

Dear authors,

first of all, thank you for the opportunity to review your article "Result Assessment Tool (RAT): Empowering search engine data analysis". Overall, the article presents a clear and concise exploration of the RAT and its functionalities. However, several areas required clarification and enhancement to ensure the comprehensive and effectiveness of the study.

Please find below some comments.

Introduction

Overall the introduction provides a clear overview of the motivation behind developing the RAT software and its intended functionalities. The background information provided about the lack of access to search engine data for researchers helps contextualise the need for RAT. However, the section could be enhanced by providing more context at the beginning, references to support the need of developing such software and the limitations researchers face in accessing search engine data. The introduction briefly mentions existing tools that cover some of RAT's functionalities but highlights the limitations of these tools, such as being outdated or lacking flexibility. Providing a more detailed comparison with existing tools, including their strengths and weaknesses, could strengthen the argument for the necessity of RAT and further highlight its contributions to the field. Furthermore, the authors should emphasize RAT's flexibility and discuss more effectively how the tool address the diverse needs of researchers.

To ensure clarity and consistency, it is important for the authors to introduce and define the concept of search systems earlier in the introduction, especially since they initially refer to platforms like Google, Twitter, and Facebook (line 32) and later discuss them to them as search systems (line 39). Clarification of difference between search engines (like Google) and search systems is discussed in lines 41 - 43. Given that Facebook and Twitter primarily function as social networks but also incorporate search functionality, it is essential to delineate these concepts clearly, at the beginning. Furthermore, looking also at the case studies presented in Table 4, it seems that the only systems analyses are search engines. Could you please clarify?

The Result Assessment Tool

This section provides an informative overview of the main components and capabilities of RAT.

However, the concise manner in which the tool is discussed may present challenges for readers seeking a more detailed understanding of its functionalities.

For example, while the RAT evaluation interface is mentioned, further elaboration on its features and usability would improve clarity for readers unfamiliar with the tool's interface and functionality.

In line 104, the flexibility of RAT is highlighted. To clarify this aspect, an example could be provided to illustrate how RAT can adapt to different research needs. For instance, researchers conducting an IR study may require different features compared to those conducting a qualitative content analysis.

User journey in RAT

Please note that there is an error to reference source not found.

- Create study: This section primarily focuses on search engines. Is it also applicable to other search systems like Facebook, Twitter, at mentioned earlier, in the Introduction? Additionally, if APIs are available for platforms like Facebook and Twitter, why would web scraping be necessary? Could the authors address these questions?
- Collect search results: the authors have not mentioned any potential legality of collecting search engine results via web scraping. Does RAT comply with terms of service of the systems/web sites being scraped? Furthermore, since the authors mention collecting the full document of search results (lines 138 to 140), do they consider any legal implication or restrictions associated with capturing and storing result documents, particularly in terms of copyright or intellectual property rights?
- Evaluate and analyse results: Could the authors provide an example of evaluation applied in practice?

Model Structure of RAT

error source reference not found

RAT Frontend error references

line 181 figure 5 (not sequential - issue with numbering)

Researcher view

- line 188: Reference source not found.

- line 188 references Figure 3D, then 3B (line 194), and finally 3A (line 202). Could the authors please review the numbering?
- line 198: “bootstrap forms providE” (removed the ending S in provide)
- line 211: how the overlap between the results from different search engines is measured?

Evaluation View

- lines 225 to 227: how this approach can guarantee anonymity?

RAT Backend

- lines 256-257: how can the classification processes be automated within RAT?
- line 259: The triggers for modules may be defined based on specific conditions or events within the system. For example, a module may be triggered to execute when a certain task is completed, when new data is available, or when a predefined time interval elapses. These triggers are typically set during the system design phase and can be adjusted as needed based on the requirements of the application. Is it correct?
- line 265: how do the authors handle jobs failure? Regarding the too many requests from the same server, is this caused by a limit in submitting the requests? What type of options have the authors considered to avoid this?
- The dependency on HTML and CSS structures highlights a potential limitation of web scraping. How did the authors address such challenges and ensure reliable data extraction? For example, by implementing mechanisms for rotating IP addresses or introducing delays between requests to mitigate the risk of being blocked.
- regarding line 265, when the authors mention “it is mandatory to update the search engine”, they likely mean that it is necessary to update the web scraping code used to collect data from search engines?
- changes in SERP structure can significantly impact the performance of tools like RAT. Can the authors discuss strategies for detecting and adapting to changes in SERP structure, such as regular monitoring and updates to the scraping algorithms?
- line 267: RAT is designed to automatically reset and restart the jobs after a configurable time up to a maximum number of attempts. However, in the case of the scraper’s IP address being blocked, the authors should clarify how RAT handles this situation and whether additional measures, such as rotating IP addresses or using proxy servers, are implemented to address IP blocking.

Search Engine Scrapers

- line 276: in the reference to search engines only, could the authors clarify why social platforms such as Facebook and Twitter were cited at the beginning, especially since they have APIs available?
- please see error reference not found.
- could the authors provide more statistical information regarding characteristics of the search engines considered for this study, and limitations? Perhaps, characteristics of search engines could be included in a table.
- line 290: could the authors clarify which type of adjustments are needed to collect search results beyond the first result page?
- line 294: error source not found. Please address it.
- line 299: can the authors clarify the approach of search engines in delivering search results? Providing an example or table would enhance understanding.
- line 307: it is mentioned that search engine scrapers need frequent updates, which represents a limitation. How does this affect the performance of the tool?

Source scraper

- line 321: could the authors clarify why they need the entire HTML source code with its content for classification purposes?

RAT Extensions

- line 364 error
- handling of research data generated by RAT
- line 280 error

Results

- line 433 error
- retrieval effectiveness studies:
- (...)

Classification studies

- line 468-469. While manual classification is provided in Mazzeo et al. for assigning a label (fake not fake), data collection was not performed manually, as referenced in their article. Please rectify.
-

Discussion and conclusion

- line 550: throughout the paper, only search engines were mentioned, not search systems.

References

line 663: <https://gs.statcounter.com/search-engine-> not found

line 672: published on May 2020 (not 2022)

line 635-636: DOI NOT FOUND

line 621-622: DOI NOT FOUND