

A DATA-BASED PERSPECTIVE ON TRANSFER LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

It is commonly believed that in transfer learning including more pre-training data translates into better performance. However, recent evidence suggests that *removing* data from the source dataset can actually help too. In this work, we take a closer look at the role of the source dataset’s composition in transfer learning and present a framework for probing its impact on downstream performance. Our framework gives rise to new capabilities such as pinpointing transfer learning brittleness as well as detecting pathologies such as data-leakage and the presence of misleading examples in the source dataset. In particular, we demonstrate that removing detrimental datapoints identified by our framework indeed improves transfer learning performance from ImageNet on a variety of target tasks.

1 INTRODUCTION

Transfer learning enables us to adapt a model trained on a *source* dataset to perform better on a downstream *target* task. This technique is employed in a range of machine learning applications including radiology (Wang et al., 2017; Ke et al., 2021), autonomous driving (Kim & Park, 2017; Du et al., 2019), and satellite imagery analysis (Xie et al., 2016; Wang et al., 2019). Despite its successes, however, it is still not clear what the drivers of performance gains brought by transfer learning actually are.

So far, a dominant approach to studying these drivers focused on the role of the source *model*—i.e., the model trained on the source dataset. The corresponding works involve investigating the source model’s architecture (Ke et al., 2021), accuracy (Kornblith et al., 2019), adversarial vulnerability (Salman et al., 2020; Utrera et al., 2020), and training procedure (Jang et al., 2019; Kumar et al., 2022). This line of work makes it clear that the properties of the source model has a significant impact on transfer learning. There is some evidence, however, that the source *dataset* might play an important role as well (Huh et al., 2016; Ngiam et al., 2018; Kolesnikov et al., 2019). For example, several works have shown that while increasing the size of the source dataset generally boosts transfer learning performance, *removing* specific classes can help too (Huh et al., 2016; Ngiam et al., 2018; Kolesnikov et al., 2019). All of this motivates a natural question:

How can we pinpoint the exact impact of the source dataset in transfer learning?

Our Contributions. In this paper, we present a framework for measuring and analyzing the impact of the source dataset’s composition on transfer learning performance. To do this, our framework provides us with the ability to investigate the counterfactual impact on downstream predictions of including or excluding datapoints from the source dataset, drawing inspiration from classical supervised learning techniques such as influence functions (Cook & Weisberg, 1982; Koh & Liang, 2017; Feldman & Zhang, 2020) and datamodels (Ilyas et al., 2022). Using our framework, we can:

- Pinpoint what parts of the source dataset are most utilized by the downstream task.
- Automatically extract granular subpopulations in the target dataset through projection of the fine-grained labels of the source dataset.
- Surface pathologies such as source-target data leakage and mislabelled source datapoints.

We also demonstrate how our framework can be used to find detrimental subsets of ImageNet (Deng et al., 2009) that, when removed, give rise to better downstream performance on a variety of image classification tasks.

2 A DATA-BASED FRAMEWORK FOR STUDYING TRANSFER LEARNING

In order to pinpoint the role of the source dataset in transfer learning, we need to understand how the composition of that source dataset impacts the downstream model’s performance. To do so, we draw inspiration from supervised machine learning approaches that study the impact of the training data on the model’s subsequent predictions. In particular, these approaches capture this impact via studying (and approximating) the counterfactual effect of excluding certain training datapoints. This paradigm underlies a number of techniques, from influence functions (Cook & Weisberg, 1982; Koh & Liang, 2017; Feldman & Zhang, 2020), to datamodels (Ilyas et al., 2022), to data Shapley values (Kwon & Zou, 2021; Karlaš et al., 2022; Ghorbani & Zou, 2019).

Now, to adapt this paradigm to our setting, we study the counterfactual effect of excluding datapoints from the *source* dataset on the downstream, *target* task predictions. In our framework, we will focus on the inclusion or exclusion of entire *classes* in the source dataset, as opposed to individual examples¹. This is motivated by the fact that, intuitively, we expect these classes to be the ones that embody whole concepts and thus drive the formation of (transferred) features. We therefore anticipate the removal of entire classes to have a more measurable impact on the representation learned by the source model (and consequently on the downstream model’s predictions).

Once we have chosen to focus on removal of entire source classes, we can design counterfactual experiments to estimate their influences. A natural approach here, the *leave-one-out* method (Cook & Weisberg, 1982; Koh & Liang, 2017), would involve removing each individual class from the source dataset separately and then measuring the change in the downstream model’s predictions. However, in the transfer learning setting, we suspect that removing a single class from the source dataset won’t significantly change the downstream model’s performance. Thus, leave-one-out methodology may be able to capture meaningful influences only in rare cases. This is especially so as many common source datasets contain highly redundant classes. For example, ImageNet contains over 100 dog-breed classes. The removal of a single dog-breed class might thus have a negligible impact on transfer learning performance, but the removal of all of the dog classes might significantly change the features learned by the downstream model. For these reasons, we adapt the *subsampling* (Feldman & Zhang, 2020; Ilyas et al., 2022) approach, which revolves around removing a random collection of source classes at once.

Computing transfer influences. In the light of the above, our methodology for computing the influence of source classes on transfer learning performance involves training a large number of models with random subsets of the source classes removed, and fine-tuning these models on the target task. We then estimate the influence value of a source class \mathcal{C} on a target example t as the expected difference in the transfer model’s performance on example t when class \mathcal{C} was either included in or excluded from the source dataset:

$$\text{Infl}[\mathcal{C} \rightarrow t] = \mathbb{E}_S [f(t; S) \mid \mathcal{C} \subset S] - \mathbb{E}_S [f(t; S) \mid \mathcal{C} \not\subset S], \quad (1)$$

where $f(t; S)$ is the softmax output² of a model trained on a subset S of the source dataset. A positive influence value indicates that including the source class \mathcal{C} helps the model predict the target example t correctly. On the other hand, a negative influence value suggests that the source class \mathcal{C} actually hurts the model’s performance on the target example t . We outline the overall procedure in Algorithm 1, and defer a detailed description of our approach to Appendix A.

A note on computational costs. In order to compute transfer influences, we need to train a large number of source models, each on a fraction of the source dataset. Specifically, we pre-train 7,540 models on ImageNet, each on a randomly chosen 50% of the ImageNet dataset. This pre-training

¹In Section 4.3, we adapt our framework to calculate more granular influences of individual source examples too.

²We experiment with other outputs such as logits, margins, or correctness too. We discuss the corresponding results in Appendix B.

Algorithm 1 Estimation of source dataset class influences on transfer learning performance.

Require: Source dataset $\mathcal{S} = \cup_{k=1}^K \mathcal{C}_k$ (with K classes), a target dataset $\mathcal{T} = (t_1, t_2, \dots, t_n)$, training algorithm \mathcal{A} , subset ratio α , and number of models m

- 1: Sample m random subsets $S_1, S_2, \dots, S_m \subset \mathcal{S}$ of size $\alpha \cdot |\mathcal{S}|$:
- 2: **for** $i \in 1$ to m **do**
- 3: Train model f_i by running algorithm \mathcal{A} on S_i
- 4: **end for**
- 5: **for** $k \in 1$ to K **do**
- 6: **for** $j \in 1$ to n **do**
- 7: $\text{Infl}[C_k \rightarrow t_j] = \frac{\sum_{i=1}^m f_i(t_j; S_i) \mathbb{1}_{C_k \subset S_i}}{\sum_{i=1}^m \mathbb{1}_{C_k \subset S_i}} - \frac{\sum_{i=1}^m f_i(t_j; S_i) \mathbb{1}_{C_k \not\subset S_i}}{\sum_{i=1}^m \mathbb{1}_{C_k \not\subset S_i}}$
- 8: **end for**
- 9: **end for**
- 10: **return** $\text{Infl}[C_k \rightarrow t_j]$, for all $j \in [n], k \in [K]$

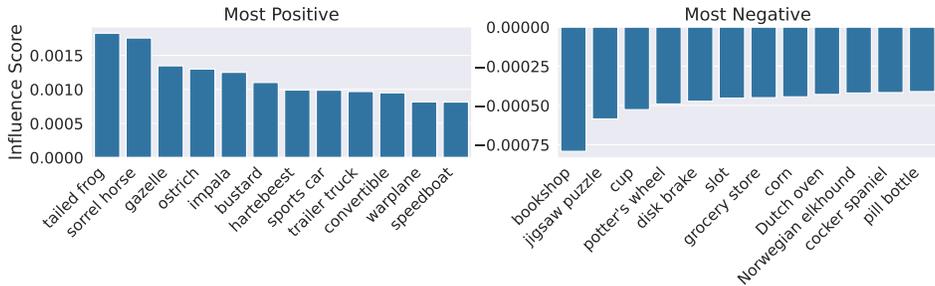


Figure 1: Most positive and negative ImageNet classes ordered based on their overall influence on the CIFAR-10 dataset. The top source classes (e.g., tailed frog and sorrel horse) turn out to be semantically relevant to the target classes (e.g., frog and horse).

step needs to be performed only once though: these same models can then be used to fine-tune on each new target task. Overall, the whole process (training the source models and fine-tuning on target datasets) takes less than 20 days using 8 V100 GPUs³.

Are so many models necessary? In Section A.5, we explore computing transfer influences with smaller numbers of models. While using the full number of models provides the best results, training a much smaller number of models (e.g., 1000 models, taking slightly over 2.5 days on 8 V100 GPUs) still provides meaningful transfer influences. Thus in practice, one can choose the number of source models based on noise tolerance and computational budget. Further convergence results can be found in Appendix A.5.

3 IDENTIFYING THE MOST INFLUENTIAL CLASSES OF THE SOURCE DATASET

In Section 2, we presented a framework for pinpointing the role of the source dataset in transfer learning by estimating the influence of each source class on the target model’s predictions. Using these influences, we can now take a look at the classes from the source dataset that have the largest positive or negative impact on the overall transfer learning performance. We focus our analysis on the fixed-weights transfer learning setting (and defer results for full model fine-tuning to Appendix E).

As one might expect, not all source classes have large influences. Figure 1 displays the most influential classes of ImageNet with CIFAR-10 as the target task. Notably, the most positively influential source classes turn out to be directly related to classes in the target task (e.g., the ImageNet label “tailed frog” is an instance of the CIFAR class “frog”). This trend holds across all of the target datasets and transfer learning settings we considered (see Appendix C). Interestingly, the source dataset also contains classes that are overall negatively influential for the target task (e.g., “bookshop” and “jigsaw

³Details are in Appendix A.

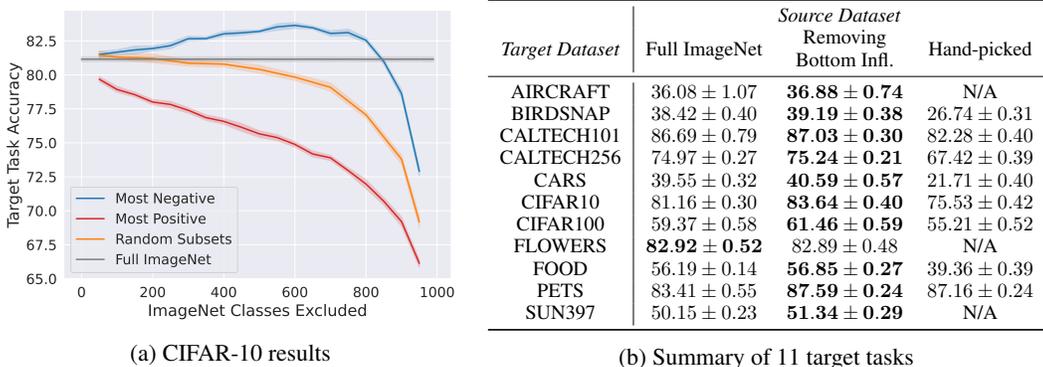


Figure 2: Target task accuracies after removing the K most positively or negatively influential ImageNet classes from the source dataset. Mean/std are reported over 10 runs. (a) Results with CIFAR-10 as the target task after removing different numbers of classes from the source dataset. We also include baselines of using the full ImageNet dataset and removing random classes. One can note that, by removing negatively influential source classes, we can obtain a test accuracy that is 2.5% larger than what using the entire ImageNet dataset would yield. Results for other target tasks can be found in Appendix C. (b) Peak performances when removing the most negatively influential source classes across a range of other target tasks. We compare against using the full ImageNet dataset or a relevant subset of classes (hand-picked, see Appendix A for details).

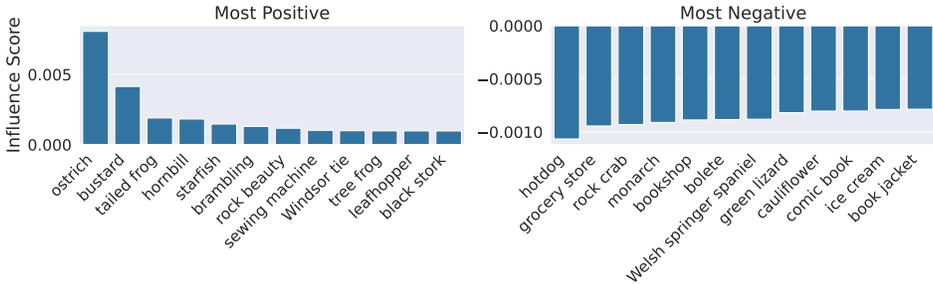


Figure 3: Most positive and negative influencing ImageNet classes for the CIFAR-10 class “bird”. These are calculated by averaging the influence of each source class over all bird examples. We find that the most positively influencing ImageNet classes (e.g., “ostrich” and “bustard”) are related to the CIFAR-10 class “bird”. See Appendix E for results on other CIFAR-10 classes.

puzzle” classes). (In Section 4, we will take a closer look at the factors that can cause a source class to be negatively influential for a target prediction.)

How important are the most influential source classes? We now remove each of the most influential classes from the source dataset to observe their actual impact on transfer learning performance (Figure 2a). As expected, removing the most positively influential classes severely degrades transfer learning performance as compared to removing random classes. This counterfactual experiment confirms that these classes are indeed important to the performance of transfer learning. On the other hand, removing the most negatively influential classes actually improves the overall transfer learning performance *beyond what using the entire ImageNet dataset provides* (see Figure 2b).

4 PROBING THE IMPACT OF THE SOURCE DATASET ON TRANSFER LEARNING

In Section 3, we developed a methodology for identifying source dataset classes that have the most impact on transfer learning performance. Now, we demonstrate how this methodology can be extended into a framework for probing and understanding transfer learning, including: (1) identifying granular target subpopulations that correspond to source classes, (2) debugging transfer learning failures, and

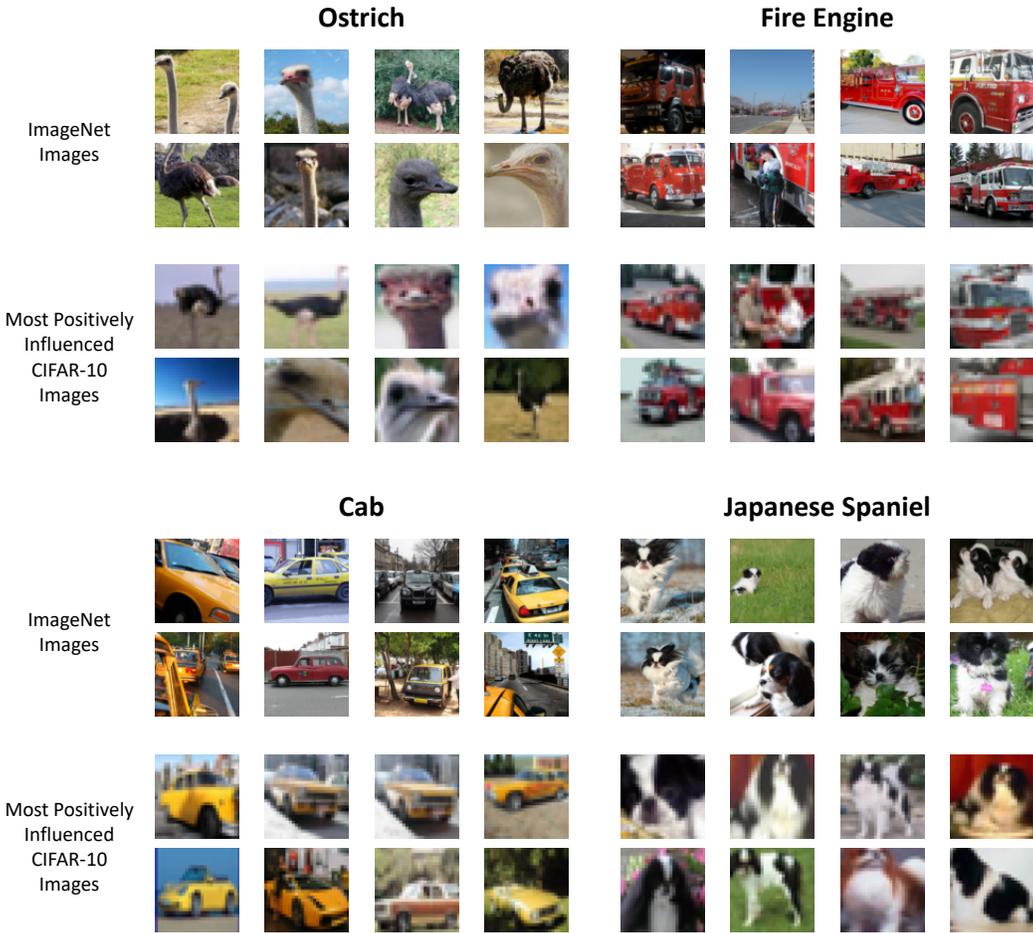


Figure 4: **Projecting source labels onto the target dataset.** The CIFAR-10 images that were most positively influenced by the ImageNet classes “ostrich”, “fire engine”, “cab”, and “Japanese Spaniel.” We find that these images look similar to the corresponding images in the source dataset.

(3) detecting data leakage between the source and target datasets. We focus our demonstration of these capabilities on a commonly-used transfer learning setting: ImageNet to CIFAR-10 (experimental details are in Appendix A).

4.1 CAPABILITY 1: EXTRACTING TARGET SUBPOPULATIONS BY PROJECTING SOURCE CLASS LABELS

Imagine that we would like to find all the ostriches in the CIFAR-10 dataset. This is not an easy task as CIFAR-10 only has “bird” as a label, and thus lacks sufficiently fine-grained annotations. Luckily, however, ImageNet *does* contain an ostrich class! Our computed influences enable us to “project” this ostrich class annotation (and, more broadly, the fine-grained label hierarchy of our source dataset) to find this subpopulation of interest in the target dataset.

Indeed, our examination from Section 3 suggests that the most positively influencing source classes are typically those that directly overlap with the target classes (see Figure 1). In particular, for our example, “ostrich” is highly positively influential for the “bird” class (see Figure 3). To find ostriches in the CIFAR-10 dataset, we thus need to simply surface the CIFAR-10 images which were most positively influenced by the “ostrich” source class (see Figure 4).

It turns out that this type of projection approach can be applied more broadly. Even when the source class is not a direct sub-type of a target class, the downstream model can still leverage salient

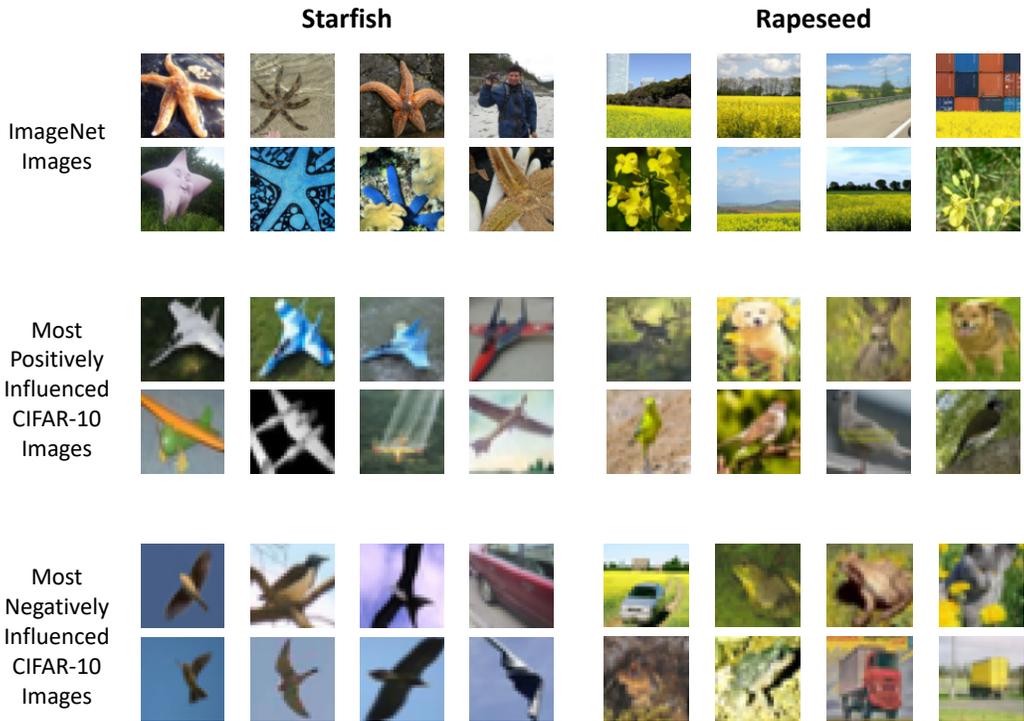


Figure 5: The CIFAR-10 images that were most positively (or negatively) influenced by the ImageNet classes “starfish” and “rapeseed.” CIFAR-10 images that are highly influenced by the “starfish” class have similar shapes, while those influenced by “rapeseed” class have yellow-green colors.

features from this class — such as shape or color — to predict on the target dataset. For such classes, projecting source labels can extract target subpopulations which share such features. To illustrate this, in Figure 5, we display the CIFAR-10 images that are highly influenced by the classes “starfish” and “rapeseed” (both of which do not directly appear in the CIFAR-10 dataset). For these classes, the most influenced CIFAR-10 images share the same shape (“starfish”) or color (“rapeseed”) as their ImageNet counterparts. More examples of such projections can be found in Appendix E.

4.2 CAPABILITY 2: DEBUGGING THE FAILURES OF A TRANSFERRED MODEL

Our framework enables us to also reason about the possible mistakes of the transferred model caused by source dataset classes. For example, consider the CIFAR-10 image of a dog in Figure 6, which our transfer learning model often mispredicts as a horse. Using our framework, we can demonstrate that this image is strongly negatively influenced by the source class “sorrel horse.” Thus, our downstream model may be misusing a feature introduced by this class. Indeed, once we remove “sorrel horse” from the source dataset, our model predicts the correct label more frequently. (See Appendix E for more examples, as well as a quantitative analysis of this experiment.)

4.3 CAPABILITY 3: DETECTING DATA LEAKAGE AND MISLEADING SOURCE EXAMPLES

Thus far, we have focused on how the *classes* in the source dataset influence the predictions of the transferred model on target examples. In this section, we extend our analysis to the *individual* datapoints of the source dataset. We do so by adapting our approach to measure the influence of each individual source datapoint on each target datapoint. Further details on how these influences are computed can be found in Appendix D.

Figure 7 displays the ImageNet training examples that have highly positive or negative influences on CIFAR-10 test examples. We find that the source images that are highly positively influential are often instances of *data leakage* between the source training set and the target test set. On the other hand,

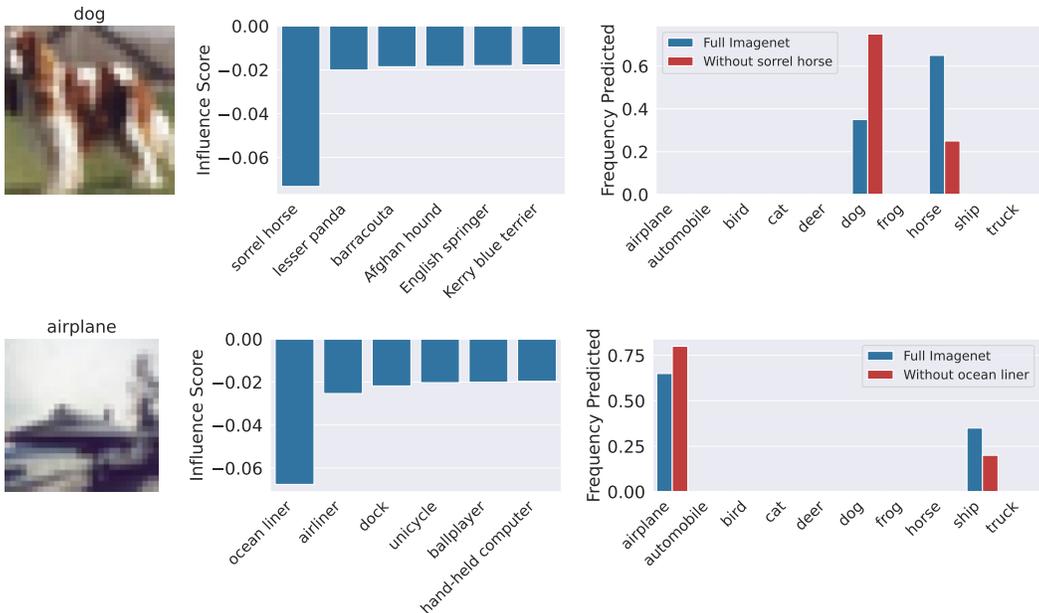


Figure 6: Pinpointing highly negatively influential source classes can help explain model mistakes. **Left:** For two CIFAR-10 images, we plot the most negatively influential source classes. **Right:** Over 20 runs, the fraction of times that our downstream model predicts each label for the given CIFAR-10 image. When the most negatively influential class is removed, the model predicts the correct label more frequently. More examples can be found in Appendix E.

the ImageNet images that are highly negatively influential are typically mislabeled, misleading, or otherwise surprising. For example, the presence of the ImageNet image of a flying lawnmower hurts the performance on a CIFAR-10 image of a regular (but similarly shaped) airplane (see Figure 7).

5 RELATED WORK

Transfer learning. Transfer learning is a technique commonly used in domains ranging from medical imaging (Mormont et al., 2018; Ke et al., 2021), language modeling (Conneau & Kiela, 2018), to object detection (Ren et al., 2015; Dai et al., 2016; Girshick et al., 2014; Chen et al., 2017). Therefore, there has been considerable interest in understanding the drivers of transfer learning’s success. For example, by performing transfer learning on block-shuffled images, Neyshabur et al. (2020) demonstrate that at least some of the benefits of transfer learning come from low-level image statistics of source data. There is also an important line of work studying transfer learning by investigating the relationship between different properties of the source model and performance on the target task (Ke et al., 2021; Kornblith et al., 2019; Salman et al., 2020; Utrera et al., 2020).

The works that are the most relevant to ours are those which studied how modifying the source dataset can affect the downstream performance. For example, Kolesnikov et al. (2019) showed that pre-training with an enormous source dataset (approximately 300 million) of noisily labeled images can outperform pre-training with ImageNet. Azizpour et al. (2015); Huh et al. (2016) investigated the importance of the number of classes and the number of images per class in transfer learning. Finally, Ngiam et al. (2018) demonstrated that more pre-training data does not always help, and transfer learning can be sensitive to the choice of pre-training data. They also presented a framework for reweighting the source datapoints in order to boost transfer learning performance.

Influence functions and datamodels. Influence functions are well-studied statistical tools that have been recently applied in machine learning settings (Hampel et al., 2011; Cook & Weisberg, 1982; Koh & Liang, 2017). For a given model, influence functions analyze the effect of a training input on the model’s predictions by estimating the expected change in performance when this training input

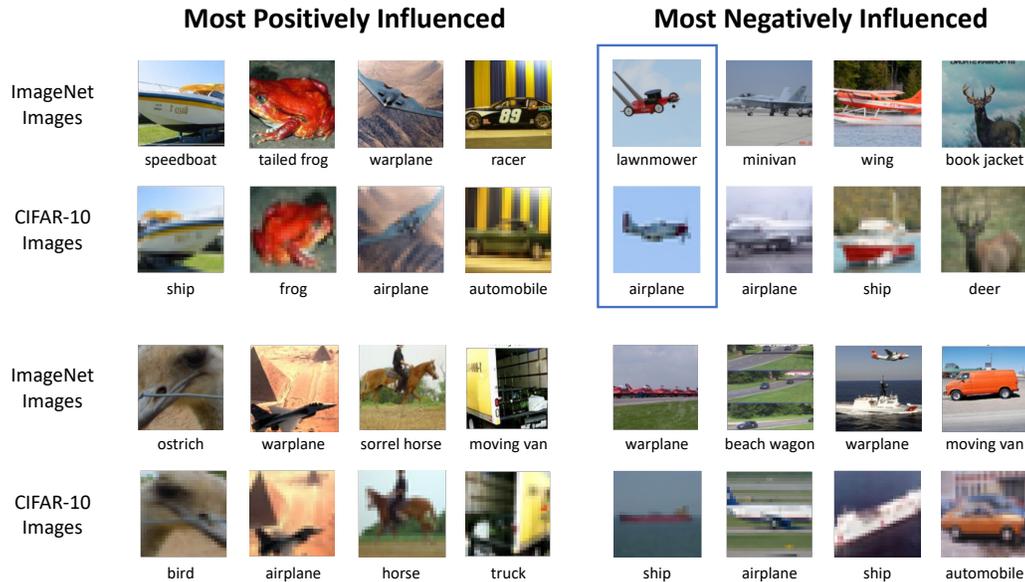


Figure 7: ImageNet training images with highest positive (**left**) or negative (**right**) example-wise (average) influences on CIFAR-10 test images. We find that ImageNet images that are highly positively influential often correspond to data leakage, while ImageNet images that are highly negatively influential are often either mislabeled, ambiguous, or otherwise misleading. For example, the presence of a flying lawnmower in the ImageNet dataset hurts the downstream performance on a similarly shaped airplane (boxed).

is added or removed. In order to apply this tool in machine learning, [Koh & Liang \(2017\)](#) propose estimating the influence functions using the Hessian of the loss function. A recent line of work estimates this quantity more efficiently by training on different subsets of the training set ([Feldman & Zhang, 2020](#)). In a similar vein, [Ghorbani & Zou \(2019\)](#) proposed running a Monte Carlo search to estimate the effect of every training input via Shapley values. More recently, [Ilyas et al. \(2022\)](#) proposed datamodeling framework as an alternative way to estimate the effect of a training input on the models’ prediction. Datamodels are represented using parametric functions (typically, linear functions) that aim to map a subset of the training set to the model’s output.

6 CONCLUSIONS

In this work, we presented a new framework for examining the impact of the source dataset in transfer learning. Specifically, our approach estimates the influence of a source class (or datapoint) that captures how including that class (or datapoint) in the source dataset impacts the downstream model’s predictions. Leveraging these estimates, we demonstrate that we can improve the transfer learning performance on a range of downstream tasks by identifying and removing detrimental datapoints from the source dataset. Furthermore, our framework enables us to identify granular subpopulations in the target dataset by projecting fine-grained labels from the source dataset, better understand model failures on the downstream task and detect potential data-leakages from the source to the downstream dataset. We believe our framework provides a new perspective on transfer learning: one that enables us to perform a fine-grained analysis of the impact of the source dataset.

REFERENCES

- Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. Factors of transferability for a generic convnet representation. *IEEE transactions on pattern analysis and machine intelligence*, 2015.
- Thomas Berg, Jiongxin Liu, Seung Woo Lee, Michelle L Alexander, David W Jacobs, and Peter N Belhumeur. Birdsnap: Large-scale fine-grained visual categorization of birds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *European conference on computer vision*, 2014.
- Emmanuel Candes, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: model-x knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(3):551–577, 2018.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- Alexis Conneau and Douwe Kiela. Senteval: An evaluation toolkit for universal sentence representations. *Language Resources and Evaluation Conference (LREC)*, 2018.
- R Dennis Cook and Sanford Weisberg. *Residuals and influence in regression*. New York: Chapman and Hall, 1982.
- Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems (NeurIPS)*, 2016.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition (CVPR)*, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Shuyang Du, Haoli Guo, and Andrew Simpson. Self-driving car steering angle prediction based on image recognition. *arXiv preprint arXiv:1912.05440*, 2019.
- Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *2004 conference on computer vision and pattern recognition workshop*, pp. 178–178. IEEE, 2004.
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pp. 2881–2891, 2020.
- Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In *International Conference on Machine Learning (ICML)*, 2019.
- Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *computer vision and pattern recognition (CVPR)*, pp. 580–587, 2014.
- Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: the approach based on influence functions*, volume 196. John Wiley & Sons, 2011.
- Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.

- Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Data-models: Predicting predictions from training data. In *International Conference on Machine Learning (ICML)*, 2022.
- Saachi Jain, Hadi Salman, Eric Wong, Pengchuan Zhang, Vibhav Vineet, Sai Vemprala, and Aleksander Madry. Missingness bias in model debugging. In *International Conference on Learning Representations*, 2022.
- Yunhun Jang, Hankook Lee, Sung Ju Hwang, and Jinwoo Shin. Learning what and where to transfer. In *International Conference on Machine Learning*, pp. 3030–3039. PMLR, 2019.
- Bojan Karlaš, David Dao, Matteo Interlandi, Bo Li, Sebastian Schelter, Wentao Wu, and Ce Zhang. Data debugging with shapley importance over end-to-end machine learning pipelines. *arXiv preprint arXiv:2204.11131*, 2022.
- Alexander Ke, William Ellsworth, Oishi Banerjee, Andrew Y Ng, and Pranav Rajpurkar. Chextransfer: performance and parameter efficiency of imagenet models for chest x-ray interpretation. In *Proceedings of the Conference on Health, Inference, and Learning*, pp. 116–124, 2021.
- Jiman Kim and Chanjong Park. End-to-end ego lane estimation based on sequential transfer learning for self-driving cars. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 30–38, 2017.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, 2017.
- Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. *arXiv preprint arXiv:1912.11370*, 2019.
- Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *computer vision and pattern recognition (CVPR)*, 2019.
- Jonathan Krause, Jia Deng, Michael Stark, and Li Fei-Fei. Collecting a large-scale dataset of fine-grained cars. 2013.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. In *Technical report*, 2009.
- Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. *arXiv preprint arXiv:2202.10054*, 2022.
- Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework for machine learning. *arXiv preprint arXiv:2110.14049*, 2021.
- Guillaume Leclerc, Andrew Ilyas, Logan Engstrom, Sung Min Park, Hadi Salman, and Aleksander Madry. ffcv. <https://github.com/libffcv/ffcv/>, 2022.
- Kaleel Mahmood, Rigel Mahmood, and Marten Van Dijk. On the robustness of vision transformers to adversarial examples. 2021.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 1995.
- Romain Mormont, Pierre Geurts, and Raphaël Marée. Comparison of deep transfer learning strategies for digital pathology. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018.
- Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? *Advances in neural information processing systems*, 33:512–523, 2020.
- Jiquan Ngiam, Daiyi Peng, Vijay Vasudevan, Simon Kornblith, Quoc V Le, and Ruoming Pang. Domain adaptive transfer learning with specialist models. *arXiv preprint arXiv:1811.07056*, 2018.

- Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, 2008.
- Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems (NeurIPS)*, 2015.
- Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Francisco Utrera, Evan Kravitz, N. Benjamin Erichson, Rajiv Khanna, and Michael W. Mahoney. Adversarially-trained deep nets transfer better. In *ArXiv preprint arXiv:2007.05869*, 2020.
- Sherrie Wang, George Azzari, and David B Lobell. Crop type mapping without field-level labels: Random forest transfer and unsupervised clustering techniques. *Remote sensing of environment*, 222:303–317, 2019.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, 2017.
- Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Computer Vision and Pattern Recognition (CVPR)*, 2010.
- Michael Xie, Neal Jean, Marshall Burke, David Lobell, and Stefano Ermon. Transfer learning from deep features for remote sensing and poverty mapping. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

A EXPERIMENTAL SETUP

A.1 IMAGENET MODELS

In this paper, we train a large number of models on various subsets of ImageNet in order to estimate the influence of each class of ImageNet on the model’s transfer performance for multiple downstream tasks. We focus on the ResNet-18 architecture from PyTorch’s official implementation found here <https://pytorch.org/vision/stable/models.html>⁴.

Training details. We fix the training procedure for all of our models. Specifically, we train our models from scratch using SGD to minimize the standard cross-entropy multi-class classification loss. We use a batch size of 1024, momentum of 0.9, and weight decay of 5×10^{-4} . The models are trained for 16 epochs using a Cyclic learning rate schedule with an initial learning rate of 0.5 and learning rate peak epoch of 2. We use standard data-augmentation: *RandomResizedCrop* and *RandomHorizontalFlip* during training, and *RandomResizedCrop* during testing. Our implementation and configuration files are attached to the submission.

A.2 IMAGENET TRANSFER TO CLASSIFICATION DATASETS

Dataset	Classes	Train Size	Test Size
Birdsnap (Berg et al., 2014)	500	32,677	8,171
Caltech-101 (Fei-Fei et al., 2004)	101	3,030	5,647
Caltech-256 (Griffin et al., 2007)	257	15,420	15,187
CIFAR-10 (Krizhevsky, 2009)	10	50,000	10,000
CIFAR-100 (Krizhevsky, 2009)	100	50,000	10,000
FGVC Aircraft (Maji et al., 2013)	100	6,667	3,333
Food-101 (Bossard et al., 2014)	101	75,750	25,250
Oxford 102 Flowers (Nilsback & Zisserman, 2008)	102	2,040	6,149
Oxford-IIIT Pets (Parkhi et al., 2012)	37	3,680	3,669
SUN397 (Xiao et al., 2010)	397	19,850	19,850
Stanford Cars (Krause et al., 2013)	196	8,144	8,041

Table 1: Image classification datasets used in this paper.

Datasets. We consider the transfer image classification tasks that are used in (Salman et al., 2020; Kornblith et al., 2019), which vary in size and number of classes. See Table 1 for the details of these datasets. We consider two transfer learning settings for each dataset: *fixed-feature* and *full-network* transfer learning.

Fixed-feature transfer. For this setting, we *freeze* the layers of the ImageNet source model⁵, except for the last layer, which we replace with a random initialized linear layer whose output matches the number of classes in the transfer dataset. We now train only this new layer for using SGD, with a batch size of 1024 using cyclic learning rate.

Full-network transfer. For this setting, we *do not freeze* any of the layers of the ImageNet source model, and all the model weights are updated. We follow the exact same hyperparameters as the fixed-feature setting.

A.3 COMPUTE AND TRAINING TIME.

We leveraged the FFCV data-loading library for fast training of the ImageNet models (Leclerc et al., 2022)⁶. Our experiments were run on two GPU clusters: an A100 and a V100 cluster.

⁴Our framework is agnostic to the choice of the model’s architecture.

⁵For all of our experiments, we do not freeze the batch norm statistics. We only freeze the weights of the model, similar to Salman et al. (2020).

⁶Using FFCV, we can train a model on the ImageNet dataset in around 1 hour, and reach $\sim 63\%$ accuracy

Training ImageNet models and influence calculation. We trained 7,540 ImageNet models on random subsets of ImageNet, each containing half of ImageNet classes. On a single V100, training a single model takes around 30 minutes. After training these ImageNet models, we compute the influence of each class as outlined in Algorithm 1. Computing the influences is fast, and takes few seconds on a single V100 GPU.

A.4 HANDPICKED BASELINE DETAILS

In our counterfactual experiments in Section 3, we automatically selected, via our framework, the most influential subsets of ImageNet classes for various downstream tasks. We then removed the classes that are detrimental to the transfer performance, and measured the transfer accuracy improvement after removing these classes. The results are summarized in Table 2b.

What happens if we hand-pick the source dataset classes that are relevant to the target dataset? Indeed, Ngiam et al. (2018) found that hand-picking the source dataset classes can sometimes boost transfer performance. We validate this approach for our setting using the WordNet hierarchy (Miller, 1995). Specifically, for each class from the target task, we look up all the ImageNet classes that are either children or parents of this target class. The set of all such ImageNet classes are used as the handpicked most influential classes. Following this manual selection, we train an ImageNet model on these classes, then apply transfer learning to get the baseline performance that we report in Table 2b.

A.5 CONVERGENCE ANALYSIS

We compute our class influence values using 7,540 source models, each of which were trained using 500 randomly chosen ImageNet classes. How many models do we actually need to compute our transfer influences?

Counterfactual Experiment To first analyze this question, we re-run our counterfactual experiment in Section 3 when using a smaller number of models to compute transfer influences (Figure 8). While using the full number of models performs the best, we get meaningful transfer influences when computing with both 4000 and 1000 models. In both cases, removing the most negatively influential classes boosts transfer learning performance over using the entire source dataset, while removing the most positively influential classes drops transfer learning performance over the random baseline.

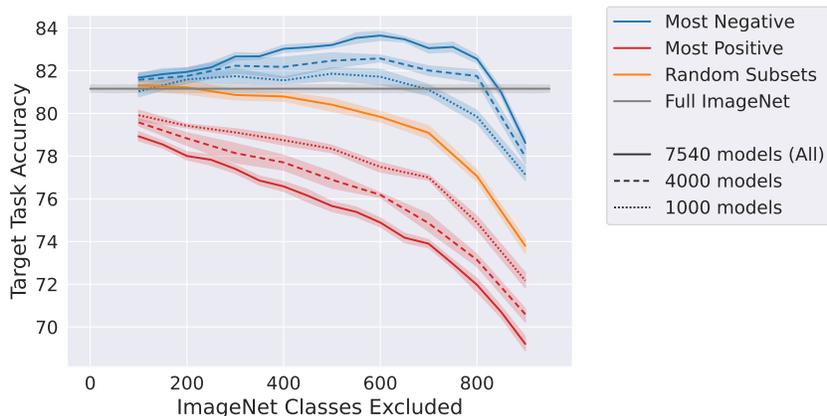


Figure 8: We replicate the counterfactual experiment in Section 3 Figure 2a using 1000 and 4000 source models for computing the transfer influences.

Bootstrap Analysis In order to analyze the convergence of the transfer influences, we track the standard deviation of the influence values after bootstrap resampling.

We consider the ImageNet \rightarrow CIFAR-10 transfer setting with fixed-feature fine-tuning. Given N models, we randomly sample, with replacement, N models to recompute our transfer influences. Specifically, we evaluate the overall transfer influences (i.e., the influence value of each ImageNet class averaged over all target examples). We perform this resampling 500 times, and measure the standard deviation of the computed overall transfer influence value for each class over these 500 resamples.

Below, we plot this standard deviation (averaged over the 1000 classes) for various number of models N . We find that the standard deviation goes down as more models are used, indicating that our estimate of the influence values has less variance. This metric roughly plateaus by the time we are using 7000 models.

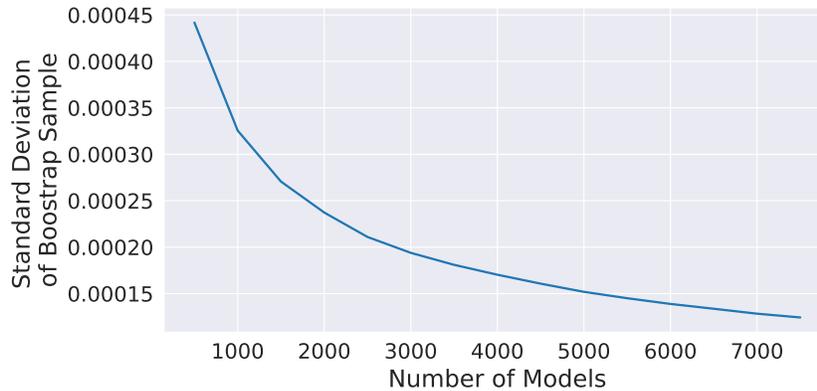


Figure 9: Standard deviation of the overall influence values (averaged over classes) after bootstrap resampling for various numbers of models.

B VARIANTS OF COMPUTING INFLUENCES

B.1 VARIATIONS OF TARGETS FOR COMPUTING TRANSFER INFLUENCES

In our paper, we used the softmax output of the groundtruth class as the target for our influence calculation. What happens if we use a different target? We compare using the following types of targets.

- Softmax Logits: the softmax output of the groundtruth class
- Is Correct: the binary value of whether the image was predicted correctly
- Raw Margins: the difference in raw output between the correct class and the most confidently predicted incorrect class
- Softmax Margins: the same as raw margins, but use the output after softmax

In Figure 10, we replicate the counterfactual experiment from the main paper in Figure 2b using these different targets. Specifically (over 2 runs), we rank the overall influence of the ImageNet classes on CIFAR-10 for fixed-feature transfer. We then remove the classes in order most most or least influence.

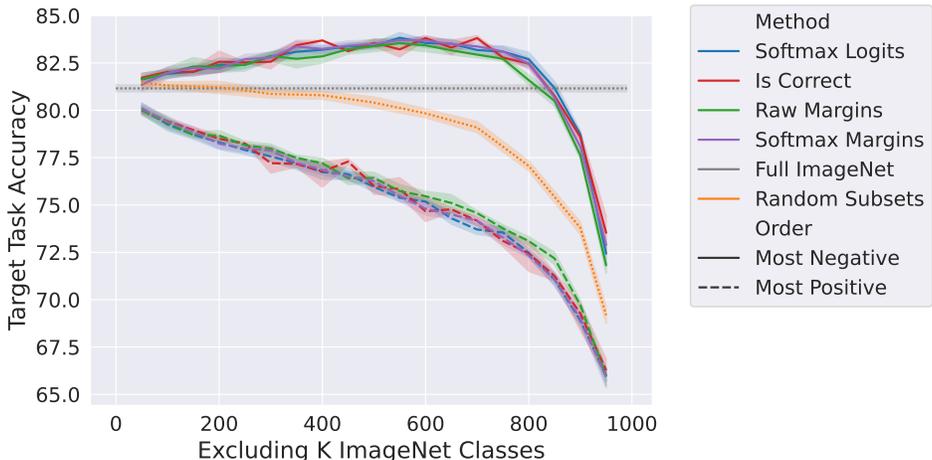


Figure 10: Target task accuracies after removing the most positively or negatively influential ImageNet classes from the source dataset with various influence targets.

We find that our method is not particularly sensitive to the individual target used. We found that using the softmax logits provided the highest benefit when removing negative influencers, and thus used that target for the reset of our experiments.

Datamodels vs. Influences. Datamodels (Ilyas et al., 2022) is another method that, similar to influences, seeks to compute the importance of a training point on a test set prediction. Specifically, instead of computing the difference in the expected accuracy of the model when a training point is removed, the method fits a linear model that, given a binary vector that denotes the composition of the training dataset, predicts the raw margin (i.e., the difference in raw output between the correct class and the most confidently predicted incorrect class). The importance of each training point is then the coefficient of the linear model for that particular example.

We adapt this method to our framework by training a linear model with ridge regression to predict the softmax output of the transfer model on the target images given a binary vector that denotes which source classes were included in the source dataset. However, we find that datamodels were more effective for computing example-based values (see Appendix D).

In Figure 11, we compare using influence values (as described in the main paper) to using these adapted datamodels. Specifically (over 5 runs), we rank the overall importance of the ImageNet

classes on CIFAR-10 for fixed-feature transfer using influences or datamodels. We then remove the classes in order most most or least influence. We find that our framework is not sensitive to the choice of datamodels or influences. However, influences performed marginally better in this counterfactual experiment, so we used influences for all other experiments in this paper.

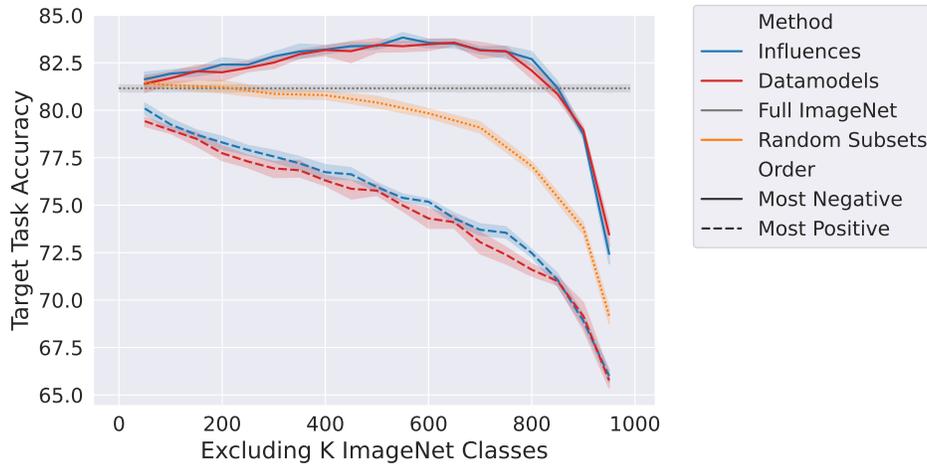


Figure 11: Target task accuracies after removing the most positively or negatively influential ImageNet classes from the source dataset using datamodels or influences.

C FULL COUNTERFACTUAL EXPERIMENT

In this section, we display the full results for the counterfactual experiment in the main paper (Figure 2b). Specifically, for each target task, we display the target task accuracies after removing the most positive (top) and negative (bottom) influencers from the dataset over 10 runs. We find that our results hold across datasets.

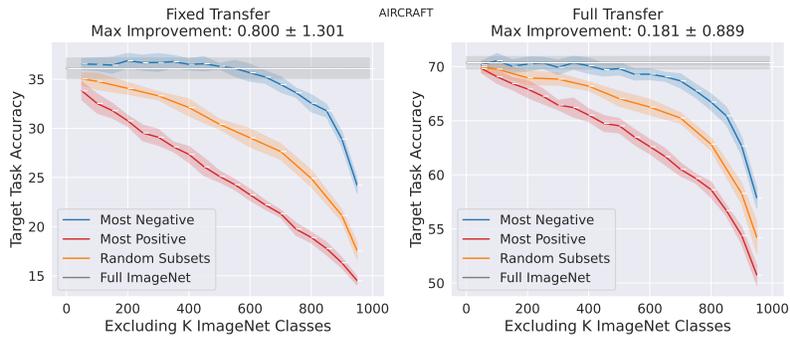


Figure 12: AIRCRAFT

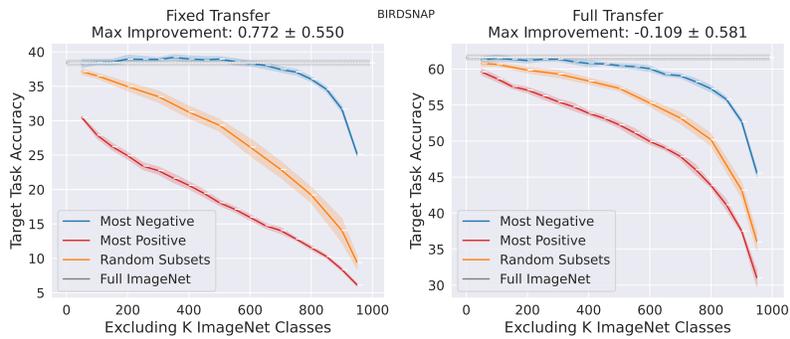


Figure 13: BIRDSNAP

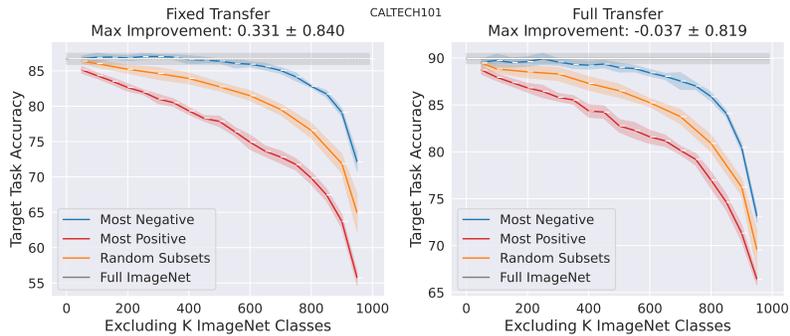


Figure 14: CALTECH101

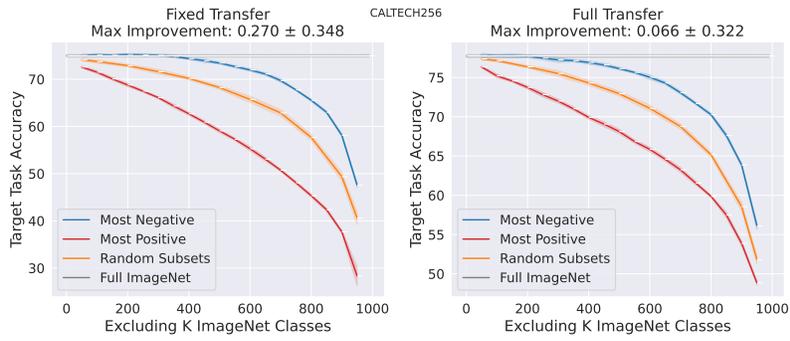


Figure 15: CALTECH256

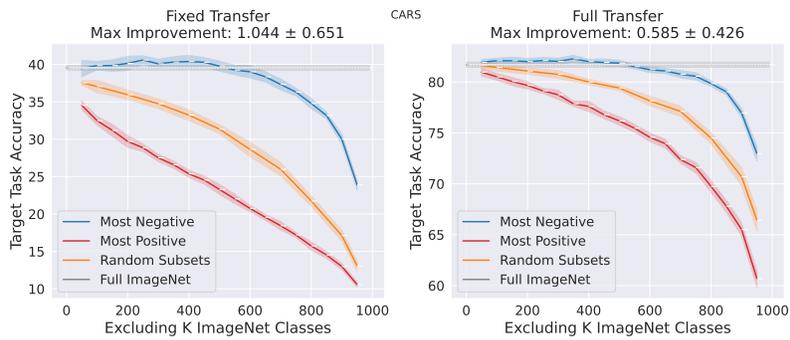


Figure 16: CARS

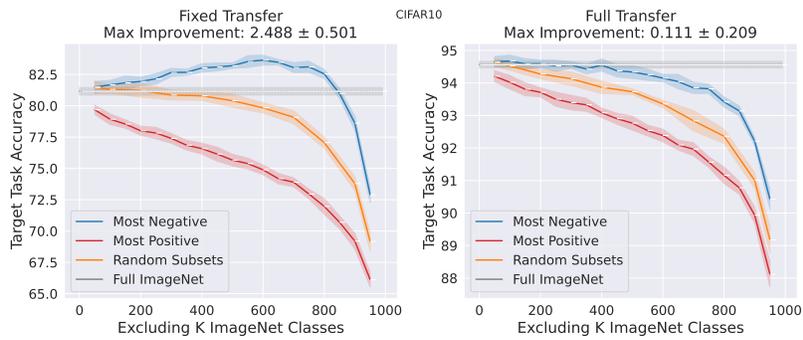


Figure 17: CIFAR10

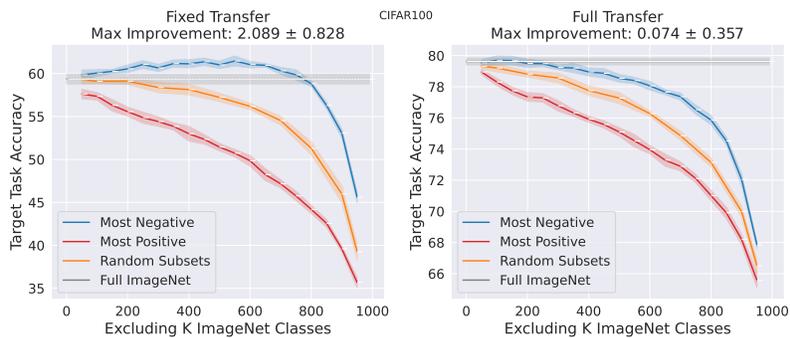


Figure 18: CIFAR100

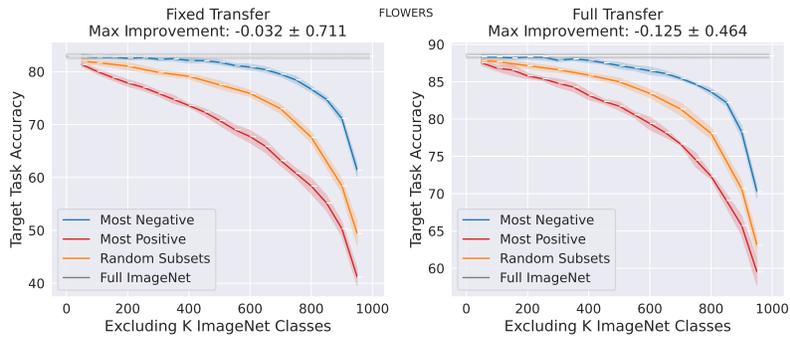


Figure 19: FLOWERS

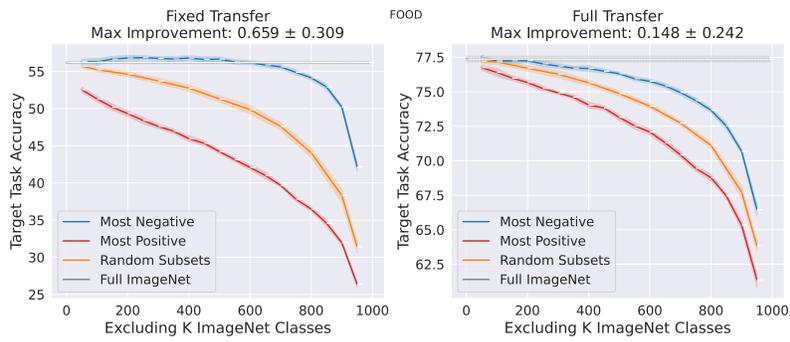


Figure 20: FOOD

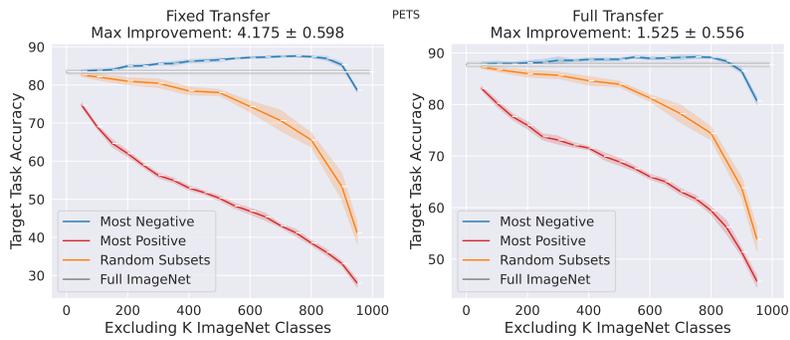


Figure 21: PETS

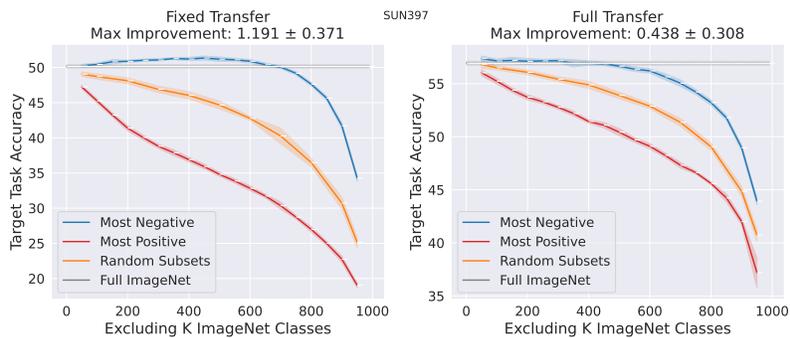
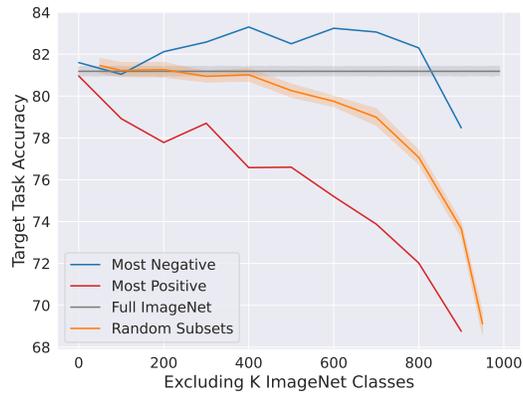


Figure 22: SUN397

C.1 PERFORMING COUNTERFACTUAL EXPERIMENTS WITH A VALIDATION SET

In this section, we replicate the experiment in Figure 2a of the main paper. Here, we split the original CIFAR-10 test set into half, where the first half is the validation set and the second is the test set.

We compute our overall class influences on the validation set, and evaluate the accuracy of removing the top and bottom source classes on the held out test set. The results match Figure 2a. Thus, the transfer influences generalize beyond the specific target examples used to compute the influences.



D ADAPTING OUR FRAMEWORK TO COMPUTE THE EFFECT OF EVERY SOURCE DATAPPOINT ON TRANSFER LEARNING

We have presented in the main paper how to compute the influences of every class in the source dataset on the predictions of the model on the target dataset. In that setup, we demonstrated multiple capabilities of our framework, such as improving overall transfer performance, detecting particular subpopulations in the target dataset, etc. Given the wide range of capabilities class-based influences provide, one natural question that arises: Can we compute the influence of every source datapoint on the predictions of the model on the target dataset? Furthermore, what do these influences tell us about the transfer learning process?

Mathematically, the computation of example-based influences (i.e., the influence of every source datapoint) is very similar to the computation of class-based influences. Specifically, to compute example-based influences, we start by training a large number of models on different subsets of the source datapoints (as opposed to source classes for class-based influences). Next, we estimate the influence value of a source datapoint s on a target example t as the expected difference in the transfer model’s performance on example t when datapoint s was either included or excluded from the source dataset:

$$\text{Infl}[s \rightarrow t] = \mathbb{E}_S [f(t; S) \mid s \in S] - \mathbb{E}_S [f(t; S) \mid s \notin S] \quad (2)$$

where $f(t; S)$ is the softmax output of a model trained on a subset S of the source dataset. Similar to class-based influences, a positive (resp. negative) influence value indicates that including the source datapoint s improves (resp. hurts) the model’s performance on the target example t .

While example-based influences provide some insights about the transfer process, we found that—in this regime—datamodels (Ilyas et al., 2022) provide cleaner results and better insights. Generally, influences and datamodels measure similar properties: the effect of the source datapoints on the target datapoints. For a particular target datapoint t , we measure the effect of every source datapoint s with datamodels by solving a regression problem. Specifically, we train a large number of models on different subsets of the source dataset. For every model f_i , we record 1) a binary mask $\mathbb{1}_{S_i}$ that indicates which source datapoints were included in the subset S_i of the source dataset, and 2) the transfer performance $f_i(t; S_i)$ of the model f_i on the target datapoint t after fine-tuning on the target dataset. Following the training and the fine-tuning stages, we fit a linear model g_w that predicts the transfer performance $f(t; S)$ from a random subset S of the source dataset as follows: $f(t; S) \approx g_w(\mathbb{1}_S) = w^\top \mathbb{1}_S$. Given this framework, $w = (w_1, w_2, \dots, w_L)$ measures the effect of every source datapoint s on the target datapoint t ⁷. We present the overall procedure in Algorithm 2.

Algorithm 2 Example-based datamodels estimation for transfer learning.

Require: source dataset $\mathcal{S} = \cup_{l=1}^L s_l$ (with L datapoints), a target dataset $\mathcal{T} = (t_1, t_2, \dots, t_n)$, training algorithm \mathcal{A} , subset ratio α , number of models m

- 1: Sample m random subsets $S_1, S_2, \dots, S_m \subset \mathcal{S}$ of size $\alpha \cdot |\mathcal{S}|$:
- 2: **for** $i \in 1$ to m **do**
- 3: Train model f_i by running algorithm \mathcal{A} on S_i
- 4: **end for**
- 5: Fine-tune f_i on the training target dataset
- 6: **for** $j \in 1$ to n **do**
- 7: Collect datamodels training set $\mathcal{D}_j = \{(\mathbb{1}_{S_i}, f_i(t_j, S_i))\}_{i=1}^m$
- 8: Compute w_j by fitting LASSO on \mathcal{D}_j
- 9: **end for**
- 10: **return** $w_j \forall j \in [n]$

⁷To estimate the datamodels, we train 71,828 models on different subsets of the source dataset.

E OMITTED RESULTS

E.1 PER-CLASS INFLUENCERS

We display for the ImageNet \rightarrow CIFAR-10 the top (most positive) and bottom (most negative) influencing classes for each CIFAR-10 class. This is the equivalent to the plot in Figure 3 in the main paper.

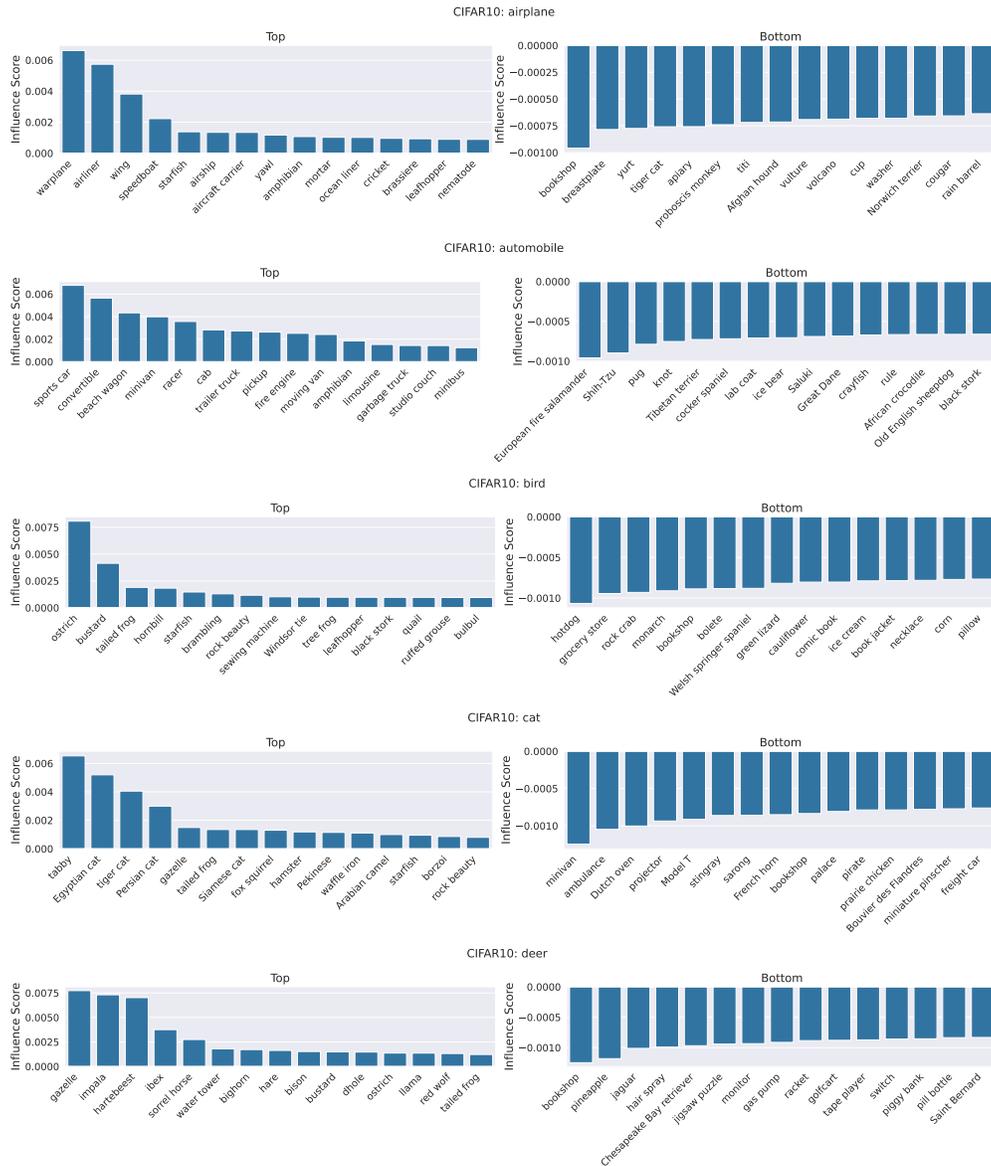


Figure 23: Top and bottom influencing ImageNet classes for all CIFAR-10 classes.

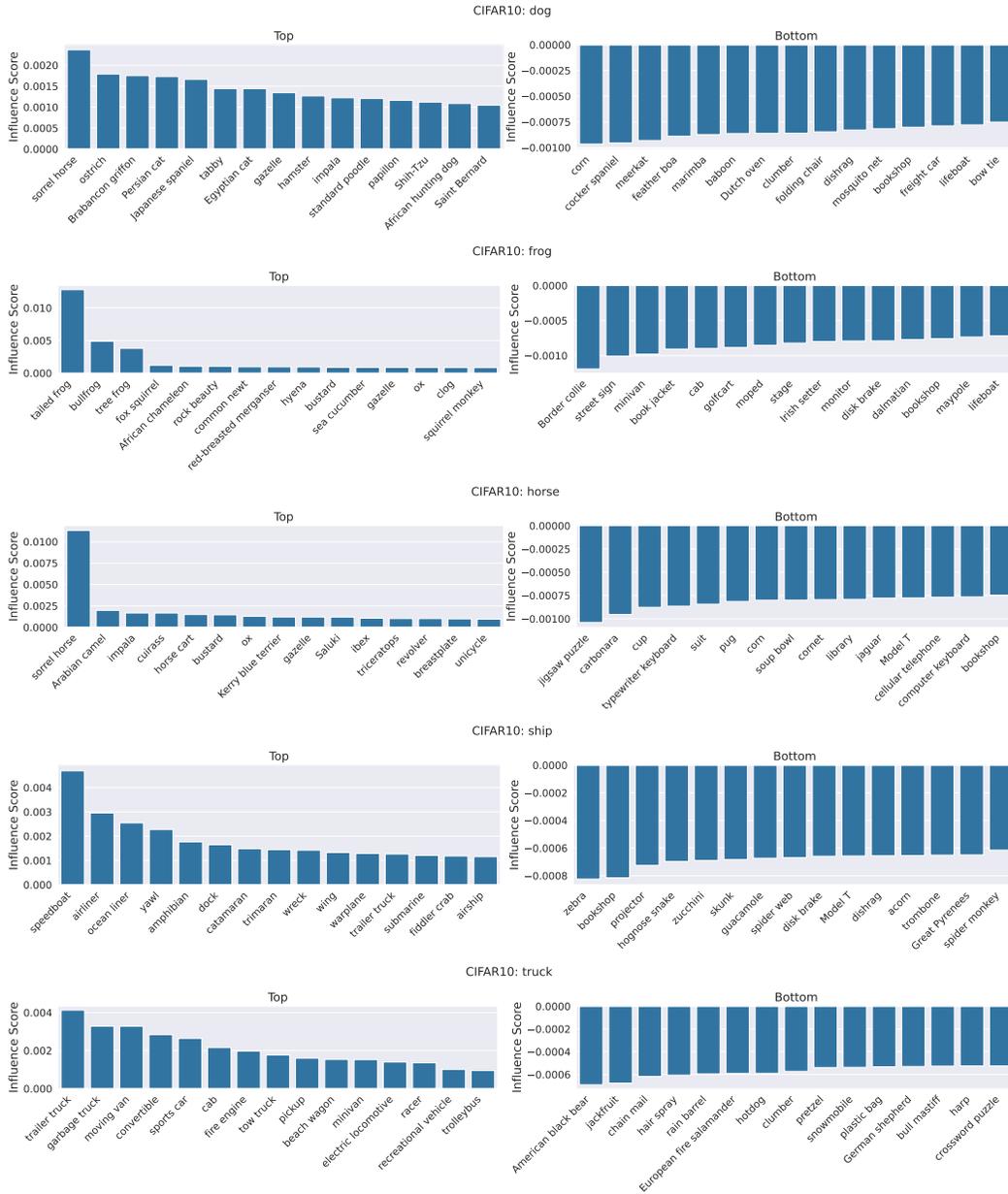


Figure 24: Top and bottom influencing ImageNet classes for all CIFAR-10 classes.

E.2 MORE EXAMPLES OF EXTRACTED SUBPOPULATIONS FROM THE TARGET DATASET

Here, we depict more examples of extracting subpopulations from the target dataset (as in Figure 4 of the main paper).



Figure 25: For each ImageNet class, we show the CIFAR-10 examples which were most positively influenced by that ImageNet class.

E.3 MORE EXAMPLES OF TRANSFER OF SHAPE AND TEXTURE FEATURE

We depict more examples of ImageNet influencers which transfer shape or texture features (as in Figure 5).



Figure 26: For each ImageNet class, we show the CIFAR-10 examples which were most highly influenced by that ImageNet class.

E.4 MORE EXAMPLES OF DEBUGGING MISTAKES OF TRANSFER MODEL USING INFLUENCERS

We display more examples of how our influences can be used to debug the mistakes of the transfer model, as presented in Figure 6 in the main paper. We find that, in most cases (Figure 27a, 27b, 27c), removing the top negative influencer improves the model’s performance on the particular image. There are a few examples where removing the top negative influencer hurts the model’s performance on the image (Figure 27d).

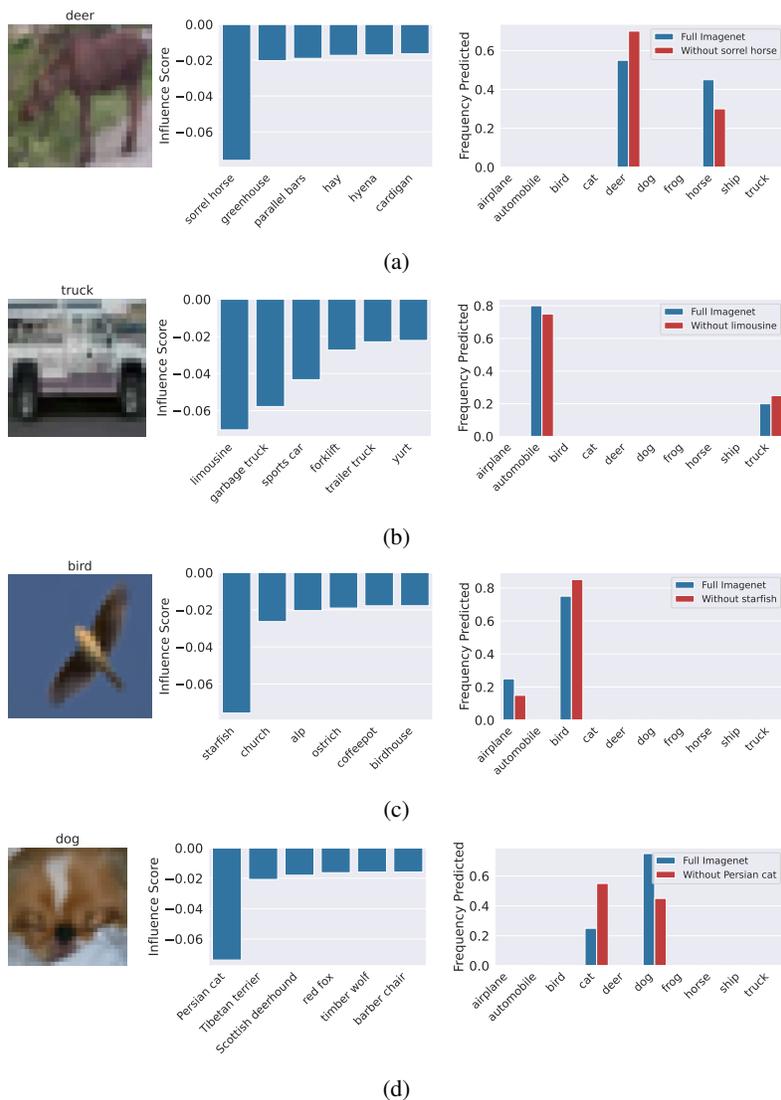


Figure 27: More examples of debugging transfer mistakes through our framework (c.f. Figure 6 in the main paper). For each CIFAR-10 image (left), we plot their most negative influencers (middle). On the right, we plot for each image the fraction (over 20 runs) of times that our transfer model predicts each class with and without the most negative influencer. While in most cases (a, b, c) we find that removing the most negative influencer helps the model predict the image correctly, in some cases removing the negative influencer does not outweigh the impact of removing images from the source dataset (d).

Quantitative analysis. How often does removing the most negative influencer actually improve the prediction on an image? For each of the following 14 classes, we run 20 runs of the ImageNet \rightarrow CIFAR-10 fixed transfer pipeline while excluding that single class from the source dataset: [“sorrel horse”, “limousine”, “minivan”, “fireboat”, “ocean liner”, “Arabian camel”, “Persian cat”, “ostrich”, “gondola”, “pool table”, “starfish”, “rapeseed”, “tailed frog”, “trailer truck”]. We compare against running the pipeline with 20 runs of the entire ImageNet dataset. Then, we look at individual CIFAR-10 images which were highly negatively influenced by one of these ImageNet classes, and check whether the images were predicted correctly more or less often when the top negative influencers were removed from the source dataset.

Of the 30 most negatively influenced ImageNet class/CIFAR-10 image pairs, 26 of them had the most negative ImageNet influencer in the above classes. Of those, 61.5% were predicted correctly more often when the negatively influential ImageNet class was removed, 34.6% were predicted incorrectly more often, and 3.9% were predicted correctly the same number of times.

We then examine the top 8 most influenced CIFAR-10 images for each of the above 14 ImageNet classes. Of those 112 images, 53% were predicted correctly more often when the image was removed, 34% were predicted incorrectly more often, and 14% were predicted correctly the same number of times.

We thus find that, for the most influenced CIFAR-10 images, removing the top negative influencer usually improves that specific image’s prediction (even though we are removing training data from the source dataset).

E.5 DO INFLUENCES TRANSFER?

E.5.1 TRANSFER ACROSS DATASETS

In this section, we seek to understand how much task-specific information is in the transfer influences that we compute. To do so, we use the transfer influences computed for CIFAR-10 in order to perform the counterfactual experiments for other datasets. We find that while using the CIFAR-10 influence values for other target datasets is more meaningful than random, they do not provide the same boost in performance when removing bottom influencers as using the task-specific influences. We thus conclude that the influence values computed by our framework are relatively task-specific.

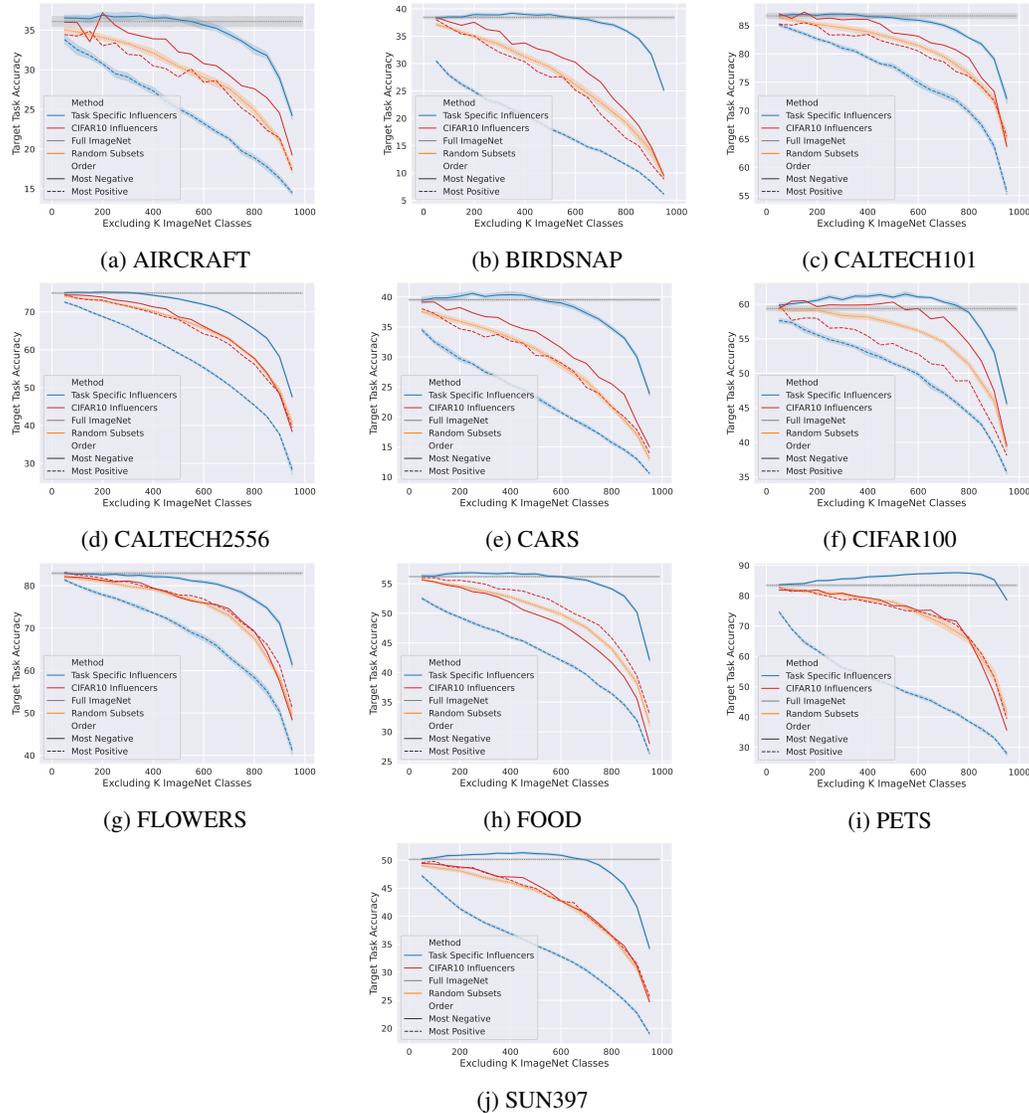


Figure 28: We repeat the counterfactual experiments (c.f. Figure 2b from the main paper) but using the influence values computed for CIFAR-10 on other target datasets.

E.5.2 TRANSFER ACROSS ARCHITECTURES

How well do our transfer influences work across architectures? Recall that we computed our transfer influences using a ResNet-18. We now repeat the counterfactual experiment from Figure 2b, using the ResNet-18 influences to remove classes from the source dataset when training a ResNet-50. We find that these influences transfer relatively well: we can thus use a smaller architecture when computing transfer influences in lieu of a larger model.

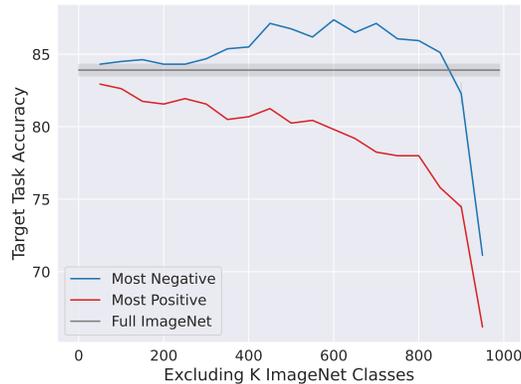


Figure 29: We repeat the counterfactual experiments (c.f. Figure 2b from the main paper) but use our influence values computed using a ResNet-18 on a ResNet-50.

We now perform the same experiment with vision transformers (ViT) (Dosovitskiy et al., 2021), which use self-attention instead of convolutions (and thus behave differently than standard CNNs (Mahmood et al., 2021; Jain et al., 2022)). Specifically, we evaluate our transfer influences on ViT-T and ViT-S (otherwise using the same hyperparameters as the ResNet-18). The ResNet-18 influences are still meaningful for the vision transformers, as removing the top classes degrades accuracy far more than the bottom classes. However, the ResNet-18 influences do not generalize as well to the ViTs as they did for the ResNet-50, as expected. In this setting, it may be more effective to use a small vision transformer (e.g., ViT-T) when computing the transfer influences.

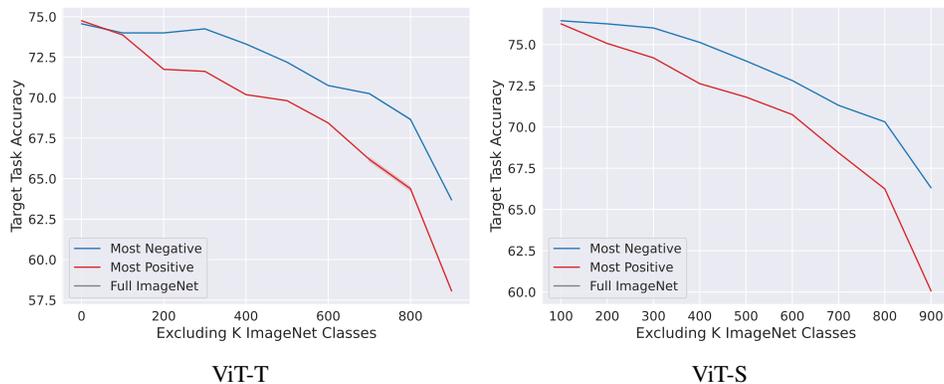


Figure 30: We repeat the counterfactual experiments (c.f. Figure 2b from the main paper) but use our influence values computed using a ResNet-18 on a ViT-T and ViT-S.

F FURTHER CONVERGENCE ANALYSIS

In this section, we analyze the sample complexity of our influence estimation in more detail. In particular, our goal is to understand how the quality of influence estimates improves with the number of models trained. We also look at the impact of different choices of model outputs.

Instead of directly measuring our downstream objective (transfer accuracy on target dataset after removing the most influential classes), which is expensive, we design two proxy metrics to gauge the convergence of our estimates:

Rank correlation. As Ilyas et al. (2022) shows, we can associate influences with a particular linear regression problem: given features $\mathbb{1}_{S_i}$, an indicator vector of the subset of classes in the source dataset, predict the labels $f(t; S_i)$, the model’s output after it is finetuned on target dataset \mathcal{T} . In fact, we can interpret influences as weights corresponding to these binary features for presence of each class. That is, the influence vector $w_t = \{\text{Infl}[C_i \rightarrow t]\}_i$ defines a linear function that, given a subset of classes that the source model is trained on, predicts the corresponding model’s output when trained on that subset. Here, we focus on analyzing average model accuracy, so in fact we consider the aggregate output $\sum_i f(t; S_i)$ and the corresponding aggregated influences $w = \{\text{Infl}[C_i]\}_i$.

Given this view, we can measure the quality of the influence estimates by measuring their performance on the above regression problem on a held-out⁸ set of examples $\{S_i, \sum_i f(t; S_i)\}$. In order to make different choices of model outputs (logit, confidence, etc.) comparable, we measure performance with spearman rank correlation between the ground truth model outputs and the predictions of the linear model (whose weights are given by the influences).

We measure this correlation while varying both the number of trained models used in the influence estimation and the choice of the model output, and the results are shown in Figure 31.

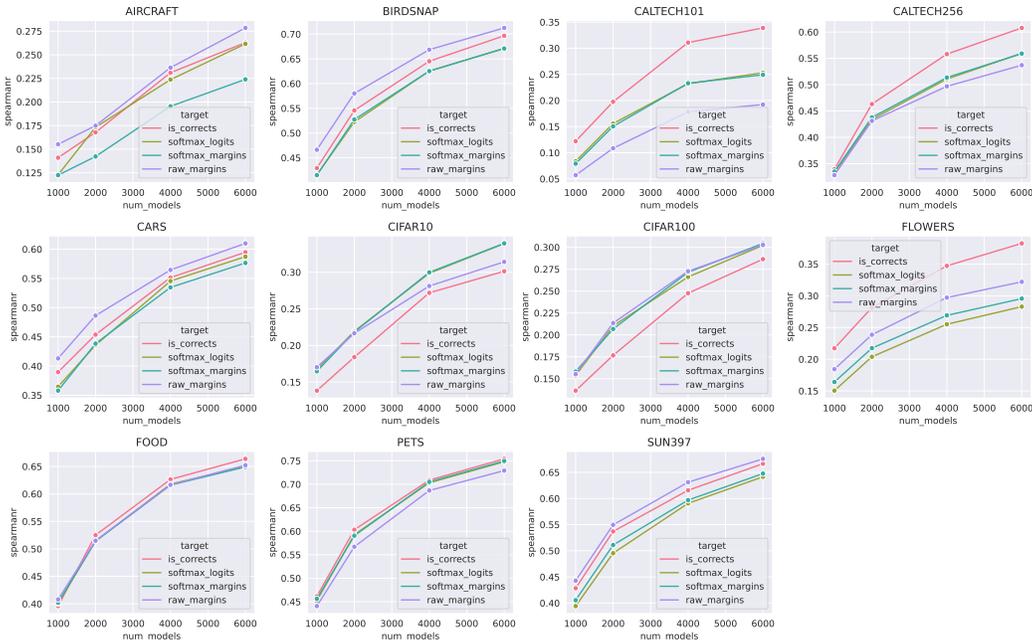


Figure 31: Measuring improvement in influence estimates using the **rank correlation** metric. The rank correlation here measures how well influence estimates perform in the underlying regression problem of predicting target accuracy from the subset of classes included in the training set. We evaluate on a held-out set of subsets independent from those used to estimate influences. Across all datasets, correlation improves significantly with more trained models.

⁸We split the 7,540 models into a training set of 6,000 and a validation set of the remainder.

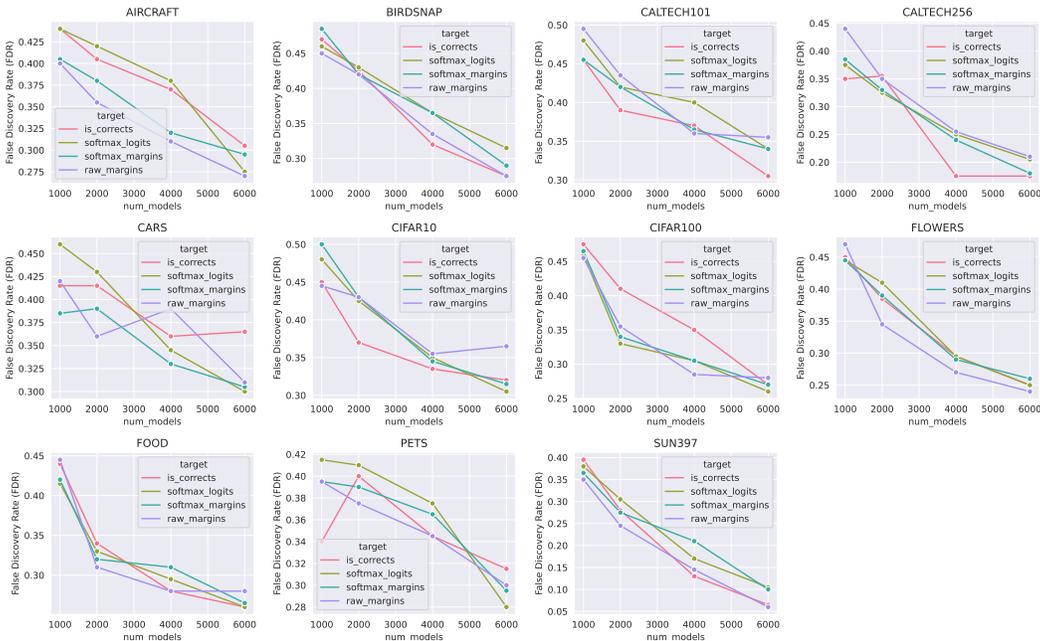


Figure 32: Measuring improvement in influence estimates using the **False Discovery Rate** heuristic. Using a procedure (loosely) based on the Knockoffs framework, we estimate the proportion of false discoveries within the top 100 features ranked by estimated influences. Across all datasets, FDR decreases generally with more trained models.

False discovery rate. Above we considered a measure based on predictive performance. Here, we focus on a more parameter-centric notion of False Discovery Rate (FDR). Intuitively, FDR here quantifies the following: how often are the top influencers actually just due to noise?

The Knockoffs (Candes et al., 2018) framework allows one to perform feature selection and also estimate the FDR. At a high level, it consists of two steps: First, one constructs “knockoff” versions of the original features which are distributed “indistinguishably” (more formally, exchangeable) from the original features, and at the same time are independent of the response by design. Second, one applies an estimation algorithm of choice (e.g., OLS or LASSO) to the augmented data consisting of both the original and the knockoff features. Then, the relative frequency at which a variable X_i has higher statistical signal than its knockoff counterpart \tilde{X}_i indicates how likely the algorithm chooses true features, and this can be used to estimate the FDR (intuitively, if X_i is independent of the response y , X_i is indistinguishable from its knockoff \tilde{X}_i and both are equally likely to have higher score).

We adapt this framework here (particularly, the version known as model-X knockoffs) as follows:

1. Sample an independent knockoff matrix \tilde{X} consisting of 1,000 binary features from the same distribution as the original mask matrix X (namely, each instance has 500 active features).⁹
2. Estimate influences for both original and knockoff features using the difference-in-means estimator.
3. Consider the top $k = 100$ features by positive influence, and count the proportion of features that are knockoff. This yields an estimate of FDR among the top 100 features. An FDR of 0.5 indicates chance-level detection.¹⁰

⁹Technically, this procedure is not exactly valid in the original FDR framework as X_i and \tilde{X}_i are exchangeable due to dependencies in the features. Nonetheless, it is accurate up to some approximation.

¹⁰This is different from the usual manner of controlling the FDR, but we look at this alternative metric for simplicity.

As with the previous metric, we measure the above FDR for each target dataset while varying the number of trained models and the target output type (Figure 32).

Discussion. We observe the following from the above analyses using our two statistics:

- There are significant gains (higher correlation and lower FDR) with increasing number of trained models.
- But neither metric appears to have plateaued with 6,000 models, so this indicates that we can improve the accuracy of our influence estimates with more trained models, which may in turn improve the max improvement in transfer accuracies (Section 3), among other results.
- The choice of the target type does not appear to have a significant or consistent impact across different datasets, which is also consistent with the results of our counterfactual experiments (Appendix B.1).