

---

# SymmetricDiffusers: Learning Discrete Diffusion on Finite Symmetric Groups

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Finite symmetric groups  $S_n$  are essential in fields such as combinatorics, physics,  
2 and chemistry. However, learning a probability distribution over  $S_n$  poses signif-  
3 icant challenges due to its intractable size and discrete nature. In this paper, we  
4 introduce SymmetricDiffusers, a novel discrete diffusion model that simplifies the  
5 task of learning a complicated distribution over  $S_n$  by decomposing it into learning  
6 simpler transitions of the reverse diffusion using deep neural networks. We identify  
7 the riffle shuffle as an effective forward transition and provide empirical guidelines  
8 for selecting the diffusion length based on the theory of random walks on finite  
9 groups. Additionally, we propose a generalized Plackett-Luce (PL) distribution for  
10 the reverse transition, which is probably more expressive than the PL distribution.  
11 We further introduce a theoretically grounded "denoising schedule" to improve  
12 sampling and learning efficiency. Extensive experiments show that our model  
13 achieves state-of-the-art or comparable performances on solving tasks including  
14 sorting 4-digit MNIST images, jigsaw puzzles, and traveling salesman problems.

## 15 1 Introduction

16 As a vital area of abstract algebra, finite groups provide a structured framework for analyzing symme-  
17 tries and transformations which are fundamental to a wide range of fields, including combinatorics,  
18 physics, chemistry, and computer science. One of the most important finite groups is the *finite*  
19 *symmetric group*  $S_n$ , defined as the group whose elements are all the bijections (or permutations)  
20 from a set of  $n$  elements to itself, with the group operation being function composition.

21 Classic probabilistic models for finite symmetric groups  $S_n$ , such as the Plackett-Luce (PL) model  
22 [35, 27], the Mallows model [28], and card shuffling methods [9], are crucial in analyzing preference  
23 data and understanding the convergence of random walks. Therefore, studying probabilistic models  
24 over  $S_n$  through the lens of modern machine learning is both natural and beneficial. This problem is  
25 theoretically intriguing as it bridges abstract algebra and machine learning. For instance, Cayley's  
26 Theorem, a fundamental result in abstract algebra, states that every group is isomorphic to a subgroup  
27 of a symmetric group. This implies that learning a probability distribution over finite symmetric  
28 groups could, in principle, yield a distribution over any finite group. Moreover, exploring this problem  
29 could lead to the development of advanced models capable of addressing tasks such as permutations  
30 in ranking problems, sequence alignment in bioinformatics, and sorting.

31 However, learning a probability distribution over finite symmetric groups  $S_n$  poses significant  
32 challenges. First, the number of permutations of  $n$  objects grows factorially with  $n$ , making the  
33 inference and learning computationally expensive for large  $n$ . Second, the discrete nature of the data  
34 brings difficulties in designing expressive parameterizations and impedes the gradient-based learning.

35 In this work, we propose a novel discrete (state space) diffusion model over finite symmetric groups,  
36 dubbed as *SymmetricDiffusers*. It overcomes the above challenges by decomposing the difficult

37 problem of learning a complicated distribution over  $S_n$  into a sequence of simpler problems, *i.e.*,  
 38 learning individual transitions of a reverse diffusion process using deep neural networks. Based on  
 39 the theory of random walks on finite groups, we investigate various shuffling methods as the forward  
 40 process and identify the riffle shuffle as the most effective. We also provide empirical guidelines  
 41 on choosing the diffusion length based on the mixing time of the riffle shuffle. Furthermore, we  
 42 examine potential transitions for the reverse diffusion, such as inverse shuffling methods and the  
 43 PL distribution, and introduce a novel generalized PL distribution. We prove that our generalized  
 44 PL is more expressive than the PL distribution. Additionally, we propose a theoretically grounded  
 45 "denoising schedule" that merges reverse steps to improve the efficiency of sampling and learning.  
 46 To validate the effectiveness of our SymmetricDiffusers, we conduct extensive experiments on three  
 47 tasks: sorting 4-Digit MNIST images, solving Jigsaw Puzzles on the Noisy MNIST and CIFAR-10  
 48 datasets, and addressing traveling salesman problems (TSPs). Our model achieves the state-of-the-art  
 49 or comparable performance across all tasks.

## 50 2 Related Works

51 **Random Walks on Finite Groups.** The field of random walks on finite groups, especially finite  
 52 symmetric groups, have been extensively studied by previous mathematicians [37, 11, 4, 38]. Tech-  
 53 niques from a variety of different fields, including probability, combinatorics, and representation  
 54 theory, have been used to study random walks on finite groups [38]. In particular, random walks on  
 55 finite symmetric groups are first studied in the application of card shuffling, with many profound  
 56 theoretical results of shuffling established. A famous result in the field shows that 7 riffle shuffles are  
 57 enough to mix up a deck of 52 cards [4], where a riffle shuffle is a mathematically precise model that  
 58 simulates how people shuffle cards in real life. The idea of shuffling to mix up a deck of cards aligns  
 59 naturally with the idea of diffusion, and we seek to fuse the modern techniques of diffusion models  
 60 with the classical theories of random walks on finite groups.

61 **Diffusion Models.** Diffusion models [40, 41, 16, 42] are a powerful class of generative models that  
 62 typically deals with continuous data. They consist of forward and reverse processes. The forward  
 63 process is typically a discrete-time continuous-state Markov chain or a continuous-time continuous-  
 64 state Markov process that gradually adds noise to data, and the reverse process learn neural networks  
 65 to denoise. Discrete (state space) diffusion models have also been proposed to handle discrete data  
 66 like image, text [3], and graphs [45]. Existing discrete diffusion models are applicable for learning  
 67 distributions of permutations. However, they focused on cases where the state space is small or has a  
 68 special (*e.g.*, decomposable) structure and are unable to deal with intractable-sized state spaces like  
 69 the symmetric group. In particular, [3] requires an explicit transition matrix, which has size  $n! \times n!$   
 70 in the case of finite symmetric groups and has no simple representations or sparsifications.

71 **Differentiable Sorting and Learning Permutations.** A popular paradigm to learn permutations  
 72 is through differentiable sorting or matching algorithms. Various differentiable sorting algorithms  
 73 have been proposed that uses continuous relaxations of permutation matrices [13, 8, 5], or uses  
 74 differentiable swap functions [33, 34, 20]. The Gumbel-Sinkhorn method [29] has also been proposed  
 75 to learn latent permutations using the continuous Sinkhorn operator. Such methods often focus on  
 76 finding the optimal permutation instead of learning a distribution over the finite symmetric group.  
 77 Moreover, they tend to be less effective as  $n$  grows larger due to their high complexities.

## 78 3 Learning Diffusion Models on Finite Symmetric Groups

79 We first introduce some notations. Fix  $n \in \mathbb{N}$ . Let  $[n]$  denote the set  $\{1, 2, \dots, n\}$ . A *permutation*  
 80  $\sigma$  on  $[n]$  is a function from  $[n]$  to  $[n]$ , and we usually write  $\sigma$  as  $\begin{pmatrix} 1 & 2 & \dots & n \\ \sigma(1) & \sigma(2) & \dots & \sigma(n) \end{pmatrix}$ . The  
 81 *identity permutation*, denoted by  $\text{Id}$ , is the permutation given by  $\text{Id}(i) = i$  for all  $i \in [n]$ . Let  
 82  $S_n$  be the set of all permutations (or bijections) from a set of  $n$  elements to itself, called the *finite*  
 83 *symmetric group*, whose group operation is the function composition. For a permutation  $\sigma \in S_n$ ,  
 84 the permutation matrix  $Q_\sigma \in \mathbb{R}^{n \times n}$  associated with  $\sigma$  satisfies  $e_i^\top Q_\sigma = e_{\sigma(i)}^\top$  for all  $i \in [n]$ . In  
 85 this paper, we consider a set of  $n$  distinctive objects  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where the  $i$ -th object is  
 86 represented by a  $d$ -dimensional vector  $\mathbf{x}_i$ . Therefore, a ranked list of objects can be represented as  
 87 a matrix  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ , where the ordering of rows corresponds to the ordering of  
 88 objects. We can permute  $X$  via permutation  $\sigma$  to obtain  $Q_\sigma X$ .

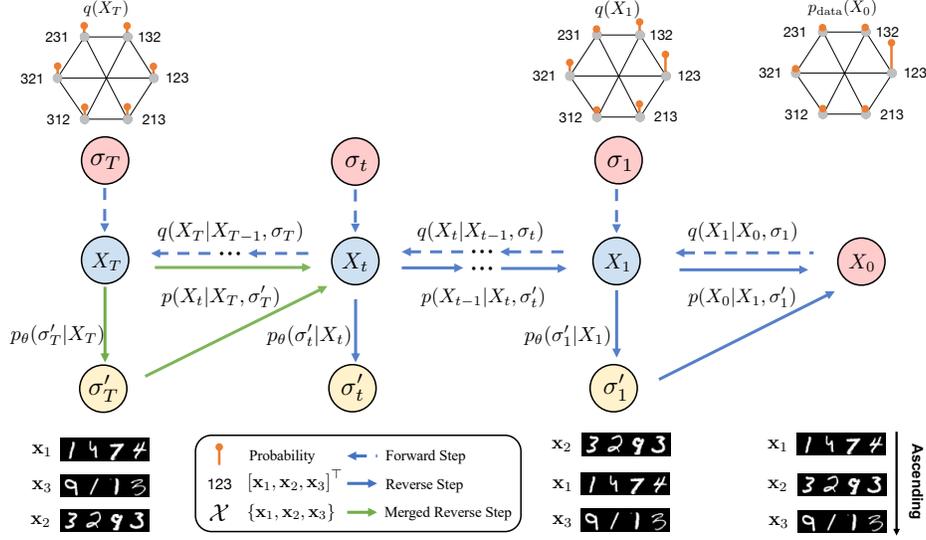


Figure 1: This figure illustrates our discrete diffusion model on finite symmetric groups. The middle graphical model displays the forward and reverse diffusion processes. We demonstrate learning distributions over the symmetric group  $S_3$  via the task of sorting three MNIST 4-digit images. The top part of the figure shows the marginal distribution of a ranked list of images  $X_t$  at time  $t$ , while the bottom shows a randomly drawn list of images.

89 Our goal is to learn a distribution over  $S_n$ . We propose learning discrete (state space) diffusion  
90 models, which consist of a *forward process* and a *reverse process*. In the forward process, starting  
91 from the unknown data distribution, we simulate a random walk until it reaches a known stationary  
92 “noise” distribution. In the reverse process, starting from the known noise distribution, we simulate  
93 another random walk, where the transition probability is computed using a neural network, until it  
94 recovers the data distribution. Learning a transition distribution over  $S_n$  is often more manageable  
95 than learning the original distribution because: (1) the support size (the number of states that can be  
96 reached in one transition) could be much smaller than  $n!$ , and (2) the distance between the initial and  
97 target distributions is smaller. By doing so, we break down the hard problem (learning the original  
98 distribution) into a sequence of simpler subproblems (learning the transition distribution). The overall  
99 framework is illustrated in Fig. 1. In the following, we will introduce the forward card shuffling  
100 process in Section 3.1, the reverse process in Section 3.2, the network architecture and training in  
101 Section 3.3, denoising schedule in Section 3.4, and reverse decoding methods in Section 3.5.

### 102 3.1 Forward Diffusion Process: Card Shuffling

103 Suppose we observe a set of objects  $\mathcal{X}$  and their ranked list  $X_0$ . They are assumed to be generated  
104 from an unknown data distribution in an IID manner, *i.e.*,  $X_0, \mathcal{X} \stackrel{\text{iid}}{\sim} p_{\text{data}}(X, \mathcal{X})$ . One can construct a  
105 bijection between a ranked list of  $n$  objects and an ordered deck of  $n$  cards. Therefore, permuting  
106 objects is equivalent to shuffling cards. In the forward diffusion process, we would like to add  
107 “random noise” to the rank list so that it reaches to some known stationary distribution like the  
108 uniform. Formally, we let  $\mathcal{S} \subseteq S_n$  be a set of permutations that are realizable by a given shuffling  
109 method in one step.  $\mathcal{S}$  does not change across steps in common shuffling methods. We will provide  
110 concrete examples later. We then define the *forward process* as a Markov chain,

$$q(X_{1:T}|X_0, \mathcal{X}) = q(X_{1:T}|X_0) = \prod_{t=1}^T q(X_t|X_{t-1}), \quad (1)$$

111 where  $q(X_t|X_{t-1}) = \sum_{\sigma_t \in \mathcal{S}} q(X_t|X_{t-1}, \sigma_t)q(\sigma_t)$  and the first equality in Eq. (1) holds since  $X_0$   
112 implies  $\mathcal{X}$ . In the forward process, although the set  $\mathcal{X}$  does not change, the rank list of objects  $X_t$   
113 changes. Here  $q(\sigma_t)$  has the support  $\mathcal{S}$  and describes the permutation generated by the underlying  
114 shuffling method. Note that common shuffling methods are time-homogeneous Markov chains, *i.e.*,  
115  $q(\sigma_t)$  stays the same across time.  $q(X_t|X_{t-1}, \sigma_t)$  is a delta distribution  $\delta(X_t = Q_{\sigma_t}X_{t-1})$  since the  
116 permuted objects  $X_t$  are uniquely determined given the permutation  $\sigma_t$  and  $X_{t-1}$ . We denote the  
117 *neighbouring states* of  $X$  via one-step shuffling as  $N_{\mathcal{S}}(X) := \{Q_{\sigma}X | \sigma \in \mathcal{S}\}$ . Therefore, we have,

$$q(X_t|X_{t-1}) = \begin{cases} q(\sigma_t) & \text{if } X_t \in N_{\mathcal{S}}(X_{t-1}) \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

118 Note that  $X_t \in N_{\mathcal{S}}(X_{t-1})$  is equivalent to  $\sigma_t \in \mathcal{S}$  and  $X_t = Q_{\sigma_t} X_{t-1}$ .

### 119 3.1.1 Card Shuffling Methods

120 We now consider several popular shuffling methods as the forward transition, *i.e.*, *random transpo-*  
 121 *sitions*, *random insertions*, and *riffle shuffles*. Different shuffling methods provide different design  
 122 choices of  $q(\sigma_t)$ , thus corresponding to different forward diffusion processes. Although all these  
 123 forward diffusion processes share the same stationary distribution, *i.e.*, the uniform, they differ in  
 124 their mixing time. We will introduce stronger quantitative results on their mixing time later.

125 **Random Transpositions.** One natural way of shuffling is to swap pairs of objects. Formally, a  
 126 *transposition* or a *swap* is a permutation  $\sigma \in S_n$  such that there exist  $i \neq j \in [n]$  with  $\sigma(i) = j$ ,  
 127  $\sigma(j) = i$ , and  $\sigma(k) = k$  for all  $k \notin \{i, j\}$ , in which case we denote  $\sigma = (i \ j)$ . We let  $\mathcal{S} =$   
 128  $\{(i \ j) : i \neq j \in [n]\} \cup \{\text{Id}\}$ . For any time  $t$ , we define  $q(\sigma_t)$  by choosing two indices from  $[n]$   
 129 uniformly and independently and swap the two indices. If the two chosen indices are the same, then  
 130 this means that we have sampled the identity permutation. Specifically,  $q(\sigma_t = (i \ j)) = 2/n^2$   
 131 when  $i \neq j$  and  $q(\sigma_t = \text{Id}) = 1/n$ .

132 **Random Insertions.** Another shuffling method is to insert the last piece to somewhere in the middle.  
 133 Let  $\text{insert}_i$  denote the permutation that inserts the last piece right before the  $i^{\text{th}}$  piece, and let  
 134  $\mathcal{S} := \{\text{insert}_i : i \in [n]\}$ . Note that  $\text{insert}_n = \text{Id}$ . Specifically, we have  $q(\sigma_t = \text{insert}_i) = 1/n$   
 135 when  $i \neq n$  and  $q(\sigma_t = \text{Id}) = 1/n$ .

136 **Riffle Shuffles.** Finally, we introduce the riffle shuffle, a method similar to how serious card players  
 137 shuffle cards. The process begins by roughly cutting the deck into two halves and then interleaving the  
 138 two halves together. A formal mathematical model of the riffle shuffle, known as the *GSR model*, was  
 139 introduced by Gilbert and Shannon [11], and independently by Reeds [37]. The model is described  
 140 as follows. A deck of  $n$  cards is cut into two piles according to binomial distribution, where the  
 141 probability of having  $k$  cards in the top pile is  $\binom{n}{k}/2^n$  for  $0 \leq k \leq n$ . The top pile is held in the  
 142 left hand and the bottom pile in the right hand. The two piles are then riffled together such that, if  
 143 there are  $A$  cards left in the left hand and  $B$  cards in the right hand, the probability that the next card  
 144 drops from the left is  $A/(A+B)$ , and from right is  $B/(A+B)$ . We implement the riffle shuffles  
 145 according to the GSR model. For simplicity, we will omit the term ‘‘GSR’’ when referring to riffle  
 146 shuffles hereafter.

147 There exists an exact formula for the probability over  $S_n$  obtained through one-step riffle shuffle.  
 148 Let  $\sigma \in S_n$ . A *rising sequence* of  $\sigma$  is a subsequence of  $\sigma$  constructed by finding a maximal  
 149 subset of indices  $i_1 < i_2 < \dots < i_j$  such that permuted values are contiguously increasing, *i.e.*,  
 150  $\sigma(i_2) - \sigma(i_1) = \sigma(i_3) - \sigma(i_2) = \dots = \sigma(i_j) - \sigma(i_{j-1}) = 1$ . For example, the permutation  
 151  $\begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ \color{red}{1} & \color{red}{4} & \color{red}{2} & \color{red}{5} & \color{red}{3} \end{pmatrix}$  has 2 rising sequences, *i.e.*, 123 (red) and 45 (blue). Note that a permutation  
 152 has 1 rising sequence if and only if it is the identity permutation. Denoting by  $q_{\text{RS}}(\sigma)$  the probability  
 153 of obtaining  $\sigma$  through one-step riffle shuffle, it is shown in [4] that

$$q_{\text{RS}}(\sigma) = \frac{1}{2^n} \binom{n+2-r}{n} = \begin{cases} (n+1)/2^n & \text{if } \sigma = \text{Id} \\ 1/2^n & \text{if } \sigma \text{ has two rising sequences} \\ 0 & \text{otherwise,} \end{cases} \quad (3)$$

154 where  $r$  is the number of rising sequences of  $\sigma$ . The support  $\mathcal{S}$  is thus the set of all permutations with  
 155 at most two rising sequences. We let the forward process be  $q(\sigma_t) = q_{\text{RS}}(\sigma_t)$  for all  $t$ .

### 156 3.1.2 Mixing Times and Cut-off Phenomenon

157 All of the above shuffling methods have the uniform distribution as the stationary distribution.  
 158 However, they have different mixing times (*i.e.*, the time until the Markov chain is close to its  
 159 stationary distribution measured by some distance), and there exist quantitative results on their mixing  
 160 times. Let  $q \in \{q_{\text{RT}}, q_{\text{RI}}, q_{\text{RS}}\}$ , and for  $t \in \mathbb{N}$ , let  $q^{(t)}$  be the marginal distribution of the Markov  
 161 chain after  $t$  shuffles. We describe the mixing time in terms of the total variation (TV) distance  
 162 between two probability distributions, *i.e.*,  $D_{\text{TV}}(q^{(t)}, u)$ , where  $u$  is the uniform distribution.

163 For all three shuffling methods, there exists a *cut-off phenomenon*, where  $D_{\text{TV}}(q^{(t)}, u)$  stays around  
 164 1 for initial steps and then abruptly drops to values that are close to 0. The *cut-off time* is the time  
 165 when the abrupt change happens. For the formal definition, we refer the readers to Definition 3.3 of

166 [38]. In [38], they also provided the cut-off time for random transposition, random insertion, and  
 167 riffle shuffle, which are  $\frac{n}{2} \log n$ ,  $n \log n$ , and  $\frac{3}{2} \log_2 n$  respectively. Observe that the riffle shuffle  
 168 reaches the cut-off much faster than the other two methods, which means it has a much faster mixing  
 169 time. Therefore, we use the riffle shuffle in the forward process.

### 170 3.2 The Reverse Diffusion Process

171 We now model the *reverse process* as another Markov chain conditioned on the set of objects  $\mathcal{X}$ . We  
 172 denote the set of realizable *reverse permutations* as  $\mathcal{T}$ , and the neighbours of  $X$  with respect to  $\mathcal{T}$  as  
 173  $N_{\mathcal{T}}(X) := \{Q_{\sigma} X : \sigma \in \mathcal{T}\}$ . The conditional joint distribution is given by

$$p_{\theta}(X_{0:T}|\mathcal{X}) = p(X_T|\mathcal{X}) \prod_{t=1}^T p_{\theta}(X_{t-1}|X_t), \quad (4)$$

174 where  $p_{\theta}(X_{t-1}|X_t) = \sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) p_{\theta}(\sigma'_t|X_t)$ . To sample from  $p(X_T|\mathcal{X})$ , one simply  
 175 samples a random permutation from the uniform distribution and then shuffle the objects accordingly  
 176 to obtain  $X_T$ .  $p(X_{t-1}|X_t, \sigma'_t)$  is again a delta distribution  $\delta(X_{t-1} = Q_{\sigma'_t} X_t)$ . We have

$$p_{\theta}(X_{t-1}|X_t) = \begin{cases} p_{\theta}(\sigma'_t|X_t) & \text{if } X_{t-1} \in N_{\mathcal{T}}(X_t) \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

177 where  $X_{t-1} \in N_{\mathcal{T}}(X_t)$  is equivalent to  $\sigma'_t \in \mathcal{T}$  and  $X_{t-1} = Q_{\sigma'_t} X_t$ . In the following, we will  
 178 introduce the specific design choices of the distribution  $p_{\theta}(\sigma'_t|X_t)$ .

#### 179 3.2.1 Inverse Card Shuffling

180 A natural choice is to use the inverse operations of the aforementioned card shuffling operations in  
 181 the forward process. Specifically, for the forward shuffling  $\mathcal{S}$ , we introduce their inverse operations  
 182  $\mathcal{T} := \{\sigma^{-1} : \sigma \in \mathcal{S}\}$ , from which we can parameterize  $p_{\theta}(\sigma'_t|X_t)$ .

183 **Inverse Transposition.** Since the inverse of a transposition is also a transposition, we can let  
 184  $\mathcal{T} := \mathcal{S} = \{(i \ j) : i \neq j \in [n]\} \cup \{\text{Id}\}$ . We define a distribution of inverse transposition (IT) over  
 185  $\mathcal{T}$  using  $n + 1$  real-valued parameters  $\mathbf{s} = (s_1, \dots, s_n)$  and  $\tau$  such that

$$p_{\text{IT}}(\sigma) = \begin{cases} 1 - \phi(\tau) & \text{if } \sigma = \text{Id} \\ \phi(\tau) (\psi(\mathbf{s}, \pi_{ij})_1 \psi(\mathbf{s}, \pi_{ij})_2 + \psi(\mathbf{s}, \pi_{ji})_1 \psi(\mathbf{s}, \pi_{ji})_2) & \text{if } \sigma = (i \ j) \text{ where } i \neq j, \end{cases} \quad (6)$$

186 where  $\psi(\mathbf{s}, \pi)_i = \exp(s_{\pi(i)}) / (\sum_{k=i}^n \exp(s_{\pi(k)}))$  and  $\phi(\cdot)$  is the sigmoid function.  $\pi_{ij}$  is any  
 187 permutation starting with  $i$  and  $j$ , *i.e.*,  $\pi_{ij}(1) = i$  and  $\pi_{ij}(2) = j$ .  $\pi_{ji}$  is any permutation starting  
 188 with  $j$  and  $i$ , *i.e.*,  $\pi_{ji}(1) = j$  and  $\pi_{ji}(2) = i$ .

189 **Inverse Insertion.** For the random insertion, the inverse operation is to insert some piece to the  
 190 end. Let  $\text{inverse\_insert}_i$  denote the permutation that moves the  $i^{\text{th}}$  component to the end, and  
 191 let  $\mathcal{T} := \{\text{inverse\_insert}_i : i \in [n]\}$ . We define a categorical distribution of inverse insertion (II)  
 192 over  $\mathcal{T}$  using parameters  $\mathbf{s} = (s_1, \dots, s_n)$  such that,

$$p_{\text{II}}(\sigma = \text{inverse\_insert}_i) = \exp(s_i) / \left( \sum_{j=1}^n \exp(s_j) \right). \quad (7)$$

193 **Inverse Riffle Shuffle.** In the riffle shuffle, the deck of card is first cut into two piles, and the two piles  
 194 are riffled together. So to undo a riffle shuffle, we need to figure out which pile each card belongs to,  
 195 *i.e.*, making a sequence of  $n$  binary decisions. We define the Inverse Riffle Shuffle (IRS) distribution  
 196 using parameters  $\mathbf{s} = (s_1, \dots, s_n)$  as follows. Starting from the last (the  $n^{\text{th}}$ ) object, each object  $i$   
 197 has probability  $\phi(s_i)$  of being put on the top of the left pile. Otherwise, it falls on the top of the right  
 198 pile. Finally, put the left pile on top of the right pile, which gives the shuffled result.

#### 199 3.2.2 The Plackett-Luce Distribution and Its Generalization

200 Other than specific inverse shuffling methods to parameterize the reverse process, we also consider  
 201 general distributions  $p_{\theta}(\sigma'_t|X_t)$  whose support are the whole  $S_n$ , *i.e.*,  $\mathcal{T} = S_n$ .

202 **The PL Distribution.** A popular distribution over  $S_n$  is the Plackett-Luce (PL) distribution [35, 27],  
 203 which is constructed from  $n$  real-valued scores  $\mathbf{s} = (s_1, \dots, s_n)$  as follows,

$$p_{\text{PL}}(\sigma) = \prod_{i=1}^n \exp(s_{\sigma(i)}) / \left( \sum_{j=i}^n \exp(s_{\sigma(j)}) \right), \quad (8)$$

204 for all  $\sigma \in S_n$ . Intuitively,  $(s_1, \dots, s_n)$  represents the preference given to each index in  $[n]$ . To  
 205 sample from  $\text{PL}_s$ , we first sample  $\sigma(1)$  from  $\text{Cat}(n, \text{softmax}(s))$ . Then we remove  $\sigma(1)$  from the  
 206 list and sample  $\sigma(2)$  from the categorical distribution corresponding to the rest of the scores (logits).  
 207 We continue in this manner until we have sampled  $\sigma(1), \dots, \sigma(n)$ . By [7], the mode of the PL  
 208 distribution is the permutation that sorts  $s$  in descending order.

209 **The Generalized PL (GPL) Distribution.** We also propose a generalization of the PL distribution,  
 210 referred to as *Generalized Plackett-Luce (GPL) Distribution*. Unlike the PL distribution, which uses  
 211 a set of  $n$  scores, the GPL distribution uses  $n^2$  scores  $\{s_1, \dots, s_n\}$ , where each  $s_i = \{s_{i,1}, \dots, s_{i,n}\}$   
 212 consists of  $n$  scores. The GPL distribution is constructed as follows,

$$p_{\text{GPL}}(\sigma) := \prod_{i=1}^n \exp(s_{i,\sigma(i)}) / \left( \sum_{j=i}^n \exp(s_{i,\sigma(j)}) \right). \quad (9)$$

213 Sampling of the GPL distribution begins with sampling  $\sigma(1)$  using  $n$  scores  $s_1$ . For  $2 \leq i \leq n$ , we  
 214 remove  $i-1$  scores from  $s_i$  that correspond to  $\sigma(1), \dots, \sigma(i-1)$  and sample  $\sigma(i)$  from a categorical  
 215 distribution constructed from the remaining  $n-i+1$  scores in  $s_i$ . It is important to note that the  
 216 family of PL distributions is a strict subset of the GPL family. Since the GPL distribution has more  
 217 parameters than the PL distribution, it is expected to be more expressive. In fact, when considering  
 218 their ability to express the delta distribution, which is the target distribution for many permutation  
 219 learning problems, we have the following result.

220 **Proposition 1.** *The PL distribution cannot exactly represent a delta distribution. That is, there does*  
 221 *not exist an  $s$  such that  $p_{\text{PL}} = \delta_\sigma$  for any  $\sigma \in S_n$ , where  $\delta_\sigma(\sigma) = 1$  and  $\delta_\sigma(\pi) = 0$  for all  $\pi \neq \sigma$ .*  
 222 *But the GPL distribution can represent a delta distribution exactly.*

### 223 3.3 Network Architecture and Training

224 We now briefly introduce how to use neural networks to parameterize the above distributions used  
 225 in the reverse process. At any time  $t$ , given  $X_t \in \mathbb{R}^{n \times d}$ , we use a neural network with parameters  
 226  $\theta$  to construct  $p_\theta(\sigma'_t | X_t)$ . In particular, we treat  $n$  rows of  $X_t$  as  $n$  tokens and use a Transformer  
 227 architecture along with the time embedding of  $t$  and the positional encoding to predict the previously  
 228 mentioned scores. For example, for the GPL distribution, to predict  $n^2$  scores, we introduce  $n$  dummy  
 229 tokens that correspond to the  $n$  permuted output positions. We then perform a few layers of masked  
 230 self-attention ( $2n \times 2n$ ) to obtain the token embedding  $Z_1 \in \mathbb{R}^{n \times d_{\text{model}}}$  corresponding to  $n$  input  
 231 tokens and  $Z_2 \in \mathbb{R}^{n \times d_{\text{model}}}$  corresponding to  $n$  dummy tokens. Finally, the GPL score matrix is  
 232 obtained as  $S_\theta = Z_1 Z_2^\top \in \mathbb{R}^{n \times n}$ . Since the aforementioned distributions have different numbers of  
 233 scores, the specific architectures of the Transformer differ. We provide more details in Appendix B.

234 To learn the diffusion model, we maximize the following variational lower bound:

$$\mathbb{E}_{p_{\text{data}}(X_0, \mathcal{X})} \left[ \log p_\theta(X_0 | \mathcal{X}) \right] \geq \mathbb{E}_{p_{\text{data}}(X_0, \mathcal{X}) q(X_{1:T} | X_0, \mathcal{X})} \left[ \log p(X_T | \mathcal{X}) + \sum_{t=1}^T \log \frac{p_\theta(X_{t-1} | X_t)}{q(X_t | X_{t-1})} \right]. \quad (10)$$

235 In practice, one can draw samples to obtain the Monte Carlo estimation of the lower bound. Due to  
 236 the complexity of shuffling transition in the forward process, we can not obtain  $q(X_t | X_0)$  analytically,  
 237 as is done in common diffusion models [16, 3]. Therefore, we have to run the forward process to  
 238 collect samples. Fortunately, it is efficient as the forward process only involves shuffling integers. We  
 239 include more training details in Appendix E.

### 240 3.4 Denoising Schedule via Merging Reverse Steps

241 If one merges some steps in the reverse process, sampling and learning would be faster and more  
 242 memory efficient. The variance of the training loss could also be reduced. Specifically, at time  $t$  of the  
 243 reverse process, instead of predicting  $p_\theta(X_{t-1} | X_t)$ , we can predict  $p_\theta(X_{t'} | X_t)$  for any  $0 \leq t' < t$ .  
 244 Given a sequence of timesteps  $0 = t_0 < \dots < t_k = T$ , we can now model the reverse process as

$$p_\theta(X_{t_0}, \dots, X_{t_k} | \mathcal{X}) = p(X_T | \mathcal{X}) \prod_{i=1}^k p_\theta(X_{t_{i-1}} | X_{t_i}). \quad (11)$$

245 To align with the literature of diffusion models, we call the list  $[t_0, \dots, t_k]$  the *denoising schedule*.  
 246 After incorporating the denoising schedule in Eq. (10), we obtain the loss function:

$$\mathcal{L}(\theta) = \mathbb{E}_{p_{\text{data}}(X_0, \mathcal{X})} \mathbb{E}_{q(X_{1:T} | X_0, \mathcal{X})} \left[ -\log p(X_T | \mathcal{X}) - \sum_{i=1}^k \log \frac{p_\theta(X_{t_{i-1}} | X_{t_i})}{q(X_{t_i} | X_{t_{i-1}})} \right]. \quad (12)$$

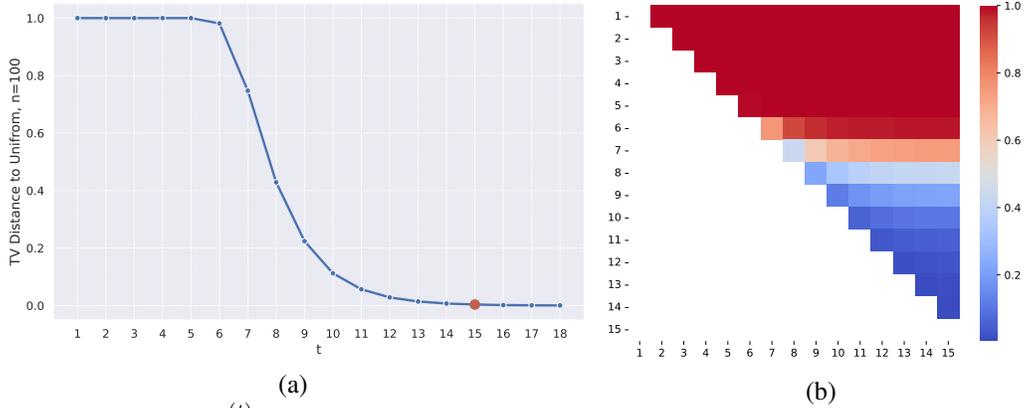


Figure 2: (a)  $D_{\text{TV}}(q_{\text{RS}}^{(t)}, u)$  computed using Eq. (14). We choose  $T = 15$  (red dot) based on the threshold 0.005. (b) A heatmap for  $D_{\text{TV}}(q_{\text{RS}}^{(t)}, q_{\text{RS}}^{(t')})$  for  $n = 100$  and  $1 \leq t < t' \leq 15$ , computed using Eq. (13). Rows are  $t$  and columns are  $t'$ . We choose the denoising schedule  $[0, 8, 10, 15]$ .

247 Note that although we may not have the analytical form of  $q(X_{t_i}|X_{t_{i-1}})$ , we can draw samples  
 248 from it. Merging is feasible if the support of  $p_{\theta}(X_{t_{i-1}}|X_{t_i})$  is equal or larger than the support  
 249 of  $q(X_{t_i}|X_{t_{i-1}})$ ; otherwise, the inverse of some forward permutations would be almost surely  
 250 unrecoverable. Therefore, we can implement a non-trivial denoising schedule (*i.e.*,  $k < T$ ), when  
 251  $p_{\theta}(\sigma'_t|X_t)$  follows the PL or GPL distribution, as they have whole  $S_n$  as their support. However,  
 252 merging is not possible for inverse shuffling methods, as their support is smaller than that of the  
 253 corresponding multi-step forward shuffling. To design a successful denoising schedule, we first  
 254 describe the intuitive principles and then provide some theoretical insights. 1) The length of forward  
 255 diffusion  $T$  should be minimal so long as the forward process approaches the uniform distribution. 2)  
 256 If distributions of  $X_t$  and  $X_{t+1}$  are similar, we should merge these two steps. Otherwise, we should  
 257 not merge them, as it would make the learning problem harder.

258 To quantify the similarity between distributions shown in 1) and 2), the TV distance is commonly  
 259 used in the literature. In particular, we can measure  $D_{\text{TV}}(q^{(t)}, q^{(t')})$  for  $t \neq t'$  and  $D_{\text{TV}}(q^{(t)}, u)$ ,  
 260 where  $q^{(t)}$  is the distribution at time  $t$  in the forward process and  $u$  is the uniform distribution. For  
 261 riffle shuffles, the total variation distance can be computed exactly. Specifically, we first introduce  
 262 the *Eulerian Numbers*  $A_{n,r}$  [32], *i.e.*, the number of permutations in  $S_n$  that have exactly  $r$  rising  
 263 sequences where  $1 \leq r \leq n$ .  $A_{n,r}$  can be computed using the following recursive formula  $A_{n,r} =$   
 264  $rA_{n-1,r} + (n-r+1)A_{n-1,r-1}$  where  $A_{1,1} = 1$ . We then have the following result.

265 **Proposition 2.** *Let  $t \neq t'$  be positive integers. Then*

$$D_{\text{TV}}(q_{\text{RS}}^{(t)}, q_{\text{RS}}^{(t')}) = \frac{1}{2} \sum_{r=1}^n A_{n,r} \left| \frac{1}{2^{tn}} \binom{n+2^t-r}{n} - \frac{1}{2^{t'n}} \binom{n+2^{t'}-r}{n} \right|, \quad (13)$$

266 *and*

$$D_{\text{TV}}(q_{\text{RS}}^{(t)}, u) = \frac{1}{2} \sum_{r=1}^n A_{n,r} \left| \frac{1}{2^{tn}} \binom{n+2^t-r}{n} - \frac{1}{n!} \right|. \quad (14)$$

267 Note that Eq. (14) was originally given in [19]. We restate it here for completeness. Once the  
 268 Eulerian numbers are precomputed, the TV distances can be computed in  $O(n)$  time instead of  $O(n!)$ .  
 269 Through extensive experiments, we have the following empirical observation. For the principle 1),  
 270 choosing  $T$  so that  $D_{\text{TV}}(q_{\text{RS}}^{(T)}, u) \approx 0.005$  yields good results. For the principle 2), a denoising  
 271 schedule  $[t_0, \dots, t_k]$  with  $D_{\text{TV}}(q_{\text{RS}}^{(t_i)}, q_{\text{RS}}^{(t_{i+1})}) \approx 0.3$  for most  $i$  works well. We show an example on  
 272 sorting  $n = 100$  four-digit MNIST images in Fig. 2.

### 274 3.5 Reverse Process Decoding

275 We now discuss how to decode predictions from the reverse process at test time. In practice, one is  
 276 often interested in the most probable state or a few states with high probabilities under  $p_{\theta}(X_0|\mathcal{X})$ .  
 277 However, since we can only draw samples from  $p_{\theta}(X_0|\mathcal{X})$  via running the reverse process, exact  
 278 decoding is intractable. The simplest approximated method is greedy search, *i.e.*, successively finding  
 279 the mode or an approximated mode of  $p_{\theta}(X_{t_{i-1}}|X_{t_i})$ . Another approach is beam search, which

Method	Metrics	Noisy MNIST					CIFAR-10		
		2 × 2	3 × 3	4 × 4	5 × 5	6 × 6	2 × 2	3 × 3	4 × 4
Gumbel-Sinkhorn Network [29]	Kendall-Tau ↑	0.9984	0.6908	0.3578	0.2430	0.1755	0.8378	0.5044	<b>0.4016</b>
	Accuracy (%)	99.81	44.65	00.86	0.00	0.00	76.54	6.07	0.21
	Correct (%)	99.91	80.20	49.51	26.94	14.91	86.10	43.59	25.31
	RMSE ↓	<b>0.0022</b>	0.1704	0.4572	0.8915	1.0570	0.3749	0.9590	1.0960
	MAE ↓	0.0003	0.0233	0.1005	0.3239	0.4515	0.1368	0.5320	0.6873
DiffSort [34]	Kendall-Tau ↑	0.9931	0.3054	0.0374	0.0176	0.0095	0.6463	0.1460	0.0490
	Accuracy (%)	99.02	5.56	0.00	0.00	0.00	59.18	0.96	0.00
	Correct (%)	99.50	42.25	10.77	6.39	3.77	75.48	27.87	12.27
	RMSE ↓	0.0689	1.0746	1.3290	1.4883	1.5478	0.7389	1.2691	1.3876
	MAE ↓	0.0030	0.4283	0.6531	0.8204	0.8899	0.2800	0.8123	0.9737
Error-free DiffSort [20]	Kendall-Tau ↑	0.9899	0.2014	0.0100	0.0034	-0.0021	0.6604	0.1362	0.0318
	Accuracy (%)	98.62	0.82	0.00	0.00	0.00	60.96	0.68	0.00
	Correct (%)	99.28	32.65	7.40	4.39	2.50	75.99	26.75	10.33
	RMSE ↓	0.0814	1.1764	1.3579	1.5084	1.5606	0.7295	1.2820	1.4095
	MAE ↓	0.0041	0.5124	0.6818	0.8424	0.9041	0.2731	0.8260	0.9990
Symmetric Diffusers (Ours)	Kendall-Tau ↑	<b>0.9992</b>	<b>0.8126</b>	<b>0.4859</b>	<b>0.2853</b>	<b>0.1208</b>	<b>0.9023</b>	<b>0.8363</b>	0.2518
	Accuracy (%)	<b>99.88</b>	<b>57.38</b>	<b>1.38</b>	0.00	0.00	<b>90.15</b>	<b>70.94</b>	<b>0.64</b>
	Correct (%)	<b>99.94</b>	<b>86.16</b>	<b>58.51</b>	<b>37.91</b>	<b>18.54</b>	<b>92.99</b>	<b>86.84</b>	<b>34.69</b>
	RMSE ↓	0.0026	<b>0.0241</b>	<b>0.1002</b>	<b>0.2926</b>	<b>0.4350</b>	<b>0.3248</b>	<b>0.3892</b>	<b>0.8953</b>
	MAE ↓	<b>0.0001</b>	<b>0.0022</b>	<b>0.0130</b>	<b>0.0749</b>	<b>0.1587</b>	<b>0.0651</b>	<b>0.0977</b>	<b>0.5044</b>

Table 1: Results (averaged over 5 runs) on solving the jigsaw puzzle on Noisy MNIST and CIFAR10.

Method	Metrics	Sequence Length							
		3	5	7	9	15	32	52	100
DiffSort [34]	Kendall-Tau ↑	0.930	0.898	0.864	0.801	0.638	0.535	0.341	0.166
	Accuracy (%)	93.8	83.9	71.5	52.2	10.3	0.2	0.0	0.0
	Correct (%)	95.8	92.9	90.1	85.2	82.3	61.8	42.8	23.2
Error-free DiffSort [20]	Kendall-Tau ↑	0.974	<b>0.967</b>	<b>0.962</b>	<b>0.952</b>	<b>0.938</b>	<b>0.879</b>	0.170	0.140
	Accuracy (%)	97.7	95.3	92.9	89.6	<b>83.1</b>	<b>57.1</b>	0.0	0.0
	Correct (%)	98.4	<b>97.7</b>	<b>97.2</b>	<b>96.3</b>	<b>95.1</b>	<b>90.1</b>	24.2	20.1
Symmetric Diffusers (Ours)	Kendall-Tau ↑	<b>0.976</b>	<b>0.967</b>	0.959	0.950	0.932	0.858	<b>0.786</b>	<b>0.641</b>
	Accuracy (%)	<b>98.0</b>	<b>95.5</b>	<b>92.9</b>	<b>90.0</b>	82.6	55.1	<b>27.4</b>	<b>4.5</b>
	Correct (%)	<b>98.5</b>	97.6	96.8	96.1	94.5	88.3	<b>82.1</b>	<b>69.3</b>

Table 2: Results (averaged over 5 runs) on the four-digit MNIST sorting benchmark.

280 maintains a dynamic buffer of  $k$  candidates with highest probabilities. Nevertheless, for one-step  
 281 reverse transitions like the GPL distribution, even finding the mode is intractable. To address this, we  
 282 employ a hierarchical beam search that performs an inner beam search within  $n^2$  scores at each step  
 283 of the outer beam search. Further details are provided in Appendix C.

## 284 4 Experiments

285 We now demonstrate the general applicability and effectiveness of our model through a variety of  
 286 experiments, including sorting 4-digit MNIST numbers, solving jigsaw puzzles, and addressing  
 287 traveling salesman problems. Additional details are provided in the appendix due to space constraints.

### 288 4.1 Sorting 4-digit MNIST Images

289 We first evaluate our SymmetricDiffusers on the four-digit MNIST sorting benchmark, a well-  
 290 established testbed for differentiable sorting [5, 8, 13, 20, 33, 34]. Each four-digit image in this  
 291 benchmark is obtained by concatenating 4 individual images from MNIST. For evaluation, we  
 292 employ several metrics to compare methods, including Kendall-Tau coefficient (measuring the  
 293 correlation between rankings), accuracy (percentage of images perfectly reassembled), and correctness  
 294 (percentage of pieces that are correctly placed).

295 **Ablation Study.** We conduct an ablation study to verify our design choices for reverse transition and  
 296 decoding strategies. As shown in Table 3, combining PL with either beam search (BS) or greedy  
 297 search yields good results in terms of Kendall-Tau and correctness metrics. In contrast, the IRS  
 298 (inverse riffle shuffle) method, along with greedy search, performs poorly across all metrics, showing  
 299 the limitations of IRS in handling complicated sorting tasks. Finally, combining GPL and BS achieves  
 300 the best accuracy in correctly sorting the entire sequence of images. Given that accuracy is the most

	GPL + BS	GPL + Greedy	PL + Greedy	PL + BS	IRS + Greedy
Kendall-Tau $\uparrow$	0.786	<b>0.799</b>	<b>0.799</b>	0.797	0.390
Accuracy (%)	<b>27.4</b>	24.4	26.4	26.4	0.6
Correct (%)	82.1	81.6	<b>83.3</b>	83.1	44.6

Table 3: Ablation study on transitions of reverse diffusion and decoding strategies. Results are averaged over three runs on sorting 52 four-digit MNIST images.

Method	OR Solvers				Learning-Based Models		
	Gurobi [14]	Concorde [1]	LKH-3 [15]	2-Opt [25]	GCN* [18]	DIFUSCO* [43]	Ours
Tour Length $\downarrow$	<b>3.842</b>	3.843	<b>3.842</b>	4.020	3.850	3.883	<b>3.849</b>
Optimality Gap (%) $\downarrow$	0.00	0.00	0.00	4.64	0.21	1.07	0.18

Table 4: Results on TSP-20. \* means we remove the post-processing heuristics for a fair comparison.

301 challenging metric to improve, we selecte GPL and BS for all remaining experiments. More ablation  
302 study (*e.g.*, denoising schedule) is provided in Appendix E.2.

303 **Full Results.** From Table 2, we can see that Error-free DiffSort achieves the best performance in  
304 sorting sequences with lengths up to 32. However, its performances drop significantly with long  
305 sequences (*e.g.*, length of 52 or 100). Meanwhile, DiffSort performs the worse due to the error  
306 accumulation of its soft differentiable swap function [20, 33]. In contrast, our method is on par with  
307 Error-free DiffSort in sorting short sequences and significantly outperforms others on long sequences.

## 308 4.2 Jigsaw Puzzle

309 We then explore image reassembly from segmented "jigsaw" puzzles [29, 31, 39]. We evaluate the  
310 performance using the MNIST and the CIFAR10 datasets, which comprises puzzles of up to  $6 \times 6$  and  
311  $4 \times 4$  pieces respectively. We add slight noise to pieces from the MNIST dataset to ensure background  
312 pieces are distinctive. To evaluate our models, we use Kendall-Tau coefficient, accuracy, correctness,  
313 RMSE (root mean square error of reassembled images), and MAE (mean absolute error) as metrics.

314 Table 1 presents results comparing our method with the Gumbel-Sinkhorn Network[29], DiffSort  
315 [34], and Error-free DiffSort [20]. DiffSort and Error-free DiffSort are primarily designed for sorting  
316 high-dimensional ordinal data which have clearly different patterns. Since jigsaw puzzles on MNIST  
317 and CIFAR10 contain pieces that are visually similar, these methods do not perform well. The  
318 Gumbel-Sinkhorn performs better for tasks involving fewer than  $4 \times 4$  pieces. In more challenging  
319 scenarios (*e.g.*,  $5 \times 5$  and  $6 \times 6$ ), our method significantly outperforms all competitors.

## 320 4.3 The Travelling Salesman Problem

321 At last, we explore the travelling salesman problem (TSP) to demonstrate the general applicability of  
322 our model. TSPs are classical NP-complete combinatorial optimization problems which are solved  
323 using integer programming or heuristic solvers [2, 12]. There exists a vast literature on learning-based  
324 models to solve TSPs [22, 23, 18, 17, 6, 24, 10, 36, 21, 43, 30]. They often focus on the Euclidean  
325 TSPs, which are formulated as follows. Let  $V = \{v_1, \dots, v_n\}$  be points in  $\mathbb{R}^2$ . We need to find some  
326  $\sigma \in S_n$  such that  $\sum_{i=1}^n \|v_{\sigma(i)} - v_{\sigma(i+1)}\|_2$  is minimized, where we let  $\sigma(n+1) := \sigma(1)$ . Further  
327 experimental details are provided in Appendix B.

328 We compare with operations research (OR) solvers and other learning based approaches on TSP  
329 instances with 20 nodes. The metrics are the total tour length and the optimality gap. Given the ground  
330 truth (GT) length produced by the best OR solver, the optimality gap is given by  $(\text{predicted length} -$   
331  $(\text{GT length})) / (\text{GT length})$ . As shown in Table 4, SymmetricDiffusers achieves comparable results  
332 with both OR solvers and the state-of-the-art learning-based methods.

## 333 5 Conclusion

334 In this paper, we introduce a novel discrete diffusion model over finite symmetric groups. We identify  
335 the riffle shuffle as an effective forward transition and provide empirical rules for selecting the  
336 diffusion length. Additionally, we propose a generalized PL distribution for the reverse transition,  
337 which is provably more expressive than the PL distribution. We further introduce a theoretically  
338 grounded "denoising schedule" to improve sampling and learning efficiency. Extensive experiments  
339 verify the effectiveness of our proposed model. In the future, we are interested in generalizing our  
340 model to general finite groups and exploring diffusion models on Lie groups.

341 **References**

- 342 [1] David Applegate, Robert Bixby, Vasek Chvatal, and William Cook. Concorde tsp solver, 2006.
- 343 [2] Sanjeev Arora. Polynomial time approximation schemes for euclidean traveling salesman and  
344 other geometric problems. *J. ACM*, 45(5):753–782, Sep 1998.
- 345 [3] Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg.  
346 Structured denoising diffusion models in discrete state-spaces, 2023.
- 347 [4] Dave Bayer and Persi Diaconis. Trailing the dovetail shuffle to its lair. *The Annals of Applied  
348 Probability*, 2(2):294 – 313, 1992.
- 349 [5] Mathieu Blondel, Olivier Teboul, Quentin Berthet, and Josip Djolonga. Fast differentiable  
350 sorting and ranking. In *International Conference on Machine Learning*, pages 950–959. PMLR,  
351 2020.
- 352 [6] Xavier Bresson and Thomas Laurent. The transformer network for the traveling salesman  
353 problem, 2021.
- 354 [7] Zhe Cao, Tao Qin, Tie-Yan Liu, Ming-Feng Tsai, and Hang Li. Learning to rank: from pairwise  
355 approach to listwise approach. In *Proceedings of the 24th International Conference on Machine  
356 Learning*, ICML '07, pages 129–136, New York, NY, USA, 2007. Association for Computing  
357 Machinery.
- 358 [8] Marco Cuturi, Olivier Teboul, and Jean-Philippe Vert. Differentiable ranking and sorting using  
359 optimal transport. *Advances in neural information processing systems*, 32, 2019.
- 360 [9] Persi Diaconis. Group representations in probability and statistics. *Lecture notes-monograph  
361 series*, 11:i–192, 1988.
- 362 [10] Zhang-Hua Fu, Kai-Bin Qiu, and Hongyuan Zha. Generalize a small pre-trained model to  
363 arbitrarily large tsp instances, 2021.
- 364 [11] E. N. Gilbert. Theory of shuffling. *Bell Telephone Laboratories Memorandum*, 1955.
- 365 [12] Teofilo F. Gonzalez. *Handbook of Approximation Algorithms and Metaheuristics (Chapman &  
366 Hall/Crc Computer & Information Science Series)*. Chapman & Hall/CRC, 2007.
- 367 [13] Aditya Grover, Eric Wang, Aaron Zweig, and Stefano Ermon. Stochastic optimization of sorting  
368 networks via continuous relaxations. In *International Conference on Learning Representations*,  
369 2018.
- 370 [14] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2023.
- 371 [15] Keld Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling  
372 salesman and vehicle routing problems, Dec 2017.
- 373 [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
- 374 [17] Chaitanya K Joshi, Quentin Cappart, Louis-Martin Rousseau, and Thomas Laurent. Learning  
375 tsp requires rethinking generalization. In *International Conference on Principles and Practice  
376 of Constraint Programming*, 2021.
- 377 [18] Chaitanya K. Joshi, Thomas Laurent, and Xavier Bresson. An efficient graph convolutional  
378 network technique for the travelling salesman problem, 2019.
- 379 [19] Shihan Kanungo. Mixing time estimates for the riffle shuffle. *Euler Circle*, 2020.
- 380 [20] Jungtaek Kim, Jeongbeen Yoon, and Minsu Cho. Generalized neural sorting networks with error-  
381 free differentiable swap functions. In *International Conference on Learning Representations  
382 (ICLR)*, 2024.
- 383 [21] Minsu Kim, Junyoung Park, and Jinkyoo Park. Sym-nco: Leveraging symmetry for neural  
384 combinatorial optimization, 2023.

- 385 [22] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional  
386 networks, 2017.
- 387 [23] Wouter Kool, Herke van Hoof, and Max Welling. Attention, learn to solve routing problems!,  
388 2019.
- 389 [24] Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai  
390 Min. Pomo: Policy optimization with multiple optima for reinforcement learning, 2021.
- 391 [25] Shen Lin and Brian W Kernighan. An effective heuristic algorithm for the travelling-salesman  
392 problem. *Operations research*, 21(2):498–516, 1973.
- 393 [26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- 394 [27] R. D. Luce. *Individual Choice Behavior*. John Wiley, 1959.
- 395 [28] Colin L Mallows. Non-null ranking models. i. *Biometrika*, 44(1/2):114–130, 1957.
- 396 [29] Gonzalo Mena, David Belanger, Scott Linderman, and Jasper Snoek. Learning latent permuta-  
397 tions with gumbel-sinkhorn networks. In *International Conference on Learning Representations*,  
398 2018.
- 399 [30] Yimeng Min, Yiwei Bai, and Carla P. Gomes. Unsupervised learning for solving the travelling  
400 salesman problem, 2024.
- 401 [31] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving  
402 jigsaw puzzles. In *European conference on computer vision*, pages 69–84. Springer, 2016.
- 403 [32] OEIS Foundation Inc. The eulerian numbers, entry a008292 in the On-Line Encyclopedia of  
404 Integer Sequences, 2024. Published electronically at <http://oeis.org/A008292>.
- 405 [33] Felix Petersen, Christian Borgelt, Hilde Kuehne, and Oliver Deussen. Differentiable sorting  
406 networks for scalable sorting and ranking supervision. In *International conference on machine  
407 learning*, pages 8546–8555. PMLR, 2021.
- 408 [34] Felix Petersen, Christian Borgelt, Hilde Kuehne, and Oliver Deussen. Monotonic differentiable  
409 sorting networks. In *International Conference on Learning Representations (ICLR)*, 2022.
- 410 [35] R. L. Plackett. The analysis of permutations. *Applied Statistics*, 24(2):193 – 202, 1975.
- 411 [36] Ruizhong Qiu, Zhiqing Sun, and Yiming Yang. Dimes: A differentiable meta solver for  
412 combinatorial optimization problems, 2022.
- 413 [37] J. Reeds. Theory of shuffling. *Unpublished Manuscript*, 1981.
- 414 [38] Laurent Saloff-Coste. *Random Walks on Finite Groups*, pages 263–346. Springer Berlin  
415 Heidelberg, Berlin, Heidelberg, 2004.
- 416 [39] Rodrigo Santa Cruz, Basura Fernando, Anoop Cherian, and Stephen Gould. Deeppermnet:  
417 Visual permutation learning. In *Proceedings of the IEEE Conference on Computer Vision and  
418 Pattern Recognition*, pages 3949–3957, 2017.
- 419 [40] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsu-  
420 pervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei,  
421 editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of  
422 *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 07–09 Jul 2015.  
423 PMLR.
- 424 [41] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data  
425 distribution, 2020.
- 426 [42] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and  
427 Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- 428 [43] Zhiqing Sun and Yiming Yang. Difusco: Graph-based diffusion solvers for combinatorial  
429 optimization, 2023.

- 430 [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
431 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- 432 [45] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal  
433 Frossard. Digress: Discrete denoising diffusion for graph generation, 2023.
- 434 [46] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforce-  
435 ment learning. *Machine Learning*, 8(3-4):229–256, 1992.

436 **A Additional Details of the GSR Riffle Shuffle Model**

437 There are many equivalent definitions of the GSR riffle shuffle. Here we also introduce the *Geometric*  
 438 *Description* [4], which is easy to implement (and is how we implement riffle shuffles in our experi-  
 439 ments). We first sample  $n$  points in the unit interval  $[0, 1]$  uniformly and independently, and suppose  
 440 the points are labeled in order as  $x_1 < x_2 < \dots < x_n$ . Then, the permutation that sorts the points  
 441  $\{2x_1\}, \dots, \{2x_n\}$  follows the GSR distribution, where  $\{x\} := x - \lfloor x \rfloor$  is the fractional part of  $x$ .

442 **B Details of Our Network Architecture**

443 We now discuss how to use neural networks to produce the parameters of the distributions discussed  
 444 in Section 3.2.1 and 3.2.2. Fix time  $t$ , and suppose  $X_t = (\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_n^{(t)})^\top \in \mathbb{R}^{n \times d}$ . Let  $\text{encoder}_\theta$   
 445 be an object-specific encoder such that  $\text{encoder}_\theta(X_t) \in \mathbb{R}^{n \times d_{\text{model}}}$ . For example,  $\text{encoder}_\theta$  can be  
 446 a CNN if  $X_t$  is an image. Let

$$Y_t := \text{encoder}_\theta(X_t) + \text{time\_embd}(t) = (\mathbf{y}_1^{(t)}, \dots, \mathbf{y}_n^{(t)})^\top \in \mathbb{R}^{n \times d_{\text{model}}}, \quad (15)$$

447 where  $\text{time\_embd}$  is the sinusoidal time embedding. Then, we would like to feed the embeddings into  
 448 a Transformer encoder [44]. Let  $\text{transformer\_encoder}_\theta$  be the encoder part of the Transformer  
 449 architecture. However, each of the distributions we discussed previously has different number of  
 450 parameters, so we will have to discuss them separately.

451 **Inverse Transposition.** For Inverse Transposition, we have  $n + 1$  parameters. To obtain  $n + 1$   
 452 tokens from  $\text{transformer\_encoder}_\theta$ , we append a dummy token of 0's to  $Y_t$ . Then we input  
 453  $(\mathbf{y}_1^{(t)}, \dots, \mathbf{y}_n^{(t)}, 0)^\top$  into  $\text{transformer\_encoder}_\theta$  to obtain  $Z \in \mathbb{R}^{(n+1) \times d_{\text{model}}}$ . Finally, we apply  
 454 an MLP to obtain  $(s_1, \dots, s_n, k) \in \mathbb{R}^{n+1}$ .

455 **Inverse Insertion, Inverse Riffle Shuffle, PL Distribution.** These three distributions all require  
 456 exactly  $n$  parameters, so we can directly feed  $Y_t$  into  $\text{transformer\_encoder}_\theta$ . Let the output  
 457 of  $\text{transformer\_encoder}_\theta$  be  $Z \in \mathbb{R}^{n \times d_{\text{model}}}$ , where we then apply an MLP to obtain the scores  
 458  $\mathbf{s}_\theta \in \mathbb{R}^n$ .

459 **The GPL Distribution.** The GPL distribution requires  $n^2$  parameters. We first append  $n$  dummy  
 460 tokens of 0's to  $Y_t$ , with the intent that the  $j^{\text{th}}$  dummy token would learn information about the  
 461  $j^{\text{th}}$  column of the GPL parameter matrix, which represents where the  $j^{\text{th}}$  component should be  
 462 placed. We then pass  $(\mathbf{y}_1^{(t)}, \dots, \mathbf{y}_n^{(t)}, 0, \dots, 0)^\top \in \mathbb{R}^{2n \times d_{\text{model}}}$  to  $\text{transformer\_encoder}_\theta$ . When  
 463 computing attention, we further apply a  $2n \times 2n$  attention mask

$$M := \begin{bmatrix} 0 & A \\ 0 & B \end{bmatrix}, \text{ where } A \text{ is an } n \times n \text{ matrix of } -\infty, B = \begin{bmatrix} -\infty & -\infty & \dots & -\infty \\ 0 & -\infty & \dots & -\infty \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & -\infty \end{bmatrix} \text{ is } n \times n.$$

464 The reason for having  $B$  as an upper triangular matrix of  $-\infty$  is that information about the  $j^{\text{th}}$   
 465 component should only require information from the previous components. Let

$$\text{transformer\_encoder}_\theta(Y_t, M) = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix},$$

466 where  $Z_1, Z_2 \in \mathbb{R}^{n \times d_{\text{model}}}$ . Finally, we obtain the GPL parameter matrix as  $S_\theta = Z_1 Z_2^\top \in \mathbb{R}^{n \times n}$ .

467 For hyperparameters, we refer the readers to Appendix E.4.

468 **C Additional Details of Decoding**

469 **Greedy Search.** At each timestep  $t_i$  in the denoising schedule, we can greedily obtain or approx-  
 470 imate the mode of  $p_\theta(X_{t_{i-1}}|X_{t_i})$ . We can then use the (approximated) mode  $X_{t_{i-1}}$  for the next  
 471 timestep  $p_\theta(X_{t_{i-2}}|X_{t_{i-1}})$ . Note that the final  $X_0$  obtained using such a greedy heuristic may not  
 472 necessarily be the mode of  $p_\theta(X_0|\mathcal{X})$ .

473 **Beam Search.** We can use beam search to improve the greedy approach. The basic idea is that,  
 474 at each timestep  $t_i$  in the denoising schedule, we compute or approximate the top- $k$ -most-probable  
 475 results from  $p_\theta(X_{t_{i-1}}|X_{t_i})$ . For each of the top- $k$  results, we sample top- $k$  from  $p_\theta(X_{t_{i-2}}|X_{t_{i-1}})$ .  
 476 Now we have  $k^2$  candidates for  $X_{t_{i-2}}$ , and we only keep the top  $k$  of the  $k^2$  candidates.

477 However, it is not easy to obtain the top- $k$ -most-probable results for some of the distributions. Here  
 478 we provide an algorithm to approximate top- $k$  of the PL and the GPL distribution. Since the PL  
 479 distribution is a strict subset of the GPL distribution, it suffices to only consider the GPL distribution  
 480 with parameter matrix  $S$ . The algorithm for approximating top- $k$  of the GPL distribution is another  
 481 beam search. We first pick the  $k$  largest elements from the first row of  $S$ . For each of the  $k$  largest  
 482 elements, we pick  $k$  largest elements from the second row of  $S$ , excluding the corresponding element  
 483 picked in the first row. We now have  $k^2$  candidates for the first two elements of a permutation, and  
 484 we only keep the top- $k$ -most-probable candidates. We then continue in this manner.

## 485 D Proofs

486 **Proposition 1.** *The PL distribution cannot exactly represent a delta distribution. That is, there does*  
 487 *not exist an  $\mathbf{s}$  such that  $p_{\text{PL}} = \delta_\sigma$  for any  $\sigma \in S_n$ , where  $\delta_\sigma(\sigma) = 1$  and  $\delta_\sigma(\pi) = 0$  for all  $\pi \neq \sigma$ .*  
 488 *But the GPL distribution can represent a delta distribution exactly.*

489 *Proof.* Assume for a contradiction that there exists some  $\sigma \in S_n$  and  $\mathbf{s}$  such that  $\text{PL}_\mathbf{s} = \delta_\sigma$ . Then  
 490 we have

$$\prod_{i=1}^n \frac{\exp(s_{\sigma(i)})}{\sum_{j=i}^n \exp(s_{\sigma(j)})} = 1.$$

491 Since each of the term in the product is less than or equal to 1, we must have

$$\frac{\exp(s_{\sigma(i)})}{\sum_{j=i}^n \exp(s_{\sigma(j)})} = 1 \tag{16}$$

492 for all  $i \in [n]$ . In particular, we have

$$\frac{\exp(s_{\sigma(1)})}{\sum_{j=1}^n \exp(s_{\sigma(j)})} = 1,$$

493 which happens if and only if  $s_{\sigma(j)} = -\infty$  for all  $j \geq 2$ . But this contradicts (16).

494 We then show that the GPL distribution can represent a delta distribution exactly. To see this, we fix  
 495  $\sigma \in S_n$ . For all  $i \in [n]$ , we let  $s_{i,\sigma(i)} = 0$  and  $s_{i,j} = -\infty$  for all  $j \neq \sigma(i)$ . Then  $\text{GPL}_{(s_{ij})} = \delta_\sigma$ .  $\square$

496 **Proposition 2.** *Let  $t \neq t'$  be positive integers. Then*

$$D_{\text{TV}}(q_{\text{RS}}^{(t)}, q_{\text{RS}}^{(t')}) = \frac{1}{2} \sum_{r=1}^n A_{n,r} \left| \frac{1}{2^{tn}} \binom{n+2^t-r}{n} - \frac{1}{2^{t'n}} \binom{n+2^{t'}-r}{n} \right|, \tag{13}$$

497 and

$$D_{\text{TV}}(q_{\text{RS}}^{(t)}, u) = \frac{1}{2} \sum_{r=1}^n A_{n,r} \left| \frac{1}{2^{tn}} \binom{n+2^t-r}{n} - \frac{1}{n!} \right|. \tag{14}$$

498

499 *Proof.* Let  $\sigma \in S_n$ . It was shown in [4] that

$$q_{\text{RS}}^{(t)}(\sigma) = \frac{1}{2^{tn}} \cdot \binom{n+2^t-r}{n},$$

500 where  $r$  is the number of rising sequences of  $\sigma$ . Note that if two permutations have the same number  
 501 of rising sequences, then they have equal probability. Hence, we have

$$\begin{aligned} D_{\text{TV}}(q_{\text{RS}}^{(t)} - q_{\text{RS}}^{(t')}) &= \frac{1}{2} \sum_{\sigma \in S_n} |q_{\text{RS}}^{(t)}(\sigma) - q_{\text{RS}}^{(t')}(\sigma)| = \frac{1}{2} \sum_{r=1}^n A_{n,r} |q_{\text{RS}}^{(t)}(\sigma) - q_{\text{RS}}^{(t')}(\sigma)| \\ &= \frac{1}{2} \sum_{r=1}^n A_{n,r} \left| \frac{1}{2^{tn}} \binom{n+2^t-r}{n} - \frac{1}{2^{t'n}} \binom{n+2^{t'}-r}{n} \right|, \end{aligned}$$

502 as claimed. For (14), replace  $q_{\text{RS}}^{(t')}(\sigma)$  with  $u(\sigma) = \frac{1}{n!}$  in the above derivations.  $\square$

## 503 E Additional Details on Experiments

### 504 E.1 Datasets

505 **Jigsaw Puzzle.** We created the Noisy MNIST dataset by adding *i.i.d.* Gaussian noise with a mean  
506 of 0 and a standard deviation of 0.01 to each pixel of the MNIST images. No noise was added to the  
507 CIFAR-10 images. The noisy images are then saved as the Noisy MNIST dataset. During training,  
508 each image is divided into  $n \times n$  patches. A permutation is then sampled uniformly at random  
509 to shuffle these patches. The training set for Noisy MNIST comprises 60,000 images, while the  
510 CIFAR-10 training set contains 10,000 images. The Noisy MNIST test set, which is pre-shuffled, also  
511 includes 10,000 images. The CIFAR-10 test set, which shuffles images on the fly, contains 10,000  
512 images as well.

513 **Sort 4-Digit MNIST Numbers.** For each training epoch, we generate 60,000 sequences of 4-digit  
514 MNIST images, each of length  $n$ , constructed dynamically on the fly. These 4-digit MNIST numbers  
515 are created by concatenating four MNIST images, each selected uniformly at random from the entire  
516 MNIST dataset, which consists of 60,000 images. For testing purposes, we similarly generate 10,000  
517 sequences of  $n$  4-digit MNIST numbers on the fly.

518 **TSP.** We take the TSP-20 dataset from [17]<sup>1</sup>. The train set consists of 1,512,000 graphs with 20  
519 nodes, where each node is an *i.i.d.* sample from the unit square  $[0, 1]^2$ . The labels are optimal TSP  
520 tours provided by the Concorde solver [1]. The test set consists of 1,280 graphs with 20 nodes, with  
521 ground truth tour generated by the Concorde solver as well.

### 522 E.2 Ablation Studies

523 **Choices for Reverse Transition and Decoding Strategies.** As demonstrated in Table 5, we have  
524 explored various combinations of forward and inverse shuffling methods across tasks involving  
525 different sequence lengths. Both GPL and PL consistently excel in all experimental scenarios,  
526 highlighting their robustness and effectiveness. It is important to note that strategies such as random  
527 transposition and random insertion paired with their respective inverse operations, are less suitable  
528 for tasks with longer sequences. This limitation is attributed to the prolonged mixing times required  
529 by these two shuffling methods, a challenge that is thoroughly discussed in Section 3.1.2.

530 **Denosing Schedule.** We also conduct an ablation study on how we should merge reverse steps. As  
531 shown in Table 6, the choice of the denosing schedule can significantly affect the final performance.  
532 In particular, for  $n = 100$  on the Sort 4-Digit MNIST Numbers task, the fact that  $[0, 15]$  has 0  
533 accuracy justifies our motivation to use diffusion to break down learning into smaller steps. The  
534 result we get also matches with our proposed heuristic in Section 3.4.

### 535 E.3 Latent Loss in Jigsaw Puzzle

536 In the original setup of the Jigsaw Puzzle experiment using the Gumbel-Sinkhorn network [29],  
537 the permutations are latent. That is, the loss function in Gumbel-Sinkhorn is a pixel-level MSE  
538 loss and does not use the ground truth permutation label. However, our loss function (12) actually  
539 (implicitly) uses the ground truth permutation that maps the shuffled image patches to their original  
540 order. Therefore, for fair comparison with the Gumbel-Sinkhorn network in the Jigsaw Puzzle  
541 experiment, we modify our loss function so that it does not use the ground truth permutation. Recall  
542 from Section 3.2 that we defined

$$p_{\theta}(X_{t-1}|X_t) = \sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) p_{\theta}(\sigma'_t|X_t). \quad (17)$$

543 In our original setup, we defined  $p(X_{t-1}|X_t, \sigma'_t)$  as a delta distribution  $\delta(X_{t-1} = Q_{\sigma'_t} X_t)$ , but this  
544 would require that we know the permutation that turns  $X_{t-1}$  to  $X_t$ , which is part of the ground truth.  
545 So instead, we parameterize  $p(X_{t-1}|X_t, \sigma'_t)$  as a Gaussian distribution  $\mathcal{N}(X_{t-1}|Q_{\sigma'_t} X_t, I)$ . At the  
546 same time, we note that to find the gradient of (12), it suffices to find the gradient of the log of (17).

<sup>1</sup><https://github.com/chaitjo/learning-tsp?tab=readme-ov-file>

		Sequence Length		
		9	32	52
RS (forward) + GPL (reverse) + greedy	Denoising Schedule	[0, 3, 5, 9]	[0, 5, 7, 12]	[0, 5, 6, 7, 10, 13]
	Kendall-Tau $\uparrow$	0.948	0.857	0.779
	Accuracy (%)	89.4	54.8	24.4
	Correct (%)	95.9	88.1	81.6
RS (forward) + PL (reverse) + greedy	Denoising Schedule	[0, 3, 5, 9]	[0, 5, 7, 12]	[0, 5, 6, 7, 10, 13]
	Kendall-Tau $\uparrow$	0.953	0.867	0.799
	Accuracy (%)	90.9	56.4	26.4
	Correct (%)	96.4	89.0	83.3
RS (forward) + PL (reverse) + beam search	Denoising Schedule	[0, 3, 5, 9]	[0, 5, 7, 12]	[0, 5, 6, 7, 10, 13]
	Kendall-Tau $\uparrow$	0.955	0.869	0.797
	Accuracy (%)	91.1	57.2	26.4
	Correct (%)	96.5	89.2	83.1
RS (forward) + IRS (reverse) + greedy	$T$	9	12	13
	Kendall-Tau $\uparrow$	0.947	0.794	0.390
	Accuracy (%)	88.6	24.4	0.6
	Correct (%)	95.9	82.5	44.6
RT (forward) + IT (reverse) + greedy	$T$ (using approx. $\frac{n}{2} \log n$ )	15	55	105
	Kendall-Tau $\uparrow$	0.490		
	Accuracy (%)	18.0	Out of Memory	
	Correct (%)	59.5		
RI (forward) + II (reverse) + greedy	$T$ (using approx. $n \log n$ )	25	110	205
	Kendall-Tau $\uparrow$	0.954		
	Accuracy (%)	91.1	Out of Memory	
	Correct (%)	96.4		

Table 5: More results on sorting the 4-digit MNIST dataset using different combinations of forward process methods and reverse process methods. Results averaged over 3 runs with different seeds. RS: riffle shuffle; GPL: generalized Plackett-Luce; IRS: inverse riffle shuffle; RT: random transposition; IT: inverse transposition; RI: random insertion; II: inverse insertion.

Denoising Schedule	[0, 15]	[0, 8, 9, 15]	[0, 7, 8, 9, 15]	[0, 7, 8, 10, 15]	[0, 8, 10, 15]
Kendall-Tau $\uparrow$	0.000	0.316	0.000	0.000	<b>0.646</b>
Accuracy (%)	0.0	0.0	0.0	0.0	<b>4.5</b>
Correct (%)	1.0	39.6	1.0	1.0	<b>69.8</b>

Table 6: Results of sorting 100 4-digit MNIST images using various denoising schedules with the combination of RS, GPL and beam search consistently applied.

547 We use the REINFORCE trick [46] to find the gradient of  $\log p_\theta(X_{t-1}|X_t)$ , which gives us

$$\begin{aligned}
& \nabla_\theta \log p_\theta(X_{t-1}|X_t) \\
&= \frac{1}{\sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) p_\theta(\sigma'_t|X_t)} \cdot \sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) \nabla_\theta p_\theta(\sigma'_t|X_t) \\
&= \frac{1}{\sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) p_\theta(\sigma'_t|X_t)} \cdot \sum_{\sigma'_t \in \mathcal{T}} p(X_{t-1}|X_t, \sigma'_t) p_\theta(\sigma'_t|X_t) (\nabla_\theta \log p_\theta(\sigma_t|X_t)) \\
&= \frac{\mathbb{E}_{p_\theta(\sigma_t|X_t)} \left[ p(X_{t-1}|X_t, \sigma'_t) \nabla_\theta \log p_\theta(\sigma_t|X_t) \right]}{\mathbb{E}_{p_\theta(\sigma_t|X_t)} \left[ p(X_{t-1}|X_t, \sigma'_t) \right]} \\
&\approx \sum_{n=1}^N \frac{p(X_{t-1}|X_t, \sigma_t^{(n)})}{\sum_{m=1}^N p(X_{t-1}|X_t, \sigma_t^{(m)})} \cdot \nabla_\theta \log p_\theta(\sigma_t^{(n)}|X_t),
\end{aligned}$$

548 where we have used Monte-Carlo estimation in the last step, and  $\sigma_t^{(1)}, \dots, \sigma_t^{(N)} \sim p_\theta(\sigma_t|X_t)$ . We  
 549 further add an entropy regularization term  $-\lambda \cdot \mathbb{E}_{p_\theta(\sigma_t|X_t)} [\log p_\theta(\sigma_t|X_t)]$  to each of  $\log p_\theta(X_{t-1}|X_t)$ .  
 550 Using the same REINFORCE and Monte-Carlo trick, we obtain

$$\nabla_\theta \left( -\lambda \cdot \mathbb{E}_{p_\theta(\sigma_t|X_t)} [\log p_\theta(\sigma_t|X_t)] \right) \approx \sum_{n=1}^N -\lambda \log p_\theta \left( \sigma_t^{(n)}|X_t \right) \nabla_\theta \log p_\theta \left( \sigma_t^{(n)}|X_t \right),$$

551 where  $\sigma_t^{(1)}, \dots, \sigma_t^{(N)} \sim p_\theta(\sigma_t|X_t)$ . Therefore, we have

$$\begin{aligned} & \nabla_\theta \left( \log p_\theta(X_{t-1}|X_t) - \lambda \cdot \mathbb{E}_{p_\theta(\sigma_t|X_t)} [\log p_\theta(\sigma_t|X_t)] \right) \\ & \approx \sum_{n=1}^N \left( \underbrace{\frac{p \left( X_{t-1}|X_t, \sigma_t^{(n)} \right)}{\sum_{m=1}^N p \left( X_{t-1}|X_t, \sigma_t^{(m)} \right)}}_{\text{weight}} - \lambda \log p_\theta \left( \sigma_t^{(n)}|X_t \right) \right) \cdot \nabla_\theta \log p_\theta \left( \sigma_t^{(n)}|X_t \right), \quad (18) \end{aligned}$$

552 where  $\sigma_t^{(1)}, \dots, \sigma_t^{(N)} \sim p_\theta(\sigma_t|X_t)$ . We then substitute in

$$p \left( X_{t-1}|X_t, \sigma_t^{(n)} \right) = \mathcal{N} \left( X_{t-1}|Q_{\sigma_t^{(n)}} X_t, I \right)$$

553 for all  $n \in [N]$ . Finally, we also subtract the exponential moving average `weight` as a control variate  
 554 for variance reduction, where the exponential moving average is given by `ema`  $\leftarrow$  `ema_rate`  $\cdot$  `ema` +  
 555  $(1 - \text{ema\_rate}) \cdot \text{weight}$  for each gradient descent step.

#### 556 E.4 Training Details and Architecture Hyperparameters

557 **Hardware.** The Jigsaw Puzzle and Sort 4-Digit MNIST Numbers experiments are trained and  
 558 evaluated on the NVIDIA A40 GPU. The TSP experiments are trained and evaluated on the NVIDIA  
 559 A40 and A100 GPU.

560 **Jigsaw Puzzle.** For the Jigsaw Puzzle experiments, we use the AdamW optimizer [26] with weight  
 561 decay  $1e-2$ ,  $\varepsilon = 1e-9$ , and  $\beta = (0.9, 0.98)$ . We use the Noam learning rate scheduler given in [44]  
 562 with 51,600 warmup steps for Noisy MNIST and 46,000 steps for CIFAR-10. We train for 120  
 563 epochs with a batch size of 64. When computing the loss (12), we use Monte-Carlo estimation for the  
 564 expectation and sample 3 trajectories. For REINFORCE, we sampled 10 times for the Monte-Carlo  
 565 estimation in (18), and we used an entropy regularization rate  $\lambda = 0.05$  and an `ema_rate` of 0.995.  
 566 The neural network architecture and related hyperparameters are given in Table 7. The denoising  
 567 schedules, with riffle shuffles as the forward process and GPL as the reverse process, are give in Table  
 568 8. For beam search, we use a beam size of 200 when decoding from GPL, and we use a beam size of  
 569 20 when decoding along the diffusion denoising schedule.

Layer	Details
Convolution	Output channels 32, kernel size 3, padding 1, stride 1
Batch Normalization	—
ReLU	—
Max-pooling	Pooling 2
Fully-connected	Output dimension $(\text{dim\_after\_conv} + 128)/2$
ReLU	—
Fully-connected	Output dimension 128
Transformer encoder	7 layers, 8 heads, model dimension ( $d_{\text{model}}$ ) 128, feed-forward dimension 512, dropout 0.1

Table 7: Jigsaw puzzle neural network architecture and hyperparameters.

Number of patches per side	Denoising schedule
$2 \times 2$	[0, 2, 7]
$3 \times 3$	[0, 3, 5, 9]
$4 \times 4$	[0, 4, 6, 10]
$5 \times 5$	[0, 5, 7, 11]
$6 \times 6$	[0, 6, 8, 12]

Table 8: Denoising schedules for the Jigsaw Puzzle task, where we use riffle shuffle in the forward process and GPL in the reverse process.

570 **Sort 4-Digit MNIST Numbers.** For the task of sorting 4-digit MNIST numbers, we use the exact  
571 training and beam search setup as the Jigsaw Puzzle, except that we do not need to use REINFORCE.  
572 The neural network architecture is given in Table 9. The denoising schedules, with riffle shuffles as  
573 the forward process and GPL as the reverse process, are give in Table 10.

Layer	Details
Convolution	Output channels 32, kernel size 5, padding 2, stride 1
Batch Normalization	—
ReLU	—
Max-pooling	Pooling 2
Convolution	Output channels 64, kernel size 5, padding 2, stride 1
Batch Normalization	—
ReLU	—
Max-pooling	Pooling 2
Fully-connected	Output dimension $(\text{dim\_after\_conv} + 128)/2$
ReLU	—
Fully-connected	Output dimension 128
Transformer encoder	7 layers, 8 heads, model dimension ( $d_{\text{model}}$ ) 128, feed-forward dimension 512, dropout 0.1

Table 9: Sort 4-digit MNIST numbers neural network architecture and hyperparameters.

Sequence Length $n$	Denoising schedule
3	[0, 2, 7]
5	[0, 2, 8]
7	[0, 3, 8]
9	[0, 3, 5, 9]
15	[0, 4, 7, 10]
32	[0, 5, 7, 12]
52	[0, 5, 6, 7, 10, 13]
100	[0, 8, 10, 15]

Table 10: Denoising schedules for the Sort 4-Digit MNIST Numbers task, where we use riffle shuffle in the forward process and GPL in the reverse process.

574 **TSP.** For solving the TSP, we perform supervised learning to train our SymmetricDiffusers to solve  
575 the TSP. Let  $\sigma^*$  be an optimal permutation, and let  $X_0$  be the list of nodes ordered by  $\sigma^*$ . We note  
576 that any cyclic shift of  $X_0$  is also optimal. Thus, for simplicity and without loss of generality, we  
577 always assume  $\sigma^*(1) = 1$ . In the forward process of SymmetricDiffusers, we only shuffle the second  
578 to the  $n^{\text{th}}$  node (or component). In the reverse process, we mask certain parameters of the reverse  
579 distribution so that we will always sample a permutation with  $\sigma_t(1) = 1$ .

580 The architecture details are slightly different for TSP, since we need to input both node and edge  
 581 features into our network. Denote by  $X_t$  the ordered list of nodes at time  $t$ . We obtain  $Y_t \in \mathbb{R}^{n \times d_{\text{model}}}$   
 582 as in Eq. (15), where  $\text{encoder}_\theta$  is now a sinusoidal embedding of the 2D coordinates. Let  $D_t \in \mathbb{R}^{n \times n}$   
 583 be the matrix representing the pairwise distances of points in  $X_t$ , respecting the order in  $X_t$ . Let  
 584  $E_t \in \mathbb{R}^{\binom{n}{2}}$  be the flattened vector of the upper triangular part of  $D_t$ . We also apply sinusoidal  
 585 embedding to  $E_t$  and add  $\text{time\_embd}(t)$  to it. We call the result  $F_t \in \mathbb{R}^{\binom{n}{2} \times d_{\text{model}}}$ .

586 Now, instead of applying the usual transformer encoder with self-attentions, we alternate between  
 587 cross-attentions and self-attentions. For cross-attention layers, we use the node representations from  
 588 the previous layer as the query, and we always use  $K = V = F_t$ . We also apply an attention mask  
 589 to the cross-attention, so that each node will only attend to edges that it is incident with. For self-  
 590 attention layers, we always use the node representations from the previous layer as input. We always  
 591 use an even number of layers, with the first layer being a cross-attention layer, and the last layer  
 592 being a self-attention layer structured to produce the required parameters for the reverse distribution  
 593 as illustrated in Appendix B. For hyperparameters, we use 16 alternating layers, 8 attention heads,  
 594  $d_{\text{model}} = 256$ , feed-forward hidden dimension 1024, and dropout rate 0.1.

595 For training details on the TSP-20 task, we use the AdamW optimizer [26] with weight decay  $1e-4$ ,  
 596  $\varepsilon = 1e-8$ , and  $\beta = (0.9, 0.999)$ . We use the cosine annealing learning rate scheduler starting from  
 597  $2e-4$  and ending at 0. We train for 50 epochs with a batch size of 512. When computing the loss  
 598 (12), we use Monte-Carlo estimation for the expectation and sample 1 trajectory. We use a denoising  
 599 schedule of  $[0, 4, 5, 7]$ , with riffle shuffles as the forward process and GPL as the reverse process.  
 600 Finally, we use beam search for decoding, and we use a beam size of 256 both when decoding from  
 601 GPL and decoding along the denoising schedule.

## 602 E.5 Baselines Implementation Details

603 **Gumbel-Sinkhorn Network.** We have re-implemented the Gumbel-Sinkhorn Network [29] for  
 604 application on jigsaw puzzles, following the implementations provided in the official repository<sup>2</sup>. To  
 605 ensure a fair comparison, we conducted a thorough grid search of the model’s hyper-parameters. The  
 606 parameters included in our search space are as follows,

Hyperparameter	Values
Learning Rate (lr)	$\{10^{-3}, 10^{-4}, 10^{-5}\}$
Batch Size	$\{50\}$
Hidden Channels	$\{64, 128\}$
Kernel Size	$\{3, 5\}$
$\tau$	$\{0.2, 0.5, 1, 2, 5\}$
Number of Sinkhorn Iterations (n_sink_iter)	$\{20\}$
Number of Samples	$\{10\}$

Table 11: Hyperparameter Search Space for the Gumbel-Sinkhorn Network

607 **Diffsort & Error-free Diffsort** We have implemented two differentiable sorting networks from  
 608 the official repository<sup>3</sup> specific to error-free diffsort. For sorting 4-digit MNIST images, error-free  
 609 diffsort employs TransformerL as its backbone, with detailed hyperparameters listed in Table 12.  
 610 Conversely, Diffsort uses a CNN as its backbone, with a learning rate set to  $10^{-3.5}$ ; the relevant  
 611 hyperparameters are outlined in Table 13.

612 For jigsaw puzzle tasks, error-free diffsort continues to utilize a transformer, whereas Diffsort employs  
 613 a CNN. For other configurations, we align the settings with those of tasks having similar sequence  
 614 lengths in the 4-digit MNIST sorting task. For instance, for  $3 \times 3$  puzzles, we apply the same  
 615 configuration as used for sorting tasks with a sequence length of 9.

616 **TSP.** For the baselines for TSP, we first have 4 traditional operations research solvers. Gurobi [14]  
 617 and Concorde [1] are known as exact solvers, while LKH-3 [15] is a strong heuristic and 2-Opt [25]

<sup>2</sup>[https://github.com/google/gumbel\\_sinkhorn](https://github.com/google/gumbel_sinkhorn)

<sup>3</sup><https://github.com/jungtaekkim/error-free-differentiable-swap-functions>

Sequence Length	Steepness	Sorting Network	Loss Weight	Learning Rate
3	10	odd even	1.00	$10^{-4}$
5	26	odd even	1.00	$10^{-4}$
7	31	odd even	1.00	$10^{-4}$
9	34	odd even	1.00	$10^{-4}$
15	25	odd even	0.10	$10^{-4}$
32	124	odd even	0.10	$10^{-4}$
52	130	bitonic	0.10	$10^{-3.5}$
100	140	bitonic	0.10	$10^{-3.5}$

Table 12: Hyperparameters for Error-Free Diffsort on Sorting 4-Digit MNIST Numbers

Sequence Length	Steepness	Sorting Network
3	6	odd even
5	20	odd even
7	29	odd even
9	32	odd even
15	25	odd even
32	25	bitonic
52	25	bitonic
100	25	bitonic

Table 13: Hyperparameters for Diffsort on Sorting 4-Digit MNIST Numbers

618 is a weak heuristic. For LKH-3, we used 500 trials, and for 2-Opt, we used 5 random initial guesses  
619 with seed 42.

620 For the GCN model[18], we utilized the official repository<sup>4</sup> and adhered closely to its default  
621 configuration for the TSP-20 dataset. For DIFUSCO[43], we sourced it from its official repository<sup>5</sup>  
622 and followed the recommended configuration of TSP-50 dataset, with a minor adjustment in the batch  
623 size. We increased the batch size to 512 to accelerate the training process. For fair comparison, we  
624 also remove the post-processing heuristics in both models during the evaluation.

## 625 F Limitations

626 Despite the success of this method on various tasks, the model presented in this paper still requires a  
627 time-space complexity of  $\mathcal{O}(n^2)$  due to its reliance on the parametric representation of GPL and the  
628 backbone of transformer attention layers. This complexity poses a significant challenge in scaling up  
629 to applications involving larger symmetric groups or Lie groups.

<sup>4</sup><https://github.com/chaitjo/graph-convnet-tsp>

<sup>5</sup><https://github.com/Edward-Sun/DIFUSCO>

630 **NeurIPS Paper Checklist**

631 **1. Claims**

632 Question: Do the main claims made in the abstract and introduction accurately reflect the  
633 paper's contributions and scope?

634 Answer: [Yes]

635 Justification: Our abstract and Section 1 accurately summarize the paper's contributions and  
636 scope.

637 Guidelines:

- 638 • The answer NA means that the abstract and introduction do not include the claims  
639 made in the paper.
- 640 • The abstract and/or introduction should clearly state the claims made, including the  
641 contributions made in the paper and important assumptions and limitations. A No or  
642 NA answer to this question will not be perceived well by the reviewers.
- 643 • The claims made should match theoretical and experimental results, and reflect how  
644 much the results can be expected to generalize to other settings.
- 645 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
646 are not attained by the paper.

647 **2. Limitations**

648 Question: Does the paper discuss the limitations of the work performed by the authors?

649 Answer: [Yes]

650 Justification: We discuss the limitations of the work in Appendix F.

651 Guidelines:

- 652 • The answer NA means that the paper has no limitation while the answer No means that  
653 the paper has limitations, but those are not discussed in the paper.
- 654 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 655 • The paper should point out any strong assumptions and how robust the results are to  
656 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
657 model well-specification, asymptotic approximations only holding locally). The authors  
658 should reflect on how these assumptions might be violated in practice and what the  
659 implications would be.
- 660 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
661 only tested on a few datasets or with a few runs. In general, empirical results often  
662 depend on implicit assumptions, which should be articulated.
- 663 • The authors should reflect on the factors that influence the performance of the approach.  
664 For example, a facial recognition algorithm may perform poorly when image resolution  
665 is low or images are taken in low lighting. Or a speech-to-text system might not be  
666 used reliably to provide closed captions for online lectures because it fails to handle  
667 technical jargon.
- 668 • The authors should discuss the computational efficiency of the proposed algorithms  
669 and how they scale with dataset size.
- 670 • If applicable, the authors should discuss possible limitations of their approach to  
671 address problems of privacy and fairness.
- 672 • While the authors might fear that complete honesty about limitations might be used by  
673 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
674 limitations that aren't acknowledged in the paper. The authors should use their best  
675 judgment and recognize that individual actions in favor of transparency play an impor-  
676 tant role in developing norms that preserve the integrity of the community. Reviewers  
677 will be specifically instructed to not penalize honesty concerning limitations.

678 **3. Theory Assumptions and Proofs**

679 Question: For each theoretical result, does the paper provide the full set of assumptions and  
680 a complete (and correct) proof?

681 Answer: [Yes]

682 Justification: We provide complete proof for Proposition 1 and Proposition 2 in Appendix  
683 D.

684 Guidelines:

- 685 • The answer NA means that the paper does not include theoretical results.
- 686 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
687 referenced.
- 688 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 689 • The proofs can either appear in the main paper or the supplemental material, but if  
690 they appear in the supplemental material, the authors are encouraged to provide a short  
691 proof sketch to provide intuition.
- 692 • Inversely, any informal proof provided in the core of the paper should be complemented  
693 by formal proofs provided in appendix or supplemental material.
- 694 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 695 4. Experimental Result Reproducibility

696 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
697 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
698 of the paper (regardless of whether the code and data are provided or not)?

699 Answer: [Yes]

700 Justification: In Appendix E, we fully disclose all information to reproduce our experimental  
701 results, including dataset preparation, training details, and choices of hyperparameters as  
702 well as baselines' implementation details.

703 Guidelines:

- 704 • The answer NA means that the paper does not include experiments.
- 705 • If the paper includes experiments, a No answer to this question will not be perceived  
706 well by the reviewers: Making the paper reproducible is important, regardless of  
707 whether the code and data are provided or not.
- 708 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
709 to make their results reproducible or verifiable.
- 710 • Depending on the contribution, reproducibility can be accomplished in various ways.  
711 For example, if the contribution is a novel architecture, describing the architecture fully  
712 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
713 be necessary to either make it possible for others to replicate the model with the same  
714 dataset, or provide access to the model. In general, releasing code and data is often  
715 one good way to accomplish this, but reproducibility can also be provided via detailed  
716 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
717 of a large language model), releasing of a model checkpoint, or other means that are  
718 appropriate to the research performed.
- 719 • While NeurIPS does not require releasing code, the conference does require all submis-  
720 sions to provide some reasonable avenue for reproducibility, which may depend on the  
721 nature of the contribution. For example
  - 722 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
723 to reproduce that algorithm.
  - 724 (b) If the contribution is primarily a new model architecture, the paper should describe  
725 the architecture clearly and fully.
  - 726 (c) If the contribution is a new model (e.g., a large language model), then there should  
727 either be a way to access this model for reproducing the results or a way to reproduce  
728 the model (e.g., with an open-source dataset or instructions for how to construct  
729 the dataset).
  - 730 (d) We recognize that reproducibility may be tricky in some cases, in which case  
731 authors are welcome to describe the particular way they provide for reproducibility.  
732 In the case of closed-source models, it may be that access to the model is limited in  
733 some way (e.g., to registered users), but it should be possible for other researchers  
734 to have some path to reproducing or verifying the results.

#### 735 5. Open access to data and code

736 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
737 tions to faithfully reproduce the main experimental results, as described in supplemental  
738 material?

739 Answer: [Yes]

740 Justification: We've included codes to reproduce the main results in the supplemental  
741 material. We also attach a detailed README file that provides sufficient instructions.

742 Guidelines:

- 743 • The answer NA means that paper does not include experiments requiring code.
- 744 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
745 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 746 • While we encourage the release of code and data, we understand that this might not be  
747 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not  
748 including code, unless this is central to the contribution (e.g., for a new open-source  
749 benchmark).
- 750 • The instructions should contain the exact command and environment needed to run to  
751 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.  
752 cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 753 • The authors should provide instructions on data access and preparation, including how  
754 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 755 • The authors should provide scripts to reproduce all experimental results for the new  
756 proposed method and baselines. If only a subset of experiments are reproducible, they  
757 should state which ones are omitted from the script and why.
- 758 • At submission time, to preserve anonymity, the authors should release anonymized  
759 versions (if applicable).
- 760 • Providing as much information as possible in supplemental material (appended to the  
761 paper) is recommended, but including URLs to data and code is permitted.

## 762 6. Experimental Setting/Details

763 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
764 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
765 results?

766 Answer: [Yes]

767 Justification: All experimental settings and details are specified in Appendix E.4.

768 Guidelines:

- 769 • The answer NA means that the paper does not include experiments.
- 770 • The experimental setting should be presented in the core of the paper to a level of detail  
771 that is necessary to appreciate the results and make sense of them.
- 772 • The full details can be provided either with the code, in appendix, or as supplemental  
773 material.

## 774 7. Experiment Statistical Significance

775 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
776 information about the statistical significance of the experiments?

777 Answer: [Yes]

778 Justification: All reported experimental results are averaged over at least three runs with  
779 different random seeds.

780 Guidelines:

- 781 • The answer NA means that the paper does not include experiments.
- 782 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
783 dence intervals, or statistical significance tests, at least for the experiments that support  
784 the main claims of the paper.
- 785 • The factors of variability that the error bars are capturing should be clearly stated (for  
786 example, train/test split, initialization, random drawing of some parameter, or overall  
787 run with given experimental conditions).

- 788
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- 789
- The assumptions made should be given (e.g., Normally distributed errors).
- 790
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- 791
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- 792
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 793
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 794
- 795
- 796
- 797
- 798
- 799
- 800

## 801 8. Experiments Compute Resources

802 Question: For each experiment, does the paper provide sufficient information on the com-  
803 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
804 the experiments?

805 Answer: [Yes]

806 Justification: We provide information on the computation resources used for our experiments  
807 in Appendix E.4.

808 Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 817 9. Code Of Ethics

818 Question: Does the research conducted in the paper conform, in every respect, with the  
819 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

820 Answer: [Yes]

821 Justification: We preserve anonymity with the NeurIPS Codes of Ethics.

822 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 828 10. Broader Impacts

829 Question: Does the paper discuss both potential positive societal impacts and negative  
830 societal impacts of the work performed?

831 Answer: [NA]

832 Justification: There is no societal impact of the work performed.

833 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

834

835

836

- 837
- 838
- 839
- 840
- 841
- 842
- 843
- 844
- 845
- 846
- 847
- 848
- 849
- 850
- 851
- 852
- 853
- 854
- 855
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 856 11. Safeguards

857 Question: Does the paper describe safeguards that have been put in place for responsible  
858 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
859 image generators, or scraped datasets)?

860 Answer: [NA]

861 Justification: The paper poses no such risks.

862 Guidelines:

- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- The answer NA means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
  - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 873 12. Licenses for existing assets

874 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
875 the paper, properly credited and are the license and terms of use explicitly mentioned and  
876 properly respected?

877 Answer: [Yes]

878 Justification: We have cited the original paper of our reference code and datasets.

879 Guidelines:

- 880
- 881
- 882
- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- The answer NA means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
  - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 891
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 892
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 893
- 894

895 **13. New Assets**

896 Question: Are new assets introduced in the paper well documented and is the documentation  
897 provided alongside the assets?

898 Answer: [NA]

899 Justification: The paper does not release new assets.

900 Guidelines:

- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908

909 **14. Crowdsourcing and Research with Human Subjects**

910 Question: For crowdsourcing experiments and research with human subjects, does the paper  
911 include the full text of instructions given to participants and screenshots, if applicable, as  
912 well as details about compensation (if any)?

913 Answer: [NA]

914 Justification: The paper does not involve crowdsourcing nor research with human subjects.

915 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 916
- 917
- 918
- 919
- 920
- 921
- 922
- 923

924 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
925 Subjects**

926 Question: Does the paper describe potential risks incurred by study participants, whether  
927 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
928 approvals (or an equivalent approval/review based on the requirements of your country or  
929 institution) were obtained?

930 Answer: [NA]

931 Justification: The paper does not involve crowdsourcing nor research with human subjects.

932 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
  - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942