

Identification of Lifelike Characteristics of Human Crowds Through a Classification Task

Jamie Webster and Martyn Amos

Department of Computer and Information Sciences, Northumbria University, Newcastle upon Tyne, United Kingdom.
Corresponding author: martyn.amos@northumbria.ac.uk

Abstract

Crowd simulations are used extensively to study the dynamics of human collectives. Such studies are underpinned by specific movement models, which encode rules and assumptions about how people navigate a space and handle interactions with others. These models often give rise to macroscopic simulated crowd behaviours that are statistically valid, but which lack the noisy microscopic behaviours that are the signature of believable “real” crowds. In this paper, we use an existing “Turing test” for crowds to identify “lifelike” features of real crowds that are generally omitted from simulation models. Our previous study using this test established that untrained individuals have difficulty in classifying movies of crowds as “Real” or “Simulated”, and that such people often have an idealised view of how crowds move. In this follow-up study (with new participants) we perform a second trial, which now includes a training phase (showing participants movies of real crowds). We find that classification performance significantly improves after training, confirming the *existence* of features that allow participants to identify real crowds. High-performing individuals are able to *identify* the features of real crowds that should be incorporated into future simulations if they are to be considered “lifelike”.

Introduction

A significant amount of artificial life research is concerned with studying the collective dynamics of *mobile agents* operating in a spatially-explicit environment. Relevant domains include the flocking behaviour of birds and other “animats” (“boids” being the archetypal example (Reynolds, 1987)), the power of distributed swarm robotics (Brambilla et al., 2013), and the engineering of biological cell populations (Gorochoowski, 2016). In all such cases, agents (whether simulated or physically realised) are situated in Cartesian space, and may interact both with one another and with their environment.

One specific area of growing interest is the study of *crowd dynamics* (Adrian et al., 2019); that is, the behaviour of large numbers of human individuals moving through and interacting in a given environment. The need to understand collective human behaviour in physical space is pressing, as

it has significant implications for events planning and management (Crociani et al., 2016), urban design (Feng et al., 2016), and incident response and analysis (Harding et al., 2011; Pretorius et al., 2015). During and after the COVID pandemic, with potentially long-lasting and profound structural and behavioural changes being made, the need to understand the crowd will persist (Pouw et al., 2020).

Due to the inherent difficulty of performing large-scale experiments with human participants, *crowd simulations* (Thalmann and Musse, 2013) (usually using an agent-based approach) are often used to investigate collective behaviour and the impact of physical or behavioural interventions on crowd dynamics. Two features of simulations are of interest; *validity* and *believability*. Validity describes how closely the output of the model matches data obtained from the real world (Klüpfel, 2007; Pettré et al., 2009; Seer et al., 2014). *Believability* is subtly different, and concerns the human perception of whether or not a crowd’s behaviour is *life-like*, or plausible. We are not concerned with “cinematic”, photo-realistic believability of the rendering of a crowd, but whether or not observers are able to detect characteristic *patterns of behaviour* in real crowds which are absent in simulated crowds. Fundamentally, we assume that a simulation is valid, and are interested in whether or not it also *looks realistic*.

The rest of the paper is organised as follows; we give some background motivation, outline our hypothesis, and describe our crowd Turing test framework for its investigation. We then describe our experimental method for the current study, and describe our results. We conclude with a discussion of the implications of our findings, and suggest possible future work.

Background and Motivation

Crowd simulations are now used extensively in a wide range of application domains, from urban planning (Aschwenden et al., 2011), emergency response (Mahmood et al., 2017), games and training simulations (Mckenzie et al., 2008), and the CGI generation of Hollywood movie scenes (a classic example being the large-scale battle scenes in *The Lord*

of the *Rings* series) (Ricks, 2013). Most crowd simulations are underpinned by a behavioural/movement model, which makes simplifying assumptions about individuals, and which is used by agents to determine their trajectories through the simulated space.

The Social Forces Model (SFM) (Helbing and Molnar, 1995) lies at the heart of many scientific and commercial crowd simulation packages, such as FDS+EVAC (Korhonen et al., 2010), PedSim (Gloor, 2016), SimWalk (Kimura et al., 2003) and MassMotion (Oasys, 2019). This is a microscopic, continuous model which uses “attractive” and “repulsive” force fields between individuals (and between individuals and their environment) to guide movement. However, there are well-established deficiencies in this and other existing movement models. As (Lerner et al., 2007) argue, “While such approaches may capture the broad overall behaviour of the crowd, they often miss the subtle details displayed by the individuals. The range of individual behaviours that may be observed in a real crowd is typically too complex for a simple behavioural model... Simple things such as walking in pairs, stopping to talk to someone, changing one’s mind and heading off in a different direction or aimlessly wandering about, are just a few examples which are difficult to capture.” The emphasis here is less on the locomotion model of avatars or the cosmetic appearance of the agents, and more on the *patterns* and “quirks” of movement that distinguish a real crowd from a simulated one.

Why is this important? After all, emergency planners (to take one significant user group) will generally be satisfied if the overall outcome of a simulation (in terms of the time required to evacuate a stadium, for example) is broadly valid, and will usually not concern themselves with micro-level “turbulence” and other localised phenomena. However, as (Fuchsberger et al., 2017) argue, crowd simulations still meet with resistance from decision makers in some significant industrial and societal domains, and this may be due to a lack of trust in their outputs (caused, in turn, by a lack of “realism”). Specific concerns identified of relevance to the current paper include “unnatural motion paths”, so if we can go some way towards addressing this, then it may lead to increased acceptance and uptake of these techniques.

As we argue in (Webster and Amos, 2020), there is still a need for more realistic behavioural/movement models in crowd simulation, and “This is motivated by a widely-acknowledged need for crowd simulations to include more “lifelike” features derived from individual and social psychology (such as group-level behaviours, indecision, etc.) (Lemerrier and Auberlet, 2016; Seitz et al., 2017; Templeton et al., 2015), which are generally not included in software packages, and which give rise to rather unrealistic or “robotic” patterns of behaviour at the population level”.

Much work has already been done on making crowd simulations more realistic; here we highlight some representative contributions. (Lerner et al., 2007) describes the con-

struction of a database of behavioural “motifs” which may be incorporated into an agent’s behaviour. (Peters and Ennis, 2009) used manual annotation of observations to extract information about group-level behaviours that were then incorporated into simulations (this study also included human trials of perception of realism). More recently, (Wei et al., 2018; Yao et al., 2020) used machine learning to extract features of observed crowds, which were then incorporated into a crowd simulation, but neither study assessed whether or not these modifications actually made the overall crowd behaviour more realistic.

Fundamentally, what passes for “lifelike” is inherently subjective. To our knowledge, until we performed this study no extensive work had been done on capturing the “essence” of what makes a crowd lifelike *from the perspective of human observers*.

Our previous work (Webster and Amos, 2020) showed that crowd simulations that employ the most commonly-used movement model are valid (in terms of their outputs having the same statistical properties as observed crowds), but they still possess a “signature” that allows them to be distinguished from real crowds. Simply put, to human observers, simulated crowds are still perceived differently to real crowds. Importantly, though, we also found that although people are able to reliably *partition* crowds into “Real”/“Simulated”, *they are unable to tell which is which*. That is, individuals are able to separate crowd movies into two categories, but they are unable to reliably label the real crowds. We found that individuals tend to have an idealised view of the behaviour of real crowds, which is often at odds with reality. These findings confirm the observation that real and simulated crowds have different microscopic features that allow them to be partitioned, if not classified.

To summarise, our previous work established the *existence* of features that are present in real crowds but not in simulated crowds; the aim of the current paper is to *identify* those features. In (Webster and Amos, 2020) we argue that “Our results suggest a possible framework for establishing a minimal set of collective behaviours that should be integrated into the next generation of crowd simulation models.” Here, we use the “Turing test” classification task to identify that specific set of features that allow trained viewers to reliably *classify* (not just partition) “Real” and “Simulated” crowds. Our results show that classification performance over a population of observers increases significantly after an initial training phase, and that individuals are able to identify a core set of “lifelike” behaviours that are present in real crowds, but which are absent in simulated crowds. This immediately suggests new features that must be incorporated into future crowd simulations if they are to be considered “lifelike”.

Hypothesis

In a landmark paper (Turing, 1950), Alan Turing proposed a method to investigate what would become known as “artificial intelligence”. Rather than directly answering the somewhat ambiguous question “Can machines think?”, Turing preferred to reframe the issue in terms of an “imitation game”, in which an interrogator engaged in conversation with two agents via “teletypes”. One of the agents (A) is a man, and the other (B) a woman, and the interrogator’s objective is to decide which is which by asking questions of both and assessing their responses. The task of A is to cause the interrogator to guess *incorrectly* (that is, persuade them that he is a woman), and the task of B is to “help” the interrogator to guess correctly, generally by giving truthful answers. We may, therefore, interpret the imitation game (commonly referred to as the “Turing test”) more generally, with the role of A being played by an artificial system that seeks to persuade a human observer that it is the “genuine article”, and B being played by an actual “real world” example of the system under study. Importantly, the test does not seek to establish the “truth” of A’s outputs (that is, their validity), but simply whether or not A could be said to represent a reasonable facsimile of the system represented by B.

This conceptual framework has been proposed for biological modelling (Harel, 2005) and artificial life (Cronin et al., 2006) as a way of investigating the lifelike properties of artificial systems. We previously used the same approach to investigate crowd simulations, basing our approach on a related Turing test for collective motion in fish (Herbert-Read et al., 2015). In (Webster and Amos, 2020), we describe the results of initial experiments, using a total of 540 in-person participants. The first set of trials presented individuals with a sequence of paired movies, using a side-by-side representation. In each pair, one of the movies represented the movement of a real crowd, and the other represented a computer simulation of the same scenario (the ordering was randomised). All observations were of the same physical space, and both movies were generated using the same custom rendering engine. For each pair (over six pairs in total), participants were asked to specify which of the pair they thought was the real crowd (that is, they had to *identify* the real crowd). For the second set of trials, participants were presented with the movies individually, and this time they were asked to *classify* each movie as either “Real” or “Simulated”.

We found that participants performed better when they were asked to *classify* crowds rather than having to choose between the two, but a striking feature of our results was that neither mode allowed participants to perform better than random guessing. A simplistic interpretation of this result could be that existing simulations are good enough to “pass” the crowd Turing test, as human observers are unable to distinguish between them, but here we emphasise that the imita-

tion game, as originally described by Turing, requires the interrogator to be able to specify *which* agent is the man.

Strikingly, the most common score in the first trial was zero, meaning that a significant proportion of participants (36.46%) failed to identify a single real crowd. That is, their entire perception of what constitutes a real crowd was perfectly “flipped” compared to reality. This sizeable group of participants were able to perfectly partition movies into “Real” or “Simulated”, but were utterly unable to say which was which. This confirmed the existence of a set of real crowd behaviours (informally described by participants in terms of “standing around” and “moving with purpose”) that allowed individuals to separate real from simulated, but which were incorrectly ascribed to the simulation as generating “unrealistic” crowd behaviour. Our conclusion was that participants had an idealised view of real crowd behaviour, and preferred to think that it was much less “messy” and unpredictable than observations would suggest.

Our hypothesis, therefore, is that participants in a crowd Turing test will improve their classification performance after being trained by viewing real crowds, as a result of being able to identify and ascribe *only to real crowds* the lifelike features that are manifested in the training set.

Experimental Methods

Our protocol was largely modelled on that of (Webster and Amos, 2020), but limitations imposed by the COVID pandemic required us to perform our trials online, as opposed to in-person. We do not believe that this modification had any significant impact on our results; indeed, it actually allowed us to recruit a more diverse range of participants, rather than using only University students (which was a possible criticism of the original study).

We performed two sets of Turing test experiments; the first (Test 1) was an online-only repetition of the second (classification) test from (Webster and Amos, 2020), with entirely new participants. We attracted 232 participants, who were recruited via social media. This first test allowed us to assess the ability of each untrained participant to classify crowds as either “Real” or “Simulated”, thus assigning each one a baseline score. We allowed an appropriate period of time to pass (4 months) in order to ensure that the tests were independent (that is, any learning effects from the first test would not be carried over to the second). We then contacted every Test 1 participant who supplied an email address to invite them to participate in the follow-up Test 2 (they were each offered a £10 gift card as an incentive); 50 participants accepted our invitation. Test 2 participants were then “trained” by asking them to first watch six rendered movies of crowds that were explicitly described as real. Participants then performed a second version of the classification task (as in Test 1), using a different set of real and simulated clips to those used previously (in order to avoid effects induced by familiarity with the clips).



Figure 1: Single movie frame of the Edinburgh Informatics Forum, taken from (Majecka, 2009).

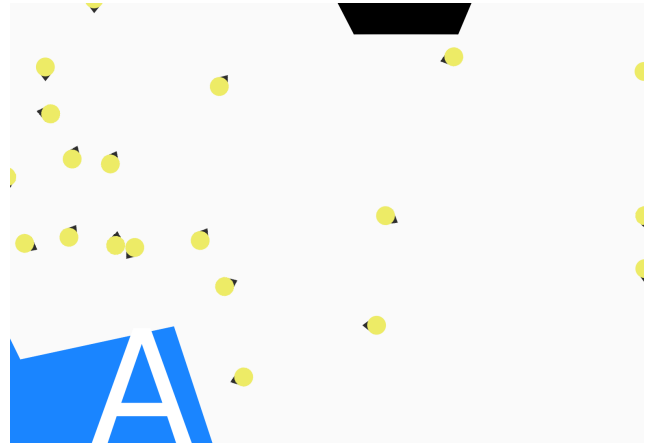


Figure 3: Example rendering of a crowd scene.

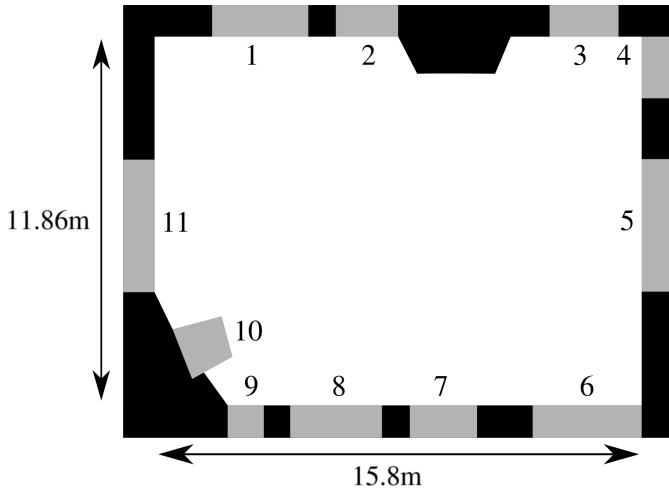


Figure 2: Diagram of Edinburgh Informatics Forum (ingress and egress points numbered).

Given that each participant had a known baseline score from Test 1, we were able to establish whether or not the training phase had a significant effect on classification ability. Participants were specifically asked to identify features that they thought allowed them to distinguish between real and simulated crowds.

Test 1 was performed at the end of June-start of July 2020, and Test 2 was performed in December 2020. Our trial protocol was approved by the Northumbria University Faculty of Engineering and Environment Ethics Committee, application number 24623. We now describe each component of the trial in more detail.

Real pedestrian motion dataset

As we employed the same dataset used in our previous study, we take our description of it from (Webster and Amos,

2020). We used data on real pedestrians from the University of Edinburgh School of Informatics (Majecka, 2009). This public dataset, captured in 2010, contains over 299,000 individual trajectories corresponding to the movement of individuals through the School Forum, and is one of the largest open datasets of its type. A photo of the Forum space is shown in Figure 1, and a diagram is shown in Figure 2. The Forum is rectangular in shape (measuring approximately 15.8×11.86 metres), has eleven ingress/egress points, and is generally clear of obstructions. Images were captured (9 per second) by a camera suspended 23m above the Forum floor, from which individual trajectories were extracted and made available (extraction was performed by the author of (Majecka, 2009)). We note that only the *trajectories* have been made publically available, and not the original video recordings, for ethical and practical reasons (these files require several terabytes of storage). This dataset has been used in several studies of pedestrian movement/tracking, including (Fernando et al., 2018; Lovreglio et al., 2017; Wang and O’Sullivan, 2016). Importantly, none of the individuals whose trajectories were captured were actively participating in movement studies; the trajectories, therefore, are as close to “natural” as possible (i.e., they have “behavioural ecological validity” (Lovreglio et al., 2017)).

In what follows, we use the term “clip” to specifically refer to a time-limited sequence of trajectory data (whether taken from the Edinburgh dataset or from the output of a simulation), as opposed to a movie visualisation. We wrote a utility to search the Edinburgh dataset and extract clips of a specific duration containing a specific number of individuals. This allowed us to ensure that the “real” and “simulated” crowds contained the same number of individuals for any single comparison.

Simulation construction and validation

Each test required participants to classify a number of clips of pedestrian movement as either “Real” or “Simulated”. We began by selecting, at random, a number of clips (30s duration) from the Edinburgh dataset, and extracting information about the number of individuals visible and the entry/exit point distribution. This information was then used to “seed” a simulation. In this way, we obtained both “Real” and “Simulated” versions of the same scenario; the real version was a rendered version of the actual observations, and the simulated version was a rendered version of the output of the model. Uniform rendering was performed by our own custom Java program, which produced “top down” visualisations of both real and simulated clips that were identical in appearance, with individuals represented as filled circles, and headings depicted by an arrow (see Figure 3).

In order to model the scenarios captured in each real Edinburgh clip, we simulated pedestrian movement using the Vadere package (Kleinmeier et al., 2019). This is an open-source package, which means that (unlike commercial software) its movement models are open to inspection. Importantly, it also allows for easy exporting of simulated pedestrian trajectories, which is necessary for rendering.

In Test 1 we used only the SFM movement model; in Test 2, we divided the simulations between the SFM and an alternative movement model, the Gradient Navigation Model (GNM) (Dietrich and Köster, 2014) in order to test whether different movement models have unique movement “signatures”. We used the default Vadere parameter values for each model.

It is important to ensure that simulations (regardless of the movement model) produce outputs that are valid, so we first calculated several statistical properties for a set of simulations and the Edinburgh observations on which they were based. As in (Webster and Amos, 2020), we used two metrics (Herbert-Read et al., 2015); *polarization* and *nearest neighbour distance* (NND). The first metric is particularly useful for describing the existence of large groups who might be moving together along the same heading (e.g., from a lecture towards an exit), while the second metric is used for estimating overall crowd density (detailed descriptions of each are supplied in (Webster and Amos, 2020)).

We selected 20 random Edinburgh clips with varying crowd sizes, and then simulated each scenario 20 times with each movement model. Results are presented in Figure 4; these confirm that both movement models produce high-level outputs that are comparable to the real-world scenarios, and that there are no significant differences between the outputs of each movement model.

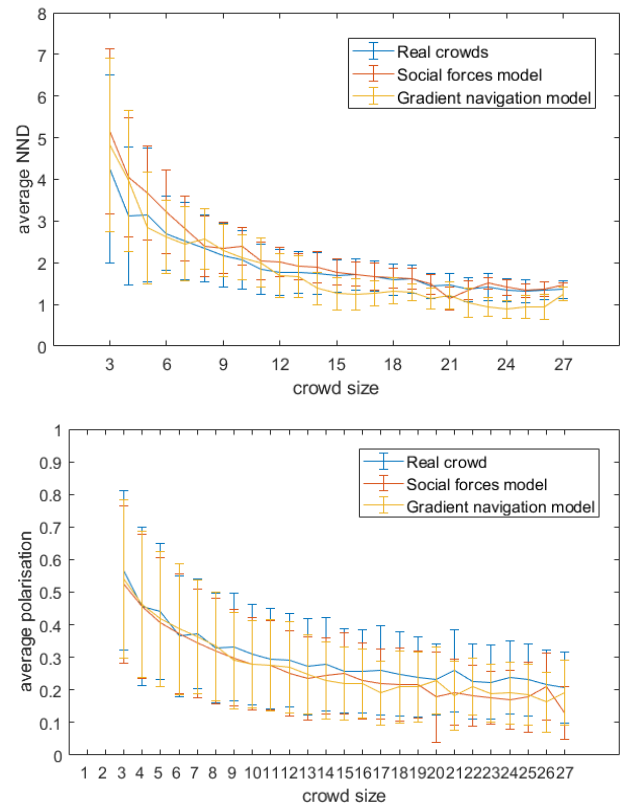


Figure 4: Movement models/real crowd statistical comparisons: Nearest Neighbour Distance (NND) (top) and polarisation (bottom) as a function of crowd size. The outputs of both movement models have properties that are close to those of the real crowds.

Classification tests

For both tests, we constructed a web-based application¹ which presented users with an information screen, asked them to click to confirm their consent to participate, and then presented participants with a randomised sequence of movies. For each movie, participants were asked to click either a “Real” or “Simulated” button, according to their own perception and opinion. At the end of the sequence, users were asked in a free text box to supply short notes on any features that they thought allowed them to identify the real crowd, to specify their level of expertise in crowd science (“High”, “Medium” or “Low”), and to supply their email address (this was used as a participant ID to allow for tracking across the two tests). Once the user submitted their information, their responses were stored on the server, and they were told how many real crowds they had correctly identified (this may have inadvertently helped with recruitment, as some particularly high-scoring participants shared screenshots of their success on social media...)

¹Available at <http://www.martynamos.com/TF2/>

Set	Test 1	s.d.
$P_1 - P_2$	31.21%	20.19
P_2	27%	19.31

Table 1: Test 1 average scores for $P_1 - P_2$ and P_2 . Scores are presented as “% correctly classified”, as the number of movies differed between tests. Analysis confirms that P_2 is representative.

For Test 1, we showed participants a sequence of 12 movies, 6 of which were based on real trajectories, and 6 of which were generated using the SFM-based simulation. Each movie was 30s in duration (in all cases, participants were free to choose “early”, before the end of the movie, and move on to the next one).

For Test 2, we first required participants to undertake a training phase, in which they were shown 6 representative clips generated from Edinburgh observations. Participants were made explicitly aware that they were watching “real” crowds. They were then shown 18 movies in total; 6 based on observations, 6 derived from SFM-based simulations, and 6 from GNM-based simulations.

Results

In this Section we present our trial results. In what follows, we adopt the following notation for participant groups; P_1 is the initial set of 232 participants who took the first Turing test, and P_2 is the subset of 50 participants in P_1 who went on to take the second test.

Classification accuracy

We first consider whether or not group P_2 is representative of the larger set of participants. In both Test 1 and Test 2, participants were scored according to their ability to correctly classify movies, and received 1 point for every correct classification. We calculate the average Test 1 scores for both $P_1 - P_2$ (that is, participants who only took Test 1) and P_2 (participants who took both Tests), and present them in Table 1 (scores are presented as % due to the fact that the number of movies differed between tests).

A Lilliefors test confirms that neither dataset is normally distributed, so we use a two-sided Wilcoxon rank sum test to confirm that data in $P_1 - P_2$ and P_2 are samples from continuous distributions with equal medians ($p = 0.0724$). We conclude, therefore, that P_2 is a representative group.

We then calculate the average Test 1 and Test 2 classification scores for P_2 only; these are shown in Table 2. This reveals a *significant* improvement in overall correct classification score after training (from 27% to 60%). In Trial 2, participants correctly identified SFM-derived movies 63% of the time, and GNM-derived movies 59% of the time, suggesting that there is no significant difference between the two models in terms of the overall characteristics of their outputs.

Test 1	s.d.	Test 2	s.d.
27%	19.31	60.22%	26.35

Table 2: Test 1 and Test 2 average scores for P_2 only.

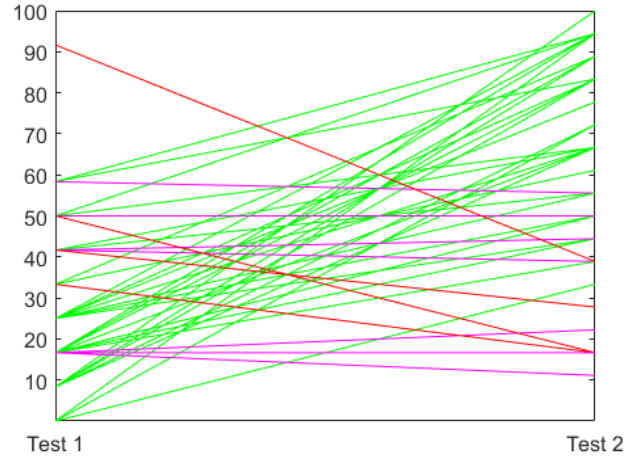


Figure 5: Slopegraph plot of changes in individual classification performance between Test 1 and Test 2 (50 individuals shown in total). Green lines show significant improvements, purple lines show small changes, and red lines show significant reductions in performance.

In Figure 5 we depict the individual changes in performance for the 50 members of P_2 ; visual inspection confirms that the vast majority of participants showed a marked improvement in classification performance after training.

These results confirm the first part of our hypothesis; that suitably trained individuals are able to improve their classification performance after viewing movies of real crowds.

Narrative findings

We now move on to consider the free text supplied by members of P_2 , and extract common themes that enable us to identify specific features of real crowds that allow them to be identified as such. We performed an informal version of this analysis in (Webster and Amos, 2020), but extracted only a small number of general themes, and did not correlate them with classification performance (as we do here).

All 50 participants supplied feedback, so this provides useful additional context to explain the general uplift in performance. Given the relatively small amount of text, we performed manual thematic analysis to extract the predominant features highlighted in the supplied corpus. Each line of free text was broken down into thematic “atoms”, which were then semantically mapped onto over-arching themes. These are summarised in Table 3, partitioned into those features ascribed to real crowds, and those to simulated crowds. We also give the relative frequency of each feature/theme (a link

Real	Freq. %	Simulated	Freq. %
Heterogenous/diverse paths/speeds (R1)	9.21	Homogeneous behaviour (S1)	5.26
Chaotic/unpredictable/erratic movement - rapid changes (R2)	21.05	Rapid direction/speed changes (S2)	3.95
Decisiveness/purposefulness - direct movement (R3)	6.56	Goal-driven (S3)	3.95
Stop-start movement (R4)	7.89	Smooth/continuous movement (S4)	15.79
Static individuals/groups (R5)	2.63	Clusters (S5)	1.32
Groups/flocking/close proximity/collisions (R6)	7.89	Long interactions/collisions and close proximity (S6)	6.58
Collision avoidance (R7)	5.26	Collision avoidance (S7)	2.63

Table 3: Themes identified in narrative comments (labels given in brackets), and their observed frequencies. Related themes are presented alongside one another, although there may not always be an *exact* correlation.

to the full dataset is supplied at the end of the paper). We label each feature for ease of presentation/discussion.

We immediately notice two dominant features; R2 (*real* crowds exhibit chaotic or unpredictable movement, sometimes with rapid changes in speed/direction) accounted for 21% of thematic atoms, and S4 (*simulated* crowds show smooth/continuous movement) accounted for nearly 16% of all atoms. These observations are clearly complementary, in that (after training) observers believe that real crowds are more unpredictable than simulated crowds, which move more smoothly.

However, it is not sufficient to simply analyse the *frequency* of themes, since dominant features may not necessarily correlate with good classification performance in the participants who identify them. We also need to extract the features that have been identified *by the participants who perform best* (or who show the best relative improvement) in the classification task. We first consider *relative* changes in scores, and then look at the *absolute* changes, as each perspective yields insights.

In Figure 6 we plot each theme against both their frequency of mentions and the average relative change in classification performance of participants who specifically mention that theme. All scores are expressed in terms of the *percentage* of movies that were correctly classified, not the “raw” score (as previously stated, the number of movies differs between tests). For participant, p , in test i , relative change is calculated as $((score_{p,2} - score_{p,1}/score_{p,1}) * 100)$, where $score_{p,1} > 0$. For example, a participant who scored 3/12 (25%) in Test 1 and 15/18 (83%) in Test 2 would have their relative change calculated as $((83 - 25)/25) * 100 = 232\%$. When calculating the average relative change, we discard 4 participants with a Test 1 score of zero, as the notion of relative change is not defined for a zero reference value (however, these participants are still included in the discussion of actual score differences, below).

We notice, from inspection, a cluster of themes that are relatively infrequently mentioned ($< 10\%$), but which are

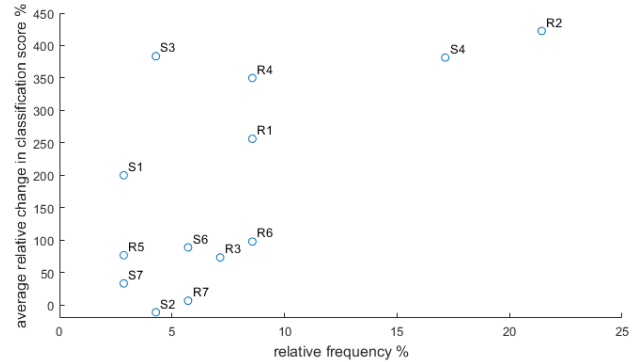


Figure 6: Thematic frequency versus average *relative* change in classification performance. The upper-right quadrant shows two themes (S4 and R2) which both appear frequently and which are correlated with significant positive relative change in classification performance in those participants who mention those themes.

associated with significant improvements in classification performance. However, we see that the two themes that are mentioned with frequency $> 15\%$ - S4 (smooth/continuous movement in simulated crowds) and R2 (unpredictable movement in real crowds) - are both also associated with performance improvements of around 400%. As noted earlier, these themes are complementary.

This finding is entirely consistent with our earlier informal narrative results (Webster and Amos, 2020), where participants who had “flipped” the real and simulated crowds believed that erratic movement was characteristic of “fake” (simulated) crowds, and that real crowds moved smoothly and predictably. After training on real crowds, however, the participants in this second trial correctly identified that real crowds are actually more noisy and unpredictable, and that overwhelmingly smooth, predictable trajectories are a characteristic of simulations.

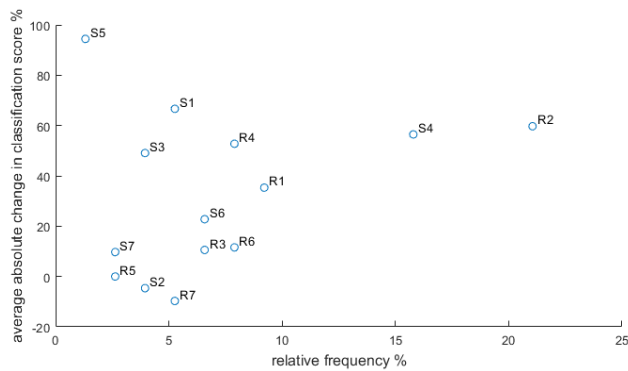


Figure 7: Thematic frequency versus average *absolute* change in classification performance. S2 and R7 are low-frequency themes that are nonetheless associated with reductions in classification performance.

We now consider *absolute* changes in classification score between tests. We see roughly the same clustering of labels as before (S5: presence of clusters in simulated crowds is an outlier, in that it was mentioned only by a single person, albeit one who saw a significant improvement in their classification score). Here we draw particular attention to the (albeit infrequently mentioned) themes that are correlated with *negative* shifts in performance. That is, the features that are mentioned by participants whose classification performance got worse after training. The two features to which this applies are S2 (rapid direction/speed changes in simulated crowds) and R7 (collision avoidance in real crowds).

Again, these findings are entirely consistent with both the current results and our previous study. If high-performing participants correctly spot that simulated crowds move smoothly, then it is entirely to be expected that low-performing participants will (incorrectly) ascribe S2 to them. Collision avoidance in real crowds (R7) is also specifically mentioned in our previous study; participants who performed badly assumed that individuals in real crowds would naturally avoid one another. As we observe in (Webster and Amos, 2020), “In reality, the opposite is true, as the real dataset contains multiple instances of individuals coming into close proximity. Moreover, the social forces model explicitly tries to keep individuals apart unless close proximity is unavoidable, so the behaviour (distance keeping) that participants attributed to real people was actually an in-built feature of the simulation.”

We conclude, therefore, that the primary feature of real crowds that allows trained individuals to correctly distinguish them from simulated crowds is their higher degree of unpredictability in terms of individual trajectories. A secondary feature is collision avoidance (specifically, proximity). Based on this work, our main suggestion (if what we seek is “lifelike” believability in crowd simulations) is

that models should include the facility to add “noise” to the movement of individual agents (surprisingly, this feature is not generally provided). Models might also benefit from a relaxation of collision detection radii to allow for closer proximity of agents. In this way, we might easily replicate the appearance of at least some of the micro-level behaviours referenced by (Lerner et al., 2007).

Discussion and Conclusions

In this paper we report the results of a human trial to identify the “signature” characteristics of real crowds that allow them to be distinguished from simulated crowds. We find that unpredictability in terms of individual trajectories is by far the best discriminator, and proximity in collision detection is also relevant. We note some limitations of our study; the underlying crowd dataset is based on a relatively small physical space which is quite regular in nature, but we point out that it is actually much larger than the arenas used for artificial crowd experiments. Moreover, the observations have a higher level of ecological validity, as the recorded pedestrians were not consciously aware of being participants in an experiment. Our second test used a relatively small number of participants, but we have established that they were representative of a larger set. Finally, our findings are only applicable to “routine” crowds (that is, where people are going about their everyday business), and not to “emergency” or “evacuation” crowds, where behaviours will be very different. However, there is still significant value in updating simulation of such routine crowds to render them more “lifelike”, especially if important policy or design decisions are to be made based on how they are perceived. This study has provided empirical evidence to support the inclusion of relatively straightforward modifications to any and all of the movement models underpinning both scientific and commercial crowd simulation packages. Importantly, the addition of noise to individual trajectories and the relaxation of collision detection radii are entirely generic updates, but ones that could significantly improve the believability of crowd simulations across a range of applications.

Future work may include the automatic detection of features of real crowds from larger and more complex datasets, consideration of the impact of changing movement model parameters, and the integration of identified features into commercial crowd simulation packages in order to test their impact on believability (thus “closing the circle”).

Materials

All code (simulations and analysis scripts) and datasets generated are available at <http://doi.org/10.6084/m9.figshare.c.5280902>

Acknowledgements

JW was supported by a Ph.D. studentship from the Faculty of Engineering and Environment, Northumbria University. We thank Gerta Köster and her research team for useful discussions, and all of the trial participants for their contributions.

References

- Adrian, J., Amos, M., Baratchi, M., et al. (2019). A glossary for research on human crowd dynamics. *Collective Dynamics*, 4(A19):1–13.
- Aschwanden, G., Haegler, S., Bosché, F., Van Gool, L., and Schmitt, G. (2011). Empiric design evaluation in urban planning. *Automation in Construction*, 20(3):299–310.
- Brambilla, M., Ferrante, E., Birattari, M., and Dorigo, M. (2013). Swarm robotics: a review from the swarm engineering perspective. *Swarm Intelligence*, 7(1):1–41.
- Crociani, L., Lämmel, G., and Vizzari, G. (2016). Multi-scale simulation for crowd management: a case study in an urban scenario. In *International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 147–162. Springer.
- Cronin, L., Krasnogor, N., Davis, B., Alexander, C., Robertson, N., Steinke, J., Schroeder, S., Khlobystov, A., Cooper, G., Gardner, P., et al. (2006). The imitation game - a computational chemical approach to recognizing life. *Nature Biotechnology*, 24(10):1203.
- Dietrich, F. and Köster, G. (2014). Gradient navigation model for pedestrian dynamics. *Physical Review E*, 89(6):062801.
- Feng, T., Yu, L.-F., Yeung, S.-K., Yin, K., and Zhou, K. (2016). Crowd-driven mid-scale layout design. *ACM Transactions on Graphics*, 35(4):132–1.
- Fernando, T., Denman, S., Sridharan, S., and Fookes, C. (2018). Tracking by prediction: A deep generative model for multi-person localisation and tracking. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1122–1132. IEEE.
- Fuchsberger, A., Tahmasbi, N., and Ricks, B. (2017). A framework for achieving realism in agent-based pedestrian crowd simulations. In *Semantics, Ontologies, Intelligence and Intelligent Systems (SIGODIS)*.
- Gloor, C. (2016). PedSim: Pedestrian crowd simulation. <http://pedsim.silmaril.org>.
- Gorochowski, T. E. (2016). Agent-based modelling in synthetic biology. *Essays in Biochemistry*, 60(4):325–336.
- Harding, P., Gwynne, S., and Amos, M. (2011). Mutual information for the detection of crush. *PLoS ONE*, 6(12):e28747.
- Harel, D. (2005). A Turing-like test for biological modeling. *Nature Biotechnology*, 23(4):495.
- Helbing, D. and Molnar, P. (1995). Social force model for pedestrian dynamics. *Physical Review E*, 51(5):4282.
- Herbert-Read, J., Romensky, M., and Sumpter, D. (2015). A Turing test for collective motion. *Biology Letters*, 11:20150674.
- Kimura, T., Sekine, H., Sano, T., Takeichi, N., Yoshida, Y., and Watanabe, H. (2003). Pedestrian simulation system SimWalk. In *Summaries of Technical Papers of Annual Meeting Architectural Institute of Japan, E-1*, pages 915–916.
- Kleinmeier, B., Zönnchen, B., Gödel, M., and Köster, G. (2019). Vadere: An open-source simulation framework to promote interdisciplinary understanding. *arXiv preprint arXiv:1907.09520*.
- Klüpfel, H. (2007). The simulation of crowd dynamics at very large events - calibration, empirical data, and validation. In *Pedestrian and Evacuation Dynamics (PED) 2005*, pages 285–296. Springer.
- Korhonen, T., Hostikka, S., Heliövaara, S., and Ehtamo, H. (2010). FDS+Evac: an agent based fire evacuation model. In *Pedestrian and Evacuation Dynamics (PED) 2008*, pages 109–120. Springer.
- Lemercier, S. and Auberlet, J. (2016). Towards more behaviors in crowd simulation. *Computer Animation And Virtual Worlds*.
- Lerner, A., Chrysanthou, Y., and Lischinski, D. (2007). Crowds by example. In *Computer Graphics Forum*, volume 26, pages 655–664. Wiley Online Library.
- Lovreglio, R., Dias, C., Song, X., and Ballerini, L. (2017). Towards microscopic calibration of pedestrian simulation models using open trajectory datasets: the case study of the Edinburgh Informatics Forum. In *Conference on Traffic and Granular Flow, Washington DC, USA*.
- Mahmood, I., Haris, M., and Sarjoughian, H. (2017). Analyzing emergency evacuation strategies for mass gatherings using crowd simulation and analysis framework: Hajj scenario. In *Proceedings of the 2017 ACM SIGSIM Conference on Principles of Advanced Discrete Simulation*, pages 231–240.
- Majecka, B. (2009). Statistical models of pedestrian behaviour in the forum. Master's thesis, School of Informatics, University of Edinburgh.
- Mckenzie, F. D., Petty, M. D., Kruszewski, P. A., Gaskins, R. C., Nguyen, Q.-A. H., Seevinck, J., and Weisel, E. W. (2008). Integrating crowd-behavior modeling into military simulation using game technology. *Simulation & Gaming*, 39(1):10–38.
- Oasys (2019). Mass motion product page. <https://www.oasys-software.com/products/pedestrian-simulation/massmotion/>.
- Peters, C. and Ennis, C. (2009). Modeling groups of plausible virtual pedestrians. *IEEE Computer Graphics and Applications*, 29(4):54–63.
- Petré, J., Ondřej, J., Olivier, A.-H., Cretual, A., and Donikian, S. (2009). Experiment-based modeling, simulation and validation of interactions between virtual walkers. In *Proceedings of the 2009 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 189–198. ACM.
- Pouw, C. A., Toschi, F., van Schadewijk, F., and Corbetta, A. (2020). Monitoring physical distancing for crowd management: Real-time trajectory and group analysis. *PLoS ONE*, 15(10):e0240963.
- Pretorius, M., Gwynne, S., and Galea, E. (2015). Large crowd modelling: an analysis of the Duisburg Love Parade disaster. *Fire and Materials*, 39(4):301–322.

- Reynolds, C. W. (1987). Flocks, herds and schools: A distributed behavioral model. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques*, pages 25–34.
- Ricks, B. C. (2013). *Improving Crowd Simulation with Optimal Acceleration Angles, Movement on 3D Surfaces, and Social Dynamics*. PhD thesis, Brigham Young University.
- Seer, S., Rudloff, C., Matyus, T., and Brändle, N. (2014). Validating social force based models with comprehensive real world motion data. *Transportation Research Procedia*, 2:724–732.
- Seitz, M., Templeton, A., Drury, J., Köster, G., and Philippides, A. (2017). Parsimony versus reductionism: How can crowd psychology be introduced into computer simulation? *Review of General Psychology*, 21(1):95–102.
- Templeton, A., Drury, J., and Philippides, A. (2015). From mindless masses to small groups : conceptualizing collective behavior in crowd modeling. *Review of General Psychology*, 19(3):215–229.
- Thalmann, S. and Musse, S. (2013). *Crowd Simulation*. Springer.
- Turing, A. (1950). Computing machinery and intelligence. *Mind*, 59(236):433.
- Wang, H. and O’Sullivan, C. (2016). Globally continuous and non-Markovian crowd activity analysis from videos. In *European Conference on Computer Vision*, pages 527–544. Springer.
- Webster, J. and Amos, M. (2020). A Turing test for crowds. *Royal Society Open Science*, 7(200307).
- Wei, X., Lu, W., Zhu, L., and Xing, W. (2018). Learning motion rules from real data: Neural network for crowd simulation. *Neurocomputing*, 310:125–134.
- Yao, Z., Zhang, G., Lu, D., and Liu, H. (2020). Learning crowd behavior from real data: A residual network method for crowd simulation. *Neurocomputing*, 404:173–185.