

# A Comprehensive Nomenclature System for Cyclodextrins

Amelia Anderson<sup>a</sup>, Matthew S. O'Connor<sup>a</sup>, James Pipkin<sup>b</sup>, Milo Malanga<sup>c</sup>, Tamas Sohajda<sup>c</sup>, Thorsteinn Loftsson<sup>d</sup>, Lajos Szente<sup>e</sup>, Rebeca García-Fandiño<sup>f\*</sup> and Ángel Piñeiro<sup>g†</sup>

<sup>a</sup>Cyclarity Therapeutics, 8001 Redwood Blvd Novato, CA 94945, USA

<sup>b</sup>Ligand Pharmaceuticals Incorporated, 3911 Sorrento Valley Boulevard, San Diego, CA, USA

<sup>c</sup>CarboHyde, Budapest, Berlini u. 47-49, 1045, Hungary

<sup>d</sup>Faculty of Pharmaceutical Sciences, University of Iceland, Hofsvallagata 53, IS-107 Reykjavik, Iceland

<sup>e</sup>CycloLab Cyclodextrin R&D Laboratory Ltd., Illatos u. 7., Budapest, H-1097, Hungary

<sup>f</sup>Department of Organic Chemistry, Center for Research in Biological Chemistry and Molecular Materials, University of Santiago de Compostela, CIQUS, Spain

<sup>g</sup>Department of Applied Physics, Faculty of Physics, University of Santiago de Compostela, Spain

\*Corresponding author: [rebeca.garcia.fandino@usc.es](mailto:rebeca.garcia.fandino@usc.es)

†Corresponding author: [Angel.Pineiro@usc.es](mailto:Angel.Pineiro@usc.es)

## Abstract

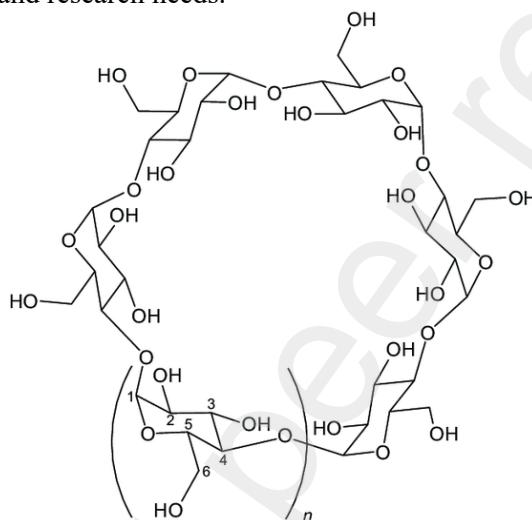
Modified cyclodextrins (CDs) are cyclic oligosaccharides with many applications in drug delivery, catalysis, and as active pharmaceutical ingredients. In general, they exist as distributions of structurally diverse molecules rather than single-isomer compounds. Their performance depends on the number of glucopyranose units (GPUs), and the type, number, and position of chemical substitutions in their hydroxyl groups. Effectively targeting individual species within these distributions is essential for optimizing CDs for specific applications. Computational techniques can generate large datasets to AI-driven structural optimization, but the absence of a standardized nomenclature system for modified CDs presents a major barrier to progress in this direction. This lack of consensus limits effective communication, data sharing, automation, and collaboration. To address this, a clear and extensible nomenclature for modified CDs is proposed. In this framework, GPUs are treated like amino-acid residues, with unsubstituted GPUs as reference building-blocks and substituted ones considered as mutations. This approach precisely defines substitution types and patterns, resolves cyclic permutation ambiguities, and offers versatility for both simple and complex modifications, including chiral center alterations and covalently linked CD oligomers. By introducing this standardized nomenclature, we aim to enhance molecular design, improve reproducibility, and streamline both experimental and computational research in the CD field.

**Note:** To facilitate the adoption of the proposed nomenclature, we have developed a web-based tool named **CyDexID Generator**, which automatically assigns standardized names to modified CD monomers and generates the corresponding 3D structures: <https://cydexid-app.lm.r.appspot.com>.

**Keywords:** cyclodextrins; nomenclature; mutations; chemical substitutions; cyclic permutations; sequences.

## 49 1.- Introduction

50 Cyclodextrins (CDs) are cyclic oligosaccharides composed of  $\alpha$ -D-glucopyranose units (GPUs)  
51 linked by  $\alpha$ -(1 $\rightarrow$ 4) glycosidic bonds (Fig. 1). Since their discovery by Villiers in 1891[1], CDs  
52 have attracted significant interest across various fields—including pharmaceuticals, food science,  
53 cosmetics, and chemistry—due to their unique ability to form inclusion complexes with a wide  
54 range of guest molecules[2–4]. This property, combined with their biocompatibility and low  
55 toxicity, has made CDs invaluable in drug delivery systems, solubility enhancement, stabilization  
56 of volatile substances, and applications in supramolecular catalysis[2,4,5]. The most common  
57 natural CDs are  $\alpha$ -,  $\beta$ -, and  $\gamma$ -cyclodextrins, consisting of 6, 7, and 8 GPUs, respectively. Chemical  
58 modifications at positions C-2, C-3, and C-6 of the glucose units (Fig. 1) have significantly  
59 expanded the utility of CDs, allowing for modulation of their physicochemical properties such as  
60 solubility, specificity to recognize individual molecules, complexation capacity, and reactivity[6].  
61 These modifications have broadened the spectrum of CD applications, enabling tailored solutions  
62 for particular industrial and research needs.



63 **Figure 1.-** Chemical structure of an  $\alpha$ -cyclodextrin molecule ( $n=1$ ), where each  
64 glucopyranose unit is linked via  $\alpha$ -1,4-glycosidic bonds, forming a cyclic oligosaccharide.  
65 The hydroxyl groups attached to the C-2, C-3, and C-6 positions of each glucose unit are  
66 highlighted, representing potential sites for chemical substitutions.  
67

68  
69 As the complexity and diversity of modified CDs have increased, the opportunity to use specific  
70 structures for targeted applications, including as active pharmaceuticals, has become more  
71 evident. CDs are already prevalent in food, drug delivery, and environmental applications, where  
72 they enhance stability, solubility, and controlled release. However, synthesis typically yields  
73 mixtures of isomers, underscoring the need for precise and unambiguous structural descriptions.  
74 A standardized CD nomenclature would improve identification, quality control, and cross-  
75 industry communication, as well as facilitate intellectual property protections for specific isomers.  
76 Current naming systems cannot capture complex CD modifications or oligomeric forms [8,9],  
77 leading to inconsistencies, communication barriers, and challenges in building databases or  
78 conducting systematic computational studies. Emerging computational techniques like molecular  
79 docking, molecular dynamics, and machine learning further highlight the need for such a naming  
80 system to maximize research and application potential [8,9].  
81

82 In this context, we propose a new comprehensive nomenclature and representation system for  
83 modified CDs and their oligomers. Drawing inspiration from protein nomenclature, concepts such  
84 as *residues*, *mutations*, and *hierarchical levels of structure* (primary, secondary, tertiary and  
85 quaternary)[7–9] are adapted to the structural specifics of CDs. The main objectives and  
86 applications of this proposal are: (i) to provide a unique and unambiguous representation for each  
87 modified CD structure; (ii) to facilitate precise communication among researchers and regulators  
88 as well as the creation of coherent databases; (iii) to enable automatic structure generation for

89 computational studies; (iv) to establish a *building-block* approach for representation and even for  
90 parameterization in MD simulations, similar to proteins; and (v) to improve the reproducibility  
91 and comparability of computational, and eventual wet-lab, studies involving CDs. The proposed  
92 system encodes modified CD structures in a concise character string, specifying substitution type,  
93 position, chirality, and linker presence in oligomers while resolving cyclic permutation  
94 ambiguities. A key feature of this approach is its parallelism with protein bioinformatics systems.  
95 This not only facilitates understanding and adoption by researchers familiar with bioinformatics,  
96 but also paves the way for adapting established algorithms and methodologies from protein  
97 studies to the CD domain. Alongside this string-based nomenclature, a structured data format is  
98 proposed to simplify CD information processing, and an alternative, less precise notation reduces  
99 structures to a lower-dimensional space by indicating only substitution type and count. While this  
100 compact notation is less precise than the expanded version, it remains useful for certain  
101 comparisons, especially between computational and experimental studies where full structural  
102 details are not available.

103  
104 The following sections detail the proposed nomenclature systems, including descriptions, the  
105 algorithm for ensuring uniqueness, and the structured data file format. Several examples  
106 illustrating the application of the system to various modified CDs are presented—including  
107 complex but theoretically viable structures to push its limits. The system simplifies CD  
108 parametrization for computational modeling, facilitating improvements to molecular docking,  
109 MD simulations, and machine learning predictions for CDs. Finally, we discuss its limitations  
110 and future directions for development. We hope that this nomenclature and representation system  
111 will serve as a significant step toward standardization in the field of modified CDs, enhancing  
112 communication, computational research, and ultimately the rational and automated design of new  
113 CDs with optimized properties for specific applications.

## 114 2.- Background

### 115 2.1.- Structure and properties of cyclodextrins

117 Modified CDs are currently referred to simply by the number of GPUs in the ring, the type of  
118 substitution, and the average degree of substitution (DS) in the sample, however, this description  
119 is incomplete. Each GPU in a CD possesses three hydroxyl groups available for modification: one  
120 primary hydroxyl at the C-6 position and two secondary hydroxyls at the C-2 and C-3 positions  
121 (Fig. 1) [10]. These hydroxyl groups are the primary sites for chemical modifications that can  
122 significantly alter the physicochemical properties of CDs[11,12]. In this work, these substitutions  
123 are interpreted as mutations in the native structure of CDs, extrapolating the terminology typically  
124 employed for proteins, as will be explained later in more detail. Knowing the number of unique  
125 molecular configurations arising from mutations at specific positions is important for  
126 understanding the structural diversity, and so the expected variety of functional properties, of  
127 CDs. Ignoring the probability of each structure and focusing just on preventing the violation of  
128 chemical laws, the number of different possible structures can be determined using *Burnside's*  
129 *Lemma*[13,14]. This lemma allows for consideration of the rotational symmetries of CD  
130 monomers and the fact that, due to their chirality, mirror images are not equivalent. According to  
131 this proposition, and considering CDs composed of  $n$  GPUs, where each GPU can undergo one  
132 of  $k$  type of mutations, the number of unique structures ( $N$ ) is given by:

$$133 N = \frac{1}{|G|} \sum_{d|n} \phi(d) \times k^{\frac{n}{d}} \quad [1]$$

134  
135 where  $|G|$  is the order of the symmetry group, representing the total number of rotational  
136 symmetries (for a cyclic molecule with  $n$  GPUs,  $|G| = n$ ),  $d$  denotes each divisor of  $n$ ,  $\phi(d)$  is  
137 Euler's Totient function, which counts the number of integers up to  $d$  that are coprime with  $d$ , and  
138  $k$  is the number of possible *mutations* per GPU. Eight different possibilities for each GPU unit ( $k$   
139 = 8) will be considered here: native GPUs, GPUs with single modifications at C-2, C-3 or C-6,  
140 GPUs with double mutations (at positions C-2 and C-3, C-2 and C-6 or C-3 and C-6), and GPUs  
141 with all three hydroxyl groups substituted. This latter configuration, while chemically feasible

142 and relatively accessible[15,16], is less frequently observed in experimental samples. Under these  
143 considerations, the number of unique possible configurations for different CD types can be  
144 determined as follows:

145  $\alpha$ -CD,  $n = 6$ :

146

$$147 \quad N_{\alpha\text{-CD}} = \frac{1}{6}(\phi(1) \times 8^6 + \phi(2) \times 8^3 + \phi(3) \times 8^2 + \phi(6) \times 8^1)$$

148  $= \frac{1}{6}(1 \times 262,144 + 1 \times 512 + 2 \times 64 + 2 \times 8) = \frac{1}{6} \times 262,800 = 43,800$

149

150

151  $\beta$ -CD,  $n = 7$ :

152 
$$N_{\beta\text{-CD}} = \frac{1}{7}(\phi(1) \times 8^7 + \phi(7) \times 8^1)$$

153  $= \frac{1}{7}(1 \times 2,097,152 + 6 \times 8) = \frac{1}{7} \times 2,097,200 = 299,600$

154

155

156  $\gamma$ -CD,  $n = 8$ :

157 
$$N_{\gamma\text{-CD}} = \frac{1}{8}(\phi(1) \times 8^8 + \phi(2) \times 8^4 + \phi(4) \times 8^2 + \phi(8) \times 8^1)$$

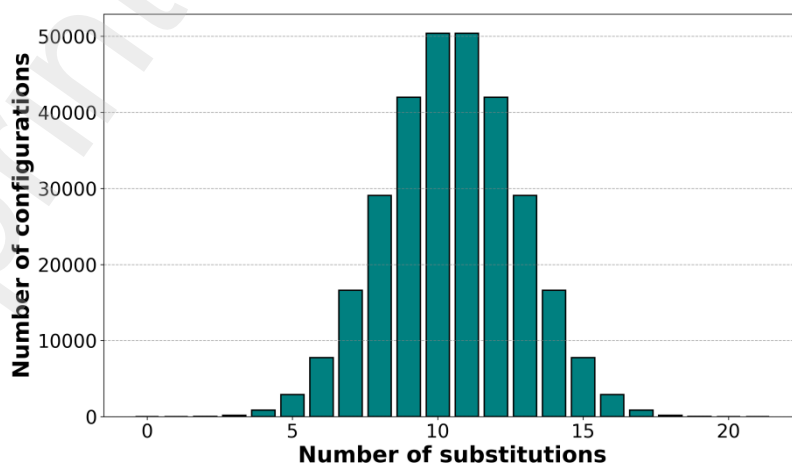
158  $= \frac{1}{8}(1 \times 16,777,216 + 1 \times 4,096 + 2 \times 64 + 4 \times 8) = \frac{1}{8} \times 16,781,472 = 2,097,684$

159

160

161 Thus, for monomeric  $\alpha$ ,  $\beta$ , and  $\gamma$ -CDs with a single substitution type we could, eventually, have  
162 up to 43,800, 299,600, and 2,097,684 unique structures, respectively. To illustrate this, the  
163 number of possible structures as a function of the total number of substitutions in  $\beta$ -CD, regardless  
164 of their location, are presented in Fig. 2, while all possible configurations for several different  
165 combinations of substitutions in  $\alpha$ -CD are shown in Fig. 3. It is important to note that the  
166 probability of each configuration and the distribution of structures with a specific number of  
167 substitutions can vary significantly, depending on the specific chemical synthesis pathway and  
168 the subsequent purification processes. Note also that the physicochemical properties of each of  
169 these structures might significantly depend on the number and location of the substitutions, so  
170 identifying them unambiguously is key to optimizing specific applications. Equation [1] has been  
171 previously applied to determine the number of distinct substitution patterns in C6 cycles[17].  
172 However, to the best of our knowledge, the total number of possible configurations for CD  
173 molecules, accounting for all substitution patterns, has not been reported yet.

174



175

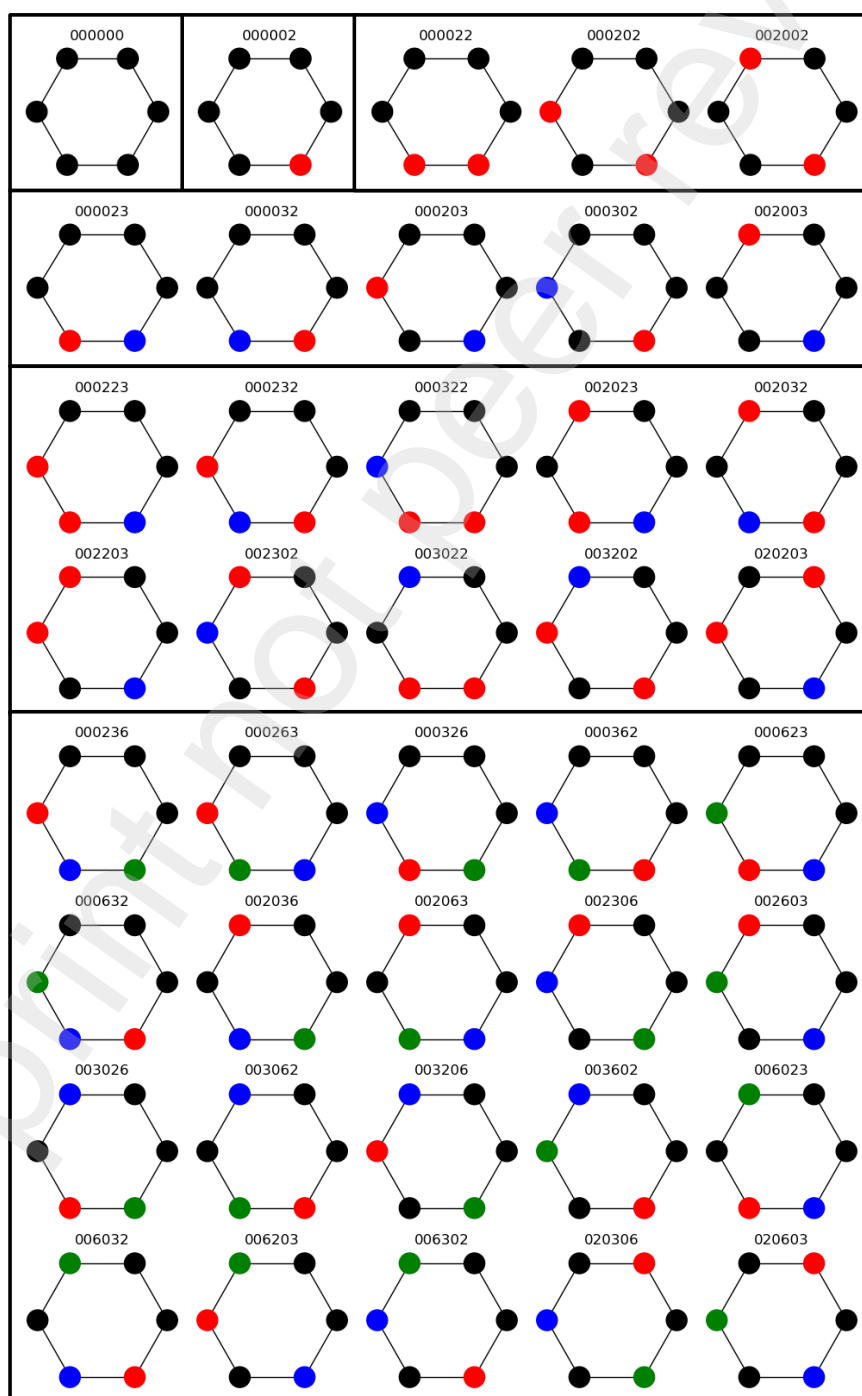
176

177

**Figure 2.-** Number of possible different structures for a  $\beta$ -CD (7 GPUs) as a function of the total number of substitutions at positions C-2, C-3 and C-6.

178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191

The previous results apply to CD monomers with a single substitution type. The number of different configurations for covalently linked CD dimers and higher-order oligomers[18,19] can be determined similarly, with a significant increase in structural diversity by several orders of magnitude due to the possible mutation combinations in each CD subunit. It is important to note that the rotational symmetry observed in monomers is lost in oligomers due to the presence of chemical linkers, which disrupt the equivalence of identical substitution configurations upon rotation of individual monomers. However, dimers with symmetric linkers where the attachment points on both CD monomers are equivalent (for instance a saturated alkyl chain bound to two CD units at C-2), lead to a new type of symmetry. In this case, certain configurations become equivalent when the two monomers are exchanged. It is worth mentioning that hybrid modifications or even conjugation of CDs with specific molecules or entities can also be present in monomers or dimers[20], increasing the dimensions of the structural space even more.



192

193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245

**Figure 3.-** All possible substitution patterns for  $\alpha$ -CD based on specific examples of increasing complexity. Each vertex of the hexagon represents a GPU, with native (unsubstituted) GPUs depicted in black, and substitutions at positions C-2, C-3, and C-6 represented by red, blue, and green spheres, respectively. The native CDs without any substitutions or with a single substitution at one specific position have only one structure each. For the case with two identical substitutions, there are 3 possible structures. The CD with two different substitutions yields 5 structures. For CDs with three substitutions, where two may be identical and one different, or all three are different, the number of possible structures increases to 10 and 20, respectively. Altogether, the figure contains 40 distinct configurations out of the 43,800 possible unique structures for  $\alpha$ -CD. The sequence representing the substitution pattern for each structure is indicated above the corresponding diagram (see description of the notation in the “*Description of the nomenclature system*” section). It is worth noting that if all the GPUs had different substitution patterns, considering rotational but not mirror symmetry, there would be 120 unique combinations for  $\alpha$ -CD. For  $\beta$ -CD, with distinct mutations on all GPUs, the number of unique structures would increase to 720.

The take-home message of the previous calculations and discussion is that, even in the simplest cases, modified CDs can present at least several tens of thousands of structures, and, in more complex scaffolds, this number can reach into the millions or even billions. Each of these structures is expected to display different physicochemical properties that are crucial for their function, affecting their ability to encapsulate specific compounds as well as their self-assembly behavior[18,21]. One of the overarching goals of computational modeling is to enable efficient studies of these structures individually [19,22], with the results ultimately guiding the extraction of the most suitable fractions from real solutions to optimize practical applications.

The extensive number of potential substitution patterns requires a robust and systematic nomenclature framework. Current nomenclature systems for modified CDs vary widely and often lack consistency. Common approaches include abbreviated structural representations and systematic or semi-systematic names. Additionally, a modified CD may obtain patent protection and a measure of commercial use, often under a branded or trademarked name. Such is the case for CAPTISOL<sup>®</sup>, a distribution of  $\beta$ -cyclodextrin molecules with a variety of sulfobutylether substitutions (<https://www.captisol.com/>) and ADVASEP<sup>™</sup>. CAPTISOL<sup>®</sup> is also referred to as SBE- $\beta$ -CD, which is an example of a semi-systematic name and abbreviated structural representation[23]. A suffix with the average number of substitutions in the sample, such as SBE- $\beta$ -CD\_DS6.5, is typically added to the semi-systematic names, but they are widely ambiguous because, as explained above, they may contain many thousands or even millions of different structures. The classical IUPAC names are also systematic chemical names[24]. They are precise but difficult to read and to parse into building blocks for computational codes. In 1997, a systematic nomenclature system was proposed for modified CDs[23], to describe their structural modifications, focusing on the type, position, and average number of substituents. This system uses the base CD structure (e.g.,  $\alpha$ -,  $\beta$ -, or  $\gamma$ -CD) as a starting point, followed by abbreviations for the substituents and their positions when known. For example, HP4- $\beta$ -CD represents a  $\beta$ -CD molecule with an average of four 2-hydroxypropyl substituents, while 6-SBE1- $\beta$ -CD specifies a sulfobutylether substituent attached to the 6th position. This nomenclature facilitates quick identification and categorization of modified CDs but cannot fully specify exact substitution patterns or describe the precise distribution of substitution patterns within a sample.

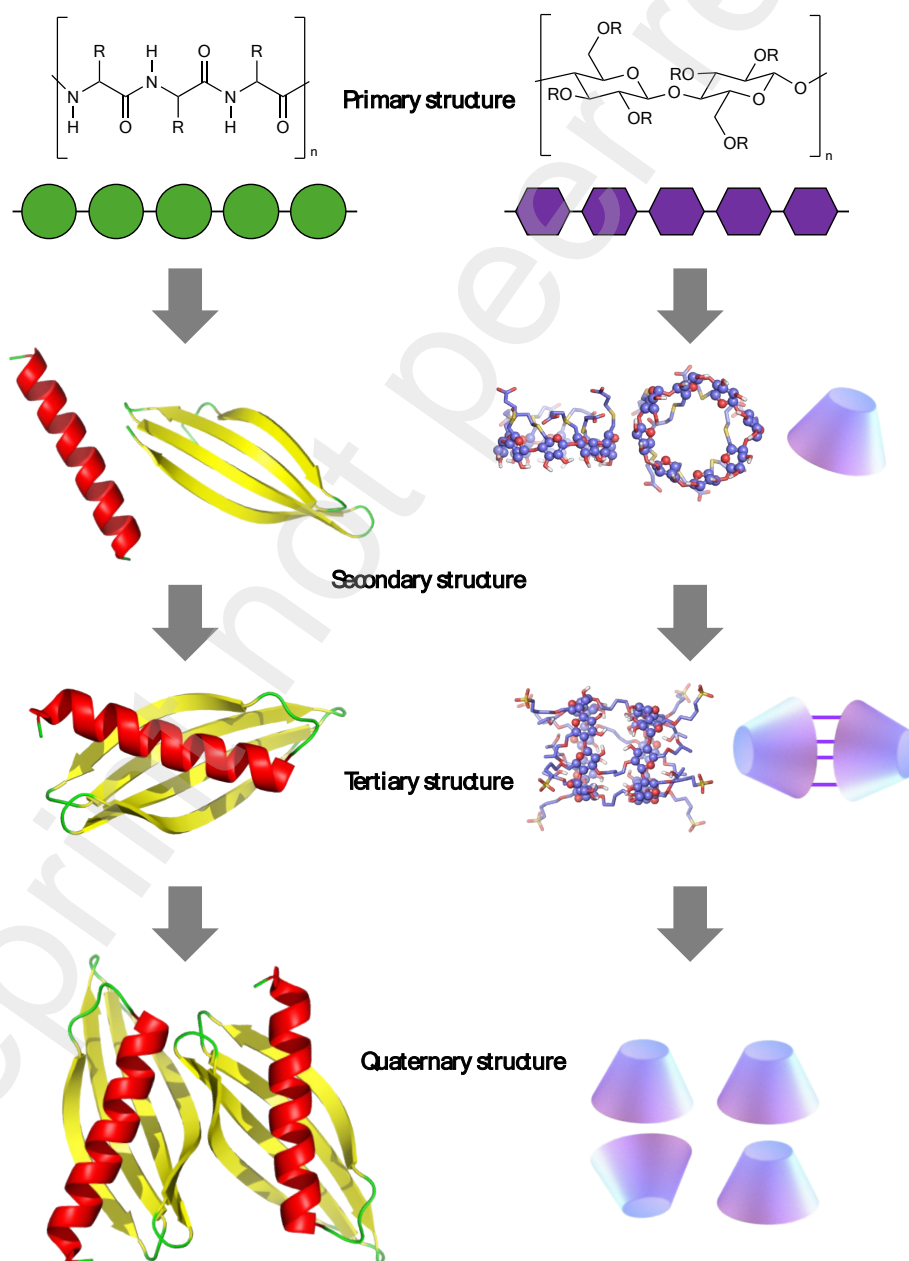
In general, these nomenclature systems suffer from the following limitations: (i) lack of uniformity across different modification types; (ii) inability to precisely represent the position of modifications on specific GPUs; (iii) difficulty in representing complex or hybrid modifications; (iv) challenges in computational processing and database storage; and (v) ambiguity in

246 representing dimeric or polymeric CD structures. These shortcomings highlight the need for a  
247 more comprehensive and standardized nomenclature system[25].

248

## 249 2.2.- Analogies between oligosaccharides and proteins

250 There are significant similarities between CDs and other biomolecules, particularly peptides and  
251 proteins, which are made from polymerized amino acids while CDs consist of repeating glucose  
252 units. The sequence of modifications in CDs can be compared to the amino acid sequence in  
253 peptides, suggesting the introduction of a *primary structure* for CDs, defined by the list of GPU  
254 units specifying a given substitution pattern. Moreover, the spatial arrangement of these  
255 modifications influences their function, much like *secondary structure* in protein folding. CD  
256 dimers and higher-order oligomers also parallel protein oligomers, which opens the possibility of  
257 defining a *tertiary structure* for CDs based on specific interactions between subunits, such as tail-  
258 to-tail, tail-to-head, or head-to-head orientations, among others (Fig. 4). These analogies point to  
259 the potential advantages of adopting concepts from protein nomenclature and bioinformatics to  
260 develop a more robust system for CDs.  
261



262

**Figure 4.-** Structural organization levels in proteins (left column) and cyclodextrins (right column), highlighting the similarities between the two molecular systems. From top to bottom: primary structure (linear chain of amino acids for proteins and GPUs for cyclodextrins), secondary structure ( $\alpha$ -helix and  $\beta$ -sheet for proteins, and the arrangement of GPUs into cyclodextrin monomers), and tertiary structure (folded 3D structure of proteins and cyclodextrin dimers or higher order oligomers). A quaternary structure can also be defined for both as the specific but not covalent aggregation of several subunits.

While the similarities are significant and they can be exploited to define a precise nomenclature system taking proteins as a reference, the most significant difference between general protein systems and CDs is the cyclization of the chain. Although cyclic peptides exist in nature[26], they do not represent the most common structural pattern. Of course, non-cyclic oligosaccharides also exist, but here we will focus exclusively on CDs.

The main motivation to propose a new nomenclature is the growing ability to perform large-scale computational studies on CDs, including individual molecules in solution, inclusion complexes, and aggregates of homogeneous or heterogeneous CD systems[11]. While many steps of methods like docking, molecular dynamics, and quantum calculations are automated for proteins, applying them to CDs is challenging due to the lack of automated tools for CD structure building and parameterization. [7–9]. Importantly, computational approaches in the field of CDs are expected to yield more accurate results compared to protein systems, given the lower degrees of freedom of cyclic oligosaccharides[27]. A standardized, building block-based nomenclature for CDs would significantly advance computational research, aiding the design of optimized structures for specific applications through machine learning. Existing protein-focused tools, such as AutoDock Vina[28–36], AMBER[37], and GROMACS[30–36], could be adapted for CDs by incorporating ring-specific constraints and tailored force fields that better describe the unique stereochemistry and flexibility of CDs. Based on this building block scheme, some of us have recently adapted Autodock Vina to specifically consider flexibility just in the hydroxyl and substituted groups of CDs and developed specific software for automatic parameterization of modified CDs and CD dimers for further MD simulations[27]. These adaptations would ensure accurate modeling of CD systems, particularly in scenarios involving host-guest interactions or derivatized CDs with varied substitution patterns.

### 3.- Description of the nomenclature system

The proposed nomenclature system aims to address the limitations of existing alternatives by providing a standardized, unambiguous, and computationally tractable representation of modified CDs and their dimers. The key principles of our proposal are:

- **Uniqueness and completeness:** Each particular CD structure should have a unique identifier and the method should be able to represent any CD structure.
- **Comprehensiveness:** The CD expression should capture all relevant structural information, including the type, number, and position of modifications.
- **Scalability:** It should accommodate simple modifications, complex hybrid structures with multiple substitution types, and even oligomeric CDs with several branches or cross-links.
- **Computational friendliness:** The format should be easily parsed and processed by computer algorithms.
- **Human readability:** The system should also be easily interpretable by researchers with minimal training at a glance.
- **Analogy to protein nomenclature:** Leveraging established protein nomenclature conventions to adapt bioinformatics concepts and tools. This includes a building-block strategy for nomenclature that can also be used for digital synthesis methods, i.e. building 3D structures and obtaining parameterizations for computational experiments such as molecular dynamics simulations.



317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370

Next, the proposal will be described step by step, starting with the simplest structures and demonstrating how it meets the prerequisites outlined in the previous list of key principles.

### 3.1.- Encoding single and multiple substitutions of a single type for CD monomers

Each CD structure is represented by a string formed by four concatenated fields:

[Type]\_[#Substitutions]\_[#GPUs]x[Positions]

where

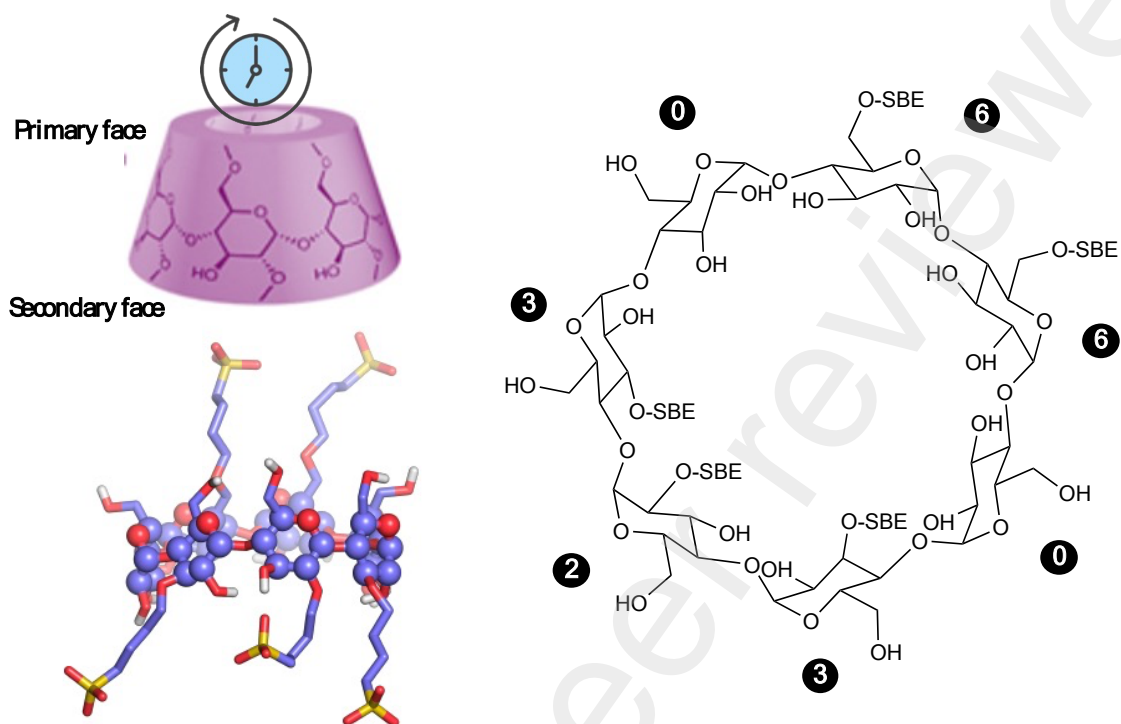
- **Type**: it represents an abbreviation for the substitution type (e.g., SBE for sulfobutylether, HP for 2-hydroxypropyl, ME for methyl, etc).
- **#Substitutions**: it is the total count of this modification type in the CD, the minimum number is zero for any native CDs, and the maximum would be 18, 21, and 24 for  $\alpha$ ,  $\beta$ , and  $\gamma$ -CD monomers, respectively.
- **#GPUs**: it is the number of glucose units in the CD (e.g., 6 for  $\alpha$ -CD, 7 for  $\beta$ -CD, 8 for  $\gamma$ -CD). Nomenclature for CDs with less than 6 or more than 8 GPUs can be easily created with this system.
- **Positions**: it is a sequence indicating substitution positions on each GPU. By convention, when the primary face of the CD is oriented upwards, the sequence is traversed in a clockwise direction (Fig. 5). In order to indicate the substitution positions, the following encoding is proposed (Fig. 6):
  - **0**: for native GPUs in the corresponding position.
  - **2, 3, 6**: for substitutions at positions C-2, C-3, or C-6.
  - **5 (2+3), 8 (2+6), 9 (3+6)**: for double substitutions at the positions indicated by the numbers between parentheses.
  - **1**: for the triple substitution (2+3+6) within the same GPU.
- To indicate the stereochemistry of chiral positions, additional strings can be introduced if needed. By default, all GPU residues are assumed to have natural chirality, and their stereochemistry is not specified.
  - Chiral modifications at positions 2, 3, and/or 5 will be marked with a string starting with RS236, where the “chiral mutations” are identified by the same numbers as those used for the substitutions at positions 2, 3 and/or 6. Although the chiral center is formally at position **5**, we use **6** to denote modifications in its stereochemistry. This choice ensures consistency with the substitution notation while avoiding ambiguity, as **2+3=5** is already used to represent a different type of substitution. Additionally, since carbon 5 is covalently bound to carbon 6, this representation remains chemically intuitive (Fig. 6).
  - Chiral modifications at positions 1 and/or 4 will be marked with an extra string starting by RS14, followed by a sequence of digits 0, 1, 4, or 5, depending on the number and location of chiral mutations at positions 1, 4 or both in each residue.

Note that the fields **#Substitutions** and **#GPUs** are redundant, as they can be easily computed from the **Positions** field; however, they do not significantly increase the length of the string, as only one character is needed to include this information, and they enhance human readability as well as comprehension. Additionally, the code used to label the substitution position at each GPU is clear, concise, and easy to process. A single character represents both the number of substitutions and their locations at each GPU, while the character's position in the sequence indicates the corresponding GPU's order within the cyclic structure. The digits used to indicate the number and location of substitutions within each GPU are intuitive: **0** represents no substitution, while **2, 3, and 6** represent single substitutions at the corresponding positions. Double substitutions are indicated by digits resulting from the sum of the substituted positions (**5, 8, 9**) and, by convention, **1** is used for triple substitutions since the sum of the three digits is 11, and it is not practical to use two digits to represent a single GPU in the sequence.

371  
372  
373  
374  
375

SBE\_5\_7x0323066

which represents a  $\beta$ -CD with five SBE substitutions at specific positions indicated by the sequence 0323066 (Fig. 5).



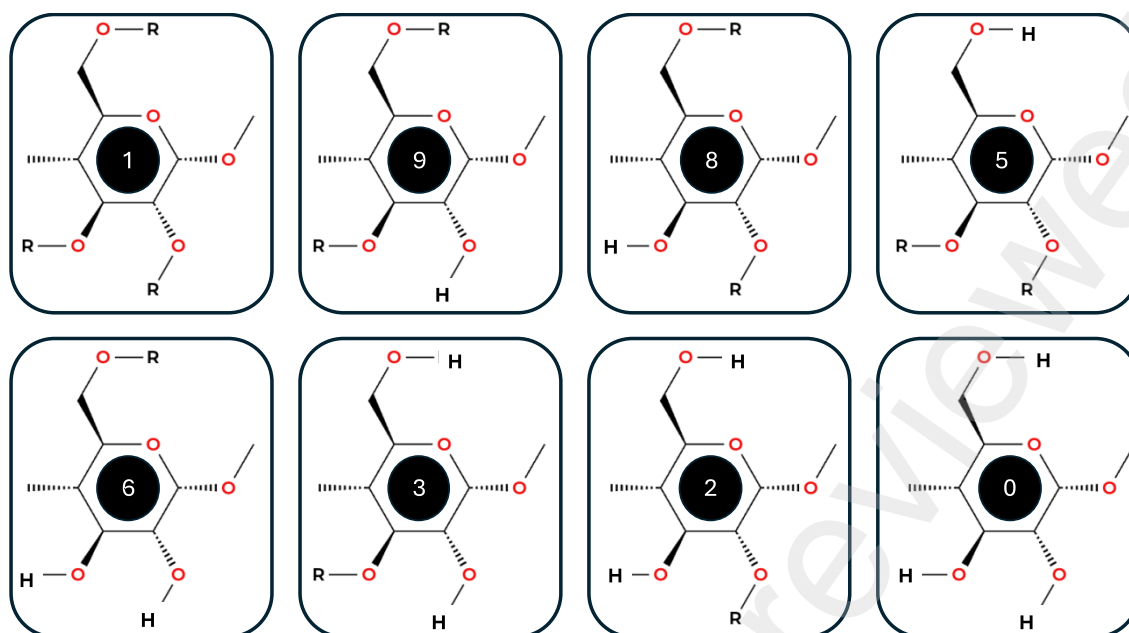
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387

**Figure 5.-** Structure corresponding to SBE\_5\_7x0323066.

Eventual chiral modifications at different positions of this structure would be represented by a longer string:

SBE\_5\_7x0323066\_RS236\_6\_7x0015200\_RS14\_2\_7x0001004

In this last example, although the string is relatively long for a monomer, it provides a complete and unambiguous specification of the structure, detailing the type and location of all substitutions as well as the stereochemical configurations of all chiral centers.



388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399

**Figure 6.-** Schematic representation of the encoding system for chemical modifications on a glucopyranoside ring. Each ring is labeled with a single-digit code displayed in a black circle, representing specific substitution patterns at positions 2, 3, and 6. A ring without modifications (all positions with hydrogen) is assigned the number 0. Single substitutions are encoded as 2 (position 2), 3 (position 3), or 6 (position 6). Double substitutions are summed as follows: 5 for positions 2 and 3 (2+3), 8 for positions 2 and 6 (2+6), and 9 for positions 3 and 6 (3+6). A triple substitution at positions 2, 3, and 6 is represented by the number 1, ensuring single-digit representation despite the sum  $2+3+6 = 11$ . An equivalent numerical code is used for chiral modifications around the same positions.

### 3.2.- Representation of hybrid modifications for CD monomers

400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415

For CD monomers with multiple modification types, each is represented separately, ordered by the number of substitutions (descending) and alphabetically if such numbers are equal. For example:

SBE\_6\_7x2306608\_HP\_4\_7x0030010

denotes a  $\beta$ -CD with both SBE and HP modifications. In the given example, all GPUs have at least one substitution. GPUs 1, 2, 4, 5, and 7 have SBE substitutions at positions C-2, C-3, C-6, C-6, and C-2+C-6, respectively, while GPUs 3 and 6 have HP substitutions at positions C-3 and C-2+C-3+C-6, respectively. This notation allows for hybrid mutations even within the same GPU, provided that the two substitution types do not occur at exactly the same position. Furthermore, this encoding is not only easy to validate computationally but also easy for a human to interpret. Another example with just one substitution per GPU is given by (Fig. 7):

SBE\_5\_7x2306206\_HP\_2\_7x0030060

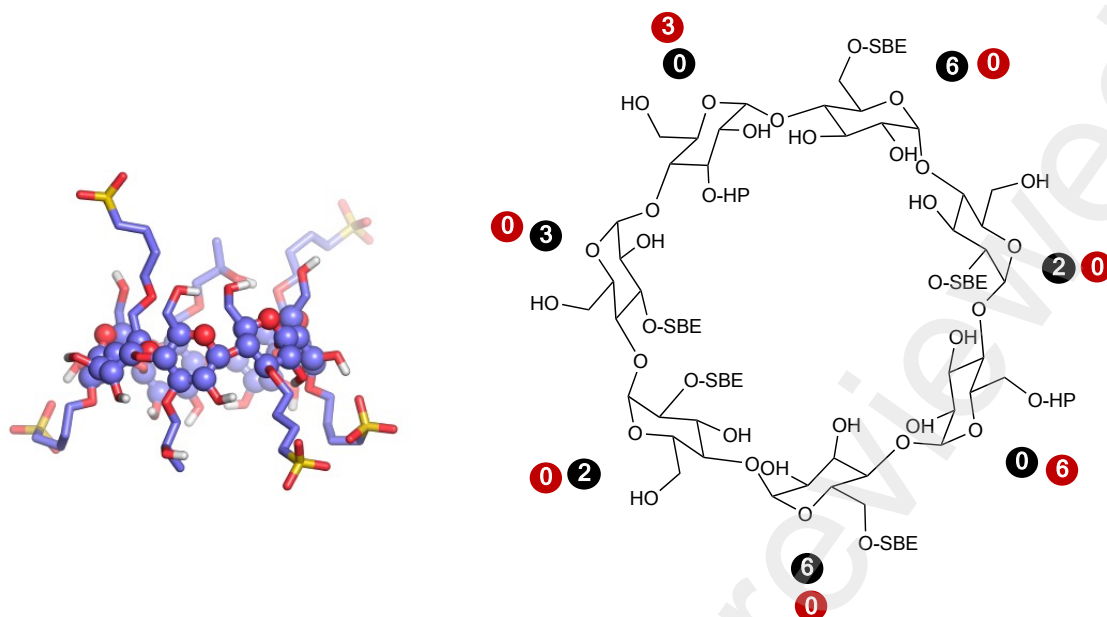


Figure 7.- Structure corresponding to SBE\_5\_7x2306206\_HP\_2\_7x0030060.

### 3.3.- Algorithm for ensuring uniqueness for CD monomers

As previously discussed, CD monomers have cyclic symmetry, meaning that for a CD composed of  $n$  GPUs, there are  $n$  possible sequences that describe the same structure. This applies to CDs with either homogeneous or heterogeneous substitution patterns. To eliminate ambiguities caused by cyclic permutations, a straightforward algorithm that can be easily implemented in a computational code has been developed. The algorithm generates all possible cyclic permutations of a given sequence and assigns a unique numerical value to each permutation. This value is simply the number obtained by concatenating the digits describing each permutation. In mathematical terms, it would be obtained as follows:

$$P_i = \sum_{j=1}^n d_j \times 10^{(n-j)} \quad [2]$$

where  $d_j$  is the digit at position  $j$  in the permutation. The sequence with the smallest resulting value is then selected as the canonical representation for the structure. When several substitution types are present, this process is applied to the first substitution type, and the same permutation is then used for the remaining modifications.

**Example:** Consider the following sequence of substitutions for a CD monomer with 6 GPUs:

**230660**

The cyclic permutations and the corresponding values of the metric given by equation [2] are:

P<sub>1</sub>: 230660  
 P<sub>2</sub>: 306602  
 P<sub>3</sub>: 066023  
 P<sub>4</sub>: 660230  
 P<sub>5</sub>: 602306  
 P<sub>6</sub>: 023066

452 The smallest value corresponds to  $P_6 = 023066$ . Therefore, the canonical representation for  
453 the sequence is:

454  
455 **023066**

456  
457  
458 **3.4.- Description of oligomers and linkers**

459 For dimeric, and potentially multimeric, CD structures covalently linked together, the complete  
460 string consists of the independent contributions from each CD monomer, along with an additional  
461 segment representing one or more linkers between subunits. These three components are  
462 concatenated together. The definition of each monomer follows the approach outlined in the  
463 previous section, which already accounts for any possible substitution pattern in each sequence.  
464 The substring corresponding to the linker(s) is placed between the substrings of the monomers it  
465 connects. The format for this linker substring is as follows:

466 `_[Residue ID][Position]_[Linker Name]_[Residue ID][Position]_`

467  
468 where the first [Residue ID][Position] refers to the GPU number and the position it is  
469 linked to in the first monomer, while the second equivalent chain provides the same information  
470 for the second monomer. Thus, for a general oligomer of CDs the nomenclature would be:

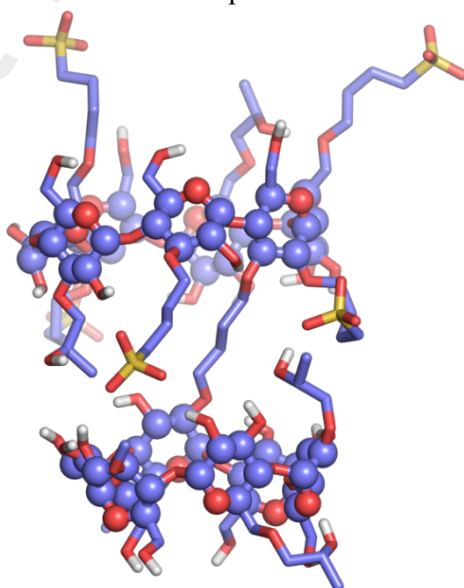
471  
472 `[CD1]_[Linker(s) 1-2]_[CD2]_[Linker(s) 2-3]_[CD3]_...`

473  
474 where [CD $i$ ] represents the string corresponding to the  $i$ -th CD, and [Linker(s)  $i$ - $j$ ]  
475 represents the string corresponding to the linker(s) joining CD $i$  with CD $j$ .

476  
477  
478 **Example:**

479 `SBE_6_7x2306603_HP_2_7x0030060`  
480 `13_BUT_32_`  
481 `HP_2_6x200006_ME_2_6x006030`

482  
483 This represents a highly complex dimeric structure, where the first CD is connected to the second  
484 via a BUT (butyl) linker, attached at residue 1, position 3 of the first CD, and residue 3, position  
485 2 of the second CD (Fig. 8). The first CD has the substitution pattern described in section 3.2,  
486 while the second subunit is a CD composed of six GPUs with two HP substitutions (2-  
487 hydroxypropyl), one at position 2 of the first GPU and another at position 6 of GPU 6.  
488 Additionally, there are two ME substitutions at positions 6 and 3 of GPUs 3 and 5, respectively.



489

490 **Figure 8.-** Structure corresponding to SBE\_6\_7x2306603\_HP\_4\_7x0030060\_  
491 13\_BUT\_32\_HP\_2\_6x200006\_ME\_2\_6x006030.

### 492 493 494 **3.5.- Complementary short notation**

495 In addition to the full or expanded version of the nomenclature, a more compact version is  
496 introduced for cases where detailed structural information is not required. This compact version  
497 of the notation simplifies the description of modified CDs by focusing just on the type and number  
498 of substitutions without specifying their exact positions on the GPUs. Although this approach  
499 results in a loss of precision compared to the expanded notation, it remains useful in experimental  
500 contexts where the exact locations of substitutions are either unknown or considered irrelevant.  
501 This situation is common when using CDs as excipients, where specific recognition of a single  
502 ligand is less important than the ability to increase the availability and solubility of different  
503 ligands. The full version of the notation will be referred to as **CyDexID-E**, while the compact or  
504 short version will be referred to as **CyDexID-S**. The compact notation condenses the structural  
505 information into a single string per subunit, indicating the total number of substitutions and the  
506 type(s) of chemical groups involved, but omitting details about their distribution across the CD  
507 ring.

508  
509 For instance, the expanded or full notation for a SBE-substituted  $\beta$ -CD with five modifications at  
510 specific positions on the GPUs might be expressed as SBE\_5\_7x2306603 (CyDexID-E). In  
511 the short form, this could simply be written as SBE5\_7u (CyDexID-S), where SBE5 indicates  
512 the number and type of substitutions, 7 denotes the number of glucopyranose units, and u serves  
513 as a placeholder to indicate that the exact positions are unspecified. Another example with two  
514 different types of substitutions would be SBE\_5\_7x2306206\_HP\_2\_7x0030060 in the  
515 expanded form (CyDexID-E), which becomes SBE5HP2\_7u in the compact form (CyDexID-S).  
516 Finally, the heterogeneous dimer given by  
517 SBE\_6\_7x2306603\_HP\_4\_7x0030060\_13\_BUT\_32\_HP\_2\_6x200006\_ME  
518 \_2\_6x006030 in the expanded form (CyDexID-E) is written as  
519 SBE6HP4\_7u\_BUT1\_HP4ME2\_6u in the compact format (CyDexID-S).

520  
521 The compact notation is especially useful when comparing distributions of modified CDs, such  
522 as in experimental samples, where the average DS is more relevant than the precise substitution  
523 pattern, or in cases where experimental (distribution of isomers) and computational (specific  
524 isomers) results are being compared. Despite its ambiguity, the short notation retains enough  
525 information to facilitate comparison to some extent between samples and computational  
526 predictions, while providing a simpler and more manageable format. One could even imagine  
527 defining mixtures of CDs for specific uses by their percentages in an actual sample and  
528 concatenating the names together either for clear labeling and/or more robust computational  
529 modeling of mixtures.

### 530 531 **4.- Structured data file format**

532 Although the defined string satisfies the requirements of an ideal nomenclature, it remains  
533 difficult to generate manually. To simplify this process, a structured, human-readable data file in  
534 a YAML-like format[38] was developed to define the CD structure. This file can be easily parsed  
535 by a Python script to automatically generate the full string. This structured format not only aids  
536 in generating the string but is also designed for computational analysis, providing clear and easily  
537 parsed sections for both human and software interpretation. The data file for a CD dimer has the  
538 following format:

```
539 CD1: [Modification Details for First Monomer]  
540 CD2: [Modification Details for Second Monomer (if applicable)]  
541 Linker_name: [Name(s) of Linker(s)]  
542
```

543 ini\_res/pos: [Initial Residue/Position for Each Linker]  
544 end\_res/pos: [End Residue/Position for Each Linker]

#### 545 546 **4.1.- Encoding CD sequences**

547 The sequence of each CD monomer is detailed on a separate line, with modifications listed as  
548 space-separated entries for each residue:

549  
550 [Type][Position] [Type][Position] [Type][Position]...

551  
552 where Type represents again the type of modification and position the number and location of  
553 the substitutions in the corresponding GPU.

554  
555  
556 CD1: SBE6 SBE6 SBE6 SBE6 SBE6 SBE6 SBE6

557  
558 represents a  $\beta$ -CD with SBE modifications at position 6 on all seven residues. The extension for  
559 hybrid modifications is trivial:

560  
561 CD1: SBE6 HP8 SBE2 ME1 HP5 ME9

562  
563 while hybrid modifications within the same residue can also easily be specified:

564  
565 CD1: SBE6 SBE3\_HP8 SBE2\_HP6 ME1 HP5 ME9

566  
567 which can be translated in the following CyDexID-E code:

568  
569 HP\_5\_6x086050\_ME\_5\_6x000109\_SBE\_3\_6x632000

570  
571 which collapses to:

572  
573 HP5ME5SBE3\_6u

574  
575 when using the CyDexID-S nomenclature.

#### 576 577 578 **4.2.- Encoding linker information**

579 Linkers are only present in dimers and higher order oligomers. As in the case of substitution types,  
580 different chemical groups can be employed to covalently join multiple CD subunits. As explained  
581 above, the chemical group and the position to which it is bound to both CD monomers must be  
582 described. Additionally, it is possible to have several linkers joining the same subunits. This  
583 information is supplied by the Linker\_name, the ini\_res/pos and the end\_res/pos  
584 fields of the file.

585  
586 **Example:**

587 Linker\_name: BUT AMD HEP  
588 ini\_res/pos: 1/2 4/3 6/2  
589 end\_res/pos: 1/2 2/3 6/6

590  
591 where we would have three different linkers: a butyl (BUT) linker connecting residue 1,  
592 position 2 of CD1 to residue 1, position 2 of CD2; an amide (AMD) linker connecting residue  
593 4, position 3 of CD1 to residue 2, position 3 of CD2; and a heptyl (HEP) linker connecting  
594 residue 6, position 2 of CD1 to residue 6, position 6 of CD2. Writing/reading this structure  
595 file as well as parsing it to the string-based nomenclature, and vice versa, is straightforward.  
596

597 Let's see a complete example that could even be interesting from the practical point of view  
598 (Fig. 9):  
599

600 **File Format Representation:**

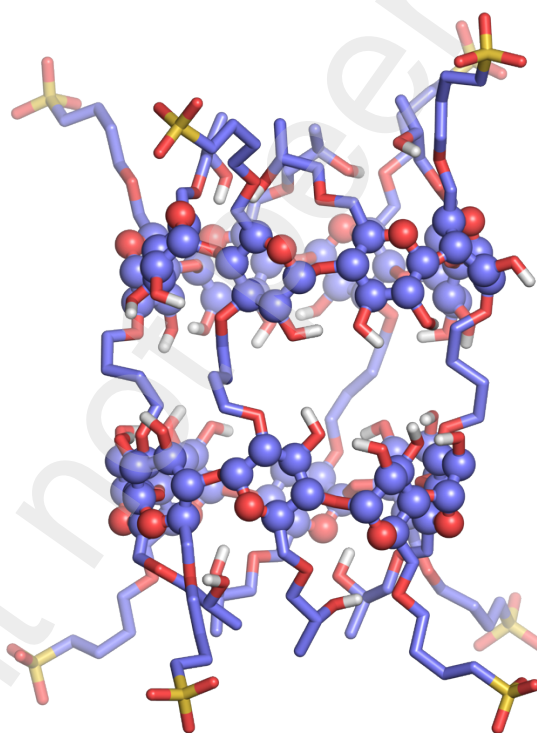
601  
602 CD1: SBE6 HP6 SBE6 HP6 SBE6 HP6 SBE6 HP6  
603 CD2: SBE6 HP6 SBE6 HP6 SBE6 HP6 SBE6 HP6  
604 Linker\_name: BUT BUT BUT BUT  
605 ini\_res/pos: 1/3 3/2 5/3 7/2  
606 end\_res/pos: 6/2 4/3 2/2 8/3  
607

608 which can be mapped to the following CyDexID-E string:

609  
610 HP\_4\_8x06060606\_SBE\_4\_8x60606060  
611 \_13\_32\_53\_72\_BUT\_BUT\_BUT\_BUT\_62\_43\_22\_83\_  
612 HP\_4\_8x06060606\_SBE\_4\_8x60606060  
613

614 and to the more compact CyDexID-S string:

615  
616 HP4SBE4\_8u\_BUT4\_HP4SBE4\_8u  
617



618  
619 **Figure 9.-** Structure corresponding to HP\_4\_8x06060606\_SBE\_4\_8x60606060\_13\_  
620 32\_53\_72\_BUT\_BUT\_BUT\_BUT\_62\_43\_22\_83\_HP\_4\_8x06060606\_SBE\_4\_8x6  
621 0606060 in the CyDexID-E format, which can be projected to the CyDexID-S notation as  
622 HP4SBE4\_8u\_BUT4\_HP4SBE4\_8u.  
623

624  
625 The structure shown in Figure 9 is elaborate, as it includes four linkers between the two subunits  
626 and eight substitutions of two different types (SBE and HP) on each monomer. Both the structure  
627 file and the CyDexID-E code describe the structure unambiguously, including the specific  
628 connections between the CD subunits for each linker. This structure could potentially be useful  
629 as a host for a variety of molecules. Notably, the CyDexID-S code for this molecule applies to  
630 many structures with the same components but different linker and substitution distributions.



631 Changes, such as placing SBE or HP substitutions at different positions or altering the symmetry  
632 of sequences, would significantly affect the molecule's properties, including cavity dynamics and  
633 encapsulation specificity. This highlights the critical importance of specifying the precise  
634 locations of all substitutions and linkers for functional accuracy.  
635

## 636 5.- Functional groups

637 One of the challenges in implementing the proposed nomenclature is the need for a concise and  
638 standardized code to represent functional groups. A unique three-letter code for each functional  
639 group would greatly enhance the clarity and usability of the nomenclature, allowing for the  
640 systematic, universal, and compact description of modified CDs. This is particularly important  
641 given the diversity of functional groups that can be introduced in CDs, as well as the need to  
642 maintain consistency and brevity in both experimental and computational contexts. While existing  
643 molecular notations, such as SMILES and SMARTS, provide robust frameworks for describing  
644 entire molecular structures or querying substructures, they do not offer standalone, standardized  
645 representations for functional groups. A three-letter coding system would streamline  
646 classification and enhance data interoperability across software, databases, and theoretical  
647 chemistry workflows.  
648

649 A review of the literature highlights active efforts in computational functional group recognition  
650 and annotation using rule-based systems, hierarchical classifications, and semantic annotations  
651 (e.g., Toxtree, CheckMol, FGO, OWL-DL) [39–41] [42,43]. However, these approaches lack  
652 standardized alphanumeric codes. Tools like BiSSCat [44] [45,46] and OCHEM [39] provide  
653 useful frameworks for specific contexts, but they rely on database-specific or numerical identifiers  
654 rather than interoperable codes. This fragmentation underscores a significant opportunity to  
655 develop standardized alphanumeric coding systems, tailored to facilitate interoperability and  
656 streamline functional group applications across computational chemistry, cheminformatics tools,  
657 and broader chemical databases. We propose a preliminary set of two and three-letter  
658 abbreviations to encode commonly used functional groups in modified CDs (Table 1). The table  
659 below lists the functional groups along with their proposed abbreviations:  
660

661 **Table 1.-** List of representative functional groups along  
662 with their proposed abbreviations.

Substitution	Abbreviation
2-Hydroxypropyl	HP
4-Sulfobutyl	SBE
Carboxythioether	CTE
Carboxymethyl	CME
Quaternary Ammonium	QA
Acetyl	AC
Tosyl	TOS
Succinyl	SUC
Benzoyl	BZ
Methyl	ME
Ethyl	ETH
1-Propyl	PRP
1-Butyl	BUT

663 These codes serve as a foundation for further refinement and standardization, addressing the gap  
664 in the literature and enabling seamless integration into both experimental workflows and  
665 computational tools.  
666

667

## 668 6.- Conclusions

669 There is much evidence indicating that the properties of modified CD molecules depend on the  
670 number of GPUs, as well as on the type, number, and even the specific location of substituted  
671 groups[47]. However, standard nomenclature systems, which only account for the number of  
672 GPUs, the type of substitution, and the average number of substitutions, can represent thousands  
673 or even millions of different structures. This is particularly limiting when designing and  
674 optimizing new CD applications. This work provides detailed calculations and examples  
675 showcasing the variety of structures that can arise in a sample with a given DS, including cases  
676 where the same number and types of substitutions occur at different locations. Because these  
677 details may be more or less relevant depending on the context, two forms of notation are proposed:  
678 an expanded version (CyDexID-E), which provides detailed information and can unambiguously  
679 describe any CD structure, and a compact version (CyDexID-S), which is useful in contexts where  
680 the exact position of substitutions is either unknown or unnecessary. This dual approach allows  
681 researchers to choose the level of detail appropriate for their specific application, ensuring both  
682 accuracy and simplicity where needed. CyDexID-E codes can be paired with structured data files  
683 for enhanced clarity and computational workflows, enabling seamless representation and  
684 scalability for diverse CD modifications, including stereochemical and polymerizable groups.  
685

686 Several examples of CD structures with increasing complexity are provided to demonstrate the  
687 method's ability to describe them unambiguously. Some examples intentionally push the  
688 boundaries of structural complexity—while remaining within the bounds of chemical laws—to  
689 rigorously test and challenge the proposed nomenclature system. Furthermore, to maximize  
690 standardization, a list of three-letter abbreviations for functional groups typically employed in  
691 modified CDs is proposed. In the present work, our discussion is limited to native glucose-based  
692 CD units. Thus, 3,6-monoanhydro-CD or per-anhydro forms are not considered within the CD  
693 scaffold. Furthermore, the nomenclature proposed here only considers structures where all units  
694 are glucose-based. The current proposal should be extended in order to explicitly consider these  
695 groups. A similar strategy used to label substitutions could be employed for this aim.  
696

697 The expanded version of the proposed nomenclature system for CD derivatives satisfies all the  
698 essential requirements aimed at solving the problems of current nomenclature systems:  
699 uniqueness, completeness, comprehensiveness, scalability, computational efficiency, human  
700 readability, and alignment with established protein nomenclature systems. A key strength of this  
701 nomenclature system is its flexibility, enabling it to accommodate any kind of CD modifications,  
702 from simple single substitutions in CD monomers to complex hybrid structures and covalently  
703 linked oligomers with any number and location of the linkers between subunits. Additionally, the  
704 proposed nomenclature considers stereochemical modifications on any of the chiral centers of the  
705 glucose units. Functional groups that can polymerize are denoted by a suffix indicating the  
706 polymerization degree. Its scalability ensures that it can adapt to the growing diversity of CD-  
707 based applications, whether in drug delivery, specific molecular recognition, catalysis, or  
708 supramolecular chemistry.  
709

710 The proposed system is highly suitable for computational applications, offering a format that is  
711 easy to parse and manipulate. It enables automatic generation of CD structures, parameterization  
712 for molecular docking and dynamics simulations, and the creation of databases to support  
713 machine learning models for predicting CD properties. The structured data file format  
714 complements the string-based nomenclature, streamlining workflows and facilitating efficient  
715 analysis of complex CD structures.  
716

717 One of the standout features of this system is its parallelism with protein nomenclature, allowing  
718 the use of familiar terms such as residues, mutations, and sequences. Moreover, it facilitates the  
719 classification of structural hierarchies, from primary sequences to secondary intramolecular  
720 patterns, tertiary oligomerization through linkers, and even quaternary non-covalent aggregates.

721 Just as the universe of proteins is described by the term proteome[48], that of lipids by  
722 lipidome[49], and that of saccharides by the glycome[50], the entire collection of cyclodextrins  
723 (as a unique subset of oligosaccharides within the glycome) can be distinguished as the  
724 cyclodextrinome, with each compound unambiguously identified by its CyDexID-E.  
725

726 To encourage wide adoption of this nomenclature, we have developed an interactive web-based  
727 tool named CyDexID Generator, which automatically assigns standardized names to modified  
728 CD monomers. The tool is accessible at: <https://cydexid-app.lm.r.appspot.com>. This tool can also  
729 convert between SMILES, IUPAC, and CyDexID E/S nomenclatures, and even produce  
730 corresponding 3D structures. We hope the clear advantages of our proposal presented in this work  
731 and the availability of an automated tool contribute to rapid and widespread adoption of our  
732 proposed nomenclature. Additionally, we commit to using this nomenclature in our future  
733 publications, which will help to demonstrate its utility and facilitate its broader acceptance.  
734 Finally, although we believe our proposal is solid, we cannot disregard future optimizations. Thus,  
735 in order to ensure future adaptability, a version identifier could be included as a prefix at the start  
736 of the CyDex-ID string, allowing for backward compatibility and seamless integration of potential  
737 extensions without disrupting existing representations.  
738

739 In summary, the proposed nomenclature system addresses key limitations in the current landscape  
740 of CD research while opening new opportunities for both experimental and computational studies.  
741 Its ability to facilitate communication, reproducibility, and computational efficiency, as well as  
742 its integration with AI-driven drug design and other predictive modeling techniques, opens  
743 exciting possibilities for future research. By providing a more systematic and standardized  
744 approach, this system is expected to accelerate the discovery of novel CD-based compounds with  
745 optimized properties for various applications. Despite the potential challenges in gaining  
746 acceptance, we believe its long-term impact will lead to significant advancements in the rational  
747 design and optimization of CDs for diverse applications.  
748

## 749 **7. Acknowledgments**

750 This work was supported by the European Union's Horizon Europe Research and Innovation  
751 Programme (Marie Skłodowska-Curie grant agreement Bicyclos N 101130235), the Spanish  
752 Agencia Estatal de Investigación (AEI) and the ERDF (PID2019-111327GB-I00, PDC2022-  
753 133402-I00, PID2022-141534OB-I00 and CNS2023-144353), by Xunta de Galicia (ED431C  
754 2021/21 and Centro de investigación do Sistema universitario de Galicia accreditation 2023-2027,  
755 ED431G 2023/03) and the European Union (European Regional Development Fund – ERDF).  
756 We acknowledge CESGA for providing computational support.  
757

758 During the preparation of this work the author(s) used ChatGPT 4 in order to improve readability  
759 and flow of the language. After using this tool/service, the author(s) reviewed and edited the  
760 content as needed and take(s) full responsibility for the content of the publication.  
761

762 **References**

- 763 [1] A. Villiers, Sur la fermentation de la fécule par l'action du ferment butyrique, *Compte*  
764 *Rendus Des Séances de l'Académie Des Sciences (France) CXII* (1891) 536–538.
- 765 [2] T. Loftsson, M.E. Brewster, *Pharmaceutical Applications of Cyclodextrins. 1. Drug*  
766 *Solubilization and Stabilization, J Pharm Sci* 85 (1996) 1017–1025.  
767 <https://doi.org/10.1021/JS950534B>.
- 768 [3] P.C. Manor, W. Saenger, *Topography of Cyclodextrin Inclusion Complexes. III. Crystal*  
769 *and Molecular Structure of Cyclohexaamylose Hexahydrate, the (H<sub>2</sub>O)<sub>2</sub> Inclusion*  
770 *Complex, J Am Chem Soc* 96 (1974) 3630–3639. <https://doi.org/10.1021/JA00818A042>.
- 771 [4] E.M.M. Del Valle, *Cyclodextrins and their uses: a review, Process Biochemistry* 39 (2004)  
772 1033–1046. [https://doi.org/10.1016/S0032-9592\(03\)00258-9](https://doi.org/10.1016/S0032-9592(03)00258-9).
- 773 [5] Fenyvesi, M. Vikmon, L. Szente, *Cyclodextrins in Food Technology and Human*  
774 *Nutrition: Benefits and Limitations, Crit Rev Food Sci Nutr* 56 (2016) 1981–2004.  
775 <https://doi.org/10.1080/10408398.2013.809513>.
- 776 [6] G. Crini, *Review: A history of cyclodextrins, Chem Rev* 114 (2014) 10940–10975.  
777 <https://doi.org/10.1021/CR500081P>.
- 778 [7] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K.  
779 Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, A. Bridgland, C. Meyer, S.A.A.  
780 Kohl, A.J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back,  
781 S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T.  
782 Berghammer, S. Bodenstein, D. Silver, O. Vinyals, A.W. Senior, K. Kavukcuoglu, P.  
783 Kohli, D. Hassabis, *Highly accurate protein structure prediction with AlphaFold, Nature*  
784 2021 596:7873 596 (2021) 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- 785 [8] E.J. Stollar, D.P. Smith, *Uncovering protein structure, Essays Biochem* 64 (2020) 649–  
786 680. <https://doi.org/10.1042/EBC20190042>.
- 787 [9] H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N.  
788 Shindyalov, P.E. Bourne, *The Protein Data Bank, Nucleic Acids Res* 28 (2000) 235–242.  
789 <https://doi.org/10.1093/NAR/28.1.235>.
- 790 [10] J. Szejtli, *Introduction and general overview of cyclodextrin chemistry, Chem Rev* 98  
791 (1998) 1743–1753. <https://doi.org/10.1021/CR970022C>.
- 792 [11] Á. Piñeiro, J. Pipkin, V. Antle, R. Garcia-Fandino, *Aggregation versus inclusion*  
793 *complexes to solubilize drugs with cyclodextrins. A case study using sulphobutylether-β-*  
794 *cyclodextrins and remdesivir, J Mol Liq* 343 (2021) 117588.  
795 <https://doi.org/10.1016/J.MOLLIQ.2021.117588>.
- 796 [12] M.E. Davis, M.E. Brewster, *Cyclodextrin-based pharmaceuticals: past, present and future,*  
797 *Nature Reviews Drug Discovery* 2004 3:12 3 (2004) 1023–1035.  
798 <https://doi.org/10.1038/nrd1576>.
- 799 [13] T. Hao, *Burnside's Lemma and Its Applications in Combinatorics Problems, Highlights*  
800 *in Science, Engineering and Technology* 47 (2023) 126–130.  
801 <https://doi.org/10.54097/HSET.V47I.8175>.
- 802 [14] W. Burnside, *Title: Theory of Groups of Finite Order, (2012).*  
803 <https://www.gutenberg.org/ebooks/40395> (accessed December 4, 2024).
- 804 [15] T. Kraus, M. Buděšínský, J. Závada, *General approach to the synthesis of persubstituted*  
805 *hydrophilic and amphiphilic β-cyclodextrin derivatives, Journal of Organic Chemistry* 66  
806 (2001) 4595–4600. <https://doi.org/10.1021/jo010046q>.
- 807 [16] X.J. Chen, G.L. Yang, X.D. Xu, J.J. Sheng, J. Shen, H.X. Dong, *Preparation and*  
808 *chromatographic evaluation of β-cyclodextrin derivative CSPs bearing substituted*  
809 *phenylcarbamate groups for HPLC, J Liq Chromatogr Relat Technol* 39 (2016) 647–657.  
810 <https://doi.org/10.1080/10826076.2016.1227993>.
- 811 [17] B. Wang, E. Zaborova, S. Guieu, M. Petrillo, M. Guitet, Y. Blériot, M. Ménand, Y. Zhang,  
812 M. Sollogoub, *Site-selective hexa-hetero-functionalization of α-cyclodextrin an*  
813 *archetypical C<sub>6</sub>-symmetric concave cycle, Nat Commun* 5 (2014).  
814 <https://doi.org/10.1038/ncomms6354>.

- 815 [18] Z. Fülöp, S. Kurkov, T. Nielsen, K. Larsen, T. Loftsson, Self-assembly of cyclodextrins:  
816 formation of cyclodextrin polymer based nanoparticles, 22 (n.d.) 2012.
- 817 [19] A.M. Anderson, I. Manet, M. Malanga, D.M. Clemens, K. Sadrafi, Á. Piñeiro, R. García-  
818 Fandiño, M.S. O'Connor, Addressing the complexities in measuring cyclodextrin-sterol  
819 binding constants: A multidimensional study, *Carbohydr Polym* 323 (2024) 121360.  
820 <https://doi.org/10.1016/J.CARBPOL.2023.121360>.
- 821 [20] P.F. Garrido, M. Calvelo, A. Blanco-González, U. Veleiro, F. Suárez, D. Conde, A.  
822 Cabezón, Á. Piñeiro, R. Garcia-Fandino, The Lord of the NanoRings: Cyclodextrins and  
823 the battle against SARS-CoV-2, *Int J Pharm* 588 (2020) 119689.  
824 <https://doi.org/10.1016/J.IJPHARM.2020.119689>.
- 825 [21] P. Saokham, A. Sá Couto, A. Ryzhakov, T. Loftsson, The self-assemble of natural  
826 cyclodextrins in aqueous solutions: Application of miniature permeation studies for  
827 critical aggregation concentration (cac) determinations, *Int J Pharm* 505 (2016) 187–193.  
828 <https://doi.org/10.1016/j.ijpharm.2016.03.049>.
- 829 [22] A. Anderson, Á. Piñeiro, R. García-Fandiño, M.S. O'Connor, Cyclodextrins: Establishing  
830 building blocks for AI-driven drug design by determining affinity constants in silico,  
831 *Comput Struct Biotechnol J* 23 (2024) 1117–1128.  
832 <https://doi.org/10.1016/J.CSBJ.2024.02.011>.
- 833 [23] D.O. Thompson, Cyclodextrins-Enabling Excipients: Their Present and Future Use in  
834 Pharmaceuticals, *Critical Reviews in Therapeutic Drug Carrier Systems* 14 (1997) 1-104.  
835 <https://doi.org/10.1615/CritRevTherDrugCarrierSyst.v14.i1.10>.
- 836 [24] H.A. Favre, W.H. Powell, *Nomenclature of Organic Chemistry: IUPAC*  
837 *Recommendations and Preferred Names 2013*, Royal Society of Chemistry, 2014.
- 838 [25] C.J. Easton, S.F. Lincoln, *Modified Cyclodextrins: Scaffolds And Templates For*  
839 *Supramolecular Chemistry*, (1999).  
840 <https://books.google.es/books?id=cM42DwAAQBAJ&printsec=frontcover&hl=es#v=onepage&q&f=false> (accessed October 25, 2024).
- 841 [26] N.L. Daly, D.T. Wilson, Plant derived cyclic peptides, *Biochem Soc Trans* 49 (2021)  
842 1279–1285. <https://doi.org/10.1042/BST20200881>.
- 843 [27] A. Anderson, R. García-Fandiño, Á. Piñeiro, M.S. O'Connor, Unraveling the molecular  
844 dynamics of sugammadex-rocuronium complexation: A blueprint for cyclodextrin drug  
845 design, *Carbohydr Polym* 334 (2024) 122018.  
846 <https://doi.org/10.1016/J.CARBPOL.2024.122018>.
- 847 [28] O. Trott, A.J. Olson, AutoDock Vina: Improving the speed and accuracy of docking with  
848 a new scoring function, efficient optimization, and multithreading, *J Comput Chem* 31  
849 (2010) 455–461. <https://doi.org/10.1002/jcc.21334>.
- 850 [29] J. Eberhardt, D. Santos-Martins, A.F. Tillack, S. Forli, AutoDock Vina 1.2.0: New  
851 Docking Methods, Expanded Force Field, and Python Bindings, *J Chem Inf Model* 61  
852 (2021) 3891–3898. <https://doi.org/10.1021/acs.jcim.1c00203>.
- 853 [30] D. Van Der Spoel, E. Lindahl, B. Hess, G. Groenhof, A.E. Mark, H.J.C. Berendsen,  
854 GROMACS: Fast, flexible, and free, *J Comput Chem* 26 (2005) 1701–1718.  
855 <https://doi.org/10.1002/jcc.20291>.
- 856 [31] H.J.C. Berendsen, D. van der Spoel, R. van Drunen, GROMACS: A message-passing  
857 parallel molecular dynamics implementation, *Comput Phys Commun* 91 (1995) 43–56.  
858 [https://doi.org/10.1016/0010-4655\(95\)00042-E](https://doi.org/10.1016/0010-4655(95)00042-E).
- 859 [32] S. Páll, A. Zhmurov, P. Bauer, M. Abraham, M. Lundborg, A. Gray, B. Hess, E. Lindahl,  
860 Heterogeneous parallelization and acceleration of molecular dynamics simulations in  
861 GROMACS, *Journal of Chemical Physics* 153 (2020). <https://doi.org/10.1063/5.0018516>.
- 862 [33] E. Lindahl, B. Hess, D. van der Spoel, GROMACS 3.0: A package for molecular  
863 simulation and trajectory analysis, *J Mol Model* 7 (2001) 306–317.  
864 <https://doi.org/10.1007/S008940100045>.
- 865 [34] B. Hess, C. Kutzner, D. van der Spoel, E. Lindahl, GROMACS 4: Algorithms for Highly  
866 Efficient, Load-Balanced, and Scalable Molecular Simulation, *J Chem Theory Comput* 4  
867 (2008) 435–447. <https://doi.org/10.1021/ct700301q>.
- 868

- 869 [35] S. Pronk, S. Páll, R. Schulz, P. Larsson, P. Bjelkmar, R. Apostolov, M.R. Shirts, J.C.  
870 Smith, P.M. Kasson, D. Van Der Spoel, B. Hess, E. Lindahl, GROMACS 4.5: A high-  
871 throughput and highly parallel open source molecular simulation toolkit, *Bioinformatics*  
872 29 (2013) 845–854. <https://doi.org/10.1093/bioinformatics/btt055>.
- 873 [36] M.J. Abraham, T. Murtola, R. Schulz, S. Páll, J.C. Smith, B. Hess, E. Lindahl, Gromacs:  
874 High performance molecular simulations through multi-level parallelism from laptops to  
875 supercomputers, *SoftwareX* 1–2 (2015) 19–25.  
876 <https://doi.org/10.1016/j.softx.2015.06.001>.
- 877 [37] D.A. Case, H.M. Aktulga, K. Belfon, D.S. Cerutti, G.A. Cisneros, V.W.D. Cruzeiro, N.  
878 Forouzes, T.J. Giese, A.W. Götz, H. Gohlke, S. Izadi, K. Kasavajhala, M.C. Kaymak, E.  
879 King, T. Kurtzman, T.S. Lee, P. Li, J. Liu, T. Luchko, R. Luo, M. Manathunga, M.R.  
880 Machado, H.M. Nguyen, K.A. O’Hearn, A. V. Onufriev, F. Pan, S. Pantano, R. Qi, A.  
881 Rahnamoun, A. Risheh, S. Schott-Verdugo, A. Shajan, J. Swails, J. Wang, H. Wei, X. Wu,  
882 Y. Wu, S. Zhang, S. Zhao, Q. Zhu, T.E. Cheatham, D.R. Roe, A. Roitberg, C. Simmerling,  
883 D.M. York, M.C. Nagan, K.M. Merz, AmberTools, *J Chem Inf Model* 63 (2023) 6183–  
884 6191. <https://doi.org/10.1021/acs.jcim.3c01153>.
- 885 [38] Yet Another Markup Language (YAML) 1.0, (n.d.). [https://yaml.org/spec/history/2001-  
886 12-10.html](https://yaml.org/spec/history/2001-12-10.html) (accessed October 25, 2024).
- 887 [39] E.S. Salmina, N. Haider, I. V. Tetko, Extended functional groups (EFG): An efficient set  
888 for chemical characterization and structure-activity relationship studies of chemical  
889 compounds, *Molecules* 21 (2016). <https://doi.org/10.3390/molecules21010001>.
- 890 [40] R. Benigni, O. Tcheremenskaia, A. Worth, Computational Characterisation of Chemicals  
891 and Datasets in Terms of Organic Functional Groups-a New Toxtree Rulebase, (n.d.).  
892 <https://doi.org/10.2788/33281>.
- 893 [41] I. Sushko, E. Salmina, V.A. Potemkin, G. Poda, I. V. Tetko, ToxAlerts: A web server of  
894 structural alerts for toxic chemicals and compounds with potential adverse reactions, *J*  
895 *Chem Inf Model* 52 (2012) 2310–2316. <https://doi.org/10.1021/ci300245q>.
- 896 [42] N. Villanueva-Rosales, M. Dumontier, Describing chemical functional groups in OWL-  
897 DL for the classification of chemical compounds, in: *CEUR Workshop Proc*, 2007.  
898 <http://ontology.dumontierlab.com/cfg-owled-2007>.
- 899 [43] P. Kumar Varadwaj, T. Lahiri, Hypothesis FGO: A novel ontology for identification of  
900 ligand functional group, *Bioinformation* 2 (2007) 113–118. [www.bioinformation.net](http://www.bioinformation.net).
- 901 [44] M. Kotera, A.G. McDonald, S. Boyce, K.F. Tipton, Functional group and substructure  
902 searching as a tool in metabolomics, *PLoS One* 3 (2008).  
903 <https://doi.org/10.1371/journal.pone.0001537>.
- 904 [45] M. Korichi, V. Gerbaud, P. Floquet, A.H. Meniai, S. Nacef, X. Joulia, Computer aided  
905 aroma design I–Molecular knowledge framework, *Chemical Engineering and Processing:  
906 Process Intensification* 47 (2008) 1902–1911. <https://doi.org/10.1016/J.CEP.2008.02.008>.
- 907 [46] D. Qv, J. Su, M. Muraki, T. Hay Aka Wa, A Decoding System for a Group Contribution  
908 Method, 1992. <https://pubs.acs.org/sharingguidelines>.
- 909 [47] G. Kali, S. Haddadzadegan, A. Bernkop-Schnürch, Cyclodextrins and derivatives in drug  
910 delivery: New developments, relevant clinical trials, and advanced products, *Carbohydr*  
911 *Polym* 324 (2024). <https://doi.org/10.1016/j.carbpol.2023.121500>.
- 912 [48] M.R. Wilkins, J.C. Sanchez, A.A. Gooley, R.D. Appel, I. Humphery-Smith, D.F.  
913 Hochstrasser, K.L. Williams, Progress with proteome projects: Why all proteins expressed  
914 by a genome should be identified and how to do it, *Biotechnol Genet Eng Rev* 13 (1996)  
915 19–50. <https://doi.org/10.1080/02648725.1996.10647923>.
- 916 [49] M.R. Wenk, The emerging field of lipidomics, *Nature Reviews Drug Discovery* 2005 4:7  
917 4 (2005) 594–610. <https://doi.org/10.1038/nrd1776>.
- 918 [50] A. Varki, R.D. Cummings, J.D. Esko, P. Stanley, G.W. Hart, M. Aebi, D. Mohnen, T.  
919 Kinoshita, N.H. Packer, J.H. Prestegard, R.L. Schnaar, P.H. Seeberger, *Essentials of*  
920 *Glycobiology*, Cold Spring Harbor (NY) (2022) 892.  
921 <https://doi.org/10.1101/9781621824213>.
- 922