

International Chemical Identifier for Reactions (RInChI)

The key to effectively managing reaction databases

Guenter Grethe¹, Gerd Blanke², Jonathan M. Goodman³, Hans Kraut⁴

Background

Since its inception, the IUPAC International Chemical Identifier (InChI) has found wide acceptance as a standard in the chemical community. In order to widen the applicability of the identifier, the IUPAC Division VIII Subcommittee and the InChI Trust have initiated several projects to extend the usage of the identifier. Among these is the development of a non-proprietary, international identifier for reactions (RInChI) to describe chemical reactions in a unique, machine-readable, character string based on the InChI algorithm suitable for data storage and indexing. Prototype versions of the RInChI supported by the IUPAC and the University of Cambridge have been available since 2011. The first official release (RInChI-V1.00), funded by the InChI Trust is now available for download.

Reactions are more difficult to codify than molecules. A molecule can usually be uniquely and precisely defined by a structural drawing, which shows the atom types, the bond connectivity, and the stereochemistry. A reaction includes molecules, but also times, temperatures, concentrations, rates of mixing, yields, and numerous other quantities that are recorded with varying levels of precision.

The goal of the development of RInChI is to define a reaction identifier that has enough detail to be useful and yet not so much that practically identical processes would be given different labels as described in the following diagram:

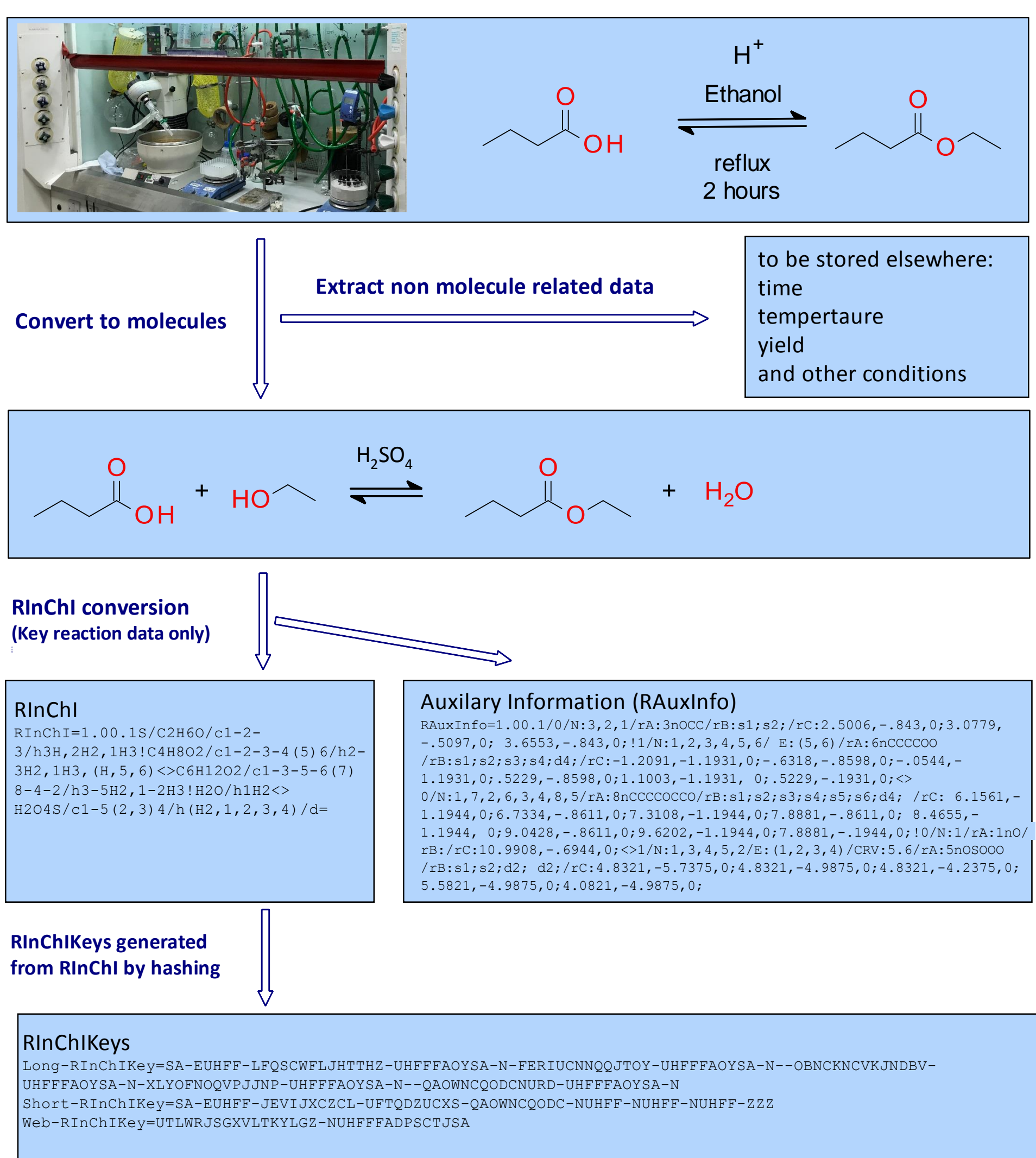
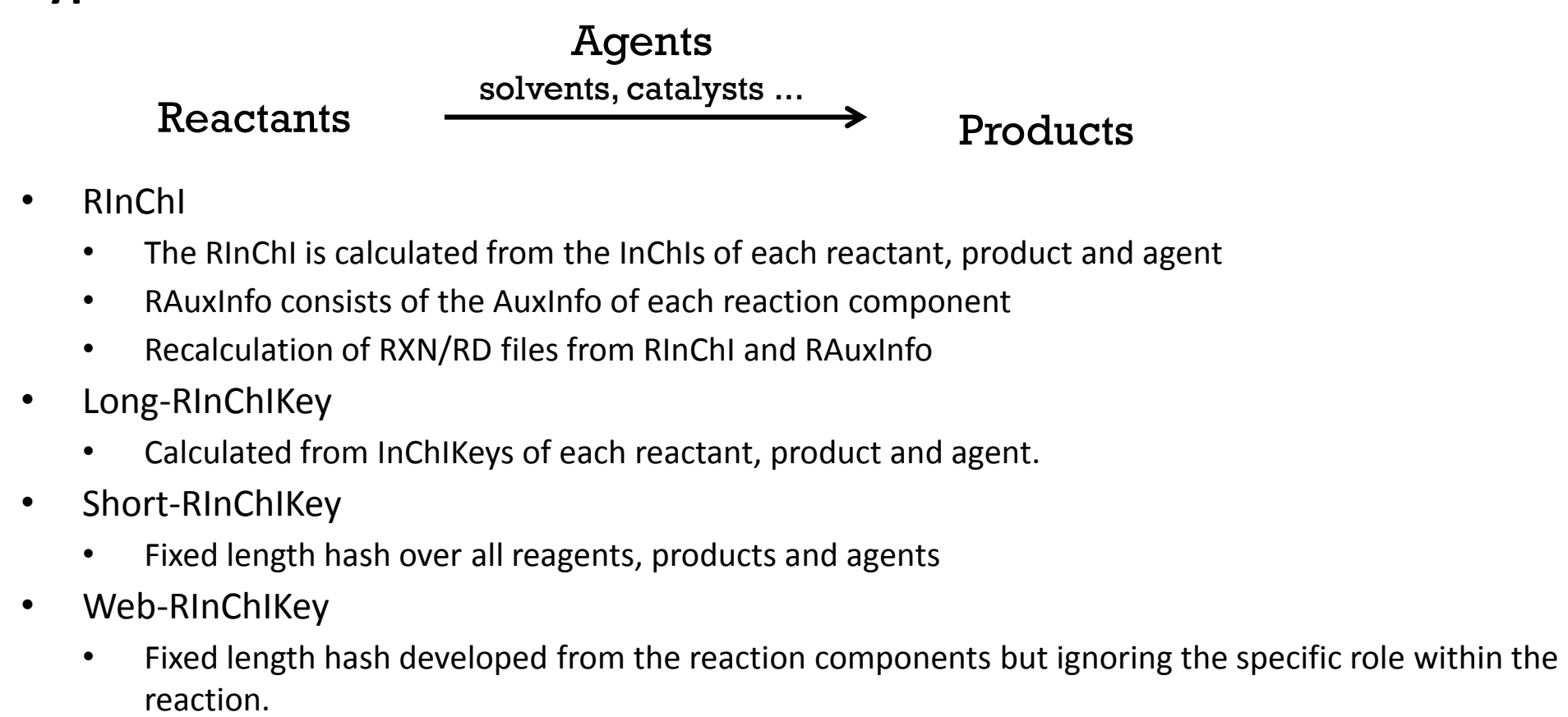


Fig 1: RInChI creation diagram

RInChI as new type of identifier is extremely useful in the creation and analysis of reaction databases as well as in Web searches where reactions must be found in searches over multiple databases with different data models.

While identifiers for molecule exist for a long time (e.g. CAS-Number, MFCDNumber, InChI) there is no public identifier for reactions available up to now. In fact, the Web-RInChIKey introduced here is the first reaction identifier that is suitable for the special conditions of Web searches.

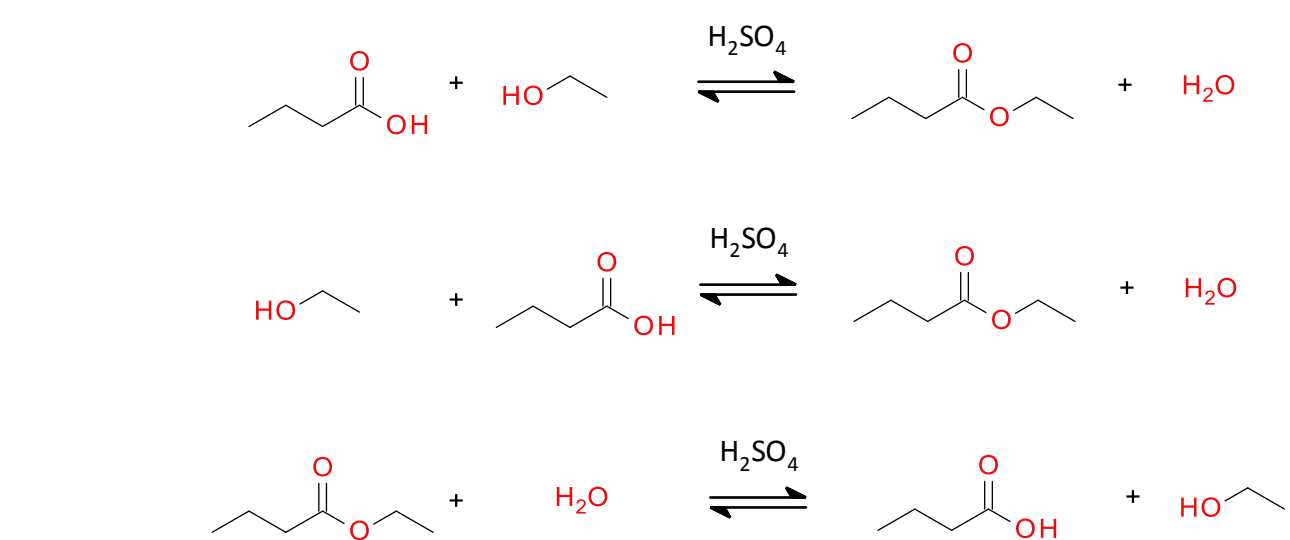
Types of RInChIs



Format of RInChI

- RInChI=1.00.1S/layer2<layer3>layer4/d(+,-,=)/u#2-#3-#4
- Layer 1 identifying the RInChI version 1.00 using Standard InChI version 1
 - Layer 2 / 3: InChIs of starting materials and products
 - Layer 4: InChIs of all agents (catalysts, solvents, etc.)
 - Layer 5: Directional identifier with $d=$ forward, $d=-$ backward and $d=$ equilibrium reaction
 - Layer 6: No-structure flag layer identified by /u with #2 as number of no-structures in layer 2, #3 number of no-structures in layer 3 and #4 number of no-structures in layer 4
 - Additional rules for building the RInChI (string)
 - The "InChI=1S/" is omitted from InChIs listed in the RInChI as it is already specified in layer 1.
 - Separators
 - Use "." as divider between main layers
 - Use "<" as divider between the reactant, product and agent layers
 - Use "/" as divider between molecules
 - Use "-" as divider for the occurrence number of no-structures
 - Separators are omitted if they have nothing to separate
 - Make RInChI representation unique by alphabetical ordering
 - Order InChIs within each layer alphabetically
 - Order layer 2 and 3 alphabetically.
 - In case that layer 2 and 3 are re-ordered the directional identifiers $d=$, $d=-$ must be exchanged to keep the reaction direction, i.e. $d=$ or $d=-$ or vice versa $d=->d=$ (i.e. implies an equilibrium, and we intend to introduce $d=$ for reactions which do not work)
 - The same molecule may not appear in both layer two and layer three, and no molecule should be repeated within each layer. Anything that is in both layers should be removed from both and added to layer four.
 - Empty layers are permitted.

Example Esterification (Reaction 1)

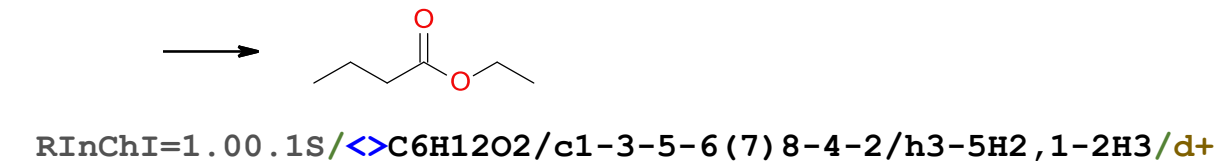


RInChI=1.00.1S/C2H6O/c1-2-3/h3h,2H,1H3!C4H8O2/c1-2-3-4(5)6/h2-3H2,1H3,(H,5,6)<<C6H12O2/c1-3-5-6(7)8-4-2/h3-5H2,1-2H3!H2O/h1H2>>H2O4S/c1-5(2,3)4/h(H2,1,2,3,4)/d=

Special cases: "half" reactions and no structures

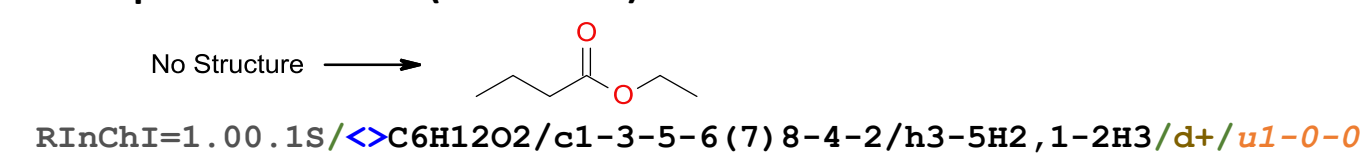
Half reactions where only the starting materials or the products are described are represented using an empty layer for the unknown part of the reaction; empty strings are not added to the RInChI.

Example half reaction (Reaction 2)



The RInChI for no-structures is an empty string as well. To distinguish no-structures in reactions from half reactions, no-structures are ignored within their group but are added as an additional flag at the right end of the RInChI string using the format u#1-#2-#A with # as integer representing the total number of No-Structures in group1, group2 and the agent group as 3rd group

Example no-structure (Reaction 3)



RAuxInfo

The RAuxInfo contains the Auxinfos of the InChIs participating in the reaction, i.e. it contains the atom coordinates, the information the original connectivity including stereochemistry, tautomeric characteristics and other molecule properties that are not directly handled by InChIs.

The RAuxInfo is necessary to re-build the original RXN or RD file. Without RAuxInfo the standard InChI routines return molfiles with all coordinates set to 0, so that the molecule coordinates must be recalculated. The recalculation may end in other atom locations of the same molecule. Although the stereocenters are kept by InChI the re-calculation may lead into another distribution of the bound ligands around the stereocenter so that the stereo parity may flip.

Example Esterification (Reaction 1)

RAuxInfo=1.00.1/0/N:3,2,1/rA:3n00C/rB:s1;s2;/rC:2.5006,-.843,0;3.0779,-.5097,0;3.6553,-.843,0;11/N:1,2,3,4,5,6/E:(5,6)/rA:6n0000C/rB:s1;s2;s3;s4;d4;/rC:-1.2091,-1.1931,0,-.6318,-.8598,0,-.0544,-1.1931,0;-5.229,-.8598,0;1.1003,-1.1931,0;-5.229,-.1931,0;0/N:1,7,2,6,3,4,8,5/rA:8n0000C/rB:s1;s2;s3;s4;s5;s6;d4;/rC:6.1561,-1.1944,0;6.7334,-.8611,0;7.3108,-1.1944,0;7.8881,-.8611,0;8.4655,-1.1944,0;9.0428,-.8611,0;9.6202,-1.1944,0;7.8881,-1.1944,0;10/N:1/rA:1n0/rB:/rC:10.9908,-.6944,0;<1/N:1,3,4,5,2/E:(1,2,3,4)/CRV:5.6/rA:5n0000C/rB:s1;s2;d2;/rC:4.8321,-5.7375,0;4.8321,-4.9875,0;4.8321,-4.2375,0;5.8212,-4.9875,0;4.8321,-4.2375,0

Example Half Reaction (Reaction 2)

RAuxInfo=1.00.1/<0/N:1,7,2,6,3,4,8,5/rA:8n0000C/rB:s1;s2;s3;s4;s5;s6;d4;/rC:5.7129,-7.3792,0;6.2903,-7.0459,0;6.8676,-7.3792,0;7.445,-7.0459,0;8.0223,-7.3792,0;8.5997,-7.0459,0;9.1777,-7.3792,0;7.445,-6.3792,0

Example No-Structure (Reaction 3)

RAuxInfo=1.00.1/<0/N:1,7,2,6,3,4,8,5/rA:8n0000C/rB:s1;s2;s3;s4;s5;s6;d4;/rC:6.8243,-5.2486,0;7.3225,-4.961,0;7.8165,-5.2486,0;8.3105,-4.961,0;8.8046,-5.2486,0;9.3027,-4.961,0;9.7968,-5.2486,0;8.3105,-4.3903,0

Reaction InChIKeys (RInChIKey)

Like the InChIKey the RInChIKeys are created by hashing the InChIs ordered in RInChI format. Depending on the purpose three different types of RInChIKeys have been developed: Long-RInChIKeys, Short-RInChIKeys and Web-RInChIKeys

Long-RInChIKeys

Long-RInChIKeys are created by building the InChIKeys for each of the components participating in the reaction based on the order of RInChI

Long-RInChIKey=SA-(F,B,E,U)UHFF-layer3-layer4-layer5

- Layer 1 identifying the InChI version SA as Standard InChI 1 (=A)
- Layer 2:
 - The first letter describes the direction of the reaction with F(forward) = 'd', B(ackward) = 'd-', E(quilibrium) = 'd=' and U(nspecified)
 - The remaining 4 characters "UHFF" are currently not used and will handle reaction conditions in a later version of RInChI
- Layer 3 to 5: InChIKeys related to the InChIs in the RInChI layers 2 to 4 using the order of InChIs in the corresponding RInChI layers
 - Multiple InChIKeys in the same layer are separated by a single hyphen "-"
 - Layers 3, 4, and 5 are separated by a double hash "-"
 - Each No-structure or one of the related pseudotoms (R, X, A and *-atom) is represented by the hash value for an empty string ("MOSFUXAXDLML-UHFFFAOYSA-N")
 - Layers are only displayed if they contain at least one InChIKey. If the last layer of the string is empty the separator "-" is omitted as well.

Long-RInChIKeys are a valuable tool for the database storage of reactions. Beside uniqueness checks, they allow the identification of each reaction component by simple text searches based on Standard InChIKeys.

Example Esterification (Reaction 1)

Long-RInChIKey=SA-EUHFF-LFQSCWFLJHTTSE-UHFFFAOYSA-N-FERIUCNNQJTOY-UHFFFAOYSA-N--OBNCNKVKJNDVB-UHFFFAOYSA-N-XLYOQVJUNP-UHFFFAOYSA-N--QAOWNCQODCNURD-UHFFFAOYSA-N

Example Half Reaction (Reaction 2)

Long-RInChIKey=SA-FUHFF---OBNCNKVKJNDVB-UHFFFAOYSA-N

Example No-Structure (Reaction 3)

Long-RInChIKey=SA-FUHFF-MOSFUXAXDLML-UHFFFAOYSA-N--OBNCNKVKJNDVB-UHFFFAOYSA-N

Because any molecule within a reaction can be identified by a text search by its InChIKey the Long-RInChIKeys are a good tool to provide reaction component searches and can be used to build synthesis trees where e.g. a starting material is identified being a product within another reaction.

On the other hand the size of Long-RInChIKeys depends on the number of reaction components and is not fixed. Therefore, it is not the preferred tool for database actions like reaction indexing.

Short-RInChIKeys

Short-RInChIKeys have a fixed length of 55 letters plus 8 hyphens as separators resulting in a total of 63 characters. They are built by hashing each of the layers of the RInChI to a fixed length using the following format

- Short-RInChIKey=SA-(F,B,E,U)UHFF - hash over all major layers in RInChI 2 - hash over all major layers in RInChI 3 - hash over all major layers in RInChI 4 - hash over all minor layers plus sum of protonation states in RInChI 2 - hash over all minor layers plus sum of protonation states in RInChI 3 - hash over all minor layers plus sum of protonation states in RInChI 4 - ### with # = Z, A, B...
- Layer 1 and 2 correspond to layer 1 and 2 of the Long-RInChIKey (see above)
 - Layer 3, 4, and 5 are hashes over all major layers of the InChIs from the layers 2, 3, and 4 of the RInChI keeping the first 10 letters.
 - Layer 6, 7, and 8 are hashes over all minor layers of the InChIs from the layers 2, 3, 4 of the RInChI keeping the first 4 letters plus one letter representing the sum over all protonation states of the related InChIs
 - Layer 9 consists of three letters, each of them representing the count of no-structures found in the sixth layer of RInChI. To simplify reading, "Z" corresponds to a count of zero no-structures (i.e. no no-structure), while "A" stands for one no-structure, B for two and so forth. The format is given by ### with # = Z, A, B
 - All layers are separated by an hyphen "-"

Example Esterification (Reaction 1)

Short-RInChIKey=SA-EUHFF-JEVIJXCZCL-UFTQDZUCXS-QAOWNCQODC-UHFFFAOYSA-N--OBNCNKVKJNDVB-UHFFFAOYSA-N

Example Half Reaction (Reaction 2)

Short-RInChIKey=SA-FUHFF-UHFFFAOYSA-N--OBNCNKVKJNDVB-UHFFFAOYSA-N

Example No-Structure (Reaction 3)

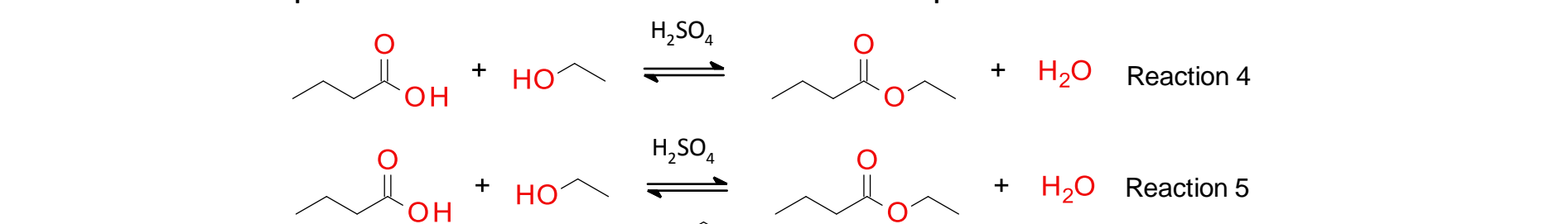
Short-RInChIKey=SA-FUHFF-UHFFFAOYSA-N--OBNCNKVKJNDVB-UHFFFAOYSA-N

The fixed length of the Short-RInChIKey makes it especially suitable for

- exact searches of reactions in databases (and in the WEB)
- indexing reactions in databases
- Identifying identical reactions in different databases.

Web-RInChIKeys

The depiction of a chemical reaction is not uniquely defined. The following 3 depictions are alternative descriptions of the esterification used in example 1



Special acknowledgements are due to the RInChI Working Group for their contributions. We are grateful to Alan McNaught and Steve Heller from the InChI Trust for initializing and supporting the project. Financial support from the IUPAC Division VIII Subcommittee and the Royal Society of Chemistry is very much appreciated. Chad Allen, Matthew Morton, James Davies, Rudi Pisa, Ben Hammond, Nicholas Parker, James Athorp, Bryn Reinstadler and Duncan Hampshire are thanked for their contributions to the program. We are also grateful to NextMove Software, FIZ Chemie (Wiley) Berlin, InfoChem and Elsevier for providing trial datasets to test the program. Working group: Colin Batchelor (Royal Society of Chemistry), Gerd Blanke (StructurePendium), Jonathan Goodman (University of Cambridge), Guenter Grethe, Jan Holst Jensen (BioChemFusion), Hans Kraut (InfoChem), Alexander Lawson, Henry Matuszczyk (InfoChem), Martin Schmidt (FIZ Chemie), Keith Taylor (Lederac Consulting). Authors: ¹ Dr. Guenter Grethe, Poway, CA 92064, US; ² StructurePendium Technologies GmbH, Essen, Germany; ³ University of Cambridge, Department of Chemistry, Cambridge, UK; ⁴ InfoChem Gesellschaft für chemische Information mbH, Munich, Germany. Funding: InChI Trust, IUPAC, Royal Society of Chemistry. Downloads: The Reaction InChI (RInChI) is available on the InChI Trust website: <http://www.inchi-trust.org/> Accessed May 2018