# Data Augmentation for Object Detection: A Comparative Study of Transformation Techniques

Aqsa Mohsin, Nasrul Huda

Neural Networks Course Project

10 July 2025

### Abstract

In recent years, object detection has become an essential component of many computer vision applications. However, the success of detection models heavily depends on the availability of large, diverse, and annotated datasets. This research explores the efficacy of data augmentation techniques to enhance model performance in scenarios where training data is limited. Using the Pascal VOC 2012 dataset, we conducted controlled experiments by training a Faster R-CNN model with and without various augmentation strategies. Due to hardware limitations, we restricted our study to a single object class with the smallest sample size: *cow*. Through this experiment, we aim to highlight the importance of augmentation in achieving better generalization with scarce data and provide empirical evidence on the comparative performance of different augmentation methods.

## 1 Introduction

Object detection is a core task in computer vision, with applications ranging from autonomous vehicles to medical imaging and surveillance systems. Deep learning models like Faster R-CNN and YOLO have significantly advanced detection accuracy, primarily fueled by large-scale annotated datasets such as COCO and Pascal VOC. However, these models typically require thousands of labeled images per category to generalize effectively.

In practice, collecting such extensive datasets is often infeasible due to limitations in time, resources, or domain-specific constraints. A common and effective strategy to mitigate these constraints is data augmentation. By applying a variety of transformations to existing images, augmentation helps simulate dataset diversity without requiring additional data collection.

This research focuses on evaluating the impact of data augmentation in the context of limited data availability. Initially, we intended to train our model using all classes in the Pascal VOC dataset. However, repeated system crashes due to the volume of data forced us to reconsider our strategy. We analyzed the class-wise distribution of training images and selected *cow* as the object category with the least training images (151 for training and 152 for validation). This reduction allowed us to run extensive experiments within the

constrained computational environment while still providing meaningful insights into the effectiveness of various augmentation techniques.

# 2 Methodology

Our objective was to compare the effectiveness of different augmentation strategies on a standard object detection model using minimal training data. The pipeline consisted of the following stages:

## 2.1 Dataset Selection

We used the Pascal VOC 2012 dataset and filtered it to include only the *cow* class. This decision was driven by the need to minimize memory consumption and avoid runtime failures in Google Colab. Figure 1 demonstrates that the *cow* class has the smallest sample size among all Pascal VOC classes, justifying our selection for this constrained study. Figure 2 illustrates the quality and consistency of bounding box annotations in our filtered dataset.
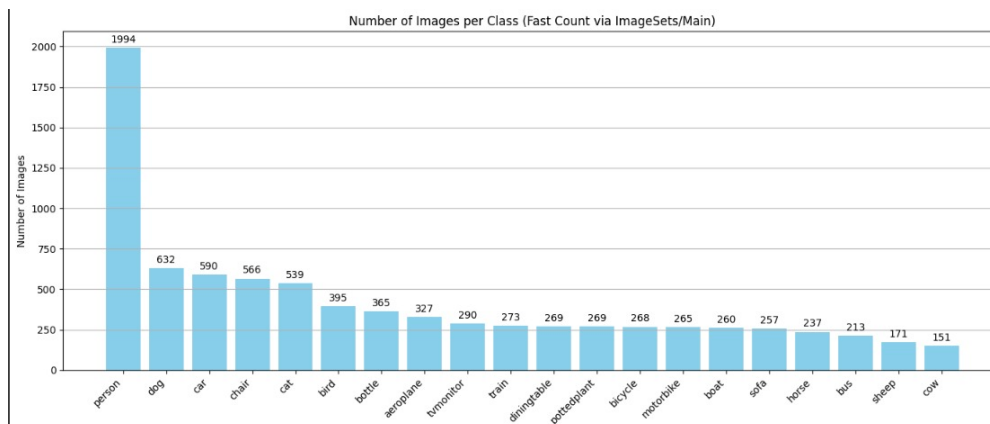


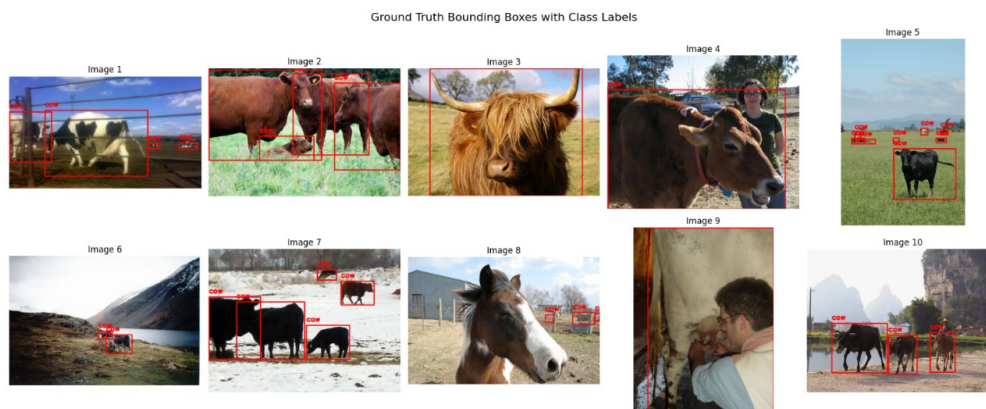Figure 1: Number of images per class in the Pascal VOC 2012 dataset



Figure 2: Ground truth boxes for the *cow* class in the Pascal VOC 2012 dataset

## 2.2 Model Architecture

We employed the `fasterrcnn_mobilenet_v3_large_320_fpn` model pre-trained on COCO. This variant of Faster R-CNN provides a lighter and more memory-efficient architecture suited for low-resource environments.

## 2.3 Training Setup

All models were trained for 30 epochs. Due to the small size of our dataset and memory limitations in the Colab environment, we used a batch size of 2. The learning rate was set to 0.001 to allow gradual convergence without overfitting. Since our dataset was very small, these parameters were chosen to ensure a balance between learning and stability.

## 2.4 Augmentation Techniques

To evaluate the impact of different augmentation strategies on object detection, five distinct augmentation setups were explored. The baseline model was trained on the original dataset without any augmentation, serving as a reference point for comparison. In the horizontal flip experiment, each image had a 50% probability of being flipped along the vertical axis, aiming to improve the model's ability to detect objects regardless of orientation. The random crop setup involved selecting random regions from the image while ensuring that the objects remained within the cropped area, which encouraged the model to learn from partial views and different object scales. The rotation experiment applied small-angle rotations to simulate various viewing perspectives, enhancing spatial generalization. Finally, a combined augmentation strategy was introduced, where a sequence of transformations—including horizontal flip, cropping, and rotation—was applied to each image to maximize diversity. All transformations were implemented using the Albumentations library, and each setup was treated as an independent experiment resulting in a separately trained model.

## 2.5 Evaluation Metrics

To evaluate the effectiveness of different augmentation strategies on object detection, we employed a set of quantitative metrics. First, the Average Loss per Epoch was recorded throughout training to monitor convergence behavior and stability. This metric reflects how well the model fits the training data over time.

Additionally, we calculated Precision, defined as the proportion of predicted bounding boxes whose Intersection over Union (IoU) with ground truth exceeds a threshold of 0.5. This metric measures the model's accuracy in predicting correct object locations while avoiding false positives.

The Average IoU across all predictions further served as a global measure of localization accuracy, indicating how closely predicted bounding boxes matched the true object boundaries.

# 3   Results and Analysis

Across all models, we observed distinct differences in how augmentations affected learning dynamics and detection performance. As shown in Figure 3, the training loss curves reveal distinct convergence patterns for each augmentation strategy.
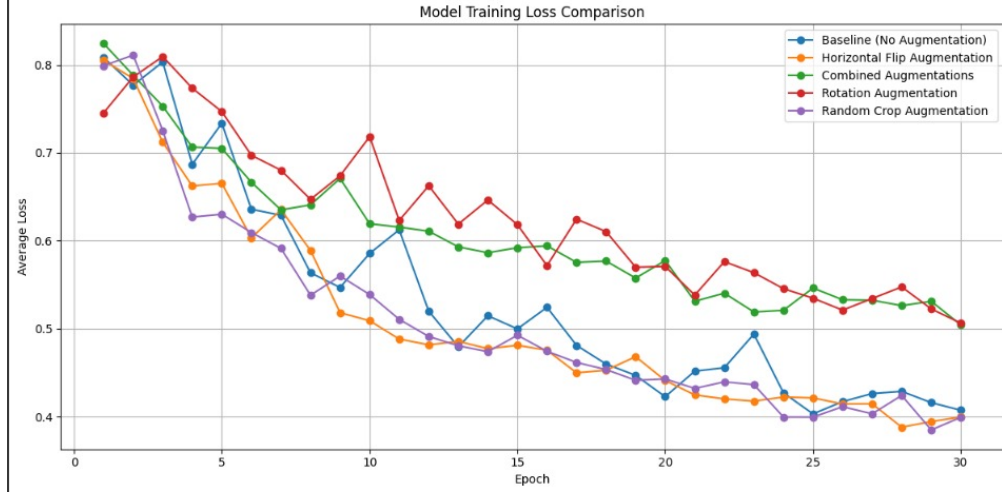


Figure 3: Training loss comparison for different augmentation strategies over 30 epochs.

The baseline model achieved an average IoU of approximately 0.38 and average loss converging to 0.407. The Horizontal Flip augmentation improved both IoU and precision, with an average IoU of 0.406 and loss of 0.399.

The Crop Augmentation strategy achieved comparable results, with similar IoU and slightly better loss. The Rotation Augmentation showed inconsistent performance, while the Combined Augmentation plateaued in IoU and resulted in higher loss (0.504). Figure 4 visually confirms these quantitative results, clearly showing that Flip and Crop augmentations achieve the highest average IoU values.

| Model | Precision | Avg IoU | Avg Loss |
|---|---|---|---|
| Baseline | 0.433 | 0.381 | 0.407 |
| Flip Augmentation | 0.446 | 0.406 | 0.399 |
| Mixed Augmentation | 0.416 | 0.380 | 0.504 |
| Rotation Augmentation | 0.442 | 0.397 | 0.506 |
| Crop Augmentation | 0.435 | 0.406 | 0.399 |

Table 1: Performance metrics for different augmentation strategies

The comprehensive performance metrics are summarized in Table 1, which quantifies the trade-offs between different augmentation approaches.
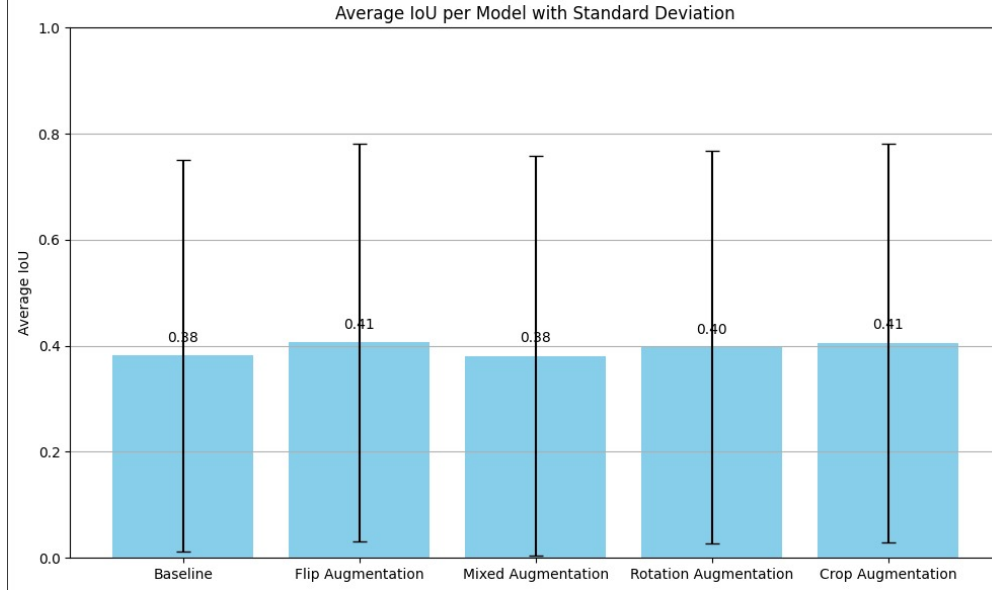
Figure 4: Average IoU comparison across different augmentation strategies.

# 4 Discussion

The results suggest that simple augmentations like horizontal flips and random cropping can significantly enhance model performance in data-scarce scenarios. Both strategies achieved identical IoU improvements (0.406), while maintaining lower training losses (0.399) compared to more complex approaches. Mixed augmentations introduced excessive variability, resulting in worse performance (0.380 IoU) and higher loss (0.504), while rotation showed modest improvements (0.397 IoU) but with elevated loss (0.506).

# 5 Conclusion

This study provides empirical evidence for the effectiveness of targeted data augmentation in object detection under severe data constraints. Our systematic comparison of five augmentation strategies reveals several key insights:

The quantitative results demonstrate that simple geometric transformations yield the most significant improvements. Horizontal Flip and Crop augmentations both achieved an average IoU of 0.406, representing a 6.6% improvement over the baseline (0.381). These same strategies also achieved the lowest training losses (0.399), indicating better model convergence.

Conversely, our findings reveal that complex augmentation combinations can be counterproductive in data-scarce scenarios. The Mixed Augmentation approach not only failed to improve upon individual techniques (0.380 IoU) but also resulted in the highest training loss (0.504), suggesting that excessive transformations may introduce noise that hampers learning with limited data.

These results have practical implications for practitioners working with limited datasets,

suggesting that focused, semantically-preserving augmentations are more effective than complex transformation pipelines.

# 6 Future Work

Building upon the findings of this research, several promising directions can be pursued to further understand and enhance the role of data augmentation in object detection. Firstly, extending the study to include other object classes from the Pascal VOC dataset can help verify whether the observed effects of augmentations generalize across categories with different visual characteristics and object complexities.

In addition, experimenting with more diverse augmentation techniques, such as color jittering, brightness/contrast adjustments, and synthetic image generation, may introduce useful variations that enhance the model's robustness. Modern augmentation policies like AutoAugment and RandAugment, which automatically search for optimal transformation strategies, could also be integrated to replace the manual design of augmentation pipelines.

Future work should also aim to scale this study to larger datasets and utilize more powerful hardware resources. This would enable a thorough evaluation of the augmentation methods under realistic conditions, reflecting the demands of production-grade detection systems.

Furthermore, analyzing the training dynamics of each augmentation type in detail, such as how quickly models converge and the stability of their loss functions, can provide insights into optimal training durations and early stopping criteria.

Collectively, these extensions would build a more comprehensive understanding of the interplay between data augmentation and object detection performance, especially in resource-constrained environments.

# Acknowledgments

# References

[1] Lin, T. Y., et al. *Microsoft COCO: Common Objects in Context.* ECCV, 2014.

[2] Everingham, M., et al. *The Pascal Visual Object Classes (VOC) Challenge.* IJCV, 2010.

[3] Ren, S., et al. *Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks.* NeurIPS, 2015.

[4] Buslaev, A., et al. *Albumentations: Fast and Flexible Image Augmentations.* Information, 2020.

[5] Paszke, A., et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library.* NeurIPS, 2019.