

NOVEMBER 11, 2023

# LEVERAGING HUMAN-MACHINE INTERACTIONS FOR COMPUTER VISION DATASET QUALITY ENHANCEMENT

Esla Timothy Anzaku, Hyesoo Hong, Jin-Woo Park, Wonjun Yang, Kangmin Kim, JongBum Won, Deshika Vinoshani Kumari Herath, Arnout Van Messem, Wesley De Neve



FACULTY OF ENGINEERING  
AND ARCHITECTURE



# OVERVIEW

- 1 THE IMAGENET DATASET
- 2 THE SINGLE-LABEL ASSUMPTION
- 3 WHY REVISIT THE SINGLE-LABEL ASSUMPTION?
- 4 OUR FRAMEWORK – MULTILABELFY
- 5 FINAL THOUGHTS

# THE IMAGENET DATASET

# INTRODUCING IMAGENET AND IMAGENET-1k

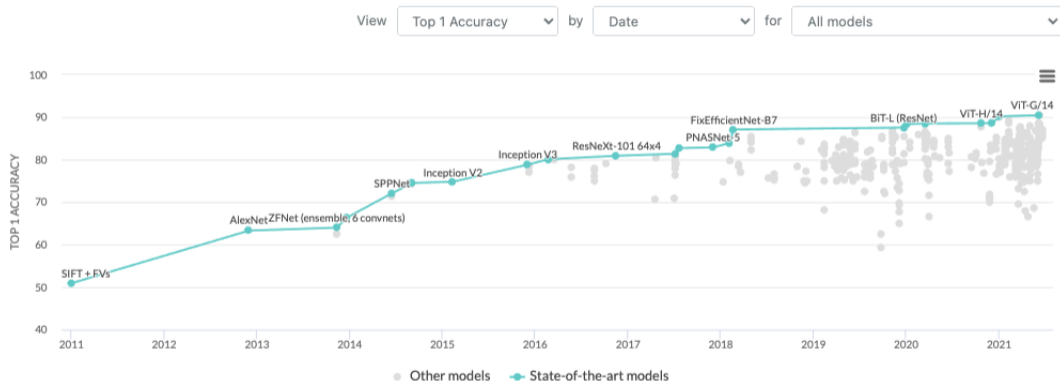
- **ImageNet**: Largest visual dataset for object recognition.
- Over 14 million images across approximately 22k categories.
- **ImageNet-1k**: A subset with 1k categories and over 1 million images.
  - Used for the ImageNet Large Scale Visual Recognition Challenge (ILSVRC).
  - Spans categories from 'dogs' and 'plants' to 'building' and 'vehicles'
  - Central to major deep learning breakthroughs.
    - Example: Transfer Learning
  - Benchmark for model evaluation in computer vision.
    - Example: Supervised and Self-supervised Benchmarking





# INTRODUCING IMAGENET AND IMAGENET-1k

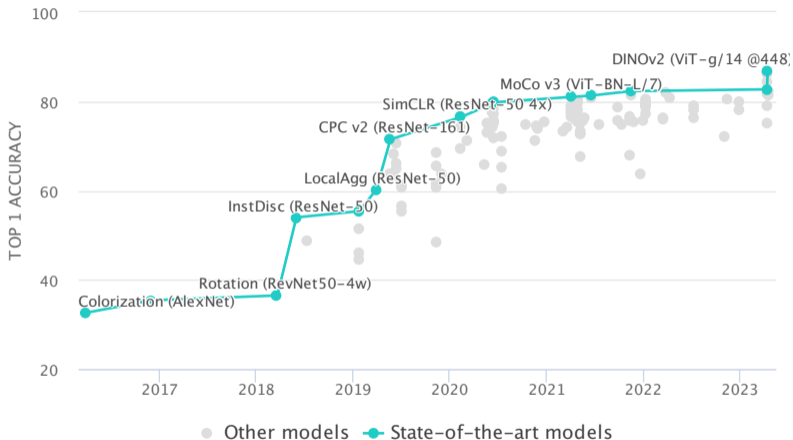
## BENCHMARKING SUPERVISED IMAGE MULTI-CLASS CLASSIFICATION<sup>1</sup>



<sup>1</sup> Source: <https://paperswithcode.com>

# INTRODUCING IMAGENET AND IMAGENET-1k

## BENCHMARKING SELF-SUPERVISED IMAGE MULTI-CLASS CLASSIFICATION<sup>2</sup>



<sup>2</sup>Source: <https://paperswithcode.com/sota/self-supervised-image-classification-on>

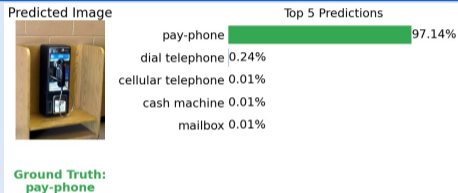
# THE SINGLE-LABEL ASSUMPTION

# THE SINGLE-LABEL ASSUMPTION IN IMAGENET-1k

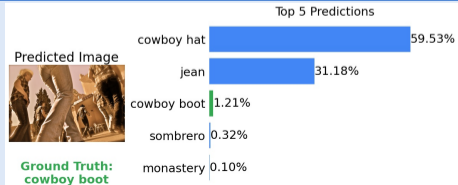
## IMPLICATIONS FOR METRIC ACCURACY AND MODEL EVALUATION

- **Single-label Assumption:** Each image in ImageNet-1k is annotated with single label.
- Common metrics: *Top-1* and *Top-5* accuracies.
  - *Top-1 Accuracy:* The model's prediction matches the ground truth.
  - *Top-5 Accuracy:* The true label is among the model's top 5 predictions.
- **Assuming single-label correctness could skew evaluations, impacting not just top-1 and top-5 metrics but also Precision, Recall, ROC AUC, Negative Log Likelihood, ECE, and more.**

**Top-1 correctness:** Ground truth = topmost prediction



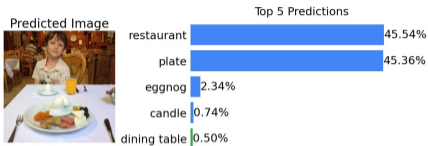
**Top-5 correctness:** Ground truth among topmost 5 predictions



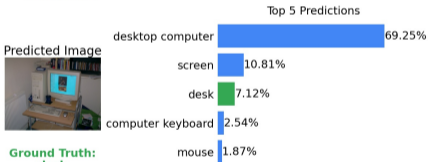
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION?

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



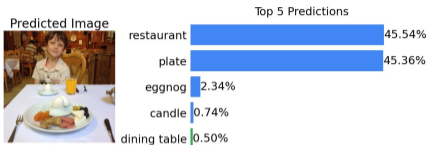
Ground Truth:  
dining table



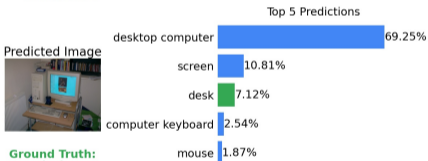
Ground Truth:  
desk

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

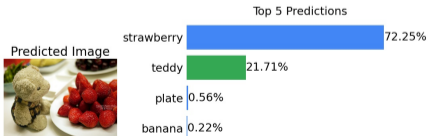
## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



Ground Truth:  
dining table



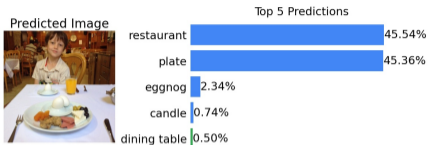
Ground Truth:  
desk



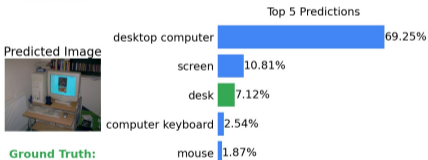
Ground Truth:  
teddy

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

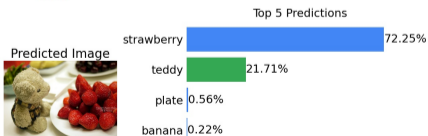
## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



Ground Truth:  
dining table



Ground Truth:  
desk

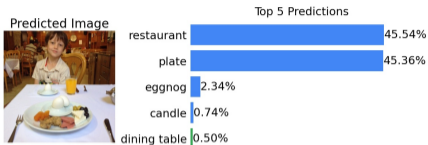


Ground Truth:  
teddy

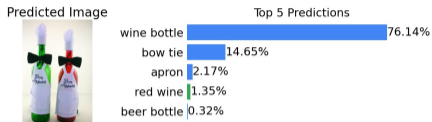


# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

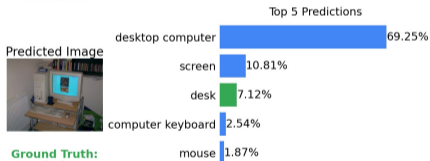
## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



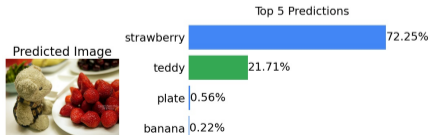
Ground Truth:  
dining table



Ground Truth:  
red wine



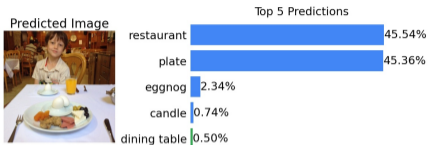
Ground Truth:  
desk



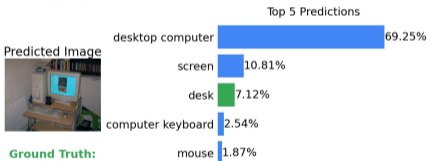
Ground Truth:  
teddy

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

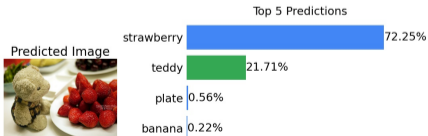
## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



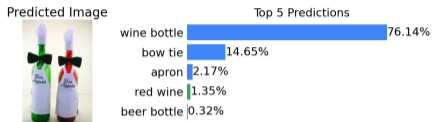
Ground Truth:  
dining table



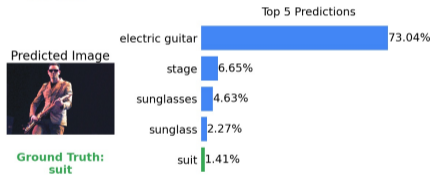
Ground Truth:  
desk



Ground Truth:  
teddy



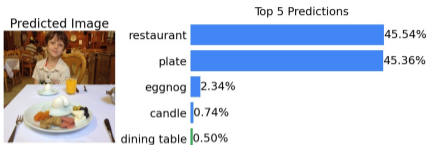
Ground Truth:  
red wine



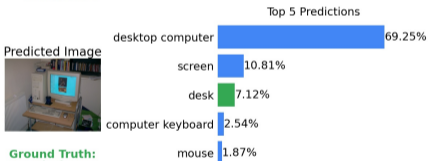
Ground Truth:  
suit

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

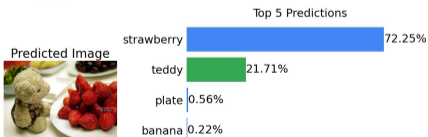
## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES



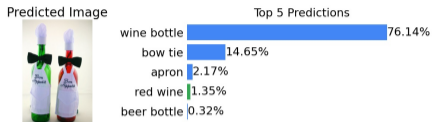
Ground Truth:  
dining table



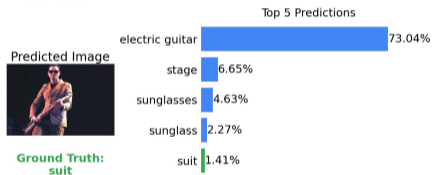
Ground Truth:  
desk



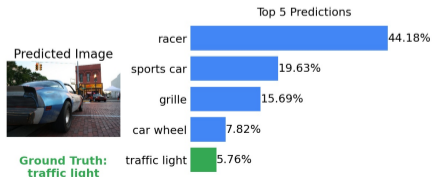
Ground Truth:  
teddy



Ground Truth:  
red wine



Ground Truth:  
suit



Ground Truth:  
traffic light

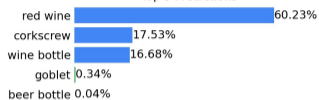
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions



Ground Truth:  
goblet

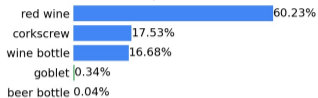
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions

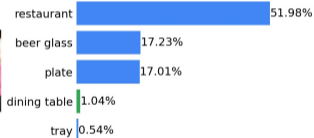


**Ground Truth:**  
goblet

Predicted Image



Top 5 Predictions



**Ground Truth:**  
dining table

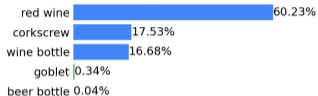
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions

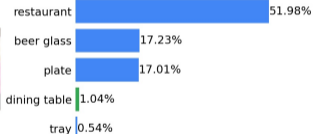


Ground Truth:  
goblet

Predicted Image



Top 5 Predictions

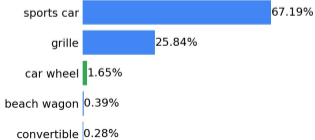


Ground Truth:  
dining table

Predicted Image



Top 5 Predictions



Ground Truth:  
car wheel

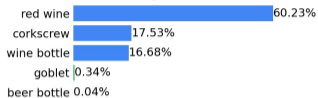
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions

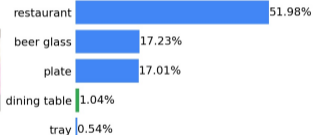


Ground Truth:  
goblet

Predicted Image



Top 5 Predictions

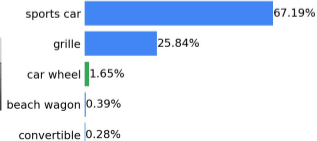


Ground Truth:  
dining table

Predicted Image



Top 5 Predictions

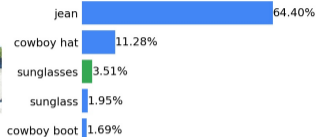


Ground Truth:  
car wheel

Predicted Image



Top 5 Predictions



Ground Truth:  
sunglasses

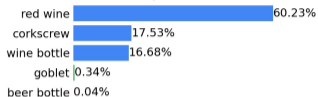
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions

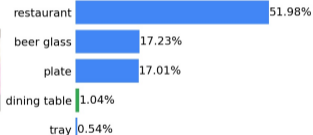


Ground Truth:  
goblet

Predicted Image



Top 5 Predictions

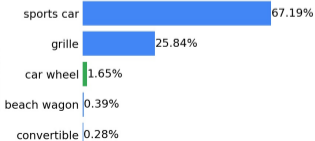


Ground Truth:  
dining table

Predicted Image



Top 5 Predictions

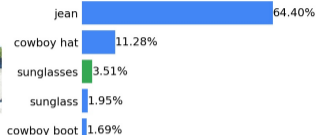


Ground Truth:  
car wheel

Predicted Image



Top 5 Predictions

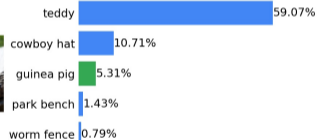


Ground Truth:  
sunglasses

Predicted Image



Top 5 Predictions



Ground Truth:  
guinea pig



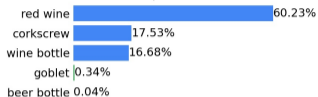
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (1/4)

## QUALITATIVE OBSERVATIONS: CONTRARY EXAMPLES

Predicted Image



Top 5 Predictions

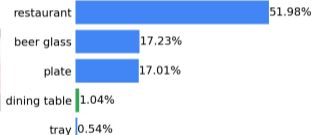


Ground Truth:  
goblet

Predicted Image



Top 5 Predictions

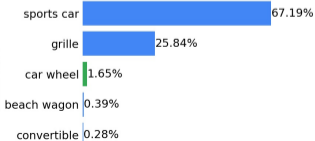


Ground Truth:  
dining table

Predicted Image



Top 5 Predictions

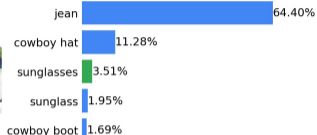


Ground Truth:  
car wheel

Predicted Image



Top 5 Predictions

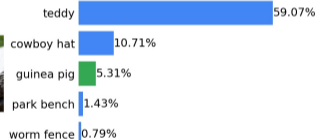


Ground Truth:  
sunglasses

Predicted Image



Top 5 Predictions

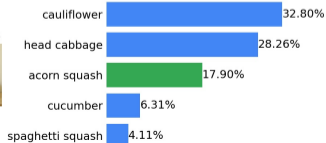


Ground Truth:  
guinea pig

Predicted Image



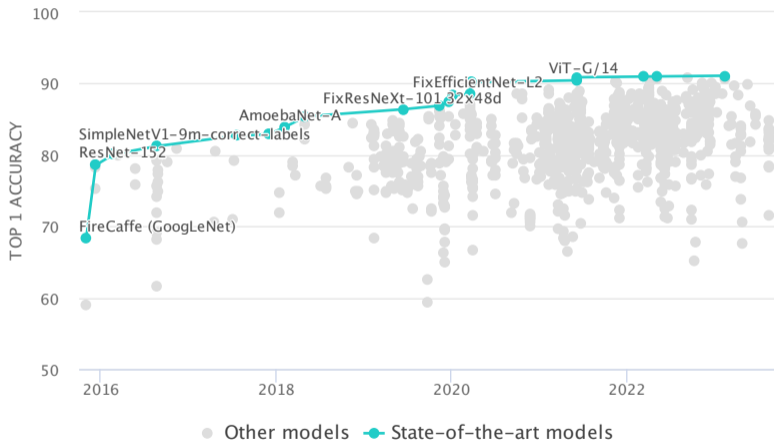
Top 5 Predictions



Ground Truth:  
acorn squash

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (2/4)

ACCURACY SATURATION: IS SOMETHING WRONG WITH THE DATA?<sup>3</sup>



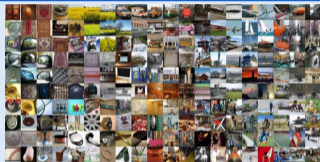
**Regardless of model architecture, training technique, dataset, and model size**

<sup>3</sup> Source: <https://paperswithcode.com/sota/image-classification-on-imagenet>

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (3/4)

UNEXPECTED ACCURACY DEGRADATION ON IMAGENET V2 DATASET

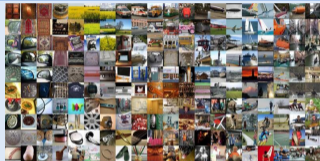
ImageNet validation set (50k Images)



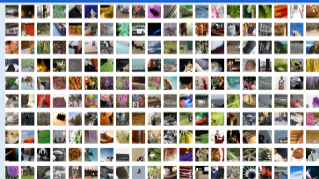
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (3/4)

## UNEXPECTED ACCURACY DEGRADATION ON IMAGENET V2 DATASET

### ImageNet validation set (50k Images)



### ImageNet V2<sup>a</sup> (10k Images)

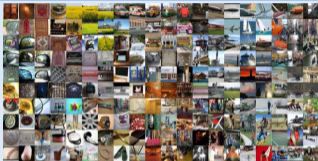


<sup>a</sup>Recht et. al., Do ImageNet Classifiers Generalize to ImageNet? (2019)

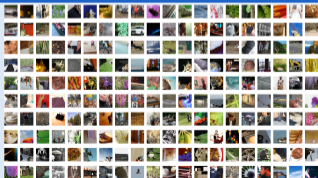
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (3/4)

## UNEXPECTED ACCURACY DEGRADATION ON IMAGENET V2 DATASET

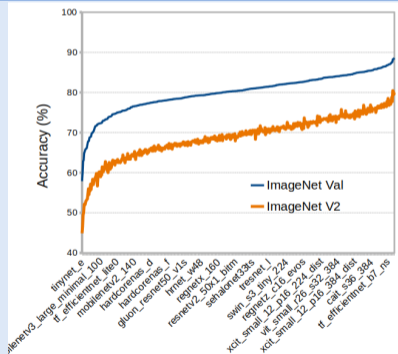
### ImageNet validation set (50k Images)



### ImageNet V2<sup>a</sup> (10k Images)



### Degradation Consistent Across 591 Models



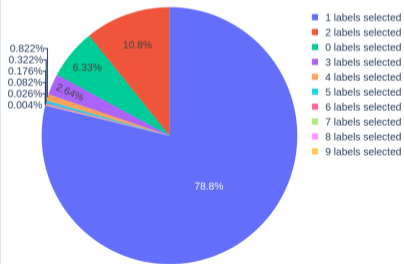
<sup>a</sup>Recht et. al., Do ImageNet Classifiers Generalize to ImageNet? (2019)

# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (4/4)

## PUBLISHED WORK<sup>4</sup> ON THE MULTI-LABEL NATURE OF IMAGENET VALIDATION SET

- Reassessed ImageNet validation labels (50k images)
- **Task:** Identify all distinct objects in each image

### Multi-label Proportions



<sup>4</sup> Source: Tsipras et. al., From ImageNet to ImageNet Classification: Contextualizing Progress on Benchmarks (2020).

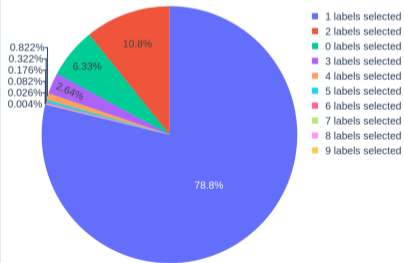
# WHY REVISIT THE SINGLE-LABEL ASSUMPTION? (4/4)

## PUBLISHED WORK<sup>4</sup> ON THE MULTI-LABEL NATURE OF IMAGENET VALIDATION SET

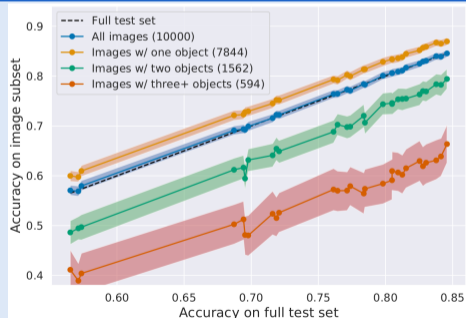
- Reassessed ImageNet validation labels (50k images)
- **Task:** Identify all distinct objects in each image

- Five annotators re-labeled the ImageNet-1k val. set
- **Full test set:** 50k images of the ImageNet validation set
- **All images:** 10k randomly selected images from the full val. set

### Multi-label Proportions



### Subset Accuracy



<sup>4</sup> Source: Tsipras et. al., From ImageNet to ImageNet Classification: Contextualizing Progress on Benchmarks (2020).

# OUR FRAMEWORK – MULTILABELFY



# WHY THE NEED FOR MULTILABELFY?

## DATASET ENHANCEMENT CHALLENGES & OPPORTUNITIES

- Annotation is labor-intensive and prone to errors
- Platforms like Mechanical Turk are often out of reach for smaller research groups
- A demand exists for open-sourced and rigorously reviewed dataset enhancement frameworks
- Available pre-trained models can be efficiently leveraged
- A user-friendly interface can greatly improve human-machine synergy



**Our framework aims to leverage the opportunities while mitigating the challenges presented.**

# MULTILABELFY USER INTERFACE



Current Image Index: 7369

Leave your comments here...

20 most likely labels -> **1-5** 6-10 11-15 16-20

cucumber, cuke



bell pepper



zucchini, courgette



grocery store, grocery,  
food market, market




acorn squash



# MULTILABELFY FRAMEWORK OVERVIEW

## Label Proposal Generation




**Label 1**  
tennis ball

**Label 2**  
racket, racquet

**Label 3**  
radiator


⋮

**Label 20**  
solar dish, solar collector, solar furnace



## Human Multi-Label Annotation


Image to Annotated




Current Image Index: 7402

20 Most likely labels --> 1-5 6-10 11-15 16-20

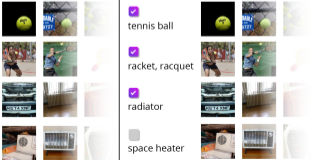
tennis ball racket, racquet radiator space heater pot, flowerpot



## Human Annotation Refinement



Before Refinement	After Refinement
<input checked="" type="checkbox"/> tennis ball	<input checked="" type="checkbox"/> tennis ball
<input checked="" type="checkbox"/> racket, racquet	<input checked="" type="checkbox"/> racket, racquet
<input checked="" type="checkbox"/> radiator	<input checked="" type="checkbox"/> radiator
<input checked="" type="checkbox"/> space heater	<input type="checkbox"/> space heater



## Annotation Disagreement Analysis

Annotator 1	Annotator 2	Disagreement
<input checked="" type="checkbox"/> tennis ball <input checked="" type="checkbox"/> racket, racquet	<input checked="" type="checkbox"/> tennis ball <input checked="" type="checkbox"/> racket, racquet	<b>X</b>
<input checked="" type="checkbox"/> tennis ball <input checked="" type="checkbox"/> racket, racquet <input checked="" type="checkbox"/> space heater	<input checked="" type="checkbox"/> tennis ball <input checked="" type="checkbox"/> racket, racquet	<b>O</b>

# DATA ENHANCEMENT CASE STUDY WITH MULTILABELFY

## RE-LABELING IMAGENET V2: EXPERIMENT SUMMARY

### ■ Stages1: Label Proposal Generation (Automated)

- Pre-trained Model Used: EVA-02<sup>5</sup> (**Top-1**: 90.05%; **Top-5**: 99.05%)
  - DNN Architecture: Vision Transformer
  - Trained on 38 million images
  - First fine-tuned on ImageNet-22k then fine-tuned on ImageNet-1k
- Generates top 20 candidate labels per image

---

<sup>5</sup> Source: Sun et. al., A Visual Representation for Neon Genesis (2023)

# DATA ENHANCEMENT CASE STUDY WITH MULTILABELFY

## RE-LABELING IMAGENET V2: EXPERIMENT SUMMARY

### ■ **Stages1: Label Proposal Generation (Automated)**

- Pre-trained Model Used: EVA-02<sup>5</sup> (**Top-1**: 90.05%; **Top-5**: 99.05%)
  - DNN Architecture: Vision Transformer
  - Trained on 38 million images
  - First fine-tuned on ImageNet-22k then fine-tuned on ImageNet-1k
- Generates top 20 candidate labels per image

### ■ **Stage 2: Human Multi-Label Annotation**

- 14 annotators of various experience levels with computer vision and ImageNet dataset
- All underwent training on the nuances of the task
- Each image was annotated by two annotators

---

<sup>5</sup> Source: Sun et. al., A Visual Representation for Neon Genesis (2023)

# DATA ENHANCEMENT CASE STUDY WITH MULTILABELFY

## RE-LABELING IMAGENET V2: EXPERIMENT SUMMARY

### ■ **Stages1: Label Proposal Generation (Automated)**

- Pre-trained Model Used: EVA-02<sup>5</sup> (**Top-1**: 90.05%; **Top-5**: 99.05%)
  - DNN Architecture: Vision Transformer
  - Trained on 38 million images
  - First fine-tuned on ImageNet-22k then fine-tuned on ImageNet-1k
- Generates top 20 candidate labels per image

### ■ **Stage 2: Human Multi-Label Annotation**

- 14 annotators of various experience levels with computer vision and ImageNet dataset
- All underwent training on the nuances of the task
- Each image was annotated by two annotators

### ■ **Stage3: Annotation Disagreement Analysis (Automated)**

- 6,425 images were selected for the next refinement stage

---

<sup>5</sup> Source: Sun et. al., A Visual Representation for Neon Genesis (2023)

# DATA ENHANCEMENT CASE STUDY WITH MULTILABELFY

## RE-LABELING IMAGENET V2: EXPERIMENT SUMMARY

### ■ **Stages1: Label Proposal Generation (Automated)**

- Pre-trained Model Used: EVA-02<sup>5</sup> (**Top-1**: 90.05%; **Top-5**: 99.05%)
  - DNN Architecture: Vision Transformer
  - Trained on 38 million images
  - First fine-tuned on ImageNet-22k then fine-tuned on ImageNet-1k
- Generates top 20 candidate labels per image

### ■ **Stage 2: Human Multi-Label Annotation**

- 14 annotators of various experience levels with computer vision and ImageNet dataset
- All underwent training on the nuances of the task
- Each image was annotated by two annotators

### ■ **Stage3: Annotation Disagreement Analysis (Automated)**

- 6,425 images were selected for the next refinement stage

### ■ **Stage 4: Human Annotation Refinement**

- 5 annotators participated; 4 of them had participated in Stage 2.
- Only 129 of 10k images remained unlabeled after this stage.

---

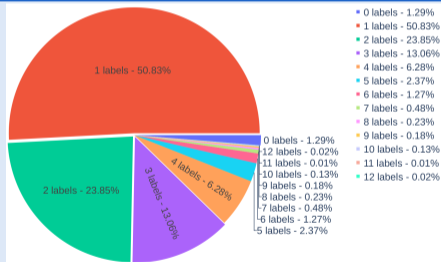
<sup>5</sup> Source: Sun et. al., A Visual Representation for Neon Genesis (2023)

# DATA ENHANCEMENT CASE STUDY WITH MULTILABELFY

## RE-LABELING IMAGE NET V2: KEY RESULTS

- **About 50% images have more than one valid label**

### Multi-label Proportions



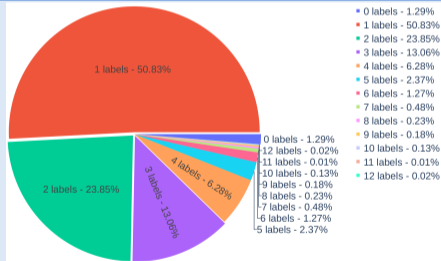


# DATA ENHANCEMENT CASESTUDY WITH MULTILABELFY

## RE-LABELING IMAGE NET V2: KEY RESULTS

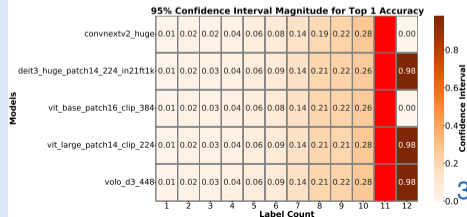
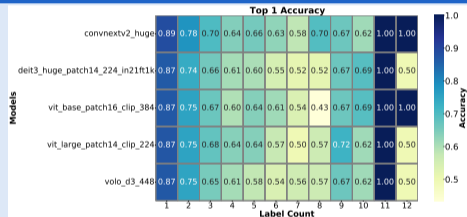
- About 50% images have more than one valid label

### Multi-label Proportions



- Label count negatively correlates with top-1 accuracy

### Top-1 Accuracy versus Label Count



# FINAL THOUGHTS



# REEVALUATING THE SINGLE-LABEL ASSUMPTION

## WHY EMBRACING MULTI-LABEL REALITIES MATTERS

- **To Reflect Real-World Complexities**

- Ensure future labeling reflects real-world complexities
- Our DNN models are already hinting at the disconnect

# REEVALUATING THE SINGLE-LABEL ASSUMPTION

## WHY EMBRACING MULTI-LABEL REALITIES MATTERS

- **To Reflect Real-World Complexities**

- Ensure future labeling reflects real-world complexities
- Our DNN models are already hinting at the disconnect

- **To Enhance Model Evaluation**

- Capture true model capabilities without bias
- Prevent unfair penalization of models for valid alternative predictions

# REEVALUATING THE SINGLE-LABEL ASSUMPTION

## WHY EMBRACING MULTI-LABEL REALITIES MATTERS

### ■ **To Reflect Real-World Complexities**

- Ensure future labeling reflects real-world complexities
- Our DNN models are already hinting at the disconnect

### ■ **To Enhance Model Evaluation**

- Capture true model capabilities without bias
- Prevent unfair penalization of models for valid alternative predictions

### ■ **To Inform Data Collection and Labeling**

- Advocate for datasets that allow DNNs to demonstrate their full potential
- Encourage the incorporation of a broader spectrum of labels

# REEVALUATING THE SINGLE-LABEL ASSUMPTION

## WHY EMBRACING MULTI-LABEL REALITIES MATTERS

### ■ **To Reflect Real-World Complexities**

- Ensure future labeling reflects real-world complexities
- Our DNN models are already hinting at the disconnect

### ■ **To Enhance Model Evaluation**

- Capture true model capabilities without bias
- Prevent unfair penalization of models for valid alternative predictions

### ■ **To Inform Data Collection and Labeling**

- Advocate for datasets that allow DNNs to demonstrate their full potential
- Encourage the incorporation of a broader spectrum of labels

### ■ **To Fuel Progress in the Field**

- Foster innovation with more accurate and holistic model evaluations

# REEVALUATING THE SINGLE-LABEL ASSUMPTION

## WHY EMBRACING MULTI-LABEL REALITIES MATTERS

### ■ **To Reflect Real-World Complexities**

- Ensure future labeling reflects real-world complexities
- Our DNN models are already hinting at the disconnect

### ■ **To Enhance Model Evaluation**

- Capture true model capabilities without bias
- Prevent unfair penalization of models for valid alternative predictions

### ■ **To Inform Data Collection and Labeling**

- Advocate for datasets that allow DNNs to demonstrate their full potential
- Encourage the incorporation of a broader spectrum of labels

### ■ **To Fuel Progress in the Field**

- Foster innovation with more accurate and holistic model evaluations

### ■ **To Boost Reliability and Trust**

- Promote rigorous validation for consistent real-world performance
- Establish more reliable benchmarks to inspire stakeholder confidence

## FUTURE RESEARCH INTERESTS

- **What are the costs of the single-label assumption?**
  - How does this assumption contribute to the **surprising brittleness** of DNN models?
  - What are the costs of utilizing **powerful models** on **simplified assumptions**?
  - To what extent does the single-label assumption foster **overfitting** to **dataset idiosyncrasies**?
  - Could challenging the single-label assumption stimulate a renewed discussion on **nuanced model evaluation**?



**THANK YOU!**

Esla Timothy Anzaku  
eslatimothy.anzaku@ugent.be