# A Human and Group Behaviour Simulation Evaluation Framework utilising Composition and Video Analysis

Rob Dupre
Kingston University
R.Dupre@kingston.ac.uk

Vasileios Argyriou
Kingston University
Vasileios.Argyriou@kingston.ac.uk

## Abstract

*In this work we present the modular Crowd Simulation Evaluation through Composition framework (CSEC) which provides a quantitative comparison between different pedestrian and crowd simulation approaches. Evaluation is made based on the comparison of source footage against synthetic video created through novel composition techniques. The proposed framework seeks to reduce the complexity of simulation evaluation and provide a platform from which the comparison of differing simulation algorithms as well as parametric tuning can be conducted to improve simulation accuracy or providing measures of similarity between crowd simulation algorithms and source data. Through the use of features designed to mimic the Human Visual System (HVS), specific simulation properties can be evaluated relative to sample footage. Validation was performed on a number of popular crowd datasets and through comparisons of multiple pedestrian and crowd simulation algorithms.*

## Introduction

Pedestrian and crowd simulation has applications in a wide range of industries including pedestrian facility suitability and capacity [1], computer graphics and gaming [2], the social sciences [3] and engineering [4]. This broad range of uses has led to extensive research into how crowds and pedestrians move around and interact with their environment.

The problem of how to evaluate these simulation algorithms is a developing area of research. One of the most prominent issues with crowd and pedestrian simulation research is the lack of a simple and suitable form of comparison between different simulation and modelling approaches. This often means that a given methodology is developed and evaluated for a specific purpose, with its wider abilities and properties left unconfirmed. This task is made more difficult as the developed approaches cover a huge range of applications, where evaluation techniques for one are not always applicable to the others.

Generally the evaluation techniques utilised can be split into qualitative [5] and quantitative measures [2, 1]. The former including assessments made by experts in the field or context of the intended application [3], as well as category based rating systems [6] designed to define the capabilities of an algorithm (such as emergent behaviours). These assess whether the simulation *looks* natural and that the agents within the simulation are not acting in an unusual fashion.

A number of quantitative measures have been suggested to provide a numeric measure of accuracy for a simulation, which include but are not limited to: speed, pedestrian density, number of steps taken to destination and duration. These evaluation techniques tend to be data driven, and as such require some kind of ground truth data from which to test against. The concept of an evaluation framework has been suggested before [7, 8, 9, 10, 11, 12]; with most deducing various metrics based on a simulation in an effort to rate simulation algorithms or tune parameters. Often these frameworks evaluate the quality of a simulation based on data driven methods i.e. how closely does simulated agent A's track match pedestrian A's track in the source data. This relies on the assumption that a good simulation must mimic captured source data exactly, however, humans moving through the same environments on a regular basis will look similar but have slightly different properties, rendering this assumption flawed.

Many of these evaluation frameworks have merit in their given context, however to make a comparison, often a number of requirements are imposed on their source data. Most commonly this pertains to tracks for the pedestrians in source data, this introduces issues to the data collection process pertaining to cost, time, ethics and suitability for large outdoor environments. Additionally the focus of these frameworks is a statistical analysis of the simulations, specifically on individual agents rather that of the simulation as a whole meaning the way the simulation appears is often overlooked.

Figure 1: Source CCTV footage and generated composition video with computer controlled agents.

To address these gaps in the existing research a new modular evaluation framework is proposed allowing comparison of a simulation algorithm to source video footage using limited ground truth data. Using novel compositing techniques, a video comprised of the 2D background of a source video and superimposed 3D agents (controlled by the simulation algorithm), can be created (Figure 1). Using this simulated video, a direct comparison against the source video can be performed using Human Visual System (HVS) features [13], analysing crowd properties such as density, speed and track. Due to the modular nature of the framework, both the composition techniques as well as the video analysis methods are interchangeable as required. This ensures the method is future proofed to new and improved methodologies.

Due to the modular nature of the framework, any simulation algorithm can be used, providing the facility to compare the performance of different simulation algorithms relative to source material. The framework also provides the functionality for model tuning, by creating a fast feedback loop which allows the adjustment of model parameters to improve simulation accuracy. Furthermore as the simulated composite videos produced as part of the proposed framework have a known groundtruth, they are very suitable for the evaluation of pedestrian tracking algorithms.

This framework and the associated features are motivated by the work in [13] and looks to improve on the ideas proposed. As such, the following novel framework is suggested which reduces the complexity of crowd and pedestrian simulation evaluation, providing a quantitative comparison regardless of context. Simulated video sequences are created using sample video data and crowd simulation algorithms, combined using novel compositing and visualisation techniques. Through the use of background subtraction and scene composition to generate simulated video sequences, lengthy scene reconstruction steps are eliminated simplifying the overall process. The paper will continue as follows: firstly, an overview of the research in this field is given, next, the proposed modular framework and compositing methodology are introduced and an overview of HVS features presented. Finally the evaluation process is reviewed and conclusions drawn.

## Literature Review

Within this section a brief overview of current simulation algorithms is given, followed by a review of existing evaluation metrics and frameworks.

### Simulation Algorithms

Simulating behaviour virtually has been an area of intense research in recent years, one of the first agent-based simulation algorithms was proposed by Reynolds in [14, 15] which focused on birds flocking, and was later developed and better defined for gaming applications. Later Helbing et al [16] introduces the Social Force Model (SFM). This became widely successful due to its ability to emulate common attributes seen in pedestrian movements and emergent behaviours in crowds such as line forming in tight areas. Another key advantage of this model was the use of variables that related to physical principles in our world. The use of these parameters allowed the application of other forms of research to drive the simulation and formed a basis for evaluation. Another example of this type of approach is proposed by Xi et al [4]. Here a dense model is proposed integrating extended decision field theory, the social force model and a dynamic planning algorithm involving AND/OR graphs. Extensive testing is done on potential profit for a shopping mall where factors of an agents AI, for example group dynamics, visual field or intention to buy, are changed. However no real validation based on real results is presented. Survey and observation data is used in the setting of these model parameters.

With the advent of so much research in the area of pedestrian simulation, the ability to determine what is a good simulation has been well researched and is very much a topic for debate. As such work has been done to survey the current research with a number of different simulation algorithms and models being put through the same test environments [6, 17, 18, 19, ?, ?]. Evaluations are based on properties such as computational performance as well as the presence of emergent behaviours and a model's abilities, such as route choice and agent strategy. Additionally numerical metrics are introduced spanning areas such as speed and acceleration, distance covered, time and collisions. However there is no real analysis of the quality of simulated visual experience, i.e if the methodology produces visually similar or natural behaviour, which is often a key requirement of the design.

In more individual cases it is seen that often the evaluation technique is limited by the context. This of course makes sense but often means that other key aspects of a simulation implementation are not analysed. For example in evacuation simulations, the ground truth data often does not exist from which to compare, leaving only qualitative

assessments by experts in the field [3] leaving any other properties that these algorithms have un-explored.

Portez et al [5] focuses on the simulation of crowds around bottle necks and looks for events. Here density matching against recorded video footage is used as a quantitative measurement, more specifically number of people per square meter. This is backed up using visual checks against the original footage to ensure the simulated crowds resemble those in the captured data. Pettre et al [20] introduces an agent interaction method and uses density plots based on aspects of the simulation, such as reaction times, as well as a defined likelihood function to perform their evaluation. The likelihood function is based on simulation variables and the assessed difference from captured source data, however these are only implemented on interaction scenarios comprising of two people and as such their validity as metrics for larger scale crowd simulation remains untested.

Kim et al [2] manages a quantitative assessment of modelled scenarios by looking at psychological studies of specific situations. I.e looking for a simulation to repeat observations made. For example, a study of pedestrians average crossing speeds relative to the start of the crossing signal [21]. This is then compared with the outputs of the simulation to provide a similarity measure. This is a more abstract approach as they are not comparing specific paths but rather duration, focusing on the properties of the simulation rather than specific paths.

## Simulation Evaluation Frameworks

Charalambous et al [7] looks at the creation of an analysis tool that characterises outlying behaviour in simulations. Two processes are suggested; outlier detection which searches for odd behaviour within that simulation and novelty detection, which finds trends or actions that differ from some reference data. However for a good analysis, reference data must be very similar to the simulated. Additionally as the comparison made is purely a data driven approach the concept of something looking similar is not addressed.

Guy et al [9] uses a computed entropy score to compare simulated data to captured real world data. The metric is defined as the entropy of the distribution of errors between the evolution of a crowd predicted by a simulator and the source data. Using three differing datasets ranging from simplistic to dense crowds, evaluations are done to produce simulations that closely resemble the source data. This statistical analysis again relies on the need for position information from both the simulation and comparable real world examples. Additionally a number of assumptions are made, the most noticeable being that the crowd simulator is not systemically more accurate for some agents within a crowd than for others, however this is not always accurate as their will always be aspects of a simulation that are more accurate than others [7].

Wang et al [22] present the Stochastic Variational Dual Hierarchical Dirichlet Process (SV-DHDP) model in which groups of similar trajectories (trending paths) can be combined to generate an overall path pattern for an environment. The path patterns created are therefore the result of local dynamics and global factors allowing differing insights based on the simulation environment. The resultant visualisations allow for detailed qualitative analysis and the introduction of an inference based similarity metric allows for the comparison of extracted path patterns from differing data sources. Providing a good generalised view of a given scene. However analysis is done on defined paths for source and test data which requires complex post processing techniques or data captured in a specific format which can create inaccuracies [1], more accurate systems can be used [23] but are not always cost effective due to the requirement of specialized equipment.

Lerner et al [24, 25] address the concept of look and feel of a crowd by assigning a similarity metric by comparing an agent's actions at a given moment in time to a database of observed actions. The database is taken from annotated frames of video, defining path vectors for pedestrians in both sparse and dense crowds. A state-action pair for each frame is defined using firstly a state (set of recorded variables such as trajectory, speed and position) and an action (a density measure representing local density changes over time). The similarity between a state-action pair from the database and a test state-action pair is defined as the similarity between the actions (differences in positions along the trajectories) and the distance between the states (differences in densities for the surrounding regions).

Musse et al [12] also address the issue of tracking generalised paths in crowds using four dimensional histograms to describe movement within a crowd. By applying the Bhattacharyya distance as a form of measurement between the defined features similarity is assessed on criteria such as speed, spacial occupancy and orientation. The results produce similarity measurements for aspects of orientation and speed but fail to take into account the density of the crowds. Additionally no analysis of what is visually similar is given which highlights a systemic problem with many of the similarity features suggested such that evaluation is given based on similarity to extracted trajectories but not on if the results look the same to a group of people.

Additionally the work in [11, 26] proposes interesting similarity metrics and [10, 27] are also worth mentioning.

The proposed framework is focusing on providing a statistical analysis of the realism in different simulations. In order to evaluate the realism of a crowd or pedestrian simulation algorithm, vision based features are utilised. Motion estimation [28, 29] and tracking are a few of the vision based steps applied during the process of pedestrian

and crowd behavior analysis for visual surveillance in dynamic scenes [30, 31].

The majority of today's optical flow methods strongly resemble the original formulation provided by Horn and Schunck [32] as well as the work by Lucas and Kanade [33]. The accuracy and robustness of optical flow estimation algorithms has seen significant improvement over the last decade [34, 35, 36]. Additionally work in tracking pedestrians and crowds has also seen much work [37, 38]. Specifically techniques have been developed for estimating the flux of people in public areas, such as stores or travel sites, which can then automatically provide congestion analysis assisting in management of crowds and pedestrians [39, 40, 41].

A technique that incorporates optical flow for accuracy evaluation in crowd simulation was proposed in [42]. In this work a solution is proposed which allows the relationship of optical flow to physical velocity to be defined. The main issue of this approach is that it requires manual annotation and performs well only for specific relative orientations of the camera and pedestrians.

All these approaches for pedestrian and crowd simulation evaluation either do not consider the quality of simulated visual experience, or fail to compare their given metric with human validation. Additionally the requirement for data with a defined groundtruth introduces issues pertaining to cost, time, ethics and suitability for large outdoor environments. All key components when considering the accessibility and accuracy of crowd simulation evaluation.

## Methodology

### Simulation Evaluation using Compositing Techniques

The following section will explain in detail the various aspects of the proposed modular Crowd Simulation Evaluation through Composition (CSEC) framework. The framework provides a method of simulation algorithm evaluation which rates how realistic simulated human walking behaviours look compared to sample footage. Evaluation can be done on a frame by frame basis or on a sequence as a whole, providing flexibility in how the simulation is evaluated. Additionally the proposed methodology requires no track or path information for the source material, allowing any pedestrian video footage captured from a static viewpoint to be used as source material. Evaluation of an algorithm's performance is key to defining how realistic the simulation outputs are.

Comparison is made using the original source footage and a video created using composition techniques. This utilises background subtraction techniques as well as methods to extract 3D data from 2D images. This allows the construction of a 3D space in which virtual agents can navigate around. Through the use of composition, a final visualisation combining this background and 3D space is generated to form the simulated video sequence in which simulated agents are superimposed into the background of the source video data. Analysis of these generated videos is through the use of features designed to evaluate the visual similarity of the two videos to provide a quantifiable similarity metric. These are designed to emulate the way the Human Visual System (HVS) perceives motion, and include the principles of Weber's Law [43] to better match the metrics to the way humans see.

Fundamentally the framework is made up of two components; simulation visualisation and video similarity. Importantly the modular nature of the framework supporting inputs of any simulation algorithm or video analysis techniques, depending on application, whilst still retaining the ability to produce a quantifiable similarity score (Figure 2). The HVS features [13] provide a good generalisation of simulation evaluation requirements in a broad range of situations, however additional feature descriptors could be developed and inserted into the framework to determine other crowd statistics such as lane formation and direction of travel. Figure 3 provides a more specific overview of the CSEC framework as it is utilised in this work. Further detail of each section will follow.

As the framework compares video data to derive a similarity value, firstly a simulated video must be constructed. Initially, using the source video sequence, the background is extracted. Next, a two dimensional plane is extracted representing a top down view of the given scene. Simulations are run to produce paths for virtual agents to follow based on the extracted plane. The visualisation component is used to composite the extracted 2D background image and 3D rendered agents as they follow the simulated paths. Frames are output from the visualisation into a final simulated video sequence (Figure 4). Once both a simulated and source video are available, the similarity can be evaluated. Optical flow and tracklet analysis are run and features extracted from the subsequent data. Finally a distance measure is used to evaluate the difference in features to give the final similarity metric.

#### Background and Plane Extraction

To allow the composition of the simulated video to be created, the background of the source video sequence is required. For this work the mean value of each pixel in a video sequence is used to create a compound background image. Other methods based on Gaussian mixture models could also be used in order to obtain more accurate results [44].

Once the background image has been subtracted the process of defining the perspective grid is applied. The per-
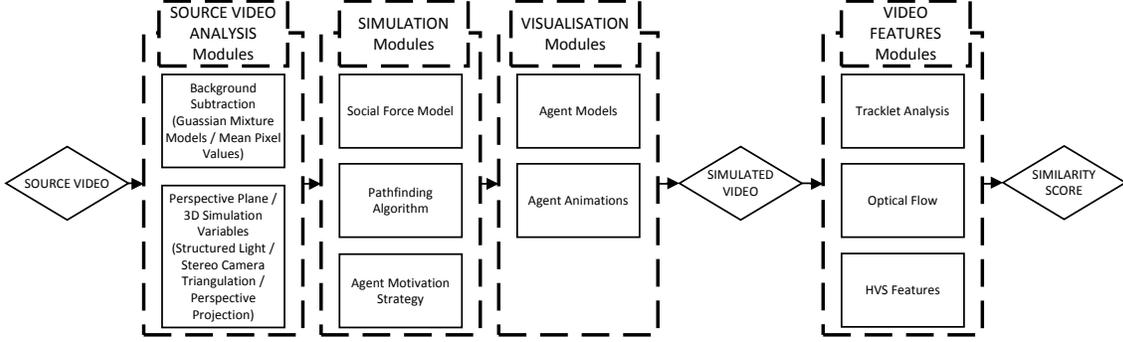
Figure 2: Outline of the modular nature of the Crowd Simulation Evaluation through Composition (CSEC) framework. Using an initial source video, an analysis stage produces the inputs for simulation and visualisation. Simulation generates the paths that the agents will follow and the visualisation modules combine these aspects to create the simulated video. The video features module is concerned with the comparative aspect of the framework and the similarity score created as an output.
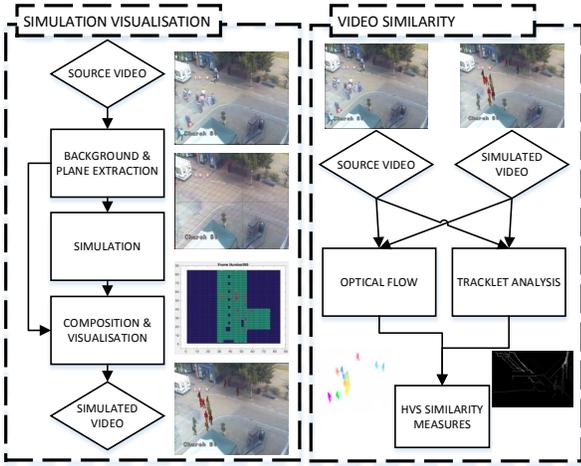


Figure 3: Overview of the Crowd Simulation Evaluation through Composition (CSEC) framework.



Figure 4: Frames of source CCTV footage and generated video using the composition techniques.

spective grid allows scale mapping of an environment from the viewpoint of the source video camera pose. The resultant grid represents a top down environment map of the viewable area and is used during agent simulation. Using the concept of perspective scale along a line we can, through the definition of two parallel lines that run to the vanishing point of an image, estimate distance in arbitrary units of measure within this perspective space (Figure 5b). This unit can be based upon an object in the scene with known dimensions or using pedestrians [45].

Initially the user defines the points $\mathbf{i}$ and $\mathbf{j}$, in the 2D image space, forming a line along an edge that leads to the vanishing point of the image. A second line is defined by the points $\mathbf{k}$ and $\mathbf{l}$, such that it runs 'parallel', relative to the vanishing point in the 3D space of the captured image, to the line defined by points $\mathbf{i}$ and $\mathbf{j}$ (Figure 5a).

At a location along the line $\mathbf{ij}$ the user defines another point $\mathbf{u}_1$, such that the line $\mathbf{iu}_1$ represents the unit of distance $m$ from which all further perspective points are defined. An additional point $\mathbf{u}_2$ is defined on top of the line $\mathbf{ik}$ which represents the same relative distance in 3D space as $m$.

For the next step of the proposed algorithm the reference points $\mathbf{T}_{vanish}$, $\mathbf{R}$, $\mathbf{R}_0$ and $\mathbf{T}_{n-1}$ are initialised automatically (Figure 5a).
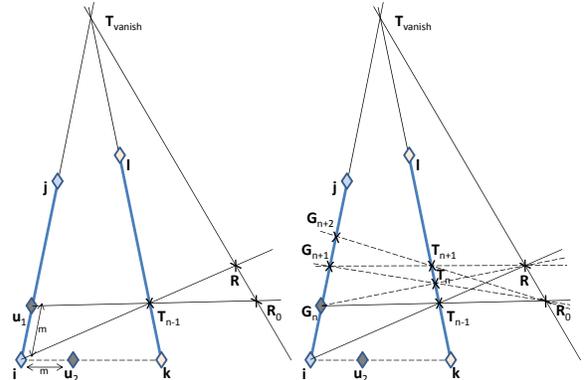


Figure 5: (a) User defined points and initialisation. (b) The first two iterations of the recursive algorithm.

In more detail, the vanishing point $\mathbf{T}_{vanish}$ is defined as

5

the point at which the lines **ij** and **kl** intersect, this may well be at a position outside of the image plane. As such $\mathbf{T}_{vanish}$ is defined as

$$\mathbf{T}_{vanish} = f(\mathbf{i}, \mathbf{j}, \mathbf{k}, \mathbf{l}) \qquad (1)$$

An arbitrary point $\mathbf{R}$ is selected at a random location outside the triangle $\mathbf{i}\mathbf{T}_{vanish}\mathbf{k}$. The point $\mathbf{T}_{n-1}$ is defined as the point of intersection of the lines $\mathbf{i}\mathbf{R}$ and $\mathbf{k}\mathbf{T}_{vanish}$

$$\mathbf{T}_{n-1} = f(\mathbf{i}, \mathbf{R}, \mathbf{k}, \mathbf{T}_{vanish}) \qquad (2)$$

Finally the point $\mathbf{R}_0$ is defined.

$$\mathbf{R}_0 = f(\mathbf{u}_1, \mathbf{T}_{n-1}, \mathbf{R}, \mathbf{T}_{vanish}) \qquad (3)$$

With these points initialised a recursive algorithm is applied to calculate equidistant points along the line $\mathbf{i}\mathbf{T}_{vanish}$ in 3D space. As the user has already defined the first of these points $\mathbf{u}_1$, for the purposes of the recursive step, these will be relabeled as $\mathbf{G}_n$. This is a two-step iterative process, with the point $\mathbf{T}_n$ being defined as the intersection of the lines $\mathbf{G}_n\mathbf{R}$ and $\mathbf{k}\mathbf{T}_{vanish}$.

$$\mathbf{T}_n = f(\mathbf{G}_n, \mathbf{R}, \mathbf{k}, \mathbf{T}_{vanish}) \qquad (4)$$

and during the second step the next equidistant point $\mathbf{G}_{n+1}$ on the line $\mathbf{i}\mathbf{T}_{vanish}$ is expressed as a function of

$$\mathbf{G}_{n+1} = f(\mathbf{R}_0, \mathbf{T}_n, \mathbf{i}, \mathbf{T}_{vanish}) \qquad (5)$$

This process is repeated until $\mathbf{G}_{n+1}$ is no longer within the borders of the original background image.

The grid is initially defined using all the equidistant points on the line **ik** using the distance $\mathbf{i}\mathbf{u}_2$ as a unit. Lines are defined between each of these points and the vanishing point $\mathbf{T}_{vanish}$ of the image. The scale points $\mathbf{G}$ are plotted along each of these newly defined lines forming the grid. Additionally if required, the recursive process can be inverted to create points moving away from the vanishing point. This ensures that the entire image plane is encapsulated by the defined grid, regardless of where the user has defined their points.

The resultant grid represents the perspective plane of the source image. On that grid the areas (cells) with obstacles (i.e. cells where pedestrians cannot walk) are annotated as is information about entrance/exit locations. In order to help the user; the obtained grid is superimposed on the extracted background image (Figure 6). Here red cells indicate areas where agents can walk, white represent obstacles and green marks an entrance or exit. This annotated version of the perspective plane is then used as the ground plane by the simulation algorithms.
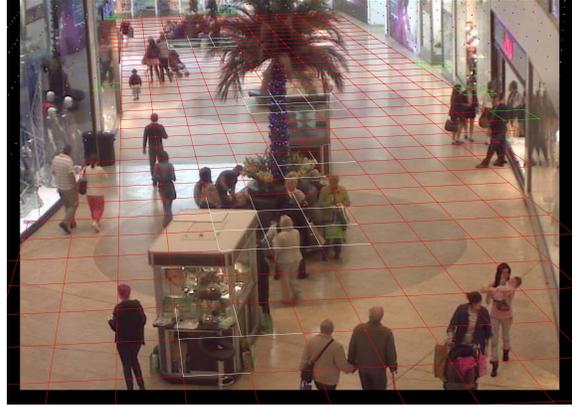


Figure 6: Resultant perspective grid overlayed on the original source image. Red cells indicate areas where agents can walk, white represent obstacles and green marks an entrance or exit.

**Pedestrian and Crowd Simulation Algorithm**

To simulate the agent movement through a given scene an algorithm based on a combination of simulation methodologies is used. Firstly a steering simulator based on the work of [15] in which the concepts of simple crowd behaviours such as separation, object avoidance and agent collision detection are utilised. These have been implemented with the social forces model structure in which each of these elements produce a force applied to the agent to adjust their movement vector. The magnitude of these forces is based on distance. An additional step, using a planning simulation methodology, based on the work of Karamouzas et al [46] is used as a predictive collision detection algorithm to produce natural agent avoidance within the simulations, this again is implemented by the application of a force upon the simulated agent. As outlined in (6).

$$F_a = g_a - p_a + \sum_{i=1}^{n} f(a, b_i) + \sum_{j=1}^{m} f(o_j) + \sum_{k=1}^{o} f(a, b_k) \qquad (6)$$

where $g_a$ is the current destination along the path of the agent $a$ to its final goal, with $p_a$ being the agent's current position. The forces for separation, $f(a, b_i)$, object avoidance $f(o_j)$ and the predicative agent avoidance $f(a, b_k)$, is calculated for any relevant entity within a defined neighbourhood. For overall path planning, an agent performs a route plan using the A* algorithm and the perspective plane obtained previously to estimate the most direct course from their start location to their destination. The variables associated with defining an agent and their respective movements are based on existing work used by Asano et al. [1, 47] who in turn derive their values from the existing studies and from

Figure 7: Example composition of the Kvan scene with test agents and perspective floor plan.

observational data from their datasets.

## Composition and Visualisation

The visualisation stage of the framework performs the composition of a scene utilising the extracted background obtained from the source video and the generated perspective plane. The key to a visually similar composition is the positioning of a virtual camera at the same location as in the original scene. By using layers the camera can have the source image as a background and the visualised 3D agents controlled by the simulation superimposed. Due to this, it is important to line up the perspective plane with the background to give the illusion of the agents walking through the scene. This alignment can be performed manually using the position and orientation of the camera or automatically using camera calibration techniques [48, 49].

Sample agents are then placed in the scene at various locations to ensure that the perspective and scaling parameters of the agents are appropriate to the scene. Figure 7 demonstrates this with agents in blue positioned at different locations in the scene. These values can be adjusted manually or calculated automatically using the methods provided in [50]. In Figure 7 it can also be seen that the imported floor plan is coloured according to each cell's defined values, green meaning areas the agents can walk and blue defines entrance exit points and red represents obstructions that will obscure the agents when they are located behind. To control the agents in the scene, position and orientation information is required for each frame. This is obtained using the desired simulation algorithm and the same perspective grid map outlined in the Simulation section.

As the goal is to create videos with similar crowds, a number of parameters from the original video are required. Using the source video sequence, an analysis is made of the pedestrians in the scene, outlining paths and estimated crowd density. For this work the information was extracted manually or provided by the datasets used, however work exists to help automate this process [22, 11]. Once all the required parameters for the agents are defined the simulation is run and outputs recorded at the frame rate of the original source video.

Finally, with the composition completed and the simulations run the visualisation of the scene is performed. Agent models are created and sized according to the obtained parameters. For each frame of the simulation the agent location and rotation is updated based on the simulation algorithm output and a composite frame is captured. Once visualisation is completed the individual frames are compiled into a video sequence. Importantly the resolution and the number of frames in the new composite video should be equal to that of the source video.

## Simulation Similarity Metrics

Once visualisation of the composite video is complete, the source and the simulated video sequences are used to extract features in order to measure their level of similarity. These features are based on the optical flow and tracklets of the moving objects in both sequences.

An optical flow method tries to calculate the motion between two image frames at times $t$ and $t + \delta t$ at each pixel position [33]. The solution as given by Lucas and Kanade is a non-iterative method, which assumes a locally constant flow in a small window. Black and Anandan in [34], describe how the single motion assumption, as well as the constant brightness constraint are not always valid. They discuss how these assumptions can be relaxed in order to develop a more robust estimation framework.

Tracklet estimation is a well researched topic with many algorithms available in the literature. These can be based on motion or other features and utilise particle and Kalman filters [37, 38, 51, 52]. Specifically, the problem of motion based tracking can be split into detecting moving objects in each frame and the association of those moving elements to a continuous corresponding object over time.

In the case of Kalman filters, the track's location in each frame is predicted and a likelihood of a detection is assigned to each track. The Kalman filter is a recursive estimator, meaning that only the estimated state from the previous time step and the current measurement are needed for computation of the current state. The Kalman filter has two distinctive features; firstly its mathematical model is described in terms of state-space concepts; Secondly, the solution is computed recursively. Usually the Kalman filter is described by a system state model and a measurement model.

In this system the optical flow algorithm proposed in [34] and the tracking method presented in [51] were utilized,

however the system is designed in such a way that allows the incorporation of multiple motion estimation or tracking methods as plugins. Based on this system architecture the proposed evaluation framework is dynamic and capable of utilizing current and future state of the art tracking methods.

**Motion and Tracklet Flux Similarity Metrics**

In order to evaluate the similarity level of the simulated and source videos a new metric is required that will allow an objective comparison incorporating the Human Visual System (HVS) based similarity features. Weber's Law [43] and the work in [53, 54] states that a human's ability to define motion, as the point when the signal-to-noise ratio is regarded as at a stimulus intensity. Therefore, the minimum motion contrast $dV$ as a function of background motion $V$, required for the human visual system to notice a change is expressed as:

$$dm = L\frac{dV}{V} \qquad (7)$$

where $dm$ is the differential change in motion perception, $dV$ is the differential increase in the velocity and $V$ is the velocity. The parameter $L$ is to be estimated using experimental data. The proposed measure includes Fechners Law, which relates velocity $V$, to perceived motion, $\mathbf{M}$, as seen by the human visual system, as follows:

$$\mathbf{M} = Lln(\frac{V}{V_{max}}) \qquad (8)$$

where $V_{max}$ is the upper threshold of the human eye. The proposed metric is based on the motion and tracklet flux histograms obtained from the perceived motion $\mathbf{M}$ utilizing standard computer vision algorithms.

Let us assume that $I_R(\vec{u},t)$ and $I_S(\vec{u},t)$ are the image frames of a real and the correspondent simulated scene, respectively. The motion vectors for each pixel location in each frame are estimated using the optical flow techniques shown in (9) and (10).

$$M_R(\vec{u},t) = f(I_R(\vec{u},t), I_R(\vec{u},t-1)) \qquad (9)$$
$$M_S(\vec{u},t) = f(I_S(\vec{u},t), I_S(\vec{u},t-1)) \qquad (10)$$

The estimated tracklets are obtained using motion information and Kalman filters.

$$T_R(n_R,\vec{u},t) = f(M_R, I_R) \qquad (11)$$
$$T_S(n_S,\vec{u},t) = f(M_S, I_S) \qquad (12)$$

Since the motion vectors and the tracklets are available the Histogram of Oriented Optical Flow (HOOF) [55] is calculated both for the real and simulated scenes.

$$f_R^{HOOF} = HOOF(M_R) \qquad (13)$$
$$f_S^{HOOF} = HOOF(M_S) \qquad (14)$$

Also, a 2D histogram of the motion parameters is obtained using (15) and (16).

$$f_R^{H2D}(r_{ij}) = m_{ij}(M_R) \qquad (15)$$
$$f_S^{H2D}(r_{ij}) = m_{ij}(M_S) \qquad (16)$$

where $r_{ij}$ is the $i^{th}$ and $j^{th}$ motion level in an interval, and $m_{ij}$ is the number of pixels in all the given frames whose motion level is $r_{ij}$. Regarding the tracklets, the time parameter in (11) and (12) is removed by superimposing all of them at the same time instance. The similarity metric here can be applied on any given time interval, which can be the whole sequence or a small time fragment. In the same way as in (15) and (16) we obtain:

$$f_R^T(r_{ij}) = m_{ij}(T_R) \qquad (17)$$
$$f_S^T(r_{ij}) = m_{ij}(T_S) \qquad (18)$$

Finally, the flux of the features in (13) - (18) is represented by the surface integral of the given vector field.

$$\Phi(\vec{u},t) = \Sigma_{\vec{u}}\Sigma_t f\, du\, dt \qquad (19)$$

Based on (19), we obtain $\Phi_R^{HOOF}$, $\Phi_S^{HOOF}$, $\Phi_R^{H2D}$, $\Phi_S^{H2D}$, $\Phi_R^T$ and $\Phi_S^T$ that correspond to the proposed HVS features. All the features can be applied either on the whole sequence or on smaller blocks allowing specio-temporal adaptation of the proposed features and metrics. In order to measure the similarity and rank the algorithms, the Bhattacharyya distance is utilised due to its use in similar work [12].

## Results

To evaluate the proposed Crowd Simulation Evaluation through Composition (CSEC) framework, a total of five different scenes were used from various crowd datasets (Mall Dataset [56], PETS2009 [57] and RBK [58]) and captured crowd and pedestrian videos sequences. Scenes of different environments including both indoor and outdoor spaces, with a large range of camera orientations and crowd configurations. Additionally the frame rates of the videos varied from less than 10fps up to 24fps providing a challenging

Figure 8: Example source and simulated frames. Top Row (Source): Road, Mall, Krad2. Bottom Row (Simulated): Road, Mall, Krad2.

| # Agents | Agents Speed | | | |
|---|---|---|---|---|
| | Very Slow | Slow | **Same** | Fast |
| Few | 9.32 | 9.07 | 8.08 | 8.01 |
| **Same** | 9.42 | 8.41 | <u>7.33</u> | 7.91 |
| Many | 9.41 | 8.83 | 8.95 | 9.47 |

and diverse set of scenes from which to evaluate the effectiveness of the evaluation framework.

Composite simulated videos for each of the tested scenes were created using four different levels of agent speed and three different population levels, totalling 12 simulations and therefore 12 composite videos per scene. This demonstrates the framework's ability to evaluate source footage against a single simulation algorithm. An additional 12 composite videos were created using Reynolds' flocking algorithm [14], providing a comparative test between simulation algorithms and demonstrating a comprehensive assessment of the relative features of the framework. Figure 8 presents example source and simulated frames for a few scenes.

For each scene the background image was extracted and the perspective grid defined. Simulations were performed for each of the cases mentioned previously and the outputs used to create composite visualisations for each. Simulated videos were then created with the same frame rate, length and resolution as the source videos. The simulations themselves are run at a set frame rate (50 frames per second), a desired frame rate is also specified allowing for the movement of the agents to be visualised at the same frame rate as the source video without having to adjust the agent properties between simulations.

Each simulated video, and its respective source video, had the optical flow and tracklets estimated. Finally the HVS features [13] were used to compare each visualisation

against its source. Three features were used in the comparison, the tracklets ($\Phi_R^T$ (11) and $\Phi_S^T$ (12)), Histogram of Orientated Optical Flow per frame ($\Phi_R^{HOOF}$ (13), $\Phi_S^{HOOF}$ (14)) and the Histogram of Orientated Optical Flow for the sequence ($\Phi_R^{H2D}$ (15), $\Phi_S^{H2D}$ (16)).

The HOOF features and 2D Histograms used a window size of $64 \times 64$ pixels per frame. Using these features, a generalised statistical measure of the differences in movement from the source human behaviour to the simulated agents is defined. The distance metric used to compare the features is the Bhattacharyya Distance [59] due to its use in similar work. For these experiments no pre-defined groundtruth is required, instead each scene has the number of agents and their speed estimated. It is expected that simulations that have a similar number of agents and relative speed to the source video will produce the lowest distance measure.

Tables 1 - 4 (left side) contain the average distance measures, after applying equation (19), across all tested scenes for the 12 composite videos using the algorithm outlined in the Pedestrian and Crowd Simulation Algorithm section. As expected the lowest distance values are seen when the simulation parameters closely match those of the source material. Also in Table 4 the average feature combination results are given for the composite videos generated using the Reynolds flocking algorithm [14], in these simulations the same initialisation values and constraints were applied as in the previous results. The cells of the table are coloured green-yellow-red, whereby green is a low distance and therefore a close match to the source footage and red represents a large measured distance and therefore less accurate simulation. Therefore Table 4 demonstrates a direct comparison where it can be seen that the distance measures from the source footage for the Reynolds [14] algorithm are consistently higher than seen in the comparative results from the proposed simulation algorithm. This is logical as modern alternatives to Reynolds [14] flocking algorithm should present more realistic movement. This demonstrates that the suggested framework provides a useful comparison tool from which to analyse simulation similarity.

To further evaluate the methodology, a group of ten people were asked to rate the 12 simulated visualisations created using the simulation algorithm outlined in the Pedes-

Table 2: Average Bhattacharyya distance between source and each simulated video sequence using the simulation algorithm outlined in the Simulation section, across all five Scenes for $\Phi^{H2D}$ features.

| # Agents | Agents Speed | | | |
|---|---|---|---|---|
| | Very Slow | Slow | **Same** | Fast |
| Few | 2.50 | 2.19 | 2.10 | 2.24 |
| **Same** | 2.63 | 2.52 | 2.44 | 2.55 |
| Many | 2.86 | 2.69 | 2.71 | 2.88 |

Table 3: Average Bhattacharyya distance between source and each simulated video sequence using the simulation algorithm outlined in the Simulation section, across all five Scenes for $\Phi^T$ features.

| # Agents | Agents Speed | | | |
|---|---|---|---|---|
| | Very Slow | Slow | **Same** | Fast |
| Few | 4.90 | 3.33 | 2.49 | 3.34 |
| **Same** | 4.59 | 2.91 | 1.45 | 3.48 |
| Many | 6.19 | 4.33 | 4.95 | 3.89 |

trian and Crowd Simulation Algorithm Section, against their respective source material. Focus was given to evaluating the speed, number and track of the agents in each video compared to the source. Using the Mean Opinion Score (MOS) method, the participants were asked to provide a rating of one to five where five represented a high similarity to the source material and one a strong dissimilarity. The values for all the participants were averaged to give a score for each scene. This evaluation technique demonstrates how similarly the proposed features and metrics react compared with the participants Human Visual Systems. The results are contained in Table 5, here the same green-yellow-red colour scheme is employed to allow comparison with Table 4, in these cases it is expected that a similar colour distribution should be seen between the tables.

The left three columns of Table 6 is a breakdown of individual features against the human participant's ability to evaluate video properties. It can be seen that in certain instances the correlation between human and specific feature types is reasonably high. However by using the weighted sum of all three proposed metrics, and again comparing to the MOS, a more robust methodology is seen. This is not surprising as it is often observed that humans have difficulty distinguishing the difference between large amounts of slow moving agents versus a smaller amount of agents moving faster. As a result the combination of distance metrics from all three features more closely matches the Human Visual System's ability to evaluate motion. The weighting of the combination in this case is equal, however the optimal combination will be application dependant. Some metrics will perform better on different types of scenario. For



Figure 9: Side by side tracklet comparison for Road and Kvan. (a-b) Still from source Road video and tracklet. (c-d) Still from simulation and tracklet. (e-f) Still from source Kvan video and tracklet. (g-h) Still from simulation and tracklet.

example videos recorded from a lower point of view may not return descriptive tracklet information. To better match the HVS feature outputs with human perception Weber's Law is incorporated (7), the right three columns of Table 6 demonstrate the improvement seen to the MOS correlations by doing so. In all cases the combination of metrics better correlates to the human perception of movement in the videos.

By using the average distance from all of the three proposed features a robust system is demonstrated. However each of the individual features provides a unique insight into the simulation accuracy. For example, evaluation of the tracklets allows an insight into how accurately the simulation model replicates the movements of the source material. As such in complex scenarios where the source agents change direction a number of times, a strong dissimilarity is expected, likewise in more simplistic scenes where the

Table 4: The average Bhattacharyya distance between source and each simulated video sequence using the simulation algorithm outlined in the Simulation section and, for comparison Reynolds [14], across all five scenes for the feature combination. Results have been colour coded to provide a heat map of similarity, where green represents low distance from source video and red high.

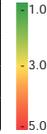| # Agents | Proposed Method | | | | Reynolds [14] | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Agents Speed | | | | | | | |
| | V. Slow | Slow | **Medium** | Fast | V. Slow | Slow | **Medium** | Fast |
| Few | 3.15 | 2.45 | 2.44 | 2.47 | 2.82 | 3.38 | 3.60 | 2.97 |
| **Same** | 3.11 | 2.63 | 2.22 | 2.74 | 3.49 | 3.66 | 3.80 | 3.27 |
| Many | 4.01 | 3.06 | 3.29 | 3.09 | 3.73 | 3.67 | 3.85 | 3.66 |

Table 5: Mean Opinion Score (MOS) of human observations of similarity. Results have been colour coded to provide a heat map of similarity, where green represents high MOS for observed similarity to source video and red low. Provides a visual comparison between the MOS scores and the frameworks distance metrics.

| # Agents | Speed of Agents | | | |
| --- | --- | --- | --- | --- |
| | Very Slow | Slow | Medium | Fast |
| Few | 1.02 | 2.95 | 3.35 | 2.32 |
| Same | 1.92 | 3.27 | 4.36 | 2.87 |
| Many | 1.53 | 3.18 | 3.45 | 2.72 |

Table 6: Correlation (Pearson) between combination features distance and MOS, with and without Weber's Law applied.

| Metric | Without Weber | | | With Weber | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Avg | Agts | Spd | Avg | Agts | Spd |
| $\Phi^{HOOF}$ | 0.67 | 0.49 | 0.70 | 0.54 | 0.31 | 0.62 |
| $\Phi^{T}$ | 0.59 | 0.46 | 0.60 | 0.63 | 0.59 | 0.57 |
| $\Phi^{H2D}$ | 0.24 | 0.06 | 0.41 | 0.28 | 0.02 | 0.44 |
| Comb | 0.55 | 0.36 | 0.60 | 0.61 | 0.44 | 0.65 |

simulation agents closely follow the source tracks a low dissimilarity is expected. A visual example is given in Figure 9, where it can be seen in the first scene (a-d) that there is an obvious visual difference between the source and the simulation, whereas in the second scene (e-g) the similarity is much higher. This is visualised using the tracklet plots which represent a compound image of the tracklets over the duration of the video.

Utilising the HOOF feature per frame and per sequence, an analysis of the amount of movement and magnitude of the optical flow can be made. Visualised examples of these two features are presented in Figures 10 & 11. Figure 10 is the compounded HOOF features for an entire sequence. The Figure 10 (a-b) represents the source material, with (c-d) being the simulation with similar values for number of agents and their speed. Figure 10 (e-f), (g-h) represent low
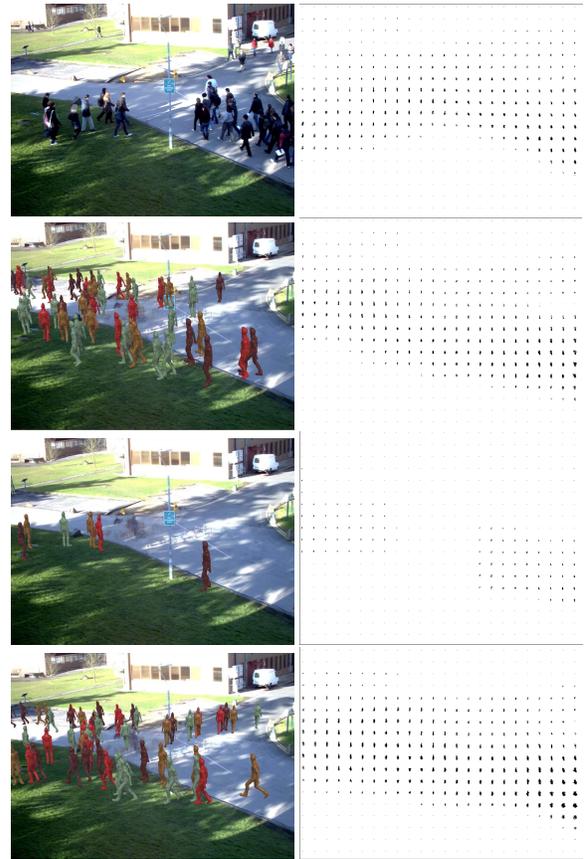


Figure 10: Histogram of Orientated Optical Flow per sequence using Road (Left example still from the video sequence and right, HOOF visualisation). (a-b) Source image, (c-d) medium, (e-f) low and (g-h) high speed and number of agent examples.

and high levels of movement respectively. Figure 11 is the HOOF features using an individual frame. As before (a-b) is the source with (c-d), (e-f) and (g-h) being simulations with the previously mentioned parameters. In both cases its clear to see how the adjustment of speed and number of agents affects the output. Additionally effects on the track-
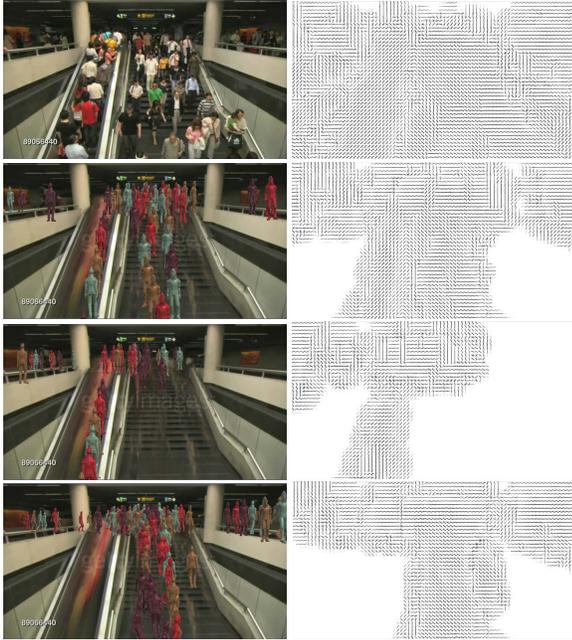
Figure 11: Histogram of Orientated Optical Flow per frame using Krad2 (Left: example still from the video sequence. Right: HOOF visualisation).(a-b) Source image, (c-d) medium, (e-f) low and (g-h) high speed and number of agent examples.

lets can be seen. In the examples where the agent's speed is very low, parts of the scene are left unchanged by agent movement.

## Conclusion

A novel Crowd Simulation Evaluation through Composition (CSEC) framework was presented which reduces the complexity of simulation evaluation and provides tangible and relevant metrics that can be used for comparison and parametric tuning. To extract and produce simulated video, a semi-automated perspective plane extraction process is introduced which allows the conversion of source material into a composited video with controlled agents (3D models) replacing those of the humans. Through the use of a modular system, any crowd or pedestrian simulation model can be evaluated and compared by generating agent motion for use in the final visual simulation. Additionally, any video analysis feature can be utilised to evaluate similarity. Through evaluation on a large range of challenging and diverse scenes, it has been shown that the methodology presents quantifiable measures of video properties such as speed and number of agents. Utilising HVS features to replicate the human's ability to perceive movement, the framework outputs correlate well with human participant analysis of the same videos showing that the system closely

emulates the Human Visual System.

The Crowd Simulation Evaluation through Composition (CSEC) framework introduces a number of key benefits over existing methods. As the framework only takes in a source and simulated video as an input, a number of the time consuming ground truth and annotation steps, such as pedestrian tracks, are reduced. The framework also allows researchers who wish to compare their algorithm against others a quick and efficient way of doing so, either by using the same well-known source material and datasets in the field or simply rerunning the framework with other pedestrian or crowd simulation algorithms to compare with. Additionally for model tuning; the proposed method can create a fast feedback loop that allows the modification of parameters to improve simulation accuracy. As the ground truth data for any simulated visualisations is already intrinsically known, and as specific testing scenarios and behaviours can be simulated, the methodology is also very suitable for the evaluation of pedestrian tracking algorithms on video data.

## Acknowledgements

## References

[1] Miho Asano, Takamasa Iryo, and Masao Kuwahara. A pedestrian model considering anticipatory behaviour for capacity evaluation. *Transportation and Traffic Theory*, 18:28, 2009.

[2] Sujeong Kim, Stephen J. Guy, Dinesh Manocha, and Ming C. Lin. Interactive simulation of dynamic crowd behaviors using general adaptation syndrome theory. *ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games - I3D*, 1(212):55, 2012.

[3] Franziska Klugl, Georg Klubertanz, and Guido Rindsfuser. Agent-based pedestrian simulation of train evacuation integrating environmental data. In *Lecture Notes in Computer Science*, volume 5803, pages 631–638, 2009.

[4] Hui Xi, Seungho Lee, and Young-jun Son. An integrated pedestrian behavior model based on extended decision field theory and social force model. In *Human-in-the-Loop Simulations*, pages 69–95. 2011.

[5] Andrea Portz and Armin Seyfried. Analyzing stop-and-go waves by experiment and modeling. In *Pedestrian and Evacuation Dynamics*, pages 577–586. 2010.

[6] Dorine C. Duives, Winnie Daamen, and Serge P. Hoogendoorn. State-of-the-art crowd motion simulation models. *Transportation Research Part C: Emerging Technologies*, 37:193–209, 2013.

[7] P.a Charalambous, I.b Karamouzas, S.J.b Guy, and Y.a Chrysanthou. A data-driven framework for visual crowd analysis. *Computer Graphics Forum*, 33(7):41–50, 2014.

[8] D. Wolinski, S. Guy, A.H. Olivier, M. Lin, D. Manocha, and J. Pettr. Parameter estimation and comparative evaluation of crowd simulations. *Computer Graphics Forum*, 2(33):303–312, 2014.

[9] Stephen J Guy, Jur van den Berg, Wenxi Liu, Rynson Lau, Ming C Lin, and Dinesh Manocha. A statistical similarity measure for aggregate crowd dynamics. *ACM Transactions on Graphics*, 31(6):1, 2012.

[10] Mubbasir Kapadia, Matt Wang, Shawn Singh, Glenn Reinman, and Petros Faloutsos. Scenario space: Characterizing coverage, quality, and failure of steering algorithms. *Proceedings of the 2011 ACM SIGGRAPH/Eurographics Symposium on Computer Animation - SCA '11*, 1:53, 2011.

[11] Mikel Rodriguez, Josef Sivic, Ivan Laptev, and Jean-Yvesl Audibert. Data-driven crowd analysis in videos. In *International Conference on Computer Vision*, pages 1235–1242, 2011.

[12] Soraia R. Musse, Vinicius J. Cassol, and Cláudio R. Jung. Towards a quantitative approach for comparing crowds. *Computer Animation and Virtual Worlds*, 23(1):49–57, 2012.

[13] K Jablonski, V Argyriou, and D Greenhill. Crowd Simulation for Dynamic Environments based on Information Spreading and Agents' Personal Interests. *Transportation Research Procedia*, 2:412–417, 2014.

[14] Craig W Reynolds. Flocks, herds and schools: A distributed behavioral model. In *ACM SIGGRAPH Computer Graphics*, volume 21, pages 25–34. ACM, 1987.

[15] C. W Reynolds. Steering behaviors for autonomous characters. In *Game Developers Conference*, volume 1999, pages 763–782, 1999.

[16] D Helbing and P Molnar. Self-organization phenomena in pedestrian crowds. *Condensed Matter*, pages 569–577, 1998.

[17] Eleonora Papadimitriou, George Yannis, and John Golias. A critical assessment of pedestrian behaviour models. *Transportation Research Part F: Traffic Psychology and Behaviour*, 12(3):242–255, 2009.

[18] Shawn Singh, Mubbasir Kapadia, Petros Faloutsos, and Glenn Reinman. Steerbench: A benchmark suite for evaluating steering behaviors. *Computer Animation and Virtual Worlds*, 20(5-6):533–548, 2009.

[19] Beibei Zhan, Dorothy N. Monekosso, Paolo Remagnino, Sergio A. Velastin, and Li Qun Xu. Crowd analysis: A survey. *Machine Vision and Applications*, 19(5-6):345–357, 2008.

[20] Julien Pettré, Jan Ondrej, Anne-hélène Olivier, Armel Cretual, and Stéphane Donikian. Experiment-based modeling, simulation and validation of interactions between virtual walkers. *Proceedings of the 2009 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, 2009:189, 2009.

[21] D Crompton. Pedestrian delay, annoyance and risk: pre-liminary results from a 2 years study. In *In Proceedings of PTRC Summer Annual Meeting*, pages 275–299, 1979.

[22] He Wang, Jan Ondej, and Carol O'Sullivan. Path patterns: Analyzing and comparing real and simulated crowds. In *Proceedings of the 20th ACM SIGGRAPH Symposium on Interactive 3D Graphics and Games - I3D '16*, number February, pages 49–57, 2016.

[23] Qiang Wang, Yan Liu, and Juan Chen. Accurate indoor tracking using a mobile phone and non-overlapping camera sensor networks. *2012 IEEE International Instrumentation and Measurement Technology Conference Proceedings*, pages 2022–2027, 2012.

[24] Alon Lerner, Yiorgos Chrysanthou, Ariel Shamir, and Daniel Cohen-Or. Data driven evaluation of crowds. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 5884 LNCS:75–83, 2009.

[25] Alon Lerner, Yiorgos Chrysanthou, and Ariel Shamir. Context-dependent crowd evaluation. *Computer Graphics Forum*, 29(7):2197–2206, 2010.

[26] Bikramjit Banerjee and Landon Kraemer. Evaluation and comparison of multi-agent based crowd simulation systems. In *Agents for games and simulations II*, pages 53–66. Springer, 2011.

[27] Francesco Zanlungo, Dražen Brščić, and Takayuki Kanda. Pedestrian group behaviour analysis under different density conditions. *Transportation Research Procedia*, 2:149–158, 2014.

[28] V. Argyriou and T. Vlachos. Quad-tree motion estimation in the frequency domain using gradient correlation. *IEEE Transactions on Multimedia*, 9(6):1147–1154, Oct 2007.

[29] V. Argyriou. Sub-hexagonal phase correlation for motion estimation. *IEEE Transactions on Image Processing*, 20(1):110–120, Jan 2011.

[30] Victoria Bloom, Vasileios Argyriou, and Dimitrios Makris. G3di: A gaming interaction dataset with a real time detection and evaluation framework. In Lourdes Agapito, Michael M. Bronstein, and Carsten Rother, editors, *Computer Vision - ECCV 2014 Workshops*, pages 698–712, Cham, 2015. Springer International Publishing.

[31] V. Bloom, D. Makris, and V. Argyriou. Clustered spatio-temporal manifolds for online action recognition. In *2014 22nd International Conference on Pattern Recognition*, pages 3963–3968, Aug 2014.

[32] Berthold K P Horn and Brian G. Schunck. Determining optical flow. *Artificial Intelligence*, 17(1-3):185–203, 1981.

[33] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *7th International Joint Conference on Artificial intelligence*, volume 2, pages 674–679, 1981.

[34] Deqing Sun, Stefan Roth, and Michael J. Black. Secrets of optical flow estimation and their principles. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2432–2439, 2010.

[35] Deqing Sun, Stefan Roth, and Michael Black. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *International Journal of Computer Vision*, 106(2):115–137, 2013.

[36] D Makris and T Ellis. Path detection in video surveillance. *Image and Vision Computing*, 20(12):895–903, 2002.

[37] Stefan Munder, Christoph Schnörr, and Dariu M Gavrila. Pedestrian detection and tracking using a mixture of view-based shape texture models. *IEEE Transactions on Intelligent Transportation Systems*, 9(2):333–343, 2008.

[38] Michalis Raptis and Stefano Soatto. Tracklet descriptors for action modeling and video analysis. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 6311 LNCS(PART 1):577–590, 2010.

[39] Pierre Allain, Nicolas Courty, and Thomas Corpetti. Crowd flow characterization with optimal control theory. *Asian Conference on Computer Vision (ACCV)*, pages 279–290, 2009.

[40] W Hu, T Tan, L Wang, and S Maybank. A survey on visual surveillance of object motion and behaviors. *IEEE Transactions on Systems, Man and Cybernetics, Part C*, 34(3):334–352, 2004.

[41] T. I Lakoba, D. J Kaup, and N. M. Finkelstein. Modifications of the Helbing-Molnár- Farkas- Vicsek social force model for pedestrian evolution. *Simulation*, 81(5):339–362, 2005.

[42] T L Clarke, DJ Kaup, Linda Malone, Rex Oleson, and Mario Rosa. Crowd model verification using video data. *Proceedings of EMSS 2007*, pages 4–6, 2007.

[43] Elke U Weber. De Pulsu, Resorptione, Auditu et Tactu. *Annotationes anatomicae et physiologicae*, pages 44–174, 1834.

[44] Zoran Zivkovic. Improved adaptive gaussian mixture model for background subtraction. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 2, pages 28–31, 2004.

[45] Antoni B. Chan, Zhang Sheng John Liang, and Nuno Vasconcelos. Privacy preserving crowd monitoring: Counting people without people models or tracking. *26th IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2008.

[46] Ioannis Karamouzas, Peter Heil, Pascal Van Beek, and Mark H. Overmars. A predictive collision avoidance model for pedestrian simulation. *Lecture Notes in Computer Science*, 5884:41–52, 2009.

[47] Miho Asano, Takamasa Iryo, and Masao Kuwahara. Microscopic pedestrian simulation model combined with a tactical model for route choice behaviour. *Transportation Research Part C: Emerging Technologies*, 18(6):842–855, 2010.

[48] Shuqiang Guo, Zhaoyang Qu, and Liqun Wang. Camera pose estimation using frequency analysis. In *2014 International Conference on Information Science and Applications (ICISA)*, pages 3–6, 2014.

[49] Bernhard Zeisl, Torsten Sattler, and Marc Pollefeys. Camera pose voting for large-scale image-based localization. *Proceedings of the IEEE International Conference on Computer Vision*, pages 2704–2712, 2015.

[50] Yongtae Do. On the neural computation of the scale factor in perspective transformation camera model *.

In *IEEE International Conference on Control and Automation (ICCA)*, pages 712–714, 2013.

[51] W. L. Khoong, W. Y. Kow, H. T Tan, H. P Yoong, Kenneth Teo, and Kin Tze. Kalman filtering based object tracking in surveillance video system. In *Proceedings of the 3rd CUTSE International Conference*, 2011.

[52] Min Hu, Weiming Hu, and Tieniu Tan. Tracking people through occlusions. *Proceedings - International Conference on Pattern Recognition*, 2:724–727, 2004.

[53] E. Wharton, K. Panetta, and S. Agaian. Human visual system based similarity metrics. In *IEEE International Conference on Systems, Man and Cybernetics*, pages 685–690, 2008.

[54] Johannes M Zanker. Does motion perception follow Weber's law? *Perception*, 24(4):363–372, 1995.

[55] Rizwan Chaudhry, Avinash Ravichandran, Gregory Hager, and Ren?? Vidal. Histograms of oriented optical flow and Binet-Cauchy kernels on nonlinear dynamical systems for the recognition of human actions. *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2009*, pages 1932–1939, 2009.

[56] Chen Change Loy, Shaogang Gong, and Tao Xiang. From semi-supervised to transfer counting of crowds. *2013 IEEE International Conference on Computer Vision*, pages 2256–2263, 2013.

[57] J. Ferryman and A. Ellis. PETS2010: Dataset and challenge. *Proceedings - IEEE International Conference on Advanced Video and Signal Based Surveillance, AVSS 2010*, pages 143–148, 2010.

[58] Damien Simonnet, Sergio A. Velastin, James Orwell, and Esin Turkbeyler. Selecting and evaluating data for training a pedestrian detector for crowded conditions. *2011 IEEE International Conference on Signal and Image Processing Applications, ICSIPA 2011*, pages 174–179, 2011.

[59] A Bhattachayya. On a measure of divergence between two statistical population defined by their population distributions. *Bulletin Calcutta Mathematical Society*, 35:99–109, 1943.