
Data Attribution for Segmentation Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The quality of segmentation models is driven by their training datasets labeled with
2 detailed segmentation masks. How does the composition of such a training dataset
3 contribute to the performance of the resulting segmentation model? In this work,
4 we take a step towards attaining such an understanding by applying the lens of data
5 attribution to it. To this end, We first identify specific behaviors of these models
6 to attribute, and then provide a method for computing such attributions efficiently.
7 We validate the resulting attributions, and leverage them to both identify harmful
8 labeling errors and curate a 50% subset of the MS COCO training dataset that leads
9 to a $2.79\% \pm 0.49\%$ increase in mIOU over the full dataset.

10 1 Introduction

11 Semantic segmentation is a fundamental task in computer vision. Indeed, segmentation models
12 classify individual pixels in an image, enabling a detailed understanding of complex scenes with
13 diverse applications such as autonomous driving [1, 2], medical image analysis [3, 4], and agricultural
14 field monitoring [5]. However, the quality of segmentation models is impacted by the difficulty of
15 collecting high-quality training data [6].

16 In particular, the process of manually creating the whole segmentation masks for an individual
17 image is far more expensive than assigning a single label to it (as is done in the classification
18 setting). Because of this complexity of creating pixel-wise annotations, the resulting masks often
19 contain labeling errors [7, 8]. To what extent are these label errors detrimental to the performance of
20 segmentation models though? And can we identify and remove problematic training examples in
21 order to curate better training datasets?

22 More broadly, we might wonder about the relationship between training data and the resulting
23 behavior of segmentation models. For example, we might want to understand the possible pitfalls of
24 training on synthetic data [9, 10], an increasingly popular alternative to manually annotated data.

25 **Our Contributions** To answer such questions, in this work we study *data attribution* for segmen-
26 tation models—that is, the task of tracing such models’ behavior to individual training examples.

27 Performing data attribution in the segmentation setting requires addressing a number of complications.
28 First, a segmentation model might aim to segment multiple distinct objects that can be of different
29 classes, making it less clear how to choose a specific target of attribution (i.e., *what* to attribute). To
30 this end, we begin by identifying specific behaviors (i.e., parts of the output) of segmentation models
31 that we wish to attribute to the training dataset. Additionally, since the output of segmentation models
32 is high-dimensional, computing gradients of this model output with respect to training inputs (a key
33 step in most data attribution methods) might be difficult. However, by leveraging recent work in
34 data attribution within classification settings [11], we provide an efficient method for computing data
35 attributions for segmentation models, and validate the faithfulness of the resulting attributions.

36 Finally, we demonstrate that the attributions identified by our method surface harmful labeling errors
37 in segmentation datasets, and leverage these attributions to curate a 50% subset of the MS COCO

What impact does *each* training example have on the model’s ability to segment...

...the class cats?



...that specific shoe?



...background from foreground?



Figure 1: **What do we want to attribute?** In the segmentation setting, there are many possible targets for attribution. Here, we visualize three possible targets of interest: a model’s predictions for a particular class (left), a specific object (middle), or separation of background from foreground (right).

38 training dataset that achieves a $2.79\% \pm 0.49\%$ increase in mIOU over the full training dataset.
39 Related works are in Appendix A and a conclusion is in Appendix B.

40 2 Preliminaries

41 **Data Attribution** Data attribution has been closely studied for classification models, and within
42 this setting, previous works have employed data attribution for a number of tasks that can provide
43 great value in the segmentation setting. For example, mislabeled or otherwise poisoned training
44 examples can be identified by examining negatively influential training examples [12–17]. More
45 broadly, evaluating the overall impact of individual training examples can be useful for curating
46 training datasets [13, 14, 18]. Data attribution can also be useful for debugging model behavior [19]
47 and comparing learning algorithms [20]. Below, we formalize the task of data attribution.

48 Consider a learning algorithm \mathcal{A} (e.g., the training process for a neural network) paired with an
49 n -element training dataset $S \in \mathcal{Z}^n$ from input space \mathcal{Z} . Broadly, data attribution aims at identifying
50 how the behavior of models trained with algorithm \mathcal{A} is impacted by each training point $z_i \in S$. In
51 particular, given some held-out example $z \in \mathcal{Z}$, we can quantify the behavior of the model on this
52 example via a *model output function* $f(z, \theta(S)) : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathbb{R}$, where $\theta(S) \in \mathbb{R}^d$ denotes the model
53 parameters resulting from running the algorithm \mathcal{A} on the dataset S .

54 A valuable primitive that captures many of the underlying goals of data attribution is *datamodel-*
55 *ing* [19]: the task of predicting the model output function $f(z, \theta(S'))$ that results from running
56 algorithm \mathcal{A} on arbitrary subsets $S' \subset S$. In fact, despite the complex dynamics of modern non-
57 convex neural networks, prior works [19, 15, 11] have demonstrated that this prediction task can be
58 well approximated by a *linear* mapping $\{0, 1\}^{|S'|} \rightarrow \mathcal{R}$. So, following Park et al. [11], we can formal-
59 ize a *data attribution method* as a function $\tau : \mathcal{Z} \times \mathcal{Z}^n \rightarrow \mathbb{R}^n$ that assigns a score $\tau(z, S)_i \in \mathbb{R}$ to
60 each training example $z_i \in S$, indicating the change in $f(z, \theta(S))$ induced by removing z_i from S . In
61 particular, we can interpret these scores as weights of a *linear datamodel*: $f(z, \theta(S')) \approx \tau(z, S)^T \mathbf{1}_{S'}$,
62 where $\mathbf{1}_{S'}$ is the indicator vector of S' in S .

63 **Segmentation Models** We focus on the task of *semantic segmentation*, in which each pixel in an
64 image is assigned a class. Each training example $z \in \mathcal{Z}$ consists of an image x of size $H \times W \times N_{ch}$,
65 where N_{ch} is the number of channels, and a label y of size $H \times W$. Each pixel x_{ij} has a label
66 $y_{ij} \in \{0, 1, \dots, N_{cls}\}$, where N_{cls} is the number of classes (and 0 represents the background class).

67 3 Attributing Segmentation Models

68 In order to perform data attribution for segmentation models, we first need to identify a specific target
69 of attribution (and pair it with a model output function that formalizes this attribution target). So,
70 we begin this section by identifying two such targets, towards which we will focus our attributions.
71 Then, we introduce a method for computing such attributions efficiently by leveraging TRAK [11], a
72 recent work in data attribution for supervised models.

73 3.1 What do we want to attribute?

74 Classification models aim to predict a single discrete value (the class label). So, in this setting, the
75 target of attribution is generally the correctness of the model in predicting this label. To represent

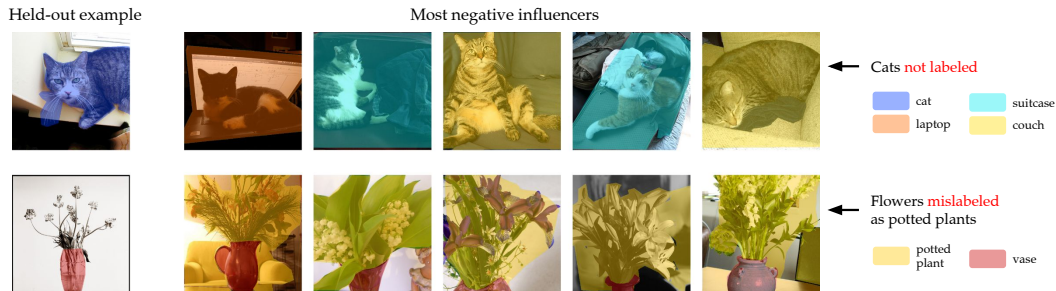


Figure 2: **Investigating negative influencers uncovers labeling errors.** For two MS COCO validation examples, we display the five most negatively influential training images according to our attribution method. Note that these negative influencers reflect consistent mislabeling errors: in the top row, each training example includes a cat that is not labeled, and in the bottom row, the flowers in each image are mislabeled as potted plants. See Appendix D for further examples.

76 correctness, a common choice for the model output function is the *correct-class margin* [19, 11], that
 77 is, the difference between the logit for the correct class and the highest incorrect logit.

78 For segmentation models, however, the question of *what* to attribute is less clear. In particular,
 79 segmentation models output a high-dimensional predicted mask with a separate prediction for each
 80 pixel of the input image. This output mask for a single image often predicts the presence of *multiple*
 81 objects, which may be of different classes. Thus, given a specific image of interest, there are a number
 82 of possible questions we might want to ask about the behavior of segmentation models on this image.

83 For instance, consider the image (and corresponding ground-truth masks) in Figure 1. Beyond
 84 identifying the training examples that impact the model’s overall ability to segment the objects in
 85 this image, there might be a number of more fine-grained questions we might want to ask in order
 86 to *isolate* attributions to specific “subtasks” of the segmentation task. For example, we might be
 87 interested in identifying the training examples that impact the model’s ability to segment specifically
 88 the cats in this image, or the model’s ability to generally separate the background from foreground
 89 objects. Clearly, there are many possible such subtasks of interest. In this work, we limit our focus to
 90 the following two possible targets of attribution:

- 91 1. **Full-image Attribution:** As a natural initial choice, we can directly attribute with respect
 92 to the full segmentation of a given image of interest. Specifically, we attribute the average
 93 of the per-pixel class predictions (each of which is its own classification task).
- 94 2. **Class-Specific Attribution:** As one avenue of identifying more fine-grained attributions,
 95 we additionally consider attributing the model’s ability to segment a specific class. Specif-
 96 ically, given a class c we attribute the per-pixel binary classification task of identifying
 97 whether each pixel belongs to class c .

98 3.2 Adapting TRAK to Segmentation Models

99 One common strategy for assigning attribution scores $\tau(z)_i$ for a given input z is to compute the
 100 *leave-one-out* influence of removing that particular example on the model output function:

$$\tau(z)_i = f(z, \theta(S)) - f(z, \theta(S \setminus \{z_i\})).$$

101 In the linear regression setting, we can compute this influence directly, as there exists a closed form
 102 solution for the parameters $\theta(S')$ given a training subset $S' \subset S$. However, estimating the resulting
 103 model parameters when leaving out an example is difficult in the non-convex deep learning setting.

104 To overcome this difficulty, TRAK [11] first *linearize* the model and then apply classical methods for
 105 data attribution in the linear setting [21]. In particular, recent work in studying the neural tangent
 106 kernel (NTK) has shown that linearizing neural networks with their first-order Taylor expansion can
 107 closely approximate training dynamics [22–24]. This expansion allows us to view the model output
 108 function $f(z, \theta(S))$ as a linear model acting on inputs $\nabla_{\theta} f(z, \theta(S))$ (see Section 3.2 of Park et al.
 109 [11] for a more detailed explanation). Now, applying TRAK to a given attribution task requires two
 110 steps: (1) computing gradients of the training examples with respect to the training loss, and (2)



Figure 3: **Examples of positive influencers.** For two MS COCO validation examples, we show the most positive influencers identified by our method, along with the most similar training examples according to CLIP similarity.



Figure 4: **Per-class positive influencers.** An example MS COCO validation example that includes segmentations for the *person*, *car*, and *traffic light* classes, as well as the top two most positively influential training images according to class-specific TRAK scores for each of these classes.

121 computing gradients of each held-out example of interest with respect to the model output functions
 122 (see Section 3.4 of Park et al. [11] for more details).

123 At first glance, it may appear that since the output of a segmentation model is high-dimensional
 124 (with a multi-class prediction for each pixel in the image), computing the gradient with respect to
 125 each training example (the first step) might require tracking a large number of independent gradients.
 126 However, note that each training example impacts the final model parameters *only* through its
 127 contribution to the gradient of the training loss. Thus, we only need to compute the gradient of each
 128 training example with respect to this loss. Finally, to compute the second step, below we define a
 129 model output function for each of the two attribution targets introduced in Section 3.1.

130 **Full-image Attribution** In the classification setting, a common choice of model output function
 131 (as we introduced in Section 3.1), is the *correct-class margin* between the correct-class logit and
 132 the largest incorrect logit [19]. However, prior theoretical work [25] motivates the use of the model
 133 output function

$$f(z; \theta) = \log \left(\frac{p(z; \theta)}{1 - p(z; \theta)} \right),$$

134 where $p(z; \theta)$ is the correct-class probability; this model output can be viewed as a “soft” version of
 135 the correct-class margin.

136 Since the training loss takes the average of individual pixel-wise cross-entropy losses over all pixels
 137 in an image, a natural choice of model output function for the semantic segmentation setting is to
 138 adapt the above model output in the same manner. Specifically, in the full-image attribution setting,
 139 we calculate the above model output function for each pixel-wise classification problem, then average
 140 over pixels as follows:

$$f(z; \theta(S)) = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \log \left(\frac{p(z_{ij}; \theta)}{1 - p(z_{ij}; \theta)} \right), \quad (1)$$

131 where $p(z_{ij}; \theta)$ is the correct-class probability predicted for pixel x_{ij} .

132 **Class-Specific Attribution** To attribute a model’s segmentation predictions for a specific class c ,
 133 we treat the segmentation task as a binary classification task with respect to c . That is, we adapt the

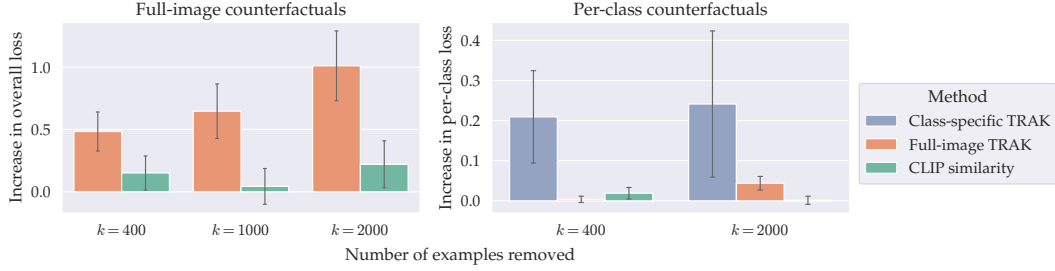


Figure 5: **Counterfactual experiments.** We quantify the counterfactual impact of removing the highest scoring training examples and retraining according to TRAK scores targeted at either full-image or per-class attribution, or according to CLIP similarity. (Left) We apply this intervention to 15 validation examples and plot the average increase in cross-entropy loss of the segmentation task. (Right) We apply this intervention to 15 pairs of (validation example, class) and plot the average increase in cross-entropy loss over the pixel-wise binary classification task of predicting the presence of that class, with per-class attributions targeted at the this class. Error bars represent standard error.

134 above model output such that all classes that other than c are treated as the same "not- c " class. Our
 135 model output for an example $z = (x, y)$ then becomes:

$$f(z; \theta(S)) = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \left(\mathbf{1}\{y_{ij} = c\} \log \left(\frac{p_c(x_{ij}; \theta)}{1 - p_c(x_{ij}; \theta)} \right) + \mathbf{1}\{y_{ij} \neq c\} \log \left(\frac{1 - p_c(x_{ij}; \theta)}{p_c(x_{ij}; \theta)} \right) \right), \quad (2)$$

136 where $p_c(x_{ij}; \theta)$ refers to the predicted probability that pixel x_{ij} belongs to class c .

137 4 Experiments

138 We now evaluate our attribution method on segmentation models trained on the MS COCO
 139 dataset [26]. After visually inspecting the computed attributions, we validate their counterfactual
 140 significance, and then demonstrate their value for curating training datasets.

141 **Experimental Setup** We use DeepLabV3 [27] segmentation models trained on segmentations
 142 from the MS COCO dataset [26]. When validating our attribution scores, we use image similarity
 143 according to a pretrained CLIP [28] ViT-L/14 [29] model as a baseline.

144 4.1 Visually Inspecting Our Attributions

145 In Figures 2 and 3, we visualize the most negative and most positive influencers, respectively, for
 146 images from the MS COCO validation set. We find that positive influencers identify semantically
 147 similar (or even nearly-duplicated) training examples, while negative influencers often surface
 148 consistent labeling errors. In Figure 4, we additionally visualize an example of top positive class-
 149 specific influencers. See Appendix D for further examples of attributions identified by our method.

150 4.2 Validating Our Attributions

151 To validate our full-image attributions, we measure the counterfactual impact of the following
 152 intervention: removing the most positive influencers for a given held-out example, retraining, and
 153 then measuring the change in loss between the original and new models. If the attributions are
 154 meaningful, we would expect this intervention to cause a significant increase in loss. We repeat this
 155 process over 25 randomly selected validation examples and compare to the baseline of removing the
 156 most visually similar training examples according to CLIP (see Figure 5, left). Our results suggest
 157 that that the full-image TRAK scores indeed identify influencers that are counterfactually significant.

158 To validate our class-specific attributions, we repeat the above intervention with top influencers
 159 identified by either the class-specific attributions for a given class within an image or full-image

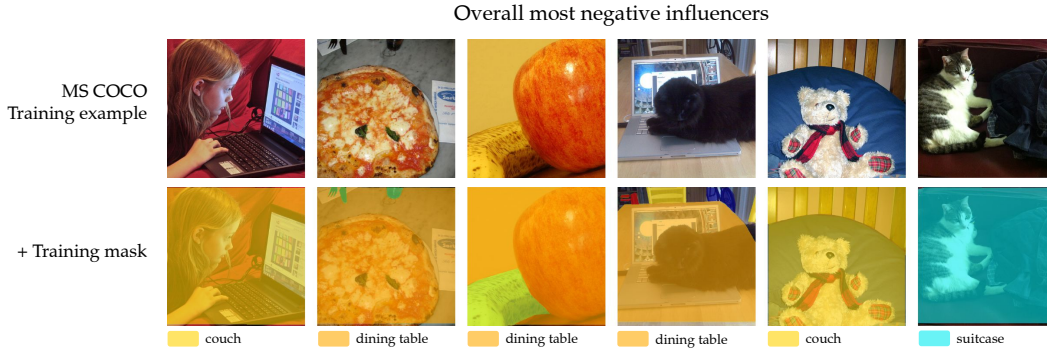


Figure 6: **Overall Most Negative Influencers.** Six of the twenty overall most negatively influential training examples across the full MS COCO validation set, according to our attribution method. Note that *person*, *cat*, *pizza*, *laptop*, *apple*, and *teddy bear* are all MS COCO classes, but labels for these classes are missing in the above training masks.

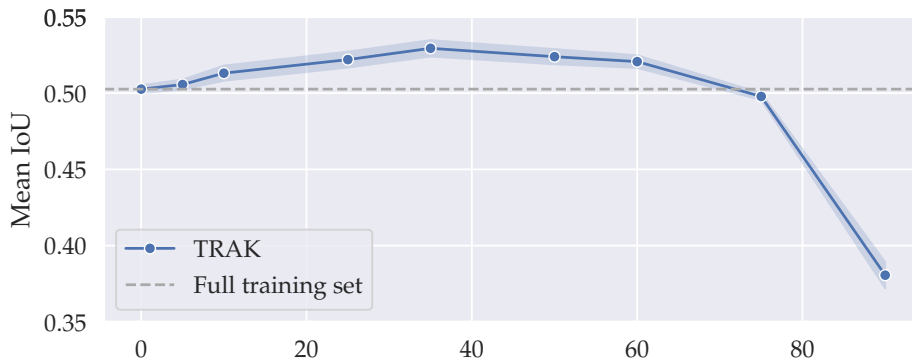


Figure 7: **Dataset curation.** Mean IoU of models trained on MS COCO after removing the most negatively influential examples of the training set according to our attribution scores. The shaded region represents standard error.

160 attributions, and measure the impact of this intervention on the loss of the binary classification task
 161 with respect the corresponding class. We find that when targeting attributions towards a specific class,
 162 the impact of the intervention is greater than that of full-image attributions (see Figure 5, right).

163 4.3 Curating Datasets for Segmentation

164 In Section 4.1, we saw that the most negative influencers often surface training examples with labeling
 165 errors. Intuitively, such training examples should be detrimental to the segmentation model. By
 166 removing such problematic examples from the training dataset, can training on the resulting *curated*
 167 dataset improve the quality of the resulting model?

168 As proxy for each individual training image’s effect on overall model’s performance, we average its
 169 attribution scores across the entire validation set. In Figure 6, we visualize samples among the most
 170 negative “overall” influencers according to this proxy and find that the surfaced training examples
 171 include notable cases of mislabeling, often missing multiple objects within a single image. Now,
 172 for various values of N , we remove the bottom $N\%$ of training examples (according to averaged
 173 attribution scores) and re-train new models on the resulting smaller, curated dataset. To prevent
 174 leakage, we randomly split the validation set into two halves, using one half to calculate average
 175 influence scores and the other half to evaluate our re-trained models. We evaluate models via mean
 176 intersection-over-union (mIoU) of predicted masks versus ground-truth masks. We find that our
 177 curated datasets are able to outperform the full MS COCO training set (see Figure 7), and in particular,
 178 we achieve a $2.79\% \pm 0.49\%$ increase in mean IoU with a 50% subset of the training set.

References

- 179
180 [1] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proceedings of the*
181 *IEEE conference on computer vision and pattern recognition*, pages 3061–3070, 2015.
- 182 [2] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson,
183 Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding.
184 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223,
185 2016.
- 186 [3] Nathan Moon, Elizabeth Bullitt, Koen Van Leemput, and Guido Gerig. Automatic brain and tumor
187 segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2002: 5th*
188 *International Conference Tokyo, Japan, September 25–28, 2002 Proceedings, Part I 5*, pages 372–379.
189 Springer, 2002.
- 190 [4] Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a
191 self-configuring method for deep learning-based biomedical image segmentation. *Nature methods*, 18(2):
192 203–211, 2021.
- 193 [5] Furkat Safarov, Kuchkorov Temurbek, Djumanov Jamoljon, Ochilov Temur, Jean Chamberlain Chedjou,
194 Akmalbek Bobomirzaevich Abdusalomov, and Young-Im Cho. Improved agricultural field segmentation
195 in satellite imagery using tl-resunet architecture. *Sensors*, 22(24):9784, 2022.
- 196 [6] Șerban Vădineanu, Daniël Maria Pelt, Oleh Dzyubachyk, and Kees Joost Batenburg. An analysis of
197 the impact of annotation errors on the accuracy of deep learning for cell segmentation. In *International*
198 *Conference on Medical Imaging with Deep Learning*, pages 1251–1267. PMLR, 2022.
- 199 [7] Matthias Rottmann and Marco Reese. Automated detection of label errors in semantic segmentation datasets
200 via deep learning and uncertainty quantification. In *Proceedings of the IEEE/CVF Winter Conference on*
201 *Applications of Computer Vision*, pages 3214–3223, 2023.
- 202 [8] Vedang Lad and Jonas Mueller. Estimating label quality and errors in semantic segmentation data via any
203 model. *arXiv preprint arXiv:2307.05080*, 2023.
- 204 [9] Yuhua Chen, Wen Li, Xiaoran Chen, and Luc Van Gool. Learning semantic segmentation from syn-
205 thetic data: A geometrically guided input-output adaptation approach. In *Proceedings of the IEEE/CVF*
206 *conference on computer vision and pattern recognition*, pages 1841–1850, 2019.
- 207 [10] Virginia Fernandez, Walter Hugo Lopez Pinaya, Pedro Borges, Petru-Daniel Tudosiu, Mark S Graham,
208 Tom Vercauteren, and M Jorge Cardoso. Can segmentation models be trained with fully synthetically
209 generated data? In *International Workshop on Simulation and Synthesis in Medical Imaging*, pages 79–90.
210 Springer, 2022.
- 211 [11] Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak:
212 Attributing model behavior at scale. In *Arxiv preprint arXiv:2303.14186*, 2023.
- 213 [12] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In
214 *International Conference on Machine Learning (ICML)*, 2019.
- 215 [13] Rajiv Khanna, Been Kim, Joydeep Ghosh, and Sanmi Koyejo. Interpreting black box predictions using
216 fisher kernels. In *The 22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- 217 [14] Ruoxi Jia, Fan Wu, Xuehui Sun, Jiachen Xu, David Dao, Bhavya Kailkhura, Ce Zhang, Bo Li, and Dawn
218 Song. Scalability vs. utility: Do we have to sacrifice one for the other in data importance quantification?
219 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- 220 [15] Jinkun Lin, Anqi Zhang, Mathias Lecuyer, Jinyang Li, Aurojit Panda, and Siddhartha Sen. Measuring
221 the effect of training data on deep learning predictions via randomized experiments. *arXiv preprint*
222 *arXiv:2206.10013*, 2022.
- 223 [16] Zayd Hammoudeh and Daniel Lowd. Identifying a training-set attack’s target using renormalized influence
224 estimation. In *arXiv preprint arXiv:2201.10055*, 2022.
- 225 [17] Shuming Kong, Yanyan Shen, and Linpeng Huang. Resolving training biases via influence-based data
226 relabeling. In *International Conference on Learning Representations (ICLR)*, 2022.
- 227 [18] Zhuoming Liu, Hao Ding, Huaping Zhong, Weijia Li, Jifeng Dai, and Conghui He. Influence selection for
228 active learning. In *ICCV*, 2021.

- 229 [19] Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Datamodels:
230 Predicting predictions from training data. In *International Conference on Machine Learning (ICML)*, 2022.
- 231 [20] Harshay Shah, Sung Min Park, Andrew Ilyas, and Aleksander Madry. Modeldiff: A framework for
232 comparing learning algorithms. In *arXiv preprint arXiv:2211.12491*, 2022.
- 233 [21] Daryl Pregibon. Logistic regression diagnostics. In *The Annals of Statistics*, 1981.
- 234 [22] Philip M Long. Properties of the after kernel. In *arXiv preprint arXiv:2105.10585*, 2021.
- 235 [23] Alexander Wei, Wei Hu, and Jacob Steinhardt. More than a toy: Random matrix models predict how
236 real-world neural representations generalize. In *ICML*, 2022.
- 237 [24] Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based view of
238 language model fine-tuning. In *arXiv preprint arXiv:2210.05643*, 2022.
- 239 [25] Nikunj Saunshi, Arushi Gupta, Mark Braverman, and Sanjeev Arora. Understanding influence functions
240 and datamodels via harmonic analysis. 2023.
- 241 [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
242 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on
243 computer vision (ECCV)*, 2014.
- 244 [27] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab:
245 Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.
246 *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- 247 [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
248 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
249 natural language supervision. In *arXiv preprint arXiv:2103.00020*, 2021.
- 250 [29] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
251 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is
252 worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning
253 Representations (ICLR)*, 2021.
- 254 [30] R Dennis Cook and Sanford Weisberg. *Residuals and influence in regression*. New York: Chapman and
255 Hall, 1982.
- 256 [31] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: the
257 approach based on influence functions*, volume 196. John Wiley & Sons, 2011.
- 258 [32] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
259 *International Conference on Machine Learning*, 2017.
- 260 [33] Pang Wei Koh, Kai-Siang Ang, Hubert HK Teo, and Percy Liang. On the accuracy of influence functions
261 for measuring group effects. In *Neural Information Processing Systems (NeurIPS)*, 2019.
- 262 [34] Andrea Schioppa, Polina Zablotskaia, David Vilar, and Artem Sokolov. Scaling up influence functions. In
263 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8179–8186, 2022.
- 264 [35] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li,
265 Ce Zhang, Dawn Song, and Costas J. Spanos. Towards efficient data valuation based on the shapley value.
266 In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*,
267 2019.
- 268 [36] Bojan Karlaš, David Dao, Matteo Interlandi, Bo Li, Sebastian Schelter, Wentao Wu, and Ce Zhang.
269 Data debugging with shapley importance over end-to-end machine learning pipelines. *arXiv preprint
270 arXiv:2204.11131*, 2022.
- 271 [37] Chih-Kuan Yeh, Joon Sik Kim, Ian E. H. Yen, and Pradeep Ravikumar. Representer point selection for
272 explaining deep neural networks. In *Neural Information Processing Systems (NeurIPS)*, 2018.
- 273 [38] Garima Pruthi, Frederick Liu, Mukund Sundararajan, and Satyen Kale. Estimating training data influence
274 by tracing gradient descent. In *Neural Information Processing Systems (NeurIPS)*, 2020.

275 **A Related Work**

276 **Data Attribution** The task of data attribution has been studied closely, with early work focusing on the
277 linear regression setting [30, 21]. In the modern machine learning settings, there is a long line of work in
278 *influence functions* [31–34], as well as works that perform data attribution through Shapley values [12, 35, 36]
279 or other heuristic approaches [37, 38]. Recently, Park et al. [11] proposed a method that significantly improves
280 upon tradeoffs between computational efficiency and accuracy by leveraging the empirical kernel structure of
281 neural networks (see Section 3 for further discuss on this work).

282 **Detecting Labeling Errors in Segmentation Models.** Recent methods for detection labeling error [8, 7]
283 compare each training example’s mask label to predictions from a segmentation model train trained without
284 when excluding training example. For example, Lad and Mueller [8] propose an efficient strategy for label
285 error detection using a label quality score for segmentation models based on taking a soft minimum over
286 pixel-wise correct-class probabilities estimated by a trained segmentation model. Rather than directly processing
287 pixel-wise predictions, Rottmann and Reese [7] instead identify labeling errors by comparing the predicted
288 “connected components” (i.e., pixel-wise connected regions in an image with the same class label) from a trained
289 segmentation model to the connected components in the ground truth mask to identify mismatches.

290 **B Conclusion**

291 In this work, we study the task of data attribution for segmentation models. We first identify two specific
292 behaviors of these models that we hope to attribute: their overall ability to segment an image (i.e., full-image
293 attribution), and their ability to properly segment a specific class within an image (i.e., class-specific attribution).
294 We then provide a method for computing such attributions and instantiate this method on the MS COCO dataset.
295 We validate the counterfactual significance of our computed attributions, and leverage them to both uncover
296 labeling errors and curate subsets of the MS COCO training dataset that lead to increased accuracy over training
297 on the full dataset.

298 C Experimental details

299 **Dataset and models.** We train and test all models on segmentations from the MS COCO dataset (2017
300 version) [26]. We center-crop images to a size of 224×224 and use random horizontal flip as an augmentation.
301 For our models, we train DeepLabV3 [27] models for semantic segmentation for 120 epochs with the Adam
302 optimizer with initial learning rate 0.0001, batch size 64, weight decay 0.0005, and a learning rate scheduler that
303 reduces when loss plateaus. We use cross-entropy loss with no smoothing applied.

304 **TRAK hyperparameters.** For all uses of TRAK, we use a projection dimension of $k = 2048$. We find that
305 attributions using a projection dimension of $k = 8192$ did not appear much stronger visually, and yielded very
306 similar results when running preliminary data curation experiments. We also do not apply soft-thresholding to
307 improve the quality of TRAK’s attributions, though this is discussed in [11].

308 **D Additional Results**

309 We show additional top positive influencers for four randomly MS COCO validation examples in Figure 8. We
 310 also show the most negative influencers for two MS COCO validation examples in Figure 9, and highlight the
 311 labeling errors reflected in these influencers.

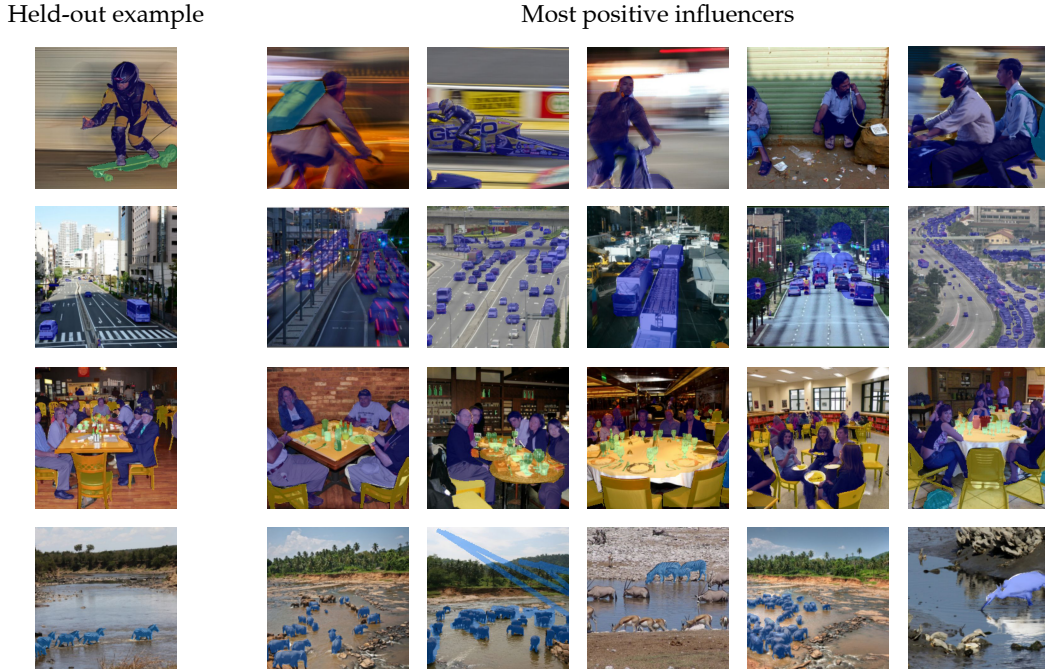


Figure 8: **More examples of top influencers.** For four MS COCO validation examples, we show the most positive influencers identified by TRAK.



Figure 9: **More examples of negative influencers.** For two MS COCO validation examples, we show the most negative (lowest-scoring) influencers identified to TRAK.