

Optimal Universal Black-Box Privacy Preservation and Tight Adversarially-Adaptive Composition for Entropic Secret

Abstract. Privacy entails a fundamental tradeoff between information propagation and preservation, and it is generally *not* free. A central goal of privacy research is to determine and, when possible, attain the minimum utility loss necessary to meet a desired privacy guarantee. Unfortunately, in the classic regime of *Input-Independent Indistinguishability* (III), as the foundation to protect *intractable a priori* secret, strong impossibility results are known that a privacy solution cannot be even validated in general. However, in the regime of randomly-generated secret with entropy, the problem, either from possibility and impossibility side, largely remains open. In this paper, we provide an affirmative answer in a more restrictive *black-box* setting: we assume no structural knowledge of the leakage-generation procedure, and the user is granted only limited access via end-to-end simulation. Stemmed from *PAC Privacy* and α -*Mutual Information*, we establish a universal framework that *automatically* and *asymptotically* approaches the optimal randomization strategy for an arbitrary *black-box* leakage function, given a Bayesian posterior budget against an *arbitrary* adversarial inference procedure.

We further generalize our results to the composition scenario, in which a random secret is reused across a sequence of potentially adversarially selected leakage applications. Still in a *black-box* setting, with respect to both adversarial strategy and leakage functions, we present a *locally optimal* randomization strategy and an asymptotically tight characterization of the cumulative privacy risk. We also provide a mechanism-agnostic variant by establishing a general composition upper bound using only marginal Bayesian risk measurements. Our results have strong practical potential for mitigating side-channel leakage, especially from physical channels that lack tractable analytical models. We include experiments on mitigating leakage from power consumption during Advanced Encryption Standard (AES) secret-key generation and from timing measurements in RSA encryption. The code and the full version with appendix can be found in the following anonymous Github link https://anonymous.4open.science/r/Optimal_Blackbox_Privacy-812F/.

Keywords: Black-Box Privacy Analysis, Adversarial Adaptive Composition, Secret Entropy, PAC Privacy, α -Mutual Information

1 Introduction

Any processing of a secret can potentially leak information. As one of the most extensively-studied security problems, privacy essentially asks the following: *given the information revealed through leakage, how much of the underlying*

secret can an adversary recover? Despite its intuitive interpretation, rigorous formalization and provable protection are highly non-trivial for many practical applications. At a high level, modern privacy research can be organized along three main directions: ❶ *definition-level*—how to characterize privacy risk and privacy guarantees, both *semantically* and *mathematically*; ❷ *measurement/accounting-level*—how to *tightly* quantify the privacy risk induced by a given leakage; and ❸ *operation-level*—how to construct *optimal* mitigations for a given leakage to ensure a required privacy guarantee.

1.1 Privacy Definition

On ❶, existing privacy definitions can be largely categorized by two types of secrets to protect: (A) entropic secrets sampled from an *accessible distribution*¹, and (B) secrets for which such entropy is not tractable. A canonical example of a Type-(A) secret is a uniformly-generated secret key in a cryptographic protocol. For Type-(A) secrets, a natural notion of privacy is Bayesian: the adversary’s optimal posterior success probability for recovering the secret [44]. In particular, when the adversary’s objective is *identification* (i.e., the success criterion is exact recovery of the secret), this risk coincides with the classic notion of (conditional) *min-entropy* studied in [7,34,22]. Related notions include *guessing entropy* [26], which captures the expected number of guesses required for successful identification.

A key prerequisite of the Bayesian viewpoint above is the ability to specify a prior distribution over the secret x , which is often inappropriate for Type-(B) secrets where we may not know, or may not have access to, the secret-generation mechanism². To obtain a rigorous notion in a *prior-free* setting, Shannon [36], in 1949, initiated the concept of *input-independent indistinguishability* (III), which essentially considers the worst case over all priors. Rather than directly measuring the adversary’s posterior knowledge, III evaluates the worst-case *posterior advantage*: the maximal change from prior to posterior when distinguishing between two secret hypotheses under an arbitrary prior distribution. Goldwasser and Micali [15] further introduced a computational form of III against computationally-bounded adversaries, laying the foundation of modern cryptography. They also showed that negligible III for binary hypotheses in the worst case implies negligible advantage for *any* efficient inference task. As a prior-free worst-case notion, III admits a clear semantic interpretation: regardless of the adversary’s (subjective) prior belief about the secret, observing the leakage should not significantly change their posterior belief. Beyond cryptography, many information-theoretic privacy definitions for Type-(B) secrets, including *Differential Privacy* (DP) [10,9], *Pufferfish Privacy* [21], and *Maximal Leakage* [18], are also rooted in the III philosophy.

¹ An *accessible distribution* refers to a secret-generation procedure from which we can repeatedly sample instances of the secret.

² Even if one could posit that all possible data is generated by some complicated stochastic process, in many real scenarios the observation is effectively unique and cannot be resampled at the relevant moment.

1.2 Universal Privacy Risk Measurement and Leakage Mitigation

The ❷ measurement problem and the ❸ mitigation problem require a deeper study of the leakage itself. Broadly speaking, leakage can be any quantity statistically correlated with a secret x that contains sensitive information, and it can be generally represented as the output of a function $\mathcal{F}(x)$. In practice, leakage $\mathcal{F}(x)$ may arise from the release of any algorithmic processing of sensitive data x , for instance, a machine learning model trained on health data [2]. Leakage can also arise as a physical signal related to the processing of x : a large body of work has identified many side channels, such as timing [6,48], power [33], memory [3] consumption, or network traffic patterns [41], that may enable successful attacks.

To mitigate leakage (❸), most existing approaches introduce additional randomization into $\mathcal{F}(x)$. A widely used strategy is perturbation, releasing a noisy version $\mathcal{F}(x) + e$. When $\mathcal{F}(x)$ is a statistic (e.g., a mean salary estimate or neural-network weights), e may be independent Gaussian [9] or Laplace [10] noise, as is standard in the DP literature. For physical signals $\mathcal{F}(x)$, e may also be introduced physically, e.g., by sending dummy queries in anonymous communication systems [41]. In any case, privacy mechanisms typically trade utility for protection. The appropriate utility loss measure depends on the application: in anonymous communication, for example, the loss may be the delay introduced by dummy messages, with the relevant metric determined by latency sensitivity or protocol constraints.

Given the complexity of practical leakage functions $\mathcal{F}(\cdot)$ and the diversity of utility metrics, it is desirable to have a *universal* framework that addresses ❷ (tight measurement) and ❸ (optimal mitigation) simultaneously. Unfortunately, in the Type-(B) regime, neither tight measurement nor optimal mitigation can be efficiently determined under III-style definitions in general. Indeed, Xiao and Tao [46] show that computing sensitivity, the maximal outcome difference

$$\sup_{\bar{x}, \bar{x}'} \|\mathcal{F}(\bar{x}) - \mathcal{F}(\bar{x}')\|$$

over arbitrary secret hypotheses (\bar{x}, \bar{x}') , is NP-hard in the worst case. This suggests that even *verifying* whether a leakage function $\mathcal{F}(\cdot)$ satisfies a given III-style guarantee can be computationally intractable.

As a consequence, known III-style privacy analysis for Type-(B) secrets typically apply only in specific *white-box* settings, and one of most representative methodologies is *decompose-then-compose*. In particular, over the last two decades the DP literature has developed a rich collection of such tools, including *subsample-then-aggregate* [29], the *Sparse Vector Technique* [16,25], *Differentially Private SGD* (DP-SGD) [1], PATE [30], and *private evolution* [23]. In these frameworks, a complex processing pipeline is decomposed into simpler components³, each of which is separately privatized and then composed. However,

³ Ideally, each component in a *decompose-then-compose* framework ensures bounded sensitivity $\sup_{\bar{x}, \bar{x}'} \|\mathcal{F}(\bar{x}) - \mathcal{F}(\bar{x}')\| \leq c$ for some threshold c under an appropriate

this methodology can lead to overly conservative (loose) noise requirements, especially for high-dimensional tasks with intricate dependencies among sensitive inputs, such as clustering [14] and outcomes from machine learning models [32]. These limitations can be even more severe for settings without tractable algorithmic structure, including many of the physical side channels [6,33,3] mentioned earlier.

In contrast, for Type-(A) entropic secrets, the landscape can be starkly different for both ② and ③. While tightness and optimality remained open, recent work PAC Privacy [44] demonstrates the possibility of universal privacy analysis even in a more restricted *black-box* setting. Here, *black-box* means that we assume no particular algorithmic structure for $\mathcal{F}(\cdot)$; instead, the user only has oracle access to end-to-end evaluations of $\mathcal{F}(\cdot)$. In the context of ② and ③, this means that both risk measurement and mechanism design must be derived from a finite number m of evaluations of $\mathcal{F}(\cdot)$ on selected inputs $\mathbf{X} = \{\bar{x}_i : i \in [m]\}$. Not surprisingly, the follow-up works [45,39] show that the provable Bayesian upper bounds derived from PAC Privacy can be significantly sharper than the prior-free worst case in many applications, after the secret's inherent entropy is exploited which makes inference intrinsically harder.

Fundamentally, this difference in addressing ② and ③ with respect to Type-(A) and Type-(B) secrets arises because an III-style guarantee in Type-(B) scenario demands divergence control between leakages from *arbitrary* pairs of secret hypotheses.⁴ This strict "arbitrary" requirement is the root cause of the computational impossibility [46]. For Type-(A) secrets, however, it suffices to only analyze Bayesian posteriors with respect to the *ground-truth* secret distribution \mathcal{D} [20], since any mis-specified or biased prior always weakens the adversary's inference performance relative to the optimal Bayesian adversary. Thus, even if extreme worst cases exist and are hard to identify, they may be negligible under \mathcal{D} , enabling near-optimal Bayesian analysis. Moreover, higher entropy in \mathbf{x} (more uncertainty) should allow substantially cheaper privacy mechanisms than conservative prior-free worst-case approaches. We formalize ② and ③ for Type-(A) entropic secrets in the black-box setting as follows.

Problem 1 (Tight black-box privacy analysis with entropy). Given a prior distribution \mathcal{D} of an entropic secret \mathbf{x} and an arbitrary leakage function $\mathcal{F}(\cdot) : \mathcal{X} \rightarrow \mathcal{O}$, we want to design a (possibly randomized) perturbation-selection mechanism $\text{Alg} : (\mathcal{X} \times \mathcal{O})^m \rightarrow \mathbb{P}$ that outputs a distribution \mathcal{D}_e of perturbation \mathbf{e} with *minimal overhead* (measured by some metric $\kappa(\cdot)$), such that for an arbitrary inference task captured by a criterion ρ and target posterior success rate δ_ρ , the following experiment is *impossible* for an informed adversary:

metric $\|\cdot\|$. In the DP context, a common structure is *aggregation*: each data point can be clipped to a bounded range and contributes *independently* to the released output [1].

⁴ This is necessary if one wants to upper bound the posterior advantage under *arbitrary* priors [15].

1. The user selects m strings $\mathbf{X} = \{\bar{x}_i \in \mathcal{X} : i \in [m]\}$, evaluates $\mathcal{F}(\bar{x}_i)$, and applies Alg to the simulation data $\mathbf{S} = \{(\bar{x}_i, \mathcal{F}(\bar{x}_i)) : i \in [m]\}$. The mechanism outputs a perturbation distribution $\text{Alg}(\mathbf{S}) = \mathcal{D}_e$.
2. The user samples an $\mathbf{x} \sim \mathcal{D}$, samples perturbation $\mathbf{e} \sim \mathcal{D}_e$, and sends the noisy leakage $\mathcal{F}(\mathbf{x}) + \mathbf{e}$ to the adversary.
3. The adversary, who has full knowledge of $\mathcal{F}(\cdot)$ and the user's mechanism Alg , outputs $\hat{\mathbf{x}}$ such that $\rho(\hat{\mathbf{x}}, \mathbf{x}) = 1$ with probability greater than δ_ρ .

In Problem 1, we follow PAC Privacy [44] and model an arbitrary adversarial inference task via a criterion function $\rho(\cdot, \cdot)$, measuring privacy risk by the Bayesian posterior success probability δ_ρ . The criterion ρ can be chosen so that the preimage of $\rho(\cdot, \mathbf{x}) = 1$ captures all *unacceptable* reconstructions of \mathbf{x} . For instance, given an l -bit secreting string \mathbf{x} , $\rho(\hat{\mathbf{x}}, \mathbf{x}) = 1$ iff $\hat{\mathbf{x}} = \mathbf{x}$ captures exact identification; $\rho(\hat{\mathbf{x}}, \mathbf{x}) = 1$ iff the Hamming distance between $\hat{\mathbf{x}}$ and \mathbf{x} is at most c captures approximate reconstruction with at most c -bit errors. As for the utility-loss metric κ , a common choice is the variance of the injected noise \mathbf{e} ; in general, one may determine a score for each output instance in \mathcal{O} and selects κ as the expected quality of the outcome from the perturbed mechanism $\mathcal{M}(\mathbf{x}) = \mathcal{F}(\mathbf{x}) + \mathbf{e}$.

On the other hand, it is worthwhile noting that the protection strategy \mathcal{D}_e , the distribution of the perturbation/modification \mathbf{e} imposed to leakage $\mathcal{F}(\mathbf{x})$ can be dependent on the secret input \mathbf{x} , which, strictly speaking, should be written as $\mathcal{D}_e(\mathbf{x})$, although in most existing works [9,1], \mathbf{e} is simply an *independent* (e.g., Gaussian) noise. For notation simplicity we still use \mathcal{D}_e in the following. We also emphasize the adversarial model in the following remark.

Remark 1 (White-Box Adversarial View). Throughout the paper, the *black-box* assumption is from the user's perspective: the user lacks structural knowledge of how leakage is incurred. The adversary, in contrast, is always assumed to be fully informed of both the leakage function \mathcal{F} and the user's privacy-preserving strategy Alg . The only uncertainty from the adversary's perspective is the randomness in secret generation (e.g., $\mathbf{x} \sim \mathcal{D}$) and in the modified leakage process (e.g., $\mathbf{e} \sim \mathcal{D}_e$), over which the posterior success probability δ_ρ is defined.

1.3 Privacy Composition

So far, we have focused on a static leakage function $\mathcal{F}(\cdot)$. In practice, leakage may arise in a potentially adversarial interactive environment, and when the secret is reused, we must track its cumulative leakage. To model leakage in such a dynamic setting, we consider a more general function $\mathcal{F}(\mathbf{x}, \mathbf{v})$, where an environment parameter \mathbf{v} can be adaptively selected by an adversary to maximize information gained through sequential leakages. For example, given a secret key \mathbf{x} , a user may be asked to encrypt multiple messages $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T\}$, which can be adversarially chosen; each encryption produces leakage $\mathcal{F}(\mathbf{x}, \mathbf{v}_t)$. Intuitively, an adversary can adaptively choose the next query \mathbf{v}_t based on accumulated information about \mathbf{x} , potentially increasing their posterior success probability.

Existing composition results primarily address Type-(B) secrets under III-style guarantees, and have largely been developed in the DP literature [11,19,5,27,8]. Tight composition accounting has also motivated new privacy definitions, evolving from pure ϵ -DP to approximate (ϵ, δ) -DP, to concentrated DP [5] and Rényi DP [27], and most recently to f -DP [8]. Typically, these results rely on two key prerequisites: (i) a *global* worst-case III guarantee for all possible leakage functions $\mathcal{F}(\cdot, \mathbf{v})$ (independent of auxiliary information and adversarial adaptivity), and (ii) independence of the randomness used at each step of the composed mechanism: the global worst-case guarantee (i) bounds the posterior advantage increment regardless of the adversary’s query, while independence in (ii) enables aggregation via concentration techniques.

For Type-(A) entropic secrets in the *black-box* regime, however, the situation changes dramatically and existing III-style composition analyses do *not* directly apply. First, in a black-box setting, we cannot generally obtain non-trivial input-independent guarantees from finitely many end-to-end simulations. Second, and more importantly, the entropic secret \mathbf{x} is generated once and reused: even if fresh randomization is applied at each step, the resulting leakages remain correlated through the shared secret. We formalize the composition version of Problem 1 as follows.

Problem 2 (Noise mechanism for tight composition). Given a prior secret distribution \mathcal{D} , a privacy budget in terms of posterior success rate δ_ρ for an inference task captured by ρ , and a family of (black-box) leakage functions $\mathcal{F}(\cdot, \cdot) : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{O}$, we want to design a joint noise mechanism Alg such that there does *not* exist an adversarial query algorithm \mathcal{Q}_{adv} for which the following experiment between a user (the secret holder) and the adversary is possible:

1. The user samples a secret $\mathbf{x} \sim \mathcal{D}$.
2. At iteration $t \in [T]$, from the adversary side, possibly based on all previously observed leakages $\{\mathbf{o}_j = \mathcal{F}(\mathbf{x}, \mathbf{v}_j) + \mathbf{e}_j : j \in [t-1]\}$ and past queries $\{\mathbf{v}_j : j \in [t-1]\}$, the adversary chooses a new query $\mathbf{v}_t \leftarrow \mathcal{Q}_{\text{adv}}(\{\mathbf{o}_j, \mathbf{v}_j\}_{j < t})$ and sends \mathbf{v}_t to the user.
3. At iteration $t \in [T]$, from the user side, after receiving the adversarial query \mathbf{v}_t , the user performs end-to-end simulations of the leakage functions $\{\mathcal{F}(\cdot, \mathbf{v}_j) : j \in [t]\}$ on a selected input set $\{\bar{\mathbf{x}}_i^{(t)} : i \in [m]\}$, and applies Alg to the resulting evaluations

$$\mathbf{S}^{(t)} = \{(\bar{\mathbf{x}}_i^{(t)}, \mathcal{F}(\bar{\mathbf{x}}_i^{(t)}, \mathbf{v}_j)) : i \in [m], j \in [t]\},$$

together with the previously sampled noises $\mathbf{e}^{(t-1)} = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{t-1})$, to determine a noise distribution $\mathcal{D}_e^{(t)}$. The user then samples $\mathbf{e}_t \sim \mathcal{D}_e^{(t)}$ and returns the noisy leakage $\mathbf{o}_t = \mathcal{F}(\mathbf{x}, \mathbf{v}_t) + \mathbf{e}_t$ to the adversary.

4. After T iterations, the adversary outputs an estimate $\hat{\mathbf{x}}$ based on $\{(\mathbf{o}_t, \mathbf{v}_t)\}_{t \in [T]}$ such that $\rho(\hat{\mathbf{x}}, \mathbf{x}) = 1$ with probability greater than δ_ρ .

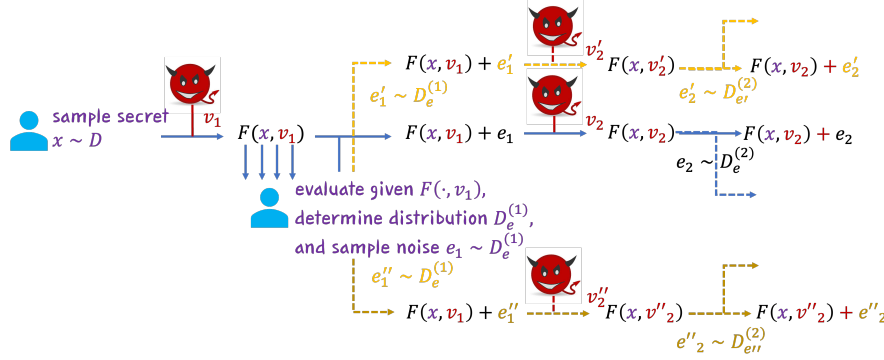


Fig. 1. Illustration of the *Parallel Universe Challenge* in adversarially adaptive composition for entropic secrets. The solid blue line corresponds to the specific (local) realization observed by the user, while dashed lines represent alternative (unobserved) realizations induced by randomness in both the secret and the mechanism.

Figure 1 illustrates the technical challenge behind Problem 2, which we refer to as the *parallel universe* restriction. With an adaptive adversary, the joint distribution of the composed leakages $\{\mathcal{F}(x, v_t)\}_{t=1}^T$ is jointly determined by the (unknown) adversarial strategy \mathcal{Q}_{adv} and the user's actions (secret sampling and noise injection). Under the black-box model, the user cannot infer \mathcal{Q}_{adv} or the induced global leakage distribution required for a direct Bayesian analysis. The user's only visibility into \mathcal{Q}_{adv} is through the adversary's reaction *conditional* on the particular randomness realized in the user's own execution. Concretely, even conditioning on a first query v_1 and a first-round noise distribution $\mathcal{D}_e^{(1)}$, the user observes only the adversary's next query v_2 corresponding to the specific sampled noise instance e_1 . Different noise realizations $e'_1 \sim \mathcal{D}_e^{(1)}$ could induce different adversarial branches—analogue to "parallel universes" inaccessible to the user. Problem 2 thus asks the user to infer and mitigate leakage from the *local* universe they observe, while ensuring that aggregating over all (unobserved) branches yields the desired global posterior guarantee.

Stemmed from Problem 2, a further generalization is to develop composition theory that is agnostic to the specific privacy-preserving mechanism used at each step. A natural application setting is one where multiple parties share a secret but are processing the secret separately, and each party lacks knowledge of the leakage functions and defenses used by others. That is, given T black-box leakage functions $\{\mathcal{F}(x, v_t) : t \in [T]\}$ that share the same entropic secret x , and given only that each marginal leakage satisfies certain posterior inference hardness with respect to x , how can we upper bound the cumulative Bayesian privacy risk? This motivates the following problem.

Problem 3 (Mechanism-agnostic composition). Given a random secret $x \sim \mathcal{D}$ and an arbitrary sequence of T leakage functions $\{\mathcal{F}(\cdot, v_t) : t \in [T]\}$, suppose that, conditioned on the shared secret input x , all additional random-

ness in each $\mathcal{F}(\cdot, \mathbf{v}_t)$ is independent across $t \in [T]$. Under what types of *marginal* measurements for each $\mathcal{F}(\cdot, \mathbf{v}_t)$ can we derive a global bound on the cumulative Bayesian privacy risk?

1.4 Techniques Overview

In this section, we provide a roadmap of techniques to address Problem 1 - 3. We begin with a special case of Problem 1 where the inference task ρ is fixed to identification (exact reconstruction). In this case, the adversary's optimal inference rule $\mathcal{A}(o)$ given an observation $o = \mathcal{F}(\mathbf{x})$ is the classical *Maximum A Posteriori* (MAP) estimator [47]:

$$\mathcal{A}(o) \in \arg \max_{\hat{\mathbf{x}}} \Pr(\mathbf{x} = \hat{\mathbf{x}} \mid \mathcal{F}(\mathbf{x}) = o), \quad (1)$$

i.e., the adversary outputs the hypothesis with maximal posterior probability.

Accordingly, the optimal identification success probability, denoted by δ_I , can be written as an aggregate of weighted 0-1 losses:

$$\begin{aligned} \delta_I &= \Pr_{\mathbf{x} \sim \mathcal{D}} [\mathcal{A}(\mathcal{F}(\mathbf{x})) = \mathbf{x}] \\ &= \sum_{\hat{\mathbf{x}}} \Pr_{\mathbf{x} \sim \mathcal{D}} (\mathbf{x} = \hat{\mathbf{x}}) \cdot \int_o \mathbf{1} \left[\hat{\mathbf{x}} \in \arg \max_z \Pr(\mathbf{x} = z \mid \mathcal{F}(\mathbf{x}) = o) \right] \cdot \Pr(\mathcal{F}(\hat{\mathbf{x}}) = o) \, do \\ &= \mathbb{E}_{\hat{\mathbf{x}} \sim \mathcal{D}} \left[\int_o \mathbf{1} \left[\hat{\mathbf{x}} \in \arg \max_z \Pr(\mathbf{x} = z \mid \mathcal{F}(\mathbf{x}) = o) \right] \cdot \Pr(\mathcal{F}(\hat{\mathbf{x}}) = o) \, do \right]. \end{aligned} \quad (2)$$

Here, $\mathbf{1}(C) = 1$ iff condition C holds. Equation (2) shows that δ_I is the probability mass of observation events for which the *true* secret happens to be MAP-optimal given the observed leakage.

Unfortunately, in the black-box setting, directly computing (2) is typically infeasible because the secret space \mathcal{X} may be exponentially large. However, the expectation form in (2) suggests a route toward **2**: if there exists a *bounded* function $\mathcal{G}(\cdot)$ such that

$$\delta_I = \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{D}} [\mathcal{G}(\bar{\mathbf{x}}, \mathcal{F}(\bar{\mathbf{x}}))], \quad (3)$$

then the empirical average $\frac{1}{m} \sum_{i=1}^m \mathcal{G}(\bar{\mathbf{x}}_i, \mathcal{F}(\bar{\mathbf{x}}_i))$ over samples $\mathbf{X} = \{\bar{\mathbf{x}}_i\}_{i=1}^m$ is an unbiased estimator of δ_I , and concentration inequalities (e.g., Hoeffding's) yield non-asymptotic error bounds.

To incorporate optimal mitigation (**3**), we can represent a privacy mechanism by a perturbation term $\mathbf{e} \sim \mathcal{D}_e$ and formulate a constrained optimization problem: under a utility loss measure κ and a privacy requirement $\delta_I \leq \epsilon$,

$$\min_{\mathcal{D}_e} \kappa(\mathcal{D}_e) \quad \text{s.t.} \quad \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{D}} [\mathcal{G}(\bar{\mathbf{x}}, \mathcal{F}(\bar{\mathbf{x}}))] \leq \epsilon. \quad (4)$$

As in statistical learning, because the population expectation constraint in (4) is unknown, one may replace it with an empirical proxy:

$$\min_{\mathcal{D}_e} \kappa(\mathcal{D}_e) \quad \text{s.t.} \quad \frac{1}{n} \sum_{i=1}^m \mathcal{G}(\bar{\mathbf{x}}_i, \mathcal{F}(\bar{\mathbf{x}}_i)) \leq \epsilon'. \quad (5)$$

To ensure that solutions satisfying (5) also satisfy (4) with high probability (i.e., to control generalization error), a standard approach is to choose a stricter threshold $\epsilon' < \epsilon$ and bound the gap $\epsilon - \epsilon'$ via uniform convergence or complexity measures (e.g., Rademacher complexity) [43]. In black-box settings, one may alternatively adopt a *propose-verify* strategy that gradually increases noise ϵ until an independent test (using fresh samples) is passed.

All of the above depends on constructing such a function \mathcal{G} , which is generally challenging. Building on PAC Privacy [44] and α -mutual information [37,42,18], we establish a new family of estimators $\{\tilde{\mathcal{G}}_\alpha : \alpha \in \mathbb{R}^+\}$ (termed *reference-prior-weighted α -information*), which yields an unbiased representation of δ_I in a *doubly asymptotic* sense:

$$\delta_I = \mathbb{E}_{\bar{x} \sim \mathcal{D}} \left[\lim_{\alpha \rightarrow \infty} \tilde{\mathcal{G}}_\alpha(\bar{x}, \mathcal{F}(\bar{x})) \right]. \quad (6)$$

More importantly, for any fixed α , we have the upper bound

$$\delta_I \leq \mathbb{E}_{\bar{x} \sim \mathcal{D}} \left[\tilde{\mathcal{G}}_\alpha(\bar{x}, \mathcal{F}(\bar{x})) \right],$$

so any non-asymptotic analysis based on a specific $\tilde{\mathcal{G}}_\alpha$ provides a provable high-probability guarantee, with larger α yielding tighter bounds. Building on the identification case, in Section 3.2 we further show that analyzing an arbitrary inference criterion ρ can be reduced to an identification problem, which completes the solution to Problem 1.

For composition, as illustrated in Fig. 1, the central challenge is that intermediate randomness and adversarial adaptivity induce a global leakage distribution—a mixture over potentially infinitely many conditional branches ("parallel universes"); but the user observes only a single realized path (the "local universe"). Rather than performing a direct Bayesian analysis of this global mixture, we propose a randomization mechanism constrained to the local universe, ensuring that the cumulative risk for the realized leakage instance *matches* the privacy budget exactly. Intuitively, while each instance operates locally, if all alternate instances across parallel universes satisfies their respective local guarantees, the collection of local controls yields a tight global composition bound.

To capture this, we consider decomposing the overall privacy risk accounting into a sequence of smaller subproblems, each solvable within its appropriate "universe". Once composed, these local controls provide a global privacy guarantee. We consider two composition setups, distinguished by whether the per-round mechanism may depend on the realized transcript of previous rounds.

In the *history-accessible* setting, the user is allowed to condition each round's mechanism on the full interaction history, including previously selected mechanisms, realized outputs, and adversarial queries. This allows one to track a running "leakage score" for each candidate secret along the current transcript. Consequently, this enables a *branch-dependent* decomposition: per-round subproblems may depend on the current history, allowing the user to adaptively calibrate the randomization at each round.

Conversely, in the *mechanism-agnostic* setting, the decomposition must be determined *independently* of the realized history. In this scenario, the adversary may concentrate on the most vulnerable region of the secret space to maximize inference success, and the user is therefore compelled to enforce a worst-case (uniform) guarantee over all possible, unseen histories with carefully-selected separate measurement of each marginal leakage function.

2 Preliminaries

2.1 Statistical Divergence

We first introduce a family of useful distributional-distance measurements.

Definition 1 ((Generalized) f-Divergence). *Let $f : (0, +\infty) \rightarrow \mathbb{R}$ be a convex function. Let P and Q be two probability distributions defined on the same measurable space, with P absolutely continuous with respect to Q . The f -divergence between P and Q is defined as*

$$D_f(P \parallel Q) := \mathbb{E}_{X \sim Q} \left[f \left(\frac{dP}{dQ}(X) \right) \right]. \quad (7)$$

Here, dP (resp., dQ) denotes the probability density function in the continuous case or the probability mass function in the discrete case.

Classic definition of f -divergence [35] further requires $f(1) = 0$ to ensure the divergence between two identical distributions equals 0, which is not necessary for the analysis in this paper. It is also noted that many commonly-used statistical divergences are special cases of f -divergence. For example, *KL-divergence* and *total variation* correspond to the selection of $f(z) = z \log(z)$ and $f(z) = |z - 1|$ in Definition 1, respectively. In particular, one important variant we will constantly use in the paper is the α -divergence, D_α , by selecting $f(z) = z^\alpha$.

Definition 2 (α -Divergence). *Let P and Q be two probability distributions defined on the same measurable space. For any $\alpha > 1$,*

$$D_\alpha(P \parallel Q) := \mathbb{E}_{x \sim Q} \left[\left(\frac{dP}{dQ}(x) \right)^\alpha \right].$$

As a side note, the well-known Rényi-divergence [27] can be viewed as a logarithmic transformation of the α -divergence, $\frac{1}{\alpha-1} \log D_\alpha(P \parallel Q)$. An important property of f -divergence is their joint convexity, as shown below.

Lemma 1 (Joint convexity of f-Divergence). *Let $f : (0, +\infty) \rightarrow \mathbb{R}$ be a convex function. Then the function*

$$(z, v) \mapsto v f\left(\frac{z}{v}\right)$$

is jointly convex in (z, v) over $z \geq 0$ and $v > 0$. Consequently, f -divergence is jointly convex with respect to the pair of mixture distributions $(c \cdot P_1 + (1 - c)P_2, cQ_1 + (1 - c)Q_2)$ for any $c \in [0, 1]$, such that

$$D_f(c \cdot P_1 + (1 - c) \cdot P_2 \parallel c \cdot Q_1 + (1 - c) \cdot Q_2) \leq cD_f(P_1 \parallel Q_1) + (1 - c)D_f(P_2 \parallel Q_2).$$

2.2 Privacy Definitions

In the following, we cover the background of privacy definitions. As discussed earlier, depending on the type of secret to protect, there are two main philosophies: *Input-Independent Indistinguishability* (III) for Type-(B) secret and Bayesian for Type-(A) secret.

Definition 3 (ϵ input-independent indistinguishability [36,15,10,21]). Given a leakage function $\mathcal{F}(\cdot) : \mathcal{X} \rightarrow \mathcal{O}$ and for an arbitrary pair of adjacent secret hypotheses $\{\bar{x}, \bar{x}'\}$ such that \bar{x} and \bar{x}' only differ in the sensitive information that desires protection, then $\mathcal{F}(\cdot)$ satisfies ϵ -III if the following experiment is impossible:

A user randomly selects x from $\{\bar{x}, \bar{x}'\}$ and sends $\mathcal{F}(x)$ to the adversary; the (computationally bounded or unbounded) adversary can return a guessing \hat{x} such that with probability at least ϵ , taken over the randomness of x and \mathcal{F} , $\hat{x} = x$.

As for Type-(A) entropic secret, the main focus of this paper, a general and formal definition of Bayesian risk is *PAC Privacy* [44].

Definition 4 ($(\mathcal{D}, \rho, \delta_\rho)$ PAC Privacy [44]). For a leakage function $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{O}$, let \mathcal{D} be the prior distribution of secret x over \mathcal{X} , and $\rho : \mathcal{X} \times \mathcal{X} \rightarrow \{0, 1\}$ be an inference criterion: $\rho(\hat{x}, x) = 1$ iff an adversarial guess \hat{x} is successful. Then \mathcal{F} satisfies $(\mathcal{D}, \rho, \delta_\rho)$ PAC Privacy if the following experiment is impossible:

A user samples $x \sim \mathcal{D}$ and releases $\mathcal{F}(x)$ to an informed adversary who knows both \mathcal{D} , \mathcal{F} and ρ . The adversary returns an estimation \hat{x} of x such that with probability greater than δ_ρ , $\rho(\hat{x}, x) = 1$.

In the following, when there is no contextual ambiguity, we will simply call the $(\mathcal{D}, \rho, \delta_\rho)$ PAC Privacy as δ_ρ privacy.

2.3 Bayesian Analysis Tools

In this section, we introduce some information-theoretical tools for the objective Bayesian analysis. One operationally-efficient result, which we term *f-divergence Fano Inequality*, as a slightly-generalized version of the main theorem of PAC Privacy (Theorem 1 in [44]), presents an upper bound of target δ_ρ in an expectation form.

Lemma 2 (f-divergence Fano Inequality [44]). *For any selected f-divergence D_f , an arbitrary leakage function $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{O}$ satisfies $(\delta_\rho, \rho, \mathcal{D})$ PAC Privacy where the adversarial posterior success rate δ_ρ can be selected as the minimal number satisfying either one of the following inequalities:*

$$\begin{aligned} D_f(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{o,\rho}}) &\geq \inf_{\mathbb{P}_W} D_f(\mathbb{P}_{\mathbf{x}, \mathcal{F}(\mathbf{x})} \parallel \mathbb{P}_{\mathbf{x}} \otimes \mathbb{P}_W) = \inf_{\mathbb{P}_W} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_f(\mathbb{P}_{\mathcal{F}(\mathbf{x})} \parallel \mathbb{P}_W)] ; \\ D_f(\mathbf{1}_{\delta_{o,\rho}} \parallel \mathbf{1}_{\delta_\rho}) &\geq \inf_{\mathbb{P}_W} D_f(\mathbb{P}_{\mathbf{x}} \otimes \mathbb{P}_W \parallel \mathbb{P}_{\mathbf{x}, \mathcal{F}(\mathbf{x})}) = \inf_{\mathbb{P}_W} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_f(\mathbb{P}_W \parallel \mathbb{P}_{\mathcal{F}(\mathbf{x})})] . \end{aligned} \quad (8)$$

Here, $\delta_{o,\rho}$ is the optimal prior success rate of recovering \mathbf{x} before observing the leakage; $\mathbf{1}_{\delta_\rho}$ and $\mathbf{1}_{\delta_{o,\rho}}$ are two Bernoulli distributions of parameters δ_ρ and $\delta_{o,\rho}$, respectively; $\mathbb{P}_{\mathbf{x}, \mathcal{F}(\mathbf{x})}$ and $\mathbb{P}_{\mathbf{x}}$ are the joint distribution and the marginal distribution of $(\mathbf{x}, \mathcal{F}(\mathbf{x}))$ and \mathbf{x} , respectively; and \mathbb{P}_W represents the distribution of an arbitrary random variable W defined over \mathcal{O} .

It has been proved (Lemma 1 in [44]) that the left hand sides of (8), $D_f(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{o,\rho}})$ and $D_f(\mathbf{1}_{\delta_{o,\rho}} \parallel \mathbf{1}_{\delta_\rho})$, are both monotone to δ_ρ . Thus, Lemma 2 suggests that through optimizing a *reference* marginal distribution \mathbb{P}_W over the leakage domain, the expectation of an arbitrary f-divergence between leakage likelihood $\mathbb{P}_{\mathcal{F}(\mathbf{x})}$ conditional on specific secret \mathbf{x} and the \mathbb{P}_W can form an upper bound to target δ_ρ . In particular, when the f-divergence D_f is selected to be the KL-divergence, (8) reduces to the well-known Fano's inequality [13], which connects the Shannon entropy and the Bayesian. Another concept related to (8) is α -mutual information [37] by selecting the f-divergence D_f to be α -divergence (Definition 2).

Definition 5 (Sibson's α -mutual information [37]). *Let \mathbb{P}_{XY} be some joint distribution of random variables (X, Y) whose marginal distribution with respect to X is \mathbb{P}_X . For any $\alpha > 1$, Sibson's α -mutual information is defined by the Rényi information radius for arbitrary marginal distribution Q of Y ,*

$$I_\alpha^{sib}(X; Y) := \inf_Q \frac{1}{\alpha - 1} \cdot \log D_\alpha(\mathbb{P}_{XY} \parallel \mathbb{P}_X \otimes Q). \quad (9)$$

When $\alpha \rightarrow \infty$, $I_\infty^{sib}(X; Y)$ forms the foundation of *maximal leakage* [18]. In a *prior-free* setup, taken over the worst case of prior distribution \mathbb{P}_X , $\sup_{\mathbb{P}_X} I_\infty^{sib}(X; Y)$ captures the logarithm of the maximal multiplicative posterior gain [18].

3 Towards Tight Privatization of Static Leakage Function

In this section, we systemically address Problem 1 for static leakage function. Provided the simulatable nature of the expectation form in (8), a natural question stemmed from Lemma 2 is that whether the equality of (8) is achievable. That is to say, through optimizing the selection of f-divergence and reference distribution \mathbb{P}_W , whether the optimal posterior rate δ_ρ can be tightly characterized from (8). In Section 3.1, we start to identify the possibility of a special

case: the uniform-prior identification setting, for which we develop progressively tighter upper bounds on δ_I through α -divergence and translate them into concrete, asymptotically optimal randomization rule \mathcal{D}_e (Algorithm 1). However, for general priors and inference settings, the application of Lemma 2 is *not* tight anymore and thus in Section 3.2, we establish a generalized structure, termed *Prior-(Reference)-Weighted α -Information*, and show how to tightly determine and control general δ_ρ .

3.1 Tight α -Information to Identification under Uniform Priors

We start by revisiting f-divergence Fano Inequality (Lemma 2) in a setting when \mathcal{D} is a uniform distribution and studying the conditions that can lead to the equality in (8). As discussed in Section 1.4, after observing a leakage $\mathbf{o} = \mathcal{F}(\mathbf{x})$, the adversary's optimal decision rule \mathcal{A} is MAP decoding:

$$\mathcal{A}(\mathbf{o}) = \arg \max_{\hat{\mathbf{x}}} \Pr(\mathbf{x} = \hat{\mathbf{x}} \mid \mathcal{F}(\mathbf{x}) = \mathbf{o}) = \arg \max_{\hat{\mathbf{x}}} \Pr(\mathbf{x} = \hat{\mathbf{x}}) \Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o}).$$

Under a uniform prior, $\Pr(\mathbf{x} = \hat{\mathbf{x}})$ is constant, so the success probability is driven by the largest likelihood $\max_{\hat{\mathbf{x}}} \Pr(\mathcal{F}(\hat{\mathbf{x}}) = \mathbf{o})$. This leads to the 0-1-loss expression of δ_I in (2), where only these most likely hypotheses determine the optimal posterior rate. This intuitively contrasts with the expectation structure in the right-hand side of (8)

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} D_f(\mathbb{P}_{\mathcal{F}(\mathbf{x})} \parallel \mathbb{P}_W) = \int_{\mathcal{O}} \Pr(W = \mathbf{o}) \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[f\left(\frac{\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})}{\Pr(W = \mathbf{o})}\right) \right] d\mathbf{o}, \quad (10)$$

which averages contributions from all \mathbf{x} . This suggests that the choice of f matters. To approach equality, we expect the f in (10) grows steeply for large arguments $\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})$ so that the expectation can be dominated by the largest ratio terms, and (10) behaves more like a worst-case control rather than an average-case one, yielding a tighter bound.

Motivated by this, we focus on the family $f(z) = z^\alpha$, which induces the α -divergence (Definition 2). Compared to the KL choice $f(z) = z \log z$ adopted in PAC Privacy [44], which grows only slightly super linearly, the power function t^α penalizes large likelihood ratios much more aggressively as α increases. Indeed, it can be proved that as $\alpha \rightarrow \infty$, (10) provides tight characterization of δ_I .

Proposition 1 (Asymptotic tightness of α -divergence approximation).

Let $\mathbf{x} \sim \mathcal{D}$ be a secret input and let $\mathcal{F}(\cdot) : \mathcal{X} \rightarrow \mathcal{O}$ be an arbitrary leakage function. Define a reference distribution \mathbb{P}_W on the \mathcal{O} by

$$\Pr(W = \mathbf{o}) = \frac{(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})^\alpha])^{1/\alpha}}{C},$$

where the normalizing constant C is given by

$$C = \int_{\mathcal{O}} (\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})^\alpha])^{1/\alpha} d\mathbf{o}.$$

Then, as $\alpha \rightarrow \infty$, the upper bound of posterior success rate obtained from the Lemma 2 with respect to W is exactly the constant C , which converges to the optimal posterior success rate δ_I .

Proposition 1 shows that, by choosing an appropriate reference distribution \mathbb{P}_W , the α -divergence bound becomes asymptotically tight as $\alpha \rightarrow \infty$. It also tells us that the optimal reference W in (8) is proportional to the α -norm aggregation of the output likelihoods, namely

$$\Pr(W = \mathbf{o}) \propto \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})^\alpha] \right)^{1/\alpha}, \quad (11)$$

which is exactly the choice that makes Hölder's inequality tight in the proof of Proposition 1.

In practice, exact computation of \mathbb{P}_W and the expectation (10) can be computationally infeasible, given potentially exponentially-large input domain \mathcal{X} . Fortunately, both \mathbb{P}_W and (10) can be asymptotically approached from finite samples. Given a i.i.d. sample set $\mathbf{X} = \{\bar{x}_i, i \in [m]\}$, $\bar{\mathbf{x}} \sim \mathcal{D}$, W can be instead selected as

$$\Pr(W = \mathbf{o}) \propto \left(\sum_{i=1}^m \frac{1}{m} \cdot [\Pr(\mathcal{F}(\bar{x}_i) = \mathbf{o})^\alpha] \right)^{1/\alpha}, \quad (12)$$

and notably from Lemma 2, any (even if non-optimal) selection of W always leads to a safe upper bound of target δ_I .

As for the estimation of expectation (10), concentration inequalities allow us to make high-probability statement on the estimation error from sampled instance \mathbf{X} . However, we need to be careful here since the failure rate of a small estimation error *cannot* be straightforwardly merged into the upper bound of δ_I : by Lemma 2, instead we need a global, *deterministic* upper bound on (10), the right hand side of Lemma 2. To address this, our solution is to treat the randomness of \mathbf{X} as part of the entire privacy-preserving mechanism's randomness, and view our final randomization strategy \mathcal{D}_e as a function of the sampled instance \mathbf{X} as well. To be formal, let $\mathcal{M}(\cdot) = \mathcal{F}(\cdot) + \mathbf{e}$, $\mathbf{e} \sim \mathcal{D}_e$, where \mathcal{D}_e is determined by some mechanism $\text{Alg}(\mathbf{X})$, which takes a randomly-sampled set \mathbf{X} as input. Thus, the output distribution of $\mathcal{M}(\cdot)$ is a mixture over the sample-conditioned mechanisms $\{\mathcal{M}_{\mathbf{X}}\}$ induced by drawing \mathbf{X} . To ensure the overall budget (10) is at most r , we set a stricter target budget $r' < r$ that holds on "good" sample set \mathbf{X} with small enough estimation error; the fraction of "good" sample set \mathbf{X} can be lower bounded by some threshold $(1 - \gamma)$ from standard concentration inequality. For the remaining, up to γ -fraction (exponentially decaying in the number of samplings m), "bad" sample set \mathbf{X} , we may switch to some conservative worst-case bound R . By joint convexity (Lemma 1) of f -divergences, we can show (10) is deterministically upper bounded by $(1 - \gamma)r' + \gamma R$. Therefore it suffices to choose parameters (r', γ, R) such that $(1 - \gamma)r' + \gamma R \leq r$. We put the above ideas together to establish Algorithm 1, a general framework to optimize privacy-preserving randomization \mathcal{D}_e in the black-box setting with formal guarantees shown in Theorem 1.

Algorithm 1 Randomization \mathcal{D}_e Solver

-
- 1: **Input:** Leakage function \mathcal{F} ; privacy-loss function $\mathcal{G}(\mathcal{D}_e, \cdot)$; input distribution D ; utility loss metric κ , privacy budget r ; sample size m ; global bound $B(\mathcal{D}_e)$.
 - 2: Sample $\mathbf{X} = \{\bar{x}_1, \dots, \bar{x}_m\}$ and $\mathbf{X}' = \{\bar{x}'_1, \dots, \bar{x}'_m\}$, all i.i.d. from D .
 - 3: Choose a target budget r' , a failure probability $\gamma \in (0, 1)$, and a fallback, universal bound R such that

$$(1 - \gamma)r' + \gamma R \leq r.$$

- 4: Solve the constrained optimization problem

$$\min_{\mathcal{D}_e} \kappa(\mathcal{D}_e), \quad \text{such that} \quad \frac{1}{m} \sum_{i=1}^m \mathcal{G}(\mathcal{D}_e, \bar{x}_i) \leq r' \quad (13)$$

to obtain an initial noise distribution \mathcal{D}_e .

- 5: Based on \mathbf{X}' , compute

$$L \leftarrow \frac{1}{m} \sum_{i=1}^m g(\mathcal{D}_e, \bar{x}'_i), \quad \beta \leftarrow B(\mathcal{D}_e) \sqrt{\frac{\log(1/\gamma)}{2m}}.$$

- 6: **while** $L > r' - \beta$ **or** $B(\mathcal{D}_e) > R$ **do**
 - 7: Update $\mathcal{D}_e \leftarrow \mathcal{D}_e + \mathcal{N}(0, c^2 \cdot \mathbf{I})$, by adding an independent Gaussian noise $z \sim \mathcal{N}(0, \sigma^2 \cdot \mathbf{I})$ to the current perturbation solution $\mathbf{e} \sim \mathcal{D}_e$.
 - 8: Recompute L and β as in step 5 for the updated \mathcal{D}_e .
 - 9: **end while**
 - 10: **Output:** Mechanism $\mathcal{M}_{\mathcal{D}_e}$.
-

Theorem 1 (Provable guarantees of Algorithm 1). *Let $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{O}$ be a given leakage function, and $\mathcal{M}_{\mathcal{D}_e}(\cdot) = \mathcal{F}(\cdot) + \mathbf{e}$ denote the privatized mechanism of \mathcal{F} based on strategy $\mathbf{e} \sim \mathcal{D}_e$. For any given \mathcal{D}_e , define the privacy-loss function*

$$\mathcal{G}(\mathcal{D}_e, \mathbf{x}) := \mathcal{D}_\alpha(\mathbb{P}_{\mathcal{M}_{\mathcal{D}_e}(\mathbf{x})} \parallel \mathbb{P}_{W_{\mathcal{D}_e}}),$$

where $\mathbf{x} \in \mathcal{X}$ is the secret input and $\mathbb{P}_{W_{\mathcal{D}_e}}$ is some reference distribution on \mathcal{O} , potentially depending on \mathcal{D}_e . Suppose that for every $\mathbf{x} \in \mathcal{X}$ and every admissible \mathcal{D}_e , we have $\mathcal{G}(\mathcal{D}_e, \mathbf{x}) \leq B(\mathcal{D}_e)$ for some finite bound $B(\mathcal{D}_e)$. Let \mathcal{D}_e be the mechanism returned by Algorithm 1, then we have

$$\mathbb{E}_{\mathbf{x} \sim D}[\mathcal{G}(\mathcal{D}_e, \mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim D}[\mathcal{D}_\alpha(\mathbb{P}_{\mathcal{M}_{\mathcal{D}_e}(\mathbf{x})} \parallel \mathbb{P}_{W_{\mathcal{D}_e}})] \leq r.$$

Algorithm 1 uses two independent sample sets \mathbf{X} and \mathbf{X}' . The first set \mathbf{X} is used to initialize a candidate \mathcal{D}_e (step 4), which with high chance can be asymptotically close to the optimal strategy \mathcal{D}_e^* when we replace the empirical constraint $\frac{1}{m} \sum_{i=1}^m \mathcal{G}(\mathcal{D}_e, \bar{x}_i)$ in (13) by the true population constraint $\mathbb{E}_{\mathbf{x} \sim D} \mathcal{G}(\mathcal{D}_e, \bar{x}_i)$. The second set \mathbf{X}' is used to validate whether this candidate satisfies the desired privacy budget r , which has a bijection relationship to the posterior δ_I from an equality in (8). If the proposed solution \mathcal{D}_e does not pass the test, we will continuously add additional independent noise. This separate application on disjoint

\mathbf{X} and \mathbf{X}' is essential: once \mathcal{D}_e has been chosen based on \mathbf{X} , the random variables $\{\mathcal{G}(\mathcal{D}_e, \bar{x}_i)\}_{i=1}^m$ are no longer independent of the choice of \mathcal{D}_e , so applying concentration inequality (e.g., Hoeffding's) to the same samples would be invalid. Using an independent set \mathbf{X}' restores the required independence and enables a standard concentration bound, thereby establishing the relationship between the confidence level γ and the sample size m .

3.2 Prior-Reference-Weighted (PRW) α -Information and General Inference under Arbitrary Prior

The previous section analyzed identification attacks in the special case where the secret \mathbf{x} is drawn uniformly. We now turn to the general case for arbitrary ρ and prior \mathcal{D} . However, when α -divergence is used as the underlying f -divergence, the bound in Lemma 2 is no longer tight in general. With a closer look at the proof of Proposition 1, we see that when the prior distribution \mathcal{D} is not uniform, the factor $\Pr_{\mathbf{x} \sim \mathcal{D}}(\mathbf{x} = \hat{\mathbf{x}})$ does not appear in the definition of the optimal reference distribution \mathbb{P}_W , which creates a strict gap between the optimal posterior success rate δ_ρ and the guarantee provided by Lemma 2. It is unclear whether there exists another choice of f for which the resulting bound converges to the ground truth for arbitrary priors \mathcal{D} ; we leave this as an interesting direction for future work.

Although the α -divergence version of Lemma 2 is not tight, the MAP rule still holds for an arbitrary inference criteria ρ , where the adversary's optimal strategy to maximize success rate is still to return the guess $\hat{\mathbf{x}}$ with the highest posterior probability, formalized in the following lemma.

Lemma 3 (Maximal success rate under a general inference criterion and prior distribution). *Let \mathbf{x} be the secret with prior $\Pr_{\mathbf{x} \sim \mathcal{D}}(\mathbf{x} = \hat{\mathbf{x}})$, and let \mathcal{F} be a randomized mechanism with output $\mathbf{o} = \mathcal{F}(\mathbf{x})$. Fix an inference criterion $\rho : \mathcal{X} \times \mathcal{X} \rightarrow \{0, 1\}$, where $\rho(\mathbf{x}, \hat{\mathbf{x}}) = 1$ means that the guessing $\hat{\mathbf{x}}$ is counted as success when the true secret is \mathbf{x} . Define the joint weights*

$$H(\mathbf{x}, \mathbf{o}) := \sum_{\hat{\mathbf{x}}' : \rho(\mathbf{x}, \hat{\mathbf{x}}')=1} \Pr(\mathbf{x} = \hat{\mathbf{x}}') \Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o} \mid \mathbf{x} = \hat{\mathbf{x}}').$$

Then, the maximal success probability over all decision rules is

$$\delta_\rho = \int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}) \, d\mathbf{o}, \quad (14)$$

To tightly approach (14), similarly we seek some function taking $\Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o})$ as a building block that converges to $\Pr(\mathbf{x} = \hat{\mathbf{x}}, \mathcal{F}(\mathbf{x}) = \mathbf{o})$ in the appropriate limit. It is observed that the max over x can be viewed as an ℓ_∞ -norm, so, still, it is natural to approximate it by an ℓ_α -norm as $\alpha \rightarrow \infty$. The subtlety is that we need an *upper* bound on $\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})$, whereas the standard normalized ℓ_α expression $(\frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha)^{1/\alpha}$ is a *lower* bound on the $\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})$.

To obtain an upper bound, we instead use

$$\left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha},$$

which clearly upper-bounds $\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})$. Notice that

$$\left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \geq \left(\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} = \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}), \quad (15)$$

since the sum includes the largest term $H(\mathbf{x}^*, \mathbf{o})^\alpha$ for $\mathbf{x}^* \in \arg \max_{\mathbf{x}} H(\mathbf{x}, \mathbf{o})$. Moreover, as $\alpha \rightarrow \infty$, the expression $\left(\sum_{\mathbf{x}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}$ converges to $\max_{\mathbf{x}} H(\mathbf{x}, \mathbf{o})$, as summarized in the following lemma.

Lemma 4 (ℓ_α -relaxation of the max posterior over x). *For any $\alpha > 1$, we have*

$$\int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}) \, d\mathbf{o} \leq \int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \, d\mathbf{o}. \quad (16)$$

Moreover, for every fixed $\mathbf{o} \in \mathcal{O}$,

$$\left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \geq \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}),$$

and the left-hand side converges to $\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})$ as $\alpha \rightarrow \infty$. Hence the right-hand side above is an ℓ_α -norm-based upper bound that becomes asymptotically tight as α grows.

Moreover, for operational efficiency, we further seek an upper bound of (16) above that pulls the summation over x outside the integral, so that the summation can be normalized by $\Pr(\mathbf{x} = \hat{\mathbf{x}})$ and recognized as an expectation over \mathbf{x} . This will allow us to approximate the resulting quantity using a finite sample drawn from \mathcal{D} . Let \mathbb{P} be an arbitrary reference distribution on \mathcal{O} . By Hölder's inequality, we have

Lemma 5 (Prior-reference-weighted α -norm upper bound). *For any reference distribution \mathbb{P}_W on \mathcal{O} and any $\alpha > 1$,*

$$\int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \, d\mathbf{o} \leq \left(\int_{\mathcal{O}} \frac{\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} \, d\mathbf{o} \right)^{1/\alpha} \quad (17)$$

$$= \left(\mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}} \frac{H(\mathbf{x}, \mathbf{o})^\alpha / \Pr(\mathbf{x})}{\Pr(W = \mathbf{o})^{\alpha-1}} \, d\mathbf{o} \right)^{1/\alpha}. \quad (18)$$

Moreover, the equality of (17) holds if and only if the behind Hölder's inequality is tight, i.e., there exists a constant $\lambda > 0$ such that

$$\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha = \lambda \Pr(W = \mathbf{o})^\alpha \quad \text{for almost every } \mathbf{o}.$$

Equivalently, the bound is tight when the reference distribution \mathbb{P}_W is chosen whose probability density is proportional to the α -norm of $H(\cdot, \mathbf{o})$,

$$\Pr(W = \mathbf{o}) \propto \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}.$$

Lemma 4 and Lemma 5 provide two constructions, both of which yield an asymptotically-tight upper bound of the optimal posterior success rate δ_ρ for an arbitrary inference criterion ρ and a prior secret distribution \mathcal{D} . With respect to (16) and (17), we formalize the two constructions as follows, respectively.

Definition 6 (Prior-reference-weighted α -information). Let \mathbf{x} be a random variable taking values in \mathcal{X} , \mathcal{M} be a randomized mechanism with output in \mathcal{O} , and ρ be an inference criterion. For each $\mathbf{x} \in \mathcal{X}$ and $\mathbf{o} \in \mathcal{O}$, define

$$H(\mathbf{x}, \mathbf{o}) := \sum_{\hat{\mathbf{x}}' \in \mathcal{X}: \rho(\mathbf{x}, \hat{\mathbf{x}}')=1} \Pr(\mathbf{x} = \hat{\mathbf{x}}') \Pr(\mathcal{F}(\mathbf{x}) = \mathbf{o} \mid \mathbf{x} = \hat{\mathbf{x}}'),$$

and let W be any reference distribution on \mathcal{O} that is independent of \mathbf{x} . For $\alpha > 1$, the prior-reference-weighted α -information with reference W is defined by

$$I_{\alpha, \rho}^W(\mathbf{x}; \mathcal{F}) := \left(\sum_{\mathbf{x} \in \mathcal{X}} \int_{\mathcal{O}} \frac{H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right)^{\frac{1}{\alpha}}.$$

Equivalently, in an expectation form,

$$I_{\alpha, \rho}^W(\mathbf{x}; \mathcal{F}) = \left(\mathbb{E}_{\mathbf{x}} \left[\frac{1}{\Pr(\mathbf{x})} \int_{\mathcal{O}} \frac{H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right] \right)^{\frac{1}{\alpha}}.$$

Definition 7 (Prior-weighted α -information). We further define the prior-weighted α -information by optimizing over all reference distributions W that are independent of \mathbf{x} :

$$I_{\alpha, \rho}(\mathbf{x}; Y) := \inf_{W \perp \mathbf{x}} I_{\alpha, \rho}^W(\mathbf{x}; Y).$$

As a corollary from Lemma 4 and Lemma 5, we have the following theorems:

Theorem 2 (General upper bound on posterior success). Let \mathbf{x} be a secret, $\mathcal{F}(\mathbf{x})$ an output mechanism, and ρ an inference criterion as in Definition 4. Then, for any $\alpha > 1$ and any reference W ,

$$\delta_\rho \leq I_{\alpha, \rho}^W(\mathbf{x}; \mathcal{F}(\mathbf{x})),$$

and by taking the infimum of W on the both sides of the above inequality, we obtain

$$\delta_\rho \leq \mathsf{I}_{\alpha,\rho}(\mathbf{x}; \mathcal{F}(\mathbf{x})).$$

Since the choice of W is arbitrary, we aim to identify the optimal reference distribution to derive the tight characterization of δ_ρ . Looking back to Lemma 5, we know that the optimal W should satisfy the condition under which the Hölder's inequality behind (17) is tight. This leads to the following result.

Theorem 3 (Optimal reference selection). *Let \mathbf{x} be the secret, a leakage function \mathcal{F} with output space \mathcal{O} , and $H(\mathbf{x}, \mathbf{o})$ as in Definition 6. For $\alpha > 1$, consider $\mathsf{I}_{\alpha,\rho}^W(\mathbf{x}; \mathcal{F})$ as \mathbb{P}_W ranges over all reference distributions on \mathcal{O} independent of \mathbf{x} . Then the infimum over such W is attained (by Hölder's inequality) at*

$$\Pr(W = \mathbf{o}) \propto \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha},$$

i.e.,

$$\Pr(W = \mathbf{o}) = \frac{\left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}}{\int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}')^\alpha \right)^{1/\alpha} d\mathbf{o}'}.$$

For this optimal choice W^* , the prior-reference-weighted α -information satisfies

$$\mathsf{I}_{\alpha,\rho}(\mathbf{x}; \mathcal{F}) = \int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} d\mathbf{o},$$

where $\mathsf{I}_{\alpha,\rho}(\mathbf{x}; \mathcal{F})$ denotes the infimum of $\mathsf{I}_{\alpha,\rho}^W(\mathbf{x}; \mathcal{F})$ over all such reference distributions W .

Finally, we prove the tightness of prior-weighted α -information, which asymptotically approaches the target δ_ρ .

Theorem 4 (Asymptotic tightness of weighted α information).

$$\lim_{\alpha \rightarrow \infty} \mathsf{I}_{\alpha,\rho}(\mathbf{x}; \mathcal{F}(\mathbf{x})) = \int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathcal{F}(\mathbf{x}) = \mathbf{o}) d\mathbf{o} = \delta_\rho$$

Finally, it is observed that

$$\mathsf{I}_{\alpha,\rho}^W(\mathbf{x}; \mathcal{F})^\alpha = \left(\mathbb{E}_{\mathbf{x}} \left[\frac{1}{\Pr(\mathbf{x})} \int_{\mathcal{O}} \frac{H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right] \right)$$

is still an expectation with respect to the secret $\mathbf{x} \sim \mathcal{D}$. Consequently, it can be approximated by Monte Carlo estimation using finitely many i.i.d. samples from \mathcal{D} . This leads to a concrete sample-based procedure for enforcing the target posterior success rate, which we formalize in the following theorem.

Theorem 5 (Sample based tight privacy guarantee). *Given a function $\mathcal{F} : \mathcal{X} \rightarrow \mathcal{O}$, let $\mathcal{M}_{\mathcal{D}_e}$ denote the mechanism that applies \mathcal{F} and then adds noise drawn from a distribution \mathcal{D}_e on \mathcal{O} . For each noise distribution \mathcal{D}_e and each $\mathbf{x} \in \mathcal{X}$, define the measure*

$$H_{\mathcal{D}_e}(\mathcal{M}_{\mathcal{D}_e}(\mathbf{x}) = \mathbf{o}) := \sum_{\rho(\mathbf{x}, \hat{\mathbf{x}}')=1} \Pr(\mathbf{x} = \hat{\mathbf{x}}') \Pr(\mathcal{M}_{\mathcal{D}_e}(\hat{\mathbf{x}}') = \mathbf{o}), \quad (19)$$

and let $W_{\mathcal{D}_e}$ be a reference distribution on \mathcal{O} . Define the privacy-loss function

$$\mathcal{G}(\mathcal{D}_e, \hat{\mathbf{x}}) := \frac{1}{\Pr(\mathbf{x} = \hat{\mathbf{x}})} D_{\alpha}(P_{H_{\mathcal{D}_e}(\hat{\mathbf{x}})} \| P_{W_{\mathcal{D}_e}}),$$

Suppose that for every $\mathbf{x} \in \mathcal{X}$ and every admissible \mathcal{D}_e we have $\mathcal{G}(\mathcal{D}_e, \hat{\mathbf{x}}) \leq B(\mathcal{D}_e)$ for some finite bound $B(\mathcal{D}_e)$, and an optimizer searching over \mathcal{D}_e to minimize some metric $\kappa(\mathcal{D}_e)$ subject to the constraint

$$\frac{1}{m} \sum_{i=1}^m \mathcal{G}(\mathcal{D}_e, \bar{\mathbf{x}}_i) \leq r,$$

for samples $\{\bar{\mathbf{x}}_i\}_{i=1}^m \sim D^m$. To make the mechanism satisfy δ_{ρ} privacy, set the privacy budget to

$$r = (\delta_{\rho})^{\alpha}.$$

Then the privatization mechanism $\mathcal{M}_{\mathcal{D}_e}$ returned by Algorithm 1 guarantees that the true posterior success rate under the inference criterion ρ is at most δ_{ρ} .

In the limit where the sample size $m \rightarrow \infty$, the order $\alpha \rightarrow \infty$, and $W_{\mathcal{D}_e}$ is chosen optimally, the budget choice $r = \delta_{\rho}^{\alpha}$ becomes asymptotically tight: the bound converges to the true posterior success rate under ρ , and \mathcal{D}_e is the optimal noise distribution which minimizes $\kappa(\mathcal{D}_e)$ under the privacy budget constraint.

Before concluding this section, we point out a difference between Theorem 5 and Lemma 2. For an identification problem with a uniformly generated secret \mathbf{x} , both bounds are asymptotically tight, but for finite $\alpha < \infty$ the bound given by Lemma 2 is slightly better than proposed Theorem 5. Notice that

$$D_{\alpha}(\mathbf{1}_{\delta_{\rho}} \| \mathbf{1}_{\delta_{o,\rho}}) = (1 - \delta_{o,\rho}) \left(\frac{1 - \delta_{\rho}}{1 - \delta_{o,\rho}} \right)^{\alpha} + \delta_{o,\rho} \left(\frac{\delta_{\rho}}{\delta_{o,\rho}} \right)^{\alpha} > \delta_{o,\rho} \left(\frac{\delta_{\rho}}{\delta_{o,\rho}} \right)^{\alpha}$$

When \mathbf{x} is uniform, we have $\delta_{o,\rho} = \Pr(\mathbf{x} = \hat{\mathbf{x}}) = 1/|\mathcal{X}|$, so the inequality (20) above directly implies

$$\delta_{\rho} \leq \left(\mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}} \frac{\Pr(\mathbf{x} = \hat{\mathbf{x}})^{\alpha-1} \Pr(\mathcal{M}(\mathbf{x}) = \mathbf{o})^{\alpha}}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right)^{\frac{1}{\alpha}}, \quad (20)$$

which is exactly the bound we obtain in Theorem 5. In the high-entropy regime with large α , the second term $\delta_{o,\rho} \left(\frac{\delta_{\rho}}{\delta_{o,\rho}} \right)^{\alpha}$ generally dominates and thus non-asymptotically the difference between Lemma 2 and Theorem 5 is minor. However, our result (Theorem 5) applies to *arbitrary* priors and general inference criteria ρ , while remaining asymptotically tight.

4 History-Dependent Adaptive Mechanisms

In this section, we address Problem 2 to control cumulative privacy risk in the presence of an *adaptive* adversary. We assume the user has access to the full interaction history, including all previous outputs, the adversary's queries, and the mechanisms realized in earlier rounds. We focus on the identification problem and our results are generalizable to general inference ρ using a similar reasoning as shown in Section 3.2.

At round t , the adversary chooses the next query based on the entire interaction history, i.e., its past queries and the observed outputs. Such adaptivity can cause significantly more privacy leakage than in the static (non-adaptive) model. Instead of committing to all queries in advance, the adversary can repeatedly probe the most *vulnerable* aspects of the secret revealed by earlier outputs, and thereby tries to extract more information from the interaction.

As discussed previously, the black-box nature of both the leakage function and the adversarial strategy \mathcal{Q}_{adv} presents a "parallel universe" challenge (Fig. 1). Because the user observes only a single realized transcript among infinite possibilities, we cannot simply simulate the combination of leakage function $\mathcal{F}(\cdot, \mathbf{v})$ and \mathcal{Q}_{adv} and apply the results from Section 3. However, one key insight is that *before* the interactive query-response process begins, once a privacy-preserving mechanism Alg is determined by the user, although it is infeasible for the user to either predict or simulate *all* possible future transcripts, a same rule determined by Alg will be consistently applied to privatize the leakage in every branch of the interaction tree. As a result, it suffices to prove a *history-dependent* guarantee that holds at each step conditioned on the past; once such a guarantee is established, it automatically lifts to a global guarantee over the entire adaptive interaction. Before proceeding, we first formalize the adaptive interaction model introduced in Problem 2.

Definition 8 (Adaptive Adversarial Composition).

Spaces and transcripts. Given a total iteration number T and three spaces: \mathcal{X} is the secret space, \mathcal{V} is the query space from which the adversary selects the query at each iteration, and \mathcal{O} is the leakage space. A secret $\mathbf{x} \sim \mathcal{D}$ is sampled once at the beginning and reused throughout all T rounds. For each round t , define the pre-round- t transcript space as

$$\mathcal{T}_t := \mathcal{V}^{t-1} \times \mathcal{O}^{t-1}.$$

Accordingly, we write the realized pre-round- t transcript as

$$\tau_t := (\mathbf{v}_{<t}, \mathbf{o}_{<t}) \in \mathcal{T}_t, \quad \text{here} \quad \mathbf{v}_{\leq t} = (\mathbf{v}_1, \dots, \mathbf{v}_{t-1}), \quad \mathbf{o}_{<t} = (\mathbf{o}_1, \dots, \mathbf{o}_{t-1}),$$

to record the adversary's past queries and the user's past outputs. Let $\rho : \mathcal{X} \times \mathcal{X} \rightarrow \{0, 1\}$ be an inference criterion, where $\rho(\hat{\mathbf{x}}, \mathbf{x}) = 1$ indicates that the adversary's estimate $\hat{\mathbf{x}}$ is deemed successful for the true secret \mathbf{x} .

Leakage and randomization. Given T processing functions

$$\mathcal{F}_t : \mathcal{X} \times \mathcal{V} \rightarrow \mathcal{O}, \quad t \in [T],$$

and T transcript-dependent randomization rules (corresponding to each step of Alg in Problem 1 and 2)

$$\Sigma_t : \mathcal{V} \times \mathcal{T}_t \rightarrow \mathcal{O}, \quad t \in [T],$$

At round t , the user combines \mathcal{F}_t and Σ_t to form a randomized mechanism

$$\mathcal{M}_t = \mathcal{F}_t + \Sigma_t.$$

Adaptive interaction experiment. An adversary's strategy \mathcal{Q}_{adv} is a sequence of algorithm which can be either deterministic or randomized.

$$\mathcal{Adv}_t : \mathcal{T}^t \rightarrow \mathcal{V}, \quad t \in [T].$$

The interaction proceeds for $t = 1, 2, \dots, T$ as follows:

1. Given \mathcal{T}_t , the adversary outputs $\mathbf{v}_t \leftarrow \mathcal{Adv}_t$.
2. Given the secret $\hat{\mathbf{x}}$, the current query \mathbf{v}_t , and the transcript τ_t , the user and returns $\mathbf{o}_t \leftarrow \mathcal{M}_t$ to the adversary.

At the end of the interaction, the adversary outputs an estimate $\hat{\mathbf{x}} \in \mathcal{X}$ of the secret \mathbf{x} , and we use \mathcal{M}^{Adv} to represent the aggregation of T adaptively-dependent leakage functions.

In Definition 8, the per-round mechanism \mathcal{M}_t may look complicated, but it admits a simple interpretation. Conditioned on the current query \mathbf{v}_t and the current transcript τ_t , the user's response is obtained by applying \mathcal{F}_t and then adding noise chosen as a function of (\mathbf{v}_t, τ_t) . In particular, if the randomization rule $\Sigma_t(\mathbf{v}_t, \tau_t)$ outputs a Gaussian distribution, then *conditionally on* (\mathbf{v}_t, τ_t) the mapping $\mathbf{x} \mapsto \mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t)$ is simply a Gaussian mechanism.

Our next goal is to design transcript-dependent randomization rules $\{\Sigma_t\}_{t=1}^T$ so that, for every adversary, the resulting composed interaction guarantees that the posterior success rate is bounded by the target level. To do so, we need a bridge between the induced output distributions and the adversary's posterior success rate. Lemma 2 provides one, but as discussed earlier, its α -divergence specialization is not tight under non-uniform priors. Therefore, we instead work with Theorem 2.

At first glance, Theorem 2 appears to apply only to a *non-interactive* mechanism, since it does not explicitly model adversarially chosen queries. The key observation is that this is without loss of generality: we may equivalently regard the adversary's strategy Adv as part of the overall randomized mechanism. Concretely, imagine that the adversary hands the user a "strategy machine" implementing $\{\mathcal{Adv}_t\}_{t=1}^T$ as a black box. At round t , the user first runs this machine on the current transcript to generate the next query \mathbf{v}_t , and then runs \mathcal{M}_t on

$(\mathbf{x}, \mathbf{v}_t, \tau_t)$ to generate the response \mathbf{o}_t . At the end, the adversary receives the full transcript $(\mathbf{v}_{1:T}, \mathbf{o}_{1:T})$ and outputs an estimate $\hat{\mathbf{x}}$.

This reformulation is equivalent to the original interaction: the joint distribution of the transcript $(\mathbf{v}_{1:T}, \mathbf{o}_{1:T})$ is unchanged, and the adversary has access to exactly the same information when forming $\hat{\mathbf{x}}$. Consequently, any posterior-success bound that holds for the resulting randomized mapping from \mathbf{x} to the final transcript applies equally to the original adaptive interaction.

For the reference distribution associated with the adversary's strategy $\{\mathcal{A}dv_t\}$, we may simply take an independent copy of the same strategy, since $\mathcal{A}dv_t$ is independent of the secret \mathbf{x} , which is the only requirement we impose on the reference. This leads to the theorem below.

Theorem 6 (Posterior success bound for adaptive adversary). *Under the setting of Definition 8, for any $\alpha > 1$, we have*

$$(\delta_\rho)^\alpha \leq \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \Pr(\mathbf{x})^{\alpha-1} \cdot \prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} \\ \cdot \prod_{t=1}^T \Pr(\mathcal{A}dv_t(\tau_t) = \mathbf{v}_t) \, d\mathbf{o}_{[1:T]} \, d\mathbf{v}_{[1:T]}$$

Here, we may take the reference distribution \mathbb{P}_W for the adversary's queries to be the true law of $\mathcal{A}dv_t(\tau_t)$ itself, since it is independent of \mathbf{x} conditionally on τ_t .

There is still a gap between Theorem 6 and an implementable procedure. As discussed above, the user observes only a single realized interaction (one “universe”) induced by the adversary's adaptivity and thus lacks the information of $\Pr(\mathcal{A}dv_t(\tau_t) = \mathbf{v}_t)$. So the right-hand side of Theorem 6 is still not computable. We now return to the key insight mentioned at the beginning of this section. We consider further decomposing the global control problem in Theorem 6 into a collection of *local* subproblems, each of which depends only on the current branch (i.e., only on the currently observed (\mathbf{v}_t, τ_t) and the local mechanism at round t), then each “self” in each branch can solve its own local subproblem using only the information available in that branch. Since the same local rule is executed in all branches, combining these locally enforced constraints yields a global control on the full adaptive interaction. The next Theorem 7 shows a way to decompose and since the local control can lead to a global privacy guarantee.

Theorem 7 (Adaptive Composition of α -divergence). *Under the setting of Definition 8 with a single secret $\mathbf{x} \sim D$ reused across all T rounds. For any possible query sequence, consider an auxiliary reference process*

$$W_{1:T}(\mathbf{v}_{1:T}, \tau_{1:T}) := (W_1(\mathbf{v}_1, \tau_1), \dots, W_T(\mathbf{v}_T, \tau_T)) \in \mathcal{O}^T$$

which is independent of \mathbf{x} . For each round $t \in \{1, \dots, T\}$ and transcript $\tau_t = (\mathbf{v}_{<t}, \mathbf{o}_{<t})$, define the prefix coefficient

$$A_{t-1}(\hat{\mathbf{x}}) := \Pr(\mathbf{x} = \hat{\mathbf{x}})^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{\mathbf{x}}, \mathbf{v}_s, \tau_s) = \mathbf{o}_s)}{\Pr(W_s(\mathbf{v}_s, \tau_s) = \mathbf{o}_s)} \right)^\alpha$$

and

$$D_t(\hat{x}) := \mathcal{D}_\alpha(\mathcal{M}_t(\hat{x}, \mathbf{v}_t, \tau_t) \parallel W_t(\mathbf{v}_t, \tau_t)).$$

If for every possible transcript τ_t there exists such a choice of reference process $W_{1:T}(\mathbf{v}_{1:T})$ such that for all $t = 1, \dots, T$,

$$\mathbb{E}_{\mathbf{x} \sim D} [A_{t-1}(\mathbf{x}) D_t(\mathbf{x})] \leq r_t \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_{t-1}(\mathbf{x})], \quad (21)$$

then the final interactive mechanism \mathcal{M}^{Adv} satisfies for any \mathcal{Q}_{adv} adversarial strategy, the optimal posterior success rate δ_ρ satisfies

$$(\delta_\rho)^\alpha \leq \prod_{t=1}^T r_t. \quad (22)$$

Theorem 7 makes adaptive randomization implementable in practice. Once the adversary's past queries and the user's past outputs have been realized, the transcript τ_t is fixed. The user therefore only needs to enforce the per-round condition (21) along the single interaction path that actually occurs, rather than over exponentially many counterfactual transcripts.

Note that the left-hand side of (21) is a *weighted* average of per-secret α -divergences. Intuitively, $D_t(\hat{x})$ quantifies how much information about the particular secret \hat{x} could be revealed at round t , while the prefix coefficient $A_{t-1}(\hat{x})$ reweighs secrets according to how plausible they remain under the realized transcript. Consequently, secrets with larger $A_{t-1}(\hat{x})$ contribute more to the cumulative leakage and must be protected more aggressively.

Finally, if \mathcal{D}_e is simply some unrestricted noise distribution, for example, a Gaussian, the existence of a randomization solution is also guaranteed. As the variance of \mathbf{e} grows, the output distribution of $\mathcal{M}_t(\cdot, \mathbf{v}_t, \tau_t)$ becomes increasingly insensitive to the secret, and hence $D_t(\hat{x})$ decreases toward its minimum value. In the limit of infinite variance, the left-hand side of (21) approaches the right-hand side, so the inequality can always be met by choosing sufficiently large variance whenever $r_t > 1$.

Remark 2 (Multiplicative vs. additive budgets for α -divergence). In our setting we work exclusively with α -divergence, and the per-round budget *can be expressed* either multiplicatively or additively. When the α -divergence bound is unrolled across T rounds, the integrand contains a product of per-round likelihood ratios raised to the power α . For this reason, it is often convenient to allocate per-round factors $\{r_t\}_{t=1}^T$ so that the overall budget takes the multiplicative form

$$r_{\text{tot}} = \prod_{t=1}^T r_t.$$

At the same time, an additive accounting is equally valid: defining $\varepsilon_t := \log r_t$ yields

$$\log r_{\text{tot}} = \sum_{t=1}^T \log r_t = \sum_{t=1}^T \varepsilon_t,$$

This log-additive view mirrors the familiar additive composition rule for Rényi-type privacy losses.

In the following, we further show the *tightness* of the privatization protocol derived by Theorem 7 for adversarially-adaptive composition: with proper selection of the reference process $W_{1:T}$, the equality of (22) is achievable. That is to say, given the iterative budget $\{r_1, r_2, \dots, r_T\}$, our privacy solution *exactly* reaches the global budget $\prod_{t=1}^T r_t$.

Since the user can only choose the reference distribution W_t *locally* at each round t based on the realized transcript, given a per-round budget allocation $r_{1:T}$, a natural construction of the reference process $W_{1:T}$ is through a greedy design: W_t is optimized to tighten the bound for the current step. Apparently, this greedy solution may not be jointly optimal for the *entire* interactive transcript but somewhat surprisingly, such locally-optimal choice indeed realizes the globally optimal reference process, as formalized in the following theorem.

Theorem 8 (Local selection realizes the optimal reference process W). *With the same notation of Theorem 7, define the per-round locally optimal reference W_t (for each fixed realized (\mathbf{v}_t, τ_t)) by*

$$\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t) \propto \left(\mathbb{E}_{\mathbf{x} \sim D} \left[\Pr_D(\mathbf{x})^{\alpha-1} \left(A_{t-1}(\mathbf{x}) \cdot \Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t) \right) \right] \right)^{1/\alpha}. \quad (23)$$

Assume that the cumulative condition (21) holds with equality at every round. Then the resulting composed reference process $W_{1:T}$ is optimal which means it minimizes the upper bound in Lemma 2 among all admissible reference processes independent of \mathbf{x} . Consequently, when $\alpha \rightarrow \infty$, the resulting posterior-success bound is asymptotically tight.

Finally, to further make Theorem 7 implementable from finite sampling, we establish Algorithm 2 by iteratively running Algorithm 1 based on locally-constructed privacy loss \mathcal{G}_t defined in (95). We provide the provable guarantees of Algorithm 2 in Theorem 9.

Theorem 9 (Provable guarantees of Algorithm 2). *Define the prefix coefficient*

$$A_{t-1}(\hat{\mathbf{x}}) := \Pr(\mathbf{x} = \hat{\mathbf{x}})^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{\mathbf{x}}, \mathbf{v}_s, \tau_s) = \mathbf{o}_s)}{\Pr(W_s(\mathbf{v}_s, \tau_s) = \mathbf{o}_s)} \right)^\alpha,$$

$$D_t(\hat{\mathbf{x}}) := \mathcal{D}_\alpha(\mathcal{M}_t(\hat{\mathbf{x}}, \mathbf{v}_t, \tau_t) \parallel W_t(\mathbf{v}_t, \tau_t)),$$

and the privacy-loss function

$$\mathcal{G}_t(\mathcal{D}_e, \hat{\mathbf{x}}) := (A_{t-1}(\hat{\mathbf{x}}) (D_t(\hat{\mathbf{x}}) - r_t)) \quad (24)$$

Algorithm 2 Randomization Optimizer under Adversarial Composition

- 1: **Input:** number of rounds T ; leakage functions $\{\mathcal{F}_t : \mathcal{X} \times \mathcal{V} \rightarrow \mathcal{O}\}_{t=1}^T$; input distribution D ; per-round α -divergence budgets $\{r_t\}_{t=1}^T$; sample size m ; global bound $B(\mathcal{D}_e)$.
- 2: Sample a secret input $\hat{x} \sim D$.
- 3: Choose per-round targets $\{r'_t\}_{t=1}^T$, failure probabilities $\{\gamma_t\}_{t=1}^T$, and a fallback bound R such that for every $t \in [T]$,

$$(1 - \gamma_t) r'_t + \gamma_t R \leq r_t.$$

- 4: **for** $t = 1, 2, \dots, T$ **do**
- 5: Receive adversarial parameter v_t .
- 6: Define the *prefix coefficient*

$$A_{t-1}(\hat{x}) := \Pr(x = \hat{x})^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{x}, v_s, \tau_s) = o_s)}{\Pr(W_s(v_s, \tau_s) = o_s)} \right)^\alpha$$

$$D_t(\hat{x}) := \mathcal{D}_\alpha(\mathcal{M}_t(\hat{x}, v_t, \tau_t) \parallel W_t(v_t, \tau_t)),$$

and

$$\mathcal{G}_t(\mathcal{D}_e, \hat{x}) := A_{t-1}(\hat{x}) D_t(\hat{x})$$

- 7: Take $\mathcal{F}(\cdot, v_t), g_t(\mathcal{D}_e, \hat{x}), D, r_t, B(\mathcal{D}_e)$ as the input for Algorithm 1
 - 8: Sample noise e_t from the calibrated distribution.
 - 9: Output $o_t \leftarrow \mathcal{F}_t(\hat{x}, v_t) + e_t$.
 - 10: **end for**
-

Suppose that for every $\hat{x} \in \mathcal{X}$ and every admissible \mathcal{D}_e we have $\mathcal{G}_t(\mathcal{D}_e, \hat{x}) \leq B(\mathcal{D}_e)$ for some finite bound $B(\mathcal{D}_e)$ given per-round budget 0 (It's because the the second part of the privacy loss function minus r_t), Let $\mathcal{M}_{[1:T]}$ be the interactive mechanism produced by Algorithm 2. Then for any Adv , the optimal posterior success rate δ_ρ satisfies

$$\delta_\rho^\alpha \leq \prod_{t=1}^T r_t.$$

Finally, we give an example of how $B(\mathcal{D}_e)$ can be computed. Assume that $\|\mathcal{F}\|_2 \leq b$ for all $\hat{x} \in \mathcal{X}$ and $v \in \mathcal{V}$. Consider adding Gaussian noise, and define the (empirical) local reference

$$W \propto \left(\frac{1}{m} \sum_{i=1}^m \Pr(x = \bar{x}_i)^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\bar{x}_i, v_s, \tau_s) = o_s)}{\Pr(W_s(v_s, \tau_s) = o_s)} \cdot \Pr(\mathcal{M}_t(\bar{x}_i, v_t, \tau_t) = o_t) \right)^\alpha \right)^{\frac{1}{\alpha}}.$$

The quantity A_{t-1} can be bounded along the realized transcript: the denominator terms $\Pr(W_s(v_s, \tau_s) = o_s)$ are fixed once the transcript is fixed, and the numerator terms admit an upper bound under the assumption $\|\mathcal{F}\|_2 \leq b$.

Therefore the normalizing factor is bounded. In particular,

$$\begin{aligned}
& \int_{\mathcal{O}} \left(\frac{1}{m} \sum_{i=1}^m \Pr(\mathbf{x} = \bar{\mathbf{x}}_i)^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\bar{\mathbf{x}}_i, \mathbf{v}_s, \tau_s) = \mathbf{o}_s)}{\Pr(W_s(\mathbf{v}_s, \tau_s) = \mathbf{o}_s)} \cdot \Pr(\mathcal{M}_t(\bar{\mathbf{x}}_i, \mathbf{v}_t, \tau_t) = \mathbf{o}) \right)^\alpha \right)^{\frac{1}{\alpha}} \text{ do} \\
& \leq \frac{\max_{\hat{\mathbf{x}} \in \mathcal{X}} A_{t-1}(\hat{\mathbf{x}})}{\min_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathbf{x} = \hat{\mathbf{x}})} \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathbf{x} = \hat{\mathbf{x}}) \Pr(\mathcal{M}_t(\hat{\mathbf{x}}, \mathbf{v}_t, \tau_t) = \mathbf{o}) \text{ do} \\
& \leq \frac{\max_{\hat{\mathbf{x}} \in \mathcal{X}} A_{t-1}(\hat{\mathbf{x}})}{\min_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathbf{x} = \hat{\mathbf{x}})}. \tag{25}
\end{aligned}$$

Hence the normalizing constant is bounded. Next, by Hölder's inequality we have

$$\frac{1}{\left(\frac{1}{m} \sum_{i=1}^m a_i^\alpha \right)^{\frac{\alpha-1}{\alpha}}} \leq \sum_{i=1}^m \frac{1}{m a_i^{\alpha-1}}. \tag{26}$$

Omitting (\mathbf{v}_t, τ_t) since they are fixed, then for any $\hat{\mathbf{x}} \in \mathcal{X}$, we obtain

$$\int_{\mathcal{O}} \Pr(\mathcal{M}_t(\hat{\mathbf{x}}) = \mathbf{o})^\alpha \Pr(W_t = \mathbf{o})^{1-\alpha} \text{ do} \leq C^\alpha \cdot \max_{\hat{\mathbf{x}}' \in \mathcal{X}} \mathcal{D}_\alpha(\mathcal{M}_t(\hat{\mathbf{x}}) \parallel \mathcal{M}_t(\hat{\mathbf{x}}')), \tag{27}$$

where C is a normalization constant, and therefore the privacy-loss function is bounded.

5 Mechanism Agnosticism Adaptive Composition

In this section, we further study the composition in the *mechanism-agnostic* setup. In this setting, across each round $t \in [T]$, a mechanism \mathcal{M}_t must determine its internal randomization using only the information restricted to the privacy budget allocated to that round and the specifically observed leakage function $\mathcal{F}(\cdot, \mathbf{v}_t)$. Formally speaking, the difference between the model in this section and Definition 8 is that τ_t is not allowed to be an input to the mechanism.

A key question we want to address in this section is that under what kind of separate measurement of marginal leakage, they can be finally aggregated form an upper bound of the cumulative privacy risk, independent of any potential adversarial strategy that controls these leakage functions behind. In the following we present two kinds of measurements established by KL-divergence and α -divergence that enables the *mechanism-agnostic* composition, as characterized in the following two theorems.

Theorem 10 (KL agnostic composition). *For each $t \in [T]$, any \mathbf{v}_t , and an arbitrary $W_t(\mathbf{v}_t)$ as reference, if both the following worse-case KL divergence and the average-case KL divergence are bounded:*

$$\mathcal{D}_{KL}(P_{\mathcal{M}_t(\hat{\mathbf{x}}, \mathbf{v}_t)} \parallel P_{W_t(\mathbf{v}_t)}) \leq R_t, \quad \forall \hat{\mathbf{x}} \in \mathcal{X}, \mathbf{v}_t \in \mathcal{V}_t \tag{28}$$

$$\mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t)} \parallel P_{W_t(\mathbf{v}_t)}) \leq r_t, \quad \forall \mathbf{v}_t \in \mathcal{V}_t \tag{29}$$

then the optimal posterior success rate δ_ρ of the final T -aggregated interactive mechanism \mathcal{M}^{Adv} for any Adv satisfies

$$\mathcal{D}_{KL}(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{o,\rho}}) \leq B_T,$$

where B_t is defined in a recursive form

$$B_1 = r_1, \quad (30)$$

$$B_t = B_{t-1} + \min\{\sqrt{2B_{t-1}}R_t + r_t, R_T\}, t \geq 2. \quad (31)$$

In Theorem 10, R_t behaves as a universal, input-independent upper bound. From (31), the increment between $B(t)$ and $B(t-1)$ in two consecutive iterations is upper bounded by either $\sqrt{2B(t-1)}R_t + r_t$, that is, the average increment r_t plus a worst-case increment R_t scaled by a factor $\sqrt{2B(t-1)}$ depending on the current accumulated risk, or simply by the worst-case increment R_t . Asymptotically, B_T increases quadratically in T .

We also show the result through α -divergence in the following theorem:

Theorem 11 (α agnostic composition). *For each $t \in [T]$, any \mathbf{v}_t , and an arbitrary $W_t(\mathbf{v}_t)$ as reference, if the α -divergence between $\mathbb{P}_{\mathcal{M}_t(\cdot, \mathbf{v}_t)}$ and $\mathbb{P}_{W_t(\mathbf{v}_t)}$ is bounded as*

$$\mathbb{E}_{\mathbf{x} \sim D} \Pr(\mathbf{x})^{p_t-1} D_{p_t}(\mathbb{P}_{\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t)} \parallel \mathbb{P}_{W_t(\mathbf{v}_t)}) \leq r_t, \quad \forall \mathbf{v}_t \in \mathcal{V}_t, \quad (32)$$

then the optimal posterior success rate δ_ρ of the final T -aggregated interactive mechanism \mathcal{M}^{Adv} for any Adv satisfies

$$\delta_\rho^\alpha \leq \prod_{t=1}^T r_t^{\frac{1}{p_t}}.$$

Here $p_{[1:t]}$ are positive real numbers satisfying the following conditions:

$$\sum_{i=1}^T \frac{1}{p_i} = 1 \quad (33)$$

Compared with KL-agnostic composition in Theorem 10 which requires a global worst-case guarantee R_t , in a Hölder's inequality style α -agnostic composition in Theorem 11 does not explicitly require such a worst-case guarantee but only normal marginal α -divergence. However, when the iteration number T is large, the corresponding α -divergence control likewise becomes effectively dominated by the $\max_t r_t$. This is because $\sum_{i=1}^T \frac{1}{p_i} = 1$, which implies that there must exist some $p_t \geq T \rightarrow \infty$. As expected, compared to the history-dependent composition in Section 4, the mechanism-agnostic one forces the user to add larger noise per-iteration, since the adversary can concentrate its attack on the most vulnerable region of the secret space, while the user cannot adapt the mechanism to counteract that vulnerability along the realized branch.

6 Experiments and Additional Related Work

To demonstrate the practical utility of our theory and algorithms, we evaluate the protection of secret keys in two realistic side-channel scenarios: power consumption leakage in an AES key generator and timing leakage in RSA encryption. We consider both static and composition settings. In our experiments, we model the privatization mechanism \mathbf{e} as independent Gaussian noise. We analyze how the optimal Gaussian distribution \mathcal{D}_e (characterized by minimal variance) is influenced by entropy (varying secret bit-lengths), prior success rate (across different inference tasks, including identification and approximate reconstruction), and the composition number T . Detailed experimental results are provided in Appendix P.

In Appendix Q, we derive analogous results based on Kullback-Leibler (KL) divergence, serving as variants to our α -divergence and prior-reference-weighted α -information bounds. While KL divergence is a standard metric in information privacy, we show that these variants, despite their simpler forms, yield bounds on the posterior rate δ_ρ that can be exponentially looser. This gap is illustrated in Figures 2–5 in Appendix P and is particularly pronounced in high-entropy regimes.

Finally, we discuss additional related work.

Looseness of KL-divergence and Mutual Information. Mutual Information (MI) is often the de facto metric for quantifying the leakage of entropic secrets. However, independent of our findings regarding the looseness of posterior inference bounds, existing literature also identifies other limitations of MI. For instance, MI generally captures average leakage rather than a high-probability sense [38] or the worst case [4]. Furthermore, it may not accurately reflect the empirical success rate of practical attackers [40].

Leakage-Resilient Cryptography. In contrast to our approach, which applies a universal randomization mechanism to any (potentially black-box) leakage, Leakage-Resilient Cryptography [31] focuses on redesigning cryptographic primitives to mitigate leakage directly. These methods typically rely on specific assumptions about the leakage channel, such as the t -probing model (where an adversary observes at most t wires) [17] or the bounded-leakage model [12]. Common techniques include masking and periodic key updates. Hardware- and software-level countermeasures like constant-time encryption [24] also exist, they are often computationally expensive or platform-specific.

7 Conclusion and Future Prospects

In this paper, we focus on entropic secrets and establish an asymptotically-tight Bayesian analysis framework to optimize privacy-preserving solutions for both static and adversarially-adaptive composition settings given black-box leakage functions. In contrast to the strong impossibility results for universal privatization in the regime of *Input-Independent Indistinguishability* (III) for the prior-free worst case, we demonstrate that secret entropy can be tightly exploited to

amplify privacy, enhancing adversarial inference hardness even in a black-box manner with probable guarantees.

A straightforward and promising future direction for our results is software/hardware co-design for side-channel leakage mitigation. Given a specific leakage channel (e.g., power or timing consumption from a hardware environment) and an objective task (e.g., encryption), tight black-box privatization enables the free exploration of different processing functions (e.g., cryptographic protocols) to identify the sharpest utility-privacy tradeoff.

Several technical questions also remain open within our framework. For example, it is unclear whether a more powerful f -function exists, improving upon the polynomial $f(z) = z^\alpha$, that can yield faster and more stable convergence to δ_ρ . Additionally, from an operational perspective, we have not systematically explored the construction but the just assumed the existence of numerical solvers for the constrained optimization problem (13) in Algorithm 1.

Finally, regarding mechanism-agnostic composition, we note that computing tight composition for III-style worst-case guarantees is known to be $\#P$ -hard [28]. A natural open problem is whether this complexity result also holds for entropic secrets. Furthermore, beyond the KL- and α -divergence-based marginal privacy risk measurements (Section 5), exploring other constructions that enable such mechanism-agnostic composition remains an interesting avenue for future work.

References

1. Abadi, M., Chu, A., Goodfellow, I., McMahan, H.B., Mironov, I., Talwar, K., Zhang, L.: Deep learning with differential privacy. In: Proceedings of the 2016 ACM SIGSAC conference on computer and communications security. pp. 308–318 (2016)
2. Abouelmehdi, K., Beni-Hessane, A., Khaloufi, H.: Big healthcare data: preserving security and privacy. *Journal of big data* **5**(1), 1–18 (2018)
3. Aga, S., Narayanasamy, S.: Invisimem: Smart memory defenses for memory bus side channel. *ACM SIGARCH Computer Architecture News* **45**(2), 94–106 (2017)
4. Alvim, M.S., Andrés, M.E., Chatzikokolakis, K., Degano, P., Palamidessi, C.: Differential privacy: on the trade-off between utility and information leakage. In: International workshop on formal aspects in security and trust. pp. 39–54. Springer (2011)
5. Bun, M., Steinke, T.: Concentrated differential privacy: Simplifications, extensions, and lower bounds. In: Theory of cryptography conference. pp. 635–658. Springer (2016)
6. Coppens, B., Verbauwhede, I., De Bosschere, K., De Sutter, B.: Practical mitigations for timing-based side-channel attacks on modern x86 processors. In: 2009 30th IEEE symposium on security and privacy. pp. 45–60. IEEE (2009)
7. Dodis, Y., Smith, A.: Entropic security and the encryption of high entropy messages. In: Theory of Cryptography Conference. pp. 556–577. Springer (2005)
8. Dong, J., Roth, A., Su, W.J.: Gaussian differential privacy. *Journal of the Royal Statistical Society Series B: Statistical Methodology* **84**(1), 3–37 (2022)
9. Dwork, C., Kenthapadi, K., McSherry, F., Mironov, I., Naor, M.: Our data, ourselves: Privacy via distributed noise generation. In: Advances in Cryptology-EUROCRYPT 2006: 24th Annual International Conference on the Theory and

- Applications of Cryptographic Techniques, St. Petersburg, Russia, May 28–June 1, 2006. Proceedings 25. pp. 486–503. Springer (2006)
10. Dwork, C., McSherry, F., Nissim, K., Smith, A.: Calibrating noise to sensitivity in private data analysis. In: Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4–7, 2006. Proceedings 3. pp. 265–284. Springer (2006)
 11. Dwork, C., Rothblum, G.N., Vadhan, S.: Boosting and differential privacy. In: 2010 IEEE 51st annual symposium on foundations of computer science. pp. 51–60. IEEE (2010)
 12. Dziembowski, S., Pietrzak, K.: Leakage-resilient cryptography. In: 2008 49th Annual IEEE Symposium on Foundations of Computer Science. pp. 293–302. IEEE (2008)
 13. Fano, R.M.: Fano inequality. *Scholarpedia* **3**(10), 6648 (2008)
 14. Ghazi, B., Kumar, R., Manurangsi, P.: Differentially private clustering: Tight approximation ratios. *Advances in Neural Information Processing Systems* **33**, 4040–4054 (2020)
 15. Goldwasser, S., Micali, S.: Probabilistic encryption. *Journal of computer and system sciences* **28**(2), 270–299 (1984)
 16. Hardt, M., Rothblum, G.N.: A multiplicative weights mechanism for privacy-preserving data analysis. In: 2010 IEEE 51st annual symposium on foundations of computer science. pp. 61–70. IEEE (2010)
 17. Ishai, Y., Sahai, A., Wagner, D.: Private circuits: Securing hardware against probing attacks. In: Annual International Cryptology Conference. pp. 463–481. Springer (2003)
 18. Issa, I., Wagner, A.B., Kamath, S.: An operational approach to information leakage. *IEEE Transactions on Information Theory* **66**(3), 1625–1657 (2019)
 19. Kairouz, P., Oh, S., Viswanath, P.: The composition theorem for differential privacy. In: International conference on machine learning. pp. 1376–1385. PMLR (2015)
 20. Kasiviswanathan, S.P., Smith, A.: On the ‘semantics’ of differential privacy: A bayesian formulation. *Journal of Privacy and Confidentiality* **6**(1) (2014)
 21. Kifer, D., Machanavajjhala, A.: Pufferfish: A framework for mathematical privacy definitions. *ACM Transactions on Database Systems (TODS)* **39**(1), 1–36 (2014)
 22. Köpf, B., Basin, D.: Automatically deriving information-theoretic bounds for adaptive side-channel attacks. *Journal of Computer Security* **19**(1), 1–31 (2011)
 23. Lin, Z., Gopi, S., Kulkarni, J., Nori, H., Yekhanin, S.: Differentially private synthetic data via foundation model APIs 1: Images. In: The Twelfth International Conference on Learning Representations (2024), <https://openreview.net/forum?id=YEhQs8P0Io>
 24. Liu, J., Jager, T., Kakvi, S.A., Warinschi, B.: How to build time-lock encryption. *Designs, Codes and Cryptography* **86**(11), 2549–2586 (2018)
 25. Lyu, M., Su, D., Li, N.: Understanding the sparse vector technique for differential privacy. *Proceedings of the VLDB Endowment* **10**(6), 637–648 (2017)
 26. Massey, J.L.: Guessing and entropy. In: Proceedings of 1994 IEEE International Symposium on Information Theory. p. 204. IEEE (1994)
 27. Mironov, I.: Rényi differential privacy. In: 2017 IEEE 30th computer security foundations symposium (CSF). pp. 263–275. IEEE (2017)
 28. Murtagh, J., Vadhan, S.: The complexity of computing the optimal composition of differential privacy. In: Theory of Cryptography Conference. pp. 157–175. Springer (2015)

29. Nissim, K., Raskhodnikova, S., Smith, A.: Smooth sensitivity and sampling in private data analysis. In: Proceedings of the thirty-ninth annual ACM symposium on Theory of computing. pp. 75–84 (2007)
30. Papernot, N., Song, S., Mironov, I., Raghunathan, A., Talwar, K., Erlingsson, U.: Scalable private learning with PATE. In: International Conference on Learning Representations (2018), <https://openreview.net/forum?id=rkZB1XbRZ>
31. Pietrzak, K.: A leakage-resilient mode of operation. In: Annual International Conference on the Theory and Applications of Cryptographic Techniques. pp. 462–482. Springer (2009)
32. Ponomareva, N., Hazimeh, H., Kurakin, A., Xu, Z., Denison, C., McMahan, H.B., Vassilvitskii, S., Chien, S., Thakurta, A.G.: How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research* **77**, 1113–1201 (2023)
33. Randolph, M., Diehl, W.: Power side-channel attack analysis: A review of 20 years of study for the layman. *Cryptography* **4**(2), 15 (2020)
34. Russell, A., Wang, H.: How to fool an unbounded adversary with a short key. *IEEE Transactions on Information Theory* **52**(3), 1130–1140 (2006)
35. Sason, I., Verdú, S.: f -divergence inequalities. *IEEE Transactions on Information Theory* **62**(11), 5973–6006 (2016)
36. Shannon, C.E.: Communication theory of secrecy systems. *The Bell system technical journal* **28**(4), 656–715 (1949)
37. Sibson, R.: Information radius. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* **14**(2), 149–160 (1969)
38. Smith, G.: On the foundations of quantitative information flow. In: International Conference on Foundations of Software Science and Computational Structures. pp. 288–302. Springer (2009)
39. Sridhar, M., Xiao, H., Devadas, S.: Pac-private algorithms. In: 2025 IEEE Symposium on Security and Privacy (SP). pp. 3839–3857. IEEE (2025)
40. Standaert, F.X., Malkin, T.G., Yung, M.: A unified framework for the analysis of side-channel key recovery attacks. In: Annual international conference on the theory and applications of cryptographic techniques. pp. 443–461. Springer (2009)
41. Van Den Hooff, J., Lazar, D., Zaharia, M., Zeldovich, N.: Vuvuzela: Scalable private messaging resistant to traffic analysis. In: Proceedings of the 25th Symposium on Operating Systems Principles. pp. 137–152 (2015)
42. Verdú, S.: α -mutual information. In: 2015 Information Theory and Applications Workshop (ITA). pp. 1–6. IEEE (2015)
43. Wainwright, M.J.: High-dimensional statistics: A non-asymptotic viewpoint, vol. 48. Cambridge university press (2019)
44. Xiao, H., Devadas, S.: Pac privacy: Automatic privacy measurement and control of data processing. In: Annual International Cryptology Conference. pp. 611–644. Springer (2023)
45. Xiao, H., Suh, G.E., Devadas, S.: Formal privacy proof of data encoding: The possibility and impossibility of learnable encryption. In: Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security. pp. 1834–1848 (2024)
46. Xiao, X., Tao, Y.: Output perturbation with query relaxation. *Proceedings of the VLDB Endowment* **1**(1), 857–869 (2008)
47. Young, G.A., Smith, R.L., Smith, R.L.: Essentials of statistical inference, vol. 16. Cambridge University Press (2005)
48. Zhang, J., Chen, C., Cui, J., Li, K.: Timing side-channel attacks and countermeasures in cpu microarchitectures. *ACM Computing Surveys* **56**(7), 1–40 (2024)

A Proof of Proposition 1

Taking the $\frac{1}{\alpha}$ power of both sides from Lemma 2 and let $\alpha \rightarrow \infty$, then the left-hand-side converges to $\frac{\delta_\rho}{\delta_{o,\rho}}$ and the right-hand-side converges to $\exp\{I_\infty^{sib}(\mathbf{x}; \mathcal{F}(\mathbf{x}))\}$. Here $I_\infty^{sib}(\mathbf{x}; \mathcal{F}(\mathbf{x}))$ is Sibson's mutual information of order infinity (Definition 5). Notice that for any $\hat{\mathbf{x}} \in \mathcal{X}$, $\Pr_{\mathbf{x} \sim \mathcal{D}}(\mathbf{x} = \hat{\mathbf{x}}) = \frac{1}{|\mathcal{X}|}$

Therefore,

$$\delta_\rho = \delta_{o,\rho} \cdot \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathcal{F}(\hat{\mathbf{x}}) = \mathbf{o}) \, d\mathbf{o} \quad (34)$$

$$= \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathbf{x} = \hat{\mathbf{x}}) \Pr(\mathcal{F}(\hat{\mathbf{x}}) = \mathbf{o}) \, d\mathbf{o} \quad (35)$$

$$= \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} \Pr(\mathbf{x} = \hat{\mathbf{x}}, \mathcal{F}(\hat{\mathbf{x}}) = \mathbf{o}) \, d\mathbf{o}. \quad (36)$$

B Proof of Lemma 3

Let $g : \mathcal{O} \rightarrow \mathcal{X}$ be any deterministic decision rule, where $g(\mathbf{o})$ is the guess when the observation is \mathbf{o} . The success event is $\{\rho(\mathbf{x}, g(\mathbf{o})) = 1\}$, so

$$\begin{aligned} \Pr(\rho(\mathbf{x}, g(\mathbf{o})) = 1) &= \int_{\mathcal{O}} \sum_{\bar{\mathbf{x}} \in \mathcal{X}} \mathbf{1}_{\{\rho(\bar{\mathbf{x}}, g(\mathbf{o})) = 1\}} \Pr(\mathbf{x} = \bar{\mathbf{x}}, \mathbf{o}) \, d\mathbf{o} \\ &= \int_{\mathcal{O}} \sum_{\bar{\mathbf{x}} \in \mathcal{X}} \mathbf{1}_{\{\rho(\bar{\mathbf{x}}, g(\mathbf{o})) = 1\}} \Pr(\mathbf{x} = \bar{\mathbf{x}}) \Pr(\mathbf{o} \mid \mathbf{x} = \bar{\mathbf{x}}) \, d\mathbf{o}. \end{aligned}$$

For each fixed \mathbf{o} , the inner sum is exactly $H(g(\mathbf{o}), \mathbf{o})$ by the definition of H , hence

$$\Pr(\rho(\mathbf{x}, g(\mathbf{o})) = 1) = \int_{\mathcal{O}} H(g(\mathbf{o}), \mathbf{o}) \, d\mathbf{o} \leq \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} H(\hat{\mathbf{x}}, \mathbf{o}) \, d\mathbf{o},$$

since $H(g(\mathbf{o}), \mathbf{o}) \leq \max_{\hat{\mathbf{x}}} H(\hat{\mathbf{x}}, \mathbf{o})$ for every \mathbf{o} .

Conversely, for each $\mathbf{o} \in \mathcal{O}$ we can choose $g^*(\mathbf{o}) \in \arg \max_{\hat{\mathbf{x}} \in \mathcal{X}} H(\hat{\mathbf{x}}, \mathbf{o})$ (choosing any maximizer when there are ties). For this g^* we have $H(g^*(\mathbf{o}), \mathbf{o}) = \max_{\hat{\mathbf{x}}} H(\hat{\mathbf{x}}, \mathbf{o})$ for all \mathbf{o} , and thus

$$\Pr(\rho(\mathbf{x}, g^*(\mathbf{o})) = 1) = \int_{\mathcal{O}} \max_{\hat{\mathbf{x}} \in \mathcal{X}} H(\hat{\mathbf{x}}, \mathbf{o}) \, d\mathbf{o}.$$

Combining the upper bound with the equality achieved by g^* proves the claim.

C Proof of Lemma 4

Notice that

$$\left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \geq \left(\max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} = \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}), \quad (37)$$

take integral over \mathcal{O} we directly get

$$\int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} d\mathbf{o} > \int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}) d\mathbf{o}$$

D Proof of Lemma 5

Fix $\alpha > 1$ and define, for each $\mathbf{o} \in \mathcal{O}$,

$$G(\mathbf{o}) := \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}.$$

Then the left-hand side of the inequality is simply $\int_{\mathcal{O}} G(\mathbf{o}) d\mathbf{o}$. We insert the reference density W and apply Hölder's inequality with exponents α and $\frac{\alpha}{\alpha-1}$:

$$\begin{aligned} \int_{\mathcal{O}} G(\mathbf{o}) d\mathbf{o} &= \int_{\mathcal{O}} G(\mathbf{o}) \Pr(W = \mathbf{o})^{\frac{\alpha-1}{\alpha}} \Pr(W = \mathbf{o})^{-\frac{\alpha-1}{\alpha}} d\mathbf{o} \\ &\leq \left(\int_{\mathcal{O}} \left(G(\mathbf{o}) \Pr(W = \mathbf{o})^{\frac{\alpha-1}{\alpha}} \right)^\alpha d\mathbf{o} \right)^{1/\alpha} \left(\int_{\mathcal{O}} \left(\Pr(W = \mathbf{o})^{-\frac{\alpha-1}{\alpha}} \right)^{\frac{\alpha}{\alpha-1}} d\mathbf{o} \right)^{\frac{\alpha-1}{\alpha}} \\ &= \left(\int_{\mathcal{O}} \frac{G(\mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right)^{1/\alpha} \left(\int_{\mathcal{O}} \Pr(W = \mathbf{o}) d\mathbf{o} \right)^{\frac{\alpha-1}{\alpha}}. \end{aligned}$$

Since W is a probability distribution, $\int_{\mathcal{O}} \Pr(W = \mathbf{o}) d\mathbf{o} = 1$, so we obtain

$$\int_{\mathcal{O}} G(\mathbf{o}) d\mathbf{o} \leq \left(\int_{\mathcal{O}} \frac{G(\mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right)^{1/\alpha} = \left(\int_{\mathcal{O}} \frac{\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} \right)^{1/\alpha},$$

which is the first inequality.

For the expectation form, note that

$$\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha = \sum_{\mathbf{x} \in \mathcal{X}} \Pr(\mathbf{x}) \frac{H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(\mathbf{x})},$$

so

$$\int_{\mathcal{O}} \frac{\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o} = \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}} \frac{H(\mathbf{x}, \mathbf{o})^\alpha / \Pr(\mathbf{x})}{\Pr(W = \mathbf{o})^{\alpha-1}} d\mathbf{o},$$

which yields the second equality.

Finally, Hölder's inequality is tight if and only if there exists a constant $\lambda > 0$ such that

$$G(\mathbf{o})^\alpha = \lambda \Pr(W = \mathbf{o})^\alpha \quad \text{for almost every } \mathbf{o},$$

which is equivalent to

$$\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha = \lambda \Pr(W = \mathbf{o})^\alpha \quad \text{a.e. } \mathbf{o}.$$

Rewriting this gives the stated condition that the bound is tight when W is chosen proportional to the α -norm of $H(\cdot, \mathbf{o})$:

$$\Pr(W = \mathbf{o}) \propto \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}.$$

E Proof of Theorem 1

For any admissible noise distribution \mathcal{D}_e , write

$$\mathcal{G}(\mathcal{D}_e, \mathbf{x}) = \mathcal{D}_\alpha(P_{\mathcal{M}_{\mathcal{D}_e}(\mathbf{x})} \parallel P_{W_{\mathcal{D}_e}}), \quad \mu(\mathcal{D}_e) := \mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\mathcal{D}_e, \mathbf{x})].$$

Algorithm 1 draws two i.i.d. samples $\mathbf{X} = \{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_m\}$ and $\mathbf{X}' = \{\bar{\mathbf{x}}'_1, \dots, \bar{\mathbf{x}}'_m\}$ from D , and then:

1. uses only \mathbf{X} to solve the empirical constraint (13) and obtain an initial candidate \mathcal{D}_e ; 2. uses only the independent sample \mathbf{X}' to test the proposed distribution \mathcal{D}_e . Let

$$L(\mathcal{D}_e) = \frac{1}{m} \sum_{i=1}^m \mathcal{G}(\mathcal{D}_e, \bar{\mathbf{x}}'_i), \quad \beta(\mathcal{D}_e) = B(\mathcal{D}_e) \sqrt{\frac{\log(1/\gamma)}{2m}},$$

and runs the while-loop until both $L(\mathcal{D}_e) \leq r' - \beta(\mathcal{D}_e)$ and $B(\mathcal{D}_e) \leq R$ hold. Let $\widehat{\mathcal{D}}_e$ denote the final noise distribution returned by the algorithm.

We shall show that $\mu(\widehat{\mathcal{D}}_e) \leq r$, which is exactly

$$\mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\widehat{\mathcal{D}}_e, \mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim D} [\mathcal{D}_\alpha(P_{\mathcal{M}_{\widehat{\mathcal{D}}_e}(\mathbf{x})} \parallel P_{W_{\widehat{\mathcal{D}}_e}})] \leq r.$$

Step 1: Condition on the training sample and define the ideal parameter. Fix the first sample set \mathbf{X} and condition on it. Under this conditioning, the candidate produced by the empirical optimization (13) depends only on \mathbf{X} . We further assume that there exists $\sigma^* \geq 0$ and a corresponding noise distribution \mathcal{D}_e^* , given by the law of $\mathbf{e} + \sigma^* Z$ with $\mathbf{e} \sim \mathcal{D}_e$ and $Z \sim \mathcal{N}(0, \mathbf{I})$ independent, that *saturates* the constraint, i.e.

$$\mu(\mathcal{D}_e^*) = r'.$$

If no such σ^* exists (that is, $\mu(\mathcal{D}_e) < r'$ for every admissible \mathcal{D}_e), then in particular the final output $\widehat{\mathcal{D}}_e$ of Algorithm 1 satisfies $\mu(\widehat{\mathcal{D}}_e) \leq r' \leq r$, and the desired guarantee already holds. Hence, without loss of generality, we may assume that such a σ^* (and \mathcal{D}_e^*) exists. This \mathcal{D}_e^* is a function of the fixed training sample \mathbf{X} , but crucially it is *independent* of the second, validation sample \mathbf{X}' .

Step 2: Hoeffding bound on the independent validation sample. Now condition on \mathbf{X} and on the choice of \mathcal{D}_e^* . Under this conditioning, the random variables

$$Y_i := \mathcal{G}(\mathcal{D}_e^*, \bar{\mathbf{x}}'_i), \quad i = 1, \dots, m,$$

are i.i.d. and bounded: $0 \leq Y_i \leq B(\mathcal{D}_e^*)$, and

$$\mathbb{E}[Y_i] = \mu(\mathcal{D}_e^*) = r'.$$

Their empirical average is

$$L^* := L(\mathcal{D}_e^*) = \frac{1}{m} \sum_{i=1}^m Y_i.$$

By Hoeffding's inequality, for any $\beta > 0$,

$$\Pr_{\mathbf{X}'} \left(L^* \leq r' - \beta \mid \mathbf{X} \right) \leq \exp \left(- \frac{2m\beta^2}{B(\mathcal{D}_e^*)^2} \right).$$

Set

$$\beta^* := B(\mathcal{D}_e^*) \sqrt{\frac{\log(1/\gamma)}{2m}},$$

then the right-hand side is exactly γ , and we obtain

$$\Pr_{\mathbf{X}'} \left(L^* \leq r' - \beta^* \mid \mathbf{X} \right) \leq \gamma. \quad (38)$$

Equivalently,

$$\Pr_{\mathbf{X}'} \left(L^* > r' - \beta^* \mid \mathbf{X} \right) \geq 1 - \gamma.$$

Define the “good” event

$$\mathbf{E} := \{L^* > r' - \beta^*\}.$$

On \mathbf{E} we have

$$L^* > r' - \beta^*,$$

so L^* does *not* satisfy the validation condition $L(\mathcal{D}_e) \leq r' - \beta(\mathcal{D}_e)$, and hence the ideal benchmark \mathcal{D}_e^* fails the validation test on \mathbf{E} .

Step 3: Monotonicity of the update step. Whenever the validation test fails, Algorithm 1 updates the current noise distribution by adding an independent Gaussian:

$$\mathbf{e}_{\text{new}} = \mathbf{e}_{\text{old}} + z, \quad z \sim \mathcal{N}(0, \sigma^2 \mathbf{I}),$$

so that the new privatized mechanism is a post-processing of the old one. By the data-processing inequality for f -divergences (in particular, for α -divergence), such a common post-processing cannot increase the privacy loss:

$$\mathcal{G}(\mathcal{D}_e^{\text{new}}, \mathbf{x}) \leq \mathcal{G}(\mathcal{D}_e^{\text{old}}, \mathbf{x}), \quad \forall \mathbf{x}.$$

Consequently, the population risk is nonincreasing along the updates:

$$\mu(\mathcal{D}_e^{\text{new}}) := \mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\mathcal{D}_e^{\text{new}}, \mathbf{x})] \leq \mu(\mathcal{D}_e^{\text{old}}).$$

Moreover, our global bound $B(\mathcal{D}_e)$ also decreases under this update: we have

$$B(\mathcal{D}_e^{\text{new}}) \leq B(\mathcal{D}_e^{\text{old}}),$$

and hence the Hoeffding deviation term

$$\beta(\mathcal{D}_e) := B(\mathcal{D}_e) \sqrt{\frac{\log(1/\gamma)}{2m}}$$

is also nonincreasing along the updates:

$$\beta(\mathcal{D}_e^{\text{new}}) \leq \beta(\mathcal{D}_e^{\text{old}}).$$

Equivalently, the validation threshold $r' - \beta(\mathcal{D}_e)$ is *nondecreasing* as the algorithm adds more noise:

$$r' - \beta(\mathcal{D}_e^{\text{new}}) \geq r' - \beta(\mathcal{D}_e^{\text{old}}).$$

Thus, for a fixed training set X , as the algorithm iteratively updates \mathcal{D}_e using only X' , we obtain a sequence of mechanisms whose population risks $\mu(\mathcal{D}_e)$ form a nonincreasing sequence, while the validation threshold $r' - \beta(\mathcal{D}_e)$ is nondecreasing. On the event E from Step 2, the benchmark \mathcal{D}_e^* satisfies

$$L(\mathcal{D}_e^*) > r' - \beta(\mathcal{D}_e^*),$$

so it fails the validation test and the algorithm cannot stop at \mathcal{D}_e^* . Every subsequent update strictly increases the noise level and hence yields a mechanism $\mathcal{D}_e^{\text{new}}$ with

$$\mu(\mathcal{D}_e^{\text{new}}) \leq \mu(\mathcal{D}_e^*) = r',$$

while the corresponding threshold $r' - \beta(\mathcal{D}_e^{\text{new}})$ is no smaller than $r' - \beta(\mathcal{D}_e^*)$. When the algorithm eventually stops at some $\widehat{\mathcal{D}}_e$ that passes the test $L(\widehat{\mathcal{D}}_e) \leq r' - \beta(\widehat{\mathcal{D}}_e)$, we necessarily have

$$\mu(\widehat{\mathcal{D}}_e) \leq \mu(\mathcal{D}_e^*) = r' \quad \text{on the event } E, \quad (39)$$

because $\widehat{\mathcal{D}}_e$ arises from \mathcal{D}_e^* by a sequence of such a post-processing updates.

Step 4: Handling the complement event and using the fallback bound. On the complement event E^c , we do not control the deviation of $L(\mathcal{D}_e^*)$ from r' . However, by the stopping rule of Algorithm 1, the final output $\widehat{\mathcal{D}}_e$ always satisfies

$$B(\widehat{\mathcal{D}}_e) \leq R,$$

and hence, using $\mathcal{G}(\mathcal{D}_e, x) \leq B(\mathcal{D}_e)$ for all x , we get

$$\mu(\widehat{\mathcal{D}}_e) = \mathbb{E}_{x \sim D} [\mathcal{G}(\widehat{\mathcal{D}}_e, x)] \leq B(\widehat{\mathcal{D}}_e) \leq R \quad \text{on } E^c. \quad (40)$$

Step 5: Final bound. Now take expectation over both \mathbf{X} and \mathbf{X}' . Combining (39) and (40) and using $\Pr(\mathbf{E}) \geq 1 - \gamma$ and $\Pr(\mathbf{E}^c) \leq \gamma$, we obtain

$$\mathbb{E}[\mu(\widehat{\mathcal{D}}_e)] \leq (1 - \gamma)r' + \gamma R \leq r,$$

where the last inequality is exactly the requirement imposed in step 3 of Algorithm 1.

Finally, by the definition of $\mu(\widehat{\mathcal{D}}_e)$, this is equivalent to

$$\mathbb{E}_{\mathbf{x} \sim D}[\mathcal{G}(\widehat{\mathcal{D}}_e, \mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim D}[\mathcal{D}_\alpha(P_{\mathcal{M}_{\widehat{\mathcal{D}}_e}(\mathbf{x})} \| P_{W_{\widehat{\mathcal{D}}_e}})] \leq r.$$

This proves the theorem.

F Proof of Theorem 2

Apply Lemma 3, Lemma 4 and Lemma 5, we know

$$\delta_\rho = \int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o}) \, d\mathbf{o} \tag{41}$$

$$\leq \int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} \, d\mathbf{o} \tag{42}$$

$$\leq \left(\int_{\mathcal{O}} \frac{\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha}{\Pr(W = \mathbf{o})^{\alpha-1}} \, d\mathbf{o} \right)^{1/\alpha} \tag{43}$$

$$= l_{\alpha, \rho}^W(\mathbf{x}; \mathcal{F}(\mathbf{x})) \tag{44}$$

since this hold for arbitrary W , the next inequality also holds:

$$\delta_\rho \leq l_{\alpha, \rho}(\mathbf{x}; \mathcal{F}(\mathbf{x})).$$

G Proof of Theorem 3

Applying Lemma 5, we see that choosing

$$\Pr(W = \mathbf{o}) \propto \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha}$$

is exactly the condition under which Hölder's inequality becomes tight.

H Proof of Theorem 4

$$\lim_{\alpha \rightarrow \infty} \mathsf{I}_{\alpha, \rho}(\mathbf{x}; \mathcal{F}(\mathbf{x})) \quad (45)$$

$$= \int_{\mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} H(\mathbf{x}, \mathbf{o})^\alpha \right)^{1/\alpha} d\mathbf{o}, \quad (46)$$

$$= \int_{\mathcal{O}} \max_{\mathbf{x} \in \mathcal{X}} H(\mathcal{F}(\mathbf{x}) = \mathbf{o}) d\mathbf{o} \quad (47)$$

$$= \delta_\rho \quad (48)$$

I Proof of Theorem 5

First we need to make it clear that δ_ρ is a parameter in this theorem we use another notation δ_ρ^* to denote the optimal posterior success rate. By applying Theorem 1, we have

$$\mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\mathcal{D}_e, \mathbf{x})] \leq r.$$

Use $\mathcal{G}(\mathcal{D}_e, \hat{\mathbf{x}}) := \frac{1}{\Pr(\mathbf{x} = \hat{\mathbf{x}})} \mathcal{D}_\alpha(P_{H_{\mathcal{D}_e}(\hat{\mathbf{x}})} \| P_{W_{\mathcal{D}_e}})$ into the above equation we have

$$\mathbb{E}_{\mathbf{x} \sim D} \frac{1}{\Pr(\mathbf{x} = \hat{\mathbf{x}})} \mathcal{D}_\alpha(P_{H_{\mathcal{D}_e}(\hat{\mathbf{x}})} \| P_{W_{\mathcal{D}_e}}) \leq r$$

Notice that

$$\mathbb{E}_{\mathbf{x} \sim D} \frac{1}{\Pr(\mathbf{x} = \hat{\mathbf{x}})} \mathcal{D}_\alpha(P_{H_{\mathcal{D}_e}(\hat{\mathbf{x}})} \| P_{W_{\mathcal{D}_e}}) = \mathsf{I}_{\alpha, \rho}^W(\mathbf{x}; \mathcal{M})^\alpha$$

and

$$r = (\delta_\rho)^\alpha.$$

Apply Theorem 2 we know

$$\delta_\rho^* \leq \mathsf{I}_{\alpha, \rho}^W(\mathbf{x}; \mathcal{M}) \leq \delta_\rho$$

So the optimal poster success rate δ_ρ^* is no more than δ_ρ .

J Proof of Theorem 6

We use $\mathcal{M}^{Adv} : \mathcal{X} \rightarrow \mathcal{O}^t \times V^t$ to denote the composed mechanism including the query output from *Adv*. For example, when $T = 2$

$$\begin{aligned} & \Pr(\mathcal{M}^{Adv} = (\mathbf{v}_1, \mathbf{o}_1, \mathbf{v}_2, \mathbf{o}_2)) \\ &= \Pr(Adv_1(\tau_1) = \mathbf{v}_1) \Pr(\mathcal{M}_1(\mathbf{x}, \mathbf{v}_1, \tau_1) = \mathbf{o}_1) \Pr(Adv_2(\tau_2) = \mathbf{v}_2) \Pr(\mathcal{M}_2(\mathbf{x}, \mathbf{v}_2, \tau_2) = \mathbf{o}_2) \end{aligned}$$

Apply Theorem 2, we have

$$\delta_\rho \leq \mathsf{I}_{\alpha, \rho}^W(\mathbf{x}; \mathcal{M}^{Adv})$$

For the reference distribution we define for a family of distribution:

$$W_t : \mathcal{V} \times \mathcal{T}_t \rightarrow \mathcal{O}$$

and

$$Adv'_t : \mathcal{T}_t \rightarrow \mathcal{O}$$

Apply Theorem 2 then we have the optimal posterior success rate δ_ρ satisfies:

$$\begin{aligned} \delta_\rho^\alpha &\leq l_{\alpha, \rho}^W(\mathbf{x}; \mathcal{M}^{Adv}) \\ &= \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \Pr(\mathbf{x})^{\alpha-1} \cdot \prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} \\ &\quad \cdot \prod_{t=1}^T \frac{(\Pr(Adv_t(\tau_t) = \mathbf{v}_t))^\alpha}{(\Pr(Adv'_t(\tau_t) = \mathbf{v}_t))^{\alpha-1}} d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \end{aligned} \quad (49)$$

We choose Adv' has the exact same distribution with Adv , then

$$\prod_{t=1}^T \frac{(\Pr(Adv_t(\tau_t) = \mathbf{v}_t))^\alpha}{(\Pr(Adv'_t(\tau_t) = \mathbf{v}_t))^{\alpha-1}} = \prod_{t=1}^T \Pr(Adv_t(\tau_t) = \mathbf{v}_t)$$

As a result

$$\delta_\rho^\alpha \leq (50) = \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \Pr(\mathbf{x})^{\alpha-1} \cdot \prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} \quad (51)$$

$$\cdot \prod_{t=1}^T \Pr(Adv_t(\tau_t) = \mathbf{v}_t) d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \quad (52)$$

K Proof of Theorem 7

Apply Theorem 5, we have

$$\begin{aligned} (\delta_\rho)^\alpha &\leq \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \Pr(\mathbf{x})^{\alpha-1} \cdot \prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} \\ &\quad \cdot \prod_{t=1}^T \Pr(Adv_t(\tau_t) = \mathbf{v}_t) d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \end{aligned} \quad (53)$$

For simplicity but without loss of generality, suppose $T = 2$, notice that inequality (53) can be rewritten as

$$\begin{aligned}
(\delta_\rho)^\alpha &\leq \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = \mathbf{v}_2) \\
&\quad \cdot \Pr(\mathbf{x})^{\alpha-1} \frac{(\Pr(\mathcal{M}_1(\mathbf{x}, \mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}{(\Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha} \int_{\mathcal{O}} \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} d\mathbf{o}_2 d\mathbf{v}_2 d\mathbf{o}_1 d\mathbf{v}_1 \\
&\quad (54) \\
&= \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = \mathbf{v}_2) \\
&\quad \cdot \underbrace{\mathbb{E}_{\mathbf{x} \sim D} \Pr(\mathbf{x})^{\alpha-1} \frac{(\Pr(\mathcal{M}_1(\mathbf{x}, \mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}{(\Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}}_{(A)} \underbrace{\int_{\mathcal{O}} \frac{(\Pr(\mathcal{M}_2(\mathbf{x}, \mathbf{v}_2, \tau_2) = \mathbf{o}_2))^\alpha}{(\Pr(W_2(\mathbf{v}_2, \tau_2) = \mathbf{o}_2))^{\alpha-1}} d\mathbf{o}_2 d\mathbf{v}_2 d\mathbf{o}_1 d\mathbf{v}_1}_{(B)} \\
&\quad (55)
\end{aligned}$$

Notice that term (A) is $A_1(\mathbf{x})$ and term (B) is $D_2(\mathbf{x})$, the condition in the theorem, tells us that

$$\mathbb{E}_{\mathbf{x} \sim D} [A_1(\mathbf{x}) D_2(\mathbf{x})] \leq r_2 \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_1(\mathbf{x})], \quad (56)$$

so (55) \leq

$$\begin{aligned}
&\int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = \mathbf{v}_2) \\
&\quad \cdot r_2 \cdot \underbrace{\mathbb{E}_{\mathbf{x} \sim D} \Pr(\mathbf{x})^{\alpha-1} \frac{(\Pr(\mathcal{M}_1(\mathbf{x}, \mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}{(\Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}}_{(A)} d\mathbf{v}_2 d\mathbf{o}_1 d\mathbf{v}_1 \\
&\quad (57)
\end{aligned}$$

$$\begin{aligned}
&= r_2 \cdot \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \\
&\quad \underbrace{\int_{\mathcal{O}} \Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1) \cdot \mathbb{E}_{\mathbf{x} \sim D} \Pr(\mathbf{x})^{\alpha-1} \frac{(\Pr(\mathcal{M}_1(\mathbf{x}, \mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha}{(\Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1))^\alpha} d\mathbf{o}_1}_{(D)} \\
&\quad \underbrace{\int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = \mathbf{v}_2) d\mathbf{v}_2 d\mathbf{v}_1}_{(C)} \\
&\quad (58)
\end{aligned}$$

$$\leq r_2 \cdot \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \cdot r_1 \cdot 1 = r_1 r_2 \quad (59)$$

Equality (58) follows by straightforward algebraic rearrangement. Inequality (59) is because term (D) is smaller than r_1 and term C equals to 1.

L Proof of Theorem 8

We focus on this ratio

$$\left(\prod_{t=1}^T \frac{\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t)}{\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t)} \right)^\alpha \quad (60)$$

From the equal condition of Holder's inequality as is stated in Theorem 3, it's sufficient and necessary that ratio (60) is a constant for any selection of $\mathbf{o}_{[1:T]}, \mathbf{v}_{[1:T]}$. Given our selection of W , we have

$$\mathbb{E}_{\mathbf{x} \sim D} \left[\left(\prod_{t=1}^T \frac{\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t)}{\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t)} \right)^\alpha \right] \quad (61)$$

$$= r_T \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_{T-1}(\mathbf{x})] \quad (62)$$

$$(63)$$

Use this proposition iteratively

$$\mathbb{E}_{\mathbf{x} \sim D} \left[\left(\prod_{t=1}^T \frac{\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t)}{\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t)} \right)^\alpha \right] \quad (64)$$

$$= \mathbb{E}_{\mathbf{x} \sim D} \left[A_{T-1}(\mathbf{x}) \left(\frac{\Pr(\mathcal{M}_T(\mathbf{x}, \mathbf{v}_T, \tau_T) = \mathbf{o}_T)}{\Pr(W_T(\mathbf{v}_T, \tau_T) = \mathbf{o}_T)} \right)^\alpha \right] \quad (65)$$

$$= r_T \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_{T-1}(\mathbf{x})] \quad (66)$$

$$= r_T r_{T-1} \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_{T-2}(\mathbf{x})] \quad (67)$$

$$\vdots \quad (68)$$

$$= \left(\prod_{t=1}^T r_t \right) \cdot \mathbb{E}_{\mathbf{x} \sim D} [A_0(\mathbf{x})] \quad (69)$$

$$= \prod_{t=1}^T r_t, \quad (70)$$

So ratio (60) is a constant for any selection of $\mathbf{o}_{[1:T]}, \mathbf{v}_{[1:T]}$. We emphasize again this condition is necessary and sufficient for the optimal choice of W . Hence, selecting the locally optimal W_t at each round yields a globally optimal reference W .

M Proof of Theorem 9

By applying Theorem 1, we have

$$\mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\mathcal{D}_e, \mathbf{x})] \leq 0.$$

Use $\mathcal{G}_t(\mathcal{D}_e, \hat{x}) := (A_{t-1}(\hat{x})(D_t(\hat{x}) - r_t)$ into the above equation for every t and every selection of we have

$$\mathbb{E}_{x \sim D} [A_{t-1}(x)(D_t(x) - r_t)] \leq 0 \quad (71)$$

which is equals to

$$\mathbb{E}_{x \sim D} [A_{t-1}(x) D_t(x)] \leq r_t \cdot \mathbb{E}_{x \sim D} [A_{t-1}(x)], \quad (72)$$

Consequently the condition 21 is held so we apply Theorem 7 we have:

$$\delta_\rho^\alpha \leq \prod_{t=1}^T r_t. \quad (73)$$

N Proof of Theorem 10

By Lemma 2, to bound $D_{\text{KL}}(\mathbf{1}_\rho \| \mathbf{1}_{o,\rho})$, it suffices to consider the entire interactive mechanism \mathcal{M} , which merges both the adversary's strategy Adv and the user's implementations across T time slots, and

$$D_f(\mathbf{1}_\rho \| \mathbf{1}_{o,\rho}) \leq \mathbb{E}_{x \sim D} \inf_{P_W} D_f(P_{\mathcal{M}(x)} \| P_W). \quad (74)$$

For simplicity but without loss of generality, suppose $T = 2$ and select $f(z) = z \log z$. Since the choice of the reference distribution P_W can be arbitrary, set $W = (W_1, W_2)$ such that the marginal distribution P_{W_1} equals $P_{W^{(1)}}$ and W_2 conditional on $W_1 = \mathbf{o}_1$ is the mixture of $\{P_{W_2(v)}\}_{v \in \mathcal{V}}$, i.e.,

$$\Pr(W_2 = \mathbf{o}_2 \mid W_1 = \mathbf{o}_1) = \sum_{v \in \mathcal{V}} \Pr(\text{Adv}(\mathbf{o}_1) = v) \Pr(W_2(v) = \mathbf{o}_2). \quad (75)$$

Then under this choice we have

$$\mathbb{E}_{x \sim D} D_f(P_{\mathcal{M}(x)} \| P_W) \quad (76)$$

$$\begin{aligned} &= \int_{\mathbf{o}_1 \in \mathcal{O}^{(1)}} \sum_{v \in \mathcal{V}} \Pr(\text{Adv}(\mathbf{o}_1) = v) \int_{\mathbf{o}_2 \in \mathcal{O}^{(2)}} \Pr(\mathcal{M}_1(x) = \mathbf{o}_1) \Pr(\mathcal{M}_2(x, v) = \mathbf{o}_2) \\ &\quad \cdot \log \frac{\Pr(\mathcal{M}_1(x) = \mathbf{o}_1) \Pr(\mathcal{M}_2(x, v) = \mathbf{o}_2)}{\Pr(W_1 = \mathbf{o}_1) \Pr(W_2(v) = \mathbf{o}_2)} \\ &= \underbrace{\int_{\mathbf{o}_1 \in \mathcal{O}^{(1)}} \Pr(\mathcal{M}_1(x) = \mathbf{o}_1) \log \frac{\Pr(\mathcal{M}_1(x) = \mathbf{o}_1)}{\Pr(W_1 = \mathbf{o}_1)}}_A \cdot \underbrace{\sum_{v \in \mathcal{V}} \Pr(\text{Adv}(\mathbf{o}_1) = v) \int_{\mathbf{o}_2 \in \mathcal{O}^{(2)}} \Pr(\mathcal{M}_2(x, v) = \mathbf{o}_2)}_B \\ &\quad + \underbrace{\sum_{v \in \mathcal{V}} \int_{\mathbf{o}_1 \in \mathcal{O}^{(1)}} Q(x, \mathbf{o}_1, v) \int_{\mathbf{o}_2 \in \mathcal{O}^{(2)}} \Pr(\mathcal{M}_2(x, v) = \mathbf{o}_2) \log \frac{\Pr(\mathcal{M}_2(x, v) = \mathbf{o}_2)}{\Pr(W_2(v) = \mathbf{o}_2)}}_C. \end{aligned} \quad (77)$$

where

$$Q(\mathbf{x}, \mathbf{o}_1, \mathbf{v}) = \Pr(\mathcal{M}_1(\mathbf{x}) = \mathbf{o}_1) \Pr(\text{Adv}(\mathbf{o}_1) = \mathbf{v}). \quad (78)$$

Term B equals 1. Define

$$D^{(1)}(\mathbf{x}) = \int_{\mathbf{o}_1 \in \mathcal{O}^{(1)}} \Pr(\mathcal{M}_1(\mathbf{x}) = \mathbf{o}_1) \log \frac{\Pr(\mathcal{M}_1(\mathbf{x}) = \mathbf{o}_1)}{\Pr(W_1 = \mathbf{o}_1)}, \quad (79)$$

$$D_{\mathbf{v}}^{(2)}(\mathbf{x}) = \int_{\mathbf{o}_2 \in \mathcal{O}^{(2)}} \Pr(\mathcal{M}_2(\mathbf{x}, \mathbf{v}) = \mathbf{o}_2) \log \frac{\Pr(\mathcal{M}_2(\mathbf{x}, \mathbf{v}) = \mathbf{o}_2)}{\Pr(W_2(\mathbf{v}) = \mathbf{o}_2)}. \quad (80)$$

From (77),

$$\mathbb{E}_{\mathbf{x} \sim D} D_{\text{KL}}(P_{\mathcal{M}(\mathbf{x})} \| P_W) \leq \mathbb{E}_{\mathbf{x} \sim D} \left[D^{(1)}(\mathbf{x}) + \sum_{\mathbf{v} \in \mathcal{V}} \Pr(\mathbf{v} \mid \mathbf{x}) D_{\mathbf{v}}^{(2)}(\mathbf{x}) \right]. \quad (81)$$

We now bound $\mathbb{E}_{\mathbf{x} \sim D} |\Pr(\mathbf{v} \mid \mathbf{x}) - \Pr(\mathbf{v})|$. Choose

$$P_{W_1} = \sum_{\mathbf{x}'} \Pr(\mathbf{x} = \mathbf{x}') P_{\mathcal{M}_1(\mathbf{x}')}. \quad (82)$$

Under the assumption $\mathbb{E}_{\mathbf{x}} \text{TV}(P_{\mathcal{M}_1(\mathbf{x})} \| P_{W_1}) \leq r_1$, Pinsker's inequality gives

$$\mathbb{E}_{\mathbf{x}} \text{TV}(P_{\mathcal{M}_1(\mathbf{x})} \| P_{W_1}) \leq \sqrt{2r_1}. \quad (83)$$

Furthermore,

$$\sum_{\mathbf{v} \in \mathcal{V}} \mathbb{E}_{\mathbf{x}} |\Pr(\mathbf{v} \mid \mathbf{x}) - \Pr(\mathbf{v})| = \mathbb{E}_{\mathbf{x}} \text{TV}(P_{\mathcal{M}_1(\mathbf{x})} \| P_{W_1}) \leq \sqrt{2r_1}. \quad (84)$$

Since each $D_{\mathbf{v}}^{(2)}(\mathbf{x}) \leq R_2$, the second term in (81) is bounded by

$$r_2 + R_2 \sqrt{2r_1}. \quad (85)$$

Thus,

$$\mathbb{E}_{\mathbf{x} \sim D} D_{\text{KL}}(P_{\mathcal{M}(\mathbf{x})} \| P_W) \leq r_1 + r_2 + R_2 \sqrt{2r_1}. \quad (86)$$

By induction across all T rounds, we recover the general upper bound $B(t)$ shown in (4.8).

O Proof of Theorem 11

Apply Theorem 6, to bound the optimal posterior success rate δ_ρ , we only need to bound

$$\mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \frac{\Pr(\mathbf{x})^\alpha}{D} \cdot \prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{\alpha-1}} \cdot \prod_{t=1}^T \Pr(\text{Adv}_t(\tau_t) = \mathbf{v}_t) d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \quad (87)$$

$$= \sum_{\mathbf{x} \in \mathcal{X}} \int_{\mathcal{O}^T \times \mathcal{V}^T} \frac{\Pr(\mathbf{x})^\alpha}{D} \cdot \prod_{t=1}^T \Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t) \prod_{t=1}^T \Pr(\text{Adv}_t(\tau_t) = \mathbf{v}_t) \cdot \quad (88)$$

$$\prod_{t=1}^T \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^\alpha} d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \quad (89)$$

$$\leq \prod_{t=1}^T \left(\sum_{\mathbf{x} \in \mathcal{X}} \int_{\mathcal{O}^T \times \mathcal{V}^T} \frac{\Pr(\mathbf{x})^\alpha}{D} \cdot \prod_{t=1}^T \Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t) \prod_{t=1}^T \Pr(\text{Adv}_t(\tau_t) = \mathbf{v}_t) \cdot \frac{(\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{p_t \alpha}}{(\Pr(W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t))^{p_t \alpha}} d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \right)^{\frac{1}{p_t}} \quad (90)$$

The inequality is from Holder's inequality. So we have

$$\delta_\rho^\alpha \leq \prod_{t=1}^T r_t^{\frac{1}{p_t}}. \quad (91)$$

P Experiment

We run two groups of experiments.

P.1 AES

Leakage model. We consider a single-round leakage experiment based on the AES S-box with Hamming-weight leakage. The secret \mathbf{x} is a binary string drawn from either a 256-bit keyspace or a 128-bit keyspace, always under the uniform prior:

$$\mathbf{x} \sim \text{Unif}(\{0, 1\}^{256}) \quad \text{or} \quad \mathbf{x} \sim \text{Unif}(\{0, 1\}^{128}).$$

We write $\mathbf{x} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$ as a sequence of bytes, where $n = 32$ in the 256-bit case and $n = 16$ in the 128-bit case.

Let $S(\cdot)$ denote the AES S-box and $\text{HW}(\cdot)$ the Hamming weight of an 8-bit value. For a given secret X , we model the one-round “raw” leakage score as an aggregate of per-byte S-box Hamming weights.

Inference criteria ρ . Given the noisy observation Y , the attacker outputs a guess \hat{x} for the secret. We evaluate two Identification criteria:

- *Identification:* the attack is counted as successful if $\hat{x} = x$.
- *Miss bits less than 16 bit:* the attack is counted as successful if the Hamming distance between the guess and the true secret is at most 16 bits.

Let $m(X)$ be the deterministic function of the secret. For each choice of noise variance σ^2 , we fix the Gaussian mechanism $Y = m(X) + e$ with $e \sim \mathcal{N}(0, \sigma^2)$ and evaluate the identification criterion ρ on the resulting channel. We first compute the *ground-truth* optimal failure level δ_ρ by explicitly evaluating the posterior over the secret space and taking the best possible adversarial decision rule. We then compare δ_ρ with three information-theoretic upper bounds on the same quantity: (i) a mutual-information-based bound obtained from Fano’s inequality; (ii) a PAC-privacy bound with $\alpha = 20$ using a Gaussian reference distribution W_{Gauss} ; and (iii) our prior-reference-weighted α -information bound also with $\alpha = 20$ instantiated with the optimal reference distribution W^* . This setup allows us to assess how tightly each analytic bound approximates the true posterior success rate.

In Figures 2–4, the x -axis is the standard deviation σ of the Gaussian noise, and the y -axis is \log_2 of the corresponding upper bound on the posterior success rate, computed under each of the three methods above.

Two qualitative trends are evident:

1. The mutual-information-based bound via Fano is the loosest, while the PAC-privacy bound and our prior-reference-weighted α -information bound are substantially tighter. The gap is exponential in the success probability, i.e., linear in the \log_2 -scale of the plots.
2. As the entropy of the secret increases, the gap between the mutual-information bound and the α -divergence-based bounds widens. In this regime, the α -divergence approach provides a much stronger guarantee than the mutual-information bound.

P.2 RSA

We call Algorithm 2 to instantiate our mechanism to mitigate leakage from the running time in RSA encryption. In this experiment, the attacker observes variations in runtime caused by modular exponentiation and attempts to reconstruct the secret key from the leakage. Following our implementation, the timing oracle is simulated by measuring the end-to-end runtime of an RSA pipeline: deterministic key generation from a secret seed plus encryptions the messages given by the adversary.

The secret \mathbf{x} is a binary string drawn uniformly at random from $\{0,1\}^n$, where the key size n varies from 62 to 64 bits. For each round t , we report the average accumulated noise variance in the form

$$\frac{1}{t} \sum_{i=1}^t \sigma_i^2.$$

The inference criterion ρ is identification: the attack is successful if and only if $\hat{\mathbf{x}} = \mathbf{x}$. We fix the target posterior success rate to $\delta_\rho = 2^{-60}$ and compute the required noise level in our algorithm. We evaluate our composition theorem for $\alpha \in \{50, 75, 100\}$ and compute the per-round average noise variance under each setting. All reported noise values are averaged over 100 independent trials and the composition number T is set to 8, as shown in Fig. 6, Fig. 7, and Fig. 8.

We have two key observations:

1. The required noise grows approximately linearly (after averaging across rounds). This contrasts with Theorem 10, where the required noise grows quadratically. This difference suggests that allowing the mechanism to adaptively calibrate noise based on previous outputs and adversarial queries can significantly reduce overall noise accumulation.
2. As α increases, the marginal reduction in required noise becomes progressively smaller. This indicates that the noise level converges to a minimal value as α grows. Empirically, this supports that the bound becomes increasingly tight as $\alpha \rightarrow \infty$.

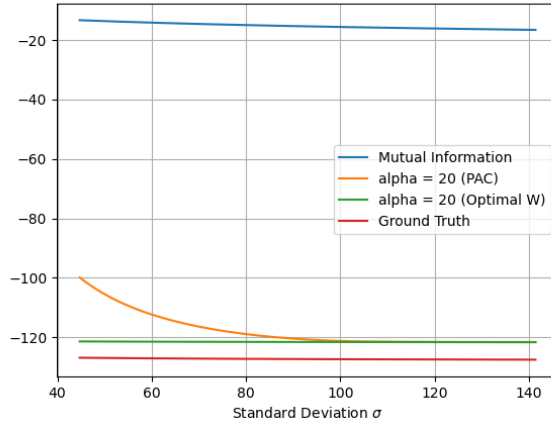


Fig. 2. \log_2 Posterior Success Bound vs. standard deviation σ (128-bit Secret, Identical Inference)

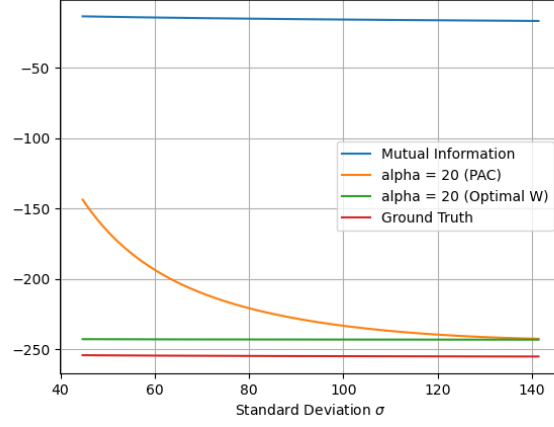


Fig. 3. \log_2 Posterior Success Bound vs. standard deviation σ (256-bit Secret, Identical Inference)

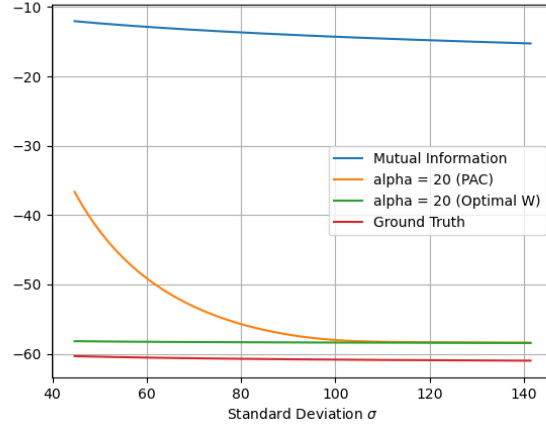


Fig. 4. \log_2 Posterior Success Bound vs. standard deviation σ (128-bit Secret, ≤ 16 -Bit Error)

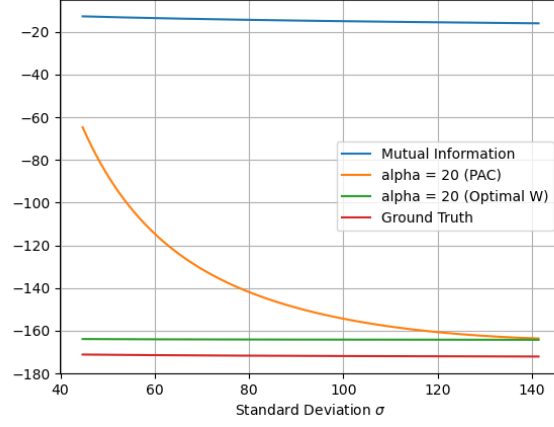


Fig. 5. \log_2 Posterior Success Bound vs. standard deviation σ (256-bit Secret, ≤ 16 -Bit Error)

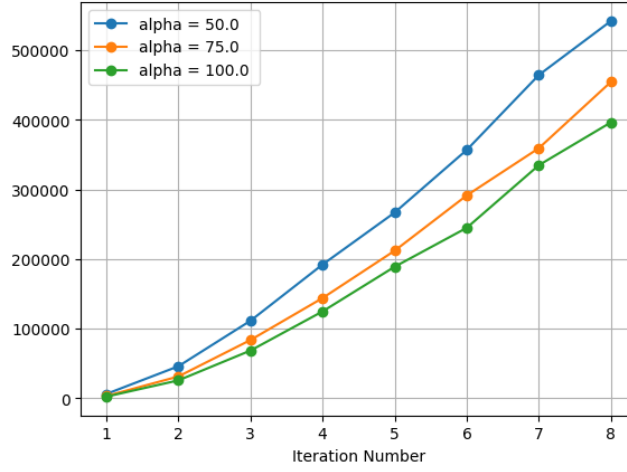


Fig. 6. Per-Round Average Noise vs. Iteration Number (64-bit Secret, Target Posterior Success $\leq 2^{-60}$)

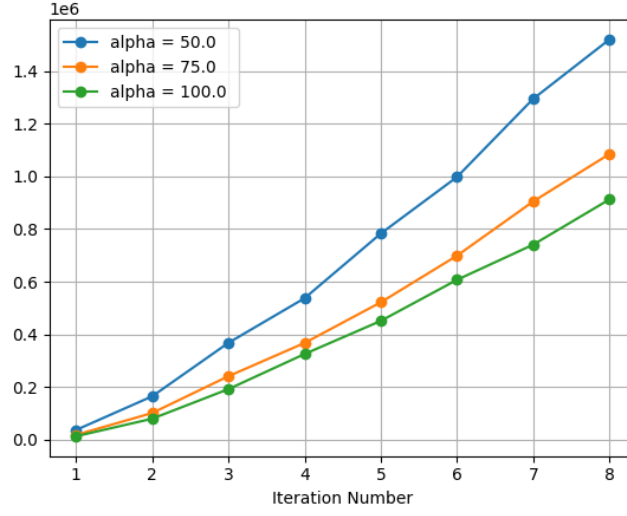


Fig. 7. Per-Round Average Noise vs. Iteration Number (63-bit Secret, Target Posterior Success $\leq 2^{-60}$)

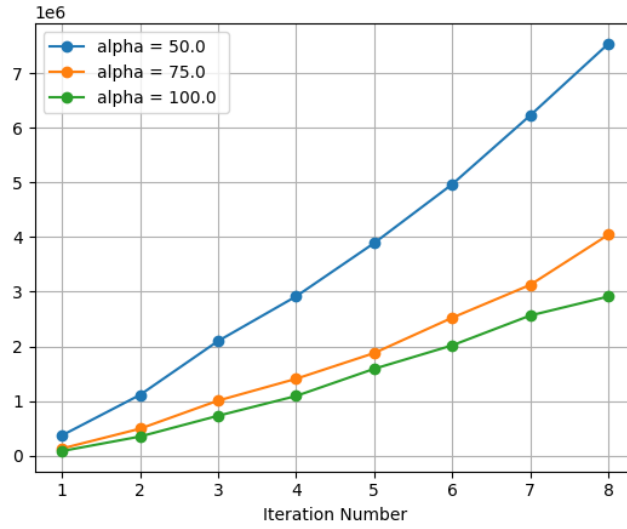


Fig. 8. Per-Round Average Noise vs. Iteration Number (62-bit Secret, Target Posterior Success $\leq 2^{-60}$)

Q Analogous Results based on KL Divergence

For the KL divergence, we apply Lemma 2 with $f(u) = u \log u$ to bound δ_ρ . We use the similar technique: consider further decomposing the global control problem into a collection of *local* subproblems, each of which depends only on the current branch and can be solved on the local branch.

Theorem 12 (Adaptive Composition of KL-divergence).

Under the setting of Definition 8 with a single secret $\mathbf{x} \sim \mathcal{D}$ reused across all T rounds. For any possible query sequence, consider an auxiliary reference process

$$W_{1:T}(\mathbf{v}_{1:T}, \tau_{1:T}) := (W_1(\mathbf{v}_1, \tau_1), \dots, W_T(\mathbf{v}_T, \tau_T)) \in \mathcal{O}^T$$

which is independent of \mathbf{x} .

Define the prefix coefficient

$$A_{t-1}(\hat{\mathbf{x}}) := \prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{\mathbf{x}}, \mathbf{v}_s, \tau_s) = \mathbf{o}_s)}{\Pr(W_s(\mathbf{v}_s, \tau_s) = \mathbf{o}_s)}$$

and

$$D_t(\hat{\mathbf{x}}) := \mathcal{D}_{KL}(\mathcal{M}_t(\hat{\mathbf{x}}, \mathbf{v}_t, \tau_t) \parallel W_t(\mathbf{v}_t, \tau_t)).$$

If for every possible transcript τ_t there exists such a choice of reference process $W_{1:T}(\mathbf{v}_{1:T})$ such that for all $t = 1, \dots, T$,

$$\mathbb{E}_{\mathbf{x} \sim D} [A_{t-1}(\mathbf{x}) D_t(\mathbf{x})] \leq r_t \quad (92)$$

then the final interactive mechanism \mathcal{M}^{Adv} satisfies for any Adv , the optimal posterior success rate δ_ρ satisfies

$$D_{KL}(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{o,\rho}}) \leq \sum_{t=1}^T r_t. \quad (93)$$

The KL condition above admits an intuitive weighted-sum interpretation. For each round t , the integrand is the conditional KL divergence between the round- t output law $\Pr(\mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \cdot)$ and the corresponding reference conditional law $\Pr(W_t(\mathbf{v}_t, \tau_t) = \cdot)$. This conditional KL term is multiplied by the prefix likelihood ratio $A_{t-1}(\mathbf{x})$ which can be viewed as an importance weight induced by the past transcript.

Just like the case for α -divergence, implementing Theorem 12 in practice requires specifying a concrete reference process $W_{1:T}(\mathbf{v}_{1:T})$. For KL, the best choice would be the marginal output law induced by the mechanism, since it yields the tightest divergence bound. To be more specific, we have the following theorem:

Theorem 13 (Optimal selection of reference for KL divergence). *For any $\mathbf{x} \sim D$ and mechanism $\mathcal{M}(\mathbf{x})$, let W denote the marginal distribution of $\mathcal{M}(\mathbf{x})$ over $\mathbf{x} \sim D$, then for arbitrary distribution Q , we have*

$$\mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathcal{M}(\mathbf{x})} \parallel P_W) \leq \mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathcal{M}(\mathbf{x})} \parallel P_Q) \quad (94)$$

So we know the optimal reference distribution W for KL divergence is the marginal distribution of $\mathcal{M}(\mathbf{x})$.

Finally, to further make Theorem 12 implementable from finite sampling, we establish Algorithm 3 (similar to Algorithm 2) by iteratively running Algorithm 1.

We provide the provable guarantees of Algorithm 2 in Theorem 9.

Theorem 14 (Provable guarantees of Algorithm 3). *Define the prefix coefficient*

$$A_{t-1}(\hat{\mathbf{x}}) := \prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{\mathbf{x}}, \mathbf{v}_s, \tau_s) = \mathbf{o}_s)}{\Pr(W_s(\mathbf{v}_s, \tau_s) = \mathbf{o}_s)}$$

$$D_t(\hat{\mathbf{x}}) := \mathcal{D}_{KL}(\mathcal{M}_t(\hat{\mathbf{x}}, \mathbf{v}_t, \tau_t) \parallel W_t(\mathbf{v}_t, \tau_t)),$$

and the privacy-loss function

$$\mathcal{G}_t(\mathcal{D}_e, \hat{\mathbf{x}}) := A_{t-1}(\hat{\mathbf{x}}) D_t(\hat{\mathbf{x}}) \quad (95)$$

Suppose that for every $\hat{\mathbf{x}} \in \mathcal{X}$ and every admissible \mathcal{D}_e we have $\mathcal{G}_t(\mathcal{D}_e, \hat{\mathbf{x}}) \leq B(\mathcal{D}_e)$ for some finite bound $B(\mathcal{D}_e)$ given per-round budget r_t . Let $\mathcal{M}_{[1:T]}$ be the interactive mechanism produced by Algorithm 3. Then for any Adv , the optimal posterior success rate δ_ρ satisfies

$$D_{KL}(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{\mathbf{o}, \rho}}) \leq \sum_{t=1}^T r_t. \quad (96)$$

This result is similar to Theorem 9. The difference here is the one we select α divergence to control while for the other we select Kl divergence to control. Since when $\alpha \rightarrow 1$, α divergence will converges to KL, Theorem 14 can viewed as a special case of Theorem 9. Notice when $\alpha \rightarrow 1$ the bound is tight so this intuitively suggests that Kl divergence is a very loose control since "1" is much smaller than " ∞ ".

R Proof of Theorem 12

Apply Lemma 2, and we have

$$D_{KL}(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{\mathbf{o}, \rho}}) \leq \mathbb{E}_{\mathbf{x} \sim D} \int_{\mathcal{O}^T \times \mathcal{V}^T} \cdot \prod_{t=1}^T \Pr(\text{Adv}_t(\tau_t) = \mathbf{v}_t) \quad (97)$$

$$(\Pr \mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t) \cdot \prod_{t=1}^T \log\left(\frac{\Pr \mathcal{M}_t(\mathbf{x}, \mathbf{v}_t, \tau_t) = \mathbf{o}_t}{\Pr W_t(\mathbf{v}_t, \tau_t) = \mathbf{o}_t}\right) d\mathbf{o}_{[1:T]} d\mathbf{v}_{[1:T]} \quad (98)$$

Algorithm 3 Randomization Optimizer under Adversarial Composition for KL divergence

- 1: **Input:** number of rounds T ; leakage functions $\{\mathcal{F}_t : \mathcal{X} \times \mathcal{V} \rightarrow \mathcal{O}\}_{t=1}^T$; input distribution D ; per-round α -divergence budgets $\{r_t\}_{t=1}^T$; sample size m ; global bound $B(\mathcal{D}_e)$.
- 2: Sample a secret input $\hat{x} \sim D$.
- 3: Choose per-round targets $\{r'_t\}_{t=1}^T$, failure probabilities $\{\gamma_t\}_{t=1}^T$, and a fallback bound R such that for every $t \in [T]$,

$$(1 - \gamma_t) r'_t + \gamma_t R \leq r_t.$$

- 4: **for** $t = 1, 2, \dots, T$ **do**
- 5: Receive adversarial parameter v_t .
- 6: Define the *prefix coefficient*

$$A_{t-1}(\hat{x}) := \Pr(x = \hat{x})^{\alpha-1} \left(\prod_{s=1}^{t-1} \frac{\Pr(\mathcal{M}_s(\hat{x}, v_s, \tau_s) = o_s)}{\Pr(W_s(v_s, \tau_s) = o_s)} \right)^\alpha$$

$$D_t(\hat{x}) := \mathcal{D}_\alpha(\mathcal{M}_t(\hat{x}, v_t, \tau_t) \parallel W_t(v_t, \tau_t)),$$

and

$$\mathcal{G}_t(\mathcal{D}_e, \hat{x}) := A_{t-1}(\hat{x}) D_t(\hat{x})$$

- 7: Take $\mathcal{F}(\cdot, v_t), g_t(\mathcal{D}_e, \hat{x}), D, r_t, B(\mathcal{D}_e)$ as the input for Algorithm 1
 - 8: Sample noise e_t from the calibrated distribution.
 - 9: Output $o_t \leftarrow \mathcal{F}_t(\hat{x}, v_t) + e_t$.
 - 10: **end for**
-

For simplicity but without loss of generality, suppose $T = 2$, notice that inequality (98) can be rewritten as

$$D_{KL}(\mathbf{1}_{\delta_\rho} \parallel \mathbf{1}_{\delta_{o,\rho}}) \tag{99}$$

$$\begin{aligned} &\leq \mathbb{E}_{x \sim D} \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = v_1) \Pr(W_1(v_1, \tau_1) = o_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = v_2) \\ &\quad \cdot \frac{\Pr(\mathcal{M}_1(x, v_1, \tau_1) = o_1)}{\Pr(W_1(v_1, \tau_1) = o_1)} \int_{\mathcal{O}} \Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2) \left(\log \left(\frac{\Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2)}{\Pr(W_2(v_2, \tau_2) = o_2)} \right) + \log \left(\frac{\Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2)}{\Pr(W_2(v_2, \tau_2) = o_2)} \right) \right) \text{d}o_2 \\ &\tag{100} \\ &= \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = v_1) \Pr(W_1(v_1, \tau_1) = o_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = v_2) \\ &\quad \cdot \mathbb{E}_{x \sim D} \underbrace{\frac{\Pr(\mathcal{M}_1(x, v_1, \tau_1) = o_1)}{\Pr(W_1(v_1, \tau_1) = o_1)}}_{(A)} \underbrace{\int_{\mathcal{O}} \Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2) \log \left(\frac{\Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2)}{\Pr(W_2(v_2, \tau_2) = o_2)} \right) \text{d}o_2}_{(B)} \text{d}v_2 \text{d}o_1 \text{d}v_1 \\ &+ \underbrace{\int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = v_1) \mathbb{E}_{x \sim D} \Pr(\mathcal{M}_1(x, v_1, \tau_1) = o_1) \log \left(\frac{\Pr(\mathcal{M}_1(x, v_1, \tau_1) = o_1)}{\Pr(W_1(v_1, \tau_1) = o_1)} \right)}_{(C)} \\ &\quad \underbrace{\int_{\mathcal{O} \times \mathcal{V}} \Pr(\mathcal{M}_2(x, v_2, \tau_2) = o_2) \Pr(\text{Adv}_1(\tau_2) = v_2) \text{d}o_2 \text{d}v_2 \text{d}o_1 \text{d}v_1}_{(D)} \tag{101} \end{aligned}$$

Notice that term (A) is $A_1(\mathbf{x})$, term (B) is $D_2(\mathbf{x})$, term (C) $\leq r_1$ and term (D) equals to 1, given the condition in the theorem, we know

$$\mathbb{E}_{\mathbf{x} \sim D} [A_1(\mathbf{x}) D_2(\mathbf{x})] \leq r_2 \quad (102)$$

so

$$\begin{aligned} (101) &\leq r_2 \cdot \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) \Pr(W_1(\mathbf{v}_1, \tau_1) = \mathbf{o}_1) \int_{\mathcal{V}} \Pr(\text{Adv}_1(\tau_2) = \mathbf{v}_2) d\mathbf{v}_2 d\mathbf{o}_1 d\mathbf{v}_1 \\ &\quad + r_1 \cdot \int_{\mathcal{O} \times \mathcal{V}} \Pr(\text{Adv}_1(\tau_1) = \mathbf{v}_1) d\mathbf{v}_1 \end{aligned} \quad (103)$$

$$= r_1 + r_2 \quad (104)$$

S Proof of Theorem 13

For each $\mathbf{x} \in \mathcal{X}$, write $P_{\mathbf{x}} := P_{\mathcal{M}(\mathbf{x})}$ for the output distribution of \mathcal{M} given input \mathbf{x} . Let W be the marginal (mixture) of these conditionals under D :

$$P_W(o) = \sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) P_{\mathbf{x}}(o), \quad o \in \mathcal{O}.$$

Let Q be any other reference distribution on \mathcal{O} . For simplicity we assume \mathcal{O} is finite or countable; the general case follows by replacing sums with integrals.

We start by expanding both sides:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathbf{x}} \| P_Q) &= \sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) \sum_{o \in \mathcal{O}} P_{\mathbf{x}}(o) \log \frac{P_{\mathbf{x}}(o)}{Q(o)}, \\ \mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathbf{x}} \| P_W) &= \sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) \sum_{o \in \mathcal{O}} P_{\mathbf{x}}(o) \log \frac{P_{\mathbf{x}}(o)}{W(o)}. \end{aligned}$$

Subtracting the second from the first gives

$$\begin{aligned} &\mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathbf{x}} \| P_Q) - \mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathbf{x}} \| P_W) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) \sum_{o \in \mathcal{O}} P_{\mathbf{x}}(o) \left(\log \frac{P_{\mathbf{x}}(o)}{Q(o)} - \log \frac{P_{\mathbf{x}}(o)}{W(o)} \right) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) \sum_{o \in \mathcal{O}} P_{\mathbf{x}}(o) \log \frac{W(o)}{Q(o)} \\ &= \sum_{o \in \mathcal{O}} \left(\sum_{\mathbf{x} \in \mathcal{X}} \Pr_D(\mathbf{x}) P_{\mathbf{x}}(o) \right) \log \frac{W(o)}{Q(o)} \\ &= \sum_{o \in \mathcal{O}} W(o) \log \frac{W(o)}{Q(o)} = \mathcal{D}_{KL}(P_W \| P_Q) \geq 0, \end{aligned}$$

where the last inequality is Gibbs' inequality for KL divergence.

Rearranging,

$$\mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathcal{M}(\mathbf{x})} \| P_W) \leq \mathbb{E}_{\mathbf{x} \sim D} \mathcal{D}_{KL}(P_{\mathcal{M}(\mathbf{x})} \| P_Q),$$

which proves that the marginal W minimizes the average KL divergence over all choices of Q .

T Proof of Theorem 14

By applying Theorem 1, we have

$$\mathbb{E}_{\mathbf{x} \sim D} [\mathcal{G}(\mathcal{D}_e, \mathbf{x})] \leq r_t.$$

Use $\mathcal{G}_t(\mathcal{D}_e, \hat{\mathbf{x}}) := A_{t-1}(\hat{\mathbf{x}}) D_t(\hat{\mathbf{x}})$ into the above equation for every t and every selection of we have

$$\mathbb{E}_{\mathbf{x} \sim D} [A_{t-1}(\mathbf{x}) D_t(\mathbf{x})] \leq r_t \tag{105}$$

Consequently the condition 92 is held so we apply Theorem 12 we have:

$$D_{KL}(\mathbf{1}_{\delta_\rho} \| \mathbf{1}_{\delta_{o,\rho}}) \leq \sum_{t=1}^T r_t. \tag{106}$$