

An interface-based molecular generative framework for protein-protein interaction inhibitors

Jianmin Wang¹, Jiashun Mao¹, Chunyan Li⁴, Hongxin Xiang², Xun Wang^{3,5}, Shuang Wang³, Zixu Wang⁶, Yangyang Chen⁶, Yuquan Li⁷, Heqi Sun⁸, Kyoung Tai No^{1,*}, Tao Song^{3,*}, Xiangxiang Zeng^{2,*}

¹ The Interdisciplinary Graduate Program in Integrative Biotechnology and Translational Medicine, Yonsei University, Incheon 21983, Republic of Korea

² College of Computer Science and Electronic Engineering, Hunan University, Changsha, Hunan, 410082, China

³ School of Computer Science and Technology, China University of Petroleum, Qingdao, 266580, Shandong, China

⁴ School of Informatics, Yunnan Normal University, Kunming, China

⁵ High Performance Computer Research Center, University of Chinese Academy of Sciences, Beijing, 100190, China

⁶ Department of Computer Science, University of Tsukuba, Tsukuba, 3058577, Japan

⁷ College of Chemistry and Chemical Engineering, Lanzhou University, Lanzhou, China

⁸ School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, Shanghai, 200240, China

*To whom correspondence should be addressed: Xiangxiang Zeng, Ph.D.

College of Computer Science and Electronic Engineering, Hunan University, Changsha, Hunan, China

Email: xzeng@hnu.edu.cn

29 **Abstract**

30 Protein-protein interactions (PPIs) play a crucial role in many biochemical processes
31 and biological processes. Recently, many structure-based molecular generative models
32 have been proposed. However, PPI sites and compounds targeting PPIs have
33 distinguished physicochemical properties compared to traditional binding pockets and
34 drugs, it is still a challenging task to generate compounds targeting PPIs by considering
35 PPI complexes or interface hotspot residues. In this work, we propose a specifically
36 molecular generative framework based on PPI interfaces, named GENiPPI. We
37 evaluated the framework and found it can capture the implicit relationship between the
38 PPI interface and the active molecules, and can generate novel compounds that target
39 the PPI interface. Furthermore, the framework is able to generate diverse novel
40 compounds with limited PPI interface inhibitors. The results show that PPI interface-
41 based molecular generative model enriches structure-based molecular generative
42 models and facilitates the design of inhibitors based on PPI structures.

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57 **Main**

58 A vast network of genes is inter-linked through protein-protein interactions and is
59 critical component of almost every biological process under physiological conditions,
60 and can be ubiquitous in many living organisms and biological pathways¹⁻⁴. Modulation
61 of PPIs expands the drug target space and has enormous potential in drug discovery. In
62 homo sapiens, it is estimated that the entire interactome comprises between 130,000 to
63 930,000 binary PPIs⁵⁻⁷. Despite significant efforts in developing modulators of PPIs,
64 drug design and development for PPI targets, especially targeting the PPI interfaces,
65 remains challenging^{6,8-12}. Structure-based rational design serves as an important tool
66 for the discovery of lead compounds in drug discovery¹³⁻¹⁷. Traditional drug targets and
67 PPIs targets have different bio-chemical features (**Table 1**)^{11,18-22}, so conventional drugs
68 and PPIs inhibitors have different physicochemical properties and drug-like properties
69 (**Table 1**)^{11,23-30}. Given their differences, developing molecular generative models of
70 different paradigms are essential for the drug design of different target types^{11,19,31}.

71
72 Generative artificial intelligence(AI) is enable to model the distribution of training
73 samples and generate novel samples^{32,33}. In drug discovery, generative AI can accelerate
74 drug discovery by generating novel molecules with desired properties. Numerous
75 excellent review articles have summarized the development in this field^{16,17,34-41}.
76 Molecular generative models in drug design can be roughly divided into three
77 categories: ligand-based molecular generative (LBMG) models, structure-
78 based(pockets or binding sites) molecular generative models (SBMG), and fragment-
79 based molecular generative models (FBMG), among which SBMG models have
80 received much attention.^{17,39,42}. Currently, some significant methods in structure-based
81 molecular generative models can be found in⁴³⁻⁵¹, molecular generative models for PPI
82 structures or PPI interfaces have been rarely reported in the literature. In recent years,
83 classical machine learning⁵²⁻⁵⁴, active learning⁵⁵, and deep learning-assisted methods⁵⁶
84 is better screening and design of PPIs inhibitors have been explored, and ligand-based
85 molecular generative models of PPI inhibitors have been reported⁵⁷. There are few

86 structure-based molecular generative models for PPIs targets have not been sufficiently
87 explored.

88

89 In this study, we developed a conditional molecular generative framework based on
90 protein-protein interaction interfaces (named GENiPPI) for the design of PPI interface
91 inhibitors. The framework was developed by a conditional Wasserstein generative
92 adversarial network (cWGAN) with convolutional neural networks (CNNs), integrated
93 graph attention networks (GATs) and long short-term memory (LSTM). It was designed
94 to efficiently capture the relationship between PPI interface with active/inactive
95 compounds to train conditional molecular generative models (**Fig. 1**). As demonstrated
96 by the conditional evaluation, GENiPPI is an effective architecture for capturing the
97 implicit relationships between the PPI interface and active compounds. In summary,
98 GENiPPI represents a potent deep learning framework for structure-based design of
99 PPI inhibitors.

100

101 **Results**

102 **Generation of molecules targeting the PPI interface**

103 Here, we introduce GENiPPI a modular deep learning framework for the design of
104 structure-based PPIs inhibitors (**Fig. 1**). GENiPPI is composed of four main modules:
105 GATs module⁵⁸⁻⁶⁰ for representation learning of the protein complex interface, CNNs
106 module for molecular representation learning, cWGAN module⁶¹ for conditional
107 molecular generation, and molecular captioning network module for SMILES strings
108 decoding (**as shown in Supplementary Figs.1, Figs.2, Figs.3 and Figs.4,**
109 **respectively**).

110

111 Our framework undergoes four steps to accomplish the generation of molecules
112 targeting the PPI interface. In the first step, we use GATs module designed for the
113 protein complex interface is to effectively capture the nuanced atomic-level interaction
114 characteristics inherent to the protein complex interface region. Next, we use CNN

115 module to provide a representation of the compound that contains voxel and electronic
116 density information in three-dimension space ⁶². And, the cGAN module is designed
117 to generate compounds that target PPI interfaces using features from the protein
118 complex interface region to regulate the inputs ⁶³. The cGAN module consists of a
119 generator, a discriminator, and a conditional network. The generator takes a Gaussian
120 random noise vector, and the protein complex interface features to generate a vector in
121 the molecular embedding space, the discriminator evaluates whether the generated
122 molecule embedding corresponds to a real or generated molecule, and the conditional
123 network evaluates whether the molecule embedding matches the protein complex
124 interface features. Finally, we use the molecular captioning network, which is made by
125 a 3D CNNs and a recurrent LSTM ⁶⁴to decode molecular representations. The
126 molecular representation generated by the generator is fed as input to the 3D
127 convolutional network with the LSTM subsequently decoding the SMILES strings.

128

129 **Conditional evaluation**

130 First, we verified the validity of the conditions that act as conditional molecular
131 generative models for the protein complex interfaces. For this purpose, we selected
132 three PPI targets: MDM2(mouse double minute 2)/p53, Bcl-2(B-cell lymphoma 2)/Bax
133 (Bcl-2 associated X), and BAZ2B(Bromodomain adjacent to zinc finger domain protein
134 2B)/H4(histone) for conditional evaluation. We generated 10,000 validated molecules
135 each by the GENiPPI framework and calculated the drug-like metrics of the generated
136 compounds: QED²⁷, QEPPI^{28,29} and Fsp3(fraction of sp3 carbon atoms)⁶⁵. We
137 compared the QED, QEPPI, and Fsp3 distributions of the active compounds and the
138 generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4 (**Fig. 2**). As shown,
139 the distributions of drug-like properties were similar between the generated compounds
140 and the active compounds for the three PPI interface targets (**Fig.2a, Fig.2b, and**
141 **Fig.2c**), while different distributions of drug-like properties were observed between the
142 generated compounds based on different targets (**Fig.2d, Fig.2e, and Fig.2f**). The
143 results demonstrate the effectiveness of the PPI interface in conditioning the molecular

144 generative model. The drug-like properties of the framework generated compounds
145 migrate relative to those of the compounds in the training dataset, indicating that the
146 framework captures the distributions of the training dataset and generates novel
147 compounds.

148

149 **Model performance**

150 In order to gain insight into the performance of the GENiPPI framework and to compare
151 it with other molecular generative models. We benchmarked our method by the MOSES
152 platform⁶⁶, a leading benchmark platform of molecular generation. We trained all
153 models on the full training dataset and randomly sampled 30,000 molecules. We
154 utilized models and hyperparameters provided by the MOSES platform, such as an
155 Adversarial Autoencoder(AAE)⁶⁷, character-level recurrent neural networks
156 (CharRNN)⁶⁸, Variational Autoencoder(VAE)⁶⁹, LatentGAN⁷⁰ and ORGAN⁷¹. To
157 validate the higher quality of the molecules generated by the conditioned model, we
158 compared them with molecules sampled from the GENiPPI framework and the
159 GENiPPI-noninterface framework without the conditioned module. We found that
160 molecules generated by the conditioned GENiPPI framework were superior to other
161 models in novelty and diversity.

162

163 As shown in **Table 2**, the GENiPPI framework has advantages in terms of uniqueness,
164 novelty, and diversity over the GENiPPI-noninterface. The GENiPPI framework
165 performs better overall in molecular generation. Compared with LatentGAN and
166 ORGAN, GENiPPI offers more benefits in terms of validity and diversity. While all
167 molecular generative models have their unique advantages in various performance
168 comparison. However, the molecular generative models tailored to specific tasks,
169 especially those based on PPI structure, have more advantages and inspirations from
170 the GENiPPI framework. To understand the similarities and differences between the
171 molecular distributions generated by the GENiPPI framework and other models. We
172 compared the distribution of molecular properties of the Testset, iPPI-DB inhibitor, and

173 the generated molecular datasets of AAE, CharRNN, VAE, LatentGAN,
174 GENiPPI(noninterface) and GENiPPI(**Supplementary Figs.4**). The generated
175 compounds have similar distributions of physicochemical properties to the compounds
176 from the training set. While most of the iPPI-DB inhibitors have QED values lower
177 than 0.5, most of them have QEPPI values higher than 0.5.

178

179 **Chemical space exploration**

180 To better obtain an estimate of the chemical space distribution of the model generated
181 molecules with the active compounds in the training datasets, we evaluated the
182 chemical drug-like space of the generated compounds by calculating t-distributed
183 random neighbourhood embedding (t-SNE) maps of MACCS fingerprint⁷². The t-SNE
184 is a dimensionality reduction method used for data points visualization in two or three-
185 dimensional space by mapping high-dimensional data to a lower dimension^{73,74}. By
186 this method, similar compounds are clustered to visualize the high-dimensional
187 chemical space of the compounds. The distribution of the generated compounds and
188 active compounds in chemical drug-like space by t-SNE visualization (**Fig.3a, Fig.3b,**
189 **and Fig.3c**). The generated drug-like compounds not only share the chemical space
190 with the active compounds, but are also homogeneously mixed in the two-dimensional
191 space. The generated compounds show a similar chemical drug-like space to that of the
192 active compounds under 2D topological fingerprint. Adding the three dimensions of
193 compounds contributes to the design of promising drug-like compounds^{30,75}. We
194 performed PMI shape analysis on the generated compounds and compared them with
195 drug-like compounds from DrugBank and iPPI-DB(**Fig.3d**). Many of the approved
196 compounds are either rod or disk shaped, and the generated drug-like compounds
197 library has a similar three-dimensional space. The PBF distribution of the library of
198 generated drug-like compounds is about 0~2 Å(**Fig.3e**). The results show that many of
199 the generated drug-like compounds are derived from relatively planar molecular
200 scaffolds. Moreover, we evaluated the ability of the model to generate target-specific
201 compounds by chemical space maps. To assess the overlapping of drug-like chemical

202 space, we utilized Tree MAP (TMAP)⁷⁶ to create the 2D projection(**Fig.3f**). Each point
203 corresponds to a compound and is colored by its target label. The dark and light colors
204 denote the generated compounds and the active compounds in the training set. These
205 results suggest that our GENiPPI model can generate compounds that are similar to the
206 active compounds in the training set and have novel structures. The results show that
207 the framework enriches and expands the chemical space of PPI-targeted drug-likeness
208 compounds.

209

210 **Few-shot molecular generation**

211 Because of the huge consumptive costs involved in data collection, only a small amount
212 of labeled biomedical data are usually available. The process of drug design and
213 optimization often faces the problem of low data⁷⁷. The lack of effectively labeled data
214 tends to diminish the practical performance of most deep learning frameworks for drug
215 design. To perform generalized molecular generative design with limited labeled data,
216 it has been a trending topic in the few-shot generative community^{78,79}. The GENiPPI
217 model was applied to generate a virtual compound library for the heat shock protein 90
218 - cell division cycle 37(Hsp90-Cdc37) interaction interface. By training the model on
219 the PPI structure of Hsp90-Cdc37 (PDB ID: 1US7) and seven disruptors, we sampled
220 500 valid compounds. The similarity between active disruptors and generated
221 compounds of Hsp90/Cdc37 in the chemical space was visualized by t-SNE projection
222 maps(**Fig.4a**). After few-shot learning, the generated compounds were mostly
223 distributed around the active disruptor, which demonstrated the effectiveness of few-
224 shot learning in navigating through the targeted chemical space. We performed
225 pharmacophore-based matching by considering DCZ3112(a novel triazine derivative
226 that disrupts Hsp90-Cdc37 interactions) as a reference molecule⁸⁰. The top 5 generated
227 molecules have similar pharmacophore and shape features with DCZ3112(**Fig.4b**),
228 demonstrating the potential of the model to be applied to low-data PPI targets. **Fig. 4c**,
229 shows the hot spot amino acid residues at the PPI interface of the Hsp90-Cdc37 protein
230 complex(PDB ID: 1US7). We performed molecular docking for prediction of the

231 binding poses(**Fig.4e**) of DCZ3112 with the Hsp90-Cdc37 complex by the UCSF
232 DOCK6.9 program⁸¹. The structure of the Hsp90-Cdc37 complex with DCZ3112
233 highlights the hydrogen bond interactions with amino acid residues:Arg32A, Glu33A,
234 Ser36A , Ser115A, Gly118A, Gln119A, and Arg167B(**Fig.4e**), which may be the major
235 energy contributors to protein-ligand interactions. The generated compounds were
236 performed molecular docking together with DCZ3112, and selected compounds with
237 reasonable binding modes and higher binding affinity by visual inspection for
238 interaction pattern analysis. The generated compounds of GENiPPI not only obtained
239 the better docking score than the active compounds, but also reproduced the interactions
240 with the key residues of the PPI interface. The generated compounds also formed
241 halogen bonds, salt bridges and π -cation interactions to improve the binding affinity of
242 the generated compounds to the target interface(**Fig.4f**). In conclusion, by analyzing
243 the interaction patterns between the generated compounds and the PPI interface,
244 GENiPPI learned the implicit interaction rules between the active compounds and the
245 PPI interface.

246

247 **Discussion**

248 We developed the GENiPPI framework, which combines protein-protein interaction
249 (PPI) interfaces features and conditional molecular generative model to generate novel
250 modulators for PPI interfaces. We validated the ability of GENiPPI framework to learn
251 the implicit relationship between PPI interfaces and active molecules through
252 conditional evaluation experiments. GENiPPI used GATs to extract key features of PPI
253 interfaces, and searched molecules by conditional wGAN with specific constraints. We
254 compared GENiPPI with various evaluation settings and benchmarks to demonstrate
255 its practical potential.

256

257 Despite the promising results, our framework has some limitations that can be
258 addressed in future work. We have not tested the model on a large number of receptor-
259 ligand pairs of PPIs, which may affect its the generalization ability. The reason is that

260 the current PPI has relatively little data of drug-PPI target complexes than the traditional
261 dataset of drug-target complexes. Furthermore, the current framework does not
262 incorporate the 3D structural information of ligand-receptor interactions of PPIs. and
263 There are still many ways to improve representation learning, balance training speed of
264 molecular generative models and the diversity of generated molecules. Several
265 potential directions could further improve GENiPPI: (1) collecting and cleaning higher
266 quality data pairs for model development and testing; (2) fusing of molecular chemical
267 language models and pre-trained models of protein-protein structural features to fine-
268 tune the datasets of receptor-ligand of PPIs to enhance the model generalization,
269 novelty and diversity of the generated compounds; (3) incorporating structural
270 information of PPIs into fragment-based molecular generative models is also a
271 promising direction; (4) change the architecture of the model or combining it with deep
272 reinforcement learning to generate novel compounds with better binding affinity.
273 Therefore, we will collect more data to further develop an enhanced version of
274 GENiPPI by combining novel representation learning methods and deep generative
275 approaches. In summary, the GENiPPI framework brings encouraging advances in PPI
276 structure-based molecular generative tasks and presents a tool for rational drug design
277 in finding modulators of macromolecule-macromolecule interactions.

278

279 **Methods**

280 **Datasets**

281 We first investigated PPI targets that were annotated with sufficient compound
282 bioactivity data for training and evaluation of our model ⁸². For this study, we selected
283 10 validated PPI drug targets that cover the binding interface(**Supplementary Table**).
284 These targets are E3 ubiquitin-protein ligase Mdm2, apoptosis regulator Bcl-2, BAZ2B,
285 apoptosis regulator Bcl-xL, BRD4 bromodomain 1 BRD4-1, CREB-binding protein
286 (CREBBP), ephrin type-A receptor 4 (EphA4), induced myeloid leukemia cell
287 differentiation protein Mcl-1, and menin. In addition, we randomly selected a subset of
288 250,000 compounds as additional inactive compounds from the ChEMBL ⁸³ dataset

289 that was used as part of the training datasets. A detailed data preprocessing can be found
290 in Supplementary Note A.

291

292 **Model strategy and training**

293 **Graph attention networks of protein-protein interaction interface**

294 In this section, the representation learning of protein-protein complexes interfaces is
295 inspired by the work in protein docking model evaluation⁶⁰, which designed a double-
296 graph representation to capture the interface features and interactions of protein-protein
297 complexes (**Supplementary Figs.1**). The extracted interface region is constructed as
298 two graphs (G^1 and G^2) for representing the interfacial information and the residues
299 involved in the two proteins participating in the interaction. A graph G can be defined
300 as $G = (V, E, \text{ and } A)$, where V is the set of nodes, and E is the set of edges between
301 them, and A is the adjacency matrix for mapping the association between the nodes of
302 the graph, which numerically denotes the connectivity of the graph. If the graph G has
303 N nodes, the dimension of the adjacency matrix A of the graph is $N * N$, where $A_{ij} > 0$
304 if the i -th node is connected to the j -th node, and $A_{ij} = 0$ otherwise. The graph G^1
305 describes the coding of the atomic types of all residues in the interface region, and its
306 adjacency matrix A^1 describes the classification of interatomic bonding types for all
307 residues at the interface region, which only considers the covalent bonds between atoms
308 of interface residues within each subunit as edges. Therefore, it is defined as follows:

$$309 \quad A_{ij}^1 = \begin{cases} 1 & \text{if atom } i \text{ and atom } j \text{ are connected by a covalent bond or if } i=j \\ 0 & \text{otherwise} \end{cases}$$

310 The graph G^2 links both covalent bonds (thus including G^1) and non-covalent residue
311 interactions as edges. The adjacency matrix A^2 for G^2 describes both covalent bonds
312 and non-covalent interactions between atoms within the range of 10.0 Å to each other.
313 The non-covalent atom pairs are defined as those which are closer than 10.0 Å to each
314 other. It is defined as follows:

$$A_{ij}^2 = \begin{cases} A_{ij}^1, & \text{if } i, j \in \text{receptor or } i, j \in \text{ligand} \\ e^{-\frac{(d_{ij}-\mu)^2}{\sigma}}, & \text{if } d_{ij} \leq 10\text{\AA} \text{ and } i \in \text{receptor and } j \in \text{ligand;} \\ & \text{or if } d_j \leq 10\text{\AA} \text{ and } j \in \text{receptor and } i \in \text{ligand} \\ 0, & \text{otherwise} \end{cases}$$

where d_{ij} represents the distance between the i -th and the j -th atoms of all atoms of all residues in the interaction region. μ and σ are learnable parameter with initial values of 0.0 and 1.0, respectively. The formula $e^{-(d_{ij}-\mu)^2/\sigma}$ decays with increasing distance between atoms.

320

The graph representation is more flexible and natural to encode interactive information and adjacent(local) relationships. For the node features of the graph, we considered the physicochemical properties of the atoms. We used the same features from the previous work^{60,84,85}. Then, the feature vector of the nodes is 23 in length and was embedded into 140 features by a one-layer fully connected (FC) network.

326

The constructed graphs are used as the input for GATs. The graph consists of adjacency matrices A^1 , A^2 , node matrices N_{mn}^1 , N_{pq}^2 , and the node features, $x^{in} = \{x_1^{in}, x_2^{in}, \dots, x_N^{in}\}$ and $x \in \mathbb{R}^F$, where F is the dimensionality of the node features. For the input graph of x^{in} , the pure graph attention coefficients are defined as follows, which represent the relative importance between the i -th and j -th nodes:

332

$$e_{ij} = x_i^T E x_j' + x_j^T E x_i',$$

where x_i' and x_j' are the transformed feature representations defined by $x_i' = W x_i^{in}$ and

$x_j' = W x_j^{in}$. W , $E \in \mathbb{R}^{F \times F}$ are learnable matrices in the GATs. e_{ij} and e_{ji} become

identical to satisfy the symmetrical property of the graph by adding $x_i^T E x_j^T$ and $x_i^T E x_j'$.

The coefficient will only be computed for i and j where $A_{ij} > 0$.

337

The attention coefficients will also be calculated for the elements in the adjacency

339 matrix. For the elements (i, j) , they are defined in the following form:

$$340 \quad a_{ij} = \frac{\exp(e_{ij})}{\sum_{j \in N_i} \exp(e_{ij})} A_{ij},$$

341 where a_{ij} represents the normalized attention coefficient between the i -th and j -th node

342 pairs, while e_{ij} is the computed symmetric graph attention coefficient. N_i denotes the

343 set of neighbors for the i -th node, which includes the interacting node j with $A_{ij} > 0$.

344 The purpose here is to define attention by considering both the physical structure A_{ij}

345 and the normalized attention coefficient e_{ij} of the interactions simultaneously.

346

347 Based on the attention mechanism, the new node features of each node are updated in

348 consideration of its neighboring nodes, which is a linear combination of the neighboring

349 node features and the final attention coefficient a_{ij} :

$$350 \quad x_i'' = \sum_{j \in N_i} a_{ij} x_j',$$

351 Making the use of the GATs mechanism described previously, we applied four layers

352 of GATs to process the node embedding information of the neighboring nodes and

353 output the updated node embedding. For two adjacency matrices A^1 and A^2 , we use a

354 shared GAT. the initial input to the network is the atomic feature. Working with two

355 matrices A^1 and A^2 , we have $x_1 = \text{GAT}(x^{in}, A^1)$ and $x_2 = \text{GAT}(x^{in}, A^2)$. In order to

356 focus only on the intermolecular interactions at the interface of the input protein-protein

357 complex, we obtain the final node embedding by subtracting the embeddings of the two

358 graphs. By subtracting the updated embedding x_1 from x_2 , we can capture aggregated

359 information on intermolecular interactions from only the other nodes in the protein-

360 protein complex interface. The output node feature is therefore defined as:

$$361 \quad x^{out} = x^2 - x^1,$$

362 After that, the updated x^{out} became x^{in} to iteratively increase the information through

363 the three following GATs layers. After the four GATs layers updated the node

364 embeddings, the node embedding of the entire graph was summed up as the overall
365 intermolecular interaction representation of the protein-protein complex:

$$366 \quad x_{graph} = \sum_{k \in G} x_k.$$

367 Finally, the FC layers were applied to the x_{graph} to obtain a [4,4,4] vector as features
368 of the protein-protein interface.

369

370 **Molecular representation**

371 For each SMILES string, a 3D conformer is generated using RDKit⁸⁶ and optimized
372 using the default settings of the MMFF94 force field. The molecular structure
373 information is then extracted into a 35Å grid centered at the geometric center of the
374 molecule using the HTMD package⁸⁷. The atoms of the molecule are discretized into a
375 1 Å cubic grid, and eight channels are considered to compute voxelized information.
376 Finally, the electronic density of the molecules 9th channel is calculated using the
377 original molecule method in Multiwfn(Supplementary Figs.2)⁸⁸.

378

379 **Conditional Wasserstein generative adversarial networks**

380 The generator takes a conditional vector and a noise vector sampled from a Gaussian
381 distribution as inputs. The PPI interface features([1,4,4,4], vector shape) are
382 concatenated with a noise vector of size [9, 4, 4, 4] and input to a 4-layer transposed
383 convolutional neural network (CNNs) with 256, 512, 1024, and 1024 filters,
384 respectively. The first three layers downsample the array size using concatenation
385 convolution (s=2). For all convolutions, we use a kernel size of 4, and the Leaky ReLU
386 is used as an activation function after convolution. BatchNorm3d is applied between
387 convolution and activation operations to normalize the values of each channel of each
388 sample.

389

390 The discriminator consists of a 4-layer sequential convolutional neural network (CNNs)
391 with 256, 512, 1024, and 1024 filters, respectively. The first three layers downsample
392 the array size using concatenation convolution (s=2). For all convolutions, we use a

393 kernel size of 4, and the Leaky ReLU ($\alpha=0.2$) is used as an activation function after
394 convolution. InstanceNorm3d is applied between convolution and activation operations
395 to normalize the values of each channel of each sample.

396

397 The physical and spatial features of the compounds are derived from the molecular
398 representation learning module, and the PPI interface features are obtained from the
399 GATs module of the protein complex interface. They are used to estimate the matching
400 probability between molecules and PPI interface features(**Supplementary Figs.3**).

401

402 **Molecular captioning network**

403 In this section, we will describe how to decode the generated molecular representation
404 into a SMILES strings. Our work is inspired by shape-based molecular generation^{89,90},
405 which designs a combination network of convolutional neural networks (CNNs) and
406 Long Short-Term Memory (LSTM)⁶⁴ to generate SMILES strings. In brief, the
407 molecular captioning network consists of a 3D CNNs and a recurrent LSTM. The
408 molecular representation generated by the generator is fed as input to the 3D CNNs,
409 and the output of the 3D CNN is fed into the LSTM to decode the SMILES strings
410 (**Supplementary Figs.4**).

411

412 **Model training**

413 The conditional generative adversarial network is trained with Wasserstein loss. The
414 loss functions for the generator ($G_{(0(z,c))}$) and discriminator ($D_0(x)$) are:

415

$$416 \begin{aligned} L_{x_0} &= E_{i y_{xx}} [-D_y(x)] + E_{z_{xx}, i y_{yx}} [D_{y y_{yx}}(G_{zy}(z, c))] + \lambda E_{i y_1} [(\|\nabla_{zx} D_y(\hat{x})\|_z - 1)^2], \\ I_{xxx} &= E_{z_{xx}, i x_{yx}} [-D_z(G_{zy}(z, c)) - \alpha \log(f_u(G_u(z, c), c))] \end{aligned}$$

417

418 where x and c are molecular representations and PPI interface features, respectively,
419 sampled from the true data distribution p_{real} , z is a random noise vector sampled from
420 a Gaussian distribution (p_z), and f_0 is a function that evaluates the probability that a

421 PPI interface feature corresponds to a molecular representation. λ and α terms are
422 regularization parameters, both empirically set to 10. λ term weighs the effect of the
423 gradient penalty on discriminator loss. α term weighs the effect of the effect of f_0 on
424 the loss of the generator.

425

426 The model was trained for 50,000 calendar hours with a batch size of 8 (65 steps per
427 calendar hour). The discriminators were updated after each step, while the generators
428 were updated every 30 steps. The network was trained using the RMSprop optimizer
429 with a learning rate of 1×10^{-4} for the generator and discriminator. during training, we
430 monitored the similarity between real and generated molecular representations using
431 Fréchet distances. The weights of the conditional networks were pre-trained on a binary
432 cross-entropy loss and frozen during GAN training. Training was performed on a single
433 NVIDIA A40 GPU, and all neural networks were built and trained using Pytorch 1.7.1
434 ⁹¹ and Tensorflow 2.5 ⁹².

435

436 **Molecular generation**

437 After the model has been trained, the embedding information of the protein-protein
438 complex interface is used to guide the model to generate novel molecules from the
439 latent space. The maximum sampling strategy was used in the LSTM, meaning that the
440 SMILES strings are generated by selecting the next token based on the highest
441 prediction probability⁸⁹.

442

443 **Evaluation settings**

444 **Conditional evaluation metrics**

445 In this study, the key is to evaluate the effectiveness of the proposed framework of
446 protein-protein interaction interface-based conditional molecular generation. Therefore,
447 we sampled the same number of valid molecules for the three PPI targets. Then we
448 calculated the QED values and Fsp3 by RDKit⁸⁶ and calculated the QEPPI values by
449 the QEPPI package (<https://github.com/ohuelab/QEPPI>) for the generated compounds

450 and others, and plotted the density distribution for comparing the differences of drug-
451 likeness.

452

453 **MOSES evaluation metrics**

454 To evaluate the performance of our proposed conditional molecule generation
455 framework, we used the evaluation metrics of validity, uniqueness, novelty and
456 diversity provided by the MOSES platform⁶⁶, which are defined as follows:

457 Validity: Molecules defined as valid in the generated molecules.

$$458 \quad \text{Validity} = \frac{N_{\text{valid}}}{N_{\text{generated}}}$$

459 Uniqueness: The proportion of unique molecules found among the generated valid
460 molecules.

$$461 \quad \text{Uniqueness} = \frac{N_{\text{unique}}}{N_{\text{valid}}}$$

462 Novelty: The generated molecules are not to be covered in the training set.

$$463 \quad \text{Novelty} = \frac{N_{\text{novel}}}{N_{\text{unique}}}$$

464 FCD(Fréchet ChemNet Distance): To detect whether the generated molecules are
465 diverse and whether they have chemical and biological properties that are similar with
466 the real molecules⁹³.

467

468 **Molecular shape**

469 To evaluate the shape space of molecules, we used two widely adopted molecular
470 descriptors to represent the three dimensions of molecular structure: principal moment
471 of inertia (PMI)⁹⁴ and the best-fit plane (PBF)⁹⁵. The PMI descriptor classifies the
472 geometric shape of molecules into the degree of rod-shaped (linear shape, such as
473 acetylene), disk-shaped (planar shape, such as benzene), and sphere (spherical shape,
474 such as adamantane). The normalized PMI ratios (NPRs) are plotted in two-
475 dimensional triangle and then used to compare the shape space covered by different sets
476 of molecules, evaluating and visualizing the diversity of the molecular shape associated
477 with a given set of molecules³⁰. PBF is a three-dimensional descriptor that represents

478 the deviation of a molecule from a plane. The PBF descriptor is the mean distance of
479 each heavy atom from the best-fit plane passing via all heavy atoms⁹⁵.

480

481 **Tree MAP**

482 To explore and explain the chemical space by unsupervised visualization of high-
483 dimensional data⁷⁶, we calculated MinHash fingerprint⁹⁶ vectors for active compounds
484 and generated compounds. Then tmap⁷⁶ and faerun⁹⁷ were utilized to construct two-
485 dimensional projections of Tree MAP (TMAP).

486

487 **Protocol for few-shot generation**

488 Targeting the Hsp90-Cdc37 PPI interface is recognized as an important option for
489 cancer therapy. The crystal structure of the Hsp90-Cdc37 protein complex (PDB ID:
490 1US7) is available for molecular docking⁹⁸. In addition, known Hsp90-Cdc37 PPI
491 disruptors were collected for training of few-shot generative. They are DCZ3112,
492 Celastrol , FW-04-804, Sulforaphane, Withaferin A, Platycodin D, Kongensin A⁹⁹.
493 OpenPharmacophore(<https://github.com/uibcdf/OpenPharmacophore>) was utilized to
494 create pharmacophore models and virtual screening. The protein structures were
495 processed by using UCSF Chimera¹⁰⁰, the program DOCK6.9 was used for semiflexible
496 docking⁸¹, and PyMOL¹⁰¹ was used to create the figures. A detailed docking protocol
497 can be found in Supplementary Note B.

498

499 **Data Availability**

500 The datasets are available at Github (<https://github.com/AspirinCode/GENiPPI>). The
501 data implementation will be provided upon acceptance of the manuscript for publication.

502

503 **Code Availability**

504 All the codes are freely available at Github (<https://github.com/AspirinCode/GENiPPI>).
505 The code implementation will be provided upon acceptance of the manuscript for
506 publication.

507

508 **Acknowledgments**

509

510 **Ethics declarations**

511 The authors declare no competing interests.

512

513 **Competing Interests Statement**

514 The authors have declared no competing interests.

515

516 **References**

- 517 1 Stelzl, U. *et al.* A human protein-protein interaction network: a resource for annotating the
518 proteome. *Cell* **122**, 957-968 (2005).
- 519 2 Rual, J.-F. *et al.* Towards a proteome-scale map of the human protein-protein interaction
520 network. *Nature* **437**, 1173-1178 (2005).
- 521 3 Titeca, K., Lemmens, I., Tavernier, J. & Eyckerman, S. Discovering cellular protein-protein
522 interactions: Technological strategies and opportunities. *Mass spectrometry reviews* **38**, 79-111
523 (2019).
- 524 4 Rhys, G. G. *et al.* De novo designed peptides for cellular delivery and subcellular localisation.
525 *Nature Chemical Biology* **18**, 999-1004 (2022).
- 526 5 Venkatesan, K. *et al.* An empirical framework for binary interactome mapping. *Nature methods*
527 **6**, 83-90 (2009).
- 528 6 Nero, T. L., Morton, C. J., Holien, J. K., Wielens, J. & Parker, M. W. Oncogenic protein
529 interfaces: small molecules, big challenges. *Nature Reviews Cancer* **14**, 248-262 (2014).
- 530 7 Oughtred, R. *et al.* The BioGRID database: A comprehensive biomedical resource of curated
531 protein, genetic, and chemical interactions. *Protein Science* **30**, 187-200 (2021).
- 532 8 Wells, J. A. & McClendon, C. L. Reaching for high-hanging fruit in drug discovery at protein-
533 protein interfaces. *Nature* **450**, 1001-1009 (2007).
- 534 9 Ivanov, A. A., Khuri, F. R. & Fu, H. Targeting protein-protein interactions as an anticancer
535 strategy. *Trends in pharmacological sciences* **34**, 393-400 (2013).
- 536 10 Ashkenazi, A., Fairbrother, W. J., Leverson, J. D. & Souers, A. J. From basic apoptosis
537 discoveries to advanced selective BCL-2 family inhibitors. *Nature reviews drug discovery* **16**,
538 273-284 (2017).
- 539 11 Shin, W.-H., Christoffer, C. W. & Kihara, D. In silico structure-based approaches to discover
540 protein-protein interaction-targeting drugs. *Methods* **131**, 22-32 (2017).
- 541 12 Shin, W.-H., Kumazawa, K., Imai, K., Hirokawa, T. & Kihara, D. Current challenges and
542 opportunities in designing protein-protein interaction targeted drugs. *Advances and*
543 *Applications in Bioinformatics and Chemistry*, 11-25 (2020).

- 544 13 Anderson, A. C. The process of structure-based drug design. *Chemistry & biology* **10**, 787-797
545 (2003).
- 546 14 Wang, X., Song, K., Li, L. & Chen, L. Structure-based drug design strategies and challenges.
547 *Current Topics in Medicinal Chemistry* **18**, 998-1006 (2018).
- 548 15 Batool, M., Ahmad, B. & Choi, S. A structure-based drug discovery paradigm. *International*
549 *journal of molecular sciences* **20**, 2783 (2019).
- 550 16 Danel, T., Łęski, J., Podlewska, S. & Podolak, I. T. Docking-based generative approaches in the
551 search for new drug candidates. *Drug Discovery Today*, 103439 (2022).
- 552 17 Isert, C., Atz, K. & Schneider, G. Structure-based drug design with geometric deep learning.
553 *Current Opinion in Structural Biology* **79**, 102548 (2023).
- 554 18 Rakers, C., Bermudez, M., Keller, B. G., Mortier, J. & Wolber, G. Computational close up on
555 protein-protein interactions: how to unravel the invisible using molecular dynamics simulations?
556 *Wiley Interdisciplinary Reviews: Computational Molecular Science* **5**, 345-359 (2015).
- 557 19 Ni, D., Lu, S. & Zhang, J. Emerging roles of allosteric modulators in the regulation of protein-
558 protein interactions (PPIs): A new paradigm for PPI drug discovery. *Medicinal research reviews*
559 **39**, 2314-2342 (2019).
- 560 20 Janin, J. & Chotia, C. The structure of protein-protein recognition sites. *The Journal of*
561 *biological chemistry (Print)* **265**, 16027-16030 (1990).
- 562 21 Smith, M. C. & Gestwicki, J. E. Features of protein-protein interactions that translate into potent
563 inhibitors: topology, surface area and affinity. *Expert reviews in molecular medicine* **14**, e16
564 (2012).
- 565 22 Wang, Z.-Z., Shi, X.-X., Huang, G.-Y., Hao, G.-F. & Yang, G.-F. Fragment-based drug
566 discovery supports drugging ‘undruggable’ protein-protein interactions. *Trends in Biochemical*
567 *Sciences* (2023).
- 568 23 Mignani, S. *et al.* Present drug-likeness filters in medicinal chemistry during the hit and lead
569 optimization process: how far can they be simplified? *Drug discovery today* **23**, 605-615 (2018).
- 570 24 Lipinski, C., Lombardo, F., Dominy, B. & Feeney, P. In vitro models for selection of
571 development candidates experimental and computational approaches to estimate solubility and
572 permeability in drug discovery and development settings. *Adv Drug Deliv Rev* **23**, 3-25 (1997).
- 573 25 Lipinski, C. A. Lead-and drug-like compounds: the rule-of-five revolution. *Drug discovery*
574 *today: Technologies* **1**, 337-341 (2004).
- 575 26 Morelli, X., Bourgeas, R. & Roche, P. Chemical and structural lessons from recent successes in
576 protein-protein interaction inhibition (2P2I). *Current opinion in chemical biology* **15**, 475-481
577 (2011).
- 578 27 Bickerton, G. R., Paolini, G. V., Besnard, J., Muresan, S. & Hopkins, A. L. Quantifying the
579 chemical beauty of drugs. *Nature chemistry* **4**, 90-98 (2012).
- 580 28 Kosugi, T. & Ohue, M. in *2021 IEEE Conference on Computational Intelligence in*
581 *Bioinformatics and Computational Biology (CIBCB)*. 1-8.
- 582 29 Kosugi, T. & Ohue, M. Quantitative Estimate Index for Early-Stage Screening of Compounds
583 Targeting Protein-Protein Interactions. *International Journal of Molecular Sciences* **22**, 10925
584 (2021).
- 585 30 Wang, J., Mao, J., Wang, M., Le, X. & Wang, Y. Explore drug-like space with deep generative
586 models. *Methods* (2023).
- 587 31 Qiu, Y. *et al.* Computational methods-guided design of modulators targeting protein-protein

- 588 interactions (PPIs). *European Journal of Medicinal Chemistry* **207**, 112764 (2020).
- 589 32 Stokel-Walker, C. & Van Noorden, R. What ChatGPT and generative AI mean for science.
590 *Nature* **614**, 214-216 (2023).
- 591 33 Urbina, F., Lentzos, F., Invernizzi, C. & Ekins, S. Dual use of artificial-intelligence-powered
592 drug discovery. *Nature Machine Intelligence* **4**, 189-191 (2022).
- 593 34 Bilodeau, C., Jin, W., Jaakkola, T., Barzilay, R. & Jensen, K. F. Generative models for molecular
594 discovery: Recent advances and challenges. *Wiley Interdisciplinary Reviews: Computational
595 Molecular Science* **12**, e1608 (2022).
- 596 35 Cheng, Y., Gong, Y., Liu, Y., Song, B. & Zou, Q. Molecular design in drug discovery: a
597 comprehensive review of deep generative models. *Briefings in bioinformatics* **22**, bbab344
598 (2021).
- 599 36 Tong, X. *et al.* Generative models for De Novo drug design. *Journal of Medicinal Chemistry*
600 **64**, 14011-14027 (2021).
- 601 37 Wang, M. *et al.* Deep learning approaches for de novo drug design: An overview. *Current
602 opinion in structural biology* **72**, 135-144 (2022).
- 603 38 Meyers, J., Fabian, B. & Brown, N. De novo molecular design and generative models. *Drug
604 Discovery Today* (2021).
- 605 39 Thomas, M., Bender, A. & de Graaf, C. Integrating structure-based approaches in generative
606 molecular design. *Current Opinion in Structural Biology* **79**, 102559 (2023).
- 607 40 Zeng, X. *et al.* Deep generative molecular design reshapes drug discovery. *Cell Reports
608 Medicine*, 100794 (2022).
- 609 41 Martinelli, D. Generative machine learning for de novo drug discovery: A systematic review.
610 *Computers in Biology and Medicine*, 105403 (2022).
- 611 42 Özçelik, R., van Tilborg, D., Jiménez-Luna, J. & Grisoni, F. Structure-based drug discovery
612 with deep learning. *arXiv preprint arXiv:2212.13295* (2022).
- 613 43 Ma, B. *et al.* Structure-based de novo molecular generator combined with artificial intelligence
614 and docking simulations. *Journal of Chemical Information and Modeling* **61**, 3304-3313 (2021).
- 615 44 Luo, S., Guan, J., Ma, J. & Peng, J. A 3D generative model for structure-based drug design.
616 *Advances in Neural Information Processing Systems* **34**, 6229-6239 (2021).
- 617 45 Drotár, P., Jamasb, A. R., Day, B., Cangea, C. & Liò, P. Structure-aware generation of drug-like
618 molecules. *arXiv preprint arXiv:2111.04107* (2021).
- 619 46 Li, Y., Pei, J. & Lai, L. Structure-based de novo drug design using 3D deep generative models.
620 *Chemical science* **12**, 13664-13675 (2021).
- 621 47 Long, S., Zhou, Y., Dai, X. & Zhou, H. Zero-Shot 3D Drug Design by Sketching and Generating.
622 Peng, X. *et al.* in *International Conference on Machine Learning*. 17644-17655 (PMLR).
- 623 49 Wang, M. *et al.* Relation: A deep generative model for structure-based de novo drug design.
624 *Journal of Medicinal Chemistry* **65**, 9478-9492 (2022).
- 625 50 Chan, L., Kumar, R., Verdonk, M. & Poelking, C. A multilevel generative framework with
626 hierarchical self-contrasting for bias control and transparency in structure-based ligand design.
627 *Nature Machine Intelligence*, 1-13 (2022).
- 628 51 Zhang, O. *et al.* ResGen is a pocket-aware 3D molecular generation model based on parallel
629 multiscale modelling. *Nature Machine Intelligence*, 1-11 (2023).
- 630 52 Neugebauer, A., Hartmann, R. W. & Klein, C. D. Prediction of protein– protein interaction
631 inhibitors by chemoinformatics and machine learning methods. *Journal of medicinal chemistry*

- 632 **50**, 4665-4668 (2007).
- 633 53 Gupta, P. & Mohanty, D. SMMPPPI: a machine learning-based approach for prediction of
634 modulators of protein–protein interactions and its application for identification of novel
635 inhibitors for RBD: hACE2 interactions in SARS-CoV-2. *Briefings in Bioinformatics* **22**,
636 bbab111 (2021).
- 637 54 Díaz-Eufracio, B. I. & Medina-Franco, J. L. Machine Learning Models to Predict Protein–
638 Protein Interaction Inhibitors. *Molecules* **27**, 7986 (2022).
- 639 55 Reker, D., Schneider, P. & Schneider, G. Multi-objective active machine learning rapidly
640 improves structure–activity models and reveals new protein–protein interaction inhibitors.
641 *Chemical science* **7**, 3919-3927 (2016).
- 642 56 Mallet, V. *et al.* InDeep: 3D fully convolutional neural networks to assist in silico drug design
643 on protein–protein interactions. *Bioinformatics* **38**, 1261-1268 (2022).
- 644 57 Wang, J. *et al.* De novo molecular design with deep molecular generative models for PPI
645 inhibitors. *Briefings in Bioinformatics* **23**, doi:10.1093/bib/bbac285 (2022).
- 646 58 Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M. & Monfardini, G. The graph neural
647 network model. *IEEE transactions on neural networks* **20**, 61-80 (2008).
- 648 59 Velickovic, P. *et al.* Graph attention networks. *stat* **1050**, 10.48550 (2017).
- 649 60 Wang, X., Flannery, S. T. & Kihara, D. Protein docking model evaluation by graph neural
650 networks. *Frontiers in Molecular Biosciences* **8**, 647915 (2021).
- 651 61 Mirza, M. & Osindero, S. Conditional generative adversarial nets. *arXiv preprint*
652 *arXiv:1411.1784* (2014).
- 653 62 Wang, Y., Wu, S., Duan, Y. & Huang, Y. ResAtom system: protein and ligand affinity prediction
654 model based on deep learning. *arXiv preprint arXiv:2105.05125* (2021).
- 655 63 Zapata, P. A. M. *et al.* Cell morphology-guided de novo hit design by conditioning GANs on
656 phenotypic image features. *Digital Discovery* (2023).
- 657 64 Schmidhuber, J. & Hochreiter, S. Long short-term memory. *Neural Comput* **9**, 1735-1780
658 (1997).
- 659 65 Wei, W., Cherukupalli, S., Jing, L., Liu, X. & Zhan, P. Fsp3: A new parameter for drug-likeness.
660 *Drug Discovery Today* **25**, 1839-1845 (2020).
- 661 66 Polykovskiy, D. *et al.* Molecular sets (MOSES): a benchmarking platform for molecular
662 generation models. *Frontiers in pharmacology* **11**, 1931 (2020).
- 663 67 Mamoshina, P., Vieira, A., Putin, E. & Zhavoronkov, A. Applications of deep learning in
664 biomedicine. *Molecular pharmaceuticals* **13**, 1445-1454 (2016).
- 665 68 Segler, M. H., Kogej, T., Tyrchan, C. & Waller, M. P. Generating focused molecule libraries for
666 drug discovery with recurrent neural networks. *ACS central science* **4**, 120-131 (2018).
- 667 69 Kingma, D. P. & Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*
668 (2013).
- 669 70 Prykhodko, O. *et al.* A de novo molecular generation method using latent vector based
670 generative adversarial network. *Journal of Cheminformatics* **11**, 1-13 (2019).
- 671 71 Guimaraes, G. L., Sanchez-Lengeling, B., Outeiral, C., Farias, P. L. C. & Aspuru-Guzik, A.
672 Objective-reinforced generative adversarial networks (organ) for sequence generation models.
673 *arXiv preprint arXiv:1705.10843* (2017).
- 674 72 Durant, J. L., Leland, B. A., Henry, D. R. & Nourse, J. G. Reoptimization of MDL keys for use
675 in drug discovery. *Journal of chemical information and computer sciences* **42**, 1273-1280

- 676 (2002).
- 677 73 Van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning*
678 *research* **9** (2008).
- 679 74 Hinton, G. E. & Roweis, S. Stochastic neighbor embedding. *Advances in neural information*
680 *processing systems* **15** (2002).
- 681 75 Meyers, J., Carter, M., Mok, N. Y. & Brown, N. On the origins of three-dimensionality in drug-
682 like molecules. *Future medicinal chemistry* **8**, 1753-1767 (2016).
- 683 76 Probst, D. & Reymond, J.-L. Visualization of very large high-dimensional data sets as minimum
684 spanning trees. *Journal of Cheminformatics* **12**, 1-13 (2020).
- 685 77 Altae-Tran, H., Ramsundar, B., Pappu, A. S. & Pande, V. Low data drug discovery with one-
686 shot learning. *ACS central science* **3**, 283-293 (2017).
- 687 78 Moret, M., Friedrich, L., Grisoni, F., Merk, D. & Schneider, G. Generative molecular design in
688 low data regimes. *Nature Machine Intelligence* **2**, 171-180 (2020).
- 689 79 Wang, J., Zheng, S., Chen, J. & Yang, Y. Meta learning for low-resource molecular optimization.
690 *Journal of Chemical Information and Modeling* **61**, 1627-1636 (2021).
- 691 80 Chen, X. *et al.* DCZ3112, a novel Hsp90 inhibitor, exerts potent antitumor activity against
692 HER2-positive breast cancer through disruption of Hsp90-Cdc37 interaction. *Cancer Letters*
693 **434**, 70-80 (2018).
- 694 81 Allen, W. J. *et al.* DOCK 6: Impact of new features and current docking performance. *Journal*
695 *of computational chemistry* **36**, 1132-1156 (2015).
- 696 82 Singh, N., Chaput, L. & Villoutreix, B. O. Fast rescoring protocols to improve the performance
697 of structure-based virtual screening performed on protein-protein interfaces. *Journal of*
698 *chemical information and modeling* **60**, 3910-3934 (2020).
- 699 83 Gaulton, A. *et al.* ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic acids*
700 *research* **40**, D1100-D1107 (2012).
- 701 84 Torng, W. & Altman, R. B. Graph convolutional neural networks for predicting drug-target
702 interactions. *Journal of chemical information and modeling* **59**, 4131-4149 (2019).
- 703 85 Lim, J. *et al.* Predicting drug-target interaction using a novel graph neural network with 3D
704 structure-embedded graph representation. *Journal of chemical information and modeling* **59**,
705 3981-3988 (2019).
- 706 86 Landrum, G. RDKit: Open-source cheminformatics. <https://www.rdkit.org/>,
707 doi:<https://www.rdkit.org/> (2006).
- 708 87 Doerr, S., Harvey, M., Noé, F. & De Fabritiis, G. HTMD: high-throughput molecular dynamics
709 for molecular discovery. *Journal of chemical theory and computation* **12**, 1845-1852 (2016).
- 710 88 Lu, T. & Chen, F. Multiwfn: A multifunctional wavefunction analyzer. *Journal of computational*
711 *chemistry* **33**, 580-592 (2012).
- 712 89 Skalic, M., Jiménez, J., Sabbadin, D. & De Fabritiis, G. Shape-based generative modeling for
713 de novo drug design. *Journal of chemical information and modeling* **59**, 1205-1214 (2019).
- 714 90 Gaulton, A. *et al.* A large-scale crop protection bioassay data set. *Scientific Data* **2**, 150032,
715 doi:10.1038/sdata.2015.32 (2015).
- 716 91 Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library.
717 *Advances in neural information processing systems* **32** (2019).
- 718 92 Abadi, M. *et al.* TensorFlow: Large-scale machine learning on heterogeneous systems, software
719 available from tensorflow.org (2015). URL <https://www.tensorflow.org> (2015).

- 720 93 Preuer, K., Renz, P., Untertiner, T., Hochreiter, S. & Klambauer, G. n. Fréchet ChemNet
721 distance: a metric for generative models for molecules in drug discovery. *Journal of chemical*
722 *information and modeling* **58**, 1736-1741 (2018).
- 723 94 Sauer, W. H. & Schwarz, M. K. Molecular shape diversity of combinatorial libraries: a
724 prerequisite for broad bioactivity. *Journal of chemical information and computer sciences* **43**,
725 987-1003 (2003).
- 726 95 Firth, N. C., Brown, N. & Blagg, J. Plane of best fit: a novel method to characterize the three-
727 dimensionality of molecules. *Journal of chemical information and modeling* **52**, 2516-2525
728 (2012).
- 729 96 Probst, D. & Reymond, J.-L. A probabilistic molecular fingerprint for big data settings. *Journal*
730 *of cheminformatics* **10**, 1-12 (2018).
- 731 97 Probst, D. & Reymond, J.-L. FUN: a framework for interactive visualizations of large, high-
732 dimensional datasets on the web. *Bioinformatics* **34**, 1433-1435 (2018).
- 733 98 Roe, S. M. *et al.* The mechanism of Hsp90 regulation by the protein kinase-specific cochaperone
734 p50cdc37. *Cell* **116**, 87-98 (2004).
- 735 99 Dike, P. P. *et al.* In silico identification of small molecule modulators for disruption of Hsp90–
736 Cdc37 protein–protein interaction interface for cancer therapeutic application. *Journal of*
737 *Biomolecular Structure and Dynamics* **40**, 2082-2098 (2022).
- 738 100 Pettersen, E. F. *et al.* UCSF Chimera—a visualization system for exploratory research and
739 analysis. *Journal of computational chemistry* **25**, 1605-1612 (2004).
- 740 101 Schrödinger, L. (November, 2015).

741

742 **Figure legends**

743 **Fig. 1. Generation of molecules targeting PPI.** 3D structural information of the
744 protein-protein complex interface is represented as a graph. Feature representation of
745 the interface region is captured by using a graph attention neural networks. The
746 representation of the voxel and electron density of the compound is encoded by a 3D
747 convolutional neural networks (CNNs). A conditional Wasserstein generative
748 adversarial networks is trained to generate molecular embeddings with interface
749 features as conditions. Generator: takes interface features and random noise vectors to
750 generate molecular embeddings for the input features. Discriminator: calculates the
751 probability that a molecule is from a real or a fake molecule. Condition: controls or
752 regulates the generation of molecules constrained by a specific protein-protein interface.
753 Finally, a long short-term memory (LSTM) networks parse SMILES strings from
754 molecular representation.

755

756 **Fig. 2:** Results of conditional evaluation. (a) The QED, QEPPI and Fsp3 distribution of
757 active compounds and compounds generated by the GENiPPI framework for
758 MDM2/p53; (b) The QED, QEPPI and Fsp3 distribution of active compounds and
759 compounds generated by the GENiPPI framework for Bcl-2/Bax; (c) The QED, QEPPI
760 and Fsp3 distribution of active compounds and compounds generated by the GENiPPI
761 framework for BAZ2B/H4; (d) The QED distribution of generated compounds for
762 MDM2/p53, Bcl-2/Bax and BAZ2B/H4; (e) The QEPPI distribution of generated
763 compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4; (f) The Fsp3 distribution of
764 generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4;

765

766 **Fig. 3:** Chemical space exploration. (a) The t-SNE visualization of active compounds
767 and generated compounds for MDM2/p53; (b) The t-SNE visualization of active
768 compounds and generated compounds for Bcl-2/Bax; (c) The t-SNE visualization of
769 active compounds and generated compounds for BAZ2B/H4; (d) The PMI ternary
770 density plots of generated compounds, small molecule drugs of DrugBank, and iPPI-
771 DB inhibitors. Top left: propyne, bottom: benzene, and the top right: adamantane; (e)
772 The molecular three-dimensionality distribution of the generated molecules was
773 visualized with NPR descriptors and PBF descriptors. (f) TMAP visualization of active
774 compounds and generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4.

775

776 **Fig. 4:** Few shot molecular generation analysis. (a) The t-SNE visualization of the
777 distribution of active compounds and generated compounds for Hsp90/Cdc37; (b)
778 Comparison of the pharmacophore of the generated molecules with the reference
779 molecule(DCZ3112); (c) PPI interface region(in green) of the Hsp90(in
780 palecyan)/Cdc37(in lightpink) complex; (d) The complex structure of DCZ3112(in
781 green) and Hsp90(in palecyan)-CDC37(in lightpink) modeled by molecular docking
782 (PDB ID: 1US7); (e) The binding poses of generated compounds(in green) and
783 Hsp90(in palecyan)-CDC37(in lightpink) modeled by molecular docking (PDB ID:
784 1US7). Hydrogen bonds are displayed as blue dotted lines. π -cation Interactions are
785 displayed as orange dotted lines.

Figure 1

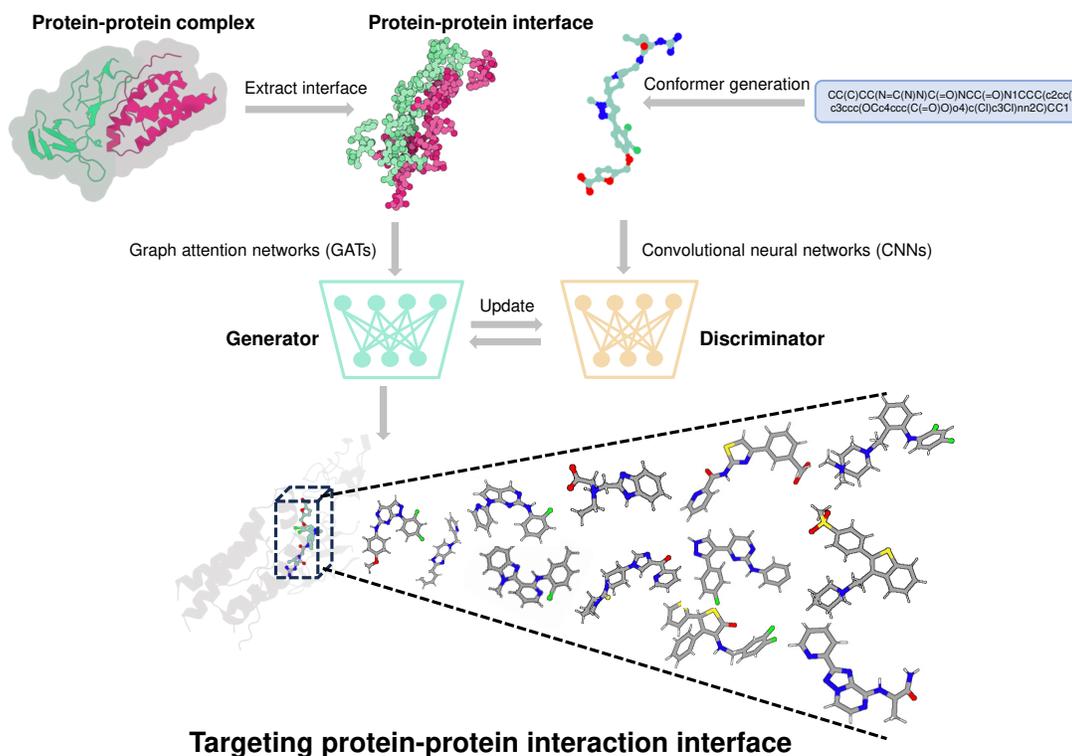


Fig. 1. Generation of molecules targeting the PPI interface. 3D structural information of the protein-protein complex interface is represented as a graph. Feature representation of the interface region is captured by using a graph attention neural networks. The representation of the voxel and electron density of the compound is encoded by a 3D convolutional neural networks (CNNs). A conditional Wasserstein generative adversarial networks is trained to generate molecular embeddings with interface features as conditions. Generator: takes interface features and random noise vectors to generate molecular embeddings for the input features. Discriminator: calculates the probability that a molecule is from a real or a fake molecule. Condition: controls or regulates the generation of molecules constrained by a specific protein-protein interface.

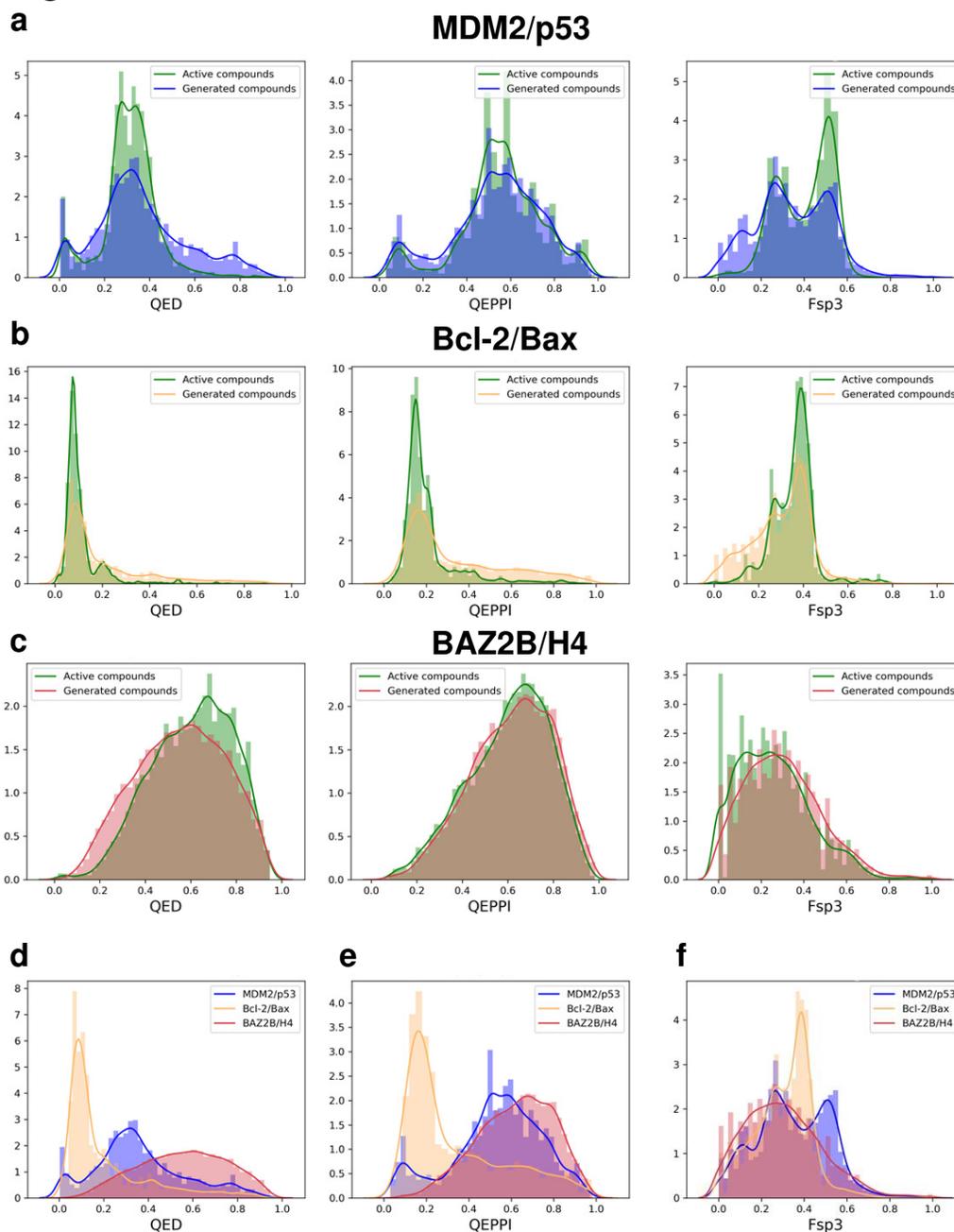
Figure 2

Fig. 2. Results of conditional evaluation. (a) The QED, QEPI and Fsp3 distribution of active compounds and compounds generated by the GENiPPI framework for MDM2/p53; **(b)** The QED, QEPI and Fsp3 distribution of active compounds and compounds generated by the GENiPPI framework for Bcl-2/Bax; **(c)** The QED, QEPI and Fsp3 distribution of active compounds and compounds generated by the GENiPPI framework for BAZ2B/H4; **(d)** The QED distribution of generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4; **(e)** The QEPI distribution of generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4; **(f)** The Fsp3 distribution of generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4.

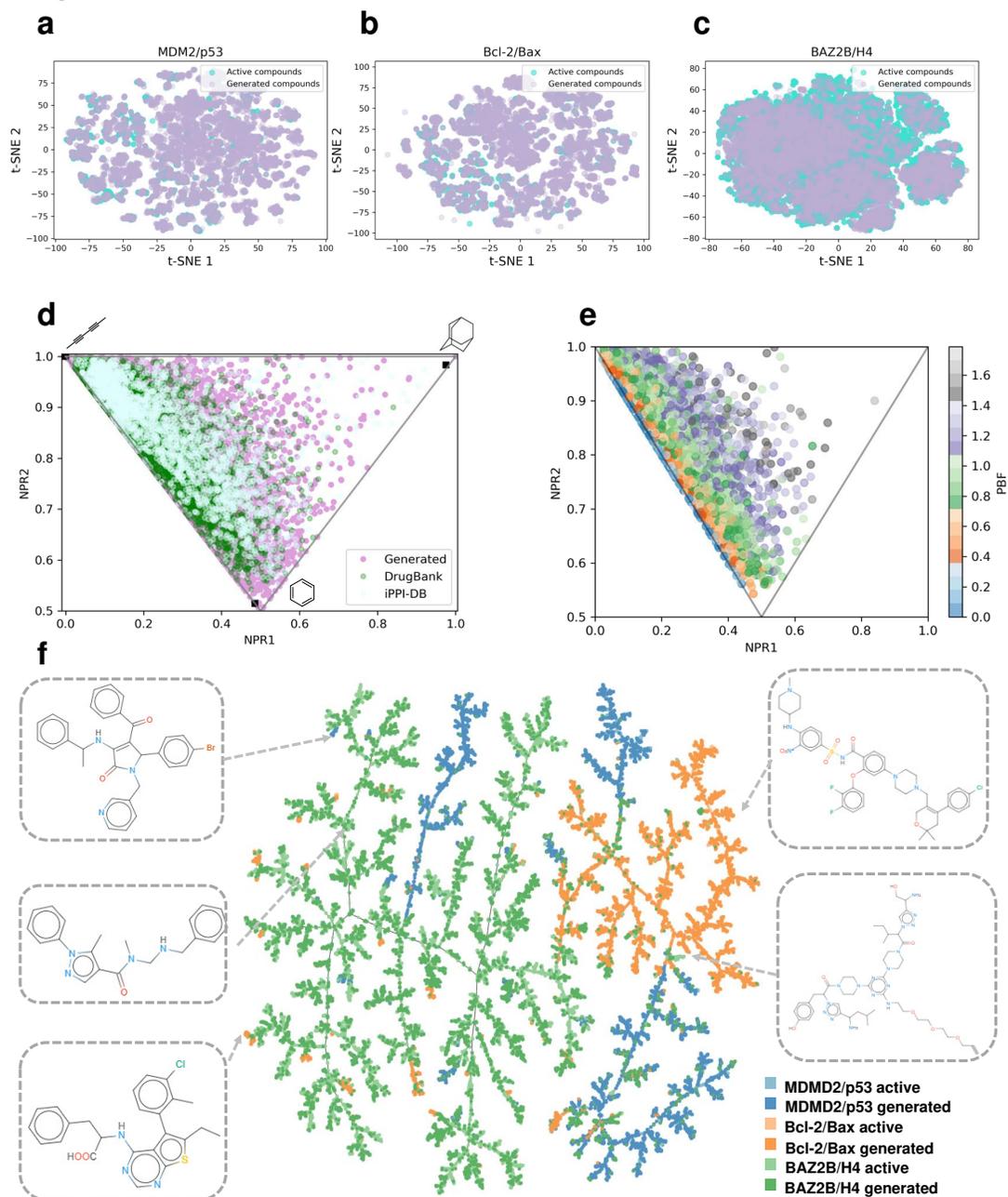
Figure 3

Fig. 3. Chemical space exploration. (a) The t-SNE visualization of active compounds and generated compounds for MDM2/p53; (b) The t-SNE visualization of active compounds and generated compounds for Bcl-2/Bax; (c) The t-SNE visualization of active compounds and generated compounds for BAZ2B/H4; (d) The PMI ternary density plots of generated compounds, small molecule drugs of DrugBank, and iPPI-DB inhibitors. Top left: propyne, bottom: benzene, and the top right: adamantane; (e) The molecular three-dimensionality distribution of the generated molecules was visualized with NPR descriptors and PBF descriptors. (f) TMAP visualization of active compounds and generated compounds for MDM2/p53, Bcl-2/Bax and BAZ2B/H4.

Figure 4

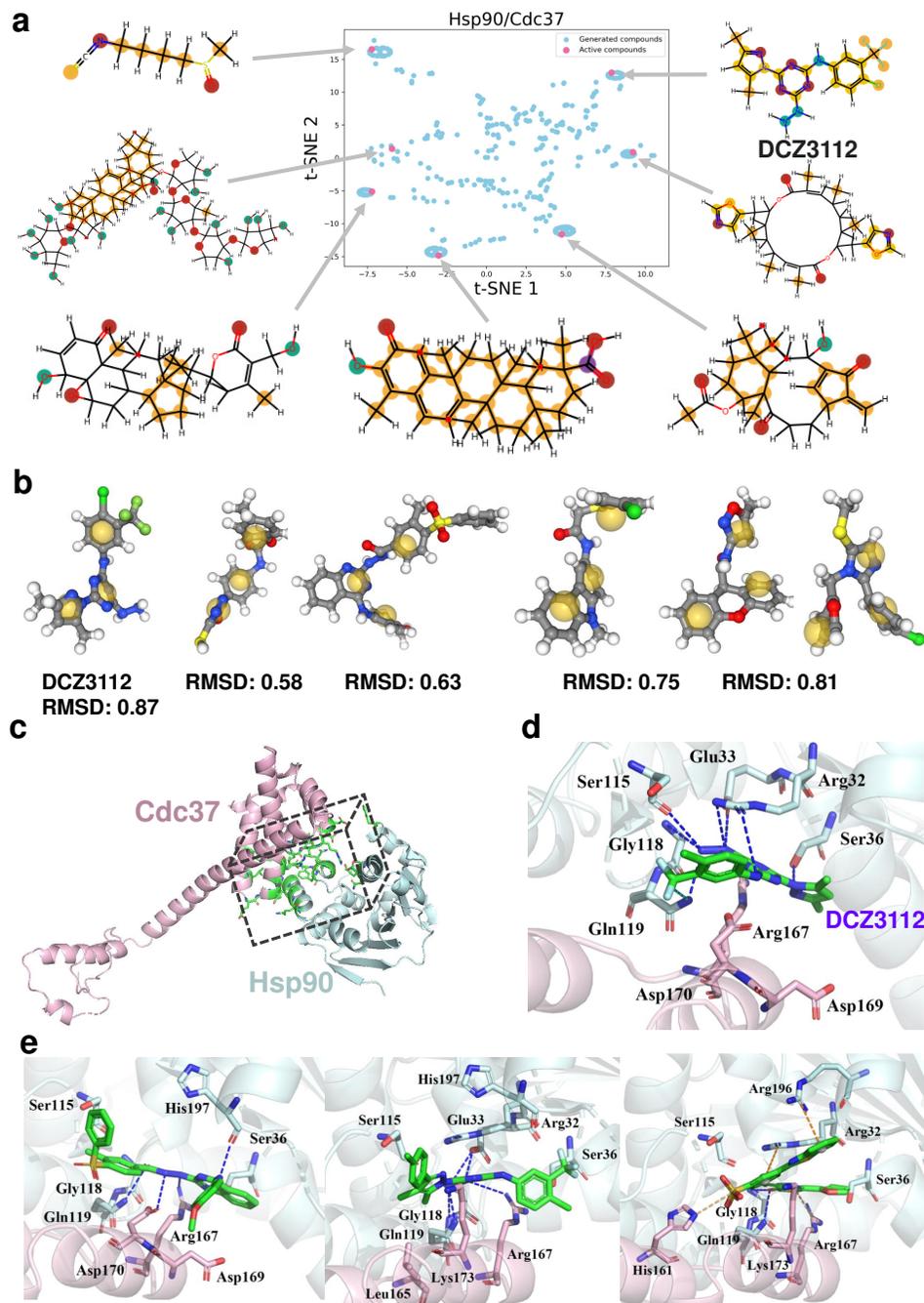


Fig. 4. Few shot molecular generation analysis. (a) The t-SNE visualization of the distribution of active compounds and generated compounds for Hsp90/Cdc37; (b) Comparison of the pharmacophore of the generated molecules with the reference molecule (DCZ3112); (c) PPI interface region (in green) of the Hsp90 (in palecyan) and Cdc37 (in lightpink) complex; (d) The complex structure of DCZ3112 (in green) and Hsp90 (in palecyan) and CDC37 (in lightpink) modeled by molecular docking (PDB ID: 1US7); (e) The binding poses of generated compounds (in green) and Hsp90 (in palecyan) and CDC37 (in lightpink) modeled by molecular docking (PDB ID: 1US7). Hydrogen bonds are displayed as blue dotted lines. π -cation interactions are displayed as orange dotted lines.

Table 1.**Table 1.** Comparisons between PPI interfaces and binding sites

PPI interfaces	Binding sites
Target properties	
Large surface area (1000-6000 Å ²)	Small surface (300-1000 Å ²) Hydrophobic
Preference for Trp (W), Tyr (Y), and Arg (R) as PPI hotspot residues; subpockets	Large volume (~260 Å ³)
Shallow, flat, flexible	Pocket, cliff
Hydrophobic, featureless, undruggability	Diverse properties
Chemical space	
MW ≥ 400	MW ≤ 500
LogP ≥ 4	LogP ≤ 500
HBA ≥ 4	HBA ≤ 10
number of rings: ≥ 4	HBD ≤ 5
Ro4 Morelli's rules	Lipinski's Rule of 5 (Ro5)
Quantitative estimate of drug-likeness scores	
QEPPi	QED

Table 2.**Table 2.** Valid, unique, novelty and FCD of sampling SMILES after training. We sampled 30,000 SMILES each time.

Model	valid	Unique@1k	Unique@10k	novelty	FCD	
					Test	TestSF
AAE	0.881	1.000	0.995	0.995	8.573	9.117
CharRNN	0.985	0.999	0.988	0.994	8.7564	8.952
VAE	0.834	1.000	0.996	0.994	7.703	8.141
LatentGAN	0.724	1.000	0.999	0.998	7.595	8.160
ORGAN	0.609	0.996	0.994	0.999	39.800	41.158
GENiPPI(noninterface)	0.999	0.997	0.975	0.997	7.653	8.132
GENiPPI	0.999	0.998	0.977	0.998	7.450	7.884