

000  
001  
002052  
053  
054

# Stream-Based Active Distillation for scalable model deployment

003  
004  
005  
006  
007  
008  
009  
010  
011055  
056  
057  
058  
059  
060  
061  
062  
063

Anonymous CVPR submission

Paper ID \*\*\*\*\*

## Abstract

This paper proposes a scalable technique for developing lightweight yet powerful models for object detection in videos using self-training with knowledge distillation. This approach involves training a compact student model using pseudo-labels generated by a large and powerful teacher model, which can help to reduce the need for massive amounts of data and computational power. However, model-based annotations in large-scale applications may propagate errors or biases, and incur high costs. To address these issues, the paper introduces Stream-Based Active Distillation (SBAD) to endow students with effective, query-efficient methods that are robust to teacher imperfections. The proposed pipeline considers on-the-fly decision-making to query or not an image from a video to avoid storing data. Simple yet effective strategies are compared, demonstrating 1) the effectiveness of distillation for low annotation budget and 2) the importance of querying images with the highest confidence.

## 1. Introduction

Deep Neural Networks (DNNs) are effective for object detection in images, but their predictive power comes at a high cost. The training of highly performant DNNs requires high-performance cloud servers with a large-scale dataset which requires (*i*) a large workforce to prepare the dataset or training implementation (*ii*) as well as a significant investment of time and money. These data, time, and hardware costs create a barrier for most practitioners in terms of going from theory to practice [5]. Furthermore, a one-shot investment in resources to create large general-purpose models, regardless of their size, is no longer sufficient. With no retraining, these models cannot be robust with respect to **stochastic** and **ever-**

**evolving** environments. In the example of Closed-Circuit Television (CCTV) monitoring traffic at city-scale, there is no dataset large enough to cover all aspects of every urban landscape [33]. Therefore, a *scalable, efficient, and recurrent* retraining is necessary to reduce costs and avoid **under-performing** systems.

Knowledge Distillation (KD) is a promising technique that enables the creation of lightweight, yet powerful models. The process assumes that for the same dataset, large models (i.e., *teachers*) have higher knowledge capacity than smaller models (i.e., *students*). The teacher, typically a pre-trained or very large generic model (e.g., YOLOv8x6), can transfer its knowledge (i.e., pattern recognition mechanisms) to students without significant model degradation. However, the recourse to other models for labelling could lead to confirmation bias, a phenomenon referring to the noise accumulation when the model is trained using incorrect predictions for semi-supervised or unsupervised learning [2]. Furthermore, an immediate rebound effect of the scheme is the multiplication, at scale, of the number of models to be trained. The inference costs could become significant. Additionally, if the teacher model is running on a cloud-based platform, there may be additional costs associated with its usage, such as hourly usage fees or data transfer costs. This could be mitigated by using Active Learning (AL), which aims to identify the most informative examples for labelling. The importance of sampling has been first formulated in [4] as the problem of developing KD methods that are both query-efficient, and robust to labelling inaccuracies due to the imperfection of the teacher (i.e., *confirmation bias*). The method developed in [4] was designed for a pool-based setting, which represents an offline scenario where a pool of unlabeled data points is made available to the learner. We claim, in many real-world applications, a large number of unlabeled samples arrive in a streaming manner, mak-

104 **ing it impossible to maintain all the data in a can-**  
 105 **didate pool.** At the best of our knowledge, there ex-  
 106 ist not a framework supporting the development of  
 107 AL methods that are both query-efficient, and ro-  
 108 bust to labelling inaccuracies in stream-based set-  
 109 tings. This paper contribution are the following :  
 110

- 111 1. Formulate Stream-Based Active Distillation  
   112 (SBAD) as the problem of developing AL meth-  
   113 ods that are both query-efficient, and robust to  
   114 labelling inaccuracies, in stream-based settings.  
   115
- 116 2. Demonstrate the scheme applicability and bene-  
   117 fits for large-scale video-based object detections  
   118 on a public dataset [24].  
   119
- 120 3. Establish simple but effective baselines to train  
   121 a YOLOv8n student from a YOLOv8x6 teacher.  
   122 For that purpose, a GIT repository is provided.  
   123

## 124 2. Related Work

### 125 2.1. Knowledge Distillation

126 KD is a method that involves training a smaller  
 127 model to imitate the performance of a larger model.  
 128 The main objectives of this technique are to prevent  
 129 a decrease in the model’s performance when it op-  
 130 erates on a dataset that is distributed differently than  
 131 the source domain, referred to as Unsupervised Do-  
 132 main Adaptation (UDA), and to produce lightweight  
 133 models suitable for the storage and computational ca-  
 134 pacities of miniaturized devices, referred to as Model  
 135 Compression (MC) applications. In this study, we  
 136 utilize a technique called *Self-training with knowl-*  
 137 *edge distillation*, which was introduced by [26]. This  
 138 technique trains a student model using pseudo-labels  
 139 generated by a teacher model, which is beneficial  
 140 when the labelled data is limited. However, it may  
 141 also propagate errors or biases. Besides, we will dis-  
 142 cuss two additional techniques of interest in the fol-  
 143 lowing paragraphs: online distillation and context-  
 144 aware distillation.  
 145

146 **Online Distillation.** This approach involves train-  
 147 ing a smaller student model to mimic the output of a  
 148 larger teacher model on a per-example basis. In [11],  
 149 the authors designed an online knowledge distillation  
 150 scheme to perform real-time human segmentation in  
 151 sports videos. Experiments show the model’s ability  
 152 to adapt to variations occurring in the context. Online  
 153 distillation is also employed in [22] to adapt a low-  
 154 cost semantic segmentation model to a target video

155 where the data distribution is not necessarily station-  
 156 ary.  
 157

158 **Context-aware Distillation.** The researches [17,  
 159 26] attempt to exploit the contextual features of the  
 160 scene to develop effective KD. They directly worked  
 161 on the distillation scheme to develop more special-  
 162 ized students. For instance, in [17], they added a  
 163 temporal dimension such that the student learns the  
 164 variations in the teacher’s intermediate features over  
 165 time. They thus take into account frame redundan-  
 166 cies within a CCTV stream.  
 167

### 168 2.2. Active Learning

169 AL is a sampling approach that selects the most  
 170 informative data points to minimize the number of  
 171 labels required for model training [31]. AL can be di-  
 172 vided into three macro scenarios: membership query  
 173 synthesis, pool-based AL, and stream-based AL [6].  
 174 The majority of approaches in deep AL have focused  
 175 on the pool-based scenario, where the learner selects  
 176 the most useful data points from a closed set of un-  
 177 labeled observations. The stream-based AL scenario  
 178 for object detectors has not been investigated. More-  
 179 over, AL assumes the availability of a perfect oracle,  
 180 where the true label of a data point is revealed when  
 181 queried. However, this assumption does not hold in a  
 182 Knowledge Distillation (KD) framework, where the  
 183 pseudo-labels provided by the teacher may be incor-  
 184 rect.  
 185

186 **Active Learning for Image Classification.** AL  
 187 strategies for pool-based classification can be cate-  
 188 gorized into uncertainty-based or diversity-based ap-  
 189 proaches [34]. Uncertainty-based strategies estimate  
 190 model uncertainty using techniques such as Monte  
 191 Carlo dropout [16] or ensemble networks [21], while  
 192 entropy and margin-based sampling strategies are  
 193 also widely employed [27]. Task-agnostic methods  
 194 like Learn loss [36] use a loss prediction module to  
 195 estimate data points that are likely to be wrongly  
 196 predicted. Among diversity-based strategies, Core-  
 197 set [30] is one of the most popular, using a K-center  
 198 Greedy algorithm to locate a set of representative  
 199 data points. Cluster-Margin [12] combines uncer-  
 200 tainty and diversity, while DRMRS [14] takes margin  
 201 and diversity into account. BADGE [3] balances un-  
 202 certainty and diversity using a  $k$ -MEANS++ seeding  
 203 algorithm on gradients obtained from the last layer  
 204 of the network. CDAL [1] replaces Euclidean dis-  
 205 tance with pairwise contextual diversity in the greedy  
 206

208 k-center algorithm used in Core-set. CLUE [23] performs  
 209 uncertainty-weighted clustering to identify target  
 210 instances that are both uncertain according to the  
 211 model and diverse in feature space. VAAL [32] uses  
 212 a Variational Autoencoder (VAE) to map instances  
 213 into a latent space, which is then fed into a discrimi-  
 214 nator that learns to differentiate between labeled data  
 215 and unlabeled samples.  
 216

217 **Active Learning for Object Detection.** AL ap-  
 218 proaches for object detection can be classified into  
 219 black-box and white-box methods [28]. Black-box  
 220 methods do not depend on the underlying network ar-  
 221 chitecture and use informativeness scores such as the  
 222 confidence obtained from the softmax layer, while  
 223 white-box methods are dependent on the architecture  
 224 of the underlying network. The minmax approach,  
 225 which selects the least confident images among the  
 226 unlabeled pool, is a popular black-box method [28].  
 227 Ensemble methods have also been used for object  
 228 detection-oriented AL [15, 29]. Query strategies  
 229 based on localization tightness and stability [19],  
 230 mixture density networks [10], and a unified box re-  
 231 gression and classification metric [37] has also been  
 232 proposed. MIAL [38] is a multiple-instance frame-  
 233 work that filters out noisy instances to bridge the  
 234 gap between instance-level and image-level uncer-  
 235 tainty. PPAL [35] is a two-stage algorithm that in-  
 236 cludes difficulty-calibrated uncertainty sampling and  
 237 category-conditioned matching similarity. [18] pro-  
 238 posed to cluster the unlabeled observations in groups  
 239 based on the frequency domain values and to use dif-  
 240 ferent sampling rates for each group.  
 241

### 2.3. Challenges of Stream-based Active Distil- 242 lation

244 The importance of sampling has been first formu-  
 245 lated in [4] as the problem of developing KD meth-  
 246 ods that are both query-efficient and robust to la-  
 247 beling inaccuracies due to the imperfection of the  
 248 teacher (i.e., *confirmation bias*). Their methods pro-  
 249 vide a theoretical guarantee that the scheme leads  
 250 to queries where the teacher provides correct labels.  
 251 However, this approach has been developed in a  
 252 pool-based setting where the student has access to the  
 253 entire information pool. Therefore, techniques such  
 254 as diversity-based strategies, clustering, or pairwise  
 255 distance matrices may not be feasible, especially in  
 256 contexts with spatio-temporal correlation among the  
 257 data that could be exploited. Another aspect is that,  
 258 due to the complexity of the student model, uncer-  
 259 tainty techniques relying on Monte Carlo dropout or

260 Learn loss modules may not be viable options.  
 261

## 3. Problem Statement

264 Let  $\theta_{student}^{base}$  define a compact general pre-trained  
 265 model learning the distribution  $\mathcal{D}$  of a data stream  
 266  $\mathcal{X}$ . We assume a spatio-temporal correlation among  
 267 the data. The student is endowed with SELECT ( $I_t$ ),  
 268 a rule that decides whether to query the image for  
 269 pseudo-labelling by a powerful but imperfect model  
 270  $\theta_{teacher}^{general}$ . The objective is to train a high-performing  
 271 student by querying the minimum number of teacher  
 272 soft-labels. We assume a large-scale setting (e.g.,  
 273 city-scale deployment of CCTV, monitoring of large  
 274 construction sites) and affordable hardware. There-  
 275 fore, the pseudo-labels constituting a training set  $\mathcal{L}$   
 276 must not exceed a maximal budget per student  $B$ ,  
 277 i.e.,  $|\mathcal{L}| \leq B$ . It is important that the SELECT strate-  
 278 gies used in SBAD are computationally efficient for  
 279 two reasons. Firstly, to prevent adding unnecessary  
 280 complexity to the student models. Secondly, to allow  
 281 for quick sampling decisions. If the decision time  
 282 for a selection rule is higher than the frame rate used  
 283 by the CCTV system, it would require the system to  
 284 have a buffer for temporarily storing recently seen  
 285 images until the sampling decision is made. This  
 286 would inevitably result in the system requiring addi-  
 287 tional resources to store and process the data, which  
 288 is not feasible in large-scale settings where hardware  
 289 resources are limited.  
 290

---

### Algorithm 1 SBAD Framework

---

291 **Require:** a pre-trained student model  $\theta_{student}^{base}$ , a  
 292 general purpose teacher model  $\theta_{teacher}^{general}$ , a budget  
 293  $B$ , a SELECT strategy.  
 294 **Ensure:**  $B \geq 1$

$\mathcal{L} \leftarrow \emptyset$	▷ Selected frames
$t \leftarrow 0$	▷ Timestamp
<b>while</b> $ \mathcal{L}  \leq B$ <b>do</b>	
Observe current frame $I_t$	
<b>if</b> SELECT( $I_t$ ) <b>is TRUE then</b>	
$\{b_i\} \leftarrow \theta_{teacher}(I_t)$	▷ Pseudo-labels
$\mathcal{L} \leftarrow \mathcal{L} \cup (I_t, \{b_i\}_t)$	▷ Pseudo-labels
<b>end if</b>	
$t \leftarrow t + 1$	▷ Timestamp
<b>end while</b>	
<b>return</b> update( $\theta_{student}^{base}$ , $\mathcal{L}$ )	

---

308 Figure 1. provides a visual illustration of the  
 309 SBAD framework. The sampling phase is a critical  
 310 step in which frames are selected based on instance  
 311

312 selection criteria to obtain the most informative sam-  
 313 ples. The selected frames are then queried by the  
 314 teacher model to generate pseudo-labels that are sub-  
 315 sequently used for fine-tuning the student models. In  
 316 the model development phase, the fine-tuned models  
 317 can be evaluated using hard-labels, but in real-world  
 318 deployment scenarios, there are typically few or no  
 319 labels available, making this step impractical.  
 320

## 321 4. Methodology

322 In the context of stream-based active learning,  
 323 the single-pass evaluation of data points is often ad-  
 324 dressed by applying a threshold to certain informa-  
 325 tiveness scores [7–9, 13, 25]. However, this approach  
 326 has not been tested in online active distillation tasks  
 327 for object detection. In this paper, we investigate the  
 328 effectiveness of thresholding algorithms based on the  
 329 confidence of the base student model  $\theta_{student}^{base}$  for the  
 330 SBAD framework. At round  $t$ , when the  
 331 student model  $\theta_{student}^{base}$  observes an image  $I_t$ ,  
 332  $n \geq 0$  objects are detected, which are defined by the  
 333 bounding boxes  $b_{it}$  and confidence scores  $c_{it}$ . Ac-  
 334 cording to [28], a unique confidence score  $C_t$  for  $I_t$   
 335 can be obtained using:  
 336

$$338 C(I_t) := \max_i c_{it}$$

339 This means that the confidence of each image is  
 340 approximated by the highest confidence score among  
 341 the objects detected in that image. By using this con-  
 342 fidence metric, we can then apply a threshold  $\Delta$  to  
 343 the confidence scores of the incoming frames. The  
 344 overall structure of the top confidence threshold sam-  
 345 pling scheme is presented in Algorithm 1. To esti-  
 346 mate the threshold  $\Delta$  for selecting the most infor-  
 347 mative frames, we introduce a warm-up phase where  
 348 the student model  $\theta_{student}^{base}$  observes incoming  
 349 frames for a period of length  $w$  without querying  
 350 any image and without storing anything other than  
 351 a single scalar representing the image-level confi-  
 352 dence scores  $C(I_t)$ , where  $t = 1, \dots, w$ . At the end  
 353 of the warm-up phase, the student model estimates  
 354 an  $(1 - \alpha)$ -upper percentile on the distribution of  
 355 the confidence scores, where  $\alpha$  represents the desired  
 356 sampling rate. In other words, we aim to select the  
 357 frames to pseudo-label and fine-tune  $\theta_{student}^{base}$  with a  
 358 ratio of  $\alpha$  frames out of the total number of  
 359 frames. While in traditional AL, the focus is on  
 360 querying images that the student model is least con-  
 361 fident about, this approach may not be optimal for  
 362 stream-based KD scenarios. The least confident im-  
 363 ages often correspond to very hard examples that

364 may not be informative enough for the student model  
 365 in the early rounds of AL when it has not been fine-  
 366 tuned for the specific scene. Additionally, select-  
 367 ing images with high uncertainty for pseudo-labeling  
 368 may lead to confirmation bias as the pseudo-labels  
 369 may not align with the ground truth due to the im-  
 370 perfection of the teacher model  $\theta_{teacher}^{general}$  as an oracle.  
 371 To address these issues, we propose to let the stu-  
 372 dent model  $\theta_{student}^{base}$  query the most confident frames  
 373 based on a threshold  $\Delta$  given by:  
 374

$$375 \mathbb{P}(C(I_t) \geq \Delta) = \alpha$$

376 ”Distillation” Ideally, using this threshold, the student  
 377 will sample informative examples that the teacher  
 378 model would properly pseudo-label and that will  
 379 contribute to faster fine-tuning while avoiding frames  
 380 that are too uncertain to be used in the initial stages  
 381 of AL.

## 382 5. Experiments

### 383 5.1. Experimental Settings

384 **Dataset** In order to evaluate the effectiveness of  
 385 the SBAD approach, we utilized the WALT dataset  
 386 [24], which consists of multi-camera recordings. The  
 387 dataset provides spatial and temporal diversity, in-  
 388 cluding various points of view and lighting condi-  
 389 tions (day and night).

390 **Distillation implementation** In line with the prin-  
 391 ciples of data distillation proposed by Rivas et al.  
 392 [26], we employ a large and complex teacher model,  
 393 YOLOv8x6 (261.1 GFLOPs), to generate pseudo-  
 394 labels. These labels are then utilized to train several  
 395 smaller student models, YOLOv8n (8.7 GFLOPs),  
 396 with less architectural complexity. Both networks are  
 397 initially pre-trained on the COCO dataset [20]. The  
 398 student models are re-trained for 100 epochs with a  
 399 batch size of 16. The budget of the SBAD frame-  
 400 work is determined by the number of pseudo-labels  
 401 used for fine-tuning, which ranges from 25 to 250 in  
 402 our experiments.

403 **Methods.** As the SBAD problem for object detec-  
 404 tion has not been explored before, there are not many  
 405 state-of-the-art methods that can be used. To explore  
 406 the effectiveness of the confidence-based threshol-  
 407 ding algorithm, we used different baselines. First,  
 408 a naïve *N-First* approach has been implemented,  
 409 where the student models are fine-tuned by simply

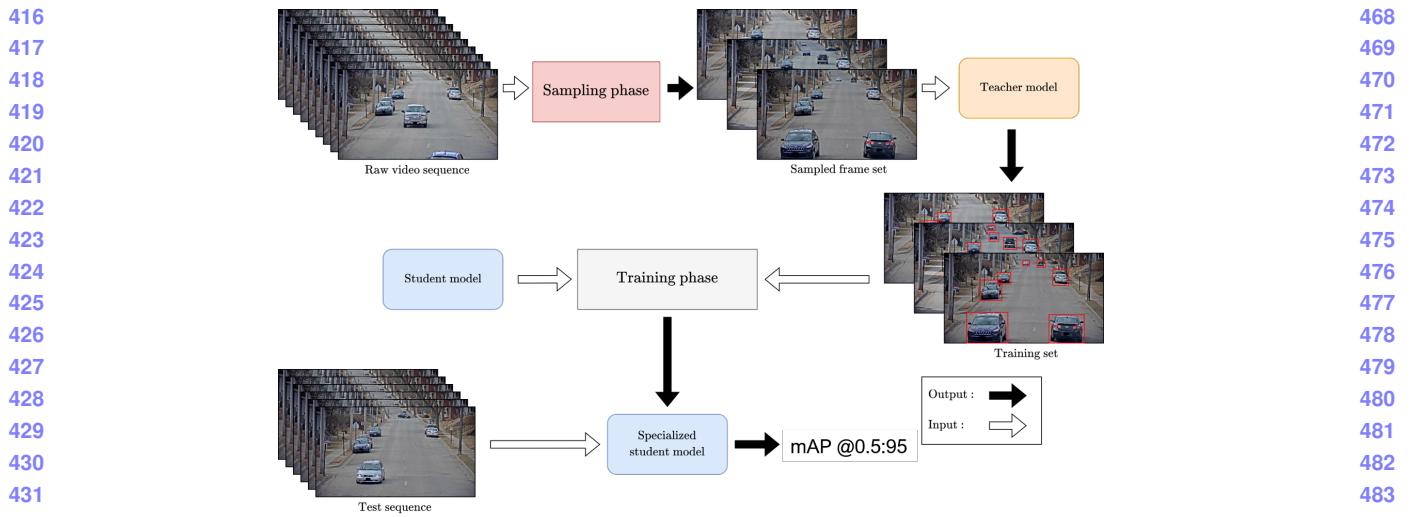


Figure 1. SBAD pipeline: sampling, fine-tuning, and evaluation.

taking the first  $N$  images observed from each camera. A second baseline is given by a *random* sampling approach, where a number  $s \sim U(0, 1)$  is generated for each incoming frame, which is queried only if  $s \geq 1 - \alpha$ . A third baseline is given by a more active learning-oriented *least confidence* approach, where similarly to the top confidence case, we impose a threshold on the image-level confidence score. The main difference is that the threshold  $\Delta$  is estimated by taking the  $\alpha$ -lower percentile from the warm-up set  $\mathcal{W}$ .

## 5.2. Experimental Results

Figures 2. and 3. shows the learning curves obtained using stream-based active learning techniques on the WALT dataset. Our analysis can be approached from two perspectives. Firstly, from a knowledge distillation standpoint, we observed how the student model's performance improves as we use more frames for fine-tuning. In particular, we found that the mAP<sub>50-95</sub> score approaches that of the teacher model when 250 pseudo-labeled frames are used. However, we also noticed that the student's performance deteriorates significantly when only a small number of frames are used for fine-tuning, which could be attributed to overfitting due to the limited number of images presented to the network. In addition, if the model is fine-tuned on images biased towards a specific time of day, such as only night or day, it may perform poorly on the balanced test set used for evaluation. Furthermore, as depicted in Figure 4., choosing highly uncertain im-

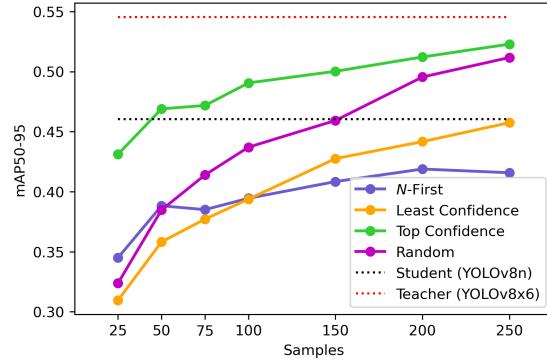


Figure 2. Learning curves obtained on the first two cameras of WALT. Results show that increasing the number of frames used for fine-tuning improves the student model's performance, approaching that of the teacher model with 250 frames. However, using only a small number of frames may lead to overfitting and poor performance on balanced evaluation sets. Top confidence thresholding is more effective than least confidence-based methods for stream-based active learning, highlighting the importance of avoiding highly uncertain images during fine-tuning.

ages for pseudo-labeling may lead to incorrect labels due to the teacher's own bad prediction.

From an active learning perspective, the performance achieved with the *top confidence threshold* algorithm is significantly better than that obtained using the least confidence-based method. This highlights the importance of fine-tuning the model with highly certain images, especially when the model has not yet been specialized for the scene.

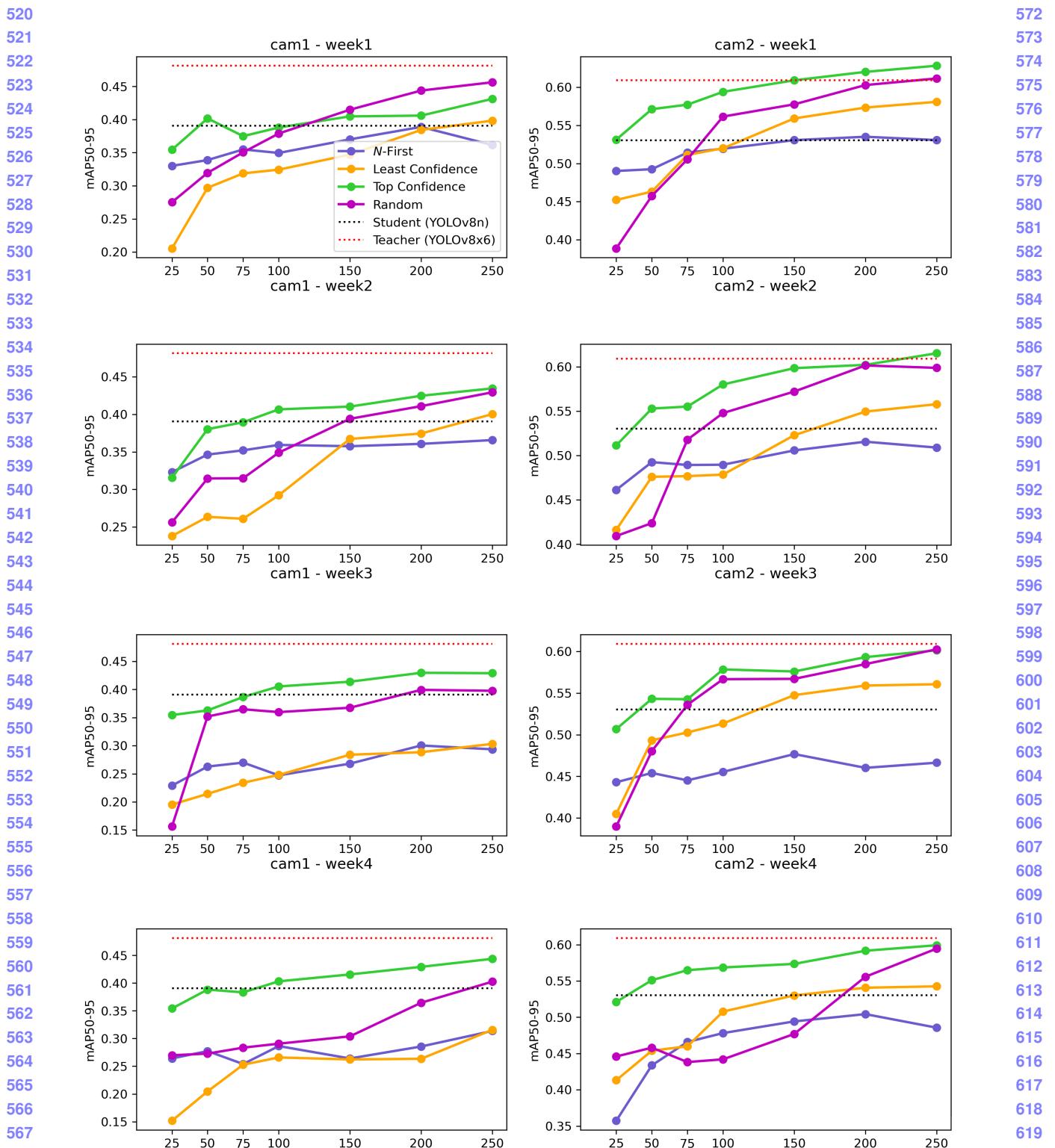


Figure 3. Weekly analysis on the first two cameras of WALT.



Figure 4. Two difficult examples (one for each camera) that lead to *confirmation bias*: when the student requests highly uncertain images based on its predictions (in yellow), wrong pseudo labels are revealed (in red).

## 6. Conclusion

This paper proposes SBAD to bridge the gap between large-scale and affordable deep learning models while adapting to changing environments. This framework enables the scalable deployment of deep learning models under tight budget constraints.

The framework evaluates the informativeness of each frame, accounting for teacher imperfections in a KD scheme. Experiments demonstrate that traditional AL strategies may not be optimal for KD. Future research could explore alternative sampling strategies and distillation mechanisms to improve performance.

## References

- [1] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *European Conference on Computer Vision (ECCV) 2020*, 8 2020. 2
- [2] Eric Arazo, Diego Ortego, Paul Albert, Noel E O’Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2020. 1
- [3] Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *2020 International Conference on Learning Representations*, 6 2019. 2
- [4] Cenk Baykal, Khoa Trinh, Fotis Iliopoulos, Gaurav Menghani, and Erik Vee. Robust active distillation. *arXiv preprint arXiv:2210.01213*, 2022. 1, 3
- [5] Lucas Beyer, Xiaohua Zhai, Amélie Royer, Larisa Markeeva, Rohan Anil, and Alexander Kolesnikov.

Knowledge distillation: A good teacher is patient and consistent. *CoRR*, abs/2106.05237, 2021. 1

- [6] Davide Cacciarelli and Murat Kulahci. A survey on online active learning. <https://arxiv.org/abs/2302.08893>, 2023. 2
- [7] Davide Cacciarelli, Murat Kulahci, and John Tyssedal. Online active learning for soft sensor development using semi-supervised autoencoders. In *ICML 2022 Workshop on Adaptive Experimental Design and Active Learning in the Real World*, 12 2022. 4
- [8] Davide Cacciarelli, Murat Kulahci, and John Sølve Tyssedal. Stream-based active learning with linear models. *Knowledge-Based Systems*, 254:109664, 10 2022. 4
- [9] Andrea Castellani, Sebastian Schmitt, and Barbara Hammer. Stream-based active learning with verification latency in non-stationary environments. In *Artificial Neural Networks and Machine Learning 2022*, 4 2022. 4
- [10] Jiwong Choi, Ismail Elezi, Hyuk-Jae Lee, Clément Farabet, and Jose M. Alvarez. Active learning for deep object detection via probabilistic modeling. *CoRR*, abs/2103.16130, 2021. 3
- [11] Anthony Cioppa, Adrien Deliege, Maxime Is-tasse, Christophe De Vleeschouwer, and Marc Van Droogenbroeck. Arthus: Adaptive real-time human segmentation in sports through online distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 2
- [12] Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Ros-tamizadeh, and Sanjiv Kumar. Batch active learning at scale. In *Conference on Neural Information Processing Systems*, 7 2021. 2

- 728 [13] Sanjoy Dasgupta, Adam Kalai, and Claire Mon-  
729 teleoni. Analysis of perceptron-based active learning.  
730 In *Lecture Notes in Computer Science*, volume 10, 12  
731 2005. 4
- 732 [14] Ehsan Elhamifar, Guillermo Sapiro, Allen Yang, and  
733 S. Shankar Sasry. A convex optimization framework  
734 for active learning. In *Proceedings of the IEEE Interna-*  
735 *tional Conference on Computer Vision*, pages 209–  
736 216. Institute of Electrical and Electronics Engineers  
737 Inc., 2013. 2
- 738 [15] Di Feng, Xiao Wei, Lars Rosenbaum, Atsuto Maki,  
739 and Klaus Dietmayer. Deep active learning for effi-  
740 cient training of a lidar 3d object detector. In *30th*  
741 *IEEE Intelligent Vehicles Symposium*, 2019. 3
- 742 [16] Yarin Gal, Riashat Islam, and Zoubin Ghahramani.  
743 Deep bayesian active learning with image data. In  
744 *Proceedings of the 34th International Conference on*  
745 *Machine Learning*, 2017. 2
- 746 [17] Amirhossein Habibian, Haitam Ben Yahia, Davide  
747 Abati, Efstratios Gavves, and Fatih Porikli. Delta dis-  
748 tillation for efficient video processing, 2022. 2
- 749 [18] Wei Huang, Shuzhou Sun, Xiao Lin, Dawei Zhang,  
750 and Lizhuang Ma. Deep active learning with weight-  
751 ing filter for object detection. *Displays*, page 102282,  
752 1 2022. 3
- 753 [19] Chieh-Chi Kao, Teng-Yok Lee, Pradeep Sen, and  
754 Ming-Yu Liu. Localization-aware active learning for  
755 object detection. In *Asian Conference on Computer*  
756 *Vision (ACCV) 2018*, 1 2018. 3
- 757 [20] Tsung-Yi Lin, Michael Maire, Serge J. Belongie,  
758 Lubomir D. Bourdev, Ross B. Girshick, James Hays,  
759 Pietro Perona, Deva Ramanan, Piotr Dollár, and  
760 C. Lawrence Zitnick. Microsoft COCO: common ob-  
761 jects in context. *CoRR*, abs/1405.0312, 2014. 4
- 762 [21] Salman Mohamadi, Gianfranco Doretto, and Don-  
763 ald A Adjeroh. Deep active ensemble sampling for  
764 image classification. In *16th Asian Conference on*  
765 *Computer Vision (ACCV 2022)*, 2022. 2
- 766 [22] Ravi Teja Mullapudi, Steven Chen, Keyi Zhang,  
767 Deva Ramanan, and Kayvon Fatahalian. Online  
768 model distillation for efficient video inference. In  
769 *Proceedings of the IEEE/CVF International Confer-  
770 ence on Computer Vision*, pages 3573–3582, 2019.  
771 2
- 772 [23] Viraj Prabhu, Arjun Chandrasekaran, Kate Saenko,  
773 and Judy Hoffman. Active domain adaptation via  
774 clustering uncertainty-weighted embeddings. In *Inter-  
775 national Conference on Computer Vision (ICCV)*  
776 2021, 2020. 3
- 777 [24] N. Dinesh Reddy, Robert Tamburo, and Srinivasa G.  
778 Narasimhan. Walt: Watch and learn 2d amodal repre-  
779 sentation from time-lapse imagery. In *Proceedings of*  
780 *the IEEE/CVF Conference on Computer Vision and*  
781 *Pattern Recognition (CVPR)*, pages 9356–9366, June  
782 2022. 2, 4
- 783 [25] Carlos Riquelme, Ramesh Johari, and Baosen Zhang.  
784 Online active linear regression via thresholding. In  
785 *31st AAAI Conference on Artificial Intelligence*,  
786 2017. 4
- 787 [26] Daniel Rivas, Francesc Guim, Jordà Polo, Pubudu M  
788 Silva, Josep Ll Berral, and David Carrera. Towards  
789 automatic model specialization for edge video ana-  
790 lytics. *Future Generation Computer Systems*, 2022.  
791 2, 4
- 792 [27] Dan Roth and Kevin Small. Margin-based active  
793 learning for structured output spaces. In *European*  
794 *Conference on Machine Learning (ECML)*, 2006. 2
- 795 [28] Soumya Roy, Asim Unmesh, and Vinay P Nambood-  
796 iri. Deep active learning for object detection. *29th*  
797 *British Machine Vision Conference(BMVC)*, 2018. 3,  
798 4
- 799 [29] Sebastian Schmidt, Qing Rao, Julian Tatsch, and  
800 Alois Knoll. Advanced active learning strategies for  
801 object detection. In *2020 IEEE Intelligent Vehicles*  
802 *Symposium*, 2020. 3
- 803 [30] Ozan Sener and Silvio Savarese. Active learning for  
804 convolutional neural networks: A core-set approach.  
805 In *ICLR 2018*, 8 2017. 2
- 806 [31] Burr Settles. Active learning literature survey. *Com-  
807 puter Sciences Technical article, University of Wis-  
808 consin–Madison*, 2009. 2
- 809 [32] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell.  
810 Variational adversarial active learning. In *Proceed-  
811 ings of the IEEE/CVF International Conference on*  
812 *Computer Vision*, pages 5972–5981, 2019. 3
- 813 [33] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and  
814 Nicu Sebe. Curriculum self-paced learning for cross-  
815 domain object detection. *Computer Vision and Image*  
816 *Understanding*, 204:103166, 2021. 1
- 817 [34] Jiaxi Wu, Jiaxin Chen, and Di Huang. Entropy-based  
818 active learning for object detection with progressive  
819 diversity constraint. In *Proceedings of the IEEE/CVF*  
820 *International Conference on Computer Vision*, 2022.  
821 2
- 822 [35] Chenhongyi Yang, Lichao Huang, and Elliot J. Crow-  
823 ley. Plug and play active learning for object detection.  
824 <http://arxiv.org/abs/2211.11612>, 2022. 3
- 825 [36] Donggeun Yoo and In So Kweon. Learning loss for  
826 active learning. In *IEEE/CVF Conference on Com-  
827 puter Vision and Pattern Recognition*, 5 2019. 2
- 828 [37] Weiping Yu, Sijie Zhu, Taojinnan Yang, and Chen  
829 Chen. Consistency-based active learning for object  
830 detection. In *IEEE Conference on Computer Vision*  
831 and *Pattern Recognition (CVPR) Workshops*, 2021. 3
- 832 [38] Tianning Yuan, Fang Wan, Mengying Fu, Jianzhuang  
833 Liu, Songcen Xu, Xiangyang Ji, and Qixiang Ye.  
834 Multiple instance active learning for object detection.  
835 In *IEEE Conference on Computer Vision and Pattern*  
836 *Recognition (CVPR)*, 4 2021. 3