# Biologically-oriented mud volcano database: muddy_db

**Alexei Remizovschi** [1], **Rahela Carpa** [Corresp. 2]

[1] Department of Molecular Biology and Biotechnology, Faculty of Biology and Geology, Babes-Bolyai University, Cluj-Napoca, - Select -, Romania

[2] Department of Molecular Biology and Biotechnology, Faculty of Biology and Geology, Babes Bolyai University, Cluj-Napoca, - Select -, Romania

Corresponding Author: Rahela Carpa
Email address: rahela.carpa@ubbcluj.ro

Mud volcanoes (MVs) are naturally occurring hydrocarbon hotbeds with continuous methane discharge, contributing to global warming. They host microbial communities adapted to hydrocarbon oxidation. Given their research value, MVs still represent a niche topic in microbiology and are neglected by hydrocarbon-oriented research. All the data regarding MVs is sporadic and decentralized. To mitigate this problem, we built a custom Natural Language Processing pipeline (muddy_mine), and collected all the available MV data from open-access articles. Based on this data, we built the muddy_db database. The muddy_db represents the first biologically oriented database rendered as a user-friendly web app. This database includes all the relevant MV data, ranging from microbial taxonomy to hydrocarbon occurrence and geology. The muddy_mine and muddy_db tools are licensed under the GPLv3.

muddy_db R Shiny web app: https://muddy-db.shinyapps.io/muddy_db/

muddy_db R package: https://github.com/TracyRage/muddy_db

muddy_mine Conda package: https://github.com/TracyRage/muddy_mine

# Biologically-oriented mud volcano database: muddy_db

**Remizovschi Alexei**[1] **and Rahela Carpa**[1]

[1]**Department of Molecular Biology and Biotechnology, Babeș-Bolyai University, Cluj-Napoca, Cluj, Romania**

Corresponding author:
Rahela Carpa[1]

Email address: rahela.carpa@ubbcluj.ro

## ABSTRACT

Mud volcanoes (MVs) are naturally occurring hydrocarbon hotbeds with continuous methane discharge, contributing to global warming. They host microbial communities adapted to hydrocarbon oxidation. Given their research value, MVs still represent a niche topic in microbiology and are neglected by hydrocarbon-oriented research. All the data regarding MVs is sporadic and decentralized. To mitigate this problem, we built a custom Natural Language Processing pipeline (muddy_mine), and collected all the available MV data from open-access articles. Based on this data, we built the muddy_db database. The muddy_db represents the first biologically oriented database rendered as a user-friendly web app. This database includes all the relevant MV data, ranging from microbial taxonomy to hydrocarbon occurrence and geology. The muddy_mine and muddy_db tools are licensed under the GPLv3.

muddy_db R Shiny web app: https://muddy-db.shinyapps.io/muddy_db/

muddy_db R package: https://github.com/TracyRage/muddy_db

muddy_mine Conda package: https://github.com/TracyRage/muddy_mine

## INTRODUCTION

Mud volcanoes (MVs) represent hydrocarbon discharging landforms (Mazzini and Etiope, 2017). They are distributed worldwide in both marine and terrestrial environments (Milkov, 2000). The most distinctive feature of MVs is recurrent methane emission. Due to methane emissions, MVs contribute extensively to global warming (Etiope et al., 2009).

MV genesis is mainly caused by a naturally mediated process - kerogen maturation (Vandenbroucke and Largeau, 2007). Therefore, the surrounding area of MVs can provide valuable data regarding both aerobic and anaerobic hydrocarbon microbial oxidation (Cheng et al., 2012).

Over the years, MVs research has been mainly focusing on anaerobic oxidation of methane (AOM) and the implicit interaction between sulfate-reducing bacteria and methane oxidizing archaea (ANME) (Bose et al., 2013; Cui et al., 2014). In addition to AOM research, MVs were also investigated in the context of hydrocarbon research. A myriad of MV studies discussed the thermogenic and biogenic origin of the methane (Etiope et al., 2009; Sano et al., 2017).

Despite these studies, the biological aspect of MVs is still a niche and unexplored topic. The biological data regarding MVs are sporadic and mostly biased towards AOM. Even worse, the already available data is not centralized. Due to the lack of the MV-dedicated database, researchers still analyze MV as a singular entity and avoid comparative studies between MVs with geographically distinct locations.

Meanwhile, mainstream biomedical fields have extensively employed natural language processing (NLP) techniques to mine meaningful data from the research articles (Wang et al., 2020). Simultaneously, the number of databases related to biomedical fields is considerable (Luo et al., 2016). Niche environmental science fields have not caught up. Lacking the possibility to mine environmental-oriented articles and build field-specific databases delays the publication of the meta-analyses or any comparative studies.

Fortunately, democratic NLP models and tools have been published over the last years. Some of them can be easily used by environmental scientists with limited computer science (CS) experience,

47 for example, the spaCy library, ScispaCy models, and S2ORC database (Honnibal and Johnson, 2015;
48 Neumann et al., 2019; Lo et al., 2020).
49 Cumulatively, the latest advancements in NLP can provide opportunities for consolidating and
50 promoting niche environmental topics such as MV microbiology.
51 Given these facts, we built the first biologically oriented mud volcano database, muddy_db, a niche
52 database that consolidates all the relevant biological data, which will be of great use for researchers
53 specializing in bacterial hydrocarbon oxidation or MV microbiology. Collaterally, our custom pipeline
54 can serve as a methodological blueprint for research collectives interested in NLP and building their own
55 specialized databases.

## METHODS

57 To collect all the available data regarding the biological aspects of MVs, we had to exclusively rely on
58 open-access articles. Having these articles, we could freely mine all the biologically oriented tokens,
59 including taxonomy-, chemicals-, geology-, and MV-specific terms. Additionally, we built a custom
60 mining pipeline - muddy_mine (Fig. 1). The scope of muddy_mine is to provide and enrich the muddy_db
61 database with relevant MV data.

### Data collection

63 We used the S2ORC (20200705v1) database to collect open-access articles. S2ORC represents a central-
64 ized database that includes 12.7 million articles with a fully preserved paper structure. S2ORC is quite
65 comprehensive and includes niche environment science articles (Lo et al., 2020). Given these facts, we
66 extracted all the available MV-related titles (N=118 total, N=115 deduplicated) from the S2ORC.

### Token extraction and muddy_mine pipeline

68 Having MV articles, we proceeded with token extraction (i.e. the extraction of the terms of interest) using
69 the muddy_mine pipeline.
70 Taxonomy extraction represented a difficult challenge due to the fact that we intended to collect
71 as many tokens as possible. To overcome this problem, we used the spaCy library (2.3.2), third-party
72 ScispaCy NLP models (0.3.0), and the most recent NCBI Taxonomy database (20 November, 2020)
73 (Honnibal and Johnson, 2015; Neumann et al., 2019; Schoch et al., 2020). First, we extracted all the taxon
74 tokens using en_core_sci_sm ScispaCy model (2.2.5). Second, we checked those tokens against a local
75 NCBI Taxonomy database. Third, we counted the extracted tokens. The higher the counting number, the
76 more likely the token was explicitly discussed in the article.
77 By interating this third-step algorithm, we managed to centralize MV-specific taxonomy on all the
78 possible levels: phylum, class, order, family, and genus.
79 The other non-taxonomy tokens were also extracted with the above-mentioned algorithm. We extracted
80 and counted the tokens related to the following categories: chemistry (inorganic ions, hydrocarbons),
81 geology (geological periods, minerals), MV terminology (ANME, methanogenesis type), and experimental
82 methods (PCR types, amplified genes, chromatography). The comprehensive list of categories can be
83 consulted by visiting the muddy_db repository.
84 The raw output of the muddy_mine pipeline represents a set of csv tables with MV data.

### Building muddy_db database

86 By obtaining muddy_mine raw output, we can advance to the next step - building a user-friendly database.
87 To create this kind of database, we created a Shiny web app, entitled muddy_db. In order to build it,
88 we used the following R packages: shiny (1.5.0), semantic.dashboard (0.2.0), and golem (0.2.1) (Filip
89 and Igras, 2021; Chang et al., 2021; Fay et al., 2021). This app includes all the output generated by the
90 muddy_mine pipeline. Specifically, it displays the counted tokens extracted both from the integral article
91 bodies (N=57) and abstracts (N=115). Additionally, we added an annotated map, which displays the
92 geographical distribution of MVs and their affiliated research metadata.

### System requirements

94 We would like to mention that muddy_mine pipeline was designed to run on systems with modest memory
95 requirements. We achieved this feature by using Python generators. Intel Core i3 (3rd Gen) 3217U / 1.8
96 GHz processor (Intel, USA) and 4GB RAM system was used to build muddy_db. The mining process
97 lasted 24 h.

## RESULTS

The ~~scope~~ of the muddy_db is to gather all the available MV biologically relevant data and include it in a user-friendly database (Fig. 2). First, we collected all the known taxa associated with MVs. The muddy_db includes data regarding archaeal and bacterial taxonomy on all the possible taxonomy levels. This particularity can facilitate the detection of microbial consortia patterns. Second, we gather information regarding metabolic pathways, geology, hydrocarbon availability, and experimental methods performed on MV sediments (Fig. 3). This information can guide specialists to implement appropriate research strategies.

## DISCUSSION

MVs are considered to be one of the settings where ~~the~~ early life evolved (Pons et al., 2011). They sustain a plethora of bacterial metabolic pathways, ranging from methane oxidation and synthesis to sulphate reduction (Kleindienst et al., 2014; Cheng et al., 2012). These pathways and their affiliated microbial communities could provide valuable data regarding (1) origin of microbial life, (2) the effect of the naturally occurring methane discharging systems on global warming, (3) the contribution of microbial consortia to oil souring (Gieg et al., 2011; Etiope et al., 2009; Pons et al., 2011) ~~As we can see, the~~ accumulation of MVs data could enhance ~~our~~ knowledge regarding topics that range from fundamental studies to ecology and engineering.

Given these facts, MVs should be the main focus of hydrocarbon-oriented research and ecology. Unfortunately, data regarding the biological aspects of MVs are scarce. Additionally, the data already gathered are not combined in a dedicated database. The lack of a specialized MV database determines mud volcano microbiology to be a niche and neglected topic in environmental science.

Biomedical fields have always represented the cutting-edge subset of natural science, which actively implement CS techniques, and are tightly intertwined with the big data term (Luo et al., 2016). Simultaneously, the implementation of CS methods in niche environmental fields lags. To both apply CS methods in an environmental context and chronically mitigate the data deficient field of MV microbiology, we created a muddy_mine NLP pipeline and muddy_db database.

The creation of muddy_db tools aims to create a platform, that would provide sufficient data to perform meta-analyses or comparative studies of the MVs. Specifically, we hope that muddy_db would facilitate the discovery of atypical taxonomic patterns and point out the influence of geography on MVs characteristics. Simultaneously, muddy_mine represents a reproducible example of a mining technique applied in the context of environmental studies. As a result, muddy_mine could encourage researchers to mine their data of interest, being free from any field of study constraints.

In addition to positive sides of our projects, we would like to address the evident limitations of muddy_mine pipeline. We were limited to only to open-access articles which were found in the S2ORC. Currently, we are constricted to this corpus due to the following facts:

1. Journals use various article xml encoding standards such as TEI-XML and JATS-XML. Therefore, it is difficult to design a universal xml to json parser. As a result, it is challenging to manually create an exhaustive list of structured texts, which are appropriate for the mining process.
2. Unlike the open-access articles, the mining of articles behind the paywall might represent copyright infringement risks.

## CONCLUSIONS

The muddy_db represents the first biologically oriented mud volcano database. It was designed to provide a comprehensive data corpus that can facilitate mud volcano research and shed light on the topic as a whole. The muddy_db contains data ranging from taxonomy to geology and experimental methods. Simultaneously, the muddy_mine NLP pipeline can serve as an example of accessible implementation of NLP techniques in environmental sciences.

## REFERENCES

Bose, A., Rogers, D., Adams, M., Joye, S., and Girguis, P. (2013). Geomicrobiological linkages between short-chain alkane consumption and sulfate reduction rates in seep sediments. *Frontiers in Microbiology*, 4.

148 Chang, W., Cheng, J., Allaire, J., Sievert, C., Schloerke, B., Xie, Y., Allen, J., McPherson, J., Dipert, A.,
149     and Borges, B. (2021). *shiny: Web Application Framework for R*. R package version 1.6.0.
150 Cheng, T., Chang, Y., Tang, S., Tseng, C., Chiang, P., Chang, K., Sun, C., Chen, Y., Kuo, H., Wang,
151     C., Chu, P., Song, S., Wang, P., and Lin, L. (2012). Metabolic stratification driven by surface and
152     subsurface interactions in a terrestrial mud volcano. *ISME J*, 6(12):2280–2290.
153 Cui, M., Ma, A., Qi, H., Zhuang, X., and Zhuang, G. (2014). Anaerobic oxidation of methane: an "active"
154     microbial process. *MicrobiologyOpen*, 4(1):1–11.
155 Etiope, G., Feyzullayev, A., and Baciu, C. (2009). Terrestrial methane seeps and mud volcanoes: A global
156     perspective of gas origin. *Mar Pet Geol*, 26(3):333–344.
157 Fay, C., Guyader, V., Rochette, S., and Girard, C. (2021). *golem: A Framework for Robust Shiny*
158     *Applications*. R package version 0.3.1.
159 Filip, S. and Igras, K. (2021). *semantic.dashboard: Dashboard with Fomantic UI Support for Shiny*. R
160     package version 0.2.0.
161 Gieg, L., Jack, T., and Foght, J. (2011). Biological souring and mitigation in oil reservoirs. *Applied*
162     *Microbiology and Biotechnology*, 92(2):263–282.
163 Honnibal, M. and Johnson, M. (2015). Proceedings of the 2015 conference on empirical methods in
164     natural language processing. Association for Computational Linguistics.
165 Kleindienst, S., Herbst, F., Stagars, M., von Netzer, F., von Bergen, M., Seifert, J., Peplies, J., Amann,
166     R., Musat, F., Lueders, T., and Knittel, K. (2014). Diverse sulfate-reducing bacteria of the desulfos-
167     arcina/desulfococcus clade are the key alkane degraders at marine seeps. *ISME J*, 8(10):2029–2044.
168 Lo, K., Wang, L., Neumann, M., Kinney, R., and Weld, D. (2020). Proceedings of the 58th annual meeting
169     of the association for computational linguistics. Association for Computational Linguistics.
170 Luo, J., Wu, M., Gopukumar, D., and Zhao, Y. (2016). Big data application in biomedical research and
171     health care: A literature review. *Biomed Inform Insights*, 8:BII.S31559.
172 Mazzini, A. and Etiope, G. (2017). Mud volcanism: An updated review. *Earth Sci Rev*, 168:81–112.
173 Milkov, A. (2000). Worldwide distribution of submarine mud volcanoes and associated gas hydrates.
174     *Mar Geol*, 167(1-2):29–42.
175 Neumann, M., King, D., Beltagy, I., and Ammar, W. (2019). Proceedings of the 18th bionlp workshop
176     and shared task. Association for Computational Linguistics.
177 Pons, M., Quitte, G., Fujii, T., Rosing, M., Reynard, B., Moynier, F., Douchet, C., and Albarede, F. (2011).
178     Early archean serpentine mud volcanoes at isua, greenland, as a niche for early life. *Proc. Natl. Acad.*
179     *Sci. U.S.A.*, 108(43):17639–17643.
180 Sano, Y., Kinoshita, N., Kagoshima, T., Takahata, N., Sakata, S., Toki, T., Kawagucci, S., Waseda, A.,
181     Lan, T., Wen, H., Chen, A., Lee, H., Yang, T., Zheng, G., Tomonaga, Y., Roulleau, E., and Pinti, D.
182     (2017). Origin of methane-rich natural gas at the west pacific convergent plate boundary. *Scientific*
183     *Reports*, 7(1).
184 Schoch, C., Ciufo, S., Domrachev, M., Hotton, C., Kannan, S., Khovanskaya, R., Leipe, D., Mcveigh,
185     R., O'Neill, K., Robbertse, B., Sharma, S., Soussov, V., Sullivan, J., Sun, L., Turner, S., and Karsch-
186     Mizrachi, I. (2020). Ncbi taxonomy: a comprehensive update on curation, resources and tools.
187     *Database*, 2020.
188 Vandenbroucke, M. and Largeau, C. (2007). Kerogen origin, evolution and structure. *Org Geochem*,
189     38(5):719–833.
190 Wang, J., Deng, H., Liu, B., Hu, A., Liang, J., Fan, L., Zheng, X., Wang, T., and Lei, J. (2020). Systematic
191     evaluation of research progress on natural language processing in medicine over the past 20 years:
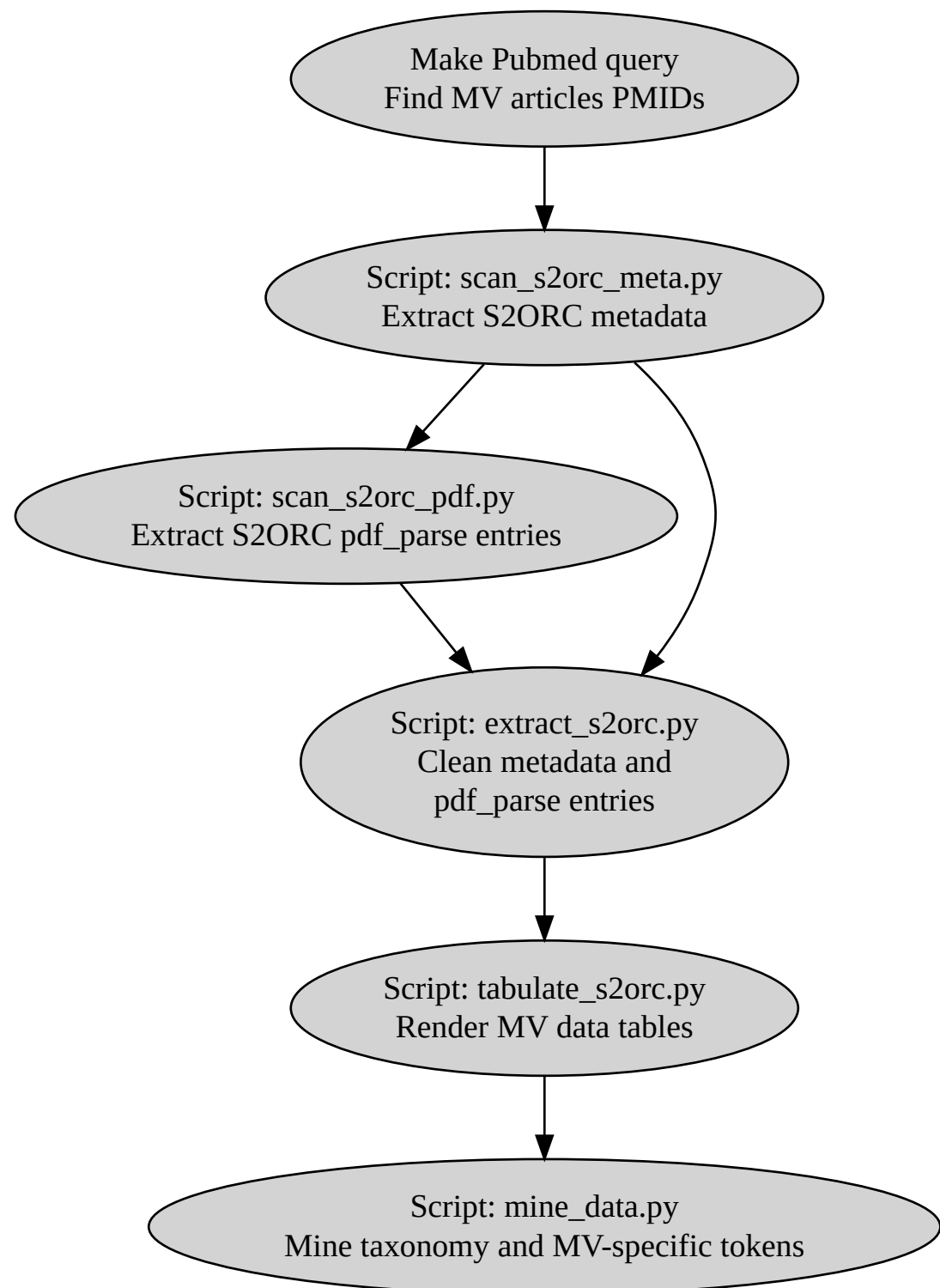192     Bibliometric study on pubmed. *J Med Internet Res*, 22(1):e16816.

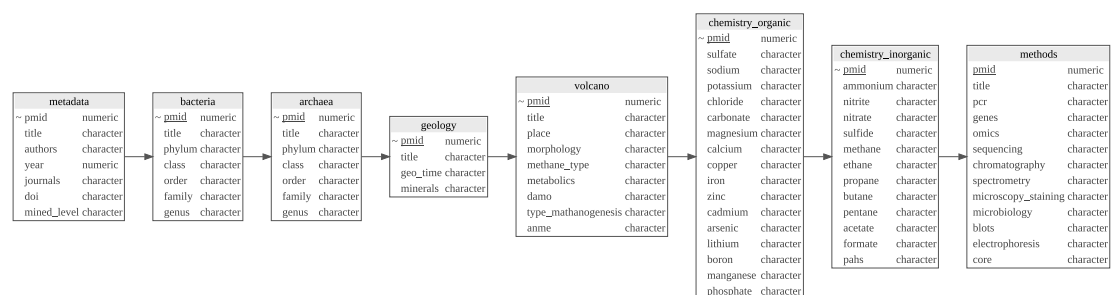**Figure 1.** muddy_mine - pipeline used to build the muddy_db database. MV - mud volcano, PMIDs - Pubmed

**Figure 2.** muddy_db general appearance



**Figure 3.** muddy_db schema