

# Chemical Science



Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: L. A. Gonzalez-Ortiz, L. Noriega, F. Ortiz-Chi, G. Vidales-Ayala, E. Soberanis-Cáceres, A. Meneses, A. Aspuru-Guzik and G. Merino, *Chem. Sci.*, 2025, DOI: 10.1039/D5SC05004A.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the <u>Information for Authors</u>.

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard <u>Terms & Conditions</u> and the <u>Ethical guidelines</u> still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.



View Article Online DOI: 10.1039/D5SC05004A

# **Grammar-Driven SMILES**

# Standardization with TokenSMILES

Luis Armando Gonzalez-Ortiz,<sup>1,\*</sup> Lisset Noriega,<sup>1</sup> Filiberto Ortiz,<sup>2</sup> Gabriela Vidales-Ayala,<sup>1</sup> Emmanuel Soberanis-Cáceres,<sup>1</sup> Amilcar Meneses,<sup>3,\*</sup> Alan Aspuru-Guzik,<sup>4,\*</sup> and Gabriel Merino.<sup>1,\*</sup>

<sup>1</sup> Departamento de Física Aplicada, Centro de Investigación y de Estudios Avanzados,
Unidad Mérida, km 6 Antigua Carretera a Progreso, Apdo. Postal 73, Cordemex 97310,
Mérida, Yucatán, México.

<sup>2</sup> Secihti-Departamento de Física Aplicada, Cinvestav-IPN, Antigua Carretera a Progreso km 6, Mérida, Yucatán, 97310, México.

<sup>3</sup> Departamento de Computación, Centro de Investigación y de Estudios Avanzados, Unidad Zacatenco, Av. IPN No. 2508, Apdo. Postal 07000, Col. San Pedro Zacatenco, CDMX, México.

<sup>4</sup> Department of Chemistry, University of Toronto, DB 421D, Lash Miller Chemical Laboratories, 80 St. George Street, Toronto, ON, M5S 3H6, Canada.

e-mails: gmerino@cinvestav.mx

luis.gonzalezo@cinvestav.mx

alan@aspuru.com

DOI: 10.1039/D5SC05004A

**Abstract** 

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

The redundancy of *SMILES* notation, where multiple strings can describe the same molecule, remains a challenge in computational chemistry and cheminformatics. To mitigate this issue, we introduce TokenSMILES, a grammatical framework that standardizes SMILES into structured sentences composed of context-free words. By applying five syntactic constraints (including branch limitations, balanced parentheses, and aromaticity exclusion), *TokenSMILES* 

minimizes redundant SMILES enumerations for alkanes while maintaining valence and octet

compliance through semantic parsing rules.

TokenSMILES does not replace SMILES but rather formalizes its syntax into a standardized, machine-interpretable form. This grammatical structure enables controlled generation and manipulation of valid SMILES strings, ensuring syntactic and semantic

consistency while substantially reducing redundancy.

Implemented into SmilX, an open-source tool, TokenSMILES generates valid SMILES with accuracy comparable to existing computational implementations for molecules with low hydrogen deficiency (HDI  $\leq$  4). Its applicability extends beyond alkanes through stoichiometric modifications such as bond insertion, cyclization, and heteroatom substitution. Nevertheless, challenges remain for highly unsaturated systems, where canonicalization artifacts highlight the need for dynamic feasibility checks.

By integrating linguistic principles with cheminformatics, *TokenSMILES* establishes a scalable framework for systematic chemical space exploration, supporting applications in drug

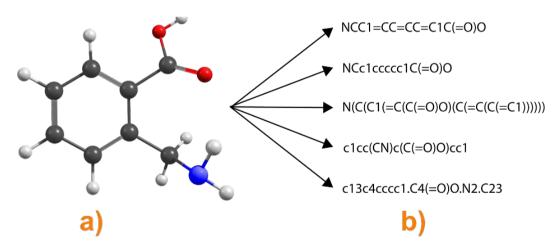
discovery, materials design, and machine learning-driven molecular innovation.

**Keywords**: SMILES, grammar, constraints, isomers

View Article Online DOI: 10.1039/D5SC05004A

#### Introduction

The analysis of chemical space using artificial intelligence relies on comprehensive and standardized molecular representations. The Simplified Molecular Input Line Entry System (SMILES), introduced by Weininger and co-workers in the 1980s, encodes molecular structures via a context-sensitive grammar using ASCII characters. <sup>1–3</sup> Initially, SMILES was designed to represent atoms, bonds, branching, cycles, aromaticity, charges, and hydrogen counts. Over time, its syntax was expanded to include stereochemistry, isotopic labeling, and hybridization. Despite its versatility, SMILES is not unique: several distinct strings can describe the same molecule (Figure 1). This redundancy arises from permissible syntactic variations within the language. Furthermore, grammatical extensions have been proposed to include more complex systems as polymers and crystal structures.<sup>4,5</sup>



**Figure 1.** 2-aminomethylbenzoic acid: (a) molecular model and (b) selected *SMILES* strings representing the same molecule, illustrating the standardization problem addressed by *TokenSMILES*.

Table 1 shows syntactic variations,<sup>6</sup> including Kekulé, aromatic, branching, and syntax and syntax aromatic variations, and syntax are syntax aromatic variations.

**Table 1.** Syntactic variations in *SMILES* notation.

Syntax style	Representative SMILES
Kekulé syntax	NCC1=CC=CC=C1C(=O)O
Aromatic syntax	NCc1cccc1C(=O)O
Branching syntax	N(C(C1(=C(C(=O)O)(C(=C(C(=C1))))))
Ring numbers and dot bond syntax	c13c4cccc1.C4(=O)O.N2.C23

Recent advances have sought to overcome these limitations by introducing alternative representations with improved syntactic control. For instance, *DeepSMILES* simplifies parenthesis handling by adopting postfixed ring-numbering rules,<sup>7</sup> while *t-SMILES* encodes functional groups explicitly, avoiding parentheses and ring indices.<sup>8</sup> *BigSMILES* extends the notation to polymers through the use of braces to denote stochastic binding patterns,<sup>8,9</sup> whereas *CurlySMILES* embeds annotations within braces { } to describe noncovalent or coordinated structures, while preserving the core *SMILES* grammar.<sup>10</sup>

Beyond these structural adaptations, new languages such as *SELFIES*,<sup>11</sup> *Group SELFIES*,<sup>12</sup> and *JAM*<sup>13,14</sup> have emerged. *SELFIES* guarantees that every token sequence corresponds to a chemically valid structures, thereby minimizing parsing and valency errors. *Group SELFIES* refines this idea by representing rings or functional groups as single tokens, simplifying substructure encoding. In contrast, *JAM* adapts *SMILES*-like syntax to describe stacking sequences in crystalline or layered materials, combining chemical and geometric information. Table 2 summarizes these languages using 2-(aminomethyl)benzoic acid as an example, highlighting improvements in grammatical clarity and robustness.

View Article Online DOI: 10.1039/D5SC05004A

Table 2. Grammatical representations of 2-(aminomethyl)benzoic acid in different notations.

Notation	Representation
SMILES	NCc1cccc1C(=O)O
DeepSMILES	NCccccc6C=O)O
t-SMILES	c1([1*])c([2*])cccc1^[1*]C(=O)O^[2*]CN
SELFIES	[N][C][C][=C][C][=C][C][=C][Ring1][=Branch1][C][=Branch1][C][=O][O]
Group SELFIES	[:benzene][Branch][:CH2NH2][pop][Branch][:COOH][pop]

In this work, we introduce *TokenSMILES*, a grammatical and graph-theoretical framework that formalizes the *SMILES* language, together with *SmilX*, its open-source implementation. *SmilX* applies the *TokenSMILES* grammar to generate and validate molecular structures based on explicitly defined syntactic and semantic rules. Both the conceptual framework and its software implementation are presented here for the first time.

The proposed grammar provides a standardized representation for organic molecules that retains the descriptive capacity of *SMILES* while minimizing ambiguity. By treating molecular strings as complete sentences governed by grammatical rules rather than collections of discrete symbols, *TokenSMILES* enables systematic computational analysis and exhaustive isomer enumeration (Figure 2). Implemented in *SmilX*, this framework departs from matrix-based approaches (e.g., *MOLGEN*, <sup>15</sup> *MAYGEN* <sup>16</sup>) and block-based methods (e.g., *SMILIB* <sup>17,18</sup>), offering a more structured and linguistically consistent paradigm for molecular representation.

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

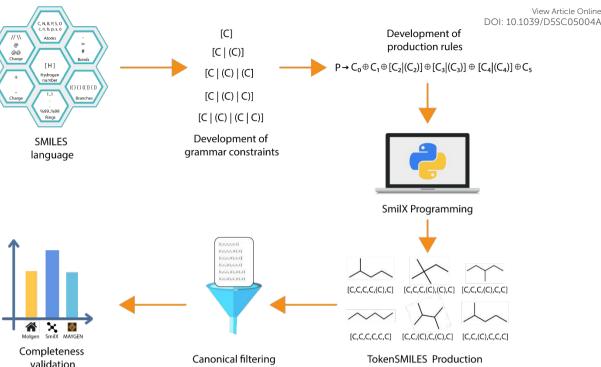


Figure 2. General workflow diagram for the development of the *TokenSMILES* grammatical framework and the validation of the number of structures.

### **Theoretical Development**

#### Kekulé Syntax for SMILES Generation

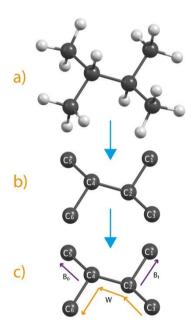
Let us start by constructing *SMILES* strings for saturated hydrocarbons. The procedure defines transversal paths across all atoms and bonds in a molecule, forming the basis for grammatical constraints in the Kekulé SMILES syntax. Using 2,3-dimethylbutane as example (Figure 3a), the process follows these steps:

- 1) **Hydrogen removal.** Generate a hydrogen-free molecular graph (Figure 3b).
- 2) **Atom labeling.** Assign unique numerical labels to non-hydrogen atoms (Figure 3b).
- 3) Path definition. Define an ordered set W representing the transversal path. In Figure 3c, the path  $P = \{(C_3, C_2), (C_2, C_4), (C_4, C_6)\}$  corresponds to  $W = \{C_3, C_2, C_4, C_6\}$ .

- 4) **Branch identification.** Identify branches  $B_m$ , where each branch contains  $S_m$  identifies  $S_m$  maximum possible number of connected atoms. The ordered set of branches is  $B = \{B_0, B_1, ..., B_m\}$ . For the system in Figure 3c,  $B_0 = \{C_4, C_5\}$  and  $B_1 = \{C_2, C_1\}$ .
- 5) **Branch insertion**. Insert the branches  $B = \{B_0, B_1\}$  into W, omitting the first atom of each branch since it already appears in W. Parentheses mark the branch boundaries:

$$\{C_3, C_2, (, C_1,), C_4, (, C_5,), C_6\}$$
 (1)

6) Symbol replacement and concatenation. Replace the atomic labels in (1) with the corresponding atomic symbols and concatenate them to obtain the string CC(C)C(C)C.



**Figure 3.** (a) Molecular model for 2,3-dimethylbutane; (b) hydrogen-free molecular graph with labeled carbon atoms; (c) traversal pathway (orange and purple arrows) covering all atoms.

The resulting *SMILES* string, CC(C)C(C)C, represents 2,3-dimethylbutane with implicit hydrogen atoms. Since *SMILES* grammar is non-commutative, the steps must be performed in the specified order.

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

View Article Online DOI: 10.1039/D5SC05004A

Our method transforms the Kekulé syntax into a standardized form that equalizes string lengths and isolates chemical information by assigning individual tokens to each atom and symbol. For example, the *SMILES* representation of 2,3-dimethylbutane, CC(C)C(C)C, can be rewritten as [C, C, (C), C, (C), C], where "(" and ")" denote the beginning and end of branches, respectively. So, in this tokenized form, branches occur at positions 1 and 3 (zero-indexed). The resulting sequence of tokens constitutes a standardized-length representation referred to as *TokenSMILES*.

This tokenization follows two sequential rules: First, the original string is parsed into individual characters, each enclosed in square brackets. For CC(C)C(C)C, this yields [C, C, (, C, ), C, (, C, ), C], maintaining the exact order of symbols in the original notation. In this representation, all symbols are enclosed within square brackets to form an ordered sequence. Although conventional set notation implies unique elements, here repeated tokens are intentionally preserved to retain positional information.

Second, the tokens are categorized according to their syntactic context. Left-context symbols, [, (, =, and #, are placed immediately before atomic symbols, while right-context symbols, @, ), %, ], and a digit n, are placed immediately after them. Applying these rules to the example yields the standardized *TokenSMILES* form [C, C, (C), C, (C), C].

#### **Grammar Constraints for TokenSMILES**

The production rules used to generate *SMILES* strings from molecular paths allow multiple valid representations for a single structure, making exhaustive *SMILES* enumeration a non-deterministic polynomial-time hard (NP-hard) problem. <sup>19,20</sup> To mitigate this redundancy, and following previous work on grammatical constraints in formal languages, <sup>21</sup> we introduce five grammar rules to reduce the number of equivalent strings representing organic isomers:

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

**Constraint 1** 

View Article Online DOI: 10.1039/D5SC05004A

For a molecular model G without ring number and dot bonds, TokenSMILES strings are constructed from traversal paths W that include all non-hydrogen atoms and bonds:

$$W = \{a_0, a_1, ..., a_{n-1}\} \text{ or } W = [a_0, a_1, ..., a_{n-1}],$$

where each  $a_i \in W$  denotes a non-hydrogen atom.

#### **Constraint 2**

Branch symbols "(" and ")" are not permitted at terminal positions  $(a_0, a_1, \text{ and } a_{n-1})$  since these locations cannot generate new expressions.

#### **Constraint 3**

To control branching, the second and penultimate atoms  $(a_2 \text{ and } a_{n-2})$  may optionally include branch symbols, *i.e.*  $[a_2 \mid (a_2)]$  and  $[a_{n-2} \mid (a_{n-2})]$ . These optional insertions produce distinct yet grammatically valid *SMILES* variants.

#### **Constraint 4**

Parentheses introduced between  $a_3$  and  $a_{n-3}$  must remain balance, with every opening parenthesis "(" matched by a corresponding closing parenthesis ")". Valid expressions include  $[a_i \mid (a_i \mid a_i) \mid (a_i)]$ , ensuring structural consistency across the entire sequence.

#### **Constraint 5**

Aromaticity symbols (c, n, b, p, s, and o) are excluded to retain the grammar to uppercase atomic symbols (C, N, B, P, S, and O).

DOI: 10.1039/D5SC05004A

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

#### Production Rules for Alkane TokenSMILES

This section presents a systematic procedure for generating TokenSMILES representations of alkanes under the grammatical constraints defined previously. Using the alphabet  $\{C, (C), (C, C)\}$ , a token dictionary is defined as  $\pi = \{C, (C), (C, C)\}$ . To construct SMILES for all  $C_nH_{2n+2}$  isomers, we define a molecular path  $W = [a_0, a_1, ..., a_{n-1}]$ , where each  $a_i \in W$ is replaced by tokens from  $\pi$  via the production rule P:

$$P \to a_0 \oplus a_1 \oplus \dots \oplus a_{n-2} \oplus a_{n-1} \tag{2}$$

Here, the concatenation operator ( $\oplus$ ) joints the strings in W, and the arrow ( $\rightarrow$ ) indicates the production process.

Applying Constraint 2,  $a_0$ ,  $a_1$  and  $a_{n-1}$  are replaced by "C" without branch symbols, resulting in:

$$P \rightarrow C \oplus C \oplus a_2 \oplus ... \oplus a_{n-2} \oplus C$$
 (3)

Next, Constraint 3 specifies that  $a_2$  and  $a_{n-2}$  may take the forms [C|(C)], yielding:

$$P \to \mathcal{C} \oplus \mathcal{C} \oplus [\mathcal{C}|(\mathcal{C})] \oplus a_3 \oplus \dots \oplus a_{n-3} \oplus [\mathcal{C}|(\mathcal{C})] \oplus \mathcal{C}$$
 (4)

For simplify, the inner sequence  $a_3 \oplus ... \oplus a_{n-3}$  is replaced with a variable  $\Omega$ , rewriting (4) as:

$$P \to C \oplus C \oplus [C|(C)] \oplus \Omega \oplus [C|(C)] \oplus C \tag{5}$$

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence.

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19:25:59.

To define the instances in  $\Omega$ , Constraint 4 is applied. Each element in  $\{a_3$ ,  $\alpha_{n,133}\}$  has co5004A four possible forms:  $[C \mid (C \mid C) \mid (C)]$ . Balanced parentheses are maintained through recursive rules.

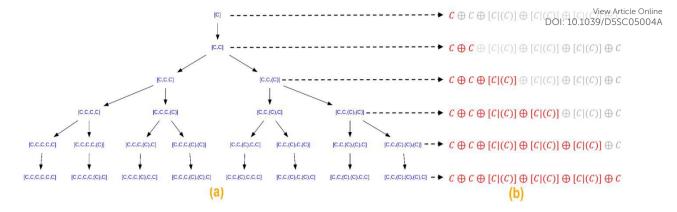
The first rule,  $q_0 \rightarrow C \mid C \mid q_0$ , generates chains such as "CC", "CCC", and "CCCC". The second,  $q_1 \rightarrow (C) \mid (C) \mid q_1$ , introduces terminal branches represented by parentheses. Combinations of the two are obtained using  $q_2 \rightarrow q_0 \mid q_1 \mid q_0 q_1 \mid q_1 q_0$ , allowing permutations of linear and branched fragments. Nested branches are introduced through  $q_3 \rightarrow (CC) \mid (Cq_3C) \mid (Cq_2C) \mid (Cq_3C) \mid (Cq$ 

$$P \to \mathcal{C} \bigoplus \mathcal{C} \bigoplus [\mathcal{C}|(\mathcal{C})] \bigoplus [q_{\mathcal{A}}] \bigoplus [\mathcal{C}|(\mathcal{C})] \bigoplus \mathcal{C} \tag{6}$$

Using (6), the number of *SMILES* representations for  $C_nH_{2n+2}$  isomers decreases drastically. For example, Figure 4 shows the construction of *SMILES* for  $C_6H_{14}$  using (6) with  $[q_4] = [C|(C)]$ :

$$P \to C \oplus C \oplus [C|(C)] \oplus [C|(C)] \oplus [C|(C)] \oplus C \tag{7}$$

However, (7) does not account for atomic equivalence or valence restrictions, which may result in redundant (e.g., [C, C, C), C, C, C]  $\equiv$  [C, C, C, C, C, C], or chemically invalid strings (e.g., [C, C, (C), (C), (C), C]). To address these limitations, the next section introduces the semantic parsing of *TokenSMILES*, which filters out chemically inconsistent structures and enforces the octet rule.



**Figure 4.** (a) A generator tree for the *SMILES* strings of the  $C_6H_{14}$  system and (b) the progress of the production rule (7) to obtain the atomic symbols in the *SMILES*.

#### **Chemical Context Definition**

Semantic parsing<sup>22,23</sup> was applied to determine the connectivity encoded in each *TokenSMILES* string. Following the conventions of Weininger et al., each atom  $(\alpha)$  is represented by its atomic symbol, and the symbol is linked to a set of chemical properties collectively referred to as the context of the atom, Context $(\alpha)$ . For our analysis, the context is defined as:

$$Context(\alpha) = (p_0, p_1, p_2)$$

where  $p_0$  is the atomic symbol,  $p_1$  corresponds to the valence, and  $p_2$  is the atomic connectivity.

As an example, the *TokenSMILES* [C, C, (C), C, (C), C] can be indexed as  $[C_0, C_1, (C_2), C_3, (C_4), C_5]$ . The chemical context for each substring is summarized in Table 3, where the first element identifies the atomic symbol ("C"), the second indicates the valence (4), and the third lists the corresponding connectivity set.

**Table 3.** Chemical context of  $[C_0, C_1, (C_2), C_3, (C_4), C_5]$ .

View Article Online DOI: 10.1039/D5SC05004A

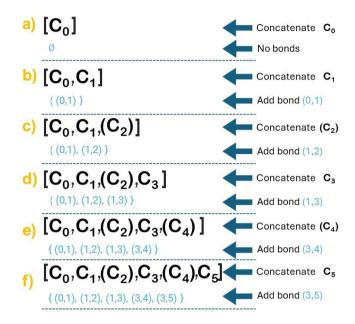
String	Chemical context
C <sub>0</sub>	( C, 4, {(0,1)})
$C_1$	( C, 4, {(0,1),(1,2),(1,3)})
$(C_2)$	(C, 4, {(1,2)})
$C_3$	(C, 4, {(1,3), (3,4), (3,5)})
$(C_4)$	(C, 4, {(3,4)})
C <sub>5</sub>	( C, 4, {(3,5)})

#### Connectivity Extraction from TokenSMILES

In *SMILES* notation, connectivity is often implicit and must be derived through semantic parsing to reconstruct chemical relationships. Figure 5 shows the procedure used to extract connectivity from the *TokenSMILES* sequence  $[C_0, C_1, (C_2), C_3, (C_4), C_5]$ :

- a) "C<sub>0</sub>" is not concatenated with any other token, so its connectivity set remains empty.
- b) " $C_1$ " is concatenated with " $C_0$ ", creating a new bond represented by the tuple (0, 1).
- c) The same procedure is applied to the remaining tokens. Figure 5c shows how  $(C_2)$  is concatenated, using parentheses to indicate the start and end of a branch emerging from " $C_1$ ", adding the tuple (1, 2).
- d) " $C_3$ " is then concatenated, forming the tuple (1, 3) rather than (2, 3), since " $(C_2)$ " is a closed branch.
  - e) " $(C_4)$ " is concatenated next, adding the tuple (3, 4).
  - f) Finally, "C<sub>5</sub>" is concatenated, producing the tuple (3, 5).

After processing all tokens, the resulting bond set is  $\{(0, 1), (1, 2), (1, 3), (3, 4), (3, 3, 7)\}$  described corresponding to [C, C, (C), C, (C), C]. Each tuple represents the connectivity in the *TokenSMILES* notation.



**Figure 5.** Procedure for extracting bonds from the *TokenSMILES* of 2,3-dimethylbutane.

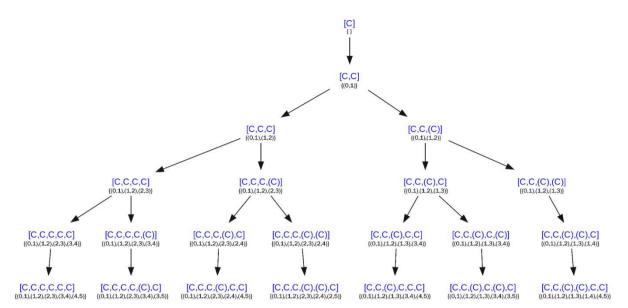
This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

As shown in Figure 5a, the bond contexts initially consist of empty sets. As the strings are concatenated according to the *SMILES* grammar and defined constraints, their contexts expands as new bonds are introduced. From this stage onward, the strings listed in Table 3 are treated as context-free strings or simply words.<sup>22</sup> These words lack explicit connectivity to other atoms, excluding hydrogen. Consequently, a *TokenSMILES* can be interpreted as a sequence of such words forming a chemically meaningful sentence.

To show this process, we return to the example of the  $C_6H_{14}$  isomers. Using production rule (7), connectivity was assigned to each *TokenSMILES* string, and those that violated the octet rule were removed, yielding seven valid *SMILES* candidates (Figure 6). To eliminate duplicates, each *TokenSMILES* was converted to a canonical *SMILES* form using the canonicalization algorithm implemented in the *RDkit* module.<sup>24</sup> Strings sharing the same

canonical form were identified, and only unique SMILES were retained. After filtering  $^{\text{virtagicle Online}}_{550004A}$  distinct SMILES remained, corresponding precisely to the five constitutional isomers of  $C_6H_{14}$ .



**Figure 6**. Production of valid *SMILES* for the  $C_6H_{14}$  isomers. Unlike Figure 5, the invalid *TokenSMILES* [C, C, (C), (C), (C), C] is omitted here because it does not satisfy the octet rule.

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19:25:59.

String lengths normalized to 6 words	String length varies from 6 to 14 symbols.
[C,C,C,C,C]	CCCCC, CCCC(C), CCC(CC), CCC(CCC), C(CC)CCC, C(CC)(CCC), C(CC)CC(C), C(CCC)C(CC), C(CCC)(CC), C(CCC)CC, C(CCC)C(C), C(CCCC)C, C(CCCC)(C), C(CCCCC), C(C)(CCC), C(C)(CCC), C(C)CC(CC), C(C)(CCCC), C(CCCC), C(C)(CCC)
[C,C,C,C,(C),C] [C,C,(C),C,C,C]	CCCC(C)C, CCCC(C)(C), CCC(C(C)C), CC(C)CCC, CC(C)CC(C), C(CC)(C(C)C), C(C(C)C), C(C(CCC)C), C(C(C)CCC), C(C(C)C)C(C), C(C(C)C)CC, C(CCC(C)C), C(CCC)(C)C, C(CCC)(C)(C), C(CC(C)C)C, C(CC(C)C)(C), C(C)(CCC)(C), C(C)(CCC)C, C(C)(C)(CCC), C(C)(C)C(C), C(C)(CCC)(C), C(C)(C)CC(C), C(C)(C)C(C)C, C(C)(C)CCC, CC(CCC)C, CC(CC(C)C), CC(CCC)(C), C(C)(CCC), CC(CCC)C, CC(CC(C)C, C(C)(C)C(C)C), CC(C)(CCC), CC(C)C(C), C(C)(C)C, C(C)(C)C(C)C)
[C,C,C,(C),C,C]	CCC(C)CC, CCC(C)C(C), CCC(C)(CC), CCC(CC)C, CCC(CC)(C), C(CC)(C)CC, C(CC)(CC)C, C(CC)(CC)(C), C(C(CC)CC), C(CC)(C)(CC), C(CC)(CC)C, C(C(C)CC)C, C(C(C)CC)(C), C(CC(C)CC), C(CC(CC)C), C(C)(CC)(CC), C(C)(C(CC)C), C(C)(CC)CC, C(C)(CC)C(C), C(C)(C(C)CC), CC(CC)CC, CC(CC)C(C), C(C)(CC)(C), C(C)(C(C)CC, C(C)C(C)CC), C(C)C(CC)C, C(C)C(C)(C), CC(C(C)CC), CC(C(CC)C), C(CC)(CC), C(C)(C(C)(C), C(C)(C)(C)) CC(C(CC)C), CC(CC)(CC), CC(C(C)(C)), C(CC)(C)(C)
[C,C,(C),C,(C),C]	C(C(C(C)C)C), C(C(C)C)(C)C, C(C(C)C)(C)(C), C(C(C)C(C)C), C(C)(C(C)C)(C)C, C(C)(C)(C)C), C(C)(C)C)C, C(C)(C)C)C, C(C)C)C, C(C)C)C, CC(C)C)C, CC(C)C,
[C,C,(C),(C),C,C] [C,C,C,(C),(C),C]	CCC(C)(C)C, CCC(C)(C)(C), C(CC)(C)(C)C, C(C(CC)(C)C), C(CC)(C)(C)(C), C(C(C)(C)CC), C(C(C)(C)CC), C(CC)(C)(C)C), C(C(C)(C)CC), C(C(C)(C)CC), C(C)(C)(C)(C), C(C)(CC)(C)C, C(C)(C)C)C), C(C)(C)(C)(C)C, C(C)(C)CC C(C)(C)(C)C, C(C)(C)(C)C, C(C)(C)(C)CC, C(C)(C)(C)CC, C(C)(C)C, C(C)(C)(C)C, C(C)(C)(C)C, C(C)(C)(C), CC(C)(C)C, C(C)(C)(C)C, C(C)(C)(C)C, CC(C)(C)C, C(C)(C)C, C(C)(C)(C)C, CC(C)(C)C, CC(C)(C)(C)C, CC(C)(C)(C)C

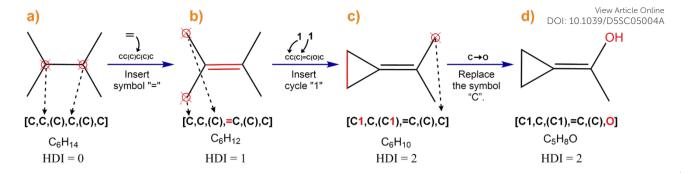
#### Modifying *TokenSMILES* Stoichiometry

This section defines the rules for modifying TokenSMILES syntax to change stoichiometry and to generate SMILES for systems beyond C<sub>n</sub>H<sub>2n+2</sub>. Three operations are introduced: two insertion operations, which add new symbols into a word, and one replacement operation, which alters atomic symbols. Each operation produces a copy of the entire sentence while recording the applied modifications. This procedure enables the systematic generation of isomers with stoichiometries different from those of alkanes.

The modification of *TokenSMILES* stoichiometry is achieved through the analysis of the online nesting levels and adjacency sets, which describe the implicit connectivities within a string. Each word in a *TokenSMILES* carries a grammar-based index that defines its nesting depth and adjacency relations. By editing these relations while preserving valence consistency, the algorithm modifies molecular topology in a rule-based manner. This approach allows direct structural transformations, such as the insertion of rings or double bonds, from the *TokenSMILES* representation without relying on IUPAC nomenclature or predefined templates (see Supplementary Information for details).

To illustrate the process, we begin with 2,3-dimethylbutane ( $C_6H_{14}$ ) and modify it to generate 1-cyclopropylideneethanol ( $C_5H_8O$ , Figure 7). First, the connectivity set is extracted as  $\{(0, 1), (1, 2), (1, 3), (3, 4), (3, 5)\}$ . Using these bonds, the algorithm identifies adjacent words representing bonds in *TokenSMILES*. Words at positions 1 and 3 (Figure 7a) are selected for double-bond insertion. To evaluate feasibility, the valence excess ( $\Delta$ ) is computed for each atom, defined as the difference between its valence and its degree (the number of incident bonds). When  $\Delta \geq I$ , a "=" symbol is inserted into the word at the higher position in the sentences (Figure 7b).

To prevent excessive "=" insertions, the hydrogen deficiency index (HDI) of the initial system ( $C_6H_{14}$ ) is used as a control variable. Each additional double bond increases the HDI by one, ensuring that the connectivity modifications remain chemically valid and stoichiometrically consistent.



**Figure 7.** Syntax transformation from 2,3-dimethylbutane *TokenSMILES* to 1-cyclopropylideneethanol *TokenSMILES*. Changes are highlighted in red, and dotted arrows indicate the modified word.

Next, to insert a cycle, two non-adjacent words are selected based on the same connectivity  $\{(0, 1), (1, 2), (1, 3), (3, 4), (3, 5)\}$ . If no bond exists between the chosen positions, the algorithm proceeds with cycle insertion. For the example in Figure 7b, words at positions 0 and 2 are chosen, and the absence of bond (0, 2) is confirmed. The valence excess for both atoms in then evaluated; if  $\Delta > 0$  for each, cycle symbols are inserted (Figure 7c). A random ring number between 1 and 9 is assigned and placed to the right of the atomic symbol (e.g., "C1"). For numbers 10-99, the cycle symbol is prefixed with "%" (e.g., "C%10"). Each cycle insertion increases the HDI by one (Figure 7c).

Finally, to substitute a carbon atom with oxygen, the condition  $degree(C) \leq valence(O)$  must be satisfied before replacement. The operation reduces the carbon and hydrogen count by one and two, respectively (Figure 7d). For heteroatoms, the allowed are B, Br, Cl, F, I, N, O, P, and S. This approach generates all structural isomers of  $C_5H_8O$  (Table SI1). Once all transformations are complete, equivalent *SMILES* are filtered using the canonicalization algorithm in *RDKit*.

SmilX Software

View Article Online
DOI: 10.1039/D5SC05004A

To automate the generation of *SMILES* under grammatical constraints, the *SmilX* program was developed as an open-source Python tool, available at <a href="https://github.com/LuisOrz/SmilX.git">https://github.com/LuisOrz/SmilX.git</a>. *SmilX* constructs *SMILES* representations of isomers while maintaining compliance with the specified stoichiometry. A user-oriented web interface was also implemented using *Streamlit*, which provides both interactive functionality and server infrastructure. The interface is included with the package and accessible at <a href="https://smilx-isogenerator.streamlit.app/">https://smilx-isogenerator.streamlit.app/</a>.

The workflow begins with a molecular formula provided as a string, such as  $C_nH_{2n+2}$ . Based on this input, SmilX generates all possible TokenSMILES corresponding to the defined stoichiometry and adjusts their syntax accordingly. The resulting words in each TokenSMILES are concatenated to form complete SMILES strings, which are then processed using a canonicalization algorithm to eliminate duplicates and retain unique structures. Finally, the software uses the RDKit module to generate stereoisomeric variants, returning a curated list of SMILES strings that satisfy the input molecular formula.

#### **Results and Discussion**

Two experiments were performed to evaluate the isomer-generation capabilities of *SmilX* and to validate the *TokenSMILES* framework. The first reproduced the data reported by Elyashberg et al.<sup>25</sup> for C–H systems (Figure 7), which serve as reference structures for more complex compositions. The resulting isomer counts from *SmilX* were compared with those obtained using *MOLGEN* (a matrix-based generator) and *MAYGEN* (an open-source, *Java*-based tool). The second experiment assessed *SmilX*'s performance in systems containing heavier elements (O, N, Cl) across a range of hydrogen deficiency indices (HDI) and atom counts.

DOI: 10.1039/D5SC05004A

omitting the "C=n" notation.

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

Isomer enumeration followed the molecular formulas reported by Elvashberg et al..<sup>25</sup> to allow direct comparison between SmilX and MAYGEN. Systems containing at least two hydrogen atoms were prioritized, while hydrogen-free cases (C<sub>6</sub> or C<sub>10</sub>) were excluded by

SmilX reproduced Elyashberg's isomer counts for most systems (Figure 8). Minor deviations occurred when the HDI approached the total heavy-atom count: SmilX occasionally yielded one missing structure (false negative, e.g., C<sub>8</sub>H<sub>4</sub>) or 1-7 additional isomers (false positives) in systems as  $C_8H_2$ ,  $C_9H_2$ ,  $C_{10}H_2$ ,  $C_9H_4$ ,  $C_{10}H_4$ ,  $C_{10}H_6$ ,  $C_9H_8$ , and  $C_{10}H_{10}$  (Figure 8).

MOLGEN accurately enumerated isomers at low HDI values but showed decreased accuracy as HDI increased (Figure 8). In contrast, MAYGEN's online version encountered memory saturation in larger systems ( $C_9H_6$ ,  $C_9H_8$ ,  $C_9H_{10}$ ,  $C_{10}H_2$ ,  $C_{10}H_4$ ,  $C_{10}H_6$ ,  $C_{10}H_8$ ,  $C_{10}H_{10}$ , C<sub>10</sub>H<sub>12</sub>, C<sub>10</sub>H<sub>14</sub>, and C<sub>10</sub>H<sub>16</sub>) but reproduced Elyashberg's data for smaller, computationally feasible cases.

MAYGEN employs a canonicalization procedure distinct from RDKit's implementation of the Weininger algorithm. 16 Consequently, canonical SMILES can struggle to distinguish equivalent atoms in molecules containing multiple nested rings, particularly when HDI is equal to or greater than half the number of heavy atoms. SmilX's reliance on RDKit likely accounts for the few observed false positives and negatives.

#### **Second Experiment**

Building on these results, the second experiment examined SmilX's robustness in systems of increasing compositional complexity: (C, H, O), (C, H, O, N), and (C, H, O, N, Cl). These compositions typically produce more isomers than pure C-H systems. For each composition, six molecular formulas were generated for HDI values ranging from 0 to 5, and the number of heavy atoms was limited to 3–10 to avoid cases where HDI  $\geq$  half the heavy atom count, previously associated with enumeration errors. Results are summarized in Tables SI1–SI3.

SmilX showed efficient performance, aided by disk-caching optimization. All three tools produced identical isomer counts for HDI  $\leq$  4, with small discrepancies emerging at higher HDI values. These deviations further support the hypothesis that RDKit's canonicalization algorithm faces difficulties in handling highly nested ring topologies.

The *TokenSMILES* framework effectively reduced both string redundancy and computational overhead through grammatical constraints and caching. While *SmilX* correctly enumerated the majority of organic systems, boundary cases where HDI approached the heavy-atom count remained problematic. These discrepancies are consistent with the theoretical limitations of canonicalization algorithms in systems containing high symmetry or multiple fused rings, rather than with specific software errors. Despite these edge-case issues, *SmilX* maintained low execution times, and the reuse of cached structures enabled scalable exploration of extensive chemical spaces.

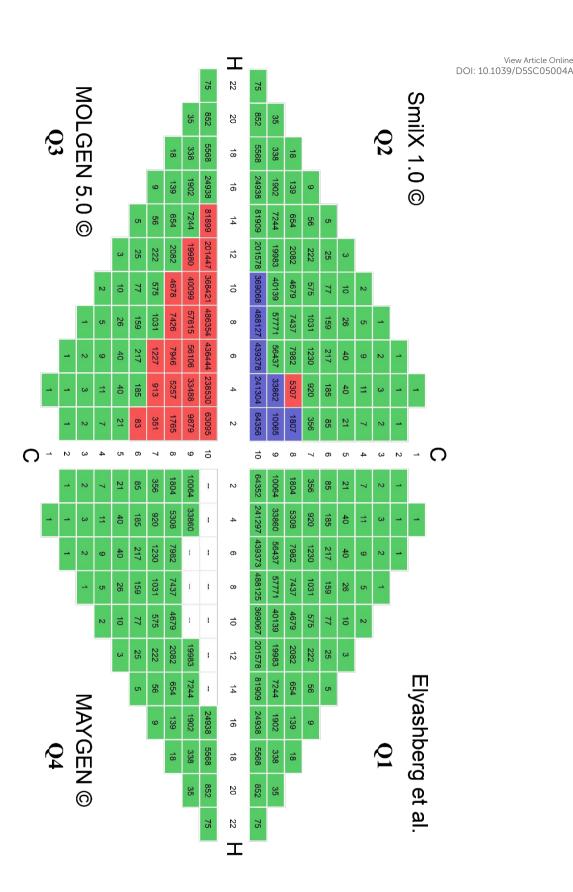


Figure 8. Isomers composed of carbon and hydrogen as reported by Elyashberg et al? (Quadram Co5004A 1) compared with those obtained using *SmilX* (Q2), *MOLGEN* (Q3), and *MAYGEN* (Q4). Green cells indicate matching isomers; blue cells represent false positives; red cells correspond to false negatives; and white cells denote missing data due to resource limitations. All results are shown relative to the reference data from Elyashberg et al.

Classical *SMILES* representations provide a foundation for cheminformatics but suffer from significant redundancy, as shown in Table 4. To overcome this limitation, the present work redefines *SMILES* not merely as atomic sequences but as grammatically structured sentences, a conceptual framework rooted in formal language theory. The *TokenSMILES* approach implements this through hierarchical syntax (word- and sentence-level organization), enforced grammatical constraints, and standardized string lengths. This reformulation reduces redundancy, ensures systematic chemical-space coverage, and facilitates new computational applications.

Unlike conventional structure generators such as *MOLGEN* or *MAYGEN*, *TokenSMILES* emphasizes formal representation rather than speed or memory efficiency, which justifies the omission of runtime benchmarks. Similarly, unlike cheminformatics toolkits such as *RDKit*, it augments rather than replaces *SMILES* syntax by operating on grammatical constructs. This strategy enables dynamic programming, partial-solution reuse, and exploration beyond the limitations of purely graph-based methods, while remaining compatible with evolutionary and machine-learning algorithms.

#### **Conclusions**

This study presents *TokenSMILES* as a grammatical framework that redefines the *SMILES* language through explicit syntactic rules. By interpreting *SMILES* as structured sentences composed of context-free words, *TokenSMILES* minimizes redundancy and enforces

grammatical consistency. Constraints on branching, parentheses balance, and paromatical reduce the number of valid SMILES variants for  $C_nH_{2n+2}$  isomers. This structured representation facilitates systematic chemical-space exploration while ensuring valence and octet compliance through semantic parsing. Integrated within SmilX, TokenSMILES performs comparably to MOLGEN and MAYGEN in generating isomers for systems with low hydrogen deficiency (HDI  $\leq 4$ ), demonstrating reliable canonical SMILES generation.

Beyond alkanes, *TokenSMILES* enables stoichiometric modifications such as bond insertion, cyclization, and heteroatom substitution, extending its applicability to broader organic systems. In high-HDI cases, minor misidentifications arise from *RDKit*'s canonicalization limitations, suggesting the need for improved feasibility checks.

Currently, *TokenSMILES* prioritizes grammatical completeness and semantic accuracy over computational efficiency. Although benchmarking was not the focus of this work, future versions could adopt optimization strategies inspired by algebraic isomer generators. The framework is inherently compatible with machine learning due to its discrete syntax, fixed-length representations, and reusable grammatical components, which enable hybrid symbolic-neural modeling and grammatical evolution.

Treating *SMILES* as grammatically structured sentences introduces a new paradigm for cheminformatics, linking linguistic theory with chemical representation. This approach supports machine-learning-based molecular design and systematic chemical-space mapping. Future extensions to polymers, organometallics, and crystalline systems may open new applications in materials and drug discovery.

## Acknowledgments

This work was supported by Cinvestav. A.G.-O. thank Conahcyt for their PhD fellowship. L.N. thanks Conahcyt for the postdoctoral fellowship.

DOI: 10.1039/D5SC05004A

References

- D. Weininger, J. Chem. Inf. Comput. Sci., 1988, 28, 31–36. DOI: 10.1021/ci00057a005.
- D. Weininger, A. Weininger and J. L. Weininger, *J. Chem. Inf. Comput. Sci.*, **1989**, *29*, 97–101, DOI: 10.1021/ci00062a008.
- 3 D. Weininger, *J. Chem. Inf. Comput. Sci.*, **1990**, *30*, 237–243, DOI: 10.1021/ci00067a005.
- 4 H. Kim, J. Na and W. B. Lee, *J. Chem. Inf. Model.*, **2021**, *61*, 5804–5814, DOI: 10.1021/acs.jcim.1c01289.
- M. Quirós, S. Gražulis, S. Girdzijauskaitė, A. Merkys and A. Vaitkus, *J. Cheminform.*,
   2018, 10, 1–17, DOI: 10.1186/s13321-018-0279-6.
- 6 C. A. James, R. Apodaca, N. O'Boyle, A. Dalke, J. van Drie, P. Ertl, G. Hutchison, G. Landrum, C. Morley, E. Willighagen, H. De Winter, T. Vandermeersch and J. May, OpenSMILES specification, Version 1.0, 2016-05-15, http://opensmiles.org/opensmiles.html, (accessed 14 February 2025).
- O'Boyle, N.; Dalke, A. DeepSMILES: An Adaptation of SMILES for Use in Machine-Learning of Chemical Structures, *ChemRxiv*, **2018**, DOI: 10.26434/chemrxiv.7097960.v1.
- 8 J.-N. Wu, T. Wang, Y. Chen, L.-J. Tang, H.-L. Wu and R.-Q. Yu, *Nat. Commun.*, **2024**, 15, 4993, DOI: 10.1038/s41467-024-49388-6.
- T.-S. Lin, C. W. Coley, H. Mochigase, H. K. Beech, W. Wang, Z. Wang, E. Woods, S. L. Craig, J. A. Johnson, J. A. Kalow, K. F. Jensen and B. D. Olsen, *ACS Cent. Sci.*, 2019, 5, 1523–1531, DOI: 10.1021/acscentsci.9b00476.
- 10 A. Drefahl, J. Cheminform., **2011**, 3, 1–7, DOI: 10.1186/1758-2946-3-1.
- M. Krenn, F. Häse, A. Nigam, P. Friederich and A. Aspuru-Guzik, *Mach. Learn. Sci. Technol.*, **2020**, *1*, 045024, DOI: 10.1088/2632-2153/aba947.

Discov., 2023, 2, 748–758, DOI: 10.1039/d3dd00012e.

12

This article is licensed under a Creative Commons Attribution-NonCommercial 3.0 Unported Licence

Open Access Article. Published on 13 2025. Downloaded on 14.11.25 19.25:59.

- J. Arcudia, F. Ortíz-Chi, A. Sánchez-Valenzuela, A. Aspuru-Guzik and G. Merino, Matter, 2023, 6, 1503–1513, DOI: 10.1016/j.matt.2023.02.014.
- J. Arcudia, F. Ortíz-Chi, J. Barroso and G. Merino, *Nanoscale*, 2025, 17, 2215–2223,
   DOI: 10.1039/D4NR03696D.
- C. Benecke, T. Grüner, A. Kerber, R. Laue and T. Wieland, *Fresenius J. Anal. Chem.*,
   1997, 359, 23–32, DOI: 10.1007/s002160050530.
- M. A. Yirik, M. Sorokina and C. Steinbeck, *J. Cheminform.*, 2021, 13, 1–14, DOI:
   10.1186/s13321-021-00529-9.
- 17 A. Schüller, G. Schneider and E. Byvatov, *QSAR Comb. Sci.*, **2003**, *22*, 719–721, DOI: 10.1002/qsar.200310008.
- 18 A. Schüller, V. Hähnke and G. Schneider, *QSAR Comb. Sci.*, **2007**, *26*, 407–410, DOI: 10.1002/qsar.200630101.
- D. S. Hochbaum, Approximation Algorithms for NP-hard Problems, *ACM Sigact News*,1997, 28, pp. 447–476, ISBN 978-0534949681.
- 20 R. M. Wharton, *Inf. Contr.*, **1977**, *33*, 253–272, DOI: 10.1016/S0019-9958(77)80005-3.
- 21 S. Kadioglu and M. Sellmann, *Constraints*, **2010**, *15*, 117–144, DOI: 10.1007/s10601-009-9073-4.
- S. L. Graham and M. A. Harrison, in Advances in Computers, eds. M. Rubinoff and M.
   C. Yovits, Elsevier, 1976, 14, pp. 77–185, DOI: 10.1016/S0065-2458(08)60451-9.
- 23 D.-Q. Zhang, K. Zhang and J. Cao, *Comput. J.*, **2001**, *44*, 186–200, DOI: 10.1093/comjnl/44.3.186.
- G. Landrum, P. Tosco, B. Kelley, R. Rodriguez, D. Cosgrove, R. Vianello, Sriniker, P.

Gedeck, G. Jones, NadineSchneider, E. Kawashima, D. Nealschneider, A: Dalke, M. Colour M. Wojcikowski, I. Swain, B. Cole, S. Turk, A. Savelev, tadhurst-cdd, A. Vaucher, M. Wojcikowski, I. Take, V. F. Scalfani, R. Walker, K. Ujihara, D. Probst, J. Lehtivarjo, H. Faara, G. Godin, A. Pahl and J. Monat, rdkit/rdkit: 2024\_09\_5 (Q3 2024) Release, Zenodo, 2025, DOI: 10.5281/zenodo.28640.

M. Elyashberg, A. Williams and K. Blinov, Contemporary Computer-assisted Approaches to Molecular Structure Elucidation, Royal Society of Chemistry, Cambridge, U.K, 2012, 1, pp. 28–30, ISBN 978-1-84973-432-5.

## Data availability statements

The data supporting this article have been included as part of the Supplementary Information.