

# BRIDGING DIVERSITY AND UNCERTAINTY IN ACTIVE LEARNING WITH SELF-SUPERVISED PRE-TRAINING

Paul Doucet, Benjamin Estermann, Till Aczel & Roger Wattenhofer

ETH Zürich

{pdoucet, estermann, taczel, wattenhofer}@ethz.ch

## ABSTRACT

This study addresses the integration of diversity-based and uncertainty-based sampling strategies in active learning, particularly within the context of self-supervised pre-trained models. We introduce a straightforward heuristic called **TCM** that mitigates the cold start problem while maintaining strong performance across various data levels. By initially applying TypiClust for diversity sampling and subsequently transitioning to uncertainty sampling with Margin, our approach effectively combines the strengths of both strategies. Our experiments demonstrate that TCM consistently outperforms existing methods across various datasets in both low and high data regimes.

## 1 INTRODUCTION

Training machine learning models are known to depend on a lot of labeled data. However, in many settings labeled data is not easy to acquire, but has to be created through expensive manual labeling. The goal of active learning is to address this challenge by providing a way to select the most informative samples for labeling. These are the samples for which training a classifier on them increases performance the fastest.

Depending on how many labeled samples the model has already been trained on, the properties of these samples have to change. In a low-budget setting, diverse and typical samples that cover the complete data distribution are most important. As the budget increases, the model does not profit from such samples as much anymore. Instead, once a model has learned the general rules of the data distribution, the most informative samples are now the ones that show where exactly the decision boundaries are. The issue is that different active learning methods work best in different data budgets and it is not clear when to switch between the methods.

In this work, we provide a new heuristic called TCM on how to combine two such methods, namely TypiClust (Hacohen et al., 2022) and Margin. TypiClust shows excellent performance in low data regimes, while Margin excels afterwards. We specifically analyze the setting where a self-supervised pre-trained backbone model is available. Such a backbone not only massively increases performance compared to training from scratch, but also the transition dynamics from low to high data regimes simplify.

We show that TCM achieves consistently strong performance, regardless of the labeling budget and the dataset (see Figure 1). It outperforms its underlying methods TypiClust and Margin during the complete training process. Thanks to the simplicity and effectiveness of our TCM, we provide clear guidelines for practitioners on how to easily use active learning in their own setting.

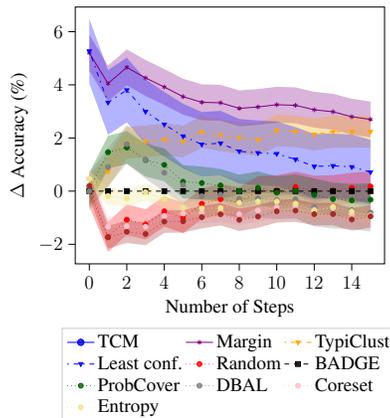


Figure 1: Accuracy improvement compared to random for all baselines and our (TCM) strategy. The accuracy improvement mean and standard deviation is computed over all budget sizes for CIFAR10, CIFAR100 and ISIC2019.

## 2 RELATED WORK

In active learning, two primary sampling strategies emerge as critical: diversity-based and uncertainty-based sampling. Diversity-based sampling aims to select a representative set of samples that span the entire feature space, thus ensuring broad coverage of the input domain. In contrast, uncertainty-based sampling focuses on querying instances for which the model exhibits the highest prediction uncertainty. In this way, uncertainty-based sampling aims to refine the model’s performance in challenging cases. Early on in the training, diversity-based methods tend to perform better as with limited samples it is harder to cover the complete data distribution. Further, in that stage, the classifier uncertainty is a weak predictor of hard samples. This is also called the “cold start problem” (Mittal et al., 2019) of uncertainty-based methods.

### 2.1 SELF-SUPERVISED PRE-TRAINING

Before we introduce any specific active learning methods, we address self-supervised learning. In self-supervised learning, a model is trained using a pretext task, allowing it to learn useful representations without explicit external labels. Self-supervised learning complements active learning by learning embeddings before any data is labeled, simplifying both the sample quivering and classifier training. Popular and commonly used models include SimCLR (Chen et al., 2020) and DINO (Caron et al., 2021). SimCLR uses a contrastive loss, where different views of an image are pushed to have a close representation and views of different images a distant representation. The DINO pretext task involves learning data transformation invariant representations, by distilling different views of the same image from a teacher to a student network, where the teacher network is an exponential moving average of the student network.

### 2.2 UNCERTAINTY-BASED ACTIVE LEARNING

We first briefly introduce relevant uncertainty-based active learning methods. Some methods try to quantify model uncertainty based on the output logits of a classifier. **Least confidence** (Lewis & Gale, 1994) selects samples where the highest class probability is the lowest. **Entropy** (Joshi et al., 2009) measures model prediction uncertainty by the classifier probability distribution and selects samples with the highest entropy. **Margin** selects samples for which the difference between the class probabilities of the most likely two classes is the lowest. There also exist methods based on a Bayesian approach. **DBAL** (Gal et al., 2017) uses Bayesian convolutional neural networks as the classifier and queries samples based on their highest entropy. **BALD** (Gal et al., 2017) also uses Bayesian convolutional neural networks, but in contrast to DBAL, it selects samples that maximize information gained about the model parameters. Our TCM strategy builds on top of Margin for uncertainty sampling due to its strong performance and simple design.

### 2.3 DIVERSITY-BASED ACTIVE LEARNING

In the realm of diversity-based active learning, **Coreset** (Sener & Savarese, 2018) queries diverse samples by selecting points that form a minimum radius cover of the remaining samples in the unlabeled pool. The minimum radius cover ensures that all remaining unlabeled sample has a nearby sample that gets labeled. **ProbCover** (Yehuda et al., 2022) improves on Coreset by building on top of self-supervised embeddings and selecting samples of high-density regions of the embedding space. While Coreset ensures that it queries samples from the whole distribution, it is likely to select outliers that do not benefit the training. In contrast, ProbCover (Yehuda et al., 2022) samples from uncovered high-density regions, selecting more representative samples. **TypiClust** (Hacohen et al., 2022) queries diverse samples by first clustering, and then selecting the most typical sample from each cluster. The number of clusters increases by the sampling size at each step, ensuring that there are enough clusters to sample new points from unexplored regions of the embedding space. Typicality is measured by the inverse average distance to other points in the cluster. The typical points queried by TypiClust enable strong performance, especially at the beginning of the training process. For this reason, our TCM strategy utilizes TypiClust as the initial sampling strategy, thereby avoiding the cold-start problem.

## 2.4 HYBRID METHODS

Some works in active learning have previously combined diversity and uncertainty-based sampling. BADGE (Ash et al., 2020) and BatchBALD (Kirsch et al., 2019) developed methods to ensure diversity within a batch of uncertain samples. SelectAL (Hacohen & Weinshall, 2023) on the other hand developed a complex algorithm that automatically detects the current data regime and selects a corresponding diversity or uncertainty-based active learning algorithm. However, these methods mostly focused on the case where the classifier is trained from scratch. As we show later in this work, when using a pre-trained backbone, there is no need for a complex switching strategy, as the transition point between diversity and uncertainty-based sampling always occurs early on in the training.

## 3 METHODOLOGY

The best querying strategy depends on the already labeled dataset size. As can be seen in Figure 1, a strong diversity-based method such as TypiClust (Blue) performs strongly in the first few steps of sampling. On the other hand, an uncertainty-based method such as Margin (Yellow) shows stronger performance the larger the cumulative budget grows. The best and most consistent performance can be achieved when utilizing both methods when they perform the strongest. We therefore propose a hybrid sampling strategy combining TypiClust and Margin and call this strategy **TCM**.

### 3.1 TRANSITION POINT

To get greater insights into the dynamics of when to transition from TypiClust to Margin, we perform an ablation starting with TypiClust and switching to Margin after  $N$  sampling steps. The results of this ablation, displayed in Figure 2, show that the optimal transition point depends on the initial budget and corresponding step size. The larger the initial budget, the better it is to quickly switch from TypiClust to Margin to achieve consistently strong performance.

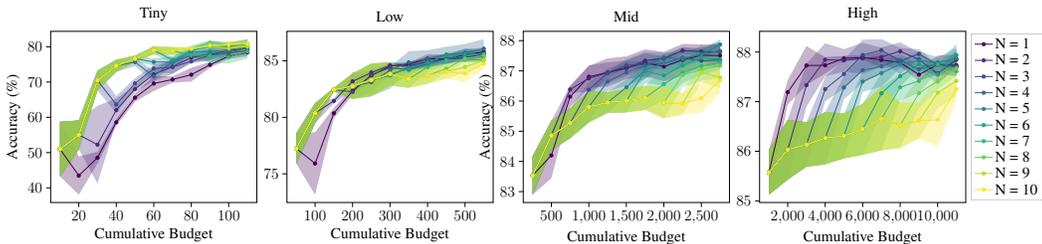


Figure 2: Transition point ablation on the CIFAR10 dataset. Switching to Margin in the last step is equal to  $N = 10$ , while only using TypiClust for the initial sampling, and switching to Margin imminently is  $N = 1$ .

### 3.2 STEP SIZE

We further analyze the effect of the step size for the Margin part of TCM. Thanks to the design of TypiClust, its performance is mostly independent of step size. For this reason, we only perform detailed analysis on Margin, after having selected the initial budget with TypiClust. By just selecting a batch of samples depending on the class probabilities, Margin might select a lot of similar samples if the step size is too big. At the same time, too small of a step size is not practical, as the model would need to be retrained too often. For this reason, we perform an ablation on the effect of the step size on the performance of TCM, shown in Figure 3. Surprisingly, the results show that overall, there is no clear difference in performance for different step sizes. A large step size for the Mid and High budget might have a slightly negative impact, however towards the end of the cumulative budget, the difference disappears. The implications of these results are strongly positive, as they indicate that the step size of TCM can be adapted to other needs such as the availability of experts for labeling data and computational resources.

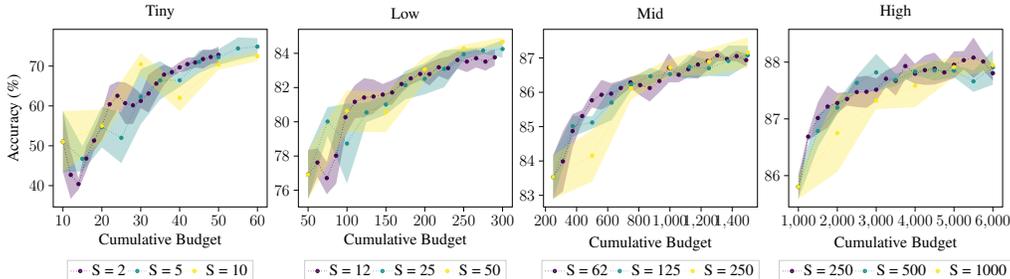


Figure 3: Step size ablation for TCM on the CIFAR10 dataset. For each regime, we evaluate three different step sizes  $S$ . Overall, there is no clear performance difference between the step sizes.

Based on our ablations, we devise the following simple heuristic for TCM. We use a step size equal to the size of the initial budget of each setting. In the Tiny and Low settings, we perform 3 steps of TypiClust before switching to Margin. In the Medium setting, we perform 2 steps of TypiClust and in the High setting, we perform a single step.

However, based on our ablations we can also provide a broader rule of thumb for an active learning practitioner. In a new setting, we suggest using a total budget of roughly 20 times the number of categories for TypiClust and then switching to Margin, with no strict constraint on step size.

## 4 EXPERIMENTAL EVALUATION

### 4.1 SETUP

Today’s active learning research has an inconsistent landscape with many contradicting results. For this reason, we follow the evaluation framework proposed by Lüth et al. (2023) to have a fair and rigorous comparison to all active learning baselines. In particular, we tune the hyperparameters on the CIFAR10 (Krizhevsky, 2009) dataset and use all other datasets as rollout datasets. We evaluate the datasets CIFAR10 and CIFAR100 (Krizhevsky, 2009) and ISIC2019 (Codella et al., 2017), but also consider the long-tail versions of the two datasets, denoted with LT (Cao et al., 2019). These long-tail versions feature class imbalance, where the number after LT indicates the ratio between the most common and the least common class, combined with an exponential decay in sample sizes for the other classes. We perform all our experiments on 4 different data budget sizes. These include small, medium, and large as defined by Lüth et al., as well as the tiny budget size with the initial sample size equal to the number of classes. Further, we consistently use a step size equal to the number of initial samples for all budgets. Some sampling methods such as ProbCover, TypiClust and our proposed TCM rely on self-supervised pre-trained representations. For all these methods, we use the same representations. Specifically, for CIFAR10, CIFAR10-LT5, CIFAR10-LT10, CIFAR100, CIFAR100-LT5, CIFAR100-LT10 we use the SimCLR features provided by Hacoen et al.. For ISIC2019, we train a DINO model from scratch and use the features of the model.

To underpin the efficacy of our proposed active learning framework, we leverage the same pre-trained representations for training our classifier. For all experiments conducted within this work, we utilize the selected backbone as a foundational feature extractor. On top, we train a linear prediction head, fine-tuning the model to the labels gained during the active learning process. This approach allows us to harness the rich representational capacity of the pre-trained models, leading to a much stronger performance of the classifier especially in a low-data regime. Furthermore, by using publicly available pre-trained model embeddings, the classifier training following each sampling step demands significantly fewer computational resources. Due to the above mentioned advantages of using a pre-trained backbone and considering that TypiClust relies on such a backbone, we do not perform any experiments where we train a classifier from scratch. We fix all training hyperparameters for classifier training for all datasets, budgets, and querying strategies. For all experiments, we plot the mean and standard deviation of runs with 3 seeds. While we have already presented an aggregation over all datasets and budget sizes in Figure 1, we present specific results for each evaluated dataset in Figure 4.

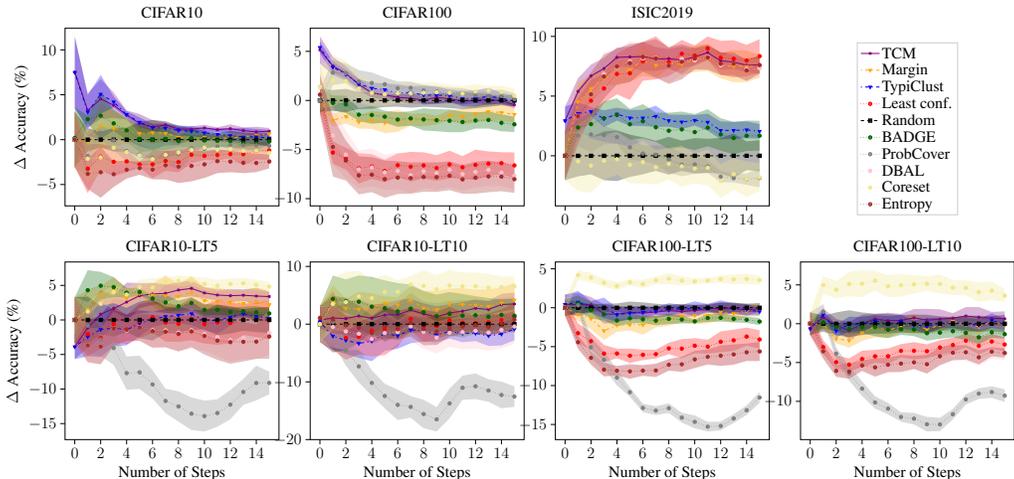


Figure 4: Accuracy improvement compared to random for all baselines and our (TCM) strategy. The accuracy improvement mean is computed over all 4 budget sizes tiny, small, medium, and large. Standard deviation is aggregated with respect to the random seed. The top row shows the main evaluated datasets, while the bottom row shows an ablation on the imbalanced versions of CIFAR10 and CIFAR100. For all imbalanced datasets, reported accuracy is balanced by computing the average of recall obtained for each class. TCM shows consistently strong performance for all datasets, even for datasets for which TypiClust or Margin on their own show suboptimal performance. Coreset shows strong performance on the LT datasets. Unfortunately, this performance does not transfer to the real-life imbalanced dataset ISIC2019.

## 4.2 RESULTS

Our results showcase the consistent and strong performance of TCM compared to other baselines. This performance gain transfers to the imbalanced long-tail versions of CIFAR10 and CIFAR100. It can be seen that TCM performs well even if either TypiClust or Margin do not perform well on their own. In CIFAR10-LT5 as well as in CIFAR100-LT10, TCM performs better than both its underlying methods, showcasing the potential of the combination. On the main datasets CIFAR10, CIFAR100 and ISIC2019, Coreset shows performance that is worse than selecting samples at random, but shows strong performance in the LT ablation. We hypothesize that this could be because, for the LT datasets, the backbone was pre-trained on the balanced version of the dataset. This is not the case for ISIC2019, where the DINO backbone was trained entirely on ISIC2019. Other methods such as ProbCover, DBAL, or Least confidence seem to struggle a lot with providing consistent performance in all data regimes and datasets. Furthermore, our results show that when using a pre-trained backbone, more complex methods such as SelectAL do not seem necessary, because the transition point occurs much earlier compared to training from scratch. This fact is exploited by TCM and allows for its strong and consistent performance for all budget sizes. Switching between diversity and uncertainty-based strategies after the few first iterations consistently outperforms other methods.

## 5 CONCLUSION

In this work, we have shown that when training a classifier using a pre-trained backbone, the transition point from diversity to uncertainty-based active learning methods occurs early. Based on these results, we present TCM, a simple, yet effective hybrid strategy that outperforms all other methods when compared on a wide range of data regimes and datasets. By using TypiClust for the first few steps and then switching to Margin, TCM selects informative instances in both low and high data regimes. Using the simple heuristics laid out by TCM, practitioners can apply active learning easily and effectively to their use case.

## REFERENCES

- Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=ryghZJBKPS>.
- Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Aréchiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 1565–1576, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/621461af90cadfdaf0e8d4cc25129f91-Abstract.html>.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 9630–9640. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00951. URL <https://doi.org/10.1109/ICCV48922.2021.00951>.
- 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>.
- Noel C. F. Codella, David A. Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin K. Mishra, Harald Kittler, and Allan Halpern. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (ISIC). *CoRR*, abs/1710.05006, 2017. URL <http://arxiv.org/abs/1710.05006>.
- 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, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1183–1192. PMLR, 2017. URL <http://proceedings.mlr.press/v70/gall17a.html>.
- Guy Hacohen and Daphna Weinshall. How to select which active learning strategy is best suited for your specific problem and budget. *CoRR*, abs/2306.03543, 2023. doi: 10.48550/ARXIV.2306.03543. URL <https://doi.org/10.48550/arXiv.2306.03543>.
- Guy Hacohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies suit high and low budgets. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 8175–8195. PMLR, 2022. URL <https://proceedings.mlr.press/v162/hacohen22a.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>.
- Andreas Kirsch, Joost van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 7024–7035, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/95323660ed2124450caaac2c46b5ed90-Abstract.html>.

- 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\_1. URL [https://doi.org/10.1007/978-1-4471-2099-5\\_1](https://doi.org/10.1007/978-1-4471-2099-5_1).
- Carsten T. Lüth, Till J. Bungert, Lukas Klein, and Paul F. Jaeger. Toward realistic evaluation of deep active learning algorithms in image classification. *CoRR*, abs/2301.10625, 2023. doi: 10.48550/ARXIV.2301.10625. URL <https://doi.org/10.48550/arXiv.2301.10625>.
- 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. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=H1aIuk-RW>.
- Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering lens. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/8c64bc3f7796d31caa7c3e6b969bf7da-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/8c64bc3f7796d31caa7c3e6b969bf7da-Abstract-Conference.html).