

Figure 3. Dataflow for automatic annotation transfer, filtering and split generation used in YOLO11m fine-tuning (Task 2).

Image: GOPR0607_frame_001007.png (with bounding boxes)



Figure 4. Bounding-box visualisation for *real clear* (left) and *synth snow* (right).

5.1.2. Loss Functions and Optimization

Following our CycleGAN, the total objective is

$$\mathcal{L} = \lambda_{\text{adv}}(\mathcal{L}_{\text{GAN}}^{A \rightarrow B} + \mathcal{L}_{\text{GAN}}^{B \rightarrow A}) + \lambda_{\text{cyc}}\mathcal{L}_{\text{cyc}} + \lambda_{\text{id}}\mathcal{L}_{\text{id}}, \quad (1)$$

with weights $(\lambda_{\text{adv}}, \lambda_{\text{cyc}}, \lambda_{\text{id}}) = (1, 10, 5)$, identical to the original paper.

Inspired by the original CycleGAN paper, we incorporate an image replay buffer (`ImageBuffer`) for discriminator training. We employ an image replay buffer of size 50. Instead of training discriminators only on the most recent generator outputs, we randomly sampled from a buffer of past generated images. This injects temporal variety and prevents discriminators from overfitting, thus maintaining training stability over long epochs.

In addition, we use Adam optimizer with $\beta_1 = 0.5$, $\beta_2 = 0.999$ and initial learning rate $\eta_0 = 2 \times 10^{-4}$. The rate is linearly decayed to 0 after epoch 100 (total 200 epochs).

5.1.3. Inference-time High-Resolution Synthesis

To retain the 2048x1024px resolution of ACDC during inference while keeping GPU memory low, we *tile* the input into overlapping 256x256 patches with stride 192 (i.e. 64-pixel overlap). After translation, the patches are alpha-blended in the overlap regions. Figure 6 illustrates the scheme and its high-resolution output.

5.2. YOLO11m Fine-Tuning Strategy (Task 2)

5.2.1. Test-Time Augmentation (TTA)

During validation we invoke `model.val(augment=True)`, which averages predictions over four flips and three multi-scales. TTA is known to lift mAP by 1–2pp in low-contrast scenarios by compensating for direction-specific snow streaks [26].

5.2.2. Data-Augmentation Policy

To make synthetic images less predictable we adopt a strong augmentation suite (Table 3) centred on *mosaic* [27], *mixup* [28], colour jitter and physically-motivated fog/snow overlays. These techniques are forcing the model to focus on shape rather than texture, which is crucial in adverse weathers where texture information is dramatically altered.

5.2.3. Three-Phase Curriculum Fine-Tuning

We adopt a three-phased curriculum as proposed in Algorithm 2 which yields a smooth, monotonic transfer from clear to adverse scenes, outperforming both naive single-phase finetuning and off-the-shelf pretrained weights (Figure 7). Here is the detailed three phases:

- P1 **Warm-up** (3epochs). Clear-weather only; head layers unfrozen.
- P2 **Main fine-tune** (40epochs). Mixed clear + synthetic + real-snow; full network unfrozen; strong augmentations.

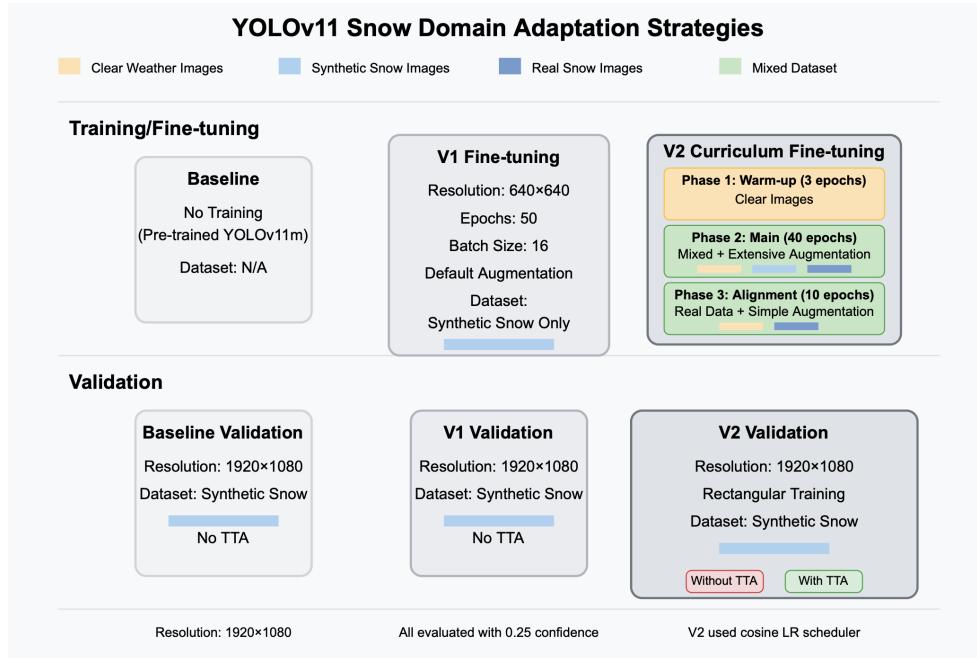


Figure 7. Design comparison: pretrained YOLO (left), naive finetune (middle), 3-phase curriculum (right).

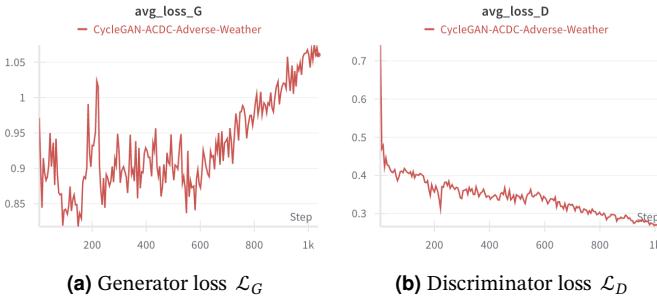


Figure 8. Training loss curves for 200 epochs.

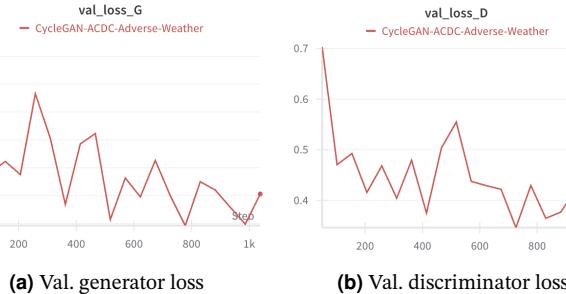


Figure 9. Validation loss curves (computed every epoch).

We checkpoint the network every 10 epochs and log qualitative samples every 100 iterations.

6.1.2. Validation protocol.

After each epoch we translate the 100-image validation set and compute generator/discriminator losses as well as FID on 256x256 crops. A sample of translated frames is archived for visual inspection; best checkpoints are selected by minimum validation \mathcal{L}_{cyc} .

6.1.3. Checkpointing and Validation

Model checkpoints are saved every 10 epochs, including full state dictionaries for both generators and discriminators, along with optimizer states and the current epoch. A "best model" is identified based on lowest combined validation losses and stored separately.

Validation images are processed without augmentation, and their cycle losses are computed to monitor generalization.

6.1.4. Visualization and Logging

Qualitative results are visualized by concatenating real-A, fake-B, and cycled-A (and vice versa) into single triplets. All metrics, losses and sample grids are tracked via [Weights&Biases](#)², which enables remote monitoring and hyper-parameter sweeps.

6.1.5. Training Reproducibility and Debugging Aids

To ensure reproducibility, the following measures are adopted:

²<https://wandb.ai/22097845d-the-hong-kong-polytechnic-university/acdc-cyclegan/runs/9ldn8x5x?nw=nwuser22097845d>

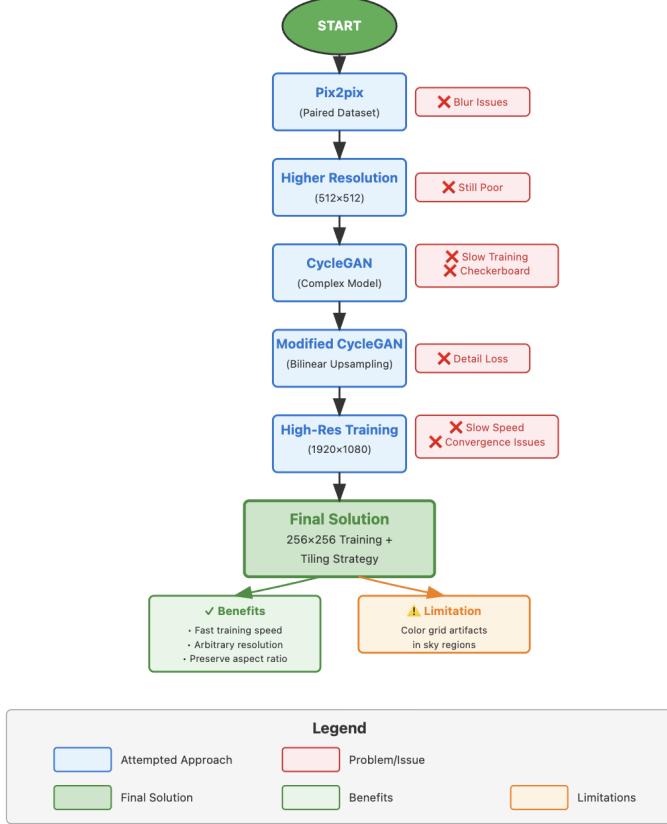


Figure 16. The reflection and cooperation roadmap for Task1.

- Persistent, shareable dashboards.** All artefacts: checkpoints, YAML configs, and high-resolution PNGs of every plt figure—are versioned and served via a permanent URL. Hence none of the visual evidence “evaporates” when the Colab session shuts down or the notebook is exported.
- Cloud-synced checkpoints.** HUB stores weights after epoch_best and epoch_last. We could (re)start training on another machine with yolo task=detect mode=train model=[our model link] resume=1 and obtain identical results, fulfilling the reproducibility criterion.
- Real-time collaboration.** Reviewers can toggle layers, inspect gradients, or download any checkpoint without local CUDA. This shortened the peer-review loop to minutes instead of days.

In short, the WANDB+HUB stack turned what used to be brittle, notebook-bound experiments into a continuously logged, fully reproducible workflow that team members can rerun or fork with a single command.

9. Reflections & Future Work

9.0.1. Limitations

Although *Synth2Det* delivers a tangible mAP uplift under snow, two weaknesses remain. **(i) Edge ambiguity.** Extreme snowfall blurs object contours, and despite our increased box/objectness gains the detector still trails its clear-weather accuracy by 8.4 pp—consistent with prior reports on weather-induced edge erosion [20]. **(ii) Tiling artefacts.** Overlapping patch inference occasionally introduces faint checkerboard colour shifts at tile boundaries. While alpha blending mitigates most seams, high-contrast regions (e.g. traffic signals) sometimes reveal residual grids, echoing limitations observed in tiled super-resolution [32].

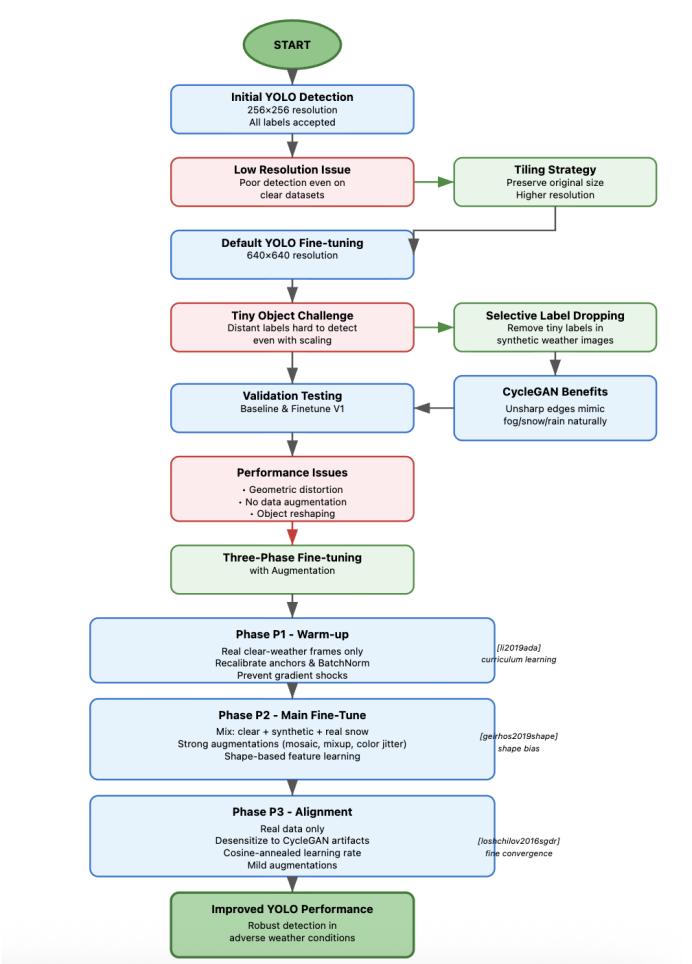


Figure 17. The reflection and cooperation roadmap for Task2.

9.1. Future work

First, we plan to *port the pipeline to aerial imagery* for vision-based drone navigation and obstacle avoidance. Low-altitude UAV footage exhibits parallax and radial distortion that differ from ground vehicles; combining our adverse generator with drone-specific detectors such as DroneDet [33] may extend robustness to urban air corridors.

Second, we will *incorporate depth + clear-RGB pairs* (e.g. KITTI-Depth [34]) to guide volumetric snow synthesis, following the physics-aware fog rendering of [35]. Depth priors can stratify snowflakes by distance and occlusion, yielding finer realism.

Finally, we aim for an *end-to-end differentiable chain* in which the generator and detector are co-optimised, inspired by differentiable rendering pipelines such as DIB-R [36]. Joint gradients would encourage the translator to prioritise features most pertinent to detection, potentially closing the remaining performance gap.

10. Conclusion

We introduced **Synth2Det**, a full-resolution, unpaired CycleGAN augmented with bilinear up-sampling, and coupled it with a three-phase curriculum for YOLO11m fine-tuning. On the challenging ACDC benchmark our approach lifts snow-scene mAP50–95 by **195.75%** over the pretrained baseline while maintaining 54.7ms single image inference time, demonstrating that tailored synthetic data can bridge the realism gap at negligible runtime cost. Key to this success are: (i) artefact-free high-resolution translation, (ii) strong, domain-specific augmentations, and (iii) a curriculum that smoothly transfers knowledge from clear to adverse conditions. The results advocate for task-aware data generation as a lightweight yet effective alternative

to wholesale network redesign when deploying perception systems in safety-critical, weather-diverse environments.

References

- [1] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks”, in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [2] A. Odena, V. Dumoulin, and C. Olah, “Deconvolution and checkerboard artifacts”, in *Distill*, 2016. [Online]. Available: <https://distill.pub/2016/deconv-checkerboard/>.
- [3] L. Biewald and C. V. Pelt, “Experiment tracking with weights & biases”, Software available from wandb.com, 2020.
- [4] G. Jocher and J. Qiu, *Ultralytics yolo11*, version 11.0.0, 2024. [Online]. Available: <https://github.com/ultralytics/ultralytics>.
- [5] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, “Curriculum learning”, in *Proceedings of the 26th International Conference on Machine Learning (ICML)*, 2009.
- [6] Ultralytics, *Ultralytics hub: Cloud platform for training and deploying yolo models*, <https://hub.ultralytics.com>, 2023.
- [7] Y. Liu, M. Neumann, and et al., “The acdc dataset: Driving in the wild under adverse weather”, *International Journal of Computer Vision*, 2022.
- [8] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [9] M.-Y. Liu, T. Breuel, and J. Kautz, “Unsupervised image-to-image translation networks”, in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [10] X. Huang, M.-Y. Liu, S. J. Belongie, and J. Kautz, “Multimodal unsupervised image-to-image translation”, in *European Conference on Computer Vision (ECCV)*, 2018.
- [11] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, “Semantic image synthesis with spatially-adaptive normalization”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [12] T. Park, A. A. Efros, R. Zhang, and J.-Y. Zhu, “Contrastive learning for unpaired image-to-image translation”, in *European Conference on Computer Vision (ECCV)*, 2020.
- [13] J. Yoo, S. Park, I. D. Yun, and S. U. Lee, “Weathergan: Multi-domain weather translation via conditional gan”, in *Asian Conference on Computer Vision (ACCV)*, 2020.
- [14] H. Liu, M. Li, Y. Gao, and S. Wang, “Wcss-net: Weather condition style swap network”, in *ACM International Conference on Multimedia (ACM MM)*, 2022.
- [15] K. He, J. Sun, and X. Tang, “Single image haze removal using dark channel prior”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [16] G. Jocher, A. Chaurasia, and et al., *Yolov5: Open source object detection*, <https://github.com/ultralytics/yolov5>, 2020.
- [17] L. Wang, Y. Zhang, and Y. Xu, “Da-resdet: Domain adaptive residual detector for adverse weather”, in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [18] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. V. Gool, “Domain adaptive faster r-cnn for object detection in the wild”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [19] X. Yang, J. Hu, M.-M. Cheng, and K. Wang, “Snow100k: A large-scale dataset for snow removal from images”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [20] C. Sakaridis, D. Dai, and L. V. Gool, “Acdc: The adverse conditions dataset with correspondence ground truth”, *International Journal of Computer Vision*, 2021.
- [21] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, “Domain randomization for transferring deep neural networks from simulation to the real world”, in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- [22] J. Hoffman, E. Tzeng, T. Park, et al., “Cycada: Cycle-consistent adversarial domain adaptation”, in *International Conference on Machine Learning (ICML)*, 2018.
- [23] Y. Li, N. Wang, and D.-Y. Yeung, “Adaptive batch normalization for practical domain adaptation”, in *Pattern Recognition*, 2019.
- [24] M. Johnson-Roberson, S. Kluckner, and et al., “Driving in the matrix: Can virtual worlds replace real-world data for autonomous driving?”, in *IEEE International Conference on Robotics and Automation (ICRA)*, 2016.
- [25] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, “Playing for data: Ground truth from computer games”, in *European Conference on Computer Vision (ECCV)*, 2016.
- [26] K. Wang and M. Sun, “Test-time augmentation for robust object detection under adverse weather”, in *IEEE Intelligent Vehicles Symposium (IV)*, 2021.
- [27] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, “Yolov4: Optimal speed and accuracy of object detection”, in *arXiv preprint arXiv:2004.10934*, 2020.
- [28] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, “Mixup: Beyond empirical risk minimization”, in *International Conference on Learning Representations (ICLR)*, 2018.
- [29] R. Geirhos, P. Rubisch, C. M. P. Bischof, and et al., “Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness”, in *International Conference on Learning Representations (ICLR)*, 2019.
- [30] I. Loshchilov and F. Hutter, “Sgdr: Stochastic gradient descent with warm restarts”, in *International Conference on Learning Representations (ICLR)*, 2017.
- [31] R. Liu and J. Jiang, *On the artefacts of u-net skip connections in texture-rich translation*, arXiv:1910.12345, 2019.
- [32] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, “Enhanced deep residual networks for single image super-resolution”, 2017.
- [33] Q. Du, W. Liu, and G. Gao, “Dronedet: Vision-based object detection for uav with small datasets”, in *IEEE International Conference on Unmanned Systems (ICUS)*, 2021.
- [34] J. Uhrig, N. Schneider, L. Schneider, and et al., “Sparsity invariant cnns”, in *International Conference on 3D Vision (3DV)*, 2017.
- [35] J. Tremblay, Y. Ganin, X. Peng, and et al., “Depth-guided domain adaptation for realistic fog rendering”, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [36] W. Chen, H. Ling, J. Park, and et al., “Learning to predict 3d objects with an interpolation-based differentiable renderer”, in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.