

Active learning driven prioritisation of compounds from on-demand libraries targeting the SARS-CoV-2 main protease

Ben Cree,^a Mateusz K. Bieniek, Siddique Amin, Akane Kawamura, and Daniel J.

1

Cole*

School of Natural and Environmental Sciences, Newcastle University, Newcastle upon Tyne

NE1 7RU, United Kingdom

E-mail: daniel.cole@ncl.ac.uk

^aB.C. and M.K.B. contributed equally to this work.

Abstract

FEgrow is an open-source software package for building congeneric series of compounds in protein binding pockets. For a given ligand core and receptor structure, it employs hybrid machine learning / molecular mechanics potential energy functions to optimise the bioactive conformers of supplied linkers and functional groups. Here, we introduce significant new functionality to automate, parallelise and accelerate the building and scoring of compound suggestions, such that it can be used for automated de novo design. We interface the workflow with active learning to improve the efficiency of searching the combinatorial space of possible linkers and functional groups, make use of interactions formed by crystallographic fragments in scoring compound designs, and introduce the option to seed the chemical space with molecules available from on-demand chemical libraries. As a test case, we target the main protease of SARS-CoV-2, identifying several small molecules with high similarity to molecules discovered by the COVID Moonshot effort, using only structural information from a fragment screen in a fully automated fashion. Finally, we order and test 19 compound designs, of which three show weak activity in a fluorescence-based Mpro assay, but work is needed to further optimise the prioritisation of compounds for purchase.

19 Introduction

20 Recent advances in structural biology, from sample preparation, to synchrotron infrastructure
21 and data analysis pipelines, have transformed the throughput of protein-ligand complexes
22 available to inform drug discovery campaigns.¹ When soaked with carefully designed com-
23 pound libraries,² the numbers of small molecule (or fragment) structural hits can reach 10s
24 or 100s against a single therapeutic target.³ A frequently employed next step is to attempt to
25 grow and/or link the hit compounds, using either custom synthesis² or ordering from cata-
26 logues of purchasable compounds.^{4,5} However, chemical space is vast such that even choosing
27 follow-up compounds for purchase from on-demand libraries, such as the readily accessible
28 (REAL) Enamine database⁶ (> 5.5 bn compounds in 2022), becomes highly non-trivial.⁷

29 As such, attention is turning to cheminformatics and machine learning based algorithms
30 for structure-based *de novo* hit expansion, linking and merging.⁸ A wide range of approaches
31 are available to build from initial structural biology data, including DeepFrag⁹ that identifies
32 promising fragments for addition to an input bound ligand, using a deep convolutional
33 neural network, and DEVELOP¹⁰ that combines 3D pharmacophoric constraints from the
34 binding pocket with a graph-based deep generative model for R-group and linker design.
35 The SILVR method enables an equivariant diffusion model to be conditioned to generate
36 molecules based on a reference structure, such as a fragment from a crystallographic screen.¹¹
37 The V-SYNTHESIS approach makes use of on-demand libraries for hit-finding by decomposing
38 compounds from purchasable databases into reactive scaffolds and synthons, and using the
39 highest scoring docked fragments as seeds for further growth.¹² One particularly noteworthy
40 example is the use of fragment merging to design hits against the nonstructural protein 3
41 (NSP3) of the severe acute respiratory syndrome-coronavirus-2 (SARS-CoV-2).¹³ Fragments
42 from a crystallographic screen were merged using the Fragmenstein package,¹⁴ ensuring
43 placement of molecular substructures onto the original fragments, and subsequently used
44 as templates for searching on-demand chemical space. In this way, fragments were rapidly
45 elaborated into a 0.4 μ M hit (representing a >400-fold improvement in affinity).

46 While extremely promising, all of the above *de novo* design approaches suffer from some
47 combination of the following issues: i) reliance on an approximate classical molecular me-
48 chanics force field or knowledge-based algorithm for generating and optimising binding poses,
49 ii) use of an approximate objective function (usually a docking score) as a surrogate measure
50 of binding affinity, iii) approximation of a rigid target receptor structure, and iv) limited syn-
51 thetic tractability of the designed compounds. We therefore developed the FEgrow software
52 as an open-source, interactive Jupyter notebook based workflow for building user-defined
53 congeneric series of ligands in protein binding pockets to start to address some of these
54 open questions (Figure 1A).¹⁵ FEgrow grows user-defined functional groups (R-groups) off
55 a constrained core of a known hit compound, thus incorporating input from structural bi-
56 ology and the expertise of the user in selecting synthetically tractable elaborations. Since
57 publication, we have added functionality for connecting R-groups to the core via a flexible
58 linker, which can be chosen from a library of those common to bioactive molecules.¹⁶ In this
59 way, users can choose from 1M+ combinations of linker and R-group from our distributed
60 libraries (or upload their own R-group modifications). The modular workflow allows for the
61 incorporation of state-of-the-art molecular modelling algorithms, such as the use of hybrid
62 machine learning / molecular mechanics potential energy functions to optimise the ligand
63 binding pose,^{17,18} and the gnina convolutional neural network scoring function to predict
64 the binding affinity.¹⁹ We plan to expand the range of available optimisation algorithms and
65 scoring functions as they become available (see Methods Section).

66 While interactive work is useful for small-scale studies, we have found it useful to au-
67 tomate the workflow for use on high performance computing (HPC) clusters, and since
68 publication have added an application programming interface (API) to FEgrow (Figure 1B).
69 This enables us to build virtual libraries with a common core, for example, using reaction-
70 based generative scaffold decoration with LibInvent²¹ or substructure searching of compound
71 libraries,²² and then rapidly build the compounds into the protein binding pocket with FE-
72 grow. However, unless the libraries are designed using information from the binding pocket,

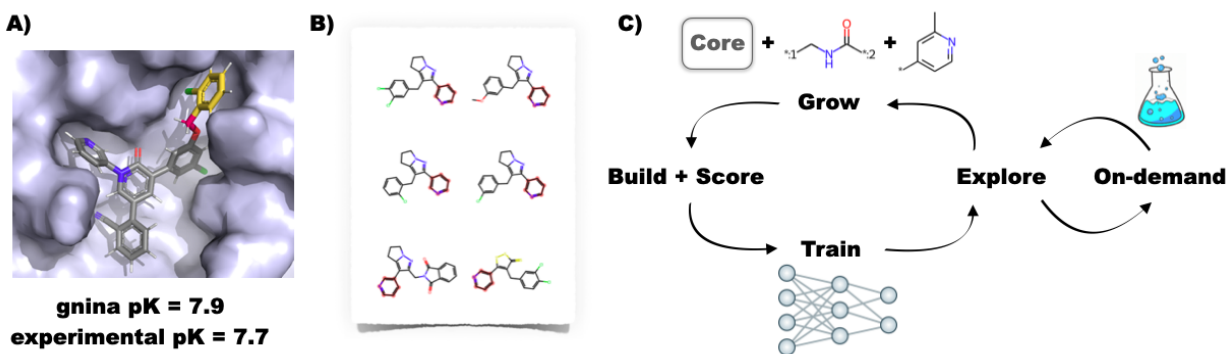


Figure 1: A) Example building and scoring of a SARS-CoV-2 inhibitor²⁰ using the interactive FEgrow workflow.¹⁵ The fixed core (grey) is extended using a user-defined, flexible linker (pink) and R-group (yellow), and scored using gnina.¹⁹ B) Compound libraries with substructures that match the rigid core can now be automatically grown and scored, treating the rest of the molecule as fully flexible. C) Proposed active learning cycle. Compounds are grown, built in the binding pocket and scored with FEgrow. The outputs are used to train a machine learning model, which is used to select the next batch of compounds. Optionally, the chemical space can be seeded using compounds available from on-demand chemical libraries.

73 time is wasted building and scoring compounds that are unlikely to be beneficial and it is
 74 still not feasible to routinely scan all possibilities.

75 Hence, rather than exhaustive or random searches of chemical space, we investigate here
 76 the use of active learning to elaborate compound design with FEgrow. The general idea
 77 behind this approach is that a subset of compounds is evaluated using an expensive design
 78 objective function (in this case the molecular growing and scoring algorithms in FEgrow)
 79 and used to train a machine learning model (Figure 1C).²³ The machine learning model
 80 then predicts the objective function for the remainder of the chemical space, and the next
 81 subset of molecules is picked for evaluation (for example, in order to optimise the objective
 82 or further explore the chemical space). By cycling through this procedure, the algorithm
 83 can iteratively make up for any lack of diversity in the initial training subset, and it has
 84 been found previously that the most promising compounds can be identified by evaluating
 85 only a fraction of the total chemical space.

86 Several studies have investigated the effects of choices such as machine learning algorithm,

87 sample selection protocol and total dataset size on active learning efficiency for experimen-
88 tal and computational affinity predictions.^{24–28} In general, active learning has been shown
89 to increase enrichment of hits compared to either random or one-shot training of a ma-
90 chine learning model, at low additional cost, and to be relatively insensitive to choices of
91 molecular representation, model hyperparameters and initial training subsets. Active learn-
92 ing has shown practical utility in prioritising compounds based on objective functions from
93 docking^{29–31} or free energy calculations.^{25,26,32,33}

94 Here, we interface FEgrow with active learning to efficiently search the chemical space of
95 linkers and R-groups from a user-defined vector. As well as using a docking score to guide
96 optimisation, we also experiment with functions that combine other molecular properties,
97 such as molecular weight, and 3D structural information, such as protein-ligand interaction
98 profiles (PLIP).³⁴ To address the issue of synthetic tractability of the compound designs, we
99 combine the workflow with regular searches of the Enamine REAL database to ‘seed’ the
100 chemical search space with promising purchasable compounds. After testing and optimising
101 the hyperparameters of the active learning models, we apply the algorithm to the prospective
102 design of inhibitors of the main protease (MPro) of SARS-CoV-2, the virus responsible for
103 the COVID-19 pandemic. This target has undergone extensive study in recent years. The
104 COVID Moonshot Consortium used open science crowd-sourced designs, in combination
105 with high-throughput structural biology and assays, free energy calculations, and machine
106 learning driven synthetic route predictions, to generate a series of potent inhibitors.⁴ Other
107 notable approaches that include biological confirmation of hits have employed, for exam-
108 ple, structure-based design starting from a drug repurposing study,²⁰ virtual screening of
109 a curated collection of commercially available compounds,³⁵ a deep reinforcement learning
110 model using pharmacophore and substructure matches with known inhibitors,³⁶ and a deep
111 generative framework using only target sequence information as input (along with priori-
112 tisation based on factors such as docking score and retrosynthetic feasibility).³⁷ Here we
113 employ active learning to prioritise compounds for purchase and testing from the Enamine

114 REAL database based only on early fragment hits. We suggest several novel designs that
115 show activity in a fluorescence-based Mpro assay, as well as automatically generating several
116 compounds that show high similarity to known Moonshot hits.

117 **Methods**

118 **Workflow Design**

119 The FEgrow software package is described in detail elsewhere.¹⁵ Briefly, FEgrow aims to
120 grow a ligand within a protein binding pocket, starting from a provided receptor structure,
121 ligand core and growth vector (Figure 1A). Libraries comprising 2000 linkers¹⁶ and around
122 500 R-groups, are provided, or users can supply their own. Merging is achieved using the
123 RDKit package,³⁸ which also generates an ensemble of ligand conformations via the ETKDG
124 algorithm,³⁹ with the atoms of the core strongly restrained to the input structure. That is,
125 the default behaviour is to allow flexibility only in the regions of the grown linkers and
126 R-groups. The ensemble of ligand structures is filtered to remove any that clash with the
127 protein, and the remaining conformers are structurally optimised in the context of a rigid
128 protein binding pocket using the OpenMM software.¹⁸ During energy minimisation, the
129 protein is treated using the AMBER FF14SB force field,⁴⁰ while intramolecular energetics
130 of the ligand are described, where possible, using the ANI-2x machine learning potential.¹⁷
131 Non-bonded interactions between the protein and ligand are described using a mechanical
132 embedding scheme, that is, they use electrostatics and Lennard-Jones terms described by
133 either the Open Force Field ‘Sage’⁴¹ or GAFF2⁴² general force fields. The goal of this
134 hybrid machine learning / molecular mechanics approach is to correct for known deficiencies
135 in potential energy surfaces of classical force fields, while ensuring that optimisations are
136 significantly faster than using full QM/MM.

137 The lowest energy structures are then output for scoring. In the first iteration of FEgrow,
138 we used the gnina convolutional neural network (CNN), which has been jointly trained on

139 binding pose and affinity prediction.^{19,43,44} We showed that the gnina ‘CNNAffinity’ scores
140 (predicted pK) correlated reasonably well with experiment for ten series of congeneric in-
141 hibitors built using FEgrow.¹⁵ Here, we add further options for scoring molecules based on
142 protein-ligand interaction profile (PLIP),³⁴ molecular properties, or a combination thereof.
143 For construction of the PLIP score, interactions formed in the available protein-fragment
144 complex crystal structures were one-hot encoded to form a reference vector of desired inter-
145 actions (here, hydrophobic, hydrogen-bonding, π -stacking, and salt bridge were all identi-
146 fied). A similar vector was constructed for the designed de novo compound, and its Tanimoto
147 similarity to the reference vector used as the objective for optimisation. It has been argued
148 that combining information from various properties can also be advantageous,⁸ for example
149 by using pharmacophore constraints in combination with docking scores, and we make use
150 here of a simple, combined score (CS):

$$CS = \left(\frac{pK}{MW} \right) \times \left(\frac{PLIP}{0.3} \right) \times 100 \quad (1)$$

151 which aims to maximise the predicted gnina affinity (pK) and the protein-ligand interaction
152 profile (PLIP) similarity to reference structures, while keeping the molecular weight (MW)
153 low.

154 Active Learning

155 Active learning²³ is a subset of machine learning that is based on iteratively labeling data
156 points from an unlabeled dataset (in our case, de novo compounds that are built into protein
157 binding pockets and scored). The aim is to pick the most useful samples for training a
158 surrogate model, whilst ultimately minimising the potentially expensive computation needed
159 to find instances that maximise an objective function. There are two main components to an
160 active learning workflow: the regression model, and the acquisition function. Every scored
161 instance is used to train a specified machine learning model, with more examples refining

162 the model accuracy, which is then used to select new molecules to be built. In this work we
163 consider and benchmark two models.

164 The first approach is gradient boosting machine (GBM), which is a random forest based
165 technique, utilising ensembles of decision trees. These trees are created from random sub-
166 sets of features (fingerprints), that are then used to make predictions. GBMs expand on
167 traditional random forests by using the gradient of the error to construct trees specifically
168 designed to minimise this error. Gradually increasing the number of relatively poor individ-
169 ual trees additively increases their predictive power (hence ‘gradient boosted’). The second
170 model is Gaussian Process (GP) regression, which is a Bayesian approach that makes predic-
171 tions by assuming observations can be modelled by the probability distribution over possible
172 reasonable (Gaussian) functions.⁴⁵ These Gaussian distributions are iteratively refined by the
173 observation of new samples. Because model prediction is performed via a probability distri-
174 bution, it natively incorporates uncertainty and other useful quantities, such as estimates of
175 expected improvement of a given new sample.⁴⁶

176 The acquisition function defines the method by which new molecules are picked at the
177 start of each active learning cycle, with the simplest example being a ‘greedy’ approach,
178 which directly selects the (currently predicted) highest scoring molecules. However, an ac-
179 quisition function has to balance picking the best compounds, with the need to further refine
180 the accuracy of the machine learning model. Picking the best scoring candidates in descend-
181 ing order might initially increase the objective function, but the algorithm will have the
182 propensity to get stuck in local maxima and to be sensitive to the initial selection of training
183 molecules.

184 There are a variety of alternatives that aim to avoid the problems of a simple greedy
185 approach, and the approach used here is the upper confidence bound (UCB) uncertainty-
186 based acquisition function.⁴⁷ UCB considers not just the value of the objective function,
187 but also the variance of the prediction (model uncertainty), effectively biasing towards the
188 selection of molecules about which the model is the least certain of the predicted score. The

189 UCB function is defined by:

$$UCB(x) = \mu(x) + \beta\sigma(x), \quad (2)$$

190 where $\mu(x)$ and $\sigma(x)$ are the mean and standard deviation of the regressor for molecule x ,
191 and β is a parameter controlling the degree of exploration (high β increases the chances
192 that a molecule with moderate score but high uncertainty will be picked). The effects of the
193 choice of machine learning model and acquisition function, as well as other active learning
194 hyperparameters, are discussed later.

195 Database Search

196 A challenge for automated growing of linkers and R-groups, and for de novo design in general,
197 is the synthetic tractability of the designed compounds. Approaches to address this limitation
198 could include a synthetic accessibility score in the objective function⁴⁸ or the expert curation
199 of libraries with known synthetic routes.³² However, we wished to fully automate the design
200 process, and be confident of acquiring compounds for rapid design-make-test-analyse cycles.
201 We therefore make use of the rapidly-growing make-on-demand compound libraries as a
202 surrogate measure of synthetic accessibility. Ideally, we might use the entire catalogue as
203 a chemical space in which to perform the active learning. Although such an approach has
204 been used as a one-off screen,⁴⁹ evaluating the regression models used here soon becomes
205 prohibitively expensive in an active learning cycle. On the other hand, highly efficient
206 methods have been developed for similarity and substructure searches of these libraries.²²
207 We therefore make use of these searches to seed the chemical space with compounds that
208 are similar to the predicted actives at each step of the active learning cycle (Figure 1(C)). In
209 this way, at the subsequent acquisition step, we enable the algorithm to pick compounds for
210 growing and scoring that are likely to be scored highly (due to similarity with other highly
211 scoring compounds) and available for purchase or synthesis (due to presence in on-demand
212 libraries).

213 In detail, the Enamine REAL database of 4.5 B compounds was searched for similarity to
214 designed molecules through the public interface to SmallWorld <https://sw.docking.org>,
215 using a graph-edit-distance space search.²² At each cycle, 100 new, top-scoring compounds
216 were searched, and up to 100 of the most similar compounds from the REAL database were
217 extracted per search query (using a maximum distance of 5 steps). This 10 K compound
218 set was filtered for substructure match with the core using RDKit,³⁸ and those compounds
219 that passed were added to the active learning search space. Active learning then selects
220 compounds for scoring following Enamine enrichment, as usual, but there is no explicit bias
221 to select compounds from the on-demand catalogue.

222 Computational Details

223 Protein input structures were taken from the set of noncovalent complexes crystallised early
224 during the COVID-19 pandemic.³ In particular, the input PDB: 5R83 was used as the
225 receptor structure for active learning design, and Chimera was used to add hydrogen atoms.⁵⁰
226 The ligand was truncated to include only the pyridyl moiety, as this appeared in other
227 available crystallised fragments in a consistent binding mode (PDB: 5RE4, 5REH, 5R84,
228 5RF3³) and with a suitable vector for growth into the binding pocket. The full set of 23
229 non-covalent complexes (that had ligands bound in areas of the pocket accessible by a growth
230 vector) was additionally used for construction of the reference PLIP³⁴ interactions.

231 For testing of the active learning protocols, the chemical space was assembled by combin-
232 ing the pyridyl moiety with 508 R-groups⁵¹ and 100 of the most common linkers¹⁶ from the
233 FEgrow library. A total of 47710 unique molecules were successfully grown into the bind-
234 ing pocket and scored using the gina CNN scoring function.¹⁹ A further 1656 molecules
235 were assigned a penalty score of $pK = 0$ as they could not be embedded due to steric clash
236 with the protein. In cases where rare errors occurred, such as a failure to assign force field
237 parameters, the molecules were discarded completely.

238 The previously tested FEgrow molecule building protocol was applied throughout.¹⁵ The

239 ETKDG algorithm³⁹ was used to generate 50 conformers, using a 0.5 Å root-mean-square
240 similarity threshold. Any conformers with an atom closer than 1 Å to any atom in the
241 protein was discarded. Energy minimisation was applied using a hybrid machine learning
242 / molecular mechanics energy function in a mechanical embedding scheme.¹⁵ The ANI-2x
243 potential¹⁷ was used for the ligand, in cases where all elements in the molecule are covered by
244 the model, or the Open Force Field Sage⁴¹ potential otherwise. The lowest energy conformer
245 was retained for scoring.

246 An active learning library based on scikit⁵² and modAL⁵³ python packages was adopted
247 from another study.²⁶ A set of molecules to initialise the active learning cycle can be se-
248 lected via RDKit's MaxMin picker³⁸ from the chemical space, or picked at random. The
249 processing was parallelised using the python library Dask,⁵⁴ which supports a diverse set of
250 technologies, including the Slurm Workload Manager that is deployed ubiquitously on high-
251 performance computing clusters. Dask is used to secure resources (scheduling workers on
252 Slurm), submitting work and retrieving results. The three major computationally-expensive
253 components were parallelised: 1) building and scoring of the molecules, 2) computing the
254 Morgan fingerprints, and 3) computing the Tanimoto similarity across the chemical space
255 for the Gaussian Process modelling.

256 Results

257 **Interfacing FEgrow with active learning enables efficient search of** 258 **chemical space.**

259 In order to investigate the performance of the active learning protocol, and the effect of
260 machine learning hyperparameters, we built a labelled 'oracle' set of 47 K compounds using
261 standard FEgrow input settings (see Computational Details). This is a larger set of com-
262 pounds than would be typically built and scored against a target, but knowing the affinities
263 of the full chemical space enables us to assess the performance of the active learning ap-

264 proach. The common core was selected to be a pyridyl fragment common to several early
265 crystal structures of the SARS-CoV-2 main protease,³ located in the S1 pocket with a vector
266 pointing into the enzyme active site (Figure 2(a)).

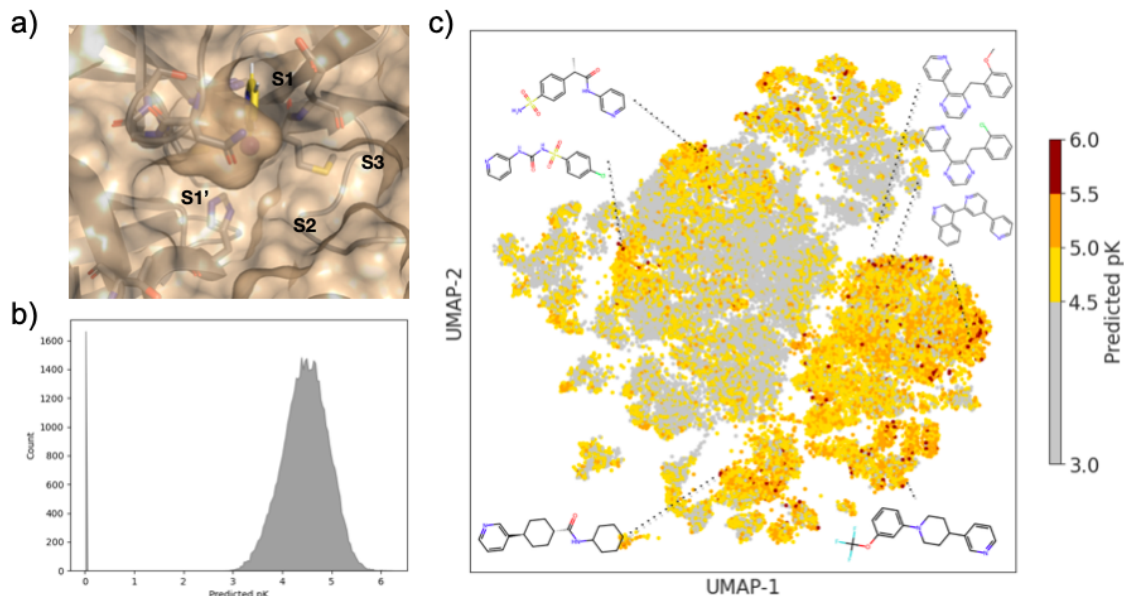


Figure 2: a) The position of the ligand core and definitions of binding pocket labels, the purple sphere is the hydrogen atom for replacement. b) Histogram of computed pK for the 47 K compound oracle dataset. c) UMAP of entire 47 K oracle chemical space, coloured by computed pK. 2D structures of representative strong binders are included.

267 Figure 2(b) shows the distribution of predicted binding affinities, computed using the
268 gnina convolutional neural network scoring function¹⁹ from FEgrow built structures. The
269 scores are symmetrically distributed around pK = 4.5, with a maximum affinity of around
270 6.0, which is indicative of a set of low molecular weight (range between 100 and 350 Da,
271 **Figure S1**), unoptimised compounds at the start of a hit finding effort. Indeed, it is at this
272 stage where the options for expansion are vast, and strategies to suggest exploration of hits
273 are particularly valuable. Note that compounds that could not be built (for example, due to
274 steric clashes with the protein) are arbitrarily assigned a pK of zero, so that this information
275 can be included in the active learning model.

276 Figure 2(c) further shows the UMAP projection of the chemical space, coloured by pre-
277 dicted pK. The visualisation shows a diverse composition of linkers and functional groups,

278 with well-spread clusters of the highest affinity binders, potentially providing a challenging
279 search space for active learning. Figure 2(b) also shows locations in the chemical space of ex-
280 ample linker and R-groups, attached to the pyridyl core, that make up the stronger predicted
281 binders. Favourable predicted linkers include amides, sulfonylurea and various 6-membered
282 ring heterocycles, and relatively bulky R-groups are feasible, which is generally expected
283 given the size and shape of the binding pocket.^{3,4} (Note that at this stage no consideration
284 is given to synthetic accessibility or stability of the compound designs).

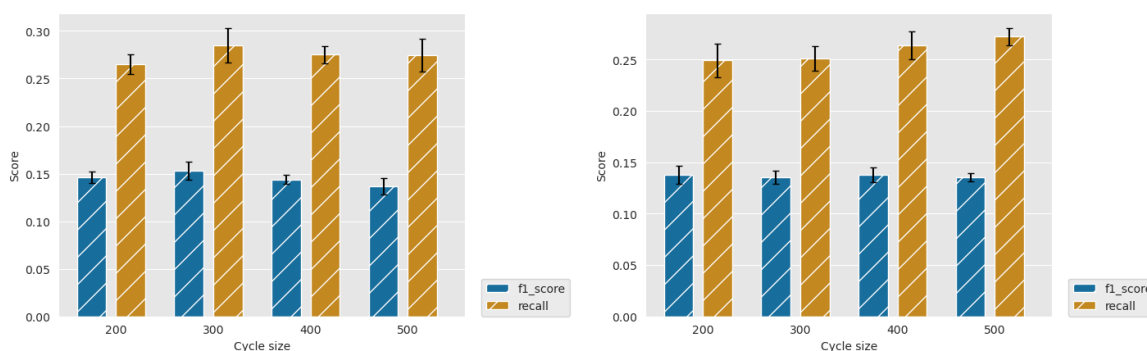


Figure 3: Recall and F1 score for diverse initial selection GBM (left) and GP (right) models, and greedy acquisition for identification of top 2 % scoring compounds for different cycle sizes. Error bars show standard errors over five runs.

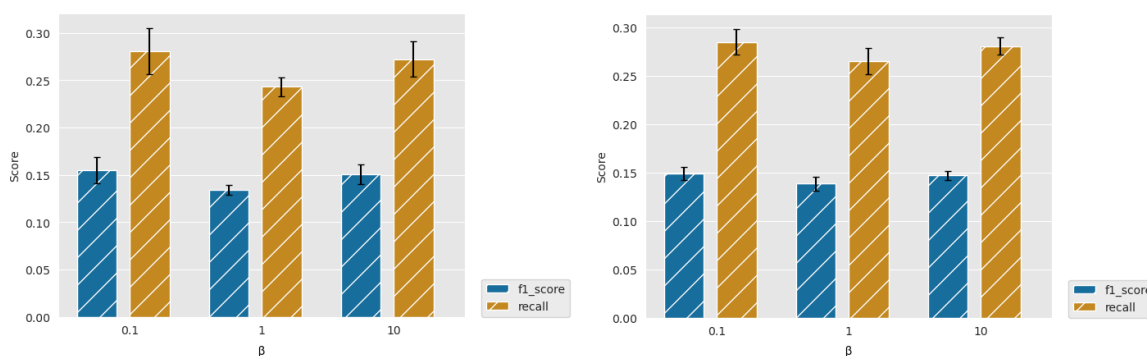


Figure 4: Recall and F1 score for diverse initial selection using GP and UCB acquisition (and varying β) with cycle sizes of 200 (left) and 400 (right) for identification of top 2 % scoring compounds. Error bars show standard errors over five runs.

285 We next sought to use active learning to accelerate the search through this chemical
286 space, using the oracle to assess the performance of model hyperparameters, and using the
287 predicted binding affinity as the optimisation target. In particular, we have investigated the

288 effects of initial compound selection (random or diverse), number of compounds picked per
289 cycle (in the range 200–500), machine learning model (GBM or GP) and acquisition method
290 (greedy or UCB). As discussed, the dependence of active learning efficiency on the choice of
291 model parameters is well documented, and so we do not devote much space to it here.

292 By way of example, Figure 3 shows the effect of the number of compounds picked per
293 cycle on model recall and precision (F1 score) for the two machine learning models (GBM
294 and GP). For a fixed total number of compounds selected (here, 2500), one might expect the
295 model to improve at small sample sizes (hence, more active learning cycles), but we find that
296 the efficiency is already well converged when picking 500 per cycle. Similarly, the choice of
297 machine learning model has little effect, with slightly higher metrics for the GBM model, but
298 both recall and precision comparisons are within the error bars. Figure 4 further shows the
299 effect of using the UCB uncertainty-based acquisition function, instead of greedy selection,
300 in conjunction with the GP machine learning model. There is some small improvement in
301 recall over greedy selection, but no significant change in the metrics used either as a function
302 of cycle size or the β parameter in eq 2.

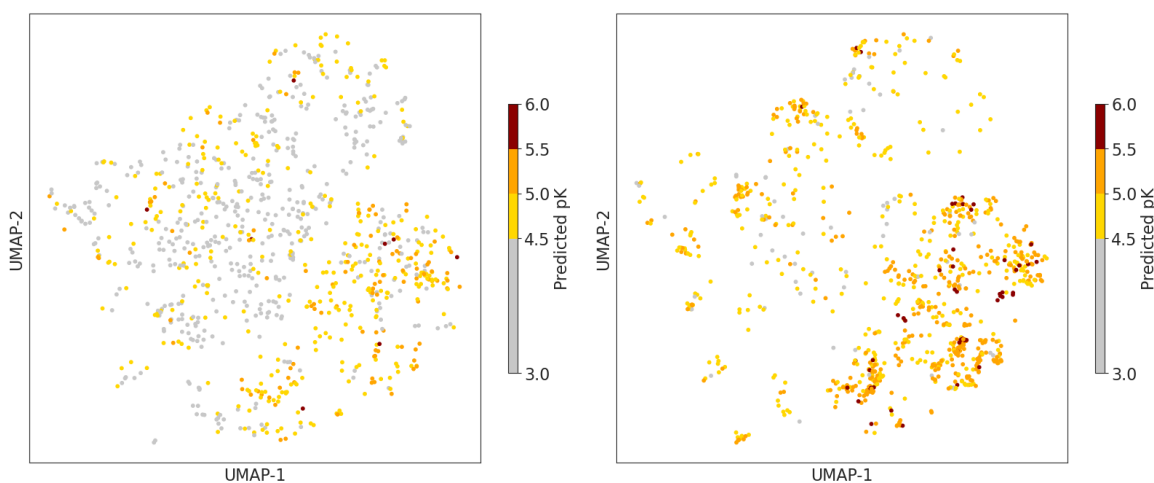


Figure 5: Difference in selection for first (left) and final (right) active learning cycles, showing a narrowing into areas predicted to be potent and avoiding unpromising areas.

303 Note that for the current dataset, random selection would give a recall of 0.05 and F1
304 score of 0.03 for identification of the top 2% of compounds. Therefore, with recall of around

305 0.25–0.30 for most of our experiments, we see efficiency improvements with active learning of
306 around a factor of 5x compared to random selection. For reference, the growth and scoring
307 of this compound set in FEgrow requires around 1000 cpuhrs, which is not prohibitive, but
308 automated acceleration at no cost is clearly worthwhile.

309 In the next section, we choose to use a GP model with UCB acquisition function, with
310 a cycle size of 200 and a diverse set of starting compounds. The overall accuracy of the
311 chosen regression model (using $\beta = 10$), following training on 5 % of the dataset, is 0.97 pK
312 units (**Figure S2**), which is competitive with typical models used in active learning with
313 fingerprint-based representations.²⁷ Figure 5 shows a similar UMAP projection as in Figure 2,
314 but now only showing compounds acquired by our chosen active learning model in the first
315 (left) and final (right) cycles. We observe both a wide exploration of the chemical space,
316 which is important to increase diversity in the final set, and a focusing of the explored regions
317 in the final cycle to compounds with a higher predicted binding affinity, which is important
318 for the use of the model to identify strong binders.

319 **Active learning driven fragment expansion identifies potential SARS-** 320 **CoV-2 MPro inhibitors.**

321 Having established that the active learning protocols tested here are able to improve the
322 efficiency of chemical space searches with FEgrow, we turn now to prospective design of
323 potential noncovalent SARS-CoV-2 MPro inhibitors. A wealth of computational and experi-
324 mental data has been generated for this target in recent years, but here we limit ourselves to
325 structural information that was available in the early months of the COVID-19 pandemic. In
326 particular, as in the previous section, we consider expansion of the pyridyl fragment (PDB:
327 5R83) along a vector into the binding pocket containing the catalytic cysteine (Cys145).³
328 We now expand the size of the chemical space to an initial 250,000 molecules, built from
329 the combination of supplied libraries of 500 linkers and 500 R-groups, such that full building
330 and scoring of the space is prohibitively expensive for routine study. To address the issue of

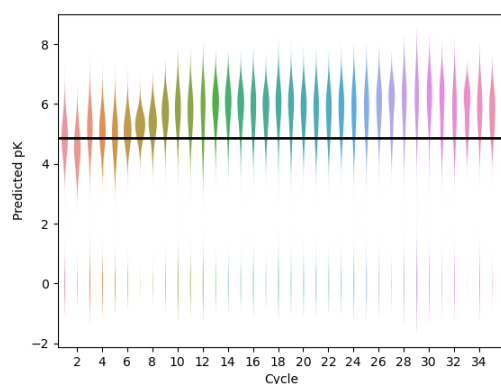


Figure 6: Active learning drives improvements in predicted binding affinity. A GP model is used, with UCB acquisition function ($\beta = 0.1$), a cycle size of 200 and a diverse set of starting compounds. The solid horizontal line shows the average score for 377 compounds randomly selected from the REAL database that were built with FEgrow.

331 synthetic feasibility of the output designs, we add an additional step in the active learning
 332 cycle (Figure 1C), whereby the chemical space is periodically seeded with compounds from
 333 the REAL database that are similar to the highest scoring compounds (see Methods). **Fig-**
 334 **ure S3** demonstrates successful incorporation of the Enamine compounds into the active
 335 learning cycles, with a significant fraction of the built and scored compounds originating
 336 from this source.

337 Figure 6 shows an example design run, optimising the compounds for predicted pK using
 338 the gnina scoring function (further examples are given in the **Supporting Information**).
 339 The distribution of predicted affinity increases over the first 10 active learning cycles then
 340 starts to saturate with a mean predicted pK close to 6 (micromolar affinity). Over the full
 341 run, 95% of the compounds were successfully built (assigned $pK > 0$) and 15% had a pre-
 342 dicted $pK > 6$. For comparison, we also extracted 1000 molecules at random that contained
 343 the pyridyl substructure from the REAL database used to seed the active learning cycles.
 344 For this set, 377 molecules (38%) could be successfully built, with an average predicted
 345 $pK = 4.9$ and only two compounds with predicted $pK > 6.0$ (0.2%).

346 Figure 7a) shows the highest scoring compound from this run, with a predicted affinity
 347 of 88 nM. The compound extends hydrophobic contacts into the S3 and S1' pockets, for ex-

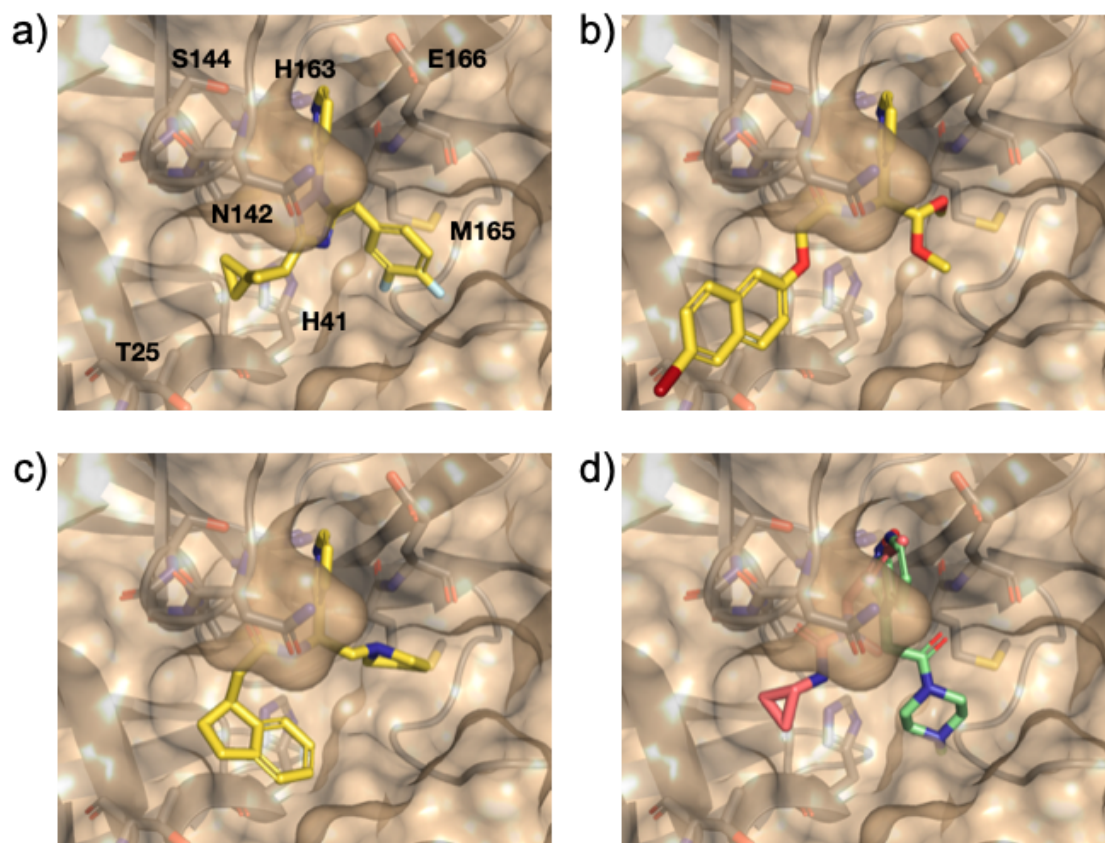


Figure 7: a) Top-scoring compounds optimised for a) predicted pK, b) protein-ligand interaction profile and c) combined scoring function. d) Fragment 5RGI shown in pink (H-bond donation by Gly143, Ser144, Cys145 and His163), and 5RF7 in green (hydrophobic and H-bond donation with Glu166).

348 ample with Met165 and Thr25, but despite this does not form any specific polar interactions
349 (other than the original core interaction with His163). Since an early fragment screen had
350 provided valuable information about the nature of potential protein–ligand interactions in
351 this binding pocket, we sought to reduce the reliance on the gnina scoring function and drive
352 the active learning towards compounds that recovered known crystallographic information
353 (see Methods). Figure 7b) shows the top-scoring compound, as defined by the Tanimoto
354 similarity to the vector of reference interactions. In this case, the grown molecule forms
355 additional hydrogen bonding interactions with Asn142, Gly143, Ser144, Cys145 and Glu166,
356 and hydrophobic interactions with Thr25 and Glu166. The majority of these interactions
357 are recapitulated by, for example, fragments PDB: 5RGI and 5RF7 (Figure 7d)).

358 Finally, we sought to combine the strengths of both docking scores and crystallographic
359 information to optimise a combined scoring function. Figure 7c) shows the top-scoring
360 compound as defined by eq 1 after 33 cycles of active learning. Although this compound is
361 scored much lower by the gnina scoring function (predicted affinity 2 μ M), it extends into
362 the S3 and S1' pockets and retains many of the interactions observed in Figure 7b) (e.g.
363 hydrogen bonding interactions with Asn142, Gly143, Ser144, Cys145 and Glu166).

364 **Analysis of hit compounds.**

365 The top 500 compounds from each of four active learning runs (two optimising predicted
366 pK, one optimising protein-ligand interactions, and one optimising the combined scoring
367 function) were checked for availability from the Enamine store. Interestingly, very few of the
368 top scored by predicted pK were available (four in total). This is likely due to an important
369 unavailable building block(s), and could be mitigated in future by increasing diversity and/or
370 including direct store queries in the search process. In any case, we focussed here on outputs
371 from the remaining two runs, and submitted the top 10 protein-ligand interaction and top
372 25 combination scoring compounds for costing. Finally, a total of 19 designed compounds
373 were purchased (of which 15 had been optimised used the combination score) based on

374 quoted price and excluding similar compounds (based on visual inspection). Two control
375 compounds were also included; one known binder from a crystallographic fragment screen
376 (Enamine ID: Z44592329; PDB: 5R83)³ and one elaborated compound from the COVID
377 Moonshot study (Enamine ID: Z4943052515 (literature IC₅₀ 0.288 μM)).⁴ The twenty one
378 purchased compounds (**Figure S11**) were evaluated in a fluorescence-based Mpro activity
379 assay at 1000, 500, 10 μM (**Figure S12**). Compounds **5** and **6** were excluded from the study
380 due to solubility issues at 1000 μM in the assay conditions. Five compounds (**8**, **10**, **12**, **14**
381 and **21** (the positive control⁴)) showed reduction of Mpro activity ≤ 50% at 1000 μM. The
382 IC₅₀ values of these compounds, except **8** which displayed background autofluorescence, were
383 further determined (Figure 8). Compounds **10**, **12** and **14** showed a concentration-dependent
384 inhibition of Mpro activity (measured pIC₅₀ 2.10, 3.01, 2.80 respectively). Nirmatrelvir, an
385 orally bioavailable antiviral drug targeting Mpro, showed inhibition (pIC₅₀ 6.01), which was
386 slightly higher than the reported IC₅₀ (0.022 μM⁵⁵), likely due to the limit of the assay (the
387 enzyme concentration was at 0.2 μM). Figure 9 shows the predicted structures of compounds
388 **12** and **14** from the active learning design runs. Both compounds form hydrogen bonding
389 interactions with the backbone of Glu166, as well as hydrophobic interactions in the S1'
390 pocket.

391 Finally, to investigate whether the relatively low affinity of designed compounds is due to
392 insufficient exploration of chemical space or the empirical objective functions used to optimise
393 molecules, we performed a retrospective analysis of the designed compound space against
394 known binders resulting from the COVID Moonshot crowd-sourced discovery campaign.⁴ In
395 particular, Figure 10 shows the three most similar compounds from the active learning runs
396 (as defined by Tanimoto similarity search between RDKit Morgan fingerprints with a radius
397 of 3 and size of 2048) to a curated set of 292 hit compounds. Considering that our FEGrow
398 runs took as input only a single PDB receptor structure and pyridyl fragment core, it is
399 clear that this fragment growing and on-demand library screening approach holds promise
400 for suggesting biologically active compounds early in hit discovery campaigns. However,

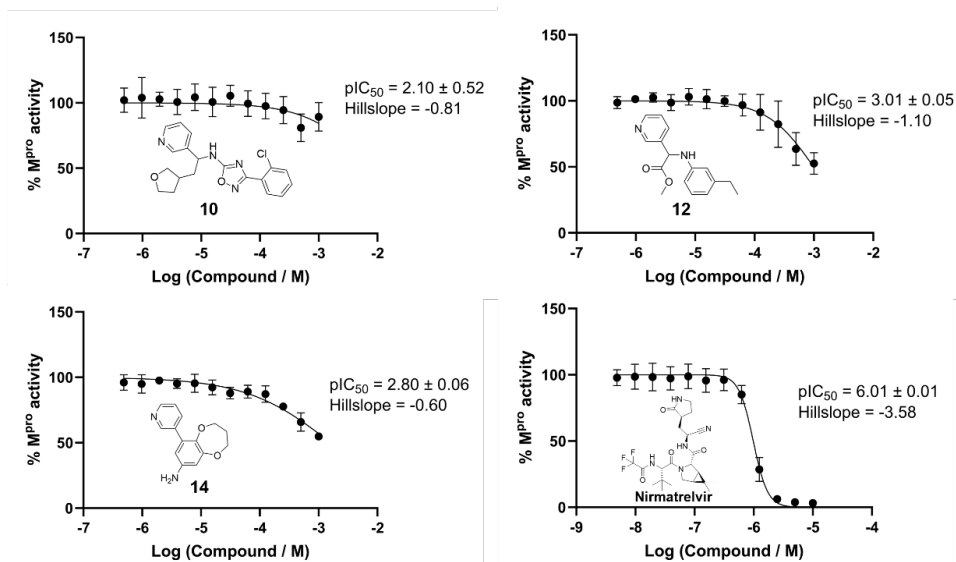


Figure 8: IC₅₀ determination of selected compounds with Mpro. Compounds **10**, **12** and **14** were tested at a top concentration of 1000 μ M. Nirmatrelvir was tested at a top concentration of 10 μ M as a positive control. Datapoints presented as mean \pm SD; pIC₅₀ presented as mean \pm SEM; two biological repeats consisting of three technical replicates. **10** consists of one biological repeat with three technical replicates. Conditions: Mpro (0.2 μ M), 12-hour pre-incubation with compounds, 20 μ M fluorescent substrate, 50 mM Tris-HCl (pH 7.3), 1 mM EDTA and temperature 25°C.

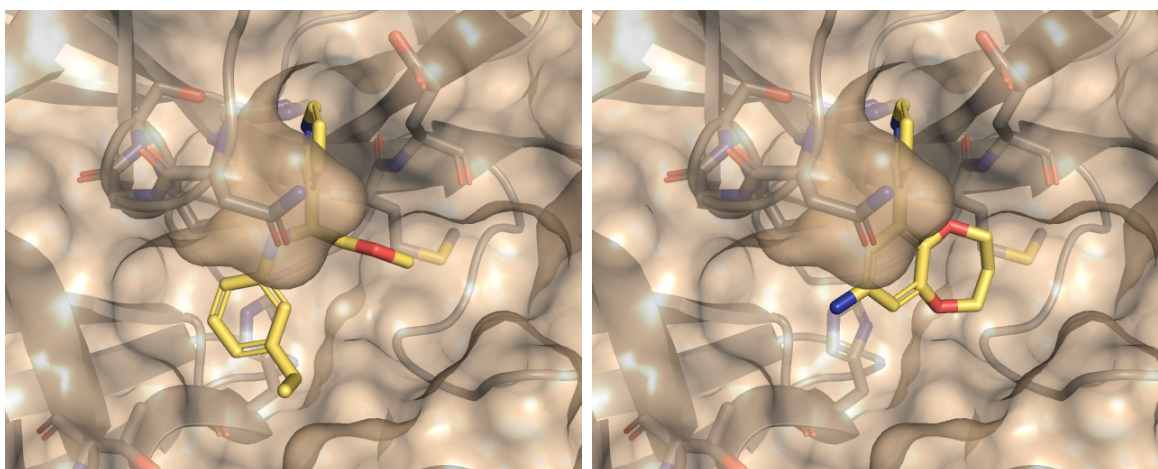


Figure 9: Predicted bound structures of compounds **12** (Z1470573089) and **14** (Z8969017446).

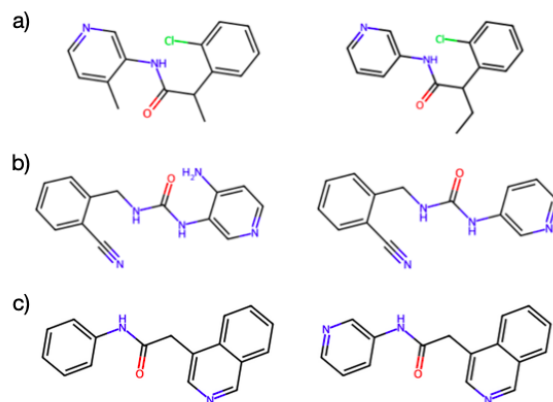


Figure 10: a) Experimental Moonshot compound (literature IC_{50} 17 μM) and most similar compound from this study, from active learning optimisation of predicted pK ($\beta=10$), b) Experimental Moonshot compound (literature IC_{50} 54 μM) and most similar compound from this study, from active learning optimisation of predicted pK ($\beta=10$), c) Experimental Moonshot compound (literature IC_{50} 57 μM) and most similar compound from this study, from active learning optimisation of combination scoring function.

401 further work is needed to ensure that the most promising compounds are located at the top
 402 of ranked lists for synthetic prioritisation and testing.

403 Discussion and Conclusions

404 In this study, we have combined the FEGrow software, an open modular workflow for building
 405 and scoring ligands in protein binding pockets, with active learning to guide and automate
 406 chemical space searches for promising binders. In agreement with numerous other studies,²⁷
 407 we have shown that search efficiency is not too dependent on the hyperparameters of the
 408 active learning model, which include the choice of regression model, the acquisition function
 409 and number of compounds picked per cycle. For this particular study, we find efficiency
 410 improvements of a factor of around 5x over random selection, which will aid throughput of
 411 future prospective design efforts.

412 With the design of FEGrow, we hope to overcome some of the current limitations of de
 413 novo drug design discussed in the Introduction. Some of these limitations are addressed in
 414 the current study, and some will be addressed in future aided by ongoing advances in molec-

415 ular modelling and machine learning. For example, we tackle the question of binding pose
416 optimisation by using a fast and accurate machine learning potential (ANI-2x¹⁷) to describe
417 the ligand energetics in a mechanical embedding scheme. However, with the flexibility of
418 the FEgrow interface with OpenMM,¹⁸ new models could be substituted in, and these are
419 now approaching sufficient speed and accuracy (including for long-ranged interactions) such
420 that the entire protein-ligand complex could be described using a single, consistent machine
421 learning potential.^{56,57} In this study, we made the approximation that the protein binding
422 pocket is rigid and used a single receptor structure for design. However, now that ligand
423 building and scoring is fully automated, future studies could use, for example, ensembles of
424 receptor structures, which may be beneficial in cases where the pocket is more flexible.

425 A limitation of this and other similar studies is the choice of objective function in the ac-
426 tive learning cycles. To demonstrate the flexibility of the FEgrow package, we demonstrated
427 four design cycles here, two optimising for predicted affinity using the gnina CNN scoring
428 function and two including a more direct optimisation of protein-ligand contacts extracted
429 from crystallographic fragment screens. While we do not have enough data to assess the
430 relative merits of these scoring functions, we expect the latter to be useful where experimen-
431 tal structural data exists, at least as part of a multi-objective optimisation in future.⁵⁸ As
432 a flexible alternative to PLIP scores trained on system-dependent crystal structures, it has
433 also been shown that transferable neural networks can be trained on the PDBbind structural
434 database to recognise favourable protein-ligand interactions.⁵⁹

435 As shown in Figure 1c), to address the issue of synthetic tractability of the de novo
436 built compounds, we inserted regular queries of the Enamine REAL database into the active
437 learning cycles. In this way, we can use the initial chemical space to train the active learning
438 regression models, and then over time seed the chemical space with compounds that are both
439 similar to predicted actives and purchasable. In this way, we were able to test the predictions
440 of the active learning workflow with a turn around time of a few weeks from order to biological
441 testing. Of the 19 designed compounds that were purchased here, three showed measurable

442 activity, but none approached the desired levels for further progression. Nevertheless, a
443 similarity search showed the presence of effective inhibitors in the built chemical space, and
444 so further investigation will focus on ranking compound designs ahead of purchase, perhaps
445 via an extra stage of physics-based free energy calculations.²⁶

446 Code Availability

447 FEgrow is freely available, with a set of tutorials, at [https://github.com/cole-group/](https://github.com/cole-group/FEgrow)
448 FEgrow.

449 Author Contributions

450 **Ben Cree:** Conceptualisation, Data curation, Formal Analysis, Investigation, Methodology,
451 Software, Validation, Visualisation, Writing - original draft.

452 **Mateusz Bieniek:** Conceptualisation, Data curation, Formal Analysis, Investigation, Method-
453 ology, Software, Validation, Visualisation, Writing - original draft.

454 **Siddique Amin:** Data curation, Formal Analysis, Investigation, Methodology, Writing -
455 original draft.

456 **Akane Kawamura:** Resources, Supervision, Writing – review & editing.

457 **Daniel Cole:** Conceptualisation, Funding acquisition, Methodology, Project administra-
458 tion, Resources, Supervision, Writing – original draft, Writing – review & editing.

459 Competing Interests Disclosure

460 The authors declare no competing interests.

461 Acknowledgement

462 D.J.C. and M.K.B. acknowledge support from a UKRI Future Leaders Fellowship (grant
463 MR/T019654/1). B.C. and S.A are grateful for support from the EPSRC Centre for Doctoral
464 Training in Molecular Sciences for Medicine (grant EP/S022791/1). This research made use
465 of the Rocket High Performance Computing service at Newcastle University. We are grateful
466 to Dr Mathew Martin, Dr Julia Hubbard, and Dr Sara Pintar for provision of the Mpro
467 expression construct and to Prof. Christopher Schofield for the provision of Nirmatrelvir.

468 Supporting Information Available

469 References

- 470 (1) Douangamath, A.; Powell, A.; Fearon, D.; Collins, P. M.; Talon, R.; Kro-
471 jer, R., T.and Skyner; Brandao-Neto, J.; Dunnett, L.; Dias, A.; Aimon, A.;
472 Pearce, N. M.; Wild, C.; Gorrie-Stone, T.; von Delft, F. Achieving Efficient Fragment
473 Screening at XChem Facility at Diamond Light Source. *J. Vis. Exp.* **2021**, *171*, e62414.
- 474 (2) Wood, D. J.; Lopez-Fernandez, J. D.; Knight, L. E.; Al-Khawaldeh, I.; Gai, C.; Lin, S.;
475 Martin, M. P.; Miller, D. C.; Cano, C.; Endicott, J. A.; Hardcastle, I. R.; Noble, M.
476 E. M.; Waring, M. J. FragLites—Minimal, Halogenated Fragments Displaying Phar-
477 macophore Doublets. An Efficient Approach to Druggability Assessment and Hit Gen-
478 eration. *Journal of Medicinal Chemistry* **2019**, *62*, 3741–3752.
- 479 (3) Douangamath, A. et al. Crystallographic and electrophilic fragment screening of the
480 SARS-CoV-2 main protease. *Nature Communications* **2020**, *11*, 5047.
- 481 (4) Bobby, M. L. et al. Open science discovery of potent noncovalent SARS-CoV-2 main
482 protease inhibitors. *Science* **2023**, *382*, eabo7201.

- 483 (5) Mackinnon, S. R.; Krojer, T.; Foster, W. R.; Diaz-Saez, L.; Tang, M.; Huber, K. V. M.;
484 von Delft, F.; Lai, K.; Brennan, P. E.; Arruda Bezerra, G.; Yue, W. W. Fragment
485 Screening Reveals Starting Points for Rational Design of Galactokinase 1 Inhibitors to
486 Treat Classic Galactosemia. *ACS Chemical Biology* **2021**, *16*, 586–595.
- 487 (6) Grygorenko, O. O.; Radchenko, D. S.; Dziuba, I.; Chuprina, A.; Gubina, K. E.; Mo-
488 roz, Y. S. Generating Multibillion Chemical Space of Readily Accessible Screening
489 Compounds. *iScience* **2020**, *23*, 101681.
- 490 (7) Warr, W. A.; Nicklaus, M. C.; Nicolaou, C. A.; Rarey, M. Exploration of Ultralarge
491 Compound Collections for Drug Discovery. *Journal of Chemical Information and Mod-
492 eling* **2022**, *62*, 2021–2034.
- 493 (8) Kuan, J.; Radaeva, M.; Avenido, A.; Cherkasov, A.; Gentile, F. Keeping pace with the
494 explosive growth of chemical libraries with structure-based virtual screening. *WIREs
495 Computational Molecular Science* **2023**, *13*, e1678.
- 496 (9) Green, H.; Koes, D. R.; Durrant, J. D. DeepFrag: a deep convolutional neural network
497 for fragment-based lead optimization. *Chem. Sci.* **2021**, *12*, 8036–8047.
- 498 (10) Imrie, F.; Hadfield, T. E.; Bradley, A. R.; Deane, C. M. Deep generative design with
499 3D pharmacophoric constraints. *Chem. Sci.* **2021**, *12*, 14577–14589.
- 500 (11) Runcie, N. T.; Mey, A. S. SILVR: Guided Diffusion for Molecule Generation. *Journal
501 of Chemical Information and Modeling* **2023**, *63*, 5996–6005.
- 502 (12) Sadybekov, A. A. et al. Synthron-based ligand discovery in virtual libraries of over 11
503 billion compounds. *Nature* **2022**, *601*, 452–459.
- 504 (13) Gahbauer, S. et al. Iterative computational design and crystallographic screening iden-
505 tifies potent inhibitors targeting the Nsp3 macrodomain of SARS-CoV-2. *Proceedings
506 of the National Academy of Sciences* **2023**, *120*, e2212931120.

- 507 (14) Ferla, M.; Sánchez-García, R.; Skyner, R.; Gahbauer, S.; Taylor, J.; von Delft, F.; Mars-
508 den, B.; Deane, C. Fragmenstein: predicting protein-ligand structures of compounds de-
509 rived from known crystallographic fragment hits using a strict conserved-binding-based
510 methodology. *ChemRxiv* **2024**,
- 511 (15) Bieniek, M.; Cree, B.; Pirie, R.; Horton, J.; Tatum, N.; Cole, D. An open-source
512 molecular builder and free energy preparation workflow. *Commun. Chem.* **2022**, *5*,
513 136.
- 514 (16) Ertl, P.; Altmann, E.; Racine, S. The most common linkers in bioactive molecules and
515 their bioisosteric replacement network. *Bioorganic & Medicinal Chemistry* **2023**, *81*,
516 117194.
- 517 (17) Devereux, C.; Smith, J. S.; Huddleston, K. K.; Barros, K.; Zubatyuk, R.; Isayev, O.;
518 Roitberg, A. E. Extending the Applicability of the ANI Deep Learning Molecular Po-
519 tential to Sulfur and Halogens. *J. Chem. Theory Comput.* **2020**, *16*, 4192–4202.
- 520 (18) Eastman, P. et al. OpenMM 8: Molecular Dynamics Simulation with Machine Learning
521 Potentials. *The Journal of Physical Chemistry B* **2024**, *128*, 109–116.
- 522 (19) McNutt, A. T.; Francoeur, P.; Aggarwal, R.; Masuda, T.; Meli, R.; Ragoza, M.; Sun-
523 seri, J.; Koes, D. R. GNINA 1.0: molecular docking with deep learning. *J. Cheminf.*
524 **2021**, *13*, 1–20.
- 525 (20) Zhang, C.-H. et al. Potent Noncovalent Inhibitors of the Main Protease of SARS-CoV-2
526 from Molecular Sculpting of the Drug Perampanel Guided by Free Energy Perturbation
527 Calculations. *ACS Cent. Sci.* **2021**, *7*, 467–475.
- 528 (21) Fialková, V.; Zhao, J.; Papadopoulos, K.; Engkvist, O.; Bjerrum, E. J.; Kogej, T.;
529 Patronov, A. LibINVENT: Reaction-based Generative Scaffold Decoration for in Silico
530 Library Design. *Journal of Chemical Information and Modeling* **2022**, *62*, 2046–2063.

- 531 (22) Irwin, J. J.; Tang, K. G.; Young, J.; Dandarchuluun, C.; Wong, B. R.; Khurelbaatar, M.;
532 Moroz, Y. S.; Mayfield, J.; Sayle, R. A. ZINC20—A Free Ultralarge-Scale Chemical
533 Database for Ligand Discovery. *Journal of Chemical Information and Modeling* **2020**,
534 *60*, 6065–6073.
- 535 (23) Yu, J.; Li, X.; Zheng, M. Current status of active learning for drug discovery. *Artificial*
536 *Intelligence in the Life Sciences* **2021**, *1*, 100023.
- 537 (24) Graff, D. E.; Shakhnovich, E. I.; Coley, C. W. Accelerating high-throughput virtual
538 screening through molecular pool-based active learning. *Chem. Sci.* **2021**, *12*, 7866–
539 7881.
- 540 (25) Khalak, Y.; Tresadern, G.; Hahn, D. F.; de Groot, B. L.; Gapsys, V. Chemical Space
541 Exploration with Active Learning and Alchemical Free Energies. *Journal of Chemical*
542 *Theory and Computation* **2022**, *18*, 6259–6270.
- 543 (26) Thompson, J.; Walters, W. P.; Feng, J. A.; Pabon, N. A.; Xu, H.; Maser, M.; Gold-
544 man, B. B.; Moustakas, D.; Schmidt, M.; York, F. Optimizing active learning for free
545 energy calculations. *Artificial Intelligence in the Life Sciences* **2022**, *2*, 100050.
- 546 (27) Gorantla, R.; Kubincová, A.; Suutari, B.; Cossins, B. P.; Mey, A. S. J. S. Benchmarking
547 Active Learning Protocols for Ligand-Binding Affinity Prediction. *Journal of Chemical*
548 *Information and Modeling* **2024**, *64*, 1955–1965.
- 549 (28) van Tilborg, D.; Grisoni, F. Traversing Chemical Space with Active Deep Learning: A
550 Computational Framework for Low-data Drug Discovery. *ChemRxiv* **2024**,
- 551 (29) Gentile, F.; Agrawal, V.; Hsing, M.; Ton, A.-T.; Ban, F.; Norinder, U.; Gleave, M. E.;
552 Cherkasov, A. Deep Docking: A Deep Learning Platform for Augmentation of Structure
553 Based Drug Discovery. *ACS Central Science* **2020**, *6*, 939–949.

- 554 (30) Yang, Y.; Yao, K.; Repasky, M. P.; Leswing, K.; Abel, R.; Shoichet, B. K.; Jerome, S. V.
555 Efficient Exploration of Chemical Space with Docking and Deep Learning. *Journal of*
556 *Chemical Theory and Computation* **2021**, *17*, 7106–7119.
- 557 (31) Sivula, T.; Yetukuri, L.; Kalliokoski, T.; Käsnänen, H.; Poso, A.; Pöhner, I. Machine
558 Learning-Boosted Docking Enables the Efficient Structure-Based Virtual Screening of
559 Giga-Scale Enumerated Chemical Libraries. *Journal of Chemical Information and Mod-*
560 *eling* **2023**, *63*, 5773–5783.
- 561 (32) Konze, K. D.; Bos, P. H.; Dahlgren, M. K.; Leswing, K.; Tubert-Brohman, I.; Bor-
562 tolato, A.; Robbason, B.; Abel, R.; Bhat, S. Reaction-Based Enumeration, Active
563 Learning, and Free Energy Calculations To Rapidly Explore Synthetically Tractable
564 Chemical Space and Optimize Potency of Cyclin-Dependent Kinase 2 Inhibitors. *Jour-*
565 *nal of Chemical Information and Modeling* **2019**, *59*, 3782–3793.
- 566 (33) Gusev, F.; Gutkin, E.; Kurnikova, M. G.; Isayev, O. Active Learning Guided Drug
567 Design Lead Optimization Based on Relative Binding Free Energy Modeling. *Journal*
568 *of Chemical Information and Modeling* **2023**, *63*, 583–594.
- 569 (34) Adasme, M. F.; Linnemann, K. L.; Bolz, S. N.; Kaiser, F.; Salentin, S.; Haupt, V. J.;
570 Schroeder, M. PLIP 2021: expanding the scope of the protein–ligand interaction profiler
571 to DNA and RNA. *Nucleic Acids Research* **2021**, *49*, W530–W534.
- 572 (35) Glaab, E.; Manoharan, G. B.; Abankwa, D. Pharmacophore Model for SARS-CoV-2
573 3CLpro Small-Molecule Inhibitors and in Vitro Experimental Validation of Computa-
574 tionally Screened Inhibitors. *Journal of Chemical Information and Modeling* **2021**, *61*,
575 4082–4096.
- 576 (36) Hazemann, J.; Kimmerlin, T.; Lange, R.; Sweeney, A. M.; Bourquin, G.; Ritz, D.;
577 Czodrowski, P. Identification of SARS-CoV-2 Mpro inhibitors through deep reinforce-

- 578 ment learning for de novo drug design and computational chemistry approaches. *bioRxiv*
579 **2024**,
- 580 (37) Chenthamarakshan, V. et al. Accelerating drug target inhibitor discovery with a deep
581 generative foundation model. *Science Advances* **2023**, *9*, eadg7865.
- 582 (38) Landrum, G. RDKit: Open-source cheminformatics. <http://www.rdkit.org/>.
- 583 (39) Riniker, S.; Landrum, G. A. Better Informed Distance Geometry: Using What We
584 Know To Improve Conformation Generation. *J. Chem. Inf. Model.* **2015**, *55*, 2562–
585 2574.
- 586 (40) Maier, J. A.; Martinez, C.; Kasavajhala, K.; Wickstrom, L.; Hauser, K. E.; Simmer-
587 ling, C. ff14SB: Improving the Accuracy of Protein Side Chain and Backbone Parame-
588 ters from ff99SB. *J. Chem. Theory Comput.* **2015**, *11*, 3696–3713.
- 589 (41) Boothroyd, S. et al. Development and Benchmarking of Open Force Field 2.0.0: The
590 Sage Small Molecule Force Field. *Journal of Chemical Theory and Computation* **2023**,
591 *19*, 3251–3275.
- 592 (42) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and
593 testing of a general amber force field. *J. Comput. Chem.* **2004**, *25*, 1157–1174.
- 594 (43) Ragoza, M.; Hochuli, J.; Idrobo, E.; Sunseri, J.; Koes, D. R. Protein–Ligand Scoring
595 with Convolutional Neural Networks. *J. Chem. Inf. Model.* **2017**, *57*, 942–957.
- 596 (44) Francoeur, P. G.; Masuda, T.; Sunseri, J.; Jia, A.; Iovanisci, R. B.; Snyder, I.;
597 Koes, D. R. Three-Dimensional Convolutional Neural Networks and a Cross-Docked
598 Data Set for Structure-Based Drug Design. *J. Chem. Inf. Model.* **2020**, *60*, 4200–4215.
- 599 (45) Deringer, V. L.; Bartók, A. P.; Bernstein, N.; Wilkins, D. M.; Ceriotti, M.; Csányi, G.
600 Gaussian Process Regression for Materials and Molecules. *Chemical Reviews* **2021**, *121*,
601 10073–10141.

- 602 (46) Wang, J. An Intuitive Tutorial to Gaussian Process Regression. *Computing in Science*
603 *& Engineering* **2023**, *25*, 4–11.
- 604 (47) Wang, A.; Liang, H.; McDannald, A.; Takeuchi, I.; Kusne, A. G. Benchmarking active
605 learning strategies for materials optimization and discovery. *Oxford Open Materials*
606 *Science* **2022**, *2*, itac006.
- 607 (48) Ertl, P.; Schuffenhauer, A. Estimation of synthetic accessibility score of drug-like
608 molecules based on molecular complexity and fragment contributions. *J. Cheminform-*
609 *atics* **2009**, *1*, 8.
- 610 (49) Alnammi, M.; Liu, S.; Ericksen, S. S.; Ananiev, G. E.; Voter, A. F.; Guo, S.; Keck, J. L.;
611 Hoffmann, F. M.; Wildman, S. A.; Gitter, A. Evaluating Scalable Supervised Learning
612 for Synthesize-on-Demand Chemical Libraries. *Journal of Chemical Information and*
613 *Modeling* **2023**, *63*, 5513–5528.
- 614 (50) Pettersen, E. F.; Goddard, T. D.; Huang, C. C.; Couch, G. S.; Greenblatt, D. M.;
615 Meng, E. C.; Ferrin, T. E. UCSF Chimera—A visualization system for exploratory
616 research and analysis. *J. Comput. Chem.* **2004**, *25*, 1605–1612.
- 617 (51) Takeuchi, K.; Kunimoto, R.; Bajorath, J. R-group replacement database for medicinal
618 chemistry. *Future Sci. OA* **2021**, *7*, 8.
- 619 (52) Pedregosa, F. et al. Scikit-learn: Machine Learning in Python. *Journal of Machine*
620 *Learning Research* **2011**, *12*, 2825–2830.
- 621 (53) Danka, T.; Horvath, P. modAL: A modular active learning framework for Python.
622 available on arXiv at <https://arxiv.org/abs/1805.00979>.
- 623 (54) Dask Development Team Dask: Library for dynamic task scheduling. 2024.
- 624 (55) Noske¹, G. D.; de Souza Silva¹, E.; de Godoy¹, M. O.; Dolci¹, I.; Fernandes¹, R. S.;
625 Guido¹, R. V. C.; Sjö, P.; Oliva¹, G.; Godoy, A. S. Structural basis of nirmatrelvir

- 626 and ensitrelvir activity against naturally occurring polymorphisms of the SARS-CoV-2
627 main protease. *Journal of Biological Chemistry* **2023**, *299*, 103004.
- 628 (56) Anstine, D.; Zubatyuk, R.; Isayev, O. AIMNet2: A Neural Network Potential to Meet
629 your Neutral, Charged, Organic, and Elemental-Organic Needs. *ChemRxiv* **2024**,
- 630 (57) Kovács, D. P.; Moore, J. H.; Browning, N. J.; Batatia, I.; Horton, J. T.; Kapil, V.;
631 Witt, W. C.; Magdau, I.-B.; Cole, D. J.; Csányi, G. MACE-OFF23: Transferable
632 Machine Learning Force Fields for Organic Molecules. 2023.
- 633 (58) Fromer, J. C.; Graff, D. E.; Coley, C. W. Pareto optimization to accelerate multi-
634 objective virtual screening. *Digital Discovery* **2024**, *3*, 467–481.
- 635 (59) Powers, A. S.; Yu, H. H.; Suriana, P.; Koodli, R. V.; Lu, T.; Paggi, J. M.; Dror, R. O.
636 Geometric Deep Learning for Structure-Based Ligand Design. *ACS Central Science*
637 **2023**, *9*, 2257–2267.