

---

# Towards Context-Agnostic Learning Using Synthetic Data

---

**Charles Jin**  
CSAIL  
MIT

Cambridge, MA 02139  
ccj@csail.mit.edu

**Martin Rinard**  
CSAIL  
MIT

Cambridge, MA 02139  
rinard@csail.mit.edu

## Abstract

We propose a novel setting for learning, where the input domain is the image of a map defined on the product of two sets, one of which completely determines the labels. We derive a new risk bound for this setting that decomposes into a bias and an error term, and exhibits a surprisingly weak dependence on the true labels. Inspired by these results, we present an algorithm aimed at minimizing the bias term by exploiting the ability to sample from each set independently. We apply our setting to visual classification tasks, where our approach enables us to train classifiers on datasets that consist entirely of a single synthetic example of each class. On several standard benchmarks for real-world image classification, we achieve robust performance in the context-agnostic setting, with good generalization to real world domains, whereas training directly on real world data without our techniques yields classifiers that are brittle to perturbations of the background.<sup>1</sup>

## 1 Introduction

We study the problem of learning to classify images of known objects when placed in context, given only a single synthetic example of each object; our empirical evaluation considers the tasks of traffic sign and handwritten character recognition.

Our methods also enable us to explore the extent to which deep neural networks trained using real-world data to perform image classification can over-rely on signals from the backgrounds of images, even when no such information is necessary for classification. Intuitively, even though contextual clues may be required in some settings, given an input that *unambiguously* contains the object of interest in the foreground, we expect a robust classifier to be invariant to the rest of the image.

We present two main contributions. First, we introduce a formal setting for our study, where the input space is decomposed into object and context spaces, and the labels are independent of contexts when conditioned on the objects. We introduce the goal of learning a *context-agnostic* classifier, i.e., a classifier whose predictions are invariant under perturbations of the context. We derive a new risk bound for this setting that is *tight up to a factor of two* but also *independent of the true labels*, which holds so long as the classifier outperforms random guessing.

Second, we present a new technique for automatically generating data to train a deep neural network for image classification. The technique works by sampling independently from the object space, which contains the transformed views of the objects, and the context space, which is designed to include challenging backgrounds that make the resulting images difficult to classify. The hypothesis is that training the network to accurately classify objects against these challenging backgrounds

---

<sup>1</sup>Code is available at [https://github.com/charlesjin/synthetic\\_data](https://github.com/charlesjin/synthetic_data).

should produce a trained network that robustly generalizes to accurately classify objects against the range of backgrounds that it may encounter when deployed in more natural settings.

We empirically validate our methods by training deep neural networks for a variety of real-world image classification tasks using only a single synthetic example of each class, obtaining robust performance in the context-agnostic setting on natural data. Conversely, we find that classifiers trained without our techniques using only natural data achieve negligible accuracy even under relatively benign perturbations that leave a well-defined object in the foreground completely untouched. These results demonstrate the ability of our technique to train accurate and robust classifiers using only small amounts of high quality synthetic data, while also highlighting the need for future work in understanding the performance of deep learning systems in the context-agnostic setting when trained on natural data.

## 2 Related work

Domain shift refers to the problem that occurs when the training set (source domain) and test set (target domain) are drawn from different data distributions. In this setting, a classifier which performs well on the source domain may not generalize well to the target domain. A standard method for addressing this challenge is domain adaptation, which leverages a small amount of data from the target domain to adapt a function that is learned over the source domain [Blitzer et al., 2006]. Approaches can be further categorized as supervised, using small amounts of labelled data from the target domain [Donahue et al., 2014, Motiian et al., 2017b, Saenko et al., 2010, Tzeng et al., 2015]; unsupervised, typically requiring large amounts of unlabelled data from the target domain [Baktashmotlagh et al., 2013, Fernando et al., 2013, Ganin and Lempitsky, 2014, Gopalan et al., 2011]; or semi-supervised, using a mixture of both types of data [Gong et al., 2012, Pan et al., 2010, Yao et al., 2015].

In the context of learning from synthetic data, the domain shift that occurs between synthetic and real world data is known as the reality gap [Jakobi et al., 1995]. State-of-the-art rendering engines, such as those used for video games, can help narrow this gap by generating photorealistic data for training [Dosovitskiy et al., 2017, Johnson-Roberson et al., 2016, Qiu and Yuille, 2016]. Another technique, particularly prevalent in the robotics community, is known as domain randomization, where the synthetic training data is generated with more variability than expected in the actual test environment (e.g., extreme lighting conditions and camera angles), so that the additional robustness also transfers across the reality gap [Tobin et al., 2017, Tremblay et al., 2018]; in particular, Torres et al. [2019] apply domain randomization to traffic sign detection and find that arbitrary natural images suffice for the task. Another body of work exploits generative adversarial networks (GANs) [Goodfellow et al., 2014a] to generate synthetic domains [Hoffman et al., 2017, Liu et al., 2017, Shrivastava et al., 2016, Taigman et al., 2016, Tzeng et al., 2017]. In particular, Shetty et al. [2019] use a GAN trained to perform in-painting and replace extraneous objects in images as a data-augmentation technique to reduce the trained model’s dependence on co-occurring classes.

A different paradigm for the low-data regime is few-shot learning. In contrast to domain adaptation, the goal of few-shot learning is to generalize to new classes given only a few examples, given the ability to train on large amounts of data containing related classes. Early approaches emphasized capturing knowledge in a Bayesian framework [Fe-Fei et al., 2003], which was later formulated as Bayesian program learning [Lake et al., 2015]. Another approach based on metric learning is to find a nonlinear embedding for objects where closeness in the geometry of the embedding generalizes to unseen classes [Koch, 2015, Snell et al., 2017, Sung et al., 2018, Vinyals et al., 2016]. Meta-learning approaches aim to extract higher level concepts which can be applied to learn new classes from a few examples [Finn et al., 2017, Munkhdalai and Yu, 2017, Nichol et al., 2018, Ravi and Larochelle, 2016]. A conceptually-related method that leverages synthetic training data is learning how to generate new data from a few examples of unseen classes; in contrast to our work, however, these methods still require a large number of samples to learn the synthesizer [Schwartz et al., 2018, Zhang et al., 2019]. Finally, some works combine domain adaptation with few-shot learning to learn under domain shift and limited samples (Motiian et al. [2017a]).

The main characteristic that differentiates our work from these approaches is that we are interested in learning classifiers that are *context-agnostic*, i.e., do not rely on background signals. As such, while we find our approach is applicable to many of the same tasks as the aforementioned works, our theoretical setting and objectives differ significantly. From a practical perspective, we demonstrate

our techniques when *the entire training set consists solely of a single synthetic image of each class*, though our techniques can certainly be applied when more data is available; however the reverse does not hold, i.e., existing approaches for domain adaptation or few-shot learning cannot be applied to our setting. Indeed, we consider this work to be complementary in that we are concerned with exploiting the additional structure that is inherent in certain domains, while the goal of domain adaptation and few-shot learning is to achieve good performance on a downstream target distribution given data from a related source distribution.

### 3 Context-agnostic learning

In this section, we introduce the formal setting of context-agnostic learning. We begin with some notation from the standard supervised learning setting. We are given an input space  $\mathcal{X}$ , an output space  $\mathcal{Y}$ , and a hypothesis space  $\mathcal{H}$  of functions mapping  $\mathcal{X}$  to  $\mathcal{Y}$ . A domain  $P_D$  is a probability distribution over  $(\mathcal{X}, \mathcal{Y})$ . Given a target domain  $P_T$  and a loss function  $\ell$ , the goal is to learn a classifier  $h \in \mathcal{H}$  that minimizes the risk, i.e., the expected loss  $R_{P_T}(h) := \mathbb{E}_{P_T}[\ell(h(x), y)]$ . The training procedure is given as input a set of  $n$  training samples  $(x_1, y_1), \dots, (x_n, y_n)$  drawn from a source domain  $P_S$ . The standard approach is empirical risk minimization, which takes the classifier that minimizes  $R_{emp}(h) = \frac{1}{n} \sum_i \ell(h(x_i), y_i)$ . This technique is known to converge to the optimal classifier over  $P_S$  as the number of samples increases; furthermore, if  $P_S$  is sufficiently “close” to  $P_T$  (e.g., if  $P_S = P_T$ , as is the case when there is no domain adaptation), then such a classifier also achieves low risk in the target domain.

#### 3.1 Formal setting

In general, we can frame the goal of classification as learning to extract reliable signals for the label  $y$  from points  $x \in \mathcal{X}$ . This task is often complicated by the presence of noise or other spurious signals. However, for input spaces generated by physical processes, such signals are generally produced by distinct physical entities and can thus be thought of as independent signals that become mixed via an observation process. We aim to capture this additional structure in our setting.

Concretely, we have an object space  $\mathcal{O}$ , a context space  $\mathcal{C}$ , and an observation function  $\gamma$  on  $\mathcal{O} \times \mathcal{C}$ . The input space  $\mathcal{X}$  is defined as the image of  $\gamma : \mathcal{O} \times \mathcal{C} \rightarrow \mathcal{X}$ . We will assume that points in  $\mathcal{O}$  are associated with a unique label in  $\mathcal{Y}$ , and points in  $\mathcal{X}$  inherit labels from  $\mathcal{O}$  via  $\gamma$ ; in particular, we require that  $\gamma$  be injective with respect to labels, i.e.,  $\gamma$  always maps objects with different labels in  $\mathcal{O}$  to distinct points in  $\mathcal{X}$ .

In this work, we will consider the special case when  $\mathcal{X} \subseteq \mathcal{C}$ . Conceptually, the context space is an “ambient space” containing not only valid inputs, but also random noise or irrelevant classes; the input space is a subset of the context space for which there exists a well-defined label. For example, in our experiments we explore such a decomposition for the task of traffic sign recognition, where the object space  $\mathcal{O}$  consists of traffic signs viewed from different angles, the context space  $\mathcal{C}$  is unconstrained pixel space, and the input space  $\mathcal{X}$  is the set of images that contain a traffic sign.

Recall that the standard objective of learning is to find a good classifier for an unknown subdomain  $\mathcal{X}_{P_T} \subseteq \mathcal{X}$ . We consider instead the task of learning a classifier on the entire input space  $\mathcal{X}$ . To sample from  $\mathcal{X}$  we are given oracle access to the observation function and draw (labelled) samples from  $\mathcal{O}$  and  $\mathcal{C}$  independently. Clearly, if this problem is realizable, i.e., there exists  $h^* \in \mathcal{H}$  for which  $R_{\mathcal{X}}(h^*) = 0$ , then we do not even need to know the target domain  $P_T$ , since

$$\mathcal{X}_{P_T} \subseteq \mathcal{X} \implies [R_{\mathcal{X}}(h^*) = 0 \implies R_{P_T}(h^*) = 0] \quad (1)$$

Hence our new goal will be to learn a classifier over  $\mathcal{X}$  which depends only on signals from  $\mathcal{O}$ ; more precisely, we have the following definitions:

**Definition 3.1.** A function  $f$  on  $\mathcal{X}$  is *context-agnostic* if

$$\Pr[f \circ \gamma(o, c) = y] = \Pr[f \circ \gamma(o, c') = y] \quad \forall c, c' \in \mathcal{C}, o \in \mathcal{O}, y \in \text{Im}(f) \quad (2)$$

**Definition 3.2.** The objective of *context-agnostic learning* is to find  $h \in \mathcal{H}$  such that  $h$  achieves the lowest risk of all context-agnostic classifiers.

**Remark.** First, note that if the true label function  $y^*$  is realizable, then the lowest risk classifier is also context-agnostic. Second, we recover the standard supervised setting for the trivial context space  $\mathcal{C} = \emptyset$ . Conversely, our setting remains well-defined even in the trivial object space  $\mathcal{O} = \{y_i\}_i$ , the set of classes; however, this pushes all the complexity to the observation function  $\gamma$ , which may be hard to define or intractable to compute. Finally, we do not preclude the existence of useful signals originating from the context for certain domains. For instance, a great deal of information can often be gleaned from the backgrounds of photos, e.g., stop signs are more often found in cities than on highways. Our theoretical setting avoids this issue by assuming realizability and uniqueness of labels; more practically, we argue that a “good” classifier should nonetheless recognize stop signs on the highway, and our experimental results provide evidence that over-reliance on such background signals leads to brittle classifiers.

### 3.2 A new risk bound for the context-agnostic setting

Our central tool in the context-agnostic setting is a new risk bound that decomposes into separate terms over the context and object spaces. We first develop a formal notion of contextual bias. For clarity we will assume a binary classification task and slightly abuse notation, denoting the classifier as  $h$  instead of  $h \circ \gamma$ , i.e.,  $h : \mathcal{O} \times \mathcal{C} \rightarrow \{-1, 1\}$ . We will denote the true label function as  $y^*$ .

**Definition 3.3.** For an object  $o \in \mathcal{O}$ , the expected classification  $\bar{o}$  and object error  $\hat{o}$  are defined as

$$\bar{o} := \mathbb{E}_{c \sim \mathcal{C}}[h(o, c)] \quad (3)$$

$$\hat{o} := |y^*(o) - \bar{o}| \quad (4)$$

**Definition 3.4.** The context bias  $B(h, c)$  of a classifier  $h$  on the context  $c$  is defined as

$$\text{sgn}(B(h, c)) := \text{sgn}(\mathbb{E}_{o \sim \mathcal{O}}[h(o, c) - \bar{o}]) \quad (5)$$

$$\|B(h, c)\| := \mathbb{E}_{o \sim \mathcal{O}}[\ell(h(o, c), \bar{o})] \quad (6)$$

where  $\ell$  is the hinge loss  $\ell(i, j) := \max(0, 1 - i * j)$ .

Intuitively, the sign of the bias corresponds to the label toward which the classifier is biased by a given context; the magnitude measures the strength of this bias. Clearly, the classifier is context-agnostic exactly when the bias is zero. We are now ready to state our main theoretical result, which gives an upper bound on the risk in terms of the context bias on  $\mathcal{C}$  and object error over  $\mathcal{O}$ .

**Theorem 3.1.** Let  $h$  be a classifier with average bias  $K$ , and denote the true risk as  $R(h)$ . Then the risk can be lower bounded as

$$K/2 \leq R(h) \quad (7)$$

with equality if and only if  $R(h) = K = 0$ . Furthermore, if the object error for all objects is bounded from above by  $\alpha < 1$ , then we may also upper bound the risk as

$$R(h) \leq K/(2 - \alpha) \quad (8)$$

with equality if and only if all object errors equal  $\alpha$ .

The proof is deferred to Appendix A. Since  $\alpha < 1$ , we can further upper bound  $K/(2 - \alpha)$  as  $K$ , which yields the following concise corollary

$$K/2 \leq R(h) < K \quad (9)$$

(under the same assumptions).

**Remark.** Both the learning objective and the standard definition of risk depend crucially on the true labels  $y^*$ ; then given the setting of context-agnostic learning, one might be tempted to focus on minimizing some objective over the object space, which determines the labels exclusively. However from the theorem we see that measuring the context bias of a classifier gives a bound on the risk that is tight up to a constant factor of two *without access to any labels*. In fact, the dependence on the object space is fairly weak in that the labels only enter in via the assumption  $\alpha < 1$ , which is achieved exactly as soon as the classifier outperforms random guessing. Note that the error bound  $\alpha$  and the bias bound  $K$  are not independent; the conclusion is rather that the risk depends more strongly on a good estimate of the bias. In particular,  $\alpha = 0$  if and only if  $K = 0$  and  $\alpha < 1$ ; observe also that when  $\mathcal{C} = \emptyset$ , then  $K = 0$  holds trivially, but in this case we can identify  $\mathcal{O}$  with  $\mathcal{X}$ , so the assumption that  $\alpha < 1$  equivalently yields that the classifier is correct on all inputs.

---

**Algorithm 1: Greedy Bias Correction**

---

**Input:** Object space  $\mathcal{O}$ , context space  $\mathcal{C}$ , observation function  $\gamma$ , number of rounds  $R$ , resample probability  $p$ , classifier update subroutine  $\text{Fit}$ , binary classifier  $h$

**Output:** Trained classifier  $h$

```
// initialize random context and label
 $c \sim \mathcal{C}$ ;
 $y \sim \{-1, 1\}$ ;
for  $r \leftarrow 1$  to  $R$  do
   $o \sim \mathcal{O}(y)$ ; // sample object
   $x \leftarrow \gamma(o, c)$ ; // observe object and context
   $h \leftarrow \text{Fit}(h, x, y)$ ; // perform classifier update
  // update context and label
   $p' \leftarrow \text{Uniform}(0, 1)$ ;
  if  $p' < p$  then
    // resample random context and label
     $c \sim \mathcal{C}$ ;
     $y \sim \{-1, 1\}$ ;
  else
     $c \leftarrow x$ ; // previous image becomes new context
     $y \leftarrow -y$ ; // flip label
  end
end
```

---

## 4 Context-agnostic learning of visual tasks using synthetic data

We now present a pair of algorithms for the setting of context-agnostic learning. The first algorithm is a generic approach to context agnostic learning that attempts to minimize the context bias in line with Theorem 3.1. The second algorithm is a specialization of our framework to the visual domain, where the observation function is given by superposition of objects over contexts. Both algorithms assume unlimited independent samples from  $\mathcal{C}$  and  $\mathcal{O}$ , along with black-box access to the true labelling function  $h^*$  and the observation function  $\gamma$ .

These assumptions allow us to learn  $h^*$  simply by taking the number of samples to infinity. Unfortunately, learning a classifier on the entire input space  $\mathcal{X}$  generally requires many more samples than learning a classifier on a smaller target domain  $\mathcal{X}_{P_T}$ . Thus we should aim to learn  $h^*$  using as few samples as possible. Our main strategy we will be to exploit the *a priori* knowledge that the true label function  $h^*$  is context-agnostic, and thus learn  $h^*$  through the decomposed object and context spaces. To that end, note that while we only need  $\max(|\mathcal{O}|, |\mathcal{C}|)$  samples to observe every object and context once, we need  $|\mathcal{O}||\mathcal{C}|$  samples to observe every object in every context. Hence, the main challenge when the number of samples is low will be avoiding *spurious signals*, i.e., statistical correlations between context and objects (and by extension, labels) which are artifacts of the sampling process and do not generalize outside the training set. This intuition corresponds to the formal notion of contextual bias introduced in the previous section.

### 4.1 Greedy bias correction

The central idea behind Theorem 3.1 is leveraging the fact that labels depend only on objects to factor the risk into separate terms for object error and context bias. This factorization enables us to exploit our ability to sample independently from the object and context spaces. More specifically, we can use samples from  $\mathcal{O}$  to minimize the object error, and samples from  $\mathcal{C}$  to minimize the context bias. Since we only need  $\alpha < 1$  (i.e., any performance that exceeds random guessing), we continue to draw objects randomly; however given an object  $o$ , we aim to observe it with the context for which the classifier has the strongest opposing bias. Intuitively, this allows the classifier to “correct” its bias and unlearn the spurious signals, thereby minimizing the bias and also the risk.

Adopting this approach without modification requires computing the bias of every context in  $\mathcal{C}$ . In most cases, however, even estimating a single bias may be prohibitively expensive. Thus, rather than

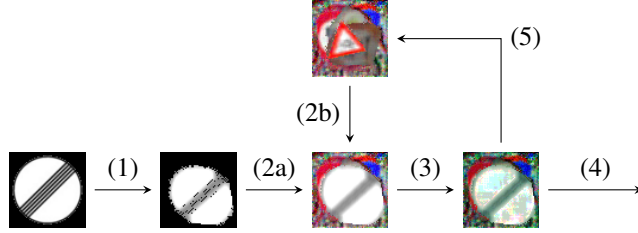


Figure 1: A graphical representation of the generative loop in Algorithm 2 using real training data. (1) Sample from object space. (2) Observe object and context. (3) Perform local refinement. (4) Add to training set. (5) Previous image becomes next context (resample from  $\mathcal{C}$  with probability  $p$ ).

solve for the maximum bias explicitly, we instead propose a heuristic for identifying contexts with large biases. Note that since  $\mathcal{X} \subseteq \mathcal{C}$ , a reasonable assumption is that the classifier learns a strong bias on recent training inputs when taken as contexts. This suggests a simple greedy approach for correcting biases by repurposing recent training inputs as contexts; we call this algorithm Greedy Bias Correction and present a description in Algorithm 1.

## 4.2 Learning visual tasks with synthetic data

We next introduce an instantiation of Greedy Bias Correction for learning visual tasks using synthetic data. We are given a function which takes a label  $y$  and outputs a rendering of the corresponding class in a random pose without any background. The context is the background of the image, on which we place no restrictions. The observation function  $\gamma$  superimposes an object over a background.

**Local refinement via robustness training** We note that our observation function  $\gamma$  is insufficient to capture the true range of possible inputs for a given object; for instance, we do not support occlusions. Because our ultimate goal will be to perform on data taken from a real-world context, we aim to capture this discrepancy using robustness training.<sup>2</sup> In particular, we assume that the image of  $\gamma$  is an  $\epsilon$ -covering of  $\mathcal{X}$ , where a set  $A$  is said to be an  $\epsilon$ -covering of another set  $B$  iff for all points  $b \in B$ , there exists a point  $a \in A$  such that  $\|a - b\| \leq \epsilon$ . Then for a given sample  $x$ , we will instead add the point in the  $\epsilon$ -neighborhood of  $x$  which maximizes the training loss, i.e., for a classifier  $h$  and a sample  $x = \gamma(o, c)$ , we instead aim to identify

$$x' = \operatorname{argmax}_{x' \in N_\epsilon(x)} \ell(h(x'), y) \quad (10)$$

This formulation is often used to train models which are robust against local perturbations. An empirically effective method for finding approximations to  $x'$  is known as Projected Gradient Descent (PGD) [Goodfellow et al., 2014b, Madry et al., 2017]. The algorithm can be summarized as

$$x_0 \leftarrow x + \delta \quad (11)$$

$$x_i \leftarrow \Pi_{N_\epsilon(x)}(x_{i-1} + \eta \cdot \operatorname{sgn}(\nabla_x \ell(h(x_{i-1}), y))), \quad i = 1, \dots, n \quad (12)$$

where  $\delta$  is a small amount of random noise,  $\Pi$  is a projection back onto the  $\epsilon$ -ball,  $\eta$  is the step size, and  $n$  is the number of iterations. As is standard for robustness training, we take the  $\epsilon$ -ball with respect to the  $\ell_\infty$ -norm, defined as  $\|(x_1, \dots, x_n)\|_\infty = \max_i x_i$ . Our choice of  $\epsilon$  will depend on the task at hand, and we also use different  $\epsilon$  for the portions of the image corresponding to the object and context.

Additionally, since we are no longer in a binary context, we sample a random permutation on labels instead of flipping the label deterministically. The full algorithm is presented as Algorithm 2 in Appendix B; Figure 1 provides a visualization of the key generative process, with images taken from a real step of training a deep neural network to perform classification of traffic signs.

From a practical standpoint, this algorithm makes concrete several benefits of our approach. First, rendering object classes, i.e. sampling from  $\mathcal{O}$ , is often relatively easy. In the case of two-dimensional

<sup>2</sup>Robustness training is more commonly referred to as adversarial training in the adversarial robustness community whence we borrow this technique. We use the nonstandard term to avoid confusion with the unrelated (generative) adversarial methods found in the few-shot learning literature.



Figure 2: Images from the training (top) and test (bottom) set for GTSRB (left) and MNIST (right).

Table 1: Performance of Algorithm 2 on various benchmarks, plus ablation studies.

| Approach          | Picto $\rightarrow$ GTSRB | Digit $\rightarrow$ MNIST | Omnifont $\rightarrow$ Omniglot |
|-------------------|---------------------------|---------------------------|---------------------------------|
| baseline          | 72.0                      | 81.9                      | 71.9                            |
| + random-context  | 72.1                      | 88.3                      | 69.8                            |
| + refinement-only | 86.4                      | 89.7                      | 90.8                            |
| + bias-correction | 87.3                      | 89.2                      | 80.5                            |
| <b>+ full</b>     | <b>95.9</b>               | <b>90.2</b>               | <b>92.2</b>                     |

rigid body objects, this can be captured using standard data augmentation such as rotations, flips, and perspective distortions. Indeed, in this setting, our work can be viewed as a form of minimal one-shot learning, where the training data consists solely of a single unobstructed straight-on shot for each object class. Second, since we allow the context space  $\mathcal{C}$  to be unconstrained, there is no requirement to perform realistic rendering of backgrounds, avoiding an additional layer of complexity.

Finally, because our approach is context agnostic, our classifiers are learned without any reference to target domains. In the formal setting, we assumed that the target domain was contained in the image of the observation function; however, synthetic images will always be subject to the reality gap. Our experiments suggest that our approach overcomes this barrier and successfully generalizes to natural images while training on synthetic data only.

## 5 Experiments

We evaluate our approach to learning visual tasks using synthetic data on three benchmarks for image recognition. Our training sets consist of a single synthetic image for each object class with no additional information about the target domain; Figure 2 shows examples of the training and test images from two of the datasets. On all three benchmarks, our models perform comparably with previous state-of-the-art results from related settings using few-shot learning and domain adaptation. Table 1 provides a summary of our results; comprehensive results and comparisons are compiled in Appendix D. Appendix C provides the full experimental setup and training details. Sample images from all datasets referenced below, including examples of rendered training data from the experiments and ablation studies, are shown in Appendix E.

### 5.1 GTSRB

For the traffic sign recognition task, our training set consists of a single, canonical pictogram of each class taken from the visualization software accompanying the target dataset, which we refer to as **Picto**. The target dataset is the German Traffic Sign Recognition Benchmark (**GTSRB**) [Stallkamp et al., 2012], which has 39,209 training and 12,630 test images of 43 classes of German traffic signs taken from the real world. We achieve 95.9% accuracy on the GTSRB test set training only on Picto, against a human baseline of 98.8%. A comprehensive comparison with existing approaches can be found in Appendix D, Table 2.

As a baseline, we also consider approaches using the **SynSign** [Moiseev et al., 2013] dataset of synthetic images designed to provide realistic training data for traffic sign recognition. The dataset comprises 100,000 synthetically generated images of signs from Sweden, Germany, and Belgium in a variety of poses, rendered against domain-appropriate backgrounds (e.g. trees, roads, sky) taken from real-world images. The dataset contains a superset of the GTSRB classes; as a result, Saito et al. [2017] report 79.2% accuracy by training directly on SynSign.

For domain adaptation, all approaches train on the full 100,000 images in SynSign plus part of the GTSRB training set. ATT [Saito et al., 2017] is the only method with better performance than ours, achieving 0.3% higher accuracy; however they use 31,367 unlabelled images from the GTSRB training set (in addition to SynSign). Methods using few-shot learning train on roughly half of the data (22 classes) from the GTSRB training set. The leading few-shot learning approach, VPE [Kim et al., 2019], adds a pictographic dataset similar to Picto, but achieves only 83.79% accuracy. In comparison, our training set consists of only 43 images, none of which are from GTSRB.

## 5.2 Handwritten character recognition

We consider two subtasks for handwritten character recognition. For the first subtask of classifying images of Arabic numerals, our training set, **Digit**, consists of a single example of each digit taken from a standard digital font. The target dataset, **MNIST** [LeCun], consists of 60,000 training and 10,000 test images of handwritten Arabic numerals in grayscale against a blank background. We achieve 90.2% accuracy on MNIST by training only on Digit, compared to human accuracy of 98%.

On MNIST, every approach using domain adaptation uses the full Street View House Numbers (**SVHN**) training set of 73,257 images of house numbers obtained from Google Street View [Netzer et al., 2011], plus varying amounts of data from MNIST. The domain shift problem faces a similar challenge as Digit, namely, handwriting exhibits different characteristics than house numbers fonts. Nevertheless, we note that SVHN contains far more examples of each digit. The only non-baseline approach to exceed our performance is CyCADA [Hoffman et al., 2017], which achieves 0.2% better performance by performing domain adaptation using 60,000 unlabelled images from the MNIST training set (in addition to training on SVHN). In contrast, we use only 10 images, none of which are from MNIST.

For the second subtask, we use the **Omniglot** [Lake et al., 2015] challenge, which consists of 1623 hand-written characters from 50 different alphabets, with 20 samples each. The samples were sourced online from 20 workers on Amazon’s Mechanical Turk, who were asked to copy each character from a single font-based example using digital input (e.g., a mouse). We obtained the original representations (one per character) for our training images, which we call **OmniFont**. We achieve 92.2% on the 20-way Omniglot classification task training only on Omnifont, compared to human accuracy of 95.5%. Tables 3 and 4 in Appendix D compare our results on the handwritten character tasks with approaches using few-shot learning and domain adaptation.

Omniglot is often described as an MNIST-transpose, where the goal is learn handwriting rather than specific symbols, and is widely used as a benchmark for few-shot learning. We reproduce the most common split given in Lake et al. [2015], which uses a predefined set of 30 alphabets, with 19,280 images for training. Test performance is reported as an average over random subsets of  $n = 5$ , 20 unseen classes for the  $n$ -way task (given one labelled example). In comparison, for each test run, we retrain a model using only the corresponding  $n$  images from OmniFont. As expected, our method finds 5-way classification easier than 20-way classification (95.8% vs 92.2%). In both cases, our performance lags behind the state-of-the-art for few-shot learning (>99%), though we emphasize that our experimental setup differs significantly in both the type and amount of training data used.

Finally, several approaches apply few-shot learning from Omniglot to MNIST, with the idea of transferring extracted features from human handwriting. We hypothesize that in comparison to Omniglot, where all the samples come from the same 20 subjects, MNIST may be particularly difficult for transfer one-shot learning, since any two examples will likely exhibit high “variance”; conversely, our approach benefits from using a canonical form which might be closer to the “mean” representation.

**Remark.** Handwritten characters and GTSRB present conceptually opposed challenges for learning: in GTSRB, the objects are rigid two-dimensional objects and backgrounds are complex settings in the natural world; in Omniglot and MNIST, backgrounds are uniform, but classes no longer have a strict specification and individual examples exhibit high variability. Thus, the main challenge of these tasks is learning how to generalize over the object class. Despite the inherent variation, a baseline model trained on Digit with plain data augmentation was able to achieve 81.9% accuracy on MNIST, exceeding many domain adaptation approaches and all the one-shot learning results; Omniglot is more difficult, with an Omnifont plus data augmentation baseline accuracy of 71.9%.



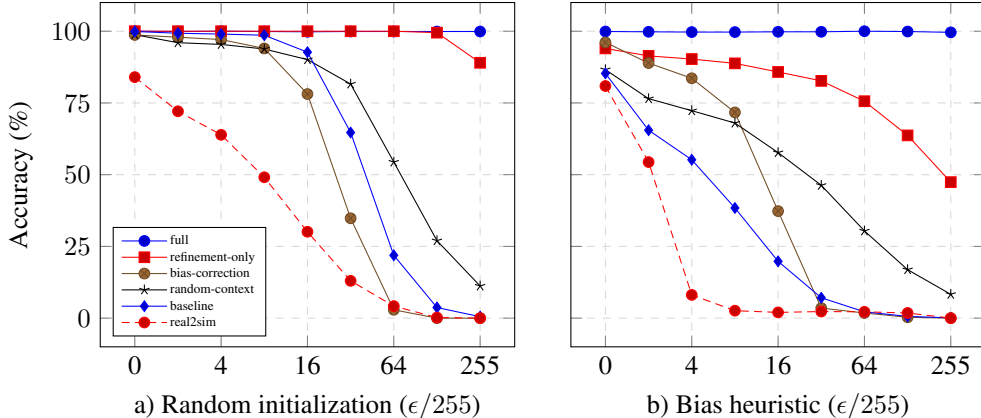


Figure 3: Context-agnostic performance on Picto using a PGD adversary on the background.

### 5.3 Ablation studies

We conduct two sets of ablation studies to better understand our approach to context-agnostic learning. The first study tests the individual components of our algorithm for their contributions to generalization over the real world dataset. All strategies employ the same data augmentation and use the following sampling procedures: **baseline** picks a fresh random background for each training point, and measures the performance of training on our synthetic dataset with plain data augmentation; **random-context** reuses random backgrounds as contexts; **bias-correction** reuses previous training images as contexts; **refinement-only** is the same as random-context with the addition of PGD-based refinement; **full** is the full algorithm as described in Algorithm 2. The results are in Table 1.

In all cases, we observe that both bias correction and local refinement contribute individually and jointly to the performance of our models. For GTSRB, a particularly interesting comparison is training on SynSign, a dataset designed to provide synthetic training data with realistic backgrounds for GTSRB, which yields 79.2% accuracy [Saito et al., 2017]. Though this is an improvement over our baseline of using random backgrounds at 72.0% accuracy, refinement-only and bias-correction achieve higher accuracy at 86.4% and 87.3%, respectively. Both methods leverage the background of training images to combat spurious signals, generating completely unrealistic backgrounds; this suggests that learning context-agnostic features is more effective than using realistic backgrounds.

The second study measures classification performance in a context-agnostic setting on the synthetic Picto dataset. By definition, the performance of a context-agnostic classifier should not degrade under perturbations of the background. We thus run an adaptive attack using a PGD adversary which fixes the foreground pixels, and ranges from fixed to unbounded on the background pixels, effectively searching the context space for a background that causes a misclassification on the given object. We also consider two initialization strategies for the PGD adversary: a standard random initialization, and initializing to the previous image, inspired by our bias heuristic.

We evaluate the same set of ablated strategies as before, plus a classifier trained directly on the GTSRB training set (discussed in the following section). Appendix E.2 contains samples of the generated images, and the results are plotted in Figure 3. Across all experiments, the models have worse (or very close) performance when using our bias heuristic for initialization. We believe this supports our usage of the bias heuristic for context-agnostic learning. Additionally, in the last column of Figure 3b, only our full method maintains passable accuracy, which suggests the gap between models is larger than performance on the GTSRB test set would indicate.

### 5.4 Comparison with training on real data

In this section, we compare a classifier trained using our method on only synthetic data with a model of the same architecture but trained directly on the GTSRB training set and achieving 98% performance on the GTSRB test set, which we refer to as **real2sim**. We first note that Figure 3 indicates the real2sim method seems to suffer from a “synthetic gap” even at  $\epsilon = 0/255$ , which is not entirely unexpected. However, in both settings, the performance of the real2sim model degrades very



Figure 4: Test images from the second ablation study using the Picto dataset. Examples of test images generated using a PGD adversary initialized randomly (top) and with the bias heuristic (bottom) at  $\epsilon = 255/255$ . Note that only backgrounds are perturbed, while foregrounds remain unambiguous.

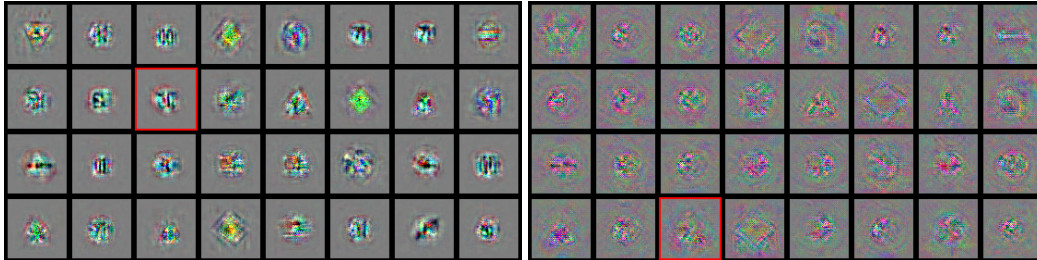


Figure 5: Guided Grad-CAM visualizations for the full (left) and real2sim (right) methods on the GTSRB test set. Misclassified images are marked with red boxes.

quickly as  $\epsilon$  increases: the effect is most pronounced when the bias heuristic is used to initialize the PGD adversary, though in both cases the accuracy eventually drops to 0. We emphasize that all of the experiments leave the foreground objects completely unperturbed (and easily human-identifiable); Figure 4 presents examples of the generated test images.

To better understand the differences between a full classifier trained using our method and the real2sim classifier trained on real data, we used the Guided Grad-CAM method [Selvaraju et al., 2019], which localizes regions of fine-grained features that are important to the classifier’s output. Figure 5 presents visualizations for the real2sim and full classifiers on images from the GTSRB test set. We see that the real2sim classifier (on the right) trained on real images of traffic signs has a more diffuse activation map with no clear interpretation, whereas the classifier trained using our method using purely synthetic data (on the left) is more focused, with more semantically-aligned features. Our results thus suggest that classifiers trained on natural images can become over-reliant on contextual signals, leading to surprisingly brittle behavior even given unambiguous foregrounds.

## 6 Conclusion

We introduce the task of context-agnostic learning, a theoretical setting for learning models whose predictions are independent of background signals. Leveraging the ability to sample objects and contexts independently, we propose an approach to context-agnostic learning by minimizing a formally defined notion of context bias. Our algorithm has a natural interpretation for training classifiers on vision-based tasks using synthetic data, with the distinct advantage that we do not need to model the background. We evaluate our methods on several real-world domains; our results suggest that our approach succeeds in learning context-agnostic classifiers that generalize to natural images using only a single synthetic image of each class, while training with natural images can lead to brittleness in the context-agnostic setting. Our performance is competitive with existing methods for learning when data is limited, while using significantly less data. More broadly, the ability to learn from single synthetic examples of each class also affords fine-grained control over the data used to train our models, allowing us to sidestep issues of data provenance and integrity entirely.

## Acknowledgments and Disclosure of Funding

We gratefully acknowledge support from DARPA Grant HR001120C0015. The views expressed are those of the authors and do not reflect the official policy or position of the Department of Defense or the U.S. Government.

## References

- Antreas Antoniou, Harrison Edwards, and Amos J. Storkey. How to train your MAML. *CoRR*, abs/1810.09502, 2018. URL <http://arxiv.org/abs/1810.09502>.
- Mahsa Baktashmotlagh, Mehrtash T Harandi, Brian C Lovell, and Mathieu Salzmann. Unsupervised domain adaptation by domain invariant projection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 769–776, 2013.
- John Blitzer, Ryan McDonald, and Fernando Pereira. Domain adaptation with structural correspondence learning. In *Proceedings of the 2006 conference on empirical methods in natural language processing*, pages 120–128, 2006.
- Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *Advances in neural information processing systems*, pages 343–351, 2016.
- G. Bradski. The OpenCV Library. *Dr. Dobb’s Journal of Software Tools*, 2000.
- Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655, 2014.
- Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. *arXiv preprint arXiv:1711.03938*, 2017.
- Li Fe-Fei et al. A bayesian approach to unsupervised one-shot learning of object categories. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pages 1134–1141. IEEE, 2003.
- Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *Proceedings of the IEEE international conference on computer vision*, pages 2960–2967, 2013.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1126–1135. JMLR. org, 2017.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. *arXiv preprint arXiv:1409.7495*, 2014.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Muhammad Ghifary, W. Bastiaan Kleijn, Mengjie Zhang, David Balduzzi, and Wen Li. Deep reconstruction-classification networks for unsupervised domain adaptation, 2016.
- Jacob Gildenblat and contributors. Pytorch library for cam methods. <https://github.com/jacobgil/pytorch-grad-cam>, 2021.
- Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2066–2073. IEEE, 2012.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014a.

- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014b.
- Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *2011 international conference on computer vision*, pages 999–1006. IEEE, 2011.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. *arXiv preprint arXiv:1711.03213*, 2017.
- Nick Jakobi, Phil Husbands, and Inman Harvey. Noise and the reality gap: The use of simulation in evolutionary robotics. In *European Conference on Artificial Life*, pages 704–720. Springer, 1995.
- Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? *arXiv preprint arXiv:1610.01983*, 2016.
- Junsik Kim, Seokju Lee, Tae-Hyun Oh, and In So Kweon. Co-domain embedding using deep quadruplet networks for unseen traffic sign recognition. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Junsik Kim, Tae-Hyun Oh, Seokju Lee, Fei Pan, and In So Kweon. Variational prototyping-encoder: One-shot learning with prototypical images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9462–9470, 2019.
- Gregory Koch. Siamese neural networks for one-shot image recognition. 2015.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
- Aoxue Li, Tiange Luo, Tao Xiang, Weiran Huang, and Liwei Wang. Few-shot learning with global class representations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9715–9724, 2019.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Boris Moiseev, Artem Konev, Alexander Chigorin, and Anton Konushin. Evaluation of traffic sign recognition methods trained on synthetically generated data. In *International Conference on Advanced Concepts for Intelligent Vision Systems*, pages 576–583. Springer, 2013.
- Saeid Motiian, Quinn Jones, Seyed Iranmanesh, and Gianfranco Doretto. Few-shot adversarial domain adaptation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6670–6680. Curran Associates, Inc., 2017a. URL <http://papers.nips.cc/paper/7244-few-shot-adversarial-domain-adaptation.pdf>.
- Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5715–5725, 2017b.
- Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2554–2563. JMLR. org, 2017.

- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2010.
- 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. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Weichao Qiu and Alan Yuille. Unrealcv: Connecting computer vision to unreal engine. In *European Conference on Computer Vision*, pages 909–916. Springer, 2016.
- Tiago Ramalho and Marta Garnelo. Adaptive posterior learning: few-shot learning with a surprise-based memory module, 2019.
- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016.
- Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pages 213–226. Springer, 2010.
- Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. Asymmetric tri-training for unsupervised domain adaptation. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, pages 2988–2997. JMLR. org, 2017.
- Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850, 2016.
- Alice Schoenauer-Sebag, Louise Heinrich, Marc Schoenauer, Michele Sebag, Lani F. Wu, and Steve J. Altschuler. Multi-domain adversarial learning, 2019.
- Eli Schwartz, Leonid Karlinsky, Joseph Shtok, Sivan Harary, Mattias Marder, Rogerio Feris, Abhishek Kumar, Raja Giryes, and Alex M. Bronstein. Delta-encoder: an effective sample synthesis method for few-shot object recognition, 2018.
- Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128(2):336–359, Oct 2019. ISSN 1573-1405. doi: 10.1007/s11263-019-01228-7. URL <http://dx.doi.org/10.1007/s11263-019-01228-7>.
- Rakshith Shetty, Bernt Schiele, and Mario Fritz. Not using the car to see the sidewalk—quantifying and controlling the effects of context in classification and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8218–8226, 2019.
- Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Josh Susskind, Wenda Wang, and Russ Webb. Learning from simulated and unsupervised images through adversarial training, 2016.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017.

- J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, (0):–, 2012. ISSN 0893-6080. doi: 10.1016/j.neunet.2012.02.016. URL <http://www.sciencedirect.com/science/article/pii/S0893608012000457>.
- Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1199–1208, 2018.
- Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016.
- 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 *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.
- Lucas Tabelini Torres, Thiago M. Paixão, Rodrigo F. Berriel, Alberto F. De Souza, Claudine Badue, Nicu Sebe, and Thiago Oliveira-Santos. Effortless deep training for traffic sign detection using templates and arbitrary natural images, 2019.
- Jonathan Tremblay, Aayush Prakash, David Acuna, Mark Brophy, Varun Jampani, Cem Anil, Thang To, Eric Cameracci, Shaad Boochoon, and Stan Birchfield. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 969–977, 2018.
- Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4068–4076, 2015.
- Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7167–7176, 2017.
- Stéfan van der Walt, Johannes L. Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D. Warner, Neil Yager, Emmanuelle Gouillart, Tony Yu, and the scikit-image contributors. scikit-image: image processing in Python. *PeerJ*, 2:e453, 6 2014. ISSN 2167-8359. doi: 10.7717/peerj.453. URL <https://doi.org/10.7717/peerj.453>.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016.
- Ting Yao, Yingwei Pan, Chong-Wah Ngo, Houqiang Li, and Tao Mei. Semi-supervised domain adaptation with subspace learning for visual recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2142–2150, 2015.
- Sung Whan Yoon, Jun Seo, and Jaekyun Moon. Tapnet: Neural network augmented with task-adaptive projection for few-shot learning, 2019.
- Jian Zhang, Chenglong Zhao, Bingbing Ni, Minghao Xu, and Xiaokang Yang. Variational few-shot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

## A Proofs

*Proof of Theorem 3.1.* Recall that we denote the true label function as  $y^*$ . First, we show the upper bound. By the assumption that  $\alpha < 1$ , we have that for all  $o \in \mathcal{O}$ , the signs of the expected classification  $\bar{o}$  and correct classification  $y^*(o)$  match, so that  $\alpha \geq \hat{\alpha} = |y^*(o) - \bar{o}| = 1 - |\bar{o}|$ . Then for all  $o \in \mathcal{O}$ ,

$$\ell(\bar{o}, y^*(o)) = 1 - \bar{o}y^*(o) \quad (13)$$

$$= 1 - |\bar{o}| \quad (14)$$

$$= \frac{1 - \bar{o}\bar{o}}{1 + |\bar{o}|} \quad (15)$$

$$= \frac{\ell(\bar{o}, \bar{o})}{1 + |\bar{o}|} \quad (16)$$

$$\leq \frac{\ell(\bar{o}, \bar{o})}{2 - \alpha}, \quad (17)$$

with equality if and only if  $\alpha = \hat{\alpha}$ . Now to bound the risk, we can write,

$$R(h) := \mathbb{E}_{o \sim \mathcal{O}, c \sim \mathcal{C}}[\ell(h(o, c), y^*(o))] \quad (18)$$

$$= \int_{\mathcal{C}} \int_{\mathcal{O}} \ell(h(o, c), y^*(o)) \quad (19)$$

$$= \int_{\mathcal{O}} \ell(\bar{o}, y^*(o)) \quad (20)$$

$$\leq \int_{\mathcal{O}} \frac{\ell(\bar{o}, \bar{o})}{2 - \alpha} \quad (21)$$

$$= \frac{1}{2 - \alpha} \int_{\mathcal{O}} \int_{\mathcal{C}} \ell(h(o, c), \bar{o}) \quad (22)$$

$$= \frac{1}{2 - \alpha} \int_{\mathcal{C}} \|B(h, c)\| \quad (23)$$

$$= \frac{K}{2 - \alpha}, \quad (24)$$

where all double integrals are exchangeable by Fubini's theorem. It also follows that equality for the upper bound holds if and only if  $\alpha = \hat{\alpha}$  for all  $o$ .

For the lower bound, notice that for all  $o \in \mathcal{O}$ , we have  $|\bar{o}| \leq 1$ , so

$$\ell(\bar{o}, y^*(o)) \geq 1 - |\bar{o}| \quad (25)$$

$$= \frac{\ell(\bar{o}, \bar{o})}{1 + |\bar{o}|} \quad (26)$$

$$\geq \frac{\ell(\bar{o}, \bar{o})}{2}, \quad (27)$$

where line 25 is an equality if and only if the signs of the expected and correct classifications match, and line 27 is an equality if and only if  $|\bar{o}| = 1$ ; together, these two conditions imply that the object is correctly classified over all contexts. Then a similar computation as for the upper bound shows that

$$R(h) \geq \int_{\mathcal{O}} \frac{\ell(\bar{o}, \bar{o})}{2} \quad (28)$$

$$= \frac{K}{2}, \quad (29)$$

where now equality holds if and only if all objects are correctly classified over all contexts, i.e.,  $R(h) = 0$ .  $\square$

## B Greedy Bias Correction for visual tasks

---

**Algorithm 2:** Visual Learning Using Context-Agnostic Synthetic Data

---

**Input:** Object space  $\mathcal{O}$ , context space  $\mathcal{C}$ , random permutations  $\Pi$ , observation function  $\gamma$ , number of rounds  $R$ , batch size  $B$ , number of classes  $N$ , resample probability  $p$ , classifier update subroutine  $\text{Fit}$ , projected gradient descent subroutine  $\text{PGD}$ , classifier  $h$

**Output:** Trained classifier  $h$

```
for  $r \leftarrow 1$  to  $R$  do
  // initialize empty training batch and random contexts
   $X \leftarrow \emptyset$ ;
  for  $n \leftarrow 1$  to  $N$  do
    |  $c_n \sim \mathcal{C}$ ;
  end
  for  $b \leftarrow 1$  to  $B$  do
    // sample random permutation
     $\pi \sim \Pi(N)$ ;
    // generate new training data
    for  $n \leftarrow 1$  to  $N$  do
      |  $o \sim \mathcal{O}(n)$ ; // sample object for class
      |  $x \leftarrow \gamma(o, c_{\pi(n)})$ ; // observe object and random (permuted) context
      |  $x' \leftarrow \text{PGD}(h, x)$ ; // perform local refinement
      |  $X \leftarrow X \cup \{(x', y)\}$ ; // add to training set
      |  $c_n \leftarrow x'$ ; // previous sample becomes next context
    end
    // resample contexts
    for  $n \leftarrow 1$  to  $N$  do
      |  $p' \leftarrow \text{Uniform}(0, 1)$ ;
      | if  $p' < p$  then
      | |  $c_n \sim \mathcal{C}$ ;
      | end
    end
  end
  // perform classifier update
   $h \leftarrow \text{Fit}(h, X)$ ;
end
```

---



## C Experimental setup

We used PyTorch 1.5.0 [Paszke et al., 2019], OpenCV 4.2.0 [Bradski, 2000], and scikit-image 0.17.2 [van der Walt et al., 2014] for all experiments. In setting the number of epochs, we did not observe any significant degradation or improvements in performance when training for longer. We use fewer epochs in the case of Omniglot due to computational constraints, as the model is retrained for each test split.

For GTSRB, we use a 5-layer convolutional neural network adapted from the official PyTorch tutorials. To train with Picto, the data augmentation consists of PyTorch transforms `RandomAffine(5, translate=(.15, .15), scale=(0.65, 1.05), shear=5)`, `RandomPerspective(0.5, p=1)`; `ColorJitter(brightness=.8, contrast=.8, saturation=.8, hue=.05)`; OpenCV box blur with a random kernel size between 1 and 6 in both dimensions (independently sampled, so not necessarily square); and a random exposure adjustment by adjusting all pixels by the same random amount between  $-30\%$  and  $50\%$ . For refinement, we used step sizes of  $\alpha = 2/255$  with 8 steps and an epsilon of  $\epsilon = 4/255$  for the foreground only. For the observation function, we superimpose the segmented foreground of the transformed pictographic sign over the context. We train for 300 epochs using the Adam optimizer (learning rate  $1e-4$ , weight decay  $1e-4$ ), with 5 examples of each class per batch and 20 batches per epoch. We report results for the model that achieves the best performance on the training set, checking every 5 epochs.

For MNIST, we use the two-layer convolutional neural network from the official PyTorch examples for MNIST, with Dropout regularization replaced with pre-activation BatchNorm. To train with Digit, the data augmentation consists of PyTorch transforms `RandomAffine(15, translate=(.15, .15), scale=(0.75, 1.05), shear=40)`, `RandomPerspective(0.5, p=1)`; OpenCV box blur with a random kernel size between 1 and 6 in both dimensions (independently sampled, so not necessarily square); then set the foreground to all pixels with value greater than 0.2. For refinement, we used step sizes of  $\alpha = 1.6/255$  with 8 iterations and no projection ( $\epsilon = \infty$ ). For the observation function, we blend the object with the context at a 2:1 ratio; this ensures that inputs have a well-defined ground truth label. We train for 300 epochs using the Adam optimizer (learning rate  $1e-4$ , weight decay  $1e-4$ ), with 5 examples of each class per batch and 20 batches per epoch. We report results for the model that achieves the best performance on the training set, checking every 5 epochs.

For Omniglot, we use the pre-activation variant of ResNet18 [He et al., 2015]. To train with Omnifont, we first preprocess with scikit-learn `skeletonize` and `dilation` to standardize stroke widths. Data augmentation consists of PyTorch transforms `RandomAffine(15, translate=(.15, .15), scale=(0.75, 1.1), shear=20)`, `RandomPerspective(0.25, p=1)`; OpenCV box blur with a random kernel size between 1 and 3 in both dimensions (independently sampled, so not necessarily square); then resize the images to 28 by 28. For refinement, we used step sizes of  $\alpha = 1.6/255$  with 8 iterations and no projection ( $\epsilon = \infty$ ). For the observation function, we blend the object with the context at a 2:1 ratio; this ensures that inputs have a well-defined ground truth label. For the  $n$ -way classification task, we randomly sample  $n$  characters from the Omniglot test set, and use the corresponding characters from the Omnifont dataset as our training set. We then train a fresh model for 150 epochs using the Adam optimizer (learning rate  $1e-4$ , weight decay  $1e-4$ ), and report performance on the all  $20n$  images in the Omniglot test set, averaged over 20 runs (10 runs for the ablation studies).

For the GradCAM visualizations, we use the public grad-cam package [Gildenblat and contributors, 2021] from the Python Package Index (PyPI), with the `target_layer` set to the last LeakyReLU layer in the encoder, and both `aug_smooth` and `eigen_smooth` set to true.

## D Full experimental results

We compare a model trained using our methods with previous state-of-the-art results from related settings using few-shot learning and domain adaptation on GTSRB (Table 2), MNIST (Table 3), and Omniglot (Table 4). When multiple experiments are reported for the same approach, we compare against both the most accurate result as well as the result using the least amount of target data. We distinguish between labelled (**L**) and unlabelled (**UL**) data; experiments for which the training data is not known are marked (?).

Table 2: GTSRB results.

| Approach          | Method   | Training Data |              | Accuracy (%) |
|-------------------|--|---------------|--------------|--------------|
|                   |  | Source        | Target       |              |
| Baselines         | Source Only (Saito et al. [2017])                | SynSign       |              | 79.2         |
|                   | Human (Stallkamp et al. [2012])                  |               |              | 98.8         |
|                   | Target Only (Ganin et al. [2016])                |               | All L        | 99.8         |
| Few-Shot Learning | VPE (Kim et al. [2019]) <sup>§</sup>             | Picto*        | 22 classes L | 83.8         |
|                   | MatchNet (Vinyals et al. [2016]) <sup>§</sup>    |               | 22 classes L | 53.3         |
|                   | QuadNet (Kim et al. [2018]) <sup>§†</sup>        |               | 22 classes L | 45.3         |
| Domain Adaptation | DSN (Bousmalis et al. [2016])                    | SynSign       | 1280 UL      | 93.0         |
|                   | ML (Schoenauer-Sebag et al. [2019]) <sup>§</sup> | SynSign       | 22 classes L | 89.1         |
|                   | MADA (Pei et al. [2018]) <sup>§‡</sup>           | SynSign       | 22 classes L | 84.8         |
|                   | DANN (Ganin et al. [2016])                       | SynSign       | 31367 UL     | 88.7         |
|                   | ATT (Saito et al. [2017])                        | SynSign       | 31367 UL     | <b>96.2</b>  |
| Context Agnostic  | baseline   | Picto         | 0            | 72.0         |
|                   | + random-context                                 | Picto         | 0            | 72.1         |
|                   | + refinement-only                                | Picto         | 0            | 86.4         |
|                   | + bias-correction                                | Picto         | 0            | 87.3         |
|                   | + full   | Picto         | 0            | <b>95.9</b>  |

<sup>§</sup>Test accuracy on remaining 21 unseen classes.

\*Kim et al. [2019] use a pictographic dataset similar to Picto.

<sup>†</sup>Reported in Kim et al. [2019].

<sup>‡</sup>Reported in Schoenauer-Sebag et al. [2019].

Table 3: MNIST results.

| Approach          | Method                            | Training Data |                | Accuracy (%) |
|-------------------|-----------------------------------|---------------|----------------|--------------|
|                   |                                   | Source        | Target         |              |
| Baselines         | Human (Netzer et al. [2011])      |               |                | 98.0         |
|                   | Target Only (Tzeng et al. [2017]) |               | All L          | 99.2         |
| Few-Shot Learning | FADA (Motiian et al. [2017a])     | SVHN          | 1 L / class    | 72.8         |
|                   | + more data                       | SVHN          | 7 L / class    | 87.2         |
|                   | SiamNet (Koch [2015])             | Omniglot      | 1 L / class    | 70.3         |
|                   | MatchNet (Vinyals et al. [2016])  | Omniglot      | 1 L / class    | 72.0         |
|                   | APL (Ramalho and Garnelo [2019])  | Omniglot      | 1 L / class    | 61.0         |
|                   | + more data                       | Omniglot      | ? <sup>‡</sup> | 86.0         |
| Domain Adaptation | DSN (Bousmalis et al. [2016])     | SVHN          | 1000 UL        | 82.7         |
|                   | DRCN (Ghifary et al. [2016])      | SVHN          | ?              | 81.9         |
|                   | DANN (Ganin et al. [2016])        | SVHN          | ?              | 73.9         |
|                   | ATT (Saito et al. [2017])         | SVHN          | ? L + 1000 UL  | 86.0         |
|                   | ADDA (Tzeng et al. [2017])        | SVHN          | 60,000 UL      | 76.0         |
|                   | CyCADA (Hoffman et al. [2017])    | SVHN          | 60,000 UL      | <b>90.4</b>  |
| Context Agnostic  | baseline                          | Digit         | 0              | 81.9         |
|                   | + random-context                  | Digit         | 0              | 88.3         |
|                   | + refinement-only                 | Digit         | 0              | 89.7         |
|                   | + bias-correction                 | Digit         | 0              | 89.2         |
|                   | + full                            | Digit         | 0              | <b>90.2</b>  |

<sup>‡</sup>Cumulative accuracy from adapting over the test set.

Table 4: Omniglot results for one-shot classification.<sup>‡</sup>

| Approach          | Method                           | Training Data | Accuracy (%)      |             |
|-------------------|----------------------------------|---------------|-------------------|-------------|
|                   |                                  |               | 5-way             | 20-way      |
| Baselines         | Human (Lake et al. [2015])       |               |                   | 95.5        |
| Few-Shot Learning | MANN (Santoro et al. [2016])     | Omniglot      | 82.2              |             |
|                   | SiamNet (Koch [2015])            | Omniglot      | 96.7 <sup>§</sup> | 92.0        |
|                   | MatchNet (Vinyals et al. [2016]) | Omniglot      | 98.1              | 93.8        |
|                   | PN (Snell et al. [2017])         | Omniglot      | 98.8              | 96.0        |
|                   | BPL (Lake et al. [2015])         | Omniglot      |                   | 96.7        |
|                   | APL (Ramalho and Garnelo [2019]) | Omniglot      | 97.9              | 97.2        |
|                   | RN (Sung et al. [2018])          | Omniglot      | 99.6              | 97.6        |
|                   | MAML++ (Antoniou et al. [2018])  | Omniglot      | 99.5              | 97.7        |
|                   | TapNet (Yoon et al. [2019])      | Omniglot      |                   | 98.1        |
|                   | GCR (Li et al. [2019])           | Omniglot      | <b>99.7</b>       | <b>99.6</b> |
| Context Agnostic  | baseline                         | Omnifont      |                   | 71.9        |
|                   | + random-context                 | Omnifont      |                   | 69.8        |
|                   | + refinement-only                | Omnifont      |                   | 90.8        |
|                   | + bias-correction                | Omnifont      |                   | 80.5        |
|                   | + full                           | Omnifont      | <b>95.8</b>       | <b>92.2</b> |

<sup>‡</sup>The exact set up of the one-shot classification task often varies between authors. We believe the broad performance numbers are still useful for contextualizing our approach, and refer the reader to the original works for details.

<sup>§</sup>As reported in Vinyals et al. [2016]

## E Training and test set visualizations

### E.1 Datasets



Figure 6: From top to bottom: samples from the GTSRB test set, Picto dataset, and SynSign dataset.

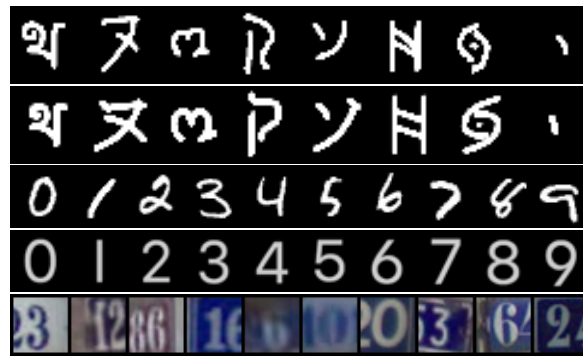


Figure 7: From top to bottom: samples from the Omniglot test set, Omnifot dataset, MNIST test set, Digit dataset, and SVHN training set.

## E.2 Ablation studies



Figure 8: Training images from the first ablation study using Picto dataset. From top to bottom: baseline, random-context, refinement-only, bias-correction, full.

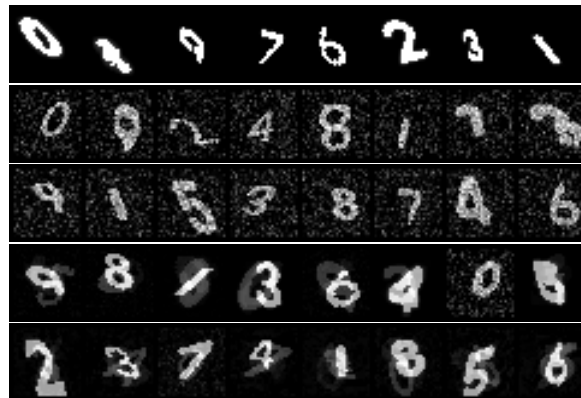


Figure 9: Training images from the first ablation study for the Digit dataset. From top to bottom: baseline, random-context, refinement-only, bias-correction, full.



Figure 10: Grad-CAM visualizations for the PGD background perturbation ablation studies (initialized using the bias heuristic). Full (left) and real2sim (right) methods as perturbations increase over  $\epsilon = 0, 2, 4, 8, 16, 32, 64, 128, 255$  (in order from top to bottom). Regions in yellow are more important to the classifier output. Misclassified images are marked with red boxes.



Figure 11: Guided Grad-CAM visualizations for the full (middle) and real2sim (bottom) methods on the GTSRB test set (original images on top). Regions with color visualize the fine-grained features that contribute to the classifier output. Misclassified images are marked with red boxes.