

# Quality Evaluation Methods for Crowdsourced Image Segmentation

Doris Jung-Lin Lee, Akash Das Sarma, Aditya Parameswaran

## Abstract

Instance-level image segmentation provides rich information crucial for scene understanding in a variety of real-world applications, such as robotics and surveillance. In this paper we propose and evaluate several crowdsourced algorithms, including novel worker-aggregation based algorithms and retrieval-based methods based on prior work, for the image segmentation problem. We also characterize the different types of worker errors observed, and present a clustering algorithm that is able to capture semantic errors and filter workers with different semantic perspectives. We demonstrate that aggregation-based algorithms attains better performance than existing retrieval-based approaches, while scaling better with increasing numbers of collected worker segmentations.

## 1 Introduction

Precise, instance-level object segmentation is crucial for identifying and/or tracking objects in a variety of real-world emergent applications of autonomy, including robotics and autonomous vehicles, surveillance, image organization and retrieval, and medicine (Irshad and et. al. 2014; Yamaguchi 2012). To this end, there has been a lot of work on employing crowdsourcing to generate training data for computer vision, including Pascal-VOC (Everingham *et al.* 2015), LabelMe (Torralba *et al.* 2010), OpenSurfaces (Bell *et al.* 2015), and MS-COCO (Lin *et al.* 2012). Unfortunately, raw data collected from crowdsourced image processing tasks are known to be noisy due to varying degrees of worker skills, attention, and motivation (Bell *et al.* 2014; Welinder *et al.* 2010b).

In order to deal with these challenges, many have employed heuristics indicative of segmentation quality (Cabezas *et al.* 2015; Sameki *et al.* 2015; Sorokin and Forsyth 2008). While this approach identifies the segmentation drawn by more diligent workers, since it simply picks one *best* bounding box as the solution, it ends up discarding the majority of the worker responses and is limited by what the best worker can do. In this paper, we introduce a novel class of aggregation-based methods, capable of incorporating portions of responses from multiple workers into a combined segmentation and compare its performance with existing retrieval-based methods. In addition, we propose a preprocessing technique that resolves different worker perspectives in multiple segmentations.

## 2 Related Work

Many large-scale efforts in image segmentation contain little to no information on the quality characterization and

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

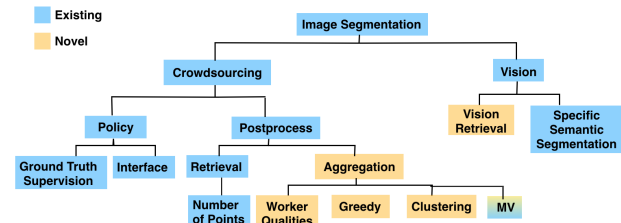


Figure 1: Flowchart summarizing the classes of existing algorithms for image segmentation (blue) and a novel class of algorithms proposed in this paper (yellow). Majority-vote (MV) is colored both blue and yellow, since a common algorithm in crowdsourcing literature, but have not been extensively applied to crowdsourced image segmentation.

evaluation of the collected dataset (Torralba *et al.* 2010; Martin *et al.* 2001; Li *et al.* 2009; Gurari *et al.* 2015), which indicate the lack of standardized approaches for quality evaluation in crowdsourced image segmentation. As shown in Figure 1, we break down the existing quality evaluation methods into several categories:

**Policy-based methods** include specialized segmentation interfaces or workflows that ensures that the data collected are of good quality, including periodic verification workflows (Lin *et al.* 2014; Everingham *et al.* 2015), specialized segmentation interfaces (Song *et al.* 2018), and vision supervision of crowdsourced segmentation (Russakovsky *et al.* 2015; Gurari *et al.* 2016).

**Retrieval-based methods** seek to pick the “best” worker segmentation based on some scoring criteria that evaluates the quality of each segmentation, including the use of vision information (Vittayakorn and Hays 2011; Russakovsky *et al.* 2015), expectation-maximization (EM) approaches for bounding box quality estimation (Welinder *et al.* 2010b), and click-stream behavior (Cabezas *et al.* 2015; Sameki *et al.* 2015; Sorokin and Forsyth 2008).

**Aggregation-based methods** use multiple worker segmentations to produce a single combined segmentation. Aggregation-based majority vote have been introduced in Sameki *et al.* (2015) as a way for aggregating expert segmentations to obtain a ground truth segmentation for evaluation purposes.

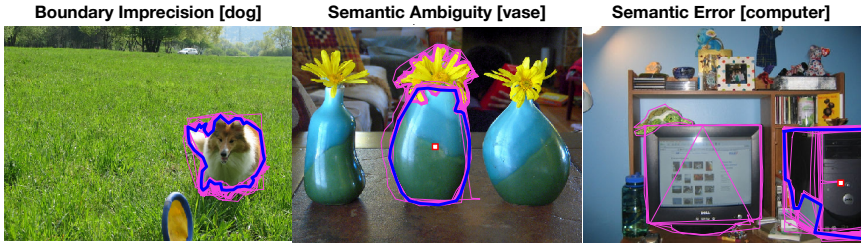


Figure 2: Pink is the segmentation from individual workers. Blue solid line delineates the ground truth. The red boxed pointer indicates the task of interest shown to users.

### 3 Preliminaries

#### 3.1 Data & Goals

We collected crowdsourced segmentation data from Amazon Mechanical Turk where each HIT consisted of one segmentation task for a specific pre-labeled object in the image. There were a total of 46 objects in 9 images from the MSCOCO dataset (Lin *et al.* 2014). For each object, we collected segmentation masks from a total of 40 workers. Each task contains a semantic keyword and a pointer indicating the object to be segmented. These tasks represent a diverse set of task difficulty (different levels of clutteredness, occlusion, lighting) and levels of task ambiguity.

#### 3.2 Evaluation Metrics

Evaluation metrics used in our experiment measures how well the final segmentation ( $S$ ) produced by these algorithms compare against ground truth ( $GT$ ). The most common evaluation metric used in literature are area-based methods which take into account the intersection,  $IA = area(S \cap GT)$ , or union,  $UA = area(S \cup GT)$ , between the user and the ground truth segmentations. Specifically, we use Precision ( $P$ ) =  $\frac{IA(S)}{area(S)}$ , Recall ( $R$ ) =  $\frac{IA(S)}{area(GT)}$ , and Jaccard ( $J$ ) =  $\frac{UA(S)}{IA(S)}$  metrics to evaluate our algorithms.

#### 3.3 Error Analysis

As shown in Figure 2 (left to right), common worker errors can be classified into three types: (1) **Boundary Imprecision**: unintentional mistakes while drawing the boundaries, either due to low precision of the image, small area of the object, or lack of drawing skills, (2) **Semantic Ambiguity**: have differing opinions about whether particular regions belong to part of an object; or (3) **Semantic Mistakes**: annotate the wrong object entirely.

Since the main focus for quality evaluation in past literature have been focused on finding worker segmentation with minimal boundary precision issues, we will first describe novel aggregation-based algorithms that we have developed and compare them with existing retrieval-based methods for addressing boundary imprecision. In the following section, we will discuss a preprocessing method that we have developed to resolve the semantic ambiguity and mistakes, which have been observed in prior work (Sorokin and Forsyth 2008; Lin *et al.* 2014; Gurari *et al.* 2018).

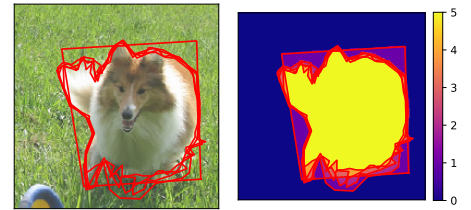


Figure 3: Left: Segmentation boundaries drawn by five workers in red. Right: Overlaid segmentation creates a masks where the color indicates the number of workers who voted for the tile region.

### 4 Precision-savvy algorithms: Aggregation v.s. Retrieval Comparison

At the heart of all our aggregation techniques is the “tile” data representation, where we logically overlay all workers’ segmentations on top of each other within the framework of the overall image, as illustrated in Figure 3, to create non-overlapping discrete tile units. The intuition here is that by splitting the image into tiles, we get finer granularity information than by looking at complete segmentations. This also allows us to aggregate data from multiple workers rather than having to choose a single worker bounding box—this allows for the potential of choosing the best partial segmentations for an object and joining them, or fixing one worker’s errors by taking the help of another worker’s segmentation. Now, we will describe several algorithms for picking a good set of tiles.

**Aggregation: Majority Vote Aggregation (MV)**: Include tiles in the output segmentation if and only if the tile is covered by at least 50% of all worker segmentations.

**Aggregation: Expectation-Maximization (EM)** Unlike MV, which assumes that all workers performs uniformly, EM approaches use worker quality models to infer the likelihood that a tile is part of the ground truth segmentation. An EM framework is used to simultaneously estimate both worker qualities and tile likelihoods as hidden variables. We detail the formal derivation and three worker quality models that we have developed in more detail in our technical report. **Aggregation: Greedy Tile Picking (greedy)** Using estimated tile probabilities to estimate intersection area between ground truth and tile. Greedily pick tiles with the largest intersection area ratio until Jaccard score begins to decrease. The Jaccard score is computed between the merged output from the selected set of tiles and MV segmentation.

**Retrieval: Number of Control Points (num pts)** Pick the worker segmentation with the largest number of control points around the segmentation boundary (i.e. most precise drawing) as the output segmentation.

#### 4.1 Retrieval v.s. Aggregation-based Comparison

**Aggregation-based methods performs significantly better than retrieval-based methods**

Figure 4 left shows that amongst the algorithms that do not make use of ground truth information, the performance of aggregation-based algorithms (greedy, EM) exceeds the best achievable through the existing retrieval-based method (num pts) and the vision-based algorithms. By making use of

Algorithm	Retrieval-based		Aggregation-based			
	num pts	best worker	MV	EM	greedy	best greedy
Worker Scaling	-6.30	2.58	1.63	1.64	2.16	5.59
Clustering Effect	5.92	-0.02	2.05	1.38	5.55	-0.06

Table 1: The first row lists the average percentage change in Jaccard between 5 workers samples and 30 workers sample. The second row lists the average percentage change between the no clustering and clustering results.

ground truth information, the best aggregation-based algorithm can achieve a close-to-perfect average Jaccard score of 0.983 as an upper bound with the 30 workers sample, far exceeding the results achievable by any single ‘best’ worker ( $J=0.91$  for 30 workers in Figure 4 right). This result demonstrates that aggregation-based methods are able to achieve better performance by performing inference at the *tile* granularity, which is guaranteed to be finer than any worker’s segmentations.



Figure 4: Jaccard performance comparison between best-performing algorithms from retrieval and aggregation-based methods with clustering as a preprocessing step where possible. We compare between the original algorithms that do not make use of ground truth information (Left) and ones that do (Right).

### Performance of aggregation-based methods scales well as more workers segmentation are added.

Intuitively, larger worker samples results in finer granularity tiles for the aggregation-based methods, resulting in an monotonically increasing relationship between number of worker segmentation used in the sample and performance evident in Table 1. However, worker scaling for retrieval-based methods are not guaranteed.

## 5 Perspective Resolution in Crowdsourced Image Segmentation

### 5.1 Worker Clustering

Our clustering-based approach is based on the intuition that workers with similar perspectives will have segmentations that are closer to each other, while workers with different perspectives from each other will have segmentations that

differ from each other. We capture the similarity between a pair of workers by computing the Jaccard score between their segmentations and perform *spectral clustering* to separate workers into clusters. Figure 5 illustrates how spectral clustering is capable of dividing the worker responses into clusters with meaningful semantic associations, reflecting the crowd’s diversity of perspectives in completing same task.

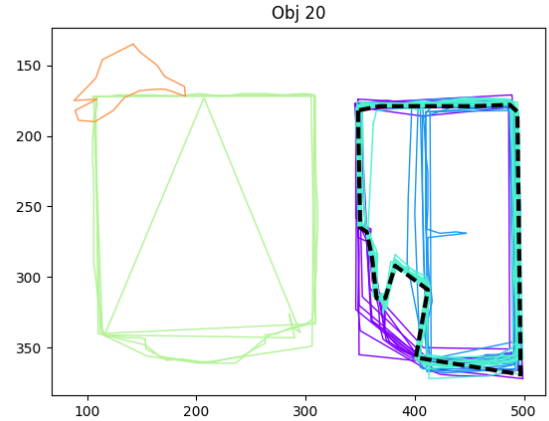


Figure 5: Example image showing clustering performed on the same object from Figure 2.

Clustering results can also be used as a preprocessing step to any of the quality evaluation algorithms by keeping only the segmentations that belong to the largest cluster, which is typically free of any semantic errors. As shown in Table 1, on average, clustering generally results in an increase the resulting algorithmic performance. Since the ground-truth supervised variants are already free of semantic ambiguity and errors, there is minimal improvement resulting from clustering.

Clustering offers an additional benefit of preserving worker’s semantic intentions in the case where there are multiple instances of different errors. For example, in Figure 5, the mistaken clusters included semantic concepts “monitor” and “turtle”. While these are considered *bad* annotations for this particular task, this cluster of annotation can provide more data for another semantic segmentation task “monitor”. A potential direction of future work includes adding an additional crowdsourcing task for semantic labelling of clusters (which is cheaper and more accurate than segmentation) to enable reuse of annotations across objects and lower the cost of data collection.

## 6 Conclusion & Future Work

Given the success of aggregation-based models, immediate future work investigates why and how majority vote —, we have preliminary studies that shows that worker qualities are good indicators of actual —. why majority vote, while simple, performs nearly as well as the advanced EM and greedy based approaches. This is because both EM and greedy have learned worker qualities that converged to MV behavior.

In this paper, we perform an extensive study of several image segmentation algorithms spanning semi-supervised vision approaches, crowdsourced retrieval approaches, and novel crowdsourced aggregation approaches. We identified three different types of errors that workers typically make on segmentation tasks, some caused by differing perspectives, and developed a clustering-based method to filter out workers that are making semantic errors. We demonstrate the strength of our worker clustering algorithm as well as the aggregation-based segmentation algorithms through extensive experiments in 1) its ability to improve as more worker segmentations are collected and 2) yield better performance than retrieval-based methods. We also found that while majority vote is a fairly simple algorithm, it performs nearly as well as the advanced EM and greedy inference approaches. Our code is open source and available for researchers to benchmark and compare techniques. Our work represents a first step in understanding and comparing the different types of algorithms available for image segmentation tasks. It opens a number of exciting directions for exploration, for instance: (a) Studying the effect of task difficulty, or worker qualities across different objects, and (b) Designing better hybrid algorithms that combine the different types of algorithms described in this paper.

## References

- [Bell *et al.* 2014] Sean Bell, Kavita Bala, and Noah Snavely. Intrinsic images in the wild. *ACM Trans. on Graphics (SIGGRAPH)*, 33(4), 2014.
- [Bell *et al.* 2015] Sean Bell, Paul Upchurch, Noah Snavely, and Kavita Bala. Material recognition in the wild with the materials in context database. *Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [Cabezas *et al.* 2015] Ferran Cabezas, Axel Carlier, Vincent Charvillat, Amaia Salvador, and Xavier Giro-I-Nieto. Quality control in crowdsourced object segmentation. *Proceedings of International Conference on Image Processing, ICIP*, 2015-Decem:4243–4247, 2015.
- [Everingham *et al.* 2015] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, January 2015.
- [Gurari *et al.* 2015] Danna Gurari, Diane Theriault, Mehrnoosh Sameki, Brett Isenberg, Tuan A. Pham, Alberto Purwada, Patricia Solski, Matthew Walker, Chentian Zhang, Joyce Y. Wong, and Margrit Betke. How to collect segmentations for biomedical images? A benchmark evaluating the performance of experts, crowdsourced non-experts, and algorithms. *Proceedings of 2015 IEEE Winter Conference on Applications of Computer Vision (WACV) 2015*, pages 1169–1176, 2015.
- [Gurari *et al.* 2016] Danna Gurari, Mehrnoosh Sameki, Zheng Wu, and Margrit Betke. Mixing Crowd and Algorithm Efforts to Segment Objects in Biomedical Images. *Medical Image Computing and Computer Assisted Intervention Interactive Medical Image Computation Workshop*, pages 1–8, 2016.
- [Gurari *et al.* 2018] Danna Gurari, Kun He, Bo Xiong, Jianming Zhang, Mehrnoosh Sameki, Suyog Dutt Jain, Stan Sclaroff, Margrit Betke, and Kristen Grauman. Predicting Foreground Object Ambiguity and Efficiently Crowdsourcing the Segmentation(s). *International Journal of Computer Vision (IJCV)*, 2018.
- [Irshad and et. al. 2014] H Irshad and Montaser-Kouhsari et. al. Crowdsourcing Image Annotation for Nucleus Detection and Segmentation in Computational Pathology: Evaluating Experts, Automated Methods, and the Crowd. *Biocomputing 2015*, pages 294–305, 2014.
- [Li *et al.* 2009] Li Jia Li, Richard Socher, and Li Fei-Fei. Towards total scene understanding: Classification, annotation and segmentation in an automatic framework. *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2009*, pages 2036–2043, 2009.
- [Lin *et al.* 2012] Christopher H Lin, Mausam, and Daniel S Weld. Crowdsourcing control : Moving beyond multiple choice. *AAAI Conference on Human Computation and Crowdsourcing (HCOMP)*, pages 491—500, 2012.
- [Lin *et al.* 2014] Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. *European Conference on Computer Vision (ECCV)*, 8693 LNCS(PART 5):740–755, 2014.
- [Martin *et al.* 2001] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. 8th Int’l Conf. Computer Vision*, volume 2, pages 416–423, July 2001.
- [Russakovsky *et al.* 2015] Olga Russakovsky, Li-Jia Li, and Li Fei-Fei. Best of Both Worlds: Human-Machine Collaboration for Object Annotation. pages 2121–2131, 2015.
- [Sameki *et al.* 2015] Mehrnoosh Sameki, Danna Gurari, and Margrit Betke. Characterizing Image Segmentation Behavior of the Crowd. pages 1–4, 2015.
- [Song *et al.* 2018] Jean Y. Song, Raymond Fok, Alan Lundgard, Fang Yang, Juho Kim, and Walter S. Lasecki. Two Tools are Better Than One : Tool Diversity as a Means of Improving Aggregate Crowd Performance. *IUI’18: Proceedings of the International Conference on Intelligent User Interfaces*, 2018.
- [Sorokin and Forsyth 2008] Alexander Sorokin and David Forsyth. Utility data annotation with Amazon Mechanical Turk. *Proceedings of the 1st IEEE Workshop on Internet Vision at CVPR 08*, (c):1–8, 2008.
- [Torralba *et al.* 2010] Antonio Torralba, Bryan C. Russell, and Jenny Yuen. LabelMe: Online image annotation and applications. *Proceedings of the IEEE*, 98(8):1467–1484, 2010.
- [Vittayakorn and Hays 2011] Sirion Vittayakorn and James Hays. Quality Assessment for Crowdsourced Object Annotations. *Proceedings of the British Machine Vision Conference 2011*, pages 109.1–109.11, 2011.
- [Welinder *et al.* 2010b] Peter Welinder, Steve Branson, Serge

Belongie, and Pietro Perona. The Multidimensional Wisdom of Crowds. *NIPS (Conference on Neural Information Processing Systems)*, 6:1–9, 2010b.

[Yamaguchi 2012] Kota Yamaguchi. Parsing clothing in fashion photographs. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, CVPR '12, pages 3570–3577, Washington, DC, USA, 2012. IEEE Computer Society.