

# DISSECTING SAMPLE HARDNESS: A FINE-GRAINED ANALYSIS OF HARDNESS CHARACTERIZATION METHODS FOR DATA-CENTRIC AI

**Nabeel Seedat**  
University of Cambridge  
ns741@cam.ac.uk

**Fergus Imrie**  
UCLA  
imrie@ucla.edu

**Mihaela van der Schaar**  
University of Cambridge  
mv472@cam.ac.uk

## ABSTRACT

Characterizing samples that are difficult to learn from is crucial to developing highly performant ML models. This has led to numerous Hardness Characterization Methods (HCMs) that aim to identify “hard” samples. However, there is a lack of consensus regarding the definition and evaluation of “hardness”. Unfortunately, current HCMs have only been evaluated on specific types of hardness and often only qualitatively or with respect to downstream performance, overlooking the fundamental quantitative identification task. We address this gap by presenting a fine-grained taxonomy of hardness types. Additionally, we propose the Hardness Characterization Analysis Toolkit (H-CAT), which supports comprehensive and quantitative benchmarking of HCMs across the hardness taxonomy and can easily be extended to new HCMs, hardness types, and datasets. We use H-CAT to evaluate 13 different HCMs across 8 hardness types. This comprehensive evaluation encompassing over *14K setups* uncovers strengths and weaknesses of different HCMs, leading to practical tips to guide HCM selection and future development. Our findings highlight the need for more comprehensive HCM evaluation, while we hope our hardness taxonomy and toolkit will advance the principled evaluation and uptake of data-centric AI methods.

## 1 INTRODUCTION

**Data quality, an important ML problem.** Data quality is crucial to the performance and robustness of machine learning (ML) models (Jain et al., 2020; Gupta et al., 2021; Renggli et al., 2021; Sambasivan et al., 2021; Li et al., 2021). Unfortunately, challenges arise in real-world data that make samples “hard” for ML models to learn from effectively, including but not limited to mislabeling, outliers, and insufficient coverage (Chen et al., 2021; Li et al., 2021). These “hard” samples or data points can significantly hamper the performance of ML models, creating a barrier to ML adoption in practical applications (Bedi et al., 2019; West, 2020). For instance, a model trained on mislabeled samples can lead to inaccurate predictions (Krishnan et al., 2016; Gupta et al., 2021). Outliers can bias the model to learn suboptimal decision boundaries (Liu et al., 2022; Eduardo et al., 2022; Krishnan et al., 2016), harming model performance. Long tails of samples can result in poor model performance for these cases (Feldman, 2020; Hooker et al., 2019; Hooker, 2021; Agarwal et al., 2022). Consequently, “hard” samples can pose serious challenges for training and the performance of ML models, making it crucial to identify these samples. This is especially important in where manually identifying “hard” samples is expensive, impractical, or time-consuming given the scale.

**Characterizing hardness, a growing area.** Recent interest in data-centric AI, which aims to ensure and improve data quality, has led to the development of systematic methods to characterize the data used to train ML models (Liang et al., 2022; Seedat et al., 2023b; 2022b). Data characterization typically assigns scores to each sample based on its learnability and utility for an ML task, thereby facilitating the identification of “hard” samples. We collectively refer to methods that perform the characterization as Hardness Characterization Methods (HCMs).

After samples are characterized, how and for what purpose they are used can differ. For example: (1) curating datasets via sample selection to improve model performance (Maini et al., 2022; Swayamdipta et al., 2020; Seedat et al., 2022a; Northcutt et al., 2021a; Pleiss et al., 2020; Agarwal

et al., 2022; Seedat et al., 2023a), (2) sculpting the dataset to reduce computational requirements while maintaining performance (Toneva et al., 2019; Paul et al., 2021; Sorscher et al., 2022; Mindermann et al., 2022), (3) guiding acquisition of additional samples (Zhang et al., 2022), or (4) understanding learning behavior from a theoretical perspective (Baldock et al., 2021; Shwartz et al., 2022; Jiang et al., 2021).

**Challenges in definition and evaluation.** A fundamental and *overlooked* aspect is that while different HCMs tackle the issue of “hardness”, it remains a vague and ill-defined term in the literature. The lack of a clear definition of hardness types has led HCMs, seemingly tackling the same problem, to unintentionally target and assess different aspects of hardness (Table 1). The lack of clarity is further exacerbated by: (1) *qualitative evaluation*: a significant focus on post hoc *qualitative* assessment and downstream improvement, instead of the fundamental hardness identification task, and (2) *narrow and unrepresentative scope*: even when quantitative evaluation has been performed, it has typically focused on a single hardness type, neglecting different manifestations of hardness. The lack of comprehensive and quantitative evaluation means we do not know how different HCMs perform on different hardness types and whether they indeed identify the correct samples of interest.

*Can we define sample hardness manifestations and then comprehensively and systematically evaluate the capabilities of different HCMs to correctly detect the hard samples?*

**Unified taxonomy and benchmarking framework.** To answer this question, we begin by defining a taxonomy of hardness types across three broad categories: (a) Mislabeling, (b) OoD/Outlier, (c) Atypical. We then introduce the **Hardness-Characterization Analysis Toolkit (H-CAT)**, which, to the best of our knowledge, is the first unified data characterization benchmarking framework focused on hardness. Using H-CAT, we comprehensively and quantitatively benchmark 13 state-of-the-art HCMs across various hardness types. In doing so, we address recent calls for more rigorous benchmarking (Guyon, 2022) and understanding of existing ML methods (Lipton & Steinhardt, 2019; Snoek et al., 2018). We make the following contributions:

**Contributions:** ① *Hardness taxonomy*: we formalize a systematic taxonomy of sample-level hardness types, addressing the current literature’s ad hoc and narrow scope. By defining the different dimensions of hardness, our taxonomy paves the way for a more rigorous evaluation of HCMs. ② *Benchmarking framework*: we propose H-CAT, which is both (i) a *benchmarking standard* to evaluate the strengths of different HCMs across the hardness taxonomy and (ii) a *unified software tool* integrating 13 different HCMs. With extensibility in mind, H-CAT can easily incorporate new HCMs, hardness types, and datasets, thus enhancing its utility for both researchers and practitioners. ③ *Systematic & Quantitative HCM evaluation*: we use H-CAT to comprehensively benchmark and evaluate 13 different HCMs across 8 different hardness types, comprising *over 14K* experimental setups. ④ *Insights*: our benchmark provides novel insights into the capabilities of different HCMs when dealing with different hardness types and offers practical usage tips for researchers and practitioners. The variability in HCM performance across hardness types underscores the importance of multi-dimensional evaluation, exposing gaps and opportunities in current HCMs. We hope H-CAT will promote rigorous HCM evaluations and inspire new advances in data-centric AI.

## 2 HARDNESS CHARACTERIZATION AND TAXONOMY

We now outline the hardness characterization problem and formalize a hardness taxonomy for HCMs.

**Learning problem.** Consider the typical supervised learning setting, with  $\mathcal{X}$  and  $\mathcal{Y}$  input and output spaces, respectively. We assume a  $k$ -class classification problem, i.e.  $\mathcal{Y} = [k]$ , where  $[k] = \{1, \dots, k\}$ , with a training dataset  $\mathcal{D} = \{(x_i, y_i) \mid i \in [N]\}$  with  $N \in \mathbb{N}^+$  samples, where  $x_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y}$ . The goal is to learn a model  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $\theta \in \Theta$ .

**Hardness problem.** To understand the intricacies of data hardness, we start by providing a broad definition of the hardness problem common to all HCMs. As a starting point, let us frame the hardness problem in general terms: specifically, some samples or data points are easier for a model to learn from, whilst others may either be harder to learn from or harm the model’s performance.

Formally, this assumes that the training dataset can be decomposed as  $\mathcal{D} = \mathcal{D}_e \cup \mathcal{D}_h$ , where  $\mathcal{D}_e$  are easy samples and  $\mathcal{D}_h$  are hard samples. We denote the corresponding joint distributions as  $\mathcal{P}_{XY}$ ,  $\mathcal{P}_{XY}^e$  and  $\mathcal{P}_{XY}^h$ . Going beyond this general definition to different manifestations of hardness requires a more rigorous characterization. However, what constitutes a “hard” sample has not been rigorously

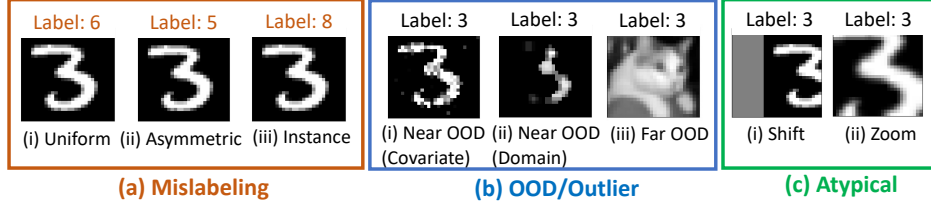


Figure 1: Examples of the hardness types included within our taxonomy and supported by H-CAT. HCMs need to be comprehensively assessed across different dimensions to quantify their ability to handle different hardness types we might expect in practice. See Sec. 2.1 for precise definitions of hardness types.

defined in the literature. To address this gap, we first formalize a taxonomy of hardness in Sec. 2.1, providing a systematic and formal definition of different types of sample-level hardness.

**Data characterization.** The goal of data characterization is to assign a scalar “hardness” score to each sample in  $\mathcal{D}$ , thereby allowing us to order samples in  $\mathcal{D}$  according to their scores. Typically, a selector function applies a threshold  $\tau$  and then assigns a hardness label  $g \in \mathcal{G}$ , where  $\mathcal{G} = \{Easy, Hard\}$ , to each sample  $(x_i, y_i)$ , i.e. assign samples in  $\mathcal{D}$  to either  $\mathcal{D}_e$  or  $\mathcal{D}_h$ . We refer to these methods as **Hardness Characterization Methods (HCM)**, having the following input-output paradigm, where methods differ based on their scoring mechanism:

- **Inputs:** (i) Dataset  $\mathcal{D} = \{(x_i, y_i)\}$  drawn from both  $\mathcal{P}_{XY}$ . (ii) Learning algorithm  $f_\theta$ .
- **Outputs:** Assign a score  $s_i$  to sample  $(x_i, y_i)$ . Apply threshold  $\tau$  to assign a hardness label  $g \in \mathcal{G}$ , where  $\mathcal{G} = \{Easy, Hard\}$ , which is then used to partition  $\mathcal{D} = \mathcal{D}_e \cup \mathcal{D}_h$ .

We group the HCMs into broad classes (see Table 1) determined based on the core metric or approach each method uses to characterize example hardness, namely (1) Learning dynamics-based: relying on metrics computed during the training process itself to characterize example hardness; (2) Distance-based: using the distance or similarity of examples in an embedding space; (3) Statistical-based: using statistical metrics computed over the data to characterize example hardness.

## 2.1 TAXONOMY OF HARDNESS

**Hardness taxonomy formalism.** The term “hardness” is broad and can manifest in various ways, as demonstrated by the different types of hardness previously examined (see Table 1). Therefore, we must first formalize a taxonomy representative of the different types of hardness we might expect in practice. Each type of hardness is characterized by a latent or unseen hardness perturbation function  $h$  that creates hard samples  $\mathcal{D}_h$  from samples in  $\mathcal{D}$ . We denote hardness perturbations on  $X$  as  $X \rightarrow X^*$  and  $Y$  as  $Y \rightarrow Y^*$ . The effect of each hardness perturbation is explained in terms of the relationship between the joint probability distributions  $P_{XY}^e(x, y)$  and  $P_{XY}^h(x, y)$ .

The taxonomy deals with three broad types (and corresponding subtypes) of hardness: (1) Mislabeling, (2) OoD/Outlier, and (3) Atypical, as illustrated in Figure 1. We anchor the different hardness manifestations with respect to relevant literature for each subtype. We define the various types next.

**Mislabeling:** Samples where the true label is replaced with an incorrect label, such that sample  $(x, y) \rightarrow (x, y^*)$ . The main distinction between easy and hard samples lies in the label space, leading to different conditional probability distributions:  $P_{Y|X}^e(y|x) \neq P_{Y|X}^h(y|x)$ . Note, the marginal probability distributions are the same  $P_X^e(x) = P_X^h(x)$ .

We consider three subtypes of mislabeling: (i) Uniform, (ii) Asymmetric, and (iii) Instance. The HCM literature primarily focuses on evaluation with a *uniform* noise model (Paul et al., 2021; Swayamdipta et al., 2020; Pleiss et al., 2020; Toneva et al., 2019; Maini et al., 2022; Jiang et al., 2021; Mindermann et al., 2022; Baldock et al., 2021), with equal probability of mislabeling across classes. However, in reality, mislabeling is often *asymmetric*, where mislabeling is label-dependent (Northcutt et al., 2021a; Sukhbaatar et al., 2015) or *instance-specific*, where certain mislabeling is more likely given the sample (Jia et al., 2022; Han et al., 2020; Hendrycks et al., 2018; Song et al., 2022). For example, mislabeling an image of a car as a truck is more likely than mislabeling a car as a dog.

Formally, the subtypes differ in their noise models to perturb the true labels, defined by the probabilities  $P(Y^* = j|Y = i)$ , where  $i$  and  $j$  are the true and perturbed labels, respectively.

- *Uniform:*  $P(Y^* = j|Y = i) = 1/(k-1) \quad \forall i, j \in [k], i \neq j$
  - *Asymmetric:*  $P(Y^* = j|Y = i) = p_{ij} \quad \forall i, j \in [k], i \neq j$
  - *Instance:*  $P(Y^* = j|Y = i, X = x) = p_{ij}(x) \quad \forall i, j \in [k], i \neq j$
- } Instance Independent  
} Instance Dependent

Table 1: HCMs, even in similar classes, are often (1) *not* quantitatively evaluated and (2) assess different hardness types. The tick and cross denote if quantitative and in brackets qualitative evaluation was performed, e.g.  $\checkmark(\checkmark)$  - *not* quantitative (qualitative). For HCM descriptions, see Appendix A.

HCM Class	HCM Name	Mislabeling			OoD/Outlier		Atypical
		Uniform	Asymmetric	Instance	Near OoD	Far OoD	
Learning-based (Margin)	AUM (Pleiss et al., 2020)	$\checkmark(\checkmark)$	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$
Learning-based (Uncertainty)	Data Maps (Swayamdipta et al., 2020)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
	Data-IQ (Seedat et al., 2022a)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$
Learning-based (Loss)	Small-Loss (Xia et al., 2021)	$\times(\times)$	$\times(\checkmark)$	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
	Action scores (Arriaga et al., 2023)	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$
	RHO-Loss (Mindermann et al., 2022)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$
Learning-based (Gradient)	GraNd (Paul et al., 2021)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\checkmark)$
	VoG (Agarwal et al., 2022)	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\checkmark)$
Learning-based (Forgetting)	Forgetting Scores (Toneva et al., 2019)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\checkmark)$
	SSFT (Maini et al., 2022)	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$
Learning-based (Statistics)	Detector (Jia et al., 2022)	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
	EL2N (Paul et al., 2021)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\checkmark)$
Distance-based	Prototypicality (Sorscher et al., 2022)	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$
Information theory	PVI (Ethayarajh et al., 2022)	$\times(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
Statistical-based	Cleanlab (Northcutt et al., 2021a)	$\times(\times)$	$\checkmark(\checkmark)$	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
	ALLSH (Zhang et al., 2022)	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\times)$
	Agreement (Carlini et al., 2019)	$\times(\times)$	$\times(\times)$	$\times(\times)$	$\times(\checkmark)$	$\times(\times)$	$\times(\checkmark)$
	Data Shapley (Ghorbani & Zou, 2019)	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$	$\checkmark(\checkmark)$	$\times(\times)$	$\times(\times)$

**OoD/Outlier:** Samples where the covariates undergo a transformation/shift, such that sample  $(x, y) \rightarrow (x^*, y)$ . The distinction between easy and hard samples lies in the feature space, leading to different marginal probability distributions  $P_X^e(x) \neq P_X^h(x)$ . Further, for any subset  $S$  within the support of  $P_X$ ,  $P_X^h(S) = 0$ . Note, conditional probability distributions remain consistent where  $P_{Y|X}^e(y|x) = P_{Y|X}^h(y|x)$ . We consider two subtypes differing in degree of shift, for clarity denoted by an arbitrary distance measure between distributions  $\text{dist}(\cdot, \cdot)$ .

- *Near OoD* (Anirudh & Thiagarajan, 2023; Mu & Gilmer, 2019; Hendrycks & Dietterich, 2018; Sun et al., 2023b; Yang et al., 2023; Tian et al., 2021): samples which have their features transformed or shifted such that they remain proximal to the original samples in  $\mathcal{D}$ , e.g. introducing noise, pixelating an image, or adding subtle texture changes. In this case,  $\text{dist}(P_X^h, P_X)$  is positive but relatively small, indicating the nearness of the perturbed distribution to the original. We can represent this bounded distance such that  $0 < \text{dist}(P_X^h, P_X) \leq \epsilon$ , with  $\epsilon > 0$ .
- *Far OoD* (Mukhoti et al., 2022; Winkens et al., 2020; Graham et al., 2022; Yang et al., 2023): samples which are distinct and likely unrelated to samples in  $\mathcal{D}$ , often not belonging to the same data-generating process. This could be by sampling from a different dataset or a perturbation s.t.  $\text{dist}(P_X^h, P_X)$  is significantly large, i.e.  $\text{dist}(P_X^h, P_X) \gg \epsilon$ . For example, for a dataset of digits, images of dogs or cats are distinctly different and unrelated. They are not just rare occurrences but images of dogs or cats represent a different data generation process compared to the digits.

**Atypical:** Samples that, although rare, are still valid instances deviating from common patterns (Yuksekgonul et al., 2023). Atypical samples are inherently part of the primary data distribution, but located in its long tail or less frequent regions (Feldman, 2020; Hooker et al., 2019; Hooker, 2021; Agarwal et al., 2022). The distinction between easy and hard samples leads to different marginal probability distributions  $P_X^e(x) \neq P_X^h(x)$ , where  $P_X^h(x)$  is very small, highlighting their rarity or infrequency. Here, for any subset  $S$  within the support of  $P_X$ ,  $P_X^h(S) > 0$ . This signifies that the long-tail samples, though rarer in occurrence, are still within the bounds of the primary data-generating process. For example, these could be images with atypical variations or vantage points compared to the standard pattern (Agarwal et al., 2022).

**Contrasting OoD/Outlier and Atypical.** Both have different marginal probability distributions  $P_X^e(x) \neq P_X^h(x)$ . The difference is that OoD/Outliers come from a shifted or completely different distribution than the original data, falling outside the support of  $P_X$ . In contrast, atypical samples are rare samples from the tails and could naturally arise, falling *within* the support of  $P_X$ .

From a practitioner’s standpoint, these distinctions can dictate different courses of action. OoD/outliers represent likely anomalies or errors that we should detect for potential sculpting or filtering from the dataset. In contrast, atypical samples are rare cases deviating from the “norm”, with no or limited similar examples. They are still valid points and should not necessarily be discarded; in fact, there might be a need to gather more of such samples. The goal of surfacing atypical samples is both for dataset auditing and understanding edge cases. We provide additional example images to provide further intuition of the difference between OoD and Atypical in Appendix A, Fig. 3.

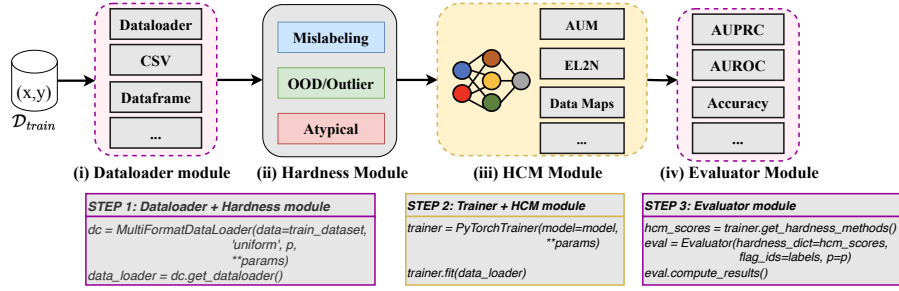


Figure 2: H-CAT facilitates comprehensive benchmarking of HCMs for multiple hardness types. Examples of the single-chain API workflow is shown below — with example usage in Appendix C. The modules described in Sec. 3 are easily extended to new HCMs, datasets, or evaluation metrics

## 2.2 GAPS AND LIMITATIONS OF CURRENT HCMs WITHIN OUR TAXONOMY

To provide a more comprehensive perspective, we critically analyze various HCMs within our proposed taxonomy. Alarming, we find inconsistent and diverse definitions of “hardness”, even within the same class of methods (refer to Table 1). Such discrepancies, as evidenced by their experimental evaluations, point to a challenge: HCMs, despite appearing to measure similar constructs, in fact, evaluate different “hardness” dimensions. Table 1 highlights two critical deficiencies discussed below, underscoring the urgent need for a systematic evaluation framework that accurately captures the scope and applicability of each HCM across different hardness types.

**Issue 1: Qualitative or indirect measures.** Many HCMs limit their evaluation to (i) *qualitative* analyses: merely showcasing flagged samples, or (ii) *indirect measures*: for example, showing downstream performance improved when removing flagged samples without directly quantifying what the HCM captures. This overlooks the necessity for an objective, quantitative evaluation.

**Issue 2: Narrow and unrepresentative scope.** HCMs that conduct quantitative evaluation focus on a single hardness type and often target the simplest manifestation. An example of this shortfall is seen in the handling of mislabeling. Many HCMs only focus on uniform mislabeling, thereby failing to account for the more realistic and complex scenarios of asymmetric or instance-wise mislabeling. Beyond this, many HCMs only test on mislabeling, overlooking other types of hardness.

These limitations emphasize the value of our fine-grained taxonomy in categorizing and evaluating various hardness types. By laying a solid foundation for the systematic evaluation and comparison of HCMs, our taxonomy enables the design and selection of HCMs that are tailored to specific hardness challenges, thereby promoting the development of more robust ML systems.

## 3 H-CAT: A BENCHMARKING FRAMEWORK FOR HCMs

To address the aforementioned limitations and facilitate benchmarking of HCMs, we propose the Hardness Characterization Analysis Toolkit (H-CAT). H-CAT serves two purposes: (1) *empirical benchmarking standard*: supporting comprehensive and quantitative benchmarking of HCMs on different hardness types within the taxonomy across multiple data modalities and (2) *software toolkit*: H-CAT unifies multiple HCMs under a single unified interface for easy usage by practitioners.

**H-CAT Design.** H-CAT has four core modules as described below, which are called sequentially (Figure 2). The framework follows widely adopted object-oriented paradigms with fit-predict interfaces (here, update-score). The workflow is simple with single-chain API calls. The stepwise composability aims to facilitate easy benchmarking and allows H-CAT to be used outside of benchmarking as a data characterization tool by practitioners.

- **Dataloader module:** loads a variety of data types for ease of use, including Torch datasets, NumPy arrays, and Pandas DataFrames, allowing users to easily use H-CAT with their chosen datasets.
- **Hardness module:** generates controllable hardness for different hardness types in the taxonomy.
- **HCM module:** provides a unified HCM interface with 13 HCMs implemented. It wraps the trainer module which is a conventional PyTorch training loop.
- **Evaluator module:** computes the HCM evaluation metrics to correctly identify the ground-truth hard samples using the scores provided by each HCM. User-specified metrics can easily be included by operating on the raw HCM scores.

**Extensibility.** H-CAT is easily extendable to include new HCMs, hardness types, or datasets by defining a simple wrapper class. For details and step-by-step code examples, refer to Appendix C.



## 4 BENCHMARKING FRAMEWORK SETUP

We now describe three key aspects of the benchmarking setup — the implementation of the Hardness Module, HCM Module and Evaluator Module.

### 4.1 HARDNESS MODULE

We describe the implementation of our hardness perturbations  $h$ . We focus on image data here, as most HCMs (10/13) have been developed for this modality. However, we discuss the corresponding hardness perturbations for tabular data in the Appendix.

**Mislabeling:** (i) *Uniform*: random mislabeling, drawn uniformly from all possible class labels. (ii) *Asymmetric*: random mislabeling, drawn asymmetrically via a Dirichlet distribution (Zhang et al., 2021; Bae et al., 2022; Zhu et al., 2022). We denote a special case of asymmetric applicable to ordinal labels, namely *adjacent*. Here, the mislabeling is to the nearest numerical class, e.g. an MNIST digit 3 mislabeled as 2 or 4. (iii) *Instance*: mislabeling probability is conditioned on an instance (reflecting human mislabeling). This can often be determined by domain/user knowledge, e.g. for MNIST, 1 is likely mislabeled as 7; for CIFAR-10, an automobile could be mislabeled as a truck.

**OoD/Outlier:** (i) *Near OoD*: perturbed data is different but related. *Covariate Shift* via Gaussian noise as performed in MNIST-C (Mu & Gilmer, 2019) or *Domain Shift*: image texture from the original photographic image is changed via edge detection and smoothing (Median filter). (ii) *Far OoD*: perturbed data is distinctly different and unrelated. We replace a subset of data with unrelated data, e.g. for MNIST replace with CIFAR-10 images and vice versa.



Figure 3: Examples providing intuition on the difference between OoD and Atypical.

**Atypical:** (i) *Shift*: translate and shift the image, causing portions to be cut off in an atypical manner. (ii) *Zoom*: create an atypical perspective by magnifying (X2) features usually seen at a smaller scale.

### 4.2 HCM MODULE

We include 13 widely used HCMs applicable to supervised classification under a unified interface. The HCMs span the range of HCM classes, at least one per class from Table 1: ■ **Learning-based<sup>1</sup> (Uncertainty)**: Data Maps (Swayamdipta et al., 2020) and Data-IQ (Seedat et al., 2022a); ■ **Learning-based (Loss)**: Sample-Loss (Xia et al., 2021; Arriaga et al., 2023); ■ **Learning-based (Margin)**: Area-under-the-margin (AUM) (Pleiss et al., 2020); ■ **Learning-based (Gradient)**: GraNd (Paul et al., 2021) and VoG (Agarwal et al., 2022); ■ **Learning-based (Statistics)**: EL2N (Paul et al., 2021) and Noise detector (Jia et al., 2022); ■ **Learning-based (Forgetting)**: Forgetting scores (Toneva et al., 2019); ■ **Statistical measures**: Cleanlab (Northcutt et al., 2021a), ALLSH (Zhang et al., 2022), Agreement (Carlini et al., 2019); ■ **Distance-based**: Prototypicality (Sorscher et al., 2022).

Specifically, we focus on HCMs that plug into the training loop and *do not*: (i) alter training, (ii) require repeated training, (iii) need additional datasets beyond the training set, or (iv) require training additional models. Consequently, we exclude the following: RHO-Loss (Mindermann et al., 2022) requires an additional irreducible loss model, SSFT (Maini et al., 2022) requires fine-tuning on a validation dataset, PVI (Ethayarajh et al., 2022) requires training a Null model. We also do not consider Data Shapley (Ghorbani & Zou, 2019) and variants (e.g. Beta Shapley (Kwon & Zou, 2022)), which have been shown to be computationally infeasible with numerical instabilities for higher dimensional data such as MNIST and CIFAR-10 with  $> 1000$  samples (Wang & Jia, 2023).

### 4.3 EVALUATOR MODULE

We directly assess the HCM’s capability to detect the hard samples. Recall that HCMs assign a score  $s$  to each sample  $(x, y)$  and then apply a threshold  $\tau$  to assign samples a group  $g \in \mathcal{G}$ , where  $\mathcal{G} = \{Easy, Hard\}$ . Many HCMs do not explicitly state how to define  $\tau$ ; hence to account for this we compute two widely used metrics: *AUPRC* (Area Under Precision-Recall Curve) and *AUROC* (Area Under Receiver Operating Curve) — for hard sample detection performance, which we denote *D-AUPRC* and *D-AUROC*<sup>2</sup>. User-specified metrics are easily computed on raw HCM scores.

<sup>1</sup>Learning-based generally refers to learning/training dynamics based HCMs

<sup>2</sup>to distinguish them from the typical downstream performance metrics.

## 5 COMPREHENSIVE EVALUATION OF HCMs USING H-CAT

We evaluate 13 different HCMs (spanning a range of techniques) across 8 distinct hardness types. To the best of our knowledge, this represents the first comprehensive HCM evaluation, encompassing over **14K** experimental setups (specific combination of HCM, hardness type, perturbation proportion, dataset, model, and seed).

We primarily focus on image datasets, as this is the modality for which the majority of the HCMs (10/13) have been developed. We use the MNIST (LeCun et al., 2010) and CIFAR-10 (Krizhevsky et al.) datasets, as their use is well-established in the HCM literature (Paul et al., 2021; Swayamdipta et al., 2020; Pleiss et al., 2020; Toneva et al., 2019; Maini et al., 2022; Jiang et al., 2021; Mindermann et al., 2022; Baldock et al., 2021). Importantly, they are realistic images yet contain almost no/little mislabeling ( $<0.5\%$ ) (Northcutt et al., 2021b). This contrasts other common image datasets like ImageNet that contain significant mislabeling (over 5%), hence we cannot perform controlled experiments. Furthermore, to show generalizability across modalities, we also evaluate on tabular datasets “Covertime” and “Diabetes130US” benchmark datasets (Grinsztajn et al., 2022) from OpenML (Vanschoren et al., 2014) (see Appendix D.6).

To assess HCM sensitivity to the backbone model, we use two different models for our image experiments with different degrees of parameterization, LeNet and ResNet-18. All experiments are repeated with 3 random seeds and for varying proportions  $p$  of hard samples.

We present aggregated results in Figs. 4-7, with more granular results in Appendix D — along with additional experiments. The main paper shows results for 6 out of 8 hardness types. We include results for other sub-types not covered in the main paper in Appendix D including: Domain shift (a type of Near OoD), Zoom shift (a type of Atypical), and Adjacent (a special case of Asymmetric mislabeling), offering similar conclusions. We investigate three aspects of HCMs (A-C), distilling the results into **benchmarking takeaways** and **practical tips**.

**A. Hardness detection performance.** Directly evaluate HCM capabilities to detect the hard samples for different hardness types, for varying perturbation proportions  $p$  – see Fig. 4.

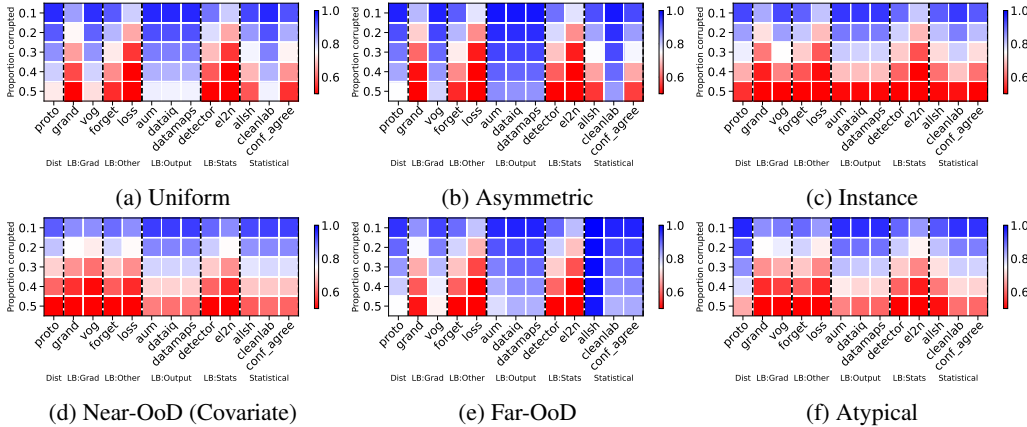


Figure 4: D-AUPRC for different HCMs for different hardness types aggregated across setups. We vary the proportion perturbed, i.e. the proportion of hard examples. Blue is better, red worse. We see variability of HCM capabilities across hardness types and proportions.

**Takeaway A1: Comprehensive testing is vital.** Many HCMs are assessed on a single hardness type. The reality is hardness manifests in many ways. We show HCM performance varies across hardness types and across proportions of hardness, with some more challenging than others, e.g. instance is harder than uniform. This result shows the critical need for comprehensive HCM evaluation.

**Takeaway A2: Hardness types vary in difficulty.** We find that different types of hardness are easier or harder to characterize. For instance, uniform mislabeling or Far-OoD are much easier than data-specific hardness like instance and Atypical. Given the performance differences of HCMs on different hardness types, it becomes important to understand the hardness type expected in practice.

**Takeaway A3: Learning dynamics-based methods with respect to output confidence are effective general-purpose HCMs.** In selecting a general-purpose HCM, we find that HCMs that characterize samples using learning dynamics on the confidence — uncertainty-based methods, which use probabilities (DataMaps, Data-IQ) or logits (AUM), are the best performing in terms of AUPRC across the board.

**Takeaway A4: HCMs typically used for computational efficiency are surprisingly uncompetitive.** We find that HCMs that leverage gradient changes (e.g. GraNd), typically used for selection to improve computational efficiency, fare well at low  $p$ . However, at higher  $p$ , they become notably less competitive compared to simpler and computationally cheaper methods.

**Practical Tip A1: HCMs should only be used when hardness proportions are low.** In general, different HCMs have significantly reduced performance at higher proportions of hardness. This is expected as we get closer to 0.5 since it’s harder to identify a clear difference between samples.

**B. Rankings and Statistical Significance.** Compare the ranking of methods, as well as assess statistical significance of performance differences using critical difference diagrams (CD diagram) (Demšar, 2006) based on the Siegel-Friedman method ( $p \leq 0.05$ ) — see Figs. 5 and 6.

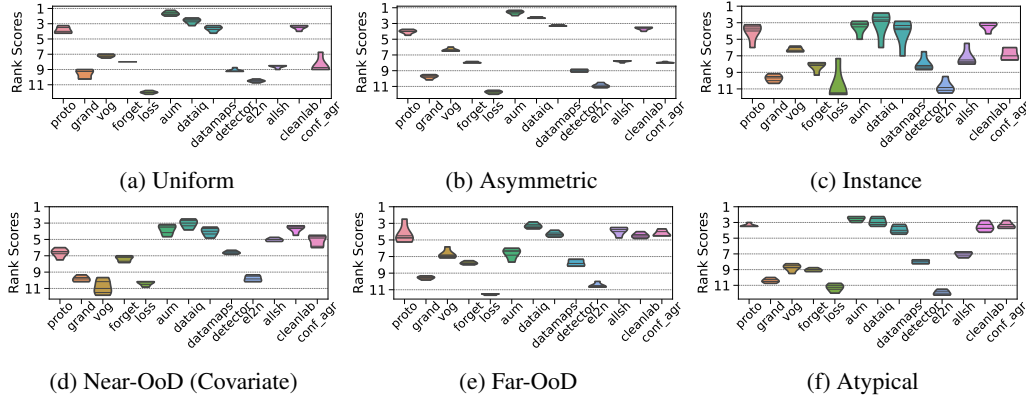


Figure 5: Performance rankings of HCMs vary depending on the hardness type.

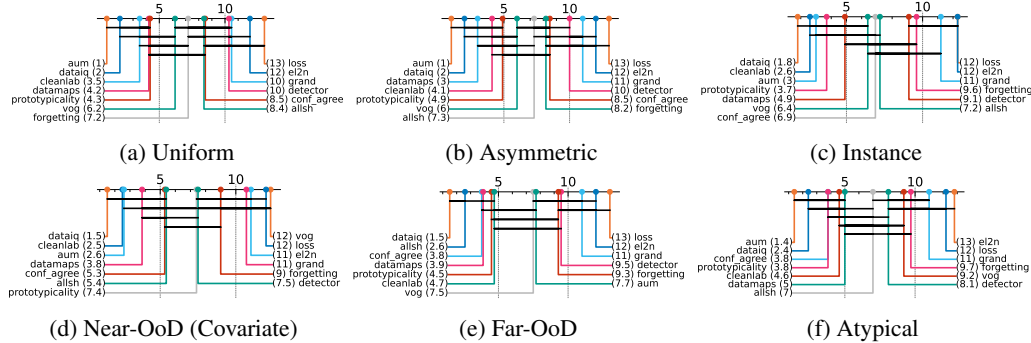


Figure 6: Critical difference diagrams highlight that similar categories/classes of HCMs *do not* have a statistically significant difference in their performance, indicated by the horizontal black lines linking HCMs which are not statistically different. The numbers in brackets denote mean rank.

**Takeaway B1: Individual HCMs within a broad “class” of methods are NOT statistically different.** We find from the critical difference diagrams that methods falling into the same class of characterization are not statistically significant from one another (based on the black connected lines), despite the minor performance differences between them. Hence, practitioners should select an HCM within the broad HCM class most suitable for the application.

**Practical Tip B1: Selecting an HCM based on the hardness is useful.** We find that confidence is a good general-purpose tool if one does not know the type of hardness. However, if one knows the hardness, one can better select the HCM. For example, Prototypicality, as expected, is very strong on instance hardness as we are able to match samples via similarity of classes in embedding space.



**Practical Tip B2: Divergence and distance-based methods are suitable primarily for distributional changes.** Divergence and distance-based methods such as ALLSH and Prototypicality should primarily be used if the hardness is with respect to a shift of the data itself rather than mislabeling.

**C. Stability/Consistency.** The rank ordering of samples is important in data characterization (Maini et al., 2022; Wang & Jia, 2023; Seedat et al., 2022a). Hence, we desire HCM scores to be stable to ensure consistent insights. As is standard (Maini et al., 2022; Seedat et al., 2022a), we compute the Spearman rank correlation across multiple runs — see Fig. 7. We also assess HCM stability/consistency across backbone models and parameterizations in Appendix D.

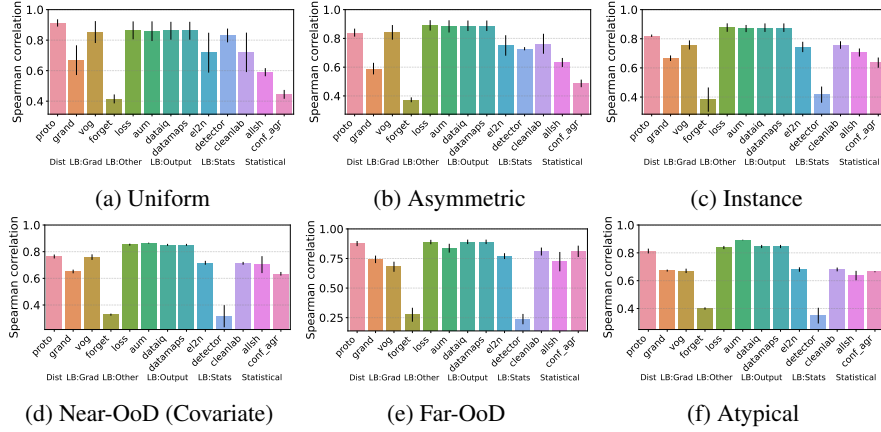


Figure 7: Certain classes of HCMs are more stable and consistent than others (retaining rank order), with higher Spearman rank correlations for multiple runs of the HCM on the same data

**Takeaway C1: Learning dynamics-based HCMs using output metrics or distance-based HCMs are most stable and consistent.** We find across all hardness types that the Spearman rank correlation is highest for HCMs that use learning dynamics on the outputs or use distance measures — specifically Uncertainty (Probabilities), Margins (Logits), Loss, or Prototypicality. The low correlation for other HCMs highlights sensitivity to the run itself when characterizing data, which would lead to inconsistent ordering.

**Takeaway C2: Insights consistent across backbone models & parameterizations.** As shown in Appendix D, we find similar results (as above) for the Spearman rank correlation of HCM scores across different backbone models and parameterizations. As the scores are consistent, this indicates that the findings and insights computed on those scores will also be consistent, i.e. those HCMs which were the most stable and consistent remain the most stable and consistent.

**Practical Tip C1: Select a stable HCM.** Certain HCMs are more sensitive to randomness. Hence, it is advised to select stable and consistent HCMs (higher Spearman correlation) to avoid such effects.

## 6 DISCUSSION

We introduce H-CAT, a comprehensive benchmarking framework for hardness characterization in data-centric AI. We analyzed 13 HCMs spanning a range of techniques for a variety of setups - over 14K. We hope our framework and insights addressing calls for rigorous benchmarking (Guyon, 2022) and understanding of existing ML methods (Lipton & Steinhardt, 2019; Snoek et al., 2018) will spur advancements in data-centric AI and that H-CAT helps practitioners better evaluate and use HCMs.

**Limitations & Future Work.** No benchmark can exhaustively test all the possible hardness manifestations and this work is no different. However, we cover multiple instances from the three fundamental types of hardness, which is significantly broader than any prior work. Building on this work, future research could investigate cases where multiple types of “hardness” manifest simultaneously (e.g. Near-OoD and mislabeling together) or where hardness is continuous rather than binary. To spur this, we provide an example of simultaneous hardness in Appendix D. From a usage perspective, future work could also look into the best way hardness scores could be used to guide better model training (e.g. data curriculum). Finally, we highlight that HCMs, and by extension H-CAT, cannot tell you which hardness type exists in a dataset; rather, H-CAT serves to benchmark the capabilities of HCMs or as a unified HCM interface.

## ACKNOWLEDGEMENTS

The authors would like to thank Nicolas Huynh, Alicia Curth and the anonymous reviewers for their helpful comments and discussions on earlier drafts of this paper and Robert Davis for discussions on the code. NS and FI gratefully acknowledge funding from the Cystic Fibrosis Trust and NSF grant (1722516), respectively. This work was supported by Azure sponsorship credits granted by Microsoft’s AI for Good Research Lab.

## ETHICS AND REPRODUCIBILITY STATEMENTS

**Ethics Statement.** HCMs that accurately identify hard samples can help make models more robust and reliable. This paper highlights the importance of rigorously evaluating HCMs to better understand their capabilities and guide better usage. This paper aims to enable the community to conduct a more systematic hardness characterization through the proposed taxonomy and benchmarking framework.

**Reproducibility Statement.** We include implementation details about our benchmark in Sec. 3 and 4, as well as in Appendix B. The code for the H-CAT framework can be found at <https://github.com/seedatnabeel/H-CAT> or <https://github.com/vanderschaarlab/H-CAT>. Step-by-step code examples can be found in the repository and in Appendix C along with a guide on how to extend H-CAT to new HCMs, hardness types, and datasets.

## REFERENCES

- Chirag Agarwal, Daniel D’souza, and Sara Hooker. Estimating example difficulty using variance of gradients. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10368–10378, 2022.
- Rushil Anirudh and Jayaraman J Thiagarajan. Out of distribution detection via neural network anchoring. In *Asian Conference on Machine Learning*, pp. 32–47. PMLR, 2023.
- Octavio Arriaga, Sebastian Palacio, and Matias Valdenegro-Toro. Difficulty estimation with action scores for computer vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 245–253, 2023.
- HeeSun Bae, Seungjae Shin, Byeonghu Na, JoonHo Jang, Kyungwoo Song, and Il-Chul Moon. From noisy prediction to true label: Noisy prediction calibration via generative model. In *International Conference on Machine Learning*, pp. 1277–1297. PMLR, 2022.
- Robert Baldock, Hartmut Maennel, and Behnam Neyshabur. Deep learning through the lens of example difficulty. *Advances in Neural Information Processing Systems*, 34:10876–10889, 2021.
- Chris Bedi, Edward Cone, Brian Foster, Riva Froymovich, Colleen Haikes, Michael Hubbard, Brian Klapac, Alan Marks, Debbie O’Brien, Kathy O’Connell, Gaby Pacini-Pinto, Kevin Palmer, Rob Pickering, Chris Pope, Dan Rogers, Julia Smith, and Dave Wright. *The Global CIO Point of View The New Agenda for Transformative Leadership: Reimagine Business for Machine Learning*. Oxford Econometric & Service Now, 2019.
- Nicholas Carlini, Ulfar Erlingsson, and Nicolas Papernot. Distribution density, tails, and outliers in machine learning: Metrics and applications. *arXiv preprint arXiv:1910.13427*, 2019.
- Haihua Chen, Jiangping Chen, and Junhua Ding. Data evaluation and enhancement for quality improvement of machine learning. *IEEE Transactions on Reliability*, 70(2):831–847, 2021.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7:1–30, 2006.
- Simao Eduardo, Kai Xu, Alfredo Nazabal, and Charles Sutton. Repairing systematic outliers by learning clean subspaces in VAEs. *arXiv preprint arXiv:2207.08050*, 2022.
- Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. Understanding dataset difficulty with  $v$ -usable information. In *International Conference on Machine Learning*, pp. 5988–6008. PMLR, 2022.

- Vitaly Feldman. Does learning require memorization? A short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 954–959, 2020.
- Amirata Ghorbani and James Zou. Data Shapley: Equitable valuation of data for machine learning. In *International Conference on Machine Learning*, pp. 2242–2251. PMLR, 2019.
- Mark S Graham, Petru-Daniel Tudosi, Paul Wright, Walter Hugo Lopez Pinaya, U Jean-Marie, Yee H Mah, James T Teo, Rolf Jager, David Werring, Parashkev Nachev, et al. Transformer-based out-of-distribution detection for clinically safe segmentation. In *International Conference on Medical Imaging with Deep Learning*, pp. 457–476. PMLR, 2022.
- Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *Advances in Neural Information Processing Systems*, 35: 507–520, 2022.
- Nitin Gupta, Hima Patel, Shazia Afzal, Naveen Panwar, Ruhi Sharma Mittal, Shanmukha Guttula, Abhinav Jain, Lokesh Nagalapatti, Sameep Mehta, Sandeep Hans, et al. Data quality toolkit: Automatic assessment of data quality and remediation for machine learning datasets. *arXiv preprint arXiv:2108.05935*, 2021.
- Isabelle Guyon. The data-centric era: How ML is becoming an experimental science. *Keynote: Advances in Neural Information Processing Systems*, 2022.
- Bo Han, Quanming Yao, Tongliang Liu, Gang Niu, Ivor W Tsang, James T Kwok, and Masashi Sugiyama. A survey of label-noise representation learning: Past, present and future. *arXiv preprint arXiv:2011.04406*, 2020.
- Lasse Hansen, Nabeel Seedat, Mihaela van der Schaar, and Andrija Petrovic. Reimagining synthetic tabular data generation through data-centric AI: A comprehensive benchmark. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2018.
- Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. Using trusted data to train deep networks on labels corrupted by severe noise. *Advances in neural information processing systems*, 31, 2018.
- Elisabeth RM Heremans, Nabeel Seedat, Bertien Buyse, Dries Testelmans, Mihaela van der Schaar, and Maarten De Vos. U-PASS: An uncertainty-guided deep learning pipeline for automated sleep staging. *Computers in Biology and Medicine*, pp. 108205, 2024.
- Sara Hooker. Moving beyond “algorithmic bias is a data problem”. *Patterns*, 2(4):100241, 2021.
- Sara Hooker, Aaron Courville, Gregory Clark, Yann Dauphin, and Andrea Frome. What do compressed deep neural networks forget? *arXiv preprint arXiv:1911.05248*, 2019.
- Abhinav Jain, Hima Patel, Lokesh Nagalapatti, Nitin Gupta, Sameep Mehta, Shanmukha Guttula, Shashank Mujumdar, Shazia Afzal, Ruhi Sharma Mittal, and Vitobha Munigala. Overview and importance of data quality for machine learning tasks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 3561–3562, 2020.
- Mohammad Hossein Jarrahi, Ali Memariani, and Shion Guha. The principles of data-centric AI (DCAI). *arXiv preprint arXiv:2211.14611*, 2022.
- Qingrui Jia, Xuhong Li, Lei Yu, Jiang Bian, Penghao Zhao, Shupeng Li, Haoyi Xiong, and Dejing Dou. Learning from training dynamics: Identifying mislabeled data beyond manually designed features. *arXiv preprint arXiv:2212.09321*, 2022.
- Ziheng Jiang, Chiyuan Zhang, Kunal Talwar, and Michael C Mozer. Characterizing structural regularities of labeled data in overparameterized models. In *International Conference on Machine Learning*, pp. 5034–5044. PMLR, 2021.

- Sanjay Krishnan, Jiannan Wang, Eugene Wu, Michael J Franklin, and Ken Goldberg. ActiveClean: Interactive data cleaning for statistical modeling. *Proceedings of the VLDB Endowment*, 9(12): 948–959, 2016.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. CIFAR-10 (Canadian Institute for Advanced Research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- Yongchan Kwon and James Zou. Beta Shapley: A unified and noise-reduced data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 8780–8802. PMLR, 2022.
- Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. CleanML: A study for evaluating the impact of data cleaning on ML classification tasks. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, pp. 13–24. IEEE, 2021.
- Weixin Liang, Girmaw Abebe Tadesse, Daniel Ho, L Fei-Fei, Matei Zaharia, Ce Zhang, and James Zou. Advances, challenges and opportunities in creating data for trustworthy AI. *Nature Machine Intelligence*, 4(8):669–677, 2022.
- Zachary C Lipton and Jacob Steinhardt. Research for practice: Troubling trends in machine-learning scholarship. *Communications of the ACM*, 62(6):45–53, 2019.
- Zifan Liu, Zhechun Zhou, and Theodoros Rekatsinas. Picket: Guarding against corrupted data in tabular data during learning and inference. *The VLDB Journal*, pp. 1–29, 2022.
- Pratyush Maini, Saurabh Garg, Zachary Chase Lipton, and J Zico Kolter. Characterizing datapoints via second-split forgetting. In *Advances in Neural Information Processing Systems*, 2022.
- Anna Majkowska, Sid Mittal, David F Steiner, Joshua J Reicher, Scott Mayer McKinney, Gavin E Duggan, Krish Eswaran, Po-Hsuan Cameron Chen, Yun Liu, Sreenivasa Raju Kalidindi, et al. Chest radiograph interpretation with deep learning models: assessment with radiologist-adjudicated reference standards and population-adjusted evaluation. *Radiology*, 294(2):421–431, 2020.
- Sören Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch, Winnie Xu, Benedikt Hölten, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, et al. Prioritized training on points that are learnable, worth learning, and not yet learnt. In *International Conference on Machine Learning*, pp. 15630–15649. PMLR, 2022.
- Norman Mu and Justin Gilmer. MNIST-C: A robustness benchmark for computer vision. *arXiv preprint arXiv:1906.02337*, 2019.
- Jishnu Mukhoti, Tsung-Yu Lin, Bor-Chun Chen, Ashish Shah, Philip HS Torr, Puneet K Dokania, and Ser-Nam Lim. Raising the bar on the evaluation of out-of-distribution detection. *arXiv preprint arXiv:2209.11960*, 2022.
- Afshan Nabi, Berke Dilekoglu, Ogun Adebali, and Ozgur Tastan. Discovering misannotated Incrnas using deep learning training dynamics. *Bioinformatics*, 39(1):btac821, 2023.
- Andrew Ng, Lora Aroyo, Cody Coleman, Greg Diamos, Vijay Janapa Reddi, Joaquin Vanschoren, Carole-Jean Wu, and Sharon Zhou. NeurIPS data-centric AI workshop, 2021. URL <https://datacentricai.org/>.
- Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021a.
- Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021b.
- Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *Advances in Neural Information Processing Systems*, 34: 20596–20607, 2021.

- Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33: 17044–17056, 2020.
- Cedric Renggli, Luka Rimanic, Nezihe Merve Gürel, Bojan Karlas, Wentao Wu, and Ce Zhang. A data quality-driven view of MLOps. *IEEE Data Engineering Bulletin*, 2021.
- Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. “Everyone wants to do the model work, not the data work”: Data cascades in high-stakes AI. In *proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pp. 1–15, 2021.
- Nabeel Seedat, Jonathan Crabbé, Ioana Bica, and Mihaela van der Schaar. Data-IQ: Characterizing subgroups with heterogeneous outcomes in tabular data. In *Advances in Neural Information Processing Systems*, 2022a.
- Nabeel Seedat, Jonathan Crabbé, and Mihaela van der Schaar. Data-SUITE: Data-centric identification of in-distribution incongruous examples. In *International Conference on Machine Learning*, pp. 19467–19496. PMLR, 2022b.
- Nabeel Seedat, Jonathan Crabbé, Zhaozhi Qian, and Mihaela van der Schaar. TRIAGE: Characterizing and auditing training data for improved regression. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
- Nabeel Seedat, Fergus Imrie, and Mihaela van der Schaar. Navigating data-centric artificial intelligence with DC-Check: Advances, challenges, and opportunities. *IEEE Transactions on Artificial Intelligence*, 2023b.
- Daniel Schwartz, Uri Stern, and Daphna Weinshall. The dynamic of consensus in deep networks and the identification of noisy labels. *arXiv preprint arXiv:2210.00583*, 2022.
- Jasper Snoek, Alex Wiltschko, and Ali Rahimi. Winner’s curse? On pace, progress, and empirical rigor. In *International Conference on Learning Representations - Workshop*, 2018.
- Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S Morcos. Beyond neural scaling laws: Beating power law scaling via data pruning. In *Advances in Neural Information Processing Systems*, 2022.
- Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Hao Sun, Alex J Chan, Nabeel Seedat, Alihan Hüyük, and Mihaela van der Schaar. When is off-policy evaluation useful? a data-centric perspective. *arXiv preprint arXiv:2311.14110*, 2023a.
- Hao Sun, Boris van Breugel, Jonathan Crabbé, Nabeel Seedat, and Mihaela van der Schaar. What is flagged in uncertainty quantification? latent density models for uncertainty categorization. *Advances in Neural Information Processing Systems*, 36, 2023b.
- Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9275–9293, 2020.
- Junjiao Tian, Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Exploring covariate and concept shift for detection and confidence calibration of out-of-distribution data. *arXiv preprint arXiv:2110.15231*, 2021.



- Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In *International Conference on Learning Representations*, 2019.
- Joaquin Vanschoren, Jan N Van Rijn, Bernd Bischl, and Luis Torgo. OpenML: Networked science in machine learning. *ACM SIGKDD Explorations Newsletter*, 15(2):49–60, 2014.
- Jiachen T Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 6388–6421. PMLR, 2023.
- Amanda West. *Machine Learning Models for the Future: The rise of the data scientist*. Refinitiv Labs, 2020.
- Steven Euijong Whang, Yuji Roh, Hwanjun Song, and Jae-Gil Lee. Data collection and quality challenges in deep learning: A data-centric AI perspective. *The VLDB Journal*, pp. 1–23, 2023.
- Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.
- Xiaobo Xia, Tongliang Liu, Bo Han, Mingming Gong, Jun Yu, Gang Niu, and Masashi Sugiyama. Sample selection with uncertainty of losses for learning with noisy labels. In *International Conference on Learning Representations*, 2021.
- Jingkang Yang, Kaiyang Zhou, and Ziwei Liu. Full-spectrum out-of-distribution detection. *International Journal of Computer Vision*, pp. 1–16, 2023.
- Mert Yuksekogonul, Linjun Zhang, James Zou, and Carlos Guestrin. Beyond confidence: Reliable models should also consider atypicality. *arXiv preprint arXiv:2305.18262*, 2023.
- Shujian Zhang, Chengyue Gong, Xingchao Liu, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. ALLSH: Active learning guided by local sensitivity and hardness. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 1328–1342, 2022.
- Yivan Zhang, Gang Niu, and Masashi Sugiyama. Learning noise transition matrix from only noisy labels via total variation regularization. In *International Conference on Machine Learning*, pp. 12501–12512. PMLR, 2021.
- Zhaowei Zhu, Jialu Wang, and Yang Liu. Beyond images: Label noise transition matrix estimation for tasks with lower-quality features. In *International Conference on Machine Learning*, pp. 27633–27653. PMLR, 2022.