

This document is confidential and is proprietary to the American Chemical Society and its authors. Do not copy or disclose without written permission. If you have received this item in error, notify the sender and delete all copies.

Modular Software for Generating and Modelling Diverse Polymer Databases

Journal:	<i>Journal of Chemical Information and Modeling</i>
Manuscript ID	ci-2023-00081z
Manuscript Type:	Article
Date Submitted by the Author:	19-Jan-2023
Complete List of Authors:	Santana-Bonilla, Alejandro; King's College London, Department of Physics Lopez-Rios de Castro, Raquel; King's College London, Department of Chemistry Sun, Peike; King's College London, Department of Physics Ziolek, Robert; King's College London, Department of Physics Lorenz, Christian; King's College London, Department of Physics

SCHOLARONE™
Manuscripts

Modular Software for Generating and Modelling Diverse Polymer Databases

Alejandro Santana-Bonilla,^{*,†} Raquel López-Ríos de Castro,^{‡,¶} Peike Sun,[¶]
Robert M. Ziolek,[¶] and Christian D. Lorenz^{*,¶}

[†]*Department of Physics, King's College London, London, WC2R 2LS, United Kingdom*

[‡]*Department of Chemistry, King's College London, London, SE1 1DB, United Kingdom*

[¶]*Biological Physics and Soft Matter Group, Department of Physics, King's College
London, London, WC2R 2LS, United Kingdom*

E-mail: alejandro.santana_bonilla@kcl.ac.uk; chris.lorenz@kcl.ac.uk

Abstract

Machine learning methods offer the opportunity to design new functional materials on an unprecedented scale however building the large, diverse databases of molecules on which to train such methods remains a daunting task. **Automated computational chemistry modelling workflows** are therefore becoming **essential tools in this data-driven hunt for new materials with novel properties**, since they offer a workflow by which to create and curate molecular databases without requiring significant levels of user input. This ensures well-founded concerns regarding data provenance, reproducibility and replicability are mitigated. We have developed a versatile and flexible software package, **PySoftK** (**Python Soft Matter at King's College London**), that provides flexible, automated computational workflows to create, model, and curate libraries of polymers with a minimal user intervention. **PySoftK is available as an efficient, fully-tested, and easily installed Python package.** Key features of the software include the wide

range of different polymer topologies that can be automatically generated and fully parallelized library generation tools. It is anticipated that PySoftK will support the **generation, modelling and curation of large polymer libraries** to support functional materials discovery in nano- and bio-technology.

Introduction

The diverse **chemical, mechanical and electronic properties** of polymers underlie their application in relevant technological areas **spanning structural materials, cosmetics, pharmaceutical formulation, electronics, and biotechnology**.^{1–10} In order to optimize their aforementioned properties for this range of applications, the role of polymer constitution and topology has been widely explored, with a range of architectures including **homopolymers, block copolymers, branched and ring polymers**.^{11–17} As a result of contemporary research activities, the chemical and structural domains of synthetic polymers are continuously growing.^{18–20}

Polymeric materials have been the focus of a significant amount of scientific research for many decades. In more recent years, access to ever-growing computational power has allowed the materials modelling community to carry out increasingly complex investigations of polymeric materials.^{21–28} Modern computer simulation studies provide predictive understanding of molecular interactions and mechanisms that determine the bulk-scale properties of polymers of interest. In tandem with experimental data, this combined insight yields rational design principles for new polymers with enhanced target properties.

While there are well-established computer simulation techniques to support the investigation of polymers, often the most technically difficult part of a polymer simulation workflow is actually generating models of the polymer(s) of interest. As a result, multiple computational platforms have been developed in recent years with the aim of making the generation of polymer models more straightforward.^{29–33} Several of these packages have been developed to allow users to build molecular models of polymers for use in molecular dynamics simulations. These packages generally require the user to provide a detailed description of the underlying

chemistry of, and connectivity between, constituent monomers to output a topology of the desired polymer. Some allow the user to input forcefield information to be used in generating input files for classical molecular dynamics simulations. All of the referenced codes allow for the user to build homopolymers while some allow the user to build heteropolymers with different architectures and topologies.^{29–32}

More recently, Polymer Structure Predictor (PSP)³⁴ has reduced the amount of information required from the user to describe the chemistry of the monomers by allowing a SMILES string to be used as the input. The code can assign forcefield parameters from the CHARMM, AMBER and OPLS generalized forcefields^{35–37} and integrates with the LAMMPS classical simulation engine³⁸ and the pysimm package²⁹ for optimization of polymer structures. While the amount of the predefined information required to build the polymers is less in PSP than is required by other packages, the polymers that can currently be built are limited to homopolymers with either linear or cyclic topologies.

Here, we present PySoftK, a modular and versatile code to model polymer structures with different topologies. We present the various modules that currently exist within PySoftK to build different polymer structures, and show how they can be combined in order to build highly complex polymer topologies in an automated way. We subsequently demonstrate the various functionalities that are found within the code to facilitate high-throughput molecular modelling calculations. PySoftK can also be used in the parameterization of dihedral terms in conjugated polymers, which are commonly poorly defined by standard classical forcefields.^{28,39} Finally, we briefly review the steps that have been taken when developing the code to ensure it can be successfully used in a broad range of applications.

Software overview

PySoftK is a modular Python package that generates molecular models of polymers with a diverse range of topologies and chemistries with minimal input from the user. It also

1
2
3 has various tools to facilitate high-throughput simulations of polymers. The code utilises
4 RDKit to generate the molecular models of the various polymers.⁴⁰ The modular nature of
5 this package will allow the scientific community to easily expand the code base to include
6 polymer architectures that are not included in this first release of the library.
7
8
9
10
11
12

13 Generating diverse polymer architectures

14

15
16 **Inputting monomers.** A key functionality of PySoftK is the automated generation of
17 a diverse range of polymer architectures. The user-inputted monomers, in which the sites
18 that link each monomer together are predefined, are the main topological descriptor used
19 by PySoftK. The chemical structure and connectivity of the monomer(s) that make up the
20 polymer can be inputted using all formats supported by RDKit, such as **SMILES strings**
21 and **.pdb**, **.mol** and **.xyz files**, as shown in Listing 1. Examples of inputs to input a furan
22 monomer are shown in Figure 1a. In order to construct a polymer from its constituent
23 monomers, a placeholder atom is used to indicate where the bond formation takes place
24 during polymerization. In the examples presented in this manuscript, bromine (Br) atoms
25 or platinum (Pt) are used as these placeholder atoms (see Fig. 1b). It is important to note
26 that this choice is entirely arbitrary: any atom can be used as the placeholder atom.
27
28
29
30
31
32
33
34
35
36
37

```
38 from rdkit import Chem
39 from rdkit.Chem import AllChem
40 # SMILES, MOL and PDB formats
41 mol_1=Chem.MolFromSmiles('c1(ccc(o1)Br)Br')
42 mol_2=Chem.MolFromPDBFile('furan_pysoftk.pdb')
43 mol_3=Chem.MolFromMolFile('furan_pysoftk.mol')
```

44
45
46
47
48 Listing 1: Code snippet showing the different input formats that PySoftK accepts.
49
50

51 **Generating polymers with different monomer distributions.** Once the monomers
52 are defined, then PySoftK can be used to describe the distribution of the monomers within
53 the polymer and the overall architecture of the polymer. The code is structured such that
54
55
56
57
58

59
60

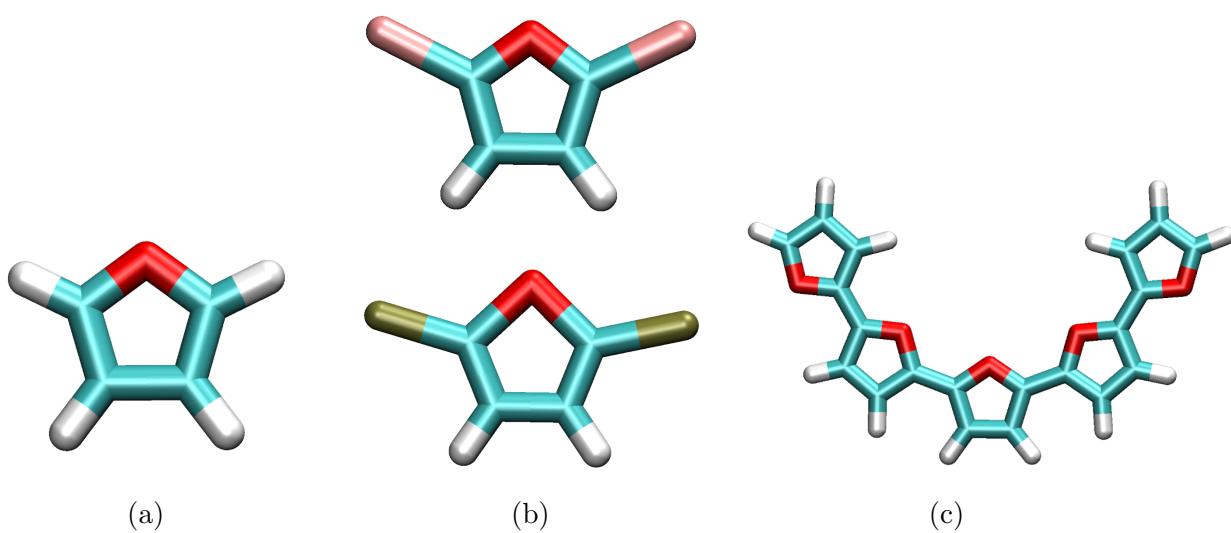


Figure 1: Creation of a linear polymer. (a) Undecorated single furan molecule. (b) Decorated furan molecule (using a Br or Pt atom) highlighting the region where intermolecular bonds are formed between monomers. (c) Linear polymer with a 5 monomers and optimized atomic positions at the RDKit **MMFF** force-field level.⁴⁰

specific modules are provided to generate homopolymers, diblock copolymers and random polymers. Block and random copolymers with more than two types of monomers, as well as alternating, periodic and statistical copolymers of any number of monomers, can be built using the patterned module within PySoftK. The modular structure of the software allows users to develop their own codes to design additional polymer topologies. In the following sections, examples of the different modules to build the various polymer topologies with different distributions of the monomers are provided.

Generating polymers with different topologies. Users can generate models of polymers with three different general topologies: (i) linear, (ii) cyclic and (iii) branched. In the following sections, examples showing how to use the various modules to build each block and random copolymers with each topology will be presented. Finally, we will present the patterned module, which provides the user more flexibility in defining the distribution of the monomers within the polymer and can be used in combination with the various topology modules to build any of these different polymer architectures.

Linear polymers. Linear polymers consist of individual monomers which are joined

1
2
3 together end-to-end forming a single molecule. In order to build a linear homopolymer in
4 PySoftK, the command `Lp(mol,atom,n_copies,shift).linear_polymer(FF,iter_ff)` is used to
5 initiate a two-step process. First, a monomer (defined by the variable `mol`) is copied and
6 translated along a predefined axis, where the distance between each monomer is defined by
7 the `shift` variable in the module call, which results in an unconnected chain of `n_copies`
8 monomers. Subsequently, a merging step between end-to-end monomers is performed em-
9 ploying a user-designated atomic placeholder (defined by the variable `atom`) indicating to
10 PySoftK the selected site on the monomer where a bond is formed. Finally, to provide a
11 realistic initial structure, PySoftK utilizes the MMFF or UFF force-field parameterizations
12 (as chosen with the `FF` parameter) as implemented in RDKit⁴⁰ in order to generate an en-
13 ergy minimized molecular conformation after `iter_ff` steps of minimisation. This process is
14 displayed in Figure 1, where a single furan molecule is decorated with a placeholder atom
15 (in this case bromine) shown in pink and extended to form a linear polymer containing five
16 monomers. The full block of code required to generate this homopolymer is found in Listing
17 2.
18
19

20 The `Db` module within PySoftK allows users to build linear diblock copolymers. The
21 command `Db(ma,mb,atom).diblock_copolymer(len_block_A,len_block_B,FF,iter_ff)` is used to
22 generate a diblock copolymer that consists of a block containing `len_block_A` monomers of `ma`
23 and a block containing `len_block_B` monomers of `mb`. Again the monomers are inputted with
24 atomic placeholders (`atom`) that identify the polymerisation site on each monomer, and the
25 resultant structure undergoes `iter_ff` steps of energy minimisation using the `FF` forcefield.
26 Listing 2 shows an example of how to practically apply this function.

```
47  
48 1 from pysftk.linear_polymer.linear_polymer import *  
49 2 from pysftk.topologies.diblock import *  
50 3  
51 4 # Build linear homopolymer  
52 5 # Input monomer with SMILES format  
53 6 mol_1=Chem.MolFromSmiles('c1(ccc(o1)Br)Br')
```

```

7 AllChem.EmbedMolecule(mol_1)
8 new=Lp(mol_1,"Br",5,shift=1.25).linear_polymer("MMFF",1)
9
10 # Build linear diblock copolymer
11 # Input ethylene oxide monomer with PDB format
12 mol_1 = Chem.MolFromPDBFile('eo.pdb')
13 # Input propylene oxide monomer with PDB format
14 mol_2 = Chem.MolFromPDBFile('po.pdb')
15 # Build a polymer with 5 eo monomers & 7 po monomers which is minimised
16 # using the MMFF forcefield over 10 iterations
17 poly = Db(mol_1,mol_2,"Br").diblock_copolymer(5,7,'MMFF',10)

```

Listing 2: Code snippet showing the creation of a linear homopolymer and a linear diblock copolymer.

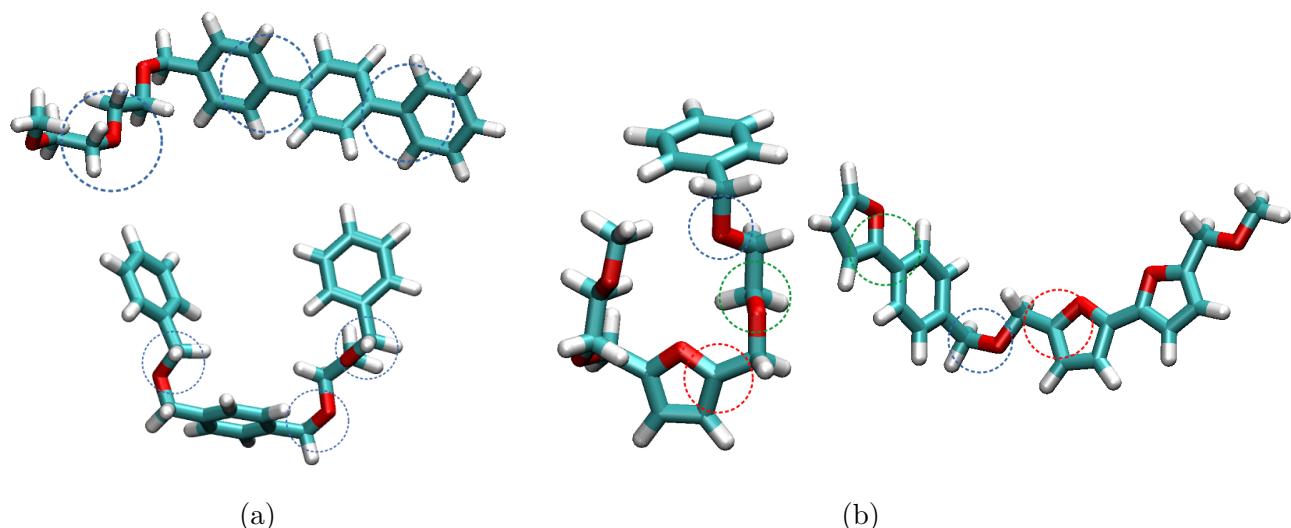


Figure 2: Creation of a random polymer with user defined probabilities. (a) Diblock copolymer architecture, where a single probability (P_1) of being linked is 0.4. The linkages are highlighted with a blue circle. Two different polymeric topologies are obtained using the same probability. (b) Triblock copolymer architecture employing a user-defined probabilities P_1 (represented in green) and P_2 (displayed in red). Two different moieties are obtained based on user defined probabilities (0.4 and 0.2), respectively.

Additionally, PySoftK supports the construction of random linear copolymers which consist of two or three different monomers. The module `Rnp(mol_1,mol_2,atom).random_ab_copolymer`

(`len`,`pA`,`iter_ff`,`FF`) is used to build a random copolymer of length `len` that contains two monomers, `mol_1` and `mol_2` which are decorated with atoms of type `atom` to indicate where the polymerisation occurs. The number of `mol_1` and `mol_2` monomers in the resulting polymer are determined from `pA*len` and `(1-pA)*len`, respectively. Similarly PySoftK can be used to generate random linear copolymers consisting of three different monomers via the module `Rnp(mol_1,mol_2,atom).random_abc_copolymer(mol_3, len, pA, pB, iter_ff, FF)`. In this case, the resulting polymer will contain `len` monomers, such that the number of `mol_1`, `mol_2` and `mol_3` monomers is `pA*len`, `pB*len` and `(1-pA-pB)*len`, respectively. A snippet showing the usage of the functions is displayed in Listing 3, and examples of the polymers produced by that code are presented in Figure 2.

```
1 from pysftk.topologies.ranpol import *
2
3 mol_1=Chem.MolFromSmiles('BrCOCBr')
4 mol_2=Chem.MolFromSmiles('c1(ccc(cc1Br)Br)Br')
5 mol_3=Chem.MolFromSmiles('c1cc(oc1Br)Br')
6
7 #diblock random polymer
8 dia = Rnp(mol_1, mol_2, "Br").random_ab_copolymer(5, 0.4, 10)
9 Fmt(dia).xyz_print("dia.xyz")
10
11 #triblock random polymer
12 tri = Rnp(mol_1, mol_2, "Br").random_abc_copolymer(mol_3, 5, 0.4, 0.2, 10)
13 Fmt(tri).xyz_print("tri.xyz")
```

Listing 3: Code snippet showing the usage of the random copolymer module functions. The functions `random..ab_copolymer` and `random...abc_copolymer` allowing the definition of the user-supplied linking probabilities.

Ring polymers. Current experimental techniques allow a precise control of polymer architectures enabling the creation of new topologies. One interesting example is ring, or cyclic, polymers, which can be described as closed macromolecular structures with no beginning or

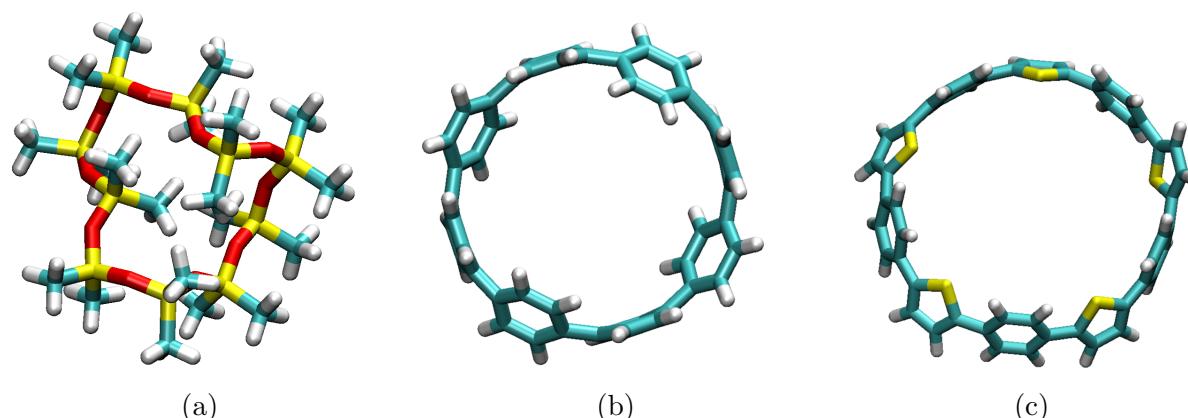


Figure 3: Creation of a ring polymers using (a) silanol⁴² and (b) benzene ring. (c) Alternating copolymer (benzene and thiol monomers combined via the `s_m` module). This example illustrates the capabilities of PySoftK to utilize previously developed algorithms for creating new architectures.

end.⁴¹ PySoftK enables the creation of ring homopolymers with the module `Rn(mol_1,atom)`.`pol_ring(len,FF,iter_ff)` that generates a ring polymer with `len mol_1` monomers by first employing the linear polymer module to generate a polymer chain and then initial and final atomic placeholders (`atom`) are used to create a bond, which generates a ring topology, as displayed in Figure 3(a,b). A snippet of code utilising this command is shown in Listing 4. It is worth mentioning that PySoftK enables the combination of modules to create new structures. This is demonstrated in the bottom part of Listing 4, where the `sm` module is used to create a linear diblock polymer that is then converted to a ring polymer topology using the `Rn` module, with the resultant ring polymer shown in Figure 3(c).

```

1 from pysoftk.topologies.ring import *
2 from pysoftk.linear_polymer.super_monomer import *
3
4 # Silanol definition using SMILES and PySoftK ring function
5 mol_1 = Chem.MolFromSmiles('C[Si](OBr)(C)Br')
6 rn2=Rn(mol_1,'Br').pol_ring(10, "MMFF", 250)
7
8 #Using the Sm module of PySoftK to generate a diblock cyclic polymer
9 mol_1=Chem.MolFromSmiles('c1cc(sc1Br)Br')

```

```
1  
2  
3 mol_2=Chem.MolFromSmiles('c1(ccc(cc1)Br)Br')  
4  
5  
6 a=Sm(mol_1,mol_2,"Br").mon_to_poly()  
7  
8 cyc=Rn(a,'Br').pol_ring(5, FF="UFF", iters=2000)  
9
```

Listing 4: Code snippet showing the creation of a cyclic polymer

Branched Polymers. Branched polymers are another topology that has been enabled in PySoftK. Generally, branched polymers are a set of secondary polymer chains linked to a primary backbone. PySoftK builds this topology by employing the user-supplied atomic placeholder as an indicator of the number of branch points from the primary backbone present in this structure. Thus, for instance, a backbone with four placeholders will generate a model where four branches (arms) are attached at the regions on the backbone indicated by the user, as shown in Figure 4. In PySoftK, the module `Bd(core,arm,atom).branched_polymer(FF, ff_iter)` is used to generate a branched polymer with a backbone of `core` and arms described by `arm` (Listing 5). The placeholder atoms are defined by `atom`. The inputted structures for `core` and `arm` can be simple monomers and therefore inputted as previously discussed or they can be resultant structures from any of the other commands and then inputted into the branched function. Therefore, the branched architecture not only illustrates the capabilities of PySoftK to create new architectures from scratch but also the versatility of PySoftK to be extended. In this case, this has been done by creating the module `topologies` where all the previous functions have been organized.

```
1 from pysoftk.topologies.branched import *  
2  
3 # Core and Arm molecules  
4 core=Chem.MolFromSmiles('BrN(Br)CCN(Br)Br')  
5 arm=Chem.MolFromSmiles('C(CC(=O)OCCOC(=O)CCBr)Br')  
6 final = Bd(core, arm, "Br").branched_polymer()
```

Listing 5: Code snippet showing the creation of a branched polymer

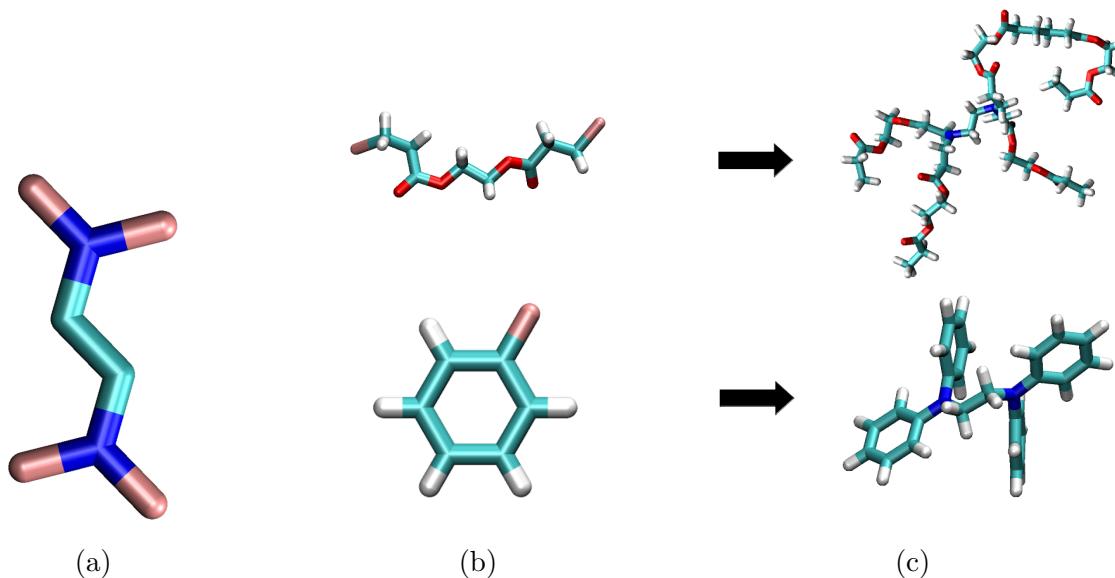


Figure 4: Creation of a branched polymers using (a) ethylenediamine as the primary backbone and (b) benzene or polyester oligomer as branches. (c) Branched polymer with provided core and arms moieties. In this case, only one placeholder atom is used to signal PySoftK the region where a bond would be formed.

Specifying absolute monomer sequences. Any linear polymer can be described by listing the combination of monomers in a specific pattern (e.g. ABBACCABBACC, ABCABABCBA). Therefore, PySoftK offers the option to build polymers with a specific pattern of monomers using an alphabet-based pattern expressed in a single string followed by a list of RDKit⁴⁰ molecular objects as presented in Listing 6. In Figure 5, thiophene, furan and benzene monomers are used to present the different possibilities provided by this function to construct all unique permutations.

```
1 from pysoftk.topologies.diblock import *
2
3 # Core and Arm molecules
4 mols=[Chem.MolFromSmiles('c1(ccc(cc1)Br)Br'),
5
6     Chem.MolFromSmiles('c1cc(oc1Br)Br'),
7
8     Chem.MolFromSmiles('c1cc(sc1Br)Br')]
9
10 patt="ABC"
```

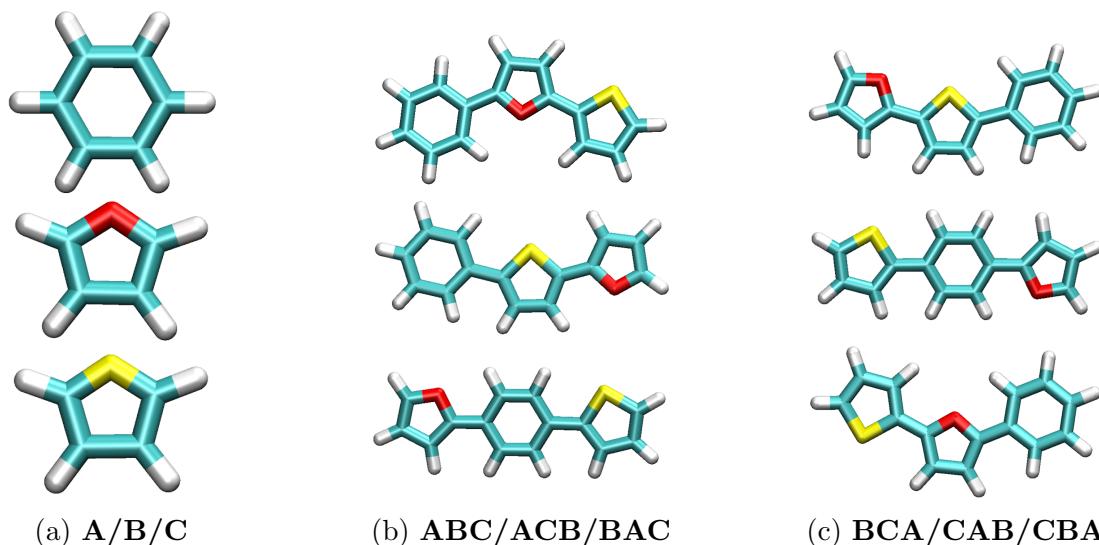


Figure 5: Figure (4) Patterned construction of a polymer based on a user-defined pattern. (a) An user created list is mapped to an alphabetic string representing the position of the molecule in the list. (b) Possible unique permutations for an arbitrary list containing 3 elements.

```
9 a=Pt(patt, mols, "Br").pattern_block_poly()
```

Listing 6: Code snippet showing the creation of a patterned polymer

This function can also be used in conjunction with the `linear_polymer` module of PySoftK to construct polymeric macromolecules that have other polymers as their monomeric units.

Facilitating High-Throughput Calculations

Folder creation and organization. The module `pysoftk.folder_manager.folder_creator` has been designed to create a user-defined number of folders with unique names. In the case of high-throughput calculations (HTC), usually many different systems are created in a single folder and then relocated to enable calculations or post-processing analyses to be conducted. As shown in Listing 7, this workflow can be carried out utilizing the function `Fld().file_to_dir` implemented in PySoftK. The result of the code in Listing 7 is that two new, uniquely named directories would be created and each of the two `.smi` files (`mol_1.smi` & `mol_2.smi`) would be moved into one of these new directories, such that each new directory

would contain one .smi file.

```
1 from pysoftk.folder_manager.folder_creator import *
2
3 print('c1cc(sc1Br)Br', file=open("mol_1.smi", 'w'))
4 print('c1(ccc(cc1Br)Br', file=open("mol_2.smi", 'w'))
5
6 # Test 3: Seek for .smi files and create the needed folders.
7 Fld().file_to_dir("smi", 2)
```

Listing 7: Code snippet showing the automatic creation of folder to organize files based on an user-defined extension.

The aforementioned function is able to search for files with a given file extension, and relocate them within individually created folders. Likewise, for cases where many files are present, this command can be executed in parallel. The creation of automated workflows (as displayed in Figure 6) can be achieved by combining many different functions in PySoftK enabling the modeling of thousands of polymers using a single script. An example demonstrating how this can be done can be found in the SI. Alternatively, these modules can be used separately as part of another workflow that generates other types of directory structures or simply a folder where many different files are located.

Automatic torsional angle detection for conjugated polymers. PySoftK offers a tool to perform analysis of the torsion angles within conjugated polymers that connect the ring subunits that are commonly found in their backbone. These dihedrals are generally poorly captured by the existing parameters within classical forcefields.^{28,39} Thus having the capability to rapidly identify those dihedrals then allows the user to easily set up the necessary *ab initio* simulations required to determine the potential energy landscape of those dihedrals, which then can be used to reparameterize the forcefield for those dihedrals.

The module `pysoftk.torsional` automatically detects and reports the atoms involved in the torsional angles found within planar conjugated polymers. PySoftK also enables the creation of molecular sketches highlighting and reporting the atom numbers that form a

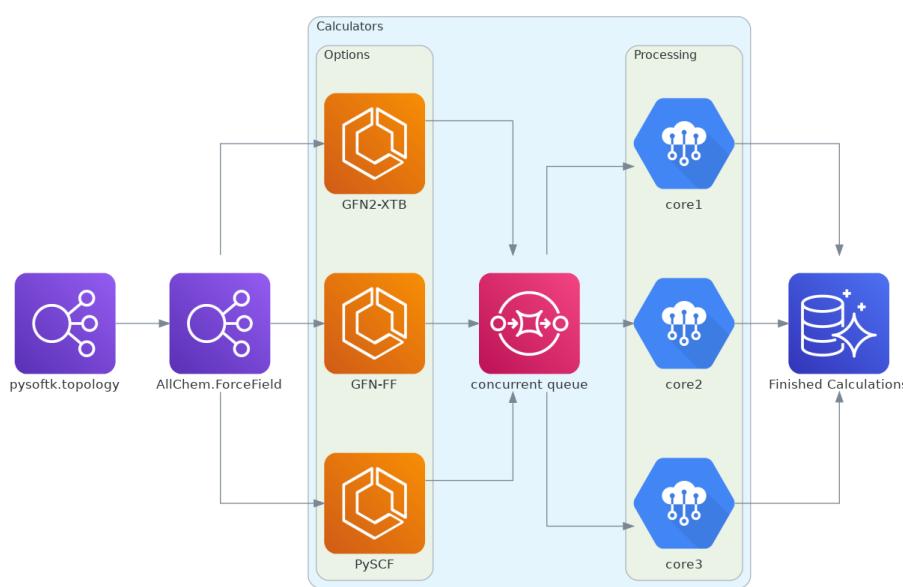


Figure 6: PySoftK workflow generation. Topology objects (polymers) are created and automatically parsed to RDKit force-field optimization.⁴⁰ External programs such as GFN2-XTB, GFN-FF or PySCF are linked using internal modules which facilitates the automation process.^{43,44} Concurrent parallelization allow the creation of events in which subsequent parallelization strategies can be used. File organization is performed.

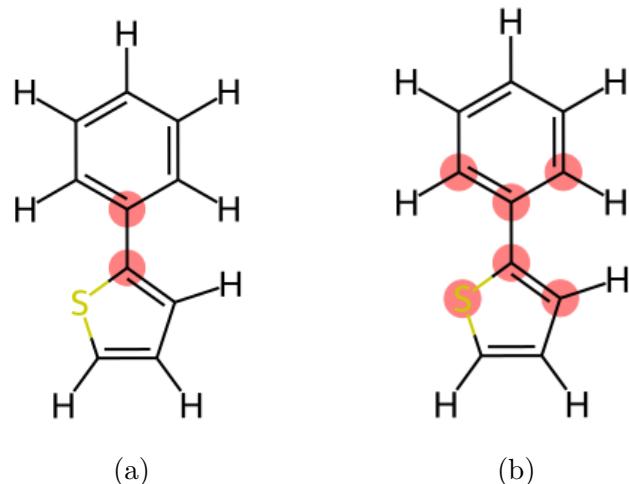


Figure 7: Graphical explanation of the **P-Graph** description and algorithm used for automatically detecting intermolecular torsional angles. (a) Bond linking atoms are detected. (b) All connecting paths are computed where previously detected atoms are used as starting and end points.

torsional angle as shown in Figure 7. The topological fingerprint used in this module relies on the description of molecules as graphs where the atoms are nodes and the covalent bonds

1
2
3 are the edges.⁴⁵ Based on the idea to convert a molecular Graph (**M-Graph**) into a path
4 Graph (**P-Graph**), we have been able to label connections between edges in the **M-Graph**
5 providing a nomenclature to name the **P-Graph**.⁴⁶
6
7

8 In order to identify the torsions between ring subunits within the backbone of a conju-
9 gated polymer, we have identified two conditions that the atoms (e.g. nodes) making up the
10 central bond of the torsion must meet. First, an atom must have three neighbours which
11 are from the inner structure (such as rings) and one provided from the next monomer (as
12 displayed in Fig. 7a). However, this condition can be also be found in atoms that are embed-
13 ded within aromatic rings. To avoid this, we have used the ring detection function provided
14 by RDKit to remove potential bonds within a ring, and therefore identify the bonds which
15 connect the ring moieties in the polymer (Fig. 7a). After identifying all of these important
16 bonds within the polymer, we then identify all of the torsions that include these bonds as the
17 central bond and therefore can provide the relevant atomic labels involved in these torsional
18 angles as shown in Figure 7b.⁴⁷
19
20

21 This module then reports all of these important torsional angles within the molecule, and
22 generates 2D molecular sketches of each one. Figure 8 shows an example of this output for
23 a boroxine polymer.⁴⁸
24
25

39 Data Provenance and Software Development

40
41

42 One of our main objectives while creating PySoftK is to provide tools to the community that
43 allow users to utilize high-performance computational facilities to automate the creation of
44 databases for polymeric structures as easily as possible. To do so, we have greatly simplified
45 the installation process of PySoftK by employing the pip command strategy. PySoftK has
46 been tested in Linux and Mac operating systems based on python 3.6 or more recent versions.
47 Parallel strategies have been developed in parts of the code to enhance the scalability in tasks
48 such as HTC or general organization of molecular databases. In this sense, we have developed
49
50
51
52
53
54
55
56
57
58

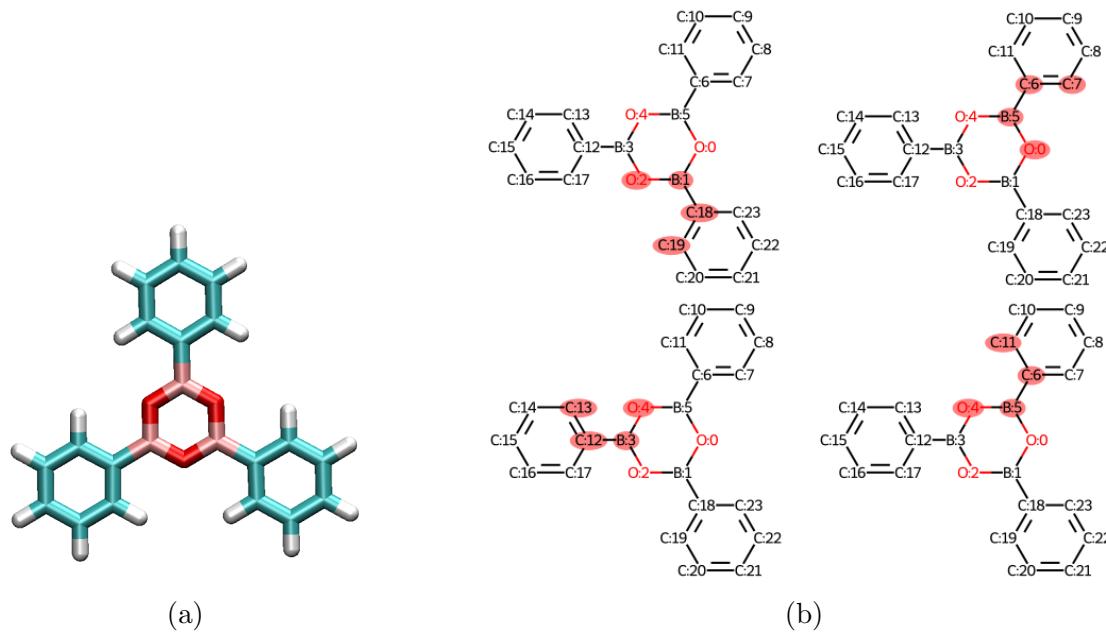


Figure 8: Automatic recognition of intermolecular torsional angles within a planar polymer. (a) Boroxine planar polymer. (b) 2D depiction of selected torsional angles highlighting the atoms and its numerical labels.

modules of PySoftK using concurrency where independently executing tasks are created and queued to use the available resources. This approach ensures that the `calculator` module can run utilizing all available resources combining task assignment and core usage for parallel codes such as GFN2-xTB⁴³ while efficiently avoiding bottlenecks, as displayed in Figure 9.⁴⁹ In this sense, the user can maximize the parallel performance of the code by tuning variables such as number of available processors against simultaneous task management.

PySoftK has been constructed using principles of modern code development practices. Thus, continuous integration (CI) strategies have been incorporated to probe the code integrity against suggested changes. In our case, these tests have been designed not only as a showcase of all of the capabilities of PySoftK, but also with the aim of maximum code coverage achieved with our current test set. This design ensures that new commits are compatible with the majority of the modules before a new version is integrated into the main branch and released.

Continuous deployment (CD) is achieved once all previous tests have passed enabling

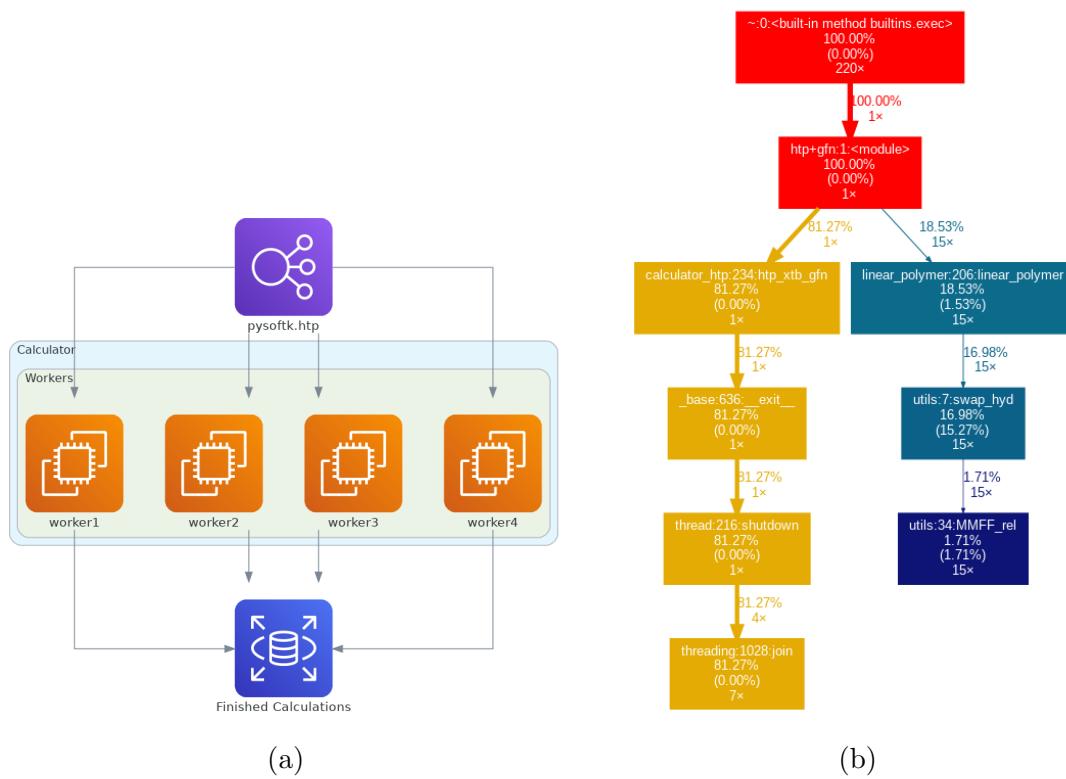


Figure 9: Parallelization strategies implemented in PySoftK. (a) Schematic diagram presenting the distribution of user defined workers (threads). Each worker performs a calculation using the `pysoftk.calculator` object where further parallelization schemes can be used (cores). (b) Usage report of a calculation performed using an Intel(R) Core(TM) i9-900 CPU 3.10 GHz computer employing 16 polymers created on-the-fly in a 4×4 scheme where 4 threads are performing task management while 4 cores are used to perform *ab-initio* calculations at the GFN2-xTB level of theory.^{43,49,50}

an agile updating of PySoftK. At this point, we would like to comment, that a successful CD is only achieved when the corresponding explanatory examples are included in the documentation. This process is greatly facilitated since the documentation is also part of PySoftK building process. Finally, a detailed documentation (which is also part of our CI/CD workflow) has been developed showing working examples (based on our tests) alongside of tutorials for all the different features enabled in PySoftK. They have been tested by members of our community and feedback has been incorporated into the latest version of the code, ensuring a clear and concise approach for new users.

Conclusions

PySoftK is a modular and versatile software that has been designed to facilitate high-throughput calculations of polymeric systems. The code contains a range of modules that allow its users to generate polymers with a uniquely broad range of topologies and compositions. Additionally, the code has unique tools to assist in the file and directory structure management which is inherent in high-throughput calculations. All of the functionality within PySoftK has been developed to be embedded as part of user-defined workflows enabling steps such as modeling or computing using the specifically defined modules. In a complementary fashion, PySoftK allows the user to keep a track-record of the data produced, facilitating post-processing analysis that can be also part of the same workflow. An ample set of testing scripts has been created aiming at a high-coverage of the code which ensures that future algorithm developments are preserving the code structure. Further documentation and tutorials are available on the PySoftK website (<https://alejandrosantanabonilla.github.io/pysoftk/#>). The modular nature of the code allows for the user community to easily contribute code that complements the existing functionality of this published version of the code.

Associated Content

Data Availability Statement

All of the data and scripts necessary to utilise PySoftK are provided within the text of this manuscript and/or within the website for PySoftK (<https://alejandrosantanabonilla.github.io/pysoftk/#>). Within the website of PySoftK we also have additional tutorials for using the various aspects of the code. Finally, the code is freely available from the following github page: <https://github.com/alejandrosantanabonilla/pysoftk>.

Acknowledgement

We are grateful to the UK Materials and Molecular Modelling Hub, which is partially funded by EPSRC (EP/P020194/1 and EP/T022213/1) and the UK HPC Materials Chemistry Consortium, which is also funded by EPSRC (EP/R029431) for providing us access to computational resources. This work also benefitted from access to the King's Computational Research, Engineering and Technology Environment (CREATE) at King's College London.⁵¹ R. L.-R. D. C. acknowledges the support by the Biotechnology and Biological Sciences Research Council (BB/T008709/1) via the London Interdisciplinary Doctoral Programme (LIDo). R. M. Z. and C. D. L. acknowledge the Engineering and Physical Sciences Research Council (EPSRC) for funding (EP/V049771/1).

For the purpose of open access, the author has applied a Creative Commons Attribution (CC BY) licence (where permitted by UKRI, 'Open Government Licence' or 'Creative Commons Attribution No-derivatives (CC BY-ND) public copyright licence' may be stated instead) to any Author Accepted Manuscript version arising.

Supporting Information Available

Results of a computational benchmark for testing the Htp module of PySoftK altogether with the pictures of all bonds detected by Boroxine are displayed in the electronic supporting information.

References

- (1) Foglia, F. et al. Disentangling water, ion and polymer dynamics in an anion exchange membrane. *Nat. Mater.* **2022**, *21*, 555–563.
- (2) Foglia, F.; Frick, B.; Nania, M.; Livingston, A. G.; Cabral, J. T. Multimodal confined

- 1
2
3 water dynamics in reverse osmosis polyamide membranes. *Nat. Commun.* **2022**, *13*,
4 2809.
5
6
7 (3) Getzler, Y. D. Y. L.; Mathers, R. T. Sustainable polymers: Our evolving understanding.
8 *Acc. Chem. Res.* **2022**, *55*, 1869–1878.
9
10 (4) Wang, S.; Zuo, G.; Kim, J.; Sirringhaus, H. Progress of conjugated polymers as emerging
11 thermoelectric materials. *Prog. Polym. Sci.* **2022**, *129*, 101548.
12
13 (5) Wang, Q.; Huang, L.; Wang, Z.; Zheng, J.; Zhang, Q.; Qin, G.; Li, S.; Zhang, S. High
14 conductive anion exchange membranes from all-carbon twisted intrinsic microporous
15 polymers. *Macromolecules* **2022**, *55*, 10713–10722.
16
17 (6) Yang, X.; Bai, R.; Zhang, Z.; Liu, Y.; Yan, X. Mechanically tunable supramolecular
18 polymer networks with different triblock backbones. *J. Polym. Sci.* **2022**, doi:
19 10.1002/POL.20220615.
20
21 (7) Li, Z.; Shen, Z.; Pei, Y.; shuang chao,; Pei, Z. Covalently bridged pillararene-
22 based polymers: structures, synthesis, and applications. *Chem. Comm.* **2023**, doi:
23 10.1039/D2CC05594E.
24
25 (8) Coscia, B. J.; Shelley, J. C.; Browning, A. R.; Sanders, J. M.; Chaudret, R.; Rozot, R.;
26 Léonforte, F.; Halls, M. D.; Luengo, G. S. Shearing friction behaviour of synthetic
27 polymers compared to a functionalized polysaccharide on biomimetic surfaces: models
28 for the prediction of performance of eco-designed formulations. *Phys. Chem. Chem.*
29 *Phys.* **2023**, doi: 10.1039/D2CP05465E.
30
31 (9) Gouveia, M. G.; Wesseler, J. P.; Ramaekers, J.; Weder, C.; Scholten, P. B. V.; Bruns, N.
32 Polymersome-based protein drug delivery – quo vadis? *Chem. Soc. Rev.* **2023**, doi:
33 10.1039/D2CS00106C.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 (10) Zhang, Y.; Chen, J.; Shi, L.; Ma, F. Polymeric nanoparticle-based nanovaccines for
4 cancer immunotherapy. *Mater. Horiz.* **2023**, doi: 10.1039/D2MH01358D.
5
6 (11) Oliveira, A. S.; Pereira, P.; Mendonça, P. V.; Fonseca, A. C.; Simões, S.; Serra, A. C.;
7 Coelho, J. F. Novel degradable amphiphilic 4-arm star PLA-b-POEOA and PLGA-b-
8 POEOA block copolymers: synthesis, characterization and self-assembly. *Polym. Chem.*
9 **2023**, *14*, 161–171.
10
11
12
13
14 (12) Wang, Y.; Su, H.; Wang, Y.; Cui, H. Discovery of Y-shaped supramolecular polymers in
15 a self-assembling peptide amphiphile system. *ACS. Macro. Lett.* **2022**, *14*, 1355–1361.
16
17
18 (13) Li, W.; Hadjigol, S.; Mazo, A. R.; Holden, J.; Lenzo, J.; Shirbin, S. J.; Barlow, A.;
19 Shabani, S.; Huang, T.; Reynolds, E. C.; Qiao, G. G.; O'Brien-Simpson, N. M. Star-
20 peptide polymers are multi-drug-resistant hram-positive bacteria killers. *ACS Appl.*
21
22 *Mater. Interfaces* **2022**,
23
24
25
26
27
28
29
30 (14) Grest, G. S.; Ge, T.; Plimpton, S. J.; Rubinstein, M.; O'Connor, T. C. Entropic mixing
31 of ring/linear polymer blends. *ACS Polym. Au* **2022**, doi: 10.1021/ACSPOLYMER-
32 SAU.2C00050.
33
34
35
36 (15) Zeng, X. C.; Zhao, H.; Chen, C.; Weil, T. Cyclic polymers: synthesis, characteristics,
37 and emerging applications. *Nanoscale Horiz.* **2022**, *2*, 1–66.
38
39
40
41 (16) Liu, M.; Blankenship, J. R.; Levi, A. E.; Fu, Q.; Hudson, Z. M.; Bates, C. M. Miktoarm
42 star polymers: Synthesis and applications†. *Chem. Mater.* **2022**, *34*, 6188–6209.
43
44
45
46 (17) Kostka, L.; Kotrchová, L.; Randárová, E.; Ferreira, C. A.; Malátová, I.; Lee, H. J.;
47 Olson, A. P.; Engle, J. W.; Kovář, M.; Cai, W.; Šírová, M.; Etrych, T. Evaluation of
48 linear versus star-like polymer anti-cancer nanomedicines in mouse models. *J. Control.*
49
50 *Release* **2023**, *353*, 549–562.
51
52
53
54
55
56
57
58
59
60

- (18) Gauthier-Jaques, M.; Mutlu, H.; Theato, P. Cage-shaped polymers synthesis: A comprehensive state-of-the-art. *Macromol. Rapid Commun.* **2022**, *32*, 2100760.
- (19) Jing, Y.; Zhu, X.; Maier, S.; Heine, T. 2D conjugated polymers: exploiting topological properties for the rational design of metal-free photocatalysts. *Trends Chem.* **2022**, *4*, 792–806.
- (20) Muramatsu, Y.; Takasu, A. Synthetic innovations for cyclic polymers. *Polym. J.* **2021**, *54*, 121–132.
- (21) Ziolek, R. M.; Omar, J.; Hu, W.; Porcar, L.; González-Gaitano, G.; Dreiss, C. A.; Lorenz, C. D. Understanding the pH-directed self-assembly of a four-arm block copolymer. *Macromolecules* **2020**, *53*, 11065–11076.
- (22) Xue, Y. L.; Huang, J.; Lau, C. H.; Cao, B.; Li, P. Tailoring the molecular structure of crosslinked polymers for pervaporation desalination. *Nat. Commun.* **2020**, *11*, 1461.
- (23) Liu, J.; Zheng, Y.; Zhao, Z.; Yuan, M.; Tsige, M.; Wang, S.-Q. Investigating nature of stresses in extension and compression of glassy polymers via stress relaxation. *Polymer* **2020**, *202*, 122517.
- (24) Ziolek, R. M.; Smith, P.; Pink, D. L.; Dreiss, C. A.; Lorenz, C. D. Unsupervised learning unravels the structure of four-arm and linear block copolymer micelles. *Macromolecules* **2021**, *54*, 3755–3768.
- (25) Scherillo, G.; Mensiteri, G.; Baldanza, A.; Loianno, V.; Musto, P.; Pannico, M.; Correa, A.; Nicola, A. D.; Milano, G. Weak interactions between poly(ether imide) and carbon dioxide: A multiscale investigation combining experiments, theory and simulations. *Macromolecules* **2022**, *55*, 10773–10787.
- (26) Mohottalalage, S. S.; Senanayake, M.; Clemmer, J. T.; Perahia, D.; Grest, G. S.;

- O'Connor, T. C. Nonlinear elongation flows in associating polymer melts: From homogeneous to heterogeneous flow. *Phys. Rev. X* **2022**, *12*, 021024.
- (27) Wang, J.; O'Connor, T. C.; Grest, G. S.; Ge, T. Superstretchable elastomer from cross-linked ring polymers. *Phys. Rev. Lett.* **2022**, *128*, 237801.
- (28) Ziolek, R. M.; Santana-Bonilla, A.; Castro, R. L.-R. D.; Kuhn, R.; Green, M.; Lorenz, C. D. Conformational heterogeneity and interchain percolation revealed in an amorphous conjugated polymer. *ACS Nano* **2022**, *16*, 14432–14442.
- (29) Fortunato, M. E.; Colina, C. M. *pysimm*: A python package for simulation of molecular systems. *SoftwareX* **2017**, *6*, 7–12.
- (30) Demidov, A. G.; Perera, B. L. A.; Fortunato, M. E.; Lin, S.; Colina, C. M. Update 1.1. to "pysimm: A python package for simulation of molecular systems". *SoftwareX* **2021**, *15*, 100749.
- (31) Girard, M.; Ehlen, A.; Shakya, A.; Bereau, T.; de la Cruz, M. O. Hoobas: A highly object-oriented builder for molecular dynamics. *Comput. Mater. Sci.* **2019**, *167*, 25–33.
- (32) Cummings, P. T. et al. Open-source molecular modeling software in chemical engineering focusing on the Molecular Simulation Design Framework. *AIChE J.* **2021**, *67*, e17206.
- (33) Yan, X.; Chaudhuri, S. An interactive polymer building toolkit for molecular dynamics simulations: PolyMAPS. *arXiv* **2022**, *2204*, 14218.
- (34) Sahu, H.; Shen, K.-H.; Montoya, J. H.; Tran, H.; Ramprasad, R. Polymer Structure Predictor (PSP): A Python toolkit for predicting atomic-level structural models for a range of polymer geometries. *J. Chem. Theory Comput.* **2022**, *18*, 2737–2748.

- 1
2
3 (35) Vanommeslaeghe, K.; A. D. MacKerell, J. Automation of the CHARMM General Force-
4 field (CGenFF) I: Bond perception and atom typing. *J. Chem. Inf. Model.* **2012**, *52*,
5 3144–3154.
6
7
8 (36) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and
9 testing of a general Amber force field. *J. Comput. Chem.* **2004**, *25*, 1157–1174.
10
11 (37) Dodda, L. S.; de Vaca, I. C.; Tirado-Rives, J.; Jorgensen, W. L. LigParGen Web Server:
12 An automatic OPLS-AA parameter generator for organic ligands. *Nucleic Acids Res.*
13 **2017**, *45*, W331–W336.
14
15 (38) Plimpton, S. J. Fast parallel algorithms for short-range molecular dynamics. *J. Comput.*
16
17 *Phys.* **1995**, *117*, 1–19.
18
19 (39) Wolf, C. M.; Guio, L.; Scheiwiller, S.; Pakhnyuk, V.; Luscombe, C.; Pozzo, L. D.
20 Strategies for the development of conjugated polymer molecular dynamics force fields
21 validated with neutron and x-ray scattering. *ACS Polym. Au* **2021**, *1*, 134–152.
22
23 (40) RDKit: Open-source cheminformatics. <https://www.rdkit.org>.
24
25 (41) You, Y.; Chen, J.; Zheng, A.; Wei, D.; Xu, X.; Guan, Y. Effect of silanol on the
26 thermal stability of poly [methyl (trifluoropropyl) siloxane]. *J. Appl. Polym. Sci.* **2020**,
27 *137*, 49347.
28
29 (42) Liu, Y.; Chaiprasert, T.; Ouali, A.; Unno, M. Well-defined cyclic silanol derivatives.
30 *Dalton Trans.* **2022**,
31
32 (43) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-xTB—An accurate and broadly
33 parametrized self-consistent tight-binding quantum chemical method with multipole
34 electrostatics and density-dependent dispersion contributions. *J. Chem. Theory Com-*
35 *put.* **2019**, *15*, 1652–1671.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 (44) Sun, Q.; Zhang, X.; Banerjee, S.; Bao, P.; Barbry, M.; Blunt, N. S.; Bogdanov, N. A.;
4 Booth, G. H.; Chen, J.; Cui, Z.-H., et al. Recent developments in the PySCF program
5 package. *J. Chem. Phys.* **2020**, *153*, 024109.
6
7
8
9 (45) Gutman, I.; Estrada, E. Topological indices based on the line graph of the molecular
10 graph. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 541–543.
11
12
13 (46) Hanser, T.; Jauffret, P.; Kaufmann, G. A new algorithm for exhaustive ring perception
14 in a molecular graph. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 1146–1152.
15
16
17 (47) Hagberg, A.; Conway, D. NetworkX: Network Analysis with Python. *URL:*
18 <https://networkx.github.io> **2020**,
19
20
21 (48) Cheng, F.; Jäkle, F. Boron-containing polymers as versatile building blocks for func-
22 tional nanostructured materials. *Polym. Chem.* **2011**, *2*, 2122–2132.
23
24
25 (49) Kwon, M. Diagrams: Diagrams as code. <https://diagrams.mingrammer.com/>, 2022.
26
27
28 (50) J. Fonseca, e. gprof2dot. <https://github.com/jrfonseca/gprof2dot>, 2022.
29
30
31 (51) King's College London (2022). King's Computational Research, Engineering and Tech-
32 nology Environment (CREATE). <https://doi.org/10.18742/rnfv-m076>.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

TOC Graphic

