

## Universal Chemical Markup (UCM) - A new format for common chemical data

**Background:** We wish to introduce a new chemical format called UCM (Universal Chemical Markup). The format is based on XML (Extensible Markup Language) and its first version focuses on recording chemical structures and their properties. **Results:** UCM currently supports structures containing isotopes, ions and various types of bonding including delocalized bonds. Properties can be expressed by combining UCM with UnitsML (Units Markup Language). Using UnitsML one defines quantities with scientific units, and then refers to them in UCM when recording property values. Users can also add literature references with BibTeXML (BibTeX Markup Language) and annotate the recorded data using plain text or XHTML (Extensible Hypertext Markup Language) descriptions. In contrast to presently available general-purpose chemical formats, UCM offers built-in validation, which combines both grammar and pattern-based XML schema languages. Thus, all recorded data can be precisely validated by UCM schemas in standard XML validators. **Conclusions:** We developed the structure for UCM from scratch on the basis of an analysis described in our previous article. Starting from scratch allowed us to integrate BibTeXML, UnitsML and XHTML as well as chemical line notations and identifiers into UCM. It also helped us to avoid unnecessary redundant parts and create the implementation that aims to minimize ambiguity and is designed to be easily extensible in the future.

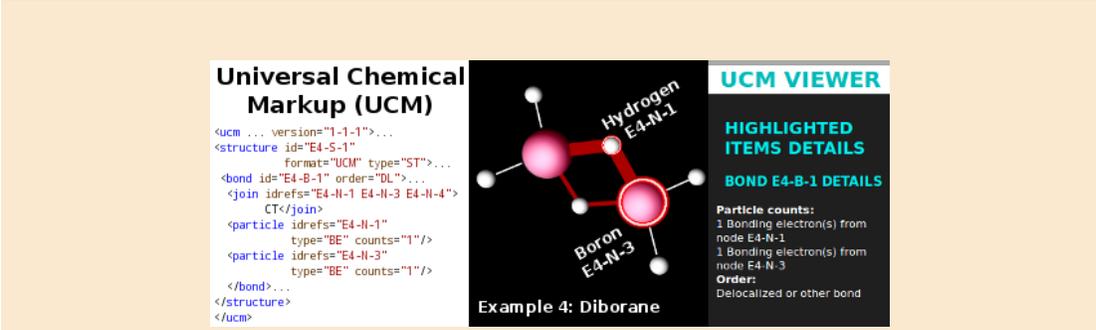
# Universal Chemical Markup (UCM) - A new format for common chemical data

Jan Mokry<sup>1</sup> and Miloslav Nic<sup>2</sup>

<sup>1</sup>Department of Inorganic Chemistry, University of Chemistry and Technology Prague, Technická 5, 166 28, Prague 6, Czech Republic

<sup>2</sup>Department of Software Engineering, Czech Technical University in Prague, Thákurova 9, 160 00, Prague 6, Czech Republic

## ABSTRACT



The abstract graphic is divided into three main sections. On the left, under the heading 'Universal Chemical Markup (UCM)', there is a snippet of XML code: 

```
<ucm ... version="1-1-1">...
<structure id="E4-S-1">...
  format="UCM" type="ST">...
  <bond id="E4-B-1" order="DL">...
  <join idrefs="E4-N-1 E4-N-3 E4-N-4">
    CT</join>
  <particle idrefs="E4-N-1"
    type="BE" counts="1"/>
  <particle idrefs="E4-N-3"
    type="BE" counts="1"/>
  </bond>...
</structure>
</ucm>
```

 In the center, a 3D ball-and-stick model of a diborane molecule (B<sub>2</sub>H<sub>4</sub>) is shown, with two boron atoms (pink) and four hydrogen atoms (white). Labels point to 'Hydrogen E4-N-1' and 'Boron E4-N-3'. Below the model is the text 'Example 4: Diborane'. On the right, a dark box titled 'UCM VIEWER' contains 'HIGHLIGHTED ITEMS DETAILS' for 'BOND E4-B-1 DETAILS'. It lists 'Particle counts: 1 Bonding electron(s) from node E4-N-1, 1 Bonding electron(s) from node E4-N-3' and 'Orders: Delocalized or other bond'.

**Background**  
We wish to introduce a new chemical format called UCM (Universal Chemical Markup). The format is based on XML (Extensible Markup Language) and its first version focuses on recording chemical structures and their properties.

**Results**  
UCM currently supports structures containing isotopes, ions and various types of bonding including delocalized bonds. Properties can be expressed by combining UCM with UnitsML (Units Markup Language). Using UnitsML one defines quantities with scientific units, and then refers to them in UCM when recording property values. Users can also add literature references with BibTeXML (BibTeX Markup Language) and annotate the recorded data using plain text or XHTML (Extensible Hypertext Markup Language) descriptions. In contrast to presently available general-purpose chemical formats, UCM offers built-in validation, which combines both grammar and pattern-based XML schema languages. Thus, all recorded data can be precisely validated by UCM schemas in standard XML validators.

**Conclusions**  
We developed the structure for UCM from scratch on the basis of an analysis described in our previous article. Starting from scratch allowed us to integrate BibTeXML, UnitsML and XHTML as well as chemical line notations and identifiers into UCM. It also helped us to avoid unnecessary redundant parts and create the implementation that aims to minimize ambiguity and is designed to be easily extensible in the future.

Keywords: Universal Chemical Markup, UCM, UCM XML structure, UCM built-in validation, UCM examples, UCM VIEWER, recording chemical structures with properties, combining XML formats, combining XML schema languages

## 1 BACKGROUND

- 2 Although many chemical formats currently exist, most are tailored to specific areas of chemistry. Relatively
- 3 few formats provide the general-purpose functionality, which is not limited to specific software or a
- 4 specific area of chemistry, and supports the recording of common chemical data (i.e. structures, reactions
- 5 and properties). This is apparent from our previous analysis of the widely used formats for common

1 chemical data.<sup>1</sup> In that analysis the strengths and weaknesses of these formats were examined. Acquired  
2 knowledge were further utilized when we designed the XML structure for UCM using the most suitable  
3 concepts from explored formats.

4 We decided to develop UCM, because our earlier analysis suggests no widely used general-purpose  
5 chemical format offers future extensibility together with the effective processing and precise built-  
6 in validation of chemical data.<sup>1</sup> The analysis results also confirmed XML technology has potentially  
7 significant benefits for chemical formats.<sup>1</sup> Thus, the development of UCM started from scratch with  
8 the aim of creating a new general-purpose XML format for common chemical data. This as well as the  
9 strengths and weaknesses of analyzed formats is described more thoroughly in our previous article.<sup>1</sup>

10 Before we explain the advantages of UCM, it is first necessary to describe at least some issues we  
11 have found in the CTfile (Chemical Table File), CDXML (ChemDraw Exchange Markup Language) and  
12 CML (Chemical Markup Language) formats. These formats, despite their weaknesses, represent the  
13 existing XML and non-XML solutions with rich general-purpose functionality for common chemical data,  
14 as can be seen from the benefits presented in our analysis.<sup>1</sup>

15 The popular CTfile formats can together record data about chemical structures, reactions, and proper-  
16 ties,<sup>2,3,4</sup> but the built-in validation is missing in all original non-XML CTfile formats (i.e. Molfile, RGfile,  
17 Rxnfile, SDfile and RDfile).<sup>1</sup> In the case of XDfile (XML Data File), which is the only XML-based CTfile  
18 format, its schema can be used by standard XML validators to check the format's basic structure.<sup>1</sup> But  
19 because the XDfile schema provides just simple grammar-based validation, and most chemical data in  
20 XDfile are stored in the embedded non-XML formats (e.g. Molfile, Rxnfile, etc.), the precise validation  
21 of chemical data is not achieved.<sup>1</sup> Furthermore, our analysis suggests that implementing the precise  
22 validation for CTfile formats (including XDfile) may require additional effort, as the original CTfile syntax  
23 is not based on XML or similar standard, which offers the validation infrastructure automatically.<sup>1</sup> Also  
24 in the non-XML CTfile formats the stored data items are less readable, as they often lack descriptions and  
25 there are only limited annotation possibilities.<sup>1</sup>

26 CDXML closely follows binary CDX (ChemDraw Exchange) specifications<sup>5</sup> and is this causes at  
27 least two potential problems. The first, similarly as in CDX, is the mixing of chemical information  
28 with embedded binary objects and data describing content visualization. It makes the XML structure of  
29 CDXML more difficult to understand and implement.<sup>1</sup> The second is that the dependency on changes  
30 in CDX specifications negatively affects the future extensibility of CDXML.<sup>1</sup> Another issue is in the  
31 validation functionality. Because CDXML uses DTD (Document Type Definition), a schema language  
32 with limited expressive power,<sup>6</sup> its basic XML structure can be validated easily, but chemical information  
33 is not validated precisely.<sup>1</sup> In addition DTD does not support XML namespaces.<sup>6</sup> Thus, CDXML will  
34 need to be redefined with a newer schema language to add a unique namespace, so that the format can be  
35 combined directly with other XML formats.<sup>1</sup>

36 CML has undergone a long evolution<sup>7</sup> and provides rich and flexible functionality focused on  
37 extensibility,<sup>8,9</sup> but also has its issues. Starting with version 2 it was redefined using XSD (W3C XML  
38 Schema Definition Language).<sup>10,11</sup> The flexibility of the format was increased by defining minimal  
39 restrictions and lax validation for many parts (see the CML schema 2.4 or 3<sup>12</sup>). In CML schema 3 any  
40 mandatory tree structures were abandoned to offer even more flexibility.<sup>13</sup> However, the fact that more than  
41 one way of recording the particular chemical information usually exists makes the format confusing.<sup>1</sup> The  
42 validation of such a format then becomes complicated, because the schema uses sometimes deliberately  
43 fuzzy concepts and enables all CML elements to contain each other.<sup>1</sup> To address these problems the authors  
44 of CML implemented a mechanism for denoting standard conventions for various CML elements.<sup>8,10,14</sup>  
45 The main issue was that until version 3 no official list of conventions was composed and CML users,  
46 including developers of chemical software, had to use their own conventions.<sup>1</sup> While with version 3  
47 CML users can still create and use custom conventions, the documentation now lists some official  
48 conventions aimed at adding additional rules and bringing more order to the format.<sup>1</sup> Part of this effort  
49 is also the CMLLite validator service intended for flexible validation based on conventions.<sup>15</sup> But the  
50 listed conventions appear to be mostly unfinished. All are marked as draft recommendations and some are  
51 missing (e.g. the convention dealing with CML support for recording crystallographic information or  
52 reactions etc.).<sup>1</sup> Until the necessary official conventions are developed and implemented at least in the  
53 CMLLite validator service, CML will remain a very variable format that is difficult to process reliably.<sup>1</sup>

1 With the exception of CML, the general-purpose formats we explored also cannot directly record  
 2 electrons participating in chemical bonds and do not properly enforce associating scientific units with the  
 3 values of properties.<sup>1</sup> This leads to a less precise description of properties and structures, especially the  
 4 structures with more complicated delocalized bonds.

5 UCM, the new XML format we propose, aims to avoid the issues identified in previously analyzed  
 6 XML and non-XML formats.<sup>1</sup> Among the main benefits of UCM is the precise validation, which uses  
 7 software from the standard XML tool chain and combines both grammar and pattern-based XML schema  
 8 languages. The grammar-based UCM schema validates the format's tree structure, while the pattern-based  
 9 schema validates constraints between UCM attributes and elements. These constraints can even check  
 10 UCM data for chemical mistakes like a carbon atom participating in five different single bonds, or a  
 11 formal charge of a structure not corresponding with the number of protons and electrons in it. The first  
 12 version of UCM focuses on recording chemical structures and their properties. Because UCM has a  
 13 unique namespace it can be directly combined with other XML formats to include complex scientific data  
 14 from various domains in a single file.

## 15 IMPLEMENTATION AND OVERALL DESIGN OF UCM

16 In the iterative development of the XML structure for UCM we used and improved the most suitable  
 17 concepts found in chemical XML formats from our earlier analysis.<sup>1</sup> This illustrates our aim of reusing  
 18 the current knowledge and available solutions where possible, and improving upon these when necessary.  
 19 Thus, the overall design philosophy of UCM is not just about creating a new format without the previously  
 20 described issues. It is also about integrating the existing knowledge and solutions into UCM to avoid  
 21 reinventing the wheel. Consequently, we designed UCM in a way that ensures various XML and non-XML  
 22 formats can be integrated in it and provide additional specialized functionality.

23 The resulting XML structure for UCM 1-1-1 is in tree structure schemes 1 and 2, while our previous  
 24 article describes how UCM implements various concepts compared to other XML formats we analyzed.<sup>1</sup>  
 25 Tree structure schemes 1 and 2 follow the complete format structure defined by UCM validation schemas,  
 26 which are available in additional file 1. We used very simple syntax to make UCM tree structure easily  
 27 readable:

- 28 • `element` – Denotes an UCM element.
- 29 • `@attribute` – Denotes an UCM attribute.
- 30 • `[CONTEXT]` – Specifies the context of an UCM element (e.g. `node [C1]` means the *node*  
 31 element that uses its first context).
- 32 • `[CONTEXT: XPATH]` – Defines the context of an UCM element (e.g. `define [C1:`  
 33 `@format = 'UCM']` means the first context for the *define* element is the *define* element with  
 34 the value "UCM" in its *format* attribute). The XPath (XML Path Language) expression only further  
 35 describes the context, which is defined by the attributes and child elements enabled for the given  
 36 UCM element.
- 37 • `(ATTRIBUTES)` – Specifies the enabled attributes of an UCM element (e.g. `point (@id,`  
 38 `@x, @y, @z)` means the *point* element with the *id*, *x*, *y* and *z* attributes).
- 39 • Quantifiers "?", "\*" and "+" are used to express 0 or 1, 0 or more, and 1 or more respectively.
- 40 • Keywords "OR", "IF" and "IF NOT" have their literal meaning.
- 41 • Element contents are indented by four spaces.
- 42 • Ellipsis means the attributes and contents of the element are in its definition.

43 As can be seen from tree structure schemes 1 and 2, some UCM elements may be utilized in more  
 44 than one context. Contexts are defined in UCM Schematron schema by additional restrictions placed on  
 45 the particular UCM element when it is used at the specific position in UCM tree structure, or when it has  
 46 some specific attribute or a certain value of some attribute. Take for example the UCM *node* element.  
 47 If it occurs inside a *define* element the following restrictions from UCM Schematron schema apply to

1 the *node* element: the mandatory *id* attribute is enabled, but other attributes must be omitted; the *stereo*  
2 child element must be omitted, but there must be at least one *particle* child element (for *description* and  
3 *property* child elements occurrence remains as defined in the main UCM schema). The restrictions define  
4 the first context for UCM *node* element. Source code snippet 3 demonstrates how that context is used for  
5 storing reusable UCM node definitions. This way UCM has relatively few elements with flexible and easy  
6 to learn functionality, but at the same time precise validation is ensured by exactly specifying all useful  
7 contexts for these elements.

## 8 Integrated XML formats

9 Using XML namespaces we integrated BibTeXML, UnitsML and XHTML into UCM. Unique XML  
10 namespaces are essential for combining two or more XML formats together. The reason is that unique  
11 XML namespaces make it possible to distinguish the attributes and elements from different formats even if  
12 these attributes or elements have the same local names.<sup>16,17</sup> Moreover, XML namespaces are also used by  
13 NVDL (Namespace-based Validation Dispatching Language),<sup>18,19</sup> which we use for the [Implementation  
14 of data validation in UCM](#).

15 The purpose of XHTML integration is simply to enhance the annotation capabilities of UCM *de-*  
16 *scription* elements. XHTML enables hyperlinks and other useful formatting features in annotations, as is  
17 demonstrated in our [Results with examples of UCM usage](#). In addition it ensures that all markup inside  
18 the *description* elements is compatible with web browsers without any additional transformations.

19 With BibTeXML integrated in UCM *define* elements one can define literature references for the  
20 *description* elements. Both the default and extended version of BibTeXML is supported in UCM. Default  
21 BibTeXML 1.0 offers the functionality of the original BibTeX format, which is widely adopted by the  
22 scientific community and is used with L<sup>A</sup>T<sub>E</sub>X document preparation system.<sup>20</sup> BibTeXML 1.0 Extended  
23 adds functionality for recording additional data items more precisely.<sup>20</sup> This is useful for modern literature  
24 references with Digital Object Identifiers and other data items, as shown in source code snippet 1. Because  
25 software tools for automatic conversion of BibTeX bibliographies into BibTeXML are freely available  
26 at the format website,<sup>20</sup> users can easily add their literature references to UCM annotations from most  
27 reference managers.

28 Similarly, UnitsML is also integrated in UCM *define* elements, as it enables defining various quantities  
29 with their scientific units. Source code snippet 2 gives an example of how it works. At first necessary  
30 scientific units are defined using UnitsML *Unit* elements, which are encapsulated in the *UnitSet* element.  
31 Then, quantities are defined inside the *QuantitySet* element by UnitsML *Quantity* elements. Each *Quantity*  
32 element usually refers to the appropriate *Unit* element that defines the scientific unit (see the UnitsML  
33 *UnitReference* element in the *Quantity* element). Because UnitsML offers a set of root units, it is possible  
34 to easily define most scientific units in a precise way.<sup>21</sup> Further information about UnitsML are available  
35 at the format website.<sup>21</sup> We especially recommend the format documentation, as it can be easily navigated  
36 thanks to cross-references and clear formatting.

## 37 Integrated non-XML formats

38 UCM also integrates non-XML formats for chemical line notations and identifiers. As apparent from tree  
39 structure scheme 1, these formats are integrated in the UCM *structure* elements. The functionality of  
40 integrated chemical line notations and identifiers nicely complements the capabilities of UCM markup.  
41 The *structure* elements that use UCM format can be validated in detail and are intended for precisely  
42 recorded structures with coordinates, properties or annotations. On the other hand sometimes one may  
43 need to store just a concise chemical identifier or query, which is usually without coordinates and more or  
44 less ambiguous. In such cases the *structure* elements can use one of the integrated formats to record the  
45 chemical structure or query using:

- 46 • Preferred IUPAC Name, General IUPAC Name, Chemical Abstracts Index Name,
- 47 • CAS RN (Chemical Abstracts Service Registry Number), Reaxys Registry Number, ChemSpider  
48 ID Number, PubChem Compound ID Number, PubChem Substance ID Number,
- 49 • InChI (International Chemical Identifier), InChIKey, Standard InChI, Standard InChIKey,
- 50 • SMILES (Simplified Molecular Input Line Entry System), SMILES Arbitrary Target Specification  
51 or SYBYL Line Notation.

1 Examples showing how these integrated formats may be used are presented in our [Results with](#)  
2 [examples of UCM usage](#). Further information are also available in UCM documentation provided in  
3 additional file 2.

#### 4 IMPLEMENTATION OF DATA VALIDATION IN UCM

5 During the development of UCM we have set the precise built-in validation as one of the main design  
6 goals, because current general-purpose chemical formats often have only limited built-in validation  
7 capabilities.<sup>1</sup> As described in our previous article, XML technology offers significant benefits one of  
8 which is the available validation infrastructure.<sup>1</sup> Virtually any formally defined XML format has a schema  
9 that provides at least some basic validation capabilities. In contrast to most XML formats, we decided  
10 to further enhance the validation capabilities of the main UCM schema. Thus, the validation in UCM  
11 employs a combination of XML schema languages.

12 The main schema of UCM uses RELAX NG (Regular Language for XML Next Generation) compact  
13 syntax. We also offer its versions translated with Trang (a converter for common XML schema  
14 languages<sup>22</sup>) into RELAX NG XML syntax and XSD. This ensures very high compatibility especially  
15 when the main schema is used independently with various XML validators. The main schema provides  
16 relatively precise grammar-based validation of UCM tree structure. The values of all UCM attributes  
17 and elements are checked for correctness. In other words each value must conform to its data type and  
18 a possible set of additional restrictions. An example is the *counts* attribute value, which must be one  
19 or more whitespace separated non-negative integer numbers. For UCM elements the main schema also  
20 checks whether they have the correct attributes. Additionally, for UCM elements that should contain  
21 other elements the correct sequence and occurrence of child elements is checked. Such validation could  
22 be considered sufficient, as many current XML formats, both chemical (e.g. CDXML, XDfile, etc.)  
23 and non-chemical (e.g. BibTeXML, UnitsML, XHTML, etc.), are defined and validated by a single  
24 grammar-based schema.<sup>1,20,21,23</sup> However, validation using a single schema written in one of the widely  
25 popular grammar-based XML schema languages (e.g. RELAX NG, XSD, etc.) has at least two limitations.  
26 The first limitation becomes apparent when one considers validation of XML files combining two or  
27 more XML formats with unique namespaces together. The grammar-based schema for one XML format  
28 cannot check the content from other XML formats unless it somehow includes their schemas too (such  
29 combinations are for example described in Mathematical Markup Language specification<sup>24</sup>). The second  
30 limitation is that the grammar-based schema cannot express validation constraints between two or more  
31 attributes, elements or values (for these constraints Schematron and its pattern-based approach is usually  
32 utilized<sup>25</sup>).

33 To validate the content from UCM namespace as well as the content from BibTeXML, UnitsML  
34 and XHTML namespaces we added a NVDL (Namespace-based Validation Dispatching Language)  
35 schema to UCM. NVDL makes it possible to load the appropriate schema for the content from each  
36 namespace.<sup>18,19</sup> This addresses the first limitation we mentioned previously for the validation that relies on  
37 a single grammar-based schema. In addition, NVDL enables the content from one namespace to be easily  
38 validated by more than one schema.<sup>19,26</sup> We plan to utilize NVDL for validating the content from UCM  
39 namespace by both RELAX NG and Schematron UCM schemas as soon as there is mature support for ISO  
40 Schematron in NVDL-capable validators. NVDL-capable validators are widely available<sup>19</sup> and except for  
41 the limited ISO Schematron support we observed no issues with common XML schema languages (e.g.  
42 RELAX NG, XSD or DTD).

43 With the Schematron schema we eventually achieved a very precise and extensible validation of  
44 more complex constraints between two or more UCM attributes, elements or values. The usage of  
45 Schematron obviously solved the second limitation of the validation based on just a single grammar-  
46 based schema. While the direct implementations of Schematron exist, the official Schematron website  
47 also offers the reference "Skeleton" implementation based on XSLT (Extensible Stylesheet Language  
48 Transformations).<sup>27,28</sup> We found the "Skeleton" implementation to be most mature especially since the  
49 UCM Schematron schema uses the current ISO version of Schematron. This way UCM Schematron  
50 validation may be used separately (i.e. we chose not to integrate it in RELAX NG schema though it is  
51 possible<sup>29</sup>) and is widely compatible as the only required software tool is the XSLT 2.0 processor.

52 Among other things the Schematron schema for UCM can validate various chemical rules, which to  
53 the best of our knowledge cannot be checked by built-in validation in current general-purpose chemical  
54 formats.<sup>1</sup> Two main examples from UCM 1-1-1 Schematron schema include: the validation of formal

1 charges on chemical structures and chemical nodes (i.e. nodes from a chemical graph of the structure),  
2 and the validation of bonding electrons (i.e. electrons to be used in UCM *bond* elements). To enable  
3 the easier implementation of these validation tasks we designed UCM around the principle of precisely  
4 defining all chemical nodes that represent the monoatomic particles in chemical structures. Thus, each  
5 monoatomic particle (e.g. an atom or a monoatomic ion) in UCM is defined as composed of protons,  
6 neutrons and electrons. In other words each UCM node is defined by a *node* element containing *particle*  
7 child elements, which specify the number of protons, neutrons and electrons. For an example see the  
8 definitions of hydrogen cation, deuterium and carbon nodes in source code snippet 3. Note how isotopic  
9 composition can be recorded with *fractions* attributes and how the number of bonding electrons can be  
10 distinguished by the "BE" value of the *type* attribute on UCM *particle* elements.

11 Using the information about the number of protons and electrons for each UCM *node* element our  
12 Schematron schema checks if the *charge* attributes on *node* and *structure* elements have the correct value.  
13 Of course in practice electrons from some *node* element may be shared with other *node* elements in  
14 a way that leads to decimal values of the *charge* attribute. An example is the structure of dihydrogen  
15 cation in source code snippet 6 or the structure of ozone in source code snippet 10. The validation then  
16 becomes more complicated, but it may be still implemented in Schematron, as one can use additional  
17 XSLT functions in Schematron validation patterns. Source code snippet 4 demonstrates this with our  
18 XSLT function and validation pattern for the *charge* attribute on an UCM *node* element. The validation  
19 pattern calculates the *charge* attribute value simply by subtracting the number of electrons (i.e. the variable  
20 "NEGATIVE-CHARGE") from the number of protons (i.e. the variable "POSITIVE-CHARGE") for  
21 the given *node* element. In order to take into account the possible shared electrons the XSLT function  
22 "If:COUNT-NODE-NEGATIVE-CHARGE" is used. This function subtracts the number of shared  
23 electrons from the total electron count obtained from the *node* element definition. Then only the correct  
24 fractions of shared electrons are added back to get the actual negative charge for the particular *node*  
25 element. The correct fractions of shared electrons are obtained from the values of the *fractions* attributes  
26 on relevant UCM *share* elements.

27 Similarly the Schematron schema also verifies whether each UCM *node* element has enough bonding  
28 electrons for all its chemical bonds. Source code snippet 5 contains the core parts of that validation  
29 procedure. At first UCM node definitions are searched for the number of electrons the *node* element  
30 should provide to *bond* elements. Then the XSLT function "If:COUNT-NODE-BONDING-ELECTRONS"  
31 calculates the number of electrons required by all UCM *bond* elements that refer to the particular *node*  
32 element. If both numbers are not equal the error occurs and the particular chemical node either participates  
33 in too many bonds (i.e. it has too few bonding electrons) or does not participate in all its bonds (i.e. it has  
34 too many bonding electrons).

35 In practice UCM validation should help users to determine if the given structure seems to be chemically  
36 correct or incorrect. Of course UCM 1-1-1 schemas cannot detect all possible chemical errors, but in  
37 future versions validation can be further enhanced and extra rules may be added. When compared to  
38 current general-purpose formats such as CDX, CDXML, CML and CTfile formats (including XDfile) we  
39 believe UCM validation brings significant improvements.

## 40 RESULTS WITH EXAMPLES OF UCM USAGE

41 To test UCM during the development we prepared examples, which use the key format functionality in  
42 practice. Complete and quite easily readable UCM files with all examples are in additional file 1. For  
43 readers of the following sections we further explain the selected examples to better illustrate how UCM  
44 can record various chemical structures and properties. In each example the chemical graph of the given  
45 structure was manually rewritten to UCM in a plain text editor. The chemical graph was constructed  
46 using information about the structure in scientific literature and chemical databases. Most examples  
47 also include 3-dimensional coordinates from available online sources and chemical databases. For some  
48 examples 3-dimensional coordinates were obtained from a skeleton structure we created with Avogadro  
49 (an open source chemical editor<sup>30,31</sup>) and optimized in Nwchem (an open source computational chemistry  
50 software<sup>32,33</sup>). We mention the original source of 3-dimensional coordinates and other information about  
51 the structure in each example, mainly to demonstrate the annotation functionality of UCM.

1 A simple UCM VIEWER, which we developed using JavaScript, XHTML and CSS (Cascading  
2 Style Sheets), helps to visualize the UCM examples in the following sections. In its first version UCM  
3 VIEWER supports the key functionality of UCM 1-1-1. To run UCM VIEWER one only needs a  
4 modern web browser with a built-in XSLT processor and support for JavaScript, or other ECMAScript  
5 implementation. UCM VIEWER utilizes JavaScript for parsing and rendering the content from UCM  
6 files. We use two open source JavaScript libraries in UCM VIEWER, three.js and jQuery. With its  
7 easy-to-use API (Application Programming Interface) compatible across a multitude of web browsers,  
8 jQuery is a well known library that simplifies various tasks in JavaScript (e.g. document traversal and  
9 manipulation, event handling, etc.).<sup>34</sup> However, equally important for UCM VIEWER is also the three.js  
10 project, which aims to create a lightweight library for 3-dimensional rendering with a very low level  
11 of complexity.<sup>35</sup> We have found the API of the three.js library intuitive and highly usable, while both  
12 quality and performance were sufficient for our use case even with the basic canvas renderer (the library  
13 provides WebGL and other renderers too<sup>35</sup>). Overall we think the source code of UCM VIEWER,  
14 available from <http://www.universalchemicalmarkup.org> under an open source license, demonstrates how  
15 the software support for UCM can be implemented relatively easily.

### 16 Examples of simple structures and structure fragments

17 In the first example we included small structures and structure fragments without 3-dimensional coordi-  
18 nates just to introduce the basic overall structure of an UCM file with chemical data. The complete  
19 example contains: dihydrogen cation, heavy water, ozone resonance hybrid with its resonance structures,  
20 and fragments from sodium chloride and hydrogen fluoride structures. The *ucm* root element in source  
21 code snippet 6 is used as a container for other UCM elements and it also specifies the namespaces and  
22 UCM version. Mandatory *version* attribute is intended to prevent future problems with versioning. Source  
23 code snippet 6 shows only the first two UCM *structure* elements with data about dihydrogen cation and  
24 heavy water. The *node* elements store chemical nodes, which represent the monoatomic particles (e.g.  
25 atoms or monoatomic ions) in the structures, while *bond* elements record bonds. Each *structure*, *node*  
26 or *bond* element has a mandatory *id* attribute to enable reliable cross-referencing using *idrefs* attributes.  
27 One can see how UCM enables recording and utilizing information about subatomic particles such as  
28 protons, neutrons and electrons when chemical bonds and nodes are described. The *idrefs* attribute on a  
29 *node* element provides cross-reference to the reusable node definition that specifies the number of protons,  
30 neutrons and electrons (see source code snippet 3). On a *bond* element the *idrefs* attribute provides  
31 cross-reference to *node* elements participating in the bond. In addition it is possible to describe electrons  
32 and how they are shared inside a bond using UCM *particle* and *share* elements. This is demonstrated  
33 inside the first *bond* element in source code snippet 6 and we explain it further when describing [Examples  
34 of more complex structures with delocalized bonds](#).

35 Before we expound the remaining examples, it is necessary to mention how XInclude (XML Inclu-  
36 sions) mechanism may be utilized in UCM and how this affects the basic structure of UCM files. XInclude  
37 is a mechanism for general purpose inclusion accomplished by merging a number of XML information  
38 sets into a single composite infoset.<sup>36</sup> As can be seen in source code snippet 6, after XInclude namespace  
39 prefix is defined (on the *ucm* root element) it is possible to use the *include* elements from XInclude  
40 namespace to incorporate the content from external XML files. XInclude processing occurs at a low  
41 level, often by a generic XInclude processor which makes the resulting information set available to higher  
42 level applications.<sup>36</sup> For UCM *define* elements such mechanism is especially advantageous, because one  
43 may reuse precise UCM, BibTeXML and UnitsML definitions simply by including them in UCM files  
44 that require them. Thus, usually the main UCM file with chemical data will reference other UCM files  
45 with definitions. We assume some definitions (e.g. UCM node definitions or UnitsML quantity and units  
46 definitions) could be offered as a part of future UCM versions, and so our UCM examples already use  
47 separate files for each group of similar definitions.

48 Now let us continue with the second example. Although optional in UCM, the coordinates of the *node*  
49 elements positions in 3-dimensional space are essential for an accurate description of the given structure  
50 and its precise visualization. Figure 1 and source code snippet 7 with the second example demonstrate this  
51 basic UCM functionality on a simple structure of urea. All positions of the *node* elements, which represent  
52 atoms in the urea structure, are stored using attributes *x*, *y*, *z* denoting coordinates in 3-dimensional space  
53 expressed in nanometers.

1 The second example also contains chemical identifiers for the urea structure and an annotation with  
2 XHTML markup inside the UCM *description* element. As depicted in figure 1, XHTML enables the user  
3 to easily include links to additional online resources and to enhance the formatting of annotation text. The  
4 main part of figure 1 shows how UCM VIEWER displays the urea structure and its chemical identifiers  
5 stored in UCM *structure* elements that have the "STID" value in their *type* attributes (see source code  
6 snippet 7).

### 7 Example of structure with properties and literature references

8 To record literature references and measured or calculated properties UCM *define* elements support  
9 BibTeXML and UnitsML, as we already explained (see source code snippets 1 and 2). When values of  
10 properties are recorded using UCM *property* elements, as in source code snippet 8, a *quantity* attribute  
11 on the given *property* element references UnitsML quantity via its *xml:id* attribute. Similarly the *litrefs*  
12 attribute on the UCM *description* element references BibTeXML literature reference via its *id* attribute.  
13 The figure 2 depicts UCM VIEWER displaying properties loaded from *property* elements together with  
14 annotations from *description* elements, which also include literature references. In source code snippet 8  
15 one can see different uses of the UCM *property* element. The *property* element with the "CN" value in its  
16 *type* attribute can describe a condition for the parent *property* element. If the *type* attribute has the value  
17 "ER", the *property* element describes errors in the *values* of the parent *property* element. This enables  
18 precise recording of errors in measured or calculated values using well defined statistical quantities like  
19 standard deviation, as can be seen from source code snippet 8. Although our simple example does not  
20 include it, *property* elements may be also utilized inside UCM *bond*, *node*, *particle* and *point* elements.  
21 Thus, it is possible to describe various properties related to these elements (e.g. bond length, bond  
22 dissociation energy, atomic weight, atomic radius, particle spin, etc.) as easily as in the case of UCM  
23 *structure* elements.

### 24 Examples of more complex structures with delocalized bonds

25 Previous examples included structures where chemical bonds could be described quite simply. However,  
26 there are many structures with more complex bonding, which current general-purpose chemical formats  
27 fail to describe in detail. It is because formats such as CDX, CDXML, CML and CTfile formats (including  
28 XDfile) do not support bond types that enforce specifying and validating electrons for more complicated  
29 bonds (e.g. aromatic and other delocalized or otherwise special bonds).<sup>1,37,12,3</sup> In UCM, the *bond*, *join*  
30 and *particle* elements can be used together to record both simple and complex bonding. The *particle*  
31 elements enable the usage of subatomic particles in UCM. Thus, it is possible to record and validate  
32 electrons participating in the given bond.

33 Source code snippet 9 and figure 3 present the example of diborane with two 3-center-2-electron bonds  
34 between the bridging hydrogen and boron atoms. These bonds are stored by two UCM *bond* elements.  
35 Each contains the *join* element with information about how to join the *node* elements that participate  
36 in the bond. The "CT" value inside both *join* elements stands for the centered interpretation of their  
37 *idrefs* attribute values. It means that the first id reference refers to the *node* element which is bonded,  
38 by the bond expressed in the parent *bond* element, to all *node* elements specified by the remaining id  
39 references. Therefore, in the case of diborane, the first id reference (in the *idrefs* attribute of the *join*  
40 element) references the *node* element that represents the bridging hydrogen atom, while the remaining  
41 id references refer to *node* elements representing the boron atoms. The id reference inside each *idrefs*  
42 attribute on the *particle* elements points to the *node* element which provided the electron. This enables  
43 the precise validation of bonding electrons, as we already described (see source code snippet 5).

44 Some advanced usage of UCM *bond* element can be seen also in the *structure* elements that describe  
45 ozone. Source code snippet 10 shows the relevant part of the first example here. The *structure* element  
46 with the "E1-S-3" value in its *id* attribute stores the resonance hybrid of classical ozone Lewis structures,  
47 which are presented in the next two *structure* elements. The most interesting is the first *bond* element  
48 (with the "E1-B-3-1" *id* attribute value) that records the 3-center-2-electron bond formed by two electrons  
49 from the pi bonding orbital covering all three oxygen atoms. Notice how *share* elements precisely describe  
50 sharing of electrons between the *node* elements and how it leads to non-zero formal charge values on  
51 these *node* elements. Sharing ratios for the *node* elements referenced in the *idrefs* attributes on each  
52 *share* element are stored in the *fractions* attributes. Of course the ozone structure is still a topic of new  
53 research. For example one of the recent papers suggests the classical Lewis structures contribute 82% to  
54 the resonance hybrid of ozone, while remaining 18% is contributed by the biradical resonance structure<sup>38</sup>

1 (see the last *structure* element in source code snippet 10). In future UCM versions we plan to add support  
2 for mixed *structure* elements capable of recording such resonance hybrids as well as other chemical  
3 substances and mixtures.

4 Another example is the structure of ferrocene in figure 4. Interesting bonds in this example include the  
5 aromatic bonds in cyclopentadienyl rings and bonds between the central iron node and cyclopentadienyl  
6 rings. Details of these bonds are in source code snippet 11. The *bond* element with the "E5-B-1" value  
7 in its *id* attribute represents the aromatic bonding in the first cyclopentadienyl ring. The "CC" value  
8 of the *join* element stands for the cyclic interpretation of its *idrefs* attribute value. This means that all  
9 aromatic bonds in the first cyclopentadienyl ring are recorded in UCM using single *bond* element, without  
10 sacrificing any details. It can be still seen that fifteen electrons from carbon atoms and one extra electron  
11 are participating in the aromatic bonding. We could even distinguish the sigma and pi electrons from  
12 carbon atoms by adding one more *particle* element.

13 Because the UCM *structure* element can be nested, it is possible to encapsulate any useful substructure  
14 and then reference it later. This is demonstrated in source code snippet 11 by the *structure* element  
15 with the "E5-S-1-3C" value in its *id* attribute. All carbon atoms of the first cyclopentadienyl ring are  
16 encapsulated in that element. Thus, just one *particle* element may be used to reference these carbon atoms  
17 in the *bond* element representing the aromatic bonding.

18 Ferrocene example also shows a straightforward usage of *charge* attribute and sharing of electrons  
19 among bonds. The *charge* attribute value must correspond to the number of proton and electron particles  
20 used inside the given *node* or *structure* element. The sharing of electrons between bonds can be seen in  
21 the *bond* element with "E5-B-13" value in its *id* attribute. That *bond* element denotes the bond between  
22 the central iron *node* and the first cyclopentadienyl ring. As can be seen in source code snippet 11, three  
23 electrons in this *bond* are from the central iron *node*, but the remaining six are pi electrons shared with the  
24 *bond* element representing the aromatic bonding. Observe how the *bond* element with the *id* attribute  
25 value of "E5-B-13" connects the central iron *node* and the centroid of the first cyclopentadienyl ring  
26 denoted by the *point* element. When the example is rendered by UCM VIEWER it emphasizes the fact  
27 that both cyclopentadienyl rings rotate.<sup>39</sup>

28 Finally there is an example of trichloro(ethene)platinate(II) anion (commonly known as Zeise's salt  
29 anion). Here the interesting bond is between the central platinum atom and ethylene ligand. Figure 5  
30 depicts both our approaches at describing this bond, while source code snippet 12 shows the relevant  
31 UCM markup. In the first approach we use the 3-center-2-electron bond to connect the central platinum  
32 atom and ethylene ligand (see the UCM *structure* element with the "E6-S-1" value in its *id* attribute).  
33 As a result only a single bond remains between the ligand carbon atoms. On the other hand our second  
34 approach emphasizes that the order of the bond between ethylene ligand carbon atoms is not completely  
35 reduced to a single bond.<sup>40,41</sup> Therefore, the central platinum atom and ethylene ligand in the structure  
36 "E6-S-2" are connected using a delocalized bond, which shares an electron with the partial double bond  
37 between the ligand carbon atoms.

### 38 Examples of structures with stereochemical configuration

39 The following examples illustrate how to describe the stereochemistry in UCM using the *stereo* element.  
40 An example in figure 6 depicts serine amino acid with chirality centre on the highlighted carbon atom.  
41 The *stereo* element in the source code snippet 13 describes the chirality centre on the *node* element with  
42 the *id* attribute value "E7-N-5". Substituents are represented by *node* elements with the *id* attribute values  
43 "E7-N-4", "E7-N-6", "E7-N-7" and "E7-N-8". References in the *idrefs* attribute on the *stereo* element  
44 are ordered by descending priority of those substituents. Because the sense of rotation from the highest  
45 to lower priority substituents is counterclockwise, when the lowest priority substituent is pointed away  
46 from the observer (as in figure 6), the *sense* attribute has the "-" value. Substituents priority was assigned  
47 according to the Cahn-Ingold-Prelog system of priority rules, so the "-" value of the *sense* attribute  
48 corresponds to the "S" configuration of the chirality centre. However, for easier implementation of UCM  
49 in software tools other algorithms can be used to assign the substituents priority.

50 Another basic example in figure 7 deals with the stereochemistry of a double bond. The UCM *stereo*  
51 element in source code snippet 14 stores the stereochemical configuration of the double bond represented  
52 by the *bond* element with the *id* attribute value "E8-B-1". Substituents are represented by *node* elements  
53 with the *id* attribute values "E8-N-3", "E8-N-4", "E8-N-5" and "E8-N-6". References in the *idrefs*  
54 attribute on the *stereo* element are ordered by descending priority of those substituents (assigned using

1 Cahn-Ingold-Prelog system of priority rules). For a double bond, the reference plane contains nodes  
2 participating in the bond and is perpendicular to the plane containing these nodes and the nodes directly  
3 bonded to them. Because both higher priority substituents are on the same side of this reference plane, the  
4 *sense* attribute has the value "+", which corresponds to the "Z" configuration of the double bond.

5 The last two examples in figure 8 demonstrate that the UCM *stereo* element can be used to describe  
6 even more complicated stereochemistry. Both examples contain metal tris chelate structures, where  
7 the central metal is attached to three bidentate ligands. The source code of these examples is available  
8 in additional file 1. In that source code we included *description* elements with annotations, which  
9 explain how the *stereo* element stores the absolute configuration of the given structure or the twist ligand  
10 conformation.

11 Further details about recording the stereochemistry in UCM are provided by the documentation in  
12 additional file 2. There the documentation of the *stereo* element and its *sense* attribute also explains how  
13 to describe the stereochemistry of the square planar or octahedral complex and the chiral axis.

## 14 DISCUSSION COMPARING UCM WITH OTHER CHEMICAL FORMATS

15 Compared to other general-purpose chemical formats UCM offers an alternative approach, as it also  
16 focuses on precise data validation. Current general-purpose formats such as CDX, CDXML, CML and CT-  
17 file formats (including XDfile) use chemical elements predefined as a set of enabled text symbols.<sup>37,13,3,4</sup>  
18 UCM on the other hand relies on reusable definitions of chemical nodes composed from protons, neutrons  
19 and electrons. This enables the detailed description of chemical bonds and the precise validation of  
20 chemical content by checking various chemical rules with the pattern-based UCM Schematron schema.  
21 Additional grammar-based schemas validate UCM tree structure including the integrated formats, which  
22 can be utilized in UCM. Our built-in validation combines grammar and pattern-based XML schema  
23 languages and similarly to the other aspects of UCM functionality it is designed with future extensibility  
24 in mind. Available general-purpose chemical formats we analyzed do not offer comparable built-in  
25 validation capabilities.<sup>1</sup> Our previous analysis suggests non-XML formats often lack any form of built-in  
26 validation, while XML formats usually provide some basic validation of the format structure. Of course  
27 some non-XML formats for chemistry include validation services by providing specialized software  
28 tools – an example is Crystallographic Information File<sup>42,43</sup> although it is not a general-purpose chemical  
29 format.<sup>1</sup> Among chemical XML formats, CML and its CMLLite validator is an exception, which aims to  
30 provide additional validation capabilities beyond what is possible to achieve with a basic grammar-based  
31 schema for the given XML format.<sup>15</sup> However, in its current state custom XSLT and Java-based CMLLite  
32 validation has significant issues described in our previous article.<sup>1</sup> Therefore, UCM validation uses  
33 standard XML schema languages instead of some custom-built mechanisms and because of that it is much  
34 more compatible with the available XML validation infrastructure.

35 Besides precise built-in validation, the clear and highly readable UCM structure provides good  
36 extensibility and flexible functionality. This was achieved by utilizing the most suitable concepts found  
37 during our earlier analysis and by implementing these concepts using relatively few UCM attributes  
38 and elements.<sup>1</sup> Because of that some UCM elements can be used in more than one context. Contexts,  
39 defined by UCM Schematron schema, further restrict (and also validate) the attributes and child elements  
40 enabled for the given UCM element when it is used in the specific way or at the specific position in  
41 UCM tree structure. Other XML-based formats like CDXML or CML tend to have more attributes and  
42 elements for the given functionality and these attributes and elements are sometimes even redundant, as  
43 we discovered through our analysis.<sup>1</sup> In addition UCM is designed to enable the integration of various  
44 XML and non-XML formats, which provide some useful specialized functionality. Our approach is quite  
45 different compared to CML and its mechanism of conventions. Conventions are available in CML from  
46 version 1 and are expressed by the *convention* attribute with originally unrestricted value.<sup>8,10,14,44</sup> In the  
47 current versions of CML this attribute may still have virtually any value conforming to a pattern specifying  
48 a string similar to the XML qualified name, but with required prefix (see the CML schema 2.4 or 3<sup>12</sup>). It  
49 effectively means that for example the "my\_convention:CustomBonds-1" value denoting a convention  
50 described by someone on some website specifies a legitimate CML convention according to documentation  
51 annotations in schema version 3 or 2.4.<sup>12</sup> Such unrestricted approach adds flexibility for new use cases,  
52 but together with the content model removal in CML schema version 3 it also leads to a highly variable  
53 and difficult to process format.<sup>1</sup> Currently it remains to be seen whether the CMLLite validation and  
54 official CML conventions, which appear to be unfinished and are marked as draft recommendations, will

1 solve at least some CML issues we described in our preceding article.<sup>1</sup> UCM approach, on the other hand,  
2 is to integrate well defined XML and non-XML formats with useful functionality into predefined UCM  
3 elements. By carefully choosing formats that complement the capabilities of selected UCM elements, we  
4 believe UCM can still be sufficiently flexible for most practical use cases without disrupting its validation  
5 and other benefits.

6 While UCM avoids many issues we identified in widely used general-purpose chemical formats,<sup>1</sup> it  
7 also has its limitations as a newly developed format. These limitations could be divided into two categories.  
8 The first concerns the format functionality. UCM 1-1-1 provides promising functionality for recording  
9 chemical structures and properties. In order to become a real general-purpose chemical format UCM will at  
10 least need the support for recording chemical reactions. We are preparing a new UCM version that is going  
11 to have this functionality together with the support for mixed structures to describe chemical substances  
12 and mixtures. With additional future UCM versions, we plan to add the support for recording polymers  
13 and crystallographic information. The second category of limitations is about available software. UCM  
14 can already benefit from the standard XML tool chain, which includes the infrastructure in programming  
15 languages and software for XML processing (especially for parsing, navigation, transformation and  
16 validation of XML documents), and can be used to process and implement UCM more easily. However, it  
17 is obviously too early to expect support for UCM in specialized chemical software. To help addressing this  
18 we are now working on utilizing Open Babel (an open-source chemical toolbox for conversion between  
19 chemical formats<sup>45</sup>) to develop a software tool for automatically converting data from current chemical  
20 formats into UCM. We also plan to extend our UCM VIEWER, so that it becomes even better reference  
21 example of how to read and display all data from various UCM files.

## 22 CONCLUSIONS

23 We have developed UCM as an alternative to the current formats for common chemical data. To avoid  
24 potential issues we identified in XML and non-XML formats for chemistry, the design of UCM reflects  
25 the knowledge from our previous analysis.<sup>1</sup> Among the main benefits of UCM is the precise grammar  
26 and pattern-based validation compatible with the standard XML tool chain; in other words one can use  
27 UCM schemas to validate not only the format's tree structure, but also various chemical constraints (e.g.  
28 whether the formal charge of a structure corresponds to the total number of protons and electrons in it).  
29 Unlike most currently available chemical formats UCM supports additional non-chemical functionality  
30 by integrating dedicated XML formats using the mechanism of XML namespaces. This way UCM is  
31 not cluttered with non-chemical functionality, while at the same time, such functionality is provided  
32 in a complete and unrestricted form by a dedicated XML format: BibTeXML for literature references,  
33 UnitsML for quantities and scientific units, and XHTML for web-compatible annotations with hyperlinks  
34 and other formatting. The first version of UCM focuses on recording chemical structures with their  
35 properties. It supports structures with isotopes, ions and various types of bonds, in which participating  
36 electrons can be recorded and validated. Besides that, chemical line notations (e.g. SLN, SMILES, etc.)  
37 and identifiers (e.g. CAS RN, InChI, InChIKey, etc.) are also supported.

38 Overall we believe UCM provides very promising functionality for the core of a modern open source  
39 chemical format. It is easily extensible and can be implemented even in web-based applications as we  
40 demonstrated with UCM VIEWER, which uses JavaScript for wide compatibility across modern web  
41 browsers. In future UCM versions we plan to add support for recording chemical reactions, mixed  
42 structures (to describe substances and mixtures), polymers and crystallographic information.

## 43 AVAILABILITY AND REQUIREMENTS

44 **Project name:** UCM 1-1-1

45 **Project home page:** <http://www.universalchemicalmarkup.org/#UCM--1-1-1>

46 **Operating system(s):** platform independent

47 **Programming language:** XML, RELAX NG, Schematron, XSLT, NVDL

48 **Other requirements:** standard XML validators (e.g. xmllint and jing), standard XSLT 2.0 processor  
49 (e.g. saxon)

50 **License:** GNU GPL 3

51 **Any restrictions to use by non-academics:** None

- 1 **Project name:** UCM VIEWER 1-1-1  
2 **Project home page:** <http://www.universalchemicalmarkup.org/#UCMV--1-1-1>  
3 **Operating system(s):** platform independent  
4 **Programming language:** JavaScript, XSLT, XHTML, CSS  
5 **Other requirements:** modern web browser with a JavaScript/ECMAScript support and built-in XSLT  
6 processor (e.g. Mozilla Firefox)  
7 **License:** GNU GPL 3  
8 **Any restrictions to use by non-academics:** None

## 9 SUPPLEMENTAL INFORMATION

### 10 Additional file 1 – UCM examples and schemas

11 All UCM examples and schemas with listed validation and processing commands.

### 12 Additional file 2 – UCM documentation

13 Complete documentation for all UCM attributes and elements.

## 14 ACKNOWLEDGMENTS

15 The authors thank Professor David Sedmidubský for helpful comments and critical review of the  
16 manuscript.

## 17 COMPETING INTERESTS

18 The authors declare that they have no competing interests.

## 19 AUTHOR CONTRIBUTIONS

20 Both authors developed UCM and prepared the examples. Each file related to the format contains  
21 information about its authors.

22 **Jan Mokrý** wrote the manuscript, developed all additional software tools and prepared the website  
23 (<http://www.universalchemicalmarkup.org>).

24 **Miloslav Nič** reviewed drafts of the manuscript, provided supervision and advice.

## 25 FUNDING

26 There was no external funding for this project. All financial resources needed to carry out this research  
27 were contributed by one of the authors, Jan Mokrý.

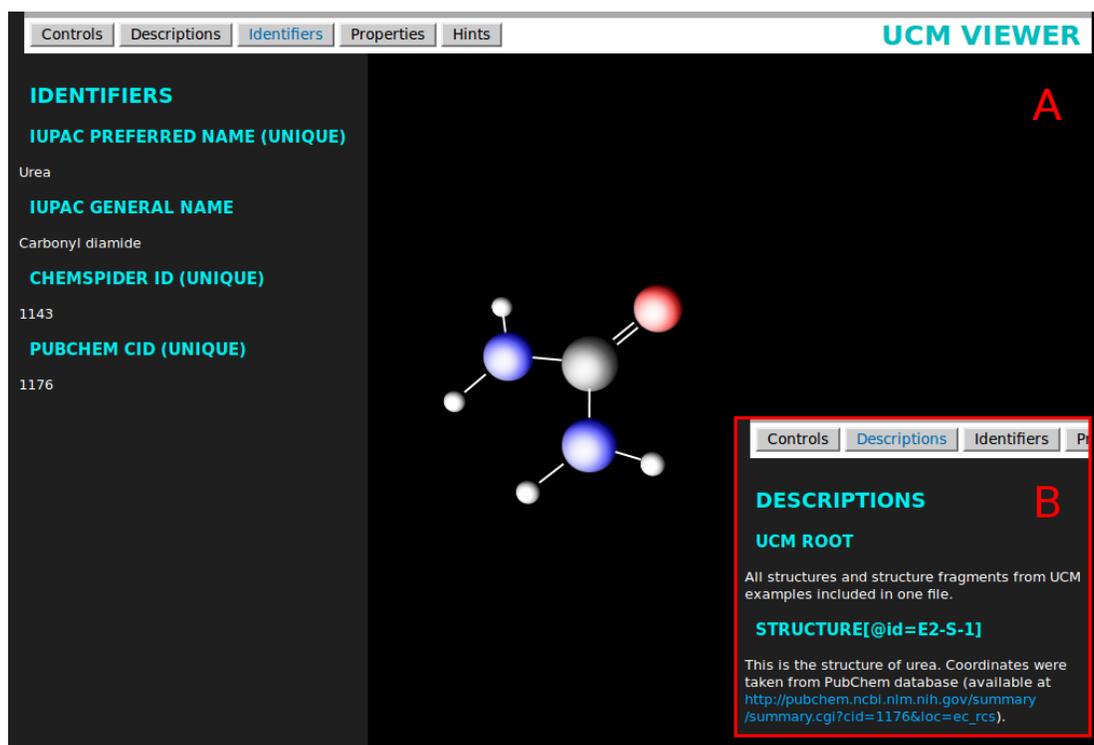
## REFERENCES

- 1 MOKRÝ, J.; NIČ, M. Designing Universal Chemical Markup (UCM) through the reusable methodology based on analyzing existing related formats. *Journal of Cheminformatics*. In press.
- 2 DALBY, A. et al. Description of several chemical structure file formats used by computer programs developed at Molecular Design Limited. *J. Chem. Inf. Comput. Sci.* 1992, vol. 32, pp. 244–255. Available also at: (<http://dx.doi.org/10.1021/CI00007A012>).
- 3 *CTfile Formats*. Accelrys. Available at: (<http://accelrys.com/products/informatics/cheminformatics/ctfile-formats/no-fee.php>).
- 4 *Symyx - CTFile Formats*. 2007. Available at: (<http://web.uni-plovdiv.bg/ksx/students/ISIS/ctfile.pdf>).
- 5 *The CDXML text-based file format*. CambridgeSoft. Available at: (<http://www.cambridgesoft.com/services/documentation/sdk/chemdraw/cdx/IntroCDXML.htm>).
- 6 LEE, D.; CHU, W. W. *Comparative Analysis of Six XML Schema Languages*. Department of Computer Science University of California, 2000. Available at: (<http://www.cobase.cs.ucla.edu/tech-docs/dongwon/ucla-200008.html>).

- 7 MURRAY-RUST, P.; RZEPA, H. S. CML: Evolution and design. *Journal of Cheminformatics*. 2011, vol. 3, pp. 44. Available also at: (<http://dx.doi.org/10.1186/1758-2946-3-44>).
- 8 MURRAY-RUST, P.; RZEPA, H. S. Chemical Markup, XML, and the Worldwide Web. 1. Basic Principles. *J. Chem. Inf. Comput. Sci.* 1999, vol. 39, pp. 928–942. Available also at: (<http://dx.doi.org/10.1021/CI990052B>).
- 9 MURRAY-RUST, P.; RZEPA, H. S. *What is CML?* Available at: (<http://www.ch.ic.ac.uk/rzepa/chimeral/documents/FAQ/FAQ.html>).
- 10 MURRAY-RUST, P.; RZEPA, H. S. Chemical Markup, XML, and the World Wide Web. 4. CML Schema. *J. Chem. Inf. Comput. Sci.* 2003, vol. 43, pp. 757–772. Available also at: (<http://dx.doi.org/10.1021/CI0256541>).
- 11 *Schemas for CML and Related Languages*. 2005. Available at: (<http://cml.sourceforge.net/schema>).
- 12 MURRAY-RUST, P.; RZEPA, H. S. *Chemical Markup Language Schema*. Available at: (<http://www.xml-cml.org/schema/>).
- 13 MURRAY-RUST, P.; RZEPA, H. S. *Chemical Markup Language Specifications*. Available at: (<http://www.xml-cml.org/spec>).
- 14 MURRAY-RUST, P.; RZEPA, H. S.; WRIGHT, M. Development of chemical markup language (CML) as a system for handling complex chemical content. *New J. Chem.* 2001, vol. 25, pp. 618–634. Available also at: (<http://dx.doi.org/10.1039/B008780G>).
- 15 TOWNSEND, J. A.; MURRAY-RUST, P. CMLite: a design philosophy for CML. *Journal of Cheminformatics*. 2011, vol. 3, pp. 39. Available also at: (<http://dx.doi.org/10.1186/1758-2946-3-39>).
- 16 BRAY, T. et al. *Namespaces in XML 1.0 (Third Edition)*. The World Wide Web Consortium (W3C), 2009. Available at: (<http://www.w3.org/TR/2009/REC-xml-names-20091208>).
- 17 BRAY, T. et al. *Namespaces in XML 1.1 (Second Edition)*. The World Wide Web Consortium (W3C), 2006. Available at: (<http://www.w3.org/TR/2006/REC-xml-names11-20060816>).
- 18 PAWSON, D. *An introduction to NVDL, ISO 19757-4*. 2008. Available at: (<http://www.dpawson.co.uk/nvdl/index.html>).
- 19 *ISO/IEC 19757-4 NVDL (Namespace-based Validation Dispatching Language)*. 2009. Available at: (<http://nvdl.org>).
- 20 GUNDERSEN, V. B.; HENDRIKSE, Z. W. *BibTeX as XML markup*. 2007. Available at: (<http://bibtexml.sourceforge.net>).
- 21 *Units Markup Language (UnitsML)*. National Institute of Standards and Technology (NIST), 2011. Available at: (<http://unitsml.nist.gov>).
- 22 CLARK, J. *Trang – Multi-format schema converter based on RELAX NG*. Thai Open Source Software Center Ltd, 2008. Available at: (<http://www.thaiopensource.com/relaxng/trang.html>).
- 23 PEMBERTON, S. et al. *XHTML 1.0: The Extensible HyperText Markup Language (Second Edition)*. The World Wide Web Consortium (W3C), 2002. Available at: (<http://www.w3.org/TR/2002/REC-xhtml1-20020801>).
- 24 CARLISLE, D.; ION, P.; MINER, R. *Mathematical Markup Language (MathML) Version 3.0 2nd Edition*. The World Wide Web Consortium (W3C), 2014. Available at: (<http://www.w3.org/TR/2014/REC-MathML3-20140410>).
- 25 PAWSON, D.; COSTELLO, R.; GEORGES, F. *ISO Schematron tutorial*. 2007. Available at: (<http://www.dpawson.co.uk/schematron/index.html>).
- 26 *ISO/IEC 19757-4:2006 Information technology – Document Schema Definition Languages (DSDL) – Part 4: Namespace-based Validation Dispatching Language (NVDL)*. International Organization for Standardization (ISO), 2006. Available at: ([http://standards.iso.org/ittf/PubliclyAvailableStandards/c038615\\_ISO\\_IEC\\_19757-4\\_2006\(E\).zip](http://standards.iso.org/ittf/PubliclyAvailableStandards/c038615_ISO_IEC_19757-4_2006(E).zip)).
- 27 *Schematron – A language for making assertions about patterns found in XML documents*. Available at: (<http://www.schematron.com/index.html>).

- 28 ROBERTSSON, E. *An Introduction to Schematron*. O'Reilly XML.com, 2003. Available at: <http://www.xml.com/pub/a/2003/11/12/schematron.html>.
- 29 ROBERTSSON, E. *Combining RELAX NG and Schematron*. O'Reilly XML.com, 2004. Available at: <http://www.xml.com/pub/a/2004/02/11/relaxtron.html?page=1>.
- 30 HANWELL, M. D. et al. Avogadro: An advanced semantic chemical editor, visualization, and analysis platform. *Journal of Cheminformatics*. 2012, vol. 4, pp. 17. Available also at: <http://dx.doi.org/10.1186/1758-2946-4-17>.
- 31 *Avogadro: An open-source molecular builder and visualization tool*. 2012. Available at: <http://avogadro.openmolecules.net>.
- 32 VALIEV, M. et al. NWChem: A comprehensive and scalable open-source solution for large scale molecular simulations. *Computer Physics Communications*. 2010, vol. 181, pp. 1477–1489. Available also at: <http://dx.doi.org/10.1016/J.CPC.2010.04.018>.
- 33 *NWChem: Delivering high-performance computational chemistry*. 2014. Available at: <http://www.nwchem-sw.org>.
- 34 RESIG, J. et al. *jQuery*. The jQuery Foundation, 2014. Available at: <http://jquery.com>.
- 35 CABELLO, R. et al. *three.js - JavaScript 3D library*. Available at: <http://threejs.org>.
- 36 MARSH, J.; ORCHARD, D.; VEILLARD, D. *XML Inclusions (XInclude) Version 1.0 (Second Edition)*. The World Wide Web Consortium (W3C), 2006. Available at: <http://www.w3.org/TR/2006/REC-xinclude-20061115>.
- 37 *CDX Format Specification*. CambridgeSoft. Available at: <http://www.cambridgesoft.com/services/documentation/sdk/chemdraw/cdx/index.htm>.
- 38 MILIORDOS, E.; XANTHEAS, S. S. On the Bonding Nature of Ozone (O<sub>3</sub>) and Its Sulfur-Substituted Analogues SO<sub>2</sub>, OS<sub>2</sub>, and S<sub>3</sub>: Correlation between Their Biradical Character and Molecular Properties. *J. Am. Chem. Soc.* 2014, vol. 136, pp. 2808–2817. Available also at: <http://dx.doi.org/10.1021/JA410726U>.
- 39 ABEL, E. W. et al. Dynamic NMR studies of ring rotation in substituted ferrocenes and ruthenocenes. *Journal of Organometallic Chemistry*. 1991, vol. 403, pp. 195–208. Available also at: [http://dx.doi.org/10.1016/0022-328X\(91\)83100-1](http://dx.doi.org/10.1016/0022-328X(91)83100-1).
- 40 BLACK, M.; MAIS, R. H. B.; OWSTON, P. G. The Crystal and Molecular Structure of Zeise's Salt, KPtCl<sub>3</sub>.C<sub>2</sub>H<sub>4</sub>.H<sub>2</sub>O. *Acta Crystallographica Section B Structural Crystallography and Crystal Chemistry*. 1969, vol. 25, pp. 1753–1759. Available also at: <http://dx.doi.org/10.1107/S0567740869004699>.
- 41 KATÔ, H. The Electronic Structure of Zeise's Salt, [PtCl<sub>3</sub>C<sub>2</sub>H<sub>4</sub>]-. *Bulletin of the Chemical Society of Japan*. 1971, vol. 44, pp. 348–354. Available also at: <http://dx.doi.org/10.1246/BCSJ.44.348>.
- 42 *Software for CIF and STAR*. The International Union of Crystallography (IUCr). Available at: <http://www.iucr.org/resources/cif/software>.
- 43 *IUCr checkCIF*. The International Union of Crystallography (IUCr). Available at: <http://checkcif.iucr.org>.
- 44 MURRAY-RUST, P. et al. The semantics of Chemical Markup Language (CML): dictionaries and conventions. *Journal of Cheminformatics*. 2011, vol. 3, pp. 43. Available also at: <http://dx.doi.org/10.1186/1758-2946-3-43>.
- 45 O'BOYLE, N. M. et al. Open Babel: An open chemical toolbox. *Journal of Cheminformatics*. 2011, vol. 3, pp. 33. Available also at: <http://dx.doi.org/10.1186/1758-2946-3-33>.

## FIGURES



**Figure 1.** The screenshot of UCM VIEWER with loaded UCM example 2 shows the structure of urea and its chemical identifiers (A). An annotation loaded from a description element is depicted in bottom right corner (B).

The screenshot displays the UCM VIEWER interface for example 3, which is ethane. The interface is divided into two main panels, A and B.

**Panel A (Properties):** This panel lists various thermodynamic and physical properties of ethane. The properties listed are:
 

- BOILING POINT** [id=E3-P-1]: Values [Kelvin]: 184.52
- ENTHALPY OF FORMATION** [id=E3-P-2]: Values [joule per mole]: -84741
- MEAN VALUE OF ENTHALPY OF COMBUSTION AT THE MEAN TEMPERATURE** [id=E3-P-3]: Values [joule per mole]: -1560665
- MEASURED VALUES OF ENTHALPY OF COMBUSTION AT THE MEAN TEMPERATURE** [id=E3-P-4]: Values [joule per mole]: -1560700, -1560740, -1560630, -1560830, -1560650, -1560440
- STANDARD DEVIATION OF THE MEAN FOR THE ENTHALPY OF COMBUSTION [ERROR]** [id=E3-P-5]: Values [joule per mole]: 48.92
- REFERENCED PROPERTIES**
  - TEMPERATURE** [id=P-CN1]: Values [Kelvin]: 298.15

 A ball-and-stick model of the ethane molecule (C<sub>2</sub>H<sub>6</sub>) is shown in the center of the interface.

**Panel B (Descriptions):** This panel provides a detailed description of the structure and associated literature references.
 

- UCM ROOT:** All structures and structure fragments from UCM examples included in one file.
- STRUCTURE** [id=E3-S-1]: This is the structure of ethane with properties. Coordinates were obtained from a skeleton structure created with Avogadro and optimized using Nwchem. Source articles for properties E3-P-1 and E3-P-2 were found using Reaxys. The source article for the property E3-P-3 was found using NIST Chemistry WebBook.
- PROPERTY** [id=E3-P-1]: Literature references: E3-REF-1
- PROPERTY** [id=E3-P-2]: Literature references: E3-REF-2
- PROPERTY** [id=E3-P-3]: Note that the values of the standard deviation of the mean and mean value are rounded in the source article. Literature references: E3-REF-3
- LITERATURE REFERENCES**
  - E3-REF-1 (ARTICLE):**

**AUTHOR:** Smith, R. D.; Udseth, H. R.  
**TITLE:** Mass spectrometry with direct supercritical fluid injection  
**JOURNAL:** Analytical Chemistry  
**YEAR:** 1983  
**VOLUME:** 55  
**PAGES:** 2266-2272  
**DOI:** 10.1021/AC00264A016  
**URL:** <http://dx.doi.org/10.1021/AC00264A016>

**Figure 2.** The screenshot of UCM VIEWER with loaded UCM example 3 shows the structure of ethane and its properties (A). Annotations with literature references loaded from description elements are depicted on right side (B).

Controls
Descriptions
Identifiers
Properties
Hints

### DESCRIPTIONS

#### UCM ROOT

All structures and structure fragments from UCM examples included in one file.

#### STRUCTURE[@id=E4-S-1]

This is the structure of diborane. Coordinates were taken from ChemSpider database (available at <http://www.chemspider.com/Chemical-Structure.17215804.html>).

#### BOND[@id=E4-B-1]

Description references: E4-D-1

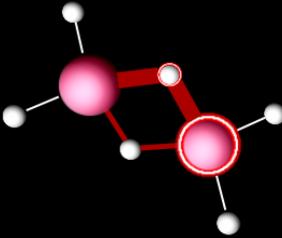
#### BOND[@id=E4-B-2]

Description references: E4-D-1

### REFERENCED DESCRIPTIONS

#### E4-D-1 (EXTERNAL DESCRIPTION)

The 3-center-2-electron bond between the node elements that represent the bridging hydrogen atom and boron atoms.



### UCM VIEWER

#### HIGHLIGHTED ITEMS DETAILS

##### BOND E4-B-1 DETAILS

**Particle counts:**  
1 Bonding electron(s) from node E4-N-1  
1 Bonding electron(s) from node E4-N-3  
**Order:**  
Delocalized or other bond

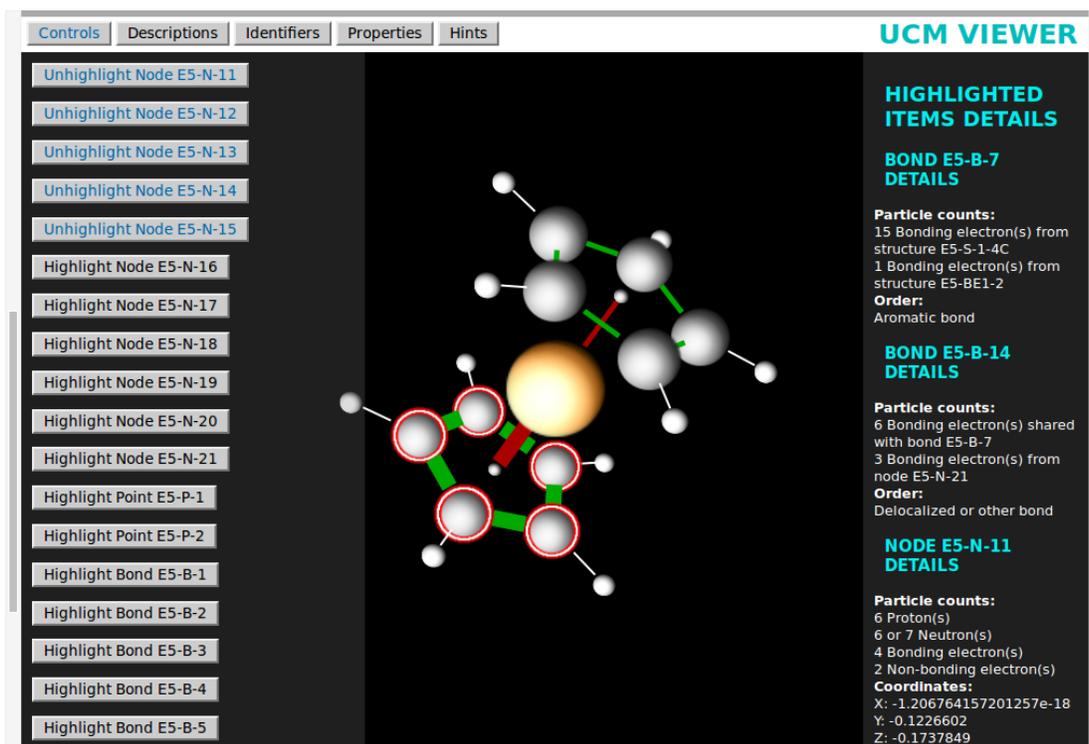
##### NODE E4-N-1 DETAILS

**Particle counts:**  
1 Proton(s)  
0 or 1 Neutron(s)  
1 Bonding electron(s)  
**Coordinates:**  
X: -0.05959875  
Y: -0.08543  
Z: -0.00022999999999999827

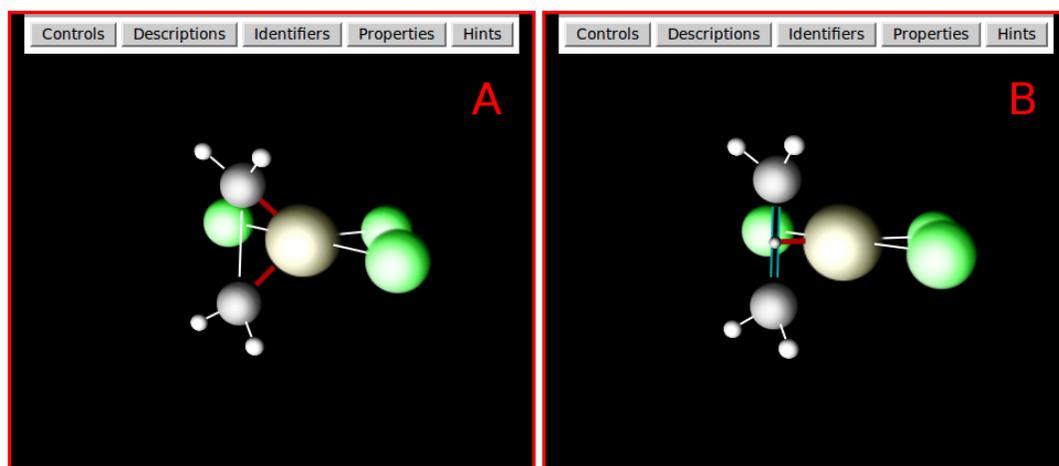
##### NODE E4-N-3 DETAILS

**Particle counts:**  
5 Proton(s)  
5 or 6 Neutron(s)  
3 Bonding electron(s)  
2 Non-bonding electron(s)  
**Coordinates:**  
X: 0.075141249999999999  
Y: -0.052389999999999999  
Z: 0.0015700000000000017

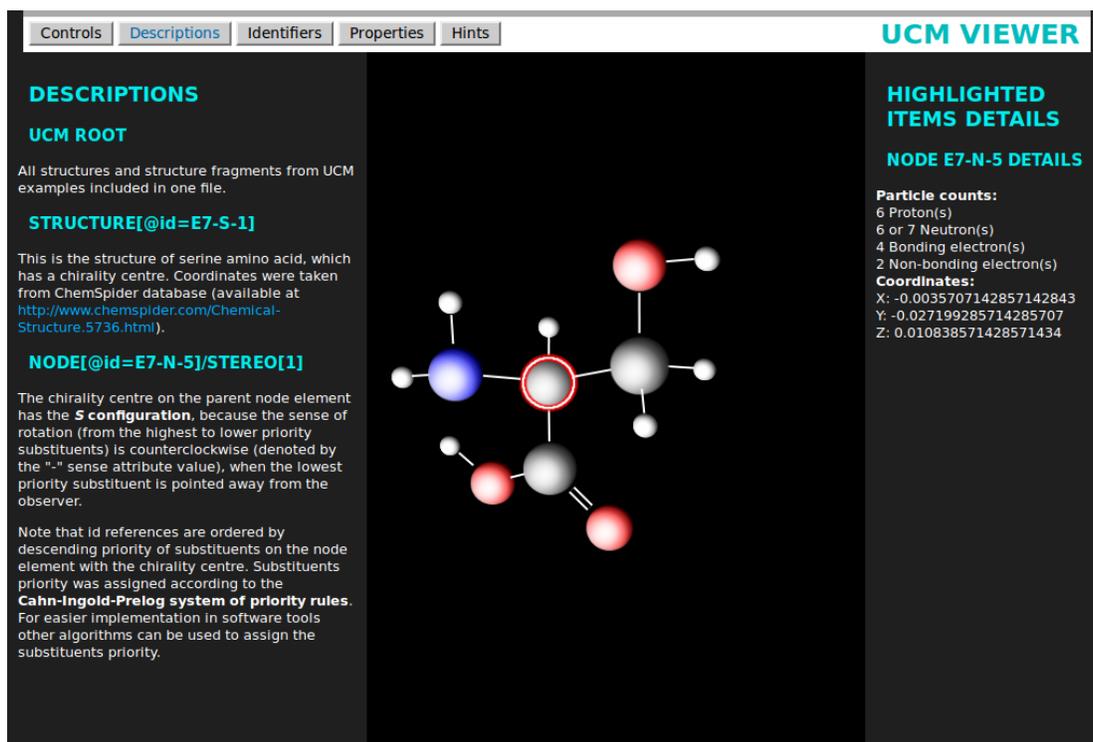
**Figure 3.** The screenshot of UCM VIEWER with loaded UCM example 4 shows the structure of diborane with 3-center-2-electron bonds.



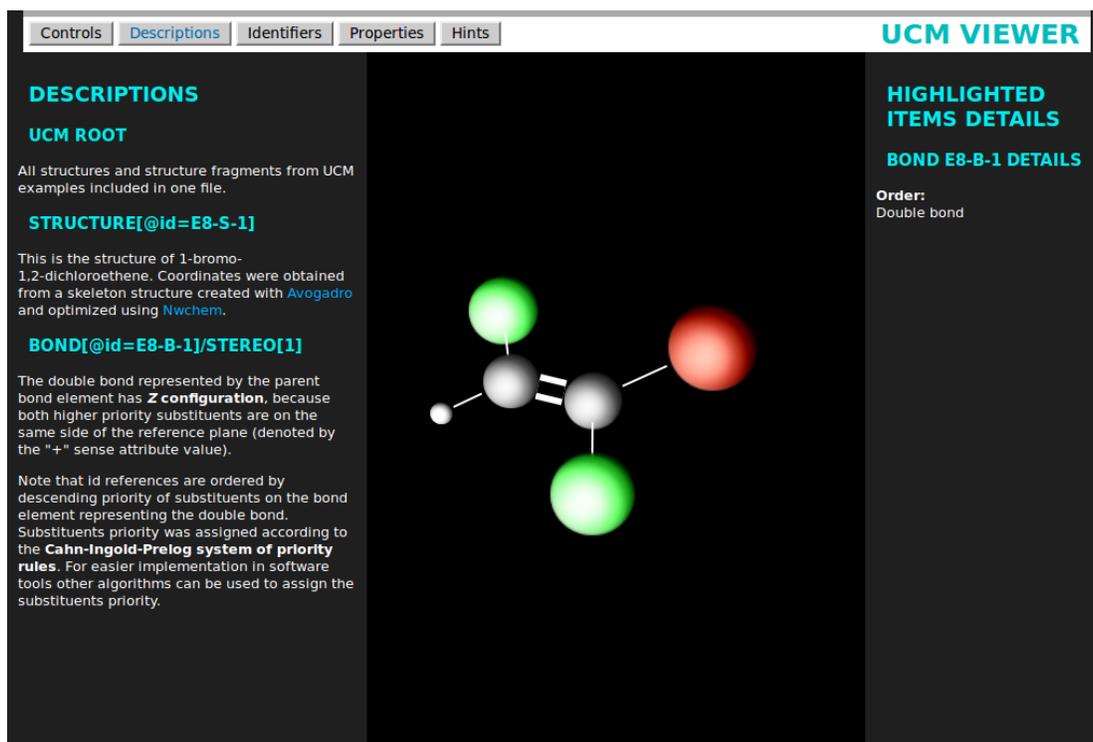
**Figure 4.** The screenshot of UCM VIEWER with loaded UCM example 5 depicts the structure of ferrocene with aromatic and delocalized bonds.



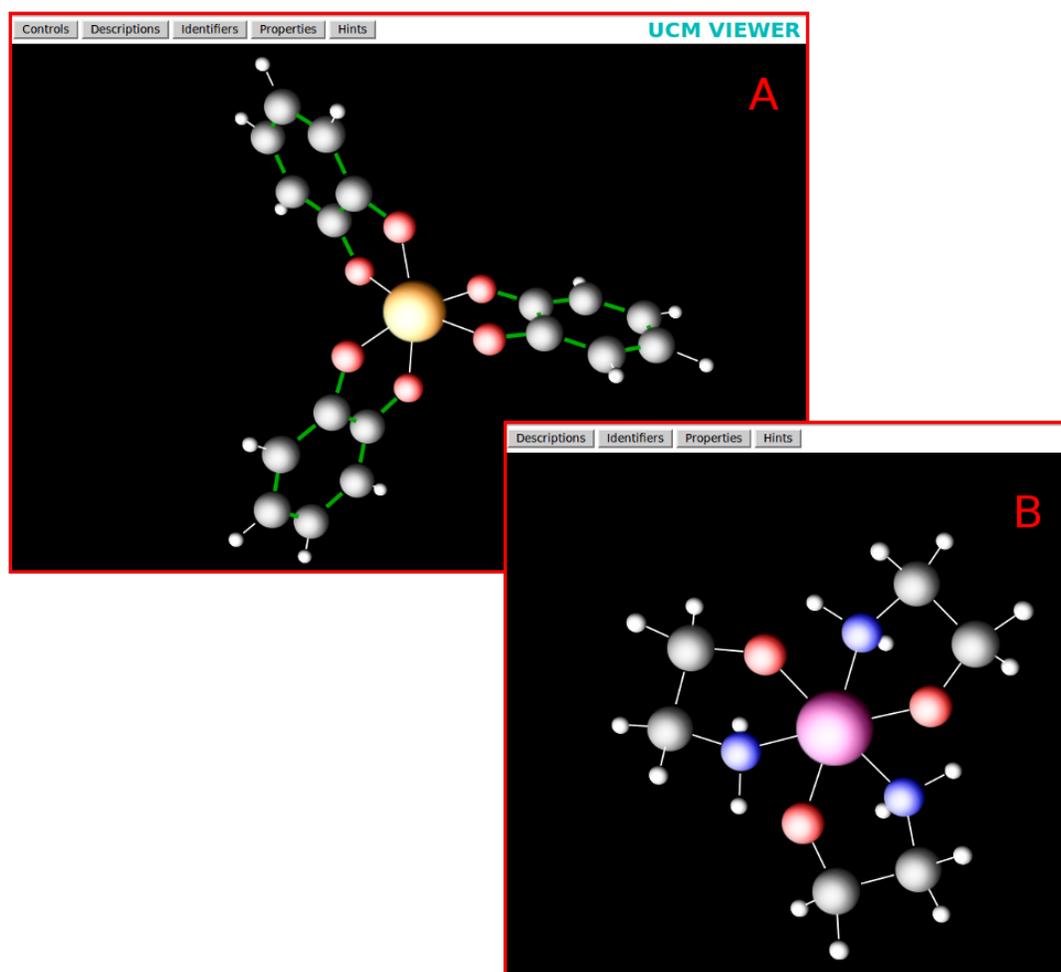
**Figure 5.** The screenshot of UCM VIEWER with loaded UCM example 6 depicts both structures of trichloro(ethene)platinate(II) anion (Zeise's salt anion) we recorded in UCM. In (A) the 3-center-2-electron bond connects the central platinum atom and ethylene ligand (a single bond is between the ligand carbon atoms). In (B) the central platinum atom and ethylene ligand are connected using a delocalized bond, which shares an electron with the partial double bond between the ligand carbon atoms.



**Figure 6.** The screenshot of UCM VIEWER with loaded UCM example 7 shows the stereochemistry of the chirality centre in serine amino acid.



**Figure 7.** The screenshot of UCM VIEWER with loaded UCM example 8 shows the stereochemistry of the double bond in 1-bromo-1,2-dichloroethene.



**Figure 8.** Screenshots of UCM VIEWER with loaded UCM examples 9 and 10 show the structure of Lambda tris(catecholato)ferrate(III) anion (A) and Lambda tris(1-hydroxy-2-aminoethane)cobalt(III) complex (B).

## TREE STRUCTURE SCHEMES

**Tree Structure Scheme 1** The first part of overall structure for UCM 1-1-1 shows the *ucm* root element and its possible content: zero or one *description* element as the first child followed by zero or more *define* and *structure* child elements.

---

```

ucm (@id?, @version)
  description [C1]? ...

  define [C1: @format = 'UCM'] (@id?, @format)*
    description [C2]* ...
    property [C1]* ...
    node [C1]* ...

  OR

  define [C2: @format = 'BIBTEXML'] (@id?, @format)*
    BIBTEXML*

  OR

  define [C3: @format = 'UNITSML'] (@id?, @format)*
    UNITSML*

  structure [C1: not(@format = 'UCM')] (@id, @format, @type)*
    IUPAC-PREFERRED-NAME-U OR IUPAC-GENERAL-NAME OR CA-INDEX-NAME
    OR CAS-RN-U OR REAXYS-RN-U
    OR CHEMSPIDER-ID-U OR PUBCHEM-CID-U OR PUBCHEM-SID
    OR INCHI OR INCHI-KEY OR S-INCHI-U OR S-INCHI-KEY
    OR SMILES OR SMARTS OR SLN

  OR

  structure [C2: @format = 'UCM'] (@id, @format, @type, @charge IF NOT 0)*
    description [C1]? ...
    structure [C1]* ... OR structure [C2]* ...
    property [C1]* ... OR property [C2]* ...
    node [C2]* ... OR node [C3]* ...
    bond [C1]* ... OR bond [C2]* ...
    point (@id, @x, @y, @z)*
      description [C1]? ...
      property [C1]* ... OR property [C2]* ...
    stereo* ...

```

---

---

**Tree Structure Scheme 2** The second part of overall structure for UCM 1-1-1 shows the remaining UCM elements with their contents.
 

---

```

description [C1: not(parent::define)] (@id?, @idrefs?, @litrefs?)
  XHTML* OR PLAINTEXT*

description [C2: parent::define] (@id, @idrefs?, @litrefs?)
  XHTML* OR PLAINTEXT*

property [C1: not(@idrefs)] (@id, @type, @quantity)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  values (@id?)?

property[C2: @idrefs] (@id, @idrefs)
  description [C1]? ...

node [C1: parent::define] (@id)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  particle [C1]+ ...

node [C2: parent::structure and @idrefs]
  (@id, @idrefs, @charge IF NOT 0, @x?, @y?, @z?)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  stereo? ...

node [C3: parent::structure and not(@idrefs)]
  (@id, @charge IF NOT 0, @x?, @y?, @z?)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  particle [C1]+ ...
  stereo? ...

bond [C1: not(@idrefs)] (@id, @order)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  join (@id?, @idrefs)+
  particle [C2]* ...
  stereo? ...

bond [C2: @idrefs] (@id, @idrefs, @order)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  particle [C2]* ...
  stereo? ...

particle [C1: parent::node] (@id?, @type, @counts, @fractions? IF @type 'N')
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...

particle [C2: parent::bond] (@id?, @idrefs, @type, @counts)
  description [C1]? ...
  property [C1]* ... OR property [C2]* ...
  share (@id?, @idrefs, @fractions)*

stereo (@id?, @idrefs, @sense)
  description [C1]? ...

```

---

## SOURCE CODE SNIPPETS

**Source Code Snippet 1** The simplified source code of BibTeXML literature references definitions. Omitted parts are indicated by ellipses.

---

```
<define xmlns="http://www.universalchemicalmarkup.org"
  xmlns:bibtexml="http://bibtexml.sf.net/"
  format="BIBTEXML">

  <bibtexml:file>
    <bibtexml:entry id="REF-0-1">
      <bibtexml:article>
        <bibtexml:author>Berglund, M. and Wieser, M. E.</bibtexml:author>
        <bibtexml:title>Isotopic compositions of the elements...</bibtexml:title>
        <bibtexml:journal>Pure and Applied Chemistry</bibtexml:journal>
        <bibtexml:year>2011</bibtexml:year>
        <bibtexml:volume>83</bibtexml:volume>
        <bibtexml:pages>397-410</bibtexml:pages>
        <bibtexml:doi>10.1351/PAC-REP-10-06-02</bibtexml:doi>
        <bibtexml:url>http://dx.doi.org/10.1351/PAC-REP-10-06-02</bibtexml:url>
      </bibtexml:article>
    </bibtexml:entry>
    ...
  </bibtexml:file>
</define>
```

---

**Source Code Snippet 2** The simplified source code of UnitsML quantities and scientific units definitions. Omitted parts are indicated by ellipses.

---

```
<define xmlns="http://www.universalchemicalmarkup.org"
  xmlns:unitsml="urn:oasis:names:tc:unitsml:schema:xsd:UnitsMLSchema-1.0"
  format="UNITSML">

  <unitsml:UnitSet>
    <unitsml:Unit xml:id="JoulePerMole">
      <unitsml:UnitName xml:lang="en-US">Joule per mole</unitsml:UnitName>
      <unitsml:RootUnits>
        <unitsml:EnumeratedRootUnit unit="joule"/>
        <unitsml:EnumeratedRootUnit powerNumerator="-1" unit="mole"/>
      </unitsml:RootUnits>
    </unitsml:Unit>
    ...
  </unitsml:UnitSet>

  <unitsml:QuantitySet>
    ...
    <unitsml:Quantity xml:id="EnthalpyOfFormation">
      <unitsml:QuantityName xml:lang="en-US">Enthalpy of formation</unitsml:QuantityName>
      <unitsml:UnitReference url="#JoulePerMole"/>
    </unitsml:Quantity>
    ...
  </unitsml:QuantitySet>
</define>
```

---

---

**Source Code Snippet 3** The simplified source code of UCM chemical nodes definitions. Omitted parts are indicated by ellipses.

---

```
<define xmlns="http://www.universalchemicalmarkup.org" format="UCM">
  ...
  <node id="H-PLUS1-BE0">
    <particle type="P" counts="1"/>
    <particle type="N" counts="0 1" fractions="0.999885 0.000115"/>
  </node>

  <node id="H2-BE1">
    <particle type="P" counts="1"/>
    <particle type="N" counts="1"/>
    <particle type="BE" counts="1"/>
  </node>
  ...
  <node id="C-BE4">
    <particle type="P" counts="6"/>
    <particle type="N" counts="6 7" fractions="0.9893 0.0107"/>
    <particle type="BE" counts="4"/>
    <particle type="NBE" counts="2"/>
  </node>
  ...
</define>
```

---

**Source Code Snippet 4** The simplified source code for the Schematron validation of UCM *charge* attributes. Omitted parts are indicated by ellipses.

```

...
<pattern name="UCM:*@CHARGE">
  ...
  <rule context="//ucm:node[parent::ucm:structure and @idrefs]">
    <let name="ID" value="./@id"/>
    <let name="IDREFS" value="./@idrefs"/>
    <let name="NODE-DEFINITION" value="//ucm:node[parent::ucm:define and @id = $IDREFS]"/>
    <let name="DEFINED-ELECTRON-COUNT"
        value="xs:integer(sum($NODE-DEFINITION/ucm:particle[contains(@type,
        'E')]/@counts))"/>
    <let name="SHARED-ELECTRON-COUNT"
        value="xs:integer(sum(//ucm:particle[ucm:share and @idrefs = $ID]/@counts))"/>
    <let name="RELEVANT-SHARE-ELEMENTS" value="//ucm:share[contains(@idrefs, $ID)]"/>
    <let name="NEGATIVE-CHARGE"
        value="lf:COUNT-NODE-NEGATIVE-CHARGE($ID, $DEFINED-ELECTRON-COUNT,
        $SHARED-ELECTRON-COUNT, $RELEVANT-SHARE-ELEMENTS)"/>
    <let name="POSITIVE-CHARGE"
        value="sum($NODE-DEFINITION/ucm:particle[@type = 'P']/@counts)"/>
    <let name="TOTAL-CHARGE" value="$POSITIVE-CHARGE - $NEGATIVE-CHARGE"/>
    <let name="CURRENT-CHARGE" value="if (./@charge) then ./@charge else 0"/>
    <assert test="$TOTAL-CHARGE = $CURRENT-CHARGE">
      <value-of select="lf:OUTPUT-ERROR-MESSAGE(.,
        (concat('The charge attribute value on this element must be ',
        $TOTAL-CHARGE, ', because it contains'),
        concat($POSITIVE-CHARGE, ' proton(s) and ', $NEGATIVE-CHARGE,
        ' electron(s).')))" />
    </assert>
  </rule>
  ...
</pattern>
...
<xsl:function name="lf:COUNT-NODE-NEGATIVE-CHARGE" as="xs:decimal">
  ...
  <xsl:param name="NODE-ID" as="xs:string"/>
  <xsl:param name="DEFINED-ELECTRON-COUNT" as="xs:integer"/>
  <xsl:param name="SHARED-ELECTRON-COUNT" as="xs:integer"/>
  <xsl:param name="RELEVANT-SHARE-ELEMENTS" as="node()*"/>
  <xsl:variable name="FRACTIONS-OF-SHARED-ELECTRON-COUNTS">
    <xsl:for-each select="$RELEVANT-SHARE-ELEMENTS">
      <xsl:value-of select="xs:decimal(subsequence(tokenize(./@fractions, '\s+'),
        index-of(tokenize(./@idrefs, '\s+'), $NODE-ID), 1))
        * xs:decimal(./parent::ucm:particle/@counts)"/>
      <xsl:if test="position() != last()">
        <xsl:text> </xsl:text>
      </xsl:if>
    </xsl:for-each>
  </xsl:variable>
  <xsl:value-of select="$DEFINED-ELECTRON-COUNT - $SHARED-ELECTRON-COUNT + sum(
    for $i in tokenize($FRACTIONS-OF-SHARED-ELECTRON-COUNTS, ' ')
    return xs:decimal($i))"/>
  ...
</xsl:function>
...

```

**Source Code Snippet 5** The simplified source code for the Schematron validation of bonding electrons in UCM. Omitted parts are indicated by ellipses.

```

...
<pattern name="UCM:NODE-2">
  ...
  <rule context="//ucm:node[parent::ucm:structure and @idrefs]">
    <let name="IDREFS" value="./@idrefs"/>
    <let name="NODE-BONDING-ELECTRONS"
        value="sum(//ucm:node[parent::ucm:define and @id = $IDREFS]/ucm:particle[
            @type = 'BE']/@counts)"/>
    <let name="REQUIRED-NODE-BONDING-ELECTRONS"
        value="1f:COUNT-NODE-BONDING-ELECTRONS(.,
            //ucm:bond[not(ucm:particle)],
            //ucm:particle[@type = 'BE'])"/>
    <assert test="$NODE-BONDING-ELECTRONS = $REQUIRED-NODE-BONDING-ELECTRONS">
      <value-of select="1f:OUTPUT-ERROR-MESSAGE(.,
          (concat('The node definition of this element contains ',
              $NODE-BONDING-ELECTRONS, ' bonding electron(s), but all'),
              concat('bonds of this element require ',
                  $REQUIRED-NODE-BONDING-ELECTRONS, ' bonding electron(s).')))" />
    </assert>
  </rule>
  ...
</pattern>
...
<xsl:function name="1f:COUNT-NODE-BONDING-ELECTRONS" as="xs:integer">
  ...
  <xsl:param name="NODE-ELEMENT" as="node()" />
  <xsl:param name="RELEVANT-BOND-ELEMENTS" as="node()*" />
  <xsl:param name="RELEVANT-PARTICLE-ELEMENTS" as="node()*" />

  <xsl:variable name="PARENT-STRUCTURE-ELEMENT"
      select="$NODE-ELEMENT/parent::ucm:structure"/>
  <xsl:variable name="PARENT-STRUCTURE-ELEMENT-ID-REGEX"
      select="concat('^|)', $PARENT-STRUCTURE-ELEMENT/@id, '( |$)'" />
  <xsl:variable name="NODE-ELEMENT-ID-REGEX"
      select="concat('^|)', $NODE-ELEMENT/@id, '( |$)'" />
  <xsl:variable name="NODE-COUNT" select="count($PARENT-STRUCTURE-ELEMENT/ucm:node)" />

  <xsl:variable name="SINGLE-BONDS-ELECTRON-COUNT"
      select="count($RELEVANT-BOND-ELEMENTS[matches(@idrefs,
          $NODE-ELEMENT-ID-REGEX) and @order = 'S'])" />
  <xsl:variable name="DOUBLE-BONDS-ELECTRON-COUNT"
      select="count($RELEVANT-BOND-ELEMENTS[matches(@idrefs,
          $NODE-ELEMENT-ID-REGEX) and @order = 'D']) * 2" />
  <xsl:variable name="TRIPLE-BONDS-ELECTRON-COUNT"
      select="count($RELEVANT-BOND-ELEMENTS[matches(@idrefs,
          $NODE-ELEMENT-ID-REGEX) and @order = 'T']) * 3" />
  <xsl:variable name="QUADRUPLE-BONDS-ELECTRON-COUNT"
      select="count($RELEVANT-BOND-ELEMENTS[matches(@idrefs,
          $NODE-ELEMENT-ID-REGEX) and @order = 'Q']) * 4" />
  <xsl:variable name="OTHER-BONDS-ELECTRON-COUNT"
      select="sum($RELEVANT-PARTICLE-ELEMENTS[matches(@idrefs,
          $NODE-ELEMENT-ID-REGEX)]/@counts)
          + (sum($RELEVANT-PARTICLE-ELEMENTS[matches(@idrefs,
              $PARENT-STRUCTURE-ELEMENT-ID-REGEX)]/@counts) div $NODE-COUNT)" />

  <xsl:value-of select="$SINGLE-BONDS-ELECTRON-COUNT + $DOUBLE-BONDS-ELECTRON-COUNT
      + $TRIPLE-BONDS-ELECTRON-COUNT + $QUADRUPLE-BONDS-ELECTRON-COUNT
      + $OTHER-BONDS-ELECTRON-COUNT" />
</xsl:function>
...

```

**Source Code Snippet 6** The simplified source code of UCM example 1 with small structures and structure fragments without 3-dimensional coordinates. Omitted parts are indicated by ellipses.

```
<ucm xmlns="http://www.universalchemicalmarkup.org"
  xmlns:xi="http://www.w3.org/2001/XInclude" version="1-1-1">
...
<xi:include href="Definitions-1-1-UCM-NODES.xml"/>
...
<structure id="E1-S-1" format="UCM" type="ST" charge="1">
...
  <node id="E1-N-1-1" idrefs="H-BE1" charge="0.5"/>
  <node id="E1-N-1-2" idrefs="H-PLUS1-BE0" charge="0.5"/>

  <bond id="E1-B-1-1" idrefs="E1-N-1-1 E1-N-1-2" order="PS">
    <particle idrefs="E1-N-1-1" type="BE" counts="1">
      <share idrefs="E1-N-1-1 E1-N-1-2" fractions="0.5 0.5"/>
    </particle>
  </bond>
</structure>

<structure id="E1-S-2" format="UCM" type="ST">
...
  <node id="E1-N-2-1" idrefs="O-BE2"/>
  <node id="E1-N-2-2" idrefs="H2-BE1"/>
  <node id="E1-N-2-3" idrefs="H2-BE1"/>

  <bond id="E1-B-2-1" idrefs="E1-N-2-1 E1-N-2-2" order="S"/>
  <bond id="E1-B-2-2" idrefs="E1-N-2-1 E1-N-2-3" order="S"/>
</structure>
...
</ucm>
```

**Source Code Snippet 7** The simplified source code of UCM example 2 with the structure of urea and its chemical identifiers. Omitted parts are indicated by ellipses.

```
<ucm xmlns="http://www.universalchemicalmarkup.org"
  xmlns:xhtml="http://www.w3.org/1999/xhtml"
  xmlns:xi="http://www.w3.org/2001/XInclude" version="1-1-1">
...
<structure id="E2-S-1" format="UCM" type="ST">

  <description><xhtml:p>This is the structure of urea. Coordinates were taken from PubChem
  database (available at <xhtml:a href="http://...">...</xhtml:a>.</xhtml:p></description>

  <structure id="E2-S-1-1" format="IUPAC-PREFERRED-NAME-U" type="STID">Urea</structure>
  ...
  <structure id="E2-S-1-4" format="PUBCHEM-CID-U" type="STID">1176</structure>

  <node id="E2-N-1" idrefs="N-BE3" x="0.06903" y="-0.11479" z="0.00001"/>
  <node id="E2-N-2" idrefs="C-BE4" x=... y=... z=.../>
  <node id="E2-N-3" idrefs="N-BE3" x=... y=... z=.../>
  <node id="E2-N-4" idrefs="O-BE2" x=... y=... z=.../>
  <node id="E2-N-5" idrefs="H-BE1" x=... y=... z=.../>
  <node id="E2-N-6" idrefs="H-BE1" x=... y=... z=.../>
  <node id="E2-N-7" idrefs="H-BE1" x=... y=... z=.../>
  <node id="E2-N-8" idrefs="H-BE1" x=... y=... z=.../>

  <bond id="E2-B-1" idrefs="E2-N-1 E2-N-2" order="S"/>
  <bond id="E2-B-2" idrefs="E2-N-1 E2-N-5" order="S"/>
  <bond id="E2-B-3" idrefs="E2-N-1 E2-N-6" order="S"/>
  <bond id="E2-B-4" idrefs="E2-N-2 E2-N-3" order="S"/>
  <bond id="E2-B-5" idrefs="E2-N-2 E2-N-4" order="D"/>
  <bond id="E2-B-6" idrefs="E2-N-3 E2-N-7" order="S"/>
  <bond id="E2-B-7" idrefs="E2-N-3 E2-N-8" order="S"/>

</structure>

</ucm>
```

**Source Code Snippet 8** The simplified source code of UCM example 3 with the recorded properties of ethane. Omitted parts are indicated by ellipses.

```

...
<xi:include href="Definitions-..." />
...
<structure id="E3-S-1" format="UCM" type="ST">
  ...
  <property id="E3-P-1" type="PR" quantity="#BoilingPoint">
    <description litrefs="REF-1"/>
    <values>184.52</values>
  </property>

  <property id="E3-P-2" type="PR" quantity="#EnthalpyOfFormation">
    <description litrefs="REF-2"/>
    <property id="E3-P-CN1-1" idrefs="P-CN1"/>
    <values>-84741</values>
  </property>

  <property id="E3-P-3" type="PR" quantity="#MeanValueOfEnthalpyOfCombustion...">
    <description litrefs="REF-3">...</description>
    <property id="E3-P-4" type="CN" quantity="#MeasuredValuesOfEnthalpy...">
      <property id="E3-P-CN1-2" idrefs="P-CN1"/>
      <values>-1560700 -1560740 -1560630 ...</values>
    </property>
    <property id="E3-P-5" type="ER" quantity="#StandardDeviationOfMean...">
      <values>48.92</values>
    </property>
    <values>-1560665</values>
  </property>
  ...
</structure>
...

```

**Source Code Snippet 9** The simplified source code of UCM example 4 with 3-center-2-electron bonds in diborane. Omitted parts are indicated by ellipses.

```

...
<structure id="E4-S-1" format="UCM" type="ST">
  ...
  <node id="E4-N-1" idrefs="H-BE1" .../>
  <node id="E4-N-2" idrefs="H-BE1" .../>
  <node id="E4-N-3" idrefs="B-BE3" .../>
  <node id="E4-N-4" idrefs="B-BE3" .../>
  ...
  <bond id="E4-B-1" order="DL">...
    <join idrefs="E4-N-1 E4-N-3 E4-N-4">CT</join>
    <particle idrefs="E4-N-1" type="BE" counts="1"/>
    <particle idrefs="E4-N-3" type="BE" counts="1"/>
  </bond>
  <bond id="E4-B-2" order="DL">...
    <join idrefs="E4-N-2 E4-N-3 E4-N-4">CT</join>
    <particle idrefs="E4-N-2" type="BE" counts="1"/>
    <particle idrefs="E4-N-4" type="BE" counts="1"/>
  </bond>
  ...
</structure>
...

```

---

**Source Code Snippet 10** The simplified source code of UCM example 1 with ozone resonance hybrid and its resonance structures without 3-dimensional coordinates. Omitted parts are indicated by ellipses.
 

---

```

...
<structure id="E1-S-3" format="UCM" type="ST">
...
  <node id="E1-N-3-1" idrefs="O-BE1" charge="-0.5"/>
  <node id="E1-N-3-2" idrefs="O-BE4" charge="1"/>
  <node id="E1-N-3-3" idrefs="O-BE1" charge="-0.5"/>

  <bond id="E1-B-3-1" order="DL">
    <description>The 3-center-2-electron bond between ... all oxygen atoms.</description>
    <join idrefs="E1-N-3-1 E1-N-3-2 E1-N-3-3">SQ</join>
    <particle idrefs="E1-N-3-2" type="BE" counts="1">
      <share idrefs="E1-N-3-1 E1-N-3-2" fractions="0.5 0.5"/>
    </particle>
    <particle idrefs="E1-N-3-2" type="BE" counts="1">
      <share idrefs="E1-N-3-2 E1-N-3-3" fractions="0.5 0.5"/>
    </particle>
  </bond>
  <bond id="E1-B-3-2" idrefs="E1-N-3-1 E1-N-3-2" order="S"/>
  <bond id="E1-B-3-3" idrefs="E1-N-3-2 E1-N-3-3" order="S"/>

</structure>

<structure id="E1-S-3-1" format="UCM" type="ST">
...
  <node id="E1-N-3-1-1" idrefs="O-BE0" charge="-1"/>
  <node id="E1-N-3-1-2" idrefs="O-BE4" charge="1"/>
  <node id="E1-N-3-1-3" idrefs="O-BE2"/>

  <bond id="E1-B-3-1-1" idrefs="E1-N-3-1-1 E1-N-3-1-2" order="S">
    <particle idrefs="E1-N-3-1-2" type="BE" counts="2">
      <share idrefs="E1-N-3-1-1 E1-N-3-1-2" fractions="0.5 0.5"/>
    </particle>
  </bond>
  <bond id="E1-B-3-1-2" idrefs="E1-N-3-1-2 E1-N-3-1-3" order="D"/>

</structure>

<structure id="E1-S-3-2" format="UCM" type="ST">...</structure>

<structure id="E1-S-3-3" format="UCM" type="ST">
...
  <node id="E1-N-3-3-1" idrefs="O-BE1"/>
  <node id="E1-N-3-3-2" idrefs="O-BE2"/>
  <node id="E1-N-3-3-3" idrefs="O-BE1"/>

  <bond id="E1-B-3-3-1" idrefs="E1-N-3-3-1 E1-N-3-3-2" order="S"/>
  <bond id="E1-B-3-3-2" idrefs="E1-N-3-3-2 E1-N-3-3-3" order="S"/>

</structure>
...

```

---

**Source Code Snippet 11** The simplified source code of UCM example 5 shows the aromatic bonding in the first cyclopentadienyl ring and also the bond between the central iron node and the first cyclopentadienyl ring. Omitted parts are indicated by ellipses.

```

...
<structure id="E5-S-1" format="UCM" type="ST">
  ...
  <structure id="E5-S-1-3" format="UCM" type="SBST" charge="-1">
    ...
    <structure id="E5-S-1-3C" format="UCM" type="SBST">
      <node id="E5-N-1" idrefs="C-BE4" .../>
      <node id="E5-N-2" idrefs="C-BE4" .../>
      <node id="E5-N-3" idrefs="C-BE4" .../>
      <node id="E5-N-4" idrefs="C-BE4" .../>
      <node id="E5-N-5" idrefs="C-BE4" .../>
    </structure>
    ...
    <node id="E5-BE1-1" charge="-1">...
      <particle type="BE" counts="1"/>
    </node>

    <bond id="E5-B-1" order="A">...
      <join idrefs="E5-N-1 E5-N-2 E5-N-3 E5-N-4 E5-N-5">CC</join>
      <particle idrefs="E5-S-1-3C" type="BE" counts="15"/>
      <particle idrefs="E5-BE1-1" type="BE" counts="1"/>
    </bond>
    ...
    <point id="E5-P-1" x="0" y="0" z="0.173016"/>

  </structure>

<structure id="E5-S-1-4" format="UCM" type="SBST" charge="-1">...</structure>

<node id="E5-N-21" idrefs="Fe-PLUS2-BE6" charge="2" ...>...</node>

<bond id="E5-B-13" idrefs="E5-N-21 E5-P-1" order="DL">...
  <particle idrefs="E5-B-1" type="BE" counts="6"/>
  <particle idrefs="E5-N-21" type="BE" counts="3"/>
</bond>
...
</structure>
...

```

**Source Code Snippet 12** The simplified source code of UCM example 6 shows both our approaches at describing the structure of trichloro(ethene)platinate(II) anion (Zeise's salt anion). Omitted parts are indicated by ellipses.

```

...
<structure id="E6-S-1" format="UCM" type="ST" charge="-1">
  ...
  <structure id="E6-S-1C" format="UCM" type="SBST">
    <node id="E6-N-1-1" idrefs="C-BE4" .../>
    <node id="E6-N-1-2" idrefs="C-BE4" .../>
  </structure>
  ...
  <node id="E6-N-1-10" idrefs="Pt-BE2" .../>
  ...
  <bond id="E6-B-1-8" idrefs="E6-N-1-1 E6-N-1-2" order="S">...</bond>
  <bond id="E6-B-1-9" order="DL">...
    <join idrefs="E6-N-1-10 E6-N-1-1 E6-N-1-2">CT</join>
    <particle idrefs="E6-S-1C" type="BE" counts="2"/>
  </bond>
</structure>

<structure id="E6-S-2" format="UCM" type="ST" charge="-1">
  ...
  <structure id="E6-S-2C" format="UCM" type="SBST">
    <node id="E6-N-2-1" idrefs="C-BE4" .../>
    <node id="E6-N-2-2" idrefs="C-BE4" .../>
  </structure>
  ...
  <node id="E6-N-2-10" idrefs="Pt-BE2" .../>
  ...
  <bond id="E6-B-2-8" idrefs="E6-N-2-1 E6-N-2-2" order="PD">...
    <particle idrefs="E6-S-2C" type="BE" counts="3"/>
  </bond>
  <bond id="E6-B-2-9" idrefs="E6-N-2-10 E6-P-2-1" order="DL">...
    <particle idrefs="E6-B-2-8" type="BE" counts="1"/>
    <particle idrefs="E6-S-2C" type="BE" counts="1"/>
  </bond>

  <point id="E6-P-2-1" x="0.4755" y="0.35813" z="-0.10135"/>
</structure>
...

```

**Source Code Snippet 13** The simplified source code of UCM example 7 with the recorded stereochemical configuration of the chirality centre in serine amino acid. Omitted parts are indicated by ellipses.

```

...
<structure id="E7-S-1" format="UCM" type="ST">
  ...
  <node id="E7-N-1" idrefs="O-BE2" .../>
  <node id="E7-N-2" idrefs="O-BE2" .../>
  <node id="E7-N-3" idrefs="O-BE2" .../>
  <node id="E7-N-4" idrefs="N-BE3" .../>
  <node id="E7-N-5" idrefs="C-BE4" ...>
    <stereo idrefs="E7-N-4 E7-N-7 E7-N-6 E7-N-8" sense="-">...</stereo>
  </node>
  <node id="E7-N-6" idrefs="C-BE4" .../>
  <node id="E7-N-7" idrefs="C-BE4" .../>
  <node id="E7-N-8" idrefs="H-BE1" .../>
  <node id="E7-N-9" idrefs="H-BE1" .../>
  <node id="E7-N-10" idrefs="H-BE1" .../>
  <node id="E7-N-11" idrefs="H-BE1" .../>
  <node id="E7-N-12" idrefs="H-BE1" .../>
  <node id="E7-N-13" idrefs="H-BE1" .../>
  <node id="E7-N-14" idrefs="H-BE1" .../>

  <bond id="E7-B-1" idrefs="E7-N-1 E7-N-6" order="S"/>
  <bond id="E7-B-2" idrefs="E7-N-1 E7-N-13" order="S"/>
  <bond id="E7-B-3" idrefs="E7-N-2 E7-N-7" order="S"/>
  <bond id="E7-B-4" idrefs="E7-N-2 E7-N-14" order="S"/>
  <bond id="E7-B-5" idrefs="E7-N-3 E7-N-7" order="D"/>
  <bond id="E7-B-6" idrefs="E7-N-5 E7-N-4" order="S"/>
  <bond id="E7-B-7" idrefs="E7-N-4 E7-N-11" order="S"/>
  <bond id="E7-B-8" idrefs="E7-N-4 E7-N-12" order="S"/>
  <bond id="E7-B-9" idrefs="E7-N-5 E7-N-6" order="S"/>
  <bond id="E7-B-10" idrefs="E7-N-5 E7-N-7" order="S"/>
  <bond id="E7-B-11" idrefs="E7-N-5 E7-N-8" order="S"/>
  <bond id="E7-B-12" idrefs="E7-N-6 E7-N-9" order="S"/>
  <bond id="E7-B-13" idrefs="E7-N-6 E7-N-10" order="S"/>

</structure>
...

```

**Source Code Snippet 14** The simplified source code of UCM example 8 with the recorded stereochemical configuration of the double bond in 1-bromo-1,2-dichloroethene. Omitted parts are indicated by ellipses.

```

...
<structure id="E8-S-1" format="UCM" type="ST">
  ...
  <node id="E8-N-1" idrefs="C-BE4" .../>
  <node id="E8-N-2" idrefs="C-BE4" .../>
  <node id="E8-N-3" idrefs="Cl-BE1" .../>
  <node id="E8-N-4" idrefs="Cl-BE1" .../>
  <node id="E8-N-5" idrefs="Br-BE1" .../>
  <node id="E8-N-6" idrefs="H-BE1" .../>

  <bond id="E8-B-1" idrefs="E8-N-1 E8-N-2" order="D">
    <stereo idrefs="E8-N-5 E8-N-4 E8-N-3 E8-N-6" sense="+">...</stereo>
  </bond>
  <bond id="E8-B-2" idrefs="E8-N-1 E8-N-3" order="S"/>
  <bond id="E8-B-3" idrefs="E8-N-2 E8-N-4" order="S"/>
  <bond id="E8-B-4" idrefs="E8-N-1 E8-N-5" order="S"/>
  <bond id="E8-B-5" idrefs="E8-N-2 E8-N-6" order="S"/>

</structure>
...

```