# Progress in RoboCup Revisited:
# The State of Soccer Simulation 2D

Thomas Gabel, Egbert Falkenberg, Eicke Godehardt

Faculty of Computer Science and Engineering
Frankfurt University of Applied Sciences
60318 Frankfurt am Main, Germany
{tgabel|falken|godehardt}@fb2.fra-uas.de

**Abstract.** A remarkable feature of RoboCup's soccer simulation leagues is their ability to quantify and prove the exact progress made over years. In this paper, we present and discuss the results of an extensive empirical study of the progress and the currently reached state of 2D soccer simulation. Our main finding is that the current decade has witnessed a continuous and statistically significant improvement of the overall level of play, but that the magnitude of the progress made has dropped clearly when compared to the previous decade. In accordance to this, we envision possible future prospects for the 2D league that might respond to our empirical findings.

## 1 Introduction

At RoboCup 2015, two-dimensional simulated soccer players competed with one another at a world championship tournament for the 20th time. The long history of this competition as well as the continued interest of the community in the 2D soccer simulation league made us ask the question what progress has been made in this league throughout the years. We answered this question quantitatively at the RoboCup Symposium 2010 [3] by presenting the results of an extensive empirical evaluation with which we measured the progress of playing performance within the time window from 2003 to 2007.

With eight years having passed since the end time of the interval considered in the study mentioned, we think it is time to revisit the 2D simulation league and to pose the same as well as further questions. For those questions to be answerable, however, we need a stable soccer simulation platform. In [3], we had addressed the time period from 2003 to 2007 during which there was such a stable platform since the 2D league's simulation software, the Soccer Server [6], underwent no changes. In the paper at hand, our focus is on analyzing and assessing the developments during the more recent time interval (2010 to 2015) where also no modifications were made to the platform. Beyond this, we intend to compare our current findings to the results of the first stable period (2003-2007), identify similarities and differences and draw corresponding conclusions. Finally, we are going to address the question what are the recommendations,

avenues, and prospects for the further development of the soccer simulator and the 2D league as a whole given the experience and the results reported.

In Section 2, we provide necessary background information on soccer simulation and its course of development during the two recent decades. We also raise a number of questions that shall be addressed with our studies and outline our experimental setup. Section 3 presents in depth the results of our evaluations and in Section 4 we aim at drawing conclusions from our findings.

## 2 Background

Researchers and students who have been active in the RoboCup domain for more than a single year, will easily come to the conclusion that the overall level of play at RoboCup tournaments is increasing gradually. While this observation will certainly be shared by both, participants as well as spectators, it is difficult to quantitatively prove its correctness. To this end, the soccer simulation leagues adopt a special role because no hardware development and maintenance are necessary, but instead soccer-playing agents as well as belonging coaches are merely some pieces of software [1]. This allows for repetitive and detailed evaluations as well as for forming quantitative statements of a team's strengths and weaknesses [2]. But it also allows for analyses of simulated soccer teams across various years. Given these circumstances, we are in the lucky position to derive empirically grounded statements regarding the quantitative playing strength of teams from different years and, in so doing, come up with evidence for or against a significant progress of RoboCup's soccer simulation branch.

### 2.1 Periods of Stability

In the 2D Soccer Simulation League all competitions are based on the Soccer Server software [6] which implements soccer playing as a completely distributed multi-agent system in a two-dimensional plane while adhering to the official soccer rules to the largest degree possible. During the 20 years of its existence this simulator has gone through various extensions and changes (see [1, 3] for an overview), but it has also experienced periods of stability, i.e. a number of successive years where the technical and maintenance committees in agreement with the soccer simulation community decided to introduce no changes (apart from bug fixes) and, hence, to keep the simulation platform stable.

These periods of stability are of special interest in the scope of this paper, since a stable platform is a fundamental prerequisite for doing analyses, experiments, and evaluations with published soccer team binaries from different years. Stated differently, a non-stable platform (e.g. due to the introduction of a new feature into the simulation) prevents us from performing a meaningful and fair comparison of teams that were developed for different versions of the simulator.

Figure 1 shows the alternating periods of stability and further-development of the Soccer Server starting from the Pre-RoboCup event at IROS 1996 in Osaka. After seven years of intensive development, in 2003 the Soccer Server went into
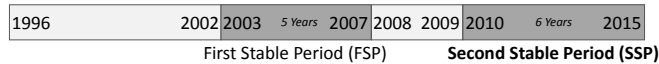
| 1996 | | 2002 | 2003 | *5 Years* | 2007 | 2008 | 2009 | 2010 | *6 Years* | 2015 |

First Stable Period (FSP)      **Second Stable Period (SSP)**

**Fig. 1.** Alternating Periods of Frozen and Continued Soccer Server Development

its *first stable period* (FSP) covering five years. This happened in parallel to the establishment of the 3D Soccer Simulation League. When, however, the 3D league moved to modeling humanoid robots instead of simple spheres in 2007, the 2D committees decided to end the stable period and introduce new features and significant changes to the 2D simulation, which resulted in new simulator versions with a changed feature set for the 2008, 2009, and 2010 competitions. Since 2010, however, the simulation has been kept stable once again such that we can now speak of a *second stable period* (SSP, six years so far, 2010-2015 and lasting) with full compatibility of all team binaries released with these years.

### 2.2 Experimental Goals and Setup

*Goals* Given the fact that the 2D Soccer Simulation League has seen two stable periods, within which a meaningful and fair evaluation across years is possible, we are going to answer the following questions in the remainder of this paper:

1. We address the question whether further progress has been made in soccer simulation 2D throughout the years of the SSP. This is similar to our analyses for the FSP (2003-2007) published in [3].
2. We are going to compare the progress indicators of both stable periods and draw corresponding conclusions.
3. We want to analyze the current state of the 2D Soccer Simulation League more thoroughly by focusing on individual teams' further development during the recent years.

*Platform* Within the SSP, Soccer Server versions from 14.0.3 to 15.2.2 were used. As already said, any changes introduced in this period were targeted solely at bug fixing and installation support. In total, there have been only two minor modifications that might interfere with older teams and hence might have an impact on our experiments. First, the drop ball time (e.g. maximal time to execute a corner kick) was cut to half in 2011 in order to push back time wasting. Second, a defect in the coach language implementation has been corrected which could be exploited by older teams to their advantage. In order to compensate these corrections and to enforce maximal compatibility with all teams from 2010 to 2015 we conducted all our experiments using server version 15.0.0.

*Team Selection* As in our 2010 study for the FSP [3], we proceed on the assumption that a year's joint level of play can be read from the quality of its top representatives. In Section 3.3 we will empirically underpin the validity of this assumption. In order to enforce comparability with the results of the mentioned study for the FSP, we again let any year from 2010 to 2015 be represented by its top four teams, i.e. those teams that made it to the semi-finals in the

respective year's RoboCup world championships. Thus, in total we utilized 24 officially released team binaries from the six years of the SSP. When addressing the first two of the three questions posed above, we refrain from naming teams by their actual names, but use "year_place" identifiers for better readability. The matching from identifiers to verbose team names can be found in the Appendix.

*Experiments* We retrieved all considered historic team binaries from the web archive[1], made them work in our evaluation setting, and let them play multiple matches against one another using Soccer Server version 15.0.0. All matches were performed on a cluster of identical machines. Since we allowed any team to face any other team at least fifty times, we had in total more than 17.000 matches which corresponds to approximately 120 days of simulated soccer.

## 3 Empirical Results

In accordance to the research questions posed above we divided our experiments into three parts. We start by presenting the results of our general progress analysis that focuses on soccer simulation's further development within the SSP. We then compare these findings to the progress that the 2D league had made within its FSP and, finally, we more critically question how to characterize its currently reached state and by which driving factors it has emerged.

### 3.1 General Progress Analysis

In this part of our study we allowed each of the 24 representatives of the SSP to face each other team 50 times under regular settings. We determined for each team the average score (number of goals shot vs. number of goals received) from these $23 \times 50 = 1150$ matches and plotted this information in the left chart of Figure 2 (standard deviations are omitted for readability). Apparently, teams from later years score on average more goals while they receive less. Furthermore, polygons formed by interconnecting data points of the same year seem to shift slightly to the bottom right of the diagram. While this observation is less obvious than in the similar chart for the FSP from 2003-2007 [3], it nevertheless represents a first indication of the fact that there has been some progress in soccer simulation 2D during the last six years.

The accompanying bar chart in the right of Figure 2 shows the share of points each team obtained within this set of experiments, where as in human soccer a victory is awarded with three points, a draw with one and a defeat with zero points. Thus, winning all of the 1150 matches would yield 100% (3450) while drawing all matches would yield 33.3%. With 82.1% of the maximal number of points the 2015 champion turned out to be the strongest team. Interestingly, teams originating from later years (bars with lighter shades of gray) place predominantly in the top half of the ranking whereas older binaries (darker shades of gray) are to be found mainly on the rear ranks.
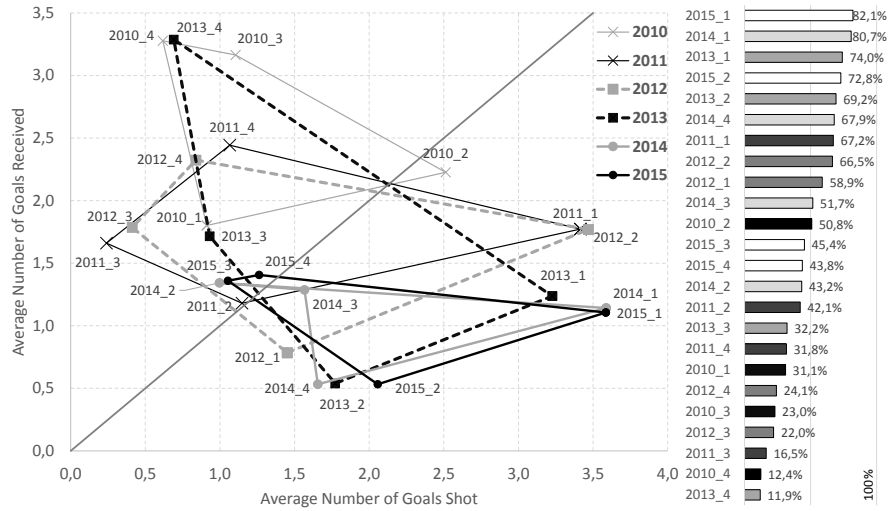
---

[1] http://chaosscripting.net/

**Fig. 2.** Left: Average scores for each of the 24 representatives from the second stable period when playing repeatedly against all other teams. Right: Corresponding share of points achieved by each team within this empirical study.

While we have concentrated on individual teams so far, we now focus on a comparison of entire years where, as stressed in Section 2, we let the joint strength of a year be represented by its four top teams. In Figure 3, we visualize the results of this year-vs-year comparison in a matrix-like representation. Each cell of this matrix shows the results of a large-scale tournament (of 800 matches each) where all representatives from the year indicated by the row played against all teams from the year indicated by the column. The bars within each cell show both, the distribution of points among both years as well as the average scores. For all matrix elements $(y_1, y_2)$ with data it holds that $y_1 < y_2$, i.e. we have the "newer" teams in the columns and we visualize it with a darker color. Therefore, we can easily read from the chart, that in any constellation $y_2$ performed superior to $y_1$. The larger the difference $y_2 - y_1$, the clearer the dominance of the newer team. For $y_1 = 2014$ and $y_2 = 2015$, however, we find that both years are nearly equally strong (points are distributed as 49.7 : 50.3). So, the first general conclusion from the mentioned evaluations are that (a) there has been substantial progress within soccer simulation 2D during the six years of the SSP and (b) that this kind of progress has slowed down recently.

### 3.2 Comparative Progress Analysis

The decision to freeze the Soccer Server development in 2003 and, thus, to enter a stable period had not been taken recklessly back in 2002 (see [3] for background information). The empirical study [3] on what happened within this period of stability has emphasized to what extent the 2D community has drawn benefits in
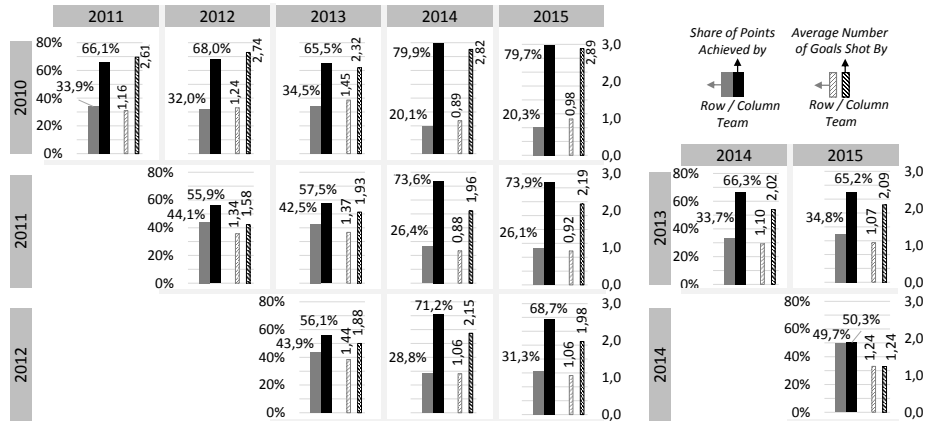
**Fig. 3.** Representatives from one year (row index) played multiple times against all teams from any other years (column index) of the SSP. Distribution of points (3/1/0) as well as average scores (over 800 matches) are plotted for each combination.

terms of increasing overall playing performance. Given the observations reported in Section 3.1, a highly interesting question to ask is whether the impact of the SSP has been as substantial as the impact of the first one.

To answer this question we now no longer target individual top teams or year-vs-year comparisons, but instead let the joint set of representatives of any year play repeatedly against *all* teams from all other years. For example, we made each of the four 2010 representatives play 50 times against all top teams from 2011 through 2015, which made a total of 4000 matches for each year considered. In the right part of Figure 4 we present the averaged scores over these 4000 matches for the six years from 2010 to 2015 including the score quotient (average number of goals scored divided by the average number of goals received). The belonging chart in the left of this figure contrasts these results with the results of the FSP. Apparently, the progress made in that time period has been much clearer because the steps from year to year were larger in terms of scoring more goals while receiving less on average. Furthermore, when comparing the starting point of the FSP (2003) with the average results achieved four years later (2007), we have to acknowledge that in this stable period the average number of goals shot has been increased by a factor of 3.27 and the goals received has been reduced by factor 8.56. By contrast, from 2010 to 2014 (also after four years of time) these increment/reduction factors are only 1.78 and 2.59, respectively. These facts are indicative of a clear decline in the progress being made in soccer simulation 2D even under stable conditions.

With 4000 samples per year group, we can quite reliably state whether a change (measured by the average score) from one year $y_1$ to a second one $y_2$ is statistically significant or not, assuming a nearly normally distributed number of goals scored and received. To this end, we interpreted each average score as a two-dimensional data point and applied a multivariate analysis of variance,
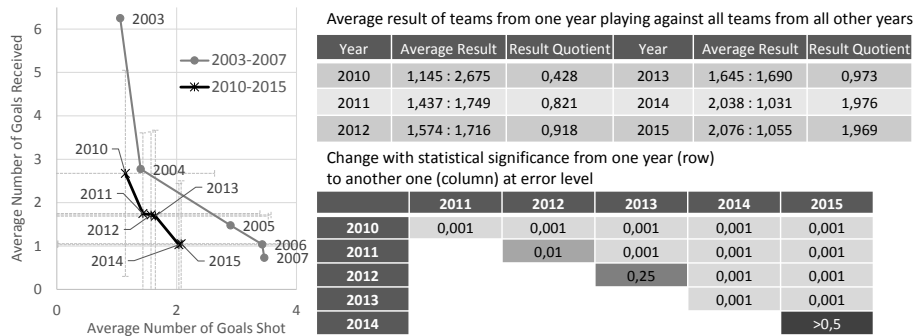
**Fig. 4.** Left: Comparison of the development of the average scores in both stable periods over years. Right top: Average scores when representatives of one year face teams from all other years. Right bottom: Error levels of a test for statistical significance in the change of the average score from the row-indicated year to the column-indicated one.

Average result of teams from one year playing against all teams from all other years

| Year | Average Result | Result Quotient | Year | Average Result | Result Quotient |
|---|---|---|---|---|---|
| 2010 | 1,145 : 2,675 | 0,428 | 2013 | 1,645 : 1,690 | 0,973 |
| 2011 | 1,437 : 1,749 | 0,821 | 2014 | 2,038 : 1,031 | 1,976 |
| 2012 | 1,574 : 1,716 | 0,918 | 2015 | 2,076 : 1,055 | 1,969 |

Change with statistical significance from one year (row) to another one (column) at error level

| | 2011 | 2012 | 2013 | 2014 | 2015 |
|---|---|---|---|---|---|
| 2010 | 0,001 | 0,001 | 0,001 | 0,001 | 0,001 |
| 2011 | | 0,01 | 0,001 | 0,001 | 0,001 |
| 2012 | | | 0,25 | 0,001 | 0,001 |
| 2013 | | | | 0,001 | 0,001 |
| 2014 | | | | | >0,5 |

given the empirically measured match results and group means (average scores). In so doing, we employed Wilks' $\Lambda$ test statistic, which for large $n$ (as in our case) is nearly $\chi^2$-distributed [4], and determined the significance levels at which the null hypothesis has to be rejected (null hypothesis: mean vector (goals shot, goals received) is in every year the same). The results are listed in the matrix in the bottom right part of Figure 4. The most important finding is that there has been a change at significance level of 0.1% from any year to any other one, if $y_2 - y_1 \geq 2$ (all non-diagonal entries). As far as changes from a year to its immediate successor year are considered, we do also find significance at a level of 0.1%, however, with three exceptions. From 2011 to 2012 the test statistic allows us to infer a change of the average score at a significance level of 1%, only. By contrast, for the transition from 2012 to 2013 and from 2014 to 2015, no statistically significant changes can be attested (error levels of 25% and >50%, respectively).

While our analyses so far has concentrated on average scores in conjunction with the law of large numbers, we now compare both stable periods from the win-draw-lose point of view. Part (a) of Figure 8 visualizes the dominance of a year's ($y_2$) representatives over the teams from the predecessor year ($y_1$) by showing the share of points ($\frac{pts(y_2)}{pts(y_1)+pts(y_2)}$) achieved when all representatives face one another multiple times (number of matches: 800 for the SSP, 240 for the FSP). The interesting point to observe here is that for the SSP we arrive at much smaller levels of dominance over preceding year. Since a value of 50% means equality (i.e. no progress made) the low numbers between 50 and 56% for 2011/13/14, respectively, do also correspond to the statistically not significant levels of progress in terms of average scores reported above.

Here, a considerable difference to the FSP is distinctive which even becomes clearer in part (b) of Figure 8 where the shares of matches won / drawn / lost are visualized, when teams play against representatives from all other years (data from 4000 matches for each year). Although the FSP lasted shorter than the SSP,
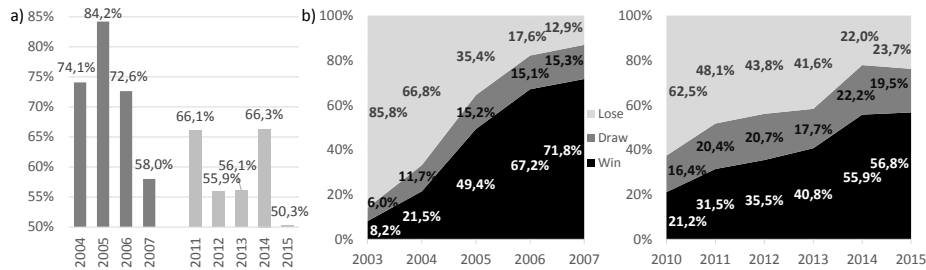
**Fig. 5.** (a) Representatives from the year ($y$) shown in the abscissa played repeatedly against teams from the immediate predecessor year ($y-1$). Share of points achieved by teams from $y$ are shown. (b) Distribution of matches won/drawn/lost by teams from one year when repeatedly facing all representatives from all other years of the same stable period.

a more distinguished slope in the performance over the years can be observed which also hints to the fact that the SSP started out from a significantly more saturated starting state which, in turn, allowed for much smaller improvements over the years.

### 3.3 Team-Focussed Progress Analysis

The goal of this section is to characterize and circumscribe certain special aspects of the current state of the 2D Soccer Simulation League as they have emerged during the SSP.

*Dominance of Top Teams* All our experimental investigations presented so far are based on the assumption that the strength of one year can be read validly from the level of play of its best representatives. We validated this claim by having a separate set of tournaments where we made ranks 1-4 from the RoboCup world cup play against ranks 5-8 repeatedly in any combination. This way we found that in such a setting the top four teams carve out 89.1% of the points whereas the ranks 5-8 yield 10.9%, only. The corresponding average score is 2.38:0.34 from the point of view of the top four. While such a kind of dominance may hold even for real soccer (perhaps less pronounced), we also have to acknowledge that the top places have been achieved more and more often by the same teams throughout the years. This fact is visualized in part (a) of Figure 6 (see the Appendix for plain team names) where the sums of the points of those teams are compared that (1) have/had the same affiliation (institution) and (2) made it to the top four for at least one year within a stable period (again: all teams played vs. all others, awarding 3/1/0 points as usual).

It is interesting to observe that the best team of the FSP (Team A) yielded 26.7% of all points, whereas the dominating team of the SSP (Team B) carves out 37.6%. The same holds, if we consider the best two / best four teams of each stable period: In the FSP they jointly obtained 49.8% (top 2) or 76.4% (top 4), in the SSP, by contrast, even 68.1% (top 2) or 91.6% (top 4). Moreover, the
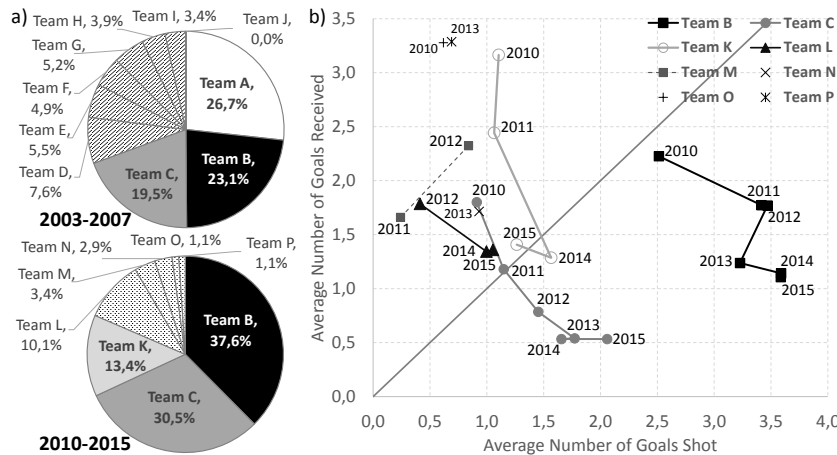
**Fig. 6.** Dominance of Established Teams in the 2D League: (a) Point distribution over successive years partitioned by teams from identical institutions. (b) Evolution of these teams throughout the second stable period in terms of average scores.

overall number of teams from different affiliations that made it to the semi-finals has decreased from 10 within the 5 years of the FSP to only 8 within the 6 years of the SSP. A conclusion to be drawn from these facts is that the performance gap between leading teams and the rest is more and more increasing.

Part (b) of Figure 6 shows the same data points as the scatter plot in Figure 2 (identical set of matches as raw data), however, this time teams from the same institution, i.e. with the same group of human developers, are now connected. This way, the dominance of the current top teams becomes even more transparent, since only teams B and C happen to position themselves below the identity function, hence scoring on average more goals than they receive. Moreover, this visualizations highlights to which extent teams improved over time.

*(Un)Worthy Champions* One of the still appealing and exciting features of 2D soccer simulation is the randomness in the simulation introduced by the Soccer Server. As in human soccer, this may result in that a team can (with a low probability) beat a stronger one and, hence, kick the stronger one out of the tournament unexpectedly. From the point of view of the actually stronger team this situation is probably very "annoying", specifically if it occurs during a final match. Figure 7 stresses the fact that the 2D Soccer Simulation League indeed had this situation once within the SSP (in 2010). In all years following, however, the provably stronger team became world champion indeed. This issue has also been addressed by Budden et al. [2] who proposed a tournament format which is aimed at minimizing fluctuations from true team performance.

*Aging Binaries* A dangerous development that might arise in a situation where a group is dominated by a small fraction of its members (as delineated at the beginning of this section) is that some kind of over-specialization is generated. Translated to the soccer simulation domain this means that most, if not all,
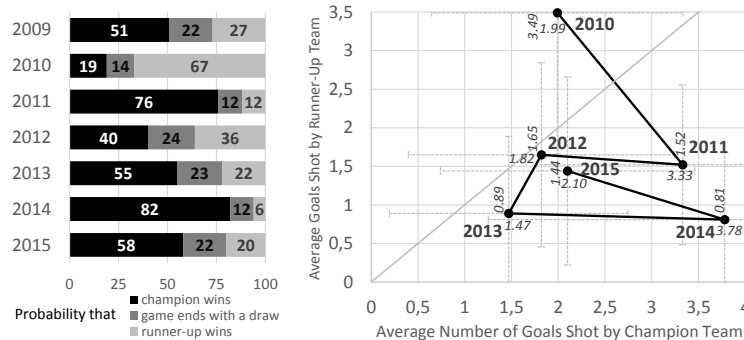
Fig. 7. Average results of 100 replays of the final matches throughout the second stable period: Share of matches won/drawn/lost (left) and average scores from the champion's perspective (right).

participants start optimizing their teams' strategies specifically against the most recent versions of last year's winner or runner-up teams. If this kind of incestuous overfitting takes place, then we should expect to see that (much) older team binaries start to come off better the older their date of publishing. In order to investigate this issue, we determined for all team binaries considered in this study (a) how well they played when they made it to the top four (and thus became a representative for its year) by playing against all other representatives from its year and (b) how much better or worse (relative to (a)) they performed when playing against the representatives of the following year. We continued this analysis for all remaining years of the SSP denoting the time elapsed as the "age" of the respective team binary.
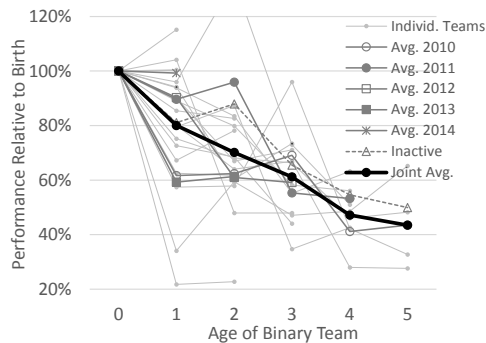


Fig. 8. No Sign of Over-specialization: Team binaries do, in general, perform worse the older they get.

Although there are minor exceptions, Figure 8 shows clearly that there is no indication that – with a team binary's increasing age – its performance starts increasing again due to reasons of over-specialization of the other teams against the most recent top teams. This finding is also confirmed by looking at the top team (Team A, binary version originating from 2010) of the FSP which has been inactive throughout the SSP (data series "Inactive"). This team's performance, though slightly above the average, decreases a little from year to year, and there is no sign of a turn-around in terms of improving performance with increased age. To sum up, on average a team binary published and used in year $y$ loses about

10-20% of its strength per year and will, for example when playing against the representatives from year $y+3$, yield only about 60% of the number of points it had carved out against the top teams from year $y$.

## 4   Conclusions

In this paper, we have presented the empirical results of a large-scale analysis of the recent progress and the current state of soccer simulation 2D. The most important conclusions from these experiments are:

– The 2D soccer simulation league has made significant progress during the recent years which we were able to quantify numerically thanks to fact that the software platform has been kept stable from 2010 to 2015.
– When compared to an earlier stable period (2003-2007) the magnitude of the progress made has been much smaller.
– The league is currently strongly dominated by a very small set of teams making it rather difficult for new teams to catch up. However, so far we do not find evidence for incestuous overfitting team strategies.

Having a stable period is extremely useful for the league and the community as a whole since it allows for studies and analyses as the one at hand. So, the policy of fixing the simulator for a couple of years should definitely be continued. With respect to the changes between the first and the second stable period this paper has shown that the league entered the SSP at a relatively mature and saturated level which did not allow for as large jumps in performance as in the FSP. Speaking about the future and a possible 2D soccer simulation agenda, we, therefore, envision three possible target directions.

– Enforce and implement finer-grained analyses of the progress of simulated soccer in terms of analyzing team play and strategies. This is certainly a shortcoming of the study at hand, since we focused solely on match outcomes without validating the actual team strategies, their maturity in multi-agent cooperation, or their level of adaptiveness.
– Enforce more focus on 2D soccer-related research benchmarks. Keepaway [8] and half-field offense [5] as well as the former coach competition are excellent examples. They might be complemented by some new learning task that becomes part of the official competitions. These points are certainly striking with regard to RoboCup's ambitious 2050 vision.
– Have a really bold, yet useful extension of the simulator. An example of such a move would be to add parts of the third dimension to the simulation (e.g. flying balls leaving the ground) and, hence, extend it from 2D to (more or less) $2\frac{1}{2}$D. This idea is not new, it had been under discussion in the technical committee already a few years ago. Among its advantages are the fact that it would *not* interfere with the 3D league which has moved to modelling humanoid robots. Furthermore, it would still allow for focusing primarily on issues of multi-agent cooperation and team play, also in conjunction with

learning approaches. Also, it might at least partially level the ground making it more attractive to new teams to enter the competition, as every team will have to adapt to the changes. Moreover, it might allow for bridging the gap to research in (human) soccer analysis where the 2D league's missing third dimension for the ball is a key obstacle [7]. Finally, it might address several of the findings brought up in this paper by allowing for a future third stable period starting out from a less saturated starting point.

## Appendix

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 2010 | HELIOS (JPN) | WrightEagle (CHN) | Oxsy (ROM) | ESKILAS (IRN) |
| 2011 | WrightEagle (CHN) | HELIOS (JPN) | Marlik (IRN) | Oxsy (ROM) |
| 2012 | HELIOS (JPN) | WrightEagle (CHN) | Marlik (IRN) | Gliders (AUS) |
| 2013 | WrightEagle (CHN) | HELIOS (JPN) | YuShan (CHN) | Axiom (IRN) |
| 2014 | WrightEagle (CHN) | Gliders (AUS) | Oxsy (ROM) | HELIOS (JPN) |
| 2015 | WrightEagle (CHN) | HELIOS (JPN) | Gliders (AUS) | Oxsy (ROM) |

**Table 1.** RoboCup World Championships' Top Teams in the 2nd Stable Period

| Id | Plain Team Name | Id | Plain Team Name | Id | Plain Team Name | Id | Plain Team Name |
|---|---|---|---|---|---|---|---|
| A | Brainstormers (GER) | E | OPU-Hana (JPN) | I | Mersad (IRN) | M | Marlik (IRN) |
| B | WrightEagle (CHN) | F | TsinghuAeolus (CHN) | J | Everest (CHN) | N | YuShan (CHN) |
| C | HELIOS (JPN) | G | AmoyNQ (CHN) | K | Oxsy (ROM) | O | ESKILAS (IRN) |
| D | STEP (RUS) | H | UvATriLearn (NED) | L | Gliders (AUS) | P | Axiom (IRN) |

**Table 2.** List of All Teams Among Top Four from 2003-07 and 2010-15

## References

1. Akiyama, H., Dorer, K., Lau, N.: On the Progress of Soccer Simulation Leagues. In: R. Bianchi, H. Akin, S. Ramamoorthy, K. Sugiura, editors, RoboCup 2014: Robot Soccer World Cup XVIII, LNCS. pp. 599–610. Springer, Joao Pessoa, Brazil (2014)
2. Budden, D., Wang, P., Obst, O., Prokopenko, M.: RoboCup Simulation Leagues: Enabling Replicable and Robus Investigation of Complex Robotic Systems. IEEE Robotics & Automation Magazine 3(22), 140–146 (2015)
3. Gabel, T., Riedmiller, M.: On Progress in RoboCup: The Simulation League Showcase. In: J. Ruiz-del-Solar, E. Chown, P. Plger, editors, RoboCup 2010: Robot Soccer World Cup XIV, LNCS. pp. 36–47. Springer, Singapore (2010)
4. Johnson, R., Wichern, D.: Applied Multivariate Statistical Analysis. Prentice Hall, New Jersey, USA (1998)
5. Kalyanakrishnan, S., Liu, Y., Stone, P.: Half Field Offense in RoboCup Soccer: A Multiagent Reinforcement Learning Case Study. In: G. Lakemeyer and E. Sklar and D. Sorenti and T. Takahashi, editors, RoboCup 2006: Robot Soccer World Cup X, LNCS. pp. 72–85. Springer, Bremen (2006)
6. Noda, I.: Soccer Server: A Simulator of RoboCup. In: Proceedings of the AI Symposium 1995. pp. 29–34. Japanese Society for Artificial Intelligence (1995)
7. Perl, J., Grunz, A., Memmert, D.: Tactics Analysis in Soccer – An Advanced Approach. International Journal of Computer Science in Sport (12), 33–44 (2013)
8. Stone, P., Sutton, R., Kuhlmann, G.: Reinforcement Learning for RoboCup-Soccer Keepaway. Adaptive Behavior 3(13), 165–188 (2005)