

---

# The Limits of Post-Selection Generalization

---

**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

1 While statistics and machine learning offers numerous methods for ensuring generalization,  
2 these methods often fail in the presence of *post selection*—the common  
3 practice in which the choice of analysis depends on previous interactions with the  
4 same dataset. A recent line of work has introduced powerful, general purpose  
5 algorithms that ensure a property called *post hoc generalization* (Cummings *et*  
6 *al.*, COLT’16), which says that no person when given the output of the algorithm  
7 should be able to find any statistic for which the data differs significantly from the  
8 population it came from.

9 In this work we show several limitations on the power of algorithms satisfying post  
10 *hoc generalization*. First, we show a tight lower bound on the error of any algorithm  
11 that satisfies post *hoc generalization* and answers adaptively chosen statistical  
12 queries, showing a strong barrier to progress in post selection data analysis. Second,  
13 we show that post *hoc generalization* is not closed under composition, despite many  
14 examples of such algorithms exhibiting strong composition properties.

15 

## 1 Introduction

16 Consider a dataset  $X$  consisting of  $n$  independent samples from some unknown population  $\mathcal{P}$ . How  
17 can we ensure that the conclusions drawn from  $X$  *generalize* to the population  $\mathcal{P}$ ? Despite decades  
18 of research in statistics and machine learning on methods for ensuring generalization, there is an  
19 increased recognition that many scientific findings do not generalize, with some even declaring this  
20 to be a “statistical crisis in science” [12]. While there are many reasons a conclusion might fail to  
21 generalize, one that is receiving increasing attention is *post-selection*, in which the choice of method  
22 for analyzing the dataset depends on previous interactions with the same dataset. Post-selection can  
23 arise from many common practices, such as *variable selection*, *exploratory data analysis*, and *dataset*  
24 *re-use*. Unfortunately, post-selection invalidates traditional methods for ensuring generalization,  
25 which assume that the method is independent of the data.

26 Numerous methods have been devised for statistical inference after post selection (e.g. [14, 16, 10,  
27 11, 20]). These are primarily *special purpose* procedures that apply to specific types of simple post  
28 selection that admit direct analysis. A more limited number of methods apply where the data is reused  
29 in one of a small number of prescribed ways (e.g. [2, 3]).

30 A recent line of work initiated by Dwork *et al.* [7] posed the question: Can we design *general-*  
31 *purpose* algorithms for ensuring generalization in the presence of post-selection? These works  
32 (e.g. [7, 6, 17, 1]) identified properties of an algorithm that ensure generalization under post-selection,  
33 including *differential privacy* [8], information-theoretic measures, and compression. They also  
34 identified many powerful general-purpose algorithms satisfying these properties, leading to algorithms  
35 for post-selection data analysis with greater statistical power than all previously known approaches.

36 Each of the aforementioned properties give incomparable generalization guarantees, and allow for  
37 qualitatively different types of algorithms. However, Cummings *et al.* [5] identified that the common  
38 thread in each of these approaches is to establish a notion of *post hoc generalization* (which they

39 originally called *robust generalization*), and initiated a general study of algorithms satisfying this  
40 notion. Informally, an algorithm  $\mathcal{M}$  satisfies post hoc generalization if there is no way, given only the  
41 output of  $\mathcal{M}(X)$ , to identify any *statistical query* [15] (that is, a bounded, linear, real-valued statistic  
42 on the population) such that the value of that query on the dataset is significantly different from its  
43 answer on the whole population.

**Definition 1.1** (Post Hoc Generalization [5]). An algorithm  $\mathcal{M} : \mathcal{X}^n \rightarrow \mathcal{Y}$  satisfies  $(\varepsilon, \delta)$ -post hoc  
generalization if for every distribution  $\mathcal{P}$  over  $\mathcal{X}$  and every algorithm  $\mathcal{A}$  that outputs a bounded  
function  $q : \mathcal{X} \rightarrow [-1, 1]$ , if  $X \sim \mathcal{P}^{\otimes n}$ ,  $y \sim \mathcal{M}(X)$ , and  $q \sim \mathcal{A}(y)$ , then

$$\mathbb{P}[|q(\mathcal{P}) - q(X)| > \varepsilon] \leq \delta,$$

44 where we use the notation  $q(\mathcal{P}) = \mathbb{E}[q(X)]$  and  $q(X) = \frac{1}{n} \sum_i q(X_i)$ , and the probability is over  
45 the sampling of  $X$  and any randomness of  $\mathcal{M}, \mathcal{A}$ .

46 Post hoc generalization is easily satisfied whenever  $n$  is large enough to ensure *uniform convergence*  
47 for the class of statistical queries. However, uniform convergence is only satisfied in the unrealistic  
48 regime where  $n$  is much larger than  $|\mathcal{X}|$ . Algorithms that satisfy post hoc generalization are interesting  
49 in the realistic regime where there will *exist* queries  $q$  for which  $q(\mathcal{P})$  and  $q(X)$  are far, but these  
50 queries cannot be *found*. The definition also extends seamlessly to richer types of statistics than  
51 statistical queries. However, restricting to statistical queries only strengthens our negative results.

52 Since all existing general-purpose algorithms for post-selection data analysis are analyzed via post  
53 hoc generalization, it is crucial to understand what we can achieve with algorithms satisfying post hoc  
54 generalization. In this work we present several strong limitaitons on the power of such algorithms.  
55 Our results identify natural barriers to progress in this area, and highlight important challenges for  
56 future research on post-selection data analysis.

## 57 1.1 Our Results

58 **Sample Complexity Bounds for Statistical Queries.** Our first contribution is strong new lower  
59 bounds on any algorithm that satisfies post hoc generalization and answers a sequence of adaptively  
60 chosen statistical queries—the setting introduced in Dwork *et al.* [7] and further studied in [1, 13, 18].  
61 In this model, there is an underlying distribution  $\mathcal{P}$ . We would like to design an algorithm  $\mathcal{M}$   
62 that holds a sample  $X \sim \mathcal{P}^{\otimes n}$ , takes statistical queries  $q$ , and returns accurate answers  $a$  such  
63 that  $a \approx q(\mathcal{P})$ . To model post-selection, we consider a *data analyst*  $\mathcal{A}$  that issues a sequence  
64 of queries  $q^1, \dots, q^k$  where each query  $q^j$  may depend on the answers  $a^1, \dots, a^{j-1}$  given by the  
65 algorithm in response to previous queries. The simplest algorithm  $\mathcal{M}$  would return the empirical  
66 mean  $q^j(X) = \frac{1}{n} \sum_i q^j(X_i)$  in response to each query, and one can show that this algorithm answers  
67 each query to within  $\pm \varepsilon$  if  $n \geq \tilde{O}(k/\varepsilon^2)$  samples.

68 Surprisingly, we can improve the sample complexity to  $n \geq \tilde{O}(\sqrt{k}/\varepsilon^2)$  by returning  $q(X)$  perturbed  
69 with carefully calibrated noise [7, 1]. The analysis of this approach uses post hoc generalization: the  
70 noise is chosen so that  $|a - q(X)| \leq \varepsilon/2$  and the noise ensures  $|q(\mathcal{P}) - q(X)| \leq \varepsilon/2$  for every query  
71 the analyst asks. Our main result shows that the sample complexity  $n = \tilde{O}(\sqrt{k}/\varepsilon^2)$  is essentially  
72 optimal for *any* algorithm that uses the framework of post hoc generalization. Our construction refine  
73 the techniques in [13, 18]—which yield a lower bound of  $n = \Omega(\sqrt{k})$  for  $\varepsilon = 1/3$ .

74 **Theorem 1.2** (Informal). *If  $\mathcal{M}$  takes a sample of size  $n$ , satisfies  $(\varepsilon, \delta)$ -post hoc generalization,  
75 and for every distribution  $\mathcal{P}$  over  $\mathcal{X} = \{\pm 1\}^{k+O(\log(n/\varepsilon))}$  and every data analyst  $\mathcal{A}$  who asks  $k$   
76 statistical queries,  $\mathbb{P}[\exists j \in [k], |q^j(\mathcal{P}) - a| > \varepsilon] \leq \delta$  then  $n = \Omega(\sqrt{k}/\varepsilon^2)$ , where the probability is  
77 taken over  $X \sim \mathcal{P}^{\otimes n}$  and the coins of  $\mathcal{M}$  and  $\mathcal{A}$ .*

78 Independently, Wang [21] proved a quantitatively similar bound to Theorem 1.2. However, their  
79 bound only applies to algorithms  $\mathcal{M}$  that receive only the empirical mean  $q(X)$  of each query. Their  
80 bound also applies for a slightly different (though closely related) class of statistics.

81 The dimensionality of  $\mathcal{X}$  required in our result is at least as large as  $k$ , which is somewhat necessary.  
82 Indeed, if the support of the distribution is  $\{\pm 1\}^d$ , then there is an algorithm  $\mathcal{M}$  that takes a sample  
83 of size just  $\tilde{O}(\sqrt{d} \log(k)/\varepsilon^3)$  [7, 1], so the conclusion is simply false if  $d \ll k$ . Even when  $d \ll k$ ,  
84 the aforementioned algorithms require running time at least  $2^d$  per query. [13, 18] also showed that  
85 any *polynomial time* algorithm that answers  $k$  queries to constant error requires  $n = \Omega(\sqrt{k})$ . We  
86 improve this result to have the optimal dependence on  $\varepsilon$ .

87 **Theorem 1.3** (Informal). Assume one-way functions exist and let  $c > 0$  be any constant. If  $\mathcal{M}$  takes  
 88 a sample of size  $n$ , has polynomial running time, satisfies  $(\varepsilon, \delta)$ -post hoc generalization, and for  
 89 every distribution  $\mathcal{P}$  over  $\mathcal{X} = \{\pm 1\}^{k^c + O(\log(n/\varepsilon))}$  and every data analyst  $\mathcal{A}$  who asks  $k$  statistical  
 90 queries,  $\mathbb{P}[\exists j \in [k], |q^j(\mathcal{P}) - a| > \varepsilon] \leq \delta$ , then  $n = \Omega(\sqrt{k}/\varepsilon^2)$ , where the probability is taken  
 91 over  $X \sim \mathcal{P}^{\otimes n}$  and the coins of  $\mathcal{M}$  and  $\mathcal{A}$ .

92 We prove the information-theoretic result (Theorem 1.2) in Section 2. Due to space restrictions, we  
 93 defer the computational result (Theorem 1.3) to the full version of this work.

94 **Negative Results for Composition.** One of the motivations for studying post hoc generalization is  
 95 to allow for exploratory data analysis and dataset re-use. In these settings, the same dataset may be  
 96 analyzed by many different algorithms, each satisfying post hoc generalization. Thus it is important  
 97 to understand whether the *composition* of these algorithms also satisfies post hoc generalization. We  
 98 show that this is not the case.

99 **Theorem 1.4.** For every  $n \in \mathbb{N}$  there is a collection of  $\ell = O(\log n)$  algorithms  $\mathcal{M}_1, \dots, \mathcal{M}_\ell$  that  
 100 take  $n$  samples from a distribution over  $\mathcal{X} = \{0, 1\}^{O(\log n)}$  such that (1) each of these algorithms are  
 101  $(\varepsilon, \delta)$ -post hoc generalizing for every  $\delta > 0$  and  $\varepsilon = O(\sqrt{\log(n/\delta)/n^{.999}})$ , but (2) the composition  
 102  $(\mathcal{M}_1, \dots, \mathcal{M}_\ell)$  is not  $(1.999, .999)$ -post hoc generalizing.

103 Theorem 1.4 states that there is a set of  $O(\log n)$  algorithms that have almost optimal post hoc  
 104 generalization, but whose composition does not have any non-trivial post hoc generalization.

105 If we consider a relaxed notion of *computational post hoc generalization*, then we show that compo-  
 106 sition can fail even for just two algorithms. Informally, computational post hoc generalization means  
 107 that Definition 1.1 is satisfied when the algorithm  $\mathcal{A}$  runs in polynomial time.

108 **Theorem 1.5.** Assume one-way functions exist. For every  $n \in \mathbb{N}$  there are two algorithms  $\mathcal{M}_1, \mathcal{M}_2$  that  
 109 take  $n$  samples from a distribution over  $\mathcal{X} = \{0, 1\}^{O(\log n)}$  such that (1) both algorithms are  
 110  $(\varepsilon, \delta)$ -computationally post hoc generalizing for every  $\delta > n^{-O(1)}$  and  $\varepsilon = O(\sqrt{\log(n/\delta)/n^{.999}})$ , but (2) the composition  
 111  $(\mathcal{M}_1, \mathcal{M}_2)$  is not  $(1.999, .999)$ -computationally post hoc generalizing.

112 We prove the information-theoretic result (Theorem 1.4) in Section 3. Due to space restrictions, we  
 113 defer the computational result (Theorem 1.5) to the full version of this work.

## 114 2 Lower Bounds for Statistical Queries

### 115 2.1 Post Hoc Generalization for Adaptive Statistical Queries

116 We are interested in the ability of *interactive* algorithms satisfying post hoc generalization to answer a  
 117 sequence of statistical queries. Definition 1.1 applies to such algorithms via the following experiment.

---

#### Algorithm 1: $\text{AQ}_{\mathcal{X}, n, k}[\mathcal{M} \rightrightarrows \mathcal{A}]$

---

$\mathcal{A}$  chooses a distribution  $\mathcal{P}$  over  $\mathcal{X}$   
 $X \sim \mathcal{P}^{\otimes n}$  and  $X$  is given to  $\mathcal{M}$  (but not to  $\mathcal{A}$ )  
**For**  $j = 1, \dots, k$   
 ┌  $\mathcal{A}$  outputs a statistical query  $q^j$  (possibly depending on  $q^1, a^1, \dots, q^{j-1}, a^{j-1}$ )  
 ┌  $\mathcal{M}(X)$  outputs  $a^j$

---

118  
 119 **Definition 2.1.** An algorithm  $\mathcal{M}$  is  $(\varepsilon, \delta)$ -post hoc generalizing for  $k$  adaptive queries over  $\mathcal{X}$  given  
 120  $n$  samples if for every adversary  $\mathcal{A}$ ,  $\mathbb{P}_{\text{AQ}_{\mathcal{X}, n, k}[\mathcal{M} \rightrightarrows \mathcal{A}]}[\exists j \in [k] \mid |q^j(X) - q^j(\mathcal{P})| > \varepsilon] \leq \delta$ .

### 121 2.2 A Lower Bound for Natural Algorithms

122 We begin with an information-theoretic lower bound for a class of algorithms  $\mathcal{M}$  that we call *natural*  
 123 *algorithms*. These are algorithms that can only evaluate the query on the sample points they are given.  
 124 That is, an algorithm  $\mathcal{M}$  is *natural* if, when given a sample  $X = (X_1, \dots, X_n)$  and a statistical query

125  $q : \mathcal{X} \rightarrow [-1, 1]$ , the algorithm  $\mathcal{M}$  returns an answer  $a$  that is a function only of  $(q(X_1), \dots, q(X_n))$ .  
126 In particular, it cannot evaluate  $q$  on data points of its choice. Many algorithms in the literature have  
127 this property. Formally, we define natural algorithms via the game  $\text{NAQ}_{\mathcal{X}, n, k}[\mathcal{M} \Leftarrow \mathcal{A}]$ . This game  
128 is identical to  $\text{AQ}_{\mathcal{X}, n, k}[\mathcal{M} \Leftarrow \mathcal{A}]$  except that when  $\mathcal{A}$  outputs  $q^j$ ,  $\mathcal{M}$  does not receive all of  $q^j$ , but  
129 instead receives only  $q_X^j = (q^j(X_1), \dots, q^j(X_n))$ .

**Theorem 2.2** (Lower Bound for Natural Algorithms). *There is an adversary  $\mathcal{A}_{\text{NAQ}}$  such that for every natural algorithm  $\mathcal{M}$ , and for universe size  $N = 8n/\varepsilon$ , if*

$$\mathbb{P}_{\text{NAQ}_{[N], n, k}[\mathcal{M} \Leftarrow \mathcal{A}_{\text{NAQ}}]} \left[ \exists j \in [k] \quad |q^j(X) - q^j(\mathcal{P})| > \varepsilon \bigvee |a^j - q^j(\mathcal{P})| > \varepsilon \right] \leq \frac{1}{100}$$

130 then  $n = \Omega(\sqrt{k}/\varepsilon^2)$

131 The proof uses the analyst  $\mathcal{A}_{\text{NAQ}}$  described in Algorithm 2. For notational convenience,  $\mathcal{A}_{\text{NAQ}}$   
132 actually asks  $k + 1$  queries, but this does not affect the final result.

---

**Algorithm 2:  $\mathcal{A}_{\text{NAQ}}$**

---

**Parameters:** sample size  $n$ , universe size  $N = \frac{8n}{\varepsilon}$ , number of queries  $k$ , target accuracy  $\varepsilon$

Let  $\mathcal{P} \leftarrow \text{U}_{[N]}$ ,  $A^1 \leftarrow [N]$ , and  $\tau \leftarrow 9\varepsilon\sqrt{2k \log(\frac{96}{\varepsilon})} + 1$

**For**  $j \in [k]$

    Sample  $p^j \sim \text{U}_{[0,1]}$

**For**  $i \in [N]$

        Sample  $\tilde{q}_i^j \sim \text{Ber}(p^j)$  and let  $q^j(i) \leftarrow \begin{cases} \tilde{q}_i^j & i \notin A^j \\ 0 & i \in A^j \end{cases}$

    Ask query  $q^j$  and receive answer  $a^j$

**For**  $i \in [N]$

        Let  $z_i^j \leftarrow \begin{cases} \text{trunc}_{3\varepsilon}(a^j - p^j) \cdot (q_i^j - p^j) & i \notin A^j \\ 0 & i \in A^j \end{cases}$

        where  $\text{trunc}_{3\varepsilon}(x)$  takes  $x \in \mathbb{R}$  and returns the nearest point in  $[-3\varepsilon, 3\varepsilon]$  to  $x$ .

        Let  $A^{j+1} \leftarrow \left\{ i \in [N] : \left| \sum_{\ell=1}^j z_i^\ell \right| > \tau - 1 \right\}$  (N.B. By construction,  $A^j \subseteq A^{j+1}$ .)

**For**  $i \in [N]$

    Define  $z_i \leftarrow \sum_{j=1}^k z_i^j$  and  $q_i^* \leftarrow \frac{z_i}{\tau}$

Let  $q^* : [N] \rightarrow [-1, 1]$  be defined by  $q^*(i) \leftarrow q_i^*$

---

133 In order to prove Theorem 2.2, it suffices to prove that either the answer  $a^j$  to one of the initial queries  
134  $q^j$  fails to be accurate (in which case  $\mathcal{M}$  is not accurate, or that the final query  $q^*$  gives significantly  
135 different answers on  $X$  and  $\mathcal{P}$  (in which case  $\mathcal{M}$  is not robustly generalizing). Formally, we have the  
136 following proposition.

137 **Proposition 2.3.** *For an appropriate choice of  $k = \Theta(\varepsilon^4 n^2)$  and  $n, \frac{1}{\varepsilon}$  sufficiently large, for any  
138 natural  $\mathcal{M}$ , with probability at least  $2/3$ , either (1)  $\exists j \in [k] \quad |a^j - q^j(\mathcal{P})| > \varepsilon$ , or (2)  $q^*(X) -$   
139  $q^*(\mathcal{P}) > \varepsilon$ . where the probability is taken over the game  $\text{NAQ}_{\mathcal{X}, n, k}[\mathcal{M} \Leftarrow \mathcal{A}_{\text{NAQ}}]$  and  $\mathcal{A}_{\text{NAQ}}$  is  
140 specified by Algorithm 2.*

141 We prove Proposition 2.3 using a series of claims. The first claim states that none of the values  $z_i$  are  
142 ever too large in absolute value, which follows immediately from the definition of the set  $A^j$  and the  
143 fact that each term  $z_i^j$  is bounded.

144 **Claim 2.4.** *For every  $i \in [N]$ ,  $|z_i| \leq \tau$ .*

145 The next claim states that, no matter how the mechanism answers, very few of the items not in the  
146 sample get “accused” of membership, that is, included in the set  $A^j$ .

147 **Claim 2.5** (Few Accusations).  $\Pr(|A_k \setminus X| \leq \varepsilon N/8) \geq 1 - e^{-\Omega(\varepsilon n)}$ .

148 *Proof.* Fix the biases  $p^1, \dots, p^k$  as well as the all the information visible to the mechanism (the query  
149 values  $\{q_i^j : i \in X, j \in [k]\}$ , as well as the answers  $a^1, \dots, a^k$ ). We prove that the probability of  $F$  is  
150 high conditioned on any setting of these variables.

151 The main observation is that, once we condition on the biases  $p^j$ , the query values at  $\{q_i^j : i \notin X, j \in$   
152  $[k]\}$  are independent with  $q_i^j \sim \text{Ber}(p^j)$ . This is true because  $\mathcal{M}$  is a natural algorithm (so it sees  
153 only the query values for points in  $X$ ) and, more subtly, because the analyst's decisions about how  
154 to sample the  $p^j$ 's, and which points in  $X$  to include in the sets  $A^j$ , are independent of the query  
155 values outside of  $X$ . By the principle of deferred decisions, we may thus think of the query values  
156  $\{q_i^j : i \notin X, j \in [k]\}$  as selected after the interaction with the mechanism is complete.

Fix  $i \notin X$ . For every  $j \in [k]$  and  $i \notin X$ , we have

$$\mathbb{E}[z_i^j] = \mathbb{E}[\text{trunc}_{3\varepsilon}(a^j - p^j) \cdot (q_i^j - p^j)] = \mathbb{E}[\text{trunc}_{3\varepsilon}(a^j - p^j)] \cdot \mathbb{E}[q_i^j - p^j] = 0.$$

157 By linearity of expectation, we also have  $\mathbb{E}[z_i] = \mathbb{E}\left[\sum_{j=1}^k z_i^j\right] = 0$ .

Next, note that  $|z_i^j| \leq 3\varepsilon$ , since  $\text{trunc}_{3\varepsilon}(a^j - p^j) \in [-3\varepsilon, 3\varepsilon]$  and  $q_i^j - p^j \in [0, 1]$ . The terms  $z_i^j$  are not independent, since if a partial sum  $\sum_{j=1}^{\ell} z_i^j$  ever exceeds  $\tau$ , then subsequent values  $z_i^j$  for  $j > \ell$  will be set to 0. However, we may consider a related sequence given by sums of the terms  $\tilde{z}_i^j = \text{trunc}_{3\varepsilon}(a^j - p^j) \cdot (\tilde{q}_i^j - p^j)$  (the difference from  $z_i^j$  is that we use values  $\tilde{q}_i^j \sim \text{Ber}(p^j)$  regardless of whether item  $i$  is in  $A^j$ ). Once we have conditioned on the biases and mechanism's outputs,  $\sum_{j=1}^k \tilde{z}_i^j$  is a sum of bounded independent random variables. By Hoeffding's Inequality, the sum is bounded  $O(\varepsilon\sqrt{k \log(1/\varepsilon)})$  with high probability, for every  $i \notin X$   $\mathbb{P}\left[\left|\sum_{j=1}^k \tilde{z}_i^j\right| > \varepsilon\sqrt{18k \ln\left(\frac{96}{\varepsilon}\right)}\right] \leq \frac{\varepsilon}{48}$ . By Etemadi's Inequality, a related bound holds uniformly over all the intermediate sums:

$$\forall i \notin X \quad \mathbb{P}\left[\exists \ell \in [k] : \left|\sum_{j=1}^{\ell} \tilde{z}_i^j\right| > \underbrace{3\varepsilon\sqrt{18k \ln\left(\frac{96}{\varepsilon}\right)}}_{\tau-1}\right] \leq 3 \cdot \mathbb{P}\left[\left|\sum_{j=1}^k \tilde{z}_i^j\right| > \varepsilon\sqrt{18k \ln\left(\frac{96}{\varepsilon}\right)}\right] \leq \frac{\varepsilon}{16}$$

158 Finally, notice that by construction, the real scores  $z_i^j$  are all set to 0 when an item is added to  $A^j$ , so  
159 the sets  $A^j$  are nested ( $A^j \subseteq A^{j+1}$ ), and a bound on partial sums of the  $\tilde{z}_i^j$  applies equally well to the  
160 partial sums of the  $z_i^j$ . Thus,  $\forall i \notin X \quad \mathbb{P}\left[\exists \ell \in [k] : \left|\sum_{j=1}^{\ell} z_i^j\right| > \tau - 1\right] \leq \frac{\varepsilon}{16}$

161 Now, the scores  $z^i$  are independent across players (again, because we have fixed the biases  $p^j$  and  
162 the mechanism's outputs). We can bound the probability that more than  $\frac{\varepsilon N}{4}$  elements  $i$  are "accused"  
163 over the course of the algorithm using Chernoff's bound:  $\mathbb{P}[|A^k \setminus X| > \frac{\varepsilon}{8}N] \leq e^{-\varepsilon N/64} \leq e^{-\Omega(n)}$   
164 The claim now follows by averaging over all of the choices we fixed.  $\square$

165 The next claim states that the sum of the scores over all  $i$  not in the sample is small.

166 **Claim 2.6.** *With probability at least  $\frac{99}{100}$ ,  $\sum_{i \in [N] \setminus X} z_i = O(\varepsilon\sqrt{Nk})$ .*

167 *Proof.* Fix a choice of  $(p^1, \dots, p^k) \in [0, 1]^k$ , the in-sample query values  $(q_X^1, \dots, q_X^k) \in \{0, 1\}^{n \times k}$ ,  
168 and the answers  $(a^1, \dots, a^k) \in [0, 1]^k$ . Conditioned on these, the values  $z_i$  for  $i \notin X$  are independent  
169 and identically distributed. They have expectation 0 (see the proof of Claim 2.5) and are bounded by  $\tau$   
170 (by Claim 2.4). By Hoeffding's inequality, with probability at least  $\frac{99}{100}$   $\sum_{i \in [N] \setminus X} z_i = O(\tau\sqrt{N}) =$   
171  $O(\varepsilon\sqrt{Nk})$  as desired. The claim now follows by averaging over all of the choices we fixed.  $\square$

172 **Claim 2.7.** *There exists  $c > 0$  such that, for all sufficiently small  $\varepsilon$  and sufficiently large  $n$ , with  
173 probability at least  $\frac{99}{100}$ , either  $\exists j \in [k] : |a^j - q^j(\mathcal{P})| > \varepsilon$  (large error), or  $\sum_{i \in [N]} z_i \geq ck$  (high  
174 scores in sample).*

175 The proof of Claim 2.7 relies on the following key lemma. The lemma has appeared in various  
176 forms [18, 9, 19]; the form we use is [4, Lemma 3.6] (rescaled from  $\{-1, +1\}$  to  $\{0, 1\}$ ).

**Lemma 2.8** (Fingerprinting Lemma). *Let  $f : \{0, 1\}^m \rightarrow [0, 1]$  be arbitrary. Sample  $p \sim \mathbb{U}_{[0,1]}$  and sample  $x_1, \dots, x_m \sim \text{Ber}(p)$  independently. Then*

$$\mathbb{E} \left[ (f(x) - p) \cdot \sum_{i \in [m]} (x_i - p) + \left| f(x) - \frac{1}{m} \sum_{i \in [m]} x_i \right| \right] \geq \frac{1}{12}.$$

*Proof of Claim 2.7.* To make use of the fingerprinting lemma, we consider a variant of Algorithm 2 that does not truncate the quantity  $a^j - p^j$  to the range  $\pm 2\epsilon$  when computing the score  $z_i^j$  for each element  $i$ . Specifically, we consider scores based on the quantities

$$\tilde{z}_i^j = \begin{cases} (a^j - p^j) \cdot (q_i^j - p^j) & \text{if } i \notin A^j, \\ 0 & \text{if } i \in A^j; \end{cases} \quad \text{and} \quad \tilde{z}_i = \sum_{j=1}^k \tilde{z}_i^j.$$

177 We prove two main statements: first, that these untruncated scores are equal to the truncated ones  
178 with high probability as long as the mechanism's answers are accurate. Second, that the expected  
179 sum of the untruncated scores is large. This gives us the desired final statement.

180 To relate the truncated and untruncated scores, consider the following three key events:

- 181 1. (“Few accusations”): Let  $F$  the event that, at every round  $j$ , set of “accused” items outside  
182 of the sample is small:  $|A_k \setminus X| \leq \epsilon N/8$ . Since the  $A^j$  are nested, event  $F$  implies the  
183 same condition for all  $j$  in  $[k]$ .
- 184 2. (“Low population error”): Let  $G$  be the event that at every round  $j \in [k]$ , the mechanism’s  
185 answer satisfies  $|a^j - p^j| \leq 3\epsilon$ .
- 186 3. (“Representative queries”): Let  $H$  be the event that  $|\tilde{q}^j(\mathcal{P}) - p^j| \leq \epsilon$  for all rounds  $j \in [k]$ —  
187 that is, each query’s population average is close to the corresponding sampling bias  $p^j$ .

188 **Sub-Claim 2.9.** *Conditioned on  $F \cap G \cap H$ , the truncated and untruncated scores are equal.  
189 Specifically,  $|a^j - p^j| \leq 3\epsilon$  for all  $j \in [k]$ .*

*Proof.* We can bound the difference  $|a^j - p^j|$  via the triangle inequality:

$$|a^j - p^j| \leq |a^j - q^j(\mathcal{P})| + |q^j(\mathcal{P}) - \tilde{q}^j(\mathcal{P})| + |\tilde{q}^j(\mathcal{P}) - p^j|.$$

190 The first term is the mechanism’s sample error (bounded when  $G$  occurs). The second is the distortion  
191 of the sample mean introduced by setting the query values of  $i \in A^j$  to 0. This distortion is at most  
192  $|A_j|/N$ . When  $F$  occurs,  $A^j$  has size at most  $|X| + |A^j \setminus X| \leq n + \epsilon N/8 = \epsilon N/4$ , so the second  
193 term is at most  $\epsilon/4$ . Finally, the last term is bounded by  $\epsilon$  when  $H$  occurs, by definition. The three  
194 terms add to at most  $3\epsilon$  when  $F$ ,  $G$ , and  $H$  all occur.  $\square$

195 We can bound the probability of  $H$  via a Chernoff bound: The probability of that a binomial random  
196 variable deviates from its mean by  $\epsilon N$  is at most  $2 \exp(-\epsilon^2 N/3)$ .

197 The technical core of the proof is the use of the fingerprinting lemma to analyze the difference  
198  $D$  between the sum of untruncated scores and the summed population errors:  $D := \sum_{i=1}^N \tilde{z}_i -$   
199  $\sum_{j=1}^k |a^j - q^j(\mathcal{P})| - k \mathbb{E} \left[ \frac{|A^j|}{N - |A^j|} \right]$

200 **Sub-Claim 2.10.**  $\mathbb{E}[D] = \Omega(k)$

201 *Proof.* We show that for each round  $j$ , the expected sum of scores for that round  $\sum_i \tilde{z}_i^j$  is at least  
202  $1/12 - \mathbb{E} \left[ |a^j - q^j(\mathcal{P})| - \frac{|A^j|}{N - |A^j|} \right]$ . This is true even when we condition on all the random choices  
203 and communication in rounds 1 through  $j - 1$ . Adding up these expectations over all rounds gives  
204 the desired expectation bound for  $D$ .

First, note that summing  $z_i^j$  over all elements  $i \in [N]$  is the same as summing over that round’s  
unaccused elements  $i \in [N] \setminus A^j$  (since  $\tilde{z}_i^j = 0$  for  $i \in A^j$ ). Thus,

$$\sum_{i=1}^N \tilde{z}_i^j = \sum_{i \in [N] \setminus A^j} \tilde{z}_i^j = (a^j - p^j) \sum_{i \in [N] \setminus A^j} (q_i^j - p^j).$$

We can now apply the Fingerprinting Lemma, with  $m = N - |A^j|$ ,  $p = p^j$ ,  $x_i = \tilde{q}_i^j$  for  $i \notin A^j$ , and  $f((x_i)_{i \notin A^j}) = a^j$  (note that  $f$  depends implicitly on  $A_j$ , but since we condition on the outcome of previous rounds, we may take  $A^j$  as fixed for round  $j$ ). We obtain

$$\mathbb{E} \left[ \sum_{i=1}^N \tilde{z}_i^j \right] \geq \frac{1}{12} - \mathbb{E} \left[ \left| a^j - \frac{1}{N - |A^j|} \cdot \sum_{i \notin A^j} q_i^j \right| \right]$$

205 Now the difference between  $\frac{1}{N - |A^j|} \sum_{i \notin A^j} q_i^j$  and the actual population mean  $\frac{1}{N} \sum_{i=1}^N q_i^j$  is at  
206 most  $N \cdot \left( \frac{1}{N} - \frac{1}{N - |A^j|} \right) = \frac{|A^j|}{N - |A^j|}$ . Thus we can upper-bound the term inside the right-hand side  
207 expectation above by  $|a^j - q^j(\mathcal{P})| + \frac{|A^j|}{N - |A^j|}$ .  $\square$

208 A direct corollary of Sub-Claim 2.10 is that there is a constant  $c' > 0$  such that, with probability at  
209 least  $199/200$ ,  $D \geq c'k$ . Let's call that event  $I$ .

210 Conditioned on  $F \cap G \cap H$ , we know that each  $\tilde{z}_i$  equals the real score  $z_i$  (by the first sub-claim  
211 above), that  $|a^j - q^j(\mathcal{P})| \leq 3\varepsilon$  for each  $j$ , and that  $|A^k| \leq \varepsilon N/8$ . If we also consider the intersection  
212 with  $I$ , then we have  $D \geq c'k - 3k\varepsilon - k \frac{\varepsilon/8}{1-\varepsilon/8} \geq k(c' - 4\varepsilon)$  (for sufficiently small  $\varepsilon$ ). By a union  
213 bound, the probability of  $\neg(F \cap H \cap I)$  is at most  $1/200 + \exp(-\Omega(\varepsilon^2 n)) \leq 1/100$  (for sufficiently  
214 large  $n$ ). Thus we get  $\mathbb{P} \left[ (\neg G) \text{ or } \left( \sum_{i=1}^N z_i \geq ck \right) \right] \geq \frac{99}{100}$ , where  $c = c' - 4\varepsilon$  is positive for  
215 sufficiently small  $\varepsilon$ . This completes the proof of Claim 2.7.  $\square$

216 To complete the proof of the proposition, suppose that  $|a^j - q^j(\mathcal{P})| \leq \varepsilon$  for every  $j$ , so that we can  
217 assume  $\sum_{i \in X} z_i = \Omega(k)$ . Then, we can show that, when  $n$  is sufficiently large and  $k \gtrsim \varepsilon^4 n^2$ , the  
218 final query  $q^*$  will violate robust generalization. A relatively straightforward calculation (omitted for  
219 space) shows that for the query  $q^*$  that we defined,  $q^*(X) - q^*(\mathcal{P}) = \Theta(\varepsilon \sqrt{k})$ . Now, we choose an  
220 appropriate  $k = \Theta(\varepsilon^4 n^2)$  we will have that  $q^*(X) - q^*(\mathcal{P}) > \varepsilon$ . By this choice of  $k$ , the first term  
221 in the final line above will be at least  $2\varepsilon$ . Also, we have  $N \geq n = \Theta(\sqrt{k}/\varepsilon^2)$ , so when  $k$  is larger  
222 than some absolute constant, the  $O(1/\sqrt{N})$  term in the final line above is  $\Theta(\varepsilon/\sqrt[4]{k}) \leq \varepsilon$ . Thus, by  
223 Claims 2.6 and 2.7, either  $\mathcal{M}$  fails to be accurate, so that  $\exists j \in [k] |a^j - q^j(\mathcal{P})| > \varepsilon$ , or we find a  
224 query  $q^*$  such that  $q^*(X) - q^*(\mathcal{P}) > \varepsilon$ .

### 225 2.3 Lower Bounds for All Algorithms via Random Masks

226 We prove Theorem 1.2 by constructing the following transformation from an adversary that defeats  
227 all natural algorithms to an adversary that defeats all algorithms. The main idea of the reduction is to  
228 use random masks to hide information about the evaluation of the queries at points outside of the  
229 dataset, which effectively forces the algorithm to behave like a natural algorithm because, intuitively,  
230 it does not know where to evaluate the query apart from on the dataset. The reduction is described in  
231 Algorithm 3. Due to space restrictions, we omit its analysis due to space.

---

#### Algorithm 3: $\mathcal{A}_{\text{AQ}}$

Parameters: sample size  $n$ , universe size  $N = \frac{8n}{\varepsilon}$ , number of queries  $k$ , target accuracy  $\varepsilon$ .

Oracle: an adversary  $\mathcal{A}_{\text{NAQ}}$  for natural algorithms with sample size  $n$ , universe size  $N$ , number of  
queries  $k$ , target accuracy  $\varepsilon$ .

Let  $\mathcal{X} = \{(i, y)\}_{i \in [N], y \in \{\pm 1\}^k}$

For  $i \in [N]$

Choose  $m_i = (m_i^1, \dots, m_i^k) \sim \mathcal{U}(\{\pm 1\}^k)$

Let  $\mathcal{P}$  be the uniform distribution over pairs  $(i, m_i)$  for  $i \in [N]$

For  $j \in [k]$

Receive the query  $\hat{q}^j : [N] \rightarrow [\pm 1]$  from  $\mathcal{A}_{\text{NAQ}}$

Form the query  $q^j(i, y) = y^j \oplus m_i^j \oplus \hat{q}^j(i)$  (NB:  $q^j(i, m_i) = \hat{q}^j(i)$ )

Send the query  $q^j$  to  $\mathcal{M}$  and receive the answer  $a^j$

Send the answer  $a^j$  to  $\mathcal{A}_{\text{NAQ}}$

---

232 **3 Post Hoc Generalization Does Not Compose**

233 In this section we prove that post hoc generalization is not closed under composition.

234 **Theorem 3.1.** For every  $n \in \mathbb{N}$  and every  $\alpha > 0$  there is a collection of  $\ell = O(\frac{1}{\alpha} \log n)$  algorithms  
235  $\mathcal{M}_1, \dots, \mathcal{M}_\ell : (\{0, 1\}^{5 \log n})^n \rightarrow \mathcal{Y}$  such that (1) for every  $i = 1, \dots, \ell$  and  $\delta > 0$ ,  $\mathcal{M}_i$  satisfies  
236  $(\varepsilon, \delta)$ -post hoc generalization for  $\varepsilon = O(\sqrt{\log(n/\delta)/n^{1-\alpha}})$ , but (2) the composition  $(\mathcal{M}_1, \dots, \mathcal{M}_\ell)$   
237 is not  $(2 - \frac{2}{n^4}, 1 - \frac{1}{2n^3})$ -post hoc generalizing.

238 The result is based on an algorithm that we call **Encrypermute**. Before proving Theorem 3.1, we  
239 introduce **Encrypermute** and establish the main property that it satisfies.

---

**Algorithm 4: Encrypermute**

---

**Input:** Parameter  $k$ , and a sample  $X = (x_1, x_2, \dots, x_n) \in (\{0, 1\}^d)^n$  for  $d = 5 \log n$ .

**If**  $X$  contains  $n$  distinct elements

Let  $\pi$  be the permutation that sorts  $(x_1, \dots, x_k)$  and identify  $\pi$  with  $r \in \{0, 1, \dots, k! - 1\}$   
Let  $\alpha \in [0, 1]$  be the largest number such that  $k \geq n^\alpha$  and let  $t \leftarrow \alpha k/20$  (NB:  $2^{dt} \leq k!$ )  
Identify  $(x_{k+1}, \dots, x_{k+t}) \in (\{0, 1\}^d)^t$  with a number  $m \in \{0, 1, \dots, k! - 1\}$   
**Return**  $c = m + r \bmod k!$

**Else**

**Return** a random number  $c \in \{0, 1, \dots, k! - 1\}$

---

240 The key facts about **Encrypermute** are as follows.

241 **Claim 3.2.** Let  $\mathcal{D}$  be any distribution over  $(\{0, 1\}^d)^n$ . Let  $D \sim \mathcal{D}$ , let  $X$  be a random permutation  
242 of  $D$ , and let  $C \leftarrow \text{Encrypermute}(X)$ . Then  $D$  and  $C$  are independent.

243 Intuitively, the claim follows from the fact that  $r$  is uniformly random and depends only on the  
244 permutation, so it is independent of  $D$ . Therefore  $m + r \bmod k!$  is random and independent of  $m$ .

245 **Lemma 3.3.**  $\forall \delta > 0$ , **Encrypermute** satisfies  $(\varepsilon, \delta)$ -post hoc generalization for  $\varepsilon = \sqrt{2 \ln(2/\delta)/n}$ .

246 Intuitively the lemma follows from the fact that  $C$  is independent of  $D$ . We omit the proof of both of  
247 these claims due to space restrictions.

248 *Proof of Theorem 3.1.* Fix  $\alpha \in (0, 1)$ , and let  $\mathcal{M}_1$  denote the mechanism that takes a database of  
249 size  $n$  and outputs the first  $n^\alpha$  elements of its sample. As  $\mathcal{M}_1$  outputs a sublinear portion of its input,  
250 it satisfies post hoc generalization with strong parameters. Specifically, by [5, Lemma 3.5],  $\mathcal{M}_1$  is  
251  $(\varepsilon, \delta)$ -post hoc generalizing for  $\varepsilon = O\left(\sqrt{\log(n/\delta)/n^{1-\alpha}}\right)$ .

252 Now consider composing  $\mathcal{M}_1$  with  $O(\frac{1}{\alpha} \log n)$  copies of **Encrypermute**, with exponentially growing  
253 choices for the parameter  $k$ , where for the  $i$ th copy we set  $k = (1 + \frac{\alpha}{20})^i \cdot n^\alpha$ . By Lemma 3.3, each  
254 of these mechanisms satisfies post hoc generalization for  $\varepsilon = O(\sqrt{\log(1/\delta)/n})$ , so this composition  
255 satisfies the assumptions of the theorem.

256 Let  $\mathcal{P}$  be the uniform distribution over  $\{0, 1\}^d$ , where  $d = 5 \log n$ , and let  $X \sim \mathcal{P}^{\otimes n}$ . By a standard  
257 analysis,  $X$  contains  $n$  distinct elements with probability at least  $(1 - \frac{1}{2n^3})$ . Assuming that this  
258 is the case, we have that the first copy of **Encrypermute** outputs  $c = m + r \bmod k!$ , where  $m$   
259 encodes the rows of  $X$  in positions  $n^\alpha + 1, \dots, (1 + \frac{\alpha}{20})n^\alpha$ , and where  $r$  is a deterministic function  
260 of the first  $n^\alpha$  rows of  $X$ . Hence, when composed with  $\mathcal{M}_1$ , these two mechanism reveal the first  
261  $(1 + \frac{\alpha}{20})n^\alpha$  rows of  $X$ . By induction, the output of the composition of all the copies of **Encrypermute**  
262 with  $\mathcal{M}_1$  reveals all of  $X$ . Hence, from the output this composition, we can define the predicate  
263  $q : \{0, 1\}^d \rightarrow \{\pm 1\}$  that evaluates to 1 on every element of  $X$ , and to -1 otherwise. This predicate  
264 satisfies  $q(X) = 1$  but  $q(\mathcal{P}) \leq -1 + 2n/2^d = -1 + 2/n^4$ .  $\square$

265 **References**

266 [1] Raef Bassily, Kobbi Nissim, Adam D. Smith, Thomas Steinke, Uri Stemmer, and Jonathan  
267 Ullman. Algorithmic stability for adaptive data analysis. In *STOC*, 2016.

268 [2] Richard Berk, Lawrence Brown, Andreas Buja, Kai Zhang, Linda Zhao, et al. Valid post-  
269 selection inference. *The Annals of Statistics*, 41(2):802–837, 2013.

270 [3] Andreas Buja, Richard Berk, Lawrence Brown, Edward George, Emil Pitkin, Mikhail Traskin,  
271 Linda Zhao, and Kai Zhang. Models as approximations: A conspiracy of random regressors  
272 and model deviations against classical inference in regression. *Statistical Science*, 1460, 2015.

273 [4] Mark Bun, Thomas Steinke, and Jonathan Ullman. Make up your mind: The price of online  
274 queries in differential privacy. In *SODA*, 2017.

275 [5] Rachel Cummings, Katrina Ligett, Kobbi Nissim, Aaron Roth, and Zhiwei Steven Wu. Adaptive  
276 learning with robust generalization guarantees. In *COLT*, 2016.

277 [6] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron  
278 Roth. Generalization in adaptive data analysis and holdout reuse. In *NIPS*, 2015.

279 [7] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron  
280 Roth. Preserving statistical validity in adaptive data analysis. In *STOC*, 2015.

281 [8] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to  
282 sensitivity in private data analysis. In *TCC*, 2006.

283 [9] Cynthia Dwork, Adam Smith, Thomas Steinke, Jonathan Ullman, and Salil Vadhan. Robust  
284 traceability from trace amounts. In *FOCS*, 2015.

285 [10] Bradley Efron. Estimation and accuracy after model selection. *Journal of the American  
286 Statistical Association*, 109(507):991–1007, 2014.

287 [11] William Fithian, Dennis Sun, and Jonathan Taylor. Optimal inference after model selection.  
288 *arXiv preprint arXiv:1410.2597*, 2014.

289 [12] Andrew Gelman and Eric Loken. The statistical crisis in science. *Am Sci*, 102(6):460, 2014.

290 [13] Moritz Hardt and Jonathan Ullman. Preventing false discovery in interactive data analysis is  
291 hard. In *FOCS*, 2014.

292 [14] Clifford M Hurvich and Chih-Ling Tsai. The impact of model selection on inference in linear  
293 regression. *The American Statistician*, 44(3):214–217, 1990.

294 [15] Michael J. Kearns. Efficient noise-tolerant learning from statistical queries. In *STOC*, 1993.

295 [16] Benedikt M Pötscher. Effects of model selection on inference. *Econometric Theory*, 1991.

296 [17] Daniel Russo and James Zou. Controlling bias in adaptive data analysis using information  
297 theory. In *AISTATS*, 2016.

298 [18] Thomas Steinke and Jonathan Ullman. Interactive fingerprinting codes and the hardness of  
299 preventing false discovery. In *COLT*, 2015.

300 [19] Thomas Steinke and Jonathan Ullman. Tight lower bounds for differentially private selection.  
301 In *FOCS*, 2017.

302 [20] Jonathan Taylor and Robert J Tibshirani. Statistical learning and selective inference. *Proceedings  
303 of the National Academy of Sciences*, 112(25):7629–7634, 2015.

304 [21] Yu-Xiang Wang. *New Paradigms and Optimality Guarantees in Statistical Learning and  
305 Estimation*. PhD thesis, Carnegie Mellon University, 2017.