# Handling new target classes in semantic segmentation with domain adaptation

Maxime Bucher[a], Tuan-Hung Vu[a], Matthieu Cord[a,b], Patrick Pérez[a]

[a]*Valeo.ai, Paris, France*
[b]*Sorbonne University, Paris, France*

## ABSTRACT

In this work, we define and address a novel domain adaptation (DA) problem in semantic scene segmentation, where the target domain not only exhibits a data distribution shift w.r.t. the source domain, but also includes novel classes that do not exist in the latter. Different to "open-set" [1] and "universal domain adaptation" [2], which both regard all objects from new classes as "unknown", we aim at explicit test-time prediction for these new classes. To reach this goal, we propose a framework that leverages domain adaptation and zero-shot learning techniques to enable "boundless" adaptation in the target domain. It relies on a novel architecture, along with a dedicated learning scheme, to bridge the source-target domain gap while learning how to map new classes' *labels* to relevant visual representations. The performance is further improved using self-training on target-domain pseudo-labels. For validation, we consider different domain adaptation set-ups, namely synthetic-2-real, country-2-country and dataset-2-dataset. Our framework outperforms the baselines by significant margins, setting competitive standards on all benchmarks for the new task. Code and models are available at: https://github.com/valeoai/buda.

arXiv:2004.01130v2 [cs.CV] 16 Feb 2021

## 1. Introduction

Detailed interpretation of operating environments is crucial for autonomous systems like self-driving cars or delivery droids. Aiming at such an exhaustive scene understanding, most recent vision systems conduct semantic segmentation, the task of predicting semantic classes for all scene pixels. As for other perception modules, significant shifts between train (source domain) and test (target domain) distributions drastically degrade the segmentation performance. Many works, *e.g.*, [3, 4, 5], propose unsupervised domain adaptation (UDA) techniques to address such *domain gaps*, while alleviating the taxing need for data labeling in the target domain. Besides these domain adaptation (DA) works, where the same categories are assumed

in both domains, few recent works consider more general set-ups in which source and target label sets differ. For example, partial DA [6, 7] is the case where the source label set contains the target one. Differently, open-set DA [1] assumes that each domain holds a set of "private" classes besides the ones they share. More challenging, universal DA [2] considers having a completely unknown target label set. These preceding efforts contribute all to make domain adaptation closer to real-world applications. In open-set and universal DA settings, objects from that classes that are absent in the source domain are all classified as "unknown". Nevertheless, for real-life scenarios where the target label set is indeed *boundless*, one could expect the final system to predict new classes explicitly. For example, one might require that the perception model, once trained on European driving scenes, behaves well on Indian streets where

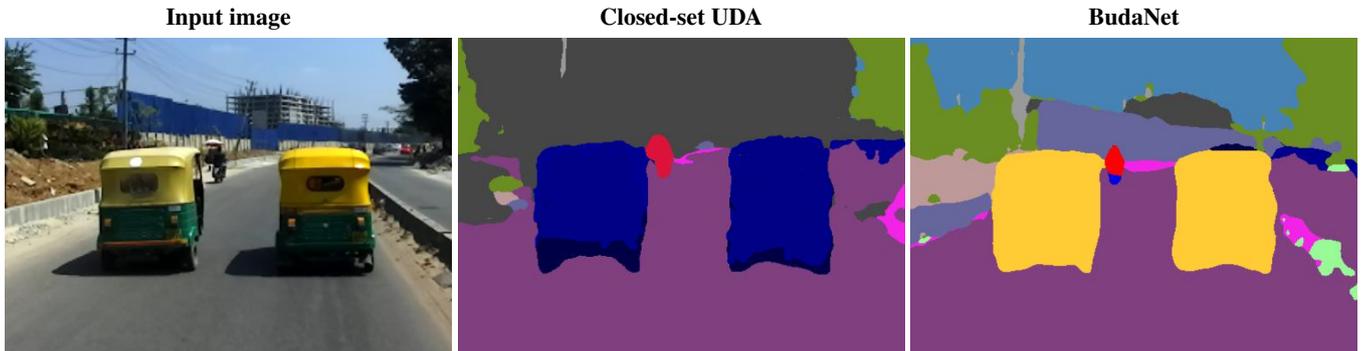| Input image | Closed-set UDA | BudaNet |
|---|---|---|



Fig. 1: **Illustration of BudaNet, the proposed approach to the new problem of boundless UDA for semantic segmentation**. At test time, for an input image (left), we display the closed-set UDA segmentation result (middle) as well as BudaNet's segmentation result (right). Both "tuk-tuk" vehicles, which are from a new class only appearing in the target domain, have been correctly identified by BudaNet: Our approach is able to deal with new classes for which no annotated images have been provided during training. The model trained following the closed-set UDA setting wrongly predicts these new vehicles as a spatial mix of car and truck.

vehicles similar to European ones cohabit with specific ones like "tuk-tuk". Zhuo *et al.* [8] address this realistic set-up for classification and propose a unified framework to handle new classes at test time while minding the domain gap.

All the above-mentioned DA works take image classification as the only test-bed.

Different to those works, we address here semantic segmentation with UDA, *i.e.*, the task of exhaustively classifying all pixels in input images in presence of a domain gap. We aim at pixel-wise recognition of objects from both shared and new classes in the target domain but, as each image may contain multiple objects, this UDA task is much more challenging than classification. As required in several real-life applications, the detailed analysis of urban scenes is an example where such complex spatial arrangements of multiple object categories are common place.

In this segmentation set-up, pixels of shared and new classes can co-occur in the same scene, requiring dedicated strategies to carefully handle source-target domain alignments; for example, as new-class objects only exist in the target domain, a naive alignment might undesirably force their visual features close to those from shared-class objects.

We call this framework "boundless" unsupervised domain adaptation", BUDA in short, and we propose ways to attack it in the present work. In particular, we introduce a full semantic segmentation pipeline, BudaNet, that jointly addresses two

main challenges: (1) Bridging the gap between source and target domains; (2) Learning discriminant visual representations for new-class objects in the target domain. Within a standard semantic segmentation framework, we propose a UDA strategy to mitigate the distribution gap for shared classes across domains, without causing unwanted misalignment to new ones. We then introduce a novel domain-aware model to generate, for all classes, pixel-level visual features for both domains. By using this generative model, we collect features from the new target-domain classes, along with features from shared classes in the source domain, to learn the last classification layer of BudaNet. Last, we refine BudaNet with a step of self-training using pseudo-labels on the target domain. The resulting behavior of our system is illustrated in Figure 1.

With this new form of DA task and the proposed BudaNet approach to solve it, we advocate for more practical and ubiquitous domain adaptation in which the target label set is boundless.

## 2. Related work and positioning

In this section, we briefly go through previous UDA works and position our setting with respect to them. We also review some zero-shot learning techniques that play an important role in the proposed framework.

*Unsupervised domain adaptation.* Most existing UDA works consider the classic setting in which source and target label sets
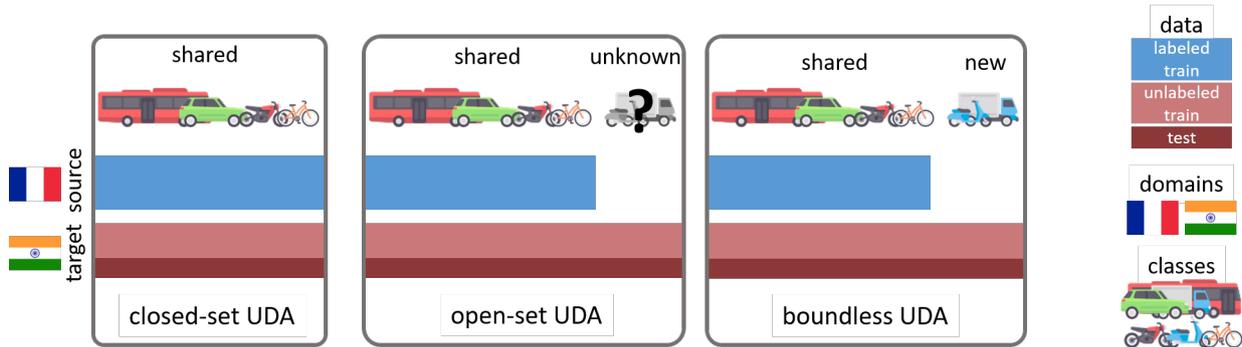
Fig. 2: **Different UDA settings**. In contrast with classic UDA ("closed-set" hypothesis), boundless UDA (BUDA) assumes that source and target label sets differ. In addition, unlike the "open-set" setting where objects from the private target classes are classified in a single *unknown* category, BUDA allows the target classes to be explicitly predicted.

are the same. Though approaching UDA via different angles, these works are in the same vein of learning task-dependent domain-invariant features, *i.e.*, minimizing inter-domain discrepancy of feature distributions. Popular techniques include regularizing the maximum mean discrepancy [9], matching deep activation correlation [10], aligning source-target distributions via adversarial training [11, 3], self-training [12, 13, 14] or self-ensembling [15, 16]. Recently, UDA for complex recognition tasks like detection [17] and semantic segmentation [18] have received more attention. Such tasks often require special techniques to handle as well the spatial layout [4, 5] or class proportions [14]. While the standard UDA setting facilitates investigations and helps gain fundamental insights, it is still far from the real-life scenarios.

Some recent works propose techniques in more relaxed UDA settings. In those cases, differences between source and target label sets make approaches by direct distribution alignment less effective. For partial DA, Cao *et al.* [6] align class-wise distributions using multiple domain discriminators, while Zhang *et al.* [7] adopt an auxiliary domain classifier to estimate source-domain sample importance. In open-set DA, where source and target domains can hold private label sets, a common practice is to use an "unknown" class gathering all target domain's private objects [1, 19]. In [2], the authors introduce a more extreme setting coined "universal" DA, which imposes no prior knowledge on the target label set.

*Dealing with news classes.* Zero-Shot Learning (ZSL) aims to recognize new classes based solely on their semantic associations with previously seen classes. To this end, attributes [20], description [21] or even word embedding [22] have been shown to be shared semantic representations that allow transferring knowledge from seen to unseen classes. In this work, we use word2vec [23] as it is extracted at minimal cost and it alleviates the need to define and assign hundreds of attributes.

ZSL for image classification has been actively studied in the literature [24, 25, 26, 27, 28, 29, 30, 31, 32] and existing methods can be generally categorized into the following two groups. The first group [24, 25, 28, 29, 30, 32] addresses ZSL as an embedding problem and maps image data and class descriptions into a common space where semantic similarity translates into spacial proximity. The second group [26, 27, 31] of methods learns from visual and semantic features of seen classes and produces generators that can synthesize visual features based on class semantic descriptions. The synthetic features are then used to train a standard classifier for object recognition. This second type of methods has been shown to be more effective than embedding approaches as it reduces the inherent bias toward seen classes. In this work, we leverage this type of generative techniques.

More and more ZSL works [14, 33, 34] address the generalized setting (GZSL) where both seen and unseen classes are jointly evaluated at test time. From a practical point of view, GZSL's evaluation protocol better reflects the real-life perfor-

mance of zero-shot models.

Very recently, generalized zero-shot learning has been extended to semantic segmentation [35, 36, 37]. Bucher *et al.* [35] introduce the task by combining a deep visual segmentation model with an approach to generate visual representations from semantic word embeddings. Kato *et al.* [36] propose a variational mapping of the class labels' embedding vectors to the visual space. Xian *et al.* [37] introduce a model with a visual-semantic embedding module to encode images in the word embedding space and a semantic projection layer that produces class probabilities. In this work, we evaluate our models using the GZSL evaluation protocol for semantic segmentation proposed in [35].

*Transductive zero-shot learning.* Transductive ZSL solves ZSL in a semi-supervised learning manner, *i.e.*, all data including test ones are available at train time and, as with classic ZSL, there is no domain gap between base-class labelled samples and new-class unlabelled test ones. The latter can be utilized to form clearer decision boundaries for both base and new classes. Rohrbach *et al.* [38] explore the manifold structure of new classes by graph-based label propagation. Fu *et al.* [39] extend the label propagation with a multi-view hyper-graph. Several approaches adopt a joint learning framework to train on labeled and unlabeled data separately [40, 41, 42]. Such a training can be in the semantic space [40], the visual space [41] or a latent space [42]. Other efforts attempt to refine visual-semantic embeddings iteratively with unlabeled unseen data [43]. A domain-invariant projection is learnt in [44], which maps visual features to semantic embeddings and then reconstructs back the same visual features. Recently, [45] described a transductive unbiased embedding to improve generalized ZSL performance.

*Positioning.* In this paper, we consider for semantic segmentation the "open UDA" setting described in [8] for classification. In this setting, the system is expected to handle explicitly objects from new classes in the target domain. Here, the target label set is thus unbounded. In Figure 2, we illustrate this boundless UDA setting alongside other UDA settings. Opposite to the open-set setting in [19], the source label set in BUDA is a subset

of the target label set. BUDA is different from both open-set and universal DA as it explicitly requires prediction for all labels at test time. Also, as briefly mentioned in Section 1 and later elaborated in Section 3, the natural co-occurrence of objects from both shared and private classes in train and test target-domain images makes BUDA more challenging.

The BUDA framework differs from classic ZSL since training images with new classes are available but are unlabeled, whereas they would not be available at all during training in a pure ZSL setting. Our scenario is also different from transductive learning as: (1) training and test sets of target-domain samples are disjoint and (2) test images stem from a domain that differs from the labeled training domain, opening the door to the domain-shift problem.

BUDA raises jointly zero-shot problems (handling private classes) and domain adaptation ones (mitigating domain shifts); traditional ZSL and DA approaches are thus insufficient here. One the one hand, a direct application of zero-shot techniques fails to close the distribution gaps between source and target shared classes. As one must leverage shared-class features and semantic relations between shared and private classes to "generate" private-class features, source/target misalignment of the shared-class features results in an erroneous mapping to private-class features in the target domain. On the other hand, a direct distribution alignment across domains should be avoided to not undesirably mix-up features of shared and private classes. These two points are discussed in Sections 3.3 and 3.4.

## 3. Proposed framework: BudaNet

We first formulate the new task of semantic segmentation with boundless UDA. Given a set $C_s$ of class labels in the source domain, a training set $\mathcal{D}^s$ is available in this domain. It is composed of pairs $(x^s, y^s)$ formed by a color image $x^s \in \mathbb{R}^{H \times W \times 3}$ of size $H \times W$ and the associated $C_s-$class ground-truth one-hot segmentation map $y^s \in \{0, 1\}^{H \times W \times C_s}$, where $C_s = |C_s|$. In addition, a set $\mathcal{D}^t$ of unlabelled color images from the target domain is also available at train time.

Aiming for a more practical DA, we argue that, given a mini-

mal semantic prior, the final system should be capable of explicitly predicting new classes, *i.e*, which are not part of $C_s$, hence with no instances visible in the annotated source-domain images. In search of such minimal and practically accessible prior, we opted for the *names* of these new classes. These names are assumed to be known beforehand, certainly something effortless to achieve. In order to exploit the semantic relations between the classes of interest, we adopt the 'word2vec' model [23] learned on a dump of the Wikipedia corpus.

In BUDA, the set $C_s$ of classes in the source domain is fully shared with the target domain – the source domain does not hold privates classes –, while the latter hold a private set of new classes we denote $C_p$ and $C_p$ its size. At run-time, we want to classify each pixel of target-domain images as one of the categories in $C_s \cup C_p$. In BUDA, we know beforehand all these labels.

### 3.1. Introduction to our strategy

In this work, we propose a multi-step framework, coined as BudaNet, with dedicated strategies to address the aforementioned concerns. We start from an existing semantic segmentation model trained with a supervised loss on source-domain data and an unsupervised loss on target-domain train data (Step 1 in Figure 3). The unsupervised criterion mitigates the source-target distribution gap over the shared classes without causing unwanted misalignment to the private ones.

Once trained, this first semantic segmentation model is limited to trained categories and, hence, unable to recognize new classes. To allow the model to recognize both shared and private categories, we propose to generate synthetic training data for private classes. This is obtained with a generative model conditioned on the semantic representation of these classes. The generator's outputs attempt to mimic the visual features computed by the current semantic segmentation model (Step 2 – blue part in Figure 3). As no existing UDA technique guarantees perfect and universal alignment between domains, we decided to take domain information into account when training the generator. This is done by conditioning not only on the class embedding but also on the domain. In addition, we introduce an additional

adversarial discriminator that tries to distinguish source-target generated features (Step 2 – red part in Figure 3).

Once the generator is trained, many synthetic features can be produced for private classes, and combined with real samples from shared classes. This new set of training data is used to retrain the classifier of the segmentation network so that it can now handle both shared and private classes (Step 3 – blue part in Figure 3). The classifier is used to extract pseudo-labels on target-domain training set. The whole segmentation network is then fine-tuned again with these pseudo-labels (Step 3 – red part in Figure 3).

Section 3.2 overviews the base zero-shot pipeline which helps to handle never-seen classes in semantic segmentation. Section 3.3 introduces our UDA strategy to align shared classes across domains with minimal negative alignment effects on private-class features. In Section 3.4, we present the novel domain-aware generative model to synthesize private-class features for the target domain. Lastly, Section 3.5 details the final self-training procedure for the classifier.

### 3.2. Base zero-shot pipeline

We now revisit the pure adaptation-free zero-shot semantic segmentation (ZS3) pipeline from [35], trained only on source-domain images in our context. This pipeline is built on top of an existing semantic segmentation network $F$, *i.e.*, DeepLabv3+ [46]. To facilitate explanations, we decompose $F$ into two consecutive parts: $F_{feat}$ as the feature extractor, which outputs pixel-wise features $f \in \mathbb{R}^{d_f}$, and $F_{cls}$ as the final $1 \times 1$ convolutional classification layer which associates class-scores to each pixel-wise feature (its output size thus depends on the number of considered classes).

The base zero-shot pipeline consists of three steps:

1. Training of a segmentation network $F$ using source-domain supervision on shared classes. Once trained, the feature extractor $F_{feat}$ can extract visual features from shared-class objects in source-domain images.

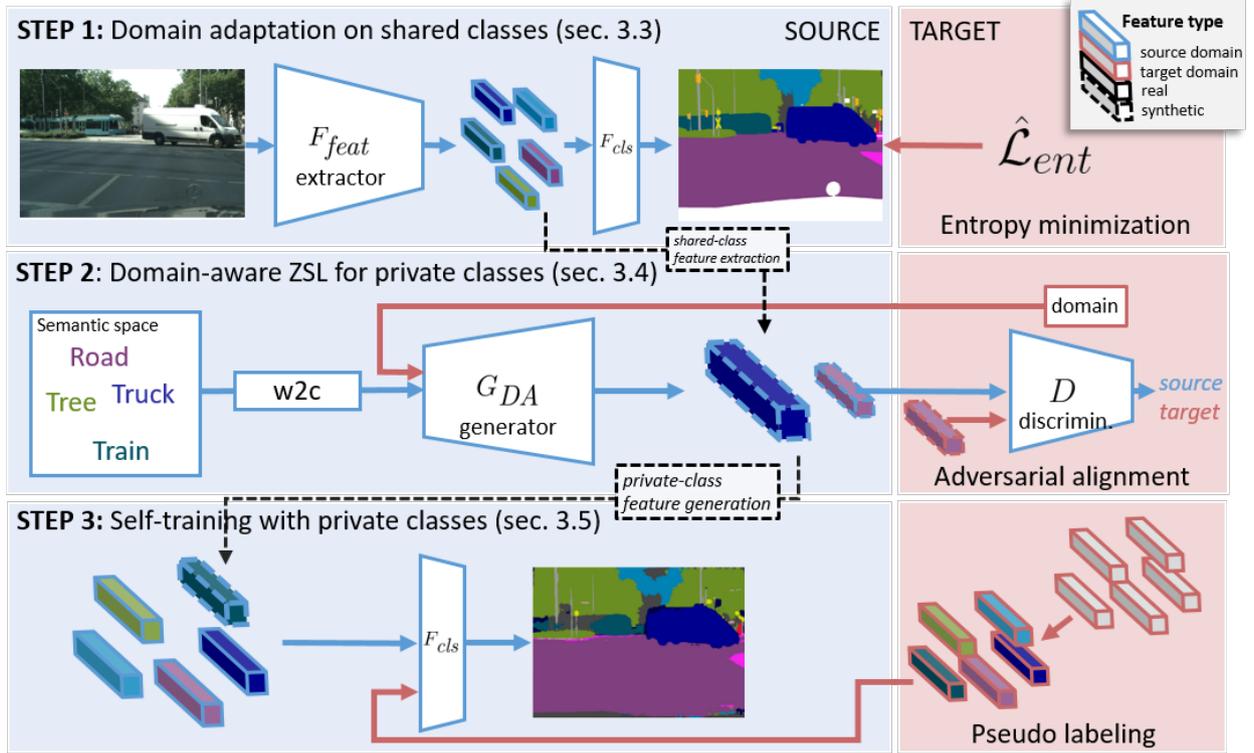2. Training of a generative network $G$, conditioned on shared-class embeddings, to generate corresponding source-

Fig. 3: **BudaNet's architecture and learning framework.** Our BudaNet learning strategy consists of three steps: Train the semantic segmentation network with domain alignment only for shared classes; Extract image features and use them as supervision for domain-aware generative model training; Combine the generated features with real ones to fine-tune the classification layer. The red-background pane (right) corresponds to the three proposed strategies for domain alignment, zero-shot learning on target domain and self-training with private classes. Black-dash arrows indicate connections between the steps. Colored arrows distinguish source-domain from target-domain flows.

domain visual features. The ground-truths are the shared-class features extracted after Step 1. The high-level idea is that geometric relations among classes in the embedding space are transferred to the generated feature space, which will help $G$ to handle private classes.

3. Training of the last classification layer $F_{cls}$ of $F$, now for both shared and private classes, with private-class features (generated by $G$ after Step 2) and shared-class features (computed by $F_{feat}$ after Step 1).

The final model is composed of $F_{feat}$ from Step 1 and $F_{cls}$ from Step 3. In detail, the segmentation network $F$ and the classifier $F_{cls}$ are trained using a standard segmentation cross-entropy loss $\mathcal{L}_{seg}$. Given a source sample $(x^s, y^s)$ and $P^{x^s}$ its soft-segmentation map predicted by $F(x^s)$, this loss reads:

$$\mathcal{L}_{seg}(x^s, y^s) = -\sum_{h=1}^{H}\sum_{w=1}^{W}\sum_{c \in C} y^s[h, w, c] \log P^{x^s}[h, w, c]. \quad (1)$$

For the training of $F$ in Step 1 (resp. in Step 3), we set $C$ to $C_s$ (resp. $C_s \cup C_p$). We follow [35] and adopt a generative moment-matching network (GMMN) as $G$, trained with a maximum mean discrepancy loss $\mathcal{L}_{GMMN}$. Given a random sample $\mathbf{z} \in \mathbb{R}^{d_z}$ from a fixed multivariate Gaussian distribution, the semantic embedding $\mathbf{a} \in \mathbb{R}^{d_a}$ and the domain indicator $d$ (set as 1 for source and 0 for target), new pixel-level feature sets are generated as $\hat{f} = G(\mathbf{a}, d, \mathbf{z}; \theta_G)$, where $\theta_G$ are the parameters of $G$. With two random populations $\mathcal{F}(\mathbf{a}, d)$ of real features and $\widehat{\mathcal{F}}(\mathbf{a}, d; \theta_G)$ of generated ones, we have:

$$\mathcal{L}_{GMMN}(\mathbf{a}, d) = \sum_{f, f' \in \mathcal{F}(\mathbf{a},d)} k(f, f') + \sum_{\hat{f}, \hat{f}' \in \widehat{\mathcal{F}}(\mathbf{a},d;\theta_G)} k(\hat{f}, \hat{f}')$$
$$- 2 \sum_{f \in \mathcal{F}(\mathbf{a},d)} \sum_{\hat{f} \in \widehat{\mathcal{F}}(\mathbf{a},d;\theta_G)} k(f, \hat{f}), \quad (2)$$

where $k$ is the Gaussian kernel, $k(f, f') = \exp(-\frac{1}{2\sigma^2}\|f - f'\|^2)$ with bandwidth parameter $\sigma$.

Figure 3, to the left, illustrates the ZS3 pipeline. Our proposed

BudaNet is built upon such a base three-steps paradigm.

### 3.3. BudaNet's Step 1: Domain adaptation on shared classes

A pure zero-shot model like ZS3 is able to handle simultaneously base and new classes from the same domain but fails to address the domain gap exhibited in BUDA. One straightforward solution for BUDA is to use UDA techniques while training $F$ in Step 1, with the hope that $F_{feat}$ can extract domain-invariant features. Additional unlabeled target-domain images are used for this purpose.

To this end, we study two recent UDA techniques in segmentation, MinEnt (self-training) and AdvEnt (adversarial training), both introduced in [5]. With MinEnt, the unsupervised entropy loss $\mathcal{L}_{ent}$, applied on target-domain samples, is jointly optimized with the supervised segmentation loss $\mathcal{L}_{seg}$ on source-domain samples. This loss is the sum of pixel-wise predictive entropies in target domain. Given a target sample $x^t$ with predicted soft-segmentation map $P^{x^t}$, the entropy loss is:

$$\mathcal{L}_{ent}(x^t) = \frac{-1}{\log(C_s)} \sum_{h=1}^{H} \sum_{w=1}^{W} \sum_{c \in C_s} P^{x^t}[h, w, c] \log P^{x^t}[h, w, c]. \quad (3)$$

Differently, AdvEnt approaches global distribution alignment via adversarial training on the weighted self-information space. We refer readers to the original work in [5] for more details.

As mentioned earlier, the mismatch between source and target label sets prevents the direct alignment between the visual distributions in source and target domains. Figure 4-(a) illustrates the problem caused by brute-force adoption of existing UDA techniques in the presence of private classes. Indeed by design, AdvEnt, using global output-space alignment, and MinEnt, using global entropy aggregation, do not differentiate between shared and private classes, which eventually results in undesirable "clusters" of mixed shared-class and private-class features.

We propose a simple yet effective strategy based on MinEnt to mitigate the above concern. A segmentation network $F^{pre}$ is first pre-trained only on the shared classes in the source domain. We then use $F^{pre}$ to produce shared-class predictions on the target-domain training images, which may include both shared-class and private-class objects. We argue that, as $F^{pre}$ has never observed private classes in the pre-training phase, top-confident

predictions coming from $F^{pre}$ should mostly constitute of shared-class pixels. Effectively, applying entropy minimization solely on such top-confident pixels minimizes the risk of misalignment in the presence of private-class target-domain features. We then fine-tune $F$ initialized with $F^{pre}$, using the optimization objective:

$$\min_{\theta_F} \frac{1}{|\mathcal{D}^s|} \sum_{x^s \in \mathcal{D}^s} \mathcal{L}_{seg}(x^s, y^s) + \frac{\lambda_{ent}}{|\mathcal{D}^t|} \sum_{x^t \in \mathcal{D}^t} \hat{\mathcal{L}}_{ent}(x^t), \quad (4)$$

where $\theta_F$ are the parameters of $F$ and $\hat{\mathcal{L}}_{ent}$ is derived from $\mathcal{L}_{ent}$ by restricting the sum of the predictive entropies to the top-$k\%$ most-confident pixels according to $F^{pre}$. The parameter $\lambda_{ent}$ controls the weight of this entropy term. Figure 3 illustrates Step 1 operating on both source and target domains.

### 3.4. BudaNet's Step 2: Domain-aware ZSL for private classes

The domain shift affects all classes, both shared and private ones. Although for BUDA we assume the absence of private classes in the source domain, if there still appeared a private-class instance in a source-domain scene, its appearance should be notably different from target-domain instances. For example, similar to 'car' images, 'tuk-tuk' images if taken in Paris should look different from those taken in India, due to distinct weather, illumination and overall appearance of each city. We thus find it crucial to take domain information into account in Step 2 where the goal is to synthesize target-domain private-class features. One may argue that after Step 1, the two domains are already aligned and the domain-invariant feature extractor $F_{feat}$ should be ready as is for feature synthesis. However, no existing UDA technique guarantees complete alignment across domains, which we thus do not expect either in our approach.

We propose a novel domain-aware generative model, denoted $G_{DA}$, which is trained to synthesize visual features for both source and target domains. This generative model is modified such that it is conditioned not only on the class embedding but also on the domain (source or target). Now, the modified model can output shared-class features for both source and target domains. We note that, although our objective is to generate domain-aware features, those stemming from the same classes should be rather close. Intuitively, we want our generator to
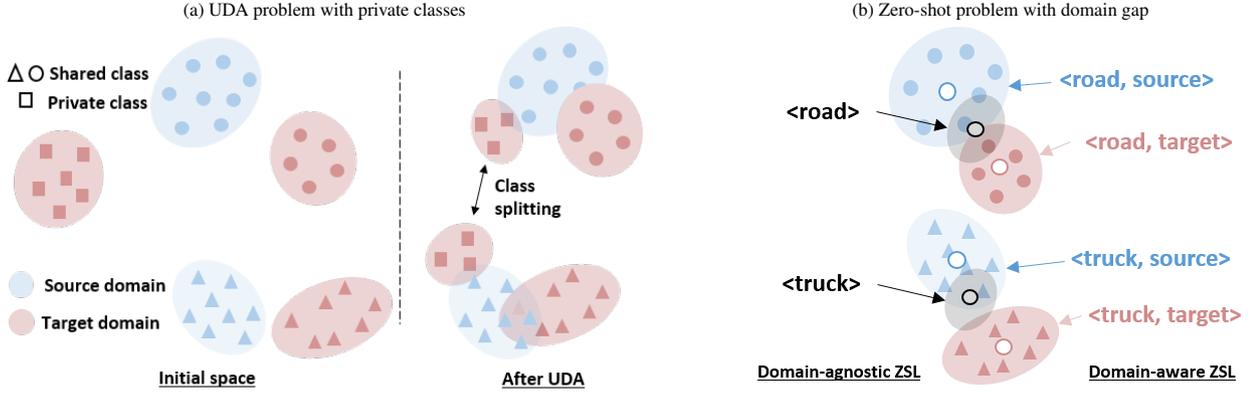
Fig. 4: **Challenges of boundless UDA.** (a) Existing UDA techniques, ignoring the presence of target private classes, tend to shift target-domain features to shared-class clusters. (b) Because of the domain gap, learning domain-agnostic mappings (black arrows) from word embedding space to visual space is sub-optimal as only mean modes (black circles) are captured. Instead, our domain-aware generative model can map to domain-specific modes (blue or red circles).

mimic well the feature generator $F_{feat}$, which produces close source-domain and target-domain feature distributions (as a result of the adaptation Step 1) yet still exhibiting certain discrepancies (as discussed above).

To this end, we leverage adversarial training. In details, we introduce an additional discriminator $D$ trying to distinguish source-target generated features. At train time, $D$ minimizes the binary cross-entropy classification loss; Given $\hat{f}^s$ and $\hat{f}^t$, two generated features for source (label 1) and target (label 0) domains, this loss reads:

$$\mathcal{L}_D(\hat{f}^s, \hat{f}^t) = \mathcal{L}_{BCE}(\hat{f}^s, 1) + \mathcal{L}_{BCE}(\hat{f}^t, 0), \qquad (5)$$

with $\mathcal{L}_{BCE}$ being the binary cross-entropy loss. Meanwhile, the generative model trained with an additional adversarial loss $\mathcal{L}_{adv}$ tries to confuse $D$. Training is now based on the two losses: $\mathcal{L}_{GMMN}$ as in Section 3.2 and the adversarial loss $\mathcal{L}_{adv}$. Given $\hat{f}^t$, a generated feature for target domain, this loss reads:

$$\mathcal{L}_{adv}(\hat{f}^t) = \mathcal{L}_{BCE}(\hat{f}^t, 1). \qquad (6)$$

To train $G_{DA}$, we use a weighting factor $\lambda_{adv}$ for $\mathcal{L}_{adv}$.

Generative training is supervised by real shared-class features coming from both source and target domains. While in the source domain we can easily assign class labels to shared-class features using the ground-truth maps, in the target domain we must opt for a heuristic pseudo-labeling strategy. Specifically, we run $F$ (pre-trained in Step 1) on the unlabeled target-domain

training set. For each pixel, we assign the class with the highest prediction probability as its pseudo-label. However, this labelling is error prone: it is not guaranteed to be correct on shared-class pixels, and cannot be correct on private-class pixels whose true labels are unknown to it. To mitigate such negative effects, only the top $p\%$ of the most confident pseudo-labels are kept; the rest are ignored during training. Doing that, we improve pseudo-label quality for shared classes and reduce the number of retained private-class pixels. In Figure 3, the strategy is illustrated as Step 2.

In Figure 4-(b), we illustrate different outcomes of our proposed model as opposed to one without domain awareness. As $G_{DA}$ can capture domain-specific modes in the visual-feature space, we expect better generalization to private classes in the target domain.

### 3.5. BudaNet's Step 3: Self-training with private classes

We start by only pre-training the final classification layer $F_{cls}$ with the shared-class source-domain features used in Step 2 and the generated private-class features coming from $G_{DA}$ – see Figure 5 (left). Once $F_{cls}$ is trained, the network can now handle both shared and private objects at once, which effectively can be used to extract pseudo-labels on target-domain training set. The whole segmentation network $F$ is then fine-tuned again with pseudo-labels – similar to above, we only train on top confident ones. Differently to the previous steps, the pseudo-labels now

cover both shared and private classes. Figures 5 and 3-(Step 3) illustrate this final self-training strategy. In Figure 5-(right), the decision boundaries, learned after pre-training $F_{cls}$ (blue lines) are shifted to new positions that are better adapted to the target domain after fine-tuning $F$ (black lines).

## 4. Experiments

In Section 4.1, we introduce our set-ups as well as some implementation details. Section 4.2 presents the main results, followed by ablation studies in Section 4.3.

### 4.1. Experimental details
*BUDA scenarios.* To evaluate our approach, we define three BUDA scenarios: *synthetic-2-real*, *country-2-country* and *dataset-2-dataset*. These very different set-ups allow investigating boundless adaptation from multiple angles of practical interest.

In the *synthetic-2-real* set-up, available at train time are labelled synthetic data and unlabeled real images. The zero-cost (and zero-risk) source-domain ground-truth acquisition makes this configuration especially appealing for semantic segmentation task. We use the SYNTHIA-RAND-CITYSCAPES split of the SYNTHIA dataset [47] containing 9,400 synthesized images as for source domain. The target-domain images are from the Cityscapes dataset [48] with 2,975 images for training and 500 images for test. In this set-up, we consider 16 shared classes, previously used in some closed-set UDA works [3, 4, 5]. Differently, our target domain holds a private label-set of 3 classes: 'terrain', 'truck' and 'train'.

The *country-2-country* set-up addresses the very practical use-case where a model trained using data from one region is deployed in other parts of the world. In details, the source domain is defined by the Cityscapes data acquired in 50 German cities while, for the target domain, we use images from the India driving dataset (IDD) [49]. The geographic gap translates into a large visual gap on both the vehicles and the overall scenes' layout. While common vehicles like cars, trucks or buses exist almost everywhere in the two countries, only in India we observe auto-rickshaws. It is also worth mentioning

that class distributions are very different, *i.e.* in Indian streets, the most frequent vehicles are motorbikes while in Germany, cars are more prevalent. Such unique signatures of different countries can be challenging for applications like self-driving cars. More into details, we use for training 2,975 real Cityscapes images with annotations along with 6,993 unlabeled IDD images. There are 19 shared classes between the two datasets and 2 private classes only in IDD: 'auto-rickshaw' (a.k.a. 'tuk-tuk') and 'animal'. We evaluate models on 981 IDD images for the total 21 classes.

Our last BUDA set-up is *dataset-2-dataset*: Pascal-VOC [50] and MS-COCO [51] as source and target domains. As both were collected from the Internet, one may expect that the two datasets are very similar. The big difference lies in the scene complexity and the level of annotation done in each: while Pascal-VOC images mostly have a single centered object with a limited 20-class label-set, MS-COCO images are more complex with multiple objects exhaustively annotated. Still, as the visual gap is small between the two datasets, we do not expect very significant performance drop in the shared classes. The challenge clearly appears once we consider a target's private label set. In detail, there are $1,464$ VOC images and 118k COCO images used in training. We consider, in the main experiment, 20 shared classes and 5 target private classes: 'truck', 'bench', 'zebra', 'giraffe' and 'laptop'. For validation, we use 5000 COCO images which contain at least one of the 25 classes. For the ablation study in Section 4.3, we also increased in this set-up the number of private classes from 5 to 20, something the two considered datasets permit while respecting the BUDA assumption that private classes do not appear in the source domain.

*Model implementation.* We adopt the DeepLabv3+ framework [46] built upon a ResNet-101 [52] backbone. The SGD [53] optimizer is used with polynomial learning rate decay with the base learning rate of $10e-2$, weight decay $1e-4$ and momentum 0.9.

Generative models $G$ and $G_{DA}$ are multi-layer perceptrons having a single hidden layer with leaky-ReLU non-linearity [54] and dropout [55]. We fix the number of hidden neurons as 256. Both
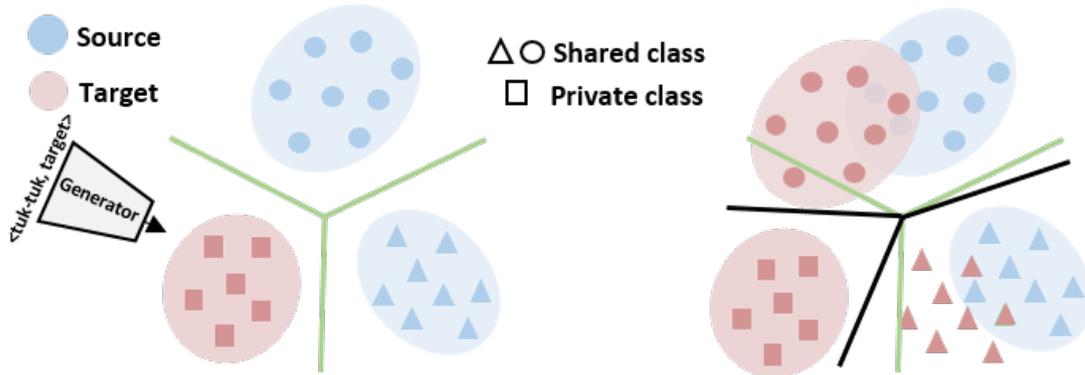
Fig. 5: **Self-training with private classes**. (Left) Pre-training $F_{cls}$ using aligned source-domain features of shared classes used in Step 2 and features generated for private classes, *e.g.*, 'tuk-tuk' indicated with squares. Solid lines illustrate the decision boundaries. (Right) Fine-tuning the whole segmenter $F$ using only pseudo-labeled target-class features – blue points are there to contrast with the left figure but are not used to train $F$.

generative models take as input a semantic classes' embedding of dimension $d_a = 300$ and a Gaussian noise vector of dimension $d_z = 300$. $G_{DA}$ additionally accepts a 1−dim vector specifying the domain (source or target). Regarding the generative loss $\mathcal{L}_{GMMN}$, we chose kernel bandwidths in $\{2, 5, 10, 20, 40, 60\}$.

For adversarial training of $G_{DA}$, we introduce an additional binary classification discriminator $D$: a fully connected layer takes 256−dim feature vectors and predicts the corresponding domain. All generative models are trained using Adam optimizer [56] with a learning rate of 2e−4.

*Evaluation metrics.* We want to assess performance of both shared and private classes at test time. Therefore, we adopt the GZSL evaluation protocol for semantic segmentation used in [35]. The protocol considers three traditional base metrics in segmentation: pixel accuracy (PA), mean accuracy (MA) and mean intersection-over-union (mIoU), separately computed for shared and private sets. Also reported are the harmonic means (hPA, hMA and hIoU) of shared and private results. The harmonic means are used because the common metrics suffer from the performance bias of shared classes, while our aim is to achieve high accuracy for both shared and private classes [57].

## 4.2. Results

Tables 1, 2 and 3 report our results in different setups.

In each experiment, we report performance for shared, private, and all classes. "Supervised" stands for the model trained with full supervision on the target domain. "Oracle" corresponds to

the zero-shot framework trained on target-domain data using only shared labels; "ZS3Net (only source)" is a similar network but only trained on source-domain data (based on ZS3Net in [35]). One straight-forward baseline is to directly apply a vanilla UDA technique, *i.e.* MinEnt [5], in Step 1 without considering the private classes. We denote "ZS3Net + UDA" this baseline. We note that in the new BUDA setting, existing UDA techniques compare differently than in the classic setting, as later reported in Section 4.3; Indeed "MinEnt" performs better than a state-of-the-art adversarial training approach. In addition, we consider the proposed strategy introduced in Section 3.3 as a stronger baseline, denoted "ZS3Net + Adaptation" in the result tables.

*Synthetic-to-Real.* Table 1 reports the segmentation performance of the models trained on SYNTHIA data and evaluated on the 19 classes of the Cityscapes validation set. Not surprisingly, the pure zero-shot segmentation method with no adaptation (ZS3Net) produces unfavourable results. The straight-forward baseline (ZS3Net + UDA) does introduce some improvement over ZS3Net. With the shared-class alignment strategy proposed in Section 3.3, our stronger baseline (ZS3Net + Adaptation) performs better on all classes. Finally, BudaNet outperforms all baselines by significant margins for both shared and private classes. Figure 6 illustrates the merit of our approach. On shared classes, ZS3Net produces noisy predictions, mistaking 'road' pixels with 'sidewalk', while BudaNet produces correct predictions for most pixels. On private classes, *i.e.* 'train' and

Table 1: **Semantic segmentation performance on SYNTHIA→Cityscapes.** Pixel accuracy (PA), mean accuracy (MA) and mean intersection-over-union (mIoU) on shared and private class sets, and their harmonic averages over the two class sets (hPA,hMA,hIoU). See text for method's discussion.

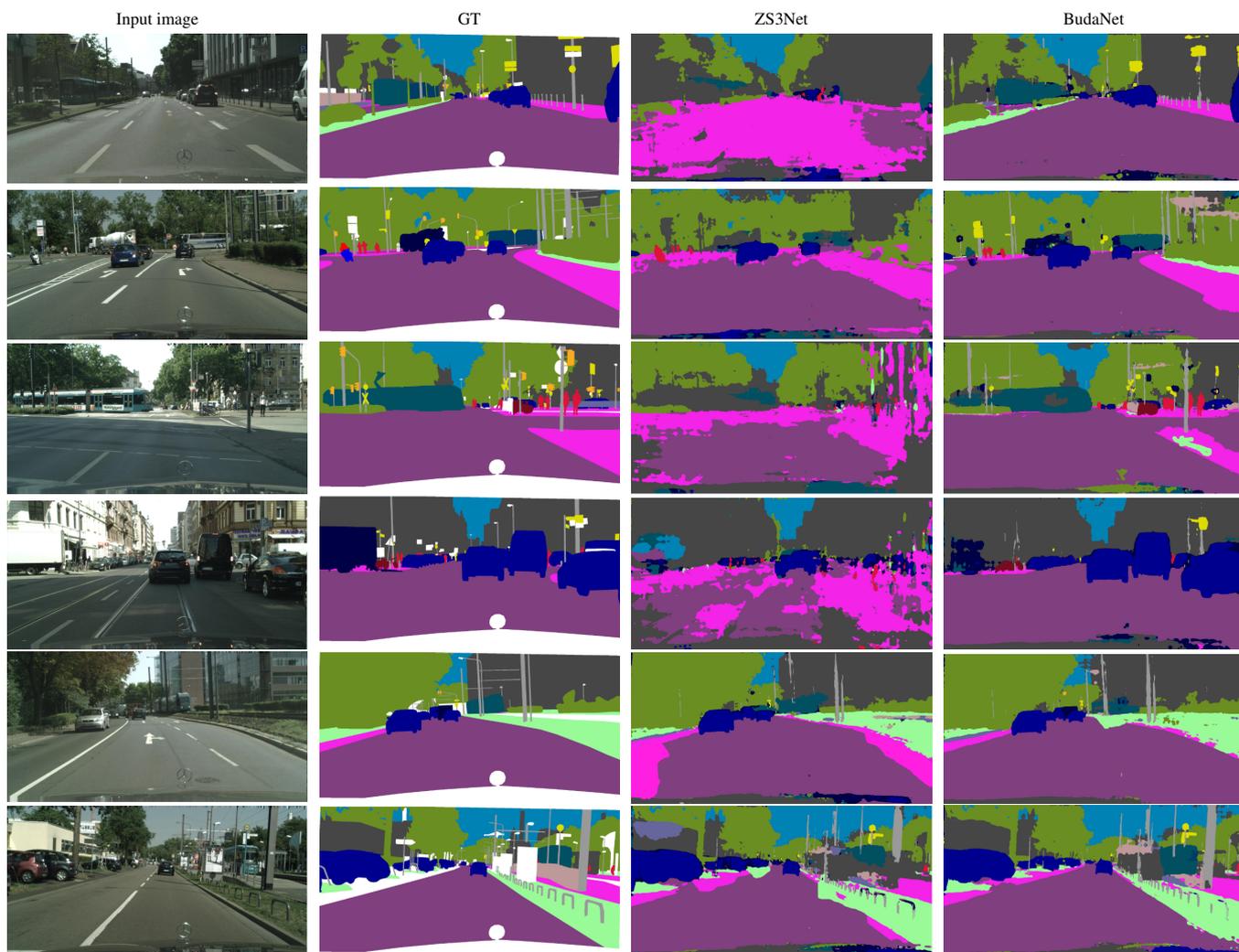| | Shared | | | Private | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | PA | MA | mIoU | PA | MA | mIoU | hPA | hMA | hIoU |
| Supervised | 95.1 | 68.1 | 61.7 | 58.1 | 60.0 | 48.2 | 72.1 | 63.8 | 54.1 |
| Oracle | 94.3 | 65.7 | 56.4 | 50.1 | 45.9 | 31.5 | 65.4 | 54.0 | 40.4 |
| ZS3Net (source only) [35] | 80.3 | 43.1 | 28.1 | 9.0 | 40.7 | 6.9 | 16.2 | 41.9 | 11.1 |
| ZS3Net + UDA | 85.3 | 42.4 | 30.1 | 12.5 | 46.8 | 8.2 | 21.8 | 44.9 | 12.8 |
| ZS3Net + Adaptation | 89.9 | 42.6 | 35.0 | 18.6 | 51.1 | 8.9 | 30.8 | 46.5 | 14.2 |
| BudaNet | **93.0** | **46.0** | **36.2** | **26.9** | **58.7** | **17.0** | **41.7** | **51.6** | **23.1** |



Fig. 6: **Qualitative results on Cityscapes**. The first and second columns show input images and corresponding segmentation ground truth. The third and fourth columns visualize results produced by ZS3Net and BudaNet. Private target classes: terrain, truck, train. Some shared classes: road, side walk, car, person, motorbike, tree, building.

Table 2: **Semantic segmentation performance on Cityscapes→IDD.** Pixel accuracy (PA), mean accuracy (MA) and mean intersection-over-union (mIoU) on shared and private class sets, and their harmonic averages over the two class sets (hPA,hMA,hIoU). See text for method's discussion.

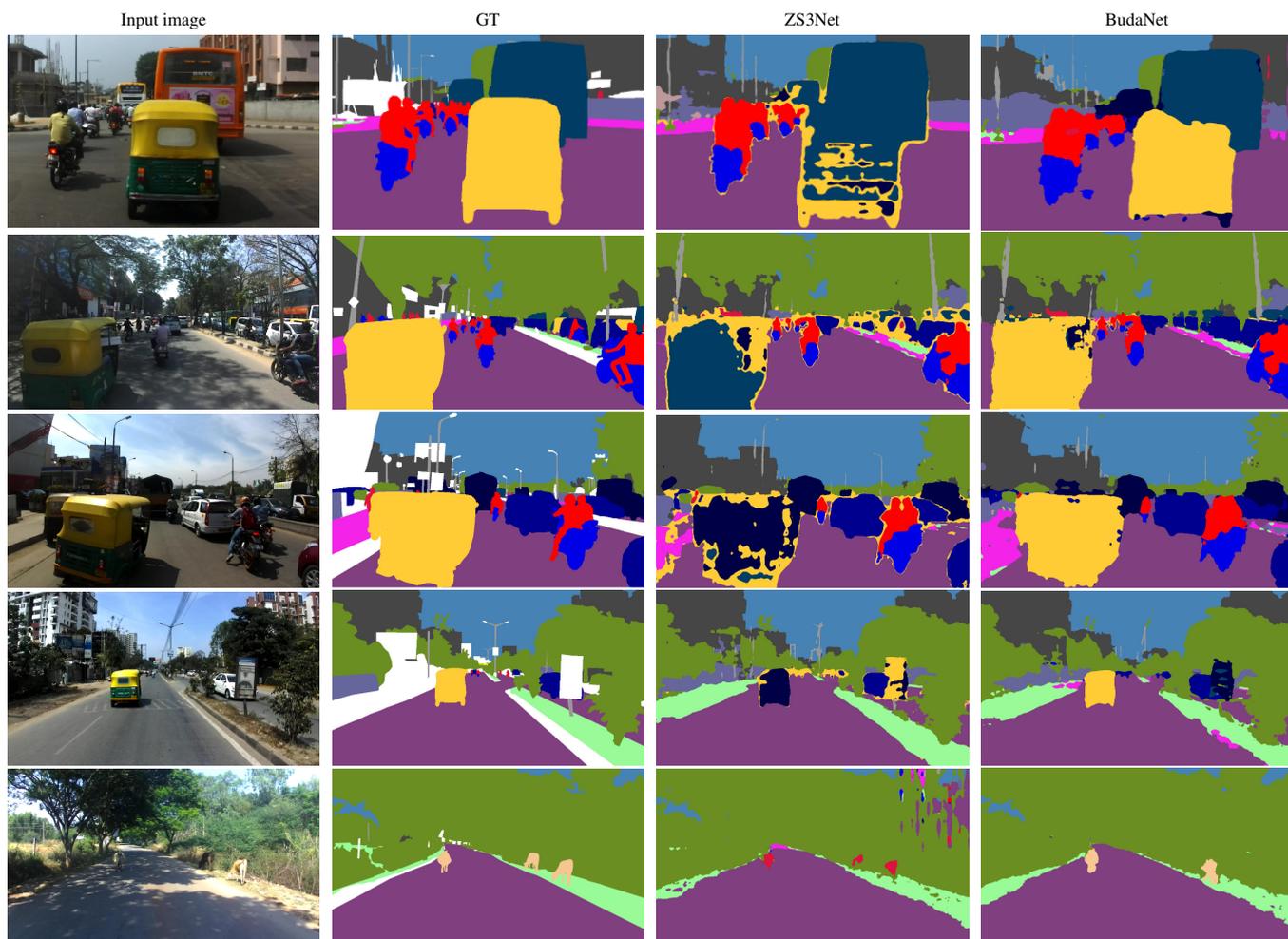| | Shared | | | Private | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|
| Method | PA | MA | mIoU | PA | MA | mIoU | hPA | hMA | hIoU |
| Supervised | 92.8 | 65.1 | 56.9 | 57.3 | 61.9 | 48.0 | 70.9 | 63.5 | 52.1 |
| Oracle | 91.7 | 63.0 | 53.1 | 47.2 | 41.8 | 28.8 | 62.3 | 50.3 | 37.3 |
| ZS3Net (source only) [35] | 81.0 | 43.8 | 29.2 | 9.3 | 41.0 | 7.9 | 16.7 | 42.4 | 12.4 |
| ZS3Net + UDA | 86.8 | 40.0 | 32.4 | 13.8 | 45.9 | 8.1 | 23.8 | 42.7 | 13.0 |
| ZS3Net + Adaptation | 88.3 | 40.5 | 32.7 | 15.9 | 47.0 | 8.6 | 26.9 | 43.5 | 13.6 |
| BudaNet | **92.0** | **47.2** | **37.3** | **28.6** | **58.9** | **18.5** | **43.6** | **52.4** | **24.7** |



Fig. 7: **Qualitative results on IDD**. The first and second columns show input images and corresponding segmentation ground truth. The third and fourth columns visualize results produced by ZS3Net and BudaNet. Private target classes: tuk-tuk, animal. Some shared classes: truck, road, side walk, car, person, motorbike, tree, building.

'truck', while ZS3Net can only localize small regions, BudaNet provides more complete areas with better contours. ZS3Net produces noisy segmentation with, for example, the shared class 'person' missing in rows 2 and 3. In the same way, the areas of

Table 3: **Semantic segmentation performance on Pascal-VOC→MS-COCO.** Pixel accuracy (PA), mean accuracy (MA) and mean intersection-over-union (mIoU) on shared and private class sets, and their harmonic averages over the two class sets (hPA,hMA,hIoU). See text for method's discussion.

| Method | Shared | | | Private | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|
| | PA | MA | mIoU | PA | MA | mIoU | hPA | hMA | hIoU |
| Supervised | 94.8 | 68.8 | 70.1 | 70.5 | 61.0 | 64.3 | 80.9 | 64.7 | 67.1 |
| Oracle | 93.8 | 68.6 | 68.0 | 61.6 | 57.7 | 39.9 | 74.4 | 62.8 | 50.3 |
| ZS3Net (source only) [35] | 92.0 | 67.3 | 63.3 | 29.9 | 51.4 | 17.3 | 45.1 | 58.3 | 27.2 |
| ZS3Net + UDA | 92.3 | 68.0 | 63.3 | 32.3 | 50.3 | 19.5 | 52.0 | 47.9 | 29.8 |
| ZS3Net + Adaptation | 92.3 | 68.2 | 63.8 | 36.2 | 54.1 | 20.9 | 52.0 | 60.3 | 31.5 |
| BudaNet | **93.1** | **68.4** | **65.0** | **38.5** | **56.5** | **23.8** | **54.5** | **61.9** | **34.8** |

the private class 'truck' are hardly detected and wrongly localized elsewhere in rows 1 and 3. Our BudaNet provides better predictions on these areas.

*Country-to-Country.* Table 2 reports the results on the 21 classes of the IDD validation set. We observe similar behaviors for the baselines and for BudaNet. While ZS3Net trained only on the source domain produces poor results, ZS3Net + UDA baseline performs better thanks to the straight-forward domain alignment. Using our adaptation strategy in Step 1, the ZS3 + Adaptation baseline enhances recognition performances on all classes. BudaNet introduces gains in both shared and private classes with a significant improvement of +11.1% hIoU in comparison to ZS3 + Adaptation baseline. Figure 7 provides some segmentation examples. Qualitatively, BudaNet's results look much better. Both ZS3Net and BudaNet produce reasonable segmentation for the shared classes. However, in rows 1 to 5, ZS3Net hallucinates 'auto rickshaw' parts on all vehicles. Further, most of 'auto rickshaw' pixels are partially predicted or incorrectly classified as 'truck' (row 3) or 'bus' (row 1, 2, 3 and 4). The last row shows the second private class 'animal'. While ZS3Net completely misses the animals on the road, BudaNet correctly predicts two of them.

*Dataset-to-Dataset.* Table 3 reports the semantic segmentation performances on the MS-COCO validation set. We only report results for the 20 shared + 5 private classes and do not consider

others. Unlike the two other settings, Pascal-VOC and MS-COCO exhibit a small domain gap since they were both collected from the Internet. Effectively with ZS3Net (no adaptation), the mIoU drop on shared classes is only 4.7% compared to the Oracle. Still, we demonstrate here similar improvements after addressing adaptation and zero-shot challenges. Private-class segmentation results remain much lower due to the absence of training samples for these categories.

Domain alignment with the proposed strategy in Step 1 (ZS3Net + Adaptation) brings a +3.6% mIoU improvement on private classes; the scores on shared classes remain comparable (63.3% *vs.* 63.8% mIoU). BudaNet helps boosting the performance in all classes with a hIoU increase of +7.6% in comparison to ZS3Net, most of which accounts to private classes. Figure 8 shows some outputs of ZS3Net and BudaNet. ZS3Net wrongly classifies private objects' parts as background or shared classes. BudaNet produces better predictions with accurate object contours. In row 1, while ZS3Net wrongly segments most pixels of the private class 'bench' as 'chair', BudaNet can recognize more complete objects. In the 2nd and 3rd examples, ZS3Net confuses the pixels of private 'giraffe' and 'zebra' classes as 'background'. By contrast, BudaNet delivers good predictions for these classes. In the last row, BudaNet is shown to improve performance for the 'truck' class.

*Transductive domain adaptation.* While the focus of the present work is on inductive learning (test data is unknown at train time),

Fig. 8: **Qualitative results on MS-COCO**. The first and second columns show input images and corresponding segmentation ground truth. The third and fourth columns visualize results produced by ZS3Net and BudaNet. Private target classes: bench, giraffe, zebra, truck. Shared classes: person, cow, sheep, background.

transductive learning could also be considered with the proposed tools. In that very specific setup, there is no distinction between train and test data. BudaNet could thus be further refined using pseudo-labeling of the target-domain test set in Step 3. This actually boosts the performance on all classes, with a large gain of 24.9% in harmonic IoU on MS-COCO. Note however that such a transductive scenario is extremely contrived and is not suited for systems to be deployed in open environments.

### 4.3. Ablation studies

*Effect of the proposed strategies.* We report in Table 4 an ablative study that analyses the impact of four key ingredients of BudaNet: Entropy minimization while training the segmenter $F$ ("MinEnt"); Entropy minimization on top-confident pixels only ("Domain Adaptation on shared classes" ); Generator training with adaptation ("Domain-aware ZSL on private classes" ); fine-tuning of the semantic segmentation model with pseudo-labels ("Self-training"). To this end, five different models are compared: A model trained only on source, hence without

Table 4: **Ablation experiments**. Performance in hIoU on the validation set of Cityscapes for five variants of the approach, adding one key component at a time on top of the adaptation-free zero-shot semantic segmentation. See text for details.

| Setup | MinEnt | Domain Adaptation on shared cls. | Domain-aware ZSL on private cls. | Self-training | hIoU (%) |
|---|:---:|:---:|:---:|:---:|---:|
| ZS3Net | | | | | 11.1 |
| ZS3Net + UDA | ✓ | | | | 12.8 |
| BudaNet - Step 1 | ✓ | ✓ | | | 14.2 |
| BudaNet - Step 1+2 | ✓ | ✓ | ✓ | | 22.0 |
| BudaNet - Step 1+2+3 | ✓ | ✓ | ✓ | ✓ | **23.1** |

Table 5: **Analysis of the adaptation strategy on shared classes.** Target-domain test performance in mIoU on shared ("Shar.") and private ("Priv.") classes and in hIOU on both sets ("All") for the three BUDA settings, and for different adaptation strategies in Step 1. We compare the proposed stategies "ZS3Net + UDA" and "ZS3Net + Adaptation" to the ZS3Net variant adopting the adversarial DA technique AdvEnt [5]. Note that "ZS3Net + Adaptation" amounts to "BudaNet - Step 1" in Table 4.

| Method | Cityscapes | | | IDD | | | MS-COCO | | |
|---|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | Shar. | Priv. | All | Shar. | Priv. | All | Shar. | Priv. | All |
| ZS3Net + AdvEnt | 30.1 | 7.5 | 12.0 | 32.3 | 8.0 | 12.8 | 63.4 | 19.5 | 29.8 |
| ZS3Net + UDA | 30.0 | 8.2 | 12.8 | 32.4 | 8.1 | 13.0 | 63.3 | 19.5 | 29.8 |
| ZS3Net + Adaptation | **35.0** | **8.9** | **14.2** | **32.7** | **8.6** | **13.6** | **63.8** | **20.9** | **31.5** |

adaptation ("ZS3Net"); Our baseline with entropy minimization ("ZS3Net + UDA"); The full-fledged model ("BudaNet - Step 1+2+3"); Two stripped-down versions of it ("BudaNet - Step 1" and "BudaNet - Step 1+2"). We first observe that the brute force entropy minimization contributes to improving the performance by 1.7% ("ZS3Net" *vs.* "ZS3Net + UDA"), which demonstrates the need for domain alignment prior to generative model training. The table indicates that the pseudo-label strategy selection for shared-class pixels in the target domain ("ZS3Net + UDA" *vs.* "BudaNet - Step 1") significantly increases the performance by 1.4%. For the training of the generator, the combination of source-domain training data with pseudo-labeled target-domain data and the adversarial loss ("BudaNet - Step 1" *vs.* "BudaNet - Step 1+2") further improves the performance, from 14.2 to 22.0. Finally, the complete BudaNet benefits from self-training (up to 1.1%), as this step of pseudo-labeling makes it possible to exploit the information related to the private-class visual features.

*Number of private classes.* To analyse the impact of the number of private classes on the performance of BudaNet, we made it vary in the *dataset-2-dataset* setup. More precisely, we defined five nested sets with 2, 5, 10, 15 and 20 private classes respectively, among the 25 classes in Pascal-VOC and MS-COO. These sets are:

$C_p^2 = \{\text{'truck', 'zebra'}\},$

$C_p^5 = C_p^2 \cup \{\text{'bench', 'giraffe', 'laptop'}\},$

$C_p^{10} = C_p^5 \cup \{\text{'elephant', 'umbrella', 'bear', 'snowboard', 'toilet'}\},$

$C_p^{15} = C_p^{10} \cup \{\text{'laptop', 'fridge', 'blanket', 'napkin', 'stone'}\},$

$C_p^{20} = C_p^{15} \cup \{\text{'tent', 'skyscraper', 'salad', 'river', 'pillow'}\}.$

The corresponding sets of shared labels have size 23, 20, 15, 10 and 5 respectively. Note that $C_p^5$ is the private label set used in the main dataset-to-dataset experiments.
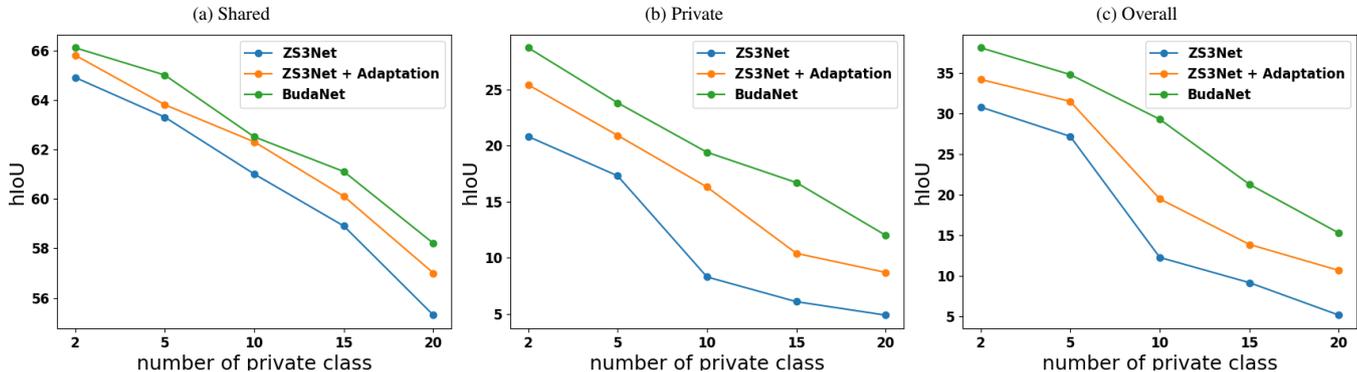
Fig. 9: **Influence of the number of target private classes**. Performance in mIoU on shared classes (a), mIoU on private classes (b) and hIoU on both (c) when training BudaNet on Pascal-VOC→MS-COCO with private class set $C_p^2, C_p^5, C_p^{10}, C_p^{15}$ and $C_p^{20}$, respectively. The number $C_s$ of shared classes is 23, 20, 15, 10 and 5 respectively.
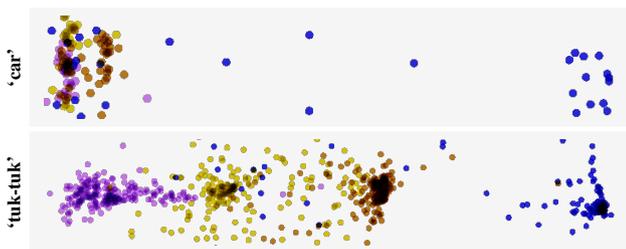


Fig. 10: **Domain shifts in visual-feature distributions for shared and private classes**. Visualization with tSNE, for 'car' (shared) and 'tuk-tuk' (private) classes, of the distributions of (from left to right) Real, BudaNet, ZS3Net + UDA and ZS3Net features. Must be viewed in color.

Figure 9 shows hIoU plots as a function of the number of private classes. We observe a decreasing trend in performance for all methods when more private classes (respectively fewer shared classes) are used. In all set-ups, BudaNet (green curves) outperforms the baselines by a significant margin.

*Domain adaptation strategy on shared classes.* In Table 5, we compare the proposed DA strategies on shared classes to the state-of-the-art adversarial DA technique AdvEnt [5]. We here notice the negative effects caused by global alignment techniques like AdvEnt in BUDA; such a drawback was similarly observed in partial DA and open-set DA works.

These experiments echo the arguments developed in Section 3.3 to justify the relevance of the proposed adaptation strategy.

*Feature visualization.* Using the models in Table 2, we visualize in Figure 10 domain shifts between real and generated pixel-wise features from different approaches for both shared and private classes. For the shared class 'car', as real features are used in training, we observe equally good alignment from BudaNet and "ZS3Net + UDA", as opposed to the adaptation-free ZS3Net. Regarding the private class 'auto-rickshaw', we clearly see a smaller gap between real and BudaNet-generated features, compared to the baseline and to the adaptation-free model. As no real private-class features are provided during training, we do not expect a "perfect" alignment in the latter case. The visualization is inline with the results in Table 2.

Table 6: **Impact of the fraction of high-scoring predictions retained as pseudo-labels**. Target-domain performance in hIoU on the SYNTHIA→Cityscapes setup as a function of $p(\%)$. "GT": see text for explanation.

| $p\%$ of high-scoring | 10% | 30% | 50% | 70% | 100% | GT |
|---|---|---|---|---|---|---|
| Cityscapes hIoU | 20.4 | 21.6 | **22.0** | 21.8 | 17.8 | 38.9 |

*Proportion of retained top-scoring predictions as pseudo-labels in Step 2.* Our aim here is to compare different $p\%$ values for the pseudo-labeling strategy on target features. Table 6 reports hIoU performance on Cityscapes as this parameter varies. Keeping only the highest-scoring predictions, at $p = 10\%$, results in a drop of performance of 1.6%. One possible explanation is that because the number of training features becomes too low, the model lacks generalization capacity. Conversely, when all predictions are kept ($p = 100\%$), a significant decrease of

performance can be observed, which confirms the need for an effective pseudo-label selection strategy. Finally, we observe that the performance is better when $p \in [30, 70]$, specifically when $p = 50\%$ with a recognition score of 22.0 hIoU [1]. In the last column of Table 6, "GT", we report the recognition score when the $p = 50\%$ highest-scoring predictions are replaced by their ground-truth labels. We observe a performance improvement of 16.9%, the pixel selection technique is therefore an essential element of the method and must be investigated in future work. In all other steps with pseudo-labeling, we also fix the percentage of retained high-scoring predictions as 50%.

## 5. Conclusion

In the context of semantic segmentation, we have explored boundless UDA, the realistic domain adaptation setting where new classes can emerge in the target test domain. Compared to existing DA settings, BUDA is more challenging because (i) objects from private new classes must be explicitly recognized at test time and (ii) objects from both shared and private classes frequently coexist in the test scenes. To address this new problem, we proposed BudaNet, a unified pipeline that leverages domain adaptation and zero-shot learning techniques. BudaNet consists of three steps with dedicated strategies to (i) mitigate the domain gap, (ii) handle the private classes and (iii) adjust the decision boundaries in the presence of new-class objects. On generalized ZSL evaluation metrics, BudaNet outperforms the baselines by significant margins and sets competitive standard in the BUDA setting. This first step towards boundless UDA will hopefully foster future research on this type of domain adaptation and, more generally, on new DA settings motivated by real-world applications.

---

[1]The result corresponds to "BudaNet - Step 1+2" in Table 4.

# References

[1] P. Panareda Busto, J. Gall, Open set domain adaptation, in: ICCV, 2017. 1, 3

[2] K. You, M. Long, Z. Cao, J. Wang, M. I. Jordan, Universal domain adaptation, in: CVPR, 2019. 1, 3

[3] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, Cycada: Cycle-consistent adversarial domain adaptation, in: ICML, 2018. 1, 3, 9

[4] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, M. Chandraker, Learning to adapt structured output space for semantic segmentation, in: CVPR, 2018. 1, 3, 9

[5] T.-H. Vu, H. Jain, M. Bucher, M. Cord, P. Pérez, Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation, in: CVPR, 2019. 1, 3, 7, 9, 10, 15, 16

[6] Z. Cao, M. Long, J. Wang, M. I. Jordan, Partial transfer learning with selective adversarial networks, in: CVPR, 2018. 1, 3

[7] J. Zhang, Z. Ding, W. Li, P. Ogunbona, Importance weighted adversarial nets for partial domain adaptation, in: CVPR, 2018. 1, 3

[8] J. Zhuo, S. Wang, S. Cui, Q. Huang, Unsupervised open domain recognition by semantic discrepancy minimization, in: CVPR, 2019. 2, 4

[9] M. Long, Y. Cao, J. Wang, M. I. Jordan, Learning transferable features with deep adaptation networks, in: ICML, 2015. 3

[10] B. Sun, K. Saenko, Deep coral: Correlation alignment for deep domain adaptation, in: ECCV, 2016. 3

[11] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by backpropagation, in: ICML, 2015. 3

[12] Y. Li, L. Yuan, N. Vasconcelos, Bidirectional learning for domain adaptation of semantic segmentation, in: CVPR, 2019. 3

[13] S. Xie, Z. Zheng, L. Chen, C. Chen, Learning semantic representations for unsupervised domain adaptation, in: ICML, 2018. 3

[14] Y. Zou, Z. Yu, B. Kumar, J. Wang, Domain adaptation for semantic segmentation via class-balanced self-training, in: ECCV, 2019. 3

[15] G. French, M. Mackiewicz, M. Fisher, Self-ensembling for visual domain adaptation, in: ICLR, 2018. 3

[16] X. J. Zhu, Semi-supervised learning literature survey, Tech. rep., University of Wisconsin-Madison Department of Computer Sciences (2005). 3

[17] Y. Chen, W. Li, C. Sakaridis, D. Dai, L. Van Gool, Domain adaptive faster r-cnn for object detection in the wild, in: CVPR, 2018. 3

[18] J. Hoffman, D. Wang, F. Yu, T. Darrell, FCNs in the wild: Pixel-level adversarial and constraint-based adaptation, arXiv:1612.02649. 3

[19] K. Saito, S. Yamamoto, Y. Ushiku, T. Harada, Open set domain adaptation by backpropagation, in: ECCV, 2018. 3, 4

[20] C. H. Lampert, H. Nickisch, S. Harmeling, Attribute-based classification for zero-shot visual object categorization, TPAMI. 3

[21] S. Reed, Z. Akata, H. Lee, B. Schiele, Learning deep representations of fine-grained visual descriptions, in: CVPR, 2016. 3

[22] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al., Devise: A deep visual-semantic embedding model, in: NIPS, 2013. 3

[23] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: NIPS, 2013. 3, 5

[24] Z. Akata, S. Reed, D. Walter, H. Lee, B. Schiele, Evaluation of output embeddings for fine-grained image classification, in: CVPR, 2015. 3

[25] M. Bucher, S. Herbin, F. Jurie, Improving semantic embedding consistency by metric learning for zero-shot classification, in: ECCV, 2016. 3

[26] M. Bucher, S. Herbin, F. Jurie, Generating visual representations for zero-shot classification, in: ICCV, 2017. 3

[27] R. Gao, X. Hou, J. Qin, L. Liu, F. Zhu, Z. Zhang, A joint generative model for zero-shot learning, in: ECCV, 2018. 3

[28] E. Kodirov, T. Xiang, S. Gong, Semantic autoencoder for zero-shot learning, in: CVPR, 2017. 3

[29] B. Romera-Paredes, P. Torr, An embarrassingly simple approach to zero-shot learning, in: ICML, 2015. 3

[30] R. Socher, M. Ganjoo, C. D. Manning, A. Ng, Zero-shot learning through cross-modal transfer, in: NIPS, 2013. 3

[31] Y. Xian, T. Lorenz, B. Schiele, Z. Akata, Feature generating networks for zero-shot learning, in: CVPR, 2018. 3

[32] Z. Zhang, V. Saligrama, Zero-shot learning via semantic similarity embedding, in: ICCV, 2015. 3

[33] V. Kumar Verma, G. Arora, A. Mishra, P. Rai, Generalized zero-shot learning via synthesized examples, in: CVPR, 2018. 3

[34] E. Schonfeld, S. Ebrahimi, S. Sinha, T. Darrell, Z. Akata, Generalized zero-and few-shot learning via aligned variational autoencoders, in: CVPR, 2019. 3

[35] M. Bucher, T.-H. Vu, M. Cord, P. Pérez, Zero-shot semantic segmentation, in: NeurIPS, 2019. 4, 5, 6, 10, 11, 12, 13

[36] N. Kato, T. Yamasaki, K. Aizawa, Zero-shot semantic segmentation via variational mapping, in: ICCV Workshops, 2019. 4

[37] Y. Xian, S. Choudhury, Y. He, B. Schiele, Z. Akata, Semantic projection network for zero-and few-label semantic segmentation, in: CVPR, 2019. 4

[38] M. Rohrbach, S. Ebert, B. Schiele, Transfer learning in a transductive setting, in: NIPS, 2013. 4

[39] Y. Fu, T. M. Hospedales, T. Xiang, S. Gong, Transductive multi-view zero-shot learning, TPAMI. 4

[40] Y. Guo, G. Ding, X. Jin, J. Wang, Transductive zero-shot recognition via shared model space learning, in: AAAI, 2016. 4

[41] E. Kodirov, T. Xiang, Z. Fu, S. Gong, Unsupervised domain adaptation for zero-shot learning, in: ICCV, 2015. 4

[42] Y. Yu, Z. Ji, X. Li, J. Guo, Z. Zhang, H. Ling, F. Wu, Transductive zero-shot learning with a self-training dictionary approach, IEEE transactions on cybernetics. 4

[43] Y. Yu, Z. Ji, J. Guo, Y. Pang, Transductive zero-shot learning with adaptive structural embedding, IEEE transactions on neural networks and learning systems. 4

[44] A. Zhao, M. Ding, J. Guan, Z. Lu, T. Xiang, J.-R. Wen, Domain-invariant projection learning for zero-shot recognition, in: NIPS, 2018. 4

[45] J. Song, C. Shen, Y. Yang, Y. Liu, M. Song, Transductive unbiased embedding for zero-shot learning, in: CVPR, 2018. 4

[46] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-decoder with atrous separable convolution for semantic image segmentation, in: ECCV, 2018. 5, 9

[47] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, A. M. Lopez, The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes, in: CVPR, 2016. 9

[48] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic urban scene understanding, in: CVPR, 2016. 9

[49] G. Varma, A. Subramanian, A. Namboodiri, M. Chandraker, C. Jawahar, Idd: A dataset for exploring problems of autonomous navigation in unconstrained environments, in: WACV, 2019. 9

[50] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, IJCV. 9

[51] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, C. L. Zitnick, Microsoft coco: Common objects in context, in: ECCV, 2014. 9

[52] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: CVPR, 2016. 9

[53] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: COMPSTAT, 2010. 9

[54] A. L. Maas, A. Y. Hannun, A. Y. Ng, Rectifier nonlinearities improve neural network acoustic models, in: ICML, 2013. 9

[55] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: a simple way to prevent neural networks from overfitting, JMLR. 9

[56] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, In ICLR. 10

[57] Y. Xian, C. H. Lampert, B. Schiele, Z. Akata, Zero-shot learning-a comprehensive evaluation of the good, the bad and the ugly, TPAMI. 10