

---

# When does dough become a bagel?

## Analyzing the remaining mistakes on ImageNet

---

Vijay Vasudevan<sup>1</sup> Benjamin Caine<sup>1</sup> Raphael Gontijo-Lopes<sup>1</sup> Sara Fridovich-Keil<sup>1,2</sup> Rebecca Roelofs<sup>1</sup>

### Abstract

Image classification accuracy on the ImageNet dataset has been a barometer for progress in computer vision over the last decade. Several recent papers have questioned the degree to which the benchmark remains useful to the community (Stock & Cissé, 2017; Beyer et al., 2020; Shankar et al., 2020; Yun et al., 2021; Tsipras et al., 2020), yet innovations continue to contribute gains to performance, with today’s largest models achieving 90%+ top-1 accuracy. To help contextualize progress on ImageNet and provide a more meaningful evaluation for today’s state-of-the-art models, we manually review and categorize every remaining mistake that a few top models make and provide insights into the long-tail of errors on one of the most benchmarked datasets in computer vision. We focus on the multi-label subset evaluation of ImageNet, where today’s best models achieve upwards of 97% accuracy. Our analysis reveals that nearly half of the supposed mistakes are not mistakes at all, and we uncover new valid multi-labels, demonstrating that, without careful review, we are significantly underestimating the performance of these models. On the other hand, we also find that today’s best models still make a significant number of mistakes (40%) that are obviously wrong to human reviewers. To calibrate future progress on ImageNet, we provide an updated multi-label evaluation set, and we curate **ImageNet-Major**: a 68-example "major error" slice of the obvious mistakes made by today’s top models—a slice where models should achieve near perfection, but today are far from doing so.

## 1. Introduction

<sup>1</sup>Google Research, Brain Team <sup>2</sup>University of California, Berkeley. Correspondence to: Vijay Vasudevan <vrv@google.com>, Rebecca Roelofs <rofls@google.com>.

Computer vision models often evaluate their performance on the ImageNet classification dataset (Deng et al., 2009; Russakovsky et al., 2015) and many variants (Recht et al., 2019; Hendrycks & Dietterich, 2019; Hendrycks et al., 2019; 2020; Wang et al., 2019), as a signal of capability for visual understanding. As performance on the standard sets have reached diminishing returns to top-1 and top-5 accuracy, much recent work (Stock & Cissé, 2017; Beyer et al., 2020; Recht et al., 2019; Shankar et al., 2020; Tsipras et al., 2020; Kornblith et al., 2019) has focused on understanding what is left for the computer vision community to solve, and where the community should be driving toward. Prior studies of ImageNet errors have identified issues stemming from lack of multi-labels, label noise, under-specified classes, and more (Stock & Cissé, 2017; Shankar et al., 2020; Beyer et al., 2020; Tsipras et al., 2020; Lee et al., 2017).



Label: dough; Model: bagel.  
When does dough become a bagel?

Label errors and label noise affect the evaluation of any model (Northcutt et al., 2021b), and ImageNet is no exception. Many studies above have spent effort to correct and improve these labels, showing that while ImageNet performance improvements are approaching diminishing returns, the dataset can remain useful to the community, but only if we collectively continue to shepherd it. As the best models improve, however, it is becoming increasingly challenging to assess the often novel predictions these models make. For example, should we penalize models for being the first to predict that a pre-baked bagel may be a bagel, as one of the models we review in this work does?

Machine learning models tend to make mistakes with varying severity and importance (a function of both the prediction as well as the label definitions), and prior work (Stock & Cissé, 2017; Beyer et al., 2020) has shown that non-experts find it challenging to determine the correctness of a model’s prediction on ImageNet. Our own experience, highlighted by the doughy-bagel, is that many of the remaining mistakes these top models make are quite reasonable and

probably should not be considered mistakes—understanding the severity and type of these remaining mistakes can help calibrate progress.

Indeed, to our knowledge, there has not been an expert-review, categorization, and severity assessment of the remaining long-tailed mistakes, which becomes particularly important at these margins. Our experience working with production teams on deployed applications has suggested that manual triage and assessing individual failures provides a useful indicator of model performance that aggregate measures fail to capture. Thus, in this work we attempt to analyze (as expert reviewers) *every remaining mistake* that a few state-of-the-art models make to better understand (a) which of the remaining mistakes remain egregious errors, (b) what error category they might fall in, and (c) what evaluations might capture the most important long-tail failures.

In this paper we analyze the ImageNet multi-label validation subsets (Shankar et al., 2020), in which expert labelers were used to assess the correctness of model predictions through the year 2020, and on which a 1000-image human-evaluated subset provides a direct comparison to expert human performance. By analyzing the mistakes of two large 2022-era ImageNet models, we found that:

- **Nearly half of each model’s mistakes were deemed correct** under a careful, expert multi-label re-evaluation, halving the error rate. Had we not analyzed the models’ mistakes, we would be severely underestimating the models’ actual performance.
- **Approximately 40% of the remaining mistakes can be classified as ‘major’ errors:** errors that most humans would likely not make, suggesting that many of the long-tailed mistakes aren’t simply label noise, but legitimate mistakes that leave room for improvement.

What do these lessons portend for the future of ImageNet evaluation? Top-1 will become increasingly noisy as our best models get better (though we have not yet completely saturated top-1). Our work shows that multi-label accuracy, while better at capturing "true" errors compared to top-1, suffers from a lack of a comprehensive, accurately labeled large evaluation set, which is expensive to procure and challenging to maintain.

We therefore propose ImageNet-M, a 68-example evaluation split composed of "major" mistakes that several top-performing models make; we believe this subset is one that future image classification models should achieve near perfect accuracy on, and provides three clear benefits: (1) we attempt to comprehensively-label all examples for multi-label annotations to prevent the need to review novel correct predictions, (2) we endeavor to maintain and provide a way for the public to add new correct predictions; (3) the evaluation set is small enough to encourage completeness and allow the community to inspect their own errors.

**Dataset release.** To evaluate on the `imagenet2012_multilabel` and the ImageNet-M subset, please visit [https://www.tensorflow.org/datasets/catalog/imagenet2012\\_multilabel](https://www.tensorflow.org/datasets/catalog/imagenet2012_multilabel). We also include a notebook with example code for evaluating on these splits, including multi-label accuracy and the ImageNet-M split using pre-computed logits for the ViT-3B and Greedy Soups models used in the paper.

## 2. Mistake analysis method and taxonomy

To obtain an initial set remaining mistakes, we used a standard ViT (Dosovitskiy et al., 2021) model scaled to 3B parameters (ViT-3B) that was pre-trained on JFT-3B (Sun et al., 2017) and fine-tuned on ImageNet-1K (Deng et al., 2009), achieving a top-1 accuracy of 89.5% (details in Appendix E). We also later reviewed mistakes made by the Greedy Soups model (Wortsman et al., 2022). Using the `imagenet2012_multilabel` dataset (Shankar et al., 2020), we measured the initial *multi-label accuracy* (MLA) of the ViT-3B model to be 96.3%. We describe the procedure we followed to review mistakes using in Appendix B.1 and the severity and category definitions in Appendix B.2.

## 3. Analyzing the remaining mistakes

After review of all original 676 mistakes (comprising both novel predictions and previously reviewed mistakes), we found that a total of 298 were either correct or unclear, or determined the original groundtruth incorrect or problematic. Our evaluation of the ViT-3B model on this re-labeled dataset is shown in Table 1, with the model making a total of 378 mistakes on the dataset. In other words, approximately 44% of the initial mistakes made by this model were determined to be correct!

### 3.1. Mistake category and severity

Each of the 378 remaining mistakes was assigned both a mistake category and severity (Table 2) by the expert panel.

**Category:** 78.3% of errors were assessed to be fine-grained in nature (either in-vocabulary or OOV), and 13.8% categorized as a spurious correlation. To measure whether these categories are meaningful, we measured the "hierarchy distance" between the groundtruth label and the model’s prediction using the WordNet hierarchy (Miller, 1995). For example, a hierarchy distance of 1 means that the groundtruth and model prediction share the same parent; a distance of 2 means they share the same grandparent. We found that 80.9% of errors with a hierarchy distance of 1 were assessed as fine-grained. In contrast, 54.3% of errors with a hierarchy distance of 3 were fine-grained, matching intuition that predictions close in the hierarchy are very likely to be

	All		Organisms		Objects	
	MLA	MLA Re-labeled	MLA	MLA Re-labeled	MLA	MLA Re-labeled
ViT-3B	96.3%	97.9%	96.3%	97.8%	96.4%	98.0%

Table 1: Multi-label accuracy (MLA) of ViT-3B model before and after our re-labeling.

fine-grained errors, and predictions far in the hierarchy are more likely to be spurious correlations, fine-grained with out-of-vocabulary, or non-prototypical examples. We note that WordNet does not provide a perfect (automated) category function: "goblet" and "vase" are a hierarchy distance of 4 apart, and we encountered one model mistake for this pair that we nonetheless assessed as fine-grained.

**Severity:** We determined that around 40% of errors were assessed to be "major" errors, indicating that this model still appears to make mistakes that a human familiar with the class definitions would not make, despite the fact that the model on average performs better than an expert human. We return to 'major' errors later in Section 4.1, as we believe that a subset of these errors can be a useful evaluation slice for future ImageNet benchmarking.

### 3.2. Generalization to new models

As models produce higher top-1 accuracy, how do the types of mistakes they make and improve upon change? We use the Greedy Soups model that obtains 90.9% top-1 accuracy on ImageNet validation (Wortsman et al., 2022), measuring its MLA after our initial re-labeling at 98.1%, and yielding 341 total remaining (and partially unreviewed) errors.

The Soups model corrected 209 mistakes that the ViT-3B model made, while the model made 170 mistakes where the ViT-3B model was correct, yielding an overall accuracy improvement; 28 mistake examples were common with ViT-3B but with a different prediction. In total there were 198 novel predictions made by this model that needed to be reviewed; upon review using the same panel method, we found 46.5% (92/198) were problematic (10) or actually correct (82), showing with a second model that model predictions on mistakes need to be reviewed, and that the single label expected by top-1 is often insufficient. The Soups model in the end made only 249 errors, for an MLA of 98.6%, a 0.5% absolute increase compared to unreviewed mistakes. The categorization and severity of these errors is shown in Table 2. A chi-square test of independence shows that there is no significant difference between either the categorization or severity ( $\chi^2(3, N = 629) = 2.41, p = .49$ ) or mistake severity ( $\chi^2(1, N = 629) = 1.61, p = .20$ ).

**Additional results:** We include additional results analyzing generalization of our mistake analysis to ImageNetV2, analyzing class confusions, comparing model performance to human performance, and analyzing model errors through the lens of training data in Appendix D.

## 4. Recommendations and Discussion

In this section we provide recommendations for future evaluation, starting with ImageNet-M: our curated multi-label evaluation set of "major mistakes" that we suggest should be reported on in addition to metrics such as top-1 and multi-label accuracy.

### 4.1. ImageNet-M: A "major mistakes" evaluation split

Studies over the last few years seeking to understand whether ImageNet remains an informative benchmark have typically concluded that aspects of ImageNet remain useful but top-1 accuracy less so. These works encourage alternative related metrics for ImageNet such as multi-label accuracy (Beyer et al., 2020; Shankar et al., 2020) or object vs. organism breakdowns (Shankar et al., 2020). In many cases, stronger, more generalizable models often continue to incrementally but inevitably improve on these metrics.

With an emphasis on understanding which long-tail errors are unambiguously errors, we suggest that a benchmark focusing on the most egregious long-tail errors could provide a useful additional signal about whether the improvements are meaningful. In particular, we desire a long-tailed benchmark where we believe that 100% accuracy is achievable. Because the minor mistakes are more subject to interpretation and discussion than the major mistakes, we believe a benchmark focused on the latter will help the community judge what is meaningful improvement on ImageNet.

To that end, we leverage our expert-reviewed analysis to produce a small slice of the ImageNet multi-label set where (1) today's *best* top-1 models are still more wrong than right, and (2) the mistakes are largely unambiguous to a human given a reasonable understanding of the ImageNet label set. We call this evaluation slice **ImageNet-Major**.

**ImageNet-M example selection method.** The ViT-3B model made 155 "major" mistakes, for which we analyzed whether each example was labeled correctly for three additional models: (1) the Greedy Soups model, (2) a model pre-trained on Instagram data but fine-tuned on ImageNet that achieves 85.4% top-1 (Mahajan et al., 2018), and (3) A zero-shot evaluation (Radford et al., 2021; Jia et al., 2021; Pham et al., 2021) using a CoCa (Yu et al., 2022) model pretrained on JFT and noisy image-text data. In order to maximize prediction diversity, we purposefully selected models with varying pre-training data and training methodologies, including a zero-shot model that does not see ImageNet image-label associations directly (Gontijo-Lopes

## When does dough become a bagel? Analyzing the remaining mistakes on ImageNet

Model	Dataset	Categories				Severities	
		Fine-grained (FG)	FG w/ OOV	Spurious	Non-prototypical	Major	Minor
ViT-3B	ImageNet	64.0%	14.3%	13.8%	7.9%	40.7%	59.3%
ViT-3B	ImageNetV2	66.0%	9.4%	15.1%	9.4%	41.5%	58.5%
Greedy Soups	ImageNet	69.1%	10.4%	12.9%	7.6%	35.7%	64.3%

Table 2: **Mistake Category and Severity:** We classified the majority of the ViT-3B mistakes as fine-grained or a variant of fine-grained; many of the mistakes were considered "major" mistakes; these distributions held on ImageNetV2 as well as for the Greedy Soups model.

et al., 2021).

From this suite of four models, we assembled a subset where three or more of the models make a major mistake, yielding 68 such major mistakes. This process is similar in spirit to ImageNet-A (Hendrycks et al., 2019), except we use four high-performing but diverse models, and we restrict the set to ImageNet images and the corresponding model prediction that were rated as "major" errors. We analyzed the predictions of all models on these examples (including any novel predictions made by these additional models) and verified that none of them were correct new multi-labels, and that any model’s mistakes were major mistakes. In addition, we attempted to comprehensively manually label additional labels that no model has yet predicted but would be correct, in an attempt to reduce the likelihood that future models are penalized for making novel but unreviewed correct predictions.

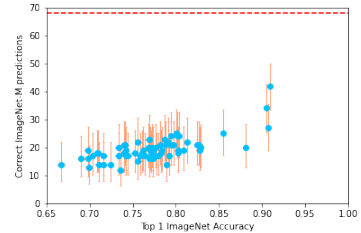
We design the ImageNet-M 68-example subset as an additional split of the validation set that we believe has the following properties: (a) many top performing models trained in different ways make mistakes on this set; (b) the mistakes are all major mistakes as determined by expert-reviewers; (c) the example set is small enough to permit manual inspection by model evaluators; (d) strives to be comprehensively labeled with respect to the ImageNet label set; (e) in theory provides a subset that *future models could achieve perfect accuracy on* without worrying about underspecified class definitions. We anticipate stronger models will make correct unreviewed predictions on this slice, so we will endeavor to update the set of multi-labels as needed.

**Evaluation.** By construction, our ViT-3B model achieves 0% accuracy on ImageNet-M; the Instagram-pretrained model gets 9 of the 68 correct, while the Greedy Soups model gets 19 correct. The zero-shot model gets the best performance with 24 correct, even though the zero-shot model overall achieves lower multi-label accuracy (94.2%) than any of the other models. Because we use these four models to help choose the mistakes, these specific models comparatively will perform poorly on this benchmark. Like with ImageNet-A (Hendrycks et al., 2019), a dataset chosen using models may bias selection in such a way that future models may easily get very high accuracy on this subset, though we try to mitigate this effect by using multiple differently-trained models in our selection criteria.

How do models that were not used to select this dataset perform? We evaluate the suite of 70 models from Shankar et al. (Shankar et al., 2020) on this dataset, in addition to four recent top models

not directly used to help filter the ImageNet-M set: a ViT-G/14 model (Zhai et al., 2021) (90.5% top-1), a BASIC model (Pham et al., 2021) fine-tuned on ImageNet (90.7% top-1), an ALIGN model (Jia et al., 2021) fine-tuned on ImageNet (88.1% top-1), and a CoCa model (Yu et al., 2022) fine-tuned on ImageNet (91.0% top-1). The plot shown here shows that most models as far back as AlexNet through ResNets get between 10-25 examples correct, but recent high accuracy models such as ViT-G/14, BASIC-FT, and CoCa-FT are starting to solve more of these ‘major’ mistakes: CoCa-FT gets 42 of the 68 examples correct. We reviewed the mistakes made by these four models, which yielded a total of 5 novel predictions; 4 of them were verified to be wrong (and major), and 1 additional new valid prediction, for which we updated the label set accordingly. We note that future models may predict classes on these examples that are "minor" mistakes, since the definition of severity is linked to the (prediction, example) pair; should it be useful, the dataset slice can be augmented with a ‘minor\_wrong\_multi\_label’ attribute to provide more fine-grained signals.

Although we did not find statistical evidence that stronger models solve major errors first, we hope that progress on image classification can be evaluated by whether improvements help reduce egregious mistakes before focusing on nebulous ones. Overall, we encourage reporting on several slices, including ImageNet-M, to give a better sense of the strengths and weaknesses of various models. For a more thorough discussion of the limitations of our analysis see Appendix C.



Trend of models on top-1 ImageNet vs. ImageNet-M, using Clopper-Pearson intervals. Red dashed line indicates total number of images in ImageNet-M as an upper-bound.



## References

- Andreassen, A., Bahri, Y., Neyshabur, B., and Roelofs, R. The evolution of out-of-distribution robustness through-out fine-tuning, 2021. <https://arxiv.org/abs/2106.15831>.
- Baker, D. Datasets have worldviews. <https://pair.withgoogle.com/explorables/dataset-worldviews>. Accessed: 2022-04-28.
- Barz, B. and Denzler, J. Do we train on test data? purging cifar of near-duplicates. *Journal of Imaging*, 6(6):41, 2020.
- Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X., and Oord, A. v. d. Are we done with imagenet? *CVPR*, 2020. <https://arxiv.org/abs/2006.07159>.
- Ciregan, D., Meier, U., and Schmidhuber, J. Multi-column deep neural networks for image classification. In *2012 IEEE conference on computer vision and pattern recognition*, pp. 3642–3649. IEEE, 2012.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009. <https://ieeexplore.ieee.org/document/5206848>.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Friedler, S. A., Scheidegger, C., and Venkatasubramanian, S. On the (im) possibility of fairness. *arXiv preprint arXiv:1609.07236*, 2016.
- Geirhos, R., Meding, K., and Wichmann, F. A. Beyond accuracy: quantifying trial-by-trial behaviour of cnns and humans by measuring error consistency. *NeurIPS*, 2020. <https://arxiv.org/abs/2006.16736>.
- Gontijo-Lopes, R., Dauphin, Y., and Cubuk, E. D. No one representation to rule them all: Overlapping features of training methods, 2021. <https://arxiv.org/abs/2110.12899>.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*, 2019.
- Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. *CoRR*, abs/1907.07174, 2019. URL <http://arxiv.org/abs/1907.07174>.
- Hendrycks, D., Basart, S., Mu, N., Kadavath, S., Wang, F., Dorundo, E., Desai, R., Zhu, T., Parajuli, S., Guo, M., Song, D., Steinhardt, J., and Gilmer, J. The many faces of robustness: A critical analysis of out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.
- Hooker, S., Courville, A., Clark, G., Dauphin, Y., and Frome, A. What do compressed deep neural networks forget? *arXiv preprint arXiv:1911.05248*, 2019.
- Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., Le, Q., Sung, Y.-H., Li, Z., and Duerig, T. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pp. 4904–4916. PMLR, 2021.
- Kolesnikov, A., Beyer, L., Zhai, X., Puigcerver, J., Yung, J., Gelly, S., and Houlsby, N. Big transfer (bit): General visual representation learning. In *European conference on computer vision*, pp. 491–507. Springer, 2020.
- Kornblith, S., Shlens, J., and Le, Q. V. Do better imagenet models transfer better? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Lee, H. S., Agarwal, A. A., and Kim, J. Why do deep neural networks still not recognize these images?: A qualitative analysis on failure cases of imagenet classification. *arXiv preprint arXiv:1709.03439*, 2017. URL <https://arxiv.org/abs/1709.03439>.
- Mahajan, D., Girshick, R., Ramanathan, V., He, K., Paluri, M., Li, Y., Bharambe, A., and Van Der Maaten, L. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 181–196, 2018.
- Mania, H. and Sra, S. Why do classifier accuracies show linear trends under distribution shift? *arXiv preprint arXiv:2012.15483*, 2020. <https://arxiv.org/abs/2012.15483>.
- Mania, H., Miller, J., Schmidt, L., Hardt, M., and Recht, B. Model similarity mitigates test set overuse. *NeurIPS*, 2019. <https://arxiv.org/abs/1905.12580>.
- Miller, G. A. Wordnet: A lexical database for english. *Communications of the ACM*, 38:39–41, 1995.
- Nguyen, T., Raghu, M., and Kornblith, S. Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. In *International Conference on Learning Representations (ICLR)*, 2021. URL <https://arxiv.org/abs/2010.15327>.

- Northcutt, C., Jiang, L., and Chuang, I. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021a.
- Northcutt, C. G., Athalye, A., and Mueller, J. Pervasive label errors in test sets destabilize machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021b. <https://arxiv.org/abs/2103.14749>.
- Paullada, A., Raji, I. D., Bender, E. M., Denton, E., and Hanna, A. Data and its (dis) contents: A survey of dataset development and use in machine learning research. *Patterns*, 2(11):100336, 2021.
- Pham, H., Dai, Z., Ghiasi, G., Liu, H., Yu, A. W., Luong, M.-T., Tan, M., and Le, Q. V. Combined scaling for zero-shot transfer learning. *arXiv preprint arXiv:2111.10050*, 2021.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Recht, B., Roelofs, R., Schmidt, L., and Shankar, V. Do imagenet classifiers generalize to imagenet? In *International Conference on Machine Learning*, pp. 5389–5400, 2019.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Li, F.-F. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 2015. <https://arxiv.org/abs/1409.0575>.
- Shankar, V., Roelofs, R., Mania, H., Fang, A., Recht, B., and Schmidt, L. Evaluating machine accuracy on imagenet. In *International Conference on Machine Learning (ICML)*, 2020. <http://proceedings.mlr.press/v119/shankar20c/shankar20c.pdf>.
- Shen, J., Nguyen, P., Wu, Y., Chen, Z., Chen, M. X., Jia, Y., Kannan, A., Sainath, T., Cao, Y., Chiu, C.-C., et al. Lingvo: a modular and scalable framework for sequence-to-sequence modeling. *arXiv preprint arXiv:1902.08295*, 2019.
- Stock, P. and Cissé, M. Convnets and imagenet beyond accuracy: Explanations, bias detection, adversarial examples and model criticism. *CoRR*, abs/1711.11443, 2017. URL <http://arxiv.org/abs/1711.11443>.
- Sun, C., Shrivastava, A., Singh, S., and Gupta, A. Revisiting unreasonable effectiveness of data in deep learning era. In *International Conference on Computer Vision (ICCV)*, 2017. <https://arxiv.org/abs/1707.02968>.
- Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tsui, P., Guo, J., Zhou, Y., Chai, Y., Caine, B., et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2446–2454, 2020.
- Torralba, A. and Efros, A. A. Unbiased look at dataset bias. In *CVPR 2011*, pp. 1521–1528. IEEE, 2011.
- Tsipras, D., Santurkar, S., Engstrom, L., Ilyas, A., and Madry, A. From imagenet to image classification: Contextualizing progress on benchmarks. *CoRR*, abs/2005.11295, 2020. URL <https://arxiv.org/abs/2005.11295>.
- Van Horn, G., Branson, S., Farrell, R., Haber, S., Barry, J., Ipeirotis, P., Perona, P., and Belongie, S. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *CVPR, 2015*. [https://vision.cornell.edu/se3/wp-content/uploads/2015/05/Horn\\_Building\\_a\\_Bird\\_2015\\_CVPR\\_paper.pdf](https://vision.cornell.edu/se3/wp-content/uploads/2015/05/Horn_Building_a_Bird_2015_CVPR_paper.pdf).
- Wang, H., Ge, S., Lipton, Z., and Xing, E. P. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pp. 10506–10518, 2019.
- Wortsman, M., Ilharco, G., Gadre, S. Y., Roelofs, R., Gontijo-Lopes, R., Morcos, A. S., Namkoong, H., Farhadi, A., Carmon, Y., Kornblith, S., et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. *arXiv preprint arXiv:2203.05482*, 2022.
- Xu, Y., Lee, H., Chen, D., Hechtman, B., Huang, Y., Joshi, R., Krikun, M., Lepikhin, D., Ly, A., Maggioni, M., et al. Gspmd: general and scalable parallelization for ml computation graphs. *arXiv preprint arXiv:2105.04663*, 2021.
- Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., and Wu, Y. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. URL <https://arxiv.org/abs/2205.01917>.
- Yun, S., Oh, S. J., Heo, B., Han, D., Choe, J., and Chun, S. Re-labeling imagenet: from single to multi-labels, from global to localized labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2340–2350, 2021. <https://arxiv.org/abs/2101.05022>.
- Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. Scaling vision transformers. *arXiv preprint arXiv:2106.04560*, 2021. URL <https://arxiv.org/abs/2106.04560>.

## A. Related Work

**Multi-label annotations on ImageNet.** As models continue to improve ImageNet top-1 accuracy, there has been an increased interest in evaluating ImageNet multi-label accuracy (Stock & Cissé, 2017; Beyer et al., 2020; Shankar et al., 2020; Yun et al., 2021; Tsipras et al., 2020). Stock et al. (Stock & Cissé, 2017) use non-expert human studies and explanations (e.g., model criticisms) on predictions of a then-SOTA model on ImageNet, finding that machine accuracy is underestimated and advocating for multi-label evaluation. Beyer et al. (Beyer et al., 2020) introduce a set of Reassessed Labels (ReaL) for the ImageNet validation set containing multi-label annotations. The researchers first collected proposal labels from model predictions using a testbed of 19 models, and then, in order to reduce the number of predictions to review, they narrowed the set of models down to the 6 models that had the highest precision and recall above 97% on a set of 256 images reviewed by 5 experts. The top predictions from these 6 models were then reviewed by human annotators from a crowdsourcing platform. In a similar vein, Shankar et al. (Shankar et al., 2020) provide a multi-label annotation dataset for 20,000 of the 50,000 ImageNet validation images and find that roughly 20% of images have more than one valid label. To generate the annotations, the researchers first collected predictions from a testbed of 70 models for each image and then reviewed the unique model predictions using a panel of three human experts. Additionally, human accuracy was evaluated on a subset of 1,000 images and a panel of 5 human experts reviewed human and model predictions on this subset. Tsipras et al. (Tsipras et al., 2020) also find that 20% of images in the validation set contain objects from multiple classes, and identify sources of ambiguous label classes. Hooker et al. (Hooker et al., 2019) use human evaluations to label a subset of examples from ImageNet, finding that they contained multiple labels 40–60% of the time. More recently, Yun et al. (Yun et al., 2021) obtained pixel-level multi-label ground truths for the ImageNet *training set* using a machine annotator, and found that training a model with the dense annotations leads to small improvements in both top-1 and multi-label ImageNet accuracy.

Our paper focuses on (and adds to) the multi-label evaluation of ImageNet using expert labelers, updating the dataset to collapse classes that are overlapping or subset relationships to better capture the remaining real errors. We adopt the labeling methodology of Shankar et al. (Shankar et al., 2020) and use a panel of human experts to determine the validity of novel predictions. In contrast to prior work, we re-visit all remaining model mistakes using human expert review, and categorize them by type and severity.

**ImageNet Label Error.** Multi-label annotation datasets for ImageNet (including our updated annotations) identify a set

of images that have no correct ground truth label (i.e. label errors). Van Horn et al. (Van Horn et al., 2015) used experts from the Cornell Lab of Ornithology to estimate that at least 4% of the bird images are misclassified, and more recently Northcutt et al. (Northcutt et al., 2021b) used MTURK workers to review algorithmically identified potential errors and found a label error rate of 5.83% in ImageNet. Lee et al. (Lee et al., 2017) sample 400 ImageNet mistakes and perform a categorization, similarly finding significant label error and label ambiguity.

**Mistake analysis.** Recent work has sought to understand or analyze model mistakes on ImageNet, mainly focusing on similarities or differences between mistake sets of independently trained classifiers. For example, Mania et al. (Mania et al., 2019) found that predictions of different models on ImageNet are more similar than one would expect if the models were making mistakes independently. Geirhos et al. (Geirhos et al., 2020) studied a 16-class ImageNet classification task and similarly found that CNNs make remarkably consistent mistakes with one another but CNNs and humans have an error consistency that is only slightly above what can be expected from chance. In follow up work, Mania et al. (Mania & Sra, 2020) studied the *dominance probabilities* for pairs of ImageNet models, which capture the probability that a higher accuracy model will make a mistake on a particular image that a lower accuracy model correctly classifies. Their empirical analysis of dominance probabilities on ImageNet implies that the mistakes of higher accuracy models are typically subsets of the mistakes of lower accuracy models. However, Gontijo-Lopes et al. (Gontijo-Lopes et al., 2021) and Andreassen et al. (Andreassen et al., 2021) found that training on larger and more diverse datasets as well as zero-shot evaluation can lead to models with more mistake diversity, which in turn can be used to build more accurate ensembles. Similarly, Nguyen et al. (Nguyen et al., 2021) found systematic differences in the errors between wide and deep ResNets on ImageNet. Ciregan et al. (Ciregan et al., 2012) demonstrated that powerful models on MNIST allowed for remaining error analysis in the single-label context, highlighting the ambiguity in many of those remaining errors.

In contrast to this prior work, we take a step towards fully calibrating SOTA model evaluations by exhaustively and visually reviewing every remaining model mistake using manual expert review in the *multi-label context*; we believe the remaining errors identified to be legitimately incorrect and tempered by severity and category ratings that we hope prove useful to future evaluations.

## B. Methods

### B.1. Panel Review

To exhaustively and accurately assess every remaining mistake, we formed a panel of five reviewers and followed a process similar to (Shankar et al., 2020) to evaluate the predictions made by this model on the 676 mistakes — we avoided using non-expert crowd-source platforms specifically because the remaining mistakes are often difficult to assess by non-expert annotators (Tsipras et al., 2020; Beyer et al., 2020). For every mistake, the panel determined: (1) Did the model make a mistake? (2) Was the original ground truth annotation correct? (3) If the model made a mistake, what is the category, type, and severity of the mistake? The `imagenet2012_multilabel` dataset contains a field for every image indicating which classes a large previous suite of models predicted that were determined to be incorrect (`wrong_multi_labels`). Of these 676 initial mistakes, 221 were novel: they were not reviewed in the original multi-label annotation process since none of the models evaluated made the same prediction. Each member of the panel reviewed all 221 novel mistakes.

Similar to (Shankar et al., 2020), we built a review tool that allowed each panelist to see a) the predicted class, b) the predicted top softmax score, c) the set of ground-truth labels, d) the set of previously incorrect labels, and e) the image. We also employed the labeling guide produced by the authors of (Shankar et al., 2020) when investigating the definition of a class, and a tool to iterate through the images of every ImageNet validation example for that class, using the validation images to help define the class boundaries (rather than the gloss or lay definition of the class). See Appendix F for screenshots of the review tools. In addition, we collapsed a small number of classes for which previous work has identified exhibited extreme overlap (Northcutt et al., 2021a), such as ‘missile’ and ‘projectile missile’. In Appendix G we provide these new collapsed class mappings.

We also used Google Image Search to help provide context to some assessments; in one interesting but not isolated case, a prediction of a taxi cab (with no obvious taxi cab indicators beyond yellow color) was present in the image; we determined the prediction to be correctly a taxi cab and not just a standard vehicle by identifying a landmark bridge in the background in order to localize the city, and a subsequent image search for taxis in that city yielded the images of the same taxi model and license plate design, validating the model’s actually correct prediction.

Each panelist rated whether these novel mistakes had a mislabeled ground truth, or whether the prediction should be added to the set of correct, unclear, or wrong multi-labels. As a group, we reviewed any image where there was no unanimous agreement, allowing those in the minority to

make their case and change minds, or highlight potential oversights. The use of a panel and discussion was important: in a non-trivial number of cases, a single panelist found unique evidence for or against a prediction that no other panelist saw that led to a different outcome. After locking in final votes, we took the majority assessment, or used ‘unclear’ for a tie. After this re-assessment, 140 of the novel predictions were deemed correct (or the ground truth deemed incorrect), leaving 536 remaining mistakes to assess.

### B.2. Mistake severity and category

The remaining mistakes then comprised images that had either previously been deemed wrong by the panel in Shankar et al. (Shankar et al., 2020), or deemed wrong by the current panel. We then began a review of the mistakes for the category and severity of mistakes. During this second phase review, we re-reviewed all images in detail, potentially overturning decisions made by the previous panel, where they missed the presence of an object that was in an example or the decision was inconsistent with the labeling guide or validation examples we relied on.<sup>1</sup>

**Severity:** We assessed each mistake’s severity with the assumption that not all mistakes are equal: some mistakes are extremely borderline, particularly because the ImageNet class definitions are imprecise or because the image itself provides ambiguous or incomplete information. For any remaining mistake, we broke down the severity levels into **major** and **minor** mistakes. A **major mistake** is a model prediction that a human who understands the class definitions would find obviously incorrect. For example Figure 1(a) shows an example image where the prediction is jigsaw puzzle but the label is dough; although the pieces are somewhat jigsaw-puzzle-shaped, untrained humans are more likely to classify this as dough than jigsaw puzzle. A **minor mistake** is a model prediction that a human who understands the class definitions would probably find to be incorrect, but in a more subtle way than a major mistake. Some minor mistakes are so subtle that even expert-trained humans might debate their correctness.

We recognize that these severities are subject to the influences of the worldview of the panelists (and the web) and should be judged accordingly (see Section C). For transparency, we provide the panel assessments of the severity for all the mistakes, and in Figure 1 we provide some examples of the severities of various mistakes made by this model for the reader’s calibration. In Appendix H we provide many additional examples of every severity and category.

<sup>1</sup>The labeling guide in Shankar et al. (Shankar et al., 2020) was constructed after the initial panel review but prior to the human accuracy assessment, resulting in some validation labels that would be inconsistent with the labeling guide we used.





(a) **Major mistake**  
 Label: dough  
 Model: jigsaw puzzle



(b) **Minor mistake**  
 Label: kuvasz  
 Model: Great Pyrenees



(c) **New multilabel**  
 Label: tape player  
 Model: cassette



(d) **Problematic**  
 Label: bee  
 Model: fly

Figure 1: **Mistake Severity.** Examples of the two mistake severities (a-b), a correct model prediction where the model identifies a previously missing multi-label (c); and a problematic example (d) where the label (bee) is incorrect (object is a bee-fly, which is a type of fly).



(a) **Fine-grained**  
 Label: wall clock  
 Model: sundial



(b) **Fine-grained w/ OOV**  
 Label: syringe  
 Model: hamster



(c) **Spurious correlation** Label: mouse, desk, monitor, screen  
 Model: desktop computer



(d) **Non-prototypical**  
 Label: stove  
 Model: hamper

Figure 2: **Mistake Category.** Examples of the four mistake categories. In the fine-grained with OOV example, the animal is a chinchilla, which is not an ImageNet class but is visually similar to a hamster, which is an ImageNet class. In the spurious correlation example, the scene contains relevant context for desktop computer, but there is no such object in the image.

**Category:** After reviewing many mistakes, we formulated four mistake types (Figure 2).

**(1) Fine-grained errors** are where the model makes a mistake between two similar types of organisms or objects, one of which is a groundtruth label. These mistakes often occur when the two confused classes are already similar (e.g. two dog breeds), or when either of them is very broad or ill-scoped (e.g. the “bake shop” class includes any baked good or bakery).

**(2) Fine-grained with out-of-vocabulary:** there is an object in the image that is not in the ImageNet class hierarchy but is similar to a predicted class that is in ImageNet. We separate out this category because it highlights the possible benefits of training models to expect new classes to appear at test time, and the importance of model uncertainty and calibration in the face of this ‘open world’.

**(3) Spurious correlations:** Either (a) the predicted object is not plausibly in the image but surrounding cues may have been used, or (b) the predicted object does not match the groundtruth. In extreme cases, there is no clear indication of the predicted object; in more subtle cases, it is more clear the model is trying to predict the class of an object in the image but the predicted object would not be considered either semantically or visually similar to the groundtruth class.

**(4) Non-prototypical labels:** The predicted label is not present but the groundtruth object is a non-prototypical example of the class that bears resemblance to the predicted label. Non-prototypical mistakes are relatively rare, and capture the ‘long tail’ of examples for each class. These mistakes highlight the internal diversity of each class, and the difficulty of modeling the long within-class tail.

## C. Limitations of Analysis

Much of our analysis on mistake categorization, severity, and data cleaning depends greatly on qualitative factors determined by the authors and experiment design, which we briefly discuss here.

**Limitations due to mistake subset.** We only reviewed multi-label annotation examples where the ViT-3B model and Greedy Soups model we chose were incorrect; while the multi-label dataset itself has undergone review of mistakes from a suite of many models, we never review any validation image whose groundtruth might be wrong, if all models evaluated also make the incorrect groundtruth prediction. We do partially review some of the mistakes made by a few other models, and build upon a dataset where mistakes made by many other models have already been reviewed, but most of the work here assesses these two specific model’s predictions.

**Definition of a mistake.** We re-iterate that we used qualitative judgments to decide whether a prediction was a mistake, and if so, its severity and categorization. Our qualitative judgments are therefore based on a biased worldview (Torralba & Efros, 2011; Friedler et al., 2016; Baker; Paullada et al., 2021) comprising the five panelists; moreover, we are not world experts on dog or animal species, though we believe our assessments are at least as good or better than the original labeling process used for the validation set, given the research effort we made on each mistake. As a mitigation against imbuing too biased a worldview on class definitions, we relied heavily on the validation data to define the boundaries of the class, even when those examples did not match up with our personal definitions of the class. Moreover, updating (and evaluating) multi-labels potentially allows for different world-views to be expressed.

## D. Additional Results

### D.1. Out-of-distribution generalization

We evaluated the ViT-3B model on the ImageNetV2 multi-label subset which produced over 900 unreviewed errors. To assess what aspects of our analysis generalize to other datasets, we sampled 100 of these errors using the same panel review system employed for ImageNetV1. We discovered that 47 of the 100 ImageNetV2 predictions were either correct or had problematic labels, leaving 53 mistakes that we reviewed for category and severity. For ImageNetV1, we had previously found that 44% (296/676) of mistakes were either problematic or correct and we found no statistically significant difference between the two datasets in this regard ( $\chi^2(1, N = 776) = 0.36, p = .55$ ). These proportions suggest that large models are frequently uncovering new correct multi-labels, suggesting that mistake analysis and label correction needs to be part of the lifecycle and maintenance of benchmark development of long-tailed errors to properly assess performance as a benchmark saturates. Table 2 compares the category and severity breakdowns between the two datasets—overall the model is making similar types and severities of mistakes on both datasets. A chi-square test of independence shows that there is no significant difference between either the mistake categorization ( $\chi^2(3, N = 434) = 0.97, p = .80$ ) or mistake severity ( $\chi^2(1, N = 434) = 0.01, p = .92$ ).

### D.2. Analyzing class confusions.

Given that most failures were fine-grained, we tried to identify any patterns present in the class confusions made by the model, but found no consistent pattern. The inline figure shows the frequency of occurrence for confused class pairs. The distribution is long-tailed in nature—the majority of class pairs occur exactly once or twice and only a handful of class pairs occur three or more times. The long-tailed nature

of the class confusions suggests that we will not be able to resolve a large fraction of model mistakes by focusing on cleaning up or adding additional data to only a small number of classes.

### D.3. Comparison to humans

We compare the performance of the ViT-3B and Greedy Soups models to the best human labeler from Shankar et al. (Shankar et al., 2020)<sup>2</sup> by evaluating on the subset of 1,000 ImageNet val images used to compute human accuracy in the prior work. To fairly compare to the models, we re-compute the human accuracy scores using the original human predictions and our updated label set. Overall, the re-labeling did not significantly change human accuracy; the best human labeler achieved 97.3% MLA on the original multi-labels and 97.4% MLA on the updated multi-labels.

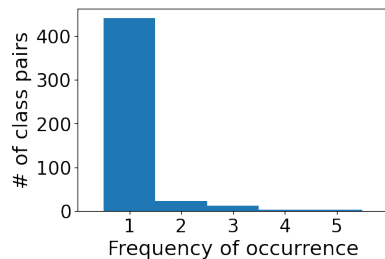


Table 3 compares the re-labeled MLA of the ViT-3B and Greedy Soup model to the best human for all ImageNet classes as well as the subset of ImageNet classes corresponding to objects and organisms. Similar to Shankar et al. (Shankar et al., 2020), we also find that the performance of the models is more uniform across the object and organism classes, but humans do substantially better on the object classes than the organism classes. However, unlike prior work, current models outperform the best human when evaluated on all ImageNet classes (though humans still achieve slightly better performance on the object classes).

### D.4. Analyzing the Training Data

Finally, we investigate how much we can understand model validation set errors through the lens of the training data. To do so, we inspect the  $K = 10$  nearest neighbor training images using JFT-pretrained (before ImageNet fine tuning) embeddings for the ViT-3B model. Doing this, we rediscover (originally documented in Sun et al. (Sun et al., 2017) and Kolesnikov et al. (Kolesnikov et al., 2020)) that 797 (1.59%) ImageNet validation images exist in the training set as **exact duplicates**. Interestingly, we are the first to notice that **every single leaked image has a different label in the training set than the validation set**. In Appendix J.1 we show removing these images has relatively little impact on the model’s performance, and detail a more pernicious leakage pattern of "near duplicates" in Appendix J.2 that is

<sup>2</sup>Best human is chosen based on MLA on all classes (as opposed to object or organism subset) and corresponds to Human E in the original work.

hard to fully quantify. Finally, we show how the training data sometimes explain spurious correlations in Appendix J.3.

## E. ViT-3B model details

The ViT model we use in this work is based on a standard Vision Transformer (Dosovitskiy et al., 2021) model scaled to nearly 3 billion parameters, using a patch size of 14, 16 heads, 64 blocks, an MLP dimension of 8192 and a hidden dimension of 2048. The model is defined and trained in Lingvo (Shen et al., 2019); we additionally employ GSPMD (Xu et al., 2021) for training. The model is pre-trained on JFT-3B (Sun et al., 2020) using training settings that optimize for performance on JFT-3B rather than for fine-tuning on ImageNet; notably, we do not use the training recipe that helps few-shot transfer performance (Zhai et al., 2021). For fine-tuning on ImageNet, we use the AdamW optimizer (beta1=0.9, beta2=0.999, epsilon=1e-8, weight\_decay=0.3) with a cosine learning rate schedule (max learning rate of 3e-2, warmup of 2k steps, final rate of 3e-4), a training batch size of 512, and fine-tune for a total of 10 epochs.

## F. Review tools

We include screenshots of the reviewing tools we built to analyze model mistakes. Figure 3 shows the UI for reviewing model predictions and Figure 4 shows the UI that displays the labeling guide and slide bar to browse images for a particular class.

When does dough become a bagel? Analyzing the remaining mistakes on ImageNet

	All Classes	Organisms	Objects
ViT-3B	97.8%	97.4%	98.1%
Greedy Soup	98.6%	98.4%	98.8%
Best Human (Shankar et al., 2020)	97.4%	95.4%	99.0%

Table 3: **Multi-label accuracy compared to humans.** Both the ViT-3B and Greedy Soup model achieve better MLA on all ImageNet classes than the best human labeler from Shankar et al. (Shankar et al., 2020). However, on the object classes, where the class boundaries are substantially less ambiguous for humans, the best human labeler still outperforms the models.

ILSVRC2012\_val\_00001327.JPEG prediction: ['jean', 'blue jean', 'denim']:

- unclear
- correct
- wrong
- is\_problematic
- unreviewed

Notes:

Reviewed by panel:

- ILSVRC2012\_val\_00001327.JPEG

Top predicted class: ['jean', 'blue jean', 'denim'] ((usually plural) close-fitting trousers of heavy denim for manual work or casual wear)  
 Top prediction score: 0.6165266633033752

Other predictions: [['seashore', 'coast', 'seacoast', 'sea-coast'], ['sandbar', 'sand bar'], ['sandal'], ['maillot']]

Correct labels:

- ['sandbar', 'sand bar'] ( a bar of sand )
- ['seashore', 'coast', 'seacoast', 'sea-coast'] ( the shore of a sea or ocean )

Wrong labels:

- ['backpack', 'back pack', 'knapsack', 'packsack', 'rucksack', 'haversack'] ( a bag carried by a strap on your back or shoulder )
- ['binoculars', 'field glasses', 'opera glasses'] ( an optical instrument designed for simultaneous use by both eyes )
- ['cowboy boot'] ( a boot with a high arch and fancy stitching; worn by American cowboys )
- ['go-kart'] ( a small low motor vehicle with four wheels and an open framework; used for racing )
- ['knee pad'] ( protective garment consisting of a pad worn by football or baseball or hockey players )
- ['motor scooter', 'scooter'] ( a wheeled vehicle with small wheels and a low-powered gasoline engine geared to the rear wheel )
- ['parachute', 'chute'] ( rescue equipment consisting of a device that fills with air and retards your fall )
- ['ballplayer', 'baseball player'] ( an athlete who plays baseball )



Figure 3: A screenshot of the UI we built to review model predictions. For each image, we determined whether the prediction was correct, wrong, or unclear. We also flagged images as problematic if the ground truth label for the image was incorrect.



## bakery, bakeshop, bakehouse

- The image does not have to show the building.
- Many **bakery** images have many baked goods, e.g., an entire tray of cupcakes.
- Cakes appear individually more often in this class. There is also no other class for cake, so we count cakes as **bakery**.
- More generally, if the image shows an individual baked item and there is no other suitable baked good class for it, the image counts as **bakery**.

### Display images

show\_n: 50

dataset: INetVal

[Show code](#)



a workplace where baked goods (breads and cakes and pastries) are produced or sold - Key: 2136211865742104015-n02776631-ILSVRC2012\_val\_00010307.JPEG, is\_problematic: NO



Figure 4: A screenshot of the class search tool we built that displays the labeling guide and a slider bar that allows users to browse validation images for a particular class.

## G. Collapsed mappings

We provide our collapsed class mappings that we employed unilaterally, based on determining that the classes exhibited

significant overlap based on the validation set images, or there was a strict superset relationship between two or more classes. For example, a prediction of ‘eskimo dog’ for a ‘siberian husky’ label would be considered correct, whereas a prediction of ‘siberian husky’ for an ‘eskimo dog’ label might not.

All siberian huskies and malamutes are also eskimo dogs. 250: 248, 249: 248

Sunglass and sunglasses are the same class (bidirectional). 836: 837, 837: 836

Indian and African elephants are also tuskers. 385: 101, 386: 101

A coffee mug is also a cup. 504: 968

Maillot and maillot, tanksuit are the same class (bidirectional). 638: 639, 639: 638

Missile and projectile missile are the same class (bidirectional). 657: 744, 744: 657

Notebook computer and laptop are the same class (bidirectional). 620: 681, 681: 620

Monitor and screen are the same class (bidirectional). 664: 782, 782: 664

A cassette player is also a tape player. 482: 848

Weasel, polecats, black-footed ferrets, and minks are all the same class. 356: [357, 358, 359], 357: [356, 358, 359], 358: [356, 357, 359], 359: [356, 357, 358],

All bathtubs are tubs, but not all tubs are bathtubs. 435: 876

## H. Mistake Examples by Severity



GT: tripod  
Pred: swing



GT: wardrobe, bucket, broom  
Pred: entertainment center



GT: cleaver  
Pred: hatchet



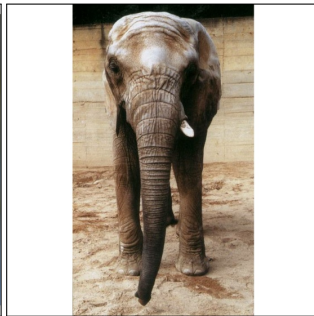
GT: agama  
Pred: frilled lizard



GT: hen  
Pred: cock



GT: jay  
Pred: magpie



GT: African elephant, tusk  
Pred: Indian elephant



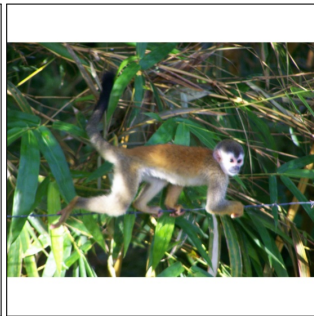
GT: Madagascar cat  
Pred: indri



GT: water snake  
Pred: ringneck snake



GT: toilet tissue, pot, window shade; Pred: paper towel



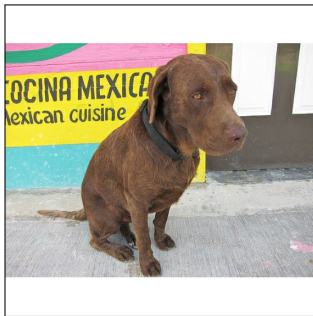
GT: squirrel monkey  
Pred: titi



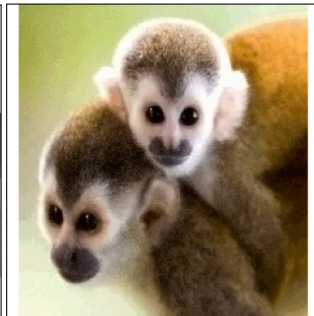
GT: reflex camera  
Pred: Polaroid camera



GT: Border collie, patio  
Pred: collie



GT: Chesapeake Bay retriever  
Pred: Labrador retriever



GT: squirrel monkey  
Pred: titi



GT: whippet  
Pred: Italian greyhound

Figure 5: **Major mistakes.** Additional examples of major mistakes. Of the correct multi-labels, the original ImageNet label is listed first.



When does dough become a bagel? Analyzing the remaining mistakes on ImageNet

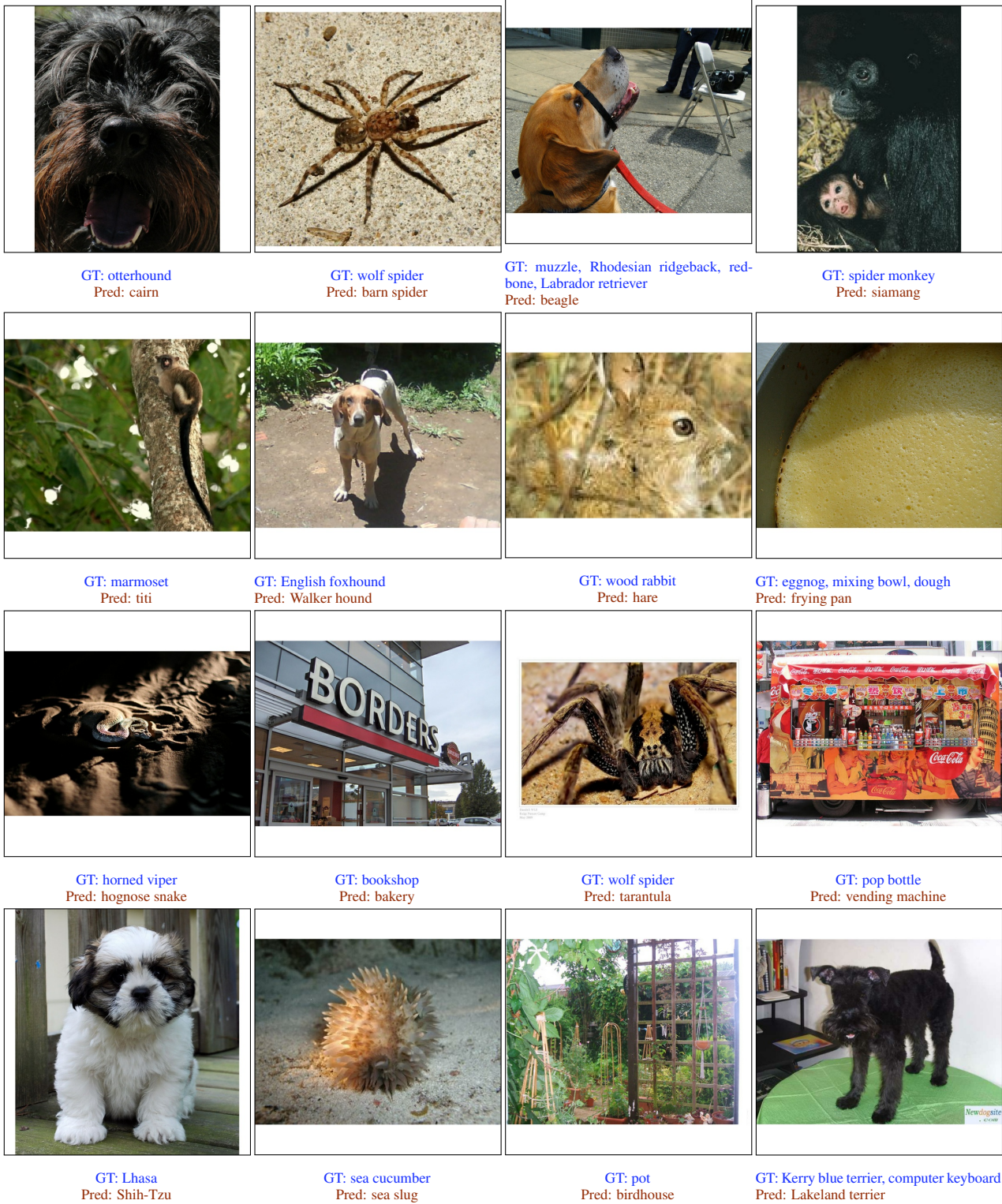
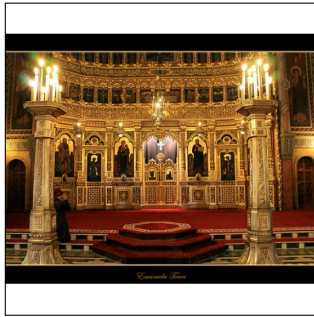


Figure 6: **Minor mistakes.** Additional examples of minor mistakes. Of the correct multi-labels, the original ImageNet label is listed first.



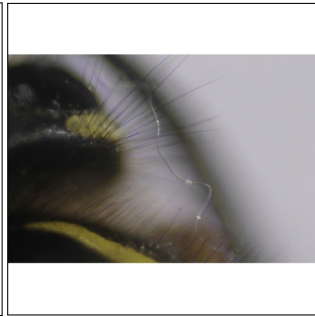
When does dough become a bagel? Analyzing the remaining mistakes on ImageNet



GT: altar, church  
Pred: church



GT: kelpie, German shepherd  
Pred: German shepherd



GT: nematode, bee  
Pred: bee



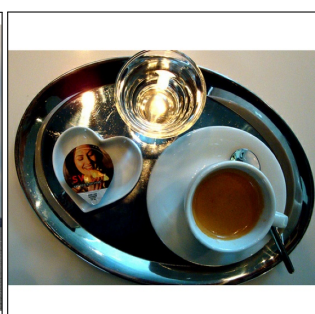
GT: wing, airliner  
Pred: airliner



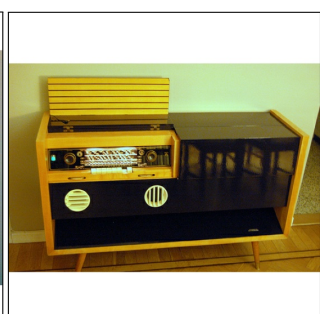
GT: cellular telephone, hand-held computer  
Pred: hand-held computer



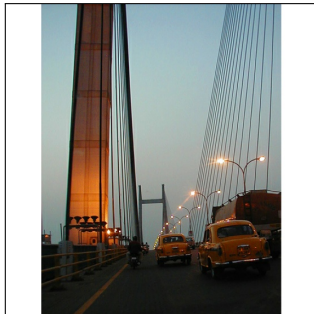
GT: suspension bridge, pier  
Pred: pier



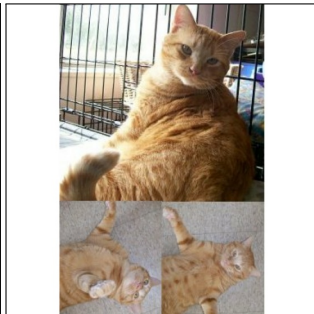
GT: tray, espresso  
Pred: espresso



GT: tape player, cassette player, radio  
Pred: radio



GT: suspension bridge, cab  
Pred: cab



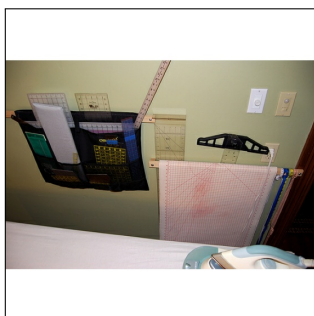
GT: tiger cat, tabby  
Pred: tabby



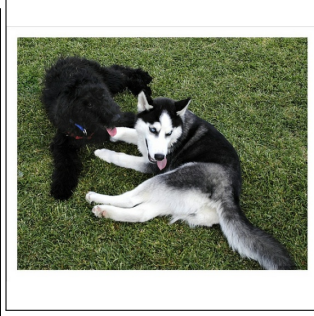
GT: ice cream, plate  
Pred: plate



GT: suspension bridge, pier  
Pred: pier



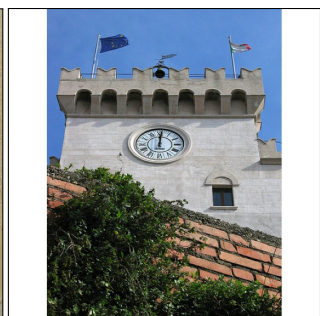
GT: rule, iron  
Pred: iron



GT: Bouvier des Flandres, Siberian husky, Eskimo dog  
Pred: Siberian husky



GT: handkerchief, velvet  
Pred: velvet



GT: wall clock, analog clock, bell cote  
Pred: bell cote

Figure 7: **Correct "mistakes"**. Additional examples where the model makes a correct prediction that we add to the original multi-label annotations. Of the original multi-labels shown, the original ImageNet label is listed first. Often the model correctly identifies a different object in the image, and in some cases a single object has ambiguous class membership and could plausibly belong to either the ground truth or the predicted class.

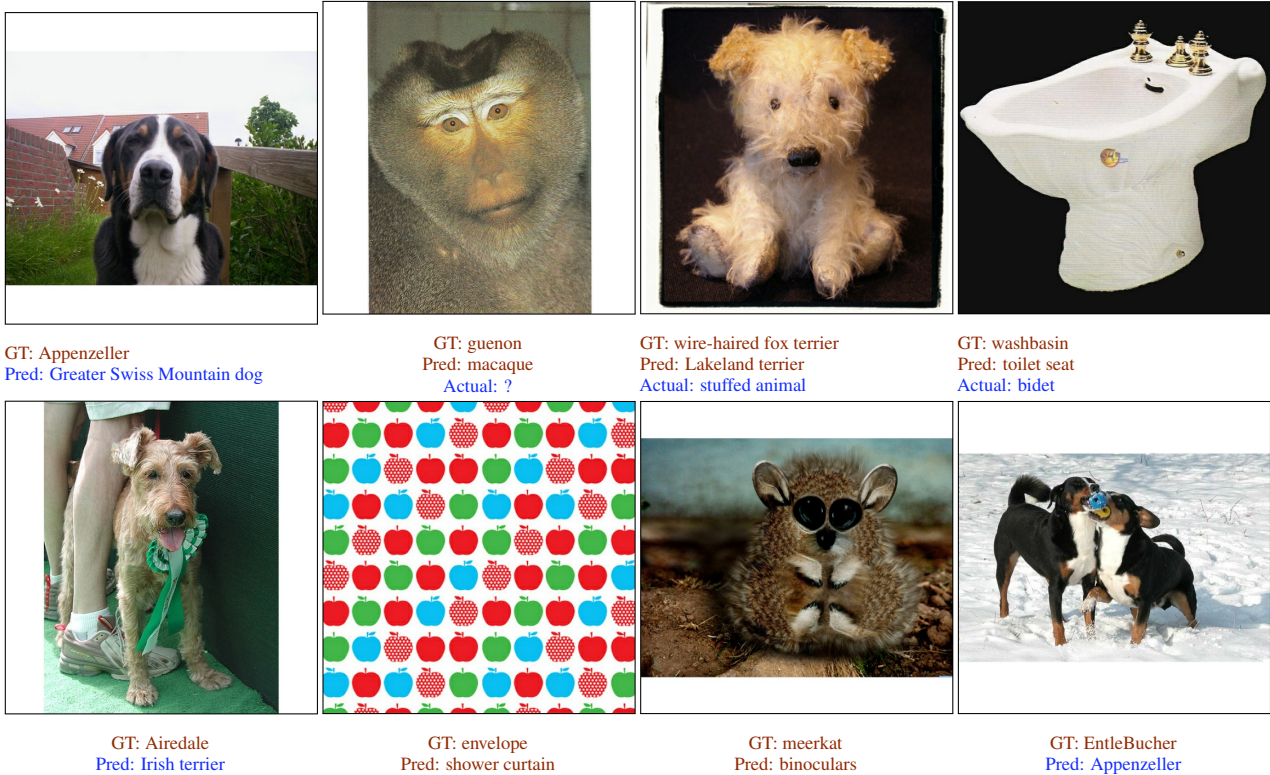
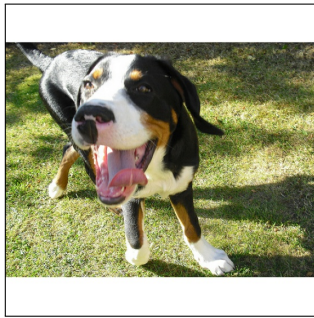


Figure 8: **Problematic “mistakes”**. Examples where our panel determined that the image or its original label was problematic (and therefore should not be in the validation set). Most problematic examples are problematic because the original ImageNet label was deemed incorrect because the prediction by the model was indeed correct.



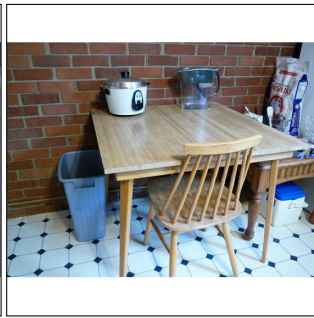
## I. Mistake Examples by Category



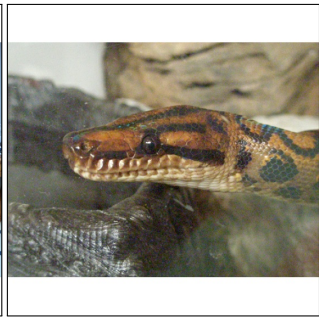
GT: Appenzeller  
Pred: EntleBucher



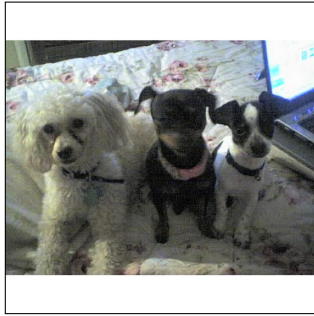
GT: measuring cup  
Pred: nipple



GT: water jug  
Pred: Crock Pot



GT: rock python  
Pred: boa constrictor



GT: miniature poodle, quilt, laptop, note-  
book  
Pred: toy poodle



GT: cleaver  
Pred: hatchet



GT: Appenzeller  
Pred: EntleBucher



GT: spotted salamander  
Pred: European fire salamander



GT: soft-coated wheaten terrier  
Pred: Irish terrier



GT: breastplate, pole  
Pred: cuirass



GT: china cabinet  
Pred: bookcase



GT: wall clock, analog clock  
Pred: sundial



GT: miniature schnauzer  
Pred: standard schnauzer



GT: wooden spoon, spatula  
Pred: ladle



GT: agama  
Pred: green lizard



GT: four-poster, crib  
Pred: cradle

Figure 9: **Fine-grained mistakes.** Additional examples of fine-grained mistakes. Of the correct multi-labels, the original ImageNet label is listed first.



When does dough become a bagel? Analyzing the remaining mistakes on ImageNet


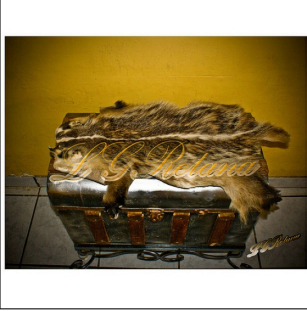














			
<p>GT: jersey Pred: crossword puzzle OOV: Calendar design</p>	<p>GT: chest Pred: badger OOV: Other animal</p>	<p>GT: sliding door, patio, studio couch, Pred: rocking chair OOV: spinning chair</p>	<p>GT: grille Pred: police van OOV: Police car</p>
			
<p>GT: cradle Pred: mosquito net OOV: Cradle cover</p>	<p>GT: cauliflower Pred: broccoli OOV: Romanesco</p>	<p>GT: pot Pred: ear OOV: Corn plants</p>	<p>GT: desktop computer, desk Pred: Angora OOV: Other rabbit/hare</p>
			
<p>GT: zucchini Pred: head cabbage OOV: Other vegetable</p>	<p>GT: strawberry Pred: strainer OOV: Other pan</p>	<p>GT: cowboy hat Pred: banded gecko OOV: Other reptile</p>	<p>GT: hay Pred: barn OOV: Other building</p>
			
<p>GT: parking meter Pred: cab OOV: Other vehicle</p>	<p>GT: teddy Pred: toyshop OOV: Other shop</p>	<p>GT: plate, soup bowl Pred: trifle OOV: Other food</p>	<p>GT: Siberian husky, Eskimo dog Pred: Norwegian elkhound OOV: Other breed</p>

Figure 10: **Fine-grained with OOV mistakes.** Additional examples of fine-grained with out-of-vocabulary mistakes. Of the correct multi-labels, the original ImageNet label is listed first.



When does dough become a bagel? Analyzing the remaining mistakes on ImageNet

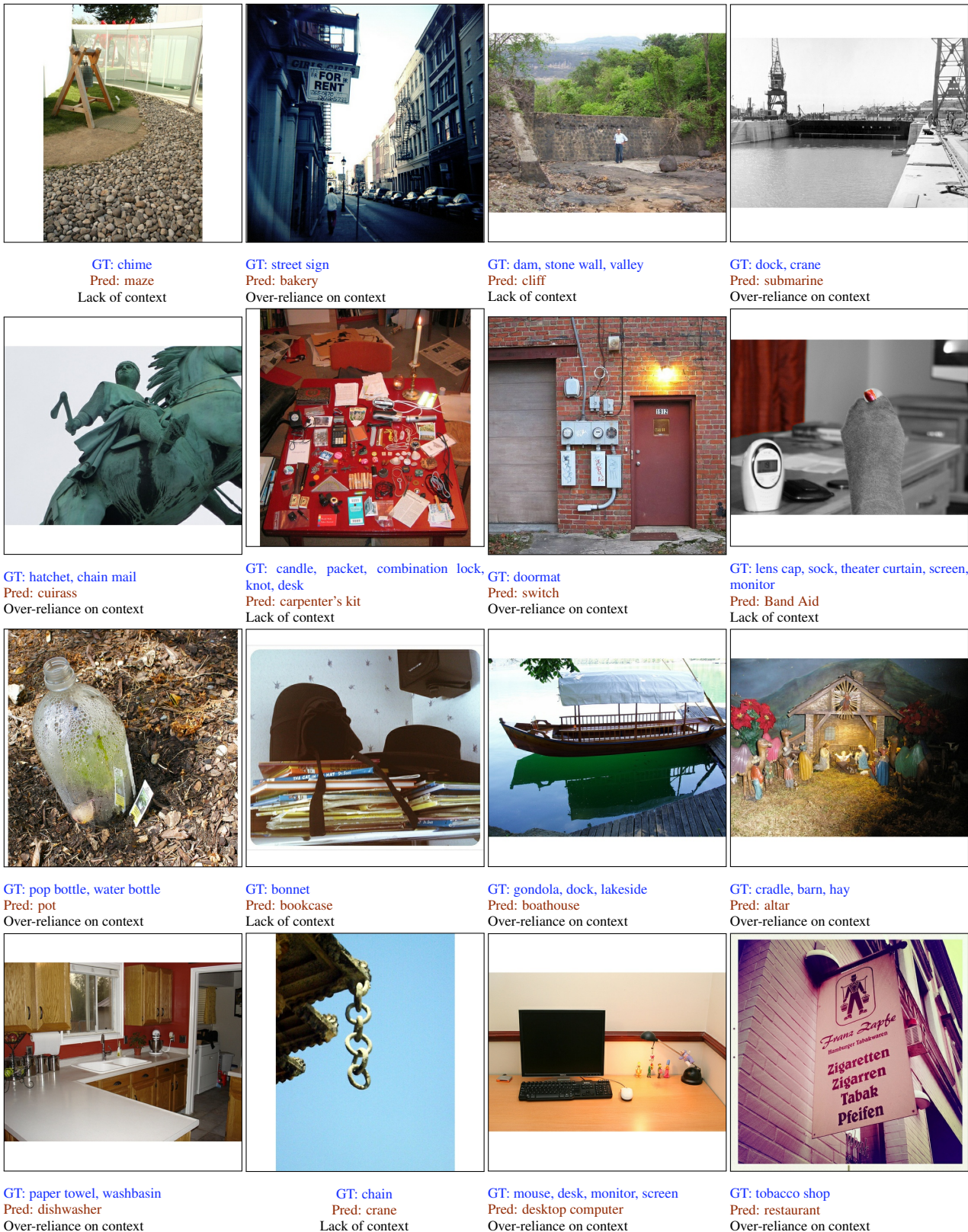


Figure 11: **Spurious correlation examples.** Of the correct multi-labels, the original ImageNet label is listed first. Over-reliance on context indicates that surrounding cues in the image correlate with the predicted class, although the predicted class is not present. Lack of context indicates that the model has failed to understand relevant context in the image, and predicts a class that is inconsistent with a holistic understanding of the image.



When does dough become a bagel? Analyzing the remaining mistakes on ImageNet

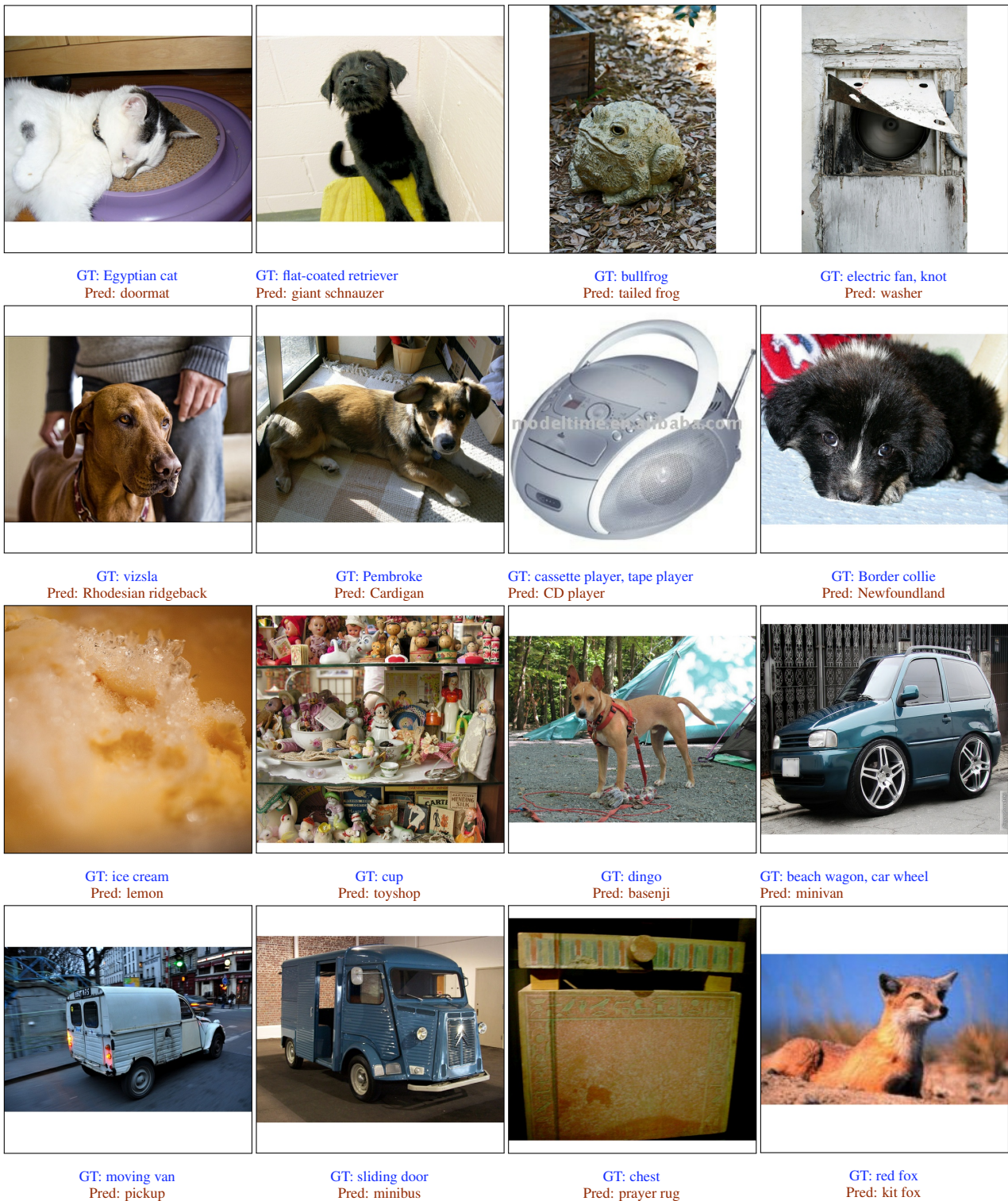
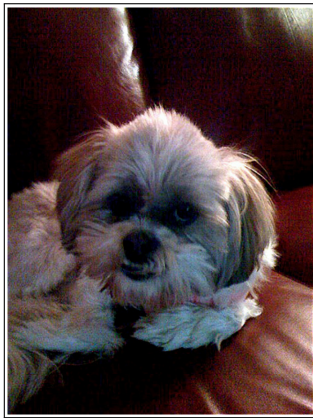


Figure 12: **Non-prototypical mistakes.** Additional examples of non-prototypical mistakes. Of the correct multi-labels, the original ImageNet label is listed first. Non-prototypical examples are typically unusual border cases of the groundtruth class, such as puppies of a dog breed, or unusual/unique versions of the class.



(a)  
GT: Lhasa



(b)  
GT: West Highland white terrier



(c)  
GT: Norwich terrier



(d)  
GT: bluetick



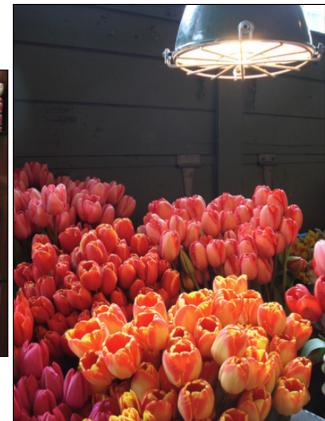
(e)  
GT: jack-o'-lantern  
Pred: ballpoint



(f)  
GT: torch  
Pred: lighter



(g)  
GT: cup  
Pred: toyshop



(h)  
GT: spotlight  
Pred: grocery store

Figure 13: **Difficult images for humans and models.** Top row: the only 4 images that all humans (Shankar et al., 2020) classify incorrectly (our model classifies these correctly). Bottom row: images that should-be-easy (all humans get correct), but the model gets incorrect.



## J. Analyzing the Training Data

In this section, we try to understand the model’s remaining mistakes by investigating the training data. We do this through the lens of looking at the nearest training examples to the ViT-3B models remaining mistakes. To do this, we generate embeddings for every training and validation example using the pretrained ViT-3B models JFT checkpoint (before fine-tuning on Imagenet), and for every error, we query the K=10 nearest neighbors using an exact nearest neighbor lookup.

### J.1. Validation Set Leakage.

One of the most interesting findings using nearest neighbors was rediscovering that 797 ImageNet validation images (1.594%) **exist in the training set as exact duplicates** (pixelwise L2 distance of 0), 34 of them more than once for a total of 831 duplicate training images. While this was previously documented in (Sun et al., 2017) and (Kolesnikov et al., 2020), we are the first to notice that every single leaked sample has a different label in the train set than in the validation set, indicating the ImageNet authors did deduplicate within a class, but not across classes. Analyzing these duplicates we find most of them represent challenging fine grained image classes (e.g. two similar dog breeds), or images where multiple annotations are appropriate. Additionally, in the Appendix J.2 we detail a second, harder to detect leak pattern we saw commonly in the training data with "near duplicates", images in the training set that are from the same photo-shoot or scene as a validation image, or are cropped or processed versions of a validation image. This phenomenon, with similar discovery methodology, has also been observed on CIFAR (Barz & Denzler, 2020).

To understand the impact of this validation set leakage, we remove all the exact duplicates from the training set and retrain both a ResNet50 from scratch, and our JFT pretrained ViT-3B. Results are shown in Table 4. Unintuitively, when we remove all the leaked validation images from the training set and retrain, we see both Top1 Accuracy and Multi-label Accuracy (MLA) actually stay the same or decrease overall, despite all leaked training images having different labels than their leaked validation image counterpart. MLA accuracy both before and after deduplication are high, which leads us to believe that the training labels may be correct under multi-label evaluation. To verify this, we find 320 of the leaked validation images were in our 20k multi-label validation set, corresponding to 331 training images. Using these multi-labels, we find that 219/331 (66.2%) of these training images labels would have been correct under multi-label evaluation. Nevertheless, we do see an increase in both Top1 and MLA accuracies on the leaked subset of the data. Finally, it seems the ViT-3B model is less sensitive to leaked validation images, which may be a function of the

fine tuning recipe we used (which was exactly the same as the fine-tuning recipe on the original non-deduped data).

### J.2. Near Duplicates

In addition to the 831 training examples that are exact duplicates of validation examples (all with different labels), there is a large collection of "near duplicates" in the Imagenet train set. "Near duplicates" are images that are either crops, augmentations, or resizes of validation images, or more unexpectedly, images from the same scene or photo shoot. These images are often visually different but semantically the same, and as a result are much harder to detect with traditional embedding distance thresholding based deduplication. Nevertheless, these can also leak test set information, introduce label noise (if the labels are different than the validation set), or if many training examples are from the same scene, reduce the effective dataset size of that class.

While we cannot provide an estimate of the prevalence of this problem, we find it often while analyzing the K-Nearest Neighbors (in JFT embedding space) of validation images. We show some of these examples in Figures 14 and 15. While we just show dog and object examples here, we find this also happens with reasonable frequency for human related classes (especially activity related classes like soccer, basketball, parallel bars, etc). The existence of these duplicates raises an interesting question: How close can training examples be to your validation set before it becomes problematic?

### J.3. Neighbors of Spurious Correlations

While fine grained errors and OOV mistakes are often somewhat intuitive to a human, spurious correlations are harder to understand. To try to understand them in more detail, we look at the neighbors of several of the ViT-3B models major spurious correlations. In Figure 16 we show two such examples, but find this to be relatively common.

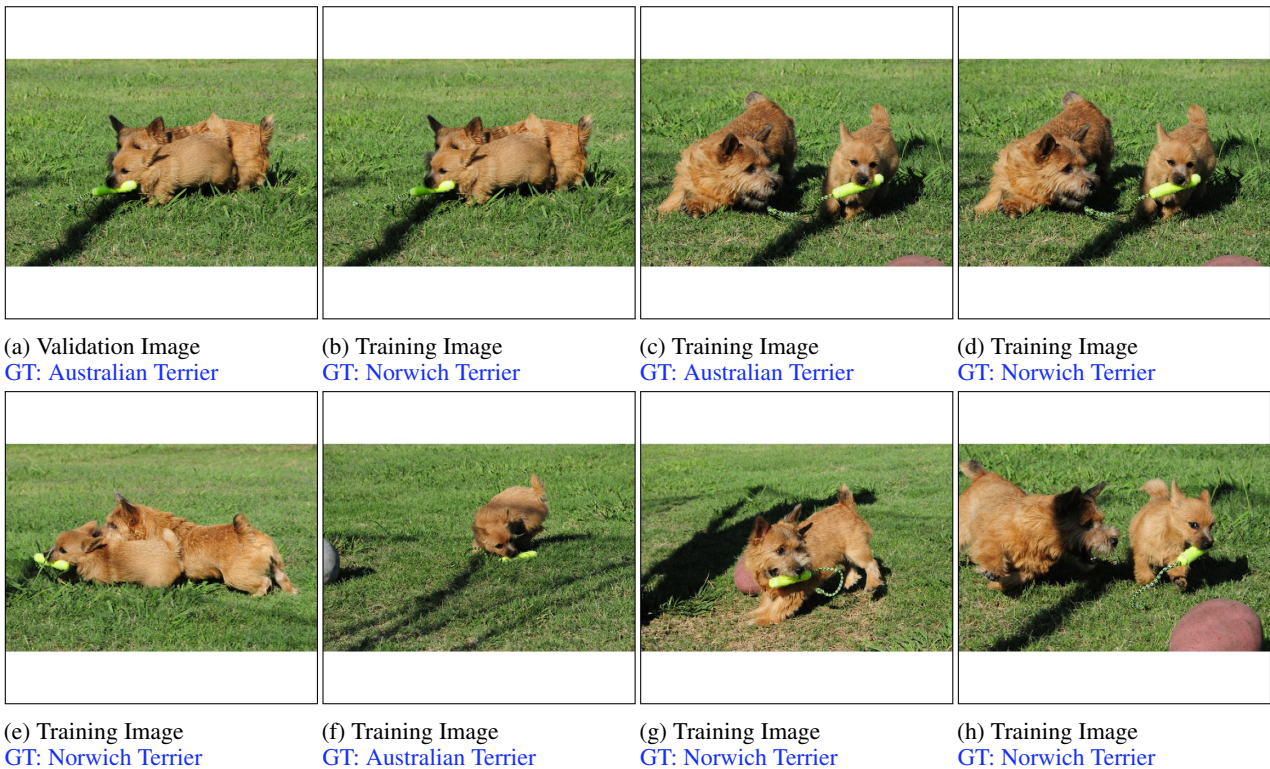


Figure 14: **Near Duplicates:** (a) We show a validation image of two dogs playing, labeled originally as an Australian Terrier. When looking at the  $K = 10$  nearest neighbors, we find all 10 of them to be of the same two dogs with one of two labels, shown as images (b) through (h), including some training examples that are duplicates of each other. Because we only retrieve the 10 nearest neighbors, there could be even more than 10 images of this scene.

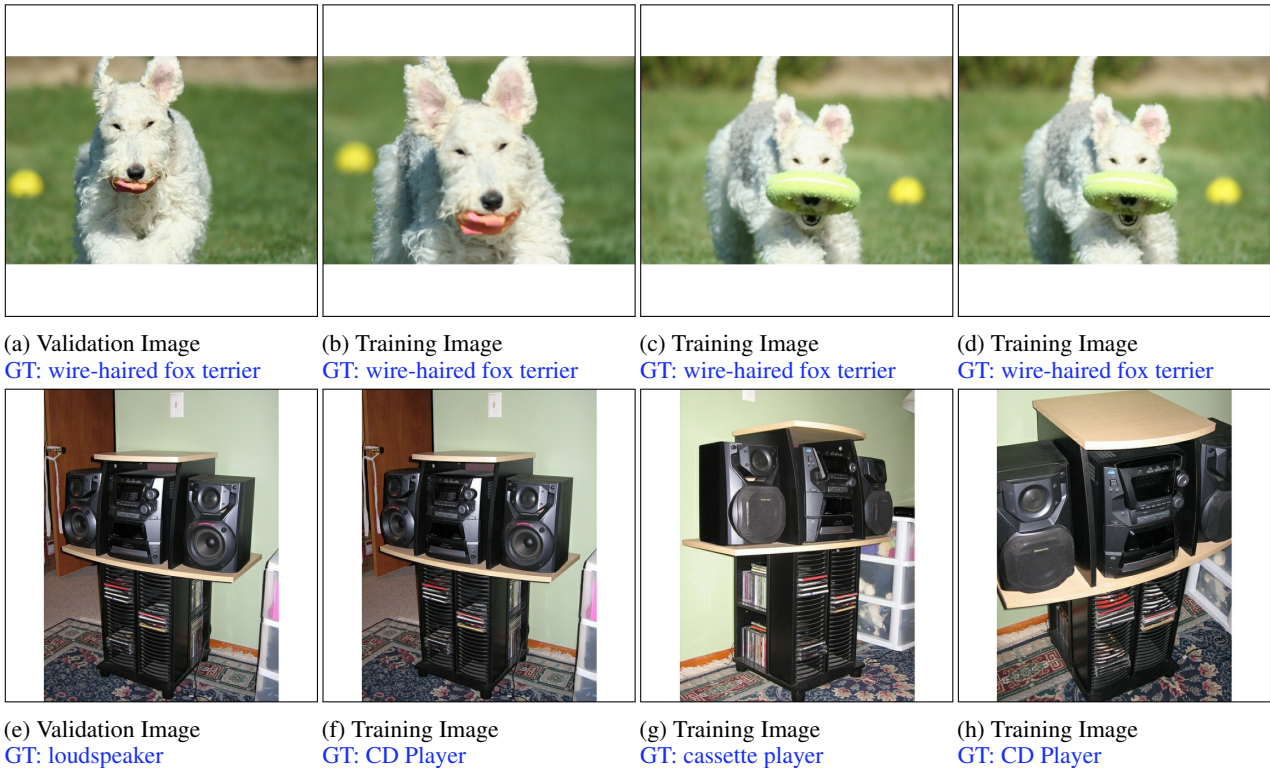


Figure 15: **Near Duplicates:** Top: (a) We show a wire-haired fox terrier with a cropped version of the validation image as a training image, and two more training images that are cropped versions of each other. We find 6/10 of the nearest neighbors are of the same dog. Bottom: We show a "loudspeaker" in the validation set with nearby images from the training set of the same speaker setup with labels of "CD Player" and "cassette player". We find 8/10 of its nearest neighbors are the same speaker setup, with several of the training images being exact duplicates of each other (or the validation set) with different labels.



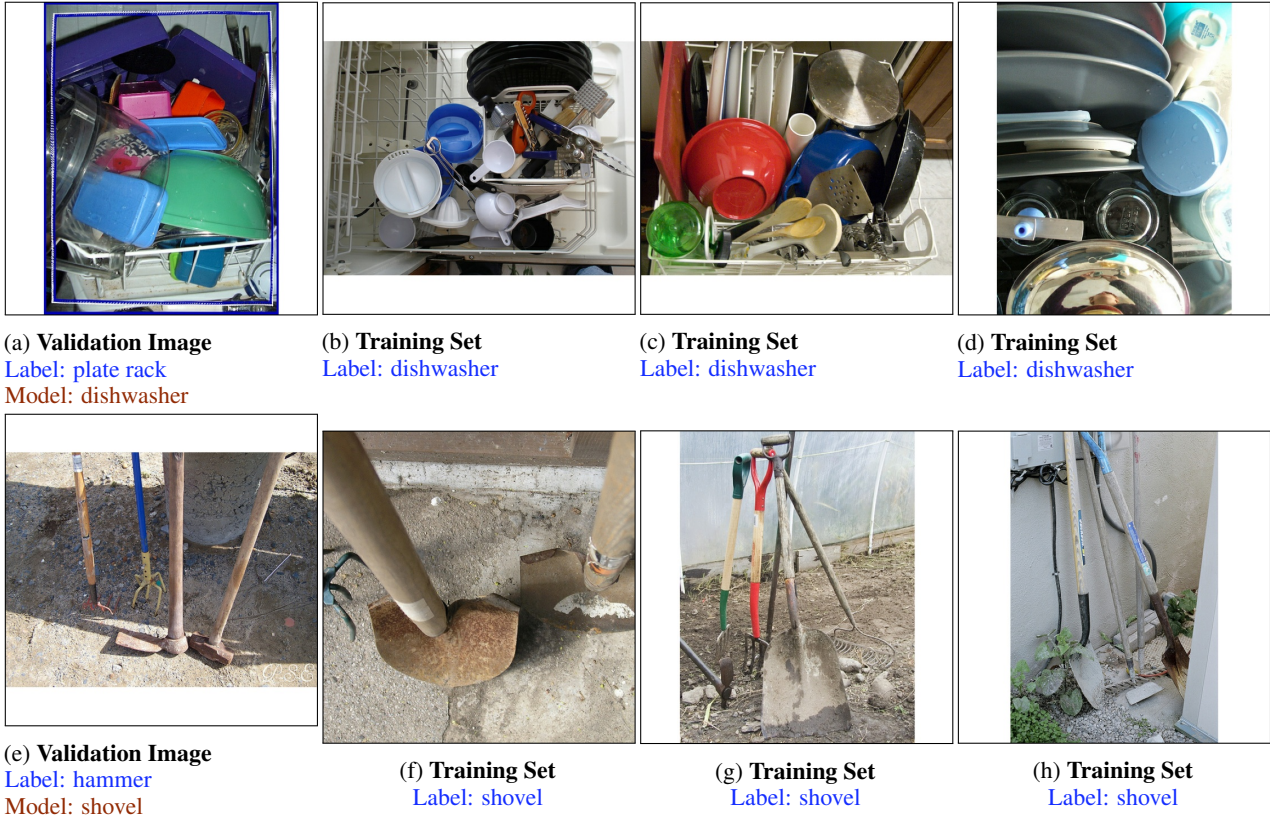


Figure 16: **Neighbors of a Major Spurious Correlation.** We show an example of two major spurious correlations and a subset of their  $K=10$  nearest neighbors (in JFT embedding space). For the dishwasher example (top), 8/10 of the nearest neighbors were pictures of cluttered dishes where the dishwasher machine was not in view. For the shovel example, we find all 10 nearest neighbors are shovels standing upright and outdoors. There are no other validation images of the hammer class with it standing upright or even outdoors.

Model	Deduped?	Top1	Top1 on Leaked	MLA	MLA on Leaked
ResNet50	-	76.0%	26.9%	84.8%	80.1%
ResNet50	✓	76.0%	40.2%	84.6%	82.1%
ViT-3B	-	89.5%	42.5%	97.8%	93.4%
ViT-3B	✓	89.4%	45.1%	97.4%	94.0%

Table 4: **Change in performance when removing leaked training examples.** We show both our ViT-3B and a ResNet50 for comparison, and report both Top1 accuracy and Multi-label Accuracy (MLA) on both the whole validation set and the leaked subset.