

Identifying Biologically Active Compound Classes Using Phenotypic Screening Data and Sampling Statistics

Justin Klekota,* Erik Brauner,* and Stuart L. Schreiber

Howard Hughes Medical Institute, Harvard Institute of Chemistry and Cell Biology, Broad Institute of Harvard and MIT, Harvard University, 12 Oxford Street, Cambridge, Massachusetts 02138

Received March 16, 2005

Scoring the activity of compounds in phenotypic high-throughput assays presents a unique challenge because of the limited resolution and inherent measurement error of these assays. Techniques that leverage the structural similarity of compounds within an assay can be used to improve the hit-recovery rate from screening data. A technique is presented that uses clustering and sampling statistics to predict likely compound activity by scoring entire structural classes. A set of phenotypic assays performed against a commercially available compound library was used as a test set. Using the class-scoring technique, the resultant activity prediction scores were more reproducible than individual assay measurements, and class scoring recovered known active compounds more efficiently than individual assay measurements because class scoring had fewer false positives. Known biologically active compounds were recovered 87% of the time using class scores, suggesting a low false-negative rate that compared well to individual assay measurements. In addition, many weak and potentially novel classes of active compounds, overlooked by individual assay measurements, were suggested.

INTRODUCTION

Chemical genetics experiments are based on the notion that small molecules can emulate the effects of genetic mutations by inducing interesting cellular phenotypes and thereby uncovering novel biological targets and functionalities. This process parallels drug discovery, which seeks to identify chemical modulators of biological targets of known therapeutic value. Notably, the advent of high-throughput assay technology and combinatorial chemistry has dramatically increased the number of compounds screened.^{1–3} Computational tools have played an increasing role in optimizing the identification of potent and bioavailable compounds,^{4,5} and these tools have potential applications in scoring phenotypic assay screening data.⁶

When computer learning (such as decision trees or clustering) is used to identify *multiple* active structural classes in one round of screening, libraries increasingly enriched in active structures can be synthesized for subsequent rounds of screening.^{7–17} Similar methods use Bayesian learning to identify active structural classes by scoring active substructures^{18–21} or inactive substructures²² using screening data. When precision measurements such as EC50 values or binding constants are available, quantitative structure–activity relationship models are used to optimize structurally the activity of an individual lead compound class (having a common structural skeleton) using linear regressions with graphical and physicochemical descriptors,^{23–30} nearest neighbor relationships,^{13–16,31–41} molecular alignment,^{42–47} and simulated binding.^{48–53}

While drug discovery technologies such as decision trees can model the activity of multiple structural classes of active

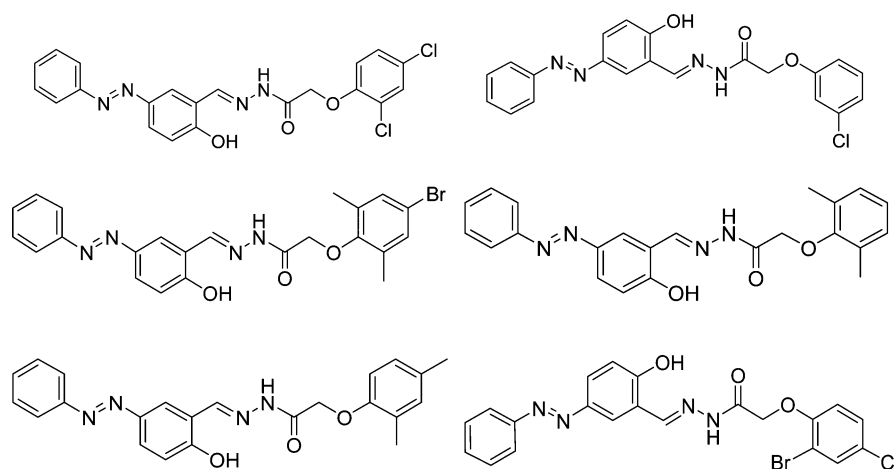
compounds using “low resolution” screening data, they are most often, but not always,⁵⁴ applied to a single protein target.^{9–12} By comparison, compounds scored in phenotypic assays can modulate multiple biological targets to induce the same phenotype. Computational technologies that can cope with the added complexity of multiple unknown targets and the multiple compound classes that act on them are required to predict activity in phenotypic assays reliably.

If the chemical similarity principle holds and structurally similar compounds have similar biological activities in phenotypic assays, then clustering should distinguish the numerous classes of compounds accessed by phenotypic screens and allow the activity of each compound class to be scored. To this end, sampling statistics can be used to measure the significance of the observed number of active compounds within each compound class. These computational tools, namely, clustering and sampling statistics, should be capable of identifying classes of active compounds with high reproducibility. In addition, the scoring of individual compounds as classes should also make more reliable predictions of activity (with fewer false positives) than the scoring of compounds using a single assay measurement, which is the most commonly employed approach. Weaker hits and active compounds not necessarily tested in an assay should also be identified as active by the class-scoring approach, suggesting a higher hit-recovery rate. Furthermore, structure–activity relationships within each class could be inferred between active and inactive molecules.

METHODS

Chemical Descriptors. Daylight fingerprints, a broadly used commercially available descriptor set,^{13–16,34,39,55} were used to represent each compound and then identify com-

* Author to whom correspondence should be addressed. E-mail: Erik_Brauner@hms.harvard.edu (E.B.), JKDklekota@aol.com (J.K.).

Chart 1. Cluster with Common Structural Backbone^a

^a This cluster demonstrates the “ideal” case when all the compounds in a cluster have a common structural backbone. The limited structural heterogeneity increases the chance that compounds in this cluster will have common biological activity.

pound classes through clustering. Daylight fingerprints containing 4096 bits were used to encode two-dimensional substructures (i.e., no specified stereochemistry) up to seven bonds in length for each compound in the Chembridge Diverse Set E library. If fewer bits were used in the fingerprint, the number of bits shared by the encoded substructures would have increased, diminishing the fingerprint's ability to differentiate compound structures and its ability to distinguish biologically active structures from inactive ones.⁵⁵ 2D fingerprints are satisfactory because combinatorial libraries often contain enantiomeric mixtures anyway and previous publications suggest that 3D fingerprints perform no better than 2D ones to distinguish active and inactive molecules.^{13,56}

All chemical libraries were represented as SD files and were subsequently converted to nonisomeric canonical SMILES using Daylight's *mol2smi* algorithm. Once in canonical SMILES format, salts and charges were removed from the compound structures using Perl scripts. A module obtained from cpan.org was modified to accommodate SMILES editing.⁵⁷ After salts and charges were removed, the 4096-bit Daylight fingerprints were then calculated directly from the nonisomeric canonical SMILES representations.

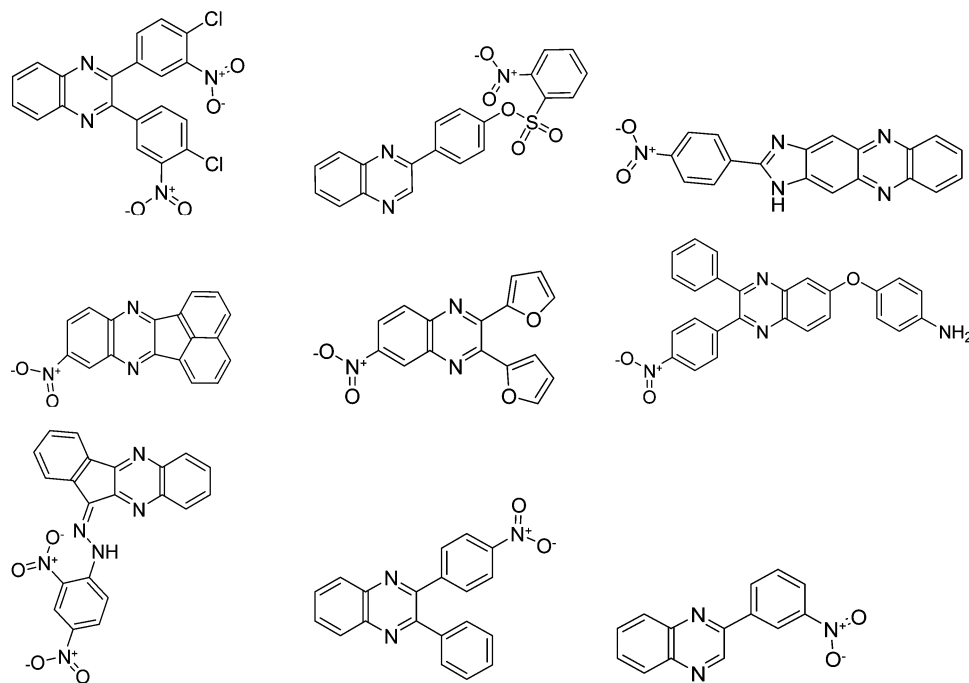
Clustering. Clustering algorithms have been widely used to identify distinct compound classes.^{13–17,32–36,38–40} These algorithms can be used to segregate compounds into distinct structural classes that induce similar biological assay phenotypes as well. Comparisons between clustering algorithms that used Daylight fingerprints^{13–16} as descriptors among others^{58,59} suggest the reported differences in the ability to distinguish structural classes or distinguish active and inactive compounds between any two widely studied clustering algorithms are rarely greater than 10%¹⁴ with no consensus on the ideal algorithm, particularly given conflicting work in genomics associating gene attributes with gene-expression clusters.⁵⁸ Nonetheless, there is a consensus that singletons are essentially uninformative; therefore, while the number of singletons is clearly dataset-dependent, partitionings with high rates of singletons are still not very useful (> 10%).^{14,15}

Using the Daylight fingerprints, each compound in the 16 320-member Chembridge Diverse Set E library was

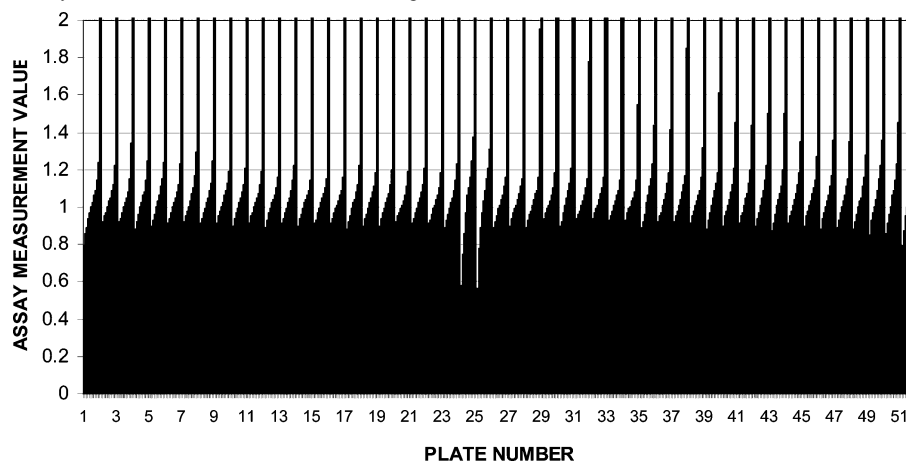
represented by a point in a 4096-dimension chemical space with each dimension corresponding to a single bit in the fingerprint. The compound points were then partitioned by *K*-modal clustering with the Euclidean distance or Tanimoto coefficients as a measure of similarity into *K* = 2, 5, 11, 21, 42, 85, 170, 340, 680, 1360, and 2720 clusters using modified open source software from Cluster 3.0.⁶⁰ Each average cluster size (*N/K*) was approximately a factor of 2 smaller than the previous one, so that the geometry of the chemical space could be probed logarithmically. Each cluster size was generated in triplicate, as *K*-modal clustering has a random component due to random initialization such that any two cluster sets rarely match perfectly. Each cluster represents an approximate structural class of compounds, and these classes may contain various degrees of structural heterogeneity depending on cluster size and the library's diversity. Two sample clusters from the Chembridge library are shown (Charts 1 and 2).

After clustering the 4096-bit Daylight fingerprints into cluster sets with various *K* values, the distribution of cluster sizes among the different cluster sets was examined (see the Supporting Information, Chart A). When the approximate criterion that less than 10% of the clusters be singletons was used,^{14,15} only the three cluster sets with *K* = 2720 derived from the Euclidean distance fell short of that criterion. The number of singletons in the cluster sets derived from Tanimoto coefficients was comparable. It is not surprising that there was a high rate of singletons in the cluster sets with *K* = 2720, since the average cluster size was only six compounds/cluster. Because the Chembridge Diverse Set E compound library represents an ostensibly diverse collection of compounds, the smooth distribution of compound sizes and the absence of any large cluster sizes straying from the distribution were consistent with our expectations.

Class Annotation with Assay Data. Compound classes can be annotated with many different types of assays that have differing numbers of biological targets (most often protein targets). Binding assays measure the binding of only one biological target of interest. Assays that measure very general phenotypes such as growth inhibition, such as NCI's publicly available panel of 70 cancer cell lines, can measure the chemical modulation of a large number of biological

Chart 2. Cluster with Common Substructures^a

^a This cluster demonstrates the case when there is no common structural backbone but there are substructures common to the cluster members. There is less structural uniformity in this cluster, and there may be less uniformity in terms of biological activity of these compounds.

Chart 3. Distribution of Assay Measurements for Each Plate in a Single Screen^a

^a By examining the empirical distribution of plate and row averages, Plates 24 and 25 were found to contain a high percentage of abnormal rows (according to each row's average assay measurement) for this particular assay and were removed from the assay data set. By ordering the assay measurements in increasing order for each plate, the defective plates were also identifiable visually from the above graph.

targets.^{61,62} Phenotypic assays ranging from cell-cycle arrest to lipid uptake have the potential to measure the chemical modulation of large numbers of biological targets while yielding a greater degree of target specificity.^{63,64} The assays considered were a balanced collection of assays including target-based fluorescence polarization assays,⁶⁵ growth-based assays,^{64,66} and phenotypic assays measuring microtubule assembly, actin polymerization, endocytosis, and acetylation among others.^{67–79}

Assays were screened singly or in duplicate using only rough estimates of compound concentration (often in the 10 μ M range). The measurements for these assays included fluorescent or luminescent signals or visually assessed phenotypes from automated microscopic images. Given the large signal variation associated with high-throughput assay

measurements, typically only two or three distinguishable states (hit/nonhit or enhancer/nonhit/suppressor) were considered possible.

Because of the potential for screening defects in high-throughput assays, any plates that had evidence of defects were removed prior to class annotation by examining the empirical distribution of plate rows and plate means and flagging as suspect any row or plate measurement average in the top or bottom 1%. Plates with a large number of suspect rows (3 or more out of 16) also appeared abnormal when their measurement distributions were visually compared to others (Chart 3), and in some cases, the original screeners also identified the plates as defective. Additionally, because many assays employed some form of fluorescent readout, compounds were assayed for fluorescence in the

fluorescein (520 nm) channel and were removed from the analysis in those assays if they were at least 75% more fluorescent than the background.

Each set of structural classes defined by the above clustering algorithms was subsequently annotated with the remaining biological assay data generated by the 48 ICCB high-throughput screens. The high-throughput screening results for each compound were then converted to binary designations (hit or nonhit) using arbitrary thresholds, top 1%, bottom 1%, top 4%, and bottom 4%. Each hit threshold was examined individually in order to determine if certain clusters (i.e., compound classes) contained more hits than others in a given assay.

Chemical Similarity Principle. Before biological assay data can be used to compare compound classes and score each class for the presence of biological activity, the quality of the assay data must first be verified. This can be accomplished by testing the chemical similarity principle, namely, that similar compounds tend to have similar biological activities.^{55,80–82} The ability of high-throughput assays to capture this phenomenon is not guaranteed given assay measurement error, few replicate measurements, and the numerous structural classes potentially active in phenotypic assays, resulting from chemical action on multiple biological targets inducing a single assay phenotype; all these factors could obscure the structural trends of active compounds. Analyses that verify the chemical similarity principle, whether for descriptor validation purposes or the validation of clustering methods, often employ some form of permutation analysis.^{13,81} Permutation analysis entails the comparison of the distribution of assay “hits” in the actual cluster set to the hit distributions in randomly assigned cluster sets using a test statistic such as class membership measures, χ^2 analysis, and configurational entropy.^{83,84}

Both χ^2 analysis and configurational entropy (explained below) statistically evaluate the selective concentration of “hits” in certain classes relative to others. If each cluster is considered a class and the “hits” and “nonhits” are counted in each, then χ^2 analysis and configurational entropy can be directly applied to these data sets and used to measure the selective concentration of hits in certain compound classes relative to others, central to satisfying the chemical similarity principle. The χ^2 analysis and configurational entropy tests must be performed using one assay, one hit threshold, and one cluster set at a time.

To perform the χ^2 test, the measured χ value must be compared to the critical χ value. Counting the number of hits h in each class (as annotated above) for a given assay and calculating the difference from the expected number of hits $E[h]$, the selective enrichment of hits in a given class relative to others can be measured

$$\chi_{\text{measured}} = \sum_{[\text{class}]} (h - E[h])^2 / E[h]$$

where [class] indicates the set of classes obtained from a single clustering routine.

If χ_{observed} is greater than χ_{critical} ($\chi_{\text{observed}} > \chi_{\text{critical}}$), then the chemical similarity principle is satisfied. The χ^2 test measures the nonrandomness with which “hit” compounds were distributed among structure-based compound clusters with an empirical χ_{critical} derived from 95% of 200 randomized data sets ($\alpha = 0.05$). Additionally the χ_{measured} z score

was calculated from the 200 randomized data sets for each combination of assay, threshold, and cluster set.

To corroborate the χ^2 analysis, configurational entropy (based on the multivariate hypergeometric probability distribution) was also employed.^{83,84} For each class of c tested compounds with h hits in a given assay, entropy was calculated:

$$\text{configurations} = C(c, h) = c! / (c - h)! h!$$

$$\text{Entropy}_{\text{observed}} = \sum_{[c]} \log(\text{configurations})$$

Only classes with mixtures of hits and nonhits add entropy to the system. The observed entropy was compared to the empirically derived fifth percentile of 200 random distributions ($\alpha = 0.05$), and z scores were calculated for each assay, threshold, and cluster set combination. In this case, lower than random entropy (with zero entropy indicating a perfectly ordered state) was indicative of the selective concentration of hits in certain compound classes over others.

The above procedure was repeated for each of the 66 cluster sets (11 average cluster sizes generated in triplicate using either the Euclidean distance or Tanimoto coefficients) for the Chembridge Diverse Set E library. For each cluster set, 48 assays were tested using the top 1%, bottom 1%, top 4%, and bottom 4% assay thresholds. By z -score normalizing each χ^2 and entropy test statistic, the observed z scores of different cluster sets can be directly compared (as long as the random distributions used to calculate the z scores were found to be normally distributed by the KS statistic).^{85–87} The KS statistic revealed that configurational entropy was more robust than χ^2 analysis to nonnormal conditions. (See Results and Discussion.)

Scoring Compound Classes. Hypergeometric probability is ideal for scoring the hit density of compound structural classes, and its use has precedence in gene cluster annotation.^{59,88–90} Hypergeometric probability exactly measures the likelihood of observing a given number of hits within a class by chance alone given library size and total hits.

For a class of c tested compounds with h hits given a library with N tested compounds and H total hits in assay a , the probability of getting h or more hits in the cluster is

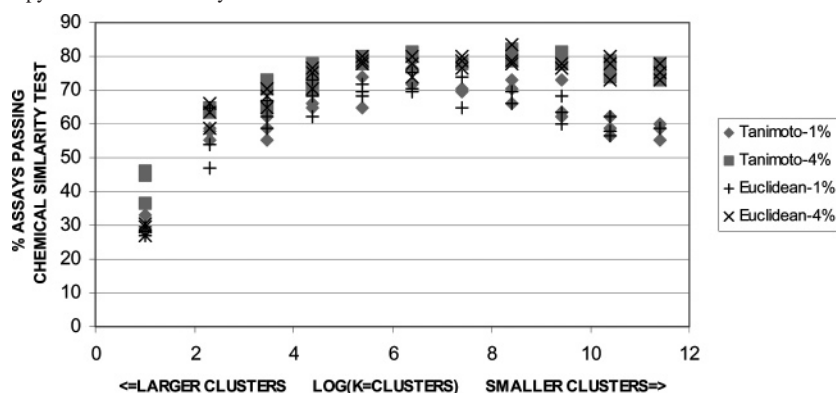
$$P_{a+}(c, h, N, H) = \sum_{i=h \dots \min(c, H)} C(H, i) C(N - H, c - i) / C(N, c)$$

where $C(h, i)$ is the number of configurations as previously defined.

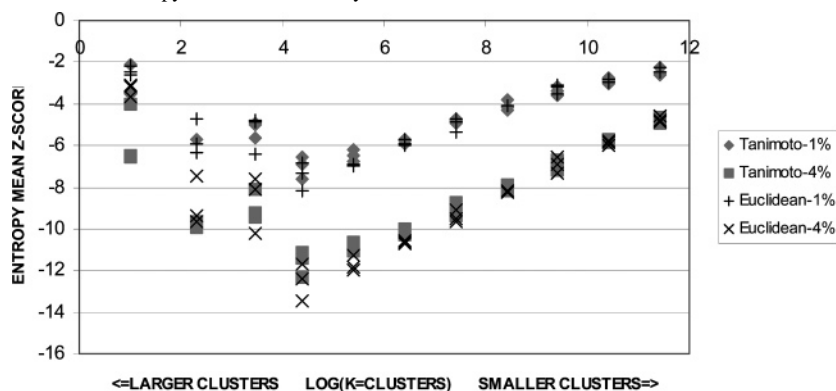
Similarly, to determine the selectivity of compound classes for certain activities, the classes were also scored for the measured *absence* of chemical activity. The probability of getting the observed number of hits or fewer within a particular assay and threshold was calculated for each class using

$$P_{a-}(c, h, N, H) = \sum_{i=\max(0, c-(N-H)) \dots h} C(H, i) C(N - H, c - i) / C(N, c)$$

Using these scores, compound classes could be prioritized for more screening for activity P_{a+} and evaluated further or penalized for lack of activity P_{a-} and never screened again in that assay. The P_{a-} statistic is similar in principle to other statistics that guide follow-up screening on the basis

Chart 4. Passage of the Entropy Chemical Similarity Test^a

^a Each cluster set is represented by a single point using either the 4% or 1% threshold, and the percentage of passing assays was determined. Assays pass the chemical similarity test at a far greater rate than would be expected by chance alone ($\geq 5\%$). Assays tend to pass the chemical similarity test at a greater rate using the top/bottom 4% thresholds rather than the top/bottom 1% thresholds. Also, the cluster sets with large clusters ($\log K < 4$) pass the chemical similarity test less often than the cluster sets with small clusters.

Chart 5. Mean Entropy z Scores of the Entropy Chemical Similarity Test^a

^a The z scores of the entropy values, measuring chemical similarity, indicate that the strongest chemical similarity trends appear around $\log K = 4$ clusters averaging 384 compounds per cluster. According to the KS statistic (not shown), the entropy distributions for $\log K < 4$ become nonnormal, so their z scores cannot be directly compared. Nonetheless, it is clear that global trends in chemical similarity inside the medium–large clusters dominate. Interestingly, every cluster size showed statistical significance ($p < 0.05$), having average z scores $|z| > 2$. Clusters derived from the Euclidean distance appear to slightly outperform clusters derived from Tanimoto coefficients given the greater magnitude of their z scores for this biological assay data set.

of the statistically demonstrated absence of activity in certain compound classes.²² (See Results and Discussion.)

RESULTS AND DISCUSSION

Chemical Similarity Principle. The chemical similarity tests using both χ^2 and configurational entropy as test statistics were passed at a higher rate than expected by chance alone ($E = 5\%$); specifically, 65–80% of threshold and assay combinations passed the chemical similarity test for most cluster sets (Chart 4). Additionally, the chemical similarity test was passed more often using the top/bottom 4% assay thresholds than the top/bottom 1% assay thresholds (Chart 4), and this pattern is corroborated by z scores as expected (Chart 5). The magnitude of the entropy z scores also suggests that clusters derived from the Euclidean distance slightly better capture the structural trends of hit compounds in our data set compared to clusters derived from Tanimoto coefficients (Chart 5). For the smallest clusters, both the percentage of assays passing the chemical similarity test (Chart 4) and the magnitude of the assay chemical similarity z scores decrease (Chart 5), indicating the presence of global trends in biological activity not captured by smaller cluster sizes.

These results suggest that many structural classes of compounds may have *weak, but statistically significant* tendencies to produce “hits” in particular assays (compounds with activities between the first and fourth percentile). This contrasts with the common practice of screeners who only follow up compounds scoring in the top 1% of their assay or higher. Expectedly, assays that failed the chemical similarity test were typically not screened against the entire Chembridge library or were phenotypes assessed visually in a low throughput format covering only a fraction of the library. Only a handful of assays screened against the *entire* library actually failed the chemical similarity test. With these observations taken together, the central premise that the chemical similarity principle would be satisfied by very noisy high-throughput assay data appears verified.

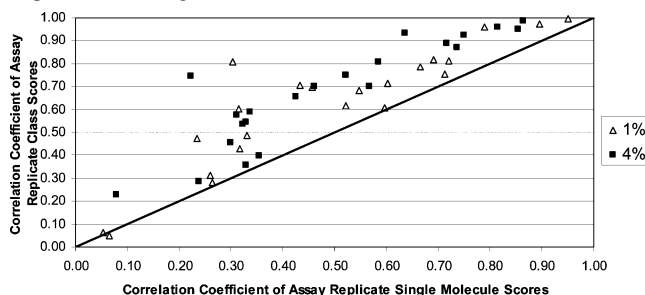
The z scores for each assay’s χ^2 and entropy values are calculated from randomized distributions, which must be normal in order for z scores to be comparable and have meaningful probability values associated with them. The KS statistic indicated that the random distributions used to calculate the χ^2 z scores became nonnormal for over half the cluster sizes tested, whereas the random distributions used to calculate the entropy z scores became nonnormal only for

the largest clusters ($K < 42$). Therefore, the entropy test statistic provided more reliable z scores (necessary to compare different cluster sets to each other) because of the greater prevalence of nonnormal conditions when calculating the χ^2 test statistic. It is important to note that random distributions for cluster sizes with $K < 42$ (regardless of the test statistic used) tend to stray significantly from the normal distribution as indicated by the KS statistic, so z scores calculated from them for $K = 2, 5, 11$, and 21 cannot be directly compared to those for $K = 42, 85, 170, 340, 680, 1360$, and 2720 .

Surprisingly, the optimal cluster size indicated by the maximum assay z scores (valid for $K > 21$) and maximum passage of the chemical similarity test (which is unaffected by nonnormal conditions) appears to be 384 compounds per cluster ($K = 42$). While chemical similarity trends are seen over the entire geometry of the chemical space beginning at $K = 2$ and ending at $K = 2720$, this finding suggests that global trends in chemical similarity, occupying large clusters in chemical space, dominate. It is not the current focus to evaluate such global trends that could be the result of differences between drug-like and nondrug-like chemical subspaces^{91–93} or even so-called privileged structures imbued with activity against multiple biological targets,⁹⁴ but it is a topic worthy of future study. Additionally, while the trends in these larger clusters may be more significant statistically, they often contain a lower percentage of active compounds, so smaller clusters (contained in the three cluster sets, with $K = 1360$ having less heterogeneity while retaining statistical significance) will be analyzed further as individual compound classes are examined and scored for biological assay activity. The analysis will also be further restricted to cluster sets derived from the Euclidean distance given their slight outperformance of cluster sets derived from Tanimoto coefficients using the entropy-based chemical similarity test for this biological assay data set.

Scoring Compound Classes. Scoring compounds as *classes* using the hypergeometric probability proved more reliable than scoring compounds *individually* for any given assay. The ability of $P+$ to distinguish compound classes with various levels of activity far outperformed the original assay measurements, which typically occupy only one of two to three states (hit/nonhit or enhancer/nonhit/suppressor). The probability of activity in a single assay, $P+$ described above, spanned 14 orders of magnitude, well separating clusters for their activity enrichment. In fact, the logarithm of the $P+$ values for each compound (which will be called class annotation scores or class scores) were actually more reproducible than the binary hit designations of individual compounds (which will be called single-molecule scores), as indicated by the respective correlation coefficients between assay duplicates (Chart 6). Points falling above the line indicate that the correlation between the $P+$ class scores of assay duplicates is greater than the correlation of single-molecules scores of those same assay duplicates. Class scores were more correlated than single-molecule scores for every single assay using either the top or bottom 4% threshold and all but one assay using the top or bottom 1% threshold (Chart 6). This suggests that class scoring may mitigate the effects of assay measurement error, allowing classes of strong hits (scoring in best 1%) and weak hits (scoring between the first

Chart 6. Class Scores Are More Reproducible between Assay Replicates than Single-Molecule Scores^a

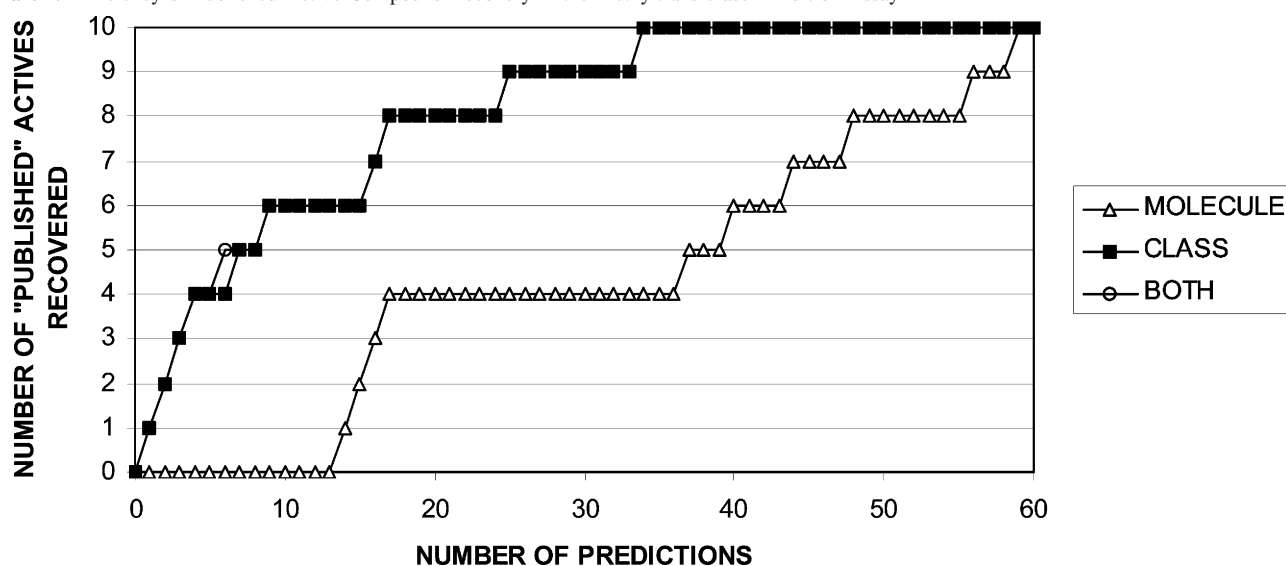


^a For nearly every assay and threshold combination (each represented by a point), the class scores between assay duplicates were more highly correlated than the original binary hit designations (indicated by points falling above the line shown). This argues that class scoring is a more reliable way to evaluate assay data, potentially mitigating the measurement error inherent in noisy phenotypic assay data. These class scores were calculated using one of the cluster sets derived from the Euclidean distance and $K = 1360$.

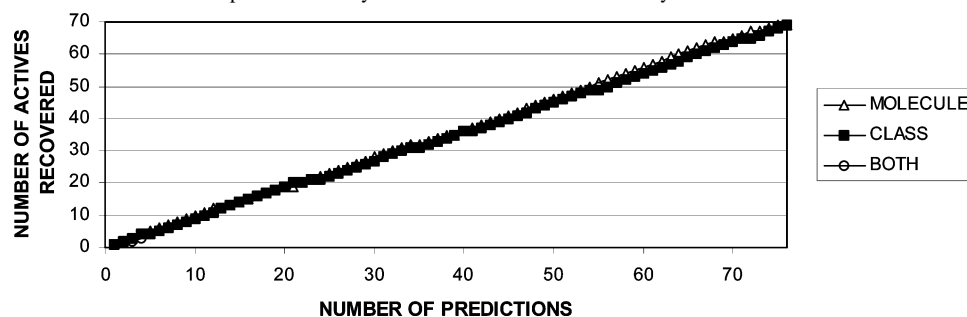
and fourth percentiles) to be pursued more reliably with fewer false positives and false negatives.

Indeed, the false-positive rate appears to be lower for class scores than single-molecule scores given the ability of class scoring to identify known actives with fewer predictions than single-molecule scores. Follow-up data available for two assays, a protein arginine methyltransferase inhibition assay⁷⁸ and a TG-3 mitotic arrest assay,⁶⁴ were examined in order to identify false positives (compounds that did not retest in the assay) and true positives (compounds that retested successfully). This information was then used to compare the two scoring methods. Compounds were ranked separately by class scores (calculated using the 4% activity threshold) and then by single-molecule scores (in this case, the primary assay measurement, not a binary hit designation). The number of known actives recovered by a given scoring method was plotted as a function of the top N scoring compounds predicted to be active. Roughly speaking, the compound scoring method that recovers the known actives using fewer predictions has the lower false positive rate.

Considering, first, the *in vitro* protein arginine methyltransferase inhibition (PRMTI) assay conducted by Bedford and co-workers⁷⁸ (discussed more below), the top scoring 60 compounds, according to the single-molecule score, were selected and ranked according to the single-molecule score and then ranked according to the class score. Within these 60 compounds, 10 were published as PRMT inhibitors, including one compound present twice. The PRMTI assay had a false-positive rate estimated as high as 57% by the original screeners,⁷⁸ indicating that single-molecule scores seemed particularly error prone. Nonetheless, this particular selection of compound test sets favors the single-molecule score since each compound was selected on the basis of its single-molecule score; all compounds not published as active were assumed to be inactive. Despite the compound test set that favored single-molecule scores, class scoring outperformed single-molecule scoring, recovering all 10 active compounds in 34 predictions, whereas single-molecule scoring required 59 to recover all the actives (Chart 7). Combining the scores by ranking compounds first by class scoring and then, secondarily, by their original single-molecule score did not significantly affect the results because most compounds in this assay test set fell into different

Chart 7. Efficiency of Published Active Compound Recovery in the Methyltransferase Inhibition Assay^a

^a Class scoring recovered all 10 active compounds using only 34 predictions of activity, compared to single-molecule scoring, which required almost twice as many predictions to identify the 10 active compounds in the test set. Ranking compounds by both class scoring and then, secondarily, by their original single-molecule score did not significantly affect the results because most compounds in this test set fell into different clusters. This particular assay was very noisy with an estimated false-positive rate of 57%. These class scores were calculated using one of the cluster sets derived from the Euclidean distance and $K = 1360$.

Chart 8. Efficiency of Published Active Compound Recovery in the TG-3 Mitotic Arrest Assay^a

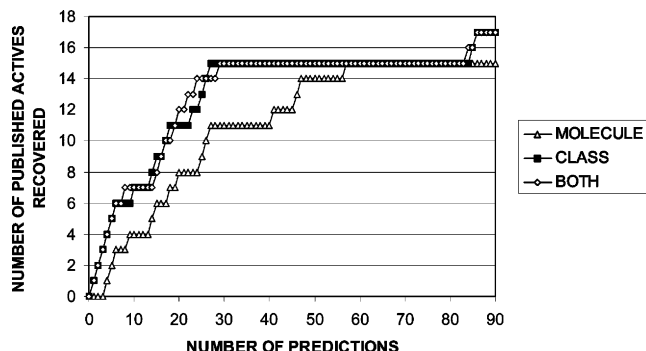
^a Class scoring and single-molecule scoring performed equally well identifying the 69 active compounds in the test set. Because this assay had a very low false-positive rate (estimated at 9%), there was little room for improvement. Ranking compounds by both class scoring and then, secondarily, by their original single-molecule score did not significantly affect the results because most compounds in this test set fell into different clusters, and there were very few reported false positives. These class scores were calculated using one of the cluster sets derived from the Euclidean distance and $K = 1360$.

classes (Chart 7). Given that class scoring recovered active compounds with significantly fewer predictions, it can be directly inferred that class scoring has a lower false-positive rate than single-molecule scoring.

The cell-based TG-3 mitotic arrest assay, conducted by Haggarty and co-workers (discussed more below),⁶⁴ was also considered, and the hit-recovery rates of class scoring and single-molecule scoring were compared. The top 76 scoring compounds according to the single-molecule score were selected and then ranked by single-molecule score, class score, and the combination of both scores. Since only 7 out of the 76 compounds considered (9%) were identified as false positives by follow-up testing,⁶⁴ there was little room for improvement in hit recovery. Given the small number of false positives, class scoring and single-molecule scoring performed equally well at recovering known actives (Chart 8).

So far, the comparisons between class scores and single-molecule scores employ compound test sets that originally favored single-molecule scoring. To offer a "more balanced" comparison, the active compounds recovered by the best class scores were compared to the known actives recovered by

the best single-molecule scores in the TG-3 mitotic arrest assay using a test set of 33 published active compounds selected for their original assay measurement (favoring single-molecule scores) or homology to known actives (favoring class scores).⁶⁴ Class scoring was able to recover known active compounds without assay data available as well as weaker structural homologues that did not score in the original primary assay at all. The combination of class scoring followed by ranking by the single-molecule score did not yield a significant improvement in hit recovery. According to the graph displaying some of the active compounds recovered (Chart 9), a larger share of compounds with the best class scores was published as active with fewer predictions compared to the share of active compounds with the best single-molecule scores. This is the result of a large number of active compounds with or without data available falling into a single class ("Class 2"; discussed below), the inability to assign single-molecule scores to compounds without primary assay data available, and the published activity of less potent compounds that did not score well in the primary assay (but were selected because of their

Chart 9. Efficiency of Recovery of Weak Hits and Untested Compounds in the TG-3 Mitotic Arrest Assay^a

^a Class scoring outperforms single-molecule scoring by recovering known active compounds with fewer predictions. This particular compound test set included compounds that did not score in the original primary assay or did not have primary assay data available and class scoring efficiently recovered them anyway. The combination of both scores did not significantly improve the hit recovery rate. These class scores were calculated using one of the cluster sets derived from the Euclidean distance and $K = 1360$.

homology to known actives). It's worth noting that many compounds predicted to be active by the class-scoring method (assumed to be inactive if they were not published) may actually contain active compounds as well, so this graph likely understates class scoring's capabilities dramatically.

With these observations taken together, class scoring appears to not only score the strongest hits more reliably (with fewer false positives) but it has the ability to recover weaker hits and identify possible activity for compounds for which assay data is not available. To select assay hits for follow-up studies using class scoring, one could choose the most potent compound (according to its original assay measurement) from each of the top scoring compound classes (essentially combining the two scoring methods) in order to get a diverse sampling of likely assay positives. The five cases studies below will provide additional evidence of class scoring's capabilities.

Recovery of Published Hits. Because class scores have fewer false positives than single-molecule scores, less time and resources would be spent testing inactive compounds. Additionally, class scoring should be able to identify structural classes of weaker hits otherwise ignored (as suggested by Chart 9), which should retest at a higher rate as well. Intuitively, selecting weak hits identified by a single assay measurement conveys less evidence of biological activity than an entire class of positively scoring compounds. Some experimental efforts have suggested that compounds belonging to enriched classes with as low as a 30% signal increase over the background may be reproducible (unpublished data), so class scoring may be a better way to pursue and optimize weak but novel structural classes. However, the evidence of weak hit recovery is largely anecdotal, so the most direct evidence of class scoring's capabilities comes from the high rate of identification of published hits as active (strong hits and their weak homologues that have activity confirmed in secondary assays)^{64,70,71,77,78} using part or all of the original high-throughput assay data.

In order for a class to be considered "enriched" in a particular assay activity, the $P_a + (c, h, N, H)$ must be less than or equal to 0.01 for the given assay and threshold combination a . Since three cluster sets (Euclidean E1360.1, E1360.2,

and E1360.3) and two thresholds (1% and 4%) were considered, any single compound had a $p < 2 \times 3 \times 0.01 = 6\%$ chance of being identified as active in a given assay under a conservative independence assumption. The actual probably is less than 6% because the cluster sets and thresholds are partially correlated. Even using the 6% figure, only 4/60 published hits would have been identified as active by chance alone compared to the 34/39 (87%) hits with primary assay data available and 11/20 (55%) hits with unavailable primary assay data successfully identified as active by class scoring. By comparison, using the compounds scoring in the best 6% according to single-molecule scores, 35/39 (90%) of the published compounds were identified as active. Compounds without available assay data have no single-molecule score associated with them, so the corresponding recovery rate is 0/20 (0%). The high recovery rate of 87% by class scoring compares well to the single-molecule scores' 90% recovery rate and directly suggests a similarly low false-negative rate (13% and 10%, respectively). (See Supporting Information, Table B.) When using only one cluster set Euclidean E1360.2 and only the 1% threshold to score the compounds with class scores, 28/39 (72%) of the published active compounds with data available and 10/20 (50%) with data unavailable were identified as active versus an expected $p = 1 \times 1 \times 0.01 = 1\%$ by chance alone. By comparison, using the best-scoring 1% of compounds according to single-molecule scores, 31/39 (79%) of the published compounds with data available were identified as active (and 0%, 0/20, with data unavailable were identified as active). Given that some cluster sets perform better than others, it is better to use multiple cluster sets to score compounds to ensure a higher hit-recovery rate. Nonetheless, considering that the published compounds were originally identified by single-molecule scores, the nearly identical hit-recovery rate (and false-negative rate) of class scoring for this biased test set is noteworthy.

It's encouraging that this method shows a 55% recovery of bioactive compounds with primary assay data unavailable (using all three cluster sets and two thresholds). This was achieved by using screening data from the other parts of the library that were tested and by inferring likely activity from those compounds' membership in otherwise active compound classes. Correspondingly, using leave-one-out elimination, the same class score was calculated for the 39 published hit compounds with data available by removing the published hit from each class (equivalent to these compounds being untested); this method recovered 22/39 (56%) of the published hits, corroborating the 55% recovery figure when the compound data is actually unavailable. This also illustrates the potential that class scoring has to use weaker hits to recover active compound classes, otherwise overlooked when the most potent homologue is not screened.

The ability of class scoring to predict active compounds also seems to outperform the application of Tanimoto coefficients calculated using Daylight fingerprints to single bioactive structures, which correctly predict that a structural homologue will share the activity of a known active compound (in the 10 μM range) only 33–49% of the time.⁵⁵ The structural homologues of the 59 published compounds determined by Daylight Tanimoto coefficients (>0.85) were active in the primary assay 55% of the time; however, over a third of the published compounds had no Tanimoto

Chart 10. Class of Actin Depolymerizers^a

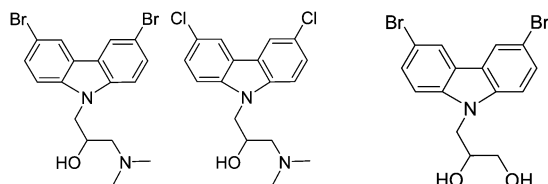
a: Scoring

b: Non-scoring but Published Active

Wiskostatin

Chlorine-substituted

Hydroxyl-Substituted



^a Class scoring recovered the entire class of actin depolymerizers (a, b) as potentially active, although only two of the published class members (a) scored positively in the primary assay. The participation of nonscoring members in this active class correctly suggested their activity at higher concentrations.

homologues at all (far greater than the 7% rate of singletons used in class scoring). This suggests that the clusters used in class scoring are more diverse structurally than Tanimoto homologues, but when scored statistically, they are potentially more informative given their higher rate of hit-recovery of 87%. Class scoring's apparent advantage is also notable because the primary assay data in the ICCB collection is of lower resolution than the IC₅₀s often used in industry.⁵⁵ Naturally, this comparison between class scoring and Tanimoto coefficients is only appropriate when evaluating the assay data of an entire chemical library, because class scoring requires assay data for an entire compound class, whereas the calculation of Tanimoto coefficients of structural homologues only requires a single active compound.

Comparing the binary hit designations to the class scores using the 40 published compounds with *available* screening data, it is not unexpected that a few compounds (5/39) were missed (in part because of structurally heterogeneous compound classes). However, the 87% recovery seems compelling enough to conclude that class scoring performs roughly on par with binary hit designations for the strongest hits. With the added bonus of fewer false positives and the possible recovery of weaker but structurally distinct hits, class scoring may become the preferred (or at least complementary) way to evaluate primary screening data.

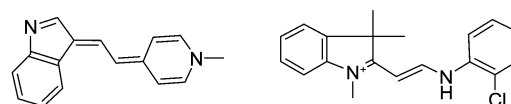
Actin Depolymerizers, Wiskostatin. Rosen and co-workers identified a class of small molecules that inhibit actin polymerization by binding neural Wiskott–Aldrich syndrome protein (*N-WASP*) in an inactive, autoinhibited conformation.⁷⁷ Wiskostatin was originally identified by a high-throughput assay that monitored the polymerization of pyrene-labeled actin monomers mixed with the cytoplasmic extracts of *Xenopus laevis* eggs and PIP-2. Upon the addition of PIP-2, control wells fluoresce as actin polymers assemble; the loss of a fluorescence signal in the presence of the compound indicated the inhibition of actin polymerization.⁷⁷ Four structural homologues were subsequently identified with various degrees of activity in follow-up assays. Only wiskostatin and a chlorine-substituted derivative (Chart 10a) scored in the bottom 1% of the primary assay, suggesting inhibition, and an active hydroxyl-substituted derivative (Chart 10b) scored in the top 1%, suggesting compound autofluorescence at 386 nm. Follow-up assays measured the EC₅₀s of wiskostatin (4.35 μ M) and its chlorine-substituted homologue (8.66 μ M) as well as those of the nonscoring homologues (59.2–242 μ M); the disparities in EC₅₀ values

Chart 11. Class of Mitochondriotoxics^a

a: Scoring

b: Non-scoring

F16



^a Class scoring recovered the mitochondriotoxic molecule F16 (a), which selectively inhibited growth of cells overexpressing *neu*++. A nonscoring homologue (b) is also shown.

potentially explain why the three weaker homologues failed to score in the primary assay conducted at ~ 10 μ M.

Using the class annotation scores, all five compounds were assigned to at least one structural class statistically enriched in actin depolymerization assay activity. All five compounds were assigned to a single active class in cluster sets E1360.1 and E1360.2 for both the bottom 4% and the bottom 1% thresholds. In the third cluster set E1360.3, only wiskostatin, its chlorinated homologue, and its unhalogenated homologue were recovered. This illustrates the importance of evaluating *more than one* cluster set and threshold combination when applying class annotation scores (with the tradeoff of a slightly higher false-positive rate).

Cluster 73 of set E1360.2 had 13/32 compounds scoring in the bottom 4% of the actin assay, including six that scored in the bottom 1%; these measurements suggested the inhibition of actin polymerization with a probability of $P_a + (32, 13, 15 \text{ } 679, 626) = 1.00 \text{ } 10^{-10}$ and $P_a + (32, 6, 16 \text{ } 579, 155) = 6.22 \text{ } 10^{-7}$ for the respective thresholds. This level of statistical significance is substantially smaller than the 0.01 cutoff described above. In addition to the published actin depolymerizers, 11 other structural homologues in this class scored positively in the primary actin assay. One possible interpretation of the other 19 compounds belonging to this class, which did not score in the primary assay, is that they are biologically less potent structural homologues active at higher concentrations; indeed, this is the case for the other three published actin depolymerizers that did not actually score in the primary assay.

Mitochondriotoxic, F16. Leder and co-workers identified a small molecule, F16 (Chart 11), that selectively kills mammary epithelial cells overexpressing the *neu* oncogene as well as related breast cancer cell lines by disrupting the crucial metabolic machinery of mitochondria.⁷¹ F16 was originally identified in a high-throughput assay that compared relative cell growths in the presence of small molecules at 10–15 μ M as measured by BrdU incorporation in two EpH4 mouse mammary epithelial cells, one wild type and the other overexpressing *neu*.⁷¹ F16 was among the most potent compounds that selectively inhibited the growth of *neu*-overexpressing cell lines (*neu*++), operating within the impressively low 0.1 μ M range; F16 acted by accumulating inside the mitochondria of those mutated cell lines and disrupting their metabolism as determined by follow-up studies.⁷¹ According to the primary assay data, having normalized the *neu*++ growth signal relative to the wild type, F16 scored *well* within the best 1% for selective inhibition of *neu*++ growth.

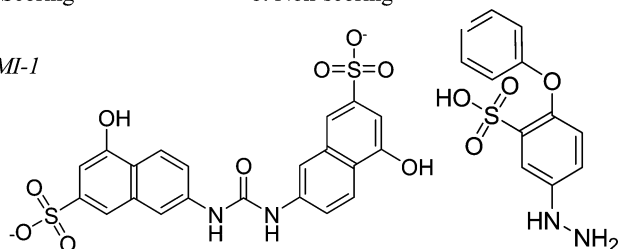
Using the class annotation scores, F16 was correctly assigned to structural classes enriched in activity for the selective inhibition of *neu*++ growth for all three Euclidean

Chart 12. Class of Protein Arginine Methyltransferase Inhibitors^a

a. Scoring

b. Non-scoring

AMI-1



^a Class scoring recovered AMI-1 (a), an inhibitor of protein arginine methyltransferase, by including weaker hits from the original assay in the score. A nonscoring class member (b) is also shown. This class contained more structural heterogeneity than the previous examples.

$K = 1360$ cluster sets and both the 1% and 4% measurement thresholds. Equivalently, to satisfy the $P_a +$ cutoff of 0.01, each of the structural classes containing F16 contained at least one other positively scoring structural homologue.

Cluster 78 of set E1360.2 had 4/9 compounds scoring in the best 4%, including F16 and a homologue that scored in the best 1%; these measurements suggested the selective killing of *neu++* cells relative to the wild type with a probability of $P_a + (9, 4, 15997, 640) = 2.72 \times 10^{-4}$ and $P_a + (9, 2, 15997, 160) = 3.42 \times 10^{-3}$ for the 4% and 1% thresholds, respectively. In addition to F16, three other structural homologues (unpublished) were identified as potentially active according to their primary assay measurements, not to mention the nonscoring homologues that may possess activity at higher concentrations.

Protein Arginine Methyltransferase Inhibitors, AMI-1. Bedford and co-workers identified multiple classes of inhibitors of protein arginine *N*-methyltransferases (PRMTs).⁷⁸ Methylation enhancers and inhibitors (such as AMI-1) were originally identified by a high-throughput assay for chemical modulation of Npl3 methylation by the yeast methyltransferase Hmt1p in the presence of AdoMet, as detected by the 1E4 antibody in vitro.⁷⁸ Subsequent follow-up assays revealed that AMI-1 (Chart 12) selectively inhibits arginine methyltransferases at 3–9 μM ; this activity is consistent with the primary assay, which was conducted at 10–15 μM , given that AMI-1 scored well within the bottom 1%.⁷⁸

Using class annotation, AMI-1 was assigned to a structural class enriched in PRMT inhibition activity for all three Euclidean cluster sets with $K = 1360$. Notably, the activity of these structural classes ($P_a + < 0.01$) was identified using the bottom 4% assay threshold, but not the bottom 1% assay threshold. AMI-1 was the only compound scoring in the bottom 1% of the primary assay for each of the compound classes to which it was assigned. Class scoring recovered AMI-1 as a potential hit only because of the participation of its *less potent* structural homologues in the class score using the 4% threshold. This observation demonstrates the potential equivalency of class scoring using *weak* hits and the scoring of compounds individually: the most potent hits tend to have weaker structural homologues in the library and, thus, can be selected by class scoring using weaker activity thresholds even if the most potent compound is not screened.

Class 45 of cluster set E1360.2 contains six compounds tested in the PRMT assay with four compounds scoring in the bottom 4%, including AMI-1, which also scored in the bottom 1%. The probability of this activity is $P_a + (6, 4, 8318, -$

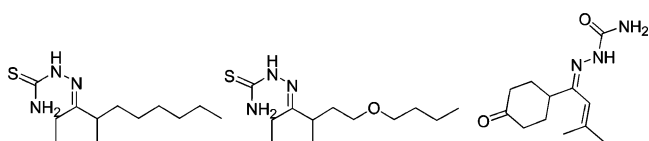
Chart 13. Class of Endocytosis Inhibitors^a

a. Scoring

b. Non-scoring

BLT-1

BLT-2



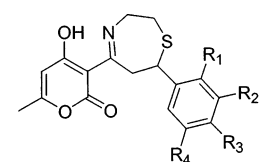
^a Class scoring recovered a class of endocytosis inhibitors BLT-1 and BLT-2 (a). A nonscoring class member is also shown (b). Only the inclusion of the strongest hits (the best scoring 1% of the assay) in the class score recovered this structural class as active.

333) = 3.55×10^{-5} for the bottom 4% threshold and $P_a + (6, 1, 8318, 84) = 0.06$ for the bottom 1%. As noted earlier, the 1% threshold did not pass the arbitrarily selected $P_a + < 0.01$ cutoff, while the 4% threshold did; however, the 1% threshold's probability of 0.06 is still notable, although frequent use of cutoffs this high does run the greater risk of false positives.

Endocytosis Inhibitors, BLT-1 and BLT-2. Kirchhausen and co-workers identified five compounds that inhibit the selective uptake and efflux of lipids by the high-density lipoprotein receptor, scavenger receptor class B type I (SR-BI).⁷⁰ SR-BI inhibitors, such as BLT-1 and BLT-2 (Chart 13), were originally identified by a high-throughput assay that measured the uptake of fluorescently labeled DiI-HDL by compound-treated IdIA-7 cells stably transfected to express high levels of murine SR-BI, relative to untreated cells. The depression of the fluorescent signal in the 5–10 μM compound-treated cells was measured relative to the signal in the same untreated cell line in order to determine if HDL endocytosis was blocked by compound treatment. Follow-up assays revealed that BLT-1 and BLT-2 both inhibit HDL uptake and efflux at the impressively low submicromolar range with IC₅₀s of 0.1 and 0.3 μM , consistent with their primary assay scores within the best 1% and 4%, respectively.

Using the class annotation scores, BLT-1 and BLT-2 were both assigned to the same structural classes enriched in endocytosis-inhibiting activity for cluster sets E1360.1 and E1360.2. BLT-1 and -2 did not score in cluster set E1360.3. For both cluster sets E1360.1 and E1360.2, BLT-1 and -2's probability of activity in the endocytosis assay met the $P_a + < 0.01$ cutoff using the 1% assay threshold but narrowly missed the $P_a + < 0.01$ cutoff using the 4% assay threshold. As opposed to the trend observed for the classes containing AMI-1 with predicted methyltransferase activity, the participation of weaker hits at the 4% threshold did not suggest activity in endocytosis inhibition, whereas the use of only the best 1% of hits did. This example shows the potential for lower activity thresholds such as 4% to underestimate the *potency* of active compounds and the presence of assay activity.

Class 391 of Cluster Set E1360.1 contains 12 compounds tested in the endocytosis assay with 3 compounds scoring with the best 4%, including 2 compounds scoring in the best 1%. BLT-1 and BLT-2 scored in the best 1% and 4%, respectively, as noted above. The probability of this activity is $P_a + (12, 3, 15675, 628) = 0.01074$ for the bottom 4% threshold (narrowly missing the $P_a + < 0.01$ cutoff) and $P_a + (12, 2, 15675, 157) = 0.00616$ for the bottom 1%.

Chart 14. Consensus Structure of Microtubule Destabilizer, "Class 2"^a

- | | |
|---|---|
| 2a R ₁ =OMe, R ₂ =OMe, R ₃ =OMe, R ₄ =H | 2b R ₁ =OMe, R ₂ =H, R ₃ =OMe, R ₄ =H |
| 2c R ₁ =OMe, R ₂ =H, R ₃ =H, R ₄ =OMe | 2d R ₁ =OMe, R ₂ =H, R ₃ =H, R ₄ =H |
| 2e R ₁ =OH, R ₂ =OMe, R ₃ =H, R ₄ =H | 2f R ₁ =H, R ₂ =OCH ₂ O-, R ₃ =OCH ₂ O-, R ₄ =H |
| 2g R ₁ =H, R ₂ =OMe, R ₃ =OMe, R ₄ =OMe | 2h R ₁ =OMe, R ₂ =H, R ₃ =OMe, R ₄ =OMe |
| 2i R ₁ =OMe, R ₂ =OMe, R ₃ =H, R ₄ =H | 2j R ₁ =H, R ₂ =H, R ₃ =OMe, R ₄ =H |
| 2k R ₁ =H, R ₂ =OMe, R ₃ =OH, R ₄ =H | 2l R ₁ =H, R ₂ =OH, R ₃ =H, R ₄ =H |
| 2m R ₁ =OH, R ₂ =Cl, R ₃ =H, R ₄ =Cl | 2n R ₁ =H, R ₂ =H, R ₃ =Cl, R ₄ =H |
| 2o R ₁ =H, R ₂ =Br, R ₃ =H, R ₄ =H | 2p R ₁ =H, R ₂ =H, R ₃ =isopropyl, R ₄ =H |
| 2q R ₁ =H, R ₂ =H, R ₃ =Me, R ₄ =H | |

^a Class scoring recovered this entire class of microtubule destabilizers, which arrest cells in mitosis, although not all compounds scored positively in the primary mitotic arrest assay. The participation of nonscoring members in this active class correctly suggested their activity at higher concentrations.

Microtubule Destabilizers, "Class 2". Schreiber and co-workers identified multiple compound classes that arrest cells in mitosis, one of which contains 17 compounds that act by destabilizing microtubules⁶⁴ (Chart 14). This class was originally identified by a high-throughput "TG-3 cytoblot" assay in A549 lung epithelial cells measuring the chemical modulation (at 20–50 μ M) of the protein nucleolin that is specifically phosphorylated during mitosis. Elevation of the phosphonucleolin signal in the assay suggests a chemically induced mitotic arrest. In addition to inducing mitotic arrest, this class was subsequently classified as microtubule destabilizing in a follow-up assay using purified bovine brain tubulin, supporting the conclusion that this class of compounds targeted tubulin α/β directly, triggering the mitotic spindle checkpoint.⁶⁴ Follow-up assays revealed that the EC₅₀ of compounds 2a–f ranged from 0.5 to 10 μ M and many of their weaker homologues 2g–q only showed similar activity at higher concentrations close to 50 μ M. Primary assay data available for 2a, 2c, 2e, 2h, 2i, 2o, and 2p revealed that each of these compounds scored in the top 1% of the assay with the exception of 2h, which was nonscoring.

Using class annotation scores, the compounds contained in "Class 2" were correctly assigned to classes statistically enriched in TG-3 mitotic arrest activity for all three cluster sets using both top 1% and top 4% assay thresholds. In cluster sets E1360.1 and E1360.3, these compounds were divided between two classes, for which one class in both cluster sets narrowly missed the $P_a+ < 0.01$ cutoff for the top 4% assay threshold but still passed using the top 1% threshold. All the compounds listed in "Class 2" (Chart 14) possess some microtubule destabilizing and mitotic arresting activity at high or low concentrations. Although some of the weaker homologues, 2g–q, did not score positively in the primary assay according to the Haggarty et al. publication, their presence in a compound class enriched in mitotic arrest activity correctly suggested activity at higher concentrations. This observation corroborates the earlier one noting that the structural homologues of wiskostatin, belonging to an active compound class but not scoring in the primary assay, possess the same activity, only at higher concentrations, too.

However, much like the case with BLT-1 and -2, the use of the more stringent top 1% assay threshold recovered all of the active classes, whereas the more liberal top 4% threshold did not. Therefore, paradoxically, the potential to infer activity for the entire class could depend on scoring *only* the most potent compounds in that class statistically. Clearly, there is no general rule for assay threshold choice, since opposite conclusions could be reached by comparing AMI-1, which was *only* recovered by the more liberal 4% threshold, to BLT-1 and -2 and some classes containing members of "Class 2," which were recovered only by the more stringent 1% threshold.

Class 229 of cluster set E1360.2 contains 14 compounds tested in the TG-3 arrest assay including 6 that scored in the top 1%. Fifteen of the "Class 2" microtubule destabilizers listed above are present in this class (structures for 2f and 2j were not identifiable), but assay data is only available for seven of them, so the other eight compounds were considered "untested" and, naturally, did not participate in the class annotation score. The probability of this activity is $P_a+(14,6,8315,328) = 8.267 \cdot 10^{-6}$ for the top 4% threshold and $P_a+(14,6,8315,80) = 1.852 \cdot 10^{-9}$ for the top 1%. These class annotation scores were well-below the cutoff of $P_a+ < 0.01$.

CONCLUSIONS

With the particular collection of small molecules used in this study, class scoring appears to be a more effective way to evaluate phenotypic primary assay data than using single-molecule measurements. This conclusion is bolstered by statistics showing class scores to be more reproducible than binary hit designations, lower false-positive rates using class scoring rather than single-molecule scores, and the validation of class scores' activity predictions using experimental evidence from previous studies emanating from the ICCB screening facility.^{64,70,71,77,78} Class scoring may also be useful in identifying classes of weakly active compounds, otherwise overlooked by screeners, which may prove valuable if those classes contain structural novelties and the ability to be optimized into better lead compounds. Within each class, structural features relevant to biological activity could also be identified by comparing hits and nonhits in order to facilitate lead optimization. Assay data is underutilized when screeners focus only on potency in the primary assay during their follow-up efforts, and class scoring may better capture the full structural diversity of an assay's biologically active compounds. Future studies with additional diverse collections of small molecules and assay data sets will be required to confirm the generality of the conclusions from the current study.

ACKNOWLEDGMENT

We thank Jason McIntosh, Jeremy Muhlich, Carol Chang, Dara Greenhouse, Andrew Lach, Nicola Tolliday, and Caroline Shamu of Harvard Institute of Chemistry and Cell Biology; Frederick Roth of Harvard Medical School; and Tudor Oprea of the University of New Mexico for their invaluable consultations. We thank the National Cancer Institute and their initiative for Chemical Genetics, which supports all screening and informatic (Chembank) activity that served as the basis for the current study. Stuart L. Schreiber is an investigator at the Howard Hughes Medical

Institute located at the Department of Chemistry and Chemical Biology, Harvard University.

Note Added after ASAP Publication. This article was released ASAP on September 1, 2005 with an incorrect notation below the summation sign in two equations. The correct version was posted September 9, 2005.

Supporting Information Available: A chart (Chart A) of the distribution of cluster sizes in cluster sets studied and a table (Table B) of published compounds analyzed and the ability of class annotation to identify them as active. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Patel, D. V.; Gordon, E. M. Applications of small molecule combinatorial chemistry to drug discovery. *Drug Discovery Today* **1996**, 1 (4), 134–144.
- (2) Baum, R. M. Combinatorial approaches provide fresh leads for medicinal chemistry. *Chem. Eng. News* **1994**, 72 (6), 20–26.
- (3) Schreiber, S. L. Target-Oriented and Diversity-Oriented Organic Synthesis in Drug Discovery. *Science* **2000**, 287 (5460), 1964–9.
- (4) Agrafiotis, D. K.; Lobanov, V. S.; Salemme, F. R. Combinatorial informatics in the post-genomics ERA. *Nat. Rev. Drug Discov.* **2002**, 1 (5), 337–46.
- (5) Agrafiotis, D. K.; Myslik, J. C.; Salemme, F. R. Advances in diversity profiling and combinatorial series design. *Mol. Divers.* **1998**, 4 (1), 1–22.
- (6) Schreiber, S. L. The small-molecule approach to biology. *Chem. Eng. News* **2003**, 81 (9), 51–61.
- (7) Agrafiotis, D. K.; Bone, R. F.; Salemme, F. R.; Soll, R. M. System and method for automatically generating chemical compounds with desired properties. U.S. Patent 5,463,564, 1995.
- (8) Graybill, T. L.; Agrafiotis, D. K.; Bone, R.; Illig, C. R.; Jaeger, E. P.; Locke, K. T.; Lu, T.; Salvino, J. M.; Soll, R. M.; Spurlino, J. C.; Subasinghe, N.; Tomczuk, B. E.; Salemme, F. R., *Molecular Diversity and Combinatorial Chemistry: Libraries and Drug Discovery*; American Chemical Society: Washington, DC, 1996; pp 16–27.
- (9) Young, S. S.; Hawkins, D. M. Analysis of a 2⁹ Full Factorial Chemical Library. *J. Med. Chem.* **1995**, 38 (14), 2784–8.
- (10) Rusinko, A., III; Farnen, M. W.; Lambert, C. G.; Brown, P. L.; Young, S. S. Analysis of a Large Structure/Biological Activity Data Set Using Recursive Partitioning. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (6), 1017–26.
- (11) Rusinko, A., III; Young, S. S.; Drewry, D. H.; Gerritz, S. W. Optimization of Focused Chemical Libraries Using Recursive Partitioning. *Comb. Chem. High Throughput Screening* **2002**, 5 (2), 125–33.
- (12) van Rhee, A. M. Use of Recursion Forests in the Sequential Screening Process: Consensus Selection by Multiple Recursion Trees. *J. Chem. Inf. Comput. Sci.* **2003**, 43 (3), 941–8.
- (13) Brown, R. D.; Martin, Y. C. Use of Structure–Activity Data To Compare Structure-Based Clustering Methods and Descriptors for Use in Compound Selection. *J. Chem. Inf. Comput. Sci.* **1996**, 36 (3), 572–584.
- (14) Brown, R. D.; Martin, Y. C. The Information Content of 2D and 3D Structural Descriptors Relevant to Ligand–Receptor Binding. *J. Chem. Inf. Comput. Sci.* **1997**, 37 (1), 1–9.
- (15) Menard, P. R.; Lewis, R. A.; Mason, J. S. Rational Screening Set Design and Compound Selection: Cascaded Clustering. *J. Chem. Inf. Comput. Sci.* **1998**, 38 (3), 497–505.
- (16) Wild, D. J.; Blankley, C. J. Comparison of 2D Fingerprint Types and Hierarchy Level Selection Methods for Structural Grouping Using Ward's Clustering. *J. Chem. Inf. Comput. Sci.* **2000**, 40 (1), 155–62.
- (17) Bajorath, J. Selected Concepts and Investigations in Compound Classification, Molecular Descriptor Analysis, and Virtual Screening. *J. Chem. Inf. Comput. Sci.* **2001**, 41 (2), 233–45.
- (18) Lagunin, A.; Stepanchikova, A.; Filimonov, D.; Poroikov, V. PASS: prediction of activity spectra for biologically active substances. *Bioinformatics* **2000**, 16 (8), 747–8.
- (19) Labute, P. Binary QSAR: a new method for the determination of quantitative structure activity relationships. *Pac. Symp. Biocomput.* **1999**, 444–55.
- (20) Labute, P.; Nilar, S.; Williams, C. A probabilistic approach to high throughput drug discovery. *Comb. Chem. High Throughput Screening* **2002**, 5 (2), 135–45.
- (21) Glick, M.; Klon, A. E.; Acklin, P.; Davies, J. W. Enrichment of Noisy High Throughput Screening Data Using a Naïve Bays Classifier. *J. Biomol. Screening* **2004**, 9 (1), 32–6.
- (22) Schreyer, S. K.; Parker, C. N.; Maggiora, G. M. Data Shaving: A Focused Screening Approach. *J. Chem. Inf. Comput. Sci.* **2004**, 44 (2), 470–9.
- (23) Hansch, C.; Li, R.; Blaney, J. M.; Langridge, R. Comparison of the Inhibition of *Escherichia coli* and *Lactobacillus casei* Dihydrofolate Reductase by 2,4-Diamino-5-(substituted-benzyl)pyrimidines: Quantitative Structure–Activity Relationships, X-ray Crystallography, and Computer Graphics in Structure–Activity Analysis? *J. Med. Chem.* **1982**, 25 (7), 777–84.
- (24) Selassie, C. D.; Li, R. L.; Poe, M.; Hansch, C. On the Optimization of Hydrophobic and Hydrophilic Substituent Interactions of 2,4-Diamino-5-(substituted-benzyl)pyrimidines with Dihydrofolate Reductase. *J. Med. Chem.* **1991**, 34 (1), 46–54.
- (25) Hasegawa, K.; Miyashita, Y.; Funatsu, K. GA Strategy for Variable Selection in QSAR Studies: GA-Based PLS Analysis of Calcium Channel Antagonists. *J. Chem. Inf. Comput. Sci.* **1997**, 37 (2), 306–10.
- (26) Loukas, Y. L. Adaptive Neuro-Fuzzy Inference System: An Instant and Architecture-Free Predictor for Improved QSAR Studies. *J. Med. Chem.* **2001**, 44 (17), 2772–83.
- (27) Viswanadhan, V. N.; Mueller, G. A.; Basak, S. C.; Weinstein, J. N. Comparison of a Neural Net-Based QSAR Algorithm (PCANN) with Hologram- and Multiple Linear Regression-Based QSAR Approaches: Application to 1,4-Dihydropyridine-Based Calcium Channel Antagonists. *J. Chem. Inf. Comput. Sci.* **2001**, 41 (3), 505–11.
- (28) Liu, S. S.; Liu, H. L.; Yin, C. S.; Wang, L. S. VSMP: A Novel Variable Selection and Modeling Method Based on the Prediction. *J. Chem. Inf. Comput. Sci.* **2003**, 43 (3), 964–9.
- (29) Beger, R. D.; Buzatu, D. A.; Wilkes, J. G.; Lay, J. O., Jr. Comparative Structural Connectivity Spectra Analysis (CoSCoSA) Models of Steroid Binding to the Corticosteroid Binding Globulin. *J. Chem. Inf. Comput. Sci.* **2002**, 42 (5), 1123–31.
- (30) Stanton, D. T. Evaluation and Use of BCUT Descriptors in QSAR and QSPR Studies. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (1), 11–20.
- (31) Cramer, R. D.; Poss, M. A.; Hermsmeider, M. A.; Caulfield, T. J.; Kowala, M. C.; Valentine, M. T. Prospective Identification of Biologically Active Structures by Topomer Shape Similarity Searching. *J. Med. Chem.* **1999**, 42 (19), 3919–33.
- (32) Lewis, R. A.; Mason, J. S.; McLay, I. M. Similarity Measures for Rational Set Selection and Analysis of Combinatorial Libraries: The Diverse Property-Derived (DPD) Approach. *J. Chem. Inf. Comput. Sci.* **1997**, 37 (3), 599–614.
- (33) Willett, P.; Winterman, V.; Bawden, D. Implementation of Nonhierarchical Cluster Analysis Methods in Chemical Information Systems: Selection of Compounds for Biological Testing and Clustering of Substructure Search Output. *J. Chem. Inf. Comput. Sci.* **1986**, 26 (3), 109–118.
- (34) Doman, T. N.; Cibulskis, J. M.; Cibulskis, M. J.; McCray, P. D.; Spangler, D. P. Algorithm5: A Technique for Fuzzy Similarity Clustering of Chemical Inventories. *J. Chem. Inf. Comput. Sci.* **1996**, 36 (6), 1195–1204.
- (35) Bayada, D. M.; Hamersma, H.; Geerestein, V. J. Molecular Diversity and Representativity in Chemical Databases. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (1), 1–10.
- (36) Stanton, D. T.; Morris, T. W.; Roychoudhury, S.; Parker, C. N. Application of Nearest-Neighbor and Cluster Analyses in Pharmaceutical Lead Discovery. *J. Chem. Inf. Comput. Sci.* **1999**, 39 (1), 21–27.
- (37) Xiao, Z.; Xiao, Y. D.; Feng, J.; Golbraikh, A.; Tropsha, A.; Lee, K. H. Modeling of Epipodophyllotoxin Derivatives Using Variable Selection *k* Nearest Neighbor QSAR Method. *J. Med. Chem.* **2002**, 45 (11), 2294–309.
- (38) Pirard, B.; Pickett, S. D. Classification of Kinase Inhibitors Using BCUT Descriptors. *J. Chem. Inf. Comput. Sci.* **2000**, 40 (6), 1431–40.
- (39) MacCuish, J.; Nicolaou, C.; MacCuish, N. E. Ties in Proximity and Clustering Compounds. *J. Chem. Inf. Comput. Sci.* **2001**, 41 (1), 134–46.
- (40) Espinosa, G.; Arenas, A.; Giralt, F. An Integrated SOM-Fuzzy ARTMAP Neural System for the Evaluation of Toxicity. *J. Chem. Inf. Comput. Sci.* **2002**, 42 (2), 343–59.
- (41) Zheng, W.; Tropsha, A. Novel Variable Selection Quantitative Structure–Property Relationship Approach Based on the *k*-Nearest-Neighbor Principle. *J. Chem. Inf. Comput. Sci.* **2000**, 40 (1), 185–94.
- (42) Cramer, R. D. Topomer CoMFA: A Design Methodology for Rapid Lead Optimization. *J. Med. Chem.* **2003**, 46 (3), 374–88.
- (43) Cramer, R. D., III; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis CoMFA. 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, 110 (18), 5959–5967.
- (44) So, S. S.; Karplus, M. Three-Dimensional Quantitative Structure–Activity Relationships from Molecular Similarity Matrices and Genetic Neural Networks. 1. Method and Validations. *J. Med. Chem.* **1997**, 40 (26), 4347–59.
- (45) Hopfinger, A. J.; Wang, S.; Tokarski, J. S.; Jin, B.; Albuquerque, M.; Madhav, P. J.; Duraiswami, C. Construction of 3D-QSAR Models

- Using the 4D-QSAR Analysis Formalism. *J. Am. Chem. Soc.* **1997**, *119* (43), 10509–10524.
- (46) Vedani, A.; Dobler, M. 5D-QSAR: The Key for Simulating Induced Fit? *J. Med. Chem.* **2002**, *45* (11), 2139–49.
- (47) Vedani, A.; Briem, H.; Dobler, M.; Dollinger, H.; McMasters, D. R. Multiple-Conformation and Protonation-State Representation in 4D-QSAR: The Neurokinin-1 Receptor System. *J. Med. Chem.* **2000**, *43* (23), 4416–27.
- (48) Vedani, A.; Dobler, M.; Zbinden, P. Quasi-Atomistic Receptor Surface Models: A Bridge between 3-D QSAR and Receptor Modeling. *J. Am. Chem. Soc.* **1998**, *120* (18), 4471–4477.
- (49) Kuntz, I. D. Structure-based strategies for drug design and discovery. *Science* **1992**, *257* (5073), 1078–82.
- (50) Taylor, R. D.; Jewsbury, P. J.; Essex, J. W. A review of protein-small molecule docking methods. *J. Comput.-Aided Mol. Des.* **2002**, *16* (3), 151–66.
- (51) Rowland, R. S. Using X-ray crystallography in drug discovery. *Curr. Opin. Drug Discovery Dev.* **2002**, *5* (4), 613–9.
- (52) Shimada, J.; Ishchenko, A. V.; Shakhnovich, E. I. Analysis of knowledge-based protein–ligand potentials using a self-consistent method. *Protein Sci.* **2000**, *9* (4), 765–75.
- (53) Ishchenko, A. V.; Shakhnovich, E. I. Small Molecule Growth 2001 (SMoG2001): An Improved Knowledge-Based Scoring Function for Protein–Ligand Interactions. *J. Med. Chem.* **2002**, *45* (13), 2770–80.
- (54) Stockfisch, T. P. Partially Unified Multiple Property Recursive Partitioning (PUMP-RP): A New Method for Predicting and Understanding Drug Selectivity. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (5), 1608–13.
- (55) Martin, Y. C.; Kofron, J. L.; Traphagen, L. M. Do Structurally Similar Molecules Have Similar Biological Activity? *J. Med. Chem.* **2002**, *45* (19), 4350–8.
- (56) Matter, H.; Pötter, T. Comparing 3D pharmacophore triplets and 2D fingerprints for selecting diverse compound subsets. *J. Chem. Inf. Comput. Sci.* **1999**, *39* (6), 1211–1225.
- (57) Tubert, I. *Chemistry–Smiles–0.13: Smile Parser*; CPAN: www.cpan.org, 2003.
- (58) Gibbons, F. D.; Roth, F. P. Judging the Quality of Gene Expression-Based Clustering Methods Using Gene Annotation. *Genome Res.* **2002**, *12* (10), 1574–81.
- (59) King, O. D.; Foulger, R. E.; Dwight, S. S.; White, J. V.; Roth, F. P. Predicting Gene Function From Patterns of Annotation. *Genome Res.* **2003**, *13* (5), 896–904.
- (60) Hoon, M. D.; Imoto, S.; Miyano, S. *The C Clustering Library*; The University of Tokyo, Institute of Medical Science, Human Genome Center: Tokyo, Japan, 2003.
- (61) Holbeck, S. L. Update on NCI in vitro drug screen utilities. *Eur. J. Cancer* **2004**, *40* (6), 785–93.
- (62) Monks, A.; Scudiero, D. A.; Shoemaker, R. H.; Paull, K. D.; Vistica, D.; Hose, C.; Langley, J.; Cronise, P.; Vaigro-Wolff, A.; Gray-Goodrich, M.; Campell, H.; Mayo, J.; Boyd, M. R. Feasibility of a high-flux anticancer screen using a diverse panel of cultured human tumor lines. *J. Natl. Cancer Inst.* **1991**, *83* (11), 757–766.
- (63) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Adv. Drug Delivery Rev.* **2001**, *46* (1–3), 3–26.
- (64) Haggarty, S. J.; Mayer, T. U.; Miyamoto, D. T.; Fathi, R.; King, R. W.; Mitchison, T. J.; Schreiber, S. L. Dissecting cellular processes using small molecules: identification of colchicine-like, taxol-like and other small molecules that perturb mitosis. *Chem. Biol.* **2000**, *7* (4), 275–86.
- (65) Degterev, A.; Lugovskoy, A.; Cardone, M.; Mulley, B.; Wagner, G.; Mitchison, T.; Yuan, J. Identification of small-molecule inhibitors of interaction between BH3 domain and Bcl-xL. *Nat. Cell Biol.* **2001**, *3* (2), 173–82.
- (66) Koeller, K. M.; Haggarty, S. J.; Perkins, B. D.; Leykin, I.; Wong, J. C.; Kao, M. C.; Schreiber, S. L. Chemical Genetic Modifier Screens: Small Molecule Trichostatin Suppressors as Probes of Intracellular Histone and Tubulin Acetylation. *Chem. Biol.* **2003**, *10* (5), 397–410.
- (67) Mayer, T. U.; Kapoor, T. M.; Haggarty, S. J.; King, R. W.; Schreiber, S. L.; Mitchison, T. J. Small molecule inhibitor of mitotic spindle bipolarity identified in a phenotype-based screen. *Science* **1999**, *286* (5441), 971–4.
- (68) Stockwell, B. R.; Haggarty, S. J.; Schreiber, S. L. High-throughput screening of small molecules in miniaturized mammalian cell-based assays involving post-translational modifications. *Chem. Biol.* **1999**, *6* (2), 71–83.
- (69) Kim, T.; Kim, T. Y.; Lee, W. G.; Yim, J.; Kim, T. K. Signaling Pathways to the Assembly of an Interferon- β Enhanceosome. *J. Biol. Chem.* **2000**, *275* (22), 16910–7.
- (70) Nieland, T. J.; Penman, M.; Dori, L.; Krieger, M.; Kirchhausen, T. Discovery of chemical inhibitors of the selective transfer of lipids mediated by the HDL receptor SR–BI. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99* (24), 15422–7.
- (71) Fantin, V. R.; Berardi, M. J.; Scorrano, L.; Korsmeyer, S. J.; Leder, P. A Novel Mitochondriotoxic Small Molecule That Selectively Inhibits Tumor Cell Growth. *Cancer Cell* **2002**, *2* (1), 29–42.
- (72) Kao, R. Y.; Jenkins, J. L.; Olson, K. A.; Key, M. E.; Fett, J. W.; Shapiro, R. A small-molecule inhibitor of the ribonucleolytic activity of human angiogenin that possesses antitumor activity. *Proc. Natl. Acad. Sci. U.S.A.* **2002**, *99* (15), 10066–71.
- (73) Yarrow, J. C.; Feng, Y.; Perlman, Z. E.; Kirchhausen, T.; Mitchison, T. J. Phenotypic screening of small molecule libraries by high throughput cell imaging. *Comb. Chem. High Throughput Screening* **2003**, *6* (4), 279–86.
- (74) Feng, Y.; Yu, S.; Lasell, T. K.; Jadhav, A. P.; Macia, E.; Chardin, P.; Melancon, P.; Roth, M.; Mitchison, T.; Kirchhausen, T. Exo1: a new chemical inhibitor of the exocytic pathway. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100* (11), 6469–74.
- (75) Haggarty, S. J.; Koeller, K. M.; Wong, J. C.; Grozinger, C. M.; Schreiber, S. L. Domain-selective small-molecule inhibitor of histone deacetylase 6 (HDAC6)-mediated tubulin deacetylation. *Proc. Natl. Acad. Sci. U.S.A.* **2003**, *100* (8), 4389–94.
- (76) Straight, A. F.; Cheung, A.; Limouze, J.; Chen, I.; Westwood, N. J.; Sellers, J. R.; Mitchison, T. J. Dissecting Temporal and Spatial Control of Cytokinesis with a myosin II Inhibitor. *Science* **2003**, *299* (5613), 1743–7.
- (77) Peterson, J. R.; Bickford, L. C.; Morgan, D.; Kim, A. S.; Querfelli, O.; Kirschner, M. W.; Rosen, M. K. A chemical inhibitor of N-WASP reveals a new mechanism for targeting protein interactions. *Nat. Struct. Mol. Biol.* **2004**, *11* (8), 747–55.
- (78) Cheng, D.; Yadav, N.; King, R. W.; Swanson, M. S.; Weinstein, E. J.; Bedford, M. T. Small molecule regulators of protein arginine methyltransferases. *J. Biol. Chem.* **2004**, *279* (23), 23892–9.
- (79) Strausberg, R. L.; Schreiber, S. L. From knowing to controlling: a path from genomics to drugs using small molecule probes. *Science* **2003**, *300* (5617), 294–5.
- (80) *Concepts and Applications of Molecular Similarity*; Johnson, M. A., Maggiora, G. M., Eds.; John Wiley & Sons: New York, 1990.
- (81) Kearsley, S. K.; Sallamack, S.; Fluder, E. M.; Andose, J. D.; Mosley, R. T.; Sheridan, R. P. Chemical Similarity Using Physicochemical Property Descriptors. *J. Chem. Inf. Comput. Sci.* **1996**, *36* (1), 118–127.
- (82) Flower, D. R. On the Properties of Bit String-Based Measures of Chemical Similarity. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (3), 379–386.
- (83) Buck, K. K.; Subramanian, V.; Block, D. E. Identification of Critical Batch Operating Parameters in Fed-Batch Recombinant *E. coli* Fermentations Using Decision Tree Analysis. *Biotechnol. Prog.* **2002**, *18* (6), 1366–76.
- (84) Agresti, A. *Categorical Data Analysis*; John Wiley & Sons: New York, 1990.
- (85) Lehmann, K. G.; Melkert, R.; Serruys, P. W. Contributions of frequency distribution analysis to the understanding of coronary restenosis. A reappraisal of the Gaussian curve. *Circulation* **1996**, *93* (6), 1123–32.
- (86) Bohning, D.; Hempfling, A.; Schelp, F. P.; Schlattmann, P. The area between curves (ABC)—measure in nutritional anthropometry. *Stat. Med.* **1992**, *11* (10), 1289–304.
- (87) Weiss, M. S. Testing correlated “EEG-like” data for normality using a modified Kolmogorov-Smirnov statistic. *IEEE Trans. Biomed. Eng.* **1986**, *33* (12), 1114–20.
- (88) Jakt, L. M.; Cao, L.; Cheah, K. S.; Smith, D. K. Assessing Clusters and Motifs from Gene Expression Data. *Genome Res.* **2001**, *11* (1), 112–23.
- (89) Berriz, G. F.; King, O. D.; Bryant, B.; Sander, C.; Roth, F. P. Characterizing gene sets with FuncAssociate. *Bioinformatics* **2003**, *19* (18), 2502–4.
- (90) Tavazoie, S.; Hughes, J. D.; Campbell, M. J.; Cho, R. J.; Church, G. M. Systematic determination of genetic network architecture. *Nat. Genet.* **1999**, *22* (3), 281–5.
- (91) Oprea, T. I.; Gottfries, J. Chemography: The Art of Navigating in Chemical Space. *J. Comb. Chem.* **2001**, *3* (2), 157–66.
- (92) Oprea, T. I.; Zamora, I.; Ungell, A. L. Pharmacokinetically Based Mapping Device for Chemical Space Navigation. *J. Comb. Chem.* **2002**, *4* (4), 258–66.
- (93) Oprea, T. I. Current trends in lead discovery: Are we looking for the appropriate properties? *J. Comput.-Aided Mol. Des.* **2002**, *16* (5–6), 325–34.
- (94) Horton, D. A.; Bourne, G. T.; Smythe, M. L. The Combinatorial Synthesis of Bicyclic Privileged Structures or Privileged Substructures. *Chem. Rev.* **2003**, *103* (3), 893–930.