

1 QSPRpred: a Flexible Open-Source Quantitative
2 Structure-Property Relationship Modelling Tool

3 Helle W. van den Maagdenberg^{1†}, Martin Šícho^{1,2†},
4 David Alencar Araripe^{1,3}, Sohvi Luukkonen^{1,4},
5 Linde Schoenmaker¹, Michiel Jaspers¹,
6 Olivier J. M. Béquignon^{1,5}, Marina Gorostiola González^{1,6},
7 Remco L. van den Broek¹, Andrius Bernatavicius^{1,7}, J. G. Coen
8 van Hasselt¹, Piet. H. van der Graaf^{1,8}, Gerard J. P. van Westen^{1*}

9 ^{1*}Leiden Academic Centre for Drug Research, Leiden University,
10 Einsteinweg 55, Leiden, 2333 CC, The Netherlands.

11 ²CZ-OPENSREEN: National Infrastructure for Chemical Biology,
12 Department of Informatics and Chemistry, Faculty of Chemical
13 Technology, University of Chemistry and Technology Prague, Technická
14 5, Prague, A-4040, Czech Republic.

15 ³Department of Human Genetics, Leiden University Medical Center,
16 Einthovenweg 20, Leiden, 2333ZC, The Netherlands.

17 ⁴ELLIS Unit Linz and LIT AI Lab, Institute for Machine Learning,
18 Johannes Kepler University, Altenberger Straße 69, Linz, 610101,
19 Austria.

20 ⁵Department of Neurosurgery, Brain Tumor Center Amsterdam,
21 Amsterdam University Medical Center, Cancer Center Amsterdam, De
22 Boelelaan 1117, Amsterdam, 1081 HV, The Netherlands.

23 ⁶Oncode Institute, The Netherlands.

24 ⁷Leiden Institute of Advanced Computer Science Leiden University,
25 Niels Bohrweg 1, Leiden, 2333 CA, The Netherlands.

26 ⁸Certara UK, University Road, Canterbury Innovation Centre, Unit 43
27 Canterbury, Kent, CT2 7FG, UK.

28 *Corresponding author(s). E-mail(s): gerard@lacdr.leidenuniv.nl;

29 Contributing authors: h.w.van.den.maagdenberg@lacdr.leidenuniv.nl;
30 m.sicho@lacdr.leidenuniv.nl; d.figueiredo.vidal@lacdr.leidenuniv.nl;
31 luukkonen@ml.jku.at; l.schoenmaker@lacdr.leidenuniv.nl;
32 m.jespers@lacdr.leidenuniv.nl; o.j.m.bequignon@lacdr.leidenuniv.nl;
33 m.gorostiola.gonzalez@lacdr.leidenuniv.nl;
34 r.l.van.den.broek@lacdr.leidenuniv.nl;
35 a.bernatavicus@liacs.leidenuniv.nl; coen.vanhasselt@lacdr.leidenuniv.nl;
36 p.vandergraaf@lacdr.leidenuniv.nl;

37 †These authors contributed equally to this work.

38 Abstract

39 Building reliable and robust quantitative structure-property relationship (QSPR)
40 models is a challenging task. First, the experimental data needs to be obtained,
41 analyzed and curated. Second, the number of available methods is continu-
42 ously growing and evaluating different algorithms and methodologies can be
43 arduous. Finally, the last hurdle that researchers face is to ensure the repro-
44 ducibility of their models and facilitate their transferability into practice. In
45 this work, we introduce QSPRpred, a toolkit for analysis of bioactivity data
46 sets and QSPR modelling, which attempts to address the aforementioned chal-
47 lenges. QSPRpred’s modular Python API enables users to intuitively describe
48 different parts of a modelling workflow using a plethora of pre-implemented com-
49 ponents, but also integrate customized implementations in a “plug-and-play”
50 manner. QSPRpred data sets and models are directly serializable, which means
51 they can be readily reproduced and put into operation after training as the
52 models are saved with all required data pre-processing steps to make predic-
53 tions on new compounds directly from SMILES strings. The general-purpose
54 character of QSPRpred is also demonstrated by inclusion of support for multi-
55 task and proteochemometric modelling. The package is extensively documented
56 and comes with a large collection of tutorials to help new users. In this paper,
57 we describe all of QSPRpred’s functionalities and also conduct a small bench-
58 marking case study to illustrate how different components can be leveraged to
59 compare a diverse set of models. QSPRpred is fully open-source and available at
60 <https://github.com/CDDLeiden/QSPRpred>.

61 Scientific Contribution

62 QSPRpred aims to provide a complex, but comprehensive Python API to conduct
63 all tasks encountered in QSPR modelling from data preparation and analy-
64 sis to model creation and model deployment. In contrast to similar packages,
65 QSPRpred offers a wider and more exhaustive range of capabilities and inte-
66 grations with many popular packages that also go beyond QSPR modelling. A
67 significant contribution of QSPRpred is also in its automated and highly stan-
68 dardized serialization scheme, which significantly improves reproducibility and
69 transferability of models.

70 **Keywords:** QSPR modelling, QSAR modelling, proteochemometrics,
71 cheminformatics, machine learning, software

72 1 Introduction

73 Quantitative Structure-Property Relationship (QSPR) modelling can be described
74 as the application of empirical methods (i.e. statistical and machine learning (ML)
75 approaches) to find the mathematical relationship between molecular structure and a
76 property of interest [1]. In the past decades, QSPR and mainly Quantitative Structure-
77 Activity Relationship (QSAR) modelling methods have established themselves as key
78 instruments in drug discovery [2–6] and beyond [1, 7]. Reliable QSPR models have the
79 potential to reduce the need for time and resource intensive experimental screening of
80 compounds by enabling effective compound selection in the drug development pipeline
81 [8]. The increasing amount of experimental data available (i.e., in ChEMBL [9] and
82 PubChem [10]) also enables the use of more advanced methods which have seen rapid
83 development [11, 12]. Therefore, given the prevalence of QSPR modelling and its
84 constant development, there is a need for software tools that can support researchers
85 not only in the tasks of developing and deploying models in practice, but also in
86 the critical assessment of new methods and validation of non-trivial computational
87 workflows that include many preliminary steps such as collection, curation and analysis
88 of data as well as model training, evaluation and deployment.

89 In the traditional sense, QSPR modelling focuses mainly on describing the
90 relationship between the compound structure and a property of interest, but pro-
91 teochemometric modelling (PCM) has emerged as an extension that also introduces
92 the protein target information into the equation [13, 14]. A PCM approach can extrap-
93 olate similarities and differences across (super)families and is therefore promising
94 in poly-pharmacology and off-target prediction [15], as well as a strategy for data
95 augmentation and relevant binding residue identification [16, 17]. Although in the tra-
96 ditional sense, the architecture is identical to that of a single-task model, it includes
97 bio-activity endpoints for multiple proteins, by featurizing each compound-protein
98 combination separately [18]. Therefore, PCM has inherent applicability challenges that

99 combine those of single-task and multi-task modelling (i.e. dataset size [19, 20], data
100 balance [21, 22], sparsity [23], but also unique featurization requirements that need to
101 take the proteins themselves into account [19, 20]).

102 No matter the underlying philosophy, both traditional and PCM QSPR models
103 can be obtained by combining various algorithms and methodologies, which often
104 need to be benchmarked against one another in a systematic and comprehensive man-
105 ner. While there is still an ongoing debate on whether and under what conditions a
106 meaningful comparison of QSPR methodologies is possible [24], benchmarking and sys-
107 tematic comparison have become an integral part of the QSPR field. Whether it is the
108 comparison of new algorithms [18, 25, 26], molecular representations [27–29], model
109 development strategies [20, 21], model validation [30, 31], the nature of data [32], or
110 to provide usage guidelines [33], one common denominator of such studies is that they
111 base their conclusions on a systematic comparison using a standardized subset of data.
112 During this task, researchers are faced with many challenges from compiling a repre-
113 sentative subset of data for the diverse set of tasks seen in QSPR modeling [29] to
114 choosing the right method to obtain statistically sound results [34–36]. However, even
115 more fundamental problems such as the dominance of median predictions [24] or com-
116 bining data sources [37] can plague benchmarking results. Therefore, even with the
117 long history of the QSPR modeling field, it is clear that benchmarking methodologies
118 are important, but not problem-free.

119 Another challenge that QSPR modellers face is the reproducibility of results and
120 model deployment. While not specific to QSPR modelling [38, 39], reproducibility is a
121 topic widely discussed in cheminformatics and computational drug discovery [40–42]
122 and pertains mainly to correct estimation of real-world performance data, but also
123 the practical transition from model building and evaluation to deployment [43]. For
124 example, the deployment phase also needs to implement crucial steps to process com-
125 pound structures before prediction to ensure equivalent compound representation as

126 during training. However, currently, there is a lack of open-source tools that would
127 sufficiently address the reproducibility of results and the transfer of models into prac-
128 tice and it is often up to the modeller to provide these, which leads to a large disparity
129 between researchers in how reproducibility and deployment of models is addressed.

130 Several open-source applications are available that support researchers in QSPR
131 modelling and help in solving some of the aforementioned challenges. A popular
132 framework is KNIME [44], which utilizes a GUI with visual workflows and has many
133 pre-implemented components. However, designing custom components in KNIME can
134 be challenging and the integration of Python extensions in the Java-based API is not
135 always straightforward [45]. With a focus on deep-learning-based models, DeepChem
136 [46], was one of the pioneering Python packages for molecular modelling. It offers a
137 wide array of different featurizers and models and a flexible and easy to understand
138 API that is modular and extensible. However, not all DeepChem models address the
139 aforementioned reproducibility and deployment issues out-of-the-box. For example, the
140 offered `SklearnModel` class does serialize the model itself, but reproducing the prepa-
141 ration workflow and creating the feature matrix is left to the users themselves, which
142 can be inconvenient or, worse, might create reproducibility and deployment problems.
143 Extending DeepChem, AMPL [47] prioritized automated machine learning for bench-
144 marking and it enables users to conveniently build and validate models. However, it
145 still lacks functionality to readily deploy and use models in practice. ZairaChem [48]
146 is another recent package, which proposes an automated cascade for training machine
147 learning models, empowering users with little knowledge in data science to train robust
148 ensemble-based models. However, one limitation of ZairaChem is that it currently only
149 supports classification models and also does not enable model serialization with prepa-
150 ration steps included. This is addressed in a recently published package, PREFER
151 [49], which wraps trained models fully, including data preprocessing. It implements a
152 pipeline based on AutoSklearn [50], covering steps such as data preparation, model

153 selection, and model evaluation. However, PREFER, offers a slightly less flexible API
154 than the previous options and combining different feature representations and split-
155 ting methods is not possible without modifying the source code of the package itself.
156 One more package that supports comprehensive serialization of data preprocessing
157 steps within produced models is Qptuna [51]. It features a modular API with variety
158 of pre-implemented algorithms and featurizers, as well as a focus on model explain-
159 ability. However, due to its focus on hyperparameter optimization and streamlining
160 the modelling process the API is not as rich and extensible as in the case of some
161 other packages such as DeepChem. For example, it might be tedious to use differ-
162 ent hyper-parameter optimization strategies not implemented by the Optuna package,
163 which Qptuna is based on. All of the aforementioned packages also lack support for
164 PCM modelling with Qptuna as a notable exception with its support for simple Z-
165 scale descriptors. Even though there exists a package with pure focus on PCM, an R
166 tool called `camb` (Chemically Aware Model Builder) [52], it has not been maintained
167 since 2017 [53]. Therefore, support for accessible and straightforward PCM modelling
168 is still lacking among contemporary open source packages. Similarly, the inclusion of
169 applicability domain of QSPR models is also not fully considered in any of the above
170 packages.

171 In this work, we present QSPRpred, an open-source package that attempts to com-
172 pile the essentials of QSPR modelling into a compact and accessible Python package.
173 As such, the package aims to address a user base with varying ranges of expertise
174 from students to well-rounded QSPR modellers interested in developing and testing
175 new approaches. With QSPRpred we try to provide high-level interfaces to accomplish
176 QSPR modelling tasks in few steps, but at the same time try to encourage writing of
177 modular Python code by making sure all variable steps are encapsulated and easily
178 replaceable by custom implementations. In addition to being customizable, all steps

179 can also be combined using the built-in benchmarking workflow, which enables stream-
180 lined model building to the likes of AMPL or Qptuna, but completely described in
181 Python with "plug-and-play" components. This provides a platform for researchers
182 to experiment and validate novel approaches quickly. In addition, QSPRpred tackles
183 the problem of reproducibility and deployment of models by streamlining setting of
184 random seeds and using a global serialization model that not only includes the model
185 itself, but also the molecule preparation and featurization steps. This means that a
186 shared model can be reloaded and directly used for predictions from SMILES strings
187 without the need to repeat any preparation steps.

188 **2 Implementation**

189 Although a wide variety of different QSPR models and descriptors have been described
190 in literature [1], there is a significant overlap in the general QSAR/QSPR modelling
191 workflow. QSPRpred leverages this by providing a modular framework that comes
192 with many out-of-the-box components while also providing clear interface definitions
193 to facilitate comprehensive extensibility features (Figure 1). The description of vari-
194 ous aspects of this workflow will be the subject of the following subsections: Section
195 2.1 includes an extensive description of all the data preparation components that
196 QSPRpred provides. Section 2.2 provides relevant details on model training and eval-
197 uation. In addition, QSPRpred also provides various visualization tools, which are
198 described in section 2.2.6.

199 **2.1 Data**

200 **2.1.1 Data collection**

201 The first step of any QSPR modelling project is the collection of data. The supplied
202 data should at least contain the molecule SMILES sequences and the property to be
203 predicted, but can contain any extra information as needed. Furthermore, data can

204 be retrieved from an external source programmatically, using the `DataSource` class,
205 which is used to describe the creation of data sets from different sources. By default
206 QSPRpred provides the `Papyrus` data source, which can be used to collect data from
207 the eponymous dataset Papyrus, a large-scale curated bioactivity dataset [54]. This
208 method of data collection is also used for creating multiple benchmarking datasets
209 dynamically, see section 3 for an example. Fetching data with a `DataSource` will
210 directly return the data in a `MoleculeTable` object or a `QSPRDataset`. `MoleculeTable`
211 is a data container that provides functionality to handle different operations for
212 molecule preparation and descriptor calculation while `QSPRDataset` is its extension
213 that provides functionality specific to QSPR modelling such as data splitting.

214 An instance of `QSPRDataset` can also be initialized directly (without a `DataSource`)
215 by the user through a tabular format (e.g. CSV/TSV) or an SDF file. The user always
216 needs to specify one or more properties to be predicted during initialization or dur-
217 ing the lifetime of the data set. These are specified as `TargetProperty`, which are
218 associated with a specific modelling `TargetTask`, such as `REGRESSION` or `MULTICLASS`.

219 In the following sections, several data pre-processing and preparation steps will be
220 described which can be applied to the `QSPRDataset`. The `prepareDataset` method
221 `QSPRDataset` streamlines the process by applying the data preparation steps in a fixed
222 order: data filtering, descriptor calculation, data splitting, feature filtering and feature
223 standardization. The user can specify which components should be used for each step,
224 which includes custom components that can be implemented by creating subclasses
225 of the abstract base class for each component. The `QSPRDataset` and `MoleculeTable`
226 are serializable to JSON (JavaScript Object Notation) and other associated files (see
227 section 2.2.7), which enable the user to save and reload prepared datasets in machine-
228 and human-readable format.

229 **2.1.2 Data pre-processing**

230 Before being used for model fitting, data samples may be removed based on the follow-
231 ing criteria: specific values (`CategoryFilter`, e.g. remove all samples from a specific
232 source or year), identical descriptor values (`RepeatsFilter`) or user-defined filters.
233 Moreover, as the underlying pandas [55] dataframe can be accessed and edited, any
234 further data analysis steps can be executed in this phase. For easy visualization and
235 exploratory analyses of the dataset, integration with Scaffviz [56] is provided, which
236 is described in more detail in the visualization section 2.2.6.

237 **2.1.3 Descriptor calculation**

238 Descriptor calculation in QSPRpred is facilitated through the `DescriptorSet` class
239 which wraps one or more types of descriptors. There are many integrated implemen-
240 tations of the `DescriptorSet` available including: RDKit descriptors [57], Mordred
241 [58], Mold² [59], PaDEL [60] and Tanimoto distance to other molecules. Additionally
242 a range of different types of molecular fingerprints are implemented, such as Mor-
243 gan and MACCs [61] fingerprints. A trained QSPRpred model can itself also be used
244 as descriptor for stacked modelling. While users can add new implementations of
245 `DescriptorSets` this may not always be practical, for example when using descrip-
246 tors from experimental measurements. In this case, the descriptors can be provided
247 as a data frame directly via the built-in `DataFrameDescriptorSet`. Moreover, a cus-
248 tom implementation of the `DescriptorSet` interface, also exists for protein descriptors
249 (`ProteinDescriptorSet`) that are commonly used for PCM modelling (i.e. z-scales
250 [62], BLOSUM [63] and VHSE [64]). This is provided via the standalone ProDEC
251 [65] package, which enables calculation of multiple sequence alignment-based descrip-
252 tors for proteins. The multiple sequence alignment is done with Clustal Omega [66] or
253 MAFFT [67], but QSPRpred also describes an API to easily integrate other alignment
254 methods.

255 **2.1.4 Feature filtering**

256 Feature filtering is commonly used to select the most informative features from the
257 calculated descriptors. Here, the feature filters currently implemented in QSPRpred
258 will be discussed, however, new feature filters can be easily added by the user. Filters
259 in QSPRpred are always calculated using only the training set or the training fold
260 of the data set to avoid data leakage (see subsection 2.1.5). The first filter that is
261 implemented is the `LowVarianceFilter`, which calculates the variance within a feature
262 on the training set and removes it if it is below a threshold specified by the user.
263 A `HighCorrelationFilter` is also available. It calculates the correlation between
264 features and removes them if it is higher than the user-specified threshold. Finally,
265 QSPRpred provides integration with BorutaPy [68], which is a Python implementation
266 of the Boruta filter [69]. Boruta filtering is an all-relevant feature selection method
267 that removes features based on their relevance compared to random features.

268 **2.1.5 Data splitting**

269 The choice of how to split data can greatly influence our impression of future per-
270 formance of the model [22, 70]. QSPRpred supports any scikit-learn-style [71] data
271 splitter class that has a method `split(X,y)` that yields for each split/fold the indices
272 for each subset. Splits can be applied at two levels, during dataset preparation to cre-
273 ate the independent test set (which will be used by `TestSetAssessor`) and during
274 model optimization/training to create (cross-)validation sets with `CrossValAssessor`
275 (described in section 2.2.3). QSPRpred also offers several integrated splits. All
276 of these methods support the creation of a single train-test split or k -folds.
277 Firstly, `SklearnRandomSplit` which wraps scikit-learn's `(Stratified)ShuffleSplit`
278 method. Three splits (`RandomSplit`, `ClusterSplit`, `ScaffoldSplit`) are included
279 that utilize the `BALANCESPLIT` [72] package to create well-balanced data splits for
280 (sparse) datasets without data leakage between different tasks. The `ManualSplit` can

281 be used to apply a predefined split. The `TemporalSplit` is used to make time series-
282 based test or cross-validation splits. `QSPRpred` also has a `BootstrapSplit` class that
283 can wrap any split and use it for repeated sampling of the dataset applying the specified
284 split in replicates. Furthermore, any of the above-mentioned splitters can be applied
285 to PCM modelling with the help of the `PCMSplit` class, which will ensure the splits are
286 balanced with respect to the protein targets. Two other PCM-specific splitters are also
287 available `LeaveTargetOut` and `TemporalPerTarget`. `LeaveTargetOut` will remove all
288 the data points for a certain target, to evaluate how well the PCM model extends to
289 new targets. The `TemporalPerTarget` split applies a temporal split that avoids data
290 leakage with multiple tests for compounds for different targets over time, by using the
291 first occurrence of the molecule in the dataset for all proteins as timepoint.

292 **2.2 Modelling**

293 **2.2.1 General**

294 `QSPRpred` inherently supports both single and multi-task variants of regression and
295 single-class and multi-class classification. All models implemented with `QSPRpred`
296 are wrapped in a `QSPRModel` subclass, which can then be applied to a `QSPRDataset`
297 instance. The model task is automatically derived from the target properties specified
298 in the dataset, i.e. two single-class `TargetTask` target properties, will result in a multi-
299 task single-class model. Model tasks in `QSPRpred` are encoded as a Python `Enum` class
300 (e.g. `REGRESSION`, `SINGLECLASS`, `MULTITASK_MULTICLASS`), which allows easy specifica-
301 tion of which tasks a model can support. Like the dataset object, a `QSPRModel` can be
302 serialized to JSON (see section 2.2.7). Currently, available model implementations are:
303 `SklearnModel`, `DNNModel`, `PyBoostModel` and `ChempropModel`. The `SklearnModel` is a
304 wrapper for all scikit-learn estimators [71]. The `DNNModel` is a PyTorch [73] implemen-
305 tation of a fully-connected neural network. `PyBoost` is a wrapper around the Py-Boost
306 [74] package, a Gradient Boosted Decision Tree toolkit. The `ChempropModel` wraps

307 the basic Chemprop [75] message passing neural network, although it does not sup-
308 port all functionality that Chemprop provides. Moreover, any new type of model can
309 be simply implemented by creating a subclass of the abstract base class `QSPRModel`,
310 requiring the user to just implement methods for fitting, predicting and serialization of
311 the model. There is also a dedicated tutorial to implementing new models (see section
312 2.5) to help QSPR practitioners with integration of novel methods.

313 **2.2.2 PCM modelling**

314 Creating a PCM model differs slightly from creating a standard QSPR model. PCM
315 models require protein featurization and can introduce the need for different splitting
316 methods, as discussed in section 2.1.3 and 2.1.5 respectively. To create a PCM dataset,
317 `QSPRpred` provides a class called `PCMDataset` which forms the alternative input for
318 a model of the class `PCMModel`, which is slightly altered from the base `QSPRModel` and
319 can be used to wrap `QSPRModel` implementations for PCM. This is mainly necessary
320 because in order to make predictions with a trained model, the protein identifier of
321 the protein to make predictions for is needed.

322 **2.2.3 Model Assessment**

323 To evaluate model performance, `QSPRpred` provides a class called `ModelAssessor`
324 that defines a structure for performance evaluation. Given a `QSPRpred` model and a
325 dataset, it will run an evaluation and return the specified metric. All scikit-learn [71]
326 scoring functions can be used or custom ones can be provided to the `ModelAssessor`.

327 By default, two `ModelAssessor` subclasses are implemented, namely the
328 `TestSetAssessor` and the `CrossValAssessor`. As the name suggests, the
329 `TestSetAssessor` evaluates the model performance on the test set of the provided
330 dataset. It will use the model to predict values for the test set and return the score
331 given by the provided metric. On the other hand, the `CrossValAssessor` can perform
332 cross-validation on the training set. The folds are determined by the user-specified

333 splitting method (see section 2.1.5). It will return a list of scores for each fold. Because
334 of the flexible splitting method, this class is not limited to cross-validation, but can
335 also perform bootstrapping, through resampling of the training set (see section 2.1.5).

336 Both described `ModelAssessor` implementations will write the predictions to the
337 model directory in TSV format, including unique molecule identifiers provided by the
338 dataset. Using the identifiers each prediction can be linked back to the original data
339 point, which allows for detailed analysis and visualization (see section 2.2.6) of the
340 model performance.

341 **2.2.4 Hyperparameter optimization**

342 Finding the right hyperparameters for a model, can be a challenging task. As
343 with previously discussed components, QSPRpred provides a flexible base class
344 `HyperparameterOptimization`. Hyperparameter optimization requires specification
345 of the search space; which model parameters to tune and for which values or within
346 which bounds. A template of how to provide the search space file can be found in the
347 documentation [76]. Furthermore, a `ModelAssessor` needs to be specified, that deter-
348 mines how a set of hyperparameters will be evaluated. If the `ModelAssessor` returns
349 multiple scores, e.g. in the case of cross-validation a score for each fold is returned,
350 the score will be aggregated by a user-specified function.

351 Two default implementations of `HyperparameterOptimization` exist: `GridSearch`
352 and `OptunaOptimization`. `GridSearch` is an exhaustive search algorithm, that eval-
353 uates all combinations of the specified hyperparameters. `OptunaOptimization` is a
354 form of Bayesian optimization using Optuna's [77] Tree-structured Parzen Estimator
355 algorithm as the acquisition function. Bayesian optimization is used for iteratively
356 proposing new hyperparameters for a machine learning model applying the Bayesian
357 principle, which allows for finding the optimal hyperparameter combinations without
358 an exhaustive search. The user needs to set the number of iterations because it is an

359 iterative process. By default the search starts with a random sample of ten parameter
360 combinations, afterwards, the acquisition function is used.

361 **2.2.5 Applicability Domain**

362 For evaluation of the applicability domain of trained models, QSPRpred provides
363 integration with MLChemAD [78]. MLChemAD implements many commonly used
364 definitions of the applicability domain for cheminformatics models, based on k-nearest
365 neighbors, the local outlier factor or on bounding approaches. It is also possible to
366 implement a custom applicability domain using the `ApplicabilityDomain` base class.
367 An `ApplicabilityDomain` object may be attached to a `QSPRDataset`. Then during
368 the dataset preparation, the applicability domain can be fit on the training set to
369 identify or remove outliers from the test set. Furthermore, the `ApplicabilityDomain`
370 object can be attached to a model and fit on the whole dataset. In production mode,
371 when predicting from SMILES the model will return whether this compound is within
372 the applicability domain of the trained model.

373 **2.2.6 Visualization**

374 The test set and cross-validation assessments write to result files in the model directory
375 as described in section 2.2.3. These result files are human readable and can be used
376 to easily generate any visualization users require. QSPRpred also has a `ModelPlot`
377 class, that provides a number of different plots that can be generated directly from
378 an instantiated model with result files present. These plots include receiver operating
379 characteristic (ROC) curves, precision-recall curves, calibration plot and barplots for
380 a range of scikit-learn [71] classification metrics. Metrics such as precision, recall and
381 Matthews Correlation Coefficient can be visualized not only for single-task models, but
382 also for multi-task and multi-class classification models. In the case of multi-class mod-
383 els metrics are calculated per class (one-vs-rest) and with different averages. Moreover,

384 correlation plots can be generated for multi and single-task regression models. In addi-
385 tion to the native `ModelPlot`, QSPRpred’s `MoleculeTable` and `QSPRModel` instances
386 can be used directly with the interactive cheminformatics visualization package Scaf-
387 fviz [56]. Scaffviz offers alternative visualization of model errors and is essentially an
388 adapter between molplotly [79] and QSPRpred. It can be used to apply dimension-
389 ality reduction methods to obtain 2D embeddings of molecules from descriptors and
390 display them in an interactive scatter plot. However, any properties from the data
391 set can be displayed on each axis as well. Points may be coloured by any property
392 in the data set, including training and test splits. It can also visualize model errors
393 for mispredicted compounds as color overlay as well, which helps identifying difficult
394 compounds to predict.

395 **2.2.7 Reproducibility & Transferability**

396 In accordance with the R(eusability) of the FAIR principles [80], almost everything in
397 the QSPRpred API is serializable to a human-readable file format. The vast major-
398 ity of objects are serializable to a JSON file, which can be read easily even without
399 QSPRpred, including model parameters of created models and workflow settings. This
400 is possible thanks to the `ml2json` package [81] and the `jsonpickle` [82] project. Since
401 `pandas.DataFrame` instances are used to represent all tabular data, they can be saved
402 to several formats, including human-readable `.csv` files. When a human-readable rep-
403 resentation is not possible (i.e. with deep learning models), sufficient metadata is
404 saved to be able to recreate it as closely as possible. Additionally, the random state of
405 QSPRpred is globally configurable, allowing for full reproducibility of results involv-
406 ing random operations. The random seed, whether newly generated or passed by the
407 user, is saved to metadata and can be re-used to get the same results. Each dataset
408 is initialized with a single random state, which is adopted by models and used in all
409 subsequent random operations. However, QSPRpred also allows further fine-grained

410 control over the random state by having the option to control the random state of
411 models and splits separately from the dataset.

412 **2.3 Architecture**

413 The structure of the QSPRpred package reflects the intended usage as outlined in
414 Figure 1 with several packages and subpackages separating the different tasks (Figure
415 2). Every data structure and functional element of the API has its own abstract def-
416 inition, which is followed in the reference implementations. For example, the main
417 `QSPRDataset` (located in `qsprpred.data.tables.qspr`) class implements several
418 interfaces defined in `qsprpred.data.tables.base`. Likewise, the `QSPRModel` abstract
419 base class (located in `qsprpred.models.base`) defines the API of a model while the
420 `qsprpred.models.sklearn` is its implementation that facilitates a compatibility layer
421 between the `scikit-learn` [71] package and QSPRpred. Such a standardized approach
422 to API development makes integration of new tools and data structures easier and
423 changes to code minimal as libraries update over time.

424 The modular architecture of QSPRpred also makes optional installation of depen-
425 dencies possible. Much of the functionality located in `qsprpred.extra` depends on
426 external packages that can sometimes pull many dependencies alongside them, but
427 with careful modular separation it is not necessary to install those dependencies
428 unless they are truly needed. Therefore, QSPRpred also supports various instal-
429 lation flags to make sure only necessary dependencies are installed for different
430 intended use cases. For example, the dependencies needed to support functionalities in
431 `qsprpred.extra.gpu` are only installed if the `[gpu]` flag is specified upon installation,
432 but without it the rest of the package will still function normally.

433 In addition to the flexible Python API, QSPRpred also offers extensive command
434 line interface (CLI). While less customizable than the API, the CLI allows the user to
435 train a wide variety of QSPR models without having to write any code. The QSPRpred

436 CLI is subdivided into three callable scripts (`data_CLI`, `model_CLI` and `predict_CLI`)
437 that cover the entire QSPR model building workflow (see figure 1).

438 With the `data_CLI` it is possible to create a range of datasets for different tasks with
439 one command, including datasets for multi-class and multi-task modelling (see section
440 2.2). Furthermore, most of the base QSPRpred data pre-processing functionality is
441 available through this CLI, including all the different descriptors sets, and data splits
442 as well as target imputation for multi-task modelling.

443 After creating one or more datasets through the CLI or the Python API, the model
444 CLI can be used for model training. With the model CLI a range of scikit-learn [71]
445 models and a PyTorch [73] fully-connected neural net can be trained. The `model_CLI`
446 provides the same main training steps as the Python API, namely hyperparameter
447 optimization, cross-validation and test set predictions. Finally, after model training
448 has completed, the trained models can be used to make predictions on new sets of
449 molecules with the `predict_CLI`.

450 2.4 Documentation

451 In QSPRpred, much effort is also devoted to guiding new users as well as potential
452 contributors. The code follows a predetermined style guide [83], which requires that
453 every functionality needs to be properly documented via a docstring that is then
454 visible on the QSPRpred documentation page [76] upon publishing a new release.
455 The style guide also requires that Python type hints are present for all methods
456 and functions to indicate which data structures are compatible and expected. Pre-
457 commit hooks are also available that can be used to check code before committing to
458 ensure compliance with the style guide. In addition to these strict API documentation
459 requirements, the documentation pages also contain guides on how to use the CLI and
460 other miscellaneous items pertaining to the usage of the package.

461 2.5 Tutorials

462 We have dedicated several Jupyter notebooks to featuring real-life QSPR modelling
463 examples. The tutorials expect basic understanding of QSPR modelling concepts and
464 are designed to mainly showcase QSPRpred functionality. While they are easily acces-
465 sible by beginners and are suitable to be used by students, they also include specialized
466 notebooks for more advanced users who want to customize behaviour and integrate
467 new methods.

468 The user progressively learns all functionalities within QSPRpred from data set
469 acquisition to model evaluation. A quick start tutorial is designed to get the user to
470 prepare a dataset and run a single-task QSAR regression model with QSPRpred as
471 quickly as possible. After the quick start tutorial, a series of one, seven and three tuto-
472 rials covering respectively the basics of benchmarking, data handling and modelling
473 within QSPRpred are available. These tutorials go over the main aspects that need
474 to be taken into account in any modelling project prior to modelling (i.e data collec-
475 tion, preparation, featurization and splitting), but also the necessary tools to build
476 and validate models (i.e. formulation of model tasks, model assessment and logging of
477 progress and results). These basic tutorials can be accessed individually or followed
478 sequentially from the quick start guide.

479 On top of the basic tutorials, a series of ten advanced data and modelling tuto-
480 rials are available. These help the user to build on top of the already acquired basic
481 knowledge of QSPRpred by teaching them how to customize functionalities and add
482 new features. These tutorials are aimed at researchers who develop new methods and
483 want to take advantage of QSPRpred to automate and standardize certain tasks. The
484 advanced tutorials showcase the possibilities to perform hyperparameter optimization
485 and add custom descriptors and data splitting methods, as well as custom models.
486 Moreover, these tutorials also dive deeper into the modelling options by showing how

487 to build deep learning, multi-task, and PCM models, and advanced model monitoring
488 via the popular Weights & Biases (W&B) framework [84].

489 **2.6 Testing**

490 In order to ensure that all functionalities of the package remain operational even as
491 the code changes, be it by internal factors (code updates) or external (dependency
492 updates), frequent testing of the code is necessary. Therefore, the inherent part of
493 the QSPRpred package is unit testing and every new functionality added needs to be
494 accompanied with testing code.

495 QSPRpred also takes advantage of continuous integration (CI) not only to run unit
496 tests, but also to frequently run the tutorial code and check consistency of models.
497 Therefore, it is always ensured that upon any modification of the code all tutorials are
498 up to date, all code is working as expected and all reference models return the same
499 expected results. The latter is especially important to guarantee model consistency
500 across different versions of the code and its dependencies by highlighting changes that
501 could lead to past results being unreproducible.

502 **3 Results**

503 QSPRpred also provides an overarching API to conduct benchmarking experiments.
504 These features are located in the `qsprpred.benchmarking` package and provide a
505 streamlined way to test various combinations of molecular descriptors, model algo-
506 rithms and even data set preparation strategies. In this section, we will show two
507 example scenarios of experiments focused on building and comparing regression mod-
508 els. The code to reproduce these experiments is available at [https://github.com/
509 CDDLeiden/qsp-bench](https://github.com/CDDLeiden/qsp-bench), but the repository can be easily extended and adapted to
510 other scenarios as well using the instructions within.

511 3.1 Experiment 1: Benchmarking Single-Task and Multi-Task 512 Regression Models

513 A multi-task modeling approach can be beneficial in modeling several endpoints at
514 once for similar biological targets [85]. However, it is not always clear if such an
515 approach will lead to an improvement over the more traditional single-task modeling
516 or what multi-task methodology would be the most optimal for the data at hand.
517 There is a plethora of methods that could be considered. As a result, researchers are
518 often confronted with a large selection of viable workflows and methods [86].

519 For this small case study, we chose a bioactivity data set of 4 adenosine receptors
520 (A1, A2A, A2B and A3). The adenosine receptors are a highly conserved family of
521 enzymes that share many similarities and selective modulators of adenosine receptors
522 are of interest in drug discovery. Therefore, many compounds that share structural
523 similarities are often tested against multiple or all of these receptors. This makes multi-
524 task modeling an eligible method to consider when creating a QSAR model for these
525 receptors. The data set used in this study was assembled by querying the Papyrus
526 data set (version 05.6) on the respective UniProt accession keys (P30542, P29274,
527 P29275 and P0DMS8) using QSPRpred's integration. Only minor modifications to
528 the adapter were made through inheritance to facilitate multi-task modeling and the
529 adapted implementation is available in the `qsp-bench` repository. Compounds in this
530 data set were represented by Morgan fingerprints with radius 3 and bit length of 2048
531 (as implemented in RDKit [57]).

532 The choice of models, in this case, study was motivated by the recently added
533 multi-task capability to the popular `xgboost` package [87]. Since the models imple-
534 mented in this package adhere to the `scikit-learn` API, they can be readily used
535 in QSPRpred with the `SklearnModel` class. For the multi-task scenario, we compared
536 the two modelling strategies currently implemented by `xgboost`: (1) Multi Output
537 Tree (MOT) and (2) One Output Per Tree (OOPT).

538 The goal of the case study was to compare these two multi-task strategies with
539 a baseline `KNeighborsRegressor` algorithm from the popular `scikit-learn` package
540 [71], but also against a simple `MedianModel` model inspired by the recent work of Janela
541 et al. [24], which predicts the median target property value for every test instance.
542 In all workflows, the target property was the median pCHEMBL value as determined
543 from the Papyrus data set [54]. All experiments were conducted in 30 replicas and
544 using either a standard random train-test splitting strategy or a cluster split strategy,
545 `RandomSplit` and `ClusterSplit` classes in the QSPRpred API, respectively. For each
546 model, we report the R^2 metric on each task separately 3.

547 Overall, we found no benefit in using a multi-task model over a single-task model
548 when using 30 repeated experiments (Figure 3). In fact, the multi-task models showed
549 slightly worse performance overall in both the random split and clustered split in all
550 tasks (Figure 3). This is likely due to the sparsity of the multi-task target variables
551 and the effect of imputation.

552 It can also be seen that the cluster split (Figures 3C and 3D) is more difficult
553 than a standard random split (Figures 3A and 3B) for all models as indicated by a
554 significant drop in performance for both single-task and multi-task models. This is
555 expected since the clustered split is designed with the intent to present the model with
556 more difficult examples in the test set.

557 In addition, the performance shows more variance across the scaffold split than
558 the random split. This is likely due to the difficulty of the test set fluctuating more
559 with the varying similarity of the test compounds to the training set, which depends
560 on the clusters created (Figure 3).

561 In all experiments the `xgboost` models were comparable to the
562 `KNeighborsRegressor` model (Figure 3). For the multi-task case the 00PT strategy
563 worked slightly worse or was comparable to MOT (Figures 3B and 3D). Contrary to

564 findings of Janela et al. [24], we did not see elevated performance of the `MedianModel`
565 baseline on this particular data set and benchmarking conditions.

566 Therefore, using multiple single task models is likely a better choice for the data at
567 hand. However, it should be noted that more models and strategies could be consid-
568 ered here and also compared in a follow up study. For example, deep learning models
569 may be better-suited for multi-task modeling than the algorithms tested herein. Their
570 architecture can be adapted in a number of ways to accommodate multi-task learn-
571 ing [88]. Another point to consider is the influence of sparsity of the data [89]. We
572 simply imputed missing labels with a median value in our experiments, which may
573 have introduced large bias to these central values in the multi-task models presented.
574 Surprisingly, this did not affect the performance of the `MedianModel`, which should
575 theoretically benefit from such imputation. Finally, it should also be noted that proper
576 statistical testing should also be conducted to determine under which conditions the
577 performance of models truly significantly differs.

578 **3.2 Experiment 2: Comparison of Regression Models of** 579 **Different Architectures**

580 One problem that QSPRpred is trying to address is how to bring models built with dif-
581 ferent algorithms and, thus, different software requirements under one roof and how to
582 run and benchmark them with as similar API as possible. Therefore, in this case study
583 we show an example that integrates and compares both the `XGBoostRegressor` model
584 and a deep learning based method `Chemprop` in one benchmarking experiment on sev-
585 eral data sets. Both models have different hardware and software requirements with
586 `XGBRegressor` being CPU-based and accepting fingerprints as input and `Chemprop`
587 requiring GPU-accelerated training and raw SMILES as input. With the exception of
588 raw standardized SMILES for `Chemprop`, the same fingerprint, replica counts, splitting

589 strategies and evaluation metrics were used in this case as well. The `MedianModel` was
590 again used as a simple baseline.

591 We also switched to four MoleculeNet benchmarking sets for this experiment to
592 show how data sets from this source can be integrated for benchmarking. The data
593 sets we chose have a diverse set of target properties. In particular we chose to evaluate
594 the predictivity of the built models on the Lipophilicity, clearance, solubility and free
595 solvation energy (Figure 4).

596 Clearance is known to be a difficult property to predict [90] and was also the
597 worst performing data set in our experiments (Figure 4). In the cluster split scenario,
598 the `XGBRegressor` did not even perform above the `MedianModel` baseline (Figure
599 4B). The `XGBRegressor` was also the inferior model in all scenarios we observed.
600 However, it should be noted we did not optimize hyperparameters for these models
601 and just used them with their defaults. Therefore, it is possible a better performance
602 could still be achieved for `XGBRegressor` as well as `ChemProp`. Again, we observed
603 that the clustered split showed degraded performance across all data sets and the
604 variance between replicas was larger (compare Figure 4A and 4B). This is consistent
605 with Experiment 1. The `XGBRegressor` also showed larger variance in predictions
606 as compared to `ChemProp`, which might indicate that `ChemProp` is more stable and
607 consistent in its predictive performance.

608 These case studies were not by all means complete or exhaustive, but they illustrate
609 how QSPRpred could be used to run various experiments for validation of novel QSPR
610 methodologies on a variety of problems. Moreover, the available code in the `qsp-`
611 `bench` repository can be quickly derived from to create other benchmarks. We think
612 that this is an especially interesting prospect for developers of new models who can
613 validate their approach within QSPRpred, but the simple act of integrating it into the
614 benchmarking workflow also makes it readily available for deployment by others.

615 4 Conclusions

616 QSPRpred is a new and versatile open-source package for QSPR modelling. QSPRpred
617 addresses a number of issues in the QSPR modelling field, including a need for tools
618 that facilitate easy comparison and validation of an ever-growing number of different
619 QSPR workflows. QSPRpred enables this by its modular Python API that simpli-
620 fies the implementation of a plethora of QSPR modelling tasks, which can be easily
621 tied together in its benchmarking workflow. Furthermore, reproducibility of results is
622 ensured through consistent serialisation of models and data in human-readable for-
623 mat. Inclusion of data pre-processing and featurization steps with the models enables
624 direct application of trained models to new compounds using only a SMILES string.
625 Moreover, to our knowledge this is the first QSPR modelling tool to support pro-
626 teochemometric modelling in Python. QSPRpred is also integrated with a number of
627 other packages developed in the Leiden Computational Drug Discovery group, notably
628 DrugEx [91] for *de novo* drug design and Papyrus [54] for collection of bio-activity
629 data. Extensive documentation and comprehensive tutorials are available.

630 In the future, we will continue to develop QSPRpred and extend its capabili-
631 ties. Most notably we intend to further extend the range of available descriptors (i.e.
632 Molfeat from datamol.io [92]) and models (i.e. by adding support for ensemble mod-
633 elling and a wider range of neural network architectures with the help of skorch [93], a
634 scikit-learn [71] compatible neural network library that wraps PyTorch [73]). Further-
635 more, we will integrate QSPRpred within GenUI [94], which provides an accessible
636 user interface for cheminformatics, QSAR modelling and AI-based molecular genera-
637 tion provided by the associated DrugEx framework [91]. Furthermore, continued efforts
638 are needed to teach and adhere to high standards for FAIR [80] research practices
639 and we hope that QSPRpred will prove to be a helpful tool to assist researchers in
640 reproducible and transferable QSPR modelling. Therefore, with an assortment of sup-
641 porting and depending packages already developed or under active development and

642 with active user base from both Leiden University and UCT Prague, we aim to keep
643 QSPRpred up to date and help researchers to perform experiments faster and in a
644 more standardized and reproducible manner as the field of QSPR modelling evolves.

645 **Supplementary information.** The source code is available from <https://github.com/CDDLeiden/QSPRpred>. The described version of the software (v3.0.1) is also
646 archived on zenodo at <https://doi.org/10.5281/zenodo.10720563>.
647

648 **Availability and requirements.**

- 649 • **Project name:** QSPRpred
- 650 • **Project home page:** <https://github.com/CDDLeiden/QSPRpred>
- 651 • **Case Study Code:** <https://github.com/CDDLeiden/qsp-bench>
- 652 • **Operating system(s)** Full support for Linux. QSPRpred is supported on Windows
653 apart from the PCM modelling, which relies on Clustal Omega [66] or MAFFT [67].
654 These require manually installation on Window. We are currently working on full
655 support for MacOS.
- 656 • **Programming Language:** Python
- 657 • **Other requirements:** Python 3.10, see [https://github.com/CDDLeiden/
658 QSPRpred/blob/main/pyproject.toml](https://github.com/CDDLeiden/QSPRpred/blob/main/pyproject.toml) for full list of requirements.
- 659 • **License:** MIT

660 **Declarations**

- 661 • Availability of data and materials Not applicable
- 662 • Competing interests Not applicable
- 663 • Funding
 - 664 – M.Š. was supported by Czech Science Foundation Grant No. 22-17367O and by
665 the Ministry of Education, Youth and Sports of the Czech Republic (project
666 number LM2023052).

667 – D.A.A. was supported by funding from the European Union’s Horizon 2020
668 research and innovation programme under the Marie Skłodowska Curie grant
669 agreement No 955879..

670 – Research reported in this publication was supported by Oncode Accelerator, a
671 Dutch National Growth Fund project under grant number NGFOP2201.

672 • Authors’ contributions **H.W.vdM.**: Methodology, Software, Writing – Original
673 Draft, Writing - Review & Editing, Visualization, Project administration **M.Š.**:
674 Methodology, Software, Benchmarking - Data Curation, code & execution, Writing
675 – Original Draft, Writing - Review & Editing, Visualization, Project administra-
676 tion **D.A.A.**: Methodology, Software, Writing – Original Draft **S.L.**: Methodology,
677 Software, Writing – Original Draft **L.S.**: Methodology, Software, Writing – Orig-
678 inal Draft, Writing - Review & Editing **M.J.**: Methodology, Software, Writing –
679 Original Draft **O.J.M.B.**: Methodology, Software, Writing – Original Draft, Writ-
680 ing - Review & Editing **M.G.G.**: Methodology, Software, Writing – Original Draft,
681 Writing - Review & Editing **R.L.vdB.**: Methodology, Software, Writing – Orig-
682 inal Draft, Benchmarking - Data Curation & code **A.B.**: Methodology, contributed
683 expertise for model integration **J.G.C.vH.**: Resources, Writing - Review & Editing,
684 Supervision, Funding acquisition **P.H.vdG.**: Resources, Writing - Review & Edit-
685 ing, Supervision, Funding acquisition **G.J.P.vW.**: Resources, Writing - Review &
686 Editing, Supervision, Funding acquisition

687 • Acknowledgements We would like to thank Dr. Wim Dehaen for contributing some
688 of his code and knowledge to the project.

689 • Author’s information **H.W.vdM.** ([orcid:0000-0002-9718-7806](https://orcid.org/0000-0002-9718-7806)); **M.Š.**
690 ([orcid:0000-0002-8771-1731](https://orcid.org/0000-0002-8771-1731)); **D.A.A.** [orcid:0000-0002-5104-1959](https://orcid.org/0000-0002-5104-1959); **S.L.** ([orcid:0000-0001-9387-1427](https://orcid.org/0000-0001-9387-1427)); **L.S.** ([orcid:0000-0001-9879-1004](https://orcid.org/0000-0001-9879-1004)); **M.J.** [orcid:0009-0003-2083-0159](https://orcid.org/0009-0003-2083-0159);
691 **O.J.M.B.** ([orcid:0000-0002-7554-9220](https://orcid.org/0000-0002-7554-9220)); **M.G.G.** ([orcid:0000-0003-1568-0881](https://orcid.org/0000-0003-1568-0881));
692 **R.L.vdB.** ([orcid:0009-0008-5661-1157](https://orcid.org/0009-0008-5661-1157)); **A.B.** ([orcid:0000-0002-0058-3678](https://orcid.org/0000-0002-0058-3678));
693

694 J.G.C.vH. ([orcid:0000-0002-1664-7314](https://orcid.org/0000-0002-1664-7314)); P.H.vdG. ([orcid.org:0000-0003-1314-3484](https://orcid.org/0000-0003-1314-3484));
695 G.J.P.vW. ([orcid:0000-0003-0717-1817](https://orcid.org/0000-0003-0717-1817))

696 References

- 697 [1] Muratov, E.N., Bajorath, J., Sheridan, R.P., Tetko, I.V., Filimonov, D., Poroikov,
698 V., Oprea, T.I., Baskin, I.I., Varnek, A., Roitberg, A., Isayev, O., Curtalolo,
699 S., Fourches, D., Cohen, Y., Aspuru-Guzik, A., Winkler, D.A., Agrafiotis, D.,
700 Cherkasov, A., Tropsha, A.: QSAR without borders. *Chemical Society Reviews*
701 **49**(11), 3525–3564 (2020) <https://doi.org/10.1039/D0CS00098A>
- 702 [2] Hansch, C., Fujita, T.: p - σ - π Analysis. A Method for the Correlation of Biologi-
703 cal Activity and Chemical Structure. *Journal of the American Chemical Society*
704 **86**(8), 1616–1626 (1964) <https://doi.org/10.1021/ja01062a035>
- 705 [3] Jiménez-Luna, J., Grisoni, F., Weskamp, N., Schneider, G.: Artificial intelligence
706 in drug discovery: recent advances and future perspectives. *Expert Opinion on*
707 *Drug Discovery* **16**(9), 949–959 (2021) [https://doi.org/10.1080/17460441.2021.](https://doi.org/10.1080/17460441.2021.1909567)
708 [1909567](https://doi.org/10.1080/17460441.2021.1909567)
- 709 [4] Alves, V.M., Bobrowski, T., Melo-Filho, C.C., Korn, D., Auerbach, S., Schmitt,
710 C., Muratov, E.N., Tropsha, A.: QSAR Modeling of SARS-CoV Mpro Inhibitors
711 Identifies Sufugolix, Cenicriviroc, Proglumetacin, and other Drugs as Candidates
712 for Repurposing against SARS-CoV-2. *Molecular Informatics* **40**(1), 2000113
713 (2021) <https://doi.org/10.1002/minf.202000113>
- 714 [5] Tejera, E., Munteanu, C.R., López-Cortés, A., Cabrera-Andrade, A., Pérez-
715 Castillo, Y.: Drugs Repurposing Using QSAR, Docking and Molecular Dynamics
716 for Possible Inhibitors of the SARS-CoV-2 Mpro Protease. *Molecules* **25**(21), 5172
717 (2020) <https://doi.org/10.3390/molecules25215172>

- 718 [6] Väitalo, P.A.J., Griffioen, K., Rizk, M.L., Visser, S.A.G., Danhof, M., Rao,
719 G., Graaf, P.H., Hasselt, J.G.C.: Structure-Based Prediction of Anti-infective
720 Drug Concentrations in the Human Lung Epithelial Lining Fluid. *Pharmaceutical
721 Research* **33**(4), 856–867 (2016) <https://doi.org/10.1007/s11095-015-1832-x>
- 722 [7] Paduszyński, K., Klebowski, K., Królikowska, M.: Predicting melting point of
723 ionic liquids using QSPR approach: Literature review and new models. *Journal of
724 Molecular Liquids* **344**, 117631 (2021) [https://doi.org/10.1016/j.molliq.2021.
725 117631](https://doi.org/10.1016/j.molliq.2021.117631)
- 726 [8] Pillai, N., Dasgupta, A., Sudsakorn, S., Fretland, J., Mavroudis, P.D.: Machine
727 Learning guided early drug discovery of small molecules. *Drug Discovery Today*
728 **27**(8), 2209–2215 (2022) <https://doi.org/10.1016/j.drudis.2022.03.017>
- 729 [9] Mendez, D., Gaulton, A., Bento, A.P., Chambers, J., De Veij, M., Félix, E.,
730 Magariños, M., Mosquera, J., Mutowo, P., Nowotka, M., Gordillo-Marañón, M.,
731 Hunter, F., Junco, L., Mugumbate, G., Rodriguez-Lopez, M., Atkinson, F., Bosc,
732 N., Radoux, C., Segura-Cabrera, A., Hersey, A., Leach, A.: ChEMBL: towards
733 direct deposition of bioassay data. *Nucleic Acids Research* **47**(D1), 930–940 (2019)
734 <https://doi.org/10.1093/nar/gky1075>
- 735 [10] Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker,
736 B.A., Thiessen, P.A., Yu, B., Zaslavsky, L., Zhang, J., Bolton, E.E.: PubChem
737 2023 update. *Nucleic Acids Research* **51**(D1), 1373–1380 (2023) [https://doi.org/
738 10.1093/nar/gkac956](https://doi.org/10.1093/nar/gkac956)
- 739 [11] Cherkasov, A., Muratov, E.N., Fourches, D., Varnek, A., Baskin, I.I., Cronin, M.,
740 Dearden, J., Gramatica, P., Martin, Y.C., Todeschini, R., Consonni, V., Kuz'min,
741 V.E., Cramer, R., Benigni, R., Yang, C., Rathman, J., Terfloth, L., Gasteiger,
742 J., Richard, A., Tropsha, A.: QSAR Modeling: Where Have You Been? Where

- 743 Are You Going To? *Journal of Medicinal Chemistry* **57**(12), 4977–5010 (2014)
744 <https://doi.org/10.1021/jm4004285>
- 745 [12] Tropsha, A., Isayev, O., Varnek, A., Schneider, G., Cherkasov, A.: Integrating
746 QSAR modelling and deep learning in drug discovery: the emergence of deep
747 QSAR. *Nature Reviews Drug Discovery*, 1–15 (2023) [https://doi.org/10.1038/
748 s41573-023-00832-0](https://doi.org/10.1038/s41573-023-00832-0)
- 749 [13] Bongers, B.J., IJzerman, A.P., Van Westen, G.J.P.: Proteochemometrics: recent
750 developments in bioactivity and selectivity modeling. *Drug Discovery Today:
751 Technologies* **32**, 89–98 (2019) <https://doi.org/10.1016/j.ddtec.2020.08.003>
- 752 [14] Cortés-Ciriano, I., Ul Ain, Q., Subramanian, V., B. Lenselink, E., Méndez-
753 Lucio, O., P. IJzerman, A., Wohlfahrt, G., Prusis, P., E. Malliavin, T., Westen,
754 G.J.P.v., Bender, A.: Polypharmacology modelling using proteochemometrics
755 (PCM): recent methodological developments, applications to target families, and
756 future prospects. *MedChemComm* **6**(1), 24–50 (2015) [https://doi.org/10.1039/
757 C4MD00216D](https://doi.org/10.1039/C4MD00216D)
- 758 [15] Burggraaff, L., Lenselink, E.B., Jespers, W., Engelen, J., Bongers, B.J., Gorosti-
759 ola González, M., Liu, R., Hoos, H.H., Vlijmen, H.W.T., IJzerman, A.P., Westen,
760 G.J.P.: Successive statistical and structure-based modeling to identify chemically
761 novel kinase inhibitors. *Journal of Chemical Information and Modeling* **60**(9),
762 4283–4295 (2020) <https://doi.org/10.1021/acs.jcim.9b01204>
- 763 [16] Gorostiola González, M., Broek, R.L., Braun, T.G.M., Chatzopoulou, M., Jespers,
764 W., IJzerman, A.P., Heitman, L.H., Westen, G.J.P.: 3DDPDs: describing protein
765 dynamics for proteochemometric bioactivity prediction. A case for (mutant) G
766 protein-coupled receptors. *Journal of Cheminformatics* **15**(1), 74 (2023) [https:
767 //doi.org/10.1186/s13321-023-00745-5](https://doi.org/10.1186/s13321-023-00745-5)

- 768 [17] Born, J., Huynh, T., Stroobants, A., Cornell, W.D., Manica, M.: Active Site
769 Sequence Representations of Human Kinases Outperform Full Sequence Repre-
770 sentations for Affinity Prediction and Inhibitor Generation: 3D Effects in a 1D
771 Model. *Journal of Chemical Information and Modeling* **62**(2), 240–257 (2022)
772 <https://doi.org/10.1021/acs.jcim.1c00889>
- 773 [18] Lenselink, E.B., Dijke, N., Bongers, B., Papadatos, G., Vlijmen, H.W.T., Kowal-
774 czyk, W., IJzerman, A.P., Westen, G.J.P.: Beyond the hype: deep neural
775 networks outperform established methods using a ChEMBL bioactivity bench-
776 mark set. *Journal of Cheminformatics* **9**, 45 (2017) [https://doi.org/10.1186/
777 s13321-017-0232-0](https://doi.org/10.1186/s13321-017-0232-0)
- 778 [19] Playe, B., Stoven, V.: Evaluation of deep and shallow learning methods in
779 chemogenomics for the prediction of drugs specificity. *Journal of Cheminformatics*
780 **12**(1), 11 (2020) <https://doi.org/10.1186/s13321-020-0413-0>
- 781 [20] Atas Guvenilir, H., Doğan, T.: How to approach machine learning-based predic-
782 tion of drug/compound–target interactions. *Journal of Cheminformatics* **15**(1),
783 16 (2023) <https://doi.org/10.1186/s13321-023-00689-w>
- 784 [21] Lopez-del Rio, A., Picart-Armada, S., Perera-Lluna, A.: Balancing Data on Deep
785 Learning-Based Proteochemometric Activity Classification. *Journal of Chemical*
786 *Information and Modeling* **61**(4), 1657–1669 (2021) [https://doi.org/10.1021/acs.
787 jcim.1c00086](https://doi.org/10.1021/acs.jcim.1c00086)
- 788 [22] Luukkonen, S., Meijer, E., Tricarico, G.A., Hofmans, J., Stouten, P.F.W., Westen,
789 G.J.P., Lenselink, E.B.: Large-Scale Modeling of Sparse Protein Kinase Activity
790 Data. *Journal of Chemical Information and Modeling* **63**(12), 3688–3696 (2023)
791 <https://doi.org/10.1021/acs.jcim.3c00132>

- 792 [23] Kanev, G.K., Zhang, Y., Kooistra, A.J., Bender, A., Leurs, R., Bailey, D.,
793 Würdinger, T., Graaf, C.d., Esch, I.J.P.d., Westerman, B.A.: Predicting the tar-
794 get landscape of kinase inhibitors using 3D convolutional neural networks. PLOS
795 Computational Biology **19**(9), 1011301 (2023) <https://doi.org/10.1371/journal.pcbi.1011301>
796
- 797 [24] Janela, T., Bajorath, J.: Rationalizing general limitations in assessing and com-
798 paring methods for compound potency prediction. Scientific Reports **13**(1), 17816
799 (2023) <https://doi.org/10.1038/s41598-023-45086-3>
- 800 [25] McElfresh, D., Khandagale, S., Valverde, J., C, V.P., Feuer, B., Hegde, C.,
801 Ramakrishnan, G., Goldblum, M., White, C.: When Do Neural Nets Outperform
802 Boosted Trees on Tabular Data? arXiv (2023). <https://doi.org/10.48550/arXiv.2305.02997>
803
- 804 [26] Grinsztajn, L., Oyallon, E., Varoquaux, G.: Why do tree-based models still out-
805 perform deep learning on tabular data? arXiv (2022). <https://doi.org/10.48550/arXiv.2207.08815>
806
- 807 [27] Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A.,
808 Hopper, T., Kelley, B., Mathea, M., Palmer, A., Settels, V., Jaakkola, T., Jensen,
809 K., Barzilay, R.: Analyzing Learned Molecular Representations for Property Pre-
810 diction. Journal of Chemical Information and Modeling **59**(8), 3370–3388 (2019)
811 <https://doi.org/10.1021/acs.jcim.9b00237>
- 812 [28] Deng, J., Yang, Z., Wang, H., Ojima, I., Samaras, D., Wang, F.: A systematic
813 study of key elements underlying molecular property prediction. Nature Commu-
814 nications **14**(1), 6395 (2023) <https://doi.org/10.1038/s41467-023-41948-6>
- 815 [29] Wu, Z., Ramsundar, B., Feinberg, E.N., Gomes, J., Geniesse, C., Pappu,

- 816 A.S., Leswing, K., Pande, V.: MoleculeNet: a benchmark for molecular
817 machine learning. *Chemical Science* **9**(2), 513–530 (2018) <https://doi.org/10.1039/C7SC02664A>
818
- 819 [30] Tossou, P., Wognum, C., Craig, M., Mary, H., Noutahi, E.: Real-World Molecular
820 Out-Of-Distribution: Specification and Investigation. *ChemRxiv* (2023). <https://doi.org/10.26434/chemrxiv-2023-q11q4-v2>
821
- 822 [31] Steshin, S.: Lo-Hi: Practical ML Drug Discovery Benchmark. *arXiv* (2023). <http://arxiv.org/abs/2310.06399> Accessed 2023-12-11
823
- 824 [32] Tilborg, D., Alenicheva, A., Grisoni, F.: Exposing the Limitations of Molecular
825 Machine Learning with Activity Cliffs. *Journal of Chemical Information and
826 Modeling* **62**(23), 5938–5951 (2022) <https://doi.org/10.1021/acs.jcim.2c01073>
- 827 [33] Boldini, D., Grisoni, F., Kuhn, D., Friedrich, L., Sieber, S.A.: Practical guidelines
828 for the use of gradient boosting for molecular property prediction. *Journal of
829 Cheminformatics* **15**(1), 73 (2023) <https://doi.org/10.1186/s13321-023-00743-7>
- 830 [34] Gramatica, P., Sangion, A.: A Historical Excursus on the Statistical Validation
831 Parameters for QSAR Models: A Clarification Concerning Metrics and Terminology.
832 *Journal of Chemical Information and Modeling* **56**(6), 1127–1131 (2016)
833 <https://doi.org/10.1021/acs.jcim.6b00088>
- 834 [35] Walters, P.: Comparing Classification Models - You're Probably Doing
835 It Wrong (2023). [https://practicalcheminformatics.blogspot.com/2023/11/
836 comparing-classification-models-youre.html](https://practicalcheminformatics.blogspot.com/2023/11/comparing-classification-models-youre.html) Accessed 2023-12-07
- 837 [36] Walters, P.: We Need Better Benchmarks for Machine Learning in
838 Drug Discovery (2023). [http://practicalcheminformatics.blogspot.com/2023/08/
839 we-need-better-benchmarks-for-machine.html](http://practicalcheminformatics.blogspot.com/2023/08/we-need-better-benchmarks-for-machine.html) Accessed 2023-12-08

- 840 [37] Landrum, G.A., Riniker, S.: Combining IC50 or Ki Values from Different Sources
841 Is a Source of Significant Noise. *Journal of Chemical Information and Modeling*
842 (2024) <https://doi.org/10.1021/acs.jcim.4c00049>
- 843 [38] Sciences, E.N.A.o., Affairs, P.a.G., Science, E., Information, B.o.R.D., , Sci-
844 ences, D.o.E.a.P., Statistics, C.o.A.a.T., Analytics, B.o.M.S., , Studies, D.o.E.a.L.,
845 Board, N.a.R.S., Education, D.o.B.a.S.S., , Statistics, C.o.N., Behavioral, C.,
846 Science, C.o.R.a.R.i.: Understanding Reproducibility and Replicability. In: Repro-
847 ducibility and Replicability in Science. National Academies Press (US), Washing-
848 ton (DC) (2019). <https://www.ncbi.nlm.nih.gov/books/NBK547546/>
- 849 [39] Hutson, M.: Artificial intelligence faces reproducibility crisis. *Science* (New York,
850 N.Y.) **359**(6377), 725–726 (2018) <https://doi.org/10.1126/science.359.6377.725>
- 851 [40] Schaduangrat, N., Lampa, S., Simeon, S., Gleeson, M.P., Spjuth, O., Nantase-
852 namat, C.: Towards reproducible computational drug discovery. *J. Cheminform.*
853 **12**(1), 9 (2020) <https://doi.org/10.1186/s13321-020-0408-x>
- 854 [41] Clark, R.D.: A path to next-generation reproducibility in cheminformatics. *J.*
855 *Cheminform.* **11**(1), 62 (2019) <https://doi.org/10.1186/s13321-019-0385-0>
- 856 [42] Hoyt, C.T., Zdrzil, B., Guha, R., Jeliaskova, N., Martinez-Mayorga, K., Nit-
857 tinger, E.: Improving reproducibility and reusability in the Journal of Chemin-
858 formatics. *Journal of Cheminformatics* **15**(1), 62 (2023) <https://doi.org/10.1186/s13321-023-00730-y>
- 860 [43] Patel, M., Chilton, M.L., Sartini, A., Gibson, L., Barber, C., Covey-Crump, L.,
861 Przybylak, K.R., Cronin, M.T.D., Madden, J.C.: Assessment and Reproducibil-
862 ity of Quantitative Structure–Activity Relationship Models by the Nonexpert.
863 *Journal of Chemical Information and Modeling* **58**(3), 673–682 (2018) <https://doi.org/10.1021/acs.jcim.8b00049>

- 864 [//doi.org/10.1021/acs.jcim.7b00523](https://doi.org/10.1021/acs.jcim.7b00523)
- 865 [44] Berthold, M.R., Cebon, N., Dill, F., Gabriel, T.R., Kötter, T., Meinl, T., Ohl, P.,
866 Sieb, C., Thiel, K., Wiswedel, B.: KNIME: The Konstanz Information Miner. In:
867 Preisach, C., Burkhardt, H., Schmidt-Thieme, L., Decker, R. (eds.) Data Analysis,
868 Machine Learning and Applications. Studies in Classification, Data Analysis, and
869 Knowledge Organization, pp. 319–326. Springer, Berlin, Heidelberg (2008). https://doi.org/10.1007/978-3-540-78246-9_38
870 [//doi.org/10.1007/978-3-540-78246-9_38](https://doi.org/10.1007/978-3-540-78246-9_38)
- 871 [45] KNIME: Create a New Python based KNIME Extension (2024). https://docs.knime.com/latest/pure_python_node_extensions_guide/index.html#introduction
872 https://docs.knime.com/latest/pure_python_node_extensions_guide/index.html#introduction
873 Accessed 2024-02-26
- 874 [46] Ramsundar, B., Eastman, P., Walters, P., Pande, V.: Deep Learning for the Life
875 Sciences: Applying Deep Learning to Genomics, Microscopy, Drug Discovery, and
876 More, 1st edition edn. O’Reilly Media, Sebastopol, CA (2019)
- 877 [47] Minnich, A.J., McLoughlin, K., Tse, M., Deng, J., Weber, A., Murad, N., Madej,
878 B.D., Ramsundar, B., Rush, T., Calad-Thomson, S., Brase, J., Allen, J.E.: AMPL:
879 A Data-Driven Modeling Pipeline for Drug Discovery. *Journal of Chemical Infor-*
880 *mation and Modeling* **60**(4), 1955–1968 (2020) [https://doi.org/10.1021/acs.jcim.](https://doi.org/10.1021/acs.jcim.9b01053)
881 [9b01053](https://doi.org/10.1021/acs.jcim.9b01053)
- 882 [48] Turon, G., Hlozek, J., Woodland, J.G., Chibale, K., Duran-Frigola, M.: First fully-
883 automated AI/ML virtual screening cascade implemented at a drug discovery
884 centre in Africa. *bioRxiv* (2022). <https://doi.org/10.1101/2022.12.13.520154>
- 885 [49] Lanini, J., Santarossa, G., Sirockin, F., Lewis, R., Fechner, N., Misztela, H., Lewis,
886 S., Maziarz, K., Stanley, M., Segler, M., Stiefl, N., Schneider, N.: PREFER: A New
887 Predictive Modeling Framework for Molecular Discovery. *Journal of Chemical*

- 888 Information and Modeling (2023) <https://doi.org/10.1021/acs.jcim.3c00523>
- 889 [50] Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., Hutter,
890 F.: Efficient and Robust Automated Machine Learning. In: Advances in Neural
891 Information Processing Systems, vol. 28. Curran Associates, Inc., Red Hook, New
892 York, USA (2015)
- 893 [51] Mervin, L., Voronov, A., Kabeshov, M., Engkvist, O.: Qptuna: an automated
894 QSAR modelling platform for molecular property prediction in drug design.
895 ChemRxiv (2024). <https://doi.org/10.26434/chemrxiv-2024-2r1k7>
- 896 [52] Murrell, D.S., Cortes-Ciriano, I., Westen, G.J.P., Stott, I.P., Bender, A., Malli-
897 avin, T.E., Glen, R.C.: Chemically Aware Model Builder (camb): an R package
898 for property and bioactivity modelling of small molecules. Journal of Cheminform-
899 matics **7**(1), 45 (2015) <https://doi.org/10.1186/s13321-015-0086-2>
- 900 [53] Murrell, D.S., Cortes-Ciriano, I., Westen, G.J.P., Stott, I.P., Bender, A., Malli-
901 avin, T.E., Glen, R.C.: cambDI/camb. cambDI (2021). [https://github.com/](https://github.com/cambDI/camb)
902 [cambDI/camb](https://github.com/cambDI/camb) Accessed 2024-02-26
- 903 [54] Béquignon, O.J.M., Bongers, B.J., Jespers, W., IJzerman, A.P., Water, B.,
904 Westen, G.J.P.: Papyrus: a large-scale curated dataset aimed at bioactivity pre-
905 dictions. Journal of Cheminformatics **15**(1), 3 (2023) [https://doi.org/10.1186/](https://doi.org/10.1186/s13321-022-00672-x)
906 [s13321-022-00672-x](https://doi.org/10.1186/s13321-022-00672-x)
- 907 [55] McKinney, W.: Data Structures for Statistical Computing in Python. Proceedings
908 of the 9th Python in Science Conference, 56–61 (2010) [https://doi.org/10.25080/](https://doi.org/10.25080/Majora-92bf1922-00a)
909 [Majora-92bf1922-00a](https://doi.org/10.25080/Majora-92bf1922-00a)
- 910 [56] Šícho, M.: martin-sicho/papyrus-scaffold-visualizer (2023). [https://github.com/](https://github.com/martin-sicho/papyrus-scaffold-visualizer)
911 [martin-sicho/papyrus-scaffold-visualizer](https://github.com/martin-sicho/papyrus-scaffold-visualizer) Accessed 2023-12-12

- 912 [57] RDKit: Open-source cheminformatics. (2024). <https://www.rdkit.org>
- 913 [58] Moriwaki, H., Tian, Y.-S., Kawashita, N., Takagi, T.: Mordred: a molecular
914 descriptor calculator. *Journal of Cheminformatics* **10**(1), 4 (2018) [https://doi.
915 org/10.1186/s13321-018-0258-y](https://doi.org/10.1186/s13321-018-0258-y)
- 916 [59] Hong, H., Xie, Q., Ge, W., Qian, F., Fang, H., Shi, L., Su, Z., Perkins, R., Tong,
917 W.: Mold2, Molecular Descriptors from 2D Structures for Chemoinformatics and
918 Toxicoinformatics. *Journal of Chemical Information and Modeling* **48**(7), 1337–
919 1344 (2008) <https://doi.org/10.1021/ci800038f>
- 920 [60] Yap, C.W.: PaDEL-descriptor: An open source software to calculate molecular
921 descriptors and fingerprints. *Journal of Computational Chemistry* **32**(7), 1466–
922 1474 (2011) <https://doi.org/10.1002/jcc.21707>
- 923 [61] Durant, J.L., Leland, B.A., Henry, D.R., Nourse, J.G.: Reoptimization of MDL
924 Keys for Use in Drug Discovery. *Journal of Chemical Information and Computer
925 Sciences* **42**(6), 1273–1280 (2002) <https://doi.org/10.1021/ci010132r>
- 926 [62] Hellberg, S., Sjoestroem, M., Skagerberg, B., Wold, S.: Peptide quantitative
927 structure-activity relationships, a multivariate approach. *Journal of Medicinal
928 Chemistry* **30**(7), 1126–1135 (1987) <https://doi.org/10.1021/jm00390a003>
- 929 [63] Georgiev, A.G.: Interpretable Numerical Descriptors of Amino Acid Space. *Jour-
930 nal of Computational Biology* **16**(5), 703–723 (2009) [https://doi.org/10.1089/
931 cmb.2008.0173](https://doi.org/10.1089/cmb.2008.0173)
- 932 [64] Mei, H., Liao, Z.H., Zhou, Y., Li, S.Z.: A new set of amino acid descriptors
933 and its application in peptide QSARs. *Peptide Science* **80**(6), 775–786 (2005)
934 <https://doi.org/10.1002/bip.20296>

- 935 [65] Béquignon, Olivier J. M.: ProDEC (2023). [https://github.com/OlivierBeq/](https://github.com/OlivierBeq/ProDEC/tree/master)
936 [ProDEC/tree/master](https://github.com/OlivierBeq/ProDEC/tree/master) Accessed 2023-12-12
- 937 [66] Sievers, F., Wilm, A., Dineen, D., Gibson, T.J., Karplus, K., Li, W., Lopez, R.,
938 McWilliam, H., Remmert, M., Söding, J., Thompson, J.D., Higgins, D.G.: Fast,
939 scalable generation of high-quality protein multiple sequence alignments using
940 Clustal Omega. *Molecular Systems Biology* **7**(1), 539 (2011) [https://doi.org/10.](https://doi.org/10.1038/msb.2011.75)
941 [1038/msb.2011.75](https://doi.org/10.1038/msb.2011.75)
- 942 [67] Katoh, K., Standley, D.M.: MAFFT Multiple Sequence Alignment Software
943 Version 7: Improvements in Performance and Usability. *Molecular Biology and*
944 *Evolution* **30**(4), 772–780 (2013) <https://doi.org/10.1093/molbev/mst010>
- 945 [68] Homola, D.: boruta_py. scikit-learn-contrib (2023). [https://github.com/](https://github.com/scikit-learn-contrib/boruta_py)
946 [scikit-learn-contrib/boruta_py](https://github.com/scikit-learn-contrib/boruta_py) Accessed 2023-11-28
- 947 [69] Kursá, M.B., Rudnicki, W.R.: Feature Selection with the Boruta Package. *Journal*
948 *of Statistical Software* **36**, 1–13 (2010) <https://doi.org/10.18637/jss.v036.i11>
- 949 [70] Rácz, A., Bajusz, D., Héberger, K.: Effect of Dataset Size and Train/Test Split
950 Ratios in QSAR/QSPR Multiclass Classification. *Molecules* **26**(4), 1111 (2021)
951 <https://doi.org/10.3390/molecules26041111>
- 952 [71] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
953 Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A.,
954 Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine
955 Learning in Python. *Journal of Machine Learning Research* **12**(85), 2825–2830
956 (2011)
- 957 [72] Tricarico, G.A., Hofmans, J., Lenselink, E.B., López-Ramos, M., Dréanic, M.-
958 P., Stouten, P.F.W.: Construction of balanced, chemically dissimilar training,

- 959 validation and test sets for machine learning on molecular datasets (2022) <https://doi.org/10.26434/chemrxiv-2022-m8l33-v2>
960
- 961 [73] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T.,
962 Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito,
963 Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chin-
964 tala, S.: PyTorch: An Imperative Style, High-Performance Deep Learning Library.
965 arXiv (2019). <http://arxiv.org/abs/1912.01703>
- 966 [74] Iosipoi, L., Vakhrushev, A.: SketchBoost: Fast Gradient Boosted Decision Tree
967 for Multioutput Problems. arXiv (2022). <https://doi.org/10.48550/arXiv.2211.12858>
968
- 969 [75] Heid, E., Greenman, K.P., Chung, Y., Li, S.-C., Graff, D.E., Vermeire, F.H.,
970 Wu, H., Green, W.H., McGill, C.J.: Chemprop: A Machine Learning Package
971 for Chemical Property Prediction. ChemRxiv (2023). <https://doi.org/10.26434/chemrxiv-2023-3zcf1-v3>
972
- 973 [76] QSPRpred: Documentation (2024). <https://cddleiden.github.io/QSPRpred/docs/index.html> Accessed 2023-12-18
974
- 975 [77] Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M.: Optuna: A Next-
976 generation Hyperparameter Optimization Framework. In: Proceedings of the 25th
977 ACM SIGKDD International Conference on Knowledge Discovery & Data Min-
978 ing, pp. 2623–2631. ACM, Anchorage AK USA (2019). <https://doi.org/10.1145/3292500.3330701>
979
- 980 [78] Béquignon, O.J.M.: MLChemAD (2023). <https://github.com/OlivierBeq/MLChemAD> Accessed 2024-02-25
981
- 982 [79] McCorkindale, W., Elijošius, R.: molplotly (2022). <https://github.com/wjm41/>

983 [molplotly](#)

- 984 [80] Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M.,
985 Baak, A., Blomberg, N., Boiten, J.-W., Silva Santos, L.B., Bourne, P.E., Bouw-
986 man, J., Brookes, A.J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S.,
987 Evelo, C.T., Finkers, R., Gonzalez-Beltran, A., Gray, A.J.G., Groth, P., Goble,
988 C., Grethe, J.S., Heringa, J., Hoen, P.A.C., Hooft, R., Kuhn, T., Kok, R., Kok,
989 J., Lusher, S.J., Martone, M.E., Mons, A., Packer, A.L., Persson, B., Rocca-
990 Serra, P., Roos, M., Schaik, R., Sansone, S.-A., Schultes, E., Sengstag, T., Slater,
991 T., Strawn, G., Swertz, M.A., Thompson, M., Lei, J., Mulligen, E., Velterop, J.,
992 Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., Mons, B.: The FAIR
993 Guiding Principles for scientific data management and stewardship. *Scientific*
994 *Data* **3**(1), 160018 (2016) <https://doi.org/10.1038/sdata.2016.18>
- 995 [81] Béquignon, O.J.M.: ml2json (2023). <https://github.com/OlivierBeq/ml2json>
996 Accessed 2023-12-12
- 997 [82] Aguilar, D.: jsonpickle (2023). <https://github.com/jsonpickle/jsonpickle> Accessed
998 2023-12-12
- 999 [83] QSPRpred: Style guide (2024). [https://github.com/CDDLeiden/QSPRpred/
1000 blob/main/docs/style_guide.py](https://github.com/CDDLeiden/QSPRpred/blob/main/docs/style_guide.py) Accessed 2023-12-18
- 1001 [84] Biewald, L.: Experiment Tracking with Weights and Biases (2020). [https://www.
1002 wandb.com/](https://www.wandb.com/)
- 1003 [85] Zhao, Z., Qin, J., Gou, Z., Zhang, Y., Yang, Y.: Multi-task learning models for pre-
1004 dicting active compounds. *Journal of Biomedical Informatics* **108**, 103484 (2020)
1005 <https://doi.org/10.1016/j.jbi.2020.103484>
- 1006 [86] Sosnin, S., Vashurina, M., Withnall, M., Karpov, P., Fedorov, M., Tetko, I.V.:

- 1007 A Survey of Multi-task Learning Methods in Chemoinformatics. *Molecular*
1008 *Informatics* **38**(4), 1800108 (2019) <https://doi.org/10.1002/minf.201800108>
- 1009 [87] Chen, T., Guestrin, C.: XGBoost: A Scalable Tree Boosting System. In: Pro-
1010 ceedings of the 22nd ACM SIGKDD International Conference on Knowledge
1011 Discovery And Data Mining, pp. 785–794. ACM, San Francisco California
1012 USA (2016). <https://doi.org/10.1145/2939672.2939785> . [https://dl.acm.org/doi/](https://dl.acm.org/doi/10.1145/2939672.2939785)
1013 [10.1145/2939672.2939785](https://dl.acm.org/doi/10.1145/2939672.2939785)
- 1014 [88] Vandenhende, S., Georgoulis, S., Van Gansbeke, W., Proesmans, M., Dai, D.,
1015 Van Gool, L.: Multi-Task Learning for Dense Prediction Tasks: A Survey. *IEEE*
1016 *Transactions on Pattern Analysis and Machine Intelligence* **44**(7), 3614–3633
1017 (2022) <https://doi.org/10.1109/TPAMI.2021.3054719> . Conference Name: IEEE
1018 *Transactions on Pattern Analysis and Machine Intelligence*
- 1019 [89] León, A., Chen, B., Gillet, V.J.: Effect of missing data on multitask prediction
1020 methods. *Journal of Cheminformatics* **10**(1), 26 (2018) [https://doi.org/10.1186/](https://doi.org/10.1186/s13321-018-0281-z)
1021 [s13321-018-0281-z](https://doi.org/10.1186/s13321-018-0281-z)
- 1022 [90] Lombardo, F., Bentzien, J., Berellini, G., Muegge, I.: Prediction of Human Clear-
1023 ance Using In Silico Models with Reduced Bias. *Molecular Pharmaceutics* (2024)
1024 <https://doi.org/10.1021/acs.molpharmaceut.3c00812>
- 1025 [91] Šícho, M., Luukkonen, S., Maagdenberg, H.W., Schoenmaker, L., Béquignon,
1026 O.J.M., Westen, G.J.P.: DrugEx: Deep Learning Models and Tools for Explo-
1027 ration of Drug-Like Chemical Space. *Journal of Chemical Information and*
1028 *Modeling* **63**(12), 3629–3636 (2023) <https://doi.org/10.1021/acs.jcim.3c00434>
- 1029 [92] Noutahi, E., Wognum, C., Mary, H., Hounwanou, H., Kovary, K.M., Gilmour, D.,

- 1030 thibaultvarin-r, Burns, J., St-Laurent, J., t, DomInvivo, Maheshkar, S., rbyrne-
1031 momatx: datamol-io/molfeat: 0.9.4. Zenodo (2023). [https://zenodo.org/records/
1032 8373019](https://zenodo.org/records/8373019)
- 1033 [93] Tietz, M., Fan, T.J., Nouri, D., Bossan, B., Developers: skorch: A scikit-learn
1034 compatible neural network library that wraps PyTorch (2017). [https://skorch.
1035 readthedocs.io/en/stable/](https://skorch.readthedocs.io/en/stable/)
- 1036 [94] Sicho, M., Liu, X., Svozil, D., Westen, G.J.P.: GenUI: interactive and extensible
1037 open source software platform for de novo molecular generation and cheminform-
1038 matics. *Journal of Cheminformatics* **13**(1), 73 (2021) [https://doi.org/10.1186/
1039 s13321-021-00550-y](https://doi.org/10.1186/s13321-021-00550-y)

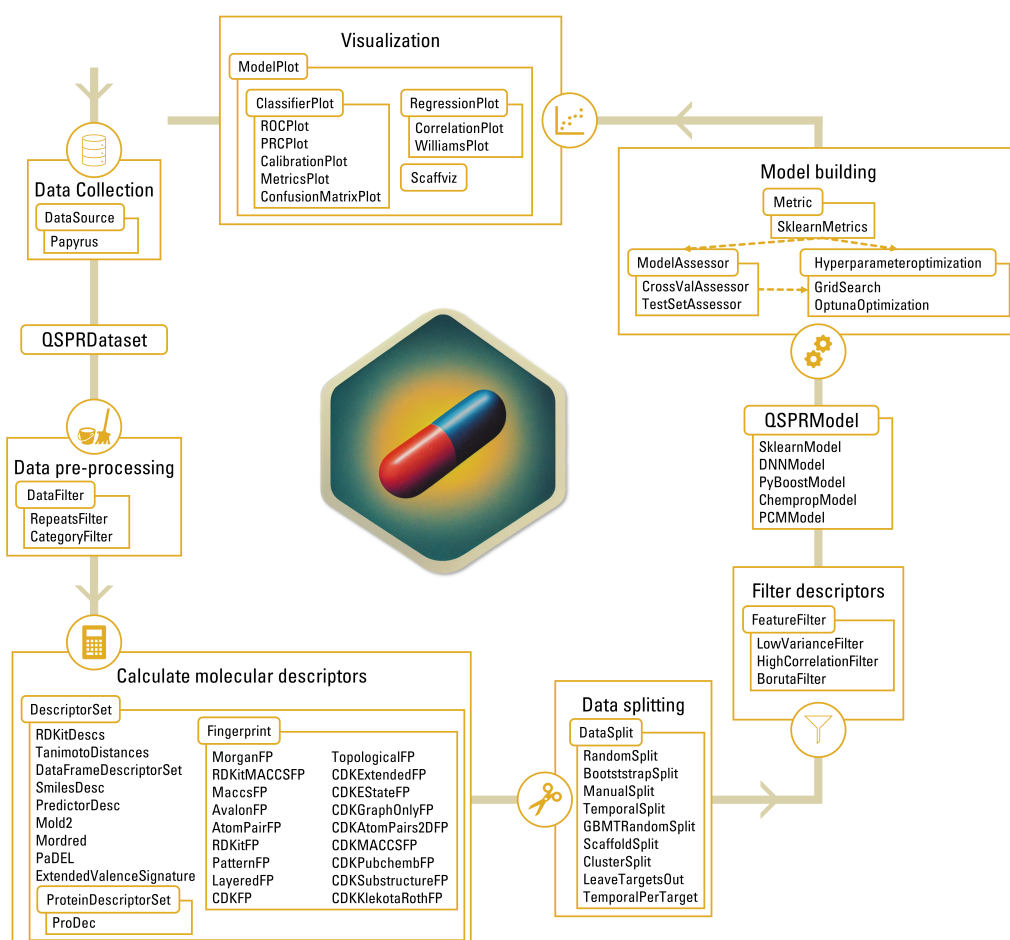


Fig. 1 Visualization of the QSPRpred workflow. Each box represents a general step in QSPR/QSAR modelling, e.g. data collection and visualization. Every rounded rectangle is an abstract base class in QSPRpred defining the interface of the respective step. Each of these classes has a number of implementations included, which are listed in the attached box, e.g. the abstract base class `DataFilter` has `RepeatsFilter` and `CategoryFilter` available as an out-of-the-box implementation. Therefore, altering the behaviour of each component can be achieved either through inheritance or by simply providing a custom implementation if the respective abstract base class.

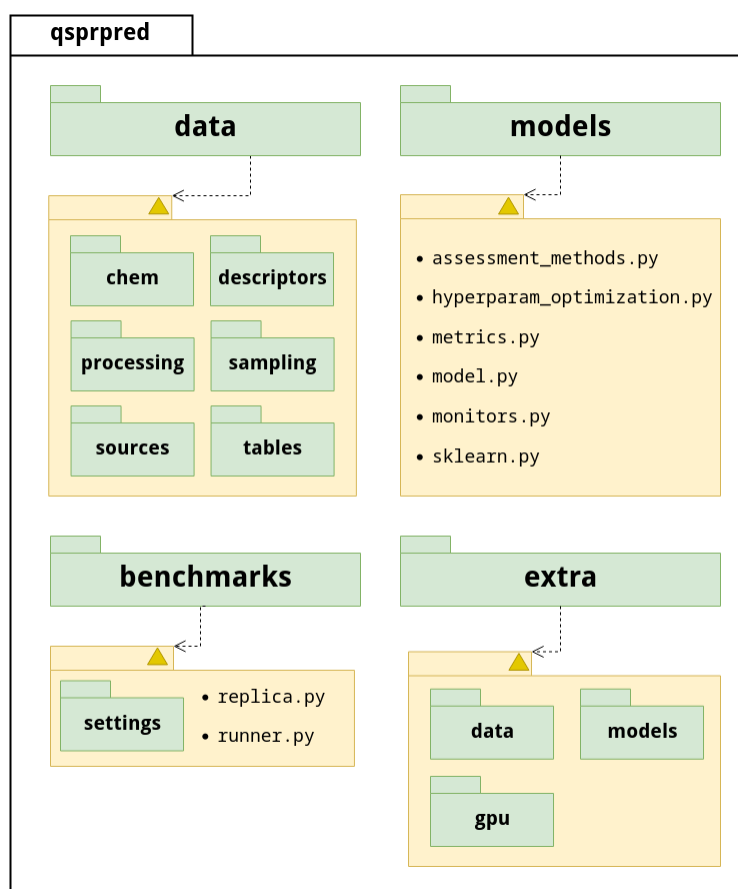
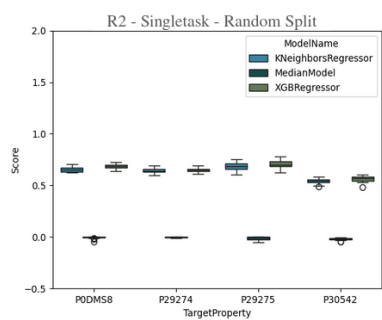
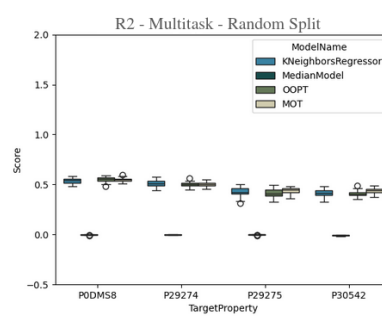


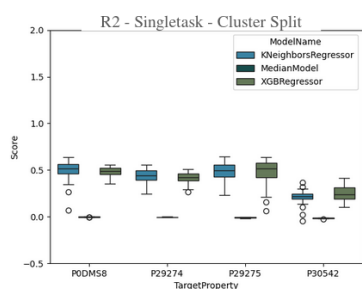
Fig. 2 Graphical overview of the QSPRpred package architecture. Only the main packages discussed in this work are included. The majority of the functionality is centered around the `qsprpred.data` package, which integrates API definitions and tools to work with chemical structures (`qsprpred.data.chem`), featurization (`qsprpred.data.descriptors`), processing tools for data and feature filtering `qsprpred.data.processing`, data splitting and resampling (`qsprpred.data.sampling`), data source adapters (`qsprpred.data.sources`) and data storage implementations (`qsprpred.data.tables`). The `qsprpred.models` package contains the base definition of a QSPRpred model API in `qsprpred.models.model`, which is implemented for scikit-learn models in `qsprpred.models.sklearn`. In addition, the `qsprpred.models` package also contains functionality needed to monitor training (`qsprpred.models.monitors`), optimize (`qsprpred.models.hyperparam_optimization`) and assess model performance (`qsprpred.models.assessment_methods` and `qsprpred.models.metrics`). The `qsprpred.extra` package extends functionality of both `qsprpred.data` and `qsprpred.models` with additional extensions that support deep neural networks and PCM modelling functionality. Finally, the `qsprpred.benchmarks` package houses data classes needed to set up (`qsprpred.benchmarks.settings`) and run benchmarking experiments (`qsprpred.benchmarks.runner`).



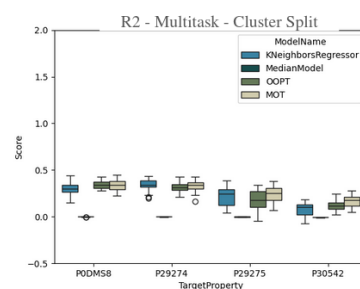
(a) R2 scores determined for single-task models on a test set determined by a random split of training data.



(b) R2 scores determined for multi-task models on a test set determined by a random split of training data.

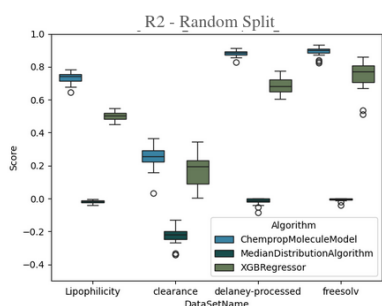


(c) R2 scores determined for single-task models on a test set determined by a cluster split of training data.

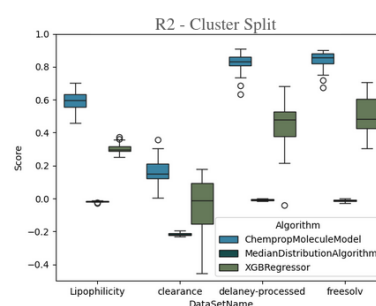


(d) R2 scores determined for multi-task models on a test set determined by a cluster split of training data.

Fig. 3 Coefficients of determination (R2) calculated for each replica in different benchmarking runs conducted in Experiment 1.



(a) R2 scores determined for regression models on the selected MoleculeNet data sets for various model architectures using a random training-test split.



(b) R2 scores determined for regression models on the selected MoleculeNet data sets for various model architectures using a cluster training-test split.

Fig. 4 Coefficients of determination (R2) calculated for each replica in different benchmarking runs conducted in Experiment 2.