

Completely Self-Supervised Crowd Counting via Distribution Matching

Deepak Babu Sam*, Abhinav Agarwalla*, Jimmy Joseph, Vishwanath A. Sindagi, R. Venkatesh Babu, *Senior Member, IEEE*, Vishal M. Patel, *Senior Member, IEEE*

Abstract—Dense crowd counting is a challenging task that demands millions of head annotations for training models. Though existing self-supervised approaches could learn good representations, they require some labeled data to map these features to the end task of density estimation. We mitigate this issue with the proposed paradigm of *complete self-supervision*, which does not need even a single labeled image. The only input required to train, apart from a large set of unlabeled crowd images, is the approximate upper limit of the crowd count for the given dataset. Our method dwells on the idea that natural crowds follow a power law distribution, which could be leveraged to yield error signals for backpropagation. A density regressor is first pretrained with self-supervision and then the distribution of predictions is matched to the prior by optimizing Sinkhorn distance between the two. Experiments show that this results in effective learning of crowd features and delivers significant counting performance. Furthermore, we establish the superiority of our method in less data setting as well. The code and models for our approach is available at <https://github.com/val-iisc/css-ccnn>.

Index Terms—Unsupervised Learning, Crowd Counting, Deep Learning



1 INTRODUCTION

THE ability to estimate head counts of dense crowds effectively and efficiently serves several practical applications. This has motivated deeper research in the field and resulted in a plethora of crowd density regressors. These CNN based models deliver excellent counting performance almost entirely on the support of fully supervised training. Such a data hungry paradigm is limiting the further development of the field as it is practically infeasible to annotate thousands of people in dense crowds for every kind of setting under consideration. The fact that current datasets are relatively small and cover only limited scenarios, accentuates the necessity of a better training regime. Hence, developing methods to leverage the easily available unlabeled data, has gained attention in recent times.

The classic way of performing unsupervised learning revolves around autoencoders ([1], [2], [3], [4]). Autoencoders or its variants are optimized to predict back their inputs, usually through a representational bottleneck. By doing so, the acquired features are generic enough that they could be employed for solving other tasks of interest. These methods have graduated to the more recent framework of self-supervision, where useful representations are learned by performing some alternate task for which pseudo labels can be easily obtained. For example, in self-supervision with

colorization approach ([5], [6], [7]), a model is trained to predict the color image given its grayscale version. One can easily generate grayscale inputs from RGB images. Similarly, there are lots of tasks for which labels are freely available like predicting angle of rotation from an image ([8], [9]), solving jumbled scenes [10], inpainting [11] etc. Though self-supervision is effective in learning useful representations, they require a final mapping from the features to the end task of interest. This is thought to be essentially unavoidable as some supervisory signal is necessary to aid the final task. For this, typically a linear layer or a classifier is trained on top of the learned features using supervision from labeled data, defeating the true purpose of self-supervision. In the case of crowd counting, one requires training with annotated data for converting the features to a density map. To reiterate, the current unsupervised approaches might capture the majority of its features from unlabeled data, but demand supervision at the end should they be made useful for any practical applications.

Our work emerges precisely from the above limitation of the standard self-supervision methods, but narrowed down to the case of crowd density estimation. The objective is to eliminate the mandatory final labeled supervision needed for mapping the learned self-supervised features to a density map output. In other words, we mandate developing a model that can be trained without using any labeled data. Such a problem statement is not only challenging, but also ill-posed. Without providing a supervisory signal, the model cannot recognize the task of interest and how to properly guide the training stands as the prime issue. We solve this in a novel manner by carefully aiding the model to regress crowd density on the back of making some crucial assumptions. The idea relies on the observation that natural crowds tend to follow certain long tailed statistics and could be approximated to an appropriate parametric prior distribution (Section 3.1). If a network trained with a self-supervised

*Manuscript received September 8, 2020; revised **** *, ***, accepted **** **, ****. This work was supported by the Science and Engineering Research Board (SERB), Department of Science and Technology (DST), Government of India under the project SB/S3/EECE/0127/2015. (D. B. Sam and A. Agarwalla contributed equally to this work.) (Corresponding authors: D. B. Sam; A. Agarwalla.)*

D. B. Sam, A. Agarwalla, J. Joseph and R. V. Babu are with the Video Analytics Lab, Department of Computational and Data Sciences, Indian Institute of Science, Bangalore, India (e-mail: deepaksam@iisc.ac.in; agarwalla-abhinav@gmail.com; jimmyj005@gmail.com; venky@iisc.ac.in).

V. A. Sindagi and V. M. Patel are with the Vision & Image Understanding Lab, Department of Electrical and Computer Engineering, Johns Hopkins University, Baltimore, USA (e-mail: vishwanathsindagi@jhu.edu; vpatel36@jhu.edu).

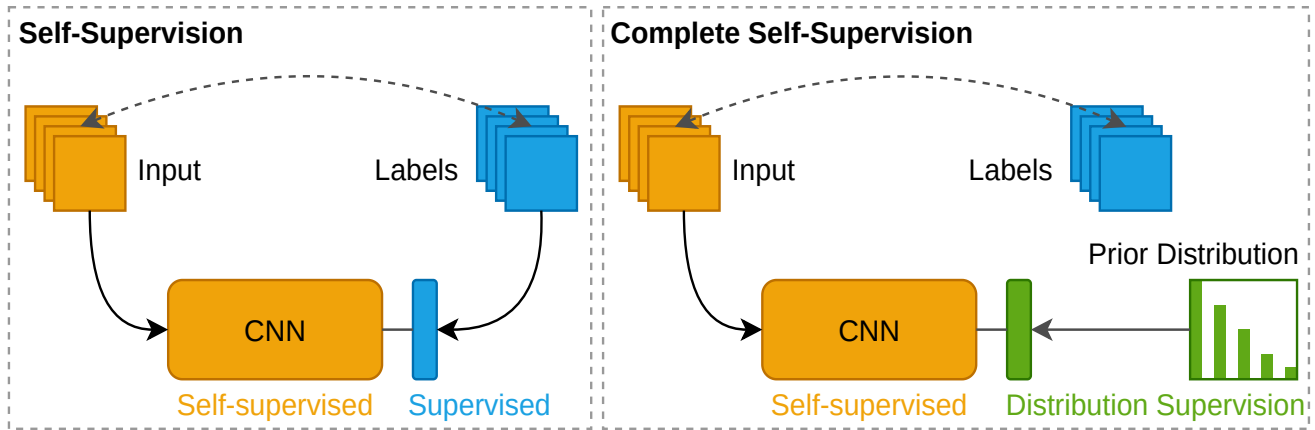


Fig. 1. Self-Supervision Vs Complete Self-Supervision: Normal self-supervision techniques has a mandatory labeled training stage to map the learned features to the end task of interest (in blue). But the proposed complete self-supervision is devoid of such an instance-wise labeled supervision, instead relies on matching the statistics of the predictions to a prior distribution (in green).

task is available (Section 3.2), its features can be faithfully mapped to crowd density by enforcing the predictions to match the prior distribution (Section 3.3). The matching is measured in terms of Sinkhorn distance [12], which is differentiated to derive error signals for supervision. This proposed framework is contrasted against the normal self-supervision regime in Figure 1, with the central difference being the replacement of the essential labeled training at the end by supervision through distribution matching. We show that the proposed approach results in effective learning of crowd features and delivers good performance in terms of counting metrics (Section 4).

In summary, our work contributes the following:

- The first *completely self-supervised* training paradigm which does not require instance-wise annotations, but works by matching statistics of the distribution of labels.
- The first *crowd counting model* that can be trained without using a single annotated image, but delivers significant regression performance.
- A detailed analysis on the distribution of persons in dense crowds to reveal the *power law nature* and enable the use of optimal transport framework.
- A novel extension of the proposed approach to *semi-supervised setting* that can effectively exploit unlabeled data and achieve significant gains.
- An efficient way to improve the Sinkhorn loss by leveraging *edge information* from crowd images.

2 RELATED WORK

The paradigm of dense crowd counting via density regression plausibly begins with [14], where hand-crafted features and frequency analysis are employed. With the advent of deep learning, many CNN based density regressors have emerged. It ranges from the initial simple models [15] to multi-network/multi-scale architectures designed specifically to address the drastic diversity in crowd images ([13], [16], [17], [18], [19]). Regressors with better, deeper and recurrent based deep models ([20], [21], [22]) are shown to improve counting performance. An alternate line of works enhance density regression by providing auxiliary information

through crowd classification ([23], [24]), scene context ([25], [26], [27]), perspective data ([28], [29]), attention ([30], [31], [32]) and even semantic priors [33]. Models designed to progressively predict density maps and perform refinement is explored in ([34], [35], [36]). Works like ([37], [38]) effectively fuse multi-scale information. A better ground truth density map should result in better regression and ([39], [40]) leverage such refinement opportunities. Some recent approaches try to bring flavours of detection to crowd counting ([41], [42], [43], [44], [45]). Developing alternate loss functions is also an area of focus with multi-task loss formulations like ([46], [47]) and probabilistic training regimes as in [48]. Some counting works employ Negative Correlation Learning [49], adversarial training [50] and divide-and-conquer approach [51]. Interestingly, all these works are fully supervised and leverage annotated data to achieve good performance. The issue of annotation has drawn attention of a few works in the field and is mitigated via multiple means. A count ranking loss on unlabeled images is employed in a multi-task formulation along with labeled data by [52]. Wang et al. [53] train using labeled synthetic data and adapt to real crowd scenario. The autoencoder method proposed in [54] optimizes almost 99% of the model parameters with unlabeled data. However, all of these models require some annotated data (either given by humans or obtained through synthetic means) for training, which we aim to eliminate.

Our approach is not only new to crowd counting, but also kindles alternate avenues in the area of unsupervised learning as well. Though initial works on the subject employ autoencoders or its variants ([1], [2], [3], [4]) for learning useful features, the paradigm of self-supervision with pseudo labels stands out to be superior in many aspects. Works like ([5], [6], [7]), learn representations through colourising a grayscale image. Apart from these, pseudo labels for supervision are computed from motion cues ([55], [56], [57]), temporal information in videos ([58], [59]), learning to inpaint [11], co-occurrence [60], spatial context ([10], [61], [62]), cross-channel prediction [63], spotting artifacts [64], predicting object rotation ([8], [9]) etc. The recent work of Zhang et al. [65] introduce the idea of auto-encoding transformations rather than data. An extensive and rigorous comparison of

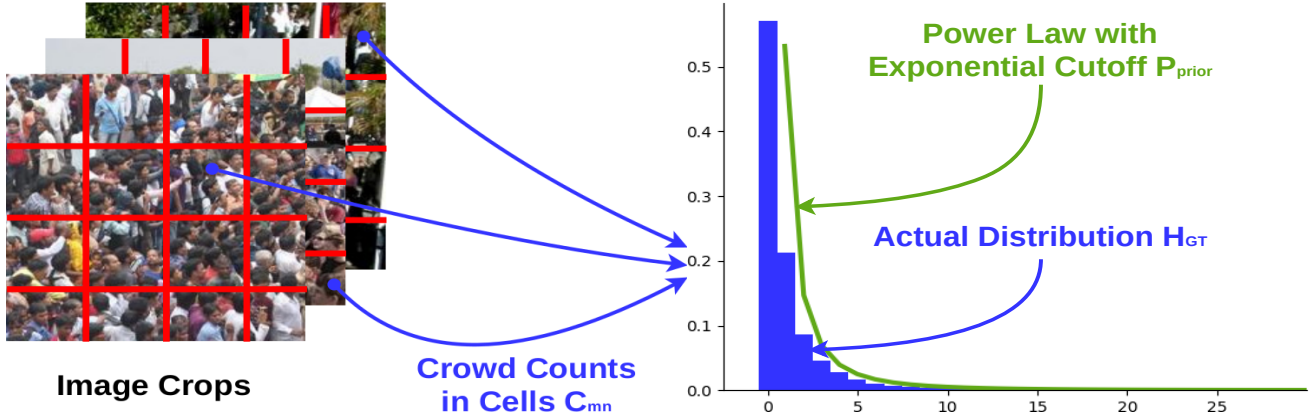


Fig. 2. Computing the distribution of natural crowds: crops from dense crowd images are framed to a spatial grid of cells and crowd counts of all the cells are aggregated to a histogram (obtained on ShanghaiTech Part_A dataset [13]). The distribution is certainly long tailed and could be approximated to a power law.

all major self-supervised methods is available in [66]. All these approaches focus on learning generic features and not the final task. But we extend the self-supervision paradigm further directly to the downstream task of interest.

3 OUR APPROACH

3.1 Natural Crowds and Density Distribution

As mentioned in Section 1, our objective of training a density regressor without using any annotated data is somewhat ill-posed. The main reason being the absence of any supervisory signal to guide the model towards the task of interest, which is the density estimation of crowd images. But this issue could be circumvented by effectively exploiting certain structure or pattern specific to the problem. In the case of crowd images, restricting to only dense ones, we deduce an interesting pattern on the density distribution. They seem to spread out following a power law. To see this, we sample fixed size crops from lots of dense crowd images and divide each crop into a grid of cells as shown in Figure 2. Then the number people in every cell is computed and accumulated to a histogram. The distribution of these cell counts is quite clearly seen to be long tailed, with regions having low counts forming the head and high counts joining the tail. The number of cell regions with no people has the highest frequency, which then rapidly decays as the crowd density increases. This resembles the way natural crowds are arranged with sparse regions occurring more often than rarely forming highly dense neighborhoods. Coincidentally, it has been shown that many natural phenomena obey a similar power law and is being studied heavily [67]. The dense crowds also, interestingly, appears conforming to this pattern as evident from multiple works ([68], [69], [70], [71] etc.) on the dynamics of pedestrian gatherings.

Moving to a more formal description, if \mathbf{D} represents the density map for the input image \mathbf{I} , then the crowd count is given by $C = \sum_{xy} \mathbf{D}_{xy}$ (please refer ([13], [14], [17] etc.) regarding creation of density maps). \mathbf{D} is framed into a grid of $M \times N$ (typically set as $M = N = 3$) cells, with C_{mn} denoting the crowd count in the cell indexed by (m, n) . Now let H^{GT} be the histogram computed by collecting the cell counts (C_{mn} s) from all the images. We

try to find a parametric distribution that approximately follows H^{GT} with special focus to the long tailed region. The power law with exponential cut-off seems to be better suited (see Figure 2). Consequently, the crowd counts in cells C_{mn} could be thought as being generated by the following relation,

$$C_{mn} \sim P_{prior}(c) \propto c^\alpha \exp(-\lambda c), \quad (1)$$

where P_{prior} is the substitute power law distribution. There are two parameters to P_{prior} with α controlling the shape and λ setting the tail length.

Our approach is to fix a prior distribution so that it can be enforced on the model predictions. Studies like ([68], [69], etc.) simulate crowd behaviour dynamics and estimate the exponent of the power law to be around 2. Empirically, we also find that $\alpha = 2$ works in most cases of dense crowds, with the only remaining parameter to fix is the λ . Observe that λ affects the length of the tail and directly determines the maximum number of people in any given cell. If the maximum count C^{max} is specified for the given set of crowd images, then λ could be fixed such that the cumulative probability density (the value of CDF) of P_{prior} at C^{max} is very close to 1. We assume $1/S$ as the probability of finding a cell with count C^{max} out of S images in the given set. Now the CDF value at C^{max} could be set to $1 - 1/S$, simply the probability for getting values less than the maximum. Note that C^{max} need not be exact as small variations do not change P_{prior} significantly. This makes it practical as the accurate maximum count might not be available in real-world scenarios. Since C^{max} is for the cells, the maximum crowd count of the full image C^{fmax} is related as $C^{max} = C^{fmax}/(MNS_{crop})$, where S_{crop} denotes the average number of crops that make up a full image (and is typically set as 4). Thus, for a given a set of highly dense images, only one parameter, the C^{fmax} is required to fix an appropriate prior distribution.

We make a small modification to the prior distribution P_{prior} as its value range starts from 1. H^{GT} has values from zero with large probability mass concentrated near the low count region. Roughly 30% of the mass is seen to be distributed for counts less than or around 1. So, that much probability mass near the head region of P_{prior} is redis-

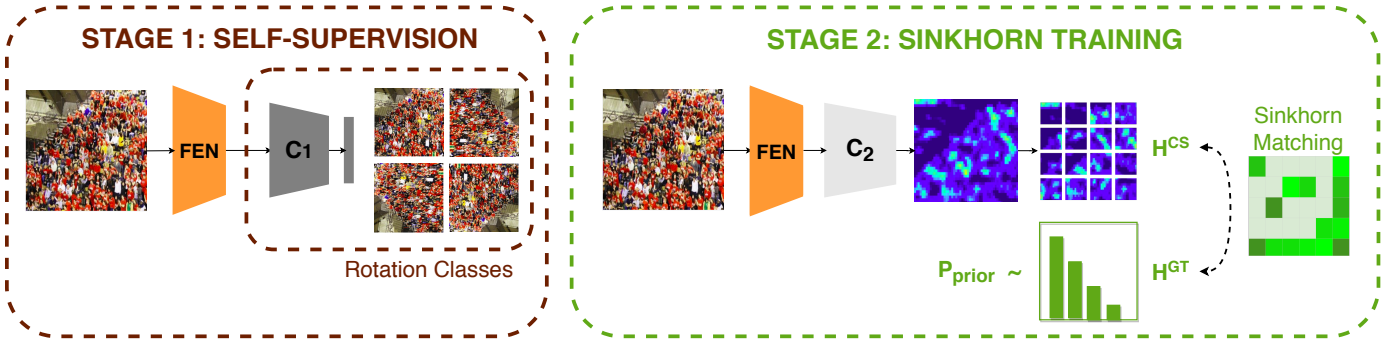


Fig. 3. The architecture of CSS-CCNN is shown. CSS-CCNN has two stages of training: the first trains the base *feature extraction network* in a self-supervised manner with rotation task and the second stage optimizes the model for matching the statistics of the density predictions to that of the prior distribution using optimal transport.

tributed to $[0, 1]$ range in a uniform manner. This is found to be better for both training stability and performance.

In short, now we have a prior distribution representing how the crowd density is being allocated among the given set of images. Suppose there exists a CNN model that can output density maps, then one could try to generate error signals for updating the parameters of the model by matching the statistics of the predictions with that of the prior. But that could be a very weak signal for proper training of the model. It would be helpful if the model has a good initialization to start the supervision by distribution matching, which is precisely what we do by self-supervision in the next section.

3.2 Stage 1: Learning Features with Self-Supervision

We rely on training the model with self-supervision to learn effective and generic features that could be useful for the end task of density estimation. That means the model has to be trained in stages, with the first stage acquiring patterns frequently occurring in the input images. Since only dense crowd images are fed, we hope to learn mostly features relevant to crowds. These could be peculiar edges discriminating head-shoulder patterns formed by people to fairly high-level semantics pertaining to crowds. Note that the model is not signaled to pick up representations explicitly pertinent to density estimation, but implicitly culminate in learning crowd patterns as those are the most prominent part of the input data distribution. Hence, the features acquired by self-supervision could serve as a faithful initialization for the second stage of distribution matching.

Regarding self-supervision, there are numerous ways to generate pseudo labels for training models. The task of predicting image rotations is a simple, but highly effective for learning good representations [66]. The basic idea is to randomly rotate an image and train the model to predict the angle of rotation. By doing so, the network learns to detect characteristic edges or even fairly high-level patterns of the objects relevant for determining the orientation. These features are observed to be generic enough for diverse downstream tasks [66] and hence we choose self-supervision through rotation as our method.

Figure 3 shows the architecture of our density regressor, named the CSS-CCNN (for *Completely Self-Supervised Counting CNN*). It has a base *Feature Extraction Network* (FEN),

which is composed of three VGG [72] style convolutional blocks with max poolings in-between. This is followed by two task heads: C_1 for the first training stage of self-supervision, and C_2 for regressing crowd density at second stage. The first stage branch has two more convolutions and a fully connected layer to finally classify the input image to one of the rotation classes. We take 112×112 crops from crowd images and randomly rotate the crop by one of the four predefined angles (0, 90, 180, 270 degrees). The model is trained with cross-entropy loss between the predicted and the actual rotation labels. The optimization runs till saturation as evaluated on a validation set of images.

Once the training is complete, the FEN has learned useful features for density estimation and the rotation classification head is removed. Now the parameters of FEN are frozen and is ready to be used in the second stage of training through matching distributions.

3.3 Stage 2: Sinkhorn Training

After the self-supervised training stage, FEN is extended to a density regressor by adding two convolutional layers as shown in Figure 3. We take features from both second and third convolution blocks for effectively mapping to crowd density. This aggregates features from slightly different receptive fields and is seen to deliver better performance. The layers of FEN are frozen and only a few parameters in the freshly added layers are open for training in the second stage of distribution matching. This particularly helps to prevent over-fitting as the training signal generated could be weak for updating large number of parameters. Now we describe the details of the exact matching process.

The core idea is to compute the distribution of crowd density predicted by CSS-CCNN and optimize the network to match that closely with the prior P_{prior} . For this, a suitable distance metric between the two distributions should be defined with differentiability as a key necessity. Note that the predicted distribution is in the form of an empirical measure (an array of cell count values) and hence it is difficult to formulate an easy analytical expression for the computing similarity. The classic Earth Mover’s Distance (EMD) measures the amount of probability mass that needs to be moved if one tries to transform between the distributions (also described as the optimal transport cost). But this is not a differentiable operation per se and cannot be used

directly in our case. Hence, we choose the Sinkhorn distance formulation proposed in [12]. Sinkhorn distance between two empirical measures is proven to be an upper bound for EMD and has a differentiable implementation. Moreover, this method performs favorably in terms of efficiency and speed as well.

Let D^{CS} represent the density map output by CSS-CCNN and C^{CS} hold the cells extracted from the predictions. To make the distribution matching statistically significant, a batch of images are evaluated to get the cell counts (C_{mn}^{CS}), which are then formed into an array H^{CS} . We also sample the prior P_{prior} and create another empirical measure H^{GT} to act as the ground truth. Now the Sinkhorn loss \mathcal{L}_{sink} is computed between H^{GT} and H^{CS} . It is basically a regularized version of optimal transport (OT) distance for the two sample sets. Designate \mathbf{h}^{GT} and \mathbf{h}^{CS} as the probability vectors (summing to 1) associated with the empirical measures H^{GT} and H^{CS} respectively. Now a transport plan \mathbf{P} could be conceived as the joint likelihood of shifting the probability mass from \mathbf{h}^{GT} to \mathbf{h}^{CS} . Define U to be the set of all such valid candidate plans as,

$$U = \{\mathbf{P} \in \mathbb{R}_+^{d \times d} \mid \mathbf{P}\mathbf{1} = \mathbf{h}^{GT}, \mathbf{P}^T\mathbf{1} = \mathbf{h}^{CS}\}. \quad (2)$$

There is a cost M associated with any given transport plan, where M_{ij} is the squared difference between the counts of i th sample of H^{GT} and j th of H^{CS} . Closer the two distribution, lower would be the cost for transport. Hence, the Sinkhorn loss \mathcal{L}_{sink} is defined as the cost pertinent to the optimal transportation plan with an additional regularization term. Mathematically,

$$\mathcal{L}_{sink}(H^{GT}, H^{CS}) = \arg \min_{\mathbf{P} \in U} \langle \mathbf{P}, M \rangle_F - \frac{1}{\beta} E(\mathbf{P}), \quad (3)$$

where $\langle \cdot \rangle_F$ stands for the Frobenius inner product, $E(\mathbf{P})$ is the entropy of the joint distribution \mathbf{P} and β is a regularization constant (see [12] for more details). It is evident that minimizing \mathcal{L}_{sink} brings the two distributions closer in terms of how counts are allotted.

The network parameters are updated to optimize \mathcal{L}_{sink} , thereby bringing the distribution of predictions close to that of the prior. At every iteration of the training, a batch of crowd images are sampled from the dataset and empirical measures for the predictions as well as prior are constructed to backpropagate the Sinkhorn loss. The value of \mathcal{L}_{sink} on a validation set of images is monitored for convergence and the training is stopped if the average loss does not improve over a certain number of epochs. Note that we do not use any annotated data even for validation. The counting performance is evaluated at the end with the model chosen based on the best mean validation Sinkhorn loss.

Thus, our Sinkhorn training procedure does not rely on instance-level supervision, but exploits matching the statistics computed from a set of inputs to that of the prior. One criticism regarding this method could be that the model need not learn the task of crowd density estimation by optimizing the Sinkhorn loss. It could learn any other arbitrary task that follows a similar distribution. The counter-argument stems from the semantics of the features learned by the base network. Since the initial training mostly captures features related to dense crowds (see Section 3.2), the Sinkhorn optimization has only limited flexibility in

what it can do other than map them through a fairly simple function to crowd density. This is especially true as there is only a small set of parameters being trained with Sinkhorn. It is highly likely and straightforward to map the frequent crowd features to its density values, whose distribution is signaled through the prior. Moreover, we show through extensive experiments in Section 4 and 5 that CSS-CCNN actually ends up learning crowd density estimation.

3.4 Improving Sinkhorn Matching

As described already, the Sinkhorn training updates the network parameters by backpropagating Sinkhorn loss \mathcal{L}_{sink} , which brings the distribution of the density predictions closer to that of the prior. But computing \mathcal{L}_{sink} relies on estimating the optimal transport plan \mathbf{P}^* (the solution to optimization in equation 3) through the Sinkhorn iterations (see [12] for more details). The quality of estimation of \mathbf{P}^* directly affects the performance of the model. Hence, it is quite beneficial to aid the computation of \mathbf{P}^* by providing additional information. Any signal that can potentially boost the transport assignments is helpful. For example, even simply grouping the prediction measures H^{CS} to a coarse sparse-dense categories and then restricting assignments within the groups from that of the prior, leads to improved performance. This is because the restricted assignments make sure that the dense samples from prediction are always mapped to dense points in the prior (similarly for sparse ones), reducing costly errors of connecting dense ones to sparse and vice versa. However, one needs to have the density category information to supplement the Sinkhorn assignments and that should be obtained in an unsupervised fashion as well.

We observe that the edge details of crowd images could serve as an indicator of density. Highly crowded regions seem to have more density of edges, while it is low for relatively sparse or non-crowds. But this is a weak signal and can have lots of false positives. The higher density of edges could arise from non-crowds such as background clutter or other patterns with more edges. Interestingly for dense crowds, we find that this weak supervisory signal is good enough for grouping regions into potential dense or very sparse categories. For any given crowd image, the standard Canny edge detector [73] is applied to extract the edge map. The map is then blurred and down-sampled to look like density maps. These pseudo maps resemble the actual ground truth crowd density in many cases, having a relatively higher response in dense region than at sparse ones. Note that the absolute values from the pseudo maps do not follow the actual crowd density and hence cannot be directly used for supervision. However, given a set of crowd patches, the relative density values are sufficient to faithfully categorize regions into two broad density groups. This is done by first sorting pseudo counts of the patches and then dividing the samples at a predetermined percentile. Crowd regions with pseudo count values above this threshold are considered to be dense while those below goes to the highly sparse or non-crowd. By employing a percentile threshold, accurate count values are not required and the pseudo counts should only need to be relatively correct across the given set of images. Since any random set of crowd patches

should follow the prior distribution (as per the assumptions and approximations in Section 3.1), the percentile threshold is fixed on the prior. We fix the threshold to be 30th percentile as there are roughly 30% samples that are non-crowds or with very low counts in the range of 1.

We modify the Sinkhorn training to incorporate the pseudo density information in the following manner: first, we compute the pseudo counts H^{CSP} corresponding to the prediction samples H^{CS} . Using pseudo counts H^{CSP} , H^{CS} is split to the sparse part H_0^{CS} and the dense H_1^{CS} . The prior samples are also grouped with the same threshold to get H_0^{GT} and H_1^{GT} . Now the Sinkhorn loss is separately found for both the categories and added. The exact loss being backpropagated is,

$$\mathcal{L}_{sink}^{++}(H^{GT}, H^{CS}) = \mathcal{L}_{sink}(H_0^{GT}, H_0^{CS}) + \mathcal{L}_{sink}(H_1^{GT}, H_1^{CS}) \quad (4)$$

By separating out the assignment of sparse and dense samples, the counting performance of the model increases as evident from the experiments in Section 4. Note that the Sinkhorn training is complete on its own without the auxiliary density information. It is a simple addendum to the method that can improve the performance.

4 EXPERIMENTS

4.1 Evaluation Scheme and Baselines

Any crowd density regressor is evaluated mainly for the standard counting metrics. There are two metrics widely being followed by the community. The first is the MAE or Mean Absolute Error, which directly measures the counting performance. It is the absolute difference of the predicted and actual counts averaged over the test set or simply expressed as $MAE = (1/S_{test}) \sum_{i=1}^N |C_i - C_i^{GT}|$, where C_i is the count predicted by the model for i th image and C_i^{GT} denotes the actual count. Note that S_{test} is the number of images in the test set. Coming to the second metric, the Mean Squared Error or MSE is defined as $MSE = \text{SQRT}((1/S_{test}) \sum_{i=1}^N (C_i - C_i^{GT})^2)$, a measure of the variance of count estimation and it represents the robustness of the model.

Our completely self-supervised framework is unique in many ways that the baseline comparisons should be different from the typical supervised methods. It is not fair to compare CSS-CCNN with other approaches as they use the full annotated data for training. Hence, we take a set of solid baselines for our model to demonstrate its performance. The *CCNN Random* experiment refers to the results one would get if only *Stage 1* self-supervision is done without the subsequent Sinkhorn training. This is the random accuracy for our setting and helpful in showing whether the proposed complete self-supervision works. Since our approach takes one parameter, the maximum count value of the dataset (C^{fmax}) as input, *CCNN Mean* baseline indicates the counting performance if the regressor blindly predicts the given value for all the images. We choose mean value as it makes for sense in this setting than the maximum (which anyway has worse performance than mean). Another important validation for our proposed paradigm is the *CCNN P_{prior}* experiment, where the model gives out a value randomly drawn from the prior distribution as its prediction for a

given image. The counting performance of this baseline tells us with certainty whether the *Stage 2* training does anything more than that by chance. Apart from these, the *CCNN Fully Supervised* trains the entire regressor with the ground truth annotations from scratch. Note that we do not initialize CCNN with any pretrained weights as is typically done for supervised counting models. *CCNN Self-Supervised with Labels*, on the other hand, runs the *Stage 1* training to learn the FEN parameters and is followed by labeled optimization for updating the regressor layers. These are not directly comparable to our approach as we do not use any annotated data for training, but are shown for completeness.

We evaluate our model on different datasets in the following sections. The results for the naive version of the Sinkhorn loss \mathcal{L}_{sink} (Section 3.3) is labeled as *CSS-CCNN*, whereas *CSS-CCNN++* represents the one with the improved \mathcal{L}_{sink}^{++} (Section 3.4). Note that only the train/validation set images are used for optimizing *CSS-CCNN* and the ground truth annotations are never used. The counting metrics are computed with the labeled data from the test after the full training. Unless otherwise stated, we use the same hyper-parameters as specified in Section 3.

4.2 Shanghaitech Dataset

The Shanghaitech Part_A [13] is a popular dense crowd counting dataset, containing 482 images randomly crawled from the Internet. It has images with crowd counts as low as 33 to as high as 3139, with an average of 501. The train set has 300 images, out of which 10% is held out for validation. There are 182 images testing. The hyper-parameter used for this is $C^{fmax} = 3000$. We compare the performance of *CSS-CCNN* with the baselines listed earlier and other competing methods in Table 1. The metrics for our method is evaluated for three independent runs with different initialization and the mean along with variance is reported. It is clear that *CSS-CCNN* outperforms all the baselines by a significant margin. This shows that the proposed method works better than any naive strategies that do not consider the input images. With the improved loss, *CSS-CCNN++* achieves around 5% less counting error than the naive version due to the more faithful Sinkhorn matching process. Moreover, the *CCNN* network with rotation self-supervision also beats the model developed in [54]. It is worthwhile to note that the performance of *CSS-CCNN* is reminiscent of the results of early fully supervised methods with the MAE being better than a few of them as well (see Table 2 of [13]). Figure 4

TABLE 1
Performance comparison of *CSS-CCNN* with other methods on Shanghaitech PartA [13]. Our model outperforms all the baselines.

Method	MAE	MSE
CCNN Fully Supervised	118.9	196.6
Sam et al. [54]	154.7	229.4
CCNN Self-Supervised with Labels	121.2	197.5
C-CNN Random	431.1	559.0
C-CNN Mean	282.8	359.9
C-CNN P_{prior}	272.2	372.5
CSS-CCNN (ours)	207.3 ± 5.9	310.1 ± 7.7
CSS-CCNN++ (ours)	195.6 ± 5.8	293.2 ± 9.3

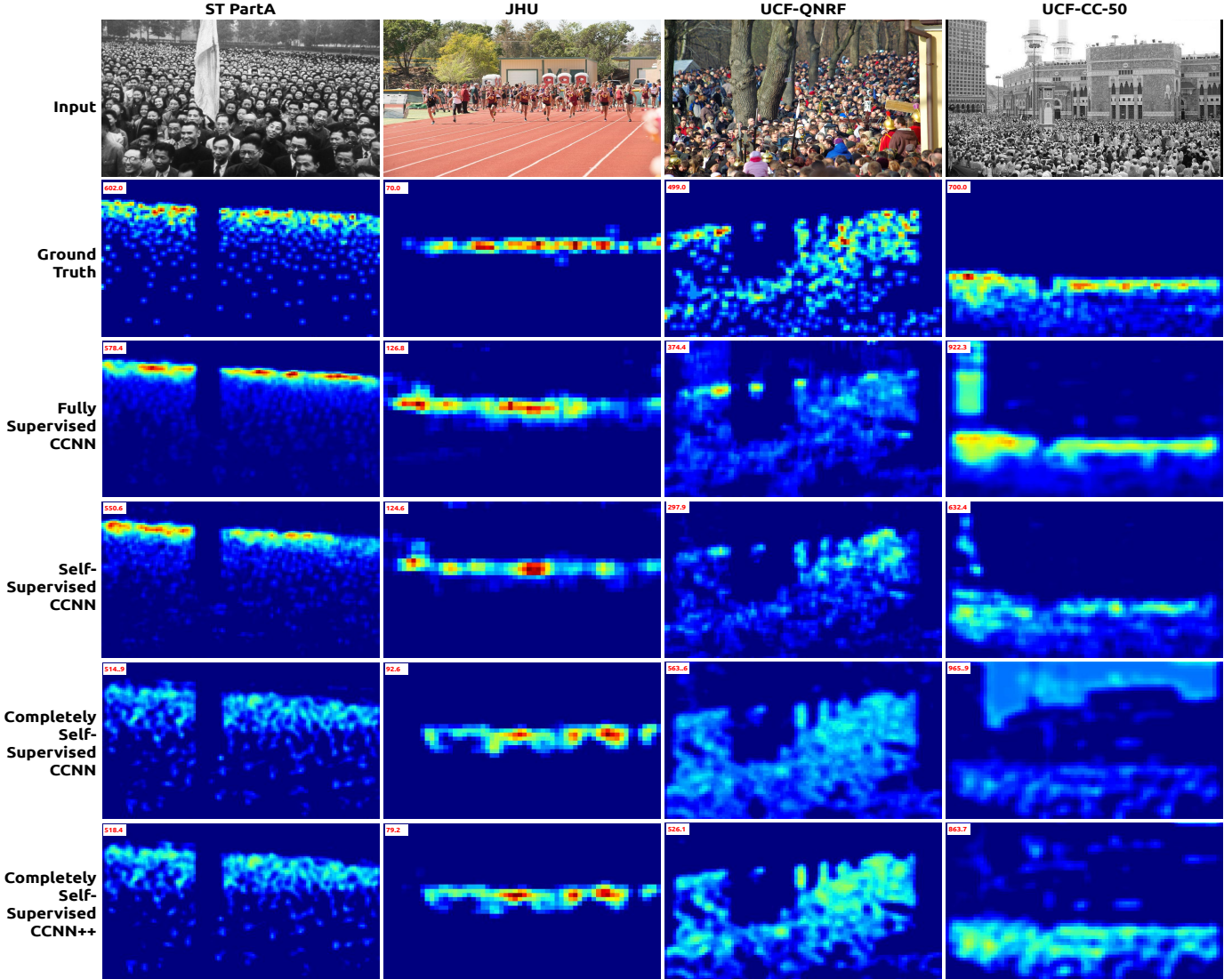


Fig. 4. Density maps estimated by CSS-CCNN along with that of baseline methods. Despite being trained without a single annotated image, CSS-CCNN is seen to be quite good at discriminating the crowd regions as well as regressing the density values.

visually compares density predictions made by CSS-CCNN and other models. The predictions of our approach are mostly on crowd regions and closely follows the ground truth, emphasizing its ability to discriminate crowds well.

4.3 UCF-QNRF Dataset

UCF-QNRF dataset [35] is a large and diverse collection of crowd images with 1.2 million annotations. There are 1535 images with crowd count varying from 49 to 12865, resulting in an average of 815 individuals per image. The dataset offers very high-resolution images with an average resolution of 2013×2902 . The count hyper-parameter provided to the model is $C^{fmax} = 12000$. We achieve similar performance trends on UCF-QNRF dataset as well. CSS-CCNN outperforms all the unsupervised baselines in terms of MAE and MSE as evident from Table 2. Since the dataset has extreme diversity in terms of crowd density, it is important to improve the Sinkhorn matching process and faithfully assign appropriate counts across density categories. Owing to the better distribution matching, CSS-CCNN++ achieves

around 9% less counting error than CSS-CCNN, despite the dataset being quite challenging.

4.4 UCF-CC-50 Dataset

UCF-CC-50 dataset [14] has just 50 images with extreme variation in crowd density ranging from 94 to 4543. The small size and diversity together makes this dataset the most challenging. We follow the standard 5-fold cross-validation scheme suggested by the creators of the dataset to report the performance metrics. Since the number of images is quite small, the assumption taken for setting the prior distribution gets invalid to certain extent. But a slightly different parameter to the prior distribution works. We set $\alpha = 1$ and $C^{fmax} = 4000$. Despite being a small and highly diverse dataset, CSS-CCNN is able to beat all the baselines. The self-supervised MAE is also better than the method in [54]. These results evidence the effectiveness of our method. CSS-CCNN++ improves upon the result significantly in terms of MSE, indicating improved performance on highly dense crowds.

TABLE 2

Benchmarking CSS-CCNN on UCF-QNRF dataset [35]. Our approach beats the baseline methods in counting performance.

Method	MAE	MSE
CCNN Fully Supervised	159.0	248.0
CCNN Self-Supervised with Labels	196.8	309.3
CCNN Random	718.7	1036.3
CCNN Mean	567.1	752.8
CCNN P_{prior}	535.6	765.9
CSS-CCNN (ours)	442.4 \pm 4.2	721.6 \pm 13.9
CSS-CCNN++ (ours)	414.0 \pm 16.3	652.1 \pm 15.6

TABLE 3

Performance CSS-CCNN on UCF-CC-50 [14]. Despite being very challenging dataset, CSS-CCNN achieves better MAE than baselines.

Method	MAE	MSE
CCNN Fully Supervised	320.6	455.1
Sam et al. [54]	433.7	583.3
CCNN Self-Supervised with Labels	348.8	484.3
CCNN Random	1279.3	1567.9
CCNN Mean	771.2	898.4
CCNN P_{prior}	760.0	949.9
CSS-CCNN (ours)	564.9	959.4
CSS-CCNN++ (ours)	557.0	737.9

TABLE 4

Evaluation of CSS-CCNN on JHU-CROWD++ [36] dataset.

Method	MAE	MSE
CCNN Fully Supervised	128.8	415.9
CCNN Self-Supervised with Labels	147.5	436.2
CCNN Random	320.3	793.5
CCNN Mean	316.3	732.3
CCNN P_{prior}	302.3	707.621
CSS-CCNN (ours)	243.6 \pm 9.1	672.4 \pm 17.1
CSS-CCNN++ (ours)	197.9 \pm 2.2	611.9 \pm 12.0

4.5 JHU-CROWD++ Dataset

JHU-CROWD++ ([36], [74]) is a comprehensive dataset with 1.51 million head annotations spanning 4372 images. The crowd scenes are obtained under various scenarios and weather conditions, making it one of the challenging dataset in terms of diversity. Furthermore, JHU-CROWD++ has a richer set of annotations at head level as well as image level. The maximum count is fixed to $C^{fmax} = 8000$. The performance trends are quite similar to other datasets, with our approach delivering better MAE than the baselines as evident from Table 4. This indicates the generalization ability of CSS-CCNN across different types of crowd datasets.

4.6 Cross Data Performance and Generalization

In this section, we evaluate our proposed model in a cross dataset setting. CSS-CCNN is trained in a completely self-supervised manner on one of the dataset, but tested on other datasets. Table 5 reports the MAEs for the cross dataset evaluation. It is evident that the features learned from one dataset are generic enough to achieve reasonable scores on the other datasets, increasing the practical utility of CSS-CCNN. The difference in performance mainly stems from the changes in the distribution of crowd density across the

TABLE 5

Cross dataset performance of our model; the reported entries are the MAEs obtained for CSS-CCNN and CSS-CCNN++ respectively.

Train \downarrow / Test \rightarrow	ST_PartA	UCF-QNRF	JHU-CROWD++
ST_PartA	207.3, 195.6	468.1, 472.4	254.0, 251.3
UCF-QNRF	251.2, 235.7	442.4, 414.0	236.5, 220.6
JHU-CROWD++	290.2, 266.3	446.2, 417.4	243.6, 197.9

TABLE 6

Evaluating CSS-CCNN in a true practical setting: the model is trained on images crawled from the web, but evaluated on crowd datasets. The counting performance appears similar to that of training on the dataset.

Train on web images	CSS-CCNN		CSS-CCNN++	
	MAE	MSE	MAE	MSE
Test on ST_PartA	208.8	309.5	184.2	268.8
Test on UCF-QNRF	450.7	755.9	422.1	699.9
Test on JHU-CROWD++	241.2	706.8	231.0	660.1

datasets. This domain shift is drastic in the case of UCF-CC-50 [14], especially since the dataset has only a few images.

4.7 CSS-CCNN in True Practical Setting

The complete self-supervised setting is motivated for scenarios where no labeled images are available for training. But till now we have been using images from crowd datasets with the annotations being intentionally ignored. Now consider crawling lots of crowd images from the Internet and employing these unlabeled data for training CSS-CCNN. For this, we use textual tags related to dense crowds and similarity matching with dataset images to collect approximately 5000 dense crowd images. No manual pruning of undesirable images with motion blur, perspective distortion or other artifacts is done. CSS-CCNN is trained on these images with the same hyper-parameters as that of ShanghaiTech Part_A and the performance metrics are computed on the datasets with annotations. From Table 6, it is evident that our model is able to achieve almost similar or better MAE on the standard crowd datasets, despite not using images from those datasets for training. This further demonstrates the generalization ability of CSS-CCNN to learn from less curated data, emphasizing the practical utility it could facilitate.

4.8 Performance with Limited Data

Here we explore the proposed algorithm along with fully supervised and self-supervised approaches when few annotated images are available for training. The analysis is performed by varying the number of labeled samples and the resultant counting metrics are presented in Figure 5. For training CSS-CCNN with data, we utilise the available annotated data to compute the optimal Sinkhorn assignments P^* and then optimize the \mathcal{L}_{sink} loss. This way both the labeled as well as unlabeled data can be leveraged for training by alternating respective batches (in a 5:1 ratio). It is clear that, at very low data, scenarios CSS-CCNN beats the supervised as well as self-supervised baselines by a significant margin. The Sinkhorn training shows 13% boost in MAE (for ShanghaiTech Part_A) by using just one labeled

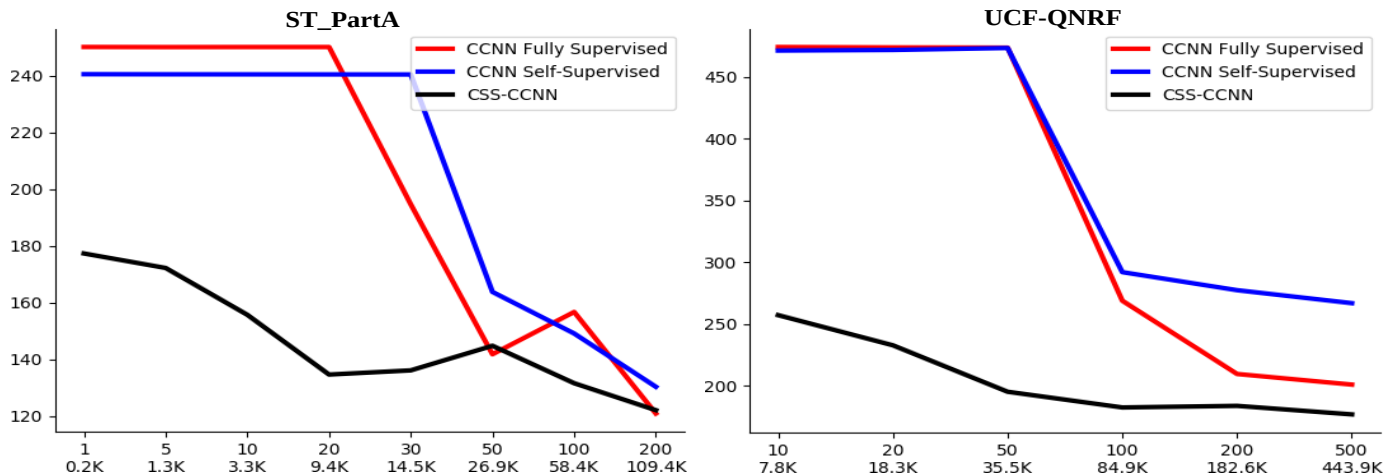


Fig. 5. Comparing our completely self-supervised method to fully supervised and self-supervised approaches under a limited amount of labeled training data. The x-axis denotes the number of training images along with the count (in thousands) of head annotations available for training, while the y-axis represents the MAE thus obtained. At low data scenarios, CSS-CCNN has significantly superior performance than others.

sample as opposed to no samples. This indicates that CSS-CCNN can perform well in extremely low data regimes. It takes about 20K head annotations for the supervised model to perform as well as CSS-CCNN. Also, CSS-CCNN has significantly less number of parameters to learn using the labeled samples as compared to a fully supervised network. These results suggests that our complete self-supervision is the right paradigm to employ for crowd counting when the amount of available annotated data is less.

5 ABLATIONS AND ANALYSIS

5.1 Ablations on Architectural Choices

In Table 7, we validate our architectural choices taken in designing CSS-CCNN. The first set of experiments ablates the *Stage 1* self-supervised training. We perform Sinkhorn training on a randomly initialized FEN (labeled as *Without Stage 1*) and receive a worse MAE. In the chosen setting of self-supervision with rotation, the input image is randomly rotated by one of the four predefined angles for creating pseudo labels. Now we analyse the effect of the number of rotation classes on the final counting metrics. As evident from the table, four angles stands to be the best in agreement with previous research on the same [66]. Self-supervision via colorization is another popular strategy for learning useful representations. The model is trained to predict the a-b color space values for a given gray-scale image (the L channel). The end performance is observed to be inferior in comparison with that of the rotation task. Another option is to load FEN with ImageNet trained weights (as this is a typical way of transfer learning) and then employ *Stage 2*. The result (*With ImageNet weights*) is worse than that of CSS-CCNN, suggesting that the self-supervised training is crucial to learn crowd features necessary for density estimation. Furthermore, the base *feature extraction network* (FEN) (see Figure 3 in main paper) is changed to ResNet blocks and CSS-CCNN is trained as well as evaluated (*with ResNet based FEN*). Simple VGG style architecture appears to be better for density regression. We also run experiments with different types of prior distributions and see that the power law with

exponential cutoff works better, justifying our design choice. The *without skip connection* experiment trains CSS-CCNN devoid of the features from the second convolutional block in FEN being directly fed to C_2 (see Figure 3 and Section 3.3). As expected, the feature aggregation from multiple layers improves the counting performance. The cell sizes used for computing count histograms (see Section 3.1) are varied (labeled *Cell Size*) to understand the effect on MAE. The metrics seem to be better with our default setting of 8×8 . CSS-CCNN employs a prior parametric distribution to facilitate the unlabeled training. We investigate the case where the prior is directly given in the form of an empirical measure derived from the ground truth annotations. For the Sinkhorn training, this *GT distribution* is sampled to get H^{GT} (see Section 3.3) instead of P_{prior} . The resultant MAE is very similar to the standard CSS-CCNN setting, indicating that our chosen prior approximates the ground truth distribution well. Lastly, we ablate the *percentile threshold* used to extract

TABLE 7

Validating different architectural design choices made for CSS-CCNN evaluated on the ShanghaiTech Part_A [13] (computed on single run).

Method	MAE	MSE
Without <i>Stage 1</i>	257.5	397.7
Rotation with 2 class	233.5	344.1
Rotation with 8 class	232.2	341.5
Colorization	242.5	363.0
With ImageNet weights	257.8	370.8
ResNet based FEN	244.8	332.4
Uniform Prior	261.8	406.0
Pareto Prior	248.3	386.2
Lognormal Prior	239.5	345.8
Without skip connection	226.8	329.1
Cell Size 2×2	243.9	374.6
Cell Size 4×4	251.6	389.4
With GT distribution	202.7	300.3
Percentile threshold 10	191.5	288.7
Percentile threshold 50	189.1	286.8
CSS-CCNN	197.3	295.9
CSS-CCNN++	187.7	280.21

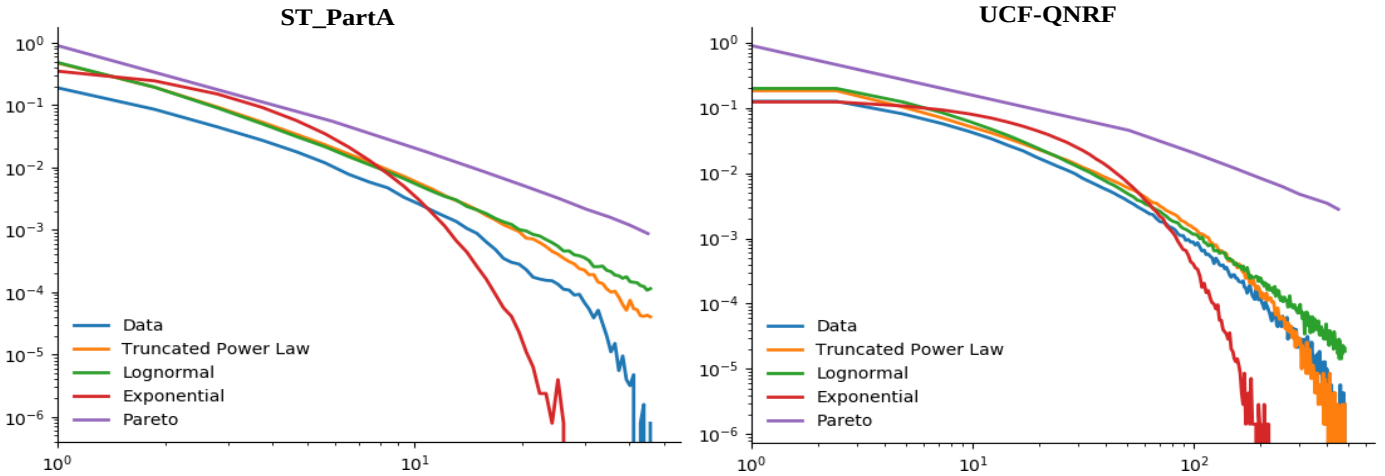


Fig. 6. Double logarithmic representation of maximum likelihood fit for the crowd counts from Shanghaitech Part_A [13] and UCF-QNRF [35].

of pseudo density category for the CSS-CCNN++ model (Section 3.4) and find that the default setting helps in better density differentiation.

5.2 Analysis of the Prior Distribution

The proposed Sinkhorn training requires a prior distribution of crowd counts to be defined and the choice of an appropriate prior is essential for the best model performance as seen from Table 7. Here we analyze the crowd data more carefully to see why the truncated power law is the right choice of prior. For this, the counts from crowd images are extracted as described in Section 3.1 and a maximum likelihood fit over various parametric distributions is performed. The double logarithmic visualization of the probability distribution of both the data and the priors are available in Figure 6. Note that the data curve is almost a straight line in the logarithmic plot, a clear marker for power law characteristic. Both truncated power law and lognormal tightly follow the distribution. But on close inspection of the tail regions, we find truncated power law to best represent the prior. This further validates our choice of the prior distribution.

5.3 Sensitivity Analysis for the Crowd Parameter

As described in Section 3.1, CSS-CCNN requires the maximum crowd count (C^{fmax}) for the given set of images as an input. This is necessary to fix the prior distribution parameter λ . One might not have the exact max value for the crowds in a true practical setting; an approximate estimate is a more reasonable assumption. Hence, we vary C^{fmax} around the actual value and train CSS-CCNN on Shanghaitech PartA [13] and UCF-QNRF [35]. The performance metrics in Table 8 show that changing C^{fmax} to certain extent does not alter the performance significantly. The MAE remained roughly within the same range, even though the max parameter is being changed in the order of 500. Note that the results are computed with single runs. These findings indicate that the our approach is insensitive to the exact crowd hyper-parameter value, increasing its practical utility. We also check the sensitivity of our approach on the power law exponent α . Varying α around 2 results in similar

TABLE 8
Sensitivity analysis for the hyper-parameters on CSS-CCNN. Our model is robust to fairly large change in the max count parameter.

ST_PartA			UCF-QNRF		
Param	MAE	MSE	Param	MAE	MSE
$C^{fmax} = 2000$	204.2	316.4	$C^{fmax} = 10000$	443.9	749.7
$C^{fmax} = 2500$	197.9	304.6	$C^{fmax} = 11000$	446.9	757.5
$C^{fmax} = 3000$	197.3	295.9	$C^{fmax} = 12000$	437.0	722.3
$C^{fmax} = 3500$	191.9	288.5	$C^{fmax} = 14000$	446.1	697.5
$\alpha = 1.9$	202.9	303.3	$\alpha = 1.9$	438.3	700.6
$\alpha = 2.0$	197.3	295.9	$\alpha = 2.0$	437.0	722.3
$\alpha = 2.1$	200.7	305.6	$\alpha = 2.1$	446.4	756.3

performances, in agreement with the findings of existing works and our design choice (see Section 3.1).

5.4 Analysis of Features

To further understand the exact learning process of CSS-CCNN, the acquired features can be compared against that of a supervised model. Figure 7 displays the mean feature map for the outputs at various convolutional blocks of CSS-CCNN along with that of the supervised baseline (see Section 3) evaluated on a given crowd image. Note that Conv4 stands for the regressor block that is trained with the Sinkhorn loss in the case of CSS-CCNN. It is clear that the self-supervised features closely follow the supervised representations, especially at the initial blocks in extracting low-level crowd details. Towards the end blocks, features are seen to diverge, with fully supervised Conv4 outputs appearing like density maps. But notice that the corresponding completely self-supervised outputs have higher activations on heads of people, which is relevant for the end task of density estimation. This clearly shows that CSS-CCNN indeed learns to extract crowd features and detect heads, rather than falling in a degenerate case of matching the density distribution without actually counting persons.

6 CONCLUSIONS AND FUTURE WORK

We show for the first time that a density regressor can be fully trained from scratch without using a single annotated

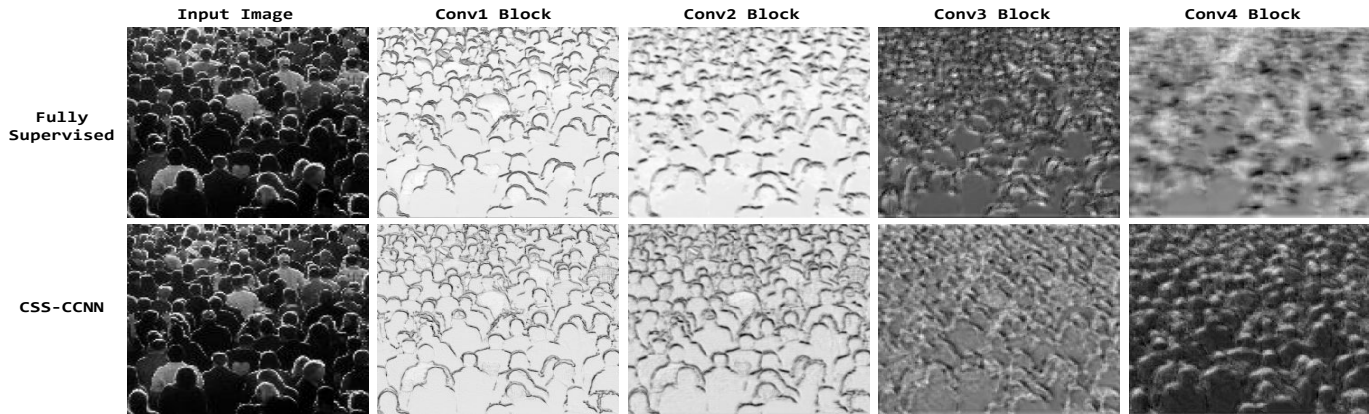


Fig. 7. Visualization of mean features extracted from different convolutional blocks of CSS-CCNN and the supervised baseline.

image. This new paradigm of complete self-supervision relies on optimizing the model by matching the statistics of the distribution of predictions to that of a predefined prior. Though the counting performance of the model stands better than other baselines, there is a performance gap compared to fully supervised methods. Addressing this issue could be the prime focus of future works. For now, our work can be considered as a proof of concept that models could be trained directly for solving the downstream task of interest, without providing any instance-level annotated data.

REFERENCES

- [1] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, 2006.
- [2] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- [3] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013.
- [4] A. Makhzani and B. J. Frey, "Winner-take-all autoencoders," in *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [5] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [6] G. Larsson, M. Maire, and G. Shakhnarovich, "Learning representations for automatic colorization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [7] —, "Colorization as a proxy task for visual understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [8] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," 2018.
- [9] Z. Feng, C. Xu, and D. Tao, "Self-supervised representation learning by rotation feature decoupling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [10] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [11] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [12] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [13] 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 (CVPR)*, 2016.
- [14] H. Idrees, I. Saleemi, C. Seibert, and M. Shah, "Multi-source multi-scale counting in extremely dense crowd images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [15] C. Zhang, H. Li, X. Wang, and X. Yang, "Cross-scene crowd counting via deep convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [16] D. Onoro-Rubio and R. J. López-Sastre, "Towards perspective-free object counting with deep learning," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [17] D. Babu Sam, S. Surya, and R. V. Babu, "Switching convolutional neural network for crowd counting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [18] D. Babu Sam, N. N. Sajjan, R. V. Babu, and M. Srinivasan, "Divide and grow: Capturing huge diversity in crowd images with incrementally growing CNN," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [19] X. Cao, Z. Wang, Y. Zhao, and F. Su, "Scale aggregation network for accurate and efficient crowd counting," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [20] 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 (CVPR)*, 2018.
- [21] X. Jiang, Z. Xiao, B. Zhang, X. Zhen, X. Cao, D. Doermann, and L. Shao, "Crowd counting and density estimation by trellis encoder-decoder networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [22] C. Liu, X. Weng, and Y. Mu, "Recurrent attentive zooming for joint crowd counting and precise localization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [23] V. A. Sindagi and V. M. Patel, "CNN-based cascaded multi-task learning of high-level prior and density estimation for crowd counting," in *Proceedings of the IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 2017.
- [24] —, "Generating high-quality crowd density maps using contextual pyramid CNNs," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [25] D. Babu Sam and R. V. Babu, "Top-down feedback for crowd counting convolutional neural network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [26] W. Liu, M. Salzmann, and P. Fua, "Context-aware crowd counting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [27] Z.-Q. Cheng, J.-X. Li, Q. Dai, X. Wu, and A. G. Hauptmann, "Learning spatial awareness to improve crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [28] Z. Yan, Y. Yuan, W. Zuo, X. Tan, Y. Wang, S. Wen, and E. Ding, "Perspective-guided convolution networks for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [29] M. Shi, Z. Yang, C. Xu, and Q. Chen, "Revisiting perspective information for efficient crowd counting," in *Proceedings of the*

- IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [30] A. Zhang, J. Shen, Z. Xiao, F. Zhu, X. Zhen, X. Cao, and L. Shao, "Relational attention network for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [31] A. Zhang, L. Yue, J. Shen, F. Zhu, X. Zhen, X. Cao, and L. Shao, "Attentional neural fields for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [32] N. Liu, Y. Long, C. Zou, Q. Niu, L. Pan, and H. Wu, "ADCrowd-Net: An attention-injective deformable convolutional network for crowd understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [33] J. Wan, W. Luo, B. Wu, A. B. Chan, and W. Liu, "Residual regression with semantic prior for crowd counting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [34] V. Ranjan, H. Le, and M. Hoai, "Iterative crowd counting," in *Proceedings of the European Conference on Computer Vision*, 2018.
- [35] H. Idrees, M. Tayyab, K. Athrey, D. Zhang, S. Al-Maadeed, N. Rajpoot, and M. Shah, "Composition loss for counting, density map estimation and localization in dense crowds," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [36] V. A. Sindagi, R. Yasarla, and V. M. Patel, "Pushing the frontiers of unconstrained crowd counting: New dataset and benchmark method," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [37] L. Liu, Z. Qiu, G. Li, S. Liu, W. Ouyang, and L. Lin, "Crowd counting with deep structured scale integration network," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [38] V. A. Sindagi and V. M. Patel, "Multi-level bottom-top and top-bottom feature fusion for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [39] J. Wan and A. Chan, "Adaptive density map generation for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [40] C. Xu, K. Qiu, J. Fu, S. Bai, Y. Xu, and X. Bai, "Learn to scale: Generating multipolar normalized density maps for crowd counting," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [41] J. Liu, C. Gao, D. Meng, and A. G. Hauptmann, "DecideNet: Counting varying density crowds through attention guided detection and density estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [42] D. Lian, J. Li, J. Zheng, W. Luo, and S. Gao, "Density map regression guided detection network for rgb-d crowd counting and localization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [43] Y. Liu, M. Shi, Q. Zhao, and X. Wang, "Point in, box out: Beyond counting persons in crowds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [44] D. Babu Sam, S. V. Peri, M. N. Sundararaman, and R. V. Babu, "Going beyond the regression paradigm with accurate dot prediction for dense crowds," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020.
- [45] D. Babu Sam, S. V. Peri, M. N. Sundararaman, A. Kamath, and R. V. Babu, "Locate, size and count: Accurately resolving people in dense crowds via detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.
- [46] M. Zhao, J. Zhang, C. Zhang, and W. Zhang, "Leveraging heterogeneous auxiliary tasks to assist crowd counting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [47] Z. Shi, P. Mettes, and C. G. Snoek, "Counting with focus for free," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [48] Z. Ma, X. Wei, X. Hong, and Y. Gong, "Bayesian loss for crowd count estimation with point supervision," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [49] Z. Shi, L. Zhang, Y. Liu, X. Cao, Y. Ye, M.-M. Cheng, and G. Zheng, "Crowd counting with deep negative correlation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [50] Z. Shen, Y. Xu, B. Ni, M. Wang, J. Hu, and X. Yang, "Crowd counting via adversarial cross-scale consistency pursuit," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [51] 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 International Conference on Computer Vision (ICCV)*, 2019.
- [52] X. Liu, J. Van De Weijer, and A. D. Bagdanov, "Exploiting unlabeled data in CNNs by self-supervised learning to rank," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [53] Q. Wang, J. Gao, W. Lin, and Y. Yuan, "Learning from synthetic data for crowd counting in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [54] D. Babu Sam, N. N. Sajjan, H. Maurya, and R. V. Babu, "Almost unsupervised learning for dense crowd counting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- [55] P. Agrawal, J. Carreira, and J. Malik, "Learning to see by moving," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [56] D. Jayaraman and K. Grauman, "Learning image representations tied to ego-motion," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [57] D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan, "Learning features by watching objects move," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [58] X. Wang and A. Gupta, "Unsupervised learning of visual representations using videos," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [59] I. Misra, C. L. Zitnick, and M. Hebert, "Shuffle and learn: unsupervised learning using temporal order verification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [60] P. Isola, D. Zoran, D. Krishnan, and E. H. Adelson, "Learning visual groups from co-occurrences in space and time," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [61] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [62] T. Nathan Mundhenk, D. Ho, and B. Y. Chen, "Improvements to context based self-supervised learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [63] R. Zhang, P. Isola, and A. Efros, "Split-brain autoencoders: Unsupervised learning by cross-channel prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [64] S. Jenni and P. Favaro, "Self-supervised feature learning by learning to spot artifacts," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [65] L. Zhang, G.-J. Qi, L. Wang, and J. Luo, "AET vs. AED: Unsupervised representation learning by auto-encoding transformations rather than data," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [66] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [67] A. Clauset, C. R. Shalizi, and M. E. Newman, "Power-law distributions in empirical data," *SIAM review*, 2009.
- [68] D. Helbing, A. Johansson, and H. Z. Al-Abideen, "Dynamics of crowd disasters: An empirical study," *Physical review E*, 2007.
- [69] M. Moussaïd, D. Helbing, and G. Theraulaz, "How simple rules determine pedestrian behavior and crowd disasters," *Proceedings of the National Academy of Sciences*, 2011.
- [70] I. Karamouzas, B. Skinner, and S. J. Guy, "A universal power law governing pedestrian interactions," *APS*, 2015.
- [71] D. Helbing, C. Kühnert, S. Lämmer, A. Johansson, B. Gehlsen, H. Ammoser, and G. B. West, "Power laws in urban supply networks, social systems, and dense pedestrian crowds," in *Complexity Perspectives in Innovation and Social Change*. Springer, 2009, pp. 433–450.
- [72] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [73] J. Canny, "A computational approach to edge detection," *IEEE Transactions on pattern analysis and machine intelligence (TPAMI)*, 1986.
- [74] V. A. Sindagi, R. Yasarla, and V. M. Patel, "JHU-CROWD++: Large-scale crowd counting dataset and a benchmark method," *Technical Report*, 2020.