

International Journal of Semantic Computing
Vol. 11, No. 2 (2017) 1–34
©World Scientific Publishing Company
DOI: 10.1142/S1793351X17002751



Enriching the Fan Experience in a Smart Stadium Using Internet of Things Technologies

Sethuraman Panchanathan*, Shayok Chakraborty†, Troy McDaniel‡,
Ramin Tadayon§ and Bijan Fakhri¶

*Center for Cognitive Ubiquitous Computing (CUBiC)
Arizona State University, Tempe, AZ 85281, USA*

**panch@asu.edu*

†*shayok.chakraborty@asu.edu*

‡*troy.mcdaniel@asu.edu*

§*rtadayon@asu.edu*

¶*bfakhri@asu.edu*

http://cubic.asu.edu

Noel O'Connor^{||}, Mark Marsden**, Suzanne Little^{††}
and Kevin Mcguinness^{‡‡}

*Insight Centre for Data Analysis, Dublin City University
Glasnevin, Dublin 9, Ireland*

^{||}*noel.oconnor@dcu.ie*

***mark.marsden@insight-centre.org*

^{††}*suzanne.little@dcu.ie*

^{‡‡}*kevin.mcguinness@dcu.ie*

David Monaghan

School of Computer Science and Statistics

Trinity College Dublin, College Green, Dublin 2, Ireland

monaghd2@tcd.ie

Rapid urbanization has brought about an influx of people to cities, tipping the scale between urban and rural living. Population predictions estimate that 64% of the global population will reside in cities by 2050. To meet the growing resource needs, improve management, reduce complexities, and eliminate unnecessary costs while enhancing the quality of life of citizens, cities are increasingly exploring open innovation frameworks and smart city initiatives that target priority areas including transportation, sustainability, and security. The size and heterogeneity of urban centers impede progress of technological innovations for smart cities. We propose a Smart Stadium as a living laboratory to balance both size and heterogeneity so that smart city solutions and Internet of Things (IoT) technologies may be deployed and tested within an environment small enough to practically trial but large and diverse enough to evaluate scalability and efficacy. The Smart Stadium for Smarter Living initiative brings together multiple institutions and partners including Arizona State University (ASU), Dublin City University (DCU), Intel Corporation, and Gaelic Athletic Association (GAA), to turn ASU's Sun Devil Stadium and Ireland's Croke Park Stadium into twinned smart stadia to investigate IoT and smart city technologies and applications.

1 *Keywords:* Internet of things; smart stadium; smart city; crowd behavior analytics; object
2 counting.

3 4 **1. Introduction**

5 People increasingly moving to urban centers is shifting the balance between rural and
6 city life. This phenomenon of rapid urbanization has brought about significant
7 changes in where the global population resides: 54% of the global population was in
8 urban in 2014, and by 2050, estimates predict that 64% of the global population will
9 be urban [41]. Rapid urbanization is exacerbating existing concerns of congestion,
10 pollution, accidents, security, and sustainability. For example, it is estimated that by
11 2050, the number of vehicles on the road will double to 2.5 billion. In 2013, the U.S.
12 spent \$124 billion due to traffic congestion, and estimates predict that by 2030, this
13 number will rise to \$186 billion with accompanying increases in “social costs” [59]. By
14 2020, \$13 billion and 1,600 premature deaths are anticipated in costs due to exposure
15 to emissions from idling vehicles during traffic jams. Traffic congestion problems are
16 a worldwide issue; as of 2014, the top 10 most congested cities [59] include Istanbul,
17 Mexico City, Rio de Janeiro, Moscow, Salvadore, Recife, St. Petersburg, Bucharest,
18 Warsaw, and Los Angeles.

19 Cities are seeking ways to reduce complexity and costs, provide better manage-
20 ment, and meet resource needs, while ensuring a high quality of life for its citizens.
21 Many cities have begun to explore open innovation frameworks and smart city
22 initiatives to address the needs of their growing populaces by targeting key priority
23 areas of health, wellness, transportation, safety, security, sustainability, and citizen
24 engagement. Cities that perform well and excel will flourish through the creation of
25 wealth and rises in productivity, paving the way for continued growth and long-term
26 success [29]. Smart city transformations rely upon not only technological and policy-
27 based advancements, but re-imagining traditional approaches to key priority areas,
28 and preparing for scalability challenges due to a city’s sheer size and heterogeneity.
29 We propose the use of a Smart Stadium as a living laboratory to more easily deploy
30 and evaluate Internet of Things (IoT) technologies and smart city solutions
31 by balancing the size and heterogeneity of a smart environment that is small enough
32 to practically trial but large and complex enough to evaluate effectiveness and
33 scalability.

34 Smart Stadium for Smart Living is an initiative developed to join institutions and
35 partners interested in IoT and smart city technologies. The initiative joins Arizona
36 State University (ASU) in Tempe, Arizona; Dublin City University (DCU) in
37 Dublin, Ireland; Gaelic Athletic Association (GAA) of Ireland; and Intel Corpora-
38 tion to turn two stadia — ASU’s Sun Devil Stadium and Ireland’s Croke Park
39 Stadium — into twinned smart stadia with the potential to be world class testbeds
40 for exploring smart city applications and IoT solutions. The projects of this initiative
41 thus far focus on two broad application areas: (i) Enriching the fan/attendee expe-
42 rience; and (ii) Enhancing stadium operation. While the application focus of these
43

1 projects is set in the context of the stadium and stadium-related events, they are
2 relevant to wider smart city application areas. The full scope of projects within this
3 initiative addresses issues of crowd management, fan engagement, event logistics,
4 stadium management, and environmental monitoring, using a variety of deployed
5 sensors such as video cameras and microphones. Given the sheer number of projects
6 within this initiative, the following discussion pertains only to projects targeting
7 enriching the fan experience.

8 Projects to enrich a fan's experience were identified by considering the entire
9 'journey' of an event attendee; that is, not only his or her interactions, behaviors, and
10 actions within the stadium, but all activities involved to attend an event. For ex-
11 ample, a fan's journey may include extensive preparation, perhaps months prior, to
12 attend an upcoming event; planning and coordination to travel to and from the
13 stadium; their involvement on social media leading up to an event as well as during
14 and after an event itself; and activities carrying over to relevant events and gath-
15 erings happening before, during, and after the stadium event itself. This work pre-
16 sents three fan-focused projects targeting efficiency/convenience, safety, and
17 engagement. These projects include: (i) *Crowd Understanding*: Improved safety via
18 vision-based and non-vision-based crowd behavior understanding and analytics; (ii)
19 *Athletic Demonstrator Platform*: Interactive serious gaming stations to support fan
20 engagement while promoting motor learning and athletic training; and (iii) *Wait*
21 *Time and Queue Estimation*: Real-time, accurate access via a mobile app to wait
22 time estimates of lines across a stadium's concession stands, souvenir stands, and
23 restrooms.

24 25 **1.1. Organization and research contributions**

26 The rest of this paper is organized as follows: Section 2 discusses the *Crowd Un-*
27 *derstanding* project. We present an efficient strategy to compute low dimensional,
28 informative features for crowd behavior understanding and anomaly detection.

29 The *Athletic Demonstrator Platform* project, outlined in Sec. 3, is a motor
30 learning environment enabling real-time motion capture, analysis, and feedback.
31 The main contributions of this work include: (1) A fusion approach for low-cost
32 Kinect-IMU motion capture and algorithms for calibration, phase detection, and
33 analysis; and (2) Insight into important research questions pertaining to the de-
34 sign of multimodal feedback including (i) What categories of performance are
35 present in real motor training feedback from a trainer to a subject, and through
36 which modalities do these interactions occur? (ii) How can a system observe these
37 metrics of performance in an individual's motion? and (iii) Does individual pre-
38 ference play a role in the assignment of modalities to feedback in a multimodal
39 environment?

40 Section 4 presents the *Wait Time and Queue Estimation* project. Our research
41 contributions in this project include: (1) A novel active learning framework to
42 identify the salient and exemplar instances from large amounts of unlabeled data to
43

1 train an object counting model and (2) Incorporating only binary (yes/no) feedback
2 into the algorithm in order to reduce the labeling burden on the user.

3 Finally, we conclude with discussions in Sec. 5.
4

5 **2. Crowd Understanding** 6

7 Sports Stadiums are multi-purpose venues within our cities where thousands gather
8 for events including sporting contests, music concerts as well as business and aca-
9 demic conferences. However, with such large gatherings of people there are signifi-
10 cant risks to public safety which must be addressed. Improving our understanding of
11 the behavior of such large crowds of people within a stadium can help maintain safety
12 and security for all involved. Early detection and a rapid response time are essential
13 in any emergency situation, especially in a highly congested public space such as a
14 stadium. To address this issue, we have developed an efficient computer vision al-
15 gorithm for detecting unusual crowd behavior in real-time on a commodity CPU.
16 Both Sun Devil Stadium (56,200 capacity) and Croke Park Stadium (82,300 ca-
17 pacity) have been fully designed to ensure the safety of all visitors, but the Smart
18 Stadium project aims to exploit visual and non-visual sensor data to gain additional
19 insight into the dynamics of crowds which will help improve the already excellent
20 safety standards.

21 The crowd understanding project uses existing CCTV camera footage from Croke
22 Park Stadium to extract scene-level holistic features and detect unusual crowd be-
23 havior at the frame level. Long-term, the system aims to learn a “steady state” of
24 what normal crowd behavior patterns look like across numerous cameras within a
25 stadium, and therefore, be able to determine when crowds don’t behave according to
26 expected patterns, and alert support staff.

27 **2.1. Crowd understanding implementation** 28

29 The objective is to design a low dimensional set of features that are quick to compute
30 and capture sufficient holistic information about objects moving in a scene to allow
31 straightforward discrimination between normal and abnormal events. The developed
32 technique for crowd behavior anomaly detection uses a set of efficiently computed,
33 easily interpretable, scene-level holistic features [40]. These features are calculated by
34 analyzing local motion patterns across a crowded scene. This low-dimensional de-
35 scriptor combines two features from the literature: crowd collectiveness [52] and
36 crowd conflict [27], with two newly developed features: mean motion speed and a
37 unique formulation of crowd density [40].

38 Crowd collectiveness is a scene-independent holistic property of a crowd system,
39 which can be defined as the degree to which individuals in a scene move in unison
40 [52]. Zhou *et al.*'s [52] method for measuring this property analyzes the tracklet
41 positions and velocities found in the current frame and constructs a weighted ad-
42 jacency matrix. The edge weights within each matrix column are summed and the
43

1 mean is calculated. This mean value corresponds to the overall collectiveness level for
2 the current frame.

3 Crowd conflict is another scene-independent holistic crowd property, which can
4 be defined as the level of friction/interaction between neighboring tracked points
5 [52]. Shao *et al.* [52] efficiently calculate this property by summing the velocity
6 correlation between each pair of neighboring tracked points in a given frame.

7 Crowd density can be defined as the level of congestion observed across a scene at
8 a given instant. The proposed approach to calculating this feature firstly divides the
9 scene into a fixed size grid (10×10) and counts the number of grid cells currently
10 occupied by one or more tracked points. Equation (1) is then used to calculate the
11 crowd density level for the current frame. A 10×10 grid was chosen to provide
12 sufficient granularity in the density calculation, with the aim being for each grid cell
13 to roughly contain one or two pedestrians in most surveillance scenarios. There are
14 obvious limitations in terms of scale invariance with this feature, however the main
15 objective is not pixel perfect accuracy but to measure a useful crowd property in a
16 highly efficient manner.

$$17 \quad \text{CrowdDensity} = \frac{\text{Occupied Grid Cells}}{\text{Total Grid Cells}} \quad (1)$$

18
19
20 Figure 1 depicts the proposed crowd density feature calculated using footage from
21 a CCTV camera covering a concession area at Croke Park Stadium during a busy
22 match. As shown, the density level increases significantly once the gates open (yel-
23 low), fall once the match begins (green), and then spike again at half time (red).

24 The heat map in Fig. 1 is taken from a video animation that illustrates how the
25 density level at various stadium locations changes over time. This was produced by
26 calculating the density level over a full match day for each camera location and
27 updating the color for each section to correspond to the density level [0.0–1.0] at that
28 time point. Using this method, the distribution of people throughout a stadium over
29 the course of a match day can be visualized. The visualization can also be sped up to
30 show the changes that take place over hours in a matter of minutes. Figure 2 shows
31 this feature being calculated on an image from the UMN dataset.

32 The mean motion speed observed within a crowded scene provides a coarse, scene-
33 level feature that can be extracted very efficiently. Our approach estimates this
34 crowd property by calculating the magnitude of each tracklet velocity vector in
35 the current frame and finding the mean. While conceptually simple, experiments
36 show that the inclusion of this feature noticeably improves the accuracy of crowd
37 behavior anomaly detection. Each of these features captures a distinct aspect of
38 crowd behavior.

39 Our holistic features are extracted for each frame in a given video sequence using
40 the following steps. Firstly, the scene foreground is segmented using the Gaussian-
41 mixture based method of KaewTraKulPong and Bowden [33] before interest points
42 are tracked using a KLT tracker [58]. These local trajectories or tracklets are then
43 analyzed to calculate four holistic features for each frame. This high-level descriptor

6 *S. Panchanathan et al.*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

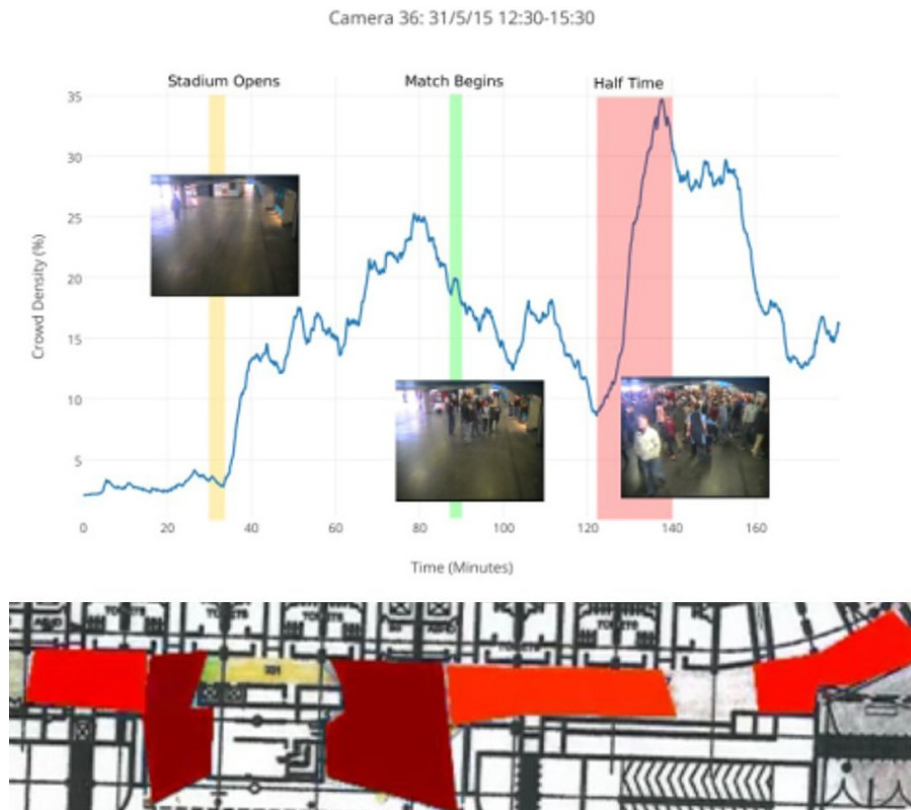


Fig. 1. (Color online) *Top*: Changes in crowd density calculated for a concession stand during a busy match day at Croke Park Stadium. *Bottom*: A heat map visualization showing differences in crowd density at different stadium locations within Croke Park.



Fig. 2. Crowd density calculation grid for a scene from the violent-flows dataset. Each green square corresponds to an occupied grid cell (crowd density in this frame = 57%).

1 of crowd behavior can be computed in real-time (30+ frames per second) even on
2 commodity hardware (e.g., an Intel i5 CPU).

3 Anomalous crowd behavior then needs to be detected using this crowd behavior
4 descriptor. We investigate two anomaly detection approaches, covering two possible
5 situations: (1) When only “Normal” behavior training data is available; and (2)
6 when both “Normal” and “Abnormal” behavior training data are available. Each
7 require the following pre-processing steps: All individual features are firstly scaled to
8 lie within the range $[0, 1]$, with respect to the range of training data values. Nor-
9 malization is then performed by dividing by the maximum magnitude vector in the
10 training set. The low-dimensional descriptor used results in almost negligible training
11 and classification times for reasonably sized datasets.

12 We use a Gaussian Mixture Model (GMM) to perform outlier detection when only
13 normal behavior training data is available. The GMM configuration (number of
14 mixture components and type of co-variance matrix) for a given experiment is se-
15 lected as the one that minimizes the Bayesian Information Criterion (BIC) value [48]
16 on the training data. The selected model is then used to calculate the log probabilities
17 for the full set of training frames, and the distribution of these log probability values
18 is used to decide upon an outlier detection threshold using Otsu’s method [43]. Test
19 frames are then classified as abnormal or normal by using the fit mixture model to
20 calculate their log probability and applying the adaptive threshold generated from
21 the training data.

22 We use a discriminative model (binary classifier) for outlier detection when
23 both normal and abnormal training data are available. Specifically, we trained a
24 Support Vector Machine (SVM) with an RBF kernel on test frames labeled as
25 normal and abnormal. The default value of 1.0 was used for the SVM regularization
26 parameter C .

27 **2.2. Crowd understanding results**

29 The proposed method is evaluated on two distinct crowd behavior anomaly datasets:
30 (i) the UMN dataset^a; and (ii) the violent-flows dataset [27]. These benchmarks
31 assess the ability of a given approach to detect unusual crowd behavior at the frame-
32 level and video-level, respectively. All experiments were carried out using MATLAB
33 2014a and Python 2.7 on a 2.8 GHz Intel Core i5 processor with 8GB of RAM.

34 The UMN dataset contains 11 sequences filmed in 3 different locations. Each
35 sequence begins with a period of normal passive behavior before a panic event/
36 anomaly occurs toward the end. The objective here is to train a classifier using frames
37 from the initial normal period and evaluate its detection performance on the sub-
38 sequent test frames. Classification is performed at the frame level and results are
39 compared in terms of the receiver operating characteristic (ROC) curve’s area under
40 the curve (AUC). For each of the three scenes, the initial 200 frames of each clip are
41 combined to form a training set, with the remaining frames used as a test set for that

42
43 ^a<http://mha.cs.umn.edu/Movies/Crowd-Activity-All.avi>.

Table 1. BIC values calculated during the GMM selection stage for the UMN dataset.

No. of mixture components	BIC
1	-20015
2	-21810
3	-22047
4	-21940

scene. This results in a roughly 1:2 split between training and test frames for each camera location and will be referred to as the single scene experiment. While this dataset is quite limited in terms of size and variation, it does provide a good means of performance evaluation during the development of a crowd anomaly detection algorithm. Since no abnormal frames are made available for training in this experiment, the GMM-based detection approach is used. Table 1 presents the BIC values calculated during the GMM selection stage, with a 3-component model ultimately used. A full co-variance matrix GMM resulted in a lower BIC value in all cases and was therefore used. Figure 3 presents the ROC curves for all three UMN scenes individually. A cross-scene anomaly detection approach is also taken, where for a

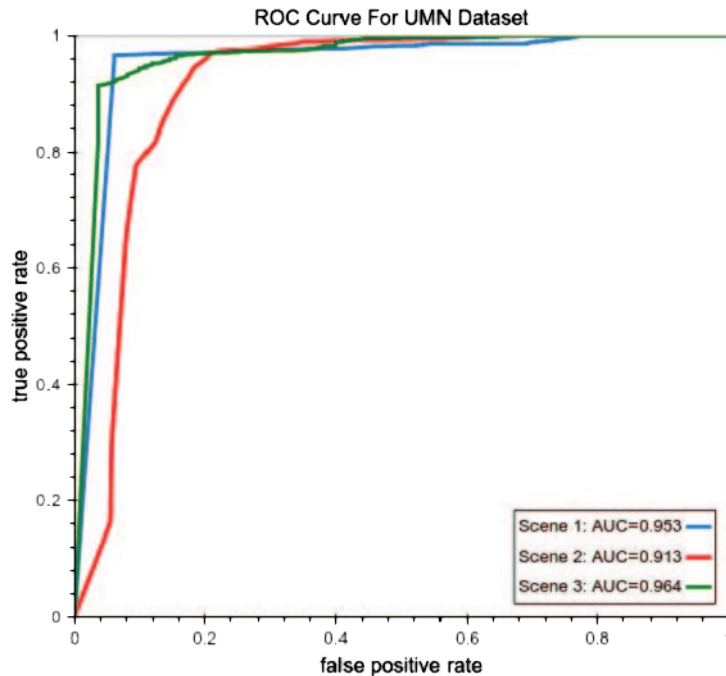


Fig. 3. Receiver operating characteristic (ROC) curve and associated area under the curve (AUC) for each UMN scene.

Table 2. ROC curve AUC performance and processing speed on the UMN dataset.

Method	AUC	Speed (FPS)
MDT	0.995	0.9
CM	0.98	5
SFM	0.97	3
Proposed Method (Single Scene)	0.929	40
Proposed Method (Cross-Scene)	0.869	40

given UMN scene, the training frames from the two other scenes are used to generate the GMM.

Table 2 compares the two variants of the proposed method with the leading approaches in terms of AUC and processing frame rate. The proposed approach achieves competitive classification performance with the state-of-the-art at just a fraction of the computational cost. The cross-scene experiment, while inferior in terms of classification performance, is noteworthy in that each scene was classified using training data only from other surveillance scenarios.

The violent-flows dataset contains 246 clips containing violent (abnormal) and non-violent crowd behavior. Classification is performed at the video level. A 5-fold cross validation evaluation approach is taken and results are compared in terms of mean accuracy. As both normal and abnormal training examples are available in this dataset, the proposed SVM-based classification approach is used. The majority classification found among the frames of a given clip is used as the overall result for that clip. An alternate approach is also taken where only the normal training examples are used, and the proposed GMM-based outlier detection approach is taken.

Table 3 presents the BIC values calculated during the GMM selection stage, with a 4-component model ultimately used. A full co-variance matrix GMM resulted in a lower BIC value in all cases and was therefore used. For this GMM-based approach the histogram of frame log probabilities for a given test clip is generated and the mode value is used to classify the overall clip by applying the Otsu threshold generated from the training data. Table 4 compares the two variations of the proposed technique with the leading approaches in terms of mean accuracy and processing

Table 3. BIC values calculated during the GMM selection stage for the violent-flows dataset.

No. of mixture components	BIC
1	-51758
2	-223161
3	-274742
4	-327545

Table 4. Mean accuracy and processing speed on the violent-flows dataset.

Method	Accuracy (%)	Speed (FPS)
SD	85.4	N/A
HOT	82.3	N/A
ViF	81.3	30
CM	81.5	5
Proposed Method (SVM)	85.53	40
Proposed Method (GMM)	65.8	40

Table 5. The contribution of each feature toward mean detection accuracy on the violent-flows dataset using proposed SVM-based detection approach.

Feature	Accuracy when excluded (%)
Crowd Collectiveness	75.2
Crowd Conflict	65.5
Crowd Density	63.5
Mean Motion Speed	81.2

frame rate. Table 5 highlights the contribution of each feature towards the achieved anomaly detection accuracy on the violent-flows dataset using the SVM-based variant. As shown, leaving out any individual feature results in a noticeable decrease in anomaly detection accuracy.

The SVM-based variant achieves state-of-the-art performance on the violent-flows dataset with a mean accuracy of $85.53 \pm 0.17\%$. The GMM-based variant achieves a very respectable $65.8 \pm 0.15\%$ accuracy, which is particularly impressive considering only half the training data, containing no violent behavior, is used in this case. The approach also achieves noticeably faster computational performance.

The proposed scene-level holistic features are easily interpretable, sensitive to abnormal crowd behavior, and can be computed in better than real-time (40 frames per second) on commodity hardware. The approach was demonstrated to improve upon the state-of-the-art classification performance on the violent-flows dataset. Future work will attempt to improve upon certain limitations of the approach such as the scale issues present in the crowd density feature, possibly using an adaptive grid cell size. Moreover, this descriptor will be used to label specific crowd behavior concepts in larger and more challenging datasets.

3. Athletic Demonstrator Platform

Modern technology has made motion sensing more accessible and prevalent than ever before, with the rise of low-cost motion-sensing hardware such as Microsoft's Kinect camera. Similarly, multimodal feedback has become increasingly ubiquitous through the introduction of haptic, visual, and audio feedback mechanisms in phones

1 and game controllers, among other devices. Thanks to this evolution of technology,
2 motor training is now more accessible to the everyday user, leading to a surge in
3 studies on motor learning in Human-Computer Interaction (HCI). With modern
4 technology, an automated system is capable of observing and reacting to a great deal
5 of information pertaining to a user's motion, and with the inclusion of expert data,
6 the system can evaluate and provide feedback on this motion in real-time, leading to
7 a new wave of "unsupervised" motor training wherein an individual interacts with a
8 system, rather than a real trainer, to gain proficiency in motor skills. This type of
9 training has a variety of applications ranging from rehabilitation [54] to sports
10 training [63]. This technology solves a critical problem in the field: a user must
11 regularly perform and receive feedback on a motor task to improve at that task at a
12 steady rate [11], but since trainer availability is limited, user compliance with this
13 training can stagnate over time [46].

14 To provide the type of feedback on motor performance that a user can consider
15 useful in comparison to a real trainer, an automated system should perform the
16 following tasks: (i) The system should accurately capture a user's motion using
17 commonly accessible technology (without the complex setup typically encountered
18 in a laboratory or clinical environment); (ii) The system should automatically recognize,
19 classify, and represent the various segments and elements of a motion; (iii)
20 The system should be able to accurately interpret motion data to form an assessment
21 of a user's performance; and (iv) The system should provide feedback on this assessment
22 that is understandable and meaningful to the user so that the user can
23 improve his or her motion in the next attempt.

24 Various aspects of the motion itself should also be considered in the provision of
25 real-time feedback including the type of motion (rehabilitation vs. sports, for example),
26 the user's proficiency level and previous experience, the complexity of the
27 motion task (typically determined by observing the number of limbs involved in the
28 motion), the type of information observed (spatial and temporal aspects of the
29 motion), the assignment of modalities to different aspects of feedback, and the timing
30 of feedback (for example, concurrent vs. terminal), among others.

31 Here we present a platform for the provision of automated multimodal feedback
32 for motor performance in a variety of motor training scenarios. The proposed Athletic
33 Demonstrator Platform implements real-time motion analysis and feedback to
34 facilitate a motor learning environment that is both useful and stimulating to enrich
35 fan engagement, excitement, and competitiveness. As part of future work, the
36 platform has potential for athlete training.

37 **3.1. Related work in motion capture**

38 The field of Motion Capture, or "MoCap", is a widely studied area in which various
39 techniques and methods have been applied toward the quantification and digital
40 representation of a human's motion in an automated system [60]. Perhaps the most
41 cutting-edge system to date for this task is the Vicon system, which uses accurate
42
43

1 and high-quality tracking of worn body markers to record and analyze complex
2 motion. However, this system is expensive and often restricted to laboratory envir-
3 onments or professional motion capture scenarios due to its complexity, making it
4 impractical for the typical user. As a result, cost-efficiency has become a recent
5 concern in the field, leading to the rise of more affordable alternative systems [21, 19]
6 which rely on computer vision [61] and depth-sensing [20, 62] to form lower-quality
7 estimates of a user’s body orientation and joint movement during a motion. Inevi-
8 tably, these mechanisms are subject to the errors caused by occlusion of body parts
9 and other issues relating to static camera sensing.

10 In addition to camera-based techniques, some wearable alternatives to Vicon
11 exist for motion capture. One popular alternative is Inertial Measurement Units
12 (IMUs), body-worn 3D motion sensors which offer an accuracy that can compete
13 with the gold standard [2, 36]. One example of IMU application in MoCap is the
14 XSens system [47]. These systems have seen limited success in practice due to the
15 accumulation of calculation errors which affect the accuracy of their measurements
16 over a time period. Furthermore, if only IMUs are used to handle motion capture, a
17 significant amount would be needed to cover all body motion, which can be very
18 costly. To address this issue, we can take a hybrid approach, which utilizes both
19 worn IMUs and Kinect depth camera sensing, and fuses the readings from these two
20 devices [18, 17] to correct for the accuracy errors of one while solving the occlusion
21 issue of the other.

22 In previous work [3], we have shown that the hybrid approach, in combination of
23 2–3 IMU sensors with real-time joint-tracking camera data and the implementation
24 of advanced algorithms for calibration, phase detection, and analysis, can provide a
25 low-cost yet accurate mechanism for capture of motor activity. Here we discuss the
26 design of a platform that utilizes the fusion approach, along with multimodal feed-
27 back in its final design, to provide learning interaction for motion of any type,
28 depending on the location of IMU sensors worn on the body.

30 **3.2. Athletic demonstrator platform implementation**

31 The athletic demonstrator platform utilizes a combination of two IMU sensors
32 (currently wrist-worn, but can be reconfigured), the Microsoft Kinect V2 depth
33 camera for joint tracking, and the Unity engine for multi-platform game develop-
34 ment, as its core components. Each IMU sensor communicates with a central com-
35 puter running the platform’s software over a Bluetooth connection at 256 Hz and
36 includes accelerometer and gyroscope output for position and orientation. The sys-
37 tem first calibrates the sensing by requiring the user to stand in a “T-pose”, thus
38 synchronizing the IMU sensors to the Kinect’s coordinate space. After this, the
39 skeletal tracking of the Kinect is fused with the readouts of the IMU sensors to
40 determine accurate joint positioning based on techniques shown in [18, 17, 3]. The
41 joint data tracked by movement of the worn IMUs are utilized to determine the
42 position and orientation of those joints, while the Kinect’s data is utilized to
43

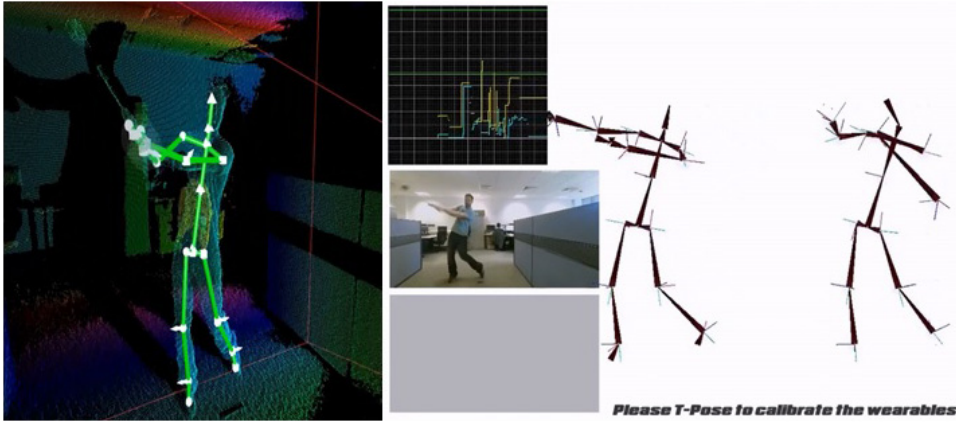


Fig. 4. Proposed low-cost Kinect-IMU motion capture fusion approach. *Left:* Demonstration of Kinect skeletal tracking. *Right:* Calibration phase and fusion of Kinect and IMU data.

determine the position and orientation of all other joints during a motion. This fusion technique is shown in Fig. 4.

To learn a motion in the current design of the platform, the user first views a demonstration of the motion by an expert through an on-screen video, which is accompanied by both an avatar representation of the fused Kinect/IMU data and a graph which depicts 3D IMU accelerometer information over time. Having viewed this demonstration, the user is then asked to attempt the motion under the same interface, with a 3D virtual avatar mirroring the user's motion as a form of concurrent visual feedback. Mechanisms are also in place for the provision of haptic feedback and audio cues at key points during this motion attempt, although the concurrent feedback used is purely visual in the initial prototype shown in Fig. 5.

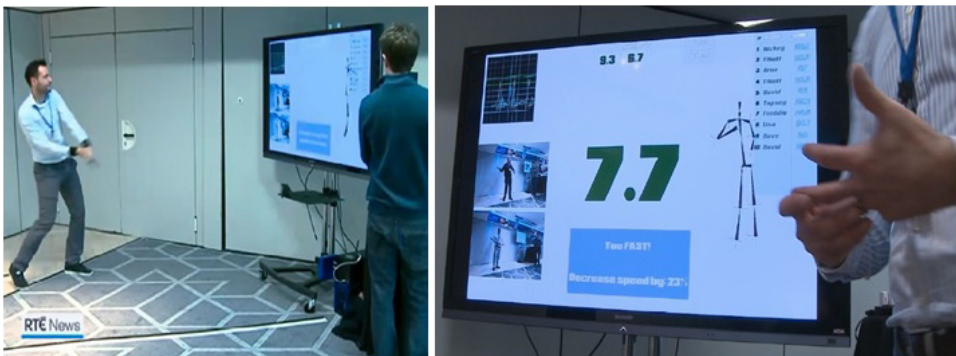


Fig. 5. Athletic Demonstrator Platform. *Left:* Live demonstration of the platform for the Irish sport of hurling. *Right:* Gamified score feedback based on expert player with top ten scoreboard to promote competitiveness.

1 Once the user completes an attempt of the motion, he or she is then provided
2 terminal feedback on performance using a scoring system that depicts the proximity
3 of the user's motion, captured with both the IMU sensors and Kinect camera, to the
4 motion sample provided by the expert. This scoring is accompanied by feedback on
5 the user's speed, specifying whether the user should slow down or speed up the rate of
6 motion on the next attempt. This terminal feedback provides an overview of the
7 individual's performance for a single attempt; after several attempts, the user is given
8 an overall score for the motion as a final measure of his or her current performance.
9 This overall score is submitted to a leaderboard indicating the best performances on
10 that motion, which can be used either by a single user to determine how his or her
11 performance is progressing over time, or by multiple users to compare their perfor-
12 mances on the same task.
13

14 **3.3. Athletic demonstrator platform: Related studies**

15 The athletic demonstrator platform was designed to be highly configurable, allowing
16 for different multimodal designs for the provision of concurrent and terminal mul-
17 timodal feedback on motor performance. It was also designed to handle a large
18 variety of motions with the fusion capture method. This flexible design has led to a
19 series of research questions on multimodal implementation which we have addressed
20 through research studies. These questions include: (i) What categories of perfor-
21 mance are present in real motor training feedback from a trainer to a subject, and
22 through which modalities do these interactions occur? (ii) How can a system observe
23 these metrics of performance in an individual's motion? (iii) Does individual pre-
24 ference play a role in the assignment of modalities to feedback in a multimodal
25 environment? Findings related to each of these questions are discussed below, and
26 together they will inform the final design of the athletic demonstrator platform.
27

28 **3.3.1. Case study on categories of feedback**

29 To address the first question, real motor training scenarios were observed as part of a
30 case study between a subject and a martial arts trainer. The goal of the first phase in
31 this case study was to determine what forms of feedback occur in real-time as the
32 subject interacts with the trainer, and in what modalities these interactions occur.
33 To achieve this, a live training session was recorded between these individuals on
34 video, and specific instances of feedback given by the trainer during the interaction
35 were noted. For these feedback instances, both the modality of feedback and the
36 category of feedback were determined.
37

38 Through this study, detailed in [55], three main categories of real-time motor
39 feedback were identified: (i) Posture: a spatial measure of feedback relating to the
40 configuration of the user's body and limbs during motion; (ii) Progression: a spatial
41 metric which relates to the range of motion and the accuracy of an individual's
42 motion trajectory compared to the ideal motion; and (iii) Pacing: a temporal
43

1 measure representing the speed of an individual's motion, its consistency, and its
2 comparison to the ideal rate of motion.

3 Together, the three categories above constitute a complete representation of
4 motor performance. While they were applied in this case to rehabilitative motion,
5 these categories can be applied toward sports motions as well. For example, a
6 football or baseball throw relies on proper configuration of the elbow and grip of the
7 ball (posture), momentum of forward motion prior to release (pacing), and release
8 of the ball at the correct moment to achieve an ideal trajectory (progression), as
9 indicated in [4].

10 The primary modalities of feedback discovered were audio (delivered as verbal
11 feedback from the trainer), visual (delivered as demonstrations of correct motion by
12 the trainer), and haptic (delivered as guiding nudges by the trainer to ensure the
13 subject reaches the desired range of motion).

14 The first design of the athletic demonstrator platform uses primarily postural
15 information to deliver score-based feedback as it is often the most important cate-
16 gory of feedback for performance in sports motion, but other categories of feedback
17 will be added to the platform to allow for a richer set of information on performance
18 with the potential to improve motor learning.

19 3.3.2. *Case study on quantification of feedback*

20 Once the categories of feedback and modalities of feedback in motor learning were
21 determined, the next step was to determine how an automated system can observe and
22 provide feedback on an individual's performance in each of these categories. In the
23 athletic demonstrator platform, the system has access to a 3-dimensional represen-
24 tation of a user's motion as a time-series dataset extracted from fused Kinect and IMU
25 data. In a similar project, "Autonomous Training Assistant" [56], we found that all
26 three categories of performance can be inferred from this data by comparing to expert
27 motion. To determine when to provide feedback, it is useful to set a threshold at which
28 an error can be identified in each category. In other words, once the user's motion
29 deviates from the expert's motion by a targeted amount, feedback can be given to
30 correct that motion. We call this method "tolerance thresholding", and it can be used
31 to refine our definition of each modality of feedback in the following ways:
32
33

34 Postural data may be described as the way in which an individual's joint angles,
35 and for the relevant joints in a motion, relate to one another and to an expert's joints
36 in 3D space. At any given point in time, the Kinect can determine the location and
37 angle of a user's shoulders, elbows, wrists, knees, and other joints for coarse postural
38 adjustment (fine postural adjustment requires more sensitive recording mechanisms
39 which may be implemented, for example, as wearable sensors). For each joint related
40 to the posture of a motion, we can define postural performance as the proximity of a
41 user's joint angle to that of the expert at that point in time, adjusted using Dynamic
42 Time Warping (DTW) methods to ensure the two are equally scaled. The tolerance
43

1 threshold for posture can then be defined as the maximum amount, in degrees, that a
2 subject's joint is allowed to deviate from the expert's joint for the motion to be
3 considered "correct". Deviation beyond this point can be considered an error and
4 feedback can be given accordingly.

5 For progression, a system can observe the trajectory of a user's motion and
6 compare it to the expert's trajectory for assessment, noting the proximity of the two
7 in time-adjusted 3D space. It would be difficult to perform this assessment for every
8 recorded data point of a motion in real-time; instead, only the most essential points,
9 i.e., "critical points", representing the shape and form of the motion may be ob-
10 served. An arc, for example, can be represented as a progression of five points in
11 space. At these points along the motion, a user's data point can be compared to an
12 expert's data point using a standard 3-dimensional distance measure. The tolerance
13 threshold can be defined for progression as the maximum allowed distance between
14 two critical points for a motion to be "correct" at that point in time.

15 Finally, for pacing, a system can observe the rate at which a user progresses from
16 the start to the end of a motion, compared to an expert. The difference between these
17 two forms the user's error in pacing. A user's motion is allowed to be slower or faster
18 than the expert's motion up to a specific tolerance threshold to be considered
19 "correct". Beyond this range, feedback is necessary. Note that in this case a system
20 must specify, as the trainer does in our first example above, whether the motion is
21 slower or faster than the desired rate so that the user can make adjustments in the
22 proper direction.

23 The final design of the athletic demonstrator platform can use the above metrics,
24 determined through the quantification of the trainer's feedback in the case study, to
25 form a detailed profile of a user's performance for a motion.

26 27 3.3.3. *Case study for individual preference*

28 To determine the effects of individual preference on the effectiveness of a modality in
29 a multimodal feedback scenario in motor learning, a study was designed with the case
30 study subject wherein a multimodal environment with the Autonomous Training
31 Assistant was presented. In this environment, the subject was asked to complete a
32 series of simple motor exercises with two feedback conditions. In the first condition,
33 modalities (haptic, audio, visual) were assigned to feedback categories (posture,
34 progression, pacing) based on the mapping suggested by the review of Sigrist *et al.*
35 [53] for concurrent multimodal feedback. In the second condition, the subject was
36 able to choose the mapping based on individual preference. The subject then com-
37 pleted a series of three basic martial arts exercises (umbrella motion, twirl motion,
38 and witik motion) assigned by the subject's martial arts trainer using the Autono-
39 mous Training Assistant interface for each condition.

40 Each exercise was completed in a 2-minute interval with breaks in-between, and a
41 longer break between the two conditions to prevent fatigue and minimize learning
42 effects. The subject's performance was measured in each category using error rate in
43

1 each of the three performance categories. It was found that in the preference con-
2 dition, the subject performed significantly better in categories that were mapped
3 differently by preference, while performance in unchanged categories of feedback
4 remained the same between conditions. This improvement held consistently across
5 all three exercises, suggesting that individual preference in multimodal feedback
6 selection may have an effect on performance in multimodal training environments.
7 Furthermore, it was observed that, in both conditions, the subject would focus on a
8 single modality of feedback over the other modalities in the presence of multimodal
9 feedback. In this case, the subject seemed to focus on haptic feedback as indicated by
10 an increased responsiveness to feedback in modality.

11 Further studies on a larger scale using the athletic demonstrator platform can
12 help determine whether these observations are generalizable across a variety of users
13 and motor training exercises. Currently, the platform is capable of providing ter-
14 minal feedback using a score system to indicate performance on a motor exercise via
15 expert data, which is purely visual. Haptic feedback will be added to the platform
16 through the introduction of wrist-worn vibrotactile motors to guide the user at
17 regular intervals through movements as initially described in [55]. Furthermore,
18 rhythmic audio cues will be added to accompany both the demonstration and at-
19 tempt screens of the platform to help the user compare the rhythm of their motion to
20 that of an expert as an additional form of evaluation.

21 22 **4. Wait Time and Queue Estimation**

23 The objective of this project is to enrich the fan experience by providing access to
24 wait times at restrooms and concession stands via a mobile app. Such a technology
25 will allow fans to maximize their time watching and enjoying a game rather than
26 waiting in long lines during the course of a game. We adopt a computer vision based
27 approach to count the number of people in a queue. We assume the presence of
28 cameras in strategic locations in the vicinity of restrooms and concession stands; the
29 video feed from these cameras is analyzed to accurately estimate the count of people
30 in the queues. Once the count is obtained, wait times can be obtained from the
31 average service time per person.

32 Counting the number of objects in an image is a problem of paramount practical
33 importance. It arises in myriads of real-world applications including crowd behavior
34 monitoring, security and surveillance, medical imaging and developing infra-
35 structures for smart cities, among others. Counting is often posed as a supervised
36 learning problem, where a regression function is learned directly from some global
37 image features to the number of objects in it. The regression-based algorithms depict
38 commendable performance in counting the number of objects in images. However,
39 they necessitate a large amount of manually annotated data from human oracles to
40 train the regression models. This is an expensive process in terms of time, labor and
41 human expertise. Further, annotating an image for object counting requires much
42 more time and effort than annotating an image for a face recognition or an object
43

**37 Pedestrians****21 Pedestrians**

Fig. 6. Two images with ground truth object counts.

recognition application, for instance. Figure 6 shows two images of pedestrians in a shopping mall and in an outdoor walkway, together with the corresponding ground truth counts. It is evident that hand-labeling such images with counts of objects is an extremely tedious task and highly prone to annotation errors. Thus, while annotating a face/object image requires only a cursory glance, counting objects is much more laborious and demands significantly more time, effort and concentration from a human oracle. It is therefore a significant challenge to obtain a large amount of labeled training images with the exact counts of the number of objects in them. In this paper, we propose a novel learning framework, with the following two features, to address this fundamental problem: (i) the first feature, binary user feedback, relaxes the requirement of exact count of objects as labels; (ii) the second feature, active sampling, aims to reduce the amount of labeled training data (and hence, the amount of manual effort) required to induce a regression model. These are detailed below:

4.1. *Binary user feedback*

We present a general learning framework which requires only binary (yes/no) feedback from the user. During each instance of interaction, the human user is presented with an image and a threshold (an integer) and he merely has to say whether the number of objects in the image is greater than the threshold or not. Providing such an input is extremely easy; it is also less prone to human errors as the number of objects in an image needs to be compared only against a given threshold every time.

In order to quantitatively compare the two types of user feedback: exact (where the exact count of the number of objects needs to be provided) and binary (where only a *yes/no* response needs to be provided about whether the count of objects is greater than a given threshold), we conducted experiments on 15 users. Each user was shown a sequence of four random images, one from each of the following

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Exact

**How many pedestrians are there
in the image?**

Binary

**Is the number of pedestrians
greater than 25?**

Fig. 7. Exact and binary annotation examples (the thresholds for the binary annotations in the experiment were computed using our algorithm).

datasets: the Mall [37], the UCSD pedestrians [10], the Fudan [57] and the TRANCOS [42]. These datasets contain images captured under challenging real-world conditions. For the first two images, the user was asked to provide exact annotations and for the next two, binary annotations (the thresholds for the binary annotations in the experiment were computed using our algorithm and is detailed in Sec. 4.4). Sample images are shown in Fig. 7.

We computed the response time (time taken to annotate an image) for both exact and binary annotations; we also requested each user to provide an overall score between 1 (extremely difficult) and 10 (extremely easy) about the ease of binary annotation over exact annotation. The results are depicted in Table 6. We note that the binary feedback requires much lesser user interaction time than the exact annotations. Moreover, as evident from the scores, users were much more comfortable with the binary annotations since it does not involve the strenuous task of counting the exact number of objects in an image. In summary, binary feedback provides an extremely appealing user interaction model for the vision based object counting application.

Table 6. User study results on exact and binary annotations.

Annotation type	Mean response time (seconds)	Mean score
Exact	20.98 ± 4.34	9.26 ± 0.88
Binary	11.32 ± 4.74	

4.2. *Active sampling*

Active learning algorithms have gained popularity in reducing human annotation effort for training machine learning models. When exposed to large amounts of unlabeled data, such algorithms automatically identify the salient and prototypical instances which can augment maximal information to the underlying models [50]. While serial-query based active learning algorithms query a single unlabeled sample at a time, batch mode active learning (BMAL) techniques query a batch of samples simultaneously for manual annotation and are effective in utilizing the presence of multiple labeling oracles. BMAL has been successfully used in a variety of computer vision applications such as face and facial expression recognition [8], image and video retrieval [31] and image clustering [22] among others. In this work, we exploit batch sampling algorithms to identify the exemplar images that need to be queried for labels, from vast amounts of unlabeled image samples. This can tremendously reduce the human annotation effort required to induce the regression learner, as only the exemplar samples identified by the algorithm need to be labeled manually. To the best of our knowledge, this is the first research effort to address the problem of active data selection with binary user feedback in the context of vision-based object counting. Although validated on object counting in this paper, the proposed algorithm is generic and can be used in any regression-based application where the exemplar instances need to be selected from large amounts of unlabeled data and a model needs to be trained based on binary user feedback.

4.3. *Related work*

In this section, we present a survey of vision based object counting methodologies as well as a brief survey of active learning.

Vision-based Object Counting: Unsupervised learning techniques have been used to address the vision-based object counting problem. They mostly rely on grouping objects based on self-similarities [1] or motion similarities [44]. However, these techniques are limited in their counting accuracy, which has paved the way for supervised learning approaches for counting. Detection-based supervised algorithms attempt to train object detectors (e.g. pedestrian detectors) to localize the individual object instances within an image; the count is then estimated as the number of localized objects. Common approaches of detection-based counting include non-maximum suppression [16], generative techniques [5] and blob tracking [26] among others. Fusion based approaches have also been explored for people counting [32] which rely on multiple sources of information (low confidence head detections, repetition of texture elements and frequency domain analysis) to estimate counts of individuals in extremely crowded images. However, all these techniques need to solve object detection, which is a challenging computer vision problem, especially for overlapping and occluded instances.

The regression-based counting techniques avoid solving the hard detection problem and attempt to learn a mapping directly from some global image feature to

1 the number of objects in it. Cho *et al.* [13] used edge features together with back-
2 ground subtraction and reported promising performance while estimating crowd
3 density using a neural network. Kong *et al.* [34] performed feature normalization in a
4 neural network model to deal with perspective projection and camera orientation and
5 proposed a viewpoint invariant approach to count pedestrians. Chen *et al.* [12] re-
6 cently proposed a scalable multi-output regression model to estimate people count in
7 spatially localized regions. Marana *et al.* [38] postulated that images of low density
8 crowds tend to present coarse textures while images of dense crowds present fine
9 textures; a self-organizing neural network was used to extract features from such
10 images for crowd density estimation. Lempitsky and Zisserman [35] proposed to
11 recover a density function F as a real function of the pixels in an image I , so that
12 integrating F over the entire image yields the count of the number of objects in it.
13 Very recently, deep learning algorithms have been exploited to count the number of
14 objects in an image [49].

15 **Active Learning:** Active learning is a well-studied problem in machine
16 learning. Several techniques have been developed over the last several years and a
17 review of these can be found in [50]. In a typical pool-based batch mode active
18 learning (BMAL) setting, the learner is exposed to a pool of unlabeled instances
19 and it iteratively queries batches of samples for annotation. Initial BMAL tech-
20 niques were largely based on heuristic measures such as maximizing the diversity of
21 the selected samples, computed as their distance from the decision hyperplane [6].
22 More recently, optimization based strategies have been proposed which have been
23 shown to outperform the heuristic approaches. Hoi *et al.* [30] used the Fisher in-
24 formation matrix as a measure of model uncertainty and proposed to query the set
25 of points that maximally reduced the Fisher information. Semi-supervised BMAL
26 algorithms have also been explored in the context of SVMs, where a kernel function
27 was first learned from a mixture of labeled and unlabeled samples, which was then
28 used to identify the informative and diverse examples through a min-max frame-
29 work [31]. Guo and Schuurmans [25] proposed a discriminative strategy that se-
30 lected a batch of points which maximized the log-likelihoods of the selected points
31 with respect to their optimistically assigned class labels and minimized the entropy
32 of the unselected points in the unlabeled pool. Guo also proposed a batch mode
33 active learning scheme which maximized the mutual information between the la-
34 beled and unlabeled sets and was independent of the classification model [24].
35 Chakraborty *et al.* [9] introduced an active matrix completion algorithm to select
36 the most informative queries to complete a low rank matrix. Researchers have also
37 explored theoretical properties of active learning and have established concrete
38 mathematical bounds on the expected number of queries to achieve a given
39 error rate [15].

40 While active learning has been extensively studied in a variety of computer vision
41 applications, it has been comparatively much less explored for object counting.
42 Loy *et al.* [37] proposed a regression based active learning algorithm (m-landmark)
43 for crowd counting, which was based on computing the normalized Graph Laplacian

1 L followed by k -means clustering. However, the algorithm did not consider binary
2 user feedback about the object count. The Elastic Net algorithm proposed by Tan
3 *et al.* [57] was based on a similar clustering strategy; however, it was more focused on
4 selecting a promising set of initial training samples for semi-supervised learning,
5 rather than active learning. In this paper, we propose a novel object counting algo-
6 rithm which can identify the exemplar unlabeled samples for manual annotation
7 and requires only binary feedback from human oracles. We now describe the pro-
8 posed framework.

9 **4.4. Proposed framework**

11 Let $\{x_{i1}, x_{i2}, \dots, x_{iN}\}$ be the set of N instances, which are labeled with their exact
12 counts $Y = \{y_1, y_2, \dots, y_N\}$ and let $\{x_{u1}, x_{u2}, \dots, x_{uM}\}$ be the set of the unlabeled
13 instances. Our objective is to select a batch containing k most informative unlabeled
14 samples from the unlabeled set, obtain their binary annotations from the human
15 oracle and use that to predict the labels of all the unlabeled samples. This task can be
16 decomposed into the following two research questions (**RQs**): (i) How can we use
17 active learning to select the k most informative samples from the unlabeled set for
18 binary user annotation? and (ii) Given the current labeled set containing the exact
19 counts and the set of k newly selected samples from the unlabeled set with binary
20 (*yes/no*) annotations, how can we predict the labels of all the unlabeled samples?

21 Conventional regression-based counting algorithms (such as the ridge regression
22 or the support vector regression) require the exact count of the number of objects in
23 each data sample and are hence unsuitable for our application. Given our problem
24 set-up, we need a framework which can incorporate inequality constraints (greater
25 or less than a given threshold) in estimating the count of objects in images. An
26 alternative strategy is to pose regression learning as the problem of completing a
27 low rank matrix [9]. Further, Marecek *et al.* [39] recently proposed a matrix com-
28 pletion algorithm under interval uncertainty, to impute the missing entries of a
29 data matrix in the presence of equality and inequality constraints. In this paper, we
30 exploit matrix completion algorithms for the problem of object counting from bi-
31 nary user feedback.

32 **4.4.1. Matrix completion**

33
34 The data collected in most computer vision/machine learning applications are
35 structured in the form of matrices. For instance, in a classification/regression
36 problem, each row represents a data sample, with corresponding label(s) and each
37 column denotes a feature; in a recommendation system, the data is represented in the
38 form of a matrix, where each row is a user, each column is an object and the cor-
39 responding entry represents the rating given by the particular user to that object.
40 Due to flaws in the feature acquisition process or the unwillingness of subjects to
41 disclose personal information, the collected data often contains missing entries,
42 which can bias results, reduce generalizability and lead to erroneous conclusions.
43

1 Matrix completion algorithms attempt to reconstruct a matrix from a set of partially
 2 observed entries and are of immense practical importance [7, 45]. Such techniques
 3 have also been exploited to address classification and regression problems [23]. The
 4 fundamental assumption is that the stacked matrix $Z = [Y^0; X^0]$ containing the
 5 label matrix Y^0 and the feature matrix X^0 is jointly low rank. The missing entries in
 6 the matrix correspond to the labels of the unlabeled samples and are estimated using
 7 matrix completion algorithms. It is posed as the following optimization:

$$\begin{aligned}
 & \min_Z \quad \text{rank}(Z) \\
 & \text{s.t.} \quad Z_{ij} = E_{ij}, \quad \forall i, j \in E
 \end{aligned} \tag{2}$$

11 where E is the set of the observed entries. Several methods have been devised to
 12 efficiently optimize this problem. The Fixed Point Continuation (FPC) method in
 13 particular, is an iterative algorithm consisting of a gradient step and a shrinkage step
 14 in each iteration with guaranteed monotonic convergence [23].

16 4.4.2. RQ1: Active sampling of the unlabeled data instances

17 Our object counting framework is based on the theory of matrix completion, nec-
 18 cessitating an active learning framework within the matrix completion paradigm.
 19 Chakraborty *et al.* [9] recently proposed the *Active Matrix Completion* algorithm to
 20 identify the missing entries in a partially observed matrix, which are the most in-
 21 formative to reconstruct the original matrix. The fundamental idea was to compute a
 22 measure of uncertainty of prediction of every missing entry in the incomplete data
 23 matrix; the top uncertain entries were then queried for manual annotation. Three
 24 strategies were presented to quantify the prediction uncertainty of each missing
 25 entry in the incomplete matrix: (i) *Conditional Gaussians*, which assumes that the
 26 set of missing entries conditioned on the set of observed entries follows a multivariate
 27 normal distribution; the mean and covariance matrix of the conditional distribution
 28 are computed from the given data and the diagonal elements of the covariance
 29 matrix quantifies the variance (uncertainty) associated with each imputation; (ii)
 30 *Query by Committee (QBC)*, which uses a committee of matrix completion algo-
 31 rithms to impute the missing entries and quantifies the prediction uncertainty of a
 32 particular entry as the level of disagreement among the committee; and (iii) *Com-
 33 mittee Stability*, which is similar to QBC and quantifies the prediction uncertainty
 34 using the regularity of predictions of a particular entry from an ensemble of
 35 predictors.

36 In this work, we used the QBC algorithm for active instance sampling due to its
 37 promising performance in matrix completion [9] and active learning in general, its
 38 strong theoretical properties [51] and ease of implementation. Specifically, a com-
 39 mittee of matrix completion algorithms were applied on the partially observed data
 40 matrix to impute the missing values. The variance of prediction (among the com-
 41 mittee members) of each missing entry was taken as a measure of uncertainty of that
 42 entry. The top k uncertain entries were then queried for manual annotation. We used
 43

1 the following three commonly used matrix completion algorithms as members of our
2 committee:

3 ***k*-NN:** The *k*-nearest neighbor algorithm identifies the *k* most similar features to
4 the current one with a missing value and uses the average of these *k* nearest
5 neighbors as an estimate for the missing entry [28].

6 **EM:** This method imputes the missing values using the Expectation Maximiza-
7 tion (EM) algorithm [28]. An iteration of the EM algorithm involves two steps. In the
8 E step, the mean and covariance matrix are estimated from the data matrix (with the
9 missing entries filled with zeros or estimates from the previous M step); in the M step,
10 the missing value of each data column is imputed with their conditional expectation
11 values based on the available entries and the estimated mean and covariance. The
12 mean and the covariance are re-estimated based on the newly completed matrix and
13 the process is iterated until convergence.

14 **SVD:** Singular value decomposition (SVD) is a standard method for matrix
15 completion based on low-rank approximation [28]. In this method, initial guesses are
16 first provided to the missing data values. SVD is then applied to obtain a low rank
17 approximation of the filled-in data matrix. The missing values are then updated
18 based on their corresponding values in the low rank estimation. SVD is applied to the
19 updated matrix again and the process is iterated until convergence.

21 4.4.3. *RQ2: Counting with binary user feedback*

22 In our framework, the user provides only binary (*yes/no*) annotations to the un-
23 labeled samples selected using active learning. This necessitates a matrix comple-
24 tion scheme that can handle inequality constraints apart from equality constraints,
25 as in Eq. (2). The *MACO* algorithm proposed by Marecek *et al.* [39] uses alternating
26 parallel co-ordinate descent to complete a matrix in the presence of equality, lower
27 bound and upper bound constraints. Let X be the $m \times n$ matrix to be recon-
28 structed. Suppose that for the elements $(i, j) \in E$, we have equality constraints, for
29 the elements $(i, j) \in L$ we have lower bounds and for the elements $(i, j) \in U$, we
30 have upper bounds. Completing the matrix can thus be posed as the following
31 optimization:
32

$$\begin{aligned}
 & \min_{X \in \mathbb{R}^{m \times n}} \quad \text{rank}(X) \\
 & \text{s.t.}: X_{ij} = X_{ij}^E, \quad \forall (i, j) \in E \\
 & \quad X_{ij} \geq X_{ij}^L, \quad \forall (i, j) \in L \\
 & \quad X_{ij} \leq X_{ij}^U, \quad \forall (i, j) \in U
 \end{aligned} \tag{3}$$

33
34
35
36
37
38
39 The problem in Eq. (3) is NP-hard, even with $U = L = \emptyset$ [39]. A popular heuristic
40 enforces low rank in a synthetic way by writing X as a product of two matrices,
41 $X = AB$, where $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$. Hence, X is of rank at most r . The alter-
42 nating parallel co-ordinate descent algorithm to solve the above optimization is
43 outlined in Algorithm 1 (please refer to [39] for more detailed derivations).

Algorithm 1. The MACO algorithm for matrix completion under equality and inequality constraints

Require: E, L, U, X^E, X^L, X^U , rank r

```

1: choose  $A \in \mathfrak{R}^{m \times r}$  and  $B \in \mathfrak{R}^{r \times n}$ 
2: for  $k = 1, 2, \dots$ , do
3:   choose a random subset  $\hat{S}_{row} \in \{1, 2, \dots, m\}$ 
4:   for  $i \in \hat{S}_{row}$  do
5:     choose  $\hat{r} \in \{1, 2, \dots, r\}$  uniformly at random
6:     update  $A_{i\hat{r}} = A_{i\hat{r}} + \delta_{i\hat{r}}$ , where  $\delta_{i\hat{r}}$  is computed using co-ordinate descent
7:   end for
8:   choose a random subset  $\hat{S}_{column} \in \{1, 2, \dots, n\}$ 
9:   for  $j \in \hat{S}_{column}$  do
10:    choose  $\hat{r} \in \{1, 2, \dots, r\}$  uniformly at random
11:    update  $B_{\hat{r}j} = B_{\hat{r}j} + \delta_{\hat{r}j}$ , where  $\delta_{\hat{r}j}$  is computed using co-ordinate descent
12:   end for
13: end for
14: return  $m \times n$  matrix  $AB$ 

```

In our object counting application, the initial training set containing the exact counts of the objects forms the set E . The MACO algorithm is used only with the set E to derive estimates of the missing labels of the unlabeled samples. These estimates are used as thresholds for binary user query; the binary user feedback on the selected unlabeled samples constitute the sets L and U . The MACO algorithm is used again with the sets E, L and U to estimate the missing labels of the unlabeled samples. The pseudo-code of our algorithm is presented in Algorithm 2.

4.5. Experiments and results

Datasets and Feature Extraction: We used four challenging datasets from different application domains to study the performance of the proposed framework: (i) the Mall dataset [37] containing video frames collected using a publicly accessible webcam for crowd counting and profiling research; (ii) the UCSD Pedestrian dataset [10], which contains videos of pedestrians on UCSD walkways, taken from a stationary camera; (iii) the Fudan Pedestrian dataset [57], which contains video frames captured at one side entrance of Guanghai Tower, Fudan University, Shanghai, China; and (iv) the TRAffic ANd COngestionS (TRANCOS) dataset [42], a novel benchmark for (extremely overlapping) vehicle counting in traffic congestion situations. All these datasets are captured under challenging real-world conditions with severe inter-object occlusions, varying crowd densities from sparse to crowded, as well as diverse activity patterns (static and moving crowds) under varying illumination conditions at different times of the day. Sample images from these datasets

Algorithm 2. The proposed active object counting algorithm with binary user feedback

Require: The labeled initial training set $\{x_{l1}, x_{l2}, \dots, x_{lN}\}$, their ground truth counts $Y = \{y_1, y_2, \dots, y_N\}$, the unlabeled set $\{x_{u1}, x_{u2}, \dots, x_{uM}\}$, batch size k , rank parameter r

- 1: Form the stacked matrix $Z = [Y; X]$; the labels of the unlabeled samples constitute the missing entries
 - 2: Form the equality set E from the given label set Y
 - 3: Apply the MACO algorithm (Algorithm 1) using the constraint set E to derive estimates of the missing labels of the unlabeled samples
 - 4: Use the QBC algorithm [9] on Z to select the k most informative unlabeled samples
 - 5: Query the binary labels of the selected k samples with respect to the thresholds given by their label estimates computed in Step 3
 - 6: Form the lower bound and upper bound constraint sets L and U from the binary user feedback
 - 7: Apply the MACO algorithm again using the constraint sets E , L and U to complete the missing entries of the matrix
 - 8: **return** The labels of the unlabeled samples
-

are shown in Figs. 6 and 7. The histogram of oriented gradients (HOG) feature [14] was used as the descriptor of each image frame due to its established performance in computer vision tasks.

4.5.1. *Experimental setup*

Each dataset was randomly divided into a labeled training set and an unlabeled set. The batch size was set at 10% of the dataset size (as detailed in Table 7). A batch of samples was queried from the unlabeled set, appended to the labeled set and the performance was evaluated on the complete unlabeled set. We studied the performance of binary annotations (where the user merely provides *yes/no* answers as to whether the number of objects in an image is greater than a given threshold) for both random selection as well as active sampling of the unlabeled samples. In random sampling, a batch of samples was selected at random from the unlabeled set

Table 7. Dataset details.

Dataset	Number of samples	Batch size
Mall	2000	200
UCSD	2000	200
Fudan	1500	150
TRANCOS	1200	120

1 for annotation while in active sampling, the proposed active learning framework
 2 was used to select the unlabeled samples for annotation. We also studied the
 3 performance of exact annotations (where the user provides the exact count of the
 4 number of objects in an image) for both random and active sampling. For exact
 5 annotations, we used only the equality constraint set E in the MACO algorithm;
 6 the lower and upper bound constraint sets L and U were empty since the user
 7 provided the exact counts. We used the mean-squared error (MSE) on the unlabeled
 8 set as the evaluation metric in our work. For each dataset, we studied the
 9 performance with different sizes of the initial training set, from 10% to 50% in steps
 10 of 10%.

12 4.5.2. Regression using matrix completion

13 Matrix completion algorithms have been used successfully to address regression
 14 problems [9]. We first studied the performance of matrix completion for the regression-based
 15 object counting problem. We used three common regression algorithms — ridge regression,
 16 kernelized ridge regression and support vector regression — as comparison baselines. The results
 17 on the four datasets are reported in Table 8 (in each experiment, 70% of the data was used for training
 18 and 30% for testing).

19 Thus, matrix completion provides comparable performance to other counting
 20 techniques. However, our method has the flexibility of incorporating binary user
 21 feedback in contrast to other methods which need the exact counts for model
 22 training.

25 4.5.3. Active counting with binary feedback

26 The results for the four datasets are reported in Tables 9–12. All the results were
 27 averaged over 5 runs (with different labeled and unlabeled sets) to rule out the effects
 28 of randomness. **BaseMSE** denotes the mean squared error using the current training
 29 data (before sample selection and annotation); **Binary Ann** denotes the MSE
 30 corresponding to binary annotations while **Exact Ann** denotes the MSE corresponding
 31 to exact annotations. **Random** denotes the case when the unlabeled samples
 32 are selected at random for user annotation while **Active** denotes the case
 33 when active sampling is used to select the unlabeled samples.

34 Table 8. Comparison of matrix completion (MC) against regression algorithms. Error metric:
 35 Mean squared error.

36 Dataset	37 MC	38 Ridge regression	39 Kernel ridge regression	40 Support vector regression
41 Fudan	2.09	1.51	3.0	1.54
42 Mall	9.26	4.60	7.20	4.69
43 TRANCOS	115.99	86.30	147.08	89.84
UCSD	30.01	6.54	7.13	6.78

Table 9. MSE comparison results on the Mall dataset. Lower values denote better performance.

Train %	BaseMSE	Binary Ann		Exact Ann	
		Random	Active	Random	Active
10	918.14	629.59	505.75	548.27	394.24
20	813.29	475.91	317.65	422.42	225.02
30	714.46	370.50	209.1	344.80	148.78
40	612.29	334.80	154.68	322.69	111.5
50	510.05	253.69	90.99	264.71	58.35

Table 10. MSE comparison results on the UCSD dataset. Lower values denote better performance.

Train %	BaseMSE	Binary Ann		Exact Ann	
		Random	Active	Random	Active
10	780.71	527.55	526.32	361.89	347.33
20	704.85	407.78	315.24	269.77	144.10
30	611.38	311.61	247.05	218.77	136.78
40	525.84	245.6	176.75	158.16	70.09
50	443.26	218.43	151.79	152.72	73.79

Table 11. MSE comparison results on the Fudan dataset. Lower values denote better performance.

Train %	BaseMSE	Binary Ann		Exact Ann	
		Random	Active	Random	Active
10	36.21	30.79	28.43	28.51	24
20	33.43	27.92	19.81	22.41	15.84
30	31.89	19.11	12.61	15.85	10.99
40	28.72	17.64	9.8	14.62	8.21
50	24.23	16.67	7.86	14.93	7.04

Table 12. MSE comparison results on the TRANCOS dataset. Lower values denote better performance.

Train %	BaseMSE	Binary Ann		Exact Ann	
		Random	Active	Random	Active
10	1.3 e+03	984.14	911.94	861.94	829.24
20	1.23 e+03	685.56	640.86	574.59	515.72
30	1.1 e+03	548.57	483.86	440.63	371.47
40	952.78	394.25	344.82	313.92	271.21
50	794.23	296	266.1	227.15	197.07

1 We first note that the MSE reduces with increasing size of the initial training
2 set, which is intuitive. We also note that the algorithm based on binary annota-
3 tions delivers much better performance compared to the baseline error. This
4 corroborates the usefulness of the proposed framework in tremendously reducing
5 the error rate by exploiting only binary feedback from the human user. Moreover,
6 active sampling successfully identifies the salient and exemplar unlabeled instances
7 and further improves the error rate over random selection in a binary user feed-
8 back setting. The same pattern is evident for different sizes of the initial training
9 set and for all the datasets, depicting the generalizability of our framework.
10 Thus, while conventional learning frameworks can operate only with data anno-
11 tated with the exact counts of objects, our framework offers more flexibility and
12 ease of interaction between the user and the machine. From these results, we
13 conclude that the proposed framework can be immensely useful to boost the ac-
14 curacy of an object counting system while minimizing the labeling burden on
15 human oracles.

16 The algorithm based on exact annotations produces better performance compared
17 to that based on binary annotations. This is intuitive, as exact annotation provides
18 more information to the underlying machine learning models. As before, active in-
19 stance sampling further reduces the error rate compared to random sampling. More
20 importantly, we note that active sampling with binary user annotations often pro-
21 vides comparable results (and sometimes, even outperforms) random sampling with
22 exact annotations, which is the conventional method to address the counting
23 problem. This depicts the merit of our algorithm in tremendously reducing human
24 annotation effort with minimal effect on the counting accuracy.

25 26 4.5.4. *Threshold study*

27 In our framework, a threshold is first computed by the algorithm and the user
28 provides a binary feedback as to whether the number of objects in the image is
29 greater or less than the threshold (the threshold is computed as the current label
30 estimate of the sample in question). Thus, the user annotation time depends on the
31 threshold computed by the MACO algorithm. If the threshold is close to the actual
32 count of objects, the annotation time will be higher and vice versa. In this experi-
33 ment, we study the thresholds computed by our algorithm on 50 random unlabeled
34 samples for 10% and 50% initial labeled training data. The results on the Mall and
35 TRANCOS datasets are depicted in Fig. 8.

36 We note that with 10% labeled training data, the thresholds computed are coarse
37 and thus, the binary annotation time will be low. As the percentage of training data
38 increases, the prediction accuracy increases and consequently, the computed
39 thresholds are much closer to the actual counts. Hence, the binary annotation time
40 will be almost similar to the absolute annotation time, since an exhaustive count of
41 all the objects will be necessary for accurate annotations. Our framework is therefore
42 most useful in the initial stages of learning, when the amount of labeled training data
43

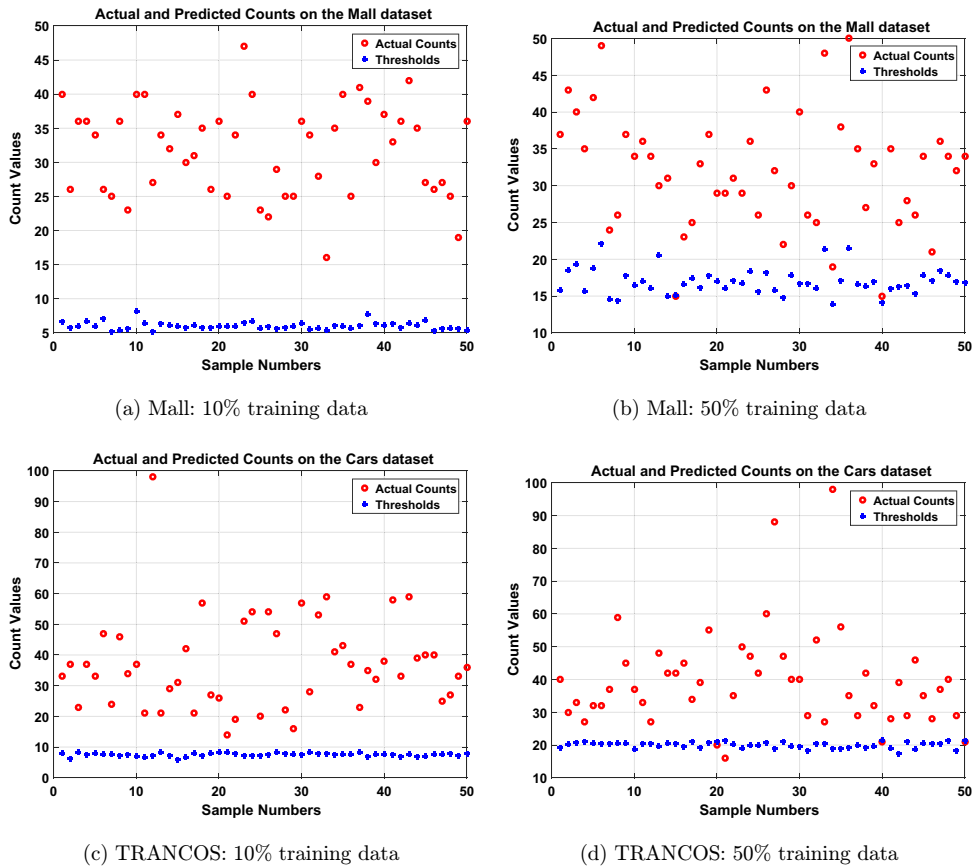


Fig. 8. Study of the threshold computed in our algorithm. Best viewed in color.

is scarce. With abundant labeled data, absolute annotations are more advantageous. A hybrid framework can be envisioned, where binary annotations are used in the initial stages and absolute annotations in the later stages of learning. This will be investigated as part of future research.

5. Discussion and Future Work

Three fan enrichment projects in the scope of the Smart Stadium for Smarter Living initiative were presented. These projects targeted improved safety (Crowd Understanding), fan engagement (Athletic Demonstrator Platform), and efficiency/convenience (Wait Time/Queue Estimation). Through use of smart stadia as testbeds, the manageable size and heterogeneity of these testbeds enabled practical trials while still providing a useful environment to explore challenges of scalability and real-timeness. Preliminary results presented here demonstrate the potential of these technologies for smart city solutions.

1 As part of future work, we are developing and deploying new projects for the
2 smart stadia. These projects include smart solutions to address issues of congestion
3 and difficulty parking during large events at the stadia; game-within-a-game halftime
4 interactions to enrich fan engagement; and projects that target important priority
5 areas of energy efficiency and sustainability. We are also investigating the use of the
6 Athletic Demonstrator Platform as a low-cost, accurate platform to augment tra-
7 ditional athlete training intended for use outside of sessions involving the trainer or
8 coach. Moreover, future studies with the platform are being planned to observe how
9 this feedback can be integrated over time to adapt to a user's proficiency level, and
10 how this integration can differ between individuals and various types of movement.
11 One such study will investigate how multimodal feedback delivery may be tuned for
12 fast sports motion interaction as opposed to slower rehabilitative movements, and
13 how the type of movement may be inferred from the nature of the expert data
14 samples.

15 Acknowledgments

16 The authors thank Intel Corporation, National Science Foundation, Arizona State
17 University, and Dublin City University for their funding support. This material is
18 partially based on work supported by: Intel Corporation under grant Joint Path
19 Finding (JPF) Proposal: Smart Stadium and Smart Living Research; and National
20 Science Foundation under Grant No. 1069125.

21 References

- 22 [1] N. Ahuja and S. Todorovic, Extracting texels in 2.1d natural textures, in *IEEE Inter-*
23 *national Conference on Computer Vision*, 2007.
- 24 [2] A. Ahmadi, E. Mitchell, F. Destelle, M. Gowing, N. E. O'Connor, C. Richter and K.
25 Moran, Automatic activity classification and movement assessment during a sports
26 training session using wearable inertial sensors, in *Proc. 11th International Conference*
27 *on Wearable and Implantable Body Sensor Networks*, 2014, pp. 98–103.
- 28 [3] A. Ahmadi, F. Destelle, D. Monaghan, K. Moran, N. E. O'Connor, L. Unzueta and M. T.
29 Linaza, Human gait monitoring using body-worn inertial sensors and kinematic model-
30 ling, in *Proc. IEEE SENSORS*, 2015, pp. 1–4.
- 31 [4] A. E. Atwater, Biomechanics of overarm throwing movements and of throwing injuries,
32 *Exercise and Sport Sciences Reviews* **7**(1) (1979) 43–86.
- 33 [5] O. Barinova, V. Lempitsky and P. Kohli, On the detection of multiple object instances
34 using hough transforms, in *IEEE Conference on Computer Vision and Pattern Recog-*
35 *nition*, 2010
- 36 [6] K. Brinker, Incorporating diversity in active learning with support vector machines, in
37 *International Conference on Machine Learning*, 2003.
- 38 [7] E. Candes and T. Tao, The power of convex relaxation: Near-optimal matrix completion,
39 in *IEEE Transactions on Information Theory* **56**(5) (2010) 2053–2080.
- 40 [8] S. Chakraborty, V. Balasubramanian, Q. Sun, S. Panchanathan and J. Ye, Active batch
41 selection via convex relaxations with guaranteed solution bounds, in *IEEE Transactions*
42 *on Pattern Analysis and Machine Intelligence* **37**(10) (2015) 1945–1958.
- 43

- 1 [9] S. Chakraborty, J. Zhou, V. Balasubramanian, S. Panchanathan, I. Davidson and J. Ye,
2 Active matrix completion, in *IEEE International Conference on Data Mining*, 2013.
- 3 [10] A. Chan, Z. Liang and N. Vasconcelos, Privacy preserving crowd monitoring: Counting
4 people without people models or tracking, in *IEEE Conference on Computer Vision and
5 Pattern Recognition*, 2008.
- 6 [11] D. K. Chan, C. Lonsdale, P. Y. Ho, P. S. Yung and K. M. Chan, Patient motivation and
7 adherence to postsurgery rehabilitation exercise recommendations: The influence of
8 physiotherapists autonomy-supportive behaviors, *Arch Phys Med Rehabil* **90**(12) (2009)
9 1977–1982.
- 10 [12] K. Chen, C. Loy, S. Gong and T. Xiang, Feature mining for localized crowd counting, in
11 *British Machine Vision Conference*, 2012.
- 12 [13] S. Cho, T. Chow and C. Leung, A neural-based crowd estimation by hybrid global
13 learning algorithm, in *IEEE Transactions on Systems, Man and Cybernetics* **29**(4)
14 (1999) 535–541.
- 15 [14] N. Dalal and B. Triggs, Histograms of oriented gradients for human detection, in *IEEE
16 Conference on Computer Vision and Pattern Recognition*, 2005.
- 17 [15] S. Dasgupta, Coarse sample complexity bounds for active learning, in *Advances of
18 Neural Information Processing Systems*, 2005.
- 19 [16] C. Desai, D. Ramanan and C. Fowlkes, Discriminative models for multi-class object
20 layout, in *IEEE International Conference on Computer Vision*, 2009.
- 21 [17] F. Destelle, A. Ahmadi, N. E. O’Connor, K. Moran, A. Chatzitofis, D. Zarpalas and P.
22 Daras, Low-cost accurate skeleton tracking based on fusion of Kinect and
23 wearable inertial sensors, in *Proc. 22nd European Signal Processing Conference*, 2014,
24 pp. 371–375.
- 25 [18] F. Destelle, A. Ahmadi, K. Moran, N. E. O’Connor, N. Zioulis, A. Chatzitofis, D. Zar-
26 palas, P. Daras, L. Unzueta, J. Goenexea and M. Rodriguez, A multi-modal 3D cap-
27 turing platform for learning and preservation of traditional sports and games, in *Proc.
28 23rd ACM International Conference on Multimedia*, 2015, pp. 747–748.
- 29 [19] J. E. Deutsch, M. Borbely, J. Filler, K. Huhn and P. Guarrera-Bowlby, Use of a low-Cost,
30 commercially available gaming console (Wii) for rehabilitation of an adolescent with
31 cerebral palsy, *Phys. Ther.* **88**(10) (2008) 1196–1207.
- 32 [20] S. Essid, D. Alexiadis, R. Tournemenne, M. Gowing, P. Kelly, D. Monaghan, P. Daras,
33 A. Drmeau and N. E. O’Connor, An advanced virtual dance performance evaluator, in
34 *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing*, 2012,
35 pp. 2269–2272.
- 36 [21] E. Farella, L. Benini, B. Ricc and A. Acquaviva, MOCA: A low-power, low-cost
37 motion capture system Based on integrated accelerometers, *Adv. MultiMedia* **2007**(1)
38 (2007) 11.
- 39 [22] C. Fu and Y. Yang, A batch-mode active learning SVM method based on semi-super-
40 vised clustering, in *Intelligent Data Analysis*, 2015.
- 41 [23] A. Goldberg, X. Zhu, B. Recht, J. Xu and R. Nowak, Transduction with matrix com-
42 pletion: Three birds with one stone, in *Advances of Neural Information Processing
43 Systems*, 2010.
- 44 [24] Y. Guo, Active instance sampling via matrix partition, in *Advances of Neural Infor-
45 mation Processing Systems*, 2010.
- 46 [25] Y. Guo and D. Schuurmans, Discriminative batch mode active learning, in *Advances of
47 Neural Information Processing Systems*, 2007.
- 48 [26] I. Haritaoglu, D. Harwood and L. Davis, W4: real-time surveillance of people and their
49 activities, in *IEEE Transactions on Pattern Analysis and Machine Intelligence* **22**(8)
50 (2000) 809–830.

- 1 [27] T. Hassner, Y. Itcher and O. Kliper-Gross, Violent flows: Real-time detection of violent
2 crowd behavior, in *Proc. IEEE Computer Society Conference on Computer Vision and*
3 *Pattern Recognition Workshops*, 2012, p. 16.
- 4 [28] T. Hastie, R. Tibshirani, G. Sherlock, M. Eisen, P. Brown and D. Botstein, Imputing
5 missing data for gene expression arrays, in Technical Report, Stanford University, 1999.
- 6 [29] S. Hodgkinson, Is your city smart enough? Digitally enabled cities and societies will
7 enhance economic, social, and environmental sustainability in the urban century, in
8 *OVUM Report*, 2011.
- 9 [30] S. Hoi, R. Jin, J. Zhu and M. Lyu, Batch mode active learning and its application to
10 medical image classification, in *International Conference on Machine Learning*, 2006.
- 11 [31] S. Hoi, R. Jin, J. Zhu and M. Lyu, Semi-supervised SVM batch mode active learning
12 for image retrieval, in *IEEE Conference on Computer Vision and Pattern Recognition*,
13 2008.
- 14 [32] H. Idrees, I. Saleemi, C. Seibert and M. Shah, Multi-source multi-scale counting in
15 extremely dense crowd images, in *IEEE Conference on Computer Vision and Pattern*
16 *Recognition*, 2013.
- 17 [33] P. K. Tra, K. Pong and R. Bowden, An improved adaptive background mixture model
18 for real-time tracking with shadow detection, in *Video-Based Surveillance Systems*,
19 eds. P. Remagnino, G. A. Jones, N. Paragios and C. S. Regazzoni (Springer, 2002),
20 pp. 135–144.
- 21 [34] D. Kong, D. Gray and H. Tao, Counting pedestrians in crowds using viewpoint invariant
22 training, in *British Machine Vision Conference*, 2005.
- 23 [35] V. Lempitsky and A. Zisserman, Learning to count objects in images, in *Neural Infor-*
24 *mation Processing Systems*, 2010.
- 25 [36] H. Liu, X. Wei, J. Chai, I. Ha and T. Rhee, Realtime human motion control with a small
26 number of inertial sensors, in *Proc. Symposium on Interactive 3D Graphics and Games*,
27 2011, pp. 133–140.
- 28 [37] C. Loy, S. Gong and T. Xiang, From semi-supervised to transfer counting of crowds, in
29 *IEEE International Conference on Computer Vision*, 2013.
- 30 [38] A. Marana, S. Velastin, L. Costa and R. Lotufo, Estimation of crowd density using image
31 processing, in *Image Processing for Security Applications*, 1997.
- 32 [39] J. Marecek, P. Richtarik and M. Takac, Matrix completion under interval uncertainty, in
33 *European Journal of Operational Research* **256**(1) (2017) 35–43.
- 34 [40] M. Marsden, K. McGuinness, S. Little and N. E. OConnor, Holistic features for real-time
35 crowd behaviour anomaly detection, in *Proc. IEEE International Conference on Image*
36 *Processing*, 2016, pp. 918–922.
- 37 [41] U. Nations, World urbanization prospects: The 2014 revision, highlights. Department of
38 Economic and Social Affairs, in *Population Division, United Nations*, 2014.
- 39 [42] R. Olmedo, B. Jimnez, R. Sastre, S. Bascn and D. Rubio, Extremely overlapping vehicle
40 counting, in *Iberian Conference on Pattern Recognition and Image Analysis*, 2015.
- 41 [43] N. Otsu, A threshold selection method from gray-level histograms, *Automatica* **11**
42 (285296) (1975) 2327.
- 43 [44] V. Rabaud and S. Belongie, Counting crowded moving objects, in *IEEE Conference on*
Computer Vision and Pattern Recognition, 2006.
- [45] B. Recht, A simpler approach to matrix completion, in *Journal of Machine Learning*
Research **12** (2011) 3413–3430.
- [46] D. J. Reinkensmeyer and S. J. Housman, “If I cant do it once, why do it a hundred
times?”: Connecting volition to movement success in a virtual environment motivates
people to exercise the arm after stroke, in *Proc. Virtual Rehabilitation 2007*, 2007,
pp. 44–48.

- 1 [47] D. Roetenberg, H. Luinge and P. Slycke, Xsens MVN: Full 6DOF human motion
2 tracking using miniature inertial sensors, in *Xsens Motion Technologies BV*, Technical
3 Report, 2009.
- 4 [48] G. Schwarz, Estimating the dimension of a model, *The Annals of Statistics* **6**(2) (1978)
5 461–464.
- 6 [49] S. Segui, O. Pujol and J. Vitria, Learning to count with deep object features, in *IEEE*
7 *Conference Computer Vision and Pattern Recognition Workshop*, 2015.
- 8 [50] B. Settles, Active learning literature survey, in Technical Report 1648, University of
9 Wisconsin-Madison, 2010.
- 10 [51] H. Seung, M. Opper and H. Sompolinsky, Query by committee, in *Workshop on*
11 *Computational Learning Theory*, 1992.
- 12 [52] J. Shao, K. Kang, C. C. Loy and X. Wang, Deeply learned attributes for crowded scene
13 understanding, in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*,
14 2015, pp. 4657–4666.
- 15 [53] R. Sigrist, G. Rauter, R. Riener and P. Wolf, Augmented visual, auditory, haptic,
16 and multimodal feedback in motor learning: A review, *Psychon Bull. Rev.* **20**(1) (2013)
17 21–53.
- 18 [54] N. Skjret, A. Nawaz, T. Morat, D. Schoene, J. L. Helbostad and B. Vereijken, Exercise
19 and rehabilitation delivered through exergames in older adults: An integrative review of
20 technologies, safety and efficacy, *International Journal of Medical Informatics* **85**(1)
21 (2016) 1–16.
- 22 [55] R. Tadayon, T. McDaniel, M. Goldberg, P. M. Robles-Franco, J. Zia, M. Laff, M. Geng
23 and S. Panchanathan, Interactive motor learning with the autonomous training assist-
24 ant: A case study, in *Human-Computer Interaction: Interaction Technologies* (Springer
25 International Publishing, 2015), pp. 495–506.
- 26 [56] R. Tadayon, T. McDaniel and S. Panchanathan, Autonomous training assistant: A
27 system and framework for guided at-home motor learning, in *Proc. 18th International*
28 *ACM SIGACCESS Conference on Computers and Accessibility*, 2016, pp. 293–294.
- 29 [57] B. Tan, J. Zhang and L. Wang, Semi-supervised elastic net for pedestrian counting, in
30 *Pattern Recognition* **44**(10–11) (2011) 2297–2304.
- 31 [58] C. Tomasi and T. Kanade, Detection and tracking of point features, in Technical Report
32 CMU-CS-91-132, 1991.
- 33 [59] Smart cities council — Transportation [Online]. Available: [http://readinessguide.
34 smartcitiescouncil.com/readiness-guide/transportation-0](http://readinessguide.smartcitiescouncil.com/readiness-guide/transportation-0).
- 35 [60] G. Welch and E. Foxlin, Motion tracking: No silver bullet, but a respectable arsenal,
36 *IEEE Computer Graphics and Applications* **22**(6) (2002) 24–38.
- 37 [61] M. Windolf, N. Gtzen and M. Morlock, Systematic accuracy and precision analysis of
38 video motion capturing systems exemplified on the Vicon-460 system, *Journal of Bio-*
39 *mechanics* **41**(12) (2008) 2776–2780.
- 40 [62] Z. Zhang, Microsoft Kinect sensor and its effect, *IEEE MultiMedia* **19**(2) (2012) 4–10.
- 41 [63] L. Zhang, J. C. Hsieh, T. T. Ting, Y. C. Huang, Y. C. Ho and L. K. Ku, A Kinect based
42 golf swing score and grade system using GMM and SVM, in *Proc. 5th International*
43 *Congress on Image and Signal Processing*, 2012, pp. 711–715.