# Learning Across Domains and Devices:
# Style-Driven Source-Free Domain Adaptation in Clustered Federated Learning

Donald Shenaj[*,1], Eros Fanì[*,2], Marco Toldo[1], Debora Caldarola[2], Antonio Tavera[2],
Umberto Michieli[†,1], Marco Ciccone[†,2], Pietro Zanuttigh[†,1], and Barbara Caputo[†,2]

[1]University of Padova, Italy          [2]Politecnico di Torino, Italy

## Abstract

*Federated Learning (*FL*) has recently emerged as a possible way to tackle the domain shift in real-world Semantic Segmentation (*SS*) without compromising the private nature of the collected data. However, most of the existing works on* FL *unrealistically assume labeled data in the remote clients. Here we propose a novel task (*FFReeDA*) in which the clients' data is unlabeled and the server accesses a source labeled dataset for pre-training only. To solve* FFReeDA*, we propose* LADD*, which leverages the knowledge of the pre-trained model by employing self-supervision with ad-hoc regularization techniques for local training and introducing a novel federated clustered aggregation scheme based on the clients' style. Our experiments show that our algorithm is able to efficiently tackle the new task outperforming existing approaches. The code is available at* `https://github.com/Erosinho13/LADD`.

## 1. Introduction

Federated Learning (FL) [49, 38, 1, 31, 52, 7, 36, 10] is a relatively new field of research that is attracting increasing interest. In FL, a learning task is solved through a collaboration among several edge devices, *i.e.*, clients, coordinated by a central server [49]. This learning paradigm is useful when data cannot be freely shared due to regulations, laws, and ethical principles: FL allows training a global model without leaking the users' data, preserving their privacy.

As an example, FL also constitutes a practical solution to tackle real-world vision tasks with data collected from multiple users in different scenarios. For instance, in the case of Semantic Segmentation (SS), it can be employed by self-driving cars for obstacle detection, and avoidance [19]. Most existing FL works assume the availability of labeled data on the client side. This assumption is clearly unrealistic due to the high cost and amount of manual work needed for dense pixel-level annotations [72].
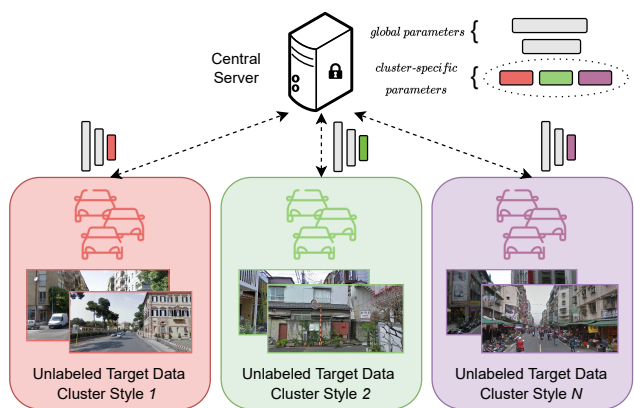


Figure 1. FFReeDA overview: clients having similar appearance are clustered together, while local learning is carried out exploiting both global and cluster-specific parameters. The clients' data is unlabeled and the source labeled dataset is kept on the server.

In this work, we focus on autonomous driving applications introducing a novel, more realistic setting for SS: **F**ederated source-**Free D**omain **A**daptation (FFReeDA). In FFReeDA, the server can pre-train the model on labeled source data. However, further accessing the source data is forbidden as in the Source-Free Domain Adaptation (SFDA) setting [41]. Clients access only their *unlabeled* target dataset, which they cannot share with other clients or the server. In particular, we consider real-world scenarios with several clients, each with a limited amount of images.

After the pre-training phase, the training setting is fully unsupervised. However, the objective of FFReeDA is not only to solve a multi-target domain adaptation problem in SS, rather to tackle specific issues arising in FL, such as *statistical* and *system heterogeneity* [38, 51], *communication bottleneck* [24], and clients' *privacy preservation* [37, 6]. To the best of our knowledge, no previous works addressed this problem and their related issues at the same time.

To address the FFReeDA problem, we propose LADD, a novel federated algorithm that assumes the presence of

---

*: Equal contribution. †: Equal supervision.

multiple distributions hidden among the clients. To exemplify, it is reasonable to assume that self-driving cars within the same city collect similar images. Indeed, the geographical proximity of two self-driving cars and different weather conditions could make the local datasets more or less similar. Therefore, LADD partitions the clients into clusters based on the styles of the images belonging to each client, trying to match them with their actual latent distribution. To minimize parameters duplication and improve communication efficiency, LADD splits the model's parameters in *shared*, globally aggregated across all clients, and *cluster-specific* which are aggregated only across clients within the same cluster, as visible in Fig. 1. Moreover, LADD takes full advantage of the source dataset during the pre-training stage with style transfer data augmentation [71], randomly loading the target styles in the source images to mimic the target distributions. Finally, LADD also leverages self-training through an ad-hoc pseudo-labeling strategy and stabilizes training with regularization techniques.

As FFREEDA is a novel setting, we adapted several baselines from other settings. LADD outperforms all baselines, showing the importance of designing specific algorithms for the proposed setting. To summarize:

- We introduce FFREEDA, a novel SS task for FL where we dropped the unrealistic assumption of dense labeled data at client side.
- We propose two realistic benchmarks for it, based on the Mapillary Vistas [53] and CrossCity [13] datasets.
- We propose LADD, a new federated algorithm tackling FFREEDA based on style-transfer and clustering.
- LADD shows excellent performance on all benchmarks with a source dataset (GTA5 [56]) and three different targets (Cityscapes [15], CrossCity, Mapillary), with diversified splits of the data across the clients.

## 2. Related Work

**Semantic Segmentation (SS)**, *i.e.*, classifying each pixel of an image with the corresponding semantic class, is an important challenge in many use cases such as self-driving cars [20]. State-of-the-art SS models rely on an encoder-decoder architectures, based on CNNs [44, 11, 12, 74, 59, 25] or transformers [17, 43, 14, 68] to generate dense predictions. These approaches typically assume a simplified, centralized setting in which the whole training dataset is available on a central server. However, this is not always possible due to privacy and efficiency constraints, and distributed training solutions must be considered.

**Domain Adaptation (DA).** Being a complex structured prediction task, SS generally requires expensive dense annotations. Recently, an increasing number of methods [63, 16] tackle this by training on synthetic data generated in virtual environments [56, 57, 2, 62]. Nonetheless, models trained on these data fail to generalize to the real world

because of the inherent domain shift between the simulated and real distributions. DA aims at reducing the performance gap between a *source domain* on which a model has been trained and a *target* one. When the target data is unlabeled, this is called Unsupervised DA (UDA). Initially, DA methods attempted to close the gap by measuring domain divergence [45, 66, 58]. Another popular direction is adversarial training [65, 48, 50], which includes the segmentation network and a domain discriminator competing in a *minimax* game. Other applications attempt to reduce domain shift by employing image-to-image translation algorithms to generate images modified with the style of the other domain [27, 55, 64]. Since this is a time-consuming technique, some non-trainable style translation algorithms, such as FDA [71], have been introduced. Modern approaches [40, 76, 4, 28] use self-learning techniques to create pseudo-labels from the target data, allowing the model to be fine-tuned even in a federated scenario in which each client observes its unlabeled domain.

**Main Challenges in FL.** Clients in FL have different hardware capabilities (*system heterogeneity*) and their data may belong to different distributions (*statistical heterogeneity*). Additionally, clients-server communication should be efficient [24] and privacy must be preserved by preventing the server to access clients' local data [37, 6].

**Vision Tasks in FL.** Thanks to its many applications in the real world and its potential in managing sensitive data, FL [49] has recently captured the interest of the research community [36, 30, 73]. However, most research papers focus on the theoretical aspects of FL [38, 31, 1, 47], neglecting its application to more complex vision tasks, *e.g.*, SS, and realistic scenarios, *e.g.*, heterogeneous domain distribution and unlabeled data observed at clients. A few exceptions are [52, 19, 8], which study FL SS and FL in the context of autonomous driving, and [61, 39, 5, 70, 69] that leverage medical images. Their main limitation is the costly assumption of having labeled data available.

In [47], the authors deal with the novel unsupervised FL setting from strong theoretical assumptions whilst only focusing on classification tasks for simple datasets such as MNIST [34], and CIFAR10 [32]. Focusing on SS, and proposing a more realistic approach, [72] introduces FMTDA (Federated Multi-target Domain Adaptation) to handle a few clients with unlabeled target local datasets belonging to different distributions while maintaining an open-access labeled source dataset on the server-side. Inspired by this work, we investigate the more complex setup of SFDA [42], in which the source dataset is only visible on the server during the pre-training phase and is not available to the clients. Moreover, we study a more realistic scenario where many more clients collaborate in the training but access much less data. As in FMTDA, we assume that clients' data may differ in terms of visual domains, *e.g.*, the scenes

collected by the autonomous vehicles in different geographical locations may have different weather or light conditions or may not show some semantic classes.

The study of DA in FL (both UDA and SFDA) is still in its early stages: [75] leverages UDA techniques for face recognition, [54] tackles domain shift via adversarial approaches, while [72] sees each client as a distinct target domain. To the best of our knowledge, this is the first work adapting SFDA to FL. Additional insights on vision tasks in FL other than SS and DA are reported in [3].

**Clustered FL (CFL).** In a real-world context, subsets of users typically share some common characteristics: for example, users in nearby geographic locations experience cities with similar architecture or weather conditions. Therefore, clients can be partitioned into clusters, each representing a specific set of conditions that we match to a corresponding *style* [33]. This approach falls under the literature of CFL [60], in which clustering is usually exploited for building personalized models that work well in a specific subdomain of interest [18, 23, 9]. Differently from these methods, we cluster clients based on the styles extracted from the unlabeled samples seen by each client.

## 3. Problem Setting

In this section, we formalize the proposed Federated Source Free Domain Adaptation (FFREEDA) setting.

Given a central server and the set of all clients $\mathcal{K}$ with $|\mathcal{K}| = K$, the input space $\mathcal{X}$, the output space $\mathcal{Y}$ and $N_p$ pixels in each image, the datasets are distinguished as follows: the source dataset $\mathcal{D}^S$ is kept on the server-side and is made of pairs of image and segmentation label $(x^S, y^S) \sim \mathscr{P}^S(\boldsymbol{x}, \boldsymbol{y})$, where $\boldsymbol{x}$ and $\boldsymbol{y}$ are the random variables following the distribution $\mathscr{P}^S$, associated with $x^S \in \mathcal{X}$ and $y^S \in \mathcal{Y}^{N_p}$ respectively; the $K$ target training datasets $\mathcal{D}_k^T = \{x_{k,i}^T \in \mathcal{X} \; \forall i \in |\mathcal{D}_k^T|\}$ are local to each client $k \in [K] := \{0, 1, ..., K-1\}$ and $x_{k,i}^T \sim \mathscr{P}_k^T(\boldsymbol{x})$.

By definition of the SFDA scenario, the source and test datasets share the same set of categories $\mathcal{Q} = \mathcal{Q}^S = \mathcal{Q}^T$. As for the federated setting, $K$ is reasonably large and the local datasets differ in terms of both size and distributions but have typically a much smaller size than the source dataset. Since users may share some common characteristics, it may happen that $\mathscr{P}_k^T(\boldsymbol{x}) = \mathscr{P}_h^T(\boldsymbol{x})$ for some $k, h \in [K]$. We assume the local datasets to be drawn from the same meta-distribution, which contains $G$ latent visual domains (*e.g.*, different cities), and each $\mathcal{D}_k^T$ contain only images from one of the $G$ latent domains. The test dataset $\mathcal{D}_{test}^T$ follows the target distribution $\mathscr{P}^T$ and is used to evaluate the final model learned across domains and devices.

Given the model $f(w) : \mathcal{X} \to \mathbb{R}^{N_p \times |\mathcal{Y}|}$ parametrized by $w$, the global objective is to obtain optimal segmentation performance on the target data distribution $\mathscr{P}^T(\boldsymbol{x})$, and it can be achieved by minimizing a suitable loss function, *i.e.*:

---

**Algorithm 1:** LADD (Learning Across Domains and Devices)

**Require:**
Source (labeled) dataset $\mathcal{D}^S$, clients $k \in \mathcal{K}$ with target (unlabeled) datasets $\mathcal{D}_k^T$, global model $f(w) = f(\{\theta, \phi\})$

**Clustering of the clients $\mathcal{K}$ and Pre-Training of $f$ on $\mathcal{D}^S$**
  Extract the styles $\mathcal{P}_k^s$ for each $k \in \mathcal{K}$
  Define the style-based clusters $\mathcal{C}$ (refer to Algorithm 2)
  Train $f(w)$ on $\mathcal{D}^S$ with style-transfer from $\mathcal{P}^s = \bigcup_{k \in \mathcal{K}} \mathcal{P}_k^s$

**Adaptation of $f$ on $\mathcal{D}^T$**
  **Initialize:**
  Cluster models $f_c(w_c) = f(w)$ and teachers $g_c(w_{g_c}) = f(w)$
  **for** *each round $t \in [T]$* **do**
    Randomly extract $\mathcal{K}^t \subset \mathcal{K}$. Let $c := \Gamma_{\mathcal{C}}(k)$.
    **for** $k \in \mathcal{K}^t$ *in parallel* **do**
      Set $f_k(w_k) = f_c(w_c)$
      $\phi_k^t, \theta_k^t \leftarrow \text{CLIENTUPDATE}(f_k, g_c, f, \mathcal{D}_k^T)$ (Sec. 4.3)
    $\phi^{t+1} \leftarrow$ Aggregate $\phi_k^t$ globally
    $\theta_c^{t+1} \leftarrow$ Aggregate $\theta_k^t$ within the cluster $c$
    **if** $t \mod \omega \equiv 0$ **then**
      **if** $t \geq t_{\text{START}}$ **then**
        $g_c^{t+1}(w_{g_c}) = \text{SWATUPDATE}(g_c^t) \; \forall c$ (Sec. 4.3)
      **else**
        $g_c(w_{g_c}) = f_c^t(w_c^t) \; \forall c \in \mathcal{C}$

---

$$w^* = \arg\min_w \sum_{k \in [K]} \frac{|\mathcal{D}_k^T|}{|\mathcal{D}^T|} \mathcal{L}_k(w) \qquad (1)$$

where $\mathcal{L}_k$ is the local loss function and $\mathcal{D}^T = \bigcup_{k \in \mathcal{K}} \mathcal{D}_k^T$.

## 4. Method

In this section, we describe in detail our FL algorithm by detailing the pre-training strategy (Sec. 4.1), the aggregation (Sec. 4.2) and the adaptation techniques (Sec. 4.3). The procedure is summarized in Fig. 2 and Algorithm 1.

### 4.1. Server Pre-training

The first step of LADD is a pre-training stage on the labeled source dataset $\mathcal{D}^S$. Before training, to bring the styles of source and target images closer together and improve the generalization of the pre-trained model, we apply the FDA [71] style transfer technique. First, the clients' style is transferred to the server and then applied to $\mathcal{D}^S$, on which the model is trained. Specifically, the $k$-th client extracts the style $s_k$ from each of its images, given by a window of width $l_s$ located at the center of the amplitude spectrum of that image [71], *i.e.* representing the amplitude of the lowest spatial frequency coefficients. Critically, these coefficients do not contain relevant information on the scene content, thus not breaking the user's privacy. The pool of styles $\mathcal{P}_k^s$ extracted from client $k$ is populated by sending the average of its extracted styles, *i.e.* $\mathcal{P}_k^s = \{\bar{s}_k\}$. On the server-side, the randomly initialized model $f(w)$ is trained on the source dataset $\mathcal{D}^S$, augmenting the source images with random styles extracted from the set $\mathcal{P}^s = \bigcup_{k \in \mathcal{K}} \mathcal{P}_k^s$. It is worth noting that these styles are never shared among the clients, and even a few images are sufficient to compute
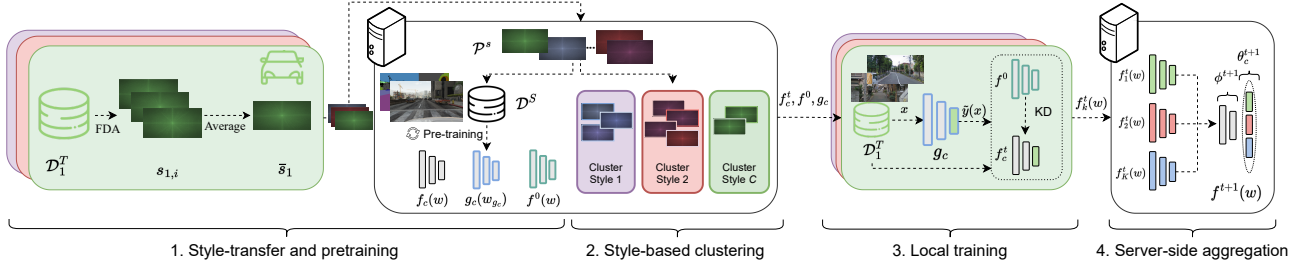
Figure 2. LADD overview (best seen in colors). 1) Each client $k$ extracts the average style $\bar{s}_k$ of its local data $\mathcal{D}_k^T$ using FDA. At server-side, the collected styles $\mathcal{P}^s$ are applied to the source dataset $\mathcal{D}_s$ during the supervised pre-training. 2) Clients are clustered according to their style. 3) At client-side, the cluster-specific teacher $g_c$ outputs the pseudo-labels, used for training $f_c^t$, leveraging KD from the pre-trained model. 4) At the server-side aggregation, we distinguish between global ($\phi^{t+1}$) and cluster-specific parameters ($\theta_c^{t+1}$).

$\bar{s}_k$ on each client. Once the pre-training stage is completed, the source dataset is no longer used.

## 4.2. Style-based Aggregation

In a realistic FL setting, different clients may observe similar samples, *e.g.* self-driving cars in the same region are likely to collect similar images and are not subject to statistical heterogeneity [38] during the server aggregation. On the other side, this premise does not apply to self-driving cars scattered throughout various distant places, which may learn conflicting information, thus affecting performances if naively aggregated. In addition, users may have access to a limited number of images, hindering clients from generalizing only from local optimization [8].

Taking these factors into account, we propose to explicitly cluster the clients to find their $G$ latent visual domains. To this end, we partition $\mathcal{K}$ in a set of non-empty clusters $\mathcal{C}$ with $|\mathcal{C}| = C$ and $\sum_{c \in \mathcal{C}} |c| = K$ basing on the clients' transferred styles. The centroid $\mu_c$ of each cluster is computed using $\bar{s}_k, \forall k \in \mathcal{C}$. We summarize our approach in Algorithm 2. We refer to H-Means instead of K-Means for symbolic convenience. First, we compute H-Means $N$ times $\forall \text{H} \in [n]_m$, with $N, n, m$ positive integers. Then, for each value of H, we select the partition $\mathcal{C}_\text{H}$ with the smallest intra-cluster distance, and compute its Silhouette Score. Finally, we select the clustering $\mathcal{C}$ with the highest Silhouette Score.

As for the server-side aggregation, instead of averaging the updates of the selected clients at each round as done by the standard FedAvg [49] algorithm, LADD introduces a clustered and layer-aware aggregation policy. We define as: (i) $w_k^t$ the weights of the model of client $k$ after $E$ local epochs of training at round $t$; (ii) $\theta_k^t$ and $\phi_k^t$ the group of *cluster-specific* and *global* parameters of the local model, such that $w_k^t = \theta_k^t \cup \phi_k^t$ and $\theta_k^t \cap \phi_k^t = \emptyset$; (iii) $\mathcal{K}^t \subset \mathcal{K}$ the subset of clients selected at round $t$. We globally aggregate the global parameters $\phi_k^t$ over all the selected clients in $\mathcal{K}^t$ to obtain the new parameter set $\phi^{t+1}$. On the other side the cluster-specific parameters $\theta_k^t$ are averaged within the clusters, resulting in $C$ specific parame-

---

**Algorithm 2:** Clustering Selection algorithm.

Let $d(\cdot, \cdot)$ be the *L2-norm* operator.

**Require:**
Clients $k \in \mathcal{K}$, target datasets $\mathcal{D}_k^T \; \forall k \in [K]$, function $\Gamma$ assigning each
  client to one of the $C$ clusters. Hyper-params $n, m, N \in \mathbb{N}_0, m < n$
**for** $\text{H} \in [n]_m := \{m, m+1, ..., n-1\}$ **do**
  **for** $n \in [N]$ **do**
    Change random seed $rs$
    $\mathcal{C}_\text{H}^{rs}$ = H-MEANS
    Compute
      $a_k(\mathcal{C}_\text{H}^{rs})$ = INTRACLUSTERDIST$(\mathcal{C}_\text{H}^{rs}, k) \; \forall k \in \mathcal{K}$
  $\mathcal{C}_\text{H} = \arg\min_{\mathcal{C}_\text{H}^{rs}} \sum_{k \in \mathcal{K}} a_k(\mathcal{C}_\text{H}^{rs})$
  Define $a_k^\text{H} := a_k(\mathcal{C}_\text{H}) \; \forall k \in \mathcal{K}$
  Compute $b_k^\text{H} := b_k(\mathcal{C}_\text{H})$ = INTERCLUSTERDIST$(\mathcal{C}_\text{H}, k) \; \forall k \in \mathcal{K}$
  Compute $\bar{\sigma}(\mathcal{C}_\text{H})$ = SILHOUETTESCORE$(a_k^\text{H}, b_k^\text{H} \; \forall k \in \mathcal{K})$
**return** $\mathcal{C} = \arg\max_\text{H} \bar{\sigma}(\mathcal{C}_\text{H})$

INTRACLUSTERDIST$(\mathcal{C}, k)$
  **return** $\frac{1}{|\Gamma_\mathcal{C}(k)| - 1} \sum_{h \in \Gamma_\mathcal{C}(k), h \neq k} d(k, h)$
INTERCLUSTERDIST$(\mathcal{C}, k)$
  **return** $\min_{c \in \mathcal{C}, c \neq \Gamma_\mathcal{C}(k)} \frac{1}{|c|} \sum_{h \in c} d(k, h)$
SILHOUETTESCORE$(\mathcal{C}, a, b)$
  $\sigma_k = \frac{b_k - a_k}{\max(a_k, b_k)}$ if $|\Gamma_\mathcal{C}(k)| > 1$, 0 otherwise, $\forall k \in \mathcal{K}$
  **return** $\frac{1}{K} \sum_{k \in \mathcal{K}} \sigma_k$

---

ter sets $\theta_c^{t+1}$ and $C$ models $f_c^{t+1}(\boldsymbol{x}; w_c^{t+1}), \forall c \in \mathcal{C}$ where $w_c^{t+1} = \phi^{t+1} \cup \theta_c^{t+1}$. Note that the server is not required to store independent models for each cluster, it is sufficient to save only the cluster-specific parameters $\theta_c^{t+1} \; \forall c \in \mathcal{C}$ and the global parameters $\phi^{t+1}$, loading them when needed. At test time, given the $i$-th target test image, we (i) extract the style $s_{\text{TEST}, i}$, (ii) compute the *L2-norm* between $s_{\text{TEST}, i}$ and all the cluster centroids $\mu_c \; \forall c \in \mathcal{C}$, (iii) select the cluster $c$ with the smallest *L2-norm*, and (iv) use the model $f_c^{t+1}(\boldsymbol{x}; w_c^{t+1})$ to evaluate the model on the $i$-th test image.

## 4.3. Client Adaptation

Here we introduce the main components of the local loss $\mathcal{L}_k(w)$ and the employed regularization techniques.

**Self-Training.** At round $t$, given an image $x$ on the $k$-th client, we train the local model $f_k^t(\boldsymbol{x}; w_k^t)$ by employing hard one-hot pseudo-labels $\tilde{y}(x) \in \mathbb{R}^{Q \times N_p}$ with the same thresholding mechanism as proposed in [71],

Table 1. Federated SS splits employed in our work.

| Split | $Q$ | $\vert\mathcal{D}^T\vert$ | $\vert\mathcal{D}^T_{test}\vert$ | $\vert\mathcal{K}\vert$ | # Img/Client (range) |
|---|---|---|---|---|---|
| Cityscapes | 19 | 2975 | 500 | 144 | [10, 45] |
| CrossCity | 13 | 12800 | 400 | 476 | [17, 37] |
| Mapillary | 19 | 17969 | 2000 | 357 | [16, 100] |
| CrossCity (split of [72]) | 13 | 12800 | 400 | 4 | 3200 |

where $Q = \vert\mathcal{Q}\vert$ is the number of classes. To reduce the computation burden on the client, we avoid having client-specific teacher networks. On the contrary, the pseudo-labels are computed using a cluster-specific teacher network $g_c^t(\boldsymbol{x}; w_{g_c}^t)$, which outputs the predictions $\hat{y}(x) := g_c^t(x; w_{g_c})$. The teacher parameters $w_{g_c}^t$ are first initialized as $w_{g_c}^0 = w$ and then updated every $\omega$ rounds as $w_{g_c}^t = w_c^t$.

**Regularization.** Pseudo-labels allow the clients to mimic the presence of the labels. However, after a few training iterations, the learning curve starts dropping [21]. At first, self-training allows to reduce the gap between the knowledge extracted by $\mathcal{D}^S$ and the one needed to perform well on the target datasets $\mathcal{D}^T_k$. Later on though the network starts being too confident on its predictions, reducing its effectiveness and making more miss-classifications. It therefore becomes of the utmost importance to try to reverse this trend. To this end, we exploit a **Knowledge Distillation** (KD) loss $\mathcal{L}_{\text{KD}}$ [26] based on the soft predictions given by the pre-trained model to prevent $f^t(w)$ from forgetting the knowledge acquired during the pre-training phase. However, our experiments showed KD on its own was not enough for avoiding overfitting, since the learning curve starts slightly dropping again during the last rounds of adaptation (see Suppl. Mat. for more details). Inspired by the recent success of Stochastic Weight Averaging (SWA) [29] in FL [8], we apply a moving average to the clients' teachers $g_c$ after a starting round $t_{\text{START}}$ as $w_{g_c}^{t+\omega} = (w_{g_c}^t n_{g_c}^t + w_c^{t+\omega})/(n_{g_c}^t + 1)$, where $n_{g_c}^t = (t - t_{\text{START}})/\omega$. We name this technique **SWA teacher** (SWAt). SWAt allows noise reduction, further stabilizes the learning curve and enables the model to better converge to the local minimum of the total loss $\mathcal{L} = \mathcal{L}_{\text{PSEUDO}} + \lambda_{\text{KD}}\mathcal{L}_{\text{KD}}$, where $\lambda_{\text{KD}}$ is an hyper-parameter to control the KD.

## 5. Experiments

### 5.1. Experimental Setup

**Datasets.** We evaluate the proposed framework in synthetic-to-real experimental setups for autonomous driving applications, which are commonly used as benchmark for domain adaptation methods. For the source domain, we opt for the synthetic GTA5 dataset [56]. It comprises 24966 highly realistic road scenes of typical US-like urban and suburban environments. As for the real (*i.e.*, target) domain, we experiment with three different datasets: Cityscapes [15], CrossCity [13] and Mapillary Vistas [53]. We use un-labeled training samples from all the datasets. Results are reported on the original validation split for Cityscapes and Mapillary, while on the test split for CrossCity.

Cityscapes provides street-view images from 50 cities in Central Europe. CrossCity includes more diverse locations and appearances, collecting driving scenes from multiple cities around the world (*i.e.*, Rome, Rio, Tokyo, and Taipei). Finally, the Mapillary Vistas dataset collects geo-localized street-view images from all around the world. We consider the largest number of overlapping classes among GTA5 and the real datasets (*i.e.*, 19 for Cityscapes and Mapillary, and 13 for CrossCity).

We propose a federated partitioning of the target datasets among the clients for each of the target dataset (*i.e.*, a *split*), as summarized in Table 1 and detailed in Suppl. Mat. For Cityscapes, we use the *heterogeneous* split from [19], where 144 clients observe images taken only from one city. We emulated the same kind of split also for the CrossCity dataset. Finally, we used the GPS information of the Mapillary dataset to discover clients with spatially near images.

**Baselines and Competitors.** To support the efficacy of the proposed approach in the unexplored FFREEDA setup, we compare with multiple methods on both centralized and federated setups. As the lower bound, we consider the naïve *source only* approach, which entails the sole use of source labeled data during training. At the other end, as the upper bound we propose the FTDA and Oracle comparisons. Both the methods assume the availability of supervised target data, either on the server side (*i.e.*, centralized framework) or on the client side (*i.e.*, federated framework). However, the FTDA method implies that target data is used to fine-tune the model after a source-only pre-training, while the Oracle simply consists in a supervised FedAvg training on the labeled version of the target dataset. Moreover, we re-implemented the Maximum Classifier Discrepancy (MCD) [58] method, and we adapted it to the federated setting as the authors of [72] did. Furthermore, we compare our method with DAFormer [28], that we regard as the current state-of-the-art UDA approach. We remark that both baselines are evaluated in the UDA setting, *i.e.*, a simpler scenario where source and target datasets are jointly available. Concerning the FL aggregation, [35, 8] show that algorithms like FedProx [38] and SCAFFOLD [31] typically do not provide an improvement for vision tasks, therefore we focused on other aspects of the training, like domain adaptation, clustering and style transfer techniques.

**Server Pre-train.** We pre-trained the model on GTA5 using a power-law decreasing learning rate $\eta$, starting from $\eta = 5.0 \cdot 10^{-3}$ with power 0.9 and using SGD optimizer with momentum equal to 0.9 and no weight decay. We pre-train the model for 15k steps. Each client computes the style on all its images using a window of size $3 \times 3$ and sends the mean style to the server, before the pre-training starts.

**Federated Adaption.** In CrossCity, we run the experiments with fixed $\eta = 1.0 \cdot 10^{-2}$, training on 4 clients per round for a number of rounds $T = 1000$, with $\lambda_{KD} = 20$; we update the pseudo-label teacher model every round ($\omega = 1$) and set $t_{START} = 400$ for SWAt. In Cityscapes, we trained on 5 client per round, with $T = 300$, fixed $\eta = 5e - 5$, $\lambda_{KD} = 10$, $\omega = 5$ and $t_{START} = 200$. For Mapillary we used $\eta = 1.0 \cdot 10^{-2}$, $\lambda_{KD} = 10$, 6 clients per round, with $T = 100$, same pseudo-label policy of Cityscapes and $t_{START} = 50$. In both settings, for all datasets, the batch size was 16. We performed data augmentation as follows: random scaling $(0.7, 2)$, random crop of $1024 \times 512$, color jitter with brightness, contrast and saturation equal to 0.5, and image normalization. For Mapillary instead of random scaling we forced a fixed rescaling with width equal to 1024.

## 5.2. Experimental Results

**GTA5→Cityscapes** Our first setup is the GTA5 → Cityscapes adaptation. Experimental results are reported in Table 2. Even though providing high quality realistic images, the GTA5 dataset still suffers a domain gap compared to real-world images, as those included in Cityscapes. We notice that simply training over supervised source data (*i.e.*, *source only*) leads to a significant performance discrepancy compared to the full target supervision. Even applying the state-of-the-art DAFormer method [28] in an UDA setting (*i.e.*, assuming joint availability of source supervised and aggregated target unsupervised data, which violates the assumptions of our setup), there is a noticeable performance drop of around $25\%$ of mIoU from the supervised oracle.

In our setup, we assume a federated learning framework with private target data distributed among multiple clients and a large-scale source dataset only available in a central server for the pre-training stage only. This introduces additional challenges not present in standard centralized domain adaptation settings. In particular, we assume that source and target data are not accessible on the same device, and that target data is available in small batches scattered among devices and not in a single place. Furthermore, target data is heterogeneously distributed among clients. The increase in the task complexity is noticeable from the performance drop of the supervised target oracle and FTDA methods (which still assume target supervision). This is also true for the MCD [58] UDA approach, which loses almost $10\%$ of mIoU when tested in a federated setting.

The proposed method is able to obtain robust results in this challenging setting achieving a mIoU of around $36.5\%$. In particular, the efficient pre-training based on domain stylization, along with the self-training optimization scheme, allows to tackle the lack of source data at the client side. We additionally improve training stability, which is hindered by the small amount of target data available within single clients, with KD and SWAt, as shown by the very small

Table 2. Results on the heterogeneous split of Cityscapes.

| Setting | Method | mIoU (%) |
|---|---|---|
| centralized | Oracle | $66.64 \pm 0.33$ |
| centralized | Source Only | $24.05 \pm 1.14$ |
| centralized | FTDA | $65.74 \pm 0.48$ |
| centralized | MCD [58] | $20.55 \pm 2.66$ |
| centralized | DAFormer [28] | $42.31 \pm 0.20$ |
| federated | Oracle | $58.16 \pm 1.02$ |
| federated | FTDA | $59.35 \pm 0.61$ |
| FL-UDA | MCD [58] | $10.86 \pm 0.67$ |
| FFREEDA | FedAvg$^\dagger$ [49] + Self-Tr. | $35.10 \pm 0.73$ |
| FFREEDA | **LADD (cls)** | $\mathbf{36.49 \pm 0.13}$ |
| FFREEDA | **LADD (all)** | $\mathbf{36.49 \pm 0.14}$ |

standard deviation of the results in Table 2. Finally, we provide an enhanced aggregation mechanism, which indirectly shares task information in a effective manner among clients sharing similar input statistics (*i.e.*, with smaller domain gap), according to our style-based client clustering. By observing results in Table 2, we notice that LADD keeps a similar performance gap w.r.t. the target oracle compared to what DAFormer achieves in a centralized UDA setting. Competitive results are provided by different variations of the proposed style-based clustering, *i.e.*, by keeping only the decoder (*i.e.*, *LADD (cls)*) or the whole network (*i.e.*, *LADD (all)*) as cluster-specific during the aggregation. Additional analyses are provided in Sec. 5.3.

**GTA5→CrossCity** We further investigate the performance of the proposed approach in the GTA5→CrossCity scenario. Quantitative results are reported in Table 3. We compare with the naïve source only baseline, as well as with MCD [58]. Due to the lack of target supervision on training images, the upper bound of the target oracle cannot be provided, nor the result of FTDA.

The diverse content and appearance of CrossCity's road scenes, due to variable geographic origin of its samples, provide a heterogeneous target distribution. The enhanced heterogeneity w.r.t. the more uniform Cityscapes dataset in turn leads to a tougher challenge for federated training. For instance, the MCD method when extended from a centralized to a federated learning framework suffers from a substantial performance reduction. Instead, LADD provides a much higher accuracy in a federated setting, with more than $17\%$ gain over federated MCD, while also not requiring reuse of source data after the initial pre-training. This is indicative of the robustness of our method w.r.t. the statistical diversity of client target data.

Finally, we note that, by allowing only a minimal amount of network parameters to be cluster-dependent, we achieve a final accuracy very close to our best result, obtained without any parameter sharing across clusters of clients. This result shows that LADD demands limited communication overhead w.r.t. standard FedAvg.

Table 3. Results on the proposed CrossCity split.

| Setting | Method | mIoU (%) |
|---|---|---|
| centralized | Source Only | 26.49 ± 1.46 |
| centralized | MCD [58] | 27.15 ± 0.87 |
| FL-UDA | MCD [58] | 24.80 ± 1.56 |
| **FFREEDA** | FedAvg[†] [49] + Self-Tr. | 33.59 ± 1.25 |
| **FFREEDA** | **LADD (cls)** | 39.87 ± 0.14 |
| **FFREEDA** | **LADD (all)** | **40.09 ± 0.19** |

Table 4. Results on the proposed Mapillary split.

| Setting | Method | mIoU% |
|---|---|---|
| centralized | Oracle | 61.46 ± 0.21 |
| centralized | Source Only | 32.40 ± 0.71 |
| centralized | MCD | 31.93 ± 1.89 |
| federated | Oracle | 49.91 ± 0.49 |
| FL-UDA | MCD | 19.15 ± 0.75 |
| **FFREEDA** | FedAvg[†] [49]+ Self-Tr. | 38.97 ± 0.21 |
| **FFREEDA** | **LADD (cls)** | **40.16 ± 1.02** |
| **FFREEDA** | **LADD (all)** | 38.78 ± 1.82 |

**GTA5→Mapillary** Finally, we provide an experimental analysis with target data from Mapillary Vistas, while the GTA5 still serves as source dataset. Table 4 contains numerical results of the evaluation. The diverse assortment of target data collected around the world, and dispersed among clients according to geographic location (see Sec. 5.1), makes client data distribution even more heterogeneous than the previous setups. We notice that with target supervision (*i.e.*, oracle method), there exists a significant performance drop of around 11.5% of mIoU from centralized to federated settings. This is even more noticeable with the MCD [58] UDA approach, which struggles when tested under the considered federated learning setup, suffering from a similar mIoU decrease of 12%. The proposed LADD framework instead provides considerable performance, with results surpassing the source only optimization by a large margin (more than 8% of mIoU) in its best version with only classifier weights kept cluster-specific (*LADD (cls)*), and approaching the federated oracle result. We further observe that LADD outperforms (in its best configuration) the simpler framework based on FedAvg and self-training. This supports the effectiveness of the proposed additional modules (concerning local training regularization and federated aggregation, see Sec. 4.2 and 4.3) in tackling the domain adaptation problem in a distributed learning setting. Finally, we remark that the lightweight aggregation scheme with a cluster-specific classifier achieves the best results.

## 5.3. Ablation Studies

**Impact of Optimization Modules.** We now study the contribution of each component of our method. In Table 5, we report the target mIoU computed with modules incrementally activated, showing the gain brought by each of

---

†: Same pretrain as LADD.

Table 5. Ablation of our optimization framework, performed on CrossCity dataset for the presented split.

| FDA | ST | KD | SWAt | Cluster Aggr | mIoU (%) |
|---|---|---|---|---|---|
| | | | | | 26.49 ± 1.46 |
| | ✓ | | | | 30.58 ± 0.59 |
| ✓ | | | | | 32.43 ± 0.61 |
| | | ✓ | ✓ | ✓ | 32.78 ± 0.09 |
| ✓ | ✓ | | | | 33.59 ± 1.25 |
| ✓ | ✓ | ✓ | | | 37.49 ± 0.14 |
| ✓ | ✓ | ✓ | ✓ | | 38.83 ± 0.12 |
| ✓ | ✓ | ✓ | | ✓ | 39.18 ± 0.24 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **40.09 ± 0.19** |

them (GTA5→CrossCity setup). We notice that the introduction of the style transfer technique during server pre-training generates an improvement of almost 6% of mIoU, where the lower bound (source only pre-training) coincides with the centralized source only experiment. The naïve federated adaptation using the proposed self-training routine (and FedAvg aggregation) allows to gain an initial improvement of 2% of mIoU, but leads to unstable training curves (see Suppl. Mat.). Adding the KD module we further boost the performance (by almost 3% mIoU) and prevent clients' optimization from undertaking unsteady behaviors, with the initial stable configuration as anchor point. By activating SWAt we get an extra boost in terms of final mIoU and training stability, and training convergence is achieved much more consistently. We remark how the std drops when enabling KD and SWAt. When further introducing cluster aggregation, but leaving SWAt disabled, we get a small mIoU increase, at the price of higher instability. Finally, by adding the cluster aggregation along with SWAt we get our complete method, for a final score of 40.09% of mIoU, achieved with stable training.

**Style-based Pre-training.** We analyze the impact of the stylization mechanism on the pre-training. We test different style extraction schemes, by varying the size of the Fourier amplitude window. Table 6 reports quantitative results of the ablation study (GTA5→CrossCity setup). It is possible to observe that a window of 1x1 is sufficient to capture and transfer useful domain-dependent information across domains. However, by increasing the dimension of the style window to $3 \times 3$ and $5 \times 5$ pixels we get an improvement of 1% mIoU. In addition, even though providing similar results to the $5 \times 5$ size in terms of final mIoU, the $3 \times 3$ window leads to a more stable pre-training (testified by lower std), due to less artifacts being introduced in the stylized images.

We remark that the style data occupies a very limited amount of memory (in the order of a few bytes), and thus requires little communication overhead to be transmitted from client to server before federated rounds start. The style window corresponds to a small portion of the Fourier Transform and, prior to its transmission, is averaged over all data samples within each client. Recall also that shape information is
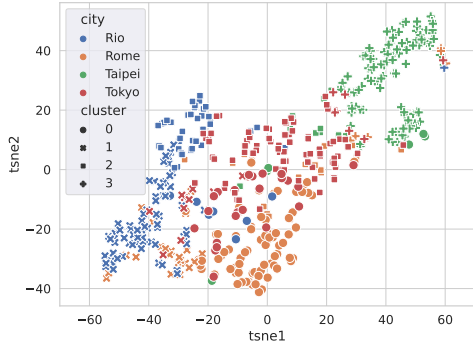
Figure 3. t-SNE of the styles. The colors represent the city ground-truth, while the symbols represent the inferred clusters. The clustering accuracy, considering each cluster a city, is equal to 0.68.

| Size | mIoU (%) | | Layers | mIoU (%) |
|------|----------|--|--------|----------|
| None | $26.49 \pm 1.46$ | | None | $38.83 \pm 0.12$ |
| $1 \times 1$ | $31.59 \pm 0.68$ | | BN | $38.72 \pm 0.20$ |
| $3 \times 3$ | $32.43 \pm 0.61$ | | Backbone | $39.31 \pm 0.13$ |
| $5 \times 5$ | $32.51 \pm 0.75$ | | Classifier | $39.87 \pm 0.14$ |
| | | | **All** | **$40.09 \pm 0.19$** |

Table 8. Results on the CrossCity split proposed in [72].

| Method | Rio | Rome | Taipei | Tokyo | **Avg** |
|--------|-----|------|--------|-------|---------|
| Source Only | 27.9 | 27.6 | 26.0 | 28.2 | 27.4 |
| Cent-MCD [58] | 31.3 | 30.6 | 28.8 | 31.6 | 30.5 |
| Fed-DAN [46] | 27.3 | 26.4 | 26.0 | 28.5 | 27.1 |
| Fed-DANN [22] | 28.6 | 26.0 | 26.6 | 28.6 | 27.5 |
| Fed-MCD [58] | 27.7 | 27.3 | 26.5 | 29.0 | 27.6 |
| DualAdapt [72] | 29.2 | 28.0 | 27.6 | 30.7 | 28.9 |
| LADD (ours) | **35.4** | **34.0** | **31.5** | **32.4** | **33.3** |

mapped to the phase data, while we transmit only amplitude information. Therefore, we argue that it encloses a negligible fraction of the overall information of the local image data, and thus does not violate data privacy limitations.

**Cluster-level Aggregation.** We study how different cluster-based aggregation schemes affect the overall adaptation performance. In particular, we select different groups of model parameters for cluster-specific and global aggregations. Results obtained in the GTA5→CrossCity setup are collected in Table 7. We notice that keeping only batch-norm parameters cluster-dependent (second row) provides similar results as the standard FedAvg (first row), where all model parameters are global. When holding per-cluster backbone and classifier blocks individually, we get improved performance. The accuracy is further slightly boosted when the entire model is kept cluster-specific, showing that backbone and classifier both enclose cluster-dependent information.

Finally, we observe that solely treating the lightweight classifier block as cluster-specific gives comparable performance to the best full-model intra-cluster aggregation. Therefore, we found a less computational and memory demanding version of the proposed LADD approach, still providing robust performance.

**Style-based Client Clustering.** We analyze the distribution of clients across the style-based clusters identified by our approach in an unsupervised fashion. The study is conducted when the CrossCity target dataset is employed, so that we can compare style-based clusters with those determined by the city of origin. In Fig. 3 we associate each cluster with a point in a 2D space according to its style tensor. In particular, each style tensor is flattened resulting into a 27-D vector, and the t-SNE [67] dimensional reduction method is used to project it into a 2D space. We can observe that clients from different cities (each of which identified by a different color) tend to be clustered together. At the same time, we notice that clients with similar styles (*i.e.*, are associated to the same style-based cluster by our approach, see Sec. 4.2) are projected in adjacent regions. Fur-

thermore, both city- and style-based partitions appear to be highly overlapping, signifying that style-based clustering is effectively able to capture domain-dependent information, which is highly correlated to the geographical location.

**Comparison in a simpler CrossCity split.** We compare with the FMTDA method [72] in the GTA5→CrossCity adaptation setup they propose, with target data distributed over 4 total clients, each containing images of one of the 4 CrossCity's cities. For a fair comparison, we use the same segmentation network adopted in [72]. In this simpler federated setting, a restricted version of LADD without the cluster-level aggregation is still effective. In Table 8, we report the mIoU computed on each city, along with the average value, for different approaches directly taken from [72]. We observe a consistent improvement, which in the average target mIoU reaches almost 5%, demonstrating the superiority of the proposed LADD approach.

## 6. Conclusion

In this work we introduced FFREEDA, a new challenging and realistic setting for Source-Free Domain Adaptation in Federated Learning for Semantic Segmentation. In FFREEDA, a server-side labeled dataset is used for pre-training the model, while local training uses only the unlabeled clients' data. We introduced LADD, an innovative algorithm to solve FFREEDA, employing (i) style-transfer, knowledge distillation and SWA teacher on the pseudo-labels for regularizing learning, and (ii) style-driven clustering for learning both global and personalized parameters. LADD has no direct competitors due to the novel setup FFREEDA, but is still able to achieve competing results compared to state-of-the-art algorithms. We also provided two new splits adapting the CrossCity and Mapillary Vistas datasets to the federated scenario as a reference for future research in FL SS.

# References

[1] Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N Whatmough, and Venkatesh Saligrama. Federated learning based on dynamic regularization. *ArXiv preprint*, 2021.

[2] E. Alberti, A. Tavera, C. Masone, and B. Caputo. IDDA: A large-scale multi-domain dataset for autonomous driving. *IEEE Robotics and Automation Letters*, 5(4):5526–5533, 2020.

[3] Mohammed Aledhari, Rehma Razzak, Reza M Parizi, and Fahad Saeed. Federated learning: A survey on enabling technologies, protocols, and applications. *IEEE Access*, 8:140699–140725, 2020.

[4] Francesco Barbato, Marco Toldo, Umberto Michieli, and Pietro Zanuttigh. Latent space regularization for unsupervised domain adaptation in semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2835–2845, 2021.

[5] Cosmin I Bercea, Benedikt Wiestler, Daniel Rueckert, and Shadi Albarqouni. Feddis: Disentangled federated learning for unsupervised brain pathology segmentation. *arXiv preprint arXiv:2103.03705*, 2021.

[6] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečnỳ, Stefano Mazzocchi, Brendan McMahan, et al. Towards federated learning at scale: System design. *Proceedings of Machine Learning and Systems*, 1:374–388, 2019.

[7] Daniel Ramage Brendan McMahan. Federated learning: Collaborative machine learning without centralized training data. ai.googleblog.com, 2017.

[8] Debora Caldarola, Barbara Caputo, and Marco Ciccone. Improving generalization in federated learning by seeking flat minima. In *European Conference on Computer Vision*, 2022.

[9] Debora Caldarola, Massimiliano Mancini, Fabio Galasso, Marco Ciccone, Emanuele Rodolà, and Barbara Caputo. Cluster-driven graph federated learning over multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2749–2758, 2021.

[10] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *ArXiv preprint*, 2018.

[11] 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 (TPAMI)*, 40(4):834–848, 2017.

[12] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.

[13] Yi-Hsin Chen, Wei-Yu Chen, Yu-Ting Chen, Bo-Cheng Tsai, Yu-Chiang Frank Wang, and Min Sun. No more discrimination: Cross city adaptation of road scene segmenters. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 1992–2001, 2017.

[14] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Neural Information Processing Systems (NeurIPS)*, 34:9355–9366, 2021.

[15] 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 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3213–3223, 2016.

[16] Gabriela Csurka, Riccardo Volpi, and Boris Chidlovskii. Unsupervised domain adaptation for semantic image segmentation: a comprehensive survey. *arXiv preprint arXiv:2112.03241*, 2021.

[17] 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.

[18] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. *Advances in Neural Information Processing Systems*, 33:3557–3568, 2020.

[19] Lidia Fantauzzo, Eros Fanì, Debora Caldarola, Antonio Tavera, Fabio Cermelli, Marco Ciccone, and Barbara Caputo. Feddrive: Generalizing federated learning to semantic segmentation in autonomous driving. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022.

[20] Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaeser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems*, 22(3):1341–1360, 2020.

[21] Francois Fleuret et al. Uncertainty reduction for model adaptation in semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9613–9623, 2021.

[22] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning (ICML)*, pages 1180–1189. PMLR, 2015.

[23] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19586–19597. Curran Associates, Inc., 2020.

[24] Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient algorithm for federated learning. In *International Conference on Machine Learning (ICML)*, pages 3973–3983. PMLR, 2020.

[25] Ali Hassani, Steven Walton, Nikhil Shah, Abulikemu Abuduweili, Jiachen Li, and Humphrey Shi. Escaping the big data paradigm with compact transformers. *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2021.

[26] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015.

[27] Judy Hoffman, E. Tzeng, T. Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A. Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *International Conference on Machine Learning (ICML)*, pages 1989–1998, 2018.

[28] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. *ArXiv preprint*, 2021.

[29] Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *ArXiv preprint*, 2018.

[30] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.

[31] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. *ArXiv preprint*, 2019.

[32] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.

[33] Jogendra Nath Kundu, Akshay Kulkarni, Amit Singh, Varun Jampani, and R Venkatesh Babu. Generalize then adapt: Source-free domain adaptive semantic segmentation. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 7046–7056, 2021.

[34] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[35] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. *ArXiv preprint*, 2021.

[36] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020.

[37] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.

[38] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *ArXiv preprint*, 2018.

[39] Wenqi Li, Fausto Milletarì, Daguang Xu, Nicola Rieke, Jonny Hancox, Wentao Zhu, Maximilian Baust, Yan Cheng, Sébastien Ourselin, M. Jorge Cardoso, and Andrew Feng. Privacy-preserving federated brain tumour segmentation. In *Machine Learning in Medical Imaging*. Springer, 2019.

[40] Qing Lian, Fengmao Lv, Lixin Duan, and Boqing Gong. Constructing self-motivated pyramid curriculums for cross-domain semantic segmentation: A non-adversarial approach. In *Proceedings of International Conference on Computer Vision (ICCV)*, October 2019.

[41] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. *Arxiv preprint*, 2021.

[42] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1215–1224, 2021.

[43] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 10012–10022, 2021.

[44] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3431–3440, 2015.

[45] Mingsheng Long, Y. Cao, J. Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. In *International Conference on Machine Learning (ICML)*, pages 97–105, 2015.

[46] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In *International Conference on Machine Learning (ICML)*, pages 2208–2217. PMLR, 2017.

[47] Nan Lu, Zhao Wang, Xiaoxiao Li, Gang Niu, Qi Dou, and Masashi Sugiyama. Federated learning from only unlabeled data with class-conditional-sharing clients. *ArXiv preprint*, 2022.

[48] Yawei Luo, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2507–2516, 2019.

[49] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics (AISTATS)*, pages 1273–1282. PMLR, 2017.

[50] Umberto Michieli, Matteo Biasetton, Gianluca Agresti, and Pietro Zanuttigh. Adversarial learning and self-teaching techniques for domain adaptation in semantic segmentation. *IEEE Transactions on Intelligent Vehicles*, 5(3):508–518, 2020.

[51] Umberto Michieli and Mete Ozay. Are all users treated fairly in federated learning systems? In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2318–2322, 2021.

[52] Umberto Michieli and Mete Ozay. Prototype guided federated learning of visual feature representations. *Arxiv preprint*, 2021.

[53] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kontschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 4990–4999, 2017.

[54] Xingchao Peng, Zijun Huang, Yizhe Zhu, and Kate Saenko. Federated adversarial domain adaptation. In *International Conference on Learning Representations*, 2020.

[55] Fabio Pizzati, Raoul de Charette, Michela Zaccaria, and Pietro Cerri. Domain bridge for unpaired image-to-image translation and unsupervised domain adaptation. In *Proceedings of Winter Conference on Applications of Computer Vision (WACV)*, pages 2990–2998, 2020.

[56] Stephan R. Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *Proceedings of European Conference on Computer Vision (ECCV)*, Lecture Notes in Computer Science, pages 102–118. Springer International Publishing, Cham, 2016.

[57] 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 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3234–3243, 2016.

[58] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3723–3732, 2018.

[59] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4510–4520, 2018.

[60] Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE transactions on neural networks and learning systems*, 32(8):3710–3722, 2020.

[61] Micah J Sheller, G Anthony Reina, Brandon Edwards, Jason Martin, and Spyridon Bakas. Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation. In *International MICCAI Brainlesion Workshop*, pages 92–104. Springer, 2018.

[62] Paolo Testolina, Francesco Barbato, Umberto Michieli, Marco Giordani, Pietro Zanuttigh, and Michele Zorzi. SELMA: SEmantic Large-scale Multimodal Acquisitions in Variable Weather, Daytime and Viewpoints. *arXiv preprint arXiv:2204.09788*, 2022.

[63] Marco Toldo, Andrea Maracani, Umberto Michieli, and Pietro Zanuttigh. Unsupervised domain adaptation in semantic segmentation: a review. *Technologies*, 8(2):35, 2020.

[64] Marco Toldo, Umberto Michieli, Gianluca Agresti, and Pietro Zanuttigh. Unsupervised domain adaptation for mobile semantic segmentation based on cycle consistency and feature alignment. *Image and Vision Computing*, 95:103889, 2020.

[65] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, and M. Chandraker. Learning to adapt structured output space for semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[66] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance, 2014.

[67] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research (JMLR)*, 9(11), 2008.

[68] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Neural Information Processing Systems (NeurIPS)*, 34:12077–12090, 2021.

[69] An Xu, Wenqi Li, Pengfei Guo, Dong Yang, Holger R Roth, Ali Hatamizadeh, Can Zhao, Daguang Xu, Heng Huang, and Ziyue Xu. Closing the generalization gap of cross-silo federated medical image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20866–20875, 2022.

[70] Xuanang Xu, Tianyi Chen, Han Deng, Tianshu Kuang, Joshua C Barber, Daeseung Kim, Jaime Gateno, Pingkun Yan, and James J Xia. Federated cross learning for medical image segmentation. *Arxiv preprint*, 2022.

[71] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4085–4095, 2020.

[72] Chun-Han Yao, Boqing Gong, Hang Qi, Yin Cui, Yukun Zhu, and Ming-Hsuan Yang. Federated multi-target domain adaptation. In *Proceedings of Winter Conference on Applications of Computer Vision (WACV)*, pages 1424–1433, 2022.

[73] Chen Zhang, Yu Xie, Hang Bai, Bin Yu, Weihong Li, and Yuan Gao. A survey on federated learning. *Knowledge-Based Systems*, 216:106775, 2021.

[74] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, July 2017.

[75] Weiming Zhuang, Xin Gan, Yonggang Wen, Xuesen Zhang, Shuai Zhang, and Shuai Yi. Federated unsupervised domain adaptation for face recognition. *arXiv preprint arXiv:2204.04382*, 2022.

[76] Yang Zou, Zhiding Yu, BVK Vijaya Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 289–305, 2018.

# Supplementary Material
# Learning Across Domains and Devices:
# Style-Driven Source-Free Domain Adaptation in Clustered Federated Learning

This document contains supporting material for the paper *Learning Across Domains and Devices: Style-driven Source-Free Domain Adaptation in Clustered Federated Learning*. Here, we include additional details on the federated splits employed in the paper along with analyses of the convergence stability of our approach when compared to competing strategies adapted to our federated setup. Finally, we show some qualitative segmentation maps.

## S1. Additional Details on Splits

In this section, we complete the description of how the federated splits used in our experiments are generated.

**Cityscapes.** We used the *heterogeneous* federated split of Cityscapes [15] proposed in [19]. The split comprises 144 clients, where each client has between 10 and 45 samples belonging to a single city from the dataset. Further details on the distribution of the number of images per client are shown in Figure S1.

**CrossCity.** We generated the federated split of the CrossCity [13] dataset by assigning $27 \pm 10$ images taken from the same city to each client, where the number of samples per client is uniformly sampled. The final distributions of the number of images per client are shown in Figure S2 both per city and overall. We observe how the distributions are balanced across the four cities.

**Mapillary.** We propose a novel split for the Mapillary Vistas [53] dataset via a clustering procedure based on the GPS coordinates of the images. We started from the original training set of 18000 images and discarded 31 of them missing the GPS coordinates. Then, we run the k-Means algorithm over the GPS coordinates six times, one per continent. The k-Means algorithm is constrained to assign every client a random number of images in the range 16 and 100. The procedure resulted in 357 clients, where each client observed samples from only one continent. The final distributions of the number of images per client are shown in Figure S3. Unlike the other scenarios, we observe a large variability across the distributions obtained in different continents due to the highly imbalanced nature of the dataset. Also, note that the two entries with higher values, 16 and 100, correspond to the extreme values of the constrained
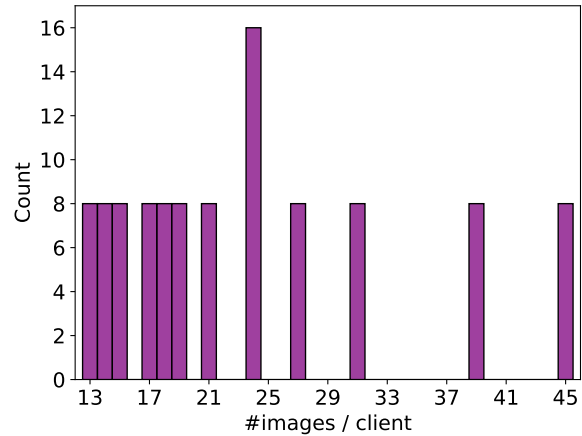


Figure S1. Histogram of images per client in the federated Cityscapes split.

k-Means process.

## S2. Additional Details on the Style-Based Client Clustering

In a realistic FL setting, different clients may observe similar samples, *e.g.* self-driving cars in the same region are likely to collect similar images, thus they are not subject to statistical heterogeneity during the server aggregation. Therefore, we proposed a style-driven client clustering as one of the foundational parts of our algorithm. During the FL optimization stage, we employed the identified communities in a clustered and layer-aware aggregation policy on the server side.

First of all, we remark that the four clusters identified by the styles extracted from the images contain mostly clients belonging to one single geographical location (*i.e.*, city). Table S1 shows the number of clients belonging to a specific city assigned to each cluster for the federated Cross-City dataset. Overall, the clustering accuracy, considering each cluster a city, is equal to 68%. Therefore, there is not a one-to-one correspondence of the clusters with the cities.

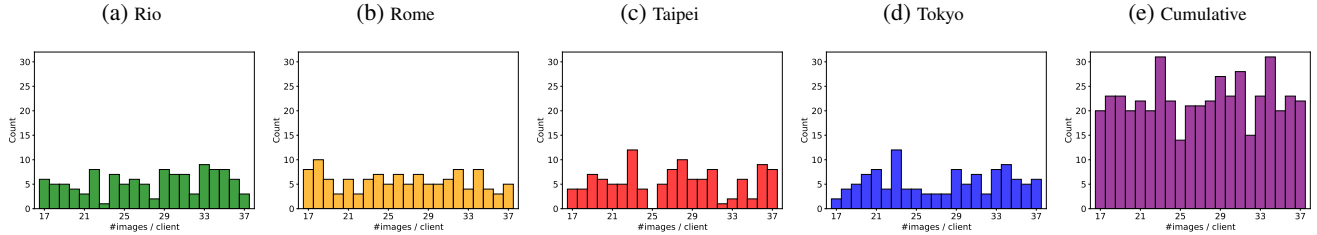To investigate this aspect, we show in Figure S4 some

(a) Rio  (b) Rome  (c) Taipei  (d) Tokyo  (e) Cumulative

Figure S2. Histogram of images per clients in the proposed federated CrossCity split.



(a) Africa  (b) Asia  (c) Europe  (g) Cumulative
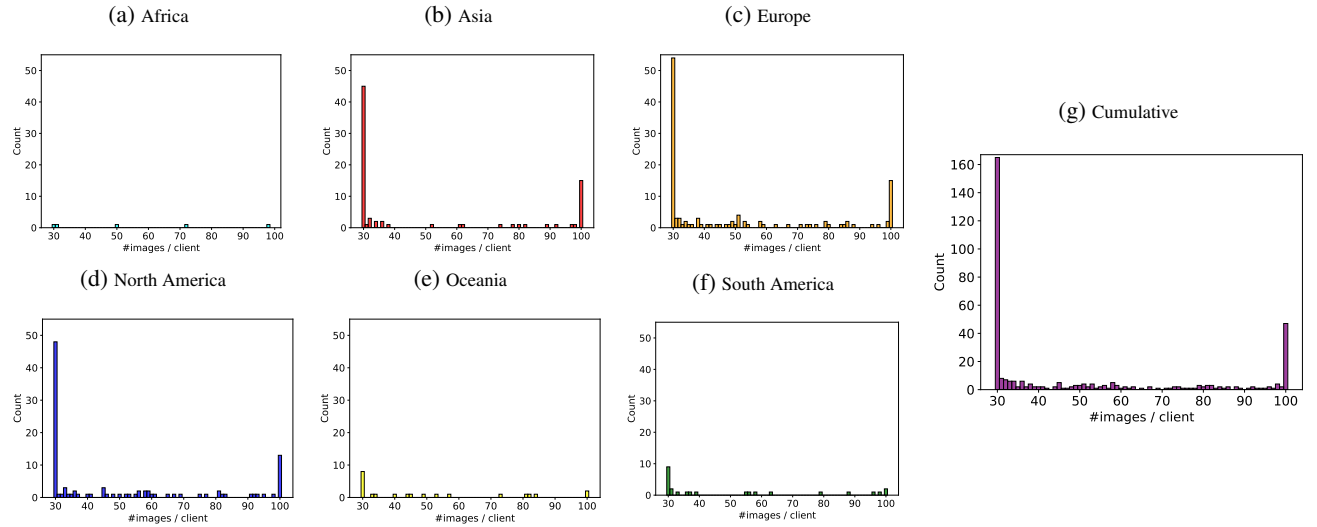(d) North America  (e) Oceania  (f) South America

Figure S3. Histogram of images per clients in the proposed federated Mapillary Vistas split.

Table S1. Number of clients belonging to a specific city assigned to each cluster for the federated CrossCity split.

|        | Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
|--------|-----------|-----------|-----------|-----------|
| **Rio**    | 7    | 1    | **70** | 38   |
| **Rome**   | **76** | 6  | 22   | 17   |
| **Taipei** | 6    | **103** | 0 | 9    |
| **Tokyo**  | 26   | 8    | 10   | **73** |

samples taken from the clients belonging to each of the four clusters in the federated CrossCity dataset. Here, we observe an interesting finding: despite being generated via style information only, the clusters tend to show scenes with similar semantics. For instance, *Cluster 1* contains clients having images of large and trafficked streets, and grayish sky. *Cluster 2* contains clients having images of narrow streets with little to no vegetation, many buildings, a few parked cars and whitish sky. *Cluster 3* contains clients having images of empty roads with green surrounding vegetation. *Cluster 4* contains clients having images from sunny weather and blue sky, narrow streets with no traffic and green vegetation.

Finally, we show in Figure S5 some samples taken from the clients belonging to each of four clusters in the federated Mapillary dataset. Unlike as for CrossCity, here we do not appreciate a clear assignment as the number of clus-

ters is different from the number of towns or continents. Therefore, we observe that here the clustering is much more appearance-related, according to the style of the images.

For instance, *Cluster 1* contains clients having cloudy and foggy images where the visual appearance is grayish. *Cluster 2* contains clients having grayish sky and yellowish buildings with some similar semantics across clients. *Cluster 3* contains clients having images at the sunset or sunrise where the light scatters yellow shadows. *Cluster 4* contains clients having images with predominant blue colors in the sky.

## S3. Implementation Details

The proposed method is implemented in PyTorch, the code and federated splits are available at `https://github.com/Erosinho13/LADD`.

The semantic segmentation network used is DeepLab-V3 [12] with Mobilenet-V2 [59] as the backbone and width multiplier equal to 1, representing a good compromise in terms of performance and lightness, important aspects to consider for real-world applications, such as self-driving cars. On each communication round, the selected clients are trained sequentially, allowing to perform the complete simulation and reproduce the results on a single GPU with

Figure S4. Sample images in each cluster for the federated CrossCity split.



Figure S5. Sample images in some clusters for the federated Mapillary split.

32GB of VRAM (we used a NVIDIA RTX 3090).

## S4. Qualitative Results

We provide some qualitative results in the form of segmentation maps of target images generated by the segmentation model subject to different adaptation schemes. Figures S6, S7 and S8 refer to the 3 adaptation setups chosen for experimental evaluations, with respectively CrossCity, Cityscapes and Mapillary as target datasets. We compare the naïve source only training (3rd columns in all the aforementioned figures) and the baseline federated adaptation strategy (4th columns), based on FedAvg[49] aggregation and local self-training, with the proposed LADD (when cluster-specific aggregation is extended to all the segmentation network layers) (last columns). For fair comparison we employ the same pretraining for FedAvg and LADD. By inspecting the segmentation maps produced by the different adaptation strategies, we notice how the *source only* maps show inconsistent and noisy predictions, where semantically similar classes are confused, such as *sidewalk* and *road* or *terrain* in all the reported samples. Local self-training and standard FedAvg aggregation at server-side partially mitigate the prediction accuracy drop caused by domain shift between source and target data. Nonetheless, we observe that the adapted model still tends to mistake semantically-similar classes such as sidewalk and road in
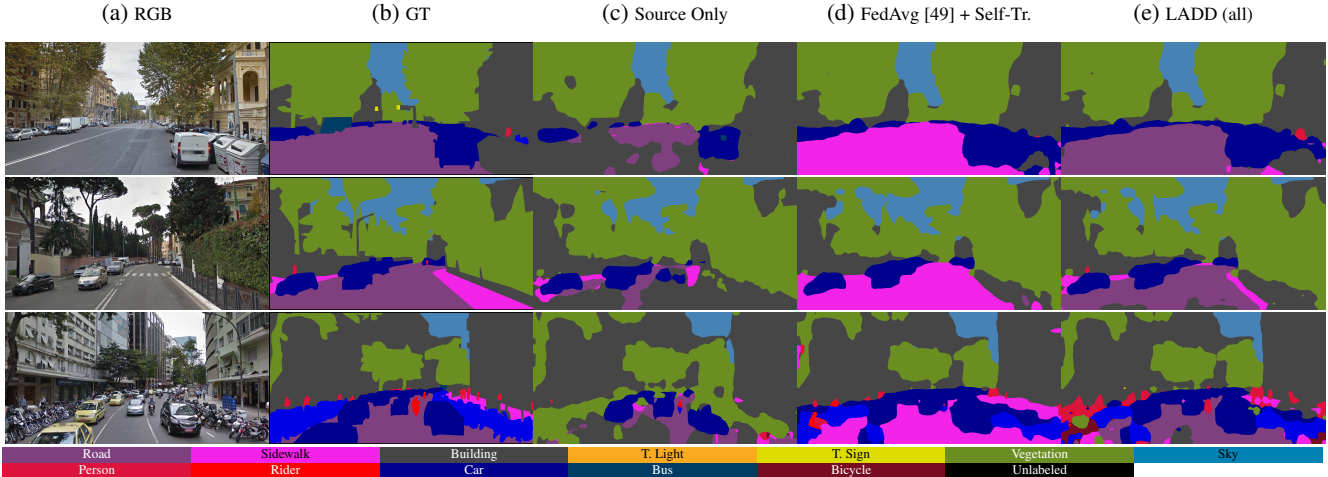
| (a) RGB | (b) GT | (c) Source Only | (d) FedAvg [49] + Self-Tr. | (e) LADD (all) |

| Road | Sidewalk | Building | T. Light | T. Sign | Vegetation | Sky |
| Person | Rider | Car | Bus | Bicycle | Unlabeled | |

Figure S6. GTA5→CrossCity qualitative results.



| (a) RGB | (b) GT | (c) Source Only | (d) FedAvg [49] + Self-Tr. | (e) LADD (all) |

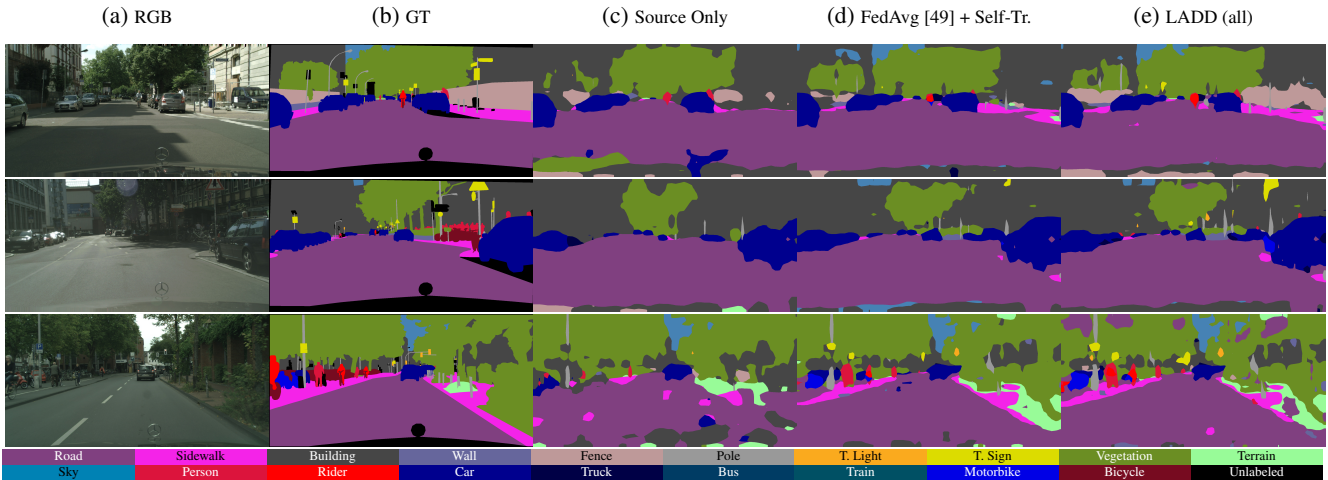| Road | Sidewalk | Building | Wall | Fence | Pole | T. Light | T. Sign | Vegetation | Terrain |
| Sky | Person | Rider | Car | Truck | Bus | Train | Motorbike | Bicycle | Unlabeled |

Figure S7. GTA5→Cityscapes qualitative results.

the first sample of Figure S6. The proposed regularized local training leads to more robust local optimization, which otherwise tends to suffer from unsteady behavior, due to the small amount of available training data and the lack of any form of supervision (even from the source domain) at the client side. This, along with the cluster-specific semantically aware aggregation mechanism, results into less noisy and more accurate predictions as we can see in the last columns of the figures.

## S5. Additional Quantitative Results

Finally, we report additional results in the form of per-class IoUs achieved when different modules of our framework are enabled. Once more, results are reported with CrossCity (Table S2), Cityscapes (Table S3) and Mapillary (Table S4) as target datasets, in terms of mean and standard deviation computed over the last 10% rounds.

When enabled, we observe that each module improves

the overall mIoU score, which is also generally shared by the individual IoU scores of the semantic classes in the different experimental setups.

In addition, in Figure S9 we report the learning curves as a result of federated optimization under different configurations of the proposed LADD method in the GTA→CrossCity setup. When only ST is employed in the client-side optimization, the training is extremely unstable, showing a small initial burst of performance followed by a rapid decrease after few rounds. When adding KD and then SWAt, the training curves become progressively more robust and stable, achieving the best results when KD and SWAt are joined by the cluster-specific aggregation, in either classifier-exclusive or full model configuration of cluster-specific parameters. We finally remark how LADD in its complete configuration is characterized by steady and converging learning curves, unaffected by diverging phenomena.
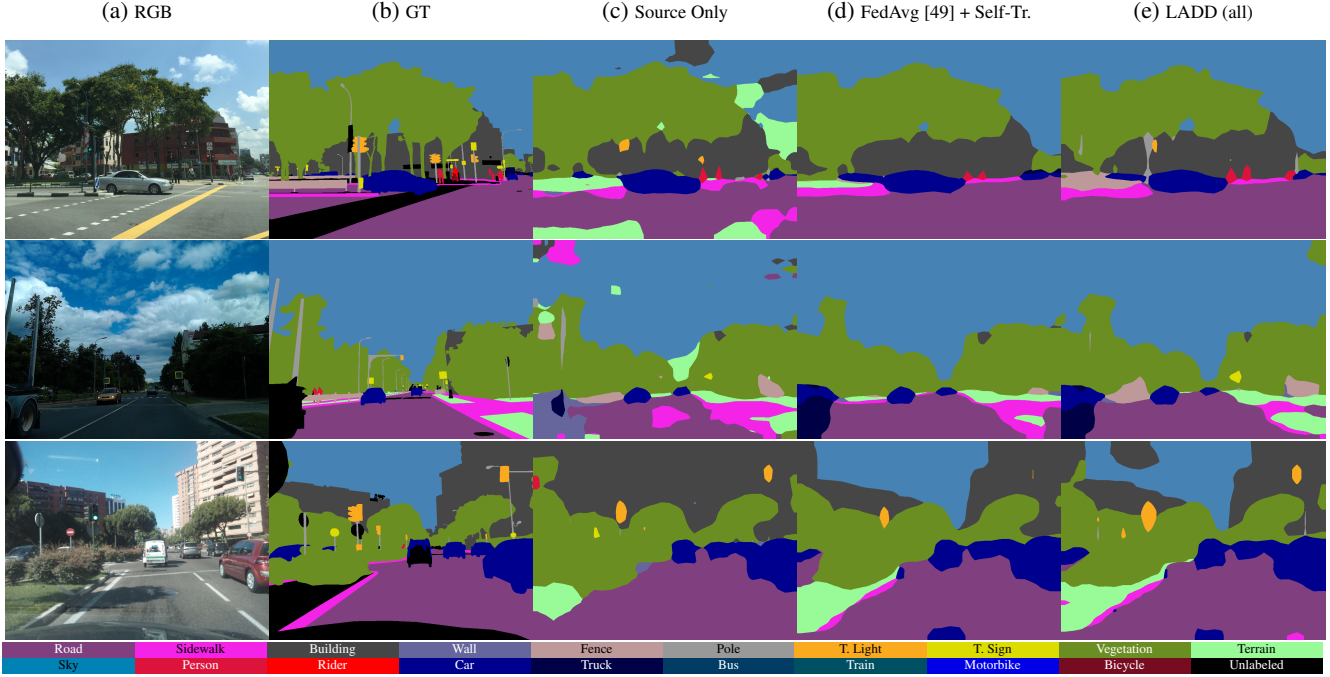
Figure S8. GTA5→Mapillary qualitative results.

Table S2. CrossCity IoU by class and mIoU (%).

| FDA | ST | KD | SWAt | Cl Aggr | road | sidewalk | building | traffic light | traffic sign | vegetation | sky | person | rider | car | bus | motorcycle | bicycle | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | 25.6 | 21.6 | 65.9 | 3.9 | 8.6 | 67.5 | 73.5 | 33.1 | 2.1 | 43.0 | 6.6 | 0.3 | 0.2 | 26.5 ± 1.5 |
| ✓ | | | | | 38.2 | 24.0 | 74.8 | 7.0 | 8.9 | 70.5 | 80.9 | 37.0 | 4.0 | 63.6 | 12.0 | 3.5 | 0.0 | 32.4 ± 0.6 |
| ✓ | ✓ | | | | 21.9 | 17.9 | 81.3 | 9.5 | 14.5 | 77.4 | 85.2 | 41.0 | 2.3 | 66.1 | 10.7 | 8.0 | 0.9 | 33.6 ± 1.3 |
| ✓ | ✓ | ✓ | | | 49.5 | 26.7 | 81.2 | 11.7 | 12.4 | 77.6 | 87.0 | 40.4 | 1.0 | 68.9 | 16.3 | 11.3 | 3.4 | 37.5 ± 0.1 |
| ✓ | ✓ | ✓ | ✓ | | 53.3 | 28.6 | 81.1 | 12.1 | 12.0 | 77.5 | 87.1 | 42.1 | 1.9 | 68.9 | 17.2 | 16.7 | 4.9 | 38.8 ± 0.1 |
| ✓ | ✓ | ✓ | | ✓ | 63.4 | 32.3 | 81.5 | 12.1 | 12.2 | 77.5 | 86.9 | 41.0 | 1.1 | 69.0 | 16.0 | 12.5 | 3.8 | 39.2 ± 0.2 |
| ✓ | ✓ | ✓ | ✓ | ✓ | 64.3 | 33.7 | 81.0 | 12.5 | 14.4 | 77.2 | 86.8 | 42.1 | 1.4 | 69.1 | 18.1 | 15.6 | 4.8 | 40.1 ± 0.2 |

Table S3. Cityscapes IoU by class and mIoU (%).

| FDA | ST | KD | SWAt | Cl Aggr | road | sidewalk | building | wall | fence | pole | traffic light | traffic sign | vegetation | terrain | sky | person | rider | car | truck | bus | train | motorcycle | bicycle | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ | | | | 84.5 | 36.8 | 77.3 | 23.9 | 11.3 | 20.5 | 29.1 | 22.6 | 76.9 | 26.5 | 68.9 | 53.4 | 13.7 | 79.0 | 15.2 | 14.0 | 1.4 | 11.0 | 5.1 | 35.1 ± 0.7 |
| ✓ | ✓ | ✓ | | | 79.3 | 34.0 | 73.6 | 22.0 | 16.4 | 24.6 | 30.3 | 31.3 | 61.7 | 23.2 | 70.1 | 51.2 | 19.3 | 73.7 | 13.6 | 17.9 | 7.3 | 12.1 | 15.3 | 35.6 ± 0.1 |
| ✓ | ✓ | ✓ | ✓ | ✓ | 80.0 | 36.1 | 74.1 | 22.8 | 18.3 | 26.3 | 30.6 | 33.0 | 65.2 | 25.4 | 69.4 | 52.3 | 19.1 | 74.5 | 13.4 | 18.0 | 7.2 | 12.6 | 14.2 | 36.5 ± 0.1 |

Table S4. Mapillary Vistas IoU by class and mIoU (%).

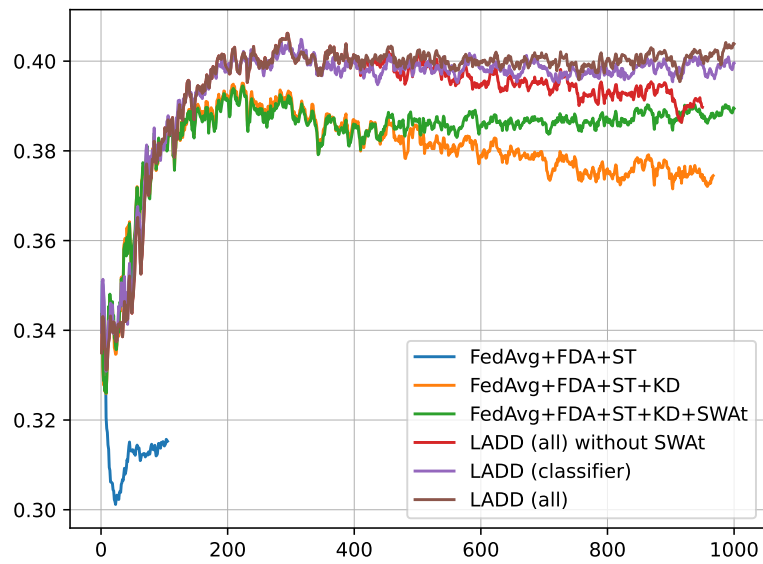| FDA | ST | KD | SWAt | Cl Aggr | road | sidewalk | building | wall | fence | pole | traffic light | traffic sign | vegetation | terrain | sky | person | rider | car | truck | bus | train | motorcycle | bicycle | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ | | | | 67.4 | 36.9 | 74.7 | 24.8 | 25.4 | 10.9 | 21.0 | 33.3 | 72.8 | 40.8 | 91.2 | 46.1 | 23.1 | 73.7 | 31.1 | 22.7 | 3.1 | 30.6 | 11.9 | 39.0 ± 0.2 |
| ✓ | ✓ | ✓ | ✓ | | 75.4 | 37.7 | 73.4 | 25.2 | 25.2 | 18.3 | 26.6 | 37.1 | 73.5 | 38.1 | 91.4 | 45.5 | 13.8 | 71.3 | 30.9 | 22.0 | 3.0 | 29.9 | 19.1 | 40.0 ± 0.1 |
| ✓ | ✓ | ✓ | ✓ | ✓ | 75.5 | 37.0 | 69.1 | 24.6 | 25.6 | 18.9 | 26.7 | 38.2 | 72.5 | 36.4 | 89.4 | 46.3 | 17.2 | 70.7 | 32.6 | 20.2 | 4.1 | 31.4 | 21.0 | 40.2 ± 1.0 |

Figure S9. Comparison of learning curves in the CrossCity federated split.