

---

# A Closed-Loop System for Improving Annotation Quality and Efficiency

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present a general system and approach to improve the quality and efficiency  
2 of interactive annotation. A specific use case based on instance segmentation  
3 of vehicles for autonomous driving is used as an illustration. Via incremental  
4 AB testing and a custom analytics pipeline, we show how to optimize human-  
5 ML interaction to systematically improve annotation efficiency, and address the  
6 shortcomings of ML models.

## 7 1 Introduction

8 Pre-annotation and interactive (few-click) annotation are now commonly used in labelling tools and  
9 platforms to accelerate annotation and therefore reduce the cost of producing labelled datasets to train  
10 Deep Neural Networks (DNNs). In this paper, we present a general analytics-driven methodology that  
11 we have found useful to optimize our Machine Learning (ML) approach and incrementally improve  
12 its integration with human annotators.

## 13 2 Interactive annotation for segmentation

14 The annotation task considered in this study consists in labelling images to train computer vision  
15 models toward autonomous driving. Human annotators are required to draw accurate polygons around  
16 vehicles on those images. This is, however, a time-consuming and costly process.

17 Given the cost involved, multiple approaches have been suggested for machine-assisted instance  
18 segmentation. These usually consist of a deep learning-based segmentation of the object(s) integrated  
19 into a human-in-the-loop system. The human can interact with the system by correcting the model  
20 output, initializing the model with one or several clicks, or a combination of those steps. Examples of  
21 such systems include Polygon-RNN++ [1], DELSE [10], DEXTR [8], and CurveGCN [7].

22 We build upon an approach presented in [2] in order to illustrate our system. In a nutshell, this  
23 approach uses a few-click segmentation model (custom DEXTR) and a post-processing procedure that  
24 converts the produced raster mask into a sparse polygon. Figure 1 shows an example of human-edited  
25 ML output, compared with a fully manual approach, on the SYNTHIA-AL synthetic automotive  
26 dataset [12].

27 The goal is to conduct a series of AB tests in order to optimize the adequacy of the system's output  
28 for human editing, thus making the annotators as efficient as possible.

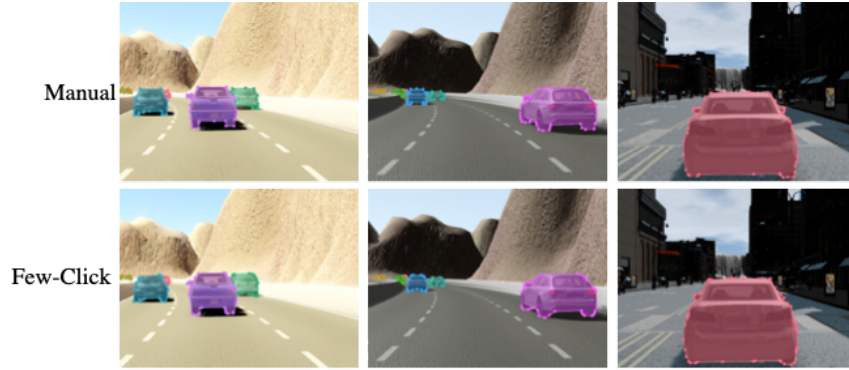


Figure 1: Comparing polygonal annotations of vehicles: (top) polygons drawn from scratch; (bottom) polygons manually edited from the output of the custom DEXTR.

### 29 3 Efficiency Metrics: a deep look into annotation analytics

#### 30 3.1 Definition of the metrics

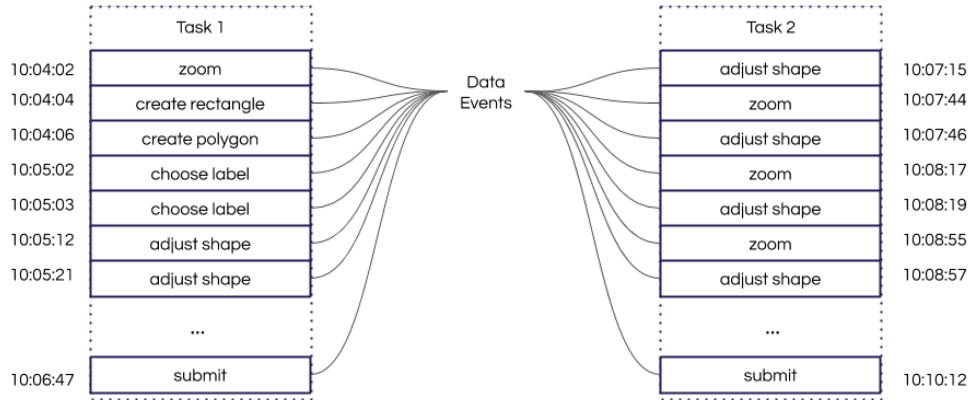
31 The performance of human annotators is dependent upon the quality of the interactions between the  
 32 person and the labelling platform. As new features (such as ML-based automation) are introduced, it  
 33 is of utmost importance to rigorously measure the impact of those features on annotation quality and  
 34 efficiency. In order to achieve this, we have introduced a series of metrics that help us gather detailed  
 35 insights into how annotation tasks are performed.

36 We have created a mapping between user tasks performed throughout the day, both on and off the  
 37 annotation tool. We assign each type of activity that annotators do, on a daily basis, to a category and  
 38 measure the amount of time devoted to them. The detailed categories are:

- 39 • **Training and feedback:** Before annotators begin work on a project, they go through a  
 40 training program. Furthermore, once the project is in production, they receive constant  
 41 feedback to improve their annotation work, based on the tasks performed so far.
- 42 • **Scanning:** Depending on the asset being labelled, finding the objects of interest can take a  
 43 significant amount of time. This can include moving around a high-resolution image, with  
 44 constant zooming in and out of regions, or navigating frames back and forth on a video.
- 45 • **Annotation creation:** This includes creating polygons that enclose objects of interest, in  
 46 the case of vector annotation, as well a marking pixels in raster-based tasks.
- 47 • **Label selection:** Depending on the complexity of the task, finding the correct label to  
 48 assign is not necessarily a quick operation. Because of that, we create a separate category  
 49 to account for the amount of time spent searching for the right label, out of a catalog of  
 50 available labels.
- 51 • **Annotation adjustment:** The process of initial polygon creation regularly requires anno-  
 52 tators to adjust the shapes drawn. We track how much time is devoted to this type of task.  
 53 In addition to this use case (i.e., create a shape and then adjust edges), there exists the  
 54 possibility of asking annotators to redo a task, when the quality criteria are not met. It is  
 55 expected that some of this feedback can be addressed by providing further adjustments to  
 56 shapes.
- 57 • **Validation:** Annotators can perform a final check on their annotations, before they submit  
 58 the work for a task to our platform. We want to assess how much of the time goes into this  
 59 type of check.

60 For the case study laid out here we focus on two of the metrics, primarily: annotation creation and  
 61 adjustment. For this, we have created a specific catalog of actions that can be performed on the user  
 62 interface, which affect each of the categories. The tool reports back those actions into our analytics  
 63 infrastructure, which allows us to create two data sources:

### A) Efficiency Metrics Detailed Log



### B) Annotator Aggregates

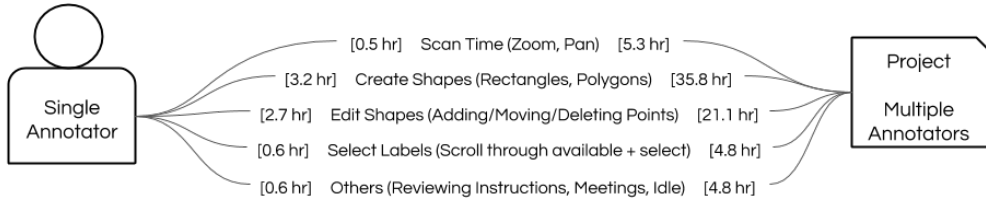
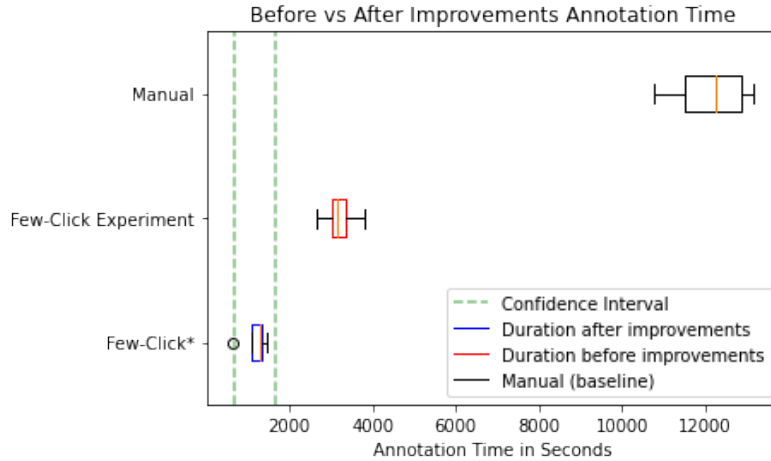


Figure 2: Efficiency Metrics data sources. **A)** The detailed log stores all events of interest processed by our internal data pipelines, which enables us to understand how each annotation task progressed. **B)** On top of the detailed data, we create daily aggregates per annotator, grouped by categories of interest, that are then fed into a statistical framework, for significance evaluation.

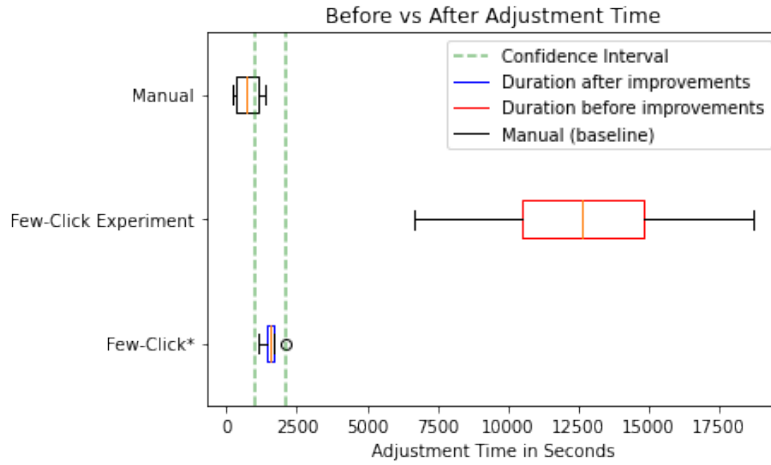
64 • **Efficiency metrics detailed log:** All user interface actions are stored sequentially for a  
 65 given task so that we can reconstruct the actions performed on the annotation tool, as  
 66 depicted in Figure 2.A. This is the finest granularity level in our analytics infrastructure and  
 67 is useful to validate patterns. For example, one agent may choose to draw all shapes first  
 68 and then do a second pass to select labels. Other agents may choose to select the label with  
 69 each shape. Some other tasks may center around adjusting shapes after they are returned to  
 70 annotators. That level of detail is only available through this source.

71 • **Annotator aggregates:** In order to power A/B testing, we need to define a granularity  
 72 level to create *populations* that can be compared against each other. Figure 2.B shows two  
 73 aggregation levels that we have found useful in our analyses. For instance, we may want  
 74 to know if we reduced the daily number of hours that annotators spent creating polygons,  
 75 due to some automation improvement. In that case we may define as *population individuals*  
 76 the daily metrics gathered per annotator and then do either before/after or side by side  
 77 comparisons between different groups. Creating this data source requires transforming  
 78 individual actions from the detailed log into accurate time aggregates. We can then query  
 79 population metrics based on the granularity of interest and assess the significance of the  
 80 changes.

81 The combined usage of these two data sources becomes a critical asset in the way we can iterate  
 82 quickly in our feature development. All of the metrics generated are updated on a near real-time basis,  
 83 which allows members in different areas of the organization, from data scientists to project managers,  
 84 to analyze the evolution of their projects and determine if the alleged improvements are verifiable.



(a) **Annotation Time**



(b) **Adjustment Time**

Figure 3: Comparison of average annotation time and adjustment time per object instance (in seconds) across all human annotators, for the baseline (*Manual*), the initial experiment (*Few-click*), and the experiment after improving the guidelines and the model output (*Few-click\**).

85 **3.2 A/B testing capabilities**

86 In an effort to conduct evidence-based process and product development, we follow best practices  
 87 in the industry [3, 5, 6, 9, 11] and adopt an experimentation approach that informs decision making.  
 88 To this end, we have developed a flexible testing infrastructure that can consume data from multiple  
 89 internal processes and is available to virtually anyone within the organization who wishes to setup  
 90 and keep track of an experiment [5].

91 The adoption or rejection of software, UI/UX, and/or procedure changes that impact our annotation  
 92 pipeline are evaluated through our AB testing framework, which measures the statistical impact  
 93 of these changes on our efficiency metrics such as drawing and adjustment times. Significance of  
 94 observed differences in a given efficiency metric is evaluated using rigorous statistical tests (*e.g.*  
 95 two-tailed hypothesis test). For the duration of an AB experiment, treatments are randomly assigned  
 96 and fixed to annotators from within our annotation tool [4].

## 97 4 Experimental setup

98 In the current experiment, we wish to compare the newly introduced interactive annotation system to  
99 a manual annotation baseline. Our setup is similar to the one described in [2]. In each group (A and  
100 B populations), each annotator is asked to label Motor Vehicles in 80 distinct images from a current  
101 customer project related to autonomous vehicles. Annotators are initially given clear guidelines along  
102 with a practice set, which is not counted towards in the metrics.

103 For this case study, efficiency metrics are gathered and stored for all annotators working on both  
104 treatments, and hypothesis tests evaluate two efficiency metrics of interest: initial polygon drawing  
105 time (from hereon dubbed annotation time) and adjustment time. We hypothesize that the new  
106 few-click tool will reduce both annotation and adjustment times.

## 107 5 Results and discussion

108 The first experiment did not yield the expected improvement (see *Manual* and *Few-click Experiment*  
109 in Figure 3). We confirmed our hypothesis concerning annotation time which saw a significant  
110 decrease, but had to reject the one pertaining to adjustment time, which was greatly increased. In  
111 order to understand what was happening, we selected images containing instances with very high  
112 adjustment time. After close inspection, we concluded the following:

- 113 • Our annotators entered sub-optimal extreme few-clicks which once fed into our custom  
114 DEXTR model along with corresponding images yielded low-quality polygons.
- 115 • Polygons on smaller objects had a much lower IoU on average when compared to established  
116 ground truths. We discovered that the source of this issue was in the parameters of the  
117 raster-to-polygon mask conversion algorithm.

118 We addressed the first issue by enhancing the training of our annotators and the labelling guidelines,  
119 and the second one by adjusting the parameters of the raster-to-polygon algorithm to maintain the  
120 quality of the output shapes at all scales. We then ran a new experiment whose results are shown  
121 as *Few-click\** on Figure 3. Even though adjustment time remained slightly higher when using the  
122 few-click tool, the total labelling time decreased by 78%, a vast improvement over the initially  
123 observed 21% increase. Furthermore, the quality of the delivered polygons was maintained with an  
124 average IoU around 95% when compared to established ground truths.

## 125 6 Conclusion

126 In this paper, we described a general methodology to incrementally improve a human-in-the-loop  
127 system based on annotation analytics and AB testing, with a specific use case in image labelling  
128 for autonomous driving. We have shown that detailed analytics can help pinpoint shortcomings in  
129 annotation guidelines and improve the quality of the output of ML models.

## 130 Broader Impact

131 Our work is concerned with improving the efficiency of human-in-the-loop annotation systems. While  
132 ML is bound to automate certain tasks (including within annotation itself), we hope that optimizing  
133 human-ML interaction in annotation systems will direct human labour towards where it is needed,  
134 thus preserving the need for these jobs while removing the more tedious aspects of it.

## 135 References

- 136 [1] David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. Efficient annotation of segmentation  
137 datasets with polygon-rnn++. In *CVPR*, 2018.
- 138 [2] M. Bertrand, F. Ratle, and L. Juillard. Human-centric efficiency improvements in image  
139 annotation for autonomous driving. In *ICML Workshop on Human In the Loop Learning (HILL)*,  
140 2020.

- 141 [3] Nikos Diamantopoulos, Jeffrey Wong, David Issa Mattos, Ilias Gerostathopoulos, Matthew  
142 Wardrop, Tobias Mao, and Colin McFarland. Engineering for a science-centric experimenta-  
143 tion platform. In *Proceedings of the ACM/IEEE 42nd International Conference on Software  
144 Engineering: Software Engineering in Practice, ICSE-SEIP '20*, page 191–200, New York, NY,  
145 USA, 2020. Association for Computing Machinery.
- 146 [4] Pavel Dmitriev, Brian Frasca, Somit Gupta, Ron Kohavi, and Garnet Vaz. Pitfalls of long-term  
147 online controlled experiments. In *Proceedings - 2016 IEEE International Conference on Big  
148 Data, Big Data 2016*, pages 1367–1376. Institute of Electrical and Electronics Engineers Inc.,  
149 2016.
- 150 [5] Raphael Lopez Kaufman, Jegar Pitchforth, and Lukas Vermeer. Democratizing online controlled  
151 experiments at booking.com. *CoRR*, abs/1710.08217, 2017.
- 152 [6] Ron Kohavi, Alex Deng, Brian Frasca, Toby Walker, Ya Xu, and Nils Pohlmann. Online  
153 controlled experiments at large scale. In *Proceedings of the ACM SIGKDD International  
154 Conference on Knowledge Discovery and Data Mining*, volume Part F128815, pages 1168–  
155 1176. Association for Computing Machinery, aug 2013.
- 156 [7] Huan Ling, Jun Gao, Amlan Kar, Wenzheng Chen, and Sanja Fidler. Fast interactive object  
157 annotation with curve-gen. In *CVPR*, 2019.
- 158 [8] K.-K. Maninis, S. Caelles, J. Pont-Tuset, and L. Van Gool. Deep extreme cut: From extreme  
159 points to object segmentation. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- 160 [9] Diane Tang, Ashish Agarwal, Deirdre O’Brien, and Mike Meyer. Overlapping experiment  
161 infrastructure: More, better, faster experimentation. In *Proceedings of the ACM SIGKDD  
162 International Conference on Knowledge Discovery and Data Mining*, pages 17–26, 2010.
- 163 [10] Zian Wang, David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. Object instance annotation  
164 with deep extreme level set evolution. In *CVPR*, 2019.
- 165 [11] Huizhi Xie and Juliette Aurisset. Improving the sensitivity of online controlled experiments:  
166 Case studies at Netflix. In *Proceedings of the ACM SIGKDD International Conference on Knowl-  
167 edge Discovery and Data Mining*, volume 13-17-August-2016, pages 645–654. Association for  
168 Computing Machinery, aug 2016.
- 169 [12] Javad Zolfaghari Bengar, Abel Gonzalez-Garcia, Gabriel Villalonga, Bogdan Raducanu,  
170 Hamed H Aghdam, Mikhail Mozerov, Antonio M Lopez, and Joost van de Weijer. Temporal  
171 coherence for active learning in videos. *arXiv preprint arXiv:1908.11757*, 2019.