
Scalable Ensemble Diversification for OOD Generalization and Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Training a diverse ensemble of models has several practical application scenarios,
2 such as model selection for out-of-distribution (OOD) generalization and the
3 detection of OOD samples via Bayesian principles. Previous approaches to diverse
4 ensemble training have relied on the framework of letting the models make the
5 correct predictions for the given in-distribution (ID) data while letting them come up
6 with different hypotheses for the OOD data. As such, they require well-separated
7 ID and OOD datasets to ensure a performant and diverse ensemble and have
8 only been verified in smaller-scale lab environments where such a separation is
9 readily available. In this work, we propose a framework, Scalable Ensemble
10 Diversification (SED), for scaling up existing diversification methods to large-scale
11 datasets and tasks (e.g. ImageNet), where the ID-OOD separation may not be
12 available. SED automatically identifies OOD samples within the large-scale ID
13 dataset on the fly and encourages the ensemble to make diverse hypotheses on
14 them. To make SED more suitable for large-scale applications, we propose an
15 algorithm to speed up the expensive pairwise disagreement computation. We verify
16 the resulting diversification of the ensemble on ImageNet and demonstrate the
17 benefit of diversification on the OOD generalization and OOD detection tasks.
18 In particular, for OOD detection, we propose a novel uncertainty score estimator
19 based on the diversity of ensemble hypotheses, which lets SED surpass all the
20 considered baselines in OOD detection task. Code will be available soon.

21 1 Introduction

22 Training a diverse ensemble of models is useful in multiple applications. Diverse ensembles are used
23 to enhance out-of-distribution (OOD) generalization, where strong spurious features learned from
24 the in-distribution (ID) training data hinder generalization [30, 31, 28, 23]. By learning multiple
25 hypotheses, the ensemble is given a chance to learn causal features that are otherwise overshadowed
26 by the prominent spurious features [39, 4]. In Bayesian machine learning, diversification of the
27 posterior samples has been studied as a means to improve the precision and efficiency of sample
28 uncertainty estimates [5, 37].

29 A common strategy to train a diverse ensemble is to introduce two objectives: one for the main
30 task and one for diversification [29, 5, 28, 23]. The main task loss, such as the cross-entropy loss
31 for classification, encourages the hypotheses to solve the task on the labeled ID training set. The
32 diversification loss encourages the hypotheses to diversify the responses on an unlabelled OOD
33 dataset [28, 23] (Figure 1). The datasets for the objectives are separated to avoid contradictory
34 objectives: prediction diversification on the ID set will encourage wrong answers if there is only one
35 correct label.

36 This strategy, however, requires a separate OOD dataset where the hypotheses may make diverse
 37 predictions without harming the main task performance on the ID training samples. Previous work
 38 has thus been tested on hypothetical lab settings where the spurious and causal features can easily be
 39 controlled to secure separate ID and OOD datasets for diverse ensemble training. It is not clear yet
 40 how one could diversify an ensemble of models for realistic, uncontrolled, and large-scale applications
 41 (e.g. ImageNet scale) where collecting a separate OOD dataset can be very costly, if not impossible.

42 To address the scalability challenge,
 43 we propose a novel diversification
 44 framework, Scalable Ensemble Diversification (SED, Figure 1). We intro-
 45 duce three ingredients. (1) OOD sam-
 46 ples are dynamically selected from the
 47 ID training samples, on which the mod-
 48 els are trained to make different predic-
 49 tions. (2) At each iteration, a subset of
 50 model pairs are stochastically selected
 51 to construct the disagreement objec-
 52 tive, rather than the full list of model
 53 pairs. (3) Deep networks are trained
 54 to diversify only a few layers at the
 55 end, rather than the full networks. This
 56 framework allows scaling up existing
 57 ensemble diversification methods. In
 58 this work, we focus on scaling up the
 59 Agree to Disagree (A2D) method [28].
 60

61 We verify that SED diversifies a model
 62 ensemble trained on ImageNet. We
 63 demonstrate the benefit of diversifica-
 64 tion on OOD generalization and OOD
 65 detection tasks. For the former, we showcase the usage of SED-diversified ensemble in three variants:
 66 (a) *vanilla ensemble* of prediction probabilities [22], (b) an average of the model weights through
 67 *model soup* [38], and (c) the *oracle selection* of the individual models for each OOD test set [23, 30].
 68 In all three cases, SED achieves a superior generalization to OOD datasets like ImageNet-A/R/C,
 69 OpenImages, and iNaturalist.

70 For OOD detection, we seek multiple ways to use the SED-diversified ensemble: (a) treating them as
 71 samples of the Bayesian posterior and (b) using our novel OODness estimate of Predictive Diversity
 72 Score (PDS) that measures the diversity of predictions from an ensemble. We show that PDS provides
 73 a superior detection of OOD samples like ImageNet-A/R/C, OpenImages, and iNaturalist.

74 Our contributions are

- 75 1. Scalable Ensemble Diversification (SED) framework that scales up existing ensemble
 76 methods;
- 77 2. Predictive Diversity Score (PDS) that computes the OODness score for samples based on
 78 ensemble prediction diversity;
- 79 3. First demonstration of the ensemble diversification and its application to OOD generalization
 80 and detection at ImageNet level.

81 The code will be released with the next versions of the manuscript.

82 2 Related work

83 In this section, we give a short overview of ensembling methods. At first, we speak about ensembles
 84 in general and the role of diversity in them (§ 2.1), then we focus on ensembling methods for neural
 85 networks and separate them into two big groups. The first group includes algorithms that use loss
 86 regularizers (§ 2.2) and the second group covers works that do not modify the training loss (§ 2.3).

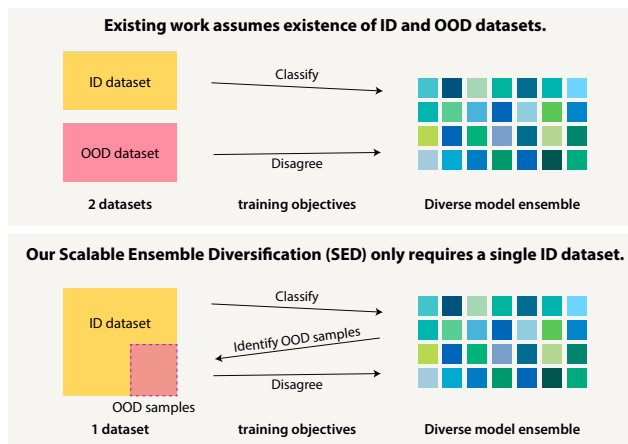


Figure 1: **Existing diversification work vs SED.** Unlike previous diversification approaches that require a separate OOD dataset on which the models are trained to diverge, our Scalable Ensemble Diversification (SED) operates on a single ID dataset where OOD samples are dynamically identified and are used to let the ensemble members diverge.

87 2.1 Ensembles as a technique

88 Ensembling is a powerful technique of aggregating the outputs of multiple models to make more
89 accurate predictions and it has been around for decades [12, 21, 18, 2, 3]. It is well known that
90 diversity in ensemble members’ outputs leads to better performance of the ensemble compared to the
91 performance of a single model [21] because ensemble members make independent errors [12, 11].
92 Therefore, one way to reduce DNNs’ reliance on spurious correlations is to train multiple models
93 on the same task and make them diverse in terms of errors they make so that their ensemble is less
94 dependent on such correlations.

95 2.2 Neural network ensembles that promote diversity through loss regularizers

96 Diversity in models can be induced by supplying training loss with a suitable regularizer.

97 Such regularizers can diversify models’ weights [5, 7, 34, 6], features [39, 4], input gradients
98 [29, 30, 31, 33] and outputs [25, 5, 28, 23].

99 Notably, in [5] authors showed that regularizer of a certain structure that repulses ensemble members’
100 weights or outputs leads to ensembles that provide a better approximation of Bayesian Model
101 Averaging. This idea was later extended by works that repulse ensemble members’ features [39] and
102 input gradients [33].

103 Since the ensemble performs better due to the diversity of errors that ensemble members make
104 [21] we want those members to give pairwise different outputs for the same inputs. Unfortunately,
105 diversity in weights space, input gradient space, or features space does not guarantee such property
106 without additional assumptions due to functional symmetry which means that models can be different
107 in terms of their weights or feature maps and input gradients they produce but still give the same
108 outputs for a given input. That is why we are focused on methods that diversify models’ outputs,
109 specifically [28, 23] which are state-of-the-art according to [1] and use regularizer of repulsive nature
110 conceptually similar to [5].

111 2.3 Neural network ensembles that promote diversity without modifying loss

112 In addition to loss regularizers, there were an uncountable number of different ways to induce diversity
113 in ensembles of neural networks that did not modify the training loss. The most straightforward
114 approach of independently training multiple models of the same architecture by changing only random
115 seeds is called Deep Ensemble [22] which was extended from the Bayesian perspective in [37].
116 Another solution is to construct an ensemble from models trained with different hyperparameters [36],
117 augmentations [24], or architectures [40]. More computationally efficient direction allows training
118 only one base model inducing diversity by ensembling either checkpoints saved in different local
119 minima along the training trajectory of this base model [19] or models produced by the base model
120 after applying dropout [10] or masking [9] to it. The mixture of experts paradigm can also be viewed
121 as an ensemble diversification technique [41] where diversification happens due to assigning different
122 training samples to different ensemble members.

123 Despite their conceptual simplicity Deep Ensembles [22] and ensembles of models trained with
124 different hyperparameters [36] are strong baselines for OOD detection [27] and OOD generalization
125 tasks, especially when combined with model souping techniques [38]. That is why we selected them
126 as baselines for our experiments.

127 3 Method

128 We present our main technical contributions, Scalable Ensemble Diversification (SED, §3.2) and the
129 Predictive Diversity Score (PDS, §3.3).

130 3.1 Preliminaries

131 We cover background materials before introducing our main technical contributions. We work with
132 a training set $\mathcal{D} := \{x_n, y_n\}_{n=1}^N$, which we refer to as the in-distribution (ID) dataset. For prior
133 diversification methods, we also assume the existence of a separate, unlabeled out-of-distribution

134 (OOD) dataset $\mathcal{D}^{\text{ood}} := \{x_n^{\text{ood}}\}_{n=1}^{N^{\text{ood}}}$. We write $f(\cdot, \theta)$ for a deep neural network classifier parametrized
 135 by θ . $f(x; \theta) \in \mathbb{R}^C$ indicates the logit outputs for C classes for input x . We write $p(x) :=$
 136 $\text{Softmax}(f(x)) \in [0, 1]^C$ for the probability outputs. We consider an ensemble $\{f^1, \dots, f^M\}$ of M
 137 models.

138 3.1.1 Existing ensemble diversification approach

139 We introduce an existing approach for diversifying an ensemble of models [28, 23]. Two objectives
 140 are imposed upon the ensemble of models: the main task loss and the diversification regularization.

141 For the main task, the community has focused on the classification task. The cross-entropy loss
 142 $-\log p_y(x; \theta)$ is used to train the model ensemble $\{f^1, \dots, f^M\}$ on the ID dataset \mathcal{D} :

$$\mathcal{L}_{\text{main}} = \frac{1}{MN} \sum_n \sum_m -\log p_{y_n}^m(x_n; \theta). \quad (1)$$

143 This encourages each member of the ensemble to behave similarly on the ID dataset.

144 Different diversification schemes use different diversification regularization loss \mathcal{L}_{div} applied on pairs
 145 (f^m, f^l) of ensemble members. The diversification objective is commonly optimized on the OOD
 146 dataset \mathcal{D}^{ood} to encourage the training of multiple hypotheses on the OOD samples while avoiding
 147 clashes with the main task objective. In this work, we focus on the Agree to Disagree [28] method.
 148 The diversification loss for a pair (p^m, p^l) is defined as:

$$\text{A2D}(p^m(x), p^l(x)) = -\log [p_{\hat{y}}^m(x) \cdot (1 - p_{\hat{y}}^l(x)) + (1 - p_{\hat{y}}^m(x)) \cdot p_{\hat{y}}^l(x)] \quad (2)$$

149 where $\hat{y} := \arg \max_c p_c^m(x)$ is the predicted class for the first model p^m . One may symmetrically
 150 define \hat{y} to be the prediction for the second model p^l ; in practice, it does not make a difference [28].
 151 Note that the diversification loss favors p^l to predict a lower likelihood for the prediction by p^m ,
 152 $p_{\hat{y}}^l(x)$, and vice versa. For M models in an ensemble, A2D is applied on the OOD dataset \mathcal{D}^{ood} for
 153 every pair of models (p^m, p^l) :

$$\mathcal{L}_{\text{div}} = \frac{1}{N^{\text{ood}} \cdot M(M-1)} \sum_n \sum_{m < l} \text{A2D}(p^m(x_n^{\text{ood}}), p^l(x_n^{\text{ood}})). \quad (3)$$

154 3.2 Scalable Ensemble Diversification (SED)

155 We present Scalable Ensemble Diversification (SED) that addresses the limitation of the existing
 156 ensemble diversification framework that requires a separate OOD dataset. We introduce two main
 157 components of SED: dynamic selection of OOD samples within the ID dataset (§3.2.1) and the
 158 stochastic selection of pairs to diverge in the optimization iterations (§3.2.2).

159 3.2.1 Dynamic selection of OOD samples

160 If only the ID training dataset is present, it is difficult to induce diversity in ensemble members,
 161 as they are uniformly incentivized to solve the main task objective: given x , predict y . Hence,
 162 previous approaches have introduced a qualitatively disjoint unlabeled set, which we refer to as
 163 the OOD dataset, where the ensemble members are encouraged to disagree with each other. The
 164 clear separation of ID and OOD datasets for the two objectives matters for ensuring a good balance
 165 between the main task performance and the diversity of hypotheses.

166 Previous works like Pagliardini et al. [28], Lee et al. [23] have performed experiments on small-scale
 167 datasets where factors are well-controlled and clean versions of OOD datasets are readily available.
 168 Examples include Waterbirds, Camelyon17, CelebA, MultiNLI, C-MNIST, and the Office-Home
 169 datasets. For example, for Waterbirds, the ID dataset is set as the cases where the bird’s habitat
 170 matches with the visual background and the OOD dataset corresponds to the complementary case.

171 While conceptually desirable, collecting a separate OOD dataset can be highly cumbersome and
 172 expensive. For a large-scale dataset like ImageNet, it is highly non-obvious how one could build a
 173 corresponding OOD dataset where the underlying feature-label correlations are different from the ID
 174 training dataset.

175 To address this challenge, we consider dynamically identifying an OOD subset of the ID dataset and
 176 letting the ensemble diverge on this subset. The desiderata for the identification of OOD samples

177 within the ID dataset are twofold: (a) we wish to discriminate samples where the ensemble members
 178 make mistakes and (b) we only trust the ensemble prediction for the OOD sample identification when
 179 the ensemble is sufficiently trained.

180 We define the sample-wise weight α_n on each ID sample $(x_n, y_n) \in \mathcal{D}$ that satisfy the two conditions:

$$\alpha_n := \frac{\text{CE}(f^1, \dots, f^M; x_n, y_n)}{\left(\frac{1}{|B|} \sum_{b \in B} \text{CE}(f^1, \dots, f^M; x_b, y_b)\right)^2} \quad (4)$$

181 where $\text{CE}(f^1, \dots, f^M; x_n, y_n) := \text{CE}(\frac{1}{M} \sum_m f^m(x_n), y_n)$ is the loss on the logit-averaged pre-
 182 diction and B is a minibatch that contains the sample (x_n, y_n) . α_n is a weight proportional to the
 183 ensemble loss on the sample; we thus meet the condition (a). The normalization is designed to handle
 184 the condition (b). To see this, consider the batch-wise weight

$$\alpha_B := \frac{1}{|B|} \sum_{b \in B} \alpha_b = \frac{1}{\frac{1}{|B|} \sum_b \text{CE}(f^1, \dots, f^M; x_b, y_b)}. \quad (5)$$

185 Note that α_B is now *inversely proportional* to the average cross-entropy loss of the ensemble on
 186 the batch B . Thus, the overall level of α_n for $n \in B$ is lower for earlier iterations of the ensemble
 187 training, where the predictions from the models are not trustworthy yet.

188 With this definition of sample-wise weight α_n for the diversification objective, we define the SED
 189 objective with the A2D loss for the diversification kernel:

$$\mathcal{L}_{\text{SED}} := \mathcal{L}_{\text{main}} + \frac{\lambda}{NM(M-1)} \sum_n \sum_{m < l} \text{stopgrad}(\alpha_n) \cdot \text{A2D}(p^m(x_n), p^l(x_n)), \quad (6)$$

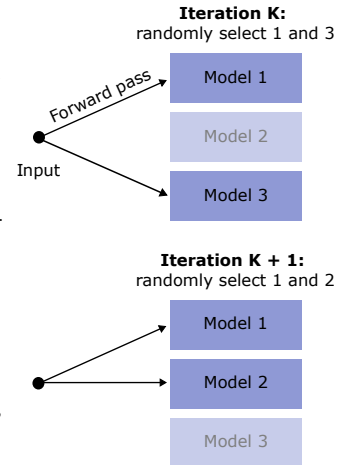
190 where $\lambda > 0$ controls the overall weight of the diversification term. Note that, compared to Equation
 191 3, this formulation does not rely on the OOD dataset \mathcal{D}^{ood} . Instead, all ID samples are treated as
 192 potential OOD samples, where their OODness is softly determined via α_n . This enables a seamless
 193 adaptation of existing ensemble diversification methods to a relaxed setting where a separate OOD
 194 dataset is unavailable.

195 3.2.2 Further tricks for scalability

196 Ensemble diversification algorithms are often based on pairwise
 197 similarities of the members. Pairwise similarity computation scales
 198 quadratically with the size of the ensemble M . The second term of
 199 Equation 6 is an example of this. This is potentially a hurdle when
 200 ensemble diversification is to be applied to $M \geq 10$, and the data
 201 and parameter sizes are in the order of millions (e.g. ImageNet).

202 We address this computational challenge by computing the summa-
 203 tion of pairwise distances as a stochastic sum. For every minibatch B
 204 of SGD iterations, we uniformly-iid sample a subset \mathcal{I} of $\{1, \dots, M\}$
 205 to compute the diversification term in Equation 6. The procedure is
 206 illustrated in the figure on the right.

207 To further speed up the SED training, we consider diversifying only
 208 a subset of layers, while freezing the other layers. In our experiments,
 209 ensemble members share the same frozen feature extractor of Deit3b
 210 [32] pretrained on ImageNet-21k [8] and we diversify only the last
 211 two layers of the models.



212 3.3 Predictive Diversity Score (PDS) for OOD Detection

213 We demonstrate several benefits of the diversified ensembles in §4. One of them is the possibility of
 214 using them for detecting OOD samples through the notion of epistemic uncertainty [13]. Given an
 215 ensemble of models, a simple baseline for OOD detection is to compute the predictive uncertainty of
 216 the Bayesian Model Averaging (BMA) by treating the ensemble members as samples of the posterior
 217 $p(\theta|\mathcal{D})$ [22, 37]:

$$\eta_{\text{BMA}} := \max_c \frac{1}{M} \sum_m p_c^m(x). \quad (7)$$

218 This notion of epistemic uncertainty does not directly exploit the potential diversity in individual
219 models of the ensemble because it averages out the predictions along the model index m .

220 We propose a novel measure for epistemic uncertainty, Predictive Diversity Score (PDS), that directly
221 measures the prediction diversity of the individual members. The formulation is given below:

$$\eta_{\text{PDS}} := \frac{1}{C} \sum_c \max_m p_c^m(x). \quad (8)$$

222 PDS is a continuous relaxation of the number of unique argmax predictions within an ensemble
223 of models. To see this, consider the special case where $p^m \in \{0, 1\}$ are one-hot vectors. Then,
224 $\max_m p_c^m(x)$ is 1 if any of m predicts c and 0 otherwise. Thus, $\sum_c \max_m p_c^m(x)$ computes the
225 number of classes that at least one of the ensemble members predicts. We show that, with our diverse
226 ensembles, PDS outperforms the DE baseline for the OOD detection task (§4.4).

227 4 Experiments

228 We verify our contributions, Scalable Ensemble Diversification (SED, §3.2) and Predictive Diversity
229 Score (PDS, §3.3), on ImageNet-scale tasks and datasets. We first verify that SED diversifies the
230 ensemble (§4.2). Then, we demonstrate the application of diversified ensemble to OOD generalization
231 (§4.3) and OOD detection (§4.4) tasks.

232 4.1 Experimental setup

233 We task the ensemble with the OOD generalization and OOD detection tasks.

234 **Training settings.** For both tasks, we train an ensemble of models with the SED framework with
235 the A2D [28] diversity regularization using AdamW optimizer [26]. We use the default settings of a
236 batch size of 16, learning rate 10^{-3} , weight decay 0.01, and the number of epochs 10. The overall
237 diversity weight λ is set to 0.1 and the stochastic pairing is done for $|\mathcal{Z}| = 2$ models for each SGD
238 batch. We use Deit3b [32] network pretrained on ImageNet21k [8] for all the experiments. Following
239 the speed-up trick in §3.2.2, we use only the last 2 layers of the network. For the in-distribution
240 (ID) dataset where the ensemble is trained to diversify, we use the training split of ImageNet with
241 $|\mathcal{D}| = 1, 281, 167$. All experiments were ran on RTX2080Ti GPUs with 12GB vRAM and 40GB
242 RAM, each experiment took from 2 to 12 hours depending on the complexity of the training.

243 **Baselines.** For naive ensemble training, we consider the *deep ensemble* [22] where each ensemble
244 member independently with different random seeds that control the weight initialization and SGD
245 batch shuffling. To match the resource usage of our SED, where we diversify only the last 2 layers
246 of the network, we consider the *shallow ensemble* variant, which is the deep ensemble where only
247 the last 2 layers are trained. We further consider a viable diversification scheme that performs deep
248 ensemble with *varying hyperparameters* [36]. In addition to that, we reimplement A2D [28] and
249 DivDis [23] algorithms and apply them without stochastic model sampling to do classification on
250 labeled samples from ImageNet-Train and disagreement on unlabeled samples from ImageNet-R.
251 For A2D we use frozen feature extractor and a parallel variant of their method which means that all
252 ensemble members are trained simultaneously and not sequentially. The computational complexity
253 of both these approaches scales quadratically with ensemble size which is why they are called Naive
254 A2D and Naive DivDis respectively.

255 **Evaluation benchmarks.** The generalization performances of the ImageNet-trained ensembles are
256 measured on multiple test datasets, ranging from the in-distribution validation split of ImageNet with
257 50,000 samples to OOD datasets like ImageNet-A (A [17], 7.5k images & 200 classes), ImageNet-R
258 (A [16], 30k images, 200 classes), ImageNet-C (C - i for corruption strength i [14], 50k images, 1k
259 classes). OpenImages-O (OI [35], 17k images, unlabeled), and iNaturalist (*iNat* [20], 10k images,
260 unlabeled). For OOD detection, we task the ensemble with the detection of the above OOD datasets
261 against the ImageNet validation split.

262 **Evaluation metrics.** For OOD generalization, we use the accuracy. For OOD detection, we use the
263 area under the ROC curve, following [15].



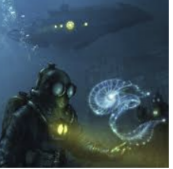
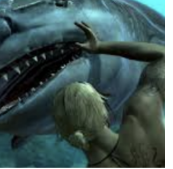

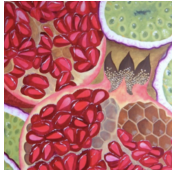




					
GT	Cowboy hat	Sea lion	Scuba diver	Great shark	Weimaraner
SED	Cowboy hat Comic book	Sea lion Otter	Scuba diver Jellyfish	Great shark Killer whale	Weimaraner Vizsla
PDS	0.300	0.300	0.294	0.292	0.292
					
GT	Pomegranate	Zebra	Pomegranate	Pomegranate	Hummingbird
SED	Pomegranate	Zebra	Pomegranate	Pomegranate	Hummingbird
PDS	0.216	0.216	0.216	0.216	0.216

Figure 2: **ImageNet-R examples leading to the greatest and least disagreement.** We show the 5 most divergent and 5 least divergent samples according to the SED ensemble. We measure the prediction diversity with the Prediction Diversity Score (PDS) in §3.3. GT refers to the ground truth category. Ensemble predictions are shown in bold, in cases when ensemble members predict classes different from the ensemble prediction we provide them on the next line with standard font.

264 4.2 Diversification

265 We start with the question of whether Scalable Ensemble Diversification (SED) truly diversify the
 266 ensemble at the ImageNet scale. To measure the diversity of the ensemble, we compute the number
 267 of unique predictions for each sample for the committee of models (#unique).

268	Method	C-1	C-5	iNat	OI
269	Deep ensemble	1.09	1.19	1.31	1.23
270	+Diverse hyperparams	1.11	1.32	1.48	1.33
271	Naive DivDis	1.04	1.14	1.19	1.16
272	Naive A2D	1.04	1.15	1.19	1.91
273	SED-A2D	5.00	5.00	4.68	4.11

274 Table 1 shows the #unique values for the IN-Val as well as multiple OOD datasets. We observe
 275 that the deep ensemble baseline does not increase the diversity dramatically (e.g. 1.09 for C-1) be-
 276 yond no-diversity values (1.0). Diversification tricks like hyperparameter diversification (1.11
 277 for C-1) or Naive A2D (1.04 for C-1) and DivDis (1.04 for C-1) do not improve the prediction di-
 278 versity dramatically. On the other hand, our SED increases the prediction diversity across the board
 (e.g. 5.00 for C-1).

Table 1: **#unique for ensembles.** We report the #unique on OOD datasets (see §4.1 for the datasets). The ensemble size M is 5 for all methods; it is the max possible #unique value.

279 Qualitative results on ImageNet-R further verify the ability of SED to diversify the ensemble (Fig-
 280 ure 2). As a measure for diversity, we use the Predictive Diversity Score (PDS) in §3.3. We observe
 281 that the samples inducing the highest diversity (high PDS scores) are indeed ambiguous: for the
 282 first image, where the “cowboy hat” is the ground truth category, we observe that “comic book” is
 283 also a valid label for the image style. On the other hand, samples with low PDS exhibit clearer
 284 image-to-category relationship.

285 4.3 OOD Generalization

286 We examine the first application of diversified ensembles: OOD generalization. We hypothesize that
 287 the superior diversification ability verified in §4.2 leads to greater OOD generalization due to the
 288 consideration of more robust hypotheses that do not rely on obvious spurious correlations.

289 **Ensemble aggregation for OOD generalization.** As a means to exploit such robust hypothe-
 290 ses, we consider 3 aggregation strategies. (1) *Oracle selection*: the best-performing individ-
 291 ual model is chosen from an ensemble [28, 30]. Final prediction is given by $f(x; \theta^{m^*})$ where

Method	M	Oracle selection					Prediction ensemble					Uniform soup				
		Val	IN-A	IN-R	C-1	C-5	Val	IN-A	IN-R	C-1	C-5	Val	IN-A	IN-R	C-1	C-5
Single model	1	85.4	37.9	44.7	75.6	38.5	85.4	37.9	44.7	75.6	38.5	85.4	37.9	44.7	75.6	38.5
Deep ensemble	5	85.4	37.9	44.9	75.7	38.6	85.4	39.9	46.3	75.7	38.6	85.3	36.7	44.6	75.5	38.3
+Diverse HPs	5	85.4	38.5	45.4	77.4	40.7	85.4	39.9	46.5	76.0	39.0	85.3	35.3	44.1	75.9	38.7
Naive DivDis	5	85.2	35.8	40.8	77.2	40.2	85.1	36.3	41.8	77.2	40.2	84.8	40.7	42.5	76.2	38.9
Naive A2D	5	85.2	36.6	44.3	77.3	40.4	85.1	37.8	45.2	77.2	40.3	84.5	39.3	45.1	75.5	39.1
SED-A2D	5	85.1	38.3	45.3	77.2	40.4	85.3	42.4	48.1	77.3	40.6	85.3	40.3	46.1	77.3	40.6
Deep ensemble	50	85.5	38.1	45.2	75.7	38.6	85.5	38.8	45.8	75.6	38.5	85.4	37.5	45.0	75.5	38.4
+Diverse HPs	50	85.5	38.5	45.6	77.5	40.8	85.5	42.5	48.5	76.0	39.0	85.4	36.4	44.8	75.9	38.8
SED-A2D	50	82.6	39.0	45.8	74.4	38.3	83.5	50.9	54.4	75.8	39.3	83.5	39.2	46.5	75.8	39.3

Table 2: **OOD generalization of ensembles.** Models are trained on the ImageNet training split. M is the ensemble size. For Naive DivDis and A2D, we use the ImageNet-R as the OOD datasets where the respective diversification objectives are applied.

292 $m^* := \arg \max_m \text{Acc}(f^m, \mathcal{D}^{\text{ood}})$. (2) *Prediction ensemble* is a vanilla prediction ensemble where
 293 the logit values are averaged: $\frac{1}{M} \sum_m f^m(x)$ [38]. (3) *Uniform soup* [38] averages the weights
 294 themselves. Final prediction is given by $f(x; \frac{1}{M} \sum_m \theta^m)$.

295 **SED improves OOD generalization for ensembles.** We show the OOD generalization performances
 296 of ensembles in Table 2, for the three ensemble prediction aggregation strategies described above. We
 297 observe that our SED framework (SED-A2D) results in superior OOD generalization performances
 298 for all three strategies. SED-A2D is particularly strong in prediction ensemble (e.g. 48.1% for $M = 5$
 299 and 54.4% for $M = 50$ on ImageNet-R) and uniform soup (e.g. 46.1% for $M = 5$ and 46.5%
 300 for $M = 50$ on ImageNet-R). We contend that the increased ensemble diversity contributes to the
 301 improvements in OOD generalization. We also remark that the SED framework (SED-A2D) envelops
 302 the performance of Naive A2D in this ImageNet-scale experiment. Together with the superiority of
 303 computational efficiency (as discussed at the end of § 4.4) of SED-A2D over the Naive A2D, this
 304 demonstrates that SED fulfills its purpose of scaling up ensemble diversification methods like A2D.

305 **Deep ensemble is a strong baseline.** We also note that deep ensemble, particularly with diverse
 306 hyperparameters, provides a strong baseline, outperforming dedicated diversification methodologies
 307 under the oracle selection strategy when $M = 5$. It also provides a good balance between ID
 308 (ImageNet validation split) and OOD generalization.

309 4.4 OOD Detection

310 We study the impact of ensemble diversifi-
 311 cation on OOD detection capabilities of an
 312 ensemble. Once an ensemble is trained, we
 313 compute the epistemic uncertainty, or like-
 314 lihood of the sample being OOD, following
 315 two schemes, η_{BMA} and η_{PDS} introduced in
 316 §3.3.

317 **SED and PDS together lead to superior**
 318 **OOD detection performances.** We show
 319 the OOD detection results in Table 3. For
 320 the BMA scores, deep ensemble remains a
 321 strong baseline. In particular, when the hy-
 322 perparameters are varied (“+Diverse HPs”),
 323 the detection AUROC reaches the maximal
 324 performances among the ensembles using
 325 the BMA scores. The quality of PDS is
 326 more sensitive to the ensemble diversity, as
 327 seen in the jump from the deep ensemble
 328 (e.g. 0.589 for OI) to the diverse-HP vari-
 329 ant (0.889). However, when the ensemble

Method	η	C-1	C-5	iNat	OI
Single model	BMA	0.615	0.833	0.958	0.909
Deep Ensemble	BMA	0.619	0.835	0.958	0.911
+Diverse HPs	BMA	0.642	0.861	0.969	0.923
Naive DivDis	BMA	0.598	0.843	0.966	0.922
Naive A2D	BMA	0.594	0.835	0.966	0.916
SED-A2D	BMA	0.641	0.845	0.960	0.915
Deep Ensemble	PDS	0.565	0.625	0.592	0.589
+Diverse HPs	PDS	0.643	0.849	0.926	0.889
Naive DivDis	PDS	0.600	0.851	0.969	0.939
Naive A2D	PDS	0.599	0.850	0.971	0.939
SED-A2D	PDS	0.686	0.896	0.977	0.941

Table 3: **OOD detection via ensembles.** For each OOD dataset (C-1, C-5, iNat, and OI), the ensembles are tasked to detect the respective OOD samples among ID samples (ImageNet validation split). We show the AUROC scores for the OOD detection task. Ensemble size is fixed at $M = 5$. η refers to the epistemic uncertainty computation framework discussed in §3.3.

330 is sufficiently diverse, such as when trained
 331 with SED-A2D, the PDS leads to high-quality OODness scores. SED-A2D with PDS achieves the
 332 best AUROC across the board, including the BMA variants.

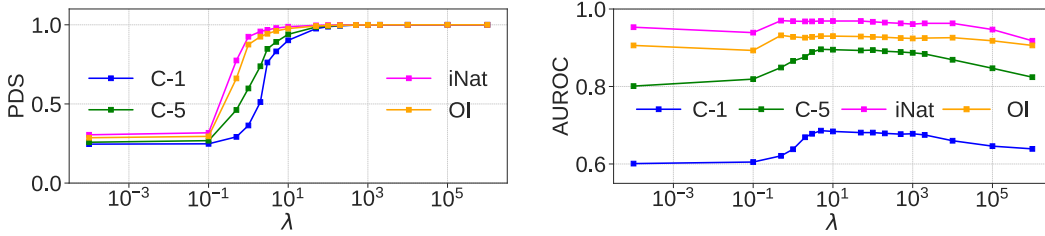


Figure 3: **Impact of diversity regulariser on OOD detection.** We show the model answer diversity, measured by PDS, and the OOD detection performance, measured by AUROC, against λ values, the loss weight for the disagreement regularizer term.

333 **Impact of diversification parameter λ .** We further study the impact of ensemble diversification
 334 on the OOD detection with the PDS estimator. In Figure 3, we observe that strengthening the
 335 diversification objective (higher λ) indeed leads to greater diversity (higher PDS), with a jump at
 336 around $\lambda \in [10^{-1}, 10^1]$. This range corresponds to the jump in the OOD detection performance
 337 (higher AUROC).

338 **Influence of ensemble size.** How ensemble size
 339 influences performance of our method? We can
 340 see that increasing ensemble size helps to im-
 341 prove AUROC for OOD detection on C-1 (Fig-
 342 ure 4). Increasing ensemble size marginally
 343 helps, but using 5 models provides already a
 344 significant improvement over the smallest pos-
 345 sible ensemble of size 2. It is also important to
 346 mention, that SED framework is computationally
 347 more efficient w.r.t. ensemble size M than Naive
 348 A2D and Naive DivDis: since we train ensembles for the fixed number of epochs, training complexity
 349 for SED is $O(1)$ thanks to stochastic model pairs selection, while for Naive A2D and Naive DivDis it
 350 is $O(M^2)$.

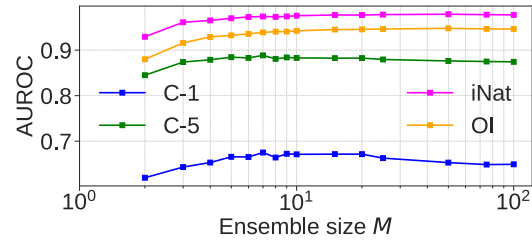


Figure 4: **Impact of ensemble size on OOD detection.**

351 5 Conclusion

352 Ensemble diversification has many implications for treating one of the ultimate goals of machine learn-
 353 ing, handling out-of-distribution (OOD) samples. By training a large number of plausible hypotheses
 354 on an in-distribution (ID) dataset, an OOD-generalizable hypothesis may appear. Moreover, the
 355 diversity of hypotheses lets us distinguish ID samples from OOD samples by measuring the degree of
 356 divergence in ensemble members’ predictions. Despite conceptual benefits, diverse-ensemble training
 357 has previously remained a lab-bound concept for several reasons. First, previous approaches required
 358 a separate OOD dataset that may nurture diverse hypotheses. Second, computational complexities of
 359 previous pairwise diversification objectives increase quadratically with the ensemble size.

360 We have addressed the challenges through the novel Scalable Ensemble Diversification (SED)
 361 framework. SED identifies the OOD-like samples from a single dataset, bypassing the need to
 362 prepare a separate OOD dataset. SED also employs a stochastic pair selection algorithm which
 363 reduces the quadratic complexity of previous approaches to a constant cost per SGD iteration. We
 364 have demonstrated good performances by SED on the OOD generalization and detection tasks, both
 365 at the ImageNet scale, a largely underexplored regime in the ensemble diversification community.
 366 In particular, for OOD detection, our novel diversity measure of Predictive Diversity Score (PDS)
 367 amplifies the benefits of diverse ensembles for OOD detection. The code to reproduce the results of
 368 our experiments will be provided with the next revision of the manuscript.

369 Limitations

370 We do not provide theoretical justification for the method. Our experiments were conducted on
 371 models with a frozen feature extractor.

References

- 373 [1] H. L. Benoit, L. Jiang, A. Atanov, O. F. Kar, M. Rigotti, and A. Zamir. Unraveling the key compo-
 374 nents of OOD generalization via diversification. In *The Twelfth International Conference on Learning*
 375 *Representations*, 2024. URL <https://openreview.net/forum?id=Lvf7GnaLru>.
- 376 [2] L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, Aug 1996. ISSN 1573-0565. doi:
 377 10.1007/BF00058655. URL <https://doi.org/10.1007/BF00058655>.
- 378 [3] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001. ISSN 1573-0565. doi: 10.1023/A:
 379 1010933404324. URL <https://doi.org/10.1023/A:1010933404324>.
- 380 [4] A. S. Chen, Y. Lee, A. Setlur, S. Levine, and C. Finn. Project and probe: Sample-efficient domain
 381 adaptation by interpolating orthogonal features. *arXiv preprint arXiv:2302.05441*, 2023.
- 382 [5] F. D’Angelo and V. Fortuin. Repulsive deep ensembles are bayesian. *Advances in Neural Information*
 383 *Processing Systems*, 34:3451–3465, 2021.
- 384 [6] A. de Mathelin, F. Deheeger, M. Mougeot, and N. Vayatis. Maximum weight entropy. *arXiv preprint*
 385 *arXiv:2309.15704*, 2023.
- 386 [7] A. de Mathelin, F. Deheeger, M. Mougeot, and N. Vayatis. Deep anti-regularized ensembles provide
 387 reliable out-of-distribution uncertainty quantification, 2023.
- 388 [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image
 389 database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.
 390 doi: 10.1109/CVPR.2009.5206848.
- 391 [9] N. Durasov, T. Bagautdinov, P. Baque, and P. Fua. Masksembles for uncertainty estimation. In *Proceedings*
 392 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13539–13548, 2021.
- 393 [10] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep
 394 learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- 395 [11] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. [http://www.](http://www.deeplearningbook.org)
 396 [deeplearningbook.org](http://www.deeplearningbook.org).
- 397 [12] L. Hansen and P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis and*
 398 *Machine Intelligence*, 12(10):993–1001, 1990. doi: 10.1109/34.58871.
- 399 [13] J. C. Helton, J. D. Johnson, and W. L. Oberkampf. An exploration of alternative approaches to the
 400 representation of uncertainty in model predictions. *Reliability Engineering & System Safety*, 85(1-3):
 401 39–71, 2004.
- 402 [14] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions
 403 and perturbations. In *International Conference on Learning Representations*, 2019. URL [https://](https://openreview.net/forum?id=HJz6tiCqYm)
 404 openreview.net/forum?id=HJz6tiCqYm.
- 405 [15] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples
 406 in neural networks. In *International Conference on Learning Representations*, 2017. URL [https://](https://openreview.net/forum?id=Hkg4TI9x1)
 407 openreview.net/forum?id=Hkg4TI9x1.
- 408 [16] D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli, M. Guo,
 409 et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings*
 410 *of the IEEE/CVF international conference on computer vision*, pages 8340–8349, 2021.
- 411 [17] D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song. Natural adversarial examples. In *Proceedings*
 412 *of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021.
- 413 [18] T. K. Ho. Random decision forests. In *Proceedings of 3rd International Conference on Document Analysis*
 414 *and Recognition*, volume 1, pages 278–282 vol.1, 1995. doi: 10.1109/ICDAR.1995.598994.
- 415 [19] G. Huang, Y. Li, G. Pleiss, Z. Liu, J. E. Hopcroft, and K. Q. Weinberger. Snapshot ensembles: Train 1, get
 416 m for free. *arXiv preprint arXiv:1704.00109*, 2017.
- 417 [20] R. Huang and Y. Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In
 418 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8710–8719,
 419 2021.
- 420 [21] A. Krogh and J. Vedelsby. Neural network ensembles, cross validation, and active learning. In
 421 G. Tesauro, D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems*,
 422 volume 7. MIT Press, 1994. URL [https://proceedings.neurips.cc/paper_files/paper/1994/](https://proceedings.neurips.cc/paper_files/paper/1994/file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf)
 423 [file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1994/file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf).
- 424 [22] B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation
 425 using deep ensembles. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
 426 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran As-
 427 sociates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf)
 428 [9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/9ef2ed4b7fd2c810847ffa5fa85bce38-Paper.pdf).

- 429 [23] Y. Lee, H. Yao, and C. Finn. Diversify and disambiguate: Out-of-distribution robustness via dis-
430 agreement. In *The Eleventh International Conference on Learning Representations*, 2023. URL
431 <https://openreview.net/forum?id=RVT0p3MwT3n>.
- 432 [24] Z. Li, I. Evtimov, A. Gordo, C. Hazirbas, T. Hassner, C. C. Ferrer, C. Xu, and M. Ibrahim. A whac-a-mole
433 dilemma: Shortcuts come in multiples where mitigating one amplifies others. In *Proceedings of the*
434 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20071–20082, 2023.
- 435 [25] Y. Liu and X. Yao. Simultaneous training of negatively correlated neural networks in an ensemble. *IEEE*
436 *Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 29(6):716–725, 1999.
- 437 [26] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *International Conference on*
438 *Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- 439 [27] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. Dillon, B. Lakshminarayanan, and
440 J. Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift.
441 *Advances in neural information processing systems*, 32, 2019.
- 442 [28] M. Pagliardini, M. Jaggi, F. Fleuret, and S. P. Karimireddy. Agree to disagree: Diversity through disagree-
443 ment for better transferability. In *The Eleventh International Conference on Learning Representations*,
444 2023. URL <https://openreview.net/forum?id=K7CbYQbyYhY>.
- 445 [29] A. Ross, W. Pan, L. Celi, and F. Doshi-Velez. Ensembles of locally independent prediction models. In
446 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5527–5536, 2020.
- 447 [30] D. Teney, E. Abbasnejad, S. Lucey, and A. van den Hengel. Evading the simplicity bias: Training a diverse
448 set of models discovers solutions with superior ood generalization. In *Proceedings of the IEEE/CVF*
449 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16761–16772, June 2022.
- 450 [31] D. Teney, M. Peyrard, and E. Abbasnejad. Predicting is not understanding: Recognizing and addressing
451 underspecification in machine learning. In S. Avidan, G. Brostow, M. Cissé, G. M. Farinella, and T. Hassner,
452 editors, *Computer Vision – ECCV 2022*, pages 458–476, Cham, 2022. Springer Nature Switzerland. ISBN
453 978-3-031-20050-2.
- 454 [32] H. Touvron, M. Cord, and H. Jégou. Deit iii: Revenge of the vit. In *European conference on computer*
455 *vision*, pages 516–533. Springer, 2022.
- 456 [33] T. Trinh, M. Heinonen, L. Acerbi, and S. Kaski. Input-gradient space particle inference for neural network
457 ensembles. In *International Conference on Learning Representations*, 2024.
- 458 [34] H. Wang and Q. Ji. Diversity-enhanced probabilistic ensemble for uncertainty estimation. In R. J. Evans
459 and I. Shpitser, editors, *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence*,
460 volume 216 of *Proceedings of Machine Learning Research*, pages 2214–2225. PMLR, 31 Jul–04 Aug
461 2023. URL <https://proceedings.mlr.press/v216/wang23c.html>.
- 462 [35] H. Wang, Z. Li, L. Feng, and W. Zhang. Vim: Out-of-distribution with virtual-logit matching. In
463 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4921–4930,
464 2022.
- 465 [36] F. Wenzel, J. Snoek, D. Tran, and R. Jenatton. Hyperparameter ensembles for robustness and uncertainty
466 quantification. *Advances in Neural Information Processing Systems*, 33:6514–6527, 2020.
- 467 [37] A. G. Wilson and P. Izmailov. Bayesian deep learning and a probabilistic perspective of generalization.
468 *Advances in neural information processing systems*, 33:4697–4708, 2020.
- 469 [38] M. Wortsman, G. Ilharco, S. Y. Gadre, R. Roelofs, R. Gontijo-Lopes, A. S. Morcos, H. Namkoong,
470 A. Farhadi, Y. Carmon, S. Kornblith, and L. Schmidt. Model soups: averaging weights of multiple
471 fine-tuned models improves accuracy without increasing inference time. In K. Chaudhuri, S. Jegelka,
472 L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference*
473 *on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 23965–23998.
474 PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/wortsman22a.html>.
- 475 [39] S. Yashima, T. Suzuki, K. Ishikawa, I. Sato, and R. Kawakami. Feature space particle inference for neural
476 network ensembles. In *International Conference on Machine Learning*, pages 25452–25468. PMLR, 2022.
- 477 [40] S. Zaidi, A. Zela, T. Elsken, C. C. Holmes, F. Hutter, and Y. Teh. Neural ensemble search for uncertainty
478 estimation and dataset shift. *Advances in Neural Information Processing Systems*, 34:7898–7911, 2021.
- 479 [41] T. Zhou, S. Wang, and J. A. Bilmes. Diverse ensemble evolution: Curriculum data-model mar-
480 riage. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Gar-
481 nett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-
482 ciates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/file/](https://proceedings.neurips.cc/paper_files/paper/2018/file/3070e6addcd702cb58de5d7897bfdae1-Paper.pdf)
483 [3070e6addcd702cb58de5d7897bfdae1-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/3070e6addcd702cb58de5d7897bfdae1-Paper.pdf).

484 **NeurIPS Paper Checklist**

485 **1. Claims**

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

488 Answer: [Yes]

489 Justification: Please refer to § 4

490 Guidelines:

- 491 • The answer NA means that the abstract and introduction do not include the claims
492 made in the paper.
- 493 • The abstract and/or introduction should clearly state the claims made, including the
494 contributions made in the paper and important assumptions and limitations. A No or
495 NA answer to this question will not be perceived well by the reviewers.
- 496 • The claims made should match theoretical and experimental results, and reflect how
497 much the results can be expected to generalize to other settings.
- 498 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
499 are not attained by the paper.

500 **2. Limitations**

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

502 Answer: [Yes]

503 Justification: Please refer to § 5

504 Guidelines:

- 505 • The answer NA means that the paper has no limitation while the answer No means that
506 the paper has limitations, but those are not discussed in the paper.
- 507 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 508 • The paper should point out any strong assumptions and how robust the results are to
509 violations of these assumptions (e.g., independence assumptions, noiseless settings,
510 model well-specification, asymptotic approximations only holding locally). The authors
511 should reflect on how these assumptions might be violated in practice and what the
512 implications would be.
- 513 • The authors should reflect on the scope of the claims made, e.g., if the approach was
514 only tested on a few datasets or with a few runs. In general, empirical results often
515 depend on implicit assumptions, which should be articulated.
- 516 • The authors should reflect on the factors that influence the performance of the approach.
517 For example, a facial recognition algorithm may perform poorly when image resolution
518 is low or images are taken in low lighting. Or a speech-to-text system might not be
519 used reliably to provide closed captions for online lectures because it fails to handle
520 technical jargon.
- 521 • The authors should discuss the computational efficiency of the proposed algorithms
522 and how they scale with dataset size.
- 523 • If applicable, the authors should discuss possible limitations of their approach to
524 address problems of privacy and fairness.
- 525 • While the authors might fear that complete honesty about limitations might be used by
526 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
527 limitations that aren't acknowledged in the paper. The authors should use their best
528 judgment and recognize that individual actions in favor of transparency play an impor-
529 tant role in developing norms that preserve the integrity of the community. Reviewers
530 will be specifically instructed to not penalize honesty concerning limitations.

531 **3. Theory Assumptions and Proofs**

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

534 Answer: [NA]

535 Justification: The paper contains no theoretical results.

536 Guidelines:

- 537 • The answer NA means that the paper does not include theoretical results.
- 538 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 539 referenced.
- 540 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 541 • The proofs can either appear in the main paper or the supplemental material, but if
- 542 they appear in the supplemental material, the authors are encouraged to provide a short
- 543 proof sketch to provide intuition.
- 544 • Inversely, any informal proof provided in the core of the paper should be complemented
- 545 by formal proofs provided in appendix or supplemental material.
- 546 • Theorems and Lemmas that the proof relies upon should be properly referenced.

547 4. Experimental Result Reproducibility

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

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

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

551 Answer: [Yes]

552 Justification: Please refer to § 4

553 Guidelines:

- 554 • The answer NA means that the paper does not include experiments.
- 555 • If the paper includes experiments, a No answer to this question will not be perceived
- 556 well by the reviewers: Making the paper reproducible is important, regardless of
- 557 whether the code and data are provided or not.
- 558 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 559 to make their results reproducible or verifiable.
- 560 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 561 For example, if the contribution is a novel architecture, describing the architecture fully
- 562 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 563 be necessary to either make it possible for others to replicate the model with the same
- 564 dataset, or provide access to the model. In general, releasing code and data is often
- 565 one good way to accomplish this, but reproducibility can also be provided via detailed
- 566 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 567 of a large language model), releasing of a model checkpoint, or other means that are
- 568 appropriate to the research performed.
- 569 • While NeurIPS does not require releasing code, the conference does require all submis-
- 570 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 571 nature of the contribution. For example
 - 572 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 573 to reproduce that algorithm.
 - 574 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 575 the architecture clearly and fully.
 - 576 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 577 either be a way to access this model for reproducing the results or a way to reproduce
 - 578 the model (e.g., with an open-source dataset or instructions for how to construct
 - 579 the dataset).
 - 580 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 581 authors are welcome to describe the particular way they provide for reproducibility.
 - 582 In the case of closed-source models, it may be that access to the model is limited in
 - 583 some way (e.g., to registered users), but it should be possible for other researchers
 - 584 to have some path to reproducing or verifying the results.

585 5. Open access to data and code

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

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

588 material?

589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639

Answer: [Yes]

Justification: Code will be available soon, please refer to § 4.1.

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: please refer to § 4.1.

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: [No]

Justification: Error bars are not reported because their magnitude was below the rounding error or roughly around it for the majority of experiments.

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).

- 640 • It should be clear whether the error bar is the standard deviation or the standard error
641 of the mean.
- 642 • It is OK to report 1-sigma error bars, but one should state it. The authors should
643 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
644 of Normality of errors is not verified.
- 645 • For asymmetric distributions, the authors should be careful not to show in tables or
646 figures symmetric error bars that would yield results that are out of range (e.g. negative
647 error rates).
- 648 • If error bars are reported in tables or plots, The authors should explain in the text how
649 they were calculated and reference the corresponding figures or tables in the text.

650 8. Experiments Compute Resources

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

654 Answer: [Yes]

655 Justification: please refer to § 4.1.

656 Guidelines:

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

665 9. Code Of Ethics

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

668 Answer: [Yes]

669 Justification: we followed the Code to the best of our knowledge.

670 Guidelines:

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

676 10. Broader Impacts

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

679 Answer: [NA]

680 Justification: We believe that this work has no societal impact.

681 Guidelines:

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

- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 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).

704 11. Safeguards

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

708 Answer: [NA]

709 Justification: We believe that our paper does not pose such risks as we train models for
710 ImageNet classification.

711 Guidelines:

- 712
- 713
- 714
- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 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.

722 12. Licenses for existing assets

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

726 Answer: [No]

727 Justification: we were unable to find the license for the dataset we used.

728 Guidelines:

- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 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 has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - 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.

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

744 **13. New Assets**

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

747 Answer: [NA]

748 Justification: the paper does not release new assets.

749 Guidelines:

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

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

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

762 Answer: [NA]

763 Justification: the paper does not involve crowdsourcing nor research with human subjects.

764 Guidelines:

- 765 • The answer NA means that the paper does not involve crowdsourcing nor research with
766 human subjects.
- 767 • Including this information in the supplemental material is fine, but if the main contribu-
768 tion of the paper involves human subjects, then as much detail as possible should be
769 included in the main paper.
- 770 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
771 or other labor should be paid at least the minimum wage in the country of the data
772 collector.

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

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

779 Answer: [NA]

780 Justification: the paper does not involve crowdsourcing nor research with human subjects.

781 Guidelines:

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