

Status of the InChI development

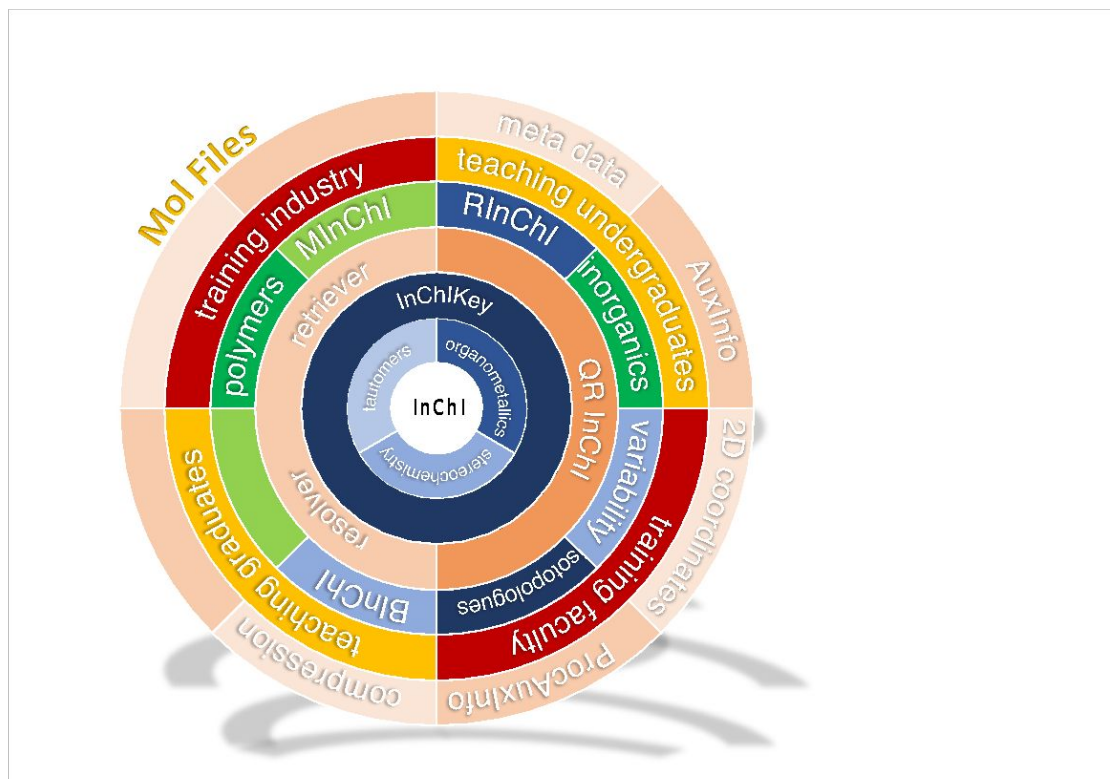
Gerd Blanke

Technical Director of the InChI Trust

Djordje Baljozovic, Felix Bänisch, Nauman
Ullah Khan, Jan Brammer, Frank Lange

24-July-2024

The International Chemical Identifier - the InChI Ecosystem



Agenda

- Status overview
 - InChI release status
 - Programming – next steps
 - Versioning
 - InChI working groups
 - RInChI
 - Organometallics
 - Stereochemistry

InChI release status

InChI release status

InChI

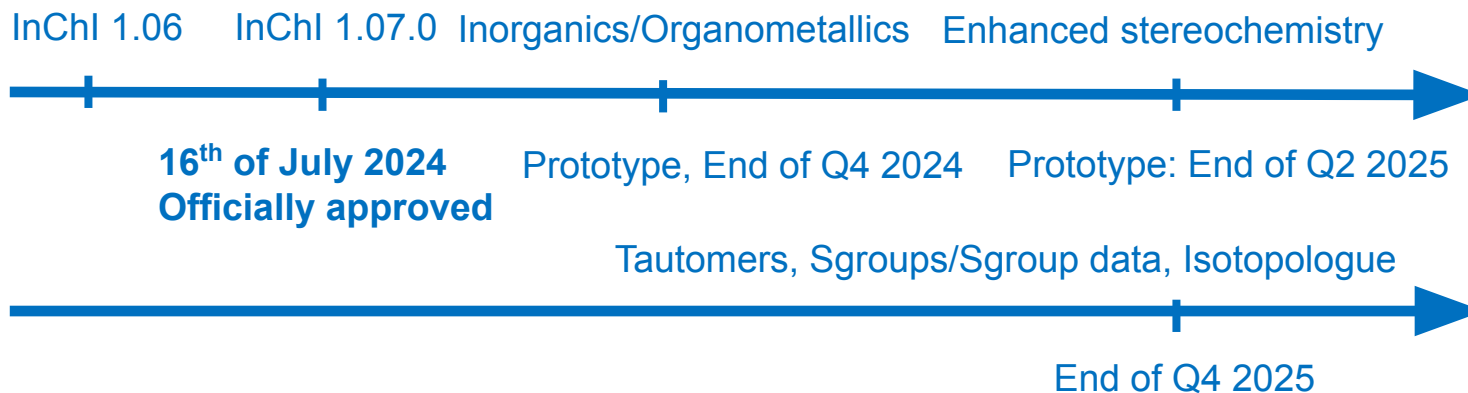
- Released **1.07.0 out!**
 - Code clean-up (about 3000 issues)
 - Polymers (in beta status)
 - 6 Tautomer transformations (for testing)
- Under development
 - Requirements (nearly) ready
 - Molecular inorganics
 - Stereochemistry
 - Tautomers
 - Sgroups and Sgroup data
 - Isotopologues
- Longer evaluations
 - Variable molecules / Markush InChIs
 - Prototypes being tested
 - <https://github.com/topics/inchi>
 - Large molecules

InChI Application Framework

- Released
 - RInChI 1.0
- Published
 - QR code
- Test phase
 - InChI Web Demo
 - <https://iupac-inchi.github.io/InChI-Web-Demo/>
 - Resolver
 - <https://github.com/inchiresolver/inchiresolver#readme>
 - RInChI (1.1)
 - <https://github.com/IUPAC-InChI/RInChI>
- Awaiting coding
 - MInChI
 - Prototype released
 - <https://github.com/cdd/mixtures>
 - RInChI (1.2)
- Evaluations
 - NanoInChI (NInChI)

InChI Roadmap

InChI



Application Framework

RInChI 1.0 , QR Code

Test: InChI Web Demo, Resolver, RInChI 1.1
Coding: MInChI, RInChI 1.2, NInChI?

Evaluations:

Large molecules
Variable molecules / Markush InChIs
NanoInChI (NInChI)

All technical developments depend on further funding.

InChI release status

- 1.07.0 released and accepted 16-July-2024
 - <https://github.com/IUPAC-InChI/InChI>
 - Availability of 6 tautomer transformation in test mode
 - Fix of a bug in InChI to structure for Buten-2
 - Fix of a buggy error message for no-structures
 - More than 3,000 issues fixed

InChI release status

- 1.07.0 released and accepted 16-July-2024
 - Google fuzz
 - Main reported issues
 - Buffer overflows
 - Memory leaks
 - Segmentation faults
 - Special local test environment in place
 - Easier and faster testing
 - Total reported number of bugs has come down
 - Further fixes in minor release versions 1.07.x
 - 1.07.1 will provide a bug fix for an issue reported by Burt Leland

InChI release status

- 1.07.0 released and accepted 16-July-2024
 - License harmonization
 - 1.07 runs under MIT license
 - License mismatch in older versions of GitHub repository
 - Issue during the GitHub set-up: The older versions under the GitHub release folder got the new MIT license.pdf. For adjustment the release folder must be re-built.

InChI release status

- Compiler versions
 - Versions 7 to 14 of GCC compilers have been used to check compatibility with older versions and the consistency of eventual warning messages.
 - 5 structures (of 300 million PubChem substances) fail but only with GCC 12, 13 and 14
 - Solved by modified setting of optimization parameters
 - Microsoft compilers
 - Open issue: MAC version
 - Up to now not supported in main InChI release code

InChI release status

- Testing
 - Test base

<https://ftp.ncbi.nlm.nih.gov/pubchem/>

	Compound	Compound 3D	Substance
download ^a	Oct 13 2023	Oct 25 2023	Oct 23 2023
size in GB (gzip) ^b	99	37	81
N SDF ^c	338	1,103	895
N structures ^d	114,726,411	23,487,296	306,711,305

- All results pertain to InChI 1.7.0 compiled with GCC 14.1.0 on Debian bookworm. All test ran on 16 physical cores.

InChI release status

- Regression tests
 - Comparison of the results of 1.07.x against 1.06

Regression Results

	Compound	Compound 3D	Substance
N structures ^e	114,726,411	23,487,296	306,711,305
N structures passed ^f	114,726,411	23,487,296	306,711,303
N structures failed ^g	0	0	2
percentage failed ^h	0	0	0.00000064
run-time total ⁱ	402 min (6 hrs, 42 min)	106 min (1 hr, 46 min)	585 min (9 hrs, 45 min)
avg run-time per structure ^j	0.21 ms	0.27 ms	0.114 ms

- Standard tests for 1.07.x versions
 - Run for each release

InChI release status

- Invariance tests
 - Renumbering of tested molfile
 - Each molfile is renumbered 10 times

Invariance Results

	Compound	Compound 3D	Substance
N structures ^k	114,726,411	23,487,296	306,711,305
N structures missing ^l	0	0	16,932,378
N structures error ^m	n/a	0	21
N structures passed ⁿ	n/a	23,487,290	289,776,775
N structures failed ^o	n/a	6	2,131
percentage failed ^p	n/a	0.000026	0.000735
run-time total ^q	n/a	389 min (6 hrs, 29 min)	4,063 min (2 days, 18 hrs, 43 min)
avg run-time per structure ^r	n/a	0.98 ms	0.84 ms

InChI release status

- Invariance tests
 - RDKit tool used instead of InChI function (easier to implement, used in other projects as well)
 - Development of unit tests continuing
 - GitHub powerful enough to run major parts of PubChem tests
 - Additional tests of InChI parameters continued
 - Because of time consumption invariance tests are only run for major releases

InChI release status

- Docker stations
 - Can be downloaded to test InChIs within own firewalls
 - Keep privateness of own data
 - Located in the test section

InChI release status

- Documentation
 - Rebuilt of documentation
 - Documentation of chemical representation into its own document (Chemical documentation)
 - Currently found in TechMan, User Guide, and FAQ
 - Additional technical InChI information about the code workflow and programming details to be added to the Technical document
 - Expected for version > 1.07.0

InChI Web Demo

- Access to 1.06 and 1.07.0 (now default!)
 - Interactive tests of actual development version

The screenshot displays the InChI Web Demo interface. At the top, the InChI TRUST logo and the title "InChI Web Demo" are visible. Below the logo, there are navigation tabs for "InChI", "RInChI", and "Funding". A toolbar contains options for "Draw structure and convert to InChI", "Convert Molfile to InChI", and "Convert InChI or AuxInfo to structure".

The main workspace shows a chemical structure of ethanol (CCO) with the labels "H₃C" and "OH". To the right of the workspace is a vertical menu of element symbols: H, C, N, O, S, P, F, Cl, Br, I, PT, [A], and ET. Below this menu is a "Select InChI version" dropdown menu, which is circled in red and currently displays "1.07.0".

On the right side of the interface, there are several configuration options:

- 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
- Tautomer options
- Treat polymers:
 - No pre-edits of original polymer structure
 - Enable CRU folding
 - Disable CRU frame shift
- Allow non-polymer Zz pseudoatoms
- [Reset InChI Options](#)

Below the configuration options, the "Select InChI version" dropdown is set to "1.07.0". At the bottom of the interface, there are several text boxes containing the following information:

- InChI: InChI=1S/C2H6O/c1-2-3/h3H,2H2,1H3
- InChIKey: Lfqscwfljhtthz-uhfffaoySA-N
- AuxInfo: AuxInfo=1/0/N:1,2,3/rA:3nCCO/rB:s1;s2;/rC:2.842,-4.25,0;3.708,-3.75,0;4.574,-4.25,0;
- Log
- InChI options:

InChI Web Demo

- Work on additional functionality
 - Display of InChI numbering in molecule
 - Display of tautomeric hydrogen zones within the molecule

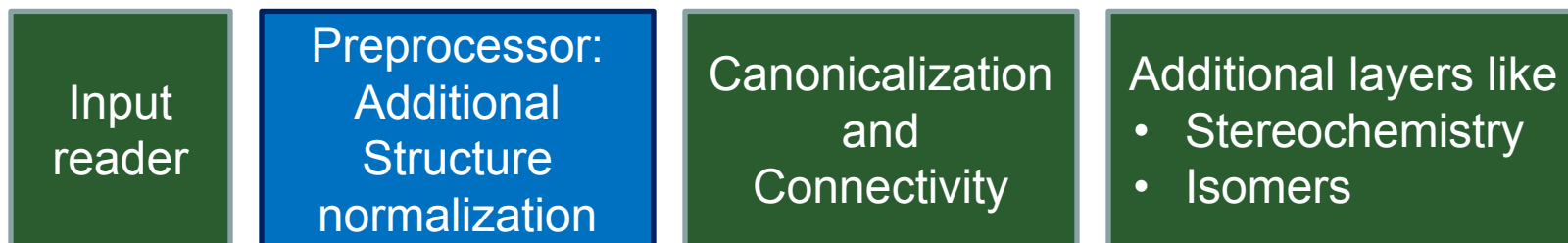
WInChI

- Windows based InChI calculator
 - Maintenance and further development will be done by the original author Dmitrii Tchekhovskoi
 - Will be added to GitHub in it sown repository

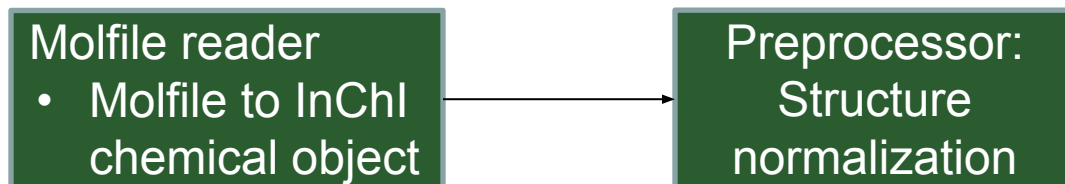
Programming – next steps

Programming – next steps

- New Architecture



- InChI internal reader: molfile to internal chemical object



QA of the input reading process

- Check of the V2 and V3 reader
 - Check of the types of values that are transferred from the molfile to the internal chemical object
 - V2 molfile reader is the most comprehensive one
 - Only expected properties fail like structures with query atoms or bonds (except aromatic bonds)
 - V3 molfiles show issues with some of the Sgroup types
 - Some of the polymer brackets are not supported
 - 3D molfile are under investigation

Programming – next steps

- Structure normalization functionality
 - Break bonds
 - Create bonds
 - Exchange bonds
 - Bonds to represent inorganics and organometal-lics are replaced by single bonds with adapted valences at the atoms
 - Coordinative bonds
 - Dative bonds
 - Haptic bonds

InChI versioning

- Discussion about InChI versioning
 - We must change to 2.x if the format is going to change
 - Modification of the InChI keys
 - Longer second section to reflect stereochemistry issues
 - Major changes in the canonicalization and connectivity string creation
 - Restricted backward compatibility
 - Current view: Keep 1.x as long as possible

InChI versioning

- Discussion about InChI versioning
 - Recognition of version 1.1 or higher
 - InChIKeys only represented by “SA”
 - “A” only points to any 1 version
 - Transfer from 1.0x to any 1.1 or higher version
 - What will be needed from the technical side will be determined during the implementation of organometallics and enhanced stereochemistry

InChI working groups

- Reaction InChI (RInChI)

RInChI

- RInChI 1.1 is prepared for release
 - Based on InChI 1.07.0
 - Additional functionality
 - Support of extended RXN format
 - Agents in third layer of RXN format
 - Fix of known technical issues
 - RInChI Mac version contributed by István Öri, member of the RInChI group
 - <https://github.com/IUPAC-InChI/RInChI>
 - RInChI became prototype for the InChI development on GitHub

RInChI

- RInChI 1.1 is prepared for release
 - Tests
 - USPTO reaction set 2008 to 2011
 - About 400,000 reactions
 - Discrepancies seen are based on differences between InChI 1.04 (used in RInChI 1.0) and InChI 1.07
 - InChI Web Demo can be used for RInChI tests
 - Developing a new RInChI tutorial (Günter Grethe)
 - To be integrated with the InChI Web Tool (if possible)
 - Development of simple search engine for (R)InChI

RInChI

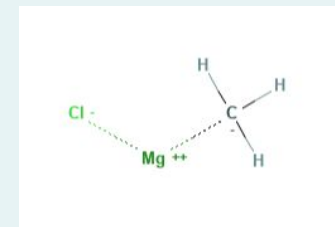
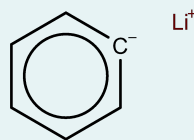
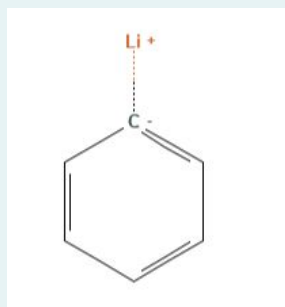
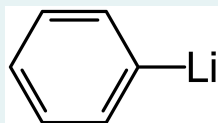
- Enhancements planned for RInChI 1.2
 - Atom mapping
 - Multithreading

InChI working groups

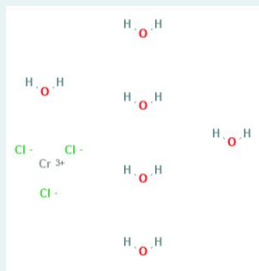
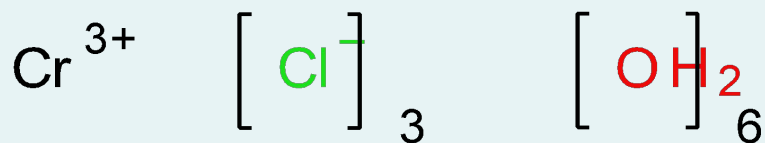
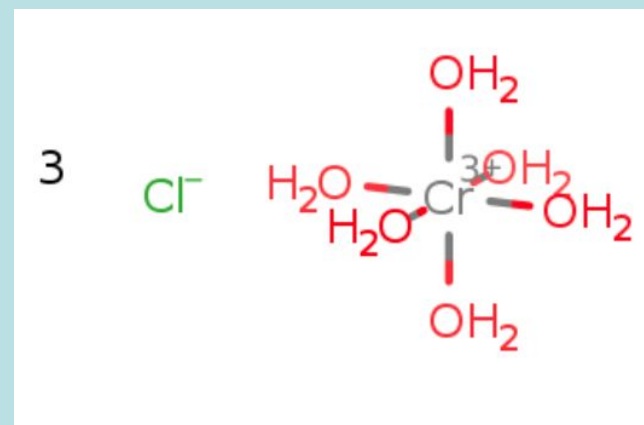
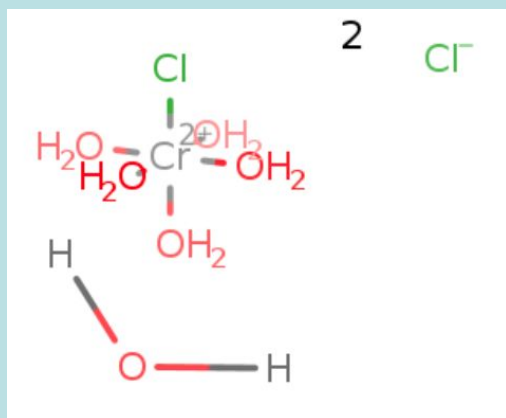
- Inorganics / Organometallics

Inorganics and organometallics

- Working group set up the general implementation guidelines



Inorganics and organometallics

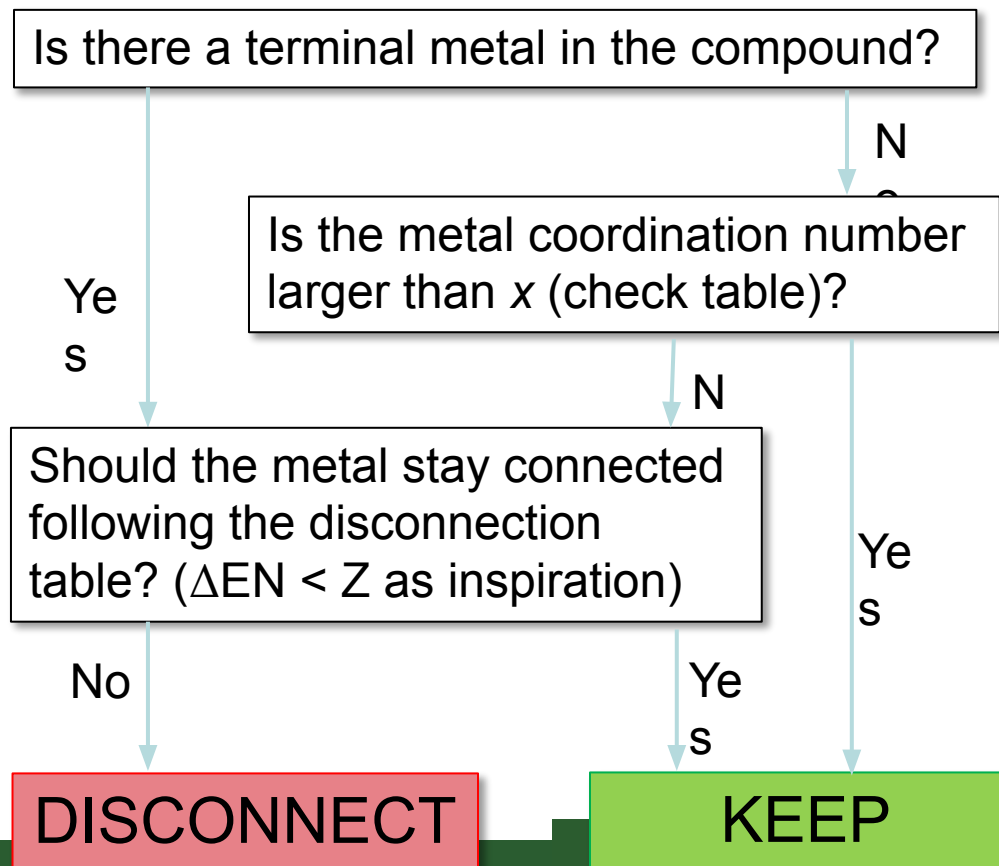


Inorganics and organometallics

- Keep connectivity for inorganics and organometallics
 - Keep explicitly drawn Hydrogens
 - Hydrogen atoms must be drawn explicitly
 - Exception: H atoms bound to atoms of organic elements like C, N, O where the H count is determined by the atom valence and the actual bond order of the atom
 - Exception of connectivity: see simple flow chart

Inorganics and organometallics

• Disconnection flow chart



To be noted:

- This flow chart has to be run for every metal.
- The „Disconnection Table“ step has to be run for every bond to a metal.
- If atoms are not connected in the input, they stay disconnected. No new bonds are formed that are not already in the input.

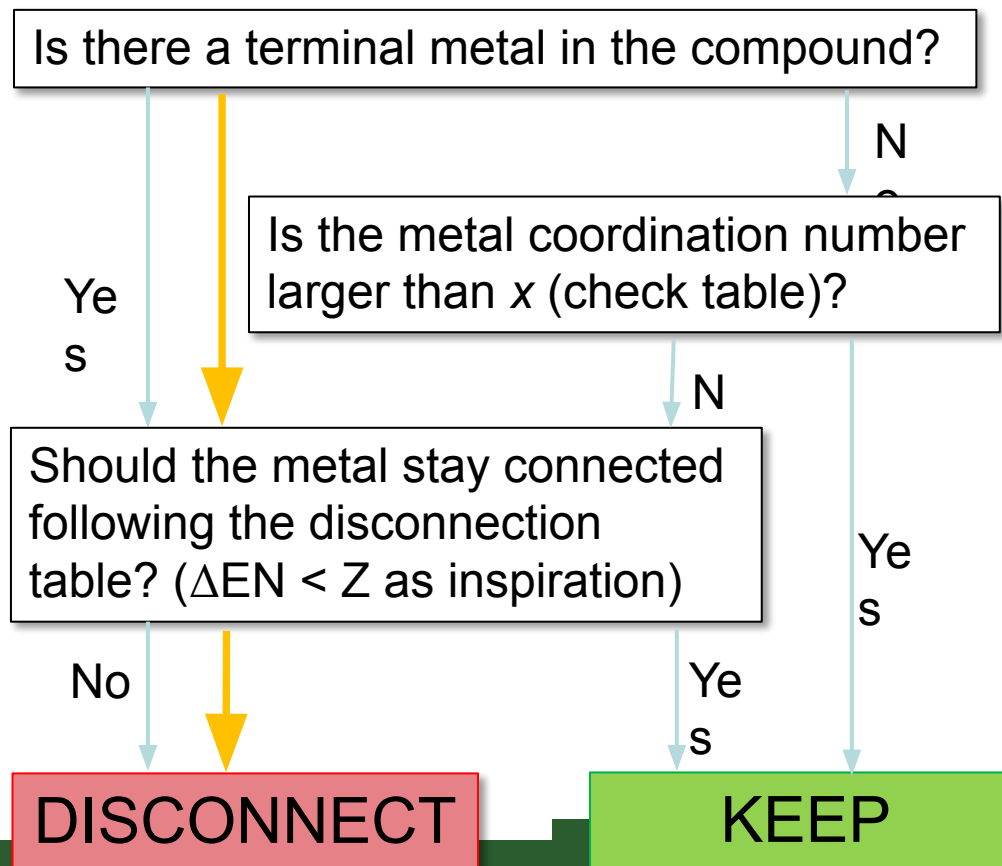
Inorganics and organometallics

- Disconnection table

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	
1	Elektronenaktivitätsdifferenzen			¹ H	² H / D	³ H / T	⁴ He	⁶ Li	⁹ Be	¹⁰ B	¹² C	¹⁴ N	¹⁶ O	¹⁸ F	²⁰ Ne	²² Na	²⁴ Mg	²⁷ Al	²⁸ Si	³² S	³⁵ Cl	³⁶ Ar	³⁹ K	⁴⁰ Ca	⁴⁵ Sc	⁴⁸ Ti	⁵¹ V	⁵² Cr		
2	Deutscher Name nach IUPAC			Wasserstoff			Helium	Lithium	Beryllium	Bor	Kohlenstoff	Stickstoff	Sauerstoff	Fluor	Neon	Natrium	Magnesium	Aluminium	Silicium	Phosphor	Schwefel	Chlor	Argon	Kalium	Calcium	Scandium	Titan	Vanadium	Chrom	
3	Gewichtliche Molmasse			1.01 u			4.00 u	6.94 u	9.01 u	10.81 u	12.01 u	14.01 u	16.00 u	19.00 u	20.18 u	22.99 u	24.31 u	26.98 u	28.09 u	30.97 u	32.07 u	35.45 u	39.95 u	39.10 u	40.08 u	44.96 u	47.87 u	50.94 u	52.00 u	
4	Elektronenaktivität nach Pauling			2.20			2.20	0.98	1.57	2.04	2.55	3.04	3.44	3.98	0.33	1.31	1.61	1.90	2.19	2.59	3.16	3.16	0.82	1.00	1.00	1.36	1.54	1.63	1.86	
5	Atomradius			73 pm			43 pm	205 pm	140 pm	117 pm	31 pm	75 pm	65 pm	57 pm	51 pm	223 pm	172 pm	182 pm	146 pm	123 pm	109 pm	97 pm	89 pm	277 pm	223 pm	209 pm	200 pm	192 pm	185 pm	
6	Kovalenzradius			32 pm			33 pm	123 pm	90 pm	82 pm	77 pm	75 pm	73 pm	72 pm	71 pm	154 pm	136 pm	118 pm	111 pm	106 pm	102 pm	99 pm	98 pm	203 pm	174 pm	144 pm	132 pm	122 pm	118 pm	
7	Ausgewählte Oxidationszahlen (nichtgatte fett)			1			1 ⁺	1 ⁺	2 ⁺	3 ⁺	-4, ..., -4	-3, ..., -5	-2, -1	-1	-	1	2	3	4	3, 5, 4	±2, 4, 6	±1, 3, 5, 7	-	1	2	3	4, 3	5, 4, 3, 2	6, 3, 2	
8	Elektronenkonfiguration im Grundzustand			1s ¹			1s ²	1s ² 2s ¹	1s ² 2s ²	1s ² 2s ² 2p ¹	1s ² 2s ² 2p ²	1s ² 2s ² 2p ³	1s ² 2s ² 2p ⁴	1s ² 2s ² 2p ⁵	1s ² 2s ² 2p ⁶	[Ne] 3s ¹	[Ne] 3s ²	[Ne] 3s ² 3p ¹	[Ne] 3s ² 3p ²	[Ne] 3s ² 3p ³	[Ne] 3s ² 3p ⁴	[Ne] 3s ² 3p ⁵	[Ne] 3s ² 3p ⁶	[Ar] 4s ¹	[Ar] 4s ²	[Ar] 4s ² 3d ¹	[Ar] 4s ² 3d ²	[Ar] 4s ² 3d ³	[Ar] 4s ² 3d ⁴	[Ar] 4s ² 3d ⁵
9	Schmelztemperatur			14.025 K			1 K (126 atm)	453.7 K	1758 K	2300 K	4700 K	63 K	50.35 K	93.48 K	24.553 K	371.0 K	922 K	933.25 K	1685 K	371.30 K	388.36 K	172.16 K	336.35 K	383.91 K	336.35 K	1122 K	1943 K	2175 K	2700 K	
10	Siedetemperatur			20.288 K			4.215 K	1615 K	2745 K	4275 K	4470 K	171.35 K	90.16 K	84.36 K	27.086 K	176 K	1763 K	2733 K	3540 K	550 K	717.75 K	239.1 K	87.30 K	1032 K	1757 K	3704 K	3562 K	3762 K	2945 K	
11	Dichte bei Raumtemperatur			0.0899 g/L			0.1787 g/L	0.53 g/cm ³	1.85 g/cm ³	2.34 g/cm ³	2.62 g/cm ³	1.251 g/L	1.429 g/L	1.696 g/L	0.901 g/L	0.97 g/cm ³	1.74 g/cm ³	2.70 g/cm ³	2.33 g/cm ³	1.82 g/cm ³	2.07 g/cm ³	3.17 g/L	1.784 g/L	0.86 g/cm ³	1.55 g/cm ³	3.0 g/cm ³	4.50 g/cm ³	5.8 g/cm ³	7.19 g/cm ³	
12	Element			H	D	T	He	Li	Be	B	C	N	O	F	Ne	Na	Mg	Al	Si	P	S	Cl	Ar	K	Ca	Sc	Ti	V	Cr	
13	IUPAC			1	2	3	4	7	9	11	12	14	16	18	20	23	24	27	28	31	32	35	40	45	48	51	52	51	52	
14	AMSL			1	2	3	4	7	9	11	12	14	16	18	20	23	24	27	28	31	32	35	40	45	48	51	52	51	52	
15	CAS			1007825	2014102	3016049	400026	7016	901219	110003	120007	1509491	180994	1909244	2200977	2300977	2300977	2609854	2709853	30097376	31097207	3409635	390624	390624	4409531	4704735	5004396	5104005	5104005	
16	Typ			non-metals (H inert)			METAL	METAL	0	0	0	0	0	0	0	METAL	METAL	METAL	0	0	0	0	METAL	METAL	METAL	METAL	METAL	METAL	METAL	
17	Elektronenaktivität (Pauling)			2.1	2.1	2.1	0	10	15	20	25	30	35	40	0	5	12	15	18	21	25	30	0	6	10	15	15	16	16	
18	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
19	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
20	CAS			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
21	Typ			non-metals (H inert)			METAL	METAL	0	0	0	0	0	0	0	METAL	METAL	METAL	0	0	0	0	METAL	METAL	METAL	METAL	METAL	METAL	METAL	
22	Elektronenaktivität (Pauling)			2.1	2.1	2.1	0	10	15	20	25	30	35	40	0	5	12	15	18	21	25	30	0	6	10	15	15	16	16	
23	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
24	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
25	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
26	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
27	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
28	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
29	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
30	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
31	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
32	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
33	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
34	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
35	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
36	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
37	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
38	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
39	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
40	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
41	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
42	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
43	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
44	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
45	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
46	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
47	IUPAC			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	
48	AMSL			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	

Inorganics and organometallics

• Disconnection flow chart



Example: Na-Cl

-> Terminal metal

-> Disconnection table says DISCONNECT

-> Na⁺ and Cl⁻

->

InChI=1S/ClH.Na/h1H;/q;+1/p-1-

Inorganics and organometallics

- Any bond is not allowed (like in organic InChI)
- Zero order bond / molfile bond type 9 / coordination bond
 - Is allowed
 - Represents connection only
 - Does not affect (alter) H-Count
 - Does not affect (alter) valence count (ligand count) ?
 - In the context of InChI calculation it is converted to single bond via setting appropriate valence count on starting atom / ending atom

Inorganics and organometallics

- Output format
 - New standard InChI format for inorganics and organometallics is based on the current reconnection layer
InChI=1S/C8H16Cu2O10/c1-5-13-9(11)15-6(2)16-10(12,14-5,18-7(3)17-9)20-8(4)19-9/h11-12H2,1-4H3
 - For consistency add a new parameter *disconnectMetal* to let the InChI calculation add the current format as additional layer

Inorganics and organometallics

- InChI (and InChIKey) do not change for organic molecules!
- Currently working out further implementation details
 - Salt handling versus metal disconnection
 - Recognition of polyhedrons
 - Stereochemistry

InChI working groups

- Stereochemistry

Stereochemistry

- Recognition and encoding of atropisomers
- Support of enhanced stereo designation
- Correct recognition of carbohydrates in Haworth and chair forms
- Recognition and encoding of configuration of long allenes
- Recognition of configuration for special 'spiro centers'

Stereochemistry

- Specific designations of configurations valid and invalid representations
- Polyhedral configurations in collaborations with inorganic/organometallics group
- Any other useful proposals related to InChI stereo

Stereochemistry

- Preparation of recommendations on enhanced stereo designation in collaboration with Division VIII.
 - Not directly connected to InChI development, but it is important for making enhanced stereo recognized by IUPAC.

Thanks