

Learning Deep Representations for Semantic Image Parsing: a Comprehensive Overview

Lili Huang, Jiefeng Peng, Ruimao Zhang, Guanbin Li, Liang Lin

School of Data and Computer Science, Sun Yat-Sen University

Abstract

Semantic image parsing, which refers to the process of decomposing images into semantic regions and constructing the structure representation of the input, has recently aroused widespread interest in the field of computer vision. The recent application of deep representation learning has driven this field into a new stage of development. In this paper, we summarize three aspects of the progress of research on semantic image parsing, i.e., category-level semantic segmentation, instance-level semantic segmentation, and beyond segmentation. Specifically, we first review the general frameworks for each task and introduce the relevant variants. The advantages and limitations of each method are also discussed. Moreover, we present a comprehensive comparison of different benchmark datasets and evaluation metrics. Finally, we explore the future trends and challenges of semantic image parsing.

1 Introduction

1.1 Semantic Image Parsing

With the development of Internet, in recent years, large-scale image and multimedia video data have increased explosively, resulting in urgent demands for advanced intelligent image analysis technology, such as semantic image parsing. As a fundamental and long-standing problem in computer vision, semantic image parsing is performed at three levels, which will be discussed below.

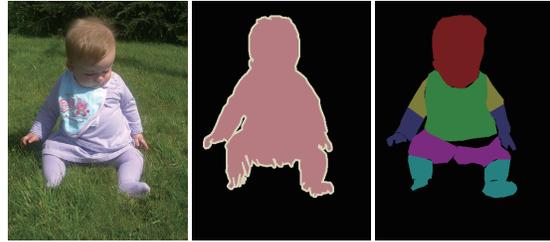


Figure 1: Illustration of the category-level semantic segmentation. The left is the original image. The middle is the basic semantic segmentation result, and the right is the semantic part segmentation result.

i. **Category-level semantic segmentation.** It attempts to assign a single category label to each pixel. Here, a category label corresponds to a specific object category or a local part of the object. Therefore, category-level semantic segmentation consists of basic semantic segmentation and semantic part segmentation (called object parsing in the literature), as illustrated in Fig. 1. The former predicts the segmentation mask and its label for the entire object, as shown in the middle of Fig. 1, while the latter refers to segmenting an object into its constituent semantic parts and predicting the segmentation mask for each local part, as shown on the right side of Fig. 1. According to the definition, part segmentation can be regarded as a special type of fine-grained category-level semantic segmentation task.

Category-level semantic segmentation is actually a pixel-wise dense prediction problem, which is supported by two key technologies: 1) classification: an

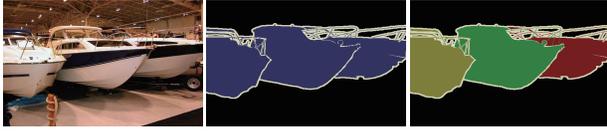


Figure 2: Comparison of category-level and instance-level semantic segmentation. The left is the original image. The middle is the category-level semantic segmentation result, and the right is the instance-level semantic segmentation result.

object is assigned a specific semantic-concept label; and 2) localization: the classification label for a pixel must match the appropriate coordinates in the output score map [1].

ii. **Instance-level semantic segmentation.** In contrast to category-level segmentation, it requires precise segmentation of each object and correct detection of all the object instances in one image [2]. In the middle of Fig. 2, three boats are segmented by assigning the same category label (i.e., "boat"). Clearly, category-level segmentation cannot distinguish the object instances belonging to the same category. In the right column of Fig. 2, the three boats are segmented by assigning different IDs with the same category label (i.e., "boat one", "boat two", and "boat three"). Thus, the instance-level segmentation requires support from both classification and detection technologies.

iii. **Beyond segmentation.** In recent years, works extending beyond semantic segmentation have also received substantial attention. This task is inspired by previous work on image parsing [2], which refers to the process of decomposing an image into its constituent visual structured configuration [3, 4, 5]. Works beyond segmentation not only semantically segment images but also predict richer and finer results, such as the structures and relations of objects and the spatial layout. Specifically, images are decomposed into semantic regions and the structures and relationships among objects are constructed. For example, in Fig. 3, the image caption is "there is one person sitting on the chair nearby the table with one monitor". Following the work in [6], the beyond segmentation method first segments all the ob-

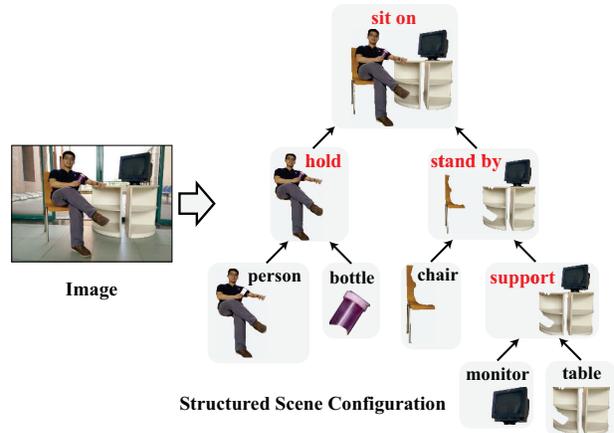


Figure 3: Illustration of beyond segmentation. (Figure extracted from [6])

jects (i.e., "person", "chair", "table", and "bottle") in the image, predicts the relations among objects (i.e., "hold", "stand by", "support", and "sit on"), and finally estimates the hierarchical structures. Intuitively, works beyond segmentation produces detailed parsing results that are consistent with human perception.

Similar to most vision problems, the discriminant features greatly affect the performance of semantic image parsing. Traditional semantic segmentation methods adopt hand-crafted features, such as SIFT [7], HOG [8], and LBP [9]. However, these hand-crafted features are not applicable to various tasks. Therefore, the automatic extraction of valuable information and effective representation of image/video data are critical. Representation learning, i.e., learning representations of data, makes it easier to extract useful information from raw data to build predictors. The representation algorithms for semantic image parsing have experienced three periods of progress in the continuous improvement of image parsing performance: 1) traditional hand-crafted methods; 2) deep learning, such as convolutional neural networks (CNNs), recurrent neural networks (RNNs); and 3) the integration of the two methods to complement each other.

Extensive experiments [10, 11, 12, 13, 2, 14, 1] have demonstrated that the representation ability of traditional hand-crafted features is insufficient. Meanwhile, deep learning currently achieves the best representation ability and has had tremendous success in many applications, such as image classification [15], object detection, and natural language understanding [16]. Therefore, we list only the main differences among the three-level semantic segmentation tasks accomplished by deep representation models, as illustrated in Table 1.

1.2 Deep Learning

Deep learning is defined as learning multiple levels of representations from the local and detailed levels in the shallow layers to the global and abstract levels in the deeper layers [18, 19, 20]. Specifically, deep neural networks consist of several simple but non-linear modules, each of which transforms the simple representation at the shallow layer (starting with the raw input) into slightly more abstract representation at the deep layer. Several well-known deep neural networks, such as the CNN, recurrent neural network, and RNN, have been reported in recent years. Moreover, abundant variants of these networks, which we discuss in the following sections, have emerged.

Convolutional Neural Networks. The CNN [21] is designed for data with grid-like structures and consists of convolutional layers, pooling layers, and non-linear rectification layers. The units in the neural network are locally connected, which results in shared weights of the local parameters and features in the deeper abstract layers being invariant to local image transformation. Despite the numerous applications of CNNs, they were not well-known until their successful application to object recognition during the ImageNet challenge in 2012. Then, CNN was quickly applied to semantic segmentation [12, 22, 2, 14, 1, 17, 23] and achieved great successes.

Recurrent Neural Networks. In contrast to CNNs, which are tailored for grid-structure data [21], recurrent neural networks are more appropriate for sequential data [24]. The principal characteristic of a recurrent neural network is that neurons (units) are connected by synaptic links to express tempo-

ral relations. To alleviate the explosion or vanishing of the backpropagated gradients in the shallow layers [25, 24], long short-term memory (LSTM) networks [26] were proposed by introducing special hidden units to memorize the observed knowledge of the previous and current inputs. The success of LSTM has demonstrated that LSTM is more effective than conventional recurrent neural networks in image captioning [27] and machine translation [28]. Additionally, many works [29, 26, 30, 31, 32, 33, 34] utilize LSTM to improve the performance of semantic image parsing.

Recursive Neural Networks. Unlike the aforementioned recurrent neural networks [35], which are designed for time sequential data, RNNs [16] are designed for hierarchical space structural data. Recurrent neural networks for chain structures by connecting hidden units, whereas RNNs recursively form a hierarchical structure because the structures of networks are similar at every level of the hierarchy. This characteristic is in line with the structures of natural language, which results in successful natural language parsing [16]. Some recent works [6, 16] proposed RNNs for structural semantic parsing.

1.3 Our Contribution to the Existing Surveys

With a unique perspective, this work comprehensively reviews deep representation learning-based semantic image parsing at three levels: category-level semantic segmentation, instance-level semantic segmentation, and beyond segmentation. Specifically, for each level of semantic segmentation, we elaborate the relative terminology and background knowledge. Furthermore, this paper reviews and compares existing models and relatively well-known datasets and evaluation metrics. To the best of our knowledge, there is no such overview of semantic image parsing in the literature.

The rest of this article is organized as follows. In Section 2, we review deep representations for semantic image parsing at three levels. Datasets and evaluation metrics are introduced in Section 3. Finally, we present the conclusions and discuss promising future research directions in Section 4.

Task	Flourishing period	Pioneering work	Key technology	Type of labels
Category-level segmentation	2015	FCN [12]	classification, localization	object, part
Instance-level segmentation	2016	FCIS [17]	classification, detection	instance
Beyond segmentation	2016	CNN-RNN [6]	classification, localization	object, part, relation, scene structure

Table 1: Comparisons of different semantic segmentation tasks performed by deep models

2 Learning Deep Representations

In previous decades, most of the successful semantic segmentation algorithms have relied on hand-crafted features combined with flat classifiers, such as boosting [36] and support vector machines [37]. Nevertheless, the performance of these algorithms is compromised by the limited feature expression.

More recently, with the emergence of big data and development of computer hardware, deep neural networks have reached their prime. In the field of computer vision, deep learning has achieved great success in image classification [12, 38, 15, 21, 39], recurrent neural networks have made tremendous achievements in expressing temporal relations [40, 41, 31, 35, 42], and RNNs have succeeded in terms of space structure relationship representation [16, 6]. The breakthroughs of deep learning in image classification are quickly repurposed to semantic image parsing. We illustrate this problem at different levels of image segmentation, i.e., category-level semantic segmentation, instance-level semantic segmentation, and beyond segmentation, in the following sections.

2.1 Category-Level Semantic Segmentation

As mentioned in Section 1, category-level semantic segmentation attempts to assign a single category label to each pixel, i.e., basic semantic segmentation and semantic part segmentation, as illustrated in Fig. 1. For convenience, we do not distinguish between these two processes. The category-level deep models for semantic segmentation are mainly divided into two types: region-based networks and fully convolu-

tional frameworks.

Region-based Networks. The previously reported deep models [43, 44, 45] are mainly region-based networks that classify each pixel by using its enclosing region for training and prediction. These methods have several limitations. First, they treat each region or pixel as a separate unit. On the one hand, this treatment ignores the importance of the context information in pixel labeling inference, while on the other hand, it ignores the spatial correlation in the image and reduces the algorithm accuracy. Secondly, the independent processing of thousands of regions results in substantial overhead and inefficiency.

Fully Convolutional Frameworks. The fully convolutional frameworks for semantic segmentation consist of two fundamental works, i.e., fully convolutional networks (FCNs) [12] and the DeepLab system [13], both of which fully utilize convolutional networks to produce spatially dense predictions.

In [12], as the fundamental application of the CNN architecture [46, 47], the authors devised FCNs for spatially dense prediction tasks by accommodating the prior advance deep networks [15, 39, 37]. Specifically, as illustrated in Fig. 4, fully connected layers in prior networks are converted into convolutional layers, and the deconvolutional layers are built by up-sampling intermediate feature maps to keep the size of the output the same as that of the input image. However, the spatial resolution of the feature maps is reduced after the consecutive combination of the max-pooling and downsampling layers in FCN, as in prior image classification models. A novel skip architecture was devised to fuse semantic information with appearance information to produce accurate and detailed segmentations [12], as shown in Fig. 5. The semantic information comes from a deep, coarse layer,

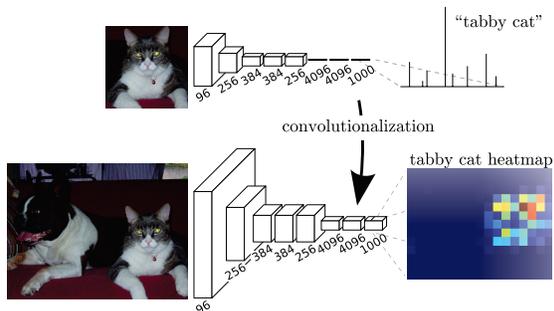


Figure 4: Illustration of the adaptation of fully connected layers into convolutional layers. (Figure extracted from [12])

while the appearance information is from a shallow, fine layer. Thus, by utilizing prior image classification models as pre-trained models, FCN is fine-tuned to learn and inference efficiently in an end-to-end manner, resulting in equivalently sized output.

Another fundamental work—DeepLab system [13]—integrated CNN with fully connected conditional random field (CRF) to expand and improve FCN [12]. As shown in Fig. 6, the responses at the final convolutional layer are fed into the fully connected CRF to capture finer details. Thus, the fully connected CRF refines the raw CNN scores, especially along object boundaries. However, the DeepLab system treats CNN and CRF as two separate components. Concretely, fully connected CRF utilizes the Gaussian CRF potentials [48] to capture long-range dependencies by treating every pixel as a CRF node to receive unary potentials.

Many subsequent variants emerges from these two fundamental works. These works generally evolve along three directions: CNN crafting tricks, integration with the random field model, and integration with recurrent neural networks. We discuss these three aspects below.

2.1.1 CNN Crafting Tricks

The majority of deep learning algorithms are based on CNNs; therefore, one intuitive fundamental idea is to design more efficient network architecture with

CNN crafting tricks, such as downsample-upsample operation, pyramid module, skip connection, and atrous convolution.

Downsample-Upsample Operations. A downsample-upsample operation is composed of two stages: downsampling and upsampling. In the downsampling stage, the feature maps are processed by convolution or unpooling and progressively shrink to smaller maps, where the receptive field of every pixel is gradually enlarged. In the upsampling stage, the object spatial dimension is recovered through deconvolution or unpooling, where the coarse-to-fine details are captured.

DeconvNet [49] treats the convolutional layers of the VGG 16-layer net as the downsampling stage, whereas the developed deconvolution network serves as the upsampling stage, which consists of deconvolution and unpooling layers to increase the resolution of small score maps with more detailed structures. Specifically, DeconvNet first generates sufficient instance-wise candidate proposals for each given image at the downsampling stage, and produces the semantic segmentation maps of each proposal at the upsampling stage. Then, the final semantic segmentation of the whole input image is obtained by assembling the maps of all proposals with non-maximum suppression. Furthermore, DeconvNet [49] is integrated with FCN [12] to improve the performance.

Similar to DeconvNet [49], SegNet [50] also introduces the unpooling operation without ReLU in the upsampling stage to recover the spatial dimensions, and the downsampling and upsampling correspond to encoder and decoder stacks, respectively. Specifically, the encoder stacks, composed of convolutions, ReLU and max-pooling, produce low-resolution feature maps while simultaneously memorizing the pooled indices. Then, the decoder stacks upsample the low-resolution maps using the pooled indices and output the semantic segmentation.

The more complicated contextualized convolutional neural network (Co-CNN) [47] is a novel downsample-upsample framework that simultaneously captures hierarchical information by seamlessly integrating three levels of context (i.e., cross-layer context, global image-level context, local super-pixel context) into a unified network, as shown in Fig.

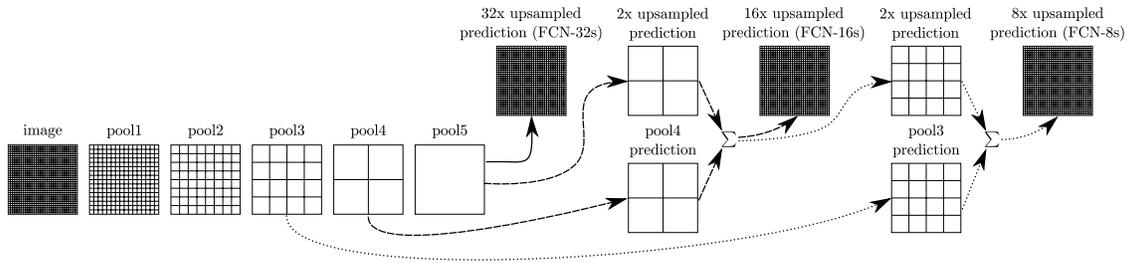


Figure 5: The "skip" architecture of FCN. (Figure extracted from [12])

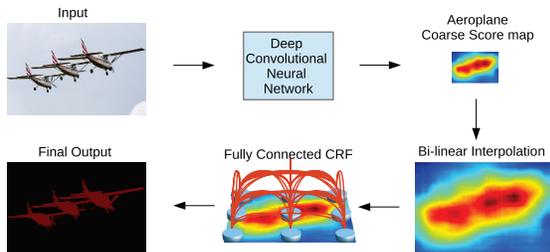


Figure 6: Illustration of the DeepLab system. (Figure extracted from [13])

7. Specifically, Co-CNN first utilizes convolutional networks to obtain the downsampled feature maps for multiple resolutions, upsamples the feature maps along with multi-level context generation, and finally produces pixel-wise predictions. Moreover, the cross-layer context, global image-level context and local super-pixel context are generated by integrating the hierarchical structure, predicting the global image-level labels, and refining super-pixels, respectively.

In general, downsampling is used to extract features from the input image, whereas upsampling produces object segmentation from the features extracted by downsampling. The seamlessly integration of downsampling with upsampling elegantly accomplishes the semantic segmentation task.

Pyramid Module. The pyramid module consists of two varieties: 1) input pyramid, where multi-scale inputs are fed into the same model with shared weights such that the large-scale inputs maintain more fine details and the small-scale inputs capture longer range information; and 2) pooling pyra-

mid, where context information is captured by spatial pyramid pooling in several ranges.

DeepLabV2 [11], the updated DeepLab system [13], employs both types of pyramid modules. On the one hand, DeepLabV2 first transforms the inputs into several scale inputs that are synchronously fed into the weight-shared CNN to produce multi-scale feature maps, which are then merged. On the other hand, DeepLabV2 segments objects at multiple scales via atrous spatial pyramid pooling (ASPP), which serves to resample features prior to convolution. Specifically, ASPP takes advantage of several parallel atrous convolutions with diverse sampling rates to capture multi-scale objects and image context.

Zhao et al. proposed a superior framework—pyramid scene parsing network (PSPNet) [1]—for scene parsing in complex scenes. PSPNet [1] adopts the pyramid pooling module to capture the global context representation and prevent the loss of context information between subregions. Specifically, the pyramid pooling module employs multiple pyramid scales to generate coarse to fine feature maps, which provide additional multi-scale contextual information from different regions. Then, the different-region-based context information is aggregated to capture the global context representation.

Apparently, the pyramid module captures multi-scale context information from local fine to global abstract to improve the performance of semantic segmentation.

Skip Connection. Similar to its first application in FCN [12], skip connection refers to the links

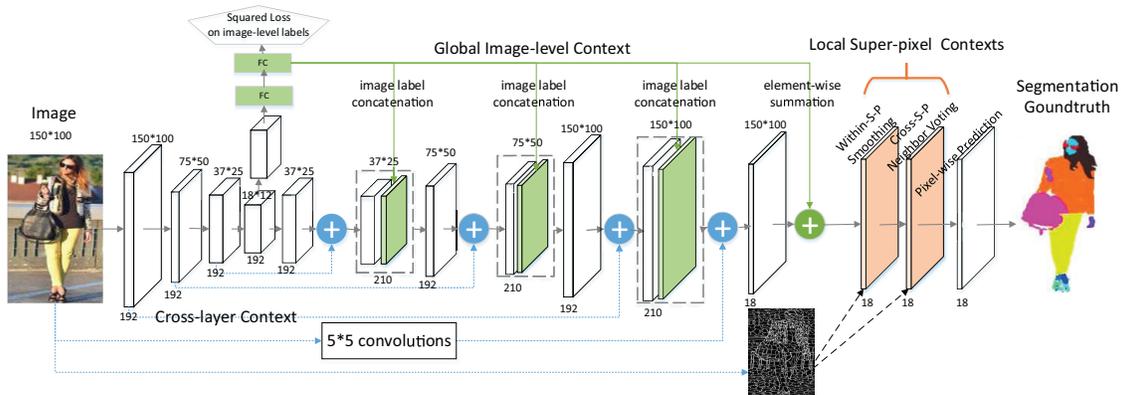


Figure 7: Illustration of the Co-CNN. (Figure extracted from [47])

between low-level layers and high-level layers at an interval of several layers; thus, detailed appearance features from shallow layers are combined with coarse semantic information from deep layers to improve the segmentation performance.

The global convolutional network (GCN) [14] with large-size kernels utilizes pretrained ResNet [38] as the feature network and FCN [12] as the segmentation framework. Specifically, the GCN and boundary refinement block are both treated as residual structures. In the feature network, each stage of the ResNet block generates different-scale feature maps, which are fed into the GCN structures to produce semantic score maps for each category. Additionally, the boundary refinement blocks are used to further refine the object boundaries. Next, outputs from the top layer of the residual structures are passed to the segmentation framework, and new high-resolution score maps are generated iteratively by skip connection [12]. Specifically, upsampled score maps in the higher layers are iteratively combined with the corresponding-resolution score maps extracted from the residual structures in the lower layers. Finally, the semantic score map, which is used to output pixel-wise semantic labels, is generated after the last upsampling.

On the basis of FCN [12], U-Net [51] proposed a u-shaped architecture composed of a contracting path and a symmetric expanding path to effectively train

deep models on small datasets. Specifically, the contracting path is similar to the typical convolution architecture used to extract and downsample the feature maps. The lowest-resolution feature maps flows into the expanding path, where the feature maps at each step are upsampled and concatenated with the same resolution feature maps cropped from the contracting path. Thus, the final segmentation maps for each category are generated after the top layer in the expanding path. The cropping step is applied to prevent the loss of border pixels during convolution operations.

Islam et al. [23] proposed the label refinement network (LRN) to improve segmentation performance by predicting segmentation labels at multiple resolutions. The LRN is formulated as an encoder-decoder framework [12, 50, 49], where the VGG16 network serves as the encoder network to extract feature maps with decreasing resolution and the decoder network predicts multi-scale coarse-to-fine label maps in several stages. The skip connection architecture combines the label maps of each stage with the corresponding feature maps in the encoder network to refine the segmentation labels. Furthermore, the LRN [23] supervises the predictions at different stages by defining a loss function for each stage.

Lin et al. [52] devised a multi-path refinement network, called RefineNet, for semantic segmentation. The cascaded architecture exploits multi-scale fea-

tures from different stages of ResNet [38] and conveys them into different stages of the RefineNet block via long-range skip connections. The RefineNet block is applied to upsample feature maps and to recover the decreased resolution through local residual connections and chained residual pooling. The long-range skip connections are used to integrate information from the coarse high-level deep layers and the fine low-level shallow layers to produce high-resolution semantic feature maps; thus, the gradient can be directly propagated to the inputs, preventing gradient vanishing and explosion.

In conclusion, skip connection structure merges hierarchical cross-layer features to improve the segmentation performance, and the gradient can be propagated backward along both the skip path and the cascaded original path to prevent gradient vanishing and explosion.

Atrous Convolution. Atrous convolution [13, 11], also called dilated convolution [53, 54], refers to convolution with an atrous rate. The rate corresponds to the stride with which the input signals are sampled. Thus, standard convolution, with a rate of 1, is a special case of atrous convolution.

The fundamental work on the DeepLab system [13, 11] first proposed the definition of atrous convolution and utilizes atrous convolution to simplify the architecture of FCN [12]. In this work, atrous convolution is constructed via convolutions with up-sampled filters. In an atrous convolution operation, the incoming input feature maps are sampled by enlarging the input stride values, resulting in enlarged field of view of filters and feature responses.

Instead of the atrous convolution with a dilated filter in DeepLab [13, 11], Yu et al. proposed a special tailored atrous convolution in Dilated-Net [54] to obtain multi-scale contextual information. Specifically, the atrous convolution in Dilated-Net is built by recomposing the convolution operator itself with dilation factors and is free from the dilated filters in DeepLab. The dilated convolution operator with different dilation factors can adopt the same filter in different ranges to capture multi-scale context. Moreover, the receptive fields are enlarged exponentially without loss of resolution, whereas the parameters in the network grow linearly.

What can be inferred from the aforementioned three works [13, 11, 54] is that, atrous convolutions can take control of the field of view of the convolution filters and feature responses without additional computation overhead.

Coarse-to-Fine Refinement. Coarse-to-fine refinement exploits cascade or supplementary structures to refine confidence maps from coarse to fine.

Active template regression (ATR) [46], which directly predicts and locates the structural masks for each label, was proposed for human parsing. The structural outputs consist of the mask template coefficients and the shape parameters. ATR builds the end-to-end relations between the input image and the structural outputs by devising two separate CNNs, i.e., a template network and a shape network. The template coefficients are predicted by the template network with max-pooling to capture the contextual correlations among all label masks. Meanwhile, the shape parameters are predicted by the shape network without max-pooling to maintain the sensitivity to the label mask position. The outputs from the two parallel CNNs provide supplementary information. Thus, the normalized mask of each semantic region is expressed as a linear combination of the learned mask templates and is then refined to a more precise mask with the shape parameters.

Li et al. [55] proposed an end-to-end deep layer cascade (LC) framework to improve the accuracy and speed of semantic segmentation. Specifically, LC treats different layers in the deep network as different stages with difficulty-aware learning. The early lower stages are trained to handle easy regions, while the challenging regions are forward propagated to the subsequent higher stages; thus, the prediction process is coarse to fine. Furthermore, dilated convolutions are used on the propagated regions to reduce the computations.

Similar to the LC framework [55], Zhou et al. [56] proposed a cascaded fixed-point model for small organ segmentation in a coarse-to-fine manner. The entire input region is fed into a coarse-scaled network to produce the coarse segmentation mask, based on which a small region is generated via a transformation function. Then, the small region serves as the input of the subsequent fine-scaled network to produce a

more accurate segmentation result. The fixed-point model is iteratively optimized by means of the strategy in [57].

Wang et al. [58] proposed a weakly supervised model, image descriptions in the wild CNN (IDW-CNN), to improve segmentation performance using object interactions and descriptions. The architecture of IDW-CNN is composed of three components, i.e., the feature extraction procedure, segmentation stream (Seg-stream) and object interaction stream (Int-stream). First, ResNet-101 is used to extract features. The Int-stream takes these features as input to predict the object interaction after producing masked features for all categories and outputting an object-presence probability vector for all categories. The Seg-stream first predicts the coarse segmentation masks for each category and further refines the segmentation results by convolving the segmentation masks with the object-presence probability vector obtained from the Int-stream as the filter.

Luo et al. [59] proposed a dual image segmentation (DIS) model to boost the segmentation performance using the image-level tags of the IDW dataset rather than using the object interactions and descriptions in IDW-CNN [58]. DIS first utilizes ResNet101 to produce the first feature map and the first feature vector for the segmentation prediction net and the tag classification net, respectively. The tag classification net outputs a tag prediction vector for all categories after two-stage refinement of the first feature vector. Meanwhile, in the segmentation prediction net, the second feature map is generated by calculating the sum of the upsampled first feature vector and the first feature map and is then further refined to obtain the initial segmentation map for all categories. The final segmentation prediction is obtained by refining the initial segmentation map with the tag prediction vector.

Essentially, these CNN crafting tricks optimize deep networks from the following perspectives: tailoring convolution or pooling operation in accordance with specific conditions, and modifying connection structure between different level layers. These tricks are universally applicable to all the three levels of semantic image parsing tasks.

2.1.2 Integration with the Random Field Model

Some recent studies [10, 13, 60] integrate random field models, such as Markov random fields (MRFs) [61] and CRFs [54], into deep learning to capture contextual information and long-term dependencies.

The DeepLab system [13] and DeepLabV2 [11] integrate fully connected CRFs into CNNs to refine the raw DCNN scores and achieve better segmentation results. Fully connected CRF utilizes the Gaussian CRF potentials [48] to capture long-range dependencies and treats every pixel as a CRF node to receive unary potentials. However, the CNN is separated from the CRF portions, so the DeepLab system and DeepLabV2 are not trained in an end-to-end manner.

Schwing et al. proposed fully connected deep-structured networks (FCDSs) [60] to jointly train the CNN and CRF. On the basis of the VGG16 network [39], the FCDS incorporates unary potentials into convolutional features and iteratively passes the error of CRF inference backward into the CNN. However, a CNN typically has millions of parameters while a CRF involves thousands of latent variables. Therefore, the simple integration of CNN with CRF is inefficient.

To alleviate this issue, Liu et al. proposed an end-to-end deep parsing network (DPN) [10] that incorporates high-order relations and a mixture of label contexts into an MRF and enables optimal computation of the MRF in a single forward pass rather than using an iterative algorithm. The DPN models unary terms and pairwise terms by the tailored VGG16 network [39] and additional designed layers, respectively.

2.1.3 Integration with Recurrent Neural Networks

Because the CNN [46, 47] can extract only neighboring context information through small convolutional filters, it obtains only local information, which limits the classification accuracy of each pixel position. Moreover, CRF can learn only the short-term dependencies of sequence data [54, 61] due to its own inner structure. Therefore, several works

[29, 26, 31, 32, 33, 34] used recurrent neural networks to simulate the graphical model for context modeling. Applications of recurrent neural network architecture range from 1D sequence data, such as speech and language, to 2D image space [62] and semantic segmentation.

Two-dimensional (2D) LSTM architecture [31] was adapted to consider the sophisticated spatial dependencies of labels for the pixel-level segmentation of large natural scene images. Specifically, 2D LSTM simultaneously performs classification, segmentation and context integration with low computational complexity by neglecting additional processing, such as multi-scale. Each local prediction is synchronously affected by its neighboring contexts and their previous spatial dependencies, which helps to efficiently capture local and global contextual information end-to-end.

Similarly to 2D LSTM [31], the long short-term memorized context fusion (LSTM-CF) model [29] was proposed to fuse 2D contextual information from photometric RGB and depth data. LSTM-CF can handle the challenges of severe occlusions and diverse appearances [63, 44, 64, 43, 65] for RGB-D indoor scene labeling. The photometric context is captured by stacking several convolutional layers, while the depth context is achieved by devising one LSTM layer that encodes both short-range and long-range spatial dependencies along the vertical direction. Moreover, another LSTM fusion layer is constructed to integrate the 2D contexts from different channels along the vertical direction to achieve true 2D global context through bi-directional propagation of the fused contexts along the horizontal direction. Finally, the 2D global contextual representation is cascaded with the RGB features extracted by convolutional layers.

Local-global LSTM (LG-LSTM) architecture [26] was developed for end-to-end embedding of local short-distance and global long-distance spatial context into the feature learning over all pixel positions for semantic part segmentation. The local short-distance spatial dependencies of each position in each LG-LSTM layer consist of one depth dimension and eight spatial dimensions (left side of Fig. 8). The former refers to the hidden cells from the same position in the previous LG-LSTM layer, whereas the spatial

dimensions refer to the hidden cells from eight neighborhood positions. Moreover, to capture the global long-distance spatial context (right side of Fig. 8), in each LG-LSTM layer, the whole hidden cell maps obtained from the previous layer are split into nine grids, each of which covers one part of the image. Then, the global context is obtained by max-pooling operations over each grid. Thus, the features at each position are greatly enhanced by stacking several LG-LSTM layers.

Furthermore, to improve the LG-LSTM architecture [26], the graph LSTM [32] network was built as the generalization of LSTM from sequential data to general graph-structured data. Traditional pixel-wise LSTM structures, e.g., row LSTM [66], grid LSTM [26] and diagonal BiLSTM [66, 67], take fixed-size pixels or patches as physical nodes and capture the context of each node by following a fixed route for different images. By contrast, for each image, graph LSTM constructs a single adaptive graph topology by viewing arbitrary-shaped superpixels as semantically consistent nodes, and the contextual information of each node is obtained along the edges, which represent the spatial relations of the adjacent superpixels.

Another extension of LG-LSTM [26], the structure-evolving LSTM model [33], was proposed to progressively and stochastically learn interpretable data representations over hierarchical graph structures via LSTM optimization. Structure-evolving LSTM is clearly distinguishable from graph LSTM [32], which processes only data with pre-fixed structures. Structure-evolving LSTM stochastically incorporates graph nodes with high compatibilities along the stacked LSTM layers, followed by progressive evolution of the multi-level graph representations from low levels to higher levels, which enables efficient propagation of long-range data dependencies. Moreover, the compatibility of two connected nodes accords with the corresponding LSTM gate outputs in each LSTM layer.

The third extended version of LG-LSTM [26] is progressively diffused networks (PDNs) [34], which unify multi-scale context modeling with deep feature learning for semantic image segmentation. Specifically, PDNs utilize multi-dimensional convolutional LSTMs to construct information diffusion layers,

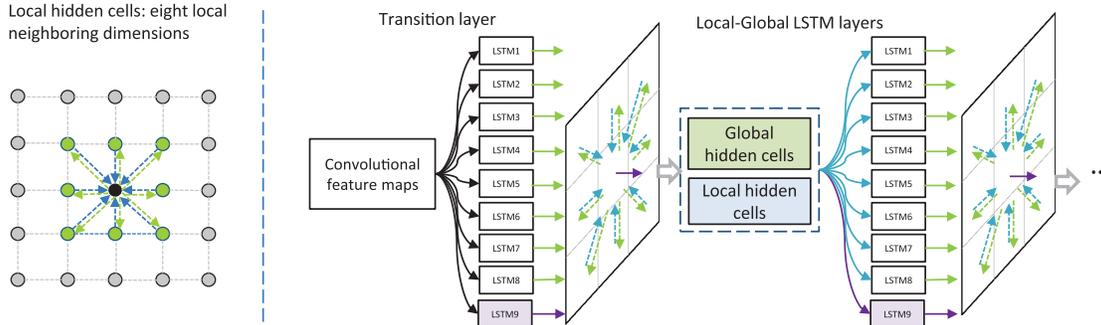


Figure 8: Illustration of the LG-LSTM layer. (Figure extracted from [26])

which contribute to diffused information over the learned feature maps. Each LSTM unit is equipped with tailored atrous filters to capture the short-range and long-range context from the neighbors of each site in the feature map.

2.2 Instance-Level Semantic Segmentation

Instance-level semantic segmentation has attracted substantial attention [68, 69, 70, 56, 71] because increasing practical applications, such as robot task planning [42] and human activity recognition [72], require different objects belonging to the same category to be distinguished. The aforementioned category-level segmentation methods cannot achieve this goal, as illustrated in Fig. 2. Instance-level semantic segmentation precisely segments each object category and correctly detects all the object instances in one image [2], seeking joint object detection and semantic segmentation. Next, we discuss this task from the three aspects: proposal-based framework, multi-task end-to-end module and metric learning embedded model.

Proposal-Based Framework. The first step in the proposal-based framework is to generate proposals, and further processing is required to produce the final segmentations. Most early deep works [68, 73, 74] first adopt a proposal generation method, extract features with tailored CNN architectures, and finally feed the intermediate results into post-

processing steps.

Typically, Hariharan et al. [68] proposed a simultaneous detection and segmentation (SDS) model. This work first generates category-agnostic candidate region proposals via bottom-up multi-scale combinatorial grouping [75] under the hypothesis that each region proposal, which consists of bounding boxes and initial segmentations, contains one object. On the basis of the proposals, features are extracted from both the bounding boxes and initial foregrounds with two separate tailored R-CNNs [64]. Support vector machines and non-maximum suppression are used to classify region proposals and to refine segmentations, respectively. This work first formulates instance-level semantic segmentation as joint object detection and semantic segmentation. But the computational cost in proposals generation phase is too expensive.

Later research [73] follows the same pipeline as that of the SDS model [68]. The differences among the methods are that 1) the refinement process in the SDS model [73] is replaced by hypercolumn-based refinement [73] to improve the segmentation accuracy; and 2) in the feature extraction step, this work [73] enlarges the bounding boxes set of detections and extracts features from just these bounding boxes without consideration of the region foreground, as in [68], which decreases the computational cost.

Recently, Li et al. [76] presented a novel salient instance segmentation approach that produces salient instance proposals by virtue of salient object contours. Specifically, This work first devised a deep

multi-scale refinement network to simultaneously detect salient region and salient object contours. Then, the salient object contours are used to generate salient object proposals, which are further filtered by subset optimization algorithm to obtain finer salient instance proposals. The final salient instance segmentation is generated by using CRF model to integrate the saliency mask with instance proposals. This work is a pioneer of joint detection of salient region and salient object contour in a unified framework, and is beneficial for the situation where multiple salient instances are spatially overlapped.

In this proposal-refinement pipeline, proposal generation precedes classification. Apparently, the deep features and large-scale training data play no role in boosting the quality of the generated proposals, therefore, the accuracy of instance segmentation is inherently limited by the quality of the initial object proposals. To resolve the issues, some newly published works [2, 17] concentrate on unifying the proposal generation and instance segmentation sub-tasks into a single end-to-end framework, and more details are discussed in the below subsection.

Multi-Task End-to-End Framework. Some methods seamlessly integrate the object segmentation of each category and the detection of all object instances into a unified framework, which is beneficial for end-to-end training without supervision in the intermediate stages.

Liang et al. [70] proposed a proposal-free network (PFN) to predict the instance numbers of different categories and each instance segmentation in end-to-end manner. This work [70] directly predicts instance-level masks through bottom-up merging, without requiring object proposals. However, PFN is not suitable for cases with small objects.

Additionally, Liang et al. proposed an alternate novel framework, called reversible recursive instance-level object segmentation (R2-IOS) [69], which recursively refines object proposals and segmentation masks. R2-IOS contains two significant sub-networks, i.e., the object proposal refinement sub-network and the instance-level object segmentation sub-network, both of which are alternately fed into each other for progressive optimization. The object proposal refinement sub-network reversibly pre-

dicts the confidences for all semantic categories and the bounding box offsets to refine the object proposals; meanwhile the instance-level object segmentation sub-network iteratively produces the foreground mask of the dominant object in each proposal. Moreover, one instance-aware denoising auto-encoder is embedded in the instance-level object segmentation sub-network, which helps R2-IOS to distinguish overlapping objects with similar appearance. This work jointly training object proposal refinement and proposal-based segmentation to complement each other, other than works in [68, 73, 76].

Dai et al. [77] presented multi-task network cascades (MNC), which dissects instance-wise segmentation into three causal sub-tasks respectively accomplished by the three sequential cascaded stages, i.e., distinguishing instances, forecasting masks and categorizing instances. Specifically, MNC first extracts the convolutional features using the stacked convolutional layers. The output is shared among the three following stages. Besides, the outputs from the early stages are also shared among the pursuant stages. This work achieves contemporary state-of-the-art accuracy by transforming complex instance-wise segmentation into three simplified sub-tasks, which, however, has its deficiencies caused by RoIPool [78, 79]: missing spatial details and repetitive computation among RoIs without sharing.

To alleviate MNC’s [77] issues, FCIS [17] provides the first fully convolutional end-to-end solution for instance-level semantic segmentation, which highly integrates FCN [12] for semantic segmentation and InstanceFCN [22] for instance mask proposal. Specifically, FCIS is divided into instance mask prediction and classification sub-task. A input image is fed into shallower convolutional layers to produce convolutional representation and further position-sensitive score maps, which are shared between subsequent two sub-networks to exploit the correlation. FCIS is fast, and preserves more spatial details without warping or resizing operations in RoIs. But FCIS has inherent drawbacks at dealing with overlapping instances [2].

Recently, a concise general framework, mask R-CNN [2], which can simultaneously detect objects in one image and generate a segmentation mask for each instance, as well as being simple to implement and

trained, was built for object instance segmentation. In concrete terms, mask R-CNN integrates one mask branch into faster R-CNN [80] so that object mask prediction is performed in parallel with the existing branch for bounding box recognition. Mask R-CNN can be generalized to other tasks, such as bounding box object detection and person keypoint detection. More importantly, mask R-CNN transcends all previous state-of-the-art results with its framework’s flexibility. However, the accuracy and speed are also restricted by the RPN and RoIPool the same as [78, 79, 80].

Metric Learning Embedded Model. Most recently, the novel research moves towards metric learning embedded deep networks for instance segmentation to measure the likelihood of different elements (e.g., pixels, detections). The distance between different elements is calculated to determine whether these elements belong to the same object instance.

Newell et al.[81] integrated associative embedding into supervised CNNs for pixel-wise predictions, which view instance segmentation as the joint detection of relevant pixels and their grouping into object instances. Here, the embeddings serve as tags to group detections with similar tags. Specifically, [81] utilizes a tailored hourglass network to simultaneously produce a detection heatmap and a grouping heatmap for each object category. The detection heatmap affords a detection score at each pixel to predict whether the pixel belongs to the foreground. Meanwhile, the grouping heatmap tags each pixel such that pixels with similar tags are grouped into the same object instance by non-maximum suppression. Besides pixel-wise embeddings [82], this work also engenders pixel-wise detection scores to reduce the output dimension of each pixel.

Coincidentally, Fathi et al. [83] also manufactured a deep metric learning method to further improve the performance of instance-wise segmentation. Specifically, a fully convolutional scoring model is first adopted to compute the seediness score of each pixel, which estimates the representativeness of the pixel comparing with other pixels in the same instance. Pixels with top seediness score serve as seed points. Then, the distance between the seed points are learned via a deep embedding model, which rep-

resents likelihood of two pixels. Thus similar pixels are grouped together into the same instance. Different from [81] using one-dimensional embedding, this work derives multi-dimensional embedding from each pixels, which makes it more appropriate for slender-shape objects.

Generally, these metric learning embedded models are trained end-to-end with fast speed and promising performance. The grouping procedure is based on pairwise constraints [84], not associated with pre-defined semantic categories. Therefore, such embedding technology maybe become a new tendency for instance-level segmentation.

2.3 Beyond Segmentation

The aforementioned segmentation research focuses on segmenting images with different-level configurations, such as category level and instance level. Each configuration assigns the label of the corresponding level for each pixel. In this section, we discuss beyond segmentation methods, which considers the implicit high-level hierarchical information in the image, such as the geometric information [30], the relations between objects [58], and the structural information [6], in addition to the aforementioned pixel-wise segmentation. This high-level information improves the image segmentation performance.

Peng et al. [30] proposed hierarchical LSTM (H-LSTM) to exploit data from the perspective of geometric attributes and geometric relations, as shown in Fig. 9. Specifically, H-LSTM simultaneously outputs the segmentation of geometric attributes (e.g., sky, ground) and geometric interaction relations (e.g., layering, supporting) through the pixel LSTM (P-LSTM) sub-network and the multi-scale super-pixel LSTM (MS-LSTM) sub-network, respectively. P-LSTM captures local contextual information to segment geometric attributes; meanwhile, MS-LSTM extracts multi-scale super-pixel representations to categorize geometric interaction relations between adjacent attributes. MS-LSTM shares basic convolutional layers with P-LSTM, which means attribute segmentation and relation categorization benefit from each other.

The major obstacles in beyond segmentation re-

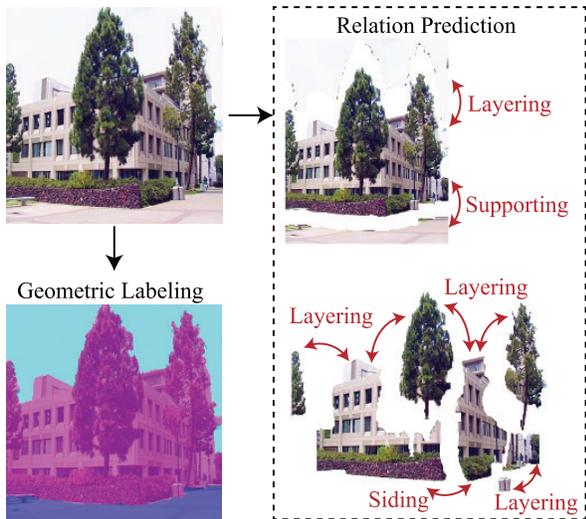


Figure 9: Illustration of the geometric scene parsing. (Figure extracted from [30])

search are the ambiguity of the image hierarchical representations and the rarity of elaborative manually annotated datasets. To alleviate these issues, some [6, 58] introduced top-down information (e.g., hierarchical object structure, object interactions) from image descriptions.

Lin et al. [6] proposed the deep-structured CNN-RNN model by integrating a CNN [12, 15] and RNN [16], which can recursively learn the representations in a semantically and structurally coherent way, as shown in Fig. 10. The CNN layer-wise extracts the feature maps of semantic objects from the input scene image (i.e., semantic segmentation results). Then, the feature maps are fed into the RNN to generate the hierarchically structured configuration (i.e., the hierarchical object structure and the object interaction relations), as shown in Fig. 3. The CNN-RNN model [6] discovers structural scene configurations from the image descriptions [27, 85] following the work of [86, 87] and is trained in a weakly supervised manner, which avoids the need for elaborate manual annotations. Furthermore, the expectation-maximization method, which alternates between latent label prediction subject to the weak annotation

constraints and optimization of the network parameter, is used to train the model.

Inspired by CNN-RNN [6], IDW-CNN [58] also exploits the image descriptions [27, 85] to capture top-down information, which further improves the image parsing performance. Wang et al. [58] designed an elaborate CNN to jointly train IDW and a subsistent image segmentation dataset. IDW dataset are raw: 1) Images are only captioned with raw sentences without pixel-wise annotation; and 2) Images and their descriptions are automatically downloaded from the Internet without subsequent manual post-processing (e.g., cleaning, refinement). IDW contributes useful object interactions to improve segmentation performance; consequently, the precise object segmentation results from the subsistent dataset also benefit object-interaction extractions in IDW. Thus, knowledge from different dataset sources can be fully explored and transferred to improve performance.

3 Datasets and Evaluation Metrics

Public datasets and relevant evaluation metrics form the foundation for improving the algorithms. The emergence of big data has driven the development of datasets and relevant evaluation metrics in the field of deep representation for semantic image segmentation, which require large-scale datasets for training. Thus, in this section, we describe these well-known public datasets and evaluation metrics in detail.

3.1 Datasets

Table 2 compares the well-known public datasets. According to the source of the image, the semantic segmentation datasets can be divided into RGB(2D) data and RGB-D(3D) data, as indicated by the "2D / 3D" term.

For category-level / instance-level semantic segmentation, the widely used datasets include PASCAL VOC, PASCAL-Part, ILSVRC 2016, MS COCO, SIFT Flow, NYUDv2, SUN RGB-D, ATR, and Fashionista.

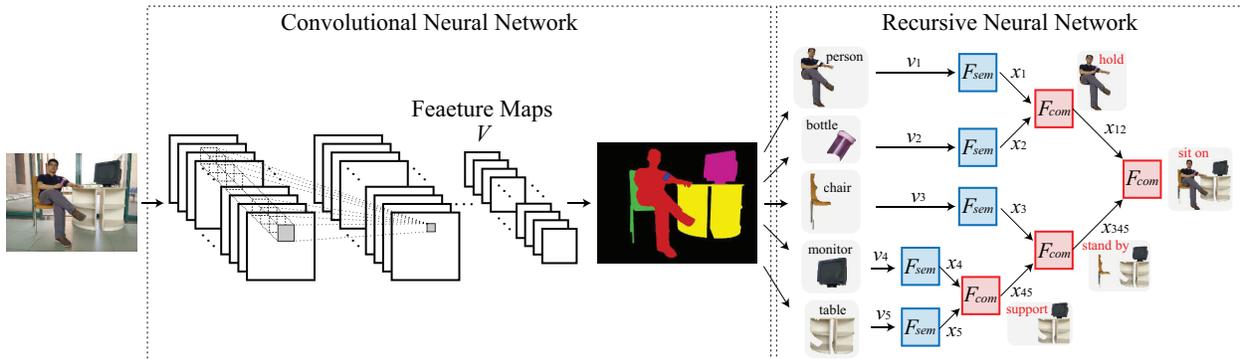


Figure 10: The detailed CNN-RNN architecture. (Figure extracted from [6])

Datasets	# of images	# of training images	# of testing images	2D/3D image	indoor/outdoor scene	# of categories
PASCAL VOC 2012	9,993	4,997	4,996	2D	both	20
PASCAL VOC 2007	7,062	3,531	3,531	2D	both	20
PASCAL Part	19,740	10,103	9,637	2D	both	88
ILSVRC 2016	25,562	20,210	3,352	2D	both	150
MS COCO	328,000	—	—	2D	both	91
SIFT Flow	2,688	—	—	2D	outdoor	33
NYUDv2	1,449	—	—	3D	indoor	40
SUN RGB-D	10,355	—	—	3D	indoor	19
Fashionista	685	456	229	2D	both	56
ATR	7,700	6,000	1,000	2D	both	18
CityScapes	5,000	3,475	1,525	2D	outdoor	30

Table 2: Comparison of the semantic segmentation datasets. “#” is short for “number”. “—” means the value cannot be found in the literature.

- PASCAL VOC.** The PASCAL VOC dataset [88] is part of the PASCAL Visual Object Classes (VOC) Challenge organized annually from 2005 to 2012. VOC data have been accepted annually for five main tasks: classification, detection, segmentation, action classification and large-scale recognition. The segmentation task was first introduced in 2007. The dataset is utilized for both category-level and instance-level segmentation. Table 2 lists VOC 2007 and VOC 2012, which are the most frequently used VOC datasets.
- PASCAL-Part Dataset.** This dataset [89] contains additional annotations for PASCAL VOC 2010, which provides segmentation masks

for each body part of an object.

- ILSVRC 2016.** The ImageNet Large Scale Visual Recognition Challenge 2016 (ILSVRC 2016) [90], organized by the MIT CSAIL Vision Group, is well-known for the image classification task, and it first introduced a scene parsing task in 2016. The dataset for this scene parsing task is the complete ADE20K Dataset [91], which contains more than 20K scene-centric images exhaustively annotated with object instances and object parts. Thus, the dataset is used for both semantic instance-level segmentation and category-level segmentation. In particular, the distribution of objects occurring in the images is non-uniform, which simulates daily real-world

scenes.

- **MS COCO.** The Microsoft Common Objects in COntext (MS COCO) dataset contains 91 common object categories in the version released in 2015 [92] and 80 categories in the 2014 version [93]. Distinct from others datasets, MS COCO contains considerably more object instances per image, which may help to exploit contextual information. MS COCO is now a widely used benchmark dataset for category-level and instance-level semantic segmentation.
- **SIFT Flow.** The SIFT Flow dataset [94] was thoroughly labeled by LabelMe users with 33 semantic categories, 3 geometric categories (i.e., ground, vertical, and sky) and 4 interaction relation labels (i.e., layering, supporting, siding and affinity). The dataset is appropriate for category-level segmentation, and it was later transformed for image geometric parsing.
- **NYUDv2.** NYUDv2 [95] is an RGB-D dataset [96] and can be used for both category-level and instance-level segmentation. Additionally, it contains labeled structural support relationships for support relation classification.
- **SUN RGB-D.** SUN RGB-D [97] is the largest RGB-D dataset currently available. The dataset combines most of the previous datasets, such as NYUDv2 [95], Berkeley B3DO [98], and SUN3D [99], as well as 3943 newly captured RGB-D images [97]. Currently, the SUN RGB-D dataset is designed for only category-level semantic segmentation.
- **Fashionista.** The Fashionista dataset [100], collected from chictopia.com, is designed for clothes parsing and contains 56 different clothing items. Thus, the dataset is tailored for category-level segmentation.
- **ATR.** The ATR dataset [46], also used for category-level segmentation, combines four human parsing datasets: Fashionista [100], Colorful Fashion Parsing Data (CFPD) [101], Daily Photos [102] and the Human Parsing in the Wild

(HPW) datasets. The labels of the Fashionista and CFPD datasets are merged into 18 categories, and the HPW dataset is newly annotated [46].

- **CityScapes.** The CityScapes dataset [103] focuses on both category-level and instance-level segmentation of urban street scenes. It provides 5,000 fine annotations, i.e., individual annotations of single instances, and 20,000 coarse annotations, which cover individual objects with marked polygons.

Semantic image parsing mainly contains structured semantic parsing [6] and geometric parsing [30]. However, to the best of our knowledge, there is no specific subsistent dataset for image parsing. Structured semantic parsing requires not only the segmentation of objects in an image but also hierarchical prediction of semantic objects with object interaction relations. Therefore, the requisite dataset is distinct from the datasets used for semantic segmentation tasks. The work in [6] constructed a dataset for this task on the basis of the existing dataset PASCAL VOC 2012. In practice, in addition to utilizing the dataset for category-level semantic segmentation, images, used for constructing the structure and relations, were selected from the PASCAL VOC 2012 segmentation dataset. Furthermore, in contrast to the pixel-wise annotations in this dataset, based on these selected images, [6] built image-level annotations by describing each image with several natural language sentences. In addition, each sentence contains objects and the hierarchy with their interaction relations in the image. Similarly, there is no specific dataset to validate the effectiveness of the proposed algorithm for geometric parsing, which simultaneously labels geometric attributes and determines the geometric interaction relations. The work in [30] transformed existing datasets (i.e., SIFT Flow, LM+SUN, and Geometric Context dataset) for use in geometric parsing.

3.2 Evaluation Metrics

The performance of pixel-wise segmentation algorithms is commonly evaluated with four metrics [12]:

pixel-wise accuracy, mean accuracy, intersection over union (IoU), and F1 score. Denote n_{ij} as the number of pixels of category i predicted to belong to category j , where there are K categories, and let $t_i = \sum_j n_{ij}$ be the total number of pixels of category i . Then,

- pixel-wise accuracy: $\sum_i n_{ii} / \sum_i t_i$
- mean accuracy: $(1/K) \sum_i n_{ii} / \sum_i t_i$
- mean IoU: $(1/K) \sum_i n_{ii} / (t_i + \sum_j n_{ji} - n_{ii})$
- F1 score: (Pixel-wise Accuracy + mean IoU) / 2

The performance of structured scene parsing algorithms is evaluated with two metrics [6]: relation accuracy and structure accuracy. Following [6], the structured scene parsing task is defined as a binary tree, and relation accuracy is computed recursively. Denote the binary tree by T and let $P = T, T_1, T_2, \dots, T_m$ be the set of enumerated subtrees (including T) of T . Each tree consists of objects and relations between different objects, while each leaf only stands for one object. A leaf T_i is considered to be correct if it is of the same object category as that in the ground truth tree. A non-leaf T_i (with two subtrees T_l and T_r) is considered to be correct if and only if object categories and relation labels in T_l and T_r are both correctly predicted. Then, the relation accuracy is calculated as (*number of correct subtrees*)/($m+1$). The structure accuracy is a simplification of the relation accuracy that ignores the relation labels when evaluating the correctness of T .

4 Conclusions and Future Work

In this work, we present a comprehensive review on deep representation learning algorithms for semantic image parsing with a unique perspective. In contrast to other surveys, we review the image parsing models in terms of the development of three-level semantic segmentation from its origins to the most recent, the relatively well-known datasets, and evaluation metrics, including 41 algorithms, 11 datasets

and 6 evaluation metrics. We believe that there are several promising research directions for semantic image parsing. The first is multi-task driven semantic parsing, such as [6], which integrates natural language understanding and image parsing. In addition, a large number of training samples are required for deep parsing models, but the collection and annotation of large-scale datasets is elaborative. Therefore, semi-supervised, weakly supervised or unsupervised learning algorithms are another direction to pursue. The third intuitive direction is to transfer image parsing ideas and technologies to the challenging video parsing task.

References

- [1] Zhao H, Shi J, Qi X, Wang X, Jia J. Pyramid scene parsing network. ArXiv Preprint ArXiv:161201105, 2016
- [2] He K, Gkioxari G, Dollár P, Girshick R. Mask r-cnn. ArXiv Preprint ArXiv:170306870, 2017
- [3] Tu Z, Chen X, et al. Image parsing: unifying segmentation, detection, and recognition. In: Proceedings of the IEEE International Conference on Computer Vision. 2003, 18–25
- [4] Tu Z, Zhu S C. Parsing images into region and curve processes. In: Proceedings of the European Conference on Computer Vision. 2002, 201–310
- [5] Han F, Zhu S C. Bottom-up/top-down image parsing by attribute graph grammar. In: Proceedings of the IEEE International Conference on Computer Vision. 2005, 1778–1785
- [6] Lin L, Wang G, Zhang R, Zhang R, Liang X, Zuo W. Deep structured scene parsing by learning with image descriptions. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. 2016,, 2276–2284
- [7] Lowe D G. Distinctive image features from scale-invariant keypoints. International Journal of Computer Vision, 2004. 60(2):91–110

- [8] Dalal N, Triggs B. Histograms of oriented gradients for human detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2005, 886–893
- [9] Ahonen T, Hadid A, Pietikäinen M. Face recognition with local binary patterns. In: Proceedings of the European Conference on Computer Vision. 2004, 469–481
- [10] Liu Z, Li X, Luo P, Loy C C, Tang X. Semantic image segmentation via deep parsing network. In: Proceedings of the IEEE International Conference on Computer Vision. 2015, 1377–1385
- [11] Chen L C, Papandreou G, Kokkinos I, Murphy K, Yuille A L. Deeplab: semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. ArXiv Preprint ArXiv:160600915, 2016
- [12] Long J, Shelhamer E, Darrell T. Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 3431–3440
- [13] Chen L C, Papandreou G, Kokkinos I, Murphy K, Yuille A L. Semantic image segmentation with deep convolutional nets and fully connected crfs. ArXiv Preprint ArXiv:14127062, 2014
- [14] Peng C, Zhang X, Yu G, Luo G, Sun J. Large kernel matters-improve semantic segmentation by global convolutional network. ArXiv Preprint ArXiv:170302719, 2017
- [15] Krizhevsky A, Sutskever I, Hinton G E. Imagenet classification with deep convolutional neural networks. In: Proceedings of Advances in Neural Information Processing Systems. 2012, 1097–1105
- [16] Socher R, Manning C D, Ng A Y. Learning continuous phrase representations and syntactic parsing with recursive neural networks. In: Proceedings of the NIPS-2010 Deep Learning and Unsupervised Feature Learning Workshop. 2010, 1–9
- [17] Li Y, Qi H, Dai J, Ji X, Wei Y. Fully convolutional instance-aware semantic segmentation. ArXiv Preprint ArXiv:161107709, 2016
- [18] Bengio Y. Deep learning of representations: looking forward. In: Proceedings of International Conference on Statistical Language and Speech Processing. 2013, 1–37
- [19] Bengio Y. Deep learning of representations for unsupervised and transfer learning. In: Proceedings of the International Conference on Machine Learning Workshop on Unsupervised and Transfer Learning. 2012, 17–36
- [20] Bengio Y, Courville A, Vincent P. Representation learning: a review and new perspectives. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013. 35(8):1798–1828
- [21] LeCun Y, Boser B, Denker J S, Henderson D, Howard R E, Hubbard W, Jackel L D. Backpropagation applied to handwritten zip code recognition. Neural Computation, 1989. 1(4):541–551
- [22] Dai J, He K, Li Y, Ren S, Sun J. Instance-sensitive fully convolutional networks. In: Proceedings of the European Conference on Computer Vision. 2016, 534–549
- [23] Islam M A, Naha S, Roohan M, Bruce N, Wang Y. Label refinement network for coarse-to-fine semantic segmentation. ArXiv Preprint ArXiv:170300551, 2017
- [24] Lipton Z C, Berkowitz J, Elkan C. A critical review of recurrent neural networks for sequence learning. Computer Science, 2015
- [25] LeCun Y, Bengio Y, Hinton G. Deep learning. Nature, 2015. 521(7553):436–444
- [26] Liang X, Shen X, Xiang D, Feng J, Lin L, Yan S. Semantic object parsing with local-global long short-term memory. In: Proceedings of

- the IEEE Conference on Computer Vision and Pattern Recognition. 2016, 3185–3193
- [27] Karpathy A, Fei-Fei L. Deep visual-semantic alignments for generating image descriptions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 3128–3137
- [28] Sutskever I, Vinyals O, Le Q V. Sequence to sequence learning with neural networks. In: Proceedings of Advances in Neural Information Processing Systems. 2014, 3104–3112
- [29] Li Z, Gan Y, Liang X, Yu Y, Cheng H, Lin L. Lstm-cf: unifying context modeling and fusion with lstms for rgb-d scene labeling. In: Proceedings of the European Conference on Computer Vision. 2016, 541–557
- [30] Peng Z, Zhang R, Liang X, Liu X, Lin L. Geometric scene parsing with hierarchical lstm. ArXiv Preprint ArXiv:160401931, 2016
- [31] Byeon W, Breuel T M, Raue F, Liwicki M. Scene labeling with lstm recurrent neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 3547–3555
- [32] Liang X, Shen X, Feng J, Lin L, Yan S. Semantic object parsing with graph lstm. In: Proceedings of the European Conference on Computer Vision. 2016, 125–143
- [33] Liang X, Lin L, Shen X, Feng J, Yan S, Xing E P. Interpretable structure-evolving lstm. ArXiv Preprint ArXiv:170303055, 2017
- [34] Zhang R, Yang W, Peng Z, Wang X, Lin L. Progressively diffused networks for semantic image segmentation. ArXiv Preprint ArXiv:170205839, 2017
- [35] Elman J L. Distributed representations, simple recurrent networks, and grammatical structure. *Machine Learning*, 1991. 7(2-3):195–225
- [36] Liu W, Rabinovich A, Berg A C. Parsenet: looking wider to see better. ArXiv Preprint ArXiv:150604579, 2015
- [37] Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A. Going deeper with convolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 1–9
- [38] He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016, 770–778
- [39] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. ArXiv Preprint ArXiv:14091556, 2014
- [40] Pinheiro P, Collobert R. Recurrent convolutional neural networks for scene labeling. In: Proceedings of International Conference on Machine Learning. 2014, 82–90
- [41] Graves A, Fernández S, Schmidhuber J. Multi-dimensional recurrent neural networks. In: Proceedings of the International Conference on Artificial Neural Networks. 2007, 549–558
- [42] Lin L, Huang L, Chen T, Gan Y, Cheng H. Knowledge-guided recurrent neural network learning for task-oriented action prediction. In: Proceedings of the IEEE International Conference on Multimedia and Expo. 2017, 625–630
- [43] Farabet C, Couprie C, Najman L, LeCun Y. Learning hierarchical features for scene labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013. 35(8):1915–1929
- [44] Gupta S, Girshick R, Arbeláez P, Malik J. Learning rich features from rgb-d images for object detection and segmentation. In: Proceedings of the European Conference on Computer Vision. 2014, 345–360

- [45] Ning F, Delhomme D, LeCun Y, Piano F, Botou L, Barbano P E. Toward automatic phenotyping of developing embryos from videos. *IEEE Transactions on Image Processing*, 2005. 14(9):1360–1371
- [46] Liang X, Liu S, Shen X, Yang J, Liu L, Dong J, Lin L, Yan S. Deep human parsing with active template regression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015. 37(12):2402–2414
- [47] Liang X, Xu C, Shen X, Yang J, Liu S, Tang J, Lin L, Yan S. Human parsing with contextualized convolutional neural network. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, 1386–1394
- [48] Krähenbühl P, Koltun V. Efficient inference in fully connected crfs with gaussian edge potentials. In: *Proceedings of Advances in Neural Information Processing Systems*. 2011, 109–117
- [49] Noh H, Hong S, Han B. Learning deconvolution network for semantic segmentation. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, 1520–1528
- [50] Badrinarayanan V, Handa A, Cipolla R. Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling. *ArXiv Preprint ArXiv:150507293*, 2015
- [51] Ronneberger O, Fischer P, Brox T. U-net: Convolutional networks for biomedical image segmentation. In: *Proceedings of International Conference on Medical Image Computing and Computer-Assisted Intervention*. 2015, 234–241
- [52] Lin G, Milan A, Shen C, Reid I. Refinenet: Multi-path refinement networks with identity mappings for high-resolution semantic segmentation. *ArXiv Preprint ArXiv:161106612*, 2016
- [53] Chen L C, Papandreou G, Schroff F, Adam H. Rethinking atrous convolution for semantic image segmentation. *ArXiv Preprint ArXiv:170605587*, 2017
- [54] Yu F, Koltun V. Multi-scale context aggregation by dilated convolutions. *ArXiv Preprint ArXiv:151107122*, 2015
- [55] Li X, Liu Z, Luo P, Loy C C, Tang X. Not all pixels are equal: difficulty-aware semantic segmentation via deep layer cascade. *ArXiv Preprint ArXiv:170401344*, 2017
- [56] Zhou Y, Xie L, Shen W, Wang Y, Fishman E K, Yuille A L. A fixed-point model for pancreas segmentation in abdominal ct scans. In: *Proceedings of International Conference on Medical Image Computing and Computer-Assisted Intervention*. 2017, 693–701
- [57] Li Q, Wang J, Tu Z, Wipf D P. Fixed-point model for structured labeling. In: *Proceedings of the International Conference on Machine Learning*. 2013, 214–221
- [58] Wang G, Luo P, Lin L, Wang X. Learning object interactions and descriptions for semantic image segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, 5859–5867
- [59] Luo P, Wang G, Lin L, Wang X. Deep dual learning for semantic image segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, 2718–2726
- [60] Schwing A G, Urtasun R. Fully connected deep structured networks. *ArXiv Preprint ArXiv:150302351*, 2015
- [61] Yang W, Luo P, Lin L. Clothing co-parsing by joint image segmentation and labeling. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2014, 3182–3189
- [62] Byeon W, Liwicki M, Breuel T M. Texture classification using 2d lstm networks. In: *Proceedings of International Conference on Pattern Recognition*. 2014, 1144–1149

- [63] Eigen D, Fergus R. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: Proceedings of the IEEE International Conference on Computer Vision. 2015, 2650–2658
- [64] Girshick R, Donahue J, Darrell T, Malik J. Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2014, 580–587
- [65] Reza M, Kosecka J, et al. Reinforcement learning for semantic segmentation in indoor scenes. ArXiv Preprint ArXiv:160601178, 2016
- [66] Oord A V d, Kalchbrenner N, Kavukcuoglu K. Pixel recurrent neural networks. In: Proceedings of International Conference on Machine Learning. 2016, 1747–1756
- [67] Kalchbrenner N, Danihelka I, Graves A. Grid long short-term memory. ArXiv Preprint ArXiv:150701526, 2015
- [68] Hariharan B, Arbeláez P, Girshick R, Malik J. Simultaneous detection and segmentation. In: Proceedings of the European Conference on Computer Vision. 2014, 297–312
- [69] Liang X, Wei Y, Shen X, Jie Z, Feng J, Lin L, Yan S. Reversible recursive instance-level object segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016, 633–641
- [70] Liang X, Wei Y, Shen X, Yang J, Lin L, Yan S. Proposal-free network for instance-level object segmentation. ArXiv Preprint ArXiv:150902636, 2015
- [71] Abtahi F, Zhu Z, Burry A M. A deep reinforcement learning approach to character segmentation of license plate images. In: Proceedings of International Conference on Machine Vision Applications, 2015 14th IAPR. 2015, 539–542
- [72] Lin L, Wang K, Zuo W, Wang M, Luo J, Zhang L. A deep structured model with radius-margin bound for 3d human activity recognition. International Journal of Computer Vision, 2016. 118(2):256–273
- [73] Hariharan B, Arbeláez P, Girshick R, Malik J. Hypercolumns for object segmentation and fine-grained localization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 447–456
- [74] Chen Y T, Liu X, Yang M H. Multi-instance object segmentation with occlusion handling. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 3470–3478
- [75] Arbeláez P, Pont-Tuset J, Barron J T, Marques F, Malik J. Multiscale combinatorial grouping. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition. 2014, 328–335
- [76] Li G, Xie Y, Lin L, Yu Y. Instance-level salient object segmentation. ArXiv Preprint ArXiv:170403604, 2017
- [77] Dai J, He K, Sun J. Instance-aware semantic segmentation via multi-task network cascades. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016, 3150–3158
- [78] Girshick R. Fast r-cnn. In: Proceedings of the IEEE International Conference on Computer Vision. 2015, 1440–1448
- [79] He K, Zhang X, Ren S, Sun J. Spatial pyramid pooling in deep convolutional networks for visual recognition. In: Proceedings of the European Conference on Computer Vision. 2014, 346–361
- [80] Ren S, He K, Girshick R, Sun J. Faster r-cnn: towards real-time object detection with region proposal networks. In: Proceedings of Advances in Neural Information Processing Systems. 2015, 91–99

- [81] Newell A, Huang Z, Deng J. Associative embedding: end-to-end learning for joint detection and grouping. In: Proceedings of Advances in Neural Information Processing Systems. 2017, 2274–2284
- [82] Harley A W, Derpanis K G, Kokkinos I. Learning dense convolutional embeddings for semantic segmentation. ArXiv Preprint ArXiv:151104377, 2015
- [83] Fathi A, Wojna Z, Rathod V, Wang P, Song H O, Guadarrama S, Murphy K P. Semantic instance segmentation via deep metric learning. ArXiv Preprint ArXiv:170310277, 2017
- [84] Yang L, Jin R. Distance metric learning: A comprehensive survey. Michigan State University, 2006. 2(2)
- [85] Xu J, Schwing A G, Urtasun R. Tell me what you see and i will show you where it is. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2014, 3190–3197
- [86] Miller G A, Beckwith R, Fellbaum C, Gross D, Miller K J. Introduction to wordnet: An online lexical database. International Journal of Lexicography, 1990. 3(4):235–244
- [87] Socher R, Bauer J, Manning C D, Ng A Y. Parsing with compositional vector grammars. In: Proceedings of Annual Meeting of the Association for Computational Linguistics. 2013, 455–465
- [88] Everingham M, Van Gool L, Williams C K, Winn J, Zisserman A. The pascal visual object classes (voc) challenge. International Journal of Computer Vision, 2010. 88(2):303–338
- [89] Chen X, Mottaghi R, Liu X, Fidler S, Urtasun R, Yuille A. Detect what you can: detecting and representing objects using holistic models and body parts. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2014, 1971–1978
- [90] Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, et al. Imagenet large scale visual recognition challenge. International Journal of Computer Vision, 2015. 115(3):211–252
- [91] Zhou B, Zhao H, Puig X, Fidler S, Barriuso A, Torralba A. Semantic understanding of scenes through the ade20k dataset. ArXiv Preprint ArXiv:160805442, 2016
- [92] Lin T Y, Maire M, Belongie S, Bourdev L, Girshick R, Hays J, Perona P, Ramanan D, Zitnick C L, Dollár P. Microsoft coco: common objects in context. ArXiv Preprint ArXiv:14050312v3, 2015
- [93] Lin T Y, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick C L. Microsoft coco: Common objects in context. In: Proceedings of the European Conference on Computer Vision. 2014, 740–755
- [94] Liu C, Yuen J, Torralba A. Nonparametric scene parsing: Label transfer via dense scene alignment. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2009, 1972–1979
- [95] Silberman N, Hoiem D, Kohli P, Fergus R. Indoor segmentation and support inference from rgb-d images. In: Proceedings of the European Conference on Computer Vision. 2012, 746–760
- [96] Gupta S, Arbelaez P, Malik J. Perceptual organization and recognition of indoor scenes from rgb-d images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2013, 564–571
- [97] Song S, Lichtenberg S P, Xiao J. Sun rgb-d: A rgb-d scene understanding benchmark suite. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015, 567–576
- [98] Janoch A, Karayev S, Jia Y, Barron J T, Fritz M, Saenko K, Darrell T. A category-level 3d

- object dataset: putting the kinect to work. In: Proceedings of IEEE International Conference on Computer Vision Workshop on Consumer Depth Cameras in Computer Vision. 2011, 1168–1174
- [99] Xiao J, Owens A, Torralba A. Sun3d: A database of big spaces reconstructed using sfm and object labels. In: Proceedings of the IEEE International Conference on Computer Vision. 2013, 1625–1632
- [100] Yamaguchi K, Kiapour M H, Ortiz L E, Berg T L. Parsing clothing in fashion photographs. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2012, 3570–3577
- [101] Liu S, Feng J, Domokos C, Xu H, Huang J, Hu Z, Yan S. Fashion parsing with weak color-category labels. *IEEE Transactions on Multimedia*, 2014. 16(1):253–265
- [102] Dong J, Chen Q, Xia W, Huang Z, Yan S. A deformable mixture parsing model with parselets. In: Proceedings of the IEEE International Conference on Computer Vision. 2013, 3408–3415
- [103] Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, Franke U, Roth S, Schiele B. The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016, 3213–3223