

1 **Chemical Proteomics-based Target Prioritization through a Residue Agnostic**  
2 **Ligandability Assessment Platform**

3 **Fettah Erdogan<sup>\*a</sup>, Raiyan Chowdhury<sup>c</sup>, Serap Beldar<sup>a</sup>, Tom Bobby Chandy<sup>d</sup>, Rebecca G.**  
4 **Allan<sup>a</sup>, Elvin D. de Araujo<sup>a</sup>, Patrick T. Gunning<sup>\*a,b</sup>**

5 <sup>a</sup> Department of Chemical and Physical Sciences, University of Toronto Mississauga, Ontario L5L  
6 1C6, Canada

7 <sup>b</sup> Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada

8 <sup>c</sup> Cheriton School of Computer Science, University of Waterloo, Ontario N2L 3G1, Canada

9 <sup>d</sup> Faculty of Health Sciences, Western University, London, ON N6A 3K7, Canada

10

11 **Corresponding Authors**

12 \*Patrick T. Gunning, Department of Chemical and Physical Sciences, University of Toronto  
13 Mississauga, 3359 Mississauga Rd N., Mississauga, Ontario L5L 1C6

14 Email: [patrick.gunning@utoronto.ca](mailto:patrick.gunning@utoronto.ca)

15 \*Fettah Erdogan, Department of Chemical and Physical Sciences, University of Toronto  
16 Mississauga, 3359 Mississauga Rd N., Mississauga, Ontario L5L 1C6

17 Email: [fettah.erdogan@alum.utoronto.ca](mailto:fettah.erdogan@alum.utoronto.ca)

18 Subject terms: covalent binding site, computational ligandability assessment, chemoproteomics,  
19 chemical proteomics, computational target prioritization, computer-aided drug discovery, machine  
20 learning

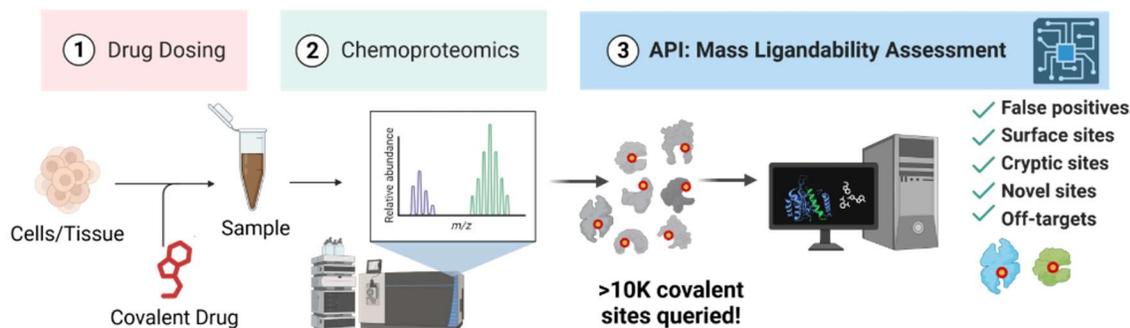
21

22

23

## 24 Graphical Abstract

### Accelerating Bottom-Up Covalent Drug Discovery



25

26

## 27 Abstract

28 The landscape of drug discovery is undergoing a transformative phase with the influx of structural  
29 biology and omics data. Identifying optimal drug targets amid this data surge presents a  
30 multifaceted challenge. Covalent inhibitors, once undervalued, now hold substantial promise,  
31 especially targeted covalent inhibitors (TCIs), effectively engaging 'undruggable' proteins and  
32 overcoming resistance mechanisms. Existing ML software can proficiently model covalent ligands  
33 but lack comprehensive utility across large chemoproteomics sites. Challenges persist in  
34 predicting and assessing cryptic ligandable sites and sites beyond cysteine, demanding advanced  
35 computational tools. As cysteine-ligandable proteins represent only ~20% of the quantifiable  
36 proteome, there is a requirement for ligandability mapping of other nucleophilic amino acids. This  
37 study introduces a pioneering computational pipeline leveraging an AI-based ligandable predictor  
38 for meticulous evaluation of chemical proteomics-based reactive sites. The pipeline offers a  
39 scalable framework to assess covalent ligandability on a large scale, filter out improbable hits and  
40 systematically evaluate potential drug targets. Our work addresses covalent drug design  
41 challenges through a pipeline that fills crucial gaps in predicting cryptic ligandable and covalent  
42 sites in addition to cysteines to foster more efficient drug discovery methodologies.

43

## 44 Introduction

45 The magnitude of structural biology and omics information arriving at the service of drug discovery  
46 and medicinal chemistry research is expediting the rate of hit-to-lead development<sup>1</sup>. However,  
47 harnessing this deluge of information to identify optimal drug targets poses a multifaceted  
48 challenge. The efficacy and success of drug development largely hinges upon precise target  
49 identification, necessitating innovative technologies to address these challenges. Covalent  
50 inhibitors, long overlooked due to concerns about reactivity and off-target effects<sup>2</sup>, are now gaining  
51 substantial attention and reverence in both academic and pharmaceutical drug design programs.  
52 Improved understanding of the factors influencing reactivity and emergence of new types of  
53 warheads in the form of targeted covalent inhibitors (TCIs), represent a promising avenue within  
54 drug discovery, particularly in targeting proteins previously deemed 'undruggable'.<sup>3-8</sup> Over 40  
55 covalent drugs are currently under clinical development<sup>2</sup> and many previously so called  
56 "undruggable" targets and mechanisms of resistance have now been effectively tackled by  
57 covalent compounds.<sup>3,9,10</sup> Recently developed and FDA approved TCIs, AMG-510 (sotorasib) and  
58 MRTX-849 (adagrasib), effectively engage Cys12 of KRas<sup>G12C</sup>, a famously difficult target  
59 implicated in 40% of lung cancers.<sup>6,7,11</sup> Additionally, ibrutinib, a first-in-class inhibitor known to  
60 successfully bind to Cyst of Bruton's tyrosine kinase (BTK), which is linked to overexpression of  
61 B cells in B cell cancers, has FDA approval for lymphoma.<sup>12-14</sup> Other FDA approved drugs such  
62 as Aspirin, Penicillin G, and Fosfomycin are all reported to engaged targets via covalent  
63 interaction<sup>15</sup>.

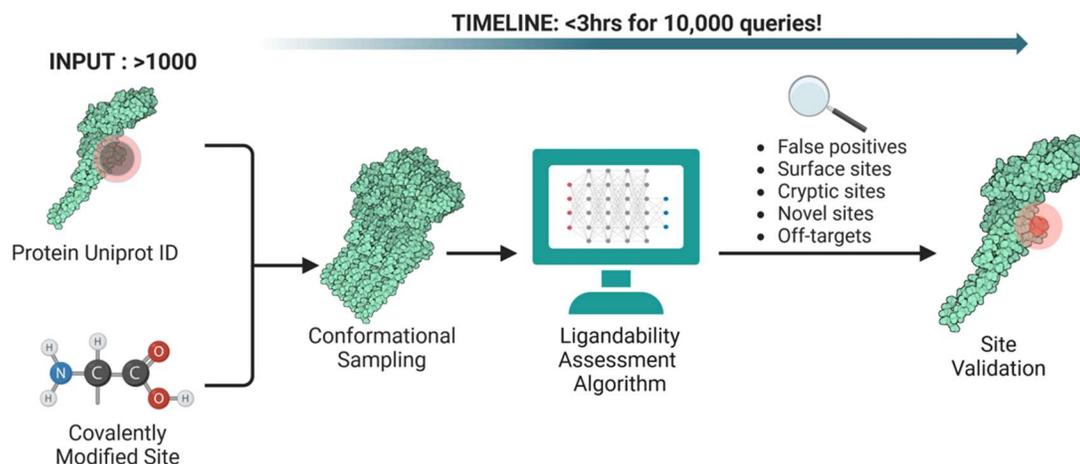
64 TCIs enact specificity and irreversible binding through a combination of covalent and non-  
65 covalent interactions at their protein's target site. Although non-covalent interactions make up  
66 the majority of contacts between a TCI and the binding pocket residues of the protein, target  
67 engagement via covalent bond formation (kinact) increases drugging efficiency (KI) by  
68 prolonging duration of action and increasing the degree of protein-drug occupancy<sup>16</sup> in some  
69 cases de-coupling PK-PD<sup>17</sup>, creating an opportunity to target shallow binding sites of  
70 challenging targets<sup>18</sup> and improving selectivity<sup>19</sup> towards target, target isoforms and disease-  
71 linked mutants. Given these appealing features, development of new tools to support the  
72 rational design of TCIs has become an indispensable task.

73 Moreover, drug discovery remains a prolonged and expensive process characterized by a  
74 notably low (13.8%) success rate<sup>20</sup>. The integration of machine learning (ML) technologies into  
75 drug discovery research emerges as a crucial avenue to address these challenges. Numerous

76 ML models focused on small-molecule design and quantitative structure-activity relationship  
77 (QSAR) have been specifically designed for virtual screening (VS) of drugs against targets.  
78 These innovations represent pioneering steps toward reducing overall costs and timelines in  
79 drug discovery research.<sup>21,22</sup> Various available computer-aided drug design (CADD) and  
80 screening programs such as Schrödinger CovDock<sup>23</sup>, DUckCov<sup>24</sup>, Cov\_DOX<sup>25</sup>, WIDOCK<sup>26</sup>,  
81 Reactive Docking<sup>27</sup>, and Bireactive<sup>28</sup>, specialize in modeling close contacts and predicting the  
82 binding mode of covalent ligands or the reactivity of warheads against a model sidechain.  
83 However, these programs efficiently serve their purpose only with prior knowledge of the binding  
84 site. As a result, they often require independent and laborious use in combination with non-  
85 /covalent binding site predictor programs such as P2Rank<sup>29</sup>, PocketFinder<sup>30</sup>, SiteMap<sup>10</sup>, and  
86 PocketMiner<sup>31</sup>, to initially identify non-covalent binding cavities, followed by calculations that  
87 suggest whether the binding pocket is amenable to covalent modification. Nevertheless, while  
88 some of these programs can be impractically coupled with nascent cysteine reactivity predictors  
89 such as DeepCoSi<sup>32</sup>, sbPCR<sup>33</sup>, and HyperCys<sup>34</sup>, many of these programs are intended for top-  
90 down proteomics discovery or have limitations. None of these platforms are designed with a  
91 utility of assessing ligandabilities across large numbers of chemoproteomics sites or are tailored  
92 for exploring sites other than cysteines. Recently, an MS-based quantitative proteolysis method  
93 (LiP-Quant) lightly integrating ML was used to prioritize true drug targets in chemoproteomics  
94 output dataset<sup>35</sup>. Although an expensive method anticipating establishment of acceptable  
95 detection rates on genuine ligandable sites and false-positives, LiP-Quant method can be  
96 promising for future chemical proteomics-based drug discovery research. In this work, we  
97 sought to develop an inexpensive method capable of analyzing existing chemoproteomics target  
98 sites in public or private repositories for proximal binding cavities and prioritize them using a  
99 ligandability score.

100 Covalent fragment-based drug discovery (FBDD) approaches using quantitative mass  
101 spectrometry (QMS) have conventionally been used to study inaccessible protein targets, reveal  
102 cryptic pockets and identify new potential targets in the proteome<sup>36-38</sup>. However, the success  
103 rate of chemical proteomics-based drug screening is comparable to crystallographic fragment-  
104 screening (5-10%).<sup>39</sup> Given the molecular dynamism of protein structures and the potential for  
105 occlusion of many transitory binding cavities, subpockets or PPI sites, an integrative *in-silico*  
106 pipeline that searches for potential druggable sites near nucleophilic residues specifically is a  
107 necessary part of the chemical proteomics-based covalent drug design pipeline.<sup>27,35,40,41</sup> Typical  
108 cell-based shotgun/bottom-up proteomics experiments can suggest >10,000 implicit protein

109 targets<sup>42</sup> however, visual inspection or *in-vitro* biophysical validation of each suggestive  
110 chemical proteomics-based hit is impractical and impossible. Given this challenge and in  
111 meeting the growing demand for TCI development<sup>5</sup>, we developed a computational pipeline that  
112 uses an AI-based ligandable predictor (DeepPocket<sup>43</sup>-developed) to scrupulously assess and  
113 mass-validate chemical proteomics-based reactive sites for apparent and cryptic ligandable  
114 cavities adjacent to each site. This comprehensive platform allows researchers to upload a  
115 proteomics-based covalent hits and perform a mass-scale ligandability evaluation of each  
116 suggestive target site for filtering out improbable hits. DeepPocket has emerged as a state-of-  
117 the-art model which builds on an established computational method, namely FPocket<sup>44</sup>, to map  
118 out the precise boundaries of ligandable sites and detect subcavities on a protein structure.  
119 Although the protein data bank (PDB) provides one of the most detailed descriptors on protein  
120 structures, it is imperative to emphasize that such structures are only simplified snap-shot  
121 models of the target macromolecules. Regeneration of these structures into conformer  
122 archetypes that contain slight perturbations in residue side chain coordinates (using tools such  
123 as CONCOORD<sup>45</sup>) can reveal transitory structural patterns for opportunistic pocket detection  
124 while preserving macromolecule's unique architectural style. CONCOORD produces protein  
125 conformers around an experimental structure using geometric restrictions. Studies of Molecular  
126 Dynamics (MD) simulations indicate that collective degrees of freedom are crucial to protein  
127 conformational changes, which are often vital to protein function. These internal constraints and  
128 configurational barriers can be used to simulate conformers without the need for more  
129 CPU intensive MD simulations<sup>45</sup>. Conformer generation using CONCOORD has been used in a  
130 similar fashion before<sup>46,47</sup> however we report its first use in covalent site ligandability  
131 assessment. Molecular structure descriptors such as PDB files are not to be thought of as a set  
132 of fixed coordinates but rather a framework for generating hypotheses based on molecular  
133 patterns to be explored<sup>48</sup>. The recent decade witnessed development of powerful generative AI  
134 models trained on omics data leading to algorithm that can recognize molecular features when  
135 faced with new data bearing similar characteristics. Our choice of DeepPocket as the pocket-  
136 predicting platform was inspired by a combination of its reliability on long-established methods,  
137 top-level performance and partitioned architecture which allowed us to further develop it towards  
138 detection of cryptic sites.



139

140 **Figure 1: Ligandability assessment pipeline architecture.** Thousands of chemoproteomics-based covalent site information is  
 141 simultaneously fed into the pipeline using the Uniprot ID and residue ID of the modified targets and their sites, respectively. Protein  
 142 structures are subject to random perturbations according to a predefined set of rules and constraints to produce conformers. Each  
 143 protein structure is passed through a geometry-based candidate pocket detection algorithm and re-ranked using an ML algorithm<sup>43,44</sup>.  
 144 Surface voxels of the pockets are then used to compute distances from the respective target site to the pocket surfaces and analyzed  
 145 by a ligandability assessment algorithm. The pipeline outputs ligandability score for each query originally input at an extraordinary  
 146 rate.

147 With diverse nucleophilic amino acids in ligandable proteins and the emerging need for  
 148 comprehensive mapping of the human proteome, covalent drug design approaches are poised  
 149 for significant advancements. The identification and targeting of cryptic ligandable sites within  
 150 proteins present immense potential for novel therapeutic interventions. However, the challenges  
 151 in predicting and assessing these sites, particularly beyond cysteine residues, underscore the  
 152 critical need for advanced computational tools. The developed pipeline (**Figure 1**), leveraging  
 153 AI-driven predictive models robustly assesses and validates chemical proteomics-based  
 154 reactive sites with exceptional speed. These advancements not only enable the detection of  
 155 cryptic pockets but also pave the way for a more efficient and systematic evaluation of potential  
 156 drug targets. By offering a scalable platform for ligandability assessment and rational hit  
 157 filtration, this work takes a significant step towards addressing the challenges in covalent drug  
 158 discovery, ultimately aiding for more efficient, precise, and successful drug development  
 159 strategies in the future.

160

161

162

163

## 164 **Methods**

### 165 **Data Curation and Preprocessing**

166 Unless otherwise mentioned, all data preparation and processing pipelines were built in-house  
167 using Python 3.10.0 and the Biopython<sup>49</sup> library. Since most PDB files contain multiple chains of  
168 proteins or macromolecules, each PDB file assess in this study was dissected for chain specific  
169 assessment, cleaned off from all ligands and/or double checked for errors with the associated  
170 covalent residue in the PDB. Any PDB chain alone with a covalent ligand sitting on a non-cavity  
171 site (absence of non-covalent interaction), missing the covalent residue in questions, or  
172 imperceptibly cryptic was omitted from analysis. Where applicable and unless otherwise  
173 mentioned, the positive covalent residue for each PDB was extracted and its respective  
174 negative was assigned by taking the most distant matching residue (matching false positive)  
175 within the chain.

176 **Holo Sites.** The initial list of 2,294 PDB coordinate files for high-resolution co-crystal structures  
177 of experimental covalent ligand-bound proteins (herein referred to as holo-protein) was scraped  
178 from the covPDB database<sup>50</sup>. A summary data file containing residue- and chain-specific  
179 annotation information was obtained directly from the covPDB website. Each annotated  
180 covalent site was subject to visual inspection (manual) for containment of the covalent ligand in  
181 a pocket-like cavity of the residue and saved as refined list (covPDBs, SupplementaryFile1).

182 **Apo sites.** Using the list of Uniprot accession codes referenced to the covPDB proteins, a list  
183 3,527 apo counterparts of the covPDBs (apo-covPDBs) was initially scraped from the RCSB  
184 database<sup>51</sup> while filtering out entities with non-polymer molecules (except waters, metal ions,  
185 and small-molecules typically used in buffers). The refined list of apo-CovPDBs used for  
186 analysis was arduously compiled by manual alignment to the covPDB counterparts and cross-  
187 checking for presence of all pocket-forming residues (apo-covPDBs, SupplementaryFile1).

188 **Cryptic sites.** The initial list of 124 apo-cryptic PDB files and their corresponding experimental  
189 holo-proteins was curated from Meller, et al<sup>31</sup> and Cimermanic, et al<sup>52</sup> and/or scraped from the  
190 RCSB repository<sup>53</sup> and refined further to a finalized list (cryptic-PDBs, SupplementaryFile1). By  
191 aligning the holo-protein counterparts of cryptic PDBs to apo-cryptic structure coordinates, a  
192 mock-covalent residue closest in distance to the ligand in the cryptic pocket of the apo-cryptic  
193 PDB was assigned to each cryptic PDB.

194

## 195 Pipeline Architecture and Data Processing

196 Except for Alphafold PDB coordinate files, each PDB chain primed for processing in this work is  
197 initially examined for potential conformational states that may reveal transitory or occluded  
198 binding pockets on protein structure. These conformational states alongside the experimental  
199 state and Alphafold predicted structures are scanned for identification of druggable binding  
200 sites. Methodologies describing protein conformer generation and pocket prediction are detailed  
201 in CONCOORD<sup>45</sup> and DeepPocket<sup>43</sup>, respectively. Briefly, each PDB chain and/or conformer  
202 structure was initially run on Fpocket<sup>44</sup> to calculate the barycenter of candidate pockets. Fpocket  
203 is a protein pocket (cavity) detection algorithm based on Voronoi tessellation and detects pocket  
204 curvatures in most protein structures with high accuracy. Subsequently, constant-sized grids are  
205 then placed at the barycenter of each candidate pocket and are scored using convolutional  
206 neural networks (CNN). A final 3D-shaped pocket structure of the top-ranked centers are then  
207 generated using a CNN segmentation model. The 3D-structure of pockets is constructed via  
208 voxelization, and the indices of the constructing voxel indices are converted to cartesian  
209 coordinates. Distance calculations are then taken between the surface of the binding pockets  
210 and the centers of the atoms of covalent residue(s) (Cys, Lys, Ser, Tyr, Thr, His, Asp or Glu).  
211 this study.

212 **Ligandability assessment.** Ligandability assessment is done to gauge the likelihood of  
213 ligandable cavity presence within rational distance of the residue within limits of standard small  
214 molecule dimensions. Whereas for the simplistic ligandability assessment, a distance threshold  
215 of 10Å was set for determining ligandability of a covalent site, to minimize false-positive rates,  
216 this threshold was set to 9Å for the advanced ligandability platform. The distance threshold is  
217 set based on the maximum of the shortest distance computed from pocket surface to covalent  
218 residue across three datasets; covalent sites of 1,647 covPDBs, 230 apo-covPDBs and mock  
219 sites of 90 cryptic-PDBs(**Figure 2**).

220

221

222

223

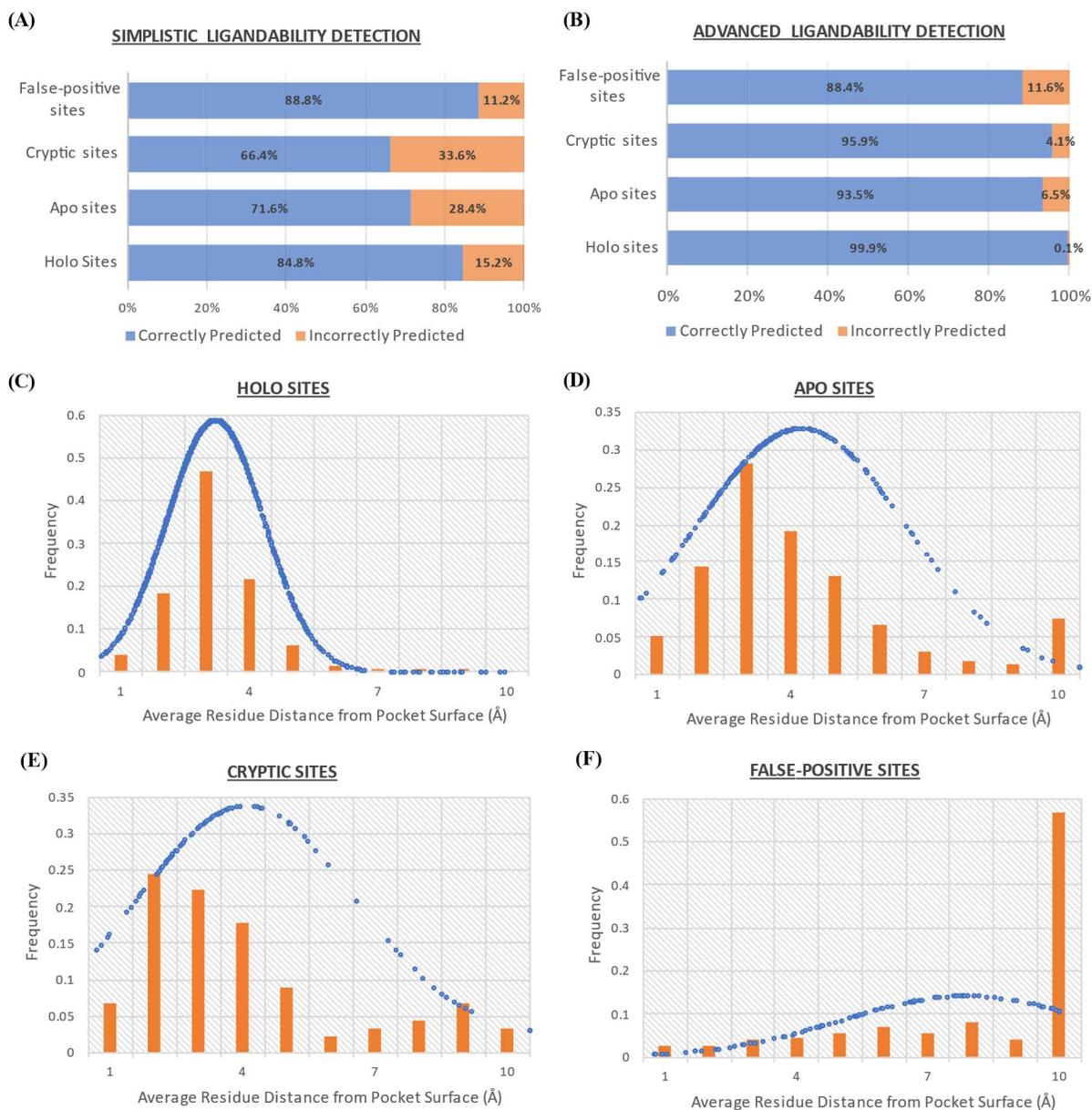
224

## 225 **Results**

### 226 **Simplistic Ligandability Assessment**

227 Using a list of empirically validated ligand binding sites, covalent site ligandability  
228 assessment platform was initially subject to performance validation using a simplistic design  
229 method of taking distance computations from a query residue to the nearby predicted pockets.  
230 Initially, four pools of PDB datasets were prepared; i) a list of empirically determined holo protein  
231 co-crystal structures with a covalent ligand bound at a ligandable sites ( "Holo sites"), ii) a list of  
232 experimentally solved apo counterpart structures of the holo protein list ("Apo sites"), iii) a list of  
233 cryptic proteins with a mock covalent site designated in their respective occluded binding cavity  
234 ("Cryptic sites") and, iv) a list of the most extreme surface residues lacking cavities or visible  
235 ligandable sites in their vicinities designates as "False-positive sites". The dataset used in  
236 ligandability detection power against holo sites comprised of 2,251 holo-protein structures from  
237 the covPDB database<sup>50</sup>, apo sites comprising of 248 apo protein structures, cryptic sites  
238 comprising of 118 cryptic protein structures from the from RCSB database<sup>51</sup>, and false-positive  
239 sites comprising of 112 false positive sites in the Cryptic-PDBs list (SupplementaryFile1).

240 Overall, covalent site ligandability assessment using a simple distance (10Å) threshold approach  
241 on rigid protein structural models demonstrated a competent detection performance of 85% when  
242 tested on holo site and satisfactorily (72%) when tested on apo site data (**Figure 2**). These scores  
243 correlate with formerly reported F-pocket and DeepPocket binding site predictions on holo and  
244 apo protein structures, respectively<sup>43,44</sup>. Generally, detection performance of many algorithms  
245 including F-pocket, Kalasanty, DeeplyTough, DeepSite, and DeepCoSI albeit considerate of  
246 signature structural features representing dynamic protein regions, exhibit suboptimal detection  
247 performance on holo sites<sup>43</sup> as demonstrated in this study. The inability of the simplistic approach  
248 to capture ligandability of apo sites could associated with multiple factors including mitigation of  
249 high false positive rates and inability to recognize unstructured features used in training the pocket  
250 detection algorithms. Apo sites are expected to possess lower solvent accessible surface area  
251 (SASA), unstructured architecture of the ligand cavity and imprecise positions of pocket-forming  
252 atoms from the ideal expected coordinates of the ligandable site. Despite their consideration of  
253 protein dynamics, these distinct features could pose challenges for algorithms trained on holo  
254 pockets, especially when identifying valid sites within apo structures.



255

256 **Figure 2: Performance metrics of the covalent ligandability assessment platforms.** Panels (A) and (B) depict the various types  
 257 of binding sites tested for correct prediction (filtering out false-positives or detecting experimental covalent binding sites) and incorrect  
 258 prediction (predicting false-positives as covalent binding sites or failing to detect experimental covalent binding sites. Predictions made  
 259 by (A) the simplistic platform designed on utilizing rigid protein structural models are compared to (B) predictions made by the  
 260 advanced platform. Panels (C-F) depict the distribution scatter of distances from the covalent residue centers to nearest pocket  
 261 surfaces in the various categories of dataset analyzed by the advanced platform. The bars in the distribution plots depict the relative  
 262 number (frequency) of sites that fall in residue-pocket distance interval bins of 1Å (i.e. 1 for 0-1Å bin, 2 for 1-2Å bin, 3 for 2-3Å bin,  
 263 etc.) in each of the dataset types. For simplicity, any distance higher than the set threshold of 9Å in the various datasets was placed  
 264 in the 9-10Å bin.

265

266

267 **Advanced Approach: Chemoproteomics Site Ligandability Assessment using an**  
268 **Amalgamated Platform**

269 The vital role of protein conformational sampling in detecting hidden or unformed covalent  
270 ligandable sites was assessed using an advanced platform designed to consider protein  
271 conformational fluctuations (**Figure 2**). When tested on the same data set of holo sites, the  
272 advanced covalent ligandability assessment platform performed remarkably by correctly detecting  
273 >99% of these sites. This distinct performance is similarly reflected on apo sites data where the  
274 advanced ligandability assessment platform was able to capture 94% of the covalent ligandable  
275 sites on apo protein structures. The binding sites on apo proteins correspond to the liganded site  
276 on their holo structure counterparts. Since applications of conformational sampling of proteins is  
277 expected to widen of existing pockets or reveal subpockets thereby decreasing average distance  
278 between pocket surface and covalently modified residue in the PDB structure (**Suppl. Figure 2**),  
279 a distance threshold of 9Å was used by the advanced ligandability assessment platform. Overall,  
280 in going from the simplistic distance computation approach (**Suppl. Figure 3**) to the advanced  
281 platform (**Figure 2**), a distance shift (0.5Å to 2.0Å) in the average distribution of residue-to-pocket  
282 distances and a shift in the relative number of sites towards the lower distance interval bins occur  
283 across the datasets. Decrease in the average residue-to-pocket distance incurred by opening or  
284 widening of pockets is more obvious when comparing apo and cryptic dataset across the two  
285 platforms.

286 The proportion of empirical holo and apo sites missed by the advanced platform was respectively  
287 0.1% and 6.5% compared to 11.2% and 28.4% using the simplistic ligandability assessment  
288 approach. Each “incorrectly predicted” holo site missed by either of the simplistic and advanced  
289 approaches was subject to visual inspection for presence of a ligandable cavity occupied by a  
290 ligand in the holo protein. For the incorrectly predicted apo sites; using the 3D atomic coordinate  
291 file of the subjective apo protein structure, the apo site in apo-covPDBs was aligned to its holo  
292 site counterpart from the covPDB pool and was visually inspected for experimental truncations  
293 and presence of residues fully forming the cavity in the ligand-complexed covPDB counterpart.  
294 Details on these corresponding covPDB IDs and their alignment analysis at the global scale as  
295 well as covalent ligand binding region can be found in SupplementaryFile1.

296

297

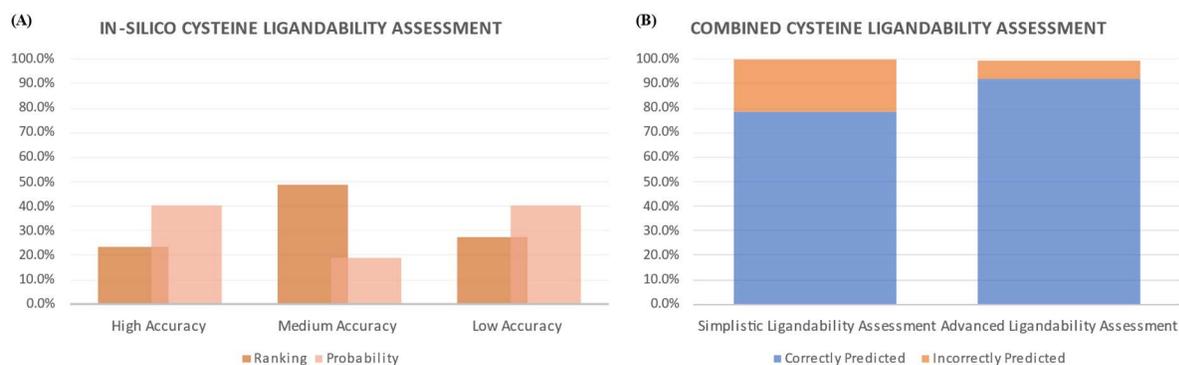
## 298 **Platform Validation on Ligandable Cryptic Sites**

299 In high probabilistic cases of ligandable covalent sites, binding of a covalent small molecule on a  
300 protein site is expected to induce structural change and potential opening or enlargement of a  
301 ligandable pocket or ostensibly a surface site. These structural transformations can give rise to  
302 discovery of hidden ligandable sites on proteins. Certainty around whether a previously  
303 unreported cavity detected in a conformational state as a false positive or, yet an undiscovered  
304 cryptic site is questionable. However, it is often possible to distinguish between superficial PPI-  
305 like sites and a buried cryptic binding site. In the past decade, many PPIs which once considered  
306 “undruggable” have now been successfully targeted by drug-like molecules<sup>54,55</sup>. Protein surface  
307 cavities and novel ligandable pockets are revealed from loop motions, secondary structure  
308 element motions and changes, and interdomain motions that increase SASA on the target. It is  
309 imaginable that such perturbations cause exposure of covalently reactive residues that may  
310 otherwise be buried as reported for most covalent cysteine sites.<sup>56–58</sup> This may give rise to  
311 indefinite estimates for setting a proximity distance threshold between covalent residue and  
312 pocket (**Figure 2, Suppl. Figure 1**). Notably, the power of the advanced covalent ligandability  
313 assessment platform is well demonstrated by its high detection scores on covalent sites in the  
314 apo-covPDB dataset including those which are buried (**Suppl. Table 1**) as well as in the cryptic-  
315 PDB dataset. Compared with a simplistic distance-based simplistic approach, the advanced  
316 ligandability assessment platform showcased an improved detection rate on missed cryptic-PDBs  
317 by 8-fold (33.6% vs 4.1%, **Figure 2**). This indicates that the advanced platform is able to capture  
318 a diversity of conformational changes that lead to formation of cryptic pockets while maintaining  
319 a comparable detection rate of false-positives sites with the simplistic ligandability assessment  
320 platform (**Suppl. Figure 3**).

## 321 **Eliminating Non-Ligandable Sites**

322 To assess the capacity of the advanced platform to successfully disqualify false-positive covalent  
323 sites and minimize the chance of a false-negatives in spite of increased positive detection rates,  
324 we selected a new mock list of false-positive sites using the apo-covPDB dataset. This list  
325 comprises of the most the superficial (solvent accessible) residue on the apo-covPDB structure  
326 whose identity matches the covalent ligand-bound residue in the CovPDB counterpart (Matching-  
327 Negative Residue column in the SupplementaryFile1). For instance, the covPDB structure 6CGE  
328 is covalent ligand-bound at HIS221 and the apo-covPDB counterpart of this protein is 1BHS. With  
329 reference to HIS221, HIS179 represents the most distant matching residue ID in 1BHS. Since

330 HIS179 remained unliganded in the experimental structure of the covalent ligand-bound structure  
 331 of the protein, it represents the most reliable matching false-positive site. Although a similar list of  
 332 matching false positives was created for the covPDBs and made available in SupplementaryFile  
 333 1, due to the large number of covPDB matching false positive sites and unfeasibility for visual  
 334 inspection to confirm correct/incorrect prediction for each site, we sufficed our discussion with  
 335 apo-covPDB matching false-positives (and Cryptic-PDB false-positives, **Figure 2**) evaluation  
 336 trials of the two platforms. Any matching false positive not found in the apo-covPDB structure or  
 337 not be superficially exposed although being most distant from the covalently modified residue was  
 338 omitted from analysis. The cryptic-PDB false positives list comprises residue sites on the most  
 339 extreme or superficial areas of the protein structure tested to ensure that false-positive detection  
 340 rates remain loyal when assessing non-ligandable surface residues (SupplementaryFile1).  
 341 Overall, although the simplistic platform demonstrates a false positive detection rate (11.2%,  
 342 **Suppl. Figure 3**) comparative to the advanced approach (11.6%), it does so at the expense of  
 343 missing real ligandable sites.



344  
 345 **Figure 3: Cysteine ligandability detection.** Covalent ligandability assessment against a list of buried cysteine residues from the  
 346 apo-covPDB dataset which are empirically confirmed to be ligandable was compared across two different methods. The first approach  
 347 (A) entirely dependent on predictions made by an ML algorithm<sup>32</sup> was able to rank 28% of the buried cysteines as first or second most  
 348 ligandable in the queried targets. The success rate of cysteine ligandability detection of the algorithm on the queried dataset is slightly  
 349 higher when taking the probability score it assigns to each cysteine (by setting 3 bins for clarity: high for probability score 1.0-1.2,  
 350 Medium for probability score of 0.8-1.0 and Low for less than 0.8). Details and exact score can be found in the SupplementaryFile1  
 351 (B) The second approach uses a combination of chemoproteomics-based (empirical) covalent sites and in-silico ligandability  
 352 assessment techniques presented in this work. In the combined approach, the simplistic ligandability approach was able to capture  
 353 78% of the ligandable cysteines compared to the advanced method which correctly judged 92% of the empirical sites.

354

355

356

## 357 **Cysteine Ligandability Detection: Computational versus Combined Methodologies**

358 Although various standalone residue reactivity prediction algorithms<sup>32,59</sup> take key ligandability  
359 features such as SASA and flexibility into account in detecting reactive residues, our experience  
360 suggests that protein conformation sampling and/or separate assessment for the presence a  
361 binding pocket in the vicinity of the predicted reactive residue site may be necessary for accurate  
362 results. Although the differentiating advantage of such an approach can be sensed from the data  
363 presented above, we further pursued a case in which we compared prediction performance of a  
364 purely computational method on a list of cysteine sites buried in an occluded binding pockets. A  
365 list of 47 apo-covPDB structures empirically confirmed to contain a buried ligandable cysteine  
366 covalent site were selected and tested for ligandability prediction by DeepCoSI<sup>32</sup> (apo-  
367 covPDBs\_buriedCys, SupplementaryFile1). DeepCoSI performs global analysis of all cysteine  
368 residues on each query PDB file and ranks the cysteines according to their “ligandability” score.  
369 Albeit considerate of SASA around the target residue, the DeepCoSI platform was able to capture  
370 at least 28% (high accuracy ranked) and at most 76% of the ligandable cysteines (the total of high  
371 and medium accuracy predictions using ranking-based scores). In contrast to ranking-based  
372 score, the proportion of top predicted ligandable cystines based on probability score of the  
373 algorithm matched that of the empirical data at least 19% and at most 40% of the time. (Figure  
374 3). Treating the same apo-covPDB samples as empirical cysteine chemoproteomics data points  
375 and assessing them for ligandability using the simplistic and advanced approaches correctly  
376 assessed 78% and 92% of the ligandable cysteines, respectively. As exemplified, covalent  
377 ligandability evaluation techniques purely dependent on computational methods may  
378 misrepresent a real ligandable site whereas that which uses the simplistic rigid protein models for  
379 site ligandability evaluation may fail to detect hidden or novel sites. By maximizing the number of  
380 correct predictions, the advanced ligandability approach which considers protein conformational  
381 fluctuations in assessing ligandable cavities around experimental sites is an ideal method to be  
382 adopted in site-directed and covalent drug discovery pipelines.

383

384

385

386

## 387 Discussion

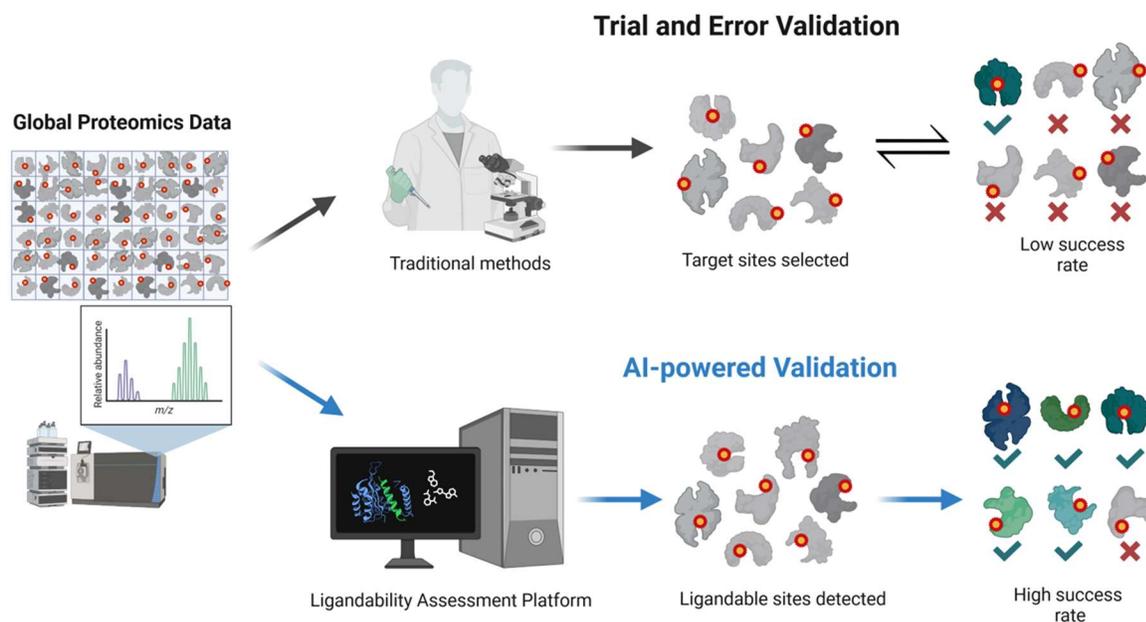
388 The identification of covalent site ligandability through machine learning presents a tremendous  
389 opportunity to expand the druggable protein space, especially those lacking well-defined binding  
390 sites. Our study aimed to introduce and validate an advanced computational platform, leveraging  
391 chemoproteomics studies, to accelerate covalent-based drug discovery processes. This platform  
392 integrates protein conformational fluctuations and structural dynamics, providing a robust scaffold  
393 for assessing and validating experimental chemoproteomics sites and identifying potential ligands  
394 for drug development.

395 In the evaluation of ligandability using a simplistic approach based on a distance threshold from  
396 pocket surfaces, our initial findings revealed a competent detection performance, particularly in  
397 detecting empirical holo sites. However, the simplistic approach exhibited limitations in identifying  
398 ligandable sites within apo structures, where structural irregularities and lower solvent accessible  
399 surface areas challenged the detection algorithm's efficacy. These results highlight the inherent  
400 challenges faced by algorithms trained on holo pockets when applied to detect valid sites within  
401 apo structures, emphasizing the need for methodologies that account for protein dynamics and  
402 conformational changes.

403 The majority of the human genome remains undrugged encoding thousands of proteins that  
404 have experimental evidence linking them to human disease. With only about 5% of the human  
405 proteome drugged and multiple targeted drugs available for a small number of driver gene  
406 targets, the vast majority of driver gene targets (>3500 for cancer driver) remain untargeted  
407 and/or inaccessible.<sup>60,61</sup> These proteins which need to be immediately explored for *de novo drug*  
408 discovery, represent a potential of multi-fold increase beyond the total number of available FDA-  
409 approved drugs existing today. Whereas membrane proteins make-up about 30% of all known  
410 proteins, over 60% of the current drug targets are membrane proteins.<sup>39</sup> The difficulty  
411 associated with targeting many of these proteins is the absence of well defined binding sites or  
412 lack of target structural information, albeit the latter is rapidly changing with recent advancement  
413 in cryo-EM techniques. The advanced covalent ligandability assessment platform, designed to  
414 consider protein conformational fluctuations, exhibited remarkable performance gains in the  
415 detection of covalent ligandable sites, surpassing the simplistic approach both in detecting holo  
416 sites and significantly outperforming in identifying apo sites. Notably, the advanced platform  
417 showcased a reduced rate of missed sites for both holo and apo structures compared to the  
418 simplistic approach, underscoring the pivotal role of considering protein structural fluctuations in

419 enhancing the accuracy of ligandability predictions. The platform also enables mass-scale and  
420 holistic analysis of proteomics data with unprecedented speed (<1 second/site) using a  
421 multidisciplinary approach aimed at expediting drug discovery research and development.  
422 Furthermore, the computational platform designed in this study can be tailored for detection of  
423 non-covalent sites by providing platform with mock chemoproteomics sites proximal to expected  
424 non-covalent binding sites.

425 Chemical proteomics-based drug discovery research typically suggests a large number (1000s)  
426 of targets of which 1 to 5 potential targets are selected based on target quality for tractability  
427 analysis using wet-lab biochemical and biophysical validation assays. Target selection and  
428 prioritization is typically done using computational biology approaches that qualify targets using  
429 curated omics data available to us. Target prioritization is a long-standing issue in drug discovery  
430 research especially due to lack of techniques that assess binding site quality and downstream  
431 functional impact of targeting each site<sup>62</sup>. Target validation is highly time consuming (est.  
432 >1year/target) and is a monetary sink (est. >\$1M/target validation) in the covalent drug discovery  
433 pipelines. Quite often, such tests eventually indicate that one or none of the elected targets are  
434 true therapeutic targets and are indeed identified as off-targets and/or false positives. Aside from  
435 being highly time and resource wasteful, this traditional approach of target validation misses on a  
436 vast opportunity that the remaining unvalidated targets pose. Ligandability evaluation techniques  
437 showcased in this work provide an opportunity to explore more than 75% of known proteome  
438 (767,580 protein structures)<sup>63,64</sup> using chemoproteomics-enabled covalent drug discovery  
439 experiments on a time and financial scale that is at least 100-fold less than the traditional  
440 approaches (Figure 4). To ensure the platform's efficacy in eliminating false-positive covalent  
441 sites, computational tools were employed to identify and disqualify superficial matching-negative  
442 residues in apo-covPDB structures. The comprehensive assessment revealed the platform's  
443 robustness in minimizing false-positive rates, crucial for ensuring the accuracy of ligandability  
444 predictions.



445

446 **Figure 4:** The traditional approach of target prioritization involves an exhaustive feedback loop of target selection and in-vitro  
 447 validation, quite often leading to pursuance of many false-positive hits. In the advanced approach powered by AI-based computational  
 448 platforms, target sites are assessed for ligandability and downstream functionality in-silico before attempting any in-vitro validation.  
 449 Although targets picked by an algorithm are still bound for in-vitro validation, the number of necessary validation and their success  
 450 rate is highly improved.

451 Application of such ligandability assessment platforms can be very broad, ranging from  
 452 development of novel QMS-based techniques and site-directed fragment-based screening to  
 453 identifying variants adjacent to novel ligandable pockets for allele-specific drug binding<sup>65</sup>. Our  
 454 results indicated that as of today, a conglomeration of empirical and computational methods  
 455 may provide the highest research throughput. We report for the first time a list of apo  
 456 counterparts of covalent-ligand bound PDB structures, which are of highly value for further  
 457 development of platforms that can accurately recognize both apo and holo covalent sites.  
 458 Although the work here is focused on a platform more suitable for a bottom-up drug discovery  
 459 approach, the data reported in this paper can be utilized in training and developing singly-  
 460 sufficient computational platforms that, through a top-down approach, accurately predict valid  
 461 covalent ligandable sites on target proteins (Suppl. Figures 3 and 4).

462 Moreover, the platform demonstrated its proficiency in identifying cryptic binding sites, which  
 463 typically remain hidden due to structural transformations induced by ligand binding. The  
 464 advanced platform's ability to detect a diverse range of conformational changes leading to the  
 465 formation of cryptic pockets was significantly higher compared to the simplistic approach,  
 466 highlighting its capability to unveil hidden ligandable sites while maintaining a low false-positive  
 467 rate.

468 In conclusion, the developed computational platform, by harnessing protein conformational  
469 sampling, represents a paradigm shift in covalent site ligandability evaluation. Its ability to  
470 consider dynamic structural changes, detect cryptic pockets, minimize false-positive rates, and  
471 enhance the accuracy of ligandability predictions stands as a significant advancement in  
472 computational biology. This platform holds immense promise in expediting drug discovery  
473 research and development by providing a holistic and multidisciplinary approach for identifying  
474 potential therapeutic targets and facilitating the design of novel covalent drugs.

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

## 490 **Acknowledgements**

491 We thank the IT department of the University of Toronto Mississauga for their support in providing  
492 infrastructure for data storage. We thank the Department of Mathematical and Computational  
493 Sciences at the University of Toronto Mississauga for providing access to computational talent  
494 and resources. We especially thank our data analysis volunteers Hrithik Kumar Advani, Sibthain  
495 Kazmi, Ahmed Nasir, and advisors Tudor B. Radu, Lisa Zhang and Sirano Dhe-Paganon for their  
496 valuable support and commitment to the development of this project. We thank Digital Research  
497 Alliance of Canada for providing access to supercomputing resources and supporting the  
498 development of this project.

499 P.T.G. is supported by research grants from NSERC (RGPIN-2014-05767), CIHR (MOP-130424,  
500 MOP-137036), Canada Research Chair (950-232042), Canadian Cancer Society (703963),  
501 Canadian Breast Cancer Foundation (705456) and infrastructure grants from CFI (33536) and  
502 the Ontario Research Fund (34876).

## 503 **Notes**

504 Visualizations are generated using Biorender.com platform.

505 The authors declare the following competing financial interest(s): F.E., is co-founder and chief  
506 scientific officer of Gene2Lead, A.D.E, and P.T.G are co-founders of Gene2Lead.

507

508

509

510

511

512

## 513 **References**

- 514 1. Westbrook, J. D. & Burley, S. K. How Structural Biologists and the Protein Data Bank  
515 Contributed to Recent FDA New Drug Approvals. *Structure* **27**, 211–217 (2019).
- 516 2. De Vita, E. 10 years into the resurgence of covalent drugs. *Future Med Chem* **13**, 193–210  
517 (2021).
- 518 3. Singh, J., Petter, R. C., Baillie, T. A. & Whitty, A. The resurgence of covalent drugs. *Nat*  
519 *Rev Drug Discov* **10**, 307–317 (2011).
- 520 4. Baillie, T. A. Targeted Covalent Inhibitors for Drug Design. *Angewandte Chemie*  
521 *International Edition* **55**, 13408–13421 (2016).
- 522 5. Abdeldayem, A., Raouf, Y. S., Constantinescu, S. N., Moriggl, R. & Gunning, P. T.  
523 Advances in covalent kinase inhibitors. *Chem. Soc. Rev.* **49**, 2617–2687 (2020).
- 524 6. Hallin, J. *et al.* The KRASG12C Inhibitor MRTX849 Provides Insight toward Therapeutic  
525 Susceptibility of KRAS-Mutant Cancers in Mouse Models and Patients. *Cancer Discov* **10**,  
526 54–71 (2020).
- 527 7. Canon, J. *et al.* The clinical KRAS(G12C) inhibitor AMG 510 drives anti-tumour immunity.  
528 *Nature* **575**, 217–223 (2019).
- 529 8. Gehringer, M. & Laufer, S. A. Emerging and Re-Emerging Warheads for Targeted Covalent  
530 Inhibitors: Applications in Medicinal Chemistry and Chemical Biology. *J Med Chem* **62**,  
531 5673–5724 (2019).
- 532 9. Bradshaw, J. M. *et al.* Prolonged and tunable residence time using reversible covalent kinase  
533 inhibitors. *Nat Chem Biol* **11**, 525–531 (2015).
- 534 10. Halgren, T. A. Identifying and Characterizing Binding Sites and Assessing Druggability. *J.*  
535 *Chem. Inf. Model.* **49**, 377–389 (2009).

- 536 11. Marx, M. A. *et al.* Abstract B30: Structure-based drug discovery of MRTX1257, a selective,  
537 covalent KRAS G12C inhibitor with oral activity in animal models of cancer. *Molecular*  
538 *Cancer Research* **18**, B30 (2020).
- 539 12. Feng, Y., Duan, W., Cu, X., Liang, C. & Xin, M. Bruton's tyrosine kinase (BTK) inhibitors  
540 in treating cancer: a patent review (2010-2018). *Expert Opinion on Therapeutic Patents* **29**,  
541 217–241 (2019).
- 542 13. da Cunha-Bang, C. & Niemann, C. U. Targeting Bruton's Tyrosine Kinase Across B-Cell  
543 Malignancies. *Drugs* **78**, 1653–1663 (2018).
- 544 14. Byrd, J. C. *et al.* Targeting BTK with ibrutinib in relapsed chronic lymphocytic leukemia. *N*  
545 *Engl J Med* **369**, 32–42 (2013).
- 546 15. Hoffer, L. *et al.* CovaDOTS: In Silico Chemistry-Driven Tool to Design Covalent Inhibitors  
547 Using a Linking Strategy. *J Chem Inf Model* **59**, 1472–1485 (2019).
- 548 16. Barf, T. & Kaptein, A. Irreversible protein kinase inhibitors: balancing the benefits and risks.  
549 *J Med Chem* **55**, 6243–6262 (2012).
- 550 17. Raouf, Y. S. *et al.* Discovery of YSR734: A Covalent HDAC Inhibitor with Cellular Activity  
551 in Acute Myeloid Leukemia and Duchenne Muscular Dystrophy. *J. Med. Chem.* (2023)  
552 doi:10.1021/acs.jmedchem.3c01236.
- 553 18. Petri, L. *et al.* A covalent strategy to target intrinsically disordered proteins: Discovery of  
554 novel tau aggregation inhibitors. *Eur J Med Chem* **231**, 114163 (2022).
- 555 19. Lonsdale, R. & Ward, R. A. Structure-based design of targeted covalent inhibitors. *Chem*  
556 *Soc Rev* **47**, 3816–3830 (2018).
- 557 20. Wong, C. H., Siah, K. W. & Lo, A. W. Estimation of clinical trial success rates and related  
558 parameters. *Biostatistics* **20**, 273–286 (2019).

- 559 21. Hu, S., Chen, P., Gu, P. & Wang, B. A Deep Learning-Based Chemical System for QSAR  
560 Prediction. *IEEE J Biomed Health Inform* **24**, 3020–3028 (2020).
- 561 22. Muegge, I. & Oloff, S. Advances in virtual screening. *Drug Discov Today Technol* **3**, 405–  
562 411 (2006).
- 563 23. Zhu, K. *et al.* Docking Covalent Inhibitors: A Parameter Free Approach To Pose Prediction  
564 and Scoring. *J. Chem. Inf. Model.* **54**, 1932–1940 (2014).
- 565 24. Rachman, M. *et al.* DUckCov: a Dynamic Undocking-Based Virtual Screening Protocol for  
566 Covalent Binders. *ChemMedChem* **14**, 1011–1021 (2019).
- 567 25. Wei, L. *et al.* Cov\_DOX: A Method for Structure Prediction of Covalent Protein–Ligand  
568 Bindings. *J. Med. Chem.* **65**, 5528–5538 (2022).
- 569 26. Scarpino, A. *et al.* WIDOCK: a reactive docking protocol for virtual screening of covalent  
570 inhibitors. *J Comput Aided Mol Des* **35**, 223–244 (2021).
- 571 27. Bianco, G. *et al.* Reactive Docking: A Computational Method for High-Throughput Virtual  
572 Screenings of Reactive Species. *J. Chem. Inf. Model.* **63**, 5631–5640 (2023).
- 573 28. Palazzesi, F. *et al.* Bireactive: A Machine-Learning Model to Estimate Covalent Warhead  
574 Reactivity. *J Chem Inf Model* **60**, 2915–2923 (2020).
- 575 29. Krivák, R. & Hoksza, D. P2Rank: machine learning based tool for rapid and accurate  
576 prediction of ligand binding sites from protein structure. *Journal of Cheminformatics* **10**, 39  
577 (2018).
- 578 30. An, J., Totrov, M. & Abagyan, R. Pocketome via comprehensive identification and  
579 classification of ligand binding envelopes. *Mol Cell Proteomics* **4**, 752–761 (2005).
- 580 31. Meller, A. *et al.* Predicting locations of cryptic pockets from single protein structures using  
581 the PocketMiner graph neural network. *Nat Commun* **14**, 1177 (2023).

- 582 32. Du, H. *et al.* Proteome-Wide Profiling of the Covalent-Druggable Cysteines with a  
583 Structure-Based Deep Graph Learning Network. *Research (Wash D C)* **2022**, 9873564  
584 (2022).
- 585 33. Wang, H. *et al.* Sequence-Based Prediction of Cysteine Reactivity Using Machine Learning.  
586 *Biochemistry* **57**, 451–460 (2018).
- 587 34. Gao, M. & Günther, S. HyperCys: A Structure- and Sequence-Based Predictor of Hyper-  
588 Reactive Druggable Cysteines. *International Journal of Molecular Sciences* **24**, 5960 (2023).
- 589 35. Piazza, I. *et al.* A machine learning-based chemoproteomic approach to identify drug targets  
590 and binding sites in complex proteomes. *Nat Commun* **11**, 4200 (2020).
- 591 36. Lanning, B. R. *et al.* A road map to evaluate the proteome-wide selectivity of covalent  
592 kinase inhibitors. *Nat Chem Biol* **10**, 760–767 (2014).
- 593 37. Browne, T., Concheiro-Guisan, M. & Prinz, M. Semi quantitative detection of signature  
594 peptides in body fluids by liquid chromatography tandem mass spectrometry (LC–MS/MS).  
595 *Forensic Science International: Genetics Supplement Series* **7**, 208–210 (2019).
- 596 38. Lanman, B. A. *et al.* Discovery of a Covalent Inhibitor of KRASG12C (AMG 510) for the  
597 Treatment of Solid Tumors. *J Med Chem* **63**, 52–65 (2020).
- 598 39. Maveyraud, L. & Mourey, L. Protein X-ray Crystallography and Drug Discovery. *Molecules*  
599 **25**, 1030 (2020).
- 600 40. Boatner, L. M., Palafox, M. F., Schweppe, D. K. & Backus, K. M. CysDB: a human cysteine  
601 database based on experimental quantitative chemoproteomics. *Cell Chem Biol* S2451-  
602 9456(23)00090–9 (2023) doi:10.1016/j.chembiol.2023.04.004.

- 603 41. Mortenson, D. E. *et al.* “Inverse Drug Discovery” Strategy To Identify Proteins That Are  
604 Targeted by Latent Electrophiles As Exemplified by Aryl Fluorosulfates. *J. Am. Chem. Soc.*  
605 **140**, 200–210 (2018).
- 606 42. Zhang, Y., Fonslow, B. R., Shan, B., Baek, M.-C. & Yates, J. R. Protein Analysis by  
607 Shotgun/Bottom-up Proteomics. *Chem Rev* **113**, 2343–2394 (2013).
- 608 43. Aggarwal, R., Gupta, A., Chelur, V., Jawahar, C. V. & Priyakumar, U. D. DeepPocket:  
609 Ligand Binding Site Detection and Segmentation using 3D Convolutional Neural Networks.  
610 *J. Chem. Inf. Model.* (2021) doi:10.1021/acs.jcim.1c00799.
- 611 44. Le Guilloux, V., Schmidtke, P. & Tuffery, P. Fpocket: an open source platform for ligand  
612 pocket detection. *BMC Bioinformatics* **10**, 168 (2009).
- 613 45. de Groot, B. L. *et al.* Prediction of protein conformational freedom from distance constraints.  
614 *Proteins* **29**, 240–251 (1997).
- 615 46. Yang, Y., Kucukkal, T. G., Li, J., Alexov, E. & Cao, W. Binding Analysis of Methyl-CpG  
616 Binding Domain of MeCP2 and Rett Syndrome Mutations. *ACS Chem Biol* **11**, 2706–2715  
617 (2016).
- 618 47. Patschull, A. O. M., Gooptu, B., Ashford, P., Daviter, T. & Nobeli, I. In Silico Assessment  
619 of Potential Druggable Pockets on the Surface of  $\alpha$ 1-Antitrypsin Conformers. *PLoS One* **7**,  
620 e36612 (2012).
- 621 48. Zheng, H., Hou, J., Zimmerman, M. D., Wlodawer, A. & Minor, W. The future of  
622 crystallography in drug discovery. *Expert Opin Drug Discov* **9**, 125–137 (2014).
- 623 49. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular  
624 biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).

- 625 50. Gao, M., Moumbock, A. F. A., Qaseem, A., Xu, Q. & Günther, S. CovPDB: a high-  
626 resolution coverage of the covalent protein-ligand interactome. *Nucleic Acids Res* **50**, D445–  
627 D450 (2022).
- 628 51. Berman, H. M. *et al.* The Protein Data Bank. *Nucleic Acids Research* **28**, 235–242 (2000).
- 629 52. Cimermancic, P. *et al.* CryptoSite: Expanding the Druggable Proteome by Characterization  
630 and Prediction of Cryptic Binding Sites. *Journal of Molecular Biology* **428**, 709–719 (2016).
- 631 53. Berman, H. M. *et al.* The Protein Data Bank. *Nucleic Acids Research* **28**, 235–242 (2000).
- 632 54. Scott, D. E., Bayly, A. R., Abell, C. & Skidmore, J. Small molecules, big targets: drug  
633 discovery faces the protein-protein interaction challenge. *Nat Rev Drug Discov* **15**, 533–550  
634 (2016).
- 635 55. Sijbesma, E. *et al.* Site-Directed Fragment-Based Screening for the Discovery of Protein-  
636 Protein Interaction Stabilizers. *J Am Chem Soc* **141**, 3524–3531 (2019).
- 637 56. Romany, A., Liu, R., Zhan, S., Clayton, J. & Shen, J. Analysis of the ERK Pathway  
638 Cysteinome for Targeted Covalent Inhibition of RAF and MEK Kinases. *J Chem Inf Model*  
639 **63**, 2483–2494 (2023).
- 640 57. Harris, R. C., Liu, R. & Shen, J. Predicting Reactive Cysteines with Implicit-Solvent-Based  
641 Continuous Constant pH Molecular Dynamics in Amber. *J Chem Theory Comput* **16**, 3689–  
642 3698 (2020).
- 643 58. Marino, S. M. & Gladyshev, V. N. Cysteine function governs its conservation and  
644 degeneration and restricts its utilization on protein surfaces. *J Mol Biol* **404**, 902–916 (2010).
- 645 59. Ustach, V. D. *et al.* Optimization and Evaluation of Site-Identification by Ligand  
646 Competitive Saturation (SILCS) as a Tool for Target-Based Ligand Optimization. *J. Chem.*  
647 *Inf. Model.* **59**, 3018–3035 (2019).

- 648 60. Warner, K. D., Hajdin, C. E. & Weeks, K. M. Principles for targeting RNA with drug-like  
649 small molecules. *Nat Rev Drug Discov* **17**, 547–558 (2018).
- 650 61. Sjöstedt, E. *et al.* An atlas of the protein-coding genes in the human, pig, and mouse brain.  
651 *Science* **367**, eaay5947 (2020).
- 652 62. K. Brown, K. *et al.* Approaches to target tractability assessment – a practical perspective.  
653 *MedChemComm* **9**, 606–613 (2018).
- 654 63. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language  
655 model. *Science* **379**, 1123–1130 (2023).
- 656 64. Suzek, B. E. *et al.* UniRef clusters: a comprehensive and scalable alternative for improving  
657 sequence similarity searches. *Bioinformatics* **31**, 926–932 (2015).
- 658 65. Nichols, C. A. *et al.* Loss of heterozygosity of essential genes represents a widespread class  
659 of potential cancer vulnerabilities. *Nat Commun* **11**, 2517 (2020).

660

661

662