

---

# Prechastic Coding: An Alternative Approach to Neural Network Description Lengths

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The minimum description length (MDL) principle has a rich history of informing  
2 neural network research and there are numerous algorithms for developing efficient  
3 neural network description lengths. Of these methods, prequential coding, based on  
4 the prequential approach to statistics, has proven to be highly successful. Despite  
5 its achievements, general prequential coding limits learning at each increment  
6 to a prefix of a given dataset - a constraint which is potentially misaligned with  
7 an effective learning process. In this paper we introduce prechastic coding, an  
8 alternative to the prequential approach which is based on a guided, noisy sequence  
9 of intermediate learning steps. In our experiments we determine that the prechastic  
10 coding can challenge prequential coding in certain scenarios, whilst also leaving  
11 significant potential for further improvement.

## 12 1 Introduction

13 Pioneered by Jorma Rissanen [1, 2, 3, 4], the MDL principle has a rich history of informing and  
14 advancing machine learning, underlying important work on topics such as variational inference for  
15 neural networks [5, 6]. At a high level, the MDL principle advocates for model selection based  
16 on measures of both model performance and model complexity. This viewpoint can be informally  
17 expressed by the notion that a good model of some data allows for the efficient transmission of both  
18 the data and the model. Intuitively, overly complex models which fit the data extremely well are not  
19 desirable as, while the model itself can achieve highly compressed lossless encodings of the data, the  
20 combined cost of communicating the data and the model is large.

21 At first, neural networks (particularly deep learning models) appear to stand in contrast to the  
22 MDL philosophy as they often demonstrate compelling performance whilst having extremely high  
23 parameter counts. This misconception stems from a naive coding scheme where parameters are passed  
24 as raw floating point numbers before the lossless transmission of data. Alternate schemes which can  
25 be used to develop far better code lengths than the naive encoding include network compression,  
26 intrinsic dimension and variational approaches [5, 6, 7, 8, 9, 10]. However, Blier and Ollivier [10]  
27 demonstrated that all of these schemes are inferior to a method known as prequential coding.

28 Prequential coding stems from the prequential approach to statistics [11] and works by sending data  
29 incrementally and updating the model after each transmission (see Figure 1 for a high-level visual  
30 diagram of an iteration). As a result, the prequential scheme leverages a model's ability to generalise  
31 from limited data. Blier and Ollivier [10]'s work established the pre-eminence of prequential coding  
32 for description lengths of deep learning models. Subsequently prequential coding has also facilitated  
33 state-of-the-art results in compression as the Large Text Compression Benchmark [12], a competition  
34 to compress one gigabyte of English Wikipedia test, is currently topped by the nncp algorithm which  
35 is largely a prequential approach with some extra features [13, 14].

36 More recently, work by Bornschein et al. [15] found that the block transmission approach to prequential coding, used by [10] to lower computational costs, could be improved using techniques from  
 37 continuous learning. Additionally there has been significant research on the use of prequential code  
 38 lengths as an evaluation metric for various criteria [16, 17, 18]. Despite this popularity, the prequential  
 39 approach is not without its drawbacks as a description length/compression mechanism. Many datasets  
 40 might not be presented in an ordering particularly conducive for learning; for example, one might  
 41 conjecture that in many scenarios incrementally learning a body of text would be better done by  
 42 increasing the level of abstract complexity of the concept at each iteration, rather than progressing  
 43 through the text one word at a time. This begets the question - can one find a better general method  
 44 of computing description lengths for neural networks? In this paper we introduce a challenger to  
 45 prequential coding termed prechastic coding (a portmanteau of predictive and stochastic). Prechastic  
 46 coding shifts the concept of intermediate training datasets from the prequential viewpoint of cumulative,  
 47 sequential partitions to noisy views of the full dataset by allowing fake labels at broadly  
 48 diminishing rates across the scheme. Rather than predicting subsequent individual labels in the data  
 49 sequence, the prechastic method uses the model to iteratively denoise the stochastic, yet curated,  
 50 intermediate datasets. In our experiments we find that in select scenarios a greedy version of the  
 51 prechastic code approaches the performance of the prequential. However, the prechastic approach as  
 52 presented herein also allows for significant future improvement as a core component of the method,  
 53 *i.e.* the selection of guiding distributions, is left as an open-ended topic of discussion.  
 54

## 55 2 Prechastic Coding

56 In this section we will describe the specifics of the general prechastic approach along with variants of  
 57 interest. Before we proceed, we first describe the standard supervised learning compression scenario.  
 58 In line with Blier and Ollivier [10]’s work on prequential coding we will use this setting throughout  
 59 the remainder of the paper in order to develop the prechastic approach (although we note that this is  
 60 by no means a necessity and is simply for pedagogical reasons). Consider a sender and a receiver  
 61 who both have a copy of a sequence of  $N$  inputs  $x_{1:N}$  and have agreed to some identically initialised  
 62 learning model. The latter agreement often involves a high level description of an architecture  
 63 and initialisation procedure along with a mutual seed for a pseudorandom number generator. The  
 64 input data is randomly ordered yet identical for both the sender and the receiver. Each  $x_i$  has a  
 65 corresponding label  $y_i \in \{1, 2, \dots, K\}$  which are, initially, only known to the sender. The sender  
 66 would like to transmit  $y_{1:N}$  to the receiver using as little information as possible.

67 Consider a sequence of probability distributions  $Q_1, Q_2, \dots, Q_T$  each defined over  $\{1, 2, \dots, K\}^N$ .  
 68 Each individual  $Q_i$  assigns probabilities to all  $K^N$  possible permutations of labels, both true or false,  
 69 to the dataset  $x_{1:N}$ . Samples from  $Q_1, Q_2, \dots, Q_T$  constitute successive, intermediate datasets for the  
 70 selected model to train on and we shall therefore refer to them as the *guiding distributions*. Since,  
 71 until the final transmission, only the sender knows all of the true labels  $y_{1:N}$ , the sender must use this  
 72 information to compute guiding distributions which lead to an efficient code length.

73 To initiate the general prechastic algorithm, the receiver creates predictions for all the labels for the  
 74 dataset from its copy of the untrained model. We will denote these predictions as  $P_1$  and note that  
 75 in many circumstances  $P_1$  is likely to be an approximate uniform distribution over all  $K^N$  possible  
 76 predictions. The sender computes an identical copy of  $P_1$  and transmits some sample  $q_1 \sim Q_1$   
 77 using  $\mathcal{O}(\text{KL}[Q_1 \parallel P_1])$  bits; possible machinery for this is discussed below in the section on relative  
 78 entropy coding.  $Q_1$  could be pre-determined or computed once  $P_1$  has been calculated. The receiver  
 79 then uses the noisy labels  $q_1$  to train the model and create an updated set of predictions  $P_2$ .

80 This process is repeated for a total of  $T$  iterations with the sender transmitting some sample from  
 81  $Q_i$  at each step to the receiver who, subsequently, uses this sample to train their model and form  
 82 updated predictions  $P_{i+1}$ . The sum cost of these transmissions, including a final lossless encoding  
 83 of  $y_{1:N}$ , constitutes the total code length for the prechastic approach. A full diagram of the process  
 84 in comparison with the prequential approach can be found in Figure 1. A summary of the general  
 85 prechastic coding algorithm is given in Algorithm 1.

86 The difference between the prequential and prechastic methods boils down to a difference in how  
 87 the intermediate training datasets are viewed: in the prequential approach such intermediate datasets  
 88 are expanding restrictions of the original dataset; in the prechastic approach they are noisy views of  
 89 the entire original dataset with an overall trend towards less noise. While deep learning models have

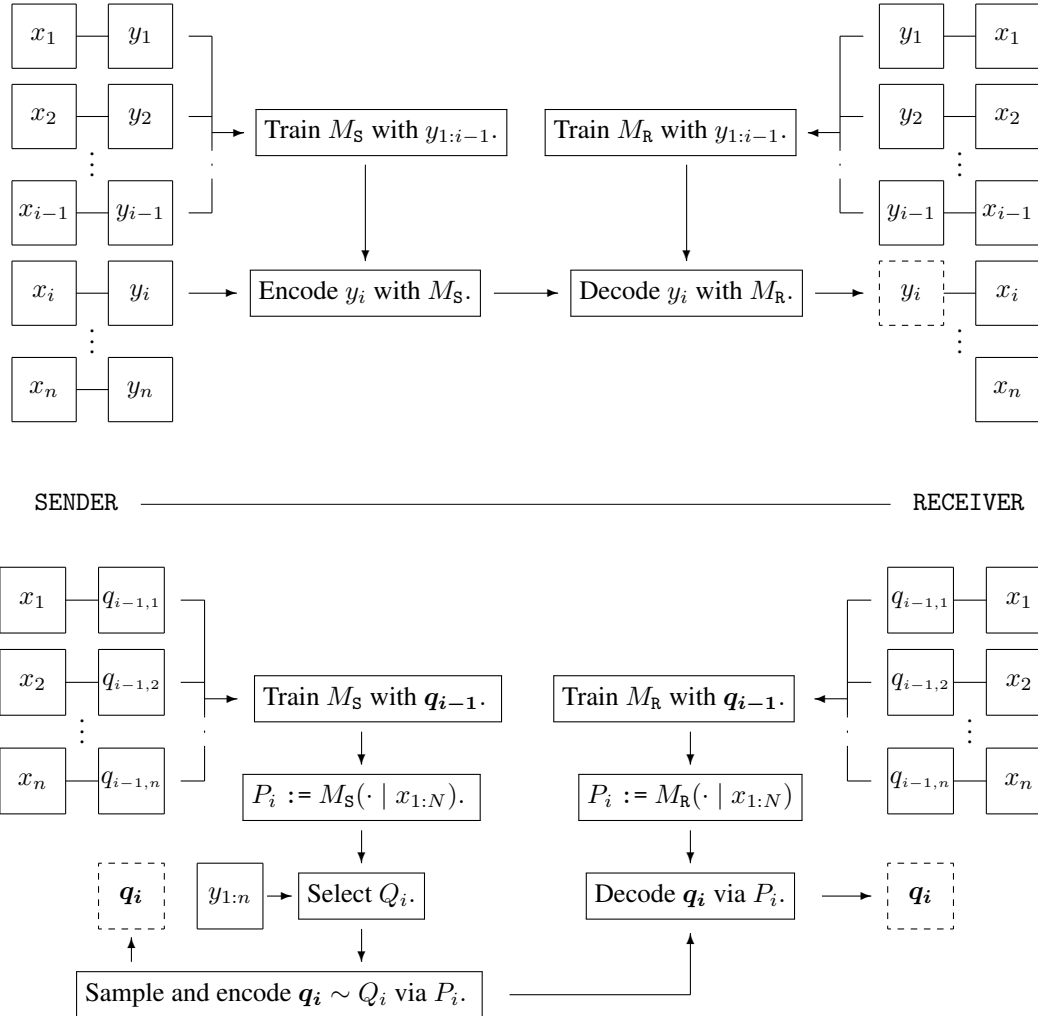


Figure 1: High level overview of an iteration of prequential coding (top) and prechastic coding (bottom). Both the sender and the receiver model are trained identically, *i.e.*  $M_S \equiv M_R$ . This diagram slightly differs from the presentation in Algorithm 1 as each iteration begins with model training.

90 the ability to reach negligible levels of training error on noisy labels, training regimens are typically  
 91 crafted for performance on unseen data. For this reason, if the guiding distributions are selected  
 92 carefully, we should in expectation see improvements between the prechastic iterations. We explore  
 93 quantitative results further in Section 3; however, we will first discuss relative entropy coding as well  
 94 as specific prechastic algorithms.

95 **Relative Entropy Coding** Consider the following communication scenario: a sender would like  
 96 to send a sample from a distribution  $Q$  to a receiver who only has access to a distribution  $P$ . The  
 97 sender does not care which particular sample is transmitted, only that it comes from the distribution  
 98  $Q$ . Relative entropy coding (REC) algorithms communicate such a sample with an expected code  
 99 length of  $\mathcal{O}(\text{KL}[Q \| P])$  [19, 20]. Initial work by Harsha et al. [21] proposed a computationally  
 100 intractable rejection sampling algorithm; later, Havasi et al. [22] used an importance sampling  
 101 approach which first generates  $M = \lceil 2^{\text{KL}[Q \| P]} \rceil$  samples from the  $P$  distribution.<sup>1</sup> Each sample  
 102  $x$  is weighted according to the ratio  $Q(x)/P(x)$  and, after normalisation, the resulting categorical  
 103 distribution is sampled to select an index from 1 to  $M$ . This index is then transmitted at a cost of  
 104  $\log_2(M) \approx \text{KL}[Q \| P]$ . Critically, Havasi et al. [22] demonstrated that, via a result from Chatterjee

<sup>1</sup>Note that in all instances in this paper, the Kullback-Leibler divergence is given in bits - many of the papers referenced in this section instead use nats.

---

**Algorithm 1** The generic prechastic coding algorithm.

---

```
1: S initialises a model  $M_S$ .
2: R initialises an identical model  $M_R$ .
3: for  $i := 1$  to  $T$  do
4:    $P_i := M(\cdot | x_{1:N})$  where  $M \equiv M_R \equiv M_S$ .
5:   S generates a sample  $q_i \sim Q_i$ .
6:   S transmits  $\mathcal{O}(\text{KL}[Q_i || P_i])$  bits which enable R to recreate  $q_i$ .
7:    $M_R$  is trained on  $q_i$ .
8:    $M_S$  is trained on  $q_i$  in an identical manner.
9: end for
10: S encodes  $y_{1:N}$  using  $M_S(y_{1:N} | x_{1:N})$ .
11: S transmits the code for  $y_{1:N}$  to R.
12: R decodes  $y_{1:N}$  using  $M_R(y_{1:N} | x_{1:N})$ .
```

---

105 and Diaconis [23], setting  $M = \lceil 2^{\text{KL}[Q || P]} \rceil$  was a sufficiently large sample size to keep the bias in  
106 the sampling low. In our experimental section we make use of the importance sampling procedure in  
107 a first-pass approach to selecting the guiding distributions. Further work on REC and REC-related  
108 methods include: Flamich et al. [19] who suggested a method of dividing the transmission process  
109 into a sequence of intermediary steps, Flamich et al. [20], who introduced approaches based on A\*  
110 sampling [24], Li and Gamal [25]’s research on Poisson functional representation, and Theis and  
111 Ahmed [26]’s Ordered Random Coding method (presented under the framework of the related reverse  
112 channel coding problem [27]). In some sense, carefully selecting guiding distributions can achieve a  
113 similar intermediary effect to the auxiliary variable method of Flamich et al. [19] as it also mitigates  
114 much of the exponential runtime effects.

## 115 2.1 The Greedy Prechastic Algorithm

116 We shall now describe an effective greedy approach to the prechastic approach which considers  
117 potential  $P_{i+1}$  values. Rather than choosing a specific  $Q_i$  from a process such as the minimisation  
118 of an optimisation problem as in Appendix A, the greedy approach generates  $G$  samples from  $P_i$   
119 then trains a model on each of these samples independently. The sample whose model, after training,  
120 minimises the cost of encoding the true values is then encoded via index at a cost of  $\mathcal{O}(\log(G))$  bits.  
121 Note that the greedy approach, which is outlined in Algorithm 2, is a slight deviation from the general  
122 prechastic scheme as it does not explicitly choose a guiding distribution  $Q_i$ .

## 123 3 Experiments

124 In the following experiments we evaluate the greedy prechastic algorithm in comparison with  
125 the prequential approach as well as two variations of the first-pass convex method described in  
126 Appendix A. We operate under the supervised learning scenario described in Section 2 and apply it to  
127 the MNIST [28] and Fashion-MNIST [29] datasets. Learning is conducted using a simple MLP with  
128 two 128 neuron hidden layers as well as a convolutional LeNet-style network [28].

129 These models were chosen in part to accommodate for the computational cost of computing a full,  
130 non-batched prequential code which requires  $\mathcal{O}(N^2)$  items of data to be processed per epoch of  
131 training. In comparison, if we consider time contributions as primarily determined by training, the  
132 general prechastic scheme is  $\mathcal{O}(TN)$  whilst the greedy scheme is  $\mathcal{O}(TNG)$ . For practical values  
133 of  $T$  and  $G$  that produce efficient code lengths, we found that our experiments had to restrict both  
134 datasets to smaller sizes of  $N = 128, 256,$  and  $512$  in order to lower runtimes. Despite this constraint,  
135 Bornschein et al. [30] found that model selection based on small dataset restrictions may give similar  
136 results to model selection which uses the entire dataset.

137 Ten trials of prequential experiments were run for each model and dataset combination using a batch  
138 size of 32. The models were trained for a total of five epochs and codelengths were computed using  
139 the final model at the end of the fifth epoch. We did not use a scheme which evaluated code lengths at  
140 the end of each epoch and subsequently took the best performing model as this would have required  
141 transmissions of the epoch index and consequently incurred a large penalty. It was determined that  
142 five epochs produced reasonably efficient prequential codelengths for the models and datasets used.

---

**Algorithm 2** The greedy prechastic coding algorithm.

---

```
1: S initialises a model  $M_S$ .
2: R initialises an identical model  $M_R$ .
3: for  $i := 1$  to  $T$  do
4:    $P_i := M(\cdot | x_{1:N})$  where  $M \equiv M_R \equiv M_S$ .
5:   R and S generate identical sets of  $G$  samples  $p_i^{(1)}, p_i^{(2)}, \dots, p_i^{(G)} \sim P_i$ .
6:   S initializes a model  $M'$ .
7:    $V := \infty, g^* := 1$ 
8:   for  $j := 1$  to  $G$  do
9:     S creates a clone of  $M'$  and trains it on  $p_i^{(j)}$ , obtaining  $M'_j$ .
10:    if  $-\log_2(M'_j(y_{1:N} | x_{1:N})) < V$  then
11:       $V := -\log_2(M'_j(y_{1:N} | x_{1:N}))$ 
12:       $g^* := j$ 
13:    end if
14:  end for
15:  S transmits  $g^*$  to R at a cost of  $\mathcal{O}(\log(G))$  bits.
16:   $M_R$  is trained on  $p_i^{(g^*)}$ .
17:   $M_S$  is trained on  $p_i^{(g^*)}$  in an identical manner.
18: end for
19: S encodes  $y_{1:N}$  using  $M_S(y_{1:N} | x_{1:N})$ .
20: S transmits the code for  $y_{1:N}$  to R.
21: R decodes  $y_{1:N}$  using  $M_R(y_{1:N} | x_{1:N})$ .
```

---

143 The prechastic experiments were conducted on the greedy prechastic algorithm along with two  
144 variants of the first pass approach from Appendix A which were iteratively solved using the CVX  
145 package [31, 32]. The first variant, FPC-Q directly sampled each  $Q_i$  whilst the second version, FPC-R  
146 used the importance sampling REC procedure of Havasi et al. [22] to indirectly sample each  $Q_i$   
147 through its respective  $P_i$ . Splitting the first pass approach into these two variants was done in order  
148 to quantify the bias affects from the importance sampling procedure. Note that because FPC-Q uses  
149 direct samples it is a thought experiment and not a practical compression algorithm. A running  
150 average of up to five of the most recently communicated samples were used as a training signal in  
151 order to improve stability.  $\beta$  was set to 7 and the cost of each iteration was logged as  $\log_2(\lceil 2^{\beta_i^*} \rceil)$   
152 (note that this does not communicate the size of  $\beta_i^*$  itself; in practice it is likely better to simply use  $\beta$   
153 and communicate it once). For the 128, 256, and 512 count dataset sizes, we used maximum iteration  
154 counts of 25, 50, and 100, respectively (the cost of transmitting the best index was included in the  
155 code lengths).

156 For the greedy prechastic experiments, the hyper-parameter  $G$  was set to 128, *i.e.* 7 bits of information  
157 was transmitted per iteration. To increase stability in the face of small datasets, multiple samples  
158 were generated from  $P_i$  for each of the  $G$  trials and the average values were used for training. Note  
159 that while there are still only  $G$  averaged options to choose from (and thus there is still only 7 bits of  
160 data transmitted per iteration) larger numbers of multiple samples drive the training signal towards  
161  $P_i$ . In order to balance this effect with the desired stability, for the 128, 256, and 512 count dataset  
162 sizes multiple sample values of 25, 5, and 2 were used, respectively. The larger datasets also required  
163 higher maximum iteration counts of 20, 40, and 60 for the 128, 256, and 512 count dataset sizes,  
164 respectively. The best result from across these iterations was taken as the code length (including the  
165 cost of transmitting this index). The results from the prequential, first pass, and greedy prechastic  
166 experiments are presented in Table 1. All code was executed on a consumer-grade build (Intel  
167 i7-4790k and an Nvidia RTX 3060) and the longer experiments typically took on the rough order of  
168 hours to a day.

## 169 4 Conclusion

170 The greedy prechastic algorithm performed well in our experiments, approaching and demonstrating  
171 comparable performance in the MNIST/LeNet testing suite. Further results across the remainder  
172 of experiments were competitive although the prequential tests consistently produced the best code

SIZE	CODING	MNIST	
		MLP	LENET
128	PREQ.	$0.896 \pm 0.008$ (380.8 $\pm$ 3.2)	$0.895 \pm 0.005$ (380.6 $\pm$ 2.1)
	FPC-Q	$0.981 \pm 0.006$ (417.0 $\pm$ 2.5)	$0.997 \pm 0.008$ (424.1 $\pm$ 3.3)
	FPC-R	$1.002 \pm 0.005$ (425.9 $\pm$ 2.1)	$1.014 \pm 0.004$ (431.3 $\pm$ 2.1)
	GREEDY	$0.948 \pm 0.006$ (402.9 $\pm$ 2.4)	$0.933 \pm 0.006$ (396.5 $\pm$ 2.4)
256	PREQ.	$0.726 \pm 0.006$ (617.7 $\pm$ 5.0)	$0.710 \pm 0.004$ (603.6 $\pm$ 3.8)
	FPC-Q	$0.879 \pm 0.012$ (747.3 $\pm$ 10.3)	$0.805 \pm 0.012$ (684.6 $\pm$ 10.1)
	FPC-R	$0.971 \pm 0.006$ (826.2 $\pm$ 4.9)	$0.935 \pm 0.010$ (795.2 $\pm$ 8.6)
	GREEDY	$0.744 \pm 0.010$ (632.6 $\pm$ 7.7)	$0.696 \pm 0.006$ (592.0 $\pm$ 5.3)
512	PREQ.	$0.540 \pm 0.006$ (917.9 $\pm$ 9.4)	$0.512 \pm 0.005$ (870.0 $\pm$ 7.8)
	FPC-Q	$0.762 \pm 0.013$ (1296.9 $\pm$ 22.1)	$0.622 \pm 0.007$ (1057.8 $\pm$ 11.1)
	FPC-R	$0.919 \pm 0.006$ (1563.2 $\pm$ 9.3)	$0.770 \pm 0.008$ (1309.0 $\pm$ 13.9)
	GREEDY	$0.620 \pm 0.006$ (1055.3 $\pm$ 9.8)	$0.517 \pm 0.004$ (878.7 $\pm$ 6.0)

SIZE	CODING	FASHION-MNIST	
		MLP	LENET
128	PREQ.	$0.718 \pm 0.011$ (305.2 $\pm$ 4.7)	$0.836 \pm 0.008$ (355.4 $\pm$ 3.2)
	FPC-Q	$0.909 \pm 0.013$ (386.5 $\pm$ 5.4)	$0.980 \pm 0.009$ (416.5 $\pm$ 3.9)
	FPC-R	$0.994 \pm 0.013$ (422.6 $\pm$ 5.3)	$1.018 \pm 0.007$ (433.1 $\pm$ 3.1)
	GREEDY	$0.851 \pm 0.007$ (361.7 $\pm$ 3.0)	$0.926 \pm 0.007$ (393.7 $\pm$ 2.9)
256	PREQ.	$0.555 \pm 0.010$ (472.4 $\pm$ 8.3)	$0.715 \pm 0.004$ (608.3 $\pm$ 3.3)
	FPC-Q	$0.746 \pm 0.010$ (634.3 $\pm$ 8.9)	$0.869 \pm 0.008$ (739.3 $\pm$ 7.0)
	FPC-R	$0.897 \pm 0.011$ (763.1 $\pm$ 9.5)	$0.952 \pm 0.009$ (809.5 $\pm$ 7.4)
	GREEDY	$0.619 \pm 0.007$ (526.7 $\pm$ 6.0)	$0.743 \pm 0.007$ (631.6 $\pm$ 5.6)
512	PREQ.	$0.447 \pm 0.005$ (761.0 $\pm$ 8.4)	$0.577 \pm 0.004$ (980.8 $\pm$ 7.2)
	FPC-Q	$0.688 \pm 0.008$ (1170.5 $\pm$ 13.2)	$0.756 \pm 0.007$ (1285.5 $\pm$ 11.2)
	FPC-R	$0.815 \pm 0.008$ (1386.1 $\pm$ 12.9)	$0.868 \pm 0.007$ (1476.6 $\pm$ 12.1)
	GREEDY	$0.494 \pm 0.004$ (840.8 $\pm$ 7.6)	$0.619 \pm 0.004$ (1053.5 $\pm$ 6.4)

Table 1: Results from prechastic and prequential experiments on restrictions of the MNIST and Fashion-MNIST datasets. The average compression ratio and the average size in bits are presented along with their corresponding standard error values.

173 lengths. However, across the largest datasets of size 512, the absolute compression ratio of the greedy  
174 approach was never more than eight percent greater than the prequential (less than fifteen percent in  
175 terms of relative performance to the prequential). Expectedly, as a first-pass at selecting the guiding  
176 distributions, the convex results fared worse than the greedy algorithm. However, the comparison  
177 between FPC-Q and FPC-R performance did yield insights into the underlying REC algorithm used  
178 in FPC-R. As there was a large drop-off from the FPC-Q results down to the FPC-R the importance  
179 sampling approach of Havasi et al. [22] clearly introduced a significant amount of sampling bias.

180 Looking beyond these experiments, the prechastic approach is a highly flexible coding scheme with  
181 potential for further development and improvement. Because the general prechastic scheme does not  
182 prescribe a specific method for selecting the guiding distributions, future work should investigate  
183 more advanced selection techniques to further improve prechastic code lengths. The greedy method  
184 could also potentially be improved by considering higher order decisions at each iteration. Future  
185 research might also consider the reduction of computational costs which partially necessitated dataset  
186 size restrictions during the experiments. Bornschein et al. [30] found that model selection based  
187 on small dataset restrictions may provide similar results to using the entire dataset, however, if one  
188 would still like to use large datasets one possible method might be to bootstrap from smaller datasets  
189 up to large ones, spreading the prechastic iterations across smaller views of the original dataset.

## References

- 190
- 191 [1] Jorma Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- 192 [2] Jorma Rissanen. Universal coding, information, prediction, and estimation. *IEEE Transactions*  
193 *on Information theory*, 30(4):629–636, 1984.
- 194 [3] Jorma Rissanen. Stochastic complexity and modeling. *The annals of statistics*, pages 1080–1100,  
195 1986.
- 196 [4] Jorma Rissanen. Stochastic complexity. *Journal of the Royal Statistical Society: Series B*  
197 *(Methodological)*, 49(3):223–239, 1987.
- 198 [5] Geoffrey E. Hinton and Drew van Camp. Keeping the Neural Networks Simple By Minimizing  
199 the Description Length of the Weights. In *Proceedings of the Sixth Annual Conference on*  
200 *Computational Learning Theory*, COLT '93, page 5–13. Association for Computing Machinery,  
201 Aug 1993.
- 202 [6] Alex Graves. Practical Variational Inference for Neural Networks. In *Advances in Neural*  
203 *Information Processing Systems*, volume 24, 2011.
- 204 [7] Song Han, Huizi Mao, and William J. Dally. Deep Compression: Compressing Deep Neural  
205 Network with Pruning, Trained Quantization and Huffman Coding. In Yoshua Bengio and Yann  
206 LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San*  
207 *Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.
- 208 [8] Ting-Bing Xu, Peipei Yang, Xu-Yao Zhang, and Cheng-Lin Liu. Margin-Aware Binarized  
209 Weight Networks for Image Classification. In *Image and Graphics: 9th International Confer-*  
210 *ence, ICIG 2017*, pages 590–601. Springer, 2017.
- 211 [9] Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the Intrinsic Di-  
212 mension of Objective Landscapes. In *6th International Conference on Learning Representations,*  
213 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.*  
214 OpenReview.net, 2018.
- 215 [10] Léonard Blier and Yann Ollivier. The description length of deep learning models. In *Advances*  
216 *in Neural Information Processing Systems 31: Annual Conference on Neural Information*  
217 *Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages  
218 2220–2230, 2018.
- 219 [11] A. P. Dawid. Present Position and Potential Developments: Some Personal Views: Statistical  
220 Theory: The Prequential Approach. *Journal of the Royal Statistical Society. Series A (General)*,  
221 147(2):278–292, 1984.
- 222 [12] Matt Mahoney. Large Text Compression Benchmark, n.d. URL [http://mattmahoney.net/](http://mattmahoney.net/dc/text.html)  
223 [dc/text.html](http://mattmahoney.net/dc/text.html).
- 224 [13] Fabrice Bellard. Lossless Data Compression with Neural Networks. Technical report, Self-  
225 published, May 2019. URL <https://bellard.org/nncp/nncp.pdf>.
- 226 [14] Fabrice Bellard. NNCP v2: Lossless Data Compression with Transformer. Technical report,  
227 Self-published, Feb 2021. URL [https://bellard.org/nncp/nncp\\_v2.1.pdf](https://bellard.org/nncp/nncp_v2.1.pdf).
- 228 [15] Jörg Bornschein, Yazhe Li, and Marcus Hutter. Sequential learning of neural networks for  
229 prequential MDL. In *The Eleventh International Conference on Learning Representations, ICLR*  
230 *2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.](https://openreview.net/pdf?id=dMMPUvNSYJr)  
231 [net/pdf?id=dMMPUvNSYJr](https://openreview.net/pdf?id=dMMPUvNSYJr).
- 232 [16] Dani Yogatama, Cyprien de Masson d’Autume, Jerome T. Connor, Tomás Kociský, Mike  
233 Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, and Phil  
234 Blunsom. Learning and evaluating general linguistic intelligence. *CoRR*, abs/1901.11373, 2019.  
235 URL <http://arxiv.org/abs/1901.11373>.
- 236 [17] Charles Lovering, Rohan Jha, Tal Linzen, and Ellie Pavlick. Predicting Inductive Biases of  
237 Pre-Trained Models. In *9th International Conference on Learning Representations, ICLR 2021,*  
238 *Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- 239 [18] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. Rissanen Data Analysis: Examining Dataset  
240 Characteristics via Description Length. In *Proceedings of the 38th International Conference on*  
241 *Machine Learning*, page 8500–8513. PMLR, Jul 2021.

- 242 [19] Gergely Flamich, Marton Havasi, and José Miguel Hernández-Lobato. Compressing images by  
243 encoding their latent representations with relative entropy coding. In *NeurIPS*, 2020.
- 244 [20] Gergely Flamich, Stratis Markou, and José Miguel Hernández-Lobato. Fast relative entropy  
245 coding with A\* coding. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*,  
246 pages 6548–6577. PMLR, 2022.
- 247 [21] Prahladh Harsha, Rahul Jain, David McAllester, and Jaikumar Radhakrishnan. The communica-  
248 tion complexity of correlation. In *Twenty-Second Annual IEEE Conference on Computational*  
249 *Complexity (CCC’07)*, pages 10–23. IEEE, 2007.
- 250 [22] Marton Havasi, Robert Peharz, and José Miguel Hernández-Lobato. Minimal random code learn-  
251 ing: Getting bits back from compressed model parameters. In *ICLR (Poster)*. OpenReview.net,  
252 2019.
- 253 [23] Sourav Chatterjee and Persi Diaconis. The sample size required in importance sampling. *The*  
254 *Annals of Applied Probability*, 28(2):1099–1135, 2018.
- 255 [24] Chris J. Maddison, Daniel Tarlow, and Tom Minka. A\* sampling. In *NIPS*, pages 3086–3094,  
256 2014.
- 257 [25] Cheuk Ting Li and Abbas El Gamal. Strong functional representation lemma and applications  
258 to coding theorems. *IEEE Trans. Inf. Theory*, 64(11):6967–6978, 2018.
- 259 [26] Lucas Theis and Noureldin Y. Ahmed. Algorithms for the communication of samples. In *ICML*,  
260 volume 162 of *Proceedings of Machine Learning Research*, pages 21308–21328. PMLR, 2022.
- 261 [27] Charles H Bennett, Peter W Shor, John A Smolin, and Ashish V Thapliyal. Entanglement-  
262 assisted capacity of a quantum channel and the reverse shannon theorem. *IEEE transactions on*  
263 *Information Theory*, 48(10):2637–2655, 2002.
- 264 [28] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
265 applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- 266 [29] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for  
267 benchmarking machine learning algorithms, 2017.
- 268 [30] Jorg Bornschein, Francesco Visin, and Simon Osindero. Small data, big decisions: Model  
269 selection in the small-data regime. In *International conference on machine learning*, pages  
270 1035–1044. PMLR, 2020.
- 271 [31] Michael Grant and Stephen Boyd. CVX: Matlab software for disciplined convex programming,  
272 version 2.1. <https://cvxr.com/cvx>, March 2014.
- 273 [32] Michael Grant and Stephen Boyd. Graph implementations for nonsmooth convex programs.  
274 In V. Blondel, S. Boyd, and H. Kimura, editors, *Recent Advances in Learning and Control*,  
275 Lecture Notes in Control and Information Sciences, pages 95–110. Springer-Verlag Limited,  
276 2008. [http://stanford.edu/~boyd/graph\\_dcp.html](http://stanford.edu/~boyd/graph_dcp.html).

## 277 A First-Pass Convex Approach

278 One of the central problems left unanswered by the general prechastic approach is how to select the  
279 guiding distributions. Choosing efficient  $Q_1, Q_2, \dots, Q_T$  is a challenging problem which essentially  
280 requires one to design an appropriately difficult curriculum for potentially complex learning models.  
281 The following is a rudimentary attempt designed largely to illustrate the difficulties of selecting the  
282 guiding distributions. Consider the convex optimization problem



$$\begin{aligned}
& \min_{\beta_i, [\mathbf{Q}_i]} \beta_i - \sum_{j=1}^N \log_2 \left( [Q_i]_{j, y_j} \right) \\
& \text{s.t.} \quad \sum_{j=1}^N \text{KL}[[\mathbf{Q}_i]_j \parallel [\mathbf{P}_i]_j] \leq \beta_i \\
& \quad \sum_{k=1}^K [Q_i]_{j,k} = 1, \quad \forall j \\
& \quad 0 \leq [Q_i]_{j,k} \leq 1, \quad \forall j, k
\end{aligned}$$

283 where  $[\mathbf{Q}_i]$  is an  $N \times K$  matrix which represents  $N$  independent categorical distributions across  
284 the label classes for each of the inputs. For an input  $x_j$ , the prediction rendered by the model  $P_i$   
285 is denoted as  $[\mathbf{P}_i]_j$ .  $[Q_i]_{j,k}$  is the probability of label  $k$  given the distribution  $[Q_i]_j$ . Note that the  
286 budget variable  $\beta_i$  is implicitly non-negative. The cost function measures the order of information  
287 that would have to be transmitted for the receiver to draw single sample from  $Q_i$  along with the  
288 cost of sending the true labels  $y_{1:N}$  if the receiver were to form predictions using  $Q_i$ . Naturally, it is  
289 unlikely that the receiver will be able to predict in a manner identical to  $Q_i$  after training on a single  
290 sample; however,  $Q_i$  is used as  $P_{i+1}$  would require model training.

291 By minimizing the cost function over  $\beta_i$  and  $[\mathbf{Q}_i]$ , a trade-off is struck between the quality of  
292 guidance and the rough cost of communicating a sample. In our experiments we also bound  $\beta_i$  to a  
293 hyper-parameter  $\beta$  by introducing the constraint  $\beta_i \leq \beta$ . This change allowed us to limit the rate of  
294 change of the guiding distributions over the iterations and also avoid intractably high computational  
295 costs from the REC importance sampling procedure.

296 **NeurIPS Paper Checklist**

297 **1. Claims**

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

300 Answer: [Yes]

301 Justification: The claims made in the paper's abstract and introduction accurately reflect its  
302 contributions and scope without embellishment.

303 Guidelines:

- 304 • The answer NA means that the abstract and introduction do not include the claims  
305 made in the paper.
- 306 • The abstract and/or introduction should clearly state the claims made, including the  
307 contributions made in the paper and important assumptions and limitations. A No or  
308 NA answer to this question will not be perceived well by the reviewers.
- 309 • The claims made should match theoretical and experimental results, and reflect how  
310 much the results can be expected to generalize to other settings.
- 311 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
312 are not attained by the paper.

313 **2. Limitations**

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

315 Answer: [Yes]

316 Justification: Yes, this is discussed in the concluding section of the paper.

317 Guidelines:

- 318 • The answer NA means that the paper has no limitation while the answer No means that  
319 the paper has limitations, but those are not discussed in the paper.
- 320 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 321 • The paper should point out any strong assumptions and how robust the results are to  
322 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
323 model well-specification, asymptotic approximations only holding locally). The authors  
324 should reflect on how these assumptions might be violated in practice and what the  
325 implications would be.
- 326 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
327 only tested on a few datasets or with a few runs. In general, empirical results often  
328 depend on implicit assumptions, which should be articulated.
- 329 • The authors should reflect on the factors that influence the performance of the approach.  
330 For example, a facial recognition algorithm may perform poorly when image resolution  
331 is low or images are taken in low lighting. Or a speech-to-text system might not be  
332 used reliably to provide closed captions for online lectures because it fails to handle  
333 technical jargon.
- 334 • The authors should discuss the computational efficiency of the proposed algorithms  
335 and how they scale with dataset size.
- 336 • If applicable, the authors should discuss possible limitations of their approach to  
337 address problems of privacy and fairness.
- 338 • While the authors might fear that complete honesty about limitations might be used by  
339 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
340 limitations that aren't acknowledged in the paper. The authors should use their best  
341 judgment and recognize that individual actions in favor of transparency play an impor-  
342 tant role in developing norms that preserve the integrity of the community. Reviewers  
343 will be specifically instructed to not penalize honesty concerning limitations.

344 **3. Theory Assumptions and Proofs**

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

347 Answer: [NA]

348 Justification: There are no theoretical results.

349 Guidelines:

- 350 • The answer NA means that the paper does not include theoretical results.
- 351 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 352 referenced.
- 353 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 354 • The proofs can either appear in the main paper or the supplemental material, but if
- 355 they appear in the supplemental material, the authors are encouraged to provide a short
- 356 proof sketch to provide intuition.
- 357 • Inversely, any informal proof provided in the core of the paper should be complemented
- 358 by formal proofs provided in appendix or supplemental material.
- 359 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 360 4. Experimental Result Reproducibility

361 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

362 perimental results of the paper to the extent that it affects the main claims and/or conclusions

363 of the paper (regardless of whether the code and data are provided or not)?

364 Answer: [Yes]

365 Justification: The paper contains all the necessary information to reproduce the experimental

366 results to the extent required to justify the claims and conclusions.

367 Guidelines:

- 368 • The answer NA means that the paper does not include experiments.
- 369 • If the paper includes experiments, a No answer to this question will not be perceived
- 370 well by the reviewers: Making the paper reproducible is important, regardless of
- 371 whether the code and data are provided or not.
- 372 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 373 to make their results reproducible or verifiable.
- 374 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 375 For example, if the contribution is a novel architecture, describing the architecture fully
- 376 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 377 be necessary to either make it possible for others to replicate the model with the same
- 378 dataset, or provide access to the model. In general, releasing code and data is often
- 379 one good way to accomplish this, but reproducibility can also be provided via detailed
- 380 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 381 of a large language model), releasing of a model checkpoint, or other means that are
- 382 appropriate to the research performed.
- 383 • While NeurIPS does not require releasing code, the conference does require all submis-
- 384 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 385 nature of the contribution. For example
- 386 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 387 to reproduce that algorithm.
- 388 (b) If the contribution is primarily a new model architecture, the paper should describe
- 389 the architecture clearly and fully.
- 390 (c) If the contribution is a new model (e.g., a large language model), then there should
- 391 either be a way to access this model for reproducing the results or a way to reproduce
- 392 the model (e.g., with an open-source dataset or instructions for how to construct
- 393 the dataset).
- 394 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 395 authors are welcome to describe the particular way they provide for reproducibility.
- 396 In the case of closed-source models, it may be that access to the model is limited in
- 397 some way (e.g., to registered users), but it should be possible for other researchers
- 398 to have some path to reproducing or verifying the results.

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

400 Question: Does the paper provide open access to the data and code, with sufficient instruc-

401 tions to faithfully reproduce the main experimental results, as described in supplemental

402 material?

403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453

Answer: [Yes]

Justification: Code for the relevant algorithms is available upon request.

Guidelines:

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

## 6. Experimental Setting/Details

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

Answer: [Yes]

Justification: The paper specifies all the necessary experimental details needed to understand the results.

Guidelines:

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

## 7. Experiment Statistical Significance

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

Answer: [Yes]

Justification: The standard error is provided for all experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 454 • It should be clear whether the error bar is the standard deviation or the standard error  
455 of the mean.
- 456 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
457 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
458 of Normality of errors is not verified.
- 459 • For asymmetric distributions, the authors should be careful not to show in tables or  
460 figures symmetric error bars that would yield results that are out of range (e.g. negative  
461 error rates).
- 462 • If error bars are reported in tables or plots, The authors should explain in the text how  
463 they were calculated and reference the corresponding figures or tables in the text.

## 464 8. Experiments Compute Resources

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

468 Answer: [Yes]

469 Justification: All computational resources are fully documented in the experiments section  
470 of the paper.

471 Guidelines:

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

## 480 9. Code Of Ethics

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

483 Answer: [Yes]

484 Justification: The research conducted for this paper fully conforms with the NeurIPS Code  
485 of Ethics.

486 Guidelines:

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

## 492 10. Broader Impacts

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

495 Answer: [NA]

496 Justification: There is no direct mechanism for societal impact.

497 Guidelines:

- 498 • The answer NA means that there is no societal impact of the work performed.
- 499 • If the authors answer NA or No, they should explain why their work has no societal  
500 impact or why the paper does not address societal impact.
- 501 • Examples of negative societal impacts include potential malicious or unintended uses  
502 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
503 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
504 groups), privacy considerations, and security considerations.

- 505 • The conference expects that many papers will be foundational research and not tied  
506 to particular applications, let alone deployments. However, if there is a direct path to  
507 any negative applications, the authors should point it out. For example, it is legitimate  
508 to point out that an improvement in the quality of generative models could be used to  
509 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
510 that a generic algorithm for optimizing neural networks could enable people to train  
511 models that generate Deepfakes faster.
- 512 • The authors should consider possible harms that could arise when the technology is  
513 being used as intended and functioning correctly, harms that could arise when the  
514 technology is being used as intended but gives incorrect results, and harms following  
515 from (intentional or unintentional) misuse of the technology.
- 516 • If there are negative societal impacts, the authors could also discuss possible mitigation  
517 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
518 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
519 feedback over time, improving the efficiency and accessibility of ML).

## 520 11. Safeguards

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

524 Answer: [NA]

525 Justification: The paper poses no such risks.

526 Guidelines:

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

## 537 12. Licenses for existing assets

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

541 Answer: [Yes]

542 Justification: All creators are properly credited in the text.

543 Guidelines:

- 544 • The answer NA means that the paper does not use existing assets.
- 545 • The authors should cite the original paper that produced the code package or dataset.
- 546 • The authors should state which version of the asset is used and, if possible, include a  
547 URL.
- 548 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 549 • For scraped data from a particular source (e.g., website), the copyright and terms of  
550 service of that source should be provided.
- 551 • If assets are released, the license, copyright information, and terms of use in the  
552 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
553 has curated licenses for some datasets. Their licensing guide can help determine the  
554 license of a dataset.
- 555 • For existing datasets that are re-packaged, both the original license and the license of  
556 the derived asset (if it has changed) should be provided.

557 • If this information is not available online, the authors are encouraged to reach out to  
558 the asset’s creators.

559 **13. New Assets**

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

562 Answer: [NA]

563 Justification: The paper does not release new assets.

564 Guidelines:

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

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

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

577 Answer: [NA]

578 Justification: The paper does not involve crowdsourcing or research with human subjects.

579 Guidelines:

- 580 • The answer NA means that the paper does not involve crowdsourcing nor research with  
581 human subjects.
- 582 • Including this information in the supplemental material is fine, but if the main contribu-  
583 tion of the paper involves human subjects, then as much detail as possible should be  
584 included in the main paper.
- 585 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
586 or other labor should be paid at least the minimum wage in the country of the data  
587 collector.

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

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

594 Answer: [NA]

595 Justification: The paper does not involve crowdsourcing or research with human subjects.

596 Guidelines:

- 597 • The answer NA means that the paper does not involve crowdsourcing nor research with  
598 human subjects.
- 599 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
600 may be required for any human subjects research. If you obtained IRB approval, you  
601 should clearly state this in the paper.
- 602 • We recognize that the procedures for this may vary significantly between institutions  
603 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
604 guidelines for their institution.
- 605 • For initial submissions, do not include any information that would break anonymity (if  
606 applicable), such as the institution conducting the review.