

Chapter 1

Image Segmentation Using Deep Learning

Shervin Minaee

Snap, Inc., Seattle, WA, USA

Yuri Boykov

University of Waterloo, Waterloo, ON, Canada

Fatih Porikli

Qualcomm, Inc., San Diego, CA, USA

Antonio Plaza

University of Extremadura, Cáceres, Spain

Nasser Kehtarnavaz

University of Texas at Dallas, Richardson, TX, USA

Demetri Terzopoulos

University of California, Los Angeles, CA, USA

Image segmentation is a key task in computer vision and image processing with important applications such as scene understanding, medical image analysis, robotic perception, video surveillance, augmented reality, and image compression, among others, and numerous segmentation algorithms are found in the literature. Against this backdrop, the broad success of Deep Learning (DL) has prompted the development of new image segmentation approaches leveraging DL models. We provide a comprehensive review of this recent literature, covering the spectrum of pioneering efforts in semantic and instance segmentation, including convolutional pixel-labeling networks, encoder-decoder architectures, multiscale and pyramid-based approaches, recurrent networks, visual attention models, and generative models in adversarial settings. We investigate the relationships, strengths, and challenges of these DL-based segmentation models and discuss promising research directions.

1. Introduction

Image segmentation has been a fundamental problem in computer vision since the early days of the field [1, Chapter 8]. An essential component of many visual understanding systems, it involves partitioning images (or video frames) into multiple segments and objects [2, Chapter 5] and plays a central role in a broad range of applications [3, Part VI], including medical image analysis (e.g., tumor boundary extraction and tissue volume measurement), autonomous vehicles (e.g., navigable surface and pedestrian detection), video surveillance, augmented reality, etc.

Image segmentation can be formulated as the problem of classifying pixels with semantic labels (semantic segmentation), or partitioning of individual objects (instance segmentation), or both (panoptic segmentation). Semantic segmentation performs pixel-level labeling with a set of object categories (e.g., human, car, tree, sky) for all image pixels; thus, it is generally a more demanding undertaking than whole-image classification, which predicts a single label for the entire image. Instance segmentation extends the scope of semantic segmentation by detecting and delineating each object of interest in the image (e.g., individual people).

Numerous image segmentation algorithms have been developed in the literature, from the earliest methods, such as thresholding,⁴ histogram-based bundling, region-growing,⁵ k-means clustering,⁶ watershed methods,⁷ to more advanced algorithms such as active contours,⁸ graph cuts,⁹ conditional and Markov random fields,¹⁰ and sparsity-based^{11,12} methods. In recent years, however, deep learning (DL) models have yielded a new generation of image segmentation models with remarkable performance improvements, often achieving the highest accuracy rates on popular benchmarks. This has caused a paradigm shift in the field.

1.1. Overview

This chapter is a shortened and slightly revised version of our 2022 survey article.¹³ It covers the recent literature in deep-learning-based image segmentation, including more than 100 such segmentation methods proposed to date. It provides a comprehensive review with insights into key aspects of these methods. The target literature is organized into the following categories:

- (1) Fully convolutional networks
- (2) Convolutional models with graphical models
- (3) Encoder-decoder based models
- (4) Multiscale and pyramid network based models
- (5) R-CNN based models (for instance segmentation)
- (6) Dilated convolutional models and DeepLab family
- (7) Recurrent neural network based models
- (8) Attention-based models
- (9) Generative models and adversarial training
- (10) Convolutional models with active contour models

Through the use of skip connections in which feature maps from the final layers of the model are up-sampled and fused with feature maps of earlier layers, the model combines semantic information (from deep, coarse layers) and appearance information (from shallow, fine layers) in order to produce accurate and detailed segmentations. Tested on PASCAL VOC, NYUDv2, and SIFT Flow, the model achieved state-of-the-art segmentation performance.

FCNs have been applied to a variety of segmentation problems, such as brain tumor segmentation,¹⁵ instance-aware semantic segmentation,¹⁶ skin lesion segmentation,¹⁷ and iris segmentation.¹⁸ While demonstrating that DNNs can be trained to perform semantic segmentation in an end-to-end manner on variable-sized images, the conventional FCN model has some limitations—it is too computationally expensive for real-time inference, it does not account for global context information in an efficient manner, and it is not easily generalizable to 3D images. Several researchers have attempted to overcome some of the limitations of the FCN. For example, Liu *et al.*¹⁹ proposed ParseNet, which adds global context to FCNs by using the average feature for a layer to augment the features at each location. The feature map for a layer is pooled over the whole image, resulting in a context vector. The context vector is normalized and unpooled to produce new feature maps of the same size as the initial ones, which are then concatenated, which amounts to an FCN whose convolutional layers are replaced by the described module.

2.2. CNNs With Graphical Models

As discussed, the FCN ignores potentially useful scene-level semantic context. To exploit more context, several approaches incorporate into DL architectures probabilistic graphical models, such as Conditional Random Fields (CRFs) and Markov Random Fields (MRFs).

Due to the invariance properties that make CNNs good for high level tasks such as classification, responses from the later layers of deep CNNs are not sufficiently well localized for accurate object segmentation. To address this drawback, Chen *et al.*²⁰ proposed a semantic segmentation algorithm that combines CNNs and fully-connected CRFs. They showed that their model can localize segment boundaries with higher accuracy than was possible with previous methods.

Schwing and Urtasun²¹ proposed a fully-connected deep structured network for image segmentation. They jointly trained CNNs and fully-connected CRFs for semantic image segmentation, and achieved encouraging results on the challenging PASCAL VOC 2012 dataset. Zheng *et al.*²² proposed a similar semantic segmentation approach. In related work, Lin *et al.*²³ proposed an efficient semantic segmentation model based on contextual deep CRFs. They explored “patch-patch” context (between image regions) and “patch-background” context to improve semantic segmentation through the use of contextual information.

Liu *et al.*²⁴ proposed a semantic segmentation algorithm that incorporates rich information into MRFs, including high-order relations and mixture of label

contexts. Unlike previous efforts that optimized MRFs using iterative algorithms, they proposed a CNN model, namely a Parsing Network, which enables deterministic end-to-end computation in one pass.

2.3. Encoder-Decoder Based Models

Most of the popular DL-based segmentation models use some kind of encoder-decoder architecture. We group these models into two categories: those for general image segmentation, and those for medical image segmentation.

2.3.1. General Image Segmentation

Noh *et al.*²⁵ introduced semantic segmentation based on deconvolution (a.k.a. transposed convolution). Their model, DeConvNet, consists of two parts, an encoder using convolutional layers adopted from the VGG 16-layer network and a multilayer deconvolutional network that inputs the feature vector and generates a map of pixel-accurate class probabilities. The latter comprises deconvolution and unpooling layers, which identify pixel-wise class labels and predict segmentation masks.

Badrinarayanan *et al.*²⁶ proposed SegNet, a fully convolutional encoder-decoder architecture for image segmentation. Similar to the deconvolution network, the core trainable segmentation engine of SegNet consists of an encoder network, which is topologically identical to the 13 convolutional layers of the VGG16 network, and a corresponding decoder network followed by a pixel-wise classification layer. The main novelty of SegNet is in the way the decoder upsamples its lower-resolution input feature map(s); specifically, using pooling indices computed in the max-pooling step of the corresponding encoder to perform nonlinear up-sampling.

A limitation of encoder-decoder based models is the loss of fine-grained image information, due to the loss of resolution through the encoding process. HRNet²⁷ addresses this shortcoming. Other than recovering high-resolution representations as is done in DeConvNet, SegNet, and other models, HRNet maintains high-resolution representations through the encoding process by connecting the high-to-low resolution convolution streams in parallel and repeatedly exchanging the information across resolutions. There are four stages: the 1st stage consists of high-resolution convolutions, while the 2nd/3rd/4th stage repeats 2-resolution/3-resolution/4-resolution blocks. Several recent semantic segmentation models use HRNet as a backbone.

Several other works adopt transposed convolutions, or encoder-decoders for image segmentation, such as Stacked Deconvolutional Network (SDN),²⁸ Linknet,²⁹ W-Net,³⁰ and locality-sensitive deconvolution networks for RGB-D segmentation.³¹

2.3.2. Medical and Biomedical Image Segmentation

Several models inspired by FCNs and encoder-decoder networks were initially developed for medical/biomedical image segmentation, but are now also being used outside the medical domain.

Ronneberger *et al.*³² proposed the U-Net for efficiently segmenting biological microscopy images. The U-Net architecture comprises two parts, a contracting path to capture context, and a symmetric expanding path that enables precise localization. The U-Net training strategy relies on the use of data augmentation to learn effectively from very few annotated images. It was trained on 30 transmitted light microscopy images, and it won the ISBI cell tracking challenge 2015 by a large margin.

Various extensions of U-Net have been developed for different kinds of images and problem domains; for example, Zhou *et al.*³³ developed a nested U-Net architecture, Zhang *et al.*³⁴ developed a road segmentation algorithm based on U-Net, and Cicek *et al.*³⁵ proposed a U-Net architecture for 3D images.

V-Net, proposed by Milletari *et al.*³⁶ for 3D medical image segmentation, is another well known FCN-based model. The authors introduced a new loss function based on the Dice coefficient, enabling the model to deal with situations in which there is a strong imbalance between the number of voxels in the foreground and background. The network was trained end-to-end on MRI images of the prostate and learns to predict segmentation for the whole volume at once. Some of the other relevant works on medical image segmentation includes Progressive Dense V-Net³⁷ for automatic segmentation of pulmonary lobes from chest CT images, and the 3D-CNN encoder for lesion segmentation.³⁸

2.4. Multiscale and Pyramid Network Based Models

Multiscale analysis, a well established idea in image processing, has been deployed in various neural network architectures. One of the most prominent models of this sort is the Feature Pyramid Network (FPN) proposed by Lin *et al.*,³⁹ which was developed for object detection but was also applied to segmentation. The inherent multiscale, pyramidal hierarchy of deep CNNs was used to construct feature pyramids with marginal extra cost. To merge low and high resolution features, the FPN is composed of a bottom-up pathway, a top-down pathway and lateral connections. The concatenated feature maps are then processed by a 3×3 convolution to produce the output of each stage. Finally, each stage of the top-down pathway generates a prediction to detect an object. For image segmentation, the authors use two multilayer perceptrons (MLPs) to generate the masks.

Zhao *et al.*⁴⁰ developed the Pyramid Scene Parsing Network (PSPN), a multiscale network to better learn the global context representation of a scene. Multiple patterns are extracted from the input image using a residual network (ResNet) as a feature extractor, with a dilated network. These feature maps are then fed into a pyramid pooling module to distinguish patterns of different scales. They are pooled at four different scales, each one corresponding to a pyramid level, and processed by a 1×1 convolutional layer to reduce their dimensions. The outputs of the pyramid levels are up-sampled and concatenated with the initial feature maps to capture both local and global context information. Finally, a convolutional layer is used to generate

the pixel-wise predictions.

Ghiasi and Fowlkes⁴¹ developed a multiresolution reconstruction architecture based on a Laplacian pyramid that uses skip connections from higher resolution feature maps and multiplicative gating to successively refine segment boundaries reconstructed from lower-resolution maps. They showed that while the apparent spatial resolution of convolutional feature maps is low, the high-dimensional feature representation contains significant sub-pixel localization information.

Other models use multiscale analysis for segmentation, among them Dynamic Multiscale Filters Network (DM-Net),⁴² Context Contrast Network and gated multiscale aggregation (CCN),⁴³ Adaptive Pyramid Context Network (APC-Net),⁴⁴ MultiScale Context Intertwining (MSCI),⁴⁵ and salient object segmentation.⁴⁶

2.5. R-CNN Based Models

The Regional CNN (R-CNN) and its extensions have proven successful in object detection applications. In particular, the Faster R-CNN⁴⁷ architecture uses a region proposal network (RPN) that proposes bounding box candidates. The RPN extracts a Region of Interest (RoI), and an RoIPool layer computes features from these proposals to infer the bounding box coordinates and class of the object. Some extensions of R-CNN have been used to address the instance segmentation problem; i.e., the task of simultaneously performing object detection and semantic segmentation.

He *et al.*⁴⁸ proposed Mask R-CNN, which outperformed previous benchmarks on many COCO object instance segmentation challenges, efficiently detecting objects in an image while simultaneously generating a high-quality segmentation mask for each instance. Essentially, it is a Faster R-CNN with 3 output branches—the first computes the bounding box coordinates, the second computes the associated classes, and the third computes the binary mask to segment the object. The Mask R-CNN loss function combines the losses of the bounding box coordinates, the predicted class, and the segmentation mask, and trains all of them jointly.

The Path Aggregation Network (PANet) proposed by Liu *et al.*⁴⁹ is based on the Mask R-CNN and FPN models. The feature extractor of the network uses an FPN backbone with a new augmented bottom-up pathway improving the propagation of lower-layer features. Each stage of this third pathway takes as input the feature maps of the previous stage and processes them with a 3×3 convolutional layer. A lateral connection adds the output to the same-stage feature maps of the top-down pathway and these feed the next stage.

Dai *et al.*⁵⁰ developed a multitask network for instance-aware semantic segmentation that consists of three networks for differentiating instances, estimating masks, and categorizing objects. These networks form a cascaded structure and are designed to share their convolutional features. Hu *et al.*⁵¹ proposed a new partially-supervised training paradigm together with a novel weight transfer function, which enables training instance segmentation models on a large set of categories, all of which have

box annotations, but only a small fraction of which have mask annotations.

Chen *et al.*⁵² developed an instance segmentation model, MaskLab, by refining object detection with semantic and direction features based on Faster R-CNN. This model produces three outputs, box detection, semantic segmentation logits for pixel-wise classification, and direction prediction logits for predicting each pixel's direction toward its instance center. Building on the Faster R-CNN object detector, the predicted boxes provide accurate localization of object instances. Within each region of interest, MaskLab performs foreground/background segmentation by combining semantic and direction prediction.

Tensormask, proposed by Chen *et al.*,⁵³ is based on dense sliding window instance segmentation. The authors treat dense instance segmentation as a prediction task over 4D tensors and present a general framework that enables novel operators on 4D tensors. They demonstrate that the tensor approach yields large gains over baselines, with results comparable to Mask R-CNN.

Other instance segmentation models have been developed based on R-CNN, such as those developed for mask proposals, including R-FCN,⁵⁴ DeepMask,⁵⁵ PolarMask,⁵⁶ boundary-aware instance segmentation,⁵⁷ and CenterMask.⁵⁸ Another promising approach is to tackle the instance segmentation problem by learning grouping cues for bottom-up segmentation, such as deep watershed transform,⁵⁹ real-time instance segmentation,⁶⁰ and semantic instance segmentation via deep metric learning.⁶¹

2.6. Dilated Convolutional Models

Dilated (a.k.a. “atrous”) convolution introduces to convolutional layers another parameter, the dilation rate. For example, a 3×3 kernel with a dilation rate of 2 will have the same size receptive field as a 5×5 kernel while using only 9 parameters, thus enlarging the receptive field with no increase in computational cost.

Dilated convolutions have been popular in the field of real-time segmentation, and many recent publications report the use of this technique. Some of the most important include the DeepLab family,⁶² multiscale context aggregation,⁶³ Dense Upsampling Convolution and Hybrid Dilated Convolution (DUC-HDC),⁶⁴ densely connected Atrous Spatial Pyramid Pooling (DenseASPP),⁶⁵ and the Efficient Network (ENet).⁶⁶

DeepLabv1²⁰ and DeepLabv2,⁶² developed by Chen *et al.*, are among the most popular image segmentation models. The latter has three key features. First is the use of dilated convolution to address the decreasing resolution in the network caused by max-pooling and striding. Second is Atrous Spatial Pyramid Pooling (ASPP), which probes an incoming convolutional feature layer with filters at multiple sampling rates, thus capturing objects as well as multiscale image context to robustly segment objects at multiple scales. Third is improved localization of object boundaries by combining methods from deep CNNs, such as fully convolutional VGG-16 or ResNet 101, and probabilistic graphical models, specifically fully-connected CRFs.

Subsequently, Chen *et al.*⁶⁷ proposed DeepLabv3, which combines cascaded and parallel modules of dilated convolutions. The parallel convolution modules are grouped in the ASPP. A 1×1 convolution and batch normalization are added in the ASPP. All the outputs are concatenated and processed by another 1×1 convolution to create the final output with logits for each pixel. Next, Chen *et al.*⁶⁸ released DeepLabv3+, which uses an encoder-decoder architecture including dilated separable convolution composed of a depthwise convolution (spatial convolution for each channel of the input) and pointwise convolution (1×1 convolution with the depthwise convolution as input). They used the DeepLabv3 framework as the encoder. The most relevant model has a modified Xception backbone with more layers, dilated depthwise separable convolutions instead of max pooling and batch normalization.

2.7. RNN Based Models

While CNNs are a natural fit for computer vision problems, they are not the only possibility. RNNs are useful in modeling the short/long term dependencies among pixels to (potentially) improve the estimation of the segmentation map. Using RNNs, pixels may be linked together and processed sequentially to model global contexts and improve semantic segmentation. However the natural 2D structure of images poses a challenge.

Visin *et al.*⁶⁹ proposed an RNN-based model for semantic segmentation called ReSeg. This model is mainly based on ReNet,⁷⁰ which was developed for image classification. Each ReNet layer is composed of four RNNs that sweep the image horizontally and vertically in both directions, encoding patches/activations, and providing relevant global information. To perform image segmentation with the ReSeg model, ReNet layers are stacked atop pre-trained VGG-16 convolutional layers, which extract generic local features, and are then followed by up-sampling layers to recover the original image resolution in the final predictions.

Byeon *et al.*⁷¹ performed per-pixel segmentation and classification of images of natural scenes using 2D LSTM networks, which learn textures and the complex spatial dependencies of labels in a single model that carries out classification, segmentation, and context integration.

Liang *et al.*⁷² proposed a semantic segmentation model based on a graph-LSTM network in which convolutional layers are augmented by graph-LSTM layers built on super-pixel maps, which provide a more global structural context. These layers generalize the LSTM for uniform, array-structured data (i.e., row, grid, or diagonal LSTMs) to nonuniform, graph-structured data, where arbitrary-shaped superpixels are semantically consistent nodes and the adjacency relations between superpixels correspond to edges, thus forming an undirected graph.

Xiang and Fox⁷³ proposed Data Associated Recurrent Neural Networks (DA-RNNs) for joint 3D scene mapping and semantic labeling. DA-RNNs use a new recurrent neural network architecture for semantic labeling on RGB-D videos. The

output of the network is integrated with mapping techniques such as Kinect-Fusion in order to inject semantic information into the reconstructed 3D scene.

Hu *et al.*⁷⁴ developed a semantic segmentation algorithm that combines a CNN to encode the image and an LSTM to encode its linguistic description. To produce pixel-wise image segmentations from language inputs, they propose an end-to-end trainable recurrent and convolutional model that jointly learns to process visual and linguistic information. This differs from traditional semantic segmentation over a predefined set of semantic classes; i.e., the phrase “two men sitting on the right bench” requires segmenting only the two people on the right bench and no others sitting on another bench or standing.

A drawback of RNN-based models is that they will generally be slower than their CNN counterparts as their sequential nature is not amenable to parallelization.

2.8. Attention-Based Models

Attention mechanisms have been persistently explored in computer vision over the years, and it is not surprising to find publications that apply them to semantic segmentation.

Chen *et al.*⁷⁵ proposed an attention mechanism that learns to softly weight multiscale features at each pixel location. They adapt a powerful semantic segmentation model and jointly train it with multiscale images and the attention model. The attention mechanism enables the model to assess the importance of features at different positions and scales, and it outperforms average and max pooling.

Unlike approaches in which convolutional classifiers are trained to learn the representative semantic features of labeled objects, Huang *et al.*⁷⁶ proposed a Reverse Attention Network (RAN) architecture for semantic segmentation that also applies reverse attention mechanisms, thereby training the model to capture the opposite concept—features that are not associated with a target class. The RAN network performs the direct and reverse-attention learning processes simultaneously.

Li *et al.*⁷⁷ developed a Pyramid Attention Network for semantic segmentation, which exploits global contextual information for semantic segmentation. Eschewing complicated dilated convolutions and decoder networks, they combined attention mechanisms and spatial pyramids to extract precise dense features for pixel labeling. Fu *et al.*⁷⁸ proposed a dual attention network for scene segmentation that can capture rich contextual dependencies based on the self-attention mechanism. Specifically, they append two types of attention modules on top of a dilated FCN that models the semantic inter-dependencies in spatial and channel dimensions, respectively. The position attention module selectively aggregates the features at each position via weighted sums.

Other applications of attention mechanisms to semantic segmentation include OCNet,⁷⁹ which employs an object context pooling inspired by self-attention mechanism, ResNeSt: Split-Attention Networks,⁸⁰ Height-driven Attention Networks,⁸¹ Expectation-Maximization Attention (EMANet),⁸² Criss-Cross Attention Network

(CCNet),⁸³ end-to-end instance segmentation with recurrent attention,⁸⁴ a point-wise spatial attention network for scene parsing,⁸⁵ and Discriminative Feature Network (DFN).⁸⁶

2.9. Generative Models and Adversarial Training

GANs have been applied to a wide range of tasks in computer vision, not excluding image segmentation.

Luc *et al.*⁸⁷ proposed an adversarial training approach for semantic segmentation in which they trained a convolutional semantic segmentation network, along with an adversarial network that discriminates between ground-truth segmentation maps and those generated by the segmentation network. They showed that the adversarial training approach yields improved accuracy on the Stanford Background and PASCAL VOC 2012 datasets.

Souly *et al.*⁸⁸ proposed semi-weakly supervised semantic segmentation using GANs. Their model consists of a generator network providing extra training examples to a multiclass classifier, acting as discriminator in the GAN framework, that assigns sample a label from the possible label classes or marks it as a fake sample (extra class).

Hung *et al.*⁸⁹ developed a framework for semi-supervised semantic segmentation using an adversarial network. They designed an FCN discriminator to differentiate the predicted probability maps from the ground truth segmentation distribution, considering the spatial resolution. The loss function of this model has three terms: cross-entropy loss on the segmentation ground truth, adversarial loss of the discriminator network, and semi-supervised loss based on the confidence map output of the discriminator.

Xue *et al.*⁹⁰ proposed an adversarial network with multiscale L1 Loss for medical image segmentation. They used an FCN as the segmentor to generate segmentation label maps, and proposed a novel adversarial critic network with a multi-scale L1 loss function to force the critic and segmentor to learn both global and local features that capture long and short range spatial relationships between pixels.

Other approaches based on adversarial training include cell image segmentation using GANs,⁹¹ and segmentation and generation of the invisible parts of objects.⁹²

2.10. CNN Models With Active Contour Models

The exploration of synergies between FCNs and Active Contour Models (ACMs)⁸ has recently attracted research interest.

One approach is to formulate new loss functions that are inspired by ACM principles. For example, inspired by the global energy formulation of Chan and Vese,⁹³ Chen *et al.*⁹⁴ proposed a supervised loss layer that incorporated area and size information of the predicted masks during training of an FCN and tackled the problem of ventricle segmentation in cardiac MRI. Similarly, Gur *et al.*⁹⁵ presented

an unsupervised loss function based on morphological active contours without edges⁹⁶ for microvascular image segmentation.

A different approach initially sought to utilize the ACM merely as a post-processor of the output of an FCN and several efforts attempted modest co-learning by pre-training the FCN. One example of an ACM post-processor for the task of semantic segmentation of natural images is the work by Le *et al.*⁹⁷ in which level-set ACMs are implemented as RNNs. Deep Active Contours by Rupprecht *et al.*,⁹⁸ is another example. For medical image segmentation, Hatamizadeh *et al.*⁹⁹ proposed an integrated Deep Active Lesion Segmentation (DALs) model that trains the FCN backbone to predict the parameter functions of a novel, locally-parameterized level-set energy functional. In another relevant effort, Marcos *et al.*¹⁰⁰ proposed Deep Structured Active Contours (DSAC), which combines ACMs and pre-trained FCNs in a structured prediction framework for building instance segmentation (albeit with manual initialization) in aerial images. For the same application, Cheng *et al.*¹⁰¹ proposed the Deep Active Ray Network (DarNet), which is similar to DSAC, but with a different explicit ACM formulation based on polar coordinates to prevent contour self-intersection.

A truly end-to-end backpropagation trainable, fully-integrated FCN-ACM combination was recently introduced by Hatamizadeh *et al.*,¹⁰² dubbed Trainable Deep Active Contours (TDAC). Going beyond their earlier work,⁹⁹ they implemented the locally-parameterized level-set ACM in the form of additional convolutional layers following the layers of the backbone FCN, exploiting Tensorflow's automatic differentiation mechanism to backpropagate training error gradients throughout the entire DCAC framework. The fully-automated model requires no intervention either during training or segmentation, can naturally segment multiple instances of objects of interest, and deal with arbitrary object shape including sharp corners.

2.11. Other Models

Other popular DL architectures for image segmentation include the following:

Context Encoding Network (EncNet)¹⁰³ uses a basic feature extractor and feeds the feature maps into a context encoding module. RefineNet¹⁰⁴ is a multipath refinement network that explicitly exploits all the information available along the down-sampling process to enable high-resolution prediction using long-range residual connections. Seednet¹⁰⁵ introduced an automatic seed generation technique with deep reinforcement learning that learns to solve the interactive segmentation problem. Object-Contextual Representations (OCR)²⁷ learns object regions and the relation between each pixel and each object region, augmenting the representation pixels with the object-contextual representation. Additional models and methods include BoxSup,¹⁰⁶ Graph Convolutional Networks (GCN),¹⁰⁷ Wide ResNet,¹⁰⁸ Exfuse¹⁰⁹ (enhancing low-level and high-level features fusion), Feedforward-Net,¹¹⁰ saliency-aware models for geodesic video segmentation,¹¹¹ Dual Image Segmentation (DIS),¹¹² FoveaNet¹¹³ (perspective-aware scene parsing), Ladder DenseNet,¹¹⁴

Bilateral Segmentation Network (BiSeNet),¹¹⁵ Semantic Prediction Guidance for Scene Parsing (SPGNet),¹¹⁶ gated shape CNNs,¹¹⁷ Adaptive Context Network (AC-Net),¹¹⁸ Dynamic-Structured Semantic Propagation Network (DSSPN),¹¹⁹ Symbolic Graph Reasoning (SGR),¹²⁰ CascadeNet,¹²¹ Scale-Adaptive Convolutions (SAC),¹²² Unified Perceptual parsing Network (UperNet),¹²³ segmentation by re-training and self-training,¹²⁴ densely connected neural architecture search,¹²⁵ hierarchical multi-scale attention,¹²⁶ Efficient RGB-D Semantic Segmentation (ESA-Net),¹²⁷ Iterative Pyramid Contexts,¹²⁸ and Learning Dynamic Routing for Semantic Segmentation.¹²⁹

Panoptic segmentation¹³⁰ is growing in popularity. Efforts in this direction include Panoptic Feature Pyramid Network (PFPN),¹³¹ attention-guided network for panoptic segmentation,¹³² seamless scene segmentation,¹³³ panoptic Deeplab,¹³⁴ unified panoptic segmentation network,¹³⁵ and efficient panoptic segmentation.¹³⁶

3. Challenges and Opportunities

We have surveyed image segmentation algorithms based on deep learning models, which have achieved impressive performance in various image segmentation tasks and benchmarks, grouped into architectural categories such as: CNN and FCN, RNN, R-CNN, dilated CNN, attention-based models, generative and adversarial models, among others. Without a doubt, image segmentation has benefited greatly from deep learning, but several challenges lie ahead. We will next discuss some of the promising research directions that we believe will help in further advancing image segmentation algorithms.

3.1. *More Challenging Datasets*

Several large-scale image datasets have been created for semantic segmentation and instance segmentation. However, there remains a need for more challenging datasets, as well as datasets of different kinds of images. For still images, datasets with a large number of objects and overlapping objects would be very valuable. This can enable the training of models that handle dense object scenarios better, as well as large overlaps among objects as is common in real-world scenarios. With the rising popularity of 3D image segmentation, especially in medical image analysis, there is also a strong need for large-scale annotated 3D image datasets, which are more difficult to create than their lower dimensional counterparts.

3.2. *Combining DL and Earlier Segmentation Models*

There is now broad agreement that the performance of DL-based segmentation algorithms is plateauing, especially in certain application domains such as medical image analysis. To advance to the next level of performance, we must further explore the combination of CNN-based image segmentation models with prominent “classical” model-based image segmentation methods. The integration of CNNs with graphical

14 *S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos*

models has been studied, but their integration with active contours, graph cuts, and other segmentation models is fairly recent and deserves further work.

3.3. Interpretable Deep Models

While DL-based models have achieved promising performance on challenging benchmarks, there remain open questions about these models. For example, what exactly are deep models learning? How should we interpret the features learned by these models? What is a minimal neural architecture that can achieve a certain segmentation accuracy on a given dataset? Although some techniques are available to visualize the learned convolutional kernels of these models, a comprehensive study of the underlying behavior/dynamics of these models is lacking. A better understanding of the theoretical aspects of these models can enable the development of better models curated toward various segmentation scenarios.

3.4. Weakly-Supervised and Unsupervised Learning

Weakly-supervised (a.k.a. few shot) learning¹³⁷ and unsupervised learning¹³⁸ are becoming very active research areas. These techniques promise to be specially valuable for image segmentation, as collecting pixel-accurately labeled training images is problematic in many application domains, particularly so in medical image analysis. The transfer learning approach is to train a generic image segmentation model on a large set of labeled samples (perhaps from a public benchmark) and then fine-tune that model on a few samples from some specific target application. Self-supervised learning is another promising direction that is attracting much attraction in various fields. With the help of self-supervised learning, many details in images can be captured in order to train segmentation models with far fewer training samples. Models based on reinforcement learning could also be another potential future direction, as they have scarcely received attention for image segmentation. For example, MOREL¹³⁹ introduced a deep reinforcement learning approach for moving object segmentation in videos.

3.5. Real-time Models for Various Applications

In many applications, accuracy is the most important factor; however, there are applications in which it is also critical to have segmentation models that can run in near real-time, or at common camera frame rates (at least 25 frames per second). This is useful for computer vision systems that are, for example, deployed in autonomous vehicles. Most of the current models are far from this frame-rate; e.g., FCN-8 takes roughly 100 ms to process a low-resolution image. Models based on dilated convolution help to increase the speed of segmentation models to some extent, but there is still plenty of room for improvement.

3.6. Memory Efficient Models

Many modern segmentation models require a significant amount of memory even during the inference stage. So far, much effort has been directed towards improving the accuracy of such models, but in order to fit them into specific devices, such as mobile phones, the networks must be simplified. This can be done either by using simpler models, or by using model compression techniques, or even by training a complex model and using knowledge distillation techniques to compress it into a smaller, memory efficient network that mimics the complex model.

3.7. Applications

DL-based segmentation methods have been successfully applied to satellite images in remote sensing,¹⁴⁰ such as to support urban planning¹⁴¹ and precision agriculture.¹⁴² Images collected by airborne platforms¹⁴³ and drones¹⁴⁴ have also been segmented using DL-based segmentation methods in order to address important environmental problems including ones related to climate change. The main challenges of the remote sensing domain stem from the typically formidable size of the imagery (often collected by imaging spectrometers with hundreds or even thousands of spectral bands) and the limited ground-truth information necessary to evaluate the accuracy of the segmentation algorithms. Similarly, DL-based segmentation techniques in the evaluation of construction materials¹⁴⁵ face challenges related to the massive volume of the related image data and the limited reference information for validation purposes. Last but not least, an important application field for DL-based segmentation has been biomedical imaging.¹⁴⁶ Here, an opportunity is to design standardized image databases useful in evaluating new infectious diseases and tracking pandemics.¹⁴⁷

References

1. A. Rosenfeld and A. C. Kak, *Digital Picture Processing*. (Academic Press, 1976).
2. R. Szeliski, *Computer Vision: Algorithms and Applications*. (Springer, 2010).
3. D. Forsyth and J. Ponce, *Computer Vision: A Modern Approach*. (Prentice Hall, 2002).
4. N. Otsu, A threshold selection method from gray-level histograms, *IEEE Transactions on Systems, Man, and Cybernetics*. **9**(1), 62–66, (1979).
5. R. Nock and F. Nielsen, Statistical region merging, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **26**(11), 1452–1458, (2004).
6. N. Dhanachandra, K. Manglem, and Y. J. Chanu, Image segmentation using K-means clustering algorithm and subtractive clustering algorithm, *Procedia Computer Science*. **54**, 764–771, (2015).
7. L. Najman and M. Schmitt, Watershed of a continuous function, *Signal Processing*. **38**(1), 99–112, (1994).
8. M. Kass, A. Witkin, and D. Terzopoulos, Snakes: Active contour models, *International Journal of Computer Vision*. **1**(4), 321–331, (1988).
9. Y. Boykov, O. Veksler, and R. Zabih, Fast approximate energy minimization via

- graph cuts, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **23** (11), 1222–1239, (2001).
10. N. Plath, M. Toussaint, and S. Nakajima. Multi-class image segmentation using conditional random fields and global classification. In *International Conference on Machine Learning*, pp. 817–824. ACM, (2009).
 11. J.-L. Starck, M. Elad, and D. L. Donoho, Image decomposition via the combination of sparse representations and a variational approach, *IEEE Transactions on Image Processing*. **14**(10), 1570–1582, (2005).
 12. S. Minaee and Y. Wang, An ADMM approach to masked signal decomposition using subspace representation, *IEEE Transactions on Image Processing*. **28**(7), 3192–3204, (2019).
 13. S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos, Image segmentation using deep learning: A survey, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **44**(7), 3523–3542 (July, 2022).
 14. J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440, (2015).
 15. G. Wang, W. Li, S. Ourselin, and T. Vercauteren. Automatic brain tumor segmentation using cascaded anisotropic convolutional neural networks. In *International MICCAI Brainlesion Workshop*, pp. 178–190. Springer, (2017).
 16. Y. Li, H. Qi, J. Dai, X. Ji, and Y. Wei. Fully convolutional instance-aware semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2359–2367, (2017).
 17. Y. Yuan, M. Chao, and Y.-C. Lo, Automatic skin lesion segmentation using deep fully convolutional networks with Jaccard distance, *IEEE Transactions on Medical Imaging*. **36**(9), 1876–1886, (2017).
 18. N. Liu, H. Li, M. Zhang, J. Liu, Z. Sun, and T. Tan. Accurate iris segmentation in non-cooperative environments using fully convolutional networks. In *International Conference on Biometrics*, pp. 1–8. IEEE, (2016).
 19. W. Liu, A. Rabinovich, and A. C. Berg, ParseNet: Looking wider to see better, *arXiv preprint arXiv:1506.04579*. (2015).
 20. L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, Semantic image segmentation with deep convolutional nets and fully connected CRFs, *arXiv preprint arXiv:1412.7062*. (2014).
 21. A. G. Schwing and R. Urtasun, Fully connected deep structured networks, *arXiv preprint arXiv:1503.02351*. (2015).
 22. S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *IEEE International Conference on Computer Vision*, pp. 1529–1537, (2015).
 23. G. Lin, C. Shen, A. Van Den Hengel, and I. Reid. Efficient piecewise training of deep structured models for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3194–3203, (2016).
 24. Z. Liu, X. Li, P. Luo, C.-C. Loy, and X. Tang. Semantic image segmentation via deep parsing network. In *IEEE International Conference on Computer Vision*, pp. 1377–1385, (2015).
 25. H. Noh, S. Hong, and B. Han. Learning deconvolution network for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 1520–1528, (2015).
 26. V. Badrinarayanan, A. Kendall, and R. Cipolla, SegNet: A deep convolutional encoder-decoder architecture for image segmentation, *IEEE Transactions on Pattern*

- Analysis and Machine Intelligence*. **39**(12), 2481–2495, (2017).
27. Y. Yuan, X. Chen, and J. Wang, Object-contextual representations for semantic segmentation, *arXiv preprint arXiv:1909.11065*. (2019).
 28. J. Fu, J. Liu, Y. Wang, J. Zhou, C. Wang, and H. Lu, Stacked deconvolutional network for semantic segmentation, *IEEE Transactions on Image Processing*. (2019).
 29. A. Chaurasia and E. Culurciello. LinkNet: Exploiting encoder representations for efficient semantic segmentation. In *IEEE International Conference on Visual Communications and Image Processing*, pp. 1–4. IEEE, (2017).
 30. X. Xia and B. Kulis, W-Net: A deep model for fully unsupervised image segmentation, *arXiv preprint arXiv:1711.08506*. (2017).
 31. Y. Cheng, R. Cai, Z. Li, X. Zhao, and K. Huang. Locality-sensitive deconvolution networks with gated fusion for RGB-D indoor semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3029–3037, (2017).
 32. O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 234–241. Springer, (2015).
 33. Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang. UNet++: A nested U-Net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pp. 3–11. Springer, (2018).
 34. Z. Zhang, Q. Liu, and Y. Wang, Road extraction by deep residual U-Net, *IEEE Geoscience and Remote Sensing Letters*. **15**(5), 749–753, (2018).
 35. Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 424–432. Springer, (2016).
 36. F. Milletari, N. Navab, and S.-A. Ahmadi. V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In *International Conference on 3D Vision*, pp. 565–571. IEEE, (2016).
 37. A.-A.-Z. Imran, A. Hatamizadeh, S. P. Ananth, X. Ding, N. Tajbakhsh, and D. Terzopoulos, Fast and automatic segmentation of pulmonary lobes from chest CT using a progressive dense V-Network, *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*. **8**(5–6), 509–518, (2020).
 38. T. Brosch, L. Y. Tang, Y. Yoo, D. K. Li, A. Traboulosee, and R. Tam, Deep 3D convolutional encoder networks with shortcuts for multiscale feature integration applied to multiple sclerosis lesion segmentation, *IEEE Transactions on Medical Imaging*. **35**(5), 1229–1239, (2016).
 39. T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2117–2125, (2017).
 40. H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2881–2890, (2017).
 41. G. Ghiasi and C. C. Fowlkes. Laplacian pyramid reconstruction and refinement for semantic segmentation. In *European Conference on Computer Vision*, pp. 519–534. Springer, (2016).
 42. J. He, Z. Deng, and Y. Qiao. Dynamic multi-scale filters for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 3562–3572, (2019).
 43. H. Ding, X. Jiang, B. Shuai, A. Qun Liu, and G. Wang. Context contrasted feature and gated multi-scale aggregation for scene segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2393–2402, (2018).

- 18 S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos
44. J. He, Z. Deng, L. Zhou, Y. Wang, and Y. Qiao. Adaptive pyramid context network for semantic segmentation. In *Conference on Computer Vision and Pattern Recognition*, pp. 7519–7528, (2019).
 45. D. Lin, Y. Ji, D. Lischinski, D. Cohen-Or, and H. Huang. Multi-scale context intertwining for semantic segmentation. In *European Conference on Computer Vision*, pp. 603–619, (2018).
 46. G. Li, Y. Xie, L. Lin, and Y. Yu. Instance-level salient object segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2386–2395, (2017).
 47. S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems*, pp. 91–99, (2015).
 48. K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *IEEE International Conference on Computer Vision*, pp. 2961–2969, (2017).
 49. S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8759–8768, (2018).
 50. J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3150–3158, (2016).
 51. R. Hu, P. Dollár, K. He, T. Darrell, and R. Girshick. Learning to segment every thing. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4233–4241, (2018).
 52. L.-C. Chen, A. Hermans, G. Papandreou, F. Schroff, P. Wang, and H. Adam. Masklab: Instance segmentation by refining object detection with semantic and direction features. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4013–4022, (2018).
 53. X. Chen, R. Girshick, K. He, and P. Dollár. Tensormask: A foundation for dense object segmentation, *arXiv preprint arXiv:1903.12174*. (2019).
 54. J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object detection via region-based fully convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 379–387, (2016).
 55. P. O. Pinheiro, R. Collobert, and P. Dollár. Learning to segment object candidates. In *Advances in Neural Information Processing Systems*, pp. 1990–1998, (2015).
 56. E. Xie, P. Sun, X. Song, W. Wang, X. Liu, D. Liang, C. Shen, and P. Luo. PolarMask: Single shot instance segmentation with polar representation, *arXiv preprint arXiv:1909.13226*. (2019).
 57. Z. Hayder, X. He, and M. Salzmann. Boundary-aware instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5696–5704, (2017).
 58. Y. Lee and J. Park. CenterMask: Real-time anchor-free instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 13906–13915, (2020).
 59. M. Bai and R. Urtasun. Deep watershed transform for instance segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5221–5229, (2017).
 60. D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee. YOLACT: Real-time instance segmentation. In *IEEE International Conference on Computer Vision*, pp. 9157–9166, (2019).
 61. A. Fathi, Z. Wojna, V. Rathod, P. Wang, H. O. Song, S. Guadarrama, and K. P. Murphy. Semantic instance segmentation via deep metric learning, *arXiv preprint arXiv:1703.10277*. (2017).
 62. L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and

- fully connected crfs, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **40**(4), 834–848, (2017).
63. F. Yu and V. Koltun, Multi-scale context aggregation by dilated convolutions, *arXiv preprint arXiv:1511.07122*. (2015).
 64. P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell. Understanding convolution for semantic segmentation. In *IEEE Winter Conference on Applications of Computer Vision*, pp. 1451–1460, (2018).
 65. M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang. DenseASPP for semantic segmentation in street scenes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3684–3692, (2018).
 66. A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, ENet: A deep neural network architecture for real-time semantic segmentation, *arXiv preprint arXiv:1606.02147*. (2016).
 67. L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, Rethinking atrous convolution for semantic image segmentation, *arXiv preprint arXiv:1706.05587*. (2017).
 68. L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *European Conference on Computer Vision*, pp. 801–818, (2018).
 69. F. Visin, M. Ciccone, A. Romero, K. Kastner, K. Cho, Y. Bengio, M. Matteucci, and A. Courville. ReSeg: A recurrent neural network-based model for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 41–48, (2016).
 70. F. Visin, K. Kastner, K. Cho, M. Matteucci, A. Courville, and Y. Bengio, ReNet: A recurrent neural network based alternative to convolutional networks, *arXiv preprint arXiv:1505.00393*. (2015).
 71. W. Byeon, T. M. Breuel, F. Raue, and M. Liwicki. Scene labeling with LSTM recurrent neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3547–3555, (2015).
 72. X. Liang, X. Shen, J. Feng, L. Lin, and S. Yan. Semantic object parsing with graph LSTM. In *European Conference on Computer Vision*, pp. 125–143. Springer, (2016).
 73. Y. Xiang and D. Fox, DA-RNN: Semantic mapping with data associated recurrent neural networks, *arXiv:1703.03098*. (2017).
 74. R. Hu, M. Rohrbach, and T. Darrell. Segmentation from natural language expressions. In *European Conference on Computer Vision*, pp. 108–124. Springer, (2016).
 75. L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille. Attention to scale: Scale-aware semantic image segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3640–3649, (2016).
 76. Q. Huang, C. Xia, C. Wu, S. Li, Y. Wang, Y. Song, and C.-C. J. Kuo, Semantic segmentation with reverse attention, *arXiv preprint arXiv:1707.06426*. (2017).
 77. H. Li, P. Xiong, J. An, and L. Wang, Pyramid attention network for semantic segmentation, *arXiv preprint arXiv:1805.10180*. (2018).
 78. J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu. Dual attention network for scene segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3146–3154, (2019).
 79. Y. Yuan and J. Wang, OCNet: Object context network for scene parsing, *arXiv preprint arXiv:1809.00916*. (2018).
 80. H. Zhang, C. Wu, Z. Zhang, Y. Zhu, Z. Zhang, H. Lin, Y. Sun, T. He, J. Mueller, R. Manmatha, et al., Resnest: Split-attention networks, *arXiv preprint arXiv:2004.08955*. (2020).
 81. S. Choi, J. T. Kim, and J. Choo. Cars can't fly up in the sky: Improving urban-scene

- segmentation via height-driven attention networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9373–9383, (2020).
82. X. Li, Z. Zhong, J. Wu, Y. Yang, Z. Lin, and H. Liu. Expectation-maximization attention networks for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 9167–9176, (2019).
 83. Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu. CCNet: Criss-cross attention for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 603–612, (2019).
 84. M. Ren and R. S. Zemel. End-to-end instance segmentation with recurrent attention. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6656–6664, (2017).
 85. H. Zhao, Y. Zhang, S. Liu, J. Shi, C. Change Loy, D. Lin, and J. Jia. PSANet: Point-wise spatial attention network for scene parsing. In *European Conference on Computer Vision*, pp. 267–283, (2018).
 86. C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. Learning a discriminative feature network for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1857–1866, (2018).
 87. P. Luc, C. Couprie, S. Chintala, and J. Verbeek, Semantic segmentation using adversarial networks, *arXiv preprint arXiv:1611.08408*. (2016).
 88. N. Souly, C. Spampinato, and M. Shah. Semi supervised semantic segmentation using generative adversarial network. In *IEEE International Conference on Computer Vision*, pp. 5688–5696, (2017).
 89. W.-C. Hung, Y.-H. Tsai, Y.-T. Liou, Y.-Y. Lin, and M.-H. Yang, Adversarial learning for semi-supervised semantic segmentation, *arXiv preprint arXiv:1802.07934*. (2018).
 90. Y. Xue, T. Xu, H. Zhang, L. R. Long, and X. Huang, SegAN: Adversarial network with multi-scale L1 loss for medical image segmentation, *Neuroinformatics*. **16**(3-4), 383–392, (2018).
 91. M. Majurski, P. Manescu, S. Padi, N. Schaub, N. Hotaling, C. Simon Jr, and P. Bajcsy. Cell image segmentation using generative adversarial networks, transfer learning, and augmentations. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, (2019).
 92. K. Ehsani, R. Mottaghi, and A. Farhadi. SegAN: Segmenting and generating the invisible. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6144–6153, (2018).
 93. T. F. Chan and L. A. Vese, Active contours without edges, *IEEE Transactions on Image Processing*. **10**(2), 266–277, (2001).
 94. X. Chen, B. M. Williams, S. R. Vallabhaneni, G. Czanner, R. Williams, and Y. Zheng. Learning active contour models for medical image segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11632–11640, (2019).
 95. S. Gur, L. Wolf, L. Golgher, and P. Blinder. Unsupervised microvascular image segmentation using an active contours mimicking neural network. In *IEEE International Conference on Computer Vision*, pp. 10722–10731, (2019).
 96. P. Marquez-Neila, L. Baumela, and L. Alvarez, A morphological approach to curvature-based evolution of curves and surfaces, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. **36**(1), 2–17, (2014).
 97. T. H. N. Le, K. G. Quach, K. Luu, C. N. Duong, and M. Savvides, Reformulating level sets as deep recurrent neural network approach to semantic segmentation, *IEEE Transactions on Image Processing*. **27**(5), 2393–2407, (2018).
 98. C. Rupprecht, E. Huaroc, M. Baust, and N. Navab, Deep active contours, *arXiv preprint arXiv:1607.05074*. (2016).

99. A. Hatamizadeh, A. Hoogi, D. Sengupta, W. Lu, B. Wilcox, D. L. Rubin, and D. Terzopoulos. Deep active lesion segmentation. In *International Workshop on Machine Learning in Medical Imaging*, vol. 11861, *Lecture Notes in Computer Science*, pp. 98–105. Springer, (2019).
100. D. Marcos, D. Tuia, B. Kellenberger, L. Zhang, M. Bai, R. Liao, and R. Urtasun. Learning deep structured active contours end-to-end. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8877–8885, (2018).
101. D. Cheng, R. Liao, S. Fidler, and R. Urtasun. DARNet: Deep active ray network for building segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7431–7439, (2019).
102. A. Hatamizadeh, D. Sengupta, and D. Terzopoulos. End-to-end trainable deep active contour models for automated image segmentation: Delineating buildings in aerial imagery. In *European Conference on Computer Vision*, pp. 730–746, (2020).
103. H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal. Context encoding for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7151–7160, (2018).
104. G. Lin, A. Milan, C. Shen, and I. Reid. RefineNet: Multi-path refinement networks for high-resolution semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1925–1934, (2017).
105. G. Song, H. Myeong, and K. Mu Lee. SeedNet: Automatic seed generation with deep reinforcement learning for robust interactive segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1760–1768, (2018).
106. J. Dai, K. He, and J. Sun. BoxSup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 1635–1643, (2015).
107. C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun. Large kernel matters — improve semantic segmentation by global convolutional network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4353–4361, (2017).
108. Z. Wu, C. Shen, and A. Van Den Hengel. Wider or deeper: Revisiting the resnet model for visual recognition, *Pattern Recognition*. **90**, 119–133, (2019).
109. Z. Zhang, X. Zhang, C. Peng, X. Xue, and J. Sun. ExFuse: Enhancing feature fusion for semantic segmentation. In *European Conference on Computer Vision*, pp. 269–284, (2018).
110. M. Mostajabi, P. Yadollahpour, and G. Shakhnarovich. Feedforward semantic segmentation with zoom-out features. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3376–3385, (2015).
111. W. Wang, J. Shen, and F. Porikli. Saliency-aware geodesic video object segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3395–3402, (2015).
112. P. Luo, G. Wang, L. Lin, and X. Wang. Deep dual learning for semantic image segmentation. In *IEEE International Conference on Computer Vision*, pp. 2718–2726, (2017).
113. X. Li, Z. Jie, W. Wang, C. Liu, J. Yang, X. Shen, Z. Lin, Q. Chen, S. Yan, and J. Feng. FoveaNet: Perspective-aware urban scene parsing. In *IEEE International Conference on Computer Vision*, pp. 784–792, (2017).
114. I. Kreso, S. Segvic, and J. Krapac. Ladder-style densenets for semantic segmentation of large natural images. In *IEEE International Conference on Computer Vision*, pp. 238–245, (2017).
115. C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang. BiSeNet: Bilateral segmentation network for real-time semantic segmentation. In *European Conference on Computer*

- 22 S. Minaee, Y. Boykov, F. Porikli, A. Plaza, N. Kehtarnavaz, and D. Terzopoulos
Vision, pp. 325–341, (2018).
116. B. Cheng, L.-C. Chen, Y. Wei, Y. Zhu, Z. Huang, J. Xiong, T. S. Huang, W.-M. Hwu, and H. Shi. SPGNet: Semantic prediction guidance for scene parsing. In *IEEE International Conference on Computer Vision*, pp. 5218–5228, (2019).
 117. T. Takikawa, D. Acuna, V. Jampani, and S. Fidler. Gated-SCNN: Gated shape cnns for semantic segmentation. In *IEEE International Conference on Computer Vision*, pp. 5229–5238, (2019).
 118. J. Fu, J. Liu, Y. Wang, Y. Li, Y. Bao, J. Tang, and H. Lu. Adaptive context network for scene parsing. In *IEEE International Conference on Computer Vision*, pp. 6748–6757, (2019).
 119. X. Liang, H. Zhou, and E. Xing. Dynamic-structured semantic propagation network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 752–761, (2018).
 120. X. Liang, Z. Hu, H. Zhang, L. Lin, and E. P. Xing. Symbolic graph reasoning meets convolutions. In *Advances in Neural Information Processing Systems*, pp. 1853–1863, (2018).
 121. B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ADE20K dataset. In *IEEE Conference on Computer Vision and Pattern Recognition*, (2017).
 122. R. Zhang, S. Tang, Y. Zhang, J. Li, and S. Yan. Scale-adaptive convolutions for scene parsing. In *IEEE International Conference on Computer Vision*, pp. 2031–2039, (2017).
 123. T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun. Unified perceptual parsing for scene understanding. In *European Conference on Computer Vision*, pp. 418–434, (2018).
 124. B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. V. Le. Rethinking pre-training and self-training, *arXiv preprint arXiv:2006.06882*. (2020).
 125. X. Zhang, H. Xu, H. Mo, J. Tan, C. Yang, and W. Ren. DCNAS: Densely connected neural architecture search for semantic image segmentation, *arXiv preprint arXiv:2003.11883*. (2020).
 126. A. Tao, K. Sapra, and B. Catanzaro. Hierarchical multi-scale attention for semantic segmentation, *arXiv preprint arXiv:2005.10821*. (2020).
 127. D. Seichter, M. Köhler, B. Lewandowski, T. Wengefeld, and H.-M. Gross. Efficient rgb-d semantic segmentation for indoor scene analysis, *arXiv preprint arXiv:2011.06961*. (2020).
 128. M. Zhen, J. Wang, L. Zhou, S. Li, T. Shen, J. Shang, T. Fang, and L. Quan. Joint semantic segmentation and boundary detection using iterative pyramid contexts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13666–13675, (2020).
 129. Y. Li, L. Song, Y. Chen, Z. Li, X. Zhang, X. Wang, and J. Sun. Learning dynamic routing for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8553–8562, (2020).
 130. A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár. Panoptic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9404–9413, (2019).
 131. A. Kirillov, R. Girshick, K. He, and P. Dollar. Panoptic feature pyramid networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6399–6408, (2019).
 132. Y. Li, X. Chen, Z. Zhu, L. Xie, G. Huang, D. Du, and X. Wang. Attention-guided unified network for panoptic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, (2019).

133. L. Porzi, S. R. Buló, A. Colovic, and P. Kotschieder. Seamless scene segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8277–8286, (2019).
134. B. Cheng, M. D. Collins, Y. Zhu, T. Liu, T. S. Huang, H. Adam, and L.-C. Chen, Panoptic-DeepLab, *arXiv preprint arXiv:1910.04751*. (2019).
135. Y. Xiong, R. Liao, H. Zhao, R. Hu, M. Bai, E. Yumer, and R. Urtasun. UPSNet: A unified panoptic segmentation network. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8818–8826, (2019).
136. R. Mohan and A. Valada, EfficientPS: Efficient panoptic segmentation, *arXiv preprint arXiv:2004.02307*. (2020).
137. Z.-H. Zhou, A brief introduction to weakly supervised learning, *National Science Review*. **5**(1), 44–53, (2018).
138. L. Jing and Y. Tian, Self-supervised visual feature learning with deep neural networks: A survey, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. (2020).
139. V. Goel, J. Weng, and P. Poupart. Unsupervised video object segmentation for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5683–5694, (2018).
140. L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson, Deep learning in remote sensing applications: A meta-analysis and review, *ISPRS Journal of Photogrammetry and Remote Sensing*. **152**, 166 – 177, (2019).
141. L. Gao, Y. Zhang, F. Zou, J. Shao, and J. Lai, Unsupervised urban scene segmentation via domain adaptation, *Neurocomputing*. **406**, 295 – 301, (2020).
142. M. Paoletti, J. Haut, J. Plaza, and A. Plaza, Deep learning classifiers for hyperspectral imaging: A review, *ISPRS Journal of Photogrammetry and Remote Sensing*. **158**, 279 – 317, (2019).
143. J. F. Abrams, A. Vashishtha, S. T. Wong, A. Nguyen, A. Mohamed, S. Wieser, A. Kuijper, A. Wilting, and A. Mukhopadhyay, Habitat-Net: Segmentation of habitat images using deep learning, *Ecological Informatics*. **51**, 121 – 128, (2019).
144. M. Kerkech, A. Hafiane, and R. Canals, Vine disease detection in UAV multispectral images using optimized image registration and deep learning segmentation approach, *Computers and Electronics in Agriculture*. **174**, 105446, (2020).
145. Y. Song, Z. Huang, C. Shen, H. Shi, and D. A. Lange, Deep learning-based automated image segmentation for concrete petrographic analysis, *Cement and Concrete Research*. **135**, 106118, (2020).
146. N. Tajbakhsh, L. Jeyaseelan, Q. Li, J. N. Chiang, Z. Wu, and X. Ding, Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation, *Medical Image Analysis*. **63**, 101693, (2020).
147. A. Amyar, R. Modzelewski, H. Li, and S. Ruan, Multi-task deep learning based CT imaging analysis for COVID-19 pneumonia: Classification and segmentation, *Computers in Biology and Medicine*. **126**, 104037, (2020).