

1 **Supplementary Information for**  
2 **A unified language model bridging *de novo* and fragment-based**  
3 **3D molecule design**

4 Han Wang<sup>1,2†</sup>, Guanglong Sun<sup>2†</sup>, Bowen Zhang<sup>3†</sup>, Yang Wang<sup>2†</sup>, Bin Xi<sup>1,2†</sup>, Minjian  
5 Yang<sup>1,2</sup>, Chuanyu Liu<sup>3</sup>, Yuyang Ge<sup>2</sup>, Fan Fan<sup>2</sup>, Wei Feng<sup>2</sup>, Yanhao Zhu<sup>2</sup>, Yang Xiao<sup>2</sup>,  
6 Xiaojian Xu<sup>4</sup>, Yuji Wang<sup>5</sup>, Zhenming Liu<sup>1,6\*</sup>, Daohua Jiang<sup>3\*</sup>, Huting Wang<sup>2\*</sup>, Wenbiao  
7 Zhou<sup>2\*</sup>, and Bo Huang<sup>1,2,5\*</sup>

8  
9 <sup>1</sup>State Key Laboratory of Natural and Biomimetic Drugs, School of Pharmaceutical Sciences, Peking  
10 University, Beijing, China.

11 <sup>2</sup>Beijing StoneWise Technology Co Ltd., Haidian Street #15, Beijing, 100080, China.

12 <sup>3</sup>Laboratory of Soft Matter Physics, Institute of Physics, Chinese Academy of Sciences, Beijing, China.

13 <sup>4</sup>Beijing Neurosurgical Institute, Capital Medical University, Beijing, China.

14 <sup>5</sup>College of Pharmaceutical Sciences, Capital Medical University, Beijing 100069, P. R. China.

15 <sup>6</sup>Key Laboratory of Xinjiang Endemic Phytomedicine Resources Ministry of Education; School of  
16 Pharmacy, Shihezi University, Shihezi 832003, Xinjiang, China.

17  
18  
19 **Contents**

20 1. Supplementary Results..... 2  
21 1.1. Ablation Studies of UniLingo3DMol..... 2  
22 1.2. Screening Pipelines for CBL-B Compound Design..... 3  
23 2. Supplementary Methods ..... 5  
24 2.1. UniLingo3DMol Development ..... 5  
25 2.2. Molecular Retrieval-Augmented Generation using UniLingo3DMol ..... 7  
26 2.3. Molde Evaluation ..... 10  
27 2.4. Compound Synthesis..... 13  
28 3. Supplementary Figures S1-S4..... 16  
29 4. Supplementary Tables S1-S9..... 21  
30 5. Supplementary Algorithms S1-S6 ..... 35  
31 Reference ..... 41

32  
33

†These authors contributed equally to this work.

\*Corresponding author(s). E-mail(s): bohuang\_011@163.com; zhouwenbiao@stonewise.cn;

wanghuting@stonewise.cn; jiangdh@iphy.ac.cn; zmliu@bjmu.edu.cn

## 34 1. Supplementary Results

### 35 1.1. Ablation Studies of UniLingo3DMol

#### 36 1.1.1. Effectiveness of Multi-stage Training

37 To investigate the impact of the pre-training, post-training, and fine-tuning stages on  
38 UniLingo3DMol's performance, we conducted ablation studies comparing its results with  
39 and without either pre-training or post-training on the *de novo* design scenario. The  
40 corresponding experimental results were presented in [Supplementary Table S3](#). As shown  
41 in [Supplementary Table S3](#), the ablation studies revealed that:

42 **Omitting the pre-training stage almost degraded all metrics.** Specifically, the %  
43 ECFP\_TS > 0.5 decreased from 74% to 8%, the min-in-place GlideSP score decreased  
44 from -7.1 to -5.0, and the NCI recovery rate decreased from 33.2% to 6.2%. These results  
45 indicated that the pre-training stage, by leveraging the diverse chemical space within its  
46 dataset, provides superior initialization parameters for ligand generation.

47 **Omitting the post-training stage weakened UniLingo3DMol's ability to generate**  
48 **drug-like molecules.** This was evidenced by a decrease in the number of drug-like  
49 molecules generated, from 81,835 to 58,583, representing a reduction of approximately  
50 30%. In addition, the quality of binding modes within pockets was also weakened, with the  
51 min-in-place GlideSP score decreased from -7.1 to -6.5. Moreover, the ability of  
52 reproducing known active compounds was also weakened, as evidenced by % ECFP\_TS >  
53 0.5 decreased from 74% to 18%. These results suggest that employing a large-scale dataset,  
54 consisting of pocket-ligand complexes with force-filed accuracy, is essential to enhance  
55 the performance of UniLingo3DMol.

56 **Omitting the fine-tuning stage improved several docking metrics**, including min-in-  
57 place GlideSP score, redocking GlideSP score, and IMP. However, these metrics are highly  
58 dependent on the force field parameters within Schrödinger's Glide module. Importantly,  
59 the fine-tuning stage utilized a dataset derived from experimentally validated pocket-ligand  
60 complexes, and some complexes in this set exhibited lower min-in-place scores.  
61 Consequently, when UniLingo3DMol was fine-tuned on this dataset, the metrics degraded.  
62 Nevertheless, UniLingo3DMol's capability to generate drug-like molecules increased, as  
63 evidenced by an increase in the drug-like molecule percentage from 73.7% to 80.7%. Due  
64 to our pursuit of generating more drug-like molecules, the fine-tuning stage is necessary  
65 for UniLingo3DMol.

66 The above results clearly highlight the critical role of each stage in the UniLingo3DMol  
67 training process. Specifically, the effectiveness of UniLingo3DMol was diminished when  
68 any stage was omitted. This finding underscores the importance of all three training stages.

### 69 **1.1.2. Effectiveness of Multi-task Training**

70 Ablation studies were performed on the *de novo* design scenario by including or excluding  
71 Task1, with results presented in [Supplementary Table S4](#).

72 [Supplementary Table S4](#) shows that removing Task1 degraded UniLingo3DMol 's  
73 performance across multiple metrics. These metrics encompassed molecular conformation  
74 quality, binding mode accuracy against target pockets, and the reproducibility of known  
75 active compounds. Specifically, the metric % ECFP\_TS > 0.5 decreased from 74% to 71%,  
76 and the min-in-place GlideSP score worsened from -7.1 to -6.9. These results highlight the  
77 critical role of incorporating Task1 into UniLingo3DMol's training scheme.

78 The benefit of Task1 stems from its objective: training UniLingo3DMol to predict NCI  
79 sites and anchor sites from a given protein pocket. This capability enables UniLingo3DMol  
80 to identify key pocket residues, such as amino acids and skeletal residues in active sites.  
81 Crucially, this residue identification directs UniLingo3DMol to generate ligands that  
82 preferentially occupy target-binding sites during molecular design.

## 83 **1.2. Screening Pipelines for CBL-B Compound Design**

### 84 **1.2.1. Screening Pipeline for Novel Scaffold Discovery**

85 The primary aim of the initial screening phase was the identification of novel scaffolds.  
86 Using the proposed UniLingo3DMol model, we generated molecules based on two retained  
87 fragments (Regions A and C; [Figure 2c](#)). In this round, the UniLingo3DMol model was  
88 used to sample 3 million times. After filtering out outputs that violated the DSMILES  
89 syntax as well as redundant SMILES entries, we obtained 697,204 unique molecules. This  
90 resulting set comprises a diverse CBL-B compound library, particularly exhibiting high  
91 variability in the B and D regions.

92 We performed virtual screening of the compound library using a customized pipeline  
93 ([Figure 3b](#)). First, molecules were filtered based on key chemical properties: QED > 0.3,  
94 SAS < 5, and number of rotatable bonds < 10. This initial filtering yielded a refined library  
95 of 186,434 unique molecules. Second, these compounds underwent molecular min-in-place  
96 docking into the binding site of the target protein (PDB: 8GCY) using Schrödinger.  
97 Subsequently, conformational filtering was applied using the TED toolkit<sup>1</sup> to retain  
98 molecules with a torsional energy score below 1, resulting in a library of 12,143 unique  
99 molecules. Third, a binding-mode-based filter was implemented to further reduce the  
100 library. Specifically, ProLIF was employed to retain only molecules forming hydrogen  
101 bond with residue F263 of the pocket, which yielded a final library of 3,356 molecules.

102 Subsequently, we computed the CSK scaffolds for all molecules in the final compound  
103 library. To prioritize novelty, we filtered out molecules that shared the same CSK scaffold  
104 as the co-crystallized ligand of protein 8GCY. The remaining molecules were then

105 clustered based on their CSK scaffolds, meaning that molecules with identical scaffolds  
106 were assigned to the same cluster. This process resulted in 344 clusters. Within each cluster,  
107 molecules were ranked in ascending order according to their min-in-place GlideSP score.  
108 Finally, the top-ranked molecule from each cluster was selected as its representative.  
109 Additionally, the visualizations for the entire set of representative molecules are provided  
110 in [Supplementary Data \(Figshare repository\)](#).

### 111 **1.2.2. Screening Pipeline for Lead Compound Optimization**

112 In contrast to the initial screening round, the second round employed the UniLingo3DMol  
113 model for the generation of molecules, specifically retaining Regions A, B, and D ([Figure](#)  
114 [2d](#)). In this round, we also employed the UniLingo3DMol model to sample 3 million times.  
115 Due to the use of a larger retained fragment compared to the first round, the number of  
116 generated molecules with identical SMILES increased. After filtering out outputs that  
117 violated the DSMILES syntax and removing redundant SMILES entries, we obtained  
118 467,920 unique molecules.

119 A virtual screening workflow ([Figure 4a](#)), analogous to that used in the first round, was  
120 subsequently applied, incorporating three distinct filter types: chemical property,  
121 molecular conformation, and binding mode. After the chemical property filter, 263,574  
122 molecules were retained, while the molecular conformation filter yielded 64,885 molecules.  
123 Notably, the binding-mode-based filter criteria were distinct in this round. Molecules were  
124 retained exclusively if they formed a hydrogen bond interaction with F263 while  
125 simultaneously establishing a salt bridge with E268. This final filter resulted in a library of  
126 1,331 molecules.

127 Finally, we computed the CSK scaffold for each molecule in this library and clustered  
128 molecules sharing identical scaffolds, resulting in 91 clusters. Molecules within each  
129 cluster were ranked in ascending order based on their min-in-place Glide SP score. The  
130 top-ranked molecule from each cluster was selected as its representative. Visualizations of  
131 all representative molecules are provided in [Supplementary Data \(Figshare repository\)](#).

132

## 133 2. Supplementary Methods

### 134 2.1. UniLingo3DMol Development

#### 135 2.1.1. Training algorithms for UniLingo3DMol

136 UniLingo3DMol implemented a hierarchical three-stage training schema that includes the  
137 pre-training, post-training, and fine-tuning stages.

138 In the pre-training stage ([Supplementary Algorithm S1](#)), UniLingo3DMol underwent  
139 training in ligand reconstruction using perturbed ligand data  $D_l^\dagger$  and DSMILES-formatted  
140 ligand data  $D_l$ . The pre-training stage was conducted for 500 epochs, with UniLingo3DMol  
141 autoregressively predicting three key components through dedicated generation heads:  
142 ligand token head, ligand pointer head, and ligand position head. The composite loss  $\mathcal{L}_{PT}$   
143 combined the token prediction loss ( $\mathcal{L}_{\text{token}}$ ), pointer loss ( $\mathcal{L}_{\text{ptr}}$ ), and positional loss ( $\mathcal{L}_{\text{pos}}$ ),  
144 with parameters updated via AdamW optimization after each epoch.

145 In the post-training ([Supplementary Algorithm S2](#)), UniLingo3DMol was initialized  
146 with pre-trained parameters to leverage transfer learning benefits, a common practice in  
147 deep learning. UniLingo3DMol interleaved three distinct tasks over 200 epochs on the  
148 pocket-ligand complex data generated through force field. The first task (Task1) predicted  
149 NCI sites and anchor sites from pockets, the second task (Task2) conducted unconstrained  
150 ligand generation, and the third task (Task3) performed constrained ligand generation with  
151 NCI/anchor awareness. Each epoch iterated through all three tasks sequentially, calculating  
152 task-specific losses ( $\mathcal{L}_\tau$ ). Crucially, parameter updates occurred after each individual task  
153 execution rather than per epoch, enabling focused gradient propagation for each objective.  
154 Both stages employed data shuffling and AdamW optimization, with the multi-task phase  
155 demonstrating a curriculum learning strategy through its alternating constraint-enabled and  
156 constraint-free generation tasks.

157 In the fine-tuning stage ([Supplementary Algorithm S2](#)), UniLingo3DMol was initialized  
158 with post-trained parameters and fine-tuned on the target dataset, which consisted of  
159 protein-ligand complex data with experimental accuracy, using the same three tasks as in  
160 the post-training stage. To prevent catastrophic forgetting and improve UniLingo3DMol's  
161 generalization ability, the fine-tuning stage was only carried out for 10 epochs.

162 The definitions of abovementioned losses are as follows:

- 163 • **Token Loss** ( $\mathcal{L}_{\text{token}}$ ): It measures the discrepancy between the predicted and ground-  
164 truth DSMILES-formatted molecular tokens.
- 165 • **Pointer Loss** ( $\mathcal{L}_{\text{ptr}}$ ): It measures the discrepancy between the predicted and ground-  
166 truth tokens' connection relationships,
- 167 • **Coordinate Loss** ( $\mathcal{L}_{\text{pos}}$ ): It consists of global and local coordinates losses:

- 168 ○  $\mathcal{L}_r^{\text{local}}$ : It measures the error between the predicted and ground-truth radial distances.
- 169 ○  $\mathcal{L}_\theta^{\text{local}}$ : It measures the discrepancy between the predicted and ground-truth bond
- 170 angles.
- 171 ○  $\mathcal{L}_\phi^{\text{local}}$ : It evaluates the difference between the predicted and ground-truth dihedral
- 172 angles.
- 173 ○  $\mathcal{L}_{\text{pos}}^{\text{global}}$ : It evaluates the difference between the predicted and ground-truth atomic
- 174 global coordinates.
- 175 • **NCI Loss** ( $\mathcal{L}_{\text{NCI}}$ ): It measures the discrepancy between the predicted and ground-truth
- 176 NCI sites.
- 177 • **Anchor Loss** ( $\mathcal{L}_{\text{Anchor}}$ ): It measures the discrepancy between the predicted and ground-
- 178 truth anchor sites.

179 Here,  $\mathcal{L}_{\text{NCI}}$  and  $\mathcal{L}_{\text{Anchor}}$  were calculated using binary cross-entropy loss functions, and all  
 180 other losses were calculated using the cross-entropy loss functions.

### 181 2.1.2. Hyperparameters of UniLingo3DMol

182 The three-stage UniLingo3DMol training experiment was performed on a distributed  
 183 computing cluster with Nvidia H20s ([Supplementary Table S8](#)). All stages consistently  
 184 employed the AdamW optimizer with identical learning rates and weight decay  
 185 regularization.

### 186 2.1.3. Sequence Augmentation Algorithms

187 The [Supplementary Algorithm S3](#) transforms an original fragment sequence  $FS_{\text{raw}}$  into a  
 188 randomized yet structurally valid sequence  $FS_{\text{trans}}$  by iteratively selecting and appending  
 189 fragments while adhering to connection constraints. Initially, the length  $n$  of  $FS_{\text{raw}}$  was  
 190 determined, and  $FS_{\text{trans}}$  was initialized as an empty list. In each iteration, if  $FS_{\text{trans}}$  is empty,  
 191 a fragment  $f$  is randomly chosen from  $FS_{\text{raw}}$ , added to  $FS_{\text{trans}}$ , and removed from  $FS_{\text{raw}}$ .  
 192 For subsequent steps, the algorithm identifies connection constraints satisfied by the  
 193 current  $FS_{\text{trans}}$ , randomly selects a compatible fragment  $f$  from  $FS_{\text{raw}}$ , transforms  $f$  into  
 194  $FS_{\text{trans}}$  using an asterisk constraint, appends  $f_{\text{trans}}$  to  $FS_{\text{trans}}$ , and removes  $f$  from  $FS_{\text{raw}}$ ,  
 195 decrementing  $n$  until all fragments are processed. Finally,  $FS_{\text{trans}}$  undergoes a legality  
 196 check to ensure compliance with structural rules, ensuring the output is both stochastically  
 197 augmented and valid.

198 The [Supplementary Algorithm S4](#) transforms an original fragment sequence ( $FS_{\text{raw}}$ ) into  
 199 a new sequence ( $FS_{\text{trans}}$ ) by retaining a specified number of fragments ( $n_r$ ) while modifying  
 200 others. First, it checks if retention is unnecessary ( $n_r \geq FS_{\text{raw}}$ ) and returns an empty  
 201 sequence. Otherwise, it initializes  $FS_{\text{trans}}$ , randomly selects a starting index, and iteratively  
 202 processes fragments in two phases:

- 203 • **Preservation phase:** Fragments are randomly selected based on unsatisfied connection  
204 constraints and added to  $FS_{\text{trans}}$  until  $n_r$  fragments are retained.
- 205 • **Conversion phase:** Remaining fragments are modified using an asterisk constraint  
206 based on satisfied connection requirements. At each step, the selected fragment is  
207 removed from  $FS_{\text{raw}}$ , and the counters are updated. The algorithm ends when all  
208 fragments are processed, followed by a legality check to validate the structural  
209 consistency of  $FS_{\text{trans}}$ . Randomized selection and constraint-driven retention ensure  
210 diversity while preserving critical sequence relationships.

## 211 2.2. Molecular Retrieval-Augmented Generation using UniLingo3DMol

212 We enhance the generative pipeline by incorporating molecular fragment retrieval, which  
213 allows the model to access a library of 3D interaction pairs between amino acids and  
214 molecular fragments ([Supplementary Figure S4](#)). Each entry in this library includes the 3D  
215 structures of an amino acid and a corresponding molecular fragment, annotated with  
216 quantum mechanically accurate interaction energies computed using the machine-learned  
217 force field AIMNet2<sup>2</sup>, as well as intuitive interaction types defined by geometric criteria  
218 via ProLIF<sup>3</sup>. All interaction pairs are derived from known protein-ligand complexes (Level  
219 0 from RComplex database), ensuring that the generation process is grounded in  
220 chemically and biologically relevant knowledge, thereby minimizing the risk of invalid or  
221 unrealistic molecular designs.

222 During retrieval, the library is queried using protein residue indices and desired  
223 interaction types—optionally focusing on backbone interactions, where all types of amino  
224 acids will be screened—to identify fragments whose amino acid partners have  
225 conformations resemble the query residue and are compatible with the protein pocket.  
226 Candidate fragments are selected based on energy and geometry criteria and then used as  
227 retained fragment for molecule generation.

### 228 2.2.1. Retrieval Library Construction

229 The retrieval library is constructed by using RComplex Level 0 data ([Methods Section 4.6](#)).  
230 For each complex, the binding pocket is defined as all residues with at least one atom within  
231 6 Å of the ligand. Ligands are decomposed into fragments by cleaving non-cyclic, non-  
232 conjugated, or non-terminal bonds. Any fragment consisting of only one heavy atom is  
233 merged into its most adjacent fragment, and hydrogens are added at the cleavage points.

234 For each amino acid residue in the binding pocket, RDKit<sup>4</sup> is used to add hydrogens to  
235 the backbone nitrogen (N) and alpha carbon (C $\alpha$ ) to neutralize unwanted charges. ProLIF<sup>3</sup>  
236 is then employed to detect interactions between the amino acid residue and each ligand  
237 fragment. For each amino acid residue–ligand fragment pair, energy minimization is  
238 performed using the GAFF2<sup>5</sup> force field and OpenMM while constraining the amino acid

239 residue conformation. The 3D conformations of the amino acid residue and ligand fragment  
240 are recorded, along with the following features: the interaction fingerprints (IFP) from  
241 ProLIF; the RMSD of ligand fragment heavy atoms before and after minimization ( $\text{min}_{\text{rmsd}}$ );  
242 and the non-covalent interaction energies (NCIE) computed with the AimNet2 force field,  
243 both before ( $\text{NCIE}_{\text{ori}}$ ) and after minimization ( $\text{NCIE}_{\text{min}}$ ).

244 After extracting all interaction pairs, they are grouped by amino acid type. For each type,  
245 an arbitrary entry is selected as a reference, and the remaining entries are aligned to it with  
246 respect to the backbone atoms (N, C $\alpha$ , C) as the reference. The resulting transformation is  
247 applied to the corresponding fragment to preserve complex geometry. A distance matrix  
248 based on side-chain heavy atoms is used to cluster the amino acid residuals into conformers,  
249 and each conformer's ID and centroid status are recorded.

250 The final library entry includes: the amino acid residue structure in PDB format; the  
251 original and minimized fragment structures in SDF format; the amino acid conformer ID;  
252 an indicator of whether the amino acid is the centroid of its cluster; RMSD of ligand  
253 fragment heavy atoms between conformation before and after minimization; NCIE before  
254 and after minimization; IFPs between the amino acid residue and ligand fragment in the  
255 original conformation; and an indicator of whether the IFP involves backbone atoms.

### 256 **2.2.2. Retrieval Method**

257 The retrieval process is query-based, using an IFP site defined by a protein residue index  
258 (chain ID and residue ID) and the desired IFP type between the target residue and the  
259 retrieved fragment. Users may optionally specify that the IFP should involve the residue's  
260 backbone; in such cases, the amino acid type is ignored, and data from all residue types are  
261 considered. The method compares the query residue conformation to those in the database  
262 to identify entries with similar amino acids conformers, while ensuring that the retrieved  
263 fragment does not clash with the protein pocket via a docking score.

264 The retrieval proceeds as follows. First, the dataset corresponding to the specified residue  
265 type is selected; if backbone IFPs are targeted, the dataset includes all residue types and is  
266 filtered for backbone interactions. The query residue conformation is compared to all  
267 cluster centroids by calculating the RMSD of matching heavy atoms after backbone  
268 alignment. Conformer clusters with the closest-matching centroids are selected. From the  
269 selected entries, those with  $\text{NCIE}_{\text{ori}} < 0$  are retained, as are those with  $\text{NCIE}_{\text{ori}} \geq 0$  that  
270 decrease to  $\text{NCIE}_{\text{min}} < 0$  after minimal adjustment (constrained minimization  $\text{RMSD} < 2$   
271 Å). Further filtering is applied based on the target IFP type if specified. The candidate  
272 amino acid is aligned to the query using backbone atoms, and the fitting RMSD is recorded  
273 with a precision of 0.05 Å. For full-residue matching, alignment is refined using all  
274 matching heavy atoms; for backbone-only cases, only backbone heavy atoms are used. The  
275 resulting transformation is applied to the fragment.

276 Potential growth sites are identified and adjusted. These sites are initially defined at  
277 original bond cleavage points. For each site, a pseudo-atom is placed 1.5 Å along the vector  
278 from the heavy atom to its connected hydrogen. The minimum distance from this pseudo-  
279 atom to all heavy atoms in the protein pocket and ligand (excluding the site itself and  
280 neighbors) is computed. If no original cleavage sites exist, other carbon atoms with  
281 attached hydrogens are considered as potential sites. Sites with a minimum distance  $\geq 2.5$   
282 Å are retained. If more than three qualify, the top three are selected based on the largest  
283 minimum distance for original cuts, or to maximize spatial diversity for new sites. If no  
284 site qualifies, the primary site is used. The maximum of the minimum distances across  
285 selected sites (max\_min\_space) is recorded and rounded to 0.1 Å. Docking scores are  
286 computed using Smina<sup>6</sup> (Vinardo scoring) to ensure no steric clashes with the protein  
287 pocket. Negative scores are rounded to the nearest multiple of 2 kcal/mol and adjusted by  
288 subtracting 0.001 to avoid zero values. The positioning direction angle is calculated as the  
289 angle between the vector from the query residue's center of mass (COM) to the fragment's  
290 COM and a reference vector from the residue COM to the pocket COM (or a reference  
291 ligand COM if provided). This angle is rounded to the nearest even degree.

292 Growth direction angles are computed for each selected growth vector (from the heavy  
293 atom to its connected hydrogen at the growth site). The angle between this vector and a  
294 reference growth vector is calculated, and the smallest angle is selected as the best growth  
295 angle (rounded to the nearest even degree). The reference vector depends on the fragment's  
296 distance to the reference COM: if the distance is less than 2 Å, the reference vector points  
297 from the reference COM to the growth site; otherwise, it points from the growth site to the  
298 reference COM.

299 Candidates are filtered to exclude those with position angles  $> 120^\circ$ , best growth angles  $>$   
300  $120^\circ$ , docking scores  $\geq 1.0$ , more than 30 heavy atoms, or no valid growth sites. Sulfur-  
301 containing fragments may be optionally excluded to avoid excessive amount of sulfonyl  
302 group. Duplicates are removed based on canonical SMILES (ignoring isotopes). The  
303 remaining fragments are ranked by docking score (ascending), fitting RMSD (ascending),  
304 position angle (ascending), number of IFPs (descending), and reference NCIE (ascending,  
305 using  $NCIE_{\min}$  if adjusted).

### 306 **2.2.3. RAG with Molecular Generation**

307 To evaluate the RAG strategy, we applied it to the DUD-E dataset ([Supplementary Table](#)  
308 [S7](#)). For each target, ProLIF was used to identify key IFP patterns (hydrogen bonds,  $\pi$ -  
309 stacking, and salt bridges) in the cocrystal structure. Protein residues involved in these  
310 interactions were treated as query sites. For each site, fragments were retrieved using the  
311 method described above. Due to resource constraints, a single fragment was selected as the

312 initial seed for molecular generation. This selection involved an additional ranking across  
313 all retrieved fragments based on docking score and the total number of IFPs formed with  
314 the protein pocket, with the top-ranked fragment chosen for subsequent generation.

## 315 **2.3. Molde Evaluation**

### 316 **2.3.1. Evaluation Sets for *De Novo* Design**

317 For the comprehensive evaluation of *de novo* molecular design strategies, the DUD-E  
318 dataset served as the foundational resource. This dataset provided 102 unique protein  
319 targets and 22,886 active compounds whose activities were previously reported in the  
320 literature. A refinement during data preparation involved the substitution of one protein  
321 structure: PDB ID 2H7L was found to be obsolete in the PDB, and as per the guidance of  
322 the RCSB PDB, the structure 4TRJ was consequently utilized.

### 323 **2.3.2. Evaluation Sets for Fragment-retained Design**

324 Three evaluation sets were constructed from the DUD-E dataset to assess model  
325 performance in the fragment-retained scenario. The core methodology involved identifying  
326 high-frequency fragments derived from co-crystallized ligands and subsequently  
327 classifying active molecules based on their presence. Specifically, both co-crystallized  
328 ligands (derived from DUD-E crystal structures) and known active molecules were  
329 systematically fragmented by breaking the non-ring single bonds. This process generated  
330 a set of molecular fragments for each molecule, and we subsequently conducted a  
331 frequency analysis of these fragments. For each target, we identified the fragments found  
332 in the co-crystallized ligands that also appeared in the active compounds set. Fragments  
333 that appeared more than five times in the active molecule set for a given target were labeled  
334 as conserved fragments. This threshold was empirically chosen to ensure the selection of  
335 well-represented motifs.

336 According to the presence or absence of these conserved fragments, three fragment-  
337 retained design evaluation sets were obtained:

- 338 • **Single-fragment-retained evaluation set:** 144 samples from 67 DUD-E proteins.  
339 Each sample contains one identified conserved fragment. For each sample, the active  
340 molecule corresponding to the conserved fragment is used as the set of known active  
341 compounds that the model needs to reproduce.
- 342 • **Two-fragment-retained evaluation set:** 79 samples from 30 DUD-E proteins. Each  
343 sample consists of two identified conserved fragments. Since the DUD-E active  
344 compounds set can hardly contain two identified conserved fragments simultaneously,  
345 we did not evaluate the active compounds reproducibility performance of models.

- 346 • **Three-fragment-retained evaluation set:** 19 samples from 9 DUD-E proteins. Each  
 347 sample contains three identified conserved fragments. Following the same rationale  
 348 applied to the two-fragment-retained set, we did not evaluate the reproducibility of the  
 349 active compounds.

### 350 2.3.3. Definition of Evaluation Metrics

351 The metrics used in this study are as follows:

- 352 • **# Generated molecules:** The sum of total number of generated molecules across 102  
 353 targets, with an upper limit of 1000 per target.
- 354 • **MW:** The averaged molecular weight of all generated molecules.
- 355 • **QED:** The averaged quantitative estimate of drug-likeness of all generated molecules.
- 356 • **SAS:** The averaged synthetic accessibility score of all generated molecules.
- 357 • **% drug-like molecules:** The percentage of molecules with QED > 0.3 and SAS < 5  
 358 over all generated molecules.
- 359 • **Bond length JSD:** The average statistical dissimilarity between the generated  
 360 molecules and the CSD<sup>7</sup> reference for 14 common bond length types.

$$\text{Bond length JSD} = \frac{\sum_i \frac{1}{2} D_{KL}(P||M) + \frac{1}{2} D_{KL}(Q||M)}{14}$$

361 where  $i$  represents one of the 14 common types of bond lengths, including C-C, C=C,  
 362 C:C, C#C, C-N, C=N, C#N, C:N, C-O, C=O, C:O, C-S, C=S and C:S. The symbol  $P$   
 363 and  $Q$  denote the bond length values from generated molecules and the CSD reference,  
 364 respectively. The term  $M$  refers to the intermediate distribution, defined as  $M = \frac{1}{2}(P +$   
 365  $Q)$ .

- 366 • **Bond angle JSD:** The average statistical dissimilarity between the generated molecules  
 367 and the CSD<sup>7</sup> reference for 11 common bond angle types.

$$\text{Bond angle JSD} = \frac{\sum_i \frac{1}{2} D_{KL}(P||M) + \frac{1}{2} D_{KL}(Q||M)}{11}$$

368 where  $i$  represents one of the 11 common types of bond angles, including CC=C,  
 369 CC=O, CCC, CCO, CN=C, CNC, COC, CSC, NC=O, NCC, and OC=O. The symbol  
 370  $P$  and  $Q$  denote the bond angles values from generated molecules and the CSD  
 371 reference, respectively. The term  $M$  refers to the intermediate distribution, defined as  
 372  $M = \frac{1}{2}(P + Q)$ .

- 373 • **Dihedral angle JSD:** The average statistical dissimilarity between the generated  
 374 molecules and the CSD<sup>7</sup> reference for 6 common dihedral angle types.

$$\text{Dihedral angle JSD} = \frac{\sum_i \frac{1}{2} D_{KL}(P||M) + \frac{1}{2} D_{KL}(Q||M)}{6}$$

375 where  $i$  represents one of the 6 common types of dihedral angles, including CCCC,  
376 cccc, CCCO, OCCO, Cccc, and CC=CC. The symbol  $P$  and  $Q$  denote the dihedral  
377 angles values from generated molecules and the CSD reference, respectively. The term  
378  $M$  refers to the intermediate distribution, defined as  $M = \frac{1}{2}(P + Q)$ .

- 379 • **Min-in-place GlideSP score:** The average of min-in-place GlideSP score using “min-  
380 in-place” method of Schrödinger Glide module. During calculation, we only considered  
381 drug-like molecules with a negative score.
- 382 • **Redocking GlideSP score:** The average of redocking GlideSP score using “docking”  
383 method of Schrödinger Glide module. During calculation, we only considered drug-  
384 like molecules with a negative score.
- 385 • **% IMP:** The percentage of generated drug-like molecules whose min-in-place GlideSP  
386 score is lower than their respective co-crystal ligand for the respective target. During  
387 calculation, we only considered drug-like molecules with a negative min-in-place score.
- 388 • **NCI recovery rate:** This metric is a measurement of the fraction of non-covalent  
389 interactions in the crystal ligand pose that are successfully replicated in the generated  
390 molecule pose using ProLIF, defined as:

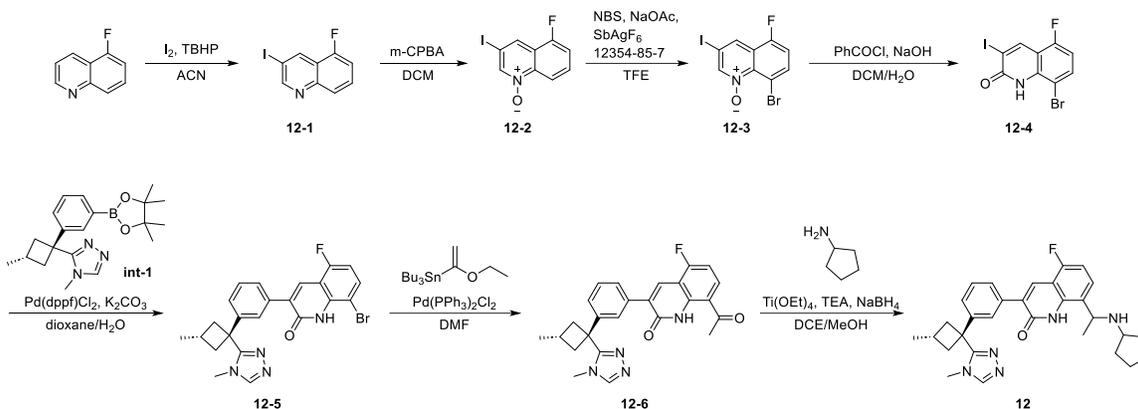
$$\text{NCI Recovery} = \frac{\sum_{i,r} \min(C_{i,r}, P_{i,r})}{\sum_{i,r} C_{i,r}}$$

391 where  $C_{i,r}$  is the count of type- $i$  interaction that crystal ligand formed with residue  $r$ ,  
392 and  $P_{i,r}$  is the count of type- $i$  interaction that a generated molecule formed with residue  
393  $r$ . Here, we considered interactions including hydrogen bond acceptor/donor, halogen  
394 bond acceptor/donor,  $\pi$ - $\pi$  stacking, cation- $\pi$ ,  $\pi$ -cation, anionic and cationic interactions  
395 were calculated by ProLIF.

- 396 • **% ECFP\_TS > 0.5:** The ratio for which at least one generated drug-like molecule has  
397 a Tanimoto similarity greater than 0.5 (based on ECFP4<sup>8</sup> fingerprints) to any known  
398 active compound for that target.

## 399 2.4. Compound Synthesis

### 400 2.4.1. Synthesis of Example 12



401

### 402 2.4.2. Synthesis of Intermediate 12-1

403 5-Fluoroquinoline (10 g, 68.03 mmol), iodine (20.7 g, 81.50 mmol), and 70 wt% *tert*-butyl  
404 hydroperoxide aqueous solution (43.7 g, 339.89 mmol) were dissolved in acetonitrile (100  
405 mL). The reaction mixture was stirred at 80°C for 16 hours. After cooling to room  
406 temperature, the mixture was diluted with saturated *aq.*  $Na_2CO_3$  until pH was adjusted to  
407 9, and extracted with ethyl acetate (200 mL  $\times$  3). The combined organics washed with brine  
408 (30 mL  $\times$  3), dried over  $Na_2SO_4$  and concentrated in vacuo to give crude product. The crude  
409 product was purified by silica gel column chromatography (eluent: PE/EA = 8:1) to give  
410 **12-1** (10.7 g, 57.6% yield) as a yellow solid, MS *m/z* (ESI): 274 [M+1].

### 411 2.4.3. Synthesis of Intermediate 12-2

412 To a mixture of **12-1** (10.7 g, 39.19 mmol) in dichloromethane (200 mL) was added *m*-  
413 chloroperbenzoic acid (13.56 g, 78.38 mmol) in batches under stirring at 0°C. The reaction  
414 mixture was stirred at room temperature for 16 hours. The mixture was diluted with H<sub>2</sub>O  
415 (100 mL) and adjusted pH to 9~10 by *aq.* NaOH (1 N) at 0°C. The mixture was extracted  
416 with dichloromethane (150 mL  $\times$  3), and the organic phases were combined, dried over  
417  $Na_2SO_4$ , evaporated and the residue was purified by silica gel column chromatography  
418 (eluent: PE/EA = 2:1) to afford **12-2** (5.0 g, 44.1% yield) as a yellow solid, MS *m/z* (ESI):  
419 290[M+1].

### 420 2.4.4. Synthesis of Intermediate 12-3

421 To a mixture of **12-2** (5.0 g, 17.30 mmol), dichloro(5-methylcyclopentadienyl)rhodium(III)  
422 dimer (536 mg, 0.87 mmol), silver hexafluoroantimonate (1.19 g, 3.47 mmol) and sodium  
423 acetate (284 mg, 3.46 mmol) in trifluoroethanol (90 mL) was added *N*-bromosuccinimide  
424 (4.0 g, 22.47 mmol) at room temperature. The reaction mixture was stirred at 50°C for 16  
425 hours. The reaction mixture was concentrated, and the residue was purified by silica gel

426 column chromatography (eluent: PE/EA = 4:1) to afford **12-3** (1.4 g, 21.9% yield) as a  
427 yellow solid, MS m/z (ESI): 368 [M+1].

#### 428 **2.4.5. Synthesis of Intermediate 12-4**

429 **12-3** (1.4 g, 3.80 mmol) was dissolved in a 30 mL mixture of dichloromethane and water  
430 (V/V = 1:1). Sodium hydroxide (304 mg, 7.60 mmol) was added, followed by the dropwise  
431 addition of benzoyl chloride (643 mg, 4.56 mmol) under stirring at 0°C. The reaction  
432 mixture was stirred at room temperature for 1 hour. After the reaction was completed, the  
433 precipitated solid was collected by filtration, and the filter cake was washed with water (5  
434 mL × 3), dried in vacuo to afford **12-4** (750 mg) as an off-white solid, MS m/z (ESI): 368  
435 [M+1].

#### 436 **2.4.6. Synthesis of Intermediate 12-5**

437 A mixture of **12-4** (600 mg, 1.63 mmol), **int-1** (692 mg, 1.96 mmol), [1,1'-  
438 bis(diphenylphosphino)ferrocene] palladium dichloride (119 mg, 0.16 mmol) and  
439 potassium carbonate (450 mg, 3.26 mmol) in 12 mL of dioxane and water (V/V = 5:1) was  
440 stirred at 80 °C under nitrogen for 2 hours. After cooling to room temperature, the reaction  
441 mixture was concentrated and the residue was purified by silica gel column  
442 chromatography (eluent: DCM/MeOH = 12:1) to afford **12-5** (560 mg, 73.5% yield) as a  
443 brown solid, MS m/z (ESI): 467 [M + 1].

#### 444 **2.4.7. Synthesis of Intermediate 12-6**

445 A mixture of **12-5** (300 mg, 0.64 mmol), tributyl(1-ethoxyvinyl)tin (462 mg, 1.28 mmol),  
446 and bis(triphenylphosphine)palladium dichloride (90 mg, 0.13 mmol) in *N,N*-  
447 dimethylformamide (5 mL) was stirred at 85 °C under nitrogen atmosphere for 2 hours.  
448 The reaction mixture was cooled to room temperature and 2N hydrochloric acid solution  
449 (10 mL) was added. The mixture was stirred at room temperature for 1 hour. The mixture  
450 was adjusted to pH = 9 by saturated *aq.* Na<sub>2</sub>CO<sub>3</sub> (10 mL × 3), extracted by ethyl acetate  
451 (20 mL × 3). The combined organics was washed by brine (10 mL × 3), dried over Na<sub>2</sub>SO<sub>4</sub>,  
452 evaporated and the residue was purified by silica gel column chromatography (eluent:  
453 DCM/MeOH = 10:1) to afford **12-6** (170 mg, 61.5% yield) as a yellow solid, MS m/z (ESI):  
454 431 [M+1].

#### 455 **2.4.8. Synthesis of Compound 12**

456 **12-6** (140 mg, 0.33 mmol), cyclopentylamine (56 mg, 0.66 mmol), tetraethyl titanate (226  
457 mg, 0.99 mmol) and triethylamine (133 mg, 1.32 mmol) were dissolved in 1,2-  
458 dichloroethane (4 mL). The reaction mixture was stirred at 80°C for 16 hours. After cooling  
459 to room temperature, methanol (4 mL) was added to the reaction mixture, followed by the  
460 addition of sodium borohydride (63 mg, 1.66 mmol) in portions while stirring at room

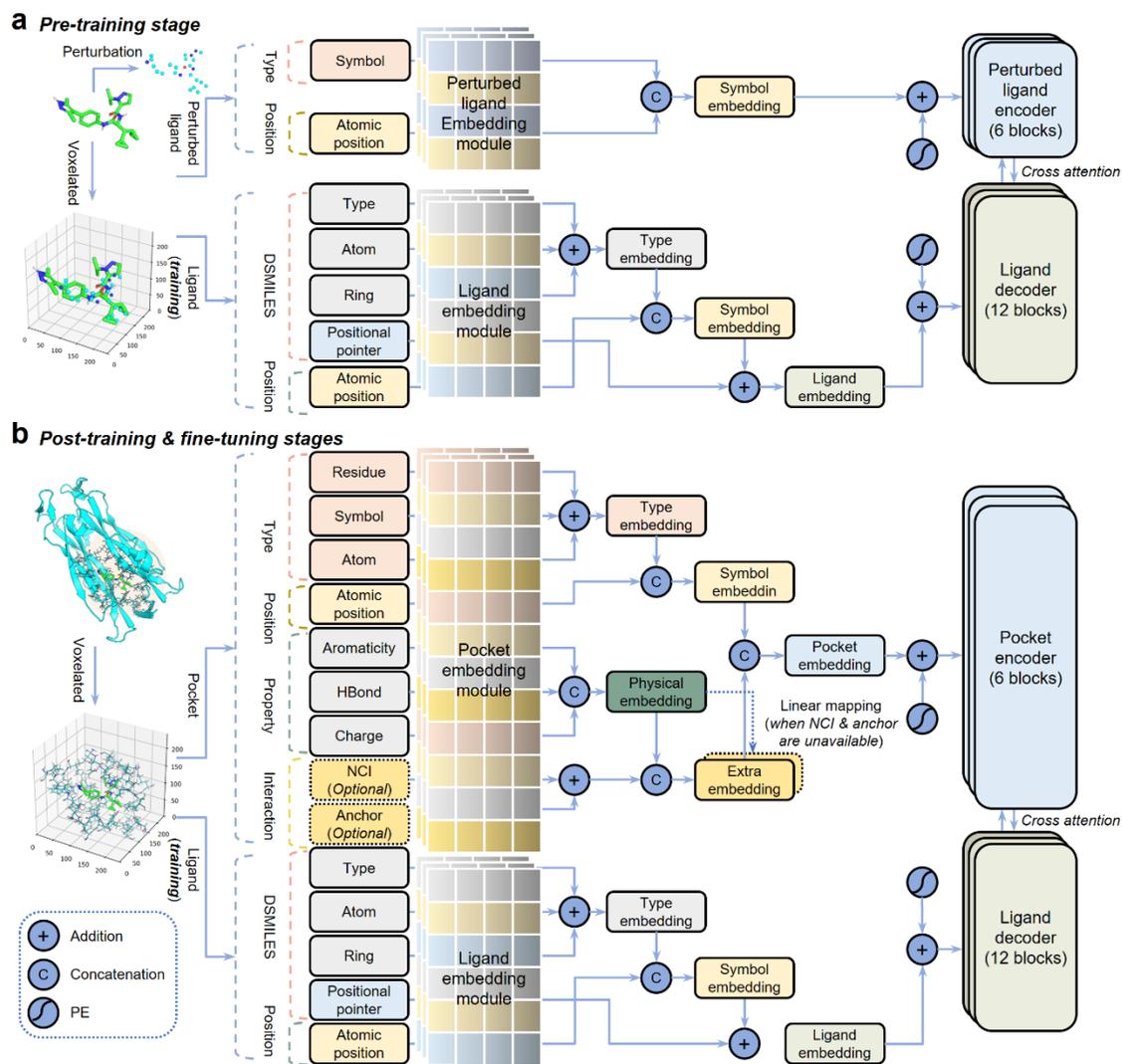
461 temperature. After the addition was completed, the reaction mixture was stirred at room  
462 temperature for another 3 hours. The mixture was diluted with H<sub>2</sub>O (20 mL) and extracted  
463 with dichloromethane (10 mL × 3). The organic phases were combined, dried over Na<sub>2</sub>SO<sub>4</sub>,  
464 evaporated and the residue was purified by silica gel column chromatography (eluent:  
465 DCM/MeOH = 9:1) to afford crude product. The crude product was purified by *Prep*-  
466 HPLC (column: Xbridge-C18 150 x 19 mm; eluent: ACN-H<sub>2</sub>O (0.05%FA); gradient: 15-  
467 40%) to afford **Cmpd. 12** (25 mg, 15.4% yield) as a white solid.

468 **Cmpd. 12** (25 mg) was separated by SFC (column: CHIRALPAK AD-H, 250 mm × 20  
469 mm, eluent: 40% IPA (NH<sub>4</sub>OH 0.2%), flow velocity: 40mL/min, RT1: 4.80 min, RT2: 5.95  
470 min, column temperature: 40°C) to afford **Cmpd. 19** (8.29 mg) as a white solid. MS m/z  
471 (ESI): 500[M+1], <sup>1</sup>H NMR (400 MHz, DMSO-*d*<sub>6</sub>) δ ppm: 8.28 (s, 1H), 8.02 (s, 1H), 7.76  
472 (s, 1H), 7.61 (d, *J* = 7.7 Hz, 1H), 7.48 - 7.29 (m, 3H), 6.98 (dd, *J* = 9.8, 8.3 Hz, 1H), 4.27  
473 - 4.18 (m, 1H), 3.22 (s, 3H), 2.92 - 2.81 (m, 3H), 2.54 (d, *J* = 6.6 Hz, 3H), 1.83 - 1.73 (m,  
474 1H), 1.72 - 1.53 (m, 3H), 1.50 - 1.26 (m, 7H), 1.09 (d, *J* = 4.9 Hz, 3H); and **Cmpd. 20**  
475 (7.47 mg) as a white solid. MS m/z (ESI): 500[M+1], <sup>1</sup>H NMR (400 MHz, DMSO-*d*<sub>6</sub>) δ  
476 ppm: 8.28 (s, 1H), 8.02 (s, 1H), 7.76 (s, 1H), 7.61 (d, *J* = 7.8 Hz, 1H), 7.48 - 7.29 (m, 3H),  
477 6.98 (dd, *J* = 9.8, 8.3 Hz, 1H), 4.27 - 4.18 (m, 1H), 3.22 (s, 3H), 2.92 - 2.82 (m, 3H), 2.54  
478 (d, *J* = 6.5 Hz, 3H), 1.83 - 1.73 (m, 1H), 1.72 - 1.53 (m, 3H), 1.50 - 1.25 (m, 7H), 1.09 (d,  
479 *J* = 4.8 Hz, 3H).

480 The compounds shown in [Supplementary Table S9](#) were synthesized using the same  
481 methods as above.

482

483 **3. Supplementary Figures S1-S4**



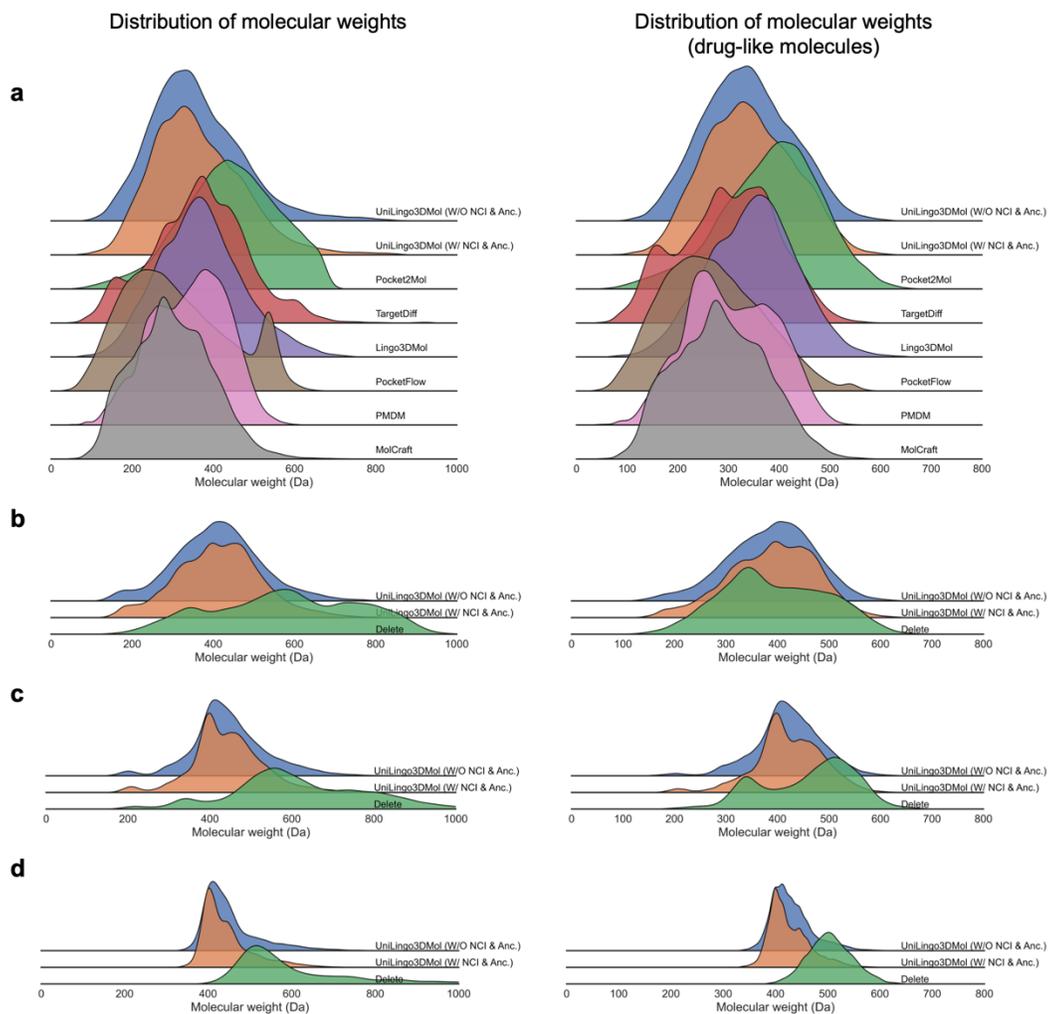
484

485 **Supplementary Figure S1. Working mechanism of UniLingo3DMol across training**  
 486 **stages.**

487 **(a) Pre-training stage.** The encoder processes perturbed ligand inputs, while the decoder  
 488 generates ligands as DSILES sequences with restored 3D conformations.

489 **(b) Post-training and fine-tuning stages.** The encoder receives protein binding pockets  
 490 as input.

491

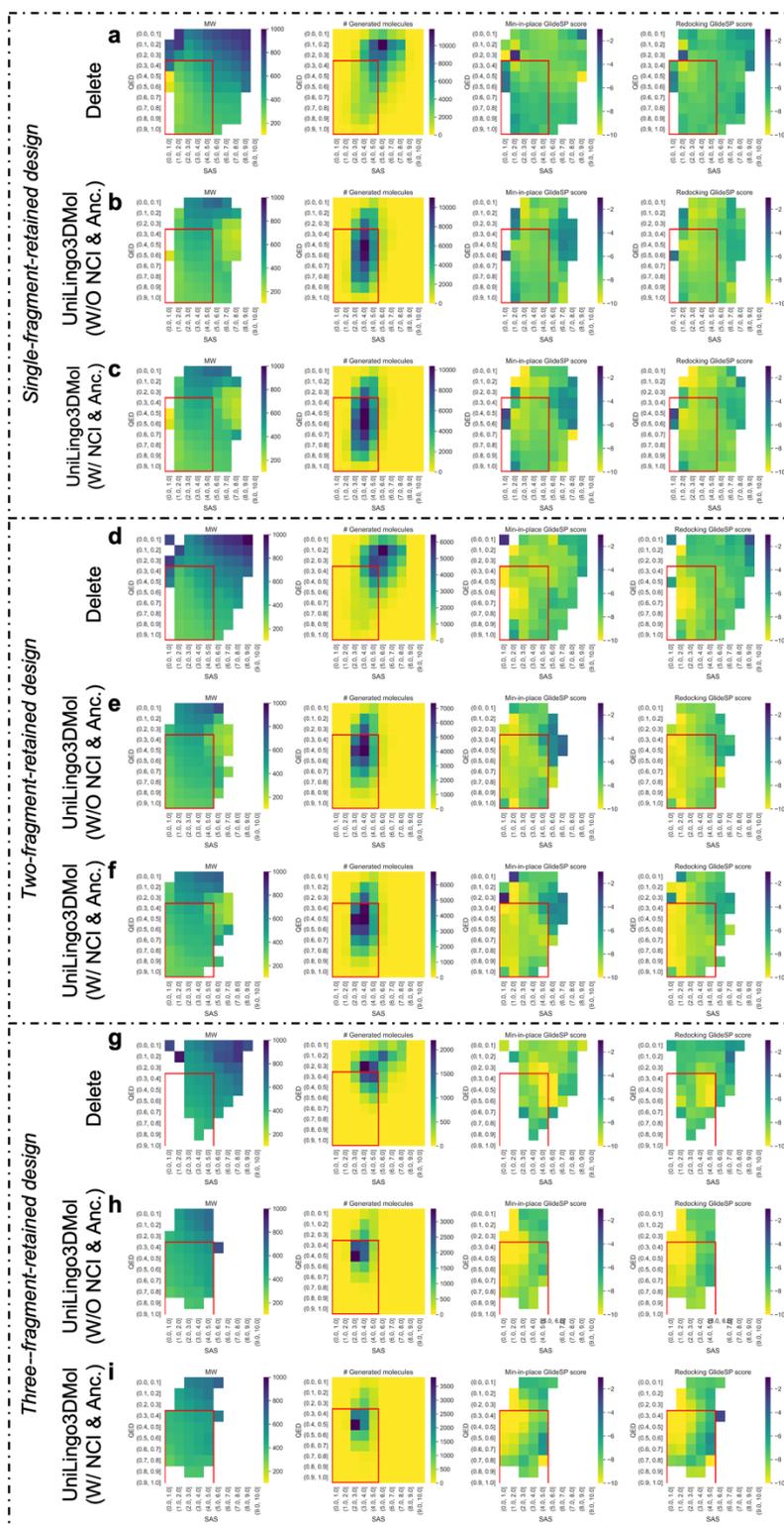


492

493 **Supplementary Figure S2. Distribution of molecular weight for molecules generated**  
 494 **by various methods on different evaluation sets.**

495 **(a-d)** show results for the *de novo* design, single-fragment-retained, two-fragment-retained,  
 496 and three-fragment-retained evaluation sets, respectively. **(a)** is supplementary to [Table 1](#),  
 497 while **(b-d)** are supplementary to [Supplementary Table S2](#).

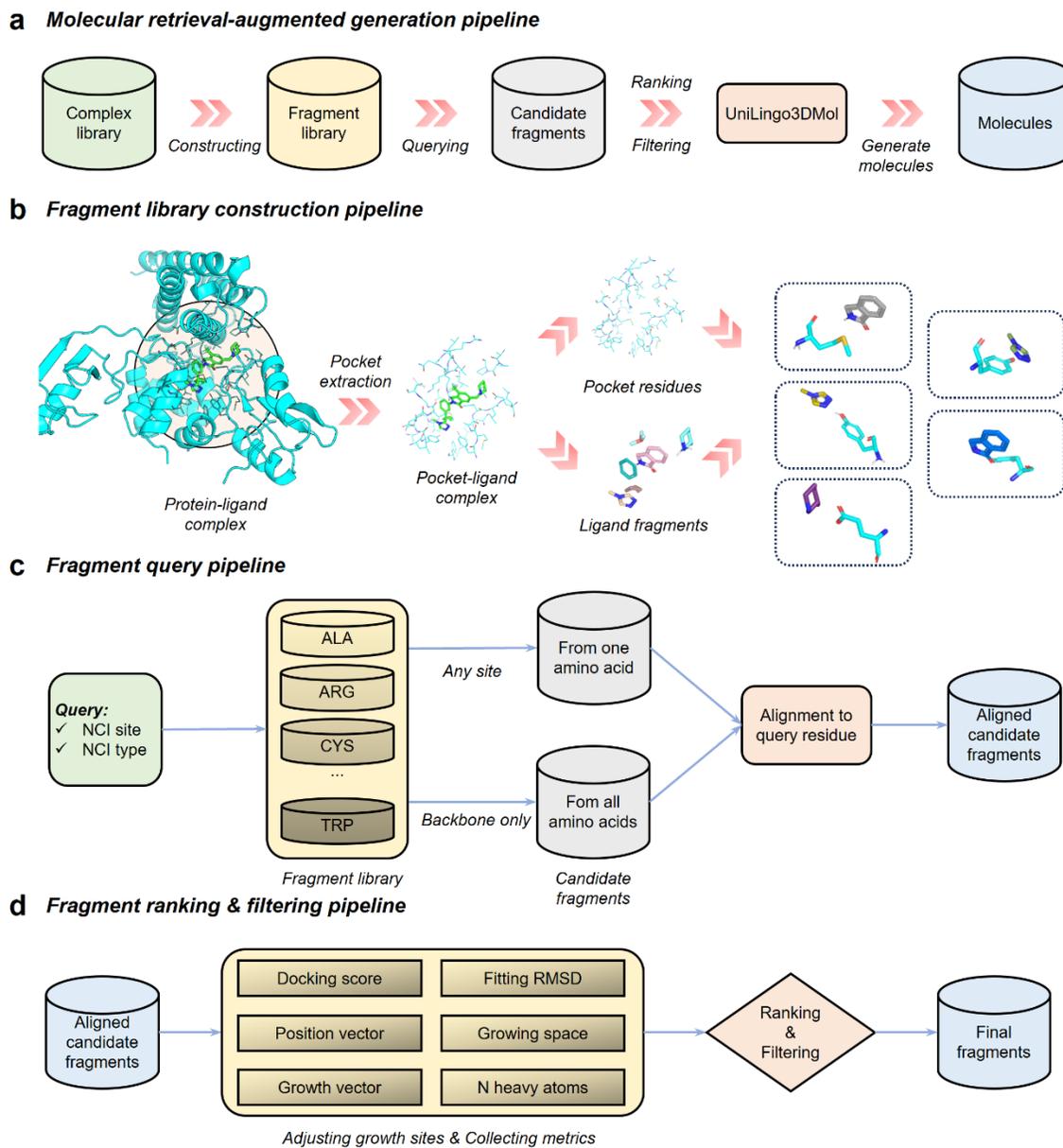
498



499  
500  
501

Supplementary Figure S3. Distributions of molecules generated by various models on the constrained design evaluation set.

502 The drug-like region with  $QED > 0.3$  and  $SAS < 5$  is indicated with red boxes. **(a-i)** show  
503 distributions on the single **(a-c)**, two **(d-f)**, and three-fragment-retained **(g-i)** evaluation sets,  
504 respectively.  
505



506

507 **Supplementary Figure S4. Overview of MRAG schemes.**

508 (a) Molecular retrieval-augmented generation pipeline.

509 (b) Fragment library construction pipeline.

510 (c) Fragment query pipeline.

511 (d) Fragment ranking and filtering pipeline.

512

513 **4. Supplementary Tables S1-S9**

514 **Supplementary Table S1. JSD of bond length, bond angle, and dihedral angle for**  
 515 **different models.** The top two results for each metric are highlighted in **bold** and  
 516 underlined, separately.

Evaluation set	Model	Bond length (↓)	Bond angle (↓)	Dihedral angle (↓)
<i>De novo</i>	Pocket2Mol	0.50	0.33	0.23
	TargetDiff	0.43	0.34	0.20
	Lingo3DMol	0.40	<u>0.23</u>	<b>0.15</b>
	PocketFlow	0.51	0.35	0.26
	PMDM	<u>0.38</u>	0.28	0.31
	MolCraft	0.39	0.27	<b>0.16</b>
	UniLingo3DMol (W/O NCI & Anc.)	<b>0.33</b>	<b>0.22</b>	<b>0.16</b>
	UniLingo3DMol (W/ NCI & Anc.)	<b>0.33</b>	<b>0.22</b>	<b>0.16</b>
	Delete	0.48	0.34	<u>0.22</u>
<b>Single-fragment-retained</b>	UniLingo3DMol (W/O NCI & Anc.)	<b>0.31</b>	<u>0.23</u>	<b>0.14</b>
	UniLingo3DMol (W/ NCI & Anc.)	<u>0.33</u>	<b>0.22</b>	<b>0.14</b>
	Delete	0.47	<u>0.34</u>	<u>0.22</u>
<b>Two-fragment-retained</b>	UniLingo3DMol (W/O NCI & Anc.)	<b>0.32</b>	<b>0.24</b>	<b>0.18</b>
	UniLingo3DMol (W/ NCI & Anc.)	<u>0.34</u>	<b>0.24</b>	<b>0.18</b>
	Delete	0.50	<u>0.34</u>	<b>0.21</b>
<b>Three-fragment-retained</b>	UniLingo3DMol (W/O NCI & Anc.)	<b>0.35</b>	<b>0.27</b>	<u>0.23</u>
	UniLingo3DMol (W/ NCI & Anc.)	<u>0.38</u>	<b>0.27</b>	<u>0.23</u>
	Delete	0.50	<u>0.34</u>	<b>0.21</b>

517 Note: 14 common bond length types include C-C, C=C, C:C, C#C, C-N, C=N, C#N, C:N, C-O, C=O,  
 518 C:O, C-S, C=S and C:S. 11 common bond angle types include CC=C, CC=O, CCC, CCO, CN=C, CNC,  
 519 COC, CSC, NC=O, NCC, and OC=O. 6 common dihedral angle types include CCCC, cccc, CCCO,  
 520 OCCO, Cccc, and CC=CC.

521

522 **Supplementary Table S2. Comparison of generated drug-like molecules on the**  
 523 **fragment-retained design scenario.** The top two results for each metric are highlighted  
 524 in **bold** and underlined, separately.

	single-fragment-retained (N=144)			two-fragment-retained (N=79)			three-fragment-retained (N=19)		
	Delete	UniLingo3DMol		Delete	UniLingo3DMol		Delete	UniLingo3DMol	
		W/O NCI & Anc.	W/ NCI & Anc.		W/O NCI & Anc.	W/ NCI & Anc.		W/ NCI & Anc.	W/O NCI & Anc.
# Generated molecules	137,302	143,294	143,148	77,249	76,343	76,883	19,000	18,373	18,075
Mean Molecular Weight (MW)	579	412	414	614	448	444	597	452	444
Mean QED (↑)	<u>0.34</u>	<b>0.50</b>	<b>0.50</b>	0.30	<u>0.44</u>	<b>0.46</b>	0.27	<u>0.39</u>	<b>0.41</b>
Mean SAS (↓)	5.0	<u>3.6</u>	<b>3.5</b>	<u>5.1</u>	<b>3.4</b>	<b>3.4</b>	4.6	<u>3.2</u>	<b>3.1</b>
# Drug-like molecules	40,829	111,799	112,944	20,783	57,338	59,892	6,128	14,491	14,769
% Drug-like molecules (↑)	29.7%	<u>78.0%</u>	<b>78.9%</b>	26.9%	<u>75.1%</u>	<b>77.9%</b>	32.3%	<u>78.9%</u>	<b>81.7%</b>
<i>The comparison below involves only drug-like molecules</i>									
Mean MW	395	388	390	459	421	421	505	430	426
Mean QED (↑)	<u>0.56</u>	<b>0.57</b>	<b>0.57</b>	0.49	<u>0.51</u>	<b>0.52</b>	0.40	<b>0.44</b>	<u>0.41</u>
Mean SAS (↓)	<u>3.7</u>	<b>3.4</b>	<b>3.4</b>	<u>3.9</u>	<b>3.3</b>	<b>3.3</b>	<u>4.0</u>	<b>3.1</b>	<b>3.1</b>
Mean Min-in-place score (↓)	<u>0.48</u>	<b>0.33</b>	<b>0.33</b>	0.46	<b>0.32</b>	<u>0.36</u>	0.47	<b>0.33</b>	<u>0.40</u>
Mean bond length JSD (↓)	0.32	<u>0.23</u>	<b>0.22</b>	<u>0.33</u>	<b>0.24</b>	<b>0.24</b>	0.33	<u>0.28</u>	<b>0.27</b>
Mean bond angle JSD (↓)	<u>0.22</u>	<b>0.16</b>	<b>0.16</b>	<u>0.23</u>	<b>0.19</b>	<b>0.19</b>	<b>0.23</b>	<u>0.24</u>	<u>0.24</u>
Mean dihedral angle JSD (↓)	-7.5	<u>-8.1</u>	<b>-8.3</b>	-8.6	<u>-8.9</u>	<b>-9.0</b>	-9.5	<b>-9.8</b>	<b>-9.8</b>
Mean Redocking score (↓)	-7.7	<u>-8.4</u>	<b>-8.5</b>	<u>-8.2</u>	<b>-9.0</b>	<b>-9.0</b>	-8.6	<b>-9.8</b>	-9.7
% IMP (↑)	15.3%	<u>15.0%</u>	<b>16.9%</b>	7.2%	<u>13.5%</u>	<b>14.6%</b>	13.1%	<b>16.0%</b>	<u>15.3%</u>
% NCI recovery (↑)	27.9%	<u>50.8%</u>	<b>58.4%</b>	33.7%	<u>57.3%</u>	<b>63.5%</b>	61.2%	<u>79.3%</u>	<b>83.8%</b>
% ECFP_TS > 0.5 (↑)	2%	<u>53%</u>	<b>55%</b>	N/A	N/A	N/A	N/A	N/A	N/A
Time cost (s,↓)	3729±2671	<u>21±4</u>	<b>23±4</b>	7307±6349	<u>31±16</u>	<b>25±14</b>	8636±8939	<u>32±20</u>	<b>27±13</b>

525 Note: For each method, we generated approximately 1000 molecules per case, using the same metrics  
 526 as in the de novo design evaluation. The % ECFP\_TS > 0.5 metric is only applicable for the single-

527 fragment-retained evaluation set, as it is not computable for the two-fragment and three-fragment-  
528 retained sets due to their specific construction ([Supplementary Section 2.3.2](#)). Additional molecular  
529 weight distribution is provided in [Supplementary Figure S2b-d](#). In addition, [Supplementary Figure S3](#)  
530 illustrates the distribution of properties of generated molecules within the entire QED-SAS chemical  
531 space.  
532  
533

534 **Supplementary Table S3. Comparative ablation study of UniLingo3DMol-generated**  
 535 **molecules under the three-stage training strategy.**

	W/O pre-training	W/O post-training	W/O fine-tuning	Standard
# Generated molecules	97,116	102,000	102,000	101,439
Mean Molecular Weight (MW)	392	445	405	357
Mean QED (↑)	<u>0.52</u>	0.40	0.49	<b>0.55</b>
Mean SAS (↓)	<b>3.3</b>	4.0	3.5	<u>3.4</u>
# Drug-like molecules	75,383	58,583	75,222	81,835
% Drug-like molecules (↑)	77.6%	57.4%	<u>73.7%</u>	<b>80.7%</b>
<i>The comparison below involves only drug-like molecules</i>				
Mean MW	359	391	370	337
Mean QED (↑)	0.59	0.54	<u>0.58</u>	<b>0.62</b>
Mean SAS (↓)	<b>3.1</b>	3.5	<u>3.2</u>	<u>3.2</u>
Mean bond length JSD (↓)	0.35	0.35	<u>0.34</u>	<b>0.33</b>
Mean bond angle JSD (↓)	<u>0.23</u>	<u>0.23</u>	<u>0.23</u>	<b>0.22</b>
Mean dihedral angle JSD (↓)	<b>0.13</b>	0.20	<u>0.16</u>	<u>0.16</u>
Mean Min-in-place score (↓)	-5.0	-6.5	<b>-7.4</b>	<u>-7.1</u>
Mean Redocking score (↓)	-7.4	-7.3	<b>-8.0</b>	<u>-7.8</u>
% IMP (↑)	7.3%	10.2%	<b>14.8%</b>	<u>11.0%</u>
% NCI recovery (↑)	6.2%	<u>25.1%</u>	<b>33.2%</b>	<b>33.2%</b>
% ECFP_TS > 0.5 (↑)	8%	18%	<u>71%</u>	<b>74%</b>

536 Note: Each of the above methods was evaluated on the *de novo* evaluation set by generating  
 537 approximately 1,000 molecules without NCI and anchor site annotations.  
 538

539 **Supplementary Table S4. Comparative ablation study of UniLingo3DMol-generated**  
 540 **molecules under the multi-task training strategy.**

	<b>W/O Task 1</b>	<b>Standard</b>
<b># Generated molecules</b>	101,745	101,439
<b>Mean Molecular Weight (MW)</b>	364	357
<b>Mean QED (↑)</b>	0.54	<b>0.55</b>
<b>Mean SAS (↓)</b>	3.5	<b>3.4</b>
<b># Drug-like molecules</b>	80,743	81,835
<b>% Drug-like molecules (↑)</b>	79.4%	<b>80.7%</b>
<i>The comparison below involves only drug-like molecules</i>		
<b>Mean MW</b>	343	337
<b>Mean QED (↑)</b>	0.61	<b>0.62</b>
<b>Mean SAS (↓)</b>	<b>3.2</b>	<b>3.2</b>
<b>Mean bond length JSD (↓)</b>	<b>0.33</b>	<b>0.33</b>
<b>Mean bond angle JSD (↓)</b>	<b>0.21</b>	0.22
<b>Mean dihedral angle JSD (↓)</b>	<b>0.15</b>	0.16
<b>Mean Min-in-place score (↓)</b>	-6.9	<b>-7.1</b>
<b>Mean Redocking score (↓)</b>	-7.7	<b>-7.8</b>
<b>% IMP (↑)</b>	10.3%	<b>11.0%</b>
<b>% NCI recovery (↑)</b>	29.0%	<b>33.2%</b>
<b>% ECFP_TS &gt; 0.5 (↑)</b>	71%	<b>74%</b>

541 Note: Each of the above methods was evaluated on the *de novo* evaluation set by generating  
 542 approximately 1,000 molecules without NCI and anchor site annotations.  
 543

544 **Supplementary Table S5. Comparative ablation study of UniLingo3DMol-generated**  
 545 **molecules across inference tasks.**

	W/ predicted NCI & Anc.	W/O NCI & Anc.
# Generated molecules	101,163	101,439
Mean Molecular Weight (MW)	418	357
Mean QED (↑)	0.48	<b>0.55</b>
Mean SAS (↓)	3.7	<b>3.4</b>
# Drug-like molecules	71,383	81,835
% Drug-like molecules (↑)	70.6%	<b>80.7%</b>
<i>The comparison below involves only drug-like molecules</i>		
Mean MW	383	337
Mean QED (↑)	0.57	<b>0.62</b>
Mean SAS (↓)	3.4	<b>3.2</b>
Mean bond length JSD (↓)	0.35	<b>0.33</b>
Mean bond angle JSD (↓)	<b>0.22</b>	<b>0.22</b>
Mean dihedral angle JSD (↓)	0.17	<b>0.16</b>
Mean Min-in-place score (↓)	-7.5	-7.1
Mean Redocking score (↓)	<b>-8.0</b>	-7.8
% IMP (↑)	<b>16.9%</b>	11.0%
% NCI recovery (↑)	<b>42.4%</b>	33.2%
% ECFP_TS > 0.5 (↑)	59%	<b>74%</b>

546 Note: Each of the above methods was evaluated on the *de novo* evaluation set. The notation “predicted  
 547 NCI and anchor” refers to a two-step process: first predicting NCI and anchor sites and then generating  
 548 molecules with NCI and anchor site annotations.  
 549

550 **Supplementary Table S6. Comparative ablation study of UniLingo3DMol-generated**  
 551 **molecules with/without fragment-retained data.**

	<b>W/O fragment-retained data</b>	<b>Standard</b>
<b># Generated molecules</b>	102,000	101,439
<b>Mean Molecular Weight (MW)</b>	351	357
<b>Mean QED (↑)</b>	<b>0.56</b>	0.55
<b>Mean SAS (↓)</b>	<b>3.4</b>	<b>3.4</b>
<b># Drug-like molecules</b>	83,333	81,835
<b>% Drug-like molecules (↑)</b>	<b>81.7%</b>	80.7%
<i>The comparison below involves only drug-like molecules</i>		
<b>Mean MW</b>	331	337
<b>Mean QED (↑)</b>	<b>0.62</b>	<b>0.62</b>
<b>Mean SAS (↓)</b>	<b>3.2</b>	<b>3.2</b>
<b>Mean bond length JSD (↓)</b>	0.34	<b>0.33</b>
<b>Mean bond angle JSD (↓)</b>	<b>0.21</b>	0.22
<b>Mean dihedral angle JSD (↓)</b>	<b>0.16</b>	<b>0.16</b>
<b>Mean Min-in-place score (↓)</b>	-6.8	<b>-7.1</b>
<b>Mean Redocking score (↓)</b>	-7.6	<b>-7.8</b>
<b>% IMP (↑)</b>	9.5%	<b>11.0%</b>
<b>% NCI recovery (↑)</b>	28.6%	<b>33.2%</b>
<b>% ECFP_TS &gt; 0.5 (↑)</b>	67%	<b>74%</b>

552 Note: Each of the above methods was evaluated on the *de novo* evaluation set by generating  
 553 approximately 1,000 molecules without NCI and anchor site annotations.  
 554

555 **Supplementary Table S7. Comparative ablation study of UniLingo3DMol-generated**  
 556 **molecules under MRAG strategy.**

	W/O MRAG		W/ MRAG	
	W/O NCI & Anc.	W/ NCI & Anc.	W/O NCI & Anc.	W/ NCI & Anc.
# Generated molecules	101,439	100,932	100,303	99,519
Mean Molecular Weight (MW)	357	362	395	393
Mean QED (↑)	<u>0.55</u>	<b>0.56</b>	0.49	0.50
Mean SAS (↓)	<b>3.4</b>	<u>3.5</u>	3.9	3.9
# Drug-like molecules	81,835	81,241	70,914	70,461
% Drug-like molecules (↑)	<b>80.7%</b>	<u>80.5%</u>	70.7%	70.8%
<i>The comparison below involves only drug-like molecules</i>				
Mean MW	337	341	365	365
Mean QED (↑)	<b>0.62</b>	<b>0.62</b>	0.57	<u>0.58</u>
Mean SAS (↓)	<b>3.2</b>	<b>3.2</b>	<u>3.6</u>	<u>3.6</u>
Mean bond length JSD (↓)	<u>0.33</u>	<u>0.33</u>	<b>0.30</b>	0.34
Mean bond angle JSD (↓)	<b>0.22</b>	<b>0.22</b>	<b>0.22</b>	<b>0.22</b>
Mean dihedral angle JSD (↓)	<b>0.16</b>	<b>0.16</b>	<u>0.17</u>	0.18
Mean Min-in-place score (↓)	-7.1	-7.3	<u>-7.6</u>	<b>-7.7</b>
Mean Redocking score (↓)	<u>-7.8</u>	<u>-7.8</u>	<b>-8.0</b>	<b>-8.0</b>
% IMP (↑)	11.0%	11.2%	<u>14.0%</u>	<b>15.4%</b>
% NCI recovery (↑)	33.2%	42.8%	<u>59.9%</u>	<b>62.6%</b>
% ECFP_TS > 0.5 (↑)	<u>74%</u>	<b>76%</b>	32%	31%

557 Note: Each of the above methods was evaluated on the *de novo* evaluation set by generating  
 558 approximately 1,000 molecules without NCI and anchor site annotations.

559 **Supplementary Table S8. Hyperparameters in each stage of UniLingo3DMol.**

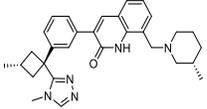
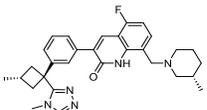
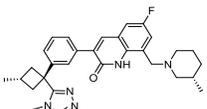
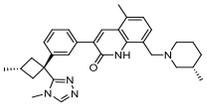
	Pre-training	Post-training	Fine-tuning
<i>Hyperparameter</i>			
Number of encoder layers	6	6	6
Number of decoder layers	12	12	12
Max sequence length of encoder	150	500	500
Max sequence length of decoder	150	150	150
Hidden size	512	512	512
Number of attention heads	8	8	8
Intermediate size	1024	1024	1024
Number of head layers	2	2	2
Head hidden size	256	256	256
<i>Training</i>			
<i>De novo</i> data (include augmentation)	~405M	~58M	~1M
Fragment-retained data (include augmentation)	NA	~54M	~1M
Optimizer	AdamW	AdamW	AdamW
$\beta_1$	0.9	0.9	0.9
$\beta_2$	0.999	0.999	0.999
$\epsilon$	1e-8	1e-8	1e-8
Learning rate	5e-5	5e-5	5e-5
weight decay	0.01	0.01	0.01
GPU type	H20	H20	H20
Number of GPUs	32	32	8
Hours to train	~226	~170	~0.5
Training epochs	500	200	10
<i>Inference</i>			
Sampling temperature (for <i>de novo</i> design)	NA	NA	0.5 (suggested)
Sampling temperature (for <i>fragment-retained design</i> )	NA	NA	1.5 (suggested)

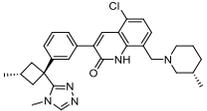
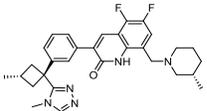
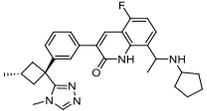
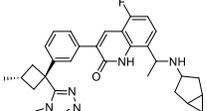
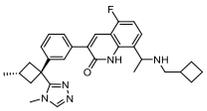
560

561

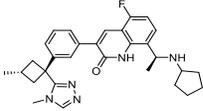
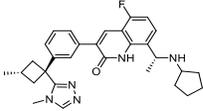
562 **Supplementary Table S9. Synthesized compounds.**

Cmpd. No	Structure	LCMS (ESI) m/z	<sup>1</sup> H NMR
Cmpd. 1		411	<sup>1</sup> H NMR (400 MHz, DMSO- <i>d</i> <sub>6</sub> ) δ 8.55 (s, 1H), 7.54 (dd, <i>J</i> = 8.7, 5.8 Hz, 1H), 7.44 - 7.39 (m, 1H), 7.34 - 7.31 (m, 1H), 6.74 (d, <i>J</i> = 1.0 Hz, 1H), 6.37 (s, 1H), 4.33 (t, <i>J</i> = 9.4 Hz, 1H), 4.04 - 3.87 (m, 1H), 3.68 - 3.61 (m, 1H), 3.43 (s, 3H), 2.27 - 2.16 (m, 4H), 2.15 - 2.04 (m, 2H), 1.60 - 1.56 (m, 1H).
Cmpd. 2		388	<sup>1</sup> H NMR (400 MHz, CDCl <sub>3</sub> ) δ ppm: 9.39 (br s, 1H), 8.28 (d, <i>J</i> = 8.3 Hz, 1H), 8.02 (s, 1H), 8.00 (s, 1H), 7.91 (d, <i>J</i> = 7.8 Hz, 1H), 7.77 (t, <i>J</i> = 7.9 Hz, 1H), 7.67 - 7.61 (m, 1H), 7.42 - 7.34 (m, 3H), 7.26 - 7.22 (m, 1H), 7.01 (d, <i>J</i> = 7.5 Hz, 1H), 2.98 (s, 3H), 2.56 (s, 3H).
Cmpd. 3		403	<sup>1</sup> H NMR (400 MHz, CDCl <sub>3</sub> ) δ ppm: 8.00 (s, 1H), 7.93 (s, 1H), 7.79 (d, <i>J</i> = 11.3 Hz, 1H), 7.62 - 7.51 (m, 2H), 7.42 - 7.34 (m, 4H), 7.26 (s, 1H), 6.90 (t, <i>J</i> = 8.0 Hz, 1H), 4.83 (s, 2H), 3.00 (s, 3H).
Cmpd. 4		371	<sup>1</sup> H NMR (400 MHz, DMSO- <i>d</i> <sub>6</sub> ) δ ppm: 11.92 (s, 1H), 8.30 (s, 1H), 8.09 (s, 1H), 7.78 - 7.71 (m, 2H), 7.66 - 7.58 (m, 1H), 7.55 - 7.46 (m, 1H), 7.42 (t, <i>J</i> = 7.7 Hz, 1H), 7.37 - 7.27 (m, 2H), 7.24 - 7.16 (m, 1H), 3.23 (s, 3H), 2.93 - 2.86 (m, 2H), 2.60 - 2.50 (m, 3H), 1.09 (d, <i>J</i> = 5.1 Hz, 3H).
Cmpd. 5		403	<sup>1</sup> H NMR (400 MHz, CDCl <sub>3</sub> ) δ ppm: 8.23 (s, 1H), 8.08 (s, 1H), 7.66 - 7.54 (m, 1H), 7.38 - 7.27 (m, 2H), 7.14 (d, <i>J</i> = 7.5 Hz, 1H), 6.92 (s, 1H), 6.87 (d, <i>J</i> = 7.1 Hz, 1H), 4.79 (t, <i>J</i> = 7.9 Hz, 2H), 3.21 (t, <i>J</i> = 8.3 Hz, 2H), 3.07 (s, 3H), 2.28 (s, 3H).

Cmpd. 6		482	<sup>1</sup> H NMR (400 MHz, DMSO- <i>d</i> <sub>6</sub> ) δ ppm: 12.08 (s, 1H), 8.28 (s, 1H), 8.14 (s, 1H), 7.80 (s, 1H), 7.69 (d, <i>J</i> = 7.2 Hz, 1H), 7.61 (d, <i>J</i> = 7.7 Hz, 1H), 7.50 - 7.25 (m, 3H), 7.16 (t, <i>J</i> = 7.6 Hz, 1H), 3.84 (s, 2H), 3.22 (s, 3H), 2.88 (d, <i>J</i> = 3.8 Hz, 2H), 2.78 (d, <i>J</i> = 9.4 Hz, 2H), 2.56 - 2.54 (m, 3H), 2.00 - 1.98 (m, 1H), 1.79 - 1.59 (m, 4H), 1.52 - 1.48 (m, 1H), 1.09 (d, <i>J</i> = 5.3 Hz, 3H), 0.95 - 0.93 (m, 1H), 0.85 (d, <i>J</i> = 6.2 Hz, 3H).
Cmpd. 7		500	<sup>1</sup> H NMR (400 MHz, CDCl <sub>3</sub> ) δ ppm: 9.35 (s, 1H), 8.11 (s, 1H), 8.09 (s, 1H), 7.82 (s, 1H), 7.58 (d, <i>J</i> = 7.6 Hz, 1H), 7.44 (t, <i>J</i> = 7.7 Hz, 1H), 7.41 - 7.35 (m, 2H), 6.89 (t, <i>J</i> = 8.8 Hz, 1H), 4.01 (s, 2H), 3.30 (s, 3H), 3.06 - 2.88 (m, 4H), 2.68 - 2.55 (m, 3H), 2.22 (t, <i>J</i> = 10.9 Hz, 1H), 1.99 - 1.83 (m, 2H), 1.83 - 1.66 (m, 3H), 1.14 (d, <i>J</i> = 5.7 Hz, 3H), 1.05 - 0.92 (m, 1H), 0.88 (d, <i>J</i> = 6.2 Hz, 3H).
Cmpd. 8		500	<sup>1</sup> H NMR (400 MHz, CDCl <sub>3</sub> ) δ ppm: 12.34 (s, 1H), 7.98 (s, 1H), 7.81 (s, 1H), 7.77 (s, 1H), 7.62 (d, <i>J</i> = 7.4 Hz, 1H), 7.44 - 7.32 (m, 2H), 7.21 (dd, <i>J</i> = 8.2, 2.1 Hz, 1H), 7.08 (s, 1H), 3.83 (s, 2H), 3.26 (s, 3H), 2.93 - 2.86 (m, 4H), 2.67 (d, <i>J</i> = 6.1 Hz, 3H), 2.10 - 1.97 (m, 1H), 1.85 - 1.72 (m, 5H), 1.14 (d, <i>J</i> = 5.2 Hz, 3H), 0.98 - 0.89 (m, 4H).
Cmpd. 9		496	<sup>1</sup> H NMR (400 MHz, DMSO- <i>d</i> <sub>6</sub> ) δ ppm: 12.17 (s, 1H), 8.29 (s, 1H), 8.08 (s, 1H), 7.83 (s, 1H), 7.61 (d, <i>J</i> = 7.7 Hz, 1H), 7.44 (t, <i>J</i> = 7.7 Hz, 1H), 7.31 (d, <i>J</i> = 8.1 Hz, 1H), 7.25 (d, <i>J</i> = 7.5 Hz, 1H), 6.99 (d, <i>J</i> = 7.5 Hz, 1H), 3.79 (s, 2H), 3.24 (s, 3H), 2.93 - 2.83 (m, 2H), 2.83 - 2.72 (m, 2H), 2.60 - 2.52 (m, 6H), 2.04 - 1.93 (m, 1H), 1.76 - 1.59 (m, 4H), 1.56 - 1.42 (m, 1H), 1.09 (d, <i>J</i> = 5.2 Hz, 3H), 1.01 - 0.90 (m, 1H), 0.85 (d, <i>J</i> = 6.2 Hz, 3H).

Cmpd. 10		516	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 12.49 (s, 1H), 8.29 (s, 1H), 8.10 (s, 1H), 7.75 (s, 1H), 7.59 (d, $J = 7.7$ Hz, 1H), 7.47 (t, $J = 7.7$ Hz, 1H), 7.41 - 7.35 (m, 2H), 7.30 (d, $J = 7.9$ Hz, 1H), 3.84 (s, 2H), 3.23 (s, 3H), 2.91 - 2.83 (m, 2H), 2.82 - 2.74 (m, 2H), 2.59 - 2.53 (m, 3H), 2.06 - 1.96 (m, 1H), 1.78 - 1.60 (m, 4H), 1.55 - 1.44 (m, 1H), 1.09 (d, $J = 5.1$ Hz, 3H), 0.99 - 0.91 (m, 1H), 0.85 (d, $J = 6.3$ Hz, 3H).
Cmpd. 11		518	$^1\text{H NMR}$ (400 MHz, $\text{CDCl}_3$ ) $\delta$ ppm: 8.03 (s, 1H), 7.99 (s, 1H), 7.86 (s, 1H), 7.58 (d, $J = 7.5$ Hz, 1H), 7.45 - 7.34 (m, 2H), 7.20 (s, 1H), 3.86 (s, 2H), 3.27 (s, 3H), 2.98 - 2.79 (m, 4H), 2.73 - 2.65 (m, 3H), 2.20 - 2.10 (m, 1H), 1.96 - 1.80 (m, 2H), 1.78 - 1.65 (m, 3H), 1.14 (d, $J = 5.1$ Hz, 3H), 1.03 - 0.95 (m, 1H), 0.87 (d, $J = 5.5$ Hz, 3H).
Cmpd. 12		500	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 8.25 (s, 1H), 7.98 (s, 1H), 7.72 (s, 1H), 7.57 (d, $J = 7.7$ Hz, 1H), 7.44 - 7.26 (m, 3H), 6.94 (t, $J = 9.0$ Hz, 1H), 4.23 - 4.14 (m, 1H), 3.19 (s, 3H), 2.88 - 2.75 (m, 3H), 2.51 (d, $J = 6.4$ Hz, 3H), 1.80 - 1.50 (m, 4H), 1.44 - 1.22 (m, 7H), 1.05 (d, $J = 4.8$ Hz, 3H).
Cmpd. 13		512	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 8.28 (s, 1H), 8.02 (s, 1H), 7.75 (s, 1H), 7.61 (d, $J = 7.7$ Hz, 1H), 7.44 (t, $J = 7.7$ Hz, 1H), 7.40 - 7.27 (m, 2H), 7.07 - 6.91 (m, 1H), 4.26 - 4.08 (m, 1H), 3.23 (s, 3H), 3.18 - 3.09 (m, 1H), 2.93 - 2.82 (m, 2H), 2.60 - 2.50 (m, 3H), 2.18 - 2.06 (m, 1H), 2.06 - 1.92 (m, 1H), 1.48 (dd, $J = 13.4$ , 4.4 Hz, 1H), 1.35 - 1.21 (m, 4H), 1.22 - 1.01 (m, 5H), 0.65 - 0.56 (m, 1H), 0.38 - 0.33 (m, 1H).
Cmpd. 14		500	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 8.24 (s, 1H), 7.99 (s, 1H), 7.70 (d, $J = 1.9$ Hz, 1H), 7.57 (d, $J = 7.7$ Hz, 1H), 7.43 - 7.26 (m, 3H),

			6.94 (dd, $J = 9.8, 8.2$ Hz, 1H), 4.12 - 4.02 (m, 1H), 3.19 (s, 3H), 2.89 - 2.79 (m, 2H), 2.65 - 2.47 (m, 4H), 2.45 - 2.33 (m, 1H), 2.36 - 2.23 (m, 1H), 2.04 - 1.86 (m, 2H), 1.86 - 1.50 (m, 4H), 1.34 (d, $J = 6.6$ Hz, 3H), 1.05 (d, $J = 5.1$ Hz, 3H).
Cmpd. 15		486	$^1\text{H}$ NMR (400 MHz, DMSO- $d_6$ ) $\delta$ ppm: 8.24 (s, 1H), 7.98 (s, 1H), 7.70 (s, 1H), 7.57 (d, $J = 7.7$ Hz, 1H), 7.43 - 7.26 (m, 3H), 6.94 (t, $J = 9.0$ Hz, 1H), 4.20 - 4.10 (m, 1H), 3.19 (s, 3H), 2.89 - 2.79 (m, 2H), 2.56 - 2.47 (m, 3H), 2.42 - 2.33 (m, 1H), 2.19 - 2.10 (m, 1H), 1.35 (d, $J = 6.6$ Hz, 3H), 1.05 (d, $J = 4.9$ Hz, 3H), 0.90 - 0.79 (m, 1H), 0.43 - 0.27 (m, 2H), 0.09 - 0.03 (m, 2H).
Cmpd. 16		502	$^1\text{H}$ NMR (400 MHz, DMSO- $d_6$ ) $\delta$ ppm: 12.97 (s, 1H), 8.28 (s, 1H), 8.02 (s, 1H), 7.73 (s, 1H), 7.61 (d, $J = 7.6$ Hz, 1H), 7.43 (t, $J = 7.7$ Hz, 1H), 7.40 - 7.28 (m, 2H), 7.07 - 6.91 (m, 1H), 4.18 - 4.03 (m, 1H), 3.22 (s, 3H), 2.95 - 2.80 (m, 2H), 2.62 - 2.52 (m, 3H), 2.39 - 2.30 (m, 1H), 2.14 - 2.00 (m, 1H), 1.43 (d, $J = 6.6$ Hz, 3H), 1.08 (d, $J = 5.1$ Hz, 3H), 0.89 (s, 9H).
Cmpd. 17		512	$^1\text{H}$ NMR (400 MHz, DMSO- $d_6$ ) $\delta$ ppm: 8.28 (s, 1H), 8.03 (s, 1H), 7.74 (s, 1H), 7.61 (d, $J = 7.7$ Hz, 1H), 7.48 - 7.38 (m, 2H), 7.33 (d, $J = 8.0, 1.4$ Hz, 1H), 6.98 (dd, $J = 9.9, 8.2$ Hz, 1H), 4.04 (s, 2H), 3.22 (s, 3H), 3.21 - 3.13 (m, 1H), 2.92 - 2.82 (m, 2H), 2.60 - 2.52 (m, 3H), 2.01 - 1.88 (m, 1H), 1.75 - 1.42 (m, 5H), 1.09 (d, $J = 5.0$ Hz, 3H), 0.53 - 0.42 (m, 2H), 0.41 - 0.30 (m, 2H).
Cmpd. 18		500	$^1\text{H}$ NMR (400 MHz, CDCl $_3$ ) $\delta$ ppm: 8.10 (s, 1H), 8.07 (s, 1H), 7.89 (s, 1H), 7.58 (d, $J = 7.4$ Hz, 1H), 7.48 - 7.30 (m, 2H), 7.24 (s, 1H), 6.85 (t, $J = 8.6$ Hz, 1H), 4.42 - 4.30 (m, 1H), 3.29 (s, 3H), 2.96 (d, $J = 5.9$ Hz, 2H), 2.74 - 2.62

			(m, 3H), 2.17 - 2.10 (m, 1H), 2.02 - 1.90 (m, 2H), 1.75 - 1.60 (m, 3H), 1.59 - 1.53 (m, 3H), 1.33 - 1.20 (m, 3H), 1.14 (d, $J = 5.4$ Hz, 3H)
Cmpd. 19		500	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 8.28 (s, 1H), 8.02 (s, 1H), 7.76 (s, 1H), 7.61 (d, $J = 7.7$ Hz, 1H), 7.48 - 7.29 (m, 3H), 6.98 (dd, $J = 9.8, 8.3$ Hz, 1H), 4.27 - 4.18 (m, 1H), 3.22 (s, 3H), 2.92 - 2.81 (m, 3H), 2.54 (d, $J = 6.6$ Hz, 3H), 1.83 - 1.73 (m, 1H), 1.72 - 1.53 (m, 3H), 1.50 - 1.26 (m, 7H), 1.09 (d, $J = 4.9$ Hz, 3H).
Cmpd. 20		500	$^1\text{H NMR}$ (400 MHz, $\text{DMSO-}d_6$ ) $\delta$ ppm: 8.28 (s, 1H), 8.02 (s, 1H), 7.76 (s, 1H), 7.61 (d, $J = 7.8$ Hz, 1H), 7.48 - 7.29 (m, 3H), 6.98 (dd, $J = 9.8, 8.3$ Hz, 1H), 4.27 - 4.18 (m, 1H), 3.22 (s, 3H), 2.92 - 2.82 (m, 3H), 2.54 (d, $J = 6.5$ Hz, 3H), 1.83 - 1.73 (m, 1H), 1.72 - 1.53 (m, 3H), 1.50 - 1.25 (m, 7H), 1.09 (d, $J = 4.8$ Hz, 3H).

564 **5. Supplementary Algorithms S1-S6**

---

**Algorithm S1** pre-training process

---

- 1: **procedure** LIGANDGENERATION( $D_l^\dagger, D_l$ )
- 2:    $\mathcal{L}_{PT} \leftarrow 0$
- 3:   Autoregressively predict with three ligand generation heads via teacher-forcing:
  - Next token in DSMILES
  - Atom pointer position
  - Local & global coordinates
- 4:   Compute loss components:

$$\mathcal{L}_{PT} = \mathcal{L}_{\text{token}} + \mathcal{L}_{\text{ptr}} + \mathcal{L}_{\text{pos}}$$

- 5:   **return**  $\mathcal{L}_{PT}$
- 6: **end procedure**
- 7:

**Require:** Perturbed ligand data  $D_l^\dagger$ , ligand data represented by DSMILES  $D_l$ , total epochs  $E$

**Ensure:** Model parameters

- 8: Shuffle data in dataloader
  - 9: **for** epoch  $e \leftarrow 1$  to  $E$  **do**  $\mathcal{L}_{PT} \leftarrow$  LIGANDGENERATION( $D_l^\dagger, D_l$ )
  - 10:   Backpropagate  $\mathcal{L}_{PT}$
  - 11:   Update parameters with AdamW optimizer
  - 12: **end for**
- 

565

566

---

**Algorithm S2** Multi-task training process

---

```
1: procedure LIGANDGENERATION( $D_p, D_l$ )
2:    $\mathcal{L}_\tau \leftarrow 0$ 
3:   Autoregressively predict with three ligand generation heads via teacher-
   forcing:
   • Next token in DSMILES
   • Atom pointer position
   • Local & global coordinates
4:   Compute loss components:
```

$$\mathcal{L}_\tau = \mathcal{L}_{\text{token}} + \mathcal{L}_{\text{ptr}} + \mathcal{L}_{\text{pos}}$$

```
5:   return  $\mathcal{L}_\tau$ 
6: end procedure
7:
8: procedure NCIANCHORPREDICTION( $D_p$ )
9:    $\mathcal{L}_\tau \leftarrow 0$ 
10:  Independently predict extra information with the pocket prediction head:
   • NCI sites prediction
   • Anchor sites prediction
11:  Compute loss components:
```

$$\mathcal{L}_\tau = \mathcal{L}_{\text{NCI}} + \mathcal{L}_{\text{Anchor}}$$

```
12:  return  $\mathcal{L}_\tau$ 
13: end procedure
14:
```

**Require:** Pocket data  $D_p$ , ligand data represented by DSMILES  $D_l$ , total epochs  $E$

**Ensure:** Model parameters

```
15: Shuffle data in dataloader
16: Load pre-trained or post-trained parameters    ▷ Depend on the training stage
17: for epoch  $e \leftarrow 1$  to  $E$  do
18:   for task  $\tau \in \{1, 2, 3\}$  do                                ▷ Multi-task iteration
19:     if  $\tau = 1$  then                                          ▷ Task1: NCI & Anc prediction
20:        $\mathcal{L}_\tau \leftarrow$  NCIANCHORPREDICTION( $D_p$ )
21:     else if  $\tau = 2$  then                                       ▷ Task2: MG without NCI & Anc constraints
22:        $\mathcal{L}_\tau \leftarrow$  LIGANDGENERATION( $D_p, D_l$ )
23:     else if  $\tau = 3$  then                                       ▷ Task3: MG with NCI & Anc constraints
24:        $\mathcal{L}_\tau \leftarrow$  LIGANDGENERATION( $D_p, D_l$ )
25:     end if
26:     Backpropagate  $\mathcal{L}_\tau$ 
27:     Update parameters with AdamW optimizer
28:   end for
29: end for
```

---

---

**Algorithm S3** Random fragment sequence augmentation algorithm

---

**Require:** Original fragment sequence  $FS_{\text{raw}}$ **Ensure:** Transformed fragment sequence  $FS_{\text{trans}}$ 

```
1:  $n \leftarrow |FS_{\text{raw}}|$  ▷ Obtain the number of fragments
2:  $FS_{\text{trans}} \leftarrow []$ 
3: while  $n > 0$  do
4:   if  $FS_{\text{trans}}$  is empty then
5:     Randomly select fragment  $f$  from  $FS_{\text{raw}}$ 
6:      $FS_{\text{trans}} \leftarrow$  Add  $f$  to  $FS_{\text{trans}}$ 
7:      $FS_{\text{raw}} \leftarrow$  Remove  $f$  from  $FS_{\text{raw}}$ 
8:   else
9:     Check satisfied connection constraints in  $FS_{\text{trans}}$ 
10:    Randomly pick fragment  $f$  from  $FS_{\text{raw}}$  based on satisfied constraints
11:     $f_{\text{trans}} \leftarrow$  Use asterisk constraint to convert  $f$ 
12:     $FS_{\text{trans}} \leftarrow$  Add  $f_{\text{trans}}$  to  $FS_{\text{trans}}$ 
13:     $FS_{\text{raw}} \leftarrow$  Remove  $f$  from  $FS_{\text{raw}}$ 
14:   end if
15:    $n \leftarrow n - 1$ 
16: end while
17: Verify the legality of the transformed fragment sequence  $FS_{\text{trans}}$ 
```

---

568

569

---

**Algorithm S4** Retained multiple fragments sequence augmentation algorithm

---

**Require:** Original fragment sequence  $FS_{\text{raw}}$ , number of retained fragments  $n_r$

**Ensure:** Transformed fragment sequence  $FS_{\text{trans}}$

```
1:  $n \leftarrow |FS_{\text{raw}}|$  ▷ Obtain the number of fragments
2:  $FS_{\text{trans}} \leftarrow []$ 
3: if  $n_r \geq n$  then ▷ Not meeting the need to preserve fragments
4:   return  $FS_{\text{trans}}$ 
5: end if
6:  $FS_{\text{trans}} \leftarrow []$ 
7: while  $n > 0$  do
8:   if  $FS_{\text{trans}}$  is empty then
9:     Randomly select fragment  $f$  from  $FS_{\text{raw}}$ 
10:     $FS_{\text{trans}} \leftarrow$  Add  $f$  to  $FS_{\text{trans}}$ 
11:     $FS_{\text{raw}} \leftarrow$  Remove  $f$  from  $FS_{\text{raw}}$ 
12:     $n \leftarrow n - 1$ 
13:     $n_r \leftarrow n_r - 1$ 
14:    while  $n > 0 \ \& \ n_r > 0$  do ▷ Preserve fragments
15:      Check unsatisfied connection constraints in  $FS_{\text{trans}}$ 
16:      Randomly pick fragment  $f$  from  $FS_{\text{raw}}$  based on unmet constraints
17:       $FS_{\text{trans}} \leftarrow$  Add  $f$  to  $FS_{\text{trans}}$ 
18:       $FS_{\text{raw}} \leftarrow$  Remove  $f$  from  $FS_{\text{raw}}$ 
19:       $n \leftarrow n - 1$ 
20:       $n_r \leftarrow n_r - 1$ 
21:    end while
22:  else
23:    Check satisfied connection constraints in  $FS_{\text{trans}}$ 
24:    Randomly pick fragment  $f$  from  $FS_{\text{raw}}$  based on satisfied constraints
25:     $f_{\text{trans}} \leftarrow$  Use asterisk constraint to convert  $f$ 
26:     $FS_{\text{trans}} \leftarrow$  Add  $f_{\text{trans}}$  to  $FS_{\text{trans}}$ 
27:     $FS_{\text{raw}} \leftarrow$  Remove  $f$  from  $FS_{\text{raw}}$ 
28:     $n \leftarrow n - 1$ 
29:  end if
30: end while
31: Verify the legality of the transformed fragment sequence  $FS_{\text{trans}}$ 
```

---

570

571



---

**Algorithm S6** Molecular generation process

---

```
1: procedure GENERATEONESTEP(Action( $i$ ))
2:   Predict the  $(i + 1)^{\text{th}}$  DSMILES token
3:   Predict the  $(i + 1)^{\text{th}}$  pointer using the  $(i + 1)^{\text{th}}$  DSMILES token
4:   Predict local coordinates  $(r, \theta, \phi)$  and global coordinates  $(x, y, z)$  using the
    $(i + 1)^{\text{th}}$  pointer and DSMILES token
5:   Define search space around predicted local coordinates for final token coordi-
   nates:
   • Bond length:  $r \pm 0.1\text{\AA}$ 
   • Bond Angle:  $\theta \pm 2^\circ$ 
   • Dihedral Angle:  $\phi \pm 2^\circ$ 
6:   Find global coordinate with highest joint probability within search space
7:   Set final predicted coordinates for the  $(i + 1)^{\text{th}}$  DSMILES token
8:   return  $(i + 1)^{\text{th}}$  DSMILES token, pointer, and coordinates
9: end procedure
10:
11: Set the maximum sequence length of the ligand to  $T$ 
12: if single/multiple retained fragments exist then
13:   Initialize Action ( $s$ ) based on single/multiple retained fragments
14:   Set the starting step  $i_{\text{start}} \leftarrow s$ 
15: else
16:   Initialize Action (0) with  $\langle \text{SOS} \rangle$ 
17:   Set the starting step  $i_{\text{start}} \leftarrow 0$ 
18: end if
19: for each generation step  $i$  in  $[i_{\text{start}}, T - 1]$  do
20:    $\text{DT}_{i+1}, \text{Ptr}_{i+1}, \text{Pos}_{i+1} \leftarrow \text{GENERATEONESTEP}(\text{Action}(i))$  ▷ Core step
21:   Update Action ( $i + 1$ ) ▷ Update ligand embeddings
22:   if  $\text{DT}_{i+1}$  is  $\langle \text{EOS} \rangle$  then
23:     break
24:   end if
25: end for
```

---

574

575

576 **Reference**

- 577 1. Fan, F. *et al.* Assessing conformation validity and rationality of deep learning-generated  
578 3D molecules. *bioRxiv* 2024–11 (2024).
- 579 2. Anstine, D. M., Zubatyuk, R. & Isayev, O. AIMNet2: A neural network potential to  
580 meet your neutral, charged, organic, and elemental-organic needs. *Chem. Sci.* **16**,  
581 10228–10244 (2025).
- 582 3. Bouysset, C. & Fiorucci, S. ProLIF: a library to encode molecular interactions as  
583 fingerprints. *Journal of cheminformatics* **13**, 72 (2021).
- 584 4. Greg Landrum *et al.* rdkit/rdkit: 2025\_09\_2 (Q3 2025) release. Zenodo  
585 <https://doi.org/10.5281/ZENODO.591637> (2025).
- 586 5. He, X., Man, V. H., Yang, W., Lee, T.-S. & Wang, J. A fast and high-quality charge  
587 model for the next generation general AMBER force field. *The Journal of Chemical*  
588 *Physics* **153**, 114502 (2020).
- 589 6. Koes, D. R., Baumgartner, M. P. & Camacho, C. J. Lessons learned in empirical scoring  
590 with smina from the CSAR 2011 benchmarking exercise. *J. Chem. Inf. Model.* **53**, 1893–  
591 1904 (2013).
- 592 7. Groom, C. R., Bruno, I. J., Lightfoot, M. P. & Ward, S. C. The cambridge structural  
593 database. *Acta Crystallogr B Struct Sci Cryst Eng Mater* **72**, 171–179 (2016).
- 594 8. Rogers, D. & Hahn, M. Extended-connectivity fingerprints. *Journal of chemical*  
595 *information and modeling* **50**, 742–754 (2010).
- 596