to moving elements. This ensures that the model not only pinpoints objects but also gauges their movement in real-time. By marrying YOLO's innate processing speed with our motion-focused enhancements, we anticipate a refined and efficient detection of moving entities across a variety of settings. To explain the proposed model that is illustrated in Figure 3, the detailed steps of the proposed detection technique are as follows:

### A. PREPARING THE VIDEO

- Decoding the Video: Isolate each frame from the video to prepare them for further steps.
- Adjusting Frame Size: Since YOLO architectures require certain input sizes like 416*416 or 608*608, modify the frame dimensions to align with YOLOv8's specific requirements.
- Adjusting Color Values: Transform pixel values to a range between [0,1] or [-1,1], based on the criteria set by the pre-trained model.

### B. REMOVING BACKGROUND

- Setting Up the Background Model: Determine a baseline background image by averaging or taking the median of a set number of frames.
- Differentiating Frames: Highlight moving objects by deducting the current frame from the background model.
- Enhance clarity by applying a binary threshold, marking significant deviations from the background as "white" and everything else as "black."
- Reducing Disturbances: Implement morphological procedures, like erosion and dilation, to eliminate distractions and bridge minor gaps.

In summary, these modifications collectively enhance the YOLO-V8 model's performance in detecting moving objects by bolstering its robustness, feature representation capabilities, and detection accuracy. The challenges introduced by moving objects, such as motion blur, alterations in appearance, and occlusion, are better addressed through these enhancements. It's important to emphasize that the actual impact of these modifications depends on factors like the quality and diversity of your training data, specific implementation details, and hyperparameter tuning. Rigorous experimentation and evaluation using a representative dataset containing moving objects are vital for accurately quantifying performance improvements.

Let's delve into the provided enhancement, named *EnhancedYOLOv8*, and dissect its components in detail.

1) **Advanced Augmentation (advanced_augmentation):**
   - Purpose: Data augmentation is crucial for object detection models to generalize well to various real-world scenarios. The better and more varied your augmentations, the more robust the model becomes.
   - Enhancement: The term *advanced_augmentation* suggests that apart from typical augmentations like simple rotations, scaling, and translations, there

---

**Algorithm 1** Enhanced Object Detection Model

**Procedure Advanced Data Augmentation**
    **Input:** Original dataset
    **Output:** Augmented dataset
    **Procedure:**
        Introduce the `advanced_augmentation` function.
        Apply advanced data augmentation techniques.
        **return** Augmented dataset

**Procedure Advanced Backbone Network**
    **Input:** Input images
    **Output:** Predicted bounding boxes
    **Procedure:**
        Employ `EnhancedYOLOv8Architecture` function to define Enhanced YOLO-V8.
        Use a sophisticated network like ResNet to capture intricate details.
        **return** Predicted bounding boxes

**Procedure Fine-Tuning**
    **Input:** Pretrained model, Augmented dataset
    **Output:** Fine-tuned model
    **Procedure:**
        Perform fine-tuning with `neural_network.fine_tune` on the augmented data.
        Adapt the model to the specific dataset and task.
        **return** Fine-tuned model.

**Procedure Advanced Post-Processing**
    **Input:** Predicted bounding boxes.
    **Output:** Refined bounding boxes.
    **Procedure:**
        Apply `advanced_postprocess` for techniques like Soft-NMS.
        Soft-NMS improves object detection accuracy
        **return** Refined bounding boxes

---

might be more sophisticated augmentations involved. Examples include:
   - Random crops and zooms simulate objects at different distances.
   - Color jittering for varying lighting conditions.
   - Implementing techniques like MixUp or CutMix to blend images.
   - Temporal augmentations if sequences of images or videos are involved.

2) **Enhanced Neural Network Architecture (EnhancedYOLOv8Architecture):**
   - Purpose: The backbone and architecture of an object detection model determine its performance and complexity.

- Enhancement: *EnhancedYOLOv8Architecture* implies modifications or advancements over the traditional YOLO architecture. Speculative improvements might include:
  - Utilizing more advanced backbones like EfficientNet or Vision Transformers.
  - Implementing more extensive Feature Pyramid Networks (FPNs) for detecting objects at different scales.
  - Incorporating attention mechanisms to focus on critical parts of an image.
  - Optimizing the architecture for specific hardware for faster inference.

3) **Loading Pretrained Weights:**
- Purpose: Pretrained weights, usually on datasets like ImageNet, provide a solid initialization point and can help the model converge faster.
- Enhancement: The ability to fine-tune these weights on augmented data suggests that the model can be optimized for specific scenarios or datasets related to moving object detection.

4) **Advanced Post-Processing (advanced_postprocess):**
- Purpose: Raw detections from a neural network usually require post-processing to be useful. This can involve Non-maximum Suppression (NMS), thresholding, or filtering out low-confidence detections.
- Enhancement: The term *advanced_postprocess* suggests steps beyond typical post-processing. For instance:
  - Soft-NMS: A variant of NMS that doesn't entirely suppress neighboring bounding boxes but down-weights them.
  - Incorporation of tracking algorithms for moving objects, ensuring consistent object identities across frames.
  - Implementing additional heuristics or logic based on the specific domain (e.g., prioritizing larger bounding boxes for vehicle detection).

5) **Overall Implications:**
- Computation: Enhanced features, especially in architecture, might add computational overhead. This can be offset by the potential increase in Frame Per Second (FPS), as suggested by the synthetic data.
- Performance: The use of advanced augmentations and post-processing, combined with sophisticated architecture, should ideally lead to better detection and tracking performance, especially in challenging scenarios.
- Generalization: Advanced augmentations can help the model generalize well to real-world conditions, making it robust against various environmental factors.

In summary, *EnhancedYOLOv8* appears to be a more advanced and optimized version of a hypothetical YOLOv8, tailored specifically for better performance in moving object detection scenarios. Its features aim to boost accuracy, maintain real-time processing capabilities, and ensure robustness in diverse conditions.

---

**Algorithm 2** Enhanced YOLO-V8 Object Detection

**function** EnhancedYOLOv8(image)
  augmented_image ← advanced_augmentation(preprocess(image))
  neural_network ← EnhancedYOLOv8Architecture()
  **if** pretrained_weights_exist **then**
    neural_network.load_weights(pretrained_weights)
    neural_network.fine_tune(augmented_image).
  **end if**
  detection_results ← neural_network.forward (augmented_image)
  detections ← advanced_postprocess(detection_results).
  **return** detections.
**end function**

---

**Algorithm 3** EnhancedYOLOv8Architecture

**Ensure:** architecture: the neural network architecture for Enhanced YOLOv8
**function** EnhancedYOLOv8Architecture
  architecture ← DefineEnhancedArchitecture()
  **return** architecture
**end function**

---

**Algorithm 4** DefineEnhancedArchitecture

**Ensure:** architecture: the defined architecture with specific layers of Enhanced YOLOv8
**function** DefineEnhancedArchitecture
  architecture ← BuildEnhancedNetworkLayers()
  **return** architecture
**end function**

---

**Algorithm 5** BuildEnhancedNetworkLayers

**Ensure:** layers: the neural network layers with advanced features
**function** BuildEnhancedNetworkLayers
  Initialize empty list *layers*
  Add advanced convolutional layers, batch normalization, activation functions, etc., to *layers*
  **return** layers
**end function**

---

### C. DISCUSSION ABOUT THE ENHANCED YOLO
#### 1) ENHANCED DATA AUGMENTATION TECHNIQUES FOR DETECTING MOVING OBJECTS

The study introduces a suite of advanced data augmentation techniques tailored to enhance the model's capability to detect moving objects. These techniques include:

**Algorithm 7** YOLOv8 Object Detection Algorithm

**function** YOLOv8(image)
    preprocessed_image ← preprocess(image)
    neural_network ← YOLOv8Architecture
    **if** pretrained_weights_exist **then**
        neural_network.load_weights(pretrained_weights)
    **end if**
    detection_results                  ←           neural_network.forward(preprocessed_image)
    detections ← postprocess(detection_results)
    **return** detections
**end function**
**function** preprocess(image)
    processed_image ← apply_preprocessing(image)
    **return** processed_image
**end function**
**function** YOLOv8Architecture
    architecture ← DefineArchitecture
    **return** architecture
**end function**
**function** DefineArchitecture
    architecture ← BuildNetworkLayers()
    **return** architecture
**end function**
**function** BuildNetworkLayers
    layers ← []             ▷ Initialize empty list of layers
    Add convolutional layers, batch normalization, activation functions, etc., to *layers*
    **return** layers
**end function**
**function** postprocess(detection_results)
    postprocessed_detections           ←  apply_postprocessing(detection_results)
    **return** postprocessed_detections
**end function**

### D. EVALUATING MODEL PERFORMANCE

- Referencing Actual Data: For model assessment, access the actual object locations for each frame from ground truth data.
- Derive IoU (Intersection over Union): Analyze the overlap between predicted and real object locations using the IoU metric for each identification.
- Deriving Precision and Recall: Given a specific confidence benchmark, compute both precision and recall values.
- Evaluate Using mAP (mean Average Precision): Determine precision for each category and compute their mean. mAP serves as a consolidated metric for assessing model accuracy across varying categories and confidence levels.

## IV. EXPERIMENT AND RESULTS

The KITTI, LASIESTA [67], PESMOD [68], and MOCS [69] datasets are utilized for the model training and evaluation processes discussed in this paper.

**TABLE 1.** labeled classes of KITTI dataset.

| Category | Description |
| --- | --- |
| Car | A typical, conventional motor vehicle. |
| Van | Various types of vehicles that are sized and shaped midway between a car and a truck. |
| Truck | The largest category of moving vehicles. |
| Pedestrian | Individuals who are walking or appear ready to walk. |
| Person (sitting) | Individuals who appear stationary within the scene, such as someone sitting on a park bench. |
| Cyclist | A person actively riding a bicycle. |
| Tram | A customary city tram used for public transportation. |
| Misc | Miscellaneous objects associated with vehicles, like trailers or Segways. |

The collection of the KITTI dataset was conducted using a specially equipped test vehicle that was outfitted with an array of imaging equipment, including both RGB and grayscale cameras, as well as a laser scanner. It also featured an inertial navigation system and varifocal lenses. The comprehensive dataset encapsulates the full spectrum of data acquired as the vehicle navigated through urban settings. For all the annotated categories, 3D tracklets are provided.

The collection of the KITTI dataset was conducted using a specially equipped test vehicle that was outfitted with an array of imaging equipment, including both RGB and grayscale cameras, as well as a laser scanner. It also featured an inertial navigation system and varifocal lenses. The comprehensive dataset encapsulates the full spectrum of data acquired as the vehicle navigated through urban settings. For all the annotated categories, 3D tracklets are provided. A description of the labeled classes available in Table 1.

The KITTI dataset encompasses a diverse range of urban environments. Detailed statistics on the tracklets are illustrated in Figure 4. As indicated by the left-side chart, 'Car' emerges as the most commonly tagged category, followed by 'Van.' Other classes, including 'Truck,' 'Pedestrian,' 'Cyclist,' and 'Misc,' display a similar frequency of labeling. Particularly scarce are the 'Tram' and 'Person (sitting)' categories, with the latter being notably challenging to locate within the dataset. The right-side chart details the distribution of tracklets across frames, revealing that a significant proportion of frames contain between two and six tracklets. For assessment purposes on the KITTI Vision Benchmark Suite, objects are categorized based on the level of difficulty in detecting them, which is determined by the degree of occlusion and truncation observed in the tracklet. The grading of difficulty spans 'Easy,' 'Moderate,' to 'Hard.' The 2D object detection benchmark page of KITTI provides a dataset with annotations for all mentioned categories. The benchmark's evaluation protocol involves testing over three difficulty levels for the 'Car' category, which requires a 70% overlap with the ground truth, and the 'Cyclist' and 'Pedestrian' categories, which necessitate a 50% overlap.

The LASIESTA dataset is an extensive and meticulously organized collection of 48 sequences, specifically designed