

AUGCAL: IMPROVING SIM2REAL ADAPTATION BY UNCERTAINTY CALIBRATION ON AUGMENTED SYNTHETIC IMAGES

Prithvijit Chattopadhyay* Bharat Goyal Boglarka Ecsedi Viraj Prabhu Judy Hoffman

Georgia Tech

{prithvijit3,bharatgoyal,becsedi3,virajp,judy}@gatech.edu

ABSTRACT

Synthetic data (SIM) drawn from simulators have emerged as a popular alternative for training models where acquiring annotated real-world images is difficult. However, transferring models trained on synthetic images to real-world applications can be challenging due to appearance disparities. A commonly employed solution to counter this SIM2REAL gap is unsupervised domain adaptation, where models are trained using labeled SIM data and unlabeled REAL data. Mispredictions made by such SIM2REAL adapted models are often associated with miscalibration – stemming from overconfident predictions on real data. In this paper, we introduce AUGCAL, a simple training-time patch for unsupervised adaptation that improves SIM2REAL adapted models by – (1) reducing overall miscalibration, (2) reducing overconfidence in incorrect predictions and (3) improving confidence score reliability by better guiding misclassification detection – all while retaining or improving SIM2REAL performance. Given a base SIM2REAL adaptation algorithm, at training time, AUGCAL involves replacing vanilla SIM images with strongly augmented views (AUG intervention) and additionally optimizing for a training time calibration loss on augmented SIM predictions (CAL intervention). We motivate AUGCAL using a brief analytical justification of how to reduce miscalibration on unlabeled REAL data. Through our experiments, we empirically show the efficacy of AUGCAL across multiple adaptation methods, backbones, tasks and shifts.

1 INTRODUCTION

Most effective models for computer vision tasks (classification, segmentation, *etc.*) need to learn from a large amount of exemplar data (Dosovitskiy et al., 2020; Radford et al., 2021; Kirillov et al., 2023; Pinto et al., 2008) that captures real-world natural variations which may occur at deployment time. However, collecting and annotating such diverse real-world data can be prohibitively expensive – for instance, densely annotating a frame of Cityscapes (Cordts et al., 2016) can take upto ~ 1.5 hours! Machine-labeled synthetic images generated from off-the-shelf simulators can substantially reduce this need for manual annotation and physical data collection (Sankaranarayanan et al., 2018; Ros et al., 2016; Blaga & Nedevschi, 2019; Savva et al., 2019; Deitke et al., 2020; Chattopadhyay et al., 2021). Nonetheless, models trained on SIM data often exhibit subpar performance on REAL data, primarily due to appearance discrepancies, commonly referred to as the SIM2REAL gap. For instance, on GTAV (SIM) \rightarrow Cityscapes (REAL), an HRDA SIM-only model (Hoyer et al., 2022b) achieves an mIoU of only 53.01, compared to ~ 81 mIoU attained by an equivalent model trained exclusively on REAL data.

While there is significant effort in improving the realism of simulators (Savva et al., 2019; Richter et al., 2022), there is an equally large effort seeking to narrow this SIM2REAL performance gap by designing algorithms that facilitate SIM2REAL transfer. These methods encompass both *generalization* (Chattopadhyay* et al., 2023; Huang et al., 2021; Zhao et al., 2022) – aiming to ensure strong out-of-the-box REAL performance of SIM trained models – and *adaptation* (Hoyer et al., 2022b;c; Vu et al., 2019; Rangwani et al., 2022) – attempting to adapt models using labeled SIM data and unlabeled REAL data. Such generalization and adaptation methods have demonstrated notable success in reducing the SIM2REAL performance gap. For instance, PASTA (Chattopadhyay* et al., 2023) (a *generalization* method) improves SIM2REAL performance of a SIM-only model from

*Correspondence to PC

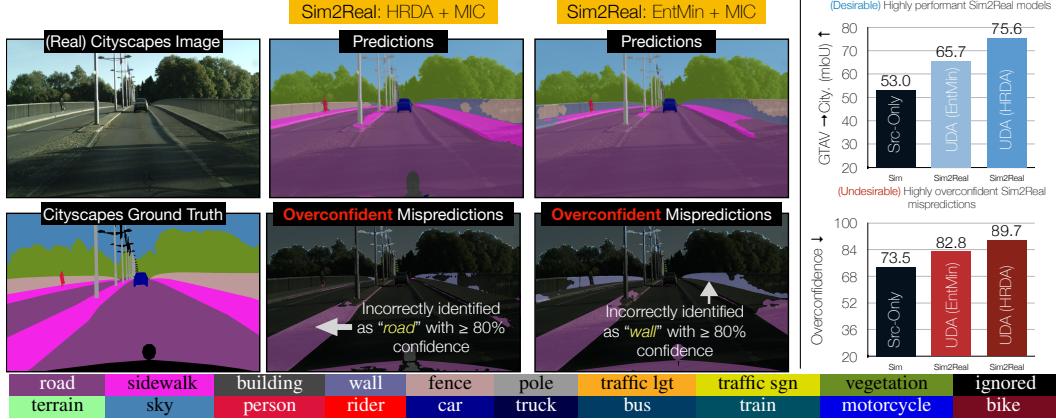


Figure 1: Overconfident SIM2REAL mispredictions. [Left] We show an example of what we mean by overconfident mispredictions. For SIM2REAL adaptation on GTAV → Cityscapes, we choose (DAFormer) HRDA + MIC (Hoyer et al., 2022c) and EntMin + MIC (Vu et al., 2019) (highly performant SIM2REAL methods) and show erroneous predictions on Cityscapes (bottom row). We can see that the model identifies *sidewalk* pixels as *road* (2nd column) and *fence* pixels as *wall* (3rd column) with very high confidence. [Right] We show how pervasive this “overconfidence” phenomena is. While better SIM2REAL adapted models – from (DAFormer) Source-Only (Hoyer et al., 2022b) to (DAFormer) EntMin + MIC (Vu et al., 2019) to (DAFormer) HRDA + MIC (Hoyer et al., 2022c) – exhibit improved transfer performance [Top, Right], they also exhibit increased overconfidence in mispredictions [Bottom, Right], affecting prediction reliability.

53.01 → 57.21 mIoU. Furthermore, HRDA + MIC (Hoyer et al., 2022c) (an *adaptation* approach) pushes performance even higher to 75.56 mIoU.

While SIM2REAL performance may increase both from generalization or adaptation methods, for safety-critical deployment scenarios, task performance is often not the sole factor of interest. It is additionally important to ensure SIM2REAL adapted models make *calibrated* and *reliable* predictions on REAL data. Optimal calibration on real data ensures that the model’s confidence in its predictions aligns with the true likelihood of correctness. Deploying poorly calibrated models can have severe consequences, especially in high-stakes applications (such as autonomous driving), where users can place trust in (potentially) unreliable predictions (Tesla Crash; Michelmore et al., 2018). We find that mistakes made by SIM2REAL adaptation methods are often associated with miscalibration caused by *overconfidence* – highly confident incorrect predictions (see Fig. 1 Left). More interestingly, we find that as adaptation methods improve in terms SIM2REAL performance, the propensity to make overconfident mispredictions also increases (see Fig. 1 Right). Our focus in this paper is to devise training time solutions to mitigate this issue.

Calibrating deep neural networks (for such SIM2REAL adaptation methods) is crucial, as they routinely make overconfident predictions (Guo et al., 2017; Gawlikowski et al., 2021; Minderer et al., 2021). While various techniques address miscalibration on “labeled data splits” for in-distribution scenarios, maintaining calibration in the face of dataset shifts, like SIM2REAL, proves challenging due to lack of labeled examples in the target (REAL) domain. To address this, we propose AUGCAL, a training-time patch to ensure existing SIM2REAL adaptation methods make *accurate*, *calibrated* and *reliable* predictions on real data. When applied to a SIM2REAL adaptation framework, AUGCAL aims to satisfy three key criteria: (1) retain performance of the base SIM2REAL method, (2) reduce miscalibration and overconfidence and (3) ensure calibrated confidence scores translate to improved reliability. Additionally, to ensure broad applicability, AUGCAL aims to do so by making two minimally invasive changes to a SIM2REAL adaptation training pipeline. First, by AUGmenting (Cubuk et al., 2020; Chattopadhyay* et al., 2023) input SIM images during training using an AUG transform that reduces distributional distance between SIM and REAL images. Second, by additionally optimizing for a CALibration loss (Hebbalaguppe et al., 2022; Liang et al., 2020; Liu et al., 2022) at training time on AUGmented SIM predictions. We devise AUGCAL based on an analytical rationale (see Sec. 3.2.1 and 3.2.2) illustrating how it helps reduce an upper bound on desired target (REAL) calibration error. Through our experiments on GTAV → Cityscapes and VisDA SIM2REAL, we demonstrate how AUGCAL helps reduce miscalibration on REAL data. To summarize, we make the following contributions:

- We propose AUGCAL, a training time patch, compatible with existing SIM2REAL adaptation methods that ensures SIM2REAL adapted models make *accurate* (measured via adaptation performance), *calibrated* (measured via calibration error) and *reliable* (measured via confidence guided misclassification detection) predictions.

- We conduct SIM2REAL adaptation experiments for object recognition (VisDA (Peng et al., 2017)) and semantic segmentation (GTAV (Sankaranarayanan et al., 2018)→Cityscapes (Cordts et al., 2016)) with three representative UDA methods (pseudo-label based self-training, entropy minimization and domain adversarial training) and show that applying AUGCAL– (1) improves or preserves adaptation performance, (2) reduces miscalibration and overconfidence and (3) improves the reliability of confidence scores.
- We show how AUGCAL improvements are effective across multiple backbones, AUG and CAL options and highlight choices that are more consistently effective across experimental settings.

2 RELATED WORK

Unsupervised Domain Adaptation (UDA). We focus on UDA algorithms to address *covariate shifts* in the SIM2REAL context (Chattopadhyay* et al., 2023; Choi et al., 2021; Zhao et al., 2022; Huang et al., 2021; Rangwani et al., 2022; Hoyer et al., 2022c; Sankaranarayanan et al., 2018; Ros et al., 2016). This involves adapting a model to an unseen target (REAL) domain using labeled samples from a source (SIM) domain and unlabeled samples from the target domain. Here, the source and target datasets share the same label space and labeling functions, but differences exist in the distribution of inputs (Farahani et al., 2021; Zhang et al., 2019). SIM2REAL UDA methods (Ganin & Lempitsky, 2014; Hoffman et al., 2018; Saenko et al., 2010; Tzeng et al., 2014) range from *feature distribution matching* (Ganin & Lempitsky, 2014; Long et al., 2018; Saito et al., 2018; Tzeng et al., 2017; Zhang et al., 2019), explicitly addressing *domain discrepancy* (Kang et al., 2019; Long et al., 2015; Tzeng et al., 2014; Rangwani et al., 2022), *entropy minimization* (Vu et al., 2019) or *pseudo-label guided self-training* (Hoyer et al., 2022b;a;c). We observe that existing SIM2REAL UDA methods usually improve performance at the expense of increasingly overconfident mispredictions on (REAL) target data (Wang et al., 2020b). Our proposed method, AUGCAL, is designed to retain SIM2REAL adaptation performance while reducing miscalibration on real data for existing methods. We conduct experiments on three representative UDA methods – Entropy Minimization (Vu et al., 2019), Self-training (Hoyer et al., 2022b) and Domain Adversarial Training (Rangwani et al., 2022).

Confidence Calibration for Deep Networks. For discriminative models, confidence calibration indicates the degree to which confidence scores associated with predictions align with the true likelihood of correctness (usually measured via ECE (Naeini et al., 2015)). Deep networks tend to be very poor at providing calibrated confidence estimates (are overconfident) for their predictions (Guo et al., 2017; Gawlikowski et al., 2021; Minderer et al., 2021), which in turn leads to less reliable predictions for decision-making in safety-critical settings. Recent work (Guo et al., 2017) has also shown that calibration worsens for larger models and can decrease with increasing performance. Several works (Guo et al., 2017; Lakshminarayanan et al., 2017; Malinin & Gales, 2018) have explored this problem for modern architectures, and several solutions have also been proposed –including temperature scaling (prediction logits being divided by a scalar learned on a held-out set (Platt et al., 1999; Kull et al., 2017; Bohdal et al., 2021; Islam et al., 2021)) and trainable calibration objectives (training time loss functions that factor in calibration (Liang et al., 2020; Karandikar et al., 2021)). Improving network calibration is even more challenging in out-of-distribution settings due to the simultaneous lack of ground truth labels and overconfidence on unseen samples (Wang et al., 2020b). Specifically, instead of methods that rely on temperature-scaling (Wang et al., 2020a; 2022) or maybe require an additional calibration split, AUGCAL explores the use of training time calibration objectives (Munir et al., 2022) to reduce micalibration for SIM2REAL shifts.

3 METHOD

3.1 BACKGROUND

Notations. Let x denote input images and y denote corresponding labels (from the label space $\mathcal{Y} = \{1, 2, \dots, K\}$) drawn from a joint distribution $P(x, y)$. We focus on the classification case, where the goal is to learn a discriminative model \mathcal{M}_θ (with parameters θ) that maps input images to the desired K output labels, $\mathcal{M}_\theta : \mathcal{X} \rightarrow \mathcal{Y}$, using a softmax layer on top. The predictive probabilities for the given input can be expressed as $p_\theta(y|x) = \text{softmax}(\mathcal{M}_\theta(x))$. We use $\hat{y} = \arg \max_{y \in \mathcal{Y}} p_\theta(y|x)$ to denote the predicted label for x and c to denote the confidence in prediction.

Unsupervised SIM2REAL Adaptation. In unsupervised domain adaptation (UDA) for SIM2REAL settings, we assume access to a labeled (SIM) source dataset $D_S = \{(x_i^S, y_i^S)\}_{i=1}^{|S|}$ and an unlabeled (REAL) target dataset $D_T = \{x_i^T\}_{i=1}^{|T|}$. We assume D_S and D_T splits are drawn from source and target distributions $P^S(x, y)$ and $P^T(x, y)$ respectively. At training, we have access to $D = D_S \cup D_T$. We operate in the setting where source and target share the same label space, and discrepancies exist

only in input images. The model \mathcal{M}_θ is trained on labeled source images using cross entropy,

$$\sum_{i=1}^{|S|} \mathcal{L}_{CE}(x_i^S, y_i^S; \theta) = -\sum_{i=1}^{|S|} y_i^S \log p_\theta(\hat{y}_i^S | x_i^S) \text{ where } \hat{y}_i^S = \arg \max_{y \in \mathcal{Y}} p_\theta(y_i^S | x_i^S) \quad (1)$$

UDA methods additionally optimize for an adaptation objective on labeled source and unlabeled target data (\mathcal{L}_{UDA}). The overall learning objective can be expressed as,

$$\min_{\theta} \underbrace{\sum_{i=1}^{|S|} \mathcal{L}_{CE}(x_i^S, y_i^S; \theta)}_{\text{Source Loss}} + \underbrace{\sum_{i=1}^{|T|} \sum_{j=1}^{|S|} \lambda_{UDA} \mathcal{L}_{UDA}(x_j^T, x_i^S, y_j^S; \theta)}_{\text{Source Target Adaptation Loss}} \quad (2)$$

Different adaptation methods usually differ in terms of specific instantiations of this objective. While AUGCAL is applicable to any SIM2REAL adaptation method in principle, we conduct experiments with three popular methods – Entropy Minimization (Vu et al., 2019), Pseudo-Label driven Self-training (Hoyer et al., 2022b) (for semantic segmentation) and Domain Adversarial Training (Rangwani et al., 2022) (for object recognition). We provide more details on these methods in Sec. A.7 of appendix.

Uncertainty Calibration. For a perfectly calibrated classifier, the confidence in predictions should match the empirical frequency of correctness. Empirically, calibration can be measured using Expected Calibration Error (ECE) (Naeini et al., 2015). To measure ECE on a test set $D = \{(x_i, y_i)\}_{i=1}^{|D|}$, we first partition the test data into B bins, $D_b = \{(x, y) \mid r_{b-1} \leq c < r_b\}$, using the confidence values c such that $b \in \{1, \dots, B\}$ and $0 = r_0 \leq r_1 \leq r_2 \leq \dots \leq r_B = 1$. Then, ECE measures the absolute differences between accuracy and confidence across instances in every bin,

$$\text{ECE} = \sum_{j=1}^B \frac{B}{|D|} \left| \frac{1}{B} \sum_{i \in D_j} \mathbf{1}_{(y_i = \hat{y}_i)} - \sum_{i \in D_j} c_i \right| \quad (3)$$

We are interested in models that exhibit high-performance and low calibration error (ECE). Note that Eqn 3 alone does not indicate if a model is overconfident. We define overconfidence (OC) as the expected confidence on mispredictions. Prior work on improving calibration in out-of-distribution (OOD) settings (Wang et al., 2022) and domain adaptation scenarios (Wang et al., 2020b) typically rely on techniques like temperature scaling. These methods often necessitate additional steps, such as employing a separate calibration split or domains (Gong et al., 2021) or training extra models (e.g., logistic discriminators for source and target features (Wang et al., 2020b)). In contrast, we consider using training time calibration objectives (Liang et al.; Hebbalaguppe et al., 2022; Liu et al., 2022) that can be optimized in addition to task-specific objectives for improved calibration.

3.2 AUGCAL

3.2.1 REDUCING MISCALIBRATION ON (REAL) TARGET

Recall that $P^S(x, y)$ and $P^T(x, y)$ denote the source (SIM) and target (REAL) data distributions. We assume $P(x, y)$ factorizes as $P(x, y) = P(x)P(y|x)$. We assume covariate shift conditions between P^S and P^T , i.e., $P^T(x) \neq P^S(x)$ while $P^T(y|x) = P^S(y|x)$ – discrepancies across distributions exist only in input images. When training a model, we can only draw “labeled samples” (x, y) from $P^S(x, y)$. We do not have access to labels from $P^T(x, y)$. Our goal is to reduce miscalibration on (unlabeled) target data using training time calibration losses. Let $\mathcal{L}_{CAL}(x, y)$ denote such a calibration loss we can minimize (on labeled data). Using importance sampling (Cortes et al., 2010), we can get an estimate of the desired calibration loss on target data as,

$$\begin{aligned} \mathbb{E}_{x, y \sim P^T(x, y)} [\mathcal{L}_{CAL}(x, y)] &= \int_x \int_y \mathcal{L}_{CAL}(x, y) P^T(x, y) dx dy \\ &= \int_x \int_y \mathcal{L}_{CAL}(x, y) \frac{P^T(x) P^T(y|x)}{P^S(x) P^S(y|x)} P^S(x, y) dx dy \\ &= \mathbb{E}_{x, y \sim P^S(x, y)} \left[\underbrace{w_S(x)}_{\text{Importance Weight}} \underbrace{\mathcal{L}_{CAL}(x, y)}_{\text{Source Loss}} \right] \end{aligned} \quad (4)$$

where $w_S(x) = \frac{P^T(x)}{P^S(x)}$ denotes the importance weight. Assuming $\mathcal{L}_{\text{CAL}}(x, y) \geq 0$ ¹, we can obtain an upper bound on step 4 (Pampari & Ermon, 2020; Wang et al., 2020b) as

$$\underbrace{\mathbb{E}_{x,y \sim P^T(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)]}_{\leq} = \mathbb{E}_{x,y \sim P^S(x,y)} [w_S(x)\mathcal{L}_{\text{CAL}}(x, y)] \quad (5)$$

$$\leq \sqrt{\mathbb{E}_{P^S(x)} [w_S(x)^2] \mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)^2]} \quad (6)$$

$$\leq \underbrace{\frac{1}{2} \left(\mathbb{E}_{P^S(x)} [w_S(x)^2] \right)}_{\text{Shift Dependent}} + \underbrace{\mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)^2]}_{\text{Source Dependent}} \quad (7)$$

where steps 6 and 7 use the Cauchy-Schwarz and AM-GM inequalities respectively. For a given model, the second RHS term in inequality 7 is computed purely on labeled samples from the source distribution and can therefore be optimized to convergence over the course of training. The gap in \mathcal{L}_{CAL} across source and target is dominated by the importance weight (first term). Following (Cortes et al., 2010), the first term can also be expressed as,

$$\mathbb{E}_{P^S} [w_S(x)^2] = d_2(P^T(x) || P^S(x)) \quad (8)$$

where $d_\alpha(P || Q) = \left[\sum_x \frac{P^\alpha(x)}{Q^{\alpha-1}(x)} \right]^{\frac{1}{\alpha-1}}$ with $\alpha > 0$ is the exponential in base 2 of the Renyi divergence (Rényi, 1960) between distributions P and Q . The calibration error gap between source and target distributions is therefore, dominated by the divergence between source and target distributions. Consequently, inequality 7 can be expressed as,

$$\underbrace{\mathbb{E}_{x,y \sim P^T(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)]}_{\text{Target Calibration Loss}} \leq \underbrace{\frac{1}{2} d_2(P^T(x) || P^S(x))}_{\text{Source and Target Divergence}} + \underbrace{\frac{1}{2} \mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)^2]}_{\text{Source Calibration Loss}} = \underbrace{U(S, T)}_{\text{Upper Bound}} \quad (9)$$

where $U(S, T)$ denotes the upper bound on target calibration loss. Therefore, to effectively reduce miscalibration on target data, one needs to reduce the upper bound, $U(S, T)$, which translates to (1) reducing miscalibration on source data (second red term in 9) and (2) reducing the distributional distance between input distributions across source and target (first blue term in 9).

3.2.2 WHY AUGCAL?

Based on the previous discussion, to improve calibration on target data, one can always invoke a training time calibration intervention (CAL) on labeled source data to reduce $\mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)]$. In practice, after training, we can safely assume that $\mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)] = \epsilon \rightarrow 0$ (for some very small ϵ). We note that while this is useful and necessary, it is not sufficient. This is precisely where we make our contribution. To reduce both (red and blue) terms in 9, we introduce AUGCAL. To do this, in addition to a training time calibration loss, \mathcal{L}_{CAL} , AUGCAL assumes that access to an additional AUG transformation that satisfies the following properties:

1. $d_2(P^T(x) || P^S(\text{AUG}(x))) \leq d_2(P^T(x) || P^S(x))$
2. After training, $\mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(x, y)] \approx \mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(\text{AUG}(x), y)] = \epsilon \rightarrow 0$

Property 1 states that the chosen AUG transformation brings transformed source data closer to target (or reduces SIM2REAL distributional distance). Property 2 states that over the course of training, irrespective of the data $\mathcal{L}_{\text{CAL}}(x, y)$ is optimized on (AUG transformed or clean source), $\mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(\cdot, \cdot)]$ can achieve a sufficiently small value close to 0. Given an AUG transformation that satisfies the above stated properties, we can claim,

$$U^{\text{AUG}}(S, T) \leq U(S, T) \quad (10)$$

where

$$U^{\text{AUG}}(S, T) = \frac{1}{2} d_2(P^T(x) || P^S(\text{AUG}(x))) + \frac{1}{2} \mathbb{E}_{P^S(x,y)} [\mathcal{L}_{\text{CAL}}(\text{AUG}(x), y)^2] \quad (11)$$

¹We make the reasonable assumption that the calibration loss function is always non-negative.

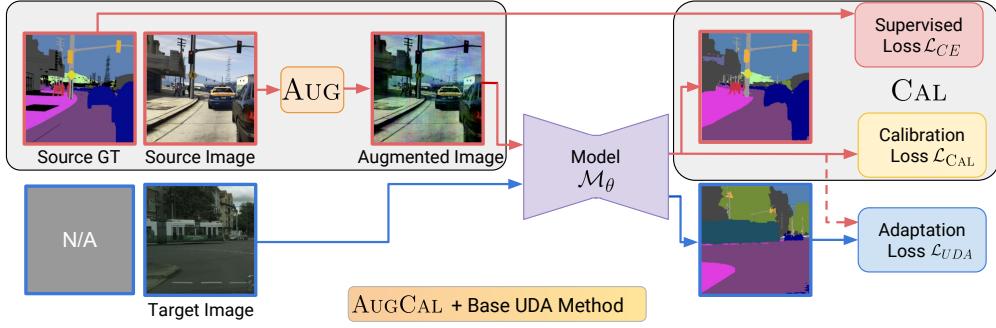


Figure 2: **AUGCAL pipeline.** AUGCAL consists of two key interventions on an existing SIM2REAL adaptation method. First source SIM images are augmented via an AUG transform. Supervised losses for SIM images are computed on the augmented image predictions. Additionally, AUGCAL optimizes for a calibration loss on AUGMENTED SIM predictions.

That is, an appropriate AUG transform, when coupled with \mathcal{L}_{CAL} , helps reduce a tighter upper bound on the target calibration error than CAL. We call this intervention – coupling AUG and CAL – AUGCAL. Naturally, the effectiveness of AUGCAL is directly dependent on the choice of AUG that satisfies the aforementioned properties.

While most augmentations can satisfy property 2, to check if an augmentation is valid according to property 1, we compute RBF Kernel based MMD distances for (SIM, REAL) and (AUGmented SIM, REAL) feature pairs using a trained model.² In Table. 1, we show how PASTA ([Chattopadhyay* et al., 2023](#)) and RandAugment ([Cubuk et al., 2020](#)), two augmentations effective for SIM2REAL transfer, satisfy these considerations on multiple shifts (read Table. 1 left to right). PASTA and RandAug are additionally (1) inexpensive when combined with SIM2REAL UDA methods, and (2) generally beneficial for SIM2REAL shifts (PASTA via SIM2REAL specific design and RandAug via chained photometric operations).

3.2.3 AUGCAL INSTANTIATION

Given a SIM2REAL adaptation method, AUGCAL additionally optimizes for improved calibration on augmented SIM source images. Since AUGCAL is applicable to any existing SIM2REAL adaptation method, we abstract away the adaptation component associated with the pipeline and denote as \mathcal{L}_{UDA} (see Eqn. 15). The steps involved in AUGCAL are illustrated in Fig. 2. Given a mini-batch, we first generate augmented views, $AUG(x^S)$, for SIM images x^S . Then, during training, we optimize \mathcal{L}_{CE} on those augmented SIM views and \mathcal{L}_{UDA} for adaptation. To improve calibration under augmentations, we optimize an additional \mathcal{L}_{CAL} loss on augmented SIM images. The overall AUGCAL optimization problem can be expressed as,

$$\min_{\theta} \underbrace{\sum_{i=1}^{|S|} \mathcal{L}_{CE}(AUG(x_i^S), y_i^S; \theta)}_{\text{Source Task Loss}} + \underbrace{\sum_{i=1}^{|T|} \sum_{j=1}^{|S|} \lambda_{UDA} \mathcal{L}_{UDA}(x_i^T, AUG(x_j^S), y_j^S; \theta)}_{\text{Source Target Adaptation Loss}} + \underbrace{\sum_{i=1}^{|S|} \lambda_{CAL} \mathcal{L}_{CAL}(AUG(x_i^S), y_i^S; \theta)}_{\text{Source Calibration Loss}} \quad (12)$$

where λ_{UDA} and λ_{CAL} denote the respective loss coefficients and the changes to a vanilla SIM2REAL adaptation framework are denoted in teal.

Choice of “AUG”. Strongly augmenting SIM images during training has proven useful for SIM2REAL transfer ([Chattopadhyay* et al., 2023](#); [Huang et al., 2021](#); [Zhao et al., 2022](#); [Cubuk et al., 2020](#)). For AUGCAL, we are interested in augmentations that satisfy the properties outlined in Sec. 3.2.2. As stated earlier, we find that RandAugment ([Cubuk et al., 2020](#)) and PASTA ([Chattopadhyay* et al., 2023](#)) empirically satisfy this criteria. We use both of them for our experiments and provide details on the operations for these AUG transforms in Sec. A.2.2 of appendix.

Choice of “CAL” (\mathcal{L}_{CAL}). AUGCAL relies on using training time calibration losses to reduce miscalibration. Prior work in uncertainty calibration has considered several auxiliary objectives

²Under bounded importance weight assumptions, MMD can be interpreted as an upper bound on KL divergence ([Wang & Tay, 2022](#)).

to calibrate a model being trained to reduce negative log-likelihood (NLL) (Hebbalaguppe et al., 2022; Kumar et al., 2018; Karandikar et al., 2021; Liang et al., 2020). For “CAL” in AUGCAL, while we consider multiple calibration losses – DCA (Liang et al.), MbLS (Liu et al., 2022) and MDCA (Hebbalaguppe et al., 2022) – and find that DCA, simple “difference between confidence and accuracy (DCA)” loss proposed in (Liang et al., 2020) is more consistently effective across experimental settings. DCA can be expressed as,

$$\mathcal{L}_{\text{CAL}} = \frac{1}{|S|} \left| \sum_{i=1}^{|S|} \mathbf{1}_{(y_i^S = \hat{y}_i^S)} - \sum_{i=1}^{|S|} p_\theta(\hat{y}_i^S | x_i^S) \right| \quad (13)$$

where $\mathbf{1}_{(y_i^S = \hat{y}_i^S)}$ and $p_\theta(\hat{y}_i^S | x_i^S)$ denote the correctness and confidence scores associated with predictions. The DCA loss forces the mean predicted confidence over training samples to match accuracy. In the following sections, we empirically validate AUGCAL across adaptation methods.

4 EXPERIMENTAL DETAILS

We conduct SIM2REAL adaptation experiments across two tasks – Semantic Segmentation (SemSeg) and Object Recognition (ObjRec). For our experiments, we train models using labeled SIM images and unlabeled REAL images. We test trained models on REAL images.

SIM2REAL Shifts. For SemSeg, we conduct experiments on the GTAV→Cityscapes shift. GTAV (Sankaranarayanan et al., 2018) consists of $\sim 25k$ densely annotated SIM ground-view images and Cityscapes (Cordts et al., 2016) consists of $\sim 5k$ REAL ground view images. We report all metrics on the Cityscapes validation split. For ObjRec, we conduct experiments on the VisDA SIM2REAL benchmark. VisDA (Peng et al., 2017) consists of $\sim 152k$ SIM images and $\sim 55k$ REAL images across 12 classes. We report all metrics on the validation split of (REAL) target images.

Models. We check AUGCAL compatibility with both CNN and Transformer based architectures. For SemSeg, we consider DeepLabv2 (Chen et al., 2017) (with a ResNet-101 (He et al., 2016) backbone) and DAFormer (Hoyer et al., 2022a) (with an MiT-B5 (Xie et al., 2021) backbone) architectures. For ObjRec, we consider ResNet-101 and ViT-B/16 (Dosovitskiy et al., 2020) backbones with bottleneck layers as classifiers. We start with backbones pre-trained on ImageNet (Deng et al., 2009).

Adaptation Methods. We consider three representative SIM2REAL adaptation methods for our experiments. For SemSeg, we consider entropy minimization (EntMin) (Vu et al., 2019) and high-resolution domain adaptive semantic segmentation (HRDA) (Hoyer et al., 2022b). For ObjRec, we consider smooth domain adversarial training (SDAT) (Rangwani et al., 2022). For both tasks, we further improve performance with masked image consistency (MIC) (Hoyer et al., 2022c) on target images during training. We use MIC (Hoyer et al., 2022c)’s implementations of the adaptation algorithms and provide more training details in Sec. A.3 of appendix.

Calibration Metrics. We use ECE to report overall confidence calibration on REAL images. Since we are interested in reducing overconfident mispredictions, we also report calibration error on incorrect samples (IC-ECE) (Wang et al., 2022) and mean overconfidence for mispredictions (OC).

Reliability Metrics. While reducing overconfidence and improving calibration on real data is desirable, this is a proxy for the true goal of improving model reliability. To assess reliability, following prior work (de Jorge et al., 2023; Malinin et al., 2019), we measure whether calibrated confidence scores can better guide misclassification detection. To measure this, we use Prediction Rejection Ratio (PRR) (Malinin et al., 2019), which if high (positive and close to 100) indicates that confidence scores can be used as reliable indicators of performance (details in Sec. A.4 of appendix).

Unless specified otherwise, we use PASTA as the choice of AUG and DCA (Liang et al.) as the choice of CAL in AUGCAL. We use $\lambda_{\text{CAL}} = 1$ for DCA.

5 FINDINGS

5.1 IMPROVING SIM2REAL ADAPTATION

Recall that when applied to a SIM2REAL adaptation method, we expect AUGCAL to – (1) retain SIM2REAL transfer performance, (2) reduce miscalibration and overconfidence and (3) ensure calibrated confidence scores translate to improved model reliability. We first verify these criteria.

▷ **AUGCAL improves or retains SIM2REAL adaptation performance.** Since AUGCAL intervenes on an existing SIM2REAL adaptation algorithm, we first verify that encouraging better calibration does not adversely impact SIM2REAL adaption performance. We find that performance is either retained or improved (*e.g.*, for EntMin + MIC in Table. 2 (a)) as miscalibration is reduced (Tables. 2(a) and (b), Perf. columns).

Table 2: AUGCAL ensures SIM2REAL adapted models make accurate, calibrated and reliable predictions. We find that applying AUGCAL to multiple SIM2REAL adaptation methods across tasks leads to better calibration (ECE, IC-ECE), reduced overconfidence (OC) and improved reliability (PRR) – all while retaining or improving transfer performance. Highlighted rows are AUGCAL variants of the base methods. For AUGCAL, we use PASTA as AUG and DCA as CAL. \pm indicates standard error.

Method	Perf. (\uparrow)		Calibration Error (\downarrow)		Reliability (\uparrow)	
	mIoU	ECE	IC-ECE	OC	PRR	
1 EntMin + MIC	65.71	5.34 \pm 0.35	77.73 \pm 0.26	82.83 \pm 0.55	45.93 \pm 0.54	
2 + AUGCAL	70.31	3.43 \pm 0.29	72.97 \pm 0.26	82.80 \pm 0.57	62.66 \pm 0.55	
3 HRDA + MIC	75.56	2.86 \pm 0.10	81.92 \pm 0.14	89.72 \pm 0.48	68.91 \pm 0.46	
4 + AUGCAL	75.90	2.45 \pm 0.09	79.09 \pm 0.16	88.26 \pm 0.49	70.35 \pm 0.51	

Method	Perf. (\uparrow)		Calibration Error (\downarrow)		Reliability (\uparrow)	
	mIoU	ECE	IC-ECE	OC	PRR	
1 SDAT + MIC	92.53 \pm 0.28	7.67 \pm 0.49	91.45 \pm 0.63	89.13 \pm 1.29	63.78 \pm 2.12	
2 + AUGCAL	92.87 \pm 0.06	6.84 \pm 0.10	89.25 \pm 0.36	85.74 \pm 0.36	67.80 \pm 0.78	

(a) GTAV \rightarrow Cityscapes. (DAFormer).

(b) VisDA SIM2REAL. (ViT-B).

Table 3: AUGCAL is better than applying AUG or CAL alone. On GTAV \rightarrow Cityscapes and VisDA, we show that AUGCAL improves over just augmented SIM training (AUG) or just optimizing for calibration on SIM data (CAL). For AUGCAL, we use PASTA as AUG and DCA as CAL. \pm indicates standard error.

Method	Perf. (\uparrow)		Calibration Error (\downarrow)		Reliability (\uparrow)	
	mIoU	ECE	IC-ECE	OC	PRR	
1 EntMin + MIC	65.71	5.34 \pm 0.35	77.73 \pm 0.26	82.83 \pm 0.55	45.93 \pm 0.54	
2 + AUG	67.58	4.30 \pm 0.33	77.59 \pm 0.25	48.05 \pm 0.53		
3 + CAL	68.70	4.04 \pm 0.26	75.86 \pm 0.26	52.52 \pm 0.54		
4 + AUGCAL	70.31	3.43 \pm 0.29	72.97 \pm 0.26	62.66 \pm 0.55		

(a) GTAV \rightarrow Cityscapes. (DAFormer).

(b) VisDA SIM2REAL. (ViT-B).

▷ **AUGCAL reduces miscalibration post SIM2REAL adaptation.** On both GTAV \rightarrow Cityscapes and VisDA, we find that AUGCAL consistently reduces miscalibration of the base method by reducing overconfidence on incorrect predictions. This is evident in how AUGCAL variants of the base adaptation methods have lower ECE, IC-ECE and OC values (AUGCAL rows, Calibration columns in Tables. 2 (a) and (b)). As an example, to illustrate the effect of improved calibration on real data, in Fig. 3, we show how applying AUGCAL can improve the proportion of per-pixel SemSeg predictions that are accurate and have high-confidence (> 0.95).

▷ **AUGCAL improvements in calibration improve reliability.** As noted earlier, we additionally investigate the extent to which calibration improvements for SIM2REAL adaptation translate to reliable confidence scores – via misclassification detection on REAL target data (see Sec. 4), as measured by PRR. We find that AUGCAL consistently improves PRR of the base SIM2REAL adaptation method (PRR columns for AUGCAL rows in Tables. 2(a) and (b)) – ensuring that predictions made AUGCAL variants of a base model are more trustworthy.

5.2 ANALYZING AUGCAL

We now analyze different aspects of AUGCAL.

▷ **Applying AUGCAL is better than applying just AUG or CAL.** In Sec. 3.2.1 and 3.2.2, we discuss how AUGCAL can be more effective in reducing target miscalibration than just optimizing for improved calibration on labeled SIM images. We verify this empirically in Tables. 3(a) and (b) for SemSeg and ObjRec. We show that while AUG and CAL, when applied individually, improve over a base SIM2REAL method, they fall short of improvements offered by AUGCAL.

▷ **AUGCAL is applicable across multiple AUG choices.** In Sec. 3.2.2 and Table. 1, we show how both PASTA (Chattopadhyay* et al., 2023) and RandAugment (Cubuk et al., 2020) are eligible for AUGCAL. In Table. 4(a), we fix DCA as CAL and find that both PASTA and RandAugment are effective in retaining or improving performance, reducing miscalibration and improving reliability.

▷ **Ablating CAL choices for AUGCAL.** For completeness, we also conduct experiments by fixing PASTA as AUG and ablating the choice of CAL in AUGCAL. We consider recently proposed training time calibration objectives – Difference of Confidence and Accuracy (DCA) (Liang et al.), Multi-class Difference in Confidence and Accuracy (MDCA) (Hebbalaguppe et al., 2022) and Margin-based Label Smoothing (MbLS) (Liu et al., 2022) – as potential CAL choices (results for SemSeg outlined in Table. 4(b)). We find that while MDCA and MbLS can be helpful, DCA is more consistently helpful across tasks and settings.

▷ **AUGCAL is applicable across multiple task backbones.** Different architectures – CNNs and Transformers – are known to exhibit bias towards different properties in images (shape, texture, etc.) (Naseer et al., 2021). Since the choice of AUG transform (which can alter such properties) is central to the efficacy of AUGCAL, we verify if AUGCAL is effective across both CNN and Transformer backbones. To do this, we conduct our SemSeg, ObjRec experiments with both transformer (DAFormer, ViT-B) and CNN (DeepLabv2-R101, ResNet-101) architectures. We find that AUGCAL

Table 4: Ablating AUG and CAL choices in AUGCAL. For a DAFormer model on GTAV→Cityscapes, AUGCAL successfully reduces miscalibration and produces reliable confidence scores for SIM2REAL adaptation using both PASTA (P) and RandAug (R) as AUG choices. We also ablate the choice of CAL in AUGCAL across DCA, MDCA and MbLS and find that DCA is more consistently effective in reducing miscalibration across tasks and settings. $\lambda_{\text{CAL}} = 1$ for MDCA and $\lambda_{\text{CAL}} = 0.1, m = 10$ for MbLS. \pm indicates standard error.

Method	AUG	Perf. (\uparrow) mIoU		Calibration Error (\downarrow) ECE		Reliability (\uparrow) PRR		Method	CAL	Perf. (\uparrow) mIoU		Calibration Error (\downarrow) ECE		Reliability (\uparrow) PRR			
		3 EntMin	4 + AUGCAL P	4 + AUGCAL R	7 HRDA	8 + AUGCAL P	8 + AUGCAL R			3 EntMin + MIC	4 + AUGCAL DCA	4 + AUGCAL MDCA	4 + AUGCAL MbLS	7 HRDA + MIC	8 + AUGCAL DCA	8 + AUGCAL MDCA	8 + AUGCAL MbLS
3 EntMin		65.71	5.34 \pm 0.35	77.73 \pm 0.26	45.93 \pm 0.54			3 EntMin + MIC	DCA	65.71	5.34 \pm 0.35	77.73 \pm 0.26	45.93 \pm 0.54				
4 + AUGCAL P	P	70.31	3.43 \pm 0.29	72.97 \pm 0.26	62.66 \pm 0.55			4 + AUGCAL DCA		70.31	3.43 \pm 0.29	72.97 \pm 0.26	62.66 \pm 0.55				
4 + AUGCAL R	R	70.65	2.34 \pm 0.14	73.77 \pm 0.21	66.65 \pm 0.46			4 + AUGCAL MDCA		69.50	3.22 \pm 0.26	72.65 \pm 0.25	59.96 \pm 0.51				
7 HRDA		75.56	2.86 \pm 0.10	81.92 \pm 0.14	68.91 \pm 0.46			4 + AUGCAL MbLS		68.77	2.90 \pm 0.24	72.53 \pm 0.23	61.57 \pm 0.48				
8 + AUGCAL P	P	75.90	2.49 \pm 0.09	79.09 \pm 0.16	70.35 \pm 0.51			7 HRDA + MIC		75.56	2.86 \pm 0.10	81.92 \pm 0.14	68.91 \pm 0.46				
8 + AUGCAL R	R	74.10	2.77 \pm 0.17	77.94 \pm 0.18	69.46 \pm 0.46			8 + AUGCAL DCA		75.90	2.45 \pm 0.09	79.09 \pm 0.16	70.35 \pm 0.51				

(a) Ablating AUG in AUGCAL. (CAL = DCA).

(b) Ablating CAL in AUGCAL. (AUG = PASTA)

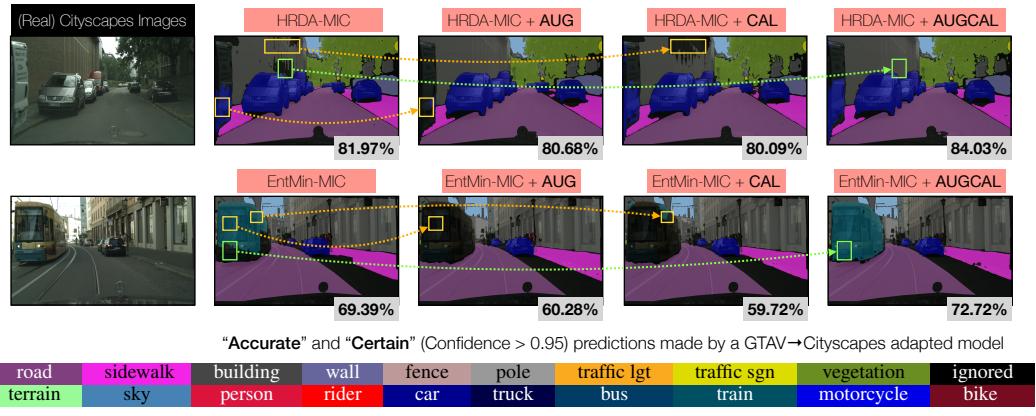


Figure 3: AUGCAL increases the proportion of “accurate” and “certain” predictions. For a (DAFormer) HRDA + MIC (row 1) and EntMin + MIC (row 2) on GTAV→Cityscapes, we show how different interventions affect the proportion of “accurate” and “certain” (confidence > 0.95) predictions (indicated in gray per column). Regions in black do not satisfy the “accurate” and “certain” filtering criteria. We see that compared to a base adaptation method, AUGCAL increases the proportion highly-confident correct predictions (green boxes). AUG and CAL applied alone can potentially reduce that proportion (yellow boxes). AUG is PASTA, CAL is DCA. (with PASTA as AUG and DCA as CAL) is effective in reducing SIM2REAL miscalibration across all settings. We discuss these results in Sec. A.5 of appendix.

▷ **How does AUGCAL compare with temperature scaling?** While we focus on “training-time” patches to improve SIM2REAL calibration, we also conduct an experiment to compare AUGCAL with “post-hoc” temperature scaling (TS) on VisDA. Specifically, we use 80% of VisDA SIM images for training models and rest (20%) for validation and temperature tuning. To ensure a fair comparison, we consider temperature tuning on both “clean” (C) and “PASTA augmented” (P) val splits. We find that irrespective of tuning on C or P, unlike AUGCAL, TS is ineffective and increases overconfidence and miscalibration. We present these results in Sec. A.5 of appendix.

▷ **AUGCAL increases the proportion of accurate and certain predictions.** In Fig. 3, we show qualitatively for GTAV→Cityscapes SemSeg how AUGCAL increases the proportion of highly-confident correct predictions. In practice, we find that this improvement is much more subtle for stronger SIM2REAL adaptation methods, such as HRDA + MIC, compared to weaker ones, such as EntMin + MIC, which have considerable room for improvement.

6 CONCLUSION

We propose AUGCAL, a method to reduce the miscalibration of SIM2REAL adapted models, often caused due to highly-confident incorrect predictions. AUGCAL modifies a SIM2REAL adaptation framework by making two minimally invasive changes – (1) augmenting SIM images via AUG transformations that reduce SIM2REAL distance and (2) optimizing for an additional calibration loss on AUGmented SIM predictions. Applying AUGCAL to existing adaptation methods for semantic segmentation and object recognition reduces miscalibration, overconfidence and improves reliability of confidence scores, all while retaining or improving performance on REAL data. AUGCAL is meant to be a task-agnostic, general purpose framework to reduce miscalibration for SIM2REAL adaptation methods and we hope such simple methods are taken into consideration for experimental settings beyond the ones considered in this paper.

Acknowledgements. This work has been partially sponsored by NASA University Leadership Initiative (ULI) #80NSSC20M0161, ARL and NSF #2144194.

7 ETHICS STATEMENT

Our proposed patch, AUGCAL, is meant to improve the reliability of SIM2REAL adapted models. We assess reliability in terms of confidence calibration (prediction scores aligning with true likelihood of correctness) and the extent to which calibrated confidence scores are useful for assessing prediction quality (measured via mis-classification detection). AUGCAL adapted models have promising consequences for downstream applications. A well-calibrated and reliable SIM2REAL adapted model can increase transparency in REAL predictions and facilitate robust decision making in safety-critical scenarios. That said, we would like to note that while AUGCAL is helpful for our specific measures of reliability, exploration along other domain specific notions of reliability remain.

REFERENCES

- Bianca-Cerasela-Zelia Blaga and Sergiu Nedevschi. Semantic segmentation learning for autonomous uavs using simulators and real data. In *2019 IEEE 15th International Conference on Intelligent Computer Communication and Processing (ICCP)*, pp. 303–310, 2019. doi: 10.1109/ICCP48234.2019.8959563.
- Ondrej Bohdal, Yongxin Yang, and Timothy Hospedales. Meta-calibration: Learning of model calibration using differentiable expected calibration error. *arXiv preprint arXiv:2106.09613*, 2021.
- Prithvijit Chattopadhyay, Judy Hoffman, Roozbeh Mottaghi, and Aniruddha Kembhavi. Robustnav: Towards benchmarking robustness in embodied navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15691–15700, 2021.
- Prithvijit Chattopadhyay*, Kartik Sarangmath*, Vivek Vijaykumar, and Judy Hoffman. Pasta: Proportional amplitude spectrum training augmentation for syn-to-real domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 19288–19300, October 2023.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- Sungha Choi, Sanghun Jung, Huiwon Yun, Joanne T Kim, Seungryong Kim, and Jaegul Choo. Robustnet: Improving domain generalization in urban-scene segmentation via instance selective whitening. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11580–11590, 2021.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.
- Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting. In *NIPS*, 2010.
- Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 702–703, 2020.
- Pau de Jorge, Riccardo Volpi, Philip Torr, and Gregory Rogez. Reliability in semantic segmentation: Are we on the right track? *arXiv preprint arXiv:2303.11298*, 2023.
- Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, et al. Robothor: An open simulation-to-real embodied ai platform. In *CVPR*, 2020.

- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 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. *arXiv preprint arXiv:2010.11929*, 2020.
- Abolfazl Farahani, Sahar Voghieri, Khaled Rasheed, and Hamid R Arabnia. A brief review of domain adaptation. *Advances in Data Science and Information Engineering: Proceedings from IC DATA 2020 and IKE 2020*, pp. 877–894, 2021.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. *arXiv preprint arXiv:1409.7495*, 2014.
- Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.
- Yunye Gong, Xiao Lin, Yi Yao, Thomas G Dietterich, Ajay Divakaran, and Melinda Gervasio. Confidence calibration for domain generalization under covariate shift. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8958–8967, 2021.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.
- Ramya Hebbalaguppe, Jatin Prakash, Neelabh Madan, and Chetan Arora. A stitch in time saves nine: A train-time regularizing loss for improved neural network calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16081–16090, 2022.
- Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning*, pp. 1989–1998, 2018.
- Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9924–9935, 2022a.
- Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Hrda: Context-aware high-resolution domain-adaptive semantic segmentation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXX*, pp. 372–391. Springer, 2022b.
- Lukas Hoyer, Dengxin Dai, Haoran Wang, and Luc Van Gool. Mic: Masked image consistency for context-enhanced domain adaptation. *arXiv preprint arXiv:2212.01322*, 2022c.
- Jiaxing Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. Fsdr: Frequency space domain randomization for domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6891–6902, 2021.
- Mobarakol Islam, Lalithkumar Seenivasan, Hongliang Ren, and Ben Glocker. Class-distribution-aware calibration for long-tailed visual recognition. *arXiv preprint arXiv:2109.05263*, 2021.

- Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. Minimum class confusion for versatile domain adaptation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*, pp. 464–480. Springer, 2020.
- Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. Contrastive adaptation network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4893–4902, 2019.
- Archit Karandikar, Nicholas Cain, Dustin Tran, Balaji Lakshminarayanan, Jonathon Shlens, Michael C Mozer, and Becca Roelofs. Soft calibration objectives for neural networks. *Advances in Neural Information Processing Systems*, 34:29768–29779, 2021.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- Meelis Kull, Telmo Silva Filho, and Peter Flach. Beta calibration: a well-founded and easily implemented improvement on logistic calibration for binary classifiers. In *Artificial Intelligence and Statistics*, pp. 623–631. PMLR, 2017.
- Aviral Kumar, Sunita Sarawagi, and Ujjwal Jain. Trainable calibration measures for neural networks from kernel mean embeddings. In *International Conference on Machine Learning*, pp. 2805–2814. PMLR, 2018.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017.
- 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, pp. 896, 2013.
- Gongbo Liang, Yu Zhang, Xiaoqin Wang, and Nathan Jacobs. Improved trainable calibration method for neural networks on medical imaging classification.
- Gongbo Liang, Yu Zhang, and Nathan Jacobs. Neural network calibration for medical imaging classification using dca regularization. In *International Conference on Machine Learning (ICML)*, 2020.
- Bingyuan Liu, Ismail Ben Ayed, Adrian Galdran, and Jose Dolz. The devil is in the margin: Margin-based label smoothing for network calibration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 80–88, 2022.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pp. 97–105. PMLR, 2015.
- Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. *Advances in neural information processing systems*, 31, 2018.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Andrey Malinin and Mark J. F. Gales. Predictive uncertainty estimation via prior networks. In *NeurIPS*, 2018.
- Andrey Malinin, Bruno Mloedeniec, and Mark Gales. Ensemble distribution distillation. *arXiv preprint arXiv:1905.00076*, 2019.
- Rhiannon Michelmore, Marta Kwiatkowska, and Yarin Gal. Evaluating uncertainty quantification in end-to-end autonomous driving control. *arXiv preprint arXiv:1811.06817*, 2018.
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances in Neural Information Processing Systems*, 34:15682–15694, 2021.

- Muhammad Akhtar Munir, Muhammad Haris Khan, M Sarfraz, and Mohsen Ali. Towards improving calibration in object detection under domain shift. *Advances in Neural Information Processing Systems*, 35:38706–38718, 2022.
- Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. *Advances in Neural Information Processing Systems*, 34:23296–23308, 2021.
- Henri J Nussbaumer. The fast fourier transform. In *Fast Fourier Transform and Convolution Algorithms*, pp. 80–111. Springer, 1981.
- Anusri Pampari and Stefano Ermon. Unsupervised calibration under covariate shift. In *arXiv:2006.16405*, 2020.
- 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 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- Nicolas Pinto, David D Cox, and James J DiCarlo. Why is real-world visual object recognition hard? *PLoS computational biology*, 4(1):e27, 2008.
- John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- 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. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Harsh Rangwani, Sumukh K Aithal, Mayank Mishra, Arihant Jain, and Venkatesh Babu Radhakrishnan. A closer look at smoothness in domain adversarial training. In *International Conference on Machine Learning*, pp. 18378–18399. PMLR, 2022.
- Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pp. 102–118. Springer, 2016.
- Stephan R Richter, Hassan Abu Alhaija, and Vladlen Koltun. Enhancing photorealism enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1700–1715, 2022.
- German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3234–3243, 2016.
- Alfréd Rényi. On measures of information and entropy. In *Proceedings of the 4th Berkeley Symposium on Mathematics, Statistics and Probability*, 1960.
- Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *European conference on computer vision*, pp. 213–226. Springer, 2010.
- Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Adversarial dropout regularization. In *International Conference on Learning Representations*, 2018.

- Swami Sankaranarayanan, Yogesh Balaji, Carlos D. Castillo, and Rama Chellappa. Generate to adapt: Aligning domains using generative adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019.
- Guardian Tesla Crash. Tesla driver dies in first fatal crash while using autopilot mode — theguardian.com. <https://www.theguardian.com/technology/2016/jun/30/tesla-autopilot-death-self-driving-car-elon-musk>. [Accessed 16-May-2023].
- Wilhelm Tranheden, Viktor Olsson, Juliano Pinto, and Lennart Svensson. Dacs: Domain adaptation via cross-domain mixed sampling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1379–1389, 2021.
- Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- 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*, pp. 7167–7176, 2017.
- Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2517–2526, 2019.
- Bailin Wang, Mirella Lapata, and Ivan Titov. Meta-learning for domain generalization in semantic parsing. *arXiv preprint arXiv:2010.11988*, 2020a.
- Chong Xiao Wang and Wee Peng Tay. Practical bounds of kullback-leibler divergence using maximum mean discrepancy. *arXiv preprint arXiv:2204.02031*, 2022.
- Deng-Bao Wang, Lei Feng, and Min-Ling Zhang. Rethinking calibration of deep neural networks: Do not be afraid of overconfidence. *Advances in Neural Information Processing Systems*, 34: 11809–11820, 2021.
- Dongdong Wang, Boqing Gong, and Liqiang Wang. On calibrating semantic segmentation models: Analysis and an algorithm. *arXiv preprint arXiv:2212.12053*, 2022.
- Ximei Wang, Mingsheng Long, Jianmin Wang, and Michael Jordan. Transferable calibration with lower bias and variance in domain adaptation. *Advances in Neural Information Processing Systems*, 33:19212–19223, 2020b.
- Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021.
- Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. Bridging theory and algorithm for domain adaptation. In *International Conference on Machine Learning*, pp. 7404–7413, 2019.
- Yuyang Zhao, Zhun Zhong, Na Zhao, Nicu Sebe, and Gim Hee Lee. Style-hallucinated dual consistency learning for domain generalized semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 289–305, 2018.

A APPENDIX

A.1 OVERVIEW

This appendix is organized as follows. In Sec. A.2.2, we provide more details on the AUG choices for AUGCAL– PASTA and RandAugment. Sec. A.3 outlines training, implementation details and choice of hyperparameters. Then, in Sec. A.4, we discuss Prediction Rejection Ratio (PRR), the metric used to assess reliability in our experiments. In Sec. A.5, we provide AUGCAL results for CNN backbones, comparisons with temperature scaling and discuss sensitivity to λ_{CAL} (coefficient of \mathcal{L}_{CAL}). Then, we provide more supporting qualitative examples in Sec. A.6. In Sec. A.7, we describe the SIM2REAL adaptation methods we use in detail. Finally, Sec. A.8 summarizes the assets used for our experiments and their associated licenses.

A.2 AUG CHOICES

A.2.1 PASTA

We use Proportional Amplitude Spectrum Training Augmentation (PASTA) ([Chattopadhyay* et al., 2023](#)) as one of the AUG choices in AUGCAL. PASTA applies structured perturbations (controlled by hyper-parameters α, β, k) to the amplitude spectra of synthetic images to generate augmented views. For a single-channel image, $x \in \mathbb{R}^{H \times W}$ (for illustration purposes), the augmentation process in PASTA is outlined below. We refer the reader to ([Chattopadhyay* et al., 2023](#)) for more details.

1. Set $\alpha = 3.0, \beta = 0.25, k = 2$ for PASTA
2. Use FFT ([Nussbaumer, 1981](#)) to obtain the Fourier spectrum of synthetic image x , as $\mathcal{F}(x) = \text{FFT}(x) \in \mathbb{C}^{H \times W}$
3. Obtain the corresponding amplitude and phase spectra as $\mathcal{A}(x) = \text{Abs}(\mathcal{F}(x))$ and $\mathcal{P}(x) = \text{Ang}(\mathcal{F}(x))$ respectively.
4. Zero-center the amplitude spectrum, as $\mathcal{A}(x) = \text{FFTShift}(\mathcal{A}(x))$, to get the lowest frequency components at the center.
5. Define perturbation strength, $\sigma \in \mathbb{R}^{H \times W}$, as $\sigma[m, n] = \left(2\alpha\sqrt{\frac{m^2+n^2}{H^2+W^2}}\right)^k + \beta$, where m, n denote the spatial frequencies.
6. Sample perturbations using the perturbation strength as, $\epsilon \sim \mathcal{N}(1, \sigma^2)$
7. Perturb the amplitude spectrum $\hat{\mathcal{A}}(x) = \epsilon \odot \mathcal{A}(x)$
8. Reset the low-frequency components, $\hat{\mathcal{A}}(x) = \text{FFTShift}(\hat{\mathcal{A}}(x))$
9. Obtain the augmented synthetic image via inverse FFT as, $\hat{x} = \text{iFFT}(\hat{\mathcal{A}}(x), \mathcal{P}(x))$

We set PASTA hyper-parameters as $\alpha = 3.0, \beta = 0.25, k = 2$ as it seems to work well across multiple SIM2REAL shifts in practice. We use Pytorch ([Paszke et al., 2019](#)) to vectorize all operations in the steps above and apply PASTA on minibatches, instead of individual synthetic images.

A.2.2 RANDAUGMENT

RandAugment ([Cubuk et al., 2020](#)) generates augmented views by applying a series of transforms sampled from a predefined vocabulary. While the vocabulary includes geometric as well as photometric transforms, we only use photometric transforms for our experiments. For SemSeg, the models we train already use geometric transforms that are optimal for SIM2REAL SemSeg performance. For ObjRec, we find it empirically beneficial to fix geometric transforms to a CenterCrop and use photometric transforms from RandAugment. These choices follow ([Hoyer et al., 2022c](#)). These operations are – AutoContrast, Equalize, Contrast, Brightness, Sharpness, Posterize, Solarize and SolarizeAdd. We use ($m = 30, n = 8$) while sampling operations from this vocabulary for RandAugment. We refer the reader to ([Cubuk et al., 2020](#)) for more details.

A.3 TRAINING AND IMPLEMENTATION DETAILS

We outline training and implementation details associated with our experiments.

Semantic Segmentation. For SemSeg, we use GTAV ([Sankaranarayanan et al., 2018](#)) as the source SIM dataset and Cityscapes ([Cordts et al., 2016](#)) as the target real dataset. We consider two

Table 5: Reproducing SIM2REAL Adaptation Baselines. We summarize results from our reproduction of the SIM2REAL adaptation baselines (from (Hoyer et al., 2022c)) used in our experiments. We find a slight difference between reported and reproduced adaptation results across models on both (a) GTAV→Cityscapes and (b) VisDA syn→real. \pm indicates standard error.

Method	Model	Status	mIoU (\uparrow)	Method	Model	Status	mAcc (\uparrow)
1 HRDA + MIC	DeepLabv2 (Hoyer et al., 2022c)	64.20	1 SDAT + MIC	ResNet-101 (Hoyer et al., 2022c)	86.90		
2 HRDA + MIC	DeepLabv2	Reproduced	64.05	2 SDAT + MIC	ResNet-101	Reproduced	81.03 ± 2.32
3 HRDA + MIC	DAFormer (Hoyer et al., 2022c)	75.90	3 SDAT + MIC	ViT-B (Hoyer et al., 2022c)	92.80		
4 HRDA + MIC	DAFormer	Reproduced	75.56	4 SDAT + MIC	ViT-B	Reproduced	92.62 ± 0.28

(a) GTAV→Cityscapes.

(b) VisDA SIM2REAL.

segmentation architectures for our experiments – DeepLabv2 (Chen et al., 2017) (with a ResNet-101 (He et al., 2016) backbone) and DAFormer (Hoyer et al., 2022a) (with an MiT-B5 (Xie et al., 2021) backbone). Both models are initialized from backbones pretrained on ImageNet (Deng et al., 2009). For EntMin (Vu et al., 2019) (with DAFormer / DeepLabv2), we use SGD as an optimizer with a learning rate of 2.5×10^{-4} and use $\lambda_{UDA} = 0.001$ as the coefficient of the “unconstrained” entropy loss. For HRDA, we use the multi-resolution self-training strategy from (Hoyer et al., 2022b) – AdamW (Loshchilov & Hutter, 2017) as the optimizer with a learning rate of 6×10^{-5} for the encoder and 6×10^{-4} for the decoder, with a linear learning rate warmup (warmup iterations 1500; warmup ratio 10^{-6}), followed by polynomial decay (to an eventual learning rate of 0). For the teacher-student self-training setup in HRDA, we additionally couple cross-domain mixing with augmentations (DACS (Tranheden et al., 2021)) to improve self-training, use the ImageNet feature distance loss (Hoyer et al., 2022a), set $\lambda_{UDA} = 1$ and use $\alpha = 0.999$ as the teacher EMA factor. For all SemSeg settings, we use a batch size of 2 (2 source images, 2 target images) and additionally use rare class sampling (based on source; following (Hoyer et al., 2022a)) to ensure consistent adaptation improvements across all classes. For all our SemSeg experiments, we feed crops of size 1024×1024 as input to the models, irrespective of the adaptation method. PASTA for AUGCAL is applied to the same crops. We train all segmentation models for 40k iterations and use the last checkpoint to report results (the SIM2REAL adaptation setting does not assume access to labels on real data which prevents selecting “best-on-target” checkpoint). For semantic segmentation, following prior work (Wang et al., 2022), for ECE (and other associated metrics), instead of pooling all pixels of different images into a set, we compute per-image numbers and then average across images. We use the same to compute standard error. We report all metrics (performance, calibration, etc.) as percentages. We use 15 bins to compute ECE.

Object Recognition. For ObjRec, we conduct experiments on the VisDA (Peng et al., 2017) SIM2REAL benchmark, using standard source-target splits (Hoyer et al., 2022c) for training. We consider ImageNet (Deng et al., 2009) pretrained ResNet-101 (He et al., 2016) and ViT-B/16 (Dosovitskiy et al., 2020) backbones with standard bottleneck (with Linear, BatchNorm and ReLU) and classifier layers. As noted earlier, we use SDAT (Rangwani et al., 2022) as the adaptation method, which relies on Conditional Domain Adversarial Adaptation (CDAN) (Long et al., 2018) and Minimum Class Confusion (MCC) (Jin et al., 2020) with SAM (Foret et al., 2020) (smoothness 0.2) – all adaptation losses are combined and λ_{UDA} is set to 1. We use SGD as the base optimizer with a learning rate of 2×10^{-4} , with a batch size of 32 (for both source and target). PASTA for AUGCAL is applied to the raw input SIM images. We train ResNet-101 backbones for 30 epochs and ViT-B/16 backbones for 15 epochs and use the last checkpoint to report results (the SIM2REAL adaptation setting does not assume access to labels on real data which prevents selecting “best-on-target” checkpoint). We report all metrics (performance, calibration, etc.) as percentages. We use 15 bins to compute ECE. For our key results in the main paper (Tables. 2(b) & 3(b)), we run experiments across 3 random seeds.

MIC Hyperparameters. We also use MIC (Hoyer et al., 2022c) with all the discussed SIM2REAL adaptation methods since it demonstrably improves performance. For MIC, following prior work (Hoyer et al., 2022c), we use a masking patch size of 64, a masking ratio of 0.7, a loss weight of 1 and an EMA factor of 0.999 for the pseudo-label generating teacher.

Compute. We conduct all object recognition experiments on RTX 6000 GPUs – every experiment requiring a single GPU. For semantic segmentation, we use one A40 GPU per experiment.

Table 6: AUGCAL ensures SIM2REAL adapted models make accurate, calibrated and reliable predictions. We find that applying AUGCAL to multiple SIM2REAL adaptation methods across tasks leads to better calibration (ECE, IC-ECE), reduced overconfidence (OC) and improved reliability (PRR) – all while retaining or improving transfer performance. Highlighted rows are AUGCAL variants of the base methods. For AUGCAL, we use PASTA as AUG and DCA as CAL. \pm indicates standard error.

Method	Perf. (\uparrow) mIoU	Calibration Error (\downarrow) ECE	Reliability (\uparrow) IC-ECE	Reliability (\uparrow) OC	PRR
1 EntMin + MIC	47.77	4.93 \pm 0.20	70.09 \pm 0.22	75.80 \pm 0.58	48.87 \pm 0.45
2 + AUGCAL	49.96	3.21 \pm 0.16	67.99 \pm 0.22	74.33 \pm 0.57	52.07 \pm 0.48
3 HRDA + MIC	64.05	3.52 \pm 0.19	78.94 \pm 0.18	86.48 \pm 0.50	63.20 \pm 0.48
4 + AUGCAL	63.95	2.55 \pm 0.11	74.71 \pm 0.21	84.13 \pm 0.52	65.37 \pm 0.50

(a) GTAV \rightarrow Cityscapes. (DeepLabv2 R-101).

(b) VisDA SIM2REAL. (ResNet-101).

Table 7: Comparing AUGCAL with Temperature Scaling for VisDA SIM2REAL. We do a controlled experiment (with an 80-20 split of the VisDA SIM2REAL split) to compare AUGCAL with Temperature Scaling (TS). B = SDAT + MIC (ViT-B). C (clean) and P (PASTA augmented) indicate the synthetic “labeled” validation splits the temperature was tuned on. For AUGCAL, we use PASTA as AUG and DCA as CAL. \pm indicates standard error.

Method	Perf. (\uparrow) mIoU	Calibration Error (\downarrow) ECE	Reliability (\uparrow) IC-ECE	Reliability (\uparrow) OC	PRR
1 SDAT + MIC (B)	92.24 \pm 0.05	8.13 \pm 0.36	91.54 \pm 0.05	89.48 \pm 0.43	63.92 \pm 1.46
1 B + TS (C)	92.24 \pm 0.05	8.19 \pm 0.26	95.41 \pm 0.08	93.97 \pm 0.20	66.58 \pm 1.23
1 B + TS (P)	92.24 \pm 0.30	8.54 \pm 0.40	93.52 \pm 0.36	92.00 \pm 0.42	64.05 \pm 1.50
2 B + AUGCAL	92.68 \pm 0.17	7.24 \pm 0.11	90.40 \pm 0.09	87.57 \pm 0.61	66.49 \pm 0.44

(b) AUGCAL vs Temperature Scaling.

(b) AUGCAL + Temperature Scaling.

Reproduced Results. We use open-sourced code for (Hoyer et al., 2022c)³ and re-run the base adaptation methods for HRDA + MIC and SDAT + MIC at our end to obtain baseline results. In Table. 5, we summarize reported and reproduced results for the same methods.

A.4 PREDICTION REJECTION RATIO (PRR)

As noted in Sec. 4 of the main paper, in addition to measuring reduced miscalibration, we also assess the extent to which such improvements in calibration are useful and lead to a more reliable model. To measure this, following prior work (de Jorge et al., 2023; Malinin et al., 2019), we measure whether confidence scores can reliably guide misclassification detection – measured via the PRR metric.

Ideally, in a real-world scenario, given a model, we would like to retrieve all (potentially) samples misclassified by the model based on confidence scores. These samples can then be skipped when the model is used for decision-making (since the model is likely to make incorrect predictions). Since confidence scores, in practice, are imperfect, measuring misclassification detection helps us assess this specific capability of the SIM2REAL adapted models. For a model that is less overconfident on incorrect samples (meaning it has reduced miscalibration), this specific ability should naturally be enhanced.

This can be measured using Rejection-Accuracy curves (de Jorge et al., 2023; Malinin et al., 2019) where we reject samples below a threshold and keep track of accuracy and the fraction of rejected samples. Since such curves are naturally biased towards models that have improved performance, the AUC for such a rejection curve can be normalized by that of an oracle. Additionally, we can subtract a score associated with a random baseline (sorting predictions for filtering in a random order) (Wang et al., 2021). Finally, we can compare this value (PRR; ranging from -100 to 100; higher is better) for multiple models to assess how reliable underlying confidence scores are.

A.5 AUGCAL RESULTS

▷ **AUGCAL results with CNN architectures.** Key results presented in Table. 2 of the main paper use transformer architectures (DAFormer for SemSeg, ViT-B for ObjRec). In Tables. 6 (a) and (b), we verify that AUGCAL improvements translate to CNN based architectures as well (also discussed in Sec. 5.2 of the paper). We use DeepLabv2 (ResNet-101) for SemSeg and a ResNet-101 based classifier for ObjRec and find that applying AUGCAL improves or retains performance, reduces miscalibration and improves reliability.

▷ **Comparing AUGCAL with Temperature Scaling (TS).** While we focus on “training-time” patches to improve SIM2REAL calibration, we also conduct an experiment to compare AUGCAL with “post-hoc” temperature scaling (TS) on VisDA. Specifically, we use 80% of VisDA SIM images for training

³<https://github.com/lhoyer/MIC>

Table 8: Sensitivity to λ_{CAL} for AUGCAL on GTAV→Cityscapes. We vary the value of $\lambda_{\text{CAL}} \in \{0.1, 0.5, 1.0, 5.0, 10.0, 20.0, 100.0\}$ and report the effect on adaptation performance, reduced miscalibration and improved reliability. We find that our choice of $\lambda_{\text{CAL}} = 1$ leads to balanced performance across desired metrics. Rows in red correspond to the baseline SIM2REAL adaptation method without the application of AUGCAL.

Value of λ_{CAL}	Perf. (\uparrow) mIoU	Calibration Error (\downarrow) ECE	Reliability (\uparrow) IC-ECE	PRR	Value of λ_{CAL}	Perf. (\uparrow) mIoU	Calibration Error (\downarrow) ECE	Reliability (\uparrow) IC-ECE	PRR
1 No AUGCAL	65.71	5.34	77.73	45.93	1 No AUGCAL	75.56	2.86	81.92	68.91
2 $\lambda_{\text{CAL}} = 0.1$	69.72	2.88	74.55	57.06	2 $\lambda_{\text{CAL}} = 0.1$	75.25	2.94	80.68	69.96
3 $\lambda_{\text{CAL}} = 0.5$	69.27	2.71	73.54	60.32	3 $\lambda_{\text{CAL}} = 0.5$	75.05	2.77	80.23	70.62
4 $\lambda_{\text{CAL}} = 1.0$	70.31	3.43	72.97	62.66	4 $\lambda_{\text{CAL}} = 1.0$	75.90	2.45	79.09	70.35
5 $\lambda_{\text{CAL}} = 5.0$	66.24	3.55	70.22	62.66	5 $\lambda_{\text{CAL}} = 5.0$	73.80	2.25	76.61	70.54
6 $\lambda_{\text{CAL}} = 10.0$	64.27	3.28	67.63	64.55	6 $\lambda_{\text{CAL}} = 10.0$	71.28	2.50	75.85	69.17
7 $\lambda_{\text{CAL}} = 20.0$	55.37	4.17	65.55	63.20	7 $\lambda_{\text{CAL}} = 20.0$	62.31	2.43	73.32	68.27
8 $\lambda_{\text{CAL}} = 100.0$	27.73	6.06	52.24	54.47	8 $\lambda_{\text{CAL}} = 100.0$	31.01	7.68	73.56	56.52

(a) EntMin + MIC (DAFormer).

(a) HRDA + MIC (DAFormer).

models and rest (20%) for validation and temperature tuning. To ensure a fair comparison, we consider temperature tuning on both "clean" (C) and "PASTA augmented" (P) val splits.⁴ We present these results in Table. 7 (a). We find that irrespective of tuning on C or P, unlike AUGCAL, TS is ineffective and increases overconfidence and miscalibration. In Table. 7, we additionally consider temperature scaling on the logits of an AUGCAL improved SIM2REAL model and find that it worsens miscalibration and overconfidence. Note that this is not entirely surprising since TS depends heavily on the "data-split" the temperature is tuned on, which in our case is SIM and not REAL.

▷ **Sensitivity to λ_{CAL} .** For an existing SIM2REAL adaptation pipeline, AUGCAL involves optimizing for an additional calibration loss \mathcal{L}_{CAL} on augmented synthetic images. In Table. 8, we vary the coefficient λ_{CAL} (values in the set $\{0.1, 0.5, 1.0, 5.0, 10.0, 20.0, 100.0\}$) for \mathcal{L}_{CAL} and note the effect on adaptation performance, calibration on real data (ECE, IC-ECE) and reliability (PRR). In Table 2 of the main paper, we already note how applying AUGCAL with $\lambda_{\text{CAL}} = 1$ improves over a baseline adaptation method (red and blue rows in Table. 8). We further note that compared to other values, our choice of $\lambda_{\text{CAL}} = 1.0$ achieves a balance between adaptation performance, reduced miscalibration and improved reliability. We find that overly high values of $\lambda_{\text{CAL}} \geq 5$ can potentially lead to reduced adaptation performance – $\lambda_{\text{CAL}} \geq 5$ significantly raises the scale of \mathcal{L}_{CAL} compared to \mathcal{L}_{CE} and \mathcal{L}_{UDA} , which leads to models optimizing for improved calibration at the expense of task performance. Based on our experiments across multiple models, shifts and tasks, we recommend restricting $\lambda_{\text{CAL}} < 5$ for SIM2REAL adaptation.

A.6 QUALITATIVE PREDICTIONS

In Figures 4 and 5, we provide more examples to demonstrate how AUGCAL improves the proportion of "accurate" and "certain" predictions. In Fig. 4, where we compare predictions for a (DAFormer) HRDA + MIC model – w/o AUGCAL 75.56 mIoU and w AUGCAL 75.90 mIoU. We find that applying AUGCAL improves the proportion of "accurate" and "certain" predictions (confidence > 0.95) – see Fig. 4, columns 3 and 5, guided by yellow arrows. For instance, we find that the base model (w/o AUGCAL) has trouble assigning high-confidence to correct "sidewalk" and "vegetation" predictions. For EntMin + MIC (w/o AUGCAL 65.71 mIoU and w AUGCAL 70.31 mIoU) in Fig. 5, we find that AUGCAL is considerably more effective in ensuring high-confidence correct predictions. Notably, we find these improvements to be more subtle as the SIM2REAL adaptation method itself improves (HRDA + MIC > EntMin + MIC).

A.7 SIM2REAL ADAPTATION METHODS

We conduct experiments with three SIM2REAL adaptation methods across two tasks. We wanted to assess the compatibility of AUGCAL with SIM2REAL adaptation methods that are competitive and are representative of the broader class of approaches used for SIM2REAL transfer. We first cover some background and describe these methods in detail.

⁴Note that temperature tuning (by definition) only affects calibration and has no impact on performance.

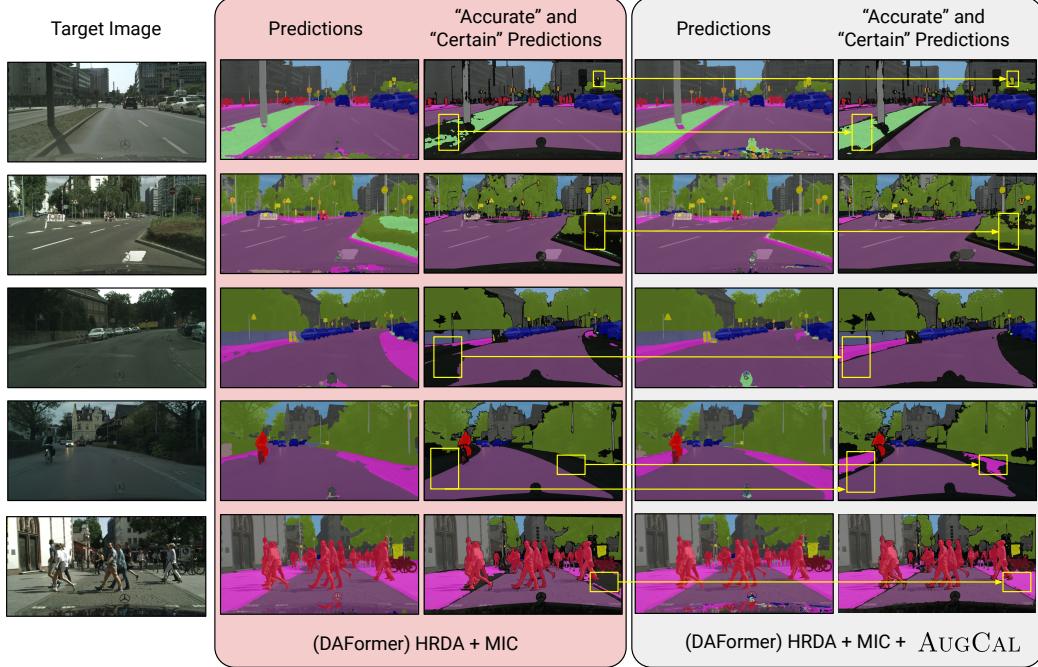


Figure 4: AUGCAL increases the proportion of “accurate” and “certain” predictions. For a base DAFormer SemSeg model trained with HRDA + MIC (State-of-the-art) on GTAV→Cityscapes, we show how applying AUGCAL at training time (right) can improve the proportion of “accurate” and “certain” (confidence > 0.95) predictions over the vanilla adaptation method (left). Regions in black do not satisfy the “accurate” and “certain” filtering criteria.

In unsupervised domain adaptation (UDA) for SIM2REAL settings, we assume access to a labeled (SIM) source dataset $D_S = \{(x_i^S, y_i^S)\}_{i=1}^{|S|}$ and an unlabeled (REAL) target dataset $D_T = \{x_i^T\}_{i=1}^{|T|}$. We assume D_S and D_T splits are drawn from source and target distributions $P^S(x, y)$ and $P^T(x, y)$ respectively. At training, we have access to $D = D_S \cup D_T$. We operate in the setting where source and target share the same label space, and discrepancies exist only in input images. The model \mathcal{M}_θ is trained on labeled source images using cross entropy,

$$\sum_{i=1}^{|S|} \mathcal{L}_{CE}(x_i^S, y_i^S; \theta) = - \sum_{i=1}^{|S|} y_i^S \log p_\theta(\hat{y}_i^S | x_i^S) \text{ where } \hat{y}_i^S = \arg \max_{y \in \mathcal{Y}} p_\theta(y_i^S | x_i^S) \quad (14)$$

For object recognition, we represent the labels corresponding to images $x \in \mathbb{R}^{H \times W \times 3}$ as $y \in \mathcal{Y}$ where $\mathcal{Y} = \{1, 2, \dots, K\}$. For semantic segmentation, since we make predictions across a vocabulary of classes for every pixel, the corresponding label can be expressed as $y \in \mathcal{Y}^{H \times W}$.

UDA methods additionally optimize for an adaptation objective on labeled source and unlabeled target data (\mathcal{L}_{UDA}). The overall learning objective can be expressed as,

$$\min_{\theta} \underbrace{\sum_{i=1}^{|S|} \mathcal{L}_{CE}(x_i^S, y_i^S; \theta)}_{\text{Source Loss}} + \underbrace{\sum_{i=1}^{|T|} \sum_{j=1}^{|S|} \lambda_{UDA} \mathcal{L}_{UDA}(x_i^T, x_j^S, y_j^S; \theta)}_{\text{Source Target Adaptation Loss}} \quad (15)$$

Different adaptation methods usually differ in terms of specific instantiations of this objective. As adaptation methods, we use HRDA (Hoyer et al., 2022b) and EntMin (Vu et al., 2019) for segmentation and SDAT (Rangwani et al., 2022) for recognition.

Entropy Minimization (EntMin). We consider the unconstrained direct entropy minimization approach from (Vu et al., 2019) as one of the adaptation methods. EntMin builds on top of the assumption the models trained only on source data tend to be under-confident (make high-entropy predictions) on target images. EntMin enforces high prediction certainty on target images by ensuring

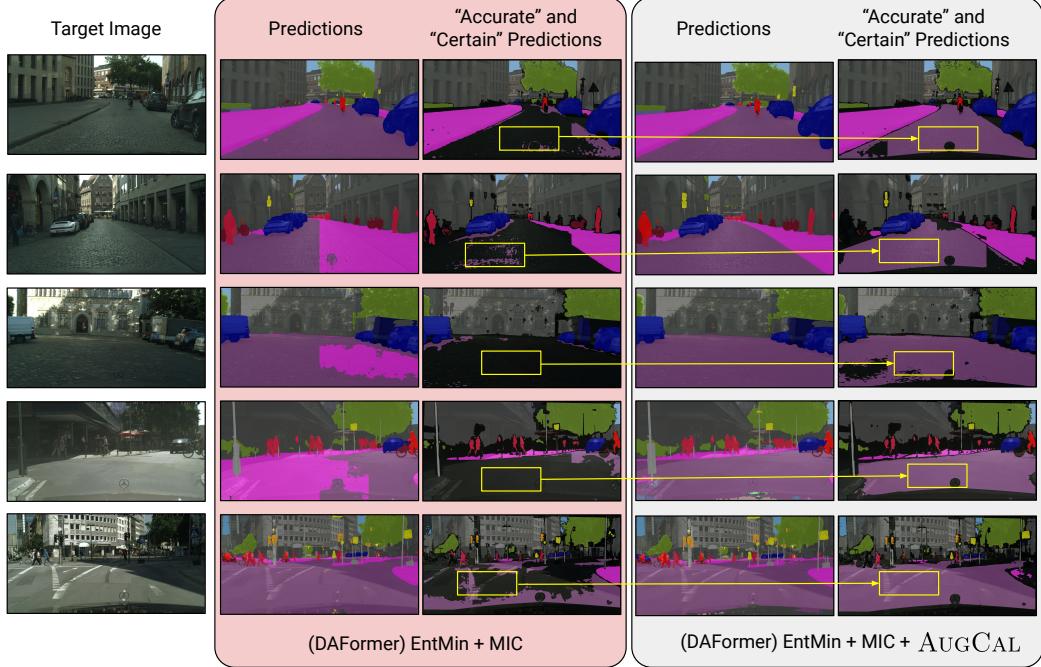


Figure 5: AUGCAL increases the proportion of “accurate” and “certain” predictions. For a base DAFormer SemSeg model trained with EntMin + MIC on GTAV→Cityscapes, we show how applying AUGCAL at training time (right) can improve the proportion of “accurate” and “certain” (confidence > 0.95) predictions over the vanilla adaptation method (left). Regions in black do not satisfy the “accurate” and “certain” filtering criteria.

that the model makes high-confidence predictions on the same. This is realized by minimizing the normalized entropy of target predictions. Specifically, given a target image x^T , if the predictive probabilities per-pixel are expressed as $\mathbf{P}_{x^T}^{(h,w,k)} \in [0, 1]^{H \times W \times K}$, the adaptation loss on target can be expressed as,

$$\mathcal{L}_{UDA}(x^T; \theta) = \frac{1}{HW} \sum_{h,w} \frac{-1}{\log K} \sum_{k=1}^K \mathbf{P}_{x^T}^{(h,w,k)} \log \mathbf{P}_{x^T}^{(h,w,k)} \quad (16)$$

EntMin can also be viewed as a “soft-assignment” version of self-training (Lee et al., 2013; Zou et al., 2018). We refer the reader to (Vu et al., 2019) for more details.

Context-Aware High-Resolution Domain-Adaptive Semantic Segmentation (HRDA). The base adaption method in HRDA (Hoyer et al., 2022b) is self-training (Zou et al., 2018; Lee et al., 2013). HRDA additionally introduces components that are beneficial for domain-adaptive semantic segmentation. Specifically, HRDA proposes a multi-resolution framework by relying on (1) a low-resolution *context* crop to learn long-range contextual dependencies and (2) a high-resolution *detail* crop to make detailed and accurate predictions. During adaptation, HRDA fuses predictions from both crops using input dependent attention. For a target image x^T , pseudo-labels are obtained from teacher network \mathcal{M}_ϕ (a moving average of \mathcal{M}_θ) as $\mathbf{P}_{x^T}^{(h,w)} = \arg \max_k \mathcal{M}_\phi(x^T)$. HRDA adapts to target by minimizing cross entropy w.r.t. $\mathbf{P}_{x^T}^{(h,w)}$ as,

$$\mathcal{L}_{UDA}(x^T; \theta) = - \sum_{h,w} q_{x^T} \mathbf{P}_{x^T}^{(h,w)} \log \mathcal{M}_\theta(x^T) \text{ with } q_{x^T} \text{ being a quality estimate for } \mathbf{P}_{x^T}^{(h,w)} \quad (17)$$

Instead of “self”-training on target images, HRDA uses cross-domain mixing (DACS (Tranheden et al., 2021)) to obtain pseudo-labels on augmented images. Additionally, components from DAFormer (Hoyer et al., 2022a) – rare class sampling and an ImageNet feature distance loss – are incorporated in the pipeline to facilitate better adaptation. The quality estimate q_{x^T} is computed as the proportion of pixels that have confidence above a specified threshold. This naturally ensures a warmup stage, where a model is first trained only on synthetic images for a few iterations, followed

by training on both synthetic and real images. We refer the reader to (Hoyer et al., 2022b) for more low-level implementation details.

Smooth Domain Adversarial Training (SDAT). SDAT (Rangwani et al., 2022) is an adaptation method for object recognition. The underlying adaptation method in SDAT is domain adversarial training (DAT), which involves reducing the discrepancy between source and target image distributions. This is realized by confusing an additional discriminator that is designed to distinguish between source and target samples. While multiple versions of DAT exist, SDAT uses CDAN (Long et al., 2018) as the default DAT method. SDAT investigates the loss-landscapes of DAT style methods and notes that smoother loss-landscapes on source data result in improved transfer to target. Consequently, SDAT proposes optimizing for smoother loss landscapes on labeled source data by modifying the supervised \mathcal{L}_{CE} loss as,

$$\sum_{i=1}^{|S|} \mathcal{L}_{CE}(x_i^S, y_i^S; \theta) = \sum_{i=1}^{|S|} y_i^S \max_{||\epsilon|| \leq \rho} \log p_{\theta+\epsilon}(\hat{y}_i^S | x_i^S) \quad (18)$$

Perturbing the weights of the network θ by ϵ in some neighborhood ρ ensures lower loss values in that neighborhood, thereby encouraging a smoother loss landscape. In practice, this is realized by using sharpness aware minimization (SAM) (Foret et al., 2020) to update model parameters. We refer the reader to (Rangwani et al., 2022; Long et al., 2018) for more details on the UDA losses used.

Masking Image Consistency (MIC). MIC (Hoyer et al., 2022c) is a general technique applicable to any existing adaptation method to improve context utilization while making predictions. MIC involves a teacher-student self-training setup where pseudo-labels are obtained from a weakly augmented target sample (via the teacher). Student predictions for a masked image are forced to match the obtained pseudo-labels. MIC relies on the recent success of masking as an auxiliary task (He et al., 2022) but instead of reconstruction, MIC sets up a prediction consistency task. This simple addition to an existing adaptation pipeline leads to considerable improvements across multiple tasks and shifts. A “masked” image for MIC is obtained by dividing the original image into patches and randomly masking a subset of the same. Similar to HRDA (Hoyer et al., 2022b), the MIC (consistency) loss is multiplied by a quality estimate of the pseudo-labels. We refer the reader to (Hoyer et al., 2022c) for more details on MIC.

A.8 ASSETS AND LICENCES

The assets used in this work can be grouped into three categories – Datasets, Code Repositories and Dependencies. We discuss source and licences for each of these below.

Datasets. For semantic segmentation, we use the GTAV (Richter et al., 2016) and Cityscapes (Cordts et al., 2016) datasets. Code used to extract densely annotated images from the GTAV game is distributed under the MIT license.⁵ The Cityscapes’ license agreement dictates that the dataset is made freely available to academic and non-academic entities for non-commercial purposes such as academic research, teaching, scientific publications, or personal experimentation and that permission to use the data is granted under certain conditions.⁶ For object recognition, we use the VisDA syn→real (Peng et al., 2017) benchmark. The VisDA-C development kit on github does not have a license associated with it, but it does include a Terms of Use, which primarily states that the dataset must be used for non-commercial and educational purposes only.⁷

Code Repositories. For our experiments, apart from code that we wrote ourselves, we build on top of the open-sourced codebase for MIC (Hoyer et al., 2022c)⁸ repository. MIC is distributed under the MIT License.

Dependencies. We use Pytorch (Paszke et al., 2019) as the deep-learning framework for all our experiments. Pytorch, released by Facebook, is distributed under a Facebook-specific license.⁹

⁵<https://bitbucket.org/visinf/projects-2016-playing-for-data/src/master/>

⁶<https://www.cityscapes-dataset.com/license/>

⁷<https://github.com/VisionLearningGroup/taskcv-2017-public/tree/master/classification>

⁸<https://github.com/lhoyer/MIC>

⁹<https://github.com/pytorch/pytorch/blob/master/LICENSE>