

1 **Small Molecule Antiviral Compound Collection (SMACC): a** 2 **database to support the discovery of broad-spectrum** 3 **antiviral drug molecules.**

4
5 Holli-Joi Martin¹, Cleber C. Melo-Filho¹, Daniel Korn¹, Richard T. Eastman², Ganesha Rai²,
6 Anton Simeonov², Alexey V. Zakharov^{2,*}, Eugene Muratov^{1,*}, Alexander Tropsha^{1,*}
7

8 ¹*UNC Eshelman School of Pharmacy, University of North Carolina, Chapel Hill, NC, USA. E-*
9 *mail: alexey.zakharov@nih.gov, murik@email.unc.edu, alex.tropsha@unc.edu*

10 ²*Division of Preclinical Innovation, National Center for Advancing Translational Sciences,*
11 *Rockville, MD, 20850*
12

13 **Abstract**

14 Diseases caused by new viruses costs thousands if not millions of human lives and trillions of
15 dollars in damage to the global economy. Despite the rapid development of vaccines for SARS-
16 CoV-2, the lack of small molecule antiviral drugs that work against multiple viral families (broad-
17 spectrum antivirals; BSAs) has left the entire world's human population vulnerable to the infection
18 between the beginning of the outbreak and the widespread availability of vaccines. Developing
19 BSAs is an attractive, yet challenging, approach that could prevent the next, inevitable, viral
20 outbreak from becoming a global catastrophe. To explore whether historical medicinal chemistry
21 efforts suggest the possibility of discovering novel BSAs, we (i) identified, collected, curated, and
22 integrated all chemical bioactivity data available in ChEMBL for molecules tested in respective
23 assays for 13 emerging viruses that, based on published literature, hold the greatest potential threat
24 to global human health; (ii) identified and solved the challenges related to data annotation accuracy
25 including assay description ambiguity, missing cell or target information, and incorrect BioAssay
26 Ontology (BAO) annotations; (iii) developed a highly curated and thoroughly annotated database
27 of compounds tested in both phenotypic (21,392 entries) and target-based (11,123 entries) assays
28 for these viruses; and (iv) identified a subset of compounds showing BSA activity. For the latter
29 task, we eliminated inconclusive and annotated duplicative entries by checking the concordance
30 between multiple assay results and identified eight compounds active against 3-4 viruses from the
31 phenotypic data, 16 compounds active against two viruses from the target-based data, and 35
32 compounds active in at least one phenotypic and one target-based assay. The pilot version of our
33 SMACC (Small Molecule Antiviral Compound Collection) database contains over 32,500 entries
34 for 13 viruses. Our analysis indicates that previous research yielded very small number of BSA
35 compounds. We posit that focused and coordinated efforts strategically targeting the discovery of
36 such agents must be established and maintained going forward. The SMACC database publicly
37 available at <https://smacc.mml.unc.edu> may serve as a reference for virologists and medicinal
38 chemists working on the development of novel BSA agents in preparation for future viral
39 outbreaks.
40

41 **Introduction**

42 Infectious diseases have had profound impacts on global human health since the beginning
43 of time. In the past two decades factors such as population growth and travel have increased the
44 rate of viral outbreaks, with a new viral threat seen nearly every year.¹ This includes the emergence
45 of severe acute respiratory syndrome coronavirus (SARS-CoV), Middle East respiratory syndrome
46 (MERS-CoV), Zika virus disease, Ebola virus disease, and variety of influenza strains (H5N1,
47 H7N9, H1N1, etc.). These viruses are just a handful of over 200 viral species annotated by the
48 International Committee for Taxonomy of Viruses as threats to human health.² The millions of
49 lives lost and trillions of dollars in damage to the global economy due to the recent pandemic
50 caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) highlight the
51 importance of medications that can offer protection against diverse viral threats that can emerge
52 in the future.³

53 It is evident from the current SARS-CoV-2 outbreak that scientific advances have
54 increased our capabilities for rapid development of new vaccines. However, none of the vaccines
55 developed thus far offers 100% protection. Furthermore, viral evolution poses a threat to further
56 decrease the vaccines efficacy, additionally, widespread hesitancy against vaccination leaves a
57 large majority of the population susceptible to the viral disease.

58 Broad-spectrum antiviral (BSA) drugs could protect against emergent viruses; however,
59 the development of such drugs has been challenging. The standard drug development and clinical
60 testing process averages 10–15 years; thus, it is impossible to see an immediate emergence of new,
61 effective drugs following an outbreak. As of today, there are only 90 approved antiviral drugs of
62 which 11 are approved to treat more than one virus.⁴ One could speculate that the ability of the 11
63 drugs to be effective against more than one virus could be explained by the conservation of their

64 targets or mechanisms of action. For example, acyclovir triphosphate competes with dGTP to
65 inhibit viral DNA polymerase activity in two human neurotropic alpha herpesviruses, herpes
66 simplex virus and varicella zoster virus.⁴ Most approved antiviral drugs are effective against
67 herpes, hepatitis, or human immunodeficiency viruses, but offer no protection against the recent
68 SARS-CoV-2 pandemic. However, the fact that medications active against more than one virus
69 do exist fuels the hypothesis that such medications can be developed in principle via concerted
70 strategic effort.

71 Despite the clear need for BSA medications, previous outbreaks have shown that the
72 interest in supporting viral research and drug discovery vanishes quickly after about a year past
73 the viral threat, leaving the work toward an effective medication unfinished.⁵ A good example is
74 the history of Paxlovid, a recently approved Pfizer medication against SARS-CoV-2. The
75 respective drug candidate was initially discovered to work against SARS in 2002-2003 by
76 inhibiting the virus' main protease (3CL-Pro) but its further development was frozen after SARS
77 vanished. When SARS-CoV-2 emerged, and it was quickly discovered that its main protease,
78 especially its active site, is almost identical to its counterpart in the original SARS, the initial drug
79 development program was restarted and Paxlovid was relatively quickly developed by Pfizer
80 through focused medicinal chemistry optimization efforts.⁶ This story clearly indicates that there
81 is a strong need for ongoing and well-funded research programs focused on the rational discovery
82 of BSA drugs.

83 More than 380 trillion different viruses exist inside the human virome, but so far only about 200
84 have been considered harmful for human health.⁷ To support focused development of BSAs and
85 learn from history, in this study we endeavored to collect, curate, and integrate all publicly
86 accessible data on compounds tested in both phenotypic and target-based assays for emerging

87 viruses of concern. To this end, we have (i) conducted a comprehensive evaluation of viruses
88 holding the greatest potential threat to global human health, (ii) used the data available in
89 ChEMBL, an online collection of bioactive molecules with drug-like properties, to build a curated,
90 annotated, and publicly available database of compounds tested in both phenotypic and target-
91 based assays for these viruses, and (iii) identified the most promising candidates with potential
92 BSA activity. We dubbed this database Small Molecule Antiviral Compound Collection (SMACC)
93 and made it publicly available online at <https://smacc.mml.unc.edu>. We expect that SMACC
94 database can support further computational and experimental medicinal chemistry studies
95 targeting rational design and discovery of novel BSAs.

96

97 **Methods**

98 **Selection of viruses of interest and initial database generation**

99 The ability of an infectious disease pathogen to cause a pandemic is impacted by several
100 intrinsic characteristics, including the mode and timing of transmission, host population
101 susceptibility, lack of effective therapeutic interventions or control measures, among others.
102 Microbial pathogens infect humans through many routes of transmission, including through animal
103 vector, fecal-oral or respiratory. Respiratory transmitted diseases are more likely to possess
104 pandemic potential, as interventions to block human-to-human transmission via aerosols are more
105 challenging to implement. The timing of disease transmission also impacts the spread of a disease,
106 if a pathogen is transmissible early in the course of disease, especially if an individual is
107 asymptomatic, this greatly facilitates potential for spread.

108

109 Viruses with a high replication rate, especially coupled with mutability of RNA and segmented
110 RNA can rapidly gain attributes, including increased transmission or evading preexisting
111 immunity, which also facilitates outbreak or pandemic spread. Viruses with high pandemic
112 potential include *Coronaviridae*, *Paramyxoviridae*, *Bunyvirales*, *Picornaviridae*, *Filoviridae*,
113 *Togaviridae*, and *Flaviviridae* virus families. Thus, we selected the following 13 viruses
114 representative of five families to query respective chemical bioactivity data in ChEMBL:
115 *Coronaviridae* (SARS-CoV-2, MERS-CoV, HCoV-229E), *Orthomyxoviridae* (H1N2, H7N7),
116 *Paramyxoviridae* (RSV, HPIV-3), *Phenuiviridae* (Sandfly Fever), and *Flaviviridae* (Dengue,
117 Zika, Yellow Fever, Powassan, West Nile).

118 All data was extracted from ChEMBL 29.⁸ The virus name, and any known alias were used as
119 keywords to extract all phenotypic and target-based assays for each virus. For the target-based
120 assays, we ran an additional search using virus and target name as the keywords to ensure no
121 respective viral data was lost. To identify drug targets for each virus we searched existing literature
122 using the keywords “ [virus_name] virus drug targets.” After extraction, the data for each virus
123 were pre-processed and curated as described below. When examining the resulting datasets, we
124 have identified a need for additional curation of assay annotations as discussed in the Results
125 section.

126

127 **Data Curation**

128 We followed protocols for chemical and biological data curation described by Fourches et
129 al.⁹⁻¹¹ In brief, specific chemotypes were normalized. Inorganic salts, organometallic compounds,
130 and mixtures were removed. Duplicate compound entries were kept to define the activity calls for
131 compounds tested against the same virus but using different assay protocols. Additionally, keeping

132 duplicates in the database can be important for analyzing the overlap of compound activity between
133 different viruses in a search for potential BSAs. However, mindful of computational modeling
134 studies that require the removal of duplicate compound entries, these entries were carefully
135 annotated in our database. All steps of data curation and integration were performed in KNIME
136 v.4.1.4¹² integrated with python v.3.7.3, RDKit v.4.2.0, and ChemAxon Standardizer v. 20.9
137 (ChemAxon, Budapest, Hungary, <http://www.chemaxon.com>). We summarized our database
138 entries after curation in **Table S1**.

139 **Identification of compounds with multiple antiviral activity**

140 A threshold of 10 μ M, irrespective of the type of activity measurement, was applied to
141 define the outcome (i.e., if a compound was active or inactive). When compound activity was
142 reported with ambiguous operators (greater than, “>”, or less than, “<”, certain value), it was
143 annotated as inconclusive. The final definition of the activity call for each compound was based
144 on the concordance of all compound replicate entries tested against the same virus but in different
145 assays (or the same viral target for the target-based dataset). Three outcomes were possible: (i) the
146 compound was active when tested in all assays; (ii) it was active in some assays and inactive in
147 others; and (iii) it was inactive in all assays. In case (i), the compound was considered active while
148 in case (iii), inactive. Any compound in case (ii), with discordant activity calls resulting from
149 different assays, i.e., with at least one activity call different from other ones for the same virus,
150 was considered inconclusive and was not used for the overlap analysis to identify compounds with
151 multiple antiviral activity. A compound was also annotated as inconclusive if the assay reported
152 the compounds activity as “Not Determined.” Finally, all compounds tested in different viruses
153 (or viral targets in the target-based dataset) were analyzed and those showing activity against two

154 or more viruses were selected as potential BSAs. Table 1 summarizes our protocols to decide on
155 the final activity calls for compounds included in SMACC database.

156

157 **Table 1. Rules for making the final activity calls for compounds in SMACC database.**

Virus	Assay	Cell Type	Activity discordance*	Standard Value	Activity Call
same	same	same	no	<10 uM or >50%	active
same	same	same	no	>10 uM or <50%	inactive
any	any	any	no	“Not Determined”	inconclusive
same	same	same	yes	Any	inconclusive
same	different	same	yes	Any	inconclusive
same	different	different	yes	Any	inconclusive

158 * At least one discordant duplicate/replicate compound

159

160 Cluster Analysis

161 The curated compound structures were submitted to hierarchical cluster analysis in KNIME
162 v.4.1.4¹² integrated with python v.3.7.3 (SciPy and Matplotlib libraries) and using RDKit
163 descriptors (RDKit v.4.2.0). The optimal number of clusters was determined by the software
164 default Euclidean distance cut-off.

165

166 Results

167 Ontological examination and curation of assays reported in ChEMBL

168 **Phenotypic assays.** While ChEMBL does an exceptional job at providing the
169 largest curated and publicly available bioactivity database, we have identified multiple issues
170 requiring additional data curation efforts to yield a clean database of antiviral activity data.
171 Specifically, we found that phenotypic assays for antiviral compounds have been annotated in
172 ChEMBL with inconsistent ontological annotation, which creates uncertainty in the data. The

173 most common finding was the inconsistent use of BioAssay Ontology annotations for the assay
174 type. As stated on the BioAssay Ontology's homepage,¹³ "The BioAssay Ontology (BAO)
175 describes chemical biology screening assays and their results including high-throughput
176 screening (HTS) data for the purpose of categorizing assays and data analysis." In practice,
177 proper BAO usage has been considered a universal best practice and highly trusted by users.
178 However, in the phenotypic assay data collection from ChEMBL, the misuse of the BAO
179 ontology was evident: for 9 of 13 viruses the assay type was recorded in ChEMBL as
180 "Organism-Based" rather than "Cell-Based". This was concerning because a virus does not meet
181 the criteria of a living organism as its life cycle relies on the host organism. Therefore, these
182 assays should be properly reported as cell-based; thus, we corrected their annotation respectively
183 in our database. The impact of this round of curation on the quality and usability of the extracted
184 data was dramatic. Indeed, in the absence of such manual analysis and correction of mis-
185 annotated data, if one were to search ChEMBL for "cell-based assays" for these viruses, 99.44%
186 (27,410 of 27,562 entries) of the data would have been uncovered. This analysis indicates a
187 critical importance of careful data processing by chemical bioactivity data curators for both the
188 accuracy of chemical structures (which has substantially improved over the years^{14,15}) and
189 correctness of activity labeling such that users can obtain the entirety of existing but effectively,
190 hidden data for which they searched.

191 Missing, i.e., absent from their designated entry field, data annotations were also extremely
192 common. For example, despite there being a distinct field for the respective entry, 13.72% of all
193 phenotypic assays results did not indicate which cell type was used. Instead, we found the records
194 of the cell type in the assay descriptions, which allowed us, in this case, to extract and properly
195 annotate this field. However, 36.73% of all missing cell types were not listed in the assay

196 description either, leaving one searching for the exact assay in the linked paper and trying to
197 identify the cell type used, which is what we had to do. This process had to be done manually,
198 which made it extremely time consuming, and, in some cases, no clear cell type could be identified.
199 If the cell type was not identified eventually, it was annotated as “unclear.” These cases are
200 reported in **Table 2** as “Cell Type Completely Missing”. Yet, this tedious work resulted in the
201 additional recovery of ~4% of all phenotypic assay results. Another issue of missing data
202 annotations was uncovered when we looked into the class of assays. Most assays were not labeled
203 to indicate whether they were primary, counter, or cytotoxicity assays. Furthermore, the assay
204 descriptions also failed to provide an appropriate level of detail. Many assay descriptions simply
205 reported “Antiviral activity against virus X.” Such descriptions are missing information on assay
206 conditions like time, substrate, equipment as well as cell type and purpose of the assay. Lacking
207 such details makes it impossible to analyze data reproducibility and prohibits meaningful
208 integration of multiple assay results. **Table 2** summarizes of our effort to procure and enrich the
209 original annotation of data found in ChEMBL.

210 **Table 2.** Curation issues of phenotypic data

Virus	# of Entries	# of Assays	Incorrect BAO Assay Annotation	Assay Type Missing in Description	Cell Type Missing	Cell Type Available in Description	Cell Type Completely Missing
SARS-CoV-2	18,190	21	18,190	2	7	0	7
MERS-CoV	49	9	49	0	49	49	0
HCoV-229E	164	11	164	7	69	65	5
Dengue	2,685	581	2,685	191	1,682	1,495	187
Yellow Fever	930	66	930	41	169	72	97
Zika	357	91	357	19	334	308	26
West Nile	514	102	514	81	308	148	160

Powassan	24	1	0	0	0	0	0
RSV	2,906	239	2,862	586	608	141	467
HPIV-3	1,632	161	1,566	435	520	96	421
H1N2	61	7	43	0	18	18	0
H7N7	26	7	26	0	18	0	18
Sandfly Fever	24	3	24	0	2	0	2

211

212 **Target based assays.** While many of the issues discussed above for phenotypic assays
213 were not present in the target-based assays, there were some cases that needed further attention.
214 One example includes 536 entries deposited as compounds tested against “genome polyprotein”
215 of West Nile or Zika viruses. However, upon closer examination of the ChEMBL records, we have
216 established that these compounds were actually tested against the NS2B-NS3 Protease, rather than
217 the entire genome polyprotein.

218 To summarize this section, when using ChEMBL as a curated¹⁶ source of data on antiviral
219 compounds, we have uncovered multiple special issues with inconsistency or mislabeling (cf.
220 Table 1) of the biological assays data. We have addressed these issues by assigning correct BAO
221 annotation to the data extracted from ChEMBL to enable the creation of a refined specialized
222 SMACC database of antiviral compounds tested in diverse antiviral assays.

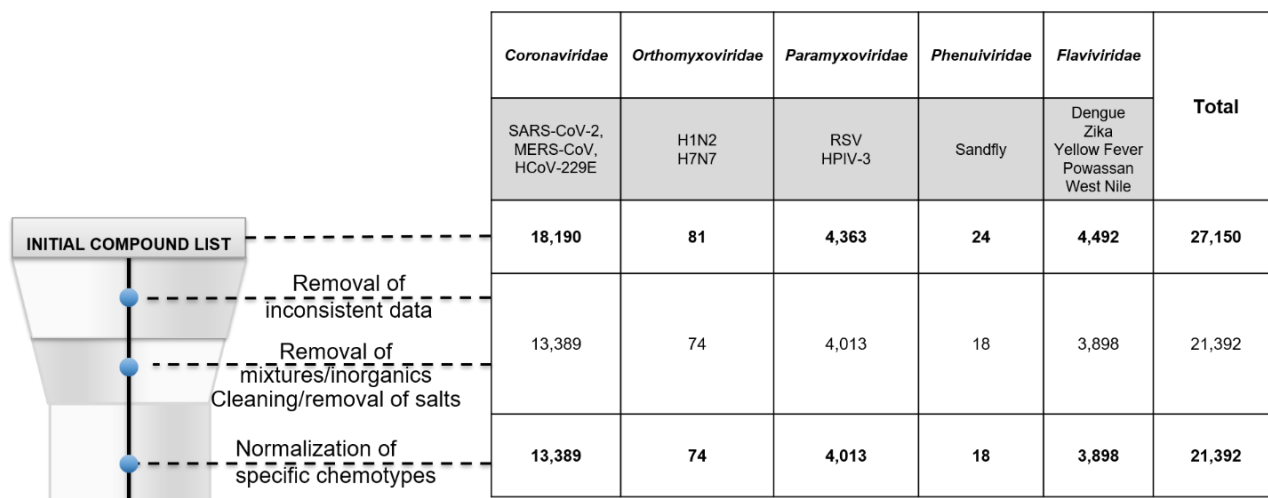
223

224 **Development of the curated data entries in the SMACC database**

225

226 Our extensive curation efforts following the protocols described in Methods led to the
227 removal of compounds from each viral family in the phenotypic data set (**Figure 1**). From the
228 initial compound list through the normalization of specific chemotypes, we removed the ~26 % of
229 compounds (n=4801) from *Coronaviridae* family, 25% (n=6) from the *Phenuiviridae*, ~21%

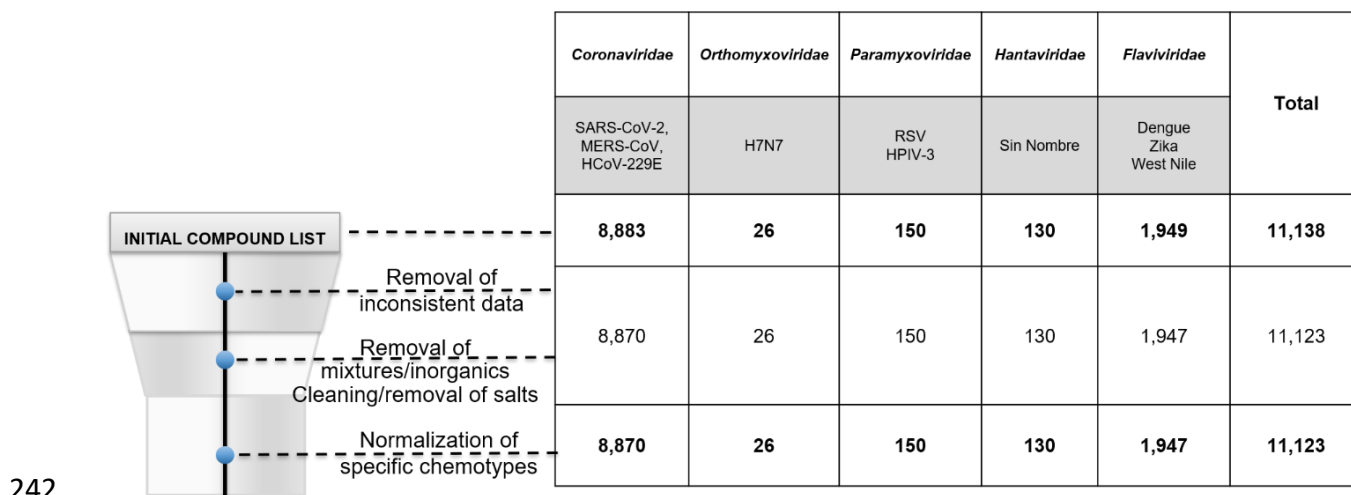
230 (n=594) from *Flaviviridae* and ~8% (n=7) from *Orthomyxoviridae* and ~8% (n=350) from
 231 *Paramyxoviridae*.



232
 233 **Figure 1.** The effect of phenotypic assay data curation on reducing the resulting dataset sizes.
 234

235
 236 For the target-based curation, there were less compounds removed due to inconsistent data
 237 and molecular cleaning (**Figure 2**). Here the family where the largest number of compounds were
 238 removed was also the *Coronaviridae* which was less than 1% of all compounds (n=13).
 239 Interestingly, our annotation efforts followed a similar pattern, where the target-based data were
 240 deposited with more annotations, therefore requiring less curation than the phenotypic data.

241



243 **Figure 2.** The effect of target-based assay data curation on reducing the resulting dataset sizes.

244 Note that at this step of data curation we intentionally kept duplicative compound records
245 reflecting our objective to check whether the same compound showed similar activity against
246 different viruses (i.e., had a potential to be a broad-spectrum agent). However, such chemically
247 duplicative entries have been annotated in SMACC to facilitate their removal prior to the
248 development of assay-specific QSAR models by the users of SMACC.

249

250 **Curated Phenotypic Data**

251

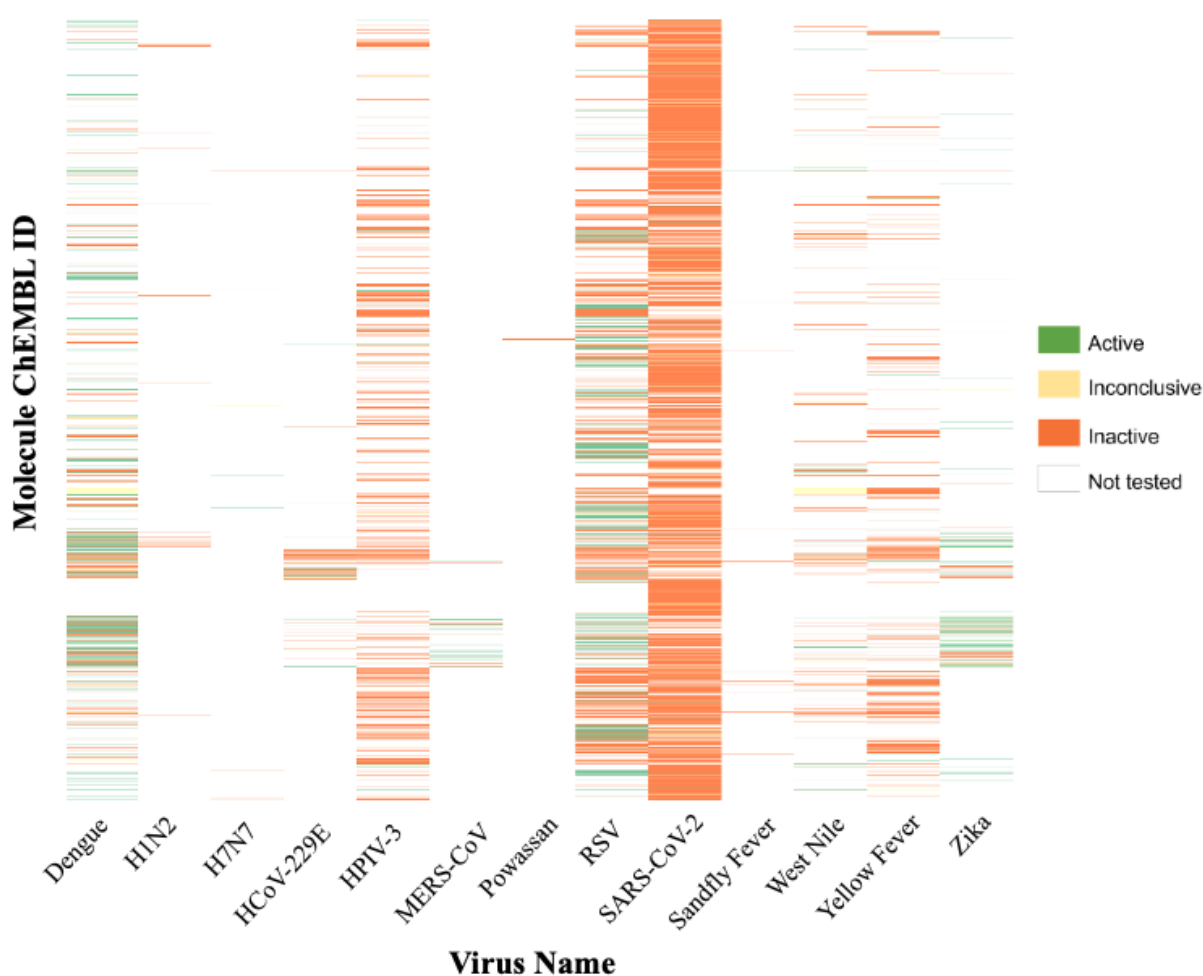
252 Curated phenotypic testing entries in our database included assay data for 13 viruses in 5
253 viral families: *Coronaviridae* (SARS-CoV-2, MERS-CoV, HCoV-229E), *Orthomyxoviridae*
254 (H1N2, H7N7), *Paramyxoviridae* (RSV, HPIV-3), *Phenuiviridae* (Sandfly Fever), and
255 *Flaviviridae* (Dengue, Zika, Yellow Fever, Powassan, West Nile).

256 **Distribution of compound activity**

257 The heatmap presented in **Figure 3** depicts the activity spectrum of all compounds tested
258 in phenotypic assays for the viruses in our database. It is evident that many compounds were either

259 inactive (80.6%) or untested. In contrast, the number of actives constituted 15.6% of the total
260 number of entries, and the fraction of “true actives” with no conflicting assay results was just
261 6.48% (1,387 compounds). Unsurprisingly, the virus with the largest number of tested compounds
262 was SARS-CoV-2 due to the many studies caused by the current pandemic, encompassing 61.6%
263 of our phenotypic assay entries. Despite the enormous testing efforts, 94.76% of compounds were
264 reported as inactive. Each virus had more inactive compounds in the dataset except Dengue
265 (Figure 3).

266



267

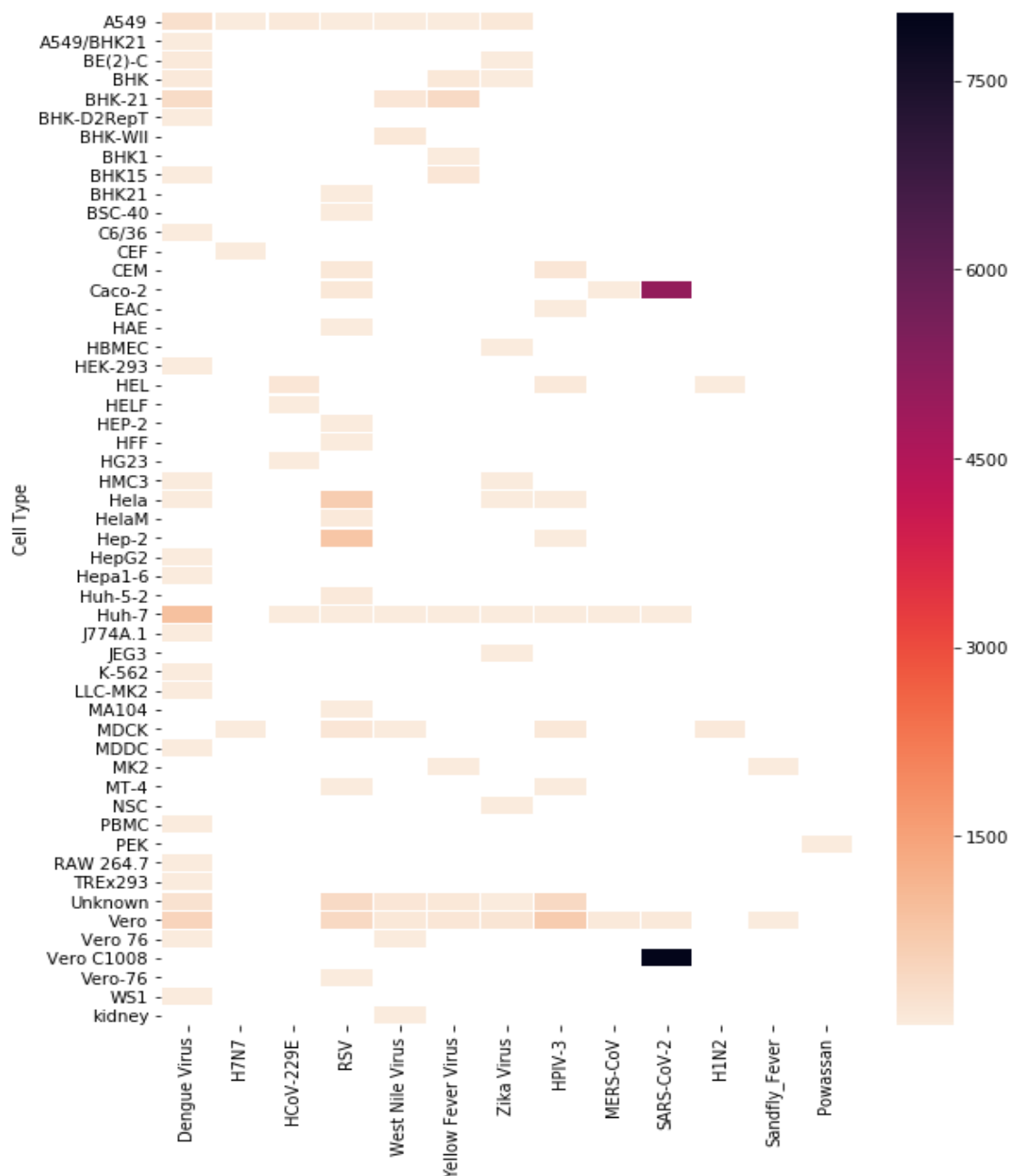
268 **Figure 3.** Activity heat map for 21,392 compounds tested in phenotypic assays for the 13 viruses

269 selected for the SMACC database.

270 **Analysis of Cell Types**

271 The 21,392 compounds integrated into our database were tested in 53 unique cell types.
272 The most common cell types were Vero C1008 (37.6% of entries), Caco-2 (24.2%), Vero (9.1%),
273 Huh-7 (4.6%), and Hep-2 (3.8%). The high propensity of testing in VeroC1008 cells is explained
274 by a single assay screening against SARS-CoV-2 (8043 entries). Other cell types, such as Caco-2,
275 were used for testing in multiple viruses and amongst various assays. Interestingly, Dengue virus
276 had the greatest number of cell types tested (26), followed by RSV (17), and Zika virus (10)
277 (**Figure 4**). Conversely, Vero cells were tested in the largest number of viruses (9) across 1,956
278 entries, followed by A549 (7 viruses), and Huh-7 cells (6 viruses).

279



280 **Figure 4.** Heat map showing distribution of compounds tested in phenotypic screens for different
281 cell types and different viruses.
282

283 **Stratifying compounds by assay type**

284 We identified 27 compounds tested in the largest number of phenotypic assays in our
285 database and further examined the effect of cell type on the resulting activity determination (Table

286 S2). As expected, there were some inconsistencies in the activities determined when stratifying
287 compounds by the virus they were tested against and the cell type for that virus and then comparing
288 their activities. The data for these 27 compounds have been recorded in a matrix with 27 respective
289 columns and 146 rows (Table S2). Unknown cell types and inconclusive activity results were
290 ignored. We identified 26 assay results when a compound was tested in the same virus and cell
291 line but showed conflicting results. In another 19 cases a compound was tested against the same
292 virus but in different cell lines, and had different results. In contrast, for 10 cases, we observed
293 completely consistent activity testing results (in 2+ entries) for a compound assayed in the same
294 virus in the same cell line. We also observed 22 cases of consistent activity when compounds were
295 tested for the same virus but in multiple cell lines. There were also many cases reporting a
296 compound tested for a single virus and multiple cell types. In this case, we only analyzed whether
297 or not the activities reported for each cell line were consistent. In summary, we observe that the
298 choice of cell type can influence the outcome of the assay, an observation reported previously¹⁷
299 and thus, the annotation of a compound as active or inactive against any virus should be always
300 reported strictly in the context of the specific underlying assay. Consequently, integration of data
301 across multiple cell line, for instance, to increase the size of the data for QSAR model
302 development, should be done with care, i.e., only when the evidence exists that compounds show
303 similar activities when tested in different cell lines.

304

305 **Identifying compounds active in multiple assays.**

306 Our efforts to extract and consolidate data on compounds tested in antiviral assays, as
307 reported in ChEMBL, revealed that identifying truly active compounds is complex and requires
308 careful curation of the available assay results. We first selected a subset of compounds from our

309 database that were tested in assays against two or more viruses. As we described in Methods, to
310 analyze the multi-viral activity we annotated each compound with one of three possible activity
311 calls. A compound was considered active if it was recorded as active when tested in all assays,
312 inactive if it was inactive in all assays, and considered inconclusive if it was active in some assays
313 and inactive in others (Table 1), if the assay result was reported with an ambiguous operator (>or
314 <), or if the compounds activity was not successfully determined by the assay. From this, we
315 created an intermediate table where each compound occupied one row, and the columns contained
316 concatenated lists of every virus the compound was reported active or inactive against. We
317 systematically analyzed this matrix to identify compounds that we considered to be true actives,
318 i.e., the compound was reported active in all assays in which it was tested.

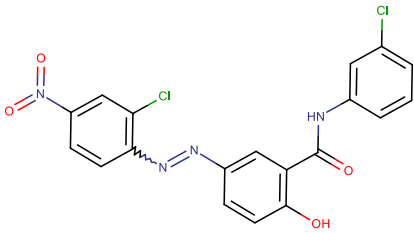
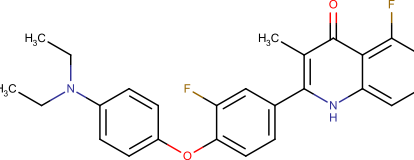
319 We report the eight most promising compounds resulting from this analysis in **Table 3**.
320 Here our top compound is CHEMBL4437334 with activity against Dengue, West Nile, Yellow
321 Fever, and Zika. It is a research compound not yet progressed into any clinical trial, which is true
322 for many compounds of this list including CHEMBL4454780 (active against RSV, MERS-CoV,
323 Dengue, and Zika), CHEMBL2016757 (active against RSV, HPIV-3, and Dengue),
324 CHEMBL4544911 and CHEMBL4562509 (active against Dengue, West Nile, and Yellow Fever).
325 Three named compounds were identified from our search: 6-azauridine (active against RSV, West
326 Nile, Dengue); amodiaquine (active against SARS-CoV-2, Dengue, Zika), which is an approved
327 drug for malaria; and brequinar (active against Dengue, West Nile, and Yellow Fever), which is
328 currently in Phase I clinical trials for treatment of acute myeloid leukemia.¹⁸⁻²⁰

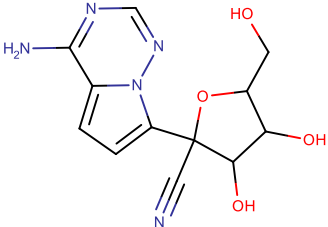
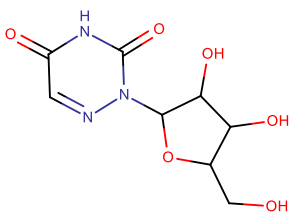
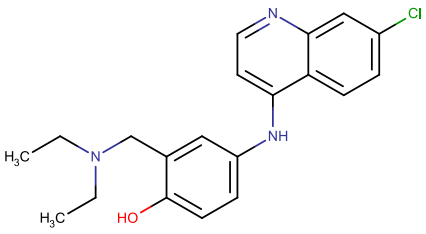
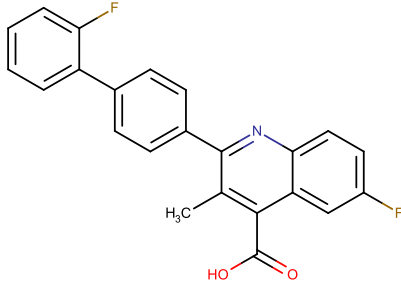
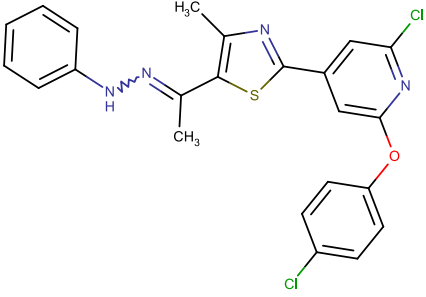
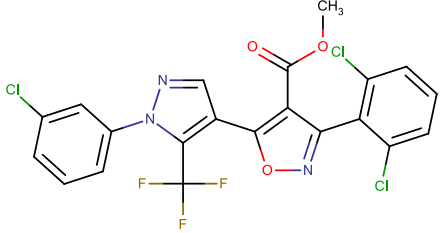
329 Interestingly, brequinar also recently underwent phase II clinical trials against SARS-CoV-
330 2.²¹ While the clinical trial was not successful using brequinar alone, research on brequinar drug
331 combinations have found this drug highly effective in combination with remdesivir or

332 molnupiravir.²² Research suggests that brequinar's antiviral activity against SARS-CoV-2 is
333 through inhibition of the host cell dihydroorotate dehydrogenase (DHODH) rather than being a
334 direct acting antiviral. The combination of a nucleobase antiviral and DHODH, or other compound
335 that impacts de novo nucleotide synthesis would in effect increase the nucleobase antiviral cellular
336 concentration thereby increasing the rate of incorporation into the viral synthesized RNA, in
337 theory. This approach has been shown to be effective against multiple viruses *in vitro*, for
338 example, brequinar has been shown to inhibit dengue, enterovirus, and Ebola viruses through this
339 same, host-targeted mechanism.²³⁻²⁵ Given the reported activity of brequinar against three
340 *flaviviruses* (Dengue, West Nile, Yellow Fever) in phenotypic assays, we hypothesize the assays
341 were detecting the human dihydroorotate dehydrogenase inhibition, rather than inhibition of a viral
342 target. As such, a future release of SMACC will include an analysis of all phenotypic assays, host-
343 target assays, and the overlap between the phenotypic assays and the host-target assays. It is our
344 hope this analysis will help develop hypotheses for, and identify, potential host-targeting broad-
345 spectrum antiviral drugs.

346

347 **Table 3.** Example of compounds active in multiple viruses.

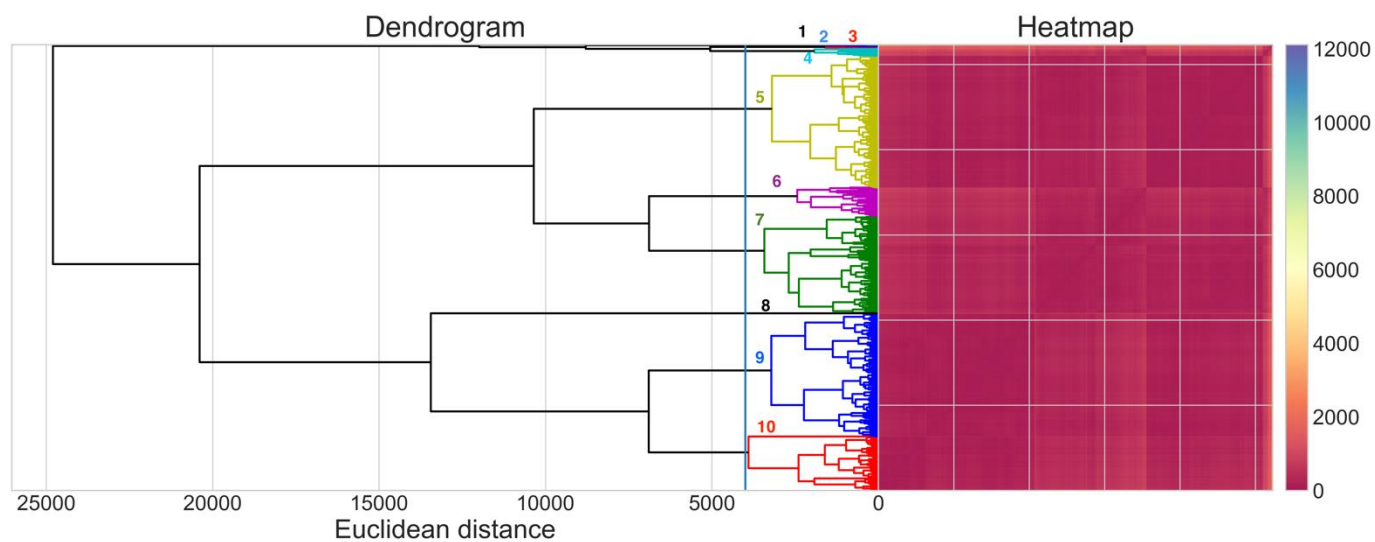
ChEMBL ID	Structure	Active against	Inactive against
CHEMBL4437334		Dengue, West Nile, Yellow Fever, Zika	
CHEMBL4454780		RSV, MERS-CoV, Dengue, Zika	

CHEMBL2016757		RSV, HPIV-3, Dengue	West Nile
CHEMBL564201 (6-Azauridine)		RSV, West Nile, Dengue	Yellow Fever
CHEMBL682 (Amodiaquine)		SARS-CoV-2, Dengue, Zika	
CHEMBL38434 (Brequinar)		Dengue, West Nile, Yellow Fever	SARS-CoV-2
CHEMBL4544911		Dengue, West Nile, Yellow Fever	
CHEMBL4562509		Dengue, West Nile, Yellow Fever	

349 There were 55 more compounds active against at least two viruses (**Table S3**). Of these,
350 seven compounds were active against different viral families: three were active against
351 *Paramyxoviridae* (RSV, HPIV-3); two were active against *Coronaviridae* (MERS-CoV, HCoV-
352 229E); and 43 against any two of our *Flaviviridae* viruses (Dengue, Zika, Yellow Fever, Powassan,
353 and West Nile). We also identified 1,324 compounds active against one virus.

354 **Cluster analysis of active compounds**

355 The structural clustering of all compounds tested in phenotypic cell-based assays revealed ten
356 clusters (**Figure 6**). The top BSA compounds, CHEMBL4437334 and CHEMBL4454780, active
357 against four different viruses, are in clusters #7 and #5, respectively. The subcluster containing
358 CHEMBL4437334 (cluster #7) has 847 compounds. Among them, some nearest neighbors of
359 CHEMBL4437334 (**Figure 7**) were active against SARS-CoV-2 and could be further tested
360 against a panel of flaviviruses (Dengue, West Nile, Yellow Fever, and Zika). The subcluster of
361 CHEMBL4454780 (cluster #5) contains 1,406 compounds; nearest neighbors of
362 CHEMBL4454780 are presented in **Figure 7**. Compounds CHEMBL1197690,
363 CHEMBL3581155, and CHEMBL7568 were active against one or two flaviviruses and could be
364 further tested against additional flaviviruses and viruses from other families like RSV and MERS-
365 CoV. Likewise, CHEMBL4303559 was only tested and active against SARS-CoV-2 and could be
366 tested against members of *Flaviviridae* and other coronaviruses such as MERS-CoV.



367

368 **Figure 6.** Clustering of compounds from the phenotypic cell-based assays by chemical structure.

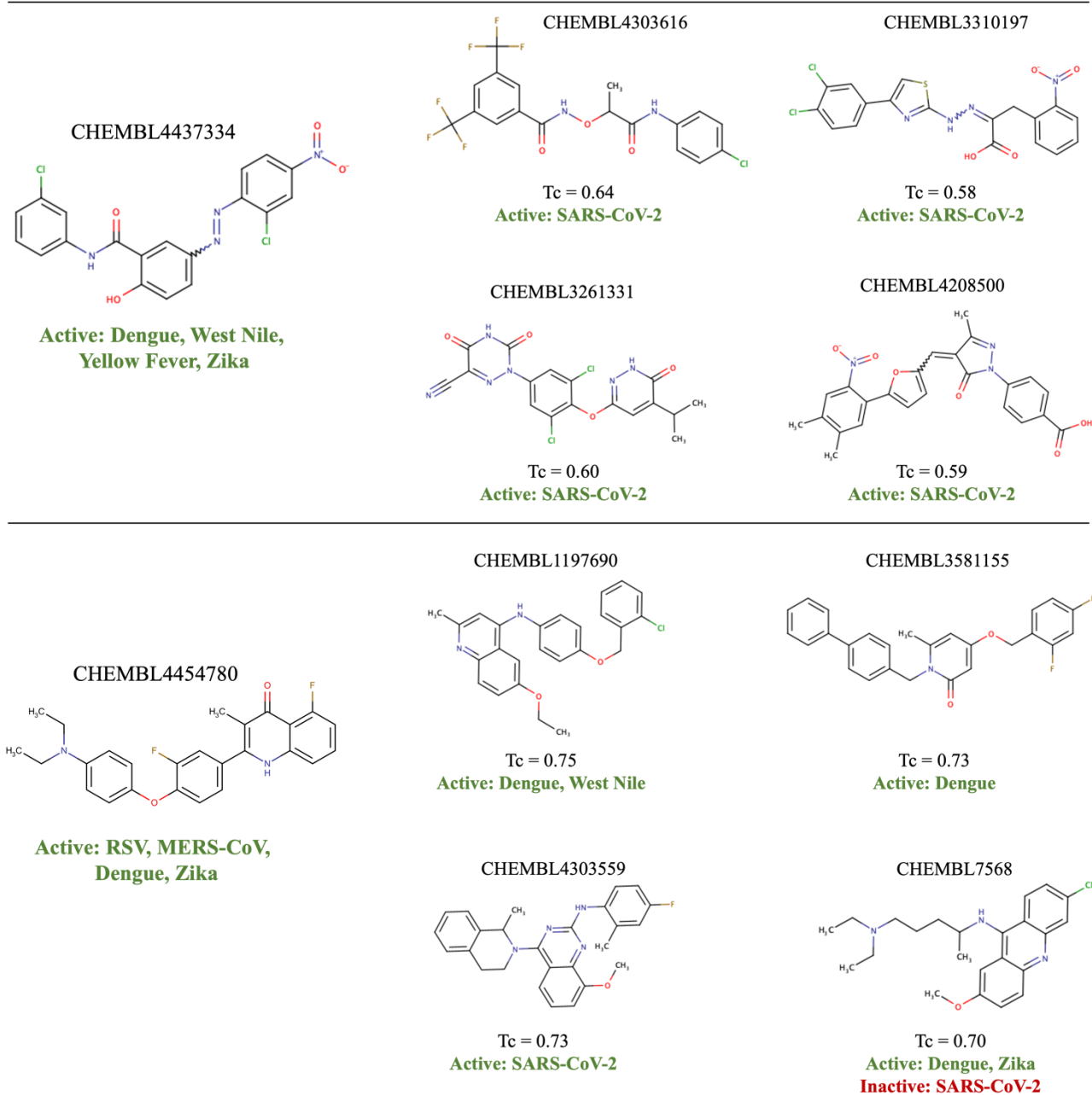
369 The colors from the heatmap are based on the Euclidean distances between compounds. Colors

370 nearer to dark red indicate a shorter distance between molecules.

371

Active against multiple viruses

Nearest neighbors inside subcluster



372

373 **Figure 7.** Examples of compounds similar to experimental broad-spectrum hits that could be

374 further tested against multiple viruses of interest.

375

376

377

378 **Target-Based Data**

379 Curated target-based testing entries (11,123) in our database include assay data for ten
380 viruses in five viral families: *Coronaviridae* (SARS-CoV-2, MERS-CoV, HCoV-229E),
381 *Orthomyxoviridae* (H7N7), *Paramyxoviridae* (RSV, HPIV-3), *Phenuiviridae* (Sin Nombre), and
382 *Flaviviridae* (Dengue, Zika, West Nile).

383

384 **Activity of Compounds**

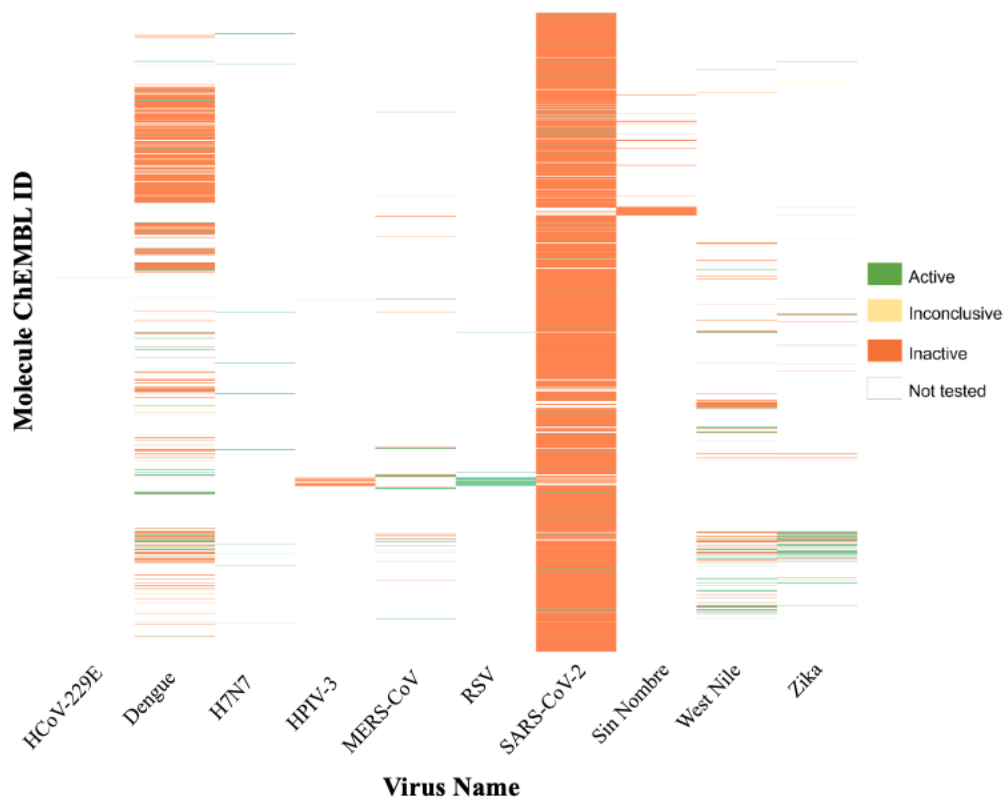
385 Using the activity calls based on the assay results as defined in Methods, most compounds (89.8%)
386 were inactive (**Figure 8**). Many compounds inactive against SARS-CoV-2 were tested because of
387 the recent multiple testing campaigns including drug repurposing screenings due to current
388 pandemic. While SARS-CoV-2 was the most tested virus, three flaviviruses were also well
389 represented in the database (Dengue, West Nile, and Zika); however, as the number of compounds
390 tested increased there was a decrease in the fraction of the active compounds for these viruses.
391 Overall, active compounds represented only ~9.9% of our total dataset where 5.78% (644
392 compounds) were “true actives”.

393

394

395

396



397
398 **Figure 8.** Activity heat map for 11,123 compounds tested in target-based assays for the 10 viruses
399 selected for the SMACC database.

400
401 **Analysis of Targets**
402 The Main Protease (3CLpro) of *Coronaviridae* was, unsurprisingly, the most studied target
403 (78.8% of the entries), followed by NS2B-NS3 Protease of *Flaviviridae* (16.2%), NS5 of
404 *Flaviviridae* (1.26%), Integrin alpha-V/beta-3 of *Hantaviridae* (1.17%), and Fusion glycoprotein
405 F0 of *Paramixoviridae* (1.1%). Interestingly, the virus with the greatest number of targets tested
406 (five) was MERS-CoV and was tested against the spike protein, RDRP, Nucleocapsid protein,
407 M^{pro}, and PL^{pro}.

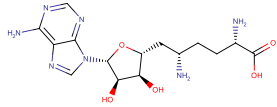
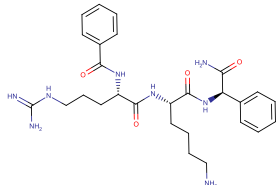
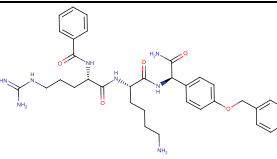
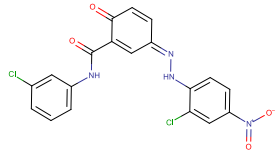
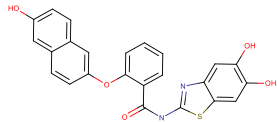
408
409

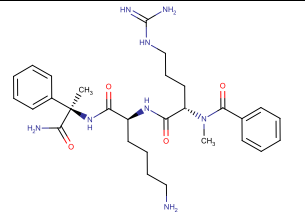
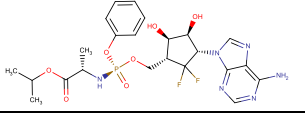
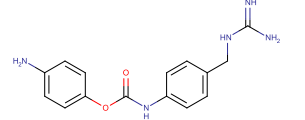
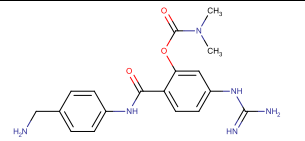
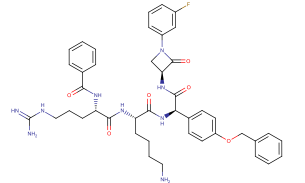
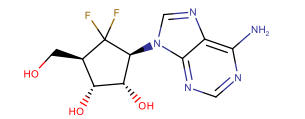
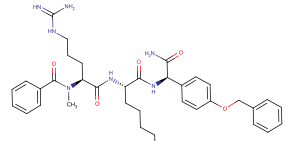
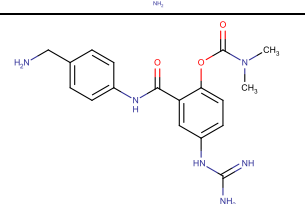
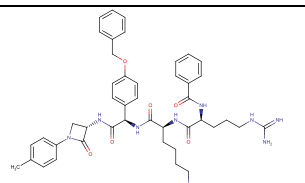
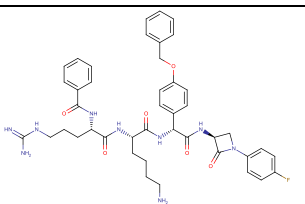
410 Analysis of Compounds

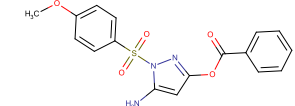
411 We followed the same approach for analyzing the target-based dataset for BSA activity, as was
412 taken for the phenotypic dataset. In this case, the intermediate table included a row for each
413 compound, and the columns were concatenated lists of every virus and target the compound was
414 reported active or inactive against. Our analysis identified 16 compounds active against two
415 viruses at the protein target level (**Table 4**). Two of these compounds (CHEMBL4544781 and
416 CHEMBL4522602) were active against targets from two different viral families (Zika's NS5 and
417 MERS-CoV's RDRP), whereas the others were active against two flaviviruses NS2B-NS3
418 Protease. We also identified 628 compounds active against one virus.

419

420 **Table 4.** Compounds active against different viruses in target-based assays.

Compound Name	Structure	Target	Virus
CHEMBL1214186 ^a		NS5	Dengue, Zika
CHEMBL3740277		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL3741422		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4437334		NS2B-NS3 Protease	Dengue, Zika
CHEMBL4440832		NS2B-NS3 Protease	Dengue, Zika

CHEMBL4474101		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4522602		NS5, RDRP	Zika, MERS-CoV
CHEMBL4531546		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4536920		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4537775		NS5, RDRP	Zika, MERS-CoV
CHEMBL4544781		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4545026		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4563372		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4568434		NS2B-NS3 Protease	Dengue, West Nile
CHEMBL4576745		NS2B-NS3 Protease	Dengue, West Nile

CHEMBL569561		NS2B-NS3 Protease	Zika, West Nile
--------------	---	-------------------	-----------------

421 ^a Inactive against SARS-CoV-2
422

423 Structural clustering of all compounds revealed 11 clusters (**Figure 9**).

424 CHEMBL4544781 was active against targets from different viral families (NS5 of Zika

425 Virus and RDRP of MERS-CoV) and is in cluster #7 along with 867 other compounds;

426 nearest neighbors of CHEMBL4544781 are presented in **Figure 10**. CHEMBL1630221

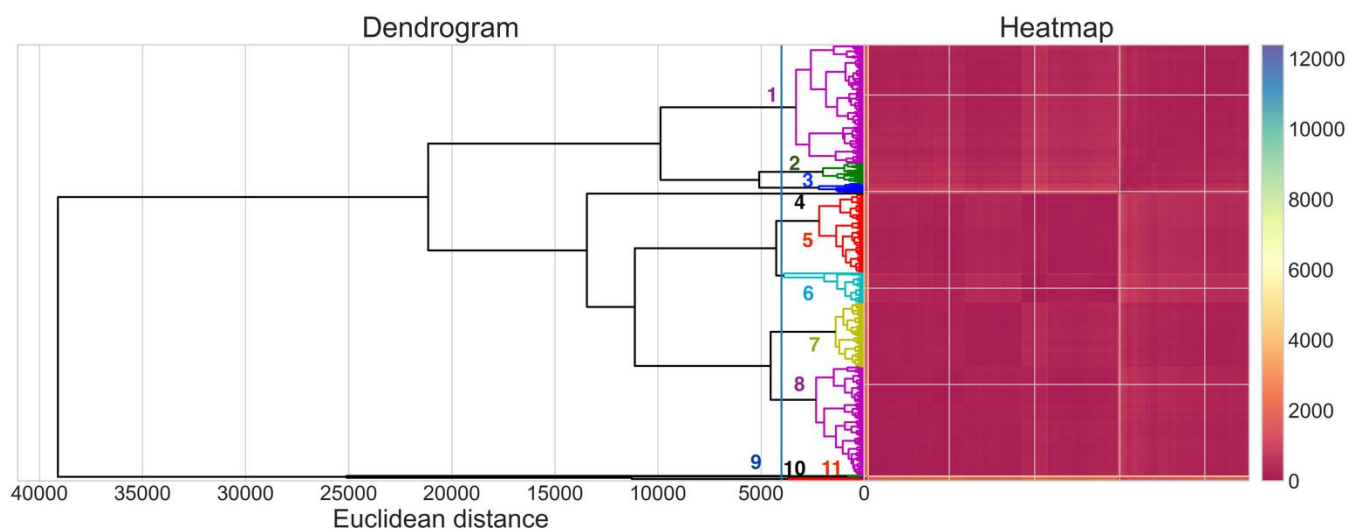
427 was only tested (and active) against the NS5 Polymerase of Zika Virus and could be further

428 tested against other polymerases from other flaviviruses and the RNA-Dependent RNA

429 Polymerase (RdRP) of MERS-CoV. Three other nearest neighbors of CHEMBL4544781

430 were only tested (and active) against SARS-CoV-2 Main Protease (M^{Pro}). These

431 compounds could be further tested against polymerases of Zika and MERS-CoV.



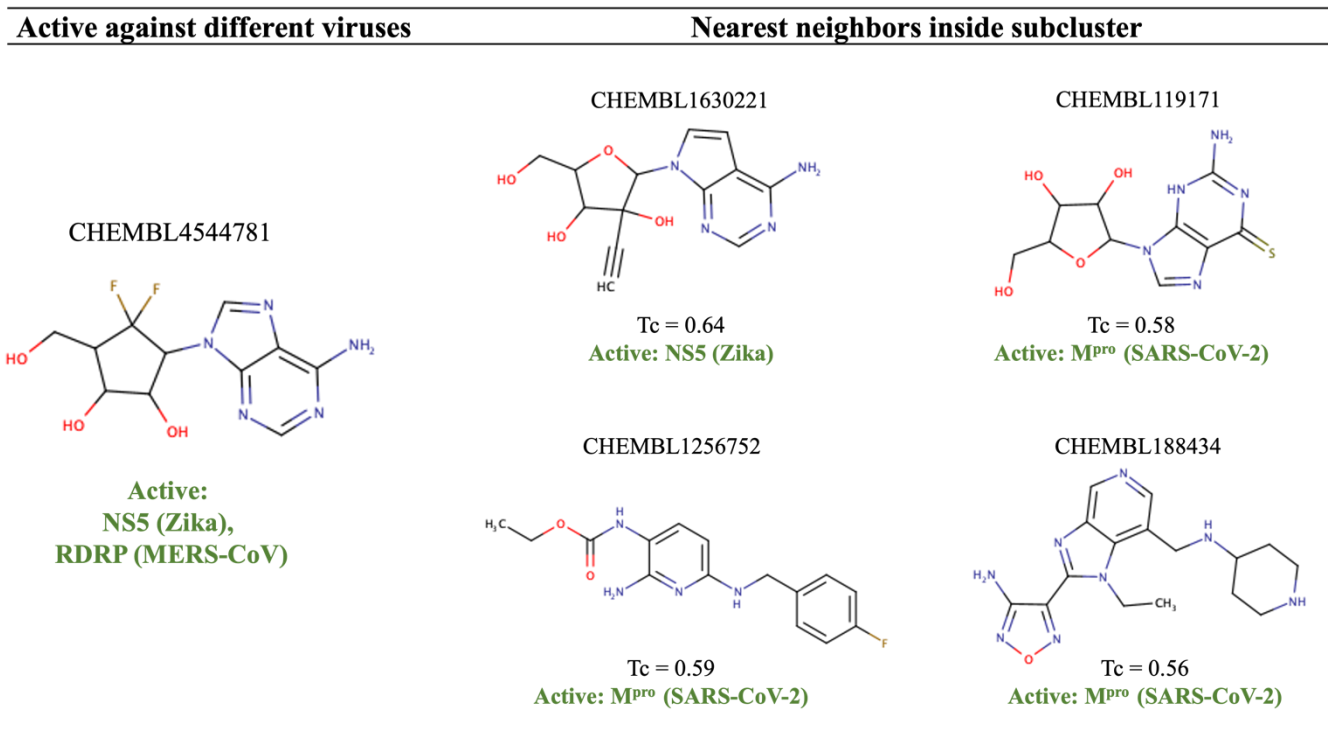
432

433 **Figure 9.** Clustering analysis of compounds from the target-based assays. The colors from the

434 heatmap are based on the Euclidean distances between compounds. Colors nearer to dark red

435 indicate a shorter distance between molecules.

436



437

438 **Figure 10.** Examples of compounds that could be further tested against different viral targets of
439 interest due to their chemical similarity to an active molecule with multiple antiviral activity.

440

441 **Concordance between the phenotypic and target-based data**

442 We analyzed the concordance between the 5,934 compounds tested in both phenotypic and
443 target-based assays to: (i) expand our list of hits by identifying potentially promising compounds
444 that may not have been tested yet, and (ii) hypothesize the mechanism of actions of compounds
445 active in a virus in a live cell and a complementary viral target. Our systematic analysis of the
446 assay results indicated that 35 compounds were active in at least one phenotypic and one target-
447 based assay (**Table S4**, Supplementary Material). In many cases, the active calls were within the
448 same viral family. For example, CHEMBL4522006 was active against Dengue Virus in a
449 phenotypic assay, and active against the Dengue NS2B-NS3 Protease in a target-based assay. Our

450 data strongly supports the hypothesis that CHEMBL4522006 is active against Dengue virus in the
451 live cell assay by inhibiting its NS2B-NS3 Protease, which supports the use of the protease assays
452 for future experimental and computational structure-activity relationship studies. Promising
453 potential BSA compounds, including CHEMBL4522006, are summarized in **Table 5**.

454 In other cases, as we observed for CHEMBL4437334, a compound was active against
455 several viruses of the same family in phenotypic assays (Dengue Virus, West Nile Virus, Yellow
456 Fever Virus, Zika Virus) and only tested and active against a subset of those viruses in the target-
457 based assays (NS2B-NS3 Protease of Dengue and Zika). In these cases, we could suggest the
458 compound be tested against the same target in the untested yet highly homologous viruses from
459 the same family, using the principle of viral protein conservation.³

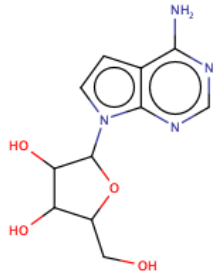
460 There were also instances of compounds, such as CHEMBL267099 (tubercidin), reported
461 active in a phenotypic assay for a virus in one family (HPIV-3) and active in a target-based assay
462 of another (Zika NS5 protein). Cases like these are particularly interesting, because after making
463 this connection one can suggest testing this compound in various HPIV-3 targets, live cell Zika
464 virus, as well as the other highly homologous *Flaviviridae* members (West Nile, Yellow Fever,
465 Dengue) in live cell assay and against the NS5 protein.

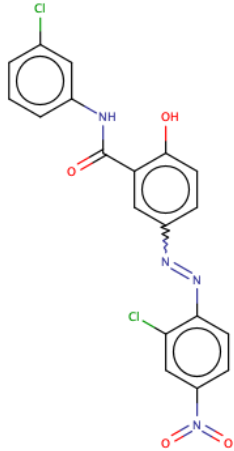
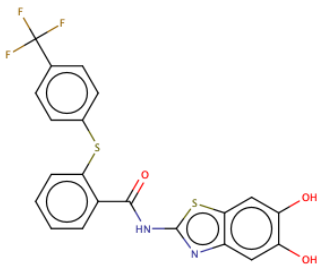
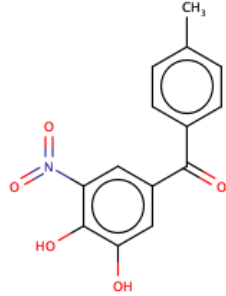
466 Our concordance analysis also revealed 52 compounds active in at least one phenotypic
467 assay and inactive in a (supposedly, relevant) target-based assay (**Table S5**, Supplementary
468 Material). In this case, we recommend compounds be tested in additional target-based assays of
469 the same viral family; it is also possible that the activity of such compounds inactive in viral
470 targeting assays but active in phenotypic assays is actually due to their host-directed mechanism
471 of action. There were also 191 compounds inactive in a phenotypic assay and active in at least one
472 target-based assay (**Table S6**, Supplementary Material). Mapping the virus and viral family of the

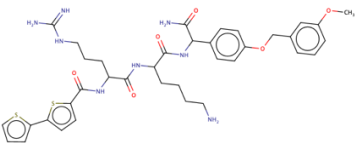
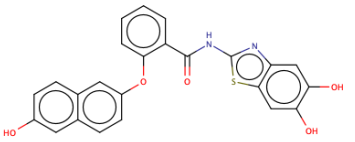
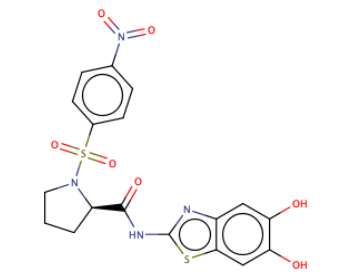
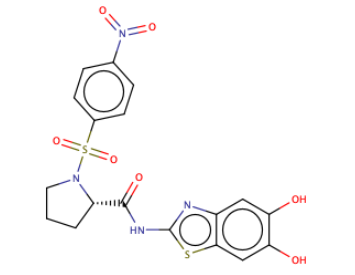
473 phenotypic result to the active result of the target helped identify potential new viral families for
474 phenotypic testing, as well as highlighted the importance whether the cell type used in the
475 phenotypic assay appropriately represented the virus and the antiviral result. Of course, due to the
476 proportion of inactive compounds in our dataset, most compounds (5,656) were concordantly
477 reported as inactive in both phenotypic and target-based assays.

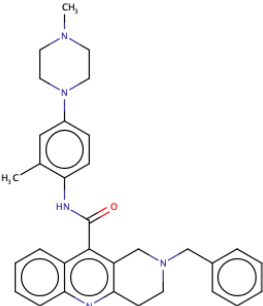
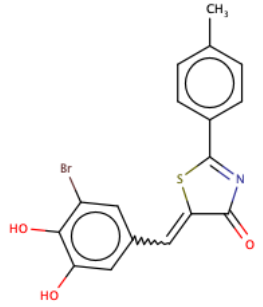
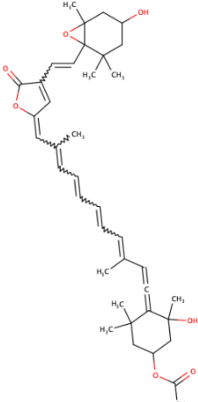
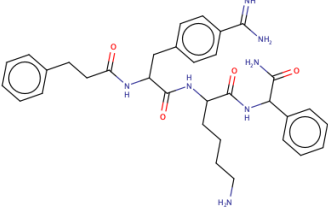
478

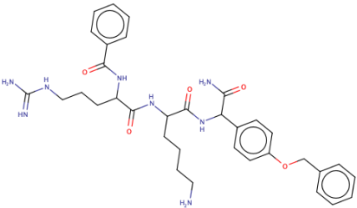
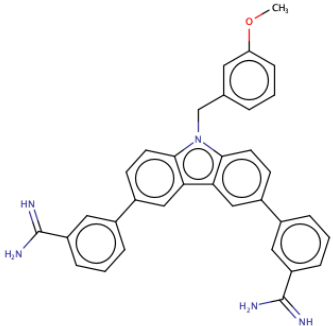
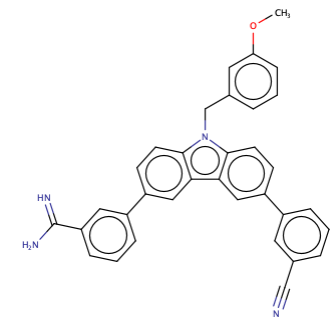
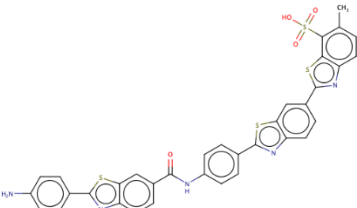
479 **Table 5.** Selection of compounds nominated for experimental testing as potential BSA agents

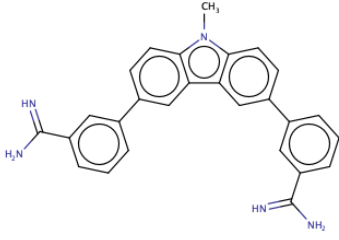
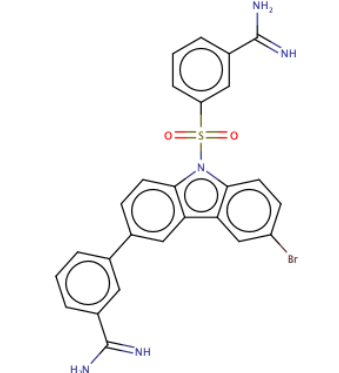
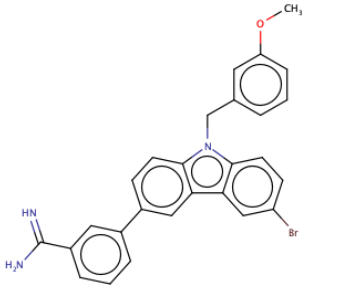
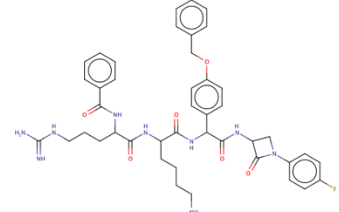
CHEMBL ID	Structure	Target active against	Virus(es) active against in cell-based assays	Suggested target-virus combination for testing	Reasoning for testing
CHEMBL267099		NS5 (Zika)	HPIV-3	Zika and other flaviviruses (cell-based assays) NS5 homologs in other flaviviruses (target-based)	Cell-based assays against Zika and other flaviviruses to confirm NS5 inhibition as a possible mechanism of action Active against NS5 of Zika that could be tested against NS5 homologs from other flaviviruses

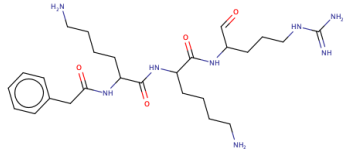
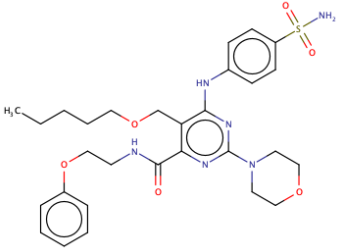
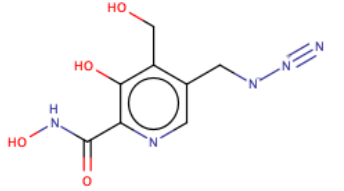
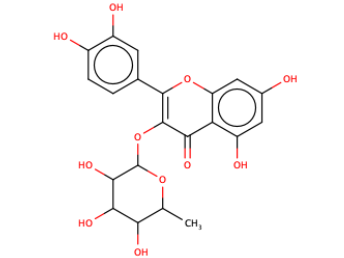
				A panel of HPIV-3 targets	Explore possible mechanism of action for HPIV-3 inhibition observed in a cell-based assay
CHEMBL4437334		NS2B-NS3 Protease (Dengue, Zika)	Dengue Virus, West Nile Virus, Yellow Fever Virus, Zika Virus	NS2B-NS3 homologs from other flaviviruses (West Nile, Yellow Fever)	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4522006		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL1324		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3

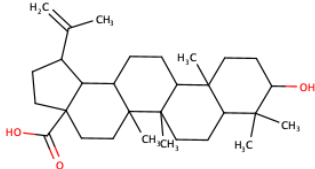
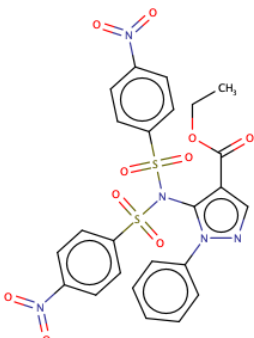
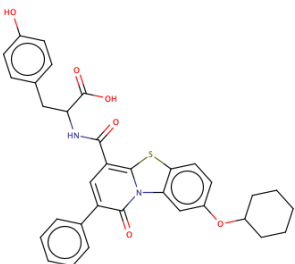
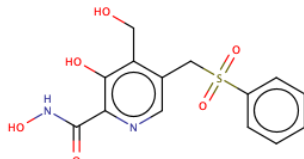
CHEMBL3741713		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4440832		NS2B-NS3 Protease (Dengue, Zika)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4462325		NS2B-NS3 Protease (Zika)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4583315		NS2B-NS3 Protease (Zika)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3

CHEMBL1370977		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL1458891		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL1980535		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL3628278		NS2B-NS3 Protease (West Nile)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3

CHEMBL3741422		NS2B-NS3 Protease (Dengue, West Nile)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4446364		NS2B-NS3 Protease (Zika)	Zika Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4447165		NS2B-NS3 Protease (Zika)	Zika Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4447800		NS2B-NS3 Protease (Dengue)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3

CHEMBL4448497		NS2B-NS3 Protease (Zika)	Zika Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4457232		NS2B-NS3 Protease (Zika)	Zika Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4532866		NS2B-NS3 Protease (Zika)	Zika Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4576745		NS2B-NS3 Protease (Dengue, West Nile)	Dengue Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3

CHEMBL522355		NS2B-NS3 Protease (West Nile)	West Nile Virus	NS2B-NS3 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS2B-NS3
CHEMBL4454990		NS5 (Zika)	Dengue Virus, Zika Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS5
CHEMBL4439416		NS5 (Zika)	Zika Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS5
CHEMBL82242		NS5 (Dengue)	Dengue Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target- based assays support the hypothesis of activity against flaviviruses by targeting NS5
CHEMBL269277		NS5 (Dengue)	Dengue Virus	NS5 homologs from other flaviviruses	Cell-based and target- based assays support the hypothesis

				Other flaviviruses in cell-based assays	of activity against flaviviruses by targeting NS5
CHEMBL4449109		NS5 (Dengue)	Dengue Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS5
CHEMBL4526128		NS5 (Dengue)	Dengue Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS5
CHEMBL4587069		NS5 (Dengue)	Dengue Virus	NS5 homologs from other flaviviruses Other flaviviruses in cell-based assays	Cell-based and target-based assays support the hypothesis of activity against flaviviruses by targeting NS5

480

481

482

Integrated, searchable SMACC Database

483

We have described above various protocols for curating data of interest to the antiviral drug

484

discovery from ChEMBL database. The resulting SMACC database currently exists as a

485

searchable Excel spreadsheet that we include for public with this paper. This spreadsheet allows

486 multiple approaches to identify compounds of interest. The approach described here in Methods,
487 i.e., removing compounds with conflicting activity calls from our final BSA activity analysis, was
488 stringent and resulted in a concise list of potential BSA compounds that we had the highest
489 confidence in. While the filtering criteria described above was appropriate to achieve our project's
490 objective, we acknowledge the value of extracting different subsets of the SMACC database using
491 different criteria depending on the study objectives, and SMACC database (even in the form of an
492 Excel spreadsheet) enables multiple analyses.

493 For instance, in contrast with the approach described above, another method for identifying
494 BSA compounds would be to consider all compounds with at least one active assay result against
495 two or more viruses. This would increase the number of compounds considered in the analysis
496 because inconclusive entries would not be removed as described above. To do this, we created a
497 subset of the SMACC database with all compound entries reported as active in phenotypic assays,
498 enumerated the number of viruses the compounds were reported active against, and removed all
499 compounds reported active against only one virus. This approach resulted in 21 new hit compounds
500 identified from phenotypic assays (**Table S7**) and 10 new hit compounds from target-target based
501 assays that were not identified in the previous approach (**Table S8**).

502 As mentioned above, there are currently only 90 antiviral drugs approved for treating nine
503 human infectious diseases.⁴ We utilized SMACCs filtering tools to enumerate their presence in
504 our database. Currently, SMACC includes chemogenomics data for RSV and two strains of human
505 influenza virus (H1N2 and H7N7), which covers only two of nine diseases with approved drugs.
506 Despite this, we identified 53 of 90 approved drugs in our phenotypic dataset and 57 of 90 in our
507 target-based dataset (**Table S9**). The compounds with reported active assay results are summarized
508 in **Table 6**. Clearly, these drugs have broader activity than they are approved for. Further

509 experimental testing based on hypotheses from this table will be extremely valuable to
 510 understanding their broad-spectrum potential.

511

512 **Table 6.** Approved drugs with active assay results found in SMACC.

Compound Information					Active Assay Results in SMACC	
Drug name	Brand name	Approved clinical use	Inhibitory MOA	ChEMBL ID	Phenotypic Assay	Target-Based Assay
Simeprevir	Olysio®	HCV	NS3/NS4B Protease	CHEMBL501849		H7N7 Matrix Protein 2
Asunaprevir	Sunvepra®	HCV	Protease	CHEMBL2105735		H7N7 Matrix Protein 2
Sofosbuvir	Sovaldi®	HCV	NS5B	CHEMBL1259059	Dengue Virus	
Ribavirin	Copegus®	HCV, RSV, fever	RdRp	CHEMBL1643	HPIV-3, Sandfly Fever, Dengue Virus, Yellow Fever Virus, RSV	
Lopinavir	Kaletra®	HIV	Protease	CHEMBL729	SARS-CoV-2	
Nelfinavir	Viracept®	HIV	Protease	CHEMBL584	Dengue Virus	
Raltegravir	Isentress®	HIV	Integrase	CHEMBL254316		H7N7 Matrix Protein 2
Elvitegravir	Vitekta®	HIV	Integrase	CHEMBL204656		SARS-CoV-2 3CLpro
Atazanavir	Reyataz®	HIV	Protease	CHEMBL1163		SARS-CoV-2 3CLpro
Rilpivirine	Edurant®	HIV-1	Nonnucleoside reverse transcriptase	CHEMBL175691	SARS-CoV-2	
Podofilox	Condylox®	HPV-related diseases	Cytotoxicity/cell division	CHEMBL61	SARS-CoV-2	
Trifluridine	Viroptic®	HSV	Viral and cellular DNA synthesis	CHEMBL1129	HPIV-3	
Idoxuridine	Dendrid®	HSV-1	Viral and cellular DNA synthesis	CHEMBL788		Zika NS5
Acyclovir	Zovirax®	HSV, VZV	Viral DNA polymerase	CHEMBL184		H7N7 Neuraminidase
Zanamivir	Relenza®	Influenza A and B	Neuraminidase	CHEMBL222813	H7N7	

513

514 Our pilot-SMACC database is currently available at <https://smacc.mml.unc.edu>. Users will

515 find freely downloadable excel sheets containing our phenotypic, target-based, and overlapping

516 datasets including tabs containing subsets of active compounds selected from the approach

517 described in Methods. These excel sheets were designed so that users could easily extract subsets
518 of the database using various filtering options. These filters include molecule (ChEMBL ID,
519 smiles, InChiKey), virus, cell or target type, activity (activity call, raw assay result), and assay
520 type. We acknowledge the widely varied objectives across antiviral research and emphasize the
521 versatility of this database.

522 **Discussion**

523 We have collected, curated, and integrated all the chemogenomic data available for a subset
524 of viruses of interest in ChEMBL to identify BSA compounds. We created a pilot version of the
525 SMACC database based on ChEMBL data. This initial data collection and curation effort can guide
526 future data collection to increase the clarity and accessibility of relevant information to a broader
527 scientific community include additional data on other emerging viruses. SMACC database adds to
528 a variety of other important datasets and databases such as SARS-CoV-2 Data Resource by
529 PubMed (<https://www.ncbi.nlm.nih.gov/sars-cov-2/>), a comprehensive COVID-19 Data Portal by
530 European Bioinformatics Institute (EBI) (<https://www.covid19dataportal.org>), a collection of over
531 20,000 screening results for compounds tested against SARS-CoV-2 in a special release of the
532 ChEMBL database (<https://www.ebi.ac.uk/chembl/>), a portal of target specific and phenotypic
533 screening results of chemical libraries in SARS-CoV-2 established by NCATS at the NIH²⁶
534 (<https://opendata.ncats.nih.gov/covid19/index.html>), and COVID-specific tools and collections
535 like CORD-19 (<https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge>),
536 COKE,²⁷ and COVID-KOP.²⁸ These collections along with several research initiatives such as the
537 Antiviral Program for Pandemics (<https://www.niaid.nih.gov/research/antivirals>) and the Rapidly
538 Emerging Antiviral Drug Development Initiative (READDI; <https://www.readdi.org>) pushed the
539 scientific community to work in an ‘open science’ format.

540 Efforts similar to ours have been made to collect antiviral data prior to the SARS-CoV-2
541 outbreak. There is a collection of antiviral activity data from ChEMBL with enhanced taxonomy
542 annotations as a tool for studying the antiviral chemical space that the authors dubbed “Viral
543 ChEMBL”.² Viral ChEMBL was compiled using information collected on compounds related to
544 many virus types (human, animal, plant) and additional curation was performed to the data by
545 mapping lists for assay and target organism data and using a dictionary of virus-related terms.
546 While this collection is quite valuable to the field, it is based on an old version of ChEMBL20
547 (released 2015, current release is ChEMBL29) and some data are not relevant to human disease.
548 Thus, our database was collected and manually curated to provide a structured, annotated
549 repository of all data available in the most current version of ChEMBL for viruses that hold the
550 greatest risk for human contraction and pandemic potential.

551 The SMACC database also has the potential to guide more informed drug repurposing
552 efforts, which was a popular strategy employed during the first year of the SARS-CoV-2
553 pandemic. Repurposing FDA approved drugs²⁹ and their combinations³⁰ quickly provided options
554 for clinical use without the need to undergo extensive toxicological testing. For instance, we have
555 identified anticancer drug brequinar as a potential antiviral agent (cf. Table 3), and the analysis of
556 additional bioactivities, including those against host targets, may reveal novel interesting
557 compounds.

558 Beyond drug repurposing, another intuitive approach used in the most recent SARS-CoV-
559 2 pandemic to identify BSA drugs across the coronaviruses was through identifying proteome
560 conservation, which was studied by Schapira et al.³¹ as well as our group.³ Schapira et al. analyzed
561 the conservation of all available PDB structures of α - and β -coronaviruses, as well as samples from
562 patients with SARS-CoV-2 by mapping druggable binding pockets onto experimental structures

563 of SARS-CoV-2 proteins. Our work complemented that of Schapira et al.,³¹ by exploring the idea
564 that similarities between homologous coronaviral proteins could be exploited for target selection
565 and the development of broad-spectrum anti-coronaviral compounds. Putting that idea into the
566 context of potential broad-spectrum inhibitors of conserved targets, we identified drugs from
567 existing literature that inhibit M^{pro}, RdRp, PL^{pro}, and nsp10-nsp16, carefully collected and
568 analyzed all known experimental data on their antiviral activity and validated our hypothesis by
569 estimating their potential as broad-spectrum drugs. These compounds are discussed extensively
570 elsewhere and are naturally included in the SMACC database. Thus, we feel exploring the
571 conservation between homologous coronaviral proteins is an extremely valuable strategy for target
572 selection and could assist the development of BSA compounds.

573 With viral protein conservation as a tool for identifying BSAs, one wonders if there may
574 be a link between protein conservation and ligand promiscuity. While in theory the framework of
575 our database would easily allow for this analysis, the unfortunate truth is that the data are not
576 available, as our collection of target-based data was already far more limited than our phenotypic
577 set. Further, it is no secret that merely collecting such data from available data sources can be
578 misleading. Errors described above depict the challenges we faced in curation and collection; for
579 example, a user looking for compounds active against NS2B-NS2 Protease would not have found
580 results due to the target being annotated generally as “genome polyprotein.” We hope that our
581 systematic analysis and enumeration of annotation deficiencies and bioactivity data curation
582 protocols could help other researchers interested in expanding our collection or creating their own
583 specialized collections. Most importantly, our efforts both identified several BSAs discovered by
584 chance without deliberate focused efforts (cf. Tables 3-4) as well as nominated several compounds
585 for additional testing (cf. Table 5). As discussed above, the SMACC database included with this

586 paper, enables user-defined filtering of the data to support the generation of specialized subsets. In
587 summary, we posit that this study provides strong motivation for continued investments into
588 research targeting the discovery and development of novel BSA agents.

589 590 **Conclusions**

591
592 We have developed a pilot version of the SMACC (Small Molecule Antiviral Compound
593 Collection) database containing over 32,500 entries for 13 emerging viruses. We followed the
594 following steps to create SMACC: (i) identification, collection, and curation of all chemical
595 bioactivity data available in ChEMBL for 13 emerging viruses holding the greatest potential threat
596 to global human health; (ii) identification and resolution of the data availability, reproducibility
597 and quality challenges; (iii) integration of curated and carefully annotated data on compounds
598 tested in both phenotypic (21,392 entries) and target-based (11,123 entries) assays for these
599 viruses; and (iv) identification of chemicals showing high potential for BSA activity. Specifically,
600 we identified eight compounds active against 3-4 viruses from the phenotypic data, 16 compounds
601 active against two viruses from the target-based data, and 35 compounds active in at least one
602 phenotypic and one target-based assay. Duplicates (phenotypic and overlap sets) and singletons
603 (all sets) were also identified and annotated. While the pilot version of SMACC has integrated all
604 chemogenomic data available in ChEMBL for these viruses, there was a large degree of sparsity
605 (93%) within the integrated data matrix. Many viruses were understudied and thus, important
606 results may be obtained by targeted testing of compounds included in SMACC against targets
607 other than those against which they were tested. In fact, we have suggested several such targeted
608 testing experiments in this paper (cf. Table 5).

609 Our analysis indicates that not many BSAs have emerged from previous disconcerted
610 studies and that special, focused efforts must be established going forward. The SMACC database

611 built in this study may serve as a reference for virologists and medicinal chemists working on the
612 development of BSA agents in preparation for future viral outbreaks. The SMACC database is
613 publicly available in the form of searchable Excel spreadsheet at <https://smacc.mml.unc.edu>.

614

615 **Conflict of interest**

616 AT and ENM are co-founders of Predictive, LLC, which develops computational methodologies
617 and software for toxicity prediction. All other authors declare they have nothing to disclose.

618

619 **Acknowledgement**

620 Authors from UNC-Chapel Hill were supported by National Institutes of Health (Grants
621 U19AI171292 and R01GM140154). This research was supported by the Intramural Research
622 Program of the National Center for Advancing Translational Sciences (NCATS), National
623 Institutes of Health (NIH).

624

625

626 **References**

627

- 628 1. Bloom, D. E., Black, S. & Rappuoli, R. Emerging infectious diseases: A proactive
629 approach. *Proc Natl Acad Sci U S A* **114**, 4055–4059 (2017).
- 630 2. Nikitina, A. A., Orlov, A. A., Kozlovskaya, L. I., Palyulin, V. A. & Osolodkin, D. I.
631 Enhanced taxonomy annotation of antiviral activity data from ChEMBL. *Database* **2019**,
632 1–18 (2019).
- 633 3. Melo-Filho, C. C. *et al.* Conserved coronavirus proteins as targets of broad-spectrum
634 antivirals. *Antiviral Research* **204**, 105360 (2022).
- 635 4. Erik, D. C. & Guangdi, L. Approved Antiviral Drugs over the Past 50 Years. *Clinical*
636 *Microbiology Reviews* **29**, 695–747 (2016).
- 637 5. Bobrowski, T. *et al.* Learning from history: do not flatten the curve of antiviral research!
638 *Drug Discovery Today* **00**, 1–10 (2020).
- 639 6. Owen, D. R. *et al.* An oral SARS-CoV-2 Mpro inhibitor clinical candidate for the
640 treatment of COVID-19. *Science (1979)* **374**, 1586–1593 (2021).
- 641 7. Mokili, J. L., Rohwer, F. & Dutilh, B. E. Metagenomics and future perspectives in virus
642 discovery. *Current Opinion in Virology* **2**, 63–77 (2012).

- 643 8. Gaulton, A. *et al.* The ChEMBL database in 2017. *Nucleic Acids Res* **45**, D945–D954
644 (2017).
- 645 9. Fourches, D., Muratov, E. & Tropsha, A. Trust, but verify: on the importance of chemical
646 structure curation in cheminformatics and QSAR modeling research. *J Chem Inf Model*
647 **50**, 1189–204 (2010).
- 648 10. Fourches, D., Muratov, E. & Tropsha, A. Trust, but Verify II: A Practical Guide to
649 Chemogenomics Data Curation. *Journal of Chemical Information and Modeling* **56**,
650 1243–1252 (2016).
- 651 11. Fourches, D., Muratov, E. & Tropsha, A. Curation of chemogenomics data. *Nature*
652 *Chemical Biology* **11**, 535–535 (2015).
- 653 12. Berthold, M. R. *et al.* KNIME: The Konstanz Information Miner. in *Studies in*
654 *Classification, Data Analysis, and Knowledge Organization (GfKL 2007)* (Springer,
655 2007).
- 656 13. The BioAssay Ontology (BAO). <http://bioassayontology.org/> (2022).
- 657 14. Papadatos, G., Gaulton, A., Hersey, A. & Overington, J. P. Activity, assay and target data
658 curation and quality in the ChEMBL database. *Journal of Computer-Aided Molecular*
659 *Design* **29**, 885–896 (2015).
- 660 15. ChEMBL Compound Curation Pipeline. [http://chembl.blogspot.com/2020/02/chembl-](http://chembl.blogspot.com/2020/02/chembl-compound-curation-pipeline.html)
661 [compound-curation-pipeline.html](http://chembl.blogspot.com/2020/02/chembl-compound-curation-pipeline.html).
- 662 16. Papadatos, G., Gaulton, A., Hersey, A. & Overington, J. P. Activity, assay and target data
663 curation and quality in the ChEMBL database. *Journal of Computer-Aided Molecular*
664 *Design* **29**, 885–896 (2015).
- 665 17. He, X. *et al.* Generation of SARS-CoV-2 reporter replicon for high-throughput antiviral
666 screening and testing. *Proceedings of the National Academy of Sciences* **118**,
667 e2025866118 (2021).
- 668 18. Park, J.-G. *et al.* Identification and Characterization of Novel Compounds with Broad-
669 Spectrum Antiviral Activity against Influenza A and B Viruses. *Journal of Virology* **94**,
670 (2020).
- 671 19. Andersen, P. I. *et al.* Novel Antiviral Activities of Obatoclox, Emetine, Niclosamide,
672 Brequinar, and Homoharringtonine. *Viruses 2019, Vol. 11, Page 964* **11**, 964 (2019).
- 673 20. Li, S. fang *et al.* Antiviral activity of brequinar against foot-and-mouth disease virus
674 infection in vitro and in vivo. *Biomedicine & Pharmacotherapy* **116**, 108982 (2019).
- 675 21. CRISIS2: A Phase 2 Study of the Safety and Antiviral Activity of Brequinar in Non-
676 hospitalized Pts With COVID-19 - Study Results - ClinicalTrials.gov.
677 <https://www.clinicaltrials.gov/ct2/show/results/NCT04575038?view=results> (2022).
- 678 22. Schultz, D. C. *et al.* Pyrimidine inhibitors synergize with nucleoside analogues to block
679 SARS-CoV-2. *Nature* **2022 604:7904** **604**, 134–140 (2022).
- 680 23. Wang, Q.-Y. *et al.* Inhibition of Dengue Virus through Suppression of Host Pyrimidine
681 Biosynthesis. *Journal of Virology* **85**, 6548–6556 (2011).
- 682 24. Fu, H., Zhang, Z., Dai, Y., Liu, S. & Fu, E. Brequinar inhibits enterovirus replication by
683 targeting biosynthesis pathway of pyrimidines. *American Journal of Translational*
684 *Research* **12**, 8247 (2020).
- 685 25. Luthra, P. *et al.* Inhibiting pyrimidine biosynthesis impairs Ebola virus replication through
686 depletion of nucleoside pools and activation of innate immune responses. *Antiviral*
687 *Research* **158**, 288–302 (2018).

- 688 26. Brimacombe, K. R. *et al.* An OpenData portal to share COVID-19 drug repurposing data
689 in real time. *bioRxiv* (2020) doi:10.1101/2020.06.04.135046.
- 690 27. Korn, D. *et al.* COVID-19 Knowledge Extractor (COKE): A Tool and a Web Portal to
691 Extract Drug - Target Protein Associations from the COVID-19 Corpus of Scientific
692 Publications on COVID-19. (2020) doi:10.26434/CHEMRXIV.13289222.V1.
- 693 28. Korn, D. *et al.* COVID-KOP: integrating emerging COVID-19 data with the ROBOKOP
694 database. *Bioinformatics* **37**, 586–587 (2021).
- 695 29. Alves, V. M. *et al.* QSAR Modeling of SARS-CoV Mpro Inhibitors Identifies Sufugolix,
696 Cenicriviroc, Proglumetacin, and other Drugs as Candidates for Repurposing against
697 SARS-CoV-2. *Molecular Informatics* **40**, 2000113 (2021).
- 698 30. Bobrowski, T. *et al.* Synergistic and Antagonistic Drug Combinations against SARS-
699 CoV-2. *Molecular Therapy* **29**, 873–885 (2021).
- 700 31. Yazdani, S. *et al.* Genetic Variability of the SARS-CoV-2 Pocketome. *Journal of*
701 *Proteome Research* **20**, 4215 (2021).
702