# Bidirectional Learning for Domain Adaptation of Semantic Segmentation

Yunsheng Li *
UC San Diego
yul554@eng.ucsd.edu

Lu Yuan
Microsoft
luyuan@microsoft.com

Nuno Vasconcelos
UC San Diego
nvasconcelos@ucsd.edu

## Abstract

*Domain adaptation for semantic image segmentation is very necessary since manually labeling large datasets with pixel-level labels is expensive and time consuming. Existing domain adaptation techniques either work on limited datasets, or yield not so good performance compared with supervised learning. In this paper, we propose a novel bidirectional learning framework for domain adaptation of segmentation. The learning process is organized as a closed loop, the image translation model and the segmentation adaptation model can be learned alternatively and promote to each other. Furthermore, we propose a self-supervised learning algorithm to learn a better segmentation adaptation model and in return improve the image translation model. Experiments show that our method is superior to the state-of-the-art methods in domain adaptation of segmentation with a big margin. The source code is available at https://github.com/liyunsheng13/BDL.*

## 1. Introduction

*(The body text of this page is rendered at a resolution too low to transcribe reliably.)*