
Pearls from Pebbles: Improved Confidence Functions for Auto-labeling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Auto-labeling techniques produce labeled data with minimal manual annotations
2 using the representations from self-supervised models and confidence scores. A
3 popular technique, threshold-based auto-labeling (TBAL) trains model using these
4 representations and manual annotations, and assigns model’s prediction as label
5 to the points where model’s confidence score is greater than certain threshold.
6 However, the model’s scores can be overconfident and lead to poor performance.
7 We show that, calibration, a common remedy for the overconfidence problem, falls
8 short in tackling this problem for TBAL. Thus, instead of using existing calibration
9 methods, we introduce a framework for optimal confidence functions for TBAL
10 and develop Colander, a method designed to maximize auto-labeling performance.
11 We perform an extensive empirical evaluation of Colander and other confidence
12 functions, using representations from CLIP and text embedding models for image
13 and text data respectively. We find Colander achieves up to 60% improvement
14 on coverage (the proportion of points labeled by model) over the baselines while
15 maintaining error level below 5% and using the same amount of labeled data.

16 1 Introduction

17 The demand for labeled data in machine learning (ML) is perpetual. Threshold-based auto-labeling
18 (TBAL) is a promising solution to obtain high-quality labeled data at low cost [41, 37, 49] using
19 model’s confidence scores and self-supervision. It powers industry products like Amazon SageMaker
20 Ground Truth [41]. The confidence function is critical to the TBAL workflow. Existing TBAL systems
21 rely on common choices like softmax outputs from neural networks [37, 49]. These functions *are*
22 *not well aligned with the objective of the auto-labeling system*. Using them results in substantially
23 suboptimal coverage (Figure 1(a)). For this reason, we ask:

24

What are the right choices of confidence functions for TBAL and how can we obtain them?

25 An ideal confidence function for auto-labeling will achieve the maximum coverage at a given auto-
26 labeling error tolerance and thus will bring down the labeling cost significantly. Finding such an ideal
27 function, however, is difficult because of the *inherent tension* between accuracy and coverage.

28 *Overconfidence* further stymies hopes of balancing accuracy and coverage. While overconfidence
29 is a challenge in general, it is exacerbated in TBAL: since models are trained on a small amount of
30 labeled data, making the problem of designing confidence functions even more challenging. Figure
31 1(a) shows that softmax scores are overconfident, resulting in poor auto-labeling performance.

32 Several methods have been introduced to address overconfidence [10]. Applying these miss out on
33 significant performance (Figure 1(b)), since the calibration goal differs from auto-labeling. From
34 the auto-labeling standpoint, we seek minimum overlap between the correct and incorrect model
35 prediction scores.

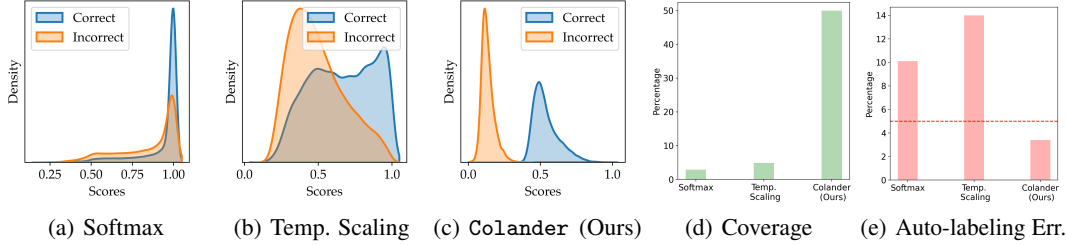


Figure 1: Scores distributions (Kernel Density Estimates) of a CNN model trained on CIFAR-10 data. (a) softmax scores of the vanilla training procedure (SGD) (b) scores after post-hoc calibration using temperature scaling and (c) scores from our Colander procedure applied on the same model. For training the CNN model we use 4000 points drawn randomly and 1000 validation points (of which 500 are used for Temp. Scaling and Colander). The test accuracy of the model is 55%. Figures (d) and (e) show the coverage and auto-labeling error of these methods. The dotted-red line corresponds to a 5% error threshold.

36 We tackle these challenges by *proposing a framework to learn suitable confidence functions* for
 37 TBAL. We summarize our contributions as follows,

- 38 1. We propose a principled framework to study the choices of confidence functions suitable for auto-
 39 labeling and provide a practical method (Colander) to learn confidence functions for efficient
 40 and reliable auto-labeling.
- 41 2. We systematically study commonly used choices of scoring functions and calibration methods
 42 and demonstrate that they lead to poor auto-labeling performance.
- 43 3. Through extensive empirical evaluation on real-world datasets, we show that using the confidence
 44 scores obtained using our procedure boosts auto-labeling performance significantly in comparison
 45 to common choices of confidence functions and calibration methods.

46 2 Background and Motivation

47 **Notation.** Let $[m] := \{1, 2, \dots, m\}$ for any natural number m . Let X_u be a set of unlabeled points
 48 drawn from some instance (feature) space \mathcal{X} . This could either be the space of raw features or the
 49 representation space from some self-supervised model. Let $\mathcal{Y} = \{1, \dots, k\}$ be the label space. There
 50 is an unknown ground truth labeling function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$. Let \mathcal{O} be a *noiseless* oracle that provides
 51 the true label for any point $\mathbf{x} \in \mathcal{X}$. Denote the model (hypothesis) class by \mathcal{H} , where each $h \in \mathcal{H}$
 52 is a function $h : \mathcal{X} \rightarrow \mathcal{Y}$. Each classifier h also has an associated *confidence function* $g : \mathcal{X} \rightarrow \Delta^k$
 53 that quantifies the confidence of the prediction by model $h \in \mathcal{H}$ on any data point $\mathbf{x} \in \mathcal{X}$. Here,
 54 Δ^k is a $(k - 1)$ -dimensional probability simplex. Let $\mathbf{v}[i]$ denote the i^{th} component for any vector
 55 $\mathbf{v} \in \mathbb{R}^d$. For any point $\mathbf{x} \in \mathcal{X}$ the prediction is $\hat{y} := h(\mathbf{x})$ and the associated confidence is $g(\mathbf{x})[\hat{y}]$.
 56 The vector \mathbf{t} denotes scores over k -classes, and $\mathbf{t}[y]$ denotes its y^{th} entry, i.e., score for class y . Table
 57 2 contains a summary of the notation.

58 **Threshold-based Auto-labeling.** Threshold-based auto-labeling (TBAL) seeks to obtain labeled
 59 datasets while reducing the labeling burden on humans. The input is a pool of unlabeled data X_u .
 60 It outputs, for each $\mathbf{x} \in X_u$, a label $\tilde{y} \in \mathcal{Y}$. The output label could be either y , from the oracle
 61 (representing a human-obtained label), or \hat{y} , from the model. Let N_u be the number of unlabeled
 62 points, $A \subseteq [N_u]$ the set of indices of auto-labeled points, and $X_u(A)$ be these points. Let N_a denote
 63 the size of the auto-labeled set A . The *auto-labeling error*, denoted by $\hat{\mathcal{E}}(X_u(A))$, and the *coverage*,
 64 denoted by $\hat{\mathcal{P}}(X_u(A))$, are defined as follows:

$$\hat{\mathcal{E}}(X_u(A)) := \frac{1}{N_a} \sum_{i \in A} \mathbb{1}(\tilde{y}_i \neq f^*(\mathbf{x}_i)), \quad \text{and} \quad \hat{\mathcal{P}}(X_u(A)) := |A|/N_u = N_a/N_u. \quad (1)$$

65 The goal of an auto-labeling algorithm is to label the dataset so that $\hat{\mathcal{E}}(X_u(A)) \leq \epsilon_a$ while maximizing
 66 coverage $\hat{\mathcal{P}}(X_u(A))$ for any given $\epsilon_a \in [0, 1]$. The TBAL algorithm proceeds iteratively. In each
 67 iteration, it queries labels for a subset of unlabeled points from the oracle. It trains a classifier from
 68 the model class \mathcal{H} on the oracle-labeled data acquired till that iteration. It then uses the model's
 69 confidence scores on the validation data to identify the region in the instance space, where the current
 70 classifier is confidently accurate and automatically labels the points in this region.

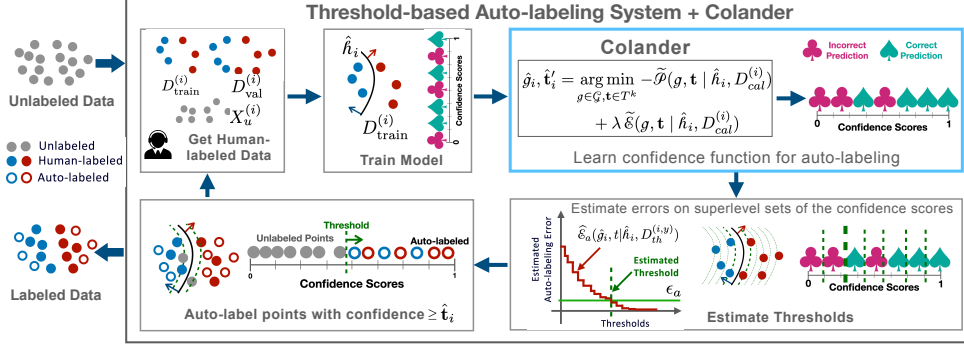


Figure 2: Threshold-based Auto-labeling with Colander: takes unlabeled data as input, selects a small subset of data points, and obtains human labels for them to create $D_{train}^{(i)}$ and $D_{val}^{(i)}$ for the i th iteration. Trains model \hat{h}_i on $D_{train}^{(i)}$. In contrast to the standard TBAL procedure, here we randomly split $D_{val}^{(i)}$ into two parts, $D_{cal}^{(i)}$ and $D_{th}^{(i)}$. Colander kicks in, takes \hat{h}_i and $D_{cal}^{(i)}$ as input and learns a coverage maximizing confidence function \hat{g}_i for \hat{h}_i . Using $D_{th}^{(i)}$ and \hat{g}_i , auto-labeling thresholds \hat{t}_i are determined to ensure the auto-labeled data has error at most ϵ_a . After obtaining the thresholds the rest of the steps are the same as standard TBAL. The whole workflow runs until all the data is labeled or another stopping criterion is met.

71 3 Proposed Method (Colander)

72 3.1 Auto-labeling optimization framework

73 In any iteration of TBAL, we have a model h trained on a subset of data labeled by the oracle. This
 74 model may not be highly accurate. However, it could be accurate in some regions of the instance
 75 space, and with the help of a confidence function g , we want to identify the points where the model is
 76 correct and auto-label them. As we saw earlier, arbitrary choices of g perform poorly on this task.
 77 Instead, we propose a framework to find the right function from a sufficiently rich family.

78 **Optimization problem.** Let $\sigma(\alpha, z) := 1/(1 + \exp(-\alpha z))$ denote the sigmoid function on \mathbb{R}
 79 with scale parameter $\alpha \in \mathbb{R}$. It is easy to see that, for any g, y and \mathbf{t} , $g(\mathbf{x})[y] \geq \mathbf{t}[y] \iff$
 80 $\sigma(\alpha, g(\mathbf{x})[y] - \mathbf{t}[y]) \geq 1/2$. Using this fact, we define the following surrogates of the auto-labeling
 81 error and coverage:

$$\tilde{\mathcal{P}}(g, \mathbf{t} | h, D_{cal}) := \frac{1}{|D_{cal}|} \sum_{(\mathbf{x}, y) \in D_{cal}} \sigma(\alpha, g(\mathbf{x})[y] - \mathbf{t}[y]), \quad (2)$$

$$\tilde{\mathcal{E}}(g, \mathbf{t} | h, D_{cal}) := \frac{\sum_{(\mathbf{x}, y) \in D_{cal}} \mathbf{1}(y \neq \hat{y}) \sigma(\alpha, g(\mathbf{x})[y] - \mathbf{t}[y])}{\sum_{(\mathbf{x}, y) \in D_{cal}} \sigma(\alpha, g(\mathbf{x})[y] - \mathbf{t}[y])}, \quad (3)$$

82 and the surrogate optimization problem as follows,

$$\arg \min_{g \in \mathcal{G}, \mathbf{t} \in T^k} -\tilde{\mathcal{P}}(g, \mathbf{t} | h, D_{cal}) + \lambda \tilde{\mathcal{E}}(g, \mathbf{t} | h, D_{cal}) \quad (\text{P1})$$

84 Here, $\lambda \in \mathbb{R}^+$ is the penalty term controlling the relative importance of the auto-labeling error and
 85 coverage. We tune it with the procedure discussed in Appendix B.4. The gap between the surrogate
 86 and actual coverage diminishes as $\alpha \rightarrow \infty$. We discuss this in the Appendix. Our framework is
 87 flexible with respect to the choice of \mathcal{G} . We use 2-layer neural nets for \mathcal{G} and use representations
 88 from the last two layers of \hat{h} as input for g . See Appendix for further details.

89 **Solving the surrogate optimization.** The optimization problem (P1) is non-convex. Nevertheless,
 90 it is differentiable and we can apply gradient-based methods. We solve for g and \mathbf{t} simultaneously
 91 using Adam [19]. Appendix B.4 lists the detailed procedure of our method in a TBAL system.

92 4 Empirical Evaluation

93 We conduct extensive empirical evaluation with variety of train-time and post-hoc calibration methods,
 94 and feature choices. We use *Vanilla* training with SGD [1, 3, 10]. *Squentropy* [16], *Correctness*

Train-time	Post-hoc	MNIST		CIFAR-10		20 Newsgroups		Tiny-ImageNet	
		Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)
Vanilla	Softmax	4.1\pm0.7	85.0 \pm 2.5	4.8 \pm 0.2	14.0 \pm 2.1	6.0 \pm 0.6	48.2 \pm 1.6	11.1 \pm 0.3	32.6 \pm 0.5
	TS	7.8 \pm 0.6	94.2 \pm 0.5	7.3 \pm 0.3	23.2 \pm 0.7	9.7 \pm 0.6	60.7 \pm 2.3	16.3 \pm 0.5	37.4 \pm 1.5
	Dirichlet	7.9 \pm 0.7	93.2 \pm 2.2	7.7 \pm 0.5	22.4 \pm 1.2	9.4 \pm 0.9	59.4 \pm 1.8	17.1 \pm 0.4	33.3 \pm 2.0
	SB	6.7 \pm 0.5	92.6 \pm 1.5	6.1 \pm 0.4	18.6 \pm 1.1	8.1 \pm 0.6	58.1 \pm 1.8	15.7 \pm 0.6	35.4 \pm 1.2
	Top-HB	7.4 \pm 1.4	93.1 \pm 3.6	6.0 \pm 0.7	15.6 \pm 1.9	9.2 \pm 1.0	59.0 \pm 2.0	16.6 \pm 0.5	37.6 \pm 2.2
	Ours	4.2 \pm 1.5	95.6\pm1.4	3.0\pm0.2	78.5\pm0.2	2.5\pm1.1	80.6\pm0.7	1.4\pm2.1	59.2\pm0.8
CRL	Softmax	4.7 \pm 0.4	86.0 \pm 4.5	5.2 \pm 0.3	15.9 \pm 0.8	5.8 \pm 0.5	48.3 \pm 0.3	10.4 \pm 0.4	32.5 \pm 0.6
	TS	8.0 \pm 0.8	94.8 \pm 0.8	6.8 \pm 0.8	20.3 \pm 1.1	9.5 \pm 1.0	61.7 \pm 1.6	15.8 \pm 0.6	37.4 \pm 1.7
	Dirichlet	8.6 \pm 0.6	93.1 \pm 1.6	7.7 \pm 0.2	20.9 \pm 1.1	8.7 \pm 0.9	58.0 \pm 1.4	16.3 \pm 0.4	33.1 \pm 1.9
	SB	7.4 \pm 0.8	93.1 \pm 2.7	5.9 \pm 0.9	17.9 \pm 1.5	8.9 \pm 1.1	57.9 \pm 3.9	15.0 \pm 0.4	35.5 \pm 1.2
	Top-HB	7.7 \pm 0.8	94.1 \pm 1.5	4.4 \pm 0.5	12.3 \pm 0.4	8.8 \pm 1.0	58.8 \pm 2.7	16.5 \pm 0.5	38.9 \pm 1.6
	Ours	4.5\pm1.4	95.6\pm1.3	2.2\pm0.6	77.9\pm0.2	1.8\pm1.2	81.3\pm0.5	2.8\pm2.1	61.2\pm1.4
FMFP	Softmax	4.8 \pm 0.8	84.2 \pm 4.1	4.9 \pm 0.4	15.6 \pm 1.7	5.4 \pm 0.7	45.4 \pm 1.9	10.5 \pm 0.3	32.4 \pm 1.4
	TS	8.0 \pm 0.6	95.3 \pm 1.6	6.5 \pm 0.3	21.0 \pm 1.5	9.5 \pm 0.5	57.7 \pm 2.2	16.2 \pm 1.1	37.7 \pm 1.8
	Dirichlet	8.2 \pm 1.3	94.0 \pm 2.2	6.9 \pm 0.4	21.7 \pm 1.2	8.9 \pm 1.0	56.6 \pm 2.4	17.4 \pm 0.8	33.0 \pm 1.8
	SB	7.2 \pm 1.1	93.1 \pm 2.3	6.1 \pm 0.5	19.5 \pm 1.0	8.6 \pm 0.4	55.8 \pm 1.3	15.5 \pm 0.6	36.1 \pm 0.5
	Top-HB	7.1 \pm 0.6	93.3 \pm 4.9	5.2 \pm 0.5	14.2 \pm 2.4	9.0 \pm 0.7	57.9 \pm 2.4	16.2 \pm 0.4	37.4 \pm 1.1
	Ours	4.6\pm0.8	95.7\pm0.2	3.0\pm0.4	77.4\pm0.2	2.5\pm0.9	80.8\pm0.6	1.8\pm2.0	60.8\pm1.4
Squentropy	Softmax	3.7\pm1.0	88.2 \pm 3.9	5.2 \pm 0.5	21.2 \pm 1.8	4.6 \pm 0.4	52.0 \pm 1.2	7.8 \pm 0.3	36.2 \pm 0.8
	TS	6.2 \pm 1.1	95.6 \pm 0.9	6.9 \pm 0.6	28.2 \pm 2.5	8.3 \pm 0.6	66.6 \pm 1.4	13.3 \pm 0.1	44.9 \pm 1.0
	Dirichlet	6.5 \pm 1.2	95.9 \pm 0.8	7.3 \pm 0.3	29.4 \pm 1.1	7.8 \pm 0.6	64.0 \pm 1.3	14.1 \pm 0.3	42.5 \pm 0.7
	SB	6.0 \pm 0.8	95.3 \pm 1.2	6.2 \pm 0.4	23.8 \pm 1.9	7.8 \pm 0.7	63.0 \pm 2.9	13.0 \pm 0.5	45.2 \pm 2.0
	Top-HB	5.3 \pm 0.4	96.4 \pm 0.9	4.3 \pm 0.5	15.8 \pm 1.4	8.2 \pm 0.8	66.5 \pm 2.2	13.7 \pm 0.1	45.9 \pm 1.4
	Ours	4.1 \pm 0.8	97.2\pm0.5	2.3\pm0.5	79.0\pm0.3	3.3\pm0.8	82.9\pm0.4	0.6\pm0.2	66.5\pm0.7

Table 1: In every round the error was enforced to be below 5%; ‘TS’ stands for Temperature Scaling, ‘SB’ stands for Scaling Binning, ‘Top-HB’ stands for Top-Label Histogram Binning. The column Err stands for auto-labeling error and Cov stands for the coverage. Each cell value is mean \pm std. deviation observed on 5 repeated runs with different random seeds.

95 *Ranking Loss (CRL)* [30] and *FMFP* [53] as train-time methods. We pipe these with *Colander*
96 and other post-hoc methods *Temperature scaling* [10], *Top-Label Histogram-Binning* [12], *Scaling-*
97 *Binning* [24] and *Dirichlet Calibration* [22]. We use raw features for *MNIST* [26] and *CIFAR-10*
98 [20], CLIP [38] embeddings for *Tiny-ImageNet* and FlagEmbedding [50] for *20 Newsgroups* [29]
99 The results are reported in Table 1.

100 TBAL with *Colander* consistently achieves significantly higher coverage compared to vanilla
101 training and softmax scores, especially on more complex datasets like *Tiny-ImageNet*, outperforming
102 baseline models across all data settings. Post-hoc calibration techniques slightly improve coverage but
103 at the cost of higher error, as their goal of mitigating overconfidence is not aligned with TBAL’s needs.
104 Importantly, *Colander* is compatible with different train-time methods and amplifies performance
105 gains, particularly when paired with *Squentropy*, where coverage increases by 6-7% and error
106 decreases, outperforming other train-time approaches. In contrast, methods focused on ordinal
107 ranking objectives, such as *CRL* and *FMFP*, perform poorly in the TBAL setting due to challenges
108 like limited training data and reduced effectiveness in differentiating between correct and incorrect
109 predictions when training error reaches zero after a few rounds.

110 5 Conclusion

111 We studied confidence scoring functions used in threshold-based auto-labeling (TBAL) and showed
112 that the commonly used choices and calibration methods can often be a bottleneck, leading to poor
113 performance. We proposed *Colander* to learn confidence functions that are aligned with the TBAL
114 objective. We evaluated *Colander* extensively against common baselines on several real-world
115 datasets and found that it improves the performance of TBAL significantly. A limitation of *Colander*
116 is that, similar to other post-hoc methods it also requires validation data to learn the scores. Reducing
117 (or eliminating) this dependence could be an interesting future work.

References

- 118
- 119 [1] S.-i. Amari. Backpropagation and stochastic gradient descent method. *Neurocomputing*, 5(4-5):
120 185–196, 1993.
- 121 [2] Y. Bai, S. Mei, H. Wang, and C. Xiong. Don’t just blame over-parametrization for over-
122 confidence: Theoretical analysis of calibration in binary classification. In M. Meila and
123 T. Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*,
124 volume 139 of *Proceedings of Machine Learning Research*, pages 566–576. PMLR, 18–24 Jul
125 2021.
- 126 [3] L. Bottou. *Stochastic Gradient Descent Tricks*, pages 421–436. Springer Berlin Heidelberg,
127 2012. ISBN 978-3-642-35289-8.
- 128 [4] Y. Chen, R. K. Vinayak, and B. Hassibi. Crowdsourced clustering via active querying: Practical
129 algorithm with theoretical guarantees. In *Proceedings of the AAAI Conference on Human
130 Computation and Crowdsourcing*, volume 11, pages 27–37, 2023.
- 131 [5] C. Corbière, N. THOME, A. Bar-Hen, M. Cord, and P. Pérez. Addressing failure prediction by
132 learning model confidence. In *Advances in Neural Information Processing Systems 32*, pages
133 2902–2913. 2019.
- 134 [6] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently
135 improving generalization. In *International Conference on Learning Representations*, 2021.
- 136 [7] D. Y. Fu, M. F. Chen, F. Sala, S. M. Hooper, K. Fatahalian, and C. Ré. Fast and three-rious:
137 Speeding up weak supervision with triplet methods. In *Proceedings of the 37th International
138 Conference on Machine Learning (ICML 2020)*, 2020.
- 139 [8] J. Gawlikowski, C. R. N. Tassi, M. Ali, J. Lee, M. Humt, J. Feng, A. Kruspe, R. Triebel,
140 P. Jung, R. Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint
141 arXiv:2107.03342*, 2021.
- 142 [9] R. G. Gomes, P. Welinder, A. Krause, and P. Perona. Crowdclustering. In *Advances in Neural
143 Information Processing Systems 24*, pages 558–566. 2011.
- 144 [10] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In
145 *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- 146 [11] C. Gupta and A. Ramdas. Distribution-free calibration guarantees for histogram binning without
147 sample splitting. In *International Conference on Machine Learning*, pages 3942–3952. PMLR,
148 2021.
- 149 [12] C. Gupta and A. Ramdas. Top-label calibration and multiclass-to-binary reductions. In
150 *International Conference on Learning Representations*, 2022.
- 151 [13] S. Hanson and L. Pratt. Comparing biases for minimal network construction with back-
152 propagation. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*,
153 volume 1. Morgan-Kaufmann, 1988.
- 154 [14] M. Hein, M. Andriushchenko, and J. Bitterwolf. Why relu networks yield high-confidence
155 predictions far away from the training data and how to mitigate the problem. 2018.
- 156 [15] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution
157 examples in neural networks. In *International Conference on Learning Representations*, 2017.
- 158 [16] L. Hui, M. Belkin, and S. Wright. Cut your losses with squentropy. In *Proceedings of the 40th
159 International Conference on Machine Learning*, pages 14114–14131, 2023.
- 160 [17] S. Hussain. Cifar 10- cnn using pytorch, Jul 2021. URL [https://www.kaggle.com/code/
161 shadabhussain/cifar-10-cnn-using-pytorch](https://www.kaggle.com/code/shadabhussain/cifar-10-cnn-using-pytorch).
- 162 [18] D. R. Karger, S. Oh, and D. Shah. Budget-optimal crowdsourcing using low-rank matrix
163 approximations. In *2011 49th Annual Allerton Conference on Communication, Control, and
164 Computing (Allerton)*, pages 284–291. IEEE, 2011.
- 165 [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint
166 arXiv:1412.6980*, 2014.
- 167 [20] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- 168 [21] A. Krogh and J. A. Hertz. A simple weight decay can improve generalization. In *Proceedings*
169 *of the 4th International Conference on Neural Information Processing Systems, NIPS'91*,
170 page 950–957, San Francisco, CA, USA, 1991. Morgan Kaufmann Publishers Inc. ISBN
171 1558602224.
- 172 [22] M. Kull, M. Perello Nieto, M. Kängsepp, T. Silva Filho, H. Song, and P. Flach. Beyond temper-
173 ature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. In
174 *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 175 [23] A. Kumar, S. Sarawagi, and U. Jain. Trainable calibration measures for neural networks from
176 kernel mean embeddings. In *Proceedings of the 35th International Conference on Machine*
177 *Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2805–2814. PMLR,
178 10–15 Jul 2018.
- 179 [24] A. Kumar, P. S. Liang, and T. Ma. Verified uncertainty calibration. *Advances in Neural*
180 *Information Processing Systems*, 32, 2019.
- 181 [25] Y. Le and X. Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- 182 [26] Y. LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- 183 [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document
184 recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 185 [28] A. Mazumdar and B. Saha. Clustering with noisy queries. In *Advances in Neural Information*
186 *Processing Systems*, volume 30, 2017.
- 187 [29] T. Mitchell. Twenty Newsgroups. UCI Machine Learning Repository, 1999.
- 188 [30] J. Moon, J. Kim, Y. Shin, and S. Hwang. Confidence-aware learning for deep neural networks.
189 In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages
190 7034–7044, 2020.
- 191 [31] J. Mukhoti, V. Kulharia, A. Sanyal, S. Golodetz, P. Torr, and P. Dokania. Calibrating deep
192 neural networks using focal loss. *Advances in Neural Information Processing Systems*, 33:
193 15288–15299, 2020.
- 194 [32] R. Müller, S. Kornblith, and G. E. Hinton. When does label smoothing help? *Advances in*
195 *neural information processing systems*, 32, 2019.
- 196 [33] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence
197 predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer*
198 *vision and pattern recognition*, pages 427–436, 2015.
- 199 [34] A. Niculescu-Mizil and R. Caruana. Predicting good probabilities with supervised learning.
200 In *Proceedings of the 22nd International Conference on Machine Learning, ICML '05*, page
201 625–632, 2005. ISBN 1595931805.
- 202 [35] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
203 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
204 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
205 *Learning Research*, 12:2825–2830, 2011.
- 206 [36] J. Platt. Probabilistic outputs for support vector machines and comparisons to regularized
207 likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- 208 [37] H. Qiu, K. Chintalapudi, and R. Govindan. MCAL: Minimum cost human-machine active
209 labeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- 210 [38] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
211 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision.
212 In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 213 [39] A. J. Ratner, C. M. D. Sa, S. Wu, D. Selsam, and C. Ré. Data programming: Creating large
214 training sets, quickly. In *Proceedings of the 29th Conference on Neural Information Processing*
215 *Systems (NIPS 2016)*, Barcelona, Spain, 2016.
- 216 [40] V. C. Raykar, S. Yu, L. H. Zhao, G. H. Valadez, C. Florin, L. Bogoni, and L. Moy. Learning
217 from crowds. *Journal of Machine Learning Research*, 11(43):1297–1322, 2010.
- 218 [41] SGT. Aws sagemaker ground truth. [https://aws.amazon.com/sagemaker/
219 data-labeling/](https://aws.amazon.com/sagemaker/data-labeling/), 2022. Accessed: 2024-05-22.

- 220 [42] C. Shin, W. Li, H. Vishwakarma, N. C. Roberts, and F. Sala. Universalizing weak supervision.
221 In *International Conference on Learning Representations*, 2022.
- 222 [43] A. Sorokin and D. Forsyth. Utility data annotation with amazon mechanical turk. In *2008 IEEE*
223 *Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages
224 1–8, 2008. doi: 10.1109/CVPRW.2008.4562953.
- 225 [44] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing
226 properties of neural networks. In *International Conference on Learning Representations*,
227 2014.
- 228 [45] R. K. Vinayak and B. Hassibi. Crowdsourced clustering: Querying edges vs triangles. In
229 *Proceedings of the 30th International Conference on Neural Information Processing Systems*,
230 NIPS’16, 2016.
- 231 [46] R. K. Vinayak, S. Oymak, and B. Hassibi. Graph clustering with missing data: Convex
232 algorithms and analysis. *Advances in Neural Information Processing Systems*, 27, 2014.
- 233 [47] R. K. Vinayak, T. Zrnic, and B. Hassibi. Tensor-based crowdsourced clustering via triangle
234 queries. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing*
235 *(ICASSP)*, pages 2322–2326. IEEE, 2017.
- 236 [48] H. Vishwakarma and F. Sala. Lifting weak supervision to structured prediction. In *Advances in*
237 *Neural Information Processing Systems*, 2022.
- 238 [49] H. Vishwakarma, H. Lin, F. Sala, and R. K. Vinayak. Promises and pitfalls of threshold-based
239 auto-labeling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- 240 [50] S. Xiao, Z. Liu, P. Zhang, and N. Muennighoff. C-pack: Packaged resources to advance general
241 chinese embedding, 2023.
- 242 [51] B. Zadrozny and C. Elkan. Learning and making decisions when costs and probabilities are both
243 unknown. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge*
244 *discovery and data mining*, pages 204–213, 2001.
- 245 [52] B. Zadrozny and C. Elkan. Transforming classifier scores into accurate multiclass probability
246 estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge*
247 *discovery and data mining*, pages 694–699, 2002.
- 248 [53] F. Zhu, Z. Cheng, X.-Y. Zhang, and C.-L. Liu. Rethinking confidence calibration for failure
249 prediction. In *European Conference on Computer Vision*, pages 518–536. Springer, 2022.

250 **Supplementary Material Organization**

251 The supplementary material is organized as follows. We provide deferred details of background and
252 motivation section in Appendix A of the method in Appendix B. Then, in Appendix C, we provide
253 additional experimental results and details of the experiment protocol and hyperparameters used for
254 the experiments. Our code with instructions to run, is uploaded along with the paper.

255 **Contents**

256	1 Introduction	1
257	2 Background and Motivation	2
258	3 Proposed Method (Colander)	3
259	3.1 Auto-labeling optimization framework	3
260	4 Empirical Evaluation	3
261	5 Conclusion	4
262	A Appendix to Section 2	9
263	A.1 Detailed Comparison with Active Learning and Self Training	9
264	A.2 Details of the motivating experiment in Section 2	9
265	B Appendix on Our Method	10
266	B.1 Detailed Algorithms	10
267	B.2 Tightness of surrogates.	10
268	B.3 Active Querying Strategy.	10
269	B.4 TBAL procedure with Colander	10
270	B.5 Glossary	12
271	C Additional Experiments and Details	15
272	C.1 Experiments on N_t , N_v and ν	15
273	C.2 Experiments on Colander input	16
274	C.3 Experiments on ϵ_a	17
275	C.4 Experiments on multiple rounds	17
276	C.5 Experiments on different architectures	19
277	C.6 Hyperparameters	19
278	C.7 Train-time and post-hoc methods	19
279	C.7.1 Train-time methods	19
280	C.7.2 Post-hoc methods	20
281	C.8 Compute resources	20
282	C.9 Detailed dataset and model	21
283	C.10 Detailed experiments protocol	21

287 A Appendix to Section 2

288 A.1 Detailed Comparison with Active Learning and Self Training

289 To illustrate the differences between TBAL and the combination of Active Learning (AL) and Self-
 290 Training for the task of data labeling, we run an experiment on the 2 concentric circles data setting as
 291 used in [49]. The details are as follows:

292 **Data setting.** We generate two concentric circles with points in the outer circle belonging to one
 293 class and the inner circle belonging to the other class. The total number of points generated is 10,000
 294 of which we use 2000 for validation.

295 **Methods.** We run TBAL, AL+Self-
 296 Training, and AL+Self-Training+SC,
 297 using logistic regression. The combina-
 298 tion of AL+Self-Training means,
 299 in each iteration, the algorithm
 300 queries human-labeled data points and
 301 pseudo-labels the points in the un-
 302 labeled data using self-training and
 303 adds both the human-labeled and
 304 pseudo-labeled points in the training
 305 pool. With this procedure, AL+Self-
 306 Training first learns the best classifier
 307 ($\hat{h}_{\text{al-st}}$) with the given budget of max-
 308 imum training points (N_t) that can be
 309 queried from humans. Then it auto-labels all the remaining unlabeled points with this classifier’s pre-
 310 dictions. For AL+Self-Training+SC, we do selective auto-labeling using $\hat{h}_{\text{al-st}}$, i.e., only auto-label
 311 the points where the classifier will have an error at most ϵ_a . We use $\epsilon_a = 1\%$ here.

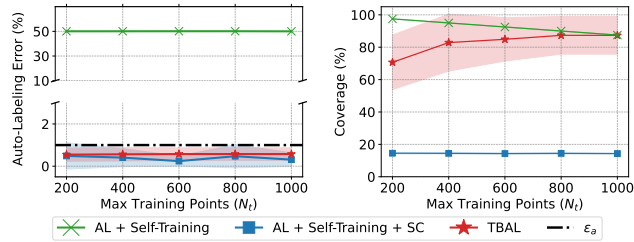


Figure 3: Results of experiment on 2-concentric circles to show the differences between TBAL, AL and ST.

312 **Results and Discussion.** The Figure 3 shows auto-labeling error and coverage achieved by these
 313 methods when run with different choices of human-labeled data budget for training. First, we can see
 314 that even with linear classifiers TBAL is able to auto-label a huge chunk of the data (high coverage)
 315 while maintaining auto-labeling error below the tolerance level of 1%. On the other hand, methods
 316 like AL+Self-Training (+SC) that try to first learn the optimal classifier in the given function class
 317 either have high auto-labeling error or very low coverage. These results are also consistent with
 318 the observations in [49] on the comparison between TBAL and AL, AL+SC. While such findings
 319 confirm the notion that there are differences—and, at least in some settings, advantages—for the
 320 TBAL approach compared to other techniques, we reiterate that our goal is to understand and improve
 321 the role of the confidence function within TBAL, rather than comparing TBAL to other techniques.

322 A.2 Details of the motivating experiment in Section 2

323 We run TBAL for a single round on the CIFAR-10 dataset with a SimpleCNN classification model
 324 with around 5.8M parameters [17]. We randomly sampled 4,000 points for training the classifier and
 325 randomly sampled 1,000 points as validation data. We train the model to zero training error using
 326 minibatch SGD with learning rate 1e-3, weight decay 1e-3 [13, 21], momentum 0.9, and batch size
 327 32. The trained model has validation accuracy around 55%, implying we could hope to get coverage
 328 around 55%. We run the auto-labeling procedure with an error tolerance of 5%.

329 B Appendix on Our Method

330 B.1 Detailed Algorithms

331 See Algorithms 1 and 2.

332 B.2 Tightness of surrogates.

333 The surrogate auto-labeling error and coverage introduced to relax the optimization problem (??) is
 334 indeed a good approximation of the actual auto-labeling error and coverage. To see this, we use a
 335 toy data setting of $x \sim \text{Uniform}(0, 1)$ with 1-dimensional threshold classifier $h_\theta(x) = \mathbb{1}(x \geq \theta)$.
 336 For any x , let true labels $y = h_{0.5}(x)$ and consider the confidence function $g_w(x) = |w - x|$. Let
 337 $\hat{y} = h_{0.25}(x)$ and consider the points on the side where $\hat{y} = 1$. We plot actual and surrogate errors in
 338 Figure 4(a) and the surrogate and actual coverage in Figure 4(a).

339 for three choices of α . As expected,
 340 the gap between the surrogates and
 341 the actual functions diminishes as we
 342 increase the α .

343 B.3 Active Querying Strategy.

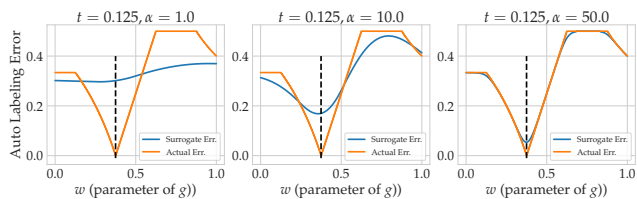
344 We employ the margin-random query
 345 approach to select the next batch of
 346 training data. This method involves
 347 sorting points based on their margin
 348 (uncertainty) scores and selecting the
 349 top Cn_b points, from which n_b points
 350 are randomly chosen. This strategy
 351 provides a straightforward and com-
 352 putationally efficient way to balance
 353 the exploration-exploitation trade-off.
 354 It’s important to acknowledge the ex-
 355 istence of alternative active-querying
 356 strategies; however, we adopt the
 357 margin-random approach as our stand-
 358 ard to maintain a focus on evaluat-
 359 ing various choices of confidence
 360 functions for auto-labeling. Note
 361 that while we use the new confidence
 362 scores computed using post-hoc methods for auto-labeling, we do not use these scores in active
 363 querying. Instead, we use the softmax scores from the model for this. We do this to avoid conflating
 364 the study with the study of active querying strategies. We use $C = 2$ for all experiments.

365 B.4 TBAL procedure with Colander

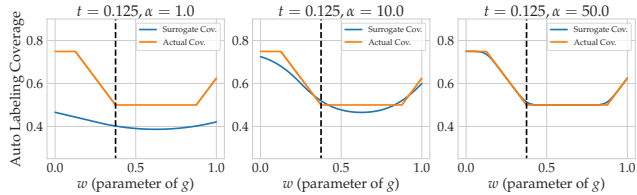
366 We take the workflow of TBAL and plugin our method Colander to learn the new confidence
 367 function and threshold. We discuss the updated workflow below and place the detailed Algorithms 1
 368 and 2 in the Appendix B due to space constraints.

369 **1. Initialization.** First select n_s points randomly from X_u and obtain human labels for them to
 370 create initial training data $D_{\text{train}}^{(1)}$. This is written as $\text{RANDOMQUERY}(X_u, n_s)$ in Algorithm 1. The
 371 procedure $\text{RANDOMQUERY}(X_u, n_s)$ selects n_s points randomly from X_u and obtains human labels
 372 for them to create $D_{\text{train}}^{(1)}$.

373 **2. Train classification model.** After obtaining human-labeled training data $D_{\text{train}}^{(i)}$ for the current
 374 round i , the procedure $\text{TRAINMODEL}(\mathcal{H}, D_{\text{train}}^{(i)})$ trains a model from model class \mathcal{H} on the training
 375 data $D_{\text{train}}^{(i)}$. Any training procedure can be used here. We use methods listed in Section ?? for model
 376 training. This step outputs a model \hat{h}_i trained on $D_{\text{train}}^{(i)}$.



(a) Auto-labeling error and surrogate error at various α .



(b) Auto-labeling coverage and surrogate coverage at various α .

Figure 4: Illustration of the tightness of surrogate error and coverage functions based on the choice of α .

377 **3. Learn new confidence function using Colander.** The model \hat{h}_i obtained in the previous step also
 378 produces softmax scores that can be used to for auto-labeling. However, as we saw earlier in Section
 379 2, using these scores may lead to poor auto-labeling performance. Thus, we plug in our procedure
 380 Colander to learn new scores that are designed to maximize the auto-labeling performance. We first
 381 randomly splits the validation data $D_{\text{val}}^{(i)}$ into $D_{\text{cal}}^{(i)}$ and $D_{\text{th}}^{(i)}$ using procedure $\text{RANDOMSPPLIT}(D_{\text{val}}^{(i)}, \nu)$.
 382 The part $D_{\text{cal}}^{(i)}$ has a fraction ν of the points from $D_{\text{val}}^{(i)}$. Then we consider problem P1 with \hat{h}_i and
 383 $D_{\text{cal}}^{(i)}$. We solve it to obtain the post-hoc confidence function \hat{g}_i .

384 **4. Threshold estimation.** The scores from the new confidence function \hat{g}_i on $D_{\text{th}}^{(i)}$ are used to
 385 estimate auto-labeling thresholds in Algorithm 2. This procedure finds thresholds for each class
 386 separately. It first splits the points in $D_{\text{th}}^{(i)}$ according to the ground truth class into subsets $D_{\text{th}}^{(i,y)}$.
 387 Then, for each class y , it finds the auto-labeling threshold $\hat{t}[y]$ by selecting the minimum threshold t
 388 such that the estimate of auto-labeling error plus a confidence interval, estimated on points in $D_{\text{th}}^{(i,y)}$
 389 having scores above t , is at most the given error tolerance ϵ_a . While we get thresholds as output
 390 from Colander, it is important to estimate them again from the held-out data $D_{\text{th}}^{(i)}$ to ensure the
 391 auto-labeling error constraint is not violated.

392 **5. Auto-labeling.** This is a simple step. We compute the scores on the remaining unlabeled data
 393 $X_u^{(i)}$ using the function \hat{g}_i and any point $\mathbf{x} \in X_u^{(i)}$ having score above $\hat{t}[\hat{y}]$ is assigned auto-label
 394 $\hat{y} = \hat{h}_i(\mathbf{x})$, and the points that did not meet this criterion remain unlabeled.

395 **6. Remove auto-labeled points.** The points that got auto-labeled in the previous steps are removed
 396 from the unlabeled pool. To make the validation data consistent with this unlabeled pool for the next
 397 round, the points in the validation data that fall into the auto-labeling region are also removed.

398 **7. Get more human-labeled data.** Lastly, it calls the procedure $\text{ACTIVEQUERY}(\hat{h}_i, X_u^{(i)}, n_b)$ to
 399 select n_b points from the remaining unlabeled pool using an active learning strategy. This newly
 400 acquired human-labeled data is added into the training data $D_{\text{train}}^{(i)}$. The details of the querying
 401 strategy are in Appendix B. The procedure then moves to step 2 and runs the loop until there are no
 402 more unlabeled points left or it has queried the stipulated number of human-labels N_t .

The notation is summarized in Table 2 below.

Symbol	Definition
$\mathbb{1}(E)$	indicator function of event E . It is 1 if E happens and 0 otherwise.
\mathcal{X}	feature space.
\mathcal{Y}	label space i.e. $1, 2, \dots, k$.
\mathcal{H}	hypothesis space (model class for the classifiers).
\mathcal{G}	class of confidence functions.
k	number of classes.
\mathbf{x}, y	\mathbf{x} is an element in \mathcal{X} and y is its true label.
h	a hypothesis (model) in \mathcal{H} .
g	confidence function $g : \mathcal{X} \rightarrow \Delta^k$.
X_u	given pool of unlabeled data points.
$X_u^{(i)}$	unlabeled data left at the beginning of i th round.
$\hat{h}^{(i)}$	ERM solution and auto-labeling thresholds respectively in i th round.
$D_{\text{query}}^{(i)}$	labeled data queried from oracle (human) in the i th round.
$D_{\text{train}}^{(i)}$	training data to learn $\hat{h}^{(i)}$ in the i th round.
$D_{\text{val}}^{(i)}$	validation data in the i th round.
$D_{\text{cal}}^{(i)}$	calibration data in the i th round to learn a post-hoc g .
$D_{\text{th}}^{(i)}$	part of validation data in the i th round to estimate threshold \mathbf{t} .
$D_{\text{auto}}^{(i)}$	part of $X_u^{(i)}$ that got auto-labeled in the i th round.
D_{out}	Output labeled data, including auto-labeled and human labeled data.
\mathbf{t}	k dimensional vector of thresholds.
$\mathbf{t}[y]$	y th entry of \mathbf{t} i.e. the threshold for class y .
$g(\mathbf{x})[y]$	the confidence score for class y output by confidence function g on data point \mathbf{x} .
\hat{y}	predicted class for data point \mathbf{x} .
f^*	unknown groundtruth labeling function.
N_u	number of unlabeled points, i.e. size of X_u .
N_t	number of manually labeled points that can be used for training h .
N_a	Total auto-labeled points in D_{out} .
ν	fraction of D_{val} that can be used for training post-hoc calibrator.
A	indices of points that are auto-labeled.
$X_u(A)$	subset of points in X_u with indices in A , i.e. the set of auto-labeled points.
\tilde{y}_i	label assigned to the i th point by the algorithm. It could be either y_i or \hat{y}_i .
y_i	groundtruth label for the i th point.
\hat{y}_i	predicted label for the i th point by classifier.
ϵ_a	auto-labeling error tolerance.
$\hat{\mathcal{E}}(g, \mathbf{t} \mid h)$	population level auto-labeling error, see eq. (??).
$\hat{\mathcal{P}}(g, \mathbf{t} \mid h)$	population level auto-labeling coverage, see eq. (??).
$\hat{\mathcal{E}}(g, \mathbf{t} \mid h, D)$	estimated auto-labeling error, see eq. (??).
$\hat{\mathcal{P}}(g, \mathbf{t} \mid h, D)$	estimated auto-labeling coverage, see eq. (??).
$\tilde{\mathcal{E}}(g, \mathbf{t} \mid h, D)$	surrogate estimated auto-labeling error, see eq. (3).
$\tilde{\mathcal{P}}(g, \mathbf{t} \mid h, D)$	surrogate estimated auto-labeling coverage, see eq. (2).

Table 2: Glossary of variables and symbols used in this paper.

Algorithm 1 Threshold-based Auto-Labeling (TBAL)

Input: Unlabeled data X_u , labeled validation data D_{val} , auto labeling error tolerance ϵ_a , N_t training data query budget, seed data size n_s , batch size for active query n_b , calibration data fraction ν , set of confidence thresholds T , coverage lower bound ρ_0 , label space \mathcal{Y} .

Output: Auto-labeled dataset D_{out}

```
1: procedure TBAL( $X_u, D_{\text{val}}, \epsilon_a, N_t, n_s, n_b, \nu, \rho_0, T, \mathcal{Y}$ )
2:    $\triangleright$  /*** Initialization. ***/
3:    $D_{\text{query}}^{(1)} \leftarrow \text{RANDOMQUERY}(X_u, n_s)$   $\triangleright$  Randomly select  $n_s$  points and get manual labels
   for them.
4:    $X_u^{(1)} \leftarrow X_u \setminus \{\mathbf{x} : (\mathbf{x}, y) \in D_{\text{query}}^{(1)}\}$   $\triangleright$  Remove the manually labeled points from the
   unlabeled pool.
5:    $D_{\text{val}}^{(1)} \leftarrow D_{\text{val}}; D_{\text{train}}^{(0)} \leftarrow \emptyset$   $\triangleright$  Validation data for the first round is full  $D_{\text{val}}$ .
6:    $D_{\text{out}} \leftarrow D_{\text{query}}^{(1)}; n_t^{(1)} \leftarrow n_s; i \leftarrow 1$   $\triangleright$  Include the manually labeled data in Step 2. in the
   output data  $D_{\text{out}}$ .
7:    $\triangleright$  /*** Run the auto-labeling loop ***/
8:    $\triangleright$  /* Until no more unlabeled points are left or the budget for manually labeled training data
   is exhausted. */
9:   while  $X_u^{(i)} \neq \emptyset$  and  $n_t^{(i)} \leq N_t$  do
10:     $D_{\text{train}}^{(i)} \leftarrow D_{\text{train}}^{(i-1)} \cup D_{\text{query}}^{(i)}$   $\triangleright$  Include the manually labeled points in the training data.
11:     $\hat{h}_i \leftarrow \text{TRAINMODEL}(\mathcal{H}, D_{\text{train}}^{(i)})$   $\triangleright$  Train a classification model.
12:     $D_{\text{cal}}^{(i)}, D_{\text{th}}^{(i)} \leftarrow \text{RANDOMSPLIT}(D_{\text{val}}^{(i)}, \nu)$   $\triangleright$  Randomly split the current validation data
   into two parts.
13:     $\triangleright$  /*** Colander block, to learn the new confidence function  $\hat{g}_i$  ***/
14:     $\hat{g}_i, \hat{\mathbf{t}}_i \leftarrow \arg \min_{g \in \mathcal{G}, \mathbf{t} \in T^k} -\tilde{\mathcal{F}}(g, \mathbf{t} \mid \hat{h}_i, D_{\text{cal}}^{(i)}) + \lambda \tilde{\mathcal{E}}(g, \mathbf{t} \mid \hat{h}_i, D_{\text{cal}}^{(i)})$   $\triangleright$  Colander
   procedure.
15:     $\triangleright$  /*** Estimate auto-labeling thresholds using  $\hat{g}_i$  and  $D_{\text{th}}^{(i)}$ . See Algorithm 2. ***/
16:     $\hat{\mathbf{t}}_i \leftarrow \text{ESTTHRESHOLD}(\hat{g}_i, \hat{h}_i, D_{\text{th}}^{(i)}, \epsilon_a, \rho_0, T, \mathcal{Y})$ 
17:     $\triangleright$  /*** Auto-label the points having scores above the thresholds. ***/
18:     $\tilde{D}_u^{(i)} \leftarrow \{(\mathbf{x}, \hat{h}_i(\mathbf{x})) : \mathbf{x} \in X_u^{(i)}\}$ 
19:     $D_{\text{auto}}^{(i)} \leftarrow \{(\mathbf{x}, \hat{y}) \in \tilde{D}_u^{(i)} : \hat{g}_i(\mathbf{x})[\hat{y}] \geq \hat{\mathbf{t}}_i[\hat{y}]\}$ 
20:     $X_u^{(i)} \leftarrow X_u^{(i)} \setminus \{\mathbf{x} : (\mathbf{x}, \hat{y}) \in D_{\text{auto}}^{(i)}\}$   $\triangleright$  Remove auto-labeled points from the unlabeled
   pool.
21:     $\tilde{D}_{\text{val}}^{(i)} \leftarrow \{(\mathbf{x}, \hat{h}_i(\mathbf{x})) : (\mathbf{x}, y) \in D_{\text{val}}^{(i)}\}$ 
22:     $D_{\text{val}}^{(i+1)} \leftarrow \{(\mathbf{x}, \hat{y}) \in \tilde{D}_{\text{val}}^{(i)} : \hat{g}_i(\mathbf{x})[\hat{y}] < \hat{\mathbf{t}}_i[\hat{y}]\}$   $\triangleright$  Remove validation points in the
   auto-labeling region.
23:     $\triangleright$  /*** Get the next batch of manually labeled data using an active querying strategy. ***/
24:     $D_{\text{query}}^{(i+1)} \leftarrow \text{ACTIVEQUERY}(\hat{h}_i, X_u^{(i)}, n_b)$ 
25:     $X_u^{(i+1)} \leftarrow X_u^{(i)} \setminus \{\mathbf{x} : (\mathbf{x}, y) \in D_{\text{query}}^{(i+1)}\}$   $\triangleright$  Remove manually labeled data from the
   unlabeled pool.
26:     $D_{\text{out}} \leftarrow D_{\text{out}} \cup D_{\text{auto}}^{(i)} \cup D_{\text{query}}^{(i+1)}$   $\triangleright$  Add the auto-labeled and manually labeled points in
   the output data.
27:     $n_t^{(i+1)} \leftarrow n_t^{(i)} + n_b$ 
28:     $i \leftarrow i + 1$ 
29:   end while
30:   return  $D_{\text{out}}$ 
31: end procedure
```

Algorithm 2 Estimate Auto-Labeling Threshold

Input: Confidence function \hat{g}_i , classifier \hat{h}_i , Part of validation data $D_{\text{th}}^{(i)}$ for threshold estimation, auto labeling error tolerance ϵ_a , set of confidence thresholds T , coverage lower bound ρ_0 , label space \mathcal{Y} .

Output: Auto-labeling thresholds $\hat{\mathbf{t}}_i$, where $\hat{\mathbf{t}}_i[y]$ is the threshold for class y .

```
1: procedure ESTTHRESHOLD( $\hat{g}_i, \hat{h}_i, D_{\text{th}}^{(i)}, \epsilon_a, \rho_0, T, \mathcal{Y}$ )
2:    $\triangleright$  /*** Estimate thresholds for each class. ***/
3:   for  $y \in \mathcal{Y}$  do
4:      $D_{\text{th}}^{(i,y)} \leftarrow \{(\mathbf{x}', y') \in D_{\text{th}}^{(i)} : y' = y\}$   $\triangleright$ Group points class-wise.
5:      $\triangleright$  /*** Only evaluate thresholds with est. coverage at least  $\rho_0$ . ***/
6:      $T'_y \leftarrow \{t \in T : \hat{\mathcal{P}}(\hat{g}_i, t | \hat{h}_i, D_{\text{th}}^{(i,y)}) \geq \rho_0\} \cup \{\infty\}$ 
7:      $\triangleright$  /*** Estimate auto-labeling error at each threshold. Pick the smallest threshold with the
        sum of estimated error and  $C_1$  times the standard deviation is below  $\epsilon_a$ .  $C_1$  is set to 0.25 here.
        ***/
8:      $\hat{\mathbf{t}}_i[y] \leftarrow \min\{t \in T'_y : \hat{\mathcal{E}}_a(\hat{g}_i, t | \hat{h}_i, D_{\text{th}}^{(i,y)}) + C_1 \hat{\sigma}(\hat{h}_i, t, D_{\text{th}}^{(i,y)}) \leq \epsilon_a\}$ 
9:   end for
10:  return  $\hat{\mathbf{t}}_i$ 
11: end procedure
```

405 **C Additional Experiments and Details**

406 **Choice of \mathcal{G} .** Our framework is flexible with respect to the choice of function class \mathcal{G} . In this work,
 407 we use neural networks with at least two layers on model class \mathcal{H} . We use representations from the last
 408 two layers of as input for the functions in \mathcal{G} . Let $\mathbf{z}^{(1)}(\mathbf{x}; h) \in \mathbb{R}^k$ and $\mathbf{z}^{(2)}(\mathbf{x}; h) \in \mathbb{R}^{d_2}$ be the outputs
 409 of the last and the second-last layer of the net h for input \mathbf{x} and let $\mathbf{z}(\mathbf{x}; h) := [\mathbf{z}^{(1)}(\mathbf{x}; h), \mathbf{z}^{(2)}(\mathbf{x}; h)]$
 410 denote the concatenation. This input is passed to network $\mathcal{G}_{nn_2} : \mathbb{R}^{k+d_2} \mapsto \Delta^k$; it outputs confidence
 411 scores for the k classes. Specifically g is defined as $g(\mathbf{x}) := \text{softmax}(\mathbf{W}_2 \tanh(\mathbf{W}_1 \mathbf{z}(\mathbf{x}; h)))$. Here
 412 $\mathbf{W}_1 \in \mathbb{R}^{(k+d_2) \times 2(k+d_2)}$ and $\mathbb{R}^{2(k+d_2) \times k}$ are the learnable weight matrices. As usual, for $\mathbf{v} \in \mathbb{R}^d$,
 413 $\text{softmax}(\mathbf{v})[i] := \exp(\mathbf{v}[i]) / (\sum_j \exp(\mathbf{v}[j]))$ and $\tanh(\mathbf{v})[i] := (\exp(2\mathbf{v}[i]) - 1) / (\exp(2\mathbf{v}[i]) + 1)$.

414 **C.1 Experiments on N_t, N_v and ν**

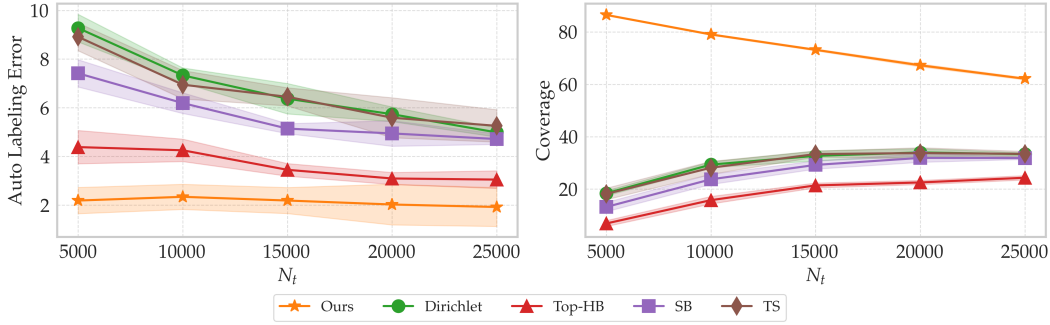


Figure 5: Autolabeling error and coverage of different post-hoc methods on CIFAR-10 for various N_t

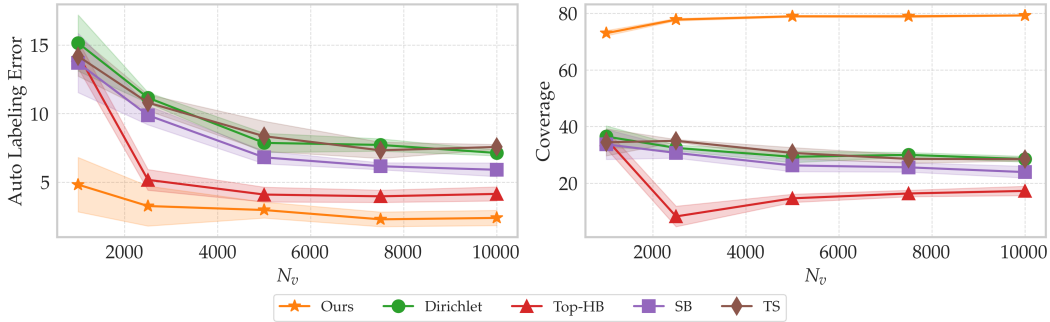


Figure 6: Autolabeling error and coverage of different post-hoc methods on CIFAR-10 for various N_v

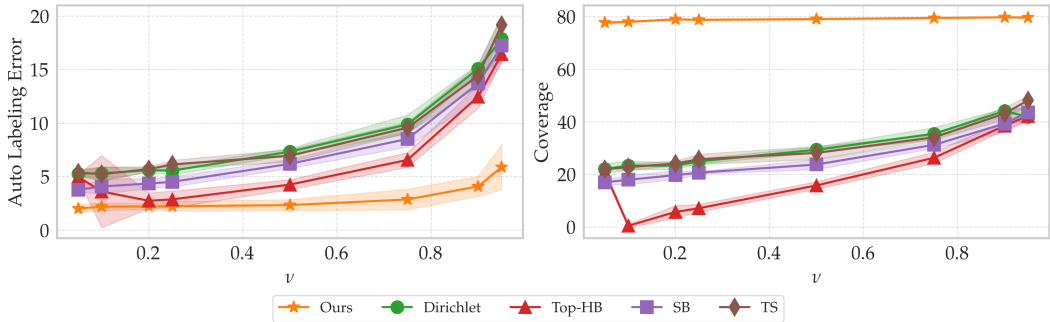


Figure 7: Autolabeling error and coverage of different post-hoc methods on CIFAR-10 for various ν

415 We need to understand the effect of training data query budget i.e. N_t , the total validation data N_v ,
 416 and the data that can be used for calibrating the model i.e. the calibration data fraction ν on the

417 auto-labeling objective. As varying these hyperparameters on each train-time method is expensive,
 418 we experimented with only Squentropy as it was the best-performing method across settings for
 419 various datasets.

420 When we vary the budget for training data N_t , we observe from Figure 5 that our method does not
 421 require a lot of data to train the base model, i.e. achieving low auto-labeling error and high coverage
 422 with a low budget. While other methods benefit from having more training data for auto-labeling
 423 objectives, it comes at the expense of reducing the available data for validation.

424 From figure 6, we observe that, while the coverage of our method remains the same across different
 425 N_v , it reduces for other methods. The cause of this phenomenon can be attributed to the fact that we
 426 are borrowing the data from the training budget as it limits the performance of the base model, which
 427 in turn limits the auto-labeling objective.

428 As we increase the percentage of data that can be used to calibrate the model, i.e., ν , we note
 429 from figure 7 that other methods improve the coverage, which can be understood from the fact that
 430 when more data is available for calibrating the model, the model becomes better in terms of the
 431 auto-labeling objective. But it’s interesting to note that even with a low calibration fraction, our
 432 method achieves superior coverage compared to other methods. It is also important to note that the
 433 auto-labeling error increases as we increase ν . This is because when ν increases, the number of data
 434 points used to estimate the threshold decreases, leading to a less granular and precise threshold.

Feature	Model	Error	Coverage
Pre-logits	Two Layer	4.6 ± 0.3	82.8 ± 0.5
Logits	Two Layer	3.2 ± 1.3	82.8 ± 0.3
Concat	Two Layer	3.3 ± 0.8	82.9 ± 0.4

Table 3: Auto-labeling error and coverage for the 3 feature representations we could use for 20 Newsgroup. As we can see, the feature representation does not lead to a significant difference in auto-labeling error and coverage.

Feature	Model	Error	Coverage
Pre-logits	Two Layer	2.1 ± 0.5	79.0 ± 0.2
Logits	Two Layer	3.1 ± 0.4	76.5 ± 0.9
Concat	Two Layer	2.3 ± 0.5	79.0 ± 0.3

Table 4: Auto-labeling error and coverage for the 3 feature representations we could use for CIFAR10 SimpleCNN. As we can see, the feature representation does not lead to a significant difference in auto-labeling error and coverage.

435 **C.2 Experiments on Colander input**

436 Figure 14 illustrates that we could use logits (last layer’s
 437 representations), pre-logits (second last layer’s representa-
 438 tions), or the concatenation of these two as the input to g .
 439 To help us decide which one we should use, we conduct a
 440 hyperparameter search for input features on the CIFAR-10
 441 and 20 Newsgroup dataset using the Squentropy train-time
 442 method. Table 3 and 4 present the auto-labeling error and
 443 coverage of using the 3 types of feature representations.
 444 As we can see, all feature representation leads to a similar
 445 auto-labeling error and coverage, and in some cases,
 446 it is better to include pre-logits as well. Therefore, we
 447 use concatenated representation (Concat), allowing more
 448 flexibility.

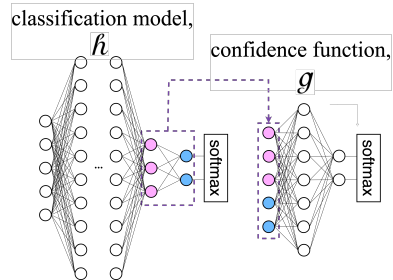


Figure 14: Our choice of g function.

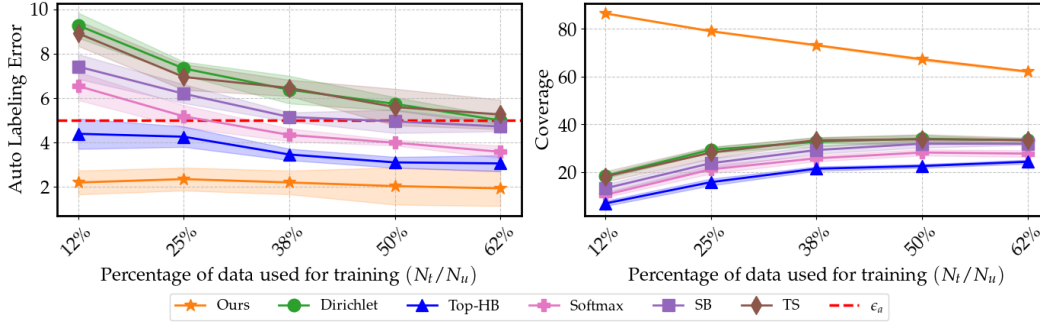


Figure 8: Auto-labeling error and coverage for different post-hoc methods on CIFAR-10 while we vary N_t . $N_u = 40,000$ is the size of the given unlabeled pool.

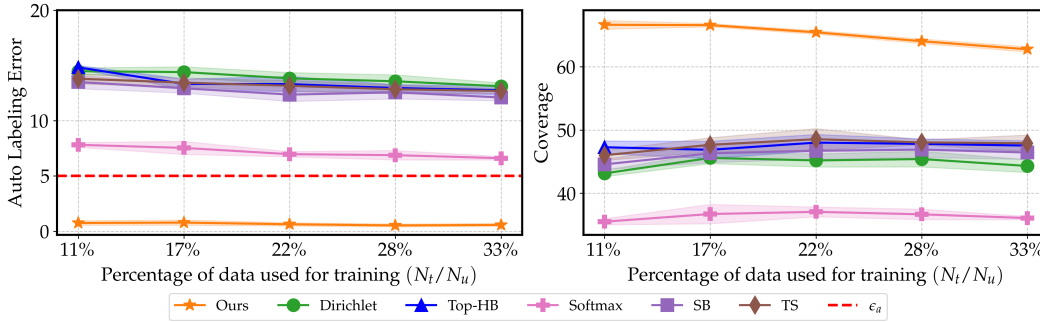


Figure 9: Auto-labeling error and coverage for different post-hoc methods on Tiny-ImageNet while we vary N_t . $N_u = 90,000$ is the size of the given unlabeled pool.

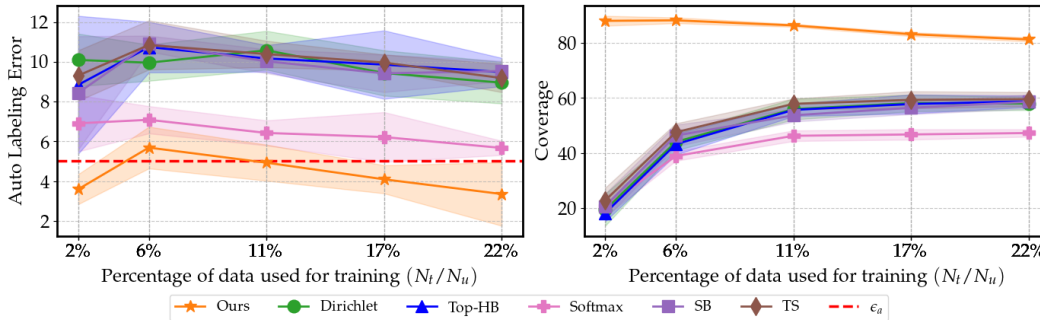


Figure 10: Auto-labeling error and coverage for different post-hoc methods on 20 Newsgroups while we vary N_t . $N_u = 9,052$ is the size of the given unlabeled pool.

449 C.3 Experiments on ϵ_a

450 We run TBAL with five values of $\epsilon_a \in \{0.01, 0.025, 0.05, 0.075, 0.1\}$ and report the results in Table
 451 5. As expected the auto-labeling error is high with larger values of and smaller with small ϵ_a .

452 C.4 Experiments on multiple rounds

453 We further demonstrate that the performance gains are due to the use of Colander, even if methods
 454 use multiple rounds. To do so, we show the evolution of coverage and error over multiple rounds
 455 in Figure 15. The effects of using Colander are visible from the first round itself, and the following
 456 rounds improve performance further. We also run a single round (passive) variant of TBAL where
 457 we sample all the human-labeled points for training (N_t) randomly at once, train a classifier, do
 458 auto-labeling, and then stop. This setting avoids confounding due to multiple rounds. We observe
 459 that using Colander yields significantly higher coverage in comparison to the baselines (see Table 6).

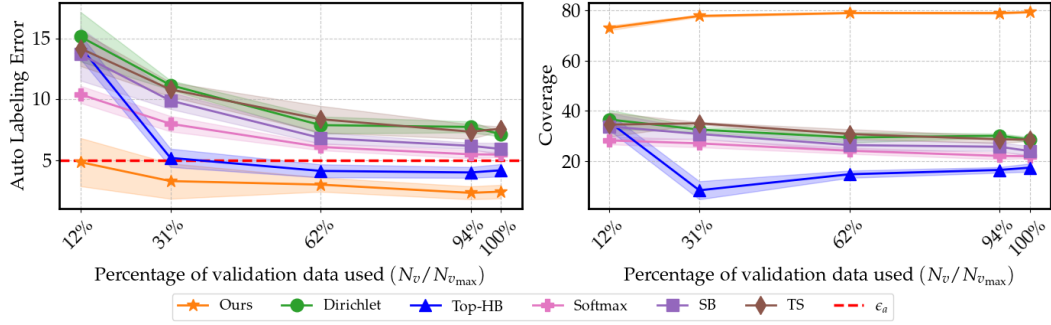


Figure 11: Auto-labeling error and coverage for different post-hoc methods on CIFAR-10 while we vary N_v . $N_{v_{\max}} = 8,000$ is the maximum number of points available for validation.

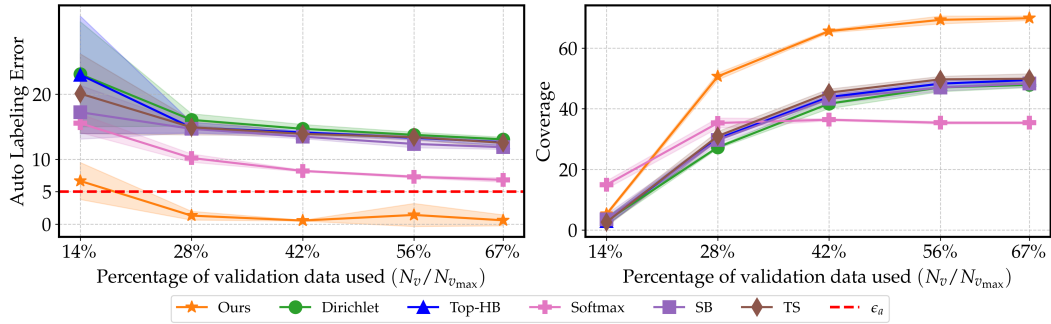


Figure 12: Auto-labeling error and coverage for different post-hoc methods on Tiny-ImageNet while we vary N_v . $N_{v_{\max}} = 18,000$ is the maximum number of points available for validation.

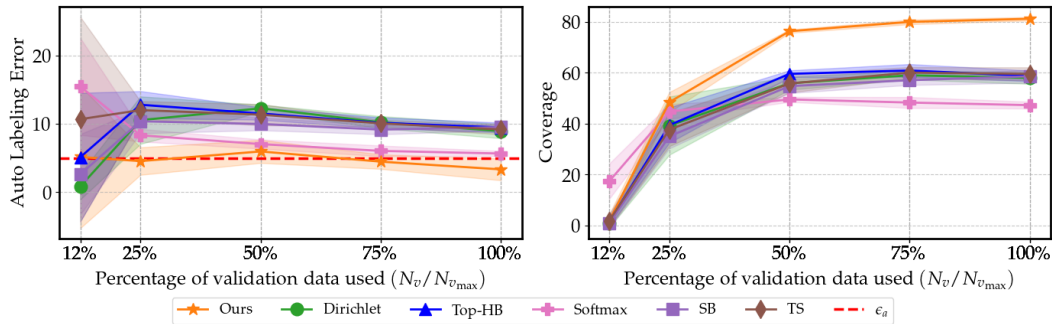
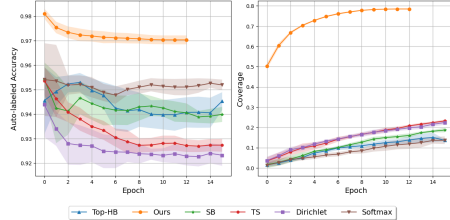


Figure 13: Auto-labeling error and coverage for different post-hoc methods on 20 Newsgroups while we vary N_v . $N_{v_{\max}} = 1,600$ is the maximum number of points available for validation.

460 This reinforces the fact that the gains in the multi-round TBAL are directly due to Colander, while
 461 multiple rounds of data selection, training, and auto-labeling are superior to doing everything in a
 462 single round.

Post-hoc Method	$\epsilon_a = 0.01$		$\epsilon_a = 0.025$		$\epsilon_a = 0.05$		$\epsilon_a = 0.075$		$\epsilon_a = 0.1$	
	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)	Err (\downarrow)	Cov (\uparrow)
Softmax	5.86 \pm 0.38	12.73 \pm 1.61	5.86 \pm 0.38	12.73 \pm 1.61	4.78 \pm 0.21	14.01 \pm 2.08	6.80 \pm 0.47	16.73 \pm 1.19	9.03 \pm 0.17	21.28 \pm 0.82
TS	8.19 \pm 0.88	19.44 \pm 1.16	8.19 \pm 0.88	19.44 \pm 1.16	7.26 \pm 0.29	23.15 \pm 0.7	9.24 \pm 0.78	22.49 \pm 0.74	11.63 \pm 0.51	25.79 \pm 1.97
Dirichlet	8.22 \pm 0.4	16.94 \pm 1.2	8.22 \pm 0.4	16.94 \pm 1.2	7.6 \pm 0.48	22.36 \pm 1.18	9.68 \pm 0.82	18.65 \pm 0.97	11.26 \pm 1.16	24.91 \pm 2.09
SB	6.15 \pm 0.52	11.74 \pm 0.57	6.15 \pm 0.52	11.74 \pm 0.57	6.09 \pm 0.35	18.58 \pm 1.13	7.81 \pm 0.65	17.37 \pm 1.3	9.13 \pm 1.08	20.52 \pm 1.11
Top-HB	5.76 \pm 0.42	9.89 \pm 0.55	5.76 \pm 0.42	9.89 \pm 0.55	5.95 \pm 0.7	15.58 \pm 1.92	7.45 \pm 0.8	13.84 \pm 0.78	8.71 \pm 1.37	17.9 \pm 0.56
Ours	1.2 \pm 0.18	78.33 \pm 0.76	1.32 \pm 0.21	78.75 \pm 0.4	2.96 \pm 0.2	78.48 \pm 0.17	4.3 \pm 0.23	78.94 \pm 0.42	6.29 \pm 0.5	78.97 \pm 0.46

Table 5: ϵ_a variation. Dataset: CIFAR-10, Train-time method: Vanilla.



463

Figure 15: Clarification on multiple rounds. Per-epoch metrics for all post-hoc methods for CIFAR10. (left) Auto-labeling accuracy (right) Coverage. Train time method is vanilla and model is medium net.

Post-hoc method	Err (\downarrow)	Coverage (\uparrow)
Softmax	2.7 \pm 0.54	11.06 \pm 1.46
TS	3.04 \pm 0.49	12.03 \pm 1.98
Dirichlet	2.98 \pm 0.32	11.22 \pm 2.1
SB	2.72 \pm 0.34	9.75 \pm 1.33
Top-HB	1.83 \pm 0.61	5.50 \pm 1.08
Ours	2.02 \pm 0.28	49.62 \pm 0.69

Table 6: Results with single round of auto-labeling. Dataset and model: CIFAR-10 setting in the paper.

464 C.5 Experiments on different architectures

465 In TBAL it is not a priori clear what model the practitioner should use. The overall system is flexible
 466 enough to work with any chosen model class. Our focus is on evaluating the effect of various training
 467 time and post-hoc methods designed to improve the confidence functions for any given model. To
 468 answer the query, we ran experiments with Resnet18 and ViT models in the CIFAR-10 setting (see
 469 Table 7). As we expected there are variations in the results in the baselines due to model choices but
 470 our method maintains high performance irrespective of the classification model used. This is due to
 471 its ability to learn confidence scores tailored for TBAL.

Post-hoc Method	Err (\downarrow)	Coverage (\uparrow)
Softmax	14.02 \pm 1.83	2.03 \pm 0.31
TS	19.32 \pm 2.51	2.54 \pm 0.33
Dirichlet	17.27 \pm 3.26	2.87 \pm 0.55
SB	9.22 \pm 10.91	0.46 \pm 0.51
Top-HB	0.00 \pm 0.00	0.00 \pm 0.00
Ours	2.62 \pm 0.32	75.56 \pm 0.15

Post-hoc Method	Err (\downarrow)	Coverage (\uparrow)
Softmax	4.48 \pm 0.23	33.24 \pm 1.14
TS	6.38 \pm 0.47	39.14 \pm 1.96
Dirichlet	6.30 \pm 0.41	37.99 \pm 1.47
SB	5.16 \pm 0.23	35.32 \pm 1.36
Top-HB	4.46 \pm 0.40	29.66 \pm 0.74
Ours	2.85 \pm 0.25	78.56 \pm 0.54

Table 7: Model variation. CIFAR-10 dataset with ViT (Left) and ResNet18 (Right), Train-time method Vanilla.

472 C.6 Hyperparameters

473 The hyperparameters and their values we swept over are listed in Table 8 and 9 for train-time and
 474 post-hoc methods, respectively.

475 C.7 Train-time and post-hoc methods

476 C.7.1 Train-time methods

477 1. *Vanilla*: Neural networks are commonly trained by minimizing the cross entropy loss using
 478 stochastic gradient descent (SGD) with momentum [1, 3]. We refer to this as the Vanilla training
 479 method. We also include weight decay to mitigate the overconfidence issue associated with this
 480 method [10].

Method	Hyperparameter	Values
Common	optimizer	SGD
	learning rate	0.001, 0.01, 0.1
	batch size	32, <u>256</u>
	max epoch	50, <u>100</u>
	weight decay	0.001, 0.01, 0.1
	momentum	0.9
CRL	rank target	softmax
	rank weight	0.7, 0.8, 0.9
FMFP	optimizer	SAM

Table 8: Hyperparameters swept over for train-time methods. Those listed next to Common are the hyperparameters for the four train-time methods: Vanilla, CRL, FMFP, and Squentropy. Therefore, we do not list those again for each method. Note that for FMFP, we used SAM optimizer instead of SGD. For each method, we swept through all possible combinations of the possible values for each hyperparameter. Underlined values are only used on TinyImageNet since it is a complicated dataset containing 200 classes.

- 481 2. *Squentropy* [16]: This method adds the average square loss over the incorrect classes to the
482 cross-entropy loss. This simple modification to the Vanilla method leads to the end model with
483 better test accuracy and calibration.
- 484 3. *Correctness Ranking Loss (CRL)* [30]: This method includes a term in the loss function of the
485 vanilla training method so that the confidence scores of the model are aligned with the ordinal
486 rankings criterion [15, 5]. The confidence functions satisfying this criterion produce high
487 scores on points where the probability of correctness is high and low scores on points with low
488 probabilities of being correct.
- 489 4. *FMFP* [53] aims to align confidence scores with the ordinal rankings criterion. It uses Sharpness
490 Aware Minimizer (SAM) [6] to train the model, with the expectation that the flat minima would
491 benefit the ordinal rankings objective of the confidence function.

492 C.7.2 Post-hoc methods

- 493 1. *Temperature scaling* [10]: This is a variant of Platt scaling [10], a classic and one of the easiest
494 parametric methods for post-hoc calibration. It rescales the logits by a learnable scalar parameter
495 and has been shown to work well for neural networks.
- 496 2. *Top-Label Histogram-Binning* [12]: Since TBAL assigns the top labels (predicted labels)
497 to the selected unlabeled points, it is appealing to only calibrate the scores of the predicted
498 label. Building upon a rich line of histogram-binning methods (non-parametric) for post-hoc
499 calibration [52], this method focuses on calibrating the scores of predicted labels.
- 500 3. *Scaling-Binning* [24]: This method combines parametric and non-parametric methods. It first
501 applies temperature scaling and then bins the confidence function values to ensure calibration.
- 502 4. *Dirichlet Calibration* [22]: This method models the distribution of predicted probability vectors
503 separately on instances of each class and assumes the class conditional distributions are Dirichlet
504 distributions with different parameters. It uses linear parameterization for the distributions,
505 which allows easy implementation in neural networks as additional layers and softmax output.

506 **Note:** For binning methods, uniform mass binning [52] has been a better choice over uniform width
507 binning. Hence, we use uniform mass binning as well.

508 C.8 Compute resources

509 Our experiments were conducted on machines equipped with the NVIDIA RTX A6000 and NVIDIA
510 GeForce RTX 4090 GPUs.

511 **C.9 Detailed dataset and model**

- 512 1. The MNIST dataset [26] consists of 28×28 grayscale images of hand-written digits across 10
513 classes. It was used alongside the LeNet5 [27], a convolutional neural network, for auto-labeling.
- 514 2. The CIFAR-10 dataset [20] contains $3 \times 32 \times 32$ color images across 10 classes. We utilized
515 its raw pixel matrix in conjunction with SimpleCNN [17], a convolutional neural network with
516 approximately 5.8M parameters, for auto-labeling.
- 517 3. Tiny-ImageNet [25] is a color image dataset that consists of 100K images across 200 classes.
518 Instead of using the $3 \times 64 \times 64$ raw pixel matrices as input, we utilized CLIP [38] to derive
519 embeddings within the \mathbb{R}^{512} vector space. We used a 3-layer perceptron (1,000-500-300) as the
520 auto-labeling model.
- 521 4. 20 Newsgroups [29, 35] is a natural language dataset comprising around 18,000 news posts
522 across 20 topics. We used the FlagEmbedding [50] to map the textual data into \mathbb{R}^{1024} embed-
523 dings. We used a 3-layer perceptron (1,000-500-30) as the auto-labeling model.

524 **C.10 Detailed experiments protocol**

525 We predefined TBAL hyperparameters for each dataset-model pair and the hyperparameters we will
526 sweep for each train-time and post-hoc method in Table 8 and Table 9 respectively. For a dataset-
527 model pair, initially, we perform a hyperparameter search for the train-time method. Subsequently,
528 we optimize the hyperparameters for post-hoc methods while keeping the train-time method fixed
529 with the previously found optimum hyperparameter for that dataset-model pair.

530 We fix the hyperparameters for the train-time method while searching hyperparameters for the post-
531 hoc method to alleviate computational budget throttle. We effectively reduce the search space to the
532 sum of the cardinalities of unique hyper-parameter combinations across the two methods instead of a
533 larger multiplicative product. Furthermore, due to the independent nature of these hyper-parameter
534 combinations, TBAL runs can be highly parallelized to expedite the search process.

535 Since TBAL operates iteratively to acquire human labels for model training, selecting hyper-
536 parameters at each round of TBAL could quickly become intractable and lose its practical significance.
537 To better align with its practical usage, we only conducted a hyperparameter search for the initial
538 TBAL round. The specific set of hyperparameters used for the search are reported in Table 9.

539 After completing the hyperparameter search for train-time and post-hoc methods, the determined
540 hyperparameter combinations are subjected to a full evaluation across all iterations of TBAL. At
541 the end of each iteration, the auto-labeled points are evaluated against their ground truth labels to
542 determine their auto-labeling error. These points are then added to the auto-labeled set, where their
543 ratio to the total amount of unlabeled data determines the coverage. This iterative process continues
544 until all unlabeled data are exhaustively labeled by either the oracle or through auto-labeling in the
545 final iteration. The auto-labeling error and coverage at the final iteration of TBAL are then recorded.

546 Since TBAL incorporates randomized components as detailed in Algorithm 1, we ran the algorithm 5
547 times, each with a unique random seed while maintaining the same hyperparameter combination. We
548 then recorded the results from the final iteration of these runs and calculated the mean and standard
549 deviation of both auto-labeling error and coverage. These figures are reported in Table 1.

550 A limitation of the grid search approach in hyper-parameter optimization becomes apparent when our
551 predefined hyper-parameter choices result in sub-optimal coverage and auto-labeling errors. Using
552 these sub-optimal hyper-parameters can adversely affect the multi-round iterative process in TBAL,
553 prompting the need for repetitive searches to find more effective hyper-parameters. When encounter-
554 ing such scenarios, TBAL users should explore additional hyper-parameter options until satisfactory
555 performance is achieved in the initial round. However, we opted for a more straightforward approach
556 to hyper-parameter selection, mindful of the computational demands of repeatedly optimizing mul-
557 tiple hyper-parameters across different methods. In scenarios expressed conditionally, we retained
558 the top-1 hyper-parameter combination for any given method if it achieved the highest coverage
559 while adhering to the specified error margin (ϵ_a). If no hyper-parameter combinations yielded an
560 auto-labeling error at most equal to the error margin (ϵ_a), we then chose the hyper-parameter combi-
561 nation with the lowest auto-labeling error, regardless of its coverage. In the case of ties, we resolved
562 them through random selection. This process results in obtaining singular values for each choice of
563 hyper-parameter after completing each method’s hyper-parameter search.

564 D Broader Impact

565 This paper contributes to the advancement of the practice of creating labeled datasets in machine
566 learning. While our work has various possible societal implications, we do not identify any specific
567 concerns that require special attention in this context.

568 E Related Work

569 **Data Labeling.** We briefly discuss prominent methods for labeling. Crowdsourcing [40, 43]
570 uses a crowd of non-experts to complete a set of labeling tasks. Works in this domain focus on
571 mitigating noise in the obtained information, modeling label errors, and designing effective labeling
572 tasks [9, 18, 28, 46, 45, 47, 4]. Weak supervision (WS), in contrast, emphasizes labeling through
573 multiple inexpensive but noisy sources, not necessarily human [39, 7, 42, 48]. Works such as [39, 7]
574 concentrate on binary or multi-class labeling, while [42, 48] extend WS to structured prediction tasks.

575 Auto-labeling occupies an intermediate position between weak supervision and crowdsourcing in
576 terms of human dependency. It aims to minimize costs to obtain human labels while generating
577 high-quality labeled data using a specific model. [37] use a TBAL-like algorithm and explore the cost
578 of training for auto-labeling with large-scale model classes. Recent work [49] theoretically analyzes
579 the sample complexity of validation data required to guarantee the quality of auto-labeled data.

580 **Overconfidence and calibration.** The issue of overconfidence [44, 33, 14, 2] is detrimental in
581 several applications, including ours. Many solutions have emerged to mitigate the overconfidence
582 and miscalibration problem. Gawlikowski et al. [8] provide a comprehensive survey on uncertainty
583 quantification and calibration techniques for neural networks. Guo et al. [10] evaluated a variety of
584 solutions ranging from the choice of network architecture, model capacity, weight decay regularization
585 [21], histogram-binning and isotonic regression [51, 52] and temperature scaling [36, 34] which
586 they found to be the most promising solution. The solutions fall into two broad categories: train-time
587 and post-hoc. Train-time solutions modify the loss function, include additional regularization terms,
588 or use different training procedures [23, 32, 31, 16]. On the other hand, post-hoc methods such as
589 top-label histogram-binning [11], scaling binning [24], Dirichlet calibration [22] calibrate the scores
590 directly or learn a model that corrects miscalibrated confidence scores.

591 **Beyond calibration.** While calibration aims to match the confidence scores with a probability
592 of correctness, it is not the precise solution to the overconfidence problem in many applications,
593 including our setting. The desirable criteria for scores for TBAL are closely related to the ordinal
594 ranking criterion [15]. To get such scores, Corbière et al. [5] add a module in the net for failure
595 prediction, Zhu et al. [53] switch to sharpness aware minimization [6] to learn the model; CRL [30]
596 regularizes the loss.

597 F NeurIPS Paper Checklist

598 1. Claims

599 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s
600 contributions and scope?

601 Answer: [Yes]

602 Justification: Our claims are backed by our novel technique in Section C and thorough empirical
603 evaluation in Section 4.

604 Guidelines:

- 605 • The answer NA means that the abstract and introduction do not include the claims made in
606 the paper.
- 607 • The abstract and/or introduction should clearly state the claims made, including the contri-
608 butions made in the paper and important assumptions and limitations. A No or NA answer
609 to this question will not be perceived well by the reviewers.
- 610 • The claims made should match theoretical and experimental results, and reflect how much
611 the results can be expected to generalize to other settings.
- 612 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are
613 not attained by the paper.

614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss them briefly.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: It is an empirical paper, it does not have theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

666 Justification: All the necessary details are provided in Section 4 and in the Appendix C. We have
667 also uploaded the code along with the submission.

668 Guidelines:

- 669 • The answer NA means that the paper does not include experiments.
- 670 • If the paper includes experiments, a No answer to this question will not be perceived well
671 by the reviewers: Making the paper reproducible is important, regardless of whether the
672 code and data are provided or not.
- 673 • If the contribution is a dataset and/or model, the authors should describe the steps taken to
674 make their results reproducible or verifiable.
- 675 • Depending on the contribution, reproducibility can be accomplished in various ways. For
676 example, if the contribution is a novel architecture, describing the architecture fully might
677 suffice, or if the contribution is a specific model and empirical evaluation, it may be
678 necessary to either make it possible for others to replicate the model with the same dataset,
679 or provide access to the model. In general, releasing code and data is often one good way
680 to accomplish this, but reproducibility can also be provided via detailed instructions for
681 how to replicate the results, access to a hosted model (e.g., in the case of a large language
682 model), releasing of a model checkpoint, or other means that are appropriate to the research
683 performed.
- 684 • While NeurIPS does not require releasing code, the conference does require all submissions
685 to provide some reasonable avenue for reproducibility, which may depend on the nature of
686 the contribution. For example
 - 687 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
688 reproduce that algorithm.
 - 689 (b) If the contribution is primarily a new model architecture, the paper should describe
690 the architecture clearly and fully.
 - 691 (c) If the contribution is a new model (e.g., a large language model), then there should
692 either be a way to access this model for reproducing the results or a way to reproduce
693 the model (e.g., with an open-source dataset or instructions for how to construct the
694 dataset).
 - 695 (d) We recognize that reproducibility may be tricky in some cases, in which case authors
696 are welcome to describe the particular way they provide for reproducibility. In the
697 case of closed-source models, it may be that access to the model is limited in some
698 way (e.g., to registered users), but it should be possible for other researchers to have
699 some path to reproducing or verifying the results.

700 5. Open access to data and code

701 Question: Does the paper provide open access to the data and code, with sufficient instructions
702 to faithfully reproduce the main experimental results, as described in supplemental material?

703 Answer: [Yes]

704 Justification: We use publicly available datasets and uploaded the code as supplementary
705 material.

706 Guidelines:

- 707 • The answer NA means that paper does not include experiments requiring code.
- 708 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/public/
709 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 710 • While we encourage the release of code and data, we understand that this might not be
711 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
712 including code, unless this is central to the contribution (e.g., for a new open-source
713 benchmark).
- 714 • The instructions should contain the exact command and environment needed to run to
715 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 716 • The authors should provide instructions on data access and preparation, including how to
717 access the raw data, preprocessed data, intermediate data, and generated data, etc.
718

- 719
- 720
- 721
- 722
- 723
- 724
- 725
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
 - At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
 - Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

726 6. Experimental Setting/Details

727 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

729 Answer: [\[Yes\]](#)

730 Justification: These details are provided in the Section 4 and Appendix C.

731 Guidelines:

- 732
- 733
- 734
- 735
- 736
- The answer NA means that the paper does not include experiments.
 - The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
 - The full details can be provided either with the code, in appendix, or as supplemental material.

737 7. Experiment Statistical Significance

738 Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

740 Answer: [\[Yes\]](#)

741 Justification: We run each setting with multiple random seeds and report the mean, standard deviations of the evaluation metrics.

743 Guidelines:

- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- 760
- 761
- 762
- The answer NA means that the paper does not include experiments.
 - The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
 - The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
 - The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
 - The assumptions made should be given (e.g., Normally distributed errors).
 - It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

763 8. Experiments Compute Resources

764 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

767 Answer: [\[Yes\]](#)

768 Justification: Provided in the Appendix C.

769 Guidelines:

- 770 • The answer NA means that the paper does not include experiments.
- 771 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or
- 772 cloud provider, including relevant memory and storage.
- 773 • The paper should provide the amount of compute required for each of the individual
- 774 experimental runs as well as estimate the total compute.
- 775 • The paper should disclose whether the full research project required more compute than
- 776 the experiments reported in the paper (e.g., preliminary or failed experiments that didn't
- 777 make it into the paper).

778 9. Code Of Ethics

779 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS
780 Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

781 Answer: [Yes]

782 Justification: We have followed the NeurIPS Code of Ethics to the best of our knowledge.

783 Guidelines:

- 784 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 785 • If the authors answer No, they should explain the special circumstances that require a
- 786 deviation from the Code of Ethics.
- 787 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration
- 788 due to laws or regulations in their jurisdiction).

789 10. Broader Impacts

790 Question: Does the paper discuss both potential positive societal impacts and negative societal
791 impacts of the work performed?

792 Answer: [Yes]

793 Justification: The paper has a brief discussion on the broader impacts.

794 Guidelines:

- 795 • The answer NA means that there is no societal impact of the work performed.
- 796 • If the authors answer NA or No, they should explain why their work has no societal impact
- 797 or why the paper does not address societal impact.
- 798 • Examples of negative societal impacts include potential malicious or unintended uses
- 799 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g.,
- 800 deployment of technologies that could make decisions that unfairly impact specific groups),
- 801 privacy considerations, and security considerations.
- 802 • The conference expects that many papers will be foundational research and not tied to
- 803 particular applications, let alone deployments. However, if there is a direct path to any
- 804 negative applications, the authors should point it out. For example, it is legitimate to point
- 805 out that an improvement in the quality of generative models could be used to generate
- 806 deepfakes for disinformation. On the other hand, it is not needed to point out that a generic
- 807 algorithm for optimizing neural networks could enable people to train models that generate
- 808 Deepfakes faster.
- 809 • The authors should consider possible harms that could arise when the technology is being
- 810 used as intended and functioning correctly, harms that could arise when the technology is
- 811 being used as intended but gives incorrect results, and harms following from (intentional or
- 812 unintentional) misuse of the technology.
- 813 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 814 strategies (e.g., gated release of models, providing defenses in addition to attacks, mecha-
- 815 nisms for monitoring misuse, mechanisms to monitor how a system learns from feedback
- 816 over time, improving the efficiency and accessibility of ML).

817 11. Safeguards

818 Question: Does the paper describe safeguards that have been put in place for responsible release
819 of data or models that have a high risk for misuse (e.g., pretrained language models, image
820 generators, or scraped datasets)?

821 Answer: [NA]

822 Justification: The paper does not release such data or models that have high risk for misuse.

823 Guidelines:

- 824 • The answer NA means that the paper poses no such risks.
- 825 • Released models that have a high risk for misuse or dual-use should be released with
826 necessary safeguards to allow for controlled use of the model, for example by requiring
827 that users adhere to usage guidelines or restrictions to access the model or implementing
828 safety filters.
- 829 • Datasets that have been scraped from the Internet could pose safety risks. The authors
830 should describe how they avoided releasing unsafe images.
- 831 • We recognize that providing effective safeguards is challenging, and many papers do not
832 require this, but we encourage authors to take this into account and make a best faith effort.

833 12. Licenses for existing assets

834 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the
835 paper, properly credited and are the license and terms of use explicitly mentioned and properly
836 respected?

837 Answer: [Yes]

838 Justification: We have appropriately credited them along with citations.

839 Guidelines:

- 840 • The answer NA means that the paper does not use existing assets.
- 841 • The authors should cite the original paper that produced the code package or dataset.
- 842 • The authors should state which version of the asset is used and, if possible, include a URL.
- 843 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 844 • For scraped data from a particular source (e.g., website), the copyright and terms of service
845 of that source should be provided.
- 846 • If assets are released, the license, copyright information, and terms of use in the package
847 should be provided. For popular datasets, paperswithcode.com/datasets has curated
848 licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 849 • For existing datasets that are re-packaged, both the original license and the license of the
850 derived asset (if it has changed) should be provided.
- 851 • If this information is not available online, the authors are encouraged to reach out to the
852 asset's creators.

853 13. New Assets

854 Question: Are new assets introduced in the paper well documented and is the documentation
855 provided alongside the assets?

856 Answer: [Yes]

857 Justification: The code is well documented along with instructions to run.

858 Guidelines:

- 859 • The answer NA means that the paper does not release new assets.
- 860 • Researchers should communicate the details of the dataset/code/model as part of their sub-
861 missions via structured templates. This includes details about training, license, limitations,
862 etc.
- 863 • The paper should discuss whether and how consent was obtained from people whose asset
864 is used.
- 865 • At submission time, remember to anonymize your assets (if applicable). You can either
866 create an anonymized URL or include an anonymized zip file.

867 14. Crowdsourcing and Research with Human Subjects

868 Question: For crowdsourcing experiments and research with human subjects, does the paper
869 include the full text of instructions given to participants and screenshots, if applicable, as well as
870 details about compensation (if any)?

871 Answer: [NA]

872 Justification: We did not use crowdsourcing or human subjects in the paper.

873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We did not use crowdsourcing or human subjects in the paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

Method	Hyperparameter	Values
Temperature scaling	optimizer	Adam
	learning rate	0.001, 0.01, 0.1
	batch size	64
	max epoch	500
	weight decay	0.01, 0.1, 1
Top-label histogram binning	points per bin	25, 50
Scaling-binning	number of bins	15, 25
	learning rate	0.001, 0.01, 0.1
	batch size	64
	max epoch	500
	weight decay	0.01, 0.1, 1
Dirichlet calibration	regularization parameter	0.001, 0.01, 0.1
Ours	λ	10, 100
	features key	concat
	class-wise	independent
	optimizer	Adam
	learning rate	0.01, 0.1
	max epoch	500
	weight decay	0.01, 0.1, 1
	batch size	64
	regularize	false
	α	0.01, 0.1, 1

Table 9: Hyperparameters swept over for post-hoc methods. For each method, we swept through all possible combinations of the possible values for each hyperparameter.