# Papyrus - A large scale curated dataset aimed at bioactivity predictions

O. J. M. Béquignon[1†], B. J. Bongers[1†], W. Jespers[1],  A. P. IJzerman[1], B. van der Water[1], G. J. P. van Westen[1*]

1 Division of Drug Discovery and Safety, Leiden Academic Centre for Drug Research, Leiden University, Leiden, The Netherlands
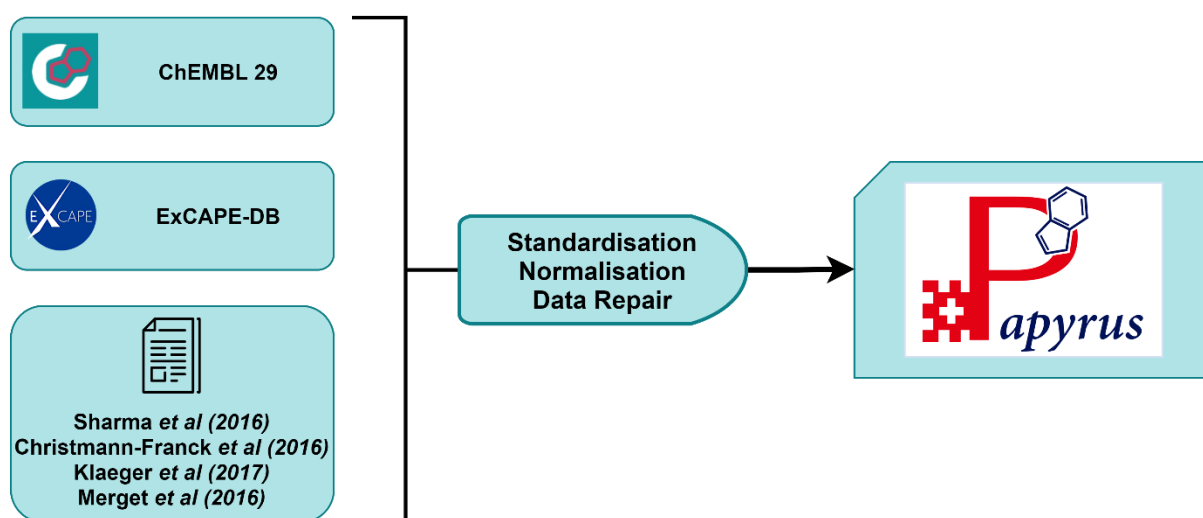

* Corresponding author

Email: gerard@lacdr.leidenuniv.nl (GJPW)


† These authors contributed equally to this work

## Abstract

With the recent rapid growth of publicly available ligand-protein bioactivity data, there is a trove of viable data that can be used to train machine learning algorithms. However, not all data is equal in terms of size and quality, and a significant portion of researcher's time is needed to adapt the data to their needs. On top of that, finding the right data for a research question can often be a challenge on its own. As an answer to that, we have constructed the Papyrus dataset (DOI: 10.4121/16896406), comprised of around 60 million datapoints. This dataset contains multiple large publicly available datasets such as ChEMBL and ExCAPE-DB combined with several smaller datasets containing high quality data. The aggregated data has been standardised and normalised in a manner that is suitable for machine learning. We show how data can be filtered in a variety of ways, and also perform some baseline quantitative structure-activity relationship analyses and proteochemometrics modeling. Our ambition is this pruned data collection constitutes a benchmark set that can be used for constructing predictive models, while also providing a solid baseline for related research.

**Keywords:** machine learning, cheminformatics, bioactivity, curated dataset, Papyrus, standardisation, normalisation.

## Introduction

Academic computational drug discovery has gained a massive boost with the growth of publicly available data[1,2]. One of the areas where this has led to improvement is the prediction of bioactivity, specifically ligand-protein affinity. Databases such as ChEMBL and BindingDB provide a wealth of information and relationships between ligands, proteins, and their interaction[3,4]. However, public data has a diverse quality range and is subject to experimental error[5,6]. In contrast to large datasets like ChEMBL, there are also smaller more focused datasets available. These typically focus on a single protein family and usually from a single set of literature such as the Klaeger clinical kinase drugs dataset[7]. Such collections contain a trove of high quality data, but are limited in their scope and are usually not viable as sole data sources in a more general study.

In previous work we compared the performance of established bioactivity prediction methods versus deep neural networks[8]. In order to publish the results from this benchmark the creation of a public dataset (to accompany the publication) was required. ChEMBL (version 20) was used and a high quality subset was extracted and made available[9]. There were several problems we ran into. Firstly there was the amount of work needed to prepare the ChEMBL dataset, leading to an inability to include the separate smaller scale datasets we were planning to add.. In addition, there was a large reduction in size due to the selection of high quality data with the final dataset being 2.5% of the total ChEMBL data.

The current research aims to address these issues and produce standardized diverse dataset. This dataset, named Papyrus[10] (in reference to Leiden Papyrus X), is created with ease-on-use and filtering in mind. We want to remove some of the limitations mentioned above and provide a dataset that does not need further curating. Aside from ChEMBL, we implemented data from the ExCAPE-DB[11] database, and added Sharma *et al*'s[12], Christmann-Franck *et al*'s[13], Klaeger *et al*'s[7] and Merget *et al*'s[14] data.

63  This current work contributes by providing a standardised and normalised dataset that
64  can be used 'out-of-the-box'. On top of that, we provide multiple sets of integrated
65  descriptors that are widely used in literature. We also show how to manipulate and
66  model the data using proteochemometrics and provide the Python scripts that we
67  used[15]. Lastly, as the focus in Papyrus is on filtering as well, we provide several Python
68  scripts for ease of querying the dataset. This includes filters for organism, activity type,
69  and accession numbers, and a link to the scripts can be found at the end of this
70  document (under 'Data Availability').

71

## Material and Methods

### Construction of Papyrus

#### ChEMBL

75  Three levels of quality were defined in the data: high, medium, and low. Data from the
76  difference sources were all classified in one of these three classifications. ChEMBL
77  version 29 (ChEMBL29) data[16] were first split between high and low quality data In
78  total 18,635,916 activity data points measured on 2,105,464 compounds and 14,554
79  targets were extracted. The following data were deemed as low quality: Data flagged
80  as potential duplicates, of questioned validity (Supplementary Table 1) - unless errors
81  were confirmed by authors, in which case they were entirely disregarded -, censored,
82  not associated with any pChEMBL value, or of questioned activity (Supplementary
83  Table 2). Remaining activity data were temporarily regarded as high quality. In the
84  ExCAPE-DB dataset[11], if the data source was PubChem[17] (source identifier 7), then the
85  data was flagged as high quality. Protein targets were retrieved along with their
86  classifications if they had Uniprot[18] accessions defined. Accession Q8MMZ4,
87  corresponding to the secondary accession for the Plasmodium falciparum (isolate
88  NF54) cGMP-dependent protein kinase was associated with a secondary UniProt
89  accession and was manually replaced by its primary accession Q8MMZ4. Using
90  accessions, protein sequences were then obtained from UniProt.  Only molecules
91  identified as small molecules and associated with a molecular registry number were
92  kept, then parsed from connection tables and standardised (see Molecular structure

93  standardisation). Activity data of high quality were reclassified as low quality if the

94  target type was other than 'single protein' or the assay confidence score was 0, 1, 2, 3,

95  4 or 6 (Supplementary Table 3). Activity data with assay confidence scores of 5 and 8

96  were reclassified of medium quality.

97

98  If low quality data were censored, inequality signs of the standard relation were

99  reversed (Supplementary Table 4) unless expressing an approximation with a tilde, in

100  which case the data were dropped. Standard values of low quality data with unassigned

101  pChEMBL values were only considered if they had case insensitive standard type of

102  either GI50, Ki, Kd, IC50, EC50 or XC50 and if standard units denoted molar or mass

103  concentrations. Scaling factors were appropriately applied to standard activity values

104  (Supplementary Table 5), mass concentrations were transformed to molar

105  concentrations and log-scale transformation applied to all concentrations. Only

106  exceptions to data with unassigned pChEMBL values were records with derivatives of

107  the following standard types: pKi, pKd, pIC50, pEC50 or pXC50 (Supplementary Table

108  6). For those records, no transformations were applied. Finally, data were flagged on

109  whether activities were derived from IC50, EC50, Ki, Kd or any other data.   The

110  preprocessed ChEMBL29 high quality data consisted of 1,097,673 activity values,

111  549,140 compounds and 4,644 targets, the medium quality data of 489,315 activity

112  values, 263,824 compounds and 2,886 targets, and low quality data of 1,510,494

113  activity values, 514,302 compounds and 4,711 targets.

114

115  *ExCAPE-DB*

116  The ExCAPE-DB dataset, consisting of 70,850,163 activity data points of 998,131

117  compounds measured on 1,667 targets, was first discarded of records originating from

118  ChEMBL version 20 or whose assay identifiers were present in the PubChem flagged

119  data of the preprocessed ChEMBL29. Gene Entrez identifiers were mapped, using the

120  identifier mapping tool of UniProt, to protein unique Swiss-Prot sequences. Only four

121  genes were manually mapped,  three genes as they resolved to multiple reviewed

122  entries and one gene as it resolved to multiple unreviewed entries(Supplementary

123    Table 7). ChEMBL29 protein classifications were then assigned to the previously

124    mapped sequences. Deposition dates of the assays were retrieved from PubChem.

125    Data with numeric activity values were considered of high quality and binary data of

126    low quality. Molecular structures failing standardisation (see Molecular structure

127    standardisation) were downloaded from PubChem and standardised afterwards. Finally

128    low quality activity data with compound-target pairs present in the high quality subset

129    were disregarded. The preprocessed ExCAPE-DB high quality data consisted of

130    278,226 activity values, 201,644 compounds and 1,535 targets, and low quality data of

131    58,445,354 activity values, 650,217 compounds and 646 targets.

132

133    *Sharma et al*

134    Sharma *et al*'s dataset[12], consisting of 258,060 activity data points of 76,017

135    compounds measured on 8 targets was considered of high quality. Gene names were

136    mapped to unique Swiss-Prot sequences using the identifier mapping tool of UniProt

137    and protein classifications retrieved from ChEMBL29. A set of 14 custom reactions

138    (Supplementary Table 8) were applied to molecular structures failing standardisation

139    (see Molecular structure standardisation), mostly fixing aromaticity-related issues.

140    Years of filing of patents were collected using Google Cloud BigQuery API patents

141    public data and manual mapping (Supplementary Table 9 & 10) after having fixed

142    erroneous patent numbers (Supplementary Table 11). Digital object identifiers or

143    PubMed identifiers of source articles were added when missing (Supplementary Table

144    12). If activity values were associated with multiple sources, only the first published

145    article or filed patent was recorded. Censored activity values or values not associated

146    with case insensitive standard types GI50, Ki, Kd, IC50 or EC50 and their

147    logarithmically-derived counterparts were disregarded. Mass concentrations were

148    transformed to molar concentrations and log-scale transformation applied to all

149    concentrations but already log-transformed. Finally infinite or activity values lower than

150    3 and higher than 14 log units were discarded. The preprocessed Sharma data

151    consisted of 77,562 activity values, 40,738 compounds and 8 targets.

152

153 *Christmann-Franck et al*

154 Christmann-Franck *et al*'s dataset[13], consisting of 344,788 activity data points of 2,065

155 compounds measured on 448 targets was considered of high quality. The wrongly

156 assigned Cryptococcus neoformans mitogen-activated protein kinase (CPK1) with

157 accession code P0CP66 was corrected to the Plasmodium falciparum calcium-

158 dependent protein kinase 1 (CDPK1) with accession code P62344. Swiss-Prot

159 sequences were retrieved using accessions and protein classifications retrieved from

160 ChEMBL29. Sequence mutations of the hepatocyte growth factor receptor (MET) and

161 the serine/threonine-protein kinase (B-raf) were corrected to M1250T and V600E

162 respectively and that of the Fibroblast growth factor receptor 1 (FGFR1) was reverted

163 to wildtype. Activity data expressed as proportion of reference activities were

164 discarded. Finally molecular structures were standardised (see Molecular structure

165 standardisation). The preprocessed Christmann-Franck data consisted of 135,948

166 activity values, 1,669 compounds and 485 targets.

167

168 *Klaeger et al*

169 Klaeger *et al*'s dataset[7], consisting of 5,916 activity data points of 229 compounds

170 measured on 520 targets was considered of high quality. Swiss-Prot sequences were

171 retrieved using HUGO Gene Nomenclature Committee (HGNC) identifiers. If multiple

172 identifiers were assigned the measurement was discarded. Protein classifications were

173 retrieved from ChEMBL29. Apparent Kd values were log-transformed and infinite

174 results disregarded. Finally molecular structures were standardised (see Molecular

175 structure standardisation), with only RDEA-436 failing for its structure was not

176 disclosed. The preprocessed Klaeger data consisted of 5,721 activity values, 228

177 compounds and 500 targets.

178

179 *Merget et al*

180 Merget *et al*'s dataset[14], consisting of 260,757 activity data points of 47,774

181 compounds measured on 241 targets was considered of high quality, except for activity

182 values originating from ChEMBL22, which were disregarded. Data originating from the

183  Published Kinase Inhibitor Set (PKIS) of GlaxoSmithKline (doi:10.1038/nbt.3374) with

184  activity values of 5 log units were considered as censored and as such reclassified as

185  low quality data. Swiss-Prot sequences were retrieved using HUGO Gene

186  Nomenclature Committee (HGNC) identifiers[19], a few of which were manually fixed

187  (Supplementary Table 13). Protein classifications were retrieved from ChEMBL29.

188  Finally molecular structures were standardised (see Molecular structure

189  standardisation). The Merget preprocessed high quality data consisted of 127,441

190  activity values, 1,666 compounds and 239 targets, and low quality data of 62,642

191  activity values, 360 compounds and 195 targets.

192

193  *Molecular Structure Standardisation*

194  During the preprocessing of each original dataset parent molecular structures were

195  gathered after a first standardisation using the ChEMBL structure pipeline[20]. Then

196  canonical tautomers were determined using the Pipeline Pilot tautomer enumerator[21]

197  with tautomerization of amides enabled. The canonical tautomers were then

198  standardised once again with the ChEMBL structure pipeline after which the parent

199  structures were obtained. Any molecule not parsable from simplified molecular input

200  line entry specification (SMILES) by the RDKit[22] at any step of the previous workflow

201  was considered failing the standardisation process.

202

203  After the individual datasets were processed and aggregated into the Papyrus dataset,

204  molecular structures were standardised. This last standardisation ensured

205  normalisation across different sources that each applied different prior

206  standardisation. For instance, tautomerization tools can alter a compound's

207  stereochemistry by removing or introducing a chiral centre. To limit the effects of

208  having bioactivity values relating to the same molecular compound having different

209  stereochemistry across sources, a set with removed stereochemistry was created and

210  deemed of higher quality than the set with conserved stereochemistry. Molecular

211  structures in the latter, after having removed stereochemistry, were first neutralized

212  with the RDKit by adding or removing hydrogen atoms. Subsequently they were

213  standardised with the ChEMBL structure pipeline after which parent structures were
214  obtained. OpenBabel[23,24] was then used to recreate dative bonds and to neutralize
215  molecules that were not during the previous step. Tetravalent negatively charged boron
216  atoms were overlooked in the latter stage, making them erroneously pentavalent. This
217  was corrected by detecting these pentavalent boron atoms, removing the newly
218  introduced implicit hydrogen atom and reassigning them a negative formal charge.
219  SMILES of molecules with the same connectivity differing by overall charge were
220  converted to InChi[25] with OpenBabel, a step that consists in incorporating the
221  normalisations after the InChI canonicalization process. For these molecules the
222  canonicalization process removed the dative bonds, these were then recreated using
223  OpenBabel. Then Dimorphite-DL[26] was used to deprotonate molecules by setting
224  minimum pH to 14.0. This ensured that after the last standardisation step, equivalent
225  to that applied to individual datasets, in which molecules are neutralised, only one
226  charge state of the same molecular species was present in the set.

227

228  *Papyrus data aggregation*

229  The processed ChEMBL29 high, medium and low quality, ExCAPE-DB high and low
230  quality, Sharma, Christmann-Franck, Klaeger and Merget datasets were aggregated
231  together. The first step consisted in ensuring that the activity of any compound-target
232  pair was contained within 3 to 14 log units. Then compound-target pairs were uniquely
233  identified by a concatenation of the compound's connectivity and of the target
234  accession along with its mutations if any. All activities relating to the same compound-
235  target pair were then filtered depending on the highest data quality available for that
236  pair. For instance, if high quality activities were identified, any data point deemed of
237  medium to low quality was filtered out. Considering censored activity values, the data
238  was filtered out if contradictory relations were identified, if not the highest recorded
239  activity was retained for lower bounds, and lowest for higher bounds. During this
240  filtering step, all patents and journal articles associated with the activity of a
241  compound-target pair were gathered whatever the quality and only the first published
242  or filed was retained. Finally, activity values were aggregated and mean averages,

9

243  medians, standard errors of the mean, standard deviations and mean average

244  distances were calculated for each unique compound–target pair.

245

246  ***Use of Papyrus***

247  *Data extraction*

248  The first subset that was extracted from Papyrus consists in adenosine receptors.

249  Using the Papyrus Python scripts, data of high quality with protein classification level

250  5 being "Adenosine receptor" was extracted. This subset, consisting of 15,941 activity

251  points, 24 protein targets and 7,967 compound structures.

252

253  *Data visualization*

254  Unique molecules of Papyrus were collected based on the uniqueness of their

255  connectivity. Each molecule was encoded using MinHash fingerprint (MHFP6)[27] and

256  then visualized using TMAP[28]. Molecules were labelled using their fraction of carbon

257  atom. Unique proteins of Papyrus were collected based on their unique target

258  identifier. Each sequence was encoded using UniRep[29] 64, 256 and 1,900 average

259  hidden states, final cell states and final hidden states. The 6660 dimensions were then

260  MinHashed and visualized with TMAP. Proteins were labelled using organisms they

261  originate from.

262

263  *Bioactivity modelling: Quantitative Structure Activity-Relationships*

264  Each protein target in the subset was modelled independently using the Papyrus

265  Python scripts. Targets for which less than 30 activity values or associated with activity

266  values spanning less than 2 log units were disregarded for modelling. Then for each

267  target, a temporal split between training and test sets was performed: datapoints

268  associated with year 2013 and above constituted the test set. If no activity data was

269  available after year 2013, then the target was disregarded. The 777 Mold2 molecular

270  descriptors[30] were calculated for each molecule and were centered and scaled to unit

271  variance. Extreme Gradient Boosting (XGBoost version 1.4.2) regressors and classifiers

272 were trained on the training set using random seed 1234 and default parameters.
273 Regressors were trained to predict mean pChEMBL values using 5-fold cross-
274 validation, while classifiers were trained to predict a binary label of activity class with
275 threshold set at 6.5 log units using 5-fold stratified cross-validation.

276

277 *Bioactivity modelling: Proteochemometrics*
278 No subsequent filtering of the subsets was carried out since proteochemometrics
279 (PCM) handles multiple targets all at once. A temporal split on year 2013 was employed
280 to split the training and test set. The 777 Mold2 molecular descriptors were calculated
281 for compounds, UniRep 64, 256 and 1,900 average hidden states, final cell states and
282 final hidden states were used as 6,660 protein descriptors and were calculated for
283 each protein. An XGBoost classifier and an XGBoost regressor were trained using the
284 same protocol as for QSAR models.

285

286 **Results and Discussion**

287 A new dataset, called Papyrus of bioactivities, resulting from the aggregation and
288 extensive standardisation of data from six sources, was created. Unless mentioned
289 otherwise, only the extensively standardised Papyrus set without stereochemistry is
290 considered in this section.

291

292 *Papyrus dataset statistics*
293 Papyrus consists of 59,763,781 compound-protein pairs, each associated with at least
294 either one activity value or activity class. Additionally, this represents the data of
295 1,268,606 unique compounds and 6,996 proteins across 496 different organisms. In
296 terms of data quality, 1,236,296 datapoints are of high quality, i.e., representing exact
297 bioactivity values measured and associated with a single protein or complex subunit.
298 335,854 datapoints are of medium quality, i.e., exact bioactivity values associated with
299 either potentially multiple proteins or a homologous single protein. 58,191,631
300 datapoints are of low quality, i.e., exact bioactivity values associated with either

301 multiple homologous proteins or homologous complex subunits, censored bioactivity
302 values and binary activity classes. When considering datapoints across all quality types,
303 2,585,248 are associated with exact bioactivity values, 354,981 with censored data
304 and 56,823,552 with binary activity classes. The repartition of data quality across the
305 ten organisms with most data (Table 1**Error! Reference source not found.**) indicates
306 a clear bias towards human, with more than 93% of the data related to it, but also
307 emphasizes the interest towards rodent targets with more than 4% of the data
308 associated with mouse and 2% with rats.

| Species | Quality | | | Total |
|---|---|---|---|---|
| | High | Medium | Low | |
| Homo sapiens (Human) | 985,579 | 246,723 | 54,363,214 | 55,595,516 |
| Mus musculus (Mouse) | 41,986 | 6,682 | 2,465,153 | 2,513,821 |
| Rattus norvegicus (Rat) | 60,374 | 32,075 | 1,151,936 | 1,244,385 |
| Escherichia coli (strain K12) | 539 | 11,283 | 54,800 | 66,622 |
| Equus caballus (Horse) | 18,326 | 32 | 27,987 | 46,345 |
| Influenza A virus (A/WSN/1933(H1N1)) | 23,813 | - | 9,143 | 32,956 |
| Trypanosoma cruzi | 5,886 | 30 | 23,927 | 29,843 |
| Schistosoma mansoni (Blood fluke) | 13,916 | - | 14,473 | 28,389 |
| Bacillus subtilis | 12,106 | - | 11,693 | 23,799 |
| Bos taurus (Bovine) | 5,923 | 5,107 | 8,913 | 19,943 |

*Table 1: Activity data of organisms in Papyrus with the most datapoints.*

309

310 When it comes to the activity types Papyrus is derived from (Table 2), most of the data
311 is either associated with untraceable data types, such as for binary data, or with types
312 derived from others, for instance the KIBA scores were derived from $IC_{50}$, $K_i$ and $K_D$
313 data[31] present in the Merget source dataset.

314

315

316

| Activity type | Original datapoints |
|:---:|:---:|
| $K_i$ | 507,821 |
| $K_D$ | 118,773 |
| $IC_{50}$ | 1,070,430 |
| $EC_{50}$ | 141,672 |
| Other | 58,315,137 |

*Table 2: Number of original datapoints in Papyrus for each activity type.*

317

318 Papyrus protein space (Figure 1A) is largely dominated by human proteins, reflecting
319 the abundance of activity values measured for these. Nevertheless, clusters of
320 homologous proteins can be observed, mostly aggregating human rat and mouse
321 protein. As a comparison, the compound space was also visualized (Figure 1B) with
322 carbon fraction evenly spread across clusters.

323 Concerning protein classification, enzymes represent nearly half of the classified and
324 annotated proteins, with more than 25 million data points, and membrane receptors
325 21% with more than 11 million (Figure 2A). Family A G protein-coupled receptors
326 represent 37% of proteins annotated with a second level class (Figure 2B), consisting
327 in more than 9 million datapoints, proteases 23%, more than 5 million, and kinases,
328 long thought undruggable targets, represent 18% of the data with more than 4.5 million
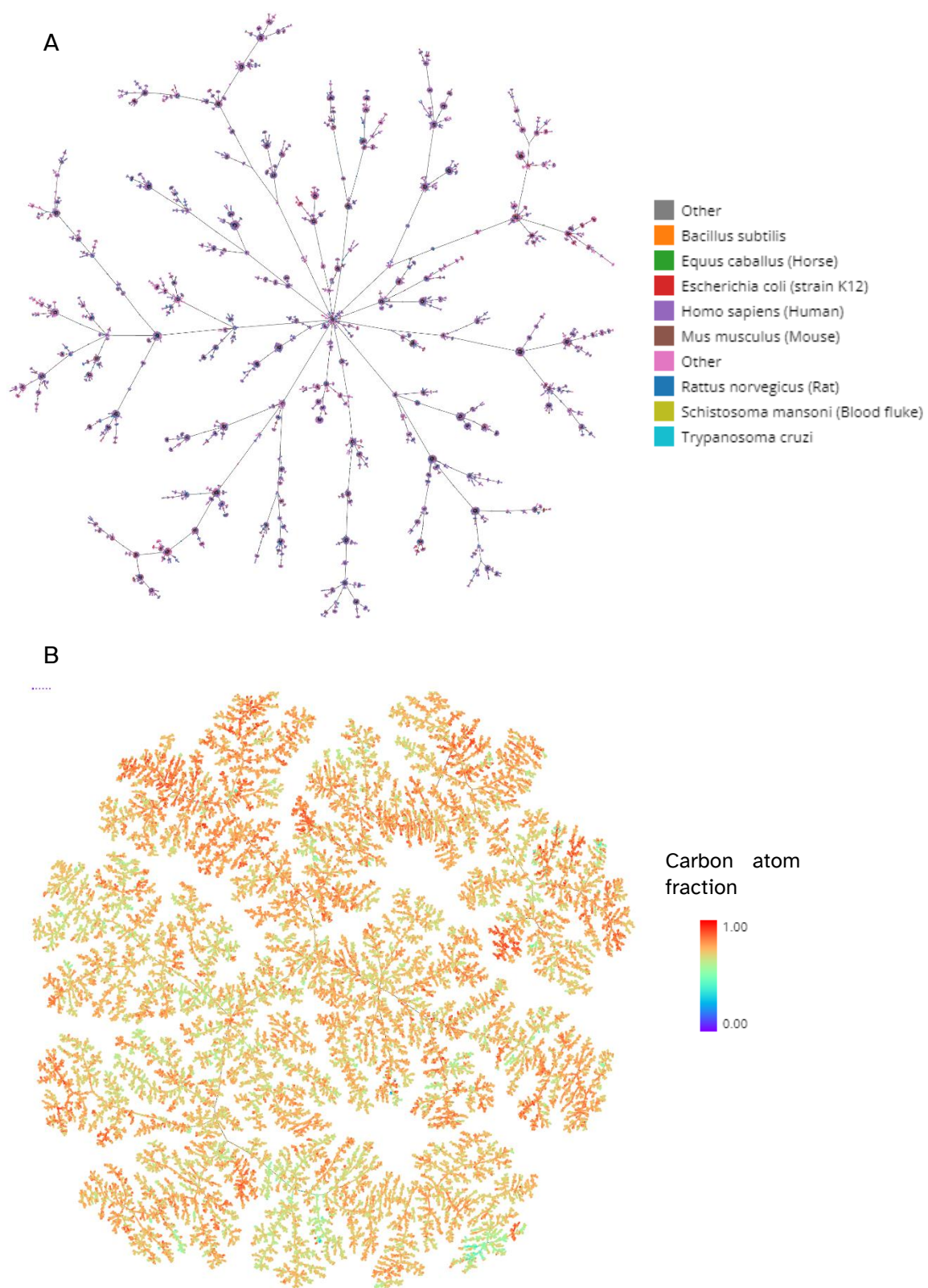329 datapoints.

330

A



| | |
|---|---|
| ■ | Other |
| ■ | Bacillus subtilis |
| ■ | Equus caballus (Horse) |
| ■ | Escherichia coli (strain K12) |
| ■ | Homo sapiens (Human) |
| ■ | Mus musculus (Mouse) |
| ■ | Other |
| ■ | Rattus norvegicus (Rat) |
| ■ | Schistosoma mansoni (Blood fluke) |
| ■ | Trypanosoma cruzi |

B



Carbon atom fraction

1.00

0.00

*Figure 1: Papyrus protein (A) and chemical (B) spaces.*

331

*Figure 2: Protein classification levels 1 (A) and 2 (B) of protein targets in Papyrus.*

*Bioactivity modeling: Results*

To exemplify the potential of Papyrus, an Adenosine receptor (AR) subset was extracted, considering only the high-quality data. Quantitative structure-activity relationships (QSAR) and proteochemometrics (PCM) regression and classification models were trained. A temporal split scheme was chosen, to better asses the prediction performance of the models[32] and minimize congeneric series being split between training and test sets.

QSAR models were trained on protein targets with sufficient data. This resulted in only 11 of 24 the adenosine receptors in the subset to be suitable for QSAR modelling. PCM models allowing the use of all related targets, all 24 adenosine receptors could be modelled.

QSAR regression models for human ADORA2b, rat ADORA1, ADORA2a and ADORA3 and mouse ADORA3 performed well with coefficient of determination $R^2$ between 0.6 and 0.7 during cross-validation (Figure 3A). Surprisingly human ADORA1, ADORA2a, mouse ADORA3 and bovine ADORA1 performed quite bad with median $R^2$ lower than 0.5 with fold performance reaching -2.3 to -2.1 for the first three. Nonetheless, the maximum error associated with these first three was significantly lower than that of other QSAR models although the root-mean-square error (RMSE) and mean absolute

15

352 error (MAE) of all models were not significantly different. With regards to external

353 validation (Figure 4B), the performance of the predictions on the temporally split test

354 set performed as expected with RMSE around 1 log unit for most models, only human

355 ADORA1 and mouse ADORA3 performing noticeably badly. The PCM regression model

356 performed quite well at cross-validation in terms of $R^2$ and RMSE with values of 0.59

357 and 0.72 but had higher median maximum error than all QSAR models. These results

358 were reflected in the temporal validation.

359 QSAR classification models of human ADORA2a performed very well with Matthews

360 correlation coefficient (MCC) ranging between 0.62 and 0.70 for cross-validation and

361 0.48 on the temporal test set (Figure 4). Except for the human ADORA3 and rat

362 ADORA3 that performed bad both during cross-validation and testing due to the

363 imbalance of the datasets (ratios of 1:7 to 2:6 of actives to inactives for human ADORA3

364 and 4:1 to 5:1 for rat ADORA3) and showing very variable sensitivity and specificity,

365 most models performed equally well at cross-validation and on the test set. The human

366 ADORA1 and ADORA2a, rat ADORA1, ADORA2a and ADORA2b, mouse ADORA1 and

367 bovine ADORA1 had balanced accuracy (BAcc) over 0.70 and area under the receiver

368 operator characteristic curve (AUC) over 0.65, which showed very good predictive

369 performance in a prospective setting. It is worth noting that the bovine ADORA1 QSAR

370 regression model $R^2$ was one of the lowest (-0.27). The PCM classification model

371 showed performance on par with well performing QSAR models during cross-validation

372 but showed lower performance on test set with MCC of 0.25 and BAcc of 0.62.

373 Overall models on the AR subset showed similar performance between regression and

374 classification. It is no surprise that the receptors that performed best, i.e. most of rat

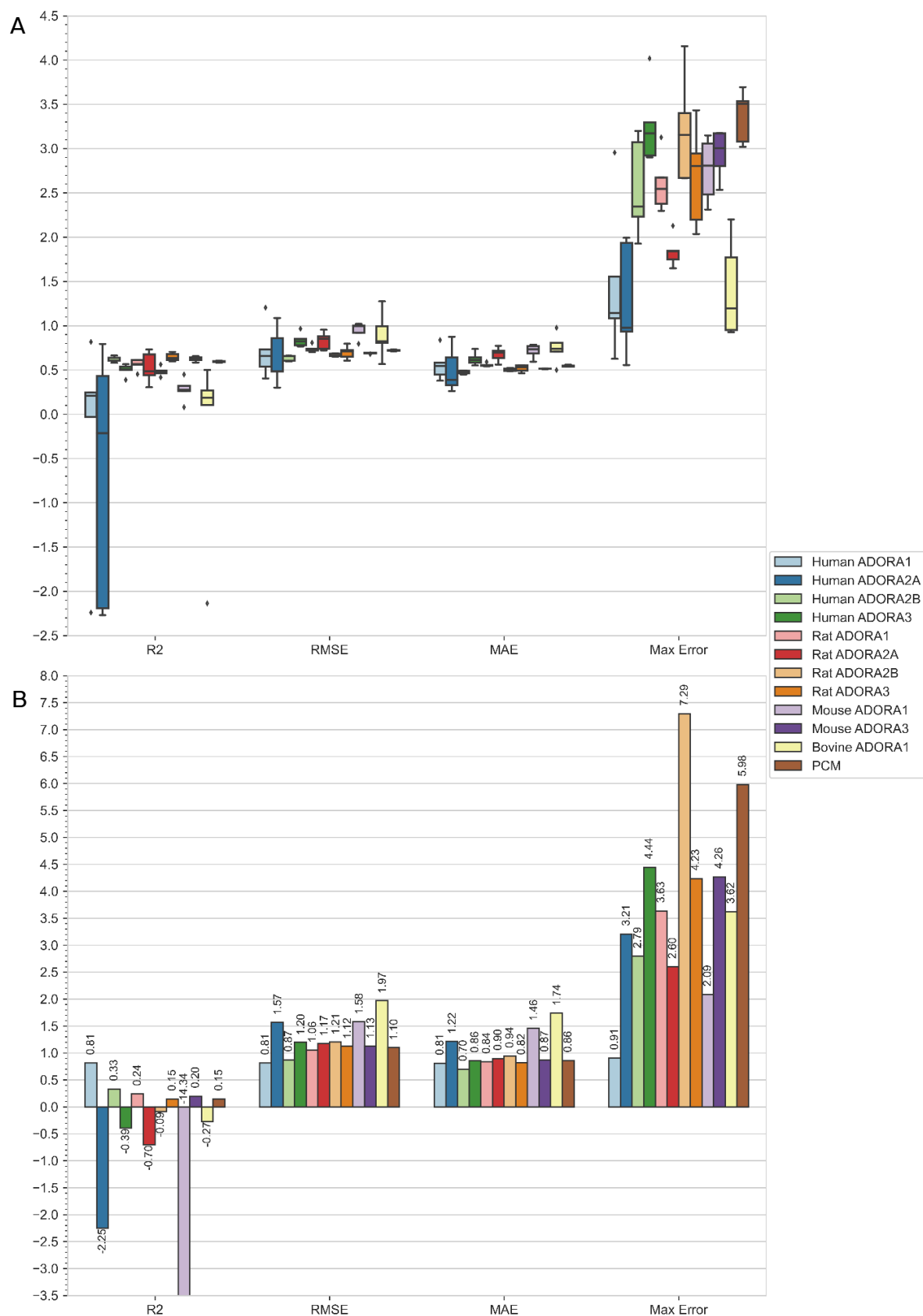375 and human receptors, were those with the most datapoints.

376

*Figure 3: Cross-validation performance (A) and temporally split test set performance (B) of regression QSAR and PCM models. $R^2$: coefficient of determination, RMSE: root-mean-square error, MAE: mean absolute error, Max Error: Maximal error.*
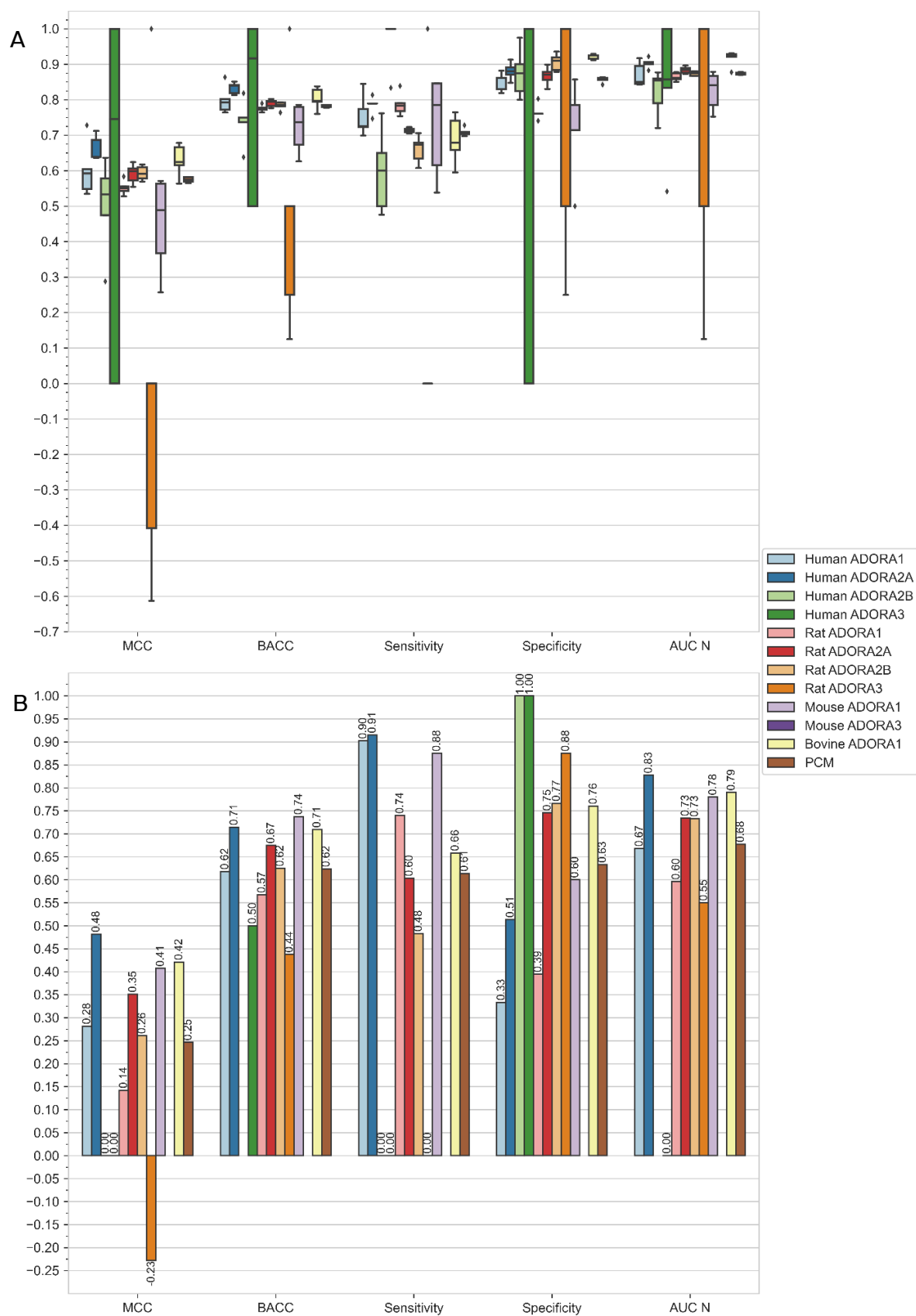
17

377

*Figure 4: Cross-validation performance (A) and temporally split test set performance (B) of regression QSAR and PCM models. R²: coefficient of determination, RMSE: root-mean-square error, MAE: mean absolute error, Max Error: Maximal error.*

8

378

## Discussion

We have shown the format of the Papyrus bioactivity dataset, as well as a few examples of baseline models that could be created from this data. While we are confident that is a reliable publicly available dataset, there are still some limitations present.

First, the most important limitation in our eyes is that although Papyrus consists of nearly 60 million activity data points, the data it contains is extremely sparse as only 0.67% of the activity data matrix is represented by the 1.25 million compounds and almost 7,000 proteins. This is unfortunately hard to avoid as many of the compounds have simply not been tested on all proteins. Only relatively popular proteins will appear in the data that is aggregated here. This makes it hard to model proteins that are understudied, however if data are available they can be added to the Papyrus dataset to create a more comprehensive set.

As Papyrus is a static dataset, updates (or corrections) are possible but are reliant on the aggregated datasets. While this is always a restriction on static datasets, there is a second degree of reliance here as all the data needs to be updated in their respective datasets. While we do not think this will pose an issue for modeling, as data freeze often occurs when a dataset is used in research, updates to Papyrus will have more time between them than the respective aggregated datasets.

Another limitation is the choice for the specific datasets that were aggregated into Papyrus in addition to ChEMBL. Firstly we have implemented ExCAPE-DB, like we did in previous work. The single article datasets are a reflection of some of the interests of our group, and we value the high quality data present. There are enough arguments to include a certain dataset that we did not mention here, which would improve the quality of the dataset even more. We do provide the full dataset and all the descriptors that are used. So if anyone wants to add a certain dataset the tools are available to do

407 so. Our goal was to create a benchmark set to set a reliable standard to perform
408 bioactivity modeling on, and we think that Papyrus meets that goal.

409

410 Additionally repetition of data in the source datasets was scrutinized and where
411 possible only the most recent bioactivity data was kept. This is for instance the case
412 of the KIBA scores of Tang *et al*[6] that used a combination of activity types of ChEMBL
413 to increase the quality of single measurements, or of ExCAPE-DB or the ChEMBL
414 subset in Merget *et al*[14] aggregating data from ChEMBL versions 20 and 21
415 respectively. While Tang's data was kept as is and subsets in ExCAPE-DB and Merget
416 *et al* were not, this phenomenon is not isolated and several activity values in Papyrus
417 might have originated from the same source. This over representation of the same
418 values would, in turn, bias the aggregated mean and standard deviation for specific
419 compound-target pairs.

420

421 Stereochemical aspects were discarded in Papyrus to ensure that differing molecular
422 standardisation processes of sources would not have an impact on the aggregation of
423 activity values. However, the procedure of removing stereochemistry completely
424 overlooked the potential of chiral molecules having opposing therapeutic or toxic
425 effects and does not allow for the modelling of activity cliffs.

426

427 Another related shortcoming of Papyrus is its disregard for peptides and nucleic acids,
428 even more so since one of the most abundant protein classes of level 3 are family A
429 GPCR peptide receptors. This means that, though many drugs and compounds have
430 been designed for these receptors, they do not have a single related data point in the
431 dataset. In turn, related peptide derived models will only show limited performance. In
432 a future version we would like to explore the possibility of increasing peptide
433 representation in the Papyrus dataset, but for now this is what we settled on.

434

435 In a similar vein as the datasets, the descriptors that were added are a selection of
436 descriptors that we frequently use. We believe that the provided set will be sufficient
437 for anyone investigating a specific protein (family), and that high quality results can be
438 obtained. However, we understand the need to tinker with all options of the process,
439 and we separated the descriptors from the main dataset instead of adding them
440 together. This gives the option for researchers to implement their own descriptors if
441 so desired, while keeping the format of the original Papyrus dataset.

442

443 We have provided several implementations of filters, to narrow down the data for use
444 in modeling (or perhaps other purposes). Using the entirety of Papyrus is not feasible
445 without adequate computational resources, and we recommend users to reduce the
446 data using the provided filters or in their own manner. It should be noted that the
447 quality annotation filter does not imply that only high-quality data should be used,
448 especially since classification models can leverage both the censored and binary data,
449 the latter constituting more than 95% of the dataset.

450

451 **Conclusion**

452 We created an extensive benchmark set named Papyrus, that contains high quality
453 data aggregated from multiple data sources. This standardised set is primarily used
454 as a reliable data source for modeling ligand-protein interactions. We have shown the
455 statistics of the Papyrus dataset and several classification and regression models
456 using QSAR and PCM, with performance on par with prior results. We anticipate that
457 the Papyrus dataset can be exploited in a myriad of ways and filtered or altered for
458 specific research questions. We believe the strength of the dataset lies in its
459 standardisation, normalisation and quality, while providing the necessary tools for
460 further manipulation to specific needs.

## Author's Contributions

OJMB, BJB and GJPvW conceived the study. OJMB, BJB and WJ performed the experimental work and analysis. APIJ, BvdW and GJPvW provided feedback and critical input. All authors read, commented on and approved the final manuscript.

## Data Availability

The Papyrus dataset can be found at https://doi.org/10.4121/16896406.v1. Python scripts can be found at https://github.com/OlivierBeq/Papyrus-scripts.
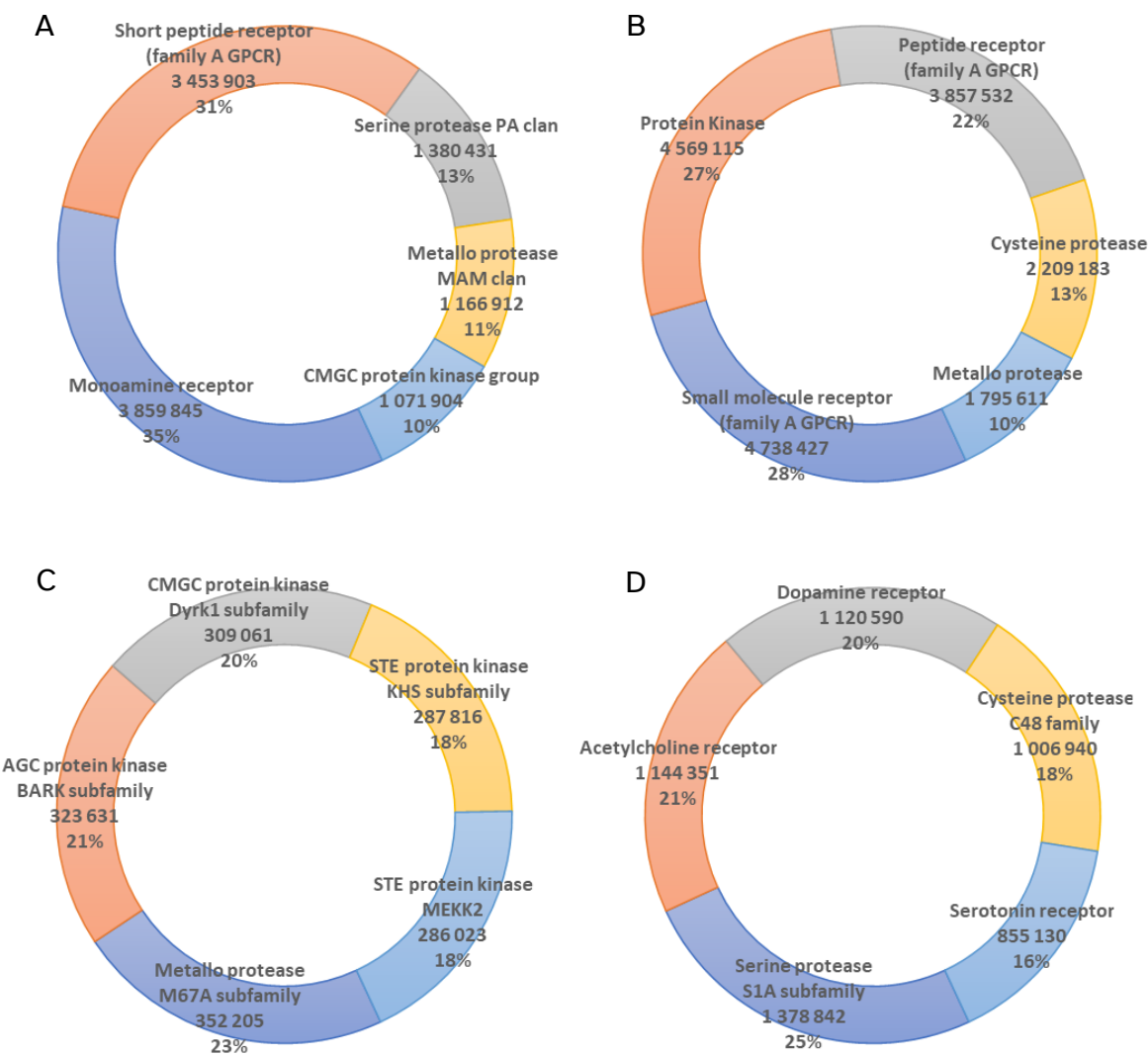
## Acknowledgements

## Funding

## Competing Interests

The authors declare that they have no competing interests.

# References

1.  Hu, Y. & Bajorath, J. Growth of Ligand–Target Interaction Data in ChEMBL Is Associated with Increasing and Activity Measurement-Dependent Compound Promiscuity. *Journal of Chemical Information and Modeling* **52**, 2550–2558 (2012).

2.  Cook, C. E. *et al.* The European Bioinformatics Institute in 2016: Data growth and integration. *Nucleic Acids Research* **44**, D20–D26 (2016).

3.  Bento, A. P. *et al.* The ChEMBL bioactivity database: an update. *Nucleic Acids Research* **42**, D1083–D1090 (2014).

4.  Gilson, M. K. *et al.* BindingDB in 2015: A public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Research* **44**, D1045–D1053 (2016).

5.  Papadatos, G., Gaulton, A., Hersey, A. & Overington, J. P. Activity, assay and target data curation and quality in the ChEMBL database. *Journal of Computer-Aided Molecular Design 2015 29:9* **29**, 885–896 (2015).

6.  Tang, J. *et al.* Making Sense of Large-Scale Kinase Inhibitor Bioactivity Data Sets: A Comparative and Integrative Analysis. *Journal of Chemical Information and Modeling* **54**, 735–743 (2014).

7.  Klaeger, S. *et al.* The target landscape of clinical kinase drugs. *Science* **358**, (2017).

8.  Lenselink, E. B. *et al.* Beyond the hype: deep neural networks outperform established methods using a ChEMBL bioactivity benchmark set. *Journal of Cheminformatics* **9**, 45 (2017).

9.  Lenselink, E. B. *et al.* Beyond the Hype: Deep Neural Networks Outperform Established Methods Using A ChEMBL Bioactivity Benchmark Set. (2019) doi:10.4121/uuid:b64986dd-3203-445e-9b93-13a5ac7ef999.

10. Béquignon, O. *et al.* Papyrus - A large scale curated dataset aimed at bioactivity predictions. (2021) doi:10.4121/16896406.v1.

512  11.  Sun, J. *et al.* ExCAPE-DB: an integrated large scale dataset facilitating Big Data
513       analysis in chemogenomics. *Journal of Cheminformatics 2017 9:1* **9**, 1–9 (2017).

514  12.  Sharma, R., Schürer, S. C. & Muskal, S. M. High quality, small molecule-activity
515       datasets for kinase research. *F1000Research* **5**, (2016).

516  13.  Christmann-Franck, S. *et al.* Unprecedently Large-Scale Kinase Inhibitor Set
517       Enabling the Accurate Prediction of Compound–Kinase Activities: A Way toward
518       Selective Promiscuity by Design? *Journal of Chemical Information and Modeling*
519       **56**, 1654–1675 (2016).

520  14.  Merget, B., Turk, S., Eid, S., Rippmann, F. & Fulle, S. Profiling Prediction of Kinase
521       Inhibitors: Toward the Virtual Assay. *Journal of Medicinal Chemistry* **60**, 474–
522       485 (2017).

523  15.  van Westen, G. J. P., Wegner, J. K., IJzerman, A. P., van Vlijmen, H. W. T. & Bender,
524       A. Proteochemometric modeling as a tool to design selective compounds and
525       for extrapolating to novel targets. *MedChemComm* **2**, 16 (2011).

526  16.  Gaulton, A. *et al.* ChEMBL: A large-scale bioactivity database for drug discovery.
527       *Nucleic Acids Research* **40**, 1100–1107 (2012).

528  17.  Kim, S. *et al.* PubChem in 2021: New data content and improved web interfaces.
529       *Nucleic Acids Research* **49**, D1388–D1395 (2021).

530  18.  The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic*
531       *Acids Research* **45**, D158–D169 (2017).

532  19.  Tweedie, S. *et al.* Genenames.org: The HGNC and VGNC resources in 2021.
533       *Nucleic Acids Research* **49**, D939–D946 (2021).

534  20.  Bento, A. P. *et al.* An open source chemical structure curation pipeline using
535       RDKit. *Journal of Cheminformatics* **12**, 51 (2020).

536  21.  Sayle, R. A. & Jack Delany. Canonicalization and Enumeration of Tautomers. in
537       *EuroMUG99* 28-29 October 1999 (1999).

538  22.  RDKit: Open-source cheminformatics (version 2021.03.5).
539       doi:10.5281/zenodo.5242603.

540 23. O'Boyle, N. M. *et al.* Open Babel: An Open chemical toolbox. *Journal of*
541 *Cheminformatics* (2011).

542 24. The Open Babel Package, version 3.0.1.

543 25. O'Boyle, N. M. Towards a Universal SMILES representation - A standard method
544 to generate canonical SMILES based on the InChI. *Journal of Cheminformatics*
545 **4**, 1–14 (2012).

546 26. Ropp, P. J., Kaminsky, J. C., Yablonski, S. & Durrant, J. D. Dimorphite-DL: An open-
547 source program for enumerating the ionization states of drug-like small
548 molecules. *Journal of Cheminformatics* **11**, 1–8 (2019).

549 27. Probst, D. & Reymond, J.-L. A probabilistic molecular fingerprint for big data
550 settings. *Journal of Cheminformatics 2018 10:1* **10**, 1–12 (2018).

551 28. Probst, D. & Reymond, J.-L. Visualization of very large high-dimensional data
552 sets as minimum spanning trees. *Journal of Cheminformatics* **12**, (2020).

553 29. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified
554 rational protein engineering with sequence-based deep representation learning.
555 *Nature Methods* **16**, 1315–1322 (2019).

556 30. Hong, H. *et al.* Mold2 , molecular descriptors from 2D structures for
557 chemoinformatics and toxicoinformatics. *Journal of Chemical Information and*
558 *Modeling* **48**, 1337–1344 (2008).

559 31. Tang, J. *et al.* Making sense of large-scale kinase inhibitor bioactivity data sets:
560 A comparative and integrative analysis. *Journal of Chemical Information and*
561 *Modeling* **54**, 735–743 (2014).

562 32. Sheridan, R. P. Time-Split Cross-Validation as a Method for Estimating the
563 Goodness of Prospective Prediction. *Journal of Chemical Information and*
564 *Modeling* **53**, 783–790 (2013).

565

566

# Additional information

A

Short peptide receptor
(family A GPCR)
3 453 903
31%

Serine protease PA clan
1 380 431
13%

Metallo protease
MAM clan
1 166 912
11%

Monoamine receptor
3 859 845
35%

CMGC protein kinase group
1 071 904
10%

B

Peptide receptor
(family A GPCR)
3 857 532
22%

Protein Kinase
4 569 115
27%

Cysteine protease
2 209 183
13%

Metallo protease
1 795 611
10%

Small molecule receptor
(family A GPCR)
4 738 427
28%

C

CMGC protein kinase
Dyrk1 subfamily
309 061
20%

STE protein kinase
KHS subfamily
287 816
18%

AGC protein kinase
BARK subfamily
323 631
21%

STE protein kinase
MEKK2
286 023
18%

Metallo protease
M67A subfamily
352 205
23%

D

Dopamine receptor
1 120 590
20%

Cysteine protease
C48 family
1 006 940
18%

Acetylcholine receptor
1 144 351
21%

Serotonin receptor
855 130
16%

Serine protease
S1A subfamily
1 378 842
25%

*Additional figure 1: Protein classification levels 3 (A), 4(B), 5 (C) and 6 (D) of targets in Papyrus.*