

# Semantic Segmentation Annotation: Common Misconceptions and Problems

# Table of Content

- Introduction
- Type of Semantic Segmentation
- Misconceptions, Problems
- Datasets and Challenges
- Use Cases
- Overcoming Challenges with Solution
- Conclusion

# Introduction

In computer vision, semantic segmentation is one of the most important components for fine-grained inference (CV). To achieve the appropriate precision levels, models must grasp the context of the environment in which they operate. As a result, through pixel accuracy, semantic segmentation supplies them with the insight.

Before we dig deep into the topic, let us understand what is Semantic Segmentation.

The goal of Semantic Segmentation is to group pixels in a meaningful way. Pixels that belong to a road, people, automobiles, or trees, for example, must be grouped individually. As a result, semantic segmentation does pixel-by-pixel categorization, such as determining if a pixel is part of a traversable road, an automobile, or a pedestrian. For self-driving automobiles and robotic navigation systems, this is critical.

Although semantic segmentation is described as the process of identifying and labeling images at the pixel level, it is sometimes mistaken for instance segmentation. The major difference is that with semantic segmentation, all pixels that belong to the same class have the same pixel value.

Another method called panoptic segmentation is used to label pictures containing items of a single class to make the object recognition process easier and more accurate. Instance and semantic segmentation are combined in this method.

## Types of Semantic Segmentation

# Region-Based Semantic Segmentation

Region-based semantic segmentation is mainly used for segmentation that incorporates region extraction and semantic-based classification. In this type of segmentation, first of all, only free-form regions are selected by the model and then these regions are transformed into predictions at a pixel level to make sure each pixel is visible to computer vision.

And this runs through the CNN, dragging features from every one of these different areas. In the end, every region is classified using a linear support vector machine specific to the chosen classes in the same class providing detail information about the subject.

The R-CNN extracts two different feature types for every region picked by the model. A frontal feature and a full region are selected. And when these two region features are joined together, resulting in the performance of the model getting improved with such segmentation.

Whereas, R-CNN models manage to utilize the discriminative CNN features and achieve improved classification performance, however, they are also limited when it comes to generating precise boundaries around the object affecting the precision.

*Actually, a specific type of framework is used to complete this in the regions through the CNN framework, or R-CNN, that uses a specific search algorithm to drag many possible section proposals from an image*

# Fully Convolutional Network-Based Semantic Segmentation

CNNs are mainly used for computer vision to perform tasks like image classification, face recognition, identifying and classifying everyday objects, and image processing in robots and autonomous vehicles. It is also used for video analysis and classification, semantic parsing, automatic caption generation, search query retrieval, sentence classification, and much more.

A Fully Convolutional Network functions are created through a map that transforms the pixels to pixels. Fully convolutional neural networks can be used to create labels for inputs for pre-defined sizes that happen as a result of fully connected layers being fixed in their inputs.

While FCNs can understand randomly sized images, and they work by running the inputs through alternating convolution and pooling layers, and often times the final result of the FCN is it predicts that are low in resolution resulting in relatively ambiguous object boundaries.

# Weakly Supervised Semantic Segmentation

This is one of the most commonly used semantic segmentation models that create a large number of images with each segment pixel-wise. Hence, manually annotating of each of the masks is not only very time consuming but also an expensive process.

Therefore, some weakly supervised methods have been proposed recently, that are dedicated to achieving the semantic segmentation by utilizing annotated bounding boxes. However, there are different methods for using bounding boxes for supervised training of the network and make the iterative improvements to the estimated positioning of the masks.

Actually, there are different methods for using bounding boxes. This technique uses the bounding boxes to supervise the training of the network and make iterative improvements to the estimated positioning of the masks. Depending on the bounding box data labeling tool the object is annotated while eliminating the noise and focusing the object with accuracy.

So, the most commonly used method for semantic segmentation is used as an FCN, as it can be also implemented by taking a pre-trained network and with the flexibility to customize the various aspects as per the network fitting in your project requirements.

Issues with the sample, errors during the collection of data, problems with the documentation, or having duplicate records would all impact data quality.

# Misconceptions and Problems

Semantic segmentation is a computer vision problem that entails putting related elements of an image into the same class. There are several problems while doing semantic segmentation. Some of them are listed below:

## 1. It's difficult and time-consuming to annotate by hand

Making semantic masks by hand is a time-consuming and difficult task. When confronted with irregular forms or locations where the boundary between items is not immediately discernible, the labeller must accurately follow the outlines of each object (see pictures below). Annotating a single frame without specialized tools is prone to mistakes, inconsistencies, and can take more than 30 minutes.

## 2. Dependency on fully automated methods

Fully automated methods are incapable of delivering high-quality results

Wouldn't it be great if we could just train a neural network to do semantic segmentation once and then have all of our annotations without having to do anything?

The reason for this is a misalignment between our perceptions of quality and how accuracy is assessed. The contour of an item is used to generate a segmentation mask, and the quality is determined by the percentage of the region that was properly detected.

## 3. It takes a long time to fix mistakes

Mistakes may be expensive in each of the aforementioned ways. Correcting an imperfect segmentation mask necessitates the correction of N additional masks, where N is the number of neighbouring masks (we'll return to this later). It takes as long to adjust the mask as it does to create it from start. As a result, human adjustment of a completely automated segmentation's output is likewise not possible. The only method to prevent this issue is to use specialized annotation software and labellers who are adequately trained.

## 4. Semantic segmentation annotation costs

As you may have seen, segmentation mask creation necessitates the use of specific annotators, equipment, and automation. This raises the price dramatically, frequently by several folds above the cost of annotating basic bounding boxes, and quickly depletes the budget.



# Datasets and Challenges



Data is perhaps one of the most crucial – if not the most critical – components of any machine learning system. When dealing with deep networks, the necessity of this is amplified. As a result, accumulating sufficient training data into a dataset is important for any deep learning-based segmentation system.

When it comes to datasets, the high quality and accurate training of these datasets are crucial. Analytics.ai is a data labeling company that provides low-cost data annotation services which can enhance the overall AI and ML models. Apart from this, Cogito Tech LLC can also assist you in the process.

Time, domain expertise to select relevant information, and infrastructure to capture and transform that data into a representation that the system can properly understand and learn are all required for gathering and constructing an appropriate dataset, which must have a large enough scale and accurately represent the system's use case. Despite its simplicity in formulation compared to advanced neural network design descriptions, this challenge is one of the most difficult to accomplish in this context.

This approach has another benefit for the community: standardized datasets allow for fair comparisons between systems; in fact, many datasets are part of a challenge that reserves some data – not provided to developers to test their algorithms – for a competition in which many methods are tested, resulting in a fair ranking of methods based on their actual performance without any data cherry-picking.

## Use Cases

Semantic segmentation for computer vision is used in a variety of fields, including:

- Recognizing people by their faces
- Recognition of handwriting
- Image search in the virtual world
- Automobiles that drive themselves
- Mapping for satellite and aerial imagery for the fashion industry and virtual try-on
- Imaging and diagnostics in medicine

In general, semantic segmentation is utilized for more complex tasks than other image annotation methods, since it allows robots to generate a higher-level judgment. For a better understanding, we'll look towards semantic segmentation common designs in the future.

# Overcoming the Challenges

Several deep networks have made such significant contributions to the industry that they are now widely accepted standards. Examples include AlexNet, VGG-16, GoogLeNet, and ResNet. They were so crucial that they're currently used in a number of segmentation systems as building blocks.

In order to get the most out of semantic segmentation, the instance-aware subtype should be used. Here are a few of the benefits behind this.

## The format is quite adaptable

With your data segmented, you can train and experiment with a variety of machine learning models, including classification, detection, and localization, picture creation, foreground/background separation, handwriting recognition, content alteration, and many others. As a result, it's employed in a variety of industries, including autonomous driving, fashion, film creation and post-production, agriculture, and so on.

## Precision unrivalled

Segmentation masks are the most exact since they only cover the position of the real item. Bounding boxes, on the other hand, frequently incorporate or connect with neighbouring territories. This is caused to non-rigid things being within or on top of other non-rigid objects.

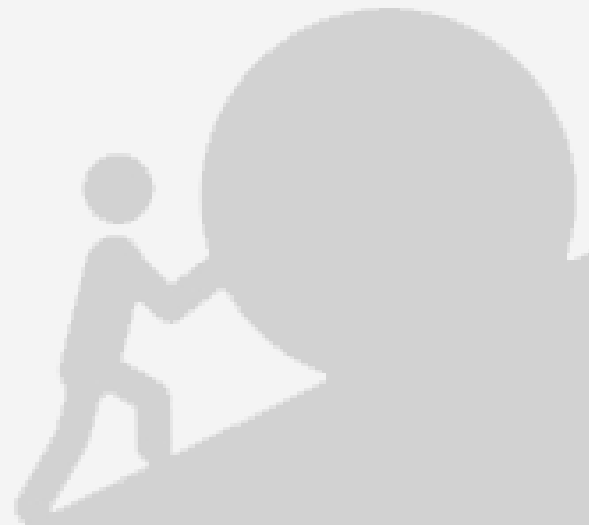
## One annotation with two annotations

Despite the fact that segmentation masks are more exact, bounding boxes are still used in many procedures. Fortunately, the surrounding bounding box can always be estimated using a segmentation mask. That's how you cover all of your bases!

## The confusion: Semantic Segmentation vs Instance segmentation

*To give a broad overview, segmentation determines which object category it belongs to, whereas instance segmentation, as the name implies, recognizes instances by assigning unique labels to them. This is not about class recognition, but about instance recognition, which means the system is seeking the same lookalike object in the scene, and any objects that seem different, even if they belong to the same class as the item in question, are disregarded.*

*This instance-level detection may be accomplished using lazy learning techniques: just store the description of the instance in a database, and during runtime, a matching score and a threshold are used to decide whether or not the instance is there.*





# Conclusion



Until date, a growing number of strategies have emerged to improve semantic segmentation accuracy or speed, or both. We hope that this overview of current advances, misconceptions, challenges and solutions in semantic segmentation will be useful to scholars working in this field.

Your segmentation challenge may be solved with the help of Cogito Tech LLC. We automate all of the hard elements of semantic segmentation, removing the complexity. Choosing Cogito as your annotation partners lowers the price of the annotation process and alleviates a lot of the set-up hassle.

# About Us

**Cogito** is a hybrid data labeling platform following model-assisted labeling (MAL) approach to cater industry's leading businesses. The MAL model leverages a human workforce to label a relevant portion of the training dataset which enables training of the AI application. Playing an important part as human-in-the-loop, our solutions encompass business verticals ranging from Retail, Manufacturing, Building, and Construction, to Medical, Food processing, E-commerce, and more.

[www.cogitotech.com](http://www.cogitotech.com)