

Parameter-free Online Test-time Adaptation

Malik Boudiaf
ÉTS Montreal *

Romain Mueller
FiveAI

Ismail Ben Ayed
ÉTS Montreal

Luca Bertinetto
FiveAI

Abstract

Training state-of-the-art vision models has become prohibitively expensive for researchers and practitioners. For the sake of accessibility and resource reuse, it is important to focus on adapting these models to a variety of downstream scenarios. An interesting and practical paradigm is online test-time adaptation, according to which training data is inaccessible, no labelled data from the test distribution is available, and adaptation can only happen at test time and on a handful of samples. In this paper, we investigate how test-time adaptation methods fare for a number of pre-trained models on a variety of real-world scenarios, significantly extending the way they have been originally evaluated. We show that they perform well only in narrowly-defined experimental setups and sometimes fail catastrophically when their hyperparameters are not selected for the same scenario in which they are being tested. Motivated by the inherent uncertainty around the conditions that will ultimately be encountered at test time, we propose a particularly “conservative” approach, which addresses the problem with a Laplacian Adjusted Maximum-likelihood Estimation (LAME) objective. By adapting the model’s output (not its parameters), and solving our objective with an efficient concave-convex procedure, our approach exhibits a much higher average accuracy across scenarios than existing methods, while being notably faster and have a much lower memory footprint. The code is available at <https://github.com/fiveai/LAME>.

1. Introduction

In recent years, training state-of-the-art models has become a massive computational endeavor for many machine learning problems (e.g. [5, 13, 38]). For instance, it has been estimated that each training of GPT-3 [5] produces an equivalent of 552 tons of CO₂, which is approximately the amount emitted in six flights from New York to San Francisco [35]. As implied in the whitepaper on “foundation

models” [4], we should expect that more and more efforts will be dedicated to the design of procedures that allow for the efficient adaptation of pre-trained large models under a variety of circumstances. In other words, these models will be “trained once” on a vast dataset and then adapted at test time to newly-encountered scenarios. Besides being important for resource reuse, being able to abstract the *pre-training stage* away from the *adaptation* is paramount in privacy-focused applications, and in any other situation in which preventing access to the training data is desirable. Towards this goal, it is important that, from the point of view of the adaptation system, there is neither access to the training data nor the training procedure of the model to adapt. With this context in mind, we are particularly interested in designing adaptation methods ready to be used in realistic scenarios, and that are suitable for a variety of models.

One aspect that many real-world applications have in common is the need to perform adaptation *online*, and with a limited amount of data. That is, we should be able to perform adaptation while the data is being received. Take for instance the vision model with which an autonomous vehicle or a drone may be equipped. At test-time, it will ingest a video stream of highly-correlated data (non-i.i.d.), which could be used for adaptation. We would like to be confident that leveraging this information will be useful, and not destructive, no matter the type of domain shift that may exist between training and test data. Such shifts could be, for instance, “low-level” (e.g. the data stream is affected by snowy weather which has never been encountered during the California-sunlit training stage), or “high-level” (e.g. the data include the particular Art Deco architecture of Miami Beach’s Historic District), or even a combination of both. To summarize, we are interested in the design of test-time adaptation systems that 1) are unsupervised; 2) can operate online and on potentially non-i.i.d. data; 3) assume no knowledge of training data or training procedure; and 4) are not tailored to a certain model, so that the progress made by the community can be directly harnessed.

This problem specification falls under the *fully test-time adaptation* paradigm studied in a handful of recent works [1, 27, 29, 56], where simple techniques like test-time

*Work done as part of a research internship at FiveAI. Corresponding author: malik.boudiaf.1@etsmtl.net

learning of batch normalization’s scale and bias parameters [56] have proven to be very effective in some scenarios, like the one represented by low-level corruptions [17]. In our experimental results, we observe that existing methods [25, 27, 29, 56] have to be used with great care in uncertain yet realistic situations because of their sensitivity to variables such as the model to adapt or the type of domain shift. As a matter of fact, we show that, when selecting their hyperparameters to maximize the average accuracy over a number of scenarios, existing methods do not outperform a non-adaptive baseline. For them to perform well, hyperparameters need to be adjusted in a scenario-specific fashion. However, this is clearly not an option when the test-time conditions are unknown in advance.

These findings suggest that, while being agnostic to both training and testing circumstances is important, it is wise to approach the problem of test-time adaptation prudently. Instead of adapting the parameters of a pre-trained model, we only adapt its *output* by finding the latent assignments that optimize a manifold-regularized likelihood of the data. The manifold-smoothness assumption has been successful in a wide range of other problems, including graph clustering [45, 46, 52], semi-supervised learning [2, 7, 19], and few-shot learning [62], as it enforces desirable and general properties on the solutions. Specifically, we embed Laplacian regularization as a corrective term, and derive an efficient concave-convex procedure for optimizing our overall objective, with guaranteed convergence. When aggregating over different conditions, this simple and “conservative” strategy significantly improves both over the non-adaptive baseline and existing test-time adaptation methods in an extensive set of experiments covering 7 datasets, 19 shifts, 3 training strategies and 5 network architectures. Moreover, by virtue of not performing model adaptation but only output correction, it reduces *by half* both the total inference time and the memory footprint compared to existing methods.

2. Related work

In general, domain adaptation aims at relaxing the assumption that “train and test distributions should match”, which is at the foundation of most machine learning algorithms. Since real-world applications rarely reflect the textbook assumption, this relaxation has generated a lot of interest and motivated a large corpus of work. Doing this topic justice would take several surveys (e.g. [10, 34, 57, 58]), and it is unfeasible given this paper format. Instead, in this section we aim at describing the overall problem setups that are more closely relevant to ours.

The applicability of early works in domain adaptation was limited, in that methods required access to the target domain [34] during training. **Unsupervised domain adaptation** [58] makes the scenario slightly more realistic by not requiring labels from the target domain. Two common gen-

eral strategies are, for instance, explicitly learning domain-invariant feature representations by minimizing some measure of divergence between source and target distributions (e.g. [20, 30, 49]); or embedding a “domain discriminator” component in the network and then penalizing its success in the loss (e.g. [14, 37]). Still, the necessity of having access, during training, to *both* source and target domains limits the usability of this class of methods.

Domain generalization (DG) foregoes the need to access the target distribution by learning a model from multiple domains, with the intent of generalizing to unseen ones [57]. Popular strategies to address this problem include: increasing the diversity of training data via either augmentations (e.g. [36, 54]), adversarial learning (e.g. [55, 61]), or generative models (e.g. [39, 47]); learning domain-invariant representations [3], and decoupling the domain-specific and domain-independent components (e.g. [18, 21, 32]). Notably, the recent work of Gulrajani & Lopez-Paz [15] showed on a large testbed that learning a vanilla classifier on a pool of datasets outperformed all modern techniques, thus sending a strong message on the importance of a carefully designed experimental protocol.

Despite the shared goal of generalizing across domains and the constraint of not having access to the target distributions in advance, one fundamental difference of DG with the setup we consider is the lack of test-time adaptability. Instead, methods falling under the **source-free domain adaptation** paradigm [9] require no access to the training data *during the process of adaptation*. Liang *et al.* [28] assume only to have access to the source dataset’s summary statistics, and relate the models fitting the source and target domains by surmising that class centroids are only moderately shifted between the two datasets. Before adaptation, Kundu *et al.* [22, 23] consider a first “vendor-side” phase, during which the target domain is not known and a model is trained on an augmented training dataset aiming at mimicking possible domain shifts and category gaps that will be encountered downstream. Li *et al.* [26] propose the Collaborative Class Conditional GAN, which integrates the output of a prediction model into the loss of the generator to produce new samples in the style of the target domain, which are in turn used to adapt the model via backpropagation. In *Test-time Training* [50], Sun *et al.* perform test-time adaptation via self-supervision by jointly optimizing two branches (one supervised and one self-supervised) during training.

While being vastly more practical than the ones addressing vanilla domain adaptation, the methods listed above are still quite limited in that they typically have an ad-hoc training procedure. As mentioned in Section 1, we would like to facilitate model reuse, so that the progress made by the community in architecture design [13], self-supervised learning [8] or multi-modal learning [38] can be directly exploited. Our setup is mostly similar to what

has been referred to in the TENT paper [56] as the **fully test-time adaptation** scenario. In this case, the intent is to perform unsupervised test-time adaptation while “*not restricting or altering model training*” [56]. In TENT, this is achieved with a simple entropy minimization loss, which informs the optimization of scale and bias parameters of batch normalization layers. As for batch normalization layers’ statistics, they are re-estimated on the test data, similarly to what is done in adaptive batchnorm (AdaBN) methods [6, 27, 31, 44], which have shown strong performance on the perturbations of ImageNet-C [17]. In similar spirit, Liang *et al.* [29] updates the parameters of the feature extractor of a given model by maximizing a mutual information objective (SHOT-IM).

Although we share many of the motivations presented in TENT and SHOT, we believe that our work differs under two main aspects. First, given our model-independence desideratum, we explicitly study the extent to which our approach works across training strategies and architectures. This analysis is missing in prior works: as we will see in Section 6, the type of model being adapted is a variable that strongly affects the effectiveness of both TENT and SHOT. Second, for the sake of usability, we are particularly focused on *online* adaptation, which leads us to also consider non-i.i.d. scenarios as an important part of our evaluation.

3. Problem Formulation

In (fully) test-time adaptation [29, 56] (TTA), we have access to a parametric model $q_{\theta}(y|\mathbf{x})$ trained on an inaccessible labelled source dataset $\mathcal{D}_s = \{(\mathbf{x}, y) \sim p_s(\mathbf{x}, y)\}$, where \mathbf{x} is an image and $y \in \mathcal{Y}$ its associated label from the set of source classes \mathcal{Y} . Additionally, we consider an unlabelled target dataset sampled from an arbitrary target distribution $\mathcal{D}_t = \{\mathbf{x} \sim p_t(\mathbf{x})\}$. We take the standard *covariate shift* assumption [48] that $p_s(y|\mathbf{x}) = p_t(y|\mathbf{x})$ and $p_s(\mathbf{x}) \neq p_t(\mathbf{x})$, which implies that shifts can only happen if there exists some class y such that $p_s(y)p_s(\mathbf{x}|y) \neq p_t(y)p_t(\mathbf{x}|y)$. This leads us to consider two types of shift throughout this work: the *prior shift*, in which $p_t(y)$ differs from $p_s(y)$, and the *likelihood shift*, in which $p_t(\mathbf{x}|y)$ differs from $p_s(\mathbf{x}|y)$.

As the target distribution shifts from the source, the parametric model $q_{\theta}(y|\mathbf{x})$ no longer necessarily well approximates the true, domain-invariant distribution $p(y|\mathbf{x})$. A toy illustration of this phenomenon can be found in Fig. 2, where the linear classifier can only properly model the true sinusoidal distribution over a limited region of the input space. Therefore, TTA methods aim at adapting $q_{\theta}(y|\mathbf{x})$ to maximize its predictive performance on the target distribution. In particular, we focus on the *online* setting, where the classifier receives a potentially non-i.i.d. stream of target samples, and must simultaneously adapt and predict.

Typical large-scale datasets contain up to tens of thousands of classes, and have been created with the purpose of

covering a large portion of the concepts that may be of interest at test-time. As such, they likely contain classes of a finer or equal (but not coarser) granularity than those required in specific TTA scenarios. Therefore, to make our setting more practical, we relax the common assumption that source classes must coincide with the target ones. Instead, we allow target classes to be superclasses, according to some pre-defined hierarchy. Authors from [53] handle this by max-pooling the softmax predictions across associated subclasses, but we empirically found average-pooling to perform slightly better, and decided to proceed with this strategy. More details in Appendix.

4. On the Risks of Network Adaptation

In order to better approximate the underlying distribution $p(z|\mathbf{x})$ at test-time, TTA methods usually propose to directly modify the parametric source model. We group such methods under the term *Network Adaptation Methods* (NAMs). Specifically, such methods [29, 56] first partition the network into *adaptable* weights θ^a and frozen weights θ^f , and proceed by minimizing an unsupervised loss $\mathcal{L}(\mathbf{x}; \theta^a \cup \theta^f)$, $\mathbf{x} \sim p_t(\mathbf{x})$ w.r.t. θ^a . TTA methods mostly differ based on their choices of partition $\{\theta^f, \theta^a\}$ and loss function \mathcal{L} . For instance, TENT [56] only adapts the scale and bias parameters (γ, β) of the batch normalization (BN) layers through entropy minimization, while SHOT [29] adapts the convolutional filters of the model through mutual information maximization.

While NAMs have the potential to substantially improve the performance of a model on the target samples, they also run the risk of dramatically degrading it. Consecutive updates of the adaptable weights θ^a on narrow portions of the target distribution can cause the model to overspecialize. Such behavior can be caused by the combination of a sub-optimal choice of hyperparameters for a specific scenario and the lack of sample diversity at the batch level. Note that the latter does not arise exclusively in video scenarios, but also in situations characterized by a high class imbalance. Moreover, adapting parameters across the network and within an iterative optimization procedure such as SGD (which spans many batches of data), can inherently lead to the degeneration of the model over time. To make this more intuitive, in Fig. 1 we showcase a failure mode of the widely used entropy minimization principle. In a low intra-batch diversity situation, entropy minimization can degenerate the model *silently*. In other words, it can fail without exhibiting any distinctive behaviour that, in the absence of labels, would allow for a clear diagnosis. An illustrative explanation of this phenomenon is conveyed in Fig. 2.

One may argue that choosing optimal hyperparameters may solve the problems mentioned above. However, tuning hyperparameters separately for each target scenario would require access to the labels. Moreover, this approach would

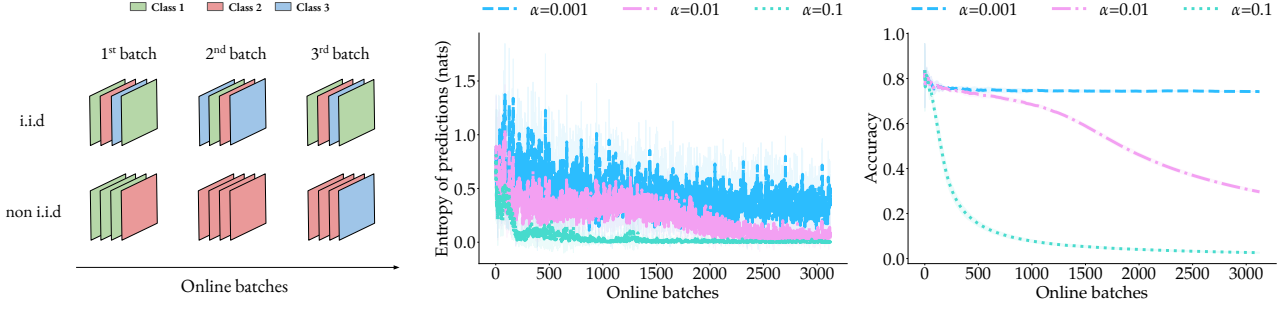


Figure 1. **Adaption through entropy minimization in a non-i.i.d. scenario may *silently* degenerate the model.** (Left) Non-i.i.d. streams are generated by batching samples according to their class. (Middle) The conditional entropy of predictions is being minimized in an online fashion on such non-i.i.d. streams. However, assessing whether the adaptation is being beneficial or detrimental solely from these curves is impractical in an unsupervised scenario. (Right) Rather, monitoring the online accuracy (which would require access to the labels) would reveal that the model is actually collapsing for two out of three learning rates considered.

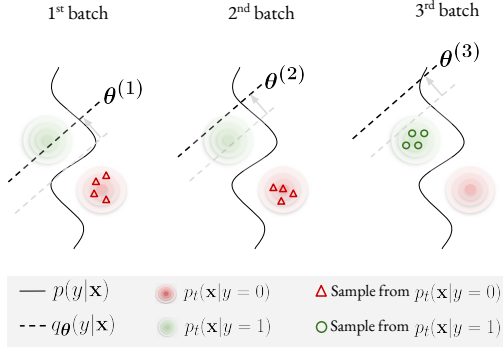


Figure 2. Minimizing the conditional entropy (as in TENT [56]) encourages the model $q_{\theta}(y|x)$ to produce high confidence predictions. Geometrically, this corresponds to increasing the margin between the decision boundary and the samples from the current batch. In the low-diversity scenario depicted above, the 1st and 2nd batches only contain red samples. This causes the boundary to move away from red samples. When green samples are finally observed in the 3rd batch, the boundary has gone past the green cluster, so that samples are (wrongly) assigned to the red class.

also require to know which scenario is going to be encountered at test time. These two points defeat the whole purpose of the TTA paradigm. Therefore, it would be desirable for NAMs’ hyperparameters to generalize well *across* scenarios. However, keeping the entropy minimization approach of TENT [56] as an example, we show on the left matrix of Fig. 3 that such generalization is, in practice, far from fulfilled. More specifically, to obtain this matrix, we created a series of 12 validation scenarios (see Section 6), providing a wide coverage of the shifts discussed in Section 3. Row i is to be read in the following way: we tune hyperparameters considering only scenario i , and then observe to which extent this choice of hyperparameters generalizes to all scenarios $j \in \{1, \dots, 12\}$. The absolute improvement

(or degradation) w.r.t. the performance of the non-adapted model is reported in the matrix. The clear trend emerging from Fig. 3 is that the entropy minimization approach is severely brittle w.r.t. its hyperparameter configuration, especially in non-i.i.d. and class-imbalanced scenarios, where a sub-optimal choice can degrade the model’s accuracy by up to an absolute 66% compared to the non-adaptive baseline. We emphasize that Fig. 3 only shows validation results obtained when using scenario-specific hyperparameters, and therefore only serves the purpose of empirically demonstrating the issue with over-specific hyperparameters. In Appendix, we show that the same trend can be observed for all NAMs we experimented with.

As an alternative, in Section 5 we propose an adaptation strategy which only affects the output of the model (not its parameters), only considers one batch of data at a time, and has only one hyperparameter to tune.

5. The LAME method

In order to address the aforementioned issues, we introduce a method that only aims at providing a *correction* of the output probabilities of a classifier instead of modifying the internal parameters of its feature extractor. On the one hand, freezing the source classifier prevents our method from accumulating knowledge across batches. On the other, it mitigates the risk of degenerating the classifier, reduces compute requirements (as gradients are neither computed nor stored), and inherently removes the need for searching over delicate hyperparameters such as learning rate or momentum of the optimizer. Overall, we empirically demonstrate that such an approach is more reliable and practical than NAMs when the test-time conditions are unknown.

Formulation. Assume we are given a batch of data sampled from the target distribution $\mathbf{X} \in \mathbb{R}^{N \times d} \sim p_t^N(\mathbf{x})$, with N the number of samples and d the feature dimension. Our method finds a latent assignment vector $\tilde{\mathbf{z}}_i =$

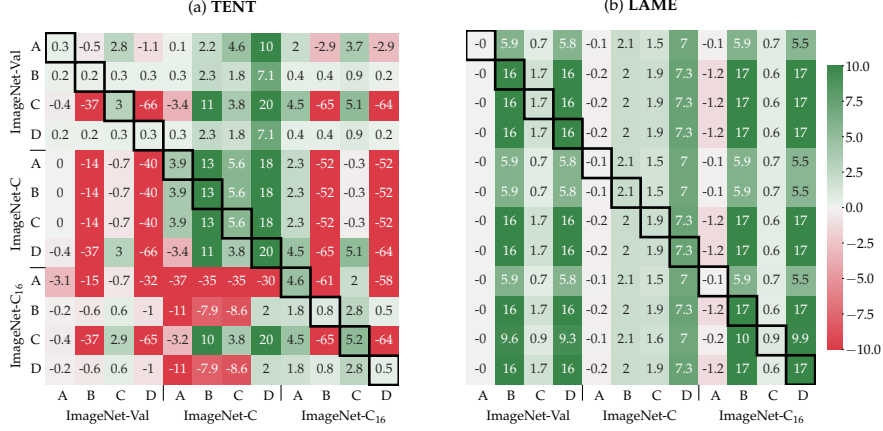


Figure 3. *Cross-shift* validation for TENT [56] (left) and our proposed LAME (right). A cell at position (i, j) shows the absolute improvement (or degradation) of the current method w.r.t. to the baseline when using the optimal hyperparameters for scenario i , but evaluating in scenario j . Legend: A = i.i.d., B = non i.i.d., C = i.i.d. + prior shift, D = non i.i.d. + prior shift. More details on the scenarios in Sec. 6

$(\tilde{z}_{ik})_{1 \leq k \leq K} \in \Delta^{K-1}$ for each data point \mathbf{x}_i , which aims to approximate the true distribution $p(z|\mathbf{x})$, with K the number of classes and $\Delta^{K-1} = \{\tilde{\mathbf{z}} \in [0, 1]^K \mid \mathbf{1}^T \tilde{\mathbf{z}} = 1\}$ the probability simplex. A principled way to achieve this is to find assignments $\tilde{\mathbf{Z}}$ that maximize the log-likelihood of the data subject to simplex constraints $\tilde{\mathbf{z}}_i \in \Delta^{K-1}, \forall i$:

$$\mathcal{L}(\tilde{\mathbf{Z}}) = \log \left(\prod_{i=1}^N \prod_{k=1}^K p(\mathbf{x}_i, k)^{\tilde{z}_{ik}} \right) \stackrel{c}{=} \sum_{i=1}^N \tilde{\mathbf{z}}_i^T \log(\mathbf{p}_i) \quad (1)$$

where $\tilde{\mathbf{Z}} \in [0, 1]^{N \times K}$ is the vector that concatenates all assignment vectors $\tilde{\mathbf{z}}_i$, $\mathbf{p}_i = (p(k|\mathbf{x}_i))_{1 \leq k \leq K} \in \Delta^{K-1}$, and $\stackrel{c}{=}$ stands for equality up to an additive constant. In order to prevent over-confident assignments, we consider a negative-entropy regularization that discourages one-hot assignments for $\tilde{\mathbf{Z}}$. Note that such regularization also acts as a barrier that restricts the domain of $\tilde{\mathbf{z}}_i$ to non-negative values, hence implicitly handling the $\tilde{\mathbf{z}}_i \geq 0$ constraint. Maximizing the regularized log-likelihood objective therefore amounts to minimizing the following Kullback–Leibler (KL) divergences subject to $\mathbf{1}^T \tilde{\mathbf{z}}_i = 1, \forall i$:

$$-\sum_{i=1}^N \tilde{\mathbf{z}}_i^T \log(\mathbf{p}_i) + \sum_{i=1}^N \tilde{\mathbf{z}}_i^T \log(\tilde{\mathbf{z}}_i) = \sum_{i=1}^N \text{KL}(\tilde{\mathbf{z}}_i \parallel \mathbf{p}_i) \quad (2)$$

Problem (2) is minimized for $\tilde{\mathbf{z}}_i = \mathbf{p}_i, \forall i$. Taking a step back, we don't have access to \mathbf{p}_i , but only to the source parametric model $\mathbf{q}_i = (q_\theta(k|\mathbf{x}_i))_{1 \leq k \leq K}$ which, recall, might be a poor approximation of the true distribution when evaluated on target samples $\mathbf{x} \sim p_t(\mathbf{x})$. In fact, simply replacing \mathbf{p}_i by \mathbf{q}_i in Eq. (2) yields the predictions from the source model as optimum: $\tilde{\mathbf{z}}_i = \mathbf{q}_i$.

To compensate for the inherent error of this approximation, we focus on Laplacian regularization, which en-

courages neighbouring points in the feature space to have consistent latent assignments. Laplacian regularization is widely used in semi-supervised learning [2, 7, 19], where it is optimized jointly with supervised losses over labelled data points, or in graph clustering [45, 46, 52], where it is optimized subject to class-balance constraints. The TTA problem is different as, unlike semi-supervised learning, cannot count on any supervision and, unlike clustering, class-balance constraints are irrelevant (or even detrimental). Hence, we introduce Laplacian Adjusted Maximum-likelihood Estimation (LAME), which minimizes the likelihood in (2) jointly with a Laplacian correction, subject to constraints $\mathbf{1}^T \tilde{\mathbf{z}}_i = 1, \forall i$:

$$\mathcal{L}^{\text{LAME}}(\tilde{\mathbf{Z}}) = \sum_i \text{KL}(\tilde{\mathbf{z}}_i \parallel \mathbf{q}_i) - \sum_{i,j} w_{ij} \tilde{\mathbf{z}}_i^T \tilde{\mathbf{z}}_j \quad (3)$$

where $w_{ij} = w(\phi(\mathbf{x}_i), \phi(\mathbf{x}_j))$, with ϕ denoting our pre-trained feature extractor and w is a function measuring the affinity between $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$. The closer the points in the feature space, the higher their affinity. Clearly, when the affinity is high (w_{ij} is large), minimizing the Laplacian term in (3) seeks the largest possible value of dot product $\tilde{\mathbf{z}}_i^T \tilde{\mathbf{z}}_j$, thereby assigning points i and j to the same class. Therefore, our model in (3) could be viewed as a graph clustering of the batch data, penalized by a KL term discouraging substantial deviations from the source-model predictions.

Efficient optimization via a concave-convex procedure.

In what follows, we show that our Problem (3) can be minimized using the Concave-Convex Procedure (CCCP) [60], which allows us to obtain a highly efficient iterative algorithm, with convergence guarantee. Each iteration updates the current solution $\tilde{\mathbf{Z}}^{(n)}$ as the minimum of a tight upper bound on the objective. This guarantees that the objective does not increase at each iteration. For the sum of a con-

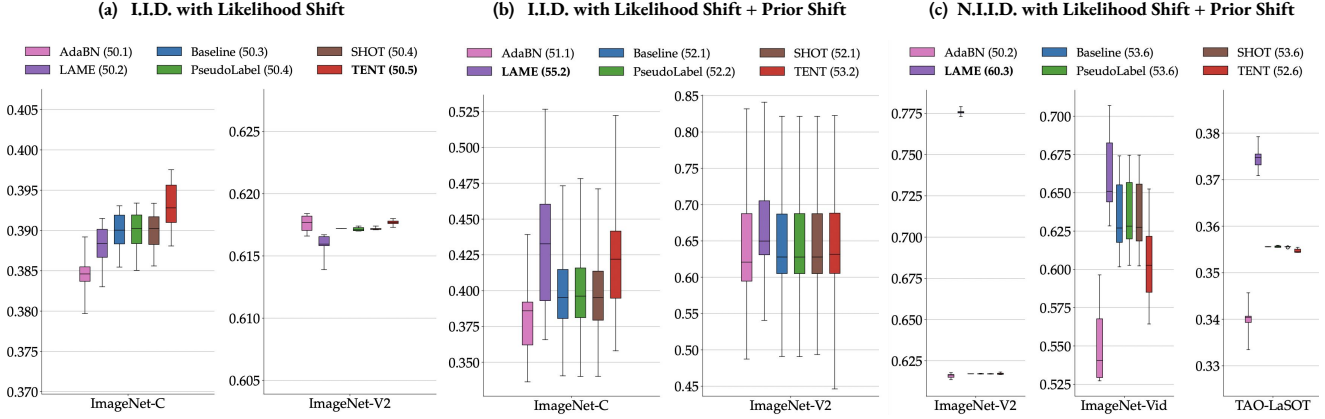


Figure 4. Results across the 7 testing scenarios, using the same Original RN-50 that was used for validation. The average for each scenario is reported in the legend. The batch size is 64. Each experiment is run 10 times with different random seeds. Experiments with *prior shift* tend to exhibit larger variance, since each random run uses new class proportions, each sampled from a Zipf distribution.

cave and a convex function, as is the case of our objective in (3), a CCCP replaces the concave part by its linear first-order approximation at the current solution, which is a tight upper bound, while keeping the convex part unchanged. In our case, the Laplacian term is concave when the affinity matrix $W = [w_{i,j}]$ is positive semi-definite, while the KL term is convex. The concavity of the Laplacian for positive semi-definite W could be verified by re-writing the term as follows¹: $-\sum_{i,j} w_{ij} \tilde{\mathbf{z}}_i^T \tilde{\mathbf{z}}_j = -\tilde{\mathbf{Z}}^T (W \otimes I) \tilde{\mathbf{Z}}$, where \otimes denotes the Kronecker product and I is the K -by- K identity matrix. We thus replace the Laplacian term in (3) by $-((W \otimes I) \tilde{\mathbf{Z}}^{(n)})^T \tilde{\mathbf{Z}}$, which yields the following tight upper bound, up to an additive constant independent of $\tilde{\mathbf{Z}}$:

$$\mathcal{L}^{\text{LAME}}(\tilde{\mathbf{Z}}) \leq \sum_i^c \text{KL}(\tilde{\mathbf{z}}_i || \mathbf{q}_i) - ((W \otimes I) \tilde{\mathbf{Z}}^{(n)})^T \tilde{\mathbf{Z}} \quad (4)$$

Solving the Karush-Kuhn-Tucker (KKT) conditions corresponding to minimizing convex upper bound (4), subject to constraints $\mathbf{1}^T \tilde{\mathbf{z}}_i = 1, \forall i$, yields the following decoupled updates of the assignment variables:

$$\tilde{z}_{ik}^{(n+1)} = \frac{q\theta(k|\mathbf{x}_i) \exp(\sum_j w_{ij} \tilde{z}_{jk}^{(n)})}{\sum_{k'} q\theta(k'|\mathbf{x}_i) \exp(\sum_j w_{ij} \tilde{z}_{jk'}^{(n)})} \quad (5)$$

which have to be iterated until convergence. The full derivation of Eq. (5) is provided in Appendix.

6. Experimental design

The design of our experimental protocol is mainly guided by the desire to assess both model and domain independence of TTA methods. For model independence, we need to evaluate the performance of methods under a variety of pre-trained models. As for domain independence,

¹ W positive semi-definite implies $W \otimes I$ positive semi-definite.

a single fixed trained model must allow to evaluate a TTA method under multiple adaptation scenarios. This implies that the source classes encoded in the pre-trained model must be able to adequately cover the classes of interest that may be encountered at test time. Note that, in practice, this is a reasonable requirement, as modern large-scale datasets span tens of thousands of classes [24, 42, 43, 59].

Networks. Because of their popularity within the community and the large number of classes covered, ImageNet-trained models represent an ideal playground for our experiments. In particular, they allow to evaluate model independence along two axis. First, with respect to the training procedure, by experimenting with the same ResNet-50 architecture (RN-50 herein), but trained in three different ways: the original release from Microsoft Research Asia (MSRA) [16], Torchvision’s [33], and using the self-supervised SimCLR [8]. Second, with respect to the architecture itself, by providing results on 5 different backbones, including RN-18, RN-50, RN-101, EfficientNet (EN-B4) [51] and the recent Vision Transformer ViT-B [13]. All models used were trained on the standard ImageNet ILSVRC-12 training set, except for ViT-B which uses an additional ImageNet-21k [12] pre-training step.

Hyperparameter search. For validation purposes, we consider 3 datasets. First, we use the original validation set of ImageNet [43]. To represent *likelihood shift*, we consider *ImageNet-C-Val*, which augments the original images with 9 realistic perturbations of varying intensity (the other 10 from the original ImageNet-C [17] are reserved for testing). Finally, we consider *ImageNet-C₁₆*, a variant of *ImageNet-C* that simulates an easier but practical scenario where a subset of ImageNet classes is mapped to 16 superclasses. By reducing the total number of classes, ImageNet-C₁₆ also reduces class diversity at the batch level, which we identi-

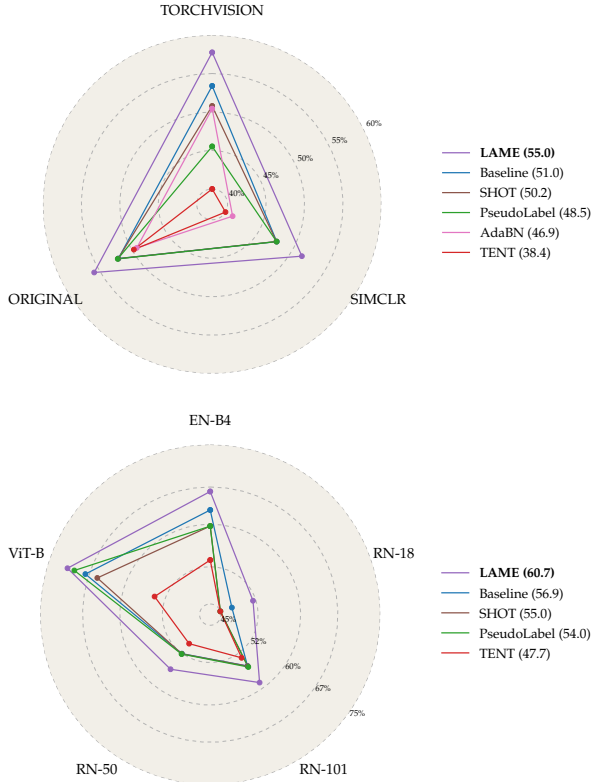


Figure 5. **Transferability of hyperparameters across models.** For each TTA method, we use the optimal set of hyperparameters obtained during validation and using the original release of RN-50 [16] as backbone. Each vertex on the chart represents the average across our 7 test scenarios for a specific architecture. The values in the legend represent the average over all the vertices. (Top): We test these hyperparameters using the same backbone but different training procedures. *Torchvision* refers to the model available in PyTorch’s model zoo, *SimCLR* to the model obtained from the self-supervised approach of [8], and *original* refers to the same model used to choose the hyperparameters. (Bottom): The *same* set of hyperparameters is used for different architectures, ranging from a RN-18 to the recent vision transformer ViT-B [13]. To allow similar setups across architectures, a batch size of 16 is used to generate the above results.

fied as a potentially critical factor for NAMs approaches in Section 4. In order to mimic realistic prior shifts, we modify the class ratios to follow a Zipf distribution [41]. Finally, to cover non-i.i.d. scenarios, we present the model with a sequence of “tasks”, where each task either represents a set of samples perturbed with the same corruption (in the case of ImageNet-C), or belonging to the same class otherwise. All the combinations of 3 datasets, 2 prior shifts (with and without Zipf-unbalanced class distribution) and 2 sampling schemes (i.i.d. or non-i.i.d.) add up to the 12 validation scenarios. For each method, a grid-search over salient hyperparameters is carried out, and the single hyperparameter set

that obtains best average performance over the 12 validation scenarios is selected, and **kept fixed for test experiments in Fig. 4 and 5**. The exact definition of the grid-search for each method is available in the Appendix.

Testing. For testing, we design 4 i.i.d. and 3 non-i.i.d. test scenarios. For the i.i.d. cases, we use the 4 combinations obtained by coupling ImageNet-C-Test and ImageNet-V2 [40] with the presence or absence of Zipf class-imbalance. As for the 3 non-i.i.d. scenarios, we use again ImageNet-V2 (with a different split), along with two video datasets: ImageNet-VID [43] and the LaSOT subset from TAO [11]. Keeping the idea of feeding the model with a sequence of tasks, video datasets allow us to evaluate realistic scenarios by simply grouping frames from the same video together. We use 10 random runs for each test experiment. More details on all datasets (and class mappings) in Appendix.

Methods. As a first baseline, we evaluate the source-trained model without any adaptation, referred to as *Baseline*. For Network Adaptation Methods (NAMs), we reproduce and evaluate four state-of-the-art TTA methods that can be run in an online fashion: *TENT* [56] based on entropy minimization, *SHOT-IM* based on mutual information maximization, *PseudoLabel* [25] based on min-entropy minimization and *AdaBN* [27] based on batch normalization statistics alignment. Finally, we evaluate LAME.

7. Experimental results

Towards domain-independent test-time adaptation. As motivated in Section 4, most scenario-sensitive hyperparameters come from the optimization of the network. By virtue of completely freezing the classifier, our LAME approach is free of such burden. Instead, LAME only tries to find optimal shallow assignments through a bound-optimization procedure that does not introduce any hyperparameter. Therefore, we are only left with the tuning of the affinity function w from Eq. (3), which is less sensitive than the optimization-related hyperparameters of NAMs. This claim is first supported by inspecting LAME’s *cross-shift* validation matrix, already used earlier to illustrate NAMs’ brittleness. Looking this time at the right plot of Fig. 3, we can see drastic improvements both in terms of average performance and worst-case degradation across all cases w.r.t. TENT.² A second empirical evidence supporting this claim comes from the results on the test scenarios, shown in Fig. 4. Consistent with the validation results in Fig. 3, Fig. 4 confirms that LAME does not help in standard i.i.d. *likelihood shifts*, and fares around 0.5% below the baseline in worst cases. However, when *prior shifts* are introduced, NAMs’ performance does not improve over the

²We speculate that introducing more hyperparameters in LAME (*e.g.* weighting the different terms of our loss) would result in worse off-diagonal terms in Fig. 3, but also higher overall performance.

baseline, whereas LAME exhibits very noticeable improvements. This is particularly evident in non-i.i.d. scenarios, where the average improvement is of (absolute) 6.7%, and goes up to 15% in the case of ImageNet-v2. Note that such improvement comes almost independently of the batch size used, as shown in Appendix.

NAMs are brittle w.r.t. the training procedure. As for model-independence, we first inspect whether methods are robust to changes to the training procedure. Such robustness is desired, for instance, in the case where the provider of the source model has released an update: in such a case, a TTA method should not require a new round of validation. As a first scenario, we investigate whether the set of hyperparameters obtained using the Original RN-50 [16] generalizes to the same methods, but when using the RN-50 provided by Torchvision. Given that both models were trained with standard supervision and minor experimental differences, one would expect the optimal set of hyperparameters to be very similar in the two cases. Results on the top chart of Fig. 5 suggest quite the opposite. While LAME preserves the same improvement w.r.t. the baseline, all NAMs lose significant ground, with TENT performing particularly poorly. We further experiment with a RN-50 trained using the self-supervised *SimCLR*, and observe that LAME once again retains its relative improvement of 4% w.r.t. the baseline, with no other method beating it.

LAME generalizes across architectures, NAMs don’t. Generalizing across different architectures should be a desirable property for any TTA method. In particular, for very large models, an exhaustive validation can become prohibitively expensive, thus making “model plug-and-play” an attractive feature. Results using five architectures (EfficientNet-B4 [51], the three ResNet variants, and the larger ViT-B [13] transformer) are shown on the bottom chart of Fig. 5. Across the board, LAME is the only method able to retain a consistently significant improvement w.r.t. the baseline, which remains a better option than any of the NAMs, especially with small backbones such as RN-18.

LAME runs twice as fast, while requiring twice less memory than NAMs. Provided that several direct applications of test-time adaptation involve real-time adaptation to data streams, the ability to run as efficiently as possible can also be a critical factor for practitioners. To measure runtimes, we divide inference into 3 stages: 1st forward, optimization (corresponding to SGD for NAMs and to the bound-optimization procedure of Section 5 for LAME), and 2nd forward (only needed for methods that modify the parameters of the model). Altogether, these three contributions account for the total runtime of each method. Results provided in Fig. 6 testify the clear advantage of LAME over the representative TENT (runtimes of other NAMs were found roughly similar to TENT). Memory-wise, LAME does not require to keep any gradient or intermediary buffer,

which roughly halves the amount of GPU memory needed w.r.t. NAMs.

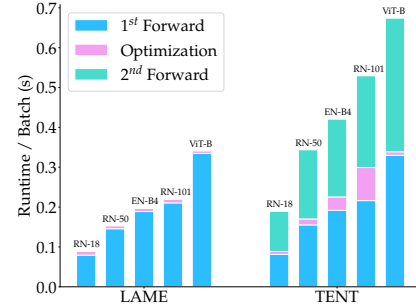


Figure 6. **Runtime per batch** of LAME vs TENT for 5 different backbones: RN-18, RN-50, EN-B4, RN-101 and ViT-B. Batch 64 is used for the RN-* family, and 16 for EN-B4 and ViT-B (as both use 380x380 images instead of 224x224). LAME provides corrected outputs without requiring a second forward pass.

8. Conclusion

Motivated by the high cost of training new models, we proposed a novel approach for online test-time adaptation (TTA) that is agnostic to both training and testing conditions. We introduced an extensive experimental protocol covering several datasets, realistic shifts and models, and evaluated existing TTA approaches by making sure that test-time domain information would not leak to inform the hyperparameters’ choice. Across the board, these methods underperform a non-adaptive baseline and can even lead to a catastrophic degradation of performance. We identified over-adaptation of the model parameters as a strong suspect for the poor performance of these methods, and opted for a more conservative approach that only corrects the output of the model. We proposed Laplacian Adjusted Maximum-likelihood Estimation (LAME), an unsupervised objective that finds the optimal set of latent assignments by discouraging deviations from the prediction of the pre-trained model, while at the same time encouraging label propagation under the manifold smoothness assumption. Averaging accuracy over the many scenarios considered, LAME outperforms all existing methods and the non-adaptive baseline, while requiring less compute and memory. Nonetheless, being restricted to the classifier’s output, LAME is also inherently limited. For one, it does not noticeably help in standard i.i.d. and class-balanced scenarios. We hope that our work will motivate further developments in this line of research. In particular, we believe that methods adopting a hybrid adaptation/correction approach, if choosing their hyperparameters under a strict regime, will have the potential to effectively tackle an even wider variety of scenarios.

References

- [1] Fatemeh Azimi, Sebastian Palacio, Federico Raue, Jörn Hees, Luca Bertinetto, and Andreas Dengel. Self-supervised test-time adaptation on video data. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022. [1](#)
- [2] Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of machine learning research*, 7(11), 2006. [2](#), [5](#)
- [3] Shai Ben-David, John Blitzer, Koby Crammer, Fernando Pereira, et al. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 2007. [2](#)
- [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. [1](#)
- [5] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. [1](#)
- [6] Collin Burns and Jacob Steinhardt. Limitations of post-hoc feature alignment for robustness. In *CVPR*, 2021. [3](#), [12](#)
- [7] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-Supervised Learning*. The MIT Press, 1st edition, 2010. [2](#), [5](#)
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. [2](#), [6](#), [7](#)
- [9] Boris Chidlovskii, Stephane Clinchant, and Gabriela Csurka. Domain adaptation in the absence of source domain data. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016. [2](#)
- [10] Gabriela Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint arXiv:1702.05374*, 2017. [2](#)
- [11] Achal Dave, Tarasha Khurana, Pavel Tokmakov, Cordelia Schmid, and Deva Ramanan. Tao: A large-scale benchmark for tracking any object. In *ECCV*, pages 436–454. Springer, 2020. [7](#), [13](#)
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. [6](#)
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. [1](#), [2](#), [6](#), [7](#), [8](#)
- [14] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*. PMLR, 2015. [2](#)
- [15] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *ICLR*, 2021. [2](#)
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. [6](#), [7](#), [8](#)
- [17] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *ICLR*, 2019. [2](#), [3](#), [6](#)
- [18] Maximilian Ilse, Jakub M Tomczak, Christos Louizos, and Max Welling. Diva: Domain invariant variational autoencoders. In *Medical Imaging with Deep Learning*. PMLR, 2020. [2](#)
- [19] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *CVPR*, 2019. [2](#), [5](#)
- [20] Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. Contrastive adaptation network for unsupervised domain adaptation. In *CVPR*, 2019. [2](#)
- [21] Aditya Khosla, Tinghui Zhou, Tomasz Malisiewicz, Alexei A Efros, and Antonio Torralba. Undoing the damage of dataset bias. In *ECCV*, 2012. [2](#)
- [22] Jogendra Nath Kundu, Naveen Venkat, R Venkatesh Babu, et al. Universal source-free domain adaptation. In *CVPR*, 2020. [2](#)
- [23] Jogendra Nath Kundu, Naveen Venkat, Ambareesh Revanur, R Venkatesh Babu, et al. Towards inheritable models for open-set domain adaptation. In *CVPR*, 2020. [2](#)
- [24] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020. [6](#)
- [25] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896, 2013. [2](#), [7](#)
- [26] Rui Li, Qianfen Jiao, Wenming Cao, Hau-San Wong, and Si Wu. Model adaptation: Unsupervised domain adaptation without source data. In *CVPR*, 2020. [2](#)
- [27] Yanghao Li, Naiyan Wang, Jianping Shi, Xiaodi Hou, and Jiaying Liu. Adaptive batch normalization for practical domain adaptation. *Pattern Recognition*, 80, 2018. [1](#), [2](#), [3](#), [7](#), [11](#)
- [28] Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan. Distant supervised centroid shift: A simple and efficient approach to visual domain adaptation. In *CVPR*, 2019. [2](#)
- [29] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *ICML*, pages 6028–6039. PMLR, 2020. [1](#), [2](#), [3](#), [12](#)
- [30] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. [2](#)
- [31] Zachary Nado, Shreyas Padhy, D Sculley, Alexander D’Amour, Balaji Lakshminarayanan, and Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate shift. *arXiv preprint arXiv:2006.10963*, 2020. [3](#)

- [32] Li Niu, Wen Li, and Dong Xu. Multi-view domain generalization for visual recognition. In *CVPR*, 2015. 2
- [33] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. pages 8024–8035, 2019. 6
- [34] Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 2015. 2
- [35] David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluís-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. Carbon emissions and large neural network training. *arXiv preprint arXiv:2104.10350*, 2021. 1
- [36] Aayush Prakash, Shaad Boochoon, Mark Brophy, David Acuna, Eric Cameracci, Gavriel State, Omer Shapira, and Stan Birchfield. Structured domain randomization: Bridging the reality gap by context-aware synthetic data. In *International Conference on Robotics and Automation (ICRA)*. IEEE, 2019. 2
- [37] Sanjay Purushotham, Wilka Carvalho, Tanachat Nilanon, and Yan Liu. Variational recurrent adversarial deep domain adaptation. In *ICLR*, 2016. 2
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 1, 2
- [39] Mohammad Mahfujur Rahman, Clinton Fookes, Mahsa Baktashmotlagh, and Sridha Sridharan. Multi-component image translation for deep domain generalization. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2019. 2
- [40] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers generalize to imagenet? In *ICML*, pages 5389–5400. PMLR, 2019. 7
- [41] William J Reed. The pareto, zipf and other power laws. *Economics letters*, 74(1):15–19, 2001. 7
- [42] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. In *NeurIPS*, 2021. 6
- [43] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 6, 7, 13
- [44] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. *NeurIPS*, 2020. 3
- [45] Uri Shaham, Kelly Stanton, Henry Li, Ronen Basri, Boaz Nadler, and Yuval Kluger. Spectralnet: Spectral clustering using deep neural networks. In *ICLR*, 2018. 2, 5
- [46] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *PAMI*, 22(8):888–905, 2000. 2, 5
- [47] Nathan Somavarapu, Chih-Yao Ma, and Zsolt Kira. Frustratingly simple domain generalization via image stylization. *arXiv preprint arXiv:2006.11207*, 2020. 2
- [48] Amos Storkey. When training and test sets are different: characterizing learning transfer. *Dataset shift in machine learning*, 30:3–28, 2009. 3
- [49] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. *arXiv preprint arXiv:1607.01719*, 2016. 2
- [50] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *ICML*. PMLR, 2020. 2
- [51] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, pages 6105–6114. PMLR, 2019. 6, 8
- [52] Meng Tang, Dmitrii Marin, Ismail Ben Ayed, and Yuri Boykov. Kernel cuts: Kernel and spectral clustering meet regularization. *IJCV*, 127:477–511, 2019. 2, 5
- [53] Rohan Taori, Achal Dave, Vaishal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. *Advances in Neural Information Processing Systems*, 2020. 3, 11
- [54] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2017. 2
- [55] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In *NeurIPS*, 2018. 2
- [56] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *ICLR*, 2021. 1, 2, 3, 4, 5, 7, 11
- [57] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Wenjun Zeng, and Tao Qin. Generalizing to unseen domains: A survey on domain generalization. *arXiv preprint arXiv:2103.03097*, 2021. 2
- [58] Garrett Wilson and Diane J Cook. A survey of unsupervised deep domain adaptation. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2020. 2
- [59] Baoyuan Wu, Weidong Chen, Yanbo Fan, Yong Zhang, Jinlong Hou, Jie Liu, and Tong Zhang. Tencent ml-images: A large-scale multi-label image database for visual representation learning. *IEEE Access*, 7:172683–172693, 2019. 6
- [60] Alan L. Yuille and Anand Rangarajan. The concave-convex procedure (CCCP). In *NeurIPS*, 2001. 5
- [61] Kaiyang Zhou, Yongxin Yang, Timothy Hospedales, and Tao Xiang. Deep domain-adversarial image generation for domain generalisation. In *AAAI*, 2020. 2
- [62] Imtiaz Ziko, Jose Dolz, Eric Granger, and Ismail Ben Ayed. Laplacian regularized few-shot learning. In *ICML*, pages 11660–11670. PMLR, 2020. 2, 12