
NTKCPL: Active Learning on Top of Self-Supervised Model by Estimating True Coverage

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 High annotation cost has driven extensive research in active learning and self-
2 supervised learning. Recent research has shown that in the context of supervised
3 learning, when we have different numbers of labels, we need to apply different
4 active learning strategies to ensure that it outperforms the random baseline. This
5 number of annotations that change the suitable active learning strategy is called the
6 phase transition point. We found, however, when combining active learning with
7 self-supervised models to achieve improved performance, the phase transition point
8 occurs earlier. It becomes challenging to determine which strategy should be used
9 for previously unseen datasets. We argue that existing active learning algorithms are
10 heavily influenced by the phase transition because the empirical risk over the entire
11 active learning pool estimated by these algorithms is inaccurate and influenced by
12 the number of labeled samples. To address this issue, we propose a novel active
13 learning strategy, neural tangent kernel clustering-pseudo-labels (NTKCPL). It
14 estimates empirical risk based on pseudo-labels and the model prediction with
15 NTK approximation. We analyze the factors affecting this approximation error and
16 design a pseudo-label clustering generation method to reduce the approximation
17 error. Finally, our method was validated on five datasets, empirically demonstrating
18 that it outperforms the baseline methods in most cases and is valid over a wider
19 range of training budgets.

20 1 Introduction

21 The boom in deep learning models in recent years stems in part from the massive amounts of
22 data [11, 17, 23]. However, the demand for large amounts of data, especially labeled data, in
23 turn, constrains the application of deep learning models, since large amounts of labels imply high
24 annotation costs [41, 1, 45]. Active learning is a path to alleviate the cost of labeling by selecting
25 informative subsets of samples to annotate.

26 However, the benefits of active learning have been increasingly questioned in recent years [25, 28].
27 One of the main concerns is that training a model initialized by self-supervised learning with randomly
28 selected labeled samples often yields results far beyond those obtained by existing active learning
29 with supervised training (randomly initialized or initialized by the last round of the active learning
30 model) [6, 8, 7, 14, 9]. Because the latter only uses labeled data to train the network, while the
31 former uses a large amount of unlabeled data to train the backbone of the network. Since most
32 existing active learning algorithms are designed in the context of supervised training, they must be
33 validated with a large number of labels compared to the number of labels required in training from a
34 self-supervised model. This means that the effectiveness of these active learning algorithms is not
35 guaranteed in the case of having access to relatively few annotations, as is the case when combining
36 with a self-supervised model. Several studies [15, 42, 4] have shown that many existing active

37 learning strategies fail to outperform the random baseline when combining them with self-supervised
38 learning. In this paper, we focus on designing an active learning strategy that works well in the
39 training method with a self-supervised model.

40 The “phase transition” phenomenon [15] is known to occur in active learning with supervised training.
41 It refers to the fact that an active learning strategy that outperforms a random baseline when the total
42 number of labels is small will be inferior to a random baseline when the total number of labels is
43 large (called the **low-budget** strategy) and vice versa (called the **high-budget** strategy). We note that
44 when combining active learning with the self-supervised model, the cut-off point between low-budget
45 and high-budget strategy occurs much earlier. For example, in the CIFAR-100 [21], the cut-off point
46 is about 10,000 labeled samples when training in the supervised learning way [16]. But, the cut-off
47 point shifts forward to about 1,500 labeled samples when training from a self-supervised model. The
48 forward-moving cut-off point means that even if the annotation budget is low (only one order of
49 magnitude above the number of classes in the dataset), it is likely to hit that cut-off point. Thus, for a
50 previously unseen dataset, it is difficult to simply determine whether a low-budget or high-budget
51 strategy should be chosen since the difficulty varies from dataset to dataset. In this paper we use this
52 problem to motivate the design of an active learning strategy with a wider effective budget range.

53 Since existing low-budget strategies are designed based on the idea of feature space coverage [24, 15,
54 42], we first analyze the problems of determining coverage based on sample feature distances in sec. 2.
55 After that, we propose that the true coverage where the empirical risk is zero, can be estimated based
56 on pseudo-labels and predictions of the model trained on the candidate set. Based on this, we propose
57 our active learning strategy, Neural Tangent Kernel Clustering-Pseudo-Labels (NTKCPL), which
58 uses the NTK [18, 27] and CPL to approximate empirical risk on active learning pool in sec. 3.2. And
59 we analyze which factor affects approximation error in sec. 3.3. Based on this analysis, we design
60 a CPL generation method in sec. 3.4. Extensive experimental results demonstrate that our method
61 outperforms state-of-the-art approaches in most cases and has a wider effective budget range. As part
62 of the results (sec. 4) we also show our method is effective for self-supervised features of different
63 quality.

64 Our contribution is summarized as follows: (1) We propose a novel active learning strategy, NTKCPL,
65 by estimating empirical risk on the whole active learning pool based on pseudo-labels. (2) We analyze
66 the approximation error of the empirical risk in the active learning pool when NTK and CPL are used
67 to approximate networks and true labels. (3) Our method outperforms both low- and high-budget
68 active learning strategies within a range of annotation quantities one order of magnitude larger than
69 traditional low-budget active learning experiments. This means that our approach can be used more
70 confidently for active learning on top of self-supervised models than existing low-budget strategies.

71 1.1 Related Work

72 Most active learning strategies are designed and validated in the high-budget scenario where network
73 weights are randomly initialized or initialized from the weights of the previous active learning
74 round. Active learning methods mainly include uncertainty-based sampling [22, 13, 19], feature
75 space coverage [32, 24, 42, 33, 5, 40], the combination of uncertainty and diversity [41, 3], learning-
76 based methods [43], and so on [34, 35]. Moreover, some recent studies explore “**look ahead**”
77 strategies [26, 38], where samples are selected based on the model trained on candidate training sets.
78 However, with the development of self-supervised training, the training approach for low-budget
79 scenarios has shifted to training based on a self-supervised pre-trained model [24]. This change in
80 the training method implies a shift in the total number of samples that need to be selected by active
81 learning. When training based on a self-supervised model, often only 0.4-6% of the total data needs
82 to be labeled to achieve similar results to training with 20-40% labeled data on a randomly initialized
83 network [4]. Recent studies have shown that there exists a phase transition phenomenon in active
84 learning strategies, whereby opposite strategies should be adopted in high-budget and low-budget
85 scenarios [15], causing many active learning strategies designed for high-budget scenarios unsuitable
86 for training based on a self-supervised model. As a result, recent studies have explored active learning
87 strategies specifically designed for low-budget scenarios [15, 42, 31, 20]. However, we find that these
88 strategies are effective only when the number of labeled data samples is extremely small, and as we
89 increase the labeled data to one order of magnitude above the number of classes of the dataset, their
90 performance falls below that of the random baseline.

91 **2 Insight: Distance is not an accurate indicator of empirical risk**

92 The goal of the active learning is to find a labeled subset, $D_C = (x_i, y_i)_{i=1}^{N_C}$, such that the model
 93 trained on that subset, f_{D_C} , has the minimized empirical risk in the entire active learning pool,
 94 $D = (x_i, y_i)_{i=1}^N$ as shown in eq. 1.

$$\operatorname{argmin}_{D_C} \frac{1}{N} \sum_{i \in D} \operatorname{Loss}(f_{D_C}(x_i), y_i) \quad (1)$$

95 Unfortunately, during active learning, we do not have the labels of the entire active learning pool, so
 96 we cannot compute this loss directly. To address this problem, current methods [32, 24, 33] covert
 97 empirical risk minimization into feature space coverage based on Lipschitz continuity. Although
 98 Lipschitz continuity guarantees that the difference between the model’s predictions is less than the
 99 product of the Lipschitz constant and the difference between inputs, it does not guarantee that their
 100 predictions fall into the same class. In practice, we cannot determine the true coverage because we
 101 do not know the distance threshold beyond which the model would change its predicted class for
 102 unlabeled samples.

103 Therefore, the current solution is to minimize the coverage radius assuming full coverage [32] or to
 104 maximize coverage based on high purity coverage [42], where purity refers to the probability that the
 105 sample has the same label within a given distance. Assuming full coverage leads to an overestimated
 106 coverage as shown in fig. 1a, i.e., some covered samples still have a large empirical risk, while high-
 107 purity coverage causes underestimated coverage as shown in fig. 1b. The overestimated coverage
 108 may cause the active learning algorithm to miss samples in areas that are not truly covered, while
 109 underestimated coverage makes active learning algorithms likely to select redundant samples. These
 110 affect the performance of active learning.

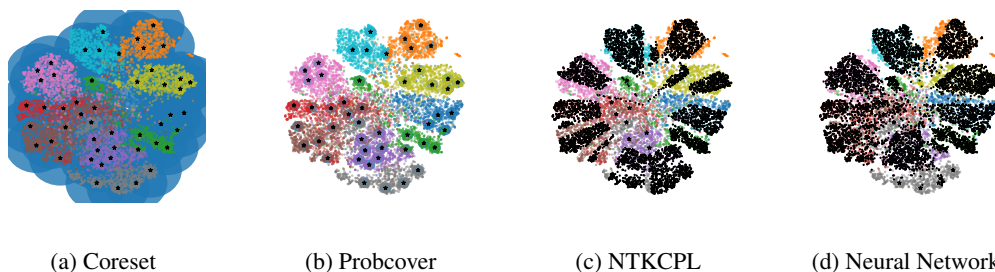


Figure 1: Coverage estimation based on sample feature distance vs. NTKCPL. Here different colors represent different categories, the black star denotes labeled samples and the blue circle represents the samples considered covered based on the feature distance approach. Coreset assumes full coverage and Probcov assumes high purity coverage. The coverage estimated by our method, NTKCPL, and true coverage based on predictions of the neural network is represented by black dots. The coverage estimated by NTKCPL is more consistent with the true coverage of the neural network than those estimated based on feature distances.

111 Additionally, estimating the empirical risk based on distance implies the assumption that model
 112 predictions are only relevant to the nearest labeled sample, which is often not the case in reality. To
 113 estimate the true coverage, we propose a new strategy, NTKCPL. It estimates the empirical risk based
 114 on the predictions of the model trained on the candidate set and pseudo-labels.

115 **3 Method: NTKCPL**

116 In sec. 3.1, we briefly review the Neural Tangent Kernel (NTK) [18] that enables active learning
 117 strategies based on the outputs of a model trained on a candidate set feasible. Then, we propose our
 118 active learning strategy, NTKCPL, in sec. 3.2 and analyze the approximation error of NTKCPL in
 119 sec. 3.3. Finally, based on the analysis, we introduce the method of generating cluster pseudo-label in
 120 sec. 3.4.

121 **3.1 Preliminaries**

122 Neural Tangent Kernel (NTK) is a powerful tool to analyze the training dynamics of neural network.
 123 Jacot et al. [18] show that the neural network is equivalent to the kernel regression with Neural
 124 Tangent Kernel when network is sufficiently wide and its weights are initialized properly [2]. The
 125 NTK, \mathcal{K} , is shown in eq. 2, where the f denotes a neural network with parameters θ and \mathcal{X} denotes
 126 train samples. When training with MSE loss, the neural network has a closed-form solution for the
 127 prediction of test sample x at iteration t as eq. 3, where \mathcal{Y} denotes labels of trainset and f_0 denotes
 128 the output of network with initialized weights.

$$\mathcal{K}(\mathcal{X}, \mathcal{X}) = \nabla_{\theta} f(\mathcal{X}) \nabla_{\theta} f(\mathcal{X})^T \quad (2)$$

$$f_t(x) = f_0(x) + \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} (\mathcal{I} - e^{-t\mathcal{K}(\mathcal{X}, \mathcal{X})}) (\mathcal{Y} - f_0(\mathcal{X})), \quad (3)$$

129 Additionally, for active learning scenarios, Mohamad [26, 27] proposes the computation time of
 130 using NTK can be further reduced by considering the block structure of the matrix, which means
 131 that look ahead type active learning strategies can be implemented in a reasonable amount of time.
 132 For example, as shown in [26], if we want to use the look ahead active learning strategy, each active
 133 learning cycle takes 3 hours to train the entire network of 15 epochs on the MNIST dataset, while it
 134 takes only 3 minutes to use NTK with a block structure.

135 **3.2 Framework**

136 We propose a look ahead strategy, NTKCPL, to approximate the empirical risk on the whole active
 137 learning pool directly. There are two challenges: (1) estimate empirical risk without labels and (2)
 138 estimate predictions of models trained with candidate sets efficiently and accurately.

139 For the first challenge, clusters on self-supervised features provide good pseudo-labels. Because most
 140 samples in the same cluster have the same label [39]. And when the number of clusters is increased,
 141 it can improve the purity of clusters, where purity refers to the probability that the sample has the
 142 same label within the same cluster. We call these clusters clustering-pseudo-labels (CPL), y_{cpl} .

143 For the second challenge, as introduced in sec. 3.1, NTK approximates the network well for random
 144 initialization and the computation time is acceptable. However, in our scenario, training on top of the
 145 self-supervised model, NTK does not approximate predictions of the whole network well. The main
 146 reason is that weights of the neural network are initialized by self-supervised learning rather than
 147 NTK initialization, i.e., drawn i.i.d. from a standard Gaussian [18]. In addition, the self-supervised
 148 initialization provides the neural network with a powerful feature representation capability that is not
 149 available in NTK. This leads to inconsistency between NTK predictions and network outputs. So, in
 150 our method, the NTK is used to approximate the classifier instead of the whole network. And the
 151 inputs of NTK are self-supervised features. Accordingly, we choose a training method following [24]
 152 that freezes the encoder initialized by self-supervised learning and trains only the MLP as a classifier.
 153 That training method achieves better or equal performance than fine-tuning the whole network in
 154 the low-budget case while its prediction is more consistent with the results of NTK. We denotes the
 155 predictions of NTK with trainset D_C as \hat{f}_{D_C} . Now, the active learning goal in eq. 1 is approximated
 156 as eq. 4.

$$\operatorname{argmin}_{D_C} \frac{1}{N} \sum_{i \in D} \operatorname{Loss}(\hat{f}_{D_C}(x_i), y_{cpl,i}) \quad (4)$$

157 The algorithm is shown in Alg. 1. For computational simplicity and without loss of generality, we
 158 use 0-1 loss to calculate empirical risk in eq. 4. In each round of active learning, after computation of
 159 NTK based on eq. 2 and generation of CPL based on the method introduced in sec. 3.4, the sample
 160 that minimizes the empirical risk on the whole active learning pool after adding labeled set is selected.

161 **3.3 NTKCPL Approximate Error**

162 In this section, we analyze what affects the accuracy of NTKCPL estimates of empirical risk on
 163 the whole active learning pool. The difference between the true empirical risk and the estimated

Algorithm 1 NTKCPL

```

1: Input: self-supervised feature  $f_{self}$ , active learning feature  $f_{al}$  labeled set  $L$ , unlabeled set  $U$ , budget  $b$ ,
   initial budget  $b_0$ , maximum cluster number  $C_{max}$ , model prediction  $Y_{pre,t-1}$  at the last active learning
   round
2: Output: labeled set  $L$ , model prediction  $Y_{pre,t}$  at this round
3: if  $L$  is  $\emptyset$  then
4:    $Y_{cpl} \leftarrow$  K-means( $f_{self}$ ,  $b_0$ )
5: else
6:    $N_{clu} = \min\{b_i/2, C_{max}\}$ 
7:    $Y_{cpl} \leftarrow$  CPL generation( $f_{al}$ ,  $Y_{pre,t-1}$ ,  $b_0$ ,  $N_{clu}$ ,  $L$ ) based on Alg. 2
8: end if
9: Initialize classifier, MLP, compute  $f_0$  and  $ker$  based on eq. 2
10: for  $itr = 1$  to  $b$  do
11:    $Emp\_risk = []$ 
12:   for  $(x_i, y_{cpl,i})$  in  $U$  do
13:     Compute  $Y_{NTK} = \hat{f}(ker, f_0, L \cup (x_i, y_{cpl,i}), U)$  based on eq. 3
14:      $Emp\_risk += [0-1Loss(Y_{NTK}, Y_{cpl})]$ 
15:   end for
16:    $i' = \operatorname{argmin} Emp\_risk$ 
17:    $L = L \cup (x_{i'}, y_{cpl,i'}), U = U \setminus x_{i'}$ 
18: end for
19: Query label  $y_{i',1,\dots,b}$  of  $x_{i',1,\dots,b}$ 
20:  $L = L \cup (x_{i',1,\dots,b}, y_{i',1,\dots,b}), U = U \setminus x_{i',1,\dots,b}$ 
21: Train classifier  $f_t$  on  $L$ 
22: model prediction  $Y_{pre,t} = f_t(U)$ 

```

164 empirical risk using NTK and CPL is shown in eq. 5. The approximation error can be divided into
165 two terms, the first one is the difference between NTK and neural network prediction, $error_{NTK}$,
166 and the second one is the difference caused by CPL during NTK estimation, $error_{CPL}$. For the
167 $error_{NTK}$, as we mentioned in sec. 3.2, NTK is used to approximate the classifier only to obtain
168 better consistency. To analyze $error_{CPL}$, we start with the relationship between the predictions of
169 NTK trained with the ground truth, $\hat{f}_y(x_i)$, and CPL, $\hat{f}_{cpl}(x_i)$.

$$\begin{aligned}
& \frac{1}{N} \sum_{i \in D} \left| Loss(f(x_i), y_i) - Loss(\hat{f}(x_i), y_{cpl,i}) \right| \\
& \leq \frac{1}{N} \sum_{i \in D} \left(\left| Loss(f(x_i), y_i) - Loss(\hat{f}(x_i), y_i) \right| + \left| Loss(\hat{f}(x_i), y_i) - Loss(\hat{f}(x_i), y_{cpl,i}) \right| \right) \quad (5)
\end{aligned}$$

170 **Definition** Denotes the j^{th} output of \hat{f}_{cpl} as \hat{f}_{cpl}^j . Label mapping function g converts NTK's
171 predictions about CPL classes, $\hat{f}_{cpl}(x_i)$, into predictions about true classes, $\hat{f}_{y_{map}}(x_i)$, based on
172 dominant labels within corresponding CPL classes as shown in eq. 6, where D_{dom} is a set of index k ,
173 where j is the dominant true label classes within CPL class, $y_{cpl,k}$.

$$\hat{f}_{y_{map}}^j(x_i) = \sum_{k \in D_{dom}} \hat{f}_{cpl}^k(x_i) \quad (6)$$

174 **Proposition** If the true labels of labeled samples are the dominant labels in their corresponding
175 CPL clusters, $\hat{f}_y(x_i) = g(\hat{f}_{cpl}(x_i))$. We defer the proof to appendix 1.

$$error_{CPL} = P_{nff} + P_{fnf} \quad (7)$$

176 As mentioned in sec. 3.2, we use 0-1 loss to calculate empirical risk. We can expand $error_{CPL}$ as
177 eq. 7, where we denote the probability that the NTK prediction agrees with the y but not with y_{cpl} as
178 P_{fnf} , and the probability that the NTK prediction does not agree with y but agrees with y_{cpl} as P_{nff} .
179 According to the proposition, we argue $\operatorname{argmax} \hat{f}_y(x_i)$ is most likely equal to $g(\operatorname{argmax} \hat{f}_{cpl}(x_i))$.

180 P_{fnf} refers to the case where different CPL classes correspond to the same true label class, i.e.,
 181 over-clustering. P_{nff} means that the true label of a sample is different from the dominant true label
 182 within its CPL class, i.e., the CPL class includes samples from different true label classes, which is
 183 called impurity. Detailed explanations and empirical evidence can be found in appendix 1.

184 3.4 Cluster Pseudo-Labels

185 As shown by eq. 7, the effect of CPL on the approximation error comes from the purity of the clusters
 186 and over-clustering. To improve clustering purity, we take two approaches: (1) clustering on the
 187 active learning feature, i.e., the output of the penultimate layer of the classifier, and (2) increasing
 188 the number of clusters. However, increasing the number of clusters may cause the labeled samples
 189 not to cover all classes of the CPL (under-coverage) and also increase the over-clustering error. For
 190 example, a group of samples with the same true label is clustered into K different classes. Even
 191 though NTK incorrectly predicts some samples as other CPL classes, their true empirical risk is zero.

192 To improve the under-coverage, we set the number of clusters to half of the total number of labels,
 193 i.e., each cluster includes two labeled samples on average. To improve the over-clustering,
 194 we manually set the maximum number of clusters and design a clustering-splitting approach
 195 instead of directly increasing the number of clusters. It splits the low-purity clusters and keeps
 196 the high-purity ones to reduce the extra over-clustering errors within samples located in the
 197 high-purity clusters. Specifically, we use the prediction of the neural network in each round of
 198 active learning to estimate the number of confusing samples within each cluster, i.e., the number
 199 of samples from classes that are different from the dominant class. The clusters that contain the
 200 largest number of confusing samples are split sequentially until a predefined number of clusters
 201 is reached. The cluster splitting algorithm is shown in Alg. 2, where we adopt the constrained
 202 K-Means [37] to improve the clusters from labeled sample constraints.
 203
 204
 205
 206
 207
 208
 209
 210

Algorithm 2 CPL generation

Input: active learning feature f_{al} , model predictions Y_{pre} , initial cluster number C_0 , cluster number C_{max} , labeled set L
Output: CPL Y_{cpl}
 $Clu_{1,\dots,C_0} \leftarrow$ Constrained K-means(f_{al}, L, C_0)
for $itr = 1$ **to** $(C_{max} - C_0)$ **do**
 $i' = argmax_i$ number of Confusing samples(Clu_i, Y_{pre})
 $f_{al,i'} \leftarrow$ f_{al} of samples within $Clu_{i'}$
 $Clu_{i',C_0+1} \leftarrow$ K-means($f_{al,i'}, 2$)
 $C_0 \leftarrow C_0 + 1$
end for
 $Y_{cpl} \leftarrow Clu_{1,\dots,C_{max}}$

211 4 Experiment Results

212 Our approach is validated on five datasets with various qualities of self-supervised features. Datasets
 213 with good self-supervised features, such as CIFAR-10 [21], CIFAR-100 [21], and ImageNet-100
 214 (a subset of ImageNet [11], following splitting in [36]), are included. SVHN [29] with poor self-
 215 supervised features is also included. Additionally, we consider practical scenarios where the total
 216 number of samples in the trainset is insufficient to support effective self-supervised training, such as
 217 Oxford-IIIT Pet dataset [30]. In this case, we evaluated the effectiveness of our method based on the
 218 model pre-trained on ImageNet [11].

219 **Baseline** We compare our proposed method with representative active learning strategies: (1)
 220 Random, (2) Entropy (uncertainty sampling, maximum entropy of output) [22], (3) Coreset (diversity
 221 active learning strategy, greedy solution of minimum coverage radius) [32], (4) BADGE (combination
 222 of uncertainty and diversity, kmeans++ sampling on grad embedding) [3], where the scalable version
 223 [10, 12], badge partition, is used in ImageNet-100, CIFAR-100 and Oxford-IIIT Pet because the
 224 huge dimension of grad embedding (5) Typiclust (designed for low-budget case) [15], (6) Lookahead
 225 (maximum output change based on NTK) [26].

226 **Implementation** Our method focuses on the low-budget regime, we followed the training method
 227 in [24], freezing weights of backbone initialized with self-supervised learning and then training a
 228 MLP as the classifier. The hyperparameters for training are set following [15] and can be found in
 229 appendix 3. For the self-supervised model, we adopt simsiam [9] for CIFAR-10, CIFAR-100 and
 230 SVHN and BYOL [14] for ImageNet-100 and Oxford-IIIT Pet. Resnet-18 [17] is used in CIFAR-10

Table 1: Comparison of accuracy of different active learning strategies on CIFAR-10. All results are averages over 5 runs. The best results are shown in red and the second-best results are shown in blue.

# Labels	Random	Entropy	Coreset(self)	BADGE	TypiClust	LookAhead	NTKCPL(self)	NTKCPL(al)
20	41.80±3.82	38.58±2.86	20.08±2.75	39.85±3.91	46.38±1.61	40.93±4.04	54.31±3.74	52.67±3.70
40	57.52±3.34	51.10±4.21	36.67±6.29	54.99±3.43	66.18±2.45	58.55±2.71	68.60±2.50	63.55±2.89
60	65.88±3.07	64.46±3.42	46.39±7.41	65.23±1.40	72.93±1.77	66.96±2.90	75.09±1.69	72.22±2.11
80	69.35±3.31	70.49±3.05	58.96±6.15	70.76±1.86	76.98±1.04	72.71±1.94	78.51±1.61	75.32±0.92
100	74.11±1.16	74.34±1.92	62.64±5.07	75.40±0.99	78.24±1.28	75.97±2.04	80.30±1.17	78.45±1.19
200	80.90±0.90	79.86±1.77	76.93±3.56	82.20±1.14	83.16±0.61	81.89±1.31	83.77±1.04	81.87±1.02
300	82.80±0.93	81.43±2.23	82.64±1.42	84.53±0.46	84.16±0.25	83.29±0.89	85.00±0.54	83.78±1.05
400	84.04±0.49	83.37±1.31	84.56±1.15	84.75±0.40	85.13±0.27	84.59±0.59	85.64±0.38	84.73±0.85
500	84.97±0.78	84.24±0.89	85.23±0.59	85.57±0.51	85.37±0.15	85.31±0.12	85.72±0.22	85.48±0.65
1000	86.26±0.38	84.94±0.48	86.75±0.36	86.06±0.31	86.07±0.14	85.69±0.47	86.83±0.33	87.15±0.57
1500	86.95±0.27	85.85±0.39	87.03±0.13	87.05±0.36	86.37±0.11	86.82±0.23	87.18±0.41	87.58±0.29
2000	87.30±0.37	86.92±0.15	87.34±0.27	87.31±0.47	86.55±0.21	87.16±0.19	87.34±0.41	87.87±0.39

231 and SVHN, WRN28-8 [44] is used in CIFAR-100 and Resnet-50 [17] is used in ImageNet-100 and
 232 Oxford-IIIT Pet.

233 The number of clusters in our method is set according to three rules, in the initial selection, it is
 234 set to the number of query samples, after that it is set to half of the query samples until the number
 235 of clusters reaches the maximum number of clusters. For CIFAR-10, CIFAR-100, ImageNet-100,
 236 SVHN, and Oxford-IIIT Pet, the maximum number of clusters is 100, 500, 300, 100, and 150,
 237 respectively. We followed [26] to sample a subset of the unlabeled set as the candidate set to select
 238 samples and estimate coverage. The candidate set includes 10,000 samples.

239 For the query step, most of the experiments (those on CIFAR-100, SVHN and Oxford-IIIT Pet)
 240 following the active learning literature by drawing a fixed number of samples from the unlabeled
 241 dataset to the oracle. Specifically, 500 for CIFAR-100, 20 for SVHN, and 40 for Oxford-IIIT Pet.
 242 We empirically found that fixed active learning query steps lead to much faster growth of classifier
 243 accuracy in the early stages of active learning (the amount of labels is about 10 times than the number
 244 of class) than in the later stages, so it is difficult to clearly observe the differences between different
 245 active learning strategies. For this reason, we empirically set varying query steps in our experiments
 246 with CIFAR-10 and ImageNet-100. Smaller query steps were used in the early stage of active learning
 247 and switched to larger query steps in the later stage. Specifically, for CIFAR-10, 20 samples are
 248 queried before 100 labels are available, 100 samples are queried before 500 labels and 500 labels
 249 are queried before 2000 labels. For ImageNet-100, 200 samples are selected before 1000 labels are
 250 available and 500 samples are queried before 2000 labels.

251 4.1 Main Results

252 All experiments were run 5 times and the avg. and std. are reported. Considering that the experiments
 253 are conducted for the scenario with a low annotation budget, it is not practical to construct a validation
 254 set to select the best checkpoints (the benefits of constructing a validation set are much less than
 255 using these labeled samples as training samples). Therefore, we report the final checkpoint accuracy,
 256 not the accuracy of the checkpoints determined by the validation set. The results are shown in fig. 2
 257 and table 1. The detailed results are in appendix 5.

258 **NTKCPL outperforms SOTA.** As shown in table 1, fig. 2. In most cases, our proposed method
 259 outperforms the baseline methods. For the few cases with only a small number of labels, our method
 260 shows comparable performance with the low-budget strategy, TypiClust, such as in CIFAR-100 with
 261 500 and 1000 labeled samples, and Oxford-IIIT Pet with 80 and 100 labeled samples.

262 **NTKCPL still shows good performance when the self-supervised features do not correspond**
 263 **well to the label classes.** Since the loss of self-supervised training is different from that of image
 264 classification, self-supervised features do not always correspond well with label classes. In SVHN
 265 dataset, self-supervised features of different classes are mixed together because the images include
 266 some irrelevant digits on both sides of the digit of interest [29]. Our method is similar to other

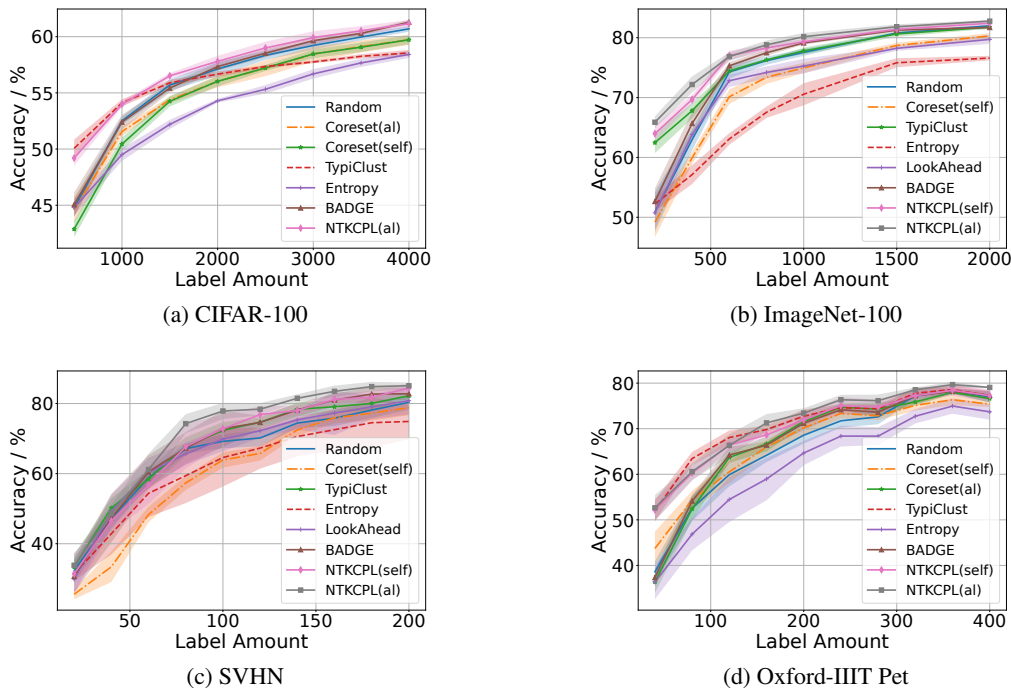


Figure 2: Performance of different active learning strategies. The shaded area represents std.

267 baseline strategies at the beginning of active learning, but it shows better results than baselines after
 268 several active learning rounds as shown in fig. 2c.

269 Another common scenario is the lack of sufficient samples to support effective self-supervised training.
 270 To evaluate in this context, we choose the Oxford-IIIT Pet dataset with the self-supervised model
 271 trained on ImageNet. The result is shown in fig. 2d. Our method has similar accuracy in the first three
 272 rounds as the TypiClust and outperforms all baseline methods afterward.

273 **NTKCPL has a wider effective budget range than**

274 **SOTA.** Active learning based on self-supervised models exhibits an intensified phase transition phenomenon.
 275 We plot the active learning gain of our method and baselines on different datasets in fig. 3. The average accuracy
 276 of our method, NTKCPL(al), outperforms the random baseline at all quantities of labels. In contrast, both the
 277 typical high-budget strategy, BADGE, and low-budget strategy, TypiClust, appear to be worse than the random
 278 baseline over a range of annotation quantities. We show the effective budget range of our method, NTKCPL, as well as the typical high-budget strategy,
 279 BADGE, and the typical low-budget strategy, TypiClust, across all experiments in table 2. The effective
 280 budget ratio refers to the proportion of the effective annotation quantity to the total annotation
 281 quantity, where the effective annotation quantity refers to the number of annotations at which active
 282 learning accuracy exceeds the random baseline (avg. + std.).
 283
 284
 285
 286
 287

Table 2: Comparison of the effective budget ratio of different active learning strategies.

	Effective Budget Ratio
TypiClust	40.8%
BADGE	42.0%
NTKCPL(al)	92.7%

288 **4.2 Ablation Study**

289 In this section, we evaluate the coverage estimation of our method and the effect of the maximum
 290 cluster number on NTKCPL. Also, we compare the effect of generating CPL on self-supervised
 291 features as well as on the active learning feature on the performance of NTKCPL.

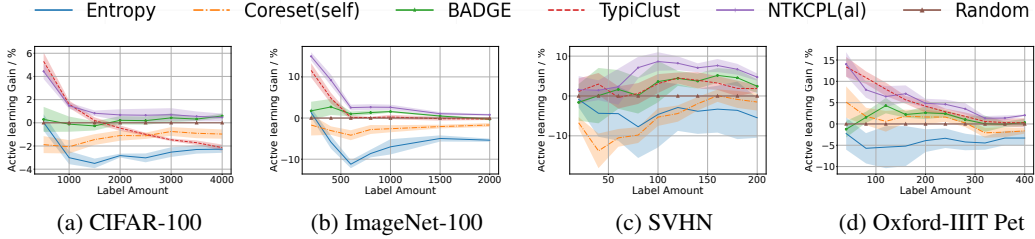


Figure 3: Active learning gain of different active learning strategies.

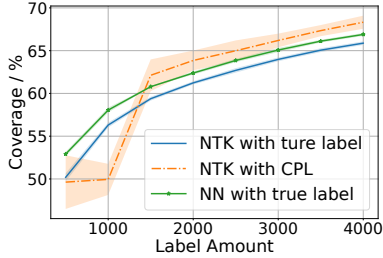


Figure 4: Coverage estimation on CIFAR-100.

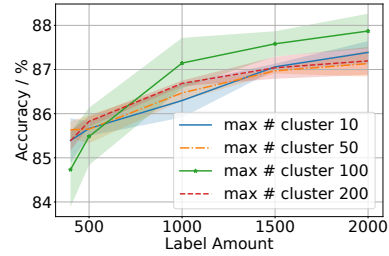


Figure 5: Effect of the maximum number of clusters on active learning performance on CIFAR-10.

292 **Coverage Estimation** We conducted experiments on CIFAR-100, where the coverage indicates the
 293 proportion of samples that are correctly predicted. The estimated coverage of NTK with true label
 294 and with CPL is shown in fig. 4. Our method approximates the true coverage well for most cases.

295 **Effect of the Maximum Number of CPL** The ablation experiments are conducted on CIFAR-10.
 296 We plot the accuracy when the number of annotations selected by active learning is greater than
 297 400 as shown in fig. 5. In this range, the number of classes of CPL is fixed at 10, 50, 100, and
 298 200, respectively. The experimental results support our analysis in sec. 3.4 that too many or too few
 299 clusters will increase the approximation error, which affects the performance of active learning.

300 **Effect of self-supervised feature-based and active learning feature-based clustering-pseudo-**
 301 **labels on NTKCPL.** We denote NTKCPL based on active learning features as NTKCPL(al) and
 302 NTKCPL based on self-supervised learning feature as NTKCPL(self). The results are shown in
 303 table 1 and fig. 2. From these experiments, we found that clustering on active learning features yields
 304 better results except for the case where the number of annotations is very small. Also, NTKCPL(self)
 305 is better than NTKCPL(al) in a wide range of annotation quantities (no more than 500), when
 306 self-supervised features are good such as experiment in the CIFAR-10.

307 5 Conclusion

308 We study the active learning problem when training on top of a self-supervised model. In this case,
 309 an intensified phase transition is observed and it influences the application of active learning. We
 310 propose NTKCPL that approximates empirical risk on the whole pool more directly. We also analyze
 311 the approximation error and design a CPL generation method based on the analysis to reduce the
 312 approximation error. Our method outperforms SOTA in most cases and has a wider effective budget
 313 range. The comprehensive experiments show that our method can work well on self-supervised
 314 features with different qualities.

315 Our approach is limited to the fixed training approach, i.e., training the classifier on top of a frozen
 316 self-supervised training encoder, which is restricted to the low-budget scenario because the fine-
 317 tuning training approach provides higher accuracy in the high-budget case. Therefore, (1) how to
 318 accurately approximate the fine-tuning model initialized with self-supervised weights using NTK and
 319 (2) whether the samples selected by our current method have good transferability for the fine-tuning
 320 would be interesting future directions.

References

- 321
- 322 [1] Inigo Alonso, Matan Yuval, Gal Eyal, Tali Treibitz, and Ana C Murillo. Coralseg: Learning coral
323 segmentation from sparse annotations. *Journal of Field Robotics*, 36(8):1456–1477, 2019.
- 324 [2] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact
325 computation with an infinitely wide neural net. *Advances in neural information processing systems*, 32,
326 2019.
- 327 [3] Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch
328 active learning by diverse, uncertain gradient lower bounds. In *International Conference on Learning
329 Representations*, 2020.
- 330 [4] Javad Zolfaghari Bengar, Joost van de Weijer, Bartłomiej Twardowski, and Bogdan Raducanu. Reducing
331 label effort: Self-supervised meets active learning. In *Proceedings of the IEEE/CVF International
332 Conference on Computer Vision*, pages 1631–1639, 2021.
- 333 [5] Zalan Borsos, Marco Tagliasacchi, and Andreas Krause. Semi-supervised batch active learning via bilevel
334 optimization. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal
335 Processing (ICASSP)*, pages 3495–3499. IEEE, 2021.
- 336 [6] Yao-Chun Chan, Mingchen Li, and Samet Oymak. On the marginal benefit of active learning: Does
337 self-supervision eat its cake? In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech
338 and Signal Processing (ICASSP)*, pages 3455–3459. IEEE, 2021.
- 339 [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
340 contrastive learning of visual representations. In *International conference on machine learning*, pages
341 1597–1607. PMLR, 2020.
- 342 [8] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-
343 supervised models are strong semi-supervised learners. *Advances in neural information processing systems*,
344 33:22243–22255, 2020.
- 345 [9] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the
346 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.
- 347 [10] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Ros-
348 tamizadeh, and Sanjiv Kumar. Batch active learning at scale. *Advances in Neural Information Processing
349 Systems*, 34:11933–11944, 2021.
- 350 [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
351 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.
352 Ieee, 2009.
- 353 [12] Zeyad Ali Sami Emam, Hong-Min Chu, Ping-Yeh Chiang, Wojciech Czaja, Richard Leapman, Micah
354 Goldblum, and Tom Goldstein. Active learning at the imagenet scale. *arXiv preprint arXiv:2111.12880*,
355 2021.
- 356 [13] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In
357 *International Conference on Machine Learning*, pages 1183–1192. PMLR, 2017.
- 358 [14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya,
359 Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own
360 latent—a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*,
361 33:21271–21284, 2020.
- 362 [15] Guy Hachohen, Avihu Dekel, and Daphna Weinshall. Active learning on a budget: Opposite strategies suit
363 high and low budgets. In *International Conference on Machine Learning*, pages 8175–8195. PMLR, 2022.
- 364 [16] Guy Hachohen and Daphna Weinshall. Misal: Active learning for every budget. 2023.
- 365 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
366 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 367 [18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization
368 in neural networks. *Advances in neural information processing systems*, 31, 2018.
- 369 [19] Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition
370 for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.

- 371 [20] Seo Taek Kong, Soomin Jeon, Dongbin Na, Jaewon Lee, Hong-Seok Lee, and Kyu-Hwan Jung. A neural
372 pre-conditioning active learning algorithm to reduce label complexity. *Advances in Neural Information*
373 *Processing Systems*, 35:32842–32853, 2022.
- 374 [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 375 [22] David D Lewis and Jason Catlett. Heterogeneous uncertainty sampling for supervised learning. In *Machine*
376 *learning proceedings 1994*, pages 148–156. Elsevier, 1994.
- 377 [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin
378 transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF*
379 *International Conference on Computer Vision*, pages 10012–10022, 2021.
- 380 [24] Rafid Mahmood, Sanja Fidler, and Marc T Law. Low-budget active learning via wasserstein distance: An
381 integer programming approach. In *International Conference on Learning Representations*, 2022.
- 382 [25] Sudhanshu Mittal, Maxim Tatarchenko, Özgün Çiçek, and Thomas Brox. Parting with illusions about deep
383 active learning. *arXiv preprint arXiv:1912.05361*, 2019.
- 384 [26] Mohamad Amin Mohamadi, Wonho Bae, and Danica J Sutherland. Making look-ahead active learning
385 strategies feasible with neural tangent kernels. In *Advances in Neural Information Processing Systems*,
386 2022.
- 387 [27] Mohamad Amin Mohamadi and Danica J Sutherland. A fast, well-founded approximation to the empirical
388 neural tangent kernel. *arXiv preprint arXiv:2206.12543*, 2022.
- 389 [28] Prateek Munjal, Nasir Hayat, Munawar Hayat, Jamshid Sourati, and Shadab Khan. Towards robust and
390 reproducible active learning using neural networks. In *Proceedings of the IEEE/CVF Conference on*
391 *Computer Vision and Pattern Recognition*, pages 223–232, 2022.
- 392 [29] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits
393 in natural images with unsupervised feature learning. 2011.
- 394 [30] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE*
395 *conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.
- 396 [31] Kossar Pourahmadi, Parsa Nooralinejad, and Hamed Pirsiavash. A simple baseline for low-budget active
397 learning. *UMBC Student Collection*, 2022.
- 398 [32] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach.
399 In *International Conference on Learning Representations*, 2018.
- 400 [33] Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled
401 method for query and training. In *International Conference on Artificial Intelligence and Statistics*, pages
402 1308–1318. PMLR, 2020.
- 403 [34] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings*
404 *of the IEEE/CVF International Conference on Computer Vision*, pages 5972–5981, 2019.
- 405 [35] Toan Tran, Thanh-Toan Do, Ian Reid, and Gustavo Carneiro. Bayesian generative active deep learning. In
406 *International Conference on Machine Learning*, pages 6295–6304. PMLR, 2019.
- 407 [36] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool.
408 Scan: Learning to classify images without labels. In *Computer Vision–ECCV 2020: 16th European*
409 *Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X*, pages 268–285. Springer, 2020.
- 410 [37] Kiri Wagstaff, Claire Cardie, Seth Rogers, Stefan Schrödl, et al. Constrained k-means clustering with
411 background knowledge. In *Icml*, volume 1, pages 577–584, 2001.
- 412 [38] Haonan Wang, Wei Huang, Ziwei Wu, Hanghang Tong, Andrew J Margenot, and Jingrui He. Deep
413 active learning by leveraging training dynamics. *Advances in Neural Information Processing Systems*,
414 35:25171–25184, 2022.
- 415 [39] Ziting Wen, Oscar Pizarro, and Stefan Williams. Active self-semi-supervised learning for few labeled
416 samples fast training. *arXiv e-prints*, pages arXiv–2203, 2022.
- 417 [40] Yichen Xie, Han Lu, Junchi Yan, Xiaokang Yang, Masayoshi Tomizuka, and Wei Zhan. Active finetuning:
418 Exploiting annotation budget in the pretraining-finetuning paradigm. *arXiv preprint arXiv:2303.14382*,
419 2023.

- 420 [41] Lin Yang, Yizhe Zhang, Jianxu Chen, Siyuan Zhang, and Danny Z Chen. Suggestive annotation: A deep
421 active learning framework for biomedical image segmentation. In *International conference on medical*
422 *image computing and computer-assisted intervention*, pages 399–407. Springer, 2017.
- 423 [42] Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a covering lens.
424 In *Advances in Neural Information Processing Systems*, 2022.
- 425 [43] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF*
426 *conference on computer vision and pattern recognition*, pages 93–102, 2019.
- 427 [44] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *British Machine Vision Conference*
428 *2016*. British Machine Vision Association, 2016.
- 429 [45] Wenqiao Zhang, Lei Zhu, James Hallinan, Shengyu Zhang, Andrew Makmur, Qingpeng Cai, and Beng Chin
430 Ooi. Boostmis: Boosting medical image semi-supervised learning with adaptive pseudo labeling and
431 informative active annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
432 *Pattern Recognition*, pages 20666–20676, 2022.