1  **Improving prediction of compound function from chemical structure using**
2  **chemical-genetic networks**

3  Hamid Safizadeh[1,2], Scott W. Simpkins[3], Justin Nelson[3], Chad L. Myers [1,2,3§]

4  1. University of Minnesota-Twin Cities, Department of Electrical and Computer Engineering,
5  Minneapolis, Minnesota, USA

6  2. University of Minnesota-Twin Cities, Department of Computer Science and Engineering,
7  Minneapolis, Minnesota, USA

8  3. University of Minnesota-Twin Cities, Bioinformatics and Computational Biology,
9  Minneapolis, Minnesota, USA

10  § Correspondence to chadm@umn.edu

11    **ABSTRACT**

12    The drug discovery process can be significantly improved through understanding how the
13    structure of chemical compounds relates to their function. A common paradigm that has been
14    used to filter and prioritize compounds is ligand-based virtual screening, where large libraries of
15    compounds are queried for high structural similarity to a target molecule, with the assumption
16    that structural similarity is predictive of similar biological activity. Although the chemical
17    informatics community has already proposed a wide range of structure descriptors and similarity
18    coefficients, a major challenge has been the lack of systematic and unbiased benchmarks for
19    biological activity that covers a broad range of targets to definitively assess the performance of
20    the alternative approaches.

21    We leveraged a large set of chemical-genetic interaction data from the yeast *Saccharomyces*
22    *cerevisiae* that our labs have recently generated, covering more than 13,000 compounds from the
23    RIKEN NPDepo and several NCI, NIH, and GlaxoSmithKline (GSK) compound collections.
24    Supportive of the idea that chemical-genetic interaction data provide an unbiased proxy for
25    biological functions, we found that many commonly used structural similarity measures were
26    able to predict the compounds that exhibited similar chemical-genetic interaction profiles,
27    although these measures did exhibit significant differences in performance. Using the chemical-
28    genetic interaction profiles as a basis for our evaluation, we performed a systematic
29    benchmarking of 10 different structure descriptors, each combined with 12 different similarity
30    coefficients. We found that the All-Shortest Path (ASP) structure descriptor paired with the
31    Braun-Blanquet similarity coefficient provided superior performance that was robust across
32    several different compound collections.

33    We further describe a machine learning approach that improves the ability of the ASP metric to
34    capture biological activity. We used the ASP fingerprints as input for several supervised machine
35    learning models and the chemical-genetic interaction profiles as the standard for learning. We
36    found that the predictive power of the ASP fingerprints (as well as several other descriptors)
37    could be substantially improved by using support vector machines. For example, on held-out
38    data, we measured a 5-fold improvement in the recall of biologically similar compounds at a
39    precision of 50% based upon the ASP fingerprints. Our results generally suggest that using high-
40    dimensional chemical-genetic data as a basis for refining chemical structure descriptors can be a
41    powerful approach to improving prediction of biological function from structure.

42    **Keywords: chemical-genetic network, chemical structure, molecular descriptor, structural**
43    **similarity, virtual screening, and machine learning**

## INTRODUCTION

Discovery, design, and development of new drugs that reveal desired and reproducible biochemical behavior against a particular biomolecular target with minimal side effects are challenging. Despite the scientific and technological advances in drug discovery during the past 60 years, the number of drugs approved per billion US dollars that were spent for the development of novel drugs has halved roughly every 9 years since 1950 (Eroom's Law in contrast to Moore's Law) [1]. Following the similar property principle (SPP) [2], Ligand-based virtual screening (LBVS) has been commonly used as an a priori step to high-throughput screening (HTS) [3,4,5] to rank compounds of a large database in the decreasing order of their similarity to a reference or target molecule with known biological activity (**Fig. 1**). According to the similar property principle, structurally similar molecules are more likely to represent similar biological activities and physicochemical properties. Although there are limitations to the similar property principle [6], such as the case of activity cliffs where a very small modification in the structure of a molecule may drastically alter its biological properties [7], this structure-activity relationship is broadly consistent throughout the larger flat regions of activity landscapes [8,9]. Hence, the need for high performance structural similarity that extracts structurally analogous compounds from a database is inevitable.

To accelerate the retrieval of the compounds of a desired class that are active against a protein target, the chemical informatics community has suggested a wide range of structure descriptors and similarity coefficients that are able to extract candidates with similar biological activity from structural compound libraries. The most widely used representation of molecular graphs in these expanding databases of two- and three-dimensional molecular structures is based upon chemical fingerprints [10,11,12], where a molecular graph is represented by a fixed-length bit-vector that enumerates all the bounded-length paths in the graph and encodes the presence or absence of substructural fragments. The degree of similarity of two structural vectors describing two different compounds is usually measured by similarity coefficients, among which the well-known Tanimoto coefficient has still remained the coefficient of choice to capture the highest level of intermolecular similarity and thus biological activity [11,13]. The Tanimoto coefficient, which is formulated as the number of features shared between two molecules divided by the total number of features presented in both molecules, offers a suitable degree of chemical similarity between compounds, although this coefficient suffers from an intrinsic bias towards selection of smaller compounds [14,15]. However, the research community has lacked a systematic benchmark that assesses the performance of these structure descriptors and similarity coefficients over a broad range of protein targets in an unbiased manner.

Chemical genomic approaches, which focus on the systematic mapping of chemical-genetic interactions, offer a valuable new source of data to connect structure to function. These chemical-genetic maps take advantage of the massive wealth of chemical-genetic interaction profiles. The yeast Saccharomyces cerevisiae is a well-characterized eukaryotic system for which the genome-wide gene deletion project has identified ~5000 viable deletion mutants [16].

83 Testing each one of these viable mutants for hypersensitivity to a bioactive compound generates
84 a chemical-genetic interaction profile in which the relative fitness of a selected group of mutant
85 strains with defined genetic perturbations in response to the bioactive compound is quantified
86 [17,18]. These chemical-genetic interaction profiles provide functional information for a
87 compound that can be interpreted through the global genetic interaction network mapped for the
88 yeast [19]. If a bioactive compound inhibits a target protein, the loss-of-function mutations in a
89 gene that encodes the protein models the primary effects of the compound, and the genetic
90 interaction profile of the target gene resembles the chemical-genetic interaction profile of the
91 compound that inhibits the target pathway. Consequently, the chemical-genetic interaction
92 profiles of bioactive compounds can link those compounds to their cellular target pathways in an
93 unbiased manner. These profiles can be annotated to specific biological processes to predict the
94 general mechanisms of action for the bioactive compounds and can serve as an unbiased
95 genome-wide measure for their biological activity [20].

96 We generated a systematic benchmark based upon the similar property principle to assess
97 the performance of several structure descriptors and similarity coefficients in prediction of the
98 chemical-genetic interaction profiles of hundreds of compounds, with the assumption that these
99 profiles provided an unbiased genome-wide measure of biological activity. We generated and
100 annotated the yeast chemical-genetic interaction profiles for more than 13,000 compounds from
101 the RIKEN NPDepo as well as several NCI/NIH/GSK compound collections [20]. We
102 systematically benchmarked 10 different structure descriptors, each combined with 12 different
103 similarity coefficients, to identify the pair with the superior prediction of biological activity by
104 using the chemical-genetic interaction profiles as a basis for the biological activity of our
105 compounds. We further developed several supervised machine learning models to improve our
106 prediction of the biological activity of compounds from chemical structures, gaining higher
107 predictive power that was not in the ability of similarity coefficients. We found that support
108 vector machines (SVMs) [21] can significantly enhance the power of our chemical fingerprints
109 for predicting the biological activity of compounds.

110

111 **RESULTS AND DISCUSSION**

112

113 To evaluate the performance of the commonly used structure descriptors and similarity
114 coefficients in predicting the biological activity of compounds, we exhaustively searched for all
115 the pairwise candidates (i.e., one structure descriptor and one similarity coefficient) that provided
116 high predictive power over a wide range of protein targets using our chemical-genetic interaction
117 data. We generated, in our labs, the chemical-genetic interaction profiles for 13,524 compounds
118 from several diverse compound collections [20]. We used a subset of these screened compounds
119 that exhibited high confidence predictions for the annotated processes and biological pathways
120 based upon our chemical-genetic interaction profiles [20,22]. We included in our benchmarking
121 system two compound sets independently: (1) 826 compounds from the RIKEN Natural Product

122     Depository (NPDepo) compound collection, which we called RIKEN high confidence collection.

123     (2) 659 compounds from several NCI/NIH/GSK collections, which we called the NCI/NIH/GSK

124     high confidence collection (see **Material**).

125     **Establishing a systematic benchmark for chemical similarity measures**

126

127     We aimed at predicting the compound pairs that exhibited the most similar function

128     based upon our chemical-genetic interaction profiles, which served as an unbiased genome-wide

129     measure of biological activity. We labeled only 10% of the compound pairs with the highest

130     (cosine) profile similarity as our gold standard for true positives, which were highly prioritized in

131     our systematic benchmark. Following the similar property principle, a large number of these true

132     positives should be identified from the structural similarity of our compounds.

133     Our benchmarking system consisted of two main components: structure descriptors and

134     similarity coefficients. We evaluated both components through our systematic benchmark to find

135     the best performer of each component for prediction of the biological activity of our compounds.

136     We used jCompoundMapper [23] to describe all our compounds in 10 different structure spaces

137     (**Table 1**), where a compound was described with a fixed-length bit-vector that indicated the

138     presence or absence of a certain number of fingerprints. The number of features required for the

139     description of the compounds in a structure space varied based upon the space definition and

140     collection properties (e.g., only 9 of the predefined features in the MACCS keys were found in

141     the RIKEN high confidence collection, while RAD2D fingerprints generated 91082 features to

142     describe this compound collection). We used 12 widely used similarity coefficients (**Table 2**) to

143     measure the degree of similarity of two compounds described by a given structure descriptor.

144     **Evaluating the performance of chemical similarity measures**

145

146     We exhaustively searched for the best-performing chemical similarity measure (i.e., one

147     structure descriptor and one similarity coefficient) in predicting the biological activity of our

148     compounds. We described all our compounds in 10 different structure spaces and measured the

149     structural similarity of two compounds described in a given structure space by using a coefficient

150     of structural similarity, generating compound similarities for all the combinations of structure

151     descriptors and similarity coefficients. We ranked all these scores of structural similarity for the

152     prediction of the chemical-genetic profile similarity of our compounds and measured precision at

153     many recalls to evaluate the performance of alternative models (**Suppl. Table 1**). To isolate the

154     winning structure descriptor, we looked at the precision of all the prediction models at several

155     (lower) recalls for our RIKEN high confidence collection, which distinguished ASP, LSTAR,

156     and RAD2D fingerprints as the structure descriptors with superior predictive power (**Fig. 2a**).

157     The wide range of precision values achieved by different structure descriptors, assuming that a

158     single similarity coefficient was used, showed that our chemical-genetic interaction profiles

159     could highly separate structure descriptors in terms of their efficiency for the prediction of the

160     biological activity of compounds. We compared the relative performance of our distinguished

161   models with the predictive power of extended-connectivity fingerprints (ECFP) [24], which has
162   recently been one of the most common descriptors to represent molecular graphs. However,
163   ECFP did not generally outperform ASP, LSTAR, and RAD2D fingerprints, except at very few
164   recalls (**Fig. 2b**). Because this finding might simply be a result of our collection property, we
165   used our NCI/NIH/GSK high confidence collection to validate the predictive performance of the
166   ASP, LSTAR, and RAD2D fingerprints, which strongly confirmed the superiority of these
167   descriptors over ECFP at many recalls (**Fig. 3**).

168   ASP fingerprints encoded a graph traversal over all atoms in a molecular graph but stored
169   only the shortest paths between atoms, whereas LSTAR and RAD2D fingerprints described the
170   radial environment of all atoms in the molecular graph [23]. As a result, the ASP encoding that
171   described our compound collections needed fewer features than LSTAR and RAD2D encodings
172   (**Table 1**) although the predictive performance of the ASP fingerprints was higher or comparable
173   with that of LSTAR and RAD2D fingerprints at several recalls. Moreover, LSTAR fingerprints
174   generally exhibited higher performance than RAD2D fingerprints at several recalls (**Figs. 2-3**),
175   which could be justified by additional information that LSTAR fingerprints collected from the
176   radial environment of the atoms by definition. Therefore, we determined ASP and LSTAR
177   fingerprints as the winning structure descriptors to describe compounds and predict the ones with
178   the highest biological similarity to a target molecule using a similarity coefficient.

179   We systematically benchmarked 12 similarity coefficients (**Table 2**) by using our RIKEN
180   high confidence collection and measured the predictive performance of every coefficient over all
181   structure descriptors to find the winning similarity coefficient. We found that several (8 out of
182   12) coefficients were able to exhibit consistent high performance across all structure descriptors,
183   although 4 similarity coefficients (Asymmetric, Russel/Rao, Euclidean, and Dot-product) failed
184   in some structure descriptor spaces because their precision significantly dropped at lower recalls.
185   In other words, precision of several models (each corresponding to one structure descriptor) was
186   substantially low for each of these 4 similarity coefficients at many lower recalls (**Suppl. Table
187   1**), which indicated that these 4 coefficients were unable to predict the biological similarity of
188   our compounds across different structure descriptors consistently. We, as a result, removed these
189   4 coefficients from our analysis and focused only on 8 remaining similarity coefficients. We
190   found that the Braun-Blanquet similarity coefficient [25] resulted in the higher precision at many
191   recalls compared to all other coefficients, including Tanimoto and cosine coefficients (**Fig. 2a**),
192   which have been widely used by the chemical informatics community. For the Braun-Blanquet
193   similarity coefficient, the average precision and the average area under the receiver operating
194   characteristic (ROC) curves across all structure descriptors were slightly higher at many recalls
195   compared to those of the Tanimoto and cosine coefficients (**Fig. 2a**; columnar green values),
196   suggesting that this simple coefficient of structural similarity could confidently be used in place
197   of the traditional Tanimoto coefficient for the ranking of database compounds in the decreasing
198   order of biological similarity to a target molecule. The Braun-Blanquet coefficient, which was
199   simply formulated as the number of features common between two molecules divided by the

200 total number of features presented by the larger molecule, determined the degree of contribution
201 of the smaller molecular graph to the larger one. We further measured the performance of our
202 predictive models using our NCI/NIH/GSK high confidence collection to validate the superiority
203 of the Braun-Blanquet coefficient over other similarity coefficients (**Fig. 3a** and **Suppl. Table 2**).
204 We, therefore, paired the Braun-Blanquet similarity coefficient with the ASP and LSTAR
205 structure descriptors as our predictive models for ligand-based virtual screening.

**Optimizing the depth of structure descriptors**

207

208 One major parameter involved in the structural description of compound collections was
209 the describing depth; a high depth generated numerous features to cover the global environment
210 of each atom, whereas a low depth only focused on describing the local neighborhood of atoms
211 in the molecular graph. We assessed the impact of the depth of 5 structure descriptors (ASP, DFS,
212 ECFP, LSTAR, and RAD2D) in predicting the biological activity of our compounds using the
213 Braun-Blanquet similarity coefficient (**Fig. 4**). Describing our RIEKN high confidence collection
214 at a high depth generally resulted in strong predictions at lower recalls but moderate outcome at
215 higher recalls because a high depth was able to predict the compound pairs that were structurally
216 and therefore functionally very similar according to the similar property principle (SPP) but also
217 pushed undesired pairs such as activity cliffs to the top of our predictions. On the other hand, a
218 low describing depth was able to capture the similarity of two compounds in the local chunks of
219 the two molecular graphs that were essential for functional similarity, which eventually resulted
220 in drawing reasonable predictions at lower recalls. We, furthermore, evaluated these results using
221 our NCI/NIH/GSK high confidence collection (**Suppl. Fig. 1**), which confirmed similar general
222 trends that impacted our predictions at several recalls using 10 different describing depths.
223 Therefore, the structural description of a compound collection at high depths not only was
224 unnecessary and inefficient but also increased the computation time and space complexity. We
225 selected depth 8 for our evaluations although other neighboring depths were also justifiable.

**Improving the prediction performance via SVM models**

227

228 To increase the ability of chemical structures in predicting the biological activity of our
229 compounds, we designed several supervised machine learning models and took advantage of the
230 great wealth of chemical-genetic interaction maps for supervision. Moreover, we used supervised
231 principal component analysis [26] via chemical-genetic interaction maps to extract a number of
232 features from the more informative compound substructures that highly related structural data to
233 the biological activity of our compounds. We found that support vector regression (SVR) models
234 [21] were able to boost the prediction performance of the functional activity of our compounds
235 from their chemical structures by weighting supervised principal components, where chemical-
236 genetic interaction profiles were also input to the learning models for supervision.

237

238        We designed a learning pipeline (**Fig. 5a**) to predict chemical-genetic profile similarities
239    by creating bootstraps [27] and pairwise structural encodings (see **Methods**). We implemented
240    support vector regression models in LibSVM [28], a popular open source support vector machine
241    learning library developed at National Taiwan University, and used Radial Basis Function (RBF)
242    kernels for building our epsilon support vector regression models. We developed precision-recall
243    (PR) curves to evaluate the performance of our models for different structure descriptors, where
244    only 10% of the compound pairs from our RIKEN high confidence collection with the highest
245    chemical-genetic profile similarity were labeled as the gold standard for true positives. We found
246    that a subset of our structure descriptors (SD1-SD6 and SD10) were able to achieve significantly
247    higher performance in the prediction of the functional similarity of our compounds than the best-
248    performing chemical similarity measures (i.e., ASP or LSTAR fingerprints along with the Braun-
249    Blanquet similarity coefficient). The learning curves for the ASP and LSTAR fingerprints (**Fig.
250    5b**) exhibited that we were able to gain a 5-fold improvement in the recall of biologically similar
251    compounds at a precision of 50%. However, the degree of improvement was dependent on the
252    functional diversity of datasets, which could result in modestly higher performance for particular
253    collections with higher diversity; for instance, we improved our predictions for the
254    NCI/NIH/GSK high confidence collection by only about 2 folds in the recall of biologically
255    similar compounds at a precision of 50% (**Fig. 5d-e**). This relatively poor improvement
256    (compared to that of the RIKEN high confidence collection) was explained by the higher
257    functional diversity of our NCI/NIH/GSK high confidence collection (score of ~25.3, against
258    ~14.6 for the RIKEN high confidence collection) although the two collections exhibited similar
259    structural diversity (score of ~62) (see **Methods**). This high functional diversity of the
260    NCI/NIH/GSK high confidence collection was due to the presence of 6 functionally different
261    sub-collections, which consequently affected the ability of our models to learn chemical-genetic
262    similarities at a high performance for this collection. Although model performance was disturbed
263    by the higher diversity of our NCI/NIH/GSK high confidence collection, we still measured more
264    distinct learning curves from the baseline while labeling 20% of functionally most similar
265    compound pairs as true positives (**Fig. 5e**), indicating that functionally similar pairs were
266    eventually pushed up to the top of the ranked lists by our learning models. Furthermore, we
267    combined the two collections, which added not only more diversity but also more compounds to
268    the RIKEN high confidence collection, and made predictions for the combined dataset, resulting
269    in about a 4.5-fold improvement in the recall of biologically similar compounds at the precision
270    of 50% (**Fig. 5c**). To accomplish higher prediction performance, we, therefore, would need a
271    larger training set (compounds with known chemical-genetic interaction profiles) to compensate
272    for the high functional diversity of compound collections and facilitate the learning process.
273

274    **Predictive power of structural similarity and SVM models**
275

276        To investigate the compounds driving our prediction models and the underlying function,
277    we clustered our compound collections into 10 functional as well as 10 structural clusters using

278     K-means and K-medoids, respectively, and mapped only the true positives at the top of our PR
279     curves to their corresponding functional and structural clusters (**Fig. 6**). We found that a large
280     group of compounds generating high prediction scores at the top of our learning curves belonged
281     to the same functional clusters (**Fig. 6b**), whereas the baseline curve for the ASP fingerprints and
282     Braun-Blanquet similarity coefficient included several functional clusters even at lower recalls
283     (**Fig. 6c**). Therefore, our learning models placed more emphasis on the learning of a few certain
284     functional clusters and boosted our prediction performance for those clusters. The first functional
285     cluster that appeared on the learning curve of ASP fingerprints for the RIKEN high confidence
286     collection (blue bar in **Fig. 6b**) was enriched for cell cycle processes based upon the predictions
287     at the MOSAIC database [29] (**Suppl. Tables 3-4** for enrichment of functional clusters). Despite
288     this functional tendency that the learning models showed, several structural clusters contributed
289     to the predicted pairs with high scores for the learning models (**Fig. 6e**), which could lead us to
290     discovering structurally diverse compounds that would exhibit similar biological activities. The
291     discovery of such compounds was of crucial importance since exploring functional analogs with
292     dissimilar structures was entirely out of the capacity of the similar property principle. Hence, our
293     learning models were able to extract compounds with similar function but distinct structures for a
294     target drug/compound, which was far beyond the scope of structural similarity coefficients. For
295     instance, our learning model for the RIKEN high confidence collection assigned a high score to
296     the biologically similar compounds NPD2186 and NPD3120 (chemical-genetic profile similarity
297     of 0.862), while the Braun-Blanquet similarity coefficient for this pair was as low as 0.027 (**Fig.**
298     **8c; Suppl. Table 5**). This property of our learning models was achieved by design, where we
299     extracted only a subset of the supervised principal components that were significantly related to
300     the biological activity of compounds, and we then weighted this subset of supervised principal
301     components using support vector regression models. This model property existed in our learning
302     models for the NCI/NIH/GSK high confidence collection as well but in a weaker manner due to
303     the high functional diversity of this collection (**Fig. 7** and **Fig. 8d**).
304

305         Furthermore, we assessed the predictive power of structural similarity (using the Braun-
306     Blanquet coefficient) against that of chemical-genetic profile similarity for our collections. For
307     the former, we predicted the chemical-genetic profile similarity of our compounds from chemical
308     structures, whereas, for the latter, we used the chemical-genetic profile similarity of compounds
309     to predict their structural similarity. The PR curves revealed that the structural similarity of our
310     compounds had higher predictive power of the chemical-genetic similarity than the latter of the
311     former (**Fig. 8a-b**) since compounds with similar biological activity could represent completely
312     different chemical structures. On the other hand, compounds with similar structures were highly
313     expected to exhibit similar biological activity; therefore, our results were a strong confirmation
314     for the similar property principle. Moreover, the degree of superiority of the predictive power of
315     structural similarity over functional similarity was an indicator of the amount of substructures
316     (i.e., compounds with similar biological activity but distinct structures) that existed in the
317     collection. The wide gap between the curves for the RIKEN high confidence collection (**Fig. 8a**),

318    therefore, represented a large number of substructures in this collection (see **Fig. 6**), while the
319    narrow gap between the curves for the NCI/NIH/GSK high confidence collection (**Fig. 8b**)
320    served as a signal that this collection, which was composed of several sub-collections, included
321    more of one-to-one correspondence between structural and functional profiles. Since the high
322    power of our learning models was to discover compounds of various structures (in addition to
323    compounds with similar structures) that exhibited similar biological activity to a target
324    drug/compound, our learning method showed enormous superiority over the baseline approach
325    (ASP fingerprints paired with the Braun-Blanquet similarity coefficient) for our RIKEN high
326    confidence collection. Although our learning method improved predictions for the functionally
327    diverse collections (such as our NCI/NIH/GSK high confidence collection) moderately, this
328    method exhibited strong predictions for the larger collections representing certain biological
329    functions with structurally diverse compounds.

330

331    **CONCLUSION**

332

333          The chemical informatics community has adopted a broad range of structure descriptors
334    and similarity coefficients for ligand-based virtual screening where the similar property principle
335    has been the basis for ranking of compounds with similar biological activity to a target molecule
336    from chemical structures. However, the research community has lacked a systematic, unbiased
337    benchmark for biological activity that would cover a wide range of targets to definitively assess
338    the performance of alternatives. We generated chemical-genetic interaction profiles from yeast in
339    our labs, covering 13,431 compounds from the RIKEN NPDepo and several NCI/NIH/GSK
340    compound collections, and used these profiles as an unbiased standard for the biological activity
341    of our compounds. Using these chemical-genetic interaction profiles as the basis for the function
342    of our compounds, we systematically benchmarked 10 different structure descriptors and 12
343    different similarity coefficients. We found that the ASP (and LSTAR) fingerprints paired with
344    the Braun-Blanquet similarity coefficient revealed as the superior choice for ranking of
345    compounds with similar biological activity to a target molecule. The ASP fingerprints encoded
346    all shortest paths between atoms obtained through an exhaustive depth-first search of the
347    molecular graph (up to a predefined depth), and the Braun-Blanquet coefficient represented the
348    number of features shared between two molecules divided by the number of features presented in
349    the larger one. Moreover, we devised a machine learning model that boosted the predictive
350    power of several fingerprints, although the degree of improvement was correlated to the
351    functional diversity of our compound collections. We found that structural similarity had a
352    higher predictive power in prediction of functional similarity than the latter of the former
353    because several substructures contributed to the similar biological activity. Although similarity
354    coefficients predicted the compounds that had both similar function and similar structure, our
355    learning models assigned higher predictive scores to most compounds with similar function by
356    weighting the supervised principal components that were strongly correlated to the chemical-
357    genetic profiles. Therefore, our learning models were able to predict compounds from a library

358    with similar biological activity but diverse structures to a target molecule, which significantly
359    improved performance relative to simple similarity coefficients applied to structure descriptors.

360    **ACKNOWLEDGEMENTS**

369
370

371    **MATERIAL AND METHODS**

372
373    **Data Collections**

374
375    We used two different compound collections independently: Our RIKEN high confidence
376    collection, as a subset of the RIKEN NPDepo, was composed largely of purified natural products
377    or natural product derivatives, whereas our NCI/NIH/GSK high confidence collection was a
378    diverse set of several sub-collections: 4 collections from the National Cancer Institute's Open
379    Chemical Repository (natural products, approved oncology drugs, and structural and mechanistic
380    diversity sets), a library of compounds from the National Institutes of Health Small Molecule
381    Repository with a history of use in human clinical trials (NIH Clinical Collection), and the
382    Glaxo-Smith-Kline kinase inhibitor collection (GSK).

383
384    **Designing the support vector machine learning pipeline**

385
386    We proposed support vector regression (SVR) models for the prediction of the functional
387    activity of our compounds based upon their chemical structures. We used LibSVM [28] for the
388    implementation of our models and bootstrapping [27] for generating our training and test sets. To
389    generate these training and test sets, we drew N (the total number of compounds in a collection)
390    samples uniformly random from the collection but with replacement, assigning ~0.632N unique
391    compounds to the training and the rest to the test set. We used supervised principal component
392    analysis [26] with adaption of chemical-genetic interaction profiles, assuming that these profiles
393    were known for the training set but unknown for the test set, to lower the dimension of structure
394    spaces. We normalized each structural vector that described a compound in the low-dimensional
395    space by its Euclidean length and multiplied each pair of the normalized vectors (both from the
396    training set or both from the test set) in the element-wise manner to create a new structure space,
397    called "pairwise structural encodings", for the representation of compound pairs (**Suppl. Fig.**

398 **2a**). We devised a pipeline (**Fig. 5a**) to predict our chemical-genetic profile similarities by using

399 pairwise structural encodings, feeding our regression models with $\binom{m}{2}$ pairwise structural

400 vectors and chemical-genetic profile similarities, where m was the number of compounds in the

401 training set, to predict the chemical-genetic profile similarities for the $\binom{N-m}{2}$ compound

402 pairs that were corresponding to the test set. We used Radial Basis Function (RBF) kernels to

403 build up epsilon support vector regression models and input a number of bootstraps to the

404 models to evaluate the average performance of our models across all bootstraps (**Suppl. Fig. 2b**).

405 To measure the prediction of our pipeline for a newly seen input, we needed to take the average

406 over all the model outputs resulted from different bootstraps, where the new input was treated as

407 a test data in the test sets; the higher the number of bootstraps the more accurate the prediction

408 value. We generated a large number of bootstraps (200 bootstraps for the RIKEN and

409 NCI/NIH/GSK high confidence collections as well as 100 bootstraps for the combined

410 collection) for our evaluations although the performance of our learning models was constant

411 after meeting a certain number of bootstraps.

412

413 **Estimating the diversity of compound collections**

414

415 We assigned all the compounds in a collection to a single cluster and split up the cluster

416 recursively to form clusters of more similar compounds. At any step of recursion, we determined

417 the cluster with the lowest average within-cluster chemical-genetic profile similarity (to compute

418 the functional diversity) or structural similarity (to compute the structural diversity) and divided

419 the cluster into two new clusters using K-means or K-medoids clustering. We stopped generating

420 new clusters right before our algorithm would generate at least two individual clusters exceeding

421 our predefined hard limit for the maximum average between-cluster chemical-genetic similarity

422 (cosine similarity of 0.3) or structural similarity (Braun-Blanquet similarity of 0.3). We repeated

423 the algorithm many times (1000 times for the functional diversity and 100 times for the structural

424 diversity) and computed the mean diversity score as the average exponentiation of the Shannon

425 entropy indices over all the instances:

426 $$D = mean\left(2^{\left(-\sum_i p_i \log_2(p_i)\right)}\right)$$

427 where $p_i$ was the proportional abundance of compounds in the $i^{th}$ cluster of the final clustering.

428

429 **References**

430

1. Scannell JW, Blanckley A, Boldon H, Warrington B. Diagnosing the decline in pharmaceutical R&D efficiency. Nat. Rev. Drug Discov. 2012; 11: p. 191–200.

2. Johnson MA, Maggiora GM. Concepts and applications of molecular similarity: John Wiley; 1990.

3. Bajorath J. Integration of virtual and high-throughput screening. Nat. Rev. Drug Discov. 2002; 1: p. 882-894.

4. Tanrikulu Y, Krüger B, Proschak E. The holistic integration of virtual screening in drug discovery. Elsevier Drug Discov. Today. 2013; 18: p. 358–364.

5. Stahura FL, Bajorath J. Virtual Screening Methods that Complement HTS. Combinatorial Chemistry & High Throughput Screening. 2004; 7: p. 259-269.

6. Eckert H, Bajorath J. Molecular similarity analysis in virtual screening: foundations, limitations, and novel approaches. Elsevier Drug Discov. Today. 2007; 12: p. 225–233.

7. Guha R, Van Drie JH. Structure-Activity Landscape Index: Identifying and Quantifying Activity Cliffs. J. Chem. Inf. Model. 2008; 48: p. 646-658.

8. Bajorath J, Peltason L, Wawer M, Guha R, Lajiness MS, Van Drie JH. Navigating structure–activity landscapes. Elsevier Drug Discov. Today. 2009; 14: p. 698-705.

9. Wassermann AM, Wawer M, Bajorath J. Activity Landscape Representations for Structure-Activity Relationship Analysis. J. Med. Chem. 2010; 53: p. 8209–8223.

10. Duanb J, Dixona SL, Lowriea JF, W. S. Analysis and comparison of 2D fingerprints: Insights into database screening performance using eight fingerprint methods. J. Molecular Graphics and Modelling. 2010; 29: p. 157–170.

11. Willett P. Similarity-based virtual screening using 2D fingerprints. Elsevier Drug Discov. Today. 2006; 11: p. 1046-1053.

12. Raymond JW, Willett P. Effectiveness of graph-based and fingerprint-based similarity measures for virtual screening of 2D chemical structure databases. J. Computer-Aided Molecular Design. 2002; 16: p. 59–71.

13. Willett P. Similarity searching using 2D structural fingerprints. In.: Springer; 2011. p. 133-158.

14. Fligner MA, Verducci JS, Blower PE. A Modification of the Jaccard–Tanimoto Similarity Index for Diverse Selection of Chemical Compounds Using Binary Strings. Technometrics. 2002; 44: p. 110- 119.

15. Holliday JD, Salim N, Whittle M, Willett P. Analysis and Display of the Size Dependence of Chemical Similarity Coefficients. J. Chem. Inf. Comput. Sci. 2003; 43: p. 819-828.

16. Giaever G, et al. Functional profiling of the Saccharomyces cerevisiae genome. Nature. 2002; 418: p. 387–391.

17. Parsons A, et al. Exploring the mode-of-action of bioactive compounds by chemical-genetic profiling in yeast. Cell. 2006; 126: p. 611–625.

18. Giaever G, et al. Chemogenomic profiling: Identifying the functional interactions of small molecules in yeast. Proc. Natl. Acad. Sci. U. S. A. 2004; 101: p. 793–798.

19. Parsons A, et al. Integration of chemical-genetic and genetic interaction data links bioactive compounds to cellular target pathways. Nat. Biotech. 2004; 22: p. 62–69.

20. Piotrowski JS, et al. Functional annotation of chemical libraries across diverse biological processes. Nat. Chem. Bio. (under review). .

21. Vapnik VN. Statistical Learning Theory: John Wiley; 1998.

22. Simpkins SW, et al. in preparation.

23. Hinselmann G, Rosenbaum L, Jahn A, Fechner N, Zell A. jCompoundMapper: An open source Java library and command-line tool for chemical fingerprints. J. Cheminformatics. 2011; 3: p. 1-14.

24. Rogers D, Hahn M. Extended-connectivity fingerprints. J. Chem. Inf. Model. 2010; 50: p. 742-754.

25. Hayek LC. Measuring and monitoring biological diversity: standard methods for amphibians. In Heyer WR, et al. , editors..: Smithsonian Books, Washington, D.C.; 1994.

26. Barshan E, Ghodsi A, Azimifar Z, Zolghadri Jahromi M. Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds. Pattern Recognition. 2011; 44: p. 1357–1371.

27. Efron B, Tibshirani R. Improvements on cross-validation: the 632+ bootstrap method. J. American Statistical Association. 1997; 92: p. 548-560.

28. Chang CC, Lin CJ. LIBSVM: a library for support vector machines. ACM Transactions on Intelligent Systems and Technology (TIST). 2011; 2: p. 1-27.

29. Nelson J, et al. in preparation.

432

433 **Figure/Table legends:**

434

435 **Figure 1. Ligand-based virtual screening of a target (e.g., NPD2186 from RIKEN NPDepo).**
436 All the compounds of the MOSAIC database (http://mosaic.cs.umn.edu) [29] were ranked in the
437 decreasing order of structural similarity to the target molecule based upon the similar property
438 principle (SPP). In this ranked list, NPD4974 had a very distinct chemical-genetic profile from
439 NPD2186, appearing as a false positive generated by SPP. Here, we described all the compounds
440 using ASP fingerprints (depth 8) and measured the structural similarity using the Braun-Blanquet
441 similarity coefficient.

442

443 **Figure 2. Model precision for all structure descriptors paired with the Cosine, Tanimoto, or**
444 **Braun-Blanquet similarity coefficient using our RIKEN high confidence data collection. (a)**
445 Precision at several recalls and the area under the ROC curve for each model. The red and green
446 values represented the highest precision achieved at a given recall and the average precision over
447 all the structure descriptors for a given similarity coefficient at a recall, respectively. **(b)** Relative
448 performance of ASP (teal), LSTAR (gold), and RAD2D (magenta) fingerprints to ECFP. For all
449 the structure descriptors that required a depth of description, precision was measured at depth 8.

450

451 **Figure 3. Model precision for all structure descriptors paired with the Cosine, Tanimoto, or**
452 **Braun-Blanquet similarity coefficient using our NCI/NIH/GSK high confidence data**
453 **collection. (a)** Precision at several recalls and the area under the ROC curve for each model. The
454 red and green values represented the highest precision achieved at a given recall and the average
455 precision over all the structure descriptors for a given similarity coefficient at a recall,
456 respectively. **(b)** Relative performance of ASP (teal), LSTAR (gold), and RAD2D (magenta)
457 fingerprints to ECFP. For all the structure descriptors that required a depth of description,
458 precision was measured at depth 8.

459

460 **Figure 4. Depth impact of structure descriptors on the performance of our prediction**
461 **models.** We measured the precision of our prediction models at 10 consequent molecular depths
462 for 5 different structure descriptors, each paired with the Braun-Blanquet similarity coefficient,
463 using our RIKEN high confidence collection.

464

465 **Figure 5. Prediction performance of learning models. (a)** Learning pipeline for one bootstrap
466 using pairwise structural encodings (see **Methods**). **(b)** Model performance for our RIKEN high
467 confidence collection. The blue curve was the prediction performance of ASP fingerprints paired
468 with the Braun-Blanquet similarity coefficient, whereas the teal and gold curves represented the
469 performance of ASP and LSTAR fingerprints using our learning models, respectively. **(c)** Model
470 performance for the combined RIKEN and NCI/NIH/GSK high confidence collections. **(d)**
471 Performance of the learning models for the NCI/NIH/GSK high confidence collection, where
472 true positives were only 10% (default for this paper) or **(e)** 20% of the compound pairs with the
473 highest chemical-genetic profile similarity.

474

475 **Figure 6. Functional and structural clustering for our RIKEN high confidence collection.**
476 **(a)** Distribution of 10 functional clusters generated by k-means using chemical-genetic profiles.
477 Contribution of these functional clusters to the top TP pairs extracted by **(b)** our learning model
478 and **(c)** the Braun-Blanquet similarity coefficient using the ASP fingerprints. **(d)** Distribution of

479   10 structural clusters generated by k-medoids using the ASP fingerprints. Contribution of these
480   structural clusters to the top TP pairs that were introduced by **(e)** our learning model and **(f)** the
481   Braun-Blanquet similarity coefficient.
482
483   **Figure 7. Functional and structural clustering for our NCI/NIH/GSK high confidence**
484   **collection. (a)** Distribution of 10 functional clusters generated by k-means using chemical-
485   genetic profiles. Contribution of these functional clusters to the top TP pairs extracted by **(b)** our
486   learning model and **(c)** the Braun-Blanquet similarity coefficient using the ASP fingerprints. **(d)**
487   Distribution of 10 structural clusters generated by k-medoids using the ASP fingerprints.
488   Contribution of these structural clusters to the top TP pairs that were introduced by **(e)** our
489   learning model and **(f)** the Braun-Blanquet similarity coefficient.
490
491   **Figure 8. Predictive power of structural similarity as a result of chemical substructures.** By
492   using the **(a)** RIKEN and **(b)** NCI/NIH/GSK high confidence collections, we measured that
493   structural similarity showed higher predictive power of chemical-genetic profile similarity than
494   the latter of the former largely because of substructures. Our learning models extracted
495   biologically similar but structurally very dissimilar compounds for **(c)** NPD2186 from our
496   RIKEN high confidence collection and for **(d)** NSC745750 from our NCI/NIH/GSK high
497   confidence collection.
498
499   **Suppl. Figure 1. Depth impact of structure descriptors on the performance of prediction**
500   **models.** We measured the precision of our prediction models at 10 consequent molecular depths
501   for 5 different structure descriptors, each paired with the Braun-Blanquet similarity coefficient,
502   using our NCI/NIH/GSK high confidence collection.
503
504   **Suppl. Figure 2. Pairwise structural encodings and bootstrapping. (a)** Pairwise structural
505   features created by the element-wise multiplication of the normalized, low-dimensional
506   structural vectors. We reduced dimension of descriptors using a supervised principal component
507   analysis method. **(b)** Smoothing average over bootstraps. At each bootstrap, the chemical-genetic
508   profile similarity of the test pairs (represented by "X") was predicted, and all the predicted values
509   for a test pair at different bootstraps were averaged to smoothen the prediction. For example, the
510   compound pair 1 was a test pair in bootstraps 1, 2, and 4 (In bootstrap 3, it might be a training
511   pair or an invalid pair where one compound belonged to the training set and the other to the test
512   set).
513
514   **Table 1. Structure descriptors.** 10 different topological, fingerprint-based structure descriptors
515   generated by jCompoundMapper for the description of each compound in our datasets. The right
516   column represented the total number of features needed to describe our high confidence RIKEN
517   collection (826 compounds).
518
519   **Table 2. Similarity coefficients.** 12 different similarity coefficients (several of these coefficients
520   were collected by Raymond and Willett [12]) for measurement of the degree of similarity of two
521   compounds described by a given fingerprint-based structure descriptor. Here, $x$ = number of bits
522   set in both fingerprints, $y$ = number of bits set in the first fingerprint, $z$ = number of bits set in the

523    second fingerprint, and $w$ = number of bits in the bit string. For the Tullos similarity coefficient,

524    $X = \log\left(1 + \dfrac{\min(y,z)}{\max(y,z)}\right)/\log(2)$, $Y = \left(\log\left(2 + \dfrac{\min(y,z)-x}{x+1}\right)/\log(2)\right)^{-\frac{1}{2}}$, $Z = \dfrac{\log\left(1 + \dfrac{x}{y}\right).\log\left(1 + \dfrac{x}{z}\right)}{\log^2(2)}$.

525

526    **Suppl. Table 1. Precision at several recalls for all combinations of structure descriptors and**
527    **similarity coefficients using our RIKEN high confidence collection.**

528

529    **Suppl. Table 2. Precision at several recalls for all combinations of structure descriptors and**
530    **similarity coefficients using our NCI/NIH/GSK high confidence collection.**

531

532    **Suppl. Table 3. Functional enrichment of clusters for the 1000 top TP pairs of our learning**
533    **model using our RIKEN high confidence collection.**

534

535    **Suppl. Table 4. Functional enrichment of clusters for the 1000 top TP pairs of our learning**
536    **model using our NCI/NIH/GSK high confidence collection.**

537

538    **Suppl. Table 5. The list of compounds with similar chemical-genetic interaction profiles but**
539    **dissimilar chemical structures for the largest functional cluster of 1000 top TP pairs (blue**
540    **bars in Figs. 6b-7b).**

**Table 1**

| ID | Name | Description | # Features |
|----|------|-------------|-----------|
| SD1 | AP2D | Topological atom pairs | 1210 |
| SD2 | ASP | All-shortest path encodings | 26114 |
| SD3 | AT2D | Topological atom triplets | 56900 |
| SD4 | DFS | All-path encodings | 48267 |
| SD5 | ECFP | Extended connectivity fingerprints | 42131 |
| SD6 | LSTAR | Local path environments | 84450 |
| SD7 | MACCS | Predefined pharmacophores | 9 |
| SD8 | PHAP2POINT2D | Topological pharmacophore pair encodings | 17 |
| SD9 | PHAP3POINT2D | Topological pharmacophore triplet encodings | 302 |
| SD10 | RAD2D | Topological molprint-like fingerprints | 91082 |

## Table 2

| ID | Name | Measurement | Range |
|----|------|-------------|-------|
| SC1 | Cosine | $\dfrac{x}{\sqrt{yz}}$ | 0 to 1 |
| SC2 | Tanimoto | $\dfrac{x}{y+z-x}$ | 0 to 1 |
| SC3 | Kulczynski | $\dfrac{x(y+z)}{2yz}$ | 0 to 1 |
| SC4 | Dice | $\dfrac{2x}{y+z}$ | 0 to 1 |
| SC5 | Sokal/Sneath | $\dfrac{x}{2y+2z-3x}$ | 0 to 1 |
| SC6 | Tullos | $XYZ$ | 0 to 1 |
| SC7 | McConnaughey | $\dfrac{x(y+z)-yz}{yz}$ | -1 to 1 |
| SC8 | Asymmetric | $\dfrac{x}{\min(y,z)}$ | 0 to 1 |
| SC9 | Braun-Blanquet | $\dfrac{x}{\max(y,z)}$ | 0 to 1 |
| SC10 | Russel/Rao | $\dfrac{x}{w}$ | 0 to 1 |
| SC11 | Euclidean | $\dfrac{1}{1+\sqrt{y+z-2x}}$ | 0 to 1 |
| SC12 | Dot-product | $x$ | 0 to $\infty$ |

RIKEN NPDepo

NPD2186

Structural Similarity Search

NPD2285

NPD2104

NPD4974

NPD2265

Figure 1

a

| | Recall = 0.002 | | | Recall = 0.005 | | | Recall = 0.02 | | |
|---|---|---|---|---|---|---|---|---|---|
| | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** |
| **SD1** | 0.850 | 0.850 | 0.861 | 0.684 | 0.683 | 0.676 | 0.355 | 0.357 | 0.370 |
| **ASP** | **0.919** | **0.919** | 0.912 | 0.829 | 0.819 | 0.830 | 0.557 | 0.567 | **0.586** |
| **SD3** | 0.849 | 0.850 | 0.866 | 0.785 | 0.785 | 0.770 | 0.515 | 0.515 | 0.523 |
| **SD4** | 0.877 | 0.888 | 0.871 | 0.765 | 0.767 | 0.740 | 0.540 | 0.544 | 0.543 |
| **SD5** | 0.855 | 0.854 | 0.861 | 0.743 | 0.749 | 0.777 | 0.555 | 0.556 | 0.556 |
| **LSTAR** | 0.883 | 0.883 | 0.883 | 0.833 | 0.833 | 0.828 | 0.565 | 0.566 | 0.560 |
| **SD7** | 0.102 | 0.102 | 0.102 | 0.074 | 0.074 | 0.074 | 0.091 | 0.091 | 0.091 |
| **SD8** | 0.138 | 0.138 | 0.138 | 0.044 | 0.044 | 0.044 | 0.049 | 0.049 | 0.049 |
| **SD9** | 0.056 | 0.056 | 0.056 | 0.081 | 0.081 | 0.081 | 0.108 | 0.108 | 0.108 |
| **RAD2D** | 0.875 | 0.866 | 0.861 | 0.810 | 0.814 | **0.834** | 0.535 | 0.536 | 0.548 |
| | **0.640** | **0.641** | **0.641** | **0.565** | **0.565** | **0.565** | **0.387** | **0.389** | **0.393** |

| | Recall = 0.05 | | | Recall = 0.2 | | | Area Under the Curve (AUC) | | |
|---|---|---|---|---|---|---|---|---|---|
| | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** |
| **SD1** | 0.221 | 0.222 | 0.230 | 0.154 | 0.154 | 0.156 | 0.590 | 0.592 | 0.593 |
| **ASP** | 0.288 | 0.287 | **0.294** | 0.154 | 0.156 | 0.158 | 0.575 | 0.579 | 0.584 |
| **SD3** | 0.254 | 0.259 | 0.256 | 0.152 | 0.153 | 0.152 | 0.581 | 0.583 | 0.584 |
| **SD4** | 0.267 | 0.270 | 0.279 | 0.154 | 0.156 | 0.158 | 0.572 | 0.576 | 0.580 |
| **SD5** | 0.263 | 0.268 | 0.277 | 0.146 | 0.149 | 0.153 | 0.574 | 0.578 | 0.584 |
| **LSTAR** | 0.288 | 0.290 | **0.294** | 0.168 | 0.170 | **0.171** | 0.587 | 0.590 | **0.596** |
| **SD7** | 0.100 | 0.100 | 0.100 | 0.116 | 0.116 | 0.116 | 0.526 | 0.526 | 0.526 |
| **SD8** | 0.062 | 0.062 | 0.062 | 0.091 | 0.091 | 0.091 | 0.499 | 0.499 | 0.501 |
| **SD9** | 0.109 | 0.109 | 0.109 | 0.098 | 0.098 | 0.101 | 0.486 | 0.488 | 0.492 |
| **RAD2D** | 0.289 | 0.291 | **0.294** | 0.161 | 0.162 | 0.164 | 0.579 | 0.582 | 0.588 |
| | **0.214** | **0.216** | **0.219** | **0.139** | **0.140** | **0.142** | **0.557** | **0.559** | **0.563** |

b



Figure 2

a

| | Recall = 0.002 | | | Recall = 0.005 | | | Recall = 0.02 | | |
|---|---|---|---|---|---|---|---|---|---|
| | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** |
| **SD1** | 0.818 | 0.813 | 0.781 | 0.594 | 0.609 | 0.644 | 0.312 | 0.319 | 0.332 |
| **ASP** | 0.839 | 0.839 | **0.932** | 0.789 | 0.810 | 0.813 | 0.535 | **0.541** | 0.529 |
| **SD3** | 0.841 | 0.828 | 0.884 | 0.819 | 0.819 | 0.801 | 0.518 | 0.495 | 0.461 |
| **SD4** | 0.828 | 0.828 | 0.916 | 0.772 | 0.789 | 0.807 | 0.491 | 0.495 | 0.490 |
| **SD5** | 0.897 | 0.897 | 0.916 | 0.783 | 0.783 | 0.795 | 0.373 | 0.395 | 0.400 |
| **LSTAR** | 0.897 | 0.897 | 0.916 | 0.856 | 0.870 | **0.871** | 0.533 | 0.527 | 0.520 |
| **SD7** | 0.190 | 0.190 | 0.190 | 0.199 | 0.199 | 0.199 | 0.155 | 0.155 | 0.155 |
| **SD8** | 0.110 | 0.110 | 0.110 | 0.144 | 0.144 | 0.144 | 0.185 | 0.185 | 0.185 |
| **SD9** | 0.280 | 0.280 | 0.280 | 0.234 | 0.234 | 0.234 | 0.175 | 0.175 | 0.175 |
| **RAD2D** | 0.861 | 0.861 | 0.878 | 0.851 | 0.851 | 0.858 | 0.456 | 0.462 | 0.472 |
| | **0.656** | **0.654** | **0.680** | **0.604** | **0.611** | **0.616** | **0.373** | **0.375** | **0.372** |

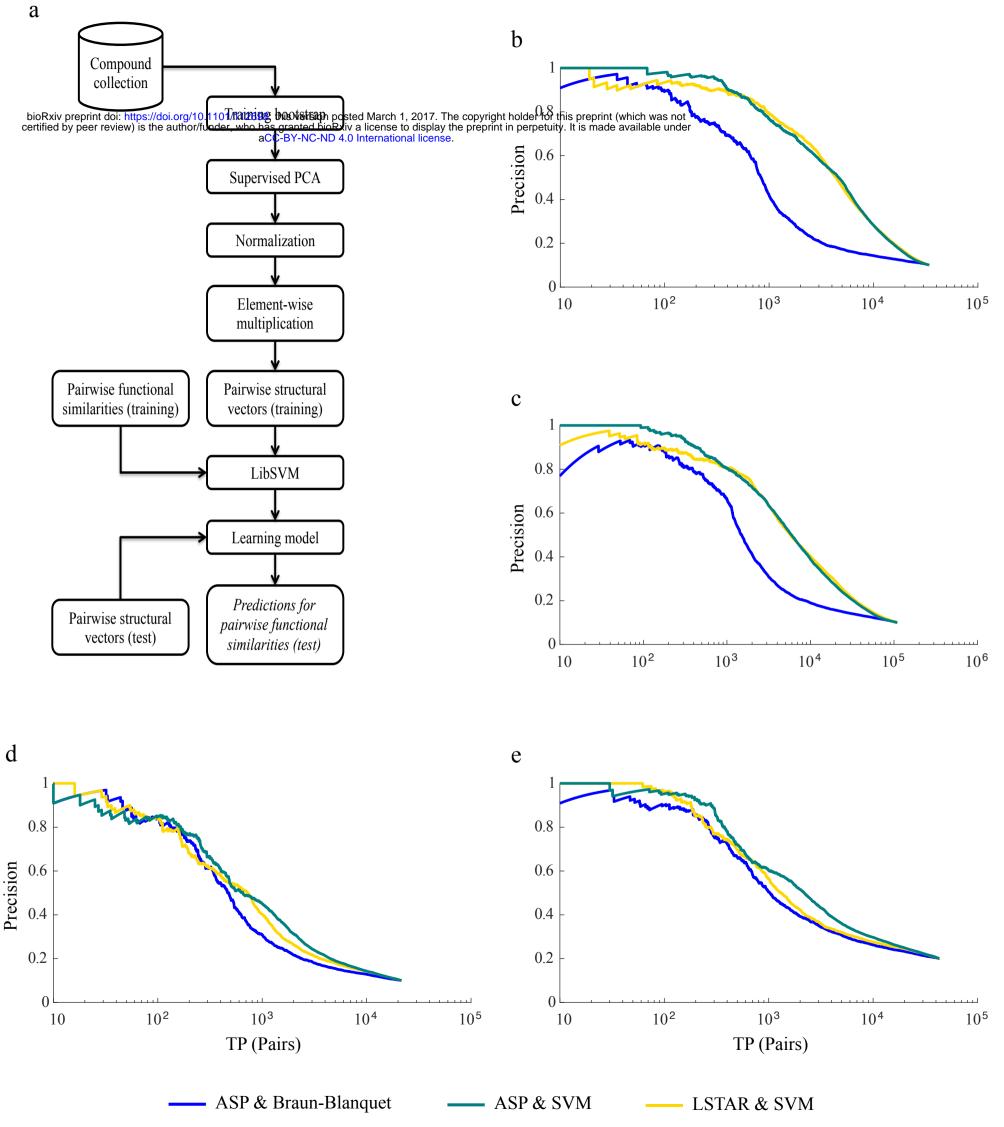| | Recall = 0.05 | | | Recall = 0.2 | | | Area Under the Curve (AUC) | | |
|---|---|---|---|---|---|---|---|---|---|
| | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** | **Cos** | **Tan** | **BB** |
| **SD1** | 0.233 | 0.235 | 0.220 | 0.153 | 0.156 | 0.156 | 0.572 | 0.575 | **0.577** |
| **ASP** | 0.272 | 0.282 | 0.290 | 0.156 | 0.161 | 0.162 | 0.567 | 0.572 | 0.576 |
| **SD3** | 0.273 | 0.273 | 0.265 | 0.154 | 0.161 | **0.165** | 0.567 | 0.571 | 0.573 |
| **SD4** | 0.282 | **0.293** | 0.291 | 0.156 | 0.161 | 0.162 | 0.566 | 0.571 | 0.575 |
| **SD5** | 0.202 | 0.208 | 0.216 | 0.128 | 0.132 | 0.137 | 0.544 | 0.550 | 0.558 |
| **LSTAR** | 0.272 | 0.279 | 0.284 | 0.150 | 0.154 | 0.160 | 0.563 | 0.568 | 0.575 |
| **SD7** | 0.129 | 0.129 | 0.129 | 0.118 | 0.119 | 0.121 | 0.519 | 0.520 | 0.520 |
| **SD8** | 0.164 | 0.164 | 0.164 | 0.136 | 0.136 | 0.136 | 0.550 | 0.550 | 0.550 |
| **SD9** | 0.148 | 0.148 | 0.148 | 0.129 | 0.129 | 0.129 | 0.546 | 0.545 | 0.542 |
| **RAD2D** | 0.243 | 0.247 | 0.270 | 0.144 | 0.147 | 0.156 | 0.561 | 0.566 | 0.574 |
| | **0.222** | **0.226** | **0.228** | **0.142** | **0.146** | **0.148** | **0.556** | **0.559** | **0.562** |

b



Figure 3

Figure 4

Figure 5

Figure 6

Figure 7

a

b

Structural similarity predicts
chemical-genetic similarity

Chemical-genetic similarity
predicts structural similarity

c

RIKEN

Machine learning model

NPD2186

NPD3120

NPD686

NPD1897

NPD450

d

NCI/NIH/GSK

Machine learning model

NSC745750

GW407323A

GSK586581B

CPD000469148

GW683134A

Figure 8