

# Quality Evaluation Methods for Crowdsourced Image Segmentation

Doris Jung-Lin Lee  
University of Illinois,  
Urbana-Champaign  
jlee782@illinois.edu

Akash Das Sarma  
Facebook, Inc.  
akashds@fb.com

Aditya Parameswaran  
University of Illinois,  
Urbana-Champaign  
adityagp@illinois.edu

## ABSTRACT

Instance-level image segmentation provides rich information crucial for scene understanding in a variety of real-world applications. In this paper, we evaluate multiple crowdsourced algorithms for the image segmentation problem, including novel worker-aggregation-based methods and retrieval-based methods from prior work. We characterize the different types of worker errors observed in crowdsourced segmentation, and present a clustering algorithm as a preprocessing step that is able to capture and eliminate errors arising due to workers having different semantic perspectives. We demonstrate that aggregation-based algorithms attain higher accuracies than existing retrieval-based approaches, while scaling better with increasing numbers of worker segmentations.

## ACM Reference Format:

Doris Jung-Lin Lee, Akash Das Sarma, and Aditya Parameswaran. 1997. Quality Evaluation Methods for Crowdsourced Image Segmentation. In *Proceedings of ACM Woodstock conference (WOODSTOCK'97)*, Jennifer B. Sartor, Theo D'Hondt, and Wolfgang De Meuter (Eds.). ACM, New York, NY, USA, Article 4, 4 pages. [https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

## 1 INTRODUCTION

Precise, instance-level object segmentation is crucial for identifying and tracking objects in a variety of real-world emergent applications of autonomy, including robotics [11], image organization and retrieval [19], and medicine [8]. To this end, there has been a lot of work on employing crowdsourcing to generate training data for segmentation, including Pascal-VOC [5], LabelMe [16], OpenSurfaces [2], and MS-COCO [9]. Unfortunately, raw data collected from the crowd is known to be noisy due to varying degrees of worker skills, attention, and motivation [1, 18].

To deal with these challenges, many have employed heuristics indicative of crowdsourced segmentation quality to pick the best worker-provided segmentation [15, 17]. However, this approach ends up discarding the majority of the worker segmentations and is limited by what the best worker can do. The contributions of this paper is as follows:

- We introduce a novel class of *aggregation-based* methods that incorporates portions of segmentations from multiple workers into a combined one described in Section 4. By overlaying worker segmentations on top of each other, we can decompose

the image into non-overlapping tiles, where each tile has some workers who believe this tile belongs to the object, and others who do not. Each tile can be treated as an independent boolean question, deriving an answer from a worker—does this tile belong to the object or not, following which we may be able to apply Expectation-Maximization (EM) [4] to derive maximum likelihood tiles and worker accuracies, a greedy approach for tile picking based on worker fraction votes, and simple majority vote aggregation.

- To our surprise, despite the intuitive simplicity of aggregation-based methods, we have not seen this class of algorithms described or evaluated in prior work. We evaluate this class of algorithms against existing methods in Section 6 and found that it performs much better than existing approaches
- We formally characterize the types of worker error in crowdsourced image segmentation in Section 3 and describe a well-known multiple perspective issue in crowdsourced image segmentation [10, 15? ], where workers often segment the wrong objects or erroneously include or exclude large semantically-ambiguous portions of an object in the resulting segmentation. To address this issue, in Section 5, we develop a clustering-based solution which can be applied as a preprocessing step to any quality evaluation methods.

## 2 RELATED WORK

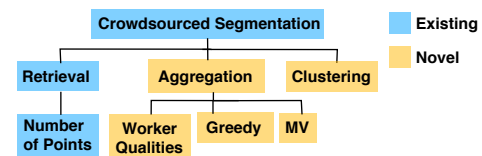


Figure 1: Taxonomy of existing algorithms for crowdsourced segmentation (blue) and a novel class of algorithms proposed in this paper (yellow).

As shown in Figure 1, quality evaluation methods for crowdsourced segmentation can be classified into the following categories: **Retrieval-based methods** pick the “best” worker segmentation based on some scoring criteria that evaluates the quality of each segmentation, including vision information [12, 17], and click-stream behavior [3, 13, 15].

**Aggregation-based methods** combine multiple worker segmentations to produce a final segmentation that is not restricted to any single worker segmentation. An aggregation-based majority vote approach was employed in Sameki et al. (2015) to create an expert-established gold standard for characterizing their dataset and algorithmic accuracies, rather than for segmentation quality evaluation as described here.

**Orthogonal methods** to improve segmentation quality include periodic verification [5, 10], specialized interfaces [14], and vision-based supervision [7, 12]. These methods could be used for quality improvement on top of any of the algorithms in this paper. Since these policy-based methods are often interface-dependent or require expensive expert-drawn ground-truth annotations or vision information, their results are not easily reproducible.

### 3 ERROR ANALYSIS

On collecting and analyzing a number of crowdsourced segmentations (described in Section 6), we found that common worker segmentation errors can be classified into three types:

- **Semantic Ambiguity:** workers have differing opinions on whether particular regions belong to an object (Figure 2 left: annotations around ‘flower and vase’ when ‘vase’ is requested);
- **Semantic Mistake:** workers annotate the wrong object entirely (Figure 2 right: annotations around ‘turtle’ and ‘monitor’ when ‘computer’ is requested.)
- **Boundary Imperfection:** workers make unintentional mistakes while drawing the boundaries, either due to low image resolution, small area of the object, or lack of drawing skills (Figure 3 left: imprecision around the ‘dog’ object).

Semantic ambiguity and mistakes have also been observed in prior work [6, 10, 15], which noted that disagreement in worker responses can come from questions that are ambiguous or difficult to answer, such as segmenting a individual person from a crowd. Since there are multiple workers annotating each object, each object can suffer from multiple types of error: we found that out of the 46 objects in our dataset, 9 objects suffered from type one error and 18 objects from type two error. Almost all objects suffer from some form of type three error of varying degrees of imprecision around the object boundary. The main evaluation methods highlighted in Section 4 focuses on resolving the imprecise, “sloppy” bounding box errors. In Section 5, we discuss a preprocessing method eliminates semantic ambiguities and errors.

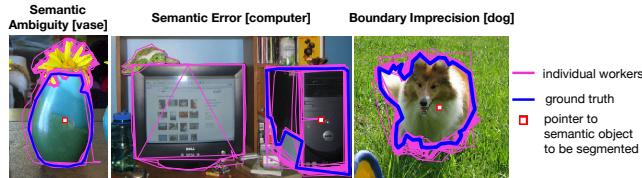


Figure 2: Examples of common worker mistakes.

### 4 FIXING BOUNDARY IMPERFECTIONS

At the heart of our aggregation techniques is the *tile data representation*. A tile is the smallest non-overlapping discrete unit created by overlaying all of the workers’ segmentations on top of each other. The tile representation allows us to aggregate segmentations from multiple workers, rather than being restricted to a single worker’s segmentation—allowing us to fix one worker’s errors with help from another. In Figure 3 (right), we display three worker segmentations for a toy example with 6 resulting tiles. Any subset of these tiles can contribute towards the final segmentation.

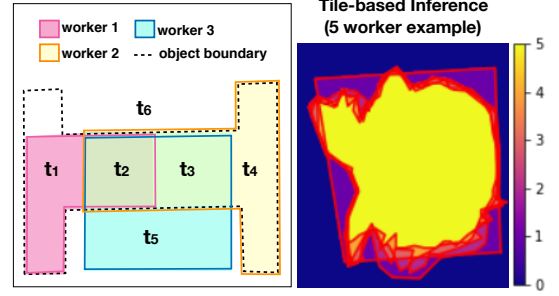


Figure 3: Left: Toy example demonstrating tiles created by three workers’ segmentations around a dumbbell object delineated by the black dotted line. Right: Segmentation boundaries drawn by five workers shown in red. Overlaid segmentation creates a mask where the color indicates the number of workers who voted for the tile region.

This simple but powerful idea of tiles also allows us to reformulate our problem from one of “generating a segmentation” to a setting that is much more familiar to crowdsourcing researchers. Since tiles are the lowest granularity units created by overlaying all workers’ segmentations on top of each other, each tile is either completely contained within or outside a given worker segmentation. Specifically, we can regard a worker segmentation as multiple boolean responses where they have voted ‘yes’ or ‘no’ to every tile independently. Intuitively, a worker votes ‘yes’ for every tile that is contained in their segmentation, and ‘no’ for every tile that is not. As shown in Figure 3 (right), tile  $t_2$  is voted ‘yes’ by worker 1, 2, and 3; tile  $t_3$  is voted ‘yes’ by worker 2 and 3. The goal of our aggregation algorithms is to pick an appropriate set of tiles that effectively trades off precision versus recall.

Now that we have modeled segmentation as a collection of worker votes for tiles, we can now develop familiar variants of standard quality evaluation algorithms for this setting.

#### Aggregation: Majority Vote Aggregation (MV)

This simple algorithm includes a tile in the output segmentation if and only if the tile has ‘yes’ votes from at least 50% of all worker segmentations.

#### Aggregation: Expectation-Maximization (EM)

Unlike MV, which assumes that all workers perform uniformly, EM approaches use worker quality models to infer the likelihood that a tile is part of the ground truth segmentation. While simultaneously estimating worker qualities and tile likelihoods as hidden variables, our basic worker quality model that we evaluate in Section 6 assumes a fixed probability for a correct vote. Details of the formal derivation and other more fine-grained worker quality models can be found in our technical report.

#### Aggregation: Greedy Tile Picking (greedy)

The greedy algorithm picks tiles in descending order based on the ratios of overlap area to non-overlap area (both with respect to ground truth), for as long as the estimated Jaccard similarity of the resulting segmentation continue to increase. Since the tile overlap and non-overlap against ground truth are unknown, we use tile-inclusion probabilities from EM to estimate these areas as a heuristic. Furthermore, since we cannot compute the actual Jaccard similarity against the unknown ground truth, we use a heuristic baseline such as MV as a proxy for the ground truth. Intuitively, tiles that have a high overlap area and low non-overlap area contribute to high

recall, at the cost of relatively little precision error. We include a proof in our technical report showing that picking tiles in such an order maximizes the Jaccard similarity of the resulting segmentation locally at every step.

#### Retrieval: Number of Control Points (num pts)

This algorithm picks the worker segmentation with the largest number of control points around the segmentation boundary (i.e., the most precise drawing) as the output segmentation [15, 17].

## 5 PERSPECTIVE RESOLUTION

As discussed in Section 3, disagreements often arise in segmentation due to differing worker perspectives on large tile regions. We developed a clustering-based preprocessing approach to resolve this issue. Based on the intuition that workers with similar perspectives will have segmentations that are close to each other, we compute the Jaccard similarity between each pair of segmentations and perform spectral clustering to separate the segmentations into clusters. Figure 2 (bottom) illustrates how spectral clustering divides the worker segmentations into clusters with meaningful semantic associations, reflecting the diversity of perspectives for the same task. Clustering results can be used as a preprocessing step for any quality evaluation algorithm by keeping only the segmentations that belong to the largest cluster, which is typically free of semantic errors.

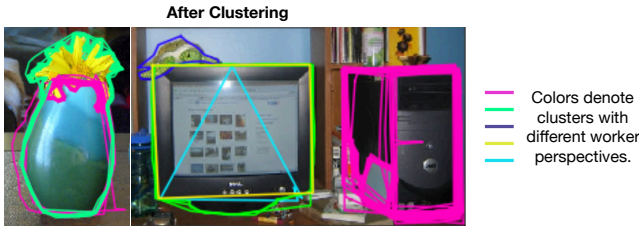


Figure 4: Example image showing clustering performed on the same object from Figure 2 left and middle.

In addition, clustering offers the additional benefit of preserving worker’s semantic intentions. For example, while the green cluster in Figure 2 (bottom right) would be considered *bad* segmentations for the particular task (‘computer’), this cluster can provide more data for another segmentation task corresponding to ‘monitor’. A potential future work direction would be to crowdsource the semantic labels for the computed clusters to enable the reuse of segmentations across multiple objects to lower costs.

## 6 EXPERIMENTAL EVALUATION

### Dataset Description

We collected crowdsourced segmentations from Amazon Mechanical Turk; each HIT consisted of one segmentation task for a specific pre-labeled object in an image. There were a total of 46 objects in 9 images from the MSCOCO dataset [10] segmented by 40 different workers each, resulting in a total of 1840 segmentations. Each task contained a keyword for the object and a pointer indicating the object to be segmented. Two of the authors generated the ground truth segmentations by carefully segmenting the objects using the same task and interface.

### Evaluation Metrics

Evaluation metrics used in our experiments measure how well the final segmentation (S) produced by these algorithms compare against

ground truth (GT). The most common evaluation metrics used in the literature are area-based methods that take into account the intersection area,  $IA = \text{area}(S \cap GT)$ , or union area,  $UA = \text{area}(S \cup GT)$  between the worker and ground truth segmentations, including Precision ( $P = \frac{IA(S)}{\text{area}(S)}$ ), Recall ( $R = \frac{IA(S)}{\text{area}(GT)}$ ), and Jaccard ( $J = \frac{IA(S)}{UA(S)}$ ).

### Experiment 1: Aggregation-based methods perform significantly better than retrieval-based methods

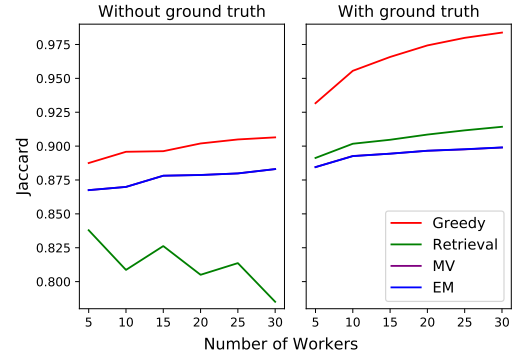


Figure 5: Performance of the original algorithms that do not make use of ground truth information (Left) and ones that do (Right). MV and EM results are so close that they overlay on each other.

In Figure 5, we vary the number of worker segmentations along the x-axis and plot the average Jaccard score on the y-axis across different worker samples of a given size across different algorithms. Figure 5 (left) shows that the performance of aggregation-based algorithms (greedy, EM) exceeds the best-achievable through existing retrieval-based method (Retrieval). Then, in Figure 5 (right), we estimate the upper-bound performance of each algorithm by assuming that the ‘full information’ based on ground truth was given to the algorithm. For greedy, the algorithm is aware of all the actual tile overlap and non-overlap areas against ground truth, and does not need to approximate these values. For EM, we consider the performance of the algorithm if the true worker quality parameter values (under our worker quality model) are known. For retrieval, the full information version directly picks the worker with the highest Jaccard similarity with respect to the ground truth segmentation. By making use of ground truth information (Figure 5 right), the best aggregation-based algorithm can achieve a close-to-perfect average Jaccard score of 0.98 as an upper bound, far exceeding the results achievable by any single ‘best’ worker ( $J=0.91$ ). 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 grained than any individual worker segmentation.

### The performance of aggregation-based methods scale well as more worker segmentations are added.

Intuitively, larger numbers of worker segmentations result in finer granularity tiles for the aggregation-based methods. The first row in Table 1 lists the average percentage change in Jaccard between 5-workers and 30-workers samples, demonstrating a monotonically increasing relationship between number of worker segmentations used and the performance. However, retrieval-based methods do not benefit from more segmentations.

## Experiment 2: Clustering as preprocessing improves algorithmic performance.

The average percentage change between the no clustering and clustering results is shown in Table 1. Clustering generally results in an accuracy increase. Since the ‘full information’ variants are already free of semantic ambiguity and errors, clustering does not assist with further improvement.

Algorithm	Retrieval-based		Aggregation-based			
	num pts	worker*	MV	EM	greedy	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: Percentage change due to worker scaling and clustering. Algorithms with \* makes use of ground truth information.**

## 7 CONCLUSION AND FUTURE WORK

We identified three different types of errors for crowdsourced image segmentation, developed a clustering-based method to capture the semantic diversity caused by differing worker perspectives, and introduced novel aggregation-based methods that produce more accurate segmentations than existing retrieval-based methods.

Our preliminary studies show that our worker quality models are good indicators of the actual accuracy of worker segmentations. We also observe that the greedy algorithm is capable of achieving close-to-perfect segmentation accuracy with ground truth information. Given the success of aggregation-based methods, including the simple majority vote algorithm, we plan to use our worker quality insights to improve our EM and greedy algorithms. We are also working on using computer vision signals to further improve our algorithms.

## REFERENCES

- [1] Sean Bell, Kavita Bala, and Noah Snavely. 2014. Intrinsic Images in the Wild. *ACM Trans. on Graphics (SIGGRAPH)* 33, 4 (2014).
- [2] Sean Bell, Paul Upchurch, Noah Snavely, and Kavita Bala. 2015. Material Recognition in the Wild with the Materials in Context Database. *Computer Vision and Pattern Recognition (CVPR)* (2015).
- [3] Ferran Cabezas, Axel Carlier, Vincent Charvillat, Amaia Salvador, and Xavier Giro-I-Nieto. 2015. Quality control in crowdsourced object segmentation. *Proceedings of International Conference on Image Processing, ICIP 2015-Decem* (2015), 4243–4247. <https://doi.org/10.1109/ICIP.2015.7351606> arXiv:arXiv:1505.00145v1
- [4] A. P. Dawid and A. M. Skene. 1979. Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm. 28, 1 (1979), 20–28.
- [5] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. 2015. The Pascal Visual Object Classes Challenge: A Retrospective. *International Journal of Computer Vision* 111, 1 (Jan. 2015), 98–136.
- [6] Danna Gurari, Kun He, Bo Xiong, Jianming Zhang, Mehrnoosh Sameki, Suyog Dutt Jain, Stan Sclaroff, Margrit Betke, and Kristen Grauman. 2018. Predicting Foreground Object Ambiguity and Efficiently Crowdsourcing the Segmentation(s). *International Journal of Computer Vision (IJCV)* (2018). <https://doi.org/10.1007/s11263-018-1065-7> arXiv:1705.00366
- [7] Danna Gurari, Mehrnoosh Sameki, Zheng Wu, and Margrit Betke. 2016. Mixing Crowd and Algorithm Efforts to Segment Objects in Biomedical Images. *Medical Image Computing and Computer Assisted Intervention Interactive Medical Image Computation Workshop* (2016), 1–8.
- [8] H Irshad and Montaser-Kouhsari et. al. 2014. Crowdsourcing Image Annotation for Nucleus Detection and Segmentation in Computational Pathology: Evaluating Experts, Automated Methods, and the Crowd. *Bioinformatics* 29, 2 (2014), 294–305. [https://doi.org/10.1142/9789814644730\\_0029](https://doi.org/10.1142/9789814644730_0029)
- [9] Christopher H Lin, Mausam, and Daniel S Weld. 2012. Crowdsourcing control : Moving beyond multiple choice. *AAAI Conference on Human Computation and Crowdsourcing (HCOMP)* (2012), 491–500. arXiv:arXiv preprint arXiv:1210.4870.
- [10] Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. *European Conference on Computer Vision (ECCV)* 8693 LNCS, PART 5 (2014), 740–755. [https://doi.org/10.1007/978-3-319-10602-1\\_48](https://doi.org/10.1007/978-3-319-10602-1_48)
- [11] E. Natonek. 1998. Fast range image segmentation for servicing robots. In *Proceedings. 1998 IEEE International Conference on Robotics and Automation (Cat. No.98CH36146)*, Vol. 1. 406–411 vol.1. <https://doi.org/10.1109/ROBOT.1998.676445>
- [12] Olga Russakovsky, Li-Jia Li, and Li Fei-Fei. 2015. Best of Both Worlds: Human-Machine Collaboration for Object Annotation. (2015), 2121–2131.
- [13] Mehrnoosh Sameki, Danna Gurari, and Margrit Betke. 2015. Characterizing Image Segmentation Behavior of the Crowd. (2015), 1–4.
- [14] Jean Y. Song, Raymond Fok, Alan Lundgard, Fang Yang, Juho Kim, and Walter S. Lasecki. 2018. Two Tools are Better Than One : Tool Diversity as a Means of Improving Aggregate Crowd Performance. *Proceedings of the International Conference on Intelligent User Interfaces* (2018).
- [15] Alexander Sorokin and David Forsyth. 2008. Utility data annotation with Amazon Mechanical Turk. *Proceedings of the 1st IEEE Workshop on Internet Vision at CVPR 08 c* (2008), 1–8. <https://doi.org/10.1109/CVPRW.2008.4562953>
- [16] Antonio Torralba, Bryan C. Russell, and Jenny Yuen. 2010. LabelMe: Online image annotation and applications. *Proc. IEEE* 98, 8 (2010), 1467–1484. <https://doi.org/10.1109/JPROC.2010.2050290>
- [17] Sirion Vittayakorn and James Hays. 2011. Quality Assessment for Crowdsourced Object Annotations. *Proceedings of the British Machine Vision Conference* (2011), 109.1–109.11. <https://doi.org/10.5244/C.25.109>
- [18] Peter Welinder, Steve Branson, Serge Belongie, and Pietro Perona. 2010. The Multidimensional Wisdom of Crowds. *NIPS (Conference on Neural Information Processing Systems)* 6 (2010), 1–9. <https://doi.org/10.1.1.231.1538>
- [19] Kota Yamaguchi. 2012. Parsing Clothing in Fashion Photographs. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (CVPR '12)*. IEEE Computer Society, Washington, DC, USA, 3570–3577. <http://dl.acm.org/citation.cfm?id=2354409.2355126>