

IUCLID-Decoder: Python code to extract the chemical property data for the substances registered under REACH

Paulina Körner¹ and Juliane Glüge¹

¹ Institute of Biogeochemistry and Pollutant Dynamics, ETH Zurich, 8092 Zurich, Switzerland

DOI: [10.26434/chemrxiv-2024-12345](https://doi.org/10.26434/chemrxiv-2024-12345)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

In partnership with



AMERICAN
ASTRONOMICAL
SOCIETY

This article and software are linked with research article DOI [10.3847/xxxxx](https://doi.org/10.3847/xxxxx) <- [update this with the DOI from AAS once you know it.](#), published in the Astrophysical Journal <- The name of the AAS journal..

Summary

The IUCLID-Decoder is a package in python that enables users to extract chemical property data from registration dossiers that have been submitted to the EU Regulation on Registration, Evaluation, Authorization and Restriction of Chemicals (REACH). The data is extracted from files that the European Chemicals Agency (ECHA) offers for download. These files, the so called 'REACH Study Results' contain more than 4 million IUCLID 6 format (i6d) files with study results. The 'decrypt_reach_results' function in the IUCLID Decoder package extracts and decodes the data and standardizes the units. The output is an SQLite database and (if set to 'True') an excel file with the database.

Statement of need

ECHA publishes on its website ([ECHA, 2024b](#)) the non-confidential substance data that have been submitted by the registrants under REACH. However, ECHA reserves the right to block systematic automated data collection activities including scraping, data mining, and extraction and re-utilization of the whole or a substantial part of the website and the ECHA databases, where justified and subject to applicable law. To be still able to access the data, ECHA offers to download the 'REACH Study Results' via the IUCLID website ([ECHA, 2024a](#)). The REACH Study Results contain results from studies that relate to physicochemical properties, environmental fate and pathways, and ecotoxicology and toxicological information. Data form free text fields are not provided, but most of the other data. As the data of the currently more than 20 000 substances come in over 4 million i6d-files, with a structure similar to xml-files, a script is needed to import them into a database. The IUCLID-Decoder package provided here enables the user to extract the information from the i6d-files and to compile them into a database or optionally an excel file. In addition, the information that is available in the REACH Study Results as numerical code including information on units, methods, reliability or study types is converted into text and the units of the study results are standardized as far as possible.

Description

The REACH Study Results come as zipped i6z-file. When using Linux or Mac, the path to the unzipped i6z-file can be given to the decrypt_reach_results function and the function can be run without considering where to place it. Under Windows, all i6z-files should be unzipped before running the function und the decrypt_reach_results should be place in the same folder as the REACH Study Results. More information on the expected document structure is provided in the GitHub repository.

39 The REACH Study Results contain one main folder per registered substance. This folder then
40 contains the so-called 'manifest' as well as individual study results and the files that are needed
41 to decode the results. The `decrypt_reach_results` function extracts the information property
42 by property. To do this, it first opens the manifest of each substance and looks up the file ID
43 for the property of interest. This property file is then opened and the data for the property
44 of interest are read out. This is done with a recursive function which calls itself up until all
45 the information has been retrieved. Additionally, data on the substance identity such as the
46 European Community (EC) number, the Chemical Abstract Service Registry Number® (CAS
47 RN®) and the substance name are extracted. SMILES are not available and would need to be
48 added later on separately. Before saving the data into a dataframe, the data are decoded, and
49 the units are standardized.

50 It is recommended to check the IUCLID website regularly, as updated REACH Study Results
51 are uploaded once or twice a year. This data can be easily transferred to a new database using
52 the IUCLID decoder package.

53 The `decrypt_reach_results` function has already been used in our previous work where we ana-
54 lyzed the physicochemical property data (Glüge & Scheringer, 2023) and the bioconcentration
55 data (Glüge, Escher, et al., 2023) in the ECHA database. We have also used it in Körner et
56 al. (2024) to check certain REACH study results on ready-biodegradation. Information on
57 SMILES and curated SMILES for most of the organic mono-constituent substances registered
58 under REACH are available in Glüge, McNeill, et al. (2023).

59 Acknowledgements

60 We thank Martin Scheringer for the initial idea of the project and Stefan Glüge for support in
61 between. PK and JG acknowledge funding from the Swiss Federal Office for the Environment.

62 References

- 63 ECHA. (2024a). *IUCLID6 - REACH study results*. [https://iuclid6.echa.europa.eu/
64 reach-study-results](https://iuclid6.echa.europa.eu/reach-study-results)
- 65 ECHA. (2024b). *REACH dissemination platform - registered substances factsheets*. [https:
66 //echa.europa.eu/de/information-on-chemicals/registered-substances](https://echa.europa.eu/de/information-on-chemicals/registered-substances)
- 67 Glüge, J., Escher, B. I., & Scheringer, M. (2023). How error-prone bioaccumulation experiments
68 affect the risk assessment of hydrophobic chemicals and what could be improved. *Integr.*
69 *Environ. Assess. Manag.*, 19(3), 792–803. <https://doi.org/10.1002/ieam.4714>
- 70 Glüge, J., McNeill, K., & Scheringer, M. (2023). Getting the SMILES right: identifying
71 inconsistent chemical identities in the ECHA database, PubChem and the CompTox
72 Chemicals Dashboard. *Environ. Sci. Adv.*, 2(4), 612–621. [https://doi.org/10.1039/
73 D2VA00225F](https://doi.org/10.1039/D2VA00225F)
- 74 Glüge, J., & Scheringer, M. (2023). Evaluation of Physicochemical Property Data in the
75 ECHA Database. *J. Phys. Chem. Ref. Data*, 52(4). <https://doi.org/10.1063/5.0153030>
- 76 Körner, P., Glüge, J., Glüge, S., & Scheringer, M. (2024). Critical insights into data curation
77 and label noise for accurate prediction of aerobic biodegradability of organic chemicals.
78 *Environ. Sci. Adv.*, under revi.