Q) \leq \Delta_{x_i}(\rho_{r_0})$ . This inequality is tight because there exists a pair  $(D^*, D'^*)$  with  $D'^* = D^* + x_i$  such that  $\|\rho_{r_0}(D'^*) - \rho_{r_0}(D^*)\|_1 = \Delta_{x_i}(\rho_{r_0})$ , and for this pair,  $\|Q(D^*) - Q(D'^*)\|_1 = \Delta_{x_i}(\rho_{r_0})$ .  $\square$

## B Global Filter Management (Gap 2)

### B.1 Algorithm

Gordon manages per-site and global privacy filters using the quota mechanisms described in §4.2. Alg. 2 presents an abstract model of Gordon's operation, capturing how it answers beneficiary queries sequentially. The function `AnswerQuery` corresponds to the implementation's `measureConversion()`.

Alg. 3 depicts the functionality triggered on receiving a report request (i.e., `getReport()`). Gordon checks and consumes budget from the relevant filters and prunes the resulting report based on filter status.

First, Gordon ensures all filters are initialized: the global filter, relevant per-site filter filters, and the impression-site quota and conversion-site quota filters. Second, it computes the privacy losses incurred—both at the epoch level (§C[8]) and at the site level (Alg. 5). Third, it checks whether all filters have sufficient budget and attempts to consume it. Because consumption must be atomic at the report level, Gordon uses a two-phase commit protocol to deduct privacy losses from multiple filters. Alg. 4 formalizes this: if any filter cannot afford its share of the privacy loss, the report is zeroed out, and no budget is consumed from any filters for that report.

Finally, Gordon returns the (possibly empty) report as the response to the query at the current time step.

### B.2 Privacy proofs

**Mechanisms.** Alg. 2 defines two types of interactive mechanisms. First, for each beneficiary site  $b$  we can denote by  $\mathcal{M}^b$  the interactive mechanism that only interacts with  $b$ . Second,  $\mathcal{M}$  is the interactive mechanism that interacts with all the beneficiary sites concurrently. Remark that the database  $D$  is fixed upfront for simplicity, but a reasoning identical to [8, Alg. 2] generalizes to adaptively generated data. Another simplification compared to [8] is that public events are never relevant events for our attribution functions (i.e., the output of a conversion report can only depend on impressions, not on other conversion). This is enforcing the constraint on queries mentioned in Case 1 of [8, Thm. 1].

**Levels of accounting.** `EpochImpSiteBudget` computes privacy loss at a different granularity than `EpochBudget`, using Def. 5.

**Definition 5** (Device-epoch-site neighborhood relation). *Consider a device  $d \in \mathcal{D}$ , an epoch  $e \in \mathcal{E}$ , an impression site  $i \in \mathcal{S}$  and a set of impression events happening on  $i$ :  $F_i \in \{i\} \times \mathcal{I}$ .<sup>8</sup>*

*We say  $D \sim_{d,e,i,F_i} D'$  if there exists  $D_0$  and  $F \in \mathcal{S} \setminus \{i\} \times \mathcal{I}$  such that:*

$$\{D, D'\} = \{D_0 + (d, e, F), D_0 + (d, e, F \cup F_i)\} \quad (24)$$

**Theorem 10.** *Consider  $x \in \mathcal{X}$  on device  $d$ , with global filter capacity  $\epsilon_{\text{global}}$  and per-site filter capacity  $\epsilon_{\text{per-site}}$ . The two following properties hold simultaneously:*

**Algorithm 2** Gordon Notations and Setup

---

```

1: Input
2: Database  $D$ 
3: Stream of adaptively chosen queries
4: function  $\mathcal{M}(D)$ 
5:    $(S_b)_{b \in \mathcal{S}} = (\emptyset)_{b \in \mathcal{S}}$ 
6:   for  $(d, e, F) \in D$  do
7:     for  $f \in F : f = (c, b, \text{conv})$  do
8:       Generate report identifier  $r \xleftarrow{\$} U(\mathbb{Z})$ 
9:       // Save mapping from  $r$  to the device that generated it
10:       $d_r \leftarrow d$ 
11:       $S_b \leftarrow S_b \cup \{(r, f)\}$ 
12:  // Each beneficiary receives its public events and corresponding report identifiers
13:  for  $b \in \mathcal{S}$  do
14:    output  $S_b$  to  $b$ 
15:  // Beneficiaries ask queries interactively
16:  for  $t \in [t_{\max}]$  do
17:    receive  $Q_t^{b_t}$  from beneficiary site  $b_t$ .
18:    output  $\text{AnswerQuery}(Q_t^{b_t})$  to  $b_t$ 
19:  // Collect, aggregate and noise reports to answer  $Q$ 
20:  function  $\text{AnswerQuery}(Q)$ 
21:    Get report identifiers  $R$  and noise parameter  $\sigma$  from  $Q$ 
22:    for  $r \in R$  do
23:      Read  $Q$  to get conversion site  $c$ , beneficiary site  $b$ , impression sites  $i$ , target epochs  $E$ , attribution function  $A$  for report  $r$ .
24:       $\rho_r \leftarrow \text{GenerateReport}(d, c, b, i, E, A, \sigma)$ 
25:      Sample  $X \sim \mathcal{L}(\sigma)$ 
26:  return  $\sum_{r \in R} \rho_r + X$ 

```

---

1.  $\mathcal{M}$  satisfies individual device-epoch  $\epsilon_{\text{global-DP}}$  for  $x$  under public information  $C$ .
2. For each beneficiary site  $b \in \mathcal{S}$ ,  $\mathcal{M}^b$  satisfies individual device-epoch  $\epsilon_{\text{per-site-DP}}$  for  $x$  under public information  $C_b$ .

*Proof.* Take a device-epoch  $x = (d, e, F) \in \mathcal{X}$  and a database  $D$  that doesn't contain  $(d, e)$ . Denote by  $x_C = (d, e, F \cap C)$  the device-epoch obtained by keeping only public events  $C$  from  $x$ , where public events are the set of all conversions. Take  $v \in \text{Range}(\mathcal{M})$ .

**1. Our first goal is equivalent to showing that:**

$$\left| \ln \left( \frac{\Pr[\mathcal{M}(D + x_C) = v]}{\Pr[\mathcal{M}(D + x) = v]} \right) \right| \leq \epsilon_{\text{global}}. \quad (25)$$

Consider any database  $D'$ .  $v$  is the vector of values returned at different points in Alg. 2. We split it into  $v = (v_{\text{pub}}, v_1, \dots, v_{t_{\max}})$  where  $v_{\text{pub}}$  is a value for the output from Line 14 of Alg. 2, and  $v_t$  is the output for query  $t$  at Line 18. We denote by  $\mathcal{M}_{\text{pub}}(D')$  the random variable of the output at Line 14, and  $\mathcal{M}_t(D')$  the random variable of the output at Line 18. By con-

**Algorithm 3** Gordon Algorithm (on device)

---

```

1: Input
2: Filter and quota capacities  $\epsilon_{\text{global}}, \epsilon_{\text{per-site}}, \epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}}$ 
3: // Generate report and update on-device budget
4: function  $\text{GenerateReport}(d, c, b, i, E, A, \sigma)$ 
5:   for  $e \in E$  do
6:      $x \leftarrow (d, e, D_d^e)$ 
7:     if  $\mathcal{F}_x$  is not defined then
8:        $\epsilon_{\text{global}}, \epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}} \leftarrow$ 
       Capacities( $\epsilon_{\text{per-site}}, N, M, r$ )
9:        $\mathcal{F}_x \leftarrow \text{InitializeFilters}(\epsilon_{\text{global}}, \epsilon_{\text{per-site}}, \epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}})$ 
10:       $F_e \leftarrow D_d^e$ 
11:      // Compute individual device-epoch-level privacy losses, same as "ComputeIndividualBudget" defined in appendix D of [8]
12:       $\epsilon_x^t \leftarrow \text{EpochBudget}(x, d, E, A, \mathcal{L}, \sigma)$ 
13:       $\epsilon_x^{i,t} \leftarrow \{\}$  // Initialize empty map
14:      for  $i \in i$  do
15:        // Compute individual device-epoch-site-level privacy losses, Alg. 5
16:         $\epsilon_x^i \leftarrow \text{EpochImpSiteBudget}(x, i, d, E, A, \mathcal{L}, \sigma)$ 
17:         $\epsilon_x^{i,t}[i] \leftarrow \epsilon_x^i$ 
18:        // Atomic filter check and update, Alg. 4
19:        if  $\text{AtomicFilterCheckAndConsume}(\mathcal{F}_x, b, c, i, \epsilon_x^t, \epsilon_x^{i,t}) = \text{FALSE}$  then
20:           $F_e \leftarrow \emptyset$  // Empty the report if any filter check fails
21:           $\rho \leftarrow A((F_e)_{e \in E})$  // Clipped attribution report
22:      return  $\rho$ 

```

---

ditioning over past outputs  $(v_{\text{pub}}, v_1, \dots, v_{t-1})$  at each time step  $t \in [t_{\max}]$  we get:

$$\Pr[\mathcal{M}(D') = v] = \Pr[\mathcal{M}_{\text{pub}}(D') = v_{\text{pub}}] \cdot \prod_{t=1}^{t_{\max}} \Pr[\mathcal{M}_t(D') = v_t | v_{<t}]. \quad (26)$$

Take  $t \in [t_{\max}]$ . By Algorithm 2 we have:

$$\Pr[\mathcal{M}_t(D') = v_t | v_{<t}] = \Pr[\text{AnswerQuery}(Q_t; D', \mathcal{F}_t) = v_t], \quad (27)$$

where the query  $Q_t$  and the state of the privacy filters (across all device-epochs)  $\mathcal{F}_t$  are functions of past results  $v_{<t}$ , and  $\text{AnswerQuery}$  is defined in Alg. 2. In particular, the behaviors of  $Q_t$  and  $\mathcal{F}$  are fixed given the views from previous steps, as  $v_{<t}$ .

Finally, if we denote by  $\rho_r(D'; \mathcal{F})$  the filtered report returned by Alg. 3 we get:

$$\Pr[\mathcal{M}_t(D') = v_t | v_{<t}] = \Pr \left[ \sum_{r \in R_t} \rho_r(D'; \mathcal{F}_{t,r}) + X_t = v_t \right], \quad (28)$$



**Algorithm 4** 2-Phase Commit Subroutine**Input:**

- 1:  $\epsilon_x^t$ : epoch-level privacy loss for a particular query
- 2:  $\epsilon_x^{i,t}$ : epoch-site-level privacy loss for a particular query
- 3: canConsume: function as is defined in Def. 1
- 4: tryConsume: function as is defined in Def. 1

**Output:**

- 5: Boolean function if all filters have enough budget for the privacy loss  $\epsilon_x^t$  or not.
- 6: **function** AtomicFilterCheckAndConsume( $\mathcal{F}_x, b, c, i, \epsilon_x^t, \epsilon_x^{i,t}$ )
- 7: *// Phase 1: Prepare - check if all filters can consume*
- 8: **if**  $\mathcal{F}_x^{\text{per-site filter}[b]}$ .canConsume( $\epsilon_x^t$ ) = FALSE **then**
- 9:   **return** FALSE
- 10: **if**  $\mathcal{F}_x^{\text{global filter}}$ .canConsume( $\epsilon_x^t$ ) = FALSE **then**
- 11:   **return** FALSE
- 12: **if**  $\mathcal{F}_x^{\text{conv-quota}[c]}$ .canConsume( $\epsilon_x^t$ ) = FALSE **then**
- 13:   **return** FALSE
- 14: **for**  $i \in i$  **do**
- 15:   **if**  $\mathcal{F}_x^{\text{imp-quota}[i]}$ .canConsume( $\epsilon_x^{i,t}[i]$ ) = FALSE **then**
- 16:     **return** FALSE
- 17: *// Phase 2: Commit - consume from all filters*
- 18: *// Privacy filters under no collusion*
- 19:  $\mathcal{F}_x^{\text{per-site filter}[b]}$ .tryConsume( $\epsilon_x^t$ )
- 20: *// Privacy filter under collusion*
- 21:  $\mathcal{F}_x^{\text{global filter}}$ .tryConsume( $\epsilon_x^t$ )
- 22: *// Quota filter for conversion site<sup>6</sup>*
- 23:  $\mathcal{F}_x^{\text{conv-quota}[c]}$ .tryConsume( $\epsilon_x^t$ )
- 24: *// Privacy filters for impression sites*
- 25: **for**  $i \in i$  **do**
- 26:   *// Individual device-epoch-impression site loss.<sup>7</sup>*
- 27:    $\mathcal{F}_x^{\text{imp-quota}[i]}$ .tryConsume( $\epsilon_x^{i,t}[i]$ )
- 28: **return** TRUE

where  $X_t$  is the Laplace noise added at time  $t$ . This equality is a direct quantification of “QueryAnswer” in Alg. 2, as for each  $r \in R_t$ , “QueryAnswer” generates filtered reports and sums over these reports with the Laplace noise  $X_t$  added to them.

Now, we instantiate  $D'$  to be a neighboring dataset of  $D$ , specifically  $D' = D + x_C$ , then we have:

$$\Pr[\mathcal{M}(D + x_C) = v] = \prod_{t=1}^{t_{\max}} \Pr \left[ \sum_{r \in R_t} \rho_r(D + x_C) + X_t = v_t \right] \quad (29)$$

$$= \prod_{t=1}^{t_{\max}} \Pr \left[ \sum_{r \in R_t} \rho_r(D) + X_t = v_t \right], \quad (30)$$

where the last equality is by Def. 2 where  $\rho_r(D + x_C) = \rho_r(D + x_C \cap F_A) = \rho_r(D \cap F_A)$  since  $F_A \subset \mathcal{I}$  and  $x_C \cap \mathcal{I} = \emptyset$ . Since the filter is unambiguous, we also abused the notation by letting  $\rho_r(\cdot) := \rho_r(\cdot, \mathcal{F}_{t,r})$ .

**Algorithm 5** Compute epoch-site level privacy budget**Input:**

- 1:  $x$ : device-epoch record ( $d, e, F$ )
- 2:  $i$ : impression site
- 3:  $d$ : device
- 4:  $E$ : set of epochs, from which we can extract set of impression sites in these epochs relevant to  $A$
- 5:  $A$ : attribution function
- 6:  $\mathcal{L}$ : parameterized noise distribution
- 7:  $\sigma$ : noise scale

**Output:**

- 8: Returns the epoch-site-level individual privacy loss for impression site  $i$  at device-epoch  $x$
- 9: **function** EPOCHIMPSITEBUDGET( $x, i, d, E, A, \mathcal{L}, \sigma$ )
- 10: Extract epoch  $e$  from  $x = (d, e, F)$
- 11:  $i_x \leftarrow \{i \in S \mid x = (d, e, F) \wedge (i, \text{imp}) \in F\}$
- 12: *// Relevant epoch-site level events for current epoch*
- 13:  $E_{\text{relevant}} \leftarrow \{f \in F \mid f = (i, \text{imp}) \wedge \|A(f)\|_1 > 0\}$
- 14: *// Compute epoch-site-level privacy losses*
- 15: **if**  $E_{\text{relevant}} = \emptyset$  **then**
- 16:   *// Case 2 of equations (13) or (14): No relevant epoch-site events in epoch*
- 17:    $\Delta \leftarrow 0$
- 18:   *// Check if we have single or multiple epoch-sites in E*
- 19:   **if**  $|E| = 1$  AND  $|\{\text{site } i \text{ relevant to } A \text{ in epoch } e \in E\}| = 1$  **then**
- 20:     *// Case 1 of equation (13): Single epoch, and only one impression site with relevant events to A in this epoch*
- 21:      $\Delta \leftarrow \Delta_{x_i}(\rho)$ , where  $\Delta_{x_i}(\rho)$  is the upper bound on the value with the same name specified in theorem 8
- 22:   **else**
- 23:     *// Case 1 of equation (14): Request either touches multiple epochs or multiples epoch-sites with relevant events*
- 24:      $\Delta \leftarrow \Delta(\rho)$ , where  $\Delta(\rho)$  is the upper bound on the value with the same name specified in theorem 4
- 25:
- 26: **return**  $\Delta / \sigma$

Similarly, we get:

$$\Pr[\mathcal{M}(D + x) = v] = \prod_{t=1}^{t_{\max}} \Pr \left[ \sum_{r \in R_t} \rho_r(D + x) + X_t = v_t \right] \quad (31)$$

$$= \prod_{t=1}^{t_{\max}} \Pr \left[ \sum_{r \in R_t^x} \rho_r(D + x) + \sum_{r \in R_t \setminus R_t^x} \rho_r(D) + X_t = v_t \right], \quad (32)$$

where  $R_t^x := \{r \in R_t \mid d_r = d, e \in E_r\}$  is the set of reports where  $x$  is queried.

Take a report  $r \in R_t^x$ . Recall that  $\text{pass}_r$  denotes whether  $r$  passed the atomic filter check in Alg. 4.

- If  $\text{pass}_r = 0$ , we have  $\rho_r(D + x) = \rho_r(D)$  because of Alg. 3, Line 20. Note that  $\text{pass}_r = 0$  can happen even if  $\mathcal{F}^{\text{global filter}}$  has enough budget, for instance if the per-site filter or a quota is out of budget.
- If  $\text{pass}_r = 1$ , we have  $\|\rho_r(D + x) - \rho_r(D)\|_1 \leq \Delta_x \rho_r \leq \sigma \cdot \epsilon_r \text{pass}_r$ , where  $\epsilon_r \text{pass}_r$  is the individual device-epoch loss for  $x$  successfully deducted from  $\mathcal{F}^{\text{global filter}}$  for report  $r$ . Since  $\text{pass}_r$  implies that  $r$  passes  $\mathcal{F}^{\text{global filter}}$ , by definition of  $\mathcal{F}^{\text{global filter}}$ , the accumulated loss over all reports is below the filter capacity, i.e.,

$$\sum_{t=1}^{t_{\max}} \sum_{r \in R_t} \epsilon_r \text{pass}_r = \sum_{t=1}^{t_{\max}} \sum_{r \in R_t^x} \epsilon_r \text{pass}_r \leq \epsilon_{\text{global}}, \quad (33)$$

where the first equality is because  $\Delta_{x'} \rho_{x'} = 0$  for  $x \neq x' \in R_t^x$  not queried.

Finally, by property of the Laplace distribution, for  $t \in [t_{\max}]$  we have:

$$\left| \ln \left( \frac{\Pr[\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\Pr[\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \leq \Delta_x Q_t / \sigma \quad (34)$$

$$\leq \sum_{r \in R_t} \Delta_x \rho_r / \sigma \quad (35)$$

$$\leq \sum_{r \in R_t} \epsilon_r \text{pass}_r, \quad (36)$$

where  $\sigma$  is the noise scale for the Laplace distribution.

Let's put everything together. First, by definition of  $x_C$ , we observe that all the outputs  $v_{\text{pub}}$  at Line 14 are identical across both worlds. For the query outputs, we apply the triangle inequality across all time steps (Ineq. (40)):

$$\left| \ln \left( \frac{\Pr[\mathcal{M}(D + x_C) = v]}{\Pr[\mathcal{M}(D + x) = v]} \right) \right| \quad (37)$$

$$= \left| \ln \left( \frac{\prod_{t=1}^{t_{\max}} \Pr[\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\prod_{t=1}^{t_{\max}} \Pr[\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \quad (38)$$

$$= \left| \sum_{t=1}^{t_{\max}} \ln \left( \frac{\Pr[\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\Pr[\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \quad (39)$$

$$\leq \sum_{t=1}^{t_{\max}} \left| \ln \left( \frac{\Pr[\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\Pr[\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \quad (40)$$

$$\leq \sum_{t=1}^{t_{\max}} \sum_{r \in R_t} \epsilon_r \text{pass}_r = \sum_{t=1}^{t_{\max}} \sum_{r \in R_t^x} \epsilon_r \text{pass}_r \quad (41)$$

$$\leq \epsilon_{\text{global}}, \quad (42)$$

where Ineq. (41) follows from Ineq. (36) and Ineq. (42) follows from Ineq. (33). This establishes the first property:  $\mathcal{M}$  satisfies individual device-epoch  $\epsilon_{\text{global}}$ -DP for  $x$  under public information  $C$ .

**2. We now consider the view of a single beneficiary  $b \in \mathcal{S}$ .**

Our goal is equivalent to showing that:

$$\left| \ln \left( \frac{\Pr[\mathcal{M}_b(D + x_C) = v]}{\Pr[\mathcal{M}_b(D + x) = v]} \right) \right| \leq \epsilon_{\text{per-site}}. \quad (43)$$

Let  $\mathcal{M}$  be the global mechanism that processes queries from all beneficiaries. Partition its queries into those from the single beneficiary  $b$  and those from other parties (the “adv”):

$$Q = Q_b \cup Q_{\text{adv}} \quad \text{where} \quad Q_b = \{t \in [t_{\max}] : b_t = b\}. \quad (44)$$

Define

$$\mathcal{M}_b(Q_b) = (\mathcal{M}(Q_b \cup Q_{\text{adv}}))|_{\text{indices in } Q_b}, \quad (45)$$

such that  $\mathcal{M}_b$  behaves exactly like  $\mathcal{M}$  on the queries issued by  $b$ , but treats the adversary's queries as fixed. Concretely, from  $b$ 's viewpoint, the queries in  $Q_{\text{adv}}$  and their outputs are independent arbitrary events that  $b$  has no control over (i.e. the adversary sees and makes the same queries in both worlds).

Hence, for any database  $D \in \mathbb{D}$  and row  $x \in \mathcal{X}$ , we have the equality of distributions w.r.t. the outputs that  $b$  observes:

$$\Pr[\mathcal{M}_b(D + x) = v] = \Pr[(\mathcal{M}(D + x))|_{Q_b} = v]. \quad (46)$$

This justifies focusing on  $\mathcal{M}_b$  rather than the full mechanism  $\mathcal{M}$  when bounding the privacy loss for  $b$ .

Recall that, as before, the public information is identical across both worlds, so we can focus on queries only —more specifically, queries that are observed by  $b$ :

$$\Pr[\mathcal{M}_b(D + x) = v] = \prod_{t \in [t_{\max}] : b_t = b} \Pr \left[ \sum_{r \in R_t} \rho_r(D + x; \mathcal{F}_{t,r}) + X_t = v_t \right]. \quad (47)$$

Secondly, we can bound the privacy loss  $\epsilon_t$  at any given time  $t \in [t_{\max}]$  by the property of Laplace distribution:

$$\left| \ln \left( \frac{\Pr[\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\Pr[\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \leq \sum_{r \in R_t} \Delta_x \rho_{b_t}^t / \sigma \quad (48)$$

$$\leq \sum_{r \in R_t} \epsilon_r \text{pass}_r \quad (49)$$

Finally, we use the fact that  $\text{pass}_r$  implies that  $r$  passes  $\mathcal{F}^{\text{per-site filter}}$  successfully, and thus:

$$\sum_{t \in [t_{\max}] : b_t = b} \sum_{r \in R_t} \epsilon_r \text{pass}_r \leq \epsilon_{\text{per-site}}. \quad (50)$$

Analogous to the first part, we conclude with a triangle

inequality:

$$\left| \ln \left( \frac{\Pr [\mathcal{M}_b(D + x_C) = v]}{\Pr [\mathcal{M}_b(D + x) = v]} \right) \right| \quad (51)$$

$$\leq \sum_{t \in [t_{\max}]: b_t = b} \left| \ln \left( \frac{\Pr [\sum_{r \in R_t} \rho_r(D) + X_t = v_t]}{\Pr [\sum_{r \in R_t} \rho_r(D + x) + X_t = v_t]} \right) \right| \quad (52)$$

$$\leq \sum_{t \in [t_{\max}]: b_t = b} \epsilon_t \leq \sum_{t \in [t_{\max}]: b_t = b} \sum_{r \in R_t} \epsilon_r \text{pass}_r \quad (53)$$

$$\leq \epsilon_{\text{per-site}} \quad (54)$$

□

### B.3 DoS resilience proofs

This section proves our main resilience result for Gordon's quota-based online algorithm: Thm. 1. First, Lemma 1 shows that the 2-PC check (Alg. 4) ensures (1) atomic consumption across all filters relevant to a query, and (2) when all filters have sufficient budget, each consumes an amount proportional to its level-specific sensitivity—either at epoch or at epoch-site level. Next, Lemma 2 bounds the total privacy budget the adversary can consume from global filter at any qualified time, using the atomic consumption guarantees of Lemma 1. This final bound directly implies Thm. 1.

**Definition 6.** We define the following notations for privacy loss accounting:

- $\epsilon_{\text{global}}^{\leq t}$ : the cumulative privacy loss charged to the global filter up to step  $t$  before 2PC is triggered at step  $t$ .
- $D^{\leq e}$ : The subset of database  $D$  containing records with epochs up to and including  $e$ .

We first formalize the atomicity property of the 2-PC algorithm for consuming privacy budgets from relevant filters when Gordon answers a query at any time step  $k$  (Alg. 4)

**Lemma 1** (2-phase commit filter guarantees). *For query  $k$ , let:*

$$\text{pass}(k) \triangleq \begin{cases} 1 & \text{if AtomicFilterCheckAndConsume returns} \\ & \text{TRUE for query } k \\ 0 & \text{otherwise} \end{cases} \quad (55)$$

The AtomicFilterCheckAndConsume function in Algorithm 4 guarantees the following properties:

For any query  $k$  processed by AtomicFilterCheckAndConsume, if  $\text{pass}(k) = 1$ , then

1. **Epoch-level Consistency Property:** the per-site filter, global filter, and conversion-site quota-filter all consume exactly the same amount of budget  $\epsilon_x^t$  for that query.
2. **Epoch-site-level Consistency Property:** the impression-site quota filter consumes exactly  $\epsilon_x^{i,t}[i]$ , which represents the device-epoch-impressionsite-level individual privacy loss.

*Proof.* We can prove both properties at the same time. Fix an arbitrary query  $k$ , for which  $\text{pass}(k) = 1$ . From Algorithm 4, we observe that AtomicFilterCheckAndConsume returns TRUE if and only if: (1) all canConsume checks in Phase 1 pass, and (2) All tryConsume operations in Phase 2 are executed. For a query from conversion site  $c$  with beneficiary site  $b$  and impression sites  $i$ , the function calls:

- $\mathcal{F}_x^{\text{per-site filter}[b]}$ .tryConsume( $\epsilon_x^t$ )
- $\mathcal{F}_x^{\text{global filter}}$ .tryConsume( $\epsilon_x^t$ )
- $\mathcal{F}_x^{\text{conversion-site quota}[c]}$ .tryConsume( $\epsilon_x^t$ )
- For each  $i \in i$ :  $\mathcal{F}_x^{\epsilon_{\text{imp-quota}}[i]}$ .tryConsume( $\epsilon_x^{i,t}[i]$ )

Note that  $\epsilon_x^t$  is computed once at EpochBudget in Line 12 of Algorithm 3 and represents the device-epoch-level individual privacy loss. Similarly, each  $\epsilon_x^{i,t}[i]$  is computed once via EpochImpSiteBudget and represents the device-epoch-impressionsite-level individual privacy loss. Therefore, when  $\text{pass}(k) = 1$ , the per-site filter, global filter, and conversion-site quota filters all consume exactly the same amount  $\epsilon_x^t$ , while each impression-site quota filter consumes its specific amount  $\epsilon_x^{i,t}[i]$ , which is proportional to its sensitivity at the impression site level. □

With such atomic guarantees for every step  $k$  up to some time  $t$ , we can show the following upper bounds for how much an adversary can deplete the global filter budget by the end of time  $t$ . For notations, at step  $t$  in Line 16, suppose that beneficiary site  $b$  requests a report  $\rho_{r,E,A}$  with noise  $\sigma$  through conversion site  $c$  for impression sites  $i$ . Consider a device-epoch  $x$ , with individual budget  $\epsilon_x^t$  computed at Line 12 in Alg. 3. Denote by  $N^{\leq t, \text{adv}}$  the number of conversion sites in  $\text{bad}_c$  with respect to  $x$  that were queried with non-zero budget by step  $t$ . Denote by  $M^{\leq t, \text{adv}}$  the number of impression sites in  $\text{bad}_i$  with respect to  $x$  that were queried with non-zero budget by step  $t$ .

**Lemma 2.** *At the end of time  $t$ , given that the adversary has created at most  $M^{\text{adv}}$  and  $N^{\text{adv}}$  imp-quota and conv-quota filters, respectively, we have the following upper bounds on the global filter filter budget that the adversary can consume:*

- The adversary consumes at most  $M^{\text{adv}} \epsilon_{\text{imp-quota}}$  budget from the global filter.
- The adversary consumes at most  $N^{\text{adv}} \epsilon_{\text{conv-quota}}$  budget from the global filter.

*Proof (first upper bound by  $N^{\text{adv}}$ ).* The total privacy loss in the global filter incurred by the attackers for the global filter by step  $t$ , not inclusive, is:

$$\epsilon_{\text{used}}^{\text{bad}} = \epsilon_{\text{global}}^{\leq t-1, \text{bad}} = \sum_{k=[t-1]: c_k \in \text{bad}_c} \epsilon_{\text{global}}^k \cdot \text{pass}(k). \quad (56)$$

By basic composition under a pure DP filter and the definition of individual privacy loss, the total privacy loss equals:

$$= \sum_{c \in \text{bad}_c} \sum_{k < t: c_k = c} \Delta_x \rho_{b_k}^k \cdot \text{pass}(k), \quad (57)$$

where  $\Delta_x \rho_{b_k}^k$  represents how much the attribution report generated for query  $k$  from beneficiary  $b_k$  can change with respect to  $x$ .

By the consistency property of lemma 1, Alg. 4 ensures that privacy loss is only incurred on  $k$  where  $\text{pass}(k) = 1$ . So, for each conversion site  $c$ , the filter,  $\epsilon_{\text{conv-quota}}[c]$ , precisely tracks the privacy losses incurred, so

$$\epsilon_{\text{conv-quota}}[c]^{\leq t-1} = \sum_{k < t: c_k = c} \Delta_x \rho_{b_k}^k \cdot \text{pass}(k). \quad (58)$$

Substitute this equality into equation (57), we get:

$$\epsilon_{\text{global}}^{\leq t-1, \text{bad}} = \sum_{c \in \text{bad}_c} \epsilon_{\text{conv-quota}}[c]^{\leq t-1}. \quad (59)$$

This sum can be restricted to conversion sites with non-zero privacy loss, i.e.:

$$= \sum_{c \in \text{bad}_c: \epsilon_{\text{conv-quota}}[c]^{\leq t-1} > 0} \epsilon_{\text{conv-quota}}[c]^{\leq t-1} \quad (60)$$

$$\leq \sum_{c \in \text{bad}_c: \epsilon_{\text{conv-quota}}[c]^{\leq t-1} > 0} \epsilon_{\text{conv-quota}}, \quad (61)$$

where  $\epsilon_{\text{conv-quota}}$  is the capacity of each  $\epsilon_{\text{conv-quota}}$  filter. It follows that the number of conversion sites with non-zero privacy loss is precisely  $N^{\leq t, \text{adv}}$ , so:

$$\leq |\{c \in \text{bad}_c : \epsilon_{\text{conv-quota}}[c]^{\leq t-1} > 0\}| \cdot \epsilon_{\text{conv-quota}} \quad (62)$$

$$= N^{\leq t-1, \text{adv}} \cdot \epsilon_{\text{conv-quota}}. \quad (63)$$

Now, during the 2-PC for time  $t$ , we have the following cases:

- Suppose  $\epsilon_x^t$  is a reasonable value, in the sense that it's bounded by the capacity  $\epsilon_{\text{conv-quota}}$ . Then,

$$\epsilon_{\text{global}}^{\leq t, \text{bad}} = \epsilon_{\text{used}}^{\text{bad}} + \epsilon_x^t \leq N^{\leq t, \text{adv}} \cdot \epsilon_{\text{conv-quota}}. \quad (64)$$

- Otherwise,  $\epsilon_x^t$  is unreasonable, in which case  $\epsilon_x^t$  exceeds the capacity  $\epsilon_{\text{conv-quota}}$ . In this case,

$$\epsilon_{\text{conv-quota}}[c_t]^{\leq t-1} + \epsilon_x^t \geq \epsilon_{\text{conv-quota}}, \quad (65)$$

causing  $\mathcal{F}_x^{\epsilon_{\text{conv-quota}}[c]}$ .  $\text{canConsume}(\epsilon_x^t)$  to return **FALSE** by definition, so no budget is spent at all. In such a case,

$$\epsilon_{\text{global}}^{\leq t, \text{bad}} = \epsilon_{\text{used}}^{\text{bad}} + 0 = \epsilon_{\text{used}}^{\text{bad}} \leq N^{\leq t-1, \text{adv}} \cdot \epsilon_{\text{conv-quota}}, \quad (66)$$

by equation (63).

Since the adversary has created at most  $N^{\text{adv}}$  by the end of time  $t$ , it must be the case that  $N^{\leq t-1, \text{adv}} \leq N^{\leq t, \text{adv}} \leq N^{\text{adv}}$ . This means that, in either case, the attackers can consume at most  $N^{\leq t, \text{adv}} \epsilon_{\text{conv-quota}} \leq N^{\text{adv}} \epsilon_{\text{conv-quota}}$  of the global filter budget by the end of time  $t$ , as desired  $\square$

*Proof (second upper bound by  $M^{\text{adv}}$ ).* By basic composition under a pure DP filter, we know  $\epsilon_{\text{global}}$  is by definition the sum of global filter consumption at each time up to the end

of time  $t-1$ :

$$\epsilon_{\text{global}}^{\leq t-1, \text{bad}} = \sum_{k \in [t-1]: c_k \in \text{bad}_c} \epsilon_x^k \cdot \text{pass}(k) \quad (67)$$

$$= \sum_{k \in [t-1]: c_k \in \text{bad}_c} \Delta_x \rho_{b_k}^k \cdot \text{pass}(k), \quad (68)$$

where  $\Delta_x \rho_{b_k}^k$  is how much the attribution report generated for query  $k$  from beneficiary  $b_k$  can change with respect to  $x$ , which we substitute next. Let  $\vec{x} = (x_{i_1}, \dots, x_{i_m})$  be the vector of  $|i_k| = m$  device-epoch-sites, where  $x_i = (d, e, F_i)$  has all its events on site  $i$ . Then, we let the corresponding neighboring dataset of  $D \in \mathbb{D}$  be  $D + \vec{x}$ , so that:

$$\rho(D + \vec{x}) - \rho(D) = \sum_{j=1}^m \rho(D + x_{i_j}) - \rho(D + x_{i_{j-1}}), \quad (69)$$

where  $x_{i_0} = 0$ . Thus, we substitute by the following decomposition of a query's sensitivity across the impression sites relevant to that query, by how it is assigned in Alg. 3 Line 16:

$$\Delta_x \rho_{b_k}^k = \max_{D, D' \in \mathcal{D}: D' = D + x} \|\rho_{b_k}^k(D') - \rho_{b_k}^k(D)\|_1 \quad (70)$$

$$= \max_{D \in \mathcal{D}} \|\rho_{b_k}^k(D + x) - \rho_{b_k}^k(D)\|_1 \quad (71)$$

$$\leq \max_{D \in \mathcal{D}} \sum_{j \in [m]: i_j \in \text{bad}_i} \|\rho_{b_k}^k(D + x_1 + \dots + x_{i_j}) - \rho_{b_k}^k(D + x_1 + \dots + x_{i_{j-1}})\|_1 \quad (72)$$

$$\leq \sum_{j \in [m]: i_j \in \text{bad}_i} \max_{D \in \mathcal{D}} \|\rho_{b_k}^k(\hat{D} + x_{i_j}) - \rho_{b_k}^k(\hat{D})\|_1 \quad (73)$$

$$= \sum_{j \in [m]: i_j \in \text{bad}_i} \Delta_x^j \rho_{b_k}^k = \sum_{i \in \text{bad}_i} \Delta_x^i \rho_{b_k}^k \quad (74)$$

$$= \sum_{i \in \text{bad}_i} \epsilon_x^{i,k}[i], \quad (75)$$

$$= \sum_{i \in \text{bad}_i} \epsilon_x^{i,k}[i], \quad (76)$$

where the inequality in equation 72 is because of the following. First, by the restriction in the sum, we know  $k$  satisfies  $c_k \in \text{bad}_c$ . Second, recall that, for a conversion site to incur privacy losses on impression sites, the conversion site must register these impression sites, meaning that if  $c_k \in \text{bad}_c$ , then  $i_k \subseteq \text{bad}_i$ . Now, plug the the inequality  $\Delta_x \rho_{b_k}^k \leq \sum_{i \in \text{bad}_i} \epsilon_x^{i,k}[i]$  in, we get

$$\epsilon_{\text{global}}^{\leq t-1, \text{bad}} \leq \sum_{k \in [t-1]: c_k \in \text{bad}_c} \sum_{i \in \text{bad}_i} \epsilon_x^{i,k}[i] \cdot \text{pass}(k) \quad (77)$$

$$= \sum_{i \in \text{bad}_i} \sum_{k \in [t-1]: c_k \in \text{bad}_c, i \in i_k} \epsilon_x^{i,k}[i] \cdot \text{pass}(k), \quad (78)$$

$$\leq \sum_{i \in \text{bad}_i} \sum_{k \in [t-1]: i \in i_k} \epsilon_x^{i,k}[i] \cdot \text{pass}(k), \quad (79)$$

by changing order of summation, and the last inequality by relaxing the " $c_k \in \text{bad}_c$ " condition. But note that

$$\epsilon_{\text{imp-quota}}^{\leq t-1}[i] = \sum_{k \in [t-1]: i \in i_k} \epsilon_x^{i,k}[i] \cdot \text{pass}(k), \quad (80)$$

because, by epoch-site-level consistency property in lemma 1, we know that only relevant site  $i$  at time  $k$ , where every filter has enough budget to pass the 2-PC check, will have epoch-site level privacy losses incurred. Substituting this equality into equation (79), we get:

$$\epsilon_{\text{global}}^{\leq t-1, \text{bad}} \leq \sum_{i \in \text{bad}_i} \epsilon_{\text{imp-quota}}^{\leq t-1} [i] \quad (81)$$

$$= \sum_{i \in \text{bad}_i: \epsilon_{\text{imp-quota}}^{\leq t-1} [i] > 0} \epsilon_{\text{imp-quota}}^{\leq t-1} [i] \quad (82)$$

$$\leq \sum_{i \in \text{bad}_i: \epsilon_{\text{imp-quota}}^{\leq t-1} [i] > 0} \epsilon_{\text{imp-quota}} \quad (83)$$

$$= \left| \left\{ i \in \text{bad}_i : \epsilon_{\text{imp-quota}}^{\leq t-1} [i] > 0 \right\} \right| \cdot \epsilon_{\text{imp-quota}}, \quad (84)$$

because only non-zero privacy losses that were incurred contribute meaningfully to the composition. Finally, we note that  $\left| \left\{ i \in \text{bad}_i : \epsilon_{\text{imp-quota}}^{\leq t-1} [i] > 0 \right\} \right| \leq M^{\leq t-1, \text{adv}}$  by definition and:

$$\epsilon_{\text{global}}^{\leq t-1, \text{bad}} \leq M^{\leq t-1, \text{adv}} \cdot \epsilon_{\text{imp-quota}}. \quad (85)$$

Following this result, similar to the proof for part 1:

- Suppose  $\epsilon_x^t \leq \epsilon_{\text{imp-quota}}$ ,

$$\epsilon_{\text{global}}^{\leq t, \text{bad}} \leq M^{\leq t, \text{adv}} \cdot \epsilon_{\text{imp-quota}}. \quad (86)$$

- Else,  $\epsilon_x^t > \epsilon_{\text{imp-quota}}$ , then  $\epsilon_{\text{imp-quota}}$  will be exceeded, causing `canConsume` to return `FALSE`, so,

$$\epsilon_{\text{global}}^{\leq t, \text{bad}} = \epsilon_{\text{global}}^{\leq t-1, \text{bad}} + 0 = \epsilon_{\text{global}}^{\leq t-1, \text{bad}} \leq M^{\leq t-1, \text{adv}} \cdot \epsilon_{\text{imp-quota}}, \quad (87)$$

by equation (85).

Since by the end of time  $t$ , the adversary has created at most  $M^{\text{adv}}_{\text{imp-quota}}$  filters, we know  $M^{\leq t-1, \text{adv}} \leq M^{\leq t, \text{adv}} \leq M^{\text{adv}}$ , which means that in both cases we have:

$$\epsilon_{\text{global}}^{\leq t, \text{bad}} \leq M^{\leq t, \text{adv}} \epsilon_{\text{imp-quota}} \leq M^{\text{adv}} \epsilon_{\text{imp-quota}}, \quad (88)$$

as desired.  $\square$

Finally, we use the preceding lemma to show the overall DoS depletion resilience guarantee over Gordon's lifetime. Consider an execution of Alg. 2. A report at time step  $k$  concerns with *one* conversion site  $c_k$ , and some subset of impression sites  $i_k \subseteq S$ . The union of attackers, in particular, can control an arbitrary subset of conversion sites  $\text{bad}_c \subseteq S$ , which may or may not contain  $c_k$  at a given  $k$ . We denote by  $N^{\text{adv}}$  the size of  $|\text{bad}_c|$  over the entire lifetime. Similarly, the adversary can control an arbitrary subset of impression sites  $\text{bad}_i \subseteq S$ , which may or may not intersect with  $i_k$ . We denote by  $N^{\text{adv}}$  the size of  $|\text{bad}_i|$  over the entire lifetime. We let  $\text{bad} = \text{bad}_c \cup \text{bad}_i$  and  $\text{good} = S \setminus \text{bad}$ .

**Theorem 1 (Resilience to DoS depletion).** *Consider an adversary who manages to create  $M^{\text{adv}}$  and  $N^{\text{adv}}_{\text{imp-quota}}$  and  $\text{conv-quota}$  filters, respectively. The maximum budget  $\epsilon_{\text{global}}^{\text{adv}}$  that the adversary can consume from the global filter*

*on a device  $d$  is such that:*

$$\epsilon_{\text{global}}^{\text{adv}} \leq \min(M^{\text{adv}} \epsilon_{\text{imp-quota}}, N^{\text{adv}} \epsilon_{\text{conv-quota}}).$$

*Proof.* The result follows directly from Lemma 2. At any time  $t$ , if the adversary controls at most  $M^{\leq t, \text{adv}} \leq M^{\text{adv}}_{\text{imp-quota}}$  and  $N^{\leq t, \text{adv}} \leq N^{\text{adv}}_{\text{conv-quota}}$  filters, then by Lemma 2:

$$\epsilon_{\text{global}}^{\text{adv}} \leq M^{\leq t, \text{adv}} \epsilon_{\text{imp-quota}} \leq M^{\text{adv}} \epsilon_{\text{imp-quota}} \quad (89)$$

$$\epsilon_{\text{global}}^{\text{adv}} \leq N^{\leq t, \text{adv}} \epsilon_{\text{conv-quota}} \leq N^{\text{adv}} \epsilon_{\text{conv-quota}}. \quad (90)$$

Therefore:

$$\epsilon_{\text{global}}^{\text{adv}} \leq \min(M^{\text{adv}} \epsilon_{\text{imp-quota}}, N^{\text{adv}} \epsilon_{\text{conv-quota}}).$$

$\square$

## C Batched Algorithm to Improve Utilization

### C.1 Algorithm

Alg. 6 describes the batched algorithm on a single device. Instead of executing `GenerateReport` as soon as a request comes, as in Alg. 3, requests are accumulated in a batch. Each epoch is divided into  $k$  scheduling intervals, and since a request can request older epochs (up to a maximum attribution window length, *a.k.a.* data lifetime) we release budget progressively over  $T$  intervals. For instance, if requests have attribution window of at most 2 epochs, we can divide this data lifetime into  $T = 4$  releases, with  $k = 2$  releases happening inside each epoch. We can also do  $T = 2$  releases with  $k = 1$  interval per epoch.

Budget release and unlocked budget semantics are defined as in [5].  $\mathcal{F}_{\text{global}}^{\text{unlock}}$  becomes a no-op after  $T$  releases, when the unlocked budget reaches the filter capacity  $\epsilon_{\text{global}}$ .

We define  $\mathcal{A}, \mathcal{U} \leftarrow \text{TryAllocate}(\mathcal{R})$  as follows. `TryAllocate` takes a set of report requests  $\mathcal{R}$ . For each request, it executes a heuristic that estimates whether Alg. 3's `GenerateReport` will succeed the request (*i.e.*, the whether the filters will return `TRUE` at Line 19). It then calls `GenerateReports` on the requests that were predicted to be allocatable, and returns two sets:  $\mathcal{A}$  the reports for requests that were executed, and  $\mathcal{U}$  the remaining unallocated requests. `TryAllocateOne` behaves like `TryAllocate`, except that it stops after the first executed request. Our heuristic is a variant of `GenerateReport` that relies purely on public information — using proxy filters  $\tilde{\mathcal{F}}$  with IDP optimizations turned off — allowing the scheduler to make decisions without leaking privacy across epochs.

`SortBatch` attaches a weight  $(b_r, \epsilon_r)$  to each request  $r$  in a batch, and then sorts by smallest weight first (in lexicographic order). The weights are defined as follows, using the proxy filters  $\tilde{\mathcal{F}}$  define for `TryAllocate`.  $\epsilon_r$  is the global epsilon requested by  $r$  (either available as a request parameter, or computed as  $\epsilon_r = \Delta \rho_r / \sigma$  for a Laplace noise scale  $\sigma$ ).  $b_r$  is the smallest budget consumed by any impression site  $i \in i_r$  requested by  $r$ , summed over the set  $E$  of all the epochs



**Algorithm 6** Batched Algorithm (On-Device)**Input:**


---

```

1:  $\epsilon_{\text{global}}, \epsilon_{\text{per-site}}, \epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}}$ : same parameters as Alg. 3.
2:  $k$ : number of scheduling intervals per epoch.
3:  $T$ : number of scheduling intervals to release the full budget.

```

---

```

4: function MAIN
5:   for  $e \in \mathbb{N}$  do
6:     // Initialize new epoch with its own filters
7:      $\mathcal{F}_e \leftarrow \text{InitializeFilters}(\epsilon_{\text{global}}, \epsilon_{\text{per-site}}, \epsilon_{\text{imp-quota}}, \epsilon_{\text{conv-quota}})$ 
8:     // Initially no global budget available
9:      $\mathcal{F}_e^{\text{global}}.\text{unlocked} \leftarrow 0$ 
10:     $\mathcal{R}_{\text{batch}} \leftarrow \emptyset$  // Requests for the batch phase
11:    for  $t \in [k]$  do
12:       $\mathcal{R}_{\text{new}} \leftarrow \text{ReceiveNewRequests}()$ 
13:       $\mathcal{A}, \mathcal{R}_{\text{batch}} \leftarrow \text{ScheduleBatch}(\mathcal{R}_{\text{new}}, \mathcal{R}_{\text{batch}})$ 
14:       $\text{SendReportsForRelease}(\mathcal{A})$ 
15:
16: function SCHEDULEBATCH( $\mathcal{R}_{\text{new}}, \mathcal{R}_{\text{batch}}$ )
17:   // 1. Initialization phase
18:   for  $e' \in [e]$  do
19:      $\mathcal{F}_e^{\text{global}}.\text{unlocked} \leftarrow \mathcal{F}_e^{\text{global}}.\text{unlocked} + \epsilon_{\text{global}}/T$ 
20:     for  $i \in \mathcal{S}$  do
21:       // Reinitialize impression-site quota.
22:        $\mathcal{F}_e^{\text{imp-quota}}[i] \leftarrow \text{InitializeFilter}(\epsilon_{\text{imp-quota}}/T)$ 
23:        $\mathcal{A}_{\text{init}}, \mathcal{U}_{\text{init}} \leftarrow \text{TryAllocate}(\mathcal{R}_{\text{batch}})$ 
24:        $\mathcal{A} \leftarrow \mathcal{A}_{\text{init}}$ 
25:   // 2. Online phase
26:    $a_{\text{online}}, u_{\text{online}} \leftarrow \text{TryAllocate}(\mathcal{R}_{\text{new}})$ 
27:    $\mathcal{A} \leftarrow \mathcal{A} \cup \mathcal{A}_{\text{online}}$ 
28:   // 3. Batch phase
29:   for  $e' \in [e], i \in \mathcal{S}$  do
30:     // Deactivate impression-site quota.
31:      $\mathcal{F}_e^{\text{imp-quota}}[i].\text{capacity} \leftarrow \infty$ 
32:    $\text{batch} \leftarrow \mathcal{U}_{\text{init}} \cup \mathcal{U}_{\text{online}}$ 
33:   do
34:      $\text{sorted} \leftarrow \text{SortBatch}(\text{batch})$ 
35:      $(a, u) \leftarrow \text{TryAllocateOne}(\text{sorted})$ 
36:      $\text{batch} \leftarrow u$ 
37:      $\mathcal{A} \leftarrow \mathcal{A} \cup a$ 
38:   while  $\mathcal{A} \neq \emptyset$ 
39:   return  $\mathcal{A}, \text{batch}$ 

```

---

considered in the batch:

$$b_r := \min_{i \in \mathcal{I}_r} \sum_{e \in E} \tilde{\mathcal{F}}^{\text{imp-quota}}[i].\text{consumed} \quad (91)$$

Pierre: This min comes from our first idea, where we compute the “total consumed budget” for each impression site,

and then start allocating requests site by site, starting with the site that received the least budget so far. The motivation was to try to maximize the minimum allocation across impression sites. But it sounds very easy to game this metric: if a request also touches a fresh impression site that received zero budget so far, then it will be able to run first. So an impression site with large allocation can piggyback on an impression site with small allocation to receive even more budget. What if we use a max instead? It has a nice property in that a site is not incentivized to request for impression sites it doesn’t need.

If all the requests have the same single-epoch attribution window  $e$ , this gives:

$$b_r := \max_{i \in \mathcal{I}_r} \tilde{\mathcal{F}}_e^{\text{imp-quota}}[i].\text{consumed} \quad (92)$$

Potentially useful fact for later: if we have a proxy that is always approximating consumption, for any  $i$  we have:

$$b_r \geq \mathcal{F}_e^{\text{imp-quota}}[i].\text{consumed} \quad (93)$$

Suppose you have an adversary that asked for a lot of budget  $\epsilon_0$  already for a malicious impression site  $i_0$ . Under what conditions can they submit another query that runs before honest queries? If the honest queries only request impression sites  $i$  that have been requested with less budget  $\epsilon_i \leq \epsilon_0$  than  $i_0$ , the honest queries will be scheduled first. But they might still not succeed if other filters fail... {Mathias: my issue with max is that an adversary can block other requests by consuming from a very common impression site, and then since everyone lists it they are all blocked, and by avoiding it the attacker can run first. I had that in mind and came to check your order, and liked the min :). I guess it has the opposite bad effect you describe though...}

Finally, SendReportsForRelease prepares the reports from allocated requests to be sent at the right time, depending on the duration specified by each request.

## C.2 DoS resilience under batching

Under our mixed online and batch algorithm, the effort in terms of on device user interactions  $U^{\text{adv}}$  required from an adversary to consume global filter budget depends on the overall workload of the system. As we saw in §4.3, in the worst case (e.g., when there are no legitimate queries in the system) this can lead to a weaker upper-bound on the budget consumption by an adversary compared to Thm. 2, with the following result:

$$\epsilon_{\text{global}}^{\text{adv}} \leq (1 + r)\epsilon_{\text{per-site}} \times \text{quota-count} \times (U^{\text{adv}} - 1). \quad (94)$$

Intuitively, this is because the attacker can batch conversion queries that all request the same impression: user interactions are only needed to create one impression under the adversary’s control, as well as  $U^{\text{adv}} - 1$  conversions that can be used to deplete global filter budget when the imp-quota filters are disabled (l.9 in Algorithm 1, ll. 29-31 in Algorithm 6) {Mathias: TODO proper algo refs}.

In practice however, we expect the benign workload to contain online queries (configured to return instantly, with no batching). To deny service to those queries requires a higher number of on device interactions for the adversary, so it is relevant to ask for a lower-bound on user interactions  $U^{\text{adv}}$  required by an adversary to prevent  $K$  legitimate online queries from being allocated. Consider  $K$  legitimate queries, with total requested budget summing to  $\epsilon_{\text{global}}^{\text{good}}$ . **{Mathias: WIP but need to go to a meeting}**

In addition, during the batch allocation phase **{Mathias: TODO II. in algos}**, the adversary would still need overcome the scheduler's sorting mechanism, to be scheduled before waiting legitimate requests. Since the sorting mechanism favors low-budget and underrepresented impression sites, the adversary would likely require more than one user interaction ( $U^{\text{adv}} \gg 1$ ) to mount an attack, making Equation 94 pessimistic. The time dynamics and workload dependency of the batching phases make the analysis of such guarantees challenging though, and we leave a proper formal treatment of any guarantees related to sorting for future work.

### C.3 Pareto efficiency proof

We prove the batched algorithm's pareto efficiency property under a restricted setting. First, we only consider requests that request budget from a single device epoch. Second, we prove pareto efficiency of the batched phase of the Gordon batched algorithm (Algorithm 6), which operates on the global filter. We do not prove Pareto efficiency for Gordon as a whole, which also manages per-site filters and quotas. Similar to prior work [5], we also assume all-or-nothing utility for DP requests, wherein their utility is equal to one if their entire global filter requested budget is granted, and is equal to zero if it is not.

**Theorem 3 [Pareto efficiency.]** *The batched global filter scheduling algorithm is Pareto efficient with regards to the global filter budget within a device epoch.*

*Proof.* Consider the batch phase with requests that request budget from a single device-epoch. We prove that no allocation from global budget can increase a request's utility without decreasing another request's utility. Consider request  $i$  that is in the queue of the batch phase. If  $i$  is allocated by Algorithm 6, its utility is one and cannot be further improved. If  $i$  has not yet been scheduled after the batch stage, this means it cannot be allocated from any released global filter budget, as Algorithm 6 grants requests until no requests can be allocated. Allocating request  $i$  would require extra global filter budget, which can only come from another allocated request. Since each allocated request has exactly its requested global filter budget this would decrease its utility from one to zero, which is not Pareto-improving.  $\square$

## D Cross-report Privacy Loss Optimization

### D.1 Example from Fig. 2

We illustrate Gordon's *cross-report* privacy loss optimization using our running example. This optimization is orthogonal to Cookie Monster's *per-report* individual-DP-based strategies (§2.3), and instead leverage structure *across* reports, often requested by different intermediaries for the same conversion.

In Fig. 2, *r1.ex*, *r2.ex*, and *r3.ex* request single-advertiser reports from `attributionObject` (a), all on behalf of client *shoes.ex*; separately, *r1.ex* and *r2.ex* request cross-advertiser reports from (b) for their own purposes. All five reports operate on the same attribution histogram, assigning \$30 to each of two impressions (epochs  $e_1$ ,  $e_2$ ). Cookie Monster computes a base epoch-level privacy loss of 0.3 per report (§2.3). Naïvely, one would expect a cumulative deduction of 0.9 from *shoes.ex*'s filters (three reports) and 1.5 from the global filter (five reports). Yet the Privacy Filters table shows only deductions of 0.6 and 0.9, respectively.

The discrepancy arises because some reports *shard* the histogram into non-overlapping pieces—enabling sensitivity-based optimizations. *r1.ex* and *r2.ex*'s single-advertiser reports from (a) each include a disjoint portion:  $\{1:30\}$  and  $\{2:30\}$ , respectively. Since both are funded by the same per-site filter (of *shoes.ex*), their combined release leaks no more than a single full histogram toward *shoes.ex*, incurring only 0.3 privacy loss. They likewise count as one deduction against the shared global filter. In contrast, *r3.ex*'s report includes the full histogram (to give *shoes.ex* a complete view across intermediaries; see §3.2), overlapping with both *r1.ex* and *r2.ex* and adding another 0.3 of loss to both *shoes.ex*'s filter and the global filter. A similar optimization applies to cross-advertiser reports from (b). These are funded from separate filters (those of *r1.ex* and *r2.ex*), so each incurs 0.3 loss. But against the global filter, they again count as one, bringing the total global filter deduction to 0.9 instead of the unoptimized 1.5.

We next formalize this optimization, whose logic we encapsulate in the `attributionObject`. This object dynamically optimizes budget deduction across the per-site, global, and quota filters on each `getReport()` call, on the basis of prior invocations and deductions.

### D.2 Privacy proof

**TODO(Alison): Please capitalize names of things you are referencing: "proposition 1" becomes "Proposition 1" (it's viewed as its name).** **Pierre: This is automated if you use macros like \Thm, Prop, etc. And actually let's be consistent across appendix sections, we've been using Lem. instead of Prop in previous sections.**

We formalize the cross-report optimization with a series of definitions, after which we analyze its sensitivity properties through a set of intermediary propositions that ultimately lead to our main correctness result in Thm. 11.

**Definition 7** (Histogram attribution function). *Let  $A$  be an*

attribution function  $A : \mathcal{P}(I)^k \rightarrow \mathbb{R}^m$ , for a fixed  $k$  and  $m$ .  $A$  is a histogram attribution function, if for any vector of a set of impressions  $\mathbf{F}$ , the following conditions all hold:

1. There exists  $a_F : I \rightarrow \mathbb{R}_+$  a function that attributes a value to each impression  $f \in I$  depending on all the other impressions  $\mathbf{F}$ .
2. There exists a one-hot encoding function  $H$  that maps each event  $f$  to one of  $m$  buckets. That is,  $H : I \rightarrow \{0, 1\}^m$  such that  $\forall f \in I, \|H(f)\|_1 = 1$ .
3.  $A(\mathbf{F}) = \sum_{f \in \mathbf{F}} a_F(f) \cdot H(f)$

**Definition 8** (Query partitions for beneficiary sites). Let  $A$  be a histogram attribution function  $A : \mathcal{P}(I)^k \rightarrow \mathbb{R}^m$ , for a fixed  $k$  and  $m$  and with an associated one-hot encoding function  $H : \mathbb{R}^m \rightarrow \{0, 1\}^m$ . Fix  $n \in \mathbb{N}_+$ . Let  $\mathbf{F}$  be a fixed vector of event sets. We define the following:

- $H_j$  is an encoding partition of  $H$  for some  $j \in [n]$  if  $H_j : I \rightarrow \{0, 1\}^m$  and for some subset of impressions  $S_j \subseteq I$  we have that

$$H_j(f) = \begin{cases} H(f) & \text{if } f \in S_j \\ 0 & \text{if } f \in I \setminus S_j \end{cases}$$

- *Pierre: We don't need to support any possible isomorphism, so I'd define  $\Psi$  once and for all, and call it "the concatenation isomorphism, or just "the concatenation map" (we'll use the isomorphism later but it sounds complicated when we mention this upfront).*  $\Psi : \{0, 1\}^m \times \dots \times \{0, 1\}^m \rightarrow \{0, 1\}^m$  is a concatenation map for a set of encoding partitions  $H_j$  of  $H$  if

$$\Psi(H_1(f), H_2(f), \dots, H_n(f)) = H(f) \quad (95)$$

for all  $f \in I$  and we have that  $\sum_{j \in [n]} H_j(f) = H(f)$

- A set of encoding partitions  $H_1, H_2, \dots, H_n$  is a valid partition set of encoding  $H$  if there exists a concatenation isomorphism  $\Psi$  such that  $\Psi(H_1, H_2, \dots, H_n) = H$  and for all  $f \in I$  we have that  $\|H_j(f)\|_1 = 1$  for exactly one  $j \in [n]$ .
- $A_j$  is the attribution partition of  $A$  for some  $j \in [n]$  where

$$A_j(\mathbf{F}) := \sum_{f \in \mathbf{F}} a_F(f) \cdot H_j(f) \quad (96)$$

- $\rho^j$  is the report partition for  $j \in [n]$  for a fixed report identifier  $r \in \mathbb{Z}$  where

$$\rho_r^j(D) = A_j(D_r^E) \quad (97)$$

Let  $H_1, H_2, \dots, H_n$  be a valid encoding partition of encoding  $H$ . Let  $B = \{b_1, \dots, b_n\}$  be a set of beneficiary sites, where each  $b_j$  has a corresponding encoding partition  $H_j$ . Consider a set of report identifiers  $R$ . We define  $Q_j$  as the query partition that results by using the attribution partition  $A_j$  and the

corresponding report partitions  $\rho^j$ . We have:

$$Q_j(D) := \sum_{r \in R} \rho_r^j(D) \quad (98)$$

**Proposition 1** (Recovering encoding partitions from a concatenated encoding). Let  $H_1, H_2, \dots, H_n$  be a valid encoding partition of encoding  $H$ , with concatenation map  $\Psi$ . By definition, each  $H_j$  has a corresponding subset of impressions  $S_j \subseteq I$ . Given  $S_1, S_2, \dots, S_n$  and  $H$ , we can recover their corresponding encoding partitions and

$$\Psi^{-1}(H) = (H_1, H_2, \dots, H_n) \quad (99)$$

*Pierre: So here we need one more step saying that given noisy results for  $Q$ , we can recover noisy results for each  $Q_i$ , when the noise is the multidimensional Laplace mechanism with some parameter  $\text{Lap } b$ . This is where we need  $\Psi$  to be an isomorphism and not just a bijection: we want the additive noise to go through (would be nice to be a bit extra verbose here about how the noise also gets decomposed into coordinates). Then, add another point saying that the mechanism  $\mathcal{M}_1, \dots, \mathcal{M}_n$  that jointly releases the  $n$  queries can be obtained by post-processing  $\mathcal{M}$  with  $\Psi^{-1}$ . And by post-processing property of DP, if  $\mathcal{M}$  is  $\epsilon$ -DP, then  $\Psi^{-1}(\mathcal{M}) = \mathcal{M}_1, \dots, \mathcal{M}_n$  is also  $\epsilon$ -DP.*

*Proof.* By definition, each  $H_j$  can be constructed by outputting  $H(f)$  for each  $f \in S_j$  and 0 for all  $f \in I \setminus S_j$ .  $\square$

**Proposition 2** (Individual sensitivity of a report partition). Fix a report identifier  $r$ , a device  $d_r$ , a set of epochs  $E_r$ , a beneficiary  $b_j \in B$ , an attribution partition function  $A_j$  with encoding partition  $H_j$ , impression set  $S_j$  and the corresponding report partition  $\rho^j : D \rightarrow A_j(D_{d_r}^{E_r})$ . We define

$$A^{\max} := \max_{\mathbf{F} \in \mathcal{P}(I)^k} \sum_{i=1}^k \sum_{f \in \mathbf{F}} a_F(f) \cdot H(f) \quad (100)$$

For a fixed device-epoch record  $x = (d, e, F) \in \mathcal{X}$ , we have that the individual sensitivity of  $\rho_j$  is

$$\Delta_x(\rho^j) \leq \begin{cases} 0 & \text{if } d \neq d_r, e \notin E_r \text{ or } F \cap S_j = \emptyset \\ \|A_j(F)\|_1 & \text{if } d = d_r \text{ and } E_r = \{e\} \\ 2A^{\max} & \text{if } d = d_r, e \in E_r \text{ and } F \cap S_j \neq \emptyset \end{cases}$$

*Proof.* Follows directly from theorem 4 of [8], with the upper bound on histogram report sensitivity from theorem 18 of [8].  $\square$

**Proposition 3** (Individual sensitivity of a concatenated report). Fix a report identifier  $r$ , a device  $d$ , a set of epochs  $E_r$ , a histogram attribution function  $A$  and a set of beneficiaries  $B = \{b_1, b_2, \dots, b_n\}$ . Suppose each beneficiary has a corresponding report partition that results from histogram attribution function  $A$  with valid encoding partitions  $H_1, H_2, \dots, H_n$  (i.e.  $\rho_j$  is the report partition for beneficiary  $b_j$  with encoding partition  $H_j$ ). For a fixed device-epoch record  $x = (d, e, F) \in \mathcal{X}$ , the

individual sensitivity of the concatenated report  $\rho$  is

$$\Delta_x(\rho) \leq \begin{cases} 0 & \text{if } d \neq d_r, e \notin E_r \\ \|A(F)\|_1 & \text{if } d = d_r \text{ and } E_r = \{e\} \\ 2A^{\max} & \text{if } d = d_r, e \in E_r \end{cases}$$

*Proof.* Follows directly from theorem 4 of [8], with the upper bound on histogram report sensitivity from theorem 18 of [8].  $\square$

**Proposition 4** (Individual sensitivity of query partitions  $Q_j$ ). Fix an impression set  $S_j \subseteq \mathcal{I}$ . Let  $Q_j$  be the query partition of beneficiary site  $j \in [n]$  with partition encoding  $H_j$ . Fix a device-epoch record  $x = (d, e, F) \in \mathcal{X}$ .

$$\Delta_x(Q_j) \leq \sum_{r \in R} \Delta_x(\rho_r^j) \quad (101)$$

*Proof.* Take  $D, D' \in \mathcal{D}$  such that  $D' = D + x$ . We have that

$$\Delta_x(Q_j) = \max_{D'=D+x \in \mathcal{D}} \|Q_j(D) - Q_j(D')\|_1 \quad (102)$$

$$= \max_{D'=D+x \in \mathcal{D}} \left\| \sum_{r \in R} \rho_r^j(D) - \rho_r^j(D') \right\|_1 \quad (103)$$

$$\leq \sum_{r \in R} \max_{D'=D+x \in \mathcal{D}} \left\| \rho_r^j(D) - \rho_r^j(D') \right\|_1 \quad (104)$$

$$\leq \sum_{r \in R} \Delta_x(\rho_r^j) \quad (105)$$

by the definition of individual sensitivity and the triangle inequality.  $\square$

**Proposition 5** (Sensitivity of concatenated query  $Q$ ). Let  $H_1, H_2, \dots, H_n$  be a valid encoding partition of  $H$ . Let  $B = \{b_1, \dots, b_n\}$  be a set of beneficiary sites, where each  $b_j$  has a corresponding encoding partition  $H_j$ . Consider a set of report identifiers  $R$ . Let each beneficiary site  $b_j$  also have a corresponding query partition  $Q_j$  and the corresponding report partitions  $\rho_r^j$  for  $r \in R$  and attribution partition  $A_j$ . Let  $Q = Q_1, Q_2, \dots, Q_n$  be the query that results from jointly processing all queries of each beneficiary site. Fix a device-epoch record  $x = (d, e, F) \in \mathcal{X}$ . We have that the sensitivity of  $Q$  is such that

$$\Delta_x(Q) \leq \sum_{r \in R} \Delta_x(\rho_r) \quad (106)$$

*Proof.* We first notice that for a fixed  $r$  we have that  $\sum_{j=1}^n \rho_r^j(D) = \rho_r(D)$ . This follows directly from our definition of a valid encoding partition where  $H_j(f) = H(f)$  for all  $f \in F$  and for

all  $j \in [n]$ .

$$\sum_{j=1}^n \rho_r^j(D) = \sum_{j=1}^n A_j(D_d^E) \quad (107)$$

$$= \sum_{j=1}^n \sum_{f \in F} a_F(f) \cdot H_j(f) \quad (108)$$

$$= \sum_{f \in F} a_F(f) \cdot H(f) \quad (109)$$

$$= \rho_r(D) \quad (110)$$

Now, take  $D, D' \in \mathcal{D}$ , such that  $D' = D + x$ . We have that

$$\Delta_x(Q) = \max_{D'=D+x \in \mathcal{D}} \|Q(D) - Q(D')\|_1 \quad (111)$$

$$= \max_{D'=D+x \in \mathcal{D}} \left\| \sum_{j=1}^n Q_j(D) - Q_j(D') \right\|_1 \quad (112)$$

$$= \max_{D'=D+x \in \mathcal{D}} \left\| \sum_{j=1}^n \sum_{r \in R} \rho_r^j(D) - \rho_r^j(D') \right\|_1 \quad (113)$$

$$= \max_{D'=D+x \in \mathcal{D}} \left\| \sum_{r \in R} \sum_{j=1}^n \rho_r^j(D) - \sum_{j=1}^n \rho_r^j(D') \right\|_1 \quad (114)$$

$$= \max_{D'=D+x \in \mathcal{D}} \left\| \sum_{r \in R} \rho_r(D) - \rho_r(D') \right\|_1 \quad (115)$$

$$\leq \sum_{r \in R} \max_{D'=D+x \in \mathcal{D}} \left\| \rho_r(D) - \rho_r(D') \right\|_1 \quad (116)$$

$$\leq \sum_{r \in R} \Delta_x(\rho_r) \quad (117)$$

$\square$

**Theorem 11** (Cross-report privacy loss optimization). Consider the setting described in Proposition 5. *Pierre: Re-emphasize that we are considering the specific class of queries defined in Def. .... Also give the actual value in the theorem instead of just saying that it's "less". Processing query partitions  $Q_1, Q_2, \dots, Q_n$  jointly deducts less global filter budget than processing them individually. Pierre: Here we need to talk about DP mechanisms  $\mathcal{M}_1$  instead of the underlying queries. Also would be good to be consistent, I don't think we use "processing" too much elsewhere, maybe we've been saying "answering" a query, to check with other propositions and algorithms.*

*Proof.* Suppose queries partitions  $Q_1, Q_2, \dots, Q_n$  were processed sequentially in isolation by algorithm 3. The global filter deduction would use the individual sensitivity of each query partition as defined in proposition 4. Given a standard deviation  $\sigma$  of the Laplace distribution  $\mathcal{L}$ , each query  $Q_j$  would deduct from the budget:

$$\epsilon_j = \sum_{r \in R} \Delta_x(\rho_r) \cdot \frac{\sqrt{2}}{\sigma} \quad (118)$$

~~In total~~By sequential composition, processing all queries would deduct the following. We use proposition 2 to show an

upper bound.

$$\epsilon_{Q_1, \dots, Q_n} = \sum_{j \in [n]} \epsilon_j \cdot \frac{\sqrt{2}}{\sigma} \quad (119)$$

$$\leq n \sum_{r \in R} 2A^{\max} \cdot \frac{\sqrt{2}}{\sigma} \quad (120)$$

$$\leq 2n|R|A^{\max} \cdot \frac{\sqrt{2}}{\sigma} \quad (121)$$

Suppose now that query partitions are first jointly processed all query partitions as one query  $Q$ . By proposition 1, we can then recover each query partition  $Q_1, Q_2, \dots, Q_n$  and distribute them among the beneficiary sites. **Pierre:** Would be good to plug the beefed up version of Prop 1 with the formal post-processing argument I sketched. **Alison:** Is it clear how you can do this? Should I provide more details?. **Roxana:** I don't know, didn't get to read. But giving all the details on a thing that we do that's new is important. The global filter deduction would use individual sensitivity concatenated query

sensitivity as defined in proposition 5 as follows:

$$\epsilon_Q = \sum_{r \in R} \cdot \frac{\sqrt{2}}{\sigma} \Delta_x(\rho_r) \quad (122)$$

By proposition 3, we can bound this as follows

$$\epsilon_Q \leq 2|R|A^{\max} \cdot \frac{\sqrt{2}}{\sigma} \quad (123)$$

It clearly follows that  $\epsilon_{Q_1, \dots, Q_n} \leq \epsilon_Q$ . □

## E Prototype Screenshot

**TODO(Nikos/Giorgio):** Add here the screenshot. Describe it in text – scenario, what it shows etc.

## F Evaluation Methodology Details

§6.1 gives a high-level overview of our methodology, covering the necessary details for result assessment. Here, we give further details that may be necessary for reproduction. We also intend to release our evaluation repository as part of artifact release.