

Optimizing the Size and Configuration of Combinatorial Libraries

Trudi Wright,[†] Valerie J. Gillet,^{*,†} Darren V. S. Green,[‡] and Stephen D. Pickett[‡]

Krebs Institute for Biomolecular Research and Department of Information Studies, University of Sheffield, Western Bank, Sheffield S10 2TN, United Kingdom, and GlaxoSmithKline, Gunnels Wood Road, Stevenage, SG1 2NY, United Kingdom

Received August 16, 2002

This paper addresses a major issue in library design, namely how to efficiently optimize the library size (number of products) and configuration (number of reagents at each position) simultaneously with other properties such as diversity, cost, and drug-like physicochemical property profiles. These objectives are often in competition, for example, minimizing the number of reactants while simultaneously maximizing diversity, and thus present difficulties for traditional optimization methods such as genetic algorithms and simulated annealing. Here, a multiobjective genetic algorithm (MOGA) is used to vary library size and configuration simultaneously with other library properties. The result is a family of solutions that explores the tradeoffs in the objectives. This is achieved without the need to assign relative weights to the objectives. The user is then able to make an informed choice on an appropriate compromise solution. The method has been applied to two different virtual libraries: a two-component aminothiazole library and a four-component benzodiazepine library.

INTRODUCTION

Many approaches to combinatorial library design have been developed since the emergence of high-throughput screening and combinatorial synthesis techniques.¹ The methods continue to evolve in response to changes in the ways in which the technologies are applied. For example, the initial emphasis was on very large diverse libraries, and more recently it has shifted toward the design of smaller more targeted libraries.² Whether the primary aim is to design diverse or focused libraries, it is now recognized that additional criteria should also be optimized, for example, the compounds contained within the library should, as far as possible, have drug-like or lead-like physicochemical properties.^{3–5} In addition, the high costs associated with the drug discovery process impose financial constraints on the cost of the reactants used. Thus, combinatorial library design is now recognized as a multiobjective optimization problem, MOP.^{6–14}

Most traditional optimization methods such as simulated annealing and genetic algorithms (GAs) are designed to solve single objective optimization problems, where the optimal solution can usually be clearly defined. For example, when designing a library based solely on diversity, typically there is one solution only, that which achieves maximum diversity. However, combinatorial library design is typical of most real-world problems that often involve the simultaneous optimization of several incommensurate and often competing objectives. The difficulty in solving MOPs arises from the fact that when two or more objectives are conflicting they cannot be optimized simultaneously into a single identifiable solution.

Despite the fact that the objectives in library design are usually in conflict, most approaches to multiobjective library design are based on aggregation methods^{6,8–10,14} whereby different objectives are combined into a single function with the problem then being solved as a single objective problem. For example, in the SELECT program multiple objectives are handled via a weighted-sum fitness function and optimized within a GA.⁸ The limitations of this approach were illustrated in our earlier work¹¹ and include difficulty in determining appropriate weights for the objectives, which is usually done by trial and error. In addition, a single solution is found, representing one particular compromise in the objectives.

The MoSELECT^{11–13} program was developed to overcome the limitations of using an aggregation method for multiobjective library design. MoSELECT handles multiple objectives via a multiobjective genetic algorithm (MOGA) and produces a family of equivalent solutions where each solution represents a different compromise in the objectives. MOGAs belong to the class of Multiobjective Evolutionary Algorithms (MOEAs) known as Pareto-based methods and were developed by Fonseca and Fleming.¹⁵ In a MOGA an optimal solution is defined using the concept of Pareto optimality where a solution is Pareto optimal when no other solution is superior to it taking all the objectives into account. In a typical MOP a family of Pareto optimal solutions exists and, in MOGA, the population nature of the GA is exploited to search for such multiple solutions in parallel. As the search progresses a hypersurface or tradeoff surface within the search space is mapped out whereby solutions that lie on the surface are considered equal in terms of overall fitness and are referred to as Pareto or nondominated solutions.

We have recently reported the application of Pareto ranking to deriving a family of QSAR models that are simultaneously optimized on accuracy, complexity, and

* Corresponding author phone: +44 1142 222 652; e-mail: v.gillet@sheffield.ac.uk.

[†] University of Sheffield.

[‡] GlaxoSmithKline.

- Randomly generate a population of potential solutions.
1. For each new generation evaluate members within the population using Pareto ranking and niche induction.
 2. Select parent individuals at random but with a bias towards the best;
 3. Apply evolutionary operators to produce offspring that may involve;
 - a. A 50% chance of one point crossover
 - b. A 50% chance of mutation
 Until 10% of the chromosomes have undergone change.
 4. Start new iteration to form the next generation.
 5. Repeat steps 2-5 for a fixed number of generations.

Figure 1. Outline of MoSELECT.

interpretability.¹⁶ The only other application of Pareto methods in the field of chemoinformatics that we are aware of is the maximum common substructure method described by Handschuh et al.¹⁷

MoSELECT has been demonstrated to be effective at finding families of equivalent solutions for both diverse and focused library designs. The objectives have included diversity, similarity to a target compound, physicochemical property profiles, a measure of bioavailability, and library cost. Using MoSELECT there is no longer a need to assign relative weights to the objectives, and the end result is that the user is able to make an informed choice on an appropriate compromise solution, rather than being presented with a single somewhat arbitrary solution.

A feature of MoSELECT (and many other library design algorithms) is that library size and configuration must be specified by the end user at run-time and then remain fixed throughout. This approach assumes that optimum library size and configuration are known a-priori; however, this is not usually the case especially when optimizing over a number of objectives. For example, the number of reactants used has an obvious effect on the cost of a library, and, as will be shown later, there is usually a conflict between minimizing the number of reactants and optimizing other properties such as diversity.

The aim of this study is to investigate the effect of varying library size and configuration simultaneously with other library characteristics such as diversity and drug-like physicochemical properties. This has been achieved within a MOGA framework in the new program MoSELECT.II. The method has been applied to two different virtual libraries: a two-component aminothiazole library and a four-component benzodiazepine library.

METHODS

The original MoSELECT method is outlined before the modifications required to allow library size and configuration to be optimized are described.

MoSELECT. An outline of the MOGA that forms the heart of MoSELECT is given in Figure 1. The main way in which a MOGA differs from a GA is in the specification of the fitness function. In a GA, a single fitness value, which may be a weighted-sum of multiple objective values, is used

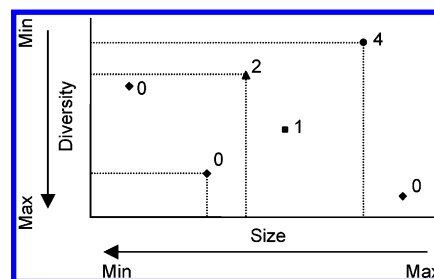


Figure 2. Pareto ranking shown for a library where the aim is to minimize size while maximizing diversity. The normal direction of the y-axis is reversed so that the direction of improvement in both objectives is toward the origin.

to rank the population following each generation. Parent selection is then biased toward the fitter individuals. In a MOGA, multiple objectives are treated independently, and the fitness ranking of a GA is replaced by Pareto ranking which is based on the concept of dominance. A nondominated solution is one where an improvement in one objective results in deterioration in one or more of the other objectives when compared with the other solutions in the population. In the Pareto ranking scheme implemented in MoSELECT, an individual is assigned a rank according to the number of times it is dominated by other individuals in the population. Thus, nondominated individuals are assigned rank 0; individuals dominated by one other individual are assigned rank 1, and so on. Parent selection is then biased toward solutions with lower ranks. The Pareto ranking method is illustrated in Figure 2 for a hypothetical library where the aim is to minimize library size while maximizing diversity. The direction of the y-axis has been reversed so that optimum values in both objectives are toward the origin and the ranks assigned to various individuals are as shown.

Niching. MOGAs can have a tendency to genetic drift where they converge toward a single solution. This effect can be minimized by ensuring that the diversity of the individuals in a population is preserved. Consequently many MOEAs incorporate niching techniques in their design strategy. In MoSELECT the preservation of diversity is maintained through the implementation of fitness sharing.¹⁸ The aim is to ensure that nondominated individuals are evenly distributed on the Pareto surface. After every generation, a niche radius is defined as a percentage of the range of values that exist for each objective within the nondomi-

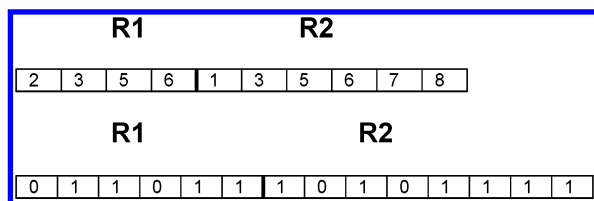


Figure 3. Chromosome encodings: (a) the integer representation of a 4×6 library used in MoSELECT; (b) the binary representation of the same 4×6 library used in MoSELECT.II.

nated set. The first solution encountered forms the center of a niche. If the next individual is within the niche radius of the first for all objective values, then the individual is penalized by increasing its rank, otherwise it forms the center of a new niche. This process continues for all nondominated individuals in the population.

Varying Library Size and Configuration. The fixed library size and configuration in MoSELECT (and SELECT) is due to the chromosome encoding scheme used in the MOGA in MoSELECT (and the corresponding GA in SELECT). The existing encoding scheme is shown in Figure 3a and involves representing a potential solution (combinatorial subset selected from the full virtual library) as an integer string. The chromosome is divided into partitions according to the number of reactant pools (variable substitution positions) in the library and the length of a partition corresponds to the number of reactants to be selected from each pool. Each element in a partition represents a unique integer value that corresponds to a particular reactant. The size of the chromosome is determined by the size and configuration of the combinatorial subset required, as specified on input, and is independent of the virtual library size.

To allow library size and configuration to vary dynamically during the search, the chromosome representation must be modified. Thus, the chromosome encoding used in MoSELECT.II is shown in Figure 3b. The encoding scheme is similar to that described for the GALOPED program.¹⁹ The entire virtual library is represented as a binary string. As before, the string is partitioned according to the number of reactant pools; however, now each bit in a partition represents a reactant available in the given reactant pool. A bit value of '1' indicates that the reactant has been selected in the combinatorial subset, whereas a bit value of '0' indicates its absence. Using this encoding scheme, the number of bits set to "1" in a chromosome can vary as a result of the genetic operators, crossover and mutation and thus the library size and configuration varies correspondingly. A chromosome is initialized by, first, generating a random number to determine how many bits should be set on and, second, generating random numbers to choose the bit positions to set to "1".

As described previously, diversity and library size are usually in conflict for example, when measuring diversity using a partitioning scheme larger libraries typically result in greater coverage of cells. Thus, when optimizing on diversity alone there will be a tendency to select very large libraries. The modified chromosome encoding scheme allows size to be included as an objective together with diversity with each objective being handled independently within the MOGA of MoSELECT.II. Hence, the tradeoff between size and diversity can thus be explored.

Handling Constraints. While it can be useful to have the ability to view the entire tradeoff surface, in practice, the user may be constrained to operate within certain limits. For example, cost considerations can mean that it is necessary to impose limits on library size, select solutions that are combinatorially efficient, or to maximize plate coverage. Each of these constraints is described in more detail below:

Plate Coverage. Combinatorial libraries are usually synthesized on 96-well plates (with one compound per well). Libraries that do not fully occupy all plates (i.e., libraries that are not multiples of 96) lead to inefficiencies with respect to the programming of the synthesis robots. Additionally, efficiency in use of materials will be maximized when the number of wells occupied by the library on a single plate is maximized.

Library Size. As already mentioned, the trend in library design has shifted away from the very large libraries that were common in the mid 1990s toward the design of smaller libraries, for example, libraries containing a few thousand compounds. Thus, having the ability to set upper and lower bounds on library size can be extremely useful. Being able to view the entire tradeoff surface between size and other library properties allows the user to choose appropriate limits and hence direct the search to restricted regions of the search space.

Combinatorial Efficiency. Combinatorial efficiency refers to the number of reactants required to generate a given number of products. The fewer the reactants required the more efficient is the library. Consider a two-component library consisting of 400 products. There are many configurations that can give rise to 400 products, for example, 400×1 ; 200×2 ; 100×4 ; ...; 20×20 ; ...; 1×400 . The total number of reactants required for each configuration is as follows: 401; 202; 104; ...; 40; ...; 401, respectively. In the absence of any other criteria, the most efficient configuration is 20×20 since this requires the minimum number of distinct reactants (40). Thus, all other criteria being equal, the most combinatorially efficient solution will be the preferred solution.

Constraints on library size and configuration have previously been applied in the GA-based GALOPED¹⁹ program and in the Monte Carlo search algorithm described by Pickett et al.²⁰ In GALOPED, constraints can be used to place upper and lower limits on library size, and potential solutions that violate the constraints are eliminated from the search. In the method described by Pickett et al. a number of parameters are used to set constraints for each library component including the maximum and minimum number of reactants and the ideal number of reactants to be included in the library design. Library solutions that fall outside the specified constraints are then penalized within the scoring function.

Many approaches have been described for the satisfaction of constraints within a MOEA.²¹ In the following, solutions that satisfy the constraints are referred to as feasible solutions to distinguish them from infeasible solutions that violate the constraints. The simplest and most efficient way to ensure that all solutions satisfy the constraints is to eliminate infeasible solutions from the search space, as done in GALOPED. However, this approach can lead to premature convergence since applying the genetic operators to infeasible individuals could result in feasible offspring.

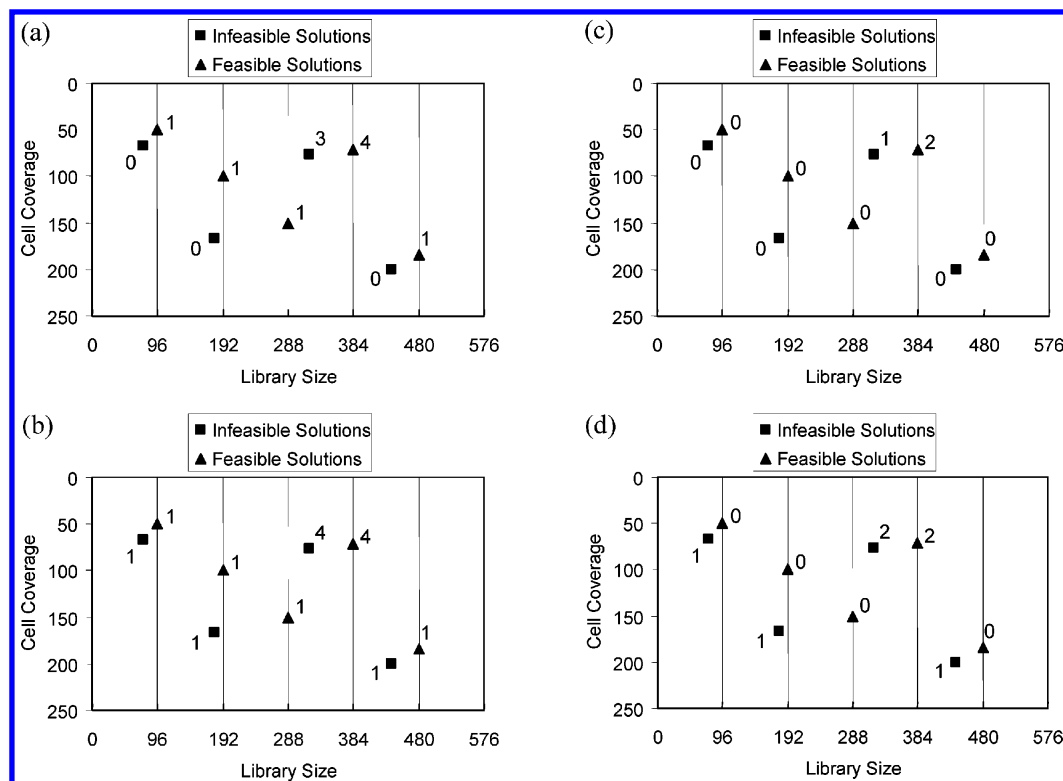


Figure 4. (a) Pareto ranking scheme implemented when no constraints are specified; (b) Simply adding a penalty to infeasible individuals can result in the entire population becoming dominated. (c) In MoSELECT.II, constraints are handled by first dividing the population into two groups: feasible individuals and infeasible individuals. Pareto ranking is carried out within each group separately. (d). The ranks of infeasible individuals are incremented by one.

The approach used by Pickett et al. in their Monte Carlo algorithm is less severe, and infeasible individuals are penalized in the fitness function with the severity of the penalty being determined by the degree to which the constraints are violated. A similar method could be implemented in the MOGA of MoSELECT.II by simply penalizing the ranks of nondominated individuals that violate the constraints so that infeasible nondominated individuals are given ranks as if they are dominated. Thus, infeasible solutions can still exist within the population and can be selected for reproduction; however, they have a lower chance of being selected than feasible, nondominated solutions. However, a potential problem with applying this method in a MOGA is that it can lead to a population in which all the individuals are dominated. This is illustrated by an example in Figure 4. Figure 4a shows Pareto ranking for a two-objective library design where the aim is optimize size and diversity simultaneously with no constraints specified, i.e., ranking is based on dominance only. Now, within this optimization scheme, assume that the aim is to constrain libraries to be multiples of 96 (shown by the vertical lines in Figure 4a) so that maximum plate coverage is achieved. Feasible solutions are shown by the triangles and fall on the vertical lines, whereas infeasible solutions are shown as squares and fall between the vertical lines. It can be seen in Figure 4a that none of the nondominated solutions satisfy the constraint. In the simple penalty scheme, these individuals will be penalized, for example, by increasing their rank by one, as shown in Figure 4b; however, now the entire population has become dominated.

Thus, in MoSELECT.II, the simple penalty method is modified to prevent the situation from occurring in which

the entire population is dominated. Prior to Pareto ranking, the population is divided into two groups: one group consisting of feasible individuals and the other consisting of infeasible individuals. Pareto ranking is then carried out within each group, that is, feasible individuals are ranked against feasible individuals, and infeasible individuals are ranked against infeasible individuals, as shown in Figure 4c. Penalties for violating the constraints are then applied in a second step where the ranks of infeasible individuals are incremented by one, as shown in Figure 4d. The effect is 2-fold: first, infeasible solutions cannot dominate feasible nondominated solutions so that there should always be at least one nondominated solution in the population; and second, since infeasible solutions are assigned higher ranks the search will favor nondominated feasible solutions over nondominated infeasible solutions.

In MoSELECT.II, the different constraints are specified in the following ways: plate constraints are specified as the minimum percentage of plate coverage that is allowed; size constraints are specified as upper and lower bounds on the final number of products in a library; and combinatorial efficiency constraints are specified as upper and lower limits on the number of reactants allowed for each library component.

Virtual Libraries. The effect of varying library size and configuration simultaneously with other library properties is investigated using two different virtual libraries: a two-component aminothiazole library and a four-component benzodiazepine library. The aminothiazole library consists of 12 850 products generated from 74 α -bromoketones coupled with 170 thioureas.¹¹ The benzodiazepine library consists of 256 036 products (constructed from $23 \times 22 \times$

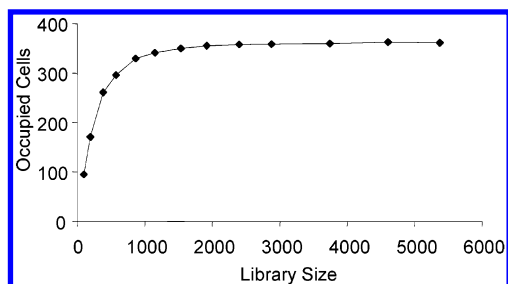


Figure 5. Cell coverage is plotted against library size for multiple runs of the SELECT program on the aminothiazole library.

23 × 22 reactants). (The benzodiazepine library is benzo_big included within Cerius² 4.5.²²).

Diversity was measured using a cell-based measure. A 3D chemistry space was constructed for each of the libraries by first calculating Cerius² default topological parameters and physicochemical properties for the virtual library and second by performing principal components analysis and selecting the top three principal components. The virtual library was then mapped into the 3D space and diversity was measured by the number of cells occupied by the library. The aminothiazole library was found to occupy 364 of the 1134 cells on the grid, and the larger benzodiazepine library was found to occupy a total of 997 cells. The experiments that are designed to achieve maximum diversity aim to select libraries that occupy the maximum number of cells in the 3D space.

RESULTS AND DISCUSSION

Optimizing Size and Diversity. The tradeoff between diversity and library size was investigated manually for the aminothiazole library by performing multiple runs of SELECT with each run configured to maximize diversity for a library of different fixed size and configuration. Thirteen runs of SELECT were carried out for libraries ranging from 96 products up to 5376 products. The results are plotted in Figure 5 where it can be seen that diversity increases steadily from 96 cells for the smallest library up to 341 cells for a library of 1152 compounds (the sixth largest library). A small gain in diversity (20 additional cells) is seen over the remaining seven libraries with the largest library occupying 361 cells with nearly four times the number of compounds. This can be compared with the virtual library of 12 580 compounds which occupies 364 cells.

The graph shows that there is no single solution to the maximum diversity versus minimum library size optimization problem. In fact, a family of different solutions exist that lie on the curve shown in the graph and that represent different compromises in the objectives. In the absence of any further information none of these libraries can be said to be better than the others (although in practice, less importance would be given to the flatter part of the curve where the tradeoff in the objectives is less pronounced).

MoSELECT.II was then used to investigate the extent to which the tradeoff surface can be explored in a single run using the MOGA approach where library size is optimized simultaneously with diversity. The program was run for 10 000 iterations with a population size of 250 with the genetic operators set at a 50% chance of crossover or mutation. A niche radius of 1% was specified. The resulting

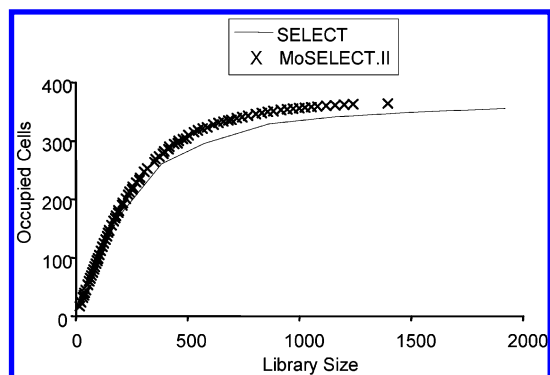


Figure 6. Solutions produced by a single run of MoSELECT.II (crosses) are compared with the multiple runs of SELECT (continuous line) for the aminothiazole library.

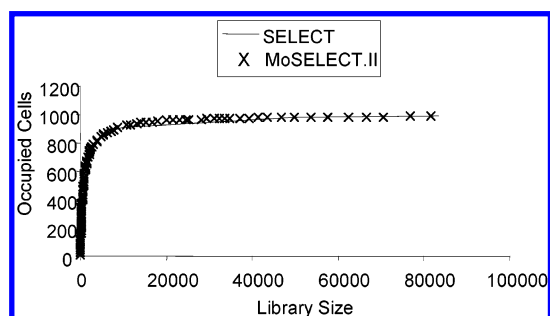


Figure 7. Solutions produced by a single run of MoSELECT.II (crosses) are compared with multiple runs of SELECT (continuous line) for the benzodiazepine library.

family of solutions is shown in Figure 6 where they are plotted as crosses and are superimposed on the results found manually using SELECT (represented by the continuous line). A total of 82 nondominated solutions were obtained.

Each MoSELECT.II solution represents a different compromise in size and diversity. The family of solutions spans the tradeoff surface found using SELECT up to a library size of 1395 compounds. In fact, this library found by MoSELECT.II occupies the same number of cells that is occupied by the entire virtual library, i.e., 364 cells. Larger libraries are not found since they could not have greater diversity and hence would be dominated by this solution. A slight increase in diversity is seen for a given library size relative to SELECT and this is most likely due to the ability of MoSELECT.II to vary configuration along with size, whereas in the SELECT runs the configuration was fixed for each library size. The MoSELECT.II run took 320s for 10 000 iterations on a R12K SG workstation running at 360 MHz. (All subsequent timings refer to the same processor.)

A similar run of MoSELECT.II was carried out for the benzodiazepine library. The population size was increased to 400 due to the increased size of the virtual library and the larger search space to be explored relative to the aminothiazole library. All other parameters were kept the same. Figure 7 shows the results. Again the MoSELECT.II solutions are shown as crosses and are superimposed on a continuous line that represents the diverse libraries found over multiple SELECT runs configured for different library sizes. In this case, the gain in diversity begins to tail off for libraries with > 10K compounds: 923 cells are occupied by the library consisting of 10692 compounds and the maximum diversity found is 988 cells occupied for a library of 81600 compounds. For comparison, the virtual library of 250K

compounds occupies 997 cells. The MoSELECT run took 2770 s (46 min) for 10 000 iterations with the increase in time relative to the aminothiazole library due to the larger population size and the larger library sizes to be explored.

Thus, MoSELECT allows the user to explore the tradeoff between size and diversity in a single run. In a typical library design scheme, a user can make an informed choice on what represents an appropriate compromise solution, for example, a library might be chosen from the region of the curve where the gain in diversity begins to flatten off with increasing size. Mapping the entire tradeoff surface is achieved in approximately the same amount of time that it takes for a single run of SELECT that will find one solution only for a library of fixed size and configuration.

Constraining the Search. As mentioned earlier, practical considerations may mean that there are constraints on library design, and the following experiments were designed to investigate the effect of constraining the search.

Plate Coverage. MoSELECT.II was run on the aminothiazole library with the aim of maximizing diversity while minimizing size with the constraint that any solution found should cover at least 75% of the final plate. The remaining MoSELECT.II parameters were as before. The solutions found are shown by the unfilled circles in Figure 8a where the vertical lines represent multiples of 96 (maximum plate coverage). The solutions found in the previous unconstrained run (illustrated in Figure 6) that happen to fall within the constraints are shown in Figure 8a by the crosses. When the solutions found in the constrained run are compared with the solutions found in the unconstrained run it can be seen that MoSELECT.II has been successful in finding a family of solutions that satisfies the plate constraint and that this is achieved without any compromise on the diversity of the solutions. In addition, a greater number of solutions that satisfy the constraints has been identified and this suggests that the constraints have influenced the search process itself as well as simply removing irrelevant solutions from the results.

Figure 8b,c shows the effect of applying a constraint that ensures all libraries produced are an exact multiple of 96 for the aminothiazole and benzodiazepine libraries, respectively, i.e., when all plates are fully occupied. Again the results found in the constrained runs are shown as unfilled circles. For the aminothiazole library, thirteen solutions are found that maximize plate coverage as compared to three solutions found in the unconstrained search (shown as crosses). In the case of the benzodiazepine library, a large number of solutions are found to satisfy the constraints.

Constraining Library Size. Runs were performed on the aminothiazole library with the aim of maximizing diversity and minimizing size with the additional upper and lower bounds on acceptable size ranges of 200–400; 400–600; 600–800; and 100–1000; respectively. The results are shown in Figure 9a–d. In each case, MoSELECT.II finds a family of solutions that satisfies the size constraints without compromising on the diversity of the solutions. The continuous line shows the corresponding solutions found for the SELECT runs shown in Figure 6.

Similar results were found when size constraints were applied to the benzodiazepine libraries.

Combinatorial Efficiency. Finally, the effect of imposing constraints on combinatorial efficiency was investigated

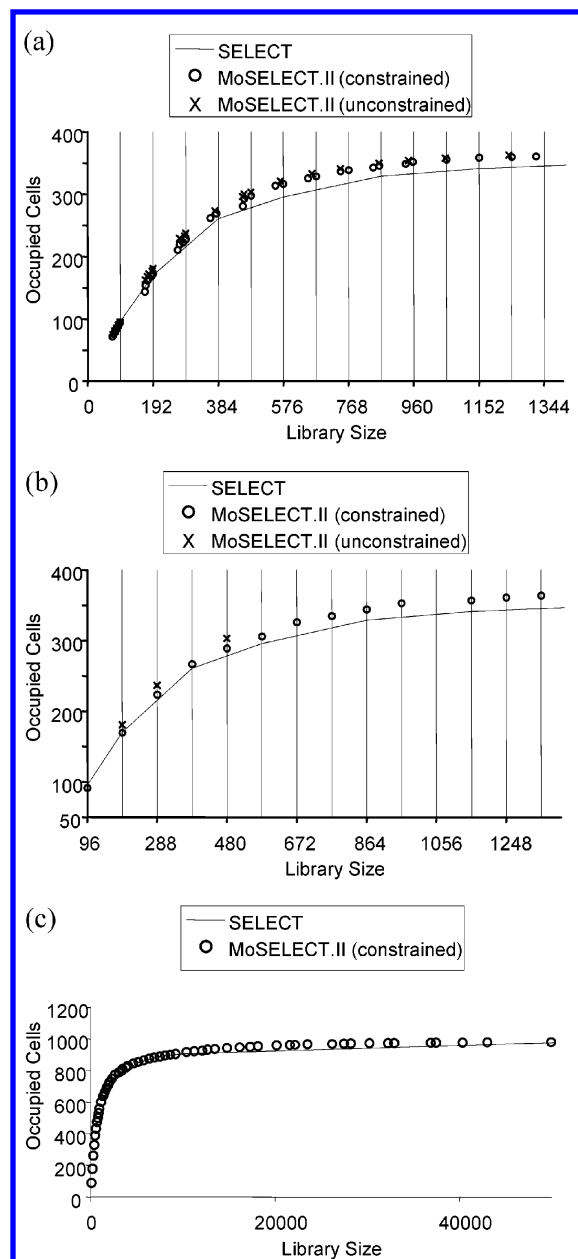


Figure 8. Solutions found when running MoSELECT.II on the aminothiazole library with the constraint that the final plate has a minimum of (a) 75% coverage; (b) 100% coverage. (c) The benzodiazepine library with the constraint that the final plate is fully occupied.

while simultaneously maximizing diversity and minimizing size. MoSELECT.II was configured to explore aminothiazole libraries within the size range of 400–600 products with the additional constraint that any solution should consist of a minimum of 20 and a maximum of 25 reactants from each pool.

The results are shown as circles in Figure 10(a) where the number of α -bromoketones selected is plotted on the y-axis and the number of thioureas is plotted on the x-axis. The crosses show the configurations of libraries found within the same size range when no constraints are placed on combinatorial efficiency. These solutions correspond to those shown in Figure 9b. It can be seen that without the constraint there is a tendency to construct libraries using a higher number of thioureas compared to α -bromoketones. Presum-

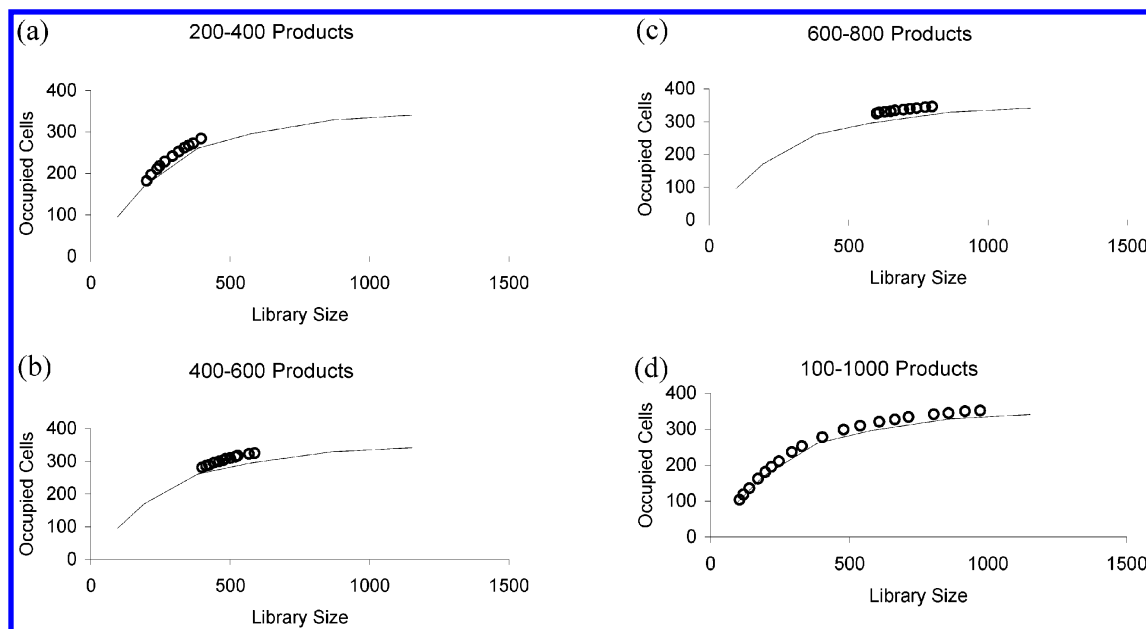


Figure 9. Aminothiazole solutions found by MoSELECT.II for different size constraints.

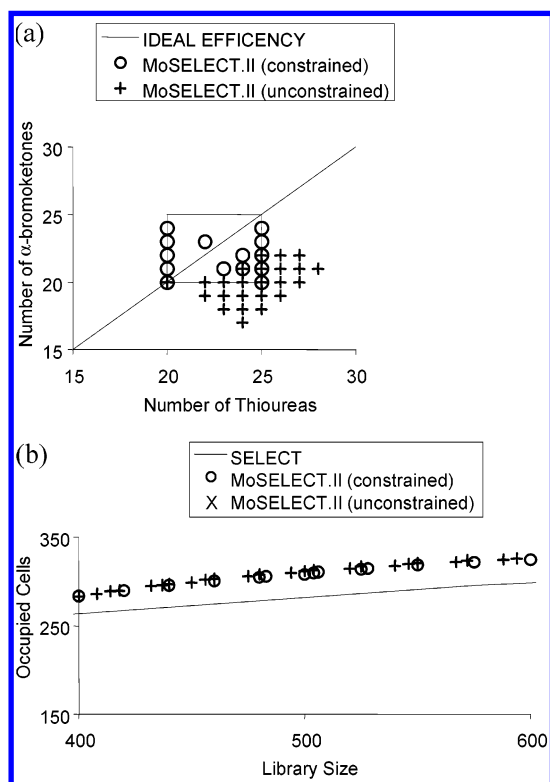


Figure 10. (a) Introducing combinatorial efficiency as a constraint on the aminothiazole libraries. The circles show results found for the constrained runs; the crosses show results found for the unconstrained runs. (b) Both sets of libraries are shown with diversity plotted against size where it can be seen that the combinatorially efficient libraries are found without compromising on diversity.

ably this is due to the fact that the virtual library consists of a higher number of thioureas (74 α -bromoketones and 170 thioureas). When the diversities of the two sets of libraries are compared, the more combinatorially efficient libraries appear on the same diversity versus size curve as the solutions previously found, Figure 10(b). Thus, better efficiency has been achieved without compromising on diversity.

Table 1. (a) Library Configurations Found for the Benzodiazepine Library When the Runs Are Constrained to Include from Four to Six Reactants for Each Position of Variability and (b) When No Constraints Are Placed on Combinatorial Efficiency

| cells occupied | library size | R1 | R2 | R3 | R4 |
|----------------|--------------|----|----|----|----|
| (a) | | | | | |
| 385 | 500 | 4 | 5 | 5 | 5 |
| 426 | 576 | 4 | 6 | 4 | 6 |
| 438 | 600 | 4 | 6 | 5 | 5 |
| 479 | 720 | 4 | 6 | 5 | 6 |
| 491 | 750 | 5 | 6 | 5 | 5 |
| 521 | 864 | 4 | 6 | 6 | 6 |
| 535 | 900 | 5 | 6 | 5 | 6 |
| 572 | 1080 | 5 | 6 | 6 | 6 |
| 600 | 1296 | 6 | 6 | 6 | 6 |
| (b) | | | | | |
| 401 | 504 | 3 | 4 | 6 | 7 |
| 415 | 525 | 3 | 5 | 5 | 7 |
| 424 | 540 | 3 | 5 | 6 | 6 |
| 431 | 576 | 3 | 4 | 6 | 8 |
| 446 | 600 | 3 | 5 | 5 | 8 |
| 471 | 648 | 3 | 6 | 6 | 6 |
| 495 | 720 | 3 | 5 | 6 | 8 |
| 507 | 756 | 3 | 7 | 6 | 6 |
| 530 | 840 | 3 | 5 | 7 | 8 |
| 537 | 864 | 3 | 6 | 6 | 8 |
| 543 | 882 | 3 | 7 | 7 | 6 |
| 557 | 945 | 3 | 5 | 7 | 9 |
| 571 | 1008 | 3 | 7 | 6 | 8 |
| 576 | 1029 | 3 | 7 | 7 | 7 |
| 597 | 1134 | 3 | 6 | 7 | 9 |
| 603 | 1176 | 3 | 7 | 7 | 8 |
| 620 | 1260 | 3 | 6 | 7 | 10 |
| 631 | 1323 | 3 | 7 | 7 | 9 |
| 633 | 1386 | 3 | 6 | 7 | 11 |
| 651 | 1470 | 3 | 7 | 7 | 10 |

The experiment was repeated for the benzodiazepine library with the constraint that there should be a minimum of four and a maximum of six reactants selected from each of the four library components. The run was also constrained to search for libraries with between 500 and 1500 products. The aim was to maximize diversity while minimizing size. Table 1(a) shows the configurations produced for the constrained MoSELECT.II runs, whereas Table 1(b) shows

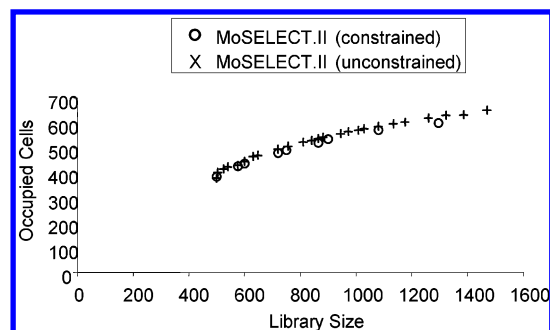


Figure 11. Solutions found when applying a combinatorial efficiency constraint to the benzodiazepine library are shown as circles. Solutions found when no constraints are placed on combinatorial efficiency are shown as crosses.

the libraries found when no constraints are placed on combinatorial efficiency. As for the aminothiazole library, the results show that greater combinatorial efficiency is achieved when the search is constrained. Finally, Figure 11 shows there is no significant decrease in the diversity of the solutions produced when the combinatorial efficiency constraint is included in the search.

Increasing the Number of Objectives. The methodology can be extended to optimize additional objectives simultaneously with library size and configuration. In the following experiments, the aminothiazole libraries were optimized on diversity, size, and distribution of molecular weight (MW) where the aim was to minimize the RMSD between the distribution of molecular weight in the libraries and the distribution of molecular weight in the World Drugs Index (WDI).²³

When there are more than two objectives it is no longer possible to represent the solutions in a simple 2D plot. Thus, solutions are presented using parallel coordinates graphs where the objectives are plotted along the *x*-axis and each solution is shown as a continuous line. Parallel coordinates graphs provide a useful way of visualizing the tradeoff in the objectives since competing objectives are shown by crossing lines. They can also be used to see if extremes in each of the objectives have been found, which gives an indication that the MOGA has converged, and they allow appropriate compromise solutions to be identified.

The first experiment was designed to select aminothiazole libraries with maximum diversity, minimum size, and drug-like molecular weight profiles. A population size of 150 was specified and the niche radius was set to 10%. No constraints were specified. A total of 38 libraries were found as shown in the parallel coordinates graph in Figure 12(a). The objective values have been scaled to allow them to be plotted on the same graph with improvement in all objectives being toward 0 on the *y*-axis. The objectives are scaled such that 0 represents the best value achieved when optimizing an objective on its own. Thus, for diversity, 0 represents a cell occupancy of 364 which is the number of cells occupied by the virtual library itself and 1 represents zero cells. For size, 0 represents a library size of 0, whereas 1 represents a library size of 2000. For MW, 0 represents an RMSD of 1 and 1 represents an RMSD of 19. The tradeoff in the objectives is clearly shown by the crossing lines in the graph. The solution libraries span a size range of 23 to 1886; the range of cell occupancies is 23 to 361; and RMSD of the molecular weight profiles varies from 2 to 10. The objective values of libraries

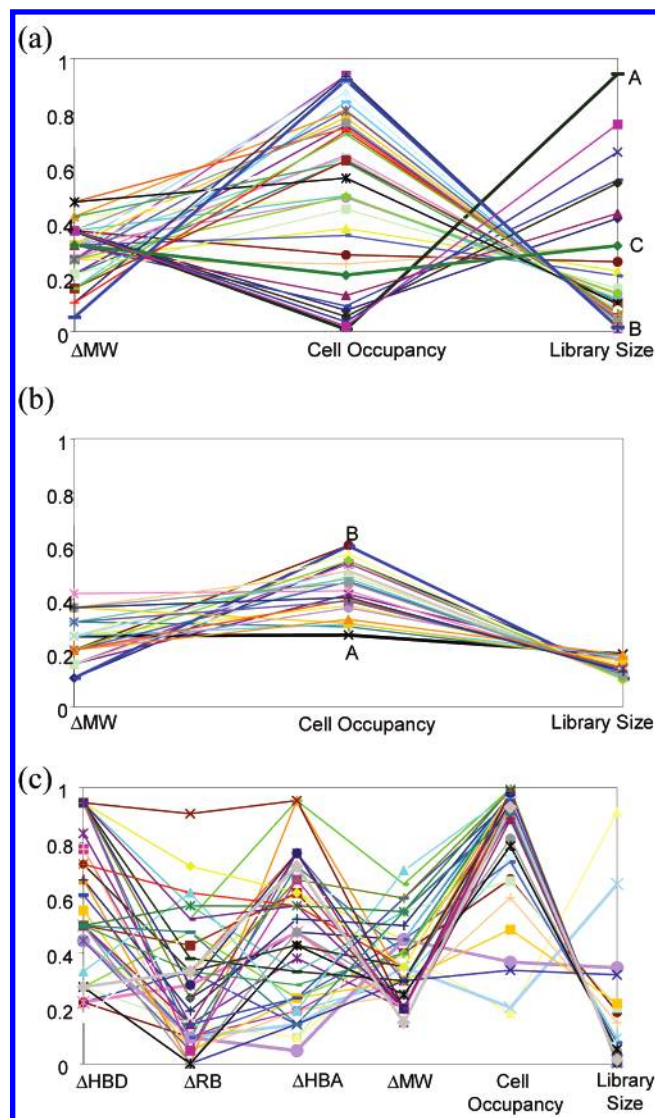


Figure 12. (a) Aminothiazole libraries found when simultaneously optimizing library size, diversity and molecular weight profile. (b) Solutions found when a size constraint (200–400) is placed on the run. (c) Aminothiazole libraries found when simultaneously optimizing library size, diversity, and profiles of molecular weight; rotatable bonds; hydrogen bond donors and acceptors.

Table 2. Example Libraries Extracted from the Solutions Shown in Figure 12(a)

| solution ^a | MW | cells occupied | library size |
|-----------------------|----|----------------|--------------|
| A | 8 | 361 | 1886 |
| B | 2 | 29 | 31 |
| C | 7 | 289 | 627 |

^a Solutions A and B represent libraries with extreme values in the objectives. Solution C represents a compromise solution.

that represent the extremes in size and diversity are shown in Table 2 as solutions A and B. Solution C shows the values of one particular compromise solution. The run took 6 min CPU time.

A further run was carried out in which library size was constrained to be within 200 to 400 products. The other parameters for the run were the same. The solutions are shown in Figure 12(b). This time 25 solutions were found that satisfy the constraints. The range of values for the objectives are 200–390 for library size; 166–267 for cell

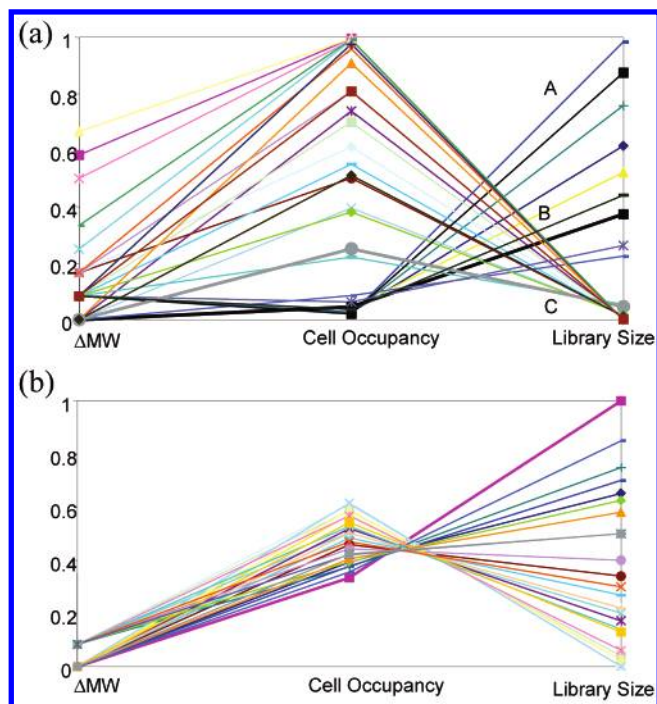


Figure 13. (a) Benzodiazepine libraries found when simultaneously optimizing library size, diversity, and molecular weight profile. (b) The search is directed to find libraries consisting of 500–1500 compounds by applying a size constraint (note that the size scale is expanded relative to (a)).

Table 3. Example Libraries Extracted from the Solutions Shown in Figure 12(b)

| solution | MW | cells occupied | library size |
|----------|----|----------------|--------------|
| A | 6 | 267 | 390 |
| B | 3 | 166 | 252 |

occupancy; and 3 to 9 for molecular weight profile. Table 3 shows the objective values for two libraries represented in Figure 12(b). The run took 3 min CPU time.

Increasing the number of objectives results in a corresponding increase in the number of solutions found. Figure 12(c) shows libraries found when simultaneously optimizing on size, diversity, and profiles of molecular weight, rotatable bonds, hydrogen bond donors and hydrogen bond acceptors, where the aim is to optimize each physicochemical property so that a drug-like distribution is achieved, using WDI as a reference. A population size of 250 was used, and a relatively high niche radius of 25% was applied in order to ensure that solutions were found that were close to the extremes in all objectives. While an unconstrained run such as this may be useful in order to see the relationships between the different objectives, in practice, it would be more useful to narrow down the search space through the use of constraints.

Finally, the effective of optimizing diversity, size, and molecular weight profile was investigated for the benzodiazepine library. Figure 13(a) shows the results found with a population of 250 and a niche radius of 10%. A total of 29 solutions were found with objectives values in the following ranges: 4 to 978 for cell occupancy (the virtual library of 250K compounds occupies 997 cells); 4 to 68 850 for library size; and 11 to 19 for molecular weight profile (the values have been scaled in a similar way to that described for Figure 12). A further run was performed with library size constrained to between 500 and 1500 with the results shown in

Figure 13(b). MW and diversity are scaled as before. This time the size scale has been expanded with a value of 1 for representing 1500 and a value of 0 representing a library of 500 products. Again, it can be seen that constraining the search can lead to a greater density of solutions within a given range.

CONCLUSIONS

This paper addresses a major issue in library design, namely how to efficiently optimize the library size (number of products) and configuration (number of reagents at each position) simultaneously with other properties such as diversity and drug-like physicochemical property profiles. This has been achieved within the MOGA framework of MoSELECT.II whereby an entire family of solutions is found in a single run without the need of an arbitrary weighting scheme to relate the competing objectives. Each solution represents a compromise in the objectives, and in the absence of any additional information they are all equally valid. MoSELECT.II has been shown to be effective in generating the entire tradeoff surface between size and diversity for two different combinatorial libraries. A user can then make an informed choice on what represents an appropriate compromise solution, for example, a library might be chosen from the region of the curve where the gain in diversity begins to flatten off with increasing size.

Practical considerations may mean that constraints are imposed on library design, for example, library size, the number of reactants used, and the efficiency with which the synthesis robot can be programmed. We have shown that such constraints can be handled effectively within MoSELECT.II.

Finally, it has been shown that the approach can be extended to include additional objectives such as drug-like physicochemical property profiles, and in all cases a range of different compromise solutions has been found. The results have been displayed using parallel coordinates graphs. While these graphs provide a useful way of visualizing tradeoffs in multiobjective space and can be used to identify appropriate compromise solutions, it can be difficult to compare the results across different runs of the MOGA. Recently, a number of methods have been suggested both for measuring the performance of MOEAs and for measuring the distribution of solutions along the tradeoff surface,^{24,25} and the application of these methods to library design is currently under investigation.

ACKNOWLEDGMENT

The authors would like to thank Peter Fleming, Peter Willett and Robin Purshouse for valuable discussions. This work is funded by GlaxoSmithKline and BBSRC via an industrial CASE studentship. The Krebs Institute for Biomolecular Research is a designated Biomolecular Sciences Centre of the Biotechnology and Biological Sciences Research Council.

REFERENCES AND NOTES

- (1) *Combinatorial Library Design and Evaluation – Principles, Software Tools and Applications in Drug Discovery*; Ghose, A. K., Viswanathan, V. N., Eds.; Marcel Dekker: New York, 2001.

- (2) Valler, M. J.; Green, D. Diversity Screening Versus Focused Screening in Drug Discovery. *Drug Discovery Today* **2000**, 5, 286–293.
- (3) Lipinski, C. A.; Lombardo, F.; Dominy, B. W.; Feeney, P. J. Experimental and Computational Approaches to Estimate Solubility and Permeability in Drug Discovery and Development Settings. *Adv. Drug Delivery Rev.* **1997**, 23, 3–25.
- (4) Oprea, T. I.; Davis, A. M.; Teague, S. D.; Leeson, P. D. Is There a Difference Between Leads and Drugs? A Historical Perspective. *J. Chem. Inf. Comput. Sci.* **2001**, 41, 1308–1315.
- (5) Hann, M. M.; Leach, A. R.; Harper, G. Molecular Complexity and its Impact on the Probability of Finding Leads for Drug Discovery. *J. Chem. Inf. Comput. Sci.* **2001**, 41, 856–864.
- (6) Good, A. C.; Lewis, R. A. New Methodology for Profiling Combinatorial Libraries and Screening Sets: Cleaning up the Design with HARPick. *J. Med. Chem.* **1997**, 40, 3926–3936.
- (7) Martin, E. J.; Crichlow, R. W. Beyond Mere Diversity: Tailoring Combinatorial Libraries for Drug Discovery. *J. Comb. Chem.* **1999**, 1, 32–45.
- (8) Gillet, V. J.; Willett, P.; Bradshaw, J.; Green, D. V. S. Selecting Combinatorial Libraries to Optimise Diversity and Physical Properties. *J. Chem. Inf. Comput. Sci.* **1999**, 39, 169–177.
- (9) Zheng, W.; Hung, S. T.; Saunders, J. T.; Seibel, G. L. PICCOLO: A Tool for Combinatorial Library Design via Multicriterion Optimization. In *Pacific Symposium on Biocomputing 2000*; Atلمان, R. B., Dunkar, A. K., Hunter, L., Lauderdale K., Klein, T. E., Eds.; World Scientific: Singapore, 2000; pp 588–599.
- (10) Brown, J. D.; Hassan, M.; Waldman, M. Combinatorial Library Design for Diversity, Cost Efficiency, and Drug-Like Character. *J. Mol. Graph. Model.* **2000**, 18, 427–437.
- (11) Gillet, V. J.; Khatib, 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.
- (12) Gillet, V. J.; Willett, P.; Fleming, P.; Green, D. V. S. Designing Focused Libraries Using MoSELECT. *J. Mol. Graph. Model.* **2002**, 20, 491–498.
- (13) UK Patent Application No. 0029361.
- (14) Agrafiotis, D. K. Multiojective Optimisation of Combinatorial Libraries. *IBM R. Res., Dev.* **2001**, 45, 545–566.
- (15) Fonseca C. M.; Fleming P. J. An Overview of Evolutionary Algorithms in Multiobjective Optimisation. *Evolutionary Comput.* **1995**, 3, 1–16.
- (16) Nicolotti, O.; Gillet, V. J.; Fleming, P. J.; Green, D. V. S. Multiobjective Optimization in Quantitative Structure–Activity Relationships: Deriving Accurate and Interpretable QSARs. *J. Med. Chem.* **2002**, 45, 5069–5080.
- (17) Handschuh, S.; Wagener, M.; Gasteiger, J. Superposition of Three-Dimensional Chemical Structures Allowing for Conformational Flexibility by a Hybrid Method. *J. Chem. Inf. Comput. Sci.* **1998**, 38, 220–232.
- (18) Goldberg, D. E. *Genetic Algorithms in Search, Optimization and Machine Learning*; Addison-Wesley Publishing Company Inc.: 1989.
- (19) Brown, R. D.; Martin, Y. C. Designing Combinatorial Library Mixtures Using a Genetic Algorithm. *J. Med. Chem.* **1997**, 40, 2304–2313.
- (20) 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.
- (21) Fonseca, C. M.; Fleming, P. J. Multiobjective Optimisation and Multiple Constraint Handling with Evolutionary Algorithms – Part I: A unified formulation. *IEEE Trans. Systems, Man, Cybernetics* **1998**, 28, 26–37.
- (22) Cerius² is available from Accelrys Inc., 9685 Scranton Road, San Diego, CA 92121.
- (23) The World Drug Index is available from Derwent Information, 14 Great Queen St., London WC2B 5DF, UK.
- (24) Fonseca, C. M.; Fleming, P. J. On the Performance Assessment and Comparison of Stochastic Multiobjective Optimizers. In *Parallel Problem Solving From Nature V*; Voigt, H. M., Ebeling, W., Rechenberg, I., Schwefel, H. P., Eds.; Springer: Berlin, Germany, 1995; pp 584–593.
- (25) Coello Coello, C. A.; van Veldhuizen, D. A.; Lamont, G. B. *Evolutionary Algorithms for Solving Multi-Objective Problems*; Kluwer Academic Publishers: New York, 2002.

CI0255836