

# Refined Plane Segmentation for Cuboid-Shaped Objects by Leveraging Edge Detection

Alexander Naumann\*, Laura Dörr\*, Niels Ole Salscheider\*, Kai Furmans†,

\*FZI Research Center for Information Technology, Karlsruhe, Germany

Email: {anaumann, doerr, salscheider}@fzi.de

†Institute for Material Handling and Logistics, Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany

Email: {kai.furmans}@kit.edu

*Abstract*—Recent advances in the area of plane segmentation from single RGB images show strong accuracy improvements and now allow a reliable segmentation of indoor scenes into planes. Nonetheless, fine-grained details of these segmentation masks are still lacking accuracy, thus restricting the usability of such techniques on a larger scale in numerous applications, such as inpainting for Augmented Reality use cases. We propose a post-processing algorithm to align the segmented plane masks with edges detected in the image. This allows us to increase the accuracy of state-of-the-art approaches, while limiting ourselves to cuboid-shaped objects. Our approach is motivated by logistics, where this assumption is valid and refined planes can be used to perform robust object detection without the need for supervised learning. Results for two baselines and our approach are reported on our own dataset, which we made publicly available. The results show a consistent improvement over the state-of-the-art. The influence of the prior segmentation and the edge detection is investigated and finally, areas for future research are proposed.

## I. INTRODUCTION

Identifying planar regions is an important task, that can be used in the context of segmentation and reconstruction for 3D scenes. Recent developments in the area of 3D plane detection from single RGB images [1], [2] open up opportunities to employ such approaches for various applications. Augmented Reality is one such application that can be used in numerous domains, ranging from logistics, manufacturing and military to education and entertainment [3]. Identified planes can be used to inpaint information and to place and simulate objects in a scene [4]. For example in manufacturing, Augmented Reality can be used to simulate, assist and improve processes [5]. In addition to that, robotics is a broad area of application, where plane segmentation can help to navigate, grasp and perform various other tasks. A common domain, where objects are confined to cuboid shapes is logistics. One benefit of using plane segmentation for object detection in such environments is its robustness compared to state-of-the-art segmentation techniques [6] that rely on knowing instances of the object categories beforehand. Thus, an accurate plane segmentation could be used for reconstruction or damage and tampering detection [7] in logistics contexts. Moreover, Augmented Reality can assist the packaging process to reduce error rates and document the process [8].

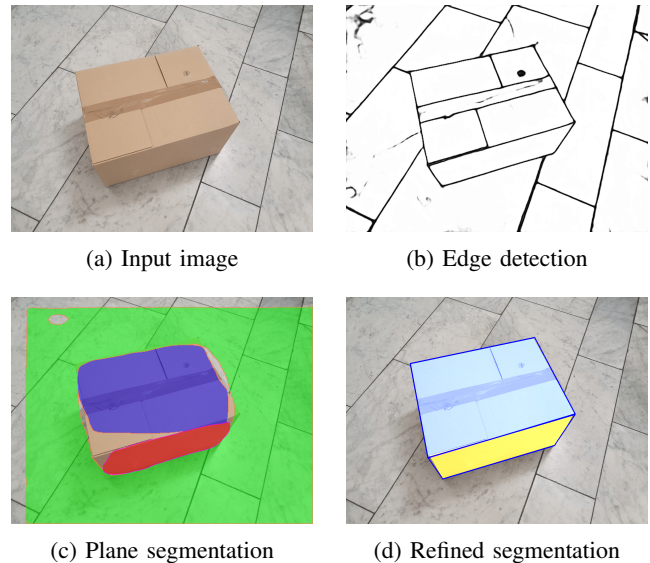


Figure 1: Overview of the segmentation pipeline. We consider an image (a) and combine edge detection (b) and plane segmentation (c) to obtain fine-grained plane segmentation results (d) close to the ground truth (blue contours).

Using state-of-the-art plane segmentation techniques alone is suitable for many applications, such as partial inpainting, however, these techniques still have difficulties in delivering fine-grained segmentation results. To overcome this issue, we propose a post-processing technique to refine the results of plane segmentation approaches such as the PlaneRCNN [2], to accurately detect the planes of cuboid-shaped objects in an image. Plane segmentation masks are refined separately and thus, not only cuboid-shaped objects, such as packages, but any plane with a rhombic shape can be rectified. In addition to that, the surfaces of those planes do not need to be completely flat, as in the case of a parcel, since also small 3D structures on the plane can be handled. Our approach uses different edge segmentation techniques [9], [10] to align the masks resulting from the plane segmentation with the edges in the image. For an overview of the pipeline, see Figure 1.

We collected a dataset of 34 images in a logistics context with ground truth plane segmentations, which we made publicly

available<sup>1</sup>. The performance of our approach, compared to the PlaneRCNN baseline [2] and a fallback routine which is also introduced in this work, is evaluated on this dataset. We show an overall improvement of the averaged Intersection over Union (IoU) of over 4.25 percentage points compared to the baseline. The quality of the refinement for individual images depends on the quality of the prior segmentation and the edge detection in the image. Several examples are provided to illustrate the capabilities of the approach even in difficult settings, but also to point out areas for further improvements. We will make our code publicly available under [https://url.fzi.de/refined\\_planeseg](https://url.fzi.de/refined_planeseg).

This work is structured as follows. In Section II we will present an overview over related literature. Section III outlines our plane segmentation mask refinement approach and Section IV evaluates our approach by comparing it to two baselines on our newly collected dataset. Section V concludes the paper.

## II. RELATED WORK

To the best of our knowledge, there has been no approach yet to improve the performance of state-of-the-art plane segmentation techniques by leveraging additional information from edge detection techniques.

Edge detection is a field thoroughly studied in computer vision [11]. Sobel's work [12] was one of the early contributions, that inspired numerous edge detection techniques. The Canny Edge Detector [9] is one of those techniques, that was developed in 1986 and is still a very common choice to date. Due to the popularity of the Canny Edge Detector, there has been a lot of research on improving it, for example by using an adaptive thresholds [13]. Recently, also Convolutional Neural Networks have been used to detect edges [14], [10] in images. These approaches often provide the advantage of being less sensitive to noise and are a promising alternative to classical techniques.

Image segmentation is a field intensively studied in computer vision that reached remarkable performance on closed-set configurations, where the objects of interests are known beforehand [15]. In the area of plane segmentation, End-to-End trained Neural Networks have only been presented recently [1], [2]. In addition to improving on the state-of-the-art in plane segmentation, these approaches also improved the state-of-the-art in single-image depth estimation.

## III. PLANE SEGMENTATION REFINEMENT

Our approach aims to improve the granularity of the plane segmentation by exploiting additional information from edge detection. We assume a prior segmentation by plane segmentation techniques, however, the procedure is independent of the source of these prior masks. In the following, we

introduce the edge detection techniques we consider in Section III-A. Afterwards, we present the plane segmentation approach whose segmentation masks are the input for our post-processing in Section III-B. Finally, we propose a method that leverages clustering and regression to find line segments along a rhombus in Section III-C and combine these ideas for the final refinement in Section III-D.

### A. Edge Detection

We consider two different techniques for edge segmentation. A recent work building upon the Canny Edge Detector [9] by automating the thresholding process [13], which we call Adaptive Canny and a machine learning based approach called DexiNed by Soria *et al.* [10]. We use two different forms of the latter approach, by applying it to the full resolution image (1280x960) and to a downsized image (640x480), since this seems to trigger a focus on important edges.

### B. Plane Segmentation

The PlaneRCNN deep neural architecture was introduced by Liu *et al.* in 2019 [2]. It improves upon earlier models [1] [16] by not requiring the maximum number of planes a priori and generalizing better to unseen scenes. The input to the model is a single RGB image, which is processed by three components. The first component is a Mask R-CNN [6] based model for plane detection. In addition to the plane segmentation, also the plane normal and depth values for each pixel are estimated. The second and third component are responsible for a refinement and the enforcement of consistency of the reconstructions. In this work, we use only the plane segmentations retrieved by the PlaneRCNN.

### C. Line Segmentation

Given a binary contour image and a segmentation of this image in planes, we consider the process of refining the given segmentation for each mask separately. We present two approaches for detecting line segments on the bounding edges of a plane belonging to a cuboid-shaped object, which will be combined in the final solution. The starting point for both approaches is to overlay the binary image with a widened contour line of the mask (See Figure 2a). For a reasonable prior segmentation, this contour should contain all or at least some of the bounding edges for the respective plane.

The first approach relies on Density-Based Spatial Clustering of Applications with Noise (DBSCAN) [17] to identify connected structures on the extract of the binary image. For a good prior segmentation, this extract of the binary image might comprise a completely connected cluster or similarly, only two or three clusters for rhombus as seen in Figure 2a. Since we try to estimate all line segments separately, we perform a corner detection [18] and break the clusters up by removing

<sup>1</sup>Dataset available under [https://url.fzi.de/dataset\\_planeseg](https://url.fzi.de/dataset_planeseg).

the detected corners from the binary mask. Thereafter, we are left with line-shaped clusters only as in Figure 2b, if all corners are detected correctly. In addition to that, we omit very small clusters and clusters with big variance in two directions, since they most likely constitute areas of noise. This leaves us with mainly line-shaped clusters of pixels. We use a RANSAC [19] linear regression to find a two point description for each of those lines. For each line, this first estimation is used to search for extensions of the line that were not captured by the widened contour line of the mask. This becomes necessary when the widened contour line only overlays with parts of the current edge, since non-overlapping contours were previously ignored. Hence, we create a mask along the estimated line across the whole image and repeat the former process. We apply clustering onto the new mask and perform a RANSAC linear regression on the dominant cluster to obtain new end point estimations for the current line.

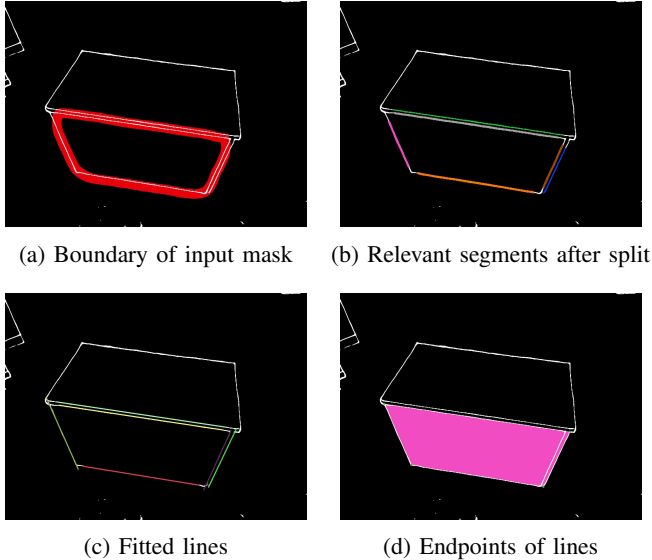


Figure 2: Overview of the pipeline for the segmentation refinement. We use the widened contour of the mask (a) to identify clusters of line segments (b). We try to fit a line for each segment (c), which can be used for the estimation of the final segmentation mask (d).

The second approach is based on the Hough Transform [20]. We cluster and average the lines resulting from applying the Hough Transform to the extract of the binary image. This approach is complementary to the first approach, since it does not rely on the connectedness within edge segments. Using the resulting lines in normal form, we estimate start and end points by leveraging their points of intersection. More precisely, for each line we compute its points of intersection with all other lines and assume that valid endpoints lie in the vicinity of the considered mask. The check for vicinity is performed by examining if a circle around the point of intersection with a 40 pixel radius intersects with the extract of the binary image.

#### D. Combined Approach

Still considering each plane separately, we combine the results from the first and the second approach from Section III-C. The result from each approach is a list of line segments described by a start and an end point, which represent approximations of the edges. We cluster the line segments to identify lines describing the same edge.

By using a k-means algorithm, we group start points into a set  $A$  and end points into a set  $B$ . These groups are each complemented by a set of ten randomly chosen points on the binary image in the vicinity of the centroids resulting from the k-means algorithm. We then identify the start point  $P_a \in A$  and the end point  $P_b \in B$  best fitting the considered edge  $e$ , for each cluster of line segments. Our cost function  $C$  describing the quality of the fit incorporates the overlap with the underlying mask  $m_e$  and the normalized length of the line with equal weights

$$C(P_a, P_b) = 0.5 \cdot I(\overline{P_a P_b}, m_e) + 0.5 \cdot \frac{\|P_a - P_b\|}{\max_{k \in A, l \in B} (\|P_k - P_l\|)},$$

where  $I(\overline{P_a P_b}, m_e)$  is the normalized intersection

$$I(\overline{P_a P_b}, m_e) = \frac{\overline{P_a P_b} \cap m_e}{\overline{P_a P_b}}.$$

The area  $\overline{P_a P_b}$  is defined by a line with one pixel thickness between  $P_a$  and  $P_b$ . Note that focusing on maximum overlap only might lead to very accurate, however, contracted line segments that do not represent the edges of the plane well.

Since we are aware of the dependence of our approach on the quality of the input information, we check the consistency of the refined mask with the prior mask by calculating their IoU. For large deviations, i.e. an IoU of less than 0.75, we resort back to a simple baseline approach. This baseline consist of calculating the prior mask's convex hull and iteratively reducing the set of points describing the mask [21] to 20 points or less.

Our approach is based on considering one mask at a time. Note that the dependencies between masks belonging to the same object can be used to further refine the segmentation results.

## IV. RESULTS AND EVALUATION

We first describe the dataset we collected in Section IV-A and comment on its separation into different classes. Subsequently, we will shortly present the baseline approaches used in the evaluation and finally discuss the results of our approach.

Dataset	PlaneRCNN	Fallback	Dexi LR	Dexi FR	Canny
Easy	83.85%	84.00%	<b>90.74%</b>	87.16%	82.03%
Medium	82.52%	82.67%	<b>84.75%</b>	82.66%	81.91%
Hard	79.88%	80.29%	<b>82.38%</b>	80.43%	79.42%
All	81.94%	82.17%	<b>85.20%</b>	83.10%	81.07%

Table I: Average IoU over all masks in each dataset for the PlaneRCNN, the fallback solution and our suggested approach with different edge segmentation techniques (Dexi LR = DexiNed with low resolution, Dexi FR = DexiNed with full resolution, Canny = Adaptive Canny).

### A. Dataset

We collected a set of 34 images, each picturing one cuboid object with different backgrounds. The backgrounds include carpet, table surfaces, floor tiles and the inside of a container. The cuboid-shaped objects include different types of parcels made from carton and plastic. The surfaces of the parcels range from almost plain to complex structures.

During the process of collecting the images, we manually discarded all images where the PlaneRCNN was not able to grasp the scene, i.e. where it did not detect all visible planes of the cuboid object of interest. This was mostly the case for flat objects and difficult camera angles, however, it also happened in some simpler scenes. We split up the dataset in three groups by manually assessing the complexity of the scenes, i.e. the quality of the prior segmentation masks and the edge detection. Of the 34 images, 7 were grouped into the category easy, 16 were grouped into the category medium and the remaining 11 were assigned the difficulty hard. Note that even the easy category contains diverse backgrounds and different types of packages. The data was labeled manually using the VGG Image Annotator [22].

### B. Baseline Approaches

The results from the PlaneRCNN [2] are used as input for our approach and thus, constitute a first baseline. In addition to that, we present the results of using our fallback solution that was described in Section III-D. Results separated by categories are presented in Table I. We use the IoU as metric for comparison as common for segmentation tasks. The IoU over all masks in an image is averaged and subsequently the average over all images in the dataset is taken.

As mentioned above, we removed images from the dataset, where the PlaneRCNN was not able to grasp the scene. If the PlaneRCNN is able to grasp the scene, it reaches 81.94% IoU with the ground truth masks. Our fallback solution shows a slight improvement on the PlaneRCNN by 0.23 percentage points for the IoU. Since it rectifies the PlaneRCNN masks, their shape is more consistent with the ground truth masks.

### C. Our Approach

The evaluation results for our approach using different edge detection techniques are reported in Table I. The evaluation shows an improvement of over 4.25 percentage points compared to the PlaneRCNN when edge detection is performed by DexiNed with a low resolution image as input. The results on the dataset classified as easy show an improvement of almost 7 percentage points. Thus, especially for reasonable prior segmentation masks and edge detections a considerable improvement over the baseline can be achieved. We exemplarily show such segmentation results in Figure 1 and Figure 3. Note that, even for complex backgrounds and feature-rich objects, our approach achieves good accuracy.



(a) Classified as easy.

(b) Classified as hard.

Figure 3: Exemplary results of the plane segmentation refinement. The blue contours constitute the ground truth.

The importance of the edge detection technique is seen by its strong influence on the evaluation results. We observe a consistent decline over all datasets moving from DexiNed with low resolution to DexiNed with full resolution and from DexiNed with full resolution to the Adaptive Canny algorithm. Using the adaptive Canny Edge Detector even yields results slightly worse than the baseline, as can be seen in Table I. This is due to feature-rich backgrounds and parcels where the bounding edges are not dominant.

By manually assessing the images, we identified reasons for bad segmentation results. The approach performs well when the prior mask and the edge detection yield reasonable inputs. Strong edges in the background as in Figure 4a and dominant lines on the parcel as the black label in Figure 4d can sometimes mislead the algorithm and cause deviations from the desired mask. In addition to that, the ability to break up the clustered segments around the mask by detecting and removing corners does affect the accuracy. In Figure 4b, for example, the vertical lines of the blue plane were not detected since they were clustered with the more dominant horizontal lines. Finally, imprecise prior masks can cause wrongfully merging two planes of an object as in Figure 4d or prohibit the algorithm from detecting the full surface of the plane as in Figure 4c.

## V. CONCLUSION

In this work, we presented an approach for the refinement of segmentation masks of cuboid-shaped objects. Existing state-of-the-art plane segmentation methods generalize well for the logistics environment we considered exemplarily and



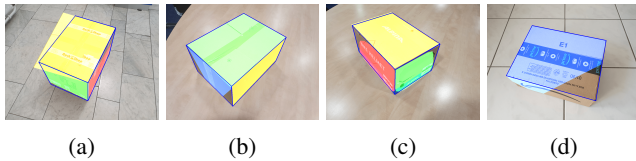


Figure 4: Visualization of the dependence of our approach on accurate prior information. Examples with strong background edges (a), failed corner detection (b), too small prior segmentation (c) and too big prior segmentation (d).

often yielded reasonable scene segmentations. However, those segmentation masks lack accuracy for fine-grained details such as corners. To enable the use of plane segmentation techniques for a wider range of applications, where this accuracy is necessary, we propose a post-processing technique. We combine edge detection and plane segmentation techniques with clustering approaches to perform a mask refinement along the edges of the object. To achieve robustness, we complement these techniques with a simple, yet effective fallback solution. Our approach improves the accuracy of state-of-the-art plane segmentation techniques by over 4.25 percentage points and generates masks that are more consistent in shape with the ground truth masks. The refined segmentation masks have several applications, for instance inpainting for Augmented Reality or object detection for cuboid-shaped objects. The further improvement of the segmentation by exploiting the fact that planes belong to the same object is left for future research.

#### REFERENCES

- [1] C. Liu, J. Yang, D. Ceylan, E. Yumer, and Y. Furukawa, “PlaneNet: Piece-Wise Planar Reconstruction from a Single RGB Image”, in *IEEE CONFERENCE on Computer Vision and Pattern Recognition*, 2018.
- [2] C. Liu, K. Kim, J. Gu, Y. Furukawa, and J. Kautz, “PlaneRCNN: 3D Plane Detection and Reconstruction From a Single Image”, in *IEEE CONFERENCE on Computer Vision and Pattern Recognition*, 2019.
- [3] D. Yu, J. S. Jin, S. Luo, W. Lai, and Q. Huang, “A Useful Visualization Technique: A Literature Review for Augmented Reality and its Application, limitation & future direction”, in *Visual Information Communication*, Springer US, 2010.
- [4] K. Karsch, K. Sunkavalli, S. Hadap, N. Carr, H. Jin, R. Fonte, *et al.*, “Automatic Scene Inference for 3D Object Compositing”, *ACM Transactions on Graphics*, vol. 33, no. 3, 2014.
- [5] S. K. Ong, M. L. Yuan, and A. Y. C. Nee, “Augmented reality applications in manufacturing: A survey”, *International Journal of Production Research*, vol. 46, no. 10, 2008.
- [6] K. He, G. Gkioxari, P. Dollar, and R. Girshick, “Mask R-CNN”, in *IEEE INTERNATIONAL CONFERENCE on Computer Vision*, 2017.
- [7] N. Noceti, L. Zini, and F. Odone, “A multi-camera system for damage and tampering detection in a postal security framework”, *EURASIP Journal on Image and Video Processing*, vol. 2018, no. 1, 2018.
- [8] M. Hochstein, J. Glöckle, T. Meyer, and K. Furmans, “Packassistent – assistenzsystem für die qualitätskontrolle während des packprozesses”, in *Logistics Journal*, Wiss. Gesellschaft für Technische Logistik, 2016.
- [9] J. Canny, “A Computational Approach to Edge Detection”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 6, 1986.
- [10] X. Soria, E. Riba, and A. D. Sappa, “Dense Extreme Inception Network: Towards a Robust CNN Model for Edge Detection”, 2020. arXiv: 1909.01955.
- [11] M. A. Oskoei and H. Hu, “A Survey on Edge Detection Methods”, University of Essex, Technical Report CES-506, 2010.
- [12] I. Sobel, “Camera Models and Machine Perception”, Technion Technical Report CS0016, 1972.
- [13] M. Fang, G. Yue, and Q. Yu, “The Study on An Application of Otsu Method in Canny Operator”, *International Symposium on Information Processing*, 2009.
- [14] S. Xie and Z. Tu, “Holistically-Nested Edge Detection”, in *IEEE International Conference on Computer Vision*, 2015.
- [15] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, “A Review on Deep Learning Techniques Applied to Semantic Segmentation”, 2017. arXiv: 1704.06857.
- [16] F. Yang and Z. Zhou, “Recovering 3D Planes from a Single Image via Convolutional Neural Networks”, in *EUROPEAN CONF. on Computer Vision*, 2018.
- [17] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, “A density-based algorithm for discovering clusters a density-based algorithm for discovering clusters in large spatial databases with noise”, in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, 1996.
- [18] C. Harris and M. Stephens, “A Combined Corner and Edge Detector”, in *Proceedings of the Alvey Vision Conference*, Alvey Vision Club, 1988.
- [19] M. A. Fischler and R. C. Bolles, “Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography”, in *Readings in Computer Vision*, Elsevier, 1987.
- [20] P. V. C. Hough, “Method and means for recognizing complex patterns”, U.S. Patent 3069654A, 1962.
- [21] D. H. Douglas and T. K. Peucker, “Algorithms for the reduction of the number of points required to represent a digitized line or its caricature”, *Cartographica: The International Journal for Geographic Information and Geovisualization*, vol. 10, no. 2, 1973.
- [22] A. Dutta and A. Zisserman, “The VIA Annotation Software for Images, Audio and Video”, in *ACM INTERNATIONAL CONFERENCE on Multimedia*, ACM, 2019.