

Supporting Information for

## Making the InChI FAIR and sustainable while moving to Inorganics

Gerd Blanke,<sup>\*[a]</sup> Jan Brammer,<sup>[b]</sup> Djordje Baljovic,<sup>[b]</sup> Nauman Ullah Khan,<sup>[b]</sup> Frank Lange,<sup>[b]</sup> Felix Bänsch,<sup>[c]</sup> Clare Tovee,<sup>[d]</sup> Ulrich Schatzschneider,<sup>[e]</sup> Richard M. Hartshorn,<sup>[f]</sup> Sonja Herres-Pawlis<sup>\*[b]</sup>

### Appendix A. Additional information on programming details

For the first time in three years, a new release of the InChI software (with version number 1.07) was published in December 2023 on a newly created *GitHub* repository.[1] All existing code fragments left by the late developer Igor Pletnev had been integrated into a new code project which served as a foundation for InChI 1.07 development. Since then, more than three thousand bugs, warnings and hints were detected and fixed in the inherited code.

The InChI version 1.07 introduces several new features that can be accessed in test/engineering mode such as:

- support for six additional types of tautomeric transformations
- support for processing of molecules with up to 32766 atoms

Particular emphasis has been given to fixing security issues and bugs, the vast majority of which had been inherited from earlier InChI versions, predominantly 1.06. Moreover, all InChI software versions have been closely monitored by the *Google*<sup>®</sup> *oss-fuzz project* team,[2] who have until now detected more than 60 bugs and security issues using various *fuzz testing* techniques.[3]

The source code has been methodically debugged using 64- and 32-bit versions of all available C/C++ compilers in *Microsoft*<sup>®</sup> *Windows* and *Ubuntu Desktop* computer operating systems. The following table (Table A1) summarizes all C/C++ compilers which were used on *Microsoft*<sup>®</sup> *Windows* and *Ubuntu Linux* platforms for testing and the InChI code development.

C/C++ compiler name	Compiler versions	64-bit	32-bit	Microsoft <sup>®</sup> Windows	Ubuntu Linux
Microsoft <sup>®</sup> Visual Studio C++	17.9.1 - 17.10.4	✓	✓	✓	-
GNU GCC	7.5 - 14.1	✓	✓	✓	✓
LLVM/Clang	14.0.3 - 18.0.4	✓	✓	✓	✓

Table A1. Summary of compiler versions used for testing and developing of InChI

The following list shows a summary of all bug and security fixes:

- 9 functions were rewritten to address issues with function arguments related to arrays of various dimensions
- 25 blocks of code were rewritten to address memory leaks, buffer overruns, various security issues, and improperly written conditional statements or bitwise operations
- 33 heap overflow issues were fixed due to use of large array dimensions
- 157 security bugs related to improper pointer dereferencing were fixed, which might cause crashes or undesired exits
- 71 memory leaks were fixed
- 530 places have been marked where optional bounds-checking functions could be used
- 2480 various bugs and other issues were addressed, such as type conversions and mismatches, removing redundant variables and/or code as well as addressing LLVM/Clang warnings
- 58 *Google*<sup>®</sup> *oss-fuzz* [2] issues were fixed, and 5 additional *Google*<sup>®</sup> *oss-fuzz* issues have been partially fixed

Finally, in order to provide a further layer of security, optional support for bounds-checking functions[4] has been provided (in accordance with *Annex K* of the C11 standard[5]).

Because of the primary focus to fix bugs and security fixes, regression tests were run to compare the results between version 1.06 and 1.07 against all PubChem "Compounds" [6], "3D Compounds" [7], and "Substances" [8]. In terms of the nearly 300 million PubChem "Substances" only 2 failed. In terms of invariance tests 2131 substances failed ( $6,4 \cdot 10^{-5}\%$ ). Most of the failures were caused by insufficient determination of H tautomerism with the structures.

All binaries/API, instructions for compiling InChI software from the source code, a comprehensive list of known issues as well as many other important technical details and the test suite as Docker container can be found on the landing page of the new InChI *GitHub* repository.[1]

## Appendix B Interactive InChI user interfaces

### *InChI's Graphical User Interface for Windows – WInChI*

An interactive tool for Microsoft Windows called *WInChI* has been part of previous releases. It calculates InChIs representations from both MDL Molfiles containing a single chemical structure or SDFfiles containing multiple structures.

After the calculation, the chemical structures are displayed together with the canonical InChI atom numbers (marked in yellow in Figure S1). Tautomeric areas within the molecules are indicated with dashed lines which depict the result of the InChI canonicalization. InChI, AuxInfo and InChIKey are displayed below the canvas with helpful information on each of the layers. *WInChI* calculates the standard InChI by default, but the InChI parameters can be changed in the option menu (Figure S2).

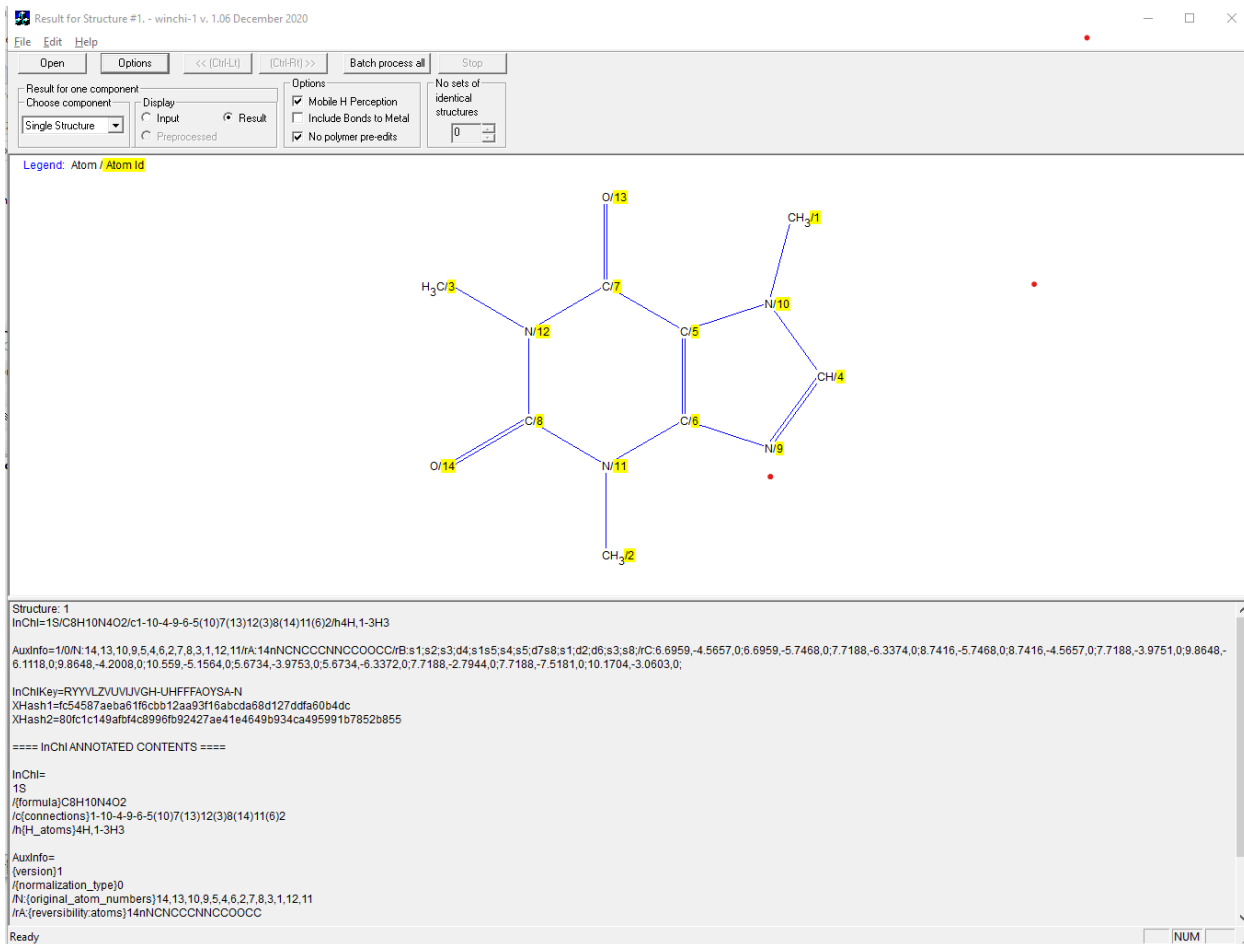


Figure S1: Screenshot of the WinChI canvas

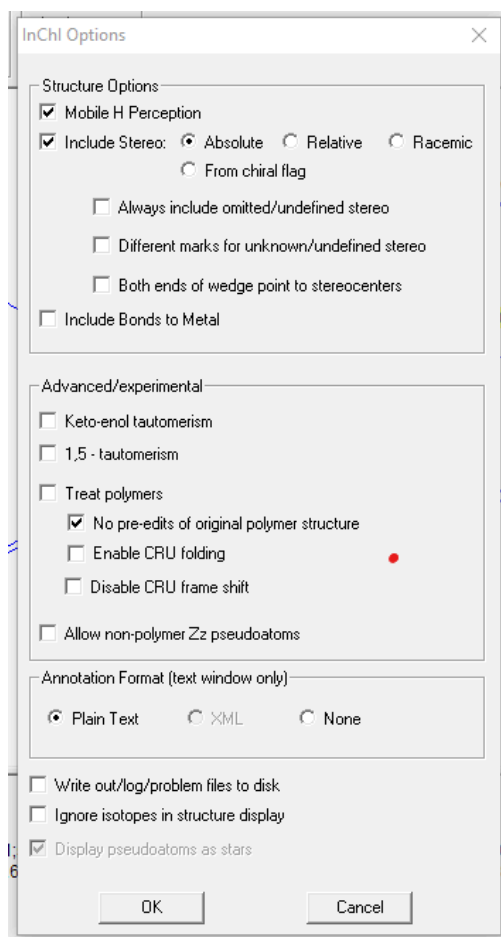


Figure S2: InChI parameter selection tool in WInChI

### *InChI web demo*

This section describes the InChI web demo[9], a new interactive and browser-based application that enables the generation of InChIs and reaction InChIs (RInChIs). Such a browser-native application has the key advantage that there is no need to install any additional software on the user's computer, such as the InChI software suite, WInChI or chemical structure editors that include the InChI library. In addition, the InChI web demo is designed such that all data remains the local browser, thus no web backend infrastructure is required.

The main use case of the application is the generation of InChIs from drawn chemical structures. Structure drawing is facilitated by the powerful Ketcher chemical structure editor[10] and InChI, AuxInfo, InChIKey, and the log of the calculation (for error messages) are shown below the drawing canvas (Figure S3). All InChI parameters and the InChI version (1.07 as well as 1.06 as legacy version) can be selected from the right part of the panel.

InChI RInChI Funding

Draw structure and convert to InChI Convert Molfile to InChI Convert InChI or Auxinfo to structure

Mobile H Perception

Include Stereo:  Absolute  
 Relative  
 Racemic  
 From chiral flag

Always include omitted/undefined stereo  
 Different marks for unknown/undefined stereo  
 Both ends of wedge point to stereocenters

Include Bonds to Metal  
 Keto-enol Tautomerism  
 1,5-tautomerism  
 Treat polymers:  
 No pre-edits of original polymer structure  
 Enable CRU folding  
 Disable CRU frame shift

Allow non-polymer Zz pseudoatoms

Reset InChI Options

Select InChI version

1.06

InChI

InChI=1S/C8H10N4O2/c1-10-4-9-6-5(10)7(13)12(3)8(14)11(6)2/h4H,1-3H3

InChIKey

RYYVLZVUVIJVGH-UHFFFAOYSA-N

AuxInfo

AuxInfo=1/0/N:14,13,10,9,3,4,2,6,7,8,5,1,12,11/rA:14nNCCCNCNCCOCC/rB:s1;s2;d3;s4;s1s5;s4;s3;d7s8;s1;d6;d2;s5;s8;/rC:7.1723,-4.9749,0;8.0385,-4.4748,0;8.9046,-4.9749,0;8.9046,-5.9751,0;8.0385,-6.4752,0;7.1723,-5.9751,0;9.8557,-6.2841,0;9.8557,-4.666,0;10.4435,-5.4751,0;6.3065,-4.475,0;6.3065,-6.475,0;8.0385,-3.4751,0;8.0385,-7.4749,0;10.1144,-3.7002,0;

Log

InChI options:

Figure S3: Screenshot of the InChI web demo demonstrating the drawing of a chemical structure, the immediate generation of InChI representation, InChIKey and AuxInfo as well as the interactive selection of the InChI parameters and version.

As an additional functionality, Molfiles can be generated from an InChI representation or AuxInfo, and the evaluation of chemical reactions to obtain RInChIs and the related Long-, Short-, and Web-RInChIKeys beside the reaction AuxInfo (RAuxInfo), which is depicted in Figure S4. Similar to the functionalities related to InChI and Molfiles, a reverse-generation of RXN/RD files from RInChI representations or RAuxInfo was implemented.

InChI RInChI Funding

Draw reaction and convert to RInChI Convert RXN/RD file to RInChI Convert RInChI to reaction Convert RInChI to RXN/RD file

RInChI

RInChI=1.00.15/C2H6O/c1-2-3/h3H,2H2,1H3!C4H8O2/c1-3(2)4(5)6/h3H,1-2H3,(H,5,6)<>C6H12O2/c1-4-8-6(7)5(2)3/h5H,4H2,1-3H3!H2O/h1H2<>H2O4S/c1-5(2,3)4/h(H2,1,2,3,4)/d+

Long-RInChIKey

Long-RInChIKey=SA-FUHFF-LFQSCWFLJHTTHZ-UHFFFAOYSA-N-KQNPFTWMSNSAP-UHFFFAOYSA-N--WDAXFOBOLVPLV-UHFFFAOYSA-N-XLYOFNOQVPJJNP-UHFFFAOYSA-N--QAQWNCQDCNURD-UHFFFAOYSA-N

Short-RInChIKey

Short-RInChIKey=SA-FUHFF-UGXVRNCBHW-QCBJAAMWEW-QAOWNCQODC-NUHFF-NUHFF-NUHFF-ZZZ

Web-RInChIKey

Web-RInChIKey=MVHGXAUVYQGEJ3NXCQ-NUHFFADPSTJSA

RAuxInfo

RAuxInfo=1.00.1/0/N:1,2,3/rA:3nCCO/rB:s1;s2;/rC:1.5529,-4.5998,0;2.4186,-4.1,0;3.2847,-4.6,0;!1/N:1,3,2,4,5,6/E:(1,2)(5,6)/rA:6nCCCCO0/rB:s1;s2;s4;d4;/rC:-2.9176,-4.5998,0;-2.0518,-4.1,0;-2.0518,-3.0999,0;-1.1857,-4.6,0;-3.196,-4.1,0;-1.1857,-5.6001,0;<0/N:8,1,3,7,2,4,6,5/E:(2,3)/rA:8nCCCCOCC/rB:s1;s2;s4;d4;s5;s7;/rC:7.6606,-4.5998,0;8.5263,-4.1,0;8.5263,-3.0999,0;9.3924,-4.6,0;10.2586,-4.1,0;9.3924,-5.6001,0;11.1246,-4.6,0;11.9907,-4.1,0;10/N:1/rA:1nO/rB:/rC:14.1176,-4.35,0;<1/N:2,3,4,5,1/E:(1,2,3,4)/CRV:5.6/rA:5nS0000/rB:d1;d1;s1;s1;/rC:5.2,-2.8,0;5.25,-3.8,0;5.175,-1.634,0;4.4,-2.816,0;6.15,-2.784,0;

Log

Figure S4: Screenshot of the RInChI view in the InChI web demo. A chemical reaction drawn in the Ketcher structure editor is transformed into a RInChI string, its respective RInChIKeys and RAuxInfo.

On the technical level, the InChI web demo utilizes the WebAssembly technology stack[11] (Wasm) by compiling InChI's C and RInChI's C++ source code to the Wasm binary instruction format using the Emscripten compiler[12]. The design and feasibility of such an application has been demonstrated previously[13,14] and the InChI algorithm has already been integrated in the web versions of major open-source cheminformatics frameworks such as RDKit.js[15] and Indigo[16] *via* the Emscripten tool chain. We would like to emphasize that the port to Wasm opens the InChI algorithm to be used on *any* processor architecture and operating system with a Wasm-compatible web browser.

The InChI web demo will follow the ongoing InChI and RInChI development. Future extensions will introduce a mode to inspect InChI's canonical atom numbering and the labeling of

tautomeric systems similar to WInChI's capabilities, bond types specific to organometallic compounds and atom mappings in reactions for RInChI. For easier use, specific tutorials such as a RInChI tutorial are being developed.

## Appendix C Testing

An important part in our effort of more open and sustainable development of the InChI is to provide a test suite, including test code, data, and documentation. Accordingly, objective and transparent quality criteria were defined that enable convenient and replicable testing. Such a test environment is indispensable for collaborative development, especially with external contributors. Our test suite consists of two kinds of tests.

### a) Regression tests

Regression tests are used to detect problems with the *stability* of the InChI across versions. A regression means that the current behavior of a software package deviates from the established proper behavior of a previous version. While developing a feature or fixing a bug, developers repeatedly run these tests to increase their confidence that no regressions are introduced. In Figure S5, the first and second run represent tests that are conducted after alterations to the codebase. During the first run, the output matches the reference state while the one results in a regression, since the output no longer matches the reference. We use the output of the last stable version as a reference for our regression tests.

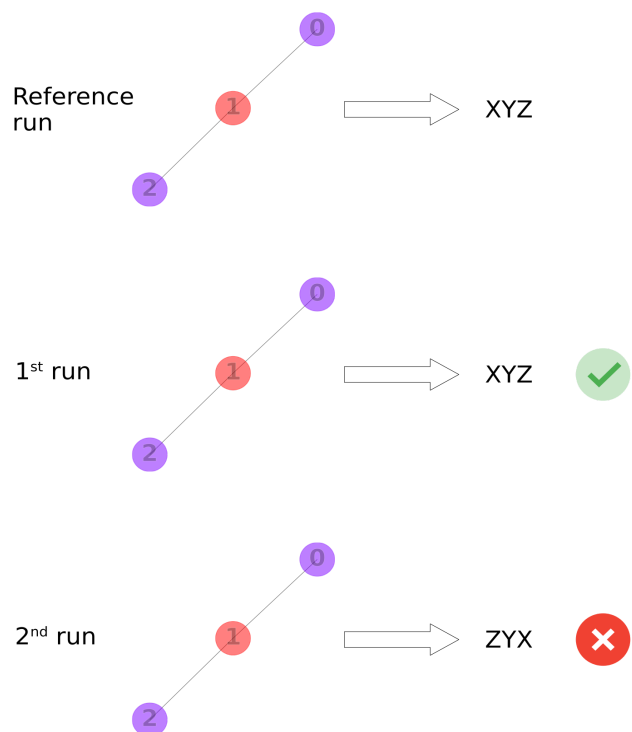


Figure S5: Schematic representation of regression tests: the first test run is successful because the output strings are identical, the second one fails because of a different output

## b) Invariance tests

Invariance tests are used to detect problems with the InChI canonicalization algorithm. As such, they increase confidence in an essential feature of the InChI, the unique and stable identification of structures. During an invariance test, the atom indices of a structure are permuted repeatedly and each permutation is expected to result in the same InChI output. In Figure S6, the output of the third permutation deviates from the previous two permutations, resulting in a failed invariance test.

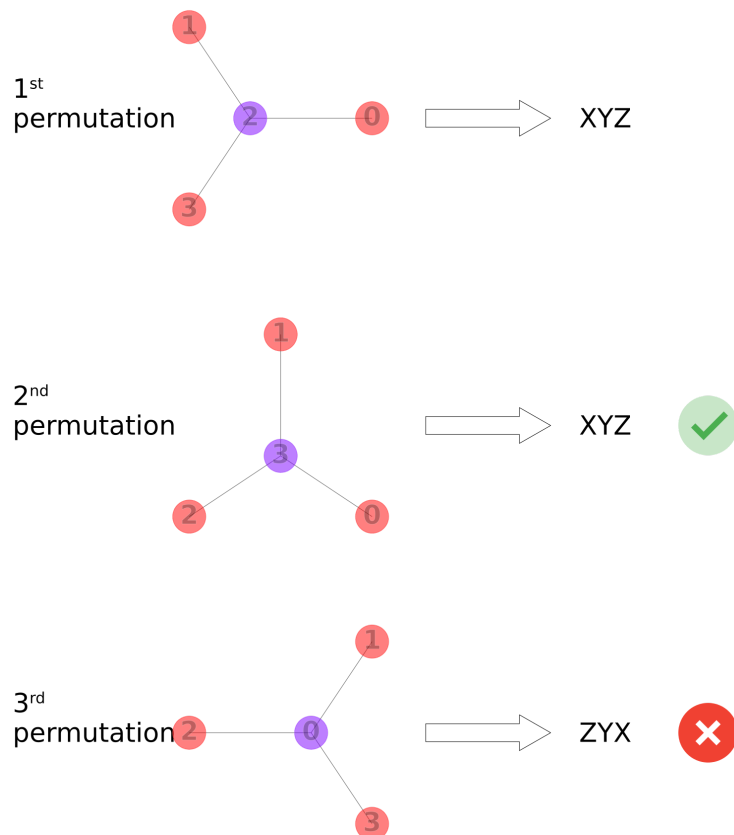


Figure S6: Permutation tests: While the input molecule features a different atom numbering in each test, the test is successful in case the output string remains identical due to a proper canonicalization algorithm. It fails in cases where one output differs from the other permutations.

We run our test suite against two data sets, each serving a unique purpose.

### Automated tests

Every contribution to our GitHub repository, by team members or external collaborators, triggers a run of the test suite. The practice of automatically running tests against each contribution is part of the continuous integration (CI).

To facilitate quick feedback, automated testing needs to be fast but thorough, to increase confidence (note that those two requirements are antagonistic). We ensure speed by limiting the size of our CI data, resulting in a test suite that runs in a matter of seconds. Furthermore, we make the data informative by including known problematic structures, which have triggered bugs in the past, as well as structures that challenge the InChI algorithmically. The CI data is



included in our GitHub repository to facilitate transparency and re-use. For example, developers can run the test suite against the CI data on their local machines, before contributing on GitHub.

### Manual release tests

Prior to each release, the test suite is run manually against the PubChem database. Specifically, tests are run against all PubChem "Compounds"[6], "3D Compounds"[7], and "Substances".[8] The PubChem data are used for testing as is, no pre-processing is applied. These comprehensive test runs are limited to releases, since they require significantly more time than the automated tests. The test results for the current stable release (1.07.0 at the time of writing) can be found here: [https://github.com/JanCBrammer/presentation\\_inchi\\_tests/blob/main/slides.pdf](https://github.com/JanCBrammer/presentation_inchi_tests/blob/main/slides.pdf). Finally, users can run our test suite against their own datasets to pilot new InChI versions before integrating them in their workflows or production environments.

### Appendix D More inorganic examples of bond handling

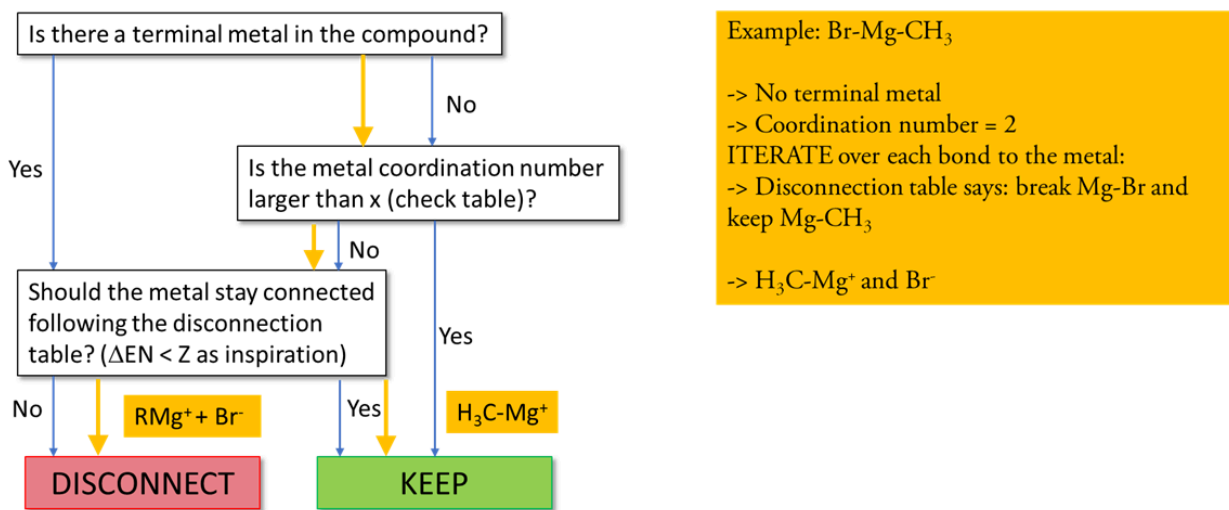
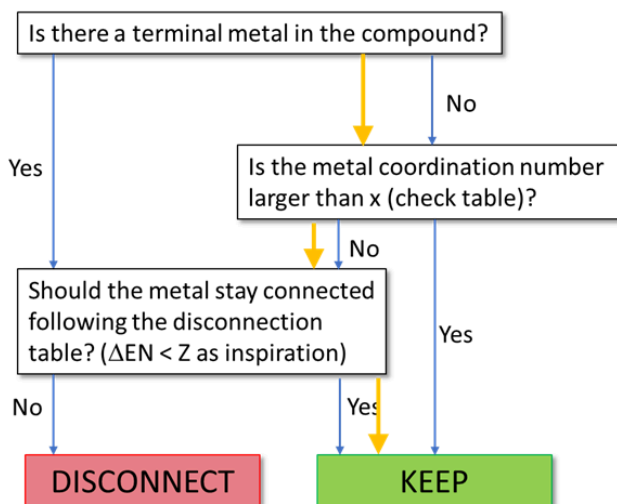


Figure S7: Bond handling for a Grignard compound



Example:  $(\text{Sn}_4\text{Bi}_4)^{2-}$

- > No terminal metal
- > Coordination number: Sn: 4, Bi: 3 (compare to standard valence: Sn: 2, Bi: 3)
- ITERATE over each bond to the metal:
- > Disconnection table says: keep all bonds

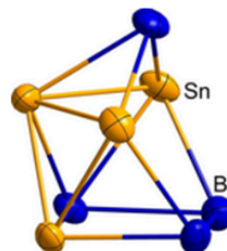
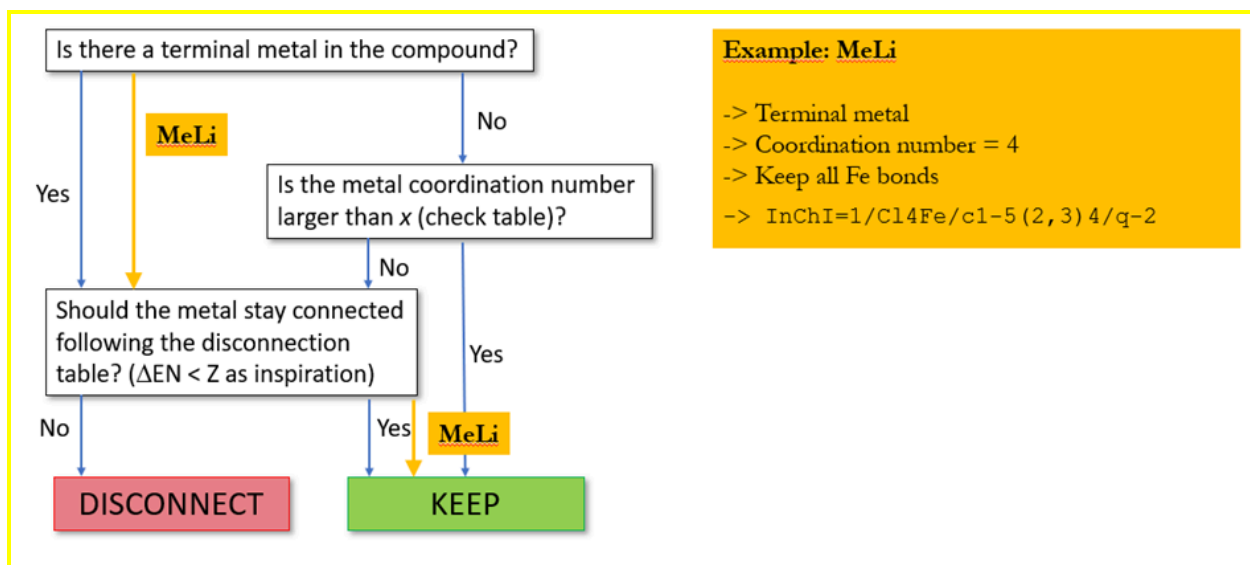


Figure S8: Bond handling for a Zintl ion synthesised by the Dehnen group[17]



Example: MeLi

- > Terminal metal
- > Coordination number = 4
- > Keep all Fe bonds
- > InChI=1/C14Fe/c1-5(2,3)4/q-2

Figure S9: Bond handling for MeLi

## References

- [1] GitHub, *Official home of the InChI*, accessed on 17/July/2024, <https://github.com/IUPAC-InChI/InChI>
- [2] Google, *OSS-Fuzz*, accessed on 17/July/2024, <https://google.github.io/oss-fuzz/>
- [3] *The Fuzzing Project*, accessed on 17/July/2024, <https://fuzzing-project.org/>
- [4] *SEI CERT C Coding Standard*, accessed on 17/July/2024, <https://wiki.sei.cmu.edu/confluence/display/c/Scope>
- [5] *C11*, accessed on 17/July/2024, <https://en.cppreference.com/w/c/11>

- [6] National Center for Biotechnology Information (NCBI), *PubChem Compound FTP site*, accessed on 17/July/2024, <https://ftp.ncbi.nlm.nih.gov/pubchem/Compound/>
- [7] National Center for Biotechnology Information (NCBI), *PubChem Compound3D FTP site*, accessed on 17/July/2024, [https://ftp.ncbi.nlm.nih.gov/pubchem/Compound\\_3D/](https://ftp.ncbi.nlm.nih.gov/pubchem/Compound_3D/)
- [8] National Center for Biotechnology Information (NCBI), *PubChem Substance FTP site*, accessed on 17/July/2024, <https://ftp.ncbi.nlm.nih.gov/pubchem/Substance/>
- [9] InChI Trust, *InChI Web Demo*, accessed on 17/July/2024, <https://iupac-inchi.github.io/InChI-Web-Demo/>
- [10] GitHub, *Ketcher*, accessed on 17/July/2024, <https://github.com/epam/ketcher>
- [11] *WebAssembly*, accessed on 17/July/2024, <https://webassembly.org/>
- [12] *Emscripten*, accessed 17/July/2024, <https://emscripten.org/>
- [13] R. L. Apodaca, *Compiling InChI to WebAssembly Part 1: Hello InChI*, accessed on 17/July/2024, <https://depth-first.com/articles/2019/05/15/compiling-inchi-to-webassembly-part-1/>
- [14] Kekule.js Lab, *InChI Demo*, accessed on 17/July/2024, <https://partridgejiang.github.io/cheminfo-to-web/demos/items/InChI/inchiDemo.html>
- [15] G. Landrum and the RDKit contributors, *RDKit.js*, accessed on 17/July/2024, <https://www.rdkitjs.com/>
- [16] GitHub, *EPAM Indigo projects*, accessed on 17/July/2024, <https://github.com/epam/Indigo>
- [17] Y. R. Lohse, B. Weinert, B. Peerless, S. Dehnen, *Z. Anorg. Allg. Chem.* 2024, 650, e202300229