

# Comparison Between CNNs and Transformers for the Occupant Detection Task in Vehicle Interior Monitoring Systems

Abdul Bari  
Matrikelnummer: 00146949  
Technische Hochschule Ingolstadt  
Email: abb0233@thi.de

## **Declaration**

I, Abdul Bari, hereby declare that the seminar report, represents my ideas in my own words and where ideas or work of others have been included, I have adequately and accurately cited and referenced the original sources.

I also declare that I have adhered to the ethics of academic honesty and integrity and have not misrepresented or fabricated any data or idea or fact or source in my submission. I understand that any violation of the above will be a cause for disciplinary action by the institute and/or the University and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been obtained. This report has not been previously formed the basis for the award of any degree, diploma or similar title of any other University.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Foundational Architectures</b>	<b>5</b>
2.1	CNNs (Convolutional Neural Networks) . . . . .	5
2.2	CNNs to YOLO . . . . .	6
2.3	Transformers . . . . .	7
2.4	Vision Transformers (ViTs) . . . . .	9
2.5	Architectures Adopted for This Study . . . . .	10
2.5.1	YOLOv5l . . . . .	10
2.5.2	Swin Transformer . . . . .	11
<b>3</b>	<b>Experiment</b>	<b>13</b>
3.1	Training Setup . . . . .	13
3.1.1	Dataset . . . . .	13
3.1.2	Model Hyperparameters . . . . .	16
3.2	Model Evaluation . . . . .	17
<b>4</b>	<b>Conclusion</b>	<b>20</b>
<b>A</b>	<b>Configuration Files</b>	<b>21</b>
A.1	YOLOv5l . . . . .	21
A.2	Swin Transformer . . . . .	21
<b>B</b>	<b>Inference Code</b>	<b>29</b>
B.1	YOLOv5l . . . . .	29
B.2	Swin Transformer . . . . .	29

## Abstract

Interior monitoring systems play a crucial role in improving the safety and comfort of passengers in modern vehicles. One key task in this domain is **occupant detection**, which involves identifying and localizing individuals inside the vehicle cabin. With the rise of deep learning, **Convolutional Neural Networks (CNNs)** gained popularity for their efficiency and hierarchical feature learning [1] while **Transformer-based models** which were originally designed for natural language processing, have demonstrated strong global reasoning capabilities in vision tasks, leading to their rapid adoption in computer vision [2, 3, 4]. Despite their success, limited comparative studies exist evaluating these architectures specifically for **interior occupant detection**. This study addresses the gap by conducting a focused analysis by comparing a **CNN-based model (YOLOv5l)** [1] and a **Transformer-based model (Swin-T)** [4], both chosen for their similar parameter sizes and specialization in object detection tasks, allowing for a fair comparison using standard object detection metrics including **mean Average Precision (mAP)** and **generalization ability** across unseen occupants and lighting conditions. The models are fine-tuned and evaluated on the **SVIRO** [5] dataset, an open-source benchmark designed for realistic vehicle interior scenes. Through this comparative study, our aim is to investigate the practical trade-offs between **CNNs** and **Transformers** in the context of vehicle interior monitoring, providing insights into their strengths, limitations, and suitability for real-world applications.

## 1 Introduction

Computer vision has seen many advances since advances in deep learning [6], particularly with the introduction of convolutional neural networks and transformers, which have significantly improved the performance of computer vision models. Transformers, which were initially designed for NLP use cases, have recently shown remarkable skills in capturing long-range dependencies in images while traditional CNNs, with their ability to extract features in a hierarchical manner, have been the fundamental in many state-of-the-art image processing pipelines for many years now. Together, these architectures have taken computer vision tasks to an entirely new level.

One of the key computer vision tasks is object detection which identifies a single or multiple objects in an image or a video and localizes it by providing both the class labels and the bounding boxes for each detected object, unlike image classification that labels the entire image. This then allows systems to recognize not only what is in an image, but also locate its exact position. This ability is important for certain usecases that require a deeper understanding of the scene, allowing systems not just to recognize objects but also to figure out their exact locations. Such detailed visual information is important for various fields, including self-driving cars, robotics and security monitoring.

One of the main object detection applications is for vehicle interior monitoring systems, this can allow us to detect and track occupants as well as their activities inside the vehicle. Furthermore, this technique can support advanced driver-assistance systems (ADAS)

by ensuring driver attentiveness, preventing distracted driving, and enabling smart features like automatic climate control based on passenger locations allowing us to improve the user experience as well as the safety of the vehicle.

As vehicle interiors increase in complexity and expectations continue to increase, the need and requirement for reliable occupant detection models is more important than ever. Selecting the right architectures whether whether it is CNNs, which are great at extracting spatial features, or transformers, which can capture broader context, can really impact the system's performance. This research aims to compare these approaches specifically for occupant and object detection in vehicles, showcasing their different strengths and weaknesses in this emerging field.

## 2 Foundational Architectures

### 2.1 CNNs (Convolutional Neural Networks)

CNNs are deep learning models that specifically deal with grid-like data (i.e. visual data) [6]. First, why did we need them instead of traditional artificial neural networks (ANN)? Because, images are made up of pixels and inputting each pixel as an input data to an ANN is very resource intensive, moreover using an ANN is in this situation can lead to overfitting because inputting each and every pixel into the network can cause it to learn the image specific data instead of the general pattern in all the images, lastly spatial info is lost because we flatten down the entire image to input it into the ANN. Because of these reasons we use a CNN as it addresses all these issues.

CNNs consist of several key layers that extract key features from the images, of which the most important one is the convolutional layer which basically applies multiple filters that slide across the entire image applying a convolution operation on that patch. A convolution operation basically computes the dot products between the filter and the region the filter is placed on top of to produce a feature map, which highlights specific visual patterns such as edges, textures, or shapes.

Following the convolution layer we have an activation function that introduces non linearity in the model, this is a common thing to add throughout deep learning. The most common activation function used in a CNN is ReLU, the main reason for choosing this is to avoid vanishing gradients in deeper networks.

Next we have the pooling layer which helps shrink the spatial dimensions (height and width) of the feature maps. The most widely used method is max pooling, where usually a  $2 \times 2$  window moves across the feature map, keeping only the highest value from each section. Pooling has several benefits: it decreases the computational load, by introducing a degree of spatial invariance (so the model is less sensitive to small translations or distortions in the image), and helps reduce overfitting by compressing the representation.

After going through multiple convolutional and pooling layers, CNNs typically finish up with one or more fully connected layers. At this point, the multi-dimensional feature maps get flattened into a single vector, which then flows through these dense layers that learn to blend the extracted features for the final prediction. In classification tasks, this process usually ends with a softmax layer, which outputs a probability distribution over the possible classes, allowing the model to make a decision [7].

## 2.2 CNNs to YOLO

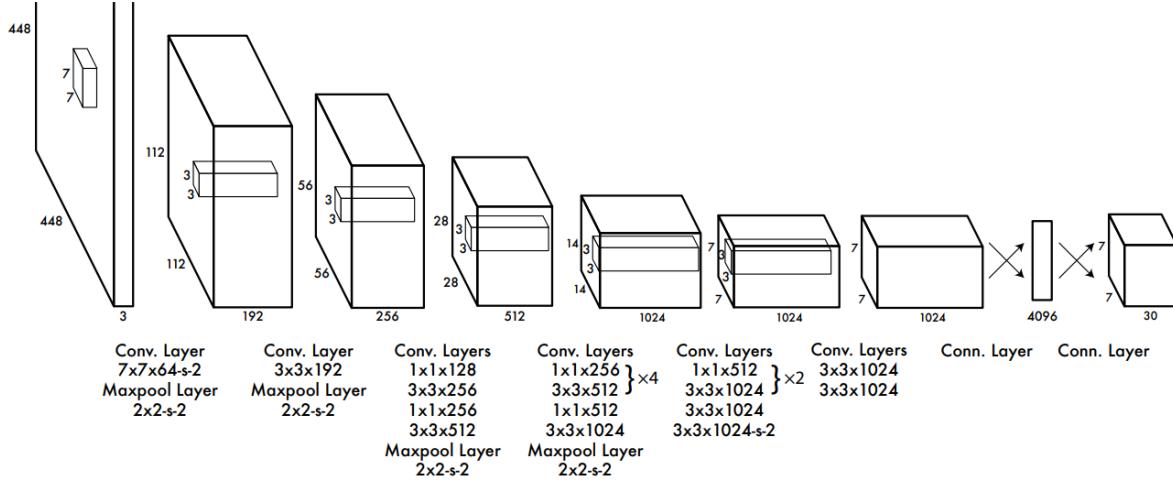


Figure 1: YOLO architecture [1]

Earlier methods of object detection used multi-stage pipelines [8] which involved separate models to predict the class labels and bounding boxes but YOLO combined them by treating object detection as a regression problem, directly predicting the bounding boxes and the class labels from an image in a single forward pass, hence the name You Only Look Once (YOLO). This resulted in significant speed improvements compared to previous methods.

The YOLO algorithm is foundationally based on CNNs. An image of fixed size is passed through multiple convolutional layers which extracts spatial features and high level patterns. These layers form the backbone of the entire model which are responsible for generating the feature maps. Instead of flattening these feature maps and using dense layers, which is commonly done in classification tasks, YOLO retains the spatial structure by applying additional convolutional layers that directly output detection results.

YOLO takes an input image and breaks it down into an  $S \times S$  grid. Each cell in this grid is tasked with predicting a set number of bounding boxes in the format  $(x, y, w, h)$  where  $x$  and  $y$  are the coordinates of the center of the box and  $w$  and  $h$  are width and height relative to the image dimensions (see Figure 1). Each prediction also includes a confidence score which tells us how likely it is that the box contains an object and class probabilities for object classification [1].

## 2.3 Transformers

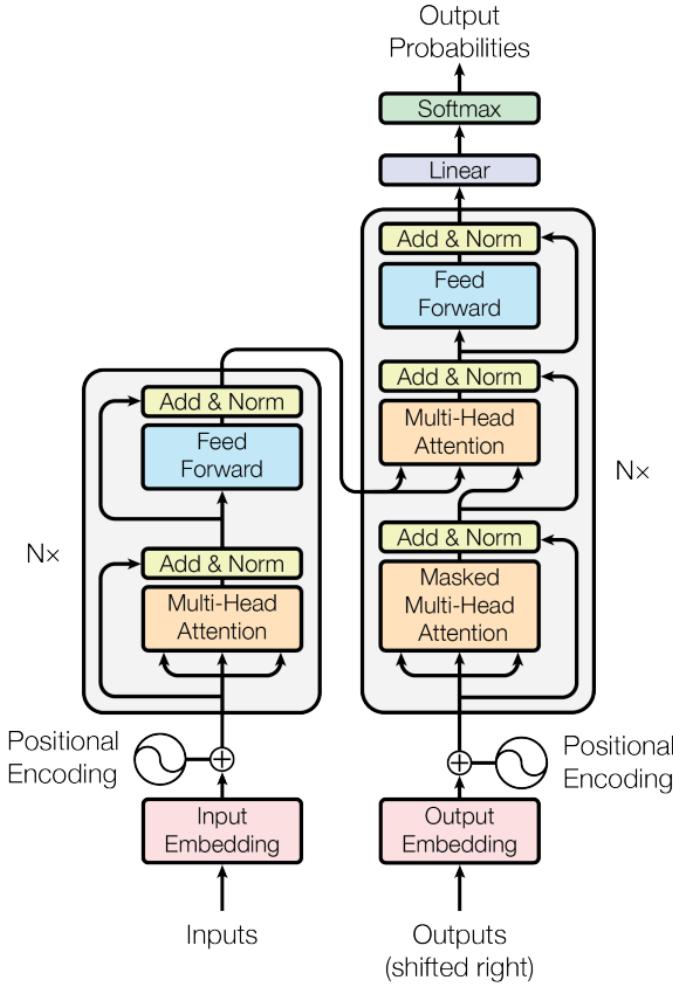


Figure 2: Transformer architecture [9]

The transformer model is a commonly used architecture in NLP applications because of its ability to learn sequential data much more efficiently than previous state of the art models such as LSTMs. Transformers basically entirely relies on the attention mechanisms particularly self attention [10], this then allows it create global dependencies with the added advantage of this being done parallelly allowing it to utilize modern GPUs thus significantly improving training speed and performance.

Transformers consist of two important parts, the encoder and the decoder with each of them composing multiple stacks of identical layers. The basic task of the encoder is to learn and produce rich contextual embeddings. However for the words to be inputted in the encoding layer, some key operations must be performed on them with the first one being performed by the embedding layer which transforms our input tokens into high dimensional vectors that preserve meaning of the token, moreover a positional encoding is added to each of the word embeddings which tell us about the relative position of the word in the sentence since transformers usually don't process data sequentially.

Processing in the encoder layer is done through mainly two layers, the multihead attention block which performs the most important operation in the entire architecture, self-attention, and the feedforward neural network. Residual connections and layer normalization is applied on the output of both of these layers (see Figure 2).

The fundamental component of the transformer architecture is the self attention mechanism [10], which transforms the previously static word embeddings into contextual embeddings. This basically tells us how much each token is related to each other token from the input sequence. This is done by deriving three further vectors from the input embedding, the query, key and the value vectors. These are derived through linear transformation which are learned from the training data. After deriving these vectors we use the following formula:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V$$

Here  $d_k$  is the dimension of the key vectors and the division by  $\sqrt{d_k}$  is done to help stabilize the gradients and improve the training performance.

Instead of using one self attention block usually multiple blocks are used, this is done to capture multiple perspectives and patterns that might be a bit nuanced, from the training data. This is done by passing our original word embedding through multiple matrices that generate multiple query, key and value vectors. Hence the name of the first block in the encoder layer being a multi-head self attention block.

After getting the contextual embeddings for each of the tokens, we pass them through the feed forward network consisting of two layers and the activation function for each of them being ReLU which introduces the non linearity in the data allowing the model to learn more complex mappings.

Coming to the decoder part, it generates the output sequence one token at a time making use of both the encoder's output and its own generated tokens, which it does similarly to the encoder by first making word embeddings and adding a positional encoding to each token in the sequence. This is then passed through a masked multihead attention which does self attention but with one change, each token in the output sequence can only attend to only previous tokens but not the future ones, this is done by applying a mask that blocks access to tokens ahead in the sequence. This is done to prevent data leakage since we are preventing the model to look ahead during training which wouldn't be possible in actual text generation.

Following this we perform cross attention using the output from the masked attention layer and the encoder. Here we try to link the input sequence to the output being generated. This is done by generating the query vector from the output of the masked attention and key and value vectors come from the output of the encoder. We then perform attention using these query, key and value vectors and the result of this effectively allows the decoder to focus on relevant parts of the input sequence while generating each token, enabling strong contextual associations between the input and output sequences.

These new contextual embeddings are then input in a feed forward network just like in the encoder consisting of two layers, each having the ReLU activation allowing the model to learn the non linearity in data. Moreover, just like in the encoder each sublayer is wrapped in a residual connection followed by a layer normalization operation.

Finally, the output layer applies a linear transformation followed by a softmax to the decoder’s final output, converting it into a probability distribution over the vocabulary (see Figure 2). This allows the model to predict the next token in the sequence during training or generation [9].

## 2.4 Vision Transformers (ViTs)

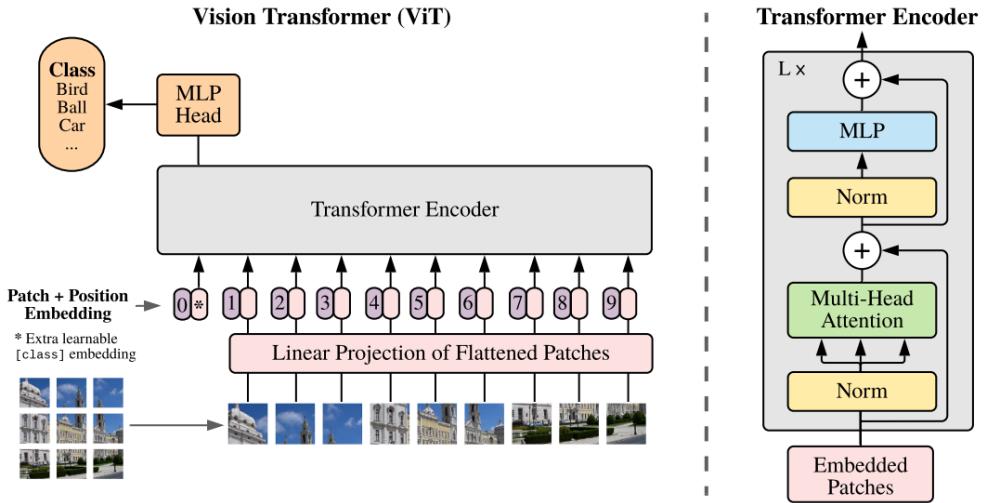


Figure 3: Vision Transformer architecture [3]

Motivated by the success of transformers in NLP tasks, researchers proposed the vision transformer [3, 11] which uses the self attention mechanism to extract feature maps from an image similar to that of CNNs. ViTs treat images in the same way Transformers treat text, by splitting them into a sequence of tokens and learning relationships among them without relying on convolutional operations (see Figure 3).

First we need to convert the input image into numbers for the model to understand it, similar to normal transformers that create word embeddings for each token, we need to create similar embeddings of our input image, this is done by first splitting our image into patches and flattening that patch into a 1D vector. The paper proposed an image of size  $224 \times 224$  pixels to be divided into  $16 \times 16$  patches, yielding 196 non-overlapping patches and the size of the flattened patch being  $16 \times 16 \times 3 = 768$  [3]. Each flattened vector is then passed through a trainable linear projection layer to map it into a fixed-dimensional embedding space. Just like the original transformer model we then add a positional encoding to each vector to tell us about the relative position of the patch since transformers usually don’t process data

sequentially. These patch embeddings with positional encoding serve as the input tokens to the Transformer model.

The Transformer encoder used in ViTs is similar to the one used in NLP models. It consists of multiple identical layers, each composed of two main sub-layers: a multi-head self-attention block and a feed-forward network. Layer normalization and residual connections are applied around both sub-layers to stabilize training (see Figure 3). The self-attention mechanism allows each image patch to attend to every other patch, capturing both local and global interactions in a single layer. Additionally, a special [CLS] token is prepended to the patch sequence, whose representation at the final encoder layer is intended to capture a global summary of the image. For classification tasks, this representation is passed through a fully connected head and a softmax layer to produce class probabilities.

Even though Vision Transformers have been quite successful [3], they do come with their own set of challenges. Unlike CNNs, they lack some built-in assumptions, such as the ability to focus on nearby pixels or effectively manage shifts in images (translation invariance). This can make them less efficient, especially when working with smaller datasets. Plus, ViTs require significantly more data and computing power to really shine. On top of that, the attention mechanism scales quadratically as the image size increases, which can pose issues for larger or more intricate tasks like object detection.

## 2.5 Architectures Adopted for This Study

### 2.5.1 YOLOv5l

YOLOv5l is one of the larger variants in the YOLOv5 family of object detection models, created and maintained by Ultralytics. Just like its siblings, YOLOv5l utilizes a fully convolutional architecture that allows it to detect objects in a single forward pass, which makes it incredibly fast and perfect for real-time applications.

It adheres to the standard YOLO workflow: an image is processed through a CSPDarknet53 backbone, based on Cross Stage Partial connections, which helps in reducing computations while maintaining a strong gradient flow resulting in an improved learning capability. This backbone is responsible for extracting high-level visual features from the input image.

Next is the PANet (Path Aggregation Network) neck, which is important for creating feature pyramids which allow the model to combine and refine features at different scales, improving its ability to detect objects of different sizes in the same image. This is especially useful in scenarios like vehicle interior monitoring, where both large (e.g., a human body) and small (e.g., hands, mobile devices) objects may need to be detected simultaneously.

Finally, we have a detection head that produces bounding box coordinates, objectness scores, and class probabilities. To refine these outputs, we apply a technique called non-maximum suppression (NMS). This helps us remove overlapping predictions and keep only the reliable ones.

We chose YOLOv5l for this study because it strikes balance between detection accuracy and inference speed. Its fully convolutional design enables real-time object detection, which is important for occupant monitoring systems that require a fast response. Moreover, when compared to the smaller YOLO variants, YOLOv5l’s deeper and wider architecture really enhances its ability to spot objects of different sizes and in various lighting conditions that are often encountered inside vehicles.

### 2.5.2 Swin Transformer

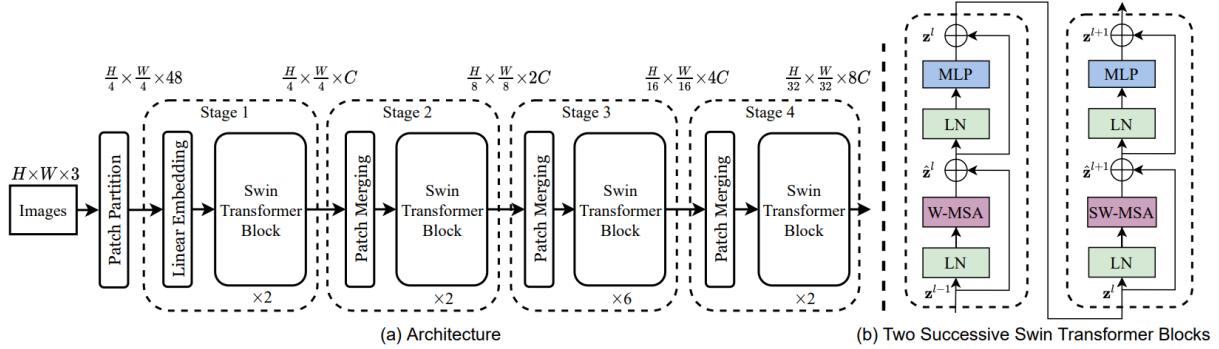


Figure 4: Swin Transformer architecture [4]

The Swin Transformer (Swin-T) addresses key drawbacks of the original Vision Transformer (ViT), such as its inability to focus on local features and high computational cost from global self-attention. ViTs perform well on classification tasks, but struggle with object detection tasks which require dense predictions, as discussed previously.

Swin-T introduces shifted window-based self-attention, which alternates across layers (see Figure 4). In the first layer of each stage, self-attention is computed within distinct, non-overlapping windows. The following layer then has the window partitioning shifted by a small offset (typically half the window size), allowing the model to connect with nearby windows. This strategy maintains the efficiency of local attention while enabling cross-window communication, effectively increasing the context without the cost of full global attention.

This architectural flexibility is further enhanced by a hierarchical design where feature maps are progressively downsampled using patch merging, similar to CNNs, making it more suitable for multi-scale vision tasks. This works well as a backbone when combined with RetinaNet or Mask R-CNN, resulting in rich, multi-scale features for accurate object detection.

The Swin Transformer, when paired with RetinaNet, was chosen for its ability to model both multi-scale and contextual information. Its hierarchical, shifted window self-attention mechanism allows for efficient extraction of global and local features, making it particularly effective at detecting partially occluded occupants or small objects in messy vehicle interiors. When combined with RetinaNet detection head, this setup provides accurate and dependable object detection even in complex situations, complementing the performance of YOLOv5l.

The following table summarizes the key differences, trade-offs, and characteristics of the two selected object detection models: YOLOv5l and Swin Transformer with RetinaNet.

Aspect	YOLOv5l	Swin Transformer (with RetinaNet)
Model Type	CNN-based single-stage detector	Transformer-based backbone with RetinaNet head (two-stage-like)
Backbone	CSPDarknet53	Swin-T (Shifted Window Transformer)
Head	YOLO detection head (anchor-based)	RetinaNet head (anchor-based, focal loss)
Feature Fusion	PANet (Path Aggregation Network)	FPN (Feature Pyramid Network)
Architecture Nature	Fully convolutional	Hierarchical transformer with local-global attention
Training Complexity	Easier and faster to train	More complex; higher memory and time requirements
Data Efficiency	Performs well on medium datasets	Requires more data for generalization
Augmentation Support	Strong built-in augmentation (Ultralytics)	Requires external pipeline (e.g., MMDetection)

Table 1: Comparison between YOLOv5l and Swin Transformer with RetinaNet

Model	mAP@[.5, .95]	Params (M)	GFLOPs	Ref
CNN (ResNet-101)	27.2	~ 44.5	7.6	[12]
ViT (DETR)	42.0	41.3	86	[2]
YOLOv3	33.0	61.9	65.9	[13]
YOLOv5l	49.0	46.5	109.1	[14]
Swin-T + Mask R-CNN	50.5	86	745	[4]

Table 2: Benchmark comparison of the discussed architectures on COCO dataset.

## 3 Experiment

The research focuses on two models: YOLOv5l and Swin Transformer paired with a RetinaNet head. The YOLOv5l model comes from the official Ultralytics GitHub repository. It has been pretrained on the COCO (Common Objects in Context) dataset [15], featuring over 200,000 labeled images across 80 different object categories. These pretrained weights enable the model to perform well on a variety of object detection tasks, especially when it's fine-tuned on a smaller, more specific dataset like SVIRO [5]. To set up the model, the yolov5 repository [14] was cloned, dependencies were installed via requirements.txt, and the training was initiated using a custom YAML config file specifying the dataset paths, class labels, and hyperparameters. The training was done using the provided train.py script.

On the other hand, the Swin Transformer, which was used alongside a RetinaNet detection head, was obtained from the MMDetection model zoo by OpenMMLab. This model has also been pretrained on the COCO dataset. Starting with strong feature representations from COCO allows the model to be fine-tuned effectively for tasks related to monitoring vehicle interiors, specifically tailored to the SVIRO dataset. To set up the Swin-T + RetinaNet model, we used the MMDetection toolbox from OpenMMLab. The process started with cloning the MMDetection repository [16] and installing all the necessary packages, including MMCV. We also had to tweak the existing Swin Transformer configuration file to match the details of the SVIRO dataset, such as the number of classes, the dataset format, and the file paths. Training was carried out using the tools/train.py script from MMDetection, and we evaluated the results with the built-in COCO evaluation tools that come with the framework.

### 3.1 Training Setup

#### 3.1.1 Dataset

The models described above are fine-tuned on the SVIRO dataset [5] (specifically for the object detection use case using the bounding box labels), which is a synthetic dataset specifically designed for vehicle interior monitoring tasks. Our training dataset contained **1732** grayscale training images, and our test set contained **500** grayscale testing images and we use it specifically for occupant and object detection in the back seat of the vehicle.

The dataset contains the following classes as shown in Figure 5.

- 1 = Infant seat
- 2 = Child seat
- 3 = Person
- 4 = Everyday object

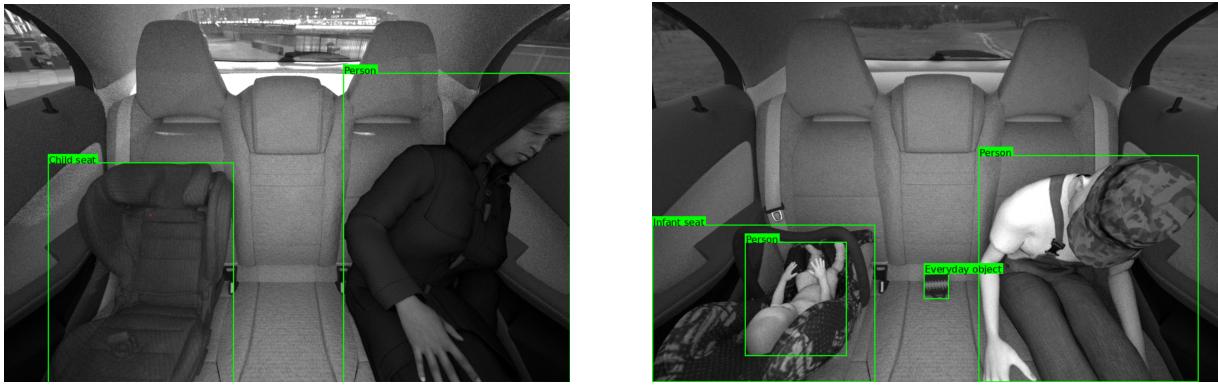


Figure 5: Example images from our dataset with their bounding box labels

## Original Format

The bounding box labels in the original dataset are in the format:

`[class_id, x_min, y_min, x_max, y_max]`

where:

- `class_id`: Integer representing the class
- `x_min`: x-coordinate of the top-left corner of the bounding box
- `y_min`: y-coordinate of the top-left corner of the bounding box
- `x_max`: x-coordinate of the bottom-right corner of the bounding box
- `y_max`: y-coordinate of the bottom-right corner of the bounding box

## YOLOv5l Format

For the YOLOv5l model, these labels were converted to the YOLO format:

`[class_id, x_center, y_center, width, height]`

where:

- `class_id`: Integer representing the class

- **x\_center**: x-coordinate of the center of the bounding box, normalized by image width
- **y\_center**: y-coordinate of the center of the bounding box, normalized by image height
- **width**: Width of the bounding box, normalized by image width
- **height**: Height of the bounding box, normalized by image height

## Swin Transformer Format (COCO)

For the Swin Transformer model, labels were converted to the COCO format:

```
{
  "image_id": int,
  "category_id": int,
  "bbox": [x_min, y_min, width, height],
  "area": float,
  "iscrowd": 0,
  "id": int
}
```

where:

- **image\_id**: Unique identifier for the image
- **category\_id**: Integer ID representing the object class
- **bbox**: Bounding box in [x\_min, y\_min, width, height] format (absolute pixel values)
- **area**: Area of the bounding box (width × height)
- **iscrowd**: Indicates whether the object is part of a crowd (set to 0)

### 3.1.2 Model Hyperparameters

Hyperparameter	YOLOv5l	Swin-T + RetinaNet
Epochs	200	50
Batch Size	8	16
Image Size	$640 \times 640$	$448 \times 448$
Optimizer	SGD	SGD
Initial Learning Rate	0.01	0.005
LR Schedule	Linear decay	Cosine annealing
Warmup Strategy	3 epochs (momentum 0.8 to 0.937, warmup_bias_lr = 0.1)	500 iters to 10% of initial LR
Momentum	0.937	0.9
Weight Decay	0.0005	0.0001
Gradient Clipping	–	Max norm = 35
Cls Loss Type	BCE with logits	Focal Loss ( $\gamma=2.0$ , $\alpha=0.25$ )
Backbone	CSPDarknet	Swin-Tiny (96-dim, [2,2,6,2], win=7)
Drop Path Rate	–	0.2
Data Augmentations	HSV jitter, scaling, translation, flipping, mosaic	Resize, flip, affine (scale, rotate, shear), photometric distortions
Dataloader Workers	8	4

Table 3: Comparison of the training setups and architectural differences between YOLOv5l and Swin-T + RetinaNet. YOLOv5l emphasizes speed and simplicity using mosaic augmentation, CIoU and objectness loss, and a CNN-based CSPDarknet backbone. In contrast, Swin-T + RetinaNet focuses on semantic depth and contextual awareness with a transformer-based backbone, cosine annealing learning rate, and focal loss to address class imbalance.

### 3.2 Model Evaluation

Metric	YOLOv5l	Swin-T+RetinaNet
Precision	0.952	0.119
Recall	0.862	0.473
mAP@0.5	0.948	0.203
mAP@0.5:0.95	0.829	0.119
Training Epochs	56	50
Parameters (Millions)	46.5	52.8

Table 4: Performance comparison of YOLOv5l and Swin-T+RetinaNet models on test data

### YOLOv5l

To assess how well the YOLOv5l model performed, we used a test set made up of 500 grayscale images of interiors. Even though according to Table 3 we set the training to run for 200 epochs, the model actually converged early and wrapped up at just 56 epochs, thanks to the early stopping rule we applied. When we evaluated the best model checkpoint on the test set, it delivered impressive results: an overall Precision (P) of 0.952, a Recall (R) of 0.862, and a mean Average Precision at IoU=0.5 (mAP@0.5) of 0.948, along with a mAP@0.5:0.95 of 0.829 (see Table 4) [17]. These results indicate that YOLOv5 performs reliably and efficiently for real-time object detection in grayscale interior environments.

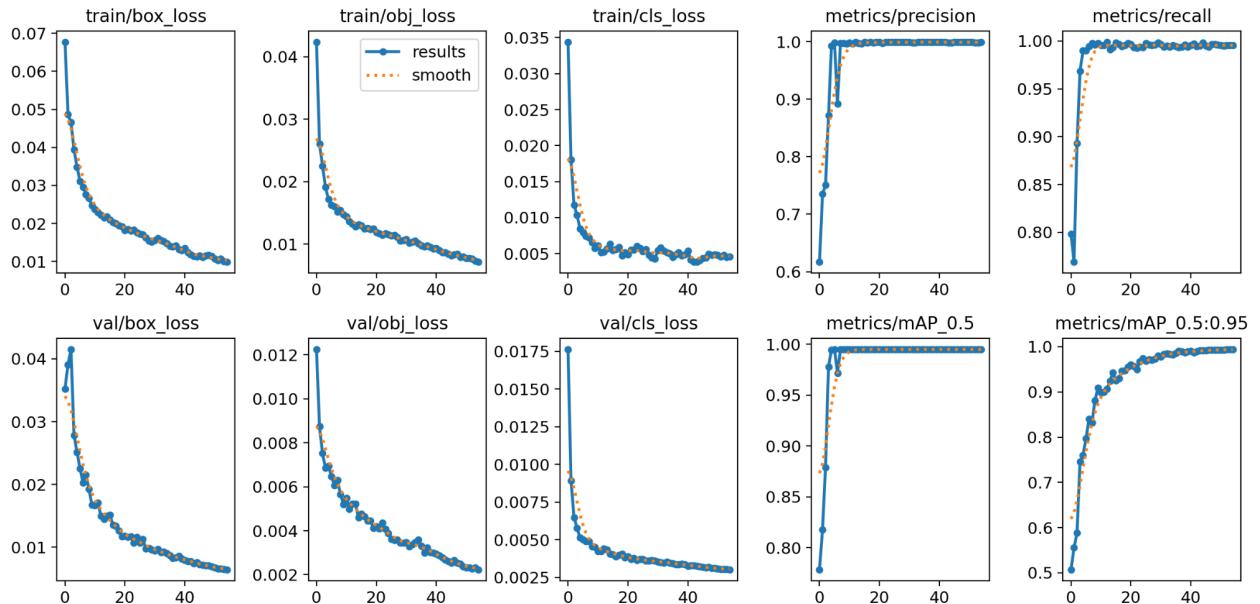


Figure 6: The results clearly show a consistent drop in all loss components, which points to stable training. At the same time, the precision, recall, and mAP metrics demonstrate a quick boost during the initial epochs before leveling off, indicating that the YOLOv5l model is effectively learning and converging.

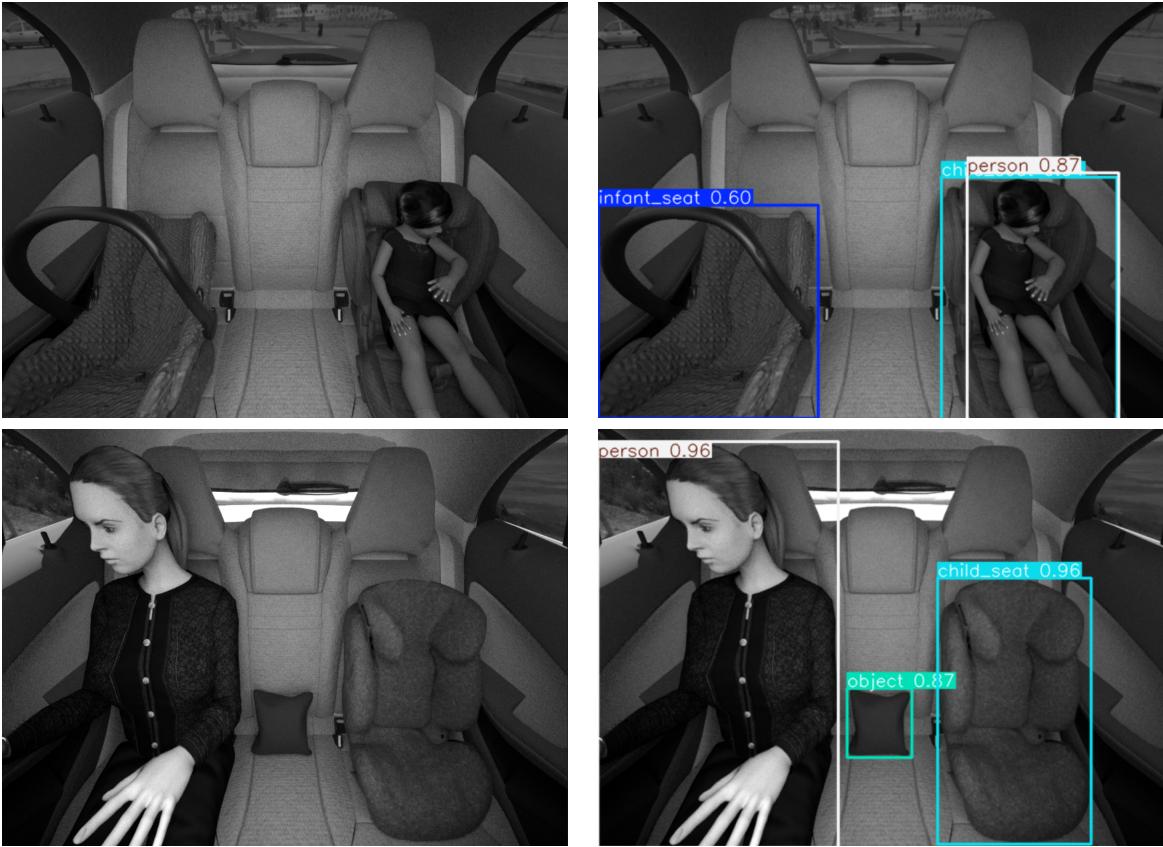


Figure 7: Detection results comparison: original images (left) with YOLO outputs (right)

Figure 7 shows the detection results from the YOLO model alongside the original input images. In the left column, we have the unannotated synthetic images of the interior, while the right column features the YOLO predictions, complete with bounding boxes and confidence scores for detected classes described above. These results really showcase the model’s impressive ability to accurately locate and classify important interior elements and occupants.

## Swin Transformer

The Swin Transformer model was trained for 50 epochs and evaluated on the same test set. It managed to achieve an Average Recall of 0.473, which means it successfully detected nearly half of the objects in the dataset. However, there was a noticeable precision-recall imbalance, leading to an overall mAP@0.5 of 0.203 and mAP@0.5:0.95 of 0.119. The model’s performance was primarily focused on large objects, with a mAP<sub>large</sub> of 0.124, and it completely missed detecting small or medium-sized objects. This aligns with the dataset’s characteristics, where most of the target objects (infant seats, child seats, persons) are relatively large within the grayscale interior scenes. While the model shows a strong capability for detection, it still needs some optimization to reduce false positives before it can be considered ready for practical deployment.

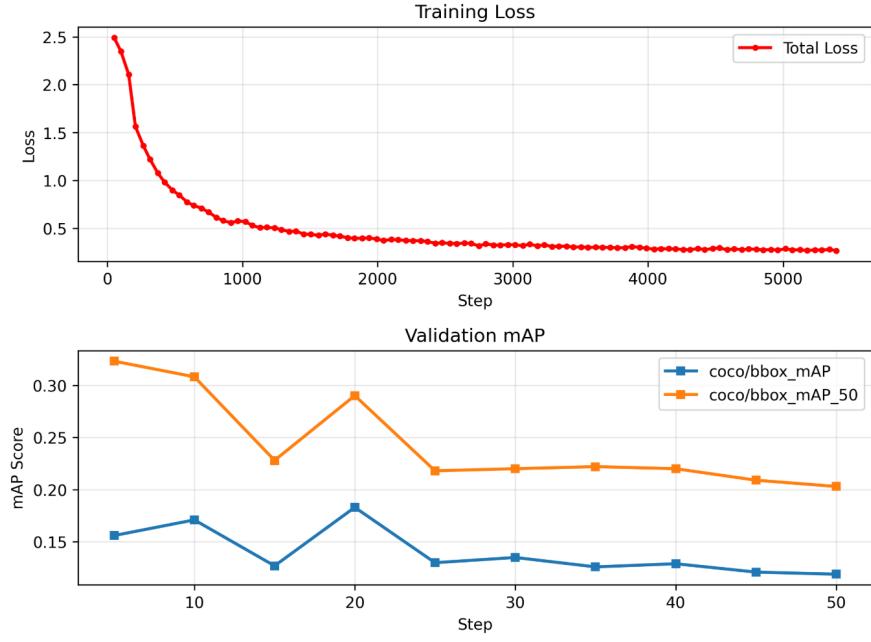


Figure 8: The results show a consistent drop in total loss, which shows that the model is converging nicely. On the other hand, the validation mAP scores are low and scattered. This implies that while the Swin-T model is good at cutting down training loss, but struggles to maintain steady detection performance when it comes to validation data.

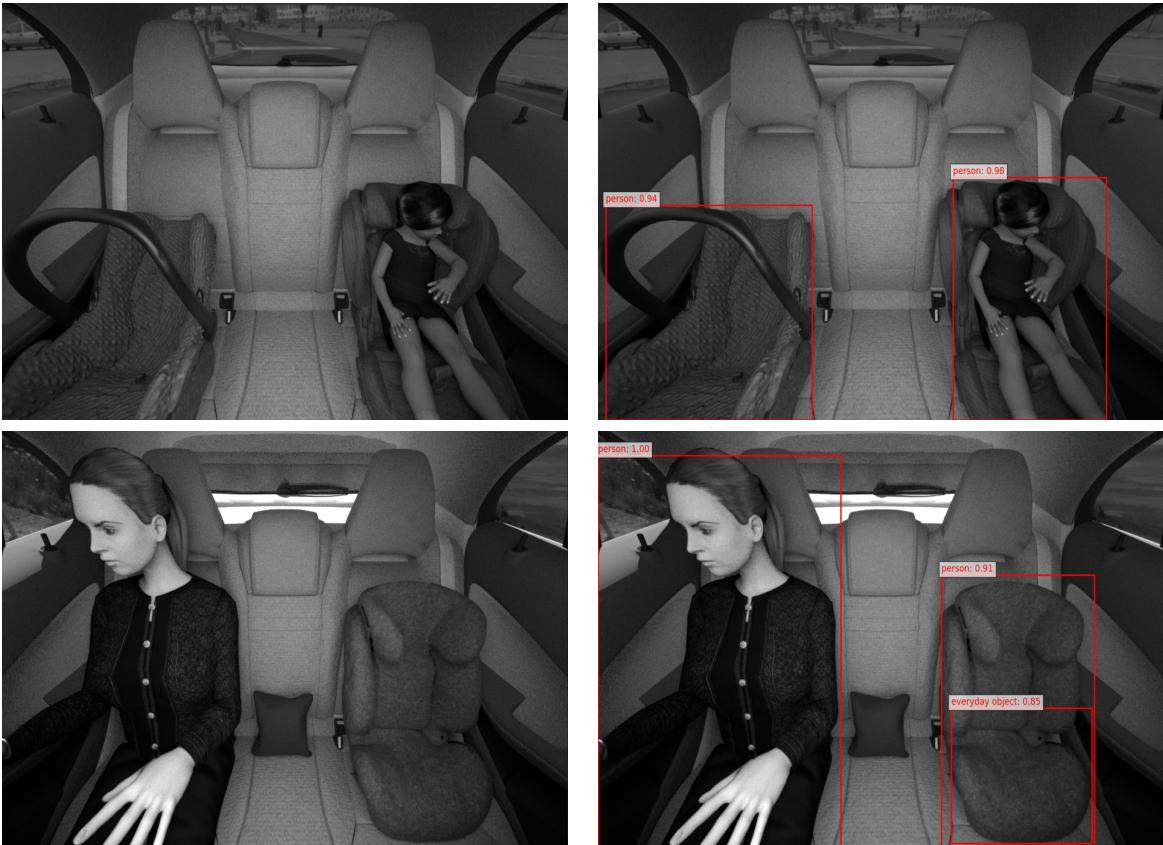


Figure 9: Detection results comparison: original images (left) with Swin-T outputs (right)

Figure 9 shows the detection results from the Swin-T model alongside the original input images. In the left column, we have the same unannotated synthetic interior images, while the right column features the Swin-T predictions, complete with red bounding boxes and confidence scores. The model shows strong performance for person detection. However, it takes a broader approach to classification, labeling objects simply as "everyday object" instead of specifying seat types like the YOLO model does. Even with this difference, the Swin-T model consistently and accurately pinpoints larger elements within the vehicle interior scenes.

## 4 Conclusion

While the Swin Transformer has shown impressive results in various vision tasks, it didn't quite match the performance of YOLOv5l in this particular case. There are a few reasons that might explain this difference, such as the relatively small training dataset and the use of grayscale images, which can restrict the model's ability to fully utilize its architectural advantages. On the other hand, YOLOv5l, with its convolutional neural network (CNN) backbone and robust data augmentation techniques, demonstrated better generalization and more efficient inference. These findings indicate that for tasks where data is limited and input domains are specialized, tried-and-true CNN-based models like YOLOv5l may still have the edge over newer transformer-based approaches, highlighting their ongoing relevance in real-world applications.

# A Configuration Files

## A.1 YOLOv5

```
1      # YOLOv5 Custom Dataset Configuration
2      path: ./yolo_dataset # dataset root dir
3      train: images/train/grayscale_wholeImage
4      val: images/test/grayscale_wholeImage
5
6      # Classes
7      nc: 4
8      names: ['infant_seat', 'child_seat', 'person', 'object']
9
10     # Class weights
11     class_weights: {1: 1.5202492211838006, 3: 2.585430463576159, 2:
12         1.0, 0: 2.0699893955461293}
```

## A.2 Swin Transformer

```
1      auto_scale_lr = dict(base_batch_size=16, enable=False)
2      custom_imports = dict(
3          allow_failed_imports=False,
4          imports=[
5              'mmdet.evaluation.metrics.coco_metric',
6          ])
7      data_root = '/home/abdulbari/uni/seminar/swin_transformer/
8          transformer_dataset/'
8      dataset_type = 'CocoDataset'
9      default_hooks = dict(
10          checkpoint=dict(
11              interval=5, max_keep_ckpts=3, save_best='auto', type='
12                  CheckpointHook'),
12          logger=dict(interval=50, type='LoggerHook'),
13          param_scheduler=dict(type='ParamSchedulerHook'),
14          sampler_seed=dict(type='DistSamplerSeedHook'),
15          timer=dict(type='IterTimerHook'),
16          visualization=dict(type='DetVisualizationHook'))
17      default_scope = 'mmdet'
18      env_cfg = dict(
19          cudnn_benchmark=False,
20          dist_cfg=dict(backend='nccl'),
21          mp_cfg=dict(mp_start_method='fork', opencv_num_threads=0))
22      launcher = 'none'
23      load_from = None
24      log_level = 'INFO'
```

```

25 log_processor = dict(by_epoch=True, type='LogProcessor', window_size
26 =50)
27 metainfo = dict(
28     classes=(
29         'infant seat',
30         'child seat',
31         'person',
32         'everyday object',
33     ))
34 model = dict(
35     backbone=dict(
36         attn_drop_rate=0.0,
37         convert_weights=True,
38         depths=[
39             2,
40             2,
41             6,
42             2,
43         ],
44         drop_path_rate=0.2,
45         drop_rate=0.0,
46         embed_dims=96,
47         init_cfg=dict(
48             checkpoint=
49                 'https://github.com/SwinTransformer/storage/releases/
50                     download/v1.0.0/swin_tiny_patch4_window7_224.pth',
51             type='Pretrained'),
52         mlp_ratio=4,
53         num_heads=[
54             3,
55             6,
56             12,
57             24,
58         ],
59         out_indices=(
60             1,
61             2,
62             3,
63         ),
64         patch_norm=True,
65         qk_scale=None,
66         qkv_bias=True,
67         type='SwinTransformer',
68         window_size=7,
69         with_cp=False),
70     bbox_head=dict(
71         anchor_generator=dict(
72             octave_base_scale=4,
73             ratios=[
```

```

72                 0.5,
73                 1.0,
74                 2.0,
75             ],
76             scales_per_octave=3,
77             strides=[
78                 8,
79                 16,
80                 32,
81                 64,
82                 128,
83             ],
84             type='AnchorGenerator'),
85             bbox_coder=dict(
86                 target_means=[
87                     0.0,
88                     0.0,
89                     0.0,
90                     0.0,
91                 ],
92                 target_stds=[
93                     1.0,
94                     1.0,
95                     1.0,
96                     1.0,
97                 ],
98                 type='DeltaXYWHBBoxCoder'),
99             feat_channels=256,
100            in_channels=256,
101            loss_bbox=dict(loss_weight=2.0, type='L1Loss'),
102            loss_cls=dict(
103                alpha=0.25,
104                gamma=2.0,
105                loss_weight=1.0,
106                type='FocalLoss',
107                use_sigmoid=True),
108            num_classes=4,
109            stacked_convs=4,
110            type='RetinaHead'),
111            data_preprocessor=dict(
112                bgr_to_rgb=True,
113                mean=[
114                    123.675,
115                    116.28,
116                    103.53,
117                ],
118                pad_size_divisor=32,
119                std=[
120                    58.395,

```

```

121             57.12,
122             57.375,
123         ] ,
124         type='DetDataPreprocessor') ,
125     neck=dict(
126         add_extra_convs='on_input',
127         in_channels=[
128             192,
129             384,
130             768,
131         ] ,
132         num_outs=5,
133         out_channels=256,
134         start_level=0,
135         type='FPN'),
136     test_cfg=dict(
137         max_per_img=100,
138         min_bbox_size=0,
139         nms=dict(iou_threshold=0.5, type='nms'),
140         nms_pre=1000,
141         score_thr=0.05),
142     train_cfg=dict(
143         allowed_border=-1,
144         assigner=dict(
145             ignore_iou_thr=-1,
146             min_pos_iou=0,
147             neg_iou_thr=0.4,
148             pos_iou_thr=0.5,
149             type='MaxIoUAssigner'),
150         debug=False,
151         pos_weight=-1,
152         sampler=dict(type='PseudoSampler')),
153         type='RetinaNet')
154 optim_wrapper = dict(
155     clip_grad=dict(max_norm=35, norm_type=2),
156     optimizer=dict(lr=0.005, momentum=0.9, type='SGD', weight_decay
157 =0.0001),
158     type='OptimWrapper')
159 param_scheduler = [
160     dict(begin=0, by_epoch=False, end=500, start_factor=0.1, type='
161     LinearLR'),
162     dict(
163         begin=0,
164         by_epoch=True,
165         end=50,
166         eta_min=1e-06,
167         type='CosineAnnealingLR'),
168 ]
169 pretrained = 'https://github.com/SwinTransformer/storage/releases/

```

```

        download/v1.0.0/swin_tiny_patch4_window7_224.pth'
168 resume = False
169 test_cfg = dict(type='TestLoop')
170 test_dataloader = dict(
171     batch_size=16,
172     dataset=dict(
173         ann_file='annotations/instances_test.json',
174         data_prefix=dict(img='test/grayscale_wholeImage'),
175         data_root=
176             '/home/abdulbari/uni/seminar/swin_transformer/
177                 transformer_dataset/',
178         metainfo=dict(
179             classes=(
180                 'infant seat',
181                 'child seat',
182                 'person',
183                 'everyday object',
184             )),
185         pipeline=[
186             dict(type='LoadImageFromFile'),
187             dict(keep_ratio=True, scale=(
188                 448,
189                 448,
190             ), type='Resize'),
191             dict(type='LoadAnnotations', with_bbox=True),
192             dict(size_divisor=32, type='Pad'),
193             dict(
194                 meta_keys=(
195                     'img_id',
196                     'img_path',
197                     'ori_shape',
198                     'img_shape',
199                     'scale_factor',
200                 ),
201                 type='PackDetInputs'),
202             ],
203             test_mode=True,
204             type='CocoDataset'),
205             drop_last=False,
206             num_workers=4,
207             persistent_workers=True,
208             sampler=dict(shuffle=False, type='DefaultSampler'))
209 test_evaluator = dict(
210     ann_file=
211         '/home/abdulbari/uni/seminar/swin_transformer/
212             transformer_dataset/annotations/instances_test.json',
213     format_only=False,
214     metric='bbox',
215     type='CocoMetric')

```

```

214 test_pipeline = [
215     dict(type='LoadImageFromFile'),
216     dict(keep_ratio=True, scale=(
217         448,
218         448,
219     ), type='Resize'),
220     dict(type='LoadAnnotations', with_bbox=True),
221     dict(size_divisor=32, type='Pad'),
222     dict(
223         meta_keys=(
224             'img_id',
225             'img_path',
226             'ori_shape',
227             'img_shape',
228             'scale_factor',
229         ),
230         type='PackDetInputs'),
231 ]
232 train_cfg = dict(max_epochs=50, type='EpochBasedTrainLoop',
233 val_interval=5)
234 train_dataloader = dict(
235     batch_sampler=dict(type='AspectRatioBatchSampler'),
236     batch_size=16,
237     dataset=dict(
238         ann_file='annotations/instances_train.json',
239         data_prefix=dict(img='train/grayscale_wholeImage'),
240         data_root=
241             '/home/abdulbari/uni/seminar/swin_transformer/
242                 transformer_dataset/',
243         filter_cfg=dict(filter_empty_gt=True, min_size=32),
244         metainfo=dict(
245             classes=(
246                 'infant seat',
247                 'child seat',
248                 'person',
249                 'everyday object',
250             )),
251         pipeline=[
252             dict(type='LoadImageFromFile'),
253             dict(type='LoadAnnotations', with_bbox=True),
254             dict(keep_ratio=True, scale=(
255                 448,
256                 448,
257             ), type='Resize'),
258             dict(prob=0.5, type='RandomFlip'),
259             dict(
260                 border=(
261                     -32,
262                     -32,

```

```

261             ),
262             max_rotate_degree=5,
263             max_shear_degree=0,
264             scaling_ratio_range=(
265                 0.8,
266                 1.2,
267             ),
268             type='RandomAffine'),
269         dict(
270             brightness_delta=32,
271             contrast_range=(
272                 0.8,
273                 1.2,
274             ),
275             hue_delta=10,
276             saturation_range=(
277                 0.8,
278                 1.2,
279             ),
280             type='PhotoMetricDistortion'),
281         dict(size_divisor=32, type='Pad'),
282         dict(type='PackDetInputs'),
283     ],
284     type='CocoDataset'),
285     num_workers=4,
286     persistent_workers=True,
287     sampler=dict(shuffle=True, type='DefaultSampler'))
288 train_pipeline = [
289     dict(type='LoadImageFromFile'),
290     dict(type='LoadAnnotations', with_bbox=True),
291     dict(keep_ratio=True, scale=(
292         448,
293         448,
294     ), type='Resize'),
295     dict(prob=0.5, type='RandomFlip'),
296     dict(
297         border=(
298             -32,
299             -32,
300         ),
301         max_rotate_degree=5,
302         max_shear_degree=0,
303         scaling_ratio_range=(
304             0.8,
305             1.2,
306         ),
307         type='RandomAffine'),
308     dict(
309         brightness_delta=32,

```

```

310     contrast_range=(
311         0.8,
312         1.2,
313     ),
314     hue_delta=10,
315     saturation_range=(
316         0.8,
317         1.2,
318     ),
319     type='PhotoMetricDistortion'),
320     dict(size_divisor=32, type='Pad'),
321     dict(type='PackDetInputs'),
322 ]
323 val_cfg = dict(type='ValLoop')
324 val_dataloader = dict(
325     batch_size=16,
326     dataset=dict(
327         ann_file='annotations/instances_test.json',
328         data_prefix=dict(img='test/grayscale_wholeImage'),
329         data_root=
330             '/home/abdulbari/uni/seminar/swin_transformer/
331                 transformer_dataset/',
332         metainfo=dict(
333             classes=(
334                 'infant seat',
335                 'child seat',
336                 'person',
337                 'everyday object',
338             )),
339         pipeline=[
340             dict(type='LoadImageFromFile'),
341             dict(keep_ratio=True, scale=(
342                 448,
343                 448,
344             ), type='Resize'),
345             dict(type='LoadAnnotations', with_bbox=True),
346             dict(size_divisor=32, type='Pad'),
347             dict(
348                 meta_keys=(
349                     'img_id',
350                     'img_path',
351                     'ori_shape',
352                     'img_shape',
353                     'scale_factor',
354                 ),
355                 type='PackDetInputs'),
356             ],
357             test_mode=True,
358             type='CocoDataset'),

```

```

358     drop_last=False,
359     num_workers=4,
360     persistent_workers=True,
361     sampler=dict(shuffle=False, type='DefaultSampler'))
362 val_evaluator = dict(
363     ann_file=
364     '/home/abdulbari/uni/seminar/swin_transformer/
365         transformer_dataset/annotations/instances_test.json',
366     format_only=False,
367     metric='bbox',
368     type='CocoMetric')
369 vis_backends = [
370     dict(type='LocalVisBackend'),
371 ]
372 visualizer = dict(
373     name='visualizer',
374     type='DetLocalVisualizer',
375     vis_backends=[
376         dict(type='LocalVisBackend'),
377     ])
378 work_dir = 'work_dirs/final'

```

## B Inference Code

### B.1 YOLOv5l

```

1 import os
2 import torch
3 from PIL import Image
4
5 model = torch.hub.load('ultralytics/yolov5', 'custom', path='yolov5/
6     runs/train/exp11/weights/best.pt', force_reload=True)
7
8 image = os.path.join('yolo_dataset', 'images', 'test', '
9     grayscale_wholeImage', 'aclass_test_imageID_489_GT_5_0_2.png')
10
11 result = model(image)
12 result.print()
13
14 Image.fromarray(result.render()[0]).show()

```

### B.2 Swin Transformer

```

1 from mmdet.apis import init_detector, inference_detector

```

```

2 import mmcv
3 import matplotlib.pyplot as plt
4 import numpy as np
5 from matplotlib.patches import Rectangle
6
7 config_file = 'mmdetection/configs/swin/retinanet_swin-t-p4-
    w7_fpn_1x_coco.py'
8 checkpoint_file = 'work_dirs/final1/epoch_50.pth'
9 device = 'cuda:0'
10
11 model = init_detector(config_file, checkpoint_file, device=device)
12
13 class_names = ['infant seat', 'child seat', 'person', 'everyday
    object']
14
15 img = 'transformer_dataset/test/grayscale_wholeImage/
    aclass_test_imageID_489_GT_5_0_2.png'
16
17 result = inference_detector(model, img)
18
19 # Visualize results
20 def visualize_detection(img_path, result, score_threshold=0.3):
21     img = mmcv.imread(img_path)
22     img = mmcv.bgr2rgb(img)
23
24     plt.figure(figsize=(12, 8))
25     plt.imshow(img)
26
27     # Access predictions from the DetDataSample object
28     pred_instances = result.pred_instances
29
30     # Get bounding boxes, scores, and labels
31     bboxes = pred_instances.bboxes.cpu().numpy()
32     scores = pred_instances.scores.cpu().numpy()
33     labels = pred_instances.labels.cpu().numpy()
34
35     # Draw detections
36     for i in range(len(bboxes)):
37         if scores[i] >= score_threshold:
38             bbox = bboxes[i]
39             label = labels[i]
40             score = scores[i]
41
42             # Create a rectangle patch
43             x1, y1, x2, y2 = bbox
44             width = x2 - x1
45             height = y2 - y1
46             rect = Rectangle((x1, y1), width, height, linewidth=2,
47                             edgecolor='r', facecolor='none')

```

```

48         plt.gca().add_patch(rect)
49
50     # Add label
51     plt.text(x1, y1-5, f'{class_names[label]}: {score:.2f}', 
52             color='red', fontsize=12, bbox=dict(facecolor='
53                                         white', alpha=0.7))
54
55     plt.axis('off')
56     plt.tight_layout()
57     plt.show()
58
59 # Print detection results
60 def print_detections(result, score_threshold=0.3):
61     # Access predictions
62     pred_instances = result.pred_instances
63
64     # Get bounding boxes, scores, and labels
65     bboxes = pred_instances.bboxes.cpu().numpy()
66     scores = pred_instances.scores.cpu().numpy()
67     labels = pred_instances.labels.cpu().numpy()
68
69     # Create a dictionary to group detections by class
70     class_detections = {}
71     for i in range(len(bboxes)):
72         if scores[i] >= score_threshold:
73             label = int(labels[i])
74             if label not in class_detections:
75                 class_detections[label] = []
76             class_detections[label].append((bboxes[i], scores[i]))
77
78     # Print results for each class
79     for class_id in sorted(class_detections.keys()):
80         print(f"\nDetections for class '{class_names[class_id]}':")
81         for i, (bbox, score) in enumerate(class_detections[class_id]):
82             x1, y1, x2, y2 = bbox
83             print(f"  Detection {i+1}: confidence={score:.2f}, bbox
84                   =[{x1:.1f}, {y1:.1f}, {x2:.1f}, {y2:.1f}]")
85
86     # Call the visualization function
87     visualize_detection(img, result, score_threshold=0.6)
88
89     # Print detection details
90     print_detections(result, score_threshold=0.6)

```

## References

- [1] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 779–788, 2016.
- [2] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European conference on computer vision*, pp. 213–229, Springer, 2020.
- [3] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- [4] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- [5] S. D. D. Cruz, O. Wasenmuller, H.-P. Beise, T. Stifter, and D. Stricker, “Sviro: Synthetic vehicle interior rear seat occupancy dataset and benchmark,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 973–982, 2020.
- [6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, vol. 25, 2012.
- [8] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, vol. 28, 2015.
- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds.), vol. 30, Curran Associates, Inc., 2017.
- [10] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” 2014.
- [11] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, *et al.*, “A survey on vision transformer,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 1, pp. 87–110, 2022.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

- [13] J. Redmon and A. Farhadi, “Yolov3: An incremental improvement,” *arXiv preprint arXiv:1804.02767*, 2018.
- [14] G. Jocher *et al.*, “Yolov5 - object detection models.” <https://github.com/ultralytics/yolov5>, 2020.
- [15] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*, pp. 740–755, Springer, 2014.
- [16] M. Contributors, “Mmdetection: Openmmlab detection toolbox and benchmark.” <https://github.com/open-mmlab/mmdetection>, 2018.
- [17] R. Padilla, S. L. Netto, and E. A. Da Silva, “A survey on performance metrics for object-detection algorithms,” pp. 237–242, 2020.
- [18] C. Shorten and T. M. Khoshgoftaar, “A survey on image data augmentation for deep learning,” vol. 6, pp. 1–48, SpringerOpen, 2019.