

1 Scaffold Generator - A Java library  
2 implementing molecular scaffold  
3 functionalities in the Chemistry  
4 Development Kit (CDK)

5

6 Jonas Schaub; Institute for Inorganic and Analytical Chemistry; Friedrich-Schiller University,  
7 Lessing Strasse 8, 07743, Jena, Germany; [jonas.schaub@uni-jena.de](mailto:jonas.schaub@uni-jena.de); ORCID: [0000-0003-  
8 1554-6666](https://orcid.org/0000-0003-1554-6666)

9

10 Julian Zander; Institute for Bioinformatics and Chemoinformatics, Westphalian University of  
11 Applied Sciences, August-Schmidt-Ring 10, 45665, Recklinghausen, Germany;  
12 [zanderjulian@gmx.de](mailto:zanderjulian@gmx.de); ORCID: [0000-0001-8197-076X](https://orcid.org/0000-0001-8197-076X)

13

14 Achim Zielesny; Institute for Bioinformatics and Chemoinformatics, Westphalian University of  
15 Applied Sciences, August-Schmidt-Ring 10, 45665, Recklinghausen, Germany;  
16 [achim.zielesny@w-hs.de](mailto:achim.zielesny@w-hs.de); ORCID: [0000-0003-0722-4229](https://orcid.org/0000-0003-0722-4229)

17

18 Christoph Steinbeck\*; Institute for Inorganic and Analytical Chemistry; Friedrich-Schiller  
19 University, Lessing Strasse 8, 07743, Jena, Germany; [christoph.steinbeck@uni-jena.de](mailto:christoph.steinbeck@uni-jena.de);  
20 ORCID: [0000-0001-6966-0814](https://orcid.org/0000-0001-6966-0814)

21

22 \*Corresponding author email: [christoph.steinbeck@uni-jena.de](mailto:christoph.steinbeck@uni-jena.de)

## 23 Abstract

24 The concept of molecular scaffolds as defining core structures of organic molecules is utilised  
25 in many areas of chemistry and cheminformatics, e.g. drug design, chemical classification, or  
26 the analysis of high-throughput screening data. Here, we present Scaffold Generator, a  
27 comprehensive open library for the generation, handling, and display of molecular scaffolds,  
28 scaffold trees and networks. The new library is based on the Chemistry Development Kit  
29 (CDK) and highly customisable through multiple settings, e.g. five different structural  
30 framework definitions are available. For display of scaffold hierarchies, the open GraphStream  
31 Java library is utilised. Performance snapshots with natural products (NP) from the COCONUT  
32 (COLleCtion of Open Natural prodUcTs) database and drug molecules from DrugBank are  
33 reported. The generation of a scaffold network from more than 450,000 NP can be achieved  
34 within a single day.

35

36 **Keywords:** cheminformatics, Chemistry Development Kit, CDK, natural products, scaffold,  
37 scaffold tree, scaffold network, fragmentation, chemical space, clustering

38

## 39 Introduction

### 40 Scaffold concept and applications

41 Molecular scaffolds, defined as the core structures of molecules and also referred to as  
42 chemotypes or frameworks in some studies, are a concept used in many areas of chemistry.  
43 In drug design, the scaffold of a molecule is considered the main structure that determines its  
44 shape and places the functional moieties into the right positions to interact with the target. For  
45 this reason, developing new drug molecules with different cores but similar biological activities  
46 has been termed “scaffold hopping” [1, 2]. Combinatorial chemistry makes use of the concept

47 in designing compound libraries by substituting a set of scaffolds with combinations of different  
48 side chains. And structures in chemical patents are often defined analogously as Markush  
49 structures [3]. The intuitive chemical scaffold concept can also be utilised for classification  
50 purposes, especially in natural product (NP) research [4-7]. In cheminformatics, scaffold-  
51 based approaches can be applied for the analysis of high-throughput screening (HTS) data  
52 [6-10], mapping and visualising chemical spaces [5, 11], or even train-test splits of molecular  
53 data sets for machine learning projects [12].

54 Another application of scaffold-based methods is identifying privileged substructures in active  
55 molecules or NP that can be used as lead structures in the development of new drugs [5, 13-  
56 19]. Within NP chemical space, macrocyclic structures or cyclic peptides are of specific  
57 interest for these medicinal chemistry purposes [20-23].

58

## 59 Scaffold approaches in cheminformatics

60 The first general definition of a molecular scaffold was the Murcko framework developed by  
61 Bemis and Murcko in 1996 [11]. According to this concept, a scaffold consists of all the rings  
62 in a molecule and the non-cyclic chains connecting them, called linkers. Excluded from the  
63 scaffold are all terminal side chains. In addition to the Murcko framework definition  
64 representing molecular properties like atomic elements and bond multiplicities, Bemis and  
65 Murcko introduced a more abstract representation that reduced each atom in the framework  
66 to a simple graph node and each bond to a simple graph vertex, called graph framework or  
67 archetype. The authors used their framework definitions to assess the structural diversity of a  
68 set of drug molecules.

69 In addition to ignoring all non-cyclic molecules, the Murcko framework has one major  
70 drawback: small changes in the ring structure or the addition of a cyclic substituent, e.g. a  
71 benzene ring in drug design or a sugar moiety in NP research, can lead to very similar  
72 molecules not being grouped together due to non-equivalent scaffolds. Therefore, multiple

73 approaches have been developed for organising molecular scaffolds in a graph-based  
74 structure to relate similar scaffolds to each other and to create a systematic scaffold hierarchy  
75 [4, 5, 7, 9, 10, 13, 24-27].

76 Early work to this end was done by Xu and Johnson, who developed multiple concepts of  
77 dissecting Murcko frameworks into constituting ring systems or abstracting them into reduced  
78 representations. They used these concepts to assign molecular equivalence numbers to  
79 molecular structures and thus classify them within chemical libraries [25].

80 Wilkens et al. in their hierarchical scaffold clustering (HierS) approach [10] use a scaffold  
81 definition similar to Murcko frameworks but additionally include all atoms that are directly  
82 attached to rings and linkers via multiple bonds. Non-cyclic molecules are taken into  
83 consideration as well and are assigned scaffolds based on their multiple bonds. To build a  
84 scaffold hierarchy, the original scaffold extracted from each molecule is dissected into its  
85 smaller parent scaffolds first. This is done by generating all smaller scaffolds that can result  
86 from the stepwise removal of ring systems, i.e. isolated single rings or fused multiple rings that  
87 share bonds or atoms, from the original scaffold. After the removal of one ring, linker atoms  
88 that have become side chains are also removed. The process is finished when only the  
89 individual ring systems are left. Via a substructure search, the scaffold hierarchy is constructed  
90 in the second step by linking parent and child scaffold if the smaller parent scaffold is a  
91 substructure of the bigger child scaffold. In the end, a tree-like hierarchy results with the  
92 individual ring systems as roots at the top, and their combinations in more complex scaffolds  
93 on the following levels. A scaffold that is not a single ring system has multiple parents in the  
94 hierarchy.

95 While HierS overcomes most limitations of the Murcko framework approach and is a good first  
96 attempt for scaffold classification, it also has some disadvantages: ring systems are not split  
97 into their constituting single rings, which can be especially problematic when studying complex  
98 ring systems of NP where the approach may be too coarse-grained. In addition, child scaffolds  
99 are linked to multiple parents in the hierarchy, which is a multi-class assignment that is often  
100 undesirable for classification tasks.

101 The latter drawback of HierS is addressed in the structural classification of natural products  
102 (SCONP) approach by Koch et al. [5] that uses the same structural scaffold definition (apart  
103 from again ignoring linear molecules) but differs in its hierarchy construction routine. One  
104 major difference is that scaffolds are not dissected here. Only the directly extracted, original  
105 scaffolds of the studied molecule set are used to construct their relations in a tree-like fashion.  
106 A more complex scaffold is linked to only a single parent scaffold that is selected from all  
107 possible parent scaffolds representing substructures of the child following a set of chemical  
108 rules. These take characteristics of the parent scaffolds into account like hetero atom count,  
109 size, and frequency in the studied dataset. This last aspect makes the approach dataset-  
110 dependent, which can lead to problems in classification tasks.

111 A combination of scaffold dissection and single-parent assignments through chemical  
112 prioritisation rules is the scaffold tree approach by Schuffenhauer et al. [7]. As a first step,  
113 scaffolds are extracted from the given molecules according to the Murcko framework definition  
114 but additionally including all atoms connected via a double-bond to ring or linker atoms in the  
115 scaffolds. These elements are included as well to preserve correct hybridisation and structural  
116 alignment of the scaffold atoms. Via an iterative removal of rings, smaller parent scaffolds are  
117 created from the original child scaffolds. Ring perception for the removal is based on a smallest  
118 set of smallest rings (SSSR) approach. This way, ring systems sharing atoms or bonds  
119 between multiple rings are not considered as one entity but dissected into their constituting  
120 rings as well. One important aspect about the scaffold tree approach is the application of 13  
121 chemical prioritisation rules at every ring removal step. Following these rules, only one  
122 terminal ring is specifically selected for removal and only one possible parent scaffold created  
123 at every scaffold dissection step. The term "terminal" indicates that the removal does not result  
124 in a disconnected scaffold structure. The specific prioritisation rules take only molecular  
125 characteristics of the rings, like size, hetero atom count, and aromaticity, into account and aim  
126 at removing the less characteristic, peripheral rings first to extract the characteristic, central  
127 parent scaffold. The scaffold dissection process continues until only one ring remains. When  
128 studying a collection of molecules, their original scaffolds and sets of created parent scaffolds

129 are arranged in a hierarchy tree, the scaffold tree. Single-ring scaffolds form the roots and  
130 more complex scaffolds are placed at the higher levels. Due to the linear scaffold dissection  
131 process using the prioritisation rules, every child scaffold in the hierarchy is exclusively  
132 assigned to only one parent scaffold. Therefore, the scaffold tree represents a hierarchical,  
133 deterministic, and unique classification of chemical scaffolds. Unlike SCONP, it is dataset-  
134 independent because it does not consider the frequency of a scaffold in the studied collection.  
135 In conclusion, the scaffold tree is a useful tool for scaffold-based classification and  
136 visualisation of large compound sets and can be successfully employed to identify active  
137 scaffolds in HTS data and promising candidates for drug development [6-8, 28-30].

138 By definition of prioritisation rules, Schuffenhauer et al. intended to create a chemically intuitive  
139 classification system which opposes a classification focussing on pharmacophoric elements  
140 [7]. Also, its capability to identify biologically active substructural motives is limited because its  
141 exploration of possible parent scaffolds is limited due to the prioritisation rules. For this reason,  
142 Varin et al. introduced the concept of scaffold networks [9], where scaffolds are extracted and  
143 dissected analogously but without the application of prioritisation rules. In this way, every  
144 possible parent scaffold is generated for a given original scaffold and the resulting hierarchy,  
145 the scaffold Network, contains multi-parent relationships between its nodes. Varin et al.  
146 generated considerably more active scaffolds in primary screening data using scaffold  
147 networks compared to the scaffold tree approach. The reason for this is the exhaustive  
148 enumeration of parent scaffolds which leads to the scaffold network containing significantly  
149 more scaffolds than a scaffold tree. Additionally, a scaffold is not linked to all parent scaffolds  
150 that are substructures of it in the scaffold tree, only to the one determined as its characteristic  
151 core. As a consequence, this scaffold may be regarded to be less active.

152 The scaffold network approach explores the scaffold space more exhaustively and supports  
153 the identification of areas that a specific compound set does not cover. In addition, more virtual  
154 scaffolds can be identified, i.e. scaffolds that are only generated as a result of scaffold  
155 dissection and do not appear directly as original scaffolds in the given molecular structures.

156 When studying a compound set linked to bioactivity data, these structures are usually of high  
157 interest when appearing frequently in active molecules [29, 30].

158 On the other hand, scaffold networks can become large and complex with a comparably small  
159 number of molecules, which makes it difficult to visualise them. When linked to bioactivity data,  
160 Varin et al. suggest to only include islands of relevant, active scaffolds in the display.

161 As a conclusion, scaffold trees are generally more suitable for a complete visualisation and  
162 overview of the defining motives and structural classes in a limited compound set. Whereas  
163 scaffold networks can be seen as more helpful for analysing compound sets linked with  
164 bioactivity data to reasonably limit the display and identify active substructural motives [9].

165 An even more extensive scaffold network approach was published recently by Manelfi et al.,  
166 named "Molecular Anatomy" [26]. While the aforementioned approaches mostly rely on one  
167 single scaffold definition, respectively, nine different scaffold types of different abstraction  
168 levels were introduced here. All of them can be dissected analogously into parent scaffolds  
169 and linked in a network representation. This way, common substructure patterns can be  
170 identified on a higher abstraction level than with scaffold networks and more relevant similar  
171 compounds determined. This may be helpful for analysing HTS data or preparing structure  
172 activity relationship (SAR) studies of scaffolds and their side chains. However, this type of  
173 scaffold network including also more abstract scaffold representations has an even stronger  
174 tendency to grow very quickly with an increasing number of included structures and hence to  
175 quickly become unfathomable without sensibly limiting the display.

176 An analogous scaffold tree-like approach to hierarchical clustering based on Xu and Johnson's  
177 more abstract scaffold representations was published by Medina-Franco et al. [24]. Here,  
178 scaffolds are not dissected into parent scaffolds but clustered in a tree structure based on  
179 scaffold representations with a lower chemical resolution at each higher level.

180 An inherently different approach to scaffold generation and clustering are methods based on  
181 analog series. Here, no *a priori* scaffold definition like the Murcko framework is applied.  
182 Instead, all structures in a given set are grouped into analog series based on methods like  
183 matched molecular pairs with additionally deriving precursor structures using the RECAP

184 (Retrosynthetic Combinatorial Analysis Procedure) [31] rules. This way, the scaffolds  
185 extracted as representatives of an analog series take synthetic accessibility into account, an  
186 important aspect in medicinal chemistry but mostly ignored in the approaches above. The  
187 analog series and their representative scaffolds can be visualised by R-group tables, mapped  
188 into coordinate-based chemical space [32], annotated with activity information to support SAR  
189 studies, or used to extract favourable lead structures for drug design campaigns [17, 18, 33-  
190 38].

191

## 192 Open implementations

193 Most of the original software tools implementing the scaffold-based approaches described in  
194 the previous section have not been published openly (Murcko frameworks, SCONP, scaffold  
195 tree) or are not findable anymore (HierS). As a result, a number of open re-implementations  
196 and more advanced, versatile software for scaffold analyses has been developed and  
197 published [28, 29, 39-44].

198 The first open software application that implemented a scaffold tree was Scaffold Hunter [28,  
199 29, 39]. Starting as a tool mainly for generating and visualising scaffold trees, it has evolved  
200 into a multi-functional cheminformatics platform for visual data analysis. By default, the  
201 prioritisation rules are applied as published by Schuffenhauer et al. [7], but they can be  
202 customised by the user or even turned off completely. Varin et al. used the latter option to  
203 generate their scaffold networks using Scaffold Hunter [9]. The rich-client application is  
204 implemented in Java and employs the Chemistry Development Kit (CDK) [45-47] for  
205 cheminformatics tasks.

206 The open command-line tool Scaffold Network Generator [41] was designed to generate both,  
207 scaffold trees and scaffold networks. It lacks the extensive visualisation functionality of  
208 Scaffold Hunter but can therefore be integrated into automated analysis workflows that do not  
209 require human interaction. Scaffold Network Generator was implemented in Java as well and



210 employs the CDK and Open Babel [48] cheminformatics toolkits. Unfortunately, it cannot be  
211 found at the internet address given in the original publication anymore.

212 The cheminformatics toolkit RDKit [49] recently integrated an extensive scaffold network  
213 functionality into its range of capabilities [40]. The module named "rdScaffoldNetwork"  
214 primarily offers the generation of scaffold networks based on a HierS-like scaffold dissection  
215 (no splitting of fused rings). Custom fragmentation rules can be added in the form of reaction  
216 SMARTS [50]. In addition, more abstract atom- and bond-generic scaffold representations can  
217 be generated. The new functionality has been employed in a study evaluating different  
218 approaches to automate chemical series classifications in medicinal chemistry [51].

219 These three open software tools for scaffold-based analyses are only a limited number of  
220 examples for many more such tools developed in the past years [42-44].

221

## 222 Motivation

223 Structural scaffold analyses are relevant in diverse areas of cheminformatics, e.g. clustering,  
224 visualisation of chemical spaces, and SAR analyses [4-10, 26, 30]. Hence, numerous open  
225 software tools for such purposes have been developed [28, 29, 39-44]. The popular  
226 cheminformatics toolkit RDKit even integrated scaffold functionalities into its core modules.  
227 For the Chemistry Development Kit, only the generation of Murcko frameworks is currently  
228 available [52]. Outside core CDK, there is no open scaffold software library exclusively based  
229 on CDK to use in workflows and software based on the toolkit. Scaffold Hunter implemented  
230 its scaffold functionalities as part of a software application, and they cannot be easily extracted  
231 from it. Scaffold Network Generator is based on CDK but on Open Babel as well and not  
232 findable anymore.

233 Here, we present Scaffold Generator, an open, stand-alone Java library for scaffold  
234 functionalities based on CDK, to fill this void. It offers the generation of scaffold trees and

235 scaffold networks with comprehensive additional scaffold-related functionalities. An integration  
236 into the main CDK modules is intended.

237

## 238 Implementation

239 The Scaffold Generator library was implemented in Java version 17 and is based on the  
240 Chemistry Development Kit (CDK) version 2.8. The openly available source code can be found  
241 on GitHub: <https://github.com/Julian-Z98/ScaffoldGenerator>. With Scaffold Generator,  
242 different scaffold representations can be extracted from given molecules, dissected into parent  
243 scaffolds in multiple ways, and organised in scaffold trees and networks. These can be  
244 visualised using the GraphStream library version 2.0 [53, 54].

245

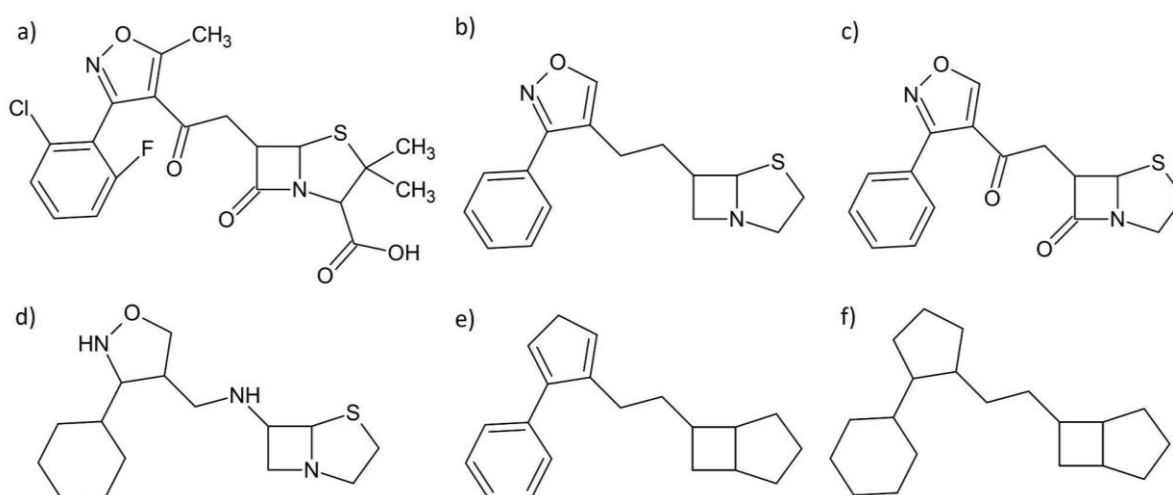
## 246 Available Functionalities

### 247 Scaffold types

248 Molecules are passed to Scaffold Generator as instances implementing the central CDK  
249 molecular structure representation, the *IAtomContainer* interface [55]. From these, molecular  
250 scaffolds can be extracted according to different scaffold definitions available. These include  
251 the Murcko framework and the scaffold definition used in most of the established approaches,  
252 like HierS or the scaffold tree. It is based on Murcko frameworks but additionally includes all  
253 atoms connected to ring or linker atoms via double-bonds [7, 10]. In Scaffold Generator, this  
254 has been extended to all atoms connected via non-single bonds to cyclic or linker atoms.  
255 Higher bond orders than 2 are considered rare in such configurations but they influence the  
256 hybridisation and structural configuration of the scaffold as strongly as exocyclic or exolinker  
257 double-bonds. Another crucial aspect to consider here is the synthetic accessibility of the  
258 represented scaffolds that is significantly influenced by the presence or absence of exocyclic

259 or exolinker multi-bonds. Additionally, two more abstract scaffold representations taken from  
260 Molecular Anatomy are available in Scaffold Generator: basic framework and basic wireframe  
261 [26]. Similar abstracted scaffold definitions have been described in earlier works as well, like  
262 the graph framework by Bemis and Murcko (analogous to basic wireframe) or the aryl cyclic  
263 system by Xu and Johnson (analogous to basic framework), but the naming was chosen here  
264 in analogy to Molecular Anatomy. A fifth scaffold type was analogously termed elemental  
265 wireframe. Here, all bonds are abstracted to single bonds, but hetero atoms are preserved  
266 (Figure 1). For the creation of scaffolds of all types, the CDK *MurckoFragmenter* class [52] is  
267 used internally and the extracted Murcko framework is post-processed according to the  
268 chosen scaffold type if necessary. If a given molecular structure has no rings, no scaffold can  
269 be extracted and an empty *IAtomContainer* instance is returned.

270



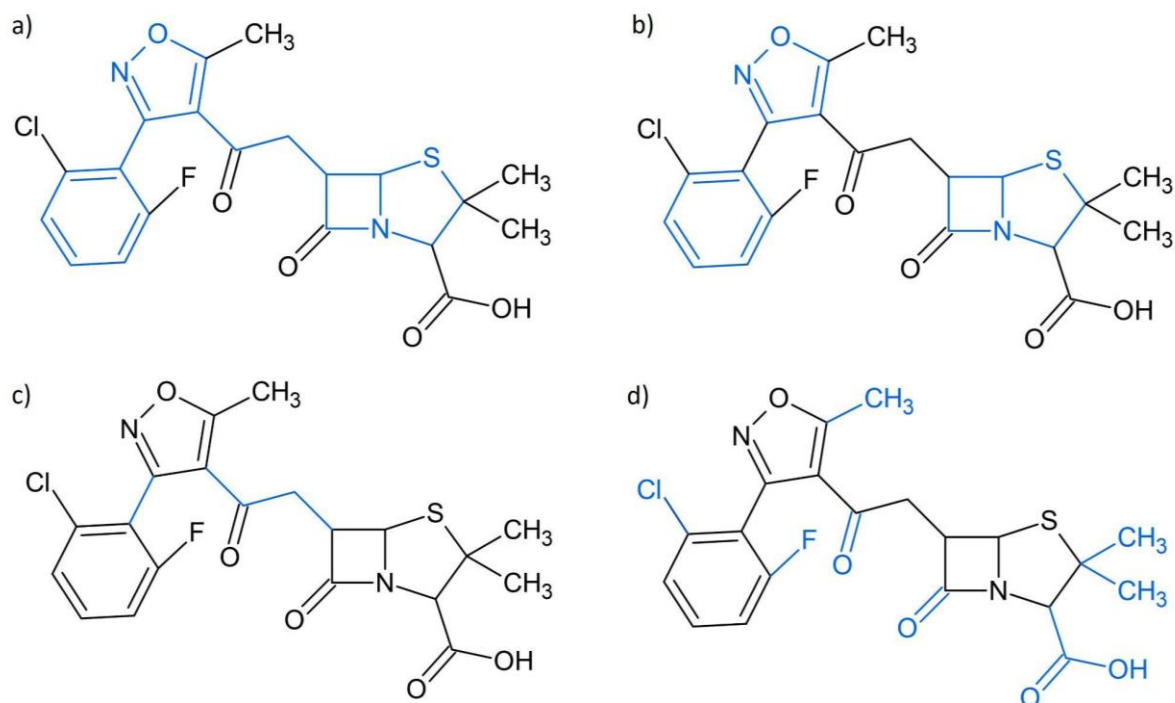
271

272 **Figure 1: Different scaffold types available in Scaffold Generator.** A) Unaltered structure  
273 of the antibiotic agent flucloxacillin (PubChem CID 21319). B) Murcko framework of  
274 flucloxacillin. C) Scaffold of flucloxacillin. D) Elemental wireframe of flucloxacillin. E) Basic  
275 framework of flucloxacillin. F) Basic wireframe of flucloxacillin.

276

277 Another functionality of Scaffold Generator is to return the building blocks of scaffolds, i.e.  
278 rings and linkers, separately. The terminal side chains excluded from the scaffold structure  
279 can also be extracted (Figure 2).

280



281

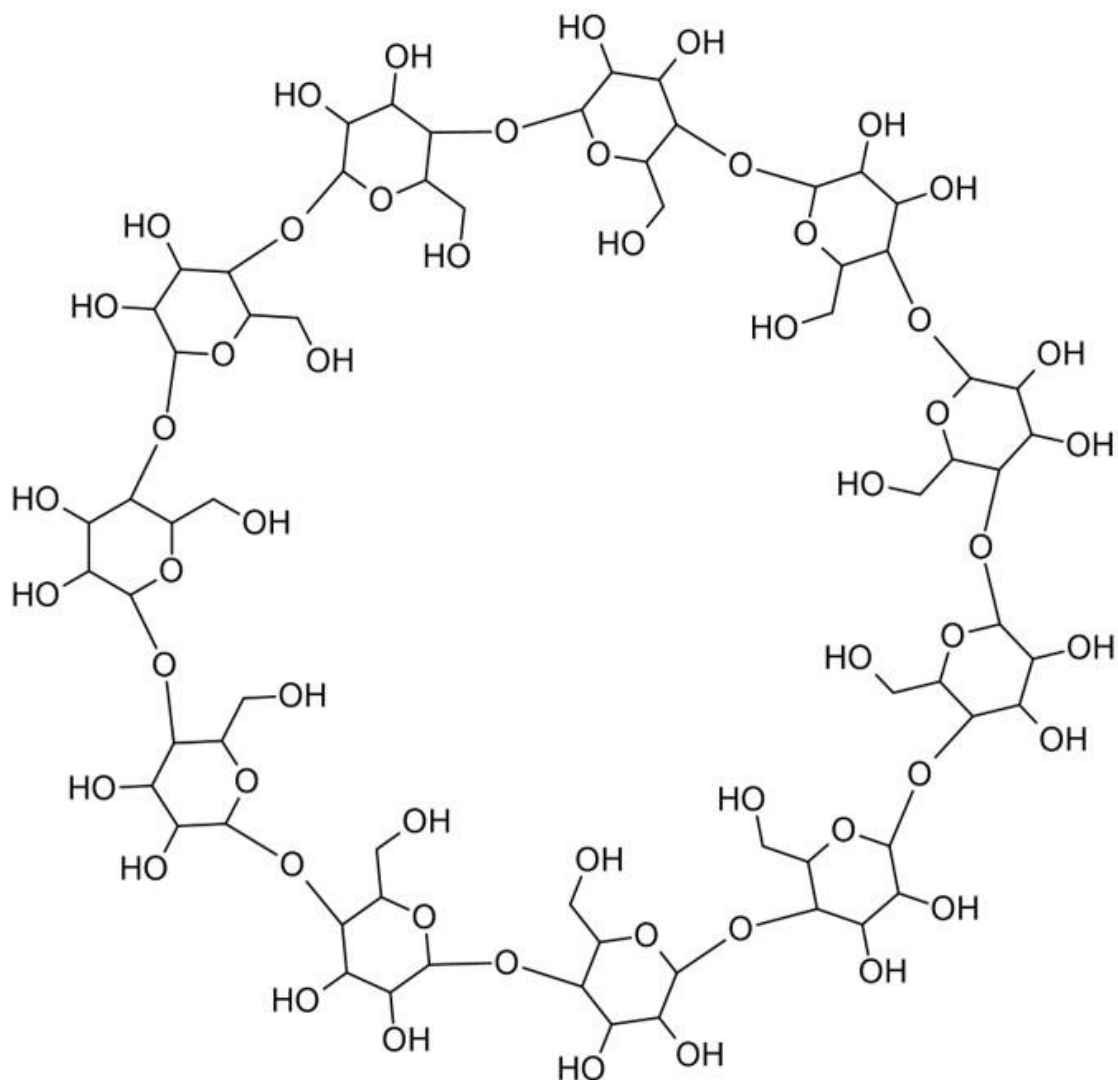
282 **Figure 2: Dissection of scaffolds into building blocks.** A) Flucloxacillin with its Murcko  
283 framework marked in blue. B) Rings of flucloxacillin marked in blue. It is important to note that  
284 the fused ring system on the right would be split into its two constituting rings in the structure  
285 set returned by the described routine of Scaffold Generator. C) Linkers of flucloxacillin marked  
286 in blue. D) Terminal side chains of flucloxacillin marked in blue.

287

## 288 Ring detection

289 Scaffold Generator dissects fused ring systems, i.e. rings that share bonds or atoms, into their  
290 constituting separate rings. This is the case not only when returning scaffold building blocks  
291 but also for the generation of parent scaffolds (see below). Internally, the CDK *Cycles.relevant*  
292 cycle finder algorithm is employed for ring detection. This algorithm detects the logical union

293 of all smallest sets of smallest rings (SSSR, also minimum cycle basis, MCB) in the given  
294 molecule [56, 57]. This way, fused ring systems are not detected as one entity, but their  
295 constituting cycles are detected separately. The *Cycles.relevant* cycle finder was chosen for  
296 Scaffold Generator to be in accordance with the original scaffold tree implementation [7]. But  
297 in rare cases, this cycle detection algorithm identifies too many rings in a given molecule,  
298 defined as more rings than there are atoms in the structure. One example is the natural product  
299 (NP) CNP0103752, taken from the COCONUT [58] database (Figure 3). Since the overarching  
300 ring connecting all 11 glycosidic rings in the structure can be detected on many different paths,  
301 *Cycles.relevant* detects 2059 rings here. In cases like this, i.e. more rings are detected than  
302 there are atoms in the molecule, Scaffold Generator uses the algorithm *Cycles.mcb* instead,  
303 which identifies one single set of SSSR/MCB instead of the logical union of all possible ones  
304 [56, 57]. In CNP0103752, it detects a more useful number for this application of 12 cycles.  
305



306

307 **Figure 3: Rings of CNP0103752 taken from COCONUT.** The CDK *Cycles.relevant* algorithm  
308 identifies 2059 rings here while *Cycles.mcb* detects 12.

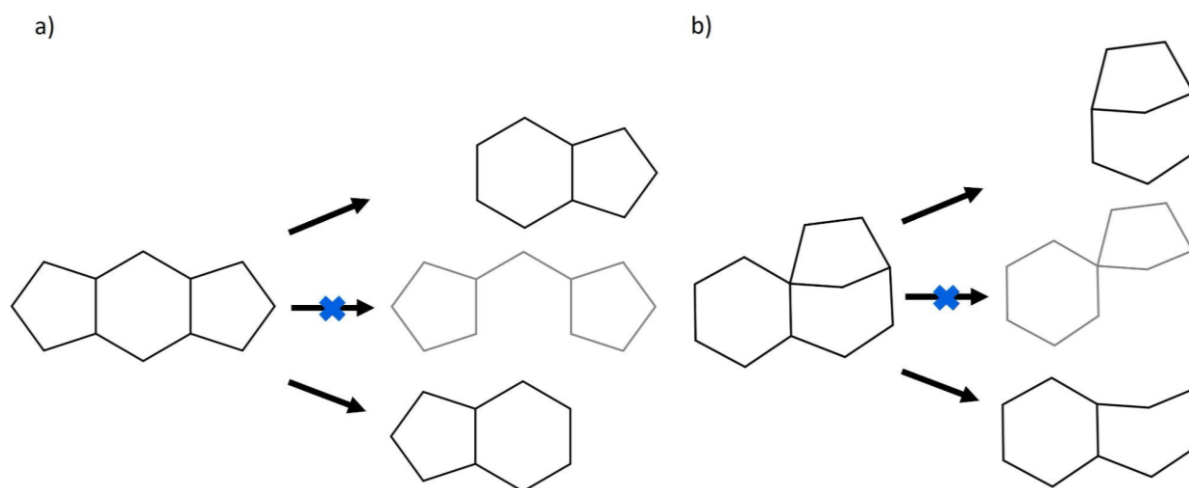
309

### 310 Ring removal

311 In the parent scaffold generation routines (see below), only rings adhering to a set of criteria  
312 are considered for removal at the individual dissection steps. The first requirement is that a  
313 ring needs to be terminal, i.e. its removal must not result in a disconnected scaffold structure.  
314 This is checked internally by removing all atoms and bonds constituting the respective ring  
315 from the scaffold, discarding potential side chains that were connected to it, e.g. when the

316 scaffold structural definition is used, and assessing whether the structure does not consist of  
317 multiple disconnected parts afterwards. If it does, the ring in question is not deemed terminal  
318 and hence not removable. This routine of checking for terminal rings has two major  
319 consequences: Internal rings that could be removed without resulting in a disconnected  
320 structure by turning some of their atoms and bonds into linker structures are still not considered  
321 terminal (Figure 4a). Secondly, the removal of rings from a scaffold cannot result in an  
322 artificially created spiro-ring system in Scaffold Generator (Figure 4b). Such cases are  
323 described in the original scaffold tree publication [7] and the fifth prioritisation rule there is  
324 intended to prevent them if other rings can be removed first. But they are possible in general  
325 and would appear in a set of all possible parent scaffolds. Because the conversion of ring  
326 atoms to linker atoms and the artificial creation of spiro-ring systems are chemically non-  
327 intuitive when generating parent scaffolds, these possibilities have been excluded in Scaffold  
328 Generator.

329



330

331 **Figure 4: Impossible parent scaffolds in Scaffold Generator.** A) Dodecahydro-s-indacene  
332 (PubChem CID 13214318) representing an example scaffold cannot be dissected in a way  
333 that turns former ring atoms into linker atoms in the created parent scaffold. B)  
334 Tricyclo[7.2.1.01,6]dodecane (PubChem CID 12758808) representing an example scaffold  
335 cannot be dissected in a way that creates a parent scaffold with a spiro-ring system which was  
336 not there in the molecule before.

337

338 Another requirement to consider a ring for removal is that it must contain at least one atom  
339 that is not part of another ring as well. This criterion is adopted from the original scaffold tree  
340 publication [7]. Here, the authors explain it with the example of adamantane. Using a ring  
341 detection algorithm that identifies the logical union of all SSSR in a structure, four rings are  
342 identified here and no atom is part of only one of them (compare Schuffenhauer et al. [7]  
343 Scheme 2). Hence, the removal of one ring is not possible because its atoms and bonds that  
344 are part of other rings as well are generally preserved in the Scaffold Generator ring removal  
345 routines. Structures like adamantane are therefore not dissected at all.

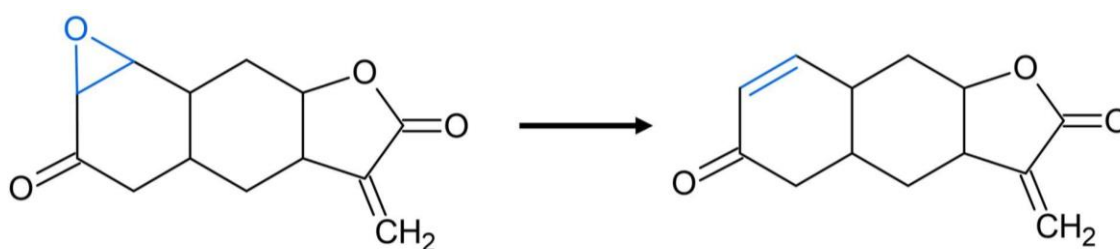
346 A similar case of structures that cannot be dissected are specific fused aromatic systems, i.e.  
347 aromatic rings that share the same atom with at least two other rings. When removing an  
348 aromatic ring sharing a bond with another ring, Scaffold Generator turns the shared bond into  
349 a double-bond to preserve the correct hybridisation of the formerly shared atoms in the  
350 remaining ring. In arrangements where the aromatic ring to remove shares an atom with at  
351 least two other rings, this double-bond insertion is not possible without violating valence rules.  
352 Such structures are not dissected as a consequence. This behaviour follows the ring removal  
353 algorithm described in the original scaffold tree publication (compare Schuffenhauer et al. [7]  
354 Scheme 3). But Scaffold Generator makes one addition here: In the original scaffold tree, this  
355 double-bond insertion is only done if an aromatic ring is fused to a non-aromatic ring and the  
356 aromatic ring is removed. In Scaffold Generator, it is also done if the remaining ring is aromatic  
357 as well. This addition has been made to preserve hybridisations and aromaticity in the  
358 remaining ring and to ensure that aromatic ring systems, if they can be dissected, are  
359 decomposed into parent scaffolds that can always be represented as valid contributing  
360 structures (as opposed to resonance hybrids). As a consequence, Scaffold Generator does  
361 not dissect most fused aromatic ring systems, e.g. pyrene. In these systems, most rings  
362 cannot be removed without altering hybridisations and bond orders in the remaining ones. And  
363 since a partial dissection does not appear reasonable because it would not produce  
364 meaningful parent scaffolds, these structures are not dissected at all. A possible future



365 extension to Scaffold Generator could be a routine that extracts meaningful parent scaffolds  
366 from fused aromatic systems, e.g. a benzene ring as root scaffold from pyrene and similar  
367 structures.

368 Another specially treated system are rings of size three containing one hetero atom that share  
369 the bond opposite to the hetero atom with another ring (Figure 5). When rings like this are  
370 removed, the shared bond is turned into a double bond to produce the precursor structure the  
371 hetero atom was most likely added to. This special case is described in the first ring removal  
372 prioritisation rule by Schuffenhauer et al. [7] but is part of the general ring removal routine of  
373 Scaffold Generator. This deviation from the original implementation does not influence the  
374 parent scaffold generation according to the scaffold tree prioritisation rules but is important to  
375 note for the enumerative generation of all possible parent scaffolds (see below).

376



377

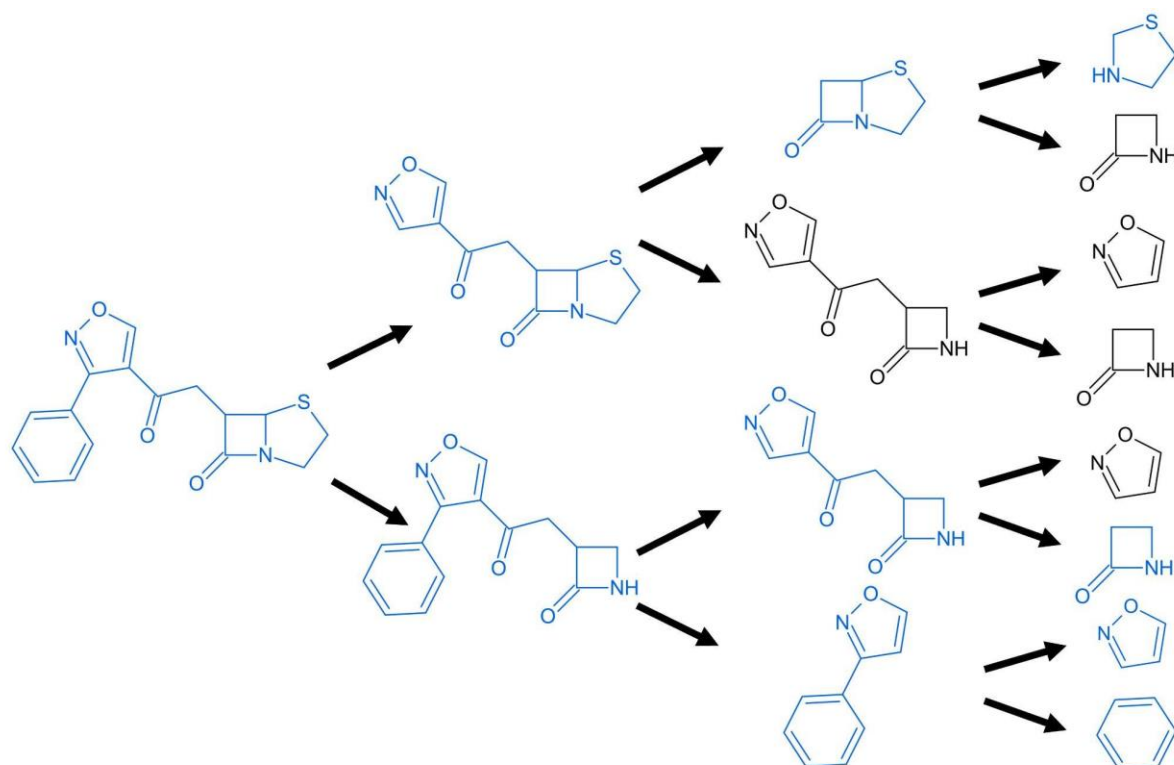
378 **Figure 5: Removal of 3-membered hetero cycles.** If the oxirane ring marked in blue is  
379 removed from himeyoshin (COCONUT CNP0151718) during parent scaffold generation, the  
380 bond shared with the cyclohexanone ring is turned into a double bond.

381

## 382 Scaffold trees and networks

383 Using Scaffold Generator, extracted molecular scaffolds can be dissected in different ways.  
384 The first one, as described above, is to decompose it into the constituting building blocks, i.e.  
385 rings and linkers. Another option is the enumerative removal that generates all possible parent  
386 scaffolds. At every iteration step, each ring adhering to the criteria listed above is removed  
387 separately to produce the resulting parent scaffold. This is repeated until only single-ring

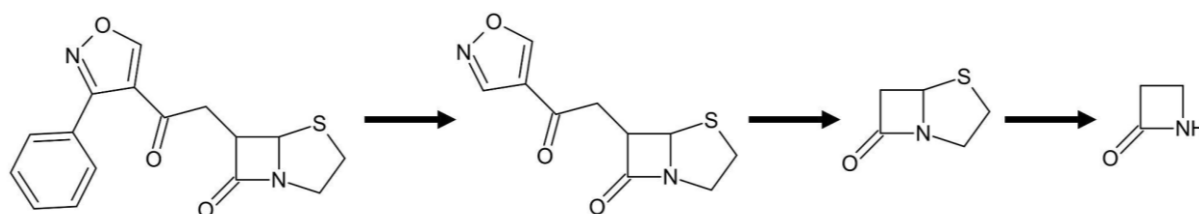
388 scaffolds remain, or no ring is removable anymore. These final scaffolds are called the root  
389 scaffolds. All generated parent scaffolds are substructures of the original scaffold. An example  
390 for the enumerative removal is shown in Figure 6. This routine can be applied to a given  
391 molecule and it returns a list with all possible parent scaffolds plus the original scaffold of the  
392 molecule. Parent scaffolds generated multiple times in the enumerative removal are returned  
393 only once. This scaffold dissection routine is the basis for generating scaffold networks. The  
394 dissection result of a single molecule can already be represented as a scaffold network by  
395 returning it as the corresponding data structure instead of a list.  
396



397  
398 **Figure 6: Enumerative parent scaffold generation of flucloxacillin.** Conceptual depiction  
399 of the enumerative parent scaffold generation routine applied to the scaffold of flucloxacillin  
400 (on the left). All possible parent scaffolds that can be created through the removal of a terminal  
401 ring are created. Marked in blue are all structures that are returned by the routine, indicating  
402 that structures occurring multiple times are still returned only once.  
403

404 Scaffold Generator implements the 13 chemical prioritisation rules that are applied in the  
405 original scaffold tree publication to specifically select only one parent scaffold at every scaffold  
406 dissection step [7]. In principle, these rules are applied to select only one ring removal path  
407 from all possible ones that are pursued in the enumerative removal (compare Figure 6). Only  
408 a few minor changes have been done to the original rules and underlying routines as reported  
409 above. Additionally, the final tie-breaking rule has been adapted to use unique SMILES  
410 representations [59, 60] as produced by the CDK, instead of canonical ones. From a given  
411 molecular structure, Scaffold Generator can generate a list of all parent scaffolds resulting  
412 from the Schuffenhauer dissection routine, plus the original scaffold (Figure 7). It produces the  
413 structures that can be used to build a scaffold tree in the second step. As with scaffold  
414 networks, a scaffold tree can already be constructed from a single molecule as well.

415



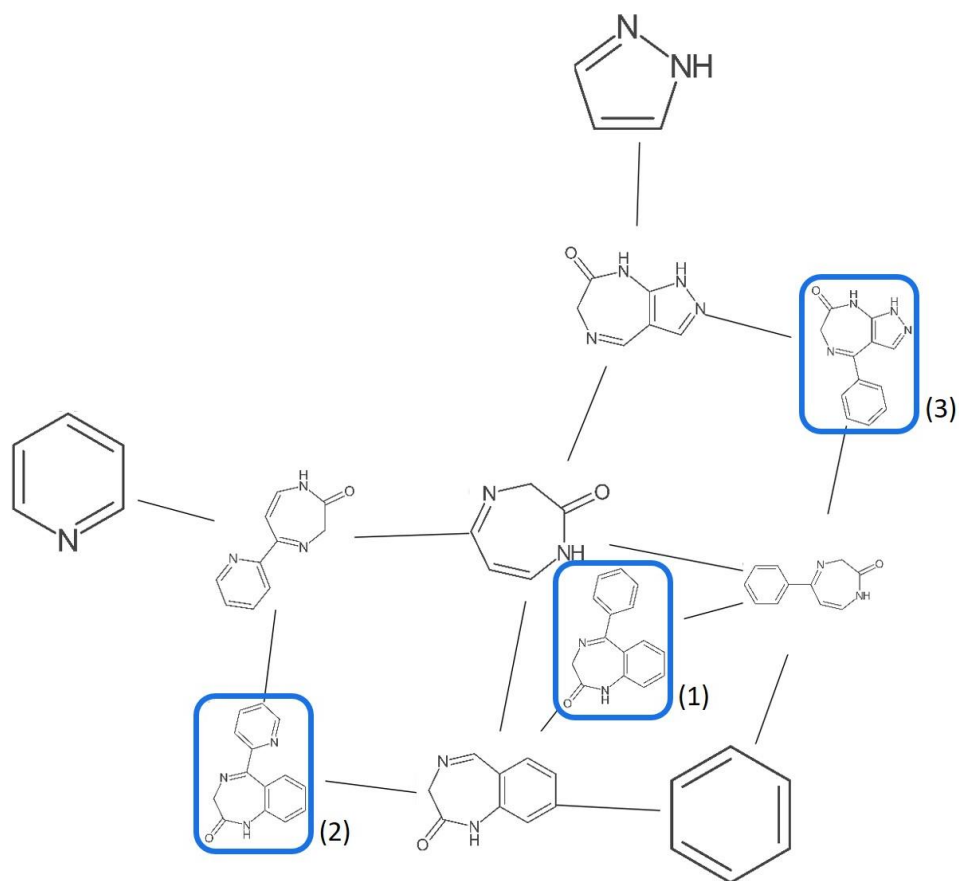
417 **Figure 7: Schuffenhauer parent scaffold generation of flucloxacillin.** Conceptual  
418 depiction of the parent scaffold generation routine employing the Schuffenhauer prioritisation  
419 rules applied to the scaffold of flucloxacillin (on the left). The rules are used to select only one  
420 parent scaffold out of all possible ones at every dissection step.

421

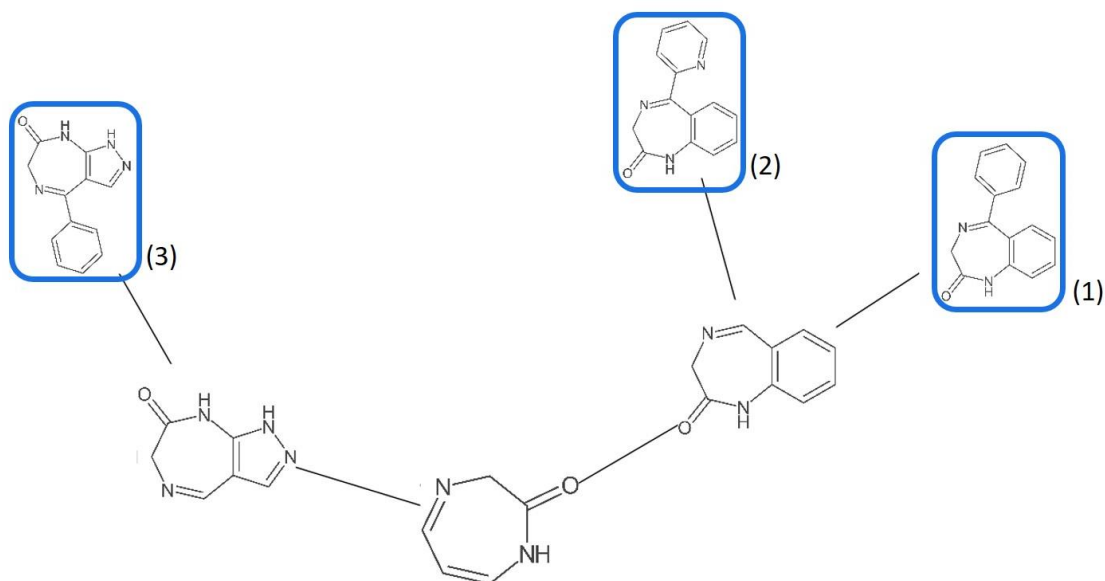
422 The main functionality of Scaffold Generator is the construction of scaffold trees and networks  
423 from given molecule collections (Figure 8). In the first step, the first molecule in the given  
424 collection is dissected into its parent scaffolds and the result is used to build the starting point  
425 of the desired structure. One by one, the remaining molecules are decomposed as well and  
426 their original scaffolds and parent scaffolds added to the tree or network if they are not already  
427 part of it. Scaffold Generator implements data structures that manage the graph nodes  
428 representing scaffolds and their parent-child connections as edges in scaffold trees and

429 networks. Both graphs are subdivided into levels with the root scaffolds on level 0 and their  
430 child scaffolds on the consecutive levels. The leaves are formed by the original scaffolds of  
431 the given molecules. But it is important to note that lower levels down to the roots can contain  
432 original scaffolds as well, e.g. when single-ring molecules are part of the given molecular set.  
433 The merging routines that are employed in the construction of a tree or network to add more  
434 scaffolds to it are also accessible after the final structures have been returned.  
435

a)



b)



436

437 **Figure 8: Scaffold network and tree depicted with the Scaffold Generator GraphStream**

438 **visualisation.** The scaffold network (a) and scaffold tree (b) of diazepam (PubChem CID

439 3016) (1), bromazepam (PubChem CID 2441) (2), and zolazepam (PubChem CID 35775) (3)

440 are displayed side-by-side for direct comparison (original scaffolds marked in blue). All three  
441 compounds are diazepinenones, a class of anxiolytics. The scaffold tree correctly identifies  
442 the diazepinenone ring as root scaffold of all three structures. But the scaffold network  
443 additionally reveals that diazepam (1) shares two-ring parent scaffolds with both the other  
444 structures, respectively. It also shows that the benzene ring is shared by all three compounds  
445 as well.

446  
447 The scaffold tree and network structures differ in some aspects: In scaffold trees, each node  
448 has only one parent node. This results from the Schuffenhauer scaffold dissection where a  
449 scaffold produces only one parent scaffold in each step. In scaffold networks, on the other  
450 hand, a node can have multiple parents since a scaffold usually produces multiple parent  
451 scaffolds in each step during the enumerative removal.

452 Another distinct aspect of scaffold trees is that only those molecules with their original  
453 scaffolds and parent scaffolds can be combined in one tree that share the same root scaffold.  
454 This is the scaffold (usually a single-ring scaffold) which results as parent scaffold in the final  
455 step of the Schuffenhauer dissection. It is unambiguously determined by the prioritisation  
456 rules. Scaffold Generator compiles the generated scaffolds of multiple molecules in one  
457 scaffold tree instance if they have the same root scaffold. If molecules with different root  
458 scaffolds are given in the molecule set, multiple scaffold tree instances will be created and  
459 returned in a list, termed scaffold forest in the nomenclature of Scaffold Generator. In the  
460 construction of scaffold networks, only one parent scaffold, i.e. at least one ring, needs to be  
461 shared between two molecules to be able to combine them in one network. But the scaffold  
462 network data structure of Scaffold Generator is also able to handle multiple disconnected  
463 graphs of scaffolds in one instance, unlike the scaffold tree structure.

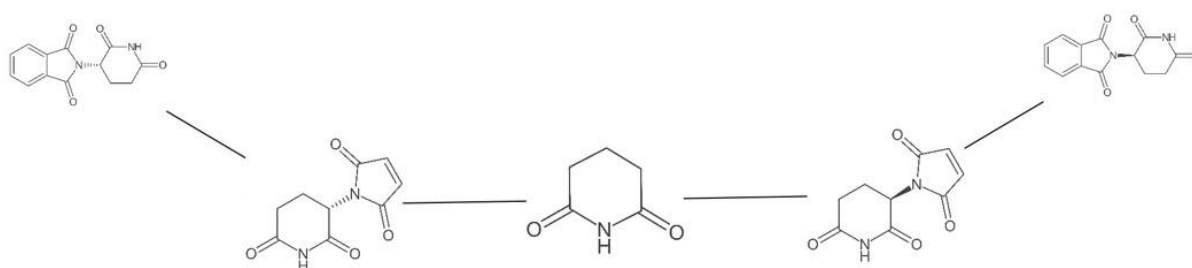
464 The tree and network data structures can generate an adjacency matrix representation of  
465 themselves that can be used for export or visualisation. Scaffold Generator offers an initial  
466 visualisation functionality for scaffold trees and networks based on the GraphStream library.  
467 The two structures can be visualised as graphs in a Java Swing application window. A layout

468 algorithm attempts to place the nodes and edges as readable as possible but modifications to  
469 the layout can be done by dragging nodes. The display can also be zoomed and moved using  
470 key commands. Some figures in this publication have been created using the Scaffold  
471 Generator GraphStream display (Figures 8 and 9). While this visualisation was helpful during  
472 the development process for visual inspection and debugging, it is not considered powerful  
473 enough for real-world use cases and will most likely not be part of a CDK integration of Scaffold  
474 Generator. A scaffold hierarchy visualisation tool that might sprout from Scaffold Generator as  
475 a separate project would have to be very interactive, i.e. zoomable, draggable, and  
476 collapsable. Especially scaffold networks tend to grow very fast with the number of included  
477 molecules. Therefore, their display needs to be limited in a comprehensive way, e.g. by only  
478 visualising islands of active scaffolds as proposed by Varin et al. [9]. Scaffold trees can  
479 become big as well, but they have the advantage that one can look at only one tree out of the  
480 forest at a time since they are disconnected.

481 When a tree or network is constructed, a crucial step is querying whether a scaffold is already  
482 part of it. This matching is done using SMILES representations of the scaffolds. The default  
483 setting is to use unique SMILES with aromaticity encoding but without stereochemical  
484 information. This can be adjusted, e.g. to include stereochemistry. Scaffold Generator  
485 generally retains given stereochemical information during scaffold creation and dissection by  
486 transferring the CDK *IStereoElement* [61] objects to the newly created structures. But this only  
487 works if all defining elements of a stereo group, i.e. atoms and/or bonds, are still present in  
488 the generated substructures. Since in the majority of cases side chains define stereochemistry  
489 and stereochemical information is often not given or incomplete in molecular data sets, the  
490 consideration of given stereochemical information in tree or network construction is turned off  
491 per default as stated above. But it can be enabled for use cases where it is relevant (Figure  
492 9). Other molecular characteristics that can generally be taken into account or not (depending  
493 on the specific use case) for the determination of equivalence between two structures in  
494 cheminformatic analyses are tautomeric forms or protonation states, for example.

495 Standardising these structures if needed has to be done in a data curation protocol that is  
496 applied to the input structures before they are passed to Scaffold Generator.

497



498

499 **Figure 9: Scaffold tree with activated stereochemistry consideration.** The Scaffold tree  
500 of (+)-thalidomide (PubChem CID 75792, on the left) and (-)-thalidomide (PubChem CID  
501 92142, on the right) with activated stereochemistry consideration is shown in the Scaffold  
502 Generator GraphStream display. If the consideration of stereochemistry in tree building was  
503 turned off, both compounds would be sharing the same two-ring scaffold as well.

504

505 The instances representing scaffold nodes in the trees and networks contain structural  
506 information about their scaffold and have references to their parents in the hierarchies.  
507 Additionally, they preserve SMILES codes of their origin molecules, i.e. structures from the  
508 data set that possess the respective scaffold. These origins are subdivided into virtual and  
509 non-virtual ones. Non-virtual origin molecules are those that have the node scaffold as their  
510 original scaffold, e.g. their Murcko framework. Virtual origins on the other hand are molecules  
511 that generate the respective scaffold only through enumerative or Schuffenhauer dissection,  
512 i.e. it is one of their parent scaffolds. This concept has been introduced in Scaffold Generator  
513 based on the definition of virtual scaffolds described in the literature [29, 30]. This term denotes  
514 scaffolds that are not directly in the data set but only identified when parent scaffolds are  
515 generated. If a scaffold node has only virtual origins, it is a virtual scaffold in Scaffold  
516 Generator. When analysing the results of a high-throughput screening (HTS) campaign, virtual  
517 scaffolds can be of particular interest if many of their child scaffolds exhibit bioactivity. A



518 promising next step can be a second screening with a smaller library based on this scaffold  
519 because the first screen might have failed to include the true active scaffold structure.  
520 An annotation of scaffold nodes in trees or networks with e.g. bioactivity data can be achieved  
521 via the stored origin molecules as well. One way to do this is to deposit the (unique) SMILES  
522 representation of the molecules in the studied data set linked to the respective annotation in  
523 a map structure. After the hierarchy is generated, its nodes can be annotated through  
524 comparing the origin molecule SMILES codes with the previously compiled annotation map.  
525 This way, e.g. scaffold nodes could be coloured according to bioactivity [7] or the hierarchy  
526 display limited to active scaffolds [9] in a more advanced visualisation tool as proposed above.  
527 During the development of Scaffold Generator, it was decided against keeping the original  
528 *IAtomContainer* instances with their structures and properties as origin references in favour of  
529 only their SMILES representations to reduce random-access memory (RAM) consumption.

530

## 531 Aromaticity handling

532 Aromaticity information and detection is relevant in multiple Scaffold Generator functionalities.  
533 As stated above, when an aromatic ring is removed, bonds it shares with other rings are turned  
534 into double bonds in some cases to preserve hybridisations and aromaticity. Since this is not  
535 possible in all configurations, aromaticity information is also relevant in the determination of  
536 possibly removable rings (see above). And many fused aromatic ring systems, e.g. pyrene,  
537 are not dissected by Scaffold Generator as a result.  
538 Aromaticity information is also significant in two of the 13 scaffold tree prioritisation rules for  
539 parent scaffold determination, namely rule 7 "A Fully Aromatic Ring System Must Not Be  
540 Dissected in a Way That the Resulting System Is Not Aromatic Any More") and rule 11 "For  
541 Mixed Aromatic/Nonaromatic Ring Systems, Retain Nonaromatic Rings with Priority") [7]. The  
542 seventh rule makes it necessary to generate all possible parent scaffolds producible by the  
543 removal of one ring at the given dissection step and apply aromaticity determination to each  
544 of them to assess whether aromaticity was lost in the remaining ring(s). Because this

545 consumes a lot of computation time and aromaticity should be conserved in most cases  
546 through the double-bond insertion, the application of the seventh prioritisation rule can be  
547 turned off individually in Scaffold Generator.

548 Aromaticity determination in CDK and hence in Scaffold Generator is carried out by  
549 *Aromaticity* instances [62] constructed from the combination of an *ElectronDonation* model  
550 [63] and a *CycleFinder* algorithm [56]. The former defines which atom types can contribute  
551 how many electrons to the aromatic system and the latter determines the cycles that can form  
552 them. All aromaticity models in CDK loosely follow the Hückel rule heuristic [62]. The specific  
553 *Aromaticity* instance used in Scaffold Generator can be configured because different models  
554 are suited for different applications.

555 Since multiple intermediate steps in scaffold dissection rely on aromaticity information of  
556 specific substructures, an initial aromaticity detection is applied at the primary scaffold  
557 generation. And again at the end of a scaffold dissection process, a final aromaticity detection  
558 is applied to all generated parent scaffolds to make sure that the aromaticity information stored  
559 on the scaffold objects is in agreement with the returned structures. This last step might lead  
560 to cases where the same ring is not detected as aromatic in a smaller parent scaffold but in  
561 the bigger child scaffold in which it is a substructure. This is due to the cycle finder algorithms  
562 usually employed for aromaticity detection that are not SSSR-/MCB-based but also take cycles  
563 into account that span multiple rings of the molecule. It should be interpreted in the way that  
564 the ring in the parent scaffold gained aromaticity in the child scaffold through combination with  
565 other rings.

566 An additional option is to turn off aromaticity detection completely in all Scaffold Generator  
567 routines. This was implemented because this process takes a lot of time and makes the results  
568 of scaffold dissection routines dependent on mostly toolkit-specific and heuristic aromaticity  
569 models. If it is disabled, initially defined aromaticity information in the input structures is  
570 preserved.

571 It must also be noted here again that all aromaticity models in CDK are based on the Hückel  
572 rule, which is the most used heuristic for aromaticity determination but not the only one and

573 has a long list of exemptions. Furthermore, it is only a heuristic determination method for the  
574 concept of aromaticity, which is itself not uniquely defined [64-67].

575

## 576 Settings and options

577 **Table 1: Settings and options of Scaffold Generator.** The settings listed in this table  
578 together with their options and default values are available in Scaffold Generator to adjust its  
579 results to specific use cases.

Setting name	Options	Default
Scaffold mode	<ul style="list-style-type: none"> <li>- Scaffold</li> <li>- Murcko framework</li> <li>- Basic wireframe</li> <li>- Basic framework</li> <li>- Elemental wireframe</li> </ul>	Scaffold
Determine aromaticity	true/false	true
Aromaticity model	All combinations of <i>CycleFinder</i> and <i>ElectronDonation</i> instances available in CDK	<i>ElectronDonation.cdk</i> and <i>Cycles.cdkAromaticSet</i>
Retain only hybridisations at aromatic bonds	true/false	false
Rule seven applied ("A Fully Aromatic Ring System Must Not Be Dissected in a Way That the Resulting System Is Not Aromatic Any More" [7])	true/false	true
SMILES generator	All <i>SmilesGenerator</i> configurations available in CDK	<i>SmiFlavor.Unique</i> and <i>SmiFlavor.UseAromaticSymbols</i>

580

581 The functionalities and routines of Scaffold Generator can be adopted for various applications

582 by a variety of settings available (Table 1). Five different structural scaffold definitions can be

583 chosen for initial scaffold extraction and scaffold dissection (Figure 1). The default setting of

584 the scaffold mode setting is to use the scaffold including all atoms directly connected to rings

585 or linkers via non-single bonds.

586 Multiple steps in scaffold dissection and the construction of Scaffold trees and networks  
587 require the testing for equivalence of molecular structures. These include the enumerative  
588 generation of all possible parent scaffolds to avoid duplicates and the identification of  
589 equivalent scaffolds when merging trees or networks. In Scaffold Generator, this is done using  
590 CDK unique SMILES codes. To allow the user the definition of structural features taken into  
591 account at these steps, e.g. stereochemistry, isotopes, or aromaticity, the CDK  
592 *SmilesGenerator* [68] instance employed can be set externally. By default, stereochemistry  
593 and atomic masses are not encoded but aromaticity is. The set *SmilesGenerator* instance is  
594 also used to create SMILES codes for origin molecules of a respective scaffold stored on  
595 nodes of scaffold trees and networks. It is important to note here that molecular characteristics  
596 of the input molecules and resulting (parent) scaffolds, like protonation states or tautomeric  
597 forms, are taken by Scaffold Generator “as is”, or rather as they are represented in the chosen  
598 SMILES encoding. The only exemption is the detection of aromatic systems which is done on  
599 input structures by default. Therefore, users have to take care of preprocessing their input  
600 data sets according to their specific needs, e.g. standardising tautomeric forms and  
601 protonation states in all input molecules, before using Scaffold Generator.

602 Another option is to exclude or include the Schuffenhauer prioritisation rule 7. This rule makes  
603 it necessary to apply aromaticity detection to different parent scaffolds created for testing  
604 purposes. This procedure is time-consuming and might not lead to a definite decision in favour  
605 of one specific parent scaffold in most cases. But by default, it is activated to be in accordance  
606 with the originally published scaffold tree implementation [7].

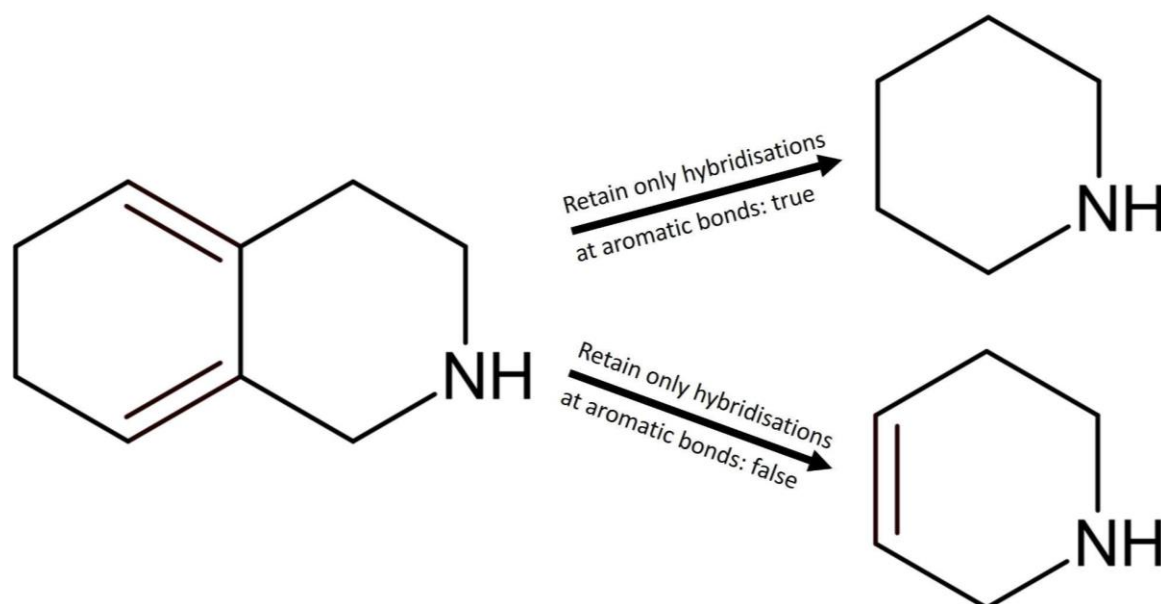
607 The aromaticity detection done in multiple steps of scaffold dissection (see above) can be  
608 configured by choosing which CDK aromaticity model is to be employed for this purpose. By  
609 default, aromaticity is determined using the *ElectronDonation.cdk* model and the  
610 *Cycles.cdkAromaticSet* cycle finder algorithm.

611 Additionally, aromaticity detection can be turned off completely in all routines to preserve initial  
612 aromaticity information of the input structures and make the results less dependent on specific

613 aromaticity models. If this is the case, rule 7 is automatically excluded from the Schuffenhauer  
614 prioritisation rules as well.

615 The fifth option of Scaffold Generator concerns post-processing after ring removal: As  
616 explained above, a double bond is inserted in some cases when an aromatic ring is removed  
617 to preserve hybridisation and aromaticity in the remaining ring(s) if possible. As an option, this  
618 insertion of double bonds can also be applied to non-aromatic systems wherever there are  
619 two  $sp^2$  hybridised atoms adjacent to a single bond that was previously shared between two  
620 rings. The bond is turned into a double bond if the two adjacent atoms would lose their  $sp^2$   
621 hybridisation because of the ring removal and if it is possible without violating valence rules  
622 (Figure 10).

623



624

625 **Figure 10: Parent scaffold of 1,2,3,4,6,7-hexahydroisoquinoline depending on the set**  
626 **value of the retain only hybridisations at aromatic bonds setting.** When the  
627 cyclohexadiene ring is removed from 1,2,3,4,6,7-hexahydroisoquinoline (PubChem CID  
628 89002720) in parent scaffold generation, the formerly shared bond with the piperidine ring is  
629 turned into a double bond if the retain only hybridisations at aromatic bonds setting is set to  
630 false. In this case, double bonds are always inserted if possible to preserve atom  
631 hybridisations in the remaining ring. If the setting is set to true, this is only done when an

632 aromatic ring is removed. In this case, no double bond is inserted in the remaining piperidine  
633 ring.

634

## 635 Software architecture

636 The central class of the Scaffold Generator library is *ScaffoldGenerator*. When instantiated,  
637 all available settings are set to their default values (Table 1) and can be adjusted using  
638 methods of the class. All main functionalities of Scaffold Generator described above can be  
639 accessed through an instance of the *ScaffoldGenerator* class, i.e. generation of scaffolds, their  
640 decomposition into building blocks, parent scaffold generation through enumerative or  
641 Schuffenhauer dissection, and the generation of scaffolds trees and networks. The two  
642 scaffold hierarchy structures are represented by a class of their own, respectively:  
643 *ScaffoldTree* and *ScaffoldNetwork*. Both extend the same base class,  
644 *ScaffoldNodeCollectionBase*, for basic functionalities and manage scaffold nodes as  
645 *TreeNode* or *NetworkNode* instances that both stem from the abstract base class  
646 *ScaffoldNodeBase*. These six classes manage scaffold structures, parent-child relationships  
647 of scaffold nodes, and origin molecule references. Trees and networks can be traversed and  
648 merged with instances of the same class, respectively. Scaffold trees can additionally be  
649 checked for validity, i.e. whether all nodes have parents, except the root node, and there is  
650 only one root node. Scaffold tree and network instances can also be exported as adjacency  
651 matrices along with scaffold structures for each represented node. This is utilised by the class  
652 *GraphStreamUtility* to display scaffold trees and networks in an interactive Java Swing  
653 application window with the GraphStream library.

654 The JUnit [69] test class *ScaffoldGeneratorTest* implements automatic tests for the basic  
655 Scaffold Generator routines, tests employing the GraphStream visualisation of scaffold trees  
656 and networks for visual inspection, and code examples for the application of Scaffold  
657 Generator. Another important set of test routines checks whether the Schuffenhauer  
658 prioritisation rules as implemented in Scaffold Generator are in accordance with the original

659 implementation, based on the examples given in the scaffold tree publication [7]. Furthermore,  
660 the COCONUT database is used to test the basic routines on a large set of natural product  
661 (NP) structures.

662 The class *PerformanceTest* represents a command-line application based on Scaffold  
663 Generator that can be used to assess its computational speed on a given structure data file  
664 (SDF). The results on COCONUT and DrugBank [70, 71] are presented in the “Results and  
665 discussion” section.

666

## 667 Results and discussion

668 A programming library for molecular scaffold functionalities named Scaffold Generator was  
669 implemented based on the Chemistry Development Kit (CDK). The openly available source  
670 code of Scaffold Generator can be found on GitHub: [https://github.com/Julian-](https://github.com/Julian-Z98/ScaffoldGenerator)  
671 [Z98/ScaffoldGenerator](https://github.com/Julian-Z98/ScaffoldGenerator). It can be utilised to extract different types of scaffolds from input  
672 molecules and dissect them further into parent scaffolds using an enumerative generation of  
673 all possible ones or a dissection according to the scaffold tree prioritisation rules. Additionally,  
674 the scaffolds and parent scaffolds can be arranged in scaffold trees and networks with these  
675 hierarchies being visualised.

676

### 677 Performance

678 Scaffold Generator can be packaged in a JAR file and used as a command-line application. It  
679 requires an SD file as input parameter and creates a performance snapshot of the main  
680 functionalities of Scaffold Generator with the given data set. First, all molecules are imported  
681 and stored in memory. From these, all structures having more than ten rings are discarded.  
682 This is done because they occur rather rarely but would influence the overall processing time  
683 disproportionately. No further filtering or preprocessing, e.g. removal of counter-ions or



684 elimination of duplicates, is done for the purpose of this performance snapshot and the  
685 following exemplary analyses. For an initial performance snapshot, all remaining molecules  
686 are processed according to the enumerative generation of parent scaffolds and the parent  
687 scaffold generation according to the scaffold tree prioritisation rules. Afterwards, the dataset  
688 is subdivided into equally large portions. The total number of fractions has to be specified in  
689 the second command-line parameter. In each following step, a growing number of created  
690 molecule subsets is combined and all included structures used to build a scaffold network and  
691 a scaffold forest, i.e. a set of scaffold trees. The number of molecules and the needed  
692 processing time is logged in every step. In the final step, all scaffolds in the network and the  
693 trees, respectively, and their frequencies determined based on their numbers of origin  
694 molecules are exported to an output file. The scaffold structures are exported as SMILES  
695 strings.

696 For this article, two performance snapshots were conducted. The first one was done on the  
697 DrugBank database containing drug molecules (DrugBank “all structures” downloaded on 8<sup>th</sup>  
698 November 2021). For comparison, the COCONUT NP database (downloaded on 1<sup>st</sup>  
699 December 2021) was analysed as well. Additionally, for some analyses, a subset of  
700 COCONUT containing 40,000 structures was compiled from the complete collection using the  
701 RDKit MaxMin algorithm implementation [49, 72]. All analyses were conducted on a  
702 workstation computer with an Intel(R) Xeon(R) Gold 6254 CPU (18 cores, 3.10 GHz) and 512  
703 GB RAM on a single core only (no multi-core parallelization). All Scaffold Generator settings  
704 were set to their default values.

705

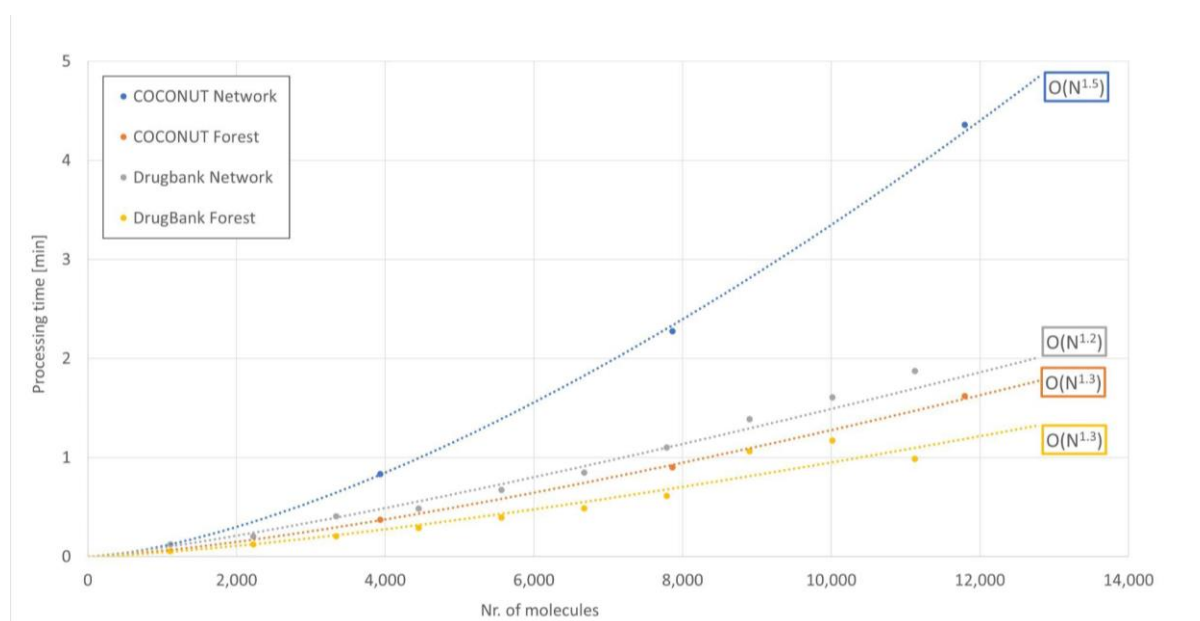
706 **Table 2: Performance snapshot of the mere parent scaffold generation routines applied**  
 707 **to COCONUT and DrugBank.**

	COCONUT	DrugBank
Initial number of molecules	406,747	11,172
Number of molecules after filtering (< 11 rings)	395,450	11,127
Schuffenhauer dissection total	1,211,063 ms (20 min)	27,656 ms (0.46 min)
Schuffenhauer dissection average per molecule	3 ms	2.5 ms
Enumerative dissection total	2,037,357 ms (34 min)	33,938 ms (0.57 min)
Enumerative dissection average per molecule	5 ms	3 ms

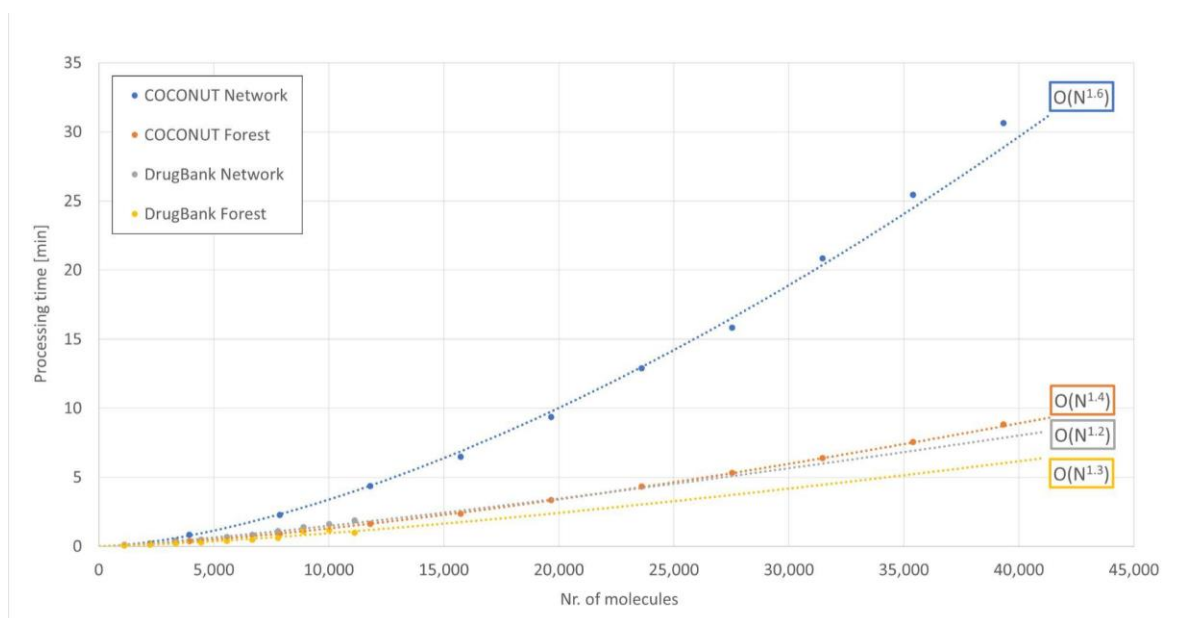
708  
 709 The complete COCONUT database contained 406,747 NP structures (Table 2). 11,297 of  
 710 these possessed 11 or more rings and were filtered. The remaining 395,450 NP were  
 711 subjected to the parent scaffold generation according to the Schuffenhauer rules, which took  
 712 1,211,063 ms (20 min). On average, the dissection of one COCONUT NP into its scaffold and  
 713 parent scaffolds according to the Schuffenhauer prioritisation rules took 3 ms. Generating all  
 714 possible parent scaffolds with the enumerative routine took 2,037,357 ms (34 min) for the  
 715 same molecule set. This is 5 ms per molecule on average.

716 The DrugBank data set of 11,172 molecules contained 45 structures with more than 10 rings  
 717 that needed to be filtered. The Schuffenhauer dissection of all structures took 27,656 ms  
 718 (0.46 min, 2.5 ms per molecule on average) and the enumerative parent scaffold generation  
 719 took 33,938 ms (0.57 min, 3 ms per molecule on average).

720 It is interesting to note that the enumeration of all possible parent scaffolds at every step  
721 required more computation time than the application of up to 13 prioritisation rules at every  
722 step. This was the case for NP as well as drug molecules which have less rings in general.  
723 The latter characteristic of drug molecules as opposed to NP is also considered the reason for  
724 the lower time it took on average to dissect the DrugBank structures. It must also be noted  
725 that these processes, the pure dissection of each molecule, scale linearly with the number of  
726 molecules and can be parallelised in multiple threads for further speed up.  
727



728 **Figure 11: Performance snapshot of scaffold forest and scaffold network construction**  
729 **in DrugBank range of molecule number.** The graph visualises the processing time it took  
730 to construct a scaffold forest or scaffold network depending on the number of input molecules  
731 taken from COCONUT or DrugBank. Exponential approximations have been applied to assess  
732 the scaling behaviour of the processes. The given range of the number of molecules is  
733 adjusted to the size of DrugBank (11,127 molecules).  
734  
735



736

737 **Figure 12: Performance snapshot of scaffold forest and scaffold network construction**

738 **in COCONUT subset range of molecule number.** The graph visualises the processing time

739 it took to construct a scaffold forest or scaffold network depending on the number of input

740 molecules taken from COCONUT or DrugBank. Exponential approximations have been

741 applied to assess the scaling behaviour of the processes. The given range of the number of

742 molecules is adjusted to the size of the curated COCONUT subset (39,324 molecules).

743

744 In a second step, it was measured how much time it took to construct scaffold forests and

745 networks from an increasing number of molecules taken from the COCONUT subset and

746 DrugBank, respectively. Figure 11 shows the results for the area of molecule number of

747 DrugBank (0 - 11,127 molecules) and Figure 12 for the area of the COCONUT subset (0 -

748 39,324 molecules). Exponential approximations show that the individual processes scaled

749 between  $O(N^{1.2})$  and  $O(N^{1.6})$ . This comparatively good scaling below a quadratic behaviour is

750 most likely due to the stepwise construction of the scaffold hierarchies that repeats the two

751 steps of scaffold dissection and integration for each molecule instead of generating all

752 scaffolds first and constructing the hierarchy later using substructure searches to establish

753 parent-child scaffold relationships.

754 Both, the generation of scaffold networks and trees from NP, scaled with higher exponents  
755 than the analogous processes for drug molecules, which can again be explained by the  
756 generally higher number of rings in the former class of compounds.

757 The generation of scaffold networks from NP structures scaled with the highest exponent.  
758 Since the number of scaffolds in a network grows faster than in a forest because more parent  
759 scaffolds are constructed for each molecule, it takes more time in network construction to  
760 integrate new molecules, i.e. their scaffolds. This traversal of the scaffold forest or network for  
761 the integration of new scaffolds is considered to be the algorithm step that dictates the scaling  
762 behaviour. In addition, this step would be more challenging to parallelise and speed up through  
763 multithreading because the same data structure would be accessed by all threads. The  
764 scaffold tree and network representations in Scaffold Generator are currently not implemented  
765 to be thread-safe, i.e. safe to use for concurrent modification.

766 According to the exponential approximation for the COCONUT subset of 40,000 NP  
767 structures, a scaffold network of up to 456,000 NP molecules could still be constructed in a  
768 single day using Scaffold Generator. The measured runtime for the complete COCONUT  
769 database of 395,450 compounds with less than 11 rings was 16.5 h (5 h for the construction  
770 of a scaffold forest). This is below the runtime of 19.2 h expected for this data set size  
771 according to the exponential function approximating the scaling behaviour of the COCONUT  
772 subset network generation. The underlying effect can be that with growing size of the network,  
773 less new scaffolds need to be integrated per newly added molecule. Here, one also has to  
774 take into account that the subset used for the performance and scaling snapshot was compiled  
775 using a diversity-preserving method [72]. This may have increased the effect even further.

776 The memory consumption of the scaffold tree and scaffold network constructed from the  
777 complete COCONUT database was below the 512 GB RAM available at all times but similar  
778 experiments on a machine with 256 GB failed.

779

780 Most frequent scaffolds in COCONUT and DrugBank

781 **Table 3: Numbers of resulting scaffolds in scaffold network and scaffold forest**  
782 **constructed from COCONUT and DrugBank.**

	COCONUT	DrugBank
Number of molecules after filtering (< 11 rings)	395,450	11,127
Number of scaffold network scaffolds	392,888	23,765
Number of scaffold trees	6,200	766
Number of scaffold tree scaffolds	173,526	10,716

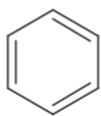
783

784 The Scaffold Generator command-line application logs the numbers of different scaffolds in  
785 network and forest built from the given data set and exports the scaffolds as SMILES  
786 representations with their frequencies as a final step. These scaffold numbers for COCONUT  
787 and DrugBank can be found in Table 3. The COCONUT scaffold network contained 392,888  
788 different (parent) scaffolds, while the DrugBank network contained 23,765. The COCONUT  
789 scaffold forest consisted of only 173,526 scaffolds distributed among 6,200 individual scaffold  
790 trees. For DrugBank, it was 10,716 scaffolds in 766 trees. According to these numbers, the  
791 enumerative parent scaffold generation produced more than twice as many scaffolds as the  
792 Schuffenhauer dissection. Using a classification by root scaffolds, the two data sets could be  
793 classified into a number of different classes according to the number of resulting scaffold trees.  
794 The 20 most frequent scaffolds in the COCONUT scaffold network and scaffold forest,  
795 respectively, as determined in this exemplary showcase analysis, are displayed in Figures 13  
796 and 14. The frequencies are given as numbers of origin molecules that produced the  
797 respective scaffold in parent scaffold generation or had it as an original scaffold. The  
798 frequencies for the network scaffolds correspond precisely to the number of molecules that  
799 possess the respective scaffold as a substructure, whereas the frequencies for the forest

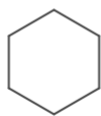
800 scaffolds correspond to the number of molecules that possess the scaffold as their most  
801 characteristic or central parent scaffold in one step of the Schuffenhauer dissection according  
802 to the prioritisation rules. Hence, 225,272 COCONUT molecules contain a benzene ring  
803 (Figure 13) but only in 29,258 molecules, it is the characteristic or central parent scaffold  
804 (Figure 14). Still, it is striking that the benzene ring is the most frequent root scaffold in the  
805 forest because some Schuffenhauer prioritisation rules explicitly assign a low relevance to it  
806 and favour its removal over that of other rings.

807 As could be expected, the first ranks in both charts are dominated by single-ring scaffolds,  
808 since they represent the final stage of scaffold dissection and have the most origin molecules,  
809 therefore. The first ranks are also dominated by 6-membered rings and parent scaffolds that  
810 are most likely resulting from the dissection of polyketides. The frequency of oxygen-  
811 containing scaffolds is higher than that of nitrogen, as can be expected for NP. The empty  
812 cells in both charts represent empty scaffolds, i.e. scaffolds of molecules that have no rings.  
813 Hence, 21,882 molecules in COCONUT do not possess any circular structures. Of 406,747,  
814 the share of linear molecules is low (5 %), but one should keep in mind that these structures  
815 are usually completely neglected in ring-based analyses like most scaffold methods.

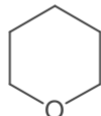
816



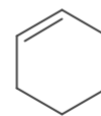
225272



78833



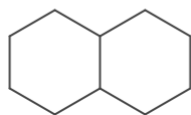
55954



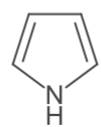
46338



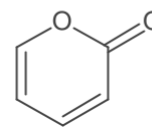
43493



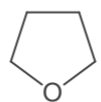
33118



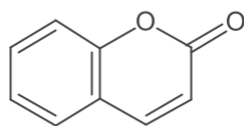
26504



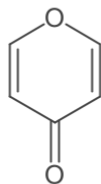
25308



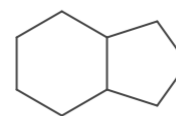
23899



23043



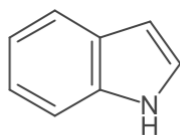
22395



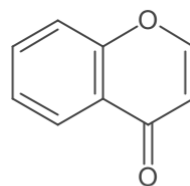
22136

No scaffold  
(linear molecules)

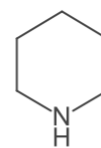
21882



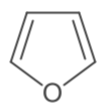
21143



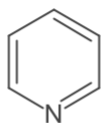
20164



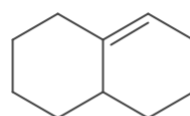
19982



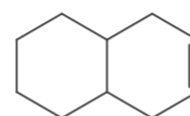
19866



18458



16250

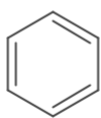


15495

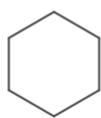
817



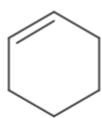
818 **Figure 13: 20 most frequent scaffold network scaffolds of COCONUT with their numbers**  
819 **of origin molecules.**  
820



29258



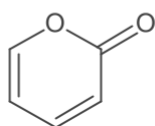
24497



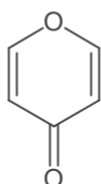
23897

No scaffold  
(linear molecules)

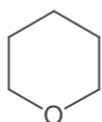
21882



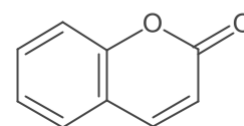
20008



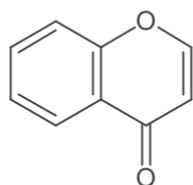
18492



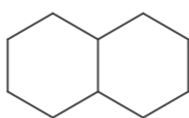
16645



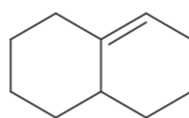
16482



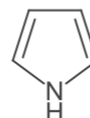
15724



12319



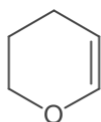
11843



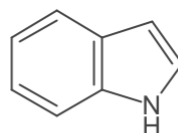
9294



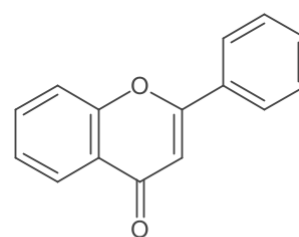
8855



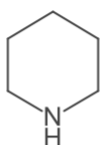
7914



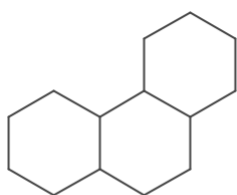
7172



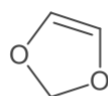
6808



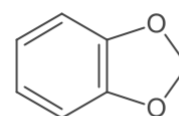
6762



6594



6307



6266

821

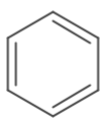
822 **Figure 14: 20 most frequent scaffold forest scaffolds of COCONUT with their numbers**  
823 **of origin molecules.**

824

825 Figures 15 and 16 analogously display the most frequent scaffolds of the created DrugBank  
826 scaffold network and scaffold forest. The first observation here is that the share of nitrogen  
827 hetero cycles is higher in these drug molecules than in NP structures. This has been reported  
828 before [73]. Also, the share of linear molecules (1,467 of 11,172, 13 %) is much higher than  
829 in NP. Benzene is again the most frequent scaffold in both analyses. But while it is by far the  
830 most frequent scaffold in the DrugBank network (6,578 origin molecules compared to 972 for  
831 the second most frequent scaffold, pyridine), its prominence is way lower in the forest (1,819  
832 origin molecules compared to 611 for pyrimidine in second place).

833 The core results of this showcase analysis comparing the most frequent NP and drug molecule  
834 scaffolds (i.e. commonness of benzene, oxygen as the dominant hetero atom in NP, nitrogen  
835 in drug molecules) are in general agreement with similar studies [15, 19, 74-76]. A significantly  
836 higher prevalence of aromatic scaffolds in drug molecules as opposed to NP that most of these  
837 studies report cannot be observed here. This stresses that the results presented here are only  
838 a proof of concept for the application of Scaffold Generator. A more detailed analysis would  
839 first of all need an extensive data curation pipeline to standardise input molecules or filter or  
840 mark duplicates between the two data sets. Furthermore, a more extensive analysis of  
841 physicochemical property distributions in the extracted scaffolds could be conducted.

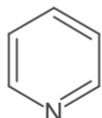
842



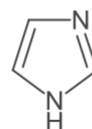
6578

No scaffold  
(linear molecules)

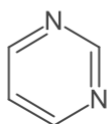
1467



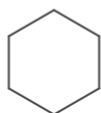
972



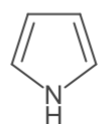
901



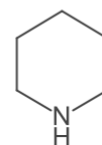
834



810



662



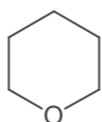
635



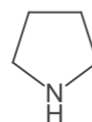
585



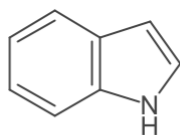
568



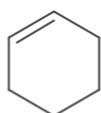
502



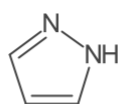
473



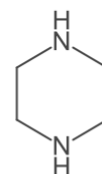
435



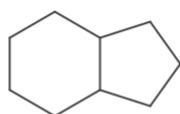
410



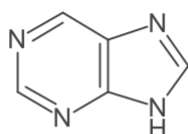
377



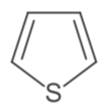
347



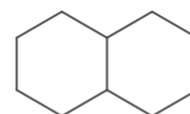
322



286



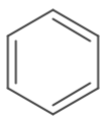
284



277

843

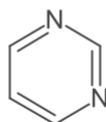
844 **Figure 15: 20 most frequent scaffold network scaffolds of DrugBank with their numbers**  
845 **of origin molecules.**  
846



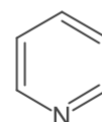
1819

No scaffold  
(linear molecules)

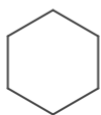
1467



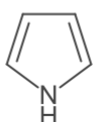
611



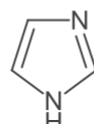
420



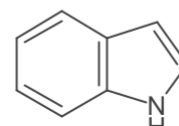
353



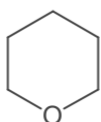
343



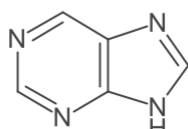
335



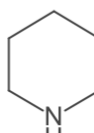
275



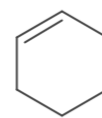
264



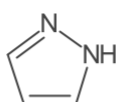
260



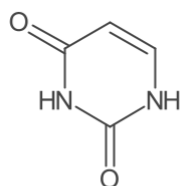
256



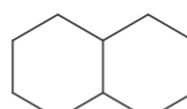
234



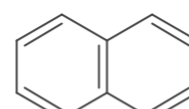
222



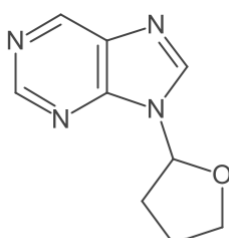
221



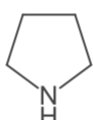
209



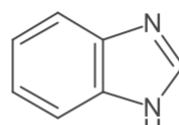
178



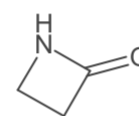
167



157



135



120

847

848 **Figure 16: 20 most frequent scaffold forest scaffolds of DrugBank with their numbers**  
849 **of origin molecules.**

850

851 This analysis of the most frequent scaffolds in COCONUT and DrugBank is only supposed to  
852 serve as a basic example for what kind of studies Scaffold Generator may be used. These  
853 results may also have been achieved through the mere dissection of scaffolds into parent  
854 scaffolds and a subsequent matching and counting of the resulting structures. With its ability  
855 to generate and represent scaffold networks and forests, Scaffold Generator may be applied  
856 to a wider variety of analyses like hierarchical classification and clustering, chemical space  
857 mapping, or HTS data interpretation. But for these, a more powerful visualisation than the  
858 existing GraphStream-based one would be very helpful.

859

## 860 Future Work

861 Scaffold Generator meets the need for an open, versatile, CDK-based library for scaffold  
862 functionalities that can be employed in software and workflows built upon this cheminformatics  
863 toolkit. To make it more accessible to potential users, an integration into the CDK core modules  
864 would be desirable since the toolkit would benefit from having more scaffold functionalities  
865 available. A corresponding request to the library maintainers has been made.

866 Another aspect that would make Scaffold Generator more applicable is a more powerful  
867 visualisation functionality than the currently available one based on the GraphStream library.

868 It should display the hierarchies in suitable layouts, i.e. a tree layout for scaffold trees and a  
869 similar layout for scaffold networks that arranges the network in its defined levels. The display  
870 should be draggable, zoomable, and collapsable. The latter aspect is especially important for  
871 scaffold networks that tend to grow very fast with the number of included molecules. For  
872 example, all scaffolds below a chosen node should be easily collapsable or only active islands  
873 of scaffolds should be displayed when bioactivity data is linked to the given molecules [9].

874 Especially the analysis of HTS data or the derivation of SAR insights would benefit from a

875 versatile scaffold hierarchy visualisation. To further support these analyses, methods to  
876 display scaffolds and their parent scaffolds hierarchically in a standardised, directly visually  
877 accessible way, like the work of Alex M. Clark [77], should be explored in future developments.  
878 Scaffold Generator can serve as core for a variety of scaffold-based functionalities.  
879 Classification, clustering, and scaffold-based fingerprints are possible applications that can be  
880 used in a second step for picking diverse training and test sets for machine learning models  
881 for example [12]. The concept of scaffolds and parent scaffolds as characteristic molecular  
882 fragments of molecules can help in the development of QSAR/QSPR models or computer-  
883 assisted structure elucidation. Applied to NP, scaffolds can serve as starting points for the  
884 creation of pseudo-NP that are regarded as promising candidates for new drug molecules [78,  
885 79]. Additionally, the study of macrocyclic structures in NP with existing scaffold  
886 methodologies and the development of new, specialised approaches for these structures are  
887 promising ways of identifying new drug candidates [20-23].  
888 Possible functional extensions of Scaffold Generator include the incorporation of more  
889 abstract scaffold representations, based on the work by Xu and Johnson [25], and the  
890 possibility to build scaffold networks or trees encompassing multiple scaffold definitions of  
891 varying chemical resolution, like in Molecular Anatomy [26] or the tree-like classification of  
892 Medina-Franco et al [24]. A major addition to the functionality of Scaffold Generator would be  
893 the inclusion of analog series based scaffold methodologies. Since these have demonstrated  
894 significant relevance in the past years, this addition must be considered.

895

## 896 Conclusion

897 An open, CDK-based, stand-alone Java library named Scaffold Generator has been  
898 developed to meet the need for scaffold functionalities in CDK-based workflows and software.  
899 It offers the extraction of different scaffolds, the dissection of scaffolds into building blocks,  
900 and the generation of parent scaffolds in two different ways. An enumerative parent scaffold



901 generation routine produces all parent scaffolds that can be created through the removal of  
902 terminal rings and forms the basis for scaffold networks. Alternatively, only characteristic or  
903 central parent scaffolds can be extracted according to the Schuffenhauer prioritisation rules  
904 that are used to build scaffold trees. Scaffold trees and networks can be internally represented  
905 as data structures and visualised in a basic display based on the GraphStream library. The  
906 generation of a scaffold network from more than 450,000 natural product structures can be  
907 achieved in a single day. A request for the integration of Scaffold Generator into the CDK core  
908 modules has been made and the process started. Scaffold Generator may serve as a starting  
909 point for diverse scaffold-based software tools, e.g. for clustering or fingerprint functionalities.  
910

## 911 List of abbreviations

912 CDK: Chemistry Development Kit  
913 CID: Compound IDentifier  
914 CNP: COCONUT Natural Product  
915 COCONUt: COLleCtion of Open Natural prodUcTs  
916 CPU: Central Processing Unit  
917 HTS: High-Throughput Screening  
918 JAR: Java ARchive  
919 MCB: Minimum Cycle Basis  
920 NP: Natural Product(s)  
921 QSAR/QSPR: Quantitative Structure Activity/Property Relationship  
922 R: Registered trademark  
923 RAM: Random-Access Memory  
924 RECAP: REtrosynthetic Combinatorial Analysis Procedure  
925 SAR: Structure Activity Relationship  
926 SD(F): Structure Data (File)

927 SMARTS: SMILES Arbitrary Target Specification

928 SMILES: Simplified Molecular Line Entry System

929 SSSR: Smallest Set of Smallest Rings

930

## 931 Availability and requirements

932 ● Project name: Scaffold Generator

933 ● Project home page: <https://github.com/Julian-Z98/ScaffoldGenerator>

934 ● Current version: v1.0.3

935 ● DOI of archived current version: <https://doi.org/10.5281/zenodo.7245473>

936 ● Operating system(s): Platform independent

937 ● Programming language: Java

938 ● Other requirements: Java v17 or higher, Maven v4 or higher, CDK v2.8 (fetched by  
939 Maven), GraphStream v2.0 (fetched by Maven), JUnit v4.13.2 (fetched by Maven)

940 ● Licence: GNU Lesser General Public Licence (LGPL) v2.1

941 ● Any restrictions to use by non-academics: None

942

## 943 Declarations

### 944 Availability of data and materials

945 Data and software are freely available under the LGPL v2.1 licence. The source code of  
946 Scaffold Generator is available on GitHub at <https://github.com/Julian-Z98/ScaffoldGenerator>.

947

## 948 Competing interests

949 AZ is co-founder of GNWI - Gesellschaft für naturwissenschaftliche Informatik mbH,  
950 Dortmund, Germany.

951

## 952 Funding

953 This work was supported by the Carl-Zeiss-Foundation.

954

## 955 Authors' contributions

956 JS designed and supervised the study. JS and JZ designed, tested, applied, and validated the  
957 features of Scaffold Generator and wrote the paper. JZ developed the Java code. CS and AZ  
958 conceived the study and acquired the funding. All authors read and approved the final  
959 manuscript.

## 960 Acknowledgements

961 The authors would like to thank the communities that created the open software libraries  
962 utilised in the development of Scaffold Generator, especially the CDK community for support  
963 during this process. Further thanks go to the authors of the scaffold tree publication,  
964 Schuffenhauer et al., for describing their prioritisation rules in a detailed way that allowed  
965 straightforward reimplementations.

966

## 967 References

968 [1] G. Schneider, W. Neidhart, T. Giller, and G. Schmid, "Scaffold-Hopping" by Topological

- 969 Pharmacophore Search: A Contribution to Virtual Screening', *Angew. Chem. Int. Ed.*, vol.  
970 38, no. 19, pp. 2894–2896, Oct. 1999, doi: 10.1002/(SICI)1521-  
971 3773(19991004)38:19<2894::AID-ANIE2894>3.0.CO;2-F. **[cito:citesAsAuthority]**
- 972 [2] H.-J. Böhm, A. Flohr, and M. Stahl, 'Scaffold hopping', *Drug Discov. Today Technol.*, vol.  
973 1, no. 3, pp. 217–224, Dec. 2004, doi: 10.1016/j.ddtec.2004.10.009.  
974 **[cito:citesAsAuthority]**
- 975 [3] E. A. Markush, 'Pyrazolone dye and process of making the same', USA101506316, Aug.  
976 26, 1924 Accessed: Jan. 07, 2022. [Online]. Available:  
977 <https://pdfpiw.uspto.gov/.piw?PageNum=USA101506316&docid=01506316&IDKey=83E682D73B35&HomeUrl=http%3A%2F%2Fpatft.uspto.gov%2Fnetacgi%2Fnph-Parser%3FSect1%3DPTO1%2526Sect2%3DHITOFF%2526p%3D1%2526u%3D%2Fnetahtml%2FPTO%2Fsrchnum.html%2526r%3D1%2526f%3DG%2526l%3D50%2526d%3DPALL%2526s1%3D1506316.PN.%2526OS%3D%2526RS%3D>  
978  
979  
980  
981  
982 **[cito:citesAsAuthority]**
- 983 [4] A. Schuffenhauer and T. Varin, 'Rule-Based Classification of Chemical Structures by  
984 Scaffold', *Mol. Inform.*, vol. 30, no. 8, pp. 646–664, Aug. 2011, doi:  
985 10.1002/minf.201100078. **[cito:citesAsAuthority]**
- 986 [5] M. A. Koch *et al.*, 'Charting biologically relevant chemical space: A structural  
987 classification of natural products (SCONP)', *Proc. Natl. Acad. Sci.*, vol. 102, no. 48, pp.  
988 17272–17277, Nov. 2005, doi: 10.1073/pnas.0503647102. **[cito:citesAsAuthority]**  
989 **[cito:discusses]**
- 990 [6] A. Schuffenhauer, N. Brown, P. Ertl, J. L. Jenkins, P. Selzer, and J. Hamon, 'Clustering  
991 and Rule-Based Classifications of Chemical Structures Evaluated in the Biological  
992 Activity Space', *J. Chem. Inf. Model.*, vol. 47, no. 2, pp. 325–336, Mar. 2007, doi:  
993 10.1021/ci6004004. **[cito:citesAsAuthority]**
- 994 [7] A. Schuffenhauer, P. Ertl, S. Roggo, S. Wetzel, M. A. Koch, and H. Waldmann, 'The  
995 Scaffold Tree – Visualization of the Scaffold Universe by Hierarchical Scaffold  
996 Classification', *J. Chem. Inf. Model.*, vol. 47, no. 1, pp. 47–58, Jan. 2007, doi:

- 997 10.1021/ci600338x. **[cito:citesAsAuthority] [cito:usesMethodIn] [cito:discusses]**
- 998 [8] T. Varin *et al.*, 'Compound Set Enrichment: A Novel Approach to Analysis of Primary  
999 HTS Data', *J. Chem. Inf. Model.*, vol. 50, no. 12, pp. 2067–2078, Dec. 2010, doi:  
1000 10.1021/ci100203e. **[cito:citesAsAuthority]**
- 1001 [9] T. Varin, A. Schuffenhauer, P. Ertl, and S. Renner, 'Mining for Bioactive Scaffolds with  
1002 Scaffold Networks: Improved Compound Set Enrichment from Primary Screening Data',  
1003 *J. Chem. Inf. Model.*, vol. 51, no. 7, pp. 1528–1538, Jul. 2011, doi: 10.1021/ci2000924.  
1004 **[cito:citesAsAuthority] [cito:usesMethodIn] [cito:discusses]**
- 1005 [10] S. J. Wilkens, J. Janes, and A. I. Su, 'HierS: Hierarchical Scaffold Clustering Using  
1006 Topological Chemical Graphs', *J. Med. Chem.*, vol. 48, no. 9, pp. 3182–3193, May 2005,  
1007 doi: 10.1021/jm049032d. **[cito:citesAsAuthority] [cito:discusses]**
- 1008 [11] G. W. Bemis and M. A. Murcko, 'The Properties of Known Drugs. 1. Molecular  
1009 Frameworks', *J. Med. Chem.*, vol. 39, no. 15, pp. 2887–2893, Jan. 1996, doi:  
1010 10.1021/jm9602928. **[cito:citesAsAuthority] [cito:usesMethodIn] [cito:discusses]**
- 1011 [12] J. Simm *et al.*, 'Splitting chemical structure data sets for federated privacy-preserving  
1012 machine learning', *J. Cheminformatics*, vol. 13, no. 1, p. 96, Dec. 2021, doi:  
1013 10.1186/s13321-021-00576-2. **[cito:citesAsAuthority]**
- 1014 [13] S. L. Schreiber, 'Target-Oriented and Diversity-Oriented Organic Synthesis in Drug  
1015 Discovery', *Science*, vol. 287, no. 5460, pp. 1964–1969, Mar. 2000, doi:  
1016 10.1126/science.287.5460.1964. **[cito:citesAsAuthority]**
- 1017 [14] D. S. Tan, 'Diversity-oriented synthesis: exploring the intersections between chemistry  
1018 and biology', *Nat Chem Biol*, vol. 1, no. 2, pp. 74–84, Jul. 2005, doi:  
1019 10.1038/nchembio0705-74. **[cito:citesAsAuthority]**
- 1020 [15] Y. Chen, C. Rosenkranz, S. Hirte, and J. Kirchmair, 'Ring systems in natural products:  
1021 structural diversity, physicochemical properties, and coverage by synthetic compounds',  
1022 *Nat. Prod. Rep.*, vol. 39, no. 8, pp. 1544–1556, 2022, doi: 10.1039/D2NP00001F.  
1023 **[cito:citesAsAuthority] [cito:agreesWith]**
- 1024 [16] S. Stone, D. J. Newman, S. L. Colletti, and D. S. Tan, 'Cheminformatic analysis of natural

- 1025 product-based drugs and chemical probes', *Nat. Prod. Rep.*, vol. 39, no. 1, pp. 20–32,  
1026 2022, doi: 10.1039/D1NP00039J. **[cito:citesAsAuthority]**
- 1027 [17] D. Dimova, D. Stumpfe, Y. Hu, and J. Bajorath, 'Analog series-based scaffolds:  
1028 computational design and exploration of a new type of molecular scaffolds for medicinal  
1029 chemistry', *Future Science OA*, vol. 2, no. 4, p. FSO149, Dec. 2016, doi: 10.4155/foa-  
1030 2016-0058. **[cito:citesAsAuthority]**
- 1031 [18] D. Stumpfe, D. Dimova, and J. Bajorath, 'Computational Method for the Systematic  
1032 Identification of Analog Series and Key Compounds Representing Series and Their  
1033 Biological Activity Profiles', *J. Med. Chem.*, vol. 59, no. 16, pp. 7667–7676, Aug. 2016,  
1034 doi: 10.1021/acs.jmedchem.6b00906. **[cito:citesAsAuthority]**
- 1035 [19] P. Ertl and T. Schuhmann, 'Cheminformatics Analysis of Natural Product Scaffolds:  
1036 Comparison of Scaffolds Produced by Animals, Plants, Fungi and Bacteria', *Mol. Inf.*, vol.  
1037 39, no. 11, p. 2000017, Nov. 2020, doi: 10.1002/minf.202000017.  
1038 **[cito:citesAsAuthority] [cito:agreesWith]**
- 1039 [20] E. Marsault and M. L. Peterson, 'Macrocycles Are Great Cycles: Applications,  
1040 Opportunities, and Challenges of Synthetic Macrocycles in Drug Discovery', *J. Med.*  
1041 *Chem.*, vol. 54, no. 7, pp. 1961–2004, Apr. 2011, doi: 10.1021/jm1012374.  
1042 **[cito:citesAsAuthority]**
- 1043 [21] P. Ermert, 'Design, Properties and Recent Application of Macrocycles in Medicinal  
1044 Chemistry', *Chimia*, vol. 71, no. 10, p. 678, Oct. 2017, doi: 10.2533/chimia.2017.678.  
1045 **[cito:citesAsAuthority]**
- 1046 [22] P. G. Dougherty, Z. Qian, and D. Pei, 'Macrocycles as protein–protein interaction  
1047 inhibitors', *Biochemical Journal*, vol. 474, no. 7, pp. 1109–1125, Apr. 2017, doi:  
1048 10.1042/BCJ20160619. **[cito:citesAsAuthority]**
- 1049 [23] C. Kramer, M. Podewitz, P. Ertl, and K. Liedl, 'Unique Macrocycles in the Taiwan  
1050 Traditional Chinese Medicine Database', *Planta Med*, vol. 81, no. 06, pp. 459–466, Apr.  
1051 2015, doi: 10.1055/s-0035-1545881. **[cito:citesAsAuthority]**
- 1052 [24] J. L. Medina-Franco, J. Petit, and G. M. Maggiora, 'Hierarchical Strategy for Identifying

1053 Active Chemotype Classes in Compound Databases', *Chemical Biology & Drug Design*,  
1054 vol. 67, no. 6, pp. 395–408, Jun. 2006, doi: 10.1111/j.1747-0285.2006.00397.x.  
1055 **[cito:citesAsAuthority]**

1056 [25] Y.-J. Xu and M. Johnson, 'Using Molecular Equivalence Numbers To Visually Explore  
1057 Structural Features that Distinguish Chemical Libraries', *J. Chem. Inf. Comput. Sci.*, vol.  
1058 42, no. 4, pp. 912–926, Jul. 2002, doi: 10.1021/ci025535l. **[cito:citesAsAuthority]**  
1059 **[cito:discusses]**

1060 [26] C. Manelfi *et al.*, "'Molecular Anatomy": a new multi-dimensional hierarchical scaffold  
1061 analysis tool', *J. Cheminformatics*, vol. 13, no. 1, p. 54, Dec. 2021, doi: 10.1186/s13321-  
1062 021-00526-y. **[cito:citesAsAuthority]** **[cito:usesMethodIn]** **[cito:discusses]**

1063 [27] Y. Hu and J. Bajorath, 'Combining Horizontal and Vertical Substructure Relationships in  
1064 Scaffold Hierarchies for Activity Prediction', *J. Chem. Inf. Model.*, vol. 51, no. 2, pp. 248–  
1065 257, Feb. 2011, doi: 10.1021/ci100448a. **[cito:citesAsAuthority]**

1066 [28] K. Klein, O. Koch, N. Kriege, P. Mutzel, and T. Schäfer, 'Visual Analysis of Biological  
1067 Activity Data with Scaffold Hunter', *Mol. Inform.*, vol. 32, no. 11–12, pp. 964–975, Dec.  
1068 2013, doi: 10.1002/minf.201300087. **[cito:citesAsAuthority]**

1069 [29] S. Wetzel *et al.*, 'Interactive exploration of chemical space with Scaffold Hunter', *Nat.*  
1070 *Chem. Biol.*, vol. 5, no. 8, pp. 581–583, Aug. 2009, doi: 10.1038/nchembio.187.  
1071 **[cito:citesAsAuthority]**

1072 [30] P. Ertl, A. Schuffenhauer, and S. Renner, 'The Scaffold Tree: An Efficient Navigation in  
1073 the Scaffold Universe', in *Chemoinformatics and Computational Chemical Biology*, vol.  
1074 672, J. Bajorath, Ed. Totowa, NJ: Humana Press, 2010, pp. 245–260. doi: 10.1007/978-  
1075 1-60761-839-3\_10. **[cito:citesAsAuthority]**

1076 [31] X. Q. Lewell, D. B. Judd, S. P. Watson, and M. M. Hann, 'RECAP-Retrosynthetic  
1077 Combinatorial Analysis Procedure: A Powerful New Technique for Identifying Privileged  
1078 Molecular Fragments with Useful Applications in Combinatorial Chemistry', *J. Chem. Inf.*  
1079 *Comput. Sci.*, vol. 38, no. 3, pp. 511–522, May 1998, doi: 10.1021/ci970429i.  
1080 **[cito:citesAsAuthority]**

- 1081 [32] J. J. Naveja and J. L. Medina-Franco, 'Finding Constellations in Chemical Space Through  
1082 Core Analysis', *Front. Chem.*, vol. 7, p. 510, Jul. 2019, doi: 10.3389/fchem.2019.00510.  
1083 **[cito:citesAsAuthority]**
- 1084 [33] J. Bajorath, 'Improving the utility of molecular scaffolds for medicinal and computational  
1085 chemistry', *Future Medicinal Chemistry*, vol. 10, no. 14, pp. 1645–1648, Jul. 2018, doi:  
1086 10.4155/fmc-2018-0106. **[cito:citesAsAuthority]**
- 1087 [34] J. J. Naveja and M. Vogt, 'Automatic Identification of Analogue Series from Large  
1088 Compound Data Sets: Methods and Applications', *Molecules*, vol. 26, no. 17, p. 5291,  
1089 Aug. 2021, doi: 10.3390/molecules26175291. **[cito:citesAsAuthority]**
- 1090 [35] J. J. Naveja, M. Vogt, D. Stumpfe, J. L. Medina-Franco, and J. Bajorath, 'Systematic  
1091 Extraction of Analogue Series from Large Compound Collections Using a New  
1092 Computational Compound–Core Relationship Method', *ACS Omega*, vol. 4, no. 1, pp.  
1093 1027–1032, Jan. 2019, doi: 10.1021/acsomega.8b03390. **[cito:citesAsAuthority]**
- 1094 [36] M. Wawer and J. Bajorath, 'Local Structural Changes, Global Data Views: Graphical  
1095 Substructure–Activity Relationship Trailing', *J. Med. Chem.*, vol. 54, no. 8, pp. 2944–  
1096 2951, Apr. 2011, doi: 10.1021/jm200026b. **[cito:citesAsAuthority]**
- 1097 [37] A. de la Vega de León and J. Bajorath, 'Matched molecular pairs derived by retrosynthetic  
1098 fragmentation', *Med. Chem. Commun.*, vol. 5, no. 1, pp. 64–67, 2014, doi:  
1099 10.1039/C3MD00259D. **[cito:citesAsAuthority]**
- 1100 [38] J. J. Naveja, B. A. Pilón-Jiménez, J. Bajorath, and J. L. Medina-Franco, 'A general  
1101 approach for retrosynthetic molecular core analysis', *J. Cheminform*, vol. 11, no. 1, p. 61,  
1102 Dec. 2019, doi: 10.1186/s13321-019-0380-5. **[cito:citesAsAuthority]**
- 1103 [39] T. Schäfer, N. Kriege, L. Humbeck, K. Klein, O. Koch, and P. Mutzel, 'Scaffold Hunter: a  
1104 comprehensive visual analytics framework for drug discovery', *J. Cheminformatics*, vol.  
1105 9, no. 1, p. 28, Dec. 2017, doi: 10.1186/s13321-017-0213-3. **[cito:citesAsAuthority]**
- 1106 [40] F. Kruger, N. Stiefl, and G. A. Landrum, 'rdScaffoldNetwork: The Scaffold Network  
1107 Implementation in RDKit', *J. Chem. Inf. Model.*, vol. 60, no. 7, pp. 3331–3335, Jul. 2020,  
1108 doi: 10.1021/acs.jcim.0c00296. **[cito:citesAsAuthority]**



- 1109 [41] M. K. Matlock, J. M. Zaretski, and S. J. Swamidass, 'Scaffold network generator: a tool  
1110 for mining molecular structures', *Bioinformatics*, vol. 29, no. 20, pp. 2655–2656, Oct.  
1111 2013, doi: 10.1093/bioinformatics/btt448. **[cito:citesAsAuthority]**
- 1112 [42] Jianxing and EX2L, *Scaffold Network Generator*. Peking University HSC. Accessed: Jan.  
1113 12, 2022. [Online]. Available:  
1114 [https://github.com/huluxiaohuowa/scaffold\\_network\\_generator](https://github.com/huluxiaohuowa/scaffold_network_generator) **[cito:citesAsAuthority]**
- 1115 [43] O. B. Scott and A. W. Edith Chan, 'ScaffoldGraph: an open-source library for the  
1116 generation and analysis of molecular scaffold networks and scaffold trees',  
1117 *Bioinformatics*, vol. 36, no. 12, pp. 3930–3931, Jun. 2020, doi:  
1118 10.1093/bioinformatics/btaa219. **[cito:citesAsAuthority]**
- 1119 [44] D. K. Agrafiotis and J. J. M. Wiener, 'Scaffold Explorer: An Interactive Tool for Organizing  
1120 and Mining Structure–Activity Data Spanning Multiple Chemotypes', *J. Med. Chem.*, vol.  
1121 53, no. 13, pp. 5002–5011, Jul. 2010, doi: 10.1021/jm1004495. **[cito:citesAsAuthority]**
- 1122 [45] C. Steinbeck, Y. Han, S. Kuhn, O. Horlacher, E. Luttmann, and E. Willighagen, 'The  
1123 Chemistry Development Kit (CDK): An Open-Source Java Library for Chemo- and  
1124 Bioinformatics', *J. Chem. Inf. Comput. Sci.*, vol. 43, no. 2, pp. 493–500, Mar. 2003, doi:  
1125 10.1021/ci025584y. **[cito:citesAsAuthority] [cito:usesMethodIn]**
- 1126 [46] C. Steinbeck, C. Hoppe, S. Kuhn, M. Floris, R. Guha, and E. Willighagen, 'Recent  
1127 Developments of the Chemistry Development Kit (CDK) - An Open-Source Java Library  
1128 for Chemo- and Bioinformatics', *Curr. Pharm. Des.*, vol. 12, no. 17, pp. 2111–2120, Jun.  
1129 2006, doi: 10.2174/138161206777585274. **[cito:citesAsAuthority]**  
1130 **[cito:usesMethodIn]**
- 1131 [47] E. L. Willighagen *et al.*, 'The Chemistry Development Kit (CDK) v2.0: atom typing,  
1132 depiction, molecular formulas, and substructure searching', *J. Cheminformatics*, vol. 9,  
1133 no. 1, p. 33, Dec. 2017, doi: 10.1186/s13321-017-0220-4. **[cito:citesAsAuthority]**  
1134 **[cito:usesMethodIn]**
- 1135 [48] N. M. O'Boyle, M. Banck, C. A. James, C. Morley, T. Vandermeersch, and G. R.  
1136 Hutchison, 'Open Babel: An open chemical toolbox', *J. Cheminformatics*, vol. 3, no. 1, p.

1137 33, Dec. 2011, doi: 10.1186/1758-2946-3-33. **[cito:citesAsAuthority]**

1138 [49] 'RDKit: Open-Source Cheminformatics Software'. <http://www.rdkit.org> (accessed Jan.

1139 14, 2022). **[cito:citesAsAuthority]**

1140 [50] 'Daylight Theory: SMARTS - A Language for Describing Molecular Patterns'.

1141 <https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html> (accessed Feb. 21,

1142 2022). **[cito:citesAsAuthority]**

1143 [51] F. Kruger, N. Fechner, and N. Stiefl, 'Automated Identification of Chemical Series:

1144 Classifying like a Medicinal Chemist', *J. Chem. Inf. Model.*, vol. 60, no. 6, pp. 2888–2902,

1145 Jun. 2020, doi: 10.1021/acs.jcim.0c00204. **[cito:citesAsAuthority]**

1146 [52] R. Guha, *MurckoFragmenter*. Accessed: Jan. 14, 2022. [Chemistry Development Kit

1147 (CDK)]. Available:

1148 <https://github.com/cdk/cdk/blob/master/tool/fragment/src/main/java/org/openscience/cd>

1149 [k/fragment/MurckoFragmenter.java](https://github.com/cdk/cdk/blob/master/tool/fragment/src/main/java/org/openscience/cdk/fragment/MurckoFragmenter.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1150 [53] A. Dutot, F. Guinand, D. Olivier, and Y. Pigné, 'GraphStream: A Tool for bridging the gap

1151 between Complex Systems and Dynamic Graphs', Dresden, Germany, Oct. 2007.

1152 [Online]. Available: <https://hal.archives-ouvertes.fr/hal-00264043> (accessed Feb. 21,

1153 2022). **[cito:citesAsAuthority] [cito:usesMethodIn]**

1154 [54] 'GraphStream - A Dynamic Graph Library', *GraphStream - A Dynamic Graph Library*.

1155 <http://graphstream-project.org/> (accessed Jan. 24, 2022). **[cito:citesAsAuthority]**

1156 **[cito:usesMethodIn]**

1157 [55] C. Steinbeck, *IAtomContainer Interface*. Accessed: Jan. 24, 2022. [Chemistry

1158 Development Kit (CDK)]. Available:

1159 <https://github.com/cdk/cdk/blob/master/base/interfaces/src/main/java/org/openscience/c>

1160 [dk/interfaces/IAtomContainer.java](https://github.com/cdk/cdk/blob/master/base/interfaces/src/main/java/org/openscience/cdk/interfaces/IAtomContainer.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1161 [56] J. May, *Cycles*. Accessed: Jan. 24, 2022. [Chemistry Development Kit (CDK)]. Available:

1162 <https://github.com/cdk/cdk/blob/master/base/core/src/main/java/org/openscience/cdk/gr>

1163 [aph/Cycles.java](https://github.com/cdk/cdk/blob/master/base/core/src/main/java/org/openscience/cdk/graph/Cycles.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1164 [57] J. W. May and C. Steinbeck, 'Efficient ring perception for the Chemistry Development

1165 Kit', *J. Cheminformatics*, vol. 6, no. 1, p. 3, Dec. 2014, doi: 10.1186/1758-2946-6-3.  
1166 **[cito:citesAsAuthority] [cito:usesMethodIn]**

1167 [58] M. Sorokina, P. Merseburger, K. Rajan, M. A. Yirik, and C. Steinbeck, 'COCONUT online:  
1168 Collection of Open Natural Products database', *J. Cheminformatics*, vol. 13, no. 1, p. 2,  
1169 Dec. 2021, doi: 10.1186/s13321-020-00478-9. **[cito:citesAsAuthority]**  
1170 **[cito:usesDataFrom]**

1171 [59] D. Weininger, 'SMILES, a chemical language and information system. 1. Introduction to  
1172 methodology and encoding rules', *J. Chem. Inf. Model.*, vol. 28, no. 1, pp. 31–36, Feb.  
1173 1988, doi: 10.1021/ci00057a005. **[cito:citesAsAuthority] [cito:usesMethodIn]**

1174 [60] D. Weininger, A. Weininger, and J. L. Weininger, 'SMILES. 2. Algorithm for generation of  
1175 unique SMILES notation', *J. Chem. Inf. Comput. Sci.*, vol. 29, no. 2, pp. 97–101, May  
1176 1989, doi: 10.1021/ci00062a008. **[cito:citesAsAuthority] [cito:usesMethodIn]**

1177 [61] E. Willighagen and J. W. Mayfield, *IStereoElement*. Accessed: Jan. 24, 2022. [Chemistry  
1178 Development Kit (CDK)]. Available:  
1179 [https://github.com/cdk/cdk/blob/master/base/interfaces/src/main/java/org/openscience/c](https://github.com/cdk/cdk/blob/master/base/interfaces/src/main/java/org/openscience/cdk/interfaces/IStereoElement.java)  
1180 [dk/interfaces/IStereoElement.java](https://github.com/cdk/cdk/blob/master/base/interfaces/src/main/java/org/openscience/cdk/interfaces/IStereoElement.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1181 [62] J. May, *Aromaticity*. Accessed: Jan. 24, 2022. [Chemistry Development Kit (CDK)].  
1182 Available:  
1183 [https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/c](https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/cdk/aromaticity/Aromaticity.java)  
1184 [dk/aromaticity/Aromaticity.java](https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/cdk/aromaticity/Aromaticity.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1185 [63] J. May, *ElectronDonation*. Accessed: Jan. 24, 2022. [Chemistry Development Kit (CDK)].  
1186 Available:  
1187 [https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/c](https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/cdk/aromaticity/ElectronDonation.java)  
1188 [dk/aromaticity/ElectronDonation.java](https://github.com/cdk/cdk/blob/master/base/standard/src/main/java/org/openscience/cdk/aromaticity/ElectronDonation.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**

1189 [64] M. Solà, 'Connecting and combining rules of aromaticity. Towards a unified theory of  
1190 aromaticity', *WIREs Comput Mol Sci*, vol. 9, no. 4, Jul. 2019, doi: 10.1002/wcms.1404.  
1191 **[cito:citesAsAuthority]**

1192 [65] Z. Chen, C. S. Wannere, C. Corminboeuf, R. Puchta, and P. von R. Schleyer, 'Nucleus-

- 1193 Independent Chemical Shifts (NICS) as an Aromaticity Criterion', *Chem. Rev.*, vol. 105,  
1194 no. 10, pp. 3842–3888, Oct. 2005, doi: 10.1021/cr030088+. **[cito:citesAsAuthority]**
- 1195 [66] K. Krämer, 'The search for the grand unification of aromaticity', *Chemistry World*, Jul. 26,  
1196 2021. Accessed: Oct. 25, 2022. [Online]. Available:  
1197 [https://www.chemistryworld.com/features/the-search-for-the-grand-unification-of-](https://www.chemistryworld.com/features/the-search-for-the-grand-unification-of-aromaticity/4013915.article)  
1198 [aromaticity/4013915.article](https://www.chemistryworld.com/features/the-search-for-the-grand-unification-of-aromaticity/4013915.article) **[cito:citesAsAuthority]**
- 1199 [67] A. Stanger, 'What is... aromaticity: a critique of the concept of aromaticity—can it really  
1200 be defined?', *Chem. Commun.*, no. 15, p. 1939, 2009, doi: 10.1039/b816811c.  
1201 **[cito:citesAsAuthority]**
- 1202 [68] O. Horlacher, S. Kuhn, and J. May, *SmilesGenerator*. Accessed: Jan. 24, 2022.  
1203 [Chemistry Development Kit (CDK)]. Available:  
1204 [https://github.com/cdk/cdk/blob/master/storage/smiles/src/main/java/org/openscience/c](https://github.com/cdk/cdk/blob/master/storage/smiles/src/main/java/org/openscience/cdk/smiles/SmilesGenerator.java)  
1205 [dk/smiles/SmilesGenerator.java](https://github.com/cdk/cdk/blob/master/storage/smiles/src/main/java/org/openscience/cdk/smiles/SmilesGenerator.java) **[cito:citesAsAuthority] [cito:usesMethodIn]**
- 1206 [69] *JUnit*. Accessed: Jan. 25, 2022. [Online]. Available: <https://junit.org/junit4/>  
1207 **[cito:usesMethodIn]**
- 1208 [70] D. S. Wishart, 'DrugBank: a comprehensive resource for in silico drug discovery and  
1209 exploration', *Nucleic Acids Res.*, vol. 34, no. 90001, pp. D668–D672, Jan. 2006, doi:  
1210 10.1093/nar/gkj067. **[cito:citesAsAuthority] [cito: usesDataFrom]**
- 1211 [71] D. S. Wishart *et al.*, 'DrugBank 5.0: a major update to the DrugBank database for 2018',  
1212 *Nucleic Acids Res.*, vol. 46, no. D1, pp. D1074–D1082, Jan. 2018, doi:  
1213 10.1093/nar/gkx1037. **[cito:citesAsAuthority] [cito: usesDataFrom]**
- 1214 [72] M. Ashton *et al.*, 'Identification of Diverse Database Subsets using Property-Based and  
1215 Fragment-Based Molecular Descriptions', *Quant. Struct.-Act. Relatsh.*, vol. 21, no. 6, pp.  
1216 598–604, Dec. 2002, doi: 10.1002/qsar.200290002. **[cito:citesAsAuthority]**
- 1217 [73] P. Ertl and T. Schuhmann, 'A Systematic Cheminformatics Analysis of Functional Groups  
1218 Occurring in Natural Products', *J. Nat. Prod.*, vol. 82, no. 5, pp. 1258–1263, May 2019,  
1219 doi: 10.1021/acs.jnatprod.8b01022. **[cito:citesAsAuthority] [cito:agreesWith]**
- 1220 [74] K. Grabowski, K.-H. Baringhaus, and G. Schneider, 'Scaffold diversity of natural products:

1221 inspiration for combinatorial library design', *Nat. Prod. Rep.*, vol. 25, no. 5, p. 892, 2008,  
1222 doi: 10.1039/b715668p. **[cito:citesAsAuthority] [cito:agreesWith]**

1223 [75] Y. Chen, M. Garcia de Lomana, N.-O. Friedrich, and J. Kirchmair, 'Characterization of the  
1224 Chemical Space of Known and Readily Obtainable Natural Products', *J. Chem. Inf.*  
1225 *Model.*, vol. 58, no. 8, pp. 1518–1532, Aug. 2018, doi: 10.1021/acs.jcim.8b00302.  
1226 **[cito:citesAsAuthority] [cito:agreesWith]**

1227 [76] P. Ertl, S. Jelfs, J. Mühlbacher, A. Schuffenhauer, and P. Selzer, 'Quest for the Rings. In  
1228 Silico Exploration of Ring Universe To Identify Novel Bioactive Heteroaromatic  
1229 Scaffolds', *J. Med. Chem.*, vol. 49, no. 15, pp. 4568–4573, Jul. 2006, doi:  
1230 10.1021/jm060217p. **[cito:citesAsAuthority] [cito:agreesWith]**

1231 [77] A. M. Clark, '2D Depiction of Fragment Hierarchies', *J. Chem. Inf. Model.*, vol. 50, no. 1,  
1232 pp. 37–46, Jan. 2010, doi: 10.1021/ci900350h. **[cito:citesAsAuthority]**

1233 [78] M. Grigalunas *et al.*, 'Natural product fragment combination to performance-diverse  
1234 pseudo-natural products', *Nat. Commun.*, vol. 12, no. 1, p. 1883, Dec. 2021, doi:  
1235 10.1038/s41467-021-22174-4. **[cito:citesAsAuthority]**

1236 [79] M. Grigalunas, A. Burhop, A. Christoforow, and H. Waldmann, 'Pseudo-natural products  
1237 and natural product-inspired methods in chemical biology and drug discovery', *Curr.*  
1238 *Opin. Chem. Biol.*, vol. 56, pp. 111–118, Jun. 2020, doi: 10.1016/j.cbpa.2019.10.005.  
1239 **[cito:citesAsAuthority]**