

# Best Practices in Active Learning for Semantic Segmentation

Sudhanshu Mittal\*  
University of Freiburg  
Freiburg, Germany

mittal@cs.uni-freiburg.de

Jörg P. Schäfer  
German Aerospace Center (DLR)  
Berlin, Germany

Joerg.Schaefer@dlr.de

Joshua Niemeijer\*  
German Aerospace Center (DLR)  
Braunschweig, Germany

Joshua.Niemeijer@dlr.de

Thomas Brox  
University of Freiburg  
Freiburg, Germany

brox@cs.uni-freiburg.de

## Abstract

Active learning is particularly of interest for semantic segmentation, where annotations are costly. Previous academic studies focused on datasets that are already very diverse and where the model is trained in a supervised manner with a large annotation budget. In contrast, data collected in many driving scenarios is highly redundant, and most medical applications are subject to very constrained annotation budgets. This work investigates the various types of existing active learning methods for semantic segmentation under diverse conditions across three dimensions - data distribution w.r.t. different redundancy levels, integration of semi-supervised learning, and different labeling budgets. We find that these three underlying factors are decisive for the selection of the best active learning approach. As an outcome of our study, we provide a comprehensive usage guide to obtain the best performance for each case. We also propose an exemplary evaluation task for driving scenarios, where data has high redundancy, to showcase the practical implications of our research findings.

## 1. Introduction

The objective of active learning is the reduction of annotation cost by selecting those samples for annotation, which are expected to yield the largest increase in the model's performance. It assumes that raw data can be collected in abundance for most large-scale data applications, such as autonomous driving, but annotation limits the use of this data. Semantic segmentation is particularly costly, as it requires pixel-level annotations. Active learning is, besides weakly supervised and semi-supervised learning, among the best-

Dataset↓	Annotation Budget			
	Low		High	
Supervision →	AL	SSL-AL	AL	SSL-AL
Diverse	✓	✓	✓	✓
Redundant	✓	✓	✓	✓

Table 1: We study current active learning (AL) methods for semantic segmentation over 3 dimensions - dataset distribution, annotation budget, and integration of semi-supervised learning (SSL-AL). Green cells denote newly studied settings in this work. Previous AL works correspond to the grey cells. This work provides a guide to use AL under all the above conditions.

known ways to deal with this situation.

In a typical deep active learning process, a batch of samples is acquired from a large unlabeled pool for annotation using an acquisition function and is added to the training scheme. This sampling is done over multiple cycles until an acceptable performance is reached or the annotation budget is exhausted. The acquisition function can be either a single-sample-based acquisition function, where a score is given to each sample individually or a batch-based acquisition function, where a cumulative score is given to the whole selected batch. Existing active learning methods for semantic segmentation assign the score to the sample either based on uncertainty [15, 35, 25] or representational value [37, 35, 34]. Most AL methods in the literature are evaluated on datasets like PASCAL-VOC [10], Cityscapes [6], and CamVid [1]. The shared attribute between these AL benchmark datasets is that they are highly diverse, as they were initially curated to provide comprehensive coverage of their corresponding domains. This curation process, however, is a sort of annotation because it is

\*These authors contributed equally to this work

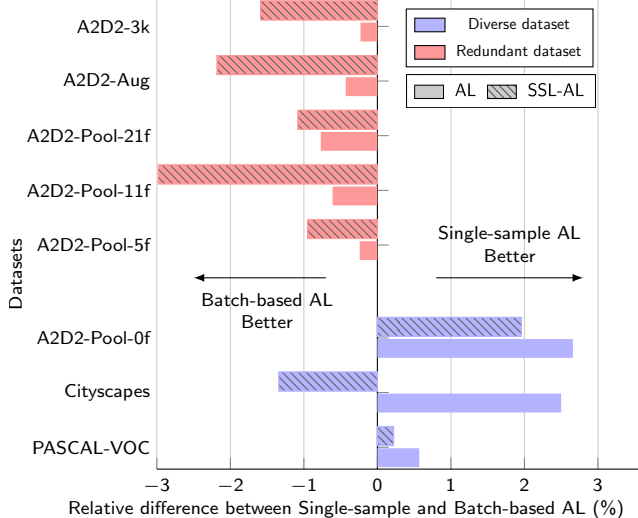


Figure 1: We analyse and compare single-sample-based AL and batch-based AL on datasets with different levels of redundancy. The figure shows the difference between the best-performing single-sample-based AL method and best performing batch-based AL method. We find that batch acquisition performs better for redundant datasets, and single-sample performs better for diverse datasets. The integration of semi-supervised learning with active learning (SSL-AL) performs well for batch-based acquisition.

typically not feasible in an entirely automated way.

State-of-the-art active learning methods for segmentation have been evaluated only in a particular experimental setup - highly diverse benchmark datasets with a comparatively large annotation budget; see Table 1. We seek answers to specific missing questions not captured by previous works.

**1. How do different active learning methods perform when the dataset has many redundant samples?** Samples with highly overlapping information are referred to as redundant samples, for example, the consecutive frames of a video. Many commonly used segmentation datasets were originally collected as videos for practical reasons, e.g., Cityscapes, CamVid, BDD100k [39]. Since active learning methods were only tested on filtered versions of these datasets, their applicability on redundant datasets is open and highly relevant.

**2. What happens when the initial unlabeled pool is also used for training along with annotated samples using semi-supervised learning (SSL)?** For image classification, many works [26, 13, 18, 28] have shown that integration of SSL into AL is advantageous. For semantic segmentation, this combination is not well studied.

**3. What happens when the annotation budget is low? Which methods scale best in such low-budget settings?**

Semantic segmentation annotations can be expensive for specific applications, especially in the medical domain. Therefore, it is critical to understand the behavior of the various active learning methods in low-budget settings.

In this work, we report the results of an empirical study designed to find answers to the above-raised questions. We study 5 existing active learning methods across the three dimensions as mentioned above - subject to different data distributions w.r.t. redundancy in the dataset, including the integration of semi-supervised learning, and under low as well as large annotation budget settings, as shown in Table 1. The outcome of this study yields new insights and provides, as the major contribution of this work, a guideline for the best selection of available techniques under the various tested conditions. Figure 1 illustrates some of the results, particularly, that the performance of acquisition functions can change depending on whether the dataset is redundant or diverse and that SSL integration plays an additional role in this. Additionally, we show that active learning in a low annotation budget setting can be particularly volatile, even nullifying the complete need for it in some cases. This further emphasizes the importance of knowing the underlying data distribution.

We also suggest a new evaluation task (A2D2-3K) for driving scenarios based on the highly redundant A2D2 dataset, which is closer to the raw data collection scheme in a driving case. The experiment outcome on this task aligns with the findings of our study for redundant dataset type with a high annotation budget setting and shows that there is a strong case for using active learning in this context.

## 2. Deep Active Learning

In this section, we briefly review the state of the art in deep active learning as relevant for our study. In particular, we review the available acquisition methods, the special considerations for segmentation, and the integration of semi-supervised learning.

The acquisition methods can be categorized into single-sample-based and batch-based approaches. They assess the value of new samples for selecting individually and collectively as a batch, respectively.

**Single sample acquisition** takes the top  $b$  samples according to the score of the acquisition function to select a batch of size  $b$ . Several methods follow this selection scheme based on either epistemic uncertainty or representation score. For example, uncertainty-based methods try to select the most uncertain samples to acquire a batch. Many methods, such as EqualAL [15], Ensemble+AT [24], and CEAL [36], estimate uncertainty based on the output probabilities. Epistemic uncertainty, estimated using Entropy [32], is often used as a strong baseline in several active learning works [15, 33, 29]. Some methods, namely BALD [17] and DBAL [12] employed a Bayesian

approach using Monte Carlo Dropout [11] to measure the epistemic uncertainty. Representation-based methods aim to select the most representative samples of the dataset that are not yet covered by the labeled samples. Numerous adversarial learning-based methods utilize an auxiliary network to score samples based on this measure, including DAAL [37], VAAL [35], and WAAL [34]. For our study, we employ Entropy, EqualAL, and BALD to represent single-sample acquisition methods due to their direct applicability to segmentation tasks. We did not include deep ensemble-based methods due to their limited scalability and adversarial methods due to their hyperparameter sensitivity. In general, single-sample acquisition approaches select individually very informative samples but do not optimize the joint improvement obtained with the whole batch.

**Batch-based acquisition** methods acquire the whole batch of size  $b$  to maximize cumulative information gain. Sener *et al.* [31] formulated the acquisition function as a core-set selection approach based on the feature representations. It is a representation-based approach that selects the batch of samples jointly to represent the whole data distribution. BatchBALD [23] is a greedy algorithm that selects a batch of points by estimating the joint mutual information between the whole batch and the model parameters. This method was also proposed to remedy the mode collapse issue, where the acquisition function collapses into selecting only similar samples (see Section 4.1 for details). However, it is limited to simple image classification datasets like MNIST [8] since its computation complexity grows exponentially with the batch size. Some more recent batch-based methods include k-MEANS++ [40], GLISTER [21], ADS [19], but these methods only evaluate on image classification tasks. For the study, we selected the Coreset method [31] to represent batch-based methods due to its effectiveness, simplicity, and easy scalability to the segmentation task.

## 2.1. Active Learning for Semantic Segmentation

When applied to semantic segmentation, active learning methods must choose which area of the image is to be considered for the acquisition: the full image [35], superpixels [2], polygons [27, 15], or each pixel [33]. There is no common understanding so far of which approach is cheaper and more effective. Thus, our study uses the straightforward image-wise selection and annotation procedure.

Most existing methods for segmentation are based on the model’s uncertainty for the input image, where the average score over all pixels in the image is used to select top-k images. **Entropy** [32] (estimated uncertainty) is a widely used active learning baseline for selection. This function computes per-pixel entropy for the predicted output and uses the averaged entropy as the final score. **EqualAL** [15] determines the uncertainty based on the consistency of the pre-

dition on the original image and its horizontally flipped version. The average value over all the pixels is used as the final score. **BALD** [17] is often used as baseline in existing works. It is employed for segmentation by adding dropout layers in the decoder module of the segmentation model and then computing the pixel-wise mutual information using multiple forward passes. **Coreset** [31] is a batch-based approach that was initially proposed for image classification, but it can be easily modified for segmentation. For e.g., the pooled output of the ASPP [4] module in the DeepLabv3+ [5] model can be used as the feature representation for computing distance between the samples. Some other methods [35, 22, 34] use a GAN model to learn a combined feature space for labeled and unlabeled images and utilize the discriminator output to select the least represented images. Our study includes Entropy, EqualAL, BALD, and Coreset approaches for the analysis, along with the random sampling baseline. In this work, these methods are also studied with the integration of semi-supervised learning.

## 2.2. Semi-supervised Active Learning

Active learning uses a pool of unlabeled samples only for selecting new samples for annotation. However, this pool can also be used for semi-supervised learning (SSL), where the objective is to learn jointly from labeled and unlabeled samples. The combination of SSL and AL has been used successfully in many contexts, such as speech understanding [20, 9], image classification [31, 13, 27, 28], and pedestrian detection [30]. Some recent works have also studied active learning methods with the integration of SSL for segmentation, but their scope is limited only to special cases like subsampled driving datasets [29] or low labeling budget [27], both cases with only single-sample acquisition methods. Our work provides a broader overview of the integration of SSL and active learning for the segmentation task. We study this integration over datasets with different redundancy levels, under different labeling budgets, and with single-sample and batch-based methods. Our findings explain when this integration is effective and boosts the active learning method.

**Integration of SSL and AL.** A successful integration can also be conceptually explained based on the underlying assumption of semi-supervised learning and the selection principle of the active learning approach. According to the *clustering assumption* of SSL, if two points belong to the same cluster, then their outputs are likely to be close and can be connected by a short curve [3]. In this regard, when labeled samples align with the clusters of unlabeled samples, the cluster assumption of SSL is satisfied, resulting in a good performance. Consequently, to maximize semi-supervised learning performance, newly selected samples must cover the unlabeled clusters that are not already

covered by labeled samples. Only acquisition functions that foster this coverage requirement have the potential to leverage the additional benefits that arise from the integration of semi-supervised learning. A batch-based method, e.g., Coreset, selects samples for annotations to minimize the distance to the farthest neighbor. By transitivity, such labeled samples would have a higher tendency to propagate the knowledge to neighboring unlabeled samples in the cluster and utilize the knowledge of unlabeled samples using a semi-supervised learning objective and help boost the model performance. Similar behavior can also be attained using other clustering approaches that optimize for coverage.

### 3. Experimental Setup

#### 3.1. Tested Approaches

In our study, we test five active learning acquisition functions as discussed in Section 2, including Random, Entropy, EqualAL, BALD, and Coreset. Here Entropy, EqualAL, and BALD approach represent single-sample, and Coreset represents the batch-based approach. All methods select the whole image for annotation. To leverage the unlabeled samples, we use the semi-supervised learning s4GAN method [26]. We pair all the used active learning approaches with SSL using this approach. This is marked by the suffix ‘-SSL’ in the experiments. In particular, we train the model using an SSL objective, which impacts the resulting model and hence the acquisition function.

#### 3.2. Datasets

Active learning methods are often evaluated on PASCAL-VOC and Cityscapes datasets, where PASCAL-VOC is naturally diverse while Cityscapes is diversified by subsampling from videos. In this work, we test on an additional driving dataset, A2D2, which is highly redundant. We evaluate the methods on these three datasets. To understand the nature of active learning methods over varying levels of redundancy in the dataset, we curate 5 smaller dataset pools from the large, original A2D2 dataset, described further below as A2D2-Pools.

**Cityscapes** [6] is a driving dataset used to benchmark semantic segmentation tasks. The dataset was originally collected as videos from 27 cities, where a diverse set of images were selected for annotation. Due to the selection, Cityscapes cannot cover the redundant data scenario in our evaluation, although it was derived from videos. As we will see in the results, the nature of the active learning method changes when considering the raw form of data in a driving scenario, and pre-filtering, as done in Cityscapes, is sub-optimal compared to directly applying active learning on the raw data (see Section 4.4).

**PASCAL-VOC** [10] is another widely used segmenta-

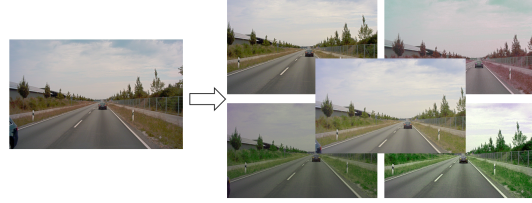


Figure 2: A2D2 Pool-Aug. Left: the original image. Right: the duplication through color augmentation and random cropping of the original image

tion dataset. We use the extended dataset [16], which consists of 10582 training and 1449 validation images. It contains a wide spectrum of natural images with mixed categories like vehicles, animals, furniture, etc. It is the most diverse dataset in this study.

**A2D2** [14] is a large-scale driving dataset consisting of 41277 annotated images with a resolution of  $1920 \times 1208$  from 23 sequences. It covers an urban setting from highways, country roads, and three cities. It contains labels for 38 categories. We map them to the 19 classes of Cityscapes for our experiments. A2D2 provides annotations for every  $\sim 10^{th}$  frame in the sequence and contains a lot of overlapping information between frames. Some consecutive frames are shown in the Appendix. We utilize 40135 frames from 22 sequences for creating our training sets and one sequence consisting of 1142 images for validation. The validation sequence ‘20180925\_112730’ is selected based on the maximum class balance. A2D2 represents the most diverse raw dataset in our study.

**A2D2 Pools.** To obtain a more continuous spectrum between diverse and redundant datasets, we created five smaller dataset pools by subsampling the large A2D2 datasets. Each pool comprises 2640 images, which is comparable in size to the Cityscapes training set. Four pools are curated by subsampling the original dataset, while the fifth pool is created by augmentation. The first four pools, denoted by Pool-Xf (where X is 0, 5, 11, and 21), were created by randomly selecting samples and X consecutive frames for each randomly selected sample from the original A2D2 dataset. Pool-0f contains only randomly selected images. We assume that the consecutive frames contain highly redundant information. Therefore, the pool with more consecutive frames has higher redundancy and lower diversity. The fifth pool, Pool-Aug, contains augmented duplicates in place of the consecutive frames. We create five duplicates of each randomly selected frame by randomly cropping 85% of the image area and adding color augmentation (see figure 2).



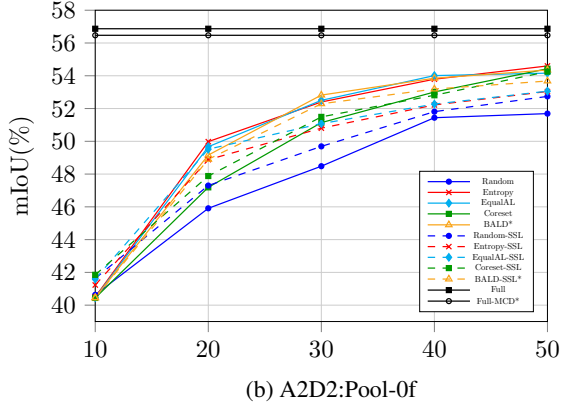
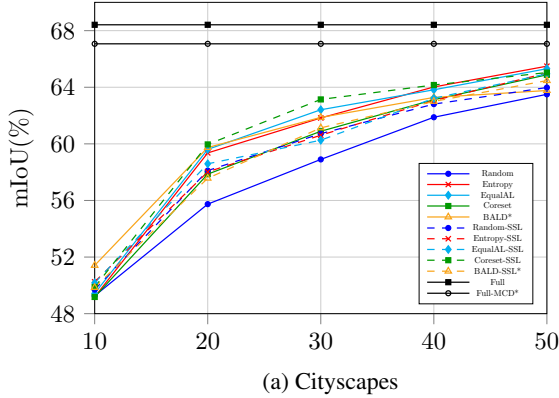


Figure 3: Results on diverse driving datasets. Active learning performance curves on Cityscapes and A2D2:Pool-Of. X-axis shows the percentage of the labeled dataset. The methods which utilize MC-dropout in their network architecture are marked with \*, and are only comparable to other methods with MC-dropout.

A	AL Method Metric →	SSL	Cityscapes		A2D2 Pool-0f	
			mIoU	AUC	mIoU	AUC
S	Random	✗	58.90	23.29	48.48	19.20
S	Entropy	✗	61.83	24.25	52.40	<b>20.37</b>
S	EqualAL	✗	62.41	24.32	<b>52.50</b>	20.35
B	Coreset	✗	60.89	23.89	51.14	19.88
S	Random-SSL	✓	60.72	23.85	49.69	19.60
S	Entropy-SSL	✓	60.61	23.93	50.80	19.90
S	EqualAL-SSL	✓	60.26	23.96	51.08	20.02
B	Coreset-SSL	✓	<b>63.14</b>	<b>24.47</b>	51.49	20.02
-	100%	✗	68.42	27.37	56.87	22.75
<i>With MC-Dropout decoder</i>						
S	BALD	✗	61.87	24.28	<b>52.82</b>	<b>20.32</b>
S	BALD-SSL	✓	61.13	23.89	52.29	20.14
-	100%-MCD	✗	67.07	26.83	56.47	22.59

Table 2: Active Learning results on Cityscapes and A2D2 Pool-Of. AUC@50 and mIoU@30 metrics are reported. A denotes Acquisition method type. S and B denotes the single-sample and batch-based acquisition, respectively.

### 3.3. Experiment Details

**Implementation details.** We used the DeepLabv3+ [5] architecture with the Wide-ResNet38 (WRN-38) [38] backbone for all experiments. The backbone WRN-38 is pre-trained using ImageNet [7]. We run all methods on 3 random seeds and report the mean performance. All other training details and hyperparameter information is included in the Appendix. Since the BALD method requires the introduction of Dropout layers into the architecture, we segregate the methods into two categories: with Monte Carlo Dropout (MCD) and without Monte Carlo Dropout layers. Random, Entropy, EqualAL, and Coreset are without MCD.

A	AL Method Metric →	SSL	VOC 5-5		VOC 10-10	
			mIoU	AUC	mIoU	AUC
S	Random	✗	70.70	13.92	72.13	28.85
S	Entropy	✗	70.38	13.94	73.72	29.17
S	EqualAL	✗	69.14	13.82	73.40	29.03
B	Coreset	✗	70.85	13.96	73.63	29.06
S	Random-SSL	✓	72.57	14.36	75.33	29.87
S	Entropy-SSL	✓	73.36	14.51	<b>76.08</b>	30.01
S	EqualAL-SSL	✓	<b>73.39</b>	<b>14.55</b>	75.89	<b>30.06</b>
B	Coreset-SSL	✓	72.88	14.46	75.91	30.03
-	100%	✗	77.00	15.40	77.00	30.80

Table 3: Active Learning results on PASCAL-VOC dataset in 5-5 and 10-10 settings. AUC@50 and mIoU@30 metric are reported. A, S and B denotes acquisition method type, single-sample and batch-based acquisition, respectively.

BALD is based on a MCD network. Since the models used in the two categories are not exactly comparable due to different architectures, we also show the fully-supervised performance for both with MCD (100% MCD) and without MCD (100%) architectures.

**Evaluation scheme.** We evaluate the methods across different data budget settings, denoted by  $\mathcal{I} - \mathcal{S}$ , where  $\mathcal{I}$  is the initial label budget,  $\mathcal{S}$  is the sampling-label budget, and  $\mathcal{I}, \mathcal{S}$  indicates the percentage of the dataset size. Images are sampled randomly to fulfill the initial label budget. For the subsequent steps, images are sampled using the AL acquisition function up to the allowed sampling-label budget. We test datasets with 10-10, 5-5, and 2-2 budget settings.

**Evaluation metrics.** We use mean Intersection over Union (mIoU) to evaluate the performance of the model at each AL cycle step. For the evaluation of the active learning method, we use two metrics: Area Under the Budget Curve

A	AL Method Metric →	SSL	Pool-5f		Pool-11f		Pool-21f		Pool-Aug	
			mIoU	AUC	mIoU	AUC	mIoU	AUC	mIoU	AUC
S	Random	✗	47.58	18.69	44.61	17.76	44.52	17.67	43.80	17.15
S	Entropy	✗	49.96	19.48	47.43	18.52	46.08	18.21	44.51	17.33
S	EqualAL	✗	49.50	19.29	47.14	18.44	46.32	18.18	44.24	17.29
B	Coreset	✗	50.08	19.44	47.72	18.69	46.68	18.38	44.70	17.54
S	Random-SSL	✓	47.92	19.03	45.25	18.02	46.27	18.19	44.17	17.29
S	Entropy-SSL	✓	48.78	19.31	47.53	18.56	46.93	18.43	44.50	17.47
S	EqualAL-SSL	✓	48.80	19.28	46.50	18.39	47.11	18.54	44.81	17.56
B	Coreset-SSL	✓	<b>50.44</b>	<b>19.69</b>	<b>48.99</b>	<b>19.01</b>	<b>47.62</b>	<b>18.69</b>	<b>45.81</b>	<b>17.74</b>
-	100%	✗	53.25	21.30	48.85	19.54	49.23	19.69	46.03	18.41
<i>With MC-Dropout decoder</i>										
S	BALD	✗	50.40	19.29	47.85	18.74	46.78	18.57	45.53	17.80
S	BALD-SSL	✓	50.33	19.62	47.34	18.61	47.06	18.57	45.16	17.72
-	100%-MCD	✗	53.82	21.53	50.86	20.34	50.43	20.17	46.62	18.65

Table 4: Active Learning results on A2D2-Pool5f, A2D2-Pool11f, A2D2-Pool-21f, and A2D2-PoolAug. AUC@50 and mIoU@30 metrics are reported. S and B denotes the single-sample and batch-based acquisition, respectively.

(AUC@B) and mean Intersection over Union at a budget B (mIoU@B). **AUC@B** is the area under the performance curves, shown in Figure 3 and 4. It captures a cumulative score of the AL performance curve up to a budget B, where B is the percentage of the labeled dataset size. For the experiments on A2D2 pools, we use B=50 in the 10-10 setting. For PASCAL-VOC, we run three experiments with B=10, 25, and 50 in 2-2, 5-5, and 10-10 settings, respectively. For Cityscapes, we experiment with B=50 in the 10-10 setting. **mIoU@B** reports the performance of the model after using a certain labeling budget B. We report performance at an intermediate labeling budget to clearly see the ranking of the AL methods.

## 4. Results

Here, we answer the three questions raised in Section 1 concerning the behavior of active learning methods w.r.t data distribution in terms of redundancy, integration of semi-supervised learning, and different labeling budgets. For each experiment, we compare random sampling, single-sample, or batch-based acquisition approaches.

### 4.1. Impact of Dataset Redundancy

Table 2 and Figure 3 show the results on Cityscapes and A2D2 Pool-0f. For both datasets, the single-sample (S) method, EqualAL, performs the best in the supervised-only setting. Table 3 shows the results obtained on the PASCAL-VOC dataset in 5-5 and 10-10 settings. Single-sample-based methods perform the best in the 10-10 setting, whereas Coreset performs the best in the 5-5 AL setting by a marginal gap w.r.t. random baseline. Table 4 and Figure 4 show the results for the redundant datasets. The

batch-based Coreset method consistently performs the best in all four datasets in the supervised-only setting.

**Diverse datasets need a single-sample method and redundant datasets need a batch-based method.** We observe that the order of best-performing models changes based on the level of redundancy in the dataset. Single-sample-based acquisition functions perform best on diverse datasets, whereas batch-based acquisition functions perform best on redundant datasets. We attribute this reversed effect to the mode collapse problem, where, for redundant datasets, single-sample acquisition methods select local clusters of similar samples. Diverse datasets are devoid of this issue as they do not possess local clusters due to high diversity across samples. Therefore, the diversity-driven acquisition is not critical for diverse datasets.

This observation is consistent for PASCAL-VOC, where single-sample-based uncertainty-type methods perform better than batch-based and random methods in the high-budget setting. The difference between the methods is only marginal here since most acquired samples add ample new information due to the highly diverse nature of the dataset. This difference further diminishes w.r.t. random baseline with a lower labeling budget (*e.g.* 5-5) since any learned useful bias also becomes weaker. The observations for the 5-5 setting tend towards a very low-budget setting which is further analysed in Section 4.3.

**Mode collapse analysis.** Here, we analyse and visualize the above mentioned model collapse issue. Mode collapse in active learning refers to the circumstance that acquisition functions tend to select a set of similar (redundant) samples when acquiring batches of data [23]. This can occur when a single-sample acquisition function gives a high score to

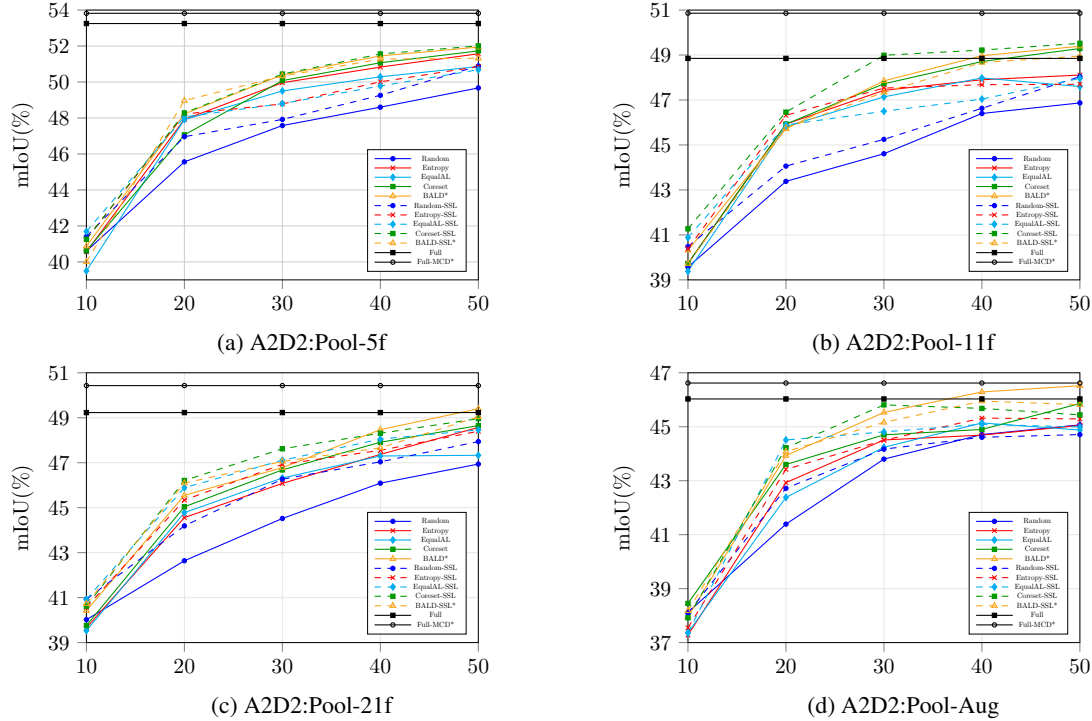


Figure 4: Results on redundant datasets. Active Learning performance curves on A2D2 dataset: Pool-5f, Pool-11f, Pool-21f, and Pool-Aug. The X-axis shows the percentage of labeled datasets. The methods which utilize MC-Dropout in their network architecture are marked with \*.

at least one of the similar samples in the set. Since similar samples have highly overlapping information, all samples in the set receive a high score. Thus, all similar samples tend to be selected, causing this collapse. Since the selected samples are all very similar, their annotation does not add much more value to the model than if a single sample was added.

We provide a qualitative analysis of the mode collapse issue on the redundant A2D2 Pool-21f. We plot the feature representations using t-SNE to show the selection process for a single-sample-based Entropy function and batch-based Coreset function, shown in Figure 5. It shows that Entropy acquisition selects many samples within local clusters, which are similar samples with overlapping information. This yields a suboptimal use of the annotation budget. In contrast, Coreset acquisition has a good selection coverage and avoids this mode collapse.

In this work, we argue that mode collapse is a common issue in many real-world datasets, containing similar samples. A good acquisition function for such datasets must be aware of the batch’s diversity to address the mode collapse issue. It is largely ignored due to the narrow scope of existing AL benchmarks like PASCAL-VOC and Cityscapes, which only cover diverse datasets.

## 4.2. Systematic Integration of SSL

For all redundant datasets, the Coreset-SSL approach consistently performs the best; see results in Table 4 and Figure 4. For diverse datasets, SSL integration is also helpful, but there is no consistent best approach. For the PASCAL-VOC dataset, single-sample based methods with SSL show the best performance, shown in Table 3. For Cityscapes, Coreset-SSL outperforms all other approaches; see Table 2 and Figure 3. For A2D2-Pool0f, Coreset-SSL improves over Coreset, but the single-sample acquisition method BALD approach shows the best performance.

**Redundant datasets favour the integration of batch-based active learning and semi-supervised learning.** The batch-based acquisition function Coreset always profits from the integration of SSL. Coreset aligns well with the SSL objective since Coreset selects samples from each local cluster, thus covering the whole data distribution. This assists SSL in obtaining maximum information from the unlabeled samples, as discussed in Section 2.2. This effect is especially strong in the redundant A2D2 pools, where Coreset-SSL always improves over Coreset and also shows the best performance. In contrast, SSL integration for single-sample methods is either harmful or ineffective, except for the PASCAL-VOC dataset. Interestingly, in Pool-

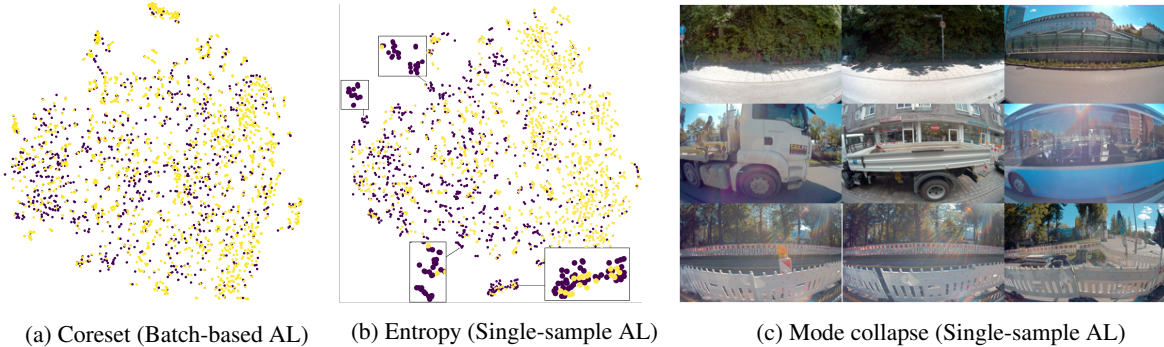


Figure 5: TSNE plots of (a) Coreset and (b) Entropy functions for A2D2 Pool-21f. The yellow points are feature representation from the unlabeled set, the violet point are the acquired points. The batch-based approach has good selection coverage, whereas the single-sample acquisition approach selects similar samples from clusters. Figure (c) shows acquired redundant samples from the violet clusters in (b).

11f, some Coreset-SSL methods even outperform the 100% baseline with less than 30% labeled data. This indicates that some labeled redundant samples can even harm the model (see Figure 4), possibly due to data imbalance. For Cityscapes, SSL with Coreset yields significant improvement, and SSL even changes the ranking of the methods. We see that EqualAL performs the best in the supervised-only setting, whereas Coreset-SSL surpasses all methods. This slight anomaly in the case of Cityscapes happens because the advantage due to the combination of SSL and batch-based method is greater than the advantage of using single-sample methods in non-redundant datasets. For diverse PASCAL-VOC, all methods align well with SSL. All methods perform well with no clear winner method since all selection criteria select samples that provide good coverage of the data distribution.

### 4.3. Low Annotation Budget

**Active learning is volatile with a low budget.** Experimenting with PASCAL-VOC in the 2-2 budget setting, Random-SSL performs the best, i.e., semi-supervised learning without active learning component. We believe that active learning fails in this setting because it fails to capture any helpful bias for selection in such a low-data regime with diverse samples. Our observations in this low-budget setting confirm and provide a stronger empirical support for similar behavior observed in [27]. For A2D2 Pool-0f and Cityscapes in the 2-2 setting, the single-sample acquisition performs the best, while its SSL integration is detrimental. These methods possibly learn some useful bias due to the specialized driving domain. For redundant datasets in low budget settings, batch-based acquisition is still the most effective way. However, SSL does not contribute any additional improvements due to insufficient labeled samples to support learning from unlabeled samples. Overall, we observe a highly volatile nature of active learning in conjunction with

A	AL Method	SSL	mIoU	AUC
B	Uniform	✗	57.75	—
S	Random	✗	56.14	5.35
S	Entropy	✗	60.16	5.53
B	Coreset	✗	60.30	5.55
S	Uniform (@5) + Entropy	✗	60.40	5.66
B	Uniform-SSL	✓	58.93	—
S	Random-SSL	✓	57.57	5.53
S	Entropy-SSL	✓	59.91	5.61
B	Coreset-SSL	✓	<b>61.13</b>	<b>5.72</b>
S	Uniform (@5) + Ent-SSL	✓	59.63	5.59
-	100%	✗	66.65	6.64

Table 5: AL results on the proposed A2D2-3k task. mIoU@7.5 and AUC@7.5 is reported. S and B denotes the single-sample and batch-based acquisition, respectively. Uniform refers to temporal subsampling selection process and (@5) means every 5<sup>th</sup> frame.

a low budget. The ideal policy transitions from random selection towards the batch-based acquisition, as the dataset redundancy goes from low to high. The result tables corresponding to this section are included in the Appendix.

### 4.4. An exemplar case study: A2D2-3K task

Previous active learning works on semantic segmentation cover only the combination of a diverse dataset and a high annotation budget. In contrast, the collected raw data can be quite redundant, like in video datasets. To study this missing redundant setting, we propose a new active learning task A2D2-3K for segmentation based on the A2D2 dataset. The aim of the new task is to select 3K images (similar size to Cityscapes) from the original A2D2 dataset (~40K images) to achieve the best performance. We select 3K images using active learning in 3 cycles with 1K images each. We



Dataset ↓	Annotation Budget			
	Low		High	
Sup. →	AL	SSL-AL	AL	SSL-AL
Diverse	Random	Random-SSL	Single	Single-SSL
Redundant	Batch	Batch	Batch	Batch-SSL

Table 6: Overview showing the best performing AL method for each scenario. Single and Batch refer to single-sample and batch-based method, and Random refers to random selection. Suffix -SSL refers to the usage of semi-supervised learning.

compare 5 acquisition functions, including Random, Entropy, and Coreset, along with SSL integration. Such video datasets are often manually subsampled based on some prior information like time or location, and then used for active learning. Therefore, we also include two such baselines - (a) where 3K samples are uniformly selected based on time information, denoted as Uniform, and (b) where every fifth sample is first selected uniformly to select  $\sim 8K$  samples and then applied with Entropy acquisition function, denoted as Uniform(@5)+Entropy. The second approach is closer to previously used active learning benchmarks in the driving context. Results are shown in Table 5. We find that the batch-based Coreset-SSL method performs the best, discussed in Section 4.2, while the subsampling-based approaches are sub-optimal. This makes an excellent case for active learning in datasets with high redundancy, as active learning filters the data better than time-based subsampling methods.

## 5. Conclusion

This work shows that active learning is indeed a useful tool for semantic segmentation. However, it is vital to understand the behavior of different active learning methods in various application scenarios. Table 6 provides an overview of the best performing methods for each scenario. Our findings indicate that single-sample-based uncertainty is a suitable measure for sample selection in diverse datasets. In contrast, batch-based diversity-driven measures are better suited for datasets with high levels of redundancy. SSL is successfully integrated with batch-based diversity-driven methods. However, it can have a detrimental impact when combined with single-sample-based uncertainty acquisition functions. Active learning with low annotation budgets is highly sensitive to the level of redundancy in the dataset. These findings have been missing in method development, which are optimized only for a few scenarios. The results of this study facilitate a broader view on the task with presumably positive effects in many applications.

## Acknowledgement

The authors would like to thank Philipp Schröppel, Jan Bechtold, and María A. Bravo for their constructive criticism on the manuscript. The research leading to these results is funded by the German Federal Ministry for Economic Affairs and Climate Action within the project “KI Delta Learning” (Forderkennzeichen 19A19013N) and “KI Wissen – Entwicklung von Methoden für die Einbindung von Wissen in maschinelles Lernen”. The authors would like to thank the consortium for the successful cooperation. Funded by the Deutsche Forschungsgemeinschaft (DFG) - 417962828.

## References

- [1] Gabriel J. Brostow, Julien Fauqueur, and Roberto Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern Recognition Letters*, 2008. 1
- [2] Lile Cai, Xun Xu, Jun Hao Liew, and Chuan Sheng Foo. Revisiting superpixels for active learning in semantic segmentation with realistic annotation costs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [3] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien, editors. *Semi-Supervised Learning*. The MIT Press, 2006. 3
- [4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 40(4):834–848, 2018. 3
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. *CoRR*, abs/1802.02611, 2018. 3, 5, 12
- [6] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 4
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009. 5, 12
- [8] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 2012. 3
- [9] Thomas Drugman, Janne Pylkkonen, and Reinhard Kneser. Active and semi-supervised learning in asr: Benefits on the acoustic and language models. In *Interspeech*, 2016. 3
- [10] Mark Everingham, Luc van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 2010. 1, 4
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning*, 2016. 3
- [12] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *ICML*, 2017. 2
- [13] Mingfei Gao, Zizhao Zhang, Guo Yu, Sercan Ö. Arık, Larry S. Davis, and Tomas Pfister. Consistency-based semi-supervised active learning: Towards minimizing labeling cost. In *ECCV*, 2020. 2, 3
- [14] Jakob Geyer, Yohannes Kassahun, Mentar Mahmudi, Xavier Ricou, Rupesh Durgesh, Andrew S. Chung, Lorenz Hauswald, Viet Hoang Pham, Maximilian Mühlegg, Sebastian Dorn, Tiffany Fernandez, Martin Jänicke, Sudesh Mirashi, Chiragkumar Savani, Martin Sturm, Oleksandr Vorobiov, Martin Oelker, Sebastian Garreis, and Peter Schuberth. A2D2: Audi Autonomous Driving Dataset. 2020. 4
- [15] Kitani Kris Golestaneh, S. Alireza. Importance of self-consistency in active learning for semantic segmentation. *BMVC*, 2020. 1, 2, 3
- [16] Bharath Hariharan, Pablo Arbeláez, Lubomir Bourdev, Subhransu Maji, and Jitendra Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998, 2011. 4
- [17] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *ArXiv*, abs/1112.5745, 2011. 2, 3
- [18] Siyu Huang, Tianyang Wang, Haoyi Xiong, Jun Huan, and Dejing Dou. Semi-supervised active learning with temporal output discrepancy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2
- [19] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li, Ce Zhang, Dawn Song, and Costas J. Spanos. Towards efficient data valuation based on the shapley value. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, 2019. 3
- [20] Stamatis Karlos, Christos Aridas, Vasileios G Kanas, and Sotiris Kotsiantis. Classification of acoustical signals by combining active learning strategies with semi-supervised learning schemes. *Neural Computing and Applications*, pages 1–18, 2021. 3
- [21] Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. Glisten: Generalization based data subset selection for efficient and robust learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021. 3
- [22] Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial

- active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [23] Andreas Kirsch, Joost van Amersfoort, and Yarín Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, 2019. 3, 6
- [24] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*. 2017. 2
- [25] Radek Mackowiak, Philip Lenz, Omair Ghori, Ferran Diego, Oliver Lange, and Carsten Rother. Cereals-cost-effective region-based active learning for semantic segmentation. *arXiv preprint arXiv:1810.09726*, 2018. 1
- [26] Sudhanshu Mittal, Maxim Tatarchenko, and Thomas Brox. Semi-supervised semantic segmentation with high-and low-level consistency. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1369–1379, 2019. 2, 4, 12
- [27] Sudhanshu Mittal, Maxim Tatarchenko, Özgün Çiçek, and Thomas Brox. Parting with illusions about deep active learning. *CoRR*, abs/1912.05361, 2019. 3, 8
- [28] Prateek Munjal, Nasir Hayat, Munawar Hayat, Jamshid Sourati, and Shadab Khan. Towards robust and reproducible active learning using neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [29] Aneesh Rangnekar, Christopher Kanan, and Matthew Hoffman. Semantic segmentation with active semi-supervised learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5966–5977, January 2023. 2, 3
- [30] Phill Kyu Rhee, Enkhbayar Erdenee, Shin Dong Kyun, Minhaz Uddin Ahmed, and Songguo Jin. Active and semi-supervised learning for object detection with imperfect data. *Cognitive Systems Research*, 45:109–123, 2017. 3
- [31] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017. 3
- [32] Claude Elwood Shannon. A mathematical theory of communication. *ACM SIGMOBILE mobile computing and communications review*, 5(1):3–55, 2001. 2, 3
- [33] Gyungin Shin, Weidi Xie, and Samuel Albanie. All you need are a few pixels: Semantic segmentation with pixelpick. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2021. 2, 3
- [34] Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled method for query and training. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, 2020. 1, 3
- [35] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 1, 3
- [36] Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang, and Liang Lin. Cost-effective active learning for deep image classification. *IEEE Trans. Cir. and Sys. for Video Technol.*, 2017. 2
- [37] Shuo Wang, Yuexiang Li, Kai Ma, Ruhui Ma, Haibing Guan, and Yefeng Zheng. Dual adversarial network for deep active learning. In *ECCV*, 2020. 1, 3
- [38] Zifeng Wu, Chunhua Shen, and Anton van den Hengel. Wider or deeper: Revisiting the resnet model for visual recognition. *Pattern Recognition*, 90:119–133, 2019. 5, 12
- [39] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. BDD100K: A diverse driving video database with scalable annotation tooling. *CoRR*, abs/1805.04687, 2018. 2
- [40] Fedor Zhdanov. Diverse mini-batch active learning. *CoRR*, abs/1901.05954, 2019. 3

## A. Datasets Visualization

Figure 6 shows examples of the A2D2 and the Cityscapes dataset. Each row shows three temporally consecutive frames in both labeled datasets. We clearly observe that the images in the A2D2 dataset have high-overlapping information, whereas images in the Cityscapes dataset are quite diverse. Therefore, to create our redundancy experiments, we chose the A2D2 dataset as the base dataset.

## B. Training details

We used the DeepLabv3+ [5] architecture with WideResNet38 (WRN-38) [38] backbone for all our experiments. The backbone WRN-38 is pre-trained using ImageNet [7]. For the supervised learning setting, the model is trained using the SGD optimizer with a base-learning rate of  $1e - 3$ , momentum of 0.9, and a weight decay of  $5e - 4$ . We utilize a polynomial learning rate scheduler with a batch size of 8 and train a model in each AL cycle for 100 epochs. The model is trained with data augmentations, including random cropping and random horizontal flipping. Input image size is  $256 \times 512$  for Cityscapes and A2D2 datasets and  $321 \times 321$  for the PASCAL-VOC dataset. We utilize the s4GAN [26] method for semi-supervised learning (SSL). We use the same training setting for the segmentation model as in the supervised setting. We use the same hyperparameters as mentioned in [26], except for the learning rate of the discriminator which is set to  $2.5e - 5$  for Cityscapes and A2D2 experiments. We add 3 dropout layers with a dropout rate of 0.1 in the decoder of the segmentation model for all the MCD-based AL methods.

## C. Evaluation Metric: AUC@B

We use the following formula to compute the Area Under the Budget Curve(AUC@B) at a total budget B, where B is the percentage of the labeled dataset:

$$AUC@B = \sum_{i=1}^{i=N} \frac{(b_{i+1} - b_i)(p_i + p_{i+1})}{2} \quad (1)$$

,where N is the number of AL acquisition steps,  $b_i$  is the percentage of labeled dataset at step  $i$ , and  $p_i$  is the performance of the model in mIoU(%) at step  $i$ .

## D. Results: AL under Different Budgets

Here, we show the performance curves and tables for different active learning methods in a low annotation budget setting. This section also contains the remaining curves and tables in a high annotation budget setting, discussed in the main paper.



Figure 6: Consecutive images from the Cityscapes and A2D2 datasets. This shows even the consecutive images in the Cityscapes dataset are different and diverse, whereas consecutive frames in the A2D2 dataset are very similar, containing redundant information.

A	AL Method Metric →	SSL	PASCAL: 2-2	
			mIoU@6	AUC@10
S	Random	✗	66.41	5.22
S	Entropy	✗	66.33	5.11
S	EqualAL	✗	65.04	5.13
B	Coreset	✗	66.24	5.19
S	Random-SSL	✓	<b>68.60</b>	<b>5.37</b>
S	Entropy-SSL	✓	67.26	5.31
S	EqualAL-SSL	✓	67.44	5.31
B	Coreset-SSL	✓	68.03	5.35
-	100%	✗	77.00	6.16

Table 7: Active Learning results on PASCAL-VOC dataset in low-budget 2-2 setting. AUC@10 and mIoU@6 metric are reported. A denotes Acquisition method type. S and B denotes the single-sample and batch-based acquisition, respectively.

### D.1. Low Budget

For diverse datasets, we evaluate active learning acquisition methods on PASCAL-VOC, A2D2-Pool-0f, and Cityscapes datasets. For redundant datasets, we show results on the A2D2-Pool-11f dataset only.

**PASCAL-VOC:** We show results for the AL methods in a low-budget setting on the PASCAL-VOC dataset. We consider the 2-2 setting as the low-budget setting, where the maximum annotation budget is 10% of the dataset size. We find that the random-SSL method performs the best. It shows that none of the AL bias is correctly learned or helpful in such a low-budget setting for such a diverse dataset.



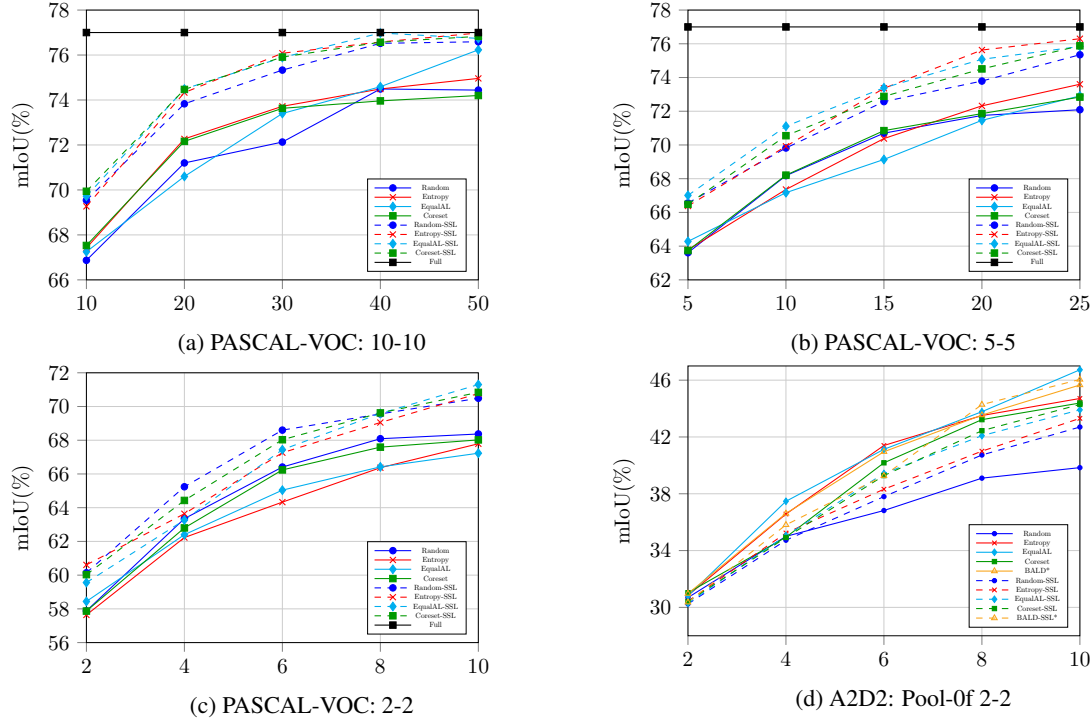


Figure 7: Active learning performance curves on PASCAL-VOC and A2D2:Pool-Of. X-axis shows the percentage of labeled dataset. The methods which utilize MC-Dropout in their network architecture are marked with \*, and are only comparable to other methods with MC-Dropout.

A	AL Method Metric →	SSL	A2D2 Pool-Of 2-2		A2D2 Pool-11f 2-2	
			mIoU@6	AUC@10	mIoU@6	AUC@10
S	Random	✗	36.82	2.92	37.74	2.93
S	Entropy	✗	<b>41.40</b>	3.18	36.37	2.92
S	EqualAL	✗	41.13	<b>3.22</b>	37.28	2.97
B	Coreset	✗	40.18	3.12	<b>39.63</b>	<b>3.10</b>
S	Random-SSL	✓	37.80	2.99	36.46	2.90
S	Entropy-SSL	✓	38.32	3.03	36.70	2.93
S	EqualAL-SSL	✓	39.43	3.07	36.31	3.06
B	Coreset-SSL	✓	39.28	3.08	39.20	3.06
-	100%	✗	56.87	4.55	48.85	3.91

Table 8: Active Learning results on A2D2 Pool-Of in 2-2 setting. AUC@10 and mIoU@6 metrics are reported. A denotes Acquisition method type. S and B denotes the single-sample and batch-based acquisition, respectively.

The integration of semi-supervised learning with random selection improves the performance over the supervised random sampling baseline. Results are shown in Table 7 and Figure 7c.

**Cityscapes:** Table 9 show results on the Cityscapes dataset in a low annotation budget, 2-2 setting. We find that the single-sample-based method performs the best. SSL integration with active learning is only useful for the batch-based Coreset approach, whereas it is detrimental for other acquisition functions.

**A2D2 Pool-Of:** We show results in the low-budget, 2-2 setting for the A2D2 Pool-Of. We find that single-sample-based methods outperform all the methods. This shows that active learning is again successful in a low-budget setting when the dataset only covers a specific domain, like the driving scenario in this case. However, semi-supervised learning does not help in this case. Results are shown in Table 8 and Figure 7d.

**A2D2 Pool-11f:** Table 8 show results on redundant datasets A2D2-pool-11f in low annotation budget, 2-2 setting. We

A	AL Method Metric →	SSL	Cityscapes: 2-2	
			mIoU@6	AUC@10
S	Random	✗	46.05	3.65
S	Entropy	✗	<b>51.24</b>	<b>4.00</b>
B	Coreset	✗	47.26	3.74
S	Random-SSL	✓	47.46	3.72
S	Entropy-SSL	✓	49.99	3.93
B	Coreset-SSL	✓	48.51	3.82
-	100%	✗	68.42	5.47

Table 9: Active Learning results on Cityscapes dataset in low-budget 2-2 setting. AUC@10 and mIoU@6 metric are reported. A denotes Acquisition method type. S and B denotes the single-sample and batch-based acquisition, respectively.

find that batch-based methods outperform all the methods. Redundant datasets still favor the batch-based acquisition method in the low-budget setting. However, SSL does not contribute any additional improvements due to insufficient labeled samples to support learning from unlabeled samples.

As discussed in the main paper, active learning methods are highly sensitive to distribution change w.r.t. levels of redundancy in the low-budget setting. The ideal policy transitions from random selection to single-sample acquisition and then to batch-based acquisition as the level of redundancy in the dataset goes from low to high.

## D.2. High-budget

**PASCAL-VOC:** Figure 7a and Figure 7b show the AL performance curves for 10-10 and 5-5 settings on the PASCAL-VOC dataset, respectively. We observe that the single-sample-based methods with semi-supervised learning perform the best. Performance tables for these settings are included in the main paper.