



# Deep Learning Algorithm using CSRNet and Unet for Enhanced Behavioral Crowd Counting in Video

B.Ganga<sup>1,2</sup>, Lata B T<sup>1</sup>, Rajshekar<sup>1</sup> and Venugopal K R<sup>3</sup>

<sup>1</sup>Department of CSE, UVCE, Bangalore, India

<sup>2</sup>Dept of CSE (Cyber Security), Sambhram Institute of Technology, Bangalore, India

<sup>3</sup>Former V.C, Bangalore University, Bangalore, India

Received 2 April 2024, Revised 30 September 2024, Accepted 3 October 2024

**Abstract:** In crowd analysis, video data incurs challenges due to occlusion, crowd densities, and dynamic environmental conditions. To address these challenges and enhance accuracy, we have proposed Behavioral Crowd Counting (BCC), which combines the Congested Scene Recognition Network (CSRNet) with Unet in video data. The CSRNet combines two networks: (1) a frontend for feature extraction and (2) backend for generating a density map. It effectively tallies individuals within densely populated regions, solving the high crowd density constraints. The Unet builds the semantic map and refines the semantic and density map of CSRNet. The Unet unravels complex patterns and connections among individuals in crowded settings, capturing spatial dependencies within densely populated scenes. It also offers the flexibility to incorporate attention maps as optional inputs to differentiate crowd regions from the background. We have also developed new video datasets, namely the Behavioral Video Dataset, from the fine-grain crowd-counting image dataset to evaluate the BCC model. Datasets include standing vs. sitting, waiting vs. non-waiting, towards vs. away, and violent vs non-violent videos, offering insights into posture, activity, directional movement, and aggression in various environments. The empirical findings illustrate that our approach is more efficient than others in behavioral crowd counting within video datasets consisting of congested scenes as indicated by metrics MSE, MAE, and CMAE.

**Keywords:** Congested Scene Recognition Network (CSRNet), Unet, Feature Extraction, Behaviour, and Crowd Analysis

## 1. INTRODUCTION

Crowd counting assesses the number of individuals in a designated area through images or video. In contrast, crowd behavior focuses on studying and interpreting individuals actions, movements, and interactions within a crowd. The study of crowd behavior has evolved significantly due to the demand for detailed crowd analysis in fields such as Retail Analysis, Law Enforcement, DineSpace Analytics, Pedestrian Flow Monitoring, Traffic surveillance, and Public Safety [1]. Vibha *et al.*, [2] explored methods for eliminating background elements to recognize moving objects in videos featuring a static background. Vibha *et al.*, [3] developed a background registration method for detecting moving vehicles. Traditionally, crowd-counting algorithms have primarily focused on quantifying the number of individuals within an image. However, this count-only approach is inadequate in providing deeper insights into crowd dynamics and behaviors, which are critical for various practical applications. Hence, there is a growing research interest in the detailed analysis of crowd videos.

Analyzing crowds using video technology is increas-

ingly important, given the vast amount of crowd-related information in video form. Traditional methods are inadequate for comprehensive understanding and interpreting the data in videos. Therefore, it is essential to focus on the immense potential of video data for in-depth crowd analysis. It goes beyond merely quantifying crowd counts in images and categorizing crowds in videos based on the action.

This work primarily analyzes crowd behavior in Retail, Law Enforcement, DineSpace Analytics, and Pedestrian Flow Monitoring applications. The transition from image to video analysis has become crucial in these applications, reflecting the current trend where most crowd data is now captured through video. Traditional retail analysis relies solely on static head counts. In contrast, video-based behavioral insights offer a dynamic perspective, providing a deeper understanding of specific sub-categories, such as individuals in queues or leisurely browsing. Additionally, in law enforcement, challenges faced in crowd management are resolved in crowd behavioral video analysis by distinguishing violent and non-violent individuals within a crowd. Similarly, in restaurant or cafeteria settings, video



analysis becomes an indispensable tool for distinguishing between standing and sitting people, enhancing the depth and precision of DineSpace analysis. In Pedestrian Flow Monitoring, our emphasis lies in leveraging video analysis to improve crowd control strategies by moving beyond traditional static observations, empowering the distinction and management of the flow of pedestrians with greater depth and precision, and addressing critical challenges in crowd dynamics.

The proposed approach has various applications, such as aiding violence detection and crowd control in intelligent city environments, ensuring operational efficiency in event management, optimizing transportation, and enhancing public space safety.

#### A. Challenges/Motivation

**Enhanced Crowd Analysis:** The growing need for in-depth crowd analysis in video data drives increasing demand across various real-world contexts, such as retail, surveillance, and public safety. Traditional approaches to crowd counting have limitations when offering thorough insights into crowd behaviors, which are vital for addressing practical challenges in these domains. The demand drives the development of more sophisticated crowd-analysis techniques that extend beyond the scope of basic crowd-counting methods to incorporate behavioral crowd-counting.

#### B. Contributions

- 1) **Development of BCC Architecture:** The introduction of the Behavioral Crowd Counting (BCC) architecture to integrate the Congested Scene Recognition Network (CSRNet) with the Unet.
- 2) **Adaptability to Congested Scenes:** The CSRNet effectively counts individuals in densely populated areas, addressing the challenge of adapting to dense crowd densities.
- 3) **Efficient Spatial Dependency Capture:** The Unet deciphers intricate patterns and relationships among individuals in crowded environments to capture spatial dependencies within crowded scenes.

#### C. Organizational Structure

The rest of the paper is as follows. Section 2 provides insights into existing research in the Related Work. Section 3 discusses the Background, Problem Definition, Objectives, and details of the BCC Architecture. Section 4 discuss experimental setup. Section 5 presents the results and discussions. Section 6 contains the Conclusion and Future Work.

## 2. RELATED WORK

Related work on Crowd Behavior Analysis Models is presented in TABLE 1. Cem *et al.*, [4] have distinguished normal and abnormal crowd behaviors in surveillance videos using Motion Information Images (MIIs) derived from optical flow data. The merit is that the optical flow data improves the accuracy of identifying panic and escape

behaviors. The demerit is that the real-time application is not explored in this work due to the resource-demanding nature of MII generation. Guo *et al.*, [5] have introduced a crowd anomaly detection method for video service robots, combining mean shift and k-means to identify abnormal behavior in crowded scenarios. The merit is that the technique classifies categories with similar motion patterns and improves anomaly detection accuracy. The demerit is that the computational parameters for domain and spatial bandwidth require precise tuning.

Junyu *et al.*, [6] have developed Multi-level Feature-aware Adaptation (MFA) and Structured Density map Alignment (SDA) to address challenges in supervised learning for crowd counting and pixel-wise density estimation. The advantages are that it overcomes data scarcity issues and outperforms existing methods in cross-domain crowd counting. The challenge arises in distinguishing between background and foreground areas with similar textures, leading to inaccuracies in the estimated crowd count.

Hyojun *et al.*, [7] have addressed wildlife monitoring issues by introducing automated multi-class object counting for endangered animal species. This work presents a fine-grained multi-class object counting dataset known as KRGRUIDAE. The advantage is that EcoCountNet's network contributes to accurate and efficient counting processes. The disadvantage is that the network requires additional computational resources, and researchers have not explored real-time applications due to its complexity. Yongtuo *et al.*, [8] have presented crowd counting model adaptability across different domains by combining two modules, Crowd Region Transfer (CRT) and Crowd Density Alignment (CDA). The merit is that the model exhibits promising performance in various adaptation scenarios. The drawback is that the point-level crowd-counting annotations for crowd images are still challenging problems and expensive.

Savchenko *et al.*, [9] have presented an efficient frame-level facial emotion analysis model that combines embeddings and scores from the EfficientNet architecture pre-trained on AffectNet. The merit is that the model outperforms the baseline on multiple tasks, including facial expression recognition and valence-arousal estimation. The demerit is that generalizing the model to real-world scenarios is a challenging problem. Justin *et al.*, [10] discussed fine-grained counting using crowd-sourced annotations to estimate individuals in crowded scenes and classify attributes. The merit is that The Seal Watch dataset contains eight fine-grained classes that advance research in animals. This work has not implemented detection-based methods for behavior analysis. Pierre *et al.*, [11] have designed crowd behavior by exploring various approaches such as microscopic, macroscopic crowd modeling, motion-based crowd behavior analysis, and optical flow utilization. Crowd analysis offers the advantage of being applicable in areas such as public safety, market analysis, urban planning, and entertainment. However, it also generates large volumes of data, which are challenging to store, process, and analyze efficiently.

Shenjian *et al.*, [12] have presented a bi-level alignment

TABLE I. Summary of Crowd Counting Models, Datasets, Advantages, and Disadvantages

Author/Year	Algorithm/Model	Dataset	Advantages	Disadvantages
Cem <i>et al.</i> , [4] 2020	MIIs from optical flow	UMN and PETS2009	Improved accuracy	Computational intensity
Shuqiang <i>et al.</i> , [5] 2019	Mean shift & k-means	UMN dataset	Improved accuracy	Parameter tuning
Junyu <i>et al.</i> , [6] 2020	MFA and SDA	ShanghaiTech Part B, WorldExpo'10, Mall, UCSD	Overcoming data scarcity	Error in similar textures
Hyojun <i>et al.</i> , [7] 2021	EcoCountNet	KR-GRUIDAE (fine-grained object counting)	Accurate counting	Additional resources and not real-time
Yongtuo <i>et al.</i> , [8] 2023	CRT and CDA	GCC, ShanghaiTech Part A	Promising adaptation	Annotation costly
Savchenko <i>et al.</i> , [9] 2022	EfficientNet	AffectNet	Outperforms baseline	Generalization challenge
Justin <i>et al.</i> , [10] 2022	Crowd-sourced fine-grained	Seal Watch dataset	Enhanced crowd management	Lack of size discussion
Pierre <i>et al.</i> , [11] 2019	Various approaches	The UCSD Anomaly Detection Dataset	Applicable in diverse areas	Large data challenge
Shenjia <i>et al.</i> , [12] 2022	Bi-level alignment	GTAS Crowd Counting (GCC) dataset	Addresses domain adaptation	Increased complexity
Sachin <i>et al.</i> , [13] 2023	MCNN	ShanghaiTech dataset	Enhanced crowd management	Size and resolution challenges
Zhikang <i>et al.</i> , [14] 2019	Count attention	ShanghaiTech dataset	Improved accuracy	Future applicability
Ye <i>et al.</i> , [15] 2019	DFFnetSeg	Test dataset from CNDnet2014	Handles modification	Increased computational complexity
Adel <i>et al.</i> , [16] 2021	U-ASD Net	Haramain, with three different scenes	Adapts to scenarios	More computational resources
Elizabeth <i>et al.</i> , [17] 2021	Integrated approach	Crowd datasets	Early detection	Data availability challenges
Xiaohegn <i>et al.</i> , [18] 2020	DANet and ASNet	ShanghaiTech Part A, UCF CC 50, UCF-QNRF, WorldExpo'10	Alleviates density differences	More accurate counting
Naveed <i>et al.</i> , [19] 2021	CNN-based model	ShanghaiTech (Part-A, Part-B), Venice	Improved accuracy	Semantic segmentation expansion
Yadi <i>et al.</i> , [20] 2023	AI-based analytics	Video records from a platform scenario	Accurate pedestrian counting	Advanced technology required
Reem <i>et al.</i> , [21] 2023	Enhanced abnormal detection	Diverse Hajj dataset	Impressive results with scalability	Model complexity

framework for enhancing synthetic-to-real Unsupervised Domain Adaptation (UDA) crowd counting. The merit is that the model addresses domain adaptation problem, while demerit is that it involves increased computational complexity. Sachin *et al.*, [13] have predicted crowd behavior using a Multicolumn Convolutional Neural Network (MCNN) on the ShanghaiTech dataset. Merit is enhanced crowd management in various real-world scenarios, and demerit is that diverse image sizes and resolutions need to be addressed in this work. Zou *et al.*, [14] have presented adaptive model to allocate different capacities to different regions in an image based on crowd density. The merit is that it improves crowd-

counting accuracy in various scenarios. The broader range of scenarios, like videos, has not been experimented with in this work.

Ye *et al.*, [15] have presented the foreground / background segmentation method DFFnetSeg for video analysis. The merit is that the model handles both unseen and changes in scene, making it suitable for diverse video scenarios. The drawback is the increase in computational complexity due to multiple frames. Adel *et al.*, [16] have proposed U-ASD Net concatenating U-Net and Adaptive Scenario Discovery to address perspective distortions and scale variations in crowd counting. It adapts to complex

scenarios, making it suitable for dense and sparse datasets. The demerit of the model is that it requires more computational resources due to the number of parameters and training time.

Elizabeth *et al.*, [17] have explored crowd behavior analysis by integrating psychology theories, IoT, and cognitive computing for predictive crowd management. The advantages are that it enables early detection of crowd disasters and the potential for real-time data processing. However, there are challenges in data availability. Xiaoheng *et al.*, [18] have addressed the challenges of crowd counting using Convolutional Neural Networks (CNNs) by combining two networks, namely Density Attention Network (DANet) and Attention Scaling Network (ASNet). It improves counting performance differences in regions with varying crowd density patterns, leading to more accurate crowd counting. Naveed *et al.*, [19] have proposed a CNN-based model comprising of three main modules viz., a backbone network for general features, Dense Feature Extraction Modules (DFEMs) with dense connections, and a Channel Attention Module (CAM) for class-specific responses. It improves counting accuracy in scenes with significant perspective variations and varying density levels. This work does not explore semantic segmentation.

Yadi *et al.*, [20] have presented a comprehensive AI-based model for crowd analytics in rail transit stations, focusing on flow volume, crowd density, and walking speed analysis. It successfully achieves accurate pedestrian counting and practical applications, requiring advanced technology and calibration equipment. Reem *et al.*, [21] have enhanced abnormal behavior detection in large crowds using the diverse Hajj dataset. While achieving impressive results, it faces challenges in real-world scalability and model complexity. Ahmed *et al.*, [22] have developed a cloud-based deep learning framework for early detection of crowding in event entrances. It demonstrates 87% accuracy but encounters real-time implementation challenges and environmental influences. Ganga *et al.*, [23] have proposed a crowd-counting method combining Unet and GAN architectures to generate crowd-density maps with minimized feature loss. Ganga *et al.*, [24] have presented AnomalyDetectNET for video anomaly detection. Ganga *et al.*, [25] have discussed crowd counting and behavioral analysis by examining CNN-based methods, particularly CSRNet and UNet, which enhance accuracy and efficiency in these tasks. Ganga *et al.*, [26] [27] have utilized a UNet GAN for accurate crowd counting and a CNN for density classification, improving stability and density evaluation for behavioral analysis and categorizes crowd behavior as violent or non-violent by analyzing density maps and individual actions.

Rongyong *et al.*, [28] have developed an enhanced MCNN for precise pedestrian head recognition and crowd mass evaluation, focusing on dynamic crowd stability analysis and accident detection. Aliyu *et al.*, [29] have investigated computer vision techniques for analyzing abnormal crowd behavior, emphasizing attribute identification and optimal detection methods in dynamic settings. Yamin *et*

*al.*, [30] have proposed SSODTL-CD2C for crowd density detection and classification, utilizing optimization and transfer learning for dense crowd environments such as HAJJ. Skander *et al.*, [31] have introduced MILP-MPCT for improved crowd tracking and UAV deployment optimization. Faisal *et al.*, [32] have developed a real-time framework using semantic segmentation and deep learning for crowd-tracking and anomaly detection in diverse outdoor environments. Maria *et al.*, [33] have developed deep convolutional neural networks for human crowd detection in drone imagery, emphasizing lightweight architectures for distinguishing crowded and non-crowded scenarios in drone protection. Manu *et al.*, [34] have introduced the IMFF framework for crowd behavior analysis, enhancing crowd management through multi-level feature fusion and local region classification.

### 3. BACKGROUND

The fine-grained crowd-counting [1] method features a two-branch architecture consisting of density-aware feature propagation and complementary attention mechanisms. In the density-aware feature propagation phase, the model iteratively propagates features to capture contextual information, explicitly focusing on high-density areas and predicting the overall crowd density map. The complementary attention mechanisms exchange information between the two branches, and individual pixels are categorized effectively. Furthermore, during training, the model combines three loss functions, counting loss, segmentation loss, and fine-grained loss, to optimize its performance. The method shows high accuracy of crowd counting in scenarios where fine-grained categorization of crowd segments is necessary. Ground-truth density maps  $Y_j$  are generated by convolving dot maps  $D_j$  with a 2D Gaussian kernel  $K$  as shown in equation (1).

$$Y_j = D_j \cdot k\sigma \quad (1)$$

Ground-truth segmentation maps  $S_j$  is obtained from the ground-truth density maps as shown in equation (2).

$$S_j = \frac{Y_j}{\eta + \sum_{j=1}^k Y_j} \quad (2)$$

Here  $\eta$  is a small number to prevent division by zero, and background segmentation  $S_{(k+1)}$  is defined by regions with low density. The soft cross-entropy loss given in equation (3) is used for segmentation.

$$\text{softCELoss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C (Y_{ij} \log \hat{Y}_{ij}) \quad (3)$$

Where  $Y_{ij}$  represents the ground truth class probability,  $\hat{Y}_{ij}$  is the predicted class probability for class,  $N$  is the total number of data points, and  $C$  is the number of classes.

#### A. Problem Definition

Given a video clip/data consisting of images of a certain length and situations, the objective is to explore Behavioral Crowd Counting (BCC) by combining CSRNet and Unet.

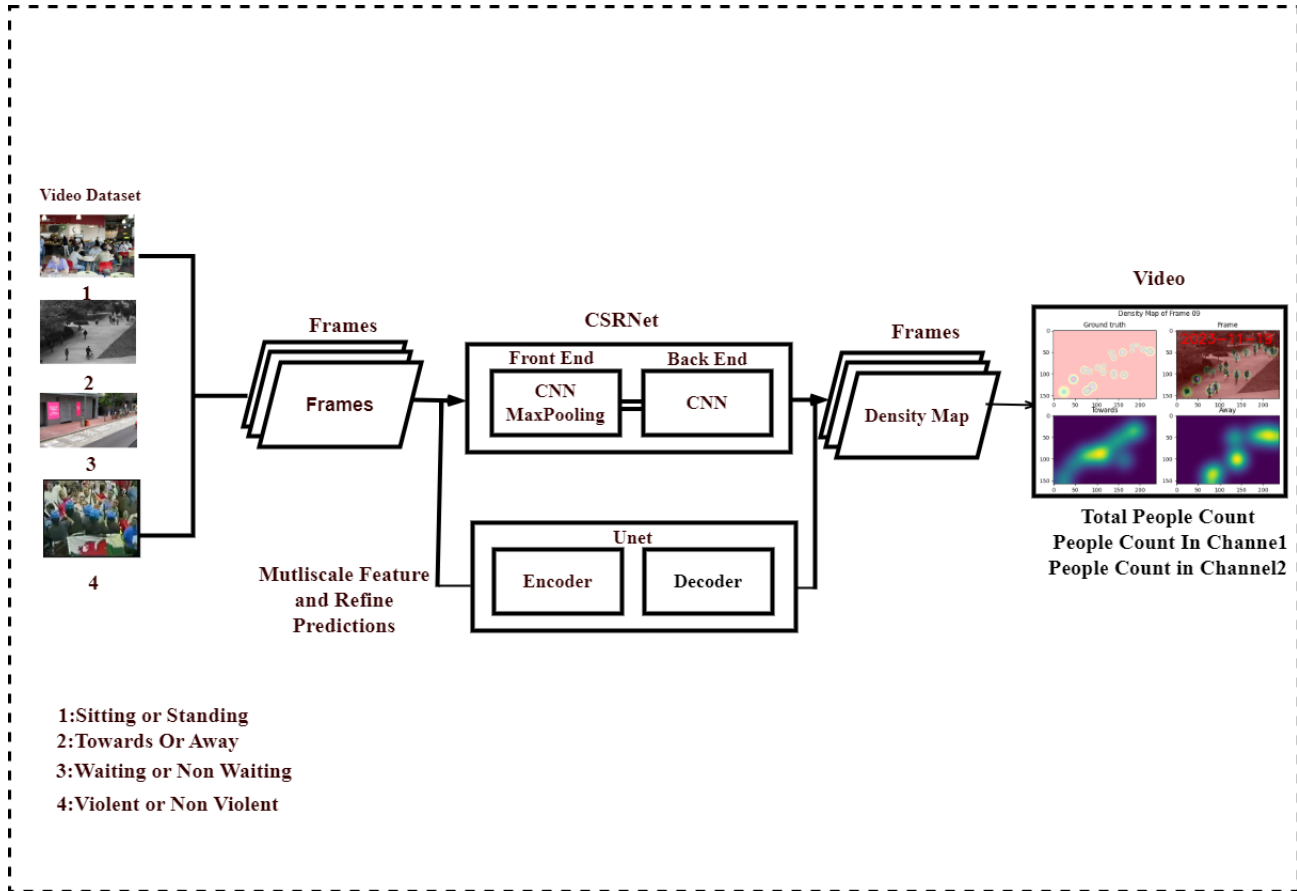


Figure 1. The Architecture of Behavioral Crowd Counting (BCC)

TABLE II. Frontend Network (Feature Extraction)

Layer	Input Channels	Output Channels	Kernel Size	Stride	Padding	Activation Function	Size
Convolutional Layer 1	1	64	5x5	1	2	Leaky ReLU (0.1)	Varies
Max-Pooling Layer 1	64	64	2x2	2	0	None	Varies
Convolutional Layer 2	64	128	5x5	1	2	Leaky ReLU (0.1)	Varies
Max-Pooling Layer 2	128	128	2x2	2	0	None	Varies
Convolutional Layer 3	128	256	5x5	1	2	Leaky ReLU (0.1)	Varies
Max-Pooling Layer 3	256	256	2x2	2	0	None	Varies
Additional Convolutional Layer	Varies	Varies	Varies	Varies	Varies	Varies	Varies

### B. Objectives

- 1) BCC Architecture: To design and construct the Behavioral Crowd Counting (BCC) architecture by integrating the CSRNet and Unet to enable behavioral crowd counting in video data.
- 2) Enhanced Accuracy and Segmentation: To improve the accuracy of crowd counting and crowd segmentation by accurately distinguishing crowd regions from the background and analyzing crowd behaviors in video data.

TABLE III. Backend Network (Density Map Estimation)

Layer	Input Channels	Output Channels	Kernel Size	Dilation	Activation Function
Convolutional Layer 1	512	512	3x3	Optional	None
Convolutional Layer 2	512	256	3x3	Optional	None
Convolutional Layer 3	256	128	3x3	Optional	None
Convolutional Layer 4	128	64	3x3	Optional	None
Output Layer	64	1	1x1	None	None

TABLE IV. Unet Architecture

Layer	Operation	Input Channels	Output Channels	Kernel Size	Stride	Padding
Encoder Conv1	Conv2d, LeakyReLU, Conv2d, LeakyReLU	i_cn	64	3x3	1	1
Max Pooling 1	MaxPool2d	64	64	2x2	-	-
Encoder Conv2	Conv2d, LeakyReLU, Conv2d, LeakyReLU	64	128	3x3	1	1
Max Pooling 2	MaxPool2d	128	128	2x2	-	-
Decoder Upconv1	ConvTranspose2d	128	64	2x2	2	0
Decoder Conv1	Conv2d, LeakyReLU, Conv2d, LeakyReLU	128	64	3x3	1	1
Decoder Upconv2	ConvTranspose2d	64	o_cn	2x2	2	0
Refinement Layer	Conv2d	64	o_cn	1x1	1	0

- 3) Optional Attention Map Integration: To integrate optional attention maps into BCC to focus on specific areas of interest within the crowd to refine density predictions.

### C. Architecture of Behavioral Crowd Counting (BCC)

The Behavioral Crowd Counting (BCC) architecture, as shown in Figure 1, concatenates Congested Scene Recognition Network (CSRNet) and Unet to achieve better crowd counting by including behavior in video data. CSRNet is renowned for its precision in estimating crowd density within congested scenes, offering a deep understanding of the intricate details in densely populated areas. The merit of CSRNet is that it excels in capturing contextual information, allowing it to comprehend the spatial relationships among individuals within a crowd. The CSRNet comprises two major components: a frontend and a backend network. The frontend network consists of several convolution and max-pooling operations for feature extraction. The backend network incorporates dilated convolutions for feature extraction and comprises a series of convolution layers that process features extracted by the frontend network, producing the estimated crowd density map. TABLES II and III show the Frontend Network (Feature Extraction) and Backend Network (Density Map Estimation) layers for their input and output channel dimensions, kernel sizes, stride values, padding, and activation functions.

The CSRNet is initialized with weights, and the front

end can be pre-trained optionally on ImageNet. The frontend has the following methods: forward method, make layers, and initialize. The forward method processes input images, extracts features to produce crowd density and semantic segmentation maps. The make layers methods generate sequential layers based on the provided configuration, including convolution layers, batch normalization, and ReLU activations. The initialize method is responsible for initializing the weights of the network modules, viz., convolution layers and batch normalization. Overall, the architecture employs VGG16-inspired features for effective feature extraction in crowd analysis, explicitly focusing on counting and semantic segmentation. The backend with dilated convolutions efficiently processes input images to produce accurate crowd density and segmentation maps.

The Unet incorporates a two-tiered structure namely, encoders serve as feature extractors and decoders for processing input features for segmentation, as shown in TABLE IV. The encoder consists of two convolution layers, LeakyReLU activation and Max-pooling. Leaky ReLU activation functions follow the convolution layers in its two blocks. Leaky ReLU activation functions initiate non-linearity in the network, which allows a slight, non-zero gradient for negative inputs, preventing dead neurons and facilitating the learning of more complex relationships in the data. Max-pooling operations follow each encoder block to down-sample and capture hierarchical and spatial information. The down-sampling is pivotal for progressive

abstracting and concentrating relevant information from the input images, enabling the network to learn hierarchical representations while maintaining computational efficiency.

The decoder has transposed convolution layers, used to up-sample the feature maps, with two such layers in each decoder block. These layers contribute to the reconstruction of spatial details lost during the down-sampling process, aiding in the precise localization of features. The transpose convolution layers refine the feature representations, maintaining symmetry with the encoder, and the final refinement occurs through a  $1 \times 1$  convolution layer in the refinement block. The Unet benefits are particularly effective in tasks such as image segmentation, where the accurate spatial delineation is essential. It also effectively performs image segmentation tasks, and its emphasis on accurate spatial delineation makes it potent in applications where precise localization of features is vital, such as medical image analysis, autonomous vehicle navigation, and satellite image processing.

#### D. Algorithm

Algorithms 1 to 6 explain Pre Processing, Segmentation of video, and Crowd Counting. The Pre Processing function is designed to handle video data. It takes a video path as input and performs the following tasks: frame extraction, saving frames to an output directory, and reading annotation data from a corresponding JSON file. The function utilizes OpenCV to read and store frames, organizing them based on the video number in a specified directory structure. The preprocessing output is the total number of frames, the path to the frame outputs extracted annotation data, and the video number. The segmentation function focuses on video segmentation. It calls the processing results function to obtain a result and then iterates through each frame. The data is prepared for each frame, and the RefineSegmentation function is applied to refine the segmentation using an Unet model. The function returns the final segmented result for each frame in the video. It takes input data containing features and an attention map, concatenates them if an attention mechanism is used, initializes the Unet model, and performs segmentation refinement. The refined segmentation result is then returned.

In the Crowd Counting function, the video undergoes crowd counting processing. It first calls the IntegratedProcessing function to obtain a processing result. The function then iterates through each frame, prepares the data, and uses the PredictCrowdCount function to estimate the crowd count. The final result represents the crowd count information for each frame in the video. The PredictCrowdCount function estimates crowd counts in images. It takes an image as input, initializes a CSRNet model for crowd counting, processes the image through its frontend and backend, and predicts the crowd count using the output layer of the model. The result is the expected crowd count for the given image.

The algorithm Enhanced Behavioral Crowd Counting (EBCC) accurately counts individuals in crowded scenes by considering their behavior, such as the direction of

movement (towards or away), posture (standing or sitting), state of people (waiting or non-waiting), and nature of the activity (violent or non-violent). It preprocesses video frames to extract and refine segmentation data and then employs a combination of Unet and CSRNet models to analyze these frames. The process counts the number of people by integrating behavioral analysis, leading to a more comprehensive understanding of crowd dynamics. The versatility of EBCC allows it to be adapted for other applications, such as monitoring public safety and enhancing surveillance systems for behavioral anomaly detection. The algorithm's time complexity is  $O(n)$ , with  $n$  representing the total frames in a video, as it processes each frame separately. The space complexity is  $O(m)$ , where  $m$  is the memory requirements of the Unet and CSRNet models during the algorithm's runtime.

#### E. Case Study

- 1) Smart city environments: Optimizing bus routes with real-time pedestrian flow data in bustling city centers improves transport efficiency and reduces congestion. Challenges include integrating with existing infrastructure and addressing data privacy concerns.
- 2) Event management: During festivals, real-time crowd monitoring improves event safety by detecting overcrowding and enabling swift responses to ensure smooth operations. Challenges include managing real-time data effectively, scaling for significant events, and maintaining attendee's privacy.
- 3) Public safety: At large protests, violence detection algorithms swiftly identify aggressive behavior, aiding law enforcement in maintaining peace and security by enhancing crowd control. Challenges include ensuring algorithm accuracy, addressing ethical concerns, and gaining public acceptance of surveillance methods.

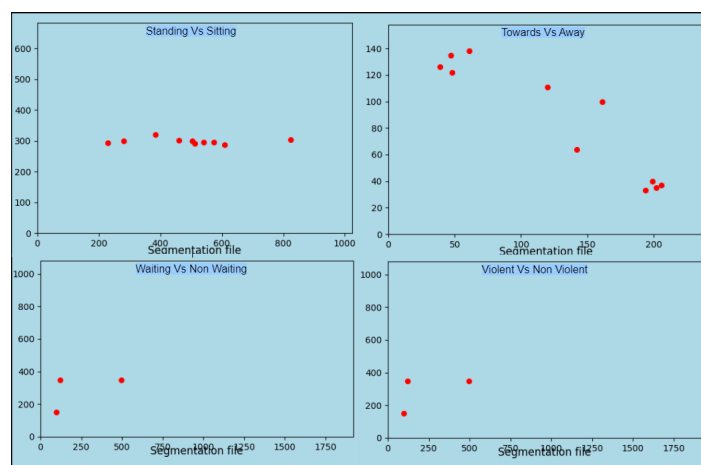


Figure 2. Segmentation Results of Four Video Dataset

## 4. EXPERIMENTS SETUP

The loss functions such as Mean Squared Error (MSE) and Mean Absolute Error (MAE) are the metrics used for

**Algorithm 1** PreProcessing

---

```

1: function PREPROCESSING(video_path, file_name)
2:   Input: video_path, file_name
3:   Output: total_frames,output_path, annotation_data, video_num
4:   dataset_path ←
5:   D:/Final_yr_project/Final_dataset_frames/
6:   file_name ← video_path
7:   video_num ← substring(1, file_name)
8:   frames_dir ← frames + video_num
9:   output_path ← Call join(dataset_path, frames_dir)
10:  if not exists(output_path) then
11:    mkdirs(output_path)
12:  end if
13:  video ← VideoCapture(video_path)
14:  total_frames ← int(video.get(CAP_PROP_FRAME_COUNT))
15:  frame_count ← 0
16:  for frame_count = 0 to total_frames do
17:    success, frame ← read_video
18:    if not success then then
19:      break
20:    end if
21:    frame_path ← output_path + frame_count
22:    Call imwrite(frame_path, frame)
23:    frame_count ← frame_count + 1
24:  end for
25:  release(video)
26:  Print "Video extraction complete!"
27:  annotation_file ← open D:/Final_yr_project/annotations/video_{video_num}.json
28:  annotation_data ← json.load(annotation_file)
29:  return total_frames, output_path, annotation_data, video_num
30: end function

```

---

**Algorithm 2** Segmentation of Video

---

```

1: function SEGMENTATION(video_path)
2:   Input: video_path
3:   Output: Segmentation_result
4:   processing_result ← video_path
5:   Extract information from processing_result
6:   for frame_index = 0 to total_frames do
7:     frame_path←join(output_path,
8:     frame_index)
9:     input_data←PrepareData(frame_path)
10:    segmentation_result←RefineSegmentation
11:    (input_data)
12:    Process segmentation_result
13:  end for
14:  return Segmentation_result
15: end function

```

---

**Algorithm 3** RefineSegmentation

---

```

1: function REFINESEGMENTATION(InputData)
2:   Input: InputData
3:   Output: RefinedSegmentation
4:   fea ← ['fea']
5:   att ← [att]  ▶ Assuming att is the attention map
6:   if att is not None then
7:     fea ← torch.cat((fea, att), 1)
8:   end if
9:   unet_model ← Unet Model(input_channels,
10:  output_channels)
11:   RefinedSegmentation ← unet_model(fea)
12:   return RefinedSegmentation
13: end function

```

---

evaluation. TABLE V shows details of four behavior video datasets, such as frame rate in seconds, length of video in seconds, and resolution of the videos.

**A. Datasets**

The Behavioral video dataset comprises four distinct datasets derived from fine-grained image datasets. Each dataset focuses on specific behavioral distinctions: standing vs. sitting, waiting vs. non-waiting, towards vs. away movement, and violent vs. non-violent actions. These videos are categorized based on criteria such as posture, activity type, and directional movement, offering a detailed analysis



---

**Algorithm 4** EBCC: Enhanced Behavioural Crowd Counting
 

---

```

1: function EBCC(video_path)
2:   Input: video_path
3:   Output: Crowd_Count
4:   Call PreProcessing(video_path)Function 1:
   Preprocessing
5:   Call Segmentation(video_path)Function 2:
   Segmentation
6:   crowd_count ← CrowdCounting(video_path)
   Function 3: Crowd Counting
7:   return crowd_count
8: end function

```

---

**Algorithm 5** Crowd Counting
 

---

```

1: function CROWDCOUNTING(video_path)
2:   Input: video_path
3:   Output: Crowd_Count
4:   processing_result ←
   IntegratedProcessing(video_path)
5:   Extract information from processing_result
6:   for frame_index = 0 to total_frames do
7:     frame_path ← join(output_path,
   frame_index)
8:     input_data ← PrepareData(frame_path)
9:     crowd_count ← Process(input_data)
10:  end for
11:  return Crowd_Count
12: end function

```

---

of human behaviors in urban settings. Each dataset includes labeled videos with statistics on the number of videos, duration, and distribution across categories, providing valuable resources for behavioral analysis and model training.

- 1) Standing vs Sitting Video Dataset: The dataset presents videos derived from a fine-grained image dataset, illustrating individuals either standing or sitting in various urban settings. It is valuable for

---

**Algorithm 6** PredictCrowdCount
 

---

```

1: function PREDICTCROWDCOUNT(x)
2:   Input: Image x
3:   Output: Predicted crowd count (dmap)
4:   i_cn ← number of input channels
5:   o_cn ← number of output channels
6:   csrnet_model ← CSRNet(i_cn, o_cn)
7:   x ← csrnet_model.frontend(x) Frontend Processing
8:   dmap_fea ← csrnet_model.backend(x) {Density
   Map Feature Extraction}
9:   dmap ← csrnet_model.output_layer(dmap_fea)Density
   Map Prediction
10:  return dmap
11: end function

```

---

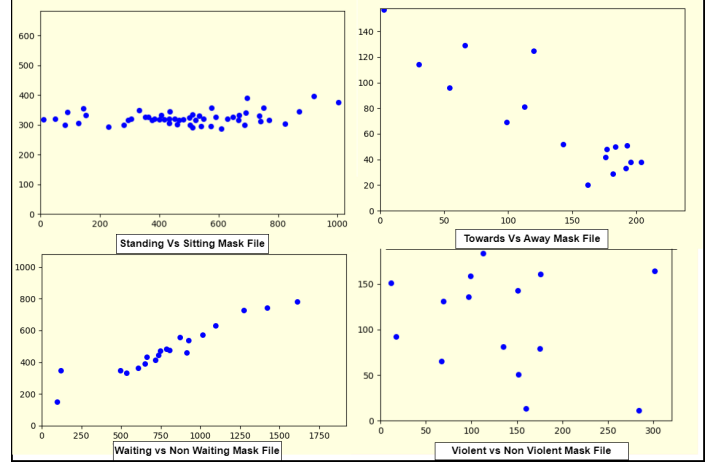


Figure 3. Mask Results of Four Video Dataset

- analyzing and distinguishing between static postures in different environmental contexts.
- 2) Waiting vs. Non-Waiting Video Dataset: The dataset presents videos built from a fine-grained image dataset capturing the people waiting (e.g., at bus stops) versus non-waiting those engaged in other activities. It's an excellent tool for studying patterns of stationary and transient behaviors in public spaces.
- 3) Towards vs. Away Video Dataset: The dataset presents videos built from a fine-grained image dataset featuring people walking towards or away from the camera. It aids in understanding directional movement and pedestrian dynamics, offering insights into approach and departure behaviors in various settings.
- 4) Violent vs Non-Violent Dataset: The dataset presents videos built from a fine-grained image dataset with violent and non-violent video scenes. This resource differentiates aggressive and non-aggressive behaviors in different contexts.

TABLE V. Details of Four Behaviour Video Dataset

Dataset Name	Frame Rate (sec)	Length (sec)	Resolution (dpi)
Violent vs Non-violent	1	1	96
Towards vs Away	1	1	96
Standing vs Sitting	1	1	96
Waiting vs Non-Waiting	1	1	96

### B. Metrics

- 1) Mean Square Error (MSE): MSE is calculated as the average error between the predicted density values and the ground-truth density values, as shown in



equation (4).

$$MSE = \frac{1}{N} \sum_{i=0}^N (Y_i - \hat{Y}_i)^2 \quad (4)$$

Where  $Y_i$  represents the ground truth density map, is the predicted density map, and  $N$  is the total number of data points.

- 2) Mean Absolute Error (MAE): MAE is calculated as the average error between the predicted density values and the ground-truth density values as shown in equation (5).

$$MAE = \frac{1}{n} \sum_i \left| \sum_j Y_{ij} - \sum_j \hat{Y}_{ij} \right| \quad (5)$$

Where  $n$  is the total number of test images,  $y_{ij}$  represents the ground-truth density map for the  $i$  and  $j$ , represents the predicted density map for the  $i_{th}$  test image and the  $j_{th}$  category.

### C. Performances

TABLE VI. Video Information

Property	Value
Frames per second	1
Total frames	6
Video created at	2023-06-29 20:39:56

TABLE VII. Frame Extraction Results

Extracted Frame	Details
Extracted frame: 1/6	Extracted frame: 4/6
Extracted frame: 2/6	Extracted frame: 5/6
Extracted frame: 3/6	Extracted frame: 6/6
Extraction complete!	

TABLE VIII. Frame Analysis Results

Frame Number	CMAE
Frame Number: 1	CMAE: 2.21
Frame Number: 2	CMAE: 1.78
Frame Number: 3	CMAE: 2.21
Frame Number: 4	CMAE: 1.64
Frame Number: 5	CMAE: 2.21
Frame Number: 6	CMAE: 2.21

Figure 2 shows the Segmentation results of the four video datasets with segmentation masks generated. A segmentation mask is a pixel-wise labeling of an image, where each pixel is assigned a category or class based on specific characteristics. The segmentation mask is generated to highlight specific regions of interest within images, guided by annotated points provided in the annotation data. The resulting segmentation mask provides a spatially detailed representation of the annotated features within the image, effectively delineating these features from the background. Figure 3 shows the mask results of the four

TABLE IX. Waiting vs Non Waiting

Metric	Value
Frames per second	1
Total frames	2
Video created at	2023-11-20 19:25:06
Average Speed	60.44 pixels per frame
Total no of people	12.50
Avg no of people in Channel1	10.50
Avg no of people in Channel2	2.00
CMAE	4.12
MAE	2.56
MSE	9.63
PATCH	$5.78 \times 10^{-6}$

TABLE X. Standing vs Non Standing

Metric	Value
Frames per second	1
Total frames	7
Video created at	2023-11-20 19:16:18
Average Speed	129.92 pixels per frame
Total no of people	56.57
Avg no of people in Standing	39.43
Avg no of people in Sitting	17.14
CMAE	5.06
MAE	3.64
MSE	15.12
PATCH	$3.85 \times 10^{-5}$

video datasets with the binary masks generated. Each mask isolates specific regions of interest within the corresponding images by assigning binary values to pixels. The regions of interest are determined by annotated points obtained from the annotation data associated with the images. The masks are created to selectively highlight and distinguish particular features within the images, as dictated by the annotated points. The grayscale intensity in the mask corresponds to the binary values assigned to pixels, where white (or lighter shades) typically represent annotated regions (binary value

TABLE XI. Violent vs Non Violent

Metric	Value
Frames per second	2
Total frames	10
Video created at	2023-11-20 19:19:40
Average Speed	0.56 pixels per frame
Total no of people	21.10
Avg no of people in Towards	12.60
Avg no of people in Away	8.50
CMAE	3.16
MAE	3.09
MSE	6.20
PATCH	0.0005573

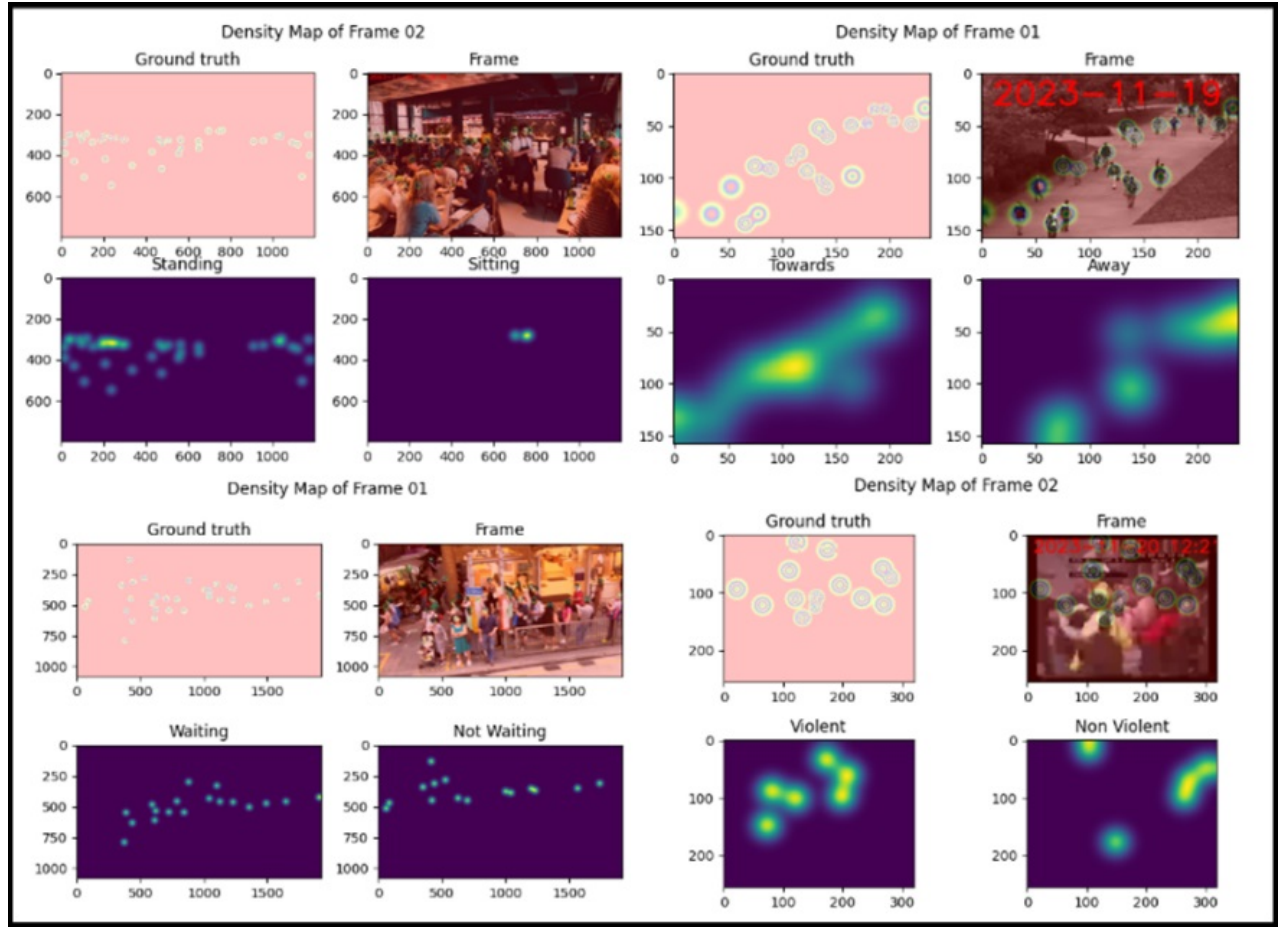


Figure 4. Results of BCC architecture with four Dataset

TABLE XII. Towards vs Away Data

Metric	Value
Frames per second	1
Total frames	6
Video created at	2023-11-20 19:21:44
Average Speed	20.18 pixels per frame
Total no of people	28.67
Avg no of people in Violent	21.67
Avg no of people in Non Violent	7.00
CMAE	4.25
MAE	3.26
MSE	10.64
PATCH	0.0003048

1), and black (or darker shades) represent non-annotated 10 areas (binary value 0). The binary masks play a vital role in the precise examination of object detection, segmentation, and feature analysis.

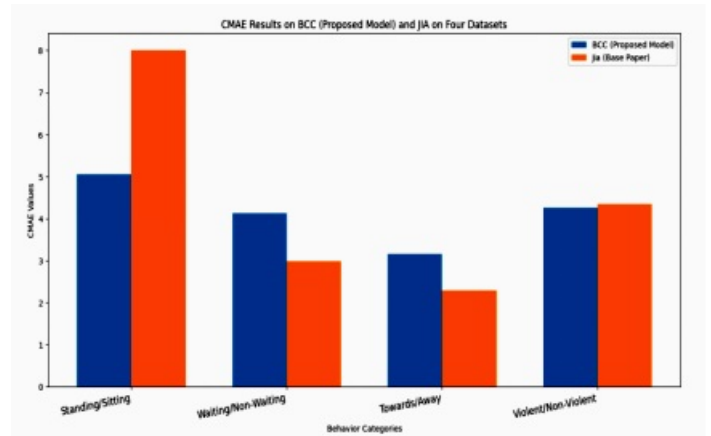


Figure 5. Bar Graph Results of CMAE for each Category

### 5. RESULTS AND DISCUSSIONS

In TABLE VI, the video information details the video, including the frames per second (1 FPS), total frames (6),



TABLE XIII. Comprehensive overview of BCC Architecture metrics across four datasets

Model	Standing/Sitting	Waiting/NonWaiting	Towards/Away	Violent/NonViolent
BCC (Proposed Model)	5.06	4.12	3.16	4.25
Jia (Base Paper) [1]	8.01	2.99	2.29	4.35

TABLE XIV. CMAE Results on four Video Datasets

Category	Frames (sec)	Total Frames	Date Created	Time Created	Avg Speed	Total Speed	Channel 1	Channel 2	CMAE	MAE	MSE	PATCH (10 <sup>-5</sup> )
Standing vs Sitting	1	7	2023-11-20	19:16:18	129.92	56.57	39.43	17.14	5.06	3.64	15.12	3.85
Waiting vs Non Waiting	1	2	2023-11-20	19:25:06	60.44	12.50	10.50	2.00	4.12	2.56	9.63	5.78
Towards vs Away	2	10	2023-11-20	19:19:40	0.56	21.10	12.60	8.50	3.16	3.09	6.20	5.57
Violent vs Non Violent	1	6	2023-11-20	19:21:44	20.18	28.67	21.67	7.00	4.25	3.26	10.64	3.04

and timestamp. The frame extraction results in TABLE VII show the extraction of each frame. In TABLE VIII, the Frame Analysis Results displays frame numbers and their respective Comparative Mean Absolute Error (CMAE) values. Lower CMAE values indicate more accurate predictions. TABLES VI to VIII offer a comprehensive overview of the video-to-frame conversion process, presenting critical metadata, analytical outcomes, and converting a video into frames. TABLES IX, X, XI and XII provide a comprehensive summary of video analysis metrics, such as processing speed (Frames per Second), total frames, video creation time, average object speed, total number of people, average count of people moving towards and away, Comparative Mean absolute Error (CMAE), Mean Absolute Error (MAE), Mean Squared Error (MSE), and a PATCH metric for all four datasets. These metrics offer insights into the dynamics, accuracy, and characteristics of crowd behavior across diverse scenarios and dataset.

Figure 4 presents the experimental results for four datasets i.e., Towards/Away, Standing/Sitting, Waiting/ Non Waiting and Violent/Non-Violent, using four separate figures for each dataset. In the Towards/Away dataset, the figures visualize crowd movement direction featuring annotated ground truth, crowd frames, and the number of people moving towards and away. For the Standing/Sitting dataset, the results focus on postures (standing or sitting), displaying annotated ground truth, crowd frames, and the number of people standing and sitting. In Waiting/Non Waiting dataset figures show the annotated ground truth, crowd frames, and number of waiting and non waiting

people. The Violent/Non-Violent dataset figures illustrate instances of violence, showcasing annotated ground truth, crowd frames, and the number of violent and non-violent actions within the crowd.

Table XIII presents the Comparative Mean Absolute Error (CMAE) results of the proposed BCC model compared with the fine grain crowd counting model of Jia [1] across diverse video datasets. Across various human behaviors such as Standing/Sitting, Waiting/Non-Waiting, Towards/Away, and Violent/Non-Violent categories, the BCC model consistently outperforms the base paper, with lower CMAE values. For instance, in the Standing/Sitting category, the BCC model achieves a CMAE of 5.06, excelling the base paper's 8.01. This trend persists across Waiting/Non-Waiting, Towards/Away, and Violent/Non-Violent categories, with the BCC model demonstrating CMAE values of 4.12, 3.16, and 4.25 respectively, compared to the base paper's 2.99, 2.29, 2.29 and 4.35.

TABLE XIV presents a comprehensive overview of key metrics with crowd behavior analysis across four datasets. Each row corresponds to a specific category: Standing vs. Sitting, Waiting vs. Non-Waiting, Towards vs. Away, and Violent vs Non-Violent. The provided metrics include frames per second, total frames, the date and time of video creation, average object speed, total people count, counts in two designated channels, and several evaluation metrics (CMAE, MAE, MSE, and PATCH). These metrics present insights into crowd behavior's characteristics, accuracy, and dynamics within diverse scenarios and behavioral categories. A contributing factor to the BCC model's enhanced

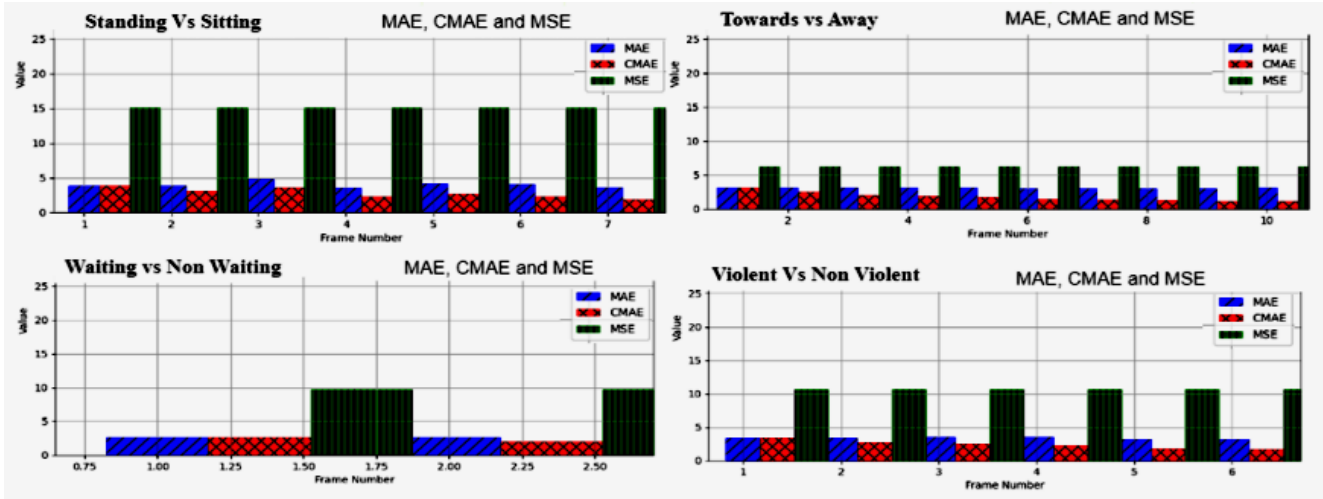


Figure 6. Results of MAE, CMAE, MSE Metric of Four Datasets

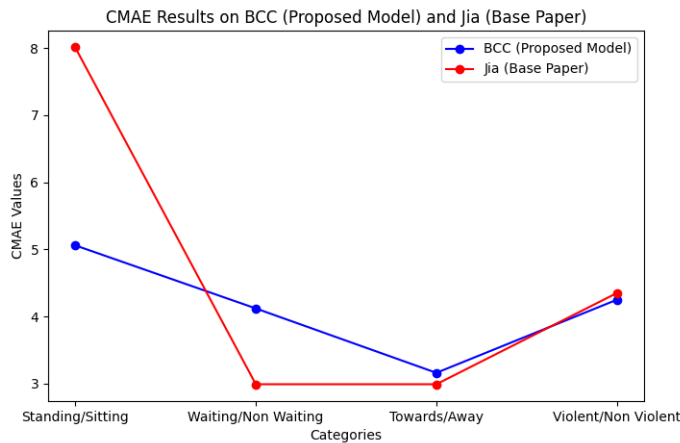


Figure 7. CMAE between BCC and JIA models

accuracy lies in incorporating the Unet architecture. The Unet captures spatial dependencies and intricate patterns in crowd behavior. Its features such as a contracting path for context capture and an expansive path for precise localization, identify even slight differences in scenes with many people. Integrating the Unet with CSRNet architecture boosts the BCC model's ability to analyze complex spatial relationships, emphasizing the importance of an advanced neural network for superior performance in crowd analysis applications.

Figure 5 shows the Comparative Mean Absolute Error (CMAE) results across four categories, each representing different aspects within the visual data. The x-axis denotes the individual categories, while the y-axis the corresponding CMAE values. Lower CMAE values on the graph signify higher accuracy in the model's predictions for each category. The graph compares and highlights the model behaviors within each category based on the calculated

CMAE values. Figure 6 illustrates the outcomes of three key metrics—Mean Absolute Error (MAE), Comparative Mean Absolute Error (CMAE), and Mean Squared Error (MSE) across four distinct video datasets. The x-axis corresponds to the individual datasets, while the y-axis corresponds to each dataset's metric values. The data points visually represent the model performance based on MAE, CMAE, and MSE, with lower values indicating more accurate predictions. The visual summary enables a comparison of model effectiveness across the various video datasets concerning these specific evaluation metrics. Figure 7 presents the performance between the BCC (Proposed Model) and Jia (Base Paper) in predicting crowd behaviors. In Standing/Sitting, BCC boasts a CMAE of 5.06 compared to Jia's 8.01. For Waiting/Non-Waiting, BCC records a CMAE of 4.12, surpassing Jia's 2.99. In Towards/Away, BCC achieves a CMAE of 3.16 against Jia's 2.29. Lastly, in Violent/Non-Violent, BCC demonstrates a CMAE of 4.25, outperforming Jia's 4.35. These numerical comparisons prove the BCC model's accuracy over the Jia model across diverse crowd behavior categories.

## 6. CONCLUSIONS AND FUTURE WORK

The Behavioral Crowd Counting (BCC) architecture combines a Congested Scene Recognition Network (CSRNet) with an Unet for enhanced behavioral crowd counting. The CSRNet consists of a frontend and a backend network for feature extraction and generation of a crowd density map. The Unet produces a density map and refines an attention-based map. It operates on video features and attention maps, refining the density map through several iterations. The refined density maps provide behavior-based crowd segmentation, separating crowd regions from the background with improved accuracy. The experimental results validate the effectiveness of the approach in behaviour crowd counting in video data consisting of congested scenes. This synergy empowers the system to perform



behavioral crowd counting, offering unprecedented insights into crowd dynamics within video datasets.

Extending BCC to recognize and analyze emotions or sentiments within the crowd enables marketing, entertainment, and event management applications. Incorporating multi-modal inputs from different data sources, such as audio, text, or social media data, provides a more comprehensive understanding of crowd behavior and improves analysis accuracy.

## REFERENCES

- [1] J. Wan, N. S. Kumar, and A. B. Chan, "Fine-grained crowd counting," *IEEE Transactions on Image Processing*, vol. 30, pp. 2114–2126, 2021.
- [2] L. Vibha, C. Hegde, P. D. Shenoy, K. R. Venugopal, and L. M. Patnaik, "Dynamic object detection, tracking and counting in video streams for multimedia mining," *IAENG International Journal of Computer Science*, vol. 35, no. 3, pp. 16–21, 2008.
- [3] —, "Moving vehicle identification using background registration technique for traffic surveillance."
- [4] C. Direkoglu, "Abnormal crowd behavior detection using motion information images and convolutional neural networks," *IEEE Access*, vol. 8, pp. 80408–80416, 2020.
- [5] S. Guo, Q. Bai, S. Gao, Y. Zhang, and A. Li, "An analysis method of crowd abnormal behavior for video service robot," *IEEE Access*, vol. 7, pp. 169577–169585, 2019.
- [6] J. Gao, Y. Yuan, and Q. Wang, "Feature-aware adaptation and density alignment for crowd counting in video surveillance," *IEEE Transactions on Cybernetics*, vol. 51, no. 10, pp. 4822–4833, 2020.
- [7] H. Go, J. Byun, B. Park, M. Choi, S. Yoo, and C. Kim, "Fine-grained multi-class object counting," in *2021 IEEE International Conference on Image Processing (ICIP)*, 2021, pp. 509–513.
- [8] Y. Liu, D. Xu, S. Ren, H. Wu, H. Cai, and S. He, "Fine-grained domain adaptive crowd counting via point-derived segmentation," in *IEEE International Conference on Multimedia and Expo (ICME)*, 2023, pp. 2363–2368.
- [9] A. V. Savchenko, "Video-based frame-level facial analysis of affective behavior on mobile devices using efficientnets," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2359–2366.
- [10] J. Kay, C. M. Foley, and T. Hart, "Fine-grained counting with crowd-sourced supervision," *arXiv preprint arXiv:2205.11398*, 2022.
- [11] P. Bour, E. Cribelier, and V. Argyriou, *Crowd Behavior Analysis from Fixed and Moving Cameras*, 2019, pp. 289–322.
- [12] S. Gong, S. Zhang, J. Yang, D. Dai, and B. Schiele, "Bi-level alignment for cross-domain crowd counting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7542–7550.
- [13] S. Bhardwaj, A. Dwivedi, A. Pandey, Y. Perwej, and P. R. Khan, "Machine learning-based crowd behavior analysis and forecasting," *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, vol. 9, no. 3, pp. 418–429, 2023.
- [14] Z. Zou, Y. Cheng, X. Qu, S. Ji, X. Guo, and P. Zhou, "Attend to count: Crowd counting with adaptive capacity multi-scale cnns," *Neurocomputing*, vol. 367, pp. 75–83, 2019.
- [15] Y. Tao, Z. Ling, and I. Patras, "Universal foreground segmentation based on deep feature fusion network for multi-scene videos," *IEEE Access*, vol. 7, pp. 158326–158337, 2019.
- [16] A. Hafeezallah, A. Al-Dhamari, and S. A. R. Abu-Bakar, "U-asd net: Supervised crowd counting based on semantic segmentation and adaptive scenario discovery," *IEEE Access*, vol. 9, pp. 127444–127459, 2021.
- [17] E. B. Varghese and S. M. Thampi, "Towards the cognitive and psychological perspectives of crowd behaviour: A vision-based analysis," *Connection Science*, vol. 33, no. 2, pp. 380–405, 2021.
- [18] X. Jiang, L. Zhang, M. Xu, T. Zhang, P. Lv, B. Zhou, X. Yang, and Y. Pang, "Attention scaling for crowd counting," in *Proceedings of The IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4706–4715.
- [19] N. Ilyas, B. Lee, and K. Kim, "Hadf-crowd: A hierarchical attention-based dense feature extraction network for single-image crowd counting," *Sensors*, vol. 21, no. 10, p. 3483, 2021.
- [20] Y. Zhu, K. Ni, X. Li, A. Zaman, X. Liu, and Y. Bai, "Artificial intelligence aided crowd analytics in rail transit station," *Transportation Research Record*, 2023.
- [21] R. Alharthi, A. Alhothali, B. Alzahrani, and S. Aldaheri, "Massive crowd abnormal behaviors recognition using c3d," in *Proceedings of the 2023 IEEE International Conference on Consumer Electronics (ICCE)*, 2023, pp. 223–226.
- [22] A. Alia, M. Maree, M. Chraibi, A. Toma, and A. Seyfried, "A cloud-based deep learning framework for early detection of pushing at crowded event entrances," *IEEE Access*, vol. 11, p. 45936–45949, 2023. [Online]. Available: <http://dx.doi.org/10.1109/ACCESS.2023.3273770>
- [23] B. Ganga, B. T. Lata, S. Admani, K. R. Venugopal, and L. M. Patnaik, "Generation of high quality density map using uskipgan," in *Proc. of the 2022 IEEE International Conference for Women in Innovation, Technology Entrepreneurship (ICWITE)*, 2022, pp. 1–6.
- [24] B. Ganga, N. Navya Shree, B. T. Lata, and K. R. Venugopal, "Anomalydetectnet: A deep learning framework for anomaly detection in video data," *IJCRT*, vol. 12, no. 1, pp. 1–7, 2024.
- [25] B. Ganga, B. T. Lata, and K. R. Venugopal, "Object detection and crowd analysis using deep learning techniques: Comprehensive review and future directions," *Neurocomputing*, May 2024.
- [26] —, "Ccsa: Crowd counting with stability analysis using adversarial network and cnn," *International Journal of Emerging Technologies and Innovative Research*, vol. 11, no. 3, pp. 170–183, 2024. [Online]. Available: <http://www.jetir.org/papers/JETIR2403B09.pdf>
- [27] B. Ganga, B. T. Lata, T. R. Pallavi, and K. R. Venugopal, "Violent behaviour analysis in crowd," in *2024 Second International Conference on Emerging Trends in Information Technology and Engineering (ICETITE)*. IEEE, Feb 2024, pp. 1–6.
- [28] R. Zhao, D. Dong, Y. Wang, C. Li, Y. Ma, and V. F. Enr iquez, "Image-based crowd stability analysis using improved multi-column

- convolutional neural network,” *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [29] A. N. Shuaibu, A. S. Malik, and I. Fay, “A comparative analysis of techniques for crowd behaviour detection in dense scenes,” *A Journal of Science and Technology*, vol. 4, no. 2, pp. 32–41, 2021.
- [30] M. Yamin, M. M. Almutairi, S. Badghish, and S. Bajaba, “Sparrow search optimization with transfer learning-based crowd density classification,” *Computers, Materials Continua*, vol. 74, no. 3, 2023.
- [31] S. Htiouech, K. Chebil, M. Khemakhem, F. Abed, and M. H. Alkiani, “An extended model for the uavs-assisted multiperiodic crowd tracking problem,” *Complexity*, vol. 2023, no. 3001812, pp. 1–14, 2023.
- [32] F. Abdullah, M. Abdelhaq, R. Alsaqour, M. H. Alatiyyah, K. Alnowaiser, S. S. Alotaibi, and J. Park, “Context aware crowd tracking and anomaly detection via deep learning and social force model,” *IEEE Access*, vol. 11, pp. 75 884–75 898, 2023.
- [33] M. Tzelepi and A. Tefas, “Graph embedded convolutional neural networks in human crowd detection for drone flight safety,” *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 5, no. 2, pp. 191–204, 2021.
- [34] Y. M. Manu, G. K. Ravikumar, and S. V. Shashikala, “An integrated multi-level feature fusion framework for crowd behaviour prediction and analysis,” *Indonesian Journal of Electrical Engineering and Computer Science*, 2023.
-