

CLIP-EBC: CLIP Can Count Accurately through Enhanced Blockwise Classification

Yiming Ma
Computer Science Department
University of Warwick
Coventry, UK
yiming.ma.1@warwick.ac.uk

Victor Sanchez
Computer Science Department
University of Warwick
Coventry, UK
V.F.Sanchez-Silva@warwick.ac.uk

Tanaya Guha
School of Computer Science
University of Glasgow
Glasgow, UK
Tanaya.Guha@glasgow.ac.uk

Abstract—The CLIP (Contrastive Language-Image Pretraining) model has exhibited outstanding performance in recognition problems, such as zero-shot image classification and object detection. However, its ability to count remains understudied due to the inherent challenges of transforming counting—a regression task—into a recognition task. In this paper, we investigate CLIP’s potential in counting, focusing specifically on estimating crowd sizes. Existing classification-based crowd-counting methods have encountered issues, including inappropriate discretization strategies, which impede the application of CLIP and result in suboptimal performance. To address these challenges, we propose the Enhanced Blockwise Classification (EBC) framework. In contrast to previous methods, EBC relies on integer-valued bins that facilitate the learning of robust decision boundaries. Within our model-agnostic EBC framework, we introduce CLIP-EBC, the first fully CLIP-based crowd-counting model capable of generating density maps. Comprehensive evaluations across diverse crowd-counting datasets demonstrate the state-of-the-art performance of our methods. Particularly, EBC can improve existing models by up to 76.9%. Moreover, our CLIP-EBC model surpasses current crowd-counting methods, achieving mean absolute errors of 55.0 and 6.3 on ShanghaiTech part A and part B datasets, respectively. The code will be made publicly available.

I. INTRODUCTION

Crowd counting is concerned with automatic estimation of the number of individuals in images or videos. In recent years, this task has garnered considerable attention due to its potential applications in critical areas, such as managing pandemics [1] and averting crowd collapses [2]. The ability to accurately quantify crowd densities is integral for enhancing public safety, urban planning, and event management.

State-of-the-art crowd-counting approaches rely on the annotated 2D coordinates of individuals’ head centers in images. These methods usually transform point annotations into binary density maps, where a value of 1 indicates that the pixel corresponds to a labeled head center while 0 indicates otherwise. A majority of these methods [3]–[9] adopts an encoder-decoder framework, aiming to directly regress the density maps. Typically, these models output density maps with reduced spatial sizes, determined by a model-specific reduction factor. Each element in the density map estimates the count value in a corresponding block of the image. However, these methods overlook the fact that the count values exhibit

a long-tail distribution, where areas with large values are severely undersampled.

To address this challenge, a few works [11], [12] reframe crowd counting as a classification task by merging count values into bins (classes), and therefore, the sample sizes of rare values can be increased. Similar to their regression-based counterparts, these methods are also based on blockwise prediction but output probability maps of reduced spatial sizes, where the vector at each spatial location represents the probability scores over the bins. During the inference phase, these methods calculate the predicted density map by aggregating mean values of bins, each weighted according to the associated probability score. The final predicted count is then derived by integrating the resulting density map. These methods, however, encounter several challenges that contribute to their underperformance. Notably, like many regression-based methods [3], [4], [8], they also incorporate Gaussian smoothing in preprocessing the ground-truth density maps, introducing a critical concern: the selection of appropriate kernel widths. Given that individuals are often depicted at different scales in images due to perspective distortion, an ideal scenario would involve matching kernel widths with head sizes, which are unfortunately not provided for counting tasks. As a result, existing classification-based methods introduce noise in the labels, leading to a degradation in performance. Moreover, Gaussian smoothing transforms the initially discrete count values into a continuous space $[0, \infty)$, necessitating the utilization of a sequence of bordering real-valued intervals as bins. This quantization policy makes samples near the borders exceptionally challenging to classify, making it difficult for models to learn optimal decision boundaries. Another limitation of current classification-based methods is their sole focus on the classification error without considering the proximity of predicted count values to the ground truth. This deficiency compromises performance in testing, as two probability distributions with identical classification errors may exhibit different expectations.

Despite the success of CLIP [1] in various recognition-related downstream tasks, such as object detection [13] and semantic segmentation [14], its capability to count has remained largely unexplored. This gap arises from two primary challenges:

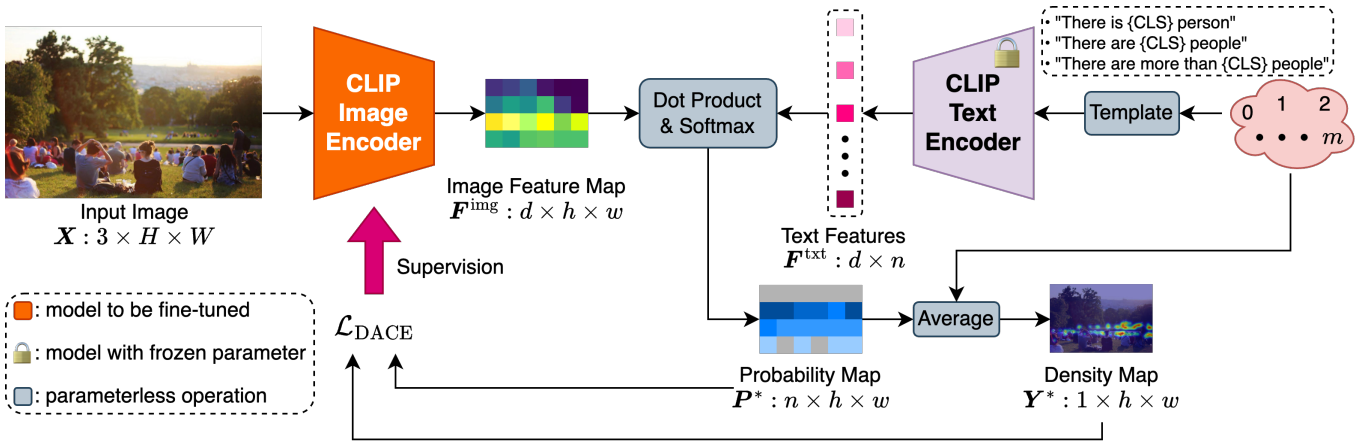


Fig. 1: Overview of our model CLIP-EBC, implemented under our proposed EBC framework with fine-grained bins. We modify the CLIP [10] image encoder by removing the pooling layer and replacing the linear projection layer with a 1×1 convolution. This modification allows the preservation of local information, which plays a vital role in crowd counting. The original text encoder is frozen and employed to extract text embeddings. Then we generate the probability map by calculating the cosine similarities between image feature vectors and text feature vectors. The density map is acquired by averaging the preset count values weighted by the probability scores. Finally, we feed the predicted probability map and the predicted density map, along with the corresponding ground truth, to the loss function $\mathcal{L}_{\text{DACE}}$, which is leveraged to fine-tune the image encoder.

1) the inherent mismatch between CLIP, originally designed for recognition, and counting, which constitutes a regression task, and 2) the limitations and suboptimal results of existing classification-based counting methods. To bridge this gap, we focus on crowd counting in this paper and propose an **Enhanced Blockwise Classification (EBC)** framework, specifically designed to address the challenges faced by current classification-based methods. With minimal modifications, existing regression-based methods can be seamlessly integrated into our EBC framework, resulting in a significant enhancement of their performances. Furthermore, building upon the EBC framework, we explore the potential of leveraging the original structure of CLIP for crowd counting and introduce **CLIP-EBC** (as shown in Figure 1). Compared to other methods [15], [16], CLIP-EBC stands out as *the first fully CLIP-based model capable of generating crowd heat maps*. The experimental results across four databases underscore the substantial improvement offered by our EBC framework over existing regression-based methods. EBC exhibits notable efficacy, achieving a significant boost of up to 76.9% reduced RMSE for CSRNet [3] on the NWPU [17] dataset. Furthermore, our proposed CLIP-EBC model surpasses the state-of-the-art crowd-counting methods, showcasing its effectiveness. Specifically, CLIP-EBC achieves a mean absolute error of 55.0 on ShanghaiTech [18] part A and 6.3 on part B. These results demonstrate that CLIP can accurately estimate crowd density maps with the support of EBC.

To summarize, our contributions are as follows:

- We propose an innovative **Enhanced Blockwise Classification (EBC)** framework that substantially improves the prior classification-based methods in three aspects:

discretization, label correction, and loss function.

- Building upon EBC, we present *the first fully CLIP-based crowd counting model CLIP-EBC*. CLIP-EBC preserves the original structure of CLIP to a maximum extent, demonstrating its capacity not only for estimating crowd sizes but also for generating detailed distribution density maps.
- We conduct extensive experiments across multiple datasets to showcase the effectiveness of EBC in enhancing existing methods and the competence of CLIP-EBC as a state-of-the-art crowd-counting method.

II. RELATED WORK

A. Regression-Based Methods

Crowd counting is dominated by a large variety of encoder-decoder-based models that regress density maps. Some models focus on tackling the scale variation problem caused by perspective distortion. Zhang *et al.*[18] introduce a multi-column CNN structure. Each branch extracts feature maps of different receptive field sizes, which are subsequently fused through concatenation. In contrast, Liu *et al.*[4] propose a VGG-16 [20]-based single-branch model, incorporating a multi-scale module to extract and fuse features across scales. Recognizing the importance of enlarging receptive fields, Li *et al.*[3] advocate for the use of dilated convolutions in generating density maps. Since ground-truth density maps are often sparse, smoothing them with Gaussian kernels is a common strategy to facilitate model optimization. Nevertheless, this approach introduces the challenge of selecting appropriate kernel widths, which ideally should match scales. Unfortunately, head sizes are usually not provided for crowd-counting tasks. Hence, Gaussian smoothing unavoidably introduces errors and noise

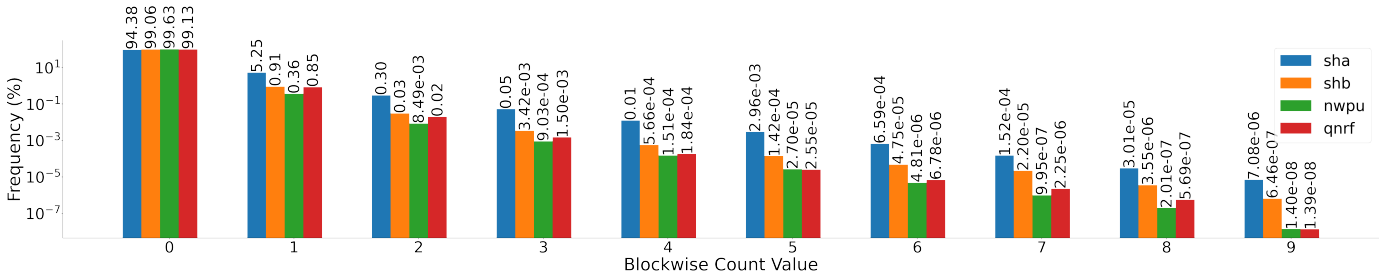


Fig. 2: The distribution of the number of people in 8×8 blocks. This block size is commonly used in many papers [3]–[6]. All datasets, namely ShanghaiTech A & B [18], NWPU-Crowd [17] and UCF-QNRF [19], demonstrate that 0 is the dominant value while other values are severely undersampled.

in labels: if the kernel size is set too small, those pixels that correspond to individual’s heads are set to 0 in the density maps; conversely, if the kernel size is too large, those pixels that correspond to the background can be mistaken as pedestrians. To address this issue, Wang *et al.*[6] introduce the DMCount loss by leveraging the discrete optimal transport theory. This loss function does not require Gaussian smoothing and models trained with it can have enhanced performance.

B. Classification-Based Methods

Existing classification-based approaches emerge from a motivation to rectify the long-tail distribution of count values (as illustrated in Figure 2), where large values are gravely under-represented, adversely affecting the performance of regression-based models. To address this issue, classification-based methods partition the support range $[0, \infty)$ into non-overlapping intervals to enhance the sample size for each class. During inference, the middle points of each interval, weighted by the probability score, are added up as the predicted count. For example, Xiong *et al.*[12] introduce DCNet, which predicts counts at multiple levels by using the same set of bins. However, this method neglects the fact that large values are less likely to appear at a local level, thereby exacerbating class imbalance. To handle this issue, Liu *et al.*[11] propose the concept of blockwise classification with models outputting probability maps, where the vector at each pixel represents the predicted probability scores. However, similar to their regression-based counterparts, these methods also smooth the ground-truth density maps with Gaussian kernels, introducing the following issues: 1) as explained in Sec. II-A Gaussian smoothing can introduce noise in labels; 2). Gaussian smoothing transforms the count from the discrete integer space into the continuous real space, and thus intervals have to border each other (e.g., $(0, 0.5]$ and $(0.5, 1]$). This quantization policy makes it difficult to classify the sample points near the boundaries. Moreover, these methods solely focus on classification results, overlooking the fact that two probability distributions can have the same classification error but different expectations, thus severely impacting the performance during testing.

C. CLIP in Crowd Counting

The Contrastive Language-Image Pre-training (CLIP) model [10] has demonstrated outstanding performance in downstream tasks such as zero-shot image classification [10], object detection [13] and [14], but there are only few studies in crowd counting. Liang *et al.*[15] propose an unsupervised method established on ranking. Their method utilizes sequences of nested image patches and a pre-defined sequence of count numbers as input, with the objective of minimizing the ranking loss of similarities. Nevertheless, one limitation of this approach is its inability to generate density maps, crucial for applications such as pandemic control and public safety monitoring. Jiang *et al.*[16] introduce a text-guided zero-shot counting method, which leverages CLIP to generate text guidance, yet it remains rooted in density-map regression, leading to suboptimal performance in crowd counting tasks. In contrast, we are the first to investigate the potential of using only CLIP to estimate crowd density maps. To this end, we introduce CLIP-EBC, the first crowd-counting method fully based on CLIP, and demonstrate that CLIP can accurately estimate crowd distributions while retaining its structure as much as possible.

III. METHODS

In this section, we first describe our **Enhanced Blockwise Classification (EBC)** framework and our contributions in three aspects: discretization, label correction, and loss design. Then, based on this backbone-agnostic framework, we introduce how to utilize CLIP to in density map estimation and present **CLIP-EBC**.

A. Enhanced Blockwise Classification

Due to the inherent noise present in point labels, predicting the density map at the pixel level becomes challenging. Consequently, our EBC framework resorts to blockwise prediction, similar to the state-of-the-art methods [3], [5], [6], [11], [12]. However, since regression-based methods suffer from undersampling of large count values (as illustrated in Figure 2), EBC groups count values into bins to increase the sample size of each bin, thereby alleviating the problem of sample imbalance. Let $\{\mathcal{B}_i \mid i = 1, \dots, n\}$ be the n pre-defined bins such that $\forall i \neq j, \mathcal{B}_i \cap \mathcal{B}_j = \emptyset$, and $\mathcal{S} \subset \cup_{i=1}^n \mathcal{B}_i$, where

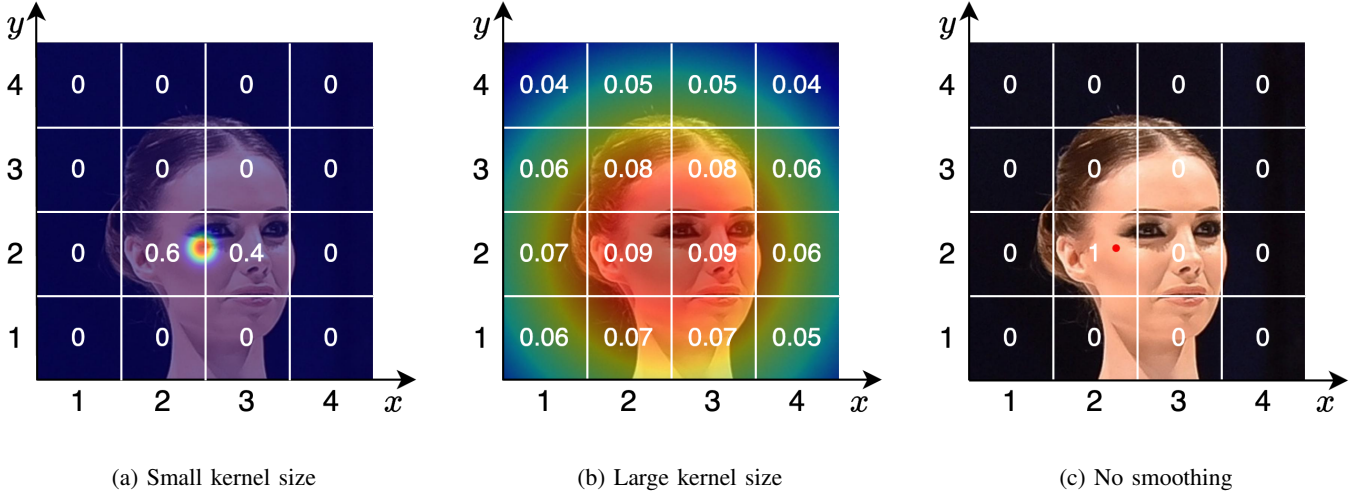


Fig. 3: Figure 3a and Figure 3b illustrate that Gaussian smoothing is sensitive to kernel width selection and can introduce noise in class labels. By comparison, as shown in Figure 3c, we bypass Gaussian smoothing and leverage a YOLO-like [21] method to create integer-based class labels.

\mathcal{S} is the support set of count values. Let $\mathbf{X} \in \mathbb{R}_+^{C \times H \times W}$ denote the input image, where C represents the number of channels, and H and W refer to the spatial height and width, respectively. EBC outputs a probability map \mathbf{P}^* with dimensions $(n, H//r, W//r)$, where $//$ represents the floor division operator and the integer r is a model-related reduction factor. For the n -dimensional vector at the spatial location (i, j) , namely $\mathbf{P}_{:,i,j}^*$, it denotes the probability scores of the bins in the region $(r(i-1) : ri, r(j-1) : rj)$ of the image. During inference, for each bin \mathcal{B}_i , let a_i be a representative count value. From the predicted probability map \mathbf{P}^* , we can obtain the predicted density map via a weighted average:

$$\mathbf{Y}_{i,j}^* = \sum_{k=1}^n a_k \cdot \mathbf{P}_{k,i,j}^*. \quad (1)$$

Summing over $\mathbf{Y}_{i,j}^*$ results in the predicted count of the whole image. Figure 4 compares regression-based methods (left) and our EBC framework (right). Existing regression-based methods can be easily tailored to fit into EBC by only changing the output dimensions.

1) *Discretization*: Following prior regression-based methods [3], [4], [18], existing classification-based methods [11], [12] smooth the ground-truth density maps with Gaussian kernels. This transforms the support set $\mathcal{S} \subset \mathbb{N}$ into a subset of \mathbb{R}_+ . Correspondingly, to cover the new support set, these methods use bordering intervals as bins: $\{0\}, (0, 0.05], (0.05, 0.1] \dots$. This strategy makes samples with labels near the boundaries (e.g., 0.05) difficult to classify. Also, Gaussian smoothing may introduce extra noise in labels as head sizes are usually not provided for crowd counting. As illustrated in Figure 3a and Figure 3b, when the kernel size is not set properly, Gaussian smoothing can create wrong class labels. To address these issues, we propose to bypass Gaussian smoothing and adopt a YOLO-like [21] approach (see Figure 3c): If an individual

is within a specific block, we compel only that block to predict the presence of this individual, while excluding other blocks from making such predictions. This strategy preserves the inherent discreteness of the count. Our support set of count values is $\mathcal{S} = \{0, 1, \dots, m\}$, where m represents the maximum allowable count value. We propose three bin strategies of varying granularity: *fine*, *dynamic*, and *coarse*. At the fine level, each bin contains only one integer; the dynamic bin policy creates bins of various sizes; and at the coarse level, each bin comprises more than one integer.

Another drawback of the previous works[11], [12] is that they employ each interval’s middle point as the representative count value, which ignores that the count values do not follow a uniform distribution. To handle this issue, we propose to use the average count values in each bin as the representative point:

$$a_i = \frac{1}{|\mathcal{B}_i|} \sum_{k=1}^M \mathbb{1}(c_k \in \mathcal{B}_i) \cdot c_k, \quad (2)$$

where $|\mathcal{B}_i|$ is the cardinality of the bin \mathcal{B}_i , M is the number of all blocks in the dataset, $\mathbb{1}$ is the indicator function, and c_k is the count value in block k .

2) *Label Correction*: All aforementioned methods overlook a crucial practical challenge: annotations within densely populated image areas can be exceedingly erroneous and noisy, leading to a significant disparity from the observable number and locations of people (as illustrated in Figure 5). This issue may arise from two factors: 1) images with low resolutions, where annotators struggle to precisely determine the head counts in congested areas; and 2) databases that are resized after labeling to optimize storage and training time. Such errors can provide crowd counting models with incorrect back-propagation signals, severely degrading their actual performance. Therefore, we propose to constrain the maximum count

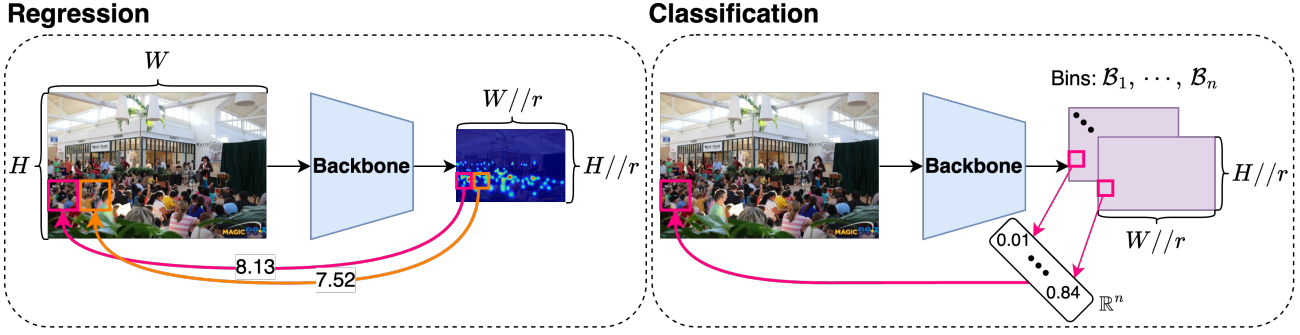


Fig. 4: Blockwise prediction. The left figure represents existing crowd-counting methods that regress density maps. For an input image, these methods usually output a density map of reduced spatial sizes. Each element of the density map predicts the count value in a corresponding block of the image. Our proposed **EBC** framework (the right figure) is also based on blockwise prediction, but it generates probability scores over the pre-defined bins instead of directly regressing the count value.

of observable people in fixed-size image patches to a small constant, solely determined by the patch size. Specifically, we posit that the minimum recognizable size for a person is $s \times s$ pixels. Hence, the maximum allowable count value can be obtained by $m = (r/s)^2$. For example, with $s = 8$, the maximum allowable count value for all 64×64 blocks, including the enclosed region in Figure 5, is restricted to $(64/8)^2 = 64$, rather than 196.

3) *Loss Design*: Among the prior blockwise classification methods, Liu *et al.*[11] solely consider the difference between the predicted probability map \mathbf{P}^* and the ground-truth \mathbf{P} in their loss formulation. Xiong *et al.* [12] introduce an additional term for divide and conquer. These methods, however, neglect the differences between the predicted count value and the ground truth. Since two probability distributions can yield the same classification error but possess different expectations, models trained with these loss functions are not guaranteed to perform well during testing. To cope with this challenge, we propose the **Distance-Aware-Cross-Entropy (DACE)** loss:

$$\begin{aligned} \mathcal{L}_{\text{DACE}} &= \mathcal{L}_{\text{class}}(\mathbf{P}^*, \mathbf{P}) + \lambda \mathcal{L}_{\text{count}}(\mathbf{Y}^*, \mathbf{Y}) \\ &= - \sum_{i=1}^{H//r} \sum_{j=1}^{W//r} \sum_{k=1}^n \mathbb{1}(\mathbf{P}_{k,i,j} = 1) \log \mathbf{P}_{k,i,j}^* \\ &\quad + \lambda \mathcal{L}_{\text{count}}(\mathbf{Y}^*, \mathbf{Y}), \end{aligned} \quad (3)$$

where $\mathbb{1}$ is the indicator function, \mathbf{P} is the one-hot encoded ground-truth probability map, \mathbf{P}^* is the predicted probability map, \mathbf{Y} is the ground-truth density map, and \mathbf{Y}^* is the predicted density map, which can be obtained from Equation (1) and Equation (2). The count loss $\mathcal{L}_{\text{count}}$ (weighted by λ) can be any function that measures the difference between two density maps, but in this paper, we majorly consider utilizing the DMCount Loss [6] since the ground truth is not smoothed.

B. The Structure of CLIP-EBC

Figure 1 depicts the structure of our CLIP-EBC model with fine-grained bins $\{0\}, \{1\}, \dots, \{m\}$. Since the mechanism of

CLIP-EBC is founded on our proposed EBC framework, here we focus solely on how to generate the predicted probability map \mathbf{P}^* . Section III-A provides additional essential details including the generation of crowd density maps, the inference process, and the training loss function.

The image encoder of CLIP-EBC comprises a feature extractor and a 1×1 convolutional layer. As CLIP-EBC is based on blockwise prediction, we remove the final pooling layer and the linear projection layer of CLIP’s image encoder and use the remaining backbone to extract the feature map \mathbf{H} of dimensions $c \times (H//r) \times (W//r)$, where c represents the number of output channels (for ResNet backbone) or the embedding dimension (for ViT backbone). Subsequently, instead of using multiple projection layers, we employ a 1×1 convolutional layer to transform \mathbf{H} into the CLIP embedding space, yielding \mathbf{F}^{img} with dimensions $d \times (H//r) \times (W//r)$, where r denotes the reduction factor (for ResNet) or the patch size (for ViT).

For text feature extraction, we commence by considering the input text prompts. Given a set of bins $\{\mathcal{B}_i \mid i = 1, \dots, n\}$, for each bin \mathcal{B}_i , we generate one text prompt according to the following rules:

- If $\mathcal{B}_i = \{b_i\}$ and $b_i < m$, where m represents the maximum allowable count, the text prompt is ‘‘There is/are b_i person/people’’, with the choice between ‘‘is/are’’ and ‘‘person/people’’ determined by whether $b_i > 1$.
- If \mathcal{B}_i contains more than one element (in this scenario, these elements must be consecutive integers and all smaller than m), then the text prompt becomes ‘‘There is/are between $\min(\mathcal{B}_i)$ and $\max(\mathcal{B}_i)$ person/people’’. The decision between ‘‘is/are’’ and ‘‘person/people’’ is again made to ensure grammatical correctness.
- If $\mathcal{B}_i = m$, then the text prompt is ‘‘There are more than m people’’.



(a) The box-annotated image

(b) Area around the box

(c) Area within the box

Fig. 5: An image containing extremely dense areas. Fig. 5a illustrates the original image, and a congested region with 196 people as the annotated count with the magenta box. Fig. 5b shows the area around the box, and people below the box can still be well recognized. The part within the box is visualized in Fig. 5c, which illustrates the difficulty of identifying the exact locations and number of pedestrians in this area. The image used in this example is IMG_120.jpg from the training split of ShanghaiTech A [18]. The box size is 64×64 , and the upper left corner of the box is located at (646, 301).

Next, the resulting n text prompts are tokenized by CLIP’s tokenizer. Subsequently, we input the tokenized text into the original CLIP text encoder, with its parameters frozen during training. This process yields text embeddings \mathbf{F}^{txt} with dimensions $d \times n$.

With the image feature maps \mathbf{F}^{img} and the text features \mathbf{F}^{txt} , we can then compute the probability map \mathbf{P}^* . First, we calculate the cosine similarity between the image feature vector $\mathbf{F}_{:,i,j}^{\text{img}}$ and the n extracted text embeddings. Subsequently, we normalize these similarities using softmax to obtain the probability $\mathbf{P}_{:,i,j}^*$, which denotes the probability scores across the n bins for block (i, j) . To obtain the predicted density map \mathbf{Y}^* , we use Equation (1), and the loss function utilized to train CLIP-EBC as defined by Equation (3).

IV. EXPERIMENTS

We conduct comprehensive experiments utilizing four publicly available crowd counting datasets: ShanghaiTech A and B [18], UCF-QNRF [19], and NWPU-Crowd [17].

1) *Model Configuration*: For fair comparisons with current approaches, we primarily focus on a block size of $r = 8$. Bilinear interpolation is leveraged to transform the feature maps’ spatial size. We set the minimum recognizable scale to $s = 4$, and hence the maximum allowable count value in each block is $m = (8/4)^2 = 4$. This configuration yields five fine-grained bins: $\{0\}, \{1\}, \{2\}, \{3\}, \{4\}$. Additionally, we explore two other block sizes: $r = 16$ and $r = 32$. For these two options, we consider three sets of bins with varying granularities (see Section IV-B for details).

2) *Training Details*: We initialize the CLIP-EBC model with the weights from CLIP [10]. For the remaining models, we initialize the encoder with pre-trained weights on ImageNet, while the decoder (if applicable) is initialized with random values drawn from a normal distribution. In the proposed

DACE loss (defined by Equation (3)), we set λ to 1, unless specified otherwise. We use the Adam optimizer [22] to train all our models with an initial learning rate of $1e-4$, which is adjusted through a cosine annealing schedule [23]. Throughout the training, we randomly crop patches of size $448u \times 448u$ with $u \sim \text{Uniform}[1, 2]$, which are subsequently resized to 448×448 . This augmentation strategy aims to increase the sample size for larger bins. For CLIP-EBC with the ViT-based image backbone, the input size is 256×256 . The batch size is fixed at 8 for all datasets.

3) *Evaluation metrics*: Following existing methods, we use the mean absolute error (MAE) and root mean square error (RMSE) to evaluate our models. These evaluation metrics are defined as follows:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |C_i - C_i^*|, \quad \text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (C_i - C_i^*)^2}, \quad (4)$$

where N is the number of images in the test set, C_i is the ground-truth global count value of image i , and C_i^* is the predicted count value by integrating over the predicted density map \mathbf{Y}^* . Notably, the lower scores indicate better results.

A. Comparison with State-of-the-Art

We compare EBC and CLIP-EBC with the state-of-the-art crowd-counting methods. Specifically, we reimplement CSRNet [3] and DMCount [6] by only changing their output dimensions to fit them into our proposed EBC framework. The modified models are denoted by CSRNet-EBC and DMCount-EBC, respectively. Table I tabulates the results and demonstrates the effectiveness of our methods. Based on the performance attained by CSRNet-EBC and DMCount-EBC and compared to their regression-based versions, we can conclude that our EBC framework can deliver nontrivial performance

Methods	Venue	SHA		SHB		QNRF		NWPU	
		MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
MCNN [18]	CVPR'16	110.2	173.2	26.4	41.3	277.0	426.0	218.5	700.6
CMTL [24]	AVSS'17	101.3	152.4	20.0	31.1	252	514	-	-
BL [5]	ICCV'19	62.8	101.8	7.7	12.7	88.7	154.8	93.6	470.3
Blockwise [11]	TCSVT'19	62.8	102.0	8.6	16.4	118	192	-	-
S-DCNet [12]	ICCV'19	58.3	95.0	6.7	10.7	104.4	176.1	-	-
DC-Reg [25]	ECCV'22	59.8	100.0	6.8	11.5	84.8	142.3	-	-
CSRNet [3]	CVPR'18	68.2	115.0	10.6	16.0	119.2	211.4	104.8	433.4
CSRNet-EBC (ours)		66.3	105.0	6.9	11.3	79.3	135.8	42.9	100.1
Improvement		2.7%	8.6%	34.9%	29.3%	33.4%	35.7%	45.5%	76.9%
DMCount [6]	NeurIPS'20	59.7	95.7	7.4	11.8	85.6	148.3	70.5	357.6
DMCount-EBC (ours)		62.3	98.9	7.0	10.9	77.2	130.4	39.6	95.8
Improvement		-4.3%	-3.3%	5.4%	7.6%	9.8%	12.0%	43.8%	73.2%
CLIP-EBC (ResNet50, ours)		55.0	88.7	6.3	10.2	80.5	136.6	38.6	90.3
CLIP-EBC (ViT/B-16, ours)		56.6	94.6	8.1	15.1	87.7	159.3	50.2	126.5

TABLE I: Comparison of **EBC** and **CLIP-EBC** with the state-of-the-art crowd counting approaches on ShanghaiTech A (SHA) & B (SHB) [18], UCF-QNRF (QNRF) [19] and NWPU-Crowd (NWPU) [17]. Methods with the “-EBC” suffix are the models slightly modified and trained under our **EBC** framework. The best results are highlighted in bold font.

Block Size	Fine		Dynamic		Coarse	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
$r = 16$	76.31	131.74	75.90	130.48	77.41	130.48
$r = 32$	76.94	130.87	76.06	127.72	79.39	132.51

TABLE II: Effect of block sizes and bin granularity. Experiments are conducted on UCF-QNRF [19]. Results indicate that the dynamic bin strategy outperforms the other due, benefited from its balance between reducing biases in representation count values and increasing class sample sizes.

improvement. Especially on the NWPU validation split, CSRNet and DMCount can be considerably boosted: 45.5% and 43.8% under MAE and 76.9% and 73.2% under RMSE, respectively. Our CLIP-EBC model can also achieve comparable results with the state-of-the-art. Particularly, our ResNet-based CLIP-EBC outperforms existing crowd-counting methods, attaining a 55.0 MAE on ShanghaiTech part A, a 6.3 MAE on part B and 38.6 on NWPU. These results demonstrate that the original CLIP model can be utilized in crowd counting with remarkable performance. On UCF-QNRF and NWPU, our DMCount-EBC has the best performance, with a 77.2 MAE and a 39.7 MAE, respectively.

B. Influence of Bin Granularity

Since the count variable can take different values in larger blocks, we explore the impact of bin granularity on the EBC framework for block sizes $r = 16$ and $r = 32$. We examine three levels of granularity in this experiment. Fine-grained bins are configured to contain only one integer each. Since the representative count value in each bin is the included integer itself, this strategy can provide bins with the lowest biases. For coarse-level bins, each bin comprises two integers, excluding 0, which forms a bin by itself. This approach acknowledges the long-tail distribution of count values and aims to increase the sample size of each bin. In constructing dynamic bins, we

Enhanced Bins	Label Correction	λ	MAE
✗	✗	0.00	140.6
✓	✗	0.00	88.3
✓	✓	0.00	85.8
✓	✓	0.01	83.9
✓	✓	0.10	83.4
✓	✓	1.00	77.9
✓	✓	2.00	82.3

TABLE III: Influence of the components of EBC, where enhanced bins and label correction indicate the use of our proposed discretization policy and label correction method in Section III-A, and λ is the weight for the count loss in Equation (3). Specifically, $\lambda = 0$ indicates that the overall loss function comprises only the classification term.

adopt a strategy that considers small count values as individual bins and combines every two for larger count values.

Table II shows the results on UCF-QNRF. For both $r = 16$ and $r = 32$, the dynamic granularity provides the best performance. This is because it can achieve a good balance between reducing the biases in representative count values and increasing sample sizes. Also, the fine granularity outperforms the coarse granularity in both scenarios, benefitting from less biased representative count values. Additionally, under MAE, using a smaller block size offers better results, since it can create more blocks, and hence, models can learn to better utilize the position information of each person during training.

C. Ablation Studies

In this section, we conduct ablation studies on the three pivotal components of EBC: discretization, label correction, and loss function. The experiments are specifically performed on the UCF-QNRF dataset [19], utilizing VGG-16 [20] as the backbone model. Models established on our discretization policy have fine-grained bins, each comprising only one integer. The

Methods	Backbone			
	VGG16	ResNet101	MobileNetV2	DenseNet201
Blockwise	140.6	127.3	135.9	118.1
EBC (ours)	77.9	79.7	83.2	82.9
Improvement	44.5%	37.3%	38.7%	29.8%

TABLE IV: Comparison between EBC and Blockwise [11] on UCF-QNRF [19] under the MAE metric. Both methods are implemented with the same backbones but different classification frameworks. Results demonstrate that the remarkably enhanced performance of EBC is irrelevant of the backbone choice.

results are summarized in Table III, where “Enhanced Bins” denotes the utilization of our proposed discretization policy, and λ represents the weight assigned to the count loss term in Equation (3). Particularly, when $\lambda = 0$, only the classification loss is employed. The baseline model, Blockwise [11], which classifies the count value into bordering continuous intervals, achieves an MAE of 140.6. However, replacing continuous intervals with integer-valued bins leads to a significant improvement, resulting in an MAE of 88.3—an impressive 37.1% enhancement. This outcome underscores that, in our scenario, the decision boundaries can be more effectively learned. Further constraining the maximum allowable count value in labels results in a decrease in MAE of 85.8. Regarding λ , increasing its value from 0.00 to 1.00 results in improved performance, with the best result achieved at $\lambda = 1$ (77.9). However, a larger value for λ (e.g., 2) compromises the model’s generalizability.

To verify that the performance improvement is independent of the backbone, we also test our EBC framework with ResNet [26], MobileNetV2 [27] and DenseNet [28]. Results in Table IV demonstrate that compared with the baseline blockwise classification method [11], EBC always achieves important performance improvements, regardless of the backbone choice.

D. Visualization

We utilize our trained CLIP-EBC model with the ResNet image backbone to generate predicted density maps on the ShanghaiTech [18] dataset. A carefully curated set of six representative images is chosen for visualization purposes—three from Part A and three from Part B. Figure 6 presents these six images in the top row, the ground-truth density maps (smoothed by Gaussian kernels for illustration) in the middle row, and the predicted density maps (resized to the images’ sizes) in the bottom row. The selection of these six images ensures comprehensive coverage across a spectrum of crowd densities, ranging from the most sparse scenarios (e.g., 15 people, left column) to the most congested cases (e.g., 1111 people, right column). These results collectively demonstrate the model’s robust performance across a diverse range of crowd density levels.

V. CONCLUSION

In this paper, we demonstrated the capability of CLIP in accurate crowd density estimation. We narrowed the gap

between CLIP and crowd counting by reformulating counting as a blockwise classification problem and proposed an Enhanced Blockwise Classification framework (EBC). Within the EBC framework, we further propose CLIP-EBC, the first fully CLIP-based crowd-counting method with density map generation. To classify local count values into pre-defined bins, CLIP-EBC compares the similarities between the corresponding local image feature with each text feature and then uses softmax on them to generate probability scores. Experiments on multiple databases showcase the effectiveness of EBC and CLIP-EBC. In the future, we will investigate using CLIP-EBC to count any objects to fully realize the potential of CLIP.

1) *Impact Statement*: Like many other crowd-counting methods, our models also analyze crowd images, which may raise privacy concerns. In real-world applications, people may feel uncomfortable knowing that they are being monitored and counted without their explicit consent. Also, some datasets may be biased, and our models trained on them may exhibit biases in predictions. This could result in certain demographic groups being underrepresented or overrepresented, leading to unfair consequences or reinforcing existing social disparities.

2) *Limitations*: CLIP has been pre-trained on millions of image-text pairs and thus should be able to count objects of any kind, but only humans are focused in this paper. In the future, we plan to explore its applications in counting other objects.

REFERENCES

- [1] I. J. C. Valencia, E. P. Dadios, A. M. Fillone, J. C. V. Puno, R. G. Baldovino, and R. K. C. Billones, “Vision-based crowd counting and social distancing monitoring using tiny-yolov4 and deepsort,” in *2021 IEEE International Smart Cities Conference (ISC2)*, IEEE, 2021, pp. 1–7.
- [2] S. A. M. Saleh, S. A. Suandi, and H. Ibrahim, “Recent survey on crowd density estimation and counting for visual surveillance,” *Engineering Applications of Artificial Intelligence*, vol. 41, pp. 103–114, 2015.
- [3] Y. Li, X. Zhang, and D. Chen, “Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1091–1100.
- [4] W. Liu, M. Salzmann, and P. Fua, “Context-aware crowd counting,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5099–5108.
- [5] Z. Ma, X. Wei, X. Hong, and Y. Gong, “Bayesian loss for crowd count estimation with point supervision,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6142–6151.
- [6] B. Wang, H. Liu, D. Samarasinghe, and M. H. Nguyen, “Distribution matching for crowd counting,” *Advances in neural information processing systems*, vol. 33, pp. 1595–1607, 2020.

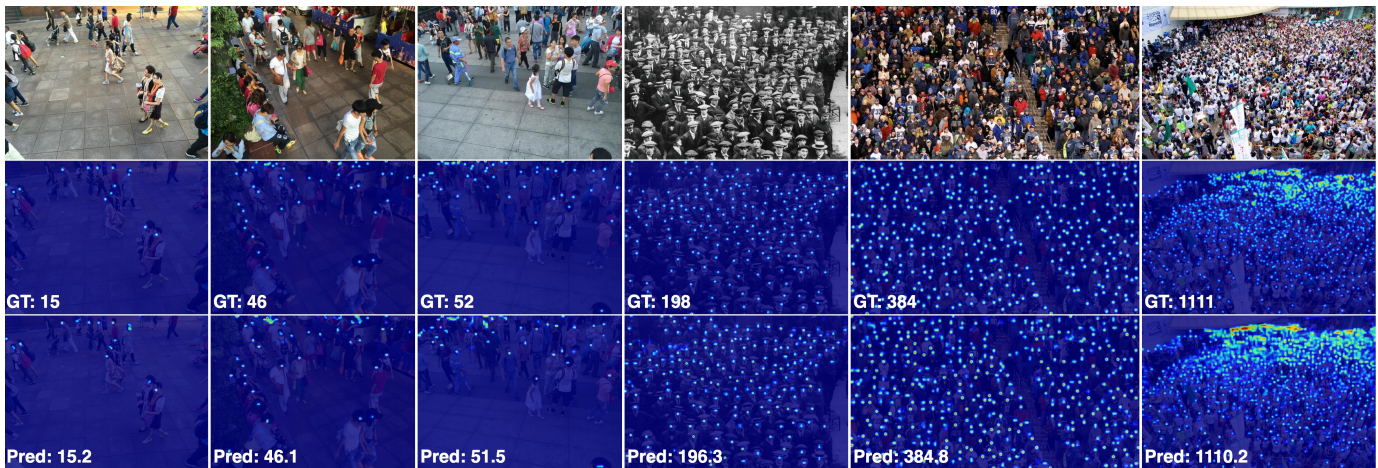


Fig. 6: Visualization on the ShanghaiTech [18] dataset. The top row shows the input crowd images; the middle row visualizes the ground-truth density maps; density maps predicted by our CLIP-EBC (ResNet50) are present in the bottom row. Top to bottom row: original images, the Gaussian-smoothed ground-truth density maps, and the predicted density maps. The left three columns represent sparsely populated cases (from part B), while the right three columns are examples of overcrowded scenarios (from part A).

- [7] P. Thanasutives, K.-i. Fukui, M. Numao, and B. Kijirikul, “Encoder-decoder based convolutional neural networks with multi-scale-aware modules for crowd counting,” in *2020 25th international conference on pattern recognition (ICPR)*, IEEE, 2021, pp. 2382–2389.
- [8] Q. Song, C. Wang, Y. Wang, *et al.*, “To choose or to fuse? scale selection for crowd counting,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, 2021, pp. 2576–2583.
- [9] Y. Ma, V. Sanchez, and T. Guha, “Fusioncount: Efficient crowd counting via multiscale feature fusion,” in *2022 IEEE International Conference on Image Processing (ICIP)*, IEEE, 2022, pp. 3256–3260.
- [10] A. Radford, J. W. Kim, C. Hallacy, *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*, PMLR, 2021, pp. 8748–8763.
- [11] L. Liu, H. Lu, H. Xiong, K. Xian, Z. Cao, and C. Shen, “Counting objects by blockwise classification,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 10, pp. 3513–3527, 2019.
- [12] H. Xiong, H. Lu, C. Liu, L. Liu, Z. Cao, and C. Shen, “From open set to closed set: Counting objects by spatial divide-and-conquer,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 8362–8371.
- [13] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, “Open-vocabulary object detection via vision and language knowledge distillation,” in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=IL3lnMbr4WU>.
- [14] T. Lüddecke and A. Ecker, “Image segmentation using text and image prompts,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7086–7096.
- [15] D. Liang, J. Xie, Z. Zou, X. Ye, W. Xu, and X. Bai, “Crowdclip: Unsupervised crowd counting via vision-language model,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 2893–2903.
- [16] R. Jiang, L. Liu, and C. Chen, “Clip-count: Towards text-guided zero-shot object counting,” *arXiv preprint arXiv:2305.07304*, 2023.
- [17] Q. Wang, J. Gao, W. Lin, and X. Li, “Nwpu-crowd: A large-scale benchmark for crowd counting and localization,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 6, pp. 2141–2149, 2020.
- [18] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma, “Single-image crowd counting via multi-column convolutional neural network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 589–597.
- [19] H. Idrees, M. Tayyab, K. Athrey, *et al.*, “Composition loss for counting, density map estimation and localization in dense crowds,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 532–546.
- [20] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *3rd International Conference on Learning Representations (ICLR 2015)*, Computational and Biological Learning Society, 2015.
- [21] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on*

computer vision and pattern recognition, 2016, pp. 779–788.

- [22] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [Online]. Available: <http://arxiv.org/abs/1412.6980>.
- [23] I. Loshchilov and F. Hutter, “SGDR: Stochastic gradient descent with warm restarts,” in *International Conference on Learning Representations*, 2017. [Online]. Available: <https://openreview.net/forum?id=Skq89Scxx>.
- [24] V. A. Sindagi and V. M. Patel, “Cnn-based cascaded multi-task learning of high-level prior and density estimation for crowd counting,” in *2017 14th IEEE international conference on advanced video and signal based surveillance (AVSS)*, IEEE, 2017, pp. 1–6.
- [25] H. Xiong and A. Yao, “Discrete-constrained regression for local counting models,” in *European Conference on Computer Vision*, Springer, 2022, pp. 621–636.
- [26] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [27] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [28] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.