

ProSAR: A New Methodology for Combinatorial Library Design

Hongming Chen,^{*,†} Ulf Börjesson,[†] Ola Engkvist,[†] Thierry Kogej,[†] Mats A. Svensson,[‡]
Niklas Blomberg,[†] Dirk Weigelt,[‡] Jeremy N. Burrows,^{£,§} and Tim Lange^{‡,§}

DECS GCS Computational Chemistry, AstraZeneca R&D Mölndal, Pepparedsleden 1, SE-43183 Mölndal, Sweden, and Medicinal Chemistry, AstraZeneca R&D Södertälje, SE-151 85 Södertälje, Sweden

Received July 9, 2008

A method is introduced for performing reagent selection for chemical library design based on topological (2D) pharmacophore fingerprints. Optimal reagent selection is achieved by optimizing the Shannon entropy of the 2D pharmacophore distribution for the reagent set. The method, termed ProSAR, is therefore expected to enumerate compounds that could serve as a good starting point for deriving a structure activity relationship (SAR) in combinatorial library design. This methodology is exemplified by library design examples where the active compounds were already known. The results show that most of the pharmacophores on the substituents for the active compounds are covered by the designed library. This strategy is further expanded to include product property profiles for aqueous solubility, hERG risk assessment, etc. in the optimization process so that the reagent pharmacophore diversity and the product property profile are optimized simultaneously via a genetic algorithm. This strategy is applied to a two-dimensional library design example and compared with libraries designed by a diversity based strategy which minimizes the average ensemble Tanimoto similarity. Our results show that by using the PSAR methodology, libraries can be designed with simultaneously good pharmacophore coverage and product property profile.

INTRODUCTION

Combinatorial and parallel synthesis technologies are powerful tools in early drug discovery. Once active compounds are identified, combinatorial or parallel synthesis libraries are a natural choice for expanding hits into potential lead series with hundreds of compounds synthesized to explore the chemical space around the identified hits. Design and synthesis of combinatorial libraries can also be used as a complement to existing screening libraries, for example a library designed for a specific protein family. The art and science of computational library design has been reviewed extensively.^{1–3} Chemical diversity^{4–6} is often used as an optimization function for library design, either on the reagent side^{7,8} or on the product side.^{9,10} Such library design strategies are often very efficient at selecting diverse compounds, but one drawback is that when the designed libraries are tested in assays, sometimes it is hard to derive a clear structure activity relationship (SAR) from the experimental results since the selected building blocks might have little or no relationship to one another.

Here we describe a new methodology to address this issue and present a library design strategy to obtain libraries for which it should be easier to derive an SAR after in vitro screening. In this method, all the reagents that are used in the library synthesis are first examined and encoded into a two-dimensional (2D) pharmacophore fingerprint by considering the pharmacophoric types and the bond distance

between each pharmacophore element and the attachment point of the reagent to the scaffold. The pharmacophore variability in the reagent set is then estimated by its Shannon entropy.¹¹ Shannon entropy was originally applied in digital communication technology to determine the amount of data that could be transmitted within a given range of frequencies. Bajorath et al.¹² have used Shannon entropy for analyzing variability of molecular descriptors in compound databases. Grootenhuys et al.^{13,14} and Miller et al.¹⁵ used the Shannon entropy of the pharmacophore distribution of the products to estimate the chemical diversity of a library. By optimizing the product Shannon entropy, a library which covered most of the pharmacophore space for the source pool could be obtained. McGregor et al.^{7,16} applied pharmacophore fingerprints in library design and proposed a strategy for maximizing the overlap between the designed library and the MDDR bioactive space with respect to the first three principal components of the pharmacophore fingerprint. Our method can be regarded as a combination of these two strategies, but with the fundamental difference that we focus on reagent pharmacophore space, whereas their work chose to optimize the coverage of product pharmacophore space. Furthermore, we generate 2-point reagent pharmacophore fingerprints where one end point is always the attachment point on the scaffold, so that the pharmacophore variability in the fingerprint is always relative to the same attachment point and can therefore be compared within the same framework to derive a SAR. Compared to libraries derived from product diversity, our strategy first filters out high complexity reagents in terms of number of functional groups and bond length to the attachment point and then selects side chains with systematically varied pharmacophore elements, helping to derive some kind of SAR for the designed library

* Corresponding author e-mail: hongming.chen@astrazeneca.com.

[†] AstraZeneca R&D Mölndal.

[‡] AstraZeneca R&D Södertälje.

[£] Present address: European Patent Office, PatentLaan 2, 2288 EE Rijswijk, The Netherlands.

[§] Present address: Medicines for Malaria Venture, Route de Pré-Bois 20, 1215 Geneva, Switzerland.

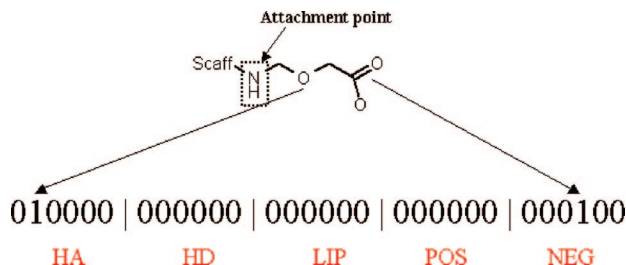


Figure 1. Encoding of a reagent pharmacophore fingerprint. The fingerprint for each reagent consists of 30 bins, corresponding to five different pharmacophore types: hydrogen bond acceptor (HA), hydrogen bond donor (HD), lipophilic group (LIP), positively charged group (POS), and negatively charged group (NEG). Each pharmacophore type has 6 bins, and the order of bins refers to the number of bonds between the pharmacophore element and the attachment point to the scaffold. In this illustrated structure, the N atom is the attachment point to scaffold, the ether oxygen atom has two bonds to the N atom, and it belongs to the HA type. So the second bin is set as 1. The carboxylic acid (NEG) is 4 bonds from the attachment point, and the 28th bin is set as 1. Features further than 6 bonds away from the attachment point are not taken into account.

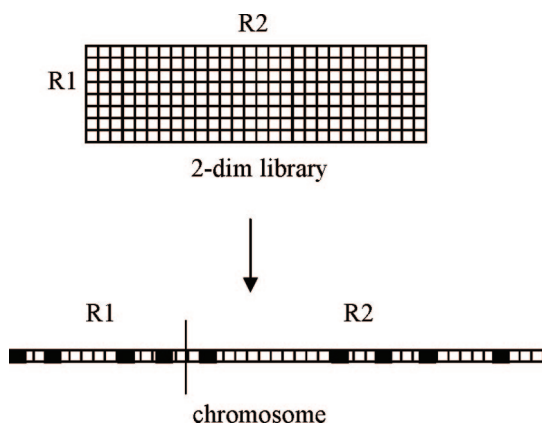


Figure 2. Encoding of a chromosome in the genetic algorithm. Each chromosome corresponds to a selected library, and it consists of an array of binary bins. Each bin refers to a reagent. Once a certain reagent is selected in the library, its bin is set as “1”, otherwise as “0”. The length of the chromosome is equal to the sum of all reagents.

as activity data become available. We have chosen to call this type of library a “ProSAR” (PSAR) library. We will also show how the PSAR strategy can be combined with product properties (for instance aqueous solubility, hERG liability, etc.) to design a library which not only will help to derive a SAR but also has an attractive property profile.

METHODS

Encoding of the Pharmacophore Fingerprint. A two-point pharmacophore is designed to encode the reagent pharmacophore fingerprint. The fingerprint consists of a single pharmacophore element and the attachment atom of the reagent as shown in Figure 1. Five standard pharmacophore types are used: hydrogen bond donor (HD), hydrogen bond acceptor (HA), positive charge center (POS), negative charge center (NEG), and lipophilic groups (LIP). Pharmacophores are defined via a set of SMARTS¹⁷ patterns. The bond distance between the pharmacophore element and the attachment atom is here limited to 6 bonds as we want

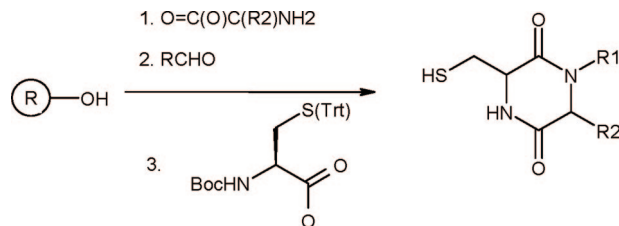


Figure 3. Combinatorial library example from Affymax.³²

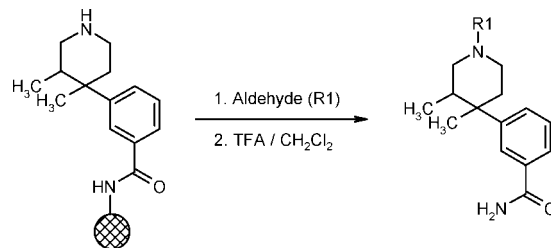


Figure 4. Combinatorial library example from Adolor.³⁵

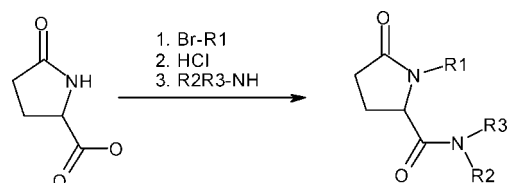


Figure 5. Library example taking into account both reagent pharmacophore entropy and library property profile as described in the text.

to keep the reagent complexity low and avoid adding long side chains to the scaffold. The total number of HD, HA, POS, and NEG functional groups on a reagent was restricted to no more than 2 to further reduce the complexity on pharmacophore elements. The total number of unique two-point pharmacophores in a reagent is therefore 30 (5×6), and a 30-bin pharmacophore fingerprint can therefore be constructed, in which each bin refers to a specific two-point pharmacophore. The value in each bin is the frequency of the specific pharmacophore in the reagent corresponding to that bin. Figure 1 shows an example of such a pharmacophore fingerprint for an amine reagent.

Optimization in Pharmacophore Space for “ProSAR”. Our library design strategy is to select reagents that cover the pharmacophore space optimally, while keeping the pharmacophore distribution as even as possible. Such a distribution can be effectively characterized by the Shannon entropy (SE), and it was therefore used to find the optimal reagent subset based on the “pharmacophore fingerprint space”. The SE for every reagent set is defined as

$$SE = - \sum_i p_i \log_2 p_i \quad (1)$$

where p_i is the probability of having a certain pharmacophore in the whole reagent set. p_i is calculated as

$$p_i = c_i / \sum_i c_i \quad (2)$$

where c_i is the population of pharmacophore i in the whole reagent set. A greater Shannon entropy value means that the pharmacophores for the selected reagent subset is more evenly distributed over the 30 bins. Hence, during the course of optimization, a set of reagent compounds is sought to maximize the Shannon entropy.

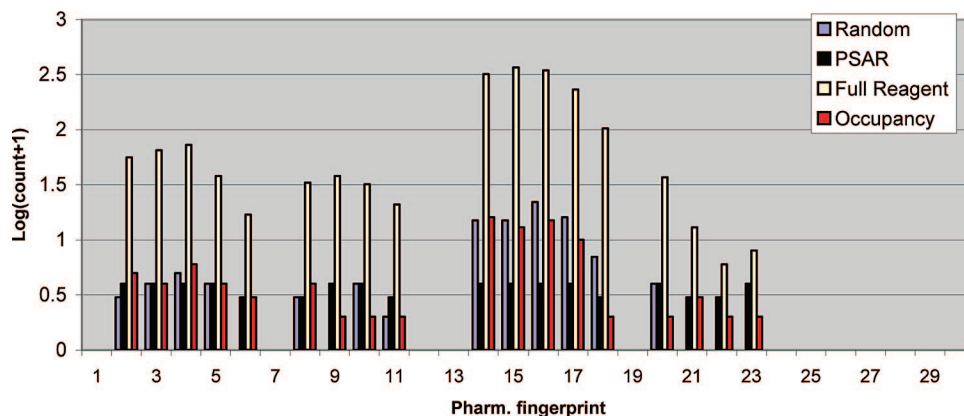
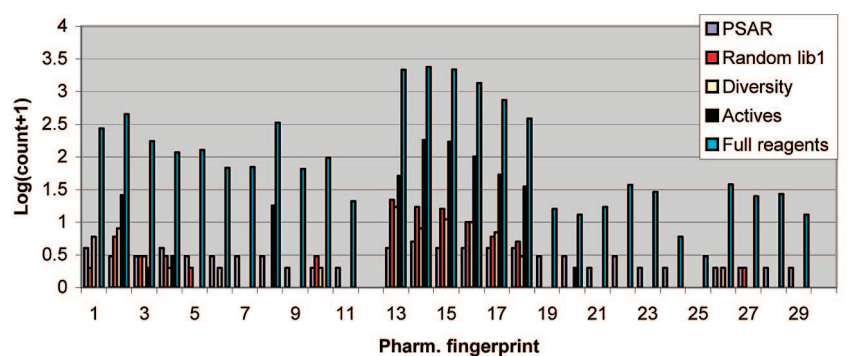
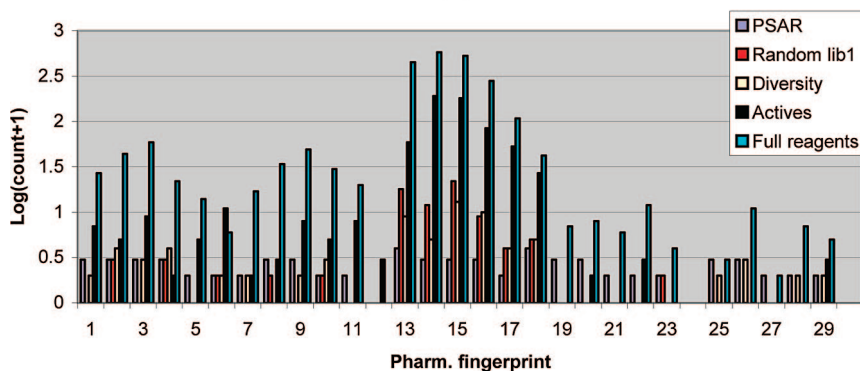


Figure 6. Pharmacophore fingerprint distribution for 20 carboxylic acids selected by using the “ProSAR” strategy, random selection (Shannon entropy value is 3.17), and fingerprint bin occupancy optimization. The X axis refers to the pharmacophore bins, and the Y axis refers to value of $\log(\text{count}+1)$; here the “count” corresponds to the frequency of each bin.



(a)



(b)

Figure 7. (a) Pharmacophore fingerprint distribution for the R1 reagents and (b) pharmacophore fingerprint distribution for the R2 reagents.

A simple “greedy” search algorithm (Miller et al.¹⁵) was chosen as the optimization engine to search for an optimal reagent subset. This was done by a first “greedy”-build up of the subset until the desired number of compounds are selected, followed by a second phase that re-evaluates each of the selected compounds in the subset to see if a better choice is available. The second stage continues until no improvement in the subset is possible.

Optimization of the Pharmacophore Entropy and the Library Property Profile. Physico-chemical properties are an important aspect to consider when a library is designed. For example, it is preferable that compounds with poor solubility or high hERG risk are not synthesized. Prediction of solubility and hERG binding have been intensively studied,^{18–20} and many articles^{21,22} consider these

properties in the library design strategy. A more realistic library design strategy would therefore be to extend the “ProSAR” concept further so as to include the library property profile. In order to calculate the properties of a library, a full enumeration has been done, and the properties are calculated at the product level. At AstraZeneca we have established a set of stringent property criteria for checking our compound collection enhancement libraries.²³ Here in this study, the calculated properties include compound novelty (compared with in-house/external compounds to check if the compound is novel), aqueous solubility,²⁴ predicted hERG liability,²⁵ and an in-house lead profile score.^{26,27} All properties were calculated by our in-house prediction tools.

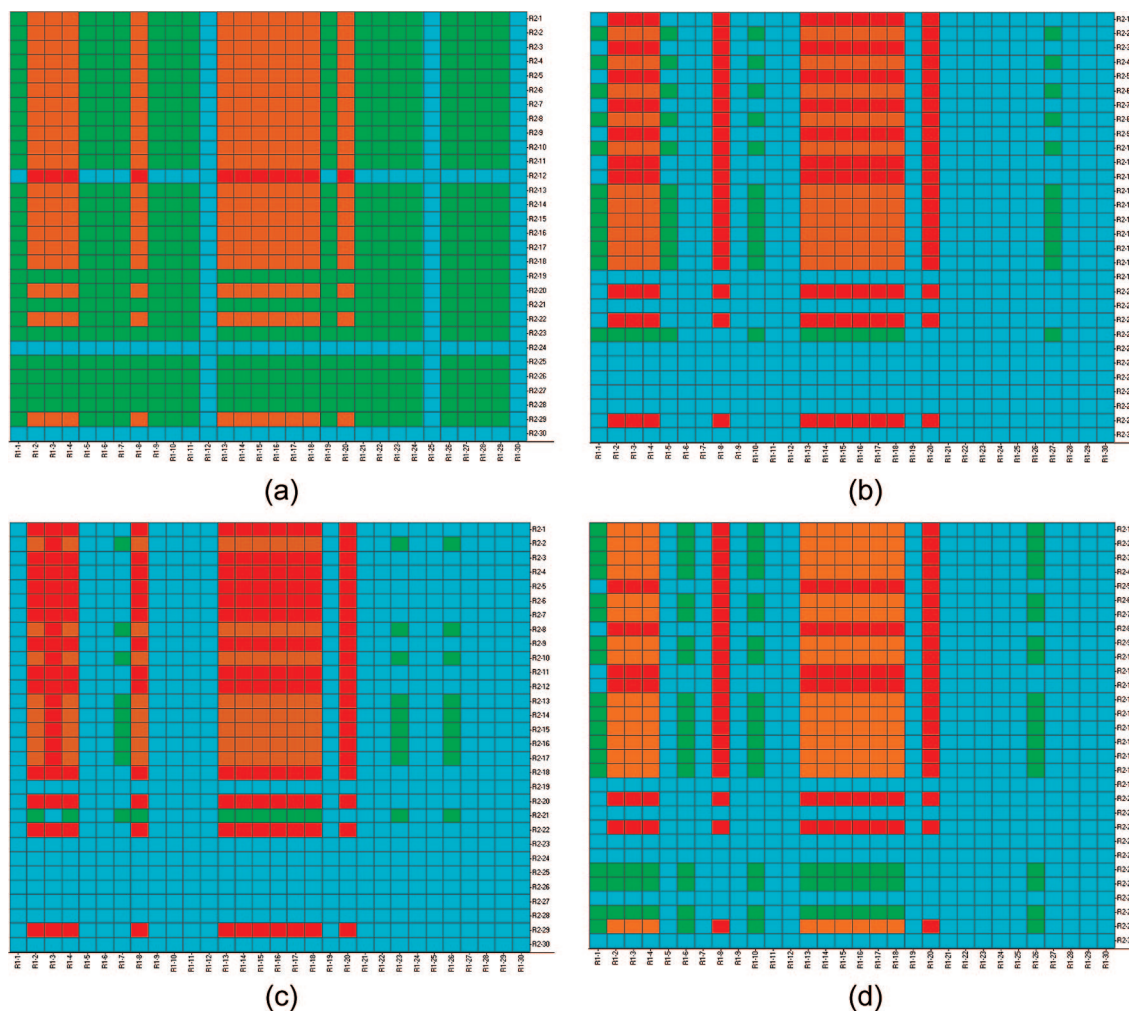


Figure 8. Comparison of the pharmacophore coverage for the Affymax library: (a) comparison of the PSAR library, (b) comparison of the random library 1, (c) comparison of the random library 2, and (d) comparison of the diversity library with the active compounds found in the Affymax library. Rows refer to pharmacophore bins of the R2 reagents, and columns refer to pharmacophore bins of the R1 reagents. The pharmacophore coverage for each library is compared with the active compounds. The red cells refer to pharmacophore bins that are present in the active compounds but not in the designed library. The green cells refer to the bins that are present in the designed library but not among the active compounds. The brown/cyano colored bins represent the pharmacophores which are present/absent in both designed library and active compounds, respectively.

Genetic algorithms (GA) or Monte-Carlo optimization is often used for library design procedure.^{5,28} An in-house library design tool GALOP was developed by using the GA optimization method to optimize the reagent pharmacophore entropy and product properties simultaneously. Figure 2 shows how the chromosome is encoded in the GA. The fitness function that the GA uses consists of two terms: one term represents the pharmacophore Shannon entropy for the reagents and the other term refers to the product properties. The fitness function formula is shown in eq 3

$$Score = w_p F + w_e \sum_j SE_j \quad (3)$$

Here, F means fraction of “good” compounds in the designed library, and SE_j refers to the Shannon entropy for the reagent set which is used for side chain j . A compound is regarded as “good” only if it meets all the specified property criteria. w_p and w_e are weighting factors for property and entropy, respectively.

Diversity Based Library Design Strategy. As comparison, a structural diversity based library design strategy was

implemented in the GALOP program and tried out in our study. The fitness function for diversity optimization is shown in eq 4

$$Score = 1 - 2 \sum_{i=1}^n \sum_{j>i}^n S_{ij} / n(n-1) \quad (4)$$

Here, S_{ij} refers to the Tanimoto similarity index between reagent i and j . So the average pairwise Tanimoto similarity for reagents will be minimized during the optimization.

COMPUTATIONAL TOOLS

Reagent pharmacophore fingerprints were generated in a two-step procedure. First, two-point pharmacophores were created by an in-house tool TRUST.²⁹ TRUST was developed based on the Open Babel toolkit.³⁰ A shell script was thereafter used to create the reagent pharmacophore fingerprints based on the TRUST output. The “greedy” search algorithm was implemented in Python³¹ to calculate and optimize reagent pharmacophore entropy. An in-house genetic algorithm driven library design tool GALOP was implemented in C++ to optimize a library in three different

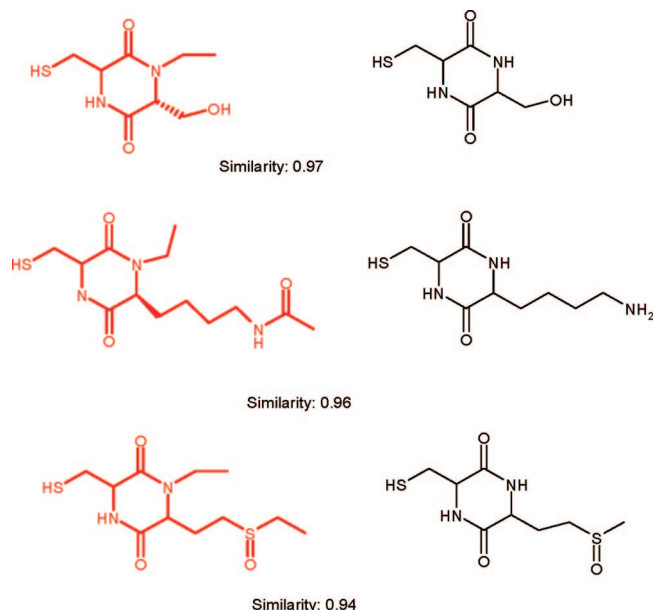


Figure 9. Some example structures for the PSAR library. The red colored structures refer to the compounds from the PSAR library, and the black colored structures refer to the Affymax active compounds (claimed as Collagenase-1, MMP-3 inhibitors³²) which have high similarity to the corresponding PSAR enumerated compounds on the left-hand side.

ways: the reagent diversity optimization, the pharmacophore entropy optimization, and the property profile minimization. Library optimization by combining diversity and property profile or by combining pharmacophore entropy with property profile can also be done in GALOP. Product properties were calculated by various in-house prediction tools. Tanimoto similarity for reagents was calculated using FOYFI fingerprint. FOYFI is an in-house developed fingerprint which is similar in spirit to the standard Daylight fingerprint.¹⁷ Database similarity searches were done by using an in-house 2D similarity search tool²⁹ and the FOYFI fingerprint. An in-house program FLUSH was used for the structure clustering. In this program, molecular distance is represented by the Tanimoto distance calculated based on FOYFI fingerprint. A threshold is set to differentiate between what one considers to be similar and different (by default is 0.3). Having set the threshold, FLUSH counts up for each molecule the number of molecules in the set that fall within this threshold. The molecule with the most neighbors is deemed to be the seed of the first cluster. All molecules in this cluster are removed from further consideration, and the process is repeated until all the molecules have been placed in a cluster. All the calculation was carried out on a Linux workstation which has two Intel Xeon cpus (2.80 GHz); the computation time for running the greedy search scripts and GALOP program on all the examples in this study is generally less than 5 min.

DATA SETS

Four test cases were investigated in this paper. The first test case is to select a 20 carboxylic acid subset from 414 carboxylic acids based on the optimal pharmacophore entropy. The second test case is to design a PSAR library based on a reaction scheme published by Affymax³² (as shown in Figure 3). The library scaffold was used as a query

for substructure search against the GVKBio database^{33,34} (a comprehensive collection of active compounds published in a variety of journals and patents). 139 active compounds were found to have the same scaffold and were therefore used as a validation set. 2518 aldehydes and 634 amino acids were used as reagents for the library enumeration. Several 20 × 20 libraries including PSAR, diversity designed and random libraries were then generated and compared with the validation set in terms of their ability to retrieve active compounds. The third case is to design a 1D library on another literature example³⁵ (shown in Figure 4) by applying PSAR strategy. The same 2518 aldehydes set were used as reagents, and several 40-compound-libraries were designed by using similar procedures as in the second test case. The last case is to design a library which has both an optimal pharmacophore entropy and a good property profile. The library reaction scheme is shown in Figure 5; it is a two-step reaction, and two types of reagent are used: a set of 112 aliphatic bromides and a set of 127 aliphatic amines. In the above studies, all reagents are from ACD.³⁶ In order to use only low complexity reagents, as previously described, selection criteria were imposed such that the maximal bond distance between a pharmacophore element and the attachment point was 6 bonds and the sum of the number of HA, HD, POS, and NEG functional groups in a reagent was no more than 2.

RESULTS AND DISCUSSION

Test Case for Carboxylic Acids. A collection of 414 carboxylic acids from ACD was selected as the first test case to demonstrate that desirable pharmacophore distribution profiles can be generated. Twenty carboxylic acids were selected using our PSAR strategy with the “greedy” search method; due to the deterministic nature of the “greedy” search algorithm, one PSAR selection was done. As a comparison, 10 randomly selected reagent collections were also generated. Jamois³⁷ et al. considered various cell-based approaches to library design, with different strategies including an entropy based one and the one that simply counts the fraction of occupied cells. In order to compare the efficiency of covering pharmacophore bins between the entropy optimized method and the occupancy optimized method, optimization of occupancy of bins is also implemented in the “greedy” search method as an alternative optimization scheme, and an occupancy optimized reagent selection was made. The total Shannon entropy of the PSAR selection is 4.15, average Shannon entropy and standard deviation of 10 random selections are 3.11 and 0.16, respectively, and Shannon entropy for occupancy based selection is 3.5. The distribution of reagent pharmacophore bins for PSAR selection, one random selection, and the occupancy selection is shown in Figure 6. The entropy optimized reagent collection can achieve the same coverage of bins with the occupancy based selection and clearly has a better coverage of the pharmacophore space than the random selection, and it has the most even distribution compared with the random selection and occupancy based selection. For example the pharmacophore corresponding to bin number 6 is missing in the random library, but the PSAR library contains this pharmacophore. For bins from 14 to 17 (corresponding to lipophilic pharmacophores), the random library’s bin count is 64 and 50 in the occupancy selection, while the frequency

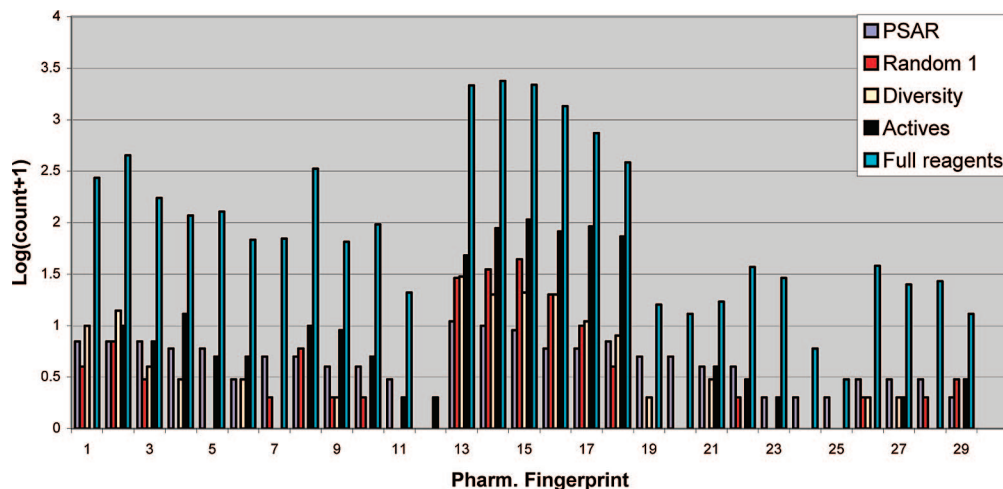


Figure 10. Pharmacophore fingerprint distribution for the Adolor library.³⁵

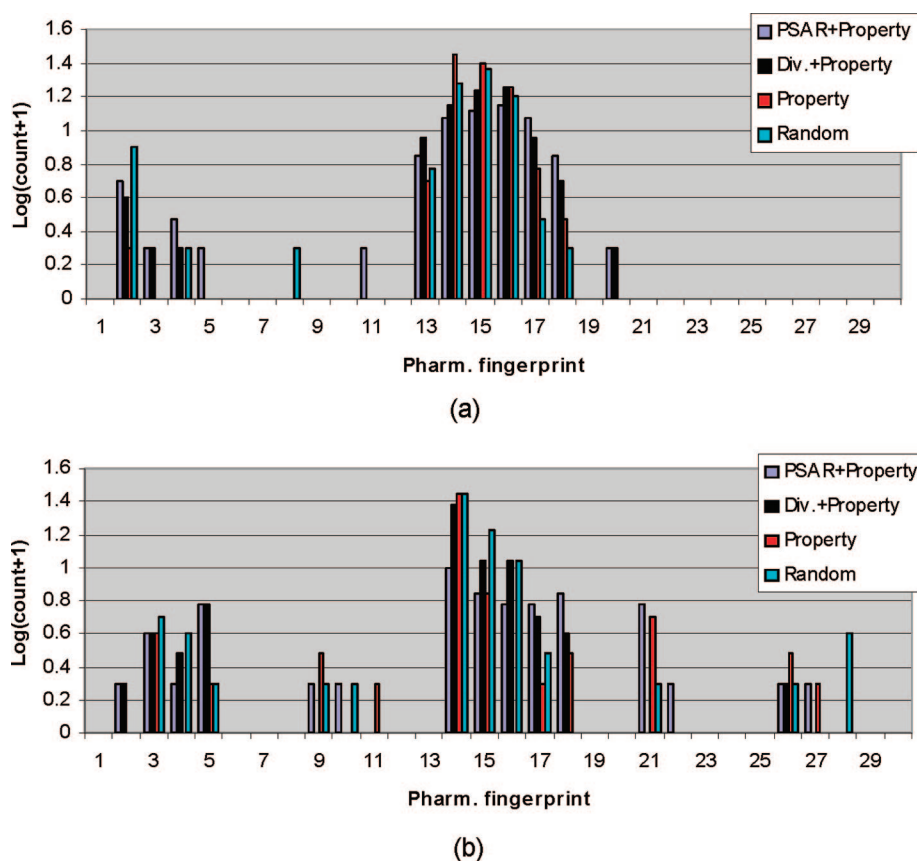


Figure 11. Comparison of pharmacophore fingerprint distribution for libraries with different design strategies: (a) pharmacophore fingerprint distribution for R1 reagents and (b) pharmacophore fingerprint distribution for R2 reagents.

of these bins has been dramatically reduced to 12 in the PSAR library. The PSAR selected reagent set does not contain bins 24 to 30, because the full reagent set does not have any reagents with these pharmacophores. These pharmacophores correspond to acids, thus the figure shows that the reagents used are monoacids, i.e. only the attachment point is an acid. This result shows that entropy optimization can achieve the same level of pharmacophore bin coverage as occupancy optimization does while having a more even distribution. So in our later studies, we make a choice of using entropy as the fitness function for pharmacophore optimization.

Affymax Library Example. As a proof-of-principle study, we compared the PSAR library with real active compounds

for a specific scaffold. A library scaffold published by Affymax was chosen as the test case for the study.³² Aldehydes and amino acids were used as building blocks to construct the library. The library has been screened for several targets,^{32,38,39} and active compounds were identified from the library. By searching the GVKBio database, 139 active compounds were found and used as a validation set. R1 (aldehyde) and R2 (amino acid) reagents were extracted from the molecular structures and encoded into pharmacophore fingerprints. Some of the reagents in the validation set have bond distances longer than 6 from the attachment point. However, when the pharmacophore fingerprints were generated, only pharmacophores up to 6 bonds from the

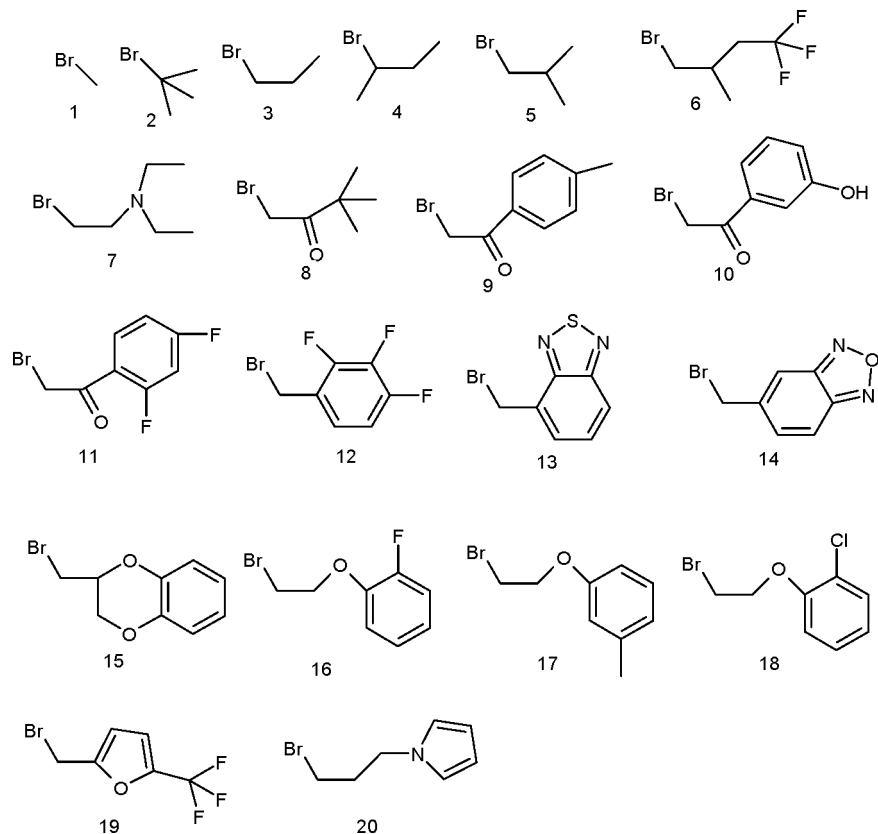


Figure 12. Selected R1 reagents for the PSAR library.

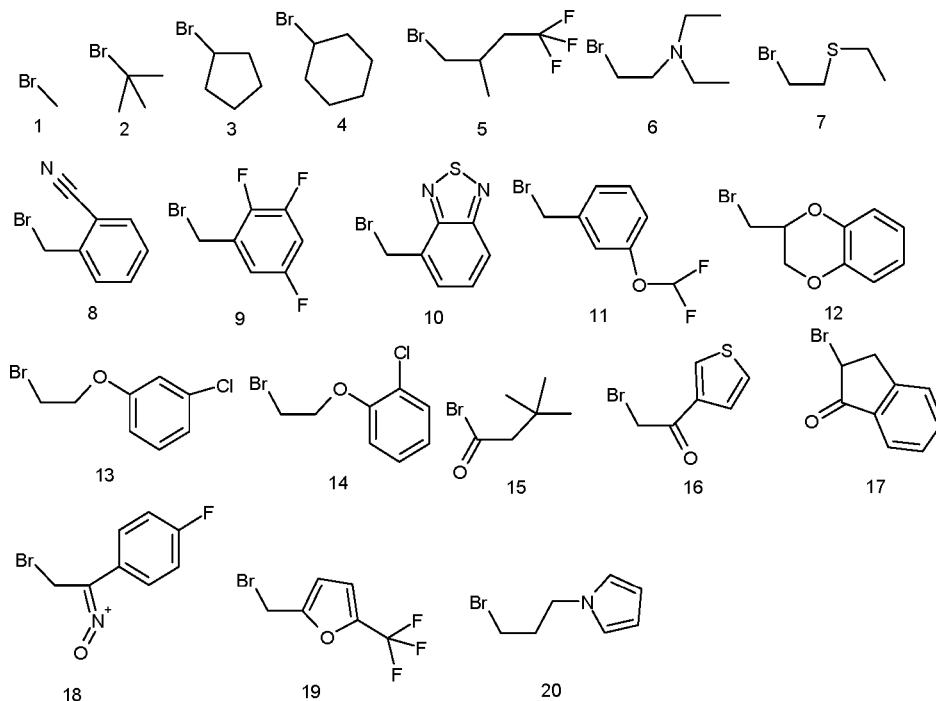


Figure 13. Selected R1 reagents for the diversity library.

attachment point were considered; pharmacophores further away than 6 bonds were ignored in this study.

Altogether 2518 aldehydes and 634 amino acids were used as the reagent pool for designing 2D libraries (20×20). A PSAR library, 10 diversity designed libraries and 10 libraries which all comprise randomly selected reagents were created from this reagent pool. The diversity library was obtained by GALOP optimization, with average Tanimoto dissimilar-

ity for the selected reagent set as the fitness function (eq 4). Figure 7 shows the distribution of pharmacophores on R1 and R2 for different reagent collections, and the results for libraries from different design strategies are summarized in Table 1. It can be seen that the PSAR selected reagent sets cover almost all of the pharmacophore bins (covered bins for R1 and R2 are both 27) and also has an even distribution among the covered bins (SE for R1 and R2 are 4.61 and

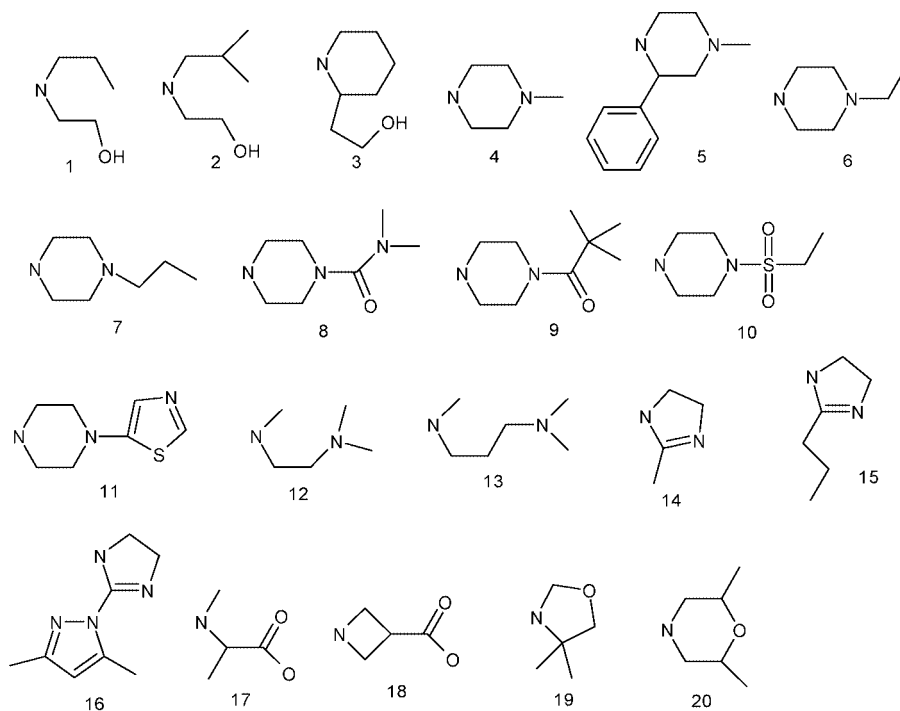


Figure 14. Selected R2 reagents for the PSAR library.

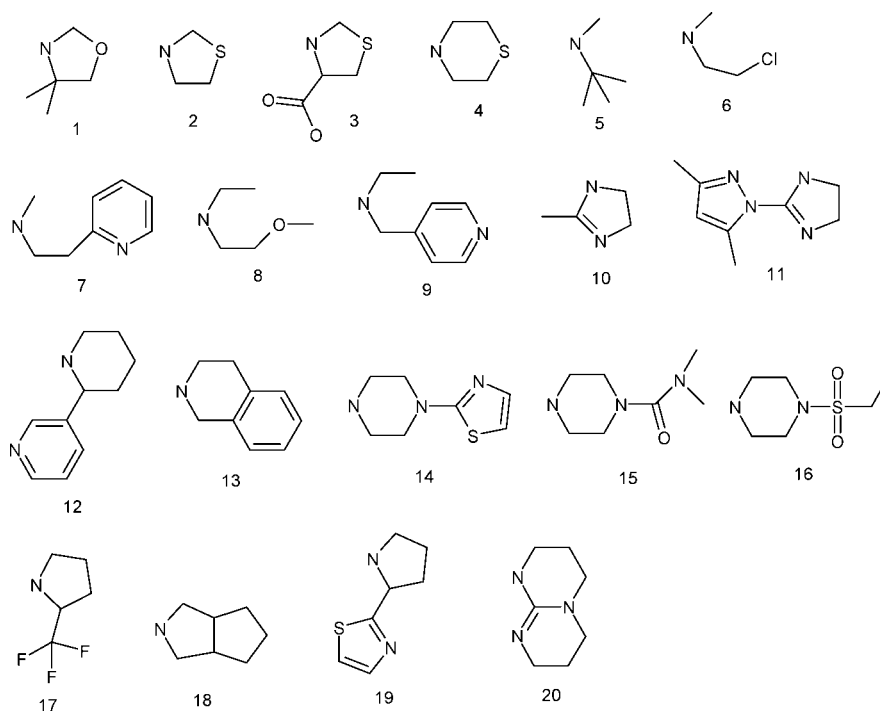


Figure 15. Selected R2 reagents for the diversity library.

4.65 respectively). Figure 8 shows the pharmacophore coverage by combining the R1 and R2 pharmacophore fingerprints into a heatmap and displaying the difference in the pharmacophore coverage between the designed library and the validation compounds. The heatmap comparison for the PSAR library and the validation set is shown in Figure 8a. Here red cells represent the pharmacophores that are present in the validation set (GVKBio actives) but are not present in the PSAR library, green cells represent the pharmacophores that are present in the PSAR library but not present in the validation set, and finally the brown or cyan cells correspond to pharmacophores which are either present

or absent in both the PSAR library and the validation set. Only 11 pharmacophore pairs (red cells) present in the validation set are not covered by the PSAR library, whereas 509 pharmacophore pairs (green cells) are only covered by the PSAR library. Figure 8b shows the comparison between random library no. 1 and the validation set. There are 132 red cells and 57 green cells here, and random library no. 2 has 159 red cells and 36 green cells (Figure 8c). Comparison between diversity library 1 and the validation set is shown in Figure 8d. It comprises 96 red cells and 99 green cells. Table 1 shows that on average random libraries miss 118.7 (54%) cells that appear in the validation set and diversity

Table 1. Results of Designed Libraries for the Affymax Example

libraries		PSAR lib.	random libraries ^d	diversity libraries ^d
no. of missed cells ^a		11	118.7	103.7
per. of missed cells		5	54	47.1
no. of additional covered cells ^b		509	72.4	81.1
no. of recovered active compounds ^c		20	11.7	1.1
Shannon	R1	4.61	3.1	3.2
entropy	R2	4.65	2.9	3.6
covered	R1	27	13.8	12.1
bins	R2	27	13.3	17.2

^a The pharmacophore cells (Figure 8) that are present in the active compounds and not present in the designed library. ^b The pharmacophore cells (Figure 8) that are present in the designed library and not present in the active compounds. ^c Retrieved GVKBio actives in similarity search, Tanimoto similarity cutoff is 0.85. ^d Average values based on 10 libraries.

libraries miss 103.7 (47.1%) cells. These results show that the PSAR library covers a vast majority of the pharmacophores that are present among the active compounds, whereas the random libraries cover the smallest number of these pharmacophores and the diversity library pharmacophore coverage lies in between. It seems that in this example the pharmacophore entropy driven PSAR library shows superior performance in covering potentially important pharmacophore elements present among the known active compounds, compared with the random library and the diversity library. Besides the pharmacophores present in actives, the PSAR library also covers many more additional pharmacophores than the other libraries (Table 1).

Although real biological activity data are not available for the compounds in our designed libraries, an estimation of the likelihood of obtaining actives from the libraries could be made by a similarity search against the GVKBio database with a high similarity cutoff (i.e., assessing the effectiveness of the libraries for generating leads). It is normally assumed that similar compounds have a higher probability of having similar bioactivity than dissimilar ones (the similar property principle),⁴⁰ and a high retrieval rate from the GVKBio database could be taken as an indication that active molecules are present in the library. Compounds in all the designed libraries were therefore used as query structures for a similarity search against the entire GVKBio database to check how many active compounds which have the same scaffold could be retrieved from the GVKBio database at a given similarity cutoff. The results are shown in Table 1. It can be seen that the PSAR library has 20 compounds that are similar to true active compounds at a similarity cutoff of 0.85. Part of results for PSAR library search are listed in Figure 9. For the random libraries, there are 11.7 similar compounds on average at a similarity cutoff of 0.85. For the diversity designed library, there are only on average 1.1 similar compounds found at this relatively high similarity cutoff. The PSAR library clearly has the best result in recovering active compounds among all libraries and, at the same time, has the highest coverage (almost 98%) of pharmacophores present in active compounds. The diversity libraries contain fewer close neighbors among known actives than random libraries do. One probable reason for this is that the average Tanimoto similarity based diversity selection algorithm, at least when not supervised by suitable physicochemical property filters, tends to select reagents which

comprise unusual or complex functional groups (however, this may not be the case for other diversity metrics^{41,42}), because the presence of “exotic” groups leads to increased Tanimoto distances when comparing structural fingerprints. Those reagents normally would not be selected in library design process; therefore, close structural neighbors to the molecules in such kind of diversity libraries are unlikely to even be synthesized and/or tested.

Adolor Library Example. A similar strategy is applied to another 1D library example published by Adolor.³⁵ The library scaffold is shown in Figure 5. Altogether 89 GVKBio compounds which comprise this scaffold are included in the validation set; they show activity against several opioid receptors.^{35,43} The same 2518 aldehydes were used as the reagents to design library, and 40 reagents were selected to enumerate library. A PSAR library, 10 random libraries, and 10 diversity libraries were designed. Distribution of pharmacophore bins on R1 position is shown in Figure 10. The PSAR library covers 28 bins in total and only misses one active bin (referring to bins covered by GVKBio active compounds), while random libraries and diversity libraries miss 7 and 7.2 bins on average. The PSAR library also shows an advantage in terms of bin distribution; its entropy value is 4.6, which is higher than that of random and diversity libraries. This means that the PSAR library is more evenly distributed on those side chain pharmacophore bins. Library performance on the likelihood of recovering active compounds was also compared in this example. The PSAR library has 21 compounds within a similarity cutoff of 0.85, while random libraries and diversity libraries have 13.9 and 9.9 instead. These results give us some confidence that the PSAR library may have advantages for finding active compounds since the compounds in this library collectively cover most of the relevant pharmacophore space.

Extension of the Optimization Strategy To Include Property Profile Optimization. So far, in the examples that are shown here, only reagent pharmacophore space diversity was considered in the library design, and compound quality was not taken into account. This is probably unrealistic for most pharmaceutical industry applications.⁴⁴ First of all, aqueous solubility is a crucial property for a compound and should be considered during library design, since it is desirable to synthesize soluble compounds. Second, risk assessment has to be done to make sure only nontoxic compounds are synthesized. For example, hERG⁴⁵ risk assessment is very important. So our strategy needs to include the compound safety profile in the optimization process. To solve this more complicated optimization problem, an in-house genetic algorithm optimizer GALOP was developed specifically to design compound libraries with multiple constraints. There are a number of published examples of this.^{46,47}

As shown in eq 3, both pharmacophore space entropy and compound property profile were included in the GA fitness function. In our experience, a weight ratio (w_e/w_p) of 2 gives reasonable and balanced libraries; this setting was used throughout the study. In the algorithm implementation, several properties were considered: (1) novelty check, (2) *in silico* predicted solubility,²⁴ (3) *in silico* predicted hERG liability,²⁵ and (4) in-house lead-like criteria.^{26,27} A “good” compound has to pass all four criteria. This extended PSAR library design strategy was applied in a hypothetical example

Table 2. Results of Designed Libraries for the Adolor Example

libraries	PSAR library	random libraries ^d	diversity libraries ^d
no. of missed bins ^a	1	7	7.2
per. of missed bins	4.8	33.3	34.3
occupied bins ^b	28	16.9	17
Shannon entropy	4.6	3.2	3.3
no. of recovered active compounds ^c	21	13.9	9.9

^a The pharmacophore bins that are present in the active compounds and not present in the designed library. ^b The pharmacophore bins that are present in the designed library. ^c Retrieved GVKBio actives in similarity search, Tanimoto similarity cutoff is 0.85. ^d Average values based on 10 libraries.

modified from an in-house reaction scheme (Figure 5). This library synthesis consists of two reaction steps: first aliphatic bromides (R1 reagents) were added to a scaffold, and the product then reacted with aliphatic amines (R2 reagents) to form a two-dimensional combinatorial library. Ten PSAR libraries were generated by running the GALOP program. In order to make a comparison, 10 libraries were assembled with randomly selected reagents, 10 diversity combining with property optimized libraries, and 10 libraries, which were only optimized by property, were also generated. For the diversity driven library design, the diversity was characterized by the Tanimoto distance of the reagents, calculated based on the in-house FOYFI fingerprint.

The final results are shown in Table 3. PSAR combined with the property calculation optimized libraries has the best reagent Shannon entropy among all other libraries on average and a high percentage of good compounds (99.7%). Diversity combining property optimized libraries generally has the best diversity; this can be seen from its average FOYFI Tanimoto distance and number of clusters, and its percentage of good compounds is also 99.7%. For the libraries, which are only optimized by property, they have both worse Shannon entropy and diversity compared with the PSAR and diversity driven libraries. As a baseline, the fully enumerated library has only 62% good compounds in total and medium entropy and diversity values. As expected, the random libraries have medium Shannon entropy and diversity and contains 60% good compounds on average, while PSAR libraries have the highest Shannon entropy values on both the R1 and R2 substituents. With respect to pharmacophore bin coverage, the PSAR libraries do not show better coverage on the R1 position than random and diversity driven libraries because compounds with desirable property profile have limited variation on the R1 position; however, on the R2 position, the bin coverage of PSAR libraries is superior to those of any other libraries.

For convenience of comparison, one of the PSAR libraries and one diversity library were randomly chosen for closer examination. The pharmacophore distribution for R1 (bromides) and R2 (aliphatic amines) reagents is shown in Figure 11, and the results for different libraries are summarized in Table 2. For the R1 reagents, the entropy of the PSAR library (3.10) is slightly better than the diversity library (2.8), and we have noted that there is a 50% overlap (10 out of 20) of the selected R1 reagents for the PSAR and diversity library; as we said before, this is mainly due to limited variation on R1 positions to compounds which has good property profile.

Compared with the PSAR library, the diversity library does not contain bins 5 and 11, which correspond to an HA (hydrogen bond acceptor) and an HD (hydrogen bond donor), 5 bonds from the attachment point, respectively. The 20 selected bromides for the PSAR library and the diversity library are shown in Figures 12 and 13. It can be seen that among the PSAR selected R1 reagents, structure **14** has an HA (aromatic N atom) element which is separated by 5 bonds from the scaffold and structure **10** has an HD group (hydroxyl group) which is 5 bonds away from the attachment point. Those pharmacophore features are missing in the diversity library. Bins 13 to 18 refer to the lipophilic elements in the reagents; for the random library and property optimized only library, there are very high occurrences of these lipophilic pharmacophores. In the PSAR library occurrences of these bins are noticeably reduced, thus making the library more balanced among the five types of pharmacophore elements.

For the R2 reagents, the entropy of the PSAR library is markedly better than any other libraries including the diversity library. Lipophilic pharmacophores corresponding to bin 14, which is over-represented in the diversity library and the property optimized library, occur much less in the PSAR library. Pharmacophores corresponding to bin numbers 9, 10, 21, 22, and 27 are missing in the diversity library compared to the PSAR library. In terms of HD functionality, the diversity library at the R2 position lacks this functionality completely, while the PSAR library has 3 compounds (structures **1**, **2**, and **3** in Figure 14) which have HD at either 3 or 4 bond distances. When considering positively charged centers (POS), the diversity library at the R2 position lacks this functionality entirely, the random library and the property optimized library occupy only one POS bin, while the PSAR library has 6 reagents which offer POS functionalities at two different bond distances to the attachment point. It seems that the PSAR library comprises a more balanced reagent set, in terms of pharmacophoric features, and has a larger variation of pharmacophore elements compared with the diversity library (see Figure 15).

Comparing with the reagent set in the diversity library, the PSAR reagent set on R1,R2 positions comprise more structurally related compounds which may serve as a good starting point for deriving SAR once bioactivity data are available. Taking a closer look at the PSAR R2 reagents in Figure 14, for structures **1**, **2**, and **3**, they are similar structures and have variations on the HD functionalities; this could potentially help to derive SAR around the HD functionality on side chains. Similarly, structures **12** and **13** (Figure 14) can provide SAR around the POS functionality, and structures **4–11** may show some SAR around the piperazine ring. For the diversity driven library design strategy that we currently use, these types of reagents will have less chance of being selected because the structural similarity among these reagents will probably decrease the fitness score during optimization. Regarding the compound properties of these designed libraries, since the property control is included in the GA optimization, all the GA optimized libraries have a high percentage of “good” compounds. The percentage of good compounds in the PSAR/property optimized library and diversity/property optimized library are 99.7% (Table 3), while the purely property optimized libraries have 100% “good” compounds.

Table 3. Results for the GA Optimized Libraries^a

libraries		PSAR+ property ^b	diversity+ property ^c	property ^d	random library ^e	full library
per. of good compounds		99.7	99.7	100	62.2	62
no. of clusters		21	46.1	14.1	23	NC
Shannon entropy	R1	3.03	2.86	2.38	2.71	2.83
	R2	3.52	2.62	2.32	2.81	2.94
diversity	R1	0.74	0.80	0.64	0.72	0.74
	R2	0.69	0.80	0.65	0.71	0.73
covered bins	R1	10.5	10.3	7	10.7	21
	R2	15.4	10.2	10.5	12	20

^a The values listed in the table are averaged based on 10 runs, except for the full library. ^b Libraries obtained by optimizing both the pharmacophore entropy and the property profile simultaneously. ^c Libraries obtained by optimizing both the diversity and the property profile simultaneously. ^d Libraries obtained by only optimizing the property profile. ^e Libraries obtained by randomly selecting reagents.

In summary, the PSAR optimized library is better than the diversity library, which is optimized by average Tanimoto similarity, in terms of reagent pharmacophore distribution, because PSAR libraries tend to include structurally related reagents that comprise systematically varied side chain pharmacophore elements. This kind of feature can be of assistance to chemists attempting to derive SARs once assay results become available. The average Tanimoto similarity derived strategy will, on the other hand, create a more structurally diverse compound set, so the choice of strategy depends on what type of diversity the library design wants to achieve.

CONCLUSIONS

We have presented a library design strategy to optimize the reagent pharmacophore space for constructing a so-called "ProSAR" library. A new way of encoding pharmacophore elements in reagents has been put forward to express systematic variation of pharmacophore elements relative to a scaffold attachment. Shannon entropy was used to represent the pharmacophore coverage of the reagent space, and optimizing this property derives a library with an optimal coverage and even distribution of the pharmacophore elements among the reagents, thus potentially making it easier for medicinal chemists to derive SARs. We illustrate this by two examples, where most of the pharmacophore features that appear in active compounds could be covered by the PSAR derived library. Furthermore, in those examples PSAR libraries include more compounds that are highly structurally similar to true active compounds than random libraries and diversity libraries which are optimized by average ensemble Tanimoto similarity. The PSAR strategy was further expanded to include compound properties to design a library which has not only good pharmacophore coverage of side chains but also desirable physicochemical properties. A GA optimization program was developed for this purpose and applied to an illustrative test case, in which the aim was to design a 400 compound two-dimensional combinatorial library. We contrast this with libraries constructed by other design strategies such as diversity (characterized by the average ensemble Tanimoto similarity in this study) driven and property driven library design. Our results demonstrate that "ProSAR" designed libraries are clearly superior in covering pharmacophore space and create more even distribution of the side chain pharmacophore elements.

ACKNOWLEDGMENT

The authors are grateful to the following colleagues at AstraZeneca: Dr. David Cosgrove for providing the FOYFI

fingerprint calculating programs and critical revision of the manuscript, diploma worker Sabbath Marchend for collecting the Adolor library example, Dr. Jens Sadowski for providing the tool to extract the R groups for library compounds, and Dr. Markus Haeberlein and Dr. Stefan Schmitt for valuable discussions.

Supporting Information Available: ACD reagents used in our examples and Affymax and Adolor active compounds from the GVKBio database. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- (1) Spellmeyer, D. C.; Grootenhuys, P. D. J. Recent Developments in Molecular Diversity: Computational Approaches to Combinatorial Chemistry. *Annu. Rep. Med. Chem. Rev.* 1999, 34, 287–296.
- (2) Beno, B. R.; Mason, J. S. The design of Combinatorial Libraries Using Properties and 3D Pharmacophore Fingerprints. *Drug Discovery Today* 2001, 6, 251–258.
- (3) Willett, P. Chemoinformatics - Similarity and Diversity in Chemical Libraries. *Curr. Opin. Biotechnol.* 2000, 11, 85–88.
- (4) Turner, D. B.; Tyrrell, S. M.; Willett, P. Rapid Quantification of Molecular Diversity for Selective Database Acquisition. *J. Chem. Inf. Comput. Sci.* 1997, 37, 18–22.
- (5) Jamois, E. A. Reagent-based and Product-based Computational Approaches in Library Design. *Curr. Opin. Chem. Biol.* 2003, 7, 326–330.
- (6) Potter, T.; Matter, H. Random or Rational Design? Evaluation of Diverse Compound Subsets from Chemical Structure Databases. *J. Med. Chem.* 1998, 41, 478–488.
- (7) McGregor, M. J.; Muskal, S. M. Pharmacophore Fingerprinting. 2. Application to Primary Library Design. *J. Chem. Inf. Comput. Sci.* 2000, 40, 117–125.
- (8) Zheng, W.; Cho, S. J.; Tropsha, A. Rational combinatorial library design. 1. Focus-2D: a new approach to targeted combinatorial chemical libraries. *J. Chem. Inf. Comput. Sci.* 1998, 38, 572–584.
- (9) Leach, A. R.; Green, D. V. S.; Hann, M. M.; Judd, D. B.; Good, A. C. Where are the gaps? A rational approach to monomer acquisition and selection. *J. Chem. Inf. Comput. Sci.* 2000, 40, 1262–1269.
- (10) Gillet, V. J.; Willett, P.; Bradshaw, J. The effectiveness of reactant pools for generating structurally diverse combinatorial libraries. *J. Chem. Inf. Comput. Sci.* 1997, 37, 731–740.
- (11) Shannon, C. E.; Weaver, W. *The Mathematical Theory of Communication*; University of Illinois Press: Urbana, IL, 1963.
- (12) Godden, J. W.; Stahura, F. L.; Bajorath, J. Variabilities of Molecular Descriptors in Compound Databases Revealed by Shannon Entropy Calculations. *J. Chem. Inf. Comput. Sci.* 2000, 40, 796–800.
- (13) Lamb, M. L.; Bradley, E. K.; Beaton, G.; Bondy, S. S.; Castellino, A. J.; Gibbons, P. A.; Suto, M. J.; Grootenhuys, P. D. J. Design of a gene family screening library targeting G-protein coupled receptors. *J. Mol. Graphics Modell.* 2004, 23, 15–21.
- (14) Bradley, E. K.; Miller, J. L.; Saiah, E.; Grootenhuys, P. D. J. Informative Library Design as an Efficient Strategy to Identify and Optimize Leads: Application to Cyclin-Dependent Kinase 2 Antagonists. *J. Med. Chem.* 2003, 46, 4360–4364.
- (15) Miller, J. L.; Bradley, E. K.; Teig, S. L. Luddite: An Information-Theoretic Library Design Tool. *J. Chem. Inf. Comput. Sci.* 2003, 43, 47–54.
- (16) McGregor, M. J.; Muskal, S. M. Pharmacophore Fingerprinting. 1. Application to QSAR and Focused Library Design. *J. Chem. Inf. Comput. Sci.* 1999, 39, 569–574.

- (17) *Daylight Theory Manual*; Daylight Chemical Information Systems, Inc. <http://www.daylight.com/dayhtml/doc/theory/> (accessed Oct. 14, 2008).
- (18) Cavalli, A.; Poluzzi, E.; De Ponti, F.; Recanatini, M. Toward a pharmacophore for drugs including the long QT syndrome: insights from a CoMFA study of HERG K(+) channel blockers. *J. Med. Chem.* **2002**, *45*, 3844–3853.
- (19) Pearlstein, R. A.; Vaz, R. J.; Kang, J.; Chen, X.-L.; Preobrazhenskaya, M.; Shchekotikhin, A. E.; Korolev, A. M.; Lysenkova, L. N.; Miroshnikova, O. V.; Hendrix, J.; Rampe, D. Characterization of HERG potassium channel inhibition using CoMSiA 3D QSAR and homology modeling approaches. *Bioorg. Med. Chem. Lett.* **2003**, *13*, 1829–1835.
- (20) Jouyban, A.; Soltanpour, S.; Soltani, S.; Chan, H. K.; Acree, W. E. Solubility prediction of drugs in water-cosolvent mixtures using Abraham solvation parameters. *J. Pharm. Sci.* **2007**, *10*, 263–77.
- (21) Egan, W. J.; Merz, K. M.; Baldwin, J. J. Prediction of drug absorption using multivariate statistics. *J. Med. Chem.* **2000**, *43*, 3867–3877.
- (22) Darvas, F.; Dorman, G.; Papp, A. Diversity measures for enhancing ADME admissibility of combinatorial libraries. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 314–322.
- (23) Muresan, S.; Kocis, P.; Chen, H.; Steele, J.; Li, J. Manuscript in preparation.
- (24) Bruneau, P. Search for Predictive Generic Model of Aqueous Solubility Using Bayesian Neural Nets. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1605–1616.
- (25) Gavaghan, C. L.; Arnby, C. H.; Blomberg, N.; Strandlund, G.; Boyer, S. Development, interpretation and temporal evaluation of a global QSAR of hERG electrophysiology screening data. *J. Comput.-Aided Mol. Des.* **2007**, *21*, 189–206.
- (26) Oprea, T. I.; Davis, A. M.; Teague, S. J.; Leeson, P. D. Is there a difference between leads and drugs? A historical perspective. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1308–1335.
- (27) Oprea, T. I. Current trends in lead discovery: Are we looking for the appropriate properties. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 325.
- (28) Reynolds, C. H.; Tropsha, A.; Pfahler, D. B.; Druker, R.; Chakravorty, S.; Ethiraj, G.; Zheng, W. Diversity and coverage of structural sublibraries selected using the SAGE and SCA algorithms. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1470–1477.
- (29) Kogej, T.; Engkvist, O.; Blomberg, N.; Muresan, S. Multifingerprint Based Similarity Searches for Targeted Class Compound Selection. *J. Chem. Inf. Model.* **2006**, *46*, 1201–1213.
- (30) The Open Babel package, version 2.0.1. http://openbabel.org/wiki/Main_Page (accessed Oct. 14, 2008).
- (31) Python Programming Language Official Website. <http://www.python.org/> (accessed Oct. 14, 2008).
- (32) Szardenings, A. K.; Antonenko, V.; Campbell, D. A.; DeFrancisco, N.; Ida, S.; Si, L.; Sharkov, N.; Tien, D.; Wang, Y.; Navre, M. Identification of Highly Selective Inhibitors of Collagenase-1 from Combinatorial Libraries of Diketopiperazines. *J. Med. Chem.* **1999**, *42*, 1348–1357.
- (33) *GVKBio database 2007*; GVK Biosciences Private Ltd.: Hyderabad 500016, India.
- (34) Southan, C.; Varkonyi, P.; Muresan, S. Complementarity Between Public and Commercial Databases: New Opportunities in Medicinal Chemistry Informatics. *Curr. Top. Med. Chem.* **2007**, *7*, 1502–1508.
- (35) Le Bourdonnec, B.; Belanger, S.; Cassel, J. A.; Stabley, G. J.; DeHaven, R. N.; Dolle, R. E. trans-3,4-Dimethyl-4-(3-carboxamidophenyl) piperidines: A Novel Class of μ -Selective Opioid Antagonists. *Bioorg. Med. Chem. Lett.* **2003**, *13*, 4459–4462.
- (36) *MDL Available Chemicals Directory database 2007*; Symyx Technologies, Inc.: Santa Clara, CA 95051.
- (37) Jamois, E. J.; Hassan, M.; Waldman, M. Evaluation of reagent-based and product-based strategies in the design of combinatorial library subsets. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 63–70.
- (38) Campbell, D. A.; Look, G. C.; Szardenings, A. K.; Patel, P. V. US6271232B1, 2001. Campbell, D. A.; Look, G. C.; Szardenings, A. K.; Patel, P. V. US5932579A, 1999. Campbell, D. A.; Look, G. C.; Szardenings, A. K.; Patel, P. V. WO97/48685A1, 1997.
- (39) Szardenings, A. K.; Harris, D.; Lam, S.; Shi, L.; Tien, D.; Wang, Y.; Patel, D. V.; Navre, M.; Campbell, D. A. Rational Design and Combinatorial Evaluation of Enzyme Inhibitor Scaffolds: Identification of Novel Inhibitors of Matrix Metalloproteinases. *J. Med. Chem.* **1998**, *41*, 2194–2200.
- (40) Martin, Y. C.; Kofron, J. L.; Traphagen, L. M. Do structurally similar molecules have similar biological activity. *J. Med. Chem.* **2002**, *45*, 4350–4358.
- (41) Agrafiotis, D. K.; Lobanov, V. S. An Efficient Implementation of Distance-Based Diversity Measures Based on k - d Trees. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 51–58.
- (42) Waldman, M.; Li, H.; Hassan, M. Novel algorithms for the optimization of molecular diversity of combinatorial libraries. *J. Mol. Graphics Modell.* **2000**, *18*, 412–426.
- (43) Mitch, C. H.; Quimby, S. J. WO03/101963A1, 2003. Armer, R. E.; Gethin, D. M.; Gibson, S. P.; Tommasini, I. EP1072592A2, 2001. Mchardy, S. F.; Liras, S.; Guediche, S. A.; Coe, J. W. WO2004/089370A1, 2004. Mchardy, S. F.; Liras, S.; Guediche, S.; Coe, J. W. US20050032837A1, 2005. Dolle, R. E.; Le Bourdonnec, B. WO2004/082623A2, 2004. Mchardy, S. F.; Liras, S.; Guediche, S.; Coe, J. W. US20040204453A1, 2004. Mitch, C. H.; Quimby, S. J. US20050222204A1, 2005. Coe, J. W.; Iredale, P. A.; Mchardy, S. F.; Mclean, S. WO2005/018670A1, 2005.
- (44) Pickett, S. D.; McLay, I. M.; Clark, D. E. Enhancing the Hit-to-Lead Properties of Lead Optimization Libraries. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 263–272.
- (45) Keating, M. T.; Sanguinetti, M. C. Molecular genetic insights into cardiovascular disease. *Science* **1996**, *272*, 681–685.
- (46) Gillet, V. J.; Khatlib, W.; Willett, P.; Fleming, P. J.; Green, D. V. S. Combinatorial library design using a multiobjective genetic algorithm. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 375–385.
- (47) Brown, R. D.; Hassan, M.; Waldman, M. Combinatorial library design for diversity, cost efficiency, and drug-like character. *J. Mol. Graphics Modell.* **2000**, *18*, 427–437.

CI800231D