

# InChI Version 1.06: Now More Than 99.99 % Reliable

**Evan Bolton**

National Center for Biotechnology Information <https://orcid.org/0000-0002-5959-6190>

**Jonathan M Goodman** (✉ [jmg11@cam.ac.uk](mailto:jmg11@cam.ac.uk))

University of Cambridge <https://orcid.org/0000-0002-8693-9136>

**Stephen R Heller**

InChI Trust <https://orcid.org/0000-0002-5538-8482>

**Igor Pletnev**

Lomonosov Moscow State University: Moskovskij gosudarstvennyj universitet imeni M V Lomonosova

<https://orcid.org/0000-0001-7175-2136>

**Paul Thiessen**

National Center for Biotechnology Information <https://orcid.org/0000-0002-1992-2086>

---

## Research article

**Keywords:** InChI, InChIKey, PubChem, RInChI

**Posted Date:** April 3rd, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-378701/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **InChI version 1.06: Now more than 99.99 % reliable**

2

3 Evan Bolton<sup>1</sup>, Jonathan M. Goodman<sup>2\*</sup>, Stephen R. Heller<sup>1,3</sup>, Igor Pletnev<sup>3,4</sup> and Paul

4 Thiessen<sup>1</sup>

5

6 1. National Center for Biotechnology Information, National Library of Medicine, National

7 Institutes of Health, Bethesda, MD 20894, USA.

8 2. Centre for Molecular Informatics, Yusuf Hamied Department of Chemistry, Lensfield

9 Road, Cambridge, CB2 1EW, UK

10 3. InChI Trust

11 4. Department of Chemistry, Lomonosov Moscow State University, 119991, Moscow, Russia

12

13 Corresponding author: Jonathan M Goodman [jmg11@cam.ac.uk](mailto:jmg11@cam.ac.uk)

14

15 ORCID IDs

16 Evan Bolton 0000-0002-5959-6190

17 Jonathan M. Goodman 0000-0002-8693-9136

18 Stephen R. Heller 0000-0002-5538-8482

19 Igor Pletnev 0000-0001-7175-2136

20 Paul Thiessen 0000-0002-1992-2086

21

22

23

24

1

## 2 **Abstract**

3 The software for the IUPAC Chemical Identifier, InChI, is extraordinarily reliable. It has been  
4 tested on large databases around the world, and has proved itself to be an essential tool in  
5 the handling and integration of large chemical databases. InChI version 1.05 was released in  
6 January 2017 and v. 1.06 in December 2020. In this paper, we report on the current state of  
7 InChI Software, the details of the improvements in the v.1.06 release , and the results of a  
8 test of the InChI run on PubChem, a database of more than a hundred million molecules.  
9 The upgrade introduces significant new features, including support for pseudo-element  
10 atoms and an improved description of polymers. We expect that few, if any, applications  
11 using the standard InChI will need to change as a result of the changes in version 1.06.  
12 Numerical instability was discovered for 0.002 % of this database, and a small number of  
13 other molecules were discovered for which the algorithm did not run smoothly. On the basis  
14 of PubChem data, we can demonstrate that InChI version 1.05 was 99.996 % accurate, and  
15 InChI version 1.06 represents a step closer to perfection. Finally, we look forward to future  
16 releases and extensions for the InChI Chemical identifier.

17

## 18 **Keywords**

19 InChI, InChIKey, PubChem, RInChI

20

## 21 **Introduction**

22 The first version of InChI was made publicly available in the spring of 2005 and further  
23 versions [1-5], including a separate InChI for Reactions (RInChI [6-7]), have been released  
24 over the years. The original version covered much of organic chemistry and it was quickly

1 adopted and used by chemists and scientific databases throughout the world. Its application  
2 in virtual databases includes SAVI-2020, a database of reactions, labelled by RInChI, with  
3 over 1.75 billion products, all of which are assigned InChI [8]. The effective coverage of  
4 chemistry had a downside: the pressure to extend and add features and capabilities was  
5 rather small. The InChI software, coded in the C language, contains about 200,000 lines of  
6 code and comments. Such a large computer program will always have bugs and  
7 imperfections. We are very grateful to everyone that has reported issues and look forward  
8 to more feedback from the large and growing community of InChI users.

9

10 The InChI is canonical. This central feature of the design requires that for every molecule  
11 there is just one InChI, and every InChI identifies just one molecule. This makes it possible  
12 for the InChI to be used as an identifier to link identical structures in databases or on the  
13 internet. It has been and is being used by many commercial and non-commercial software  
14 chemistry packages. The InChI fulfils a different role to SMILES [9], which was designed to be  
15 easily understood both by chemists and by computers. As a consequence, there are many  
16 valid SMILES representations of every molecule. This distinction makes SMILES a better tool  
17 for some applications than InChI, and InChI a better tool for others, whenever a unique  
18 identifier is helpful. There are a variety of schemes to generate canonical SMILES, but just  
19 one InChI algorithm.

20

21 The issue is more subtle than simply finding a globally-agreed canonicalization algorithm. It  
22 is also necessary to address the fundamental question: what is a molecule? InChI provides a  
23 consistent and effective answer to this question, which covers nearly all of the molecules  
24 used in medicinal chemistry. SMILES has the option of encoding molecular detail that InChI

1 does not, which makes it more flexible in its descriptions of complex objects, and the  
2 disadvantage that an object that some people consider to be one molecular system is  
3 considered as several different molecular systems by other people, with different SMILES  
4 strings. For example, 2-methylpyridine, which is also known as picoline, can be drawn in two  
5 different ways with different double bond positions (Fig. 1). Both of these representations  
6 are correct. There is just one InChI for both forms: InChI=1S/C6H7N/c1-6-4-2-3-5-7-6/h2-  
7 5H,1H3. There are many valid SMILES, including CC1=CC=CC=N1 and CC1=NC=CC=C1.

8



9

10 **Fig. 1:** these two molecules are the same as each other, and should have the same name  
11

12 The SMILES contain more information, because they record the localisation of the double  
13 bonds and the InChI does not. The localisation can lead to problems in assigning  
14 stereochemistry using the Cahn-Ingold-Prelog rules [10]. In this particular example, the  
15 extra information does not correspond to any physically observable distinction for the  
16 molecule these strings represent. Might there be molecules for which a different InChI is  
17 applied to variant representations of the same molecule, or else an important distinction in  
18 structure is missed in InChI generation? The choices that the InChI algorithm makes about  
19 what is and is not represented control how effective it is at producing useful molecular  
20 identifiers.

21

22 Tautomers represent one aspect of this question. Are two tautomers the same molecule or  
23 different molecules? There are many cases for which the answer depends on the detail of  
24 the question. Tautomers may not be separable when stored in a bottle in a storeroom but

1 the distinction between them may be important in understanding a reaction mechanism,  
2 during which they may be formed transiently, or at a very low temperature. Should they  
3 have the same InChI or different InChI? A detailed investigation of this issue is underway  
4 [11] and is beyond the scope of this paper. A few molecular features, such as  
5 atropisomerism and some aspects of organometallic structures, are not described by the  
6 current version of the InChI. Work is underway to introduce these features into future  
7 releases of the InChI.

8  
9 In this manuscript we report the details of the latest upgrade, version 1.06, which fixes a  
10 series of known issues, and the results of its extensive testing on the PubChem Compound  
11 database provided by the U.S. National Center for Biotechnology Information, containing  
12 more than one hundred million structures. [12]

13

#### 14 **Results: details of upgrade**

15 InChI v1.06 is available for download from the InChI Trust website [13]. In addition to the  
16 software and its source code, the downloads include a detailed log of changes, an API  
17 reference, and a list of known issues. Comments are welcomed on the SourceForge  
18 discussion list [14].

19

#### 20 *Upgrade to v. 1.05*

21 The key features introduced in v. 1.05 were the ability to treat polymers and the ability to  
22 generate InChI for large molecules containing up to 32767 atoms. These are still considered  
23 to be experimental features. For this reason, the InChI generated for polymers and for  
24 molecules with more than 1024 atoms are not standard InChI (InChI=1S) but are marked as

1 beta non-standard InChI, indicated by InChI=1B prefix. A partial ability to use Molfile v3000  
2 as input format was also added to v. 1.05. This facility is important in handling large  
3 molecules as molecules in the Molfile v2000 format cannot contain more than 999 atoms.

4

#### 5 *Upgrade to v. 1.06*

6 The upgrade to v1.06 includes support for pseudo element atoms, labelled "Zz" or "\*", and  
7 improved support for single-strand polymers options. All of these facilities are currently only  
8 available in the non-standard InChI, labelled with the prefix InChI=1B. In addition, there are  
9 a variety of bugfixes, updates to the API library, support for threading building blocks (TBB)  
10 scalable memory allocators and several convenience features.

11

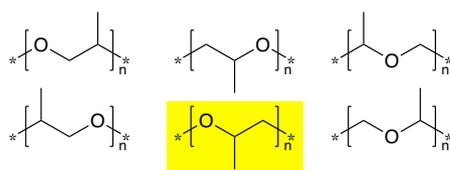
12 Support for single-strand polymers was first introduced in v. 1.05 and has been extended  
13 and improved in v. 1.06. Both structure-based and source-based representation and  
14 encoding of polymers are available. Source-based representation encodes chemical  
15 structures of the starting material(s) with an indication of polymer nature, type of polymer  
16 (block, random, and alternating) and the role and order of the components where needed.  
17 Structure-based representation of polymers is based on the structure of structural repeating  
18 units, SRU, sometimes called constitutional repeating units, CRU, enclosed in polymer  
19 brackets. Note that source- and structure-based representations and their InChI encodings  
20 are independent and, in general, no procedures are implied for algorithmic conversion and  
21 relation from one type to the other.

22

23 A pseudo-element atom is a generic placeholder designating an undefined, unknown or  
24 variable entity, and is indicated by "Zz" or "\*". The exact meaning is not defined within the

1 InChI standard. It may be used to indicate the places where constitutional repeat units in  
 2 polymers join on to each other, or, more generally, to indicate undefined univalent atoms.  
 3 Stereochemistry adjacent to Zz atoms is disabled by default, although an option is provided  
 4 to switch it on. This feature makes it possible to describe molecular systems more flexibly,  
 5 but it may be harder to describe them canonically. The central importance of canonical  
 6 representation for the InChI is the reason that these features are not a part of the standard  
 7 InChI (prefix InChI=1S) and restricted to the non-standard InChI (prefix InChI=1B).

8  
 9 This new feature improves the ability of InChI to describe polymers. For example,  
 10 polypropylene glycol (PPG, [-O-CH<sub>2</sub>-CH(CH<sub>3</sub>)-]<sub>n</sub>) can be described, in structure-based  
 11 representation, in several different but equivalent ways, Fig. 2.



12  
 13  
 14 **Fig. 2:** polypropylene glycol can be described with any of these monomers. InChI version  
 15 1.06 selects the highlighted one as the canonical version.

16  
 17 All six of these representations, which use a \* to indicate the Zz atom, are correct; the InChI  
 18 algorithm selects the highlighted one, in the middle of the bottom row, as the canonical  
 19 one: InChI=1B/C3H6OZz2/c1-3(2-5)4-6/h3H,2H2,1H3/z101-1-4(6-4,5-2)

20  
 21 Note that the canonical form here is exactly the same as one preferred by IUPAC  
 22 recommendations [15]. Special care has been taken in developing the InChI algorithm to  
 23 ensure that the basic IUPAC criteria are incorporated, but the correspondence between the

1 preferred InChI and the IUPAC recommendation (targeted to manual selection of preferred  
2 form) is not always perfect.

3

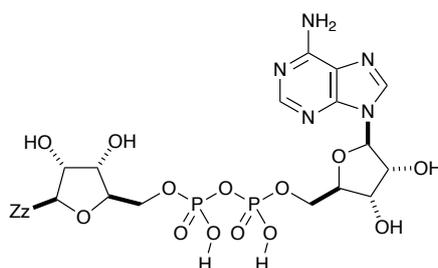
4 The InChI v1.05 preferred form for this polymer was different as the two Zz atoms,  
5 numbered 5 and 6 in this new InChI, were not available. The InChI v1.06 code can generate  
6 the older form of the polymer InChI using the Polymers105 option, but its use is deprecated.

7

8 *Pseudo element atoms*

9 Pseudo-elements use is not restricted to polymers. For example, Fig. 3, the  
10 adenosinediphosphoribosyl group is a molecular fragment, included in the EBI's directory of  
11 Chemical Entities of Biological Interest (ChEBI) as structure CHEBI:22259. InChI cannot  
12 describe molecular fragments, but the Zz atom in version 1.06 provides a non-standard-  
13 InChI for this.

14



15

16

**Fig. 3:**

17 InChI=1B/C15H22N5O13P2Zz/c16-12-7-13(18-3-17-12)20(4-19-7)14-10(23)8(21)5(31- 14)1-  
18 29-34(25,26)33-35(27,28)30-2-6-9(22)11(24)15(36)32-6/h3-6,8-11,14-15,21- 24H,1-  
19 2H2,(H,25,26)(H,27,28)(H2,16,17,18)/t5-,6-,8-,9-,10-,11-,14-/m1/s1

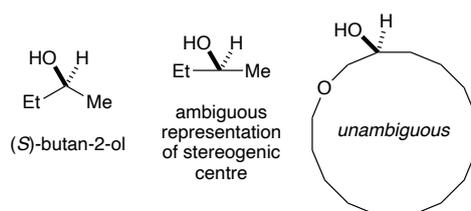
20

21 *Structure-based issues*

22 Translation, from two-dimensional representations of molecules in Molfiles to an InChI,  
23 requires interpretation. The InChI code is very effective at doing this, and version 1.06

1 contains a few improvements. One of these focusses on the representation of large rings,  
2 which are used in important classes of molecules including peptides and macrolides. In Fig.  
3 4, the left-hand structure is a clear representation of (S)-butan-2-ol. The stereochemistry of  
4 the middle structure is unclear, because the ethyl and methyl groups are represented in a  
5 straight line. The InChI code will not give a stereochemical assignment to such a structure.  
6 How straight does a line have to be, before giving a stereochemical assignment is  
7 inappropriate? The right-hand structure illustrates a molecule for which the stereochemistry  
8 can be easily interpreted, because of the large ring, and yet the chain is almost straight. In  
9 version 1.06, the InChI code allows for representations of this type to be interpreted by  
10 introducing a new option: *LooseTSACheck*. Because this option is for structure perception  
11 and not InChI generation, it can be used within the standard InChI.

12



13

14 **Fig. 4:** version 1.06 recognises the large ring and so assigns stereochemistry correctly to the  
15 right-hand structure

16

17 The InChI software is now able to identify and assign the stereochemistry of 1,7-

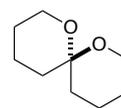
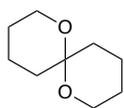
18 dioxaspiro[5.5]undecane and related structures (Fig. 5). *R* and *S* 1,7-

19 dioxaspiro[5.5]undecane are mirror images of each other, although this is not immediately

20 obvious from the diagram. InChI version 1.06 can now take structures drawn in the form on

21 the right of the figure, and assign the correct stereochemical layer.

22



1 1,7-dioxaspiro[5.5]undecane (S)-1,7-dioxaspiro[5.5]undecane (R)-1,7-dioxaspiro[5.5]undecane (R)-1,7-dioxaspiro[5.5]undecane

2 **Fig. 5:** with version 1.05, the InChI of both *R* and *S* 1,7-dioxaspiro[5.5]undecane is  
3 InChI=1S/C9H16O2/c1-3-7-10-9(5-1)6-2-4-8-11-9/h1-8H2. With version 1.06, a  
4 stereochemistry layer is added, and InChI=1S/C9H16O2/c1-3-7-10-9(5-1)6-2-4-8-11-9/h1-  
5 8H2/t9-/m1/s1 is generated from the right-hand structure.  
6

7 Two cases for which version 1.05 was unable to generate an InChI string have been fixed, as  
8 have some renumbering issues. The details are given in the description of tests with  
9 PubChem, below. Some large molecules supplied in Molfile V3000 formats caused a  
10 problem which has now been fixed, as have minor issues in AuxInfo and the non-standard  
11 InChI description of polymers.

### 12 *Other improvements*

14 Memory allocation can become a serious performance bottleneck when using the InChI  
15 library in multi-threading environment, as threads may compete for a global lock related to  
16 a single global heap. In this situation, the program's behaviour is not scalable and speed  
17 may degrade if number of processor cores increases. Intel TBB is free software package  
18 available for both Windows and Linux, licensed under the Apache License, which  
19 automatically replaces C functions for dynamic memory allocation with its own scalable  
20 memory allocators to avoid contention in most cases. This method may optionally be used  
21 with InChI Software library and, in some cases, may provide performance gain.

22  
23 A new switch, *WMnumber*, has been added which sets InChI calculation timeout in  
24 milliseconds. This complements *Wnumber* which sets the value in seconds. The finer

1 granularity is useful for managing long runs on datasets containing many millions of  
2 molecules.

3

#### 4 *Security Issues*

5 The InChI software expects valid InChI strings or Molfiles as its input. Other input should be  
6 rejected with a suitable diagnostic message. A small number of cases have been discovered  
7 for which other input files (mostly artificially corrupted) can lead to memory corruption or a  
8 crash and there is a possibility that this could cause problems for some applications. A list of  
9 people who have pointed out specific issues is in the acknowledgements section of this  
10 paper. Numerous modifications to the code have been made to catch these problems.

11

#### 12 **Results: tests on PubChem**

13 The InChI algorithm is intended to work on all molecules. Demonstrating that it does this,  
14 however, is a non-trivial task. How well does it work on molecules that are not merely un-  
15 made, but also un-imagined? A step towards answering this question is to test the algorithm  
16 on a large, curated dataset. PubChem [12] is such a database and contained 111 476 790  
17 molecules at the time of testing. Standard InChI strings have been generated for all of these  
18 molecules.

19

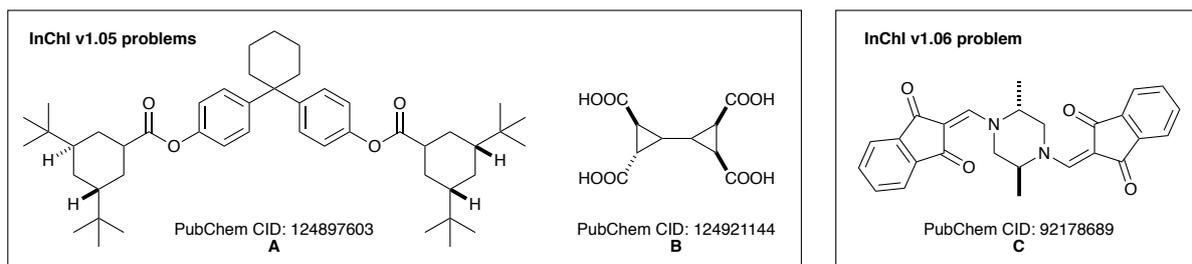
20 Several types of imperfection can be envisaged: (i) failure to create an InChI string; (ii)  
21 making the same InChI string for two different molecules; (iii) generating two different InChI  
22 strings for the same molecule; (iv) different InChI strings for v1.05 and v1.06.

23

1 The first issue, failure to create an InChI string, is serious, but easy to catch because an error  
 2 message is generated. This problem occurs for two structures for InChI v1.05 and one  
 3 structure for InChI v1.06 (Fig. 6). The first two, **A** and **B**, both have an undefined  
 4 stereocentre, and InChI v1.05 gives a STEREOCOUNT\_ERR error and does not generate an  
 5 InChI. The enantiomer of **B** is also present in PubChem (CID: 92286308) and InChI v1.05 is  
 6 able to generate an InChI for this. InChI v1.06 generates InChI strings which note, correctly,  
 7 that one of the stereocentres is undefined in each case:

8 **A:** InChI=1S/C48H72O4/c11-44(2,3)36-26-32(27-37(30-36)45(4,5)6)42(49)51-40-20-16-34(17-  
 9 21-40)48(24-14-13-15-25-48)35-18-22-41(23-19-35)52-43(50)33-28-38(46(7,8)9)31-39(29-  
 10 33)47(10,11)12/h16-23,32-33,36-39H,13-15,24-31H2,1-12H3/t32?,36-,37+,38-,39-/m0/s1  
 11 **B:** InChI=1S/C10H10O8/c11-7(12)3-1(4(3)8(13)14)2-5(9(15)16)6(2)10(17)18/h1-  
 12 6H,(H,11,12)(H,13,14)(H,15,16)(H,17,18)/t1?,3-,4+,5-,6-/m0/s1

13



14

15

**Fig. 6:** Problem structures from PubChem for InChI v1.05 and InChI v1.06

16

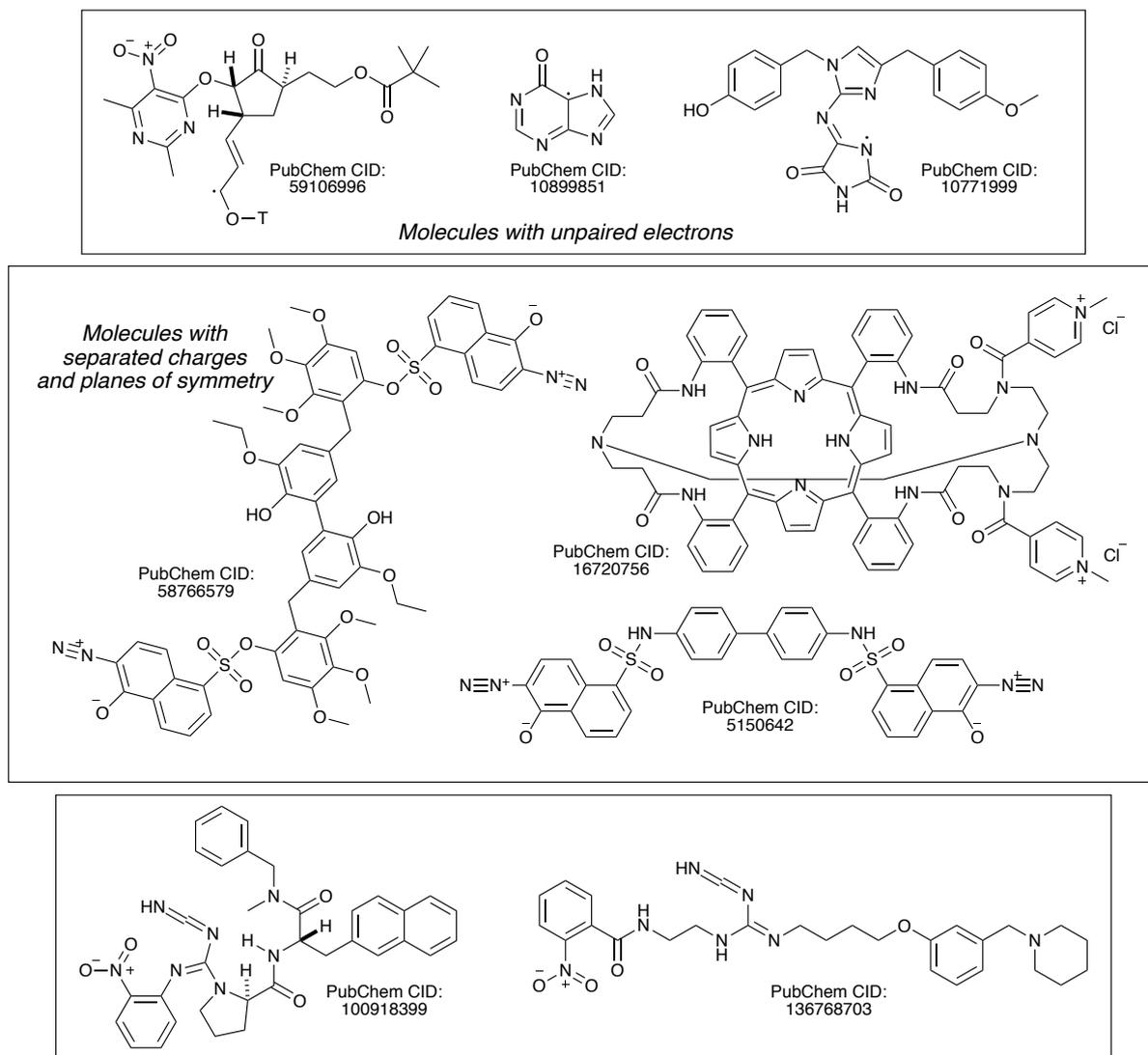
17 Neither version has any difficulty with the stereoisomers for which all the stereocentres are  
 18 defined, and both produce the same InChI strings. For molecule **C** (PubChem CID:  
 19 92178689), InChI v1.05 generates an InChI string, and InChI v1.06 may stop with a  
 20 CANON\_ERR (this appears only in some environments, e.g., Windows/inchi-1 executable but  
 21 not Linux-64bit/libinchi.so library; the behaviour is expected to be fixed in the next release).

1 There appear to be no other errors of this type for PubChem molecules, and so InChI v1.06  
2 has an error rate of about one in 111 476 790, which can be expressed as 99.999999 %  
3 reliable.

4  
5 The second issue, which would clearly be an error, is generating the same InChI string for  
6 two different molecules. When two molecules differ only by the localisation of protons on  
7 heteroatoms, they may have the same InChI and would be expected to be inseparable. With  
8 these exceptions, no such errors were observed in this survey of PubChem.

9  
10 The third issue, two different InChI strings for one molecule, is also an error. This was  
11 observed as a consequence of numbering instability. A key feature of the InChI software is  
12 that it numbers all of the atoms in a molecule, except for the hydrogens. This numbering  
13 should be the same whatever the order of the atoms in the input Molfile. A molecule with N  
14 non-hydrogen atoms can be represented in a Molfile in about N! different ways before the  
15 coordinates are considered. All of these possibilities should produce the same InChI string,  
16 and they nearly always do. However, 547 examples have been discovered in PubChem for  
17 which the numbering of the InChI string can be changed by altering the order of the atoms  
18 in the Molfile. These are all listed in the *KnownIssues* file which is a part of the InChI v1.06  
19 download. Nearly all of these molecules contain nitrogen (530 out of 547) and many contain  
20 sulfur (453 out of 547), but only 36 contain phosphorous. The majority of these molecules  
21 have a net charge (381 positive; 49 negative). Many contain metals (Na: 57; Zn: 6; Ru: 6; K:  
22 4; Fe: 3; Co: 1; Mn: 1; Ni: 1) and some contain halogens (F: 54; Cl: 92; Br: 1; I: 3). Only 84 of  
23 the 547 molecules were neutral and did not contain metals; all of these molecules contain  
24 nitrogen and three have unpaired electrons. The most distinctive thing about this list of 547

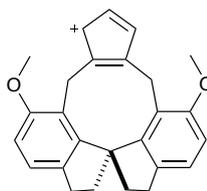
1 molecules is that 534 of them contain a “/p” layer, indicating that a proton has to be added  
 2 or removed from the formula to give the input composition. Of the remaining 13 molecules,  
 3 five have a “/q” layer, indicating a charge on the molecule. This leaves eight, which are  
 4 illustrated in Fig. 7. With a problem that affects 547 molecules out of 111 476 790, InChI  
 5 v1.06 is 99.9995% reliable.



**Fig. 7:** All molecules with numbering instability that lack “/p” and “/q” layers

The final issue, when v1.05 and v1.06 produce different strings, may be a result of improved structural perception for v1.06, but may also be inconvenient as molecules that are the

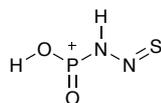
1 same are associated with two different InChI strings, one for v1.05 and another for v1.06.  
2 There are 6524 examples of this in PubChem, none of which show numbering instability.  
3 Three of these contain just carbon, hydrogen and oxygen; two of these three are Fig. 6, **A**  
4 and **B**. The third is illustrated at the top of Fig. 8. This structure has a stereocentre in a spiro-  
5 cyclic motif, which is missed by v1.05 and found by v1.06. The InChI are identical except that  
6 v1.06 adds the stereochemical layers: /t26-/m0/s1. In contrast to molecules with unstable  
7 numbering, none of the 6524 molecules in this group have a net charge indicated in the  
8 molecular formula. However, almost all of them (6388) have a “/p” layer when calculated  
9 with InChI v1.06, and 434 have a “/q” layer. 364 have both “/p” and “/q”. This contrasts with  
10 the InChI v1.05 calculations on this group, for which 6169 have a “/p” layer, 6323 have a  
11 “/q” layer, and 6154 have both “/p” and “/q”. Most of these cases are the result of a fix to a  
12 problem with molecules with acidic hydroxy groups at cationic heteroatom centres, which  
13 led to issues with numbering atoms. The details of this are in the CHANGELOG file in the  
14 release. In v1.05, such systems were often labelled with both a “/q” and a “/p” layer: “/q-  
15 1/p+2” was very common, and occurs in 5446 of the v1.05 InChI strings in this group of 6524  
16 molecules. In v1.06, the two layers are replaced by one layer: “/p+1” and the extra  
17 hydrogen atom is included in the formula of the molecule. Three examples are in Fig. 8.  
18



PubChem CID: 134958543

InChI v1.05: InChI=1S/C26H25O2/  
c1-27-22-8-6-16-10-12-26-13-11-17-7-9-23(28-2)21(25(17)26)15-19-5-3-4-18  
(19)14-20(22)24(16)26/h3-9H,10-15H2,1-2H3/q-1

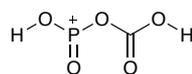
InChI v1.06: InChI=1S/C26H25O2/  
c1-27-22-8-6-16-10-12-26-13-11-17-7-9-23(28-2)21(25(17)26)15-19-5-3-4-18  
(19)14-20(22)24(16)26/h3-9H,10-15H2,1-2H3/q-1/t26-/m0/s1



PubChem CID: 154400441

InChI v1.05: InChI=1S/N2O2PS/c3-5(4)1-2-6/q-1/p+2

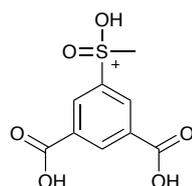
InChI v1.06: InChI=1S/HN2O2PS/c3-5(4)1-2-6/h(H-,1,3,4,6)/p+1



PubChem CID: 20681196

InChI v1.05: InChI=1S/CO5P/c2-1(3)6-7(4)5/q-1/p+2

InChI v1.06: InChI=1S/CHO5P/c2-1(3)6-7(4)5/h(H-,2,3,4,5)/p+1



PubChem CID: 147555804

InChI v1.05: InChI=1S/C9H7O6S/c1-16(14,15)7-3-5(8(10)11)2-6(4-7)9(12)13/  
h2-4H,1H3,(H-2,10,11,12,13,14,15)/q-1/p+2

InChI v1.06: InChI=1S/C9H8O6S/c1-16(14,15)7-3-5(8(10)11)2-6(4-7)9(12)13/  
h2-4H,1H3,(H2-,10,11,12,13,14,15)/p+1

1

2

3

**Fig. 8:** changes from v1.05 to v1.06

4 Only 35 molecules out of 6524 have the same "/q" layer for both InChI v1.05 and InChI  
5 v1.06. In 18 of these cases, the match arises because of the absence of a "/q" layer in both  
6 versions of the InChI string. In all 35 of these cases, the difference between v1.05 and v1.06  
7 arises in the "/t" (stereochemistry) layer. For PubChem CID: 6555836, 6589644, 6589645,  
8 11871423, 18805145, 18805148, 18805149, 18805151, 18805146, 49950537, 49950540,  
9 101988808 the stereochemistry is not clear in PubChem. This illustrates the power of the  
10 InChI string to highlight molecules within a database that would benefit from further  
11 checking. These checks demonstrate that the transition from v1.05 to v1.06 is 99.99 %  
12 consistent for the standard InChI string.

13

14 **Discussion**

1 InChI v1.06 provides substantial new functionality and corrects some issues with earlier  
2 releases, whilst retaining a very high level of backwards compatibility. Changes to the InChI  
3 can be readily incorporated in Reaction InChI (RInChI) [6-7] and into on-going work to  
4 develop an InChI-based description of mixtures [16].

5  
6 The frequency of the errors we have discovered in this survey of PubChem can be compared  
7 with the number of InChIKey collisions. It was shown by one of us (JMG, [17-18]) that InChIKey  
8 collisions can occur. The paper [19] analyzes the situation in details and demonstrates that  
9 collisions, though unavoidable, occur at a very small and expected rate. The probability of an  
10 InChIKey collision in a database of the current size of PubChem is slightly above 0.01 %. The  
11 fact that InChIKey collisions can occur is important to consider, even though their occurrence  
12 is so rare that they have been found only by deliberate searching. The observed failure rate  
13 of the InChI algorithm of one structure in 111 476 790 is slightly higher than this, and so it  
14 also needs to be kept under consideration.

## 15 16 **Conclusions**

17 That latest version of the InChI code, version 1.06, provides a significant advance in  
18 functionality and reliability, whilst retaining compatibility with previous releases. The  
19 upgraded InChI can be used now both as a molecular identifier and as a part of the Reaction  
20 InChI (RInChI) process for labelling reactions. [6-7] The InChI Trust [20] and the IUPAC InChI  
21 subcommittee [21] are always pleased to hear suggestions for new functionality and reports  
22 of issues with the current code. These should be reported either to the InChI SourceForge  
23 website and discussion list [13], or by e-mail to the secretary of the IUPAC InChI  
24 subcommittee (Jonathan Goodman, [jonathan@inchi-trust.org](mailto:jonathan@inchi-trust.org)).

1

2 We continue to work on expanding the capabilities of InChI. Organometallics and extended  
3 stereochemistry are the likely next additions to InChI capabilities. There are two  
4 fundamental issues with further InChI developments. First there is the human element with  
5 diverse opinions on what constitutes the correct answer in representing molecular  
6 structures. The second is the chemistry itself, which is not always well resolved and  
7 unambiguous. Like the InChI capabilities in the area of polymers, where the properties of a  
8 polymer are often of more concern to the chemist than the structure, which may not be  
9 known in detail. InChI is and always will be limited by how good chemists and chemistry are  
10 at representing a structure. The InChI Trust is also developing the functionality of the InChI  
11 beyond the representation of single molecules: Reactions, mixtures, Markush structures, QR  
12 codes for InChI, and other areas are all being developed and will be made available as soon  
13 as possible.

14

15 We hope that InChI v1.06 will be widely used. It retains compatibility with previous standard  
16 InChI releases, provides additional functionality, and is, on the basis of tests with PubChem,  
17 more than 99.99 % reliable.

18

## 19 **Declarations**

20

## 21 **Availability of data and materials**

22 The InChI v1.06 software is available for download from the InChI Trust website. [20]

23

## 24 **Acknowledgements**

1 We are grateful to many people and organizations who are continuously supporting InChI  
2 development, in various forms. In particular we wish to thank the InChI Trust which is the  
3 organization funding the further developments and enchantments to the InChI algorithm.  
4 Thanks are due to those who provided feedback, reported on problems and participated in  
5 discussions, namely, Google AutoFuzz team and Ian Wetherbee and Steve Boyer; Cure53  
6 team; Andrew Dalke; Burt Leland; Yulia Borodina; Marc Nicklaus; Gerd Blanke; Evan Bolton;  
7 Noel O'Boyle; Dmitry Redkin; Andrey Yerin; Daniel Lowe; John Mayfield; Igor Filippov; Wolf-  
8 Dietrich Ihlenfeldt; Roger Sayle; Egon Willighagen; Peter Linstrom; Karl Nedwed; John  
9 Salmon; Lutz Weber; Alex Clark and to many others who specifically contributed to the  
10 current release/testing, especially Dmitrii Tchekhovskoi.

11

#### 12 **Competing interests**

13 There are no competing interests.

14

#### 15 **Funding**

16 The InChI Trust is thanked for funding of the IUPAC InChI project. The work of PT and EB  
17 was supported by the Intramural Research Program of the National Library of Medicine,  
18 National Institutes of Health.

19

#### 20 **Authors' contributions**

21 IP prepared the InChI v1.06 release and ran various tests. EB and PT designed and ran the  
22 PubChem testing. JMG designed the study and performed the analysis. SRH and JMG wrote  
23 the paper. All authors read and approved the final manuscript.

24

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

## References

1. Warr WA (2015) Many InChIs and quite some feat. *J Comput-Aided Mol Des* 29, 681. doi:10.1007/s10822-015-9854-3
2. McNaught AD, Heller SR (2011) *The IUPAC Chemical Identifier (InChI), Principles of Chemical Nomenclature - A Guide to IUPAC Recommendations 2011 Edition*, pages 190-196. ISBN 978-1-94973-007-5
3. Heller S, McNaught A, Pletnev I, Stein S, Tchekhovskoi D, Pletnev I (2013) InChI - the worldwide chemical structure identifier standard. *Journal of Cheminformatics* 5, 7. doi:10.1186/1758-2946-5-7
4. Heller SR, McNaught A, Pletnev I, Stein S, Tchekhovskoi D (2015) InChI - the IUPAC International Chemical Identifier. *Journal of Cheminformatics* 7, 23. doi:10.1186/s13321-015-0068-4
5. Boucher R, Heller S, McNaught A (2017) The Status of the IUPAC InChI Chemical Structure Standard. *Chemistry International*, 48.
6. Grethe G, Goodman JM, Allen CHG (2013) International chemical identifier for reactions (RInChI). *J. Cheminformatics* 5, 45. doi: 10.1186/1758-2946-5-45
7. Grethe G, Blanke G, Kraut H, Goodman JM (2018) International chemical identifier for reactions (RInChI). *J. Cheminformatics* 10, 22. doi: 10.1186/s13321-018-0277-8
8. Patel H, Ihlenfeldt WD, Judson PN, Moroz Y, Pevzner Y, Peach M, Tarasova N, Nicklaus M (2020). *Synthetically Accessible Virtual Inventory (SAVI) (Version 2020)*. CADD Group, CBL, CCR, NCI, NIH. doi:10.35115/37N9-5738.

- 1 9. Weininger D (1988) SMILES, a chemical language and information system. 1.  
2 Introduction to methodology and encoding rules. *J Chem Inform Comput Sci*, 28(1),  
3 31-36.
- 4 10. Hanson RM, Musacchio S, Mayfield JW, Vainio MJ, Yerin A, Redkin D (2018)  
5 Algorithmic analysis of Cahn–Ingold–Prelog rules of stereochemistry: proposals for  
6 revised rules and a guide for machine implementation. *J. Chem. Inf. Model.* 58, 9,  
7 1755–1765. doi: 10.1021/acs.jcim.8b00324
- 8 11. Dhaked, DK; Ihlenfeldt, WD; Patel, H; Delanee, V; Nicklaus, MC (2020) Toward a  
9 Comprehensive Treatment of Tautomerism in Chemoinformatics Including in InChI  
10 V2. *J. Chem. Inf. Model.* 60, 1253-1275. doi:10.1021/acs.jcim.9b01080
- 11 12. Kim S, Chen J, Cheng T, Gindulyte A, He J, He S, Li Q, Shoemaker BA, Thiessen PA, Yu  
12 B, Zaslavsky L, Zhang J, Bolton EE (2019) PubChem 2019 update: improved access to  
13 chemical data. *Nucleic Acids Res.* 47(D1):D1102-1109. doi:10.1093/nar/gky1033
- 14 13. Downloads of InChI Software. <https://www.inchi-trust.org/downloads/> Accessed 1  
15 January 2021.
- 16 14. SourceForge mailing lists: inchi-discuss (InChI Facilities and Applications).  
17 <https://lists.sourceforge.net/lists/listinfo/inchi-discuss> Accessed 1 January 2021.
- 18 15. Kahovec, J.; Fox, R. B.; Hatada, K. Nomenclature of Regular Single-Strand Organic  
19 Polymers (IUPAC Recommendations 2002). *Pure and Applied Chemistry* 2002, 74,  
20 1921–1956.
- 21 16. Clark AM, McEwen LR, Geddeck P, Bunin BA (2019) Capturing mixture composition: an  
22 open machine-readable format for representing mixed substances. *J.*  
23 *Cheminformatics* 11, 33. doi:10.1186/s13321-019-0357-4

- 1 17. Goodman JM (2009) Reliable Reactions and Stable Structures. Abstracts of papers of  
2 the American Chemical Society 238, CINF18. Washington, DC
- 3 18. Goodman JM (2011) RInChIs and reactions. Abstracts of papers of the American  
4 Chemical Society 242, CINF40. Denver, CO.
- 5 19. Pletnev I, Erin A, McNaught A, Blinov K, Tchekhovskoi D, Heller S (2012) InChIKey  
6 collision resistance: an experimental testing. J. Cheminformatics 4, 39.  
7 doi:10.1186/1758-2946-4-39
- 8 20. InChI Trust. <https://www.inchi-trust.org> Accessed 1 January 2021.
- 9 21. IUPAC InChI subcommittee. [https://iupac.org/who-we-are/committees/committee-  
11 details/?body\\_code=802](https://iupac.org/who-we-are/committees/committee-<br/>10 details/?body_code=802) Accessed 1 January 2021.

## Figures



Figure 1

these two molecules are the same as each other, and should have the same name

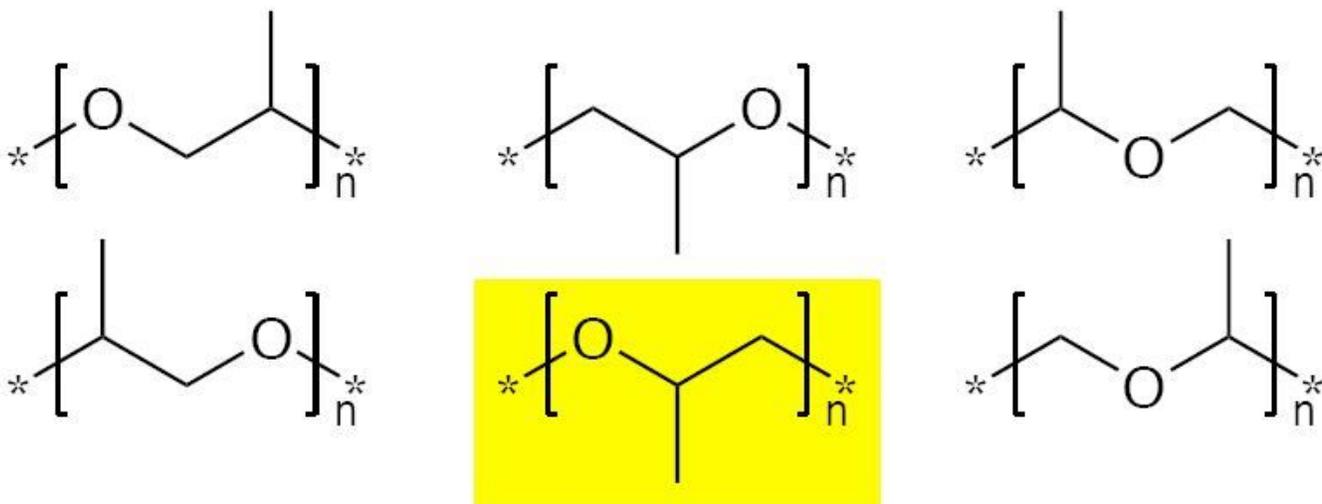
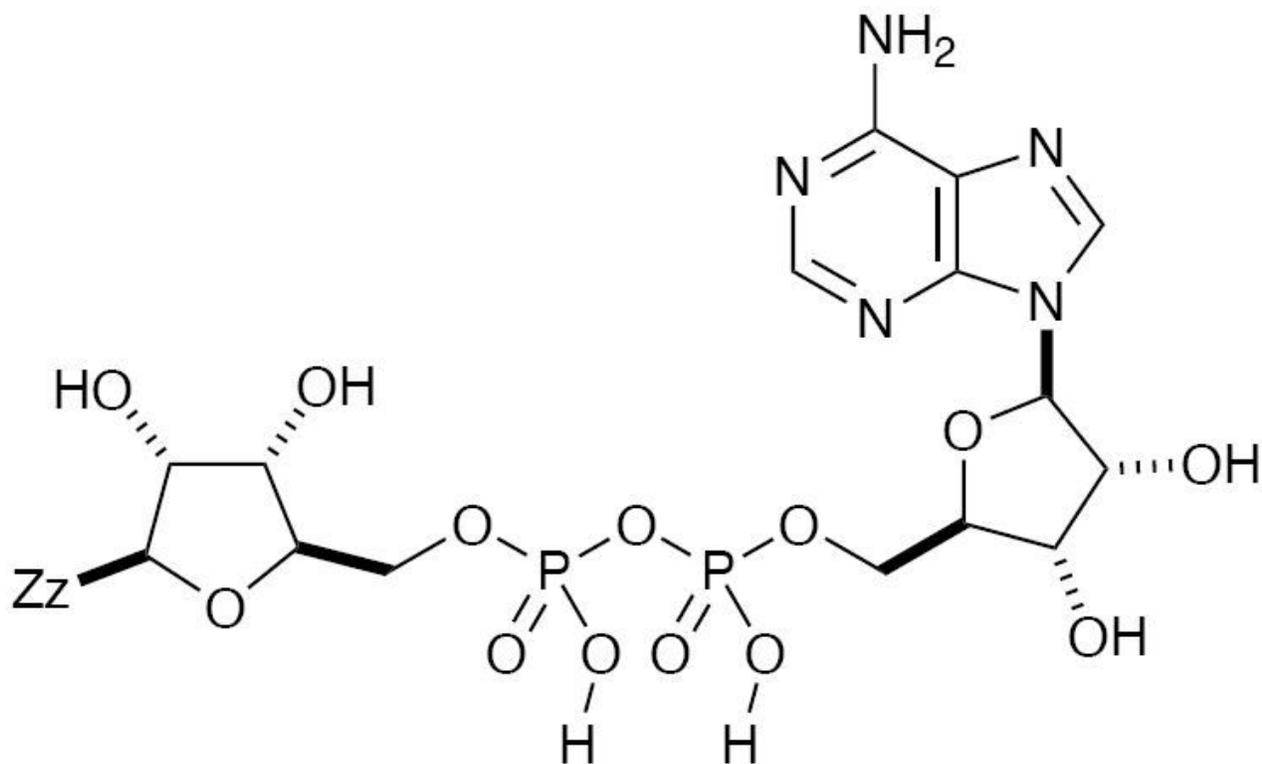


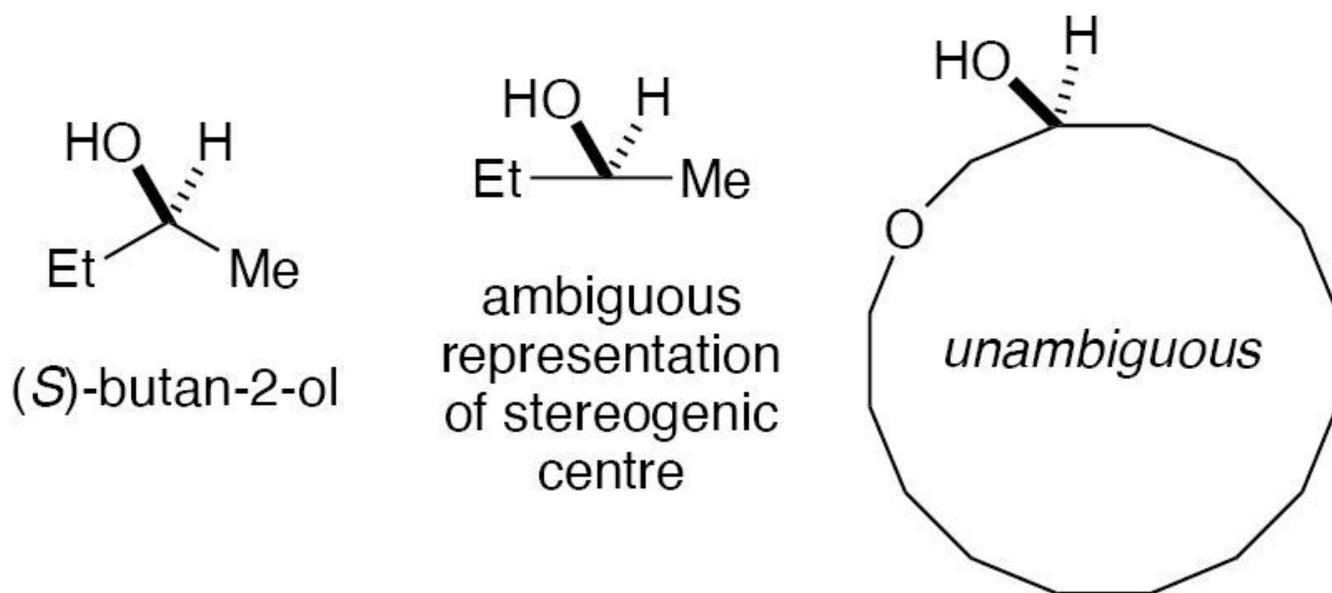
Figure 2

polypropylene glycol can be described with any of these monomers. InChI version 1.06 selects the highlighted one as the canonical version.



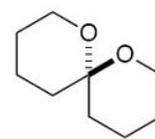
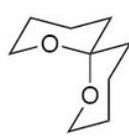
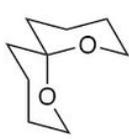
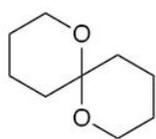
**Figure 3**

InChI=1B/C15H22N5O13P2Zz/c16-12-7-13(18-3-17-12)20(4-19-7)14-10(23)8(21)5(31-14)1-29-34(25,26)33-35(27,28)30-2-6-9(22)11(24)15(36)32-6/h3-6,8-11,14-15,21-24H,1-2H2,(H,25,26)(H,27,28)(H2,16,17,18)/t5-6,8-9,10,11,14-/m1/s1



**Figure 4**

version 1.06 recognises the large ring and so assigns stereochemistry correctly to the right-hand structure



1,7-dioxaspiro[5.5]undecane (S)-1,7-dioxaspiro[5.5]undecane (R)-1,7-dioxaspiro[5.5]undecane (R)-1,7-dioxaspiro[5.5]undecane

**Figure 5**

with version 1.05, the InChI of both R and S 1,7-dioxaspiro[5.5]undecane is InChI=1S/C9H16O2/c1-3-7-10-9(5-1)6-2-4-8-11-9/h1-8H2. With version 1.06, a stereochemistry layer is added, and InChI=1S/C9H16O2/c1-3-7-10-9(5-1)6-2-4-8-11-9/h1-8H2/t9-/m1/s1 is generated from the right-hand structure.

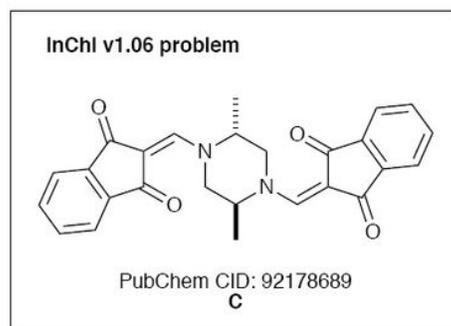
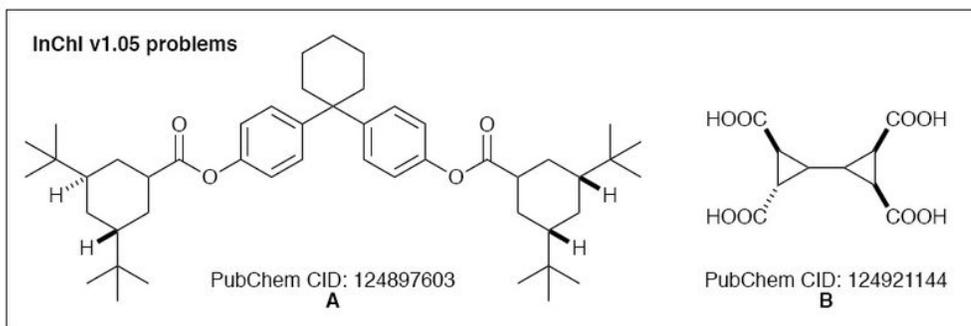


Figure 6

Problem structures from PubChem for InChI v1.05 and InChI v1.06

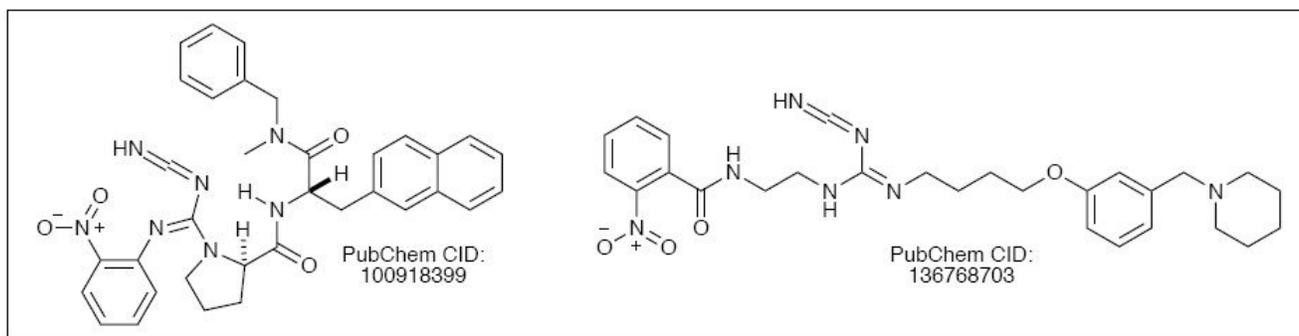
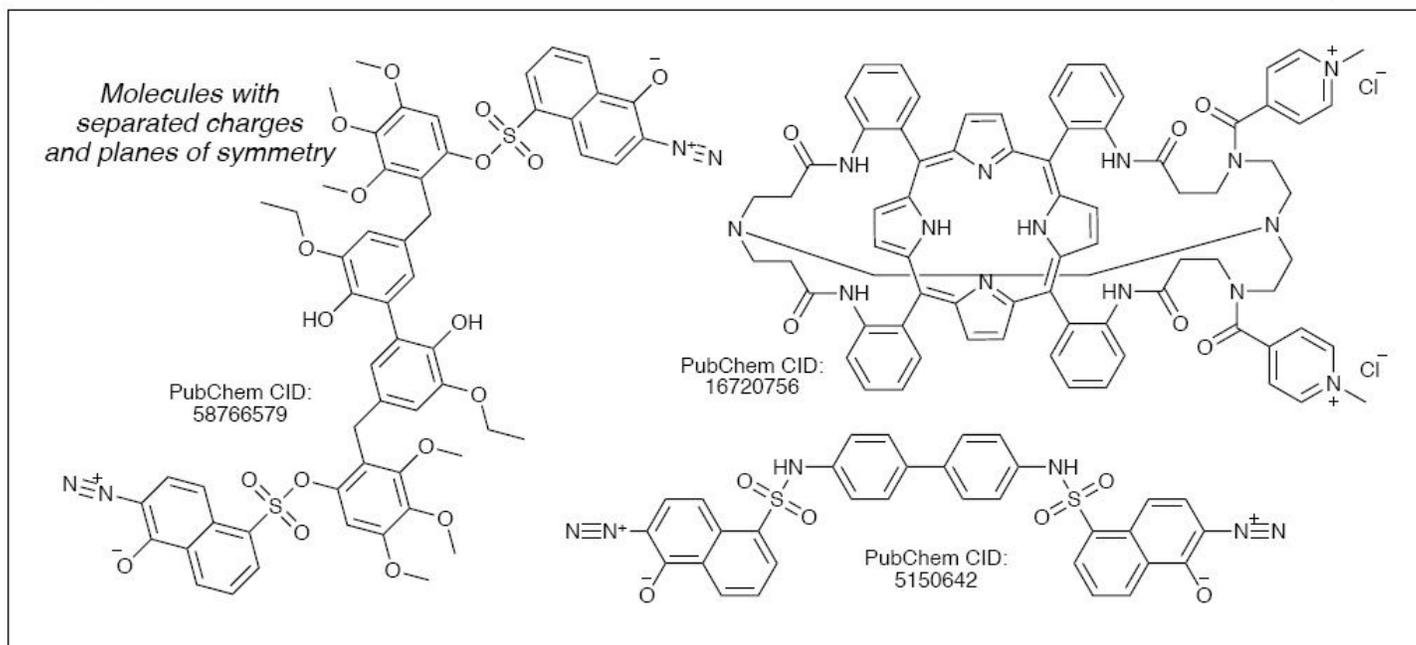
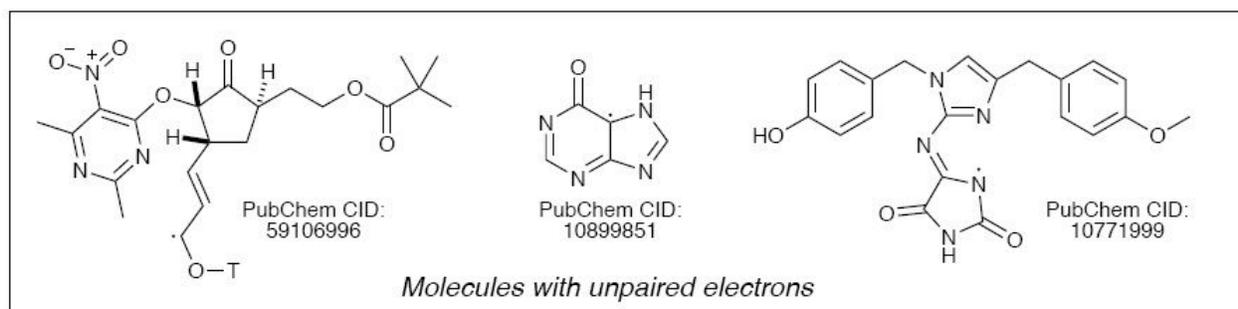
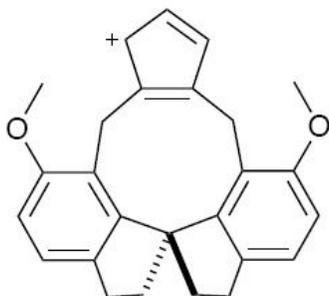


Figure 7

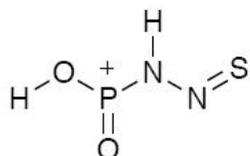
All molecules with numbering instability that lack "/p" and "/q" layers



PubChem CID: 134958543

InChI v1.05: InChI=1S/C26H25O2/  
c1-27-22-8-6-16-10-12-26-13-11-17-7-9-23(28-2)21(25(17)26)15-19-5-3-4-18  
(19)14-20(22)24(16)26/h3-9H,10-15H2,1-2H3/q-1

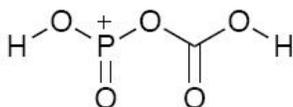
InChI v1.06: InChI=1S/C26H25O2/  
c1-27-22-8-6-16-10-12-26-13-11-17-7-9-23(28-2)21(25(17)26)15-19-5-3-4-18  
(19)14-20(22)24(16)26/h3-9H,10-15H2,1-2H3/q-1/t26-/m0/s1



PubChem CID: 154400441

InChI v1.05: InChI=1S/N2O2PS/c3-5(4)1-2-6/q-1/p+2

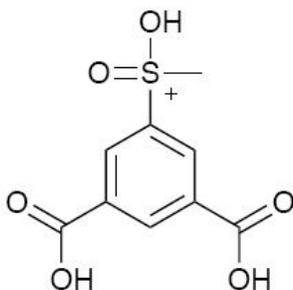
InChI v1.06: InChI=1S/HN2O2PS/c3-5(4)1-2-6/h(H-,1,3,4,6)/p+1



PubChem CID: 20681196

InChI v1.05: InChI=1S/CO5P/c2-1(3)6-7(4)5/q-1/p+2

InChI v1.06: InChI=1S/CHO5P/c2-1(3)6-7(4)5/h(H-,2,3,4,5)/p+1



PubChem CID: 147555804

InChI v1.05: InChI=1S/C9H7O6S/c1-16(14,15)7-3-5(8(10)11)2-6(4-7)9(12)13/  
h2-4H,1H3,(H-2,10,11,12,13,14,15)/q-1/p+2

InChI v1.06: InChI=1S/C9H8O6S/c1-16(14,15)7-3-5(8(10)11)2-6(4-7)9(12)13/  
h2-4H,1H3,(H2-,10,11,12,13,14,15)/p+1

## Figure 8

changes from v1.05 to v1.06