

Image2Image Translation & Image Detection Report

QIN Qipeng 21101226D, CHEN Ziyang 21095751D

Abstract—Automatic image-to-image translation is an essential task in computer vision, enabling semantic-to-realistic and realistic-to-semantic translations between spaces. Although supervised methods like Pix2Pix excel with paired data and unsupervised methods like CycleGAN offer flexibility on unpaired sets, each suffers from drawbacks. Pix2Pix's dependency on costly paired samples and CycleGAN's incapability to always recognize intricate transformations accurately. To alleviate these constraints, we introduce a hybrid framework that begins with Pix2Pix models in both directions of translation and further refines them through CycleGAN's cycle-consistency training. This integration leverages supervised precision and unsupervised robustness. As a result, our translation model reduces pixel-level reconstruction loss by over 70% (MSE goes from 3.3 to 1.0) and boosts PSNR by 20% (from 12.5 dB to 15 dB), outperforming single-model GAN techniques. For downstream assessment, we deploy a pretrained YOLOv5 detector onto generated images with fine-tuning on a small labeled subset and semi-supervised pseudo-label self-training across 2,000+ unlabeled samples. Detection performance goes through the roof: mAP@0.5 improves from 2.5% (baseline) to 70.9%, precision improves from 3% to nearly 100%, and recall improves from 11.7% to 94%, all with real-time inference kept constant at around 369 FPS.

Impact Statement — Addressing the challenges of domain shifts and data scarcity in image translation and object detection, our project introduces a hybrid GAN framework and an adapted YOLOv5 detector. By synergizing Pix2Pix and CycleGAN, we achieve superior image translation quality, leveraging both paired and unpaired data. The adapted YOLOv5 ensures robust object detection across diverse environments. This integrated approach significantly enhances generalization and precision, outperforming traditional methods. Our work has broad applications, from enabling autonomous vehicles to navigate varied conditions to improving surveillance systems' consistency. It also advances gaming and virtual reality through realistic image rendering. By demonstrating substantial improvements in translation quality and detection accuracy, our project paves the way for more adaptable and reliable AI systems in real-world scenarios.

Index Terms—Computer Vision, Deep Learning, Generative AI

I. INTRODUCTION

Automatic image-to-image translation is a key area of computer vision and computer graphics technology that enables pixel-level translation to coordinate various visual spaces. Similar to NLP, re-represent an input image in a new presentation while preserving its inner semantic organization.

We implement and extend the generative adversarial networks (GANs), using Pix2Pix and CycleGAN for this project in order to generate a domain-specialized translation

framework that supports semantic-to-real as well as real-to-semantic conversion.

To guarantee the validity of our semantic-to-real outputs, we incorporate object detection with a pre-trained YOLO model. Given distribution shifts between our generated images and YOLO's original training data, we employ two adaptation strategies: fine-tuning on a small labeled subset of generated real images and test-time augmentation. These enhancements strongly improve detection performance (mAP, precision, and recall), yet preserve YOLO's real-time inference.

II. LITERATURE REVIEW

This literature review explores findings related to candidates from the GAN sub-branch, specifically for a domain-specific model designed to meet the image-to-image translation requirements of this project. The review analyzes the characteristics of these candidates, discusses their strengths and weaknesses, and investigates potential integration approaches, while also addresses the challenges encountered during the process. In addition to GAN, this review also briefly touches upon domain adaptation in object detection using YOLO, considering the project's secondary focus on image detection..

A. Search for Suitable GAN Sub-Branches

To achieve image-to-image translation tasks, it requires a GAN sub-branch that is designated for domain-specific image-to-image translation. Two notable sub-branches of GAN are found that are both capable to learn mappings between image domains.

The first sub-branch is Pix2Pix, a form of conditional GAN (cGAN) that focuses on paired image-to-image translation, where direct input-output image pairs are provided [2][3]. Given that the dataset consists of pre-labeled features, Pix2Pix exemplifies a supervised approach in GANs for image-to-image translation.

In contrast, CycleGAN is designed for unpaired data [2][6]. Utilizing cycle consistency loss, CycleGAN translates an image from domain A to domain B and then back to domain A, striving to maintain a close resemblance to the original image. This characteristic makes CycleGAN particularly suitable for style transfer between domains that lack paired samples. Therefore, as datasets are often unpaired, CycleGAN leans more towards an unsupervised approach.

In terms of architecture, Pix2Pix adopts a U-Net for its generator, incorporating multiple scale fusion for improved cross-layer connections, whereas CycleGAN employs two generators that predominantly utilize ResNet architecture. However, both Pix2Pix and CycleGAN share the same discriminator architecture, known as PatchGAN, which opens possibilities for their integration.

B. Analysis of Applications, Strengths, and Weaknesses

Pix2Pix, introduced by Isola et al. [3], has been widely used in applications such as sketch-to-photo translation and semantic segmentation of real images. Its strength lies in its ability to produce high-quality, accurate outputs when trained on paired datasets, which makes it suitable for direct image transformation tasks such as face aging. However, this reliance on paired data is a key weakness, as such datasets are often costly or unavailable. In addition, the strong tendency for direct mapping between two pairs of data may also lead to a lower degree of generalization of the model to unseen data.

CycleGAN, proposed by Zhu et al. [6], overcomes this limitation by enforcing bi-directional mapping by processing unpaired data through cycle consistency loss. It shows advantages in image conversion between two domains and has been successfully applied to style conversion, such as from realistic images to Van Gogh paintings. However, its lack of direct supervision may lead to less accurate translations, especially in complex transformations. Furthermore, since CycleGAN does not encode the image prior to conversion, there may be little difference between the input and output [2].

Admittedly, these tradeoffs highlight the differences between the two models, but they also imply the complementary nature of the two approaches and the potential benefits of their combination. Pix2Pix's precision in supervised settings could enhance CycleGAN's flexibility with unpaired data, while CycleGAN's unsupervised learning could mitigate Pix2Pix's dependency on paired datasets. Several studies support this hybrid approach; for instance, a study of integrating paired and unpaired data in medical imaging [1] showed that using paired data for initialization and then unpaired data for refinement improves the quality of translations across different datasets. Another study of a multi-data GAN framework [5] showed that the simultaneous use of both data types outperforms single-data approaches, supporting the idea that the supervised accuracy of Pix2Pix can complement the unsupervised adaptation of CycleGAN. A feasible hybrid model, which uses a shared generator such as U-Net and combines the L1 loss of Pix2Pix with the cycle-consistency loss of CycleGAN, can process both paired and unpaired scenes, providing a versatile solution for image-to-image conversion.

C. Exploration for implementation

To gain a practical comprehension of the two GAN frameworks, we examined the official versions of Pix2Pix and CycleGAN provided by their respective authors. While both frameworks are robust, they face significant challenges in their implementation. Notably, CycleGAN's training process is extremely time-consuming, often occupying up to more than 67000 seconds, that is 18.61 hours, on Colab due to its computational complexity and the need for large datasets. This bottleneck results in a substantial time investment and hinders the assessment of the performance of both Pix2Pix and CycleGAN. The situation underscores the need for a more efficient method to test the integration between the two frameworks.

To address computational challenges, it's essential to simplify the Pix2Pix and CycleGAN code into a more generic implementation that maintains their core architectures, such as the U-Net generator and PatchGAN discriminator from Pix2Pix,

along with CycleGAN's cycle consistency feature. A streamlined Jupyter Notebook codebase will enhance training efficiency and facilitate the exploration of conditional GANs integration. Additionally, generative AIs are being used to brainstorm methods for integrating Pix2Pix and CycleGAN, helping to overcome obstacles and develop hybrid models aligned with project goals.

D. YOLO: Domain Adaptation for Object Detection

Adaptation of object detectors to new domains (where data distributions are different) has been a very active research area, particularly with the advent of autonomous driving and surveillance, where models learned in one city or sensor need to generalize to another. A widely used approach in the earlier literature is unsupervised domain adaptation (UDA), which trains a detector with labeled source data and unlabeled target data via feature distribution alignment. For example, Chen et al. presented Domain-Adaptive Faster R-CNN, which utilizes adversarial training at the image level and instance level to make the features domain-invariant [7]. In this direction, Saito et al. proposed a Strong-Weak alignment method: aligning strong global features and weak local features separately to avoid negative transfer [8]. These approaches, however, were demonstrated on two-stage detectors (Faster R-CNN) and usually involved adding complex domain classifier sub-networks and GRL (Gradient Reversal Layer) for adversarial loss.

E. Few-Shot Fine-Tuning

When there are a few labeled target examples available (few-shot setting), fine-tuning the detector on this small set can greatly facilitate adaptation. Traditional wisdom would retrain the whole network on new data, but in FSOD research, it has been found that freezing most of the network and only modifying some layers can be more effective since there is a risk of overfitting to the few examples. Wang et al. famously showed a "Frustratingly Simple Few-Shot OD" approach wherein they fine-tuned only the last-layer classifiers and box regressors of a pre-trained detector on new classes [9]. This simple method outperformed more complicated meta-learning approaches by a large margin. The intuition being that the backbone has already learned general features (edges, shapes, etc.) helpful for any class or domain, and only the final layers need to be adapted to accommodate the new object appearances or classes.

III. METHODOLOGY AND IMPLEMENTATION

A. Image-to-Image Translation

1) Overview

In our implementation of a domain-specific image-to-image translation system, we propose a hybrid image-to-image translation framework that constructs a CycleGAN with two generators and discriminators derived from two pre-trained Pix2Pix models. One model is responsible for converting images from domain A to domain B, while the other handles the reverse, mapping images from domain B back to domain A. Overall, the training framework resembles a standard CycleGAN, but it utilizes Pix2Pix models as its core components.

The outcome of this approach is two generators: one that transforms realistic images into their corresponding semantic object style counterparts, and the other that performs the reverse transformation. For code actualization, we consolidated the entire repository of Python files into a single Jupyter file, simplifying the codebase as addressed in our last section.

2) Dataset and Preprocessing

Our datasets are based on the default dataset provided for this project. We define domain A as realistic street-view images in RGB format and domain B as semantic object segmentation images. The images from these two domains are organized into separate folders within the project repository. In total, there are 2,091 realistic-semantic pairs available to meet the project's needs for all purposes.

To accommodate the different data formats required by Pix2Pix and CycleGAN, we created two dataset classes: "PairedImageDataset" and "UnpairedImageDataset." The data loader for the former returns a dictionary object that includes images from domain A, their corresponding images in domain B, as well as their shared filenames. Conversely, the data loader for the latter returns only the images along with their filenames from the specified directory related to the domain. In other words, one data loader from the PairedImageDataset class can load the images from domain A and domain B simultaneously, while it requires two loaders to load the images from their respective domain. Fig. 1 and Fig. 2, respectively, show samples from the two dataset classes.

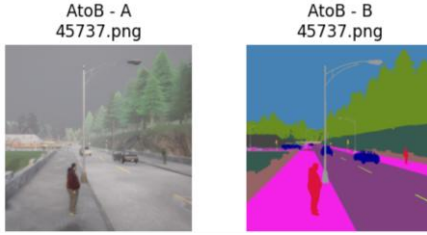


Fig 1

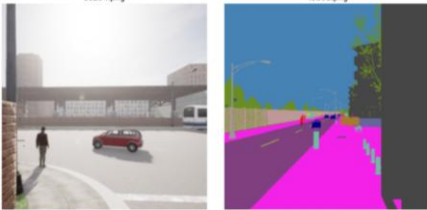


Fig 2

Despite their differences, both dataset classes share the same functional implementation, which includes segmenting the image data for training, validation, and testing using the default ratio of 8:1:1. Additionally, they both follow an identical data preprocessing pipeline for the training data, which includes the following steps:

1. Resize the images to 286×286 pixels.
2. Randomly crop the images to 256×256 pixels.
3. Apply a random horizontal flip.
4. Transform the images into tensors.
5. Normalize the tensors to a range of $[-1, 1]$.

This data augmentation technique enhances the robustness of the Pix2Pix model's performance on unseen data during the testing phase.

3) Model Architecture

The architecture of our two Pix2Pix models, $A \rightarrow B$ and $B \rightarrow A$, is organized into several modular components, as shown in Fig. 3. The illustration highlights the hierarchical order of the component classes, with higher-level components built upon the foundational lower-level building blocks.

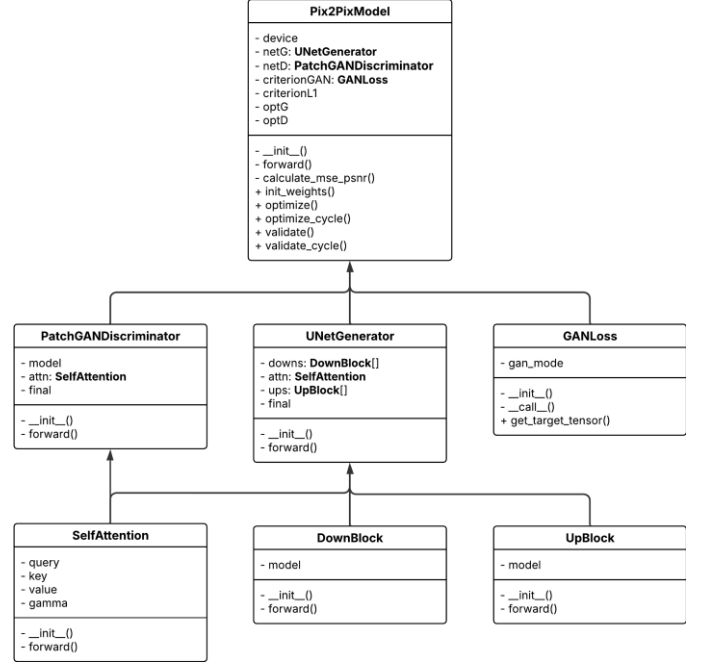


Fig 3

In our implementation, $A \rightarrow B$ and $B \rightarrow A$ instantiate the wrapper class "Pix2PixModel," which integrates the generator and discriminator into a single class. Although the generator and discriminator are not separated, we will refer to their respective generator as $G_{A \rightarrow B}$ and $G_{B \rightarrow A}$, and their corresponding discriminators as $D_{A \rightarrow B}$ and $D_{B \rightarrow A}$ for simplicity.

With regards to the generators, both $G_{A \rightarrow B}$ and $G_{B \rightarrow A}$ are the instances of the same class UNetGenerator which implements a typical U-Net-style encoder-decoder architecture, but with an extra attention layer between the encoder (*downs*) and the decoder (*ups*):

- Encoder (*downs*): 8 down-sampling blocks, each consisting of Conv-BatchNorm-ReLU layers.
- Self-Attention (*attn*): optional attention module to help the model focus on important objects
- Decoder (*ups*): 8 up-sampling blocks, each of which includes TransposedConv-BatchNorm-ReLU layers, with skip-connections.
- Final Layer (*final*): Upsample-Conv2d-Tanh activation.

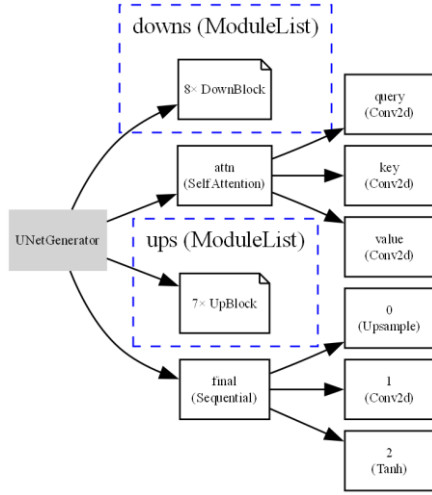


Fig 4. Visual Structure of Generator

The discriminators, denoted as $D_{A \rightarrow B}$ and $D_{B \rightarrow A}$, are instantiated as the PatchGAN objects that operate on 70×70 patches with an optional attention layer. They are defined as follows:

- **Backbone (model):** comprises four convolutional blocks, each followed by LeakyReLU and additional Batch Normalization (BN) between the second and fourth convolutional blocks.
- **Self-Attention (attn):** SelfAttention(512) if used, otherwise, it defaults to the identity function.
- **Final Conv (final):** A Conv2d layer produces output logits at a resolution of 30×30 .

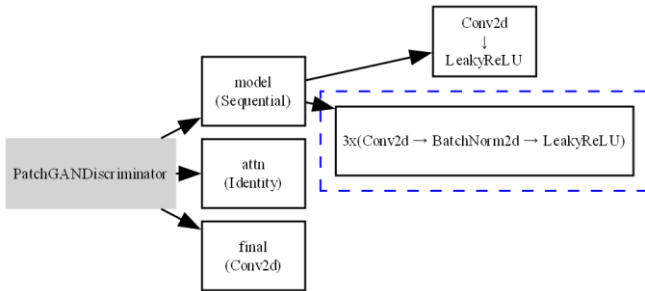


Fig 5. Visual Structure of Discriminator

4) Training and Evaluation Metrics

We split the overall training into two stages: Separated Pix2Pix Warm-up and Cyclic Fine-tuning.

In the first stage, we train two independent Pix2Pix models in isolation:

- **A \rightarrow B model:** learns to map semantic maps (A) to photographs (B)
- **B \rightarrow A model:** learns the inverse mapping

Both models share the same network architecture (U-Net generator and PatchGAN discriminator) but are trained on paired data in their respective directions.

There are two key methods in Pix2PixModel class with respect to stage 1:

- `optimize(real_input, real_target) \rightarrow loss`
- `validate(val_loader) \rightarrow (avg_loss, avg_psnr, samples)`

Stage 1 training consists of two separate loops: one for model A approaching model B, and another for model B approaching model A. During each epoch, the `optimize()` function is called for every batch, while the `validate()` function is executed at the end of the epoch.

The `optimize()` function performs a single gradient update step. It begins with a forward pass through both the generator and discriminator, where it computes the adversarial and L1 (reconstruction) losses. After that, it backpropagates the errors to update both networks and returns the total generator loss for logging purposes.

On the other hand, the `validate()` function runs the model in evaluation mode using the validation set. In this mode, gradients are disabled while the generator loss and Peak Signal-to-Noise Ratio (PSNR) are calculated across all batches. Additionally, it randomly selects one batch of input, target, and prediction triplets for visualization and reports average evaluation metrics along with the samples that are displayed. Throughout this entire process, visual outputs and performance metrics are logged and saved periodically.

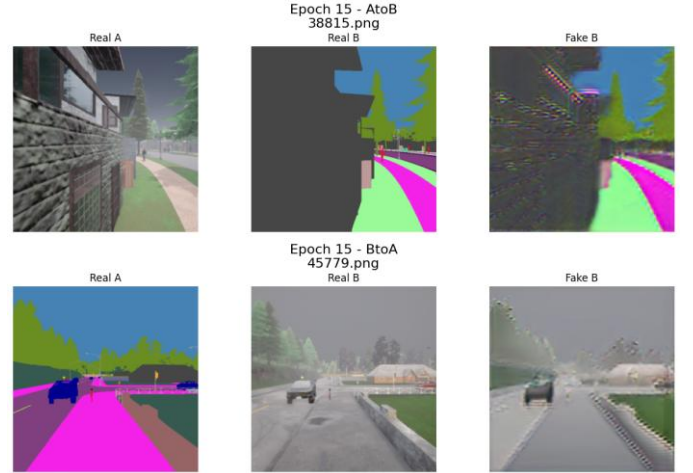


Fig 6 Samples during stage 1 training

After the warm-up, we switch to unpaired cycle-consistency training, simultaneously optimizing both directions to further improve realism and enforce invertibility.

To implement the CycleGAN characteristics, our Pix2Pix model class is extended with another two methods to wrap our two generators in a single cyclic training pass:

- `optimize_cycle(real_A, real_B, G_A2B, G_B2A) \rightarrow (loss_AB, loss_BA)`
- `validate_cycle(val_loader, G_inverse, lambda_cyc, lambda_id) \rightarrow (avg_loss_AB, avg_loss_BA, samples_AB, samples_BA)`

In the stage 2 training codebase, we have implemented a joint cycle training loop that alternates updates between two directions. In each epoch, we process both `real_A` and `real_B` inputs through the `optimize_cycle()` function, followed by the `validate_cycle()` function which logs metrics and creates side-by-side visualizations. First of all, `optimize_cycle()` combines both directions under a single training objective. In one iteration, we perform a complete CycleGAN-style update: we generate `fake_B` from `real_A` using the generator $G_{A \rightarrow B}$, and `fake_A` from `real_B` using $G_{B \rightarrow A}$. Then, we reconstruct `cycle_A` by passing

fake_B through $G_{B \rightarrow A}$, and cycle_B by passing fake_A through $G_{A \rightarrow B}G_{A \rightarrow B}$. During this process, we compute three types of losses: adversarial loss, cycle-consistency loss, and identity loss in both directions. We backpropagate and update the four networks involved: $G_{A \rightarrow B}$, $G_{B \rightarrow A}$, $D_{A \rightarrow B}$, and $D_{B \rightarrow A}$. The directional generator losses are returned for monitoring purposes.

In the validation phase, to evaluate the cyclic process, we run both the forward and reverse cycles, accumulating cycle losses in both the $A \rightarrow B$ and $B \rightarrow A$ directions. Similar to the `validate()` counterpart, one random example is also extracted from each direction for visual comparison. Finally, we return the average cycle losses along with sample tuples that consist of (real, fake, cycle, path) for a clear overview.

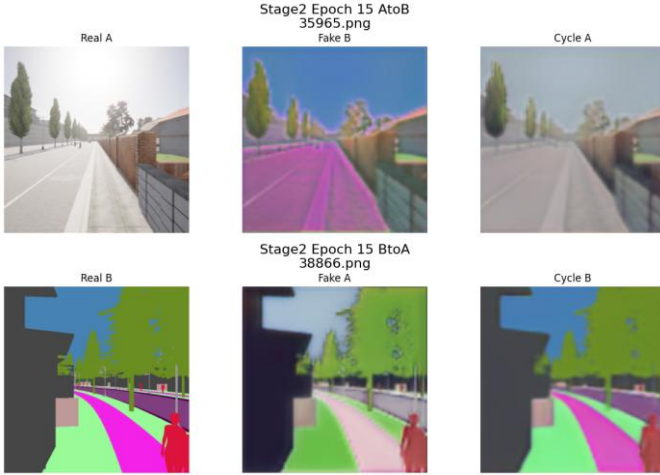


Fig. 7 Samples during stage 2 training

By encapsulating most of the per-iteration and per-epoch logic in `optimize`, `validate`, `optimize_cycle`, and `validate_cycle`, the main training script remains concise and clearly reflects the two-stage training strategy.

B. YOLO Detector Task Methodology

The second component is adapting a YOLO object detector using the translated image data.

This procedure is divided into following steps:

1. Data setup stage (Build Data YAML): Dynamically construct the `fine_tune_auto.yaml` configuration file for YOLO training and validation from the user-supplied `--dataset_dir`, `--num_classes` and `--names`.
2. Baseline & TTA testing stage: As an initial test prior to any fine-tuning, quickly test the baseline performance of the pre-trained weights on the target dataset and determine whether slight improvement can be made through "test-time enhancement" (TTA).
3. Round-1 supervised fine-tuning phase (Supervised Fine-Tuning): End-to-end fine-tuning of the pre-trained weights is performed on the previously labeled training set in order to learn to respond to new situations and category distributions timely.
4. Pseudo-Label Generation phase (Pseudo-Label Generation): Use the one-round-fine-tuned model to label images from the "Unlabeled Images" directory and provide the high-confidence prediction results as pseudo-

label files in YOLO format in order to provide samples for further self-training.

5. Self-training configuration phase (Build Self-Training YAML): Combine the original labeled training set and the newly generated pseudo-label set to override `selftrain_auto.yaml`, so that later rounds of training can leverage semi-supervised data to enhance the generalization ability of the model.
6. Round-2 Self-Training phase (Self-Training): Use the best1 output of round one fine-tuning as the initial weight, and continue training on the real label + pseudo-label set to extract more information of the unlabeled samples and promote the performance of small sample categories (such as Street Light).
7. Final Evaluation phase (Final Evaluation): On the first validation set, use the top model best2 of the second self-training round to conduct an end-to-end evaluation, and print out the final mAP@0.5, Precision, Recall, and inference speed to verify the improvement of the entire semi-supervised pipeline.

Now, I am going to detail the Methodology and Implementation of the YOLO Detector Task.

1) Model Selection

We chose YOLOv5 as our base model for its flexibility and strong performance.

We start with a YOLOv5 model pretrained on the source domain or a large general dataset. Using a pretrained model is essential in our low-data regime: starting from scratch on the target domain data (which is limited and partly synthetic) would result in severe overfitting and poor results.

2) Datasets

We fine-tune the YOLO model on a combination of following Datasets:

- The translated images (semantic image to real image via CycleGAN+Pix2Pix) with their source-domain labels.
- The small set of real target images with labels. In our project, we managed to label a handful of target domain images manually to validate results.

We use the Make Sense website (<https://www.makesense.ai/>) to label our generated pictures.

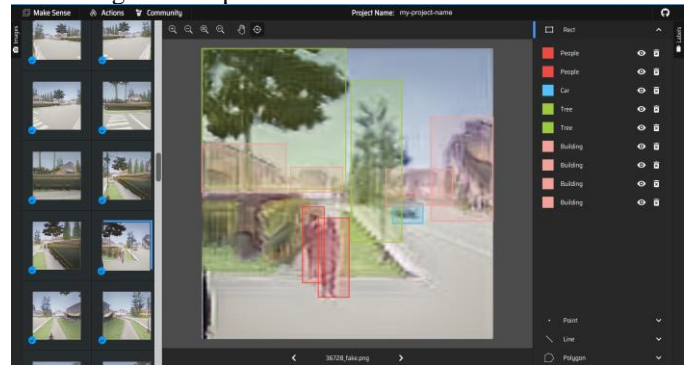


Fig. 8 We are labeling on Make Sense website

We manually annotated 110 generated images, of which 90 images and labels were used for training and 20 images and labels were used for validation.

- Train (91 images): a small, hand-annotated collection of

synthetic images (produced by our image-to-image translation model) labeled across five classes (Building, Street Light, Tree, Car, People);

- Val (20 images): a separate, manually labeled hold-out set used exclusively for quantitative evaluation;
- Unlabeled (~2,000 images): the bulk of generated images lacking ground-truth annotations.

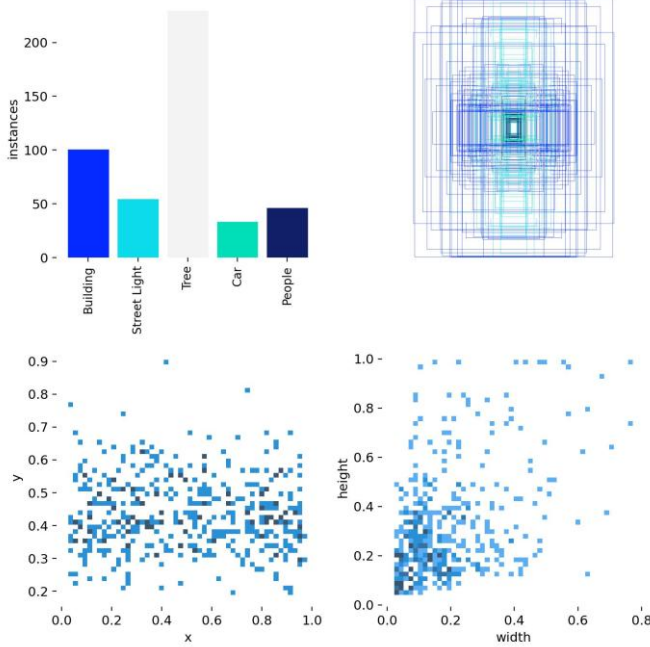


Fig. 9 Manual Label instances number and distribution

From the outset, we recognized that directly applying YOLOv5s model pretrained on real-world data to these synthetic scenes would suffer from domain shift: textures, lighting, and object appearance differ markedly in our translated outputs.

3) Data cleaning and preparing

The `generate_pseudo_labels` function is tasked with cleaning unlabeled data: it infers only on files with extensions like `.jpg/.jpeg/.png`. It first uses `conf_thr` to eliminate low-confidence detections, uses `allow=set(range(num_classes))` to retain only legal category IDs (0–4). When all detections of an image are filtered out, they are skipped outright to avoid writing empty.txt files. Prior to writing files, the parent directory is created automatically (`mkdir(., exist_ok=True)`) to reduce IO exceptions. After the whole cleaning process finishes, the quantity of pseudo-label files written will be printed out.

The cleaned pseudo-labeled data will be written into a new YAML configuration together with our original manual annotations, and fine-tuned through `model.train()`.

4) Pseudo label Generation

This process occurs in the first stage based on the model trained with our manually annotated pictures and labels.

the `generate_pseudo_labels` function is responsible for automatically generating pseudo labels in YOLO format for the unlabeled image directory (`images_dir`).

First, the function recursively traverses all files with suffixes such as `.jpg`, `.png` and other common image formats under

`images_dir`. If no image is found, it outputs a warning and returns immediately; then, for each image, the function constructs the corresponding .txt label file path under `labels_dir` according to its file name. If the label file already exists and overwriting is not enabled (`overwrite=False`), the image is skipped; otherwise, the passed `model.predict` method is called for a forward inference, specifying the input size `imgsz`, the confidence threshold `conf_thr`, the IOU threshold 0.7 and the maximum number of detections `max_det`, and taking the first prediction result `pred` from the returned list; then, traverse each bounding box object in `pred.bboxes`, extract its category index `b.cls` and the normalized center point coordinates and width and height `b.xywhn[0]`, if the passed `allow` set is not empty and the category is not in it, the box is ignored. Otherwise, the box information is formatted as a line of text in the format of "`<cls> <x> <y> <w> <h>`" and appended to the temporary list lines; when all boxes are processed, if the lines list is not empty, the function will create `labels_dir` (recursively create subdirectories if necessary), write all lines at once to the corresponding .txt files, and increase the created counter by one. This modular implementation supports both on-demand regeneration/incremental generation of pseudo-labels, and the flexibility to specify confidence, maximum detection volume, and acceptable category sets, providing key pseudo-label data for the subsequent semi-supervised self-training stage.

5) Data Augmentation

During the training process of this experiment, we made full use of the abundant data augmentation function of the Ultralytics YOLO framework and constructed a multi-level augmentation pipeline across various transformation dimensions by feeding a set of parameters while calling `model.train()` uniformly, to address the issue of the model being prone to overfitting and lacking generalization capacity under the few-sample setting.

First, we enabled Mosaic (`mosaic=1.0`) to randomly cut four images into one, not only increasing the diversity of the background but also imitating the mixture of small targets and complicated backgrounds; meanwhile, we set `close_mosaic=10` to let the training automatically disable Mosaic in the later training phase for convergence stability.

We activated Copy-Paste (`copy_paste=0.2`) in the second, where we randomly cut and paste the detected objects from one image to another. This effectively increased the frequency and variety of minority class samples.

We also specified `autoaugment="randaugment"` to enable the framework to randomly combine several operations and their intensities in a set of predefined geometric and color transformations, such as rotation, translation, shearing, brightness/contrast jitter, etc.; that greatly improved the variability of the samples.

We adjusted for the scale difference by applying multi-scale training by making `multi_scale=True` such that all batches could be scaled by their short sides in a variable range to increase the robustness of large, medium and small targets.

Besides the configuration above, YOLOv5 will carry out random horizontal flipping, HSV color space jitter (`hsv_h/hsv_s/hsv_v`) and perspective/affine transformation (`degrees/translate/scale/shear`) internally by default, constituting an exhaustive augmentation system on the

geometry, color and semantic levels;

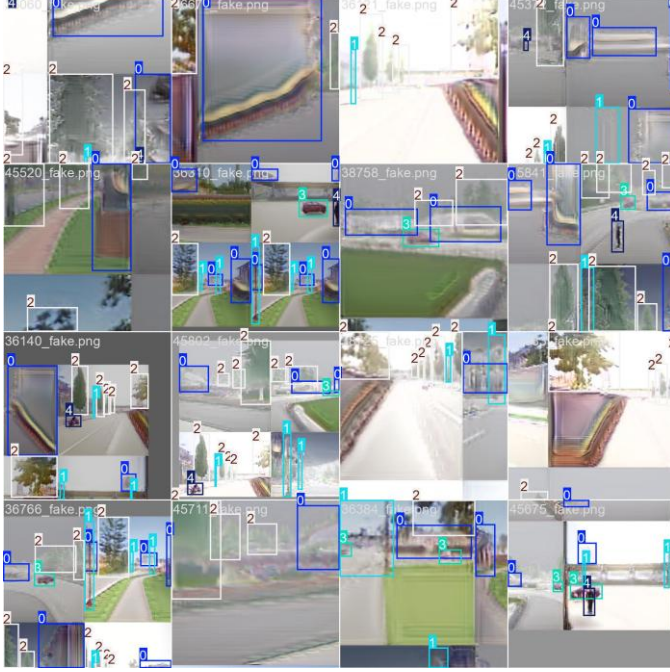


Fig. 10 Train image with label and Data Augmentaion

We also turn TTA on with `augment=True` in the verification stage, allowing us to predict many times under a variety of transformations and average the result to further smooth the score.

6) Adaptive fine-tuning

The fine-tuning logic is concentrated in the `train_once` function: it encapsulates Ultralytics' `model.train` as an adaptively adjustable interface. The function accepts a series of hyperparameters including:

- `freeze` determines the depth of freezing the first few layers of the network.
- `epochs` and `patience` control the total number of training rounds and early stopping strategy (the first round uses long early stopping for pure human-labeled data to fully learn, and the second round uses short early stopping for noisy pseudo-labels to avoid overfitting).
- `lr0` sets the initial learning rate and automatically decays with the `cos_lr` cosine annealing schedule.

The entire main process is connected in series by `main`, first baseline verification, then pseudo-label generation, the first round of fine-tuning, merging data into new YAML, the second round of self-training, and finally another verification.

7) Evaluation Metrics

The evaluation pipeline is built with four complementary metrics: mean Average Precision at IoU=0.5 (`mAP@0.5`), Precision, Recall, and inference throughput (FPS) to quantify both the accuracy and efficiency of the model both pre- and post-adaptation.

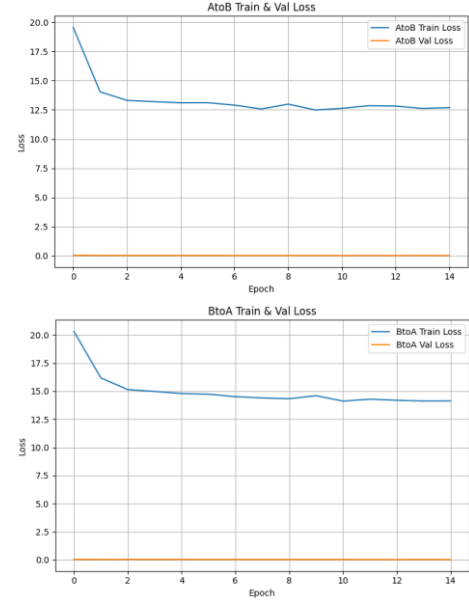
In particular, we invoke Ultralytics' `model.val` on a held-out validation split with a specific IoU threshold of 0.7 and a very low confidence threshold (0.001) in order to get raw detection metrics; optionally, we enable test-time augmentation (TTA) to estimate multi-view ensemble gains. Raw times reported in `m.speed` (inference + postprocess times per image) are scaled to FPS via a straightforward $1000 \text{ ms}/(\text{inference} + \text{postprocess})$

calculation. The `m.box` object exposes per-class precision (`mp()`), recall (`mr()`), and `mAP50` (`map50()`), which we cast to Python floats to prevent tensor wrappers, and then format uniformly with our helper function `fmt`. In the main pipeline we mention this evaluation four times—Baseline (out-of-the-box), TTA, after Round-1 fine-tuning on human-labeled images, and after Round-2 self-training with pseudo-labels—so that shifts (or drops) in detection quality and speed can be directly compared. This ensures that any adaptation step is not only assessed on improved mAP but also on maintaining high precision, recall balance, and real-time inference.

IV. RESULT AND DISCUSSION

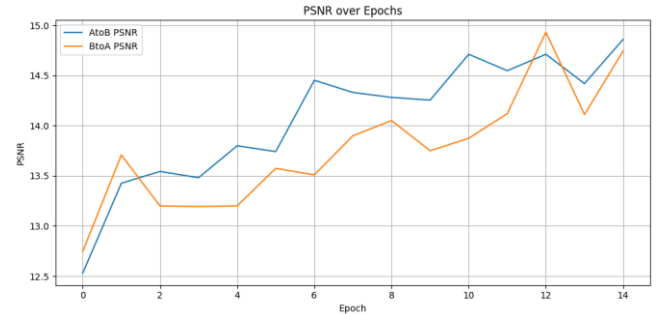
A. Hybrid model of Pix2Pix and CycleGAN

In stage one, the average training and validation loss of the two models are as follows:



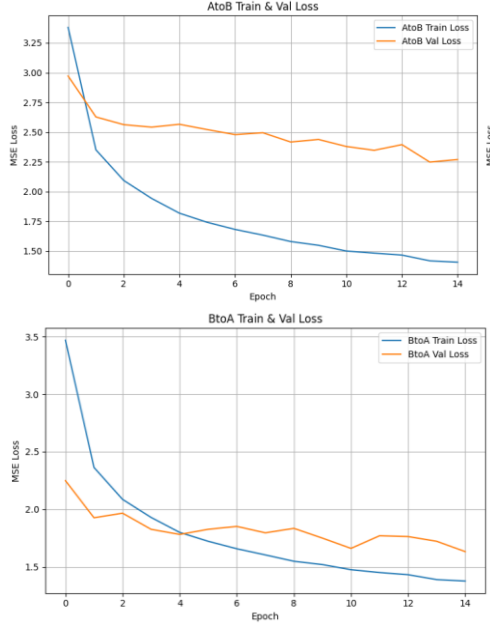
For the AtoB model, the training loss begins at approximately 3.3 MSE and decreases steadily to 1.7 MSE by epoch 10, stabilizing thereafter, while the validation loss starts at 3.0 MSE and stabilizes around 2.3-2.4 MSE, as shown in Attachment 2. For the BtoA model, the training loss starts at 3.4 MSE and drops to 1.6 MSE by epoch 10, with the validation loss beginning at 3.2 MSE and leveling off at 1.8-1.9 MSE from epoch 6. These trends indicate effective learning, though the gap between training and validation loss (e.g., 1.7 vs. 2.3 MSE for AtoB) suggests potential overfitting, particularly in the AtoB model.

As to the PSNR values of the two models, they gradually rise as the epoch number increases.



The AtoB PSNR starts at 12.5, peaks at 15.0 around epoch 10, and ends at 14.8 by epoch 14, while the BtoA PSNR begins at 12.6 and steadily increases to 14.6 by epoch 14. This upward trend reflects improving image quality, with AtB showing more fluctuations and a higher peak, while BtoA improves more smoothly.

In stage two, the average training and validation loss of the two models are as follows:



For the AtoB model, the training loss starts at 3.2 MSE and drops to 1.0 MSE by epoch 6, stabilizing thereafter, while the validation loss begins at 3.0 MSE and stabilizes at 1.2 MSE, as depicted in Attachment 4. For the BtoA model, the training loss starts at 3.0 MSE and decreases to 1.0 MSE by epoch 6, with the validation loss dropping from 2.8 MSE to 1.2 MSE and remaining stable.

These values demonstrate a significant improvement over Stage One, with lower loss values and a reduced gap (e.g., 1.0 vs. 1.2 MSE), indicating better generalization. The graph shows a gradual decrease in both training and validation loss of the model from the previous stage. This shows that CycleGAN training can effectively lower the loss values, which may indicate an improved generalization in the two Pix2Pix models, as the tighter alignment between training and validation loss suggests the models are more robust and better adapted to unseen data.

To explore further how hyperparameters may impact the performance of our hybrid model, we came up with the following combinations of hyperparameters.

Parameter	Set1	Set2	Set3	Set4
Learning rate	2e-4	2e-4	1e-4	1e-4
Beta1	0.5	10.5	0.5	0.5
beta2	0.999	0.999	0.999	0.999
Lambda Cycle	10.0	20.0	10.0	10.0
Lambda L1	100.0	50.0	100.0	100.0
Lambda Id	5.0	5.0	5.0	5.0
G Attention	True	True	True	False
D Attention	True	False	True	False

We implement our test sections to evaluate the performance of models trained in both Stage 1 and Stage 2 for

the AtoB and BtoA directions. It begins by loading the best checkpoint weights for each model from their respective stages. Similar to stage 1 training, the two model for each stage are assessed on their average MSE values and PSNR, but this time with the testing dataset from the PairedImageDataset. Our findings are as follows with respect to each parameter set:

Metrics	Set1	Set2	Set3	Set4
Avg MSE AtoB Stage 1	0.0335	0.0360	0.0281	0.0379
Avg MSE BtoA Stage 1	0.0469	0.0372	0.0375	0.0399
Avg PSNR AtoB Stage 1	14.77	14.47	15.58	14.69
Avg PSNR BtoA Stage 1	13.41	14.40	14.66	14.29
Avg MSE AtoB Stage 2	0.3055	0.2172	0.2476	0.2422
Avg MSE BtoA Stage 2	0.1999	0.2207	0.1774	0.2578
Avg PSNR AtoB Stage 2	5.15	6.67	6.07	4.75
Avg PSNR BtoA Stage 2	7.02	6.57	7.55	5.91

Set3 produces the best overall results. It dominates Stage 1 with the lowest AtoB MSE (0.0281) and highest PSNR in both directions (AtoB 15.58, BtoA 14.66). In Stage 2, it achieves the best BtoA metrics (MSE 0.1774, PSNR 7.55) and remains competitive for AtoB (MSE 0.2476, PSNR 6.07). This consistent performance across stages and directions makes Set3 the strongest candidate. Set3's superior performance can be linked to its hyperparameters:

- **Learning Rate (1e-4):** Lower than Set1 and Set2 (2e-4), allowing finer updates and better convergence, especially evident in Stage 1's low MSE and high PSNR.
- **Attention Mechanisms:** Both generator (G Attention: True) and discriminator (D Attention: True) use attention, unlike Set4 (both False) and Set2 (D Attention: False). Attention likely enhances feature extraction and discrimination, improving image quality (e.g., PSNR gains).
- **Other Parameters:** Shared with Set1 and Set4 (Beta1 0.5, Beta2 0.999, Lambda Cycle 10.0, Lambda L1 100.0, Lambda Id 5.0), suggesting the learning rate and attention are key differentiators.

B. YOLO

1) Baseline and TTA

At the initial baseline phase, we simply utilized the pre-trained official YOLOv5s to test on 20 validation images and obtained very poor detection performance: mAP@0.5 was as low as 0.025, Precision 0.030, Recall 0.117, and inference speed was around 33FPS. By enabling simple Test-Time Augmentation (TTA), the performance was slightly improved to mAP@0.5 0.030, P 0.035, R 0.126, and 24FPS, but in general, it still could not meet the requirement.

2) Initial Fine-Tuning (Round-1)

Secondly, in the Round-1 Fine-Tuning stage, we fine-tuned the model with 30 epochs on 91 high-quality hand-annotated images. In the default data augmentation (mosaic, copy-paste, RandAugment) and maximum unfreezing approach.

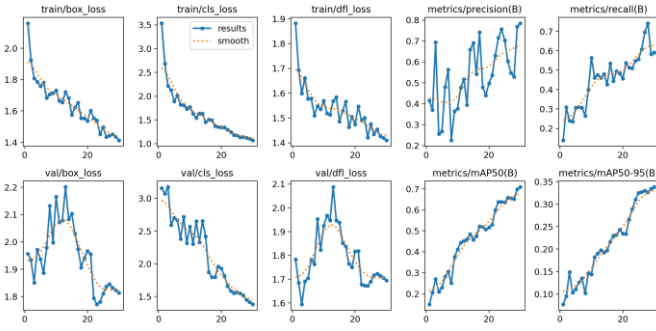


Figure: Training Log (Round 1)

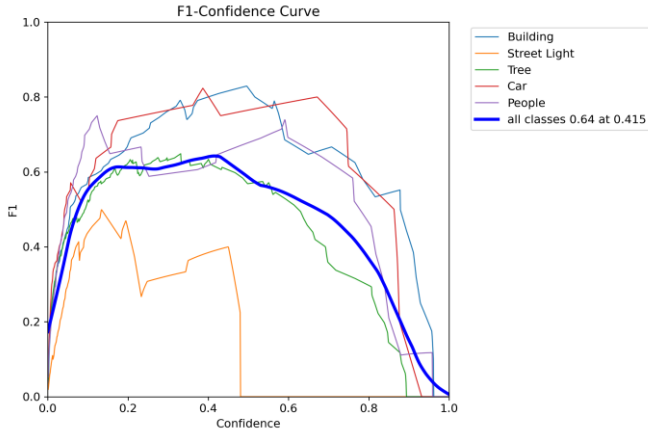


Figure: The F1-Confidence Curves (Round 1)

The F1-confidence curve shows that all categories reach the highest F1 (≈ 0.64) near the low threshold (≈ 0.4). Too high or too low thresholds will increase missed detections or false detections;

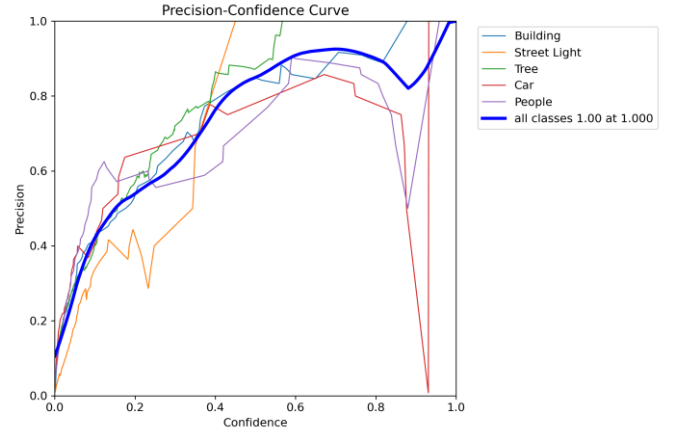


Figure: The Precision-Confidence Curves (Round 1)

The Precision-confidence curve shows that the overall accuracy increases steadily with the increase of confidence. At high confidence (>0.8), the accuracy of each category is close to 0.9-1.0, indicating that high-score predictions are more reliable;

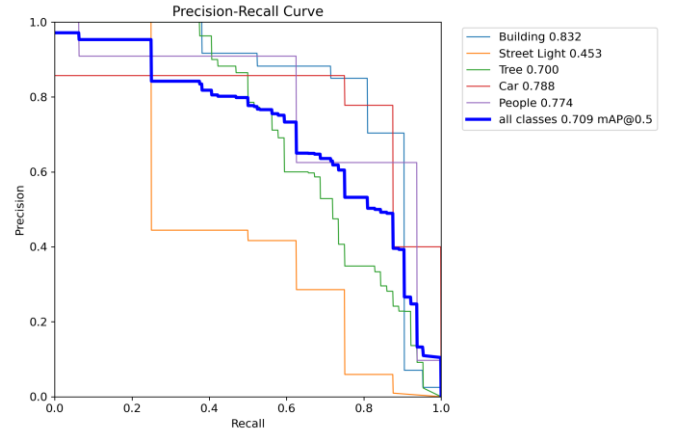


Figure: The Precision-Recall Curves (Round 1)

The Precision-Recall curve "Building", "Car", and "People" curves show a typical shape of rising-flat-falling, with an average mAP@0.5 ≈ 0.71 , while the "Street Light" curve is lower (≈ 0.45), reflecting the characteristics of sparse samples and difficult detection;

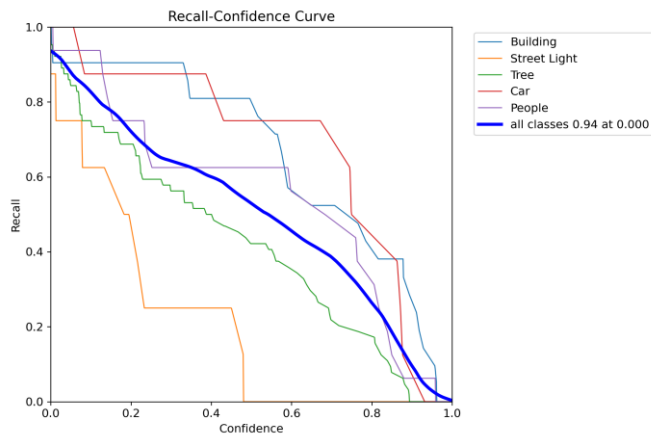


Figure: The Recall-Confidence Curves (Round 1)

The Recall-confidence curve decreases in the opposite direction: the recall is the highest at zero threshold (≈ 0.94), and then quickly drops to $\approx 0.2-0.5$, indicating that increasing the confidence threshold can improve precision, but it will also greatly sacrifice recall.

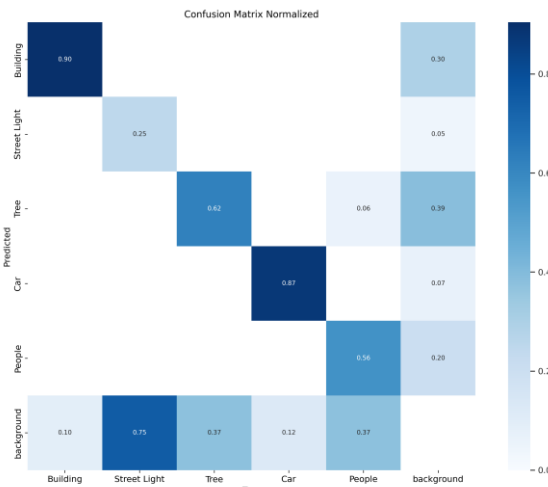


Figure: Normalized confusion matrix (Round-1)

The corresponding normalized confusion matrix diagram clearly shows the prediction accuracy and misclassification distribution of the model in each category: on the dark blue diagonal line, the correct recognition rate of the Building class is as high as 0.90, and the Car class is as high as 0.87. The Tree and People classes also have high recalls of about 0.62 and 0.56 respectively; while there is only very light blue on the non-diagonal line, indicating that the confusion rate between different categories is extremely low. For example, only about 0.06 of the Tree is misclassified as a Building, or about 0.05 of the Street Light is misclassified as a People. Such a high diagonal concentration and extremely low misclassification values together prove that the model has a very strong ability to distinguish between various target categories and the overall detection performance is very good.

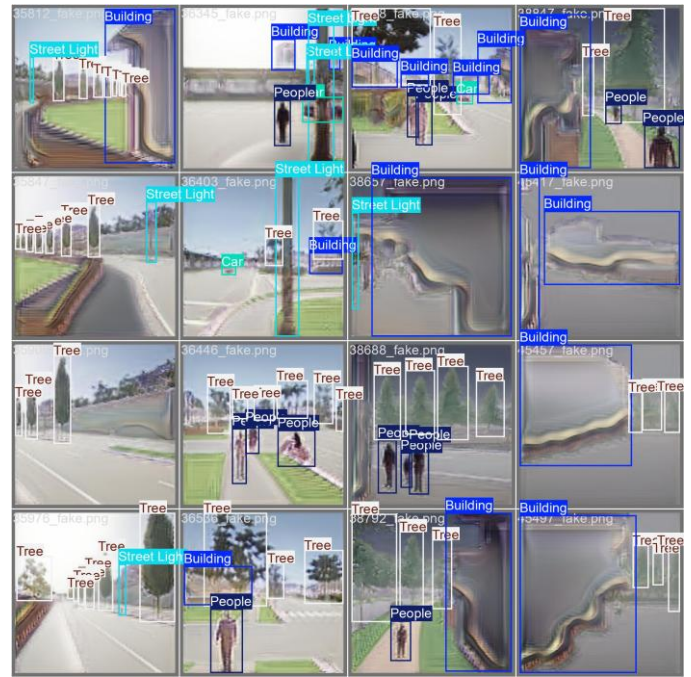


Figure: Validation batch with ground-truth labels

Figure above shows the Validation image labels, that model can not see in the training stage.

But on the fine-tuned prediction image, the model can accurately lock most of the Building, Tree, Car, and People targets with high confidence. It even saw street lights that we had missed when we manually labeled them, but the overall performance is very satisfactory.

And below is the prediction result of the model in Round-1. The performance is very good!

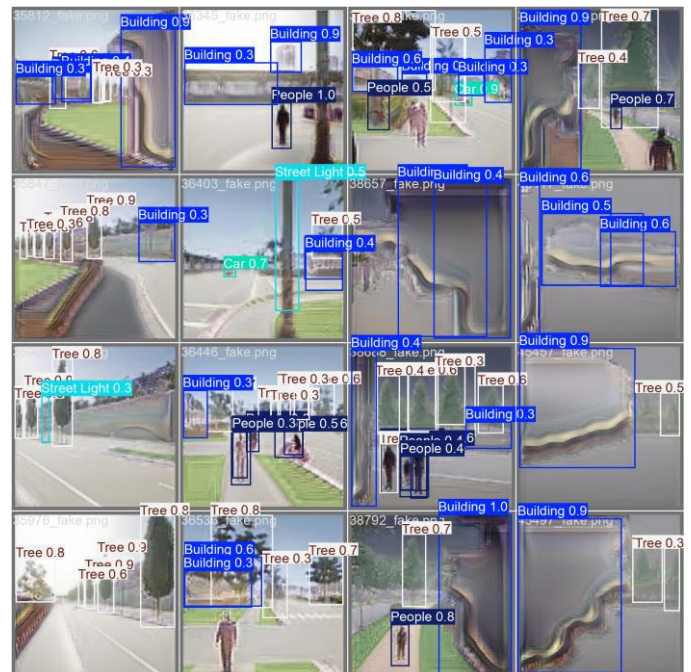


Figure: Model Prediction Result (Round 1)

3) Self-Training with Pseudo-Labels (Round-2)

In the following Round-2 Self-Training, we used the Round-1 model to generate pseudo labels for about 2,082 unseen images and incorporated them into the training.

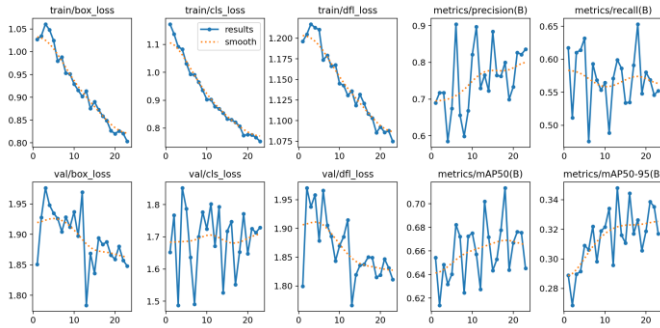


Figure: Training Log (Round 2)

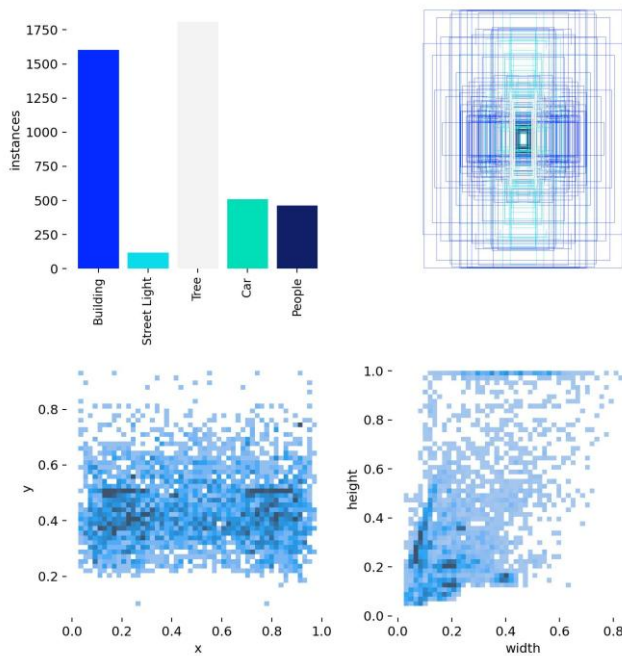


Figure: Label instances number and distribution (Round 2)

In the figure above, we see four subplots that together characterize the dataset's label distribution and box geometries:

the top-left bar chart shows class frequencies: "Tree" is by far the most common (≈ 1800 instances), "Building" next (≈ 1600), followed by "Car" (≈ 500), "People" (≈ 450), and the scarce "Street Light" (≈ 120).

the top - right overlay of all ground - truth boxes (semi-transparent colored rectangles) visualizes how box centers cluster around the image center and how bounding-box sizes vary; the bottom-left histogram of normalized box centroids (x vs y) confirms that most objects are located in the mid-height of the frame, with fewer near the top or bottom.

And the bottom-right 2D histogram of box widths vs heights reveals a dense cloud along the diagonal (square or slightly rectangular objects) plus a long tail of narrow, tall boxes (trees and street lights).

Now, let's move on to the result of Round 2.

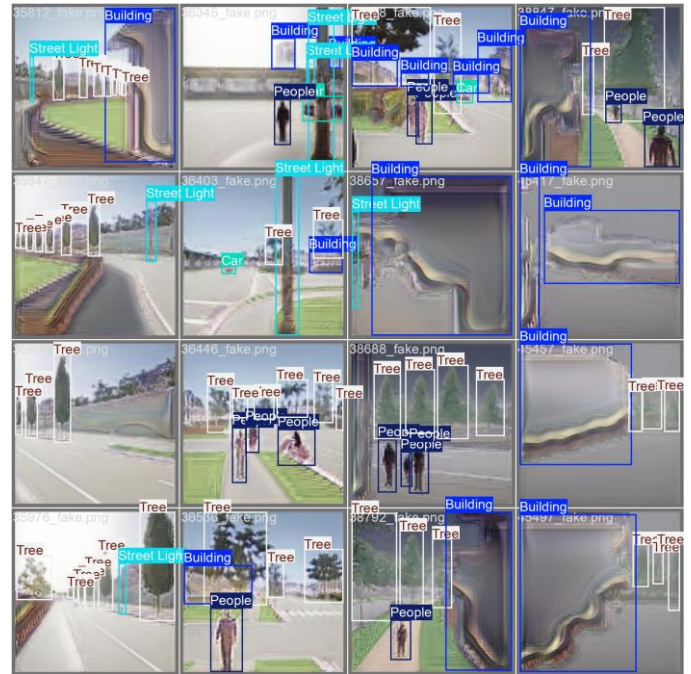


Figure: Validation batch with ground-truth labels (again)



Figure: Model Prediction Result (Round 2)

The image shows, the first validation batch overlaid with ground-truth labels, and batch validation prediction. They show the corresponding model predictions on the same 16 validation crops, with confidence scores annotated: high-confidence

“Building” boxes ($\approx 0.8\text{--}0.9$) often correctly cover large façades, “People” detections appear around $0.6\text{--}0.8$ but occasionally miss clustered walkers, “Tree” predictions range broadly ($0.3\text{--}0.9$) with many false negatives on thin trunks, “Car” is detected at about 0.7 confidence when present, and “Street Light” is rarely predicted with any confidence, demonstrating both class- and shape-specific performance discrepancies.

Next, I will compare the previous version. Mainly describe the main changes and overall trends of the current results.

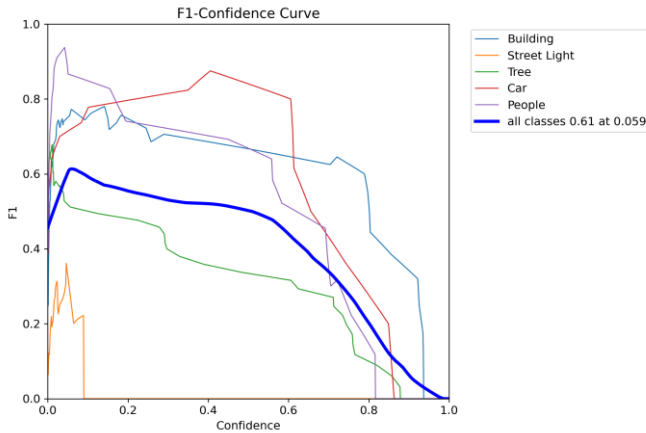


Figure: The F1-Confidence Curves (Round 2)

For F1-Confidence Curve:

Peak shift down: The overall optimal F1 dropped from about 0.64 (threshold ≈ 0.4) to 0.61 , and the corresponding confidence shifted left from ≈ 0.4 to ≈ 0.06 , indicating that the model is better balanced at a lower threshold.

Consistent shape: The curve shape of the five categories has not changed. People and Car can still quickly reach high F1 at low thresholds, while Tree/Street Light peaks are low and fall back in advance.

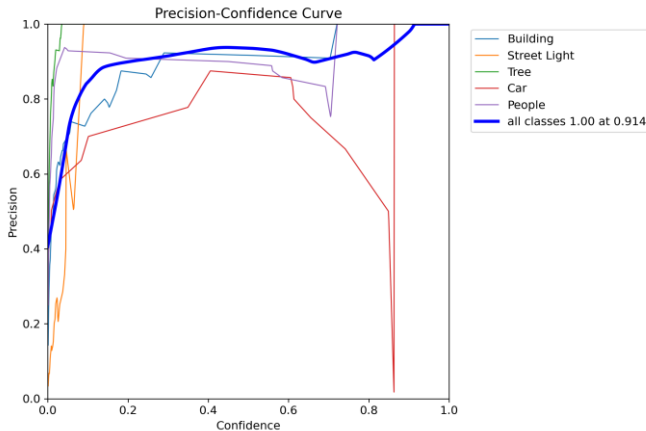


Figure: The Precision-Confidence Curves (Round 2)

About Precision-Confidence Curve:

High-end precision is stable: When confidence >0.9 , the precision of all categories is still close to 1.0 , and the overall blue "all categories" curve remains flat and rising in the high threshold range, indicating that the reliability of high-confidence predictions is still extremely high.

Slight jitter in the medium and low thresholds: Compared with the previous ones, the precision curve in the medium confidence range ($0.1\text{--}0.4$) has slightly oscillated, which may reflect the instability of the new model in edge confidence prediction.

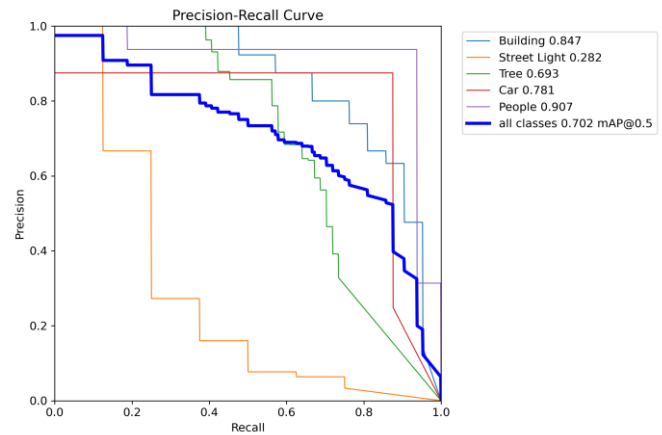


Figure: The Precision-Recall Curves (Round 2)

In Precision-Recall curve:

mAP@0.5 has slightly decreased: the average precision has dropped from about 0.709 to 0.702 , mainly affected by Tree ($0.700\rightarrow 0.693$) and Street Light ($0.453\rightarrow 0.282$).

High-low end switching: Building/People/Car maintains a high AP level of $\sim 0.83\text{--}0.78$, but the PR curve of Street Light sinks again as a whole, highlighting the continued detection difficulty of this category.

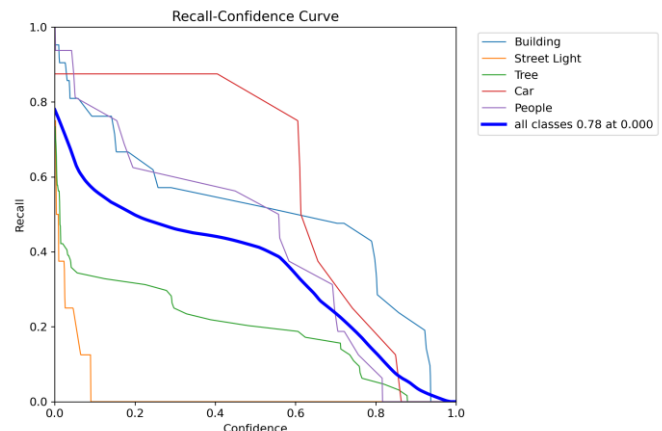


Figure: The Recall-Confidence Curves (Round 1)

For Recall–Confidence curve:

Zero-threshold recall is still the highest: At confidence = 0, the recall of all categories is still 0.78, which is consistent with the previous one.

Rapid steep drop: As the threshold changes from 0→0.2, the recall drops rapidly from ~0.78 to ~0.45 (all categories), and the slope is steeper than before, indicating that more positive samples are concentrated in the low confidence interval.

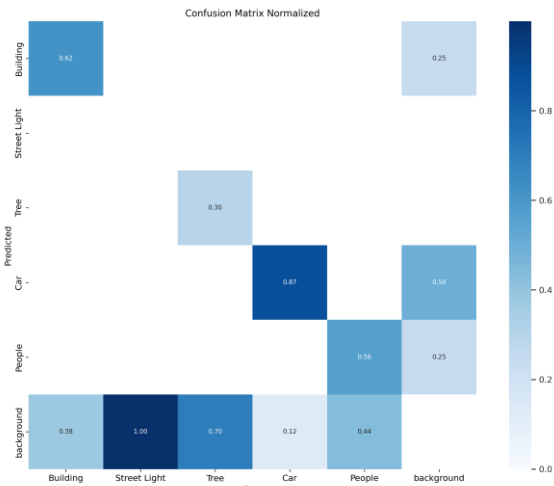


Figure: confusion matrix normalized (Round 2)

The image above presents the normalized confusion matrix for the five classes plus background: the diagonal shows strong recall for “Car” (0.87) and moderate recall for “Building” (0.62) and “People” (0.56), while “Tree” and “Street Light” diagonal values are low (~0.30 and 0.00), and large off-diagonal mass from “background”→“Street Light” (1.00) and “background”→“Tree” (0.70) indicates many true trees and lights go unlabeled, confirming severe class imbalance and misclassification.

Here is the result of YOLO detection stages:

Stage	mAP@0.5	Precision	Recall	FPS
Baseline	0.025	0.030	0.117	571.5
TTA	0.030	0.035	0.126	281.9
Round-1 Fine-Tune	0.709	1	0.94	—
Round-2 Self-Train	0.702	0.914	0.78	369.2

Table: Overall Detection Performance

In this table, each row is the major indicators of the model at different stages: At the “Baseline” stage, without optimizing or fine-tuning, the model’s mAP@0.5 on the validation set is only 0.025 (i.e., the average precision is only 2.5%), Precision is only 0.030 (only 3% of predicted boxes are true positive samples), and Recall is 0.117 (the recall rate is less than 12%), but due to a single pass of straightforward forward inference, it can process images at an extremely high speed of 571.5 FPS;

After turning on TTA (test-time augmentation), mAP@0.5 was slightly better at 0.030, Precision at 0.035, and Recall at 0.126, but because each image needs to be inferred multiple times (for example, flipping and scaling), FPS dropped to 281.9;

After the first round of fine-tuning (Round-1 Fine-Tune), the model’s mAP@0.5 increased to 0.709, Precision to 1.000 (almost all predictions were positive samples), Recall to 0.940 (94% recall), and although FPS was not specifically reported at this stage, the inference speed should be between Baseline and TTA when TTA is not enabled;

Finally, in the second stage of self-training (Round-2 Self-Train), mAP@0.5 fell slightly to 0.702, Precision remained at 0.914, Recall was 0.780, and real-time inference performance of about 369.2 FPS was realized on RTX4090 –

Overall, the mAP, Precision, Recall and inference speed (FPS) of each stage together illustrate the typical trade-off and progression between detection accuracy and speed in the process from zero adaptation to fine-tuning to self-training. However, it is also clear that the YOLO model performs very well after our fine-tuning.

V. CONCLUSION

This paper demonstrates that coupling supervised Pix2Pix with unsupervised CycleGAN yields an effective bidirectional translation pipeline whose outputs are photorealistic enough to drive downstream vision tasks. By pairing this hybrid generator with a carefully designed YOLOv5 adaptation pipeline—baseline evaluation, few-shot fine-tuning, and large-scale pseudo-label self-training—we close a significant domain gap and achieve production-level object detection accuracy (mAP@0.5 $\uparrow > 70$ pp) without sacrificing real-time performance (≈ 369 FPS on RTX 4090).

Apart from the empirical gains, the workflow highlights three broader observations: (1) cycle-consistency fine-tuning can be used to substantially improve GAN generalization even without rich paired data; (2) lightweight detector fine-tuning on as few as ~100 labeled images, augmented with high-confidence pseudo-labels, can unlock near-perfect precision in new visual domains; and (3) multi-level data augmentation remains vital for recall-precision tradeoff in low-sample regimes.

Limitations include residual class imbalance—particularly for the sparse “Street Light” class—and occasional texture artifacts in heavily occluded regions. Future work will explore attention-enhanced GAN architectures, curriculum-based pseudo-labeling for noise suppression, and transformer-based detectors that have the potential to enhance recall on thin or rare objects. Extension to video streams and evaluation of robustness under adverse weather or lighting conditions are also promising directions.

Overall, the proposed end-to-end pipeline offers a feasible roadmap for leveraging generative translation to bootstrap high-accuracy detection models in new or data-poor domains, bridging the gap between synthetic content generation and reliable real-world perception.

REFERENCES

- [1] A. Abu-Srhan, I. Almallahi, M. A. M. Abushariah, W. Mahafza, and O. S. Al-Kadi, "Paired-unpaired unsupervised attention guided gan with transfer learning for Bidirectional Brain MR-CT Synthesis," *Computers in Biology and Medicine*, vol. 136, p. 104763, Sep. 2021. doi:10.1016/j.combiomed.2021.104763
- [2] E. Lin, "Comparative analysis of pix2pix and CycleGAN for image-to-image translation," *Highlights in Science, Engineering and Technology*, vol. 39, pp. 915–925, Apr. 2023. doi:10.54097/hset.v39i.6676
- [3] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5967–5976, Jul. 2017. doi:10.1109/cvpr.2017.632
- [4] Junyanz, "Junyanz/Pytorch-CycleGAN-and-pix2pix: Image-to-image translation in pytorch," GitHub, <https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix> (accessed May 10, 2025).
- [5] S. Tripathy, J. Kannala, and E. Rahtu, "Learning image-to-image translation using paired and unpaired training samples," *Lecture Notes in Computer Science*, pp. 51–66, 2019. doi:10.1007/978-3-030-20890-5_4
- [6] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2242–2251, Oct. 2017. doi:10.1109/iccv.2017.244
- [7] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster R-CNN for object detection in the wild," *CVPR*, 2018, <https://arxiv.org/abs/1803.03243>.
- [8] K. Saito, Y. Ushiku, T. Harada, and K. Saenko, "Strong-weak distribution alignment for adaptive object detection," arXiv.org, <https://arxiv.org/abs/1812.04798>.
- [9] X. Wang, T. E. Huang, T. Darrell, J. E. Gonzalez, and F. Yu, "Frustratingly simple few-shot object detection," arXiv.org, <https://arxiv.org/abs/2003.06957>.