



1 **Facilitating the structural characterisation of non-canonical amino acids**
2 **in biomolecular NMR**

3 Sarah Kuschert¹, Martin Stroet², Yanni Ka-Yan Chin¹, Anne Clair Conibear³, Xinying Jia¹, Thomas Lee²,
4 Christian Reinhard Otto Bartling⁴, Kristian Strømgaard⁴, Peter Güntert⁵, Karl Johan Rosengren⁶, Alan
5 Edward Mark² and Mehdi Mobli¹

6 ¹Centre for Advanced Imaging, The University of Queensland, Brisbane, QLD 4072, Australia

7 ²School of Chemistry & Molecular Biosciences, The University of Queensland, Brisbane, QLD 4072, Australia

8 ³Institute of Applied Synthetic Chemistry, Technische Universität Wien, Getreidemarkt 9/163, Wien 1060, Vienna, Austria

9 ⁴Department of Drug Design and Pharmacology, University of Copenhagen, Universitetsparken 2, 2100 Copenhagen,
10 Denmark

11 ⁵Laboratory of Physical Chemistry, ETH Zürich, 8093 Zürich, Switzerland; Institute of Biophysical Chemistry, Center for
12 Biomolecular Magnetic Resonance, Goethe University Frankfurt, 60438 Frankfurt am Main, Germany; Department of
13 Chemistry, Tokyo Metropolitan University, Hachioji, Tokyo 192-0397, Japan

14 ⁶School of Biomedical Sciences, The University of Queensland, Brisbane, QLD 4072, Australia

15 *Correspondence to:* Mehdi Mobli (m.mobli@uq.edu.au)



16 **Abstract**

17 Peptides and proteins containing non-canonical amino acids (ncAAs) are a large and important class of biopolymers. They
18 include non-ribosomally synthesised peptides, post-translationally modified proteins, expressed or synthesised proteins
19 containing unnatural amino acids, and peptides and proteins that are chemically modified. Here, we describe a general
20 procedure for generating atomic descriptions required to incorporate ncAAs within popular NMR structure determination
21 software such as CYANA, CNS, Xplor-NIH and ARIA. This procedure is made publicly available via the existing Automated
22 Topology Builder (ATB) server (atb.uq.edu.au) with all submitted ncAAs stored in a dedicated database. The described
23 procedure also includes a general method for linking of sidechains of amino acids from CYANA templates. To ensure
24 compatibility with other systems, atom names comply with IUPAC guidelines. In addition to describing the workflow, 3D
25 models of complex natural products generated by CYANA are presented, including vancomycin. In order to demonstrate the
26 manner in which the templates for ncAAs generated by the ATB can be used in practice we use a combination of CYANA and
27 CNS to solve the structure of a synthetic peptide designed to disrupt Alzheimer-related protein-protein interactions.
28 Automating the generation of structural templates for ncAAs will extend the utility of NMR spectroscopy to studies of more
29 complex biomolecules, with applications in the rapidly growing fields of synthetic and chemical biology. The procedures we
30 outline can also be used to standardise the creation of structural templates for any amino acid and thus have the potential to
31 impact structural biology more generally.



32 **1 Introduction**

33 The 20 genetically encoded amino acids, together with selenocysteine and pyrrolysine, provide the basis for most proteins and
34 peptides that make up the machinery of life (Liu and Schultz, 2010; Huang et al., 2010; Bullwinkle et al., 2014). The use of
35 just 22 amino acids, however, limits the structural complexity and functional diversity that can be achieved. The chemical and
36 functional diversity of ribosomally synthesised proteins is further expanded by posttranslational modification (PTM), including
37 processes such as acylation, methylation, phosphorylation, oxidation and epimerisation (Aebersold et al., 2018; Barber and
38 Rinehart, 2018; Walsh et al., 2005). Non-ribosomal synthesis pathways also expand chemical diversity, leading to the
39 production of typically short peptides containing non-canonical amino acids (ncAAs), as well as backbone or sidechain
40 cyclisation (Link et al., 2003; Tharp et al., 2020; Caboche et al., 2008; Goodrich and Frueh, 2015; Martínez-Núñez and López,
41 2016; Strieker et al., 2010). Such non-ribosomal peptides (NRPs) are prevalent in bacteria and fungi, which produce a wide
42 range of bioactive peptides (Caboche et al., 2008; Marahiel, 2009). In addition to pathways found in nature, chemical synthesis,
43 enzymatic modification, genetic code expansion and site-selective biorthogonal transformations are increasingly used to
44 introduce novel ncAAs and PTMs into peptides and proteins (Bondalapati et al., 2016; Chuh et al., 2016; Conibear, 2020; Hoyt
45 et al., 2019; Thompson and Muir, 2020) for both investigations of peptide structure and function as well as in the design and
46 optimisation of pharmaceuticals (Noren et al., 1989; Coin, 2018; Hoesl and Budisa, 2011; Johnson et al., 2010; Wang et al.,
47 2001; Wang et al., 2020; Zou et al., 2018).

48 Despite the prevalence of ncAAs and their importance in determining the functional properties of both naturally occurring and
49 synthetic peptides, the structural characterisation of peptides and proteins containing ncAAs remains challenging. For example,
50 of the almost 200 000 structures in the Protein Data Bank (PDB), only 11 677 were annotated with a PTM (Craveur et al.,
51 2014). This paucity of protein structures bearing ncAAs also results in them being excluded from machine learning and
52 structure prediction algorithms, as there is an insufficient training set. Techniques used to determine the structure of proteins
53 and peptides (X-ray diffraction, cryo-electron microscopy and nuclear magnetic resonance (NMR)) rely heavily on modelling
54 software to transform the experimental data into structural models (Mal et al., 2002). The amount of data that can be collected
55 experimentally in NMR is, in general, insufficient to determine the structure of a given peptide or protein directly. Instead, a
56 representation of the spatial arrangement of the atoms within each amino acid, together with a description of the interactions
57 between sets of atoms, is required to translate a set of experimental restraints into a three-dimensional (3D) structure (Mal et
58 al., 2002). Most molecular modelling packages contain only the 20 canonical amino acids and a modest selection of the most
59 common ncAAs. This is true for both general molecular dynamics simulation packages (such as AMBER, CHARMM,
60 GROMACS and GROMOS) as well as software dedicated to structure refinement such as Xplor-NIH (Bermejo and
61 Schwieters, 2018), CNS (Brunger et al., 1998), ARIA (Mareuil et al., 2015; Allain et al., 2020) or CYANA (Guntert and
62 Buchner, 2015; Guntert et al., 1997). Despite the various input formats required for the different structural calculation software,
63 the internal representation of interatomic interactions is in principle the same, or closely related.

64 Peptides containing ncAAs are ideally suited to NMR structural characterisation as they are small in size, and often contain
65 sidechain links that induce local structure (Hamada et al., 2010; Weber et al., 1991). Nevertheless, their structural
66 characterisation by NMR is often restricted to measurements of a limited set of NOEs, hydrogen bonds or backbone chemical
67 shifts to support specific geometries (Mendive-Tapia et al., 2015; Umstatter et al., 2020). Alternatively, NMR analysis is
68 omitted altogether in favour of lower resolution methods such as CD spectroscopy (De Araujo et al., 2022; Wu et al., 2017).
69 The dearth of high-resolution NMR structures in this class of molecules likely stems, at least in part, from the difficulty of
70 obtaining high-quality atomic representations of the ncAAs required for computer assisted structure determination.



71 Recently the handling of ncAAs and small molecules in CYANA was addressed by Yilmaz et al. (Yilmaz and Guntert, 2015),
72 who developed CYLIB, an algorithm that enables automated template generation for ncAAs and small molecules, provided
73 that a suitable input geometry is available (CIF or Mol2 file). Here, the quality of the input geometry is critical, as the algorithm
74 does not perform any optimisation of the structure. CYLIB has internal procedures for creating the appropriate branch structure
75 and ring-closures required by CYANA, however, practical aspects of working with ncAA-containing peptides and proteins
76 such as consistent atom naming and sidechain linkages are not addressed by CYLIB.

77 For the algorithms operating in Cartesian space, topology builders have been created which take simplified representations of
78 amino acids and other molecules as input and infer topological information based on a set of internal rules (Schmid et al.,
79 2012; Van Der Spoel et al., 2005; Wang et al., 2006). The challenge for incorporating ncAAs lies in the use of specific atom
80 names to infer bonded and/or non-bonded interactions, as well as assumptions regarding how individual amino acids in a
81 peptide chain are linked or terminated. While in principle almost any molecule can be represented, users are often unaware of
82 the assumptions that have been embedded in the codes, or how to incorporate new atom types into the associated files which
83 contain the parameter definitions needed by the builders to interpret them. Furthermore, common approaches to reduce
84 complexity such as using atom names to infer interaction type, work well when dealing with a small subset of chemical space
85 (such as the 20 canonical amino acids). However, linking atom names to specific interactions rapidly leads to a combinatorial
86 problem, and an explosion in terms, as new classes of interactions are introduced. The addition of a single new atom type can
87 often require the definition of hundreds of interactions to ensure compatibility with the existing framework.

88 In attempting to simplify the input for standard (common) cases, the authors of many programs have made the treatment of
89 non-standard cases progressively more challenging. We have set out to address this problem and provide a generic mechanism
90 to generate complete topological information for amino acids and related molecules of interest in a robust and reproducible
91 manner. This can be readily translated into inputs for various simulation and structure refinement packages. Our approach
92 leverages the capability of the Automated Topology Builder (ATB), a publicly accessible webserver that provides optimised
93 geometries, validated force field parameters and topology files for a range of popular molecular dynamics (MD) simulation
94 and structure refinement packages (Stroet et al., 2018; Koziara et al., 2014; Malde et al., 2011). The ATB is provided free to
95 academic users who can both download existing molecules and submit new structures. In the case of a new structure the user
96 can submit a set of Cartesian coordinates in, Protein Data Bank (PDB) format together with the net charge on the molecule.
97 Within the ATB, the geometry of the molecule is first optimised using quantum mechanical (QM) calculations. A topology is
98 then generated by combining the results with a set of empirical rules (Malde et al., 2011). A variety of formats can be used as
99 inputs, the most basic being an initial set of Cartesian coordinates in Protein Data Bank (PDB) format and the net charge on
100 the molecule.

101 Here, we describe an extension of the ATB that allows users to directly generate template files for a number of popular NMR
102 structure determination software packages. A template recognition approach has been developed that recognises ncAAs (rather
103 than “ligands”) as forming part of a peptide chain. This allows “building block” files compatible with different software to be
104 generated, including 3D geometries and templates with atom names that adhere to IUPAC standards. We further introduce a
105 general method for linking sidechains of amino acids in CYANA and CNS. The utility of the procedure is demonstrated by
106 generating models of a number of natural products containing complex ncAAs, as well as solving the structure of a sidechain
107 cyclised synthetic peptide. We further use this method to create a publicly accessible and expanding repository of amino acids
108 within the ATB. The database currently holds templates for the 20 genetically encoded amino acids with different termini,
109 common post-translationally modified amino acids, as well as the entire content of the SwissSidechain database (230 entries
110 in both D- and L- form of amino acids) which includes a wide range of ncAAs (Gfeller et al., 2013).

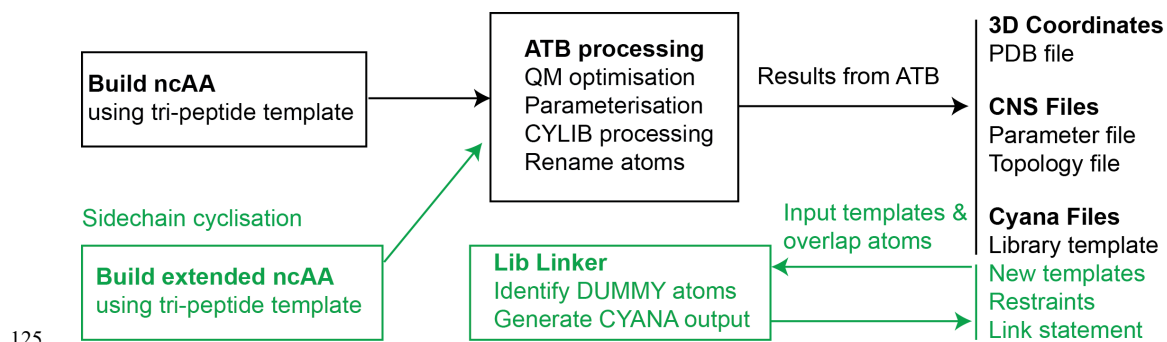


111 2 Methods

112 2.1 Outline

113 The overall workflow starts with the generation of 3D coordinates of a nCAA in a suitable amino acid template and format
114 (described below) for submission to the ATB, where geometries are optimised and a number of coordinate, parameter, topology
115 and template files are generated. The CYLIB algorithm is used internally by the ATB to generate the CYANA template files
116 from an intermediate Chemical Component Dictionary (CCD) Crystallographic Information Files (CIF) file. All atom names
117 are updated to adhere to IUPAC standards. The resulting outputs are made available to the user in CYANA or CNS format
118 and stored on the ATB server (Fig. 1). Optionally, where sidechain links are required for CYANA templates, these can be
119 generated using the CYANA Lib Linker (atb.uq.edu.au/cyana_linker). Each component of the workflow is described in detail
120 below.

121 While many pre-calculated nCAAs have been generated as described here and are stored within the ATB repository, the
122 protocol outlined below can be followed by a user to initiate the generation of parameters (templates, topologies, etc.) for any
123 new nCAA. New submissions of nCAAs to the ATB will be added to the existing database and thus the repository will become
124 progressively more complete.



125
126 **Figure 1.** Workflow for generating files for NMR structural calculation of peptides and proteins that contain nCAAs. The input
127 file must be generated using the described template format (Fig. 2) and used as input to the ATB webserver, which then
128 recognises the input as an amino acid and excises the necessary portion to generate the required input files for CYANA and
129 CNS-based structure determination softwares (shown in black on the right). To produce templates suitable for sidechain
130 linkage in CYANA (shown in green), an extended nCAA template containing atoms that extend beyond the linking bond is
131 built (Fig. 4). Two such templates are used as input to the CYANA Lib Linker interface of the ATB, where the user also
132 defines which atoms from the first template are present in the second template and vice versa. The “overlap” atoms of the
133 extension are changed to DUMMY atoms to produce a new CYANA template file (.lib file) and short upper distance restraints
134 are generated between sets of overlap atoms (.upl file). A link statement (to be added to the sequence file (.seq) in CYANA)
135 is also produced to remove the repulsion between the two atoms that are to be linked.

136

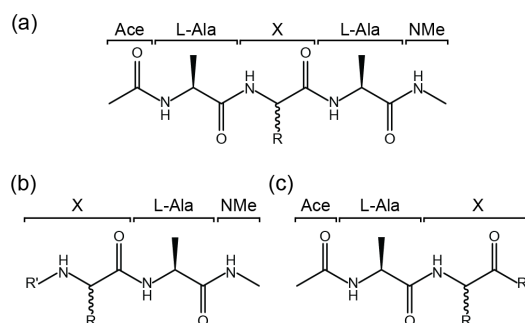


137 **2.2 Defining the group or amino acid of interest and submission to the ATB**

138 A core feature of the protocol is how the specific group or amino acid of interest is defined and automatically recognised. To
139 ensure an appropriate chemical environment for the parameterisation of the group of interest within a peptide chain, a series
140 of structural templates have been defined (Fig. 2). The general form of the peptide template is Ace-Ala-X-Ala-NMe, where
141 Ace is an N-terminal acetyl capping group, NMe is a C-terminal N-methyl amide capping group, and the central portion of the
142 structure, denoted X, is the chemical group or amino acid of interest (Fig. 2a). Similarly, structures of the form X-Ala-NMe or
143 Ace-Ala-X are recognised as representing an N-terminal or C-terminal residue respectively (Fig. 2b and 2c). This template
144 format allows the molecule to be identified as an amino acid, and the portion 'X' to be excised when generating the parameter
145 files. Each processed entry is associated with a unique molecule ID (MOLID) within the ATB database.

146 Submissions to the ATB require an unambiguous molecular representation that includes 3D coordinates (PDB format with all
147 hydrogen atoms present) along with the net charge. There is no requirement that the input geometry is optimised and thus may
148 contain clashes, non-ideal bond lengths etc. The stereochemistry, atom information and bonded connectivities, however, must
149 be specified correctly in the input file. The ATB submission page includes tools that will generate 3D coordinates from an
150 embedded 2D drawer (JSME (Bienfait and Ertl, 2013)) or a SMILES input. Existing entries within the ATB database can also
151 be loaded directly into the 2D drawing tool and modified. Documentation and ATB MOLIDs for the structural templates can
152 be found at github.com/ATB-UQ/CYANA-Examples. There is, in principle, no limit to the size of the group of interest.
153 However, for computational reasons, the default limit for density functional theory (DFT) calculations to be performed as part
154 of the parameterisation by the ATB is 50 atoms and the limit for semi-empirical QM processing is 500 atoms. Note that the
155 most significant difference between DFT and semi-empirical QM processing by the ATB is the atomic charge model, which
156 is not relevant for CYANA calculations as electrostatic interactions are not considered. While the group of interest is not
157 required to be a canonical amino acid, it must be able to be incorporated into a peptide chain via peptide bonds. The procedures
158 used by the ATB to parameterise molecules and the steps taken to validate these parameters have been described in detail
159 elsewhere (Malde et al., 2011; Stroet et al., 2018).

160 All amino acid entries are listed as such on the ATB (https://atb.uq.edu.au/index.py?tab=amino_acids), this includes entries
161 described herein, N- and C- terminal versions of all genetically encoded amino acids and all ncAAs contained within the
162 SwissSidechain DB (www.swissidechain.ch).



163

164 **Figure 2.** Template formats for submission of ncAAs to the ATB. The required templates for a) non-terminal residue, b) N-
165 terminal residue, and c) C-terminal residue in a peptide chain are shown. In each case R represents the sidechain of the ncAA
166 labelled 'X', and in the case of b) and c) R' represents possible modifications of the termini.



167 2.3 Renaming algorithm

168 When working with ncAAs, it is important to adhere to a consistent set of rules when naming the atoms on the sidechain.
169 Additionally, it is important for the workflow that atom names are unique. We therefore developed a tool to ensure both
170 requirements are fulfilled. In 1969, a set of rules were recommended by IUPAC-IUB for the representation of proteins and
171 peptides (Iupac, 1970). A section of these rules pertains to the labelling of the constituent atoms in amino acids. These rules
172 are adhered to in the presentation of NMR structures of proteins and peptides (Markley et al., 1998).

173 Non-hydrogen atoms present in the sidechain of amino acids are identified based on the lowest number of bonds that separate
174 them from the C α atom (determined from the connectivity information in the CYANA template or Mol2 format). The order of
175 atoms is indicated using the Greek alphabet (using corresponding Roman characters in files). The first atom connected to the
176 C α is C β , C γ is two bonds away and so on. While seemingly straightforward, complexities of this naming system arise for
177 branched amino acid sidechains, which are common in many ncAAs. In the event of a branch, where two atoms are the same
178 number of bonds away from the backbone and therefore are assigned the same Greek letter, chain priority is assigned based
179 on the Cahn–Ingold–Prelog priority (CIP) rules (Cahn and Ingold, 1951).

180 An algorithm was written to automatically name the sidechain of amino acids (github.com/ATB-UQ/fixnom). This involves
181 the following steps:

- 182 1. The library file containing atom coordinate and connectivity information is read as input (CYANA template format
183 or Mol2 format). The input is parsed and a reduced matrix is formed containing only the heavy atoms and their
184 connectivities (to heavy atoms).
- 185 2. The matrix is expanded to introduce dummy atoms to represent unsaturated bonds– these are required for the
186 implementation of the CIP rules. The C α atom is identified and the matrix reordered based on this information.
- 187 3. A connectivity matrix is created describing the distance (in bonds) between any two atoms. The atoms are ordered
188 based on their distance to the C α atom.
- 189 4. At this point the chains are initiated by considering all atoms connected to the C α position. All (heavy) atoms
190 connected to C α are given a position and a chain identifier. The position is simply the number of bonds away from
191 C α . The chain identifier follows the CIP rules. The chain with highest priority is assigned the lowest number and
192 additional chains are provided an incremented chain number based on priority. The priority of the chain is determined
193 as follows:
 - 194 a) The priority of that atom (atomic number)
 - 195 b) The priority of each attached atom one bond away. If all connected atoms have the same priority, then evaluate
196 all atoms two bonds away and so on.
 - 197 c) If the priority cannot be resolved by the above method, the atomic coordinates of the atoms are considered and
198 the R/S position of the branching atoms are used to assign priority. This last step is only required for cases of
199 tetrahedral atoms containing two identical chains (ignoring the preceding atom in the chain). If there are three
200 identical chains their identity is arbitrary by rotation (and no further action is taken). Practically this is
201 implemented by using the method of Cieplak and Wisniewski (Cieplak and L., 2001), where the determinant of
202 the 4×4 matrix formed by the atomic coordinates of the atoms (X, Y, Z, 1) in the stereocentre is used to evaluate
203 if an atom is clockwise or counterclockwise in position with respect to a geminal atom. The atom with the highest
204 priority is placed as the bottom row of the matrix, the two top rows are then occupied by the two atoms being
205 evaluated. As noted by Cieplak and Wisniewski the position or the identity of the third row does not affect the
206 handedness of the atoms being evaluated, and this atom can be “placed” at the position of the central atom. In



207 their implementation they do this for cases where the fourth atom is a hydrogen (often removed in databases) and
208 show that this is valid since this atom will always have a lower priority than the other atoms. Similarly, in our
209 case, we do not need to know what this atom is since it will always have a lower priority than the atom preceding
210 the branch point (according to the CIP elaboration by Markley et al. (Markley et al., 1998), the atom closest to
211 the $C\alpha$ atom always has the highest priority), thus we can simply fill the third row with the atomic coordinates
212 of the central atom.

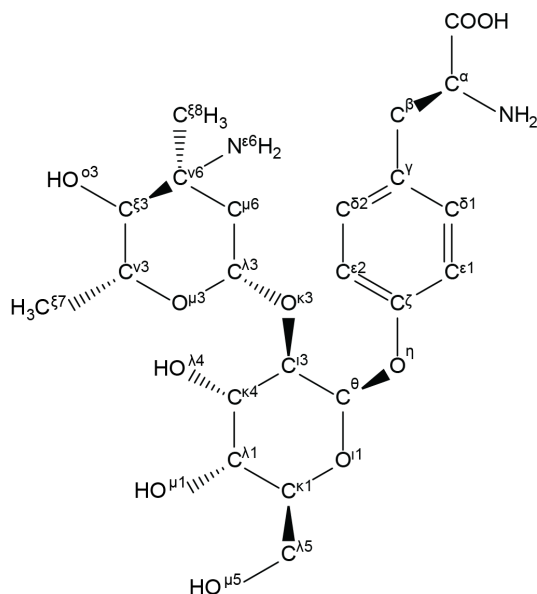
213 5. The above procedure assigns chain numbers to atoms that are one bond ahead of the atom being considered in an
214 incremental (one pass) approach. This requires that at each increment, the atoms that have already been assigned
215 priority (ahead of the atom being evaluated) are reordered based on their assigned priority in the previous step.

216 6. A situation can arise that an atom is evaluated twice if it joins two branches, as is the case for rings. In these cases,
217 the atom is given the lower chain ID between that already assigned to it, and what would be assigned to it had it not
218 been a joining atom.

219 Once all the chains have been generated, a procedure for re-evaluating identical chains (due to symmetry) in different parts of
220 the molecule is applied. Here the chain number is reordered based on the priority of the chain from which they branch. This is
221 not explicitly defined by the CIP rules but is required to ensure that if a branch contains two identical aromatic rings, the chain
222 numbers in each aromatic ring are consecutive (i.e., to avoid a ring having chain identifiers 1 and 3 or 2 and 4).

223 Once all chains have been created, the chain number and the distance from the C^α atom are used to generate the correct atom
224 names. Note, if a position exists that belongs to chain 1 and only one (heavy) atom occupies this position (number of bonds
225 away from C^α), this atom is not given a chain number, and only the Greek alphabet character corresponding to the bond position
226 is used. Finally, hydrogen (and pseudo) atoms are added. For cases where two hydrogens are attached to a tetrahedral carbon,
227 the stereochemistry of these is used to assign priority (using the same method as outlined above in step 4c).

228 This algorithm has been implemented in the ATB using an open-source standalone Perl script developed in-house
229 (github.com/ATB-UQ/fixnom) and is used to apply IUPAC atom naming to all amino acid building block outputs. This
230 algorithm is used internally in the ATB to rename atoms in all output files, including PDB files that can be used to define
231 atomic templates in NMR analysis software such as CCPN (available in v2 and planned for v3 (Skinner et al., 2016; Vranken
232 et al., 2005)). An example of a complex ncAA named by the algorithm is shown in Fig. 3.



233

234 **Figure 3.** Generation of template files following the IUPAC atom naming conventions for amino acids. As an example, the
235 atom names of the ncAA D-phenyl-Gly* (ATB MOLID 606467), present in vancomycin, were automatically generated using
236 an in-house algorithm. Greek letters are used to signify the number of bonds between a given atom and the C^α atom, with β
237 denoting one bond, γ denoting two bonds, continuing to ω denoting 14 bonds away. Chain priority is shown by the number
238 next to the Greek letter. In this diagram several important priority assignment examples can be seen. For example, at the branch
239 of C⁰ to C¹ and O¹, priority is assigned to the chain with the heavier atom, oxygen, hence the priorities reflect C³ and O¹.

240

241 2.4 ATB – CYLIB – CYANA pipeline

242 Due to the torsion angle dynamics algorithm used by CYANA to efficiently sample configurations, the parameters used to
243 describe amino acids are arranged according to a tree structure with the N-terminus at the base and the sidechains (and C-
244 terminus) as terminating branches (Guntert and Buchner, 2015). Because of the inherent complexities in representing
245 molecules in this manner, the ATB utilises the CYLIB application to produce a CYANA library file for amino acid building
246 blocks (Yilmaz and Guntert, 2015). This is achieved by excising the target group (X) from the templates outlined in Fig. 2 and
247 producing a CIF file containing the variables, amino acid backbone atoms (and atom names) to match amino acid structures
248 within the CCD (Westbrook et al., 2015). The resulting CCD compatible amino acid CIF file is passed to CYLIB with the
249 arguments required for sidechain, C- or N-terminus groups and the resulting files made available as downloadable files on the
250 ATB site. Note that in cases where sidechains or termini contain cyclic elements (within an amino acid), CYLIB also produces
251 a restraint macro to close the cycle which is also provided as a downloadable file.

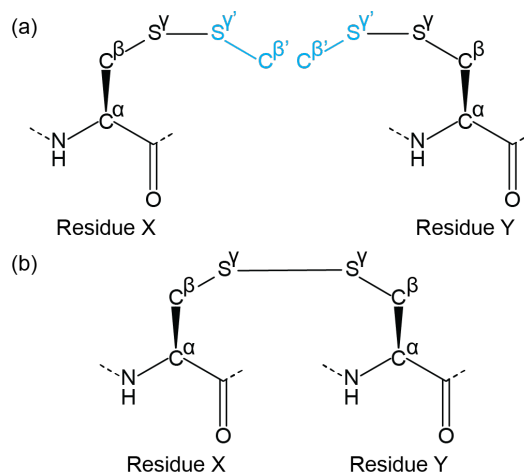
252



253 2.5 CYANA Lib Linker distance restraints

254 To enable amino acid sidechain and backbone cyclisation within CYANA, the ATB provides a tool to generate distance
255 restraints, linking statements and modified library files called CYANA Lib Linker (atb.uq.edu.au/cyana_linker).

256 The CYANA Lib Linker takes two input CYANA library files both of which contain the two atoms that will form the link
257 between the two amino acids plus at least a one atom extension beyond the linking bond (e.g. a disulfide bridge could be
258 formed by linking two sidechains of cysteines in the form $C^\alpha-C^\beta-S^\gamma-S^{\gamma'}-C^{\beta'}$) – see Fig. 4. To make the process generic for any
259 linkage, templates are required for each side of the link (in the disulfide bond example the same template is used twice as
260 input). For each template the following must then be defined: (1) “Residue index”, which is the residue number in the intended
261 peptide sequence; (2) “Linking bond”, defined by the two atoms involved in the sidechain linkage ($S^\gamma-S^{\gamma'}$); (3) “overlap atoms”
262 defined as atoms that exist in both templates (these atoms may have different names – in the disulfide bond example, this
263 would be a process of pairing atoms C^β of one template with the $C^{\beta'}$ of the other and S^γ of one with the $S^{\gamma'}$ of the other). Once
264 all of the above is satisfied, the algorithm edits the input template files by altering the atom type of the overlap atoms (and
265 attached hydrogens) in the library extension ($S^{\gamma'}-C^{\beta'}$) as “DUMMY” (excluding PSEUDO atoms which are not altered). An
266 upper distance is defined for each pair of corresponding overlap atoms with a limit of 0.04 Å (chosen arbitrarily, must be larger
267 than 0 and within experimental uncertainty) and assigned a weight of 10 (i.e. equivalent to 10 NOEs in CYANA). CYANA
268 Lib Linker also produces a link statement that removes the repulsion term between the bound atoms (between the two S^γ atoms
269 from each template) within CYANA. The link statement is included in the sequence file input of CYANA.



270

271 **Figure 4.** Example of a sidechain cyclisation using the Lib Linker for CYANA. Two extended ncAA template files are
272 generated representing either side of the linkage (extension shown in cyan), including the linking atom and neighbouring atoms
273 that define the chemical environment of the linked atoms. a) For a disulfide bond, the linked residues are the same on either
274 side and a single template can be selected as both input files for Lib Linker (Residue X and Y). Each side is, however,
275 defined as having different positions in the peptide chain. The connecting atoms (S^γ from each template) and overlap atoms (C^β and S^γ
276 overlap with $C^{\beta'}$ and $S^{\gamma'}$ of the other template, shown in black and cyan respectively) are provided as user defined input to Lib
277 Linker. b) Lib Linker generates two templates (identical in this case), where only the atoms from each side of the linking bond
278 are included. Restraint files are also generated that can be used directly for CYANA structure calculations.

279



280 **2.6 ATB-CNS Pipeline**

281 MD force fields have been incorporated into several NMR structure determination packages including XPLOR (currently
282 distributed as Xplor-NIH), ARIA and CNS (Schwieters et al., 2003; Schwieters et al., 2006; Bermejo and Schwieters, 2018).
283 For non-backbone-modified ncAAs lacking sidechain cyclisation, implementation into CNS forcefields for structure
284 calculations is straightforward. The ncAA is built and processed by the ATB as described above. To maintain compatibility
285 with the current protein force field used by CNS, the amino acid building block files produced by the ATB retain the standard
286 CNS atom types as well as the bonded and non-bonded parameters for atoms involving protein backbone definitions (N, HN,
287 CA, HA, CB, O). This allows standard linkage statements to be used for producing the molecular template. The new topology
288 and parameters can simply be added to the standard files without modification. New atom types are defined for all other atoms
289 in the ncAA allowing the new geometry and charges to be defined. In order to conserve the net charge when combining the
290 CNS protein backbone charges with the ATB charge model, any residual charge is simply added to the non-backbone atom
291 with the largest magnitude charge. Note that while the duplication of atom types is prevented between ATB outputs and the
292 standard CNS protein force field, atom types defined within separate ATB outputs are not unique. In cases where multiple and
293 distinct ATB parameterised groups are being used within a single peptide, manual resolution of atom types to ensure they are
294 unique may be required.

295 For sidechain-linked ncAAs, a manual addition of the cross link is required. As for the CYANA approach, generating building
296 blocks with overlapping atoms allows all aspects of the required geometry to be defined and topology and parameters directly
297 added to the existing forcefield. A specific linking statement that removes the extra atoms, modifies charges if required, and
298 adds the required bonds, angles and improper dihedrals (analogous to how disulfide bonds are implemented), can subsequently
299 be written by the user. Similarly, ncAAs that include backbone modifications, and thus cannot be modelled using the existing
300 standard linkage statements for creating peptide bonds, also require manual modification of the peptide bond linkage statement
301 to be incorporated and modelled correctly.

302

303 **3 Results**

304 **3.1 Disulfide bonds**

305 The handling of sidechain linkages described here is a departure from the default approach used in CYANA for linking
306 disulfide bonds. Currently, disulfide bonds are defined using a special template file called CYSS, where the template contains
307 all atoms up to the linking sulfur atom. A set of upper and lower, one- and two- bond distance limits are then imposed to
308 maintain an appropriate geometry around the introduced bond. One problem with this approach is that the χ_2 and χ_3 torsion
309 angles of disulfide bonds are not defined in the template and are therefore neither explicitly subject to the CYANA torsion
310 angle search, nor can they be easily constrained to specific values.

311 To validate the use of the new template (CYSX), we recalculated the structures of a number of disulfide-rich peptides resolved
312 in our group (e.g. 2KSL, 5LIC) using the new template (data not shown). In general, the recalculated structures are very similar
313 to those previously calculated using the CYSS template and distance restraints to define the disulfide bond. While it is beyond
314 the scope of this work to investigate in detail how the subtle differences between the two methods affect the quality of the
315 calculated structures, we can make some general observation about the process. We note that using the CYSX template instead
316 of the CYSS template in CYANA requires only trivial changes to the associated files. Specifically: (1) CYSS is replaced with



317 CYSX in the sequence file (the “link” statement which removes the repulsion between the connected sulfur atoms is used in
318 both approaches and requires no further changes); (2) CYSX is appended to the existing CYANA library; (3) the old restraints
319 for upper and lower distance limits to define the disulfide bond are removed and replaced with the new short upper distance
320 limits between “overlap atoms”; (4) after structure calculation the dummy atoms are removed and the CYSX name changed
321 to CYS. The last point can be achieved by adding four lines of commands to the end of the CYANA structure calculation script
322 (using existing CYANA functions). Warnings may occur due to a conflict between CYSX and CYSS in other restraints files
323 (i.e. dihedral angle restraints), but these can be ignored. The required template file for CYSX, the associated distance restraints
324 and CYANA commands for removing the DUMMY atoms and producing a PDB file suitable for subsequent analysis and
325 deposition to the PDB are provided in the supplementary section and our GitHub repository (github.com/ATB-UQ).

326

327 **3.2 ncAAs in complex natural products**

328 To demonstrate the capability and workflow, we have selected three examples of ncAA-containing peptides and generated the
329 additional library files required for structure calculation by CYANA. These peptides were selected to demonstrate particular
330 aspects of our pipeline and highlight potential practical applications. For each example, ncAA templates and restraints were
331 generated following the above procedure, and we demonstrate that these lead to chemically sound structures during
332 unrestrained structure calculation in CYANA.

333 Tyrocidine is a cyclic decapeptide antibiotic (Loll et al., 2014) (Fig. 5a) and contains one ncAA which is currently not present
334 in the standard CYANA library. The missing ornithine (at position 9) library file was generated by building the tripeptide:
335 Ace-Ala-Orn-Ala-NMe, in PyMOL (using the PyMOL Builder) and the resulting PDB file was saved and submitted to the
336 ATB, (MOLID 467880). The resulting ATB entry provides all necessary files for NMR structure calculation. The CYANA
337 template generated by the ATB was appended to the CYANA library. Existing procedures in CYANA were used to create a
338 backbone linkage to cyclise the peptide chain. Unrestrained CYANA calculations were performed, resulting in an ensemble
339 of chemically feasible structures (without NMR restraints). An example structure is shown in Fig. 5b.

340 Cyclosporine (Corbett et al., 2021) is an 11-residue backbone-cyclised peptide commonly used as an immunosuppressant to
341 treat rheumatoid arthritis and Crohn’s disease, and in the prevention of organ rejection in transplants. Cyclosporin contains six
342 ncAAs (Fig. 5c), five of which are not present in the standard CYANA library. The ncAAs are: D-alanine (residue 1), (4R)-4-
343 [(E)-2-butenyl]-4,N-dimethyl-L-threonine (MeBmt) (residue 5), α -aminobutyric acid (Abu) (residue 6) as well as three N-
344 methylated amino acids (sarcosine, residue 7; N-methyl-valine, residue 4; N-methyl leucine, residues 2, 3, 8, 10). Overall,
345 cyclosporine required five new library files to be defined and appended to the CYANA library (D-alanine can be generated
346 from alanine using the “library mirror” command in CYANA). A tripeptide template was generated for each and submitted to
347 the ATB. The ncAAs needed to describe cyclosporine are now available with the following MOLIDS: N-methyl leucine,
348 1175924; N-methyl-valine, 1175930; MeBmt, 1175933; Abu, 1175938; and sarcosine, 1175941. Feasible structures were
349 again obtained following unrestrained torsion angle dynamics simulations using CYANA (Fig. 5d).

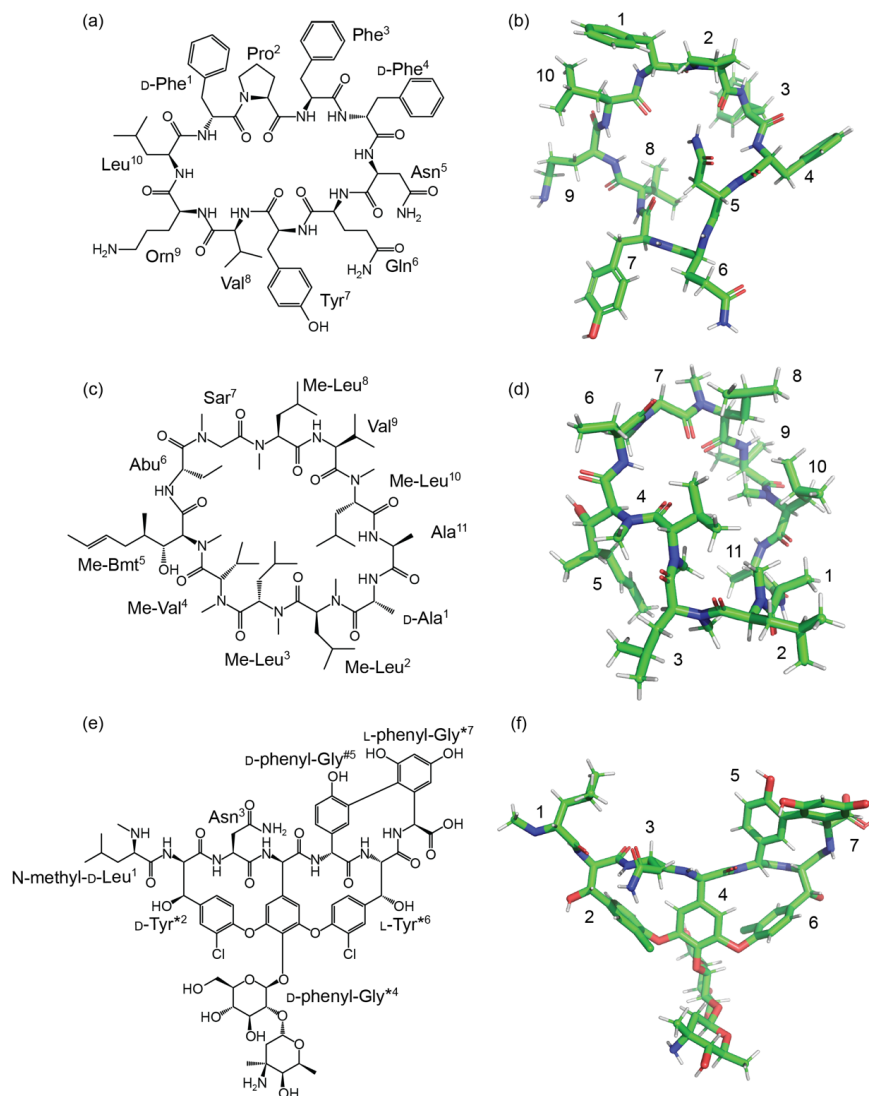
350 The final example, vancomycin, is a glycopeptide antibiotic that has been included on the WHO list of most essential medicines
351 (Schafer et al., 1996). It was, until recently, used as an antibiotic of last resort. This tricyclic heptapeptide is an example of a
352 non-ribosomal peptide. Five of the seven amino acids in vancomycin are involved in sidechain-sidechain branches. These
353 residues must be generated individually and subsequently processed through the CYANA Lib Linker algorithm to define



354 overlap atoms and distance restraint files. Additional distance restraint files are also required to close the sugar rings present
355 in residue four.

356 Three sidechain links are present in vancomycin, the first (link 1) between residues 2 and 4, the second (link 2) between
357 residues 4 and 6 and the third (link 3) between residues 5 and 7. Templates were generated with extensions beyond the linking
358 atom (including the aromatic rings on either side of the ether bond – MOLIDs: 1212202 (residue 2) and 1212203 (residue 4)
359 for link 1. The overlap atoms included the bonding atom (O^{n3} and $C^{\delta 1}$ in link 1 and $C^{\delta 2}$ and O^{n3} for link 2), as well as all atoms
360 one bond away from the linking atom ($C^{\zeta 3}$, $C^{\zeta 3}$, and $C^{\zeta 4}$ in residue 2; and $C^{\zeta 1}$, $C^{\zeta 3}$ and $C^{\gamma 1}$ in residue 4). For link 2, one additional
361 extended library file (MOLID: 1213034) was generated for residue 6, and the following overlap atom mapping was used: $C^{\gamma 2}$,
362 $C^{\epsilon 3}$ and $C^{\zeta 2}$ in residue 4; and $C^{\zeta 4}$, $C^{\zeta 3}$ and $C^{\zeta 3}$ in residue 6. Note that the modified residue 4 from link 1 is used as input when
363 creating link 2. For link 3, between residues 5 and 7, two additional library files were generated, MOLIDs 1212698 and
364 1212205 respectively. Link 3 exists between the bonding atoms $C^{\delta 1}$ and $C^{\gamma 3}$ with overlap atom mapping: $C^{\gamma 1}$, $C^{\epsilon 1}$, $C^{\zeta 3}$ and $C^{\zeta 4}$
365 in residue 5; and $C^{\epsilon 5}$, $C^{\epsilon 3}$, $C^{\delta 1}$ and C^{β} in residue 7. Finally, a template was produced for the N-terminal residue (MOLID
366 1212206). Using these template files, unrestricted torsion angle dynamics were performed in CYANA (Fig. 5f).

367 Vancomycin contains two sugar rings, and in CYANA each must be “closed” using an internal ring-closure process. CYLIB
368 automatically generates the required CYANA commands for closing rings, however, it currently cannot handle multiple ring-
369 closures, such as those for residue 4 of vancomycin. This required an additional ring-closure statement to be added manually.
370 An additional manual step involved adding a torsion angle that defines the angle between the two aromatic rings that are
371 directly fused (between residues 5 and 7).



372

373 **Figure 5.** Structures of natural products generated from un-restrained CYANA calculations using ncAA templates generated
374 by the ATB as described herein. (a) 2D and (b) 3D structure of tyrocidine containing a single ncAA (ornithine). (c) 2D and (d)
375 3D structure of cyclosporine. The sequence of cyclosporine contains five ncAA, N-methyl-leucine (Me-Leu), N-methyl-valine
376 (Me-Val), 4R)-4-[(E)-2-butenyl]-4-methyl-L-threonine (Me-Bmt), α -Aminobutyric acid (Abu), and sarcosine (Sar). (e) 2D and
377 (f) 3D structure of vancomycin. The amino acids of vancomycin have been abbreviated in some cases. D-Tyr*² is m-chloro- β -
378 hydroxy-D-Tyr. D-phenyl-Gly*⁴ is (2-[α -4-L-epi-vancosaminyll)- β -l-D-glucosyl)-D-phenyl-Gly. D-phenyl-Gly*⁵ is p-hydroxy-
379 D-phenyl-Gly. L-Tyr*⁶ is m-chloro- β -hydroxy-L-Tyr and L-phenyl-Gly*⁷ is m,m-dihydroxy-L-phenyl-Gly. The 3D structures
380 were generated by CYANA based on the templates obtained by the ATB – CYLIB – CYANA pipeline without experimental
381 restraints.

382



383 3.3 Practical application using an engineered peptide

384 Advances in peptide chemistry are rapidly expanding the chemical space accessible to high-throughput peptide synthesis
385 methods. Recently, a systematic approach was taken to explore sidechain stabilisation of a segment of the amyloid precursor
386 protein (APP; ¹NGYENPTYKFFE¹²) into the conformation it adopts when bound to the phosphotyrosine binding (PTB)
387 domain of Mint2 (Bartling et al., 2022). The method described herein was developed to solve the structure of four peptides
388 with different sidechain linkages (Bartling et al. under review). One of these peptides is sidechain-cyclised through ring closing
389 metathesis (RCM) between residues 7 and 11. The 3D structure of this peptide has been reported elsewhere (Bartling et al.
390 under review), and was solved using the methods described herein. The sidechain cyclisation following RCM was, however,
391 treated using a similar approach to that currently used in CYANA for fusing sulfur atoms in disulfide bonds (i.e. using truncated
392 template files). Here, we revisit this structure and perform the sidechain fusion using the new approach described here. We
393 further demonstrate how CNS templates, generate by the ATB, can be used for water refinement of the reported CYANA
394 structure.

395 First, Thr7 and Phe11 were replaced with the α -methyl-substituted olefin-bearing nAA named pentenyl alanine “PAL”
396 (MOLID 929126). This is an extended form of the PAL sidechain as shown in Fig. 6. The CYANA Lib Linker was then used
397 to modify the template (DUMMY and overlap atom definition) and to generate appropriate upper distance restraints. For the
398 terminal residues, templates corresponding to the N-terminal Asn and C-terminal Glu (MOLIDs 1162954 and 1159438) with
399 their respective amino and carboxy acid termini were also generated.

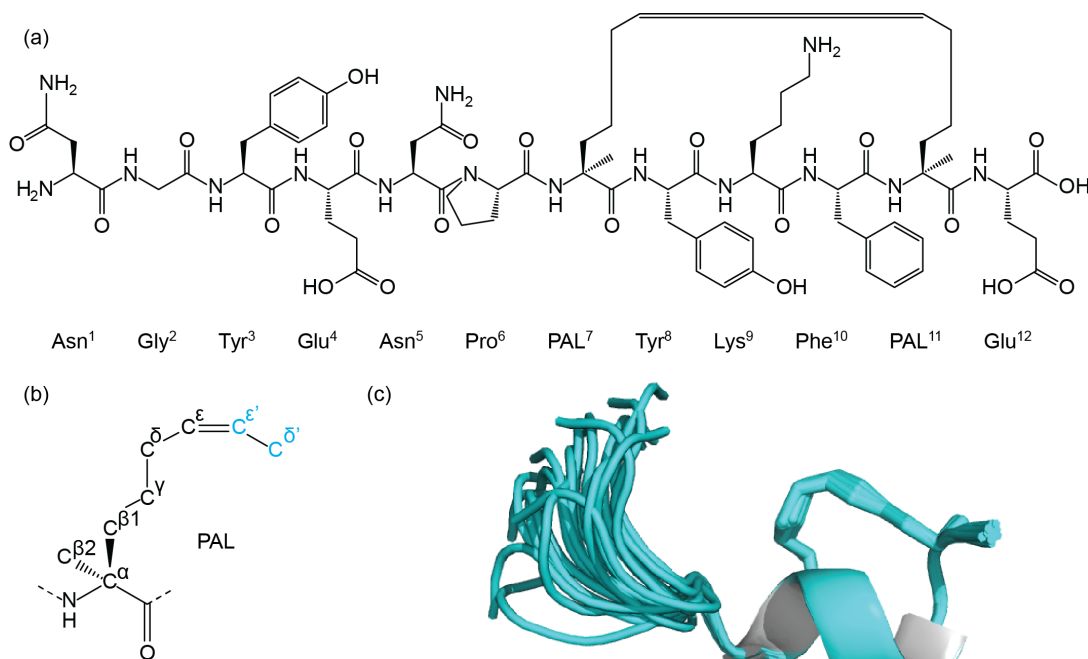
400 To assign the magnetic resonances of the peptide, CcpNmr Analysis 2.4.1 was employed (Vranken et al., 2005). The atomic
401 composition of individual nAAs in CcpNmr were defined by uploading the ATB-generated coordinate files into its molecule
402 library. Resonance assignments were obtained using a combination of 2D ¹H-¹H TOCSY, 2D ¹H-¹H NOESY and natural
403 abundance 2D ¹H-¹⁵N and ¹H-¹³C HSQC. The ¹H chemical shifts were calibrated with the reference to the water chemical shift
404 while ¹³C and ¹⁵N chemical shifts were calibrated indirectly with the reference to ¹H. Cross-peaks from the 2D ¹H-¹H NOESY
405 (mixing time of 350 ms) were manually picked to generate a list of interproton distance restraints. TALOS-N was used to
406 derive dihedral angle restraints. As TALOS-N does not recognize nAAs, we replaced the ATB generated residue codes for
407 the terminal residues with N and E respectively, in the TALOS-N input file. The angle restraint range was set to twice the
408 estimated standard deviation.

409 To perform calculations using CYANA the following files were prepared using the ATB, CcpNmr and TALOS-N: (i) a
410 sequence file listing the amino acid sequence of the peptide, (ii) a chemical shift file listing the chemical shifts of all assigned
411 atoms, (iii) a peak list of the 2D ¹H-¹H NOESY spectrum containing the chemical shifts and calibrated peak intensity (height
412 or volume) of each peak, (iv) an angle restraint file derived from TALOS-N, (v) a CYANA library file specific for the nAA
413 templates acquired from the ATB, and (vi) restraint files specific for the sidechain linkage generated by CYANA Lib Linker.
414 The sequence file and peak list (distance restraints) were directly exported from CcpNmr in the CYANA-compatible XEASY
415 format. The chemical shifts of protons in the peptide, were first exported from CcpNmr in BMRB format (any other formats
416 would omit the entries for any nAAs – a current CcpNmr limitation). The shift list in BMRB format was imported into
417 CYANA and pseudoatoms were added using an internal CYANA command. CYANA (v. 3.98.13) was then used to
418 automatically assign the peak list, extract distance restraints, and calculate 200 structures from which 20 structures with lowest
419 target function values were selected to represent the structure ensemble of the peptide.

420 The output of the CYANA structure calculation was used to perform water-refinement in CNS. The CNS topology and
421 parameter output for PAL was added to the standard forcefield and this residue included at both positions to be linked. A



422 linkage statement deleting the extra atoms was subsequently introduced to generate a complete molecular template file with
423 all the atoms and restraints required to define the geometry of the residue. Because of the additional backbone methyl group,
424 a custom peptide bond linkage statement was also constructed and used to create the residue links on either side of each ncAA
425 (both custom linkage statements are also available on github.com/ATB-UQ). Structures were calculated and minimised in
426 water using the experimental distance and angle restraints, resulting in a well-defined family of structures with excellent
427 geometry and no violations of the experimental data (Fig. 6c).



428

429 **Figure 6.** Figure showing the chemical structure of a synthetic peptide (a) that mimics the MINT-2 bound form of the amyloid
430 precursor protein (APP). The helical structure in the C-terminal tail of the peptide is stabilised through sidechain cyclisation
431 using ring closing metathesis (RCM). The RCM linkage was modelled using an extended pentenyl alanine (PAL) library entry
432 (b) generated by the ATB. The ATB generated library was further processed by the CYANA Lib Linker interface of the ATB
433 to generate a new template where the two overlap atoms C ϵ^* and C δ^* are converted to DUMMY atoms and used to constrain
434 the sidechain around the double bond. The 3D structure of the peptide was determined using CYANA and then further refined
435 in explicit water using CNS (c) with topology and parameter files generated by the ATB. The structure shows stabilisation of
436 the helical motif (cartoon) and the PAL sidechain is shown as sticks.

437

438 4 Discussion

439 A workflow has been developed that facilitates the structure determination and refinement of peptides and proteins that contain
440 ncAAs using popular NMR software, including CYANA and CNS. Methodology has been incorporated into the ATB server,
441 which allows arbitrary ncAAs to be built as required. New ncAAs are stored on the site and added to a database of entries
442 containing all of the necessary files for structure determination by NMR spectroscopy. The workflow introduces a number of



443 automated procedures to improve access to complex ncAAs without manual intervention. An algorithm has also been
444 developed that can automatically assign IUPAC atom names to ncAAs. It has already been highlighted that there have been
445 errors within the PDB when naming for diastereotopic atoms (Bottoms and Xu, 2008). The potential explosion in elaborate
446 sidechains that is possible when introducing ncAAs requires that care is taken when defining atoms, reporting and comparing
447 data.

448 We have also introduced into this workflow a stand-alone and general method for introducing sidechain-to-sidechain bonds in
449 CYANA. A graphical user interface has been incorporated into the ATB to facilitate this process using user supplied template
450 files (i.e. the files generated within the server). The new sidechain linking procedures uses overlap “DUMMY” atoms to
451 establish the missing bond in CYANA. This is a departure from the standard method currently used in CYANA for linking
452 sidechain atoms of pairs of cysteines. Tests of this procedure showed that the structures produced by both the new and old
453 methods were very similar. Importantly, the new approach allows for definition of the χ_2 and χ_3 torsion angles in cysteine
454 bridges, enabling these to be sampled by the torsion angle dynamics in CYANA, and for these angles to be defined where
455 experimental evidence is available (Armstrong et al., 2018; Ramanujam et al., 2019). In the “traditional” approach of defining
456 disulfide bonds, these torsion angles are not defined and disulfide bond geometries are instead sampled indirectly by altering
457 the χ_1 angles and the consequences of this on the imposed distance restraints and repulsion terms associated with the other
458 sidechain atoms. While the result appears to be similar using the two methods in the absence of χ_2 and χ_3 restraints, the new
459 method will certainly be favourable where such data are available. More generally, this procedure will allow all sidechain links
460 to include definition of inter-residue torsion angles in CYANA.

461 The utility of the workflow was demonstrated by building structural models of several natural products in CYANA. In most
462 cases only a single ncAA is present in a peptide or protein chain, due to incorporation of the ncAA via recombinant expression
463 methods, via peptide synthesis or chemical modification of specific AAs in the produced peptide or protein. This was captured
464 by modelling tyrocidine, a short peptide incorporating a single ncAA. In such cases, the method is very robust and little manual
465 intervention is required (beyond adding the additional templates to the CYANA library). Cases involving multiple ncAA,
466 especially those containing backbone N-methylation, are very challenging to prepare manually, however, even cyclosporin
467 poses no issues in our pipeline beyond the once-off generation of the template files. The most challenging example explored,
468 the antibiotic vancomycin, is an extreme case involving multiple interlocked ring structures formed by linking ncAAs. CYLIB
469 currently cannot automatically process groups with multiple rings, such as residue 4 in vancomycin which has two sugar
470 moieties branching from an aromatic ring. Thus, some manual steps are required to add the extra ring closure statements.
471 Similarly, the central residue 4 of vancomycin is sidechain-linked to two different residues. This requires serial iterations of
472 the CYANA Lib Linker interface to generate the required files – i.e. first templates for residues 2 and 4 are used to generate
473 modified templates, the modified template for residue 4 is then further modified when submitted as the linking residue to
474 residue 6. Even this challenging molecule could be processed using the tools developed.

475 The results described above demonstrated that the ncAA templates are compatible with existing libraries and yield high quality
476 structures. We also showed how one might use the templates in conjunction with experimental data. The APP-derived peptide
477 that includes an RCM-cyclised sidechain, was used to highlight some practical considerations associated with the pipeline
478 described in this work. Some manual processing is still required to use our templates with the popular CcpNmr software. While
479 we were able to create a working solution for v2 of CcpNmr, incorporation of ncAAs within a peptide chain in v3 is currently
480 not feasible. Other analysis packages such as POKY do not require templates and may be more suited for use with ncAAs (Lee
481 et al., 2021). We also noted that additional automation was required to appropriately import the output of CcpNmr into



482 CYANA. Although we have a working solution, it is likely that future development of CcpNmr will address these problems
483 that exist when working with ncAAs.

484 The final test involved the refinement of the APP peptide in water using the CNS software. This currently involves manual
485 creation of the linkage statements. This is relatively straightforward involving a simple modification of the standard peptide
486 bond and/or the standard disulfide bond definitions. Importantly, because the geometries are already defined in the building
487 blocks and all parameters required are included in the ATB output, the statement simply has to define which atoms are linked
488 and infer other geometrical constraints such as the planarity of the double-bond. Once written, these linkage statements can be
489 used for any variant of a given ncAA. For example, the peptide backbone link used for the “PAL” residue can be applied to
490 any residue type in which the H^α proton has been replaced with a methyl group.

491 The workflow presented here was developed to cater to the rapid growth in peptide and protein engineering in recent years,
492 examples including directed evolution mRNA display methods to generate macrocyclic peptide ligands of target receptors
493 (Goto and Suga, 2021), modified amino acids for structural studies (Mekkattu Tharayil et al., 2022) and high-throughput
494 chemical synthesis in drug development (Bartling et al., 2022). These developments have increased the interest in solving
495 structures of ncAA-containing peptides leading us to develop the described method. The protocols and tools we have developed
496 are designed to be general and to interface with a range of software. This said, it is likely that not all ncAAs that can be
497 envisaged will be able to be handled with the existing workflow. Nevertheless, the solution we presented in this work not only
498 addresses a pressing need in the cases of ncAAs but also provides a general framework that can be used to improve the
499 description of AAs in structure refinement more broadly.

500

501 **5 Conclusions**

502 Peptides containing ncAAs encompass a large pool of biologically active molecules with many potential industrial, agricultural
503 and pharmaceutical uses. The significant structural diversity found in these peptides presents significant challenges for NMR
504 spectroscopists when applying existing structure determination tools. This problem requires the development of a set of tools
505 that can automatically generate molecular representations that have suitable chemical properties. We have here provided a
506 solution to this problem in the form of an extension to the Automated Topology Builder, which now can produce template
507 files compatible with most commonly used NMR structure calculation software. We have also ensured that the ncAA templates
508 generated adhere to IUPAC naming conventions based on the Cahn-Ingold-Prelog priority rules. This extends the utility of
509 NMR structure calculation to complex natural products, synthetic peptides, and complex, natural and unnatural, post-
510 translational modifications.

511 **Code and Data availability**

512 The ATB server is available publicly at (atb.uq.edu.au). The code used in this work is available via GitHub as cited in the
513 manuscript (github.com/ATB-UQ).

514 **Author contribution**

515 SK, MS, MM and AEM conceived of the project and designed elements of the pipeline. MS implemented the code in the ATB
516 with input from AEM. SK validated the implementation with input from MM. MS, TL and MM wrote code used in the ATB
517 workflow. SK, YKYC, ACC, XJ, KJR and MM produced the models presented and analysed the NMR data. CORB and KS



518 provided the synthetic APP peptide. PG contributed to CYANA template generation and sidechain linking method. SK wrote
519 the first draft of the manuscript and all authors contributed to the final version.

520 Acknowledgements

521 This project was supported by funding from the Australian Research Council (DP DP190101177 and DP220103028 to MM.
522 DP220100896 to AEM and MS), the National Health and Medical Research Council (NHMRC APP1162597 to MM) and the
523 University of Queensland (Postgraduate Research Scholarship to SK, Research Stimulus fellowship to YC and Development
524 Fellowship to MM) and in part by the Austrian Science Fund (FWF) (Project P36101-B) to AC.

525

526 References

- 527 Aebersold, R., Agar, J. N., Amster, I. J., Baker, M. S., Bertozzi, C. R., Boja, E. S., Costello, C. E., Cravatt, B. F., Fenselau,
528 C., Garcia, B. A., Ge, Y., Gunawardena, J., Hendrickson, R. C., Hergenrother, P. J., Huber, C. G., Ivanov, A. R., Jensen, O.
529 N., Jewett, M. C., Kelleher, N. L., Kiessling, L. L., Krogan, N. J., Larsen, M. R., Loo, J. A., Ogorzalek Loo, R. R., Lundberg,
530 E., MacCoss, M. J., Mallick, P., Mootha, V. K., Mrksich, M., Muir, T. W., Patrie, S. M., Pesavento, J. J., Pitteri, S. J.,
531 Rodriguez, H., Saghatelian, A., Sandoval, W., Schluter, H., Sechi, S., Slavoff, S. A., Smith, L. M., Snyder, M. P., Thomas, P.
532 M., Uhlen, M., Van Eyk, J. E., Vidal, M., Walt, D. R., White, F. M., Williams, E. R., Wohlschlagler, T., Wysocki, V. H., Yates,
533 N. A., Young, N. L., and Zhang, B.: How many human proteoforms are there?, *Nat Chem Biol*, 14, 206-214,
534 10.1038/nchembio.2576, 2018.
- 535 Allain, F., Mareuil, F., Menager, H., Nilges, M., and Bardiaux, B.: ARIAweb: a server for automated NMR structure
536 calculation, *Nucleic Acids Res*, 48, W41-W47, 10.1093/nar/gkaa362, 2020.
- 537 Armstrong, D. A., Kaas, Q., and Rosengren, K. J.: Prediction of disulfide dihedral angles using chemical shifts, *Chem Sci*, 9,
538 6548-6556, 10.1039/c8sc01423j, 2018.
- 539 Barber, K. W. and Rinehart, J.: The ABCs of PTMs, *Nat Chem Biol*, 14, 188-192, 10.1038/nchembio.2572, 2018.
- 540 Bartling, C. R. O., Alexopoulou, F., Kuschert, S., Chin, Y., Jia, X., Sereikaite, V., Jain, P., Özcelik, D., Jensen, T. M., Nygaard,
541 M. M., Harpsøe, K., Gloriam, D. E., Mobli, M., and Strømgaard, K.: Systematic Scanning of Peptide Cyclization Structure
542 Generates Cell-Permeable Inhibitors of Protein-Protein Interactions, *Angewandte Chemie International Edition*, Submitted,
543 2022.
- 544 Bermejo, G. A. and Schwieters, C. D.: Protein Structure Elucidation from NMR Data with the Program Xplor-NIH, *Methods*
545 *Mol Biol*, 1688, 311-340, 10.1007/978-1-4939-7386-6_14, 2018.
- 546 Bienfait, B. and Ertl, P.: JSME: a free molecule editor in JavaScript, *J Cheminform*, 5, 24, 10.1186/1758-2946-5-24, 2013.
- 547 Bondalapati, S., Jbara, M., and Brik, A.: Expanding the chemical toolbox for the synthesis of large and uniquely modified
548 proteins, *Nat Chem*, 8, 407-418, 10.1038/nchem.2476, 2016.
- 549 Bottoms, C. A. and Xu, D.: Wanted: unique names for unique atom positions. PDB-wide analysis of diastereotopic atom names
550 of small molecules containing diphosphate, *BMC Bioinformatics*, 9 Suppl 9, S16, 10.1186/1471-2105-9-S9-S16, 2008.
- 551 Brunger, A. T., Adams, P. D., Clore, G. M., DeLano, W. L., Gros, P., Grosse-Kunstleve, R. W., Jiang, J. S., Kuszewski, J.,
552 Nilges, M., Pannu, N. S., Read, R. J., Rice, L. M., Simonson, T., and Warren, G. L.: Crystallography & NMR system: A new
553 software suite for macromolecular structure determination, *Acta Crystallogr D Biol Crystallogr*, 54, 905-921,
554 10.1107/s0907444998003254, 1998.
- 555 Bullwinkle, T., Lazazzera, B., and Ibba, M.: Quality control and infiltration of translation by amino acids outside of the genetic
556 code, *Annu Rev Genet*, 48, 149-166, 10.1146/annurev-genet-120213-092101, 2014.



- 557 Caboche, S., Pupin, M., Leclere, V., Fontaine, A., Jacques, P., and Kucherov, G.: NORINE: a database of nonribosomal
558 peptides, *Nucleic Acids Res*, 36, D326-331, 10.1093/nar/gkm792, 2008.
- 559 Cahn, R. S. and Ingold, C. K.: 131. Specification of configuration about tetravalent asymmetric atoms, *Journal of the*
560 *Chemical Society (Resumed)*, 612-622, 10.1039/jr9510000612, 1951.
- 561 Chuh, K. N., Batt, A. R., and Pratt, M. R.: Chemical Methods for Encoding and Decoding of Posttranslational Modifications,
562 *Cell Chem Biol*, 23, 86-107, 10.1016/j.chembiol.2015.11.006, 2016.
- 563 Cieplak, T. and L., W. J.: New Effective Algorithm for the Unambiguous Identification of the Stereochemical Characteristics
564 of Compounds During Their Registration in Databases, 10.3390/61100915, 2001.
- 565 Coin, I.: Application of non-canonical crosslinking amino acids to study protein-protein interactions in live cells, *Curr Opin*
566 *Chem Biol*, 46, 156-163, 10.1016/j.cbpa.2018.07.019, 2018.
- 567 Conibear, A. C.: Deciphering protein post-translational modifications using chemical biology tools, *Nature Reviews*
568 *Chemistry*, 4, 674-695, 10.1038/s41570-020-00223-8, 2020.
- 569 Corbett, K. M., Ford, L., Warren, D. B., Pouton, C. W., and Chalmers, D. K.: Cyclosporin Structure and Permeability: From
570 A to Z and Beyond, *J Med Chem*, 64, 13131-13151, 10.1021/acs.jmedchem.1c00580, 2021.
- 571 Craveur, P., Rebehmed, J., and de Brevern, A. G.: PTM-SD: a database of structurally resolved and annotated posttranslational
572 modifications in proteins, *Database (Oxford)*, 2014, bau041, 10.1093/database/bau041, 2014.
- 573 de Araujo, A. D., Lim, J., Wu, K. C., Hoang, H. N., Nguyen, H. T., and Fairlie, D. P.: Landscaping macrocyclic peptides:
574 stapling hDM2-binding peptides for helicity, protein affinity, proteolytic stability and cell uptake, *RSC Chem Biol*, 3, 895-
575 904, 10.1039/d1cb00231g, 2022.
- 576 Gfeller, D., Michielin, O., and Zoete, V.: SwissSidechain: a molecular and structural database of non-natural sidechains,
577 *Nucleic Acids Res*, 41, D327-332, 10.1093/nar/gks991, 2013.
- 578 Goodrich, A. C. and Frueh, D. P.: A nuclear magnetic resonance method for probing molecular influences of substrate loading
579 in nonribosomal peptide synthetase carrier proteins, *Biochemistry*, 54, 1154-1156, 10.1021/bi501433r, 2015.
- 580 Goto, Y. and Suga, H.: The RAPID Platform for the Discovery of Pseudo-Natural Macrocyclic Peptides, *Acc Chem Res*, 54,
581 3604-3617, 10.1021/acs.accounts.1c00391, 2021.
- 582 Guntert, P. and Buchner, L.: Combined automated NOE assignment and structure calculation with CYANA, *J Biomol NMR*,
583 62, 453-471, 10.1007/s10858-015-9924-9, 2015.
- 584 Guntert, P., Mumenthaler, C., and Wuthrich, K.: Torsion angle dynamics for NMR structure calculation with the new program
585 DYANA, *J Mol Biol*, 273, 283-298, 10.1006/jmbi.1997.1284, 1997.
- 586 Hamada, T., Matsunaga, S., Fujiwara, M., Fujita, K., Hirota, H., Schmucki, R., Guntert, P., and Fusetani, N.: Solution structure
587 of polytheonamide B, a highly cytotoxic nonribosomal polypeptide from marine sponge, *J Am Chem Soc*, 132, 12941-12945,
588 10.1021/ja104616z, 2010.
- 589 Hoesl, M. G. and Budisa, N.: In vivo incorporation of multiple noncanonical amino acids into proteins, *Angew Chem Int Ed*
590 *Engl*, 50, 2896-2902, 10.1002/anie.201005680, 2011.
- 591 Hoyt, E. A., Cal, P. M. S. D., Oliveira, B. L., and Bernardes, G. J. L.: Contemporary approaches to site-selective protein
592 modification, *Nature Reviews Chemistry*, 3, 147-171, 10.1038/s41570-019-0079-1, 2019.
- 593 Huang, Y., Russell, W. K., Wan, W., Pai, P. J., Russell, D. H., and Liu, W.: A convenient method for genetic incorporation of
594 multiple noncanonical amino acids into one protein in *Escherichia coli*, *Mol Biosyst*, 6, 683-686, 10.1039/b920120c, 2010.
- 595 IUPAC: IUPAC-IUB Commission on Biochemical Nomenclature. Abbreviations and symbols for the description of the
596 conformation of polypeptide chains. Tentative rules (1969), *Biochemistry*, 9, 3471-3479, 10.1021/bi00820a001, 1970.
- 597 Johnson, J. A., Lu, Y. Y., Van Deventer, J. A., and Tirrell, D. A.: Residue-specific incorporation of non-canonical amino acids
598 into proteins: recent developments and applications, *Curr Opin Chem Biol*, 14, 774-780, 10.1016/j.cbpa.2010.09.013, 2010.



- 599 Koziara, K. B., Stroet, M., Malde, A. K., and Mark, A. E.: Testing and validation of the Automated Topology Builder (ATB)
600 version 2.0: prediction of hydration free enthalpies, *J Comput Aided Mol Des*, 28, 221-233, 10.1007/s10822-014-9713-7,
601 2014.
- 602 Lee, W., Rahimi, M., Lee, Y., and Chiu, A.: POKY: a software suite for multidimensional NMR and 3D structure calculation
603 of biomolecules, *Bioinformatics*, 37, 3041-3042, 10.1093/bioinformatics/btab180, 2021.
- 604 Link, A. J., Mock, M. L., and Tirrell, D. A.: Non-canonical amino acids in protein engineering, *Curr Opin Biotechnol*, 14, 603-
605 609, 10.1016/j.copbio.2003.10.011, 2003.
- 606 Liu, C. C. and Schultz, P. G.: Adding new chemistries to the genetic code, *Annu Rev Biochem*, 79, 413-444,
607 10.1146/annurev.biochem.052308.105824, 2010.
- 608 Loll, P. J., Upton, E. C., Nahoum, V., Economou, N. J., and Cocklin, S.: The high resolution structure of tyrocidine A reveals
609 an amphipathic dimer, *Biochim Biophys Acta*, 1838, 1199-1207, 10.1016/j.bbamem.2014.01.033, 2014.
- 610 Mal, T. K., Bagby, S., and Ikura, M.: Protein structure calculation from NMR data, *Methods Mol Biol*, 173, 267-283,
611 10.1385/1-59259-184-1:267, 2002.
- 612 Malde, A. K., Zuo, L., Breeze, M., Stroet, M., Poger, D., Nair, P. C., Oostenbrink, C., and Mark, A. E.: An Automated Force
613 Field Topology Builder (ATB) and Repository: Version 1.0, *J Chem Theory Comput*, 7, 4026-4037, 10.1021/ct200196m,
614 2011.
- 615 Marahiel, M. A.: Working outside the protein-synthesis rules: insights into non-ribosomal peptide synthesis, *J Pept Sci*, 15,
616 799-807, 10.1002/psc.1183, 2009.
- 617 Mareuil, F., Malliavin, T. E., Nilges, M., and Bardiaux, B.: Improved reliability, accuracy and quality in automated NMR
618 structure calculation with ARIA, *J Biomol NMR*, 62, 425-438, 10.1007/s10858-015-9928-5, 2015.
- 619 Markley, J. L., Bax, A., Arata, Y., Hilbers, C. W., Kaptein, R., Sykes, B. D., Wright, P. E., and Wuthrich, K.:
620 Recommendations for the presentation of NMR structures of proteins and nucleic acids. IUPAC-IUBMB-IUPAB Inter-Union
621 Task Group on the Standardization of Data Bases of Protein and Nucleic Acid Structures Determined by NMR Spectroscopy,
622 *J Biomol NMR*, 12, 1-23, 10.1023/a:1008290618449, 1998.
- 623 Martínez-Núñez, M. A. and López, V. E. L. y.: Nonribosomal peptides synthetases and their applications in industry,
624 *Sustainable Chemical Processes*, 4, 13, 10.1186/s40508-016-0057-6, 2016.
- 625 Mekkattu Tharayil, S., Mahawaththa, M. C., Feintuch, A., Maleckis, A., Ullrich, S., Morewood, R., Maxwell, M. J., Huber,
626 T., Nitsche, C., Goldfarb, D., and Otting, G.: Site-selective generation of lanthanoid binding sites on proteins using 4-fluoro-
627 2,6-dicyanopyridine, *Magnetic Resonance*, 3, 169-182, 10.5194/mr-3-169-2022, 2022.
- 628 Mendive-Tapia, L., Preciado, S., Garcia, J., Ramon, R., Kielland, N., Albericio, F., and Lavilla, R.: New peptide architectures
629 through C-H activation stapling between tryptophan-phenylalanine/tyrosine residues, *Nat Commun*, 6, 7160,
630 10.1038/ncomms8160, 2015.
- 631 Noren, C. J., Anthonycahill, S. J., Griffith, M. C., and Schultz, P. G.: A General-Method for Site-Specific Incorporation of
632 Unnatural Amino-Acids into Proteins, *Science*, 244, 182-188, DOI 10.1126/science.2649980, 1989.
- 633 Ramanujam, V., Shen, Y., Ying, J., and Mobli, M.: Residual Dipolar Couplings for Resolving Cysteine Bridges in Disulfide-
634 Rich Peptides, *Front Chem*, 7, 889, 10.3389/fchem.2019.00889, 2019.
- 635 Schafer, M., Schneider, T. R., and Sheldrick, G. M.: Crystal structure of vancomycin, *Structure*, 4, 1509-1515, 10.1016/s0969-
636 2126(96)00156-6, 1996.
- 637 Schmid, N., Christ, C. D., Christen, M., Eichenberger, A. P., and van Gunsteren, W. F.: Architecture, implementation and
638 parallelisation of the GROMOS software for biomolecular simulation, *Computer Physics Communications*, 183, 890-903,
639 10.1016/j.cpc.2011.12.014, 2012.
- 640 Schwieters, C. D., Kuszewski, J. J., and Clore, G. M.: Using Xplor-NIH for NMR molecular structure determination, *Progress*
641 *in Nuclear Magnetic Resonance Spectroscopy*, 48, 47-62, 10.1016/j.pnmrs.2005.10.001, 2006.



- 642 Schwieters, C. D., Kuszewski, J. J., Tjandra, N., and Clore, G. M.: The Xplor-NIH NMR molecular structure determination
643 package, *J Magn Reson*, 160, 65-73, 10.1016/s1090-7807(02)00014-9, 2003.
- 644 Skinner, S. P., Fogh, R. H., Boucher, W., Ragan, T. J., Mureddu, L. G., and Vuister, G. W.: CcpNmr AnalysisAssign: a flexible
645 platform for integrated NMR analysis, *J Biomol NMR*, 66, 111-124, 10.1007/s10858-016-0060-y, 2016.
- 646 Strieker, M., Tanovic, A., and Marahiel, M. A.: Nonribosomal peptide synthetases: structures and dynamics, *Curr Opin Struct
647 Biol*, 20, 234-240, 10.1016/j.sbi.2010.01.009, 2010.
- 648 Stroet, M., Caron, B., Visscher, K. M., Geerke, D. P., Malde, A. K., and Mark, A. E.: Automated Topology Builder Version
649 3.0: Prediction of Solvation Free Enthalpies in Water and Hexane, *J Chem Theory Comput*, 14, 5834-5845,
650 10.1021/acs.jctc.8b00768, 2018.
- 651 Tharp, J. M., Krahn, N., Varshney, U., and Soll, D.: Hijacking Translation Initiation for Synthetic Biology, *Chembiochem*, 21,
652 1387-1396, 10.1002/cbic.202000017, 2020.
- 653 Thompson, R. E. and Muir, T. W.: Chemoenzymatic Semisynthesis of Proteins, *Chem Rev*, 120, 3051-3126,
654 10.1021/acs.chemrev.9b00450, 2020.
- 655 Umstatter, F., Domhan, C., Hertlein, T., Ohlsen, K., Muhlberg, E., Kleist, C., Zimmermann, S., Beijer, B., Klika, K. D.,
656 Haberkorn, U., Mier, W., and Uhl, P.: Vancomycin Resistance Is Overcome by Conjugation of Polycationic Peptides, *Angew
657 Chem Int Ed Engl*, 59, 8823-8827, 10.1002/anie.202002727, 2020.
- 658 Van Der Spoel, D., Lindahl, E., Hess, B., Groenhof, G., Mark, A. E., and Berendsen, H. J.: GROMACS: fast, flexible, and
659 free, *J Comput Chem*, 26, 1701-1718, 10.1002/jcc.20291, 2005.
- 660 Vranken, W. F., Boucher, W., Stevens, T. J., Fogh, R. H., Pajon, A., Llinas, M., Ulrich, E. L., Markley, J. L., Ionides, J., and
661 Laue, E. D.: The CCPN data model for NMR spectroscopy: development of a software pipeline, *Proteins*, 59, 687-696,
662 10.1002/prot.20449, 2005.
- 663 Walsh, C. T., Garneau-Tsodikova, S., and Gatto, G. J., Jr.: Protein posttranslational modifications: the chemistry of proteome
664 diversifications, *Angew Chem Int Ed Engl*, 44, 7342-7372, 10.1002/anie.200501023, 2005.
- 665 Wang, J., Wang, W., Kollman, P. A., and Case, D. A.: Automatic atom type and bond type perception in molecular mechanical
666 calculations, *J Mol Graph Model*, 25, 247-260, 10.1016/j.jm gm.2005.12.005, 2006.
- 667 Wang, L., Brock, A., Herberich, B., and Schultz, P. G.: Expanding the genetic code of *Escherichia coli*, *Science*, 292, 498-
668 500, 10.1126/science.1060077, 2001.
- 669 Wang, T., Liang, C., An, Y., Xiao, S., Xu, H., Zheng, M., Liu, L., Wang, G., and Nie, L.: Engineering the Translational
670 Machinery for Biotechnology Applications, *Mol Biotechnol*, 62, 219-227, 10.1007/s12033-020-00246-y, 2020.
- 671 Weber, C., Wider, G., von Freyberg, B., Traber, R., Braun, W., Widmer, H., and Wuthrich, K.: The NMR structure of
672 cyclosporin A bound to cyclophilin in aqueous solution, *Biochemistry*, 30, 6563-6574, 10.1021/bi00240a029, 1991.
- 673 Westbrook, J. D., Shao, C., Feng, Z., Zhuravleva, M., Velankar, S., and Young, J.: The chemical component dictionary:
674 complete descriptions of constituent molecules in experimentally determined 3D macromolecules in the Protein Data Bank,
675 *Bioinformatics*, 31, 1274-1278, 10.1093/bioinformatics/btu789, 2015.
- 676 Wu, Y., Li, Y. H., Li, X., Zou, Y., Liao, H. L., Liu, L., Chen, Y. G., Bierer, D., and Hu, H. G.: A novel peptide stapling strategy
677 enables the retention of ring-closing amino acid side chains for the Wnt/beta-catenin signalling pathway, *Chem Sci*, 8, 7368-
678 7373, 10.1039/c7sc02420g, 2017.
- 679 Yilmaz, E. M. and Guntert, P.: NMR structure calculation for all small molecule ligands and non-standard residues from the
680 PDB Chemical Component Dictionary, *J Biomol NMR*, 63, 21-37, 10.1007/s10858-015-9959-y, 2015.
- 681 Zou, H., Li, L., Zhang, T., Shi, M., Zhang, N., Huang, J., and Xian, M.: Biosynthesis and biotechnological application of non-
682 canonical amino acids: Complex and unclear, *Biotechnol Adv*, 36, 1917-1927, 10.1016/j.biotechadv.2018.07.008, 2018.
- 683