**RESEARCH ARTICLE**

# Toward computer-supported semi-automated timelines of future events

Alan de Oliveira Lyra[1*], Carlos Eduardo Barbosa[1,2], Yuri Oliveira de Lima[1], Herbert Salazar dos Santos[1], Matheus Argôlo[1] and Jano Moreira de Souza[1]

## Abstract

During a Futures Study, researchers analyze a significant quantity of information dispersed across multiple document databases to gather conjectures about future events, making it challenging for researchers to retrieve all predicted events described in publications quickly. Generating a timeline of future events is time-consuming and prone to errors, requiring a group of experts to execute appropriately. This work introduces NERMAP, a system capable of semi-automating the process of discovering future events, organizing them in a timeline through Named Entity Recognition supported by machine learning, and gathering up to 83% of future events found in documents when compared to humans. The system identified future events that we failed to detect during the tests. Using the system allows researchers to perform the analysis in significantly less time, thus reducing costs. Therefore, the proposed approach enables a small group of researchers to efficiently process and analyze a large volume of documents, enhancing their capability to identify and comprehend information in a timeline while minimizing costs.

**Keywords**  Named entity recognition, Futures Research, Machine learning, Software, System

## Introduction

The world is experiencing rapid transformations, amplifying the unpredictability and complexity of the involved processes. This situation increases the demand for studies related to planning, forecasting, and creating future visions. Futures Studies are deemed crucial for fostering the ability to establish innovation systems that align with societal needs.

In Futures Studies, finding conjectures or predictions of future events and organizing them is a process that takes a lot of time and is susceptible to failures when conducted manually, particularly when handling large-scale research that encompasses numerous documents to be examined. Furthermore, when executed by hand and in collaboration, the examination/selection can be inconsistent among individuals.

This study introduces a collaborative system called NERMAP to process various documents, and analyze and extract sentences associated with a timestamp—mainly dealing with the future and automatically assembling timelines. NERMAP uses the named entity recognition (NER) technique to identify and extract sentences that are probably related to future events. After this process, humans can then make changes to the generated Timeline.

The most significant benefit of using NERMAP is its ability to quickly process large amounts of documents, saving considerable time compared to manual processing. The increased processing speed allows for analysis that would otherwise be infeasible to perform manually.

This analysis provided by NERMAP and further refined by humans can be a complete analysis for a Roadmapping process during Futures Studies, as shown in [1–4]. We can also use it as input for other stages of future studies

*Correspondence:
Alan de Oliveira Lyra
alanlyra@cos.ufrj.br
[1] COPPE, Universidade Federal do Rio de Janeiro, Rio de Janeiro - RJ 21941-450, Brazil
[2] Centro de Análises de Sistemas Navais, Marinha do Brasil, Rio de Janeiro - RJ 20091-000, Brazil

Lyra *et al. European Journal of Futures Research*        (2023) 11:4

Page 2 of 9

workflow, such as Scanning Horizons, Futures Wheels, and Scenarios, or for identifying conflicting expectations, contributing to the Foresight area.

## Literature review

This section addresses the theoretical foundation of the fundamental concepts for NERMAP: the NER technique. NER is a component widely used in Information Extraction [5] built upon recognition tasks to extract and categorize the entities mentioned in a natural language text according to the model's predefined structure and rules [6, 7]. The entity usually refers to a place, person, institution, or date. The outcome of a NER system is an organized representation of the unstructured input text [8].

Implementations that utilize NER rely on creating patterns for recognizing entities primarily by manually intensive training models, requiring human effort and significant labor [9]. NER systems learn such labeled data patterns automatically [10]. The Conditional Random Field (CRF) [11–13], Hidden Markov Models (HMM) [14], Maximum Entropy (ME) [15], and Maximum-Entropy Markov Model (MEMM) [16–18] are statistical approaches to machine learning to execute NER.

This study utilized the CRF method for NER because the methodology applied was supervised learning, which relies on creating a training corpus. Additionally, CRF has a faster processing time than other evaluated methods [19] and yields better results in NER tasks [20].

The CRF is designed to label and segment sequential data according to the $p(Y|X)$ conditional distribution method with a corresponding graphical model. Variable $X$ is a vector of random input variables, and $Y$ is a vector of random output variables. In this way, $p(Y|X)$ is the probability of a given as input vector $X$ to get output $Y$.

The implementation of CRF for applying NER can be formulated as a problem of finding a conditional distribution. Let $X = [x_1, x_2,..., x_{n-1}, x_n]$ be a sequence of words in a text, where $x_j$ is a word located at position $j$ of the text. Let $Y = [y_1, y_2,..., y_{n-1}, y_n]$ be a sequence containing a label for each word of the vector $X$. Then, we have that $p(Y|X)$ says the probability that the text represented by $X$ is labeled $Y$ [21].

## NERMAP

We developed the NERMAP framework to aid the generation of a timeline of future events, semi-automating it, and enabling the researcher to study a large volume of analyzed documents, improving efficiency, and reducing expenses of Futures Studies.

We built the NERMAP software, a collaborative Foresight Support System [22] based on the NERMAP framework. The software supports on-site and remote (through the Internet) research of Futures Studies and Foresight fields to generate timelines of future events, as shown in [1–3]. The NERMAP model, architecture, and system are discussed in the following sections.

### NERMAP model

The detection of future events in texts can be accomplished using various methodologies. The methodology employed in this work categorizes the terms that compose a future-related text into three categories:

- Temporal space: terms that denote a date, which can be a year or a phrase of time (e.g., "next year", "decade", and "century")
- Future-indicating terms: terms that indicate a future event, always connected to a temporal context (e.g., "In 2030", "By 2030", and "Around 2030")
- Event: passages that announce future events, for example, something forecasted to occur in a specific future timeframe.

Table 1 illustrates each concept, while Table 2 displays a potential timeline for each case.
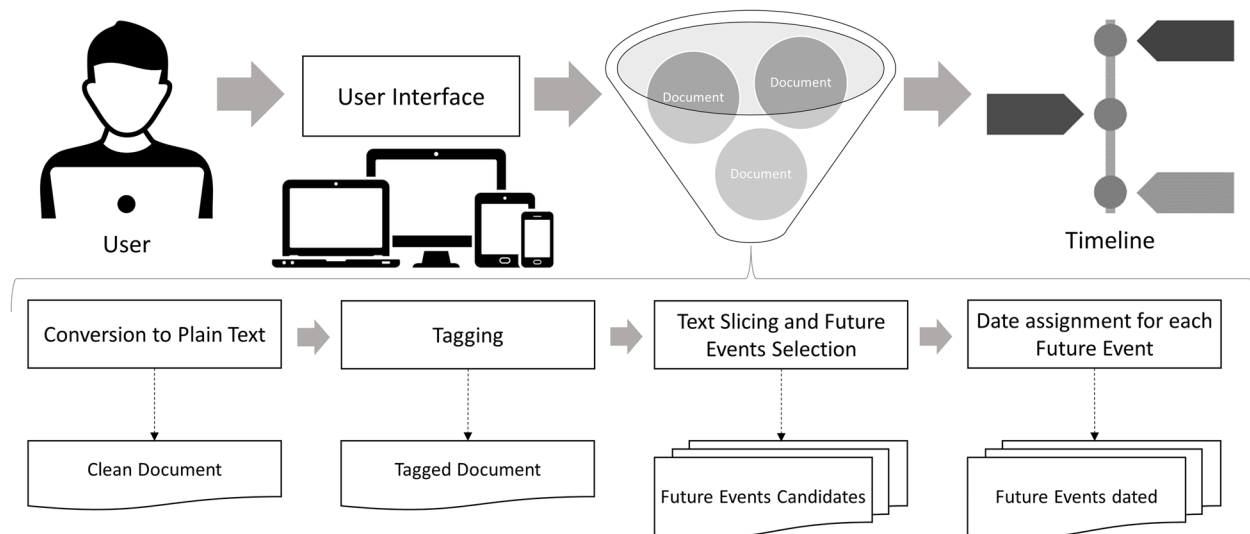
The NERMAP enables multiple people to collaborate through the use of projects. In line with this principle,

**Table 1** An example of identifying future events in a fictional text

| Type | Input |
| --- | --- |
| Temporal space | In *2040*, Brazil will be at the forefront of environmental sustainability. According to experts, by *2025*, Brazil will have reached the same level as other nations in implementing green technologies. By *2030*, Brazil will surpass Europe's efforts to reduce carbon emissions and preserve natural resources. For the year *2040*, Brazil will be at the forefront of the clean energy and conservation industry. |
| Future-indicating terms | *In* 2040, Brazil will be at the forefront of environmental sustainability. According to experts, *by* 2025, Brazil will have reached the same level as other nations in implementing green technologies. *By* 2030, Brazil will surpass Europe's efforts to reduce carbon emissions and preserve natural resources. *For the year* 2040, Brazil will be at the forefront of the clean energy and conservation industry. |
| Event | In 2040, *Brazil will be at the forefront of environmental sustainability*. According to experts, by 2025, *Brazil will have reached the same level as other nations in implementing green technologies*. By 2030, *Brazil will surpass Europe's efforts to reduce carbon emissions and preserve natural resources*. For the year 2040, *Brazil will be at the forefront of the clean energy and conservation industry*. |

**Table 2** Example of a timeline of future events

| Date | Future events |
| --- | --- |
| 2025 | Brazil will have reached the same level as other nations in implementing green technologies |
| 2030 | Brazil will surpass Europe's efforts to reduce carbon emissions and preserve natural resources |
| 2040 | Brazil will be at the forefront of environmental sustainability |
| 2040 | Brazil will be at the forefront of the clean energy and conservation industry |



**Fig. 1** NERMAP model

the process begins with the user creating a new project, as depicted in Fig. 1, which can be shared with other users via invitations.

The user loads documents—mostly articles or technical reports—in TXT or PDF formats. Users must also indicate a year that will serve as a basis to be the foresight limit of the events to be found (e.g., 2080), as well as the main area—if any—of the study (e.g., medicine, work, education, or even general). Subsequently, the documents are converted to plain text and tagged using the Named Entity Recognition method.

NERMAP slices the tagged text, selects candidates for future events, and dates each event according to the context. Dated future events can be classified into three categories: Complete, where the selected text comprises the predicted event in its entirety and encompasses its date; Partial, which displays an unfinished future event that requires manual refinement; Invalid, a section that is not a future event.

Users can view each project from the single-document perspective or the composite of all documents perspective to build the Timeline—as illustrated in Fig. 2.

NERMAP starts with document processing, in which parser libraries process the documents, converting PDFs

to plain text, followed by applying the NER technique for text tagging. The Stanford NER supports the NER used in the supervised machine-learning task of this work. We use Stanford NER due to the ease of adaptation of pre-existing models and the use of a native CRF algorithm implementation.

A training model[1]—corpus—is needed for machine learning, constructed through the manual annotation of entities in articles and technical reports, enabling general learning of the supervised learning tool.

### NERMAP architecture

The NERMAP architecture (Fig. 3) was designed with four main layers. The presentation layer (in blue) includes the User Interface, which provides a visualization of the instantiations of each module. The management layer (in brownish-red), with three modules: the User Management, which is responsible for the registration and editing of each system user; the Process Management, which

---

[1] The NERMAP corpus included the articles from the *Journal Futures* and was tagged using the BILOU [23] (B = Begin, I = Inside, L = Last, O = Outside, and U = Unit) notation.
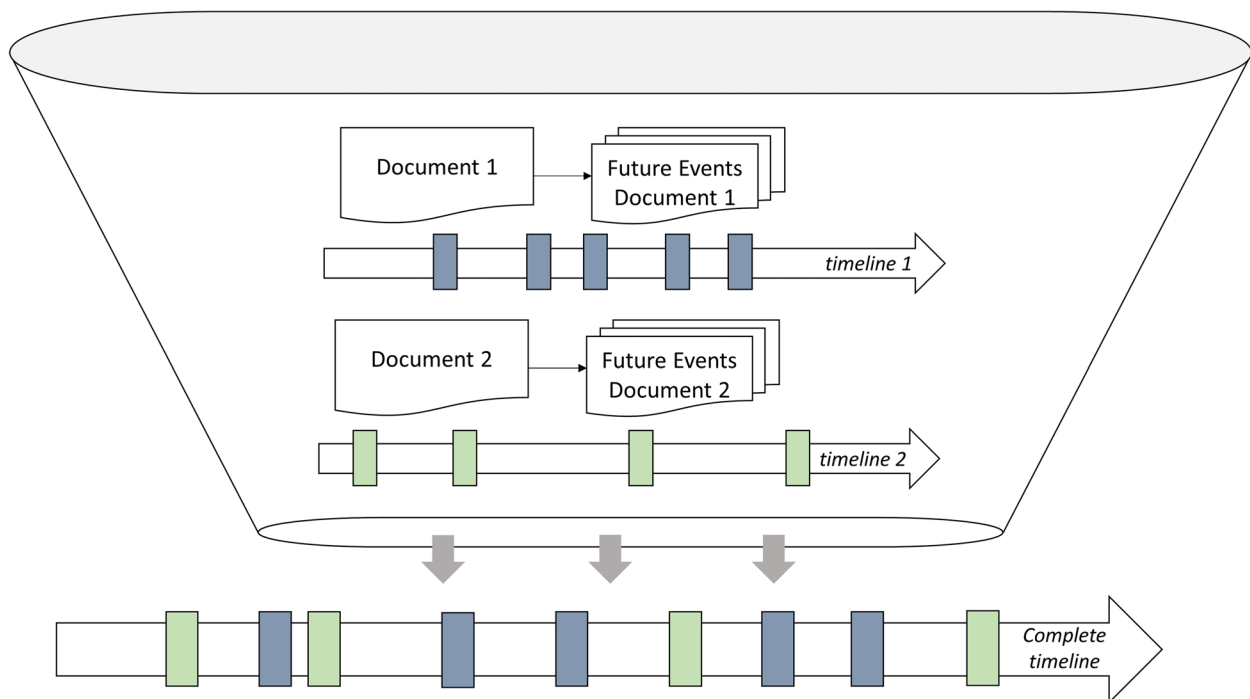
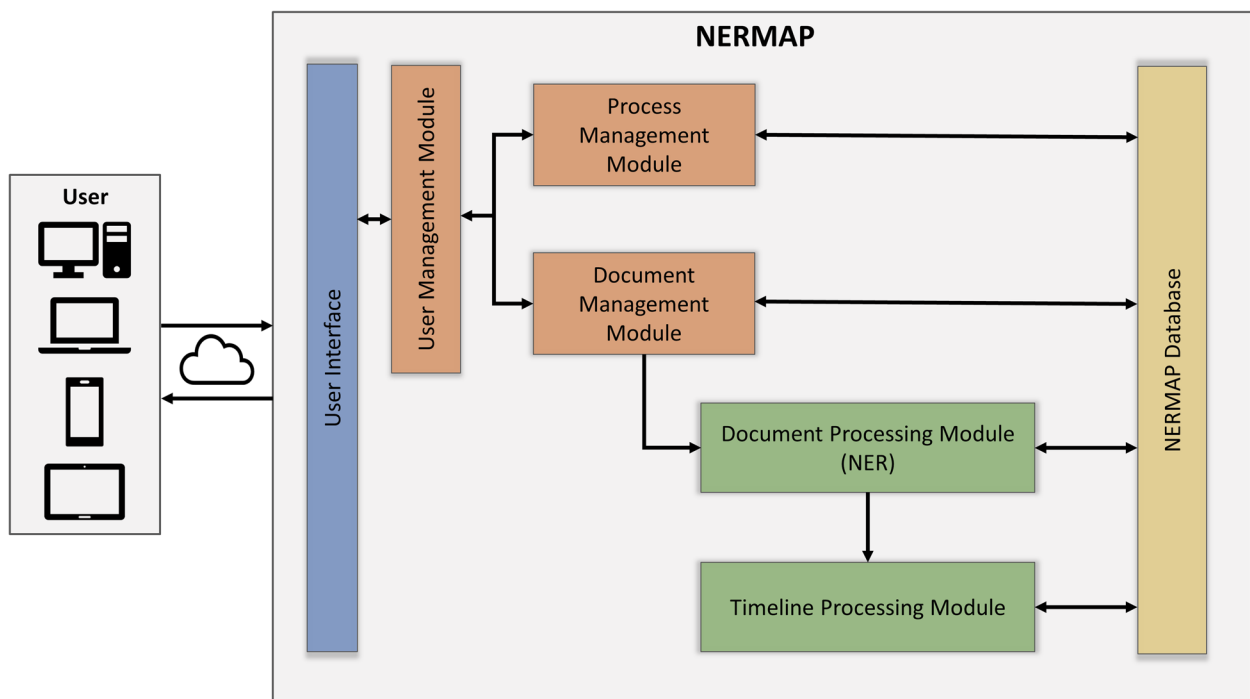**Fig. 2** NERMAP process for creating timelines



**Fig. 3** NERMAP architecture

is responsible for the creation of the process, assigning parameters and roles for each participating user, being able to guide groups; and the Document Management, that is responsible for the submission of documents to start the semi-automated process. The Processing layer (in green), with two modules: Document Processing,

responsible for converting documents to plain text and then tagging the texts based on NER, and Timeline processing, responsible for obtaining the dated future events to build the Timeline. Finally, the Storage layer (in yellow) with the NERMAP Database is responsible for storing all gathered and produced by the system.

### NERMAP system

The NERMAP system enables users to manage their projects, both created independently or shared by others on the Management screen. NERMAP includes the following functionalities: adding a metadata file with *title, authors, year of publication*, and *source reliability*; sharing projects with other users through invitations; viewing and editing added files; editing, deleting, and viewing timelines.

The Timeline—presented in Fig. 4—is the most crucial aspect of the system, as it displays the results of the analysis of submitted documents and organizes future events in a chronological format.

The following features are included:

- The ability to edit or delete future events from the Timeline
- The option to manually add a future event: when a researcher finds a future event in a publication that needs to be included in NERMAP's Timeline, uploading the entire document for processing may lead to unwanted duplicated or out-of-scope events
- The ability to filter for smaller timelines using the documents that form the complete Timeline
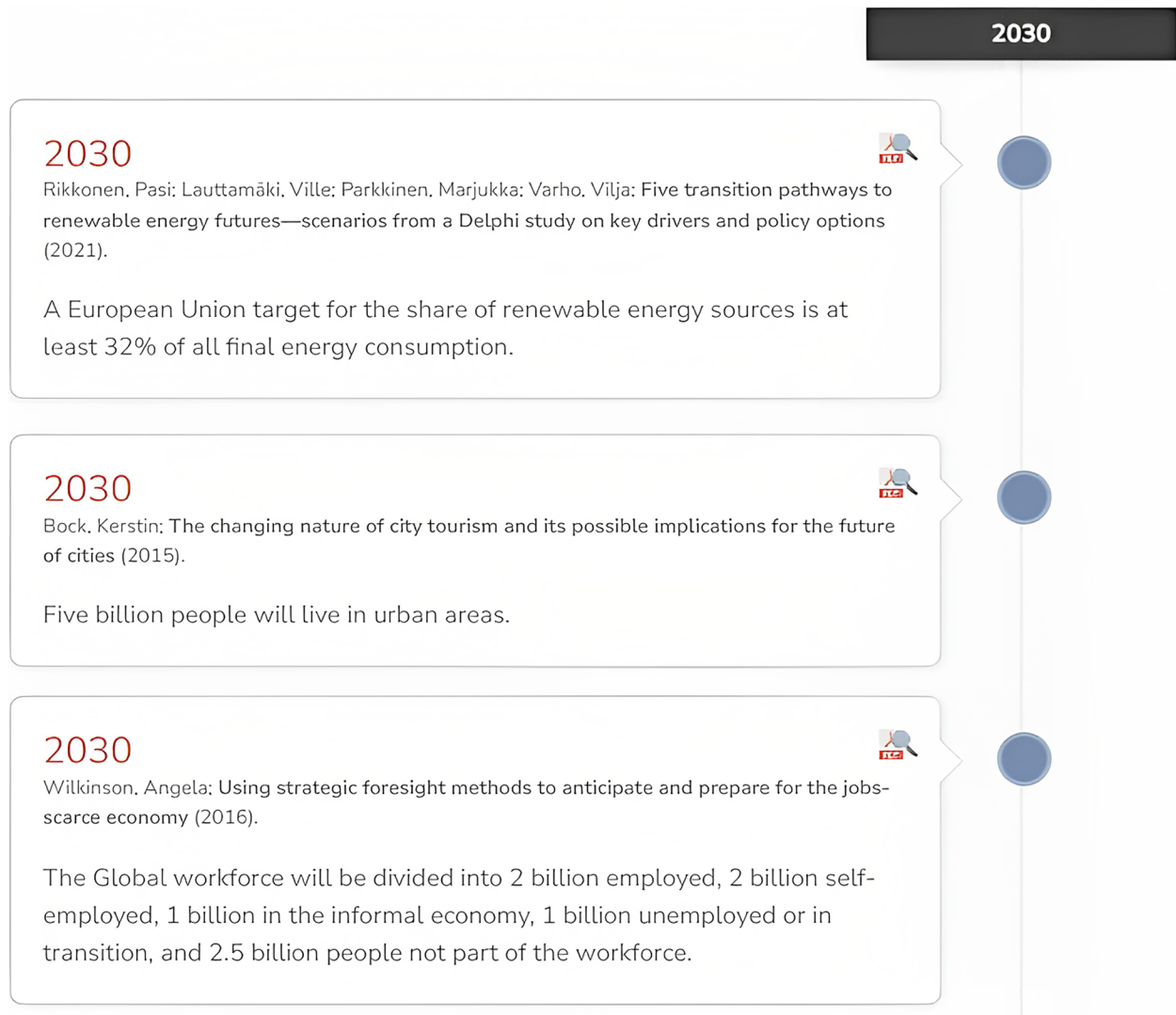


**2030**

Rikkonen, Pasi; Lauttamäki, Ville; Parkkinen, Marjukka; Varho, Vilja: Five transition pathways to renewable energy futures—scenarios from a Delphi study on key drivers and policy options (2021).

A European Union target for the share of renewable energy sources is at least 32% of all final energy consumption.

**2030**

Bock, Kerstin: The changing nature of city tourism and its possible implications for the future of cities (2015).

Five billion people will live in urban areas.

**2030**

Wilkinson, Angela: Using strategic foresight methods to anticipate and prepare for the jobs-scarce economy (2016).

The Global workforce will be divided into 2 billion employed, 2 billion self-employed, 1 billion in the informal economy, 1 billion unemployed or in transition, and 2.5 billion people not part of the workforce.

**Fig. 4** NERMAP timeline screen

- The capability to apply word or phrase filters to the Timeline to focus on a specific niche of the study: this is useful for focusing on a particular subset within a larger area of study using a previously collected set of publications
- The option to view the highlighted future event within the document: necessary function to evaluate and either include or exclude a future event by evaluating the context of the future event within the publication
- The ability to move through the years in the Timeline
- The ability to view the average timeline reliability calculated based on the credibility of each publication source in the system
- The option to export the Timeline through an API or in various formats such as CSV, JSON, TXT, PDF, and DOC.

NERMAP allows for much more collaboration in this process of Futures Study than when executed manually. Data (documents and generated timelines) in NERMAP is kept secure between projects since users cannot view or alter third-party data without the permission of the project owner.

### Evaluation

NERMAP was evaluated through a three-cycle experiment. We examined the system-generated future events and collected performance metrics from the model's tagging stage—Precision, Recall, and F-Measure.

$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}} \quad (1)$$

$$\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}} \quad (2)$$

$$\text{F-Measure} = \frac{2 * \text{Precision} * \text{Recall}}{(\text{Precision} + \text{Recall})} \quad (3)$$

The experiment cycles were rooted in a manual analysis of Futures documents published in 2020 to spot future events. We uploaded the documents to the NERMAP system using the process outlined in Fig. 2. Then, we compared the results gathered by NERMAP with those of manual analysis, as illustrated in Fig. 5, and categorizing the future events found by NERMAP as Invalid, Partial, or Complete and also determining if they had been identified in the manual analysis, as illustrated in Fig. 6.

The study utilized 1091 documents using *Journal Futures* publications from 2010 to 2020 as input. We used the period from 2010 to 2018 for training the model and the remaining period for testing.

In the first cycle, we utilized a model trained using publications from 2010 to 2012 as a benchmark for the subsequent cycles. In this cycle, we found issues such as difficulty in converting files to plain text due to the library used and new cases requiring attention by the system for future events.
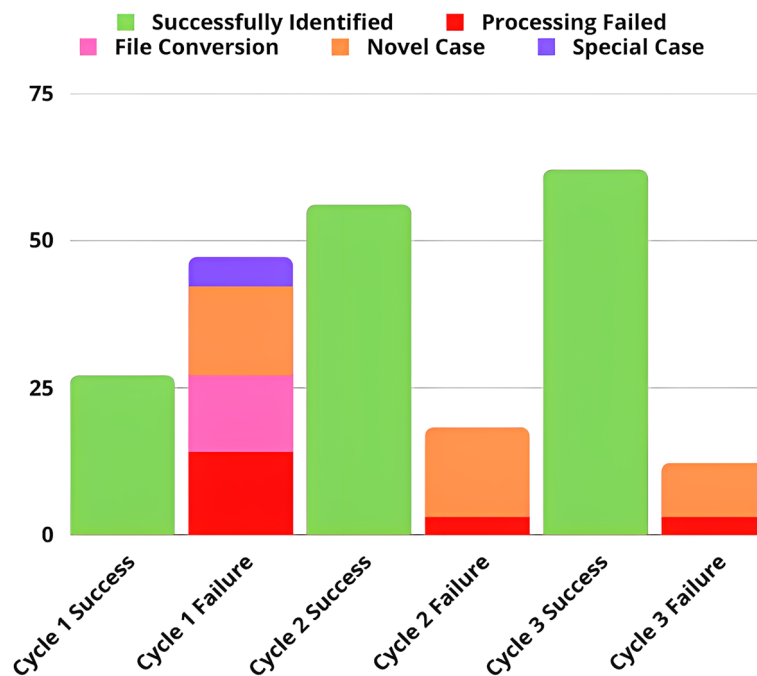


**Fig. 5** Future events discovered by the NERMAP system versus through manual means

In the second cycle, we address the problems identified in the first cycle. We manually converted the documents that had failed to convert to plain text and made the necessary adjustments for the observed cases. NERMAP's performance improved from 36 to 75% compared to the number of future events found through manual means (as shown in Fig. 5) and a total of 117 events (an improvement of 67%, as presented in Fig. 6). While the results were greatly improved, with NERMAP discovering more future events with better accuracy, the model did not identify 20% of events found through manual means as they encompass new cases not included in the training data.

In Cycle 3, our primary goal was to decrease the number of cases NERMAP could not return by expanding the corpus and upgrading the training model utilized. Publications from 2012 to 2018 were examined and added to the corpus. We note the improvement in the model in the results summarized in Table 3. We applied the process to publications from 2019 to 2020 and delivered superior results compared to the model employed in the first two cycles.

As a result of the model, NERMAP discovered 83% of the future events found through manual means, yielding a larger quantity of accurate future events.

In the third cycle, NERMAP discovered 125 future events. We found that most future events (83%) could benefit a real Futures Study. Table 4 summarizes some of the significant future events that comprise the Timeline for the third cycle.

**Table 3** Metrics

| Model | Precision | Recall | F-Measure |
|---|---|---|---|
| Cycles 1 and 2 | 0.6923 | 0.6290 | 0.6592 |
| Cycle 3 | 0.7079 (+2.25%) | 0.6774 (+7.7%) | 0.6923 (+5%) |

## Discussion

The experiment cycles showed a gradual growth of the system, being possible to identify problems and new cases of future events to improve the model. Driven by corrections in processing documents and improvement in the NER training model, NERMAP increased the number of returned events classified as complete and partial—events we must maximize—while decreasing the number of events classified as invalid—events we must minimize.

NERMAP demonstrated high efficiency, identifying 83% of events found through manual means in a large set of documents in seconds. While the system's 83% hit rate is respectable, the NERMAP model does not account for all potential use cases of future events in Futures Studies publications. Further research on other journals is required to expand the corpus, likely resulting in even higher accuracy.

One of the main challenges identified was the file processing issue, where the library chosen for plain text conversion experienced problems. The following steps
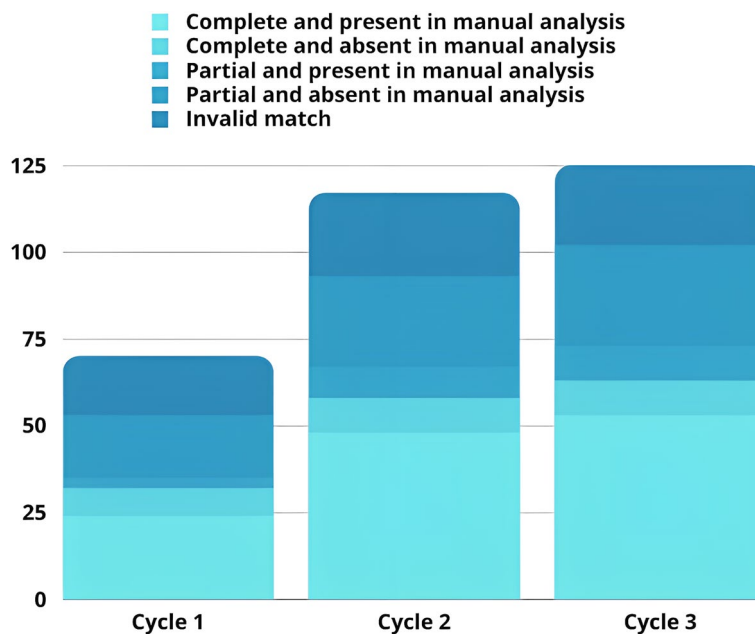


**Fig. 6** Analysis of future events discovered by NERMAP

**Table 4** The third cycle timeline

| Date | Event |
|------|-------|
| 2025 | The vaccines market increases its value to over $ 70 billion |
| 2030 | The adverse effects of climate change can increase maize the crop in southern Africa by up to 30% compared to 2008 figures |
| 2030 | Sustainable Development Goals are interlinked to ensure unification throughout all the represented sectors |
| 2030 | Annual growth rates will likely exceed 10 to 15% |
| 2033 | About 47% of the US employment is at a high risk of automation |
| 2035 | Brazil emerges as a significant producer and supplier of oil, ranking as the 6th largest energy provider globally |
| 2035 | NASA has plans to send astronauts to Mars |
| 2035 | China plans to establish a settlement on the moon, followed by a colony on Mars |
| 2040 | Global energy consumption increases by almost 80% when compared to 2016 |
| 2040 | There will be no more working days in Mexico, and everything will be assigned by projects and resolved within the designated timeframe for delivery |
| 2050 | Travel demand for motorized private transport will decline by 28% compared to 2015 as the rural demographic is expected to decline |
| 2050 | The Brazilian population stabilizes at 240 million |
| 2063 | Africa will have a more decisive role in a globalizing world, aiming to achieve the SDGs primarily through normal economic growth |
| 2068 | Technology is embedded in our daily lives, improving lifestyles globally |
| 2100 | Expansion of global per-capita growth by a factor of three to eight, regardless of biophysical limitations, compared to 2020 |
| 2300 | If $CO_2$ emissions from current fossil fuel sources continue, the highest concentration of $CO_2$ in the atmosphere could reach 1400 ppm, resulting in a global temperature increase of 8ºC or more |

for NERMAP depend on the transition to a new tool that can perform this type of file parsing.

Some studies use the NERMAP system. One of these works involved an analysis of possible trends for Industry 4.0 [1]. When used in collaboration, NERMAP can be a powerful tool, providing a platform for multiple researchers to work on the same topic and share their findings.

As seen in the previous section with the evaluation, the NERMAP tool cannot provide a 100% exact result with all possible data compared to a human reading—and interpreting—the text. Over time, NERMAP will be able to update itself with improvements that allow its model to identify a more comprehensive range of future events in texts based on broader training.

In a larger-scale study, NERMAP can save time, decrease research project or workgroup expenses, and facilitate greater collaboration among individuals by providing access to many processed documents. Automating a labor-intensive and error-prone task in Futures Studies is a significant advancement for the field. Furthermore, we can use NERMAP as a part of a Futures Study workflow. In this case, researchers can refine the output of the NERMAP's Timeline and develop many roadmaps that can be applied as input in the next steps of a Foresight.

## Conclusions

The process of forecasting and constructing a timeline of future events typically involves a thorough examination of a vast amount of data to identify future events, making it challenging for researchers to gather all the relevant information from these publications efficiently. As a result, traditional methods for creating such a timeline can be time-consuming and often require a team of experts.

In this work, we present NERMAP, a system that aims to assist researchers in semi-automating the creation of timelines for future events. The tool demonstrated an accuracy rate of 83% and quickly generated the final product using a large set of documents.

In addition to the time and cost savings in a Futures Study, NERMAP provides a much more consistent collaboration throughout the process than manual execution, centralizing the entire process, and organizing information in a single integrated environment.

In future work, we plan to integrate new journals, expand the training corpus of the NER model, and enhance the document-processing library of NERMAP. Moreover, NERMAP can incorporate new algorithms for identifying future events.

Lyra *et al. European Journal of Futures Research* (2023) 11:4

Page 9 of 9

## Declarations

### Ethics approval and consent to participate
Not applicable.

### Consent for publication
Not applicable.

### Competing interests
The authors declare that they have no competing interests.

## References
1. El-Warrak L, Nunes M, Lyra A et al (2022) Analyzing industry 4.0 trends through the technology roadmapping method. Procedia Comput Sci 201:511–518. https://doi.org/10.1016/j.procs.2022.03.066
2. Simoes RV, Parreiras MVC, Silva da ACC et al (2022) Artificial intelligence and digital transformation: analyzing future trends. p 6
3. Barbosa CE, Lima Y, Lyra A, Oliveira D (2019) Healthcare 2030: a view of how changes on technology will impact Healthcare in 2030. Laboratório do Futuro
4. Barbosa CE, de Lima YO, Costa LFC et al (2022) Future of work in 2050: thinking beyond the COVID-19 pandemic. Eur J Futures Res 10:25. https://doi.org/10.1186/s40309-022-00210-w
5. Bunescu RC (2007) Learning for information extraction: from named entity recognition and disambiguation to relation extraction. Thesis, The University of Texas at Austin, Austin
6. Doddington G, Mitchell A, rzybocki M et al (2004) The Automatic Content Extraction (ACE) Program Tasks, Data, and Evaluation. 4
7. Elloumi S, Jaoua A, Ferjani F et al (2013) General learning approach for event extraction: case of management change event. J Inf Sci 39:211–224. https://doi.org/10.1177/0165551512464140
8. Sureka A, Goyal V, Correa D, Mondal A (2009) Polarity classification of subjective words using common-sense knowledge-base. In: Sakai H, Chakraborty MK, Hassanien AE et al (eds) Rough Sets, Fuzzy Sets, Data Mining and Granular Computing. Springer, Berlin Heidelberg, Berlin, Heidelberg, pp 486–493
9. Appelt DE, Hobbs JR, Bear J et al (1993) FASTUS: a finite-state processor for information extraction from real-world text. IJCAI. pp 1172–1178
10. Ciravegna F (2001) Adaptive information extraction from text by rule induction and generalisation. 17th International Joint Conference on Artificial Intelligence 2:6
11. Amaral D, Vieira R (2014) NERP-CRF: uma ferramenta para o reconhecimento de entidades nomeadas por meio de Conditional Random Fields. Linguamática 6:41–49
12. Lafferty J, McCallum A, Pereira FCN (2001) Conditional random fields: probabilistic models for segmenting and labeling sequence data. 18th International Conference on Machine Learning (ICML) 10
13. Chatzis SP, Demiris Y (2012) The echo state conditional random field model for sequential data modeling. Expert Syst Appl 39:10303–10309. https://doi.org/10.1016/j.eswa.2012.02.193
14. Bikel DM, Miller S, Schwartz R, Weischedel R (1997) Nymble: a high-performance learning name-finder. In: Proceedings of the Fifth Conference on Applied Natural Language Processing (ANLC '97). Association for Computational Linguistics, Washington, DC
15. Chieu HL, Ng HT (2003) Named Entity recognition with a maximum entropy approach. Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003. pp 160–163
16. Bender O, Och FJ, Ney H (2003) Maximum entropy models for named entity recognition. Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003. pp 148–151
17. Curran J, Clark S (2003) Language independent NER using a maximum entropy tagge. Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003. pp 164–167
18. Finkel J, Dingare S, Manning CD et al (2005) Exploring the boundaries: gene and protein identification in biomedical text. BMC Bioinformatics 6:S5. https://doi.org/10.1186/1471-2105-6-S1-S5
19. Mota C, Santos D (2008) esafios na avaliação conjunta do reconhecimento de entidades mencionadas: O Segundo HAREM. Linguateca
20. Finkel JR, Grenager T, Manning C (2005) Incorporating non-local information into information extraction systems by Gibbs sampling. Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics - ACL '05. Association for Computational Linguistics, Ann Arbor, Michigan, pp 363–370
21 da Silva TS (2012) Reconhecimento de Entidades Nomeadas em Notícias de Governo. Dissertação de Mestrado, UFRJ
22. Barbosa CE, Lima Y, Emerick M, et al (2022) Supporting distributed and integrated execution of future-oriented technology analysis. Futures Foresight Sci. https://doi.org/10.1002/ffo2.136
23. Ratinov L, Roth D (2009) Design challenges and misconceptions in named entity recognition. Proceedings of the Thirteenth Conference on Computational Natural Language Learning. Association for Computational Linguistics, USA, pp 147–155

## Publisher's Note
Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.