
Test-Time Training with Masked Autoencoders

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Prior work has shown masked autoencoding as an effective self-supervised task
2 across many visual distributions with sufficient training data. We use masked
3 autoencoding to train on each unlabeled test sample as it arrives at test time, before
4 making a prediction. This simple method improves generalization of predictive
5 models on many visual benchmarks of unknown distributions, where input images
6 are not covered by training data. We also theoretically characterize the relationship
7 between our method and the well known bias-variance trade-off in statistics.

8 1 Introduction

9 Generalization is the central theme of supervised learning, and a hallmark of intelligence. While
10 most of the research focuses on generalization when the training and test data are drawn from the
11 same distribution, this is rarely the case for real world deployment [47], and is certainly not true for
12 environments where intelligence has emerged.

13 Most models today are fixed during deployment, even when the test distribution changes. As a
14 consequence, the trained models need to be robust to all potential distribution shifts that could happen
15 in the future [24, 14, 49, 38]. This turns out to be quite difficult because being ready for all possible
16 futures limits the model’s capacity to be good at any particular one. But only one of these futures is
17 actually going to happen.

18 This motivates an alternative perspective on generalization: instead of being ready for everything,
19 one simply adapts to the future once it arrives. Test-time training (TTT) is one line of work that
20 takes this perspective [46, 36, 20, 1]. The key insight is that each test input gives a hint about the
21 test distribution. We modify the model at test time to take advantage of this hint by setting up a
22 *one-sample learning problem*. The only issue is that the test input comes without a ground truth label.
23 But we can generate labels from the input itself thanks to self-supervised learning. At training time,
24 TTT trains on both the main task (e.g. classification) and the self-supervised task. Then at test time,
25 it adapts the model on the self-supervised task alone for each test input, before making a prediction
26 on the main task.

27 The choice of the self-supervised task is critical. On one hand, this task must be general and broad
28 enough to work on a wide range of potential test distributions. But on the other hand, it should not be
29 invariant to the changes across these test distributions, otherwise test-time training on a new sample
30 will not provide the learning signal needed to improve the main task. What self-supervised task is
31 both general and difficult enough to fit these requirements? We turn to a fundamental property shared
32 by natural visual data – spatial smoothness, i.e. the local redundancy of information in xy space. The
33 task of spatial prediction – remove parts of the data, and predict the removed content – forms the

34 basis of some of the most successful self-supervised approaches, e.g. denoising autoencoder [50],
35 context encoder [41], and masked autoencoder (MAE) [21].

36 In this paper, we should that MAE [21] turns out to be extremely well-suited for test-time training.
37 Our method leads to substantial improvements on four different object recognition datasets. We also
38 offer a theoretical justification for the use of test-time training based on a simplified linear formulation.
39 The code and the models will be made public upon acceptance.

40 **2 Related Work**

41 **2.1 Masked Autoencoders and Self-Supervised Learning**

42 Key to the success of deep learning is the formation of semantic features that transfer across many
43 tasks; the goal of self-supervised pre-training is to learn these transferable features without human
44 annotations, by making up inputs and labels to an auxiliary pretext task. Given an original input
45 image x , the pretext task of masked autoencoding [21] splits the image into many non-overlapping
46 patches (e.g. 196), randomly masks out majority of the patches (e.g. 75%), and trains a model to
47 reconstruct the missing patches, by optimizing the mean squared error between the reconstructed and
48 the original pixels.

49 This model, called a masked autoencoder (MAE), is composed of two parts: a large encoder and
50 a small decoder, both based on Vision Transformers [12]. For a downstream task, e.g. object
51 recognition, the encoder takes in all the image patches and outputs the features, which then become
52 input to a linear projection head. There are two common ways to combine the pre-trained encoder and
53 untrained head. 1) Fine-tuning: both the encoder and head are trained, end-to-end, for the downstream
54 task. 2) Linear probing: the encoder is frozen as a fixed feature extractor, and only the head is trained.
55 We will refer back to these training options in Section 3.

56 **2.2 Generalization under Distribution Shifts**

57 When training and test distributions are different, generalization is intrinsically hard. Training on
58 many distribution shifts at once causes underfitting [53, 11]. The robustness obtained by training or
59 fine-tuning on one type of distribution shift (e.g. Gaussian noise) also does not transfer to another
60 (e.g. salt-and-pepper noise) [14, 49], even for visually similar corruptions.

61 To avoid the hardness of extrapolation, researchers work on alternative problem settings that anticipate
62 the test distribution. For domain generalization [30, 29, 42, 39, 4, 15, 37, 19], a meta distribution
63 generates many distribution shifts, some of which are used for training and others for testing. For
64 domain adaptation [33, 34, 8, 5, 44, 26, 6, 18], a potentially unlabeled dataset from the test distribution
65 is available for training. Solutions motivated by these settings share the same flaw as augmenting the
66 training set: they are only effective on the particular meta / test distribution that they were trained for.

67 **2.3 Test-Time Training**

68 Test-time training (TTT) [46] is a paradigm of methods that generalizes to distribution shifts without
69 anticipating them, through self-supervised learning on each test input x . In [46], the main task is
70 object recognition, and the self-supervised task is rotation prediction [16]: x is simply rotated in the
71 image plane by multiples of 90 degrees to create a four-way classification problem.

72 The network for TTT is Y-shaped: a feature extractor f , followed simultaneous by a self-supervised
73 head g and a main task head h . During training, labels are available for both the self-supervised
74 task, using $g \circ f$, and the main task, using $h \circ f$; so both branches are trained jointly by summing
75 their losses together for gradient descent. During testing, for each x , only the self-supervised task is
76 labeled; so $g \circ f$ is trained, end-to-end, on this single input. This produces a different model, call it f_x
77 (and g_x), initialized from f (and g), locally optimized for each test input x . After making a prediction
78 on x , now with $h \circ f_x$, the algorithm never uses f_x (and g_x) again; the entire process repeats itself for
79 the next input, again initializing from the old f (and g).

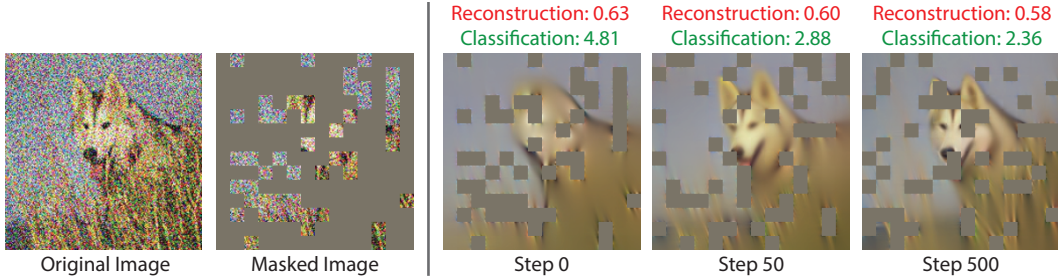


Figure 1: We train an MAE to reconstruct each test image at test time, with 75% of the input patches masked out. The three reconstructed images on the right visualize the progress of this single-input learning problem. Loss of the main task (green) – object recognition – keeps dropping even after 500 steps of gradient descent, while the network continues to optimize for reconstruction (red). The unmasked patches are not shown on the right since they are not part of the reconstruction loss.

80 TTT as explained above is not a new problem setting; it is a proposed solution to the canonical
 81 problem of generalization. Our paper stay in the same setting. Other followups have worked on new
 82 settings such as a batch [51] or dataset [32] of samples from the test distribution.

83 Before [46], single-image training has been performed for super-resolution [43]. The philosophy of
 84 TTT is also preceded by that of local learning [3] and transductive learning [48], although usually
 85 performed for support vector machines with multiple test samples [7, 27]. Since [46], TTT has
 86 been applied to many fields with domain-specific self-supervised tasks: vision-based reinforcement
 87 learning [20], model-based legged locomotion [45], tracking objects in videos [13], natural language
 88 question answering [1], and even medical imaging [28].

89 3 Method

90 At a high level, our method simply substitutes MAE [21] for the self-supervised part of TTT [46].
 91 In practice, making this work involves many design choices.

92 **Architecture.** Like in [46], our architecture is Y-shaped. The feature extractor f is exactly the
 93 encoder of MAE, and g the decoder; each is a Vision Transformer (ViT) [12]. We intentionally avoid
 94 modifying them to make clean comparisons with results in [21]. Choice of the main task (object
 95 recognition) head h presents more opportunities. [21] uses a linear projection, from the encoder
 96 feature dimension to the number of classes; the authors discuss this as mostly a historic artifact of
 97 how people have been evaluating the quality of features in unsupervised learning. Therefore we also
 98 experiment with a more expressive main task head – an entire ViT-Base – to strengthen our baseline.

99 **Training-time training.** Following standard practice, we start from the MAE model provided by
 100 the authors, with a ViT-Large encoder, pre-trained for reconstruction on ImageNet-1k [10]. There are
 101 three ways to combine it with the untrained main task head: 1) Fine-tuning: train $h \circ f$ end-to-end.
 102 2) Probing: train h only, with f frozen. 3) Joint training: train both $h \circ f$ and $g \circ f$, by summing
 103 the respective losses together. The first two are practiced by [21], only for linear h (see Subsection
 104 2.1). The third is practiced by [46], for linear h , as well as a neural network (see Subsection 2.3).
 105 We experiment with all three, and choose probing with a ViT-Base head, a.k.a. *ViT probing*, as our
 106 default setup, based on the following considerations:

107 1) Fine-tuning does not work with TTT. As discussed in [21], reconstruction and recognition require
 108 very different features. Fine-tuning makes the encoder specialize to recognition, only to find it
 109 training for reconstruction again at test time and losing the recognition-only features that the head
 110 learned to rely on. This agrees with the intuition in [46] for joint training, instead of training only for
 111 recognition, as preparation for TTT.

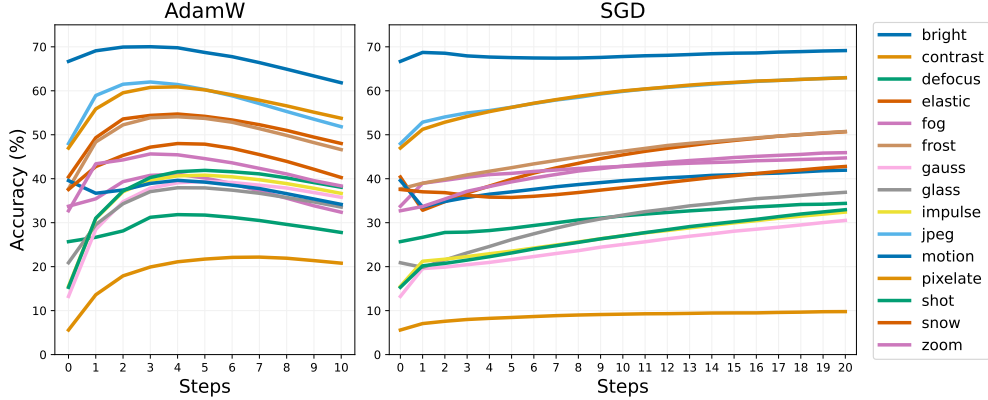


Figure 2: We experiment with two optimizers for TTT. MAE uses AdamW for pre-training, but using AdamW for TTT requires early stopping, which is unrealistic for generalization to unknown distributions. We instead use SGD, which keeps improving performance even after 20 steps.

122 2) ViT probing is much more lightweight than both fine-tuning and joint-training. It trains 3.5 times
 123 fewer parameters than even linear fine-tuning (86M vs. 306M). It also obtains higher accuracy than
 124 linear fine-tuning without aggressive augmentations on the ImageNet validation set, and accuracy
 125 similar to linear fine-tuning with augmentations, which [21] optimized heavily with many hyper-
 126 parameters that prevent overfitting. Altogether, we believe that ViT probing strikes the right balance
 127 between the highly constrained linear probing and highly expressive linear / ViT fine-tuning.

128 3) Joint training is the recommended design in [46], but its accuracy on the ImageNet validation set,
 129 to the best of our ability, has been worse than the linear / ViT probing counterparts. This is most likely
 130 because we have not explored enough of the hyper-parameter space. As discussed in 2), fine-tuning is
 131 already heavyweight and hard to train; joint training is only heavier, and has more hyper-parameters.
 132 If future work finds a reliable recipe for joint training, our method will likely benefit from its setup.
 133 Meanwhile, we consider it our contribution to find ViT probing as a lightweight substitute.

124 Formally, denote the encoder produced by MAE pre-training as f_0 (and the decoder, used later for
 125 TTT, as g_0). ViT probing produces a trained main task head h_0 :

$$h_0 = \arg \min_h \frac{1}{n} \sum_{i=1}^n l_m(h \circ f_0(x_i), y_i). \quad (1)$$

126 The summation is over the training set with n samples, each consisting of input x_i and label y_i .
 127 The main task loss l_m in our case is the cross entropy loss for classification. Note that we are only
 128 optimizing the main task head, while the pre-trained encoder f_0 is frozen.

129 **Augmentations.** Our default setup, during training-time training, only uses image cropping and
 130 horizontal flips for augmentations, follow the protocol in [21] for pre-training and linear probing.
 131 Fine-tuning in [21] (and [2, 52]), however, adds aggressive augmentations on top of the two above,
 132 including random changes in brightness, contrast, color and sharpness¹. Many of them are in fact
 133 distribution shifts in our evaluation benchmarks. Training with them is analogous to training on the
 134 test set, for the purpose of studying generalization to new test distributions. Therefore, we choose not
 135 to use these augmentations, even though they would improve our results at face value.

136 **Test-time training.** At test time, we start from the main task head h_0 produced by ViT probing, as
 137 well as the MAE pre-trained encoder f_0 and decoder g_0 . Once each test input x arrives, we optimize
 138 the following loss for TTT:

$$f_x, g_x = \arg \min_{f,g} l_s(g \circ f(\text{mask}(x)), x). \quad (2)$$

¹Other augmentations used by fine-tuning in [21] are: interpolation, equalization, inversion, posterization, solarization, rotation, shearing, random erasing, and translation. These are taken from RandAugment [9].

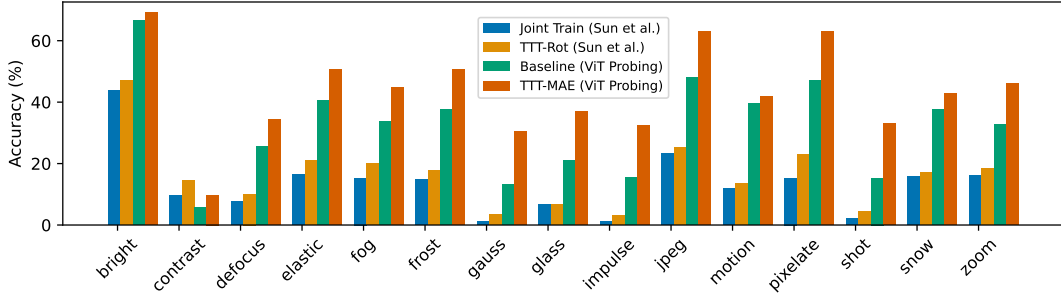


Figure 3: Accuracy (%) on ImageNet-C, level 5. Our method, TTT-MAE, significantly improves on top of our baseline, which already outperforms the method of [46]. See Subsection 4.2 for details. Numbers for our baseline and TTT-MAE can be found in the last two rows of Table 2. Numbers for Sun et al. are taken from [46].

139 The self-supervised reconstruction loss l_s computes the pixel-wise squared error of the decoded
 140 patches with the ground truth. After TTT, we make a prediction on x as $h \circ f_x(x)$. Note that
 141 gradient-based optimization for Equation 2 always starts from f_0 and g_0 . When evaluating on a test
 142 set, we always discard f_x and g_x after making a prediction on each test input x , and reset the weights
 143 to f_0 and g_0 for the next test input. By test-time training on the test inputs independently, we do not
 144 assume that they come from the same distribution.

145 The choice of optimizer for TTT is straightforward in [46]: it simply takes exactly the same optimizer
 146 setting as during the last epoch of training-time training of the self-supervised task. This choice,
 147 however, is not available for us, because the learning rate schedule of MAE reaches zero by the end
 148 of pre-training. We experiment with various learning rates for AdamW [35] and stochastic gradient
 149 descent (SGD) with momentum. Performance of both, using the best learning rate respectively, is
 150 shown in Figure 2.

151 AdamW is used in [21], but for TTT it hurts performance with too many iterations. On the other
 152 hand, more iterations with SGD consistently improve performance on all distribution shifts. Test
 153 accuracy keeps improving even after 20 iterations. This is very desirable for TTT: the single test
 154 image is all we know about the test distribution, so there is no validation set to tune hyper-parameters,
 155 or monitor performance for early stopping. With SGD, we can simply keep training.

156 4 Empirical Results

157 4.1 Implementation Details

158 In all experiments we use the MAE based on ViT-Large, configured as in [21], pre-trained on
 159 ImageNet-1k [10] for 800 epochs. Our ViT-Base head takes as input the image features from the
 160 pre-trained MAE encoder. There is a linear layer in between that resizes those features to fit as inputs
 161 to the head, just like between the encoder and the decoder in [21]. A linear layer is also appended to
 162 the final class token of the head to produce the classification logits.

163 TTT is performed with SGD, as discussed, for 20 steps, using a momentum of 0.9, weight decay of
 164 0.2, batch size of 128, and fixed learning rate of $5e-3$. The choice of 20 steps is purely computational;
 165 more steps will likely improve performance marginally, judging from the trend observed in Figure 2.
 166 Most experiments are performed on four NVIDIA A100 GPUs; hyper-parameter sweeps are ran on
 167 an industrial cluster with mostly V100 GPUs.

168 We optimize both the encoder and decoder weights during TTT, together with the class token and the
 169 mask token. We have experimented with freezing the decoder and found that the difference, even for
 170 multiple iterations of SGD, is negligible; see the appendix for detailed comparison. This is consistent
 171 with the observations in [46]. We also tried reconstruction loss both with and without normalized

172 pixels, as done in [21]. Our method works well for both, slightly better for normalized pixels because
173 the baseline is slightly better. We include these results in the appendix.

174 4.2 ImageNet-C

175 ImageNet-C [24] is a benchmark for object recognition under distribution shifts. It contains copies
176 of the ImageNet [10] validation set with 15 types of corruptions, independent across images, each
177 with 5 levels of severity. [Due to space constraints, all results in the main text are on level 5, the most
178 severe, unless stated otherwise. Results on the other four levels are in the appendix. Sample images
179 from this benchmark are also in the appendix in Figure 4.](#)

180 In the spirit of treating these distribution shifts as truly unknown, we do not use training data, prior
181 knowledge, or data augmentations derived from these corruptions. This complies with the stated
182 rule of the benchmark, that the corruptions should be used only for evaluation, and the algorithms
183 should not be corruption-specific [24]. This rule, in our opinion, is helpful for community progress:
184 as explained in Subsection 2.2, the numerous distribution shifts outside of research evaluation cannot
185 be anticipated, and algorithms that rely on information specific to the test distribution cannot scale.

186 Our main results on ImageNet-C appear in Figure 3. Following the convention in [46], we plot
187 the accuracy only on the level-5 (most severe) corruptions. Results on the other levels are in the
188 appendix. [Our baseline \(green\) in Figure 3 is the default training method throughout this paper: take
189 a pre-trained MAE encoder, perform ViT probing for object recognition on the original ImageNet
190 training set, then apply the fixed model to the corrupted test sets.](#) Performing TTT-MAE on our
191 baseline significantly improves performance.

192 **Reconstruction vs. rotation prediction.** The baseline for TTT-Rot, taken from [46], is called Joint
193 Train in Figure 3. This is a ResNet [22] with 16-layers, trained jointly for rotation prediction and
194 object recognition. Because our baseline is much more advanced than the ResNet baseline, the
195 former already outperforms the reported results of TTT-Rot, for all corruptions except contrast.² So
196 comparison to [46] is more meaningful in relative terms: TTT-MAE has higher performance gains in
197 all corruptions than TTT-Rot, on top of their respective baselines.³

198 **Rotation invariant classes.** As discussed in the introduction, many self-supervised tasks like
199 rotation prediction [16] are designed for specific domains. This can be especially limiting when the
200 goal is, in fact, to generalize to unknown domains. Here we show examples in ImageNet where the
201 assumptions behind rotation prediction do not hold.

202 We find rotation invariant classes through the following: 1) Take the ImageNet pre-trained ResNet-18
203 from the official PyTorch website [40]. 2) Run it on the original validation set, and one copy of it
204 with 90-degree rotation. 3) Filter out classes for which the original accuracy is lower than 60%. 4)
205 For each class, compare accuracy for the original and the rotate images, and choose the 5 classes for
206 which the normalized difference is the smallest.

207 Sample images from these 5 classes are shown in Table 1. Not surprisingly, these images are usually
208 taken from top-down views; rotation prediction on them can only memorize the generated labels,
209 without forming semantic features as intended. This causes TTT-Rot of [46] to harm performance, as
210 seen in Table 1. Also not surprisingly, TTT-MAE is agnostic to rotation invariance and still helps.

211 **Training-time training.** In the synonymous paragraph in Section 3, we discussed our choice of
212 ViT probing instead of fine-tuning or joint training. [The first three rows of Table 2 compare accuracy
213 of these three designs. As discussed earlier, joint training does not achieve satisfactory performance](#)

²It turns out that for the contrast corruption type, the ResNet baseline of TTT-Rot is somehow much better than our ViT baseline; TTT-MAE improves on the ViT baseline nevertheless.

³For fair comparison, we have also attempted to implement the rotation prediction task of [46] on top of our ViT probing baseline, as a simple substitution to the reconstruction task of [21]; the rotation prediction head here is the same as the main task head. However, even after joint training in the fashion of [46], TTT for this model with rotation still hurts performance on every corruption. Please see appendix for details.

	<i>Frost</i>		<i>JPEG</i>		<i>Brightness</i>		<i>Elastic</i>	
	[46]	Ours	[46]	Ours	[46]	Ours	[46]	Ours
Jellyfish	-3	0	-6	5	2	-1	-2	3
Holster	-1	4	-2	6	1	0	-1	14
Chiton	-1	3	-3	8	-1	2	-3	8
Fire salamander	1	3	-1	3	2	5	2	2
Spaghetti squash	-3	4	-2	5	2	3	1	6

Table 1: Changes after TTT in number of correctly classified images, out of 50 total for each category. On these rotation invariant classes, TTT-Rot [46] hurts performance, while TTT-MAE still helps, thanks to the generality of MAE as a self-supervised task; see details in main text. The four corruptions are selected for being the *most accurate* categories of the TTT-Rot baseline.

	brigh	cont	defoc	elast	fog	frost	gauss	glass	impul	jpeg	motn	pixel	shot	snow	zoom
Joint Train	62.3	4.5	26.7	39.9	25.7	30.0	5.8	16.3	5.8	45.3	30.9	45.9	7.1	25.1	31.8
Fine-Tune	67.5	7.8	33.9	32.4	36.4	38.2	22.0	15.7	23.9	51.2	37.4	51.9	23.7	37.6	37.1
ViT Probe	68.3	6.4	24.2	31.6	38.6	38.4	17.4	18.4	18.2	51.2	32.2	49.7	18.2	35.9	32.2
TTT-MAE	69.1	9.8	34.4	50.7	44.7	50.7	30.5	36.9	32.4	63.0	41.9	63.0	33.0	42.8	45.9

Table 2: Accuracy (%) on ImageNet-C, level 5. The first three rows are fixed models without test-time training, comparing the three design choices discussed in Section 3. The third row, ViT probing, is our default baseline throughout the paper; it has the same numbers as the baseline in Figure 3. The last row is our method: TTT-MAE on ViT probing. It achieves the best performance across all corruption types; it has the same numbers as TTT-MAE in Figure 3.

214 in most categories. While fine-tuning is initially better than ViT probing, it is not amenable to TTT.
 215 Therefore, ViT probing is chosen as the default baseline throughout the paper. All three of the above
 216 are only for training-time training, after which a fixed model is applied during testing. The last row is
 217 TTT-MAE on a model trained with ViT probing, and performs the best across all corruption types.
 218 Numbers in the last two rows are the same as for the baseline and TTT-MAE in Figure 3.

219 4.3 Other ImageNet Variants

220 We evaluated TTT-MAE on two other popular benchmarks for distribution shifts. ImageNet-A [25]
 221 contains real-world, unmodified, and naturally occurring examples where popular ImageNet models
 222 perform poorly. ImageNet-R [23] contains renditions of ImageNet classes, e.g. art, cartoons and
 223 origami. Both TTT-MAE and the baseline are the same models and algorithms, with exactly the
 224 same hyper-parameters as for ImageNet-C (Subsection 4.2). Our method continues to outperform the
 225 baseline by large margins, see Table 4.3.

226 4.4 Portraits Dataset

227 The Portraits dataset [17] contains American high school yearbook photos labeled by gender, taken
 228 over more than a century. It contains real-world distribution shifts in appearances over the years. We
 229 sort the entire dataset by year and split it into four equal parts, with 5062 images each. Between the
 230 four splits, there are visible low-level differences such as blurriness and contrast, as well as high-level
 231 change like hair-style and smile.

232 We create three experiments: for each we train on one of the first three splits and test on the fourth.
 233 Performance is in terms of accuracy in binary gender classification. As expected, performance
 234 increases slightly as the training split gets closer to test in time.

235 Our model is the same ImageNet pre-trained MAE + ViT head (with sigmoid for binary classification),
 236 and probing is performed on the training split. We use exactly the same hyper-parameters for TTT.
 237 Our results are shown in 4: our method improves on the baseline for all three experiments.

238 5 Theoretical Results

239 Why does TTT help? The theory of [46] gives an intuitive but shallow explanation: when the
 240 self-supervised task happens to propagate gradients that correlate with those of the main task. But
 241 this evasive theory only delegates one unknown – the test distribution, to another – the magical
 242 self-supervised task with correlated gradients, without really answering the question, [or showing any](#)
 243 [concrete example of such a self-supervised task.](#)

244 In this paper, we show that autoencoding, i.e. reconstruction, is a self-supervised task that makes
 245 TTT help. To do so with minimal mathematical overhead, we restrict ourselves to the linear world,
 246 where our models are linear and the distribution shifts are produced by linear transformations. We
 247 analyze autoencoding with dimensionality reduction instead of masking, as we believe the two are
 248 closely related in essence.

249 The most illustrative insight, in our opinion, is that under distribution shifts, TTT finds a better *bias-*
 250 *variance trade-off* than applying a fixed model. The fixed model is biased because it is completely
 251 based on biased training data, which do not represent the test sample under distribution shifts. The
 252 other extreme is to completely forget the training data, and train a new model from scratch on each
 253 test input. This is also undesirable because the single test input is high variance, although unbiased
 254 by definition. A sweet spot of the trade-off can be found by performing TTT while remembering the
 255 training data in some capacity: in practice by initializing with a model trained on them, and in this
 256 section by using part of their data covariance matrix.

257 It is well known that linear autoencoding is equivalent to principle component analysis (PCA).
 258 This equivalence simplifies the mathematics by giving us closed form solutions to the optimization
 259 problems both during training and test-time training. For the rest of the section, we use the term PCA
 260 instead of linear autoencoding, following convention of the theory community.

261 **Problem setup.** Let $x, y \sim P$, the training distribution, and assume that the population covariance
 262 matrix $\Sigma = \text{Cov}(x)$ is known. PCA performs spectral decomposition on $\Sigma = UDU^\top$, and takes the
 263 top k eigenvectors u_1, \dots, u_k to project $x \in \mathbb{R}^d$ to \mathbb{R}^k . Throughout this section, we denote u_i as the
 264 i th column of the matrix U , and likewise for other matrices.

265 Assume that for each x , the ground truth y is a linear function of the PCA projection with $k = 1$:

$$y = wu_1^\top x, \tag{3}$$

266 for some known weight $w \in \mathbb{R}$. For mathematical convenience, we also assume the following about
 267 the eigenvalues, i.e. the diagonal entries of D :

$$\sigma_1 > \sigma_2 = \sigma_3 = \dots = \sigma_d = \sigma. \tag{4}$$

<i>ImageNet-A</i>		<i>ImageNet-R</i>	
Baseline	Ours	Baseline	Ours
15.3	21.3	31.3	38.9

Table 3: Accuracy (%) on ImageNet-A [25] and ImageNet-R [23]; see Subsection 4.3 for details. Our baseline is ViT-Probing, as default, trained on the original ImageNet.

Split (train)	1	2	3
Baseline	76.1	76.5	78.2
TTT-MAE	76.4	76.7	79.4

Table 4: Accuracy (%) for binary classification of gender. For each column, we train on the indicated split, and test on the fourth. Our method improves on the baseline for all three training splits.

268 At test time, nature draws a new $x, y \sim P$, but we can only see \tilde{x} , a corrupted version of x . We model
269 a corruption as an unknown orthogonal transformation R , and $\tilde{x} = Rx$.

270 **Algorithms.** The fixed model baseline is to blithely apply our old procedure and predict

$$\hat{y} = wu_1^\top \tilde{x}, \quad (5)$$

271 as an estimate of y . Intuitively, this will be inaccurate if corruption is severe i.e. R is far from I . And
272 here is the PCA version of TTT. We form a new (and rather trivial) covariance matrix with \tilde{x} , the
273 only piece of information we have about the corruption. Let $\alpha \in [0, 1]$ be a hyper-parameter, and the
274 linear combination of Σ with our new covariance matrix be

$$M(\alpha) = (1 - \alpha) \cdot \Sigma + \alpha \cdot \tilde{x}\tilde{x}^\top. \quad (6)$$

275 Denote its spectral decomposition as $M(\alpha) = V(\alpha) \cdot S(\alpha) \cdot V^\top(\alpha)$. We then predict

$$\hat{y} = wv_1^\top \tilde{x}. \quad (7)$$

276 **Bias-variance trade-off.** $\alpha = 0$ is equivalent to the baseline, and $\alpha = 1$ means that we exclusive
277 use the one test input and completely forget about the training data. In statistical terms, α presents a
278 bias-variance trade-off: $\alpha \downarrow 0$ means more bias, $\alpha \uparrow 1$ means more variance.

279 **Theorem.** Define the prediction risk as $\mathbb{E}[|\hat{y} - y|]$, where the expectation is taken over the corrupted
280 test distribution. This risk is strictly dominated when $\alpha = 0$. That is, TTT with some hyper-parameter
281 choice $\alpha > 0$ is, on average, strictly better than the baseline.

282 The proof is given in the appendix.

283 **Remarks on assumptions.** The assumption in Equation 3 is mild; it basically just says that PCA is
284 at least helpful on the training distribution. It is also mild to assume that Σ and w are known in this
285 context, since the estimated covariance matrix and weight asymptotically approach the population
286 ground truth as the training set grows larger. The assumption on the eigenvalues in Equation 6 greatly
287 simplifies the math, but a more complicated version of our analysis should exist without it. Lastly,
288 we believe that our general statement should still hold for invertible linear transformation instead of
289 orthogonal transformations, since they share the same intuition.

290 6 Discussion

291 **Limitations.** Due to the additional computation at test time, our method is slower than the baseline.
292 Inference speed has not been the focus of this paper, and it can potentially be improved through better
293 optimizers and hyper-parameters. In general, test-time training has been limited by the deep learning
294 toolbox. Most tools, for architectural design, regularization and optimization, have been refined for
295 training-time training on large datasets, especially for large models like ours; this is another direction
296 for future work.

297 References

- 298 [1] Pratyay Banerjee, Tejas Gokhale, and Chitta Baral. Self-supervised test-time learning for reading compre-
299 hension. *arXiv preprint arXiv:2103.11263*, 2021.
- 300 [2] Hangbo Bao, Li Dong, and Furu Wei. Beit: BERT pre-training of image transformers. *CoRR*,
301 abs/2106.08254, 2021.
- 302 [3] Léon Bottou and Vladimir Vapnik. Local learning algorithms. *Neural computation*, 4(6):888–900, 1992.
- 303 [4] Fabio M Carlucci, Antonio D’Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain
304 generalization by solving jigsaw puzzles. In *Proceedings of the IEEE Conference on Computer Vision and*
305 *Pattern Recognition*, pages 2229–2238, 2019.
- 306 [5] Minmin Chen, Kilian Q Weinberger, and John Blitzer. Co-training for domain adaptation. In *Advances in*
307 *neural information processing systems*, pages 2456–2464, 2011.
- 308 [6] Xilun Chen, Yu Sun, Ben Athiwaratkun, Claire Cardie, and Kilian Weinberger. Adversarial deep averaging
309 networks for cross-lingual sentiment classification. *Transactions of the Association for Computational*
310 *Linguistics*, 6:557–570, 2018.
- 311 [7] Ronan Collobert, Fabian Sinz, Jason Weston, Léon Bottou, and Thorsten Joachims. Large scale transductive
312 svms. *Journal of Machine Learning Research*, 7(8), 2006.
- 313 [8] Gabriela Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint*
314 *arXiv:1702.05374*, 2017.
- 315 [9] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical data augmentation
316 with no separate search. *CoRR*, abs/1909.13719, 2019.
- 317 [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
318 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.
319 Ieee, 2009.
- 320 [11] Samuel Dodge and Lina Karam. Quality resilient deep neural networks. *arXiv preprint arXiv:1703.08119*,
321 2017.
- 322 [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
323 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth
324 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 325 [13] Yang Fu, Sifei Liu, Umar Iqbal, Shalini De Mello, Humphrey Shi, and Jan Kautz. Learning to track
326 instances without video annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
327 *Pattern Recognition*, pages 8680–8689, 2021.
- 328 [14] Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A
329 Wichmann. Generalisation in humans and deep neural networks. In *Advances in Neural Information*
330 *Processing Systems*, pages 7538–7550, 2018.
- 331 [15] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for
332 object recognition with multi-task autoencoders. In *Proceedings of the IEEE international conference on*
333 *computer vision*, pages 2551–2559, 2015.
- 334 [16] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting
335 image rotations. *arXiv preprint arXiv:1803.07728*, 2018.
- 336 [17] Shiry Ginosar, Kate Rakelly, Sarah Sachs, Brian Yin, and Alexei A. Efros. A century of portraits: A visual
337 historical record of american high school yearbooks. *CoRR*, abs/1511.02575, 2015.
- 338 [18] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain
339 adaptation. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2066–2073.
340 IEEE, 2012.
- 341 [19] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint*
342 *arXiv:2007.01434*, 2020.
- 343 [20] Nicklas Hansen, Rishabh Jangir, Yu Sun, Guillem Alenyà, Pieter Abbeel, Alexei A Efros, Lerrel Pinto, and
344 Xiaolong Wang. Self-supervised policy adaptation during deployment. *arXiv preprint arXiv:2007.04309*,
345 2020.

- 346 [21] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked
347 autoencoders are scalable vision learners. *CoRR*, abs/2111.06377, 2021.
- 348 [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
349 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 350 [23] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai,
351 Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces
352 of robustness: A critical analysis of out-of-distribution generalization. *ICCV*, 2021.
- 353 [24] Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corrup-
354 tions and perturbations. *CoRR*, abs/1903.12261, 2019.
- 355 [25] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
356 examples. *CVPR*, 2021.
- 357 [26] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A Efros, and
358 Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. *arXiv preprint arXiv:1711.03213*,
359 2017.
- 360 [27] Thorsten Joachims. *Learning to classify text using support vector machines*, volume 668. Springer Science
361 & Business Media, 2002.
- 362 [28] Neerav Karani, Ertunc Erdil, Krishna Chaitanya, and Ender Konukoglu. Test-time adaptable neural
363 networks for robust medical image segmentation. *Medical Image Analysis*, 68:101907, 2021.
- 364 [29] Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic
365 training for domain generalization. *arXiv preprint arXiv:1902.00113*, 2019.
- 366 [30] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep
367 domain generalization via conditional invariant adversarial networks. In *Proceedings of the European
368 Conference on Computer Vision (ECCV)*, pages 624–639, 2018.
- 369 [31] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro
370 Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in
371 context, 2014. cite arxiv:1405.0312Comment: 1) updated annotation pipeline description and figures; 2)
372 added new section describing datasets splits; 3) updated author list.
- 373 [32] Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre
374 Alahi. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information
375 Processing Systems*, 34, 2021.
- 376 [33] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep
377 adaptation networks. *arXiv preprint arXiv:1502.02791*, 2015.
- 378 [34] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation with
379 residual transfer networks. In *Advances in Neural Information Processing Systems*, pages 136–144, 2016.
- 380 [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint
381 arXiv:1711.05101*, 2017.
- 382 [36] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth
383 estimation. *ACM Transactions on Graphics (ToG)*, 39(4):71–1, 2020.
- 384 [37] Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain generalization using causal matching. In
385 *International Conference on Machine Learning*, pages 7313–7324. PMLR, 2021.
- 386 [38] Chengzhi Mao, Lu Jiang, Mostafa Dehghani, Carl Vondrick, Rahul Sukthankar, and Irfan Essa. Discrete
387 representations strengthen vision transformer robustness. *arXiv preprint arXiv:2111.10493*, 2021.
- 388 [39] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature
389 representation. In *International Conference on Machine Learning*, pages 10–18, 2013.
- 390 [40] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
391 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
392 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie
393 Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In
394 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in
395 Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

- 396 [41] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders:
397 Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern*
398 *recognition*, pages 2536–2544, 2016.
- 399 [42] Shiv Shankar, Vihari Piratla, Soumen Chakrabarti, Siddhartha Chaudhuri, Preethi Jyothi, and Sunita
400 Sarawagi. Generalizing across domains via cross-gradient training. *arXiv preprint arXiv:1804.10745*,
401 2018.
- 402 [43] Assaf Shocher, Nadav Cohen, and Michal Irani. “zero-shot” super-resolution using deep internal learning.
403 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3118–3126,
404 2018.
- 405 [44] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A Efros. Unsupervised domain adaptation through
406 self-supervision. *arXiv preprint*, 2019.
- 407 [45] Yu Sun, Wyatt L Ubellacker, Wen-Loong Ma, Xiang Zhang, Changhao Wang, Noel V Csomay-Shanklin,
408 Masayoshi Tomizuka, Koushil Sreenath, and Aaron D Ames. Online learning of unknown dynamics for
409 model-based controllers in legged locomotion. *IEEE Robotics and Automation Letters*, 6(4):8442–8449,
410 2021.
- 411 [46] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
412 with self-supervision for generalization under distribution shifts. In *International Conference on Machine*
413 *Learning*, pages 9229–9248. PMLR, 2020.
- 414 [47] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528.
415 IEEE, 2011.
- 416 [48] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.
- 417 [49] Igor Vasiljevic, Ayan Chakrabarti, and Gregory Shakhnarovich. Examining the impact of blur on recognition
418 by convolutional networks. *arXiv preprint arXiv:1611.05760*, 2016.
- 419 [50] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing
420 robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on*
421 *Machine Learning, ICML ’08*, page 1096–1103, New York, NY, USA, 2008. Association for Computing
422 Machinery.
- 423 [51] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time
424 adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 425 [52] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.
426 Simmim: A simple framework for masked image modeling. In *International Conference on Computer*
427 *Vision and Pattern Recognition (CVPR)*, 2022.
- 428 [53] Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural
429 networks via stability training. In *Proceedings of the IEEE conference on computer vision and pattern*
430 *recognition*, pages 4480–4488, 2016.

431 **Checklist**

- 432 1. For all authors...
- 433 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
434 contributions and scope? [Yes]
- 435 (b) Did you describe the limitations of your work? [Yes]
- 436 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 437 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
438 them? [Yes]
- 439 2. If you are including theoretical results...
- 440 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 441 (b) Did you include complete proofs of all theoretical results? [Yes]
- 442 3. If you ran experiments...
- 443 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
444 mental results (either in the supplemental material or as a URL)? [Yes] Code and data
445 will be released upon acceptance.
- 446 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
447 were chosen)? [Yes]
- 448 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
449 periments multiple times)? [N/A] Following the convention of [21] and many other
450 large-scale experiments.
- 451 (d) Did you include the total amount of compute and the type of resources used (e.g., type
452 of GPUs, internal cluster, or cloud provider)? [Yes]
- 453 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 454 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 455 (b) Did you mention the license of the assets? [Yes] In appendix
- 456 (c) Did you include any new assets either in the supplemental material or as a URL?
457 [N/A]
- 458 (d) Did you discuss whether and how consent was obtained from people whose data you're
459 using/curating? [N/A]
- 460 (e) Did you discuss whether the data you are using/curating contains personally identifiable
461 information or offensive content? [N/A]
- 462 5. If you used crowdsourcing or conducted research with human subjects...
- 463 (a) Did you include the full text of instructions given to participants and screenshots, if
464 applicable? [N/A]
- 465 (b) Did you describe any potential participant risks, with links to Institutional Review
466 Board (IRB) approvals, if applicable? [N/A]
- 467 (c) Did you include the estimated hourly wage paid to participants and the total amount
468 spent on participant compensation? [N/A]