

# *stk*: an extendable Python framework for automated molecular and supramolecular structure assembly and discovery

Lukas Turcani,<sup>1</sup> Andrew Tarzia,<sup>1</sup> Filip T. Szczypiński,<sup>1</sup> and Kim E. Jelfs<sup>1</sup>

Department of Chemistry, Molecular Sciences Research Hub, Imperial College London, White City Campus, Wood Lane, London, W12 0BZ, UK<sup>a)</sup>

(Dated: 8 March 2021)

Computational software workflows are emerging as all-in-one solutions to speed up the discovery of new materials. Many computational approaches require the generation of realistic structural models for property prediction and candidate screening. However, molecular and supramolecular materials represent classes of materials with many potential applications for which there is no go-to database of existing structures or general protocol for generating structures. Here, we report a new version of the supramolecular toolkit, *stk*, an open-source, extendable and modular Python framework for general structure generation of (supra)molecular structures. Our construction approach follows a bottom-up process and minimises the input required from the user, making *stk* user-friendly and applicable to many material classes. This version of *stk* includes metal-containing structures and rotaxanes as well as general implementation and interface improvements. Additionally, this version includes built-in tools for exploring chemical space with an evolutionary algorithm and tools for database generation and visualisation. The latest version of *stk* is freely available at [github.com/lukasturcani/stk](https://github.com/lukasturcani/stk).

## I. INTRODUCTION

Computational modelling seeks to accelerate functional material discovery and support experimental workflows by offering insights into chemical processes and structures that are not achievable through experiment. With advances in hardware and software, it is now possible to couple computational and experimental exploration of the vast array of potential materials and their properties at a much lower time and resource cost than experiment alone, and with a lower risk of wasted efforts.<sup>1–3</sup> Artificial intelligence (AI) and machine learning has the potential to assist in the efficient exploration of known and unexplored chemical space toward the optimal materials for a specific application.<sup>4</sup> For example, AI-driven computational workflows have been applied to explore the chemical space of transition metal complexes<sup>5–7</sup> and organic electronics.<sup>8,9</sup> Many of these approaches are facilitated by a strong push in the materials modelling community to develop open-source repositories of code, material structures and their properties.

There are established computational methods (of varying cost and accuracy) for calculating the properties of materials for many of the problems in materials science. Such approaches have facilitated the prediction of the properties of hypothetical and known materials for screening toward particular applications. However, the accurate calculation of many material properties requires a realistic and representative structural model. One solution to this problem is to use an existing database of structures from experimental results and screen them for their properties, which is a common approach in solid-state materials or biological materials.<sup>10–13</sup> Unfortunately, such an approach limits exploration beyond known examples. Therefore, the ideal solution for novel materials discovery is to generate structures from scratch. Re-

search groups often employ in-house scripts written for specific material types to generate structures for high-throughput materials screening. As such, scripts are often tailor-made for the groups' specific needs and are not made available to the broader scientific community. Furthermore, they are hard to maintain and difficult to generalise to a broader system set. In particular, the structure prediction of organic and supramolecular materials is currently difficult to generalise and limited to a small subset of possible chemical classes. There are programs currently available (some of which are open-source) for the generation of materials such as metal-organic polyhedra,<sup>14,15</sup> organic and inorganic molecules,<sup>16,17</sup> and polymeric systems.<sup>18–26</sup> However, it is not trivial to interface such codes for more general workflows or more complex projects.

To tackle the problem of general structure generation for materials discovery, we have previously reported developing the supramolecular toolkit (*stk*): an open-source and easily extendable Python library for the assembly of complex molecular architectures, including supramolecular assemblies.<sup>27</sup> Importantly, *stk* mainly generates molecular representations of materials and the problem of crystal structure prediction, or how molecules pack in the solid-state, is beyond the scope of this work. Here, we provide an update on the *stk* structure generation software, an extension to new material classes, and advanced capabilities including chemical space exploration. We have rewritten *stk* to focus on constructing a diverse range of structures with an easy-to-use interface and modular, modifiable functionality. The modular design of *stk* makes the addition of new structure classes and reactions as simple as possible for end-users. Also, *stk* allows users to deposit molecules, and molecular properties, easily into local or remote MongoDB databases. MongoDB databases constructed by *stk* can be viewed through *stk-vis*,<sup>28</sup> a standalone, cross-platform application providing 2D and 3D molecular rendering and molecular property tabulation. Molecules deposited by *stk* into MongoDB databases are immediately visible in *stk-vis*, allowing users to easily share computationally con-

<sup>a)</sup>Electronic mail: [k.jelfs@imperial.ac.uk](mailto:k.jelfs@imperial.ac.uk)

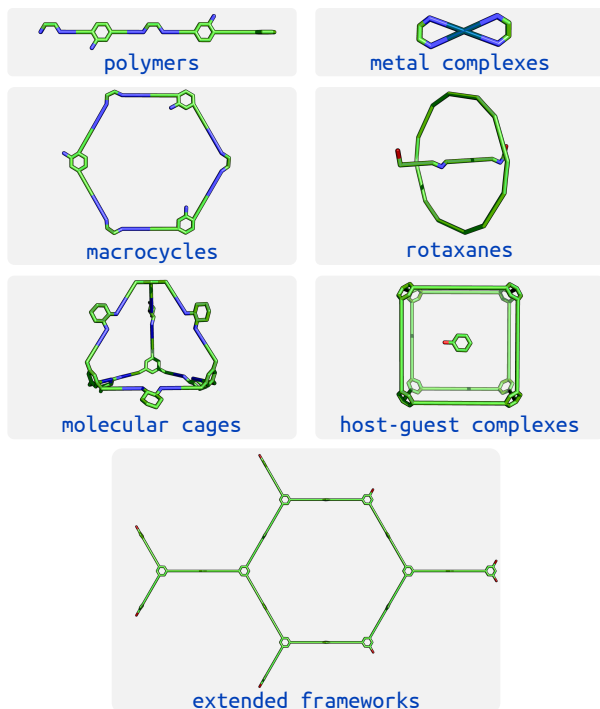


FIG. 1. Examples of buildable topology types with *stk*. Names do not match their name in the *stk* code. These structures represent the placement and alignment of building blocks on a topology graph and do not have chemically realistic bonds lengths and angles etc. between building blocks as they are not geometry optimised here.

structed molecular databases both within and across teams.

On top of the previously reported construction of covalent systems, such as linear polymers, covalent organic frameworks (COFs), and porous organic cages,<sup>27</sup> the most recent version of *stk* allows for the inclusion of metal centres, thus opening the scope to a diverse range of metal-organic cages and metal-organic frameworks (MOFs). Furthermore, *stk* now allows for the automated and custom construction of rotaxanes alongside existing supramolecular structures, such as host-guest complexes. FIG. 1 shows examples of *stk* constructed molecules from each broad topology type already implemented in *stk*. Finally, *stk* provides convenient methods for interfacing with third-party software for geometry optimisations and property calculations of *stk*-generated molecules. Geometry optimisation is beyond the scope of *stk*; for this, we have written the open-source repository *stko* that contains functions for the geometry optimisation and analysis of molecules ([github.com/JelfsMaterialsGroup/stko](https://github.com/JelfsMaterialsGroup/stko)).<sup>29</sup>

This paper describes the construction of numerous molecule types and their associated construction approaches currently implemented in *stk*. Importantly, those approaches are robust for the chemistries described here, but are also easily transferable for application to almost any chemistry type. Additionally, we describe the interfaces within *stk* for making databases of molecules and the use of an evolutionary algorithm (EA) to explore chemical space with *stk*. Further examples and more thorough documentation can be found at

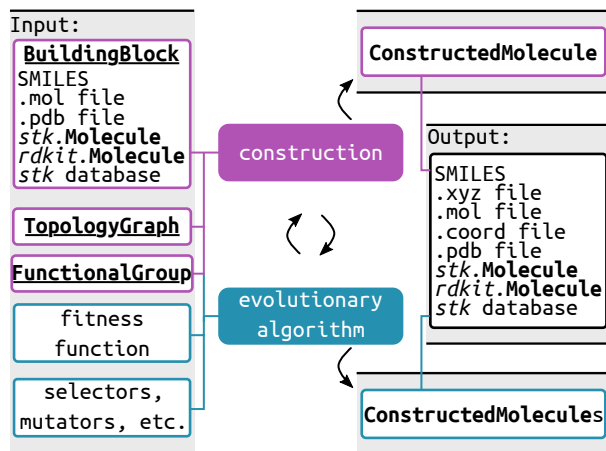


FIG. 2. The connection between the user-input, *stk* construction and evolutionary algorithms, and output. The complexity of *stk* is behind the construction and evolutionary algorithms, while the user interface is as simple as possible. We highlight the variety of input and output options, which allow for interfacing with other computational chemistry software. The construction algorithm is the crucial component behind *stk* usage.

<https://stk.readthedocs.io>; all *stk* code is freely available at [github.com/lukasturcani/stk](https://github.com/lukasturcani/stk) while *stk-vis* can be found on [github.com/lukasturcani/stk-vis](https://github.com/lukasturcani/stk-vis).

## II. SOFTWARE OVERVIEW

Here, we describe the individual software components of *stk* building up to the use of the evolutionary algorithm. *stk* is a Python library that provides users with the following capabilities: (i) automated construction of diverse molecular and supramolecular models, regardless of their complexity, (ii) automatic design of molecules with user-desired properties and (iii) creation of molecular databases. The development of *stk* focuses on a robust and straightforward user interface, with the implementation details being ultimately hidden from the user. Therefore, the above capabilities of *stk* are easily accessible with limited programming experience. For a tutorial-style introduction to using *stk*, we recommend visiting the documentation (<https://stk.readthedocs.io>). Additionally, while *stk* provides built-in functionality for each of its features, a primary design goal is that users may extend any aspect of *stk* in their code, without touching the source-code of *stk* itself. Throughout this paper, we describe the default implementation of *stk*.

The interface of *stk* handles the input, construction and output of molecules (FIG. 2). While previously *stk* provided an interface for interacting with third-party optimisation software, this functionality has since been removed, as there already exists a well-developed Python ecosystem for providing this functionality.<sup>30–32</sup> This reduction in scope allows the development of *stk* to focus on its key features. However, many of the previously implemented protocols (including Schrödinger’s MacroModel,<sup>33</sup> GULP,<sup>34,35</sup> xTB<sup>36</sup> and

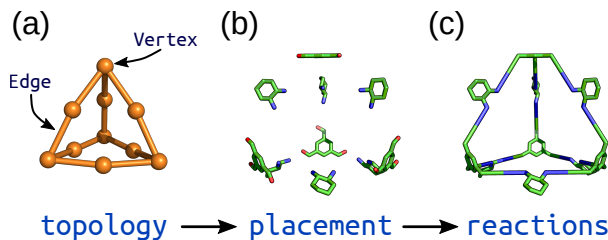


FIG. 3. Schematic of the construction process of an organic cage in *stk* starting from a (a) topology graph of vertices connected by edges. The supplied building blocks are (b) placed and aligned on the topology, and then (c) ‘reactions’ are performed between them.

RDKit<sup>37</sup>; these software packages provide access to geometry optimisations, property calculations, dynamics simulations and conformer generation algorithms) are available in our repository *stko*,<sup>29</sup> which interfaces with *stk*. Recently, we added open-source optimisation protocols to the construction process; these protocols were added into *stk* because they do not introduce significant software dependencies. Additionally, we have developed *stk-vis*,<sup>28</sup> which is a cross-platform application for the visualisation of databases created by *stk* (Section II F). *stk-vis* can also connect to remote MongoDB databases of *stk* molecules, facilitating sharing among researchers.

#### A. Construction Overview

The primary process that *stk* performs, which facilitates most of its capability, is the construction of constructed molecules from a topology graph and building blocks (represented by the **ConstructedMolecule**, **TopologyGraph** and **BuildingBlock** classes, respectively; bold text represents a class name within *stk*). The default implementation of the construction process occurs in stages: (1) **BuildingBlock** instances are placed and aligned on vertices of a **TopologyGraph**, (2) **FunctionalGroup** instances of those **BuildingBlocks** are assigned to the edges of the **TopologyGraph**, (3) ‘reactions’ are performed to connect functional groups assigned to the same edge and, as an optional final step, (4) the geometry of the structure is optimised. For example, FIG. 3 shows the placement of building blocks (separated molecules in (b)) on a topology graph to form a constructed molecule.

#### B. Topology graphs, vertices and edges

In the default implementation of *stk*, molecules are constructed by placing building blocks on topology graphs (FIG. 3(a) to (b)). Topology graphs construct molecules by first defining an underlying graph of vertices (**Vertex** instances) and edges (**Edge** instances) (FIG. 3(a)). Each vertex in a topology graph defines a position, where the building block is placed (this can be defined in Cartesian coordinates or defined relative to the vertex’s neighbours), and the series of transformations applied to the building block to align it

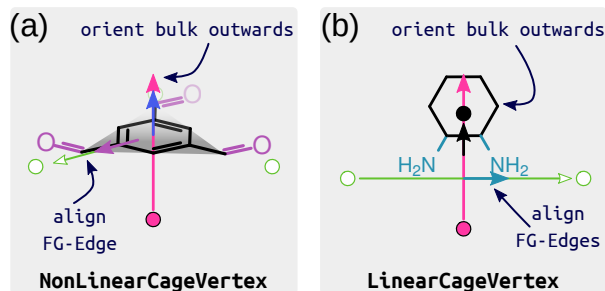


FIG. 4. Schematic of the alignment of a (a) **NonLinearCageVertex** and (b) **LinearCageVertex**, which are vertex classes specific to cage construction. These vertices use different approaches to orient the bulk (part of the chemical structure in black) of their building blocks away from the topology centre (pink circle). Nonlinear building blocks align the normal (blue arrow) to the plane defined by its functional groups (shaded triangle) with the pink arrow, while the Linear building blocks align the vector between the placer and core centroids (black arrow) with the pink arrow. Both instances use the position of the edge and functional group centroids to define the alignment of a functional group (purple in (d), cyan in (e)) with an edge (green circles).

functional groups with the neighbouring edges. Edges in a topology graph define which functional groups are joined during the reaction step of the construction process. Additionally, edges serve to provide anchors for the orientation of the building blocks, where the transformation defined by the vertex is such that functional groups are well aligned with the edges they are assigned to. The implemented alignment processes aim to align the functional groups of a building block with the edges of the vertex on which it is placed while minimising distortion and clashes in the bonds generated during construction through a series of independent transformations, such as the alignment of two vectors or the translation of a building block’s centroid to a new position (FIG. 4).

Throughout *stk*, we have implemented robust alignment processes for the built-in topology graphs that vary in complexity. In Section III, we describe the built-in molecule types that use their own topology graph, vertex and edge classes that are appropriate for their construction, where different chemical systems require different vertex transformations to achieve alignment successfully; for example the processes that orient building blocks on a cage topology graph will differ to those that orient a cycle and axle to form a rotaxane molecule. Therefore, we have made the implementation of new topology graph, vertex and edge classes straightforward such that the extension of *stk* to a user’s materials is possible. Importantly, user-defined topology graph, vertex, and edge classes can use an approach that is equivalent to one of the built-in classes or an entirely new approach; in other words, the processes used by developers of *stk* are entirely customisable.

### C. Building blocks and functional groups

The **BuildingBlock** class in *stk* represents a molecule, which is placed and aligned on the vertices of a topology graph. Building blocks will also be joined to other building blocks during the construction process, assuming an edge connects the vertices on which the building blocks are placed. The building block representation includes the atoms and bonds of a molecule, and its position matrix, which is a matrix of atomic coordinates. For placement, alignment and reactions to be carried out, building blocks also contain **FunctionalGroup** instances, which are defined by the user using the **FunctionalGroupFactory** interface (FIG. 5).

A functional group of a building block defines three sets of atoms designated for use by the default construction process: *bonders*, *deleters*, *placers*, (FIG. 5). *Bonder* and *deleter* atoms represent the atoms that will bond and be deleted, respectively, during a reaction. *Placer* atoms are used to place and align the building block. Finally, a building block also has *core* atoms, representing the bulk of the molecule, which often needs to be oriented in a specific direction. Note that *core* atoms will not get modified by reactions during the construction process as they are not part of any functional group. All building blocks have defined *core* and *placer* atom sets, while a building block can have no functional groups and, as a result, no *bonder* or *deleter* atoms. Therefore, to define *bonder* and *deleter* atoms, a functional group must be defined.

Functional groups in *stk* are defined by searching the building block for the chemical pattern that represents the desired functional group. We provide many built-in functional groups in *stk* that cover common functionalities (such as alcohols, amines, aldehydes and halogens), and the documentation covers the definition of new functional groups. Additionally, new chemical patterns can be defined very simply using SMARTS strings (a string-based representation of chemical patterns) to search for functional groups within the molecule using RDKit.<sup>37</sup> Therefore, *stk* can handle arbitrary chemical transformations of interest to the user. Importantly, this approach to defining functional groups in *stk*, through the **FunctionalGroup** classes, is as user-friendly and straightforward as possible.

### D. Reactions

Reactions are the algorithms that *stk* uses to connect building blocks during the construction process. We must emphasise that we discuss reactions on topology graphs from the perspective of the implementation within *stk*, not the associated chemical process. **Reaction** classes define algorithms that act on functional groups to either add or remove atoms or bonds. The addition of bonds between atoms of functional groups of different building blocks is what ultimately leads to building blocks being joined by the construction process. Steps two and three (FIG. 3(b) to (c)) of the default construction process define, then perform, reactions between building blocks. In the case of vertices connected by edges, the second step assigns each functional group on a **BuildingBlock** to

#### BuildingBlocks with functional groups:

```
1 bb1 = stk.BuildingBlock(
2     smiles='C1CCC(C(C1)N)N',
3     functional_groups=[stk.PrimaryAminoFactory()],
4 )
5 bb2 = stk.BuildingBlock(
6     smiles='C1=C(C=C(C=C1C=O)C=O)C=O',
7     functional_groups=[stk.AldehydeFactory()],
8 )
```

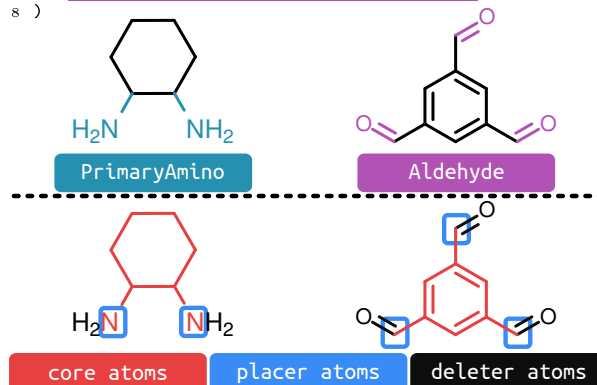


FIG. 5. Code snippet showing the generation of **BuildingBlock** instances of two molecules from their SMILES strings using built-in functional group factories. Functional group factories search the molecules for the user-requested functional groups, as an alternative to having the user specify each functional group individually. Blue and purple coloured boxes highlight the code used to request specific functional groups, shown in the same colour on the middle chemical structures. Below the dashed lines, we show the four subsets of atoms defined upon building block initialisation: *bonder* and *placer* (highlighted by blue outlines), *core* (coloured red) and *deleter* atoms (coloured black). By default, *placer* and *bonder* atoms are equivalent, which is automatically handled based on the functional groups present. However, they can also be user-defined.

an edge connected to that vertex. In step three, each edge is designated a reaction based on the set of associated functional groups, and then *stk* performs those reactions. *stk* provides three generic reactions which can be used with a wide array of functional groups: **OneOneReaction**, **OneTwoReaction** and **TwoTwoReaction**, which work on pairs of functional groups with the following combinations of *bonder* atoms: one and one, one and two, and two and two. In these cases, the algorithm is simple but generally applicable: a generic reaction between two such functional groups will delete *deleter* atoms and form bonds between *bonder* atoms on different functional groups. For a **TwoTwoReaction**, the ambiguity between which *bonders* get connected is resolved by bonding the two closest atoms and then the second nearest set.

Reactions also have the capability to create bonds with any bond order desired by the user and dative bonds for use with metal-coordinating species. As with the other parts of *stk*, the reaction process has robust defaults, but these are very customisable if need-be; *i.e.* the user can override which reactions are used to react specific pairs of functional groups when creating a **ConstructedMolecule**. By default, *stk* will select the reaction (from **OneOneReaction**, **OneTwoReaction** and **TwoTwoReaction**) that matches the number of *bonder* atoms in the functional groups. Our interface for chemical reactivity



in *stk* focuses on simplicity and generalisability. Importantly, arbitrary functional groups and reactions can be defined using the provided interface, which affords the general applicability of *stk* to user-defined problems. The default options for reactions (provided in the online documentation) are well suited to most uses of *stk*.

### E. Modular and independent construction steps

Here, we describe the default implementation of molecule construction in *stk*, which is performed at the level of vertices and edges; this means each vertex performs an independent, self-contained, operation on a single building block. Similarly, reactions between functional groups on one edge are entirely independent of those on other edges. Therefore, as long as each independent process is implemented correctly (e.g. each building block is placed and aligned correctly), the construction process is local to a single building block, or pair of building blocks for the reaction step; *i.e.* there is no inherent relationship between two vertices or edges on a topology graph. As a result of this, the construction process is trivially parallelisable, where each building block can be placed simultaneously. Additionally, a topology graph can define arbitrary vertex and edge geometries, making the structure space accessible by *stk* infinitely extendable. For example, vertices do not need to be connected by edges, which results in non-reactive topology graphs that are crucial for the study of supramolecular materials (Section III D).

The main focus of *stk* is implementing robust and general construction algorithms. To achieve this, we use the topology graph definition, which provides the placement, alignment and connection between building blocks. However, this process results in a nonphysical structure where the connectivity between building blocks is exaggerated in distance. We recognise the need for chemically reasonable structures and have implemented an interface for third-party software in our repository *stko*<sup>29</sup> and recently added two open-source geometry optimisation processes to the construction process within *stk* that are accessible as optional arguments. The newly implemented geometry optimisation protocols (part of the MCHammer package available at [github.com/andrewtarzia/MCHammer](https://github.com/andrewtarzia/MCHammer)) work to decrease the distance between building blocks on a topology graph by performing rigid translations of the building blocks after their reaction. These processes are nonphysical and, as a result, are generally applicable to any **ConstructedMolecule**. However, their lack of physical meaning suggests that they should be used with care and perhaps as the initial step in further optimisation sequences.

### F. Databases

A significant aspect of computational structure generation is developing and sharing databases of structures and their properties. As such, *stk* provides built-in support for depositing molecules and their properties into Mon-

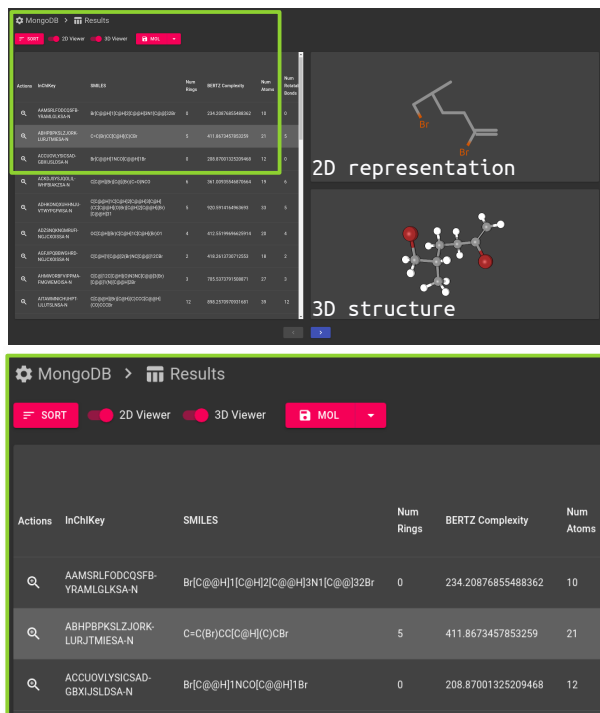


FIG. 6. Screenshot of the *stk-vis* graphical interface showing a table of molecule properties in a database and the 2D and 3D representations of the selected molecule. A zoom in on part of the database table bounded by the green box is shown at the bottom.

goDB databases; the database interface is generalisable to other database schemas. *stk* supports three database types: **MoleculeDatabase**, **ConstructedMoleculeDatabase** and **ValueDatabase**. **MoleculeDatabase** and **ConstructedMoleculeDatabase** are used for storing molecules (atoms, bonds, position matrices) and **ValueDatabase** is used for storing properties in the form of strings, numbers, lists, dictionaries and nested dictionaries thereof. Notably, the constructed molecules maintain the information about the building blocks used to construct the molecule, which are also stored in the database; these building blocks do not maintain the functional groups used.

Here, we introduce a secondary piece of software, *stk-vis*,<sup>28</sup> an open-source, cross-platform application for browsing local and remote *stk*-generated MongoDB databases. For each entry in the database, *stk-vis* provides visualisation of the properties (from a **ValueDatabase**), 2D representation and 3D structure (from the stored position matrix). Additionally, if visualising constructed molecules, *stk-vis* makes the inspection of its constituent building blocks very simple. *stk-vis* facilitates the tabulation of molecules and their properties, including sorting based on a specific property and sharing this interface through a single file transfer. FIG. 6 shows an example of the *stk-vis* graphical interface.

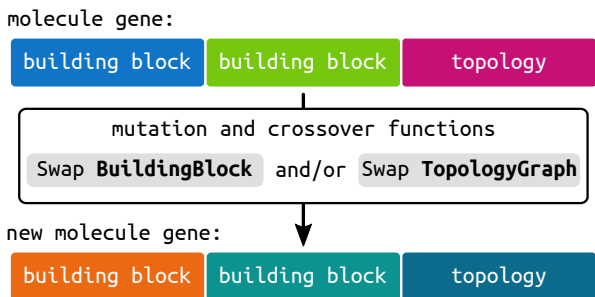


FIG. 7. Schematic of the modification of the generalised molecule gene available within *stk*, where a molecule gene is defined by the constituent building blocks and underlying topology graph, to a new molecule gene based on mutation or crossover functions. Here colours represent a change in the building block or topology graph.

### G. Evolutionary algorithm

Evolutionary algorithms (EA) are efficient approaches for exploring chemical space.<sup>38–41</sup> EAs mimic the evolutionary process by taking some population, performing mutation and crossover events on that population and then selecting survivors for the next generation. They are highly flexible algorithms, where the definition of the genome of the population, and the way the population is selected, modified and ranked are all modifiable by the user. This flexibility allows EAs to be applied in a variety of fields. The fragment-based approach treats members of the population as a collection of components that can all be modified. Such an approach is amenable to the construction approach of *stk*, where molecules are constructed from a combination of building blocks and a topology graph. We provide a general and modular implementation of the components of an EA in *stk* to automate materials discovery using the fragment-based approach.<sup>42,43</sup>

Given a population of molecules, *stk*’s EA provides a series of functions for selection, mutation, crossover, and fitness function calculation and normalisation. The EA works specifically on the **ConstructedMolecule** class and uses, by default, the building blocks and topology graph as the gene (FIG. 7). Therefore, *stk* explores mutations of that gene toward constructed molecules with the desired properties by modifying the constituent building blocks or topology graph; this process uses the *stk* construction process to generate new candidates. Finally, the EA in *stk* directly feeds its results into *stk* databases and *stk-vis* for real-time collaboration.

The entire EA process can be user-defined for a specific problem, i.e. automating the search for molecules with user-defined properties. Fitness functions, in particular, must be provided to the EA and are regular Python functions that take a **ConstructedMolecule** and return a value representing its fitness. *stk* provides multiple selection, mutation and crossover algorithms for an EA. All implemented algorithms use a fragment-based approach, where each building block is treated as a fragment of its corresponding constructed molecule. Therefore, the implemented mutation and crossover algorithms work at the building block-

level by mutating or swapping building blocks in constructed molecules. For example, the **RandomBuildingBlock** mutation will switch out the building block used in construction with a random replacement from some population of building blocks and construct a new molecule. Additionally, the user can mutate the topology graph of a constructed molecule to explore entirely distinct structures. *stk* provides a robust interface for implementing arbitrary fragment-based EAs and is continuously in development. Importantly, the documentation provides thorough examples of implementing the EA from scratch.

## III. EXAMPLES OF IMPLEMENTED SYSTEMS

In the following sections, we highlight the materials classes that *stk* can construct, including molecular materials and extended framework materials. Molecular materials are discrete and include examples of varying complexity, such as linear polymers, macrocycles, metallocycles, organic and metal-organic cages, catenanes, rotaxanes, knots and molecular machines. Extended materials are periodic and include metal-organic frameworks (MOFs) and covalent organic frameworks (COFs), which can be two- or three-dimensional. Importantly, *stk* constructs materials from a bottom-up approach, i.e. from building blocks to an assembly, and can assemble materials through covalent, coordination or noncovalent interactions. Such an approach is effectively similar to the synthetic processes used for many supramolecular materials and other molecular materials, such as cage-like molecules and crystalline frameworks. Ultimately, *stk* allows for the structure generation of molecules of arbitrary complexity.

Given a topology graph and the appropriate functional groups, *stk* can, in principle, build any structure type seen in materials chemistry, and indeed, other fields of chemistry. Crucially, the construction interface in *stk* provides the necessary control over relative building block orientation and placement on a topology graph to allow for the construction of different structural isomers of a constructed molecule. The critical distinction between material types is that their underlying topology graph will define a specific series of construction steps and vertices with specific alignment processes. Here, we describe the implemented topology graphs and corresponding molecule types, together with some examples of their use. In all cases, we show the structures directly output by *stk* without any geometry optimisation.

### A. Polymers and macrocycles

Previously we introduced the *stk* interface for constructing **Linear** polymers of any size with arbitrary repeat unit sequence and directionality.<sup>27</sup> We have added the **Macrocycle** topology graph, which allows for the construction of macrocyclic structures using a similar interface and process to the **Linear** polymer class. Both topology graphs take *repeating\_unit* and *monomer\_orientation* information as input to give the user full control over the order and orienta-

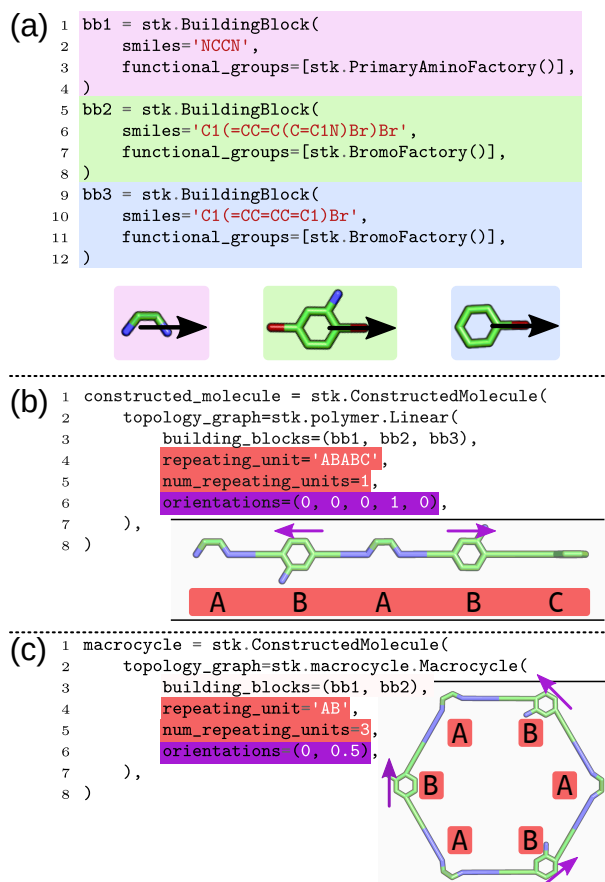


FIG. 8. (a) Definition of three building blocks in *stk*. The resulting molecules are visualised in the coloured boxes with their orientation vectors (black arrows) highlighted. Assembly of a (b) **Linear** and (c) **Macrocyclic** topology from the defined building blocks in (a). Building block ordering and repeats are highlighted in red and orientations are highlighted using purple arrows (for building blocks that are not symmetric). Coloured boxes match code snippets with their structure or effect on the structure.

tion of building blocks in the polymer or macrocycle chain (*i.e.* to control configurational isomers). In both cases, the placer atoms in the building block’s functional groups are used to align the building blocks (producing the black arrows in FIG. 8(a)). This interface allows for the user to set fractional probabilities of building block “flipping” on the topology graph (purple arrows in FIG. 8(b) and (c)). To date, the **Linear** polymer class has been used to explore very large ( $\sim 200000$  molecules) chemical spaces of organic aromatic molecules.<sup>9,44–49</sup> FIG. 8(b) and (c) shows that with a family of building blocks at hand (FIG. 8(a)), we can easily construct arbitrary **Linear** polymers and **Macrocyclic** structures.

## B. Metal-complex construction

In this release of *stk*, we have added **MetalComplex** topology graphs, which handle the placement and alignment of metal atoms and ligands on metal-complex geome-

tries. Other groups have recently implemented tools for constructing metal complexes and small molecules that encompass a broader set of possible metal-complex geometries than *stk*.<sup>5,16,17</sup> While *stk* can in principle be extended to construct any metal-complex geometry, we have focused on common geometries in supramolecular chemistry (*e.g.* variations of square planar and octahedral complexes, porphyrin and paddlewheel geometries; FIG. 9(a)). Other than the porphyrin topology, *stk* currently only handles mono- and bidentate coordination geometries. FIG. 9(b) shows a code-snippet example of the definition of a palladium(II) atom and the subsequent assembly of a bidentate square planar complex with it.

**MetalComplex** graphs have specific metal-type and ligand-type vertices, where metal-vertices are single atoms and do not undergo any orientation. All ligand orientation is based on aligning ligand binding sites (defined by **FunctionalGroup** instances) and the defined location of the edges in the **TopologyGraph**. Importantly, these topology graphs strictly define the position of **Edge** instances, which represent the ideal position of metal-coordinating atoms in the complex geometry. For bidentate ligands, the alignment process requires two idealised **Edge** positions to align the two functional groups on the ligand. This process appropriately enforces the alignment of the ligand bulk away from the metal-centre and the two ligand-binding sites inline with the two metal-binding sites (defined by the **TopologyGraph**).

To handle metal-containing systems, *stk* allows for the definition of dative bonds. Examples of defining a **Reaction** that produces dative bonds are available in the documentation (<https://stk.readthedocs.io>). Finally, our approach to metal geometries requires the strict definition of a metal-complex geometry by the user (this has been completed for the implemented examples). Therefore, distinct **TopologyGraph** classes are required to handle the following use-cases (for example): i) assembly of octahedral symmetries ( $\Lambda$  vs  $\Delta$  symmetry of tris-bidentate octahedral complexes) and ii) assembly of *cis*-protected square planar metal-complexes with free binding sites for further reaction (see FIG. 9(a)).

## C. Molecular cages

Molecular cages are a broad class of molecular systems that may have an internal cavity (*e.g.* porous molecular cages). Cages are commonly synthesised from a bottom-up building block approach and can be formed from purely organic building blocks<sup>50</sup> or building blocks that contain metal complexes.<sup>51,52</sup> They are candidate materials for solid-state and solution-phase applications in, for example, storage, separations and catalysis.<sup>50,53</sup> Of particular interest is the modularity of their design process, where specific structures or properties are targeted by choice of its constituent building blocks. Modular design processes based on constituent building blocks are well suited to *stk*. Previously, we reported functionality for constructing porous organic cages with various topology graphs,<sup>27</sup> which has since been used in the high-throughput screening of porous organic cages,<sup>42,43,54</sup> and the

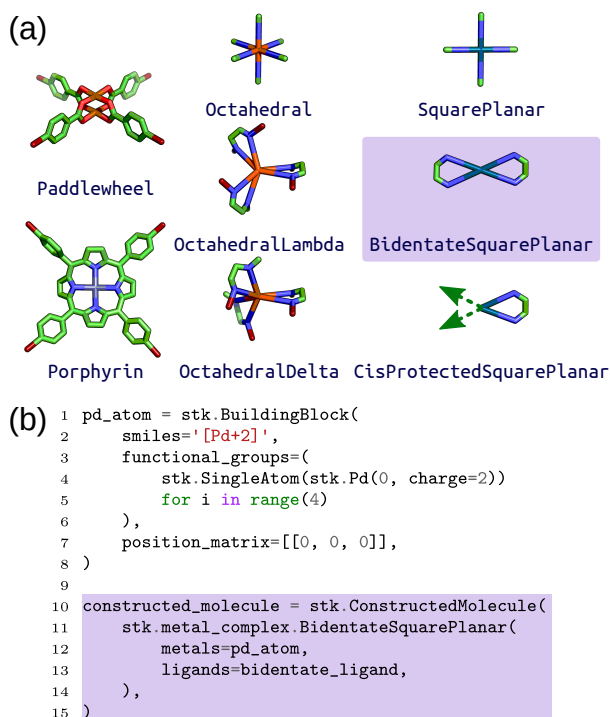


FIG. 9. (a) Implemented **MetalComplex** topology graphs. Green dashed arrows show where subsequent coordination may be performed for the **CisProtectedSquarePlanar** graph. (b) Code snippet showing the definition of a palladium (II) atom with four **SingleAtom** functional groups and the subsequent assembly of a **BidentateSquarePlanar** metal complex. Coloured boxes match up to code snippets with their structure. The code snippet in (b) is from an on-line example and is not complete because initialisation of a precursor is required.

generation of a large ( $\sim 60000$  structures) cage database for training machine learning models to predict their stability.<sup>55</sup> FIG. 10(a) shows a code snippet of cage construction using *stk*, highlighting its simplicity.

In this release, we have implemented the handling of metal-based systems toward the construction of metal-organic cages. Overall, we have implemented 31 distinct cage topologies (FIG. 10(b)) that encompass structures with diverse connectivities commonly seen in the literature. We split the topology graphs based on organic and metal-organic categories to aid the user experience and maintain domain-specific nomenclature.<sup>56</sup> However, *stk* does not make any technical distinction between them (i.e. metal-containing building blocks can be placed on an organic cage topology). The critical modification required to handle metal-containing cages is that cage topologies can now handle vertices where building blocks cannot be aligned because they are a single atom (e.g. a metal atom with **SingleAtom** functional groups). Young and co-workers recently developed *cgbind*, which is open-source software for the construction of a handful of metal-organic cage topologies.<sup>15,57</sup> In comparison, *stk*'s implementation is more general, but constructed molecules require further optimisation compared to those generated using *cgbind*. However,

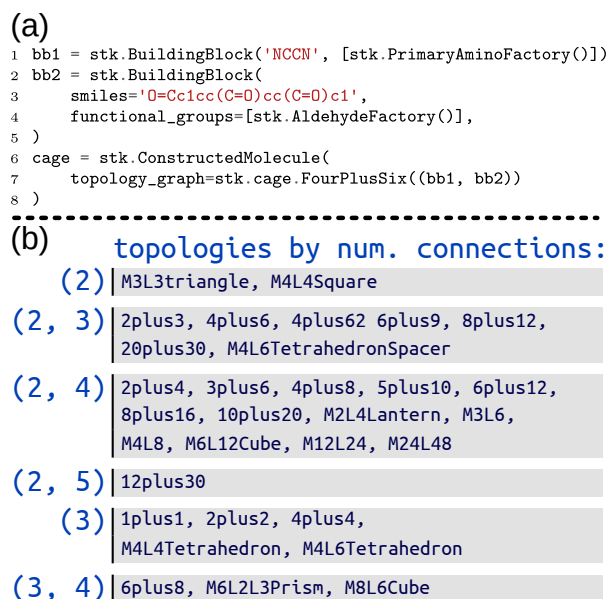


FIG. 10. (a) Code snippet showing cage construction from two building blocks. (b) Listing of all cage topology graphs built into *stk* separated by the number of connections or functional groups required by building blocks in the topology graph. Each row contains the name of topology graphs in *stk*. The column on the left shows the number of connection points of building blocks in each topology graph; i.e. (2, 3) implies that the building blocks have either two or three connection points.

we aim to overcome this using simple optimisation algorithms that result in cage structures with reasonable geometries.

Molecular cages are constructed from building blocks of different connectivity (FIG. 10(b)), including building blocks with one to five connection points. Therefore, their alignment with the topology graph's edges can be more complicated. In *stk*, we define alignment procedures based on the number of connections a building block will have in a given topology by defining the **Vertex** class (the construction approach is described in Section II E). For example, vertices with two connection points require building blocks with two functional groups and are aligned using the **LinearCageVertex** process. Similarly, vertices with three or more connections require building blocks with three or more functional groups and are aligned using the **NonLinearCageVertex** process. Additionally, we provide the **UnaligningVertex** class for building blocks that do not require alignment. The **NonLinearCageVertex** placement process happens in two steps (FIG 4): i) orientation of the building block such that the vector normal to the plane of the placer atoms aligns with the normal of the plane of edge positions, ii) orientation of the building block such that a specific structural isomer is constructed (the user can select this isomer). These steps are sufficient for avoiding collisions of building blocks and inter-building block bonds. After a building block is oriented on a vertex, functional groups are assigned to edges following a vertex-specific protocol. The assignment process specifically prevents bonds added during construction from crossing one another, which



would result in unrealistic structures. Once all building blocks are placed, and functional groups assigned to edges, reactions are performed between functional groups assigned to the same edge. So far, we have found that these processes are robust to any topology (metal-organic or organic) we implement. However, all the topology graphs we have implemented to date are concave geometries, where the bulk of the building blocks also point away from the structure’s centre.

#### D. Nonreactive topology graphs: rotaxanes and host-guest complexes

Non-reactive topology graphs construct molecules focusing on the relative spatial arrangement of the building blocks, where no bonds are created between them. Two examples discussed here are  $[n]$ rotaxanes and host-guest complexes. Rotaxanes are molecules in which a ring-shaped macrocycle is threaded on an axle with stoppers at each end of the axle.<sup>58</sup> The bulky stoppers on the axle prevent the macrocycle from slipping and, hence, the rotaxanes are mechanically interlocked, and the building blocks cannot be separated despite not being covalently bonded. The  $n$  in  $[n]$ rotaxane corresponds to the number of interconnected building blocks, so a single macrocycle on one axle would be called a  $[2]$ rotaxane. Host-guest complexes are complexes formed between a guest molecule encapsulated in the cavity of a host molecule.<sup>59</sup>

Here we describe the construction of two topology classes: **NRotaxane** and **Complex**. In both instances, no reactions are performed by *stk*, as these topology graphs focus on the relative orientation of the constituent building blocks. This differs from the synthetic process for rotaxanes, where a chemical reaction is required to form the mechanical bond holding the macrocycle on the axle. Importantly, *stk* does not attempt to model the realistic reaction processes and should be used in an alchemical way that simplifies the construction process as much as possible. We have implemented the **NRotaxane** class, which takes an axle and any number of cycles and assembles a rotaxane. FIG. 11(a) shows the formation of an **NRotaxane** from building blocks constructed using the **Linear** and **Macrocycle** classes. The macrocycles are evenly spaced along the axle, with full control over their orientation (with respect to the direction of the axle) and sequence along the axle, and are placed such that the normal of the plane of best fit of the macrocycle is parallel with the axle.

We have introduced general code to construct host-guest complexes. The method shown here is entirely generalisable to any two *stk* **BuildingBlock** instances, but focuses on the relative orientation and placement of a guest molecule to a host molecule (FIG. 11(b)). To handle the guest’s orientation relative to the host (*e.g.* to align a functional group with a specific binding site), the **Complex** can be provided with an initial and final vector, where the guest is rotated such that the initial and final vector are parallel. FIG. 11(b) shows how to use the vector of best fit through a building block’s atoms, to align the guest along a particular vector. By additionally defining the guest’s translation relative to the host’s centroid, the user can explore many host-guest conformations using *stk*.

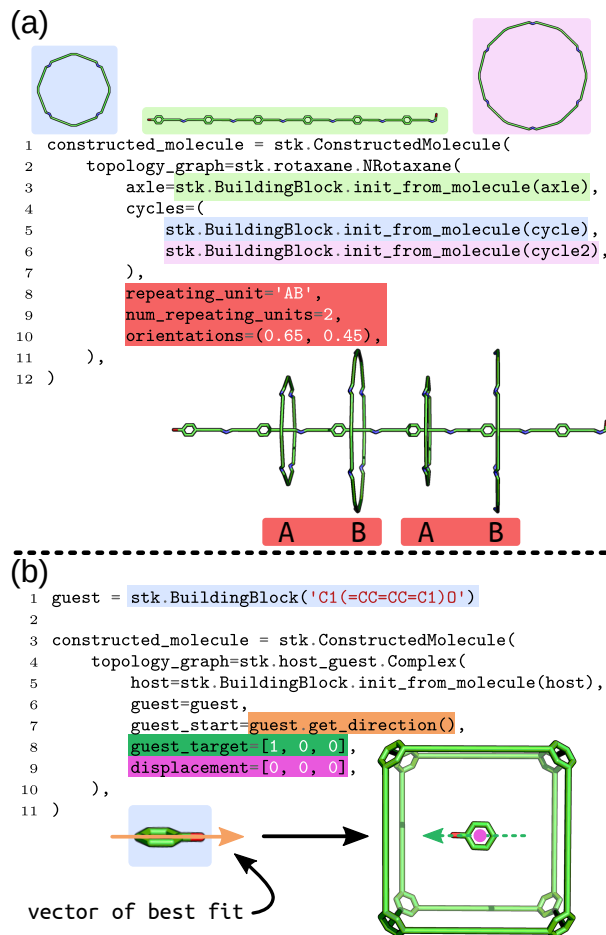


FIG. 11. (a) Assembly of an **NRotaxane** from an axle and two distinct macrocyclic building blocks. The code snippet highlights the possible modifications to the rotaxane structure by specifying the repeating unit, the number of repeating units and orientations. (b) Assembly of a **Complex** from a host and guest. The code snippet highlights the methods and vectors used for guest orientation within the host structure. Coloured boxes match up to code snippets with either their structure, impact on a structure, vector or position. Code snippets are not complete and require initialisation of precursors.

A recent example of host-guest structure generation in metal-organic cages<sup>15</sup> shows efficient ways of determining optimal guest orientation and placement, which can be automated with *stk* for any molecule class. Using low-cost simulation methods and the EA in *stk* (Section II G), we showed this on the simple case of  $C_{60}$  encapsulation in porous organic cages.<sup>43</sup> This work highlighted the benefit of fragment-based host evolution toward optimal binding that is possible within *stk*.

#### E. Constructing extended framework materials

Extended framework materials are two or three-dimensional structures that can be periodic, and hence, represented by an infinitely repeating unit cell. Covalent organic frameworks (COFs) and metal-organic frameworks

(MOFs) are two classes of extended framework materials that generally use the reticular chemistry approach<sup>60,61</sup> of constructing structures based on a given topology and systematically replacing building blocks on this topology to produce a vast potential chemical space. This process has been replicated in *stk*, and we have introduced the handling of periodicity such that *stk* can construct MOFs and COFs. Currently, we focus on 2D-COF construction, but because any **BuildingBlock** instance (with the appropriate number of functional groups) can be placed on any topology, it is possible to construct MOFs by placing metal-complexes on a COF topology.

Within *stk*, we distinguish between discrete molecules and extended materials at the topology level, where we provide a series of topology graphs defined by repeating unit cells for generating structures of extended materials. To construct an extended material, the topology graph’s unit cell can be repeated in the *x*, *y* and *z* directions to create a larger graph that the building blocks are placed on. This approach allows for the creation of infinitely repeating structures, where “periodic” bonds are created at the cell boundaries, or finite structures with unreacted groups at the cell boundaries. This distinction is required to provide an interface for generating crystal structures and “island” models of extended materials (FIG. 12(b)).

At the moment, *stk* contains four common two-dimensional topologies of COFs (hexagonal (net: **hxl**), honeycomb (two variations exist with and without a ditopic linker between three-coordinate nodes; net: **hcb**), square (net: **sql**) and kagome (net: **kgm**)) (FIG. 12(c)). However, the current implementation is extendable to three-dimensions and other extended topology graphs. For two-dimensional systems, the construction places the COF layer in the *xy* plane and orients building blocks by rotating them along the *z* direction. As with other topology graphs, unsymmetrical building blocks can be manually oriented with respect to their neighbours as desired.

#### F. Hierarchical construction and further analysis of *stk* molecules

An important feature of *stk* is the easy conversion of a **ConstructedMolecule** into a **BuildingBlock**, which allows for the simple use of a previously constructed molecule as the building block in a new construction as part of a “hierarchical” construction process. Specifically, a **BuildingBlock** can be created from an existing **BuildingBlock** object, or from a **ConstructedMolecule**, but with a different set of functional groups. Therefore, *stk* is well suited for constructing complex molecules over many steps from elementary building blocks. In particular, the **Linear** class, for example, can be used to construct combinatorial libraries of precursor molecules for further construction, e.g. to act as a family of potential rotaxane axles (FIG. 13). Additionally, we find this approach useful for assembling metallo-architectures from **MetalComplex** structures. For example, constructing an **Octahedral** metal complex, and then placing that, as a building block, on the node of a cage topology graph is much simpler than

constructing both simultaneously. By tackling the problem in a step-wise fashion, this approach simplifies the underlying topology graph (improving the code’s generalisability) and the complexity of the user input (i.e. the building blocks are simpler). A similar approach could be used to construct coordination polymers using the **MetalComplex** and **Linear** classes. Finally, all molecules generated by *stk* can be used to generate new molecules of arbitrary complexity.

In the latest version of *stk*, we also support the writing of molecules to various common file types (XYZ, MOL V3000, Protein Data Bank (PDB) and Turbomole files). The new implementation can also handle the output of periodic structures (FIG. 12(a)) for relevant file types (namely PDB and Turbomole files). Users may also straightforwardly define new functions for writing molecules to files if the built-in formats do not match their requirements. Furthermore, *stk* **Molecule** instances can be converted directly into the molecular representation of the cheminformatics software *RDKit*.<sup>37</sup> Ultimately, this allows for the interfacing of molecules generated by *stk* with many other computational chemistry software.

## CONCLUSIONS

*stk* is a Python library designed for the automated construction of structures of arbitrary complexity from their constituent building blocks. We provide a modular and open-source framework that is generally applicable and simple to extend to a user’s material of interest. Currently, constructable molecule types include linear polymers, small molecule oligomers, macrocycles, rotaxanes, metal complexes, metal-organic cages, organic cages, and extended framework materials. Importantly, we provide a robust interface to the methods required for construction: (1) functional group searching, (2) placement and alignment of building blocks on the vertices of a topology graph, and (3) reacting functional groups along edges of the topology graph; all of which can be used in any new user problem and as parts of much larger workflows. When coupled with other open-source codes in the *stk* ecosystem, including *stk-vis* and *stko*, we provide a solution to structure generation and exploration that includes a modular evolutionary algorithm and a simple interface to databasing tools (such as MongoDB) for the simplified storing and sharing of large chemical libraries. With *stk-vis*, the *stk* ecosystem is ideal for real-time collaboration between experimental and computational chemists in various materials chemistry fields. As the active developers of *stk*, we aim to provide consistent improvements and guidance for new users with example usages and tutorials; additionally, we are active in assisting new users in implementing new topology graphs or reactions. Furthermore, *stk* comes with a test-suite covering the code base, which makes extending *stk* simpler and safer. We anticipate that *stk* will provide users with a robust and general solution to structure generation and the pre- and post-processing of structures, and precursor generation.

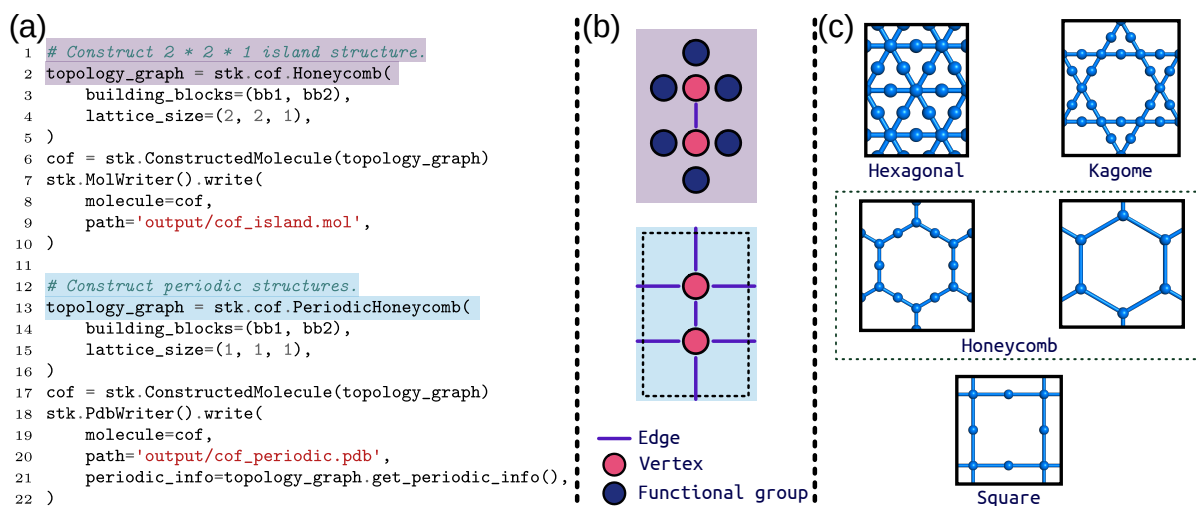


FIG. 12. (a) Assembly of a cluster (“island”) and periodic **Cof** model. The code snippet highlights the use of a “periodic” **TopologyGraph** to include periodic information, which can now be saved to PDB or Turbomole files (lines 18–22 of the code snippet). (b) Schematic showing the topological difference between the (top) nonperiodic and (bottom) periodic cases, where the main difference is the connectivity of the **FunctionalGroups** at the cell boundaries. (c) Implemented extended topology graphs showing a unit-cell containing vertices and edges. Coloured boxes match up to code snippets with their structure. Code snippets from online examples are not complete and require initialisation of precursors.

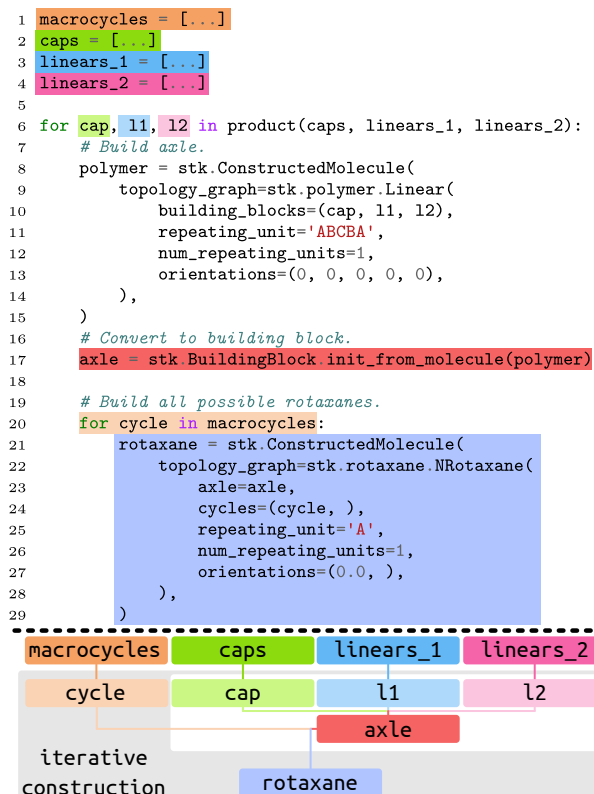


FIG. 13. Code snippet showing the hierarchical assembly of multiple rotaxane structures from a pool of macrocycles, caps and linear building blocks. Here “for” loops are used to iterate over these pools of building blocks. Coloured boxes match up to code snippets with the hierarchical process schematic below the dashed line.

## AUTHOR CONTRIBUTIONS

LT is the primary author of *stk* code. AT is the primary author of this manuscript and contributed the metal-complex and metal-organic cage code. FS contributed the rotaxane and macrocycle code. All code contributions are overseen by LT, and contributions are documented on the github history. KEJ supervised the development of *stk* and the writing of this manuscript. All authors contributed to the manuscript and oversaw its final production.

## ACKNOWLEDGMENTS

KEJ thanks the Royal Society for a University Research Fellowship and a Royal Society Enhancement Award 2018 (AT), and the ERC through Agreement Number 758370 (ERC-StG-PE5-CoMMaD). We also thank the Leverhulme Trust for a Research Project Grant (FTS). Steven Bennett is thanked for contributions to *stk* and the *stko* code. Dr. Alejandro Santana-Bonilla is thanked for assistance with implementing the xTB wrapper into the *stko* code.

## DATA AVAILABILITY STATEMENT

No data was generated for this manuscript. All code is open-source and linked to throughout the manuscript.

<sup>1</sup>R. L. Greenaway and K. E. Jelfs, “Integrating computational and experimental workflows for accelerated organic materials discovery,” *Adv. Mater.*, 2004831 (2021).

- <sup>2</sup>D. Ongari, L. Talirz, and B. Smit, "Too many materials and too many applications: An experimental problem waiting for a computational solution," *ACS Cent. Sci.* **6**, 1890–1900 (2020).
- <sup>3</sup>R. Pollice, G. dos Passos Gomes, M. Aldeghi, R. J. Hickman, M. Krenn, C. Lavigne, M. Lindner-D'Addario, A. Nigam, C. T. Ser, Z. Yao, and A. Aspuru-Guzik, "Data-driven strategies for accelerated materials design," *Acc. Chem. Res.* **54**, 849–860 (2021).
- <sup>4</sup>K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev, and A. Walsh, "Machine learning for molecular and materials science," *Nature* **559**, 547–555 (2018).
- <sup>5</sup>A. Nandy, C. Duan, J. P. Janet, S. Gugler, and H. J. Kulik, "Strategies and software for machine learning accelerated discovery in transition metal chemistry," *Ind. Eng. Chem. Res.* **57**, 13973–13986 (2018).
- <sup>6</sup>J. P. Janet, S. Ramesh, C. Duan, and H. J. Kulik, "Accurate multiobjective design in a space of millions of transition metal complexes with neural-network-driven efficient global optimization," *ACS Cent. Sci.* **6**, 513–524 (2020).
- <sup>7</sup>J. P. Janet, C. Duan, A. Nandy, F. Liu, and H. J. Kulik, "Navigating transition-metal chemical space: Artificial intelligence for first-principles design," *Acc. Chem. Res.* **54**, 532–545 (2021).
- <sup>8</sup>R. Gómez-Bombarelli, J. Aguilera-Iparraguirre, T. D. Hirzel, D. Duvenaud, D. Maclaurin, M. A. Blood-Forsythe, H. S. Chae, M. Einzinger, D.-G. Ha, T. Wu, G. Markopoulos, S. Jeon, H. Kang, H. Miyazaki, M. Numata, S. Kim, W. Huang, S. I. Hong, M. Baldo, R. P. Adams, and A. Aspuru-Guzik, "Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach," *Nat. Mater.* **15**, 1120–1127 (2016).
- <sup>9</sup>L. Wilbraham, R. S. Sprick, K. E. Jelfs, and M. A. Zwijnenburg, "Mapping binary copolymer property space with neural networks," *Chem. Sci.* **10**, 4973–4984 (2019).
- <sup>10</sup>M. Miklitz, S. Jiang, R. Clowes, M. E. Briggs, A. I. Cooper, and K. E. Jelfs, "Computational screening of porous organic molecules for Xenon/Krypton separation," *J. Phys. Chem. C* **121**, 15211–15222 (2017).
- <sup>11</sup>M. Tong, Y. Lan, Z. Qin, and C. Zhong, "Computation-ready, experimental covalent organic framework for methane delivery: Screening and material design," *J. Phys. Chem. C* **122**, 13009–13016 (2018).
- <sup>12</sup>D. Ongari, A. V. Yakutovich, L. Talirz, and B. Smit, "Building a consistent and reproducible database for adsorption evaluation in Covalent–Organic frameworks," *ACS Cent Sci* **5**, 1663–1675 (2019).
- <sup>13</sup>Z. Yao, B. Sánchez-Lengeling, N. S. Bobbitt, B. J. Bucior, S. G. H. Kumar, S. P. Collins, T. Burns, T. K. Woo, O. K. Farha, R. Q. Snurr, and A. Aspuru-Guzik, "Inverse design of nanoporous crystalline reticular materials with deep generative models," *Nat. Mach. Intell.* **3**, 76–86 (2021).
- <sup>14</sup>B. P. Hay and T. K. Firman, "HostDesigner: a program for the de novo structure-based design of molecular receptors with binding sites that complement metal ion guests," *Inorg. Chem.* **41**, 5502–5512 (2002).
- <sup>15</sup>T. A. Young, R. Gheorghe, and F. Duarte, "Cgbind: A python module and web app for automated metallocage construction and Host–Guest characterization," *J. Chem. Inf. Model.* **60**, 3546–3557 (2020).
- <sup>16</sup>E. I. Ioannidis, T. Z. H. Gani, and H. J. Kulik, "molSimplify: A toolkit for automating discovery in inorganic chemistry," *J. Comput. Chem.* **37**, 2106–2117 (2016).
- <sup>17</sup>J.-G. Sobez and M. Reiher, "Molassembler: Molecular graph construction, modification, and conformer generation for inorganic and organic molecules," *J. Chem. Inf. Model.* **60**, 3884–3900 (2020).
- <sup>18</sup>L. J. Abbott, K. E. Hart, and C. M. Colina, "Polymatic: A generalized simulated polymerization algorithm for amorphous polymers," *Theor. Chem. Acc.* **132**, 1334 (2013).
- <sup>19</sup>C. E. Wilmer, M. Leaf, C. Y. Lee, O. K. Farha, B. G. Hauser, J. T. Hupp, and R. Q. Snurr, "Large-scale screening of hypothetical metal–organic frameworks," *Nat. Chem.* **4**, 83–89 (2012).
- <sup>20</sup>P. G. Boyd and T. K. Woo, "A generalized method for constructing hypothetical nanoporous materials of any net topology from graph theory," *CrystEngComm* **18**, 3777–3792 (2016).
- <sup>21</sup>Y. J. Colón, D. A. Gómez-Gualdrón, and R. Q. Snurr, "Topologically guided, automated construction of Metal–Organic frameworks and their evaluation for energy-related applications," *Cryst. Growth Des.* **17**, 5801–5810 (2017).
- <sup>22</sup>R. L. Martin and M. Haranczyk, "Construction and characterization of structure models of crystalline porous polymers," *Cryst. Growth Des.* **14**, 2431–2440 (2014).
- <sup>23</sup>J. Keupp and R. Schmid, "TopoFF: MOF structure prediction using specifically optimized blueprints," *Faraday Discuss.* **211**, 79–101 (2018).
- <sup>24</sup>M. A. Addicoat, D. E. Coupry, and T. Heine, "AuToGraFS: Automatic topological generator for framework structures," *J. Phys. Chem. A* **118**, 9607–9614 (2014).
- <sup>25</sup>J. P. Darby, M. Arhangel'skis, A. D. Katsenis, J. M. Marrett, T. Friščić, and A. J. Morris, "Ab initio prediction of metal-organic framework structures," *Chem. Mater.* **32**, 5835–5844 (2020).
- <sup>26</sup>P. G. Boyd, Y. Lee, and B. Smit, "Computational development of the nanoporous materials genome," *Nat. Rev. Mater.* **2**, No. 17037 (2017).
- <sup>27</sup>L. Turcani, E. Berardo, and K. E. Jelfs, "Stk: A python toolkit for supramolecular assembly," *J. Comput. Chem.* **39**, 1931–1942 (2018).
- <sup>28</sup>L. Turcani, "Stk-vis," <https://github.com/lukasturcani/stk-vis>, (accessed February 18, 2021).
- <sup>29</sup>S. Bennett, A. Tarzia, and L. Turcani, "Stko:stk-optimizers," <https://github.com/JelfsMaterialsGroup/stko>, (accessed February 18, 2021).
- <sup>30</sup>A. H. Larsen, J. J. Mortensen, J. Blomqvist, I. E. Castelli, R. Christensen, Marcin Dulak, J. Friis, M. N. Groves, B. Hammer, C. Hargus, E. D. Hermes, P. C. Jennings, P. B. Jensen, J. Kermode, J. R. Kitchin, E. L. Kolsbjerg, J. Kubal, Kristen Kaasbjerg, S. Lysgaard, J. B. Maronsson, T. Maxson, T. Olsen, L. Pastewka, Andrew Peterson, C. Rostgaard, J. Schiøtz, O. Schütt, M. Strange, K. S. Thygesen, Tejs Vegge, L. Vilhelmsen, M. Walter, Z. Zeng, and K. W. Jacobsen, "The atomic simulation environment—a python library for working with atoms," *J. Phys. Condens. Matter* **29**, 273002 (2017).
- <sup>31</sup>P. Eastman, J. Swails, J. D. Chodera, R. T. McGibbon, Y. Zhao, K. A. Beauchamp, L.-P. Wang, A. C. Simmonett, M. P. Harrigan, C. D. Stern, R. P. Wiewiora, B. R. Brooks, and V. S. Pande, "OpenMM 7: Rapid development of high performance algorithms for molecular dynamics," *PLOS Comput. Biol.* **13**, e1005659 (2017).
- <sup>32</sup>S. P. Ong, W. D. Richards, A. Jain, G. Hautier, M. Kocher, S. Cholia, D. Gunter, V. L. Chevrier, K. A. Persson, and G. Ceder, "Python materials genomics (pymatgen): A robust, open-source python library for materials analysis," *Comput. Mater. Sci.* **68**, 314–319 (2013).
- <sup>33</sup>Schrödinger, "Schrödinger release 2021-1: MacroModel," Schrödinger (2021).
- <sup>34</sup>J. D. Gale, "GULP: A computer program for the symmetry-adapted simulation of solids," *J. Chem. Soc. Faraday Trans.* **93**, 629–637 (1997).
- <sup>35</sup>J. D. Gale and A. L. Rohl, "The general utility lattice program (GULP)," *Mol. Simul.* **29**, 291–341 (2003).
- <sup>36</sup>C. Bannwarth, S. Ehlert, and S. Grimme, "GFN2-xTB—An accurate and broadly parametrized self-consistent tight-binding quantum chemical method with multipole electrostatics and density-dependent dispersion contributions," *J. Chem. Theory Comput.* **15**, 1652–1671 (2019).
- <sup>37</sup>G. A. Landrum, "RDKit: Open-source cheminformatics," <http://www.rdkit.org/>, (accessed March 1, 2020).
- <sup>38</sup>N. Brown, B. McKay, F. Gilardoni, and J. Gasteiger, "A graph-based genetic algorithm and its application to the multiobjective evolution of median molecules," *J. Chem. Inf. Comput. Sci.* **44**, 1079–1087 (2004).
- <sup>39</sup>N. M. O'Boyle, C. M. Campbell, and G. R. Hutchison, "Computational design and selection of optimal organic photovoltaic materials," *J. Phys. Chem. C* **115**, 16200–16210 (2011).
- <sup>40</sup>I. Y. Kanak and G. R. Hutchison, "Rapid computational optimization of molecular properties using genetic algorithms: Searching across millions of compounds for organic photovoltaic materials," arXiv:1707.02949 [physics] (2017), [arXiv:1707.02949 \[physics\]](https://arxiv.org/abs/1707.02949).
- <sup>41</sup>J. H. Jensen, "A graph-based genetic algorithm and generative model/Monte carlo tree search for the exploration of chemical space," *Chem. Sci.* **10**, 3567–3572 (2019).
- <sup>42</sup>E. Berardo, L. Turcani, M. Miklitz, and K. E. Jelfs, "An evolutionary algorithm for the discovery of porous organic cages," *Chem. Sci.* **9**, 8513–8527 (2018).
- <sup>43</sup>M. Miklitz, L. Turcani, R. L. Greenaway, and K. E. Jelfs, "Computational discovery of molecular c 60 encapsulants with an evolutionary algorithm," *Commun. Chem.* **3**, 1–10 (2020).
- <sup>44</sup>L. Wilbraham, E. Berardo, L. Turcani, K. E. Jelfs, and M. A. Zwijnenburg, "High-throughput screening approach for the optoelectronic properties of conjugated polymers," *J. Chem. Inf. Model.* **58**, 2450–2459 (2018).



- <sup>45</sup>R. S. Sprick, C. M. Aitchison, E. Berardo, L. Turcani, L. Wilbraham, B. M. Alston, K. E. Jelfs, M. A. Zwijnenburg, and A. I. Cooper, "Maximising the hydrogen evolution activity in organic photocatalysts by copolymerisation," *J. Mater. Chem. A* **6**, 11994–12003 (2018).
- <sup>46</sup>I. Heath-Apostolopoulos, L. Wilbraham, and M. A. Zwijnenburg, "Computational high-throughput screening of polymeric photocatalysts: Exploring the effect of composition, sequence isomerism and conformational degrees of freedom," *Faraday Discuss.* **215**, 98–110 (2019).
- <sup>47</sup>Y. Bai, L. Wilbraham, B. J. Slater, M. A. Zwijnenburg, R. S. Sprick, and A. I. Cooper, "Accelerated discovery of organic polymer photocatalysts for hydrogen evolution from water through the integration of experiment and theory," *J. Am. Chem. Soc.* **141**, 9063–9071 (2019).
- <sup>48</sup>C. B. Meier, R. Clowes, E. Berardo, K. E. Jelfs, M. A. Zwijnenburg, R. S. Sprick, and A. I. Cooper, "Structurally diverse covalent triazine-based framework materials for photocatalytic hydrogen evolution from water," *Chem. Mater.* **31**, 8830–8838 (2019).
- <sup>49</sup>I. Heath-Apostolopoulos, D. Vargas-Ortiz, L. Wilbraham, K. E. Jelfs, and M. Zwijnenburg, "Using high-throughput virtual screening to explore the optoelectronic property space of organic dyes; finding diketopyrrolopyrrole dyes for dye-sensitized water splitting and solar cells," *Sustain. Energ. Fuels* (2020), 10.1039/D0SE00985G.
- <sup>50</sup>T. Hasell and A. I. Cooper, "Porous organic cages: Soluble, modular and molecular pores," *Nat. Rev. Mater.* **1** (2016), 10.1038/natrevmats.2016.53.
- <sup>51</sup>T. R. Cook and P. J. Stang, "Recent developments in the preparation and chemistry of metallacycles and metallacages via coordination," *Chem. Rev.* **115**, 7001–7045 (2015).
- <sup>52</sup>B. S. Pilgrim and N. R. Champness, "Metal-organic frameworks and metal-organic cages – a perspective," *ChemPlusChem* **85**, 1842–1856 (2020).
- <sup>53</sup>H. Vardhan, M. Yusubov, and F. Verpoort, "Self-assembled metal-organic polyhedra: An overview of various applications," *Coord. Chem. Rev.* **306**, 171–194 (2016).
- <sup>54</sup>E. Berardo, R. L. Greenaway, L. Turcani, B. M. Alston, M. J. Bennison, M. Miklitz, R. Clowes, M. E. Briggs, A. I. Cooper, and K. E. Jelfs, "Computationally-inspired discovery of an unsymmetrical porous organic cage," *Nanoscale* **10**, 22381–22388 (2018).
- <sup>55</sup>L. Turcani, R. L. Greenaway, and K. E. Jelfs, "Machine learning for organic cage property prediction," *Chem. Mater.* **31**, 714–727 (2019).
- <sup>56</sup>V. Santolini, M. Miklitz, E. Berardo, and K. E. Jelfs, "Topological landscapes of porous organic cages," *Nanoscale* **9**, 5280–5298 (2017).
- <sup>57</sup>T. A. Young, V. Martí-Centelles, J. Wang, P. J. Lusby, and F. Duarte, "Rationalizing the activity of an "artificial diels-alderase": Establishing efficient and accurate protocols for calculating supramolecular catalysis," *J. Am. Chem. Soc.* **142**, 1300–1310 (2020).
- <sup>58</sup>M. Xue, Y. Yang, X. Chi, X. Yan, and F. Huang, "Development of pseudorotaxanes and rotaxanes: From synthesis to stimuli-responsive motions to applications," *Chem. Rev.* **115**, 7398–7501 (2015).
- <sup>59</sup>S. Zarra, D. M. Wood, D. A. Roberts, and J. R. Nitschke, "Molecular containers in complex chemical systems," *Chem. Soc. Rev.* **44**, 419–432 (2015).
- <sup>60</sup>O. M. Yaghi, M. O'Keeffe, N. W. Ockwig, H. K. Chae, M. Eddaoudi, and J. Kim, "Reticular synthesis and the design of new materials," *Nature* **423**, 705–714 (2003).
- <sup>61</sup>H. Furukawa, K. E. Cordova, M. O'Keeffe, and O. M. Yaghi, "The chemistry and applications of Metal–Organic Frameworks," *Science* **341**, 1230444 (2013).
- <sup>62</sup>J. D. Crowley, S. M. Goldup, A.-L. Lee, D. A. Leigh, and R. T. McBurney, "Active metal template synthesis of rotaxanes, catenanes and molecular shuttles," *Chem. Soc. Rev.* **38**, 1530–1541 (2009).
- <sup>63</sup>P. Z. Moghadam, T. Islamoglu, S. Goswami, J. Exley, M. Fantham, C. F. Kaminski, R. Q. Snurr, O. K. Farha, and D. Fairen-Jimenez, "Computer-aided discovery of a metal–organic framework with superior oxygen uptake," *Nat. Commun.* **9**, 1378 (2018).