# Tackling Similarity Search for Soccer Match Analysis: Multimodal Distance Measure and Interactive Query Definition

Manuel Stein*
University of Konstanz

Halldor Janetzko†
University of Zurich

Tobias Schreck‡
Graz University of Technology
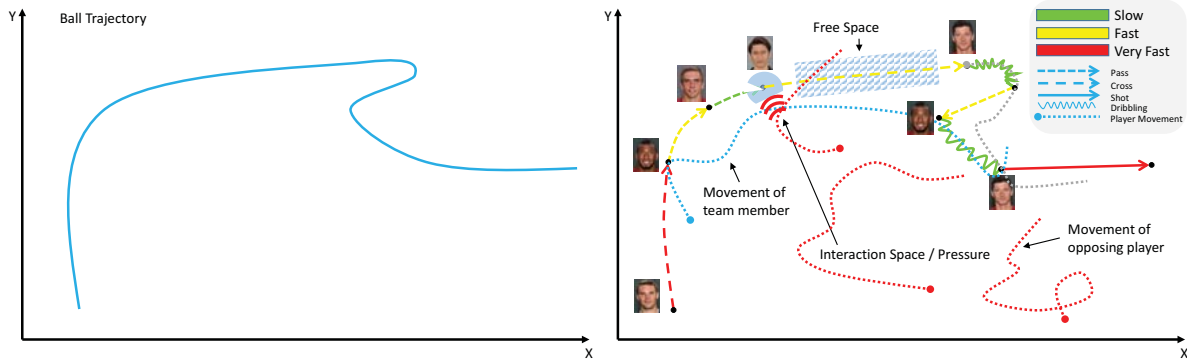
Daniel A. Keim§
University of Konstanz

Figure 1: Soccer movement trajectories are complex data. Existing trajectory similarity measures are typically based on spatio-temporal features, but lacking support for richer context. The trajectory on the left illustrates a trajectory of a soccer move consisting exclusively of x- and y-coordinates of the ball. The annotated trajectory on the right reveals the crucial movement context. Essential context data are, among others, the movement of the involved players of the ball possessing team as well as the movement of the players of the opposing team. Trajectory event data provide additional context information. We here study how to jointly take into account these data perspectives for similarity search as an important basis for soccer data exploration.

## ABSTRACT

Analysts and coaches in soccer sports need to investigate large sets of past matches of opposing teams in short time to prepare their teams for upcoming matches. Thus, they need appropriate methods and systems supporting them in searching for soccer moves for comparison and explanation. For the search of similar soccer moves, established distance and similarity measures typically only take spatio-temporal features like shape and speed of movement into account. However, movement in invasive team sports such as soccer, includes much more than just a sequence of spatial locations. We survey the current state-of-the-art in trajectory distance measures and subsequently propose an enhanced similarity measure integrating spatial, player, event as well as high level context such as pressure into the process of similarity search. We present a visual search system supporting analysts in interactively identifying similar contextual enhanced soccer moves in a dataset containing more than 60 soccer matches. Our approach is evaluated by several expert studies. The results of the evaluation reveal the large potential of enhanced similarity measures in the future.

---

*e-mail: Manuel.Stein@uni-konstanz.de
†e-mail: Halldor.Janetzko@geo.uzh.ch
‡e-mail: Tobias.Schreck@cgv.tugraz.at
§e-mail: Daniel.Keim@uni-konstanz.de

## 1 INTRODUCTION

The term *collective behavior* has been used for the first time in 1939 to describe a phenomenon observable for both humans and animals [9]. Common ground for the observation of collective behavior is that a number of entities behave in a similar, coordinated, or interdependent way. *Collective movement* as a branch of collective behavior is described as the movement of individuals in close proximity with similar speed and direction [23]. The analysis of collective behavior requires the precise recording of each individual movement, nowadays enabled by the persistent advancements and new developments of sensor and positioning technology. Resulting trajectory data consist of position coordinates of each individual as well as corresponding timestamps. A particular form of collective movement can be observed in invasive team sports such as soccer or basketball. In invasive team sports, members of a team want to reach a collective goal. Players need to make decisions and develop strategies in cooperation with their team members as well as in competition with the players of the opposing team. For training purposes supporting the players, team sport clubs employ professional (video) analysts to gain insights and to finally improve team performance. Improving the team performance requires identifying crucial and improvable tactical shortcomings of individuals or groups of players.

Nowadays, large data volumes can be captured from sport events. In soccer, in the top leagues but also elsewhere, computer vision techniques are applied to capture large trajectory and event data sets at high data rates, which need to be analyzed. The large data amount already captured in one single match or in matches of a complete season hinder manual in-depth exhaustive analyses. Therefore, assessing the performance of players is often done by computing numerical statistical features from the data sets. Simple match statistics such as a teams average ball possession, the number of shots on the opposing goal, or the amount of won tackles are nowadays mainly used to extract team performance [44, 58]. However, the

resulting statistics are challenging to interpret as aggregate statistics hardly reflect the individual or single scenes of interest, and are highly context-dependent.

An important inherent challenge of having more and detailed data is the selection of the crucial information pieces. Manual analyses are not feasible and fully automated methods can only be applied when the analyst knows the desired patterns in advance. Consequently, supporting analysts focusing on the key aspects is necessary for a successful analysis but highly context-dependent and usually ill-defined. For the identification of desired reoccurring collective movement patterns, analysts recently started to make use of conventional spatial distance measures such as DTW [8], LCSS [55] or EDR [21] during similarity search or clustering. These distance measures compare the spatial distance between two trajectories as sequences of points.

However, soccer trajectories are far more complex than current distance and similarity measures are taking into account, as demonstrated in Figure 1. The trajectory on the left of Figure 1 illustrates a trajectory of a soccer move consisting exclusively of x- and y-coordinates of the ball. Figure 1 right instead illustrates the existing movement context being lost. Essential context data are, among others, the movement of the involved players of the ball possessing team as well as the movement of the players of the opposing team. Trajectory event data provide additional context information. A passing event, e.g., indicates that a player has seen an advantage in passing the ball to another player, while a shot on goal event might indicate a successfully executed tactical behavior. The movement of players of the opposite team can also be seen by a team's interaction and free spaces [53].

In this paper, we propose an enhanced and flexible similarity measure (Section 3) to interactively integrate contextual information in addition to spatio-temporal data. We illustrate and provide examples for essential characteristics of soccer moves being of interest for analysts, and explain how to integrate and combine identified characteristics into an enhanced similarity measure. The resulting similarity measure is designated to improve our understanding of collective movement patterns. To develop such an enhanced similarity measure, we provide an overview of the commonalities in the current state-of-the-art (Section 2) including distance measures for point and movement distance. We present a user-centered visual searching system (Section 4) supporting analysts in their everyday work by enabling an efficient exploration of large amounts of soccer movement data. Our presented system allows both *query-by-example* and *query-by-description* analysis and provides several interactive visualizations enabling analysts to steer the analysis process as well as to explore the results. The resulting enhanced similarity measure is evaluated (Section 5) with the help of two experienced soccer domain experts. We discuss the results of our evaluation and conclude our paper summarizing and highlighting open research aspects (Section 6).

## 2 RELATED WORK

Distance measures are the foundation of many pattern detection algorithms as, for example, clustering, classification, or querying. To develop an enhanced similarity measure for incorporating movement context, this section provides an overview to the current state-of-the-art of trajectory distance measures as well as systems for interactive trajectory search in sport analysis. Finally, we position ourselves within the aforementioned works.

### 2.1 Trajectory Distance Measures

A trajectory consists of a sequence of points which in turn consist of position coordinates (x- and y-coordinates) as well as an optional time stamp. Most existing distance measures for trajectories, therefore, calculate the distance of two trajectories based on their spatial properties. For example, the Fréchet- and Hausdorff-metrics [25, 31] calculate the similarity of two trajectories based on the distance of

a single point combination. The Hausdorff-distance is mainly used in video surveillance and pattern recognition [33, 38]. Chen et al. developed an trajectory clustering algorithm based on Hausdorff and DBSCAN [27] to analyze hurricane trajectories [18]. The Fréchet-distance as well as its variants (e.g., weak Fréchet-distance [2], discrete Fréchet-distance [25]) are applied, for example, for optical character recognition tasks in document [52] and signature [59] collections. Furthermore, the Fréchet-distance has been applied for similarity search in traffic movement data [11, 57].

In contrast to Hausdorff- and Fréchet-metrics, which calculate the trajectory distance based on single point combinations, various algorithms exist calculating the distance based on a larger amount of point comparisons. The edit-distance, for example, calculates the number of *insert*, *delete* as well as *replace* operations required to transform a character string into another one. Originally, the edit-distance is used for text comparisons [56] and represents a whole set of distance measures based on the given edit functions [39]. For trajectory data, among the most common distance measures based on variants of the edit-distance are the *Edit Distance on Real Sequence* (EDR) [21], the *Edit Distance with Real Penalty* (ERP) [20] or the *Longest Common Subsequence* (LCSS) [55]. Both, EDR and ERP are used for trajectories with temporal offsets. Furthermore, EDR has been developed especially for trajectory data containing sensor-related outliers. Both distance measures have been used for traditional movement trajectories [60] as well as voice or music recognition [19]. LCSS is known as robust algorithm for data with measurement errors [55]. Due to this robustness, LCSS is often used for the detection and comparison of vehicle trajectories [17]. Furthermore, Sivaraman et al. [50] showed that LCSS can be applied successfully as distance measure for clustering tasks. In the given example, LCSS was used to identify movement trajectories indicating dangerous driving maneuvers based on sensor data of self-driving cars. Further works in the field of autonomous driving have been performed by Shirazi and Morris [49] as well as Choong et al. [22]. They used LCSS as similarity measure to cluster and count turning trajectories at crossings. Additionally, LCSS has been used for the analysis of time series of different lengths [30, 36].

*Dynamic Time Warping* (DTW) is explicitly taking the temporal aspects of a trajectory into account as it allows identifying similarities between trajectories containing temporal offsets through stretching operations. DTW as well as its variants (e.g., SDTW [34], PDTW [35], FastDTW [47], FTW [46]) are used in a wide variety of applications, for example, for spoken words comparison in linguistics [45, 54] or for clustering gene expression data [1, 32]. Furthermore, DTW has been applied in optical character recognition tasks [15, 16, 28, 43] such as plagiarism detection.

### 2.2 Contextual Trajectory Distance Measures

One of the first work extracting contextual movement features from trajectory data has been proposed by Dodge et al. [24]. They describe a method to extract characteristics such as speed, acceleration, direction and length from moving entities. Their approach is evaluated with an anonymized trajectory dataset containing pedestrians, car as well as bike trajectories. The extracted movement characteristics are used during evaluation to classify the given trajectories. Promising results indicate that a pure observation of the shape of a trajectory is not sufficient and that analysis can be improved by making use of additional contextual movement features. Andrienko et al. [6] extended the analysis of movement by including contextual spatio-temporal as well as event data. Their analysis focuses on the relations of moving entities and their surroundings described by the movement context. However, both works do not focus on the development of an enhanced distance measure combining these aspects.

Buchin et al. proposed the first contextual trajectory distance measure by enhancing the Fréchet-distance using several contextual

features [13,14]. For a dataset of hurricane trajectories, they used the current terrain as contextual feature while for albatross trajectories they made use of the wind direction. However, their work focuses entirely on the linking of spatial contextual features. Trajectory characteristics of other domains are not taken into account. While Buchin et al. incorporate basic geographic contextual features such as whether a tornado is moving over land or sea, they do not discuss the influence of entities on each other when moving in groups although they consider them to be very important [14]. In another work [12], Buchin et al. incorporated basic influences of moving entities on each other as, for example, only considering two entities following each other, if they are moving into one direction.

## 2.3 Interactive Trajectory Search for Sport Analysis

To date, several interactive approaches have been proposed for searching for trajectory data in different applications. In [29], sketch-based search for traffic routes within a street network is proposed. The user specifies a polyline query path by setting and adjusting control points on a map. Also, the widely-used line chart can be regarded to define trajectories. In [7], query-by-example and query-by-sketch search are proposed for exploration of patterns in time series databases.

Recently, visual analysis of sport data has received increased research attention. In [42] a comprehensive overview of current approaches is given. It is argued that as the data and analysis questions of concern typically show high variability, user interaction and search are particularly helpful for analysis. In [37], sketch-based search for team movements in rugby is proposed, supporting comparison of match situations. Specific to soccer, in [41] user search is supported by selection of time intervals or events of interest based on meta data. The system allows to compare situations by aggregation of movement elements, overlays and other rich data visualizations. In previous work [48], we have presented sketch-based search for soccer movement based on free-form trajectory sketches and specification of start-stop areas. The search in that work was based strictly on geometric trajectory features.

## 2.4 Positioning

The current state-of-the-art in trajectory distance measures reveals that a large variety of algorithms calculating the distance between trajectories exist. The focus of most distance measures, however, solely lies on the trajectories spatial coordinates. The influence of moving entities to each other is barely represented in state-of-the-art trajectory distance measures. There exists no distance measure for the comparison or aggregation of trajectories which includes various features from movement context. Furthermore, no current distance measure is able to incorporate high-level movement features such as player pressure. Our proposed enhanced similarity measure fills this gap by incorporating low and high-level movement features as well as event and player context.

## 3 INCORPORATING MOVEMENT CONTEXT IN THE SIMILARITY MEASURE

In the following section, we propose and explain our developed distance measures for soccer movement queries based on geospatial data, player pressure, and event and player annotations. Based on these individual measures, a flexible aggregate measure is defined.

### 3.1 Spatial Context

Important characteristics of a soccer trajectory with spatial context are for example position and speed. The speed of a trajectory is determined by combining points in time with their corresponding locations. We are interested in cases, in which several game moves are spatially far away, but nevertheless show similar speed patterns. Exactly opposite situations in which moves are spatially close to each other but have very different speed patterns are interesting as

well. Examples of such situations can be seen in Figure 2 (a) and Figure 2 (b). In both illustrations, a darker color of a trajectory means a higher speed. Figure 2 (a) shows two passplay moves with a larger spatial difference. A standard spatial distance measure would rather not evaluate two such trajectories as similar. However, these moves could be interesting for an analyst due to their similar speed patterns. Figure 2 (b) shows two spatially similar moves. Interestingly, the difference between the two moves can only be recognized with background knowledge from the context of trajectories. The upper fast move shows a goalkeeper kicking the ball deep into the opponent's half, which is then extended to a shot on goal shortly afterwards. The slower lower move shows the ball rolling after a low throw by the goalkeeper, followed by a longer pass combination to the opposing goal. The type of move is consequently also reflected in the speed of the move.
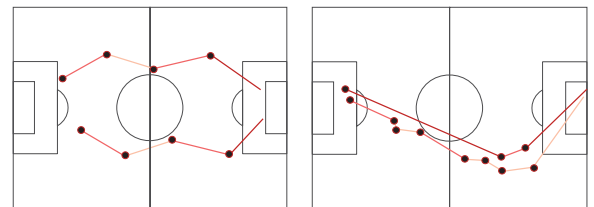
Accordingly, we combine position and speed to be able to differentiate between trajectories with similar shape but different velocity. The Euclidean Distance is a common distance measure for point data and is used in many distance measures for trajectories. However, the Euclidean Distance only uses the position feature for distance calculation. Consequently, we extend the Euclidean Distance as a distance measure for the points of a trajectory. According to the standard definition, a point of a trajectory consists of its coordinates and a time stamp. Our concept adds speed to an extended point $ep$ in addition to the standard spatial features. The speed of a point is calculated from the difference in distance and time from a previous point. Theoretically, this concept can be extended by any number of additional features. It is only important that each characteristic can be calculated at any point of a trajectory.

An extended point $ep$ consists of the coordinates $x$ and $y$ and a time stamp $t$ like a standard point. The variable $i$ shows the position within a trajectory. In addition, each point has a number of geographical features. For our play trajectories one point is extended by the characteristic speed $s$.

$$ep(i) = (x, y, t, s) | i \geq 0$$
$$ep(0) = (x, y, t, 0)$$

By combining different characteristics (features), distances from different value ranges need to be integrated, which requires normalization. Here, we normalize the feature ranges based on domain knowledge. In general, soccer pitches have a maximum length of 105 meters and a width of 68 meters. This gives a diagonal of 125.10 meters, the maximum resulting distance between two points. The speed of a move must also be limited by an upper limit in order to be normalized. We define *moves* as intervals of arbitrary length, starting with the gaining of the ball and ending with a final turnover. A study of our data showed that the average speed of more than 5000 moves is around 30 kilometers per hour. Shots at an average speed of 120 kilometers per hour are well above the average speed



(a) similar speed, dissimilar shape   (b) similar shape, dissimilar events

Figure 2: The similarity between two soccer moves depends on the characteristics taken into account. In this figure, the speed of the trajectory is encoded by color from yellow (slow) to red (fast).

|  | | | | | |
|---|---|---|---|---|---|
| $p_1(5)$ | 120 | 77 | **26** | **23** | **22** |
| $p_1(4)$ | 70 | **11** | 31 | 42 | 34 |
| $p_1(3)$ | 35 | **10** | 28 | 63 | 70 |
| $p_1(2)$ | **12** | 31 | 34 | 65 | 88 |
| $p_1(1)$ | **7** | 17 | 49 | 90 | 150 |
|  | $p_2(1)$ | $p_2(2)$ | $p_2(3)$ | $p_2(4)$ | $p_2(5)$ |

Euclidean Distance

$$\sqrt{(p_1(i) - p_2(i))^2}$$

(a)

|  | | | | | |
|---|---|---|---|---|---|
| $ep_1(5)$ | 24 | 19 | 17 | 5 | **3** |
| $ep_1(4)$ | 20 | 18 | 11 | 7 | **6** |
| $ep_1(3)$ | 13 | 17 | 8 | **3** | 11 |
| $ep_1(2)$ | 10 | 14 | **4** | 13 | 24 |
| $ep_1(1)$ | **1** | **7** | 12 | 15 | 22 |
|  | $ep_2(1)$ | $ep_2(2)$ | $ep_2(3)$ | $ep_2(4)$ | $ep_2(5)$ |

Extended Euclidean Distance
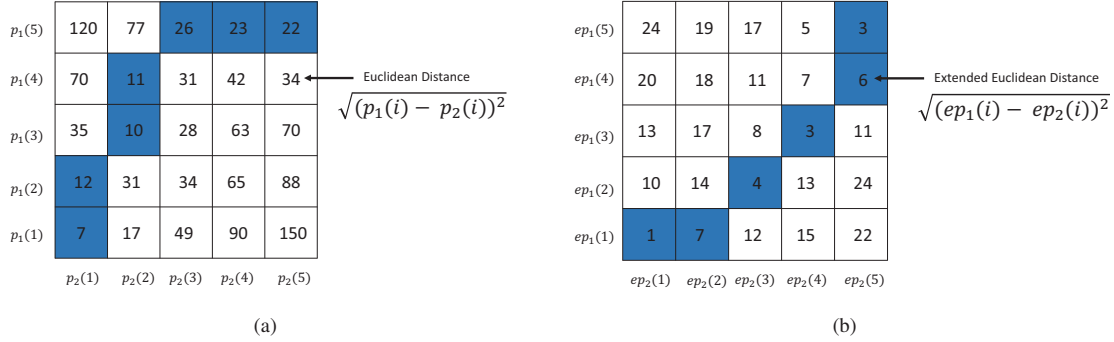
$$\sqrt{(ep_1(i) - ep_2(i))^2}$$

(b)

Figure 3: Visualization of two DTW distance matrices using the Euclidean Distance in (a) and the extended Euclidean Distance in (b). The optimal warping path is visualized in blue.

of a move. After consultation with an expert, we set the upper limit for the speed of the play train to 120 kilometers per hour, since speeds above 120 can only be shots and the speeds of dribbles and passes make up the more interesting range of values. To enable a user to adjust the combination of the characteristics position and speed to her or his analysis requirements, we link the normalization of the characteristics with a weighting. Then we extend the standard Euclidean distance by the difference in speed, the normalization of characteristics and an option for weighting. This gives us the following formula for an extended Euclidean distance measure based on our earlier defined extended points, here represented by $p$ and $q$. $w_1$ and $w_2$ enable to adjust the weighting for position and speed while $N_1$ and $N_2$ are used for the normalization.

$$D_{eECL}(p,q) = w_1 \cdot N_1 \left( \sqrt{(q(x) - p(x))^2 + (q(y) - p(y))^2} \right) + w_2 \cdot N_2 \left( |q(s) - p(s)| \right)$$

The main advantage of our advanced Euclidean metric is that it can be included in most distance measures for trajectories without any further adjustments of the similarity algorithm. After several tests of the extended Euclidean distance with the distance measurements Fréchet, EDR, LCSS and DTW, we decided to use the DTW algorithm. As described in Section 2.1, the DTW algorithm uses a distance matrix consisting of the distances of all points of two given trajectories. In our spatial distance measure, we adjust the calculation of the distance matrix by using our introduced extended Euclidean distance. Figure 3 shows the change of a distance matrix by applying the extended Euclidean distance with the characteristics position and speed. The optimal warping path is highlighted in blue. By using the extended Euclidean distance as a point distance measure, the individual point distances in the distance matrix can change significantly. As a result, the optimal warping path can also find a significantly different path through the distance matrix.

### 3.2 Player and Event Context

Event data of a move is available as a list of positions with corresponding time stamps and textual descriptions of the particular event. The description typically includes the type of event and, if applicable, the players involved. Most events relate to a movement pattern of the ball, including events such as a pass or dribbling. The incorporation of event analysis is interesting as, for example, frequent repetitions of an event sequence can indicate tactical patterns. Such tactical patterns are not necessarily high spatial similarity, which would make them undetected by conventional shape-focused similarity measures.

The player data of a soccer match is also available as a list of position data with corresponding time stamps. Additionally, player name, jersey number as well as player position such as defender or striker are stored. For analysts, for example, moves involving the same players or the same player positions are of particular interest. Figure 4 (a) and Figure 4 (b) show two situations where the move with less spatial similarity could be the more interesting move for an analyst. Move $V$ is the query trajectory for which similar moves are to be found and compared. Possible identified similar moves are trajectories $A$ and $B$. Move $A$ has less spatial similarity to $V$ than $B$. Consequently, standard trajectory distance measures would identify move $B$ as the more similar move to $V$. However, an analyst would notice the similarity of the players in Figure 4 (a) and the player positions in Figure 4 (b) and would probably identify trajectory $A$ as the more similar move.

In order to consider the sequential structure of event and player data, we propose the use of local alignment algorithms, looking for similarities in sections of sequences. In our application scenario we can thus calculate the largest sequence of identical events or players between two moves. For the local alignment calculation, we compared the algorithms Smith-Waterman [51], FASTA [40] and BLAST [4]. Smith-Waterman calculates the optimal local alignment of two sequences. However, the biggest disadvantage of the Smith-Waterman algorithm is its time-consuming calculation. The algorithms FASTA and BLAST can calculate the local alignment much faster in comparison to Smith-Waterman, but have the disadvantage that they only heuristically approach the optimal result. In the end, we chose the Smith-Waterman algorithm because we only



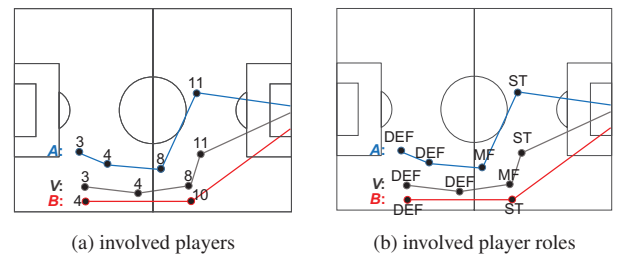(a) involved players          (b) involved player roles

Figure 4: In both displayed cases, the move with less spatial similarity (blue) could be the more interesting move for an analyst based on involved players or player roles. Move $V$ is the query trajectory for which a similar move is to be found.

Figure 5: Our visual search system supports analysts in interactively identifying similar contextual enhanced soccer moves. Displayed is an overview about the system and the interactive query definition. In the shown example, the query is enhanced by player and event information.

compare short event and player sequences and, therefore, do not expect any performance problems. The Smith-Waterman algorithm uses so-called editing functions to find local similarities between two sequences. These functions check individual sequence elements for equality or inequality and can make decisions for inserting or deleting sequence elements. Each editing is assigned a specific value, which is summed up during use (positive for *equality*, negative for *inequality* as well as *insert* and *delete*). In order to be able to combine the alignment with other characteristics and to enable comparability of the results, the result must be normalized. Here, Altschul and Erickson [3] have shown that a comparison of local alignment results is only possible by a combination of alignment value and length. The result of the normalization represents the similarity of two sequences in percent.

### 3.3 Pressure and Further High Level Context

The exertion of pressure (*pressing*) on opposing players can be a tactical pattern in soccer matches. Pressing is performed by the defending team against players of the attacking team. One aim of pressing is to prevent an attacking move of the opposing team. Especially the player with the ball is often the target of pressing to either enable a ball capture or to provoke false passes. A special form of pressing is the so-called *counter-pressing*. Counter-pressing is the direct pressing on the ball possessing player of the opposing team, shortly after the ball has been lost. Pressure is calculated by the player movement. For each player, the level of pressure can be calculated by the proximity to opposing players. The closer an opponent gets to this player, the more he or she gets under pressure. Experts an-

alyze pressing behavior as one possible basis for evaluating a move. Moves in the opposing penalty area without opposing pressure, for example, can indicate an uncoordinated defensive performance by the defending team.

Several models for calculating player pressure have been proposed. We use the method introduced by Andrienko et al. [5] which calculates pressure based on the region surrounding each player. The closer an opponent is to the player in the center, the larger is the corresponding pressure. The position of the opponent is also crucial. If an opponent stands in the running direction of a player, the pressure on this player is rated higher than if the opponent stands behind the player under pressure. To calculate the average pressing of a move, the maximum pressure on the pressing target for each point in time is accumulated and divided by the duration.

### 3.4 Resulting Similarity Measure

After a distance measure has been found for each characteristic, the similarity between two moves $P$ and $Q$ can be calculated for each move using the following formula.

$$D_{Context}(P,Q) = \sum_{i \in I} D_i(P,Q) w_i$$

$$I = \{geo, event, player, pressing\}$$

By normalizing the individual distance measures, the results have a common value scale and can be aggregated. The weighting allows the individual distance measures to be put in relation to each other

(a) Selection of several players for a filter



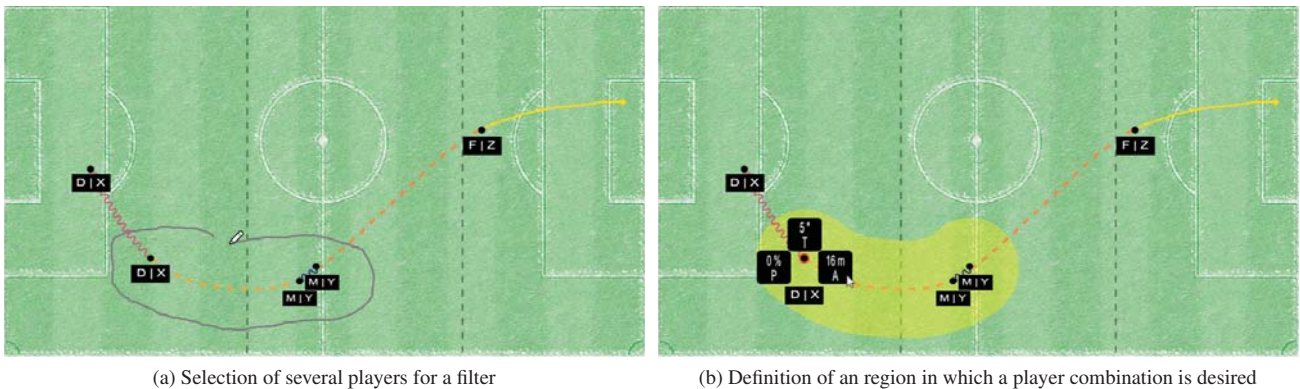(b) Definition of an region in which a player combination is desired

Figure 6: Drawing and editing filters allows an analyst to enforce the presence of specific player or event combinations or free spaces in the search results. This can be interesting for analysts, for example, when looking for a certain tactical pattern and only wanting to find results containing this pattern.

and thus to influence their semantic relevance. Due to the algorithm's modularity, further features can be added in the future.

## 4 INTERACTIVE SYSTEM FOR USER QUERIES AND RESULT VISUALIZATION

In the following, we present a system (Figure 5) for the effective search for similar moves in soccer. The system is primarily aimed at analysts of large soccer clubs in everyday operations. The analysts are enabled to search for particularly similar moves over an entire season in order to subsequently be able to test their own hypotheses about the playing behavior of their own or the opposing team. In particular, the interactive creation of the search query should be supported. In addition, the search results should be adjustable to the analysts' needs by weighting search characteristics. Our similarity measure is used to identify the most similar moves. For the results to be comprehensible for the user, they should be able to be examined and compared for each characteristic. The visualizations required for this are intended to provide the analyst with an overview and, if desired by the user, more detailed views. The system is designed to assist analysts in examining the results and suggest potentially interesting results. For an analyst to start more in-depth analyses based on already identified results, the system should be able to use game moves found as a basis for new searches. A video describing and showing our system in use can be found online (http://files. dbvis.de/stein/Enhanced_Similarity_Search.mp4).

### 4.1 Defining the Search Query

In order for analysts to be able to search for similar moves, they must first describe a move in the system. To keep the creation of a move as intuitive as possible, we decided to use a visualization of a soccer pitch comparable to a tactic board. Tactic boards are small drawing boards with a marked playing field. Coaches use tactical boards to visually explain tactics to players. The search query should include the shape of a move. An analyst should also have the possibility to extend the search query by adding features such as players involved, events or pressure. Visual variables such as color and shape are designed to help the analyst differentiate between these characteristics. Coaches already use different visualizations for the events of a move, for example a wavy line is used for dribbling. We display passes with a dashed line, dribbling with a wavy line and shots as well as kicks with a straight, solid line. The interaction concept *Query-by-Sketch* is used to create a search trajectory. The direct drawing of a move is intended to remind of the tactical boards of soccer coaches. Drawn soccer moves are, furthermore, smoothed

by applying *Smoothing Via Iterative Averaging* (SIA). To get more accurate search results, an analyst is able to manually add features such as events, players and their positions, speed as well as pressure to any position of a move. The pressure as well as the speed at different parts of the move can be changed by scrolling with the mouse wheel. Additionally, the system allows analysts to save and load created as well as identified moves so that they do not have to recreate them every time they perform recurring searches for specific moves. An example of an user created query can be seen on the soccer pitch in Figure 5.

### 4.2 Filtering Data

Our proposed system offers an analyst three filtering methods to restrict the search results. These filter methods allow an analyst to enforce the presence of specific player or event combinations or free spaces in the search results. This can be interesting for analysts, for example, when looking for a certain tactical pattern and only wanting to find results containing this pattern. Filtering by event combination, for example, allows an analyst to search only for moves with a shot. Filtering can also be linked to a position. This makes it possible, for example, to search for shots in front of the opposing goal. Other interesting event combinations are, for example, passes
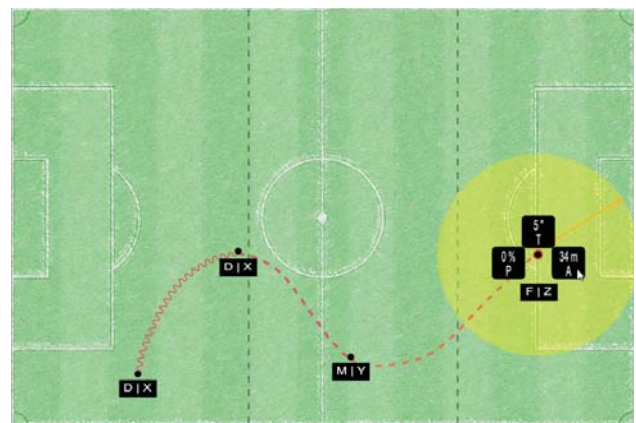


Figure 7: Drawing an event filter near the goal of the opposing team. The drawn circle allows the user to define in which area the shot on goal has to occur.
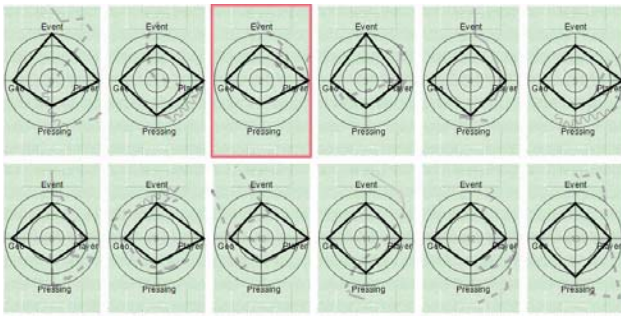
Figure 8: We use a small multiple star glyph visualization to provide a visual overview about the results of the similarity search.



Figure 10: The pressure visualization of the ball possessing players shows areas with high or low pressure for each move.

and crosses from the sides in the opponent's penalty area and shots within the penalty area allowing to explicitly examine the positional play of the attacking as well as the defending team. By filtering for player combinations, analysts can search for moves where they expect a player at a certain position. A possible application scenario for this filter is, for example, to investigate whether a defender was always at the intended position during a trained move or whether further training of a move is required. The third filter allows an analyst to filter moves for the presence of certain free spaces, which can indicate good passing options.

The interaction concept *Query-by-Sketch* is used again to enable an intuitive operation of these otherwise rather difficult to define filters. For players (Figure 6) and event (Figure 7) combinations, the corresponding object can be selected or encircled with the mouse. The selected players and events can then be given a position restriction. The area in which the desired player or event may be located can be adjusted using the mouse wheel. For a free space, the analyst can draw any polygon on the playing field and thus directly determine position and size of the filter.

### 4.3 Weighting Features

Using the weighting, users can directly interact with our proposed enhanced similarity measure and thus adapt the search process with respect to their expertise. Furthermore, the weighting allows the user to relate the distance measures to each other. For example, if analysts are only interested in the shape of a move as well as certain existing players, they can lower prioritize distance measures such as speed or events by decreasing their corresponding weight. We have chosen a discrete scale of 0 to 100 percent as the weighting factor, whereby a distance measure at a factor of 0 percent has no influence on our extended similarity measure and a weighting factor of 100 defines a distance measure as maximum important.
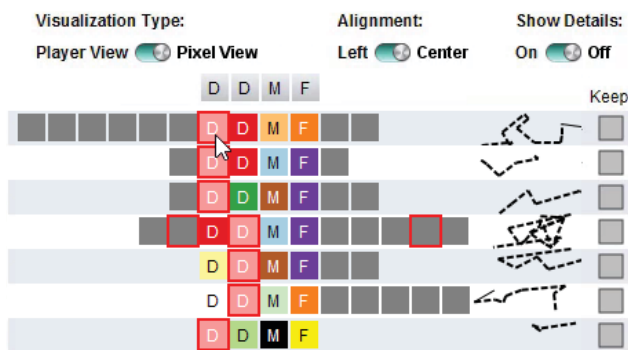


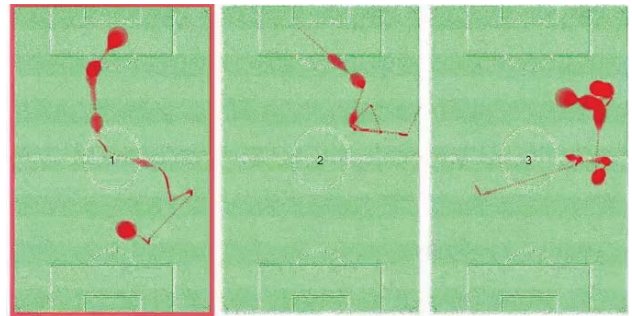Figure 9: Pixel visualization highlighting same sequence elements.

### 4.4 Visualizing the Results

The aim of the result visualizations is to support analysts investigating identified similar moves. Our proposed system provides overview as well as detail visualization for the overall similarity measure as well as for each distance measure. A visual overview about the results of our similarity search is provided via small multiple star glyphs (Figure 8). For spatial features, the overview visualization consists of an list of small multiples ordered by the calculated similarity. Each small multiple represents a single move. When focusing one small multiple, a detailed view enables the user to inspect shape, speed (by color), involved players, occurring events as well as an animation of the exact player movement. For the efficient comparison as well as pattern detection of players and events in different moves, we propose a pixel based visualization (Figure 9). All moves are visualized in a descending list. Each row contains the player or event sequence and a rough representation of the move trajectory to distinguish the moves. Pixels can either be visualized from left to right in one row or centered based on the common characteristics of the subsequent rows. The color of each pixel represents a single player or player role. To enable analysts to concentrate on similar parts of players and event sequences, the system hides all dissimilar sequence elements by default. Furthermore, users can interactively highlight and filter the displayed data. The pressure visualization shows areas with high or low pressure for each move. To display zones with high or low pressure, the amount of pressure on the ball possessing player is calculated for each point of a move and afterwards visualized via the size of the ball on top of the trajectory. Figure 10 illustrates how the size of the ball position is influenced by pressure.

### 5 Evaluation

In order to examine the usefulness of our enhanced similarity measure and our system for the visual analysis of the resulting moves we invited several domain experts. The first invited domain expert has been an active soccer player for 26 years and, additionally, has been working as a coach for 12 years in the youth sector of the German soccer club FC Bayern München. The second expert has been an active soccer player for more than 20 years and is currently working as a DFB (German Football Association)-certified coach At the beginning of our evaluation, the invited experts were given a comprehensive overview of the presented search system. Afterwards, each expert was given the opportunity to create own contextual enhanced search trajectories and search for similar moves based on a dataset containing 60 matches from two professional soccer clubs from an international first league. For every match, the data consists of the movement of every player as well as ball movement and event data. While using our system, experts were asked to express *ad hoc* comments via the thinking aloud method [10, 26]. At the end of the evaluation, we performed a semi-structured interview focus-

ing on the creation of search queries via sketching, the different developed contextual distance measures, the designed result visualizations as well as asked for fields of application and suggestions of improvement.

The results of our evaluation are very promising. Each expert used our designed system over the course of several hours. Our *query-by-sketch* approach is perceived very intuitive by both experts as it reminds them of drawings on a tactic board. Both experts also like the possibility to annotate additional features of a move as well as to restrict and define regions in which, for example, certain events have to occur. One expert explained that he would make use of such a system in order to identify and filter moves that reached the opposing goal with specific player and event combinations. Both experts emphasized the potential of our system when applied on tablet computers enabling them to use our system in the open as well as drawing and annotating moves with their fingers. The possibility to weight the several introduced distance measures was also mentioned positively. Both experts claimed that they want to adjust the importance of each features based on their current search goal and expertise. Both experts experimented with the weightings and inspected the different results.

Our designed distance measures as well as the resulting enhanced similarity measure are approved by our invited experts. The experts were not aware of a missing distance measure. The various provided overview as well as detail visualizations are seen positively as well. For example, the experts emphasized the usefulness of the provided overview small multiples enabling analysts to compare results and inspect subsets in detail visualizations. For the player and event pixel visualization, the experts pointed out the usefulness in detecting repeating patterns, especially in front of the opposing goal. As future improvement, they suggested using additional colors to display a players position. Furthermore, one expert explained that events in the first third of a soccer pitch can be irrelevant, as here the players usually act without pressure from the opponent and can therefore play the ball at will. To solve this issue, the expert suggests removing event sequences from the first third of the soccer pitch. Our proposed pressure visualization is well perceived by the invited experts. The experts complimented in particular how easy it is to recognize zones of pressure. Both experts analyzed identified patterns of pressure in various moves. One expert used the pressure weighting as well as visualization in order to determine where a team wins the ball under high pressure. Accordingly, he searched for moves with high pressure and then checked whether the players have chosen the correct playing styles. Eventually, the experts approved the possibility to inspect each identified move in detail through an interactive animation. With the help of the animation, both experts verified what they considered to be the decisive situations of identified moves. Patterns found included insufficient man coverage, a player's incorrect position, insufficient pressure on an opponent and incorrect coordination between players and their teammates.

Overall, both experts are satisfied with the opportunities of our system in order to create and search contextual enhanced moves. According to one expert, our system contains almost everything he would expect from a novel system for the search of soccer moves. Both experts believe that the proposed system will be used for the analysis of moves and that interesting tactical patterns can be discovered in similar moves. The invited experts see two application areas for the system in particular. The first area of application is the verification of ideas and hypotheses of an analyst. During a soccer match, an analyst generates hypotheses about a team's match plan based on observed and annotated moves. Using our system, they would compare the annotated moves with all the data available for the observed team in order to confirm or reject their hypotheses. The other way to use the system is to search for a new player to acquire on the transfer market. Professional coaches have a clear idea of, for example, the offensive game of their team and how the players should act in it. A coach is therefore always on the lookout for players who can play moves according to their ideas. Given the data, our system allows experts to examine similar moves to the coaches' match plan and the involvement of a potential candidate. Furthermore, the experts noted that the system allows players to be analyzed in a very simple way without having to watch the full 90 minutes of a soccer match.

## 6 DISCUSSION AND CONCLUSION

In this paper, we proposed an extended similarity measure, taking into account the context and events of a move in addition to its shape. To the best of our knowledge, there exists no comparable similarity measure for trajectories that is combining high-level contextual with spatial characteristics. Our similarity measure as well as the contained distance measures were designed very modular, thus it is possible to extend it by further trajectory characteristics. The weighting of the distance measures allows us to put them in relation to each other. This allows users to influence the relevance of the results according to their needs. Given enough training data, machine learning and recommender systems could be applied in the future to automatically detect user favorable weightings. Another interesting future development is an even stronger inclusion of player trajectories in the search process. Although player trajectories already influence the calculation of some high-level features such as pressure, the experts expressed great interest in developing and integrating more player movement related features.

Furthermore, we presented a system for the interactive sketch-based search of similar soccer moves in large amounts of match data. We put special emphasis on an intuitive and interactive operation of the system, for example, during creation of the search query. As a general future development of the system we consider the use with data from other invasive team sports such as American Football. American Football is characterized by the fact that there is only a relatively limited repertoire of different moves and a coach can call on up to 53 players. These moves must be executed by the players with utmost precision. Possible analyses are, for example, to identify exactly how well a player knows the running routes and how high his chances of success after a catch are. Conversely, it will also be very interesting to analyze which player might be the best choice for a particular move.

We evaluated our enhanced similarity measure with two domain experts in several expert studies. During the evaluation, we used match data from more than 60 matches over a complete season from two teams. The results of our evaluation indicate that experts favor our proposed enhanced similarity measure over traditional similarity measures that are limited to the shape of a trajectory. Both experts are convinced of the usability and would like to test it with data from their teams. Nevertheless, in future work a detailed quantitative evaluation of the search algorithm is desirable in addition to the qualitative evaluation. For this, however, a ground truth data set would first have to be created by an expert, in which the most similar moves for several moves are presorted.

### REFERENCES

[1] J. Aach and G. M. Church. Aligning gene expression time series with time warping algorithms. *Bioinformatics*, 17(6):495–508, 2001.
[2] H. Alt and M. Godau. Computing the fréchet distance between two polygonal curves. *International Journal of Computational Geometry & Applications*, 5(01n02):75–91, 1995.

[3] S. F. Altschul and B. W. Erickson. Locally optimal subalignments using nonlinear similarity functions. *Bulletin of mathematical biology*, 48(5-6):633–660, 1986.

[4] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.

[5] G. Andrienko, N. Andrienko, G. Budziak, J. Dykes, G. Fuchs, T. von Landesberger, and H. Weber. Visual analysis of pressure in football. *Data Mining and Knowledge Discovery*, pp. 1–47, 2017.

[6] G. Andrienko, N. Andrienko, and M. Heurich. An event-based conceptual model for context-aware movement analysis. *International Journal of Geographical Information Science*, 25(9):1347–1370, 2011.

[7] J. Bernard, D. Daberkow, D. W. Fellner, K. Fischer, O. Koepler, J. Kohlhammer, M. Runnwerth, T. Ruppert, T. Schreck, and I. Sens. Visinfo: a digital library system for time series research data based on exploratory search - a user-centered design approach. *Int. J. on Digital Libraries*, 16(1):37–59, 2015. doi: 10.1007/s00799-014-0134-y

[8] D. J. Berndt and J. Clifford. Using dynamic time warping to find patterns in time series. In *KDD workshop*, vol. 10, pp. 359–370. Seattle, WA, 1994.

[9] H. Blumer. Collective behavior. In R. E. Park, ed., *An Outline of Principles of Sociology*, pp. 219–280. Barnes & Noble, New York, 1939.

[10] T. Boren and J. Ramey. Thinking aloud: Reconciling theory and practice. *IEEE transactions on professional communication*, 43(3):261–278, 2000.

[11] S. Brakatsoulas, D. Pfoser, R. Salas, and C. Wenk. On map-matching vehicle tracking data. In *Proceedings of the 31st international conference on Very large data bases*, pp. 853–864. VLDB Endowment, 2005.

[12] K. Buchin, M. Buchin, and J. Gudmundsson. Constrained free space diagrams: a tool for trajectory analysis. *International Journal of Geographical Information Science*, 24(7):1101–1125, 2010.

[13] M. Buchin, S. Dodge, and B. Speckmann. Context-aware similarity of trajectories. In *International Conference on Geographic Information Science*, pp. 43–56. Springer, 2012.

[14] M. Buchin, S. Dodge, and B. Speckmann. Similarity of trajectories taking into account geographic context. *Journal of Spatial Information Science*, 2014(9):101–124, 2014.

[15] D. J. Burr. Elastic matching of line drawings. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-3:708–713, 1981.

[16] D. J. Burr. Designing a handwriting reader. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-5:554–559, 1983.

[17] D. Buzan, S. Sclaroff, and G. Kollios. Extraction and clustering of motion trajectories in video. In *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, vol. 2, pp. 521–524. IEEE, 2004.

[18] J. Chen, R. Wang, L. Liu, and J. Song. Clustering of trajectories based on hausdorff distance. In *Electronics, Communications and Control (ICECC), 2011 International Conference on*, pp. 1940–1944. IEEE, 2011.

[19] L. Chen. *Similarity search over time series and trajectory data*. PhD thesis, University of Waterloo, 2005.

[20] L. Chen and R. Ng. On the marriage of lp-norms and edit distance. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*, pp. 792–803. VLDB Endowment, 2004.

[21] L. Chen, M. T. Özsu, and V. Oria. Robust and fast similarity search for moving object trajectories. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pp. 491–502. ACM, 2005.

[22] M. Y. Choong, R. K. Y. Chin, K. B. Yeo, and K. Tze Kin Teo. Trajectory pattern mining via clustering based on similarity function for transportation surveillance. *International Journal of Simulation–Systems, Science & Technology*, 17(34), 2016.

[23] A. Czirók, M. Vicsek, and T. Vicsek. Collective motion of organisms in three dimensions. *Physica A: Statistical Mechanics and its Applications*, 264(1):299–304, 1999.

[24] S. Dodge, R. Weibel, and E. Forootan. Revealing the physics of movement: Comparing the similarity of movement characteristics of different types of moving objects. *Computers, Environment and Urban Systems*, 33(6):419–434, 2009.

[25] T. Eiter and H. Mannila. Computing discrete fréchet distance. Technical report, Citeseer, 1994.

[26] K. A. Ericsson and H. A. Simon. Protocol analysis, 1984.

[27] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd*, vol. 96, pp. 226–231, 1996.

[28] M. Faundez-Zanuy. On-line signature recognition based on vq-dtw. *Pattern Recognition*, 40(3):981–992, 2007.

[29] A. Godwin and J. T. Stasko. Hotsketch: Drawing police patrol routes among spatiotemporal crime hotspots. In *50th Hawaii International Conference on System Sciences, HICSS 2017, Hilton Waikoloa Village, Hawaii, USA, January 4-7, 2017*, 2017.

[30] P. Grabusts and A. Borisov. Clustering methodology for time series mining. *Scientific Journal of Riga Technical University. Computer Sciences*, 40(1):81–86, 2009.

[31] F. Hausdorff. *Mengenlehre*. Walter de Gruyter Berlin, 1927.

[32] F. Hermans and E. Tsiporkova. Merging microarray cell synchronization experiments through curve alignment. *Bioinformatics*, 23(2):e64–e70, 2007.

[33] I. N. Junejo, O. Javed, and M. Shah. Multi feature path modeling for video surveillance. In *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, vol. 2, pp. 716–719. IEEE, 2004.

[34] E. Keogh and M. Pazzani. Scaling up dynamic time warping to massive datasets. *Principles of Data Mining and Knowledge Discovery*, pp. 1–11, 1999.

[35] E. J. Keogh and M. J. Pazzani. Scaling up dynamic time warping for datamining applications. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 285–289. ACM, 2000.

[36] R. Khan, M. Ahmad, and M. Zakarya. Longest common subsequence based algorithm for measuring similarity between time series: a new approach. *World Applied Sciences Journal*, 24(9):1192–1198, 2013.

[37] P. A. Legg, D. H. S. Chung, M. L. Parry, R. Bown, M. W. Jones, I. W. Griffiths, and M. Chen. Transformation of an uncertain video search pipeline to a sketch-based visual analytics loop. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2109–2118, 2013.

[38] J. Lou, Q. Liu, T. Tan, and W. Hu. Semantic interpretation of object activities in a surveillance system. In *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, vol. 3, pp. 777–780. IEEE, 2002.

[39] G. Navarro. A guided tour to approximate string matching. *ACM computing surveys (CSUR)*, 33(1):31–88, 2001.

[40] W. R. Pearson and D. J. Lipman. Improved tools for biological sequence comparison. *Proceedings of the National Academy of Sciences*, 85(8):2444–2448, 1988.

[41] C. Perin, R. Vuillemot, and J. Fekete. Soccerstories: A kick-off for visual soccer analysis. *IEEE Trans. Vis. Comput. Graph.*, 19(12):2506–2515, 2013. doi: 10.1109/TVCG.2013.192

[42] C. Perin, R. Vuillemot, C. Stolper, J. Stasko, J. Wood, and S. Carpendale. State of the art of sports data visualization. *Computer Graphics Forum*, 2018.

[43] T. M. Rath and R. Manmatha. Word image matching using dynamic time warping. In *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 2, pp. II–II. IEEE, 2003.

[44] C. Reep and B. Benjamin. Skill and chance in association football. *Journal of the Royal Statistical Society. Series A (General)*, 131(4):581–585, 1968.

[45] H. Sakoe and S. Chiba. A dynamic programming approach to continuous speech recognition. In *Proceedings of the Seventh International Congress on Acoustics, Budapest*, vol. 3, pp. 65–69. Akadémiai Kiadó, Budapest, 1971.

[46] Y. Sakurai, M. Yoshikawa, and C. Faloutsos. Ftw: fast similarity search under the time warping distance. In *Proceedings of the twenty-fourth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 326–337. ACM, 2005.

[47] S. Salvador and P. Chan. Toward accurate dynamic time warping in linear time and space. *Intelligent Data Analysis*, 11(5):561–580, 2007.

[48] L. Shao, D. Sacha, B. Neldner, M. Stein, and T. Schreck. Visual-interactive search for soccer trajectories to identify interesting game situations. In *Proc. SPIE Conference on Visualization and Data Analysis*, 2016.

[49] M. S. Shirazi and B. Morris. Vision-based turning movement counting at intersections by cooperating zone and trajectory comparison modules. In *Intelligent Transportation Systems (ITSC), 2014 IEEE 17th International Conference on*, pp. 3100–3105. IEEE, 2014.

[50] S. Sivaraman, B. Morris, and M. Trivedi. Learning multi-lane trajectories using vehicle-based vision. In *Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on*, pp. 2070–2076. IEEE, 2011.

[51] T. Smith and M. Waterman. ªidentification of common molecular subsequences. º j. *Molecular Biology*, 147:195–197, 1981.

[52] E. Sriraghavendra, K. Karthik, and C. Bhattacharyya. Fréchet distance based approach for searching online handwritten documents. In *Document Analysis and Recognition, 2007. ICDAR 2007. Ninth International Conference on*, vol. 1, pp. 461–465. IEEE, 2007.

[53] M. Stein, H. Janetzko, T. Breitkreutz, D. Seebacher, T. Schreck, M. Grossniklaus, I. D. Couzin, and D. A. Keim. Director's cut: Analysis and annotation of soccer matches. *IEEE computer graphics and applications*, 36(5):50–60, 2016.

[54] V. Velichko and N. Zagoruyko. Automatic recognition of 200 words. *International Journal of Man-Machine Studies*, 2(3):223–234, 1970.

[55] M. Vlachos, G. Kollios, and D. Gunopulos. Discovering similar multidimensional trajectories. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pp. 673–684. IEEE, 2002.

[56] R. A. Wagner and M. J. Fischer. The string-to-string correction problem. *Journal of the ACM (JACM)*, 21(1):168–173, 1974.

[57] C. Wenk, R. Salas, and D. Pfoser. Addressing the need for map-matching speed: Localizing global curve-matching algorithms. In *Scientific and Statistical Database Management, 2006. 18th International Conference on*, pp. 379–388. IEEE, 2006.

[58] Z. Yue, H. Broich, and J. Mester. Statistical analysis for the soccer matches of the first bundesliga. *International Journal of Sports Science & Coaching*, 9(3):553–560, 2014. doi: 10.1260/1747-9541.9.3.553

[59] J. Zheng, X. Gao, E. Zhan, and Z. Huang. Algorithm of on-line handwriting signature verification based on discrete fréchet distance. In *International Symposium on Intelligence Computation and Applications*, pp. 461–469. Springer, 2008.

[60] Y. Zhou, S. Yan, and T. S. Huang. Detecting anomaly in videos from trajectory similarity analysis. In *Multimedia and Expo, 2007 IEEE International Conference on*, pp. 1087–1090. IEEE, 2007.