

# From data base to knowledge graph - using data in chemistry

Angiras Menon<sup>a,c</sup>, Nenad B. Krdzavac<sup>a,c</sup>, Markus Kraft<sup>a,b,c</sup>

<sup>a</sup>*Department of Chemical Engineering and Biotechnology, University of Cambridge, West Site, Philippa Fawcett Drive, Cambridge, United Kingdom, CB3 0AS*

<sup>b</sup>*School of Chemical and Biomedical Engineering, Nanyang Technological University, West Site, Philippa Fawcett Drive, 62 Nanyang Drive, Singapore, 637459*

<sup>c</sup>*Cambridge Centre for Advanced Research and Education in Singapore (CARES), CREATE Tower, 1 Create Way, Singapore, 138602*

---

## Abstract

Over the last couple of decades, the scientific community has made large efforts to process and store experimental and computational chemical data and information on the world wide web. This review summarizes several databases and ontologies available on the web for researchers to use. We also discuss briefly the categories of chemistry data that are stored, its main usage and how it can be accessed and understood in the framework of the Semantic Web.

*Keywords:* Knowledge Graph, Databases, Semantic Web

---

## 1. Introduction

As progress is being made in developing new and green chemical processes for a variety of industrial applications, an ever-growing amount of chemical information has been published and stored in databases online. This includes both experimental and computational chemical data. As a result, understanding how to store, access, and manipulate this vast amount of information is now key to further scientific progress. Increasingly, information science and mathematical methods such as data mining and graph theory are being used to guide various fields in chemistry and chemical engineering. Examples include analyzing organic reaction networks to understand and plan new synthetic routes for green chemistry [1, 2, 3, 4], and the use of process informatics to develop predictive chemical kinetics for combustion

13 chemistry [5]. In addition, various approaches to access and generate chem-  
14 ical knowledge are being developed using, for example, semantic web and  
15 network analysis. Semantic web technologies like knowledge graphs offer ad-  
16 ditional functionality to represent chemical knowledge. In conjunction with  
17 semantic web services the information available in chemical databases can be  
18 retrieved and changed and allows the automation of model building [6, 7, 8].  
19 The purpose of this review is to describe some of the main current databases  
20 available to researchers for data mining and review, as well as to discuss ef-  
21 forts to use ontologies as a general model for the representation of chemistry  
22 data, the improvement of the quality of these data, and the generation of  
23 resources to share consistent chemical data for a variety of purposes.

## 24 **2. Chemical Databases**

25 Several large chemical databases are available in the chemistry literature,  
26 providing a wealth of useful chemical information for researchers to use. The  
27 purpose of this section is to summarize some of the key features of such  
28 databases, for example, what information on chemical species they store and  
29 how this information can be queried. The world’s largest freely accessible  
30 database of chemical information is PubChem [9], which stores information  
31 in three primary categories: compounds, substances, and bioactivities [10, 9].  
32 Currently, PubChem has information on 97 million compounds, 242 million  
33 substances, and 280 million bioactivities [10, 9]. Information in PubChem  
34 can be queried by standard means, such as by text search, molecular formula,  
35 or chemical structure. For a common molecule, such as benzene, PubChem  
36 contains a variety of properties. This includes 2D and 3D structures as well as  
37 any crystal structures which can be downloaded in standard chemical formats  
38 such as JavaScript Object Notation (JSON), eXtensible Markup Language  
39 (XML) [11], or Common Interchange Format (CIF). PubChem also computes  
40 standard identifiers for the species in question, such as the IUPAC name, the  
41 canonical SMILES identifier [12, 13], or the InChI format [14], as well as  
42 other vendor/chemical agency identifiers. These identifiers enable identifi-  
43 cation and comparison of species between databases, so are key to linking  
44 data for the same species from different sources. Essential computed and  
45 experimental chemical and physical properties for the structure are also pro-  
46 vided by PubChem, as is any available spectral data that has been linked to  
47 the structure. PubChem also provides a large amount of information on the  
48 biological aspects of such structures, including drug information, solubility,

49 toxicity, and biological activity, which is key data for those designing drugs  
50 or green synthesis routes.

51 Another major database for chemical data is Reaxys, run by Elsevier  
52 [15, 16]. Reaxys contains much of the same information as PubChem and  
53 other chemical databases, such as structure, key identifiers, physical and  
54 chemical properties, spectral data, and biological activity for various com-  
55 pounds. What differentiates Reaxys is its focus on providing data for develop-  
56 ing synthetic routes. To this end, Reaxys has three key sets of information for  
57 a substance, namely preparations, reactions, and documents. Preparations  
58 displays key synthesis routes that can be used to prepare the substance in  
59 question. This includes the main reactions, reaction conditions, catalysts and  
60 any other information used in the synthesis routes. Each synthesis route also  
61 contains the source of the synthesis, which usually comes from the Journals  
62 and Patent databases that are linked to Reaxys via Elsevier. This enables  
63 the user to create a synthetic route for the substance of interest using Reaxys  
64 synthesis planner. Similarly, the reaction set contains the list of reactions in  
65 the Reaxys database which includes the substance the user has queried. The  
66 reactions can be filtered by structure, reagent, reaction class, solvents, cata-  
67 lysts, and yield among others, allowing the user to find reactions tailored to  
68 their application. Finally, the documents class lists the journal publications,  
69 patents, conference papers, and books that Reaxys has access to that are  
70 linked to the queried substance. This allows users of Reaxys to have access  
71 to both the data and source to analyze and select reactions.

72 Similar to Reaxys, the Chemical Abstracts Service (CAS) [17, 18] is a col-  
73 lection of databases containing information on organic and inorganic chemical  
74 substances. This information includes chemical structures, chemical names,  
75 and chemical reactions. Information stored in these databases is extracted  
76 from a wide range of literature such as patent records, journal publications,  
77 conference proceedings, Ph.D. theses, and web sources. The CAS Registry  
78 databases contain chemical structures, names, and experimental properties  
79 for more than 150 million molecules [18]. Building on the scope of the CAS  
80 Registry, the CASREACT database [19] contains several million single- and  
81 multi-step chemical reactions based on the molecules and the information  
82 stored in the CAS database. Much like Reaxys, this is provided to help users  
83 find reactions for their particular chemical application.

84 A key database for thermochemical data is the Active Thermochemical  
85 Tables (ATcT), developed by researchers at the Argonne National Labora-  
86 tory [20, 21]. The principle behind the ATcT is the thermochemical network

87 approach, which makes use of both experimental and theoretical reaction  
88 and formation enthalpies to yield estimates for the enthalpy of formation  
89 of the species in the network. The ATcT describes thermochemistry using  
90 a graph theoretic approach, with primary vertices being the enthalpies of  
91 formation of species, secondary vertices being the reaction enthalpies, and  
92 the directed edges indicating a reaction occurring between species in the net-  
93 work, with the weight determined by stoichiometry. A statistical approach  
94 is then used to analyze and solve for the optimal thermochemical values that  
95 yield a self-consistent solution. Typically, this is possible because there are  
96 multiple measurements or calculations for a given formation or reaction en-  
97 thalpy, providing the extra degrees of freedom necessary. This also means  
98 that the solution given by the ATcT can help to identify measurements that  
99 are potentially inconsistent with others in the network. Data computed by  
100 the ATcT can be found and queried online. Crucially, the reactions which  
101 contribute to the ATcT enthalpy of formation are displayed, as are uncer-  
102 tainties in the estimate of enthalpy of formation provided, making it clear  
103 which data is used and its degree of reliability.

104 On the computational chemical database side, the largest database is  
105 the Computational Chemistry Comparison and Benchmark DataBase (CC-  
106 CBDB) for thermochemical properties of species from the National Institute  
107 of Standards and Technology (NIST) [22]. Information is queried by chemi-  
108 cal name or molecular formula. The CCCBDB stores computed information  
109 in the following main categories: energy, geometry, vibrations, electrostatics,  
110 entropy and heat capacity, and reaction. All of the computed properties are  
111 displayed for the different levels of theory at which they have been calcu-  
112 lated, with the data split into categories based on the type of computational  
113 chemical method used. The CCCBDB also crucially has a comparison fea-  
114 ture, where the user can compare the results of theoretical calculations to  
115 any available experimental data in NIST’s databases, as well as look at the  
116 effect of different theoretical methods on calculated properties.

117 Other more specialized databases also exist. For example, the Alexandria  
118 library developed by van der Spoel et al. consists of molecular properties for  
119 force field development [23]. Alexandria contains molecular structures and  
120 properties for 2,704 compounds, many of which contain functional groups  
121 common to biomolecules and drugs. Alexandria contains similar informa-  
122 tion to the CCCBDB, but crucially provides more extensive multipole and  
123 polarizability calculations to guide researchers who want to develop poten-  
124 tials and force fields. Importantly, all properties in Alexandria are provided

125 at the same level of theory and the Gaussian input and output files from  
126 the calculations are also given, making reproduction of the stored informa-  
127 tion significantly easier. Even more specialized databases for computational  
128 chemists exist, such as Head-Gordon and Hait’s benchmark database specif-  
129 ically for DFT calculations on dipole moments, spanning a variety of func-  
130 tionals in the process [24]. The database from Simmie et al. is specifically  
131 for high-level enthalpies of formation for nitrogen based compounds [25].  
132 The GDB-17 database specifically enumerates small organic molecules, using  
133 graph-theoretic methods to span 166 billion such molecules with the aim of  
134 guiding new drug design [26]. Ramakrishnan et al. provide the QM9 dataset,  
135 containing DFT calculations on around 134,000 molecules for training new  
136 machine learning potentials [27]. The ANI-1 data set uniquely contains non-  
137 equilibrium DFT calculations, that is for molecules in conformers that are  
138 not their minimum energy ground state configuration [28]. ANI-1 contains  
139 around 20 million molecular conformations for 57,462 molecules taken from  
140 the GDB database. There is clearly a wide variety of chemical data, both ex-  
141 perimental and computational, that is available to researchers in a variety of  
142 fields in chemistry. This data is ever growing, and methods to store, access,  
143 and act on this data automatically are becoming more valuable for progress  
144 to be made.

### 145 **3. Ontologies for Computational Chemistry**

146 Given the variety of chemical data available, developing a consistent  
147 framework to store and access it is crucial, even more so as the amount  
148 of data available is expanding rapidly. Further data processing will increas-  
149 ingly rely on automation allowing machines to interpret, integrate, share,  
150 and perform reasoning with data of various formats.

151 One of the early efforts in storing chemical data in a standard format was  
152 the introduction of Chemical Markup Language (CML) pioneered by Murray-  
153 Rust and coworkers [29, 30, 31, 32]. The CML format is based on XML, which  
154 is suitable for storing data of any level of complexity while providing semantic  
155 information to the data stored. CML allows the representation of complex  
156 chemical objects by employing the hierarchical tree structure of XML using  
157 chemical name tags which cover different aspects of chemistry. Over the past  
158 20 years, CML has been developed to represent most aspects of chemistry,  
159 including CMLReact for chemical reactions [33], CMLSpec for spectral data

160 [34], CML for crystallography [35], and CML for polymers (PML) [36] along  
161 with the standard labels and definitions for physical properties.

162 Building on this established format for representing chemical data, Phan-  
163 dungsukanan and coworkers developed a sub-domain for storing quantum  
164 chemistry calculations data based on CML, termed CompChem [37]. The  
165 main goal of CompChem was to introduce a stricter structure into CML-  
166 based documents so that software tools know exactly how to validate and  
167 process information related to computational chemistry. To this end, the se-  
168 mantics of data stored in the CompChem based documents is modelled based  
169 on the typical nature of computational simulations or calculations, contain-  
170 ing information on the job type, input parameters, and output parameters  
171 that one would expect in these calculations. This enables the storage of a  
172 variety of output data from *ab initio* quantum chemistry calculations such  
173 as the results of geometry optimization, single point energy calculations, and  
174 frequency calculations, among others. The storage and access of this data  
175 was realized through a MolHub web service [37]. However, the original Mol-  
176 Hub did not allow for semantic inter-operability between different chemistry  
177 software tools, provide an efficient query engine, or guarantee the consistency  
178 of data.

179 To alleviate these shortcomings, a novel OntoCompChem ontology has  
180 been developed by extending the Gainesville Core (GNVC) ontology [38]  
181 while supporting the CompChem convention of CML [39]. The OntoCom-  
182 pChem ontology is currently populated by Gaussian quantum chemistry  
183 calculations through an updated version of the MolHub semantic web ser-  
184 vice (<https://como.ceb.cam.ac.uk/resources/molhub/>). The OntoCom-  
185 pChem knowledge graph forms part of a more general knowledge graph called  
186 the J-Park Simulator (JPS) [40]. This architecture supports semantic inter-  
187 operability between different domains and allows the use of propositional  
188 logic, formal query language, and Semantic Web tools such as the HerMiT  
189 [41] reasoner to check the consistency of data within the JPS knowledge  
190 graph. More recently, the OntoKin ontology [42, 43] has been developed as  
191 a component of the JPS to represent gas phase elementary reactions, which  
192 are the building block of large reaction mechanisms found in combustion  
193 and atmospheric chemistry models. The ontology allows inference engines  
194 to detect inconsistencies in chemical mechanisms and to perform semantic  
195 queries across mechanisms stored in the JPS knowledge graph. At present,  
196 both the OntoKin and MolHub frameworks are missing an intelligent system  
197 that automatically establishes semantic inter-operability between quantum

198 chemistry calculations and kinetic mechanisms. To achieve this goal, we are  
199 currently developing a formal framework that is based on reinforcement learn-  
200 ing formal tools [44], modal logic [45], and a propositional logic framework  
201 with binary metric operators [46] to provide formal language support.

202 In addition to the JPS efforts, other semantic frameworks are currently  
203 in use. The Chemical Semantics Framework (CSF) [47] stores results of  
204 quantum chemistry calculations. The core of the CSF is the GNVC ontology  
205 which forms the knowledge component of the framework. However, the ontol-  
206 ogy does not support all of CompChem’s conventions for CML features. For  
207 example, some keywords in the CML format such as geometry type are not  
208 supported. In addition, the CSF does not support semantic inter-operability  
209 between different computational chemistry tools. However, the framework  
210 allows web agents to access and, in principle, act on data stored in the CSF,  
211 representing a step towards automation of the knowledge graph. The ChEBI  
212 database stores molecular entities focused on ‘small’ chemical compounds,  
213 that is part of the Open Biomedical Ontologies effort. It uses the ChEBI  
214 ontology as a common model for classification of chemical compounds in the  
215 biomedical field. The ontology provides models for molecular structures such  
216 as hydrocarbons, common chemical roles for the molecules in the ontology,  
217 as well as for information pertaining to subatomic particles [48]. The ChEBI  
218 database can be explored using an advanced search interface, but semantic  
219 inter-operability and web agent access is currently not supported.

220 The review of ontologies for chemistry makes it clear that plenty of effort  
221 is being put towards developing methods for storing, accessing, and interpret-  
222 ing the available chemical data in an intelligent way. Key to the success of  
223 these efforts will be the development of standards for the publication and re-  
224 porting of chemical data. By having a standard format for reporting chemical  
225 data, linking this information to a semantic framework or ontology becomes  
226 substantially easier and less error prone. Efforts to this end include the work  
227 of the InChI consortium [14], the Allotrope Foundation’s work on developing  
228 a standard data format, and the work of Cronin and coworkers on developing  
229 a chemical programming language that can be used to represent experimen-  
230 tal organic chemistry [49]. These standards will help inspire the definition of  
231 classes in chemical ontologies. In conjunction with this, the development of  
232 tools for establishing semantic frameworks, as well as agents that can act on  
233 this data automatically, is still in process. This will eventually enable a self-  
234 consistent and ever-growing chemical knowledge graph based on ontologies  
235 and automated by web agents.

## 236 4. Summary and Outlook

237 In this review, we have discussed how the rapidly increasing amount of  
238 chemical information available to researchers has necessitated the develop-  
239 ment of automated methods to query, store, and share this information for a  
240 variety of applications. We have discussed some of the main databases and  
241 the usage of ontologies in the chemistry domain. Moving forward, it is hoped  
242 that more tools will be developed to provide more intelligent ways to create,  
243 update, retrieve, and maintain distributed chemical information via the Web.  
244 It is also necessary to develop tools to support more advanced community in-  
245 volvement, bridging data silos, and identifying "best" data for the solution of  
246 a particular problem. Eventually, the chemical knowledge graph will be fully  
247 automated and self-improving to provide, for example, new synthesis routes  
248 and more reliable chemical models built on the experimental and chemical  
249 data provided in the variety of databases online.

## 250 5. Acknowledgments

251 AM acknowledges Johnson Matthey for financial support. The authors  
252 also acknowledge the financial support of the Singapore National Research  
253 Foundation (NRF) through the Campus for Research Excellence and Tech-  
254 nological Enterprise (CREATE) program. MK gratefully acknowledges the  
255 support of the Alexander von Humboldt foundation.

## 256 References

- 257 [1] \*P.M Jacob, A. Lapkin, Statistics of the network of organic chemistry,  
258 *React. Chem. Eng.* 3 (2018) 102–118.  
259 Using graph-theoretic methods, the authors analyze the structure of a  
260 network of organic reactions built on chemical data mined from Reaxys.  
261 The authors show that on average most molecules can be synthesized  
262 within six steps from any other molecule, in what is the first study on  
263 such a large network.
- 264 [2] M. Fialkowski, K. J. M. Bishop, V. A. Chubukov, C. J. Campbell, B. A.  
265 Grzybowski, Architecture and evolution of organic chemistry, *Ange-  
266 wandte Chemie International Edition* 44 (2005) 7263–7269.

- 267 [3] K. J. M. Bishop, R. Klajn, B. A. Grzybowski, The core and most  
268 useful molecules in organic chemistry, *Angewandte Chemie International*  
269 *Edition* 45 (2006) 5348–5354.
- 270 [4] B. A. Grzybowski, K. J. Bishop, B. Kowalczyk, C. E. Wilmer,  
271 The 'wired' universe of organic chemistry, *Nature Chemistry* 1 (2009)  
272 31.
- 273 [5] M. Frenklach, Transforming data into knowledgeprocess informatics  
274 for combustion chemistry, *Proceedings of the Combustion Institute* 31  
275 (2007) 125 – 140.
- 276 [6] M. F. Lopez, A. Gomez-Perez, J. P. Sierra, A. P. Sierra, Building a  
277 chemical ontology using methontology and the ontology design environ-  
278 ment, *IEEE Intelligent Systems and their Applications* 14 (1999) 37–46.
- 279 [7] K. R. Taylor, R. J. Gledhill, J. W. Essex, J. G. Frey, S. W. Harris,  
280 D. C. De Roure, Bringing chemical data onto the semantic web, *Jour-  
281 nal of Chemical Information and Modeling* 46 (2006) 939–952. PMID:  
282 16711712.
- 283 [8] J. Hastings, L. Chepelev, E. Willighagen, N. Adams, C. Steinbeck,  
284 M. Dumontier, The chemical information ontology: Provenance and  
285 disambiguation for chemical data on the biological semantic web, *PLOS*  
286 *ONE* 6 (2011) 1–13.
- 287 [9] \*S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A.  
288 Shoemaker, P. A. Thiessen, B. Yu, L. Zaslavsky, J. Zhang, E. E. Bolton,  
289 PubChem 2019 update: improved access to chemical data, *Nucleic Acids*  
290 *Research* 47 (2018) D1102–D1109.  
291 The authors summarize the information available in PubChem, the  
292 world's largest open source chemical database. The authors have also  
293 expanded PubChem to include spectral information, links to scientific  
294 articles, as well as biological properties for food and agricultural chem-  
295 icals.
- 296 [10] S. Kim, P. A. Thiessen, T. Cheng, B. Yu, B. A. Shoemaker, J. Wang,  
297 E. E. Bolton, Y. Wang, S. H. Bryant, Literature information in pub-  
298 chem: associations between pubchem records and scientific articles,  
299 *Journal of cheminformatics* 8 (2016) 32.

- 300 [11] T. Bray, J. Paoli, C. M. Sperberg-McQueen, E. Maler, F. Yergeau, Ex-  
301 tensible markup language (xml) 1.0, 2000.
- 302 [12] D. Weininger, Smiles, a chemical language and information system. 1.  
303 introduction to methodology and encoding rules, *Journal of chemical*  
304 *information and computer sciences* 28 (1988) 31–36.
- 305 [13] N. M. OBoyle, Towards a universal smiles representation-a standard  
306 method to generate canonical smiles based on the inchi, *Journal of*  
307 *cheminformatics* 4 (2012) 22.
- 308 [14] S. R. Heller, A. McNaught, I. Pletnev, S. Stein, D. Tchekhovskoi, Inchi,  
309 the iupac international chemical identifier, *Journal of cheminformatics*  
310 7 (2015) 23.
- 311 [15] J. Goodman, Computer software review: Reaxys, 2009.
- 312 [16] A. J. Lawson, The making of reaxys-towards unobstructed access to  
313 relevant chemistry information, *The Future of the History of Chemical*  
314 *Information* 1164 (2014) 127–48.
- 315 [17] K. J. Meloche, J. Mears, R. J. Schenck, Intriguing Records in CAS  
316 Databases, pp. 21–40.
- 317 [18] Cas registry database, 2019. Accessed May 23rd, 2019.
- 318 [19] Casreact - cas chemical reactions database, 2019. Accessed May 23rd,  
319 2019.
- 320 [20] B. Ruscic, R. E. Pinzon, M. L. Morton, G. von Laszewski, S. J. Bittner,  
321 S. G. Nijsure, K. A. Amin, M. Minkoff, A. F. Wagner, Introduction  
322 to active thermochemical tables: Several key enthalpies of formation  
323 revisited, *The Journal of Physical Chemistry A* 108 (2004) 9979–9997.
- 324 [21] B. Ruscic, R. E. Pinzon, G. Von Laszewski, D. Kodeboyina, A. Burcat,  
325 D. Leahy, D. Montoy, A. F. Wagner, Active thermochemical tables:  
326 thermochemistry for the 21st century, in: *Journal of Physics: Confer-*  
327 *ence Series*, volume 16, IOP Publishing, p. 561.
- 328 [22] R. Johnson III, Cccbdb computational chemistry comparison and bench-  
329 mark database, NIST Standard Reference Database Number 101 (1999).

- 330 [23] \*M.M. Ghahremanpour, P. J. Van Maaren, D. Van Der Spoel, The  
331 alexandria library, a quantum-chemical database of molecular properties  
332 for force field development, *Scientific data* 5 (2018) 180062.  
333 The authors provide an open source database of quantum chemistry  
334 calculations for 2704 compounds. This establishes a key training set for  
335 the development of empirical forcefields for a variety of molecules and  
336 applications.
- 337 [24] \*D. Hait, M. Head-Gordon, How accurate is density functional theory  
338 at predicting dipole moments? an assessment using a new database of  
339 200 benchmark values, *Journal of chemical theory and computation* 14  
340 (2018) 1969–1981.  
341 The authors provide 200 benchmark dipole moments calculated using  
342 coupled cluster theory. This study then develops a hierarchy of den-  
343 sity functionals for accurately predicting dipole moments, crucial to the  
344 development of intermolecular potentials.
- 345 [25] J. M. Simmie, A database of formation enthalpies of nitrogen species by  
346 compound methods (cbs-qb3, cbs-apno, g3, g4), *The Journal of Physical  
347 Chemistry A* 119 (2015) 10511–10526.
- 348 [26] L. Ruddigkeit, R. Van Deursen, L. C. Blum, J.-L. Reymond, Enumer-  
349 ation of 166 billion organic small molecules in the chemical universe  
350 database gdb-17, *Journal of chemical information and modeling* 52  
351 (2012) 2864–2875.
- 352 [27] R. Ramakrishnan, P. O. Dral, M. Rupp, O. A. Von Lilienfeld, Quantum  
353 chemistry structures and properties of 134 kilo molecules, *Scientific data*  
354 1 (2014) 140022.
- 355 [28] J. S. Smith, O. Isayev, A. E. Roitberg, Ani-1, a data set of 20 million  
356 calculated off-equilibrium conformations for organic molecules, *Scientific  
357 data* 4 (2017) 170193.
- 358 [29] P. Murray-Rust, H. S. Rzepa, Chemical markup, xml, and the worldwide  
359 web. 1. basic principles, *Journal of Chemical Information and Computer  
360 Sciences* 39 (1999) 928–942.
- 361 [30] G. V. Gkoutos, P. Murray-Rust, H. S. Rzepa, M. Wright, Chemical  
362 markup, xml, and the world-wide web. 3. toward a signed semantic

- 363 chemical web of trust, *Journal of Chemical Information and Computer*  
364 *Sciences* 41 (2001) 1124–1130. PMID: 11604013.
- 365 [31] P. Murray-Rust, H. S. Rzepa, Chemical markup, xml, and the world  
366 wide web. 4. cml schema, *Journal of Chemical Information and Com-*  
367 *puter Sciences* 43 (2003) 757–772. PMID: 12767134.
- 368 [32] J. A. Townsend, P. Murray-Rust, Cmlite: a design philosophy for cml,  
369 *Journal of Cheminformatics* 3 (2011) 39.
- 370 [33] G. L. Holliday, P. Murray-Rust, H. S. Rzepa, Chemical markup, xml,  
371 and the world wide web. 6. cmlreact, an xml vocabulary for chemical  
372 reactions, *Journal of chemical information and modeling* 46 (2006) 145–  
373 157.
- 374 [34] S. Kuhn, T. Helmus, R. J. Lancashire, P. Murray-Rust, H. S. Rzepa,  
375 C. Steinbeck, E. L. Willighagen, Chemical markup, xml, and the world  
376 wide web. 7. cmlspect, an xml vocabulary for spectral data, *Journal of*  
377 *chemical information and modeling* 47 (2007) 2015–2034.
- 378 [35] N. Day, J. Downing, S. Adams, N. England, P. Murray-Rust, Crystaleye,  
379 URL <http://wwmm.ch.cam.ac.uk/crystaleye/>. Online (2008).
- 380 [36] N. Adams, J. Winter, P. Murray-Rust, H. S. Rzepa, Chemical markup,  
381 xml and the world-wide web. 8. polymer markup language, *Journal of*  
382 *chemical information and modeling* 48 (2008) 2118–2128.
- 383 [37] W. Phadungsukanan, M. Kraft, J. A. Townsend, P. Murray-Rust, The  
384 semantics of chemical markup language (cml) for computational chem-  
385 istry: Compchem, *Journal of cheminformatics* 4 (2012) 15.
- 386 [38] N. S. Ostlund, M. Sopek, GNVC: Gainesville core ontology - standard for  
387 publishing results of computational chemistry, ver. 0.7, 2015. Accessed  
388 October 24th, 2018.
- 389 [39] \*N. Krdzavac, S. Mosbach, D. Nurkowski, P. Buerger, J. Akroyd, J. Mar-  
390 tin, A. Menon, M. Kraft, An ontology and semantic web service for  
391 quantum chemistry calculations, *Journal of chemical information and*  
392 *modeling* 59 (2019) 3154–3165.  
393 The authors develop the OntoCompChem ontology by extending the

- 394 Gainesville Core (GNVC) ontology and establish semantic interoperabil-  
395 ity between different tools used in quantum chemistry and thermochem-  
396 istry calculations. The new ontology's use is demonstrated by querying  
397 the results from quantum chemistry calculations and using these to per-  
398 form thermodynamic data calculations for the species of interest.
- 399 [40] M. Kraft, S. Mosbach, The future of computational modelling in reaction  
400 engineering, *Philos. Trans. R. Soc., A* 368 (2010) 3633–3644.
- 401 [41] B. Glimm, I. Horrocks, B. Motik, G. Stoilos, Z. Wang, Hermit: An  
402 OWL 2 reasoner, *J. Autom. Reasoning* 53 (2014) 245–269.
- 403 [42] \* F. Farazi, J. Akroyd, S. Mosbach, P. Buerger, D. Nurkowski, M. Kraft,  
404 OntoKin: An ontology for chemical kinetic reaction mechanisms, 2019.  
405 Submitted for publication.  
406 The authors develop an ontology capable of storing data from chemical  
407 kinetics and chemical reaction mechanisms by using OWL and formal  
408 reasoning tools. The new ontology's use is demonstrated by querying  
409 and browsing different mechanism as well as modelling the atmospheric  
410 dispersion of pollutants formed in an internal combustion engine.
- 411 [43] F. Farazi, N. B. Krdzavac, J. Akroyd, S. Mosbach, A. Menon,  
412 D. Nurkowski, M. Kraft, Linking reaction mechanisms and quantum  
413 chemistry: An ontological approach (2019). Submitted for publication.
- 414 [44] R. S. Sutton, A. G. Barto, Reinforcement Learning, MIT Press,  
415 Cambridge- Massachusetts, London -England, 2nd edition, 2018.
- 416 [45] A. V. Chagrov, M. Zakharyashev, Modal Logic, volume 35 of *Oxford*  
417 *logic guides*, Oxford University Press, 1997.
- 418 [46] N. Stojanovic, N. Ikodinovic, R. Djordjevic, A propositional logic with  
419 binary metric operators, *Journal of Applied Logics - IfCoLog Journal*  
420 *of Logics and their Applications* 5 (2018) 1605–1622.
- 421 [47] N. S. Ostlund, M. Sopek, Applying the semantic web to computational  
422 chemistry, in: A. Paschke (Ed.), *Proceedings of the 6th International*  
423 *Workshop on Semantic Web Applications and Tools for Life Sciences (*  
424 *SWAT4LS 2013)*, Edinburgh, United Kingdom. Accessed February 7th,  
425 2019.

- 426 [48] J. Hastings, G. Owen, A. Dekker, M. Ennis, N. Kale, V. Muthukrish-  
427 nan, S. Turner, N. Swainston, P. Mendes, C. Steinbeck, Chebi in 2016:  
428 Improved services and an expanding collection of metabolites, *Nucleic  
429 acids research* 44 (2016) D1214–9.
- 430 [49] S. Steiner, J. Wolf, S. Glatzel, A. Andreou, J. M. Granda, G. Keenan,  
431 T. Hinkley, G. Aragon-Camarasa, P. J. Kitson, D. Angelone, et al.,  
432 Organic synthesis in a modular robotic system driven by a chemical  
433 programming language, *Science* 363 (2019) eaav2211.