

Large-Scale Analysis of Soccer Matches using Spatiotemporal Tracking Data

Alina Bialkowski^{1,2}, Patrick Lucey¹, Peter Carr¹, Yisong Yue¹, Sridha Sridharan² and Iain Matthews¹

¹Disney Research, Pittsburgh, PA, USA, ²Queensland University of Technology, Brisbane, Australia

a.bialkowski@connect.qut.edu.au, {patrick.lucey,peter.carr,yisong.yue,iainm}@disneyresearch.com, s.sridharan@qut.edu.au

Abstract—Although the collection of player and ball tracking data is fast becoming the norm in professional sports, large-scale mining of such spatiotemporal data has yet to surface. In this paper, given an entire season’s worth of player and ball tracking data from a professional soccer league ($\approx 400,000,000$ data points), we present a method which can conduct both individual player and team analysis. Due to the dynamic, continuous and multi-player nature of team sports like soccer, a major issue is aligning player positions over time. We present a “role-based” representation that dynamically updates each player’s relative role at each frame and demonstrate how this captures the short-term context to enable both individual player and team analysis. We discover role directly from data by utilizing a minimum entropy data partitioning method and show how this can be used to accurately detect and visualize formations, as well as analyze individual player behavior.

I. INTRODUCTION

Coinciding with the widespread deployment of tracking technologies, an enormous amount of trajectory data logging movements of people, vehicles, animals and weather patterns has emerged for various applications such as transportation [1], military [2], social [3], scientific studies [4] and hurricane prediction [5]. Another interesting application which has seen a recent deluge of spatiotemporal tracking data is the domain of sports analytics, where vision-based tracking systems have been deployed in professional basketball [6], soccer [7], baseball [8] tennis and cricket [9]. However, even though an enormous amount of data is being generated for visualization and umpiring consideration, no methods for large-scale analysis of the tracking data have surfaced yet.

In team sports, there are two types of analysis that can be conducted: a) individual and b) team analysis. In terms of analyzing individual player behavior, current methods plot locations of a particular player for an event or their mean position over time [10]. However, such methods lack important contextual information with regards to their team-mates. For example in soccer (see Fig. 1(a)), given we have a player who starts on the left-wing but then half way through the half he switches to the right wing, we get two distinct types of behaviors (i.e., left and right wing play). Current analysis conducts the analysis based on his original position or “role” (i.e., left-wing) which makes comparisons challenging. Ideally, we want a contextual label noting the player’s role at that specific moment (and not normal/starting position). To conduct role-specific analysis, we first need to define the set of roles. A role within a team can be defined as *a space or area that each player is assigned responsibility for relative to the other teammates*. The overall team formation therefore, can be defined as a set of roles, and these set of roles can vary

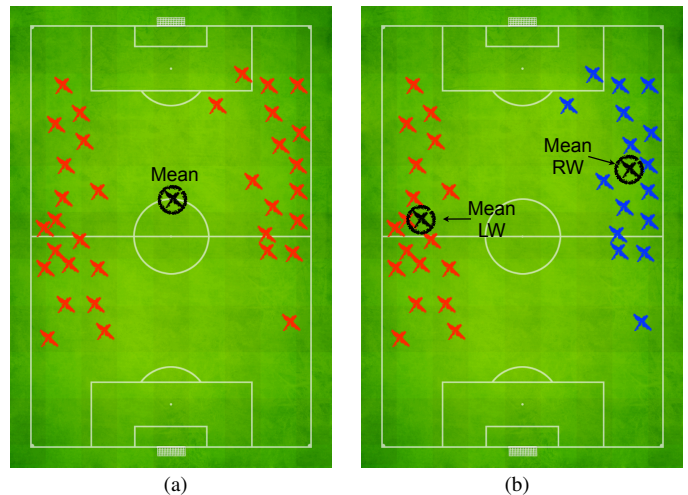


Fig. 1. (a) Shows the touches of the player who starts in the left-wing position but changes half way through the half to the right-wing. Current approaches just give the mean position which neglects the context. (b) In this paper, we utilize a role-representation which captures the context to allow for individual player analysis.

depending on the tactics/strategy of the coach as well as the behavior of the adversarial (i.e., opposing) team. In an information theoretical sense, this can be seen as partitioning the set of player positions (i.e., input data) into clusters which minimize the overlap or entropy. This is called *minimum entropy data partitioning* and in this paper we show that we can learn the roles directly from data. We show that this approach effectively aligns the multi-agent tracking data and allows for the detection and visualization of formations, as well as providing contextual information to do individual player analysis depending on their specific role (Fig. 1(b)).

A. Problem Definition

The simplest method of representing team behavior from tracking data is to use player identity. This means that given the x, y position of every player in the team, we initialize the players into some canonical ordering and these players remain *fixed* in this order throughout the entire match. So given $N = 11$ players, the representation at frame t can be given as $\mathbf{x}_t = [x_1, y_1, x_2, y_2, \dots, x_{11}, y_{11}]^T$ where the subscripts 1-11 refer to the unique identifier of the player such as jersey or player name. However, as seen in Fig. 2, the *static* assumption is not ideal as there is constant interchanging of positions making the dimensionality of the resulting subspace much higher [11]. Additionally, this assumption breaks down when there is a substitution, an expulsion of a player due to misconduct (i.e., red-card in soccer), or when comparing differ-

III. ROLE DISCOVERY USING MINIMUM ENTROPY DATA PARTITIONING

A. Minimum Entropy Data Partitioning

Given player tracking data \mathbf{D} , our goal is to estimate the underlying formation of the team, which is equivalent to finding the most probable set \mathcal{F}^* of 2D probability density functions $\mathcal{F}^* = \arg \max_{\mathcal{F}} P(\mathcal{F}|\mathbf{D})$. We begin by considering the 2D probability density function $P(\mathbf{X} = \mathbf{x})$ which models the tracking data \mathbf{D} . In other words, $P(\mathbf{x})$ represents the heat-map for an entire team. We can model the heat-map of the entire team as a linear combination of the heat maps for each role

$$\begin{aligned} P(\mathbf{x}) &= \sum_{n=1}^N P(\mathbf{x}|n)P(n) \\ &= \frac{1}{N} \sum_{n=1}^N P_n(\mathbf{x}). \end{aligned} \quad (1)$$

Strategically, a team needs to spread out its players so that the entire field is adequately covered. As a result, the probability density functions should exhibit minimal overlap (See Fig. 4 right). Equivalently, each role probability density function should exhibit minimal overlap with the team's probability density function. Following the ideas of minimum entropy data partitioning [32], [33], we employ Kullback-Lieber divergence to measure the overlap between two probability functions $P(x)$ and $Q(x)$

$$KL(P(x)||Q(x)) = \int P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx. \quad (2)$$

Since divergence is a strictly positive quantity (and completely overlapping probability density functions have zero divergence), we employ a penalty V_n based on the negative divergence value between the heat map $P_n(\mathbf{x})$ of an individual role and that of the team $P(\mathbf{x})$

$$V_n = -KL(P_n(\mathbf{x})||P(\mathbf{x})). \quad (3)$$

Computing the optimal formation \mathcal{F}^* is equivalent to determining the optimal set $\mathcal{F}^* = \{P_1(\mathbf{x}), \dots, P_N(\mathbf{x})\}^*$ of per-role probability density functions $P_n(\mathbf{x})$ that minimize the total overlap

$$\mathcal{F}^* = \arg \max_{\mathcal{F}} V. \quad (4)$$

Substituting the expressions for KL divergence into the total overlap cost illustrates the dependence on each role-specific 2D probability density function

$$\begin{aligned} V &= - \sum_{n=1}^N P(n) \int P(\mathbf{x}|n) \log P(\mathbf{x}|n) dx \\ &\quad + \sum_{n=1}^N P(n) \int P(\mathbf{x}|n) \log P(\mathbf{x}) dx. \end{aligned} \quad (5)$$

The expression for V is drastically simplified when put in terms of entropy

$$H(x) = - \int_{-\infty}^{+\infty} P(x) \log(P(x)) dx. \quad (6)$$

The total overlap cost, in terms of entropy, becomes

$$V = -H(\mathbf{x}) + \sum_{n=1}^N P(n)H(\mathbf{x}|n) \quad (7)$$

$$= -H(\mathbf{x}) + \frac{1}{N} \sum_{n=1}^N H(\mathbf{x}|n). \quad (8)$$

Substituting 8 into 4 and ignoring the constant term $H(\mathbf{x})$, the optimal formation is the set of role-specific probability density functions that minimize the total entropy

$$\mathcal{F}^* = \arg \min_{\mathcal{F}} \sum_{n=1}^N H(\mathbf{x}|n). \quad (9)$$

B. Equivalence to K-Means

As there is no way to solve this problem efficiently, we achieve an approximate solution using the expectation maximization (EM) algorithm [34]. Our approach is similar to k-means clustering. However, instead of assigning each data point to its closest cluster, we solve a linear assignment problem between identities and roles using the Hungarian algorithm [35]. Our process is the following. We arbitrarily assign each player a unique role label and assume it remains constant for the entire duration of the tracking data. As a result, each role probability density function is initialized using the tracking data of an individual player. The initial occupancy maps for each role resemble something similar to Fig. 4 (initial roles, left). We then iterate through each frame of the tracking data and assign role labels to player positions by formulating a cost matrix based on the log probability of each position being assigned a particular role label. We use the Hungarian algorithm [35] to compute the optimal assignment of role labels. Once role labels have been assigned to all frames of the tracking data, we recompute the probability density functions of each role. The process is repeated until convergence, resulting in well separated probability density functions similar to Fig. 4 (final, right). We normalize the tracking data in each frame to have zero mean in order to negate the effects of translation.

IV. DISCOVERING AND VISUALIZING ROLE

A. Spatiotemporal Tracking Data

For this work, we utilized an entire season of soccer player tracking data from Prozone. The data consisted of 20 teams who played home and away, totaling 38 games for each team or 380 games overall. Five of these games were omitted due to errors in the data files. We refer to the 20 teams using arbitrary labels $\{A, B, \dots, T\}$. Each game consists of two halves, with each half containing the (x, y) position of every player at 10 frames-per-second. This results in over 1 million data-points per game, in addition to the ball events that occurred throughout the match, consisting of 43 possible events (e.g. passes, shots, crosses, tackles etc.). Each of these ball events contains the time-stamp, location and players involved. An inventory of the data is given in Table I.

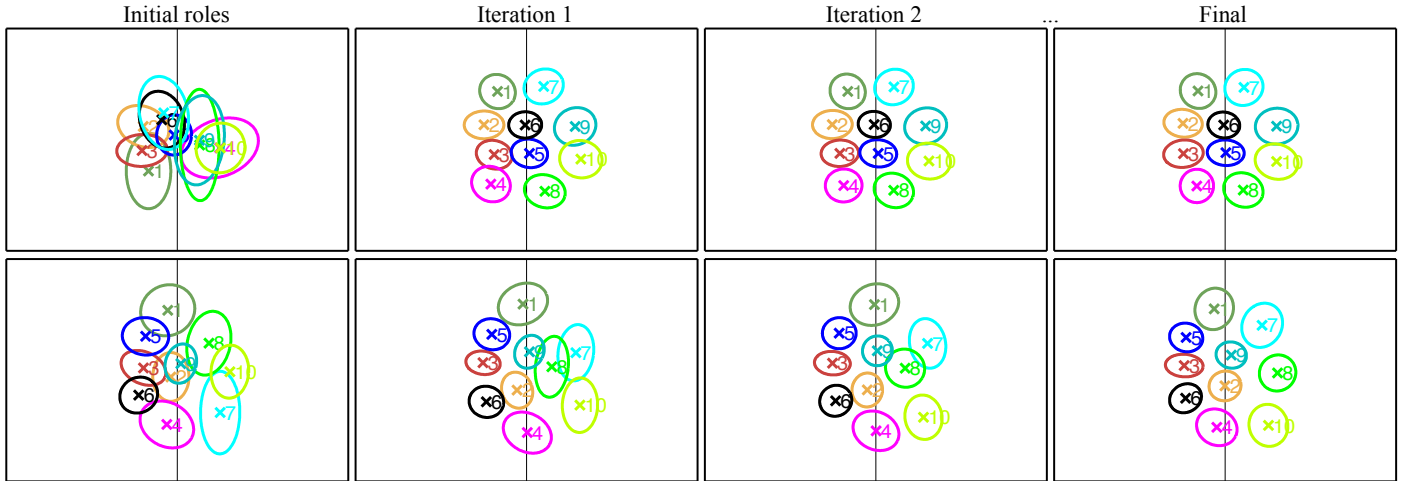


Fig. 4. Example of how our approach works at each iteration, with the role distributions drawn as Gaussians, with the mean and covariance shown.

B. Formation Detection

By using the expectation maximization procedure for minimum entropy data partitioning described in Section III-B, we simultaneously assign each player to a role at each frame of the tracking data and determine the role probability distributions, $P_n(\mathbf{x})$. We performed this procedure for each team and match half, excluding formations where players were sent off, resulting in the detection of 1411 formations. Each formation consists of a set of ten distinct role probability distributions representing the structural arrangement of the team over a half, and depicts the long-term characteristic behavior of the team.

Given these role distributions, we then automatically discovered different types of formations. We employed agglomerative clustering to group the discovered formations into clusters, using the *Earth Mover's Distance* (EMD) [36] to compute the distance between two role probability densities. The $\text{EMD}(\mathbf{a}, \mathbf{b})$ between two normalized histograms \mathbf{a} and \mathbf{b} is obtained as the solution of the transportation problem

$$\min_{f_{qt} \geq 0} \sum_{q,t=1}^D d_{qt} f_{qt} \quad \text{s.t.} \quad \sum_{q=1}^D f_{qt} = a^t, \sum_{t=1}^D f_{qt} = b^q. \quad (10)$$

The variable f_{qt} denotes a flow representing the amount transported from the q th supply to the t th demand and d_{qt} the ground distance. Using the EMD measure gave us role-to-role comparison distances, and then we set the distance between one formation and another equal to the sum of the distances between corresponding roles.

The resulting clusters are shown in Fig. 5, with the mean role positions of each formation overlaid over one another. It can be seen that clustering resulted in the discovery of distinct

Statistic	Frequency
Teams	20
Games	375
Data Points	480M
Ball Events	981K

TABLE I. INVENTORY OF DATASET USED FOR THIS WORK.

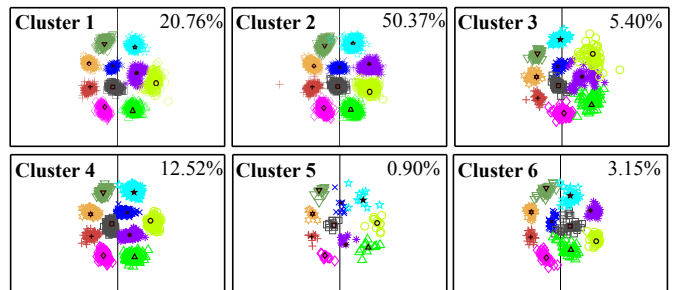


Fig. 5. Formation clustering results displaying the mean role positions of each formation assigned to each cluster, and the median formation overlaid in grey. (Note: all formations are normalized so that the team is attacking from left to right)

formation classes - e.g. Cluster 1 and 4 have only 1 striker in the front, Cluster 2 and 6 have 2 strikers, while Cluster 3 and 5 appear to have 3. Cluster 3 is the only cluster with 3 defenders at the back with the remainder all having 4. The clustering also gives an indication of which formations are more commonly adopted by teams, as given by the clustering assignment frequency (top right of each cluster in Fig. 5). We can see that Cluster 2, which appears to be a 4-4-2, is the most common with approximately 50.37% of formations being assigned to this cluster, followed by Cluster 1 (20.76%), which appears to be a 4-2-3-1. These give insight into the strategies adopted by teams (e.g. having 2 strikers instead of 1 may be considered a more attacking strategy).

To evaluate the clustering results, we compare against ground truth formation labels, where a soccer expert annotated the most frequently observed formation for each half and team according to the arrangement of players (4-4-2, 4-2-3-1, 4-3-3, 3-4-3, 4-1-4-1, or ‘other’ where the team either did not display a dominant formation or was not one of the given labels). To evaluate the results, we estimated the label of each cluster as the most frequent ground truth label within the cluster. The results are presented as a confusion matrix in Fig. 6.

It can be seen from Fig. 6 that the discovered formation clusters match the ground truth annotations quite well, with

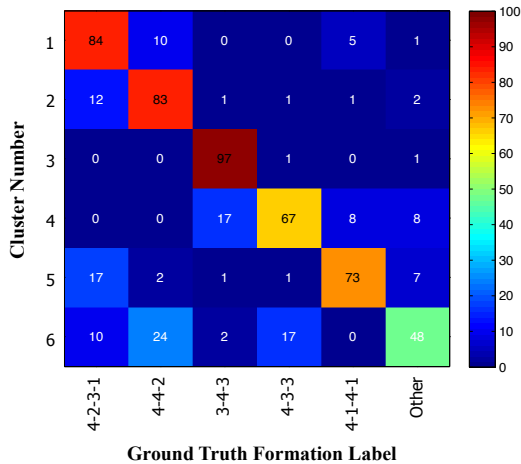


Fig. 6. Confusion matrix of formation clustering results relative to ground truth formation labels.

high within cluster label agreement and an overall correct classification rate of 75.33%. The most confusion is in Cluster 6 which appears to be a 4-1-3-2 formation (sometimes referred to as a 4-4-2 ‘diamond’), often being classified as a 4-4-2 and 4-3-3. On visual inspection of the misclassified examples, sometimes the formation appears in between two clusters, e.g. there is some confusion between the 4-4-2 and 4-2-3-1 formations when the second striker is positioned slightly behind the other.

C. Formation Visualization

In addition to representing the long-term behavior of the team in terms of formation or team structure, our method can also be used over shorter durations, to dynamically represent how a team plays throughout a match. Compared to existing statistics which only contain sparse team information (e.g. # corners, # shots, % possession), our approach can represent the spatiotemporal characteristics of the match in terms of formations and position.

One of the statistics which broadcasters present during a live-broadcast is the possession duration of both teams over the past 5 minutes which gives an indication of which team is dominating. While this is insightful, it does not give any information about where this is happening. Using a sliding window of 5 minutes on the role assigned player positions we can visualize play progression in terms of team formations and positions relative to one another, by representing the role distributions over the time window with 2D Gaussians. A film-strip of this approach is shown in Fig. 7.

D. Individual Player Analysis

Compared to existing analysis which often only looks at the mean behaviors of each player, our role assignment method dynamically assigns players to roles throughout a match and therefore, allows us to see the different characteristic behaviors of each player. First, we analyze the events that occurred within the match (See Fig 8). On the left we segmented all events by role. On the right we segmented events by player identity. In this example, the players playing left wing and right wing swap roles for part of the match, and the role representation is

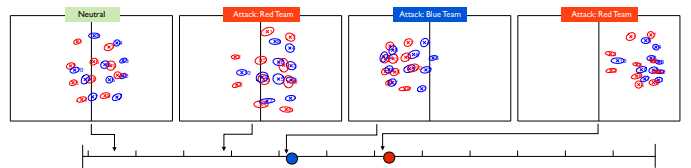


Fig. 7. By aligning the data based on role, we can quickly visualize and digest the flow of the match based on the formation of the team. Here we show snapshots of the match at different moments and highlight key events (circles=goals) with the home-team (red) attacking left-to-right. The x’s denote the mean role position across a small window of time (5 mins) and the ellipses show the variance of motion of each role.

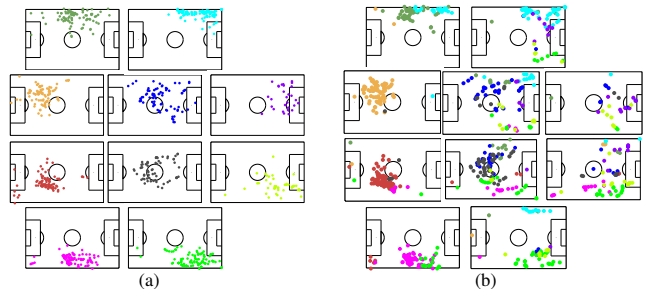


Fig. 8. Every event within a match half segmented into (a) roles, versus (b) player identity (both colored by the role of the player at the frame of the event)

effectively able to color these. If we were to simply take the mean of each players actions, we would miss this important tactical variation, and hence role provides important context in player analysis.

Next we examine the roles of each player over a match half as shown in Fig. 9. In this example, we can see the behavior of three teams and how players dynamically alternate relative positions throughout a match, essentially representing how versatile the players are within the formation. Plot (c) represents a 5 min smoothed version of the role assignments (to ignore temporary role swaps) showing dominant roles taken by each player. From this, it can be seen that in the top game, roles remain constant throughout the match, while in the 2nd game the midfielders (roles 5 and 6, shown in blue and grey) frequently swap positions, and in the 3rd game there are frequent role swaps between several of the players with only the back four players (i.e. top 4 rows in green, yellow, red and pink) remaining constant.

V. SUMMARY

In this paper, we presented a *role-based* representation to represent player tracking data, which was found by minimizing the entropy of a set of player role distributions. We showed how this could be efficiently solved using an EM approach which simultaneously assigns players to roles throughout a match and discovers the team’s overall role distributions. Using this method we show how we can perform both individual player and team analysis, providing context to player statistics and enabling large scale team analysis over a full season of player tracking data. In future work, we plan to use these methods to delve deeper into the various strategic patterns teams exhibit.

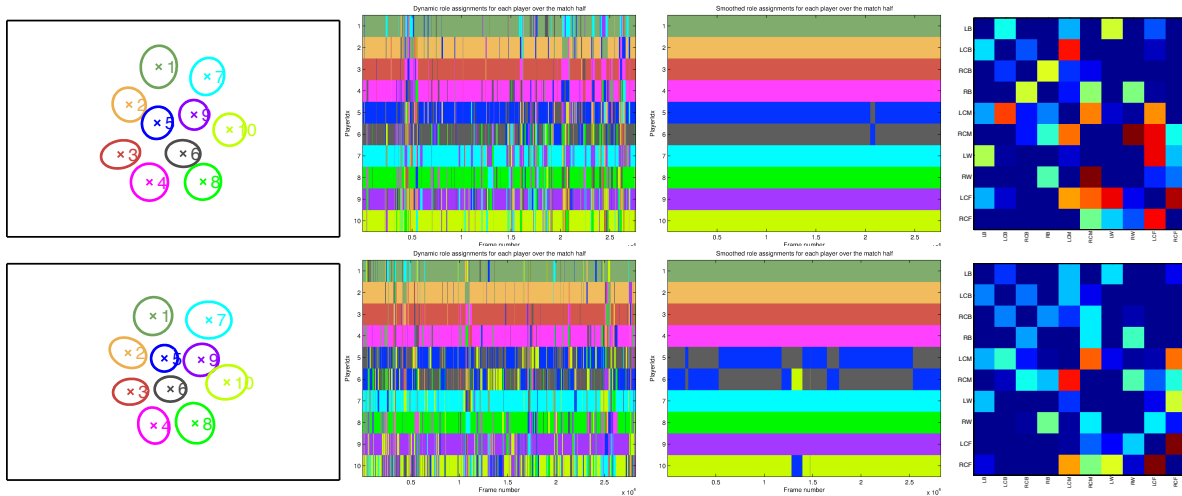


Fig. 9. The behavior of two different teams over half a match, demonstrating: (a) Their overall formation calculated using our minimum entropy data partitioning method (with roles represented as 2D Gaussians). (b) A timeline showing the role assigned to each player at each frame, colored by role. (c) A 5 min smoothed version of the role assignments (ignores temporary role swaps), (d) The role swaps across the half {1=left-back(LB), 2 = left-center-back(LCB), 3 = right-center-back(RCB), 4=right-back(RB), 5=left-centre-midfield(LCM), 6=right-centre-midfield(RCM), 7=left wing(LW), 8=right-wing(RW), 9=left forward(LF), 10=right forward(RF)}

REFERENCES

- [1] J. Yuan, Y. Zheng, C. Zhang, W. Xie, X. Xie, G. Sun, and Y. Huang, "T-Drive: Driving Directions based on Taxi Trajectories," in *GIS*, 2010.
- [2] L. Tang, X. Yu, S. Kim, J. Han, C. Hung, and W. Peng, "Tru-Alarm: Trustworthiness Analysis of Sensor Networks in Cyber-Physical Systems," in *ICDM*, 2010.
- [3] Y. Zheng, X. Xie, and W. Ma, "GeoLife: A Collaborative Social Networking Service Among User, Service," *IEEE Data Engineering Bulletin*, 2010.
- [4] J. Gudmundsson and M. Kreveld, "Computing Longest Duration Flocks in Trajectory Data," in *GIS*, 2006.
- [5] J. Lee, J. Han, and K. Whang, "Trajectory Clustering: A Partition-and-Group Framework," in *Int. Conf. Mgmt. of Data*, 2007.
- [6] STATS SportsVU, www.sportvu.com.
- [7] Prozone, www.prozonesports.com.
- [8] SportsVision, www.sportsvision.com.
- [9] Hawk-Eye, www.hawkeyeinnovations.co.uk.
- [10] E. GameCast, <http://www.espnfc.com/gamecast/392450/gamecast.html>.
- [11] P. Lucey, A. Bialkowski, P. Carr, S. Morgan, I. Matthews, and Y. Sheikh, "Representing and Discovering Adversarial Team Behaviors using Player Roles," in *CVPR*, 2013.
- [12] K. Goldsberry, "CourtVision: New Visual and Spatial Analytics for the NBA," in *MITSSAC*, 2012.
- [13] R. Masheswaran, Y. Chang, J. Su, S. Kwok, T. Levy, A. Wexler, and N. Hollingsworth, "The Three Dimensions of Rebounding," in *MITSSAC*, 2014.
- [14] D. Cervone, A. D'Amour, L. Bornn, and K. Goldsberry, "POINTWISE: Predicting Points and Valuing Decisions in Real Time with NBA Optical Tracking Data," in *MITSSAC*, 2014.
- [15] A. Miller, L. Bornn, R. Adams, and K. Goldsberry, "Factorized Point Process Intensities: A Spatial Analysis of Professional Basketball," in *ICML*, 2014.
- [16] P. Lucey, A. Bialkowski, P. Carr, E. Foote, and I. Matthews, "Characterizing Multi-Agent Team Behavior from Partial Team Tracings: Evidence from the English Premier League," in *AAAI*, 2012.
- [17] P. Lucey, D. Oliver, P. Carr, J. Roth, and I. Matthews, "Assessing team strategy using spatiotemporal data," in *ACM SIGKDD*, 2013.
- [18] J. Gudmundsson and T. Wolle, "Football Analysis using Spatiotemporal Tools," *Computers, Environment and Urban Systems*, 2013.
- [19] J. Pena and H. Touchette, "A Network Theory Analysis of Football Strategies," *arXiv preprint arXiv:1206.6904*, 2012.
- [20] X. Wei, P. Lucey, S. Morgan, and S. Sridharan, "Sweet-Spot: Using Spatiotemporal Data to Discover and Predict Shots in Tennis," in *MITSSAC*, 2013.
- [21] —, "Predicting Shot Locations in Tennis using Spatiotemporal Data," in *DICTA*, 2013.
- [22] S. Intille and A. Bobick, "Recognizing Planned, Multi-Person Action," *Computer Vision and Image Understanding*, vol. 81, pp. 414–445, 2001.
- [23] G. Zhu, Q. Huang, C. Xu, Y. Rui, S. Jiang, W. Gao, and H. Yao, "Trajectory based event tactics analysis in broadcast sports video," in *ACM Multimedia*, 2007.
- [24] M. Perse, M. Kristan, S. Kovacic, and J. Pers, "A Trajectory-Based Analysis of Coordinated Team Activity in Basketball Game," *CVIU*, 2008.
- [25] J.-C. Bricola, "Classification of multi-agent trajectories," Master's thesis, EPFL, 2012.
- [26] D. Stracuzzi, A. Fern, K. Ali, R. Hess, J. Pinto, N. Li, T. Konik, and D. Shapiro, "An Application of Transfer to American Football: From Observation of Raw Video to Control in a Simulated Environment," *AI Magazine*, vol. 32, no. 2, 2011.
- [27] K. Kim, M. Grundmann, A. Shamir, I. Matthews, J. Hodgins, and I. Essa, "Motion Fields to Predict Play Evolution in Dynamic Sports Scenes," in *CVPR*, 2010.
- [28] P. Carr, M. Mistry, and I. Matthews, "Hybrid Robotic/Virtual Pan-Tilt-Zoom Cameras for Autonomous Event Recording," in *ACM Multimedia*, 2013.
- [29] P. Lucey, A. Bialkowski, P. Carr, Y. Yue, and I. Matthews, "How to Get an Open Shot: Analyzing Team Movement in Basketball using Tracking Data," in *MITSSAC*, 2014.
- [30] A. Bialkowski, P. Lucey, P. Carr, Y. Yue, and I. Matthews, "Win at home and draw away: Automatic formation analysis highlighting the differences in home and away team behaviors," in *MITSSAC*, 2014.
- [31] X. Wei, L. Sha, P. Lucey, S. Morgan, and S. Sridharan, "Large-Scale Analysis of Formations in Soccer," in *DICTA*, 2013.
- [32] S. Roberts, R. Everson, and I. Rezek, "Minimum Entropy Data Partitioning," *IET*, pp. 844–849, 1999.
- [33] Y. Lee and S. Choi, "Minimum Entropy, K-Means, Spectral Clustering," in *International Joint Conference on Neural Networks*, 2004.
- [34] A. Dempster, N. Laird, and D. Rubin, "Maximum Likelihood from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society*, vol. 39, no. 1, pp. 1–38, 1977.
- [35] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval Research Logistics Quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [36] Y. Rubner, C. Tomasi, and L. Guibas, "A Metric for Distributions with Applications to Image Databases," in *ICCV*, 1998.