
Test-Time Adaptation to Distribution Shift by Confidence Maximization and Input Transformation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep neural networks often exhibit poor performance on data that is unlikely under
2 the train-time data distribution, for instance data affected by corruptions. Previous
3 works demonstrate that test-time adaptation to data shift, for instance using entropy
4 minimization [1], effectively improves performance on such shifted distributions.
5 This paper focuses on the fully test-time adaptation setting, where only unlabeled
6 data from the target distribution is required. This allows adapting arbitrary pre-
7 trained networks. Specifically, we propose a novel loss that improves test-time
8 adaptation by addressing both premature convergence and instability of entropy
9 minimization. This is achieved by replacing the entropy by a non-saturating surro-
10 gate and adding a diversity regularizer based on batch-wise entropy maximization
11 that prevents convergence to trivial collapsed solutions. Moreover, we propose
12 to prepend an input transformation module to the network that can partially undo
13 test-time distribution shifts. Surprisingly, this preprocessing can be learned solely
14 using the fully test-time adaptation loss in an end-to-end fashion without any target
15 domain labels or source domain data. We show that our approach outperforms
16 previous work in improving the robustness of publicly available pretrained image
17 classifiers to common corruptions on such challenging benchmarks as ImageNet-C.

18 1 Introduction

19 Deep neural networks achieve impressive performance on test data, which has the same distribution
20 as the training data. Nevertheless, they often exhibit a large performance drop on test (target) data
21 which differs from training (source) data; this effect is known as data shift [2] and can be caused for
22 instance by image corruptions. There exist different methods to improve the robustness of the model
23 during training [3, 4, 5]. However, generalization to different data shifts is limited since it is infeasible
24 to include sufficiently many augmentations during training to cover the excessively wide range of
25 potential data shifts [6]. Alternatively, in order to generalize to the data shift at hand, the model can be
26 adapted during test-time. Unsupervised domain adaptation methods such as [7] use both source and
27 target data to improve the model performance during test-time. In general source data might not be
28 available during inference time, e.g., due to legal constraints (privacy or profit). Therefore we focus
29 on the *fully test-time adaptation* setting [1]: the model is adapted to the target data given only the
30 arbitrarily pretrained model parameters and the unlabeled target data that share the same label space
31 as source data. We extend the work of Wang et al. [1] by introducing a novel loss function, using
32 a diversity regularizer, and prepending a parametrized input transformation module to the network.
33 We show that our approach outperform previous works and make pretrained models robust against
34 common corruptions on image classification benchmarks as ImageNet-C [8] and ImageNet-R [9].

35 Sun et al. [10] investigate test-time adaptation using a self-supervision task. Wang et al. [1] and
36 Liang et al. [11] use the entropy minimization loss that uses maximization of prediction confidence as
37 self-supervision signal during test-time adaptation. Wang et al. [1] has shown that such loss performs

38 better adaptation than a proxy task [10]. When using entropy minimization, however, high confidence
39 predictions do not contribute to the loss significantly anymore and thus provide little self-supervision.
40 This is a drawback since high-confidence samples provide the most trustworthy self-supervision.
41 We mitigate this by introducing two novel loss functions that ensure that gradients of samples with
42 high confidence predictions do not vanish and learning based on self-supervision from these samples
43 continues. Our losses do not focus on minimizing entropy but on minimizing the *negative log*
44 *likelihood ratio* between classes; the two variants differ in using either soft or hard pseudo-labels. In
45 contrast to entropy minimization, the proposed loss functions provide non-saturating gradients, even
46 when there are high confident predictions. We refer to Figure 1 for an illustration of the losses and the
47 resulting gradients. Using these new loss functions, we are able to improve the network performance
48 under data shifts in fully test-time adaptation.

49 In general, self-supervision by confidence maximization can lead to collapsed trivial solutions, which
50 make the network to predict only a single or a set of classes independent of the input. To overcome
51 this issue a *diversity regularizer* [11, 12] can be used, that acts on a batch of samples. It encourages
52 the network to make different class predictions on different samples. We extend the regularizer by
53 including a moving average, in order to include the history of the previous batches and show that this
54 stabilizes the adaptation of the network to unlabeled test samples. Furthermore we also introduce a
55 parametrized *input transformation module*, which we prepend to the network. The module is trained
56 in a fully test-time adaptation manner using the proposed loss function, i. e. without the need of any
57 target domain labels or source data. It aims to partially undo the data shift at hand. This helps to
58 further improve the performance on image classification benchmark with corruptions.

59 Since our method does not change the training process, it allows to use any pretrained models. This
60 is beneficial because any good performing pretrained network can be readily reused, e.g., a network
61 trained on some proprietary data not available to the public. We show, that our method significantly
62 improves performance on models that are trained on clean ImageNet data such as a ResNet50 [13],
63 as well as robust models such as ResNet50 models trained using DeepAugment+AugMix [9].

64 In summary our main contributions are as follows: we propose non-saturating losses based on the
65 negative log likelihood ratio, such that gradients from high confidence predictions still contribute to
66 test-time adaptation. We extend the diversity regularizer that acts on a batch of samples to a moving
67 average version, which includes the history of the previous batch samples. This prevents the network
68 from collapsing to trivial solutions. Furthermore we also introduce an input transformation module,
69 which partially undoes the data shift at hand. We show that the performance of different pretrained
70 models can be significantly improved on challenging benchmarks like ImageNet-C and ImageNet-R.

71 2 Related work

72 **Common image corruptions** are potentially stochastic image transformations motivated by real-
73 world effects that can be used for evaluating a model’s robustness. One such benchmark, ImageNet-C
74 [8], contains simulated corruptions such as noise, blur, weather effects, and digital image transforma-
75 tions. Additionally, Hendrycks et al. [9] proposed three data sets containing real-world distribution
76 shifts, including Imagenet-R. The ImageNet-C have been further extended to MNIST [14], several
77 object detection datasets [15], and image segmentation [16], reflecting the interest of the robustness
78 community. Most proposals for improving robustness involve special training protocols, requiring
79 time and additional resources. This includes data augmentation like Gaussian noise [17, 18, 9],
80 CutMix [19], AugMix [4], training on stylized images [3, 20] or against adversarial noise distribu-
81 tions [21]. Mintun et al. [22] pointed out that many improvements on ImageNet-C are due to data
82 augmentations which are too similar to the test corruptions, that is: overfitting to ImageNet-C occurs.
83 Thus, the model might be less robust to corruptions not included in the test set of ImageNet-C.

84 **Unsupervised domain adaptation** methods train a joint model of the source and target domain by
85 cross-domain losses, with the hope to find more general and robust features. These losses optimize
86 feature alignment [23, 24] between domains, adversarial invariance [25, 5, 26, 27], shared proxy
87 tasks [28] or adapting the entropy minimization via an adversarial loss [7]. While these approaches
88 are effective, they require explicit access to source and target data at the same time, which may not
89 always be feasible. Our approach works with any pretrained model and only needs target data.

90 **Test-time adaptation** (also termed *source free adaptation* in some literature) is a setting, when
91 training (source) data is unavailable at test-time. Several works use generative models [29, 30, 31, 32]

92 for the source free adaptation and require several thousand epochs to adapt to the target data [30, 32].
 93 Besides, there is another line of work [10, 33, 34, 35, 1] that interpret the common corruptions as
 94 data shift and aim to improve the model robustness against these corruptions with efficient test-time
 95 adaptation strategy to facilitate online adaptation. Such setting refrain the usage of generative models
 96 or methods that require larger number of adaptation steps. Our work also falls in this line of research
 97 and aims to test-time adapt the model to common corruptions with less computational overhead.

98 Sun et al. [10] update feature extractor parameters at test-time via a self-supervised proxy task
 99 (predicting image rotations). However, Sun et al. [10] alter the training procedure by including the
 100 proxy loss into the optimization objective as well, hence arbitrary pretrained models cannot be used
 101 directly for test-time adaptation. Inspired by the domain adaptation strategies [36, 37], several works
 102 [33, 34, 35] replace the estimates of Batch Normalization (BN) activation statistics with the statistics
 103 of the corrupted test images. Fully test time adaptation, studied by Wang et al. [1] (TENT) uses
 104 entropy minimization to update the channel-wise affine parameters of BN layers on corrupted data
 105 along with the batch statistics estimates. SHOT [11] also uses entropy minimization and a diversity
 106 regularizer to avoid collapsed solutions. SHOT modifies the model from the standard setting by
 107 adopting weight normalization at the fully connected classifier layer during training to facilitate their
 108 pseudo labeling technique. Hence, SHOT is not readily applicable to arbitrary pretrained models.

109 We show that pure entropy minimization [1, 11] results in vanishing gradients for high confidence
 110 predictions, thus inhibiting learning. Our work addresses this issue by proposing a novel non-
 111 saturating loss, that provides non-vanishing gradients for high confidence predictions. We show
 112 that our proposed loss function improves the network performance after test-time adaptation. In
 113 particular, performance on corruptions of higher severity improves significantly. Furthermore, we
 114 add and extend the diversity regularizer [11, 12] to avoid collapse to trivial, high confidence solutions.
 115 Note that the existing diversity regularizers [11, 12] act on a batch of samples, hence the number of
 116 classes has to be smaller than the batch size. We mitigate this problem by extending the regularizer
 117 to a running average version. Prior work [5, 38, 39] transformed inputs by an additional module
 118 to overcome domain shift, obtain robust models, and also to learn to resize. In our work, we also
 119 prepend an input transformation module to the model, but in contrast to former works, this module is
 120 trained purely at test-time to partially undo the data shift at hand and thus aids the adaptation.

121 3 Method

122 We propose a novel method for fully test-time adaption. For this, we assume that a neural network
 123 f_θ with parameters θ is available that was trained on data from some distribution \mathcal{D} , as well a set of
 124 (unlabeled) samples $X \sim \mathcal{D}'$ from a target distribution $\mathcal{D}' \neq \mathcal{D}$ (importantly, no samples from \mathcal{D} are
 125 required). We frame fully test-time adaption as a two-step process: (i) Generate a novel network g_ϕ
 126 based on f_θ , where ϕ denotes the parameters that are adapted. A simple variant for this is $g = f$ and
 127 $\phi \subseteq \theta$ [1]. However, we propose a more expressive and flexible variant in Section 3.1. (ii) Adapt the
 128 parameters ϕ of g on X using an unsupervised loss function L . We propose two novel losses L_{slr}
 129 and L_{hlr} in Section 3.2 that have non-vanishing gradients for high-confidence self-supervision.

130 3.1 Input Transformation

131 We propose to define the adaptable model as $g = f \circ d$. That is: we prepend a trainable network d
 132 to f . The motivation for the additional component d is to increase expressivity of g such that it can
 133 learn to (partially) undo the domain shift $\mathcal{D} \rightarrow \mathcal{D}'$.

134 Specifically, we choose $d(x) = \gamma \cdot [\tau x + (1 - \tau)r_\psi(x)] + \beta$, where $\tau \in \mathbb{R}$, $(\beta, \gamma) \in \mathbb{R}^{n_{in}}$ with
 135 n_{in} being the number of input channels, r_ψ being a network with identical input and output shape,
 136 and \cdot denoting elementwise multiplication. Specifically, β and γ implement a channel-wise affine
 137 transformation and τ implements a convex combination of unchanged input and the transformed input
 138 $r_\psi(x)$. By choosing $\tau = 1$, $\gamma = \mathbf{1}$, and $\beta = \mathbf{0}$, we ensure $d(x) = x$ and thus $g = f$ at initialization.
 139 In principle, r_ψ can be chosen arbitrarily. In this work, we choose r_ψ as a simple stack of 3×3
 140 convolutions, group normalization, and ReLUs (for details, we refer to the appendix). However,
 141 exploring other choices would be an interesting avenue for future work.

142 Importantly, while the motivation for d is to learn to partially undo a domain shift $\mathcal{D} \rightarrow \mathcal{D}'$, we train
 143 d end-to-end in the fully test-time adaptation setting on data $X \sim \mathcal{D}'$, without any access to samples

144 from the source domain \mathcal{D} , based on the losses proposed in Section 3.2. The modulation parameters
 145 of g_ϕ are $\phi = (\beta, \gamma, \tau, \psi, \theta')$, where $\theta' \subseteq \theta$. That is, we adapt only a subset of the parameters θ of
 146 the pretrained network f . We largely follow Wang et al. [1] in adapting only the affine parameters of
 147 normalization layers in f while keeping parameters of convolutional kernels unchanged. Additionally,
 148 batch normalization statistics (if any) are adapted to the target distribution.

149 Please note that the proposed method is applicable to any pretrained network that contains normaliza-
 150 tion layers with a channel-wise affine transformation. Even for networks that do not come with such
 151 affine transformation layers, one can add affine transformation layers into f that are initialized to
 152 identity as part of model augmentation.

153 3.2 Adaptation Objective

154 We propose a loss function $L = L_{\text{div}} + \delta L_{\text{conf}}$ for fully test-time network adaptation that consists of
 155 two components: (i) a term L_{div} that encourages predictions of the network over the adaptation dataset
 156 X that match a target distribution $p_{\mathcal{D}'}(y)$. This can help avoiding test-time adaptation collapsing
 157 to too narrow distributions such as always predicting the same or very few classes. If $p_{\mathcal{D}'}(y)$ is
 158 (close to) uniform, it acts as a diversity regularizer. (ii) A term L_{conf} that encourages high confidence
 159 prediction on individual datapoints. We note that test-time entropy minimization (TENT) [1] fits into
 160 this framework by choosing $L_{\text{div}} = 0$ and L_{conf} as the entropy.

161 3.2.1 Class Distribution Matching L_{div}

162 Assuming knowledge of the class distribution $p_{\mathcal{D}'}(y)$ on the target domain \mathcal{D}' , we propose to add a
 163 term to the loss that encourages the empirical distribution of (soft) predictions of g_ϕ on X to match
 164 this distribution. Specifically, let $\hat{p}_{g_\phi}(y)$ be an estimate of the distribution of (soft) predictions of g_ϕ .
 165 We use the Kullback-Leibler divergence $L_{\text{div}} = D_{KL}(\hat{p}_{g_\phi}(y) || p_{\mathcal{D}'}(y))$ as loss term. In a special case
 166 of $p_{\mathcal{D}'}(y)$ being a uniform distribution over the classes, this corresponds to maximizing the entropy
 167 $H(\hat{p}_{g_\phi}(y))$. Similar assumption has been made in SHOT [11] to circumvent the collapsed solutions.

168 Since the estimate $\hat{p}_{g_\phi}(y)$ depends on ϕ , which is continuously adapted, it needs to be re-estimated
 169 on a per-batch level. Since re-estimating $\hat{p}_{g_\phi}(y)$ from scratch would be computational expensive, we
 170 propose to use a running estimate that tracks the changes of ϕ as follows: let $p_{t-1}(y)$ be the estimate at
 171 iteration $t-1$ and $p_t^{\text{emp}} = \frac{1}{n} \sum_{k=1}^n \hat{y}^{(k)}$, where $\hat{y}^{(k)}$ are the predictions (confidences) of g_ϕ on a mini-
 172 batch of n inputs $x^{(k)} \sim X$. We update the running estimate via $p_t(y) = \kappa \cdot p_{t-1}(y) + (1 - \kappa) \cdot p_t^{\text{emp}}$.
 173 The loss becomes $L_{\text{div}} = D_{KL}(p_t(y) || p_{\mathcal{D}'}(y))$ accordingly. We use $\kappa = 0.9$ in the experiments.

174 3.2.2 Confidence Maximization L_{conf}

175 We motivate our choice of L_{conf} step-by-step from the (unavailable) supervised cross-entropy loss:
 176 for this, let $\hat{y} = g_\phi(x)$ be the predictions (confidences) of model g_ϕ and $H(\hat{y}, y^r) = -\sum_c y_c^r \log \hat{y}_c$
 177 be the cross-entropy between prediction \hat{y} and some reference y^r . Moreover, let the last layer of g be
 178 a softmax activation layer softmax . That is $\hat{y} = \text{softmax}(o)$, where o are the network's logits. We
 179 note that we can rewrite the cross-entropy loss in terms of the logits o and a one-hot reference y^r as
 180 follows: $H(\text{softmax}(o), y^r) = -o_{c^r} + \log \sum_{i=1}^{n_{cl}} e^{o_i}$ where c^r is the index of the 1 in y^r and n_{cl} is
 181 the number of classes.

182 In the case of labels being available for the target domain (which we do not assume) in the form of a
 183 one-hot encoded reference y_t for data x_t , one could use the *supervised cross-entropy loss* by setting
 184 $y^r = y_t$ and using $L_{\text{sup}}(\hat{y}, y^r) = H(\hat{y}, y^r) = H(\hat{y}, y_t)$. Since fully test-time adaptation assumes
 185 no label information being available, the supervised cross-entropy loss is not applicable and other
 186 options for y^r need to be used.

187 One option are (hard) *pseudo-labels*. That is, one defines the reference y^r based on the network pre-
 188 dictions \hat{y} via $y^r = \text{onehot}(\hat{y})$, where onehot creates a one-hot reference with the 1 corresponding to
 189 the class with maximal confidence in \hat{y} . This results in $L_{\text{pl}}(\hat{y}) = H(\hat{y}, \text{onehot}(\hat{y})) = -\log \hat{y}_{c^*}$, with
 190 $c^* = \arg \max \hat{y}$. One disadvantage with this loss is that the (hard) pseudo-labels ignore uncertainty
 191 in the network predictions during self-supervision. This results in large gradient magnitudes with
 192 respect to the logits $|\frac{\partial L_{\text{pl}}}{\partial o_{c^*}}|$ being generated in situations where the network is highly unconfident (see

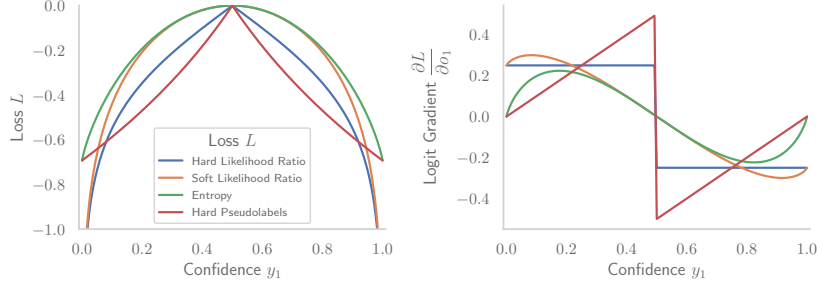


Figure 1: *Illustration of different losses for confidence maximization.* Losses (left, shifted such that maxima of all losses are at 0) and the resulting gradients with respect to the first logit (right) as a function of the first classes confidence are shown for the case of a binary classification problem. Both *entropy* and *hard pseudo-labels* have vanishing gradients for high confidence predictions. Accordingly, both have maximum gradient amplitude for low-confidence self-supervision, with this effect being stronger for the hard pseudo-labels. *Hard Likelihood Ratio* has constant gradient amplitude for any confidence and thus takes into account low- and high-confidence self-supervision equally. *Soft Likelihood Ratio* also shows non-vanishing (albeit non-maximum) gradients for high-confidence self-supervision and additionally produces small gradient amplitudes from low-confidence self-supervision. Since the likelihood ratio-based losses are unbounded, the design of the model needs to ensure that logits cannot grow unbounded.

193 Figure 1). This is undesirable since it corresponds to the network being affected most by data points
 194 where the network’s self-supervision is least reliable.

195 An alternative is to use soft pseudo-labels, that is $y^r = \hat{y}$. This takes uncertainty in network
 196 predictions into account during self-labelling and results in the *entropy minimization* loss of TENT
 197 [1]: $L_{ent}(\hat{y}) = H(\hat{y}, \hat{y}) = H(\hat{y}) = -\sum_c \hat{y}_c \log \hat{y}_c$. However, also for the entropy the logits’
 198 gradient magnitude $|\frac{\partial L_{ent}}{\partial o}|$ goes to 0 when one of the entries in \hat{y} goes to 1 (see Figure 1). For
 199 a binary classification task, for instance, the maximal logits’ gradient amplitude is obtained for
 200 $\hat{y} \approx (0.82, 0.18)$. This implies that during later stages of test-time adaptation where many predictions
 201 typically already have very high confidence, i. e. above 0.82, gradients are also dominated by
 202 datapoints with relative low confidence in self-supervision.

203 While both hard and soft pseudo-labels are clearly motivated, they are not optimal in conjunction with
 204 a gradient-based optimizer since the self-supervision from low confidence predictions dominates (at
 205 least during later stages of training). To address this issue, we propose two losses that are analogous
 206 to L_{pl} and L_{ent} , but are not based on the cross-entropy H but instead on the negative log likelihood
 207 ratios

$$R(\hat{y}, y^r) = -\sum_c y_c^r \log \frac{\hat{y}_c}{\sum_{i \neq c} \hat{y}_i} = -\sum_c y_c^r (\log \hat{y}_c - \log \sum_{i \neq c} \hat{y}_i) = H(\hat{y}, y^r) + \sum_c y_c^r \log \sum_{i \neq c} \hat{y}_i$$

208 Note that while the entropy H is lower bounded by 0, R can get arbitrary small if $y_c^r \rightarrow 1$ and the
 209 sum $\sum_{i \neq c} \hat{y}_i \rightarrow 0$ and thus $\log \sum_{i \neq c} \hat{y}_i \rightarrow -\infty$. This property will induce non-vanishing gradients
 210 for high confidence predictions.

211 The first loss we consider is the *hard likelihood ratio* loss that is defined similarly to the hard
 212 pseudo-labels loss L_{pl} :

$$L_{hlr}(\hat{y}) = R(\hat{y}, \text{onehot}(\hat{y})) = -\log\left(\frac{\hat{y}_{c^*}}{\sum_{i \neq c^*} \hat{y}_i}\right) = -\log\left(\frac{e^{o_{c^*}}}{\sum_{i \neq c^*} e^{o_i}}\right) = -o_{c^*} + \log \sum_{i \neq c^*} e^{o_i},$$

213 where $c^* = \arg \max \hat{y}$. We note that $\frac{\partial L_{hlr}}{\partial o_{c^*}} = -1$, thus also high-confidence self-supervision
 214 contributes equally to the maximum logits’ gradients. This loss was also independently proposed as
 215 negative log likelihood ratio loss by Yao et al. [40] as a replacement to the fully-supervised cross
 216 entropy loss for classification task. However, to the best of our knowledge, we are the first to motivate
 217 and identify the advantages of this loss for self-supervised learning and test-time adaptation due to its
 218 non-saturating gradient property.

219 In addition to L_{hlr} , we also account for uncertainty in network predictions during self-labelling in a
 220 similar way as for the entropy loss L_{ent} , and propose the *soft likelihood ratio* loss:

$$\begin{aligned}
 L_{slr}(\hat{y}) &= R(\hat{y}, \hat{y}) = - \sum_c \hat{y}_c \cdot \log\left(\frac{\hat{y}_c}{\sum_{i \neq c} \hat{y}_i}\right) &= - \sum_c \hat{y}_c \log\left(\frac{e^{o_c}}{\sum_{i \neq c} e^{o_i}}\right) \\
 &= \sum_c \hat{y}_c (-o_c + \log \sum_{i \neq c} e^{o_i})
 \end{aligned}$$

221 We note that as $\hat{y}_{c^*} \rightarrow 1$, $L_{slr}(\hat{y}) \rightarrow L_{hlr}(\hat{y})$. Thus the asymptotic behavior of the two likelihood
 222 ratio losses for high confidence predictions is the same. However, the soft likelihood ratio loss
 223 creates lower amplitude gradients for low confidence self-supervision. We provide illustrations of the
 224 discussed losses and the resulting logits’ gradients in Figure 1.

225 We note that both likelihood ratio losses would typically encourage the network to simply scale
 226 its logits larger and larger, since this would reduce the loss even if the ratios between the logits
 227 remain constant. However, when finetuning an existing network and restricting the layers that are
 228 adapted such that the logits remain approximately scale-normalized, these losses can provide a
 229 useful and non-vanishing gradient signal for network adaptation. We achieve this approximate
 230 scale normalization by freezing the top layers of the respective networks. In this case, normalization
 231 layers such as batch normalization prohibit “logit explosion”. However, predicted confidences can
 232 presumably become overconfident; calibrating confidences in a self-supervised test-time adaptation
 233 setting is an open and important direction for future work.

234 4 Experimental settings

235 **Datasets** We evaluate our method on image classification datasets for corruption robustness and
 236 domain adaptation. We evaluate on the challenging benchmark ImageNet-C [8], which includes a
 237 wide variety of 15 different synthetic corruptions with 5 severity levels that attribute to data shift.
 238 This benchmark also includes 4 additional corruptions as validation data. For domain adaptation, we
 239 choose ImageNet trained models to adapt to ImageNet-R proposed by Hendrycks et al. [9]. This
 240 dataset contains various naturally occurring artistic renditions of object classes from the original
 241 ImageNet. ImageNet-R comprises 30,000 image renditions for 200 ImageNet classes. Please refer
 242 Sec. A.5 for the experiments on other domain adaptation datasets VisDA-C [41], Office-Home [42].

243 **Models** Our method operates in a fully test-time adaptation setting that allows us to use any arbitrary
 244 pretrained model. We use publicly available ImageNet pretrained models ResNet50, DenseNet121,
 245 ResNeXt50, MobileNetV2 from torchvision [43]. We also test on a robust ResNet50 model trained
 246 using DeepAugment+AugMix¹ [9].

247 **Baseline for fully test-time adaptation** Since TENT from Wang et al. [1] outperformed competing
 248 methods and fits the fully test-time adaptation setting, we consider it as a baseline and compare
 249 our results to this approach. Similar to TENT, we also adapt model features by estimating the
 250 normalization statistics and optimize only the channel-wise affine parameters on the target distribution.

251 **Settings** We conduct test-time adaptation on a target distribution for 5 epochs with batch size 64 and
 252 use the Adam optimizer with cosine decay scheduler of the learning rate with initial value 0.0006.
 253 We set the weight of L_{conf} in our loss function to $\delta = 0.025$ and $\kappa = 0.9$ in the running estimate
 254 $p_t(y)$ of L_{div} (we investigate the effect of κ in the Sec. A.3). Similar to SHOT [11], we also choose
 255 the target distribution $p_{\mathcal{D}'}(y)$ in L_{div} as a uniform distribution over the available classes. We found
 256 that the models converge during 3 to 5 epochs and do not improve further.

257 For TENT, we use SGD with momentum 0.9 at constant learning rate 0.00025 with batch size 64.
 258 These values correspond to the ones of Wang et al. [1]; alternative settings of optimizer and learning
 259 rates for TENT did not improve performance. TENT is originally optimized only for 1 epoch. For a
 260 fair comparison to our method, we optimize TENT also for 5 epochs. Similar to Wang et al. [1], we
 261 also control for ordering by data shuffling and sharing the order across the methods.

262 Note that all the hyperparameter settings are tuned solely on the validation corruptions of ImageNet-C
 263 that are disjoint from the test corruptions. As discussed in Section 3.2.2, we freeze all trainable
 264 parameters in the top layers of the networks to prohibit “logit explosion”. Note that normalization

¹From <https://github.com/hendrycks/imagenet-r>. Owner permitted to use it for research/commercial purposes.

Table 1: Test-time adaptation of ResNet50 on ImageNet-C at highest severity level 5. Ground truth labels are used to adapt the model in supervised manner to obtain empirical upper bound performance.

Method	Gauss	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
No Adaptation	2.44	2.99	1.96	17.92	9.82	14.78	22.50	16.89	23.31	24.43	58.93	5.43	16.95	20.61	31.65
Pseudo Labels	2.44	2.99	1.96	17.92	9.82	14.78	22.50	16.89	23.31	24.43	58.93	5.43	16.95	20.61	31.65
Epoch 1															
TENT	32.70	35.34	35.11	32.79	31.80	47.22	53.02	51.82	43.42	60.44	68.82	27.53	58.47	61.63	55.98
TENT+	33.96	36.66	35.75	33.70	33.33	47.73	53.22	52.16	44.79	60.62	68.91	35.60	58.81	61.82	56.23
HLR (ours)	38.39	41.11	40.28	38.25	38.18	51.63	55.55	55.45	48.96	62.19	68.17	49.47	60.34	62.51	57.42
SLR (ours)	39.51	42.09	41.58	39.35	39.02	52.67	55.80	55.92	49.64	62.62	68.47	50.27	60.80	63.01	57.80
Epoch 5															
TENT	16.04	23.22	25.85	19.05	17.40	49.02	52.78	52.72	34.31	61.19	68.54	1.26	59.26	62.15	56.17
TENT+	33.97	37.95	36.93	32.69	33.36	51.42	54.33	54.55	45.80	62.09	69.03	24.08	60.36	63.10	57.21
HLR (ours)	41.37	44.04	43.68	41.74	41.09	54.26	56.43	57.03	50.81	63.05	68.29	50.98	61.15	63.08	58.13
SLR (ours)	41.52	42.90	44.07	41.69	40.78	54.76	56.59	57.35	51.01	63.53	68.72	50.65	61.49	63.46	58.32
Groundtruth	55.68	58.10	61.27	55.84	55.08	65.83	67.22	67.56	62.60	72.49	76.97	65.04	70.86	72.51	68.56

Table 2: SSIM and SLR-adapted ResNet50 accuracy without and with input transformation (IT).

Corruption	Gauss	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Bright	Contrast	Elastic	Pixel	JPEG
SSIM	0.123	0.147	0.135	0.623	0.648	0.622	0.676	0.517	0.575	0.619	0.653	0.545	0.625	0.786	0.800
SSIM+IT	0.173	0.188	0.347	0.605	0.638	0.603	0.670	0.580	0.628	0.626	0.676	0.765	0.616	0.776	0.795
SLR	41.59	43.49	43.90	41.70	41.10	54.86	56.39	57.47	50.90	63.51	68.70	51.06	61.36	63.39	58.35
SLR+IT	43.09	44.39	64.05	41.98	40.99	55.73	56.75	58.56	51.68	63.64	68.85	55.01	61.32	63.59	58.24

265 statistics are still updated in these layers. Please refer Sec. A.2 for more details regarding which
 266 layers are frozen in different networks.

267 Furthermore, we prepend a trainable input transformation module d (cf. Sec. 3.1) to the network to
 268 partially counteract the data-shift. Note that the parameters of this module discussed in Sec. 3.1 are
 269 trainable and subject to optimization. This module is initialized to operate as an identity function prior
 270 to adaptation on a target distribution by choosing $\tau = 1$, $\gamma = \mathbf{1}$, and $\beta = \mathbf{0}$. We adapt the parameters
 271 of this module along with the channel-wise affine transformations and normalization statistics in
 272 an end-to-end fashion, solely using our proposed loss function along with the optimization details
 273 mentioned above. The architecture of this module is discussed in Sec. A.1.

274 Since L_{div} is independent of L_{conf} , we also propose to combine L_{div} with TENT, i. e. $L = L_{\text{div}} + L_{\text{ent}}$.
 275 We denote this as TENT+ and also set $\kappa = 0.9$ here. Note that TENT optimizes all channel-wise
 276 affine parameters in the network (since entropy is saturating and does not cause logit explosion).
 277 For a fair comparison to our method, we also freeze the top layers of the networks in TENT+. We
 278 show that adding L_{div} and freezing top layers significantly improves the networks performance over
 279 TENT. Note that SHOT [11] is the combination of TENT, batch-level diversity regularizer, and their
 280 pseudo labeling strategy. TENT+ can be seen as a variant of SHOT but without their pseudo labeling
 281 technique. Please refer to Sec. A.4 for the test-time adaptation of pretrained models with SHOT.

282 Note that each corruption and each severity in ImageNet-C is treated as a different target distribution
 283 and in all settings we reset model parameters to their pretrained values before every adaptation. We
 284 run our experiments for three times with different random seeds (2020, 2021, 2022) in PyTorch and
 285 report the average accuracies.

286 5 Results

287 **Evaluation on ImageNet-C** We adapt different models on the ImageNet-C benchmark using TENT,
 288 TENT+, and both *hard likelihood ratio* (HLR) and *soft likelihood ratio* (SLR) losses. Figure 2
 289 (top row) depicts the mean corruption accuracy (mCA%) of each model computed across all the
 290 corruptions and severity levels. It can be observed that TENT+ improves over TENT, showcasing the
 291 importance of a diversity regularizer L_{div} . Importantly, our methods HLR and SLR outperform TENT
 292 and TENT+ across DenseNet121, MobileNetV2, ResNet50, ResNeXt50 and perform comparable
 293 with TENT+ on robust ResNet50-DeepAugment+Augmix model. This shows that the mCA% of
 294 robust DeepAugment+Augmix model can be further increased from 58% (before adaptation) to
 295 68.6% using test-time adaptation techniques. Here, the average of mCA obtained from three different

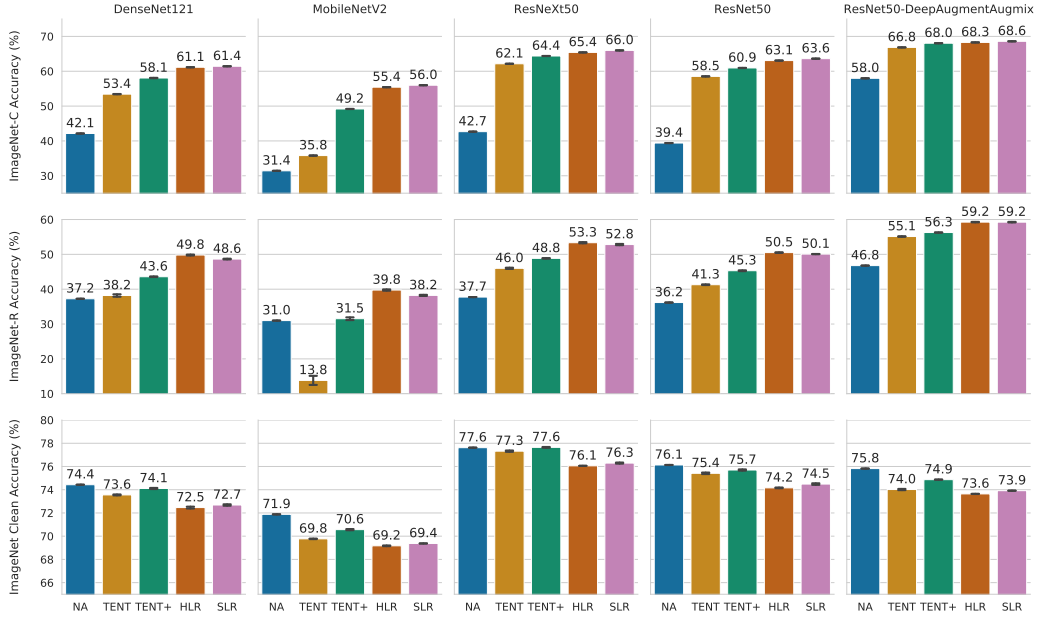


Figure 2: Test-time adaptation results on (top row) ImageNet-C, averaged across all 15 corruptions and severities, (middle row) ImageNet-R, (bottom row) clean ImageNet. NA refers to "No Adaptation".

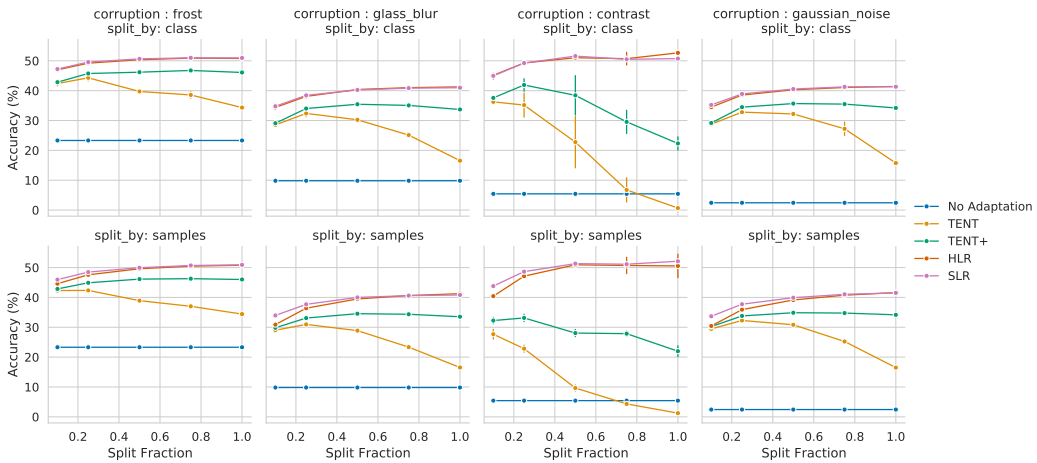


Figure 3: Test-time adaptation of ResNet50 using (top row) a subset of classes, and (bottom row) a subset of samples per class on 4 different corruptions at severity 5. Accuracy is computed based on the evaluation of adapted model on the entire target data. Note that error bars are smaller to visualize.

296 random seeds are depicted along with the error bars. These smaller error bars represent that the
 297 test-time adaptation results are not sensitive to the choice of random seed.

298 We also illustrate the performance of ResNet50 on the highest severity level across all 15 test
 299 corruptions of ImageNet-C in Table 1. Here, the adaptation results after epoch 1 and 5 are reported.
 300 It can be seen that a single epoch of test-time adaptation improves the performance significantly
 301 and makes minor improvements until epoch 5. TENT adaptation for more than one epoch result
 302 in reduced performance and TENT with L_{div} (TENT+) prevents this behavior. We note that both
 303 HLR and SLR clearly and consistently outperform TENT and TENT+ on the ResNet50. We also
 304 compare our results with the hard pseudo-labels (PL) objective and also with an oracle setting where
 305 the groundtruth labels of the target data are used for adapting the model in a supervised manner (GT).
 306 Note that this oracle setting is not of practical importance but illustrates the empirical upper bound on
 307 fully test-time adaptation performance under the chosen modulation parametrization. The reported
 308 numbers in the table are the average of three random seeds.

309 **ImageNet-R** We evaluate different adapted models on ImageNet-R and depict the results in Figure 2
310 (middle row). Results show that our methods significantly improve performance of all the models,
311 including the model pretrained with DeepAugment+Augmix. Moreover, both HLR and SLR clearly
312 outperform TENT and TENT+.

313 **Evaluation with data subsets** In the above experiments, the model is evaluated on the same data
314 that is also used for the test-time adaptation. Here, we test model generalization by adapting on a
315 subset of target data and evaluate the performance on the whole dataset, which also includes unseen
316 data that is not used for adaptation. We conduct two case studies: (i) adapt on the data from a subset
317 of ImageNet classes and evaluate the performance on the data from all the classes. (ii) Adapt only on
318 a subset of data from each class and test on all seen and unseen samples from the whole dataset.

319 Figure 3 illustrates generalization of a ResNet50 adapted on different proportions of the data across
320 different corruptions, both in terms of classes and samples. We observe that adapting a model on
321 a small subset of samples and classes is sufficient to achieve reasonable accuracy on the whole
322 target data. This suggests that the adaptation actually learns to compensate the data shift rather than
323 overfitting to the adapted samples or classes. The performance of TENT decreases as the number of
324 classes/samples increases, because L_{ent} can converge to trivial collapsed solutions and more data
325 corresponds to more updates steps during adaptation. Adding L_{div} such as in TENT+ stabilizes the
326 adaptation process and reduces this issues. Reported are the average of random seeds with error bars.

327 **Input transformation** We investigate whether the input transformation (IT) module, trained end-to-
328 end with a ResNet50 and SLR loss on data of the respective distortion *without* seeing any source
329 (undistorted) data, can partially undo certain domain shifts of ImageNet-C and also increase accuracy
330 on corrupted data. We measure domain shift via the structural similarity index measure (SSIM) [44]
331 between the clean image (unseen by the model) and its distorted version/the output of IT on the
332 distorted version. Table 2 shows that IT increases the SSIM considerably on certain distortions such as
333 Impulse, Contrast, Snow, and Frost. IT increases SSIM also for other types of noise distortions, while
334 it slightly reduces SSIM for the blur distortions, Elastic, Pixelate, and JPEG. When combined with
335 SLR, IT considerably increases accuracy on distortions for which also SSIM increased significantly
336 (for instance +20 percent points on Impulse, +4 percent points on Contrast) and never reduces
337 accuracy by more than 0.11 percent points. We provide illustrations of effect of IT in the appendix.

338 **Clean images** As a sanity check, we investigate the effect of test-time adaptation when target data
339 comes from the same distribution as training data. For this, we adapt pretrained models on clean
340 validation data of ImageNet. The results in Figure 2 (bottom row) depict that the performance of
341 SLR/HLR adapted models drops by 1.5 to 2.5 percent points compared to the pretrained model.
342 We attribute this drop to self-supervision being less reliable than the original full supervision on in-
343 distribution training data. The drop is smaller for TENT and TENT+, presumably because predictions
344 on in-distribution target data are typically highly confident such that there is little gradient and thus
345 little change to the pretrained networks by TENT. In summary, while self-supervision by confidence
346 maximization is a powerful method for adaptation to domain shift, the observed drop when adapting
347 to data from the source domain indicates that there is “no free lunch” in test-time adaptation.

348 6 Conclusion

349 We propose a method to improve corruption robustness and domain adaptation of models in a fully
350 test-time adaptation setting. Unlike entropy minimization, our proposed loss functions provide
351 non-vanishing gradients for high confident predictions and thus attribute to improved adaptation
352 in a self-supervised manner. We also show that additional diversity regularization on the model
353 predictions is crucial to prevent trivial solutions and stabilize the adaptation process. Lastly, we
354 introduce a trainable input transformation module that partially refines the corrupted samples to
355 support the adaptation. We show that our method improves corruption robustness on ImageNet-C and
356 domain adaptation to ImageNet-R on different ImageNet models. We also show that adaptation on a
357 small fraction of data and classes is sufficient to generalize to unseen target data and classes.

358 **Ethical and Societal Impact** Our non-saturating loss increases accuracy but might result in over-
359 confident predictions, which can cause harm in safety-critical downstream applications when not
360 properly calibrated. At the same time, self-supervised confidence maximization might amplify bias in
361 pretrained models. We hope that the diversity regularizer in the loss partially compensates this issue.

References

- 362
- 363 [1] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Fully
364 test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 365 [2] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence.
366 *Dataset Shift in Machine Learning*. MIT Press, 2009.
- 367 [3] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann,
368 and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias
369 improves accuracy and robustness. In *International Conference on Learning Representations*,
370 2019.
- 371 [4] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-
372 narayanan. Augmix: A simple data processing method to improve robustness and uncertainty.
373 *arXiv preprint arXiv:1912.02781*, 2019.
- 374 [5] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain
375 adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
376 pages 7167–7176, 2017.
- 377 [6] Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and
378 corruptions in natural corruption robustness. *arXiv preprint arXiv:2102.11273*, 2021.
- 379 [7] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Ad-
380 versarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings*
381 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2517–2526,
382 2019.
- 383 [8] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common
384 corruptions and perturbations. *International Conference on Learning Representations (ICLR)*,
385 2019.
- 386 [9] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
387 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin
388 Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization,
389 2020.
- 390 [10] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time
391 training with self-supervision for generalization under distribution shifts. In *International*
392 *Conference on Machine Learning*, pages 9229–9248. PMLR, 2020.
- 393 [11] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? Source
394 hypothesis transfer for unsupervised domain adaptation. In Hal Daumé III and Aarti Singh,
395 editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of
396 *Proceedings of Machine Learning Research*, pages 6028–6039. PMLR, 13–18 Jul 2020.
- 397 [12] Xiaofu Wu, Quan Zhou, Zhen Yang, Chunming Zhao, Longin Jan Latecki, et al. Entropy mini-
398 mization vs. diversity maximization for domain adaptation. *arXiv preprint arXiv:2002.01690*,
399 2020.
- 400 [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
401 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
402 pages 770–778, 2016.
- 403 [14] Norman Mu and Justin Gilmer. Mnist-c: A robustness benchmark for computer vision. *arXiv*
404 *preprint arXiv:1906.02337*, 2019.
- 405 [15] Claudio Michaelis, Benjamin Mitzkus, Robert Geirhos, Evgenia Rusak, Oliver Bringmann,
406 Alexander S Ecker, Matthias Bethge, and Wieland Brendel. Benchmarking robustness in object
407 detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484*, 2019.
- 408 [16] Christoph Kamann and Carsten Rother. Benchmarking the robustness of semantic segmenta-
409 tion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
410 *Recognition*, pages 8828–8838, 2020.
- 411 [17] Nic Ford, Justin Gilmer, Nicolas Carlini, and Dogus Cubuk. Adversarial examples are a natural
412 consequence of test error in noise. *arXiv preprint arXiv:1901.10513*, 2019.

- 413 [18] Raphael Gontijo Lopes, Dong Yin, Ben Poole, Justin Gilmer, and Ekin D Cubuk. Improving
414 robustness without sacrificing accuracy with patch gaussian augmentation. *arXiv preprint*
415 *arXiv:1906.02611*, 2019.
- 416 [19] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon
417 Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In
418 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6023–6032,
419 2019.
- 420 [20] Christoph Kamann, Burkhard Güssefeld, Robin Huttmacher, Jan Hendrik Metzen, and Carsten
421 Rother. Increasing the robustness of semantic segmentation models with painting-by-numbers.
422 *arXiv preprint arXiv:2010.05495*, 2020.
- 423 [21] Evgenia Rusak, Lukas Schott, Roland Zimmermann, Julian Bitterwolf, Oliver Bringmann,
424 Matthias Bethge, and Wieland Brendel. Increasing the robustness of dnns against image
425 corruptions by playing the game of noise. *arXiv preprint arXiv:2001.06057*, 3, 2020.
- 426 [22] Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and
427 corruptions in natural corruption robustness, 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=zbEupOtJFF)
428 [id=zbEupOtJFF](https://openreview.net/forum?id=zbEupOtJFF).
- 429 [23] Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and N Lawrence. Co-
430 variate shift and local learning by distribution matching, 2008.
- 431 [24] Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain
432 adaptation. In *Domain Adaptation in Computer Vision Applications*, pages 153–171. Springer,
433 2017.
- 434 [25] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation.
435 In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.
- 436 [26] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
437 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural
438 networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- 439 [27] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros,
440 and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International*
441 *conference on machine learning*, pages 1989–1998. PMLR, 2018.
- 442 [28] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A Efros. Unsupervised domain adaptation
443 through self-supervision. *arXiv preprint arXiv:1909.11825*, 2019.
- 444 [29] Jogendra Nath Kundu, Naveen Venkat, Rahul M V, and R. Venkatesh Babu. Universal source-
445 free domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
446 *Pattern Recognition (CVPR)*, June 2020.
- 447 [30] Rui Li, Qianfen Jiao, Wenming Cao, Hau-San Wong, and Si Wu. Model adaptation: Unsuper-
448 vised domain adaptation without source data. In *Proceedings of the IEEE/CVF Conference on*
449 *Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- 450 [31] Vinod K. Kurmi, Venkatesh K. Subramanian, and Vinay P. Namboodiri. Domain impression: A
451 source data free domain adaptation method. In *Proceedings of the IEEE/CVF Winter Conference*
452 *on Applications of Computer Vision (WACV)*, January 2021.
- 453 [32] Hao-Wei Yeh, Baoyao Yang, Pong C. Yuen, and Tatsuya Harada. Sofa: Source-data-free
454 feature alignment for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Winter*
455 *Conference on Applications of Computer Vision (WACV)*, January 2021.
- 456 [33] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias
457 Bethge. Improving robustness against common corruptions by covariate shift adaptation.
458 *Advances in Neural Information Processing Systems*, 33, 2020.
- 459 [34] Zachary Nado, Shreyas Padhy, D. Sculley, Alexander D’Amour, Balaji Lakshminarayanan, and
460 Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate
461 shift, 2021.
- 462 [35] Philipp Benz, Chaoning Zhang, Adil Karjauv, and In So Kweon. Revisiting batch normalization
463 for improving corruption robustness. In *Proceedings of the IEEE/CVF Winter Conference on*
464 *Applications of Computer Vision (WACV)*, pages 494–503, January 2021.

- 465 [36] Fabio Maria Carlucci, Lorenzo Porzi, Barbara Caputo, Elisa Ricci, and Samuel Rota Bulò.
466 Autodial: Automatic domain alignment layers. In *Proceedings of the IEEE International*
467 *Conference on Computer Vision*, pages 5067–5075, 2017.
- 468 [37] Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch
469 normalization for practical domain adaptation. *arXiv preprint arXiv:1603.04779*, 2016.
- 470 [38] Evgenia Rusak, Lukas Schott, Roland S. Zimmermann, Julian Bitterwolf, Oliver Bringmann,
471 Matthias Bethge, and Wieland Brendel. A simple way to make neural networks robust against
472 diverse image corruptions. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael
473 Frahm, editors, *Computer Vision – ECCV*, 2020.
- 474 [39] Hossein Talebi and Peyman Milanfar. Learning to resize images for computer vision tasks.
475 *arXiv preprint arXiv:2103.09950*, 2021.
- 476 [40] Hengshuai Yao, Dong-lai Zhu, Bei Jiang, and Peng Yu. Negative log likelihood ratio loss
477 for deep neural network classification. In Kohei Arai, Rahul Bhatia, and Supriya Kapoor,
478 editors, *Proceedings of the Future Technologies Conference (FTC)*, pages 276–282, Cham,
479 2020. Springer International Publishing. ISBN 978-3-030-32520-6.
- 480 [41] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko.
481 Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- 482 [42] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.
483 Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE*
484 *conference on computer vision and pattern recognition*, pages 5018–5027, 2017.
- 485 [43] Torch-Contributors. Torchvision models, 2020. URL [https://pytorch.org/vision/
486 stable/models.html](https://pytorch.org/vision/stable/models.html).
- 487 [44] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From
488 error visibility to structural similarity. *Trans. Img. Proc.*, 13(4):600–612, April 2004. ISSN
489 1057-7149.

490 Checklist

- 491 1. For all authors...
- 492 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
493 contributions and scope? [Yes]
- 494 (b) Did you describe the limitations of your work? [Yes] We discuss that our class
495 distribution matching term L_{div} requires knowledge of the class distribution $p_{\mathcal{D}'}(y)$
496 (Section 3.2.1). We also discuss that confidence maximization using a non-saturating
497 loss might result in overconfident predictions (Section 3.2.2). We also discuss the small
498 drop of accuracy when applying our method to adaptation on data from the source
499 domain (part “Clean images” in Section 4).
- 500 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Please
501 refer to the “Ethical and Societal Impact” paragraph in Section 6.
- 502 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
503 them? [Yes]
- 504 2. If you are including theoretical results...
- 505 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not
506 provide theoretical results
- 507 (b) Did you include complete proofs of all theoretical results? [N/A] We do not provide
508 theoretical results
- 509 3. If you ran experiments...
- 510 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
511 mental results (either in the supplemental material or as a URL)? [Yes] We include the
512 code and instructions as a part of supplementary material.
- 513 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
514 were chosen)? [Yes] We report training details and hyperparameters and how they are
515 chosen in Section 4, part “Settings”.

- 516 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
517 ments multiple times)? [Yes] We provide error bars in Figure 2. For the results in the
518 Tables, we report error bars in the appendix.
- 519 (d) Did you include the total amount of compute and the type of resources used (e.g., type
520 of GPUs, internal cluster, or cloud provider)? [No] We can not report the total amount
521 of compute since we did not track it. However, our organization is carbon neutral so
522 that all its activities including compute on the GPU clusters on which the experiments
523 have been performed do no longer leave a carbon footprint.
- 524 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 525 (a) If your work uses existing assets, did you cite the creators? [Yes] Our work builds upon
526 pre-trained models. We cite the creators (see Section 4, part “Models”). Our work uses
527 several openly available datasets. We cite their creators (see Section 4, part “Datasets”).
- 528 (b) Did you mention the license of the assets? [No]
- 529 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
530 We do not provide new assets.
- 531 (d) Did you discuss whether and how consent was obtained from people whose data you’re
532 using/curating? [No]
- 533 (e) Did you discuss whether the data you are using/curating contains personally identifiable
534 information or offensive content? [N/A] To the best of our knowledge data used in this
535 project does not contain any personally identifiable information or offensive content.
- 536 5. If you used crowdsourcing or conducted research with human subjects...
- 537 (a) Did you include the full text of instructions given to participants and screenshots, if
538 applicable? [N/A] No crowdsourcing was used and no research with human subjects
539 was conducted
- 540 (b) Did you describe any potential participant risks, with links to Institutional Review
541 Board (IRB) approvals, if applicable? [N/A] No crowdsourcing was used and no
542 research with human subjects was conducted
- 543 (c) Did you include the estimated hourly wage paid to participants and the total amount
544 spent on participant compensation? [N/A] No crowdsourcing was used and no research
545 with human subjects was conducted