

Predicting Oral Druglikeness by Iterative Stochastic Elimination

Anwar Rayan,* David Marcus, and Amiram Goldblum*

Molecular Modeling and Drug Design Lab and the Alex Grass Center for Drug Design and Synthesis,
Institute of Drug Research, The Hebrew University of Jerusalem 91120 Israel

Received November 8, 2009

Integration of computational methods in the early stages of drug discovery has been one of the key trends in the pharmaceutical industry. Starting with high quality drug candidates should ultimately minimize clinical attrition rates and give rise to higher success rates. In this paper, we present a novel approach for indexing oral druglikeness of compounds. With the Iterative Stochastic Elimination (ISE) Algorithm, we distinguish between orally available drugs and nondrugs by generating sets of optimized descriptors' ranges, each set constituting a "filter". We delineate in this paper how to produce an ensemble of best k-descriptor sets out of the huge number of possibilities, and how to construct a "filter bank" that retains diverse filters by clustering. Finally, we define the "orally bioavailable drug-like" character of individual molecules by combining the filters into an "Orally Bioavailable Druglike Index" (OB-DLI) which may be used to prioritize molecules in databases and discuss its uses in several potential applications. The predictive power with sets of 4–6 descriptors is high (i.e., one filter of 5 descriptors retrieved 81% true positives and >77% true negatives). Thus, OB-DLI has advantages over binary decisions (that use only one filter) not only in raising discriminative power but also in ranking drug candidates according to their chance to be successful oral drugs. We demonstrate the ability of our approach to discover molecular entities with the required property, orally bioavailable drug likeness, that are structurally dissimilar to those of the training set. Comparison of this ISE application to some of the current main methods for classification reveals that our approach has >13% improvement in the Matthews Correlation Coefficient, which measures the success of identifying true and false positives and negatives.

INTRODUCTION

Computational approaches for indexing oral druglikeness of molecules are highly important in drug discovery. Besides scientific interest, such methods would assist in the needs of reducing costs and time in drug development. For example, the number of compounds used in High Throughput Screening (HTS) experiments could be reduced. Furthermore, in combinatorial chemistry libraries,^{1–3} such an index could aid in selecting molecules to be drug candidates.

A set of simple rules based on statistics has become the landmark for predicting drug character. Lipinski's "rule of five" (ROF)^{4,5} is widely used to determine the potential bioavailability of molecules. In many cases, it is also employed to make decisions in drug development, in acquiring molecules for HTS, and in other applications. However, Lipinski's criteria are not useful for separating druglike compounds from nondruglike ones.⁶

Computational methods have been frequently applied to discern the "pharmacokinetic problem".⁷ Obviously, for orally administered drugs, the oral availability is of prime predictive interest. Drug absorption through the intestinal tract and its subsequent distribution in the body has an important role in determining its efficacy in treating diseases. Other aspects such as metabolism, excretion, and toxicity of the administered drugs are complex issues that are also frequently probed by computations^{8–11} including molecular

properties that are relevant to bioavailability such as solubility^{12–16} and lipophilicity.^{17–21}

Most predictions are "knowledge based", requiring information from databases that are not easily available.^{22–27} The information could possibly be acquired through carefully designed clinical studies. Predicting the total pharmacokinetic fate of a drug candidate in humans will remain elusive for some time, and models that combine the different characteristics cannot assume additivity of the ADME (absorption, distribution, metabolism, excretion) characteristics. However some models were developed in order to have a closer correspondence between experiments and real life;^{28–30} i.e., Caco-2 cell experiments have been used to predict intestinal absorption.^{31,32} It should however be remembered that failures of compounds are not due entirely to their inappropriate ADME profiles.³³

At this moment, it may thus be beneficial to improve our ability to predict the probability of a molecule to become a drug, in an "informatics" approach rather than a mechanistic one. An ability to distinguish between orally bioavailable drugs and others has important consequences for prioritization in synthesis and biological testing. It could be used for constructing molecular libraries, for increasing the efficiency of virtual docking and virtual screening experiments, and as a guide to design or modify drugs that will have appropriate molecular features.

A debate^{8,33–36} has developed in the past decade regarding the applicability of Lipinski's ROF for oral drug likeness, and several papers proposed to refine and reassess the

* Corresponding author e-mail: anwararrayan@gmail.com; amiram@vms.huji.ac.il.

determination of potential oral drugs.^{37–41} One of the notable contributions was made by Veber and co-workers who proposed two additional rules for predicting rat oral bioavailability:⁴² 10 or fewer rotatable bonds and a polar surface equal to or less than 140 Å² are required for a molecule to be bioavailable. However, a recent paper suggested that these rules are hardly applicable to human oral bioavailability, and the ability to predict bioavailability by simple thumb rules has been criticized.⁴³ Following that approach, many papers tried to extend the definition of oral bioavailability and to develop methods which combine that concept with “drug-likeness”.^{6,44–52} Several ADME contributions are optimized together to discriminate “drugs” from “non-drugs” by looking at different compound libraries. Other papers shift the debate from druglikeness to leadlikeness which is more practical for designing screening libraries with reduced complexity.^{53–59} Analysis shows that leadlike compounds have a lower molecular weight, less ring systems, less rotatable bonds, and a lower logP. Others suggest to use methods to describe screening libraries for other desired properties in the drug discovery process, such as CNS passage, promiscuity, and metabolites.^{60,61} Although it seems that these methods are widely applied in drug development pipelines, recent surveys show that newly developed drugs have different trends^{35,36,62} than those suggested by the informatics analysis. Still, it is important to improve the methods for discriminating between orally bioavailable drugs and nondrugs, to which we contribute in this paper. We thus present an improved method for indexing oral bioavailability of drugs and hope that it will be used for more accurate *in silico* examination of a molecule’s potential to become an oral drug.

METHODS

We investigate the distinction between orally bioavailable drugs and others by using a novel optimization method that we developed in the past few years.^{63,64} Here, we describe for the first time its application to solve a problem in the area of chemoinformatics. Iterative Stochastic Elimination (ISE) is a general optimization method that finds best solutions to complex combinatorial problems that are functions of many variables. [Abbreviations: Iterative Stochastic Elimination (ISE), Available Chemical Database (ACD), Comprehensive Medicinal Chemistry (CMC), MDL Drug Data Report (MDDR), Rule of Five (ROF), High Throughput Screening (HTS), ZINC - a free database of commercially available compounds for virtual screening (ZINC)]. The orally bioavailable character of molecules depends on variables such as solubility, acidity, molecular weight, molecular surface, types of atoms, and others. These variables or “descriptors” are useful to describe both orally bioavailable drugs and orally bioavailable nondrugs. ISE has been used to optimize the ranges of variables’ values that maximize the differences between the two sets. A function (Matthews’ Correlation Coefficient, MCC⁶) is used to score the optimizations on the basis of the numbers of true and false positives and negatives. ISE leads to the formation of a set of filters, each consisting of a set of variable value ranges. Subsequently, individual molecules may be examined and scored for their ability to pass the filter set. This score is our oral bioavailability drug-like index (OB-DLI) for individual molecules, which reflects a molecule’s chance to belong to

the database of orally bioavailable drugs. Of course, one must finally decide if a molecule will or will not be tested experimentally, which is a “binary decision”. But, the scalar presentation with OB-DLI allows one to make more elaborate decisions based on the position of a molecule along a scale, as well as on enrichment factors that can be calculated at different levels of the OB-DLI, as shown below. Our OB-DLI is presented in terms of properties of the entire molecule and is a result of accumulating results from several, sometimes many, filters.

Indexing the oral bioavailability of compounds by using data of oral drugs alone does not discriminate well between oral drugs and others. In our study, the database of orally bioavailable drugs has been assembled from CMC and MDDR, but only molecules that obey Lipinski’s rule of five and are leadlike according to Oprea’s rules^{5,54} were included. A database of nondrugs is usually a set of marketed molecules that is a less specific resource.⁶⁵ It should obey Lipinski’s rule of five and leadlikeness, according to Oprea, but it contains mostly molecules that have not been tested for treating some disease. It has been suggested that, in the nondrugs databases, between 70–90% of the molecules do not have the potential to become orally bioavailable drugs or are “true” nondrugs, while the rest could be drugs at a certain concentration.^{54,66–69} While taking account of that possibility, such databases have been widely used for distinguishing drugs from nondrugs, by diverse methods of prediction and optimization. The intricacy of biological activity and the large number of potential physical descriptors of molecular structure require special technologies that can deal with such complex problems, but it is not clear if these problems can be solved to full satisfaction.

The ISE algorithm was applied to perform simultaneous selection and optimization of the ranges of k-descriptor sets in order to distinguish between orally bioavailable drugs (OBD) and other chemicals, presumably orally bioavailable nondrugs (nOBD). Application of ISE for side chains placement, peptide/protein modeling, and docking was described elsewhere.^{63,64,70,71}

Here, we describe an application of ISE optimization for solving classification problems, a usage that requires an approach that does not depend on energy-based scoring methods. A flowchart of the basic processing steps for differentiating between OBD and others and the formation of the Orally Bioavailable Drug Like Index (OB-DLI) on the basis of 2D descriptors is presented in Figure 1.

General Description of Iterative Stochastic Elimination for Optimizing Classifications. We utilize ISE for distinguishing optimally between a database that has a specific desired property (“positives”) and another database that lacks that property (“negatives”). Problem solving proceeds in a few general steps:

The “preparations step”:

1. Calculate the values of the physicochemical descriptors of interest for all of the molecules in the two databases, using appropriate software. In the present study, MOE was used for calculating the values of descriptors for molecules.

2. Construct histograms of the descriptors of the databases of positives and negatives, in order to determine the lower and upper ranges of each variable and for deciding upon their internal division into “bins” for the subsequent computations.

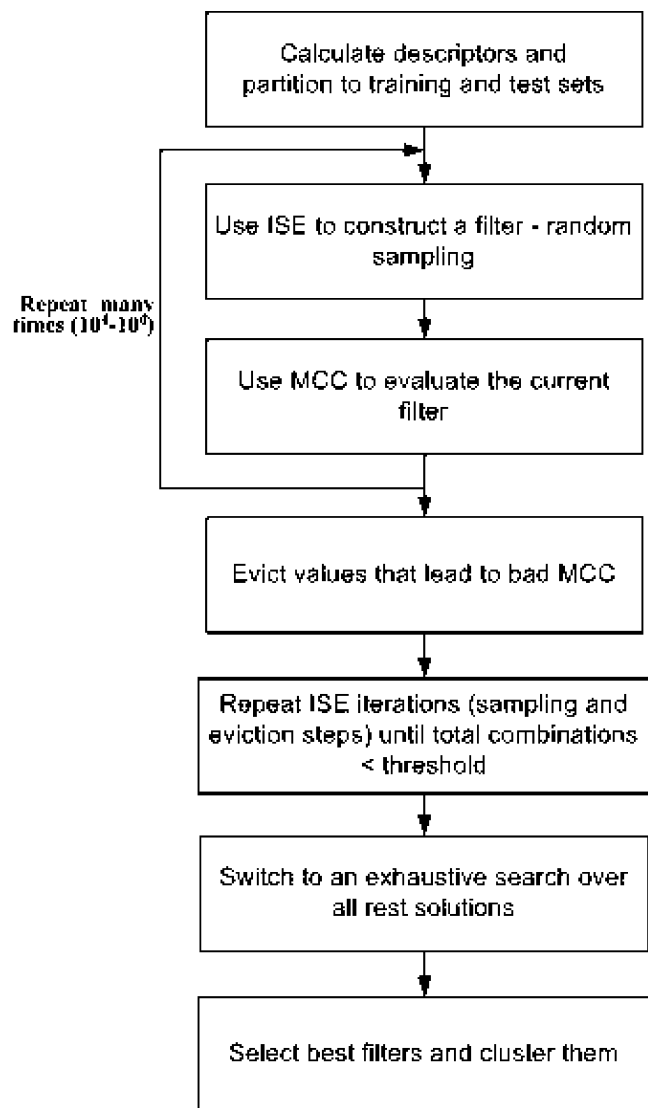


Figure 1. A flowchart of the ISE approach for differentiation between orally bioavailable drugs and nondrugs.

3. Partition the two databases of positives and negatives randomly into a training set and a test set. From this step on, only the training set is used.

The “sampling step”:

1. Pick a set of k -dimensional descriptors for optimizing their ranges.

2. Pick randomly two values for each molecular descriptor: one out of its lower range, the other out of its upper range. Together, they form a randomly chosen “range” of a descriptor. A “filter” is the set of all randomly picked descriptor ranges, for the descriptors chosen for examination, as many as is wished. The training set molecules, in the form of their descriptor values, are “passed” through such a filter. Molecules of the data set of positives must be accommodated by all values of a filter. If a molecule does not pass even a single range (of one descriptor, i.e., if a molecule’s weight is 450 while the current filter includes the range 200–380 for MW), the molecule “fails” as a positive candidate and is assigned a false negative. If it passes all ranges in a filter, it is a true positive. If the molecule’s origin is in the training set of negatives and it failed a single range or more in a filter, it is assigned a true negative. If a molecule from the negatives passes all filters, it is a false positive.

3. A scoring function, the Matthews Correlation Coefficient (MCC), is used to measure the ability to differentiate between positives and negatives in the training set, on the basis of the current filter and the results for positives and negatives from the two databases.

The “evictions step”:

After repeating sampling steps 2–3 a large number of times ($\sim 10^5$), a virtual histogram is used to examine which of the values in each variable (upper and lower limits of any descriptor are such variables) may be discarded with lowest risk. Briefly, in this process, values of the lower or upper range of a variable may be identified as contributing to results with the worst scores and should subsequently be eliminated.

The “iterations step”:

Repeat the sampling and eviction steps until a predetermined number of combinations, regularly $\sim 10^6$ – 10^7 combinations, is reached from the much larger set of initial combinations, 10^{10} – 10^{30} for sets of 4–9 descriptors.

The “solutions step”:

1. Calculate all remaining combinations exhaustively.

2. Cluster the results, and pick the top set of “best filters” (on the basis of MCC values).

The final output may be composed of hundreds of best filters.

Converting Drug/Nondrug Databases to Orally Bioavailable Databases. CMC and MDDR databases represent drug space, while the contents of the ACD database were considered as representative of the nondrugs space. From both databases, we cleaned pesticides, UV-screens, etc.^{52,69} as well as undesired atomic elements (those that are different than C, S, O, N, P, H, Si, Cl, Br, I, or F). The databases were further “washed” of counterions and solvent molecules. Despite the possible source for confusion, we used single tautomers, and the ionization state was determined on the basis of the chemical functions.^{72–74} Compounds from all databases were kept only if they obeyed Lipinski’s rule of five and Oprea’s criteria for lead-likeness. The filtered CMC and MDDR databases are called orally bioavailable databases. An overall 184 2D descriptors were computed by MOE 2008.10; 146 noncorrelated descriptors were considered.

TRAINING SETS AND TEST SETS

The orally bioavailable CMC (OB-CMC) database is employed as the basis for the training data set of orally bioavailable drugs (OBD), and the orally bioavailable ACD is employed as the training set of orally bioavailable nondrugs (nOBD). The OB-CMC, composed of 6776 molecules was randomly partitioned into two portions of 3/4 and 1/4 of the molecules, thus providing 5082 molecules in the training set and 1694 in the test set. Three random subsets from the OB-ACD were picked after cleaning and filtering as described above: each one of the three subsets includes 9128 molecule; the first was used for training while the two others served as test sets. Three random sets composed of 6070 molecules each were prepared from OB-MDDR by the procedure described above, as well as two sets composed, each, of 9300 molecules from OB-ZINC. The MDDR sets were used as OBD test sets, in addition to the OB-CMC test set, while the ZINC sets and the ACD sets were used as nOBD test sets.

Construction of Filters. Only noncorrelated descriptors were considered. At first, sets of descriptors composed of 3–6 descriptors were selected by applying ISE (using a single range per descriptor; that range was found to have highest MCC value). Subsequently, ranges were optimized by applying ISE.

Scoring. The constructed filter (comprised of ranges of n descriptors) was applied to the molecules in each of the two training sets (OBD and nOBD) to calculate the value of the scoring function, the Matthews Correlation Coefficient (MCC), eq 1.

$$\text{MCC} = \frac{(PN) - (P_f N_f)}{\sqrt{(N + N_f)(N + P_f)(P + N_f)(P + P_f)}} \quad (1)$$

P and N are the percentages of true positive and true negative predictions, while P_f and N_f are the percentages of false positives and false negatives, respectively. True positives are OBD that are identified by the current filter as OBD. False positives are nOBD identified as OBD. False negatives are OBD identified as nOBD, and true negatives are nOBD correctly identified as such.

The best possible value for MCC is 1.0 (for a perfect prediction of both P and N , with no P_f and no N_f), and the worst possible value is -1.0 (a completely erroneous prediction). An OBD molecule can either be counted as “positive” (P) or as false negative (N_f). In order to be counted as P , such a molecule must have descriptor values that fall in the range of all the descriptors’ ranges of a given filter. It is counted as N_f if even a single value does not comply with the filter’s descriptor range. nOBD can only be negatives or false positives.

Constructing the Best Filters Set by Clustering. A set of “best filters” includes the 1000 filters with best MCC, following clustering of filters. In the exhaustive phase of applying ISE, many thousands of filters are formed, many of which are quite similar to each other, differing by a small variation of the range of a single descriptor. On the basis of our interest in the diversity of the model, we cluster by determining a threshold, of $X\%$ between two filters. If a new filter changes the definitions of less than $X\%$ of the whole sample, i.e., from OBD to nOBD or from nOBD to OBD, it belongs to the same cluster to which it was compared. In clustering, however, we do not search for a “representative filter” of each cluster but use the filter that gives the best MCC value. The set of best filters may include filters having different numbers of descriptors, different kinds of descriptors, or different ranges for the same set of descriptors.

CPU Time. From an initial set of 184 2D descriptors, less than 145 remain following the elimination of redundant descriptors and of descriptors with low variance. Computation times (on a single CPU) depend on the size of the filter and training sets’ sizes. For filters of 4 descriptor ranges, the total time (up to the formation of a set of best filters) is between 1–2 h (examining filters with 4 descriptors out of ~150). It is 2–3 h for 5-descriptor filters and 3–4 h for 6 descriptors. About 12 h are required for constructing the ISE-based final prediction model.

Combined Filters and Orally Bioavailable Drug Like Index (OB-DLI). A single molecule may be examined for its “OB-DLI” by using the best filter that was found in the exhaustive stage described above. Also, molecules may

be designed so that their descriptors’ values would be in the ranges of that best filter. However, we expect that an orally bioavailable drug-like molecule, one that has a higher probability for belonging to the OBD database, will also be able to “pass” other good filters, that are less successful than the best one, the “global optimum” with the top MCC. Such filters may be regarded as “local minima” that can be quite close to the global minimum but are different. Designing a new molecule by having its descriptor values pass a few or more filters is complicated, but picking molecules that have such character is easily performed. In fact, by using a set of filters, rather than the single “best” filter, the method benefits from the larger set of good filters and increases its potential to prioritize molecules. Such prioritization based on oral druglikeness is needed in subsequent high-throughput screening (HTS) of databases, in the selection of lead compounds, and in many other situations. Equation 2 is employed for calculating the orally bioavailable drug-like index (OB-DLI) and may include as many of the “best filters” (n) as desired.

$$\text{DLI} = \frac{\sum_{i=1}^n \delta_{Di} \frac{P_{Di}}{P_{NDi}} - \delta_{NDi} \frac{N_{NDi}}{N_{Di}}}{n} \quad (2)$$

By this equation, the Orally Bioavailable Drug Like Index (OB-DLI) for a molecule is determined on the basis of a set of n filters. The number n can range from a few to thousands. The value of the δ function δ_{Di} is 0 (zero) if the molecule is nOBD according to the currently calculated filter i , and $\delta_{Di} = 1$ if it is an OBD according to that filter. Similarly, the value of the δ function δ_{NDi} is 1 if it is nOBD according to the currently calculated filter, and 0 if it is an OBD according to that filter. P_{Di} is the percentage of OBD that is predicted to be “OBD” according to filter i (“true positives”), while P_{NDi} is the percentage of false positives, i.e., nOBD that is predicted to be OBD according to filter i . N_{NDi} is the percent of nOBD identified as such by the current filter, i.e., “true negatives”, while N_{Di} is the percent of OBD identified to be nOBD according to the current filter (“false negatives”).

Thus, P_{Di}/P_{NDi} may be regarded as an “efficiency factor” of filter i for the OBD, while N_{NDi}/N_{Di} is an “efficiency factor” for identifying nOBD. OB-DLI is thus composed by combining the “successful rate” or “OBD like efficiency” of prediction of OBD with “nOBD like efficiency” of prediction, for each molecule at each filter. Thus, identification of a molecule at each filter as an OBD receives a positive contribution, while its identification as nOBD at any filter receives a negative contribution. A shift to the right side of the scale (higher OB-DLI values) means improving the OBD-like character of a molecule and giving it a better chance to be a successful OBD candidate.

As the OB-DLI for a certain molecule is larger, the confidence that this molecule could be an OBD is greater. In dealing with a set that is a mixture of OBD and nOBD, the higher the OB-DLI cutoff that is selected, the fewer OBD will “pass” that cutoff. The lower the OB-DLI cutoff, the more OBD will pass but also more nOBD. The optimal cutoff is then preferably sought, by which more OBD and fewer nOBD pass. For example, when used with HTS, the process may optionally begin by testing molecules with a higher OB-

Table 1. MDDR/CMC vs ACD—Set of Three Descriptors, Selection, and Optimization of Ranges by ISE^a

filter	MCC*	%TP/ MDDR	%TN/ ACD	BCUT_SMR_3	Q_VSA_POS	Opr_brigid
1	0.5696	84.85	71.61	2.69–3.20	≥128.4	≥7
2	0.5689	86.59	69.47	2.69–3.23	≥128.4	≥7
3	0.5687	84.89	71.47	2.69–3.22	≥126.5	≥9
4	0.5683	84.34	72.06	2.69–3.22	≥134.5	≥7
5	0.5682	87.23	68.59	2.69–3.22	≥119.8	≥7

^a MCC = Matthews correlation coefficient; TP = percent true positives; TN = percent true negatives. BCUT_SMR_3 is a descriptor of atomic contribution to molar refractivity using the Wildman and Crippen SMR method instead of partial charge. Q_VSA_POS is total positive van der Waals surface area. Opr_brigid is the number of rigid bonds.⁵⁴

Table 2. Y-Scrambling Test: OB-CMC Mixed with OB-ACD

set	MCC value
1	0.0046
2	−0.002
3	0.005

DLI (a fraction that is enriched with orally bioavailable drugs and its molecules have a higher chance to be successful oral drugs).

RESULTS AND DISCUSSION

We formed a diverse “filter bank” into which we introduced 369 diverse *k*-descriptors’ filters with $77 \geq k \geq 3$ descriptors.

Sets of Three Descriptors. Best filters with 3 descriptors are depicted below. The five unique filters described in Table 1 have different ranges, and at least 3% of the molecules in the databases change from OBD to nOBD or vice versa between filters.

Y-Scrambling Test. We scramble molecules from the OBD databases (CMC) with those from the nOBD database (ACD) by randomly picking a similar number from each of the databases, and we form two “half false” databases of OBD and nOBD from the scrambled results. Database A is considered as if it is a database of OBD, while database B is considered to represent nOBD. Each of the two databases A and B contains half of the OBD as well as half of the nOBD. Both databases (A and B) were partitioned into training sets and test sets as described in the Methods section. The process of assembling A and B was repeated three times, so that three A and B pairs of “databases” were constructed from OB-CMC + OB-ACD. There should be no difference between databases A and B, and thus a MCC value close to zero is expected. Results of Y-scrambling tests are shown in Table 2.

Orally Bioavailable Drug Like Index. Filters for the construction of OB-DLI were obtained by picking and optimizing ranges of different sets of descriptors. The best filters are clustered by a dissimilarity requirement of at least 3%. The OB-DLI is calculated according to eq 2. The results for average OB-DLI values are presented in Table 3.

Figure 2 demonstrates the ability to separate between molecules that belong to the “OBD” database and to the “nOBD” database. It seems that OB-DLI could be useful for determining the chance of any molecule to be an OBD or nOBD and ranking that chance relative to other molecules.

Table 3. Average OB-DLI Values of a Few Databases

database (# test sets)	average OB-DLI	# molecules in each set
ACD (2)	−0.91	9128
CMC (1)	0.86	1694
MDDR (3)	1.32	6070
ZINC (2)	0.46	9300

A higher OB-DLI means a better chance to be an orally bioavailable drug. This ability to transform a qualitative classification into quantification is due to the use of multiple filters that are a natural result of the ISE process. However, oral druglikeness is statistically based: the filters that constitute the basis for quantification emerge from examining a large set of molecules, and the value of MCC reflects the percentage of predictive success of each filter. Although OB-DLI characterizes a single molecule, its statistical nature and the chance for erroneous prediction for a single molecule should be kept in mind.

One way to appreciate the success of OB-DLI in distinguishing OBD from nOBD is to pick randomly molecules from each database and to position them on the same plot, as shown in Figure 3. In that figure, 300 molecules from OB-ACD and from OB-MDDR were randomly picked and positioned along the *x* axis, with values of 1–300, as random compound numbers. The *y* axis is the OB-DLI values. It is obvious from Figure 3 that most compounds of OB-ACD have lower OB-DLI values, while those from OB-MDDR have much higher values, with most of its compounds above

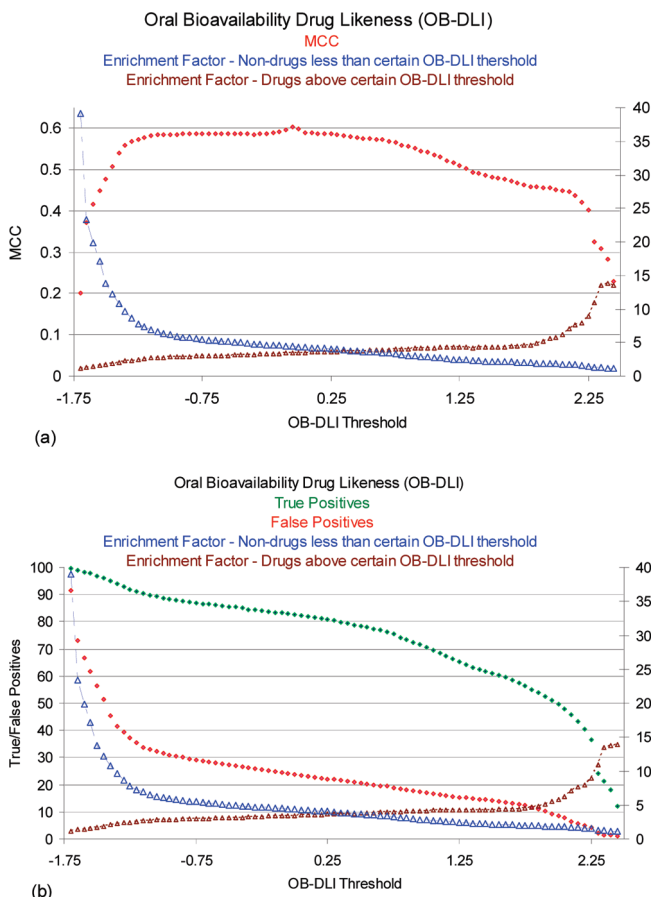


Figure 2. (a) Change of MCC along the OB-DLI axis and (b) enrichment factors along that axis. Molecules that have OB-DLI equal or above a certain threshold are considered potentially orally bioavailable drugs, while others are nondrugs.

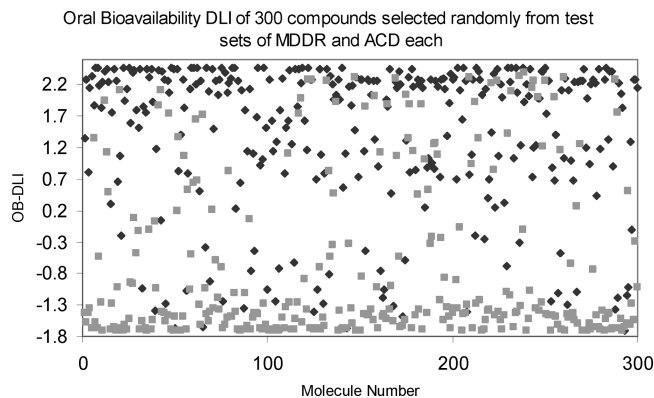


Figure 3. DLI values of 300 randomly selected molecules from two databases: ACD (gray squares) and MDDR-clinical (black rectangles).

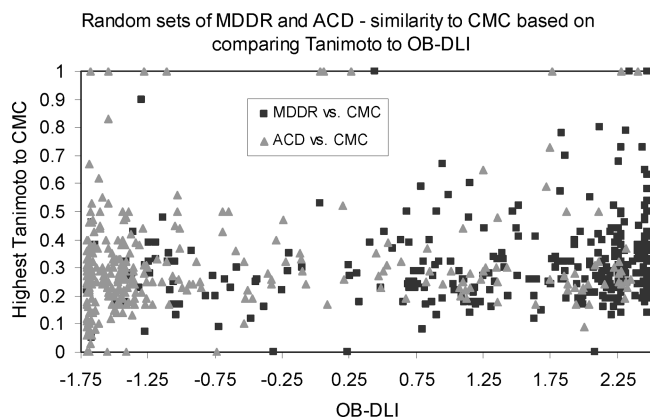


Figure 4. Similarity of MDDR (black squares) and ACD (gray triangles) to CMC vs OB-DLI values.

an OB-DLI of 0.0. It may also be seen that OB-DLI may be used as a decision making tool, in that a line drawn parallel to the x axis at some higher OB-DLI value excludes most of the nOBD and enriches the remaining set with many orally bioavailable drug candidates.

It is also worth it to indicate that some of the gray squares (molecules from the non-OBD database, ACD) that have higher OB-DLI values could be candidates for orally bioavailable drugs if they are active at some disease target.

Discovering Dissimilar Compounds by the Indexing Method. Since the indexing method is based on wide descriptor ranges as filters, rather than on direct descriptor quantities such as in other methods, we checked the performance of the ISE indexing technique for discovering structurally dissimilar compounds. Figure 4 demonstrates this ability in screened databases when using the index of oral bioavailability utilizing our OB-DLI prediction model. Black squares display the degree of similarity of 300 drug candidates, selected randomly from MDDR, to the druglike training set CMC, while gray triangles display similarity of the nondrugs, randomly selected from ACD, to the CMC training set. This figure suggests that our technique, based on ISE with multiple filters, is capable of discovering chemical entities which are significantly dissimilar in structure compared to the training set. It is also clear that molecules with a high OB-DLI, the ones that could be candidates for orally available drugs, are much different than the ones from the training set of CMC. This ability can be beneficial for finding new structural entities within a

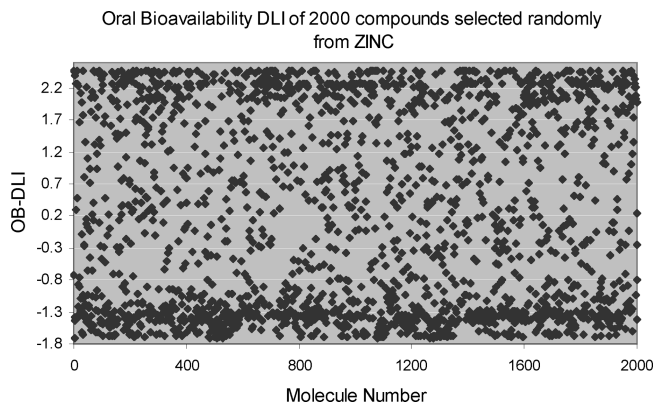


Figure 5. Indexing a random set from the ZINC database for oral druglikeness.

relatively larger chemical space defined by the filter ranges, rather than by models which are more connected to the known structures of the training set that define a relatively narrower space.

The sets of “positives” (OB-CMC) and “negatives” (OB-ACD) are not of equal size. In most applications of ISE to classification, the positives are supposed to be a small set of active molecules out of a huge set of nonactives (such as, e.g., inhibitors of an enzyme), in proportions similar to those found for hits in high throughput screening. Thus, the training set should include a number of “positives” that is much smaller than the number of “negatives”. In the scoring function, however, each of the true and false positives and negatives is represented by their percentages, thus balancing the huge difference between the data sets’ sizes.

OB-DLI: A Focusing/Enrichment Tool. The limited ability to predict the OB-DLI of a single molecule is alleviated when considering sets of molecules. Decision making on a set to be synthesized or purchased has a higher predictive stability with an increase in the number of molecules of that set. However, the enrichment factor is at its top values in the higher positive end of OB-DLI, where there are fewer molecules. The concept of “focusing” is quite parallel to the enrichment.

Thus, OB-DLI may be employed as a decision making filter by itself. Enrichment may be improved with assigning a cutoff value of OB-DLI and keeping those molecules that have higher values. As an example, let us assume a cutoff value of $\text{OB-DLI} = 1.0$. In the case of our test sets, more than 58% of the MDDR OBD passes this “filter” as true OBD, while the amount of false positives from ACD is only about 13%. Thus, at this cutoff, the enrichment factor ($\text{EF} = \text{TP}/(1 - \text{TN})$) is about 4.5.

If a very large library with yet unknown OBD and non-OBD is to be screened, what values of cutoffs should be used and what is the quality of the portion of the database that will be retained?

The enrichment factor in some fractions, mainly those with very high OB-DLI, could exceed 13-fold (see Figure 2b). Utilizing the OB-DLI indexing technique for prioritizing oral bioavailable chemicals in the ZINC database for HTS is demonstrated in Figure 5. It is preferable to discard the low fraction from biological tests since it is highly enriched in chemicals with a low chance to become successful orally bioavailable drugs. By comparing molecules from the ACD and CMC, we find that nearly 4% of the ACD molecules

Table 4. Performance of ISE and Various Computational Approaches

method	MCC	TP	TN
artificial neural network (multilayerperceptron)	0.538	81.4	72.2
support vector classifier (SMO)	0.503	66.1	83.4
KNN method (Kstar)	0.457	65.1	80.1
decision Tree (J48)	0.424	64	78
RBF network	0.317	42.3	86.2
simple filter composed of 5 descriptors ^a (ISE)	0.58	81.0	77.3
OB-DLI (ISE)	0.61	83	78

^a Best filter composed of 5 descriptors: GCUT_PEOE_3 ranges between 2.45 and 3.33; Zagreb ranges between 74.01 and 288.00; PEOE_VSA_POS ranges between 115.33 and 501.44; Opr-brigid ranges between 7 and 42; a_acc ranges between 1 and 10.

are OBD or very similar to OBD, as calculated on the basis of 2D descriptors, and those were left as part of the OB-ACD. In the set with the top OBD values, there are thus more OBD molecules than those that are expected from the comparison of CMC to ACD. In addition, ACD contains a larger percentage of molecules that have never been examined and some are expected to be potential OBD. Therefore, the enrichment factor is probably higher than the one currently calculated on the basis of OB-DLI for OBD vs nOBD. Another enrichment factor is possible when a reduction in the size of the database is required for HTS and in other experiments. If we focus on the fraction with a low index, the true negative enrichment factor of such a fraction ($EF(TN) = TN/(1 - TP)$) could exceed 100 (see Figure 3), and a low value may be used to evict many compounds from the library.

ISE versus Other Approaches. The WEKA (Waikato Environment for Knowledge Analysis) package includes many machine learning techniques for classification purposes (<http://www.cs.waikato.ac.nz/ml/weka>).⁷⁶ Five methods have been utilized to model orally bioavailable druglikeness. The employed algorithms are the state-of-the-art learning techniques for classification: neural networks, support vector machines, decision trees, and the nearest-neighbor method. Table 4 compares the performance of ISE to the other five approaches for constructing an OB-DLI prediction model using orally bioavailable drugs/nondrugs data sets coded by 184 MOE 2D descriptors (CMC and a random set selected from ACD were used for training while MDDR and a random set from ACD were used as test sets). The efficacy of each prediction model is measured by its MCC value.

The accuracy of our models for indexing oral druglikeness has >13% improvement over the best result obtained by other classification techniques in our comparative study.

Sen et al. (*J. Chem Inf Model* **2006**, *46* (3), 1394–1401) described another method for indexing oral drug likeness of chemicals. The comparison of our results with their published work was not possible since they had no access to any drug or nondrug databases; thus, they could not report the discriminating power of their approach in detail.

The screening speed of this novel indexing technique is very high. The computation time for 1000 compounds takes a few seconds with the 2D-descriptors input of the compounds.

Oral Bioavailability and Drug Likeness. There seems to still be confusion as many mentions of Lipinski's ROF are considered to characterize drug-likeness, while Lipinski's data set for assessing the ROF was composed of orally

bioavailable drugs. As there was no data set for the orally bioavailable nondrugs, the ROF does not purport to distinguish between the bioavailability of drugs and that of nondrugs, which is the main subject of our research. To diffuse some of the confusion, we optimized the ranges of ROF variables (molecular weight, numbers of H donors and acceptors, and the calculated logP) for distinguishing between the two sets of OB drugs and OB nondrugs. The optimal filter has a MCC of 0.47 corresponding to discovering 93% of the TPs and only 49% of the true negatives. This new filter is composed of only 2 descriptors: $MW \geq 240$ and number of acceptors ≥ 1 .

One of the reviewers requested to comment on the issue of discovering toxic compounds as OB drugs, if chosen by the method described herein. The issue of toxicity is, in our opinion, ambiguous. On the one hand, drugs are molecules that clearly interact with biomolecules, proteins in particular. As it is difficult to aim at a single protein, drugs interact frequently with many similar proteins and might cause toxicity by affecting pathways that are different than the one aimed for. Toxicity is related to dose, which is a main reason for searching strong binders that may be administered in a smaller amount. In our study, the CMC drugs are assumed to be given in the range of clinical doses, so that their toxicity is reduced. Therefore, distinguishing between OB-CMC and OB-ACD introduces implicitly the differences in toxicity between these two sets. Excluding the inclusion of doses, the OB drugs of CMC are evidently nontoxic, while the OB nondrugs of ACD might have toxic effects, depending on their interactions with proteins and other biomolecules and their ability to cross cell membranes and the Blood Brain Barrier (BBB). The inclusion of nondrugs in an orally bioavailable molecules' set is a result of applying Lipinski's and Oprea's rules, which do not deal explicitly with toxicity.

ISE is currently applied to issues that are related to the effectiveness of drugs, including h-Erg toxicity (manuscript in preparation), general toxicity, leadlikeness, solubility, BBB passage, P450 isoenzymes' selectivity for substrates and inhibitors, and others.

CONCLUSION

We present a new and highly efficient tool for predicting oral druglikeness of molecules on the basis of indexing. It is a tool that should be utilized to construct a set of potential candidates with high values of oral druglikeness and is proposed to help in decision making in the drug discovery process. We have shown the advantages of ISE in predicting the drug likeness of test sets, compared to the most abundant current methods. Our approach combines oral bioavailability with druglikeness and is useful for discriminating between orally bioavailable drugs and nondrugs, with better performance than many other major techniques. Finally, the ISE model presented here demonstrates again the ability to discover molecules that are substantially different than those used for the training of the model. That is a result of using nonstructural criteria of descriptors.

ACKNOWLEDGMENT

We thank the reviewers for making suggestions that helped to improve the paper. We thank Dr. Dinorah Barasch for technical assistance with running the programs and Dr. Jamal

Raiyan for preparation of Figure 1. A.R. is a member of COST Action BM0608.

REFERENCES AND NOTES

- Matter, H.; Baringhaus, K. H.; Naumann, T.; Klabunde, T.; Pirard, B. Computational approaches towards the rational design of drug-like compound libraries. *Comb. Chem. High Throughput Screening* **2001**, *4* (6), 453–475.
- Lobanov, V. S.; Agrafiotis, D. K. Scalable methods for the construction and analysis of virtual combinatorial libraries. *Comb. Chem. High Throughput Screening* **2002**, *5* (2), 167–178.
- Tropsha, A.; Zheng, W. F. Rational principles of compound selection for combinatorial library design. *Comb. Chem. High Throughput Screening* **2002**, *5* (2), 111–123.
- 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* (1–3), 3–25.
- Lipinski, C. A. Drug-like properties and the causes of poor solubility and poor permeability. *J. Pharm. Toxicol. Methods* **2000**, *44* (1), 235–249.
- Frimurer, T. M.; Bywater, R.; Naerum, L.; Lauritsen, L. N.; Brunak, S. Improving the odds in discriminating “Drug-like” from “Non Drug-like” compounds. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (6), 1315–1324.
- Dearden, J. C. In silico prediction of ADMET properties: how far have we come. *Expert Opin. Drug Metab. Toxicol.* **2007**, *3* (5), 635–9.
- van de Waterbeemd, H.; Gifford, E. ADMET in Silico modelling: Towards prediction paradise. *Nat. Rev. Drug Discovery* **2003**, *2* (3), 192–204.
- Butina, D.; Segall, M. D.; Frankcombe, K. Predicting ADME properties in Silico: methods and models. *Drug Discovery Today* **2002**, *7* (11), S83–S88.
- Ekins, S.; Waller, C. L.; Swaan, P. W.; Cruciani, G.; Wrighton, S. A.; Wikel, J. H. Progress in predicting human ADME parameters in Silico. *J. Pharm. Toxicol. Methods* **2000**, *44* (1), 251–272.
- Falah, M.; Nassar, T.; Rayan, A. A simple approach discriminating cardio-safe drugs from toxic ones. *Bioinformation* **2009**, *3* (9), 389–93.
- Huuskonen, J.; Salo, M.; Taskinen, J. Aqueous solubility prediction of drugs based on molecular topology and neural network modeling. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (3), 450–456.
- Katritzky, A. R.; Wang, Y. L.; Sild, S.; Tamm, T.; Karelson, M. QSPR studies on vapor pressure, aqueous solubility, and the prediction of water-air partition coefficients. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (4), 720–725.
- Mitchell, B. E.; Jurs, P. C. Prediction of aqueous solubility of organic compounds from molecular structure. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (3), 489–496.
- Abraham, M. H.; Le, J. The correlation and prediction of the solubility of compounds in water using an amended solvation energy relationship. *J. Pharm. Sci.* **1999**, *88* (9), 868–880.
- Jorgensen, W. L.; Duffy, E. M. Prediction of drug solubility from Monte Carlo simulations. *Bioorg. Med. Chem. Lett.* **2000**, *10* (11), 1155–1158.
- Breindl, A.; Beck, B.; Clark, T.; Glen, R. C. Prediction of the n-octanol/water partition coefficient, logP, using a combination of semiempirical MO-calculations and a neural network. *J. Mol. Model.* **1997**, *3* (3), 142–155.
- Haerberlein, M.; Brinck, T. Prediction of water-octanol partition coefficients using theoretical descriptors derived from the molecular surface area and the electrostatic potential. *J. Chem. Soc., Perkin Trans. 2* **1997**, (2), 289–294.
- Hawkins, G. D.; Liotard, D. A.; Cramer, C. J.; Truhlar, D. G. OMNISOL: Fast prediction of free energies of solvation and partition coefficients. *J. Org. Chem.* **1998**, *63* (13), 4305–4313.
- Buchwald, P.; Bodor, N. Octanol-water partition: Searching for predictive models. *Curr. Med. Chem.* **1998**, *5* (5), 353–380.
- Tetko, I. V.; Tanchuk, V. Y.; Villa, A. E. P. Prediction of n-octanol/water partition coefficients from PHYSPROP database using artificial neural networks and E-state indices. *J. Chem. Inf. Comput. Sci.* **2001**, *41* (5), 1407–1421.
- Zheng, C. J.; Sun, L. Z.; Han, L. Y.; Ji, Z. L.; Chen, X.; Chen, Y. Z. Drug ADME-associated protein database as a resource for facilitating pharmacogenomics research. *Drug Dev. Res.* **2004**, *62* (2), 134–142.
- Goshorn, J. Information resources for the twenty-first century: Web-based databases and other resources developed by the Toxicology and Environmental Health Information Program of the USA National Library of Medicine. *Toxicol. Appl. Pharmacol.* **2004**, *197* (3), 225–225.
- Klopman, G.; Chakravarti, S. K.; Zhu, H.; Ivanov, J. M.; Saiakhov, R. D. ESP: A method to predict toxicity and pharmacological properties of chemicals using multiple MCASE databases. *J. Chem. Inf. Comput. Sci.* **2004**, *44* (2), 704–715.
- Renwick, A. G. Toxicology databases and the concept of thresholds of toxicological concern as used by the JECFA for the safety evaluation of flavouring agents. *Toxicol. Lett.* **2004**, *149* (1–3), 223–234.
- Rendic, S.; Kitajima, M.; Ciloy, J. M. Human metabolic databases for drug development and drug interactions. *Drug Metab. Rev.* **2004**, *36*, 105–105.
- Hou, B. K.; Kim, J. S.; Jun, J. H.; Lee, D. Y.; Kim, Y. W.; Chae, S.; Roh, M.; In, Y. H.; Lee, S. Y. BioSilico: an integrated metabolic database system. *Bioinformatics* **2004**, *20* (17), 3270–3272.
- Grootenhuys, P. D. J.; Penzotti, J.; Miller, J.; Xu, R.; Kassel, D. Combining admet in Silico, in vitro, and in vivo for drug discovery. *Abstr. Pap. Am. Chem. Soc.* **2000**, *219*, U457–U457.
- Migeon, J. C.; Rogalski, S. L.; Krejsa, C. M.; Horvath, D.; Mao, B.; Barbosa, F.; Merrick, S. E.; Mersberg, M.; Lakehal, F. Using large in vitro ADME data sets to predict in vivo properties. *Drug Metab. Rev.* **2003**, *35*, 168–168.
- Norris, D. A.; Bremer, T.; Holme, K.; Leesman, G.; Sud, M. The failure of in vitro ADME properties to correctly determine in vivo outcomes. *Abstr. Pap. Am. Chem. Soc.* **2003**, *226*, U453–U453.
- Egan, W. J.; Merz, K. M.; Baldwin, J. J. Prediction of drug absorption using multivariate statistics. *J. Med. Chem.* **2000**, *43* (21), 3867–3877.
- Wessel, M. D.; Jurs, P. C.; Tolan, J. W.; Muskal, S. M. Prediction of human intestinal absorption of drug compounds from molecular structure. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (4), 726–735.
- Kubinyi, H. Drug research: myths, hype and reality. *Nat. Rev. Drug Discovery* **2003**, *2* (8), 665–668.
- Abad-Zapatero, C. A Sorcerer’s apprentice and the rule of five: from rule-of-thumb to commandment and beyond. *Drug Discovery Today* **2007**, *12* (23–24), 995–997.
- Keseru, G. M.; Makara, G. M. The influence of lead discovery strategies on the properties of drug candidates. *Nat. Rev. Drug Discovery* **2009**, *8* (3), 203–212.
- Leeson, P. D.; Springthorpe, B. The influence of drug-like concepts on decision-making in medicinal chemistry. *Nat. Rev. Drug Discovery* **2007**, *6* (11), 881–90.
- Bai, J. P. F.; Utis, A.; Crippen, G.; He, H. D.; Fischer, V.; Tullman, R.; Yin, H. Q.; Hsu, C. P.; Jiang, L.; Hwang, K. K. Use of classification regression tree in predicting oral absorption in humans. *J. Chem. Inf. Comput. Sci.* **2004**, *44* (6), 2061–2069.
- Bergstrom, C. A. S.; Stafford, M.; Lazorova, L.; Avdeef, A.; Luthman, K.; Artursson, P. Absorption classification of oral drugs based on molecular surface properties. *J. Med. Chem.* **2003**, *46* (4), 558–570.
- Martin, Y. C. A bioavailability score. *J. Med. Chem.* **2005**, *48* (9), 3164–3170.
- Wenlock, M. C.; Austin, R. P.; Barton, P.; Davis, A. M.; Leeson, P. D. A comparison of physicochemical property profiles of development and marketed oral drugs. *J. Med. Chem.* **2003**, *46* (7), 1250–1256.
- Yoshida, F.; Topliss, J. G. QSAR model for drug human oral bioavailability. *J. Med. Chem.* **2000**, *43* (13), 2575–2585.
- Veber, D. F.; Johnson, S. R.; Cheng, H. Y.; Smith, B. R.; Ward, K. W.; Kopple, K. D. Molecular properties that influence the oral bioavailability of drug candidates. *J. Med. Chem.* **2002**, *45* (12), 2615–2623.
- Hou, T. J.; Wang, J. M.; Zhang, W.; Xu, X. J. ADME evaluation in drug discovery. 6. Can oral bioavailability in humans be effectively predicted by simple molecular property-based rules. *J. Chem. Inf. Model.* **2007**, *47* (2), 460–463.
- Biswas, D.; Roy, S.; Sen, S. A simple approach for indexing the oral druglikeness of a compound: discriminating druglike compounds from nondruglike ones. *J. Chem. Inf. Model.* **2006**, *46* (3), 1394–401.
- Hutter, M. C. Separating drugs from nondrugs: A statistical approach using atom pair distributions. *J. Chem. Inf. Model.* **2007**, *47* (1), 186–194.
- Schneider, N.; Jackels, C.; Andres, C.; Hutter, M. C. Gradual in Silico filtering for druglike substances. *J. Chem. Inf. Model.* **2008**, *48* (3), 613–628.
- Takaoka, Y.; Endo, Y.; Yamanobe, S.; Kakinuma, H.; Okubo, T.; Shimazaki, Y.; Ota, T.; Sumiya, S.; Yoshikawa, K. Development of a method for evaluating drug-likeness and ease of synthesis using a data set in which compounds are assigned scores based on chemists’ intuition. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (4), 1269–1275.
- Wagener, M.; van Geerestein, V. J. Potential drugs and nondrugs: Prediction and identification of important structural features. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (2), 280–292.
- Xu, J.; Stevenson, J. Drug-like index: A new approach to measure drug-like compounds and their diversity. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (5), 1177–1187.

- (50) Zernov, V. V.; Balakin, K. V.; Ivaschenko, A. A.; Savchuk, N. P.; Pletnev, I. V. Drug discovery using support vector machines. The case studies of drug-likeness, agrochemical-likeness, and enzyme inhibition predictions. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (6), 2048–2056.
- (51) Zheng, S. X.; Luo, X. M.; Chen, G.; Zhu, W. L.; Shen, J. H.; Chen, K. X.; Jiang, H. L. A new rapid and effective chemistry space filter in recognizing a druglike database. *J. Chem. Inf. Model.* **2005**, *45* (4), 856–862.
- (52) Zuccotto, F. Pharmacophore features distributions in different classes of compounds. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (5), 1542–1552.
- (53) 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* (3), 856–864.
- (54) Oprea, T. I. Property distribution of drug-related chemical databases. *J. Comput.-Aided Mol. Des.* **2000**, *14* (3), 251–264.
- (55) Oprea, T. I. Current trends in lead discovery: Are we looking for the appropriate properties. *J. Comput.-Aided Mol. Des.* **2002**, *16* (5–6), 325–334.
- (56) Oprea, T. I.; Allu, T. K.; Fara, D. C.; Rad, R. F.; Ostopovici, L.; Bologa, C. G. Lead-like, drug-like or “pub-like”: how different are they. *J. Comput.-Aided Mol. Des.* **2007**, *21* (1–3), 113–119.
- (57) Oprea, T. I.; Davis, A. M.; Teague, S. J.; Leeson, P. D. Is there a difference between leads and drugs? A historical perspective. *J. Chem. Inf. Comput. Sci.* **2001**, *41* (5), 1308–1315.
- (58) Rishton, G. M. Nonleadlikeness and leadlikeness in biochemical screening. *Drug Discovery Today* **2003**, *8* (2), 86–96.
- (59) Wunberg, T.; Hendrix, M.; Hillisch, A.; Lobell, M.; Meier, H.; Schmeck, C.; Wild, H.; Hinzen, B. Improving the hit-to-lead process: data-driven assessment of drug-like and lead-like screening hits. *Drug Discovery Today* **2006**, *11* (3–4), 175–180.
- (60) Dobson, P. D.; Patel, Y.; Kell, D. B. ‘Metabolite-likeness’ as a criterion in the design and selection of pharmaceutical drug libraries. *Drug Discovery Today* **2009**, *14* (1–2), 31–40.
- (61) Ritchie, T. J.; Luscombe, C. N.; Macdonald, S. J. F. Analysis of the Calculated Physicochemical Properties of Respiratory Drugs: Can We Design for Inhaled Drugs Yet. *J. Chem. Inf. Model.* **2009**, *49* (4), 1025–1032.
- (62) Vistoli, G.; Pedretti, A.; Testa, B. Assessing drug-likeness - what are we missing. *Drug Discovery Today* **2008**, *13* (7–8), 285–294.
- (63) Glick, M.; Grant, G. H.; Richards, W. G. Docking of flexible molecules using multiscale ligand representations. *J. Med. Chem.* **2002**, *45* (21), 4639–4646.
- (64) Rayan, A.; Senderowitz, H.; Goldblum, A. Exploring the conformational space of cyclic peptides by a stochastic search method. *J. Mol. Graphics Modell.* **2004**, *22* (5), 319–33.
- (65) Sadowski, J. Optimization of the drug-likeness of chemical libraries. *Perspect. Drug Discovery Des.* **2000**, *20* (1), 17–28.
- (66) Ajay; Walters, W. P.; Murcko, M. A. Can we learn to distinguish between “drug-like” and “nondrug-like” molecules. *J. Med. Chem.* **1998**, *41* (18), 3314–3324.
- (67) Sadowski, J.; Kubinyi, H. A scoring scheme for discriminating between drugs and nondrugs. *J. Med. Chem.* **1998**, *41* (18), 3325–3329.
- (68) Anzali, S.; Barnickel, G.; Cezanne, B.; Krug, M.; Filimonov, D.; Poroikov, V. Discriminating between drugs and nondrugs by prediction of activity spectra for substances (PASS). *J. Med. Chem.* **2001**, *44* (15), 2432–2437.
- (69) Muegge, I.; Heald, S. L.; Brittelli, D. Simple selection criteria for drug-like chemical matter. *J. Med. Chem.* **2001**, *44* (12), 1841–1846.
- (70) Gorelik, B.; Goldblum, A. High quality binding modes in docking ligands to proteins. *Proteins* **2008**, *71* (3), 1373–86.
- (71) Rayan, A.; Noy, E.; Chema, D.; Levitzki, A.; Goldblum, A. Stochastic algorithm for kinase homology model construction. *Curr. Med. Chem.* **2004 Mar**, *11* (6), 675–92.
- (72) Kubinyi, H. From narcosis to hyperspace: The history of QSAR. *Quant. Struct.-Act. Relat.* **2002**, *21* (4), 348–356.
- (73) Pospisil, P.; Ballmer, P.; Scapozza, L.; Folkers, G. Tautomerism in computer-aided drug design. *J. Recept. Signal Transduction* **2003**, *23* (4), 361–371.
- (74) Tetko, I. V.; Bruneau, P. Application of ALOGPS to predict 1-octanol/water distribution coefficients, logP, and logD, of AstraZeneca in-house database. *J. Pharm. Sci.* **2004**, *93* (12), 3103–3110.
- (75) MOE. http://www.chemcomp.com/Journal_of_CCG/Features/descr.htm (accessed Jan. 2010).
- (76) Pirard, B. Computational methods for the identification and optimization of high quality leads. *Comb. Chem. High Throughput Screening* **2004**, *7* (4), 271–280.

CI9004354