

# Towards Transparency and Knowledge Exchange in AI-assisted Data Analysis Code Generation

Robert Haase<sup>1,2</sup>

<sup>1</sup>Data Science Center, Leipzig University, Humboldtstraße 25, 04105  
Leipzig, Germany.

<sup>2</sup>Center for Scalable Data Analytics and Artificial Intelligence  
(ScaDS.AI) Dresden / Leipzig.

Contributing authors: [robert.haase@uni-leipzig.de](mailto:robert.haase@uni-leipzig.de);

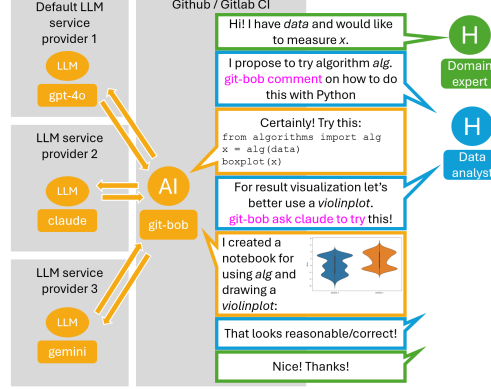
## Abstract

Todo: Enter standfirst here.

Generative artificial intelligence (AI) and Large Language Models (LLMs) in particular are changing the way we do data science. Most prominently, scientists use the technology for interacting with scientific data [1], answer data analysis questions [2, 3], generate data analysis code [4–6], and [re-]write scientific manuscripts [7]. Unfortunately, the prompts sent to LLMs are commonly not conserved, and thus, at the time of publication, it might be hard to differentiate human-made and AI-generated parts of the scientific work. A professional peer-review system, for documenting how LLM-generated code was prompted for, and which human reviewed it, is not established in contemporary scientific culture. However, such systems do exist for collaborative code editing involving multiple humans. For example, the source code repositories GitHub and GitLab are well-established in the open-source software community for discussing issues and potential solutions, building code together, and for peer-reviewing contents. As it was shown before that LLMs can solve real-world GitHub issues [8], developing an AI-assistant that interacts with humans directly within the GitHub platform is the obvious next step.

Here, I present git-bob, an implementation of an LLM-based AI-assistant that can respond to GitHub issues, discuss potential solutions with humans iteratively, write code for them, and submit it as pull-request to be reviewed by humans. It is technically similar to various online services for data analysis such as the OpenAI ChatGPT Data

Analyst or GitHub Copilot Workflows, with three major differences. First, multiple humans can interact with git-bob in one communication thread. This allows bringing together domain specialists, such as life scientists, data-analysts and the AI-assistant in one discussion, stimulating knowledge exchange on how to interact properly with the AI-assistant. Second, discussions with git-bob and resulting code modifications are conserved in an online-platform that others can read and follow, making the interaction with the AI-assistant fully transparent. Third, git-bob is completely open-source. Other developers can read its built-in system prompts and modify them to their needs.



**Fig. 1 Schematic view of the interaction with git-bob:** In one discussion thread, multiple humans (H) can interact with the AI-assistant, which is installed in GitHub or GitLab Continuous Integration (CI) infrastructure. Depending on how it is triggered, the AI-assistant may use LLMs from multiple LLM service providers in the background.

A schematic workflow involving git-bob is depicted in Figure 1: a domain expert opens an issue, a type of discussion thread, on a repository on GitHub, where git-bob is installed. A repository member can then add more information to the request and trigger git-bob to answer by writing a command such as “git-bob comment”. If externals try so, an automatic response will inform them that only repository members are allowed to trigger git-bob because running git-bob may cause costs for repository owners. Once triggered, git-bob will use an LLM to respond to the question, potentially including a code snippet and resulting plots or images. Users and the AI-assistant can then discuss back and forth until some potential solution is reached. This way, good-scientific-practice can be maintained by involving not just domain experts but also data analysis experts in the discussion. Optionally, git-bob can then be asked to implement the solution and send a GitHub pull-request, another type of discussion thread, but accompanied by file modifications to the repository, for instance, including a Jupyter Notebook containing the previously discussed code solution to a given issue. A human would need to review this pull-request and merge it into the code base of the repository. Git-bob also has the capability to review pull-requests originating from humans, but it is not allowed to merge them. This reflects established practices in

science, where eventually a scientist is responsible for data analysis code that becomes part of the project.

Common tasks git-bob is capable of are:

- Giving advice on how to solve a data analysis or data visualization task ([Supplementary Figure 1](#))
- Supporting users of open source libraries by providing advice and code examples, as shown in [Supplementary Figure 2](#). As prompt engineering techniques have the potential to decrease wrong answers and hallucinations [9], also git-bob can be instructed to forward the question to a human in case of doubt ([Supplementary Figure 3](#)). It shall be noted that there is no guarantee that the LLM makes this choice with perfect accuracy.
- Documenting code ([Supplementary Figure 4](#)). Such a task can be time-consuming when performed without AI-assistance, which can generate documentation for multiple Python functions in seconds to minutes.
- Analysing data in the repository directly, and summarizing and plotting data in CSV files ([Supplementary Figure 5](#)).
- Assisting in writing if manuscript files are stored in a GitHub repository, for example in latex format, git-bob can assist in writing. For example, the abstract for this manuscript was written by the AI-assistant and this is documented transparently as shown in [Supplementary Figure 6](#).

A highlight of git-bob is that a local installation is not required. Git-bob is implemented as GitHub workflow or GitLab pipeline, which can be installed by uploading a configuration file to a repository and setting access rights. It is compatible and was tested with the commercial LLMs OpenAI’s GPT4-omni, Anthropic’s Claude 3.5 Sonnet, Google Gemini 1.5 Pro 002, and freely available models hosted on GitHub Models Marketplace. Git-bob reports which model was used in all of its messages, as good scientific practice suggests. Obviously, the communication with the selected LLM is transmitted to the service provider, including source code files from the repository and images provided with the GitHub issue. Hence, users are recommended to not submit any personal or sensitive information. When writing data analysis code, git-bob is intrinsically limited by the capabilities of the used LLM. For example, it has been shown that state-of-the-art (SOTA) LLMs can solve bio-image analysis questions by generating functionally correct code just above 50% of tested cases [5]. This fundamental limitation may disappear when improved LLMs are published. For now, this can be evaded by the humans guiding the AI-assistant in multi-turn interactions towards a workable solution. Further technical limitations arise from prompt-length limitations of the underlying LLMs. When modifying or generating a file, these files must be below specified limits, for example GPT4-omni has 128k tokens input and 16k output tokens as limit (1 token  $\approx$  approx. 3/4 words). Also when processing data, limitations of the GitHub IT infrastructure have to be considered: Git-bob executed in public repositories runs on virtual machines with 4 CPU cores, 16 GB of RAM and 14 GB of SSD storage. In private repositories, only 2 CPU cores and 7 GB RAM are available [10]. More capable systems are available on a paid basis. By installing git-bob in an institutional GitLab server, users can setup freely chosen hardware to run git-bob on.

**Supplementary information.** Supplementary Figures are available in a separate document.

**Acknowledgements.** I would like to thank Elena Katharina Nicolay (UFZ Leipzig) for testing git-bob in its early days and for providing constructive feedback on the manuscript. I also would like to thank Volker Hilsenstein for pushing for GitLab interoperability.

## Declarations

- Funding: I acknowledge the financial support by the Federal Ministry of Education and Research of Germany and by Sächsische Staatsministerium für Wissenschaft, Kultur und Tourismus in the programme Center of Excellence for AI-research "Center for Scalable Data Analytics and Artificial Intelligence Dresden/Leipzig", project identification number: ScaDS.AI. I also acknowledge financial support from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under the National Research Data Infrastructure – NFDI 46/1 – 501864659 - NFDI4BioImage.
- Conflict of interest/Competing interests: The authors declare no conflict of interest.
- Ethics approval and consent to participate: Not applicable.
- Consent for publication: Consent.
- Data availability: Not applicable.
- Materials availability: Not applicable.
- Code availability: The complete source code of git-bob is available online: <https://github.com/haesleinhuepf/git-bob> . The manuscript is available openly too: <https://github.com/haesleinhuepf/git-bob-manuscript>
- Author contribution: RH wrote the software and the manuscript. Some portions of code and text were generated using generative artificial intelligence. The git-history of the both above mentioned repositories informs about which parts were written by RH: Parts which parts were AI-generated are authored by github-actions[bot].

# Supplementary Information

## References

- [1] Royer, L.A.: The future of bioimage analysis: a dialog between mind and machine. *Nature Methods* **20**(7), 951–952 (2023) <https://doi.org/10.1038/s41592-023-01930-y>
- [2] Lai, Y., Li, C., Wang, Y., Zhang, T., Zhong, R., Zettlemoyer, L., Yih, S.W.-t., Fried, D., Wang, S., Yu, T.: DS-1000: A Natural and Reliable Benchmark for Data Science Code Generation (2022)
- [3] Lei, W., Fuster-Barceló, C., Reder, G., *et al.*: Bioimage.io chatbot: a community-driven ai assistant for integrative computational bioimaging. *Nature Methods* **21**, 1368–1370 (2024) <https://doi.org/10.1038/s41592-024-02370-y>
- [4] Royer, L.A.: Omega — harnessing the power of large language models for bioimage analysis. *Nature Methods* **21**(8), 1371–1373 (2024) <https://doi.org/10.1038/s41592-024-02310-w>
- [5] Haase, R., Tischer, C., Hériché, J.-K., Scherf, N.: Benchmarking large language models for bio-image analysis code generation. *bioRxiv* (2024) <https://doi.org/10.1101/2024.04.19.590278>  
<https://www.biorxiv.org/content/early/2024/04/25/2024.04.19.590278.full.pdf>
- [6] Chen, M., Tworek, J., Jun, H., Yuan, Q., Oliveira Pinto, H.P., *et al.*: Evaluating large language models trained on code. *CoRR* **abs/2107.03374** (2021) [2107.03374](https://arxiv.org/abs/2107.03374)
- [7] Lu, C., Lu, C., Lange, R.T., Foerster, J., Clune, J., Ha, D.: The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292* (2024)
- [8] Jimenez, C.E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., Narasimhan, K.: SWE-bench: Can Language Models Resolve Real-World GitHub Issues? (2024). <https://arxiv.org/abs/2310.06770>
- [9] Yin, Z., Sun, Q., Guo, Q., Wu, J., Qiu, X., Huang, X.: Do Large Language Models Know What They Don’t Know? (2023). <https://arxiv.org/abs/2305.18153>
- [10] GitHub: About GitHub-hosted Runners <https://docs.github.com/en/actions/using-github-hosted-runners/using-github-hosted-runners/about-github-hosted-runners>. Accessed: 2024-10-14 (2024). <https://docs.github.com/en/actions/using-github-hosted-runners/using-github-hosted-runners/about-github-hosted-runners>

A)

[Plotting question]: Pair-plot #48

[Open](#) haesleinhuepf opened this issue last week · 1 comment

Assignees



haesleinhuepf commented last week

Owner · ...

Simulate a dataset of 100 patients and three measurements. In a dataframe, the patients are rows and the measurements are columns. There is a fourth column where there is the stage classification (1, 2 or 3). Draw a pairplot of these measurements against each other and use the stage classification as hue.

[git-bob](#) try to do this.



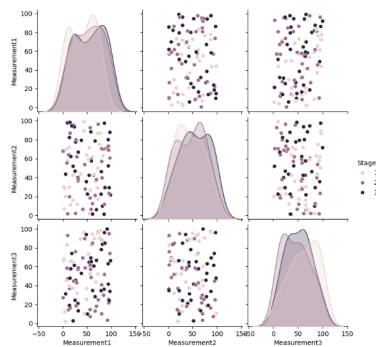
[haesleinhuepf](#) self-assigned this last week

github-actions bot commented last week

Contributor · ...

This message was generated by [git-bob](#) (version: 0.11.0, model: gpt-4o-2024-08-06, [log](#)), an experimental AI-based assistant. It can make mistakes and has [limitations](#). Check its messages carefully.

The changes made to the haesleinhuepf/git-bob-playground repository addressed Issue #48 by creating a Jupyter Notebook, which simulates a dataset of 100 patients with three measurements and a stage classification. The notebook generates a pair-plot of these measurements with the stage classification as the hue. Additionally, the pair-plot created by the notebook is saved as an image file.



- An image of the pair-plot showing the simulated data with stage classification as hue was added.

- [pairplot\\_notebook.ipynb](#) A Jupyter Notebook was created to simulate patient data and generate a pair-plot.



B)

[git-bob-playground / pairplot\\_notebook.ipynb](#)

github-actions[bot] A Jupyter Notebook for si... dafa908 · last week

227 lines (227 loc) · 231 KB

Preview Code Blame



## Pair-Plot of Simulated Patient Data

In this notebook, we'll simulate a dataset of 100 patients with three measurements and a stage classification. We'll then create a pair-plot to visualize these measurements with their stage classifications as the hue.

### Step 1: Import Libraries

...

### Step 2: Simulate the Dataset

We'll create a DataFrame with 100 patients, with three random measurements, and randomly assigned stages (1, 2, or 3).

In [2]:

```
# Seed for reproducibility
np.random.seed(0)

# Simulate measurements
measurements = np.random.rand(100, 3) * 100 # Example range from 0

# Simulate stage classification
stages = np.random.choice([1, 2, 3], size=100)

# Create DataFrame
df = pd.DataFrame(measurements, columns=['Measurement1', 'Measurement2', 'Measurement3'])
df['Stage'] = stages
```

...

### Step 3: Create the Pair-Plot

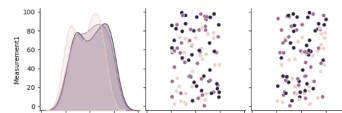
Now, we'll create a pair-plot using Seaborn, setting the `hue` parameter to the stage classification.

In [3]:

```
# Create the pairplot
pairplot = sns.pairplot(df, hue='Stage')

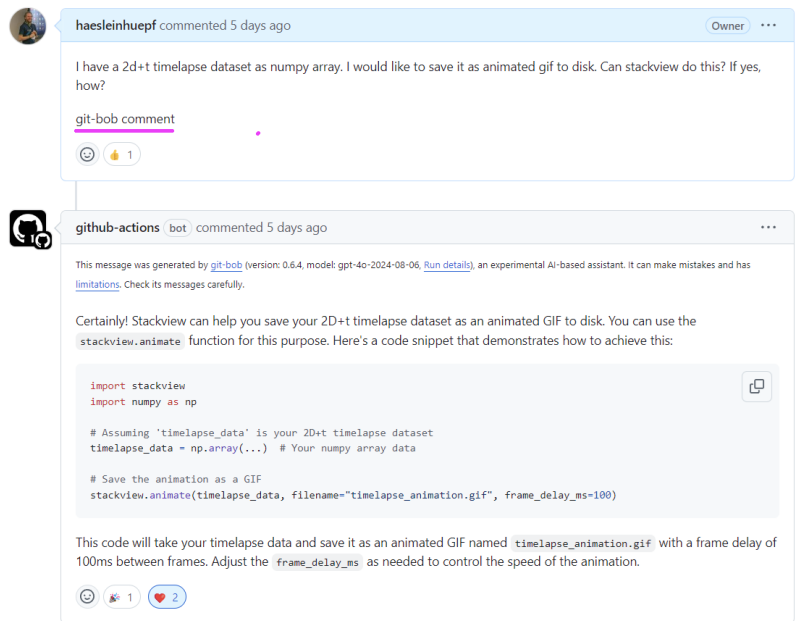
# Save the pairplot
pairplot.savefig('pairplot.png')

# Show the plot
plt.show()
```

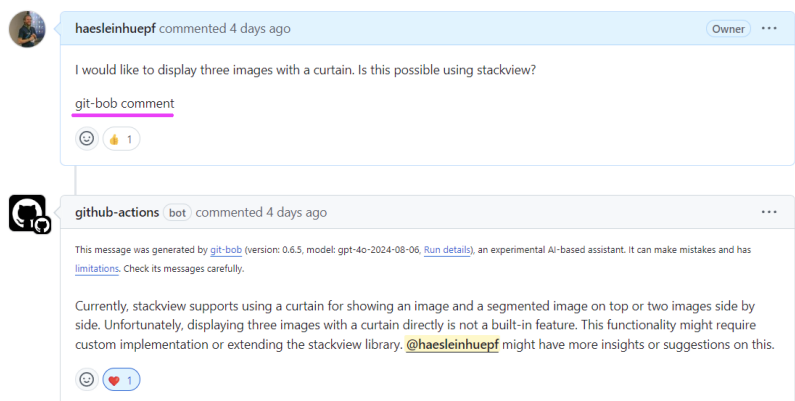


...

**Supplementary Figure 1 Use-case example for generating data analysis code:** The user explains a scenario (A) and triggers git-bob (underlined in magenta). The AI-assistant generates and executes code and visualizes the resulting plot. The user can click on the link to the generated notebook (underlined in green) to go to the notebook (B) and read the code and see intermediate results. The shown notebook is an excerpt as indicated by "...". The entire discussion and corresponding code can be read online: <https://github.com/haesleinhuepf/git-bob-playground/issues/48>



**Supplementary Figure 2 Use-case example for supporting users:** The assistant can be configured to act as expert on a specific Python library and answer user questions. Words triggering git-bob are underlined in magenta. The entire discussion and corresponding code can be read online: <https://github.com/haesleinhuepf/stackview/issues/79>



**Supplementary Figure 3 Use-case example for asking an expert:** The answer to the question shown here is "No", but this is nowhere written in the documentation or the configuration of the assistant. In this case the assistant is not sure, and it can be configured to forward a question to a maintainer of the library where the question arrived. Words triggering git-bob are underlined in magenta. The entire discussion and corresponding code can be read online: <https://github.com/haesleinhuepf/stackview/issues/80>

A)

## Add numpy-style docstrings to \_github\_utilities.py #29

Merged haesleinhuepf merged 1 commit into main from mod-ajh4CjX8wR on Aug 2

Conversation 0 Commits 1 Checks 0 Files changed 1

github-actions (bot) commented on Aug 2 Contributor ...

This comment was generated by [git-bob](#), an AI-based assistant.

Added numpy-style docstrings for all functions in `src/git_bob/_github_utilities.py`. Closes [#22](#).

Update file content d123ac1

haesleinhuepf merged commit dafef0ec into main on Aug 2 Revert

B)

Conversation 0 Commits 1 Checks 0 Files changed 1

Changes from all commits File filter Conversations Jump to Ask Copilot Review

src/git\_bob/\_github\_utilities.py

☐ Viewed ...

```
@@ -6,6 +6,18 @@
6 6  from github import Github
7 7
8 8  def add_comment_to_issue(repository, issue, comment):
9 9  + """
10 10 +     Add a comment to a specific GitHub issue.
11 11 +
12 12 +     Parameters
13 13 +     -----
14 14 +     repository : str
15 15 +         The full name of the GitHub repository (e.g., "username/repo-name").
16 16 +     issue : int
17 17 +         The issue number to add a comment to.
18 18 +     comment : str
19 19 +         The comment text to add to the issue.
20 20 +     """
9 21  print(f"-> add_comment_to_issue({repository}, {issue}, ...)")
10 22
11 23  access_token = os.getenv('GITHUB_API_KEY')
```

**Supplementary Figure 4 Use-case example for adding and revising documentation in code:** git-bob was used to partially write the code documentation of its own code. When asked to add documentation in a specific format, it sent a pull-request (A) and the human could inspect the code modifications (B, excerpt) before merging the code into the project's code base. The entire discussion and corresponding code can be read online: <https://github.com/haesleinhuepf/git-bob/pull/29>






haesleinhuepf commented 3 minutes ago

Member

...


In the folder "download\_statistics", there are multiple files named for example "20240628.csv". The filename without ending is obviously a date, the creation date of the data. Write a jupyter notebook that opens all these files using pandas and plots over time how many records are in these files. Save the notebook in the notebooks directory and make sure the relative link the download\_statistics folder is correct.

git-bob try again

1

---



github-actions bot commented 2 minutes ago

Contributor

...

This message was generated by [git-bob](#) (version: 0.8.2, model: github\_models:gpt-4o, [log](#)), an experimental AI-based assistant. It can make mistakes and has [limitations](#). Check its messages carefully.

In this update, a Jupyter notebook was created to read CSV files and plot the number of records over time, resulting in the creation of [notebooks/plot\\_zenodo\\_links\\_over\\_time.ipynb](#). Additionally, an image file was generated:



Date	Number of Records
2024-07-01	34
2024-07-15	34
2024-08-01	37
2024-08-15	40
2024-09-01	43
2024-09-15	43
2024-10-01	59



**Supplementary Figure 5 Use-case example for plotting data:** after explaining the assistant the folder structure of the project, it generates code for parsing a folder of CSV files and plotting results. Words triggering git-bob are underlined in magenta. The entire discussion and corresponding code can be read online: <https://github.com/NFDI4BIOIMAGE/training/issues/250>

A)

### Add abstract #8

Open haesleinhuepf opened this issue 5 minutes ago · 2 comments

haesleinhuepf commented 5 minutes ago

We need an abstract for the manuscript in the git-bob-manuscript.tex file. The manuscript will be submitted to a journal where life scientists are the target audience. I'd like to have an abstract focusing on the scientific challenges (transparency and knowledge exchange) when using LLMs. git-bob should be mentioned in the second half of the abstract.

git-bob think about this and answer in plain text with no markdown or tex styling.

github-actions bot commented 4 minutes ago

This message was generated by git-bob (version: 0.9.0, model: gpt-4o-2024-08-06, [log](#)), an experimental AI-based assistant. It can make mistakes and has [limitations](#). Check its messages carefully.

An abstract focusing on the scientific challenges of using LLMs in life sciences is provided below:

The integration of Large Language Models (LLMs) in scientific research presents both opportunities and challenges for life scientists. Key challenges include ensuring transparency in AI-generated content and facilitating efficient knowledge exchange among researchers. These issues arise from the opaque nature of AI-driven code generation and the informal sharing of AI insights, which may hinder reproducibility and collaboration. This paper introduces git-bob, an innovative AI-assistant designed to address these challenges by fostering an interactive and transparent collaboration platform within GitHub. By enabling seamless dialogue between humans and AI, git-bob ensures that AI contributions are transparent and verifiable. Moreover, it supports collaborative knowledge exchange, enhancing the interdisciplinary dialogue necessary for cutting-edge life sciences research. The open-source nature of git-bob further promotes accessibility and customization, positioning it as a vital tool in employing LLMs responsibly and effectively within scientific communities.

haesleinhuepf commented 2 minutes ago

Awesome! Add this abstract to the manuscript!

git-bob implement this.

github-actions bot mentioned this issue 1 minute ago

Add abstract to manuscript highlighting key challenges and git-bob solution. #9

Open

B)

### Add abstract to manuscript highlighting key challenges and git-bob solution. #9

Open

github-actions bot commented 2 minutes ago

This message was generated by git-bob (version: 0.9.0, model: gpt-4o-2024-08-06, [log](#)), an experimental AI-based assistant. It can make mistakes and has [limitations](#). Check its messages carefully.

Added an abstract focusing on the scientific challenges of using LLMs in life sciences to the git-bob-manuscript.tex file, addressing issue #8 by highlighting transparency, knowledge exchange, and introducing git-bob as a solution.

closes #8

I added an abstract to the manuscript.

C)

Conversation 0

Commits 1

Checks 0

Files changed 1

Changes from all commits

File filter

Conversations

Ask Copilot

Review

Jump to

git-bob-manuscript.tex

Viewed

@@ -102,7 +102,7 @@

102 102

103 103

104 104

105 105

106 106

107 107

108 108

\begin{abstract}

- Abstract will be added later

+ The integration of Large Language Models (LLMs) in scientific research presents both opportunities and challenges for life scientists. Key challenges include ensuring transparency in AI-generated content and facilitating efficient knowledge exchange among researchers. These issues arise from the opaque nature of AI-driven code generation and the informal sharing of AI insights, which may hinder reproducibility and collaboration. This paper introduces git-bob, an innovative AI-assistant designed to address these challenges by fostering an interactive and transparent collaboration platform within GitHub. By enabling seamless dialogue between humans and AI, git-bob ensures that AI contributions are transparent and verifiable. Moreover, it supports collaborative knowledge exchange, enhancing the interdisciplinary dialogue necessary for cutting-edge life sciences research. The open-source nature of git-bob further promotes accessibility and customization, positioning it as a vital tool in employing LLMs responsibly and effectively within scientific communities.

\end{abstract}

**Supplementary Figure 6 Use-case example for working on scientific manuscripts:** after a first draft of the manuscript was written, git-bob was asked to formulate an abstract (A). The abstract was then submitted as pull-request with a short description (B). The human can also review and potentially modify the proposed text in this online interface (C). Words triggering git-bob are underlined in magenta. The entire discussion can be read online: <https://github.com/haesleinhuepf/git-bob-manuscript/issues/8> and <https://github.com/haesleinhuepf/git-bob-manuscript/pull/9>.