

Using Temporal Convolutional Networks to estimate ball possession in soccer games

Original

Using Temporal Convolutional Networks to estimate ball possession in soccer games / Borghesi, Matteo; Lorenzo D., Costa; Morra, Lia; Lamberti, Fabrizio. - In: EXPERT SYSTEMS WITH APPLICATIONS. - ISSN 0957-4174. - STAMPA. - 223:(2023). [10.1016/j.eswa.2023.119780]

Availability:

This version is available at: 11583/2976470 since: 2023-05-03T12:22:11Z

Publisher:

Elsevier

Published

DOI:10.1016/j.eswa.2023.119780

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

1 Using Temporal Convolutional Networks to Estimate
2 Ball Possession in Soccer Games

3 Matteo Borghesi^a, Lorenzo Dusty Costa^b, Lia Morra^a, Fabrizio Lamberti^a

*^aDepartment of Control and Computer Engineering, Politecnico di Torino, Corso Duca
degli Abruzzi 24, Torino, 10129, Italy*

^bMath&Sport, Via Privata Deruta 20, Milano, 20132, Italy

Email addresses: `s268199@studenti.polito.it` (Matteo Borghesi),
`lorenzo.costa@mathandsport.com` (Lorenzo Dusty Costa), `lia.morra@polito.it` (Lia
Morra), `fabrizio.lamberti@polito.it` (Fabrizio Lamberti)

5 **Abstract**

6 The use of tracking data in the field of sport analytics has increased in the
7 last years as a starting point for in-depth tactical analyses. This work in-
8 vestigates the use of Temporal Convolutional Networks (TCNs), a powerful
9 architecture for sequential data analysis, to extract ball possession informa-
10 tion from tracking data. This task is a crucial step for many tactical analysis
11 and is nowadays carried out manually by a human operator in the stadium,
12 which is costly, difficult to implement, and prone to errors. In this work, sev-
13 eral classification approaches are explored to classify the game state as dead,
14 ball owned by the home team, or by the away team: as a single-branch,
15 ternary prediction, or as two binary predictions, first detecting whether the
16 game is dead or alive and then which team owns the ball. TCNs are ex-
17 ploited to create independent trajectory embeddings from tracking data of
18 each object; since there is no semantic ordering among the tracked objects,
19 we investigate different permutation-invariant layers to combine the embed-
20 dings, namely, an element-wise sum over the embeddings, a self-attention
21 module, and the use of 2D convolutions. Performance evaluation on tracking
22 data from professional soccer games shows that the proposed method out-
23 performs state-of-the-art rule-based methods, achieving 86.2% accuracy in
24 possession estimation (+7.3% compared to the state of the art) and 89.2%
25 accuracy in dead-alive classification (+33.2% compared to the state of the
26 art). Extensive ablation studies were conducted to investigate how different
27 input data concur to the final prediction.

28 *Keywords:* Sport analytics, Deep learning, Tracking data, Ball possession,

30 **1. Introduction**

31 In recent years, the field of sport analytics has received increasingly atten-
 32 tion, as it has been realized that the systematical analysis of the big amount
 33 of data produced by sports daily can help to develop strategies capable of
 34 increasing the chances of winning a match.

35 In this paper, we focus on the automatic extraction of ball possession in-
 36 formation from a soccer game from spatio-temporal tracking data. In typical
 37 soccer analytics pipelines, estimating ball possession and game state is the
 38 first step in understanding the events that occur in a game and their rela-
 39 tionship. Without this information, only physical quantities, e.g., on covered
 40 distance and speeds, can be measured and aggregated. In turn, the avail-
 41 ability of a ball possession estimation component opens to a wider range of
 42 analyses: besides computing simple game state statistics, it becomes possible
 43 to analyze every single pass, to split the game in actions, to classify them
 44 and to study players and team behaviors in attack and defense phases, etc.

45 The need to develop this component stems from considering that raw
 46 tracking data about players and ball positions are today commonly extracted
 47 by specialized companies and made available for the world’s top leagues such
 48 as Premier League, Bundesliga, LaLiga and Serie A, but these companies
 49 still rely on human interventions for the provision of game state information.
 50 In particular, for ball possession the soccer industry uses a definition which
 51 considers the amount of time that a team spends controlling the ball and, to
 52 extract this information, has long relied on a human operator watching the

53 game armed with a three-button timer. The buttons are used to record the
54 beginning of a new game phase, which can be either the home team having
55 possession, the away team having possession, or a stoppage (because the ball
56 is outside the pitch, or the referee has interrupted the game, e.g., after a
57 foul).

58 The reliance on a human operator is motivated by the fact that there are
59 a number of situations where it can be extremely difficult to define clearly
60 which team owns the ball (Bialik, 2014). The need to rely on human op-
61 erators watching and annotating each game, however, clearly has economic
62 and logistic impacts, which are detrimental to the implementation of an au-
63 tomatic data analytics pipeline. Furthermore, it has also been observed that
64 these annotations are prone to errors, which negatively impact the subse-
65 quent data analysis steps (Richly et al., 2017). Another important reason
66 for making the extraction of ball possession automatic is that it would allow
67 to add this information to tracking datasets regarding past matches where
68 it was originally missing, making them accessible for further analyses that
69 otherwise would be extremely time consuming.

70 Scientific literature regarding the automatic estimation of game state
71 is not particularly rich (especially from tracking data), mainly due to the
72 scarcity of public datasets. The few approaches proposed so far generally re-
73 lied on a different definition of ball possession based on the number of passes
74 completed by a team (Glasser, 2014; Sarkar et al., 2019) that is not the com-
75 monly adopted one, or leveraged handcrafted rules (Link and Hoernig, 2017;
76 Morra et al., 2020; Khaustov and Mozgovoy, 2020) that can hardly capture
77 the discrimination abilities of human operators.

78 Hence, in the present paper, we propose to:

- 79 • use deep learning to make a computer learn how to automatically es-
80 timate the state of a soccer game starting from spatio-temporal data
81 about players and ball positions, without resorting to rules defined
82 based on domain knowledge about soccer;
- 83 • output this information in the same format of the standard, three-
84 button timer mechanism.

85 In particular, we investigate the use of Temporal Convolutional Networks
86 (TCNs), which in recent years proved to be particularly effective for deal-
87 ing with classification tasks on sequential data. In this work, TCNs are
88 used to create independent trajectory embeddings from tracking data. We
89 then experimented with three architectures that combine the embeddings
90 in different ways and compared the obtained results with those achieved by
91 state-of-the-art methods leveraging predefined rules.

92 The remaining of the paper is organized as follows. Section 2 reviews
93 relevant literature pertaining sport analytics. Section 3 introduces the pro-
94 posed method, whereas Section 4 presents the protocol that has been set up
95 to evaluate it. Section 5 reports on experimental results. Finally, Section 7
96 provides conclusions and suggests possible directions for future research in
97 this field.

98 **2. Related work**

99 Over the years, sport analytics and, particularly, event detection in soccer
100 games have been addressed in different ways. The existing literature can

101 be roughly classified on the basis of the type of input used, which can be
102 represented by either visual or tracking data (or a combination of them).

103 *2.1. Visual data*

104 Videos indeed represent the most common source of input in the context
105 of sport analytics. Unsurprisingly, in the last several years, most of the
106 research explored the use of deep learning models. Deep learning has been
107 proved successful for many purposes, from players tracking (Kukleva et al.,
108 2019), (Xu et al., 2018), to video summarization (Rockson et al., 2019), (Gao
109 et al., 2020) and the generation of high-level game statistics (Fernández et al.,
110 2019), (Memmert and Rein, 2018), (Theagarajan et al., 2018).

111 Among the possible applications, two tasks that are particularly relevant
112 for the goal of this paper are action recognition and event recognition. For
113 instance, Hong et al. (2018) focused on the possibility to use transfer learn-
114 ing with state-of-the-art Convolutional Neural Network (CNN) models to
115 detect events like corner, free-kick, penalty and goal plus different types of
116 camera shots from soccer videos. Other authors, like Xu and Tasaka (2020),
117 focused on improving the accuracy and speeding up the identification of par-
118 ticular events in 4K multi-view videos of soccer games extending well-known
119 CNN-based object detection and pose estimation methods (such as YOLO
120 (Farhadi, 2016) and OpenPose (Cao et al., 2021)).

121 The most common approach to address the above task on video data is
122 known as Convolutional Recurrent Neural Network (RNN): this approach
123 extends the architecture used in previously cited works since, first, features
124 are extracted from each frame in the video using a CNN, then they are passed
125 to a RNN which produces the output.

126 An example of this setting can be found in Sorano et al. (2020), which
 127 aims at producing a graph of passes that occur in a soccer game. In the
 128 proposed architecture, video frames are processed both by a convolutional
 129 object detection network (YOLO, in this case) and by a feature extraction
 130 network (ResNet18 (He et al., 2016)). For each frame, the feature extraction
 131 module produces a vector describing the whole scene; the object detector, in
 132 turn, is responsible for detecting the players as well as the ball, and returns
 133 a vector describing the position of the ball and the players closest to it.
 134 The two vectors are then concatenated. By processing all the frames in
 135 the video, it is possible to produce a sequence of feature vectors that are
 136 fed into the sequence classification module, which consists of a bidirectional
 137 LSTM (Long Short-Term Memory), a model commonly used in this context.
 138 This module outputs a pass vector that indicates, for each frame of the
 139 original sequence, whether it is part of a pass sequence or not. Another work
 140 exploiting this methodology to detect a larger set of events is represented
 141 by Jiang et al. (2016). Here, play-break segments are first obtained. Then,
 142 semantic features are extracted from them using a CNN. Finally, four event
 143 types are classified (namely, corner, goal, goal attempt and card) using a
 144 RNN.

145 The above approach is also adopted in Roy Tora et al. (2017). In this
 146 case, the focus is on ice hockey, but the task is closer to that addressed in the
 147 present work, as the authors’ goal is to recognize puck possession events. Like
 148 in the above works, the frames are processed in parallel by two CNNs which
 149 extract frame-related features and, based on the output of an object detector,
 150 individual player-related features. The features are then concatenated and

151 passed to an LSTM, which processes them sequentially and produces the
152 output.

153 The main drawback of the methods reviewed so far lays in the fact that
154 they rely on recurrent architectures, which have been proven to be character-
155 ized by a performance bottleneck due to the use of sequential computations.
156 A way to cope with the above limitation when dealing with sequential data
157 is represented by TCNs. These networks rely on convolutional layers, whose
158 operations can be easily parallelized, thus benefiting of continuous advance-
159 ment in computing technology. Bai et al. (2018) and Guirguis et al. (2021),
160 who focused on comparing recurrent architectures with their convolutional
161 counterparts on a variety of tasks, showed the largely higher performance of
162 the latter models in terms of accuracy, as well as of training and inference
163 time.

164 In the context of sport analytics, TCNs have been largely applied to
165 action recognition (the domain they actually stemmed from). An example is
166 provided by Martin et al. (2018), where a Siamese spatio-temporal CNN is
167 used to simultaneously process color images and optical flow data associated
168 with table tennis games to this purpose.

169 These models have also been widely used for event detection, which can
170 be considered as a particular case of action recognition. Some examples in
171 this field are represented, e.g., by Liu et al. (2017), Lee et al. (2018), Yu et al.
172 (2019) and Khan et al. (2018b).

173 The approach of Khan et al. (2018b) is particularly interesting since it
174 uses C3D (Tran et al., 2015), which is basically the three-dimensional (spatio-
175 temporal) counterpart of the well-known two-dimensional VGG network (Si-

176 monyan and Zisserman, 2014) and leverages many of the state-of-the-art
177 characteristics for image classification (such as a high number of layers and
178 small kernels); moreover, the authors showed for the first time how to use
179 such an architecture not only for classifying a video, but also for creating
180 effective descriptors of it, which could be used in a transfer learning pipeline
181 for further analyses.

182 It is worth observing that the problems addressed by the above works
183 are different than that tackled by the present paper. In fact, as stated in
184 Section 1, the task of estimating the game state has not been dealt with in
185 depth by the research community yet.

186 Besides some research done in the field of game state description (ad-
187 dressed in the context of summarization), one of the few works that consid-
188 ered ball possession is represented by Khan et al. (2018a). The authors do
189 not use deep learning end-to-end, since they propose a framework in which
190 the frames of soccer videos are first processed by a Single Shot MultiBox De-
191 tector (SSD)-based object detection module (Liu et al., 2016). The output is
192 then passed to a rule-based system that uses temporal and logical operators,
193 which starts by detecting the so-called “simple” events and assembles them
194 to recognize the “complex” events.

195 2.2. *Tracking data*

196 As shown by the last work reviewed in the previous section, an alternative
197 way to deal with the problem of interest for this paper and, in general,
198 with sport analytics tasks consists of exploiting tracking data. With the
199 improvements in tracking technology, sport researchers are using them for
200 ever more complex tasks, from event detection, to statistics generation, tactic

201 effectiveness quantification, etc.

202 As for video data, the most recent works in this field rely on deep learning,
203 and leverage spatio-temporal convolutions directly on raw data to create low-
204 dimensional representations that summarize the motion of objects of interest
205 in space over time periods. In the literature, these representations are referred
206 to as trajectory embeddings.

207 An important milestone in the use of trajectory embeddings in sport ana-
208 lytics has been set by the work reported in Horton (2020). The author’s goal
209 is to learn an internal feature representation of the movements of all players
210 in a soccer game. To this purpose, a network is designed that takes as input
211 raw trajectory data and learns an internal representation of the individual
212 and coordinated movements of all players. The trajectory of a single player is
213 represented as a sequence of time-stamped frames, and each frame is a vector
214 containing the x and y coordinates for the player at that time, possibly with
215 additional information such as his or her orientation and speed. The main
216 contribution of this work stems from the consideration that most machine
217 learning methods require a predefined structure in the input format that also
218 comprehends an ordering within each input element, such as in the case of
219 an image. However, in the case of tracking data, it is often impossible to
220 define a predefined shape of the input due to the naturally variable duration
221 of a game play, and there is no intrinsic ordering of players in a given interval
222 of play that persists throughout the game or from game to game (in some
223 sports number of player can even change, e.g., due to red cards). Previously,
224 both the variable duration problem and player-ordering problem had been
225 circumvented by introducing a preprocessing step in which raw tracking data

226 are transformed into structured feature representations designed ad hoc for
 227 the task at hand. The method adopted by Horton (2020) avoids the limi-
 228 tations typically associated with feature engineering, and addresses the first
 229 problem by means of 1D convolutions and adaptive pooling mechanisms; it
 230 then deals with the second problem by using a set-based architecture (Za-
 231 heer et al., 2017) that treats the input as an unordered set, devising a model
 232 that learns the feature representation directly from raw data. The authors
 233 used the proposed method to create two models for making predictions about
 234 passes (probability of completion, length, and reception location) and tackles
 235 (probability for a player to be the first to attempt it, distance covered, and
 236 location).

237 Other ways that have been explored to achieve set-based learning in sport
 238 analytics consist in leveraging roles rather than identities for players (e.g.,
 239 when the task is to study the behavior of an entire adversarial team, like in
 240 Lucey et al. (2013)), or in identifying an object, like a player or a ball, that
 241 can be used as “anchor” and define an ordering relative to it. An example of
 242 this latter approach is given in Mehrasa et al. (2018). Like in Horton (2020),
 243 trajectory embeddings are created by 1D convolutions. Then, a permutation-
 244 invariant sorting scheme is defined based on the distance of a candidate object
 245 (a player, in this case) to the anchor, with the trajectory of the anchor being
 246 placed always in the first position, the closest object next to it, and the
 247 farthest object appended to the end. The authors applied this technique
 248 to two different tasks, i.e., event recognition in ice hockey (with six events
 249 considered, and the player carrying the puck acting as the anchor), and team
 250 classification in basketball (with the ball selected as the anchor). It is worth

251 observing that the devised approach based on trajectory data was found to
252 outperform the C3D model that uses video as input, and to be capable of
253 achieving even better performance when used in combination with video.

254 A work that is particularly interesting considering the focus of the present
255 paper is represented by Sanford et al. (2020). The authors address the detec-
256 tion of atomic actions in a soccer game (pass, shot, and reception), and focus
257 on analyzing the performance of vision-based and trajectory-based models.
258 The authors considered four vision-based models. All of them rely on an
259 inflated 3D CNN (I3D) (Carreira and Zisserman, 2017). In the first model,
260 the 3D convolution is applied to the whole image frame. In the other cases, it
261 is applied to players’ “tubelets”, i.e., sequences of bounding boxes containing
262 a single player; the features extracted from the tubelets are then processed
263 in three different ways: via max-aggregation, a Graph Convolutional Net-
264 work (GCN), and a transformer. The best performance was achieved by
265 the model that processes the whole frame, without using the tubelets. For
266 trajectory-based detection, they considered three models, namely a TCN
267 (named Wavenet (van den Oord et al., 2016)), a transformer, and a TCN
268 followed by a transformer. In all three cases, these blocks were followed by
269 a fully connected layer to predict the actual player’s activity. Experiments
270 were run both using only the ball trajectory and using the ball along with
271 the K-nearest players: the model containing both the TCN and the trans-
272 former and using the (five) K-nearest players proved be the one providing
273 the best performance on all three atomic activities. When comparing the
274 two approaches, vision- and trajectory-based models were found to provide,
275 on average, comparable results. The findings of the latter work are particu-

276 larly interesting since they confirm the effectiveness of the trajectory-based
277 method for dealing with sports analytics tasks. They also pinpoint TCNs as
278 the best candidate to address the detection of game events.

279 Works reported above are all relevant to the present paper, since they
280 provide concrete architectures that can be used to perform classification and
281 action detection tasks on tracking data. Notwithstanding, their goal still dif-
282 fers from that considered here, since they cannot be directly used for ball pos-
283 session estimation (although this limitation mainly concerns the last layers of
284 the network, which are described by the authors themselves as task-specific).

285 Like for visual data, the amount of works that focused on ball possession
286 by leveraging tracking data is still quite limited. An example is represented
287 by Link and Hoernig (2017). This work uses a rule-based system to segment
288 the game into possession phases, which are further subdivided into actions
289 and void phases (e.g., when the ball is in the air). A possession phase begins
290 when a new player starts to interact with the ball. Interactions are detected
291 looking at the spikes in the ball acceleration; when a local maximum greater
292 than 4ms^{-2} is found, the possession is assigned to the player closest to the
293 ball. The idea of looking at the derivatives of the ball position comes from
294 the fact that tracking data often do not include the z coordinate; therefore,
295 it is necessary to prevent accidental changes in possession during phases in
296 which the ball crosses intermediate players.

297 A similar idea is exploited in Morra et al. (2020). Here, the temporal
298 and logical operators that were originally used in Khan et al. (2018a) on the
299 output of a visual data processing stage are extended and applied to spatio-
300 temporal data obtained through a soccer game simulator for the detection of

301 game events, including player’s ball possession. In this case, the ball speed is
302 considered, based on the consideration that while a player controls the ball,
303 the latter should move relatively slowly. Furthermore, for a possession to be
304 valid, the distance between the player and the ball should be low for a certain
305 amount of time, and the distance between the opponents and the ball should
306 be above a given threshold.

307 A more recent example of a rule-based system that predicts ball possession
308 from spatio-temporal data is given in Khaustov and Mozgovoy (2020). As in
309 the previous cases, when a possession change occurs, the ball is assigned to
310 the closest player. Possession changes are detected in three ways: through
311 changes in ball speed, changes in ball direction, and prolonged proximity to
312 the ball.

313 In the present work, we address the problem of estimating ball possession
314 from tracking data following a different approach. First, as in the works
315 focusing on action detection reported at the beginning of this section, our
316 aim is to remove the need to rely on handcrafted rules. The objective is
317 to devise models capable to learn directly from data, without resorting to
318 domain knowledge about soccer, which humans use to explain the (possi-
319 ble ambiguous) concept of possession. Second, our expected outcome also
320 slightly deviates from that of the latter works reviewed above. In fact, they
321 actually addressed the ball possession problem in a more fine-grained way, as
322 they estimate the possession on an individual level, telling which player, not
323 only which team, owns the ball. Indeed, this fact is directly connected with
324 the nature of rule-based systems, which follow a bottom-up approach that
325 allows to extract semantic knowledge from the intermediate results. How-

326 ever, we intentionally faced the problem from the point of view of the team
 327 rather than of the player, since, as said, the actual mechanism to collect
 328 ball owner information is based on the three-button timer used by a manual
 329 operator. Data collected through this mechanism only describe the state of
 330 the game (home team controlling the ball, away team controlling the ball,
 331 game stopped), without providing information about the single player who
 332 is owning the ball. Thereby, it is convenient to start with a less fine-grained
 333 approach, and only afterwards add information about the single player on
 334 top of the data obtained for the team.

335 **3. Proposed method**

336 This section illustrates the main principles that underlie our methodol-
 337 ogy. First, the general formulation of the problem statement is presented in
 338 Section 3.1. The three architectures evaluated in this study are then intro-
 339 duced in Section 3.2. Finally, the loss function, TCN design, and aggregation
 340 functions are discussed in detail in Sections 3.3, 3.4, and 3.5, respectively.

341 *3.1. Problem statement*

342 As anticipated in Section 1, our goal is to propose a network architecture
 343 able to classify the game state for a given time window by leveraging tracking
 344 information.

345 Let us define the proposed network as a function $\mathcal{H} : \mathbf{x} \in \mathbb{R}^{n_f \times n_o \times n_c} \mapsto$
 346 $y \in \mathcal{Y}$, where \mathbf{x} is the input tensor and $\mathcal{Y} = \{\text{DEAD}, \text{HOME}, \text{AWAY}\}$ is the
 347 three-class output. The classifier is trained in a supervised fashion from a
 348 labeled dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$.

349 We assume the input tensor to be of size $n_f \times n_o \times n_c$, where n_f is the
350 number of frames in the observed time window, n_o is the number of tracked
351 objects (e.g., players, ball, referee, etc.) and n_c is the number of feature
352 channels associated to each object. For simplicity, and without loss of gen-
353 erality, we assume that the feature channels include at least the position of
354 the tracked objects with respect to the pitch and the team; however, as ex-
355 emplified later in Section 4.1, the feature vector can be extended to include
356 other features such as the player id, the object velocity, visual features, etc.

357 To simplify the explanation of the learning procedure, we decompose the
358 network \mathcal{H} as a combination of three functions

$$\mathcal{H}(\mathbf{x}) = f_c(\Lambda(f_{tcn}(\mathbf{x}))) \quad (1)$$

359 each implemented by one or more layers.

360 The *embedding* function $f_{tcn} : \mathbb{R}^{n_f \times n_o \times n_c} \rightarrow \mathbb{R}^{n_o \times l_{traj}}$ maps the trajectories
361 of each individual object to an embedding vector of length l_{traj} . As detailed
362 later, this component is based on TCNs, hence the subscript.

363 The *aggregation* function Λ combines the embeddings associated with
364 different objects into a single feature vector, which is then given as input to
365 the actual classifier f_c (last function). The order of the players within the
366 input data is based solely on the jersey number of each player within the team;
367 hence, this order does not carry any semantic meaning and needs therefore
368 to be abstracted. It is crucial that, given the same position of the objects on
369 the pitch, the network predicts the same result if two players are swapped,
370 i.e., that the output does not vary in case of a permutation in the input data:
371 hence the need to define a *permutation-invariant* function. An alternative
372 strategy, detailed in Section 3.5.3, is to order the objects according to a

predefined criterion, which bypasses the need to use a permutation-invariant function.

Finally, the last *classification* function f_c is a simple feed-forward network (FFN) that computes the output class c . Alternatively, it is possible to redefine the output space as a combination of two binary labels (y_{DA}, y_{POSS}) , where $y_{DA} \in \{\text{DEAD}, \text{ALIVE}\}$ and $y_{POSS} \in \{\text{HOME}, \text{AWAY}\}$. The peculiarity of this formulation, as discussed in Section 3.3, is that y_{POSS} is not defined when $y_{DA} = \text{DEAD}$.

3.2. High-level architecture

This section explores in detail several variants for each of these three components and their combinations, and introduces the three high-level architectures that were experimentally compared in this work.

All architectures exploit TCNs as the *embedding* function. As discussed in Section 2, according to the recent deep learning literature, TCNs applied to tracking data have proven to work well in different tasks, such as event detection, team classification, etc. They also compared favorably with respect to recurrent architectures (reported, e.g., in Bai et al. (2018) and Guirguis et al. (2021)).

In particular, the proposed architectures are based on $k \times 1$ convolutional filters in order to produce separate trajectory embeddings for each object on the field (in our case, as said, the players, the ball, and the referee).

The first proposed architecture, denoted in the following as the *single-branch* model, frames the problem as a ternary classification. The network, depicted in Fig. 1a consists of three blocks that implement the functions introduced in Section 3.1. In this formulation, $f_c(\cdot)$ is a FFN with three

398 output classes, which takes as input the trajectory embeddings obtained
 399 through spatio-temporal convolution as detailed in the problem statement.
 400 It is important to stress that the aggregation function Λ must be invariant
 401 to permutation; different architectural choices that satisfy this property are
 402 illustrated in Section 3.5.

403 The second class of architectures requires producing two binary classifi-
 404 cations: one telling if the game state is active, the other one telling which
 405 team owns the ball in an active game phase. This leads to an architecture
 406 with two parallel output layers, each responsible for one classification. At
 407 the network level, it is possible to achieve these goals in two ways, outlined
 408 respectively in Fig. 1b and Fig. 1c.

409 The first variant, illustrated in Fig. 1b and denoted in the following as
 410 the *two-branch* network, computes the trajectory embeddings once and uses
 411 them to predict both output variables. The TCN output is passed to two
 412 different Λ layers, which in turn pass their output to two separate FFNs,
 413 the Dead-Alive (DA) branch, and the Possession (POSS) branch. Each FFN
 414 produces a scalar output, representing respectively $P(Y_{DA} = \text{DEAD} \mid X)$ and
 415 $P(Y_{POSS} = \text{HOME} \mid X, Y_{DA} = \text{ALIVE})$. Alternatively, it is also possible to
 416 share both the TCN and the Λ layers, splitting only the FFN network or
 417 part of it. The choice clearly represents a trade-off between computational
 418 needs and flexibility; here, we preferred to keep the Λ layers separated, since
 419 we expect that having distinct representations may be useful to optimize
 420 each classification. It is important to note that both branches are trained
 421 in parallel, i.e., a single backpropagation is performed, and hence the TCN
 422 is trained to jointly optimize both tasks. Parallel training can be achieved

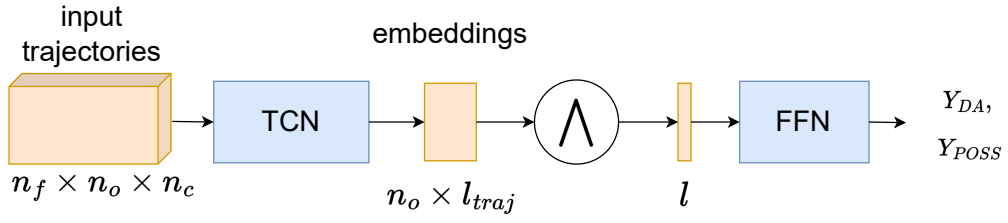
423 using a combined loss function (discussed in Section 3.3) that produces a
 424 single scalar value resulting from both branches.

425 Alternatively, it is possible to perform the classifications by two separate
 426 networks, a Dead-Alive (DA) network and a Possession (POSS) network, as
 427 shown in Fig. 1c. This variant will be denoted in the following as the *two-*
 428 *networks* configuration. In this case, each network computes its trajectory
 429 embeddings that are then passed to the Λ layers and finally to the FFNs for
 430 the binary classification. Computing separate embeddings allows the TCNs
 431 to capture those aspects of the tracking data that may be more relevant
 432 for the specific task, rather than producing a set of general-purpose feature
 433 vectors that are expected to solve both tasks at the same time. The two
 434 networks are trained separately end-to-end, with the possibility of adapting
 435 them to specific task needs, which include using different sets of hyperparam-
 436 eters. The drawback of this alternative is that training two networks requires
 437 roughly twice as much computational resources; this choice is viable only if it
 438 brings about a boost in performance that justifies such investment. Further-
 439 more, the use of separate trajectory embeddings goes against the concept
 440 of embedding as a general descriptor that effectively summarizes the data
 441 and can be used in a wide range of applications, as described in Khan et al.
 442 (2018b).

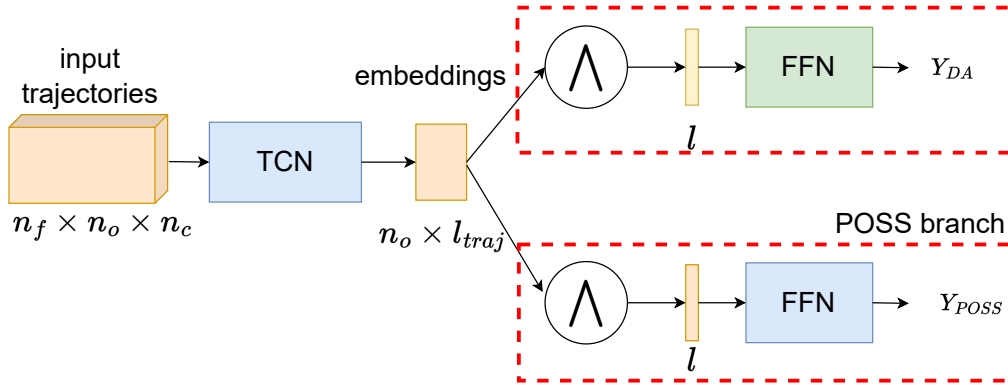
443 3.3. Loss functions

444 The single-branch, multi-class model is trained using a standard cross-
 445 entropy loss written as

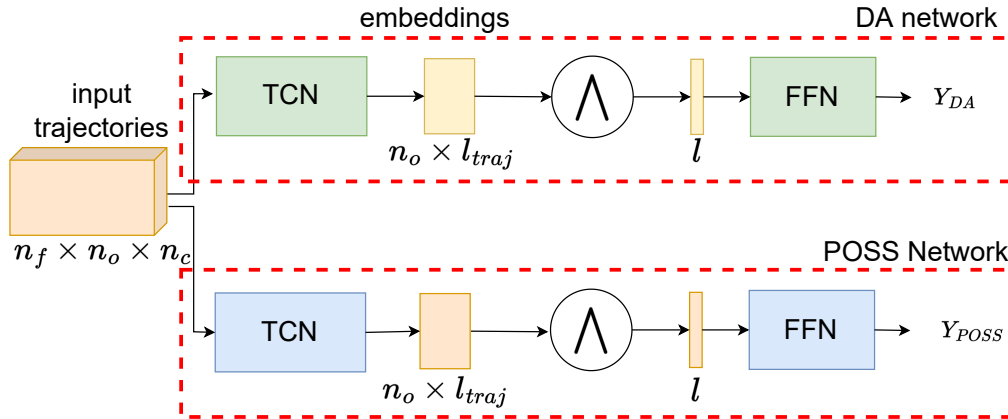
$$L(\hat{y}) = - \sum_i y_i \cdot \log(\hat{y}_i) \quad (2)$$



(a) Single-branch



(b) Two-branches



(c) Two-networks

Figure 1 (*previous page*): Comparison of the three proposed architectures. All take as input a multi-dimensional array of size $n_f \times n_o \times n_c$, where n_f is the number of frames, n_o is the number of objects (including all players, the ball and the referees), and n_c is the number of channels (i.e., features) associated with each object (including, e.g., the position, velocity, team, etc.). All architectures output two scalars representing the probabilities $P(Y_{DA} = \text{DEAD} \mid X)$ and $P(Y_{POSS} = \text{HOME} \mid X, Y_{DA} = \text{ALIVE})$. Each architecture is composed by one or more TCN computing the embeddings (one of each object), a permutation-invariant aggregation function Λ that combines the trajectories of all objects, and finally one or more FFNs f_c that computes the output probabilities. While the single-branch architecture computes both output probabilities using a single TCN and FFN (a), in the two-branch architecture two separate FFN layers are defined on top of a single shared embedding TCN (b). In the two-networks architecture, DA state and ball possession are estimated using separate embedding and classification functions (c).

446 Using a one-hot encoding, y_i is zero for all classes but the correct one; hence,
 447 the cross-entropy loss turns out to be $-\log(\hat{y}_K)$, whereby K is the true class.

448 For the two-networks model, the DA network does not differ substan-
 449 tially from the previous one, except that it performs a binary classification:
 450 however, this can be considered as a special case of multiclass classification,
 451 which allows to use a slightly different cross-entropy function that accounts
 452 for the fact that the network outputs a scalar value instead of a vector. Since
 453 the network predicts the conditional probability of $Y_{DA} = \text{DEAD}$, the true
 454 label $y_{(DA)}$ should be a scalar with value 1 if the game state is DEAD and
 455 0 otherwise. With these modifications, the binary loss function of the DA
 456 network can be expressed as:

$$L_{DA}(\hat{y}_{(DA)}) = -(y_{(DA)} \cdot \log(\hat{y}_{(DA)}) + (1 - y_{(DA)}) \cdot \log(1 - \hat{y}_{(DA)})) \quad (3)$$

457 The issue is more complex when considering the POSS network. The classifi-
 458 cation here does not only depend on the input data X , but also on the value
 459 of Y_{DA} : Y_{POSS} is meaningful only as long as the game is active, otherwise it
 460 is useless to estimate which team owns the ball. During training, this means
 461 that the network should not update its parameters if it is faced with a sample
 462 where the true label is DEAD. To obtain this result, it is possible to define
 463 the loss function as follows:

$$L_{POSS}(\hat{y}_{(POSS)}) = \begin{cases} BCE(\hat{y}_{(POSS)}), & y_{(DA)} = 0 \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

464 where BCE is the binary cross entropy:

$$BCE(\hat{y}_{(POSS)}) = -(y_{(POSS)} \cdot \log(\hat{y}_{(POSS)}) + (1 - y_{(POSS)}) \cdot \log(1 - \hat{y}_{(POSS)})) \quad (5)$$

465 Finally, the two-branch model is trained using a multi-tasking loss func-
 466 tion defined as

$$L(\hat{y}_{(DA)}, \hat{y}_{(POSS)}) = \alpha \cdot L_{DA}(\hat{y}_{(DA)}) + (1 - \alpha) \cdot L_{POSS}(\hat{y}_{(POSS)}) \quad (6)$$

467 i.e., as an average of the two loss functions described above with an addi-
 468 tional weight parameter α . During backpropagation, the derivative of L with
 469 respect to an arbitrary parameter x is given by the formula

$$\nabla_x L = \alpha \cdot \nabla_x L_{DA} + (1 - \alpha) \cdot \nabla_x L_{POSS} \quad (7)$$

470 For the parameters located in the branches, this means that the update
 471 process is the same as in the two-networks model: parameters lying in one

branch do not impact the loss function of the other branch; thus, one of the two members of the derivative above will be zero. With respect to the parameters in the TCN, the update will depend on both losses, according to the weight factor α . In particular, it can be noticed that if a sample belongs to a segment of inactive game, the function L_{POSS} and its derivatives will be zero, which means that the parameters in the TCN are updated based only on the output of the DA branch.

3.4. TCN design

In the proposed architectures, the TCN is responsible for producing trajectory embeddings, i.e., fixed-size representations of the movements on the pitch of every relevant object. This is achieved by stacking several layers of temporal convolutions, which gradually incorporate information from different points in time into a single vector. The structure of the layers defines *a priori* the size of the receptive field, i.e., the number of elements in the sequence that concur to the final prediction. As a result, the receptive field determines how many frames are needed to form an input sample, a parameter that has already been introduced as n_f . In this choice, there should be a trade-off (which has to be made at design time) between two factors: on the one hand, larger sequences allow to consider a larger portion of the game when producing an output; on the other hand, they require a deeper network, which in turn needs more time and more data to be trained.

It is important to note that the target frame, i.e., the frame for which we want to predict the game state, can be located anywhere within the sample: in case the input sequence only includes past frames, the convolution is said to be *causal*, otherwise it is called *acausal*. The choice between these

alternatives depends on how fast the ball possession prediction has to be made; however, it is important to consider that seeing how the action goes on after the target frame can help to enhance the model performance. For example, using only the tracking data, it is difficult to recognize immediately whether a foul was called: in this case, it can be helpful to consider also some frames afterwards, based on the consideration that if a foul is called, the players will probably stop running or move towards the referee. For this reason, unless the system has strict time constraints, it seems appropriate to opt for acausal convolutions.

With respect to these two concepts, it can be useful to point out two aspects pertaining TCNs. First, it is clear that the input slices in two adjacent forward passes have almost the same elements; yet, since they are in different positions, it is not possible to reuse the results of the convolutions from one pass to the other. Second, as shown in Fig. 2, the temporal convolutions are computed on all elements in the sequence, applying dilations and padding when necessary. However, only a small part of the computations (which are shown in the figure with continuous lines) effectively contribute to the output results (the orange circle in the top right). The other computations are needless, since their results are gradually discarded by the following layers.

In order to design the internal structure of the TCN, it is important to recall from the problem statement that it should apply a function f_{tcn} to the input data, such that

$$f_{tcn} : \mathbb{R}^{n_f \times n_o \times n_c} \rightarrow \mathbb{R}^{n_o \times l_{traj}} \quad (8)$$

However, since the trajectory embeddings should be computed separately for each object, f_{tcn} is equivalent to applying on n_o inputs a function f'_{tcn} ,

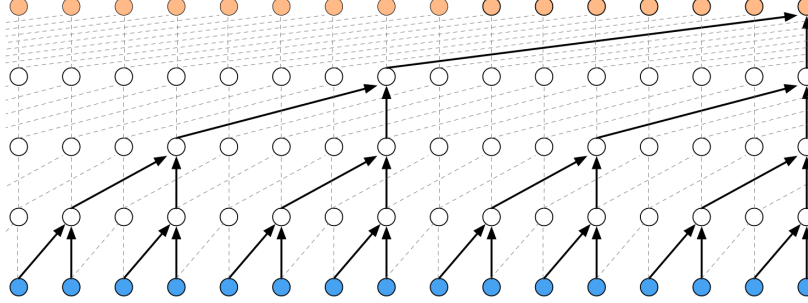


Figure 2: Scheme of dilated convolution (van den Oord et al., 2016): black lines show the convolutions that actually contribute to the result.

such that

$$f'_{tcn} : \mathbb{R}^{n_f \times n_c} \rightarrow \mathbb{R}^{l_{traj}} \quad (9)$$

A function with these characteristics can be achieved using 1D convolutions, i.e., convolutions with a filter of size k and not $k_1 \times k_2$, as in the more common 2D convolutions. At implementation time, it should be considered that filters always have one additional dimension, since they are applied over many channels at the same time; however, this aspect is usually disregarded in the definitions, which explains why they are referred to as 1D convolutions even though the input is two-dimensional. In order to apply f'_{tcn} in parallel on all objects, the most straightforward way is to arrange the operation as a 2D convolution with a $k \times 1$ filter on the whole input, which has size $n_f \times n_o \times n_c$. This technique, proposed by Horton (2020), allows at each step the filter to be convolved with a portion of the input tensor, containing k frames related to only one object. The result of the 2D convolutional layer is a matrix of size $n_o \times l_{traj}$, whose columns correspond to the output of the 1D convolution applied to the respective object. In other words, by means of a $k \times 1$ filter it is possible to compute the function f_{tcn} on the whole input tensor in a single

537 pass.

538 The final structure of the TCN is given in Table 1. The first layer is a 1×1
539 convolution, in order to adapt the third dimension of the input to the size
540 of the final embeddings, which is l_{traj} . Next, a batch normalization layer is
541 applied, as proposed in Ioffe and Szegedy (2015). After that, the architecture
542 features a block containing three layers: the first one is a convolutional layer
543 with a $k \times 1$ filter, which constitutes the most relevant part of the function
544 f_{tcn} . Then, there is a dropout layer and another 1×1 convolution. The block
545 is repeated multiple times (the exact number n_blocks is a hyperparameter
546 of the network) with an exponentially growing dilation rate: as said at the
547 beginning of this section, the number of blocks in the network determines the
548 receptive field of the TCN and, hence, the length of the subsequence consid-
549 ered at each forward pass. Finally, after having applied dropout and batch
550 normalization once again, the last sequence element is selected since, as said,
551 this element captures the whole receptive field, thus offering a summarized
552 representation of the whole temporal sequence.

553 3.5. *Permutation invariance*

554 In the architectures presented in Section 3.2, a major role is played by the
555 aggregation layer Λ , which transforms the individual trajectory embeddings
556 into a global representation of the game sequence, which in turn can then
557 be classified by an FFN. It has already been pointed out that Λ should be
558 permutation-invariant, i.e., the result should be independent of the players’
559 order in the input tensor. In this section, three possible ways to achieve this
560 goal are analyzed, with different characteristics and complexities.

Layer type	Output size	Parameters
input	$n_f \times n_o \times n_c$	-
conv	$n_f \times n_o \times l_{traj}$	filter 1×1
batch_norm	"	-
conv	"	filter $k \times 1$
dropout	"	-
conv	"	filter 1×1
dropout	"	-
batch_norm	"	-
slice	$n_o \times l_{traj}$	-

Table 1: Architecture of the TCN module.

561 3.5.1. Reduce by sum

562 Considering an input matrix A , a simple invariant operation with respect
563 to column permutation is the multiplication $A \cdot \mathbf{1}$, where $\mathbf{1}$ is a vector of all
564 ones. This operation is equivalent to computing a vector whose i -th value
565 is the sum of all the values in the i -th row of A . It is evident that if two
566 columns in the input matrix are swapped, the sum of the values across each
567 row remains unchanged; therefore, the function

$$f_{reduce_sum} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^m, A \mapsto A \cdot \mathbf{1} \quad (10)$$

568 is permutation-invariant. In the current case, the TCN outputs a tensor A_{tcn}
569 of size $n_o \times l_{traj}$. Since the position of the ball and the referee is already fixed
570 by the fact that their tracking data are placed in the first two columns of the
571 dataset, in order to make the Λ layer permutation invariant, it is sufficient to

consider the two submatrices $A_{tcn}^{(home)}$ and $A_{tcn}^{(away)}$ containing the embeddings of the home and away players. The submatrices have size $(11 + 6) \times l_{traj}$, since each soccer team has 11 starting players and up to six substitutions ¹, and can be passed to f_{reduce_sum} obtaining two vectors v_{home} and v_{away} of size l_{traj} that act as trajectory embeddings of one team each.

It is therefore possible to construct a permutation-invariant Λ layer through the linear transformation

$$f_{\Lambda} : \mathbb{R}^{n_o \times l_{traj}} \rightarrow \mathbb{R}^{4 \times l_{traj}}, A_{tcn} \mapsto \begin{pmatrix} | & | & | & | \\ v_{ball} & v_{ref} & v_{home} & v_{away} \\ | & | & | & | \end{pmatrix} \quad (11)$$

where v_{ball} and v_{ref} are the trajectory embeddings of the ball and the referee. It is possible to use this result as input to a FFN by flattening the matrix into a one-dimensional vector of size $n_o \cdot l_{traj}$, as is commonly done in CNNs designed for image classification.

3.5.2. Self-attention

A second possibility in the construction of Λ is to take advantage of recent advances in the field of attention models. In particular, the self-attention module introduced by Vaswani et al. (2017) allows to create embeddings of the original elements that not only take into consideration other elements in the tensor, but are linear combinations of those elements (or more precisely, of their *values*). Notably, in the original self-attention model, the tensor consists of different elements of a sequence, but at this point there are no

¹One additional player is encoded as an extra substitution to account for possible tracking errors

591 sequential data to work with since the temporal information is flattened
592 by the TCN into the trajectory embeddings. Thus, in this case, the self-
593 attention model is not used to process different elements within a temporal
594 sequence, but rather elements that are part of an unordered set, such as the
595 submatrices $A_{tcn}^{(home)}$ and $A_{tcn}^{(away)}$ introduced above.

596 Since the biggest part of the possession estimation is related to one single
597 object, namely the ball, it seems reasonable to think that if the self-attention
598 module is able to grasp all the interactions where the ball is involved, it is
599 possible to make a reliable prediction without considering the interactions
600 among the other objects. At the same time, since self-attention is specifically
601 designed to output a weighted representation of the interactions between the
602 input columns, extracting such a representation of the ball trajectory should
603 provide enough information for a successful classification. Based on these
604 considerations, if the ball trajectory has to be enriched by all the other
605 objects, it is evident that the self-attention layer should receive the whole
606 tensor produced by the TCN: it will be then the self-attention task to rec-
607 ognize which objects have a role in determining the possession and which
608 objects are irrelevant, such as bench players. In this sense, it is relevant to
609 note that the value of a given column in the output of a self-attention layer
610 is independent of the ordering of the other columns. This means that the
611 computation of the column related to the ball is permutation-invariant with
612 respect to the columns related to the players.

613 If we denote by $S \in \mathbb{R}^{n_o \times n_{att}}$ the matrix produced by the self-attention
614 layer and by s_i its column vectors, the operation of the Λ layer can be de-

615 scribed by the function

$$f_{\Lambda} : \mathbb{R}^{n_o \times l_{traj}} \rightarrow \mathbb{R}^{n_{att}}, A_{tcn} \mapsto s_1 \quad (12)$$

616 where n_{att} indicates the size of the *query*, *key* and *value* vectors as defined in
 617 Vaswani et al. (2017).

618 3.5.3. 2D Convolutions

619 A way that is often used to achieve permutation invariance is to impose an
 620 ordering based on an anchor object. In this case, the game state is estimated
 621 based predominantly on the ball: it is possible therefore to order the players
 622 according to their distance to this object, so that the network can operate on
 623 the data independently of how the players are arranged in the input tensor.
 624 Although this idea can also be applied to the options presented above, e.g.,
 625 by limiting the reduce or the self-attention operations to the players close to
 626 the ball, its most powerful consequence is that it allows one to structure the
 627 input tensor in a way that avoids the need to create isolated embeddings for
 628 the objects.

629 The input data can be arranged across two dimensions, which should be
 630 flattened in order to obtain a single prediction of the game state. The two
 631 dimensions are the temporal axis and the different objects, and their sizes
 632 are n_f and n_o , respectively. Structuring the input data allows to operate
 633 on them as commonly done for video streams, where the temporal and the
 634 spatial dimensions are processed in parallel. In other terms, the input data
 635 can be interpreted as two-dimensional, rather than one-dimensional.

636 Taking this fact into account, it is possible to apply a temporal convo-
 637 lution to the input by using a 2D convolutional layer. However, while in

638 Section 3.4 2D convolution was performed by means of a $k \times 1$ filter to sep-
639 arate each object, here $k_1 \times k_2$ filters are used in order to fuse together the
640 information along both axes. This third proposal to obtain a permutation-
641 invariant representation of the global features therefore does not foresee any
642 Λ layer: instead, permutation invariance is a by-product of the TCN design,
643 after introducing the additional constraint that players in the input data are
644 pre-ordered according to their distance from the ball.

645 4. Experiments

646 4.1. Dataset

647 The dataset at our disposal consists of tracking data taken from 35 games
648 during the 2019-20 season of a top professional European league. The data
649 are collected at an average rate of 16 frames per second (fps) and for each
650 frame the following information is provided:

- 651 • a *frame number*, an incremental id of the frame starting from 1;
- 652 • a *game state* label, as described in the problem statement: this is the
653 target variable of the system;
- 654 • a *timestamp*, indicating at which moment the data was collected;
- 655 • a *half* flag, indicating whether the frame was collected in the first or
656 in the second half of the game;
- 657 • the x and y coordinates of the ball;
- 658 • the x and y coordinates of the referee;

659 • for each player, the x and y coordinates, and a flag to distinguish
660 goalkeepers.

661 Ball and players coordinates are provided from a third part company
662 specialized in real time tracking technologies for the sport sector, through a
663 system of ad-hoc cameras installed directly in each venue. The coordinates
664 spaces is a rectangle corresponding to the football pitch and the coordinates
665 system is centered in the center of the pitch (kickoff point) with x-positive
666 axis pointing to the right and y-positive axis pointing up. Considering the
667 standard dimensions of a football pitch (105×68 meters), the range is $[-$
668 $52.5, 52.5]$ for x-axis and $[-34, 34]$ for y-axis both for ball and players. If the
669 tracking system could not locate an object or if the object was not on the
670 pitch (e.g., in the case of a player sitting on the bench or being expelled),
671 the corresponding x and y fields are empty.

672 Target labels DEAD, HOME and AWAY were manually assigned in real
673 time by a human operator as part of a series of services provided by a third
674 company to the league organization. The target labels in the dataset are
675 distributed as follows: about 40% of the samples belong to the class DEAD,
676 the rest of the samples are quite evenly distributed between the classes HOME
677 (29.6%) and AWAY (30.3%).

678 Since the model takes as input a tensor of size $n_f \times n_o \times n_c$, the whole
679 game is split in sequences of length n_f . Clearly, it is also possible for samples
680 to overlap with each other, as the game sequence is transformed into samples
681 following a sliding window approach with a stride of 1 (i.e., adjacent samples
682 differ only by one frame).

683 Datasets acquired during real games often have highly variable quality. A

684 simple yet effective metric to assess data quality is the percentage of samples
 685 in which the ball coordinates are missing. A slightly more informative version
 686 of this metric can be obtained by considering only the samples in which the
 687 game is active, i.e., the game state is not DEAD. The eight games with
 688 the lowest percentage of missing ball coordinates, measured according to
 689 the latter metric, were included in this study. These games were then split
 690 in three subsets of respectively four, two, and two games. From the first
 691 subset, 100K samples were randomly chosen to create the training set; from
 692 the second subset, 5K samples were randomly chosen to create the validation
 693 set; from the third subset 5K samples were randomly chosen to create the test
 694 set. By using different games in each subset, we aimed to have statistically
 695 independent data across the phases. Furthermore, the data in the three
 696 subsets were acquired in different stadiums and with different teams involved
 697 to ensure that the model generalizes well in other contexts.

698 *4.2. Implementation*

699 Each of the alternatives presented in Section 3 is defined by two inde-
 700 pendent factors, namely, the high-level architecture (i.e., single-branch, two-
 701 branch or two-networks configuration) and the permutation-invariant layers
 702 (i.e., reduce by sum, self-attention or 2D convolution). Both factors can be
 703 varied as desired even within the same architecture, e.g., it is possible to de-
 704 fine a two-networks model in which one network uses self-attention and the
 705 other one uses 2D convolution. The only constraint in this sense is that in
 706 a two-branch model, it is not possible to combine a 2D convolution with an-
 707 other permutation-invariant function: in fact, when using 2D convolutions,
 708 the TCN outputs a single vector that is passed to both branches. Therefore,

709 while it is possible in a two-branch network to use a sum layer in one branch
 710 and a self-attention layer in the other one, in the case of 2D convolution the
 711 choice affects necessarily both branches since the TCN is shared by both.

712 As said, each data sample is structured in a three-dimensional tensor
 713 of size $n_f \times n_o \times n_c$. In this work, we set $n_f = 64$, based on experimental
 714 evidence and domain knowledge. The target frame is the 48th element within
 715 the sequence, which means that the model is acausal.

716 The total number of objects n_o is equal to 1 ball+1 referee+22 starting play
 717 ers + 12 substitutions = 36. Padding columns with empty values are added
 718 when the teams did not exploit all possible substitutions. Finally, the $n_c = 11$
 719 channels are defined for each object with the following information:

- 720 • the x and y coordinates;
- 721 • the velocity in the x and y directions, computed by subtracting the
 722 coordinate vector in two adjacent time points and dividing it by the
 723 frame period;
- 724 • three channels encoding the role of the object, i.e., whether it is a ball,
 725 a referee or a goalkeeper;
- 726 • two channels encoding whether the object belongs to the home team
 727 or to the away team;
- 728 • a flag telling whether the object is located outside of the pitch;
- 729 • a flag telling whether the object is missing.

730 Before training, the data are *preprocessed* to ensure training convergence
 731 and reduce the effect of noise. The ball coordinates are first interpolated

732 using the Akima spline (Akima, 1970). Then, each x and y coordinates are
733 separately rescaled using min-max scaling so that they fall into the interval
734 $[-1, 1]$. Missing data are assigned the value -2 , and values far outside of
735 the game field are truncated before scaling in order to provide more stability
736 to normalization.

737 The network architecture has been described in detail in Section 3. The
738 FFN consists of two fully connected layers, with 64 and 32 units, respec-
739 tively. The TCN and the self-attention module are initialized according to
740 the Xavier normalized algorithm, while the FFN initialization follows the
741 Xavier uniform algorithm (Glorot and Bengio, 2010). All layers have the
742 ReLU activation function, except for the FFN layers which use an ELU ac-
743 tivation. All networks were implemented in Python based on Keras v2.4.3 e
744 Tensorflow v2.3.1. For training, the Adam (Kingma and Ba, 2014) optimizer
745 was used, with an initial learning rate of 10^{-5} and a decay rate of 0.7 after
746 each epoch.

747 4.3. Performance assessment

748 The models were evaluated on the basis of three accuracy metrics. First,
749 *global accuracy* is considered, i.e., the percentage of correct predictions among
750 all predictions made within the ternary classification setting presented in the
751 problem statement. Global accuracy can be thus expressed as

$$acc_{global} = \frac{\# \text{correct predictions}}{\# \text{all predictions}} \quad (13)$$

752 Especially for multi-branch and multi-network models, it is also interest-
753 ing to consider two additional metrics, namely *dead-alive accuracy*, acc_{DA} ,
754 and *possession accuracy*, acc_{POSS} . These measures represent the ability of

755 a model to solve one of the two sub-tasks into which the problem can be
 756 decomposed. In particular, the dead-alive accuracy represents the percent-
 757 age of samples for which the model correctly identifies whether the game is
 758 active or not and is computed as

$$acc_{DA} = \frac{\#tp_{DA} + \#tn_{DA}}{\#tp_{DA} + \#tn_{DA} + \#fp_{DA} + \#fn_{DA}} \quad (14)$$

759 where tp_{DA} are the samples for which both the true and predicted labels are
 760 not DEAD, tn_{DA} are the samples for which both the true and predicted labels
 761 are DEAD, fp_{DA} are the samples for which the true label is DEAD while their
 762 predicted label is not DEAD, and fn_{DA} is the opposite case. On the contrary,
 763 the possession accuracy represents the percentage of samples for which the
 764 game is active, and the model correctly identifies which team owns the ball.
 765 It is computed as

$$acc_{POSS} = \frac{\#tp_{POSS} + \#tn_{POSS}}{\#tp_{POSS} + \#tn_{POSS} + \#fp_{POSS} + \#fn_{POSS}} \quad (15)$$

766 where tp_{POSS} are the samples for which both the true and the predicted label
 767 are HOME, tn_{POSS} are the samples for which both the true and the predicted
 768 label are AWAY, fp_{POSS} are the samples for which the true label is AWAY
 769 while their predicted label is HOME, and fn_{POSS} is the opposite case. It is
 770 thus important to note that acc_{POSS} only considers those samples for which
 771 true label is not DEAD, i.e., those frames where the game is active.

772 For the selected architectures, the inference time (mean and standard
 773 deviation) needed to process one batch was calculated. Execution time
 774 was measured on a PC equipped with an NVIDIA 1080Ti GPU with 11Gb
 775 VRAM, 32G RAM and Intel i7-7700 CPU @ 3.60GHz.

776 5. Results

777 The goal of this section is to provide an evaluation of the presented meth-
778 ods. Thus, in Section 5.1, different design alternatives are compared in order
779 to identify the best model to solve the problem statement. This model is
780 then compared in Section 5.2 with other methods taken from the existing
781 literature on related topics. Finally, in Section 5.3, some ablation studies are
782 conducted in order to identify which parts of the model contribute most to
783 the final outcome.

784 5.1. Comparison of design alternatives

785 The results obtained by comparing different design alternatives are shown
786 in Table 2, where each row represents a different combination of architecture
787 and aggregation function. Overall, most of the results are within a small
788 range: the mean accuracy (\pm standard deviation) for all models is $82.93\% \pm$
789 1.71% . On average, the accuracy achieved with single-branch ($83.11\% \pm$
790 1.39%) and two-branch architectures ($83.79\% \pm 1.85\%$) is higher than the
791 two-networks solution ($82.40\% \pm 1.82\%$).

792 We also measured inference time for the best performing network, i.e.,
793 the two-branch network with self-attention aggregation layers. The average
794 time (\pm standard deviation) needed to process one batch is equal to 37.18
795 $\text{ms} \pm 4.66 \text{ ms}$ for a batch size of 1, $57.47 \text{ ms} \pm 5.94 \text{ ms}$ for a batch size
796 of 16, and $52.39 \text{ ms} \pm 1.26 \text{ ms}$ for a batch size of 32. Given that the data
797 is sampled at 16 frames/s, processing times are comparable with real-time
798 inference even on relatively low-performance, consumer-grade GPU.

Architecture	Aggregation function	acc_{global}
Single-branch	sum	83.42 %
	self-att.	84.32 %
	2D-conv.	81.6 %
Two-branch	sum + sum	82.74 %
	self-att. + sum	84.55 %
	sum + self-att.	83.78 %
	self-att. + self-att.	86.39 %
	2D-conv.	81.49 %
Two-networks	sum + sum	82.78 %
	self-att. + sum	84.44 %
	2D-conv. + sum	82.38 %
	sum + self-att.	82.86 %
	self-att. + self-att.	84.32 %
	2D-conv + self-att.	79.42 %
	sum + 2D-conv.	79.62 %
	self-att. + 2D-conv.	82.16 %
	2D-conv + 2D-conv.	83.64 %

Table 2: Performance (global accuracy) of different design alternatives. The first column refers to the high-level architectures, whereas the second column reports the permutation-invariant aggregation function Λ . For multi-branch/multi-network architectures, the first aggregation function refers to the DA branch/network, whereas the second one refers to its POSS counterpart.

799 5.2. Comparison with the state of the art

800 In order to fully assess the contribution of this work, it is important to
 801 provide a quantitative analysis with respect to the state of the art. Since there
 802 are no works that address the overall problem of estimating the game state,
 803 the comparison will be made separately with respect to the two subtasks of
 804 estimating the densities $P(Y_{DA} | X)$ and $P(Y_{POSS} | X, Y_{DA} = \text{ALIVE})$.

805 First of all, the classification between active and inactive game phases is
 806 considered, comparing the model presented in this work with the one from
 807 Wei et al. (2013), which uses a decision tree trained with the ball coordi-
 808 nates only. Each model is tested on 20K samples randomly selected from
 809 two games, chosen among those that were not used to train the neural net-
 810 work. As shown in Table 3 (upper part), the network greatly outperforms
 811 the baseline model, which in turn performs only 6% better than a random
 812 classifier (since it is a binary classification and the classes are relatively bal-
 813 anced, a random classifier has a 50% chance of guessing the correct label).
 814 In Section 5.3 it will be also shown that, even if the network is provided only
 815 with the ball coordinates (thus holding out the players and the referee), it is
 816 still able to achieve 83% accuracy in the dead-alive problem, which is 27%
 817 better than the decision tree.

818 Regarding possession, the current work is compared with three methods,
 819 taken respectively from Link and Hoernig (2017), Morra et al. (2020) and
 820 Khaustov and Mozgovoy (2020). These works propose rule-based systems,
 821 in which possession is estimated starting from considerations drawn from
 822 domain knowledge, regarding, e.g., the closest player to the ball, the speed
 823 and acceleration of the ball, etc. Again, each model is tested on 20K sam-

824 ples randomly selected from two games; the test set is also pruned of those
825 samples where the game is inactive, since the baseline models are designed
826 for estimating ball possession only.

827 All competing models were reimplemented based on the available infor-
828 mation from the original papers. Specifically, in Link and Hoernig (2017),
829 the ball acceleration was computed as a finite difference starting from the
830 ball coordinates. The threshold on the ball acceleration was set to 4ms^{-2} , as
831 proposed by the authors. The minimum distance T_P between the player and
832 the ball, used to discriminate if the player is interacting with the ball, is not
833 provided in the work and was set through validation to 1.5m. In Morra et al.
834 (2020), ball possession is estimated based on the distance from the closest
835 player, the movement of the player and the ball speed, each controlled by a
836 separate threshold. Hyperparameters were taken from the code released by
837 the authors and set to 1.09m, 1.19m, and 8.6ms^{-1} , respectively. As concerns
838 Khaustov and Mozgovoy (2020), the algorithm as well as its hyperparameters
839 are thoroughly listed in the paper and were kept unchanged.

840 The obtained results are listed in the lower part of Table 3, and show
841 that the best performance is achieved by the neural network, with a margin
842 of 7% in accuracy with respect to the best rule-based model, which is the
843 one from Morra et al. (2020).

844 5.3. Ablation studies

845 The goal of this section is to analyze which parts of the input data concur
846 to the final result, in order to understand what aspects are deemed as more
847 important by the network to produce the output, and what is ignored. In
848 particular, ablation studies are performed on two axes: on the one hand,

Solution	acc_{DA}	acc_{POSS}
Ours	89.2%	86.2%
(Wei et al., 2013)	56.0%	-
(Link and Hoernig, 2017)	-	64.5%
(Morra et al., 2020)	-	79.1%
(Khaustov and Mozgovoy, 2020)	-	75.4%

Table 3: Comparison of our best model (two-branch network with self-attention aggregation layers) with the state of the art on the task of dead-alive classification (acc_{DA}) and possession classification (acc_{POSS}).

we evaluate what happens when we remove *objects*, in particular players; on the other hand, we investigate the role of individual *channels*, i.e., of the information related to each object. The two directions are followed separately in an orthogonal way, i.e., when objects are removed, all the channels are considered, and vice versa.

Ablation studies report all three different metrics introduced in Section 4.3. In fact, some objects – or some channels – may be important to determine only one of the two aspects, i.e., only if the game is active or which team owns the ball. The ultimate goal of this analysis is therefore to understand *which parts* of the input are important to produce *which parts* of the output. This is particularly relevant since, as it has been shown above, it is possible to build a model using two separate branches or even two separate networks, each of which performs a binary classification. Knowing which parts of the input data are more important for each prediction enables us to finetune separately the training of each branch/network.

864 Ablation studies are performed on an extended test set which includes
865 20 games, encompassing a larger variety, in terms of acquisition settings and
866 data quality, with respect to the games included in the training set. The two-
867 branch model with self-attention, which achieves the highest global accuracy
868 as reported in Table 2, is selected as baseline.

869 5.3.1. Ablation of objects

870 The object ablation study progressively removes some of the objects from
871 the input data. The input data consist of a tensor of size $n_f \times n_o \times n_c$, where
872 n_o amounts to 36, since it includes the ball, the referee and 17 players from
873 each team (11 starting players and 6 possible incoming players). Performing
874 an ablation study on the objects thus means to cut away a slice of the input
875 on the second axis, passing to the network a tensor of size $n_f \times n'_o \times n_c$, where
876 n'_o is the number of objects that are kept.

877 The ablation is performed in two steps. First, the players far from the
878 ball are removed. The distance can be computed in different ways: here, the
879 Euclidean distance is considered at the frame in which the game state should
880 be estimated. This approach, based on the idea of the K-nearest neighbors
881 (KNN) algorithm, is rather common and can be found in several works from
882 the literature (Sanford et al., 2020) (Mehrasa et al., 2018). In the second
883 step, a more aggressive ablation is performed, and only the ball is retained:
884 the intuition behind this choice is that the ball trajectory, by itself, carries a
885 considerable part of the information.

886 The results in terms of global accuracy are shown in Fig. 3. The blue
887 dots in the figure represent the baseline, which achieves a mean accuracy of
888 81.59% on the test set, as shown in Table 4. The yellow dots refer to the

Ablation	acc_{global}	acc_{DA}	acc_{POSS}
baseline	81.59 %	88.25 %	84.95 %
5NN	79.41 %	85.8 %	84.98 %
ball only	58.35 %	83.62 %	50.54 %
$(x, y) + \text{roles}$	67.25 %	74.89 %	82.09 %
(x, y) only	56.16 %	75.81 %	51.89 %

Table 4: Mean accuracies of different ablation models.

889 model trained using the ball and its five nearest players (5NN) and performs
890 about 2% worse than the baseline. Finally, the red dots show the performance
891 when the model is trained using only the tracking data of the ball: this leads
892 to a considerable drop in the accuracy, since only 58% of the samples are
893 classified correctly on average.

894 Table 4 compares models also with respect to the additional metrics acc_{DA}
895 and acc_{POSS} . The latter presents a similar trend as the global accuracy: the
896 5NN model performs on par with the baseline, while the ball-only model
897 performs significantly worse. On the contrary, in order to estimate if the game
898 is active, it is useful to include all players, since there is a 2.5% difference
899 in acc_{DA} between the 5NN model and the baseline (which ultimately causes
900 the difference in global accuracy). Most interestingly, it can be noticed that
901 the ball trajectory alone is able to achieve a good 83.62% mean acc_{DA} .

902 5.3.2. Ablation of channels

903 Channel ablation studies aims to explain which part of the information
904 about each object are important to produce the output result. In the original

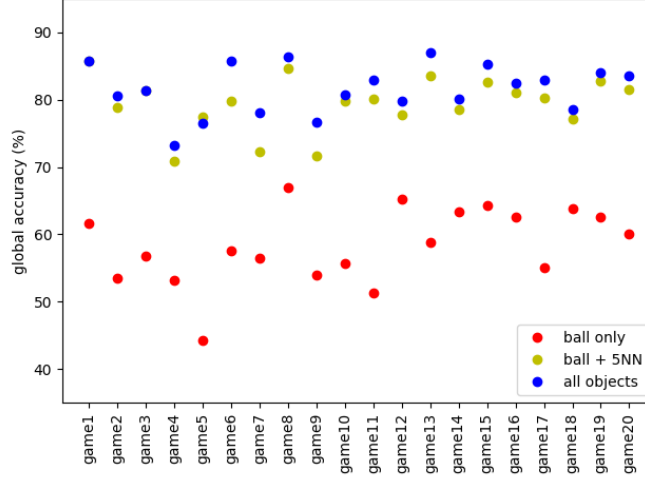


Figure 3: Results of the object ablation study. Each dot represents the global accuracy of the baseline (blue), ball + 5NN (yellow) and ball only (red) models, calculated separately on each game in the extended test set.

input data, 11 channels are passed to the network, including some hand-crafted features, such as velocity, pre-computed in the data pre-processing phase. The goal of this section is therefore to identify which information can be considered as redundant, and whether the designed architecture is capable of automatically encoding or compute features from the raw spatial coordinates, if they are indeed relevant for the classification. Since one of the most relevant characteristics of neural networks is their ability to recognize hidden patterns, which avoids the need to hand engineer the input features as it was typical of the earlier machine learning techniques, we aim to measure up to which point the network is able to fulfill this expectation, and conversely when it is better to provide some explicit information in order to improve the performance.

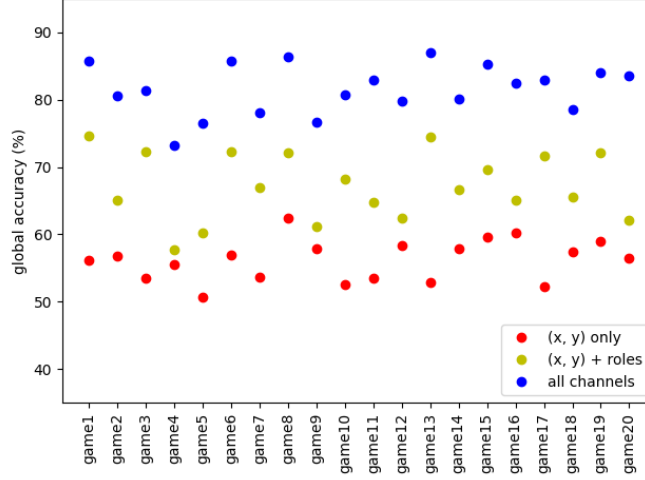


Figure 4: Results of the feature ablation study. Each dot represents the global accuracy of the baseline (blue), (x, y) coordinates + roles (yellow) and (x, y) coordinates only (red) models, calculated separately on each game in the extended test set.

As in the previous section, the ablation is done in two steps: in the first step, information about the coordinates, the roles and the team is kept (in total seven channels), whereas in the second step only the two coordinates are used. The detailed results in terms of global accuracy are shown in Fig. 4, whereas the mean accuracy across the 20 games in the test set are reported in Table 4. The results show a clear difference between the three models: the baseline achieves 81.59% mean accuracy, the model using only the roles has 67.25%, whereas the model that uses only the spatial coordinates has 56.16%. This means that, from a general point of view, all groups of channels make a significant contribution to the output.

When considering separately the accuracy on the two binary classification settings, however, it is possible to note some differences. In fact, in terms

929 of the dead-alive accuracy, the role model (i.e., the first ablation model) has
930 nearly the same performance as the coordinate model (the second ablation
931 model), which indeed performs a slight 1% better. On the contrary, with
932 respect to the possession accuracy, the role model has a performance less
933 than 3% worse than the baseline, whereas the coordinate model achieves as
934 little as 51.9% accuracy.

935 6. Discussion

936 In this paper, we have investigated different TCN architectures to esti-
937 mate the state of a soccer game starting from spatio-temporal data about
938 players and ball positions. All proposed architectures are based on common
939 principles: first, TCNs are employed to map trajectories into an embedding
940 space, and second, the architecture is designed to be permutation-invariant
941 with respect to the orders of the players. However, they differ with respect
942 to other design choices, such as the number of branches, the choice of the
943 permutation-invariant aggregation function, and the loss, which were exper-
944 imentally compared in this paper.

945 With respect to the *global architecture*, the two-networks architecture,
946 in which dead-alive classification and possession estimation are predicted by
947 two separate networks, performs on average worse than those based on a
948 single network. A possible interpretation is that in order to build effective
949 trajectory embeddings, training simultaneously on samples from both active
950 and inactive game phases is more beneficial than having a more flexible net-
951 work with a higher number of parameters. When training on related tasks,
952 multi-task learning can improve performance by promoting implicit regular-

953 ization and more robust feature representation (Ruder, 2017), (Vandenhende
954 et al., 2020). In addition, models consisting of two separate networks may
955 need significantly more resources for both training and inference.

956 Taking into account the trade-off between training time and performance,
957 as well as between memory and performance, the single-branch models achieve
958 results that are often similar or even better than more complex variants.
959 For example, when using 2D convolution, a single classification branch does
960 not perform worse than its two-branch counterpart. From a computational
961 perspective, the processing time of the dual-branch architecture with self-
962 attention is low and even compatible with real-time use. However, it must
963 be stressed that the processing time to extract players and ball tracking data
964 from sensors and/or videos was not considered in the present work. At the
965 same time, many sports analytics pipeline do not require real-time processing
966 capabilities, but rather high accuracy.

967 The choice of the *aggregation function* Λ has a moderate impact on the
968 overall performance. Most of the information can be captured by simple
969 functions, such as summing over all trajectory embeddings. Yet, the best
970 overall performance (86.39% global accuracy) is achieved by the two-branch
971 model using self-attention in both branches: self-attention is the most elab-
972 orated of the three aggregation functions, and allows to capture task-specific
973 features that cannot be recognized otherwise.

974 *Another important aspect to consider is how different input features affect*
975 *the overall performance.* In this case, the input is composed by multiple
976 objects (i.e., the players and the ball), each further characterized by several
977 features (or channels), including the (x, y) coordinates, additional features

978 related to the position (the velocity in the x and y directions, whether the
 979 object is located inside or outside the pitch, and whether it is missing), the
 980 role played by each object, and the team. Both aspects were studied through
 981 extensive ablation studies. In order to globally classify the game state, it is
 982 not possible to consider only the position of the ball, as the accuracy drops
 983 slightly above chance level ($acc_{global}=58.35\%$). However, our ablation studies
 984 show that, on average, it is sufficient to consider the five players closest to the
 985 ball at the beginning of the sequence ($acc_{global}=79.41\%$). It should be noticed
 986 that, because the distance is computed only at one point in the sequence,
 987 samples in which the ball is kicked at the beginning of the sequence could be
 988 misclassified (e.g., in the case of a long pass to an empty area of the pitch,
 989 in which the possession does not change even if the passing player is very far
 990 from the ball at the moment of the evaluation).

991 However, the input information required for each specific task is different.
 992 To determine whether the game state is active or not, the trajectory of the
 993 ball alone achieves a strong performance ($acc_{DA} = 83.62\%$), quite close to
 994 the baseline ($acc_{DA} = 88.25\%$): hence, ball tracking information accounts
 995 for 94% of the information captured by the network that allows to determine
 996 whether the game state is active or not. Removing information about all but
 997 the closest players also reduces the performance by 2.5% ($acc_{DA} = 85.8\%$).
 998 On the contrary, in order to estimate ball possession, it is sufficient to include
 999 the five nearest players ($acc_{POSS} = 84.98\%$) to achieve comparable results
 1000 to the baseline ($acc_{POSS} = 84.95\%$), whereas ball tracking information alone
 1001 cannot reach accuracy above chance level.

1002 Similar considerations apply for the features (channels) associated with

1003 each object. In terms of dead-alive classification, removing velocities and
 1004 position with respect to the external line has a large impact on accuracy.
 1005 In fact, accuracy when using only (x, y) coordinates drops significantly
 1006 ($acc_{DA} = 75.81\%$), and adding role information even slightly degrades perfor-
 1007 mance ($acc_{DA} = 74.89\%$). On the other hand, with respect to possession ac-
 1008 curacy, role information is crucial, whereas velocities and other features play
 1009 a minor role: in fact, a model that takes as input only position and role of
 1010 each object achieves accuracy comparable to the baseline ($acc_{POSS} = 82.09\%$
 1011 vs. $acc_{POSS} = 84.89\%$). Both these insights are in line with intuition: in
 1012 order to tell if the game is active, it is important to know the velocity of the
 1013 objects (e.g., to know if the ball is moving) and if they are inside the pitch,
 1014 whereas to assign the ball possession it is essential to correctly assign each
 1015 object to the proper team.

1016 The results of the ablation studies are consistent with those of the com-
 1017 parison of different architectures. In fact, a two-branch model that uses
 1018 self-attention in both branches would be able to automatically select the
 1019 most relevant features for each task. On the other hand, if a two-networks
 1020 architecture is selected, it would be advisable to tailor the input data passed
 1021 to each network in order to maximize the performance of the system. Like-
 1022 wise, in a two-branch model, since the trajectories are computed separately
 1023 for each object, it is possible to pass only a subset of the embeddings to each
 1024 branch, based on which objects are most important for the classification. For
 1025 example, if only the nearest players are needed to determine Y_{POSS} , it would
 1026 be reasonable to prune the input of the POSS branch in Fig. 1b, selecting
 1027 only the trajectory embeddings related to the objects needed.

1028 Finally, *the proposed model outperforms previously published solutions* on
 1029 both possession accuracy (+7%) (Link and Hoernig, 2017; Morra et al., 2020;
 1030 Khaustov and Mozgovoy, 2020) and game state classification (+27%) (Wei
 1031 et al., 2013). The most recent competing methods (Morra et al., 2020; Khaus-
 1032 tov and Mozgovoy, 2020) are based on rules or temporal logic techniques;
 1033 these methods do not require training, but may include provisions to tune
 1034 rule-specific hyper-parameters (Morra et al., 2020). It is worth noticing that
 1035 all previous techniques were reimplemented and tested on the same dataset to
 1036 ensure a fair comparison; however, hyper-parameters were kept to the original
 1037 values proposed by the authors, and were thus tuned on different datasets,
 1038 at least in one case leveraging synthetic datasets (Morra et al., 2020). The
 1039 comparison offers an interesting insight about the trade-offs present in rule-
 1040 based and deep learning models. On the one hand, handcrafted rules allow
 1041 to build hierarchical models, which can be expanded more easily (e.g., to
 1042 perform event detection) and often have nice by-products, such as the fact
 1043 that the possession estimation is already done at individual level. However,
 1044 this may come at a price in terms of performance, since neural networks
 1045 often present a greater flexibility that allows them to learn more difficult
 1046 mappings. In this case, it is particularly reasonable to opt for a deep learn-
 1047 ing system because the dataset is quite big, which allows to train larger and
 1048 more powerful networks with little impact on generalization.

1049 **7. Conclusions and future work**

1050 This study aimed to devise a deep learning system capable of estimating
 1051 the state of a soccer game on a frame-by-frame basis given a set of spatio-

1052 temporal tracking data. The best performing architecture is a two-branch ar-
1053 chitecture which exploits a TCN backbone to extract trajectory embeddings
1054 for each object/player, and self-attention modules to aggregate embeddings
1055 in a permutation-invariant way. Extensive experimental analysis on track-
1056 ing data from a professional soccer league show that the proposed method
1057 outperforms, by a large amount, state-of-the-art rule-based systems in both
1058 dead-alive classification and ball possession classification.

1059 The present study can be considered as a stepping stone towards automat-
1060 ing a task that presently requires constant human input and supervision. At
1061 the same time, it represents an important contribution to the state of the
1062 art, which currently lacks methods to simultaneously and reliably estimate
1063 ball possession and game state. From a technical point of view, this study
1064 proved that techniques and network architectures that have been success-
1065 fully developed in similar fields, such as event detection, can be applied in
1066 the context of ball possession as well. This work also systematically com-
1067 pares different techniques for achieving permutation invariance on set-based
1068 data, which may be of interest for other applications based on the analysis
1069 of tracking data.

1070 Ample directions for future research emerge from the results of the present
1071 study. For instance, the dataset used in this work is based on cameras pro-
1072 viding only x and y coordinates: improvements in the accuracy of the model
1073 could be achieved by leveraging more advanced systems that provides a very
1074 accurate tracking of the ball, including the z coordinate. Regarding the
1075 methodology, an interesting alternative to the approach adopted here could
1076 be to estimate the game state from a set of events *by subtraction*, i.e., by

1077 detecting all the events that determine a change in the game state, and seg-
1078 menting the game accordingly. In this way, it would be possible to exploit
1079 the large body of existing literature in the field of event detection, as well
1080 as to take one more leap in the direction of an end-to-end deep learning
1081 system capable of analyzing spatio-temporal data. Clearly, this would also
1082 require the availability of a more fine-grained annotated dataset, including
1083 information on the individual players as well as the team in the classification.

1084 Acknowledgements

1085 We would like to thank Ms. Beatrice Gaviraghi, former Math&Sport
1086 project manager, for her consistent support, assistance and organization es-
1087 pecially during the initial phases of this research work.

1088 References

- 1089 Akima, H. (1970). A new method of interpolation and smooth curve fitting
1090 based on local procedures. *Journal of the ACM (JACM)*, 17(4):589–602.
- 1091 Bai, S., Kolter, J. Z., and Koltun, V. (2018). An empirical evaluation of
1092 generic convolutional and recurrent networks for sequence modeling. *arXiv*
1093 *preprint arXiv:1803.01271*.
- 1094 Bialik, C. (2014). The people tracking every touch, pass and tackle in
1095 the world cup. [https://fivethirtyeight.com/features/the-people-](https://fivethirtyeight.com/features/the-people-tracking-every-touch-pass-and-tackle-in-the-world-cup/)
1096 [tracking-every-touch-pass-and-tackle-in-the-world-cup/](https://fivethirtyeight.com/features/the-people-tracking-every-touch-pass-and-tackle-in-the-world-cup/). Ac-
1097 cessed: 2021-06-25.

- 1098 Cao, Z., Hidalgo, G., Simon, T., Wei, S.-E., and Sheikh, Y. (2021). Openpose:
1099 Realtime multi-person 2D pose estimation using part affinity fields. *IEEE*
1100 *Transactions on Pattern Analysis and Machine Intelligence*, 43(1):172–
1101 186.
- 1102 Carreira, J. and Zisserman, A. (2017). Quo vadis, action recognition? A new
1103 model and the kinetics dataset. In *Proc. IEEE Conference on Computer*
1104 *Vision and Pattern Recognition*, pages 6299—6308.
- 1105 Farhadi, J. R. S. D. R. G. A. (2016). You Only Look Once: Unified, real-
1106 time object detection. In *Proc. IEEE Conference on Computer Vision and*
1107 *Pattern Recognition*, pages 779–788.
- 1108 Fernández, J., Bornn, L., and Cervone, D. (2019). Decomposing the immea-
1109 surable sport: A deep learning expected possession value framework for
1110 soccer. In *Proc. 13 th Annual MIT Sloan Sports Analytics Conference*.
- 1111 Gao, X., Liu, X., Yang, T., Deng, G., Peng, H., Zhang, Q., Li, H., and Liu, J.
1112 (2020). Automatic key moment extraction and highlights generation based
1113 on comprehensive soccer video understanding. In *Proc. IEEE International*
1114 *Conference on Multimedia and Expo Workshops*, pages 1—6.
- 1115 Glasser, H. (2014). The problem with possession: The inside story of soccers
1116 most controversial stat. [https://slate.com/culture/2014/06/soccer-](https://slate.com/culture/2014/06/soccer-possession-the-inside-story-of-the-games-most-controversial-stat.html)
1117 [possession-the-inside-story-of-the-games-most-controversial-](https://slate.com/culture/2014/06/soccer-possession-the-inside-story-of-the-games-most-controversial-stat.html)
1118 [stat.html](https://slate.com/culture/2014/06/soccer-possession-the-inside-story-of-the-games-most-controversial-stat.html). Accessed: 2021-06-25.
- 1119 Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training
1120 deep feedforward neural networks. In *Proceedings of the thirteenth inter-*

1121 *national conference on artificial intelligence and statistics*, pages 249–256.
1122 JMLR Workshop and Conference Proceedings.

1123 Guirguis, K., Schorn, C., Guntoro, A., Abdulatif, S., and Yang, B. (2021).
1124 SELD-TCN: sound event localization & detection via temporal convolu-
1125 tional networks. In *2020 28th European Signal Processing Conference (EU-*
1126 *SIPCO)*, pages 16–20. IEEE.

1127 He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning
1128 for image recognition. In *Proc. IEEE conference on computer vision and*
1129 *pattern recognition*, pages 770—778.

1130 Hong, Y., Ling, C., and Ye, Z. (2018). End-to-end soccer video scene and
1131 event classification with deep transfer learning. In *Proc. 2018 International*
1132 *Conference on Intelligent Systems and Computer Vision*, pages 1—4.

1133 Horton, M. (2020). Learning feature representations from football tracking.
1134 In *Proc. MIT Sloan Sports Analytics Conference*.

1135 Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep
1136 network training by reducing internal covariate shift. *arXiv preprint*
1137 *arXiv:1502.03167*.

1138 Jiang, H., Lu, Y., and Xue, J. (2016). Automatic soccer video event detection
1139 based on a deep neural network combined CNN and RNN. In *Proc. IEEE*
1140 *28th International Conference on Tools with Artificial Intelligence*, pages
1141 490—494.

1142 Khan, A., Lazzerini, B., Calabrese, G., and Serafini, L. (2018a). Soccer event

1143 detection. In *Proc. 4th International Conference on Image Processing and*
1144 *Pattern Recognition*, pages 119—129.

1145 Khan, M. Z., Saleem, S., Hassan, M. A., and Khan, M. U. G. (2018b).
1146 Learning deep C3D features for soccer video event detection. In *2018 14th*
1147 *International Conference on Emerging Technologies (ICET)*, pages 1–6.
1148 IEEE.

1149 Khaustov, V. and Mozgovoy, M. (2020). Recognizing events in spatiotempo-
1150 ral soccer data. *Applied Sciences*, 10(22):8046.

1151 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimiza-
1152 tion. *arXiv preprint arXiv:1412.6980*.

1153 Kukleva, A., Asif, M., Hafez Farazi, K., and Behnke, S. (2019). Utilizing
1154 temporal information in deep convolutional network for efficient soccer
1155 ball detection and tracking. In *Robot World Cup*, pages 112—125.

1156 Lee, J., Kim, Y., Jeong, M., Kim, C., Nam, D.-W., Lee, J., Moon, S., and
1157 Yoo, W. (2018). 3D convolutional neural networks for soccer object motion
1158 recognition. In *Proc. 20th International Conference on Advanced Commu-*
1159 *nication Technology*, pages 354—358.

1160 Link, D. and Hoernig, M. (2017). Individual ball possession in soccer. *PLoS*
1161 *One*, 12(7).

1162 Liu, T., Lu, Y., Lei, X., Zhang, L., Wang, H., Huang, W., and Wang, Z.
1163 (2017). Soccer video event detection using 3d convolutional networks and
1164 shot boundary detection via deep feature distance. In *Proc. International*
1165 *Conference on Neural Information Processing*, pages 440—449.

- 1166 Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Scott, R., Fu, C.-Y., and
1167 Berg, A. C. (2016). SSD: Single Shot MultiBox Detector. In *Proc. 14th*
1168 *European Conference on Computer Vision*, pages 21–37.
- 1169 Lucey, P., Bialkowski, A., Carr, P., Morgan, S., Matthews, I., and Sheikh,
1170 Y. (2013). Representing and discovering adversarial team behaviors using
1171 player roles. In *Proc. IEEE Conference on Computer Vision and Pattern*
1172 *Recognition*, pages 2706—2713.
- 1173 Martin, P.-E., Benois-Pineau, J., Péteri, R., and Morlier, J. (2018). Sport ac-
1174 tion recognition with siamese spatio-temporal CNNs: Application to table
1175 tennis. In *Proc. International Conference on Content-Based Multimedia*
1176 *Indexing*, pages 1—6.
- 1177 Mehrasa, N., Zhong, Y., Tung, F., Bornn, L., and Mori, G. (2018). Deep
1178 learning of player trajectory representations for team activity analysis. In
1179 *11th MIT Sloan Sports Analytics Conference*, volume 2, page 3.
- 1180 Memmert, D. and Rein, R. (2018). Match analysis, big data and tactics: Cur-
1181 rent trends in elite soccer. *German Journal of Sports Medicine*, 69(3):65–
1182 71.
- 1183 Morra, L., Manigrasso, F., Canto, G., Gianfrate, C., Guarino, E., and Lam-
1184 berti, F. (2020). Slicing and dicing soccer: Automatic detection of complex
1185 events from spatio-temporal data. In *Proc. International Conference on*
1186 *Image Analysis and Recognition*, pages 107—121.
- 1187 Richly, K., Moritz, F., and Schwarz, C. (2017). Utilizing artificial neural
1188 networks to detect compound events in spatio-temporal soccer data. In

- 1189 *Proc. 3rd SIGKDD Workshop on Mining and Learning from Time Series*,
1190 pages 13—17.
- 1191 Rockson, A., Rafiq, M., and Gyu Sang, C. (2019). Soccer video summa-
1192 rization using deep learning. In *Proc. IEEE Conference on Multimedia*
1193 *Information Processing and Retrieval*, pages 270—273.
- 1194 Roy Tora, M., Chen, J., and Little, J. J. (2017). Classification of puck
1195 possession events in ice hockey. In *IEEE Conference on Computer Vision*
1196 *and Pattern Pecognition Workshops*, page 147–154.
- 1197 Ruder, S. (2017). An overview of multi-task learning in deep neural networks.
1198 *arXiv preprint arXiv:1706.05098*.
- 1199 Sanford, R., Gorji, S., Hafemann, L. G., Pourbabae, B., and Javan, M.
1200 (2020). Group activity detection from trajectory and video data in soccer.
1201 In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
1202 *Pattern Recognition Workshops*, pages 898–899.
- 1203 Sarkar, S., Chakrabarti, A., and Mukherjee, D. P. (2019). Generation of ball
1204 possession statistics in soccer using minimum-cost flow network. In *Proc.*
1205 *IEEE Conference on Computer Vision and Pattern Recognition Work-*
1206 *shops*, pages 2515–2523.
- 1207 Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks
1208 for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- 1209 Sorano, D., Carrara, F., Cintia, P., Falchi, F., and Pappalardo, L. (2020).
1210 Automatic pass annotation from soccer video streams based on object
1211 detection and LSTM. *arXiv preprint arXiv:2007.06475*.

- 1212 Theagarajan, R., Pala, F., Zhang, X., and Bhanu, B. (2018). Soccer: Who
1213 has the ball? Generating visual analytics and player statistics. In *Proc.*
1214 *IEEE Conference on Computer Vision and Pattern Recognition Work-*
1215 *shops*, page 1749–1757.
- 1216 Tran, D., Bourdev, L., Fergus, R., Torresani, L., and Paluri, M. (2015).
1217 Learning spatiotemporal features with 3D convolutional networks. In *Proc.*
1218 *IEEE international conference on computer vision. 2015*, pages 4489–
1219 –4497.
- 1220 van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves,
1221 A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). Wavenet: A
1222 generative model for raw audio. In *9th ISCA Speech Synthesis Workshop*,
1223 pages 125–125.
- 1224 Vandenhende, S., Georgoulis, S., Proesmans, M., Dai, D., and Van Gool,
1225 L. (2020). Revisiting multi-task learning in the deep learning era. *arXiv*
1226 *preprint arXiv:2004.13379*.
- 1227 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N.,
1228 Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. *Advances*
1229 *in neural information processing systems*, 30:5998–6008.
- 1230 Wei, X., Sha, L., Lucey, P., Morgan, S., and Sridharan, S. (2013). Large-
1231 scale analysis of formations in soccer. In *2013 international conference on*
1232 *digital image computing: techniques and applications (DICTA)*, pages 1–8.
1233 IEEE.

- 1234 Xu, J., Kanokphan, L., and Tasaka, K. (2018). Fast and accurate object
1235 detection using image cropping/resizing in multi-view 4k sports videos.
1236 In *Proc. 1st International Workshop on Multimedia Content Analysis in*
1237 *Sports*, pages 97—103.
- 1238 Xu, J. and Tasaka, K. (2020). Keep your eye on the ball: Detection of
1239 kicking motions in multi-view 4K soccer videos. *ITE Transactions on*
1240 *Media Technology and Applications*, 8(2):81—88.
- 1241 Yu, J., Lei, A., and Hu, Y. (2019). Soccer video event detection based on
1242 deep learning. In *Proc. International Conference on Multimedia Modeling*,
1243 page pp. 377–389.
- 1244 Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., R., S., and J., S. A.
1245 (2017). Deep Sets. In *Proc. Annual Conference on Neural Information*
1246 *Processing Systems*.