



Universiteit
Leiden
The Netherlands

Transforming data into knowledge for intelligent decision-making in early drug discovery

Paricharak, S.A.

Citation

Paricharak, S. A. (2017, February 9). *Transforming data into knowledge for intelligent decision-making in early drug discovery*. Retrieved from <https://hdl.handle.net/1887/45874>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/45874>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/45874> holds various files of this Leiden University dissertation

Author: Paricharak, S.A.

Title: Transforming data into knowledge for intelligent decision-making in early drug discovery

Issue Date: 2017-02-09

Chapter five

Data-driven Derivation of an “Informer Compound Set” for Improved Selection of Active Compounds in High- Throughput Screening (manuscript submitted)

Shardul Paricharak, Adriaan P. IJzerman, Jeremy L. Jenkins, Andreas Bender, and Florian Nigsch

Abstract

Despite the usefulness of high-throughput screening in drug discovery, for some systems, low assay throughput or high screening cost can prohibit the screening of large numbers of compounds. In such cases, iterative cycles of screening involving active learning (AL) are employed, creating the need for smaller “informer sets” that can be routinely screened to build predictive models for selecting compounds from the screening collection for follow-up screens. Here, we present a data-driven derivation of an informer compound set with improved predictivity of active compounds in HTS, and validate its benefit over randomly selected training sets on 46 PubChem assays comprising at least 300000 compounds and covering a wide range of assay biology. The informer compound set showed improvement in BEDROC($\alpha=100$), PRAUC and ROCAUC values averaged over all assays of 0.015, 0.010 and 0.016, respectively, compared to randomly selected training sets, all with paired *t*-test *p*-values $< 10^{-15}$. A per-assay assessment showed that the BEDROC($\alpha=100$), which is of particular relevance for early retrieval of actives, improved for 40 out of 46 assays, increasing the success rate of smaller follow-up screens. Overall, we showed that an informer set derived from historical HTS activity data can be employed for routine small-scale exploratory screening in an assay-agnostic fashion. This approach led to a consistent improvement in hit rates in follow up screens without compromising on scaffold retrieval. The informer set is adjustable in size depending on the number of compounds one intends to screen, as performance gains are realized for sets with more than 3000 compounds, and this set is therefore applicable to a variety of situations. Finally, our results indicate that random sampling may not adequately cover descriptor space, drawing attention to the importance of the composition of the training set for predicting actives.

Introduction

Over the past three decades, high-throughput screening (HTS) has become a well-established method used during early drug discovery.¹⁻⁷ However, low assay throughput or high screening cost can at times prohibit the screening of large numbers of compounds.^{8,9} Given this drawback, iterative cycles of design-screen-refine involving active learning (AL) strategies can be used when only a small number of compounds can or should be screened.¹⁰⁻¹² This, in combination with recent advances in machine learning, has recently

prompted efforts to improve bioactivity modeling in order to identify active compounds *in silico*, with the aim of increasing the hit rates in compound screens.¹¹

For this purpose, a high-throughput screening fingerprint (HTS-FP) was developed by Petrone *et al.*¹³ and later by Dančák *et al.*,¹⁴ which profiles compounds according to their bioactivity across a range of HTS assays. This work was based on the idea that such fingerprints are predictive of compound affinity on targets *not* covered in the fingerprint and showed the value of HTS-FP for virtual screening and biodiverse selection of actives. This concept has previously been explored computationally on smaller datasets,¹⁵⁻¹⁸ but without large-scale experimental validation. More recently, Riniker *et al.*¹⁹ benchmarked the predictive performance of chemical fingerprints and HTS-FP in conjunction with a variety of classification methods across a large number of assays performed in Novartis and those in the public domain (available in PubChem).²⁰ It was found that random forest (RF) methods with HTS-FP often outperformed machine learning methods developed on chemical descriptors.¹⁹ On a related note, Maciejewski *et al.*²¹ explored an experimental design strategy where AL was used to enhance the chemical diversity of large training sets comprising over 50000 compounds, leading to improvement in model performance. While the mentioned studies addressed the dependence of the model on descriptor and classification method used, a comprehensive assessment of how the composition of the initially screened compound set (training set) affects model performance and early retrieval of actives from the remaining screening collection was not performed.

The effectiveness of HTS screening sets in identifying actives has been widely discussed.²² Given the possible existence of over 10^{63} drug-like molecules,⁷ it is remarkable that HTS campaigns comprising “only” 10^6 compounds succeed in finding hits at all.²²⁻²⁴ A plausible explanation for this is that screening libraries are not random, but rather biased towards biogenic compounds, likely to interact with the druggable proteome. This claim has been reinforced by studies showing the chemical similarity between metabolite space, natural product space and bioactive space.²⁵⁻²⁷ A comprehensive analysis by Klekota *et al.*²⁸ showed that certain “privileged” chemical substructures, such as benzodiazepines,²⁹ enrich for bioactivity, creating further avenues for modeling the likelihood of compounds being bioactive in *any* therapeutically relevant setting (hereafter referred to as joint bioactivity modeling), rather

than target- or phenotype-specific bioactivity modeling (also shown by Gillet *et al.*).³⁰

In this study, we harnessed bioactivity information from a large number of PubChem²⁰ HTS assays to derive an assay-agnostic “informer compound set” that, once screened, predicts bioactivity better than randomly selected sets for almost all HTS assays, improving the efficiency of subsequent screens. We used AL to iteratively derive this set. Due to the difficulty in implementing AL under extreme class imbalance³¹ as is the case for all HTS assays analyzed in his study, activities from multiple assays were combined to derive binary labels representing assay-agnostic bioactivity for each compound. This was based on the idea of joint bioactivity modeling^{28,30} and led to a class-balanced dataset suitable for AL. HTS-FPs were used as descriptors, as they showed improved performance over chemical fingerprints.¹⁹ Moreover, this informer set was constructed with the aim to facilitate routine screens, as pre-composed sets are easier to screen routinely from an infrastructure point of view.

Methods

HTS Data. The public HTS data used by Riniker *et al.*¹⁹ was used in this study (see **Tables S1** and **S2** of this reference for the list of assays used). HTS data from the NIH molecular libraries program (MLP) comprising at least 300000 compounds per assay, and submitted by the NCGC, the Scripps Research Institute Molecular Screening Center, or the Burnham Center for Chemical Genomics were extracted from PubChem.²⁰ This resulted in a total of 141 cell-based and target-based assays (mainly using fluorescence readout technologies), covering a wide range of assay biology (kinases, proteases, ion channels, GPCRs and other target classes). Assay-specific z-scores were calculated for all compounds tested based on the activity measurement used to define the PubChem activity outcome. The set of assays was subsequently split into 2 groups: 95 “group 1 assays” (comprising over 338000 compounds) and 46 “group 2 assays” (comprising 300000–338000 tested compounds, depending on operational turnover of the compound collection at the screening centers). Group 1 assays (referred to as “historical assays” by Riniker *et al.*)¹⁹ were used exclusively for the construction of HTS-FP,¹³ a fingerprint used as a descriptor for machine learning, profiling the activity of a compound across HTS assays based on z-scores (float version).¹³ Group 2 assays (referred to as “test assays” by Riniker *et al.*)¹⁹ were used for deriving

labels and for model training and testing. This distinction between assay groups ensured that there was no overlap in targets between the two groups.¹⁹

HTS-FP. For each compound, an HTS-FP was computed, in which each element corresponds to the z-score (based on activity) of the compound in one of the group 1 assays. Missing z-scores (15% of all data points; not every compound is tested in each assay) were assumed to be 0 (the mean of z-scores), as implemented earlier by Riniker *et al.*¹⁹

Workflow. In this study we tested the performance of bioactivity models developed on an informer set derived with AL. First, we evaluated the performance for predicting bioactivity independent of tested assay (**Figure 23**, joint bioactivity modeling).

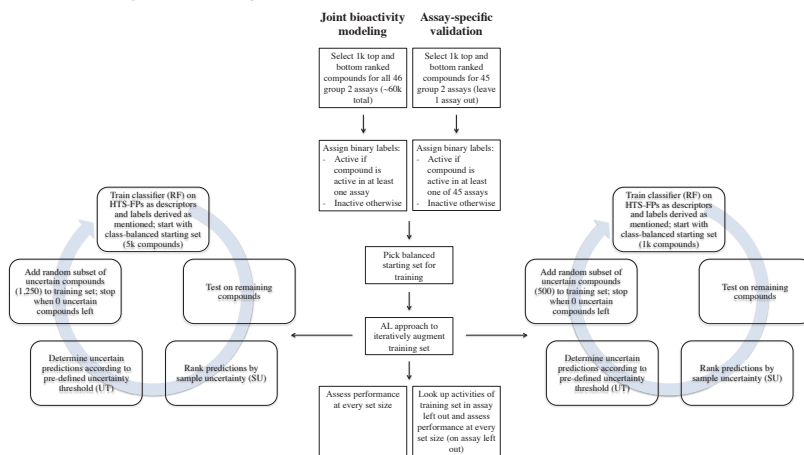


Figure 23. Overview of workflow. In this study, two analyses were performed. Firstly (left), a joint bioactivity model was developed on the 1000 top and bottom ranked (based on z-scores) compounds. An AL approach was used to iteratively augment the training set, for which model performance (ROCAUC) was assessed at every set size. The second analysis (right) involved an assay-specific validation, where a joint bioactivity model was developed on all assays except the assay left out of training. The training set was iteratively augmented with uncertain samples using AL, and at every set size, activities of these compounds were looked up in the assay left out. Subsequently, model performance (ROCAUC, PRAUC, BEDROC) for the training set was assessed on the assay left out, rather than on the joint activities dataset.

Here, activities from group 2 assays were combined to derive binary labels representing assay-agnostic bioactivity for each compound in order to construct a class-balanced dataset suitable for AL. Improved model performance at this step was considered a prerequisite for the more challenging task of predicting actives for individual assays. An assay-specific validation was performed to address the latter task: the informer set was

derived from activity data from 45 group 2 assays and predictivity was assessed on the one assay remaining (**Figure 23**, assay-specific validation). This was repeated 46 times, effectively leaving each group 2 assay out once.

Joint Bioactivity Modeling. The 1000 least and most active compounds (based on z-scores) were selected from each group 2 assay, resulting in a total of 58768 compounds. A skewed distribution of the number of assays these compounds were active in was observed, with 45%, 33%, 12% and 10% of compounds active in 0, 1, 2 and more than 2 assays, respectively. Each compound was labeled as “active” if it was active in *any* of the group 2 assays (as defined by the PubChem activity outcome) or “inactive” otherwise, resulting in a total of 32171 actives and 26597 inactives. This labeling was based on the concept of considering activities independent of the assay they were tested in (joint bioactivity). An RF model (scikit-learn)³²⁻³⁴ was developed on a randomly selected class-balanced training set of 5000 compounds, and the performance of the model was assessed on the remaining compounds. Using AL, this training set was iteratively augmented with up to 1250 uncertain samples at each iteration, with the aim to improve model performance on the remaining compounds (see “Active Learning” section for more details). The model for this informer set was benchmarked against a model developed on a randomly selected set at each set size using the area under the receiver operating characteristic curve (ROCAUC).

Assay-specific Validation. Here, the informer set was derived from activity data from 45 group 2 assays, and a model was trained on group 1 assay HTS-FPs and labels derived from the one assay left out. The starting set was a randomly selected class-balanced set of 1000 compounds, which was iteratively augmented by up to 500 compounds using AL (see “Active Learning” section for more details). The size of the training and augmentation set was kept smaller here than for the joint bioactivity modeling due to observed improvement in performance at the earlier stages of the algorithm. Performance on the assay left out was assessed at each set size using the ROCAUC, the area under the precision-recall curve (PRAUC),³⁵ Boltzmann-enhanced discrimination of ROC (BEDROC) ($\alpha=100$),^{36,37} and the retrieval of Murcko scaffolds³⁸ belonging to the active compounds. The BEDROC($\alpha=100$) was used due to its relevance in early retrieval of actives in imbalanced datasets and the PRAUC was used because it captures the effect of the large number of inactive compounds on the model’s performance.³⁵ Both these metrics were therefore considered more relevant than the ROCAUC for the

assay-specific validation (by contrast, for the joint bioactivity modeling the ROCAUC was considered an adequate metric due to class balance).

The model was benchmarked against models developed on a randomly selected set and a set comprising compounds with the highest median z-scores across the 45 assays left in (the frequent hitter set). The latter comparison was included to ensure that the performance gain for the informer set was better than when simply more actives from other assays (including more frequent hitters) were trained on.

Machine Learning. The RF parameters used were: 100 trees (no maximum depth), minimum samples to split = 2, and minimum samples for a leaf = 2.

Active Learning (AL). The AL approach consisted of three iterative steps: (1) training of an RF model, (2) model testing on the remaining compounds and (3) augmenting the training set with a randomly selected subset of uncertain labeled samples (1250 and 500 compounds for the joint bioactivity modeling and assay-specific validation, respectively); when the number of uncertain samples was smaller than the size of the subset, all uncertain samples were selected. The AL algorithm was terminated when the number of uncertain samples was zero. Sample uncertainty (SU) of a given compound c was defined as the absolute probability difference in active versus inactive class predictions:

$$SU_c = |p_c^{active} - p_c^{inactive}| \quad (4)$$

with SU_c in the range of 0–1 where 0 and 1 represent the most uncertainty and complete certainty in prediction, respectively. Only samples with an SU value smaller than the uncertainty threshold (UT) were considered uncertain. We investigated the effect of varying the UT from 0.5 (least stringent) to 0.01 (most stringent) for the joint bioactivity modeling, and used a UT of 0.1 for the assay-specific validation. The presence of uncertain samples suggests undersampling of bioactivity space. Including these samples could improve model performance over random sampling.¹⁰

Software Used. The workflow comprised Python scripts for data analysis, using scikit-learn³⁴ for machine learning and RDKit³⁹ for scaffold derivation. Tableau⁴⁰ was used for data exploration and R⁴¹ was used for the visualization of results.

Results and discussion

The development of an informer set for the prediction of joint bioactivity is presented first (see **Figure 23** – left). Prediction of joint bioactivity allowed the identification of compounds more likely to be bioactive regardless of the assay used. This was followed by a performance assessment of the informer set on individual assays (assay-specific validation; see **Figure 23** – right), and an analysis of scaffold retrieval and set composition. The assay-specific validation was performed in order to determine whether the informer set is more useful than a randomly selected set in predicting actives for novel assays one might perform.

Joint Bioactivity Modeling. The gap in ROCAUC between models developed on the AL sets and on randomly selected sets consistently widens from set sizes of ~5000 onwards (see **Figure 24** – top).

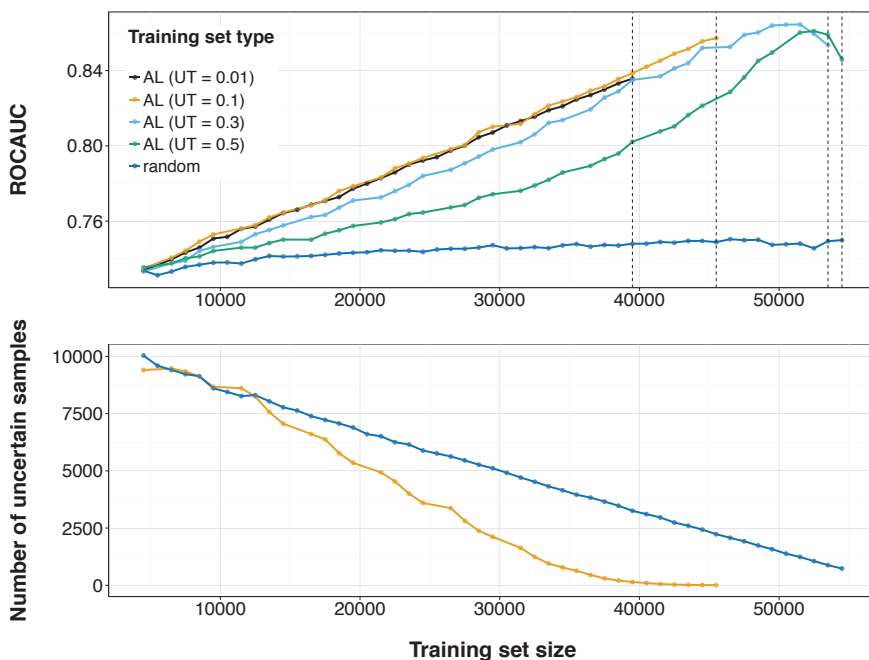


Figure 24. Comparison of model performance for the AL and randomly selected training sets. The ROCAUC (top) is shown for the models trained on AL and randomly selected sets. Performance across all set sizes is consistently better for all AL sets than it is for the randomly selected set. At a set size of 40000 an average gain in performance of 0.08 is observed. In addition, lower *UT* values led to better performance than higher *UT* values. A *UT* value of 0.1 was chosen for the assay-specific validation on the basis of a trade-off between improvement in performance and maximum training set size. For the AL set (*UT* = 0.1), the number of uncertain samples declines faster compared to the randomly selected set (bottom), indicating more efficient sampling of bioactivity space.

At a set size of 40000 an average gain in performance of 0.08 is observed for the AL sets (average ROCAUC of 0.83 compared to 0.75 for randomly selected sets). Stringent *UT* values led to sets with a greater gain in performance at the cost of maximum set size, as fewer samples are classified as uncertain, and the number of uncertain samples reduces to zero earlier in the AL process. Moreover, the number of uncertain samples declines faster for the AL (*UT* = 0.1) set than for the randomly selected set (**Figure 24** – bottom), indicating the benefit of AL in sampling relevant bioactivity space more efficiently. For example, almost all uncertain samples were exhausted for a set size of approximately 40000 using AL, whereas the random set did not exhaust the uncertain samples even at set sizes upwards of 50000. These results indicate that AL is able to consistently sample descriptor space better than random sampling, hereby improving the identification of compounds bioactive in one or more group 2 assays. For further analysis, we chose a *UT* value of 0.1 on the basis of a trade-off between gain in performance and maximum training set size.

Predictive Performance of Informer Set on Individual Assays. In an attempt to translate performance gain in predicting joint bioactivity (see previous section) to performance gain in individual large-scale assays, we performed an assay-specific validation for all group 2 assays. Improved predictive performance in this setting would corroborate the usefulness of an informer set, as no prior information about the assay left out would be required for its construction.

The BEDROC($\alpha=100$),^{36,37} PRAUC and ROCAUC were calculated for an RF classifier trained on the informer set (AL), a randomly selected set, and the frequent hitter set. These values were averaged over all 46 assay-specific validation experiments and were binned by set size (see **Figure 25**).

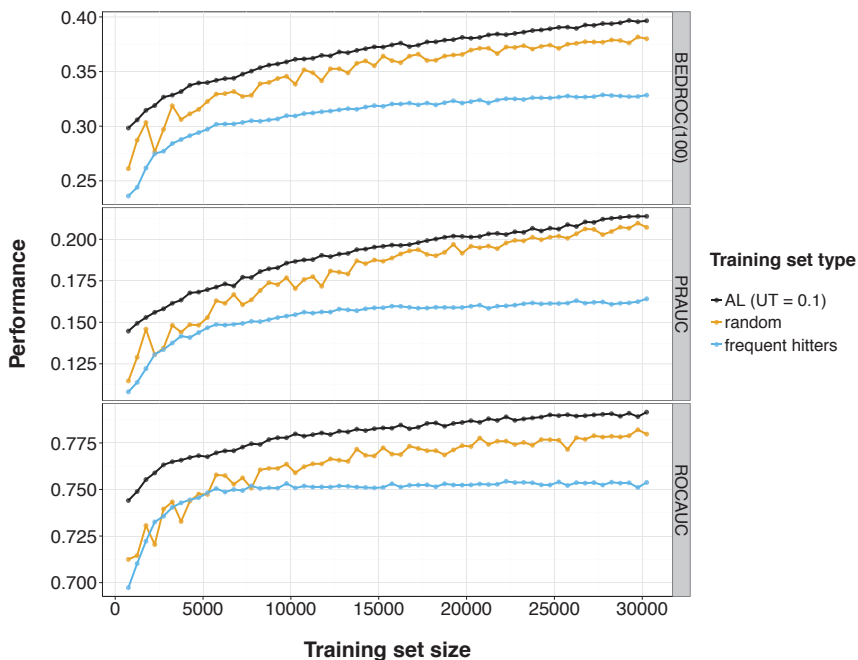


Figure 25. Comparison of model performance for the AL ($UT = 0.1$), random and frequent hitter training sets (assay-specific validation). The BEDROC($\alpha=100$)^{36,37} (top), PRAUC (middle) and ROCAUC (bottom) binned by set size are shown for all three training sets (bin width=500). The assay-averaged performance for the AL set (all metrics) is consistently better than that for the randomly selected set. For the frequent hitter set, performance is consistently worse than both the AL set and the randomly selected set for training sets larger than 5000 compounds. These results indicate that models trained on the AL set consistently retrieve more actives compared to models trained on the other sets.

The frequent hitter set was used as a benchmark, to ensure that the performance gain of the AL set was better than when simply more actives from other assays (including more frequent hitters) were trained on.

Overall, the performance for the AL set was enhanced compared to the randomly selected set, with an average increase of 0.015, 0.010 and 0.016 in average BEDROC, PRAUC and ROCAUC, respectively (all with paired t -test p -values $< 10^{-15}$). The apparent low values of the average BEDROC (0.25-0.40) can be explained by the Boltzmann enhancement, as early retrieval of actives is strongly preferred. Low values of the average PRAUC metric (0.10-0.25) can be explained by the extreme class imbalance: a random classifier would achieve a PRAUC of ~ 0.007 given the average fraction of actives is only $\sim 0.7\%$. For the frequent hitter set, performance is consistently worse for set sizes larger than 5000, indicating that simply including more actives from other

assays does not account for the performance gain observed for the informer set. This finding is in line with the results of the “weak reinforcement strategy” as described in the study by Maciejewski *et al.*²¹ Here, training sets with a large number of actives similar in descriptor space (including frequent hitters^{42,43} in our study, as the descriptor space is based on bioactivity profiles) were found to be poor at identifying the remaining small number of actives in the test set due to insufficient coverage of descriptor space. By contrast, training sets containing compounds outside the applicability domain, corresponding to uncertain samples in this study, were much better at identifying the remaining actives in the screening collection.

Next, the average improvement in performance over all set sizes of the informer set was calculated separately for each assay (see **Figure 26**).

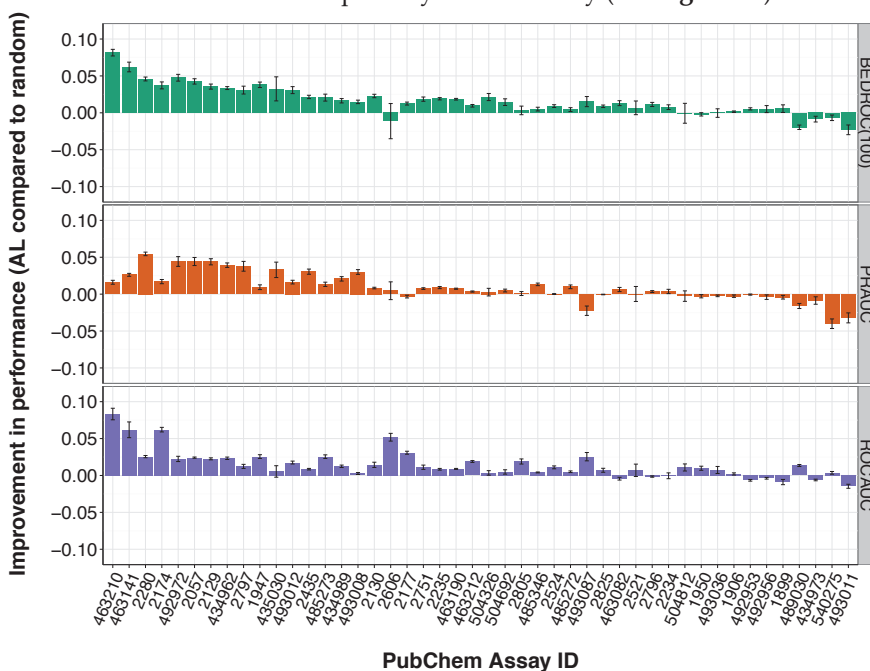


Figure 26. Improvement in model performance for the AL ($UT = 0.1$) set compared to the randomly selected set for separate assays. The average difference in BEDROC($\alpha=100$)^{36,37} (top), PRAUC (middle) and ROCAUC (bottom) between the AL set and the randomly selected set is shown for separate assays. Error bars represent standard error of the mean. For 25 out of 46 assays all three metrics improved, whereas the BEDROC($\alpha=100$), which is of most relevance for early retrieval of actives,^{36,37} improved for 40 out of 46 assays. In practice, the results indicate that if a subsequent screen were performed for each assay, more actives would be retrieved for 40 assays, compared to when random training sets would be used.

For 29 out of 46 assays, all three metrics improved by average 0.02 on average, whereas the BEDROC, which is of most relevance for early retrieval of actives,^{36,37} improved for 40 out of 46 assays by 0.02 on average. The best increase in performance was observed for assays number 463210 (caspase 7), 463141 (caspase 3), 2280 (GLD-1 protein) and 2174 (lysophospholipase 1), with BEDROC improvements of 0.08, 0.06, 0.05 and 0.04, respectively. While improvement was modest for most assays, it was consistent, as shown by the error bars representing the standard error of the mean difference in performance between the informer set and the randomly selected set across all sizes. Given the relatively small training sets, varying in size from ~0.3% to 10% of the entire screening collection, large improvements in predictive power over the remaining 90%-99.7% would be unrealistic. We attempted to investigate the cause for the performance loss for the remaining 6 assays, but could not find an explanation: there was no apparent relationship with the average performance for that assay, nor the number of actives in that assay.

Scaffold Retrieval for Individual Assays. We analyzed the scaffold retrieval rate (defined as the retrieved percentage of unique scaffolds belonging to active compounds in the test set; see **Figure 27** – top) and the median z-scores (see **Figure 27** – bottom) of actives identified in the top 5% ranked compounds in order to assess whether these actives were enriched for frequent hitters.

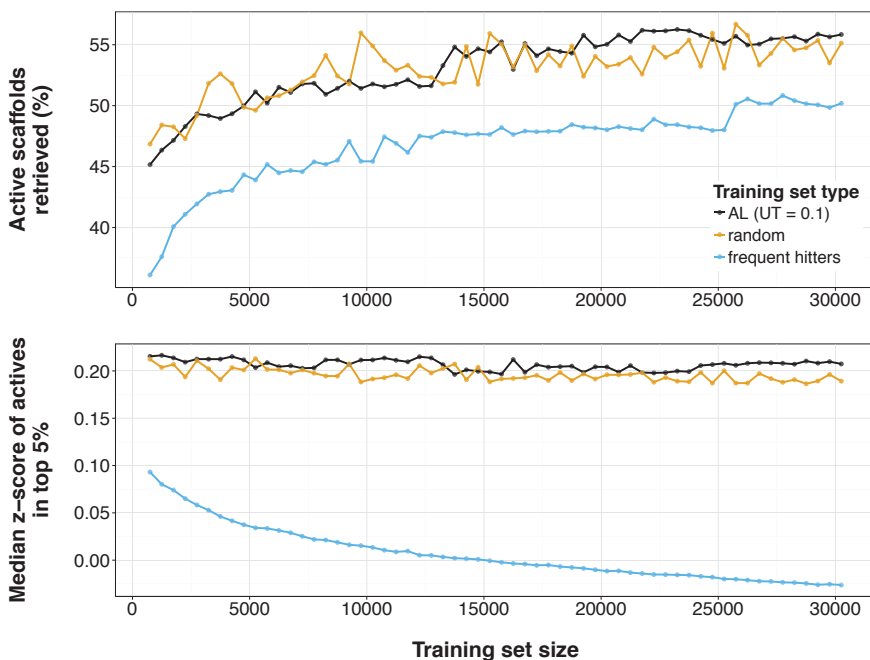


Figure 27. Active scaffold retrieval (%) and median z-scores of actives in top 5% (assay-specific validation). Similar values in scaffold retrieval and the median z-scores of actives in the top 5% ranked compounds for the AL set and the randomly selected set indicate that the AL approach does not compromise on the scaffold retrieval of active compounds, nor does it substantially enrich for frequent hitters. For the frequent hitter set, scaffold retrieval is consistently reduced, hence showing that simply including active compounds from other assays in the training set does not improve the retrieval of diverse sets of actives.

Similar values in scaffold retrieval (45%-55%) and median z-scores (~0.20) for the AL set and the randomly selected set indicate that the AL approach does not compromise on the retrieval of diverse sets of active compounds, nor does it substantially enrich for frequent hitters. Remarkably, the fluctuation in scaffold retrieval is somewhat higher for the random set. This finding is not surprising, as the number of active scaffolds in the training set (which varies in different random sets of compounds due to chance) determines scaffold retrieval in the test set. By contrast, the AL set is iteratively augmented with uncertain compounds that are more likely to have different scaffolds, and hence shows less fluctuation in scaffold retrieval. The frequent hitter set consistently shows worse performance than the other two sets in scaffold retrieval. In addition, the median z-score of the actives retrieved consistently drops from 0.09 to below 0 (**Figure 27** – bottom). The latter drop is likely caused due to fewer compounds with high median z-scores remaining in the

test set as training set size increases. Relative stability of the median z-score is observed for both the AL and random sets, indicating no enrichment for frequent hitters in the training set. In summary, we conclude that when the AL approach is used the scaffold retrieval is not impaired, frequent hitters are not enriched for and at the same time overall hit rates are improved.

Composition of informer set. In order to analyze the composition of the informer set in more detail, we calculated the fraction of the number of active compounds picked from the group 2 assays relative to the number of active compounds for each assay (see **Figure 28**).

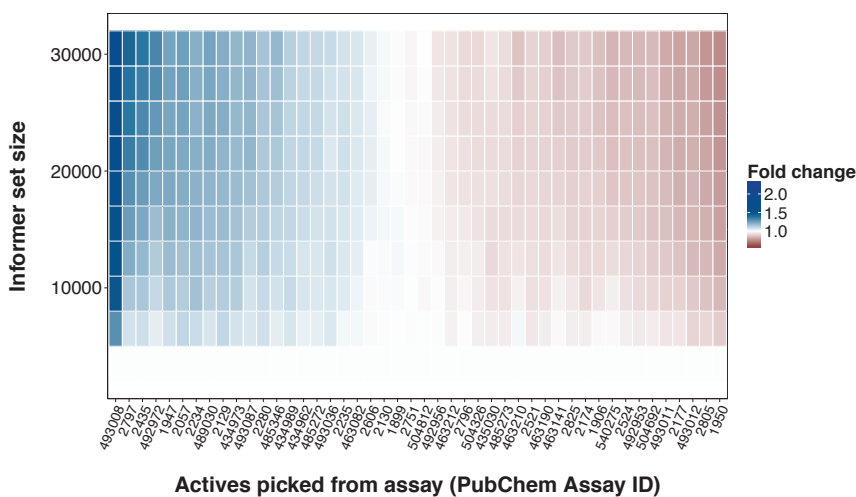


Figure 28. Composition of the informer set in terms of active compounds selected from group 2 assays. The heat map represents the composition of the informer set at varying sizes in terms of the fraction of the number of active compounds selected from group 2 assays relative to the number of active compounds for each assay. On the one hand, active compounds from assays number 493008 (troponin C type 1), 2797 (vasopressin V1a receptor), 2435 (oxytocin receptor) and 492972 (platelet-activating factor acetylhydrolase 1b subunit γ) are consistently overrepresented (fold change > 1.3 at a set size of 30000). On the other hand, active compounds from assays number 1950 (EBNA-1 protein), 2805 (intestinal alkaline phosphatase), 493012 (DNA deaminase APOBEC-3G) and 2177 (lysophospholipase 2) are underrepresented (fold change < 0.8 at a set size of 30000). While the AL approach improves performance for all assays, the average BEDROC($\alpha=100$) is much higher for the assays with overrepresented actives (0.75) than for the assays with underrepresented actives (0.20).

On the one hand active compounds from assays number 493008 (troponin C type 1), 2797 (vasopressin V1a receptor), 2435 (oxytocin receptor) and 492972 (platelet-activating factor acetylhydrolase 1b subunit γ) are consistently overrepresented in the informer set (maximum fold change > 1.3) while on the other hand active compounds from assays number 1950 (EBNA-1 protein),

2805 (intestinal alkaline phosphatase), 493012 (DNA deaminase APOBEC-3G) and 2177 (lysophospholipase 2) are underrepresented (minimum fold change < 0.8). While the AL approach improves performance for all the assays mentioned above (see **Figure 26**), interestingly, the average BEDROC is much higher for those assays of which the active compounds are *overrepresented* (0.75) than for the assays of which the active compounds are *underrepresented* (0.20). This indicates that more actives are picked from assays already exhibiting good performance. In addition, as determined by Riniker *et al.*¹⁹ the assays of which the active compounds are overrepresented share over 20% of actives with at least six group 1 assays, whereas the assays of which the active compounds are underrepresented share over 20% actives with only at most one group 1 assay (group 1 assays were used to define descriptor space), explaining the difference in BEDROC between these assays.

We attempted to investigate whether bias towards active compounds from particular assays in the informer set was related to improvement in performance over models trained on randomly selected sets for those assays, but could not find any link. We therefore conclude that this improvement in performance is due to better sampling of bioactivity space, as the AL approach iteratively augments the informer set with uncertain samples.

Conclusions

Strategies involving iterative cycles of feedback-driven compound selection and testing can be used when low assay throughput or high screening cost hinders the screening of large compound libraries. This creates the need for the exploratory screening of smaller informer sets to build predictive models for compound selection for follow-up testing. In this study, we performed a data-driven construction of an informer compound set with improved retrieval of actives in a subsequent selection round for apparently unrelated HTS assays. The benefit of this informer set was validated over randomly selected training sets on 46 PubChem²⁰ assays comprising at least 300000 compounds. Overall, we highlight that such a set – of adjustable size, depending on the number of compounds one intends to screen – can be employed for routine exploratory screening in an assay-agnostic fashion for a gain in predictive power.

Averaged over all assays, an improvement in BEDROC, PRAUC and ROCAUC (of 0.015, 0.010 and 0.016 respectively) was observed with respect to random training sets, all with paired *t*-test *p*-values < 10⁻¹⁵. The informer set

improved the BEDROC for 40 out of 46 assays, indicating better early retrieval of actives. In addition, we found that our approach did not compromise on the retrieval of diverse sets of active compounds, nor did it enrich for frequent hitters, as both scaffold retrieval and the median z-score activity of the actives retrieved were unaffected. The informer set overrepresented actives from certain assays, and underrepresented actives from other assays. Interestingly, while the informer set increased performance for both groups of assays, the BEDROC was much higher (0.75) for the assays of which the actives were overrepresented, than for assays with underrepresented actives (0.20).

We conclude that our AL approach is able to more effectively sample descriptor space, expected to improve the retrieval of active compounds in subsequent screens, thereby reducing the time and expense required to arrive at the same number of hits.

References

- (1) Macarron, R. (2006) Critical review of the role of HTS in drug discovery. *Drug Discov. Today* 11, 277–279.
- (2) Mayr, L. M., and Fuerst, P. (2008) The future of high-throughput screening. *J. Biomol. Screen.* 13, 443–448.
- (3) Phatak, S. S., Stephan, C. C., and Cavasotto, C. N. (2009) High-throughput and in silico screenings in drug discovery. *Expert. Opin. Drug Discov.* 4, 947–959.
- (4) Mayr, L. M., and Bojanic, D. (2009) Novel trends in high-throughput screening. *Curr. Opin. Pharmacol.* 9, 580–588.
- (5) Valler, M. J., and Green, D. (2000) Diversity screening versus focussed screening in drug discovery. *Drug Discov. Today* 5, 286–293.
- (6) Fox, S., Farr-Jones, S., Sopchak, L., Boggs, A., Nicely, H. W., Khoury, R., and Biros, M. (2006) High-Throughput Screening: Update on Practices and Success. *J. Biomol. Screen.* 11, 864–869.
- (7) Koutsoukas, A., Paricharak, S., Galloway, W. R. J. D., Spring, D. R., IJzerman, A. P., Glen, R. C., Marcus, D., and Bender, A. (2014) How diverse are diversity assessment methods? A comparative analysis and benchmarking of molecular descriptor space. *J. Chem. Inf. Model.* 54, 230–242.
- (8) Bajorath, J. (2002) Integration of virtual and high-throughput screening. *Nat. Rev. Drug. Discov.* 1, 882–894.
- (9) Astashkina, A., Mann, B., and Grainger, D. W. (2012) A critical evaluation of in vitro cell culture models for high-throughput drug screening and toxicity. *Pharmacol. Ther.* 134, 82–106.
- (10) Settles, B. (2010) Active Learning Literature Survey. *Mach. Learn.* 15, 201–221.
- (11) Reker, D., and Schneider, G. (2015) Active-learning strategies in computer-assisted drug discovery. *Drug Discov. Today* 20, 458–465.
- (12) Paricharak, S., IJzerman, A. P., Bender, A., and Nigsch, F. (2016) Analysis of iterative screening with stepwise compound selection based on Novartis in-house HTS data. *ACS Chem. Biol.* 11, 1255–1264.

- (13) Petrone, P. M., Simms, B., Nigsch, F., Lounkine, E., Kutchukian, P., Cornett, A., Deng, Z., Davies, J. W., Jenkins, J. L., and Glick, M. (2012) Rethinking Molecular Similarity: Comparing Compounds on the Basis of Biological Activity. *ACS Chem. Biol.* 7, 1399–1409.
- (14) Dančik, V., Carrel, H., Bodycombe, N. E., Seiler, K. P., Fomina-Yadlin, D., Kubicek, S. T., Hartwell, K., Shamji, A. F., Wagner, B. K., and Clemons, P. A. (2014) Connecting Small Molecules with Similar Assay Performance Profiles Leads to New Biological Hypotheses. *J. Biomol. Screen.* 19, 771–781.
- (15) Kauvar, L. M., Higgins, D. L., Villar, H. O., Sportsman, J. R., Engqvist-Goldstein, Å., Bukar, R., Bauer, K. E., Dilley, H., and Roche, D. M. (1995) Predicting ligand binding to proteins by affinity fingerprinting. *Chem. Biol.* 2, 107–118.
- (16) Bender, A., Jenkins, J. L., Glick, M., Deng, Z., Nettles, J. H., and Davies, J. W. (2006) “Bayes affinity fingerprints” improve retrieval rates in virtual screening and define orthogonal bioactivity space: when are multitarget drugs a feasible concept? *J. Chem. Inf. Model.* 46, 2445–2456.
- (17) Nguyen, H. P., Koutsoukas, A., Mohd Fauzi, F., Drakakis, G., Maciejewski, M., Glen, R. C., and Bender, A. (2013) Diversity Selection of Compounds Based on “Protein Affinity Fingerprints” Improves Sampling of Bioactive Chemical Space. *Chem. Biol. Drug Des.* 82, 252–266.
- (18) Givehchi, A., Bender, A., and Glen, R. C. (2006) Analysis of activity space by fragment fingerprints, 2D descriptors, and multitarget dependent transformation of 2D descriptors. *J. Chem. Inf. Model.* 46, 1078–1083.
- (19) Riniker, S., Wang, Y., Jenkins, J. L., and Landrum, G. A. (2014) Using Information from Historical High-Throughput Screens to Predict Active Compounds. *J. Chem. Inf. Model.* 54, 1880–1891.
- (20) Wang, Y., Xiao, J., Suzek, T. O., Zhang, J., Wang, J., Zhou, Z., Han, L., Karapetyan, K., Dracheva, S., Shoemaker, B. A., Bolton, E., Gindulyte, A., and Bryant, S. H. (2012) PubChem’s BioAssay Database. *Nucl. Acids Res.* 40, D400–D412.
- (21) Maciejewski, M., Wassermann, A. M., Glick, M., and Lounkine, E. (2015) An Experimental Design Strategy: Weak Reinforcement Leads to Increased Hit Rates and Enhanced Chemical Diversity. *J. Chem. Inf. Model.* 55, 956–962.
- (22) Hert, J., Irwin, J. J., Laggner, C., Keiser, M. J., and Shoichet, B. K. (2010) Quantifying Biogenic Bias in Screening Libraries. *Nat. Chem. Biol.* 5, 479–483.
- (23) Fox, S., Farr-Jones, S., Sopchak, L., Boggs, A., and Comley, J. (2004) High-throughput screening: searching for higher productivity. *J. Biomol. Screen.* 9, 354–358.
- (24) Pereira, D. A., and Williams, J. A. (2007) Origin and evolution of high throughput screening. *Br. J. Pharmacol.* 152, 53–61.
- (25) Ertl, P., Roggo, S., and Schuffenhauer, A. (2008) Natural product-likeness score and its application for prioritization of compound libraries. *J. Chem. Inf. Model.* 48, 68–74.
- (26) Gupta, S., and Aires-de-Sousa, J. (2007) Comparing the chemical spaces of metabolites and available chemicals: Models of metabolite-likeness. *Mol. Divers.* 11, 23–36.
- (27) O’Hagan, S., Swainston, N., Handl, J., and Kell, D. B. (2015) A “rule of 0.5” for the metabolite-likeness of approved pharmaceutical drugs. *Metabolomics* 11, 323–339.
- (28) Klekota, J., and Roth, F. P. (2008) Chemical substructures that enrich for biological activity. *Bioinformatics* 24, 2518–2525.
- (29) Evans, B. E., Rittle, K. E., Bock, M. G., DiPardo, R. M., Freidinger, R. M., Whitter, W. L., Lundell, G. F., Veber, D. F., and Anderson, P. S. (1988) Methods for drug discovery: development of potent, selective, orally effective cholecystokinin antagonists. *J. Med. Chem.* 31, 2235–2246.
- (30) Gillet, V. J., Willett, P., and Bradshaw, J. (1997) Identification of biological activity profiles using substructural analysis and genetic algorithm. *J. Chem. Inf. Comput. Sci.* 38, 165–179.

- (31) Attenberg, J., and Ertekin, S. (2013) Class imbalance and active learning, in *Imbalanced Learning: Foundations, Algorithms, and Applications, First Edition* (He, H., and Ma, Y., Eds.), pp 101–149. John Wiley & Sons, Inc.
- (32) Breiman, L. (2001) Random Forests. *Mach. Learn.* 45, 5–32.
- (33) Riniker, S., Fechner, N., and Landrum, G. A. (2013) Heterogeneous Classifier Fusion for Ligand-Based Virtual Screening: Or, How Decision Making by Committee Can Be a Good Thing. *J. Chem. Inf. Model.* 53, 2829–2836.
- (34) Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, É. (2011) Scikit-learn : Machine Learning in Python. *J. Mach. Learn. Res.* 12, 2825–2830.
- (35) Davis, J., and Goadrich, M. (2006) The Relationship Between Precision-Recall and ROC Curves, in *Proceedings of the 23rd International Conference on Machine learning*, pp 233–240.
- (36) Truchon, J., and Bayly, C. I. (2007) Evaluating virtual screening methods: Good and bad metrics for the “early recognition” problem. *J. Chem. Inf. Model.* 47, 488–508.
- (37) Riniker, S., and Landrum, G. A. (2013) Open-source platform to benchmark fingerprints for ligand-based virtual screening. *J. Cheminform.* 5, 26–42.
- (38) Bemis, G. W., and Murcko, M. A. (1996) The properties of known drugs. 1. Molecular frameworks. *J. Med. Chem.* 39, 2887–2893.
- (39) RDKit: cheminformatics and machine learning software (<http://www.rdkit.org/>); 2013.
- (40) Tableau Desktop, version 9.0.1; Tableau Software Inc., 2015.
- (41) Dessau, R. B., and Pipper, C. B. (2008) R–project for statistical computing. *Ugeskr. Laeger.* 170, 328–330.
- (42) Baell, J., and Walters, M. A. (2014) Chemical con artists foil drug discovery. *Nature* 513, 481–483.
- (43) Che, J., King, F. J., Zhou, B., and Zhou, Y. (2012) Chemical and Biological Properties of Frequent Screening Hits. *J. Chem. Inf. Model.* 52, 913–926.