Automatic Feasibility Study via Data Quality Analysis for ML: A Case-Study on Label Noise

Cedric Renggli^{*,†}, Luka Rimanic^{*,†}, Luka Kolar^{*,†}, Wentao Wu[§], Ce Zhang[†]

[†]ETH Zurich, [§]Microsoft Research

[†]{cedric.renggli, luka.rimanic, luka.kolar, ce.zhang}@inf.ethz.ch, [§]wentao.wu@microsoft.com

Abstract—In our experience of working with domain experts who are using today's AutoML systems, a common problem we encountered is what we call "unrealistic expectations" - when users are facing a very challenging task with a noisy data acquisition process, while being expected to achieve startlingly high accuracy with machine learning (ML). Many of these are predestined to fail from the beginning. In traditional software engineering, this problem is addressed via a feasibility study, an indispensable step before developing any software system. In this paper, we present Snoopy, with the goal of supporting data scientists and machine learning engineers performing a systematic and theoretically founded feasibility study before building ML applications. We approach this problem by estimating the irreducible error of the underlying task, also known as the Bayes error rate (BER), which stems from data quality issues in datasets used to train or evaluate ML models. We design a practical Bayes error estimator that is compared against baseline feasibility study candidates on 6 datasets (with additional real and synthetic noise of different levels) in computer vision and natural language processing. Furthermore, by including our systematic feasibility study with additional signals into the iterative label cleaning process, we demonstrate in end-to-end experiments how users are able to save substantial labeling time and monetary efforts. Index Terms-Feasibility Study for ML, Data Quality for ML

I. INTRODUCTION

Modern software development is typically guided by software engineering principles that have been developed and refined for decades [1]. Even though such principles are yet to come to full fruition regarding the development of machine learning (ML) applications, in recent years we have witnessed a surge of work focusing on ML usability through supporting efficient ML systems [2]–[4], enhancing developer's productivity [5]–[7], and supporting the ML application development process itself [8]–[15].

<u>Calls for a Feasibility Study of ML</u>: In this paper, we focus on one specific "failure mode" that we frequently witness whilst working with a range of domain experts, which we call "*unrealistic expectations*." Unlike classical software artifacts, the quality of ML models (e.g., its accuracy) is often a reflection of the *data quality* used to train or test the model. We regularly see developers that work on challenging tasks with a dataset that is *too* noisy to meet the unrealistically high expectations on the accuracy that can be achieved with ML — such a project is predestined to fail. Ideally, problems of

this type should be caught *before* a user commits significant amount of resources to train or tune ML models.

In practice, if this were done by a *human ML consultant*, she would first analyze the representative dataset for the defined task and assess the *feasibility* of the target accuracy — if the target is not achievable, one can then explore alternative options by refining the dataset, the acquisition process, or investigating different task definitions. Borrowing the term from classic software engineering, we believe that such a *feasibility study* step is crucial to the usability of future ML systems for application developers. In this paper, we ask: *Can we provide some systematic and theoretically understood guidance for this feasibility study process*?

Quantitative Understanding of "Data Quality for ML": Data quality, along with its cleaning, integration, and acquisition, is a core data management problem that has been intensively studied in the last few decades [16]-[22]. Agnostic to ML workloads, the data management community has been conducting a flurry of work aimed at understanding and quantifying data quality issues [23]-[26]. In addition to these fundamental results, the presence of an ML training procedure as a downstream task over data provides both challenges and opportunities. On the one hand, systematically mapping these challenges to ML model quality issues is largely missing (with the prominent exceptions of [27]–[29] together with some of our own previous efforts [30]-[33]). On the other hand, the ML training procedure provides a quantitative metric to measure precisely the utility of data, or its quality. In this paper, we take one of the early steps in this direction and ask: "Can we quantitatively map the quality requirements of a downstream ML task to the requirements of the data quality of the upstream dataset?"

<u>The Scope and Targeted Use Case</u>: As one of the first attempts towards understanding this fundamental problem, this paper by no means provides a complete solution. Instead, we have a very specific application scenario in mind for which we develop a deep understanding both theoretically and empirically. Specifically, we focus on the case in which a user has access to a dataset D, large enough to be representative for the underlying task. The user is facing the following question: Is my current data artefact D clean/good enough for <u>some</u> ML models to reach a target accuracy α_{target} ? If the answer to this question is "Yes", the user can start expensive AutoML runs and hopefully can find a model that can reach α_{target} ;

^{*} The first three authors contributed equally to this paper.

otherwise, it would be better for the user to improve the *quality* (via data cleaning for example) of the data artefact \mathcal{D} before starting an AutoML run which is "doomed to disappoint". We call the process of answering this question a "*feasibility study*". Our main goal is to derive a strategy for the feasibility study that is (i) *informative and theoretically justified*, (ii) *inexpensive*, and (iii) *scalable*.

Such a process can be useful in many scenarios. In this paper we develop a fundamental building block of the feasibility study and evaluate it focusing on one specific use case as follows — the dataset \mathcal{D} is noisy in its labels, probably caused by (1) the inherent noise of the data collection process such as *crowd sourcing* [34]–[37], or (2) *bugs* in the data preparation pipeline (which we actually see quite often in practice). The user has a target accuracy α_{target} and can spend time and money on two possible operations: (1) manually clean up some labels in the dataset, or (2) find and engineer suitable ML models, manually or automatically.

<u>Challenges of the Strawman</u>: There are multiple strawman strategies, each of which has its own challenge. A natural, rather trivial approach is to run a *cheap proxy model*, e.g., logistic regression, to get an accuracy α_{proxy} , and use it to produce an estimator $\alpha_{\text{est}} = c \cdot \alpha_{\text{proxy}}$, say, for some $c \in [1, 1/\alpha_{\text{proxy}}]$. The challenge of this approach is to pick a universally good constant c, which depends not only on the data but also on the cheap proxy model. It is important, but often challenging, to provide a principled and theoretically justified way of adjusting the gap between α_{proxy} and α_{est} .

An alternative approach would be to simply fire up an AutoML run that systematically looks at various configurations of ML models and potentially neural architectures. Given enough time and resources, this could converge closely to the best possible accuracy that one can achieve on a given dataset; nevertheless, this can be very expensive and time consuming, thus might not be suitable for a quick feasibility study.

Feasibility Study: Theory vs. Practice: The main challenges of the strawman approaches motivate us to look at this problem in a more principled way. From a theoretical perspective, our view on feasibility study is not new, rather it connects to a decades-old ML concept known as the Bayes error rate (BER) [38], the "irreducible error" of a given task corresponding to the error rate of the Bayes optimal classifier. In fact, all factors leading to an increase of the BER can be mapped to classical data quality dimensions (e.g., "label noise" to "accuracy", or "missing features " to "completeness") [39]. Estimation of the BER has been studied intensively for almost half a century by the ML community [40]–[46]. Until recently, most, if not all, of these works are mainly theoretical, evaluated on either synthetic and/or very small datasets of often small dimensions. Over the years, we have been conducting a series of work aimed at understanding the behavior of these BER estimators on larger scale, real-wold datasets. This paper builds on two of these efforts, notably (1) a framework to compare BER estimators on large scale real-world datasets with unknown true BER [47], and (2) new convergence bounds for a simple BER estimator on top of pre-trained transformations [48]. Guided by the insights and theoretical understanding we gained from these prior works, which we treat as preliminaries and do not see them as a technical contribution of this work, we ask the following non-trivial questions:

Q1. How to use estimations of the BER for the purpose of systematic feasibility study for ML?

Q2. How can we build a scalable system to make decades of theoretical research around the BER practical and feasible on real-world datasets?

Summary of Contributions: We present Snoopy— a fast, practical and systematic feasibility study system for machine learning. We make three technical contributions.

C1. Systems Abstractions and Designs: In Snoopy, we model the problem of feasibility study as estimating a lower bound of the BER. Users provide Snoopy with a dataset representative for their ML task along with a target accuracy. The system then outputs a binary signal assessing whether the target accuracy is realistic or not. Being aware of failures (false-positives and false-negatives) in the binary output of our system, which we carefully outline and explain in this paper, we support the users in deciding on whether to "trust" the output of our system by providing additional numerical and visual aids. The technical core is a practical BER estimator. We propose a simple, but novel approach, which consults a collection of different BER estimators based on a 1NN estimator inspired by Cover and Hart [38], built on top of a collection of pre-trained feature transformations, and aggregated through taking the minimum. We provide a theoretical analysis on the regimes under which this aggregation function is justified.

C2. System Optimizations: We then describe the implementation of Snoopy, with optimizations that improve its performance. One such optimization is the *successive-halving* algorithm [49], a part of the textbook Hyperband algorithm [7], to balance the resources spent on different estimators. This already outperforms naive approaches significantly. We further improve on this method by taking into consideration the convergence curve of estimators, fusing it into a new variant of successive-halving. Moreover, noticing the iterative nature between Snoopy and the user, we take advantage of the property of kNN classifiers and implement an incremental version of the system. For scenarios in which a user cleans some labels, Snoopy is able to provide real-time feedback (0.2 ms for 10K test samples and 50K training samples).

C3. Experimental Evaluation: We perform a thorough experimental evaluation of Snoopy on 6 well-known datasets in computer vision and text classification against the baselines that use cheap and expensive proxy models. We show that Snoopy consistently outperforms the cheap, and matches the expensive strategy in terms of predictive performance for synthetic and natural label noise, whilst being computational much more efficient than both approaches. In an end-to-end use-case, where noisy datasets are iteratively cleaned up to a fraction required to achieve the target, our system enables

large savings in terms of overall cost, especially in cheap label-cost regimes. In label-cost dominated regimes (i.e., large label costs or cheap compute costs), our system adds little to no overhead compared to the baselines. Furthermore, by exploring the regimes in which Snoopy fails to provide a correct answer, we show the benefits of additional signals given to the user.

Limitations: In this paper we focus on the challenging endeavor of estimating the irreducible error for the task definition and data acquisition process, originating from data quality issues. We focus on label noise, representing one of the most prominent source for non-zero irreducible error. The exploration of other aspects of poor data quality, such as noisy or incomplete features, are left as future work. We by no means provide a conclusive solution to prevent unrealistic or very costly endeavors of training ML models with finite data. Rather, we view our contribution as a first step towards a practical treatment of this problem, which is the key for enabling a systematic feasibility study for ML. Concretely, we focus on classification tasks which, compared to other ML tasks, benefit of a solid theoretical understanding of the irreducible error and ways of estimating it. As a result, in Section II we carefully describe limiting assumptions on the data distributions, as well as failure causes and failure examples, presented in Sections III and VI respectively, hoping that this can stimulate future research from the community.

Future Extension: The feasibility study functionality targeted in this paper is ideal for new ML projects designed to replace existing "classical" code with certain accuracy. Nothing prevents data scientists and ML engineers to use Snoopy prior to any batch trained ML model though. This is particularly appealing in the context of data-centric AI, where the signal can be used to understand the impact of data actions (e.g., cleaning labels). For stream-based or continual learning there are some extra challenges. First, the window of data should typically be small in order to have a good representation of the current distribution, which renders an accurate estimation of the BER challenging. Secondly, it is unclear how a BER estimator can be designed to cope with adversarial examples. Both aspects represent interesting lines of future research. Finally, when understanding the impact of distributional drift, the test accuracy of a fixed model is typically inspected. Designing drift-aware BER estimator could to detect such a drift for any model on a distributional level is left as future work.

II. PRELIMINARIES

In this section, we give a short overview over the technical terms and the notation used throughout this paper. Let \mathcal{X} be the feature space and \mathcal{Y} be the label space, with $C = |\mathcal{Y}|$. Let $X \in \mathcal{X}, Y \in \mathcal{Y}$ be random variables. Let p(X, Y) be their joint distribution, often simplified by p(x, y) = p(X = x, Y = y). We define $\eta_y(x) = p(y|x)$ when C > 2, and $\eta(x) = p(1|x)$ when C = 2, assuming $\mathcal{Y} = \{0, 1\}$.

<u>Bayes Error Rate:</u> Bayes optimal classifier is the classifier that achieves the lowest error rate among all possible classifiers from \mathcal{X} to \mathcal{Y} , with respect to p. Its error rate is called the *Bayes error rate (BER)* and we denote it by $R_{X,Y}^*$, often abbreviated to R_X^* when Y is clear from the context. It can be expressed as $R_X^* = \mathbb{E}_X [1 - \max_{y \in \mathcal{Y}} \eta_y(x)]$.

<u>k-Nearest-Neighbor (kNN)</u> Classifier: Given a training set $\mathcal{D}_n := \{(x_1, y_1), \dots, (x_n, y_n)\}$ and a new instance x, let $(x_{\pi(1)}, \dots, x_{\pi(n)})$ be a reordering of the training instances by their distances from x, based on some metric (e.g., Euclidean or cosine dissimilarity). The kNN classifier $h_{n,k}$ and its *n*sample error rate $(R_X)_{n,k}$ are defined by

$$h_{n,k}(x) = \operatorname*{arg\,max}_{y \in \mathcal{Y}} \sum_{i=1} \mathbf{1}_{\{y_{\pi(i)}=y\}}$$
$$(R_X)_{n,k} = \mathbb{E}_{X,Y} \mathbf{1}_{\{h_{n,k}(X)\neq Y\}},$$

respectively. The *infinite-sample error rate* of kNN is given by $(R_X)_{\infty,k} = \lim_{n \to \infty} (R_X)_{n,k}$. Cover and Hart derived the following fundamental [38], and now well-known, relationship between the nearest neighbor algorithm and the BER (under mild assumptions on the underlying probability distribution):

$$(R_X)_{\infty,1} \ge R_X^* \ge \frac{(R_X)_{\infty,1}}{1 + \sqrt{1 - \frac{C(R_X)_{\infty,1}}{C-1}}}.$$
 (1)

Determining such a bound for k > 1 and C > 2 is still an open problem, and in this work we mainly focus on k = 1.

Bayes Error Estimation: The task of estimating the BER, given a finite representative dataset, is inherently difficult and has been investigated by the ML community for decades — from Fukunaga's early work in 1975 [40] to Sekeh et al.'s work in 2020 [46]. Existing BER estimators can be divided into three groups: density estimators (KDE [42], DE-kNN [50]), divergence estimator (GHP [46]), kNN classifier accuracy (1NN-kNN [41], kNN-Extrapolation [51], 1NN inspired by [38]).

As mentioned earlier, we have been conducting a series of work in order to understand the theoretical and empirical behavior of deploying and comparing BER estimators on larger scale, real-world datasets, using powerful pre-trained embeddings. This paper builds on these efforts [47], [48] but treats them as preliminaries — they provide important insights into many decisions in our system, but they do not count as the technical contribution of this paper. We next summarize these efforts and the gained insights.

A. Evaluating Bayes Error Estimators on Real-World Datasets

Evaluating the relative performance of BER estimators on real-world dataset is far from trivial. In one of our previous endeavor [47], we proposed FeeBee, a novel framework for evaluating BER estimators on real-world data. The key insights for building such a framework lies in the realization that evaluating BER estimators on a single point for tasks with unknown true BER is infeasible. Rather, one has to construct a series of points, for which the evolution of the BER is known. We do so by injecting uniformly distributed noise over the labels for different amounts of label noise and following the evolution of the BER through the following theoretical result.

Lemma 2.1 (From [47]): Let Y_{ρ} be a random variable defined on \mathcal{Y} by setting $Y_{\rho} = Z \cdot U(\mathcal{Y}) + (1-Z) \cdot Y$, where U



Fig. 1. Overview of Snoopy: Given user's input in terms of a dataset and a target, the system consults various Bayes error estimators (yellow ticks), aggregates them by taking the minimum (green tick) and outputs its belief whether the target is realistic, together with the insights in terms of convergence plots, gap towards target accuracy, and estimated number of additional samples needed to reach the target accuracy.

is a uniform variable taking values in \mathcal{Y} , and Z is a Bernoulli variable with probability $0 \le \rho \le 1$, both independent of X and Y. Then $R_{X,Y_{\rho}}^* = R_{X,Y}^* + \rho(1 - 1/C - R_{X,Y}^*)$.

The above lemma is sufficient in determining the strength of each BER estimator [47]. However, in order to apply a BER estimator in a system for a feasibility study in realworld datasets, where human annotators typically introduce more noise for classes which are harder to distinguish, we need to be able to go beyond uniform noise. Therefore, in Section III-A we provide a generalization that does not assume uniform label noise.

As a major finding of FeeBee, we established that the *INN*-based estimator is a powerful one — on par or better than all other estimators when it comes to performance, whilst being highly scalable and insensitive to hyper-parameters [47]. In this paper, *INN*-based estimator on top of a feature transformation is our default choice. For a fixed transformation f, and n-samples, it is defined by

$$\widehat{R}_{f(X),n} = \frac{(R_{f(X)})_{n,1}}{1 + \sqrt{1 - \frac{C(R_{f(X)})_{n,1}}{C - 1}}}.$$
(2)

B. On Convergence of Nearest Neighbor Classifiers over Feature Transformations

Since our estimator will combine the 1NN algorithm with pre-trained feature transformations, also called *embeddings*, such as those publicly available, we need to understand the influence of such transformations. In our theoretical companion to this paper [48], we provide a novel study of the behavior of a kNN classifier on top of a feature transformation, in particular its convergence rates on *transformed* data, previously known only for the *raw* data.¹ We prove the following theorem, recalling that a real-valued function g is *L*-Lipschitz if $|g(x) - g(x')| \leq L||x - x'||$, for all x, x', and defining $\mathcal{L}_g(f) := \mathbb{E}_X[(g \circ f)(X) - \eta(X)]^2$.

Theorem 2.2 (From [48]): Let $\mathcal{X} \subseteq \mathbb{R}^D$ and $\widetilde{\mathcal{X}} \subseteq \mathbb{R}^d$ be bounded sets, and let (X, Y) be a random vector taking values in $\mathcal{X} \times \{0, 1\}$. Let $g: \widetilde{\mathcal{X}} \to \mathbb{R}$ be an L_g -Lipschitz function. Then for all transformations $f: \mathcal{X} \to \widetilde{\mathcal{X}}$,

$$\mathbb{E}_{n}\left[(R_{f(X)})_{n,k}\right] - R_{X}^{*}$$

$$= \mathcal{O}\left(\frac{1}{\sqrt{k}}\right) + \mathcal{O}\left(L_{g}\left(\frac{k}{n}\right)^{1/d}\right) + \mathcal{O}\left(\sqrt[4]{\mathcal{L}_{g}(f)}\right). \quad (3)$$

Motivated by the usual architecture of trained embeddings, in Theorem 2.2 one should think of g as a softmax prediction layer with weights w, which allows taking $L_g = ||w||^2$. Equation 3 shows that there is a trade-off between the improved convergence rates (in terms of L_g and d) and the bias introduced by the transformation independent of kNN.

III. DESIGN OF SNOOPY

We next present the design of Snoopy. A high-level overview of the workflow of our system is given in Figure 1.

Functionality: Snoopy interacts with users in a simple way. The user provides an input dataset that is representative for the classification task at hand, along with a target accuracy α_{target} . The system then estimates the "highest possible accuracy" that an ML model can achieve, and outputs a binary signal — REALISTIC, if the system deducts that this target accuracy is achievable; UNREALISTIC, otherwise. We note that Snoopy does not provide a model that can achieve that target, only its *belief* on whether the target is achievable, using an *inexpensive* process. Furthermore, the goal of Snoopy is not to provide a perfect answer on feasibility, but to give information that can guide and help with the decision-making process — the signal provided by the system may as well be wrong, as we will discuss later. The best possible accuracy is implicitly returned to the user in the form of the gap between target and projected accuracy (c.f., Section IV-C).

<u>Interaction Model</u>: The binary signal of Snoopy given to a user is often correct, but not always. We now dive into the user's and Snoopy's interaction upon receiving the signal.

The Case When Snoopy Reports REALISTIC: In general, one should trust the system's output when it reports the target to be realistic, and proceed with running AutoML. We note that wrongly reporting realistic can be a very costly mistake which any feasibility system should try to avoid. In theory, our system could also be wrong in that fashion, due to (1) a lower bound estimate based on the 1NN estimator by Cover and Hart [38] that is known to be not always tight, or (2) the fact that the estimators are predicting asymptotic values. However, as presented in the next section, we construct our

¹We restrict ourselves to C = 2, as usual in theoretical results about the convergence rates of a kNN classifier.

estimator in a theoretically justified way that aims at reducing such mistakes and in our experiments we do not observe this behavior. Even if this were the case, we expect (2) to be the dominating reason, in which case gathering more data for the task at hand and running AutoML on this larger dataset might very well confirm the system's prediction.

The Case When Snoopy Reports UNREALISTIC: In this case, our experiments showed that the system's output is also trustworthy, but with more caution. Under reasonable computational resources², Snoopy is often correct in preventing unrealistic expectations for a varying amount of both synthetic and natural label noise. Nonetheless, there are two possible reasons for making wrong predictions in this manner: (1) either the data is not representative enough for the task (i.e., users might need to acquire more data), or (2) the transformations applied in order to reduce the feature dimension, or to bring raw features into a numerical format in the first place (e.g., from text), increased the BER.³ We note that (1) and (2) are complementary to each other. Even though estimating the BER on raw features (if applicable) prevents (2) from happening, having "better" transformations can lower the number of samples required to accurately estimate the BER. In an ideal world, one could rule out (1) by checking whether the BER estimator converged on the given number of samples. That is why Snoopy provides insights in terms of convergence plots and finite-sample extrapolation numbers to help users understand the relation between increasing number of samples and BER estimate, giving insights into the source of predicting UNREALISTIC, and increase the confidence in the prediction.

A. Data Quality Issues and the BER

The power of the BER, and the reason that Snoopy focuses on estimating this quantity, is that it provides a link connecting data quality to the performance of (best possible) ML models. This link can be made more explicit, even in closed form, if we assume some noise model. We take one of the most prominent source of data quality issues, label noise, as an example, and illustrate this connection via a novel theoretical analysis.

<u>Noise Model</u>: We focus on a standard noise model: classdependent label noise [52]. We assume that we are given a noisy random variable Y_{ρ} through a transition matrix t with

$$t_{\widetilde{y},y} := \mathbb{P}(Y_{\rho} = \widetilde{y} \mid Y = y, X = x) = \mathbb{P}(Y_{\rho} = \widetilde{y} \mid Y = y),$$
(4)

where the equality follows from the assumption that we are in the class-dependent label noise scenario, rather than in instance-dependent one. One can think of $\rho(y) = 1 - t_{y,y}$ to be the fraction of class y that gets flipped. Let $y_x := \arg \max_{y \in \mathcal{Y}} p(Y = y|x)$. We further assume that $y_x = \arg \max_{y \in \mathcal{Y}} p_{\rho}(Y_{\rho} = y|x)$, meaning that the maximal label per sample x is preserved after flipping, albeit possibly with lower probability (which then increases the BER). Our main result is the following theorem. *Theorem 3.1:* Let Y_{ρ} be a random variable taking values in \mathcal{Y} that satisfies (4). Then

$$R_{X,Y_{\rho}}^{*} = R_{X,Y}^{*} + \mathbb{E}_{X}[\rho(y_{x})p(y_{x}|x)] - \sum_{y \neq y_{x}} \mathbb{E}_{X}\left[t_{y_{x},y}p(y|x)\right]$$

One can prove Theorem 3.1 using the law of total expectation, together with careful manoeuvring of the terms that involve the elements of the transition matrix.⁴ Setting $\rho(y) = \rho \cdot (1-1/C)$, for all $y \in \mathcal{Y}$, and $t_{y,y'} = \rho/C$, for all $y' \neq y$, recovers Lemma 2.1, and one can further deduct the following valid bounds on the evolution of the BER:

$$(1 - s_{X,Y}) \min_{y} \rho(y) - s_{X,Y} \max_{y,y':y \neq y'} t_{y,y'} \le R^*_{X,Y_{\rho}} \le s_{X,Y} + \max_{y} \rho(y),$$

where $s_{X,Y}$ denotes the error of state-of-the-art model.

Other Data Quality Dimensions: Whilst we assume that the BER for zero label noise is typically small, it does not have to be equal to zero (c.f., [39] for examples). Nonetheless, by estimating the BER, we implicitly quantify the data quality issues along *all* dimensions (e.g., missing features, or combinations of feature and label noise). Deriving alternative noise models to theoretically and empirically disentangling these factors is a challenging task and left as future work.

IV. IMPLEMENTATION

The core component of Snoopy is a BER estimator, which estimates the irreducible error of a given task. The key design decision of Snoopy is to *consult a collection of BER estimators and aggregate them in a meaningful way*. More precisely, for a collection of feature transformations \mathcal{F} , e.g., publicly available pre-trained feature transformations (or lastlayer representations of pre-trained neural networks) on platforms like TensorFlow Hub, PyTorch Hub, and HuggingFace Transformers, we define our main estimator of the BER (on *n* samples) using Equation 2 by

$$\hat{R} = \min_{f \in \mathcal{T}} \hat{R}_{f(X),n}.$$

Finally, the system's output is

REALISTIC, if
$$\widehat{R} \leq 1 - \alpha_{\text{target}}$$
,
UNREALISTIC, otherwise.

A. "Just a Lightweight AutoML System?"

At first glance, our system might seem like a "lightweight AutoML system," which runs a collection of fast models (e.g., kNN classifiers) and takes the minimum to get the best possible classifier accuracy. We emphasize the difference the accuracy of an AutoML system always corresponds to a concrete ML model that can achieve this accuracy; however, a BER estimator does *not* provide this concrete model. That is, Snoopy does *not* construct a model that can achieve \hat{R} . This key difference between AutoML and feasibility study makes the latter inherently more computationally efficient, with almost instantaneous re-running, which we will further illustrate with experiments in Section VI-B.

²For instance, reproducing the state-of-the-art model performance for wellestablished benchmark datasets is often a highly non-trivial task.

³We have shown in our theoretical companion [48] that any deterministic transformation can only increase the BER.

⁴We provide the full proof and additional discussion in the extended online version of this work [53].

B. Theoretical Analysis

Given a collection of 1NN-based BER estimators over feature transformations, Snoopy aggregates them by taking the minimum. This seemingly simple aggregation rule is far from trivial, raising obvious questions — Why can we aggregate BER estimators by taking the minimum? When will this estimator work well and when will it not?

In order to mathematically quantify different regimes, we need a few simple definitions. We define the *asymptotic tightness* of our estimator for a fixed transformation f as

$$\Delta_f = R^*_{f(X)} - \lim_{n \to \infty} \widehat{R}_{f(X),n}.$$
(5)

Equation 1 implies $\Delta_f \geq 0$. We define the corresponding *transformation bias* by

$$\delta_f = R_{f(X)}^* - R_X^*, \tag{6}$$

with $\delta_f \geq 0$ (by [48]). Finally, the *n*-sample gap (of the estimator) is given by

$$\gamma_{f,n} = \widehat{R}_{f(X),n} - \lim_{n \to \infty} \widehat{R}_{f(X),n}, \tag{7}$$

with $\gamma_{f,n} \ge 0$ in expectation (also by [48]).

The fundamental challenge lies in the fact that none of the three quantities above can be derived in practice: Δ_f is dependent on the underlying unknown distribution, δ_f is intractable for complex neural networks [48], and $\gamma_{f,n}$ relies on the convergence of the estimator, which requires the number of samples to be exponential in the input dimension [51], making it impossible to generalize to representations on real-world datasets. Nevertheless, the connection between the quantities, together with the empirical analysis from [47] and Section VI of this work, allows us to define meaningful regimes next.

<u>When is R optimal?</u> In other words, when does the transformation that yields the minimum outperform all the others? A sufficient condition is given by

$$\forall f \in \mathcal{F}: \quad \delta_f + \gamma_{f,n} - \Delta_f \ge 0. \tag{8}$$

If the sum of finite-sample gap and transformation bias (i.e., the normalized constants of the second and third terms in Equation 3) is larger than the asymptotic tightness of the estimator (i.e., the normalized constants of the first term in Equation 3) for all transformations, then all estimators yield a number larger than the true BER, and therefore the minimum can be taken. Intuitively, this means that all the curves in the convergence plot are above the true BER. We note that this has to include the identity transformation, where there is no transformation bias. If Condition 8 holds, R will not underestimate the BER. The above trivially holds if for all $f \in \mathcal{F}$ one has $\Delta_f = 0$. Note that any classifier accuracy (nonscaled, to be used as proxy) also trivially falls into this regime, although it is usually worse than R. Furthermore, the system is guaranteed to not predict YES when the target is unreachable, thus avoiding costly mistakes. If the system wrongly predicts UNREALISTIC, it is guaranteed that its predicted error is off by at most $\delta_f + \gamma_{f,n} - \Delta_f$.

We note that all empirical evidence in Section VI and in our companion work on BER evaluation framework [47] suggests that we are in this regime for reasonable label noise (e.g., less than 80%) on a wide range of datasets and transformations.



Fig. 2. Theoretical justifications: (Left) 1NN error and its estimator values for raw features and the best transformation. (Right) Scaling down the logistic regression error on the best transformation and normalizing it by plugging it into the 1NN estimator of Equations 2. Notice that the solid green and solid pink line are identical in the right plot, leading to a dark purple line.

What if R is not optimal? We distinguish two different cases in this regime. In the first one, we suppose that the suggested estimator $R_{X,\infty}$ of Cover and Hart [38] on the raw features performs well in the asymptotic regime, i.e. that Δ_{id} is small. In that case, a sufficient condition for \widehat{R} to be at least as good as $\widehat{R}_{X,\infty}$ is $\Delta_f \leq \Delta_{id}$, for all $f \in \mathcal{F}$. Intuitively, this states that if all transformations do not increase the asymptotic tightness of the estimator by transforming the underlying probability distribution with respect to the raw distribution, taking the minimum over all transformations is no worse than running the estimator with 1NN on infinite samples. This condition can be seen empirically by inspecting the linear shape of the 1NN-based BER estimator values with increasing label noise for different transformations (c.f., Figure 2 on the left). One could weaken the condition for the finite-sample regime, resulting in a sufficient condition for Rto perform better than $R_{X,\infty}$:

$$\forall f \in \mathcal{F}: \quad \delta_f + \gamma_{f,n} - \Delta_f + \Delta_{id} \ge 0. \tag{9}$$

For the second case, when \hat{R} performs poorly, we ask: What is the worst-case error of underestimation? Using the fact that 1NN error is trivially above the true BER (c.f., left inequality of Equation 1), we can bound the difference of the 1NN-Based estimator value. In fact, the estimator value is at most the scaling factor of Equation 2 away from the true BER (i.e., 1/2 for a binary classification problem). However, our analysis and empirical verification reveals that our estimator of choice rarely ends up in the worst-case scenario. In fact, \hat{R} is usually the optimal choice and, when it is not, we often end up in the regime in which $\hat{R}_{X,\infty}$ already performs well and \hat{R} outperforms it.

Downscaling classifiers other than 1NN: As the worstcase error holds for scaling down any classifier accuracy, one could be tempted to use a downscaled version (e.g., dividing by a constant c > 1, or by plugging the error value in place of the 1NN error into the estimator of Equation 2) of other classifiers as a proxy. Contrary to the 1NN-based estimators, it is easy to show that for many datasets, any scaled version of a proxy model accuracy quickly falls into this worst-case regime (c.f., Figure 2 on the right, or Figure 4b in [47]), supporting the challenges of the Strawman outlined in the introduction.

C. Additional Guidance

To support users of Snoopy in deciding whether to "trust" the output of the system, regardless of the outcome, additional information is provided. It comes in the form of (a) the estimated BER and, thus, the gap between the projected accuracy and the target accuracy, (b) the convergence plots indicating the estimated BER value with respect to increased number of training samples over all deployed BER estimators (as illustrated in Figure 1), and (c) an additional estimate of the required number of additional samples to reach the target accuracy for the minimal transformation. Such an estimate is fairly non-trivial. Although Snapp et. al. [51] suggest how to approximate the kNN error by fitting a parametrized function to sampled data, the number of samples required to attain high confidence and accuracy is exponential in the feature dimension. This method is thus not practical for either finitesample extrapolation, or estimating the BER, as shown in our companion work [47]. Instead, to support users of Snoopy beyond purely visual aids, we approximate the estimate based on the 1NN error using a simple log-linear function [54]

$$\log\left((R_X)_{n,k}\right) \approx -\alpha \log\left(n\right) + C,\tag{10}$$

for two positive constants α and C. The idea of approximating the error is motivated by recent observations of scaling laws across different deep learning modalities [55], [56]. Notice that Equation 10 should only be used to extrapolate the convergence for a small number of additional data points. The function (i.e., the exponential of the righthand side of Equation 10) is known to converge to 0, implying that regardless of the label noise or true BER, it will always underestimate the BER for a too large number of samples. We show the benefits and failures of using this approximation in Section VI-C.

V. SYSTEM OPTIMIZATIONS

The suggested and theoretically motivated estimator from the previous section relies on the 1NN classifier being evaluated on a possibly large collection of publicly available pretrained feature transformations. We present optimizations that improve the performance, making it more scalable.

Algorithm: There are five computational steps involved:

- Take user's dataset with n samples: (i) features $X_1, X_2, ..., X_n$ and labels $Y_1, Y_2, ..., Y_n$.
- (ii) For pre-defined m transformations $\mathcal{F} = \{f_1, f_2, \dots, f_m\},\$ calculate the corresponding features for every sample in the dataset by applying all the transformations in \mathcal{F} .
- (iii) For each feature transformation $j \in [m]$, calculate the 1NN classifier error $R_j = (R_{f_j(X)})_{n,1}$ on the transformed features $f_i(X_i)$, for all samples $i \in [n]$.
- (iv) Based on the 1NN classifier error, derive the lower-bound estimates $\widehat{R}_{f_j(X),n}$ using Equation 2. (v) Report the overall estimate $\widehat{R} = \min_{j \in [m]} \widehat{R}_{f_j(X),n}$.

Note that the dataset is split into training samples and test samples. The test set is only used to estimate the accuracy of the classifier and is typically orders of magnitude smaller than the training set. The quality of Snoopy depends heavily on



Fig. 3. Improved Successive-Halving: At each point for each convergence curve we construct a tangent and check whether there are more than half of the remaining curves that are better than the tangent.

the list of feature transformations that are fed into it. Since we take the minimum over all transformations in \mathcal{F} , increasing the size of the set only *improves* the estimator. On the downside, an efficient implementation is by no means trivial with an everincreasing number of (publicly) available transformations.

Computational Bottleneck: When analyzing the previously defined algorithm, we realize that the major computational bottleneck comes from transforming the features. Especially when having large pre-trained networks as feature extractors, running inference on large datasets, in order to get the embeddings, can be very time-consuming and result in running times orders of magnitude larger than the sole computation of the 1NN classifier accuracy. More concretely, given a dataset with n samples and m feature transformations, the worst case complexity is $\mathcal{O}(mn)$, which highlights the importance of providing an efficient version of the algorithm.

Multi-armed Bandit Approach: Inspired by ideas for efficient implementation of the nearest-neighbor search on hardware accelerators [57], running inference on all the training data for all feature transformations simultaneously is not necessary. Rather, we define a streamed version of our algorithm by splitting the steps (ii) to (iv) into iterations of fixed batch size per transformation. This new formulation can directly be mapped to a non-stochastic best arm identification problem, where each arm represents a transformation. The successivehalving algorithm [49], which is invoked as a subroutine inside the popular Hyperband algorithm [7], is designed to solve this problem efficiently. We can summarize the idea of successivehalving as follows: Uniformly allocate a fixed initial budget across all transformations and evaluate their performance. Keep only the better half of the transformations, and repeat this until a single transformation remains.

Improved Successive-Halving: We develop a variant of successive-halving that further improves the performance. The main idea comes from observing the convergence curve of a kNN classifier. We know that under some mild assumptions, the kNN error decreases as a function of $n^{-2/d}$, where n is the number of samples [58]. Therefore, we can assume that the convergence curve is decreasing and convex. This allows us to predict a simple lower bound for the convergence curve at the end of each step - using the tangent through the curve at the last known point, as illustrated in Figure 3. If the tangent line at the end point is worse than half of the remaining curves at the current point, the curve will not proceed to the next round. To simplify the implementation, we approximate the tangent by a line through the two last-

TABLE I DATASETS AND SOTA PERFORMANCES.

Name	Classes C	Train / Test Samples	SOTA %
MNIST	10	60K / 10K	0.16 [59]
CIFAR10	10	50K / 10K	0.63 [60]
CIFAR100	100	50K / 10K	6.49 [60]
IMDB	2	25K / 25K	3.79 [61]
SST2	2	67K / 872	3.2 [61]
YELP	5	500K / 50K	27.80 [61]

known values of the convergence curve and develop a variant of successive-halving that uses this as a stopping condition. An important property of our improvement is that the remaining transformations after each step are the same as the ones from the original successive-halving, which implies that all theoretical guarantees of successive-halving still hold.

Parameters of Successive-Halving: We eliminate the dependency on the initial budget by implementing the *doubling-trick* (cf. Section 3 in [49]). The batch size of the iterations has a direct impact on the performance and speedup of the algorithm. This is linked to properties of the underlying hardware and the fact that approximating the tangent for points that are further apart becomes less accurate. Hence, we treat the batch size as a single hyper-parameter, which we tune for all transformations and datasets.

Efficient Incremental Execution: For the specific scenario of incrementally cleaning labels until a target accuracy is reachable, we provide a simple yet effective optimization that enables re-running Snoopy almost instantly. After its initial execution, the system keeps track of the label of a single sample per test point – its nearest neighbor. As cleaning labels of test or training samples does not change the nearest neighbor, calculating the 1NN accuracy after cleaning any training or test samples can be performed by iterating over the test set exactly once, thus, providing real-time feedback.

VI. EXPERIMENTS

We now present the results of our empirical evaluation by describing the benefits of performing a feasibility study in general, and using the binary output of Snoopy over other baselines. We focus on a specific use-case scenario motivated in the introduction. We also show how the additional guidance can increase trust in the binary signal. We then analyze the generalization properties of our system on certain vision tasks and conclude this section by performing a detailed performance analysis of Snoopy. The code of Snoopy is available via https://github.com/easeml/snoopy, whereas the code to reproduce the results can be found under https: //github.com/DS3Lab/snoopy-paper.

A. Experimental Setup

<u>Datasets</u>: We perform the evaluation on two data modalities that are ubiquitous in modern machine learning and are accompanied by strong state-of-the-art (SOTA) performances summarized in Table I. Implicitly, a strong SOTA yields a low natural BER (i.e., originating from all data quality dimensions). The first group consists of visual classification tasks,

TABLE II CIFAR-N datasets statistics. The variable $t_{i,j}$ refers to an element of the noise transition matrix t (c.f., Section III-A).

Dataset	Noise	$\max_i t_{i,i}$	$\min_i t_{i,i}$	$\max_{i \neq j} t_{i,j}$
CIFAR10-Aggre	9%	17%	3%	10%
CIFAR10-Random1	17%	26%	10%	23%
CIFAR10-Random2	18%	26%	10%	23%
CIFAR10-Random3	18%	26%	10%	23%
CIFAR100-Noisy	40%	85%	8%	31%

including CIFAR10 [62], CIFAR100 [63], and MNIST [64]. The second group consists of standard text classification tasks, where we focus on IMDB, SST2, and YELP [61]. We remark that the SOTA values for SST2 and YELP are provided on slightly different sizes of training sets.

We mainly focus our study on datasets with noisy labels. The ML community usually works on high-quality, noise-free benchmark datasets. As an exception, Wei et. al. [52] published different noisy variants of the popular CIFAR datasets, called CIFAR-N. The noise levels vary between 10% and 40% (c.f., Table II). The datasets are provided with their noise transition matrix, allowing us to use the bounds derived from Theorem 3.1. The assumption therein corresponds to the each diagonal element being the maximal value per row, which is given for all datasets. Supported by our theoretical understanding of the impact of label noise on the BER and its evolution, we also synthetically inject uniform label noise for 20% and 40% of the label into all six datasets from Table I.

Feature Transformations: We compile a wide range of more than 15 different feature transformations per data modality, such as PCA and NCA [65], as well as state-of-the-art pre-trained embeddings. The pre-trained feature transformations are taken from public sources such as TensorFlow Hub, PyTorch Hub, and HuggingFace, whereas PCA and NCA are taken from scikit-learn⁵. The pre-trained embeddings can either be directly accessed via the corresponding source, or have to be extracted from the last-layer representations of pre-trained neural networks. More details about the transformations supported, for each modality individually, can be found in the extended version of this work [53].

<u>Settings of Snoopy</u>: When running Snoopy, we define the time needed to reach the lowest 1NN error across all embeddings based on multiple independent runs as described in Section VI-F. These runtimes include the 1NN computation and running inference on a single GPU, with the latter being the most costly part, particularly for large NLP models. In the end-to-end experiments, when re-running Snoopy after having restored a fixed portion of the synthetically polluted labels (set to 1% of the dataset size), we use the fact that the "best" embedding did not change and, therefore, no additional inference needs to be executed.

We compare with a diverse set of baselines that estimate the BER: (i) training a logistic regression (LR) model on top of all pre-trained transformations, (ii) running AutoKeras, and (iii)

⁵TensorFlow Hub: https://tfhub.dev, PyTorch Hub: https://pytorch.org/ hub, HuggingFace: https://huggingface.co/models/ and scikit-learn: https:// scikit-learn.org/



Fig. 4. Error Estimations vs. Time on three synthetics noise levels. The dashed horizontal lines represent the expected increase of the SOTA using Lemma 2.1.

fine-tuning a state-of-the-art (SOTA) foundation model [66] for each data modality.

B. Evaluation of BER Estimations

<u>Baseline 1: LR Models</u>: As mentioned before, when training the logistic regression models we assume that the representations for all the training and test samples are calculated in advance exactly once. In the end-to-end experiments, after having restored the same fixed portion of labels (i.e., 1% of test and train samples), re-training the LR models does not require any inference. We train all LR models on a single GPU using SGD with a momentum of 0.9, a mini-batch size of 64 and 20 epochs. We select the minimal test accuracy achieved over all combinations of learning rate in $\{0.001, 0.01, 0.1\}$ and L_2 regularization values in $\{0.0, 0.001, 0.01\}$. We calculate the average time needed to train a LR based on the best transformation, without label noise, on all possible hyperparameters. The hyper-parameter search was conducted 5 independent times for any value of randomly injected noise.

Baseline 2: AutoML Systems: To mimic the use of an AutoML systems on a single GPU without any prior datasetdependent knowledge, we run AutoKeras with the standard parameters of a maximum of 100 epochs and 2 trials on top of all datasets. We additionally run auto-sklearn with two different configurations to simulate a short execution time (max 1 hour), and a longer execution time (max 10 hours). Auto-sklearn does not natively support text as input and we therefore execute it using universal sentence embedding representations omitting the time to extract those representations. We report the mean of 5 independent executions in terms of times and accuracy, noting little variance amongst the results.

Baseline 3: Finetune: The goal of this baseline is to replicate the SOTA values achieved for all datasets. We remark that this baseline is equipped with a strong prior knowledge which is usually unavailable for performing a cheap feasibility study and it only serves as a reference point. Unfortunately, reproducing the exact SOTA values was not possible for any of the dataset involved in the study, which is mainly due to computational constraints and the lack of publicly available reproducible code. We therefore perform our best, mostly manual, efforts to train a model on the original non-corrupted data. For multi-channel vision tasks (i.e., CIFAR10 and CIFAR100, and its noisy variants), we fine-tune EfficientNet-B4 using the proposed set of hyper-parameters [67], whereas for NLP tasks, we fine-tune BERT-Base with 3 different learning rates and for 3 epochs [68], using a maximal sequence length of 512, batch size of 6 and the Adam optimizer.

We first evaluate Snoopy by comparing its BER estimation of the best achievable accuracy with other baselines and show how this benefits an end-to-end scenario.

Snoopy vs. Baselines on Synthetic Noise: In Figure 4 we present our main findings on three levels of realistically injected label noise - 0%, 20% and 40%, which we visualize by adding the increase of the SOTA (at the time of writing) in expectation as horizontal lines to indicate a proxy of the ground truth BER error. We see that Snoopy is comparable to the short execution of auto-sklean whilst producing much better estimations. Furthermore, Snoopy is much faster than all other methods, often by orders of magnitude. The only exception is YELP in which running over large models (e.g., GPT2 or XLNET) slows down Snoopy in a fashion comparable to AutoKeras, whilst still outperforming it in terms of the estimated accuracy. It also produces BER estimations that are *comparable*, *if not better* than all other approaches. In fact, it is often better than both LR and, particularly, AutoKeras. It is only slightly worse than the LR classifier on text tasks for IMDB and SST2, while being orders of magnitude faster.

<u>Snoopy vs. Baselines on Real Noise</u>: In Figure 5 we run the same set of experiments for real noisy variants of CIFAR10 and CIFAR100 from [52]. We realize that Snoopy constantly outperforms all baselines both in terms of speed and estimation accuracy. When comparing the error values to the lower and upper bounds, we realize that whilst Snoopy remains inside the bounds, there is a considerable gap between them. Nonetheless, Snoopy produces estimates close to the expected increase of the SOTA using Theorem 3.1.

Is Taking the Minimum Necessary? When analyzing the performance of the system with respect to the number of feature transformations, one might ask the question whether a single transformation always outperforms all the others and hence makes the selection of the minimal estimator obsolete. When conducting our experiments, we observed that selecting the *wrong* embedding can lead towards a large gap when compared to the optimal embedding, e.g., favoring the embedding USELARGE over XLNET on IMDB doubles the gap of the estimated BER to the known SOTA value [47], whereas favoring XLNET over USELARGE on SST2 increases the gap by $1.5 \times$, making proper selection necessary (c.f., Figure 6).

C. Usefulness of the Additional Guidance

When evaluating the 1NN estimator accuracy for varying label noise, and its convergence under different feature trans-



Fig. 5. Error Estimations vs. Time on multiple real noisy variant of CIFAR10 and CIFAR100 (c.f., Table II). The dashed horizontal lines represent bounds and the expected increase of the SOTA using Theorem 3.1.



Fig. 8. CIFAR100 - (Left) Extrapolation accuracy to the full dataset based on a fraction of the samples. (Right) Approximation based on 5% of the samples.

formations, we see that even the best transformations are constantly over-estimating the lower bound when increasing label noise, validating the key arguments for taking the minimum over all estimators. All the results indicate the median, 95% and 5% quantiles over multiple independent runs (i.e., 10 for YELP and 30 otherwise). We observe the presence of much more instability in SST2 when compared to other datasets. This is not at all surprising since SST2 has a very small test set consisting of less than one thousand samples, as seen in Table I. This naturally results in higher variance and less confidence for the 1NN classifier accuracy compared to the larger number of test samples for the other datasets.

Figure 7 illustrates a convergence plot for a fixed embedding

(EfficientNet-B5) and the clean CIFAR100 dataset injected with 20% and 40% label noise respectively. The two target accuracies visualized by a horizontal line represent exactly the noise level, and the noise level plus 10%. Note that the noise level is only reachable if the original BER of the dataset is zero. From the visualizations, the target of 0.5 on the dataset with 0.4 noise is highly likely. By using the approximation from Equation 10, we realize that less than 10K more samples should suffice to attain this accuracy. Conversely, for a target of 0.3 and 0.2 noise, although possibly realizable, the number of additional samples to verify the quality of the extrapolation is already more than 260K. Note that Equation 10 converges to zero, and therefore any target can be realizable. Targeting exactly the label noise for each of the datasets yields an extrapolated number of more than 16M and 84M, which both should be seen as not trustworthy approximations based on the much smaller number of available samples in the training set. This thus implies that the target accuracy is not achievable based on the given transformations and numbers of samples. To illustrate this fact, we subsample the low label noise dataset and the same embedding for a fixed fraction. We then extrapolate the achievable target for the full dataset (i.e., 50K training samples) and plot the difference between the extrapolated target and the true BER estimate in Figure 8 on the left. The right part of Figure 8 illustrates the extrapolation based on 5% of the samples. Notice that this provided example illustrates when to trust the estimated number of additional samples required using Equation 10 (i.e., when the number if relatively low), not the BER estimate of Snoopy. The same results can easily be shown for any other dataset.

D. End-to-end Use Case

How can we take advantage of Snoopy to help practical use cases? In this section, we focus on a specific end-to-end use case of a feasibility study in which the user's task contains a target accuracy and a representative, but noisy dataset. The goal is to reach the target accuracy. At each step, the user can perform one of the following three actions: (1) clean a portion of the labels, (2) train a high-accuracy model using AutoKeras or fine-tune a state-of-the-art pre-trained model, (3) perform a feasibility study by either using the cheap LR model or Snoopy. To simulate the cleaning process on a noisy dataset, which usually requires human interactions of an expert labeler, we focus on the manually polluted datasets with synthetic label noise, where we can simply *restore* the original label



from the dataset. Being aware of different human costs for cleaning labels in real-world scenarios (i.e., depending on the application and the required expertise), we compare the impact of different cost scenarios outlined below. We report the mean (accuracy and run-time) over at least 5 independent runs.

Different User Interaction Models: We differentiate two main scenarios in our end-to-end experimental evaluation: (1) without feasibility study and (2) with feasibility study. Without a feasibility study, users will start an expensive, high-accuracy run (i.e., running the fine-tuning baseline) using the input data. If the achieved accuracy is below the desired target, users will clean a fixed portion of the data (1%, 5%, 10%, or 50%, which we call steps) and re-run the expensive training system. This is repeated until a model reaches the desired accuracy or all labels are cleaned. With a feasibility study, users alternate between running the feasibility study system and cleaning a portion of the data (set to 1% of the data) until the feasibility study returns a positive signal or all labels are cleaned. Finally, a single expensive training run is performed. The lower bound on computation is given by training the expensive model exactly once.

<u>Different Cost Scenarios</u>: We measure the cost in hypothetical "dollar price" for different regimes, depending on the *human-labeling* costs and on the *machine* costs. For the former, we define two scenarios: 'free', 'cheap' (0.002 dollars per label, resulting in 500 labels per dollar) and 'expensive' (0.02 dollars per label, resulting in 50 labels per dollar). For the latter, we fix the *machine* cost to 0.9\$ per hour (the current cost of a single GPU Amazon EC2 instance).

<u>Key Findings</u>: We only present the results on CIFAR100 here and leave the rest to the extended online version of this work [53] —we observe similar results on all datasets for a wide range of initial noise levels and target accuracies. We show the results in Figures 9 and 10, for 2 different cost setups described above (cheap and expensive), each over 2 values of the initial noise (0.40, 0.20) and, respectively, 2 target accuracies (0.60, 0.80). Each dot represents the result of one run of the expensive training process. More results for the third cost scenario and other datasets can be found in [53].

(I) Feasibility Study Helps.: When comparing the costs of repetitively training an expensive model to those of using an efficient and accurate system that performs a feasibility study, such as Snoopy, we see significant improvements across all results (c.f., blue vs. brown lines in Figure 9). Without a system that performs a feasibility study, users are facing a dilemma. On the one hand, if one does not train an expensive model frequently enough, it might clean up more labels than necessary, to achieve the target accuracy, e.g., FineTune (step 50%), which makes it intense on the humanlabeling costs. This can be seen by the size of the vertical gap between the end point of a method and the horizontal line indicating the minimum number of samples to be cleaned to achieve the target accuracy. On the other hand, if one trains an expensive model too frequently, visible in the number of stairs for expensive fine-tune lines or the steepness of the curves for faster methods, it becomes computationally expensive, wasting a lot of computation time. With a feasibility study, the user can balance these two factors better. As running lowcost proxy models is significantly cheaper than training an expensive model, the user can get feedback more frequently (having in mind the efficient incremental implementation from Section V). Finally, notice that when we enter the label-cost dominated regime (e.g., Figure 10), one seeks at cleaning the minimum amount of labels necessary, ignoring the computational costs. Nevertheless, finding the right step size is critical, making it a difficult task.

(II) Snoopy Outperforms Baselines.: When comparing different estimators that can be used in a feasibility study, in most cases, Snoopy is more effective compared to running a cheaper model such as LR, with its accuracy as a proxy. Snoopy offers significant savings compared to LR when the labeling costs are high. The LR model will often be of a lower accuracy than an expensive approach; hence, it requires to clean more labels than necessary to reach the target. We note that there are cases (e.g., for IMDB) where the best LR model yields a lower error than the BER estimator used by Snoopy. In such cases, there exists a regime where the costs of using the LR proxy are comparable or superior to using Snoopy despite being more expensive to compute. However, we see this as an exception and Figure 9 clearly show that the LR proxy is usually significantly more costly than using Snoopy.

E. Generalization to Other Tasks

In this section we examine two potential limitations of Snoopy, when deployed on a new task: (i) *its dependence on large datasets*, and (ii) *the necessity of having "good" pre-trained feature transformations for the given task.* For this, we use the results of [15] on the popular visual task adaptation benchmark (VTAB) [69] which is known to be a



Fig. 11. Snoopy vs. Maximal fine-tune accuracies on VTAB (Left) for proprietary expert models derived from [15] and (**Right**) for public models from Huggingface: Showing Snoopy's generalization ability on small datasets and embeddings trained on different tasks.



Fig. 12. Different selection strategies for (Left) CIFAR10, (Right) CIFAR100.

diverse collection of datasets (19 different tasks), each being *small* (1K training samples), and our collection of pre-trained transformations *does not contain* any trained on these datasets. Additionally, we fine-tune the same 19 datasets on a set of 235 publicly available PyTorch models from Huggingface.

To validate that Snoopy does not suffer from the above limitations, in Figure 11 we illustrate the difference between Snoopy's predictions and the best achieved post-fine-tune accuracies. We observe that on most datasets, Snoopy produces a useful estimate of the fine-tune accuracy (except for some negative transfer results enabled by the low data regime) for both proprietary expert models and publicly available models. The estimates of the later are slightly shifted to the right as expected. Even though this is sufficient to say that the currently available embeddings are supporting Snoopy's performance, we expect this figure to improve over time as more and better embeddings become publicly available via repositories such as Huggingface, which also start including learned representation for different modalities such as tabular data.

F. Efficiency of Snoopy

We saw that the gain of using Snoopy comes from having an (i) *efficient estimator* of (ii) *high accuracy*. Those two requirements are naturally connected. While having access to more and "better" (pre-trained) transformations is key for getting a high accuracy of our estimator, it requires the implementation of our algorithm to scale with respect to the ever-increasing number of transformations.

<u>Runtime Analysis:</u> To showcase the importance of the successive-halving (SH) algorithm, with and without the tangent method presented in Section V, we compare different strategies for deploying the 1NN estimator in Figures 12. The strategies are evaluated with respect to the runtime (averaged across multiple independent runs on a single Nvidia Titan



Fig. 13. Incremental execution (log scale): On all datasets, rerunning Snoopy incrementally is several orders of magnitude faster than running from scratch.

Xp GPU) needed to reach an estimation within 1% of the best possible value using all the training samples. Running the estimator only on the transformation yielding the minimal result is referred to as the *perfect* strategy providing a lower bound, whereas we also test the uniform allocation baseline described in [49]. We report the runtime by selecting the best batch size out of 1%, 2%, or 5% of the training samples. We observe that running the entire feasibility study using Snoopy on CIFAR100 on a single GPU takes slightly more than 16 minutes, whereas the largest examined NLP dataset YELP requires almost 8.5 hours, with a clear improvement of SH with the tangent method over the one without. Putting these numbers into context, fine-tuning EfficientNet-B4 on CIFAR100 on the same GPU with one set of hyper-parameters (out of the 56 suggested by the authors [67]) requires almost 10 hours (without knowing whether other embeddings would perform better), whereas training large NLP models usually requires several hundred accelerators [70].

<u>Incremental Execution</u>: In Figure 13, not surprisingly, we see the benefit of Snoopy's ability to quickly rerun incremental evaluation described in Section V. Compared to running from scratch, we see that incremental execution is faster by several orders of magnitude on all datasets.

VII. CONCLUSION

We present Snoopy, a novel system that enables a systematic feasibility study for ML application development. By consulting a range of estimators of the Bayes error and aggregating them in a theoretically justified way, Snoopy suggests whether a predefined target accuracy is achievable. We demonstrate system optimizations that support the usability of Snoopy, and scale with the increase in the number and diversity of available pre-trained embeddings in the future.

Acknowledgements. CZ and the DS3Lab gratefully acknowledge the support from the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number MB22.00036 (for European Research Council (ERC) Starting Grant TRIDENT 101042665), the Swiss National Science Foundation (Project Number 200021_184628, and 197485), Innosuisse/SNF BRIDGE Discovery (Project Number 40B2-0_187132), European Union Horizon 2020 Research and Innovation Programme (DAPHNE, 957407), Botnar Research Centre for Child Health, Swiss Data Science Center, Alibaba, Cisco, eBay, Google Focused Research Awards, Kuaishou Inc., Oracle Labs, Zurich Insurance, and the Department of Computer Science at ETH Zurich.

REFERENCES

- H. Van Vliet, H. Van Vliet, and J. Van Vliet, Software engineering: principles and practice. John Wiley & Sons, 2008, vol. 13.
- [2] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B. Su, "Scaling distributed machine learning with the parameter server," in OSDI, 2014, pp. 583–598.
- [3] X. Meng, J. K. Bradley, B. Yavuz, E. R. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. B. Tsai, M. Amde, S. Owen, D. Xin, R. Xin, M. J. Franklin, R. Zadeh, M. Zaharia, and A. Talwalkar, "MLlib: Machine learning in Apache Spark," *Journal of Machine Learning Research*, vol. 17, pp. 34:1–34:7, 2016.
- [4] M. Zaharia, A. Chen, A. Davidson, A. Ghodsi, S. A. Hong, A. Konwinski, S. Murching, T. Nykodym, P. Ogilvie, M. Parkhe *et al.*, "Accelerating the machine learning lifecycle with mlflow." *IEEE Data Eng. Bull.*, vol. 41, no. 4, pp. 39–45, 2018.
- [5] J. Bergstra, D. Yamins, and D. D. Cox, "Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms," in *Proceedings of the 12th Python in Science Conference*, vol. 13. Citeseer, 2013, p. 20.
- [6] D. Baylor, E. Breck, H.-T. Cheng, N. Fiedel, C. Y. Foo, Z. Haque, S. Haykal, M. Ispir, V. Jain, L. Koc et al., "Tfx: A tensorflow-based production-scale machine learning platform," in *Proceedings of the 23rd* ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2017, pp. 1387–1395.
- [7] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: A novel bandit-based approach to hyperparameter optimization," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6765–6816, 2017.
- [8] M. Vartak, H. Subramanyam, W.-E. Lee, S. Viswanathan, S. Husnoo, S. Madden, and M. Zaharia, "Modeldb: a system for machine learning model management," in *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*, 2016, pp. 1–3.
- [9] T. Kraska, "Northstar: An interactive data science system," *PVLDB*, vol. 11, no. 12, pp. 2150–2164, 2018.
- [10] N. Polyzotis, M. Zinkevich, S. Roy, E. Breck, and S. Whang, "Data validation for machine learning," *Proceedings of Machine Learning and Systems*, vol. 1, 2019.
- [11] S. Nakandala, A. Kumar, and Y. Papakonstantinou, "Incremental and approximate inference for faster occlusion-based deep cnn explanations," in *Proceedings of the 2019 International Conference on Management of Data*, 2019, pp. 1589–1606.
- [12] S. Nakandala, Y. Zhang, and A. Kumar, "Cerebro: a data system for optimized deep learning model selection," *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 2159–2173, 2020.
- [13] F. A. Hubis, W. Wu, and C. Zhang, "Quantitative overfitting management for human-in-the-loop ML application development with ease. ml/meter," arXiv preprint arXiv:1906.00299, 2019.
- [14] C. Renggli, B. Karlas, B. Ding, F. Liu, K. Schawinski, W. Wu, and C. Zhang, "Continuous integration of machine learning models with ease.ml/ci: Towards a rigorous yet practical treatment," in *SysML Conference*, 2019.
- [15] C. Renggli, A. S. Pinto, L. Rimanic, J. Puigcerver, C. Riquelme, C. Zhang, and M. Lucic, "Which model to transfer? finding the needle in the growing haystack," *IEEE/CVF Conference on Computer Vision* and Pattern Recognition, 2022.
- [16] G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma, "Improving data quality: Consistency and accuracy." in *VLDB*, vol. 7, 2007, pp. 315–326.
- [17] F. Chiang and R. J. Miller, "Discovering data quality rules," *Proceedings of the VLDB Endowment*, vol. 1, no. 1, pp. 1166–1177, 2008.
- [18] S. Sadiq, N. K. Yeganeh, and M. Indulska, "20 years of data quality research: themes, trends and synergies," in *Proceedings of the Twenty-Second Australasian Database Conference-Volume 115*, 2011, pp. 153– 162.
- [19] W. Fan, "Data quality: From theory to practice," Acm Sigmod Record, vol. 44, no. 3, pp. 7–18, 2015.
- [20] Z. Abedjan, X. Chu, D. Deng, R. C. Fernandez, I. F. Ilyas, M. Ouzzani, P. Papotti, M. Stonebraker, and N. Tang, "Detecting data errors: Where are we and what needs to be done?" *Proceedings of the VLDB Endowment*, vol. 9, no. 12, pp. 993–1004, 2016.
- [21] S. Schelter, D. Lange, P. Schmidt, M. Celikel, F. Biessmann, and A. Grafberger, "Automating large-scale data quality verification," *Proceedings of the VLDB Endowment*, vol. 11, no. 12, pp. 1781–1794, 2018.
- [22] I. F. Ilyas and X. Chu, Data cleaning. Morgan & Claypool, 2019.

- [23] R. Y. Wang and D. M. Strong, "Beyond accuracy: What data quality means to data consumers," *Journal of management information systems*, vol. 12, no. 4, pp. 5–33, 1996.
- [24] D. M. Strong, Y. W. Lee, and R. Y. Wang, "Data quality in context," *Communications of the ACM*, vol. 40, no. 5, 1997.
- [25] M. Scannapieco and T. Catarci, "Data quality under a computer science perspective," *Archivi & Computer*, vol. 2, 2002.
- [26] C. Batini, C. Cappiello, C. Francalanci, and A. Maurino, "Methodologies for data quality assessment and improvement," ACM computing surveys, vol. 41, no. 3, 2009.
- [27] S. Krishnan, J. Wang, E. Wu, M. J. Franklin, and K. Goldberg, "Active-Clean: Interactive Data Cleaning for Statistical Modeling," *Proceedings* of the VLDB Endowment, vol. 9, no. 12, 2016.
- [28] A. Ghorbani and J. Zou, "Data shapley: Equitable valuation of data for machine learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2242–2251.
- [29] W. Wu, L. Flokas, E. Wu, and J. Wang, "Complaint-driven training data debugging for query 2.0," in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, 2020, pp. 1317–1334.
- [30] P. Li, X. Rao, J. Blase, Y. Zhang, X. Chu, and C. Zhang, "CleanML: A Benchmark for Joint Data Cleaning and Machine Learning [Experiments and Analysis]," arXiv preprint arXiv:1904.09483, 2019.
- [31] R. Jia, D. Dao, B. Wang, F. A. Hubis, N. Hynes, N. M. Gürel, B. Li, C. Zhang, D. Song, and C. J. Spanos, "Towards efficient data valuation based on the shapley value," in *The 22nd International Conference on Artificial Intelligence and Statistics*, 2019, pp. 1167–1176.
- [32] R. Jia, D. Dao, B. Wang, F. A. Hubis, N. M. Gurel, B. Li, C. Zhang, C. Spanos, and D. Song, "Efficient task-specific data valuation for nearest neighbor algorithms," *Proceedings of the VLDB Endowment*, vol. 12, no. 11, pp. 1610–1623, 2019.
- [33] B. Karlaš, P. Li, R. Wu, N. M. Gürel, X. Chu, W. Wu, and C. Zhang, "Nearest neighbor classifiers over incomplete information: From certain answers to certain predictions," arXiv preprint arXiv:2005.05117, 2020.
- [34] D. Compton, T. P. Love, and J. Sell, "Developing and assessing intercoder reliability in studies of group interaction," *Sociological Methodology*, vol. 42, no. 1, pp. 348–364, 2012.
- [35] U. Gadiraju, R. Kawase, S. Dietze, and G. Demartini, "Understanding malicious behavior in crowdsourcing platforms: The case of online surveys," in *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, 2015, pp. 1631–1640.
- [36] A. Checco, J. Bates, and G. Demartini, "All that glitters is gold—an attack scheme on gold questions in crowdsourcing," in *Proceedings* of the AAAI Conference on Human Computation and Crowdsourcing, vol. 6, no. 1, 2018.
- [37] D. Q. Sun, H. Kotek, C. Klein, M. Gupta, W. Li, and J. D. Williams, "Improving human-labeled data through dynamic automatic conflict resolution," arXiv preprint arXiv:2012.04169, 2020.
- [38] T. M. Cover and P. A. Hart, "Nearest neighbor pattern classification," *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 21–27, 1967.
- [39] C. Renggli, L. Rimanic, N. M. Gürel, B. Karlaš, W. Wu, and C. Zhang, "A data quality-driven view of mlops," *Data Engineering*, p. 11, 2021.
- [40] K. Fukunaga and L. Hostetler, "k-nearest-neighbor Bayes-risk estimation," *IEEE Transactions on Information Theory*, vol. 21, no. 3, pp. 285–293, 1975.
- [41] P. A. Devijver, "A multiclass, k-NN approach to Bayes risk estimation," *Pattern recognition letters*, vol. 3, no. 1, pp. 1–6, 1985.
- [42] K. Fukunaga and D. M. Hummels, "Bayes error estimation using parzen and k-NN procedures," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 9, no. 5, pp. 634–643, May 1987.
- [43] L. J. Buturovic and M. Z. Markovic, "Improving k-nearest neighbor bayes error estimates," in 11th IAPR International Conference on Pattern Recognition. Vol. II. Conference B: Pattern Recognition Methodology and Systems, vol. 1. IEEE Computer Society, 1992, pp. 470–471.
- [44] T. Pham-Gia, N. Turkkan, and A. Bekker, "Bounds for the Bayes error in classification: A Bayesian approach using discriminant analysis," *Statistical Methods & Applications*, vol. 16, no. 1, pp. 7–26, Jun. 2007.
- [45] V. Berisha, A. Wisler, A. O. Hero, and A. Spanias, "Empirically estimable classification bounds based on a nonparametric divergence measure," *IEEE Transactions on Signal Processing*, vol. 64, no. 3, pp. 580–591, 2016.
- [46] S. Y. Sekeh, B. L. Oselio, and A. O. Hero, "Learning to bound the multiclass Bayes error," *IEEE Transactions on Signal Processing*, 2020.

- [47] C. Renggli, L. Rimanic, N. Hollenstein, and C. Zhang, "Evaluating bayes error estimators on read-world datasets with feebee," Advances in Neural Information Processing Systems (Datasets and Benchmarks), vol. 34, 2021.
- [48] L. Rimanic, C. Renggli, B. Li, and C. Zhang, "On convergence of nearest neighbor classifiers over feature transformations," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [49] K. Jamieson and A. Talwalkar, "Non-stochastic best arm identification and hyperparameter optimization," in *Artificial Intelligence and Statistics*, 2016, pp. 240–248.
- [50] K. Fukunaga and D. Kessell, "Nonparametric Bayes error estimation using unclassified samples," *IEEE Transactions on Information Theory*, vol. 19, no. 4, pp. 434–440, 1973.
- [51] R. R. Snapp and T. Xu, "Estimating the Bayes risk from sample data," in Advances in Neural Information Processing Systems, 1996, pp. 232–238.
- [52] J. Wei, Z. Zhu, H. Cheng, T. Liu, G. Niu, and Y. Liu, "Learning with noisy labels revisited: A study using real-world human annotations," in *International Conference on Learning Representations*, 2022.
- [53] C. Renggli, L. Rimanic, L. Kolar, W. Wu, and C. Zhang, "Automatic feasibility study via data quality analysis for ml: A case-study on label noise," arXiv preprint arXiv:2010.08410, 2022.
- [54] T. Hashimoto, "Model performance scaling with multiple data sources," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4107–4116.
- [55] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.
- [56] J. S. Rosenfeld, A. Rosenfeld, Y. Belinkov, and N. Shavit, "A constructive prediction of the generalization error across scales," in *International Conference on Learning Representations*, 2020.
- [57] J. Johnson, M. Douze, and H. Jégou, "Billion-scale similarity search with gpus," *IEEE Transactions on Big Data*, 2019.
- [58] R. R. Snapp, D. Psaltis, and S. S. Venkatesh, "Asymptotic slowing down of the nearest-neighbor classifier," in *Advances in Neural Information Processing Systems*, 1991, pp. 932–938.
- [59] A. Byerly, T. Kalganova, and I. Dear, "A branching and merging convolutional network with homogeneous filter capsules," *arXiv preprint* arXiv:2001.09136, 2020.
- [60] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby, "Large scale learning of general visual representations for transfer," arXiv preprint arXiv:1912.11370, 2019.
- [61] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "Xlnet: Generalized autoregressive pretraining for language understanding," in *Advances in Neural Information Processing Systems*, 2019, pp. 5754–5764.
- [62] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [63] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, "Sharpness-aware minimization for efficiently improving generalization," in *International Conference on Learning Representations*, 2020.
- [64] A. Byerly, T. Kalganova, and I. Dear, "No routing needed between capsules," *Neurocomputing*, 2021.
- [65] Z. Wu, A. A. Efros, and S. X. Yu, "Improving generalization via scalable neighborhood component analysis," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 685–701.
- [66] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021.
- [67] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.
- [68] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv* preprint arXiv:1810.04805, 2018.
- [69] X. Zhai, J. Puigcerver, A. Kolesnikov, P. Ruyssen, C. Riquelme, M. Lucic, J. Djolonga, A. S. Pinto, M. Neumann, A. Dosovitskiy *et al.*, "A large-scale study of representation learning with the visual task adaptation benchmark," *arXiv preprint arXiv:1910.04867*, 2019.
- [70] J. Rasley, S. Rajbhandari, O. Ruwase, and Y. He, "Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters," in *Proceedings of the 26th ACM SIGKDD International*

Conference on Knowledge Discovery & Data Mining, 2020, pp. 3505–3506.