



Reducing redundancy in Semantic-KITTI using data augmentation within AL

Ngoc Phuong Anh Duong, Alexandre Almin, Léo Lemarié, Ravi Kiran

► To cite this version:

Ngoc Phuong Anh Duong, Alexandre Almin, Léo Lemarié, Ravi Kiran. Reducing redundancy in Semantic-KITTI using data augmentation within AL. NeurIPS 2021 Bayesian Deep Learning Workshop, Dec 2021, Paris, France. <hal-03494276>

HAL Id: hal-03494276

<https://hal.science/hal-03494276v1>

Submitted on 18 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

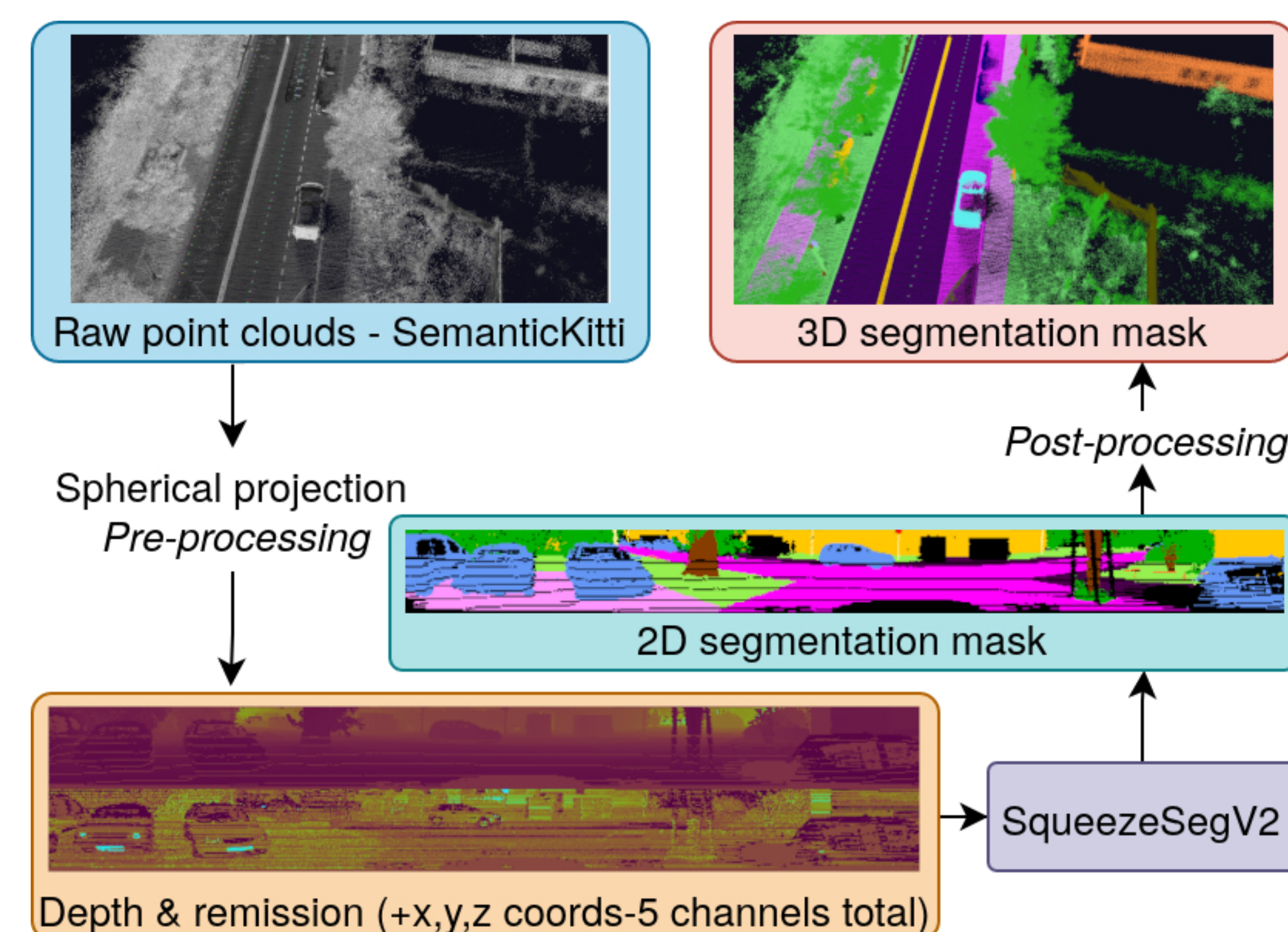


HAL Authorization

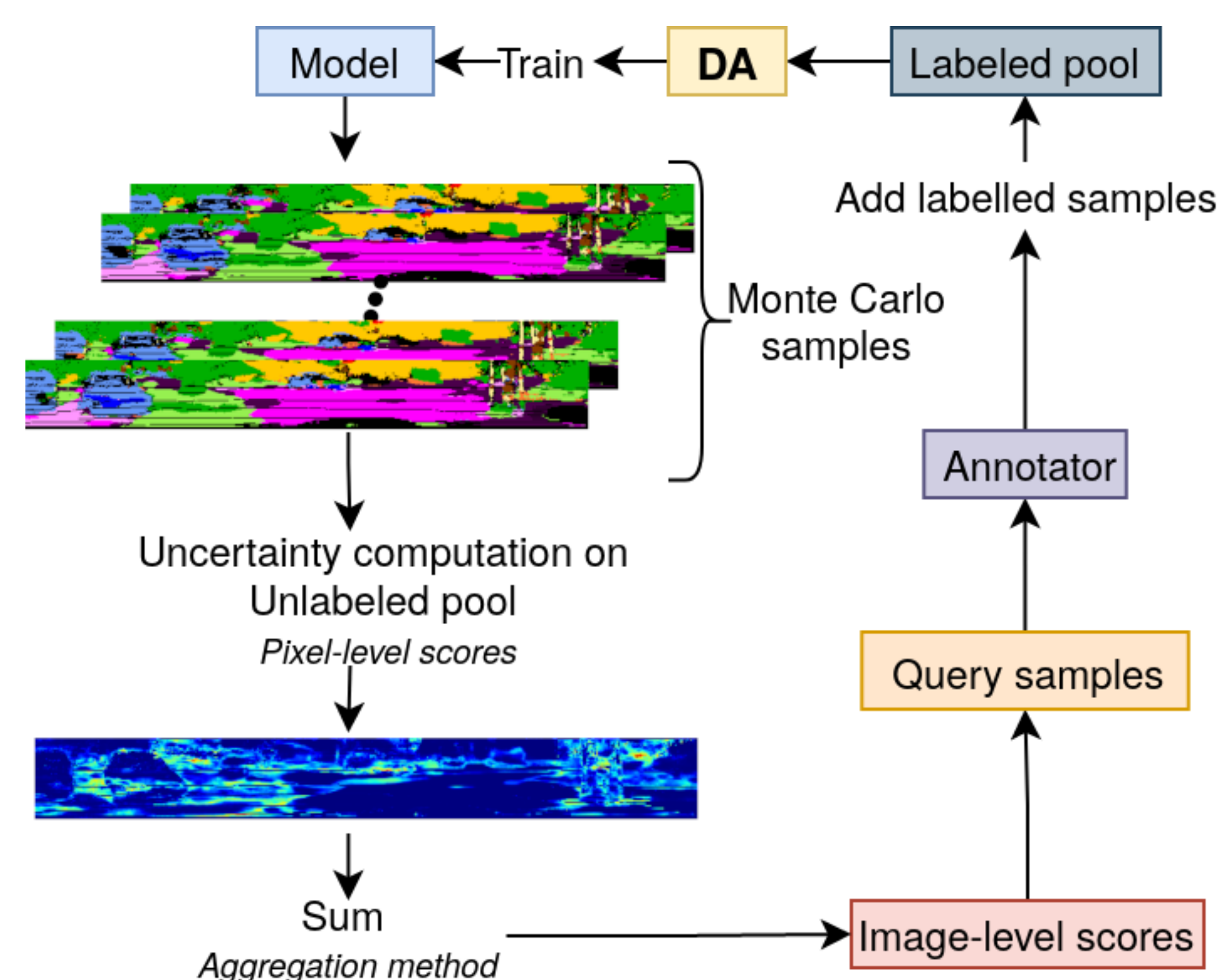
Introduction

- **Problem:** Large-scale datasets contain high redundancy.
- **Motivation:** Data augmentation (DA) models the redundant generative process that gives similar samples in autonomous driving datasets.
- **Our focus:** Evaluate the effect of DA on uncertainty sampling and labeling efficiency in active learning (AL) for large-scale point cloud semantic segmentation datasets.

Model & Dataset



Active Learning Loop



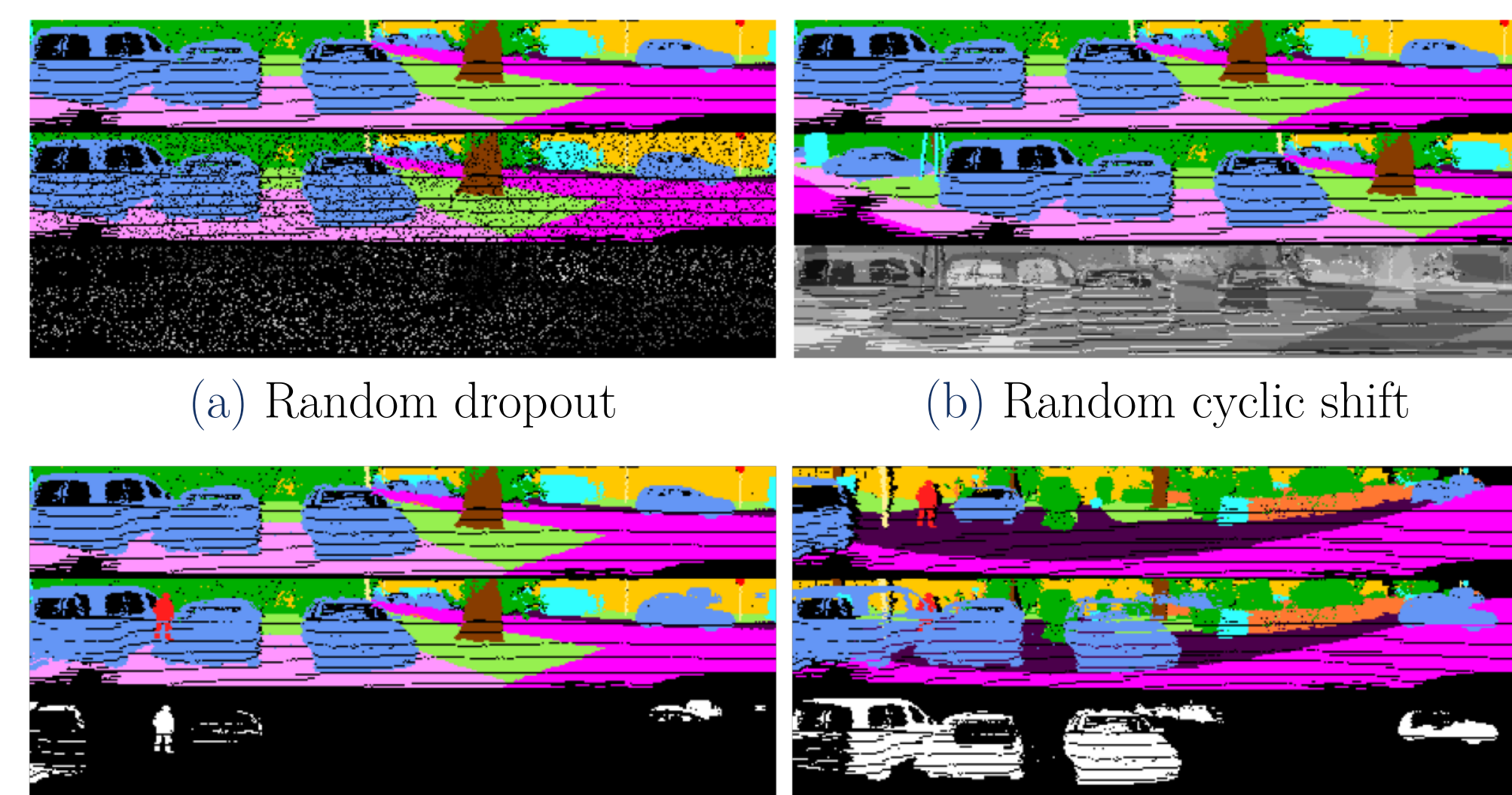
Pool	Test	Init	Budget	AL steps
6000	2000	240	240	25

Experiments settings

Heuristic functions:

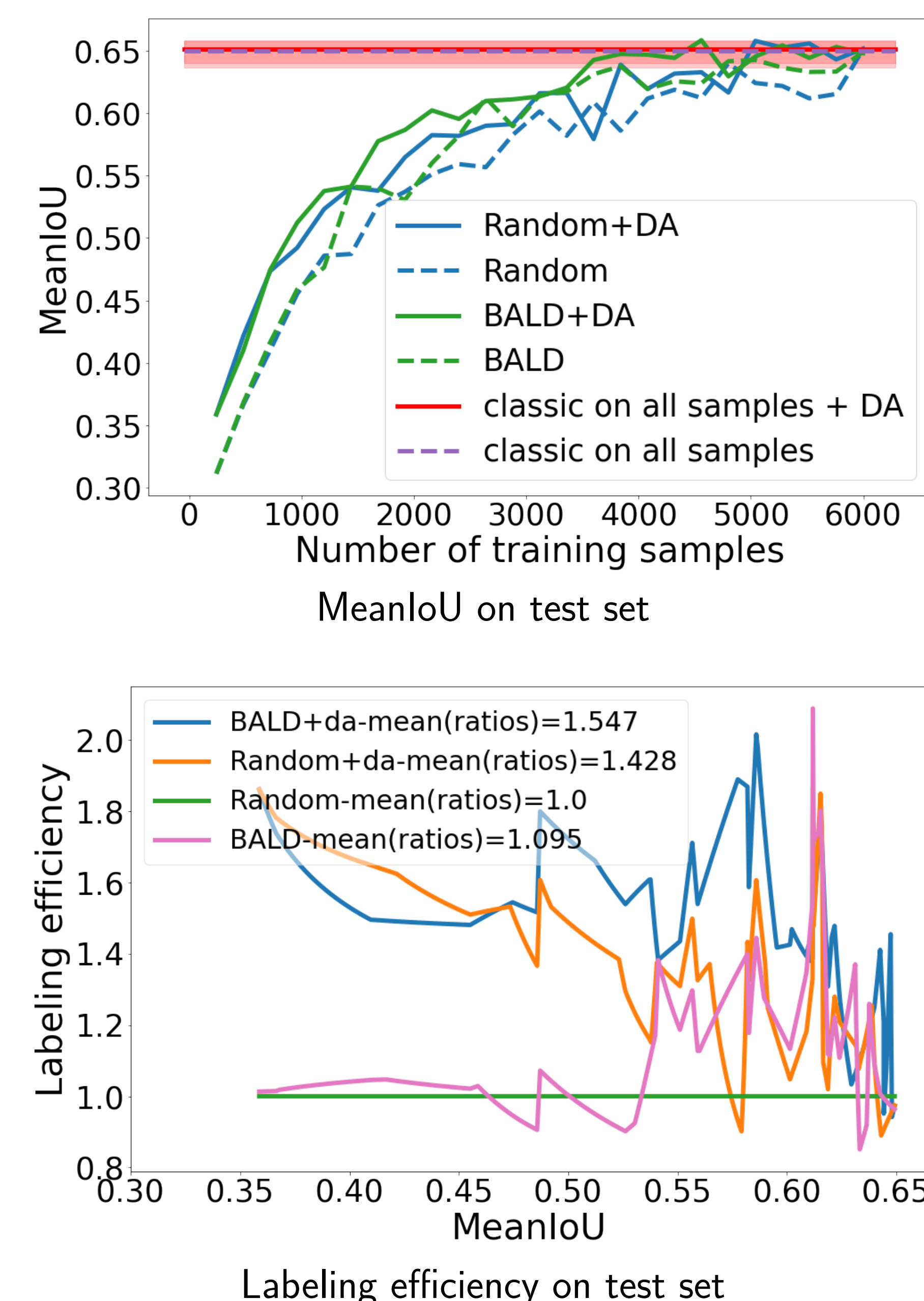
Random, BALD ($avg[entropy(preds)] - entropy[avg(preds)]$)

Examples of DA on range images



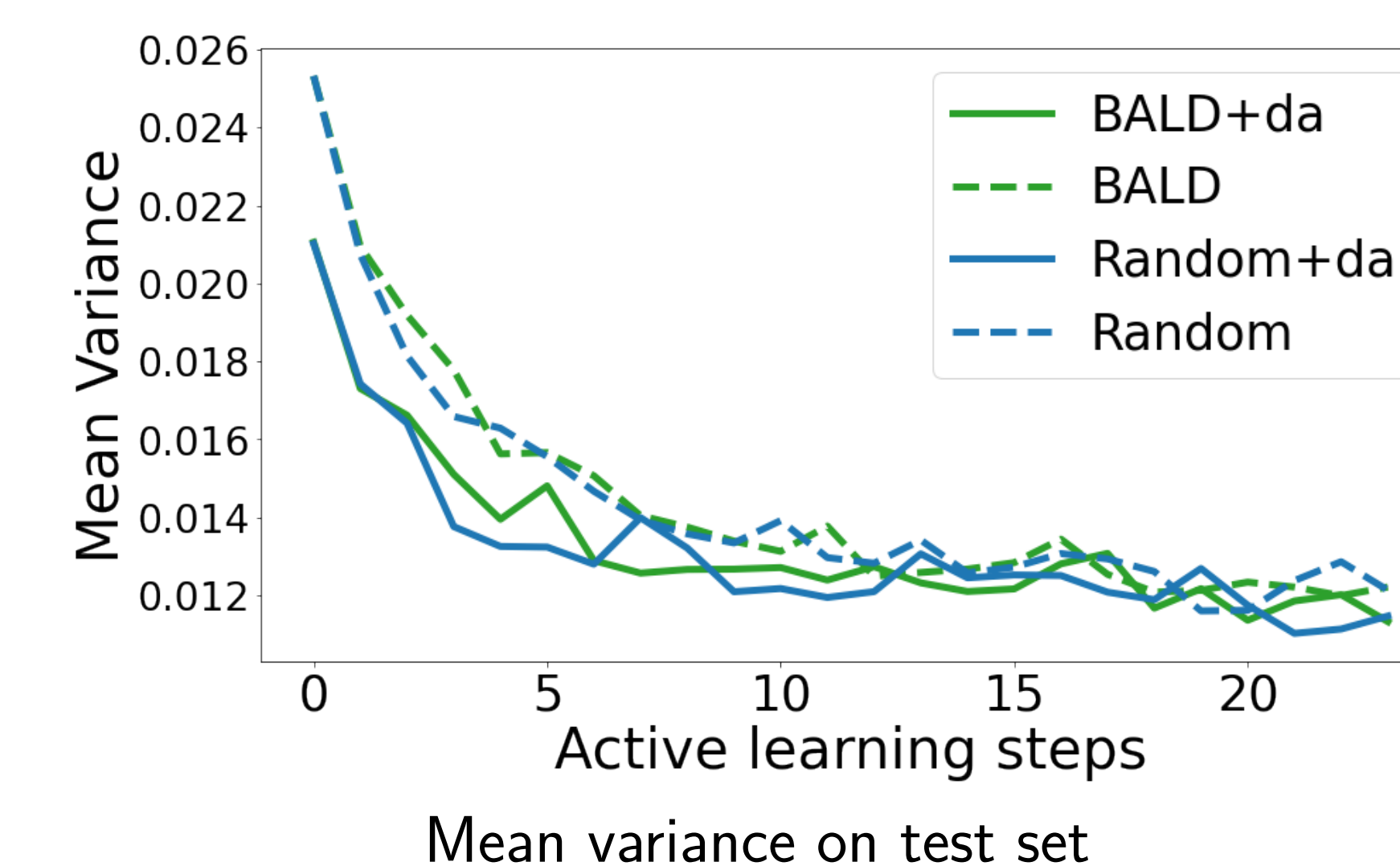
Top: Original | Middle: Augmented | Bottom: Difference

DA evaluation

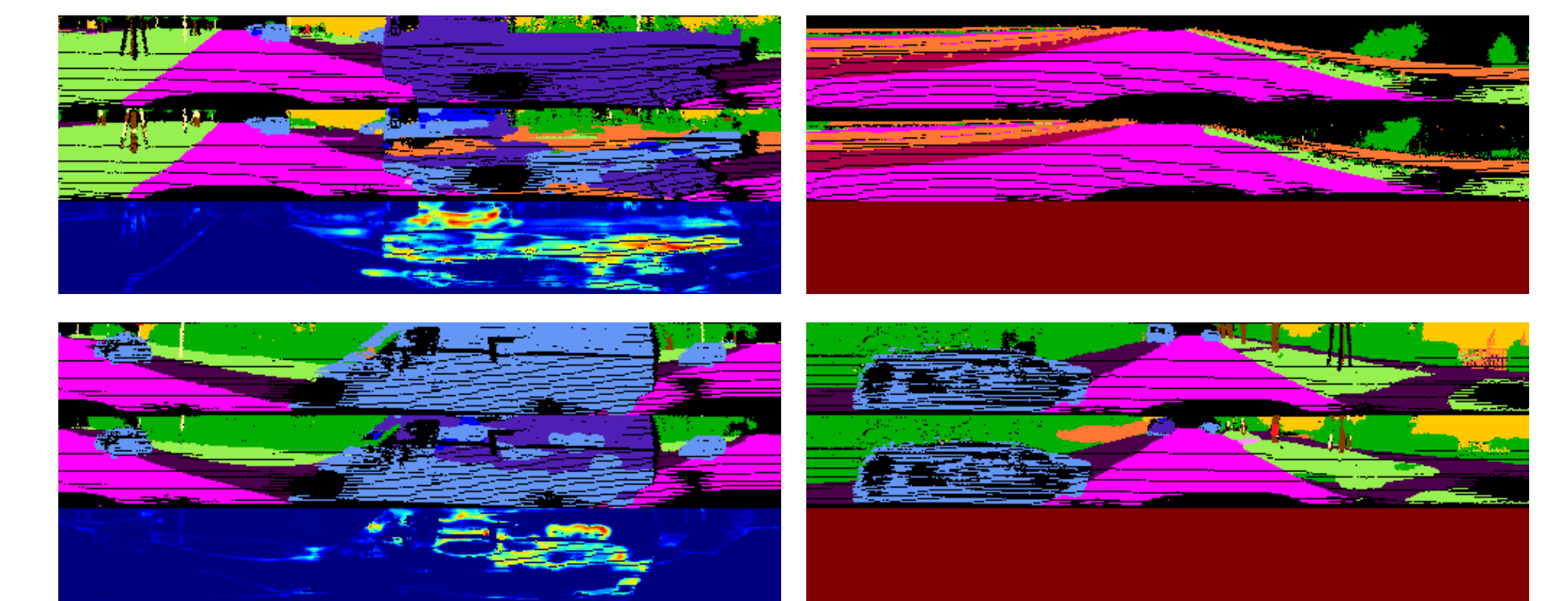


Results

- BALD with DA can achieve an important dataset compression, by using only 60% of the total sample pool.
- DA improves model generalization by regularization, making the sample selection more informative
- DA improves the stability of models.



Top hardest samples



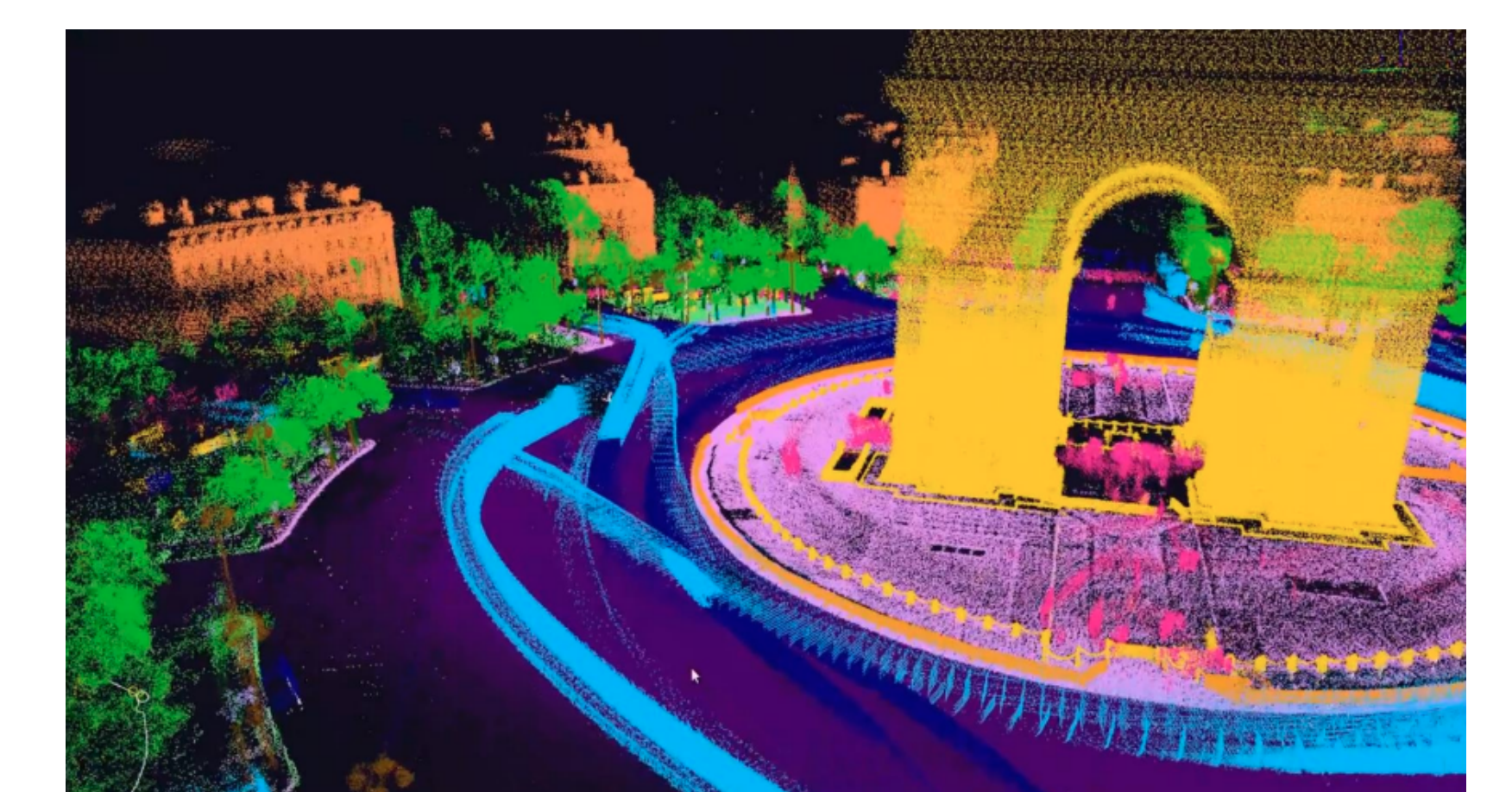
(a) BALD (b) Random
Top: Target | Middle: Prediction | Bottom: Heuristic score

Conclusion

- DA helped to select better informative samples, by improving the heuristic function stability.
- Improved labeling efficiency: With only 60% of the samples, we reached the accuracy of supervised training on full selected dataset.
- DA reduced annotation costs as well as reducing training time in production over large datasets.
- Similar gains have been observed using DA by [1] on CIFAR dataset for classification task.

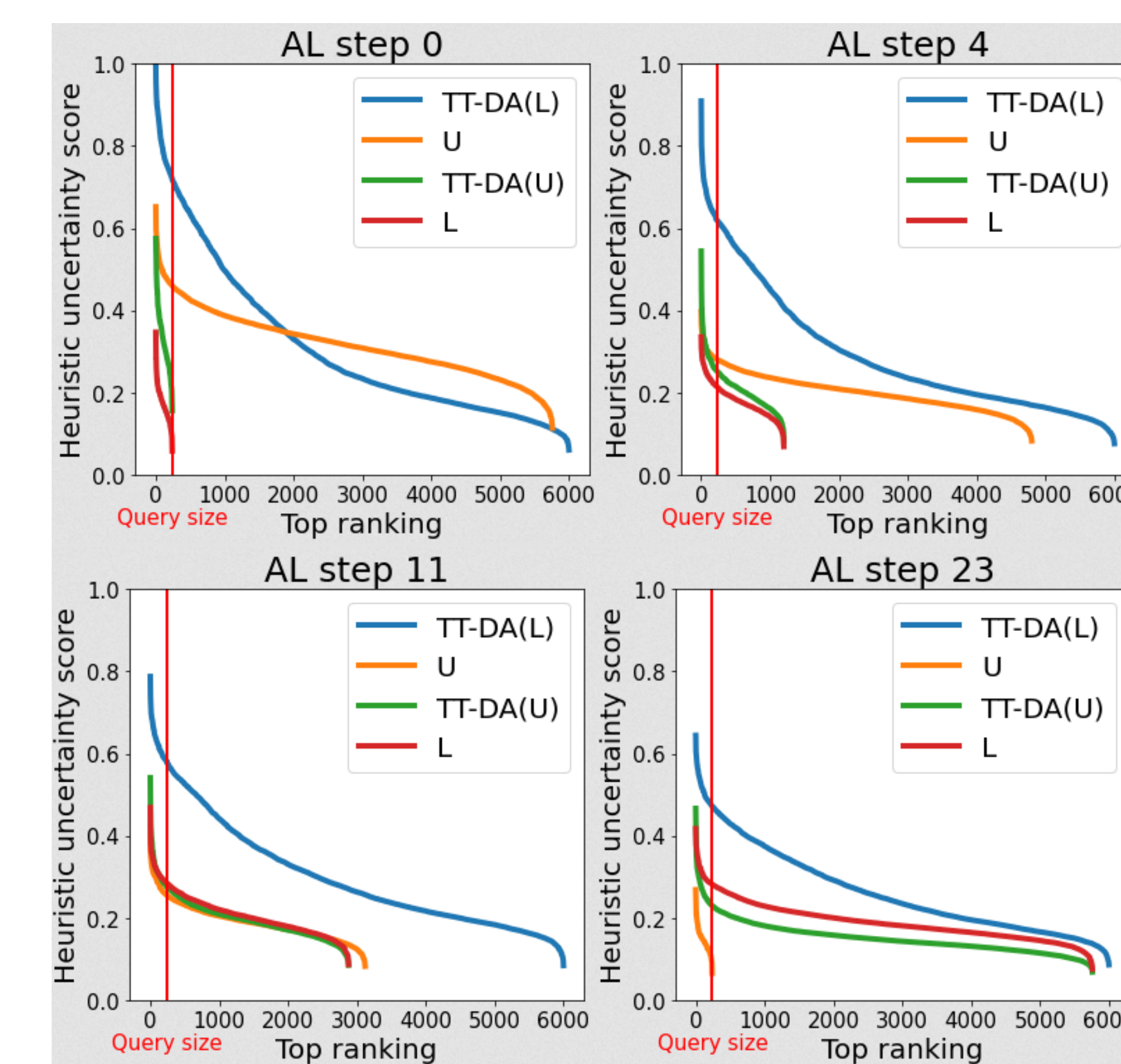
[1] Nathan Beck et al. *Effective Evaluation of Deep Active Learning on Image Classification Tasks*. 2021. arXiv: 2106.15324 [cs.CV].

Upcoming Navya dataset



Acknowledgements: This work was granted access to HPC resources of [TGCC/CINES/IDRIS] under the allocation 2021- [AD011012836] made by GENCI (Grand Equipement National de Calcul Intensif). It is part of the Deep Learning Segmentation (DLS) project financed by ADEME.

Aggregated heuristic evaluation



Aggregated heuristic score of samples sorted by decreasing value (using models w/o DA) over labeled pool(L), unlabeled pool(U), and Test-Time DA (TTDA) of L and U.

- **Goal:** Understand how the AL pipeline ranks and selects augmented vs normal samples w.r.t heuristic.
- **Result:** DA could generate high-score samples, but it could also provide transformed samples outside of the dataset distribution.