
Automatic Dataset Construction (ADC): Sample Collection, Data Curation, and Beyond

Anonymous Author(s)

Affiliation

Address

email

Abstract

Large-scale data collection is essential for developing personalized training data, mitigating the shortage of training data, and fine-tuning specialized models. However, creating high-quality datasets quickly and accurately remains a challenge due to annotation errors, the substantial time and costs associated with human labor. To address these issues, we propose Automatic Dataset Construction (ADC), an innovative methodology that automates dataset creation with negligible cost and high efficiency. Taking the image classification task as a starting point, ADC leverages LLMs for the detailed class design and code generation to collect relevant samples via search engines, significantly reducing the need for manual annotation and speeding up the data generation process. Despite these advantages, ADC also encounters real-world challenges such as label errors (label noise) and imbalanced data distributions (label bias). We provide open-source software that incorporates existing methods for label error detection, robust learning under noisy and biased data, ensuring a higher-quality training data and more robust model training procedure. Furthermore, we design three benchmark datasets focused on label noise detection, label noise learning, and class-imbalanced learning. These datasets are vital because there are few existing datasets specifically for label noise detection, despite its importance. Finally, we evaluate the performance of existing popular methods on these datasets, thereby facilitating further research in the field.

1 Introduction

In the era of Large Language Models (LLMs), the literature has observed an escalating demand for fine-tuning specialized models [1, 2, 3], highlighting the urgent need for customized datasets [4, 5, 6].

Traditional Dataset Construction (TDC) typically involves sample collection followed by labor-intensive annotation, requiring significant human efforts [7, 8, 9, 10]. Consequently, TDC is often hindered by the limitations of human expertise, leading to suboptimal design [11], data inaccuracies [12, 13, 14, 7, 15], and extensive manual labor [16, 17]. Furthermore, certain datasets are inherently challenging or risky to collect manually, such as those for fall detection in elderly individuals, dangerous activities like extreme sports, and network intrusion detection. Therefore, there is a growing need for more automated and efficient data collection methods to enhance accuracy and efficiency in dataset creation [18, 19, 20]. To address these challenges, we propose the **Automatic Dataset Construction (ADC)**, an innovative approach designed to construct customized large-scale datasets with minimal human involvement. Our methodology reverses the traditional process by starting with detailed annotations that guide sample collection. This significantly reduces the

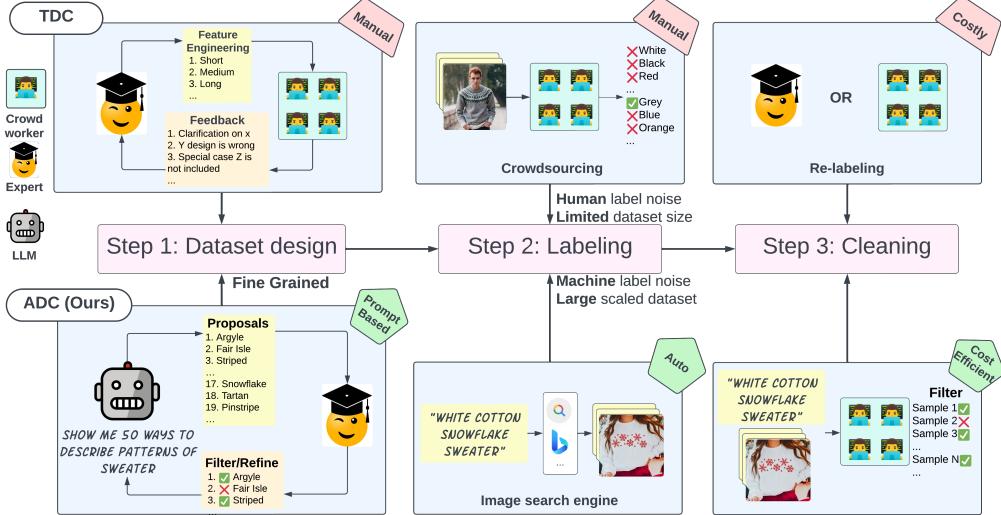


Figure 1: Comparisons of key steps in dataset construction. In **Step 1: Dataset design**, ADC utilizes LLMs to search the field and provide instant feedback, unlike traditional methods that rely on manual creation of class names and refine through crowdsourced worker feedback. In **Step 2: Labeling**, ADC reduces human workload by flipping the data collection process, using targets to search for samples. In **Step 3: Cleaning**, ADC instructs human labor to filter noisy labeled samples from previous steps, instead of relabeling.

34 workload, time, and cost associated with human annotation, making the process more efficient and
 35 targeted for LLM applications, ultimately outperforming traditional methods.

36 **Traditional-Dataset-Construction v.s. Automatic Dataset Construction** Figure 1 illustrates the
 37 difference between Traditional Dataset Construction (TDC) and Automatic Dataset Construction
 38 (ADC). TDC typically unfolds in two stages: developing classification categories and employing
 39 human labor for annotation. Creating comprehensive categories requires deep domain knowledge
 40 and experience, tasks that even expert researchers find challenging [11]. Crowdsourcing is often used
 41 to refine these categories, but it increases time and costs without necessarily improving label quality
 42 [16, 17]. Annotation by human workers introduces label noise, which impacts dataset reliability, even
 43 when multiple inputs are aggregated [21]. In contrast, ADC offers improvements at each key step.
 44 In the “Dataset design”, ADC uses LLMs to automate field searches and provide instant feedback,
 45 unlike traditional manual class and attribute creation. In the sample annotation steps, ADC reverses
 46 the labeling process by using predefined targets to search for samples, human annotators are then
 47 instructed to filter noisy labeled samples, significantly reducing the need for costly human annotation.

48 Our main contributions can be summarized as follows:

- 49 • **The Automatic-Dataset-Construction (ADC) pipeline:** We introduce Automatic-Dataset-
 50 Construction (ADC), an automatic data collection approach that requires minimal human efforts,
 51 tailored for the specialized large-scale data collection. Leveraging ADC, we developed Clothing-
 52 ADC, an image dataset containing one million images with over 1,000 subclasses for each clothing
 53 type. Our dataset offers a rich hierarchy of categories, creating well-defined sub-populations that
 54 support research on a variety of complex and novel tasks.
- 55 • **Software efforts for addressing dataset construction challenges:** We explore several challenges
 56 observed in real-world dataset construction, including detecting label errors, learning with noisy
 57 labels, and class-imbalanced learning. To improve the quality of the constructed data and model
 58 training, we provide well-written software that incorporates existing solutions to these challenges.
- 59 • **Benchmark efforts and Comprehensive Service:** To further facilitate the exploration of the
 60 aforementioned challenges (label noise detection and learning, class-imbalanced learning), we cus-
 61 tomize three benchmark subsets and provide benchmark performances of the implemented methods
 62 in our software. This offers researchers a platform for performance comparisons, enhancing the
 63 evaluation and refinement of their approaches.

64 **2 Automatic-Dataset-Construction (ADC)**

65 In this section, we discuss the detailed procedure of ADC, as well as an empirical application.

66 **2.1 The ADC pipeline**

67 The ADC pipeline generates datasets with finely-grained class and attribute labels, utilizing data
68 diagnostic software to perform data curation. Below, we provide a step-by-step guide to collecting the
69 Clothing-ADC, a clothes image dataset, along with an overview of its statistics and key information.
70 The overall Automatic-Dataset-Construction (ADC) pipeline is illustrated in Figure 1.

71 **Step 1: Dataset design with large language models (LLM)**

72 • **Customization:** The LLM is tailored to understand specific domains (i.e., fashion) via In-Context-
73 Learning based prompts to leverage its pre-existing knowledge, or fine-tuning on the customized
74 dataset. This enables the LLM to significantly simplify the design phase, making the task easier
75 for designers by allowing them to select from existing options rather than creating from scratch.
76 Given the extensive availability of "clothes"-related domain knowledge on the internet, we skip
77 this procedure for Clothing-ADC.

78 • **Detailed categories identification:** LLMs assist researchers in conducting a more thorough search
79 in the field by processing and analyzing numerous concepts simultaneously, unlike humans who
80 may overlook certain factors when faced with a large volume of concepts [11]. We utilize LLMs to
81 identify attribute types for each class. Then use a prompt of "*Show me <30-80> ways to describe*
82 *<Attribute> of <Class>*" to generate the proposed subclasses.

83 • **Iterative refinement:** The initial category list generated by the LLM undergoes review and
84 refinement either by domain experts or through self-examination by the LLM itself, ensuring
85 alignment with specific application or research needs, as shown in Figure 2. This iterative
86 refinement process enables the creation of a high-quality dataset with finely-grained class labels.
87 Additionally, this approach facilitates rapid iterative feedback during the design phase, offering
88 a significant advantage over traditional methods that rely on annotator feedback during the test
89 run annotation phase. This acceleration enables researchers to explore and refine their ideas more
90 efficiently, resulting in better dataset quality and reduced development time.

Help me pick the best color options to describe sweaters from the many choices I've found:
Ash, black, blue ... Green, Grey ... Indigo, Ivory ... <More attribute designs>
First remove the **wrongly labeled** and **repetitive** samples. Then **organize** them by categories.

[LLM]: I've removed the duplicates and categorized the colors for you:
Neutrals: Ash, Beige, Black, Camel, Charcoal, Dark Gray, ...
... <More categories>

Very Good. I would like to have a diverse range of color, help me to get a set of 10 from this list.

[LLM]: I've selected a diverse set of 10 colors for you:
Navy (Blues) Forest Green (Greens) Burgundy (Reds) Mustard (Yellows) Charcoal (Neutrals)
Teal (Greens) Plum (Purples) Light Blue (Blues) Tan (Earth Tones) Coral (Pinks)

Figure 2: Example of using LLM for iterative refinement of attribute designs.

91 **Step 2: Automated querying and downloading** For image data collection, ADC utilizes APIs
92 provided by Google Images or Bing Images for automated querying. Each category and attribute
93 identified in the first step can be used to formulate search queries.

94 **Step 3: Data curation and cleaning** Implement or use existing data curation software capable of
95 identifying and filtering out irrelevant images, such as Docta [22], CleanLab [23], and Snorkel [24],
96 etc. For example, these tools can identify when an item is mislabeled regarding its type, material, or
97 color. Finally, ADC aggregates the suggested labels recommended by the dataset curation software
98 and removes potentially mislabeled or uncertain samples. For illustration, we adopt a data-centric
99 label curation software (Docta) in Algorithm 1. The high-level idea of this algorithm is to estimate
100 the essential label noise transition matrix T_Est without using ground truth labels, achieved through

101 the consensus equations (**Part A**). Following this, Algorithm 1 identifies those corrupted instances
 102 via the cosine similarity ranking score among features as well as a well-tailored threshold based on
 103 the obtained information (i.e., T_Est), and then relabels these instances using KNN-based methods
 104 (**Part B**). For more details, please refer to work [25, 26, 27].

Algorithm 1 Data centric curation (Docta)

```

1: procedure DOCTA(noisyDataset, preTrainedModel)
2:   Part A: Encode images and estimate label noise transition matrix
3:   features  $\leftarrow$  EncodeImages(noisyDataset, preTrainedModel)
4:    $T\_Est \leftarrow$  EstimateTransitionMatrix(features, noisyLabels)
5:   Part B: Identify and relabel corrupted instances
6:   corruptedInstances  $\leftarrow$  SimiFeat-rank(features, noisyLabels,  $T\_Est$ )
7:   curedLabels  $\leftarrow$  KNN-based Relabeling(corruptedInstances)
8:   Return curedLabels
9: end procedure

```

105 **2.2 Clothing-ADC**

106 To illustrate the ADC pipeline, we present the Clothing-ADC dataset, which comprises a substantial
 107 collection of clothing images. The dataset includes 1,076,738 samples, with 20,000 allocated for
 108 evaluation, another 20,000 for testing, and the remaining samples used for training. Each image
 109 is provided at a resolution of 256x256 pixels. The dataset is categorized into 12 primary classes,
 110 encompassing a total of 12,000 subclasses, with an average of 89.73 samples per subclass. Detailed
 111 statistics of the dataset are provided in Table 1. The following subsection elaborates on the dataset
 112 construction process in comprehensive detail.

113 **Subclass design** Utilizing GPT-4, we identified numerous attribute options for each clothing type.
 114 For example, in the case of sweaters, we recognized eight distinct attributes: color, material, pattern,
 115 texture, length, neckline, sleeve length, and fit type. The language model was able to find 30-50
 116 options under each attribute. Our Clothing-ADC dataset includes the three most common attributes:
 117 color, material, and pattern, with each attribute having ten selected options. This results in 1000
 unique subclasses per clothing type. The selected attributes are detailed in Table 6 (Appendix).

Dataset Overview	
Number of Samples	1,076,738
Resolution	256 \times 256
Dataset Split	
Train set(with web noise)	1,036,738
Evaluation set (Clean)	20,000
Test set (Clean)	20,000
Classification Structure	
Main Class	12
Total Subclasses	12,000
Subclass Details	
Attribute (Color)	10
Attribute (Material)	10
Attribute (Pattern)	10
Ave. Samples per attribute	89.73

118 Table 1: Dataset information summary of Clothing-ADC Dataset.

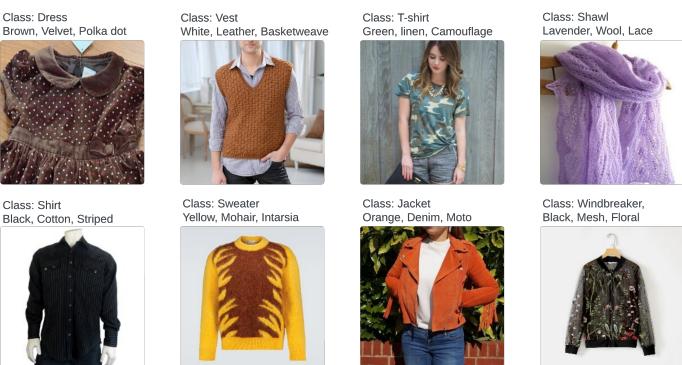


Figure 3: Samples from the collected Clothing-ADC Dataset

119 **Data collection** The ADC pipeline utilizes the Google Image API to collect clothing images by
 120 formulating queries that include attributes such as "Color + Material + Pattern + Cloth Type" (e.g.,
 121 "white cotton fisherman sweater"). Figure 3 shows examples of these queries and the corresponding
 122 images retrieved. The relevance of the search results tends to decline after a significant number of
 123 samples are gathered, leading us to set a cutoff threshold of 100 samples per query. After removing
 124 broken links and improperly formatted images, each subclass retained approximately 90 samples.
 125 These queries generated noisy, webly-labeled data for the training set.

126 **Creating a test set** Note that the collected samples may suffer from web-based label noise, where
 127 annotations might be incorrect due to mismatches provided by search engines, the traditional ap-

128 proach typically involves manually re-labeling existing annotations and aggregating multiple human
 129 votes per label to ensure a high-quality subset for testing purposes. Our ADC pipeline enhances
 130 efficiency by presenting annotators with a set of samples that share the same machine-generated label.
 131 Annotators are then tasked with selecting a subset of correctly labeled samples, choosing a minimum
 132 of four samples out of twenty. This method significantly reduces both manual effort and difficulty,
 133 encouraging annotators to critically evaluate machine-generated labels and thereby reducing the effect
 134 of human over-trust in AI answers [28, 18]. The samples selected through this process are considered
 135 “clean” labels, representing a consensus between human judgment and machine-generated labels [29].
 136 Table 2 provides an insightful comparison between existing datasets and Clothing-ADC. Briefly
 137 speaking, compared with existing datasets, the ADC pipeline is able to help humans without domain
 138 expertise to create fine-grained attributes for the dataset, and automatic annotation and label cleaning
 139 drastically eliminate human effort during label creation.

Table 2: Compare with existing datasets Our ADC pipeline creates a large-scale image classification dataset with a clean test set. Most existing datasets require human effort for labeling, whereas our pipeline can automatically annotate and clean the data. While Clothing-ADC provides fine-grained attribute labels, our dataset design does not require human expertise in the field.

Dataset	# Train/Test	# Classes	Noise Rate(%)	Has Attributes	Auto annotation	Require expert?
iNaturalist [30]	579k/279k	54k	Close to 0	✗	✗	✓
WebVision [14]	2.4M/100k	1000	20	✗	✓	✓
ANIMAL-10N [31]	50k/10k	10	8	✗	✗	✗
CIFAR-10N [9]	50k/10k	10	9.03/25.60/40.21	✗	✗	✗
CIFAR-100N [9]	50k/10k	100	25.6/40.2	✗	✗	✗
Food-101N [32]	75.75k/25.25k	101	18.4	✗	✗	✓
Clothing1M [7]	1M in all	14	38.5	✗	✗	✓
Clothing-ADC (Ours)	1M/20k	12	22.2-32.7	12k	✓	✗

140 3 Challenge one: dealing with imperfect data annotations

141 The first pervasive and critical challenge during the automatic dataset construction lies in the prevalence
 142 of noisy/imperfect labels. This issue is intrinsic to web-sourced data, which, although rich in
 143 diversity, often suffers from inaccuracies due to the uncurated nature of the internet. These errors
 144 manifest as mislabeled images, inconsistent tagging, and misclassified attributes, introducing non-
 145 negligible noise into the dataset that may adversely affect the training and performance of machine
 146 learning models. The following discussion bridges the gap between imperfect data and curated data
 147 via mining and learning with label noise, to refine data quality, enhance label accuracy, and ensure
 148 the reliability of Auto-Dataset-Construction (ADC) for high-stakes AI applications.

149 **Formulation** Let $D := \{(x_n, y_n)\}_{n \in [N]}$ represent the training samples for a K -class classification
 150 task, where $[N] := \{1, 2, \dots, N\}$. Suppose that these samples $\{(x_n, y_n)\}_{n \in [N]}$ are outcomes of the
 151 random variables $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, drawn from the joint distribution \mathcal{D} . Here, \mathcal{X} and \mathcal{Y} denote the
 152 spaces of features and labels, respectively. However, classifiers typically access a noisily labeled
 153 training set $\tilde{D} := \{(x_n, \tilde{y}_n)\}_{n \in [N]}$, assumed to arise from random variables $(X, \tilde{Y}) \in \mathcal{X} \times \tilde{\mathcal{Y}}$, drawn
 154 from the distribution $\tilde{\mathcal{D}}$. It is common to observe instances where $y_n \neq \tilde{y}_n$ for some $n \in [N]$. The
 155 transition from clean to noisy labels is typically characterized by a noise transition matrix $T(X)$,
 156 defined as $T_{i,j}(X) := \mathbb{P}(\tilde{Y} = j | Y = i, X)$ for all $i, j \in [K]$ [12, 13, 33].

157 3.1 The challenge of label noise detection

158 While employing human annotators to clean data is effective in improving label quality, it is often
 159 prohibitively expensive and time-consuming for large datasets. A practical alternative is to enhance
 160 label accuracy automatically by first deploying algorithms to detect potential errors within the dataset
 161 and then correcting these errors through additional algorithmic processing or crowdsourcing.

162 3.1.1 Existing approaches to detect label noise

163 **Learning-Centric Approaches:** Learning-centric approaches often leverage the behavior of models
 164 during training to infer the presence of label errors based on how data is learned. One effective
 165 strategy is confidence-based screening, where labels of training instances are scrutinized if the model’s

166 prediction confidence falls below a certain threshold. This approach assumes that instances with low
 167 confidence scores in the late training stage are likely mislabeled [34]. Another innovative technique
 168 involves analyzing the gradients of the training loss w.r.t. input data. Pruthi et al. [35] utilize gradient
 169 information to detect anomalies in label assignments, particularly focusing on instances where the
 170 gradient direction deviates significantly from the majority of instances. Researchers have also utilized
 171 the memorization effect of deep neural networks, where models tend to learn clean data first and only
 172 memorize noisy labels in the later stages of training. Techniques that track how quickly instances are
 173 learned during training can thus identify noisy labels by focusing on those learned last [36, 37, 38].

174 **Data-Centric Approaches:** Data-centric methods focus on analyzing data features and relationships
 175 rather than model behavior for detection. The ranking-based detection method [39] ranks instances
 176 by the likelihood of label errors based on their alignment with model predictions. An ensemble of
 177 classifiers evaluates each instance, flagging those that consistently deviate from the majority vote
 178 as noisy. Neighborhood Cleaning Rule [40] uses the k -nearest neighbors algorithm to check label
 179 consistency with neighbors, identifying instances whose labels conflict with the majority of their
 180 neighbors as potentially noisy. Zhu et al. [27] propose advanced data-centric strategies for detecting
 181 label noise without training models. Their local voting method uses neighbor consensus to validate
 182 label accuracy, effectively identifying errors based on agreement within the local feature space.

183 3.1.2 Clothing-ADC in label noise detection

184 We prepared a subset of 20,000 samples from the Clothing-ADC dataset for the label noise detection
 185 task, including both noisy and clean labels. We collected three human annotations for each image via
 186 Amazon MTurk. Annotators were instructed to classify the labels as correct, unsure, or incorrect.
 187 Each sample received three votes. Based on these annotations, we determined the noise rate to be
 188 22.2%-32.7%. Using majority vote aggregation implies uncertainty of the label correctness. By using
 189 a more stringent aggregation criterion, more samples are considered as noisy labeled. Under the
 190 extreme case where any doubts from any human annotator can disqualify a sample, our auto collected
 191 dataset still retains 61.3% of its samples. For a detailed distribution of human votes, see Table 7 in
 192 the Appendix.

193 **Benchmark efforts** Detection perfor-
 194 mance comparisons of certain existing
 195 solutions are given in Table 3. We adopt
 196 ResNet-50 [43] as the backbone model

Table 3: F_1 -Score comparisons among several label noise detection methods on Clothing-ADC.

Methods	CORES [41]	CL [34]	Deep k -NN [42]	Simi-Feat [27]
F_1 -Score	0.4793	0.4352	0.3991	0.5721

197 to extract the feature here. For each method, we use the default hyper-parameter reported in the
 198 original papers. All methods are tested on 20,000 points and predict whether the data point is
 199 corrupted or not. We follow [27] to apply the baseline methods to our scenario. In Table 3, the
 200 performance is measured by the F_1 -score of the detected corrupted instances, which is the harmonic
 201 mean of the precision and recall, i.e., $F_1 = \frac{2}{\text{Precision}^{-1} + \text{Recall}^{-1}}$. Let $v_n = 1$ indicate that the n -th label
 202 is detected as a noisy/wrong label, and $v_n = 0$ otherwise. Then, the precision and recall of detecting
 203 noisy labels can be calculated as: Precision = $\frac{\sum_n \mathbb{1}(v_n=1, \tilde{y}_n \neq y_n)}{\sum_n \mathbb{1}(v_n=1)}$, Recall = $\frac{\sum_n \mathbb{1}(v_n=1, \tilde{y}_n \neq y_n)}{\sum_n \mathbb{1}(\tilde{y}_n \neq y_n)}$.

204 3.2 The challenge of learning with noisy labels

205 Another technique is robust learning that can effectively learn from noisy datasets without being
 206 misled by incorrect labels, thus maintaining high accuracy and reliability in real-world applications.

207 3.2.1 Existing approaches to learn with label noise

208 In this subsection, we contribute to the literature with robust learning software; all covered meth-
 209 ods can be mainly summarized into the following three categories: robust loss functions, robust
 210 regularization techniques, and multi-network strategies.

211 **Robust loss designs** Loss Correction modifies the traditional loss function to address label noise by
 212 incorporating an estimated noise transition matrix, thereby recalibrating the model’s training focus
 213 [33]. Loss-Weighting strategies mitigate the impact of noisy labels by assigning lower weights to
 214 likely mislabeled instances, reducing their influence on the learning process [13, 44]. Symmetric
 215 Cross-Entropy Loss balances the contributions of correctly labeled and mislabeled instances, im-
 216 proving the model’s resilience to label discrepancies [45]. Generalized Cross-Entropy Loss, derived

217 from mean absolute error, offers enhanced robustness against outliers and label noise [46]. Peer
 218 Loss Functions form a family of robust loss functions [47, 48, 41], leveraging predictions from peer
 219 samples as regularization to adjust the loss computation, thereby increasing resistance to noise.

220 **Robust Regularization Techniques** Regularization techniques are designed to constrain or modify
 221 the learning process, thereby reducing the model’s sensitivity to label noise. Mixup [49] generates
 222 synthetic training examples by linearly interpolating between pairs of samples and their labels,
 223 enhancing model generalization and smoothing label predictions. Label Smoothing [50, 51] combats
 224 overconfidence in unreliable labels by adjusting them towards a uniform distribution. Negative Label
 225 Smoothing [52] refines this approach by specifically adjusting the smoothing process for negative
 226 labels, preserving model confidence in high-noise environments. Early-Learning Regularization
 227 tackles the issue of early memorization of noisy labels by dynamically adjusting regularization
 228 techniques during the initial training phase [37, 38].

229 **Multi-network strategies** Employing multiple networks can enhance error detection and correction
 230 through mutual agreement and ensemble techniques. In Co-teaching, two networks concurrently train
 231 and selectively share clean data points with each other, mitigating the memorization of noisy labels
 232 [53]. MentorNet [54] equips a student network with a curriculum that emphasizes samples likely
 233 to be clean, as determined by the observed dynamics of a mentor network. DivideMix leverages
 234 two networks to segregate the data into clean and noisy subsets using a mixture model, allowing for
 235 targeted training on each set to manage label noise better [55].

236 3.2.2 Clothing-ADC in label noise learning

237 We provide two versions
 238 of the Label Noise Learn-
 239 ing task, Clothing-ADC
 240 and Clothing-ADC (tiny).
 241 Specifically, Clothing-
 242 ADC leverages the whole
 243 available (noisy) training
 244 samples to construct
 245 the label noise learning
 246 task. The objective is to
 247 perform class prediction
 248 w.r.t. 12 clothes types:
 249 Sweater, Windbreaker,
 250 T-shirt, Shirt, Knitwear,
 251 Hoodie, Jacket, Suit,

Table 4: Experiment results of label noise learning methods on Clothing-ADC and Clothing-ADC (tiny). We report the model prediction accuracy on the held-out clean labeled test set for comparisons.

Methods / Dataset	Clothing-ADC	Clothing-ADC (tiny)
Cross-Entropy	74.76	67.72 ± 0.40
Backward Correction [33]	77.51	70.49 ± 0.06
Forward Correction [33]	78.45	70.60 ± 0.14
(Positive) LS [51]	81.94	70.67 ± 0.15
(Negative) LS [52]	78.65	70.14 ± 0.13
Peer Loss [47]	78.58	70.92 ± 0.17
f -Div [48]	77.43	68.98 ± 0.22
Divide-Mix [55]	77.00	71.58 ± 0.11
Jocor [56]	78.47	70.92 ± 0.09
Co-Teaching [53]	80.49	70.55 ± 0.17
LogitCLIP [57]	77.85	70.16 ± 0.14
TaylorCE [58]	81.87	70.07 ± 0.01

252 Shawl, Dress, Vest, Underwear. We also provide a tiny version of Clothing-ADC, which contains 50K
 253 training images, sharing similar size with certain widely-used ones, i.e., MNIST, Fashion-MNIST,
 254 CIFAR-10, CIFAR-100, etc.

255 **Estimated noise level of Clothing-ADC** We selected a subset of 20,000 training samples and asked
 256 human annotators to evaluate the correctness of the auto-annotated dataset. After aggregating three
 257 votes from annotators, we estimate the noise rate to be 22.2%-32.7%, which consists of 10.5% of the
 258 samples having ambiguity and 22.2% being wrongly labeled. The remaining 77.8% of the samples
 259 were correctly labeled. The detailed distribution of human votes is given in Appendix Table 7.

260 **Benchmark efforts** In this task, we aim to provide the performance comparison among various
 261 learning-with-noisy-label solutions. All methods utilize ResNet-50 as the backbone model and are
 262 trained for 20 epochs to ensure a fair comparison. We report the model prediction accuracy on the
 263 held-out clean labeled test set. For the tiny version, we conduct three individual experiments using
 264 three different random seeds and calculate the mean and standard deviation. As shown in Table 4,
 265 certain methods, such as Positive LS and Taylor CE, significantly outperform Cross-Entropy. These
 266 results underscore the importance and necessity of pairing ADC with robust learning software.

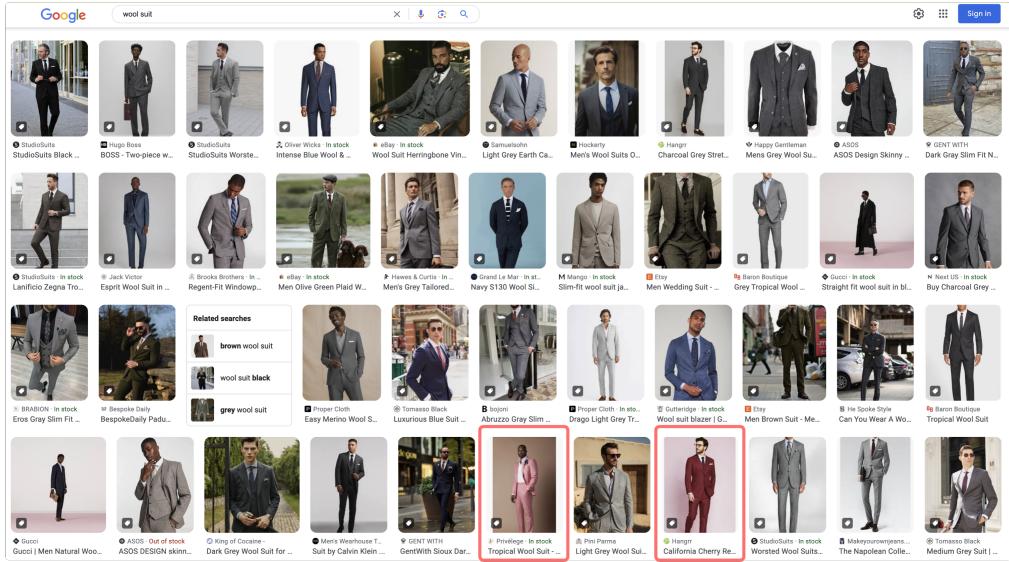


Figure 4: Long-tailed data distribution is a prevalent issue in many datasets. Searching “wool suit” in Google image results in dark wool suits, while only a few are of a light color (red/pink).

267 4 Challenge two: dealing with imbalanced data distribution

268 We now discuss another real-world challenge: when imperfect annotations meet with imbalanced
 269 class/attribute distributions. As shown in Figure 4, long-tailed data distribution is a prevalent issue
 270 in web-based datasets: to collect a dataset of wool suits without a specified target color on Google
 271 Image, the majority would likely be dark or muted shades (grey, black, navy), with few samples in
 272 brighter colors like pink or purple. This natural disparity results in most data points belonging to a
 273 few dominant categories, while the remaining are spread across several minority groups.

274 We are interested in how class-imbalance intervenes with learning. In real-world scenarios, the
 275 distribution of classes tends to form a long-tail form, in other words, the head class and the tail class
 276 differ significantly in their sample sizes, i.e., $\max_k \mathbb{P}(Y = k) \gg \min_{k'} \mathbb{P}(Y = k')$.

277 4.1 Existing approaches for class imbalance learning

278 **Data-Level Methods** Data-level methods modify training data to balance class distribution, fo-
 279 cusing on adjusting the dataset by increasing minority class instances or decreasing majority class
 280 instances. Oversampling increases the number of minority class instances to match or approach the
 281 majority class. This can be done through simple duplication [59] (e.g., random oversampling) or gen-
 282 erating synthetic data [60, 61, 62, 63]. Undersampling reduces the number of majority class instances,
 283 helping to balance class distributions but potentially discarding useful information [64, 65, 66].

284 **Algorithm-Level Methods** These methods adjust the training process or model to handle unequal
 285 class distributions better. Specifically, cost-sensitive learning assigns different costs to misclassifica-
 286 tions of different classes, imposing higher penalties for errors on the minority class [67]. It modifies
 287 the loss function to incorporate misclassification costs, encouraging the model to focus more on
 288 minority class errors [68, 69]. Thresholding adjusts the decision threshold for class probabilities to
 289 account for class imbalance. Instead of using a default threshold, different thresholds are applied
 290 based on class distribution, modifying the decision process for predicting class labels [70, 71].

291 4.2 Clothing-ADC in class-imbalanced learning

292 Note that in the label noise learning task, the class distributes with almost balanced prior. However,
 293 in practice, the prior distribution is often long-tail distributed. Hence, the combined influence of label
 294 noise and long-tail distribution is a new and overlooked challenge presented in the literature. To
 295 facilitate the exploration of class-imbalanced learning, we tried to reduce the impact of noisy labels

296 via selecting high-quality annotated samples as recognized by dataset curation software. Human
 297 estimation suggested a noise rate of up to 22.2%, and 10.5% marked as uncertain. To address this,
 298 we employed two methods to remove noisy samples: a data centric curation (Algorithm 1), which
 299 removed 26.36% of the samples, and a learning-centric curation (Appendix Algorithm 3), which
 300 removed 25%. Combined, these methods eliminated 45.15% of the samples, with an overlap of 6.21%
 301 between the two approaches. We provide Clothing-ADC CLT, which could be viewed as the long-tail
 302 (class-level) distributed version of Clothing-ADC. Denote by ρ the imbalanced ratio between the
 303 maximum number of samples per class and the minimum number of samples per class. In practice,
 304 we provide $\rho = 10, 50, 100$ (class-level) long-tail version of Clothing-ADC.

305 **Benchmark efforts** Regarding the evaluation metric, we follow from the recently proposed metric
 306 [72], which considers an objective that is based on the weighted sum of class-level performances on
 307 the test data, i.e., $\sum_{i \in [K]} g_i \text{Acc}_i$, where Acc_i indicates the accuracy of the class i :

$$\delta\text{-worst accuracy: } \min_{g \in \Delta_K} \sum_{i \in [K]} g_i \text{Acc}_i, \text{ s.t. } D(\mathbf{g}, \mathbf{u}) \leq \delta.$$

308 Here, Δ_K denotes the $(K - 1)$ -dimensional probability simplex, where K is the number of classes
 309 as previously defined. Let $\mathbf{u} \in \Delta_K$ be the uniform distribution, and $\mathbf{g} := [g_1, g_2, \dots, g_K]$ is the class
 310 weights. The δ -worst accuracy measures the worst-case \mathbf{g} -weighted performance with the weights
 311 constrained to lie within the δ -radius ball around the target (uniform) distribution. For any chosen
 312 divergence D , it reduces to the mean accuracy when $\delta = 0$ and to the worst accuracy for $\delta \rightarrow \infty$.
 313 The objective interpolates between these two extremes for other values of δ and captures our goal of
 314 optimizing for variations around target priors instead of more conventional objectives of optimizing
 315 for either the average accuracy at the target prior or the worst-case accuracy.

316 Different from the previous dataset we used in noise learning, we use a cleaner dataset for this
 317 class-imbalance learning to avoid the distractions of noisy labels. The size of this dataset consists of
 318 56,2263 images rather than 1M. The backbone model we use is ResNet-50. For the class distributions
 319 for different ρ , we include them in the Appendix. All the experiments are run for 5 times and we
 320 calculate the mean and standard deviation. With the imbalance ratio going larger, the accuracy
 321 becomes worse, which is expected for a more difficult task.

Table 5: δ -worst accuracy of class-imbalanced learning baselines on clothing-ADC CLT dataset.

Method	$\delta = 0$ Worst Accuracy			$\delta = 1$ Worst Accuracy			$\delta = \infty$ Worst Accuracy		
	$\rho = 10$	$\rho = 50$	$\rho = 100$	$\rho = 10$	$\rho = 50$	$\rho = 100$	$\rho = 10$	$\rho = 50$	$\rho = 100$
Cross Entropy	57.80 ± 0.25	33.85 ± 0.13	30.10 ± 0.22	19.79 ± 0.23	0.35 ± 0.11	0.00 ± 0.00	0.96 ± 0.26	0.00 ± 0.00	0.00 ± 0.00
Focal [73]	72.70 ± 0.19	65.17 ± 0.29	62.28 ± 0.31	49.66 ± 1.09	34.14 ± 1.05	29.12 ± 0.92	38.12 ± 1.76	19.46 ± 1.49	13.44 ± 1.73
LDAM [74]	72.50 ± 0.15	65.70 ± 0.26	63.25 ± 0.35	51.13 ± 0.78	36.86 ± 1.03	30.88 ± 1.07	40.90 ± 1.53	23.24 ± 1.69	15.69 ± 2.13
Bal-Softmax [75]	74.18 ± 0.08	70.48 ± 0.55	69.47 ± 0.44	56.57 ± 0.93	53.37 ± 2.31	44.24 ± 2.83	48.54 ± 2.27	45.64 ± 3.98	50.60 ± 1.40
Logit-Adjust [76]	74.08 ± 0.05	70.94 ± 0.24	69.44 ± 0.18	56.00 ± 1.39	53.93 ± 2.46	49.70 ± 2.64	47.45 ± 2.26	47.76 ± 4.07	43.26 ± 4.69
Post-hoc [76]	62.54 ± 0.11	54.84 ± 0.15	49.63 ± 0.71	35.67 ± 0.49	24.14 ± 1.18	19.00 ± 0.68	22.50 ± 0.78	7.15 ± 1.82	3.81 ± 0.97
Drops [72]	73.66 ± 0.29	69.14 ± 0.38	67.15 ± 0.17	58.12 ± 0.26	47.07 ± 0.74	43.42 ± 1.19	50.85 ± 0.49	36.27 ± 1.15	32.43 ± 1.90

322 Conclusion

323 In this paper, we introduced the Automatic Dataset Construction (ADC) pipeline, a novel approach
 324 for automating the creation of large-scale datasets with minimal human intervention. By leveraging
 325 Large Language Models for detailed class design and automated sample collection, ADC signifi-
 326 cantly reduces the time, cost, and errors associated with traditional dataset construction methods.
 327 The Clothing-ADC dataset, which comprises one million images with rich category hierarchies,
 328 demonstrates the effectiveness of ADC in producing high-quality datasets tailored for complex
 329 research tasks. Despite its advantages, ADC faces challenges such as label noise and imbal-
 330 anced data distributions. We addressed these challenges with open-source tools for error detection and
 331 robust learning. Our benchmark datasets further facilitate research in these areas, ensuring that ADC
 332 remains a valuable tool for advancing machine learning model training.

333 **References**

- 334 [1] Manuela Benary, Xing David Wang, Max Schmidt, Dominik Soll, Georg Hilfenhaus, Mani Nas-
335 sir, Christian Sigler, Maren Knödler, Ulrich Keller, Dieter Beule, et al. Leveraging large language
336 models for decision support in personalized oncology. *JAMA Network Open*, 6(11):e2343689–
337 e2343689, 2023.
- 338 [2] Sebastian Porsdam Mann, Brian D Earp, Nikolaj Møller, Suren Vynn, and Julian Savulescu.
339 Autogen: A personalized large language model for academic enhancement—ethics and proof of
340 principle. *The American Journal of Bioethics*, 23(10):28–41, 2023.
- 341 [3] Stanisław Woźniak, Bartłomiej Koptyra, Arkadiusz Janz, Przemysław Kazienko, and Jan Kocouć.
342 Personalized large language models. *arXiv preprint arXiv:2402.09269*, 2024.
- 343 [4] Jimmy Wu, Rika Antonova, Adam Kan, Marion Lepert, Andy Zeng, Shuran Song, Jeannette
344 Bohg, Szymon Rusinkiewicz, and Thomas Funkhouser. Tidybot: Personalized robot assistance
345 with large language models. *Autonomous Robots*, 47(8):1087–1102, 2023.
- 346 [5] Hanjia Lyu, Song Jiang, Hanqing Zeng, Yinglong Xia, and Jiebo Luo. Llm-rec: Personalized
347 recommendation via prompting large language models. *arXiv preprint arXiv:2307.15780*, 2023.
- 348 [6] Zhaoxuan Tan, Qingkai Zeng, Yijun Tian, Zheyuan Liu, Bing Yin, and Meng Jiang. Democra-
349 tizing large language models via personalized parameter-efficient fine-tuning. *arXiv preprint*
350 *arXiv:2402.04401*, 2024.
- 351 [7] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. Learning from massive
352 noisy labeled data for image classification. In *Proceedings of the IEEE conference on computer*
353 *vision and pattern recognition*, pages 2691–2699, 2015.
- 354 [8] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
355 2009.
- 356 [9] Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning
357 with noisy labels revisited: A study using real-world human annotations. *arXiv preprint*
358 *arXiv:2110.12088*, 2021.
- 359 [10] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaou Tang. Deep learning face attributes in the
360 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 361 [11] Vikram V Ramaswamy, Sunnie SY Kim, Ruth Fong, and Olga Russakovsky. Overlooked factors
362 in concept-based explanations: Dataset choice, concept learnability, and human capability. In
363 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
364 10932–10941, 2023.
- 365 [12] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning
366 with noisy labels. *Advances in neural information processing systems*, 26, 2013.
- 367 [13] Tongliang Liu and Dacheng Tao. Classification with noisy labels by importance reweighting.
368 *IEEE Transactions on pattern analysis and machine intelligence*, 38(3):447–461, 2015.
- 369 [14] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. Webvision database:
370 Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017.
- 371 [15] Jiaheng Wei, Zhaowei Zhu, Hao Cheng, Tongliang Liu, Gang Niu, and Yang Liu. Learning with
372 noisy labels revisited: A study using real-world human annotations. In *International Conference*
373 *on Learning Representations*, 2022.
- 374 [16] Joseph Chee Chang, Saleema Amershi, and Ece Kamar. Revolt: Collaborative crowdsourcing
375 for labeling machine learning datasets. In *Proceedings of the 2017 CHI conference on human*
376 *factors in computing systems*, pages 2334–2346, 2017.

- 377 [17] Todd Kulesza, Saleema Amershi, Rich Caruana, Danyel Fisher, and Denis Charles. Structured
378 labeling for facilitating concept evolution in machine learning. In *Proceedings of the SIGCHI*
379 *Conference on Human Factors in Computing Systems*, pages 3075–3084, 2014.
- 380 [18] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar,
381 Marco Tulio Ribeiro, and Daniel Weld. Does the whole exceed its parts? the effect of ai
382 explanations on complementary team performance. In *Proceedings of the 2021 CHI conference*
383 *on human factors in computing systems*, pages 1–16, 2021.
- 384 [19] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. Is the most
385 accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI*
386 *Conference on Artificial Intelligence*, volume 35, pages 11405–11414, 2021.
- 387 [20] Lei Han, Xiao Dong, and Gianluca Demartini. Iterative human-in-the-loop discovery of
388 unknown unknowns in image datasets. In *Proceedings of the AAAI Conference on Human*
389 *Computation and Crowdsourcing*, volume 9, pages 72–83, 2021.
- 390 [21] Victor S Sheng, Foster Provost, and Panagiotis G Ipeirotis. Get another label? improving
391 data quality and data mining using multiple, noisy labelers. In *Proceedings of the 14th ACM*
392 *SIGKDD international conference on Knowledge discovery and data mining*, pages 614–622,
393 2008.
- 394 [22] Docta. <https://docta.ai/>.
- 395 [23] Cleanlab. <https://cleanlab.ai/>.
- 396 [24] Snorkel. <https://snorkel.ai/>.
- 397 [25] Zhaowei Zhu, Jialu Wang, Hao Cheng, and Yang Liu. Unmasking and improving data credibility:
398 A study with datasets for training harmless language models. *arXiv preprint arXiv:2311.11202*,
399 2023.
- 400 [26] Zhaowei Zhu, Yiwen Song, and Yang Liu. Clusterability as an alternative to anchor points
401 when learning with noisy labels. In *International Conference on Machine Learning*, pages
402 12912–12923. PMLR, 2021.
- 403 [27] Zhaowei Zhu, Zihao Dong, and Yang Liu. Detecting corrupted labels without training a model
404 to predict. In *International conference on machine learning*, pages 27412–27427. PMLR, 2022.
- 405 [28] Gagan Bansal, Besmira Nushi, Ece Kamar, Daniel S Weld, Walter S Lasecki, and Eric Horvitz.
406 Updates in human-ai teams: Understanding and addressing the performance/compatibility
407 tradeoff. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages
408 2429–2437, 2019.
- 409 [29] Minghao Liu, Jiaheng Wei, Yang Liu, and James Davis. Do humans and machines have the
410 same eyes? human-machine perceptual differences on image classification. *arXiv preprint*
411 *arXiv:2304.08733*, 2023.
- 412 [30] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig
413 Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection
414 dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
415 pages 8769–8778, 2018.
- 416 [31] Hwanjun Song, Minseok Kim, and Jae-Gil Lee. SELFIE: Refurbishing unclean samples for
417 robust deep learning. In *ICML*, 2019.
- 418 [32] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative
419 components with random forests. In *European Conference on Computer Vision*, 2014.

- 420 [33] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu.
421 Making deep neural networks robust to label noise: A loss correction approach. In *Proceedings*
422 of the IEEE conference on computer vision and pattern recognition, pages 1944–1952, 2017.
- 423 [34] Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in
424 dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021.
- 425 [35] Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training
426 data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*,
427 33:19920–19930, 2020.
- 428 [36] Jiangfan Han, Ping Luo, and Xiaogang Wang. Deep self-learning from noisy labels. In
429 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5138–5147,
430 2019.
- 431 [37] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-
432 learning regularization prevents memorization of noisy labels. *Advances in neural information*
433 *processing systems*, 33:20331–20342, 2020.
- 434 [38] Xiaobo Xia, Tongliang Liu, Bo Han, Chen Gong, Nannan Wang, Zongyuan Ge, and Yi Chang.
435 Robust early-learning: Hindering the memorization of noisy labels. In *International conference*
436 *on learning representations*, 2020.
- 437 [39] Carla E Brodley and Mark A Friedl. Identifying mislabeled training data. *Journal of artificial*
438 *intelligence research*, 11:131–167, 1999.
- 439 [40] Jorma Laurikkala. Improving identification of difficult small classes by balancing class dis-
440 tribution. In *Artificial Intelligence in Medicine: 8th Conference on Artificial Intelligence in*
441 *Medicine in Europe, AIME 2001 Cascais, Portugal, July 1–4, 2001, Proceedings 8*, pages 63–66.
442 Springer, 2001.
- 443 [41] Hao Cheng, Zhaowei Zhu, Xingyu Li, Yifei Gong, Xing Sun, and Yang Liu. Learning with
444 instance-dependent label noise: A sample sieve approach. *arXiv preprint arXiv:2010.02347*,
445 2020.
- 446 [42] Nicolas Papernot and Patrick McDaniel. Deep k-nearest neighbors: Towards confident, inter-
447 pretable and robust deep learning. *arXiv preprint arXiv:1803.04765*, 2018.
- 448 [43] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
449 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
450 pages 770–778, 2016.
- 451 [44] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples
452 for robust deep learning. In *International conference on machine learning*, pages 4334–4343.
453 PMLR, 2018.
- 454 [45] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross
455 entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF international*
456 *conference on computer vision*, pages 322–330, 2019.
- 457 [46] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks
458 with noisy labels. *Advances in neural information processing systems*, 31, 2018.
- 459 [47] Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing
460 noise rates. In *International conference on machine learning*, pages 6226–6236. PMLR, 2020.
- 461 [48] Jiaheng Wei and Yang Liu. When optimizing f -divergence is robust with label noise. *arXiv*
462 *preprint arXiv:2011.03687*, 2020.

- 463 [49] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond
464 empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 465 [50] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help?
466 *Advances in neural information processing systems*, 32, 2019.
- 467 [51] Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. Does label smoothing
468 mitigate label noise? In *International Conference on Machine Learning*, pages 6448–6458.
469 PMLR, 2020.
- 470 [52] Jiaheng Wei, Hangyu Liu, Tongliang Liu, Gang Niu, Masashi Sugiyama, and Yang Liu. To
471 smooth or not? when label smoothing meets noisy labels. In *International Conference on
472 Machine Learning*, pages 23589–23614. PMLR, 2022.
- 473 [53] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
474 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.
475 *Advances in neural information processing systems*, 31, 2018.
- 476 [54] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning
477 data-driven curriculum for very deep neural networks on corrupted labels. In *International
478 conference on machine learning*, pages 2304–2313. PMLR, 2018.
- 479 [55] Junnan Li, Richard Socher, and Steven CH Hoi. Dividemix: Learning with noisy labels as
480 semi-supervised learning. *arXiv preprint arXiv:2002.07394*, 2020.
- 481 [56] Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A
482 joint training method with co-regularization. In *Proceedings of the IEEE/CVF conference on
483 computer vision and pattern recognition*, pages 13726–13735, 2020.
- 484 [57] Hongxin Wei, Huiping Zhuang, Renchunzi Xie, Lei Feng, Gang Niu, Bo An, and Yixuan Li.
485 Mitigating memorization of noisy labels by clipping the model prediction. In *International
486 Conference on Machine Learning*, pages 36868–36886. PMLR, 2023.
- 487 [58] Yipeng Chen, Ke Xu, Peng Zhou, Xiaojuan Ban, and Di He. Improved cross entropy loss for
488 noisy labels in vision leaf disease classification. *IET Image Processing*, 16(6):1511–1519, 2022.
- 489 [59] Taeho Jo and Nathalie Japkowicz. Class imbalances versus small disjuncts. *ACM Sigkdd
490 Explorations Newsletter*, 6(1):40–49, 2004.
- 491 [60] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote:
492 synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–
493 357, 2002.
- 494 [61] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. Borderline-smote: a new over-sampling
495 method in imbalanced data sets learning. In *International conference on intelligent computing*,
496 pages 878–887. Springer, 2005.
- 497 [62] Chumphol Bunkhumpornpat, Krung Sinapiromsaran, and Chidchanok Lursinsap. Safe-level-
498 smote: Safe-level-synthetic minority over-sampling technique for handling the class imbalanced
499 problem. In *Advances in Knowledge Discovery and Data Mining: 13th Pacific-Asia Conference,
500 PAKDD 2009 Bangkok, Thailand, April 27-30, 2009 Proceedings 13*, pages 475–482. Springer,
501 2009.
- 502 [63] Haibo He, Yang Bai, Edwardo A Garcia, and Shutao Li. Adasyn: Adaptive synthetic sampling
503 approach for imbalanced learning. In *2008 IEEE international joint conference on neural
504 networks (IEEE world congress on computational intelligence)*, pages 1322–1328. Ieee, 2008.
- 505 [64] Inderjeet Mani and I Zhang. knn approach to unbalanced data distributions: a case study
506 involving information extraction. In *Proceedings of workshop on learning from imbalanced
507 datasets*, volume 126, pages 1–7. ICML, 2003.

- 508 [65] Miroslav Kubat, Stan Matwin, et al. Addressing the curse of imbalanced training sets: one-sided
509 selection. In *Icml*, volume 97, page 179. Citeseer, 1997.
- 510 [66] I TOMEK. Two modifications of cnn. *IEEE Trans. Systems, Man and Cybernetics*, 6:769–772,
511 1976.
- 512 [67] Charles Elkan. The foundations of cost-sensitive learning. In *International joint conference on*
513 *artificial intelligence*, volume 17, pages 973–978. Lawrence Erlbaum Associates Ltd, 2001.
- 514 [68] Matjaz Kukar, Igor Kononenko, et al. Cost-sensitive learning with neural networks. In *ECAI*,
515 volume 15, pages 88–94. Citeseer, 1998.
- 516 [69] Zhi-Hua Zhou and Xu-Ying Liu. Training cost-sensitive neural networks with methods ad-
517 dressing the class imbalance problem. *IEEE Transactions on knowledge and data engineering*,
518 18(1):63–77, 2005.
- 519 [70] Steve Lawrence, Ian Burns, Andrew Back, Ah Chung Tsoi, and C Lee Giles. Neural network
520 classification and prior class probabilities. In *Neural networks: tricks of the trade*, pages
521 299–313. Springer, 2002.
- 522 [71] Michael D Richard and Richard P Lippmann. Neural network classifiers estimate bayesian a
523 posteriori probabilities. *Neural computation*, 3(4):461–483, 1991.
- 524 [72] Jiaheng Wei, Harikrishna Narasimhan, Ehsan Amid, Wen-Sheng Chu, Yang Liu, and Ab-
525 hishek Kumar. Distributionally robust post-hoc classifiers under prior shifts. In *The Eleventh*
526 *International Conference on Learning Representations*, 2023.
- 527 [73] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
528 object detection. In *Proceedings of the IEEE international conference on computer vision*,
529 pages 2980–2988, 2017.
- 530 [74] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced
531 datasets with label-distribution-aware margin loss. *Advances in neural information processing*
532 *systems*, 32, 2019.
- 533 [75] Jiawei Ren, Cunjun Yu, Xiao Ma, Haiyu Zhao, Shuai Yi, et al. Balanced meta-softmax for long-
534 tailed visual recognition. *Advances in neural information processing systems*, 33:4175–4186,
535 2020.
- 536 [76] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit,
537 and Sanjiv Kumar. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*,
538 2020.

539 **Checklist**

- 540 1. For all authors...
- 541 (a) Do the main claims made in the abstract and introduction accurately reflect the pa-
542 per's contributions and scope? [Yes] The abstract and introduction reflect our main
543 contribution clearly.
- 544 (b) Did you describe the limitations of your work? [Yes] We have included the limitation
545 section in the beginning of the Appendix.
- 546 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Beginning
547 of Appendix.
- 548 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
549 them? [Yes] We have read the code of ethics carefully and we are sure that we follow
550 it.
- 551 2. If you are including theoretical results...
- 552 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 553 (b) Did you include complete proofs of all theoretical results? [N/A]
- 554 3. If you ran experiments (e.g. for benchmarks)...
- 555 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
556 imental results (either in the supplemental material or as a URL)? [Yes] We have
557 included in the Supplementary.
- 558 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
559 were chosen)? [Yes] We have included in Appendix Section C.3, C.4 and C.5.
- 560 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
561 ments multiple times)? [Yes] Please refer to Table 4 and Table5
- 562 (d) Did you include the total amount of compute and the type of resources used (e.g., type
563 of GPUs, internal cluster, or cloud provider)? [Yes] We have included in the Appendix
564 SectionC.3, C.4 and C.5.
- 565 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 566 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 567 (b) Did you mention the license of the assets? [N/A]
- 568 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
569 We have included in the Supplementary.
- 570 (d) Did you discuss whether and how consent was obtained from people whose data you're
571 using/curating? [N/A]
- 572 (e) Did you discuss whether the data you are using/curating contains personally identifiable
573 information or offensive content? [N/A]
- 574 5. If you used crowdsourcing or conducted research with human subjects...
- 575 (a) Did you include the full text of instructions given to participants and screenshots, if
576 applicable? [Yes] We have included in Appendix B.3, C.1.
- 577 (b) Did you describe any potential participant risks, with links to Institutional Review
578 Board (IRB) approvals, if applicable? [N/A]
- 579 (c) Did you include the estimated hourly wage paid to participants and the total amount
580 spent on participant compensation? [Yes] We have included in Appendix B.3, C.1.

581 **Appendix**

582 The appendix is organized as follows:

- 583 • Appendix A includes additional detailed algorithms in the Automatic-Dataset-Construction
584 pipeline.
- 585 • Appendix B contains dataset statistics and more exploratory data analysis of Clothing ADC.
- 586 • Appendix C includes experiment details of our benchmark on label noise detection, label noise
587 learning, and class-imbalanced learning.

588 **Broader impacts**

589 Our paper introduces significant advancements in dataset construction methodologies, particularly
590 through the development of the Automatic Dataset Construction (ADC) pipeline:

- 591 • **Reduction in Human Workload:** ADC automates the process of dataset creation, significantly
592 reducing the need for manual annotation and thereby decreasing both the time and costs associated
593 with data curation.
- 594 • **Enhanced Data Quality for Research Communities:** ADC provides high-quality, tailored
595 datasets with minimal human intervention. This provides researchers with datasets in the fields of
596 label noise detection, label noise learning, and class-imbalanced learning, for exploration as well
597 as fair comparisons.
- 598 • **Support for Customized LLM Training:** The ability to rapidly generate and refine datasets
599 tailored for specific tasks enhances the training of customized Large Language Models (LLMs),
600 increasing their effectiveness and applicability in specialized applications.

601 Furthermore, the complementary software developed alongside ADC enhances these impacts:

- 602 • **Data Curation and Quality Control:** The software aids in curating and cleaning the collected
603 data, ensuring that the datasets are of high quality that could compromise model training.
- 604 • **Robust Learning Capabilities:** It incorporates methods for robust learning with collected data,
605 addressing challenges such as label noise and class imbalances. This enhances the reliability and
606 accuracy of models trained on ADC-constructed datasets.

607 Together, ADC and its accompanying software significantly advance the capabilities of machine
608 learning researchers and developers by providing efficient tools for high-quality customized data
609 collection, and robust training.

610 **Limitations**

611 While ensuring the legal and ethical use of datasets, including compliance with copyright laws and
612 privacy concerns, is critical, our initial focus is on legally regulated and license-friendly data sources
613 available through platforms like Google or Bing. Addressing these ethical considerations is beyond
614 the current scope but remains an essential aspect of dataset usage.

615 Besides, similar to Traditional-Dataset-Construction (TDC), Automatic-Dataset-Construction (ADC)
616 is also unable to guarantee fully accurate annotations.

617 **A Detailed algorithms in the generation of Automatic-Dataset-Construction**

618 **A.1 The algorithm of image data collection in ADC**

Algorithm 2 Image Data Collection in ADC

```
1: procedure IMAGEDATACOLLECTION
2:   Part A: Get attributes from dataset design
3:   attributes  $\leftarrow$  Step 1 Dataset Design
4:   categories  $\leftarrow$  ["sweater", "shirt", "pants", ...]            $\triangleright$  List of categories
5:   target_category  $\leftarrow$  "sweater"                                 $\triangleright$  Target category (e.g. "sweater")
6:   attributes  $\leftarrow$  attributes[target_category]                 $\triangleright$  Get attributes for target category
7:   colors, patterns, materials  $\leftarrow$  attributes["color"],          |
8:                           attributes["pattern"],                   |
9:                           attributes["material"]                  |
10:  Part B: Create search queries
11:  search_queries  $\leftarrow$  {  $c + p + m + \text{target\_category}$  |           |
12:     $c \in \text{colors}$ ,                                         |
13:     $p \in \text{patterns}$ ,                                     |
14:     $m \in \text{materials}$  }            $\triangleright$  (e.g. "beige fisherman cotton sweater")
15:  Part C: Launch distributed image search
16:  image_data  $\leftarrow$  distributed_search(search_queries,
17:                      api = Google_Images | Bing_Images,
18:                      n_process = 30)
19: end procedure
```

619 **A.2 The algorithm of learning-centric curation method in ADC**

Algorithm 3 Learning-centric curation (early-learning memorization behavior)

```
1: procedure EARLYSTOPCE(noisyDataset, percentage=25%)
2:   Part A: Train classifier over the dataset and apply early stopping
3:    $\mathcal{D} \leftarrow$  Load training data                                $\triangleright$  (images and labels)
4:   model  $\leftarrow$  Initialize neural network model                $\triangleright$  (e.g. ResNet)
5:   loss_fn  $\leftarrow$  Define loss function                          $\triangleright$  (e.g. cross-entropy)
6:   optimizer  $\leftarrow$  Choose optimizer                           $\triangleright$  (e.g. SGD, Adam)
7:   for epoch = 1 to  $E \in \{1, 2\}$  do
8:     model  $\leftarrow$  Trainer( $\mathcal{D}, loss\_fn, optimizer$ )
9:   end for
10:  Part B: Record predictions and confidence levels
11:  for batch in  $\mathcal{D}$  do
12:    images  $\leftarrow$  Get batch of images
13:    outputs  $\leftarrow$  Forward pass:  $model(images)$ 
14:    confidence  $\leftarrow$  Get confidence levels: softmax(outputs)
15:  end for
16:  Part C: Remove samples with lowest  $x\%$  confidence level
17:  threshold  $\leftarrow$  Calculate threshold: percentile(confidence,  $100 - x$ )
18:   $\mathcal{D} \leftarrow$  Filter out samples with confidence below threshold
19:  Return  $\mathcal{D}$ 
20: end procedure
```

620 **B Dataset statistics in Clothing-ADC**

621 **B.1 Collected Clothing ADC dataset**

622 Our collected Clothing-ADC dataset can be found here: Google Drive.

623 B.2 Dataset statistics in Clothing-ADC

Our automated dataset creation pipeline is capable of generating numerous designs per attribute, as shown in Table 6. This table provides a detailed list of designs generated by our pipeline, from which we selected a subset to include in our dataset.

Color			Material			Pattern					
Animal print	Gold	Pastel	Acrylic	Lace	Tulle	Abstract	Camouflage	Fishnet	Leather	Printed	Thongs
Beige	Gray	Peach	Alpaca	Leather	Tweed	Abstract Floral	Chalk stripe	Floral	Logo	Quilted	Tie-Dye
Black	Green	Pink	Angora	Lightweight	Twill	Animal Print	Check	Floral print	Low rise	Reversible	Tie-dye
Blue	Grey	Plum	Bamboo	Linen	Velvet	Animal print	Checkered	Fringe	Mesh	Ribbed	Toile
Blush Pink	Heather	Purple	Breathable	Mesha	Viscose	Aran	Chevron	G-strings	Military	Ripples	Trench
Bright Red	Ivory	Red	Cashmere	Microfiber	Water-resistant	Argyle	Color block	Galaxy	Mock turtleneck	Satin	Tribal
Brown	Khaki	Rich Burgundy	Chambray	Modal	Windproof	Aztec	Colorblock	Garter Stitch	Mosaic	Scales	Tuck stitch
Burgundy	Lavender	Royal Blue	Chiffon	Mohair	Wool	Basket check	Cotton	Garter stitch	Moss stitch	Seamless	Tweed
Burnt Orange	Light Grey	Rust	Corduroy	Neoprene	Acrylic	Basket rib	Cropped	Geometric	Moto	Seed stitch	Twill
Champagne	Maroon	Rustic Orange	Combi	Nylon	bamboo	Basket weave	Damask	Gingham	Nautilus	Shadow stripe	Vintage-inspired
Coat	Metallic	Sage	Crochet	Organza	cotton	Basketweave	Diagonal	Glen check	Nchiru	Sharkskin	Waterproof
Clouded Grey	Mustard	Silver	Denim	PVC	hemp	Batik	Diagonal grid	Giant	Nothic	Sheila	Windowpane
Cream	Yellow	Soft Pink	Down	Polyester	linen	Bikini	Diamond	Graphic	Ombre	Shibori	
Cream White	Navy	Striped	Embroidered	Rayon	lyra	Birdseye	Ditsy	Grid	Oversized	Slip Stitch	
Dark Plum	Navy Blue	Tan	Flannel	Reflective	modal	Blazer	Dogtooth	Herringbone	Oxford	Slip stitch	
Deep Blue	Neon	Teal	Fleece	Ripstop	nylon	Bomber	Embossed	High waisted	Paisley	Solid	
Deep Purple	Nude	Turquoise	Fringe	Satin	polyester	Boxer briefs	Embroidered	Honeycomb	Pearcoat	Striped	
Earthy Beige	Olive	Vibrant Turquoise	Fur	Silk	rayon	Briefs	Emoji	Houndstooth	Pin Dot	Stripes	
Forest Green	Olive Green	Warm Brown	Gore Tex	Softshell	silk	Brioché	Entrelac	Ikat	Pinstripe	Studded	
Fuchsia	Pale Yellow	Yellow	Gore-Tex	Spandex	spandex	Broken rib	Eyelet	Intarsia	Plaid	Suede	
		lilac	Hemp	Sued	tencel	Broken stripe	Fair Isle	Jacquard	Polka Dot	Tartan	
			Insulated	Synthetic	viscose	Cable	Fibonacci	Knit and Purl	Polka dot	Teddy	
			Jersey	Synthetic Blend	wool	Cable knit	Fisherman	Lace	Prince of Wales	Textured	

Table 6: The union of attributes across all clothing types in Clothing-ADC dataset.

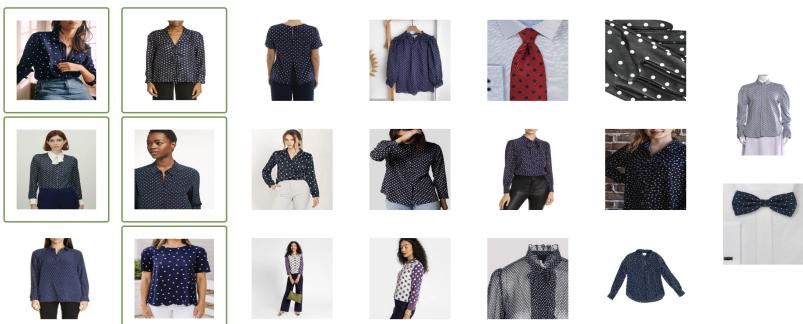
627 B.3 Test set human vote collection

Our automated dataset collection pipeline enabled us to create a large, noisy labeled dataset. We asked annotators to select the best-fitting options from a range of samples, as shown in Figure 5, with each task including at least 4 samples and workers completing 10 tasks per HIT at a cost of \$0.15 per task, totaling \$150 estimated wage of \$2.5-3 per hour, and after further cleaning the label noise, we ended up with 20,000 samples in our test set. To participate, workers had to meet specific requirements, including being Master workers, having a HIT Approval Rate above 85%, and having more than 500 approved HITs, with the distribution of worker behavior shown in Figure 6.

[Task 6 / 10] Please find 4 or more images of T-shirt with:

Color: Navy, Material: silk, Pattern: Polka dot

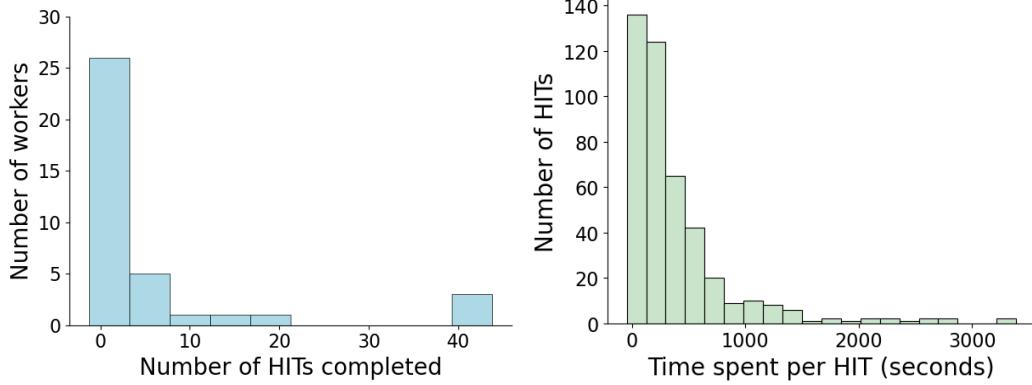
If you are not familiar with the any of these key words, feel free to search it on Google. Try your best to find the most relevant images. I am genuinely interested in your opinion and generous in accepting your answers.



Previous task

Next task

Figure 5: Collection of Clothing-ADC test set: A filtering task to the worker instead of annotation from scratch.



(a) Distribution of the HITs completed per worker (b) Distribution of work time in seconds per HIT.

Figure 6: The behaviors of workers in the creation of test set.

635 C Experiment details

636 C.1 Distribution of Human Votes for Label Noise Evaluation

637 On the annotation page, we presented the image and its original label to the worker and asked if
 638 they believed the label was correct (Figure 7). They input their evaluation by clicking one of three
 639 buttons. Note that we encouraged workers to categorize acceptable samples as "unsure". The resulting
 640 distribution is shown in Table 7. Using a simple majority vote aggregation, we found that the noise
 641 rate in our dataset is 22.15%. However, if a higher level of certainty is required for clean labels,
 642 we can apply a more stringent aggregation method, considering more samples as mislabeled. In
 643 the extreme case where any doubts from any of the three annotators can disqualify a sample, our
 644 automatically collected dataset still retains 61.25% of its samples.

645 For the label noise evaluation task, we utilized a subset of 20,000 samples from the Clothing-ADC
 646 dataset, collecting three votes from unique workers for each sample. Each Human Intelligence
 647 Task (HIT) included 20 samples and cost \$0.05. To participate, workers had to meet the following
 648 requirements: (1) be Master workers, (2) have a HIT Approval Rate above 85%, and (3) have more
 649 than 500 approved HITs. The total cost for this task was \$150, estimated wage of \$2.5-3 per hour.

650 We show the distribution of worker behavior during the noise evaluation task in Figure 8. Figure 8(a)
 651 shows the distribution of the amount of HIT completed per worker while neglecting ids with 1-2
 652 submissions. There is a total of 49 unique workers. Figure 8(b) shows the distribution of time spent
 653 per HIT.

Table 7: **Distribution of Human Votes for Label Noise Evaluation:** We employed human annotators to evaluate a subset of 20,000 samples from our collected dataset, with each sample receiving three votes from distinct annotators.

Human Votes	Percentage
Yes, Yes, Yes	61.25%
Yes, Yes, Unsure	6.10%
Yes, Yes, No	10.50%
Else	22.15%

[Task 2 / 2] The following image has been labeled as **Dress**

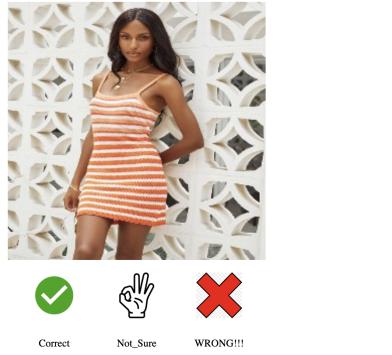
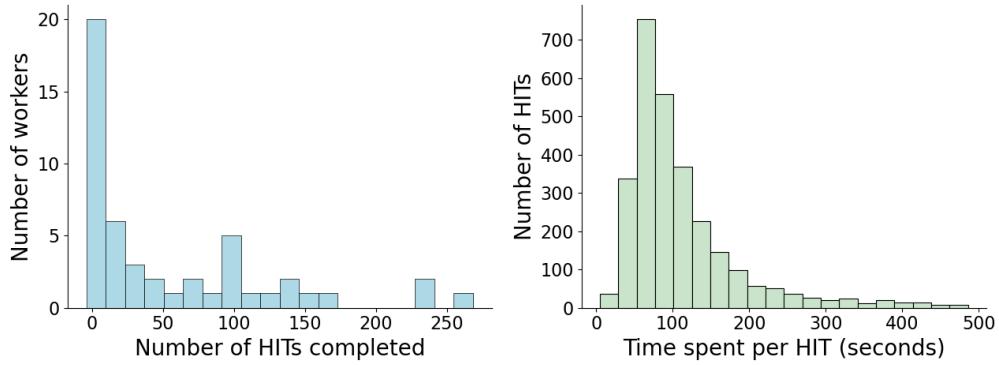


Figure 7: Label noise evaluation worker page



(a) Distribution of the HITs completed per worker (b) Distribution of work time in seconds per HIT.

Figure 8: The behaviors of workers in the collection of label noise evaluation.

654 C.2 Noisy learning and class imbalance learning benchmark implementation details

655 Our code refers to zip file in supplementary material.

```
656 1 train_set = Clothing1mPP(root, image_size, split="train")
657 2 tiny_set_ids = train_set.get_tiny_ids(seed=0)
658 3 tiny_train_set = Subset(train_set, tiny_set_ids) # Get the tiny
659     version of the dataset
660 4 val_set = Clothing1mPP(
661 5     root, image_size, split="val", pre_load=train_set.data_package
662 6 )
663 7 test_set = Clothing1mPP(
664 8     root, image_size, split="test", pre_load=train_set.data_package
665 9 )
666 10
667 11 train_loader = DataLoader(
668 12     train_set, batch_size=batch_size, shuffle=True, num_workers=
669         num_workers
670 13 )
671 14 tiny_train_loader = DataLoader(
672 15     tiny_train_set, batch_size=batch_size, shuffle=True, num_workers=
673         num_workers
```

```

67416 )
67517 val_loader = DataLoader(
67618     val_set, batch_size=batch_size, shuffle=False, num_workers=
67719     num_workers
67820 )
67921 test_loader = DataLoader(
68022     test_set, batch_size=batch_size, shuffle=False, num_workers=
68123     num_workers
68224 )

```

Listing 1: How to load data. Line 1 loads the full set of our dataset. Line 2 and Line 3 load the tiny version of our dataset. Line 4 creates the validation set. Line 5 creates the testing set. Line 11 to Line 20 create the data loader.

```

6831 python examples/main.py --config configs/Clothing1MPP/default.yaml #
6842     Run Cross Entropy
6853 python examples/main_peer.py --config configs/Clothing1MPP/default.
6864     yaml # Run Peer Loss
6875 python examples/main_jocor.py --config configs/Clothing1MPP/
6886     default_jocor.yaml # Run Jocor
6897 python examples/main_coteaching.py --config configs/Clothing1MPP/
6908     default_coteaching.yaml # Run Co-teaching
6919 python examples/main_drops.py --config configs/Clothing1MPP/
69210    default_drops.yaml # Run drops

```

Listing 2: The example of the command we use to run the algorithm in one line

```

6931 inherit_from: configs/default.yaml
6942 data: &data_default
6953   root: '/root/cloth1m_data_v3'
6964   image_size: 256
6975   dataset_name: "clothing1mpp"
6986   imbalance_factor: 1 # 1 means no imbalance
6997   tiny: False
7008
7019 train: &train
70210  num_workers: 8
70311  loss_type: 'ce'
70412  loop_type: 'default' # 'default','peer','drops'
70513  epochs: 20
70614  global_iteration: 999999999
70715  batch_size: 64
70816  # scheduler_T_max: 40
70917  scheduler_type: 'step'
71018  scheduler_gamma: 0.8
71119  scheduler_step_size: 2
71220  print_every: 100
71321  learning_rate: 0.01
71422
71523 general:
71624   save_root: './results/'
71725   whip_existing_files: True # Whip exisitng files
71826   logger:
71927     project_name: 'Clothing1MPP',
72028     frequency: 200
72129
72230 model: &model_default
72331   name: "resnet50"
72432   pretrained_model: 'IMAGENET1K_V1'
72533   cifar: False
72634
72735 test: &test_defaults

```

728:6 <<: *train

Listing 3: The example of YAML config file

729 **C.3 Label noise detection**

730 We run four baselines for label noise detection, including CORES [41], confident learning [34], deep
 731 k -NN [42] and Simi-Feat [27]. All the experiment is run for one time following [41, 27].

732 The experiment platform we run is a 128-core AMD EPYC 7742 Processor CPU and the memory is
 733 128GB. The GPU we use is a single NVIDIA A100 (80GB) GPU. For the dataset, we used human
 734 annotators to evaluate whether the sample has clean or noisy label as mentioned in Appendix C.1.
 735 We aggressively eliminates human uncertainty factors and only consider the case with unanimous
 736 agreement as a clean sample, and everything else as noisy samples. The backbone model we use
 737 is ResNet-50 [43]. For all the baselines, the parameters we use are the same as the original paper
 738 except the data loader. We skip the label corruption and use the default value from the original
 739 repository. For CORES, the cores loss whose value is smaller than 0 is regarded as the noisy sample.
 740 For confidence learning, we use the repository¹ from the clean lab and the default hyper-parameter.
 741 For deep k -NN, the k we set is 100. For SimiFeat, we set k as 10 and the feature extractor is CLIP.

742 **C.4 Label noise learning**

743 The platform we use is the same as label noise detection. The backbone model we use is ResNet-
 744 50 [43]. For the full dataset, we run the experiment for 1 time. For the tiny dataset, we run the
 745 experiments for 3 times. The tiny dataset is sampled from the full set whose size is 50. The base
 746 learning rate we use is 0.01. The base number of epochs is 20. The hyper-parameters for each
 747 baseline method are as follows. For **backward and forward correction**, we train the model using
 748 cross-entropy (CE) loss for the first 10 epochs. We estimate the transition matrix every epoch from
 749 the 10th to the 20th epoch. For the **positive and negative label smoothing**, the smoothed labels are
 750 used at the 10th epoch. The smooth rates of the positive and negative are 0.6 and -0.2. Similarly, for
 751 **peer loss**, we train the model using CE loss for the first 10 epochs. Then, we apply peer loss for the
 752 rest 10 epochs and the learning rate we use for these 10 epochs is 1e-6. The hyper-parameters for
 753 f -div is the same as those of peer loss. For **divide-mix**, we use the default hyper-parameters in the
 754 original paper. For **Jocor**, the hyper-parameters we use is as follows. The learning rate is 0.0001. λ
 755 is 0.3. The epoch when the decay starts is 5. The hyper-parameters of **co-teaching** is similar to Jocor.
 756 For logitclip, τ is 1.5. For **taylorCE**, the hyper-parameter is the same as the original paper.

757 **C.5 Class-imbalanced learning**

758 The platform we use is the same as label noise detection. The backbone model we use is ResNet-50
 759 [43]. For different imbalance ratio ($\rho = 10, 50, 100$). The class distribution is shown in Table 8. For
 760 all the methods, the base learning rate is 0.0001 and the batch size is 448. The dataset we use is
 761 not full dataset because we want to disentangle the noisy label and class imbalance learning. We
 762 use Docta and a pre-trained model trained with cross-entropy to filter the data whose prediction
 763 confidence is low. Due to the memorization effect, we fine-tune the model for 2 epochs to filter the
 764 data. We remove 45.15% data in total where Docta removes 26.36% while CE removes 25.00% with
 765 a overlap of 6.20%. Thus, the datset we use for class-imbalance learning is 54.85% of the full dataset.

imbalance ratio (ρ)	Class Distribution	Total Number
10	[39297, 31875, 25854, 20971, 17010, 13797, 11191, 9078, 7363, 5972, 4844, 3929]	191181
20	[39297, 27536, 19295, 13520, 9474, 6638, 4652, 3259, 2284, 1600, 1121, 785]	129461
100	[39297, 25854, 17010, 11191, 7363, 4844, 3187, 2097, 1379, 907, 597, 392]	114118

Table 8: The class distribution for different imbalance ratio

¹<https://github.com/cleanlab/cleanlab>