



A Community-Driven Data-to-Text Platform for Football Match Summaries

Pedro Fernandes¹, Luís António Santos², Sérgio Nunes¹

¹ INESC TEC • FEUP, University of Porto

² ICS, University of Minho



Universidade do Minho
Instituto de Ciências Sociais



Introduction

This work describes the design and deployment of Prosebot, a community-driven data-to-text platform tailored for generating textual summaries of football matches derived from match statistics.

Prosebot uses a template-based Natural Language Generation module to produce initial match reports, which are subsequently refined by the reading community.

The system enhances the visibility of lower-tier matches, traditionally accessible only through data tables, greatly increasing the production and publication of textual summaries.

domenica 28 aprile 2024 @ 12h30 | Stadio Giuseppe Meazza (ITA) (Milano)

Maria Caputi (ITA)

Serie A 2023/24 - Campionato - Giornata 34

Internazionale 2-0 Torino

56 60 (r) Hakan Çalhanoğlu

Mets Partita: 0-0

SCHEDA DELLA PARTITA CRONACA PERFORMANCE NOTIZIE STADIO ARBITRO VIDEO FOTOGRAFIE COMMENTI SINTESI

COMODA VITTORIA IL CASALINGO IN UNA GARA EMOZIONANTE

L'Internazionale ha battuto il Torino

L'[Internazionale](#) ha battuto il [Torino](#), ampliando una sequenza di 28 gare senza perdere, domenica, 2-0, nella giornata 34 della [Serie A](#). La squadra di [Simone Inzaghi](#) arrivava da una vittoria, e la squadra di Torino era arrivata da due pareggi prima di questa giornata. [Hakan Çalhanoğlu](#) in primo piano con 2 reti.

Il primo tempo è finito senza nessun gol.

Alla 49', [Adrien Tameze](#) fu espulso dopo vedere il cartellino rosso. Al minuto 56, Hakan Çalhanoğlu ha fatto il primo dell'Internazionale, da dentro l'area e con la sinistra. Il centrocampista ha definito il risultato finale con una rete al sessantesimo minuto della partita attraverso un calcio di rigore. Con questo gol, centrocampista fa il quindicesimo gol della stagione.

Dopo questo risultato l'Internazionale si incontra nel primo posto in [classifica generale](#), 89 punti, trovandosi il Torino al decimo posto, 46 punti. Quanto alla prossima giornata, la squadra di Simone Inzaghi [gioca in trasferta a casa del Sassuolo](#). A sua volta, la squadra di Ivan Juric [ospita il Bologna](#).

Questo articolo fu automaticamente generato da un algoritmo programmato dal calciocc.it

Pertinente ★★★★★ Qualità ★★★★★

NELLA COMPETIZIONE

Serie A

ITA

Serie A 2023/24

Italy

FASE
Campionato

JORNADA
Giornata 34

3



zerozero.pt

zerozero.pt is a prominent Portuguese sports portal, offering comprehensive access to match-related data and statistics, encompassing teams, players, goals, substitutions, fouls, and corners.

It also provides news coverage for premier competitions, curated by a team of journalists.

This coverage encompasses merely 1.4% of the matches cataloged in zerozero.pt's database.

To take advantage of the vast volume of data, zerozero.pt pioneered the development of Prosebot, an automated system designed to craft text-based match reports predicated on structured match data.

This work presents the Prosebot platform, highlighting its role in engaging zerozero.pt's vast community in the creation of match summaries using automatically generated texts as a starting point.



Related Work

Natural Language Generation techniques have gained significant traction in journalism.

⇒ Academia: [SumTime-Mousam](#) [weather] (2005), [Plachouras et al.](#) [finance] (2016), [PASS](#) [football] (2017).

⇒ Commercial solutions: [Arria Sports](#), [WordSmith](#), Quill.

⇒ Application examples: LA Times [[crime reports](#) + [earthquakes](#)] (2015), BBC [[elections](#)] (2019)

One important distinction of the Prosebot platform is the involvement of the community of readers in the process.



Post-Editing in Text Generation

The Prosebot platform is original in the combination of NLG and the readers in crafting sports texts.

Related work can be found in post-editing and machine-in-the-loop story generation.

⇒ [Stripada et al.](#) (2005) report on the post-editing of machine generated weather reports.

⇒ Machine translation ([Vaswani et al.](#), [Sennrich et al.](#)) and automatic grammar checking ([Dale et al.](#)) are areas where text-to-text NLG is commonly used to support user work.

⇒ Machine-in-the-loop examples include text generation systems to support creative work, such as STORIUM ([Akoury et al.](#)), Writing Buddy ([Samuel et al.](#)) and Creative Help ([Roemmele and Gordon](#)).

Journalists' Perceptions towards Natural Language Generation



Survey to Journalists

To inform the design of the platform, we conducted a survey to journalists in zerozero.pt's newsroom.

This survey also aimed to gauge their perceptions regarding automated content generation tools.

The survey was conducted in 2021 and garnered responses from fifteen (15) journalists.

⇒ the survey revealed that 86.7% of the journalists believed that summaries should undergo approval before publication.

⇒ 73.3% of the journalists concurred that both the collaborators and Prosebot should be acknowledged as co-authors. Such an approach ensures transparency, allowing readers to be informed of all contributors to the text.

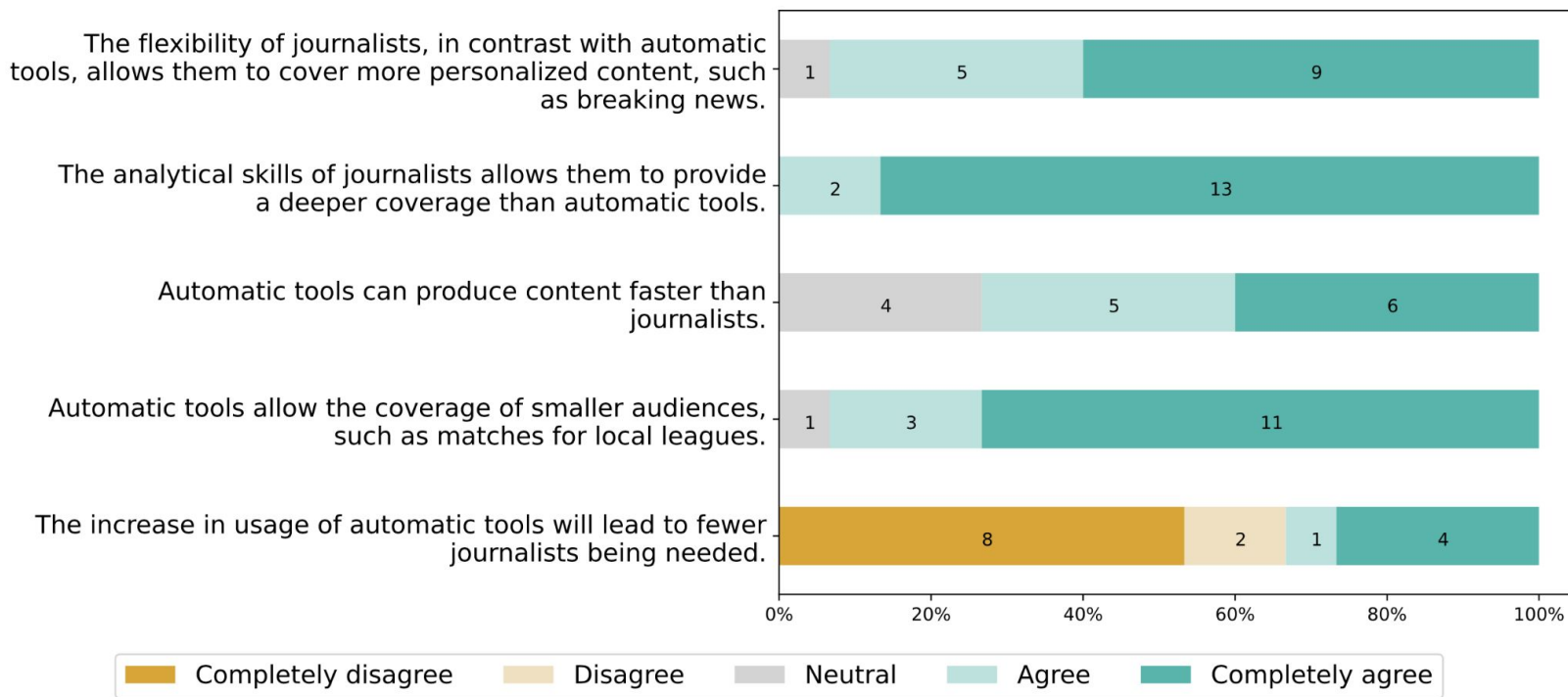


Figure 1: Overview of survey to zerozero.pt's journalists.

Prosebot Platform



Prosebot Platform

Data source ⇒ detailed match information, including participating teams, competition, final score, goal contributors, antecedent match outcomes, and player infractions from the zerozero.pt database.

Template design ⇒ developed through a collaborative effort between zerozero.pt's engineering team and its journalists. Each template encompasses both content and an associated condition.

Document structure ⇒ generated match report is organized into seven sections: (1) title, (2) subtitle, (3) small text, (4) introduction, (5) events, (6) debriefing, and (7) trivia.



Prosebot Templates

Prosebot evaluates template conditions to decide if the content should be incorporated into the generated text.

This approach has been particularly important in diversifying the template pool, thereby reducing textual monotony.

| Section | Category | Variants |
|--------------|----------------|---|
| Title | | All-purpose, 2+ goal difference, 4+ goal difference, home team won, away team won, game ended in draw, match ended with penalties |
| Subtitle | | All-purpose, favourite team lost, decisive goal in the last minutes, no goals scored, many goals scored |
| Small text | | Final score, starter and benched relevant players |
| Introduction | | Final score, previous results, best player |
| Events | Goal | All-purpose, first goal, only goal, own goal, hat-trick, poker, second goal, last goal, goal drew the match, goal increased/decreased the goal difference |
| | Substitution | Gamechanger player is subbed in, relevant player is subbed in/out |
| | Missed penalty | Goalkeeper saved, penalty taker missed |
| | Red card | Direct red card, accumulation of yellow cards |
| Debriefing | | Post-match classification, next games, match stats |
| Trivia | Stats | Best/worst result of the season for the team, best/worst overall result of the season |
| | Streaks | Increased or broke a sequence of matches |

Table 1: Available templates.

| Section | Category | Variants |
|--------------|----------------|---|
| Title | | All-purpose, 2+ goal difference, 4+ goal difference, home team won, away team won, game ended in draw, match ended with penalties |
| Subtitle | | All-purpose, favourite team lost, decisive goal in the last minutes, no goals scored, many goals scored |
| Small text | | Final score, starter and benched relevant players |
| Introduction | | Final score, previous results, best player |
| Events | Goal | All-purpose, first goal, only goal, own goal, hat-trick, poker, second goal, last goal, goal drew the match, goal increased/decreased the goal difference |
| | Substitution | Gamechanger player is subbed in, relevant player is subbed in/out |
| | Missed penalty | Goalkeeper saved, penalty taker missed |
| | Red card | Direct red card, accumulation of yellow cards |
| Debriefing | | Post-match classification, next games, match stats |
| Trivia | Stats | Best/worst result of the season for the team, best/worst overall result of the season |
| | Streaks | Increased or broke a sequence of matches |

Table 1: Available templates.

Table 1: Prosebot available templates.



NLG Algorithm (materials)

The NLG algorithm starts by parsing the template, grammar, and entity manager files.

⇒ **Template** files contain the textual content that will appear in the final result, according to the validity of their conditions given the particular match's data.

⇒ **Grammar** files cater to language-specific utilities, aiding in the representation of numbers in both ordinal and cardinal forms; and also the inclusion of articles based on the gender and plurality of the subject.

⇒ **Entity** managers primarily serve to retrieve names of entities, such as players or teams, based on stored information and previously used names



NLG Algorithm (generation)

First, the system gathers match data from internal API requests.

For each match report section, text is generated recursively by selecting and applying templates.

⇒ Template selection is based on the defined per-template conditions.

⇒ Templates are applied by replacing match data or recursively expanding nested templates.

Example template – for the first game of the match in a game where there was more than one goal:

⇒ "text": "{template.time}, {scorer.name} opened the scoring for {team.name}, {template.goal_type}, {template.assisted}."

⇒ "condition": "match_goal==1 && match_goals>1"

title

SC Braga defeated CD Tondela

Comfortable home win in goal rain

subtitle

small text

SC Braga defeated CD Tondela on Sunday, 4-2. Arsenalistas scored by João Novais, Lucas Piazón 2x and Ricardo Horta, while Tondela's team scored by João Jaquité and Souleymane Anne.

SC Braga triumphed over CD Tondela, [4-2](#), on Sunday, in the 20th round. In this competition, arsenalistas came from a win, and Tondela's team came from a win. [Lucas Piazón](#) was on fire. After 18 minutes, Lucas Piazón opened the scoring for SC Braga, with a right-foot shot inside the box, laid on by [Wenderson Galeno](#). After 40 minutes, [Ricardo Horta](#) fired home arsenalistas's second goal, laid on by Lucas Piazón. Shortly before the interval, [João Novais](#) struck for Braga's team, with a right-foot shot from outside the box, laid on by Lucas Piazón. With 50 minutes on the clock, Lucas Piazón struck for Carlos Carvalho's team, laid on by [Abel Ruiz](#). With 84 minutes on the clock, [Souleymane Anne](#) struck for CD Tondela, with a left-foot shot inside the box, laid on by [Salvador Agra](#). With 90 minutes already on the clock, [João Jaquité](#) netted the final goal of the game. After the result SC Braga are 3rd in the [table](#), 43 points, while CD Tondela occupy 12th place, 21 points. In their next fixture, arsenalistas [visit Nacional](#), while Tondela's team [will host Gil Vicente](#).

debriefing

introduction

events

trivia

Figure 2: Sample of text generated by Prosebot in English.



Saturday 20 April 2024 12h00

Jake Woodman U18 Premier League South Div First Stage 2023/24 - League - Matchweek 23



Arsenal

[U18]

Michal Rolsiak 24 (pen.), Louie Copley 45+3, Louis Zecevic-John 51, Chido Obi-Martin 90+2

4-2

Half Time: 2-2

Aston Villa

[U18]

Bradley Burrows 13, Cole Brannigan 14

MATCH REPORT LIVE COMMENTARY PERFORMANCE STADIUM REFEREE VIDEOS PHOTOS COMMENTS MATCH REPORT

ARSENAL BAG THREE IMPRESSIVE POINTS ON HOME SOIL IN A THRILLING MATCH

U18 Premier League: Arsenal defeat Aston Villa

Arsenal triumphed over Aston Villa in a fantastic win on Saturday, 4-2, in a game relating to matchweek 23. In this competition, Jack Wilshere's team came from two wins, and Gerard Nash's team came from two defeats.

The first goal of the game was scored by Bradley Burrows, with 13 minutes on the clock. In the 14th minute, Aston Villa increased their lead with a goal from Cole Brannigan. After 24 minutes, Michal Rolsiak struck for Arsenal, from the penalty spot. With 45 minutes already played, Louie Copley struck the second goal for Jack Wilshere's team.

After the break, Max Lott had a golden opportunity on his feet to score but missed a penalty. With 51 minutes on the clock, Louis Zecevic-John fired home for Arsenal. In the 82nd minute, Max Asante-Boakye left the game early after seeing a red card. Chido Obi-Martin netted the last goal in the final minutes of the game.

After the result Arsenal are 3rd in the table, with 41 points, while Aston Villa sit 11th place, with 17 points. With regard to the next fixture, Jack Wilshere's team will play away against Norwich City, while Gerard Nash's team visit Reading.



This article was automatically generated from an algorithm programmed by playmakerstats.com

Relevant? ★★★★★ Quality? ★★★★★

Ad by Refinery89

COMPETITION



Premier League

U18 Premier League South Div First Stage 2023/24

England

FASE JORNADA
League Matchweek 23



Sunday 21 April 2024 00h00

Váctor Alfonso Cáceres Hernández Liga MX Clausura 2023/24 - League - Matchweek 16



Club León

39 Federico Viñas, 47 Ángel Mena

2-0

Half Time: 1-0

Monterrey

MATCH REPORT LIVE COMMENTARY PERFORMANCE STADIUM REFEREE VIDEOS PHOTOS COMMENTS MATCH REPORT

CLUB LEÓN BAG THREE IMPRESSIVE POINTS ON HOME SOIL IN AN EMOTIONAL GAME

Club León get the better of Monterrey

Club León triumphed over Monterrey in a win on Sunday, 2-0, in a game relating to matchweek 16. In this competition, Jorge Bava's team came from two defeats, and Fernando Ortiz's team came from a draw.

The first goal was scored by Federico Viñas in the 39th minute.

Ángel Mena netted the last goal after the break.

After the result Monterrey are 3rd in the table, with 29 points, while Club León sit 11th place, with 23 points. In the next fixture, Jorge Bava's team visit Juárez. Meanwhile, Fernando Ortiz's team will play away against Necaxa.

Ad by Refinery89

COMPETITION



LIGA MX

Liga MX Clausura 2023/24

Mexico

FASE JORNADA
League Matchweek 16



This article was automatically generated from an algorithm programmed by playmakerstats.com

Relevant? ★★★★★ Quality? ★★★★★

Match report examples online.



Community-Based Platform (1)

zerozero.pt has always been a proponent of community engagement.

- 1) When a collaborator navigates to a match page, an option to craft a summary is presented.
- 2) Upon selection, they are directed to a dedicated interface.
- 3) This interface features a text editor pre-populated with content generated by Prosebot's NLG module.

For reference, the original generated text is also displayed. To guide collaborators, a brief overview of the initiative is provided, emphasizing certain guidelines. For instance, zerozero.pt advises against incorporating personal biases or subjective judgments into the summaries.



Community-Based Platform (2)

4) Upon submission of the summary, collaborators are presented with a brief survey, allowing them to share feedback on their user experience.

5) Once processed, the summary is displayed at the top of the respective match page.

In the spirit of transparency, an alert is displayed alongside the summary, crediting both Prosebot and the contributing collaborator. The platform also facilitates a collaborative review process: other community members can comment on and rate the summary based on its quality and relevance

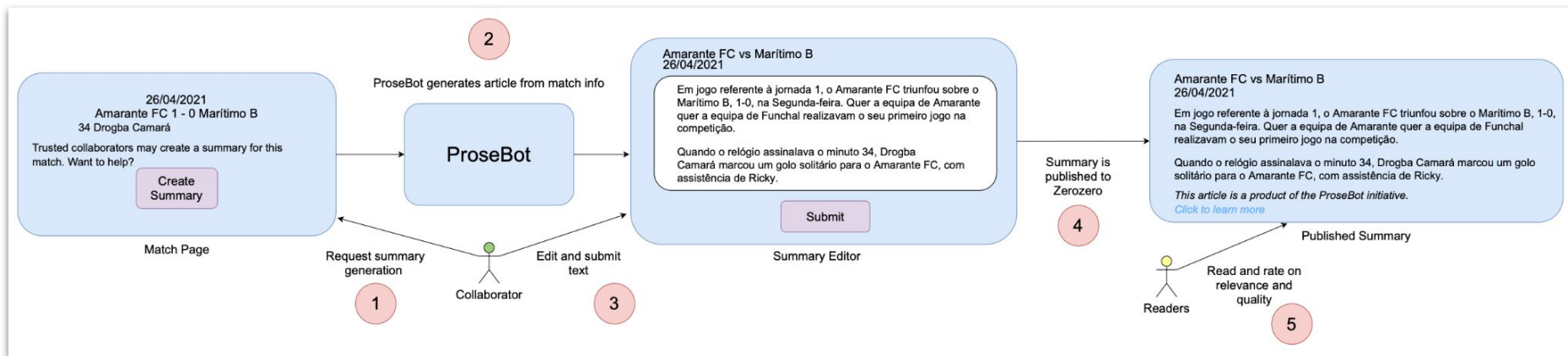


Figure 3: Overview of the Prosebot platform workflow.

Evaluation



Evaluation

To assess the efficacy and utility of the platform, we undertook a multifaceted evaluation approach.

First, we conducted an **automated comparison between the initial Prosebot drafts and the final texts** that were published on zerozero.pt.

Second, we **collected direct feedback from users through a post-submission survey**, aiming to capture their experiences and perceptions.

Third, we delved into an **analysis of the data**, examining the range of teams and competitions covered by the generated summaries and tracking the traffic these summaries attracted.



Changes to Initial Drafts

To assess the similarity between the initial drafts generated by Prosebot and the final published texts one month post-launch (174 texts), we used the Dice text similarity metric.

In a large majority of the published texts (70%), the Dice coefficient exceeds 90%, indicating a high degree of similarity between the drafts and the final texts.

Finally, we have applied a manual and automatic word frequency analysis, to understand what kind of text the collaborator usually adds or removes. Collaborators frequently added depth to the final paragraph that discusses team classifications post-match. They also highlight significant achievements, such as a team clinching the league title.



Post-Submission Survey

To gain a deeper understanding of user experiences and perceptions, we administered a survey comprising both quantitative and qualitative components.

The quantitative section consisted of seven (7) questions, each based on a 5-point Likert scale, where participants indicated their level of agreement with specific statements.

The qualitative section, on the other hand, consisted of two open-ended questions, allowing respondents to provide more detailed feedback.

We collected 46 responses, with an overall positive sentiment.

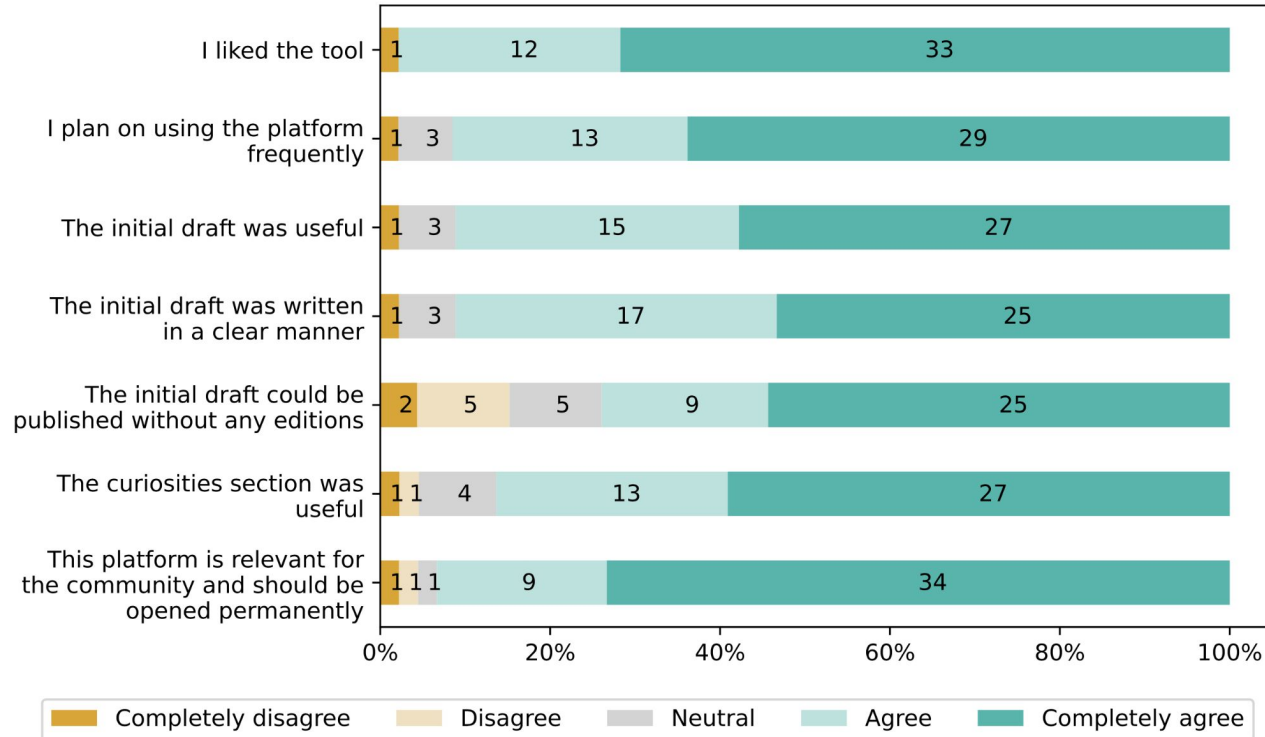


Figure 4. Results of the survey to Prosebot users.



Platform Usage and Engagement (2021)

The Prosebot platform was officially launched in early May 2021 – www.zerozero.pt/prosebot.php

- ⇒ A majority of authors produced between 1 to 4 reports.
- ⇒ The top five most active users were responsible for 40% of all reports generated.
- ⇒ Over half of the reports were published within two days of the match's conclusion.
- ⇒ A segment of users (18%) generated reports for matches from more than three months ago



Platform Usage and Engagement (2024)

In March 2024 we can observe a strong consolidation of the Prosebot platform within the community.

Over 8 thousand match summaries have been generated by more than 580 users, covering over 3 thousand distinct teams.

The top 5 users were responsible for 30% of all summaries generated and 14 users were responsible for 40%.

| Indicator | May 2021 | March 2024 |
|--------------------------------|----------|------------|
| Number of summaries | 174 | 8,079 |
| Number of authors | 43 | 587 |
| Number of teams covered | 197 | 3,116 |
| Number of competitions covered | 56 | 295 |
| Mean summaries per author | 4.0 | 13.8 |
| Mean summaries per team | 1.8 | 5.1 |
| Mean summaries per competition | 3.1 | 27.4 |
| Mean visits per summary | 53.0 | 363.8 |

Table 2: Prosebot platform key statistics.

Discussion and Conclusions



Limitations

The Prosebot platform is used to produce match summaries derived from football match statistics.

The system is currently being used in a real-world setting with a user base of hundreds of thousands of users. Since its launch it has been used to generate thousands of texts.

One of the main limitations of a template-based system such as Prosebot is the diversity of the generated text. Although a large collection of templates have already been developed, content diversity is impacted and can be improved. To tackle this problem, automatic template extraction is being explored to augment the number of available templates.



Ethical Statement

Automatically generated media, also known as synthetic media, is increasingly common.

However, its societal impact is still under intense scrutiny and debate. In the design of the Prosebot platform, actions were taken to control the impact of such systems. Most importantly, all game summaries generated have a clear disclaimer stating that the text was generated by Prosebot.

Additionally, journalists were involved in the design of the system since early prototypes.

In the design of the Prosebot system we have found that initial fears of “machines replacing humans” have been replaced by a positive perspective of “machines helping writers”, as reflected on the survey the journalists reported in the paper.



Conclusions

In this work, we present the Prosebot platform, a pioneering effort in the Portuguese domain that integrates a community of readers into the post-editing process of automatically generated sports match summaries.

A survey conducted a month post-launch revealed a favorable reception among participants.

In the span of three years post-launch, Prosebot has been instrumental in generating millions of game summaries across various languages. This has notably amplified the coverage and visibility of matches at the grassroots level, which previously might have been overlooked.



Acknowledgements

We extend our sincere appreciation to ZOS for their ongoing and productive collaboration over the years, with special thanks to Marco Sousa and Pedro Dias.

This work is financed by national funds through the Portuguese funding agency, FCT - Fundação para a Ciência e a Tecnologia, within project LA/P/0063/2020 (DOI 10.54499/LA/P/0063/2020).



A Community-Driven Data-to-Text Platform for Football Match Summaries

Pedro Fernandes¹, Luís António Santos², Sérgio Nunes¹

¹ INESC TEC • FEUP, University of Porto

² ICS, University of Minho



Universidade do Minho
Instituto de Ciências Sociais