

Betty's (Re)Search Engine

A client-based search engine for research software stored in repositories.

Vasily Seibert  ¹

Andreas Rausch  ¹

Stefan Wittek  ¹

1. Institute for Software and Systems Engineering, TU Clausthal, Clausthal-Zellerfeld.



Date Received:

2023-02-15

Licenses:

This article is licensed under: 

Keywords:

Ingrid, Data, research data management, research software

Data availability:

Data can be found here:

https://github.com/VasilySeibert/cascading_search_ecology_eval

Software availability:

Software can be found here: <https://github.com/VasilySeibert/BettysReSearchEngine>

Abstract. Promoting research, without providing the source code that was used to conduct the research, means a greater effort for every researcher down the line. Existing solutions that aim to make research software FAIR [1], fail to provide a wholesome solution, for they do not sufficiently consider already existing research software stored on platforms like GitHub or organizational GitLabs. We therefore present Betty's research engine, a client-based implementation of a cascading search process, that first finds research software stored on platforms like GitHub and then links them to corresponding publications or entries in third party databases. We evaluated 400 random search results from the domain of ecology and found that 345 out of 400 repositories made a reference to a corresponding publication / entry in third party database and therefore clearly indicating the potential of the cascading search. Betty's research engine is live and openly available under this URL: <http://nfdi4ing.rz-housing.tu-clausthal.de/>

1 Introduction

Findability and accessibility of research related software is crucial for every researcher who aims to (i) fully understand and (ii) reproduce software related publications as well as (iii) benchmark the own software related research to other solutions. However, finding fitting research software is not a trivial task. It is not obligatory for most authors of software related publications to provide a link to corresponding software repositories where the research software would then be stored. [2] annotated a stratified sample of a collection of software mentions, extracted from the CORD 19 Dataset [3] of open access publications. Their goal was to understand the status of software citation practices. Out of 295.609 software mentions in the Softcite-CORD-19 Dataset 50 % are solely names of the software without further details, 35 % provide a version, 21 % mention their publisher, and 9 % have a URL given in the Text. The authors of [4] reviewed the state of code availability in ecology using a random sample of 346 nonmolecular articles published between 2015 and 2019 under mandatory or encouraged code sharing policies. Their results show, that only 27 % of eligible articles were accompanied by code, although the percentage of ecological journals with mandatory or encouraged code sharing policies increased from 15 %

17 (2015) to 75% (2020).

18 There are existing approaches, that aim to make research software FAIR. They can be broadly
19 categorized as:

- 20 • Community driven approaches (e.g. Papers with Code [5])
- 21 • Technology driven approaches (e.g. Zenodo [6], Citation File [7])
- 22 • Standards and policies (e.g. FAIR4RS [8])

23 Papers with Code [5] is an example of a community driven approach to make research soft-
24 ware more findable and accessible. It offers references to papers from the field of mathematics,
25 physics, machine learning, astronomy and computer science as well as their corresponding data-
26 and software repositories. Generation of new content as well as quality checking is done by a
27 user community.

28 Zenodo [6] and Citation File [7] are examples of technology driven approaches. Zenodo is a
29 platform that allows users to upload research related digital artifacts (such as research software)
30 and generate a DOI (Digital Object Identifier) for them (therefore making them citable). If
31 the digital artifact is also stored on GitHub, Zenodo offers a Zenodo-Badge, a small image file
32 that can be integrated into the README to make the Zenodo DOI more visible. The citation
33 file format [7] is a file format that is human and machine readable. If a citation file is placed
34 in a Github Repository it unlocks the "cite this repository" function on GitHub, which further
35 increases visibility and the probability of getting cited.

36 Metadata standards (e.g. FAIR4RS [8]) and consequently policies (because someone has to
37 enforce these standards) are designed to uniform the way research is documented and archived,
38 with the objective of making it FAIR. Findability according to [8] can be achieved if:

- 39 1. the software is assigned a globally unique and persistent identifier,
- 40 2. the software is described with rich metadata,
- 41 3. the metadata explicitly includes the identifier of the software it describes
- 42 4. the Metadata is searchable and indexable

43 All of the described approaches rely on researchers to explicitly add additional information to
44 their publications and therefore none of the approaches can be regarded as a wholesome solution
45 for the problems described in [4] and [2] because of the following reasons:

- 46 • Finding existing research software is more important than the chance of finding (a frac-
47 tion) of future research software.
- 48 • The effort for researchers to adapt their existing software repositories to a new database,
49 technology, metadata standard is considerable and therefore a further obstacle.

50 For the above reasons there is a need for searching **already existing research software reposi-**
51 **ries** (stored on GitHub [9] or GitLab [10]) automatically and then linking them to corresponding
52 publications or databases.

53 Such a *cascading search* would (i) initiate the search process by retrieving a list of references
54 and metadata that refer to software repositories. If a given software repository contains any
55 identifiers that can be connected to a corresponding publication or the listing of that software
56 repository in a third-party database, then (ii) the search engine would access and retrieve that
57 information. The result (iii) would be a list of software repositories, enriched with additional
58 metadata and therefore sortable by the user's preferences.

59 From a technical perspective this cascading search engine faces a variety of challenges in terms
60 of privacy, networking, individualization, reliability, scalability and the complexity that arise
61 from the use of multiple APIs. A user must be able to:

- 62 • understand the search process and the logic behind it,
- 63 • perform the described cascading search fully automatically,
- 64 • perform the described cascading search in an acceptable amount of time,
- 65 • be assured that his/her search is private and not being tracked,
- 66 • access sources of information that are commercial or limited to members of an organisa-
67 tion,
- 68 • extend the cascading search with further databases.

69 In this paper we reiterate the problem of not findable, accessible research software and show
70 that existing approaches are unable to provide a wholesome solution. We motivate the need for
71 searching existing research software and linking them to publications / database entries after-
72 wards (this is what we call cascading search). We formulate top level requirements and elaborate
73 how we meet them through our architectural design and our choice of technology. We evalu-
74 ate our results by analyzing 400 repositories and reflect these results in a discussion. We then
75 conclude our work.

76 2 Cascading Search

77 [4] and [2] focused on finding research software by searching for a reference in the publications.
78 We introduce the cascading search, a search process according to which, we (i) initiate the search
79 by retrieving a list of references and metadata that refer to software repositories (from a platform
80 like GitHub). If a given software repository contains any identifiers that can be connected to a
81 corresponding publication or the listing of that software repository in a third-party database, then
82 (ii) the search engine will access and retrieve that information. The result (iii) of such a cascading
83 search is a list of software repositories, enriched with additional metadata and therefore sortable
84 by the user's preferences.

85 We find (i) and retrieve references to software repositories from GitHub by using the GitHub
86 REST API. For a given search query, GitHub will find repositories that can be associated with
87 that search query (e. g. the occurrence of that term in the name of the repository). GitHub also
88 allows the user of the REST API to modify the search query in a way that makes adjustments
89 to GitHub's internal search process. For example, the search query "seismology doi in:readme"
90 would return a list of references to repositories, that can be associated with the term "seismology"

91 and that contain the word "doi" in their README files. After retrieving a list of references,
 92 we iterate through that list to collect more information on every single repository. Among the
 93 collected information is the README file as well as the name of every other file in the main
 94 directory. We then search (ii) for any identifiers to a corresponding publication or the entry of
 95 that repository in a third party database. Searching for unique identifiers is a quality determining
 96 factor for the cascading search, since this information is most reliable for performing further
 97 search on that repository. As the results from the evaluation will show, references turn out to
 98 be highly heterogenous. We developed rules for identifying a *target reference* (a reference that
 99 refers to a corresponding publication or entry in a third party database, the counterpart of a target
 100 reference would be a *normal reference* that refers to any literature that was used to conduct or
 101 support the research).

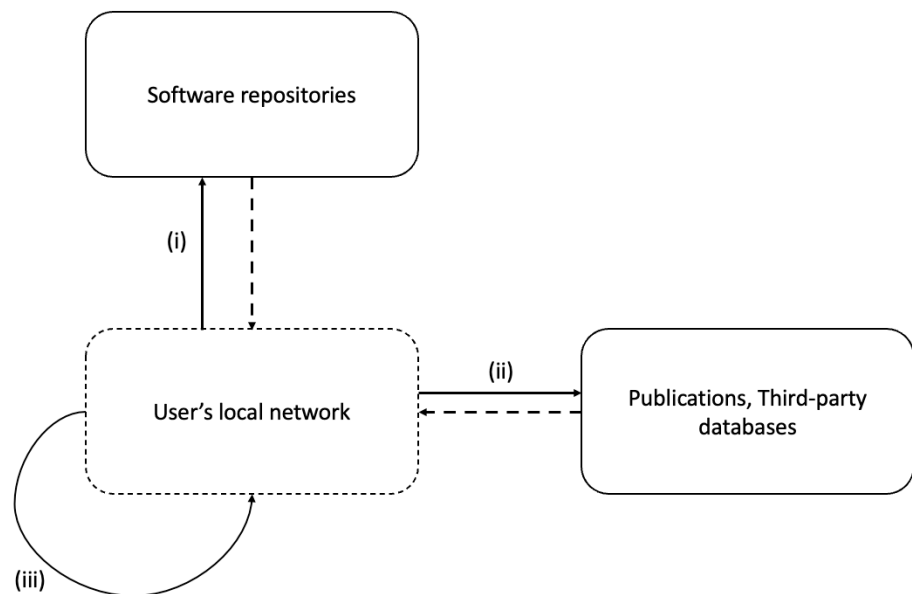


Figure 1: The cascading search (i) initiates the search process by retrieving a list of references and metadata that refer to software repositories. If a given software repository contains any identifiers that can be connected to a corresponding publication or the listing of that software repository in a third-party database, then (ii) the search engine will access and retrieve that information. The result (iii) is a list of software repositories, enriched with additional metadata and therefore sortable by the user's preferences.

102 The following rules were created with the the computational cost and number of required API
 103 calls in mind.

- 104 • If the repository contains a citation file [7], then the reference in that file is considered a
 105 target reference
- 106 • If the repository contains a Zenodo image file [6] then the DOI is considered a target
 107 reference
- 108 • If the repository contains a Zenodo DOI in plain text and the data on Zenodo is similar
 109 to the name of the repository, we consider that reference as the target reference (only

110 checking the first Zenodo DOI)

111 After retrieving the target reference of a repository, we search that reference on further platforms.

112 Current platforms include: Zenodo [6], Open Alex [11], DataCite [12] and Open citations [13].

113 The information that we retrieve from the cascading search (iii) is added to our data model.

114 With the cascading search, we enriched the repositories with additional metadata that would

115 not be available through a search on GitHub alone. The user can now, sort the repositories

116 according to personal preferences and interests (like the number of citations a repository holds).

117 The repository also holds information about what reference was considered the target reference.

118 With that, we meet the requirement of making the cascading search understandable.

119 3 Client-based search engine

120 The cascading search faces a variety of challenges. A consequence of using multiple APIs is

121 that multiple rate limitations have to be considered (one or more limitation rule per API). Since

122 rate limitations are in many cases tracked through the IP Address that makes the request, a

123 centralized backend is not possible. A centralized backend would also mean that the user could

124 not call services and databases with restricted access (commercial or organizational).

125 To face the described challenges, Betty's research engine is written solely in languages that can

126 run in the browser. We send the user everything needed to perform the cascading search from

127 the local machine. For that reason, the user is required to provide own credentials, since it is the

128 user that utilizes the different APIs. No communication is being recorded, the credentials are

129 stored on the local machine, so the user enjoys the maximum amount of privacy that is possible.

130 The architectural design follows a MVC (Model View Controller) pattern (See Fig 2.). This

131 means that the user interface (View) is strictly divided from the logic (Controller) and the

132 database (Model) of the system. The user initiates the cascading search from the view. In-

133 side the controller a *manager* holds knowledge of all available *agents*. An agent is a concrete

134 implementation of how an API is utilized. This includes making the API call, receiving the

135 returned information, processing it and passing that information through a defined interface to

136 the model. By thinking in terms of agents, the complexity that arises from using multiple APIs

137 is capsuled in a defined container. This way, the manager can orchestrate the agents by a gen-

138 eralized logic (and doesn't need to implement individual exceptions and error handlings). The

139 retrieved information is stored in the database according to a defined data model. If new infor-

140 mation is added to the model, the manager gets notified by an update mechanism. So every time

141 a data instance is added or modified, the manager gets notified and can instantiate new agents

142 that use that information.

143 Therefore the cascading search can be seen as a circular flow of data. The manager instantiates

144 a GitHub agent, who performs a call to the GitHub REST API, retrieves that information and

145 processes it before calling a model interface and adding the repositories to the database. The

146 manager gets notified about new repositories being added, and instantiates Zenodo, DataCite,

147 OpenCitations etc. agents that use the newly added information to perform searches on their

148 own, before calling another model interface and adding the information.

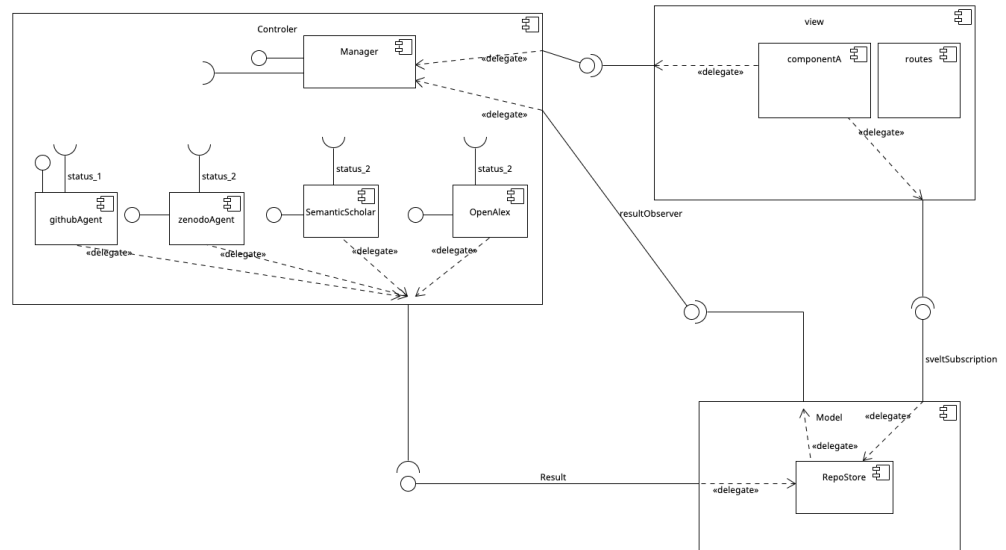


Figure 2: Betty's research engine is an implementation of the cascading search that faces the challenges of individual rate limitations and restricted access by a client based design. By handling the complexity that arises from the use of several APIs in agents, the cascading search can be performed in parallel processes while maintaining a sound, coherent structure.

149 Through the presented architectural design, the cascading search can now be conducted fully
 150 automatically and since there is no direct dependency between the agents, they can run in par-
 151 allel.

152 4 Evaluation

153 [4] analyzed 346 publications, published in ecological journals. They found that between 2015
 154 and 2019 under mandatory or encouraged code sharing policies only 27% of eligible articles
 155 were accompanied by code. To evaluate Betty's re search engine, we performed a cascading
 156 search for ecology related research software. Out of 2513 software repositories that were avail-
 157 able on GitHub, we randomly choose 400 for detailed analysis (all data and code that was used
 158 to perform this analysis can be found at [https://github.com/VasiliySeibert/cascading](https://github.com/VasiliySeibert/cascading_search_ecology_eval)
 159 [g_search_ecology_eval](https://github.com/VasiliySeibert/cascading_search_ecology_eval)). We found that 345 out of 400 (86,25%) repositories referred to a
 160 corresponding publication or the listing of the repository in a third party database.

161 With the introduced rules, we were able to correctly identify 89 (22,25%) target references
 162 (true positives). 256 (64%) repositories had a target reference, but were not identified (false
 163 negatives). 52 (13%) repositories did not refer to corresponding research and were correctly
 164 disregarded (true negatives). 3 (0,75%) repositories were falsely used for the cascading search
 165 (false positives).

166 Of further interest is how a repository refers to corresponding publications / entries in a third
 167 party database. We found that the majority of references (252, 63%) were mentions in plain
 168 text. 74 (18,5%) repositories used a Zenodo image file, 5 (1,25%) repositories had a citation
 169 file placed in their main directory. 10 (2,5%) repositories had a BibTEX text element in their
 170 README.

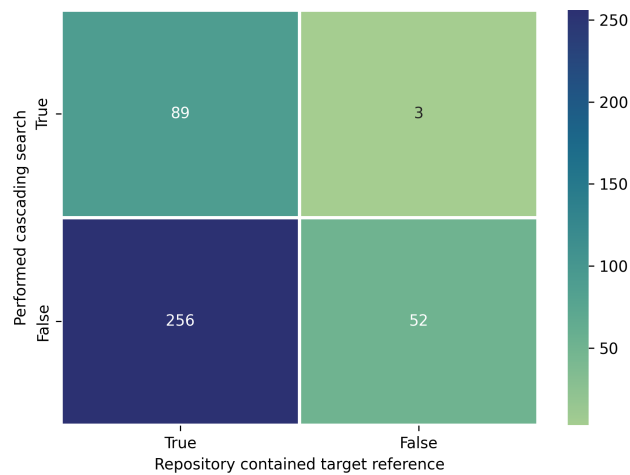


Figure 3: The introduced rules yield an accuracy of 35,25%. (all data and code that was used to perform this analysis can be found at

https://github.com/VasilySeibert/cascading_search_ecology_eval)

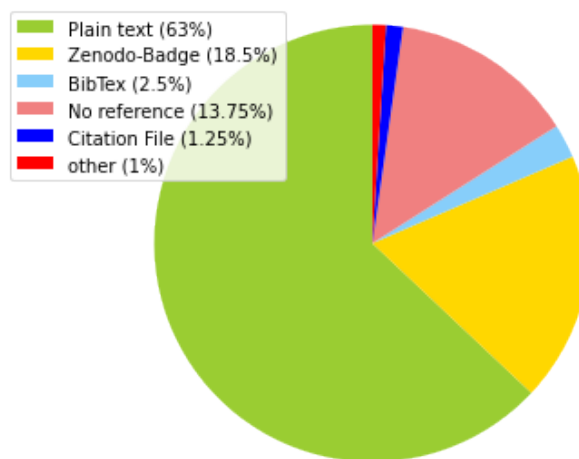


Figure 4: We found that the majority of references (252, 63%) were mentions in plain text. 74 (18,5%) repositories used a Zenodo image file, 5 (1,25%) repositories had a citation file placed in their main directory. 10 (2,5%) repositories had a BibTEX text element in their README.

171 5 Discussion

172 **Does the cascading search come up for the shortcomings of current approaches?** As the
 173 results from the evaluation chapter clearly indicate, the cascading search has the potential to
 174 compensate for the shortcomings of the approaches which were introduced in the introduction.
 175 Existing research software can be found on GitHub and furthermore connected to their respec-
 176 tive publications. The cascading search utilize technological driven approaches, in fact these
 177 approaches are most reliable when it comes to identifying target references, for their use of
 178 structured text elements. Supported by these results, the cascading search is a valid and useful

179 alternative to community- and technological driven approaches as well as metadata standards
180 and policies.

181 **How effective is Betty's research engine?** With an accuracy of merely 35,25%, Betty's re
182 search engine has plenty of room for improvement. However, in relation to the 400 analyzed
183 repositories, only 89 (22,25%) repositories made use of structured text elements Zenodo-Badges,
184 Citation Files, BibTEX. A rule-based approach might not be well suited for this type of textual
185 data. Since 63% of all references are written in plain, unstructured text, a data driven approach
186 might be more promising.

187 **In what use cases, can Betty's re search engine serve best?** Betty's research engine can
188 help researcher who aim to (i) fully understand and (ii) reproduce software related publications.
189 Unlike Papers with Code, Betty's research engine does not support any categorisation of the
190 found research. The search results do not have any information about ongoing benchmarks or
191 the state of the art for that matter. Enabling users to (iii) benchmark the own software related
192 research to other solutions would require further work.

193 6 Conclusion

194 In this paper we reiterated the problem of not findable, accessible research software and pointed
195 out how current approaches do not provide a satisfying solution, for they are not sufficiently
196 consider already existing research software repositories stored on platforms like GitHub or or-
197 ganizational GitLabs. We proposed the cascading search, a novel approach that has the potential
198 to compensate for the shortcomings of current approaches by searching for software repositories
199 first and then linking them to their corresponding publications. We presented Betty's research
200 engine, an implementation of the cascading search that faces the challenges of individual rate
201 limitations and restricted access by a client based design. We elaborated how the cascading
202 search can be performed in parallel processes while maintaining a sound, coherent structure.
203 We analyzed 400 random search results from the domain of ecology and found that (a) the cas-
204 cading search is a valid, useful alternative to current approaches, (b) the rule-based approaches
205 presented in this paper yield an accuracy of 35,25% and (c) the majority of references (63%) to
206 corresponding publications / entries in third party databases are written in plain, unstructured
207 text. Reflecting our results we found that a rule-based approach (like the one we use now)
208 might not be well suited for the problem at hand and a data-driven approach might be more
209 promising for identifying references to corresponding publications. Compared to community
210 driven approaches like Papers with Code we reflected on the importance of categorisation and
211 benchmarking of research software and concluded that further work on Betty's research engine
212 is necessary.

213 7 Acknowledgements

214 This work would not have been possible without the financial support from the German Research
215 Foundation (DFG) - project number 442146713.

216 We are grateful to all of those with whom we had the pleasure to work during this and other
217 related projects.

218 8 Roles and contributions

219 **Vasily Seibert:** Conceptualization, Writing – original draft

220 **Andreas Rausch:** Conceptualization, Writing – original draft

221 **Stefan Wittek:** Conceptualization, Writing – original draft

222 References

- 223 [1] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, *et al.*, “The fair guiding principles
224 for scientific data management and stewardship,” *Scientific data*, vol. 3, no. 1, pp. 1–9,
225 2016.
- 226 [2] C. Du, J. Cohoon, P. Lopez, and J. Howison, “Understanding progress in software citation:
227 A study of software citation in the cord-19 corpus,” *PeerJ Computer Science*, vol. 8,
228 e1022, 2022.
- 229 [3] L. L. Wang, K. Lo, Y. Chandrasekhar, *et al.*, “Cord-19: The covid-19 open research
230 dataset,” *ArXiv*, 2020.
- 231 [4] A. Culina, I. van den Berg, S. Evans, and A. Sánchez-Tójar, “Low availability of code in
232 ecology: A call for urgent action,” *PLoS Biology*, vol. 18, no. 7, e3000763, 2020.
- 233 [5] *Papers with code*. [Online]. Available: <https://paperswithcode.com/>.
- 234 [6] *Zenodo*. [Online]. Available: <https://zenodo.org/>.
- 235 [7] S. Druskat, N. C. Hong, R. Haines, and J. Baker, “Citation file format (cff)-specifications,”
236 *Zenodo. data*, 2018.
- 237 [8] M. Barker, N. P. Chue Hong, D. S. Katz, *et al.*, “Introducing the fair principles for research
238 software,” *Scientific Data*, vol. 9, no. 1, pp. 1–6, 2022.
- 239 [9] *Github*. [Online]. Available: <https://github.com>.
- 240 [10] *Gitlab*. [Online]. Available: <https://about.gitlab.com>.
- 241 [11] *Open alex*. [Online]. Available: <https://openalex.org>.
- 242 [12] *Data cite*. [Online]. Available: <https://datacite.org>.
- 243 [13] *Open citations*. [Online]. Available: <https://opencitations.net>.