

Supplementary Material: Beyond deterministic translation for unsupervised domain adaptation

Eleni Chiou¹, Eleftheria Panagiotaki¹, and Iasonas Kokkinos¹

University College London

1 Overview

We provide additional details about the training procedure and additional results including quantitative results and qualitative results.

2 Training procedure

In Sec. 2.1 we provide some additional details regarding the training procedure of the stochastic translation network while in Sec. 2.2 we provide additional details about the pseudolabeling. Finally, in Sec. 2.3 we summarize in more detail the entire training process.

2.1 Stochastic translation for unsupervised domain adaptation

The stochastic translation network is based on MUNIT [3]. We have described it in the main paper but we also describe it here in more detail. For a more extensive presentation we refer the reader to [3]. As it is illustrated in Fig. 1 the stochastic translation network consists of content encoders $\{C_s, C_t\}$, style encoders $\{S_s, S_t\}$, generators $\{G_s, G_t\}$ and domain discriminators $\{D_s, D_t\}$ for the source domain s , and the target domain t respectively.

Given a source domain image $x \in \mathcal{S}$, we start by extracting a domain-invariant content code $c = C_s(x)$ and a domain-specific style code $s_s = S_s(x)$. Then, we perform within-domain reconstruction (Fig. 1a) and cross-domain translation (Fig. 1b). We reconstruct the original image x , using the source generator G_s that takes as input at the first layer the content code c and its subsequent layers are modulated by Adaptive Instance Normalization (AdaIN) [2] driven by the style code s_s . This amounts in minimizing the following objective function:

$$L_{rec}^s = \sum_{x \in \mathcal{S}} \|x - G_s(C_s(x), S_s(x))\|.$$

We perform translation from the source domain to the target domain by passing the content code c as input to the target generator G_t , whose subsequent layers are

modulated by AdaIn driven by a random style code $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This results in the following stochastic translation function from the source domain to the target domain:

$$\mathbf{T}[x, \mathbf{v}] \doteq G_t(C_s(x), \mathbf{v}).$$

We ensure that the resulting translation matches the distribution of the target domain data by employing the following GAN objective:

$$L_{GAN}^t = \sum_{x \in \mathcal{T}} \log D_t(x) + \sum_{x \in \mathcal{S}} \log(1 - D_t(\mathbf{T}[x, \mathbf{v}])),$$

where \mathcal{T}, \mathcal{S} the target and source datasets respectively and D_t the adversarial discriminator for the the target domain t .

We ensure that the content codes c of the source image and the translated image are aligned by minimize the following objective:

$$L_{rec}^{c_s} = \sum_{x \in \mathcal{S}} \|C_t(G_t(C_s(x), \mathbf{v})) - C_s(x)\|_2.$$

Similarly, to align the target style code with the Gaussian prior distribution we use an objective of the following form:

$$L_{rec}^{s_t} = \sum_{x \in \mathcal{S}} \|S_t(G_t(C_s(x), \mathbf{v})) - \mathbf{v}\|_2.$$

As we described in detail in the main paper we ensure that the semantics are preserved during translation using the following objective function:

$$L_{sem}^s = \mathcal{L}_{ce}(F(\mathbf{T}[x, \mathbf{v}]), p),$$

where \mathcal{L}_{ce} , the cross-entropy loss, F a segmentation network trained on both source and target data, $F(\mathbf{T}[x, \mathbf{v}])$ the softmax output given the translated image $\mathbf{T}[x, \mathbf{v}]$ and $p = \text{argmax}(F(x))$ the predicted labels for the source image $x \in \mathcal{S}$.

The exact same procedure is followed for translating from the target to the source domain and the corresponding loss terms L_{rec}^t , L_{GAN}^s , $L_{rec}^{c_t}$, $L_{rec}^{s_s}$, L_{sem}^t are defined similarly.

The full objective is given by

$$\begin{aligned} \min_{\substack{C_s, S_s, G_s \\ C_t, S_t, G_t}} \max_{\substack{D_s, D_t}} & \quad \lambda_x(L_{rec}^s + L_{rec}^t) + \lambda_{GAN}(L_{GAN}^s + L_{GAN}^t) \\ & + \lambda_c(L_{rec}^{c_s} + L_{rec}^{c_t}) + \lambda_s(L_{rec}^{s_s} + L_{rec}^{s_t}) \\ & + \lambda_{sem}(L_{sem}^s + L_{sem}^t), \end{aligned} \tag{11}$$

where λ_{GAN} , λ_x , λ_c , λ_s , λ_{sem} are weights that control the importance of each term. In our experiments, we set the weights as follows: $\lambda_{GAN} = 1$, $\lambda_x = 10$, $\lambda_c = 1$, $\lambda_s = 1$, $\lambda_{sem} = 1$.

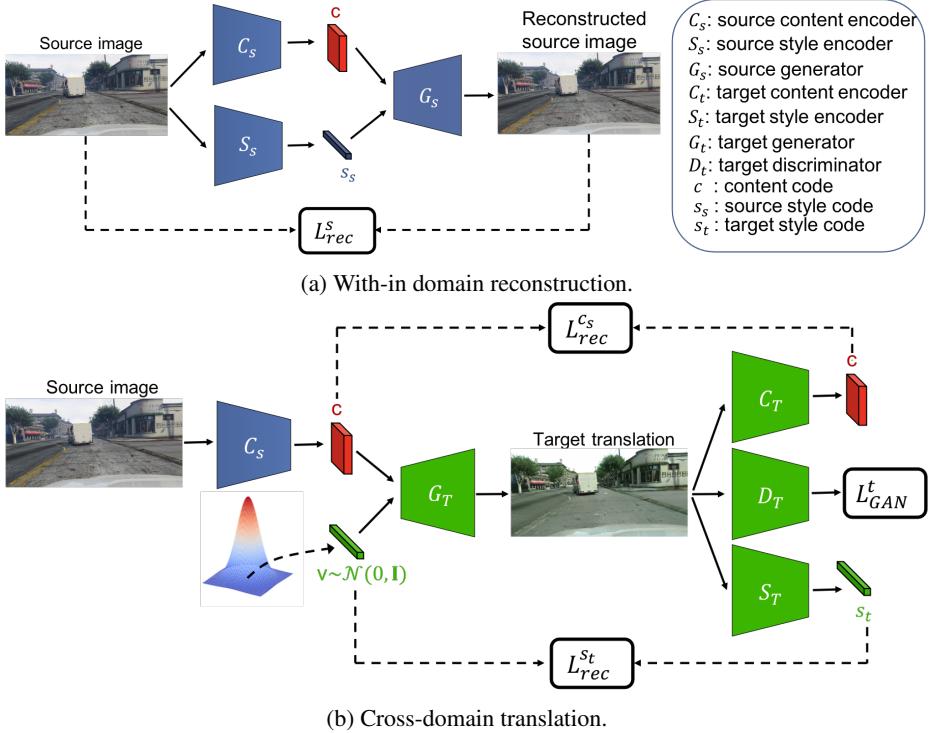


Fig. 1: Stochastic translation [3] from the source to the target domain: with-in domain reconstruction and cross-domain translation allows us to reconstruct the input and pass the content from the source image to its counterparts respectively. The target cycle is omitted for clarity.

2.2 Robust pseudolabel generation

As we mentioned in Sec. 3.4 of the main paper, we rely on three complementary networks to generate robust pseudolabels. We train two target networks $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$; one with the variance left intact and the other with the variance scaled by 10. We also train a source network, F_s , and exploit multiple samples to obtain a better estimate of the pseudo-labels as we described in Sec. 3.3 of the main paper.

The enhanced probability map used to generate pseudolabels is obtained by the weighted average of the predictions of the three networks:

$$\hat{y} = \frac{1}{3}\hat{y}_s + \frac{1}{3}\hat{y}_{t,\sigma^2=1} + \frac{1}{3}\hat{y}_{t,\sigma^2=10}, \quad (12)$$

where \hat{y}_s , $\hat{y}_{t,\sigma^2=1}$, $\hat{y}_{t,\sigma^2=10}$ indicates the pixel-level posterior distribution on labels obtained by F_s , $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$ respectively.

We assign pseudolabels to samples for which the dominant class has a score above a certain threshold θ . Similarly to [5] we use class-wise confidence thresholds to assign pseudolabels. In particular for each class c , the threshold θ_c equals the probability

ranked at $r * N_c$, where N_c is the number of pixels predicted to belong in class c and r is the proportion of pseudolabels we want to retain. In cases where $\theta_c > 0.9$, we set $\theta_c = 0.9$. In Sec. 3.1 we provide results for the selection of r based on the mIoU of the validation set.

2.3 Training process

Algorithm 1 summarizes the training process. Initially we train a segmentation network, F , operating on both domains and use this to impose a semantic consistency constraint to the stochastic translation network trained using Eq. 11. Using translated images obtained by the stochastic translation network, we train two target-domain networks and one source-domain network. We generate robust pseudolabels for the target-domain data by combining the predictions of the three models. The pseudolabels are used in the next round of training as supervision for the networks when they are driven by target-domain data. We perform two rounds (R) of pseudo-labeling and training.

Algorithm 1 Training process

Input: \mathcal{S}, \mathcal{T}
Output: $F_{t,\sigma^2=1}^{(R=2)}, F_{t,\sigma^2=10}^{(R=2)}, F_s^{(R=2)}$
 train F with Eq. 2 (main paper)
 train the stochastic translation network with Eq. 11
 train $F_{t,\sigma^2=1}^{(R=0)}$ with Eq. 4 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 train $F_{t,\sigma^2=10}^{(R=0)}$ with Eq. 4 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, 10\mathbf{I})$
 train $F_s^{(R=0)}$ with Eq. 9 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
for $i \leftarrow 1$ to 2 **do**
 generate \hat{y} with Eq. 12
 train $F_{t,\sigma^2=1}^{(R=i)}$ with Eq. 8 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 train $F_{t,\sigma^2=10}^{(R=i)}$ with Eq. 8 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, 10\mathbf{I})$
 train $F_s^{(R=i)}$ with Eq. 10 (main paper), where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
end for

3 Additional results

We report results from additional ablation studies and class-wise IoU for ablation studies already presented in the main paper. We report results on GTA-to-Cityscapes using DeepLab-V2 with ResNet-101.

3.1 Quantitative results

Training objective of the source domain network.

As we mentioned in the main paper, for the source domain network we observed experimentally that we obtained better results by adding an entropy-based adversarial loss

\mathcal{L}_{adv} , to the output of source-domain network F_s when it is driven by translated target images. In Table 1 we report results obtained with and without the entropy-based adversarial loss (Eq. 9 in the main paper). Adding the entropy-based regularization improves performance for most classes.

Loss	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
\mathcal{L}_{CE}	90.2	38.0	81.2	29.1	16.2	24.4	23.7	15.5	84.0	38.8	78.5	56.9	24.0	85.0	36.4	47.0	0.3	31.8	26.8	43.6
$\mathcal{L}_{CE} + \mathcal{L}_{adv}$	90.5	39.4	82.0	29.0	21.4	23.6	28.6	17.8	83.9	38.2	79.8	56.9	26.0	85.1	32.2	44.1	3.8	31.5	30.1	44.4

Table 1: Better performance is achieved by adding an entropy-based regularization \mathcal{L}_{adv} to the output of source-domain network F_s when it is driven by translated target images.

Selection of r for pseudolabel generation.

As we mentioned in Sec. 2.2, we select the proportion of r based on the mIoU of the validation set. In Table 2 and Table 3 we provide the per-class IoU and mean IoU (mIoU) obtained on the validation set for different values of r in the first ($R=0$) and second ($R=1$) round of pseudo-labeling respectively. In both rounds the best performance is achieved for $r = 0.6$. In the second round of pseudolabeling the networks provide more confident predictions since the performance remains the same for almost all classes when $r \leq 0.5$.

r	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
1.0	92.1	47.8	84.3	36.5	27.9	31.5	36.6	24.5	85.4	41.2	81.6	61.4	30.1	86.3	37.6	47.3	1.3	28.7	32.7	48.2
0.8	97.4	69.8	93.4	52.0	42.5	47.1	55.2	37.7	94.2	53.5	91.1	78.2	40.2	94.4	47.4	55.7	2.2	40.9	55.4	60.4
0.7	98.2	72.1	95.9	63.4	51.5	55.2	64.0	42.5	96.5	65.0	93.7	86.9	51.7	96.4	57.0	61.2	2.7	53.1	65.8	67.0
0.6	98.5	70.2	96.3	73.4	57.5	54.1	58.7	33.3	97.2	77.4	93.8	90.8	60.3	97.5	65.5	69.6	2.5	65.8	71.9	70.2
0.5	98.6	59.4	96.5	79.0	55.2	44.6	54.3	20.7	97.6	81.4	93.8	92.2	61.0	97.9	70.6	74.1	1.3	73.3	70.3	69.6
0.4	98.6	50.1	96.6	76.5	42.2	44.6	54.6	20.6	97.7	82.3	93.8	92.3	61.1	97.9	70.6	74.1	0.5	73.4	70.5	68.3
0.3	98.6	50.1	96.6	76.6	38.3	44.7	54.6	20.6	97.7	82.3	93.8	92.3	61.1	97.9	70.6	74.1	0.2	73.4	70.5	68.1

Table 2: Per-class IoU and mean (mIoU) obtained using different values of r for class-wise confidence threshold selection in the first round ($R=0$) of pseudolabeling. We observe that $r = 0.6$ gives the best results.

Class-wise IoU for ablation studies reported in the main paper.

In Table 4 we report the per-class IoU obtained from deterministic and stochastic translation (mIoU results are reported in Table 1 of the main paper). In particular the comparison builds on directly on the ADVENT baseline [4]; the first two rows compare the

<i>r</i>	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
1.0	93.0	53.3	85.8	41.2	33.1	33.4	39.1	29.7	86.4	45.4	84.5	60.0	29.3	86.9	45.8	57.7	2.7	34.6	45.8	52.0
0.8	96.8	67.4	93.6	57.3	45.2	46.7	54.2	41.9	94.2	59.2	92.4	74.8	37.1	94.4	57.1	71.3	4.6	49.2	60.8	63.1
0.7	97.2	67.3	94.2	64.5	49.4	49.0	51.9	37.1	95.3	67.6	92.5	80.8	46.3	95.9	68.6	77.9	5.4	59.7	70.5	66.9
0.6	97.3	65.4	94.5	67.8	51.3	44.0	48.5	36.0	95.7	71.9	92.5	84.4	52.5	96.9	85.3	81.2	4.9	67.9	74.7	69.1
0.5	97.3	65.4	94.5	68.1	49.3	42.8	48.6	36.3	95.7	72.2	92.5	84.7	52.7	97.0	87.5	81.6	3.0	68.6	75.1	69.1
0.4	97.3	65.4	94.5	68.1	49.3	42.8	48.6	36.3	95.7	72.2	92.5	84.7	52.7	97.0	87.5	81.6	0.4	68.6	75.1	69.0
0.3	97.3	65.4	94.5	68.1	49.3	42.8	48.6	36.3	95.7	72.2	92.5	84.7	52.7	97.0	87.5	81.6	0.4	68.6	75.1	69.0

234
235
236
237
Table 3: Per-class IoU and mean (mIoU) obtained using different values of r for class-
wise confidence threshold selection in the second round (R=1) of pseudolabeling. We
observe that $r = 0.6$ gives the best results.

238
239
240 originally published and our reproduced numbers respectively. The third row shows the
241 substantial improvement attained by training the system of ADVENT using translated
242 images. The forth row reports our stochastic translation-based result. We observe a sub-
243 stantial improvement, that can be attributed solely to the stochasticity of the translation.
244 The last row shows that imposing a semantic consistency constraint as described in Eq.
245 5 (main paper) further improves the performance. In Table 5 we report the per-class
246 IoU obtained from multiple rounds R of pseudo-labeling and training (mIoU results are
247 provided in Table 3 of the main paper). Multiple rounds of pseudolabeling and training
248 yield improved performance.

Model	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
ADVENT	89.9	36.5	81.6	29.2	25.2	28.5	32.3	22.4	83.9	34.0	77.1	57.4	27.9	83.7	29.4	39.1	1.5	28.4	23.3	43.8
ADVENT *	87.2	38.5	78.2	25.9	24.6	30.4	36.3	21.7	84.0	28.7	76.7	60.1	28.8	80.0	28.0	45.2	0.7	19.7	19.9	42.9
ADVENT *+																				
CycleGAN*	91.9	51.5	83.1	30.8	23.6	32.0	32.1	24.3	83.8	38.5	82.3	58.7	28.5	84.1	33.3	35.9	0.6	21.7	20.0	45.1
Ours	90.2	37.6	84.1	33.0	25.1	30.1	36.8	28.4	83.8	36.1	82.2	58.1	29.6	84.6	34.4	45.4	1.0	26.2	30.8	46.2
Ours w/ L_{sem}	92.1	49.9	83.5	29.1	24.7	30.3	38.3	27.2	84.8	34.4	81.1	60.4	28.1	85.2	33.0	45.7	2.5	23.8	30.4	46.6

250
251
252
253
254
255
256
257
Table 4: GTA to Cityscapes UDA using stochastic translation: We train ADVENT using
258 synthetic images obtained from deterministic translation (CycleGAN) and stochastic
259 translation (Ours). We observe a clear improvement thanks to pixel-space alignment
260 based on stochastic translation. * denotes our retrained models.

261 3.2 Qualitative results

262 Diverse translation obtained using stochastic translation.

263 Fig. 2 shows diverse translations of images from the SYNTHIA source dataset to the

Model	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU
source (R=0)	90.5	39.4	82.0	29.0	21.4	23.6	28.6	17.8	83.9	38.2	79.8	56.9	26.0	85.1	32.2	44.1	3.8	31.5	30.1	44.4
target, $\sigma^2 = 1$ (R=0)	92.1	49.9	83.5	29.1	24.7	30.3	38.3	27.2	84.8	34.4	81.1	60.4	28.1	85.2	33.0	45.7	2.5	23.8	30.4	46.6
target, $\sigma^2 = 10$ (R=0)	90.9	43.0	83.4	30.6	29.3	30.6	34.1	27.1	84.4	36.2	79.9	60.6	29.5	84.5	32.5	40.3	3.1	29.2	26.4	46.1
Ens (R=0)	92.1	47.8	84.3	36.5	27.9	31.5	36.6	24.5	85.4	41.2	81.6	61.4	30.1	86.3	37.6	47.3	1.3	28.7	32.7	48.2
source (R=1)	92.1	48.4	84.3	29.5	30.5	35.9	26.5	85.4	42.9	82.1	59.8	29.6	85.5	38.2	52.9	3.4	32.7	37.3	49.1	
target, $\sigma^2 = 1$ (R=1)	92.1	47.5	85.1	38.3	29.4	32.9	35.4	32.1	85.9	46.8	81.7	60.5	30.4	86.6	35.7	51.1	4.4	34.9	41.0	50.1
target, $\sigma^2 = 10$ (R=1)	92.9	55.2	85.1	38.1	30.6	32.8	39.8	34.8	85.9	42.2	84.0	59.0	26.1	85.4	47.9	46.3	10.1	28.4	42.8	50.9
Ens (R=1)	93.0	53.3	85.8	41.2	33.1	33.4	39.1	29.7	86.4	45.4	84.5	60.0	29.3	86.9	45.8	57.7	2.7	34.6	45.8	52.0
source (R=2)	92.3	48.2	85.1	40.7	34.3	29.8	38.5	28.2	86.5	46.7	83.3	60.9	30.2	86.9	41.3	53.1	10.4	38.4	40.5	51.3
target, $\sigma^2 = 1$ (R=2)	93.3	56.5	85.9	41.0	33.1	34.8	43.8	43.8	86.6	46.5	82.5	61.1	30.4	87.0	39.7	50.7	8.8	34.9	46.8	53.0
target, $\sigma^2 = 10$ (R=2)	93.4	56.3	85.6	40.6	33.5	35.9	43.5	41.1	85.7	43.8	84.1	60.6	29.2	87.2	44.2	53.7	13.7	33.8	39.2	52.8
Ens (R=2)	93.4	55.8	86.4	44.4	36.1	34.6	45.0	39.8	86.9	48.0	84.4	61.7	30.9	87.7	44.9	55.9	11.1	38.4	45.4	54.3

Table 5: Ablation study on GTA → Cityscapes. Averaging the predictions (Ens) of a source network F_s , and two target networks F_t trained with different degrees of stochasticity (σ^2) in the translation allows to obtain robust pseudo-labels, while using multiple rounds R of pseudo-labeling and training improves the overall performance.

Cityscapes target dataset. Fig. 3 and Fig. 4 show diverse translations of images from the Cityscapes target dataset to the SYNTHIA and GTA source datasets respectively. We observe that stochastic translation generates diverse samples that capture more faithfully the data distribution of the source domain and preserve the content of the original image allowing us to obtain more robust pseudolabels for the target data.

Stochastic versus deterministic translation.

Fig. 5 shows stochastic and deterministic translation of images from the GTA source dataset to the Cityscapes target dataset while Fig. 6 shows stochastic and deterministic translation of images from the SYNTHIA source dataset to the Cityscapes target dataset. Fig. 7 shows stochastic and deterministic translation of images from the Cityscapes target dataset to the GTA source dataset while Fig. 8 shows stochastic and deterministic translation of images from the Cityscapes target dataset to the SYNTHIA source dataset. We observe that stochastic translation generates sharp samples of noticeable diversity compared to the deterministic translation that generates a single output.

Multiple rounds of pseudolabeling.

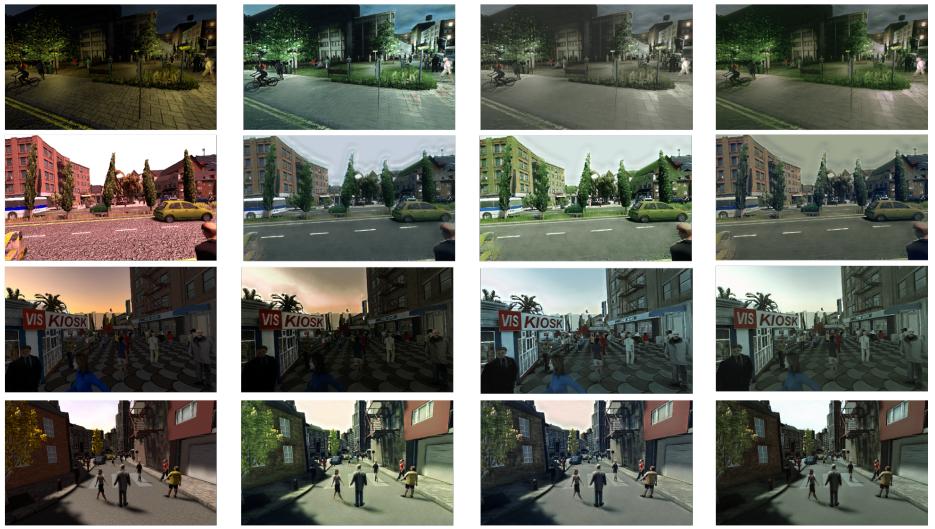
Fig. 9 shows the pseudo-labels obtained from the first (R=0) and second (R=1) round of pseudolabeling. We observe that the pseudolabels we obtained in the second round are more accurate allowing us to train more accurate models in the last round of training.

Robust pseudolabeling through ensembling.

Fig. 10 shows the pseudo-labels obtained by averaging the predictions of two target networks $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$ and a one source network F_s . Averaging the predictions allows us to generate more accurate pseudolabels.

Ensembling for improved segmentation performance.

Fig. 11 shows the predictions obtained by averaging the predictions of two target networks $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$ and a one source network F_s . Averaging the predictions allows us to further improve performance by better distinguishing similar structures (e.g., road, sidewalk) and identifying small objects.



316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
Fig. 2: Diverse translations of images from the SYNTTHIA source dataset to the
332 Cityscapes target dataset: we observe that even though the content and pixel
333 semantics stay intact, we generate diverse variants of the same scene, effectively capturing
334 more faithfully the data distribution in the target domain.
335

336 337 Qualitative comparison of the segmentation results.

338 Fig. 12 shows results segmentation results obtained by our method and DPL [1]. Our
339 method generates better predictions that are closer to the ground-truth.
340

341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359

315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359

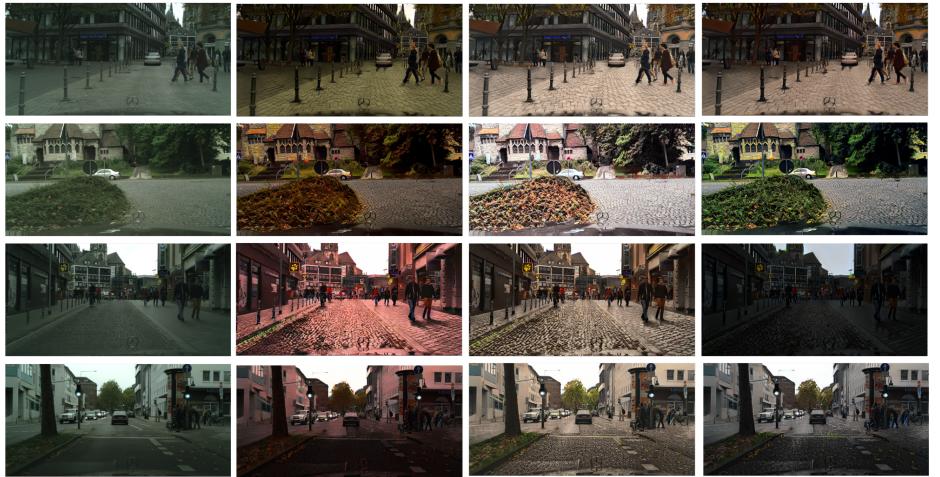


Fig. 3: Diverse translations of images from the Cityscapes target dataset to the SYNTHIA source dataset: we observe that even though the content and pixel semantics stay intact, we generate diverse variants of the same scene, effectively capturing more faithfully the data distribution of the source domain. This allows us to generate more robust pseudolabels.



Fig. 4: Diverse translations of images from the Cityscapes target dataset to the GTA source dataset: we observe that even though the content and pixel semantics stay intact, we generate diverse variants of the same scene, effectively capturing more faithfully the data distribution of the source domain. This allows us to generate more robust pseudolabels.

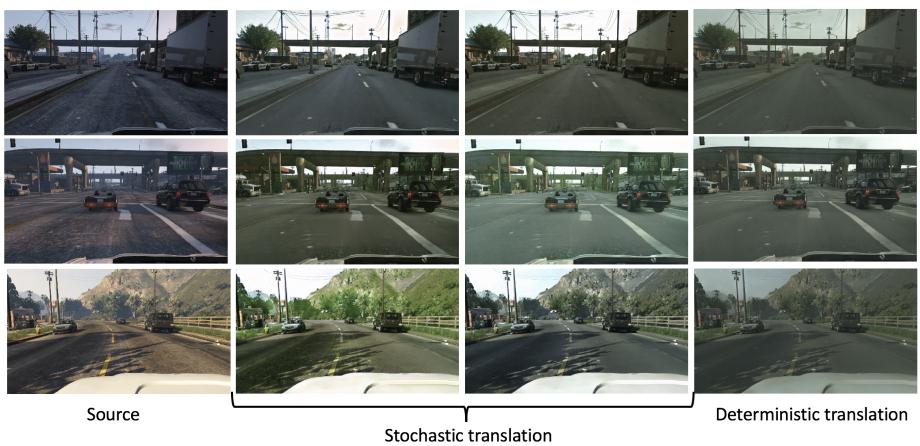


Fig. 5: Stochastic and deterministic translation of images from the GTA source dataset to the Cityscapes target dataset.



Fig. 6: Stochastic and deterministic translation of images from the SYNTHIA source dataset to the Cityscapes target dataset.

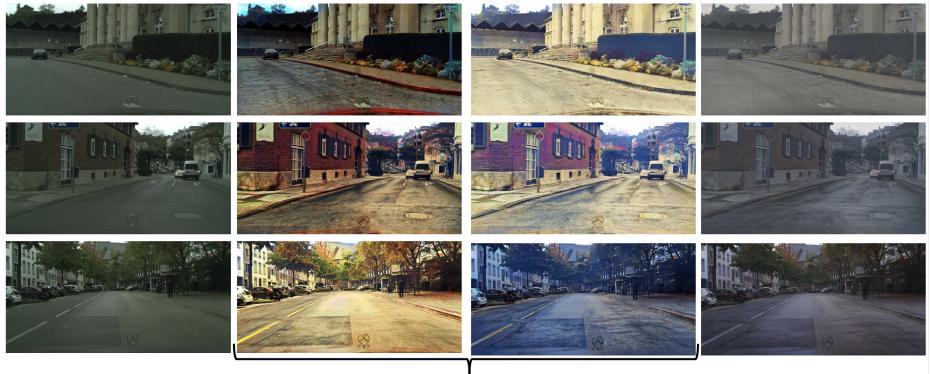


Fig. 7: Stochastic and deterministic translation of images from the Cityscapes target dataset to the GTA source dataset.



Fig. 8: Stochastic and deterministic translation of images from the Cityscapes target dataset to the SYNTHIA source dataset.



Fig. 9: Visualization of pseudolabels obtained from the first ($R=0$) and second ($R=1$) round. Pseudolabels obtained in the second round are more accurate allowing us to train more robust models in the last round of training.

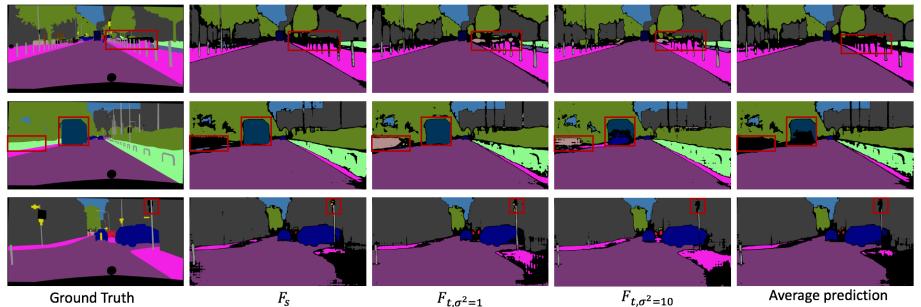


Fig. 10: Visualization of pseudolabels obtained by averaging the predictions of two target networks $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$ and a one source network F_s . Averaging the predictions allows us to generate more accurate pseudolabels.

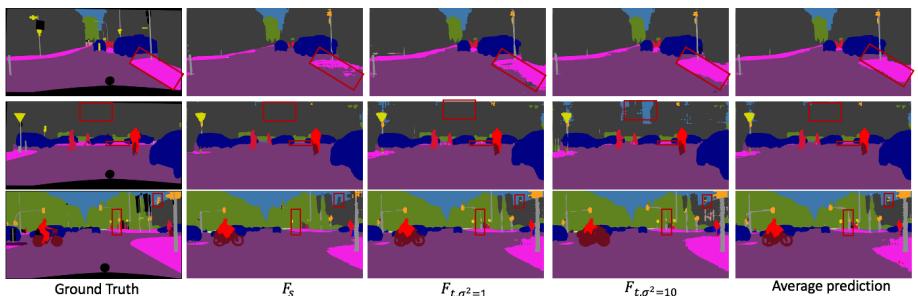


Fig. 11: Visualization of results obtained by averaging the predictions of two target networks $F_{t,\sigma^2=1}$, $F_{t,\sigma^2=10}$ and a one source network F_s . Averaging the predictions allows us to further improve performance by better distinguishing similar structures (e.g., road, sidewalk) and identifying small objects.

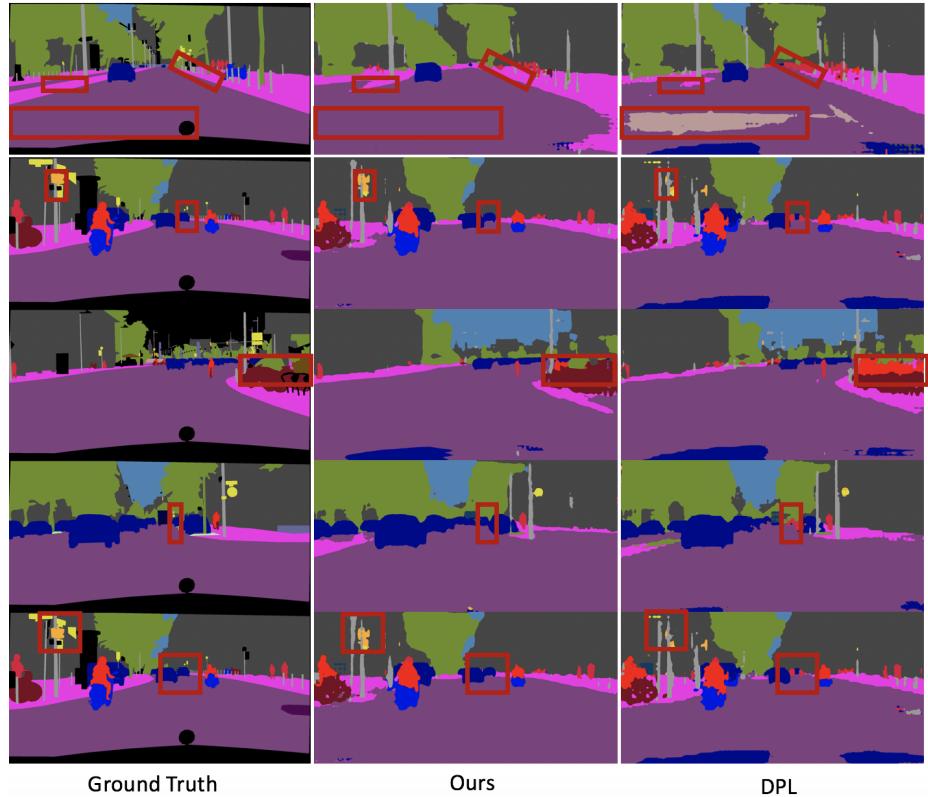


Fig. 12: Qualitative comparison of our method with DPL [1]. Our method generates better predictions that are closer to the ground-truth.

References

- 585 1. Cheng, Y., Wei, F., Bao, J., Chen, D., Wen, F., Zhang, W.: Dual path learning for domain
586 adaptation of semantic segmentation. In: ICCV (2021) 8, 13
- 587 2. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normal-
588 ization. In: ICCV (2017) 1
- 589 3. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image trans-
590 lation. In: ECCV (2018) 1, 3
- 591 4. Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P.: ADVENT: Adversarial entropy minimiza-
592 tion for domain adaptation in semantic segmentation. In: CVPR (2019) 5
- 593 5. Zou, Y., Yu, Z., Vijaya Kumar, B.V.K., Wang, J.: Unsupervised domain adaptation for seman-
594 tic segmentation via class-balanced self-training. In: ECCV (2018) 3