

University of Tripoli
Faculty of Engineering
Electrical and Electronic Engineering Department

EE569 – Deep Neural Networks
Fall 2024

Brief report for
Assignment 2
Computer Vision Task

Team members:

Mohamed Mustafa Kamil ID: 2180205565

Mohamed Abdulbaset Alshareef ID: 2180205640

Instructor: Dr. Nuri Ben Barka

1. Introduction

This report provides a detailed overview of the computer vision task performed using the LV-MHP-V2 dataset. The primary goal was to develop an instance segmentation model that detects, segments, and classifies individuals by gender using Detectron2. The report includes dataset processing, model training, evaluation, encountered challenges, and optimizations.

2. Objective

- Detect multiple people in an image.
- Generate segmentation masks for each detected person.
- Classify each detected individual by gender (i.e. male or female).

3. Dataset: LV-MHP-V2

3.1 Overview

The LV-MHP-V2 (Look into Person Multi-Human Parsing V2) dataset is a large-scale benchmark designed for human parsing, instance segmentation, and dense human-centric tasks. It contains densely packed images with multiple humans in real-world scenarios, making it a challenging dataset for segmentation and classification models.

The LV-MHP-V2 dataset contains densely packed human images with:

- Bounding box annotations.
- Parsing annotations (segmentation masks for body parts).
- No gender labels, which were manually annotated.

The dataset is structured as follows:

- Training set
 - Images
 - Parsing annotations
- Validation set
 - Images
 - Parsing annotations
- Test set
 - Images

3.3 Data Preprocessing

Parsing annotations were modified to remove specific body parts such as faces and hands. Custom scripts were implemented to refine segmentation masks and eliminate internal contours to ensure more accurate segmentation. Since the dataset did not include gender labels, manual annotations were added to classify individuals as male or female and the custom model was trained on them.

4. Phases of Workflow

Phase 1: Setup and Data Exploration

The necessary environment was configured by installing required libraries such as PyTorch, Detectron2, OpenCV, and Matplotlib. The dataset LV-MHP-V2 was downloaded and structured.

Environment Setup:

- Python 3.8+
- Detectron2
- OpenCV
- PyTorch
- Matplotlib
- Numpy

Phase 2: Object Detection

Before implementing instance segmentation, a baseline object detection experiment was conducted using pre-trained models from Detectron2 Model Zoo. The goal of this step was to evaluate the performance of existing object detection models on the LV-MHP-V2 dataset, providing a benchmark to compare against our final instance segmentation model.

Selecting a Pre-Trained Model:

To establish an object detection baseline, we explored two popular pre-trained object detection models from Detectron2 Model Zoo:

- **Faster R-CNN (R50-FPN):** A two-stage object detector known for its high accuracy.
- **RetinaNet (R50-FPN):** A one-stage detector optimized for speed while maintaining good detection performance.

The Faster R-CNN model was ultimately selected due to its better accuracy in detecting multiple persons in crowded scenes.

Dataset Preparation:

- The LV-MHP-V2 dataset provides bounding box annotations, which were used as ground truth for evaluating the model.
- The training and validation splits were loaded using the Detectron2 DatasetCatalog.

Model Configuration:

A Detectron2 configuration was setup to load the pre-trained Faster R-CNN model. The pre-trained weights from COCO detection models were loaded, and the model was evaluated on the validation set.

Evaluation:

The model was tested on a subset of validation images, and its predictions were compared against the ground truth bounding boxes.

Metric	Faster R-CNN (R50-FPN)
mAP (IoU=0.5)	76.2%
Recall	78.4%
Precision	74.6%

Table.1 Performance Metrics

The model performed well in detecting individual persons but struggled with densely packed crowds, where bounding boxes overlapped significantly.

Phase 3: Data Preparation for Instance Segmentation and Gender Classification

The LV-MHP-V2 dataset primarily focuses on multi-human parsing and instance segmentation, but it does not include gender labels. To enable gender classification alongside segmentation, a manual annotation process was required.

Gender Labeling Process:

Since no gender information was available in the dataset, a subset of images was manually annotated with labels ("male" or "female"). CVAT annotation tool was used to add these labels. These tools allow bounding box and attribute annotations to be efficiently assigned to each individual in an image.

To prevent bias in the model, a balanced dataset with approximately equal numbers of male and female objects was created, which was randomly sampled from the LV-MHP-v2 (500 image for training and 120 image for validation). Gender annotations were stored in COCO JSON format, making them compatible with the structure of the Detectron2 dataset.

Mask Generation from Parsing Annotations:

The original LV-MHP-V2 parsing annotations provide detailed body part segmentation masks, meaning each individual in an image is labeled with different regions such as head, hands, torso, legs, and accessories. A Python script was written to process the parsing annotations and convert them into instance segmentation masks.

The key modifications applied to the masks were:

- All body parts were included except for hands and face.
- The processed masks were saved in grayscale format, where the foreground (person) was white (255) and the background was black (0).

The segmentation masks were overlaid on the original images for visual inspection. Masks were color-coded: the segmentation of each detected person was highlighted, ensuring that the model could distinguish individuals correctly.

In instance segmentation preprocessing, one of the key challenges was removing internal contours (primarily faces) from segmentation masks. The problem arose because internal contours were disconnected from the outer body mask, leaving isolated regions within the segmentation masks. This required a technique to merge these internal contours with the outer boundary, ensuring a single connected mask.

To solve this, we implemented a shortest path algorithm to establish a connection between the inner contour (face) and the outer contour (body).

Concept of the Shortest Path Approach:

Instead of simply filling the internal contours (which could create inaccurate masks), a **thin connection** was created between the inner and outer contours, ensuring proper mask connectivity. The **key steps** were:

1. Extract Contours:

- Using `cv2.findContours`, we identified all contours in the segmentation mask.
- `cv2.RETR_CCOMP` was used to retrieve both external and internal contours while maintaining hierarchical relationships.

2. Determine the Nearest Connection:

- For each internal contour (identified via hierarchy), the closest point on the outer contour was computed.
- The Euclidean distance between every point on the inner and outer contour was measured.
- The pair of points with the minimum distance was selected as the optimal connection point.

3. Merge the Contours Using a Thin Connection:

- A thin black line (`cv2.line()`) was drawn between the selected inner and outer contour points.
- The thickness of the line was optimized to ensure proper merging while preventing excessive distortion of the segmentation mask.

4. Finalize the Mask:

- After merging the contours, `cv2.drawContours()` was used to fill the entire connected region, ensuring that the face (previously an internal contour) was now part of the background.

The algorithm is computationally intensive when multiple internal contours exist, as it calculates distances between all contour points. To optimize performance, we filtered small internal contours, ensuring only significant internal regions (e.g., faces) were processed.

The `thickness=3` parameter in `cv2.line()` was fine-tuned to balance merging accuracy and mask integrity.

Dataset Class Modification:

Since the original LV-MHP-V2 dataset did not match the standard Detectron2 format, modifications were required to integrate the processed dataset into Detectron2's training pipeline.

The dataset class was **modified** to load:

- Preprocessed images

- Generated segmentation masks
- Manually annotated gender labels

Each data entry was structured as follows:

- File name
- Image ID
- Height
- Width
- Annotations
 - Bounding boxes (bbox)
 - Instance segmentation masks
 - Gender labels (category_id) 0 for male, 1 for female

Handling Gender Labels in Training:

The gender labels were included as part of the dataset dictionary, enabling Detectron2 to learn gender classification from the annotated dataset.

The modified dataset was successfully registered in Detectron2, and the training pipeline was updated to include both instance segmentation and gender classification. The dataset was formatted correctly, allowing for seamless integration with Detectron2's training framework.

Before training, the dataset was visualized to validate the correctness of segmentation masks and gender annotations. Sample images were inspected to ensure that annotations aligned properly with the corresponding objects in the images. This step was crucial in identifying potential annotation errors and inconsistencies.

Phase 4: Instance Segmentation and Gender Classification Model Training

Once the dataset was prepared with cleaned segmentation masks and manually annotated gender labels, we proceeded with training the instance segmentation and gender classification model using Detectron2. The following steps were undertaken to complete this phase effectively.

Model Selection and Configuration:

To perform instance segmentation, we selected Mask R-CNN from the Detectron2 Model Zoo. This model was chosen due to its strong performance in object detection and segmentation tasks while maintaining computational efficiency.

The model was configured with the following settings:

- **Pretrained Weights:** A model pre-trained on COCO-InstanceSegmentation was used as a starting point.
- **Number of Classes:** The model was modified to detect and segment only two classes (male and female).
- **Bounding Box Head:** Outputs bounding boxes around detected individuals.
- **Segmentation Head:** Generates pixel-wise masks for each detected individual.

The following hyperparameters were fine-tuned for optimal performance:

- Learning rate: 0.00025
- Batch size: 512
- Number of iterations: 10,000
- Number of parallel data loading processes: 4
- Number of classes: 2

The model was registered with the training and validation datasets using Detectron2's DatasetCatalog.

Training the Model:

Once the model was configured, training was performed using the prepared dataset. The model was trained with:

- Training images: Processed images with segmentation masks.
- Segmentation masks: Generated masks excluding hands and faces.
- Gender labels: Assigned labels (male or female) for each individual.

The training pipeline followed these steps:

1. Load the dataset using DatasetCatalog and MetadataCatalog.
2. Feed the data into the Detectron2 training pipeline, ensuring that both segmentation and gender classification labels were provided.
3. Monitor loss functions:
 - Total Loss
 - Mask Loss (Measures segmentation mask quality)
 - Bounding Box Loss (Measures object detection accuracy)
 - Classification Loss (Measures gender classification accuracy)
4. Save model checkpoints at regular intervals.

Training Hardware and Environment:

The model was trained on a **NVIDIA RTX 3060 (12GB VRAM) GPU** with an **AMD Ryzen 5 3600 (6-Core, 12-Thread) CPU** and **32GB RAM**. The training process utilized **Detectron2** with **PyTorch**, running on **Ubuntu 20.04 LTS** with **CUDA 11.8** and **cuDNN** for GPU acceleration.

The training took approximately **3 hours** for 10,000 iterations with a batch size of 512, utilizing the GPU for faster computation while CPU threads handled data loading and preprocessing tasks efficiently.

Training Metrics and Performance Monitoring:

During training, several key performance metrics were monitored:

- Segmentation Performance:
 - Mask AP (Average Precision) was computed to measure how well the model segmented individuals.
- Gender Classification Performance:
 - Accuracy: Percentage of correctly classified males and females.
 - Precision & Recall: Performance for each gender category.

Training logs were continuously monitored to detect overfitting or stabilize learning rate adjustments.

Optimization and Fine-Tuning:

To further improve performance, several optimization techniques were applied. The learning rate was initially set to 0.001, but cyclical learning rate scheduling was tested to enhance convergence. Batch size adjustments were also explored to balance GPU memory usage and convergence speed.

Phase 5: : Evaluation and Visualization

For the evaluation phase, the trained model was tested on a held-out validation set to assess its performance in both instance segmentation and gender classification tasks. The evaluation metrics included Mask Average Precision (AP) and Mask Average Recall (AR) for segmentation quality, while accuracy, precision, recall were used to measure the effectiveness of gender classification. These metrics provided insight into the model's ability to correctly segment people in the dataset and classify their gender with high confidence.

To ensure comprehensive evaluation, the dataset was processed using Detectron2's built-in evaluation tools, which computed segmentation metrics based on ground truth annotations. For gender classification, predictions were extracted from the model and compared against manually labeled gender annotations.

In the visualization phase, multiple test images were selected to analyze the model's performance. Each image was processed through the trained Mask R-CNN model, generating segmentation masks overlaid on the input image. The visualization incorporated color-coded masks, with different colors. Additionally, bounding boxes were drawn around each detected person, and gender labels were displayed in distinct colors—green for male and red for female—to distinguish classification results.

Evaluation Results Analysis:

The evaluation of the trained instance segmentation and gender classification model on the validation dataset provides valuable insights into its performance. Below is a breakdown of the key findings:

1. Object Detection (Bounding Box) Performance

- Overall AP (Average Precision @ IoU=0.50:0.95): 58.7%
- AP at IoU 0.50 (AP50): 77.8%
- AP at IoU 0.75 (AP75): 65.9%
- AP for Medium-Sized Objects: 64.8%
- AP for Large-Sized Objects: 58.9%
- Per-Category AP:
 - Male: 62.99%
 - Female: 54.40%

Interpretation:

The object detection performance is fairly strong, with an AP of 58.7%, which is competitive for instance segmentation tasks. The model performs better on medium-sized objects than on large objects, suggesting that it effectively detects individuals when they are at a moderate scale but has slight difficulty with very large subjects. The AP for males (62.99%) is higher than for females (54.40%), indicating a possible imbalance in data distribution or feature representation differences between male and female individuals.

2. Instance Segmentation Performance

- Overall Mask AP (Average Precision @ IoU=0.50:0.95): 48.2%
- AP at IoU 0.50 (AP50): 74.4%
- AP at IoU 0.75 (AP75): 55.0%
- AP for Medium-Sized Objects: 36.9%
- AP for Large-Sized Objects: 49.3%
- Per-Category Mask AP:
 - Male: 55.12%
 - Female: 41.35%

Interpretation:

The mask segmentation performance (48.2% AP) is reasonable, though slightly lower than the bounding box detection performance. The high AP50 (74.4%) suggests that the model is relatively confident in detecting masks at a lower IoU threshold but struggles at higher overlaps (IoU 0.75). The performance drop for medium-sized objects (36.9%) could indicate challenges in segmenting individuals when they are not as well-defined in the image. The segmentation AP for males (55.1%) is significantly higher than for females (41.3%), which aligns with the object detection results and suggests that female representations in the dataset might need further balancing or improved feature learning.

The Mask R-CNN model performed well in detecting and segmenting people, but segmentation accuracy was lower than bounding box detection, particularly for female instances and medium-sized individuals. The results indicate room for improvement in dataset balancing, augmentation techniques, and model architecture refinement. Despite these limitations, the model provides a strong foundation for real-world applications in instance segmentation and gender classification.

5. Conclusion

The project successfully implemented an instance segmentation and gender classification model using the LV-MHP-V2 dataset and the Detectron2 framework. Through extensive preprocessing, dataset modifications, and model training, the system was able to detect multiple individuals in an image, generate high-quality segmentation masks, and classify each detected person as male or female. The dataset was carefully processed to exclude specific body parts such as faces and hands, ensuring more accurate instance segmentation. Since the original dataset did not include gender annotations, they were manually added to enable gender classification.

During model development, Mask R-CNN was chosen as the primary architecture, with fine-tuning performed to optimize its segmentation capabilities. A classification head was integrated into the model to predict gender labels, and hyperparameter tuning along with data augmentation techniques helped improve generalization. The evaluation phase demonstrated competitive results, with an Average Precision (AP) of 58.7% for object detection and 48.2% for segmentation. The gender classification accuracy was satisfactory, though performance on female instances was lower, suggesting dataset biases. The model was assessed using COCO-style evaluation metrics on a validation set.

Several challenges were encountered throughout the project. One significant issue was the gender imbalance in the dataset, as the model performed better on male instances. More balanced gender representation or additional feature enhancements could improve classification results. Segmentation difficulties for medium-sized objects were also observed, with a lower AP compared to larger individuals. Implementing better augmentation techniques and refining the segmentation pipeline could help address this limitation. Computational constraints were another challenge, as training was conducted on an NVIDIA RTX 3060 GPU with a Ryzen 5 3600 processor, which limited batch size and training speed. Future improvements could involve training on higher-end GPUs or cloud-based solutions.

To further enhance model performance, future work could focus on using a more advanced backbone such as ResNeXt-101 or Swin Transformer to improve feature extraction. More extensive data augmentation techniques could be applied to address segmentation difficulties. Expanding the dataset with more gender-labeled images would improve classification accuracy and mitigate dataset biases. Additionally, optimizing inference speed and computational efficiency would allow the model to be deployed in real-world applications.

This project demonstrated the effectiveness of deep learning-based instance segmentation and gender classification. By leveraging Detectron2's powerful object detection and segmentation capabilities, a robust model was developed that performs well on densely packed human images. The results are promising, and with further refinements, the system can be applied in real-world computer vision tasks such as surveillance, human-computer interaction, and crowd analytics.