# Differentially Private Fisher Randomization Tests for Binary Outcomes

Qingyang Sun[*]   Jerome P. Reiter[†]

Department of Statistical Science, Duke University

November 27, 2025

## Abstract

Across many disciplines, causal inference often relies on randomized experiments with binary outcomes. In such experiments, the Fisher randomization test provides exact, assumption-free tests for causal effects. Sometimes the outcomes are sensitive and must be kept confidential, for example, when they comprise physical or mental health measurements. Releasing test statistics or $p$-values computed with the confidential outcomes can leak information about the individuals in the study. Those responsible for sharing the analysis results may wish to bound this information leakage, which they can do by ensuring the released outputs satisfy differential privacy. In this article, we develop and compare several differentially private versions of the Fisher randomization test for binary outcomes. Specifically, we consider direct perturbation approaches that inject calibrated noise into test statistics or $p$-values, as well as a mechanism-aware, Bayesian denoising framework that explicitly models the privacy mechanism. We further develop decision-making procedures under privacy constraints, including a Bayes risk-optimal rule and a frequentist-calibrated significance test. Through theoretical results, simulation studies, and an application to the ADAPTABLE clinical trial, we demonstrate that our methods can achieve valid and interpretable causal inference while ensuring the differential privacy guarantee.

*Keywords:* Confidentiality; Decision; Experiment; Hypothesis; Privacy.

---

[*]qingyang.sun@duke.edu, 214 Old Chemistry, Box 90251, Durham, NC 27708-0251, USA
[†]jreiter@duke.edu, 214 Old Chemistry, Box 90251, Durham, NC 27708-0251, USA

# 1 Introduction

Randomization-based inference is a cornerstone of causal analysis, offering inferences that derive solely from the randomization of treatment assignments to study subjects. Among such methods, the Fisher randomization test (FRT) is perhaps the most widely used approach (Fisher 1935). In the FRT, one presumes Fisher's sharp null hypothesis: each individual's outcome is the same regardless of treatment assignment. Under this sharp null, one can use the observed outcomes to compute the value of the chosen test statistic for any possible randomization, and thereby construct a reference distribution for the test statistic. The analyst compares the observed value of the test statistic to this reference distribution, resulting in a $p$-value under the null hypothesis. In this article, we consider FRTs for binary outcomes in completely randomized experiments.

In many causal studies, the outcomes are inherently sensitive and therefore should be kept confidential. For example, the outcomes could encode personal information like disease status, substance use, or criminal recidivism. The literature on data privacy indicates that releasing results of any statistical analysis leaks information about the underlying study subjects. Even releasing summary statistics can introduce disclosure risks (Dinur & Nissim 2003, Dwork et al. 2017, Abowd et al. 2022). Thus, those responsible for sharing the analysis results of an FRT with confidential outcomes may want to limit the amount of information leakage.

Differential privacy (DP) has emerged as a leading framework for bounding the information leakage when releasing results of analyses (Dwork 2006, Dwork et al. 2006). It establishes a bound on how much the inclusion or exclusion of any single record can change the distribution of a released statistic. Thus, adversaries seeking to use the released statistic to learn confidential information cannot tell whether or not any particular individual (or value) was part of the data used to make the published output. In this sense, DP offers a mathematically rigorous guarantee of data confidentiality.

Several researchers have developed DP methods for causal inference and hypothesis testing; however, to our knowledge, there do not exist any DP algorithms to do Fisher randomization tests. Within causal inference, D'Orazio et al. (2015) develop DP estimators for paired mean differences. Lee et al. (2019) propose a differentially private inverse probability weighting estimator that privatizes both the propensity score model and the final treatment effect estimate. Subsequent advances extend these ideas to conditional average treatment effect estimation (Niu et al. 2022) and Bayesian inference under local DP (Ohnishi & Awan 2025). Guha & Reiter (2025) introduced private algorithms for binary-outcome causal inference that support a range of weighted average treatment effect (WATE) estimators and provide standard errors and confidence intervals using a subsample-and-aggregate strategy. Finally, Mukherjee et al. (2024) present DP regression modeling strategies to analyze randomized controlled trials. Within hypothesis testing, researchers have developed DP methods for $\chi^2$ tests for goodness-of-fit (Gaboardi et al. 2016), rank-based nonparametric tests (Couch et al. 2019), and uniformly most powerful tests for discrete models (Awan & Slavković 2018). More recent work extends these ideas to yield level-$\alpha$ private tests across a variety of settings (Kazan et al. 2023, Peña & Barrientos 2025), but they are designed for parametric or asymptotic contexts rather than for randomization-based causal inference.

We aim to close that gap by developing differentially private Fisher randomization tests for binary outcomes, which we abbreviate as DP-FRT. The central goal is to estimate and

release Fisher's randomization $p$-value under a specified privacy budget while maintaining statistical validity and interpretability. We first explore direct perturbation mechanisms that add calibrated noise to the $p$-value or test statistics. We then propose a mechanism-aware, Bayesian denoising approach that explicitly models the DP noise to recover a posterior distribution for the confidential $p$-value. Furthermore, we develop decision frameworks under both Bayesian and frequentist paradigms that translate privatized $p$-values into actionable conclusions. The Bayesian framework enables analysts to abstain from making a decision when the evidence is inconclusive and to refine or update their conclusions by spending additional privacy budget, thereby incorporating uncertainty and adaptivity into decision-making. The frequentist-calibrated framework focuses on maintaining valid inference through rigorous control of type I error, ensuring that privacy protection does not compromise the nominal significance level. Together, these frameworks enable reliable and interpretable causal conclusions with formal privacy guarantees.

The remainder of this article is organized as follows. Section 2 reviews the notation and formulation of the FRT for binary outcomes and key concepts of DP. Section 3 presents differentially private methods for releasing the $p$-value from a FRT, introducing both the direct perturbation and Bayesian denoising methods. Section 4 discusses decision-making under DP-FRT, offering Bayesian risk-optimal and frequentist-calibrated decision frameworks. Section 5 describes simulation studies and a genuine data analysis that assess the performance of the proposed methods and offer guidance for their implementation. Section 6 concludes with a discussion and potential extensions.

# 2   Background

In Section 2.1, we review the Fisher randomization test, particularly in the context of completely randomized experiments with binary outcomes. In Section 2.2, we review several key concepts and properties of differential privacy.

## 2.1   Randomization Test for Binary Outcomes

### 2.1.1   Potential Outcomes and Data Representation

We first introduce some notations under the potential outcome framework for causal inference (Rubin 1974). Suppose there are $n$ units in the experiment. For each unit $i \in \{1, \ldots, n\}$ in the study, let $Z_i = 1$ when unit $i$ is assigned to the treatment and $Z_i = 0$ when the unit is assigned the control, and let $(Y_i(1), Y_i(0))$ denote the pair of potential outcomes under treatment and control, respectively. Under the stable unit treatment value assumption (SUTVA), the observed outcome is

$$Y_i^{\text{obs}} = Z_i Y_i(1) + (1 - Z_i) Y_i(0) = Y_i(0) + Z_i \{Y_i(1) - Y_i(0)\}.$$

Thus, $Y_i^{\text{obs}}$ equals the potential outcome corresponding to the realized treatment: if $Z_i = 1$ then $Y_i^{\text{obs}} = Y_i(1)$, and if $Z_i = 0$ then $Y_i^{\text{obs}} = Y_i(0)$.

The full collection $\{(Y_i(1), Y_i(0))\}_{i=1}^{n}$ is often referred to as the "Science," which represents the underlying object of interest in causal inference. The Science can be equivalently

represented as a table of unit-level potential outcomes:

| Unit $i$ | $(Y_i(1),\ Y_i(0))$ |
|---|---|
| 1 | $(Y_1(1),\ Y_1(0))$ |
| 2 | $(Y_2(1),\ Y_2(0))$ |
| $\vdots$ | $\vdots$ |
| $n$ | $(Y_n(1),\ Y_n(0))$ |

For binary outcomes, it is often more convenient to summarize the Science by the joint distribution of $(Y_i(1), Y_i(0))$, which leads to the so-called Science table. Let

$$N_{jk} = \#\{i : Y_i(1) = j, Y_i(0) = k\}, \quad j, k \in \{0, 1\},$$

with row and column sums $N_{1+}, N_{0+}, N_{+1}, N_{+0}$ satisfying $\sum_{j,k} N_{jk} = n$. The Science table is then

| | $Y_i(0) = 1$ | $Y_i(0) = 0$ | Row sum |
|---|---|---|---|
| $Y_i(1) = 1$ | $N_{11}$ | $N_{10}$ | $N_{1+}$ |
| $Y_i(1) = 0$ | $N_{01}$ | $N_{00}$ | $N_{0+}$ |
| Col sum | $N_{+1}$ | $N_{+0}$ | $n$ |

Crucially, the Science table is never directly observed, since only one of $\{Y_i(1), Y_i(0)\}$ is realized for each unit $i$. What we do observe is the pair $(Z_i, Y_i^{\text{obs}})$, which can be aggregated into the following $2 \times 2$ Outcome table

| | $Y_i^{\text{obs}} = 1$ | $Y_i^{\text{obs}} = 0$ | Row sum |
|---|---|---|---|
| $Z_i = 1$ | $n_{11}$ | $n_{10}$ | $n_1$ |
| $Z_i = 0$ | $n_{01}$ | $n_{00}$ | $n_0$ |
| Col sum | $n_{+1}$ | $n_{+0}$ | $n$ |

Here, $n_{11}$ and $n_{10}$ are the numbers of treated units with $Y_i^{\text{obs}} = 1$ and $Y_i^{\text{obs}} = 0$, respectively, and $n_{01}$ and $n_{00}$ are the corresponding counts in the control group. Accordingly, $n_1$ and $n_0$ are the sizes of the treatment and control groups, and $n_{+1}$ and $n_{+0}$ are the marginal counts of each outcome.

### 2.1.2 Fisher Randomization Test

The fundamental problem of causal inference lies in the fact that only one potential outcome is observed for each unit. Within the potential outcomes framework, randomized experiments address this missing data problem by exploiting the known treatment assignment mechanism. Specifically, let $\boldsymbol{Z} = (Z_1, \ldots, Z_n)^{\text{T}}$ be the stochastic treatment assignment vector and $\boldsymbol{Y}^{\text{obs}} = (Y_1^{\text{obs}}, \ldots, Y_n^{\text{obs}})^{\text{T}}$ the observed outcome vector. The assignment mechanism links the unobserved Science table to the observed Outcome table through the distribution of $\boldsymbol{Z}$ known by design. This connection enables exact, finite-sample valid, model-free inference. One of the canonical designs is the Completely Randomized Experiment (CRE).

**Definition 2.1** (CRE). Fix $n_1$ treated units and $n_0$ controlled units with $n = n_1 + n_0$. Define the set of admissible assignments in CRE

$$\mathcal{Z} = \left\{ \boldsymbol{z} \in \{0, 1\}^n : \sum_{i=1}^{n} z_i = n_1 \right\}, \qquad |\mathcal{Z}| = \binom{n}{n_1}.$$

The assignment mechanism is uniform over all $\mathcal{Z}$, i.e.,

$$\Pr(\boldsymbol{Z} = \boldsymbol{z}) = \frac{1}{|\mathcal{Z}|}, \qquad \text{for all } \boldsymbol{z} \in \mathcal{Z},$$

so that all randomness in the observed data arises solely from $\boldsymbol{Z}$.

Under the design of CRE, Fisher (1935) proposed to test the sharp null hypothesis of no individual causal effects.

**Definition 2.2** (Fisher's Sharp Null Hypothesis).

$$H_0^{\mathrm{F}} : Y_i(1) = Y_i(0) \qquad \text{for all } i = 1, \dots, n.$$

The hypothesis $H_0^{\mathrm{F}}$ is referred to as "sharp" because it enables the imputation of all missing potential outcomes via $Y_i(1) = Y_i(0) = Y_i^{\mathrm{obs}}$, thus the entire Science table becomes observable. Although Fisher's sharp null has been criticized for being restrictive, subsequent research has extended the framework to test weak null hypotheses (Wu & Ding 2021). We will discuss possible extensions to weak nulls in Section 6. The Fisher randomization test under $H_0^{\mathrm{F}}$ is defined as follows.

**Definition 2.3** (FRT). Under CRE and $H_0^{\mathrm{F}}$, the FRT proceeds in the following steps:

1. **Choose a test statistic.** Select a statistic $T(\boldsymbol{Z}; \boldsymbol{Y}^{\mathrm{obs}})$ that reflects deviations from $H_0^{\mathrm{F}}$, such as a difference in means, $t$-statistic, or a rank-based measure.

2. **Compute the observed value.** Let $\boldsymbol{z}^{\mathrm{obs}}$ denote the realized assignment. The observed statistic is
$$T^{\mathrm{obs}} = T(\boldsymbol{z}^{\mathrm{obs}}; \boldsymbol{Y}^{\mathrm{obs}}).$$

3. **Generate the randomization distribution.** Under the known assignment mechanism of CRE, evaluate $T(\boldsymbol{z}; \boldsymbol{Y}^{\mathrm{obs}})$ for each $\boldsymbol{z} \in \mathcal{Z}$ to obtain the randomization distribution under $H_0^{\mathrm{F}}$ for reference.

4. **Compute the $p$-value.** The (one-sided) Fisher's randomization $p$-value is

$$p_{\mathrm{FRT}} = \Pr\left\{T(\boldsymbol{z}; \boldsymbol{Y}^{\mathrm{obs}}) \geq T^{\mathrm{obs}}\right\} = \frac{1}{|\mathcal{Z}|} \sum_{\boldsymbol{z} \in \mathcal{Z}} \mathbf{1}\{T(\boldsymbol{z}; \boldsymbol{Y}^{\mathrm{obs}}) \geq T^{\mathrm{obs}}\}.$$

When $|\mathcal{Z}| = \binom{n}{n_1}$ is too large to enumerate, one common strategy is to approximate $p_{\mathrm{FRT}}$ by Monte Carlo: randomly draw $\boldsymbol{z}^{(1)}, \dots, \boldsymbol{z}^{(R)} \overset{\text{i.i.d.}}{\sim} \mathrm{Unif}(\mathcal{Z})$ for some large $R$ and compute

$$\tilde{p}_{\mathrm{FRT}} = \frac{1 + \sum_{r=1}^{R} \mathbf{1}\left\{T\left(\boldsymbol{z}^{(r)}; \boldsymbol{Y}^{\mathrm{obs}}\right) \geq T^{\mathrm{obs}}\right\}}{1 + R},$$

where the add-one correction ensures a strictly positive $p$-value.

For binary outcomes, the implications of FRT are particularly straightforward. Recall from the Outcome table in Section 2.1.1 that $n_1, n_0$ denote the treatment and control group sizes, and $n_{+1}, n_{+0}$ the totals for outcomes $Y_i^{\mathrm{obs}} = 1$ and $Y_i^{\mathrm{obs}} = 0$. Under the sharp null

hypothesis $H_0^{\mathrm{F}}$, the number of treated successes $n_{11}$ follows a hypergeometric distribution with probability mass function

$$n_{11} \sim \mathrm{Hypergeometric}(n,\ n_{+1},\ n_1), \quad \Pr(n_{11} = a) = \frac{\binom{n_{+1}}{a}\binom{n_{+0}}{n_1 - a}}{\binom{n}{n_1}},$$

where $a = \max\{0,\ n_1 - n_{+0}\}, \ldots, \min\{n_1,\ n_{+1}\}$.

Moreover, as noted by Ding & Dasgupta (2016), commonly used test statistics are monotonic in $n_{11}$ and thus yield identical rejection regions. A natural choice is the difference-in-proportions statistic, $\widehat{\tau} = n_{11}/n_1 - n_{01}/n_0$. In this case, the FRT coincides numerically with Fisher's exact test for $2 \times 2$ tables, though its validity relies on the known randomization mechanism rather than any distributional assumptions.

## 2.2   Differential Privacy

Agencies and researchers often need to release summaries of sensitive data while limiting disclosure risks. Traditional statistical disclosure control techniques can be effective in practice but typically lack formal, data-agnostic guarantees. Differential privacy (DP) addresses this limitation by ensuring that the inclusion or exclusion of a single individual's data has a limited impact on the output distribution, regardless of any auxiliary information an adversary may possess (Dwork 2006).

Formally, we view a data-release procedure as a randomized algorithm $\mathcal{M}$ that takes a dataset $D$ as input and produces a randomized output. The privacy guarantee is defined with respect to neighboring datasets, i.e., datasets that differ in the data of a single individual, such as by adding, removing, or modifying one record. Intuitively, DP requires that the output distributions produced by $\mathcal{M}$ on any pair of neighboring datasets be nearly indistinguishable. Consequently, it becomes difficult for an adversary to determine whether a specific individual is present in the dataset or to infer sensitive attributes with high confidence.

**Definition 2.4** ($\epsilon$-Differential Privacy (Dwork 2006)). A randomized algorithm $\mathcal{M}$ satisfies $\epsilon$-differential privacy ($\epsilon$-DP) if, for any pair of neighboring datasets $D$ and $D'$ and for all measurable subsets $S \subseteq \mathcal{R}(\mathcal{M})$,

$$\Pr\{\mathcal{M}(D) \in S\} \leq e^\epsilon \Pr\{\mathcal{M}(D') \in S\},$$

where the probabilities are taken over the randomness in $\mathcal{M}$, and $\mathcal{R}(\mathcal{M})$ denotes its range of possible outputs.

The privacy parameter $\epsilon > 0$, referred to as the privacy budget, quantifies the worst-case multiplicative difference between the output distributions over neighboring datasets, where smaller values of $\epsilon$ indicate stronger privacy protection. Extensions such as $(\epsilon, \delta)$-DP allow a small failure probability $\delta$ in exchange for improved utility; please refer to Dwork & Roth (2014) for more details.

To ensure DP, a typical approach is to add random noise calibrated to the sensitivity of the released quantity. A common $\epsilon$-DP mechanism used in continuous domains is the Laplace mechanism (Dwork 2006). Let $f : \mathcal{D} \to \mathbb{R}^d$ be a function defined on datasets. The $\ell_1$-sensitivity of $f$ is defined as $\Delta f = \max_{D,D'} \|f(D) - f(D')\|_1$, where the maximum is taken over all pairs of neighboring datasets $D$ and $D'$.

**Definition 2.5** (Laplace Mechanism (Dwork 2006))**.** Suppose $f(D) \in \mathbb{R}$ is a real-valued function with $\ell_1$-sensitivity $\Delta f$. The Laplace mechanism releases

$$\tilde{f}(D) = f(D) + \eta, \qquad \text{where } \eta \sim \text{Lap}\left(0, \frac{\Delta f}{\epsilon}\right),$$

with probability density function

$$p_\eta(h) = \frac{\epsilon}{2\Delta f} \exp\left(-\frac{\epsilon|h|}{\Delta f}\right), \quad h \in \mathbb{R}.$$

If $f(D) \in \mathbb{R}^d$, independent Laplace noise is added to each coordinate.

In contexts involving count queries or discrete outputs, the (two-sided) Geometric mechanism (Ghosh et al. 2012) is particularly suited as it adds integer-valued noise, which can be viewed as the discrete analogue of the Laplace mechanism.

**Definition 2.6** (Geometric Mechanism (Ghosh et al. 2012))**.** Suppose $f(D) \in \mathbb{Z}$ is an integer-valued function with $\ell_1$-sensitivity $\Delta f$, and let $\rho = \exp\{-\epsilon/\Delta f\}$. The geometric mechanism releases

$$\tilde{f}(D) = f(D) + \eta, \qquad \text{where } \eta \sim \text{Geom}(\rho),$$

with probability mass function

$$\Pr(\eta = h) = \frac{1 - \rho}{1 + \rho} \cdot \rho^{|h|}, \quad h \in \mathbb{Z}.$$

If $f(D) \in \mathbb{Z}^d$, independent geometric noise is added to each coordinate.

Two fundamental properties of DP are particularly useful when designing private algorithms. Proposition 2.7 ensures that the privacy guarantee is preserved under any transformation of the output, as long as the transformation is independent of the underlying dataset. Proposition 2.8 summarizes how privacy guarantees behave under both sequential and parallel compositions of multiple DP mechanisms.

**Proposition 2.7** (Post-Processing Invariance (Dwork et al. 2006))**.** Let $\mathcal{M}$ be an $\epsilon$-DP mechanism, and let $g$ be any (possibly randomized) function that does not depend on the dataset. Then the composed mechanism $g \circ \mathcal{M}$ also satisfies $\epsilon$-DP.

**Proposition 2.8** (Composition Properties (McSherry 2009))**.** Let $\mathcal{M}_1, \mathcal{M}_2, \ldots, \mathcal{M}_k$ be mechanisms applied to datasets.

(a) Sequential Composition: If each $\mathcal{M}_i$ is applied to the same dataset $D$ and satisfies $\epsilon_i$-DP, then their sequential composition satisfies $\left(\sum_{i=1}^k \epsilon_i\right)$-DP.

(b) Parallel Composition: If each $\mathcal{M}_i$ is applied to a disjoint subset of the dataset $D$ and satisfies $\epsilon_i$-DP, then the overall mechanism satisfies $(\max_i \epsilon_i)$-DP.

We wrap up this subsection by highlighting the fundamental privacy-utility tradeoff inherent in DP data analysis. Intuitively, a smaller value of $\epsilon$ provides stronger privacy guarantees but necessitates adding larger noise, which would reduce the accuracy or utility of the released data. Given a fixed privacy budget, the key challenge lies in designing mechanisms that maximize the utility of the output for the intended analytical task. In other words, the goal is to enable valid population-level statistical inference from privatized noisy data, closely approximating the results that would be obtained without privacy constraints.

# 3 Methods for Differentially Private Estimation of the Fisher's randomization $p$-value

This section presents $\epsilon$-DP approaches for privatizing the $p$-value of FRT. Section 3.1 introduces direct perturbation methods for the $p$-values and the test statistics. Section 3.2 develops a mechanism-aware Bayesian denoising framework with uncertainty quantification for the privatized $p$-values.

## 3.1 Direct Perturbation Approaches

### 3.1.1 Perturbation of $p$-value

The most straightforward approach is to directly perturb the exact $p$-value by adding calibrated Laplace noise. Recall that under CRE, the treatment group sizes $(n_1, n_0)$ are fixed with $n = n_1 + n_0$, and the assignment space is $\mathcal{Z} = \{ \boldsymbol{z} \in \{0, 1\}^n : \sum_i z_i = n_1 \}$ with realized assignment $\boldsymbol{z}^{\mathrm{obs}}$. We use substitution adjacency, where neighbors $D \sim D'$ differ only at one coordinate $Y_j^{\mathrm{obs}} \in \{0, 1\}$, while $(n_1, n_0)$ and $\boldsymbol{z}^{\mathrm{obs}}$ are fixed by the design of CRE. For simplicity, we let $p_{\mathrm{FRT}}$ denote the Fisher's randomization $p$-value for the difference-in-proportions statistic $\widehat{\tau}$, and we abbreviate its $\ell_1$-sensitivity as $\Delta_p$.

**Lemma 3.1** (Sensitivity of the Fisher's randomization $p$-value). *Under CRE with binary outcomes and the test statistic $\widehat{\tau}$, the $\ell_1$-sensitivity of $p_{\mathrm{FRT}}$ is $\Delta_p = \max \left\{ \dfrac{n_1}{n}, \dfrac{n_0}{n} \right\}$.*

*Proof.* Define $A(\boldsymbol{z}) = \sum_{i=1}^n z_i Y_i^{\mathrm{obs}}$ and $a = A(\boldsymbol{z}^{\mathrm{obs}})$. Since $n_{01}(\boldsymbol{z}) = n_{+1} - A(\boldsymbol{z})$, we have $\widehat{\tau}(\boldsymbol{z}; \boldsymbol{Y}^{\mathrm{obs}}) = [(1/n_1) + (1/n_0)] A(\boldsymbol{z}) - n_{+1}/n_0$, hence $\widehat{\tau}$ is strictly increasing in $A(\boldsymbol{z})$. Therefore $p_{\mathrm{FRT}} = |\mathcal{Z}|^{-1} \sum_{\boldsymbol{z} \in \mathcal{Z}} \mathbf{1}\{A(\boldsymbol{z}) \geq a\}$.

Let $D$ and $D'$ differ only at unit $j$, and set $s = Y_j' - Y_j \in \{-1, +1\}$. Then for $D'$, $A'(\boldsymbol{z}) = A(\boldsymbol{z}) + s z_j$ and $a' = a + s z_j^{\mathrm{obs}}$. If $z_j = z_j^{\mathrm{obs}}$ the indicator is unchanged; if $z_j \neq z_j^{\mathrm{obs}}$ it can change by at most one in absolute value. Averaging over $\mathcal{Z}$ gives

$$|p_{\mathrm{FRT}}(D) - p_{\mathrm{FRT}}(D')| \leq \Pr_{\boldsymbol{Z} \sim \mathrm{Unif}(\mathcal{Z})} (Z_j \neq z_j^{\mathrm{obs}}) = \begin{cases} n_1/n, & z_j^{\mathrm{obs}} = 0, \\ n_0/n, & z_j^{\mathrm{obs}} = 1. \end{cases}$$

Maximizing over $j$ yields $\Delta_p \leq \max\{n_1/n, n_0/n\}$. Tightness holds in two extremal constructions. If all $Y_i^{\mathrm{obs}} = 0$ and $z_j^{\mathrm{obs}} = 1$, then $p_{\mathrm{FRT}}(D) = 1$, while after flipping $Y_j^{\mathrm{obs}}$ to 1, we obtain $p_{\mathrm{FRT}}(D') = \Pr(Z_j = 1) = n_1/n$. So the gap is $n_0/n$. Symmetrically, if all $Y_i^{\mathrm{obs}} = 1$ and $z_j^{\mathrm{obs}} = 0$, the gap is $n_1/n$. In particular, under the balanced design $n_1 = n_0 = n/2$, the sensitivity is $\Delta_p = 1/2$. $\qquad\square$

Based on Lemma 3.1, we release privatized $p$-values by adding Laplace noise calibrated to the sensitivity and clipping to the feasible range.

**Theorem 3.2** (Laplace mechanism for the Fisher's randomization $p$-value). *Fix a privacy budget $\epsilon > 0$ and let $\Delta_p$ be as in Lemma 3.1. Define $[L, U] = [|\mathcal{Z}|^{-1}, 1]$ and release*

$$\tilde{p} = \Pi_{[L, U]} \left( p_{\mathrm{FRT}} + \eta \right), \qquad \eta \sim \mathrm{Lap}\left( 0, \frac{\Delta_p}{\epsilon} \right),$$

*where $\Pi_{[L, U]}(x) = \min\{U, \max\{L, x\}\}$ is the clipping operator. Then $\tilde{p}$ satisfies $\epsilon$-DP.*

The privacy guarantee follows from the standard Laplace mechanism, and the post-processing invariance of DP ensures that clipping does not degrade privacy.

### 3.1.2 Perturbation of Test Statistic and its Reference

Apart from directly perturbing the $p$-value, one may privatize the test statistic and pair it with a privatized randomization distribution for reference. We start by calculating the $\ell_1$-sensitivity of the difference-in-proportions statistic $\hat{\tau} = n_{11}/n_1 - n_{01}/n_0$.

**Lemma 3.3** (Sensitivity of $\hat{\tau}$). *Under CRE with binary outcomes, the $\ell_1$-sensitivity of the statistic $\hat{\tau}$ is $\Delta_{\hat{\tau}} = \max\left\{\dfrac{1}{n_1}, \dfrac{1}{n_0}\right\}$.*

*Proof.* If $Z_j^{\mathrm{obs}} = 1$, flipping $Y_j^{\mathrm{obs}}$ changes $n_{11}$ by $\pm 1$ and keeps $n_{01}$ unchanged, so $\Delta_{\hat{\tau}} = 1/n_1$. If $Z_j^{\mathrm{obs}} = 0$, it changes $n_{01}$ by $\pm 1$ and keeps $n_{11}$ unchanged, so $\Delta_{\hat{\tau}} = 1/n_0$. Maximizing over $j$ yields the claim. In a balanced design, $\Delta_{\hat{\tau}} = 2/n$. $\qquad\square$

First, we perform separate perturbation on the observed statistic and its randomization distribution. We note that privatizing the randomization distribution is necessary since it depends on the sensitive total number of observed successes $n_{+1}$.

**Theorem 3.4** (Separate perturbation mechanism). *Fix a privacy budget $\epsilon > 0$ and choose $\epsilon_{\mathrm{obs}} > 0$ and $\epsilon_{\mathrm{ref}} > 0$ with $\epsilon_{\mathrm{obs}} + \epsilon_{\mathrm{ref}} = \epsilon$. Let $\Delta_{\hat{\tau}}$ be the $\ell_1$-sensitivity of $\hat{\tau}$ as in Lemma 3.3. Produce a privatized $p$-value $\tilde{p}$ in two steps:*
*Step 1: Perturbation of the observed statistic*

$$\tilde{T}^{\mathrm{obs}} = \Pi_{[-1,1]}\left(\hat{\tau} + \eta_{\mathrm{obs}}\right), \qquad \eta_{\mathrm{obs}} \sim \mathrm{Lap}\left(0, \frac{\Delta_{\hat{\tau}}}{\epsilon_{\mathrm{obs}}}\right).$$

*Step 2: Perturbation of the randomization distribution*

$$\tilde{n}_{+1} = \Pi_{\{0,1,\dots,n\}}\left(n_{+1} + \eta_{\mathrm{ref}}\right), \qquad \eta_{\mathrm{ref}} \sim \mathrm{Geom}(e^{-\epsilon_{\mathrm{ref}}}).$$

*Given $(\tilde{T}^{\mathrm{obs}}, \tilde{n}_{+1})$, report the private Fisher tail probability*

$$\tilde{p} = \Pr\left\{\frac{\tilde{n}_{11}}{n_1} - \frac{\tilde{n}_{+1} - \tilde{n}_{11}}{n_0} \geq \tilde{T}^{\mathrm{obs}}\right\},$$

*where $\tilde{n}_{11} \sim \mathrm{Hypergeometric}\left(n, \tilde{n}_{+1}, n_1\right)$. Then $\tilde{p}$ satisfies $\epsilon$-DP.*

*Proof.* For Step 1, the release of $\tilde{T}^{\mathrm{obs}}$ is $\epsilon_{\mathrm{obs}}$-DP by the Laplace mechanism with sensitivity $\Delta_{\hat{\tau}}$. The subsequent clipping to $[-1,1]$ is post-processing.

For Step 2, under $H_0^{\mathrm{F}}$ the randomization distribution depends on $(n, n_1, n_{+1})$ through $n_{11} \sim \mathrm{Hypergeometric}(n, n_{+1}, n_1)$. Here $n$ and $n_1$ are public by design, while $n_{+1}$ is sensitive. By privatizing $n_{+1}$ using the geometric mechanism with unit sensitivity, we obtain $\tilde{n}_{+1}$ which is $\epsilon_{\mathrm{ref}}$-DP. Based on $\tilde{n}_{+1}$ we form the private hypergeometric reference. Since $\hat{\tau}$ is strictly monotone in $n_{11}$, this induces a privatized randomization distribution for $T$ through the mapping $a \mapsto a/n_1 - (\tilde{n}_{+1} - a)/n_0$.

Finally, $\tilde{p}$ is computed as a measurable function of $(\tilde{T}^{\mathrm{obs}}, \tilde{n}_{+1})$, which is again post-processing. By sequential composition, the entire mechanism satisfies $\epsilon$-DP. $\qquad\square$

Alternatively, one may privatize all statistics simultaneously when the randomization distribution is approximated via $R$ Monte Carlo samples.

**Theorem 3.5** (Monte Carlo perturbation mechanism). *Fix the number of Monte Carlo replicates $R \in \mathbb{N}_+$ and a privacy budget $\epsilon > 0$, choose $\epsilon_{\text{obs}} > 0$ and $\epsilon_1, \ldots, \epsilon_R > 0$ with $\epsilon_{\text{obs}} + \sum_{r=1}^{R} \epsilon_r = \epsilon$. Let $\Delta_{\widehat{\tau}}$ be the $\ell_1$-sensitivity of $\widehat{\tau}$ as in Lemma 3.3. Release the privatized vector $\left( \tilde{T}^{\text{obs}}, \tilde{T}^{(1)}, \ldots, \tilde{T}^{(R)} \right)$, where*

$$\tilde{T}^{\text{obs}} = \Pi_{[-1,1]} \left( \widehat{\tau} + \eta_{\text{obs}} \right), \qquad \eta_{\text{obs}} \sim \text{Lap}\left( 0, \frac{\Delta_{\widehat{\tau}}}{\epsilon_{\text{obs}}} \right),$$

*and for each $r = 1, \ldots, R$,*

$$\tilde{T}^{(r)} = \Pi_{[-1,1]} \left( T^{(r)} + \eta_r \right), \qquad \eta_r \sim \text{Lap}\left( 0, \frac{\Delta_{\widehat{\tau}}}{\epsilon_r} \right),$$

*with $T^{(r)} = T(\boldsymbol{z}^{(r)}; \boldsymbol{Y}^{\text{obs}})$ computed at data-independent assignments $\boldsymbol{z}^{(r)} \sim \text{Unif}(\mathcal{Z})$. A private Monte Carlo p-value is obtained by*

$$\tilde{p} = \frac{1 + \sum_{r=1}^{R} \mathbf{1}\left\{ \tilde{T}^{(r)} \geq \tilde{T}^{\text{obs}} \right\}}{1 + R}.$$

*Then $\tilde{p}$ satisfies $\epsilon$-DP.*

*Proof.* Each coordinate in the released $(R + 1)$-dimensional vector is differentially private with its assigned budget by the Laplace mechanism with sensitivity $\Delta_{\widehat{\tau}}$, and clipping is post-processing. The random draws $\boldsymbol{z}^{(r)}$ are independent of the dataset and consume no privacy budget. The final p-value $\tilde{p}$ is a deterministic function of the privatized vector and is therefore post-processing. By sequential composition, the mechanism satisfies $\epsilon_{\text{obs}} + \sum_{r=1}^{R} \epsilon_r = \epsilon$ differential privacy. $\square$

We wrap up this subsection by discussing several limitations of the direct perturbation approaches, which motivate the mechanism-aware Bayesian denoising framework introduced in Section 3.2.

**(1) Utility loss due to DP noise and clipping.** The addition of DP noise can strongly distort the p-value. As shown in Lemma 3.1, the sensitivity of the p-value is at least $1/2$, which is substantial relative to its full range $[0, 1]$. Even with moderate privacy budgets, the added noise can greatly reduce accuracy. In the Monte Carlo perturbation mechanism (Theorem 3.5), all $R$ replicates need to be privatized, so the budget is split and the noise scale grows linearly with $R$, a well-known issue in DP hypothesis testing (e.g., Kim & Schrab 2023). Moreover, clipping the outputs to valid ranges preserves privacy but introduces bias and further harms utility.

**(2) Invalidity of privatized p-values.** The privatized $\tilde{p}$ may not retain the properties of a valid p-value, such as uniformity. For example, once noise is added, $\tilde{p}$ is no longer guaranteed to satisfy $\Pr(\tilde{p} \leq \alpha) \leq \alpha$ under $H_0^{\text{F}}$ for any prespecified level $\alpha \in [0, 1]$. This undermines the reliability of downstream decisions based on $\tilde{p}$.

**(3) Lack of uncertainty quantification.** These direct perturbation mechanisms do not provide uncertainty quantification for the privatized outputs. Only noisy summaries are released, with no means to account for the added randomness from DP noise. The absence of uncertainty measures limits the interpretability and transparency of the results for end users.

**(4) No support for broader synthetic inference.** Directly perturbing $p$-values or test statistics only enables the release of privatized test outcomes, without providing corresponding synthetic data or treatment effect estimates that are consistent with the privatized $p$-values. Consequently, users cannot perform follow-up synthetic inferences of their interest, limiting the utility of these methods in applications that require both hypothesis testing and effect estimation.

## 3.2   Mechanism-aware Bayesian Denoising Approach

We next present a Bayesian denoising framework that explicitly accounts for the DP mechanism. The strategy is to privatize the sufficient statistics, update a Bayesian model for the underlying true counts, and then map the posterior distribution onto the space of $p$-values. Multiple studies have demonstrated the importance of accounting for DP noise during inference (Karwa et al. 2017, Seeman et al. 2020, Nixon et al. 2022, Räisä et al. 2023). By modeling the noise distribution and propagating it through the pipeline, this approach provides a full posterior distribution of $p_{\mathrm{FRT}}$ rather than a single noisy point estimate.

Recall that under CRE with binary outcomes, the observed data can be summarized as a $2 \times 2$ table with cell counts $(n_{11}, n_{10}, n_{01}, n_{00})$. Since the treated and control group sizes $(n_1, n_0)$ are fixed by design, it suffices to only privatize the success counts $n_{11}$ and $n_{01}$. Specifically, given a privacy budget $\epsilon > 0$, we perturb $(n_{11}, n_{01})$ by

$$\tilde{\boldsymbol{n}} = (\tilde{n}_{11}, \tilde{n}_{01}) = (n_{11} + \eta_{11}, \ n_{01} + \eta_{01}), \qquad \eta_{11}, \eta_{01} \overset{\text{i.i.d.}}{\sim} \mathrm{Geom}\left(\exp(-\epsilon)\right). \tag{1}$$

Since the treatment assignment is fixed by design, modifying the outcome of a single individual can only change the success count within that individual's assigned group. Consequently, at most one of $n_{11}$ or $n_{01}$ can change by $\pm 1$. Therefore, the $\ell_1$-sensitivity of the pair $(n_{11}, n_{01})$ is one. By the geometric mechanism in Definition 2.6, the release of $(\tilde{n}_{11}, \tilde{n}_{01})$ satisfies $\epsilon$-DP.

To denoise these counts, we specify a data-independent prior $\pi$ on $(n_{11}, n_{01})$ with support $\{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$. By Bayes' rule, the posterior distribution is then

$$\mathrm{Pr}(n_{11} = a, n_{01} = b \mid \tilde{\boldsymbol{n}}) = \frac{w(a, b)}{\sum_{a'=0}^{n_1} \sum_{b'=0}^{n_0} w(a', b')}, \tag{2}$$

with weights

$$w(a, b) = \pi(a, b) \kappa_\rho(\tilde{n}_{11} - a) \kappa_\rho(\tilde{n}_{01} - b),$$

where $\kappa_\rho(h) = \dfrac{1 - \rho}{1 + \rho} \rho^{|h|}$ is the (two-sided) geometric kernel with $\rho = \exp(-\epsilon)$.

For each candidate pair of true counts $(a, b)$, let $K = a + b$ denote the corresponding total number of observed successes. Recall that the difference-in-proportions statistic $\hat{\tau} = n_{11}/n_1 - n_{01}/n_0$ is a strictly increasing function of $n_{11}$ given fixed total successes. Under

the sharp null $H_0^F$, we have $n_{11} \sim \text{Hypergeometric}(n, K, n_1)$, so the one-sided FRT $p$-value corresponding to $(a, b)$ is given by:

$$p(a, b) = \Pr(n_{11} \geq a) = \sum_{t=\max\{a, K-n_0\}}^{\min(n_1, K)} \frac{\binom{K}{t}\binom{n-K}{n_1-t}}{\binom{n}{n_1}}. \tag{3}$$

Denote $\gamma(a, b) = \Pr(n_{11} = a, n_{01} = b \mid \tilde{\boldsymbol{n}})$. The deterministic mapping $(a, b) \mapsto p(a, b)$ induces the following posterior distribution of $p_{\text{FRT}}$ given $(\tilde{n}_{11}, \tilde{n}_{01})$:

$$\Pr(p_{\text{FRT}} \in B \mid \tilde{\boldsymbol{n}}) = \sum_{a=0}^{n_1} \sum_{b=0}^{n_0} \gamma(a, b) \mathbf{1}\{p(a, b) \in B\}, \qquad B \subseteq [0, 1]. \tag{4}$$

Since the original data is only used in the construction of $(\tilde{n}_{11}, \tilde{n}_{01})$, with all later steps being post-processing, we obtain the $\epsilon$-DP guarantee for this framework.

**Theorem 3.6.** *The posterior distribution* $\Pr(p_{\text{FRT}} \mid \tilde{\boldsymbol{n}})$ *released by the Bayesian denoising framework satisfies $\epsilon$-DP.*

The whole procedure is briefly summarized in Algorithm 1. This approach is mechanism-aware because it explicitly models the distribution of the DP noise. Rather than treating the noisy counts as true data, it propagates the uncertainty induced by the mechanism. The resulting posterior distribution for $p_{\text{FRT}}$ enables the reporting of point estimates together with credible sets for uncertainty quantification, which also supports more principled decision-making as discussed in Section 4.

---

**Algorithm 1** DP-FRT: Mechanism-aware Bayesian Denoising

---

**Input:** group sizes $(n_1, n_0)$; success counts $(n_{11}, n_{01})$; privacy budget $\epsilon$; prior $\pi$.

**Output:** posterior distribution of $p_{\text{FRT}}$.

**Step 1 (Privatize counts):** Apply the geometric mechanism (1) to privatize the success counts as $(\tilde{n}_{11}, \tilde{n}_{01})$;

**Step 2 (Obtain posterior):** Combine the prior $\pi$ with the noise kernels of mechanism centered at $(\tilde{n}_{11}, \tilde{n}_{01})$ to obtain the posterior distribution $\gamma(a, b)$ as in (2);

**Step 3 (Map to $p$-value):** For each $(a, b) \in \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$, compute the randomization $p$-value $p(a, b)$ using (3), and induce the posterior distribution of $p_{\text{FRT}}$ as given in (4).

---

To provide an illustrative example of the proposed DP-FRT procedure, we consider a dataset with treatment and control group sizes $n_1 = n_0 = 500$, where the observed numbers of successes are $n_{11} = 260$ and $n_{01} = 250$. Using a uniform prior on $(n_{11}, n_{01})$, we apply the approach under privacy budgets $\epsilon \in \{0.1, 0.5, 1.0\}$, and visualize the resulting posterior distributions of the Fisher's randomization $p$-value. Figure 1 shows the posterior probability mass functions for each privacy level, with the red dashed line marking the non-private $p$-value $p_{\text{FRT}} = 0.2846$ computed from the original data. When $\epsilon$ is small, the posterior is

highly diffuse and exhibits a mode at $p_{\text{FRT}} = 0.5$, reflecting greater uncertainty due to DP noise. As $\epsilon$ increases, the posterior concentrates more sharply around the non-private value, demonstrating the privacy-utility tradeoff captured by the DP-FRT framework.
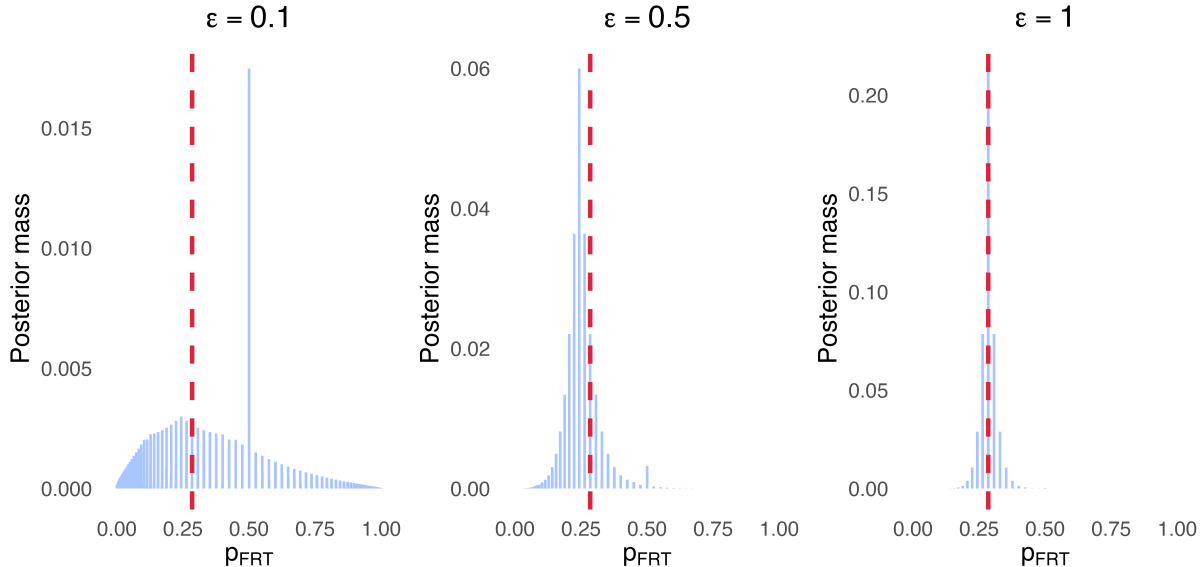


Figure 1: Posterior distributions of Fisher's randomization $p$-value under privacy budget $\epsilon \in \{0.1, 0.5, 1\}$. The red dashed line indicates the non-private $p$-value.

### 3.2.1 Prior Specification for Counts

As a key ingredient of the Bayesian denoising framework, we need to specify a data-independent prior $\pi$ on $(n_{11}, n_{01})$ to translate $(\tilde{n}_{11}, \tilde{n}_{01})$ into a posterior over the true counts via (2). When no external information is available, a natural default is the discrete uniform prior that assigns equal mass to every admissible pair:

$$\pi_{\text{unif}}(a, b) = \frac{1}{(n_1 + 1)(n_0 + 1)}, \qquad a \in \{0, \ldots, n_1\},\ b \in \{0, \ldots, n_0\}. \tag{5}$$

This choice is transparent, easy to compute with, and avoids privileging any specific configuration between the treatment and the control groups.

A closely related but slightly more structured option arises from the independent binomial formulation commonly used in clinical practice. Specifically, we independently posit $n_{11} \sim \text{Binom}(n_1, p_1)$ and $n_{01} \sim \text{Binom}(n_0, p_0)$, together with independent priors $p_1 \sim \text{Beta}(\alpha_1, \beta_1)$ and $p_0 \sim \text{Beta}(\alpha_0, \beta_0)$. Integrating out $(p_1, p_0)$ yields the following independent Beta-Binomial prior for the counts:

$$\pi_{\text{BB}}(a, b; \alpha_1, \beta_1, \alpha_0, \beta_0) = \frac{\binom{n_1}{a} B(a + \alpha_1, n_1 - a + \beta_1)}{B(\alpha_1, \beta_1)} \frac{\binom{n_0}{b} B(b + \alpha_0, n_0 - b + \beta_0)}{B(\alpha_0, \beta_0)}, \tag{6}$$

where $B(\cdot, \cdot)$ denotes the Beta function $B(x, y) = \Gamma(x)\Gamma(y)/\Gamma(x + y)$. Intuitively, this specification assumes that each group has its own baseline success rate and that the two groups are a priori independent. In particular, when taking $(\alpha_1, \beta_1) = (\alpha_0, \beta_0) = (1, 1)$,

13

it reduces to the discrete uniform prior (5). Thus, the uniform prior can be viewed as a special case of this independent Beta-Binomial construction.

An alternative that is sometimes useful encodes a common success rate across groups by setting $p_1 = p_0 = \theta$ with $\theta \sim \text{Beta}(\alpha, \beta)$. Marginalizing over $\theta$ gives a joint prior on counts that depends only on the total number of successes $K = a + b$:

$$\pi_{\text{CR}}(a, b; \alpha, \beta) = \binom{n_1}{a} \binom{n_0}{b} \frac{B\left(a + b + \alpha, n_1 + n_0 - (a + b) + \beta\right)}{B(\alpha, \beta)}. \tag{7}$$

When $(\alpha, \beta) = (1, 1)$, the induced distribution of the Fisher's randomization $p$-value is close to the uniform distribution on $(0, 1)$ up to discreteness.

When prior knowledge about the treatment effect is available from past studies or expert opinion, it can be encoded through the relationship between $(p_0, p_1)$. Two practical choices are the risk difference (RD) and the log odds ratio (log-OR):

$$\text{RD:} \quad p_0 \sim \text{Beta}(\alpha, \beta), \quad p_1 = \min\{\max\{p_0 + \tau, 0\}, 1\}, \quad \text{with } \tau \sim \mathcal{N}(\tau_0, \sigma_\tau^2),$$

$$\text{log-OR:} \quad p_0 \sim \text{Beta}(\alpha, \beta), \quad p_1 = \text{logit}^{-1}\left(\text{logit}(p_0) + \delta\right), \quad \text{with } \delta \sim \mathcal{N}(\delta_0, \sigma_\delta^2).$$

Both approaches induce flexible joint priors on $(n_{11}, n_{01})$ that reflect prior beliefs about the magnitude of the treatment effect. In practice, posterior inference typically requires MCMC sampling for these richer priors.

To conclude, we emphasize that FRT itself does not require any model assumptions. The priors above appear only in the denoising step to handle the randomness introduced by DP noise. Moreover, as Ding & Dasgupta (2016) caution, the independent binomial specification is a convenient practice rather than a consequence of randomization under the potential outcome framework. Any such specification should be regarded as a way of incorporating our prior beliefs instead of an inherent property of the design. Practically, as sample sizes grow, the likelihood dominates and the influence of $\pi$ on the posterior of $p_{\text{FRT}}$ becomes negligible.

### 3.2.2   Posterior Summaries with Uncertainty Quantification

A major advantage of the mechanism-aware Bayesian denoising framework is that it yields a full posterior distribution for $p_{\text{FRT}}$ with clear uncertainty quantification. Recall that given the posterior weights $\gamma(a, b)$ in (2) and the mapping $(a, b) \mapsto p(a, b)$ in (3), the posterior of $p_{\text{FRT}}$ is the discrete mixture in (4). Several summaries are natural and enjoy clear decision-theoretic interpretations.

For point estimation, let $w_p(u) = \sum_{(a,b):p(a,b)=u} \gamma(a, b)$ and $\mathcal{U} = \{p(a, b) : (a, b) \in \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}\}$. The posterior mean minimizes quadratic loss and is

$$\widehat{p}_{\text{mean}} = \sum_{u \in \mathcal{U}} u w_p(u) = \sum_{a=0}^{n_1} \sum_{b=0}^{n_0} \gamma(a, b) p(a, b). \tag{8}$$

The posterior median minimizes absolute loss and may not be unique because the posterior places mass at discrete points. Let $F(u) = \sum_{v \leq u} w_p(v)$ be the posterior cumulative distribution function on the support, then a valid posterior median satisfies

$$\widehat{p}_{\text{median}} \in \left\{m : F(m) \geq \frac{1}{2} \text{ and } 1 - F(m^-) \geq \frac{1}{2}\right\}. \tag{9}$$

The maximum a posteriori (MAP) estimate may also be non-unique. It is defined as the mode of the aggregated weights:

$$\widehat{p}_{\text{MAP}} \in \arg\max_{u \in \mathcal{U}} w_p(u), \tag{10}$$

which is informative when the posterior is multi-modal.

Credible sets can be constructed directly from (4) on the finite support. Let $\{u_{(1)} < \cdots < u_{(J)}\}$ be the sorted distinct support points with masses $w_p\big(u_{(j)}\big)$ and distribution function $F(u)$. The equal-tailed credible set at level $1 - \alpha$ uses posterior quantiles on this support:

$$C_{1-\alpha} = \{u_{(j)} : L \leq u_{(j)} \leq U\}, \tag{11}$$

where $L = \inf\{u : F(u) \geq \alpha/2\}$ and $U = \inf\{u : F(u) \geq 1 - \alpha/2\}$. Because $F$ is a step function, the posterior content of $C_{1-\alpha}$ is at least $1 - \alpha$ but may not be exact. When multiple $(L, U)$ achieve the same nominal level due to ties at the boundaries, it is reasonable to break ties by minimizing either the range $U - L$ or the cardinality of the included support points, and the chosen rule should be stated.

A highest-posterior-density set at level $1 - \alpha$ thresholds the weights. Let $t_\alpha$ be the smallest number such that $\sum_{u:w_p(u) \geq t_\alpha} w_p(u) \geq 1 - \alpha$. This yields

$$\text{HPD}_{1-\alpha} = \{u \in \mathcal{U} : w_p(u) \geq t_\alpha\}, \tag{12}$$

which is computed by sorting the support points in decreasing $w_p$ and accumulating mass until the target level is attained. If several points tie at the threshold, one may include all tied points to be conservative or include a minimal subset to match the nominal content, and the rule should be declared. Since many $(a, b)$ can map to the same $u$, the posterior on $p$ can be multi-modal and the HPD set need not be an interval. For interpretability, it is helpful to present both the exact finite set and its smallest enclosing interval in $[0, 1]$.

In addition to providing explicit uncertainty quantification, the Bayesian framework eliminates the need for post-processing to enforce feasible ranges on privatized counts. Specifically, the privatized counts $\tilde{n}_{11}$ and $\tilde{n}_{01}$ may fall outside the intervals $[0, n_1]$ or $[0, n_0]$, respectively, and the conventional approach is to clip these values to their feasible range. However, the following lemma shows that truncation is unnecessary for posterior inference, as the posterior update remains unchanged when the released counts are clipped under the geometric mechanism (1).

**Lemma 3.7** (Clipping invariance). *Let $n \in \{0, \ldots, M\}$ have prior $\{\pi_k\}_{k=0}^{M}$ and observe $\tilde{n} = n + \eta$ with $\eta$ drawn from a kernel $K(n, \tilde{n}) = \kappa_\rho(|n - \tilde{n}|)$ that is multiplicatively separable: $\kappa_\rho(a + b) = \kappa_\rho(a)\kappa_\rho(b)/\kappa_\rho(0)$. Define $\tilde{n}^{\text{clip}} = \min(\max(\tilde{n}, 0), M)$. Then, for all $k$,*

$$\Pr(n = k \mid \tilde{n}) = \Pr(n = k \mid \tilde{n}^{\text{clip}}).$$

*Proof.* By Bayes' rule, we have

$$\Pr(n = k \mid \tilde{n} = x) = \frac{\pi_k \kappa_\rho(|k - x|)}{\sum_{m=0}^{M} \pi_m \kappa_\rho(|m - x|)}.$$

If $x \in [0, M]$ there is nothing to prove. If $x < 0$, then for all $k \in \{0, \ldots, M\}$, $|k - x| = |k - 0| + |0 - x|$, hence

$$\kappa_\rho(|k - x|) = \frac{\kappa_\rho(|k - 0|)\kappa_\rho(|0 - x|)}{\kappa_\rho(0)}.$$

15

The factor $\kappa_\rho(|0 - x|)/\kappa_\rho(0)$ cancels between numerator and denominator, yielding

$$\Pr(n = k \mid \tilde{n} = x) = \frac{\pi_k \kappa_\rho(|k - 0|)}{\sum_{m=0}^{M} \pi_m \kappa_\rho(|m - 0|)} = \Pr(n = k \mid \tilde{n}^{\text{clip}} = 0).$$

If $x > M$, then $|k - x| = |M - k| + |x - M|$ and the same argument gives

$$\Pr(n = k \mid \tilde{n} = x) = \frac{\pi_k \kappa_\rho(|M - k|)}{\sum_{m=0}^{M} \pi_m \kappa_\rho(|M - m|)} = \Pr(n = k \mid \tilde{n}^{\text{clip}} = M).$$

Thus $\Pr(n = k \mid \tilde{n}) = \Pr(n = k \mid \tilde{n}^{\text{clip}})$ in all cases. $\qquad \square$

The argument applies coordinate-wise for $(n_{11}, n_{01})$ when independent geometric noise are used, because the kernel is multiplicatively separable and the joint kernel factorizes across coordinates.

### 3.2.3 Monte Carlo Sampling and Aggregation

To compute the exact posteriors for $p_{\text{FRT}}$, enumerating all $(a, b) \in \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$ requires $O(n_1 n_0)$ operations to compute and normalize the weights $w(a, b)$. For each pair, evaluating the tail sum in (3) takes $O(\min\{n_1, n_0\})$ time, resulting in an overall complexity of $O(n_1 n_0 \times \min\{n_1, n_0\})$.

When both $n_1$ and $n_0$ are moderate, full enumeration is feasible. For large $(n_1, n_0)$, it might be preferable to sample from the posterior rather than exhaustively evaluating the entire grid. In particular, when the prior can be factorized as $\pi(a, b) = \pi_1(a)\pi_0(b)$, as in the case of the discrete uniform prior (5) and the independent Beta Binomial prior (6), the posterior distribution also factorizes as $\gamma(a, b) = \gamma_{11}(a)\gamma_{01}(b)$, where

$$\gamma_{11}(a) \ \propto \ \pi_1(a)\kappa_\rho(\tilde{n}_{11} - a), \qquad \gamma_{01}(b) \ \propto \ \pi_0(b)\kappa_\rho(\tilde{n}_{01} - b).$$

As a result, $a^{(r)} \sim \gamma_{11}$ and $b^{(r)} \sim \gamma_{01}$ can be sampled independently for $r = 1, \ldots, R$. The cost of normalizing the two distributions is $O(n_1 + n_0)$, sampling $R$ pairs is $O(R)$, and evaluating the $p$-value for each draw via (3) costs $O(\min\{n_1, n_0\})$ per sample, resulting in a total computational cost of $O(n_1 + n_0 + R\min\{n_1, n_0\})$.

Given Monte Carlo samples $\{(a^{(r)}, b^{(r)})\}_{r=1}^{R}$, posterior summaries are obtained by mapping each pair to $u^{(r)} = p(a^{(r)}, b^{(r)})$ and aggregating over $\{u^{(r)}\}$. For example, the posterior mean can be approximated by $\tilde{p}_{\text{mean}} = \sum_{r=1}^{R} u^{(r)}/R$. The posterior distribution function can be approximated by the empirical distribution of $\{u^{(r)}\}$, from which the posterior median, credible sets, and HPD sets can be extracted as described in Section 3.2.2. The MAP estimate can also be obtained by tabulating the frequencies of the distinct values in $\{u^{(r)}\}$ and selecting the mode.

Furthermore, these Monte Carlo draws also enable posterior predictive generation of synthetic data and effect estimation. Specifically, each draw $(a^{(r)}, b^{(r)})$ determines

$$n_{11}^{(r)} = a^{(r)}, \qquad n_{01}^{(r)} = b^{(r)}, \qquad n_{10}^{(r)} = n_1 - a^{(r)}, \qquad n_{00}^{(r)} = n_0 - b^{(r)}.$$

Based on these tables, the following estimands can be computed:

$$\tau^{(r)} = \frac{a^{(r)}}{n_1} - \frac{b^{(r)}}{n_0}, \qquad \text{RR}^{(r)} = \frac{a^{(r)}/n_1}{b^{(r)}/n_0}, \qquad \text{OR}^{(r)} = \frac{a^{(r)}(n_0 - b^{(r)})}{(n_1 - a^{(r)})b^{(r)}},$$

16

where $\tau^{(r)}$ is the risk difference, $\mathrm{RR}^{(r)}$ is the risk ratio with a small continuity adjustment when needed, and $\mathrm{OR}^{(r)}$ is the odds ratio with the standard Haldane-Anscombe correction applied if any cell is zero. The empirical distributions of $\{\tau^{(r)}\}$, $\{\mathrm{RR}^{(r)}\}$, and $\{\mathrm{OR}^{(r)}\}$ provide posterior point summaries and credible intervals for treatment effects consistent with the mechanism-aware denoising. Notice that all operations are deterministic functions of $(\tilde{n}_{11}, \tilde{n}_{01})$, releasing such synthetic tables is post-processing and does not expend additional privacy budget.

# 4 Decision Making under DP-FRT

After obtaining the privatized Fisher's randomization $p$-value, one common practice is to reject the null hypothesis when it falls below a user-specified level $\alpha$. However, under differential privacy, the released statistics include randomness introduced by the privacy mechanism. A principled decision rule should also account for this uncertainty rather than rely on a single noisy realization or summary.

In this section, we base our decisions under DP-FRT on $\Pr\left(p_{\mathrm{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}}\right)$, which represents the posterior probability that the Fisher $p$-value does not exceed $\alpha$. This quantity summarizes the strength of evidence for rejection, incorporates the full posterior shape rather than a single noisy point, and connects naturally to classical summaries such as posterior quantiles and one-sided credible bounds. We introduce two frameworks for implementing this decision-making process.

## 4.1 Bayes Risk-optimal Decision Framework

We first develop a Bayes rule that minimizes posterior risk of departing from the non-private decision. The aim is to recover as faithfully as possible the significance label that would be obtained without DP noise. Fix a threshold $\alpha \in (0, 1)$ and let the decision $\delta \in \{1, 0\}$ indicate reject and not reject. Consider the loss

$$L(\delta, p_{\mathrm{FRT}}) = \begin{cases} 0, & \delta = 1, \ p_{\mathrm{FRT}} \leq \alpha, \\ \lambda_0, & \delta = 1, \ p_{\mathrm{FRT}} > \alpha, \\ \lambda_1, & \delta = 0, \ p_{\mathrm{FRT}} \leq \alpha, \\ 0, & \delta = 0, \ p_{\mathrm{FRT}} > \alpha, \end{cases} \tag{13}$$

where $\lambda_0 > 0$ and $\lambda_1 > 0$ represent the losses due to discordance with the non-private significance decision defined by $\{p_{\mathrm{FRT}} \leq \alpha\}$. More specifically, $\lambda_0$ denotes the loss of being too aggressive, such that we reject in cases where the non-private test would not. Conversely, $\lambda_1$ corresponds to the loss of being too conservative, such that we fail to reject in cases where the non-private test would. Then the corresponding posterior risks given $\tilde{\boldsymbol{n}} = (\tilde{n}_{11}, \tilde{n}_{01})$ are

$$R(\delta = 1 \mid \tilde{\boldsymbol{n}}) = \lambda_0 \Pr(p_{\mathrm{FRT}} > \alpha \mid \tilde{\boldsymbol{n}}), \tag{14}$$

$$R(\delta = 0 \mid \tilde{\boldsymbol{n}}) = \lambda_1 \Pr(p_{\mathrm{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}}). \tag{15}$$

The Bayes risk-optimal decision rule is given by

$$\delta_{\text{Bayes}}(\tilde{\boldsymbol{n}}) = \arg\min_{\delta \in \{0,1\}} \mathbb{E}[L(\delta, p_{\text{FRT}}) \mid \tilde{\boldsymbol{n}}]$$

$$= \begin{cases} 1, & R(\delta = 1 \mid \tilde{\boldsymbol{n}}) < R(\delta = 0 \mid \tilde{\boldsymbol{n}}), \\ 0, & \text{otherwise}, \end{cases} \tag{16}$$

$$= \begin{cases} 1, & \Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > \dfrac{\lambda_0}{\lambda_0 + \lambda_1}, \\ 0, & \text{otherwise}. \end{cases}$$

This rule can also be viewed as a decision based on a posterior quantile, where the quantile level is determined by the trade-off between the losses of being overly aggressive and overly conservative relative to the non-private decision. In particular, when there is no justification for assigning different losses to the two types of discordance (i.e., $\lambda_0 = \lambda_1$), the decision rule reduces to one based on the posterior median.

Although the decision rule (16) admits an intuitive interpretation, it captures only the relative penalties for overly liberal versus overly conservative decisions, and does not quantify confidence in the decision under DP noise. Therefore, it is advisable to present the decision together with $\Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$, which conveys the strength of evidence supporting the decision and clarifies how DP noise has influenced it.

### 4.1.1 Decision Confidence and Abstention under Uncertainty

When the posterior quantity $\Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$ is close to 1 or to 0, there is strong evidence to reject or not, respectively. However, with a tight privacy budget or a small sample size, the posterior for $p_{\text{FRT}}$ can be diffuse and this probability may lie near the binary Bayes threshold $\lambda_0/(\lambda_0 + \lambda_1)$, thus the decision might be unreliable.

In the spirit of Chow's rule (Chow 1957, 1970), it is natural to equip the decision procedure with an explicit abstention option to deal with these situations. Specifically, let $\delta \in \{0, 1, u\}$ with $u$ representing abstention due to uncertainty, and let $\lambda_u > 0$ be the loss incurred for abstaining, such as the operational cost of deferring the decision. The Bayes-optimal decision rule becomes

$$\delta_{\text{Bayes}}^*(\tilde{\boldsymbol{n}}) = \begin{cases} 1, & R(\delta = 1 \mid \tilde{\boldsymbol{n}}) < \min\{R(\delta = 0 \mid \tilde{\boldsymbol{n}}),\ \lambda_u\}, \\ 0, & R(\delta = 0 \mid \tilde{\boldsymbol{n}}) < \min\{R(\delta = 1 \mid \tilde{\boldsymbol{n}}),\ \lambda_u\}, \\ u, & \text{otherwise}, \end{cases}$$

$$= \begin{cases} 1, & \Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > \max\left\{ \dfrac{\lambda_0}{\lambda_0 + \lambda_1},\ 1 - \dfrac{\lambda_u}{\lambda_0} \right\}, \\ 0, & \Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) < \min\left\{ \dfrac{\lambda_0}{\lambda_0 + \lambda_1},\ \dfrac{\lambda_u}{\lambda_1} \right\}, \\ u, & \text{otherwise}. \end{cases} \tag{17}$$

In particular, when $\lambda_u$ is large relative to $\lambda_0$ and $\lambda_1$, the abstention region degenerates, and the trinary rule (17) reduces to the binary rule (16). The following lemma gives a condition that is both sufficient and tight.

**Lemma 4.1.** *If* $\lambda_u \geq \dfrac{H}{2}$, *where* $H = \dfrac{2\lambda_0\lambda_1}{\lambda_0 + \lambda_1}$ *is the harmonic mean of* $\lambda_0$ *and* $\lambda_1$, *then the abstention option degenerates.*

Guidance for specifying $\lambda_u$ can be based on the relative tolerance for indecision versus making errors. From an alternative viewpoint, $\lambda_u$ determines the width of the abstention region around the threshold $\Pr(p_{\mathrm{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}}) = \lambda_0/(\lambda_0 + \lambda_1)$. Figure 2 illustrates this rule, where the abstention region is visualized as the shaded gray band with width $\max\{0, 1 - 2\lambda_u/H\}$. A smaller $\lambda_u$ enlarges the abstention region and thereby increases the likelihood of abstention, whereas a larger $\lambda_u$ narrows the region and ultimately reduces the rule to the binary decision when $\lambda_u \geq H/2$. Therefore, $\lambda_u$ can be chosen accordingly, particularly when its direct interpretation is unclear. For instance, setting $\lambda_u = 0.025H$ yields an abstention region of width 0.95, which can be interpreted as imposing a 95% decision confidence in the Bayesian sense.
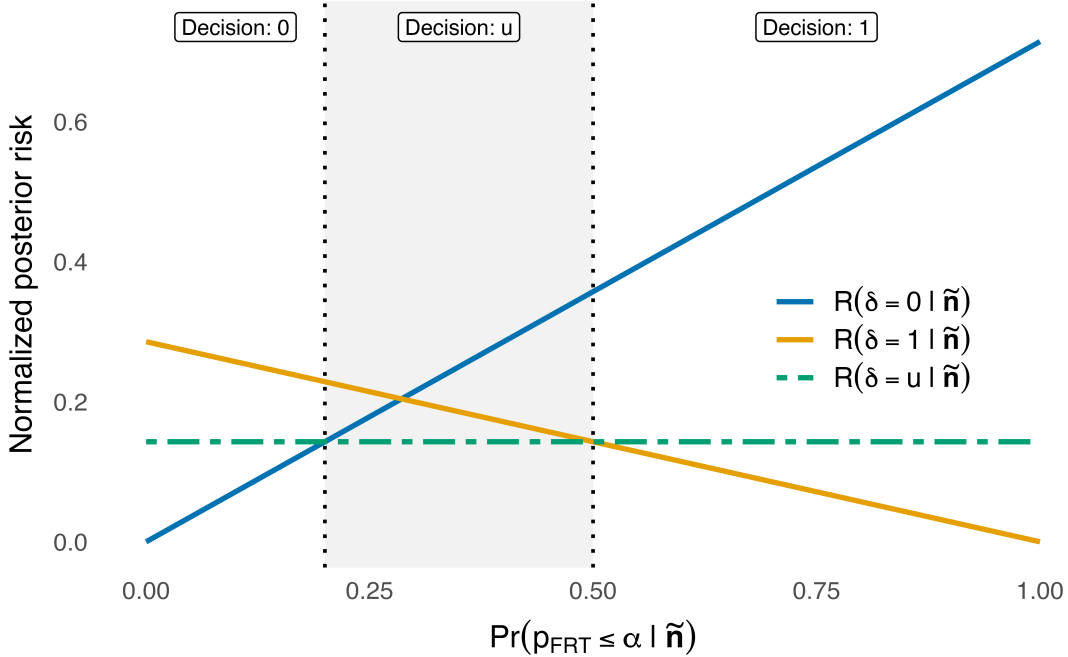


Figure 2: Illustration of the Bayes risk-optimal decision with an abstention option. The losses are specified as $\lambda_0 = 0.2$, $\lambda_1 = 0.5$, and $\lambda_u = 0.1$. The blue, orange, and green lines correspond to the posterior risks normalized by $(\lambda_0 + \lambda_1)$ for decisions $\delta = 0$, $\delta = 1$, and $\delta = u$, respectively. The shaded gray band marks the abstention region.

### 4.1.2 Sequential Decision under Additional Privacy Budget

We next consider how a Bayes risk-optimal decision rule (17) that initially abstains can be refined into a definite decision by allocating additional privacy budget.

Suppose the first privatized release $\tilde{\boldsymbol{n}} = (\tilde{n}_{11}, \tilde{n}_{01})$ is obtained with privacy budget $\epsilon > 0$, leading the trinary Bayes decision rule (17) to yield $\delta^*_{\mathrm{Bayes}} = u$. In such cases, the posterior evidence $\Pr(p_{\mathrm{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}})$ lies within the abstention interval determined by the loss parameters $(\lambda_0, \lambda_1, \lambda_u)$. To improve decision certainty without fully releasing the data,

one may adopt a sequential procedure that updates the posterior through an additional noisy observation drawn under a top-up privacy budget.

Specifically, let $\epsilon_{\text{plus}} > 0$ be the supplementary budget, we first generate an independent second noisy release

$$\tilde{\boldsymbol{n}}^+ = \left(\tilde{n}_{11}^+, \tilde{n}_{01}^+\right) = (n_{11} + \eta_{11}^+,\ n_{01} + \eta_{01}^+), \qquad \eta_{11}^+, \eta_{01}^+ \overset{\text{i.i.d.}}{\sim} \text{Geom}\left(\exp(-\epsilon_{\text{plus}})\right), \qquad (18)$$

so that the total privacy budget of the two releases is $\epsilon_{\text{tot}} = \epsilon + \epsilon_{\text{plus}}$ by the sequential composition property 2.8. Write $\rho = \exp(-\epsilon)$ and $\rho_+ = \exp(-\epsilon_{\text{plus}})$, and let $\kappa_\rho(h) = \dfrac{1-\rho}{1+\rho}\rho^{|h|}$ denote the geometric kernel. Then the sequential update of posterior follows directly from Bayes' rule as

$$\begin{aligned} \Pr(n_{11}, n_{01} \mid \tilde{\boldsymbol{n}}, \tilde{\boldsymbol{n}}^+) &\propto \Pr(\tilde{\boldsymbol{n}}^+ \mid n_{11}, n_{01}, \tilde{\boldsymbol{n}})\ \Pr(n_{11}, n_{01} \mid \tilde{\boldsymbol{n}}) \\ &\propto \Pr(\tilde{\boldsymbol{n}}^+ \mid n_{11}, n_{01})\ \Pr(\tilde{\boldsymbol{n}} \mid n_{11}, n_{01})\ \pi(n_{11}, n_{01}), \end{aligned} \qquad (19)$$

where $\pi(n_{11}, n_{01})$ denotes the prior defined in Section 3.2.1. Here, the posterior obtained from the first release serves as the prior for the second, allowing the sequential Bayesian update to combine both sources of information under the total privacy budget $\epsilon_{\text{tot}}$. Denote $\gamma^+(a, b) = \Pr(n_{11} = a, n_{01} = b \mid \tilde{\boldsymbol{n}}, \tilde{\boldsymbol{n}}^+)$. By substituting the geometric kernels, we have the following normalized posterior

$$\gamma^+(a, b) = \frac{\pi(a,b)\kappa_\rho(\tilde{n}_{11} - a)\kappa_\rho(\tilde{n}_{01} - b)\kappa_{\rho_+}(\tilde{n}_{11}^+ - a)\kappa_{\rho_+}(\tilde{n}_{01}^+ - b)}{\displaystyle\sum_{a'=0}^{n_1}\sum_{b'=0}^{n_0} \pi(a',b')\kappa_\rho(\tilde{n}_{11} - a')\kappa_\rho(\tilde{n}_{01} - b')\kappa_{\rho_+}(\tilde{n}_{11}^+ - a')\kappa_{\rho_+}(\tilde{n}_{01}^+ - b')} \qquad (20)$$

for $(a, b) \in \{0, \dots, n_1\} \times \{0, \dots, n_0\}$.

Next, we investigate how allocating additional privacy budget can improve decision certainty. Denote the abstention region as

$$A = (t_{\text{low}},\ t_{\text{high}}) = \left(\min\left\{\frac{\lambda_0}{\lambda_0 + \lambda_1},\ \frac{\lambda_u}{\lambda_1}\right\},\ \max\left\{\frac{\lambda_0}{\lambda_0 + \lambda_1},\ 1 - \frac{\lambda_u}{\lambda_0}\right\}\right),$$

which is determined solely by the loss parameters. By allocating extra privacy budget, the posterior evidence of rejection given $\alpha$ is updated from $\Pr(p_{\text{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}})$ to $\Pr(p_{\text{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}}, \tilde{\boldsymbol{n}}^+)$, which may become more concentrated around 0 or 1. Hence, although the width of the abstention region remains fixed, the probability that the updated evidence falls within this abstention region may decrease. The following theorem lower bounds this improvement on decision certainty.

**Theorem 4.2** (Lower Bound on Abstention-probability Reduction). *Denote* $\Psi = \Pr(p_{\text{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}})$, $\Psi^+ = \Pr(p_{\text{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}}, \tilde{\boldsymbol{n}}^+)$, *and let* $A = (t_{\text{low}}, t_{\text{high}})$ *be the abstention region with* $0 < t_{\text{low}} < t_{\text{high}} < 1$. *Assume* $\Pr(\Psi \in A) > 0$. *Then there exists a finite constant* $c > 0$ *such that, for all* $\epsilon_{\text{plus}} > 0$,

$$\Pr\left(\Psi^+ \notin A \mid \Psi \in A\right) \geq 1 - ce^{-\epsilon_{\text{plus}}}. \qquad (21)$$

*Proof.* The proof proceeds in three steps: (1) establishing posterior concentration on the true cell counts after the top-up release; (2) translating posterior concentration into escape from the abstention region; and (3) converting the unconditional escape bound into the desired conditional probability bound.

In the first step, we analyze the posterior concentration around the true counts. Define the event

$$\mathcal{E} = \{\eta_{11}^+ = 0, \ \eta_{01}^+ = 0\}.$$

Since $\eta_{11}^+, \eta_{01}^+ \overset{\text{i.i.d.}}{\sim} \mathrm{Geom}(\rho^+)$, we have

$$\Pr(\mathcal{E}^c) = 1 - \left(\frac{1 - \rho^+}{1 + \rho^+}\right)^2 = \frac{4\rho^+}{(1 + \rho^+)^2} \leq 4\rho^+.$$

On $\mathcal{E}$ we have $\tilde{n}_{11}^+ = n_{11}$ and $\tilde{n}_{01}^+ = n_{01}$, so

$$\kappa_{\rho^+}(\tilde{n}_{11}^+ - a)\kappa_{\rho^+}(\tilde{n}_{01}^+ - b) = \kappa_{\rho^+}(n_{11} - a)\kappa_{\rho^+}(n_{01} - b).$$

Thus, for $(a, b) = (n_{11}, n_{01})$ this factor equals $\kappa_{\rho^+}(0)\kappa_{\rho^+}(0)$, while for any $(a, b) \neq (n_{11}, n_{01})$ we have $|a - n_{11}| + |b - n_{01}| \geq 1$, so

$$\kappa_{\rho^+}(n_{11} - a)\kappa_{\rho^+}(n_{01} - b) \leq \kappa_{\rho^+}(0)^2 (\rho^+)^{|a-n_{11}|+|b-n_{01}|} \leq \kappa_{\rho^+}(0)^2 \rho^+. \tag{22}$$

Next we control the remaining factors uniformly in $\tilde{\boldsymbol{n}}$. Let $S$ denote the finite support of $(n_{11}, n_{01})$. Because the prior is strictly positive on $S$, there exist constants $0 < \pi_{\min} \leq \pi(a, b) \leq \pi_{\max} < \infty$ for all $(a, b) \in S$. For any integers $x$ and any $(a, b) \in S$,

$$\frac{\kappa_\rho(x - a)}{\kappa_\rho(x - n_{11})} = \rho^{|x-a|-|x-n_{11}|},$$

and by the triangle inequality, $\big||x - a| - |x - n_{11}|\big| \leq |a - n_{11}|$. Since $0 < \rho < 1$, this implies

$$\frac{\kappa_\rho(x - a)}{\kappa_\rho(x - n_{11})} \leq \rho^{-|a-n_{11}|}.$$

Applying this to both coordinates, we obtain, for any $(a, b) \in S$,

$$\frac{\kappa_\rho(\tilde{n}_{11} - a)\kappa_\rho(\tilde{n}_{01} - b)}{\kappa_\rho(\tilde{n}_{11} - n_{11})\kappa_\rho(\tilde{n}_{01} - n_{01})} \leq \rho^{-|a-n_{11}|-|b-n_{01}|},$$

and hence

$$\frac{\pi(a, b)\kappa_\rho(\tilde{n}_{11} - a)\kappa_\rho(\tilde{n}_{01} - b)}{\pi(n_{11}, n_{01})\kappa_\rho(\tilde{n}_{11} - n_{11})\kappa_\rho(\tilde{n}_{01} - n_{01})} \leq \frac{\pi_{\max}}{\pi_{\min}} \rho^{-|a-n_{11}|-|b-n_{01}|}.$$

Since $S$ is finite, the right-hand side admits a finite maximum over $(a, b) \neq (n_{11}, n_{01})$. Thus there exists a finite constant $C_1 > 0$, depending only on the prior and $\rho$, such that for all $\tilde{\boldsymbol{n}}$ and all $(a, b) \neq (n_{11}, n_{01})$,

$$\frac{\pi(a, b)\kappa_\rho(\tilde{n}_{11} - a)\kappa_\rho(\tilde{n}_{01} - b)}{\pi(n_{11}, n_{01})\kappa_\rho(\tilde{n}_{11} - n_{11})\kappa_\rho(\tilde{n}_{01} - n_{01})} \leq C_1. \tag{23}$$

21

Combining the bounds (22) and (23), we obtain that on $\mathcal{E}$,

$$\frac{\gamma^+(a,b)}{\gamma^+(n_{11}, n_{01})} \leq C_1 \rho^+ \qquad \text{for all } (a,b) \neq (n_{11}, n_{01}).$$

Summing over all $(a,b) \neq (n_{11}, n_{01})$ in the finite support $S$ gives

$$1 - \gamma^+(n_{11}, n_{01}) = \sum_{(a,b) \neq (n_{11}, n_{01})} \gamma^+(a,b) \leq C_2 \rho^+ \quad \text{on } \mathcal{E},$$

for some finite $C_2 > 0$ independent of $\tilde{\boldsymbol{n}}$ and $\rho^+$. Notice that on $\mathcal{E}^c$ we have the trivial bound $1 - \gamma^+(n_{11}, n_{01}) \leq 1$. Therefore,

$$1 - \gamma^+(n_{11}, n_{01}) \leq \mathbf{1}_{\mathcal{E}^c} + \mathbf{1}_{\mathcal{E}} C_2 \rho^+.$$

Taking expectations and using $\Pr(\mathcal{E}^c) \leq 4\rho^+$ yields

$$\mathbb{E}\big[1 - \gamma^+(n_{11}, n_{01})\big] \leq 4\rho^+ + C_2 \rho^+ \leq C\rho^+ = Ce^{-\epsilon_{\text{plus}}}, \tag{24}$$

for some finite constant $C > 0$ independent of $\epsilon_{\text{plus}}$.

In the second step, we translate posterior concentration into abstention probability. Let $p_{\text{FRT}} = g(n_{11}, n_{01})$ be the non-private Fisher's randomization $p$-value. Notice that the corresponding non-private decision $H = \mathbf{1}\{g(n_{11}, n_{01}) \leq \alpha\}$ is a deterministic function of $(n_{11}, n_{01})$. Then, by the definition of $\Psi^+$, we have

$$\Psi^+ = \sum_{a,b} \mathbf{1}\{g(a,b) \leq \alpha\} \gamma^+(a,b).$$

If $H = 1$, then $(n_{11}, n_{01})$ is in the rejection region and

$$\Psi^+ = \gamma^+(n_{11}, n_{01}) + \sum_{\substack{(a,b) \neq (n_{11}, n_{01}) \\ g(a,b) \leq \alpha}} \gamma^+(a,b) \geq \gamma^+(n_{11}, n_{01}),$$

hence

$$1 - \Psi^+ \leq 1 - \gamma^+(n_{11}, n_{01}).$$

If $H = 0$, then $(n_{11}, n_{01})$ is in the acceptance region and

$$\Psi^+ = \sum_{\substack{(a,b) \neq (n_{11}, n_{01}) \\ g(a,b) \leq \alpha}} \gamma^+(a,b) \leq 1 - \gamma^+(n_{11}, n_{01}).$$

Denote $d_A = \min\{t_{\text{low}}, \ 1 - t_{\text{high}}\} > 0$. If $1 - \gamma^+(n_{11}, n_{01}) \leq d_A$ and $H = 0$, then $\Psi^+ \leq d_A \leq t_{\text{low}}$, and thus $\Psi^+ \notin A$. If $1 - \gamma^+(n_{11}, n_{01}) \leq d_A$ and $H = 1$, then $1 - \Psi^+ \leq d_A \leq 1 - t_{\text{high}}$, so $\Psi^+ \geq t_{\text{high}}$ and hence $\Psi^+ \notin A$. Therefore, we have

$$\{\Psi^+ \in A\} \subseteq \{1 - \gamma^+(n_{11}, n_{01}) > d_A\}.$$

Taking probabilities and applying Markov's inequality together with (24) gives

$$
\begin{aligned}
\Pr(\Psi^+ \in A) &\leq \Pr\left(1 - \gamma^+(n_{11}, n_{01}) > d_A\right) \\
&\leq \frac{\mathbb{E}\left[1 - \gamma^+(n_{11}, n_{01})\right]}{d_A} \\
&\leq \frac{C}{d_A} e^{-\epsilon_{\mathrm{plus}}} =: C_A e^{-\epsilon_{\mathrm{plus}}}.
\end{aligned}
\tag{25}
$$

In the final step, we derive the conditional escape probability that is of interest. By Bayes' rule,

$$
\Pr(\Psi^+ \in A \mid \Psi \in A) = \frac{\Pr(\Psi^+ \in A, \Psi \in A)}{\Pr(\Psi \in A)} \leq \frac{\Pr(\Psi^+ \in A)}{\Pr(\Psi \in A)}.
$$

Since $\Psi$ depends only on the first release, $\Pr(\Psi \in A) =: p_A > 0$ is a constant independent of $\epsilon_{\mathrm{plus}}$. Combining this with (25), we obtain

$$
\Pr(\Psi^+ \in A \mid \Psi \in A) \leq \frac{C_A}{p_A} e^{-\epsilon_{\mathrm{plus}}}.
$$

Setting $c = C_A / p_A$ yields

$$
\Pr(\Psi^+ \notin A \mid \Psi \in A) = 1 - \Pr(\Psi^+ \in A \mid \Psi \in A) \geq 1 - c e^{-\epsilon_{\mathrm{plus}}},
$$

which completes the proof of Theorem 4.2. $\qquad\square$

As the additional privacy budget $\epsilon_{\mathrm{plus}}$ increases, the second-stage noise becomes negligible. In the limit $\epsilon_{\mathrm{plus}} \to \infty$, the second-stage release reveals the true cell counts with probability tending to one. Consequently, the updated posterior quantity $\Psi^+$ converges to the non-private decision $\mathbf{1}\{p_{\mathrm{FRT}} \leq \alpha\}$, and the chance of remaining in the abstention region decays exponentially in $\epsilon_{\mathrm{plus}}$.

We next turn to a complementary upper bound on how likely the refined posterior is to exit the abstention region after the top-up release. Our derivation adopts an information-theoretic viewpoint: conditioned on the first release, the top-up mechanism induces a channel whose ability to move the posterior is governed by its total variation contraction under differential privacy. By bounding this contraction and relating posterior movement to the distance from the abstention boundaries, we obtain a sharp upper bound on the reduction of abstention probability. We introduce some notation to formalize this argument. Let $S = \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$ denote the support of $(n_{11}, n_{01})$. For $h \in \{0, 1\}$, let $S_h = \{(a, b) \in S : H(a, b) = h\}$ be the subsets of $S$ corresponding to acceptance ($h = 0$) and rejection ($h = 1$). Define the $\ell_1$-distance between these two classes by $L_{\max} = \max_{(a,b) \in S_0, (a',b') \in S_1} \left(|a - a'| + |b - b'|\right)$, which can be computed without accessing any private data. We obtain the following theorem.

**Theorem 4.3** (Upper Bound on Abstention-probability Reduction). *Let $A = (t_{\mathrm{low}}, t_{\mathrm{high}})$ be the abstention region with $0 < t_{\mathrm{low}} < t_{\mathrm{high}} < 1$, and define $r(\Psi) = \min\{\Psi - t_{\mathrm{low}}, \ t_{\mathrm{high}} - \Psi\}$ whenever $\Psi \in A$. Then, for every $\epsilon_{\mathrm{plus}} > 0$,*

$$
\Pr\left(\Psi^+ \notin A \mid \Psi \in A\right) \leq 2 \min\{L_{\max} s(\epsilon_{\mathrm{plus}}), \ 1\} \, \mathbb{E}\left[\left. \frac{\Psi(1 - \Psi)}{r(\Psi)^2} \right| \Psi \in A\right],
\tag{26}
$$

*where $s(\varepsilon) = \tanh(\varepsilon/2) = (e^\varepsilon - 1)/(e^\varepsilon + 1)$.*

*Proof.* The proof proceeds in three steps: (1) a channel-based representation of posterior refinement; (2) bounding the channel's total variation contraction via DP; and (3) converting mean-square movement into an upper bound on the probability of exiting the abstention region. For simplicity, denote $H = \mathbf{1}\{p_{\mathrm{FRT}} \leq \alpha\}$. All expectations and probabilities in this proof are taken with respect to the joint law of $(H, \tilde{\boldsymbol{n}}, \tilde{\boldsymbol{n}}^+)$.

In the first step, we formulate the binary-input channel given the first release. Fix a realization $\tilde{\boldsymbol{n}}$ and condition on this event. Under this conditional law, the posterior of the non-private decision $H$ is $\Psi(\tilde{\boldsymbol{n}}) = \Pr(H = 1 \mid \tilde{\boldsymbol{n}}) \in [0, 1]$. For $h \in \{0, 1\}$, define the conditional output distributions as

$$Q_h(\cdot) := \mathcal{L}(\tilde{\boldsymbol{n}}^+ \mid H = h, \tilde{\boldsymbol{n}}).$$

Thus, conditionally on $\tilde{\boldsymbol{n}}$, $(H, \tilde{\boldsymbol{n}}^+)$ is a binary-input channel with input prior $\Pr(H = 1) = \Psi(\tilde{\boldsymbol{n}})$ and output kernel $\{Q_0, Q_1\}$.

Let $Y$ denote a generic random variable with distribution $\mathcal{L}(\tilde{\boldsymbol{n}}^+ \mid \tilde{\boldsymbol{n}})$. Then, by Bayes' rule, the refined posterior can be written as $\Psi^+(Y) = \Pr(H = 1 \mid \tilde{\boldsymbol{n}}, Y)$. Conditionally on $\tilde{\boldsymbol{n}}$, one can compute

$$\mathbb{E}\big[|\Psi^+ - \Psi| \mid \tilde{\boldsymbol{n}}\big] = 2\Psi(1 - \Psi)\mathrm{TV}(Q_1, Q_0),$$

where $\mathrm{TV}(\cdot, \cdot)$ denotes total variation distance. Since $|\Psi^+ - \Psi| \leq 1$, we also have

$$\mathbb{E}\big[(\Psi^+ - \Psi)^2 \mid \tilde{\boldsymbol{n}}\big] \leq \mathbb{E}\big[|\Psi^+ - \Psi| \mid \tilde{\boldsymbol{n}}\big] = 2\Psi(1 - \Psi)\mathrm{TV}(Q_1, Q_0).$$

Thus,
$$\mathbb{E}\big[(\Psi^+ - \Psi)^2 \mid \tilde{\boldsymbol{n}}\big] \leq 2\Psi(1 - \Psi)\mathrm{TV}(Q_1, Q_0). \tag{27}$$

In the second step, we bound $\mathrm{TV}(Q_1, Q_0)$ via $L_{\max}$ and $s(\epsilon_{\mathrm{plus}})$. Let $K(\cdot \mid a, b)$ denote the distribution of the top-up release $\tilde{\boldsymbol{n}}^+$ when the true counts are $(n_{11}, n_{01}) = (a, b)$, i.e., the geometric mechanism kernel. Conditionally on $\tilde{\boldsymbol{n}}$ and $H = h$, the posterior of the counts is supported on $S_h$ and $Q_h$ is the corresponding mixture:

$$Q_h(\cdot) = \sum_{(a,b) \in S_h} K(\cdot \mid a, b) \Pr\big(n_{11} = a, n_{01} = b \mid H = h, \tilde{\boldsymbol{n}}\big), \qquad h \in \{0, 1\}.$$

Consider two arbitrary mixtures $Q_1 = \sum_i \alpha_i P_i$, $Q_0 = \sum_j \beta_j R_j$ with weights summing to 1. By the triangle inequality, we have $\mathrm{TV}(Q_1, Q_0) \leq \sup_{i,j} \mathrm{TV}(P_i, R_j)$. Applied here, we obtain
$$\mathrm{TV}(Q_1, Q_0) \leq \sup_{(a,b) \in S_1, (a',b') \in S_0} \mathrm{TV}\big(K(\cdot \mid a, b), K(\cdot \mid a', b')\big). \tag{28}$$

Since the geometric mechanism on the counts has $\ell_1$-sensitivity one and satisfies $\epsilon_{\mathrm{plus}}$-DP with respect to the adjacency $d((a, b), (a', b')) = 1$. By Ghazi & Issa (2024), we know that for an $\epsilon$-DP mechanism $K$, any adjacent inputs $z, z'$ obey

$$\mathrm{TV}\big(K(\cdot \mid z), K(\cdot \mid z')\big) \leq s(\epsilon) := \tanh(\epsilon/2).$$

Equip $S$ with the $\ell_1$-metric $d((a, b), (a', b')) = |a - a'| + |b - b'|$. Then, for any two $(a, b), (a', b') \in S$, one can connect them by a path of $d((a, b), (a', b'))$ adjacent points in

this lattice. Applying the triangle inequality for total variation along this path and using the adjacent bound at each step gives

$$\mathrm{TV}\big(K(\cdot \mid a,b), K(\cdot \mid a',b')\big) \leq d\big((a,b),(a',b')\big)s(\epsilon_{\mathrm{plus}}).$$

Combining with (28) and the definition of $L_{\max}$, and noting that total variation is always at most one, we obtain

$$\mathrm{TV}(Q_1, Q_0) \leq \min\big\{L_{\max}s(\epsilon_{\mathrm{plus}}),\ 1\big\}. \tag{29}$$

Substituting (29) into (27), we obtain, for each fixed $\tilde{\boldsymbol{n}}$,

$$\mathbb{E}\big[(\Psi^+ - \Psi)^2 \mid \tilde{\boldsymbol{n}}\big] \leq 2\Psi(1-\Psi)\min\big\{L_{\max}s(\epsilon_{\mathrm{plus}}),\ 1\big\}. \tag{30}$$

In the final step, we establish how the mean-square movement translates into the abstention probability. Fix $\tilde{\boldsymbol{n}}$ such that $\Psi \in A$. The minimal distance from $\Psi$ to the boundaries of $A$ is $r(\Psi) = \min\{\Psi - t_{\mathrm{low}},\ t_{\mathrm{high}} - \Psi\} > 0$. If $\Psi \in A$ and $\Psi^+ \notin A$, then necessarily $|\Psi^+ - \Psi| \geq r(\Psi)$, so

$$\big\{\Psi \in A,\ \Psi^+ \notin A\big\} \subseteq \big\{\Psi \in A,\ |\Psi^+ - \Psi| \geq r(\Psi)\big\}.$$

Conditionally on this fixed $\tilde{\boldsymbol{n}}$ with $\Psi \in A$, we obtain

$$\Pr\big(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}\big) \leq \Pr\big(|\Psi^+ - \Psi| \geq r(\Psi) \mid \tilde{\boldsymbol{n}}\big).$$

Applying Markov's inequality to the nonnegative random variable $(\Psi^+ - \Psi)^2$ yields

$$\Pr\big(|\Psi^+ - \Psi| \geq r(\Psi) \mid \tilde{\boldsymbol{n}}\big) \leq \frac{\mathbb{E}\big[(\Psi^+ - \Psi)^2 \mid \tilde{\boldsymbol{n}}\big]}{r(\Psi)^2}.$$

Combining with (30) gives, for every $\tilde{\boldsymbol{n}}$ such that $\Psi \in A$,

$$\Pr\big(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}\big) \leq 2\min\{L_{\max}s(\epsilon_{\mathrm{plus}}),\ 1\}\ \frac{\Psi(1-\Psi)}{r(\Psi)^2}.$$

Finally,

$$\Pr\big(\Psi^+ \notin A \mid \Psi \in A\big) = \mathbb{E}\Big[\Pr\big(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}\big) \mid \Psi \in A\Big].$$

Using the bound from the previous step and the fact that $\Psi$ and $r(\Psi)$ are functions of $\tilde{\boldsymbol{n}}$, we obtain

$$\Pr\big(\Psi^+ \notin A \mid \Psi \in A\big) \leq 2\min\{L_{\max}s(\epsilon_{\mathrm{plus}}),\ 1\}\ \mathbb{E}\left[\frac{\Psi(1-\Psi)}{r(\Psi)^2}\ \middle|\ \Psi \in A\right],$$

which is exactly (26). This completes the proof of Theorem 4.3. $\qquad\square$

**Remark 4.4** (Interpretation of $s(\epsilon)$)**.** The quantity $s(\epsilon) = \tanh(\epsilon/2)$ is the standard conversion from a privacy budget $\epsilon$ to an upper bound on the adversary's distinguishing advantage between neighboring datasets, measured by the total variation distance between their output distributions under an $\epsilon$-DP mechanism. It thus provides an interpretable measure of how much information a DP mechanism can reveal about the underlying data (Ghazi & Issa 2024). In the context of Theorem 4.3, this quantity governs how far the posterior can move after the top-up release: the total variation distance between the top-up channels corresponding to acceptance and rejection grows at most linearly in $s(\epsilon_{\mathrm{plus}})$ up to a universal cap $L_{\max}$, so the refinement induced by additional privacy budget $\epsilon_{\mathrm{plus}}$ is necessarily limited by this contraction factor.

Based on Theorem 4.3, the following corollary quantifies a necessary lower bound on the additional privacy budget needed for obtaining a certain decision with high probability. It serves as a practical guideline for practitioners on privacy budget allocation.

**Corollary 4.5** (Budget necessary to exit the abstention region)**.** *Let* $A = (t_{\text{low}}, t_{\text{high}})$ *be the abstention region with* $0 < t_{\text{low}} < t_{\text{high}} < 1$, *and define* $r(\Psi) = \min\{\Psi - t_{\text{low}}, \ t_{\text{high}} - \Psi\}$ *whenever* $\Psi \in A$. *Fix a confidence level* $1 - \eta \in (0, 1)$. *If the refined posterior* $\Psi^+$ *satisfies*

$$\Pr(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}) \geq 1 - \eta,$$

*then it is necessary that*

$$\epsilon_{\text{plus}} \ \geq \ 2\operatorname{arctanh}\left(\frac{(1-\eta)r(\Psi)^2}{2L_{\max}\Psi(1-\Psi)}\right), \tag{31}$$

*where* $L_{\max}$ *is the* $\ell_1$-*distance between the classes as in Theorem 4.3.*

*Proof.* From the proof of Theorem 4.3, for any fixed $\tilde{\boldsymbol{n}}$ with $\Psi \in A$,

$$\Pr(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}) \leq 2\min\{L_{\max}s(\epsilon_{\text{plus}}), \ 1\} \frac{\Psi(1-\Psi)}{r(\Psi)^2}.$$

Assume that for some $\epsilon_{\text{plus}} > 0$ the desired bound $\Pr(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}) \geq 1 - \eta$ holds. Then

$$1 - \eta \leq \Pr(\Psi^+ \notin A \mid \tilde{\boldsymbol{n}}) \leq 2\min\{L_{\max}s(\epsilon_{\text{plus}}), \ 1\} \frac{\Psi(1-\Psi)}{r(\Psi)^2},$$

which implies

$$\min\{L_{\max}s(\epsilon_{\text{plus}}), \ 1\} \geq \frac{(1-\eta)r(\Psi)^2}{2\Psi(1-\Psi)}. \tag{32}$$

Since $\Psi \in A \subset (0, 1)$, we have

$$r(\Psi) = \min\{\Psi - t_{\text{low}}, \ t_{\text{high}} - \Psi\} \leq \min\{\Psi, \ 1 - \Psi\},$$

hence $r(\Psi)^2 \leq \min\{\Psi, \ 1 - \Psi\}^2 \leq \Psi(1 - \Psi)$. Combining with the fact that $L_{\max} \geq 1$, the right-hand side of (32) lies in $(0, 1)$. Therefore

$$s(\epsilon_{\text{plus}}) \geq \frac{(1-\eta)r(\Psi)^2}{2L_{\max}\Psi(1-\Psi)}.$$

Finally, since $s(\varepsilon) = \tanh(\varepsilon/2)$ and the hyperbolic arctangent function $\operatorname{arctanh}(\cdot)$ is the inverse of $\tanh(\cdot)$ on $[0, 1)$, it is equivalent to

$$\epsilon_{\text{plus}} \geq 2\operatorname{arctanh}\left(\frac{(1-\eta)r(\Psi)^2}{2L_{\max}\Psi(1-\Psi)}\right),$$

which completes the proof of Corollary 4.5. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

While Theorem 4.3 quantifies how much the additional release can at most reduce the probability of abstention, Corollary 4.5 inverts this relationship: it determines a necessary lower bound on $\epsilon_{\text{plus}}$ for ensuring, with high probability, that the refined posterior escapes the abstention region. The two results therefore provide dual perspectives on the same phenomenon: a larger privacy budget increases distinguishability, which strengthens posterior concentration and thus reduces abstention.

Intuitively, the necessary privacy budget depends on $L_{\max}$, $\eta$, and the initial posterior $\Psi$ because these quantities together capture (i) how different the acceptance and rejection hypotheses can be in the worst case ($L_{\max}$), (ii) how confident a decision the analyst requires ($\eta$), and (iii) how close the current posterior is to the abstention boundary ($\Psi$ through $r(\Psi)$). When the separation between hypotheses is smaller (smaller $L_{\max}$), when higher confidence is required (smaller $\eta$), or when the posterior lies deeper inside the abstention region so that the ratio $r(\Psi)^2/\big(\Psi(1-\Psi)\big)$ is larger, a greater privacy budget is needed to push the posterior outside the abstention region.

In practice, one may begin with a moderate initial privacy budget $\epsilon$ as a pilot release. If the decision rule (17) abstains, compute the bound above as a necessary lower bound when selecting $\epsilon_{\text{plus}}$ for a pre-specified confidence level (e.g., $\eta = 0.05$ for 95% certainty). A second release $\tilde{\boldsymbol{n}}^+$ is then generated, the posterior is updated, and the decision is recomputed. If the decision still abstains, the procedure may be repeated, or the decision may be reported as inconclusive once the total privacy budget has been fully used. This strategy spends privacy budget adaptively and only when necessary, allowing practitioners to decide at each stage whether to continue.

## 4.2  Frequentist-calibrated Decision Framework

We next calibrate a threshold to control frequentist type I error based on the posterior quantity $\Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$. Specifically, define

$$\delta_{\text{Freq}} = \begin{cases} 1, & \Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > t^*, \\ 0, & \text{otherwise}, \end{cases} \tag{33}$$

where the threshold $t^* \in [0,1]$ is chosen as the smallest value such that, under the sharp null $H_0^{\text{F}}$,

$$\Pr_{H_0^{\text{F}}}(\delta_{\text{Freq}} = 1) \le \alpha_{\text{Freq}}, \tag{34}$$

with $\alpha_{\text{Freq}} \in (0,1)$ denoting the target type I error level. Specifically, let $F_\Psi$ be the distribution of $\Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$ under $H_0^{\text{F}}$, then its $(1 - \alpha_{\text{Freq}})$-quantile

$$t^* = F_\Psi^{-1}(1 - \alpha_{\text{Freq}}) = \inf\{t \in [0,1] : F_\Psi(t) > 1 - \alpha_{\text{Freq}}\} \tag{35}$$

yields the desired calibrated threshold. Equivalently, rule (33) can also be viewed as an FRT based on the test statistic $\Psi = \Pr(p_{\text{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$, which rejects the null if its randomization $p$-value is at or below the target level $\alpha_{\text{Freq}}$.

To determine the critical value $t^*$ defined in (35), the key step is to derive $F_\Psi$, the distribution of $\Psi$ under the sharp null hypothesis. However, the construction of $F_\Psi$ depends on the total number of successes, $n_{+1} = n_{11} + n_{01}$, which is also subject to privatization and therefore unknown. Two calibration strategies are proposed below.

### 4.2.1 Worst-case Calibration

We first construct a threshold that is valid for all possible total number of successes. Recall that under the sharp null, the randomization redistributes all successes across groups. For each $K \in \{0, \ldots, n\}$, let $Q_K$ denote the probability mass function of $\tilde{\boldsymbol{n}} = (\tilde{n}_{11}, \tilde{n}_{01})$ under $H_0^{\mathrm{F}}$ with total number of successes equal to $K$, given by

$$Q_K(a, b) = \sum_{t=\max\{0, K-n_0\}}^{\min\{n_1, K\}} \frac{\binom{K}{t}\binom{n-K}{n_1-t}}{\binom{n}{n_1}} \kappa_\rho(a - t)\kappa_\rho(b - (K - t)) \tag{36}$$

for $(a, b) \in \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$. We then define $F_\Psi^{(K)}$ as the cumulative distribution function of $\Psi = \Pr(p_{\mathrm{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}})$ when $\tilde{\boldsymbol{n}} \sim Q_K$, that is,

$$F_\Psi^{(K)}(t) = \sum_{a=0}^{n_1} \sum_{b=0}^{n_0} Q_K(a, b)\mathbf{1}\{\Psi \le t\}, \qquad t \in [0, 1], \tag{37}$$

and denote $t_K = \inf\{t : F_\Psi^{(K)}(t) > 1 - \alpha_{\mathrm{Freq}}\}$ as the right-continuous $(1 - \alpha_{\mathrm{Freq}})$-quantile of $F_\Psi^{(K)}$. Finally, we set the least favorable threshold as $t_{\mathrm{LFC}}^* = \sup_{K \in \{0, \ldots, n\}} t_K$. This leads to the following decision rule as

$$\delta_{\mathrm{LFC}}(\tilde{\boldsymbol{n}}) = \begin{cases} 1, & \Pr(p_{\mathrm{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > t_{\mathrm{LFC}}^*, \\ 0, & \text{otherwise.} \end{cases} \tag{38}$$

The following theorem ensures type I error control under this rule.

**Theorem 4.6** (Type I error control with worst-case calibration)**.** *Under the sharp null* $H_0^{\mathrm{F}}$*, we have*

$$\Pr_{H_0^{\mathrm{F}}}\left(\delta_{\mathrm{LFC}}(\tilde{\boldsymbol{n}}) = 1\right) \le \alpha_{\mathrm{Freq}}, \tag{39}$$

*where the probability averages over randomization and the privacy mechanism.*

*Proof.* Fix $K$, then by construction, $\Pr\left(\Pr(p_{\mathrm{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > t_K \mid K\right) \le \alpha_{\mathrm{Freq}}$. Since $t_{\mathrm{LFC}}^* \ge t_K$, it follows that $\Pr\left(\Pr(p_{\mathrm{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > t_{\mathrm{LFC}}^* \mid K\right) \le \alpha_{\mathrm{Freq}}$. Since this holds for every $K$, including the case where $K$ equals the true total successes $n_{+1}$, the type I error is guaranteed to be no greater than $\alpha_{\mathrm{Freq}}$. $\qquad\square$

### 4.2.2 Data-adaptive Calibration with Confidence Sets

We can also restrict the least favorable search to a data-adaptive confidence set for $n_{+1}$ with high coverage. Specifically, for each $K$, define $Q_K$ as in (36), then order the lattice points $(a, b) \in \{0, \ldots, n_1\} \times \{0, \ldots, n_0\}$ by decreasing $Q_K(a, b)$ and take the smallest set $A_K$ whose total mass is at least $1 - \eta$, with $\eta \in (0, \alpha_{\mathrm{Freq}})$. Define the $(1 - \eta)$-confidence set

$$C_{1-\eta}(\tilde{\boldsymbol{n}}) = \{K \in \{0, \ldots, n\} : \tilde{\boldsymbol{n}} \in A_K\}. \tag{40}$$

Then for $\alpha' = \alpha_{\mathrm{Freq}} - \eta$, define $t_K' = \inf\{t : F_\Psi^{(K)}(t) > 1 - \alpha'\}$ as the right-continuous $(1 - \alpha')$ quantile of $F_\Psi^{(K)}$ defined in (37). By setting $t_{\mathrm{Neyman}}^*(\tilde{\boldsymbol{n}}) = \sup_{K \in C_{1-\eta}(\tilde{\boldsymbol{n}})} t_K'$, we obtain the decision rule

$$\delta_{\mathrm{Neyman}}(\tilde{\boldsymbol{n}}) = \begin{cases} 1, & \Pr(p_{\mathrm{FRT}} \le \alpha \mid \tilde{\boldsymbol{n}}) > t_{\mathrm{Neyman}}^*(\tilde{\boldsymbol{n}}), \\ 0, & \text{otherwise.} \end{cases} \tag{41}$$

Intuitively, this rule spends at most $\eta$ on potential miscoverage of $n_{+1}$ and uses the least favorable threshold within the confidence set, which yields type I error control while typically reducing conservativeness. In practice, one may choose $\eta = 0.05$ to make $C_{1-\eta}(\tilde{\boldsymbol{n}})$ a 95% confidence set. By construction, we obtain the following results.

**Lemma 4.7** (Neyman confidence set). *With $A_K$ and $C_{1-\eta}(\tilde{\boldsymbol{n}})$ defined as above, for each $K$, we have*

$$\Pr_K\{K \in C_{1-\eta}(\tilde{\boldsymbol{n}})\} = \Pr_K\{\tilde{\boldsymbol{n}} \in A_K\} \geq 1 - \eta, \tag{42}$$

*where $\Pr_K$ denotes probability under $Q_K$.*

**Theorem 4.8** (Type I error control with the data-adaptive confidence set). *Fix $\eta \in (0, \alpha_{\text{Freq}})$ and set $\alpha' = \alpha_{\text{Freq}} - \eta$. Then under the sharp null $H_0^{\text{F}}$, we have*

$$\Pr_{H_0^{\text{F}}}\left(\delta_{\text{Neyman}}(\tilde{\boldsymbol{n}}) = 1\right) \leq \alpha_{\text{Freq}}. \tag{43}$$

*Proof.* Fix $K$ and denote $\Psi = \Pr(p_{\text{FRT}} \leq \alpha \mid \tilde{\boldsymbol{n}})$. We first decompose

$$\Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}})\right) = \Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}}), \ \tilde{\boldsymbol{n}} \in A_K\right) + \Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}}), \ \tilde{\boldsymbol{n}} \notin A_K\right)$$

$$\leq \Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}}), \ \tilde{\boldsymbol{n}} \in A_K\right) + \Pr_K\left(\tilde{\boldsymbol{n}} \notin A_K\right).$$

On the event $\{\tilde{\boldsymbol{n}} \in A_K\}$ we have $K \in C_{1-\eta}(\tilde{\boldsymbol{n}})$, so by the definition of $t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}})$, $t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}}) \geq t'_K$. Hence, $\left\{\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}})\right\} \subseteq \{\Psi > t'_K\}$ on $\{\tilde{\boldsymbol{n}} \in A_K\}$. Therefore,

$$\Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}}), \ \tilde{\boldsymbol{n}} \in A_K\right) \leq \Pr_K\left(\Psi > t'_K\right) \leq \alpha'.$$

By Lemma 4.7, we have $\Pr_K\left(\tilde{\boldsymbol{n}} \notin A_K\right) \leq \eta$. Combining the two bounds gives

$$\Pr_K\left(\Psi > t^*_{\text{Neyman}}(\tilde{\boldsymbol{n}})\right) \leq \alpha' + \eta = \alpha_{\text{Freq}}.$$

Since this holds for every $K$, including the case where $K$ equals the true total successes $n_{+1}$, the type I error is bounded by $\alpha_{\text{Freq}}$. $\qquad\square$

# 5 Simulation Studies and Real Data Analysis

This section presents simulation studies that illustrate and evaluate the proposed methods. We then apply the methods to the ADAPTABLE trial data (Jones et al. 2021) to compare the effectiveness and safety of the 81 mg and 325 mg aspirin doses under privacy protection. The code for reproducing the results is available at `https://github.com/qy-sun/dp_frt`.

## 5.1 Simulation Studies

### 5.1.1 DP Studies: Assessing Differentially Private Estimation of Fisher's Randomization $p$-value

In this subsection, we evaluate the differentially private estimation of Fisher's randomization $p$-value conditional on the observed data. We consider different realized data summarized by $(n_{11}, n_{10}, n_{01}, n_{00})$, with the corresponding non-private $p$-values reported in Table

1. Under varying sample sizes and privacy budgets, we repeat the privatization procedure 1000 times to assess estimation performance.

Table 1: Simulation settings for Cases 1–4 in the DP studies. The counts $(n_{11}, n_{10}, n_{01}, n_{00})$ summarize the observed outcomes under $n \in \{50, 100, 500\}$, along with the corresponding non-private Fisher's randomization $p$-values.

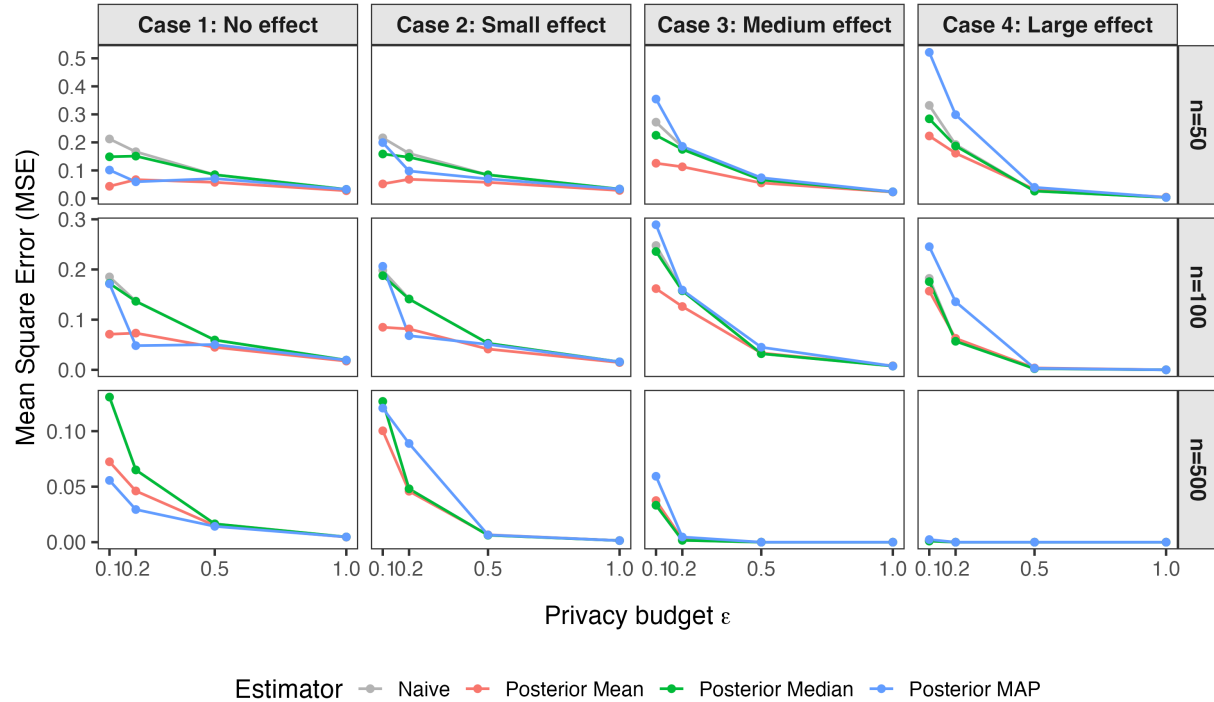| **Case** | $n$ | $n_{11}$ | $n_{10}$ | $n_{01}$ | $n_{00}$ | Non-private $p_{\mathrm{FRT}}$ |
|---|---|---|---|---|---|---|
| Case 1: No effect | 50 | 12 | 13 | 12 | 13 | 0.611 |
| | 100 | 25 | 25 | 25 | 25 | 0.579 |
| | 500 | 125 | 125 | 125 | 125 | 0.536 |
| Case 2: Small effect | 50 | 14 | 11 | 12 | 13 | 0.389 |
| | 100 | 28 | 22 | 25 | 25 | 0.344 |
| | 500 | 138 | 112 | 125 | 125 | 0.141 |
| Case 3: Medium effect | 50 | 16 | 9 | 12 | 13 | 0.197 |
| | 100 | 32 | 18 | 25 | 25 | 0.113 |
| | 500 | 162 | 88 | 125 | 125 | $5.54 \times 10^{-4}$ |
| Case 4: Large effect | 50 | 20 | 5 | 12 | 13 | $1.89 \times 10^{-2}$ |
| | 100 | 40 | 10 | 25 | 25 | $1.53 \times 10^{-3}$ |
| | 500 | 200 | 50 | 125 | 125 | $1.11 \times 10^{-12}$ |



Figure 3: MSE of the naive and Bayesian denoising estimators for Fisher's randomization $p$-value across privacy budgets $\epsilon \in \{0.1, 0.2, 0.5, 1.0\}$ in Cases 1–4 under sample sizes $n \in \{50, 100, 500\}$. The y-axis scales differ across panels.

Figure 3 presents the MSE of different estimators for the Fisher's randomization $p$-value across privacy budgets, simulation cases, and sample sizes. The posterior mean, median, and MAP estimators are obtained from the mechanism-aware Bayesian denoising framework, whereas the naive estimator directly perturbs the observed counts without accounting for DP noise. As expected, the MSE of all estimators decreases as either the privacy budget or sample size increases. The Bayesian estimators slightly improve estimation accuracy relative to the naive approach, particularly under small privacy budgets. Among the Bayesian methods, the posterior mean and posterior median achieve comparable MSE across most settings. The posterior MAP estimator occasionally exhibits larger MSE when $\epsilon$ is small, reflecting its sensitivity to multimodal posterior distributions induced by heavy DP noise. We recommend using the posterior mean estimator for its overall stability and accuracy.

Table 2: Coverage (%) and interval width of the 95% credible sets across privacy budgets $\epsilon \in \{0.1, 0.2, 0.5, 1.0\}$ in Cases 1–4 under sample sizes $n \in \{50, 100, 500\}$.

| Case | $\epsilon = 0.1$ | | $\epsilon = 0.2$ | | $\epsilon = 0.5$ | | $\epsilon = 1.0$ | |
|---|---|---|---|---|---|---|---|---|
| | Cov | Width | Cov | Width | Cov | Width | Cov | Width |
| **(a) Sample Size $n = 50$** | | | | | | | | |
| 1: No effect | 100 | 1.000 | 99.5 | 0.983 | 95.2 | 0.878 | 95.9 | 0.674 |
| 2: Small | 100 | 1.000 | 99.2 | 0.981 | 93.8 | 0.870 | 94.8 | 0.673 |
| 3: Medium | 100 | 1.000 | 99.1 | 0.986 | 95.9 | 0.825 | 96.3 | 0.564 |
| 4: Large | 100 | 1.000 | 97.9 | 0.980 | 97.1 | 0.574 | 95.5 | 0.211 |
| **(b) Sample Size $n = 100$** | | | | | | | | |
| 1: No effect | 99.8 | 0.990 | 96.3 | 0.940 | 93.5 | 0.782 | 94.4 | 0.537 |
| 2: Small | 99.6 | 0.991 | 95.7 | 0.935 | 96.3 | 0.787 | 96.7 | 0.516 |
| 3: Medium | 99.0 | 0.990 | 95.7 | 0.920 | 95.4 | 0.620 | 95.1 | 0.322 |
| 4: Large | 97.7 | 0.988 | 96.5 | 0.778 | 95.3 | 0.173 | 94.6 | 0.027 |
| **(c) Sample Size $n = 500$** | | | | | | | | |
| 1: No effect | 95.8 | 0.929 | 96.2 | 0.846 | 94.9 | 0.500 | 93.5 | 0.275 |
| 2: Small | 95.4 | 0.899 | 93.6 | 0.708 | 96.4 | 0.324 | 95.1 | 0.162 |
| 3: Medium | 95.3 | 0.585 | 94.2 | 0.145 | 94.9 | 0.009 | 94.4 | 0.002 |
| 4: Large | 95.6 | 0.062 | 94.9 | 0.000 | 96.0 | 0.000 | 97.4 | 0.000 |

More importantly, the Bayesian denoising framework provides valid uncertainty quantification for the privatized Fisher's randomization $p$-value. Table 2 reports the frequentist coverage and interval width of the 95% credible sets. The coverage of the proposed Bayesian credible sets remains close to the nominal 95% level when $\epsilon \geq 0.2$. Under very tight privacy constraints ($\epsilon = 0.1$), the credible sets are slightly conservative, with coverage approaching 100% and wider intervals. The results confirm that the credible sets produced by the mechanism-aware Bayesian framework maintain appropriate frequentist coverage.

### 5.1.2   Causal Studies: Evaluating Decision Rules under DP-FRT

In this subsection, we illustrate the decision rules based on the DP-FRT framework in finite-population randomized experiments. Specifically, we start with a finite population possessing both potential outcomes $\{Y_i(1), Y_i(0)\}_{i=1}^n$, then randomly assign $n_1$ units to the

treatment group, and finally observe the realized outcome table $(n_{11}, n_{10}, n_{01}, n_{00})$. Four representative scenarios (Cases 5–8) are reported in the Science table 3. We repeat the randomization and privatization procedures 100 times to evaluate decisions. The DP-FRT is implemented under a uniform prior with nominal significance level $\alpha = 0.05$.

Table 3: Simulation settings for Cases 5–8 in the causal studies. The counts $(N_{11}, N_{10}, N_{01}, N_{00})$ summarize the potential outcomes under $n \in \{50, 100, 500\}$, along with the corresponding treatment effects $\tau = (N_{10} - N_{01})/n$.

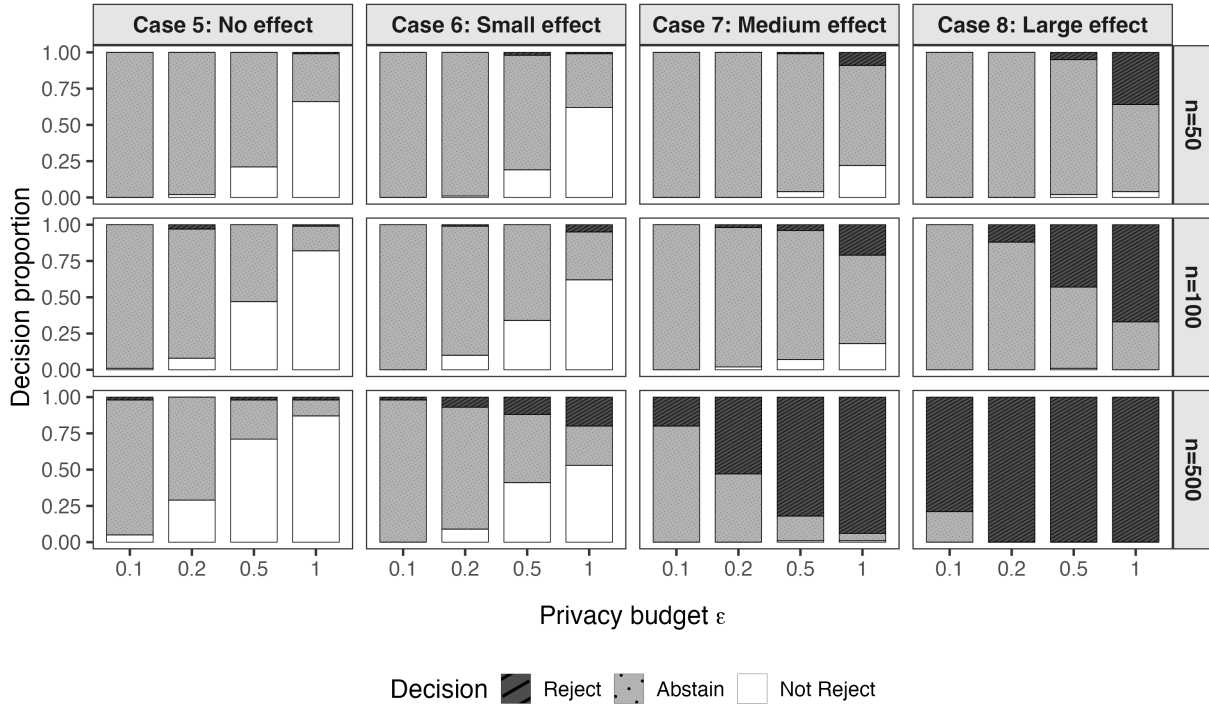| **Case** | $n$ | $N_{11}$ | $N_{10}$ | $N_{01}$ | $N_{00}$ | $\tau$ |
|---|---|---|---|---|---|---|
| | 50 | 25 | 0 | 0 | 25 | 0.00 |
| Case 5: No effect | 100 | 50 | 0 | 0 | 50 | 0.00 |
| | 500 | 250 | 0 | 0 | 250 | 0.00 |
| | 50 | 25 | 3 | 0 | 22 | 0.06 |
| Case 6: Small effect | 100 | 50 | 5 | 0 | 45 | 0.05 |
| | 500 | 250 | 25 | 0 | 225 | 0.05 |
| | 50 | 25 | 8 | 0 | 17 | 0.16 |
| Case 7: Medium effect | 100 | 50 | 15 | 0 | 35 | 0.15 |
| | 500 | 250 | 75 | 0 | 175 | 0.15 |
| | 50 | 25 | 15 | 0 | 10 | 0.30 |
| Case 8: Large effect | 100 | 50 | 30 | 0 | 20 | 0.30 |
| | 500 | 250 | 150 | 0 | 100 | 0.30 |



Figure 4: Decision proportions of the Bayes risk-optimal rule under budgets $\epsilon \in \{0.1, 0.2, 0.5, 1.0\}$ in Cases 5–8 for sample sizes $n \in \{50, 100, 500\}$ with loss parameters $(\lambda_0, \lambda_1, \lambda_u) = (1, 1, 0.025)$.

32

Figure 4 presents the decision proportions of the Bayes risk-optimal rule with abstention under different settings, with loss parameters $(\lambda_0, \lambda_1, \lambda_u) = (1, 1, 0.025)$. Each stacked bar shows the relative frequencies of the three possible decisions. Under tight privacy budgets, the proportion of abstentions is substantial, reflecting greater decision uncertainty induced by stronger privacy noise. As the privacy budget or sample size increases, abstentions gradually diminish and the proportion of correct rejections rises. For large effects (Case 8), nearly all decisions favor rejection when either $\epsilon$ or $n$ is sufficiently large, indicating that the Bayes risk-optimal rule remains sensitive to true treatment effects while effectively managing privacy-induced uncertainty through the abstention option.
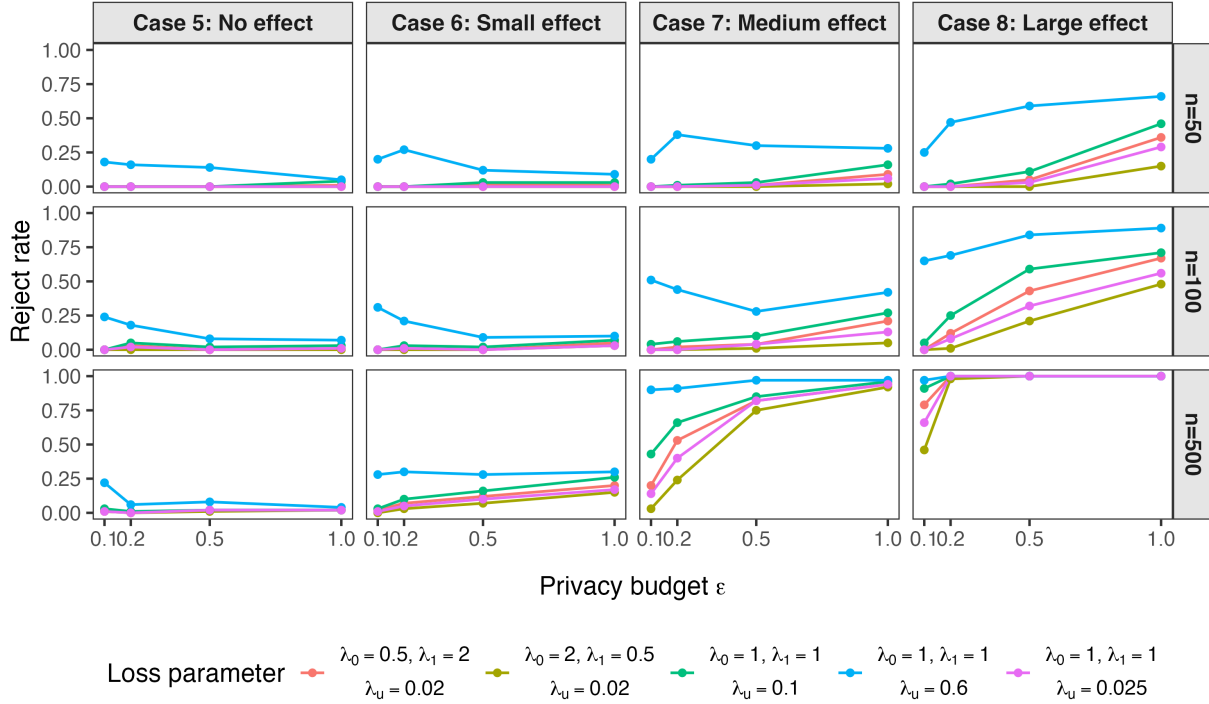


Figure 5: Rejection rates of the Bayes risk-optimal rule with abstention under varying loss parameters $(\lambda_0, \lambda_1, \lambda_u)$ across $\epsilon \in \{0.1, 0.2, 0.5, 1.0\}$ in Cases 5–8 for $n \in \{50, 100, 500\}$.

Figure 5 shows the rejection rates of the Bayes risk-optimal rule with abstention under varying loss parameters $(\lambda_0, \lambda_1, \lambda_u)$ across privacy budgets and sample sizes. The results highlight the sensitivity of decision outcomes to the relative weighting of rejection, non-rejection, and abstention losses. When $\lambda_u$ is small, the rule tends to abstain more and reject less frequently, leading to lower rejection rates across all settings. As $\lambda_u$ increases, abstention becomes less likely. Under the sharp null (Case 5), rejection rates remain near zero across all loss parameters. For cases with genuine effects (Cases 6–8), rejection rates increase with both $\epsilon$ and $n$ as expected. The results suggest that the abstention loss $\lambda_u$ provides a flexible mechanism for balancing privacy protection and inferential confidence.

## 5.2 Real Data Analysis

In this section, we illustrate the proposed DP-FRT framework using data from the ADAPT-ABLE trial (Jones et al. 2021). This is a pragmatic, randomized comparison of two aspirin

dosing strategies (81 mg vs. 325 mg daily) for secondary prevention in patients with established atherosclerotic cardiovascular disease.

Our primary binary endpoint is the composite of death from any cause, hospitalization for myocardial infarction, or hospitalization for stroke, evaluated under the intention-to-treat assignment of 81 mg ($n_0 = 7540$) versus 325 mg ($n_1 = 7536$) aspirin. Let $Y^{(1)}$ and $Y^{(2)}$ denote the primary composite endpoint and the safety endpoint, respectively, and let $n_{ij}^{(k)}$ denote the corresponding $2 \times 2$ cell counts for endpoint $k = 1, 2$. At median follow-up, 590 patients in the 81 mg group and 569 patients in the 325 mg group experienced the primary outcome, yielding the $2 \times 2$ Outcome table

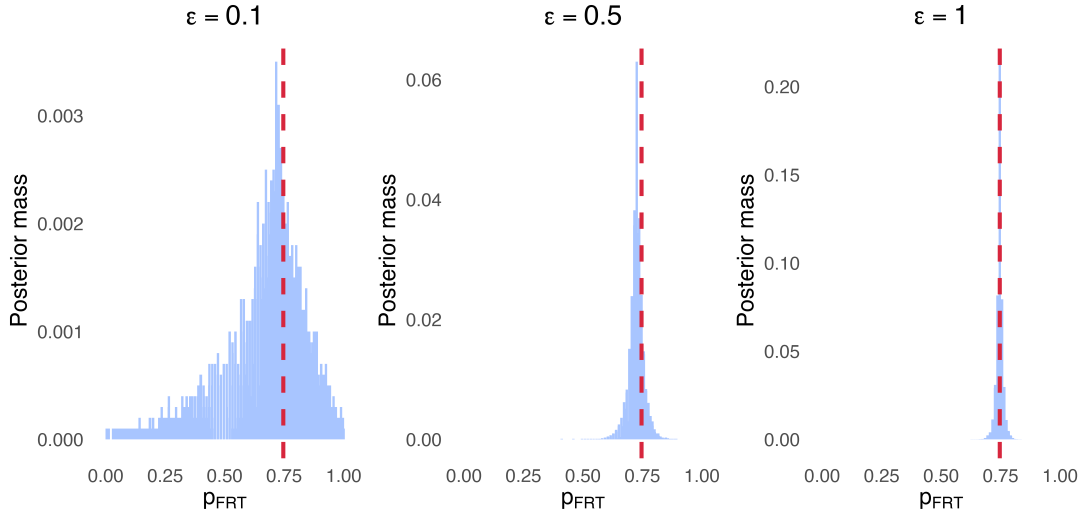|  | $Y^{(1)} = 1$ | $Y^{(1)} = 0$ | Row sum |
|---|---|---|---|
| Aspirin 325 mg | $n_{11}^{(1)} = 569$ | $n_{10}^{(1)} = 6967$ | 7536 |
| Aspirin 81 mg | $n_{01}^{(1)} = 590$ | $n_{00}^{(1)} = 6950$ | 7540 |
| Col sum | 1159 | 13,917 | |

We also consider the prespecified safety endpoint of hospitalization for major bleeding with transfusion. For this outcome, there were 53 events in the 81 mg group and 44 events in the 325 mg group, corresponding to

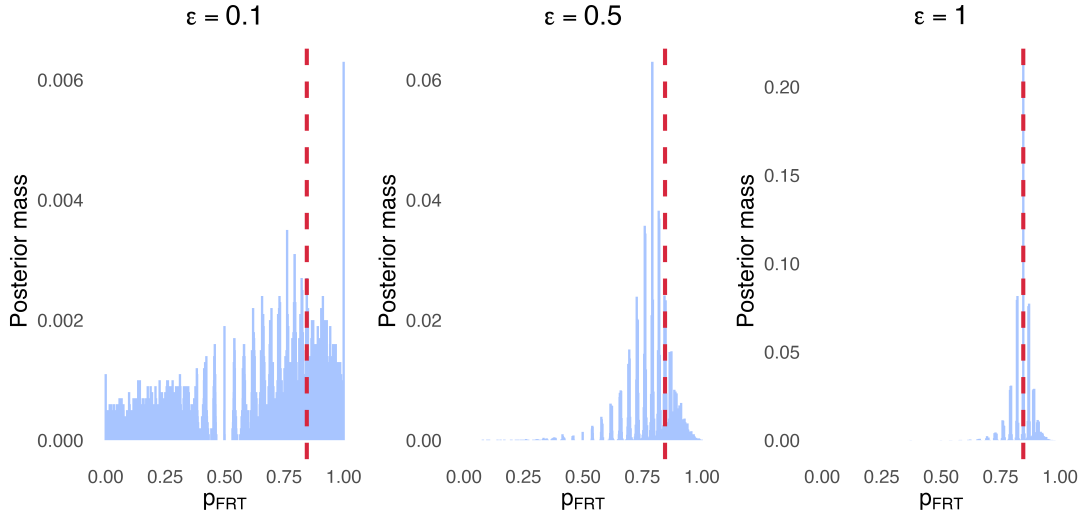|  | $Y^{(2)} = 1$ | $Y^{(2)} = 0$ | Row sum |
|---|---|---|---|
| Aspirin 325 mg | $n_{11}^{(2)} = 44$ | $n_{10}^{(2)} = 7492$ | 7536 |
| Aspirin 81 mg | $n_{01}^{(2)} = 53$ | $n_{00}^{(2)} = 7487$ | 7540 |
| Col sum | 97 | 14,979 | |

Both the primary composite endpoint and the safety endpoint are clinically sensitive outcomes, which may reveal information about severe cardiovascular and hemorrhagic events. In what follows, we treat the published $2 \times 2$ tables as the confidential input and use the proposed mechanism-aware Bayesian denoising framework to obtain posterior summaries of Fisher's randomization $p$-value and to derive the corresponding decision rules. Under the non-private FRT, the $p$-values for the primary composite endpoint and major bleeding were 0.7464 and 0.8452, respectively, providing no evidence against the sharp null hypothesis of no difference between the 325 mg and 81 mg aspirin groups. We subsequently applied the proposed DP-FRT procedure under privacy budgets $\epsilon \in \{0.1, 0.5, 1.0\}$, where $\alpha = 0.05$ and loss parameters $\lambda_0 = \lambda_1 = 1$ and $\lambda_u = 0.025$.

As shown in Figure 6 and Table 4, the posterior distributions of the differentially private Fisher's randomization $p$-values increasingly concentrate around their non-private counterparts as the privacy budget grows, with most of the posterior mass remaining well away from the rejection region. For the primary composite endpoint, the posterior rejection probabilities $\Psi$ are nearly zero across all privacy levels, yielding a Bayes decision of "not reject". For major bleeding, the procedure abstains at $\epsilon = 0.1$ due to greater posterior uncertainty ($\Psi = 0.1031$), but returns "not reject" decisions for $\epsilon \geq 0.5$.

Taken together, the DP-FRT analysis is consistent with the non-private findings from the original trial: 325 mg daily aspirin does not reduce the risk of death, myocardial infarction, or stroke compared with 81 mg, and there is no clear evidence of a difference in major bleeding risk between the two doses.

(a) Primary composite endpoint (death, myocardial infarction, or stroke).



(b) Major bleeding requiring transfusion.

Figure 6: Posterior distributions of Fisher's randomization $p$-value under privacy budget $\epsilon \in \{0.1, 0.5, 1.0\}$ for two clinical endpoints in the ADAPTABLE trial. The red dashed lines indicate the non-private $p$-values 0.7464 and 0.8452, respectively.

Table 4: Bayes risk-optimal decisions under the DP-FRT framework for the two ADAPTABLE endpoints, where $\alpha = 0.05$ and $(\lambda_0, \lambda_1, \lambda_u) = (1, 1, 0.025)$.

| Endpoint | Non-private $p_{\mathrm{FRT}}$ | Privacy budget | $\Psi$ | Decision |
|---|---|---|---|---|
| Primary composite | 0.7464 | 0.1 | 0.0012 | not reject |
| | | 0.5 | 0.0000 | not reject |
| | | 1.0 | 0.0000 | not reject |
| Major bleeding | 0.8452 | 0.1 | 0.1031 | abstain |
| | | 0.5 | 0.0000 | not reject |
| | | 1.0 | 0.0000 | not reject |

# 6 Discussion

This paper develops a framework for differentially private FRT with binary outcomes, enabling exact and interpretable causal inference under formal privacy guarantees. It introduces approaches for differentially private estimation of Fisher's randomization $p$-value and formal decision rules built upon them. These methods maintain finite-sample validity while also supporting uncertainty quantification and adaptive use of the privacy budget. While the proposed approaches focus on the most canonical setting, the framework naturally leads to several methodological extensions and theoretical directions.

**Toward bounded outcomes.** A direct extension is to randomized experiments with bounded discrete responses, such as ordinal or count data with finite support. Since the sensitivity of bounded discrete statistics remains finite, both the direct perturbation and mechanism-aware Bayesian frameworks can be adapted with minor modifications. For bounded continuous responses, direct application of the binary FRT requires discretization or binning. A principled approach is to treat the binning as a latent process and marginalize over the latent contingency tables rather than conditioning on a fixed discretization. This can be embedded within the Bayesian denoising framework, thereby propagating both binning and DP-induced uncertainties.

**Toward complex randomized designs.** The proposed methods can also be generalized beyond completely randomized experiments to stratified or block-randomized designs, paired or cluster randomization, and rerandomization procedures. In each case, the structure of the assignment mechanism changes the randomization distribution but does not affect the privacy mechanism applied to the observed outcomes. Combining DP-FRT with covariate-adjusted or rerandomized FRTs may improve efficiency, though the resulting randomization distributions will generally depend on more complex sufficient statistics requiring tailored privacy analysis.

**Toward testing weak null hypotheses.** Another promising direction is to extend the DP-FRT framework to the testing of weak null hypotheses (Ding & Dasgupta 2018, Wu & Ding 2021), which are often more relevant in practice than Fisher's sharp null. Under weak nulls, exact randomization inference is generally infeasible without additional modeling or asymptotic justification, because unit-level treatment effects are not fully specified. A natural bridge arises through the posterior predictive checking framework (Rubin 1984, Meng 1994), where one computes posterior predictive $p$-values (ppp) by averaging the randomization $p$-values over the posterior distribution of unobserved potential outcomes. As emphasized by Ding & Li (2018), the Bayesian interpretation of the Fisher randomization test is particularly valuable for weak nulls. In such settings, posterior averaging over the missing potential outcomes provides a coherent and interpretable way to account for uncertainty. Incorporating this averaging into the DP-FRT might therefore yield a unified framework for privacy-preserving inference under both sharp and weak null hypotheses.

**Toward superpopulation inference.** While our current formulation assumes a fixed finite population, many practical studies often target superpopulation inference, where the units are viewed as random draws from a larger population. Extending DP-FRT to

this context would require accounting for the randomness in the sampling process and integrating sampling design into the inference procedure.

**Toward design-based privacy efficiency.** Certain special cases naturally induce non-sensitive randomization distributions. For example, under matched-pair randomization, the paired sign test depends only on the number of discordant pairs and not on individual outcomes. Consequently, the randomization distribution itself requires no privacy noise, which can instead be applied only to the observed test statistic. Incorporating design-based structure can substantially reduce privacy cost and motivate further exploration of design-adaptive privacy mechanisms.

**Toward valid inference under data perturbation.** A broader theoretical question concerns whether the validity of hypothesis testing and decision-making can be preserved when data are stochastically perturbed, whether by privacy mechanisms, measurement error, or missingness. Our analysis suggests that explicitly modeling the noise process, as in the Bayesian denoising approach, is crucial for valid uncertainty quantification. Future research could seek to formalize this principle and characterize general conditions under which exact or asymptotically valid inference remains possible in the presence of stochastic perturbations.

# Conflict of Interest

On behalf of all authors, the corresponding author states that there is no conflict of interest.

# References

Abowd, J. M., Ashmead, R., Cumings-Menon, R. et al. (2022), 'The 2020 Census Disclosure Avoidance System TopDown Algorithm', *Harvard Data Science Review* (Special Issue 2).

Awan, J. & Slavković, A. (2018), Differentially private uniformly most powerful tests for binomial data, *in* 'Proceedings of the 32nd International Conference on Neural Information Processing Systems', Curran Associates Inc., pp. 4212–4222.

Chow, C. (1957), 'An optimum character recognition system using decision functions', *IRE Transactions on Electronic Computers* **EC–6**(4), 247–254.

Chow, C. (1970), 'On optimum recognition error and reject tradeoff', *IEEE Transactions on Information Theory* **16**(1), 41–46.

Couch, S., Kazan, Z., Shi, K. et al. (2019), Differentially private nonparametric hypothesis testing, *in* 'Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security', ACM, pp. 737–751.

Ding, P. & Dasgupta, T. (2016), 'A potential tale of two-by-two tables from completely randomized experiments', *Journal of the American Statistical Association* **111**(513), 157–168.

Ding, P. & Dasgupta, T. (2018), 'A randomization-based perspective on analysis of variance: a test statistic robust to treatment effect heterogeneity', *Biometrika* **105**(1), 45–56.

Ding, P. & Li, F. (2018), 'Causal inference: A missing data perspective', *Statistical Science* **33**(2), 214–237.

Dinur, I. & Nissim, K. (2003), Revealing information while preserving privacy, *in* 'Proceedings of the Twenty-second ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems', ACM, pp. 202–210.

D'Orazio, V., Honaker, J. & King, G. (2015), 'Differential privacy for social science inference', Sloan Foundation Economics Research Paper No. 2676160. Available at SSRN: `https://ssrn.com/abstract=2676160`.

Dwork, C. (2006), Differential privacy, *in* 'International Colloquium on Automata, Languages, and Programming', Vol. 4052, Springer, pp. 1–12.

Dwork, C., McSherry, F., Nissim, K. et al. (2006), Calibrating noise to sensitivity in private data analysis, *in* 'Theory of Cryptography Conference', Vol. 3876, Springer, pp. 265–284.

Dwork, C. & Roth, A. (2014), 'The algorithmic foundations of differential privacy', *Foundations and Trends in Theoretical Computer Science* **9**(3–4), 211–407.

Dwork, C., Smith, A., Steinke, T. et al. (2017), 'Exposed! a survey of attacks on private data', *Annual Review of Statistics and Its Application* **4**(1), 61–84.

Fisher, R. A. (1935), *The Design of Experiments*, 1st edn, Oliver and Boyd, Edinburgh.

Gaboardi, M., Lim, H., Rogers, R. et al. (2016), Differentially private chi-squared hypothesis testing: Goodness of fit and independence testing, *in* 'International Conference on Machine Learning', Vol. 48, PMLR, pp. 2111–2120.

Ghazi, E. & Issa, I. (2024), 'Total variation meets differential privacy', *IEEE Journal on Selected Areas in Information Theory* **5**, 207–220.

Ghosh, A., Roughgarden, T. & Sundararajan, M. (2012), 'Universally utility-maximizing privacy mechanisms', *SIAM Journal on Computing* **41**(6), 1673–1693.

Guha, S. & Reiter, J. P. (2025), 'Differentially private estimation of weighted average treatment effects for binary outcomes', *Computational Statistics & Data Analysis* **207**, 108145.

Jones, W. S., Mulder, H., Wruck, L. M. et al. (2021), 'Comparative effectiveness of aspirin dosing in cardiovascular disease', *New England Journal of Medicine* **384**(21), 1981–1990.

Karwa, V., Krivitsky, P. N. & Slavković, A. B. (2017), 'Sharing social network data: differentially private estimation of exponential family random-graph models', *Journal of the Royal Statistical Society Series C: Applied Statistics* **66**(3), 481–500.

Kazan, Z., Shi, K., Groce, A. et al. (2023), The test of tests: A framework for differentially private hypothesis testing, *in* 'International Conference on Machine Learning', Vol. 202, PMLR, pp. 16131–16151.

Kim, I. & Schrab, A. (2023), 'Differentially private permutation tests: Applications to kernel methods', *arXiv preprint arXiv:2310.19043* .

Lee, S. K., Gresele, L., Park, M. et al. (2019), 'Privacy-preserving causal inference via inverse probability weighting', *arXiv preprint arXiv:1905.12592* .

McSherry, F. D. (2009), Privacy integrated queries: An extensible platform for privacy-preserving data analysis, *in* 'Proceedings of the 2009 ACM SIGMOD International Conference on Management of data', ACM, pp. 19–30.

Meng, X. L. (1994), 'Posterior predictive $p$-values', *The Annals of Statistics* **22**(3), 1142–1160.

Mukherjee, S., Mustafi, A., Slavkovic, A. et al. (2024), Improving privacy for respondents in randomized controlled trials: A differential privacy approach, *in* 'Data Privacy Protection and the Conduct of Applied Research: Methods, Approaches and their Consequences', University of Chicago Press.

Niu, F., Nori, H., Quistorff, B. et al. (2022), Differentially private estimation of heterogeneous causal effects, *in* 'Conference on Causal Learning and Reasoning', Vol. 177, PMLR, pp. 618–633.

Nixon, M., Barrientos, A., Reiter, J. P. et al. (2022), 'A latent class modeling approach for differentially private synthetic data for contingency tables', *Journal of Privacy and Confidentiality* **12**(1).

Ohnishi, Y. & Awan, J. (2025), 'Locally private causal inference for randomized experiments', *Journal of Machine Learning Research* **26**(14), 1–40.

Peña, V. & Barrientos, A. F. (2025), 'Differentially private hypothesis testing with the subsampled and aggregated randomized response mechanism', *Statistica Sinica* **35**, 671–691.

Räisä, O., Jälkö, J., Kaski, S. et al. (2023), Noise-aware statistical inference with differentially private synthetic data, *in* 'International Conference on Artificial Intelligence and Statistics', Vol. 206, PMLR, pp. 3620–3643.

Rubin, D. B. (1974), 'Estimating causal effects of treatments in randomized and nonrandomized studies', *Journal of Educational Psychology* **66**(5), 688–701.

Rubin, D. B. (1984), 'Bayesianly justifiable and relevant frequency calculations for the applied statistician', *The Annals of Statistics* **12**(4), 1151–1172.

Seeman, J., Slavkovic, A. & Reimherr, M. (2020), Private posterior inference consistent with public information: A case study in small area estimation from synthetic census data, *in* 'International Conference on Privacy in Statistical Databases', Vol. 12276, Springer, pp. 323–336.

Wu, J. & Ding, P. (2021), 'Randomization tests for weak null hypotheses in randomized experiments', *Journal of the American Statistical Association* **116**(536), 1898–1913.