

# Multiobjective Optimization of Combinatorial Libraries

Dimitris K. Agrafiotis\*

3-Dimensional Pharmaceuticals, Inc., 665 Stockton Drive, Exton, PA 19341, USA

## ABSTRACT

Combinatorial chemistry and high-throughput screening have caused a fundamental shift in the way chemists contemplate experiments. Designing a combinatorial library is a controversial art that involves a heterogeneous mix of chemistry, mathematics, economics, experience, and intuition. Although there seems to be little agreement as to what constitutes an ideal library, one thing is certain: only one property or measure seldom defines the quality of the design. In most real-world applications, a good experiment requires the simultaneous optimization of several, often conflicting, design objectives, some of which may be vague and uncertain. In this paper, we discuss a class of algorithms for subset selection rooted on the principles of multiobjective optimization. Our approach is to employ an objective function that encodes all the desired selection criteria, and then use a simulated annealing or evolutionary approach to identify the optimal (or a nearly optimal) subset from among the vast number of possibilities. Virtually any conceivable design criterion can be accommodated, including diversity, similarity to known actives, predicted activity and/or selectivity as determined by some QSAR or receptor binding model, enforcement of certain property distributions, reagent cost and availability, and many others. The method is robust, convergent, and extensible, offers the user full control over the relative significance of the various objectives in the final design, and permits the simultaneous selection of compounds from multiple libraries in full or sparse array format.

---

\* Tel: (610) 458-6045, Fax: (610) 458-8249, E-mail: [dimitris@3dp.com](mailto:dimitris@3dp.com).

## **INDEX TERMS**

Combinatorial chemistry, combinatorial library, high-throughput screening, molecular diversity, molecular similarity, quantitative structure-activity relationship, QSAR, multiobjective optimization, multicriteria optimization, Pareto optimization, simulated annealing, nonlinear mapping, multidimensional scaling.

## **I. INTRODUCTION**

Historically, drug discovery has been based on a serial and systematic modification of chemical structure aimed at producing compounds which can effectively and safely alter the activity of biological targets associated with a particular disease. This process involves four major steps following the identification of a biological target: 1) hit identification, 2) lead generation, 3) lead optimization, and 4) target validation. The first part of this process is carried out by screening large compound collections, such as combinatorial libraries, natural product collections, corporate banks, etc, to identify compounds that interact with the target enzyme or receptor. Once a hit has been identified, it is chemically modified by iterative synthesis and testing of related analogs to produce leads, i.e. compounds with improved chemical characteristics that are more suitable as potential drugs. Further chemical modification optimizes the properties of these leads and converts them into drug development candidates for further preclinical and clinical development.

Prior to the advent of combinatorial chemistry, this process involved a simple prioritization of synthetic targets based on pre-existing structure-activity data, synthetic feasibility, experience, and intuition. This situation began to change as advances in parallel synthesis and high-throughput screening have enabled the simultaneous synthesis and biological evaluation of large chemical libraries containing hundreds to tens of thousands of compounds [1]. Although throughput has increased dramatically, the number of compounds that can be made and tested in a reliable manner represents a tiny fraction of all the molecules of potential pharmaceutical interest, and the process is still fundamentally based on trial-and-error. It is becoming increasingly apparent that in order to maximize the probability of identifying sustainable drug candidates, combinatorial experiments must be carefully planned and take full advantage of whatever information is available about the biological target of interest. Whether it is used for lead discovery or optimization, the design of a good library is a complex task that requires the simultaneous optimization of several, often conflicting, design objectives. In this paper, we present an overview of a general methodology for designing combinatorial and high-throughput screening experiments rooted on the principles of multiobjective optimization.

Multiobjective (MO) optimization (also known as multicriterion, vector, or Pareto optimization) extends optimization theory by permitting several design objectives to be optimized simultaneously. Although the basic theorems can be traced back to the work of Leibniz and Euler, the principle of multiobjective optimization was first formalized by Vilfredo Pareto, an Italian economist, whose theories [2] are now considered the basis of modern welfare economics exercised by socialist economic scholars. He introduced the concept of the *Pareto optimum*, a standard of judgment in which the optimum allocation of the resources of a society is not attained as long as it is possible to make at least one individual better off in his own estimation while keeping others as well off in their own estimation. In Pareto's own words: "The principal subject of our study is [economic] equilibrium. This equilibrium results from the opposition between men's tastes and the obstacles to satisfying them. Our study includes, then, three distinct parts: 1) the study of tastes; 2) the study of obstacles; and 3) the study of the way in which these two elements combine to reach equilibrium." Throughout the years, Pareto's basic theories have been extended [3] and applied to a wide spectrum of optimization problems in economics, management, engineering, and social sciences.

A multiobjective optimization problem is solved in a manner similar to the conventional single-objective (SO) problem. The goal is to find a set of values for the design variables that simultaneously optimize several objective (or cost) functions. The solutions are often referred to as Pareto optima, vector maxima, efficient points, or non-dominated solutions. In general, the solutions obtained by individual optimization of each objective (i.e. SO optimization) do not represent a feasible solution to the multiobjective problem. Several methods have been devised for constructing Pareto-optimal sets, including hierarchical optimization, weighing objectives, distance functions, goal-programming, trade-off or constraint methods, min-max optimization, and many others [3]. These methods fall into two broad categories: 1) methods which attempt to optimize each criterion in turn, subject to constraints derived from the optimization of previously optimized criteria, and 2) methods which attempt to minimize a single objective function that combines the design objectives in some prescribed functional form (often referred to as a global criterion).

Unfortunately, the vast majority of multiobjective optimization theory deals with continuous or semi-continuous design variables, and is not directly applicable to discrete, combinatorial problems such as the

application at hand. In this paper, we review the algorithmic details of a general methodology for designing combinatorial libraries that combines the flexibility of multiobjective fitness functions with the power of simulated annealing for searching vast combinatorial state spaces. The method allows traditional design objectives such as molecular similarity or molecular diversity to be combined with other selection criteria in order to enforce certain drug-like property distributions, contain the cost of the experiment, minimize potential toxicological and pharmacokinetic liabilities, etc.

Library design has evolved into a distinct sub-discipline of computational chemistry standing on the crossroads of organic synthesis, chemometrics, QSAR and structure-based design. Although it is not our intention to provide an in-depth review of this field, the reader should be aware that the approach described here is one of many available methodologies for designing effective and well-targeted combinatorial experiments, and one that has been pursued independently by several groups. While some of the underlying principles have their origin in the field of statistical experimental design, the first reports dealing specifically with combinatorial libraries were presented in 1995 by five independent groups [4, 5, 6, 7, 8]. In what is considered by many the pivotal computational work on molecular diversity, Martin and coworkers reported a rational method for selecting a set of monomers for a peptoid combinatorial library based on D-optimal design [4]. To ensure that the design would capture biologically relevant information, the group employed a wide range of molecular descriptors that captured lipophilicity, shape and branching, chemical functionality, and receptor binding. Principal component analysis and multidimensional scaling were employed to reduce the dimensionality of the original data, and a D-optimal design procedure was used to select a representative subset of reagents to optimally explore the resulting diversity space.

An alternative approach that is more closely related to the one described herein was presented by groups at Merck, Hoffman-La Roche, and Sterling Winthrop Pharmaceuticals. Sheridan and Kearsley [6] employed a genetic algorithm to generate N-substituted glycine tripeptides which maximized similarity to a known lead or predicted activity according to a trend vector developed from a series of active molecules. Weber et al [7] presented a similar genetic scheme using a fitness function that was based on the results of an experimental enzyme assay rather than a theoretical QSAR model, and Singh *et al* [8] employed a similar methodology to optimize a set of potent and selective hexapeptide stromelysin substrates.

The group at 3DP was also quick to recognize that biological data is the best guide for the design of combinatorial and high-throughput screening experiments, and that effective data management and process integration were key for the timely and cost-effective development of new therapeutic agents. Contemporaneously to the aforementioned publications, our group was being awarded the first of a series of patents describing a new drug discovery paradigm that integrates combinatorial, structural, and computational chemistry under a unifying information management system [4]. The system, known as DirectedDiversity®, is an iterative optimization process that explores combinatorial space through successive rounds of selection, synthesis and testing. Unlike traditional combinatorial approaches where the entire library is made and tested in a single conceptual step, DirectedDiversity® physically synthesizes, characterizes and tests only a portion of that library at a time. The selection of compounds is carried out by computational search engines that combine optimal exploration of molecular diversity with a directed search based on SAR information accumulated from previous iterations of the integrated machinery. The original blueprint was based on an optimization scheme that was in many respects similar to that described by the previous authors, but emphasized more strongly the need to adjust the selection criteria in the course of a discovery program and deal with the ambiguity that is inherent in the biological response data.

In the following years, these algorithms were elaborated, calibrated, and applied to a wide range of problems in library design. From an algorithmic perspective, the most notable examples include Agrafiotis' [9, 10, 11, 12] and Hassan's *et al* [13, 14] independently developed simulated annealing implementations and their subsequent variations by Good and Lewis [15] and Zheng *et al* [16], Brown and Martin's genetic scheme to generate libraries designed to minimize the effort required to deconvolute biological hits by mass-spectroscopic techniques [17], Gillet's *et al* [18], Rassokhin *et al* [19], and Brown's *et al* [20] attempts to enforce certain property distributions on the final design, and Sheridan's *et al* [21] latest use of genetic algorithms for designing targeted libraries. These advances were complemented by important validation studies, which compared the ability of several popular sets of descriptors to differentiate between active and inactive molecules [22, 23, 24, 25, 26, 27], as well as the relative merits of reagent versus product-based designs [28, 29]. These represent only a small fraction of proposed library design methodologies, which range from conventional experimental design [4, 30, 31] to clustering [32] and cluster sampling [33], conformational sampling [34], partitioning [35, 36], boolean logic [37, 38, 39],

vector analysis [40], and some more recent greedy algorithms for selecting combinatorial arrays [41, 42]. For a more extensive account, the interested reader is referred to several recently published reviews [31, 43, 44, 45].

## II. METHODS

**Architecture** – In its prototypical form, the selection problem can be stated as follows: given a collection of  $n$  compounds and a number  $k$ , find the ‘best’ set of  $k$  compounds according to some user-defined criteria. This problem is NP-complete, and the cardinality of the search space is enormous even for the most conservative cases encountered in library design. The approach taken here is to combine all the selection criteria into a single unifying function, and maximize (or minimize) that function using an efficient optimization algorithm such as simulated annealing or evolutionary programming in order to identify the optimal (or a nearly optimal) set among the vast number of possibilities. To simplify the description of the algorithm, we define the following entities: *collections* or *libraries* which represent separate pools of chemical compounds (combinatorial libraries and/or regular collections) from which the selection is to be drawn, *subsets* which represent selections of compounds from a particular compound collection, *states* or *designs* which represent a collection of subsets from one or more chemical libraries, *selection criteria* which encode the individual design objectives, *objective functions* which combine one or more selection criteria in some arbitrary functional form and provide the overall quality of a particular design (state), and *optimizers* which search through the state space associated with the problem of interest to identify the optimal (or a nearly optimal) solution.

The overall architecture is shown in Fig. 1. An optimizer (in this case a serial or parallel implementation of simulated annealing) produces a state (i.e. a collection of subsets from one or more chemical libraries), which is evaluated against all the desired selection criteria. These are combined into a unifying objective function, which measures the overall fitness of that state, that is, its ability to collectively satisfy all the specified selection criteria. This fitness value is used by the optimizer to produce a new set of compounds (i.e. a new state), which is, in turn, evaluated against the prescribed selection criteria in the manner outlined

above. The process continues until a predefined termination criterion or time limit is met, and the best state identified during the course of the simulation is reported. This general scheme can be implemented in a serial or parallel fashion; in the later case, several states are evaluated in parallel and are subsequently combined to produce a new set of states for the next iteration (see below).

The major advantage of this approach is that the search algorithm is completely independent of the performance measure, and can be applied on a wide variety of selection criteria and fitness functions. Unlike alternative algorithms such as maxmin [46], cluster analysis [22], binning [36], stepwise elimination [33], etc., which are tailored to a particular application, this approach is completely general, programmatically simple, and easily extensible. The remaining paragraphs describe in detail the various elements of this approach, using a mixed terminology borrowed from the simulated annealing and evolutionary programming literature.

**States** – A *state* or *design* represents a collection of subsets from one or more chemical libraries. The system was designed to allow the simultaneous selection of multiple subsets from multiple collections, and thus enable the design of experiments that span multiple chemistries and/or corporate files. Depending on the nature of the parent collection, three types of subsets can be defined: *sparse arrays* (or singles), *full arrays* (or simply arrays) and *plates*.

**Sparse Arrays** – A sparse array represents any conceivable subset of  $k$  compounds from an  $n$ -membered collection. This is the most general design and is not restricted by the nature or location of the compounds, or the types of reagents required for their synthesis. The term array originates from the combinatorial chemistry literature, and refers to a subset of compounds from a combinatorial library that does not necessarily represent all possible combinations of the selected building blocks. This type of selection is illustrated in Fig. 2a and is formally defined as:

$$S \subseteq C, |S| = k, |C| = n \quad (1)$$

where  $C$  denotes the parent collection. The size of the state space, i.e. the number of different  $k$ -membered subsets of an  $n$ -membered set, is given by the binomial coefficient:

$$\binom{n}{k} = \frac{n!}{(n-k)!k!} \quad (2)$$

For example, the selection of 50 from a set of 1000 plates involves  $10^{85}$  possibilities.

*Full Arrays* – A full array is applicable only to combinatorial libraries and refers to a subset of compounds that represent all the products derived by combining a given subset of building blocks in all possible combinations as prescribed by the reaction scheme. Note that in this context, the term “full array” does not necessarily refer to the physical layout and execution of the experiment. A full array for a hypothetical 2-component library is illustrated in Fig. 2b and is defined as:

$$S = S_1 \times S_2 \times \dots \times S_R, S_i \subseteq R_i, |S_i| = k_i, |R_i| = n_i \quad (3)$$

where  $R_i$  represents the pool of reagents at the  $i$ -th variation site of the combinatorial library, and  $n_i$  and  $k_i$  are the total and selected number of reagents at that site, respectively.

The combinatorics of full and sparse arrays are vastly different. For sparse arrays, the number of states that one has to consider (the number of different  $k$ -subsets of an  $n$ -set) is given by Eq. 2, whereas for full arrays the number of possibilities (i.e. the number of different  $k_1 \times k_2 \times \dots \times k_R$  arrays derived from an  $n_1 \times n_2 \times \dots \times n_R$   $R$ -component combinatorial library) is:

$$\prod_{i=1}^R \frac{n_i!}{(n_i - k_i)!k_i!} \quad (4)$$

For a  $10 \times 10$  2-component combinatorial library there are  $10^{25}$  different subsets of 25 compounds, and only 63,504 different  $5 \times 5$  arrays. For a  $10 \times 10$  selection from a  $100 \times 100$  library, those numbers increase to  $10^{241}$  and  $10^{26}$  for sparse and full arrays, respectively. Thus, full arrays are much simpler to design and easier to synthesize on robotic hardware. In fact, most combinatorial libraries are synthesized in (full) array format.

*Plates* – The final type of selection is a collection of plates. A plate represents a set of compounds grouped together according to their physical location or some other common characteristic. The primary use of this class is to enable the selection of plates from a large corporate archive for replication and screening against one or more biological targets. A plate subset is defined as:

$$S \subseteq G, |S| = k_p, |G| = n_p \quad (5)$$

where  $G$  represents the collection of plates, and  $n_p$  and  $k_p$  are the total and selected number of plates in the collection, respectively. Note that the number of compounds in the plates need not be the same; in this case the selection criteria must be defined in a manner that does not favor the selection of plates that contain either too few or too many compounds, unless it is so desired.

**Criteria** – A selection criterion is a function that encodes the ability of a given set of compounds to satisfy a particular design objective such as maximum intrinsic diversity, similarity to a set of known actives, enforcement of particular property distributions, and many others. These functions are simple numerical indices that can be combined into a single objective function that measures the overall quality of a candidate design,  $S$ :

$$f_o(S) = f(f_1(S), f_2(S), \dots, f_n(S)) \quad (6)$$

The objective function  $f_o$  can assume any desired functional form, and by convention it is maximized to produce the optimal set. The following paragraphs describe the selection criteria that are most commonly used in library design. Note that many of these functions can be defined in a multitude of ways, and some alternative definitions are discussed in section IV.

**Similarity** – The similarity of a given set of compounds,  $S$ , to a set of leads is defined as a function of the average distance of a compound to its nearest lead:

$$S(S) = \frac{1}{k} \sum_{i=1}^k f(\min_{j=1}^l (d_{ij})) \quad (7)$$

where  $k$  is the cardinality of  $S$ ,  $l$  is the number of leads,  $d_{ij}$  is the distance between the  $i$ -th compound and the  $j$ -th lead in some molecular descriptor space, and  $f$  is a user-defined function known as a *kernel*. The default kernel is the identity function. Since typically a higher similarity score indicates a collection of compounds that are more distant and therefore less similar to the leads, focused libraries are obtained by minimizing  $S$ .

*Diversity* – The intrinsic diversity of a set of compounds,  $S$ , is defined as a function of the average nearest neighbor distance [43, 44, 45]:

$$D(S) = \frac{1}{k} \sum_i f(\min_{j \neq i} (d_{ij})) \quad (8)$$

where again  $k$  is the cardinality of  $S$ ,  $d_{ij}$  is the Euclidean distance between the  $i$ -th and  $j$ -th compounds in  $S$  in some molecular descriptor space, and  $f$  is a user-defined kernel. The kernel is used to tailor the diversity of the design, and defaults to the identity function. Since typically the value of this function increases with spread, diverse libraries are obtained by maximizing  $D$ . We have found that this function is smoother than the more commonly used minimum dissimilarity and can discriminate more effectively between the various ensembles. Naively implemented, Eq. 8 requires  $k(k-1)/2$  distance computations and scales adversely with the number of compounds selected. To reduce the quadratic complexity of the problem,  $D$  is computed using the k-d tree algorithm presented in [12]. This algorithm achieves computational efficiency by first organizing all the points in  $C$  in a k-dimensional tree, and then performing a nearest neighbor search for each point using a branch-and-bound approach. For a relatively small number of dimensions, this algorithm exhibits  $k \log k$  time complexity, and scales favorably with the number of compounds selected. Other diversity functions are discussed in section IV.

*Complementarity* – This criterion is closely related to diversity and represents the ability of a particular design to fill in the “diversity voids” that exist in a preexisting collection. Its definition is similar to Eq. 8:

$$D(S, S^*) = D(S \cup S^*) = \frac{1}{k} \sum_i f(\min_{j \neq i} (d_{ij})) \quad (9)$$

where  $i$  and  $j$  are now used to index the compounds in the combined set  $S \cup S^*$ , and  $k = |S| + |S^*|$ . Just as with Eq. 8, complementary designs are obtained by maximizing  $D$ . Complementarity represents the sum of the intrinsic and extrinsic diversity of a set of compounds with respect to a reference collection.

*Confinement* – This criterion measures the degree to which the properties of a given set of compounds fit within prescribed limits, and is defined as:

$$P(S) = \frac{1}{k} \sum_i \sum_j \max(x_j^{\min} - x_{ij}, x_{ij} - x_j^{\max}, 0) \quad (10)$$

where  $x_{ij}$  is the  $j$ -th property of the  $i$ -th compound, and  $x_j^{\min}$  and  $x_j^{\max}$  are the minimum and maximum allowed limits of the  $j$ -th property, respectively. Since the value of this function increases as more and more compounds fall outside the desired property range, constrained libraries are obtained by minimizing  $P$ . When multiple properties are used, they must be normalized to allow meaningful comparisons. In the special case when the properties of interest need to attain a particular target value (i.e. in the case of a degenerate range), Eq. 10 can be rewritten as:

$$P(S) = \frac{1}{k} \sum_i \sum_j \text{abs}(x_{ij} - x_j^*) \quad (11)$$

where  $x_j^*$  represents the target value of the  $j$ -th property. The confinement criterion was introduced in reference [41].

*Distribution* – This criterion is used to create designs that obey certain predefined molecular property distributions. It is based on the Kolmogorov-Smirnov statistic, which measures how well an experimental distribution is approximated by a particular distribution function. It is applicable to unbinned distributions that are functions of a single independent variable, and is defined as the maximum value of the absolute difference between two cumulative distribution functions:

$$K^* = \max_{-\infty < x < \infty} |P(x) - P^*(x)| \quad (12)$$

where  $P(x)$  is an estimator of the cumulative distribution function of the actual probability distribution from which it is drawn, and  $P^*(x)$  is a known cumulative distribution function. For a set of  $k$  points  $x_i$ ,  $i = 1, \dots, k$ ,  $P(x)$  represents the fraction of data points to the left of a given value  $x$  (inclusive). The method is illustrated in Fig. 3.

Unlike the more commonly used  $\chi^2$  test, the Kolmogorov-Smirnov statistic does not require binning of the data, which is arbitrary and leads to loss of information. More importantly, the function is very fast to compute since it involves sorting the data in ascending order, followed by a linear scan to identify the maximum difference from the user-defined cumulative distribution function (or a simultaneous scan of two vectors in the case of two cumulative distributions). Speed of computation is particularly important in the

application at hand where the fitness function needs to be evaluated tens of thousands of times in the course of the optimization.

The significance level of a particular value of  $K^*$  as a disproof of the hypothesis that two distributions are the same is a function of  $K^*$  and the number of data points,  $k$ . This function is relatively slow to compute, but when  $k$  is constant, it is a monotonic function of  $K^*$ . Since all we want is to determine which experimental distribution is closer to the “ideal” distribution  $P^*(x)$ , the significance level need not be computed.

The Kolmogorov-Smirnov criterion as defined by Eq. 12 is a measure of dissimilarity and takes values in the interval  $[0, 1]$ . Alternatively, we can define the similarity between two probability distributions,  $K$ :

$$K = 1 - K^* \quad (13)$$

Thus, designs that obey a particular distribution function are obtained by maximizing  $K$ . This criterion was introduced by Rassokhin and Agrafiotis in reference [19].

*Activity* – A common goal in library design is to produce arrays of compounds that are predicted to be highly active against a predefined target according to some quantitative structure-activity or receptor binding model. This can be easily accomplished using the average predicted activity of the selected compounds,  $S$ :

$$Q_a(S) = \frac{1}{k} \sum_i a_i \quad (14)$$

where  $a_i$  is some measure of the predicted activity of the  $i$ -th compound in  $S$ . Since the value of  $Q_a$  increases as the compounds become more active, focused libraries are obtained by maximizing  $Q_a$ .

*Selectivity* – A similar function to Eq. 14 can be used to measure the selectivity against a set of biological targets:

$$Q_s(S) = \frac{1}{k} \sum_i (a_{iq} - \max_{j \neq q} (a_{ij})) \quad (15)$$

where  $a_{ij}$  is the predicted activity of the  $i$ -th compound against the  $j$ -th target, and  $q$  is the target that the molecules should be selective against. Since the value of  $Q_s$  increases as the compounds become more selective, selective libraries are obtained by maximizing  $Q_s$ .

**Overlap** – This criterion measures the extend of overlap of a given selection,  $S$ , to another set of compounds,  $S^*$ :

$$O(S, S^*) = 1 - \frac{|S \cap S^*|}{|S|} \quad (16)$$

where  $|S|$  denotes the cardinality of  $S$ . Since the value of  $O$  decreases with the number of duplicates, non-redundant libraries are obtained by maximizing  $O$ . The most typical use of this criterion is to enable the selection of compounds in array or plate format that have not been previously synthesized or screened, and for the selection of multiple arrays from the same library (see below).

**Optimization** – As previously mentioned, the combinatorial nature of the problem does not permit an exhaustive enumeration of every possible combination in order to identify the optimal solution. While for simple cost functions several very effective greedy algorithms can be employed [41, 42], arbitrary objective functions have unpredictable surfaces with many local minima, and require a stochastic approach that is capable of global optimization. The method chosen here is based on simulated annealing. Simulated annealing is a global, multivariate optimization technique based on the Metropolis Monte-Carlo search algorithm. The method starts from an initial random state, and walks through the state space associated with the problem of interest by a series of small, stochastic steps. In the problem at hand, a state represents a particular selection of compounds (i.e. a list of subsets from one or more virtual collections) and a step is a small change in the composition of that set (i.e. replacement of a small fraction of the compounds comprising the set). An objective function,  $f_o$ , maps each state to a real value which represents its energy or fitness. While downhill transitions are always accepted, uphill transitions are accepted with a probability that is inversely proportional to the energy difference between the two states. This probability is computed using Metropolis' acceptance criterion:

$$p = e^{-\Delta E / K_B T} \quad (17)$$

or Felsentein's function:

$$P = \frac{1}{1 + e^{-\Delta E / K_B T}} \quad (18)$$

The later ensures that the transition probability never exceeds 0.5, and thus prohibits the system from performing random walks. Boltzmann's constant,  $K_B$ , is used for scaling purposes, and  $T$  is an artificial temperature factor that controls the ability of the system to overcome energy barriers. The temperature is systematically adjusted during the course of the simulation in a manner that gradually reduces the probability of high-energy transitions. This protocol results in two optimization phases: one in which the system explores the state space relatively freely, and one in which it equilibrates around a low energy minimum.

For simulated annealing to work, it is imperative that the transition probability is properly controlled. The difficulty with the problem at hand is that the cost function is not known *a priori* and can vary dramatically from problem to problem depending on the nature and weights of the criteria involved. To circumvent the problem of selecting an appropriate value for Boltzmann's constant, we use an adaptive approach in which  $K_B$  is not a true constant, but rather it is continuously adjusted during the course of the simulation based on a running estimate of the mean transition energy. In particular, at the end of each transition, the mean transition energy is updated, and the value of  $K_B$  is adjusted so that the acceptance probability for a mean uphill transition at the final temperature is some predefined small number (usually 0.1%). The temperature is reduced using a Gaussian cooling schedule with a half-width of 5-10 deviation units. Other cooling schedules, such as linear, exponential, and Cauchy, have also been tested; in general, schedules that involve more extensive sampling at lower temperatures seem to perform best, although it is equally important that sufficient time is spent at higher temperatures so that the algorithm does not get trapped into local minima. The following sections provide a formal description of the serial and parallel implementations of this algorithm.

*Serial Implementation* – Let  $S_i$  denote the  $i$ -th subset requested,  $C_i$  the collection from which it is drawn (or any subset thereof),  $n_S$  the total number of subsets in the design,  $\mathbf{T}$  the vector of temperatures in the cooling schedule,  $n_T$  the number of temperature cycles,  $n_C$  the number of sampling steps per temperature

cycle,  $f_o(\cdot)$  the multiobjective fitness function, and  $\overline{\Delta E}$  the average uphill transition energy (fitness). Also, let  $R_{ij}$  be the pool of available building blocks at the  $j$ -th variation site if  $C_i$  is a combinatorial library,  $S_{ij}$  the selected building blocks at the  $j$ -th variation site if  $S_i$  is a full array, and  $P_i$  and  $n_p$  be the set and number of available plates in  $C_i$  (if applicable). The serial annealing algorithm involves the following steps:

1. Initialize each subset,  $S_i$ ,  $i = 1, \dots, n_S$ , at random, and set  $S = \{S_i, i = 1, \dots, n_S\}$ . If  $S_i$  is a sparse array and  $k$  is the requested number of compounds, initialize  $S_i$  with a random subset of  $k$  compounds from  $C_i$ . If  $S_i$  is a full array and  $k_j$  is the requested array size at the  $j$ -th variation site, initialize each  $S_{ij}$  with a random subset of  $k_j$  reagents from  $R_{ij}$ . If  $S_i$  is a plate selection and  $k$  is the number of plates requested, initialize  $S_i$  with a random subset of  $k$  plates from  $P_i$ .
2. Set  $f = f_o(S)$ ,  $f_{min} = f_o(S)$  and  $\overline{\Delta E} = 0$ .
3. Perform steps 4-9 for each  $t \leq n_T$ .
4. Set  $T = \mathbf{T}[t]$ .
5. Perform steps 6-9 for each  $c \leq n_C$ .
6. Select a random subset  $S_i$ ,  $i \in [1, n_S]$ , mutate it, and denote the resulting state as  $S^*$ . If  $S_i$  is a sparse array, mutate it by replacing a small fraction of randomly chosen compounds in  $S_i$  with an equal number of randomly chosen compounds in  $\overline{S_i}$ , where  $\overline{S_i}$  denotes the complement of  $S_i$  (i.e. the compounds in  $C_i$  that are not in  $S_i$ ). If  $S_i$  is a full array, mutate it by selecting a variation site,  $j$ , at random, and replacing a randomly chosen reagent in  $S_{ij}$  with a randomly chosen reagent in  $\overline{S_{ij}}$ . Finally, if  $S_i$  is a plate subset, mutate it by replacing a randomly chosen plate in  $S_i$  with a randomly chosen plate from  $\overline{S_i}$ .
7. Set  $f^* = f_o(S^*)$  and  $\Delta E = |f - f^*|$ .
8. Update  $\overline{\Delta E}$  and set  $K_B = \frac{-\overline{\Delta E}}{T \ln a}$ , where  $a$  is a predefined small number in the interval  $[0, 1]$ .
9. If  $f^* \leq f$  or if  $f^* > f$  and  $r < e^{-\Delta E / K_B T}$ , where  $r$  is a random number in the interval  $[0, 1]$ , then:
  - 9.1. Set  $S = S^*$  and  $f = f^*$ .
  - 9.2. If  $f < f_{min}$ , set  $f_{min} = f$  and  $S_{min} = S$ .
10. Output  $S_{min}$  and  $f_{min}$ .

To avoid repeated and potentially expensive memory allocation,  $S$  and  $S^*$  are implemented as a single object that is able to revert to its previous state after an unfavorable mutation. Moreover, the mean uphill transition energy  $\overline{\Delta E}$  is computed based on the last  $w$  uphill transitions, where  $w$  is a preset window (usually a few hundred steps).

*Parallel Implementation* – The parallel algorithm described in this section is known as synchronous annealing and is designed to keep inter-process communication and thread synchronization to a minimum. As in conventional annealing, the process starts with a random initial state and walks through the state space by a series of small stochastic steps. However, during each temperature cycle each execution thread is allowed to follow its own independent Monte-Carlo trajectory. The threads synchronize at the end of each cycle, and the best among the last states visited by each thread is recorded and used as the starting point for the next iteration. Given sufficient simulation time, this parallel algorithm produces results that are comparable to those obtained with the traditional serial implementation (see below). The algorithm proceeds as follows:

1. Initialize each subset,  $S_i$ ,  $i = 1, \dots, n_S$ , at random, and set  $S = \{S_i, i = 1, \dots, n_S\}$ . If  $S_i$  is a sparse array and  $k$  is the requested number of compounds, initialize  $S_i$  with a random subset of  $k$  compounds from  $C_i$ . If  $S_i$  is a full array and  $k_j$  is the requested array size at the  $j$ -th variation site, initialize each  $S_{ij}$  with a random subset of  $k_j$  reagents from  $R_{ij}$ . If  $S_i$  is a plate selection and  $k$  is the number of plates requested, initialize  $S_i$  with a random subset of  $k$  plates from  $P_i$ .
2. Set  $f = f_o(S)$ ,  $f_{min} = f_o(S)$  and  $\overline{\Delta E} = 0$ .
3. Perform steps 4-12 for each  $t \leq n_T$ .
4. Set  $T = \mathbf{T}[t]$ .
5. Perform steps 6-11 for each  $p \leq n_p$ , where  $n_p$  is the number of execution threads (processors).
6. Set  $S_p = S$ ,  $f_p = f$ ,  $S_{min}^p = S$ ,  $f_{min}^p = f$ , and  $\overline{\Delta E}_p = \overline{\Delta E}$ .
7. Perform steps 8-11 for each  $c \leq n_C$ .

8. Select a random subset  $S_i$ ,  $i \in [1, n_s]$ , mutate it, and denote the resulting state as  $S^*$ . If  $S_i$  is a sparse array, mutate it by replacing a small fraction of randomly chosen compounds in  $S_i$  with an equal number of randomly chosen compounds in  $\overline{S_i}$ . If  $S_i$  is a full array, mutate it by selecting a variation site,  $j$ , at random, and replacing a randomly chosen reagent in  $S_{ij}$  with a randomly chosen reagent in  $\overline{S_{ij}}$ . Finally, if  $S_i$  is a plate subset, mutate it by replacing a randomly chosen plate in  $S_i$  with a randomly chosen plate from  $\overline{S_i}$ .
9. Set  $f^* = f_o(S^*)$  and  $\Delta E = |f_p - f^*|$ .
10. Update  $\overline{\Delta E}_p$  and set  $K_B^p = \frac{-\overline{\Delta E}_p}{T \ln a}$ , where  $a$  is a predefined small number in the interval  $[0, 1]$ .
11. If  $f^* \leq f_p$  or if  $f^* > f_p$  and  $r < e^{-\Delta E / K_B^p T}$ , where  $r$  is a random number in the interval  $[0, 1]$ , then:
  - 11.1. Set  $S_p = S^*$  and  $f_p = f^*$ .
  - 11.2. If  $f_p < f_{min}^p$ , set  $f_{min}^p = f_p$  and  $S_{min}^p = S_p$ .
12. Set  $f_{min} = \min_p f_{min}^p$ ,  $S_{min} = S_{min}^q : S_{min}^q \leq S_{min}^p \forall p \neq q$ ,  $f = \min_p f_p$  and  $S = S_q : S_q \leq S_p \forall p \neq q$ .
13. Output  $S_{min}$  and  $f_{min}$ .

The choice of simulated annealing was based on its programmatic simplicity, the fact that the mutation function (or step in annealing terminology) can be designed in a way that guarantees the creation of valid states (something that requires extra care with genetic approaches), and in-house comparative studies which demonstrated superior convergence compared to evolutionary approaches.

**Filters** – Filters limit the selection to specific subsets of a particular collection. These subsets can be specified as reagent lists, product lists, or plate lists. For array selections, product lists are deconvoluted to the respective reagents (note that in this case, if the input list is a sparse array, it is possible that some of the products in the final selection may not be part of the specified list). Filters can be used for a variety of purposes, such as restricting the selection to compounds having a particular substructure, reagents that are

provided by reliable vendors, plates that have a good QC score or contain compounds of a particular structural class, etc.

**Computational Details** – All programs were implemented in the C++ programming language and are part of the DirectedDiversity® [5] software suite. They are based on 3-Dimensional Pharmaceuticals' Mt++ class library [47] and are designed to run on all Posix-compliant Unix and Windows platforms. Parallel execution on systems with multiple CPUs is supported through the multithreading classes of Mt++. All calculations were carried out on a Dell workstation equipped with two 800 MHz Pentium III Intel processors running Windows 2000 Professional. Selections were carried out in 30 temperature cycles using a Gaussian cooling schedule and 1,000 sampling steps per temperature cycle. Boltzmann's constant was determined in an adaptive manner so that the acceptance probability for a mean uphill transition at the final temperature was 0.1%.

### **III. DATA SETS**

Two data sets were used in this study. The first is a 2-component virtual combinatorial library based on the reductive amination reaction, and the second is a subset of 3-Dimensional Pharmaceuticals' probe library, a collection of more than 250,000 diverse compounds that represent over 30 different structural classes.

The 2-component reductive amination library is part of a synthetic strategy that exploits the pivotal imine intermediate and is utilized for the construction of structurally diverse drug-like molecules with useful pharmacological properties, particularly in the GPCR super-family [48]. The synthetic protocol is illustrated in Fig. 4. This library was used in a number of previous studies and represents an internal standard for testing new library design methodologies. A set of 300 primary and secondary amines with 300 aldehydes were selected at random from the Available Chemicals Directory [49], and were used to generate a virtual library of 90,000 products using the library enumeration classes of the DirectedDiversity® toolkit [47]. These classes take as input lists of reagents supplied in SDF or Smiles format, and a reaction scheme written in a proprietary language that is based on Smarts and an extension of

the scripting language Tcl. All chemically feasible transformations are supported, including multiple reactive functionalities, different stoichiometries, cleavage of protecting groups, stereo-specificity, and many others. The computational and storage requirements of the algorithm are minimal (even a billion-membered library can be generated in a few seconds on a personal computer) and scale linearly with the number of reagents.

Each compound in the 90,000-membered library was characterized by a standard set of 117 topological descriptors [50, 51] computed with the DirectedDiversity® toolkit [47]. These descriptors include an established set of topological indices with a long, successful history in structure-activity correlation such as molecular connectivity indices, kappa shape indices, subgraph counts, information-theoretic indices, Bonchev-Trinajstis indices, and topological state indices. It has previously been shown that these descriptors exhibit proper “neighborhood behavior” [24] and are thus well suited for diversity analysis and similarity searching [22, 23, 25, 52].

These 117 molecular descriptors were subsequently normalized and decorrelated using principal component analysis. This process resulted in an orthogonal set of 23 latent variables, which accounted for 99% of the total variance in the data. To simplify the analysis and interpretation of results, this 23-dimensional data set was further reduced to 2 dimensions using a very fast nonlinear mapping algorithm developed by our group [53, 54, 55, 56]. The projection was carried out in such a way that the pair-wise distances between points in the 23-dimensional principal component space were preserved as much as possible on the 2-dimensional map. The resulting map had a Kruskal stress of 0.187 and was used to perform and visualize the selections. The PCA preprocessing step was necessary in order to eliminate duplication and redundancy in the data, which is typical of graph-theoretic descriptors.

Finally, in addition to the 117 topological descriptors, the octanol-water partition coefficient ( $\log P$ ) of each compound was computed independently using the Ghose-Crippen approach [57] as implemented in the DirectedDiversity® toolkit [47], and was used as the target variable for property-based designs (see below). This parameter was not included in the descriptor set used for similarity and diversity assessment.

The second data set is a subset 3-Dimensional Pharmaceuticals’ probe library containing nearly 250,000 compounds arranged in 96-well plate format containing 88 compounds each plate. This library was described using a similar methodology to that used for the reductive amination library.

#### IV. RESULTS AND DISCUSSION

**Diverse Libraries** – Molecular diversity represents the most common method for designing combinatorial libraries [43, 44, 45]. Although there has been much controversy regarding the choice of metrics and descriptors and their ability to discover novel leads, it is generally accepted that a proper diversity function should measure spread with be algorithmically efficient so that it can be applied to the kinds of data sets that are encountered in combinatorial library design.

In general, diversity metrics fall into three main categories: *distance-based* methods which express diversity as a function of the pairwise molecular dissimilarities defined through measurement or computation; *cell-based* methods which define it in terms of the occupancy of a finite number of cells that represent disjoint regions of chemical space; and *variance-based* methods which quantify diversity based on the degree of correlation between the molecules' pertinent features. In their vast majority, these metrics encode the ability of a given set of compounds to sample chemical space in an even and unbiased manner, and are used to produce space-filling designs that minimize the size of unexplored regions known as “diversity voids”. A thorough discussion of the relative merits of the most commonly used diversity functions can be found elsewhere [43, 44, 45].

The selection of 100 compounds from the reductive amination library based on their average nearest neighbor Euclidean distance on the 2-dimensional nonlinear map (Eq. 8) is shown in Fig. 5a. The compounds, which were selected as singles, are nicely distributed in the data space occupied by the virtual library, and do not exhibit any significant clustering that reflects the density distribution of the parent collection. However, this selection does suffer from one significant drawback: it requires 69 amines and 70 aldehydes, i.e. a total 139 reagents. Obviously, the physical synthesis of these compounds would be extremely laborious, and this is the primary reason why sparse arrays are rarely used in combinatorial chemistry. This method is usually employed in compound retrieval and acquisition, particularly when automated, efficient cherry-picking techniques are available. Thus, the remaining discussion will be focused mostly on (full) array selections.

Fig. 5b shows the selection of an equivalent number of compounds in the form of a  $10 \times 10$  array using the same distance function and diversity metric that was employed in the previous selection. Although the array is somewhat less diverse in terms of spread, it requires only 20 reagents as compared to 139 reagents required by the singles, and therefore it is much easier to synthesize in practice. A look at the selected reagents (Fig. 6) confirms that the design is also diverse in terms of chemical structure, as it consists of building blocks containing a wide variety of atom types, connectivity patterns, ring systems and functional groups.

Some interesting aspects of the optimization algorithm are illustrated in Fig. 7-9. Fig. 7 and 8 show the diversity score and the percent of accepted uphill transitions at the end of each temperature cycle for a single annealing run. In Fig. 7, green points indicate the diversity of the last accepted state, and red points the cost of the best state discovered at the end of each temperature cycle. Within the first 15 cycles, the algorithm is able to extract the gross features of the minimum and recover most of the diversity that is accessible by this array size. The final cycles are spent refining that minimum with relatively minimal improvements in the fitness function. The asymptotic convergence of the two curves manifests the decreasing ability of the Metropolis search algorithm to perform high-energy transitions, which is also reflected in the fraction of accepted uphill transitions in Fig. 8. Indeed, at higher temperatures the algorithm is able to overcome substantial energy barriers, but this ability is diminished at lower temperatures, and the system is eventually frozen around the 20-th cycle. This graph also shows that the adaptive determination of  $K_B$  seems to work reasonably well. This parameter need not be estimated by the user; rather, it is adjusted automatically by the algorithm as the simulation progresses and as the energy landscape is more thoroughly explored.

The method is also very robust. Fig. 9 shows the mean and standard deviation of the diversity scores obtained by 50 independent optimization runs carried out with the two algorithms, plotted against the diversity of a random array as a reference. The average score obtained with the serial and parallel implementation was 0.436 and 0.433, respectively, with a standard deviation of only 0.007 in both cases. Thus, even though the algorithm does converge to different local minima depending on the starting configuration and random seed, the solutions are essentially equivalent in terms of quality, and more than sufficient for the purpose of exploring chemical space. As for the parallel algorithm, it scales favorably

with the number of CPUs (the elapsed execution time showed a nearly perfect linear relationship with respect to the number of processors used) and produces results that are virtually identical to those obtained with the serial algorithm.

A point that is worth noting relates to the computational advantages afforded by the use of the k-d tree algorithm for computing the diversity function in Eq. 8 [12]. Fig. 10 shows the ratio between CPU time required by the naïve and k-d tree algorithms in 2 dimensions as a function of the array size (the naïve algorithm involves exhaustive enumeration of all pairwise distances in the data sample). It is clear that while for very low values of  $k$ , k-d trees impose some trivial overhead, they scale superbly with  $k$ , resulting in enormous computational savings for large selections.

This margin, however, decreases in higher dimensions. For example, the scaling factor (i.e. the increase in the CPU time when the selection size is doubled) is 2.04 for 2D, 2.31 for 3D, 2.38 for 4D and 2.98 for 10D, a trend shown graphically in Fig. 11. For low dimensions ( $\leq 5$ ), the scaling factor is less than 2.5, which is consistent with the reported theoretical value of  $n \log n$ . However, for 10 dimensions the time complexity is somewhere between linear and quadratic. This performance decline is likely to continue with increasing dimensionality and the algorithm will eventually become quadratic, with the added overhead of constructing and traversing the tree. Fortunately, our experience has shown that most high-dimensional descriptor spaces can be reduced to a relatively small number of dimensions with minimal distortion through the use of efficient nonlinear mapping techniques [54, 55, 56].

Nevertheless, there is a point where Eq. 8 does become prohibitively expensive, no matter what algorithm is used to compute the nearest neighbors. For large selections containing thousands of compounds, alternative metrics must be devised to enable a more expedient estimation of molecular diversity. Recently, we presented a novel diversity function that captures the notion of spread, is fast to compute, scales favorably with the number of compounds in the design, does not fall prey to dimensionality, and can be used to compare collections of different cardinality [58]. The method is based on the fundamental assumption that an optimally diverse sample is one that is uniformly distributed in the property space it is designed to explore. Diversity is quantified by estimating the cumulative probability distribution of inter-molecular dissimilarities in the collection of interest, and then measuring the deviation of that distribution from the respective distribution of a uniform sample using the Kolmogorov-Smirnov

statistic described in section II. Departure from uniformity induces sampling redundancy and formation of clusters, and thus results in diminishing diversity (see Fig. 12). The distinct advantage of this approach is that the cumulative distribution can be easily estimated using probability sampling and does not require exhaustive enumeration of all pairwise distances in the design, resulting in an algorithm of virtually constant time complexity. While some caution must be exercised in determining the appropriate target distribution [58], the function produces results that are consistent with our notion of spread and can do so at a fraction of the time required by alternative methodologies. The selection of a 10×10 array using this algorithm is shown in Fig. 13. While the selected compounds are not as perfectly distributed in the input space, they show no preference for any particular region of the map, nor are they biased by the local density of the parent collection, i.e. they are diverse.

A common criticism of the diversity function in Eq. 8 is that it tends to favor the extremes of the feature space and produce designs containing unusual compounds of limited pharmaceutical interest. Although in theory one can exclude ‘peculiar’ reagents from the selection through the use of appropriate filters (see section II above) or even from the virtual library itself, a ‘softer’ and more dynamic approach is to adjust the kernel function. A simple way to control the volume of the design without affecting its internal spread is to request that the mean nearest neighbor distance of the selected compounds is equal to some preset value. Consider, for example, the design in Fig. 5b. This score of this array is 0.43, and represents the maximum possible nearest neighbor intermolecular dissimilarity for any design carried out in a 10×10 format. Thus, by maximizing the function:

$$D(S) = \frac{1}{k} \sum_i -abs(a - \min_{j \neq i}(d_{ij})) \quad (19)$$

where  $a$  is set to 0.25, one can design an array that occupies roughly half of the volume of the original selection, yet it is still internally diverse (Fig. 14). The user can choose alternative values of  $a$  to produce designs with lesser or greater spread. Incidentally, Fig. 14 also illustrates a known disadvantage of Eq. 8 as a general measure of molecular diversity, and that is the fact that it measures only the intra-cluster (nearest neighbor) distances and does not take into account the inter-cluster separations. As a result, designs comprised of two or more internally diverse clusters cannot be easily distinguished, regardless of the relative separation between them.

The diversity function outlined above can also be used to design libraries that fill in the diversity voids of one or more pre-existing collections. A common problem faced by many pharmaceutical companies is the enhancement of their corporate collection through the synthesis or acquisition of chemical libraries that increase the diversity of their existing archives. Eq. 9 provides a simple and straightforward solution to this problem. By maximizing the diversity of the combined set, we can ensure that the design is both internally and externally diverse. The selection of a second 10×10 array that complements the diversity of the original selection in Fig. 5b is shown in Fig. 15. The two arrays have only 1 amine and 2 aldehydes (i.e. only 2 products) in common. Note that in terms of scatter the design is not ‘perfect’; this is due to the array constraint rather than the diversity metric itself.

***Drug-Like Libraries*** – In all the preceding examples, the selections were based exclusively on molecular diversity, and no attention was paid on the drug-likeness of the resulting compounds. Experience has shown that selecting compounds from a virtual library solely on that basis, no matter what diversity measure is used, often results in combinatorial libraries with poor pharmacokinetic properties or other undesirable characteristics. Recently there have been several attempts to quantify drug-likeness and incorporate it directly into the design process. Martin *et al.* [31] presented a reagent selection algorithm based on D-optimal design, wherein the candidate reagents were assigned to categorical bins according to their properties, and successive steps of D-optimal design were performed to generate diverse substituent sets consistent with required membership quotas from each bin. This technique was later elaborated [59, 60], and a new “parallel” sampling approach was proposed in order to eliminate the order-dependence of the original algorithm. Most recently, Koehler *et al.* [61] proposed a multipass algorithm designed to facilitate addition of compounds to an existing chemical library. This method was intended to prioritize compounds that are most similar to a specified set of favorable target molecules, and, at the same time, most dissimilar to the compounds that reside in the library being augmented. The algorithm scores and ranks each compound in the external library according to the local density of similar compounds in the target and internal libraries. Density arising from the target library adds a positive contribution to the score, while density arising from the internal library adds a negative contribution. The highest-ranking

compounds are then included in the internal library and the process is repeated until the specified number of compounds is selected.

Perhaps the most simple approach is to filter the candidate designs using Lipinski's 'rule of 5', a simple heuristic which states that for compounds which are not substrates of biological transporters poor absorption and permeation are more likely to occur when there are more than 5 H-bond donors, more than 10 H-bond acceptors, the molecular weight is greater than 500, or the logP is greater than 5 [62]. Consider, for example, a maximally diverse 20×20 array from the reductive amination library shown in Fig. 16. The molecular weight and logP distributions of these compounds are shown in Fig. 16 b-c, along with the respective distributions of known drugs derived by analyzing a subset of 7,484 compounds from the World Drug Index, which was generously provided to us in electronic format by Dr. Christopher Lipinski of Pfizer, Inc. While the distribution of molecular weights matches closely that of the reference set, the logP distribution is shifted upwards by nearly 3 logP units, with 125 out of 400 compounds (more than 30% of the selection) falling outside the boundaries of the Lipinski box.

In our multiobjective paradigm, correcting this adverse physicochemical profile can be easily accomplished in one of two different ways. The first is to combine the diversity criterion of Eq. 8 with the confinement criterion of Eq. 10 in order to penalize compounds that fall beyond the boundaries of the box. This can be achieved with an objective function of the form:

$$f_o = D - 5 \cdot (0.01 \cdot P(mw) + P(\log P)) \quad (20)$$

where  $D$  represents the diversity of the ensemble, and  $P(mw)$  and  $P(\log P)$  are the values of the confinement criterion for molecular weight and logP, respectively. As with most problems of this type, the most difficult task is to assign a meaningful set of coefficients used to weight the individual objectives. In our case, the coefficients were determined empirically based on the maximum value that each criterion could assume independently for any given array size. In the case at hand, the coefficients were chosen based on the relative scale of MW and logP, and the value of the mean nearest neighbor distance,  $D$ , of the maximally diverse 20x20 array. As shown in Fig. 17, the resulting design is fully contained within the Lipinski box, while remaining sufficiently diverse.

However, given the probabilistic nature of the problem, perhaps a better approach would be to compare the actual distributions themselves, and penalize designs whose MW and logP distributions depart from

those of known drugs. As a numerical measure of the dissimilarity between two property distributions, we used the Kolmogorov-Smirnov criterion of Eq. 13. Unlike the more commonly used  $\chi^2$  test, the Kolmogorov-Smirnov statistic does not require binning of the data [18], which is arbitrary and leads to loss of information. The distributions of known drugs were essentially normal, with a mean and sigma of 314.3 and 108.3 for molecular weight, and 1.04 and 1.78 for logP, respectively.

To enforce these distributions in the final design, we combined molecular diversity with the Kolmogorov-Smirnov statistic into the following objective function:

$$f_o = D + 0.2 \cdot K(\log P) + 0.2 \cdot K(mw) \quad (21)$$

where  $D$  is the diversity criterion defined in Eq. 14, and  $K(\log P)$  and  $K(mw)$  are the Kolmogorov-Smirnov similarities between the logP and molecular weight distributions of the selected compounds and the reference WDI set, respectively. Again, the coefficients were determined on the basis of the maximum values of the respective criteria, which were 0.18, 0.9 and 0.9 for diversity, and MW and logP distribution, respectively. These values suggested that for the three criteria to be placed on an equal footing, the value of  $K$  had to be scaled down by approximately a factor of 5. In pathological cases where the energy landscapes (i.e. the distributions of scores) of the individual criteria are very different, alternative, more complex objective functions can be devised. As shown in Fig. 18, when the selection is carried out using Eq. 21, the selected compounds' MW and logP distributions approximate very nicely the respective distributions of known drugs, and this occurs without a significant impact on the diversity of the design (Fig. 18a). Note that just like any problem of this kind, there is a point where the various objectives begin to oppose each other: a step towards improving one of the objectives, increasing molecular diversity, is a step away from improving the other, increasing their drug-likeness.

The methodology outlined above is, of course, not limited to virtual collections. A common problem faced by many pharmaceutical companies is the selection of a small number of compounds from their chemical banks for screening against a new biological assay. Compounds synthesized by combinatorial methods are usually stored as DMSO solutions in 96-well plates. Extraction of individual samples from these plates can be laborious and time-consuming, and it is often much more efficient to replicate and screen the entire plate containing the compounds of interest. This, however, is a rather simplistic approach; as we demonstrated before, the quality of a particular design is determined by the collective properties of

all of its constituent compounds, and the experiment can be vastly improved if the physical location or grouping of these compounds is taken into consideration during the design phase. This is illustrated in the selection of 10 plates from a subset of 3DP's probe library, a collection of ~250,000 diverse compounds representing more than 30 distinct structural classes. These compounds are stored in nearly 3,000 96-well plates, containing 88 compounds each. The selection of 10 plates containing 880 diverse, drug-like compounds using an objective function similar to the one employed in the previous example is illustrated in Fig. 19. Note that the requirement that these compounds exhibit a normal MW distribution centered around 315 D causes a noticeable shift of the selection to the right of the nonlinear map, which represents molecules of smaller MW (size is a dominant factor in determining molecular similarity and therefore has a significant impact on the shape of the nonlinear map).

***Focused Libraries*** – Molecular diversity is typically employed in the design of exploratory libraries for use in general screening. Once a promising hit has been identified and confirmed, subsequent iterations attempt to explore the structure-activity space around that compound and refine its pharmacological profile. This process involves two phases: the first is based on very sparse information and is guided primarily by molecular similarity, whereas the second involves the use of more rigorous structure-activity models which are derived from an active series with a broad dynamic range of biological activities. This task can be easily accomplished with the similarity, activity and selectivity criteria described in section II. For reasons of brevity, the remaining discussion will focus on molecular similarity; the extension of these principles to QSAR should be straightforward.

Consider, again, the 2-component amination library and a randomly chosen member of that library (Fig. 20) as a lead. The distribution of the compounds in the 10×10 array that is maximally similar to that compound is illustrated in Fig. 21, and the respective reagents in Fig. 22. These reagents consist predominantly of di-substituted anilines bearing a nitro substituent, and di- and tri-substituted benzaldehydes containing three oxygen atoms in their substituents. This result demonstrates not only the ability of the algorithm to produce highly focused designs, but also the fact that this particular set of descriptors is fully consistent with our general perception of molecular similarity.

Finally, we turn our attention to the issue of redundancy. Mining a virtual library is an iterative process that involves successive rounds of selection, synthesis, and testing. Since chemical synthesis and biological testing may be time-consuming and costly, it is imperative that the number of duplicates (i.e. compounds that have been previously screened) is kept to a minimum. While this is trivial in the case of sparse arrays, it becomes more problematic with full arrays, particularly when the requested size is relatively large. The problem stems from the fact that eliminating the possibility of duplicates often results in the design of mediocre experiments that are far removed from their principal objective. Thus, a more practical approach is to use duplication as a selection criterion whose influence can be tailored by the medicinal chemist on a case-by-case basis. Fig. 23 illustrates how this principle can be used to design a second focused array around the putative lead in Fig. 20, in a way that minimizes the overlap with the original selection (Fig. 21). The selection was carried out using the objective function:

$$f_o = O - S \quad (22)$$

where  $O$  is the overlap criterion in Eq. 16. As shown in Fig. 24, the two arrays share 7 aromatic amines, but have no aldehydes (and therefore no products) in common. To achieve this goal, one oxygen atom is sacrificed from the benzaldehyde group, but besides this change, the products remain closely related to the original lead. Note that in the absence of this criterion, the second selection would have been identical to the first.

## V. CONCLUSIONS

This paper presented an overview of a general and flexible methodology for compound selection and combinatorial series design. The method permits the selection of multiple subsets from multiple compound collections using multiple selection criteria, in sparse array, full array, and plate format. It is rooted on the principles of multiobjective optimization, and is best viewed as a Pareto optimal process seeking a consensus in which many objectives are balanced so that the improvement of any single objective will result in a negative impact on at least one other objective. The method can accommodate any desired selection criteria, and can be employed for the design of both exploratory and focused libraries. Although the method is stochastic in nature and requires the evaluation of a relatively large number of candidate

designs, typical selections are carried out in a few seconds to a few minutes on a modern personal computer.

## **ACKNOWLEDGMENTS**

The author is thankful to Drs. Victor S. Lobanov, Dmitrii N. Rassokhin, Sergei Izrailev, Edward Jaeger, Renee DesJarlais, Richard M. Soll and Roger Bone of 3-Dimensional Pharmaceuticals, Inc. for many useful suggestions, and Dr. Raymond F. Salemme for his insightful comments and support of this work.

## REFERENCES

- [1] Thompson, L. A., and Ellman, J. A., *Chem. Rev.*, **1996**, *96*, 555-600.
- [2] Pareto, V. *Manual of Political Economy*, **1906**, 106.
- [3] Eschenauer, H. A., Koski, J., and Osyczka, A. *Multicriteria Design Optimization: Procedures and Applications*, Springer-Verlag, New York, **1986**.
- [4] Martin, E. J., Blaney, J. M., Siani, M. A., Spellmeyer, D. C., Wong, A. K. and Moos, W. H., Measuring diversity: experimental design of combinatorial libraries for drug discovery. *J. Med. Chem.*, **1995**, *38*, 1431-1436.
- [5] Agrafiotis, D.K.; Bone, R.F.; Salemme, F.R.; Soll, R.M., United States Patents 5,463,564, **1995**; 5,574,656, **1996**; 5,684,711, **1997**; and 5,901,069, **1999**.
- [6] Sheridan, R.P., and Kearsley, S.K. Using a genetic algorithm to suggest combinatorial libraries. *J. Chem. Info. Comput. Sci.*, **1995**, *35*, 310-320
- [7] Weber, L., Wallbaum, S., Broger, C., and Gubernator, K. Optimization of the biological activity of combinatorial compound libraries by a genetic algorithm. *Angew. Chem. Int. Ed. Eng.*, **1995**, *34*, 2280-2282.
- [8] Singh, J., Ator, M. A., Jaeger, E. P., Allen, M. P., Whipple, D. A., Solowej, J. E., Chowdhary, S., and Treasurywala, A. M., Application of genetic algorithms to combinatorial synthesis: a computational approach for lead identification and lead optimization, *J. Am. Chem. Soc.*, **1996**, *118*, 1669-1676.
- [9] Agrafiotis, D. K., Stochastic algorithm for maximizing molecular diversity. *3-rd Electronic Computational Chemistry Conference*, <http://hackberry.chem.niu.edu/ECCC3/paper48>, **1996**.
- [10] Agrafiotis, D. K., Stochastic algorithms for maximizing molecular diversity. *J. Chem. Info. Comput. Sci.*, **1997**, *37*, 841-851.
- [11] Agrafiotis, D. K., On the use of information theory for assessing molecular diversity. *J. Chem. Inf. Comput. Sci.*, **1997**, *37(3)*, 576-580.
- [12] Agrafiotis, D. K., and Lobanov, V. S., An efficient implementatin of distance-based diversity metrics based on k-d trees. *J. Chem. Inf. Comput. Sci.*, **1999**, *39(1)*, 51-58.

- [13] Hassan, M., Bielawski, J. P., Hempel, J. C., and Waldman, M., Optimization and visualization of molecular diversity of combinatorial libraries. *Mol. Diversity*, **1996**, *2*, 64-74.
- [14] Waldman, M., Li, H., and Hassan, M., Novel algorithms for the optimization of molecular diversity of combinatorial libraries, *J. Mol. Graphics Mod.*, Special issue on combinatorial library design, Agrafiotis, D. K., and Martin, E., Eds., **2000**, *18(4-5)*, 412-426.
- [15] Good, A. C., and Lewis, R. A., New methodology for profiling combinatorial libraries and screening sets: cleaning up the design process with HARPCik, *J. Med. Chem.*, **1997**, *40*, 3926.
- [16] Zheng, W., Cho, S. J., and Tropsha, A., Rational combinatorial library design: 1) Focus-2D: a new approach to the design of targeted combinatorial chemical libraries, *J. Chem. Info. Comput. Sci.*, **1998**, *38*, 251.
- [17] Brown, R. D., and Martin, Y. C., Designing combinatorial library mixtures using genetic algorithms. *J. Med. Chem.*, **1997**, *40*, 2304-2313.
- [18] Gillet, V. J., Willet, P., Bradshaw, J., and Green, D. V. S., Selecting combinatorial libraries to optimize diversity and physical properties. *J. Chem. Info. Comput. Sci.*, **1999**, *39*, 169-177.
- [19] Rassokhin, D. N., and Agrafiotis, D. K., Kolmogorov-Smirnov statistic and its applications in library design. *J. Mol. Graphics Mod.*, **2000**, *18(4-5)*, 370-384.
- [20] Brown, R. D., Hassan, M., and Waldman, M., Combinatorial library design for diversity, cost efficiency and drug-like character, *J. Mol. Graphics Mod.*, **2000**, *18(4-5)*, 427-437.
- [21] Sheridan, R. P., SanFeliciano, S. G., and Kearsley, S. K., Designing targeted libraries with genetic algorithms, *J. Mol. Graphics Mod.*, **2000**, *18(4-5)*, 320-334.
- [22] Downs, G. M. and Willett, P., Similarity searching and clustering of chemical-structure databases using molecular property data. *J. Chem. Inf. Comput. Sci.*, **1994**, *34*, 1094-1102.
- [23] Brown, R. D., and Martin, Y. C., *J. Chem. Inf. Comput. Sci.*, **1996**, *36*, 572-584.
- [24] Patterson, D.E.; Cramer, R.D.; Ferguson, A.M.; Clark, R.D.; Weinberger, L.E., *J. Med. Chem.* **1996**, *39*, 3049-3059.
- [25] Brown, R. D., and Martin, Y. C., The information content of 2D and 3D structural descriptors relevant to ligand-receptor binding, *J. Chem. Inf. Comput. Sci.*, **1997**, *37*, 1-9.

- [26] Matter, H., Selecting optimally diverse compounds from structure databases: a validation study of two-dimensional and three-dimensional molecular descriptors, *J. Med. Chem.*, **1997**, *40*, 1219.
- [27] Martin, Y. C., Bures, M. G., and Brown, R. D., Validated descriptors for diversity measurements and optimization, *Pharm. Pharmacol. Commun.* **1998**, *4*, 147.
- [28] Gillet, V. J., Willett, P., Bradshaw, J., *J. Chem. Inf. Comput. Sci.*; **1997**, *37(4)*, 731-740.
- [29] Jamois, E. A., Hassan, M., and Waldman, M., *J. Chem. Inf. Comput. Sci.*, in press.
- [30] Polinsky, A.; Feinstein, R. D.; Shi, S.; Kuki, A. *Molecular Diversity and Combinatorial Chemistry*, eds. Chaiken, I. M.; Janda, K. D. American Chemical Society: Washington D.C., **1996**, 219-232.
- [31] Martin, E. J., Spellmeyer, D. C., Critchlow, R. E., Blaney, J., M., Does combinatorial chemistry obviate computer-aided drug design?, in *Reviews in Computational Chemistry*, Lipkowitz, K. B., Boyd, D. B., Eds., VCH, New York, **1997**, *10*, 75-100.
- [32] Willett, P. *Similarity and Clustering in Chemical Information Systems*. Research Studies Press: Letchworth, **1987**.
- [33] Taylor, R., *J. Chem. Inf. Comput. Sci.*, **1995**, *35*, 59-67.
- [34] Chapman, D., *J. Comput.-Aided Mol. Design*, **1996**, *10*, 501-512.
- [35] Cummins, D. J., Andrews, C. W., Bentley, J. A. and Cory, M., *J. Chem. Inf. Comput. Sci.*, **1996**, *36*, 750-763.
- [36] Pearlman, R. S., and Smith, R. S., *Perspectives Drug Discovery Design*, **1998**, *9*, 339-353.
- [37] Pickett, S., Mason, J. S., and McLay, I. M., *J. Chem. Inf. Comput. Sci.*, **1996**, *36*, 1214-1223.
- [38] Davies, E. K., and Briant, C. *Network Science*, **1995**, <http://www.awod.com/netsci/issues/>.
- [39] Shemetulskis, N. E., Weininger, D., Blankley, C. J., Yang, J. J., and Humblet, C., *J. Chem. Inf. Comput. Sci.*, **1996**, *36*, 862-871.
- [40] Boyd, S. M., Beverly, M., Norskov, L., and Hubbard, R. E., *J. Comput.-Aided Mol. Design*, **1995**, *9*, 417-424.
- [41] Agrafiotis, D. K., and Lobanov, V. S., Ultrafast algorithm for designing focused combinatorial arrays. *J. Chem. Info. Comput. Sci.*, **2000**, *40*, 1030-1038.

- [42] Stanton, R. V., et al, Combinatorial library design: maximizing model fitting compounds with matrix synthesis constraints, *J. Chem. Info. Comput. Sci.*, **2000**, *40*, 701-705.
- [43] Agrafiotis, D. K., The diversity of chemical libraries. In *The Encyclopedia of Computational Chemistry*, Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer III, H. F., and Schreiner, P. R., Eds., John Wiley & Sons, Chichester, **1998**, 742-761.
- [44] Agrafiotis, D. K., Myslik, J. C., and Salemme, F. R., Advances in diversity profiling and combinatorial series design. *Mol. Diversity*, **1999**, *4(1)*, 1-22.
- [45] Agrafiotis, D. K., Lobanov, V. S., Rassokhin, D. N., and Izrailev, S., The measurement of molecular diversity, in *Virtual screening of bioactive molecules*, Böhm, H.-J., and Schneider, G., Eds., Wiley-VCH, Weinheim, **2000**.
- [46] Lajiness, M. S., In *QSAR: Rational Approaches to the Design of Bioactive Compounds*, Silipo, C., and Vittoria, A., Eds., Elsevier, Amsterdam, **1991**, 201-204.
- [47] Copyright © 3-Dimensional Pharmaceuticals, Inc., 1994-2000.
- [48] Dhanao, D. S., Gupta, V., Sapienza, A., and Soll, R. M., Poster 26, American Chemical Society National Meeting, Anaheim, CA, **1999**.
- [49] MDL Information Systems, Inc., 140 Catalina Street, San Leandro, CA 94577.
- [50] Hall L.H.; Kier, L.B. The Molecular Connectivity Chi Indexes and Kappa Shape Indexes in Structure-Property Relations. In *Reviews of Computational Chemistry*, Boyd, D.B.; Lipkowitz, K.B., Ed.; VCH Publishers, **1991**; Chapter 9, 367-422.
- [51] Bonchev, D.; Trinajstić, N., *J. Chem. Phys.* **1977**, *67*, 4517-4533.
- [52] Lobanov, V. S., and Agrafiotis, D. K., Stochastic similarity selections from large combinatorial libraries. *J. Chem. Inf. Comput. Sci.*, **2000**, *40*, 460-470.
- [53] Agrafiotis, D. K., A new method for analyzing protein sequence relationships based on Sammon maps. *Protein Science*, **1997**, *6*, 287-293.
- [54] Agrafiotis, D. K., and Lobanov, V. S., Nonlinear mapping networks. *J. Chem. Info. Comput. Sci.*, in press.

- [55] Rassokhin, D. N., Lobanov, V. S., and Agrafiotis, D. K., Nonlinear mapping of massive data sets by fuzzy clustering and neural networks. *J. Comp. Chem.*, in press.
- [56] Agrafiotis, D. K., Rassokhin, D. N., and Lobanov, V. S., Multidimensional scaling of large molecular similarity tables. *J. Comp. Chem.*, in press.
- [57] Ghose, A. K., Viswanadhan, V. N., Wendoloski, J. J., *J. Phys. Chem. A*, **1998**, *102*, 3762-3772.
- [58] Agrafiotis, D. K., Constant time algorithm for estimating the diversity of large chemical libraries. *J. Chem. Info. Comput. Sci.*, in press.
- [59] Martin, E. J., and Critchlow, R. E., Beyond mere diversity: tailoring combinatorial libraries for drug discovery. *J. Comb. Chem.*, 1999, *1(1)*, 32-45.
- [60] Martin, E. J., and Wong, A., Sensitivity analysis and other improvements to tailored combinatorial library design. *J. Chem. Info. Comput. Sci.*, **2000**, *40(2)*, 215-220.
- [61] Koehler, R. T., Dixon, S. L., and Villar, O. H., LASSOO: a generalized directed diversity approach to the design and enrichment of chemical libraries. *J. Med. Chem.*, **1999**, *42*, 4695-4704.
- [62] Lipinski, C. A., Lombardo, F., Dominy, B. W., Feeny, P. J., *Advanced Drug Delivery Reviews*, **1997**, *23*, 3-25

## **CAPTIONS TO FIGURES**

1. General algorithm for multiobjective compound selection.
2. Combinatorial subsets: a) sparse arrays, b) full arrays.
3. Kolmogorov-Smirnov statistic for computing the difference between two cumulative distribution functions.
4. Synthetic sequence for the reductive amination library.
5. Selection of 100 compounds from the reductive amination library based on maximum diversity on the 2-dimensional nonlinear map: a) sparse array, b) full array (10×10).
6. Amine and aldehyde reagents comprising the selection in Fig. 5b.
7. Diversity score as a function of time for the selection in Fig. 5b. ‘Best’ and ‘Last’ represent the best and last state found at the end of each temperature cycle, respectively. Since simulated annealing allows uphill transitions, these two states are not always the same.
8. Percent of accepted uphill transitions at the end of each temperature cycle for the selection in Fig. 5b.
9. Mean and standard deviation of the diversity scores obtained by 50 independent selections of a maximally diverse 10×10 array from the reductive amination library, using the serial and parallel simulated annealing algorithms, plotted against the score of a randomly chosen array as a reference.
10. Ratio of the CPU time required by the naïve over the k-d tree algorithm for the selection of a maximally diverse array from the reductive amination library as a function of array size.
11. Scaling of the k-d tree algorithm as a function of dimensionality. The y axis corresponds to the ratio of the CPU time required for the selection of 1,600 vs 800 compounds in array format. A value of 2 indicates  $O(n)$  complexity, whereas a value of 4 indicates an  $O(n^2)$  complexity.
12. Effect of clustering on the shape of the cumulative probability distribution of pairwise distance for an artificial 2-dimensional data set confined in the unit square. Reproduced from reference [58].
13. Selection of a maximally diverse 10×10 array from the reductive amination library based on the Kolmogorov-Smirnov measure of molecular diversity.

14. Selection of 100 compounds from the reductive amination library in 10×10 array format based on the Eq. 19. The kernel is used to limit the volume of the design, and avoid sampling the extreme regions of the feature space.
15. Selection of a 10×10 array that complements the diversity of the selection in Fig. 5b.
16. Selection of a maximally diverse 20×20 array from the reductive amination library: a) nonlinear map, b) MW distribution, and c) logP distribution. In b and c, the blue series shows the respective distributions of known drugs. Reproduced from the J. Mol. Graphics. Mod. [].
17. Selection of a maximally diverse 20×20 array that is confined within the boundaries of the Lipinski box: a) nonlinear map, b) Lipinski profile.
18. Selection of a maximally diverse 20×20 array that exhibits a drug-like MW and logP distribution: a) nonlinear map, b) MW distribution, and c) logP distribution. In b and c, the blue series shows the respective distributions of known drugs. Reproduced from the J. Mol. Graphics. Mod. [].
19. Selection of a maximally diverse set of 880 compounds from 3DP's probe library. The selection was designed to exhibit a drug-like MW and logP distribution, and represents the best set of 10 8×11 plates. a) nonlinear map, b) MW distribution, and c) logP distribution. In b and c, the blue series shows the respective distributions of known drugs.
20. Lead compound used for similarity selections from the reductive amination library.
21. Selection of 100 compounds in 10×10 array format based on maximum similarity to the 'lead' compound in Fig. 20.
22. Amine and aldehyde reagents comprising the selection in Fig. 21.
23. Selection of 100 compounds in 10×10 array format based on maximum similarity to the 'lead' compound in Fig. 20, and minimum overlap with the selection in Fig. 21.
24. Amine and aldehyde reagents comprising the selection in Fig. 23.

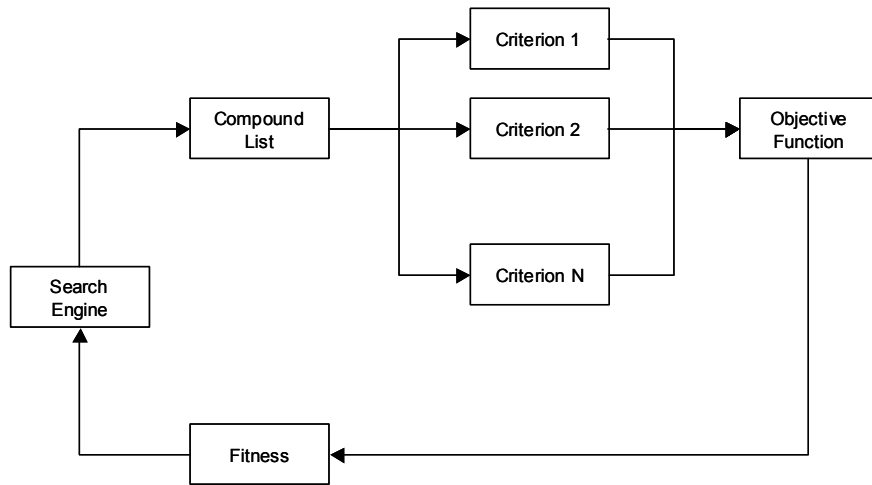


Fig. 1

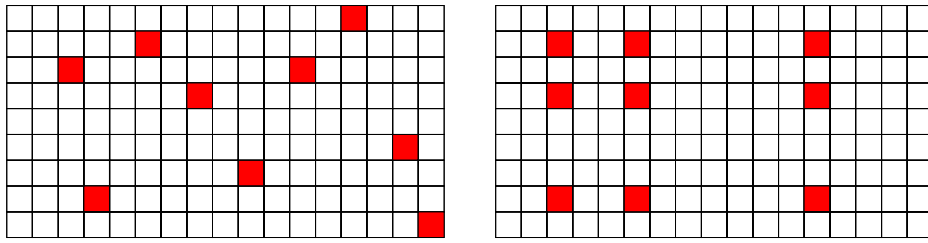


Fig. 2

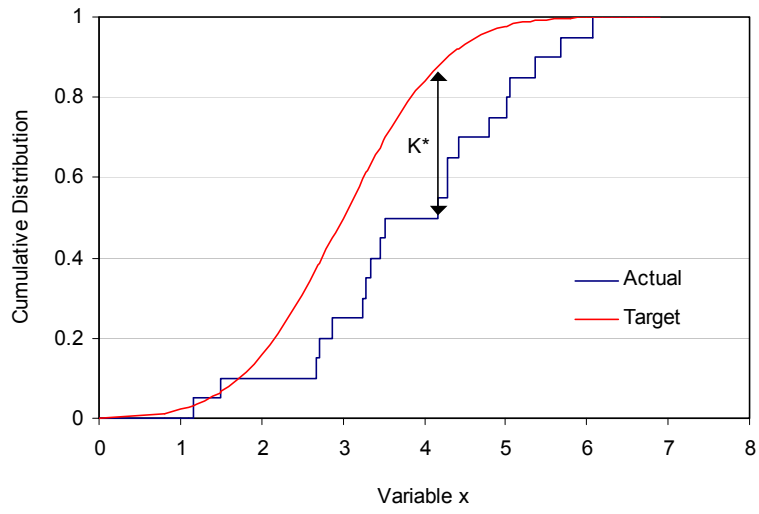


Fig. 3

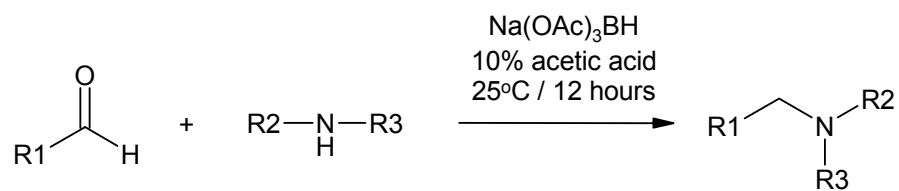
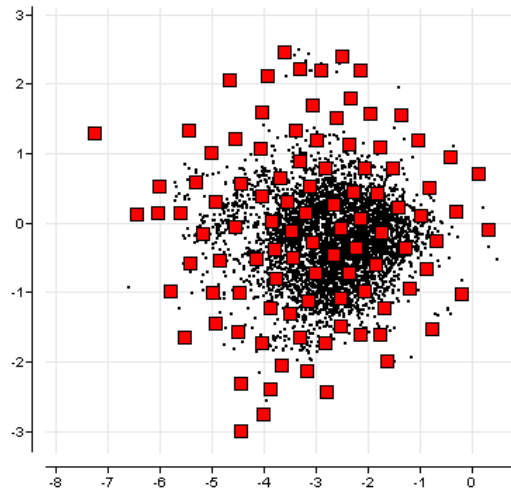
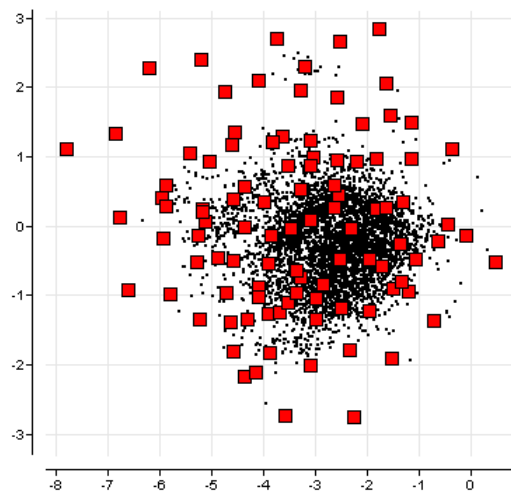


Fig. 4



a



b

Fig. 5

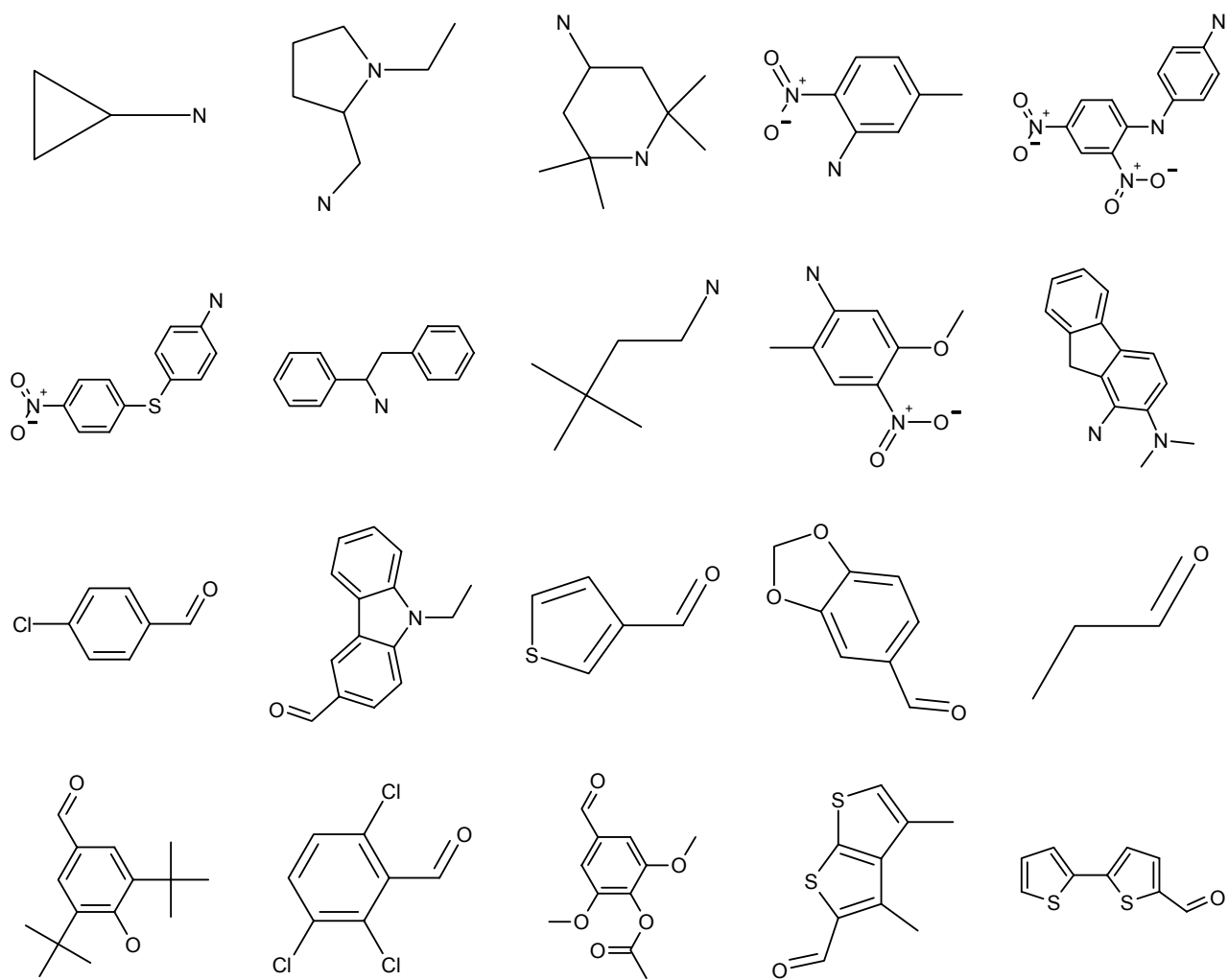


Fig. 6

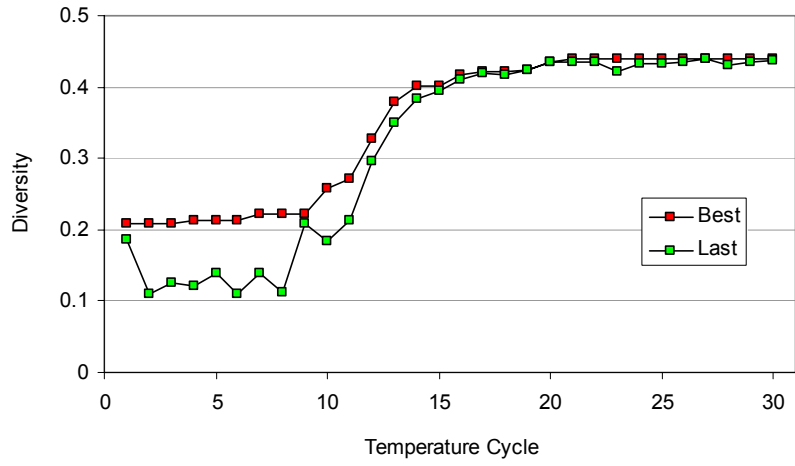


Fig. 7

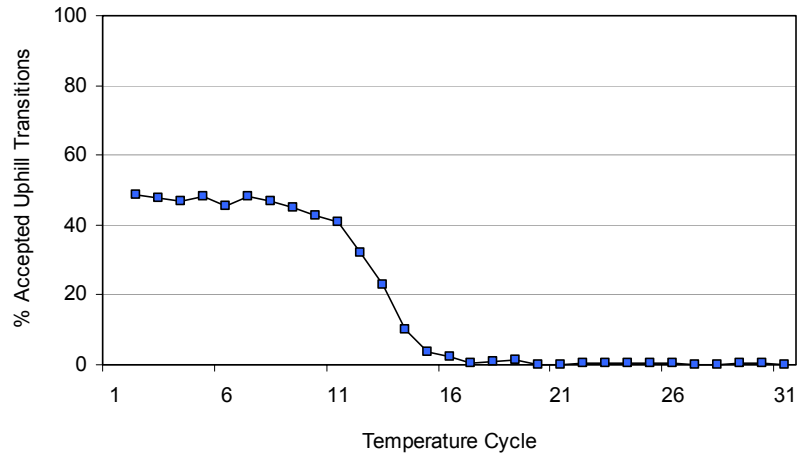


Fig. 8

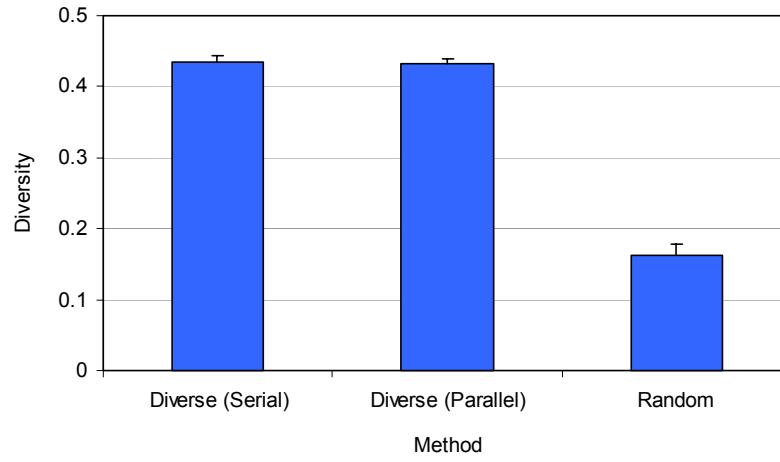


Fig. 9

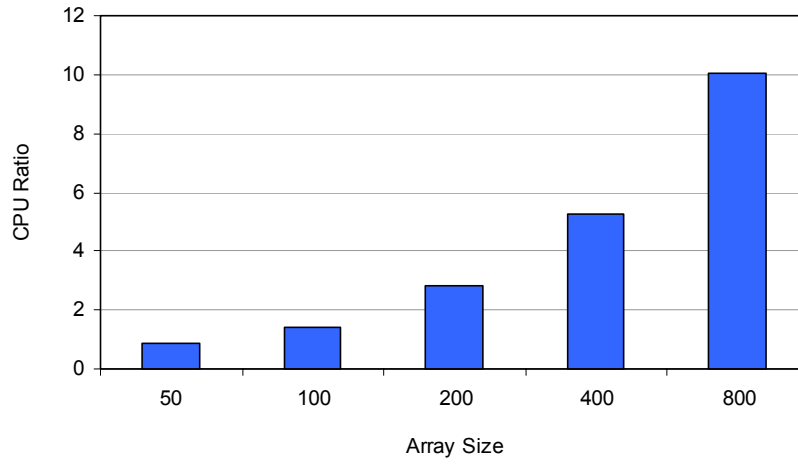


Fig. 10

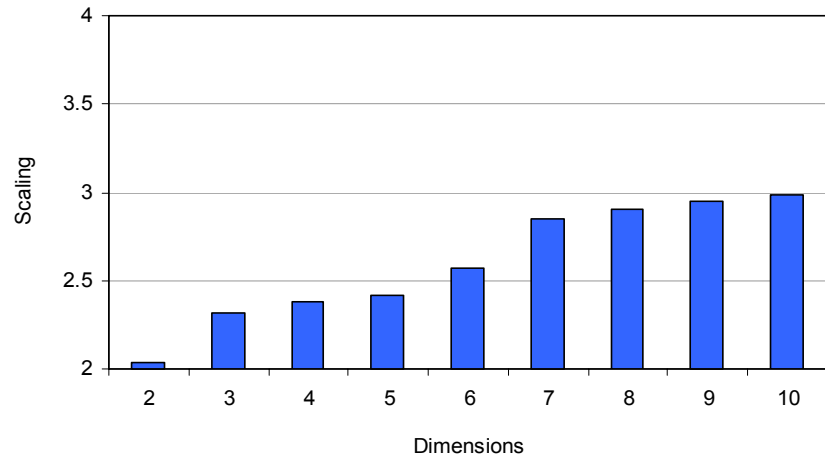


Fig. 11

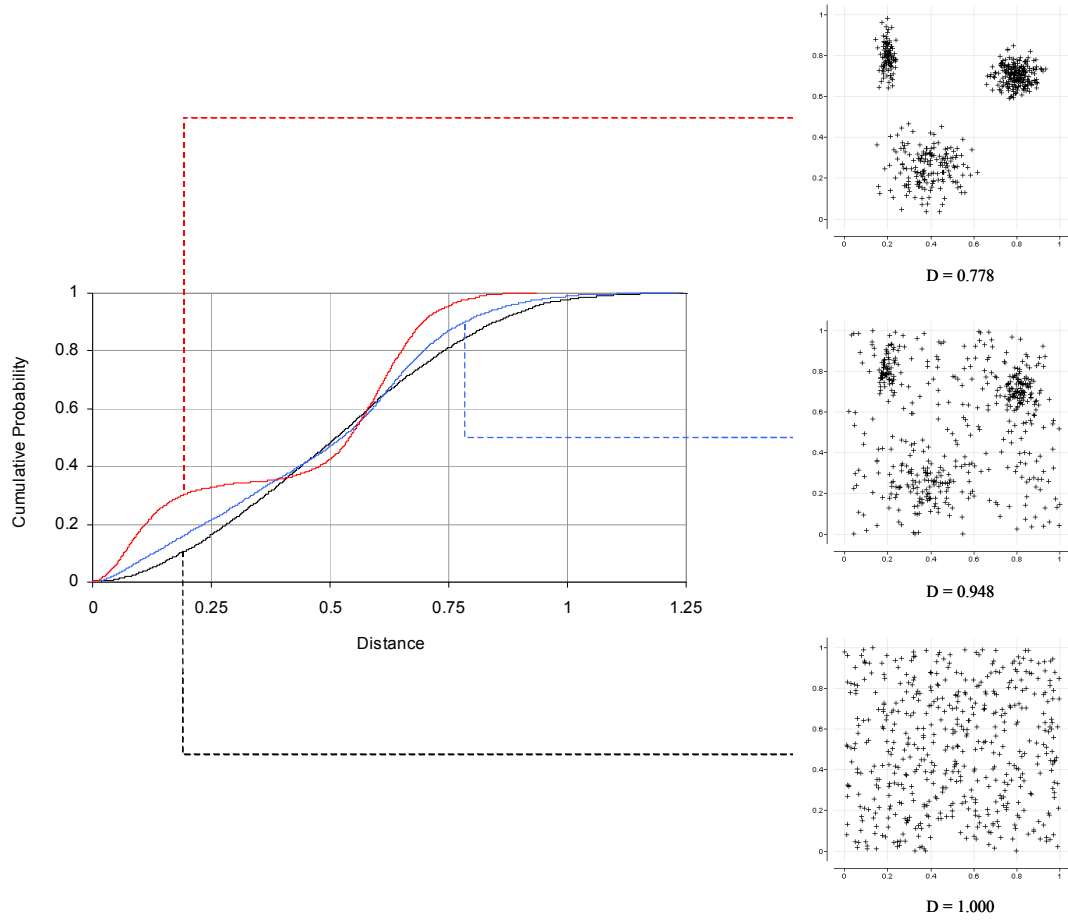


Fig. 12

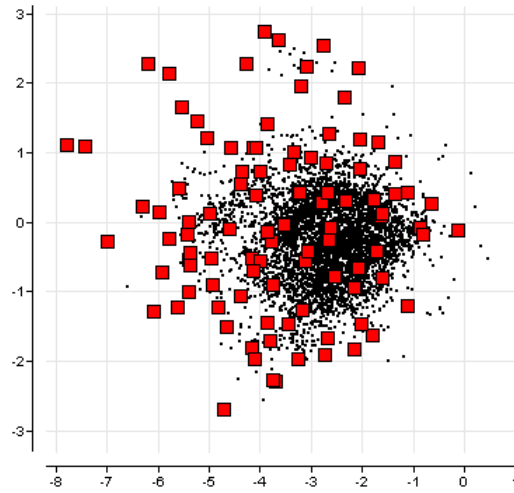


Fig. 13

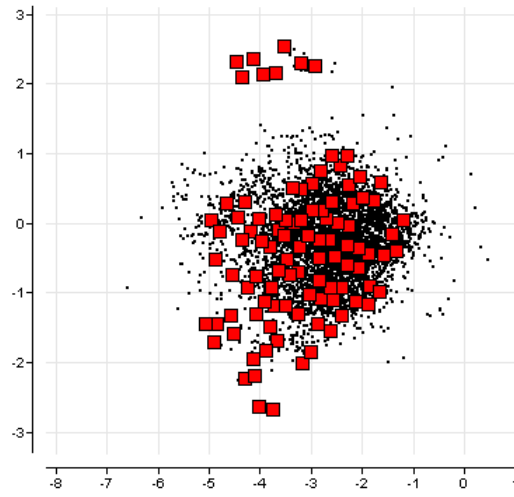


Fig. 14

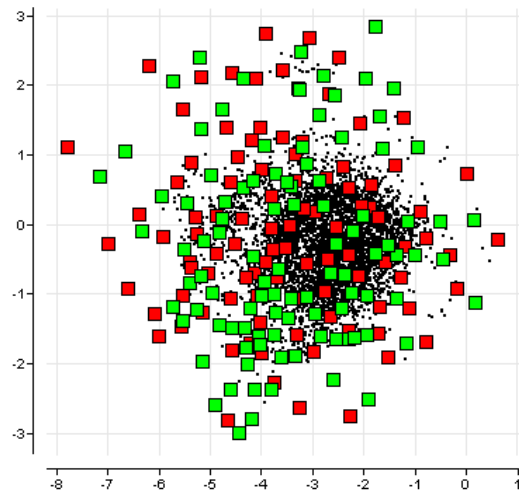


Fig. 15

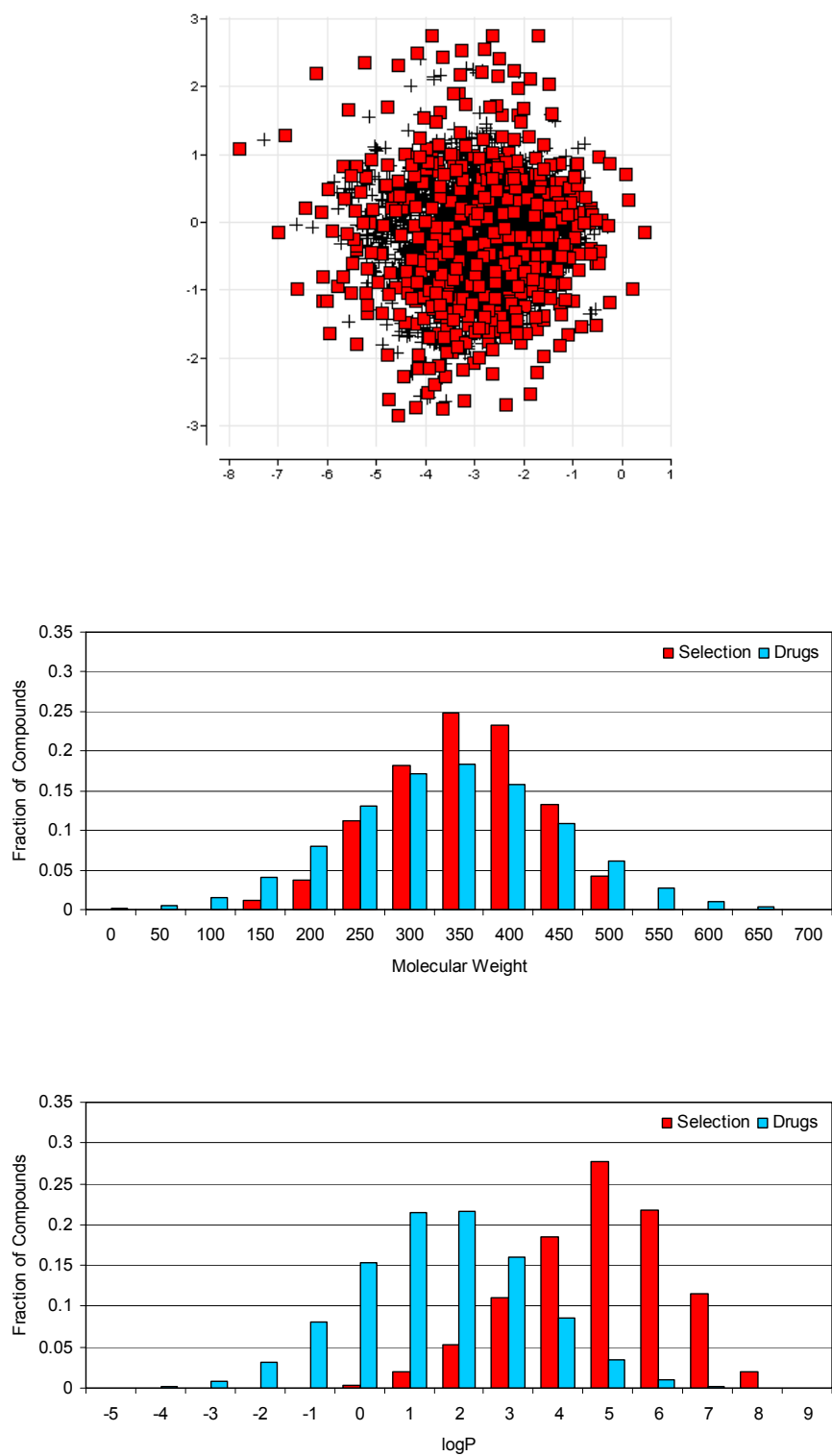


Fig. 16

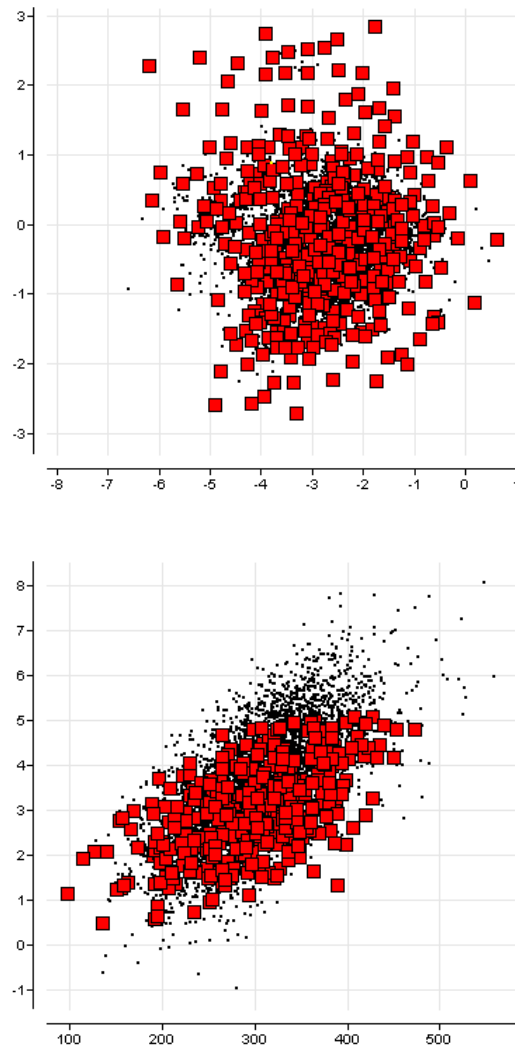


Fig. 17

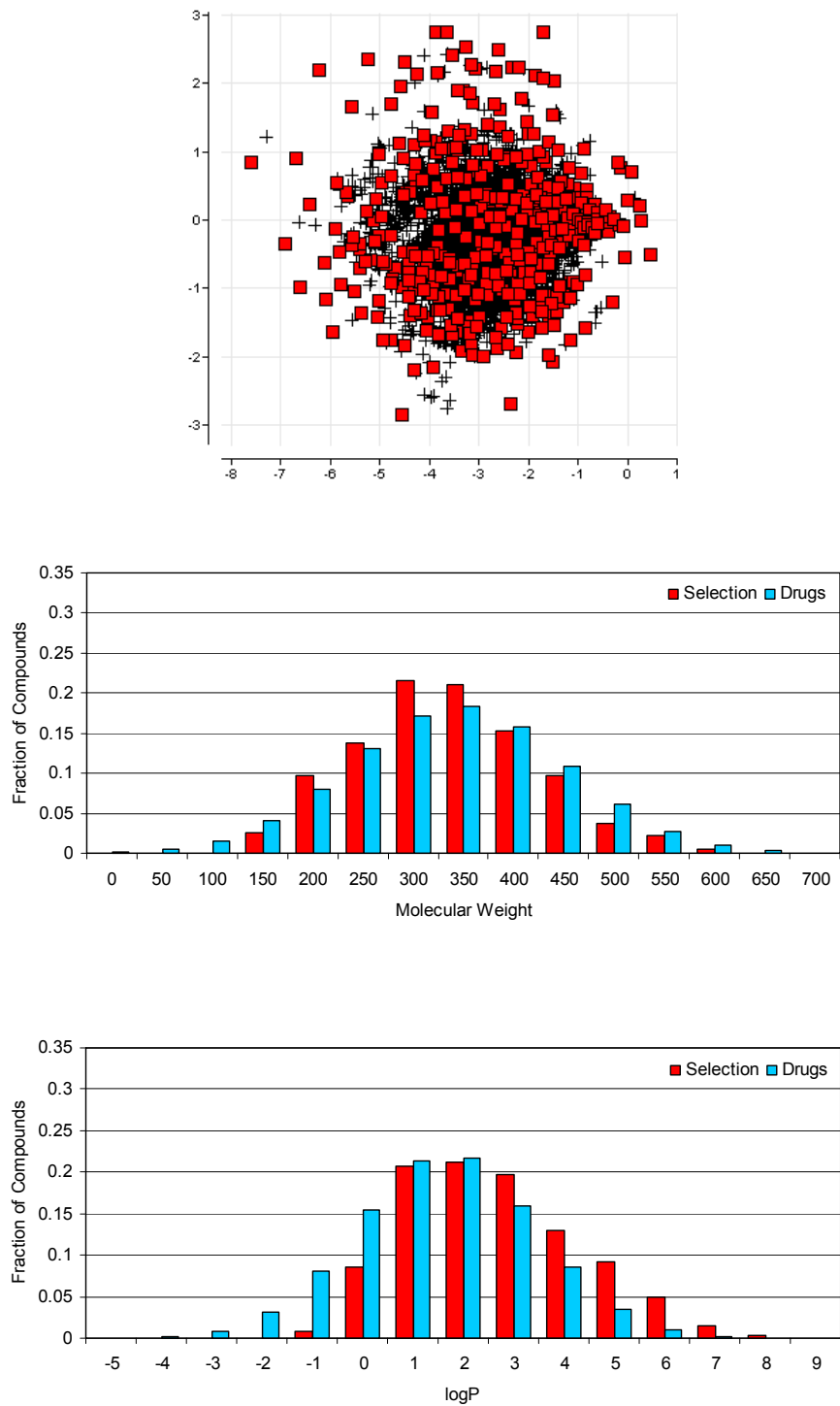


Fig. 18

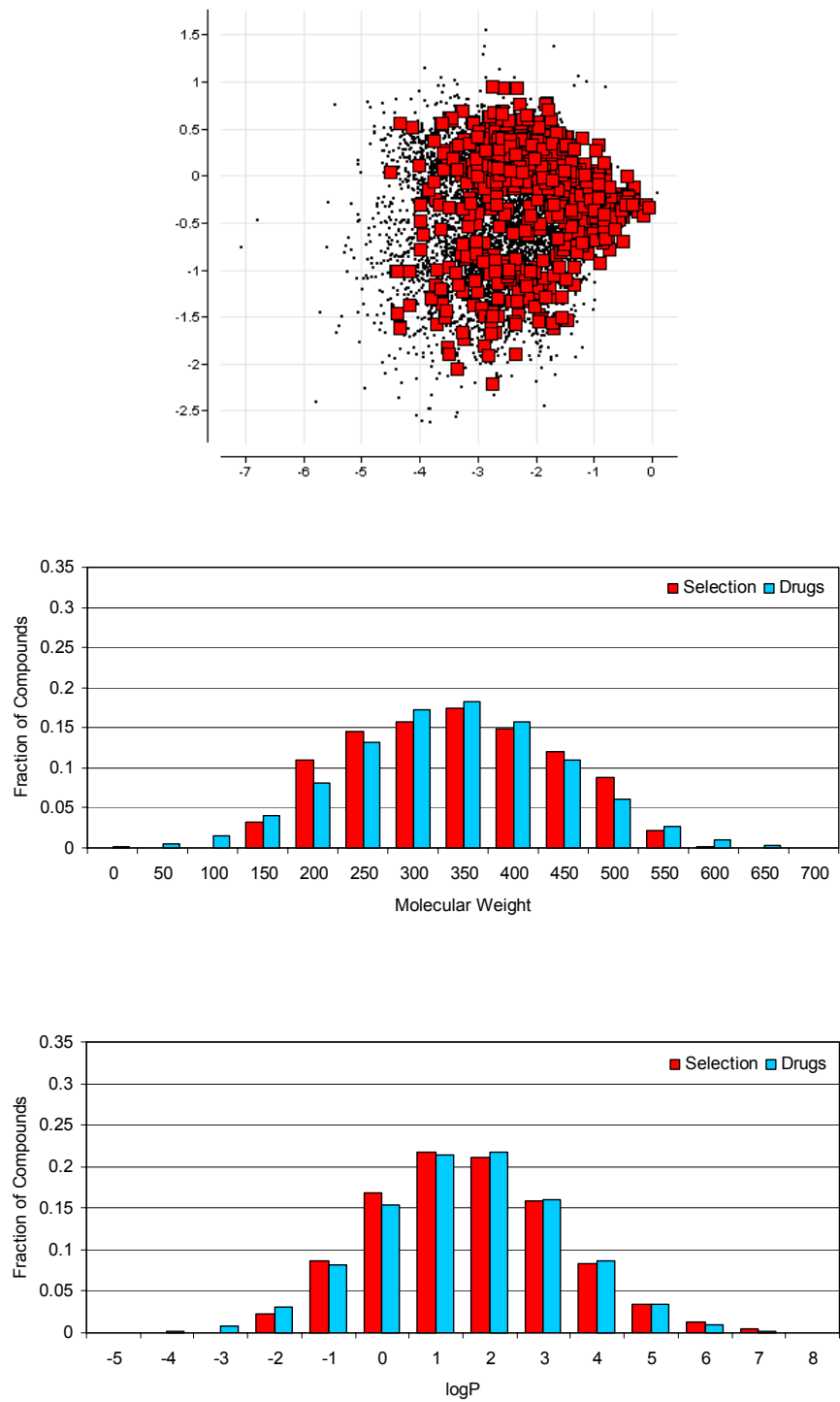


Fig. 19

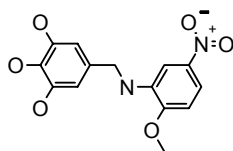


Fig. 20

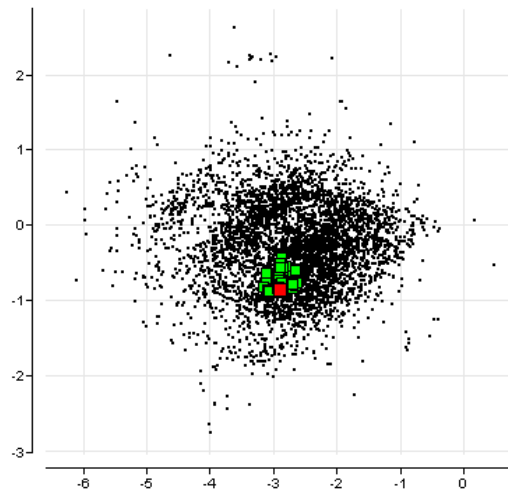


Fig. 21

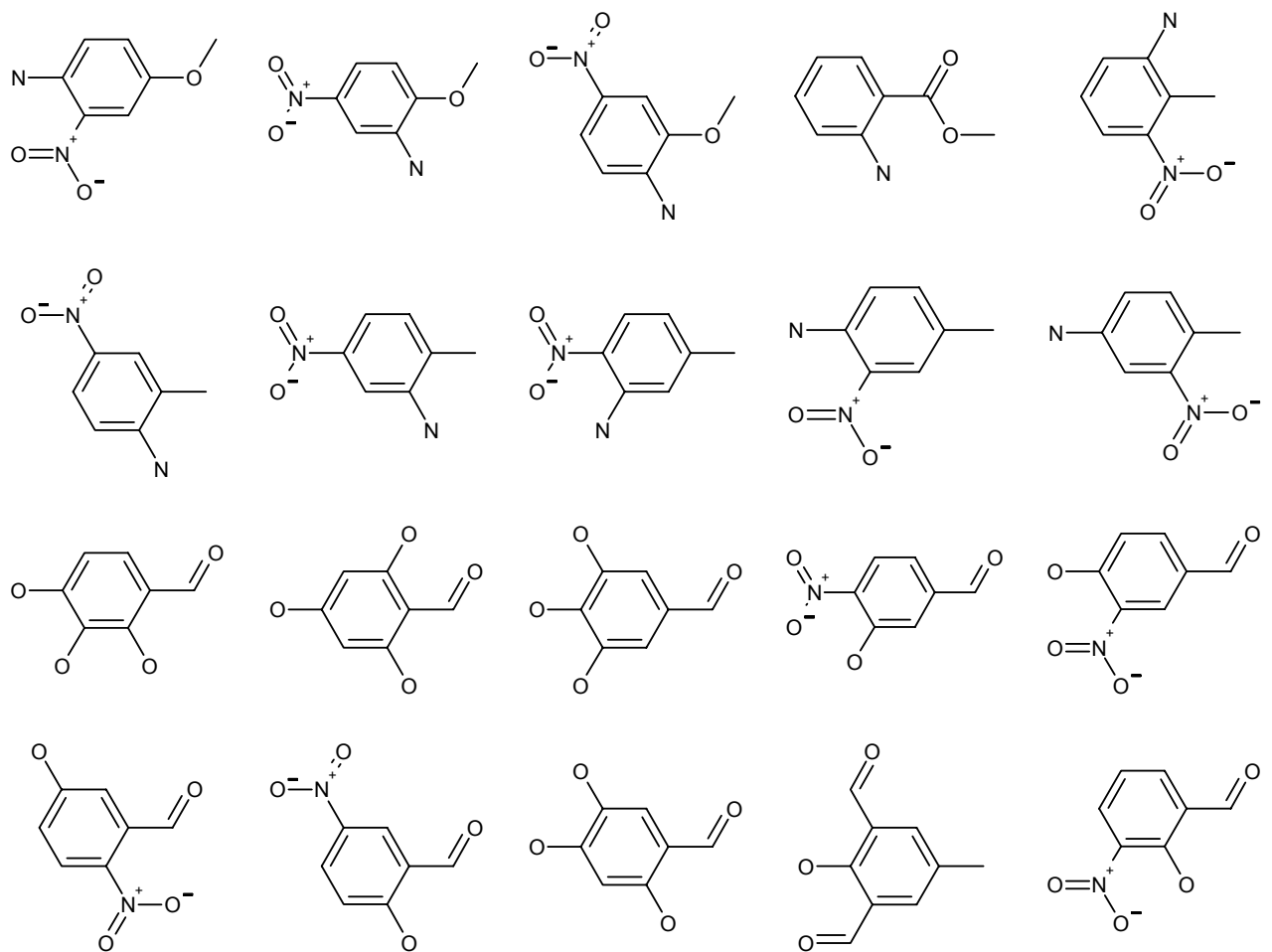


Fig. 22

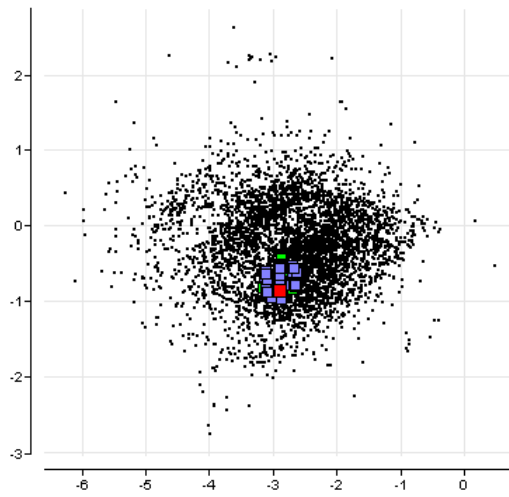


Fig. 23

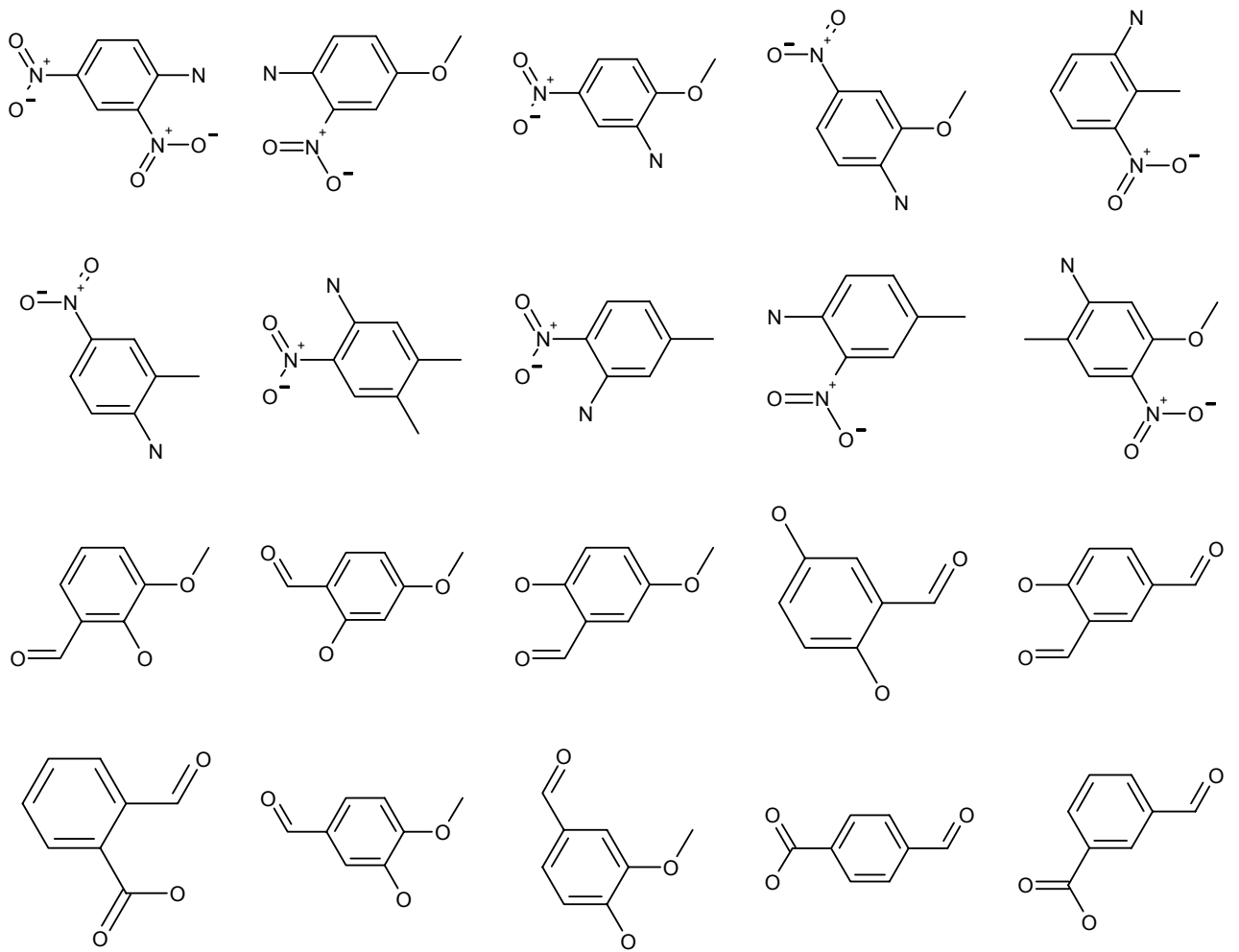


Fig. 24