












Leveraging Human-Machine Interactions for Computer Vision Dataset Quality Enhancement

Esla Timothy Anzaku^{1,2,3} , Hyesoo Hong¹ , Jin-Woo Park¹ ,
Wonjun Yang¹ , Kangmin Kim¹ , JongBum Won¹ ,
Deshika Vinoshani Kumari Herath⁵ , Arnout Van Messem⁴ , and Wesley
De Neve^{1,2,3} 

¹ Ghent University Global Campus, Incheon 21985, South Korea
eslatimothy.anzaku@ugent.be

² Center for Biosystems and Biotech Data Science, Ghent University Global Campus,
Incheon 21985, South Korea

³ IDLab, Ghent University, Technologiepark-Zwijnaarde 126, 9052 Ghent, Belgium

⁴ University of Liège, 4000 Liège, Belgium

⁵ Mediiio, Seoul, South Korea

Abstract. Large-scale datasets for single-label multi-class classification, such as *ImageNet-1k*, have been instrumental in advancing deep learning and computer vision. However, a critical and often understudied aspect is the comprehensive quality assessment of these datasets, especially regarding potential multi-label annotation errors. In this paper, we introduce a lightweight, user-friendly, and scalable framework that synergizes human and machine intelligence for efficient dataset validation and quality enhancement. We term this novel framework *Multilabelfy*. Central to *Multilabelfy* is an adaptable web-based platform that systematically guides annotators through the re-evaluation process, effectively leveraging human-machine interactions to enhance dataset quality. By using *Multilabelfy* on the *ImageNetV2* dataset, we found that approximately 47.88% of the images contained at least two labels, underscoring the need for more rigorous assessments of such influential datasets. Furthermore, our analysis showed a negative correlation between the number of potential labels per image and model top-1 accuracy, illuminating a crucial factor in model evaluation and selection. Our open-source framework, *Multilabelfy*, offers a convenient, lightweight solution for dataset enhancement, emphasizing multi-label proportions. This study tackles major challenges in dataset integrity and provides key insights into model performance evaluation. Moreover, it underscores the advantages of integrating human expertise with machine capabilities to produce more robust models and trustworthy data development.

Keywords: Computer Vision · Dataset Quality Enhancement · Dataset Validation · Human-Computer Interaction · Multi-label Annotation

1 Introduction

Deep learning, the engine behind advanced computer vision, has been largely propelled by training on expansive resources like the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) dataset [1], commonly known as ImageNet-1k. However, recent performance trends in deep neural network (DNN) models trained on these datasets have shown top-1 and top-5 accuracy stagnation across diverse DNN architectures and training techniques, irrespective of model complexity and dataset size [2,3]. This performance plateau suggests that we may be nearing the limits of model accuracy with the current ImageNet-1k dataset using the top-1 accuracy.

A potentially overlooked factor contributing to the observed stagnation may be attributed to the inherent multi-label nature of the dataset in question. It is plausible that a substantial proportion of the images in the dataset are related to more than a single ground truth label. However, the dataset only provides labels for a singular ground truth, which may impose limitations [4–6]. This single-label ground truth constraint could inadvertently lead to underestimating the performance of DNN models, particularly when utilizing the top-1 accuracy metric.

Furthermore, the performance of models significantly degrades when assessed on newer but similar datasets, such as ImageNetV2 [7]. Despite being developed following a similar protocol to the original ImageNet-1k dataset, ImageNetV2 exhibits unexplained accuracy degradation across various models, regardless of model architecture, training dataset size, or other training configurations. While efforts are being made to investigate this degradation [8,9], we found only one work that partially studied this problem [10].

Prior work has acknowledged the need for more accurate dataset labels and has published reassessed labels that reflect the multi-label nature of the ImageNet-1k validation set [4]. However, label reassessment is not a trivial task. It requires considerable resources and expertise, presenting a substantial challenge for smaller research groups. Given the vital role of the validation and test sets in DNN model selection and benchmarking, meticulous analysis of the ImageNet-1k validation set and its replicates remains indispensable. This critical importance highlights the necessity for accessible and effective frameworks to scrutinize and tackle the multi-label nature of computer vision single-label classification datasets. To address this, we propose an accessible and scalable framework, termed *Multilabelfy*, that combines human and machine intelligence to efficiently validate and improve the quality of computer vision multi-class classification datasets. *Multilabelfy* comprises four stages: (i) label proposal generation, (ii) human multi-label annotation, (iii) annotation disagreement analysis, and (iv) human annotation refinement. It is designed with two primary objectives: to strategically harness the capabilities of a diverse pool of annotators and to seamlessly blend human expertise with machine intelligence to improve the quality of a dataset. These objectives are made accessible through a user-friendly interface.

This research effort enriches existing literature by offering Multilabelify for improving the quality of computer vision multi-class classification datasets. Utilizing Multilabelify, we reassessed the labels for ImageNetV2, revealing that 47.88% of the images in this dataset could have more than one valid label. We also identified other noteworthy dataset issues. Our work accentuates the importance of recognizing and addressing the multi-label nature of ImageNet-1k and its replicates. Our ultimate goal is to contribute towards developing robust DNN models that can effectively generalize beyond their training data.

2 Related Work

2.1 Label Errors

Label errors have been identified within the test sets of numerous commonly used datasets, including a 6% error rate in the ImageNet-1k validation set [11]. The importance of tackling the issue of label errors in test partitions of datasets was further emphasized. It was found that high-capacity models are prone to mirroring these systematic errors in their predictions, potentially leading to a misrepresentation of real-world performance and distortion in model comparisons. In another work, an extensive examination of 13,450 images across 269 categories in the ImageNet-1k validation set, which predominantly includes wild animal species, was conducted [12]. Through collaboration with ecologists, it was found that many classes were ambiguous or overlapping. An error rate of 12% in image labeling was reported, with some classes being erroneously labeled more than 90% of the time. Our work further accentuates the critical role of addressing label errors in datasets used for model evaluation. It underscores the need for more precise and thorough dataset construction and assessment methodologies.

2.2 From Single-Label to Multi-Label

Single-label evaluation has traditionally served as the standard for assessing models on the ImageNet-1k dataset. However, a reassessment of the ImageNet-1k validation ground truth labels revealed that a good proportion of the images could have multiple valid labels, prompting the creation of Reassessed Labels (ReaL) [4], incorporating these multi-labels.

In a related study [6], the remaining errors that models made on the ImageNet-1k dataset were examined, focusing on the multi-label subset of *ReaL*. Nearly half of the perceived errors were identified as alternative valid labels, confirming the multi-label nature of the dataset. However, it was also observed that even the most advanced models still exhibited about 40% of errors readily identifiable by human reviewers.

2.3 ImageNet-1k Replicates

When tested on replication datasets like ImageNetV2, DNN models have been observed to demonstrate a significant, yet unexplained, drop in accuracy [7].

Despite these replication datasets, including ImageNetV2, being created by following the original datasets' creation protocols closely, the performance decline raises significant questions about the models' generalization capabilities or the integrity of the datasets. The significant performance drop on ImageNetV2, between 11% to 14% [7], was based on the conventional approach of evaluating model accuracy using all data points in the test datasets.

However, it has been argued that the conventional evaluation approach may not fully capture the behavior of DNN models and may set unrealistic expectations about their accuracy [8,9]. A more statistically detailed exploration into this unexpected performance degradation on ImageNetV2 found that standard dataset replication approaches can introduce statistical bias [8]. After correcting for this bias and remeasuring selection frequencies, the unexplained part of the accuracy drop was reduced to an estimated $3.6\% \pm 1.5\%$, significantly less than the original $11.7\% \pm 1.0\%$ earlier reported in [7]. An alternative evaluation protocol that leverages subsets of data points based on different criteria, including uncertainty-related information, provides an alternative perspective [9]. Through comprehensive evaluation leveraging the predictive uncertainty of models, the authors found that the degradation in accuracy on ImageNetV2 was not as steep as initially reported, suggesting possible differences in the characteristics of the datasets that warrant further investigation. A closely related research work studies various aspects of the ImageNet-1K and ImageNetV2 datasets using human annotators. Using a sample of 1,000 images from both datasets, the proportion of images with multiple labels was estimated to be 30.0% and 34.4% for the two datasets, respectively. This information is detailed in Section B.2 of the supplementary material in [10]. The cited work suggested that the difference in the multi-label composition between the two datasets could be a possible explanation for the accuracy degradation.

2.4 Key Modifications to Existing Approaches

Our research expands upon a previous work [4] with several essential modifications:

Model Selection for Candidate Label Proposals. In contrast to the original study's use of a hand-annotated sample of 256 images from a 50,000-image dataset to guide the selection of an optimal model ensemble, we built upon their work, utilizing their generated multi-labels and proposed *ReaL accuracy*. We selected the best-performing pre-trained model utilizing the ReaL accuracy metric, designed to evaluate multi-class classification DNN models on a multi-label test dataset. Further details on this process are provided in Sect. 3.2.

Image Pre-selection for Multi-label Annotation. The original study only utilized an ensemble of pre-trained single-label models to generate eight candidate labels. In contrast, our approach extended the candidate proposals to 20, thereby decreasing the risk of omitting valid labels and increasing selection accuracy (Sect. 3.2).

Annotation Refinement. We introduced an additional stage, wherein the top twenty model-proposed labels, alongside all human-selected labels, are presented to an additional pool of experienced annotators for further refinement (Sect. 3.5).

Open-Source Platform. Recognizing that platforms like Mechanical Turk might be inaccessible or not affordable for some research labs, we developed Multilabelfy, an open-source alternative. This platform allows in-house dataset quality improvement while maintaining a user-friendly interface.

Section 3 provides more comprehensive information regarding these contributions.

3 Proposed Framework

3.1 Overview

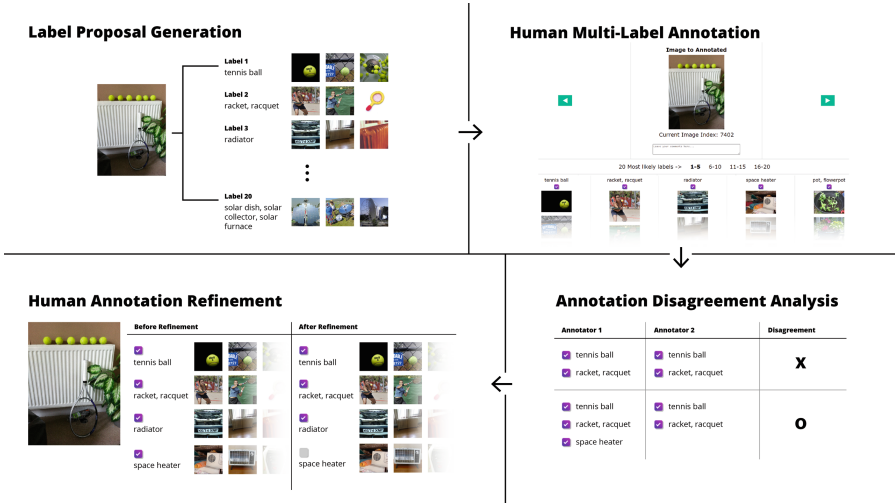


Fig. 1. Overview of the proposed framework for enhancing computer vision datasets from single-label to multi-label, enabling a more comprehensive capture of their descriptions.

The proposed multi-label dataset enhancement framework (Fig. 1) comprises four key stages: (i) label proposal generation, (ii) human multi-label annotation, (iii) annotation disagreement analysis, and (iv) human annotation refinement. The label proposal generation and annotation disagreement analysis can be automated using the appropriate algorithms while the human multi-label annotation and human annotation refinement require the involvement of human annotators.

3.2 Label Proposal Generation

Our qualitative analysis shows that pre-trained models, originally trained on single-label computer vision datasets, can effectively rank the predicted probability vector. This capability is corroborated by the near-perfect top-5 accuracies of state-of-the-art DNN classification models reaching approximately 99% [2]. For model selection, we utilize the *ReaL accuracy* metric, specifically designed to assess the performance of single-label pre-trained DNN models in multi-label scenarios. *Under this metric, an image prediction is considered correct if the prediction belongs to the set of ground truth labels assigned to the image.* The selected model is then used to generate the top-20 label proposals, an increase from the eight proposals presented in previous work, to ensure broader coverage of valid labels. Given the potential for information overload with many label proposals, we designed the annotation user interface to mitigate this concern. Additional details regarding the role of human annotators and the annotation interface are discussed in Sect. 3.3.

3.3 Multi-label Annotation by Human Annotators

Multilabelify incorporates a strategically designed web interface to alleviate the workload of human annotators, with a screenshot provided in Fig. 2. This user interface is characterized by several key features engineered to enhance the efficiency and effectiveness of the annotation process. It facilitates the display of label names, their corresponding synonyms, and representative images from the pool of twenty potential labels systematically organized into four subgroups of five labels each. In the event that the initial group does not sufficiently encompass all visible objects, annotators have the option to navigate to other label groups.

The design also incorporates a streamlined selection process facilitated by a singular checkbox assigned to each proposed label. Moreover, ten exemplar images are presented in a scrollable format for each proposed label, providing a comprehensive view without overwhelming the annotator. Further attention to detail is reflected in the feature that allows images to be clicked on, enabling annotators to inspect these images at their original resolution. These elements combined optimize the multi-label annotation process, yielding higher accuracy and efficiency.

3.4 Annotation Disagreement Analysis

Single-label multi-class classification computer vision datasets often comprise images featuring multiple objects. However, a prior research work [5] estimated that about 80% of the ImageNet-1k images contain a single object. We also expect some images with multiple labels to pose no challenges to the annotators. Considering the aforementioned observations, our framework seeks to effectively exclude such images from the pool intended for further refinement. We target

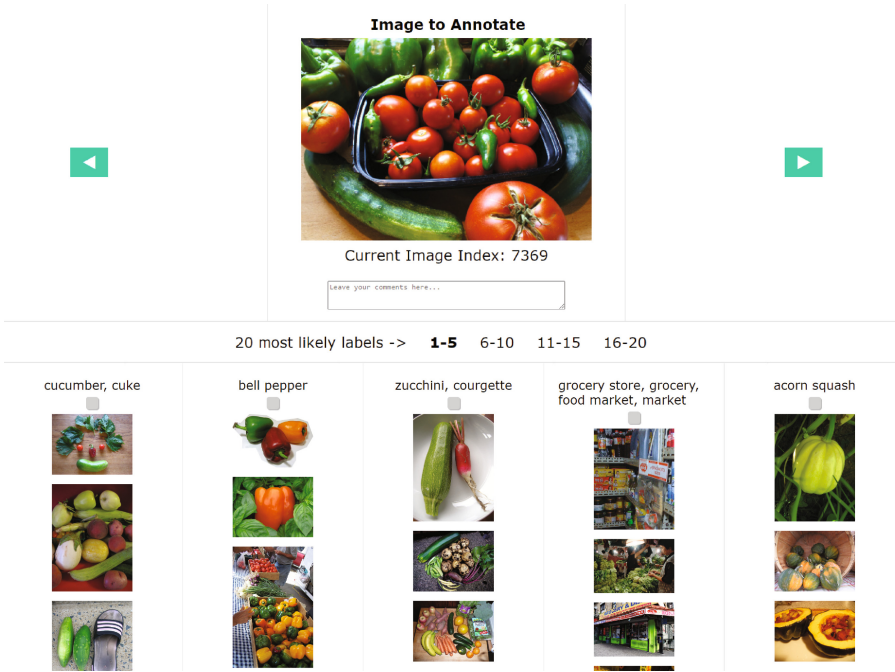


Fig. 2. The user interface of the annotation platform. It showcases key features like label presentation in groups of five, a single checkbox per proposed label, scrollable sample images, and click-to-enlarge functionality for detailed inspection of images. These features are designed to streamline the annotation process and efficiently accommodate multi-label data annotation.

images that require additional human annotation refinement during the *annotation disagreement analysis* stage, as depicted in Fig. 1. Images are selected for further annotation refinement if the labels generated by human annotators, as discussed in Sect. 3.3, fail to meet a predefined annotation agreement condition. This annotation agreement condition requires: *complete consistency across all labels identified by human annotators for a particular image and the inclusion of the originally provided ground truth label within the array of labels selected by the annotators*. This strategic condition facilitates focused refinement of annotations for the subset of images that pose more significant challenges to annotators. As a result, we minimize the misuse of annotators' time and provide an avenue for a more detailed examination of the more complex images, ultimately fostering a more thoroughly annotated dataset.

3.5 Refinement of Human Annotation

This stage follows the process described in Sect. 3.3 but with some critical distinctions. In the stage described in Sect. 3.3, annotators with varying degrees of

experience with the dataset contribute to the labeling. However, the refinement stage exclusively engages more experienced annotators. These experienced annotators are provided with the labels previously selected, which are pre-checked for the annotators to review: uncheck (to correct) or check additional missing labels. Furthermore, the annotators are instructed to document any changes they make to the labels using the comments section of the web interface. This provision ensures that a clear record is maintained for each correction, which can be invaluable in resolving potential discrepancies in the annotations. It is important to note that these annotators have undergone several tutorial sessions on the label issues of the ImageNet-1k dataset. Additionally, they reviewed and summarized related literature to ensure that they are aware of the nuanced issues that are encountered when annotating images into 1,000 categories, especially within the fine-grained categories.

4 Results

4.1 Experimental Setup

Our goal is to re-assess the labels for the ImageNetV2 dataset to accommodate and account for its multi-label nature. The four stages in Sect. 3 were carefully followed. In the label proposal generation stage, the EVA-02 [13] model was used to generate the proposal. It is one of the top performing models (90.05% top-1 accuracy [2]) on the ImageNet-1k dataset; additional details of the model can be found in the cited paper. Subsequently, in the human multi-label annotation stage, the 10,000 images of the ImageNetV2 dataset were partitioned into seven batches, and each batch was assigned to two human annotators. Fourteen human annotators having varying experience levels with the ImageNet-1k dataset and computer vision in general, participated during this stage. Upon the annotation disagreement analysis (detailed in Sect. 3.4), the annotations for 6,425 of the 10,000 images fulfilled our disagreement criteria and were selected for subsequent refinement by five more experienced annotators, four of whom were previously referenced among the group of fourteen. Each annotator refined the annotations for 1,285 images. The refined annotations were then used to generate the results presented and discussed in the following sections.

4.2 The Extent of ImageNetV2 Multi-Labeledness

Here, we provide visual statistics summarizing the multi-label nature of the labels we generated for the ImageNetV2 dataset. Specifically, we show what percentage of the dataset contains which label count, i.e., the number of ground truth labels assigned to an image. As shown in the pie chart of Fig. 3, the annotation process could not find labels for 1.29% of the images. Moreover, 50.83% of the images contain one label, 23.85% contain two labels, and 24.03% contain more than two labels.

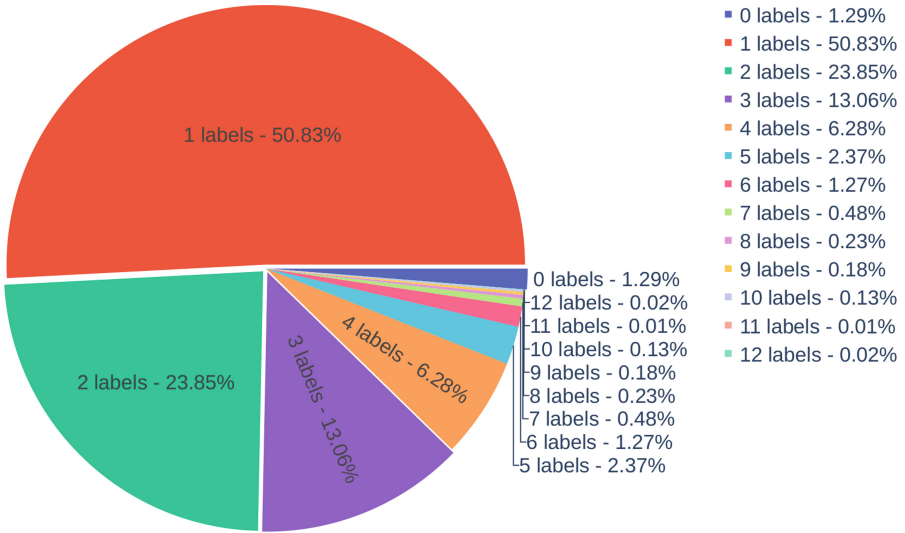


Fig. 3. The distribution of images based on the number of labels assigned to them during our annotation process.

4.3 Re-evaluation of Models on ImageNetV2 Improved Labels

Top-1 Accuracy Versus ReaL Accuracy. We provide a Scatterplot to understand the relationship between ReaL and top-1 accuracy on our generated labels (Fig. 4). Each dot in the plot represents a pre-trained model, and 57 models were evaluated on the ImageNet-1k validation set and ImageNetV2. These models are sourced from a publicly available GitHub repository [2] and represent state-of-the-art models pre-trained either exclusively on the ImageNet-1k dataset, or on additional external data. Details of these models can be found together with the paper’s code at <https://github.com/esla/Multilabelify>. The regression analysis indicates a significant correlation between the two metrics. Specifically, for every percentage point increase in top-1 accuracy, the ReaL accuracy rises by approximately 0.5788 percentage points. The coefficient of determination, R^2 , is 75.69%, suggesting that 75.69% of the variation in ReaL accuracy is explained by its linear relationship with top-1 accuracy. This result reflects a consistent positive relationship: as the top-1 accuracy of models improves, there is a proportional increase in ReaL accuracy. It is worth noting that four models visibly diverge from the regression line; these models merit additional scrutiny to identify potential model-specific quirks or underlying reasons for their divergence. A detailed investigation of these models will be addressed in future work.

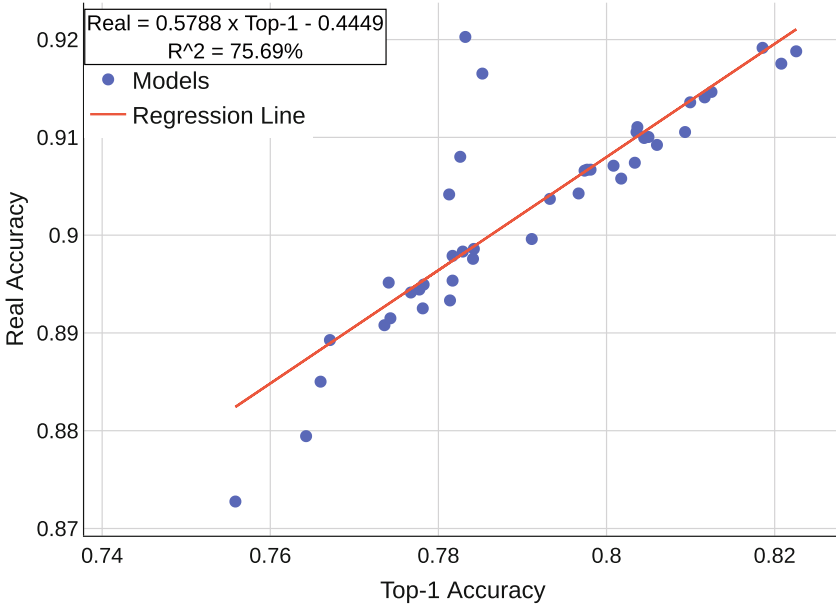


Fig. 4. Scatterplot of Real accuracy versus top-1 accuracy for 57 top-performing DNN models, pre-trained either exclusively on the ImageNet-1k dataset or additionally on external datasets.

Visual Statistics of Top-1 Accuracy Versus Image Count. We investigate the relationship between top-1 accuracy and the variability in image label assignments using heatmaps (Fig. 5). While we presented the results for 57 models in Sect. 4.2, for visual brevity, we randomly selected 5 models for the heatmaps. We determine top-1 accuracy using ground truth labels from the ImageNetV2 dataset, comparing them with our multi-label annotations. To this end, we employ a heatmap (Fig. 5, top) that presents top-1 accuracies for each evaluated model across different *label count* categories. While this heatmap is informative, it does not factor in the variability stemming from different sample sizes across label counts. For instance, images with a single label may be more prevalent than those with multiple labels, potentially leading to biases in accuracy measurements.

To enhance our understanding of accuracy computations and account for inherent uncertainties, we incorporate a secondary heatmap as shown in Fig. 5, bottom. The margin of error related to the top-1 accuracy is denoted as $U(i, j)$ and is determined using the following formula: $U(i, j) = 1.96 \times \sigma(i, j) / \sqrt{n}$. Here, $\sigma(i, j)$ stands for the standard deviation stemming from the binary outcomes of individual predictions for a specific *model* and *label count*. This standard deviation for a binary variable is expressed as $\sigma(i, j) = \sqrt{p(1-p)/n}$, where p represents the proportion of correct predictions. The variable n symbolizes the number of observations for the considered model-label count pairing. This margin

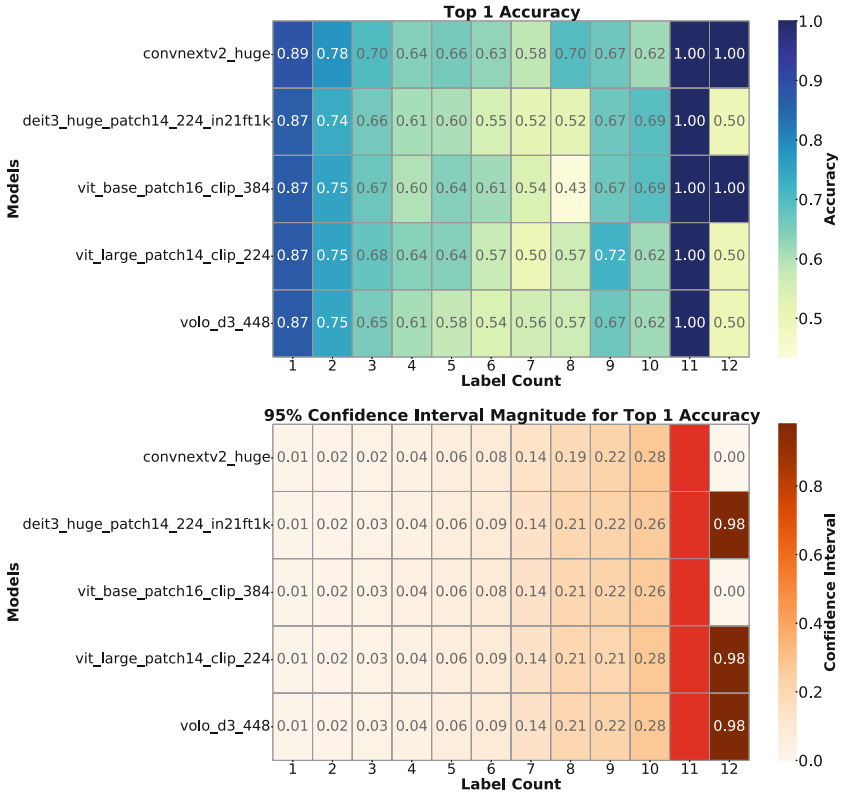


Fig. 5. Heatmaps displaying top-1 accuracy (top) for five randomly selected models evaluated on our multi-labeled ImageNetV2 dataset, and the half-width of the 95% confidence interval (bottom) associated with these accuracies. Red cells without numbers represent NaN values due to sets with one or no images for a given label count. (Color figure online)

of error, corresponding to half the width of the 95% confidence interval, offers a gauge of uncertainty for each model-label count combination. Differences in sample sizes across subsets can lead to variations in the width of the confidence interval. This variance emphasizes the significance of jointly considering both accuracy and its associated uncertainty when interpreting model performance across different label counts.

In our analysis, while results for only five models are presented for clarity, the observations are representative of numerous other models evaluated. A notable observation is that models consistently exhibit higher top-1 accuracy for images associated with a single label. However, as the number of potential labels expands, a discernible decrease in accuracy is evident. This pattern potentially indicates that models might be predicting alternative valid labels, and the top-1 accuracy metric penalizes them for such predictions. Such a negative correlation warrants

attention, as it hints at the possibility of underestimating model performance due to potentially skewed dataset assumptions.

4.4 Analysis of Images with Zero Labels

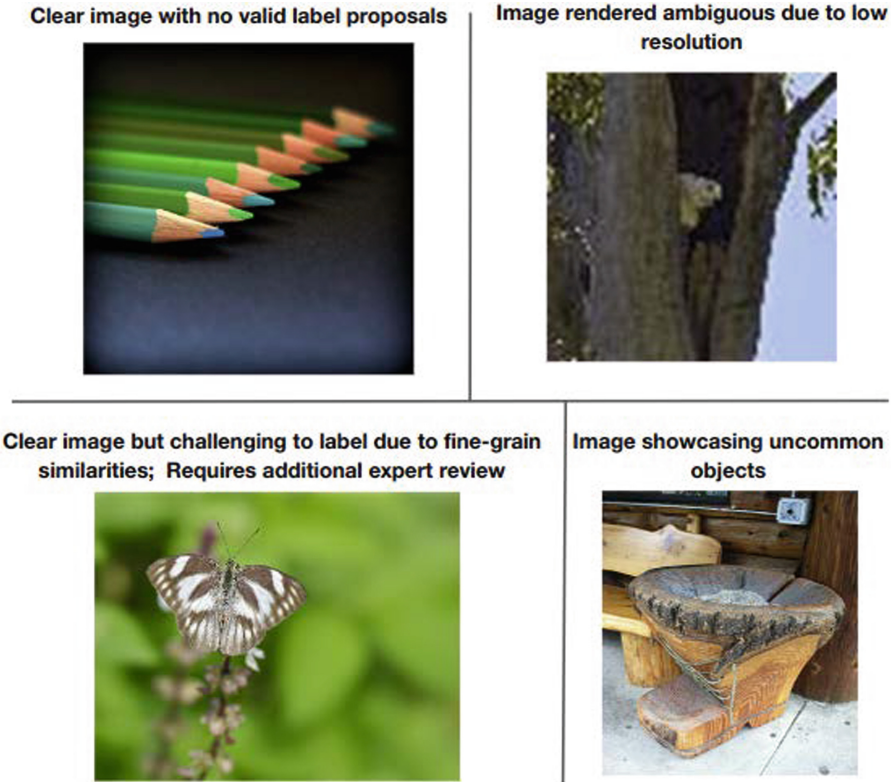


Fig. 6. Example images where annotators did not find matching labels from the 20 proposed labels. The images are categorized based on possible explanations for not finding matching labels in the labels proposed (see Sect. 4.4).

During our dataset annotation process, despite the meticulous efforts of fifteen annotators, 1.29% (129 images) had no labels assigned to them at the completion of human annotation refinement. Consequently, two of the experienced annotators further scrutinized these images. They classified the images without valid annotations into (i) clear images with no valid label proposals (21.79%), (ii) images rendered ambiguous due to low resolution (10.26%), (iii) clear images but challenging to label due to fine-grain similarities, thereby requiring additional expert review (38.46%), and (iv) images showcasing uncommon objects

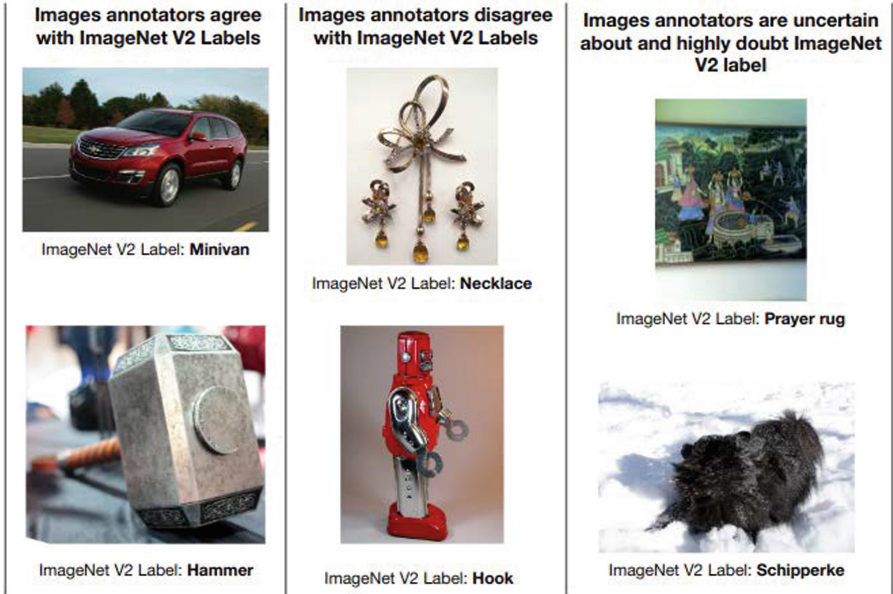


Fig. 7. Example images where annotators did not find matching labels from the 20 proposed labels. The images are categorized based on whether or not our two annotators agree with the provided ImageNetV2 ground truth label (see Sect. 4.4).

or atypical viewpoints (29.49%). One example from each of these categories is shown in Fig. 6.

While our finalized annotations did not provide labels for these images, ground truth labels from the creators of the ImageNetV2 dataset existed for reference. Using these, the annotators further categorized the images based on their alignment with the ImageNetV2 ground truth as (i) those they agree with (26.92%), (ii) those they disagree with (19.23%), and (iii) those they remain uncertain about and highly doubt (53.85%). Examples of this type of categorization are provided in Fig. 7.

5 Conclusions

Single-label multi-class classification datasets like ImageNet-1k are crucial for advancing deep learning in computer vision. However, as the demand for reliable DNN models grows, it is vital to examine these datasets for biases that could impede progress. We provide a practical framework for smaller research groups to enhance the quality of multi-class classification datasets, especially those that could contain multi-labeled images. Furthermore, we introduce new labels for the ImageNetV2 dataset to account for its multi-label nature. The purpose of our dataset enhancement platform and the provided multi-labels for ImageNetV2 is to facilitate research on the performance degradation of ImageNet-1k-trained

DNN models on the ImageNetV2 dataset. Interestingly, only about half of the 10,000 images in the ImageNetV2 dataset can be confidently categorized as having a single label, thereby underscoring the need for further investigation into the impact of the multi-labeled images on ImageNet-based benchmarks and their potential implications for downstream utilization. Such research endeavors will help us better understand how models perform on complex vision datasets.

Acknowledgment. This research was supported by Ghent University Global Campus (GUGC) in Korea. This research was also supported under the National Research Foundation of Korea (NRF), (2020K1A3A1A68093469), funded by the Korean Ministry of Science and ICT (MSIT). We want to specifically thank the following people for their contribution to the annotation process: Gayoung Lee, Gyubin Lee, Herim Lee, Hyesoo Hong, Jihyung Yoo, Jin-Woo Park, Kangmin Kim, Jihyung Yoo, Jongbum Won, Sohee Lee, Sohn Yerim, Taeyoung Choi, Younghyun Kim, Yujin Cho, and Wonjun Yang.

References

1. Krizhevsky, A., Sutskever, I., Hinton, G.: ImageNet classification with deep convolutional neural networks. *Commun. ACM* **60**(6), 84–90 (2017). <https://doi.org/10.1145/3065386>
2. Wightman, R.: Pytorch image models. GitHub (2019). <https://github.com/huggingface/pytorch-image-models/blob/main/results/results-imagenet.csv>
3. Ozbulak, U., et al.: Know your self-supervised learning: a survey on image-based generative and discriminative training. *Trans. Mach. Learn. Res.* (2023). <https://openreview.net/forum?id=Ma25S4ludQ>
4. Beyer, L., Hénaff, O., Kolesnikov, A., Zhai, X., Oord, A.: Are we done with ImageNet? arXiv preprint (2020). <http://arxiv.org/abs/2006.07159>
5. Tsipras, D., Santurkar, S., Engstrom, L., Ilyas, A., Madry, A.: From ImageNet to image classification: contextualizing progress on benchmarks. In: 37th International Conference on Machine Learning, Article no. 896, pp. 9625–9635 (2020). <https://dl.acm.org/doi/10.5555/3524938.3525830>
6. Vasudevan, V., Caine, B., Gontijo-Lopes, R., Fridovich-Keil, S., Roelofs, R.: When does dough become a bagel? Analyzing the remaining mistakes on ImageNet. In: NeurIPS (2022). <https://openreview.net/pdf?id=mowt1WNhTC7>
7. Recht, B., Roelofs, R., Schmidt, L., Shankar, V.: Do ImageNet classifiers generalize to ImageNet? In: 36th International Conference on Machine Learning (2019). <http://proceedings.mlr.press/v97/recht19a/recht19a.pdf>
8. Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D., Steinhardt, J., Madry, A.: Identifying statistical bias in dataset replication. In: 37th International Conference on Machine Learning (2020). <http://proceedings.mlr.press/v119/engstrom20a/engstrom20a.pdf>
9. Anzaku, E., Wang, H., Van Messem, A., De Neve, W.: A principled evaluation protocol for comparative investigation of the effectiveness of DNN classification models on similar-but-non-identical datasets. arXiv preprint (2022). <http://arxiv.org/abs/2209.01848>
10. Shankar, V., Roelofs, R., Mania, H., Fang, A., Recht, B., Schmidt, L.: Evaluating machine accuracy on ImageNet. In: 37th International Conference on Machine Learning, vol. 119, pp. 8634–8644 (2020). <https://proceedings.mlr.press/v119/shankar20c.html>

11. Northcutt, C., Athalye, A., Mueller, J.: Pervasive label errors in test sets destabilize machine learning benchmarks. In: Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks (2021). <https://openreview.net/pdf?id=XccDXrDNLek>
12. Luccioni, A., Rolnick, D.: Bugs in the data: how imagenet misrepresents biodiversity. In: Proceedings of the Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, Article no. 1613, pp. 14382–14390 (2023). <https://dl.acm.org/doi/10.1609/aaai.v37i12.26682>
13. Fang, Y., Sun, Q., Wang, X., Huang, T., Wang, X., Cao, Y.: Eva-02: a visual representation for neon genesis. arXiv preprint (2023). <https://doi.org/10.48550/arXiv.2303.11331>