

1 We are deeply appreciative of reviewers for their feedback amidst these trying circumstances. We are glad that reviewers  
2 appreciated the novelty of the approach for a fundamental problem, the rigor of our analysis, simplicity of the approach  
3 from an implementation perspective.

4 [R3] There seems to be a serious misunderstanding, perhaps due to certain terminology used in our paper. In this paper,  
5 we are indeed interested in comparing different sampling schemes with respect to the number of samples required. Since  
6 we are interested in discrete distributions, most standard sampling schemes, like the two-sample scheme, can be used.  
7 But such schemes suffer from an impractical requirement on the number of samples. The tests in hypothesis testing  
8 literature operate in a *black-box setting*, wherein we have strong lower bounds [4,36,37] over the required number  
9 of samples before a statistically sound conclusion can be drawn. Our setting is not black-box, as pointed out in the  
10 Abstract and Introduction. To borrow from folklore in property testing, consider two distributions  $p$  and  $q$  such that  $p$  is  
11 the uniform distribution over an  $N$  element set  $S$  and  $q$  is either equal to  $p$  or is the uniform distribution over a subset  $T$   
12 of  $S$  where  $|T| = \frac{N}{2}$ , that is for every element in  $T$  the probability of that element in  $q$  is  $\frac{2}{N}$ . If the set  $T$  is unknown,  
13 the number of samples one would need to check if  $q$  is equal to or far from the uniform distribution  $p$  over the set  $S$  is at  
14 least  $\sqrt{N}$ . Since  $N$  is exponential in the number of bits, the sample complexity has an exponential lower bound. This  
15 lower bound is exhibited in [4]. In our setting, since  $N$  is exponential in the input size, the sample complexity is also  
16 exponential. As mentioned in the Introduction, the usage of conditional sampling allows us to circumvent the hardness,  
17 and this work shows how we can handle arbitrary discrete distributions.

18 We refer the reader to citations [8] and [9] for detailed discussions about property testing literature and the limitations  
19 of the standard hypothesis testing setting. While we are aware of the literature on various sampling schemes from  
20 hypothesis testing, our paper does not exactly fit in that line of work and instead fits the property testing line of work,  
21 as pointed out in the Appendix. Therefore, we did not attempt to cast our results in the notions of null and alternate  
22 hypothesis, even though such a characterization is possible but non-standard in property testing literature. Furthermore,  
23 we are not sure about the need for such a formulation when Definition 6 and Theorem 1 are self-contained and precise.

24 [R3] We have access to an Ideal sampler in addition to the weight function (lines 80-85). An explicit description of  
25 ground truth distribution in terms of samples is intractable given the large *model count* (see Appendix E2).

26 [R1, R3] **About the Kernel and the role of randomization.** Kernel, as defined in Definition 7, allows the usage of  
27 randomization internally. As discussed on line 229, we use randomization merely to choose the literal on line 6 of  
28 Algorithm 2; we will add further discussion to this effect in the final version (We have provided source code with a  
29 detailed description). The Kernel in our paper has nothing to do with RKHS. Perhaps we should use a different name to  
30 avoid confusion.

31 [R4] As is the norm, we focused on binary variables, but the techniques work for any arbitrary discrete domain.

32 [R1, R4]: As a first step to demonstrate the working of the prototype implementation, we focused on log-linear models  
33 for which the most efficient technique is to employ the inverse transform. This also allows us to easily obtain a sampler  
34 with formal guarantees and two other samplers without guarantees. It is worth noting that STS and variants of UniGen  
35 have indeed appeared in AI/ML conferences, so we have taken the first step in addressing the distributions of interest to  
36 the AI/ML community. We do agree with [R2] (and also strongly believe) that our algorithmic framework will inspire  
37 follow-up work focused on further improvements and a deeper understanding of various complex distributions. Again,  
38 our algorithmic analysis is not restricted to log-linear models.

39 [R1] **Do sampling algorithms satisfy the ‘non-adversarial’ condition ?** Since conditioning is a fundamental opera-  
40 tion in probabilistic reasoning, one would expect that the underlying sampling algorithms support conditioning. Our  
41 non-adversariality assumption simply translates to support for conditioning. We will emphasize this further. At the  
42 same time, we do agree with the reviewer that one can claim that their sampler supports only one particular distribution,  
43 but in such a case, it is hard to see if one can do any better given the strong lower bounds [4,36,37].

44 [R1] **What does it mean for  $\hat{\varphi}$  to have "similar structure" to  $\varphi$ ?**  $\hat{\varphi}$  is constructed by conjuncting  $\varphi$  with a small  
45 formula which allows conditional sampling on the models of  $\varphi$ ; thus it can be said that  $\hat{\varphi}$  retains the structure of  $\varphi$ .

46 [R2] **What is the performance like without the constraint obfuscation heuristic, does the naive approach work?**  
47 Our preliminary studies showcased the need for constraint obfuscation methodology. We will seek to conduct a longer  
48 empirical study and add these results in the final version.

49 [R4] **The absence of a tight and practical bound.** It is worth emphasizing that Theorem 1 only implies the worst-case  
50 behavior of Fulcrum, and lack of tight upper bound does not imply that the performance of Fulcrum would be worse,  
51 a claim supported by experimental analysis as well. As a similar example, while there are no bounds better than  
52 exponential time for any algorithm to solve SAT, the SAT solvers routinely perform orders of magnitude better than  
53 a naive exponential time algorithm. In this view, the design of an efficient algorithm should be viewed as a primary  
54 contribution even though its worst-case behavior is hard to analyze.