

Casting Geometric Constraints in Semantic Segmentation as Semi-Supervised Learning

Sinisa Stekovic¹Friedrich Fraundorfer¹Vincent Lepetit^{2,1}¹Institute for Computer Graphics and Vision, Graz University of Technology, Graz, Austria²Université Paris-Est, École des Ponts ParisTech, Paris, France

{sinisa.stekovic, fraundorfer, lepetit}@icg.tugraz.at

Abstract

We propose a simple yet effective method to learn to segment new indoor scenes from video frames: State-of-the-art methods trained on one dataset, even as large as the SUNRGB-D dataset, can perform poorly when applied to images that are not part of the dataset, because of the dataset bias, a common phenomenon in computer vision. To make semantic segmentation more useful in practice, one can exploit geometric constraints. Our main contribution is to show that these constraints can be cast conveniently as semi-supervised terms, which enforce the fact that the same class should be predicted for the projections of the same 3D location in different images. This is interesting as we can exploit general existing techniques developed for semi-supervised learning to efficiently incorporate the constraints. We show that this approach can efficiently and accurately learn to segment target sequences of ScanNet and our own target sequences using only annotations from SUNRGB-D, and geometric relations between the video frames of target sequences.

1. Introduction

Semantic segmentation of images provides high-level understanding of a scene, useful for many applications such as robotics and augmented reality. Recent approaches can perform very well [23, 15, 45, 5].

In practice, however, it is difficult to generalize from existing datasets to new scenes. In other words, it is a challenging task to obtain good segmentation of images that do not belong to the training datasets. To demonstrate this, we trained a state-of-the-art segmentation method DeepLabV3+ [5] on the SUNRGB-D dataset [37], which is made of more than 5000 training images of indoor scenes. Fig. 1a shows the segmentation we obtain when we attempt to segment a new image, which does not belong to the dataset. The performance is clearly poor, showing that the



Figure 1: (a) Even the state-of-the-art method DeepLabV3+ trained with training data from SUNRGB-D makes many mistakes when segmenting an image outside the SUNRGB-D dataset. (b) After exploiting geometric constraints on an unlabeled sequence of the new scene, our semi-supervised S4-Net approach predicts much better segmentations.

SUNRGB-D dataset was not sufficient to generalize to this image, despite the size of the training dataset.

To make semantic segmentation more practical and to break this dataset bias, one can exploit geometric constraints [25, 31, 27, 14], in addition to readily available training data such as the SUNRGB-D dataset. We introduce an efficient formalization of this approach, which relies on the observation that geometric constraints can be introduced as standard terms from the semi-supervised learning literature. This results in an elegant, simple, and powerful method that can learn to segment new environments from video frames, which makes it very useful for applications such as robotics and augmented reality.

More exactly, we adapt a general technique for semi-supervised learning that consists of adding constraints on pairs of unlabeled training samples that are close to each other in the feature space, to enforce the fact that such two samples should belong to the same category [22, 39, 2]. This is very close to what we want to do when enforcing geometric constraints for semantic segmentation: Pairs of unlabeled pixels that correspond to the same physical 3D point would be labeled with the same category. In practice,

to obtain the geometric information needed to enforce the constraints, we can rely on the measurements from depth sensors or train a network to predict depth maps as well, using recent techniques for monocular image reconstruction.

In contrast to previous methods exploiting geometric constraints for semantic segmentation, our method introduces several novelties. Comparing to [25], our approach applies geometric constraints to completely unlabeled scenes. Furthermore, when compared to [31, 27, 14], which use simple label fusion to segment given target sequence, our approach can generalize from a representation of one target sequence from the target scene to segmenting unseen images of the target scene. We demonstrate this aspect further in the evaluation section.

In short, our contribution is to show that semi-supervised learning is a simple yet principled and powerful way to exploit geometric constraints in learning semantic segmentation. We demonstrate this by learning to annotate sequences of the ScanNet [7] dataset using only annotations from the SUNRGB-D dataset. We also demonstrate effectiveness of the proposed method through the semantic labeling of our own newly generated sequence unrelated to SUNRGB-D and ScanNet.

In the rest of the paper, we discuss related work, describe our approach, and present its evaluation with quantitative and qualitative experiments together with an ablation study.

2. Related Work

In this section, we discuss related work on the aspects of semantic segmentation, domain adaptation, general semi-supervised learning, and also recent methods for learning depth prediction from single images, as they also exploit geometric constraints similar to our approach. Finally, we discuss similarities and differences with other works that also combine segmentation and geometry.

2.1. Supervised Semantic Segmentation with Deep Networks

The introduction of deep learning made a large impact on performance of semantic segmentation. Fully Convolutional Networks (FCNs) [23] allow segmentation prediction for input of arbitrary size. In this setting, standard image classification task networks [36, 16] can be used by transforming fully-connected layers into convolutional ones. FCNs use deconvolutional layers that learn the interpolation for upsampling process. Other works including SegNet [3] and U-Net [33] rely on similar architectures. Such works have been applied to a variety of segmentational tasks [33, 1, 29].

Recent methods address the problem of utilizing global context information for semantic segmentation. PSPNet [45] proposes to capture global context information through a pyramid pooling module that combines features

under four different pyramid scales. DeepLabV3+ [5] uses atrous convolutions to control response of feature maps and applies atrous spatial pyramid pooling for segmenting objects at multiple scales. In our experiments, we apply our approach to both DeepLabV3+ and PSPNet to demonstrate it generalizes to different network architectures. In principle, any other architecture could be used instead.

2.2. Semi-Supervised Learning with Deep Networks

Availability of ground truth labels is often the main limitation of supervised methods in practice. In contrast, semi-supervised learning is a general approach aiming at exploiting both labeled and unlabeled or weakly labeled training data. Some approaches rely on adversarial networks to measure the quality of unlabeled data [8, 10, 21, 17]. More in line with our work are the popular consistency-based models [22, 39, 2]. These methods enforce the model output to be consistent under small input perturbations. As explained in [2], consistency-based models can be viewed as a student-teacher model: To measure consistency of model f , or the student, its predictions are compared to predictions of a teacher model g , a different trained model, while at the same time applying small input perturbations.

Π -model [22] is a recent method using a consistency-based model where the student is its own teacher, *i.e.* $f = g$. It relies on a cross-entropy loss term applied to labeled data only and an additional term that penalizes differences in predictions for small perturbations of input data. Our semi-supervised approach is closely related to the Π -model but relies on geometric consistency instead of enforcing consistent predictions for different input perturbations.

As pixel-level annotations, required for semantic segmentation tasks, are typically very time consuming to obtain, weakly-supervised methods become very interesting options for further increasing the amount of training samples. One way of obtaining more training data is through image-level annotations or bounding boxes. [30] demonstrates that a network trained with large number of such weakly-supervised samples in combination with small amount of samples with pixel-level annotations achieves comparable results to a fully supervised approach. Given image-level annotations rather than pixel-level annotations, [41] generates dense object localization maps which are then utilized in a weakly- or semi-supervised framework to learn semantic segmentation. Our geometric constraints can be seen as a form of weak supervision but instead of weak labels our approach relies only on weak constraints enforcing consistent annotations for 3D points of the scene.

2.3. Domain Adaptation For Semantic Segmentation

Domain adaptation has been studied for the field of semantic segmentation. One can argue that overcoming the

dataset bias is closely related to the field of domain adaptation. In the context of semantic segmentation, such approaches usually leverage the possibility of using inexhaustive synthetic datasets for improving performance on real data [24, 28, 38, 4]. However, as further explained in [38], due to large domain gap between real and synthetic images, such domain adaptation methods easily overfit to synthetic data and can fail to generalize to real images.

Very recently, in terms of domain adaptation approaches that rely on real data only, Kalluri *et al.* [20] proposed a unified segmentation model for different target domains that minimizes supervised loss for labeled data of the target domains and exploits visual similarity between unlabeled data of the domains. Results indicate increase in performance for all of the target domains. However, such approach still requires labeled images for all of the target domains. Here, we focus on adapting the source domain to a related target domain for which no labeling is available.

2.4. Geometric Constraints and Label Propagation

Geometry in semantic segmentation has already been considered for purpose of semantic mapping. [25] trains a CNN by propagating manual hand-labeled segmentations of frames to new frames by warping. In contrast, we do not need any manual annotations for the target sequences of the scene. SemanticFusion [27] uses a pre-trained CNN together with ElasticFusion SLAM method [42], and merges multiple segmentation predictions from different viewpoints. [31, 14] rely on using a pre-trained CNN together with 3D reconstruction methods and improve accuracy over initial segmentations. However, these approaches are applied to a CNN with fixed parameters and rely on geometric constraints during inference time. In contrast, our method uses geometric constraints to improve single-view segmentation predictions for the target scene and afterwards requires only color information for segmenting unseen images of the scene.

2.5. Single-View Depth Estimation

Because of view warping, our approach is also related to recent work on unsupervised single-view depth estimation. Both Zhou *et al.* [46] and Godard *et al.* [12] proposed an unsupervised approach for learning depth estimation from video data. This is done by learning to predict a depth map so that a view can be warped into another one. This research direction became quickly popular, and has been extended since by many authors [44, 26, 40, 13].

Our work is related to these methods as it also introduces constraints between multiple views, by using warping. We demonstrate that this type of constraints can be utilized for the task of semantic segmentation.

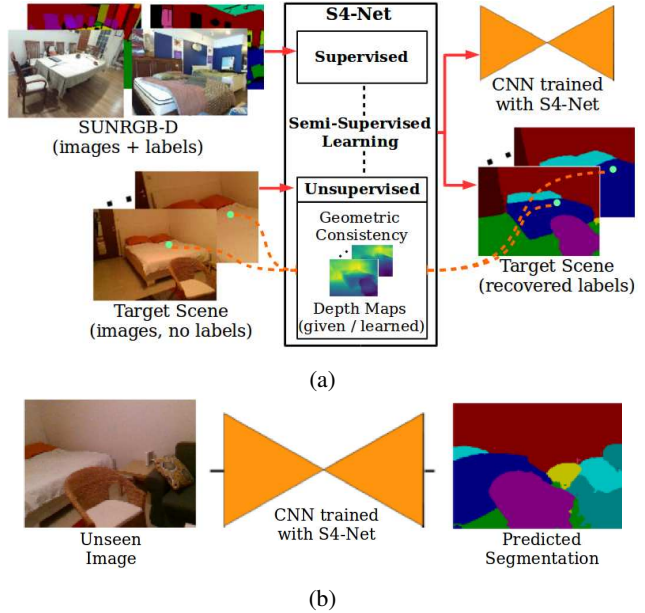


Figure 2: Method overview. (a) Our S4-Net approach combines supervised data from SUNRGB-D and an image sequence from a target scene without any annotations. By exploiting geometric constraints of the target image sequence, we obtain a network with high performance for the target scene, and labels for the target sequence. (b) After being trained by S4-Net, the segmentation network can be applied to unseen images of the target scene with much better performance.

3. Approach Overview

For the rest of the paper, we refer to our Semi-Supervised method for Semantic Segmentation as S4-Net. We assume that we are given a dataset of color images and their segmentations:

$$\mathcal{S} = \{e_i = (I_i, A_i)\}_i,$$

where I_i is a color image, and A_i is the corresponding ground truth segmentation. In practice we use the SUNRGB-D dataset. Based on these annotations, we would like to train a segmentation model $f()$ for a new scene given a sequence of registered frames, for which no labels are known *a priori*:

$$\mathcal{U} = \{e_j = (I_j, D_j, T_j)\}_j,$$

where I_j is a color image, D_j is the corresponding depth map, and T_j the corresponding camera pose. As a direct result of S4-Net, we obtain automatic annotations for sequence \mathcal{U} . Additionally, the output of S4-Net is a trained network $f()$. At test time, network $f()$ trained with S4-Net can then be used to predict correct segmentation for new images of the scene. We present the method overview in Fig. 2.

3.1. Semi-Supervised Learning and Geometric Consistency

We optimize the parameters Θ of $f()$ by minimizing the semi-supervised loss term:

$$L = L_S + \lambda L_G, \quad (1)$$

where L_S is a supervised loss term and L_G is a term that exploits geometric constraints. In practice, we set the discount factor λ for all experiments to the same value. L_S is a standard term for supervised learning of semantic segmentation:

$$L_S = \sum_{e \in \mathcal{S}} l_{\text{WCE}}(f(I(e); \Theta), A(e)), \quad (2)$$

where l_{WCE} is the weighted cross-entropy of segmentation prediction $f(I; \Theta)$ relative to manual annotation A . The class weights are calculated using median frequency balancing [11] to prevent overfitting to most common classes.

L_G exploits geometric constraints to enforce consistency between predictions for images taken from different viewpoints:

$$L_G = \sum_{e \in \mathcal{U}} l_{\text{CE}}(f(I(e); \Theta), \text{Merge}(\text{Warp}(f(I(e'); \Theta')))), \quad (3)$$

where $\mathcal{N}(e)$ is a subset of \mathcal{U} containing samples with a viewpoint that overlaps with the view point of e . $\text{Warp}(S)_{e' \rightarrow e}$

function warps segmentation S from frame e' to frame e . We give more details on this warp operation in Section 3.2. Merge function merges given neighbouring views by first $\text{Merge}_{e' \in \mathcal{N}(e)}$ summing the pixelwise probabilities and then performing argmax operation to obtain the final pixelwise labels.

We consider prediction $f(I(e'); \Theta')$ as a teacher prediction and, similarly to the Π -model [22], it is treated as a constant when calculating the update of the network parameters. Parameters Θ' are updated every 100 iterations to equal parameters Θ . We found that this step helps to further stabilize the learning process. l_{CE} is the standard cross-entropy loss function that compares the predicted segmentations. We found empirically that using weighted cross-entropy tends to converge to solutions with incorrect segmentations.

3.2. Segmentation Warping

We base our warping function Warp on the inverse warping method used in [46]. For a 2D location p in homogeneous coordinates of a target sample e , we find the corresponding location p' of the source sample e' using:

$$p' = K T_{e \rightarrow e'} d K^{-1} p, \quad (4)$$

where K is the intrinsic matrix of the camera, $T_{e \rightarrow e'}$ is the relative transformation matrix between the target and the

source samples, and d is the predicted depth value at location p . Since p' value lies in general between different image locations, we use the differentiable bilinear interpolation from [18] to compute the final projected value from the 4 neighbouring pixels. This transformation is applied to the segmentation probabilities predicted by the network.

In practice, not every pixel in the target sample has a correspondent pixel in the source sample. This can happen as depth information is not necessarily available for every pixel when using depth cameras, and since some pixels in the target sample may not be visible in the source sample, because they are occluded or, simply, because they are not in the field of view of the source sample. If the difference between the depths is larger than a threshold value, this means that the pixel is occluded and does not correspond to the same physical 3D point. We simply ignore the pixels without correspondents in the loss function. Additionally, we ignore the pixels that are located near the edges of the predicted segments: Segmentation predictions in these regions tend to be less reliable and, for such regions, insignificant errors in one view can easily induce significant errors in other views because of the different perspectives, as shown in Fig. 3

3.3. S4-Net with Depth Prediction

For sequences captured with an RGB camera, depth data is not available, and we rely on predicted depths to enforce geometric constraints. If the supervised dataset \mathcal{S} also includes ground truth depths, we can introduce additional loss terms to learn depth estimation:

$$L_D = L_{DS} + \lambda_D L_{DG}, \quad (5)$$

where L_{DS} is a supervised depth loss term and L_{DG} is a semi-supervised term that exploits geometric constraints in the depth domain. λ_D is a weighting factor. L_{DS} is the absolute difference loss term:

$$L_{DS} = \sum_{e \in \mathcal{S}} |f_d(I(e); \Theta_d) - \hat{D}(e)|, \quad (6)$$

where $f_d(I; \Theta_d)$ is the depth prediction for network parameters Θ_d and \hat{D} is the ground truth depth map. Term L_{DG} corrects noisy depth predictions for the target scene \mathcal{U} through geometric constraints only:

$$L_{DG} = \sum_{e \in \mathcal{U}} \sum_{e' \in \mathcal{N}(e)} L_{INT}((I(e), \text{Warp}(I(e')))), \quad (7)$$

where L_{INT} is a loss term comparing pixelwise intensities together with the structure similarity loss term from recent literature on monocular depth prediction [46, 12, 44, 26, 40, 13]. We apply this term only to the target image pixels where the predicted segmentations are consistent with each other and further away from segmentation borders:

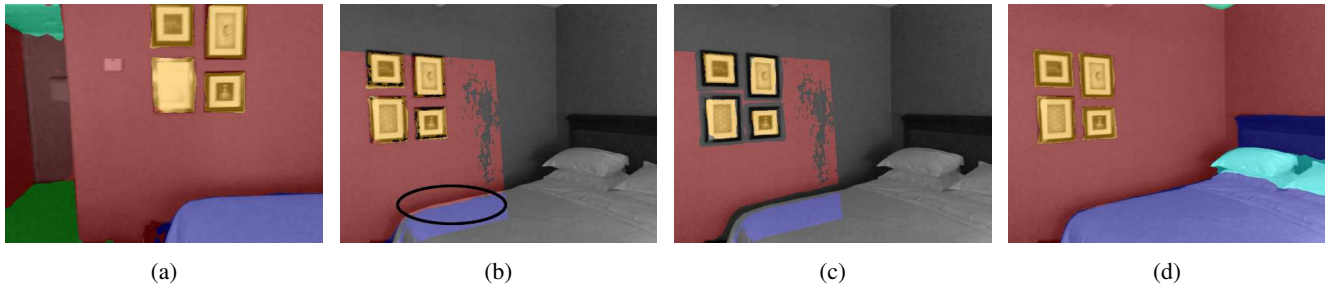


Figure 3: Process of warping source segmentation prediction in (a) to the corresponding target view. Encircled region in (b) demonstrates that warping in boundary regions of segmentation predictions can induce errors in the target view. In this case, the warped segmentation prediction falsely assigns wall labels to the edges of the bed region. Hence, we introduce a segmentation boundary mask to resolve this issue in (c). As direct result, S4-Net is able to recover quality segmentations in the affected boundary regions in (d).

We found that such mask helps regularize depth predictions for occluded regions of the image. We explain this further in supplementary material. For enforcing geometric constraints on semantic segmentation with term L_G , we consider the depth prediction as a constant when calculating the update of the network parameters.

3.4. Network Architecture

We use DeepLabV3+ [5] or PSPNet [45] as network $f()$ in our experiments. In both cases, as the base network, we use ResNet-50 [16] pre-trained on the ImageNet classification task [9]. However, S4-Net is not restricted to a specific type of architecture and could be applied to other architectures as well. When predicting depth maps, the encoder is shared between the depth network and the segmentation network. The depth decoder has the same architecture as the segmentation decoder, but they do not share any parameters. We show further details on network initialization and training procedure in supplementary material.

4. Evaluation

We evaluate S4-Net on the task of learning to predict semantic segmentation from color images for a target scene. The task for the network is to learn segmentations for a target scene without any knowledge about the ground truth labels for the scene. Hence, S4-Net requires an uncorrelated annotated dataset to obtain prior knowledge about the segmentation task that it needs to perform. Additionally, in order to learn accurate segmentations for the target scene, it utilizes frame sequences of that scene. By exploiting geometric constraints between the frames, S4-Net learns to predict segmentations across the target scene.

Datasets. In all of our evaluations, we use the SUNRGB-D dataset [37], consisting of 5285 annotated RGB-D images for training, as the supervised dataset \mathcal{S} and perform mirroring on these samples to augment the super-

vised data. The SUNRGB-D dataset is a collection made of an original dataset and additional datasets previously published [35, 19, 43]. The images are manually segmented into 37 object categories that are typical for an indoor scenario. The full list of object categories is given in supplementary materials. First, we evaluate S4-Net on scenes from the ScanNet dataset [7]. Second, we show that S4-Net is general by applying it to our own data. Finally, we show that enforcing geometric constraints through depth predictions can be used to learn quality segmentation predictions for the target scene.

4.1. Evaluation on ScanNet

As previously discussed, we use SUNRGB-D for the supervised training data only. Therefore, for our first experimental setup, we evaluate S4-Net on 6 scenes from the ScanNet dataset [7] to demonstrate the generalization aspect of S4-Net, as it is the scenario that motivates our work. These scenes represent different indoor scenarios, including apartment, hotel room, public bathroom, large kitchen space, lounge area, and study room. Even though the RGB-D sequences in the ScanNet dataset are annotated and can be mapped to our desired segmentation task, we utilize these annotations only to validate our results.

Data Split. Intentionally, we choose scenes from the ScanNet dataset which were scanned twice during the creation of the dataset. The first scan of each scene is utilized during training while the second scan is used for validation purposes only. We refer to these independently recorded scans as “Scan 1” and “Scan 2”. During training, we use the registered RGB-D sequence from “Scan 1” when applying geometric constraints. For evaluation, we additionally validate performance on “Scan 2” for which the camera follows a different pathway.

In this experiment, we show that our network, trained with geometric constraints from a target scene of ScanNet and supervised data from SUNRGB-D, notably improves

	“Scan 1” (ScanNet) (Unlabeled images during training)				“Scan 2” (ScanNet) (Excluded from training)			
	pix_acc	mean_acc	mIOU	fwIOU	pix_acc	mean_acc	mIOU	fwIOU
DeepLabV3+ network architecture								
Supervised baseline	0.765	0.634	0.533	0.692	0.772	0.651	0.544	0.697
<i>S4-Net</i>	0.803	0.687	0.59	0.733	0.803	0.691	0.593	0.732
<i>S4-Net with depth dred.</i>	0.794	0.679	0.581	0.724	0.797	0.684	0.586	0.725
PSPNet network architecture								
Supervised baseline	0.727	0.597	0.486	0.644	0.737	0.61	0.499	0.654
<i>S4-Net</i>	0.781	0.648	0.539	0.701	0.78	0.652	0.546	0.699

Table 1: Quantitative evaluation on the target scenes from ScanNet. We include results averaged over the target scenes used during experiments. The results for “Scan 1” show significant improvements for the images where we applied S4-Net. Furthermore, the segmentation accuracy for “Scan 2” indicates that our trained network brings similar improvements over the supervised baseline for the images that were not utilized by S4-Net during training. Our experiments also show that S4-Net can be applied to different segmentation networks with a significant gain in accuracy in comparison to the supervised baselines and the results with depth prediction show a comparable increase in performance over supervised baseline.

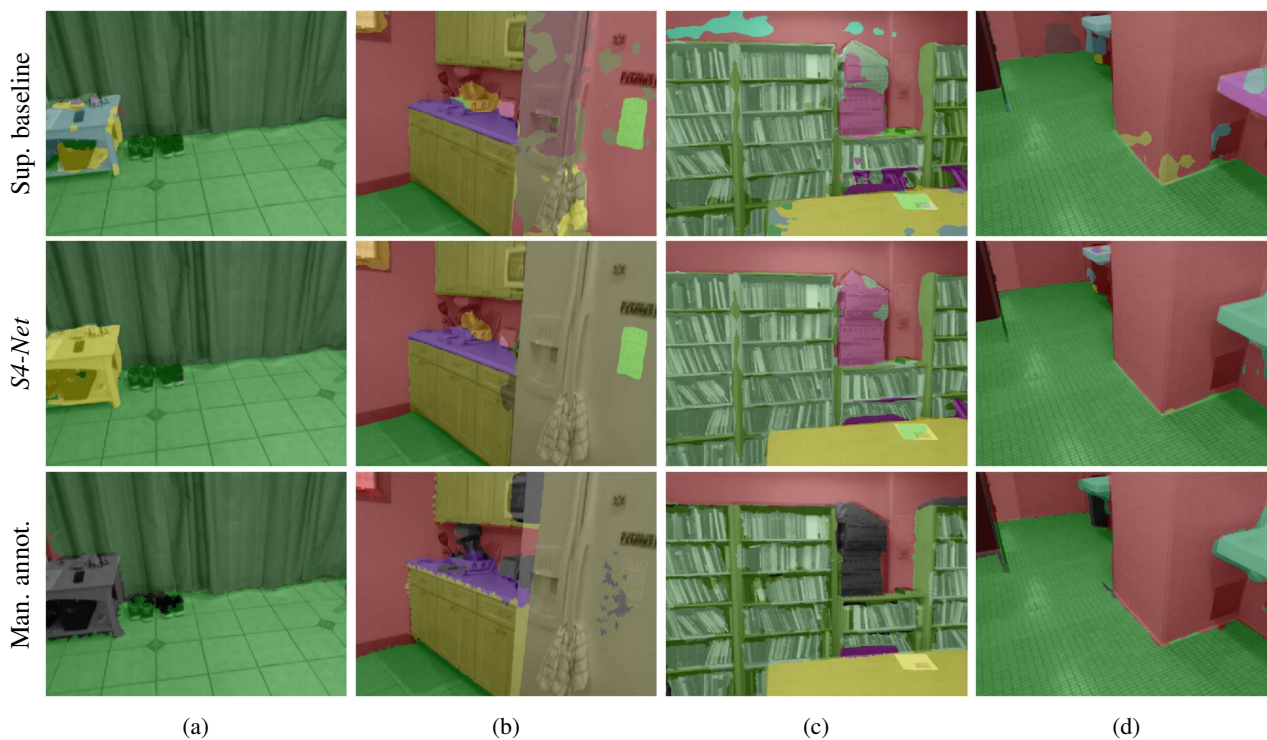


Figure 4: Qualitative results on unseen images from “Scan 2” of the target scenes from ScanNet, DeepLabV3+ network architecture. As S4-Net does not rely on manual annotations for the target scenes, it predicts segmentations that are sometimes more accurate than manual annotations. More specifically, it correctly segments the otherwise unlabelled table and box regions in (a) and (c), and in regions with wrong manual annotations it correctly predicts paper segments in (b) and (c).

performance over our supervised baseline for the target scene. This is true for all of our experimental scenes in ScanNet in regions where our supervised baseline already provides a certain level of generalization for some of the viewpoints. For “Scan 1” we measure a significant increase

in performance. To further demonstrate different use cases of S4-Net, we show that the network fine-tuned for “Scan 1” predicts high-grade segmentations also for “Scan 2” of the scene. The two scans are recorded independently which results in different camera paths for the recordings. As there

are no direct neighboring frames between the two independently recorded scans, the results on “Scan 2” demonstrate the ability of the S4-Net trained network to generalize to independently recorded scans of the same scene.

In our quantitative evaluations in Table 1, we present results averaged over all of our experimental scenes. We observe that the S4-Net approach clearly overcomes the dataset bias correlated with the supervised approach as it demonstrates superior performance on the target scene in comparison to its supervised baseline. Not only does the performance increase for “Scan 1” but we also observe that the increase in performance for the images of “Scan 2” is as significant. Furthermore, by observing performance for different network architectures, we show that S4-Net can be applied to arbitrary segmentation network architectures.

We further demonstrate the benefits of our approach in Fig. 4 where we show some qualitative results. We observe that, in areas where the supervised approach predicts very noisy predictions, our approach predicts consistent segmentations. This is the indicator that confident segmentation predictions are propagated to less confident viewpoints, and not the other way around.

Our experiments on ScanNet demonstrate that S4-Net is useful for different practical applications. First, the evaluations on “Scan 1” show that the approach is applicable for the use case of automatically labelling indoor scenes. The second application is that, once the network has converged for the target sequence, we can reliably segment new images of the scene without the need for the depth data.

4.2. Evaluation on Additional Scene

So far we have demonstrated that S4-Net works well for the chosen scenes from the ScanNet dataset. To show that S4-Net generalizes well, we also evaluate it on our own data. For this purpose, we captured and registered a living room area using an Intel[®]RealSense[™] D400 series depth camera¹ and registered the scene using an implementation of a scene reconstruction pipeline [6] from Open3D library [32]. We refer to this scene as the “Room” dataset.

Data Split. In line with the ScanNet experiments, we scanned the “Room” scene twice. “Scan 1” contains roughly 6000 training images. For evaluation purposes, we then sampled 20 images from “Scan 1” that capture different viewpoints of the scene, and we manually annotated them using the LabelMe annotation tool [34]. We annotated 20 additional images from “Scan 2” that was recorded independently of “Scan 1” to further demonstrate the aspect of generalization across the scene.

Table 2 gives the results of our quantitative evaluation. We again observe a significant increase in performance over the supervised baseline approach for images of “Scan 1” and “Scan 2”. Our qualitative evaluations in Fig. 5

¹<https://software.intel.com/en-us/realsense/d400>

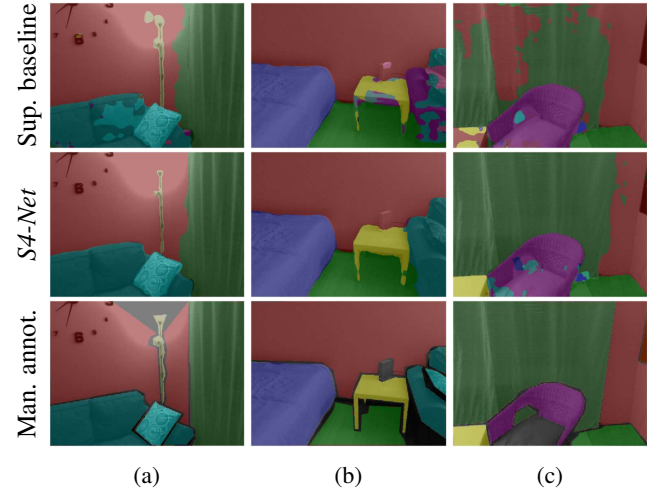


Figure 5: Qualitative evaluation on unseen images from “Scan 2” of the “Room” scene for the DeepLabV3+ network architecture. The supervised baseline already predicts high quality segmentations across the scene. However, the supervised baseline still predicts noisy or incorrect segmentations for some view points, especially due to partial visibility of objects, for example the sofa in (b) and the table in (c). S4-Net demonstrates notable improvements in these regions.

show many overall improvements. Even though our supervised baseline might predict quality segmentations for specific viewpoints, for other viewpoints it fails completely as these data samples are not presented well throughout the SUNRGB-D dataset. In contrast, the S4-Net approach preserves quality segmentations in such regions. This further proves that the usage of geometric constraints is, indeed, a very powerful method for transferring knowledge from the supervised baseline to a new scene.

4.3. Evaluation of S4-Net with Depth Prediction

Furthermore, we evaluate the aspect of using depth predictions for enforcing geometric constraints. As SUNRGB-D also contains depth ground truth data, it provides supervision for both the depth network and the segmentation network in this scenario. When enforcing geometric constraints for the target scenes, warping between different view points is performed by using depth predictions instead of the ground truth depth images. For this experiment, we found empirically that setting λ_D to 0.1 achieved satisfying quality for segmentation and depth predictions. In case of PSPNet, due to low accuracy of initial depth predictions, we excluded this part in our evaluations.

Our quantitative evaluations in Table 1 demonstrate comparable performance to S4-Net with depth ground truth for the ScanNet scenes. Similarly, in Table 2 we observe that

	“Scan 1” (“Room”) (Unlabeled images during training)				“Scan 2” (“Room”) (Excluded from training)			
	pix_acc	mean_acc	mIOU	fwIOU	pix_acc	mean_acc	mIOU	fwIOU
DeepLabV3+ network architecture								
Supervised baseline	0.89	0.757	0.699	0.827	0.817	0.726	0.66	0.76
<i>S4-Net</i>	0.934	0.846	0.799	0.884	0.938	0.75	0.719	0.906
<i>S4-Net with depth pred.</i>	0.922	0.801	0.757	0.868	0.911	0.755	0.712	0.867
PSPNet network architecture								
Supervised baseline	0.862	0.673	0.597	0.788	0.728	0.629	0.497	0.651
<i>S4-Net</i>	0.888	0.723	0.645	0.817	0.847	0.708	0.604	0.782

Table 2: Quantitative evaluation on the “Room” scene. Similarly to our experiments on ScanNet, S4-Net demonstrates significant performance increase for both the images from “Scan 1” and “Scan 2” of the scene. We observe improvements for different network architectures and also when using S4-Net with depth prediction network.

S4-Net with depth predictions shows significant improvements for the “Room” scene in comparison to the supervised baseline. In our qualitative results in Fig. 6 we visualize the results of S4-Net with depth predictions for the ScanNet scenes. Even though one would expect that noisy depth predictions considerably decrease the quality of geometric constraints, S4-Net still demonstrates quality improvements in this scenario that is comparable to our results when using ground truth depth for enforcing geometric constraints. Even though initial depth predictions for the supervised baseline are noisy, S4-Net also learns better depth predictions for the target scene. Hence, geometric constraints on semantic segmentation improve during training enabling convergence for S4-Net. For unseen “Scan 2” sequences from ScanNet, the Root Mean Square (RMS) error drops from 0.61 to 0.4 on average after applying S4-Net. For “Scan 2” images from “Room” scene, the average RMS error drops from 0.58 to 0.49. We show further quantitative and qualitative results on depth predictions in supplementary material.

5. Conclusion

We showed that semi-supervised learning is a good theoretical framework for enforcing geometric constraints and for adapting semantic segmentation to new scenes. We also investigated a potential problem which could appear with such semi-supervised constraints on non-annotated sequences. It would be possible that the learning may assign labels which are consistent among views, but wrong. Our experiments have shown that this is only very rarely the case. Instead, the semi-supervised constraints yield significant improvements, without the need for additional manual labels. This is possible because the network can learn to propagate labels from locations where it is confident to more difficult locations.

In summary, our S4-Net approach yields quality labels across given target sequences which makes it very interest-

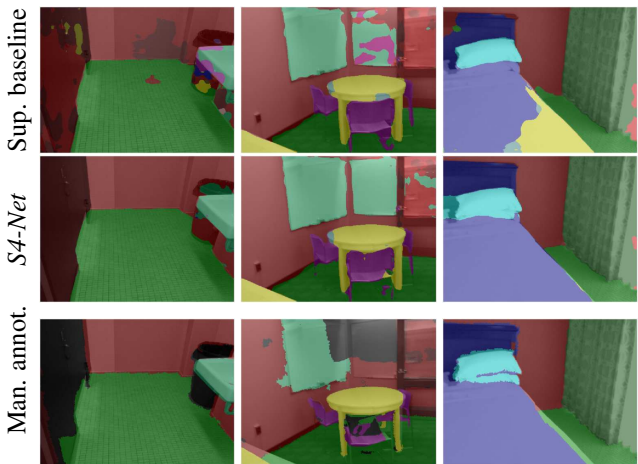


Figure 6: Qualitative results of S4-Net with depth predictions on unseen images from “Scan 2” of the target scenes from ScanNet. Similarly to our previous observations, S4-Net predicts quality segmentations in many regions which are noisy, wrongly labeled or unlabeled in the manual annotations.

ing for the task of sequence labelling. The segmentation network trained with S4-Net also generalizes nicely to unseen images of the target scene. This makes our approach useful for applications relying on semantic segmentation, for example in robotics and augmented reality. Finally, we have shown that the attractive idea of enforcing geometric constraints by means of depth predictions produces satisfying segmentations and achieves accuracy that is comparable to the accuracy when using ground truth depth information.

Acknowledgment

This work was supported by the Christian Doppler Laboratory for Semantic 3D Computer Vision, funded in part by Qualcomm Inc.

References

- [1] A. Armagan, M. Hirzer, P. M. Roth, and V. Lepetit. Learning to Align Semantic Segmentation and 2.5D Maps for Geolocalization. In *CVPR*, 2017. 2
- [2] B. Athiwaratkun, M. Finzi, P. Izmailov, and A. G. Wilson. There Are Many Consistent Explanations of Unlabeled Data: Why You Should Average. *ICLR*, 2019. 1, 2
- [3] V. Badrinarayanan, A. Kendall, and R. Cipolla. SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation. *TPAMI*, 2017. 2
- [4] W.-L. Chang, H.-P. Wang, W.-H. Peng, and W.-C. Chiu. All about Structure: Adapting Structural Information across Domains for Boosting Semantic Segmentation. In *CVPR*, 2019. 3
- [5] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation. In *ECCV*, 2018. 1, 2, 5
- [6] S. Choi, Q.-Y. Zhou, and V. Koltun. Robust Reconstruction of Indoor Scenes. In *CVPR*, 2015. 7
- [7] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Niessner. ScanNet: Richly-Annotated 3D Reconstructions of Indoor Scenes. In *CVPR*, 2017. 2, 5
- [8] Z. Dai, Z. Yang, F. Yang, W. W. Cohen, and R. Salakhutdinov. Good Semi-supervised Learning That Requires a Bad GAN. In *NIPS*, 2017. 2
- [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 5
- [10] C. N. dos Santos, K. Wadhawan, and B. Zhou. Learning Loss Functions for Semi-Supervised Learning via Discriminative Adversarial Networks. *arXiv preprint arXiv:1707.02198*, 2017. 2
- [11] D. Eigen and R. Fergus. Predicting Depth, Surface Normals and Semantic Labels With a Common Multi-Scale Convolutional Architecture. In *ICCV*, 2015. 4
- [12] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised Monocular Depth Estimation with Left-Right Consistency. In *CVPR*, 2017. 3, 4
- [13] C. Godard, O. Mac Aodha, and G. J. Brostow. Digging Into Self-Supervised Monocular Depth Estimation. *arXiv preprint arXiv:1806.01260*, 2018. 3, 4
- [14] M. Grinvald, F. Furrer, T. Novkovic, J. J. Chung, C. Cadena, R. Siegwart, and J. Nieto. Volumetric Instance-Aware Semantic Mapping and 3D Object Discovery. *arXiv preprint arXiv:1903.00268*, 2019. 1, 2, 3
- [15] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick. Mask R-CNN. In *ICCV*, 2017. 1
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016. 2, 5
- [17] W.-C. Hung, Y.-H. Tsai, Y.-T. Liou, Y.-Y. Lin, and M.-H. Yang. Adversarial Learning for Semi-Supervised Semantic Segmentation. In *BMVC*, 2018. 2
- [18] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial Transformer Networks. In *NIPS*, 2015. 4
- [19] A. Janoch, S. Karayev, Y. Jia, J. T. Barron, M. Fritz, K. Saenko, and T. Darrell. A Category-Level 3D Object Dataset: Putting the Kinect to Work. In *Consumer Depth Cameras for Computer Vision*. 2013. 5
- [20] T. Kalluri, G. Varma, M. Chandraker, and C. Jawahar. Universal semi-supervised semantic segmentation. *ICCV*, 2019. 3
- [21] A. Kumar, P. Sattigeri, and P. T. Fletcher. Improved Semi-Supervised Learning with GANs using Manifold Invariances. In *NIPS*, 2017. 2
- [22] S. Laine and T. Aila. Temporal Ensembling for Semi-Supervised Learning. In *ICLR*, 2017. 1, 2, 4
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully Convolutional Networks for Semantic Segmentation. In *CVPR*, 2015. 1, 2
- [24] Y. Luo, P. Liu, T. Guan, J. Yu, and Y. Yang. Significance-aware Information Bottleneck for Domain Adaptive Semantic Segmentation. *arXiv preprint arXiv:1904.00876*, 2019. 3
- [25] L. Ma, J. Stückler, C. Kerl, and D. Cremers. Multi-View Deep Learning for Consistent Semantic Mapping with RGB-D Cameras. In *IROS*, 2017. 1, 2, 3
- [26] R. Mahjourian, M. Wicke, and A. Angelova. Unsupervised Learning of Depth and Ego-Motion from Monocular Video Using 3D Geometric Constraints. In *CVPR*, 2018. 3, 4
- [27] J. McCormac, A. Handa, A. Davison, and S. Leutenegger. SemanticFusion: Dense 3D Semantic Mapping with Convolutional Neural Networks. In *ICRA*, 2017. 1, 2, 3
- [28] U. Michieli, M. Bassetton, G. Agresti, and P. Zanuttigh. Adversarial Learning and Self-Teaching Techniques for Domain Adaptation in Semantic Segmentation. *arXiv preprint arXiv:1909.00781*, 2019. 3
- [29] A. Milioto, P. Lottes, and C. Stachniss. Real-Time Semantic Segmentation of Crop and Weed for Precision Agriculture Robots Leveraging Background Knowledge in CNNs. In *ICRA*, 2018. 2
- [30] G. Papandreou, L.-C. Chen, K. P. Murphy, and A. L. Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, 2015. 2
- [31] Q.-H. Pham, B.-S. Hua, T. Nguyen, and S.-K. Yeung. Real-time Progressive 3D Semantic Segmentation for Indoor Scenes. In *WACV*, 2019. 1, 2, 3
- [32] Qian-Yi Zhou and Jaesik Park and Vladlen Koltun. Open3D: A Modern Library for 3D Data Processing. *arXiv:1801.09847*, 2018. 7
- [33] O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *MICCAI*, 2015. 2
- [34] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman. LabelMe: A Database and Web-Based Tool for Image Annotation. *IJCV*, 2008. 7
- [35] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor Segmentation and Support Inference from RGBD Images. In *ECCV*, 2012. 5
- [36] K. Simonyan and A. Zisserman. Very Deep Convolutional Networks for Large-Scale Image Recognition. *ICLR*, 2015. 2

- [37] S. Song, S. P. Lichtenberg, and J. Xiao. SUN RGB-D: A RGB-D Scene Understanding Benchmark Suite. In *CVPR*, 2015. 1, 5
- [38] R. Sun, X. Zhu, C. Wu, C. Huang, J. Shi, and L. Ma. Not All Areas Are Equal: Transfer Learning for Semantic Segmentation via Hierarchical Region Selection. In *CVPR*, 2019. 3
- [39] A. Tarvainen and H. Valpola. Mean Teachers are Better Role Models: Weight-Averaged Consistency targets improve Semi-Supervised Deep Learning Results. In *NIPS*, 2017. 1, 2
- [40] C. Wang, J. M. Buenaposada, R. Zhu, and S. Lucey. Learning Depth from Monocular Videos using Direct Methods. In *CVPR*, 2018. 3, 4
- [41] Y. Wei, H. Xiao, H. Shi, Z. Jie, J. Feng, and T. S. Huang. Revisiting dilated convolution: A simple approach for weakly- and semi-supervised semantic segmentation. In *CVPR*, 2018. 2
- [42] T. Whelan, R. F. Salas-Moreno, B. Glocker, A. J. Davison, and S. Leutenegger. ElasticFusion: Real-Time Dense SLAM and Light Source Estimation. *IJRR*, 2016. 3
- [43] J. Xiao, A. Owens, and A. Torralba. SUN3D: A Database of Big Spaces Reconstructed Using SfM and Object Labels. In *ICCV*, 2013. 5
- [44] Z. Yin and J. Shi. GeoNet: Unsupervised Learning of Dense Depth, Optical Flow and Camera Pose. In *CVPR*, 2018. 3, 4
- [45] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *CVPR*, 2017. 1, 2, 5
- [46] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe. Unsupervised Learning of Depth and Ego-Motion from Video. In *CVPR*, 2017. 3, 4