# SUPClust: Active Learning at the Boundaries

**Yuta Ono, Till Aczel, Benjamin Estermann & Roger Wattenhofer**
ETH Zürich
{yutono,taczel,estermann,wattenhofer}@ethz.ch

## Abstract

Active learning is a machine learning paradigm designed to optimize model performance in a setting where labeled data is expensive to acquire. In this work, we propose a novel active learning method called SUPClust that seeks to identify points at the decision boundary between classes. By targeting these points, SUPClust aims to gather information that is most informative for refining the model's prediction of complex decision regions. We demonstrate experimentally that labeling these points leads to strong model performance. This improvement is observed even in scenarios characterized by strong class imbalance.

## 1 Introduction

Progress in deep learning for classification tasks has been following an impressive pace in recent years (Ioffe & Szegedy, 2015; Dosovitskiy et al., 2021; Srivastava & Sharma, 2024). In order to achieve high classification accuracy on a target dataset, many of these methods necessitate a substantial amount of annotated data. However, in many settings, annotating data is both time-consuming and costly, posing a challenge to the application of these successful methods in scenarios with limited resources. One of the ways to mitigate this problem is active learning. Active learning aims to maximize performance by selecting the most informative and valuable data points to be annotated for model training.

But how can a model correctly classify points of different classes? Classical support vector machines (SVMs) search for a hyperplane that separates two classes with the largest possible margin (see Figure 1). The points that lie on this decision boundary are called *support vectors*. In other words, these support vectors define the boundary of all samples of a class and are critical for a model to know in order to correctly separate the classes. We hypothesize that points close to the decision boundary are similarly relevant for neural network-based models.

In this work, we propose a novel active learning method (SUPClust) that tries to identify these points so that they can be annotated. Since the labels of the points are not known a priori, we rely on self-supervised representation learning in combination with clustering in order to break down the high-dimensional input space. For each cluster, we then identify the points close to a neighboring cluster, thereby selecting potential support vector points. Thanks to selecting points from all clusters, we ensure a broad coverage of the input space. In practice, data distributions often include outliers and the decision boundary between different classes is not always clearcut. For this reason, we further constrain our points to be somewhat typical according to a typicality metric introduced by Hacohen et al. (2022).



Figure 1: Decision boundary of an SVM classifier.

Our experimental evaluation demonstrates the merit of sampling points closer to the decision boundary, underscored by the strong performance compared to baseline active learning methods. SUPClust manages to not only mitigate the "cold start problem (Mittal et al., 2019)", it also shows strong performance in datasets with strong class imbalance. In ablation experiments, we ensure that all building blocks of SUPClust are necessary and contribute to the final result.
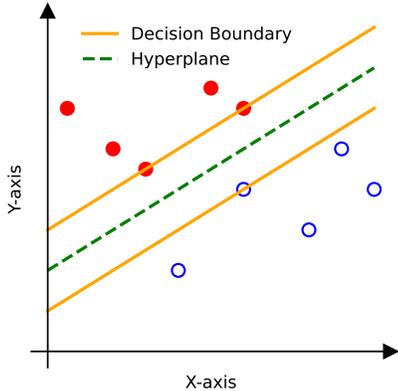
## 2   RELATED WORK

Various active learning methods have been proposed to this end, which can be categorized in uncertainty-based and diversity-based. Uncertainty-based approaches (Lewis & Gale, 1994; Joshi et al., 2009; Gal et al., 2017) leverage the prediction uncertainty of the classification model under training to select informative data samples for annotation. Diversity-based approaches (Sener & Savarese, 2017; Yehuda et al., 2022; Hacohen et al., 2022) aim to annotate a diverse range of samples spanning the complete data distribution, avoiding the selection of too similar ones. There also exist hybrid methods (Ash et al., 2019) which try to identify samples that have high uncertainty and are diverse at the same time. Some of these models rely on embeddings learned during self-supervised pre-training. Self-supervised learning involves training a model on a pretext task, allowing it to learn valuable representations without relying on explicit external labels. These representations complement the active learning task because they contain important information about the structure of the data distribution.

We give a short summary of the most used uncertainty-based approaches. Least confidence (Lewis & Gale, 1994), Entropy (Joshi et al., 2009), and Margin all select uncertain samples according to an uncertainty measure based on the output logits of the trained classifier. DBAL and BALD (Gal et al., 2017) on the other hand utilize Bayesian convolutional neural networks as a classifier and then select samples based on the highest entropy in the classifier or largest information gain. Many of these methods suffer from the "cold start problem", where their performance in low-budget regimes is worse than randomly selecting samples. This is possibly caused by the uncertainty estimates to be bad when the underlying model is not trained sufficiently due to limited labeled samples. SUPClust avoids this issue by selecting samples close to the decision border between clusters in the embedding space of a self-supervised pre-trained model.

In the realm of diversity-based methods, Coreset (Sener & Savarese, 2017) queries diverse samples through the selection of points that form a minimum radius cover of the remaining samples in the unlabeled pool. To do this, Coreset works on the embeddings generated by the penultimate layer of the classifier. In comparison, ProbCover (Yehuda et al., 2022) and TypiClust (Hacohen et al., 2022) rely on the embeddings of a self-supervised pre-trained model. ProbCover selects a maximum cover set for fixed-sized balls in this pre-trained embedding space. Typiclust builds clusters in the embedding space. From each cluster, it then selects the most typical sample. This combination ensures both broad coverage of the input space as well as selecting informative points, which shows in its strong performance in low-budget regimes. Typicality is measured in the following way:

$$Typicality(\boldsymbol{x}) = \left( \frac{1}{K} \sum_{\boldsymbol{x}_n \in K-\mathrm{NN}(\boldsymbol{x})} \|\boldsymbol{x} - \boldsymbol{x}_n\| \right)^{-1} \tag{1}$$

Here, $K - \mathrm{NN}(\boldsymbol{x})$ is a set of $K$ nearest neighbors of $\boldsymbol{x}$ in an embedding space. SUPClust also relies on typicality in order to ensure that the selected points are still representative of the cluster they come from.

## 3   SUPCLUST

SVM classifiers are defined by a few key points located at the decision boundary between the categories. Our querying strategy selects instances situated near the decision boundary, as they provide a strong signal to the learning process of neural network based models too. Traditional active learning methods have approached this problem by using model uncertainty as an indication for samples at the decision boundary. However, these methods suffer from the cold-start problem, where in low-budget scenarios, the model uncertainty is unable to identify hard instances. In this work, we introduce a novel method to find such samples by exploiting pre-trained representations. We can see in Figure 2 on the example of CIFAR-10 that similar categories are clustered together in the representation space. As category boundaries align with cluster boundaries, we use clustering to identify samples of interest. To quantify proximity to the decision boundary, we compute, for each sample, the weighted mean distance to all other cluster centers. The weights are the same for all samples within a cluster and are dependent on the distance of the cluster center to all other cluster centers. Clusters positioned at the "edge" of the data distribution select an instance that is close to
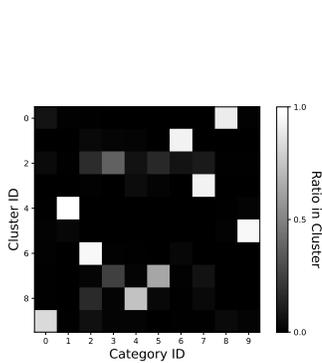
TypiClust
SUPClust



Figure 2: Distribution of classes within each cluster on SimCLR embeddings for CIFAR-10. Cluster boundaries align with category boundaries.
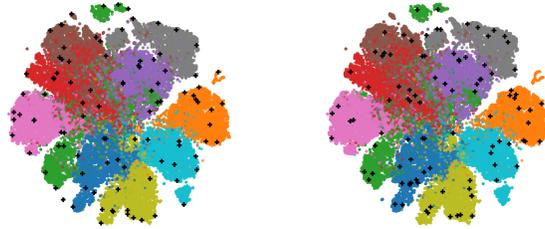
Figure 3: t-SNE plots of 100 queried instances by TypiClust and SUPClust (ours) in the CIFAR-10 embedding space. Colors represent the categories. For clusters on the "edge" of the data distribution, SUPClust tends to select samples that are closer to other clusters in the embedding space.

the nearest cluster. Conversely, clusters in the "middle" of the distribution do not select instances that are close to just one of the clusters. To normalize the weights to 1, we use the softmax function with the negative L2 distance as the logits and the temperature parameter $T$. For a point in cluster $i$, the weight to the cluster $j$ is given by

$$w_i^j = \frac{\exp\left(-\frac{\|\boldsymbol{c}_i - \boldsymbol{c}_j\|}{T}\right)}{\sum_{k \in C \setminus \{i\}} \exp\left(-\frac{\|\boldsymbol{c}_i - \boldsymbol{c}_k\|}{T}\right)}, \tag{2}$$

where $\boldsymbol{c}_i$, $\boldsymbol{c}_j$ and $\boldsymbol{c}_k$ are the centers of cluster $i$, $j$ and $k$ respectively, and $C$ is the set of all clusters. For cluster $i$, we select the sample $\boldsymbol{x}$, that has the minimum distance, or maximum $SUP$ to the decision boundary computed by Equation 3.

$$SUP(\boldsymbol{x}) = \left(\sum_{j \in C \setminus \{i\}} w_i^j \|\boldsymbol{x} - \boldsymbol{c}_j\|\right)^{-1} \tag{3}$$

Real data distributions are noisy, contain outliers, and can not be separated by a hyperplane, thus sampling only based on $SUP$ leads to subpar performance. To avoid outliers on the decision boundary we combine typicality with $SUP$. Typicality and $SUP$ are not correlated, see Figure 4, thus using both metrics for sample selection can improve performance.
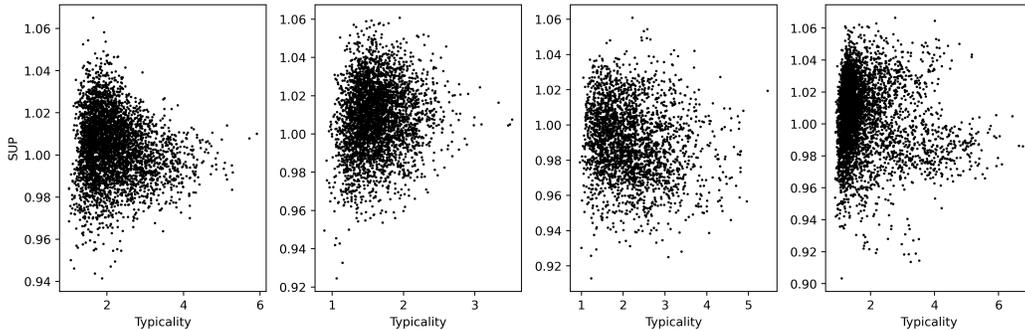


Figure 4: Relationship between typicality and SUP on CIFAR-10 on 4 randomly selected clusters, with temperature 1. Typicality and SUP have no strong correlation, using both metrics to select instances can improve the querying strategy.

Our proposed strategy **SUPClust** consists of 4 parts. 1) Train a self-supervised model on the unlabeled pool. 2) Partition the data into $N$ clusters, where $N$ is the number of labeled samples after the end of the current step. Excluding clusters that contain samples from the already labeled pool, selecting as many of the biggest clusters as samples are queried. 3) In each cluster filter the top 10% of samples based on typicality. 4) Select the sample with the highest $SUP$. In Figure 3 we compare samples queried by TypiClust and and SUPClust. Notably, TypiClust chooses more samples from the "edge" of the data distribution, whereas SUPClust prioritizes samples that lie closer to other categories.

## 4  RESULTS

### 4.1  EXPERIMENTAL SETUP

All strategies are evaluated on image classification tasks using CIFAR-10, CIFAR-100 (Krizhevsky, 2009), CIFAR-10-LT (Cao et al., 2019), and ISIC-2019 (Kassem et al., 2020). CIFAR-10 and CIFAR-100 consist of 60k natural images of size 32x32 with 10/100 classes. CIFAR-10-LT is a class-imbalanced subset of CIFAR-10. We apply an imbalanced factor of 50, meaning a 50-fold difference in the number of images between the most and least frequent class. ISIC-2019 consists of 25331 skin cancer images with 8 imbalanced classes. To standardize the image dimensions, all images are resized to 224x224 pixels. In alignment with TypiClust, we adopt tiny and small budget sizes, involving querying step sizes 1 and 5 times the number of classes respectively.

We evaluate AL strategies in the following two frameworks. 1) Fully supervised (FSL): training a deep neural network, ResNet18 (He et al., 2015), exclusively on the labeled set which is acquired by active queries. 2) Fully supervised with self-supervised embedding (SSL): training a linear classifier on the labeled embeddings obtained by active queries. These self-supervised embeddings for the classifier are obtained from a pre-trained SimCLR (Chen et al., 2020). Within these frameworks, we compare SUPClust to nine baseline strategies: Random, Margin, Least confidence, Entropy, BALD, Coreset, DBAL, TypiClust, ProbCover. For the clustering and sampling with TypiClust and SUPClust, we use SimCLR representations, namely the ResNet18 backbone for CIFAR-10, CIFAR-10-LT50 and ISIC-2019, and the ResNet34 for CIFAR-100.

### 4.2  ABLATION STUDY

To assess the significance of individual components within SUPClust, we conduct ablation experiments for each component. We display the results in Figure 5. When leaving out our SUP-based acquisition metric (SUPClust w/o SUP) and instead selecting a sample randomly from the top 10% typical samples within each cluster, the performance noticeably declines, falling below that of TypiClust. Similarly, relying solely on SUP without considering typicality for sample selection (SUPClust w/o typicality) fails to achieve the performance levels observed with other querying strategies. As a comparison, we also show the default TypiClust (typiclust-rp), which always selects the most typical sample of a cluster. Our results showcase that all components of SUPClust are necessary and contribute to its performance.
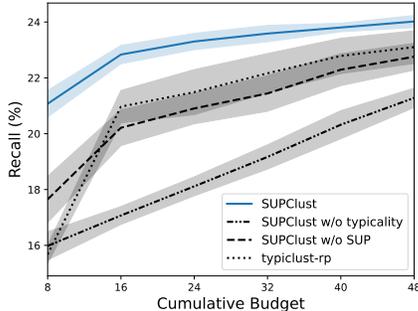


Figure 5: Ablation study on ISIC-2019 with budget=8 and with self-supervised embeddings

### 4.3  MAIN RESULTS

We present the main results of our evaluation in Figure 6 for the tiny and the small budget regime. We can see that SUPClust performs well on all evaluated datasets, especially in imbalanced settings. On CIFAR10-LT50 and ISIC2019, SUPClust demonstrates a strong performance gain compared to TypiClust. We hypothesize that by selecting points according to maximum $SUP$, SUPClust is
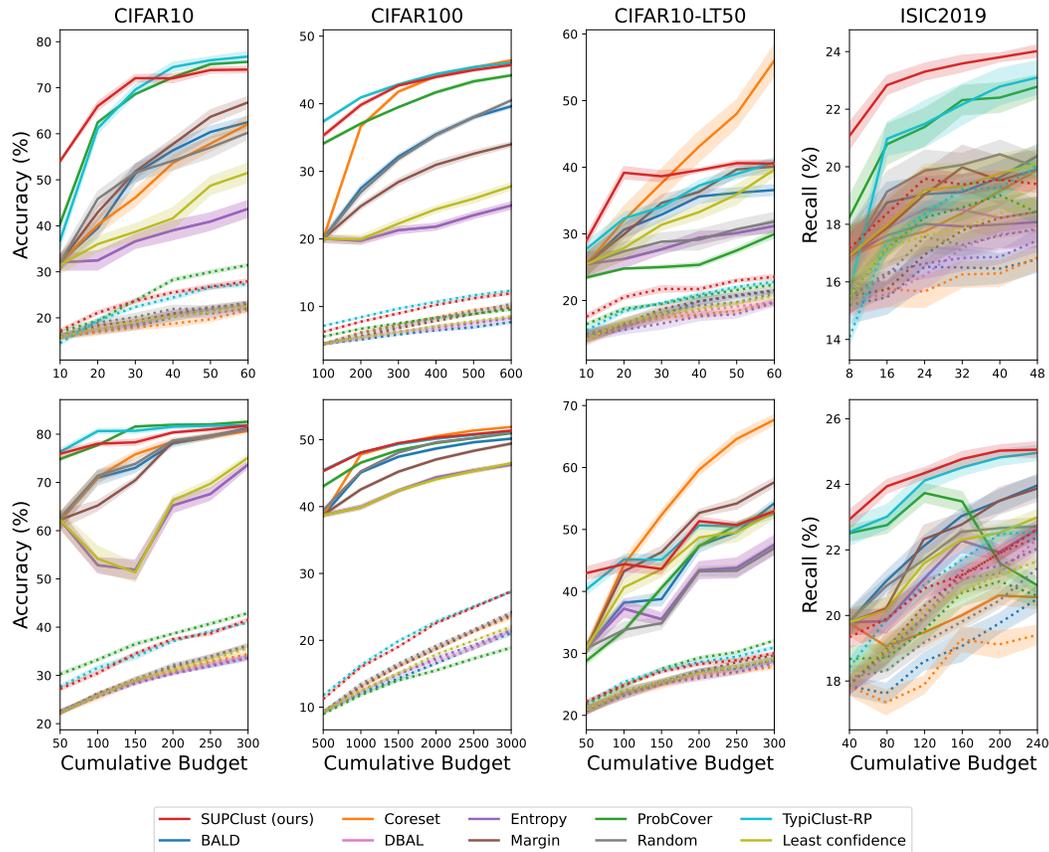
Figure 6: Results in the tiny (top) and small (bottom) budget regime. Solid lines represent results with the SSL setting, and dotted lines represent results with the FSL setting. The mean and the standard error with 10 different random seeds are shown. Our method (SUPClust) shows robust performance compared to other baselines, across all datasets and both data regimes.

able to select more informative points relevant for distinguishing the classes, irrelevant of imbalance. In our low-budget regimes, diversity-based methods such as TypiClust, Coreset and Prob-Cover generally perform better than their uncertainty-based counterparts. This is to be expected, as uncertainty-based methods bring stronger benefits only in higher budget regimes. Building on the self-supervised pre-trained embeddings improves performance across all datasets. The performance of Coreset on CIFAR10-LT50 in the SSL setting is surprising. The embeddings of the pre-trained backbone allow Coreset to select very informative samples. Unfortunately, when training in the FSL setting or on any other dataset, the performance of Coreset is diminished compared to other algorithms.

## 5 DISCUSSION

Active learning can bring performance benefits to settings where acquiring labeled data is expensive. Samples close to the decision boundary between categories provide a strong training signal. The introduction of the novel $SUP$ metric provides a non-label-based means of quantifying the distance of each sample to the decision boundary. Utilizing $SUP$, when selecting which samples to label for classifier training improves sample efficiency, especially in the low data budget regime. Our findings contribute to the broader understanding of active learning dynamics, shedding light on the relationship between the $SUP$ metric, typicality, and diversity. Exploring the changing relationship between diversity, typicality and the $SUP$ metric across various data regimes remains future work.

REFERENCES

Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. June 2019.

Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *Adv. Neural Inf. Process. Syst.*, 32, 2019.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1597–1607. PMLR, 2020. URL http://proceedings.mlr.press/v119/chen20j.html.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL https://openreview.net/forum?id=YicbFdNTTy.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1183–1192. PMLR, March 2017.

Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies suit high and low budgets. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 8175–8195. PMLR, July 2022.

Kaiming He, X Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, pp. 770–778, December 2015.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 448–456. JMLR.org, 2015. URL http://proceedings.mlr.press/v37/ioffe15.html.

Ajay J. Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. Multi-class active learning for image classification. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA*, pp. 2372–2379. IEEE Computer Society, 2009. doi: 10.1109/CVPR.2009.5206627. URL https://doi.org/10.1109/CVPR.2009.5206627.

Mohamed A Kassem, Khalid M Hosny, and Mohamed M Fouad. Skin lesions classification into eight classes for ISIC 2019 using deep convolutional neural network and transfer learning. *IEEE Access*, 8:114822–114832, 2020.

Alex Krizhevsky. Learning multiple layers of features from tiny images. pp. 32–33, 2009. URL https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.

David D. Lewis and William A. Gale. A sequential algorithm for training text classifiers. In W. Bruce Croft and C. J. van Rijsbergen (eds.), *Proceedings of the 17th Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval. Dublin, Ireland, 3-6 July 1994 (Special Issue of the SIGIR Forum)*, pp. 3–12. ACM/Springer, 1994. doi: 10.1007/978-1-4471-2099-5. URL https://doi.org/10.1007/978-1-4471-2099-5.

Sudhanshu Mittal, Maxim Tatarchenko, Özgün Çiçek, and Thomas Brox. Parting with illusions about deep active learning. December 2019.

Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A Core-Set approach. August 2017.

Siddharth Srivastava and Gaurav Sharma. Omnivec: Learning robust representations with cross modal sharing. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 1236–1248, January 2024.

Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering lens. In S Koyejo, S Mohamed, A Agarwal, D Belgrave, K Cho, and A Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 22354–22367. Curran Associates, Inc., May 2022.