

Multi-modal line-crossing using a multi-headed model for flow-enhanced density maps and a pixel-wise accumulator

Jan Erik van Woerden

30 October 2020

Contents

1	Introduction	4
1.1	Introduction	4
1.2	Contributions	5
1.3	Thesis outline	5
2	Background	6
2.1	Region of Interest	6
2.1.1	Density Map	6
2.2	Flow Estimation	7
2.3	Line of Interest	8
3	Related Work	10
3.1	Crowd Counting	10
3.1.1	CSRNet	10
3.2	Flow Estimation	11
3.2.1	PWCNet	11
3.2.2	DDFlow	11
3.3	Line of Interest	12
3.3.1	Instant LOI counting	12
3.3.2	Region-level LOI counting	14
4	Method	15
4.1	Velocity Map and Density Map	15
4.1.1	Models	16
4.1.2	Realigning	17
4.2	LOI counting	17
5	Implementation	19
5.1	Models	19
5.1.1	Baseline 1	19
5.1.2	Baseline 2	19
5.1.3	Unified model	20
5.1.4	Flow context density map	20
5.2	Environment	20
5.2.1	Maxing filter	20
5.2.2	Line Crossing	20
5.2.3	Optimizer	21
5.2.4	Augmentation	21
5.2.5	Metrics	21

6 Datasets	22
6.1 Requirements	22
6.2 Labelling	22
6.3 Datasets	23
6.3.1 UCSD	23
6.3.2 Fudan-ShanghaiTech	23
6.3.3 AI City Challenge	23
7 Results	25
7.1 Fudan-ShanghaiTech	25
7.2 UCSD	25
7.3 AI City Challenge	26
7.4 Real world performance	26
7.5 Flow estimation impact	26
8 Conclusion	27
8.1 Further research	27
A Appendix	32

Chapter 1

Introduction

1.1 Introduction

In recent years the amount of surveillance camera's has increased immensely to a point that it is hard to supervise them all manually by humans. Since the increased use of surveillance cameras a lot of research is done to automate information extraction from those cameras [Sreenu and Saleem Durai, 2019].

A widely researched area for extracting information from camera's is Crowd Counting [Chan and Vasconcelos, 2008, Wang et al., 2020, Li et al., 2018, Fang et al., 2019, Liu et al., 2019b]. In Crowd Counting the amount of pedestrians present in the frame is counted. Where pedestrians in low density areas can be easily counted by general object recognition, higher density area's need specialized methods to accurately count the amount of pedestrians [Zhang et al., 2016].

While Crowd Counting only focusses on counting the exact amount of pedestrians in a frame, it doesn't take into account the amount of pedestrians walked by over time. This is no issue when camera's are present in the whole area of interest, however with large areas this becomes a much bigger issue. In for example a shop it would be much more convenient to count the amount of customers inside the shop by only tracking the customers walking inside and outside the shop, instead of having cameras in the whole shop.

This research area of Crowd Crossing Counting is much less researched [Ma and Chan, 2013, Zhao et al., 2016, Zheng et al., 2019]. By adding Flow Estimation to the Crowd Density (Crowd Counting) the flow of pedestrians can be obtained and the amount of people entering and exiting an area can be measured.

In early papers on Crowd Crossing Counting prediction was done using key-point extraction and feeding the keypoints to a regression model [Ma and Chan, , Ma and Chan, 2013]. More recently the introduction of Convolutional Neural Networks was made into the field of both Crowd Counting [Zhang et al., 2016, Liu et al., 2019b, Li et al., 2018, Wang et al., 2020] and Flow Estimation [Sun et al., , Dosovitskiy et al., 2015]. Which sparked the research in those fields.

In Crowd Crossing Counting the amount of research done using Neural Networks is limited [Zhao et al., 2016, Cao et al.,]. New research in both Crowd Counting and Flow estimation provides lots of new opportunities to improve the State-of-the-Art of Line Crossing. New research also adds some new challenges, which we try to solve as well in this thesis.

1.2 Contributions

The latest Neural Networks for Crowd Crossing Counting [Zhao et al., 2016] learns both Crowd Counting and Flow Estimation in a supervised way. However the labeling process of both these methods are a time consuming method, especially labeling the images for Flow Estimation. To reduce the labour intensive task of labeling, a new method is introduced to use and train unsupervised Flow Estimation during Line Crossing Counting.

Additionally a multi-headed model is proposed which trains the Crowd Counting in a supervised matter and the Flow Estimation in an unsupervised matter. The model utilizes a shared encoder, while two distinct decoders are used to predict each a individual task.

To further utilize the availability of both the predictions a flow enhanced model is proposed. This model uses the Flow Estimation to further enhance the Crowd Counting prediction.

The main dataset used in Crowd Crossing Counting is UCSD. This dataset is low resolution [Ma and Chan, 2013, Ma and Chan,] (720x480) and doesn't provide much diversity in scenery. In Crowd Counting higher resolution datasets are already available. In this thesis Line Crossing labeling is provided and published for the Fudan-Shanghai dataset (1920x1080, 25FPS) [Fang et al., 2019] using self-developed labeling software.

Lastly thorough research is done on the usability of the presented system in real world scenario's. Impact on video frame rate are tested and the processing time is compared. Additionally to show the generality of the method a cross-domain dataset from the AI City challenge is used, which counts passing cars on a crossroad.

1.3 Thesis outline

The rest of this thesis is divided into the next chapters:

- **Background**, explains several fields to understand the starting point for this thesis.
- **Related work**, which explains more in depth related work which is used in this thesis.
- **Method**, presents the method of the proposed solution.
- **Implementation**, presents the hyperparameters and evaluation methods.
- **Datasets**, presents the used datasets and used approach to label the proposed new datasets.
- **Results**, discusses the results of the experiments.
- **Conclusion**, wraps it up and summarizes what we can conclude.

Chapter 2

Background

In this chapter a selection of terms is explained which gives a basis to understand the rest of this thesis. This background is created with the assumption that the reader has a basic background in Machine Learning and (Convolutional) Neural Networks.

Reframe ROI
and LOI in the
introduction or
in the rest of
the paper

2.1 Region of Interest

The Region of Interest problem is a widely studied problem in which the goal is to estimate the amount of pedestrians given a single image. Directly predicting the counted pedestrians given a Neural Network is a hard task, because of the lack of supervision, this would require a large amount of samples to accurately solve this task. All recent State-of-the-Art methods therefore use an intermediate representation to give the model enough supervision to perform Crowd Counting with a low amount of training samples.

In the early days of Crowd Counting several methods have been proposed which use *detection-based* methods to estimate the amount of pedestrians [Dalal and Triggs, 2005, Dollár et al., 2012]. Several papers were introduced which tried to detect only the head [Subbaraman et al., 2012] and others tried to focus on general part detection [Wu and Nevatia, 2007, Lin and Davis, 2010]. These methods rely on individually detecting the pedestrians. This becomes much harder when occlusion of the pedestrians start to happen. This is why the performance of these methods start to degrade when the density of the pedestrians in an image start to increase.

Later papers introduced a *regression-based* solution, which tries to predict the amount of pedestrians in crowd blobs [Chan and Vasconcelos, 2009, Idrees et al., 2013, Zheng et al., 2019]. Using SVM or other regressor methods and several features such as the amount of foregrounds pixels of the blob and detected key points the count inside crowd blobs were predicted. Regression based solutions were an improvement over the detection methods, but still lack the capabilities to estimate pedestrians counts in highly occluded areas.

2.1.1 Density Map

With the introduction of Convolutional Neural Networks in the field of Crowd Counting density maps were proposed as well to count pedestrians [Zhang et al., 2016, Liu et al., 2019b, Li et al., 2018]. A *density map* (Figure 2.1) used for Region of

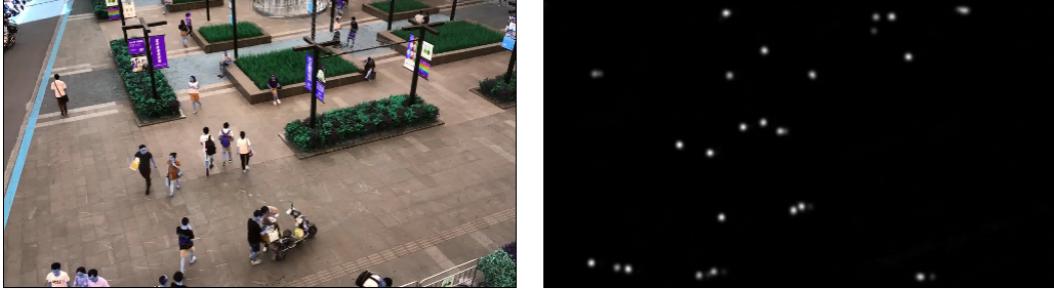


Figure 2.1: Example of generated density map on the right side, for the left image

Interest is a map which represents the density of pedestrians of each pixel. The density map is generated by taking the locations of each pedestrian ($p = \begin{bmatrix} x_p \\ y_p \end{bmatrix}$ in equation 2.1) and place those locations on the the density map.

Individual dots are very hard for a Neural Network to detect correctly and are prone to errors. To circumvent this a Gaussian shaped circle is created around this location, still with with a sum of 1. The amount of pedestrians in the frame can be extracted from the density map by taking the sum over all the pixels of the density map (Equation 2.2, where $D_t(p)$ is the density for location $p = \begin{bmatrix} x_p \\ y_p \end{bmatrix}$ for trainings frame t).

$$D_t(p) = \frac{1}{2\pi\sigma_p^2} \sum_{p \in P} e^{\frac{(x_p - x)^2 + (y_p - y)^2}{-2\sigma_p^2}} \quad (2.1)$$

$$C_t = \sum_{p \in P} D_t(p) \quad (2.2)$$

Several methods have been presented to optimize the generation of density maps [Zhang et al., 2016, Li et al., 2018, Wan and Chan, 2019]. For most medium dense frames the difference in methods is minimal. Often in benchmarks with medium dense frames a fixed sigma is used ($\sigma_p = \sigma_i$ in equation 2.1). For highly dense frames the use of different methods can have a difference, especially when the difference in size between close pedestrians and pedestrians in the background is large [Li et al., 2018]. In our own preliminary results, this has no effect on our datasets.

2.2 Flow Estimation

The research which is done on the Flow Estimation problem is widely used. Approaches on this topic can be used in a wide range of applications which makes it very interesting. Already in the early 1980's Horn and Schunck [Horn and Schunck, 1981] published the first paper which tried to predict flow. Since then lots of different approaches have been published [Mémin and Pérez, 1998, Bruhn et al., 2005, Brox et al., 2014]. Long conventional mathematical approaches have ruled the flow estimation field. Later also learnable models were introduced [Pock et al., 2008, Wedel et al., 2009].

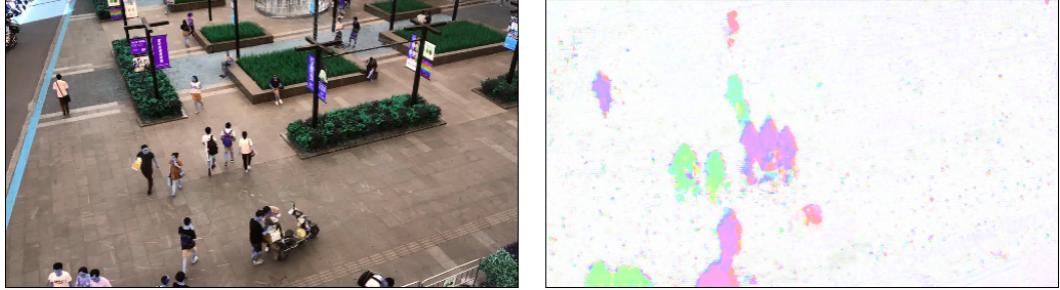


Figure 2.2: Example of generated velocity map on the right side, for the left image

Recent papers however make use of Convolutional Neural Network based models [Dosovitskiy et al., 2015, Ilg et al., 2016, Sun et al., , Ranjan and Black, 2017, Hui et al., 2018]. These models predict pixel-precise velocity maps. The *velocity map* (Figure 2.2) is a map which predict per pixel of the frame the amount of movement to another location. In equation 2.3, $V_t(p)$ shows the velocity map as a difference between the location of the pixel in the current frame (p) and the location of this pixel in the next frame ($N_t(p)$).

$$V_t(p) = N_t(p) - \begin{bmatrix} x_p \\ y_p \end{bmatrix} \quad (2.3)$$

Creating a real world dataset that utilizes the power of pixel-wise flow estimation is very hard [Dosovitskiy et al., 2015]. There are no real world devices which could capture both video and create pixel perfect ground-truths to train the flow estimation models on. Most of the flow estimation benchmarks are therefore generated videos. Computer 3D-engines make it possible to generate pixel-perfect flow estimation based on the generated videos in the engine.

However there is a large gap in domain and scene between the generated datasets and real world applications [Liu et al., 2008]. Recent supervised papers [Dosovitskiy et al., 2015, Sun et al.,] tend to overfit on these datasets and therefore perform rather poor on real world applications. One promising direction is unsupervised learning [Yu et al., 2016, Janai et al., 2018, Liu et al., 2019a, Liu et al.,]. Early papers only predicted non-occluded pixels [Yu et al., 2016, Janai et al., 2018], but recent papers use methods to estimate occluded pixels as well [Liu et al., 2019a, Liu et al.,]. Further details about these methods in related work.

2.3 Line of Interest

Line of Interest is very similar to Region of Interest. Where Region of Interest is the interest of the amount of people inside the ROI, the Line of Interest is the focus on the amount of pedestrians that cross the specified line during a certain timeframe. This LOI is defined as a single line between two points p_1 and p_2 (Top and bottom point in figure 2.3)

With the Line of Interest problem the goal is to give the amount of pedestrians crossing of each side given a set of frames (a pre-captured video or video stream).

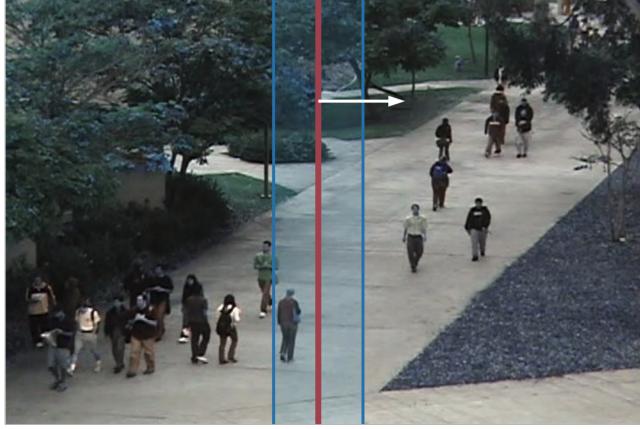


Figure 2.3: Example of LOI, red: LOI, blue: LOI area

The output of the prediction should give two numbers c_1 and c_2 which are the amount of pedestrians crossing from each side.

Only a handful of papers are published about Line of Interest. In the earlier papers [Ma and Chan, , Cao et al.,], slicing was a widely used approach to estimate the Line of Interest. With slicing a small area, called the LOI area (Blue area in figure 2.3), is taken around the LOI. Over a set of consecutive frames each slice of the frame was taken and stitched together into a single image. On the images slow walking pedestrians appear rather wide and fast walking pedestrians shallow. By counting the amount of pedestrians present on the stitched image, the total amount of pedestrians crossing the line can be counted.

The area is defined by all the pixels that have a maximum distance to the LOI of d and can be projected on the LOI. When projected, the pixels fall between p_1 and p_2 .

Recent papers discard this method [Zhao et al., 2016, Zheng et al., 2019], because it makes it hard to track pedestrian with different speeds and walking in different directions give artifacts which make it hard to track those pedestrians [Zhao et al., 2016]. The slicing method is replaced with an actual frame by frame prediction method. Using two consecutive frames the amount of pedestrians crossing the line is measured. These newer methods predict both location and direction of the pedestrian.

Based on these new papers, the problem of Line of Interest is divided into three separate problems. Locating the pedestrians (Region of Interest), estimate the direction (Flow Estimation) of the pedestrians and combining these two streams of information into the count for Line of Interest. Further details of this approach will be provided in related work.

Chapter 3

Related Work

In this chapter we try to explain a couple of key papers on which the proposed methods are build.

3.1 Crowd Counting

3.1.1 CSRNet

Since Zhang et al. [Zhang et al., 2016] lots of research was done on better predicting the density maps for Crowd Counting. Zhang et al. and following proposed several multi-column models to predict pedestrians of different sizes. However Li et al. [Li et al., 2018] changed this with the proposal of CSRNet. Instead of using multi-column models it uses dilated kernels. This method beat the multi-column models on several benchmarks.

The advantage of dilated kernels is the increase of the reception field without increasing the amount of parameters and layers in the network. As show in figure 3.1 the dilation rate is the distance between each filter input of the kernel. By increasing the dilation rate the size of area the kernel covers without increasing the amount of parameters of a small-sized kernel.

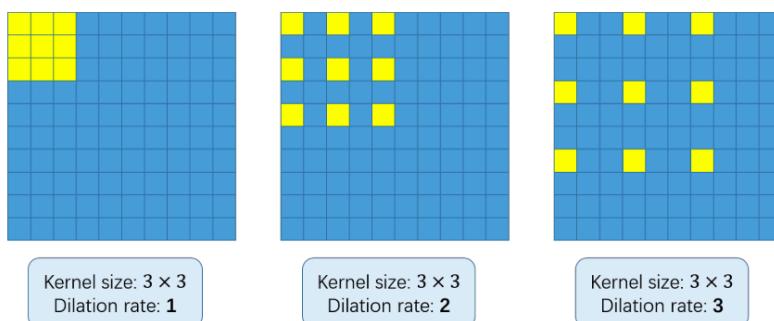


Figure 3.1: Dilation rates on a 3×3 kernel

3.2 Flow Estimation

3.2.1 PWCNet

A popular Flow Estimation network is PWCNet [Sun et al.,]. It uses the original ideas of FlowNet, but it improves FlowNet in a lot of ways. FlowNet traditionally is used to fully predict the flow with a neural network. PWCNet massively reduces the number of weights (See figure 3.2), which results in faster training and much quicker prediction. Additionally the network shows a higher accuracy on several benchmarks.

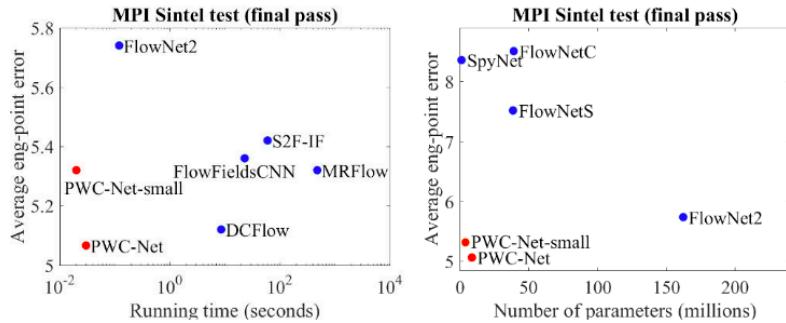


Figure 3.2: PWCNet compared to other existing flow estimation architectures

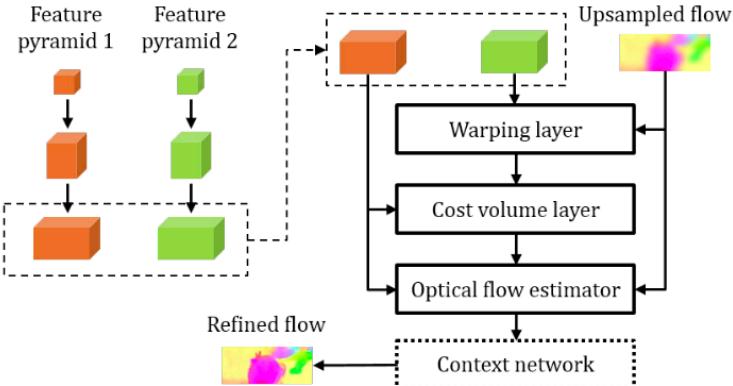


Figure 3.3: PWCNet architecture

PWCNet uses a pyramid shape architecture to predict the velocity map (See figure 3.3). The encoder encodes both input frames separately into two feature volumes. During decoding several decoding steps are taken. After the initial velocity map estimation, each decoding step upscales the so-far estimated map and further refines the map. At the start of each decoding step the estimated velocity map is used to warp the feature volume of the second frame and correlate this warped volume with the volume of the first image to focus on area of refinement in the velocity map.

3.2.2 DDFlow

In DDFlow [Liu et al., 2019a] a general method is proposed to estimate a velocity map in an unsupervised manner. Earlier methods already proposed several methods

to optimize using a photometric loss [Yu et al., 2016] and ignore occluded-pixels during loss calculation [Janai et al., 2018]. However all these earlier methods used handcrafted methods to estimate the occluded-pixels.

The paper [Liu et al., 2019a] proposes a method which learns occluded-pixels using a distillation from unlabeled data without supervision of humans to label the ground truth. Because of the generality of the method, the method can be wrapped around all existing supervised flow estimation architectures.

The method uses a teacher and a student approach to train the occluded pixels. The teacher network is trained on all the non-occluded pixels, with exclusion of the occluded pixels using only the photometric loss with occlusion-awareness. After training the teacher network, the velocity maps of the teacher network are used as ground-truth for the student network.

The student network is then trained on patches of the original predicted velocity map. On the borders of the patches, occluded pixels will appear, because moving pixels from inside the patched frame will move to outside the patched frame. These border pixels will be marked as occluded pixel by the student network, but will be non-occluded pixels according to the ground-truth of the teacher network.

$$L_p = \sum \psi(I_1 - I_2^w) \odot (1 - O_f) / \sum (1 - O_f) + \sum \psi(I_2 - I_1^w) \odot (1 - O_b) / \sum (1 - O_b) \quad (3.1)$$

The photometric loss [Yu et al., 2016] with occlusion-awareness [Janai et al., 2018] proposed in the earlier papers is defined in equation 3.1. Where I_1 and I_2 are the full input frames and I_i^w the backwarped image based on the predicted velocity map. Additionally O_b and O_f are masks for the occluded pixels. These maps are calculated by checking if the mismatch between forward flow and backward flow is too large.

For the student model the photometric loss is used together with the loss for occlusion (equation 3.2). The loss of the student model then just $L_p + L_o$. In equation 3.2, \tilde{w}_f and \tilde{w}_b are the predicted velocity map forward and backward using the student model. w_b^p and w_f^p are the cropped patches from the ground-truth velocity map predicted by the teacher model. M_f and M_b are just defined as the difference of occluded pixels between the ground-truth of the parent and the teacher ($M_f = \text{clip}(\tilde{O}_f - O_f^p, 0, 1)$). These pixels are not occluded in the full frame and occluded in the patch. Which means that the groundtruth using distillation is available.

$$L_o = \sum \psi(w_f^p - \tilde{w}_f) \odot M_f / \sum M_f + \sum \psi(w_b^p - \tilde{w}_b) \odot M_b / \sum M_b \quad (3.2)$$

3.3 Line of Interest

3.3.1 Instant LOI counting

Zhao et al. [Zhao et al., 2016] introduces a new approach to calculate the amount of pedestrians crossing the Line of Interest. Instead of slicing a range of consecutive frames around the Line of Interest and stitching them together, they present a

method that directly predicts the amount crossing pedestrians using two consecutive frames.

They use both a density map and velocity map and finally merge them together to obtain the LOI counts. Both the maps are trained fully supervised. Because annotating raw velocity maps is a very hard task, they simplify the velocity map in disk-shaped area's around the pedestrian locations. Annotating each frame is still very demanding, but lowers the amount of labeling significantly in comparing to full velocity maps.

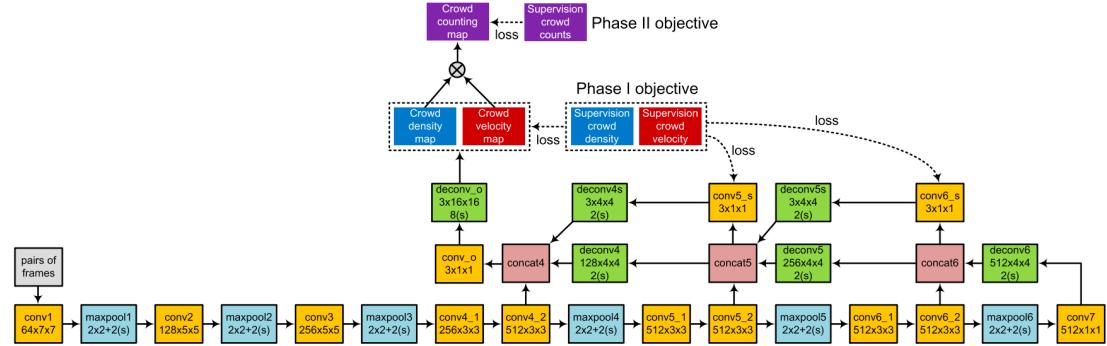


Figure 3.4: Pipeline of Zhao et al. with FlownetSimple as architecture, where in the last layer both the density map as the velocity map are predicted

They are as well the first who predict both the density map and the velocity map in a single Convolutional Neural Network. A network they use a simplified model of FlowNet. Because the velocity map and density map are so similar in their location, they predict both the velocity map and the density map in the final layer (See figure 3.4).

Training contains of two stages. In the first stage, the model is optimized for both the density map and velocity map. During the second stage the density map and velocity map are multiplied to train a directional counting map as $C_t = D_t \otimes V_t$. Which it is trained on to further align the velocity map and the density map.

To finally merge both the density map and velocity map together they use the created directional counting map. According to the paper the directional counting map gives the amount of pedestrians crossing that pixel between the time of the two frames.

$$c_{1,t} = \sum_{\{p|\cos(\theta_p) \geq 0\}} \sqrt{C_{t,x}(p)^2 + C_{t,y}(p)^2} \cdot \cos(\theta_p) \quad (3.3)$$

$$c_{2,t} = \sum_{\{p|\cos(\theta_p) < 0\}} \sqrt{C_{t,x}(p)^2 + C_{t,y}(p)^2} \cdot (-\cos(\theta_p)),$$

$$c_1 = \sum_{\{t|t \in T\}} c_{1,t}, \quad c_2 = \sum_{\{t|t \in T\}} c_{2,t} \quad (3.4)$$

To calculate per frame pair the amount of pedestrians crossed the line, a set of locations p around the LOI is taken. Then the directional counting map is normalized and summed together in equation 3.3, where θ_p is the angle between $V_t(p)$ and the LOI.

To calculate the total line crosses over a certain timeframe, the results of equation 3.3 can be summed as in equation 3.4.

3.3.2 Region-level LOI counting

Zheng et al. [Zheng et al., 2019] provides a fast method of predicting the Line of Interest. The paper outperforms [Zhao et al., 2016] in the UCSD benchmark. Zheng et al. discusses the problems with Neural Networks and the complexity the models, which makes it hard to run the models real time. It therefore introduces a non-CNN based method based on a SVM, linear regression and the Lucas-Kanade optical flow tracker.

It uses the idea of [Zhao et al., 2016] to discard the slicing method and uses pairwise prediction and uses the same method to merge density and velocity. Because the methods used in [Zheng et al., 2019] are not on a pixel-level, a method is proposed to bin on a region-level. In equation 3.5 a single velocity $v_{t,r}$ is calculated with a weighted average over all moving keypoints inside the region.

$$v_{t,r} = \frac{\sum_p w_r(p) \cdot v_{\text{towards}}^{(t)}(p)}{\sum_p w_r(p)} \quad (3.5)$$

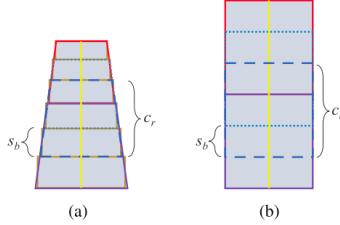


Figure 3.5: Regions around the LOI with skewness and normalized to a straight LOI

Additional to the region-level LOI counting, they introduce the method to skew the region to take into account the perspective of the camera view (See figure 3.5). This helps the SVM and linear regressor to more accurately predict the amount of pedestrians inside the region.

Preliminary results show that the binning on a region-level is not improving the Neural Network models to increase accuracy.

Chapter 4

Method

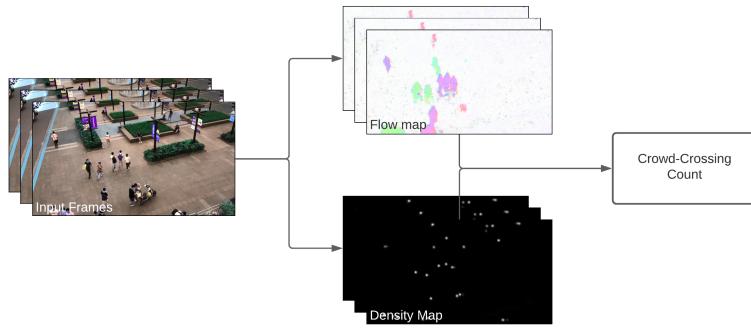


Figure 4.1: Overview of Crowd-Crossing Counting pipeline

In this chapter the proposed method for Crowd Crossing Counting is explained. As shown in figure 4.1, the full pipeline of prediction is based in two stages. First we predict both the density map and the velocity map, then we use these maps to predict the Crowd-Crossing Count. In section *Velocity Map and Density Map* are the methods explained to train/predict both the velocity map and density map. In section *LOI counting* the method is explained to predict the Crowd-Crossing Count based on those two intermediate maps.

4.1 Velocity Map and Density Map

Whereas earlier papers rely on velocity maps and density maps [Zhao et al., 2016], the amount of labeling required to train the models is high and creates an extra barrier for real world applications.

To increase the usability of the method, a supervised velocity map estimation is used. Training to predict both maps is being done in a parallel manner using a unified loss function described in equation 4.1. For the velocity loss, L_v , the photometric loss from figure 3.1 [Yu et al., 2016, Janai et al., 2018] is used, whereas for the density loss, L_c , the traditional L2 loss is used.

$$L_t = L_v + \lambda \cdot L_c \quad (4.1)$$

Predicting the Velocity Map and Density Map are both two widely researched fields, as explained in the background chapter. Two separate models would be

capable of predicting both the Velocity Map and Density Map accurately. However this gives a lot of extra overhead in processing time, due to separately fully processing the image. Additionally there is not an easy solution to enhance the other map with context to perform better.

Therefore two new unified models are proposed. Both are multi-headed networks which both predict a velocity map and a density map. It takes into account the multi-modal problem of predicting both a velocity map and density map. The second model enhances the density-head using the final velocity map as extra context information.

4.1.1 Models

Unified model

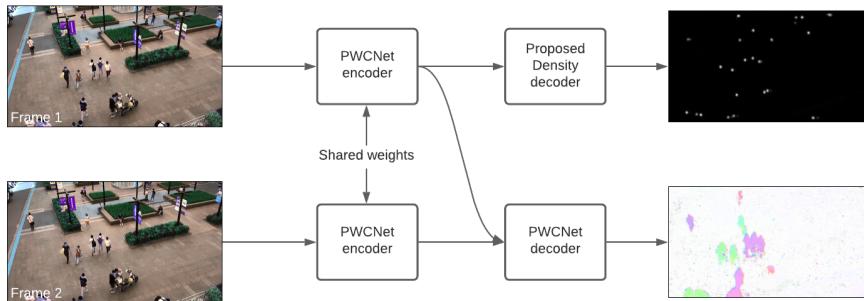


Figure 4.2: Unified model

The first proposed model is a multi-headed model with a shared encoder (Fig 4.3).

The model uses the original PWCNet network [Sun et al.,] as backbone. The proposed model shares the complete encoder and decoder of the PWCNet, but adds a second decoder to predict a density map as well, as shown in figure 4.3.

This decoder uses a decoder structure with feeding features in several stages of the decoding stage. Additionally the proposed dilation kernels in CSRNet[Li et al., 2018] are used for a larger reception field which boosts the performance of the density map prediction.

Flow context density map

In the unified model, both heads only share the encoder, but don't use the final outcomes of either of the decoders to enhance the other decoder. It would be beneficial for Line Counting when the model can be optimized for moving pedestrians. These are the ones counted during Line Crossing (Equation 4.3).

Therefore we propose a novel method to enhance the density map predictor which makes use of both of these available information streams. This model starts with the unified model (Figure 4.3) as base. By adding the output flow map as context features to the existing encoder features, the decoder could use this information to better detect moving pedestrians.

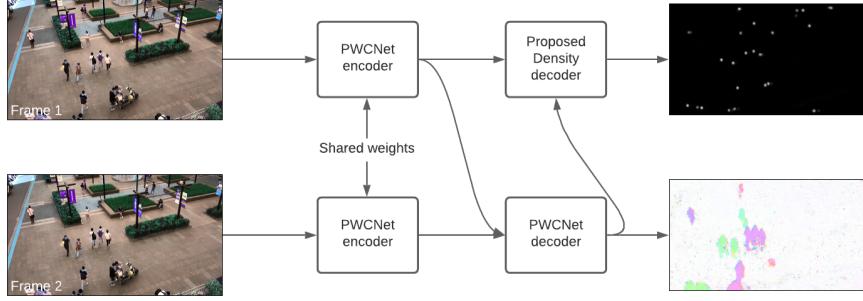


Figure 4.3: Flow enhanced model

4.1.2 Realigning

In [Zhao et al., 2016] a supervised method of aligning the velocity map and density map is used. However due to the unsupervised prediction of the velocity map in this thesis it is not possible to use this alignment. As shown in figure 4.4 this misalignment is an huge issue when multiplying the velocity map and density map in a pixel-wise matter. The blobs predicted in the density map are often much bigger than the flow of the pedestrian predicted. Additionally, for several Crowd datasets the tagging of the pedestrian is done on the head, which is why the blobs are often predicted over the head of the pedestrian.

To tackle this misalignment problem we propose an expanding method by applying a maxing filter on the flow estimation. This maxing filter takes the local maximum value in a surrounding of each pixel. Looking at figure 4.4 this helps to cover a lot of misaligned density maps and flow maps. To optimize for heads on the top side of the pedestrian, the maxing filter is focused on the bottom side of the selected pixel. Maximum values above the pixel are ignored.

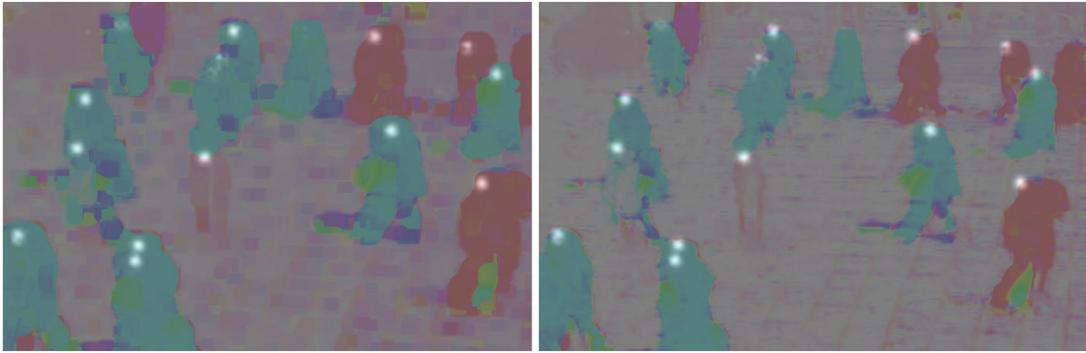


Figure 4.4: On the right a non-maxed flow estimation and on the right the flow estimation with maxing filter applied

4.2 LOI counting

For our approach the main idea of [Zhao et al., 2016] is taken (Explained in Related Work), however the current paper uses some simplifications, by reframing the pixel-

V2, better labeling the dots in the figure, separating the figure description

level counting in the following way. The approach is much more theoretical correct.

We define v_{perp} (White arrow in figure 2.3) as the normalized directional vector perpendicular to the LOI (Two solutions are perpendicular on the LOI and this defines sides 1 and 2 of the LOI counting). Then we define the collection of the pixels on the left side of the LOI and inside the LOI area as M_1 (side 1) and the pixels on the right side (side 1 and inside the LOI area) as M_2 .

The velocity towards the LOI is then defined as the dot-product of V_t and v_{perp} (Equation 4.2).

$$Q_t(p) = V_t(p) \cdot v_{perp} \quad (4.2)$$

$$\begin{aligned} c_{1,t} &= \sum_{\{p \in M_1 | Q_t(p) > 0\}} C_t(p) \cdot \frac{Q_t(p)}{d} \\ c_{2,t} &= \sum_{\{p \in M_2 | Q_t(p) < 0\}} C_t(p) \cdot \frac{-Q_t(p)}{d} \end{aligned} \quad (4.3)$$

Then the LOI count on timestep t is defined in equation 4.3. Where $\frac{Q_t(p)}{d}$ defines the percentage that the density on the specific pixel has crossed the LOI area. Lastly we can sum the count over a timespan into a single count for each side as in equation 3.4.

Add more visual description in same image as figure 2.3, explain that this assumption of left/right is without losing generality

Chapter 5

Implementation

In this chapter details about the actual implementation are explained in more depth.

5.1 Models

In the results four different models are compared: Two baselines and two proposed models.

5.1.1 Baseline 1

The first model is proposed in [Zhang et al., 2016]. Because this model assumes a supervised Flow Estimation, parts of the proposed method can't be used. Only the proposed model therefore is used as baseline, where the exact network is shown in figure 3.4.

Due to full shared nature of the network the assumption is that the network will perform very poor when predicting both the density map and the unsupervised flow map. Therefore we propose a second baseline as well.

5.1.2 Baseline 2

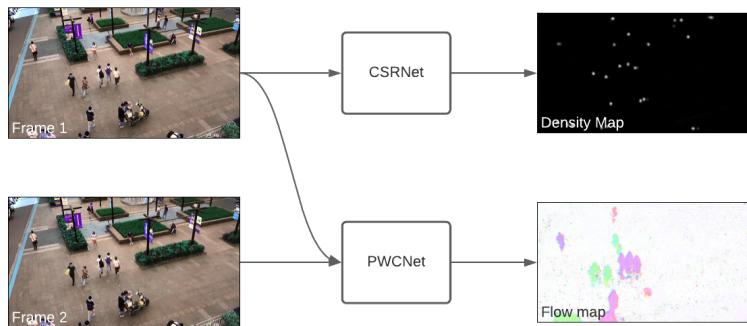


Figure 5.1: CSRNet+PWCNet baseline

This stronger baseline is a combination of two independent models (Figure 5.1). CSRNet [Li et al., 2018] for predicting the density map and PWCNet [Sun et al.,] for flow estimation. Both models perform very well in comparing to papers in their

respective field. Therefore will this baseline be a very powerful baseline to compare to. The models are trained independently of each other with no shared loss function to further optimize their performance.

In V2, a pyramid decoder is used, which intermediately sends lower resolution density maps to optimize weight optimization and reduce overfitting

5.1.3 Unified model

The first proposed model shares the decoder of the PWCNet model (Full network in appendix, figure A.1). The PWCNet makes use of a pyramide shaped architecture and the encoder provides the decoder with 6 feature maps ranging from 1/2 the size to 1/64 the size of the original image.

The decoder contains two essential processing blocks. The dilation block, which is a block with 4 conv-layers with a dilation of 2. Additionally an upscaling block is used, which first upscales the input by two and then refines by 2 conv-layers. All conv-layers use a kernel of 3x3 and a stride of 1.

Each of the four tiniest feature map are processed by a dilation block. The smallest feature map is processed first and individually upscaled using the upscaling block. The second and third feature maps are first concatenated with the earlier feature maps and then upscaled.

After processing the fourth feature map and concatenation the upscaled feature maps, the features are processed by two dilation blocks which then predicts the density map using a single output layer.

In V2 the structure is summarised in a visual graph as well, now only in appendix

5.1.4 Flow context density map

The second model proposed enhances the feature maps with the output of the flow map. The flow enhanced model uses the first proposed model as base decoder. Instead of only decoding the feature maps on each level, the output flow map is reshaped and concatenated to each feature map. Which results the information of the flow on each level available.

5.2 Environment

5.2.1 Maxing filter

During experiments the maxing filter is optimized per dataset. For the Fudan-ShanghaiTech dataset a maxing filter is applied with a distance of 12px, optimized in implementation. For the UCSD dataset a maxing filter of 4px is applied.

5.2.2 Line Crossing

In the equation 4.3 all the parameters for the Line Crossed are shown. The width of the line (parameter d) is the last parameter which is not defined. In preliminary results the difference is width is minimal, during all the experiments a width of 20px is used therefore.

5.2.3 Optimizer

For all the experiments the Adam optimizer [Kingma and Ba, 2015] is used with a learning-rate of $2 \cdot 10^{-5}$. A regularization of 10^{-4} is applied.

During training, the unified loss (Equation 4.1) is used for training using a λ of 5. Which at the start focusses mainly on the density map loss and after some initial training produces an equal loss distribution between the velocity map loss and density map loss.

5.2.4 Augmentation

To augment the dataset several augmentations will be applied on the training samples which are used in [Li et al., 2018] as well. First a crop of the image is made. Ranging from $1/3$ and $1/6$ of the total image size. Then the cropped image is resized to a size of $1/4$ of the original image (So $1/2$ the width and $1/2$ the height). Lastly the cropped image is half of the time flipped horizontally.

5.2.5 Metrics

For both the ROI and the LOI the Mean Average Error and the Mean Squared Error are used. Additionally the LOI uses the Relative Mean Average Error.

$$MAE = \frac{1}{n} \sum_{i=1}^n |C_i - P_c^{(i)}| \quad (5.1)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (C_i - P_c^{(i)})^2 \quad (5.2)$$

For the ROI the MAE and the MSE are defined as equation 5.1 and 5.2. Where $P_c^{(i)}$ is the predicted density map for the given frame.

$$MAE = \frac{1}{n} \sum_{i=1}^n |G_l^{(i)} - P_l^{(i)}| + |G_r^{(i)} - P_r^{(i)}| \quad (5.3)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (G_l^{(i)} - P_l^{(i)})^2 + (G_r^{(i)} - P_r^{(i)})^2 \quad (5.4)$$

$$RMAE = \frac{1}{n} \sum_{i=1}^n \frac{|G_l^{(i)} - P_l^{(i)}| + |G_r^{(i)} - P_r^{(i)}|}{G_l^{(i)} + G_r^{(i)}} \quad (5.5)$$

TODO, when applying no cropping, less overfitting occurs at the corners, but much slower

Rewrite, because UCSD uses a different method, which is kind of similar, but slightly different

For the LOI the MAE and MSE are defined as quation 5.3 and 5.4. The RMAE is simply defined as in equation 5.5. Where $G_l^{(i)}$ is the ground truth for sample i for side left-to-right. And $P_r^{(i)}$ is the predicted value for right-to-left.

Chapter 6

Datasets

In this chapter we explain the datasets in more depth. First we explain the requirements of the dataset to correctly train and evaluate each dataset. To compensate for lack of some required labeling a tool is written and explained. Lastly all the datasets are explained in more depth.

6.1 Requirements

For training and evaluation we need two different methods of labeling. For the density map generation the position of each pedestrian is required for the frames which are used for training. For evaluation the line crossing it is required to label the amount of pedestrians crossing the LOI from each side. Ideally the training set is purely labeled with head-tags and the evaluation set only with line crossing labeling.

6.2 Labelling

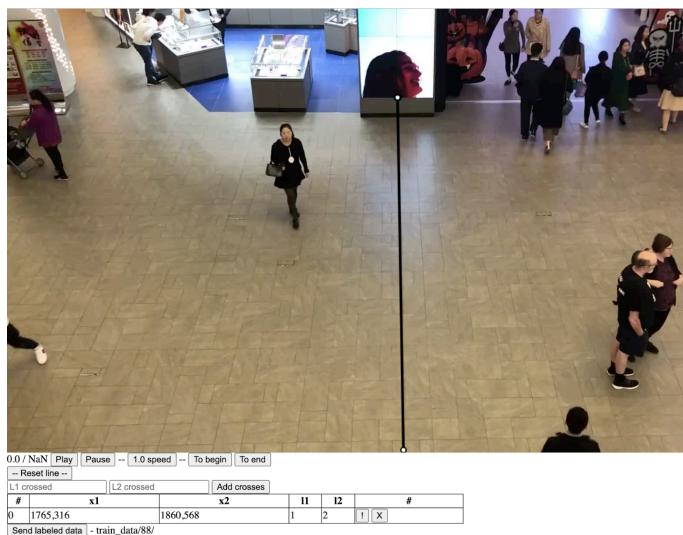


Figure 6.1: User interface of the labeler

Several Crowd Counting datasets provide sequences of frames with corresponding pedestrian labelling. However most of those datasets don't provide line crossing

labeling. Therefore I build a tool to label videos for line crossing (Figure 6.1). The labeler loads a video from the unlabeled videos. In this video the user can label multiple lines by first clicking on the video to draw a line and afterwards fill in the amount of pedestrians crossing the line during the video. Additionally the user can scroll and view through the video using multiple manipulations.

6.3 Datasets



Figure 6.2: Samples of each dataset

6.3.1 UCSD

The UCSD Pedestrian dataset [Chan and Vasconcelos, 2008] is a public dataset created in 2008. The dataset contains a video of a path where people cross each other (See figure 6.2a). The dataset is 720x480 in resolution. The original videos are in 30fps, but are downsampled to 10fps. A total of 5680 frames are captured. The first 2000 frames contain line of interest labeling [Ma and Chan, 2013]. Frames 0-599 and 1400-1999 are used for testing. The rest of the frames are used as training samples.

For the training samples all the locations of each pedestrian is provided, so no extra labeling required for density generation [Chan and Vasconcelos, 2012].

6.3.2 Fudan-ShanghaiTech

The Fudan ShanghaiTech dataset [Fang et al., 2019] is a public dataset with 100 videos of 13 different scenes. Each video contains 6 seconds of footage at 25 fps and have a resolution of 1920x1080 or 1280x720 (Sample of scene in figure 6.2b). The scenes have between 20-100 pedestrians per frame. In each frame the pedestrians in the frame are labeled with a bounding-box and a center-point of the bounding-box. The dataset contains 60 training videos and 40 testing videos.

The lack of trajectories and custom line crossing labeling requires the use of the custom build labeler (Figure 6.1). This is done on the 40 videos of the test set.

6.3.3 AI City Challenge

The AI City Challenge dataset is a huge dataset of cars crossing on a road. The dataset contains 20 different scenes, whereas 11 of the scenes are crossroads. A total of around 2.5 hour footage is captured at 10FPS (1 hour for training, 1.5 hours for testing). Each vehicle which is labeled as a bounding box, which can be used for

the centre-point. Additionally a set of Line of Interest labeling is provided. So no extra labeling for this dataset is required.

Reducing the amount of footage to create a doable benchmark, then define this here as well

Chapter 7

Results

7.1 Fudan-ShanghaiTech

Method (LOI)	MAE	RMAE
Baseline 1	5.762	0.940
Baseline 2	3.273	0.531
Baseline 2+m	1.490	0.291
Proposed+m	1.539	0.281
Proposed+f+m	1.368	0.266

Table 7.1: The LOI results for Fudan

For the Fudan-ShanghaiTech dataset both the LOI performance and the ROI performance is displayed in table 7.1 and table 7.2.

Looking at table 7.1. The original baseline 1 doesn't perform very well on the new task, especially looking at the RMAE. Which is as expected, because the buildup of the model assumes that the velocity map is trained in a supervised manner. Baseline 2 is performing much better and therefore a much stronger baseline.

By aligning the flow map and density map with the maxing filter a huge increase in performance is measured as well. With decreasing both the MAE and RMAE by a factor of 2. Which suggests that the aligning using the maxing filter has some serious benefits for Crowd Crossing Counting.

The proposed model performs slightly worse than our baseline, which is probably due to the multi modal performance of the encoder. However when the Flow map is fed to the density map decoder, the model significantly outperforms the baseline.

Looking at table 7.2, the results show that the proposed model outperforms CSRNet [Li et al., 2018] on this dataset. Adding Flow significantly decreases the ROI performance, which could mean that the model is focussed on the flow features and therefore ignores some of the non-moving pedestrians.

7.2 UCSD

Stated in table 7.3 that the proposed method doesn't produce a state-of-the-art results on the existing UCSD. But more info about UCSD soon in table 7.3.

Incorrect method of measuring the MAE/MWAE, because this is done differently than Fudan....

Method (ROI)	MAE	MSE
Baseline 1	16.402	340.435
Baseline 2	7.569	80.51
Proposed	5.066	50.37
Proposed+f	9.516	132.920

Table 7.2: The ROI results for Fudan

Method (LOI)	TMAE Left	MWAE@100 Left	TMAE Right	MWAE@100 Right
[Ma and Chan,]	0.60	0.72	0.69	0.51
[Zhao et al., 2016]	1.18	0.60	0.63	0.47
[Zheng et al., 2019]	0.46	0.41	0.50	0.41
Baseline 1	8.43	1.40	7.03	1.37
Baseline 2	2.56	0.56	4.85	0.90
Baseline 2+m	1.89	0.46	0.96	0.51
Proposed+m	1.82	0.51	0.88	0.47
Proposed+f+m	4.33	0.78	0.73	0.47

Table 7.3: The LOI results for UCSD

7.3 AI City Challenge

Method (ROI)	MAE	MSE	Method (LOI)	MAE	MSE
Baseline 1	-	-	Baseline 1	-	-
Baseline 2	-	-	Baseline 2	-	-
Proposed	-	-	Baseline 2+m	-	-
Proposed+f	-	-	Proposed+m	-	-
			Proposed+f+m	-	-

Table 7.4: The ROI results for AI City Challenge dataset

Table 7.5: The LOI results for AI City Challenge dataset

To be done about table 7.4 and table 7.5

7.4 Real world performance

To compare real world performance the models are compared on processing speed. Additionally the optical FPS is calculated. This is both done on the Fudan-ShanghaiTech dataset.

Method	FPS	ms	FPS	MAE	MSE
Baseline 2	2.7f	376ms	25	0.0	0.0
Baseline 2+m	2.3	440ms	12.5	0.0	0.0
Proposed+m	4.0	248ms	8.3	0.0	0.0
Proposed+f+m	0.0		5	0.0	0.0
Proposed+f+m+o	0.0		2.5	0.0	0.0

Table 7.6: Processing time

Table 7.7: Optimal FPS

7.5 Flow estimation impact

Qualitative comparison

Chapter 8

Conclusion

In this thesis a new method for Crowd Crossing Counting is presented which can be used multi-domain as well.

All results on the datasets show a clear benefit for a newly created unified model which is focussed on unsupervised flow estimation in comparison to [Zhao et al., 2016].

The first proposed model is performing equally well as the strong baseline. Which was predicted, due to use of the same principles the model design was based on.

The model with flow context clearly performs better on all models, which suggests that flow context indeed pushes the model indirectly to focus more on the moving pedestrians.

The Fudan-Shanghai and UCSD datasets show that the usage of realigning the density map and velocity map is crucial to perform well. The maxing filter is a good solution to artificially align the density map and the velocity map.

The AI City Challenge dataset show that the proposed method is multi-domain as well and performing seemingly well. Additionally the dataset shows that the use of realigning is of less usage for line crossing with objects where the labeling is done in the middle of the object and of sufficient size.

The results are promising for Crowd Counting for used in real world applications. Multiple streams can be processed in real time running on a single GPU. However the method is currently run on a large GPU, which is often still difficult to store. However with the increase in efficiency and performance of current GPU's, the model could soon be running on much smaller equipment.

All in all, even though the proposed model doesn't show state-of-the-art performance on the UCSD benchmark, the proposed model is an interesting new approach to Crowd Crossing Counting.

8.1 Further research

— Due to the lack of high density videos, the full potential of this network can't be shown to its full extend. For further research it would be ideal to propose such a high density map.

— Focussing on even more efficient usage of the method (Aligning is expensive now)

— No real comparison or research is done on other realigning methods. Currently some parameters need to be set, a trainable method would be interesting for further research.

Bibliography

- [Brox et al., 2014] Brox, T., Papenberg, N., and Weickert, J. (2014). High Accuracy Optical Flow Estimation based on warping - presentation. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 3024(May):25–36.
- [Bruhn et al., 2005] Bruhn, A., Weickert, J., and Schnörr, C. (2005). Lucas/Kanade meets Horn/Schunck: Combining local and global optic flow methods. *International Journal of Computer Vision*, 61(3):1–21.
- [Cao et al.,] Cao, L., Zhang, X., Ren, W., and Huang, K. Large scale crowd analysis based on convolutional neural network. 48(10):3016–3024.
- [Chan and Vasconcelos, 2009] Chan, A. and Vasconcelos, N. (2009). Bayesian Poisson regression for crowd counting Cited by me. *Computer Vision, 2009 IEEE 12th International*
- [Chan and Vasconcelos, 2008] Chan, A. B. and Vasconcelos, N. (2008). Modeling, clustering, and segmenting video with mixtures of dynamic textures. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(5):909–926.
- [Chan and Vasconcelos, 2012] Chan, A. B. and Vasconcelos, N. (2012). Counting people with low-level features and bayesian regression. *IEEE Transactions on Image Processing*, 21(4):2160–2177.
- [Dalal and Triggs, 2005] Dalal, N. and Triggs, B. (2005). Histograms of oriented gradients for human detection. *Proceedings - 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2005*, I:886–893.
- [Dollár et al., 2012] Dollár, P., Wojek, C., Schiele, B., and Perona, P. (2012). Pedestrian detection: An evaluation of the state of the art. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4):743–761.
- [Dosovitskiy et al., 2015] Dosovitskiy, A., Fischery, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Smagt, P. V. D., Cremers, D., and Brox, T. (2015). FlowNet: Learning optical flow with convolutional networks. *Proceedings of the IEEE International Conference on Computer Vision*, 2015 Inter:2758–2766.
- [Fang et al., 2019] Fang, Y., Zhan, B., Cai, W., Gao, S., and Hu, B. (2019). Locality-constrained spatial transformer network for video crowd counting. *Proceedings - IEEE International Conference on Multimedia and Expo*, 2019-July:814–819.

- [Horn and Schunck, 1981] Horn, B. K. and Schunck, B. G. (1981). Determining optical flow. *Artificial Intelligence*, 17(1-3):185–203.
- [Hui et al., 2018] Hui, T. W., Tang, X., and Loy, C. C. (2018). LiteFlowNet: A Lightweight Convolutional Neural Network for Optical Flow Estimation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 8981–8989.
- [Idrees et al., 2013] Idrees, H., Saleemi, I., Seibert, C., and Shah, M. (2013). Multi-source multi-scale counting in extremely dense crowd images. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2547–2554.
- [Ilg et al., 2016] Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., and Brox, T. (2016). FlowNet 2.0: Evolution of optical flow estimation with deep networks.
- [Janai et al., 2018] Janai, J., Güney, F., Ranjan, A., Black, M., and Geiger, A. (2018). Unsupervised Learning of Multi-Frame Optical Flow with Occlusions. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11220 LNCS:713–731.
- [Kingma and Ba, 2015] Kingma, D. P. and Ba, J. L. (2015). Adam: A method for stochastic optimization. *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, pages 1–15.
- [Li et al., 2018] Li, Y., Zhang, X., and Chen, D. (2018). Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1091–1100.
- [Lin and Davis, 2010] Lin, Z. and Davis, L. S. (2010). Shape-based human detection and segmentation via hierarchical part-template matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(4):604–618.
- [Liu et al., 2008] Liu, C., Freeman, W. T., Adelson, E. H., and Weiss, Y. (2008). Human-assisted motion annotation. *26th IEEE Conference on Computer Vision and Pattern Recognition, CVPR*.
- [Liu et al., 2019a] Liu, P., King, I., Lyu, M. R., and Xu, J. (2019a). DDFlow: Learning optical flow with unlabeled data distillation.
- [Liu et al.,] Liu, P., Lyu, M., King, I., and Xu, J. SelFlow: Self-supervised learning of optical flow.
- [Liu et al., 2019b] Liu, W., Salzmann, M., and Fua, P. (2019b). Context-aware crowd counting. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June:5094–5103.
- [Ma and Chan,] Ma, Z. and Chan, A. B. Counting people crossing a line using integer programming and local features. 26(10):1955–1969.

- [Ma and Chan, 2013] Ma, Z. and Chan, A. B. (2013). Crossing the line: Crowd counting by integer programming with local features. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2539–2546.
- [Mémin and Pérez, 1998] Mémin, E. and Pérez, P. (1998). Dense estimation and object-based segmentation of the optical flow with robust techniques. *IEEE Transactions on Image Processing*, 7(5):703–719.
- [Pock et al., 2008] Pock, T., Schoenemann, T., Graber, G., and Bischof, H. (2008). Learning optical flow. 5304(May 2014).
- [Ranjan and Black, 2017] Ranjan, A. and Black, M. J. (2017). Optical flow estimation using a spatial pyramid network. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-January:2720–2729.
- [Sreenu and Saleem Durai, 2019] Sreenu, G. and Saleem Durai, M. A. (2019). Intelligent video surveillance: a review through deep learning techniques for crowd analysis. *Journal of Big Data*, 6(1):1–28.
- [Subburaman et al., 2012] Subburaman, V. B., Descamps, A., and Carincotte, C. (2012). Counting people in the crowd using a generic head detector. *Proceedings - 2012 IEEE 9th International Conference on Advanced Video and Signal-Based Surveillance, AVSS 2012*, pages 470–475.
- [Sun et al.,] Sun, D., Yang, X., Liu, M.-Y., and Kautz, J. PWC-net: CNNs for optical flow using pyramid, warping, and cost volume.
- [Wan and Chan, 2019] Wan, J. and Chan, A. (2019). Adaptive density map generation for crowd counting. *Proceedings of the IEEE International Conference on Computer Vision*, 2019-Octob:1130–1139.
- [Wang et al., 2020] Wang, Q., Gao, J., Lin, W., and Li, X. (2020). Nwpu-crowd: A large-scale benchmark for crowd counting. *arXiv preprint arXiv:2001.03360*.
- [Wedel et al., 2009] Wedel, A., Cremers, D., Pock, T., and Bischof, H. (2009). Structure- and motion-adaptive regularization for high accuracy optic flow. *Proceedings of the IEEE International Conference on Computer Vision*, pages 1663–1668.
- [Wu and Nevatia, 2007] Wu, B. and Nevatia, R. (2007). Detection and tracking of multiple, partially occluded humans by Bayesian combination of edgelet based part detectors. *International Journal of Computer Vision*, 75(2):247–266.
- [Yu et al., 2016] Yu, J. J., Harley, A. W., and Derpanis, K. G. (2016). Back to basics: Unsupervised learning of optical flow via brightness constancy and motion smoothness. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 9915 LNCS:3–10.

- [Zhang et al., 2016] Zhang, Y., Zhou, D., Chen, S., Gao, S., and Ma, Y. (2016). Single-image crowd counting via multi-column convolutional neural network. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016-December:589–597.
- [Zhao et al., 2016] Zhao, Z., Li, H., Zhao, R., and Wang, X. (2016). Crossing-line crowd counting with two-phase deep neural networks. In Leibe, B., Matas, J., Sebe, N., and Welling, M., editors, *Computer Vision – ECCV 2016*, volume 9912, pages 712–726. Springer International Publishing.
- [Zheng et al., 2019] Zheng, H., Lin, Z., Cen, J., Wu, Z., and Zhao, Y. (2019). Cross-line pedestrian counting based on spatially-consistent two-stage local crowd density estimation and accumulation. 29(3):787–799.

Appendix A

Appendix

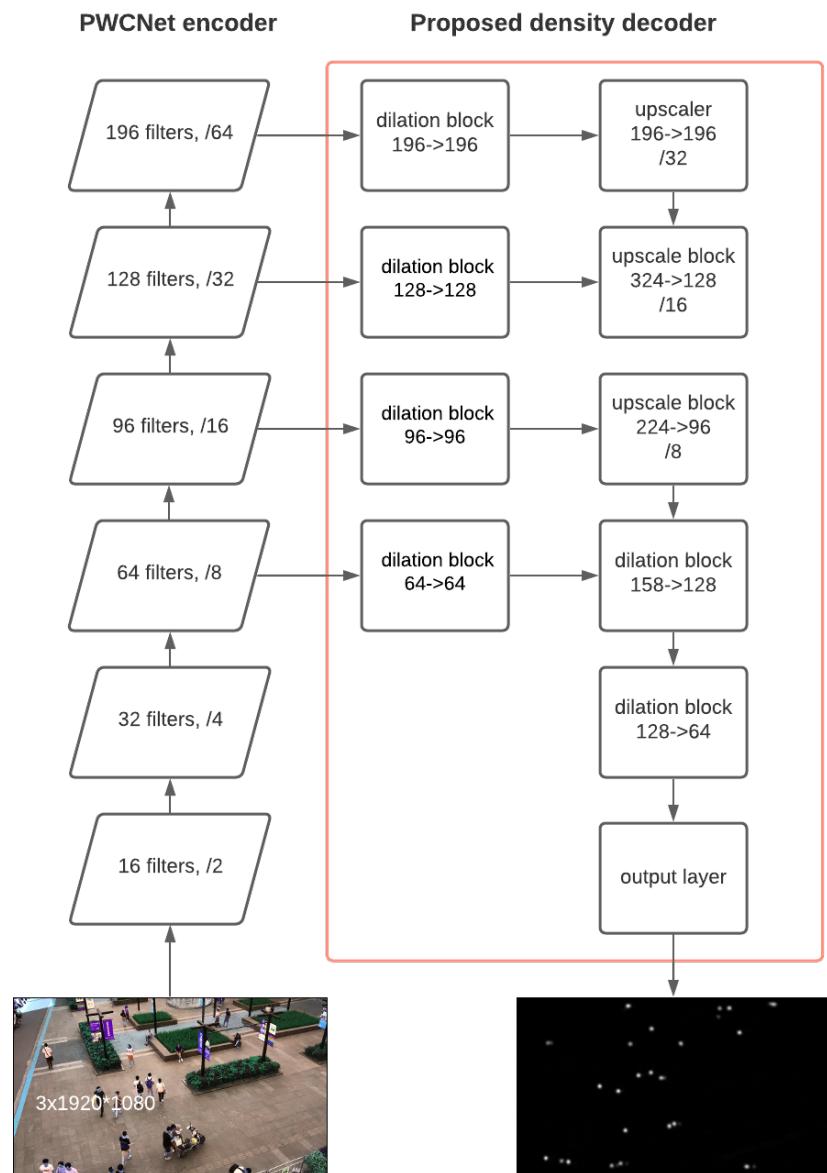


Figure A.1: Full decoder for density prediction