

1 Thanks for all the comments! We answer major comments from each reviewer below; we'll fix the minor ones.

2 **REVIEWER 1:** "This paper ranks high in novelty...The experimental results are strong, especially on Text
3 modelling. Moreover, the proposed method significantly computationally more efficient compared to competing
4 approaches...The paper is very well written and easy to understand." Thanks for the high praise!

5 "1. The authors should consider a wider array of distributions to more convincingly demonstrate the capabilities
6 of the proposed flow layers." We include many synthetic discrete problems in the Experiments section: full rank
7 discrete distributions, addition, and Potts models. We decided against including more continuous/ordinal data problems
8 as our flows don't take advantage of ordinal structure in the Categorical states, and there are many practical problems
9 without this structure.

10 "2. Some important details are unclear. E.g. what is the base distribution for sampling? Is it the factorised
11 marginal distribution? If it is how is it estimated in high dimensions (given that the number of data samples
12 needed for an accurate estimate would grow exponentially)?" We use factorized Categorical base distributions for
13 bipartite flow models and correlated Categorical base distributions parameterized by an autoregressive neural network
14 for autoregressive flow models. Note that the number of parameters for the base distribution is far less than the total
15 number of possible states (e.g., $K \cdot D$ for factorized Categorical instead of K^D).

16 **REVIEWER 2:** "Originality: This paper is the first demonstration of flow-based models to discrete data.
17 As such, the work is fairly novel...That being said, the main technical contribution amounts to...on top of the
18 existing techniques. I view this simplicity as a benefit of the approach, but some may view this a simple extension
19 of existing techniques." Thanks for the high praise! We agree about simplicity being a benefit. We explicitly designed
20 discrete flows to be natural extensions of the continuous case.

21 "Quality: The technical and experimental aspects of the paper are well-executed. Clarity: The presentation of
22 the approach is incredibly clear." Thanks!

23 "For the most part, the experiments section is also clear. Some details of the models and training set-up are un-
24 clear, particularly in the toy examples from sections 4.1 - 4.3. Additional details in the supplementary material
25 would help to clear up confusion." Great idea. We'll clean this up .

26 **Adding discrete image datasets.** Demonstrating discrete flows on image data is a good idea. So far, we focused on
27 text applications to show a domain normalizing flows hadn't yet been applied to. It's a ripe area given that both flexible
28 modeling with bidirectionality and nonautoregressive generation are of huge interest for text. We'd love to explore
29 images in future work, in particular pushing on the ordinality of pixel intensities to better handle data quantization.
30 Hooeboom et al. (2019) provide excellent complementary work in that direction.

31 "In Section 3.1, 3.2., or in the supplementary, it would be helpful to have an expanded discussion around when
32 discrete flows are invertible and what difficulties there are in ensuring this aspect. This discussion could also
33 include the invertible discrete functions alluded to in Section 5." In the revision, we'll add a section about the
34 expressivity of discrete flows, what the set of transformations are (permutation-based), and designing parameterized
35 invertible discrete functions.

36 **REVIEWER 3:** "The approach taken at this point is to ignore the problem, and employ the straight-through
37 estimator (which the authors argue work well for problems where K is not too large)." The straight-through
38 estimator is effective and commonly used in many discrete optimization problems to compute gradients through the
39 non-differentiable argmax operation. In future work we hope to investigate whether other gradient approximations
40 (including non-gradient based optimization) improve performance.

41 "The paper is clearly original, and I imagine it will be of great significance (it brings two interests together,
42 namely, flexible flow-based density estimation, and modelling discrete data)." Thanks!

43 "Since you count on backpropagation via straight-through estimator (STE), the derivative of mod K becomes
44 relevant (as it will be necessary for chain rule to update the parameters of the NN that predict flow parameters).
45 The best I could gather is that it's probably 1 everywhere except when $(u + \sigma x)$ is exactly divisible by K . Is
46 that correct?" You're correct that the operations are non-differentiable. The goal with continuous relaxations is not to
47 compute the true gradients, but rather to provide a useful signal for improving the loss. We'll add these nuances to the
48 paper.

49 "You mention STE works well if K is not too large, but is that all?" That's a great question. Depth (number of
50 flows) affects gradients, since the bias explicitly accumulates as each flow uses an approximation. This is also the case
51 for dimensionality (sequence length). We haven't found complexity of the networks parameterizing the flow to make
52 a difference, but this requires further investigate. Note an additional complication follows your above point: "true
53 gradients" are not well-defined. So instead of examining gradient bias, bias may be formalized by comparing, for
54 example, the minima from the discrete loss to the minima from the relaxed loss.