

Explicit Diversity Index (EDI): A Novel Measure for Assessing the Diversity of Compound Databases

Ákos Papp,[†] Anna Gulyás-Forró,[†] Zsolt Gulyás,[‡] György Dormán,[†] László Ürge,[†] and Ferenc Darvas^{*,§}

AMRI Hungary, Záhony u. 7, H-1031 Budapest, Hungary, ComGrid Ltd., Záhony u. 7, H-1031 Budapest, Hungary, and CompuDrug International, Inc., 115 Morgan Drive, Sedona, Arizona 86351

Received March 7, 2006

A novel diversity assessment method, the Explicit Diversity Index (EDI), is introduced for druglike molecules. EDI combines structural and synthesis-related dissimilarity values and expresses them as a single number. As an easily interpretable measure, it facilitates the decision making in the design of combinatorial libraries, and it might assist in the comparison of compound sets provided by different manufacturers. Because of its rapid calculation algorithm, EDI enables the diversity assessment of in-house or commercial compound collections.

INTRODUCTION

Molecular diversity is one of the most important characteristics of screening libraries.^{1–4} To increase the hit rate, the use of the most representative compound set that covers the chemical space relevant to the appropriate target is advised.¹ To select the most favorable screening library, a simple measure (preferably expressed as a single number) that explicitly informs the medicinal chemist about the diversity of the investigated compound set would be advantageous.

There are already numerous diversity assessing methods in the literature which measure diversity on the basis of specific structural features of the molecules combined with different mathematical methods.^{5–17} Many diverse library design procedures deal with the selection of the most appropriate reagent sets and building blocks for synthesis.^{18–21}

In the present paper, we propose a novel diversity assessment procedure, the Explicit Diversity Index (EDI), which describes diversity explicitly as a single number through combining structural and combinatorial synthesis-related dissimilarity. EDI is developed to provide assistance in designing diverse libraries primarily in the field of drug discovery.

A recent study demonstrated²² that structural dissimilarity does not contribute directly to biological activity. Surprisingly, only 30% of the purchased compounds having a Tanimoto similarity above 85% have been justified in biological tests. Thus, structural similarity or dissimilarity does not influence the hit rate as expected. A recent publication focused more on scaffold diversity, claiming that increasing the hit rate can be achieved with more scaffolds²³ (in other words, with various diverse cores or chemotypes). According to this “uniform library concept”, an equal scaffold distribution is the optimum within a library.

More recently, shape diversity²⁴ was also described, which relied on the 3D structure of the chemotype or skeleton. In this concept, common molecular skeletons display similar chemical information in the 3D space. On the basis of these assumptions, Burke and Schreiber²⁵ introduced the term skeletal diversity in connection with forward-synthetic planning and diversity-oriented synthesis. They described various examples of their “one synthesis/one skeleton” approach significantly increasing the chemical diversity.

Hogan has published an example demonstrating two magnitudes of scaffold structures:²⁶ First, they may inherently represent biological activity, which is enhanced by the substituents introduced during synthesis. In other cases, no specific biological activity can be assigned to the core structure, but it provides a skeleton for arranging the substituents into the required directions. Walters et al.²⁷ noted that the latter case is particularly important in the field of drug design and combinatorial chemistry. In conclusion, the scaffold distribution (or as we termed core representativeness) is closely related to the combinatorial realization of diversity, thus, to the practice of library design and synthesis.

In the proposed EDI approach, structural dissimilarity and the above combinatorial synthetic aspects are equally accounted for leading to a practical diversity measure. Structural dissimilarity is generally calculated by a comparison of pairwise dissimilarities of a target and a reference set. The combinatorial synthetic design aspects are considered by the involvement of “core representativeness”.

In the present paper, we discuss the calculation of the elements of EDI as well as its testing and application, examining various compound collections that are important in the practice of combinatorial drug discovery.

METHOD DESCRIPTION

Calculation of Structural Dissimilarity. To characterize the structural dissimilarity, the nearest neighbor pairwise dissimilarities of the molecules in the library are calculated. For calculation of the nearest-neighbor's average value, we

* Corresponding author phone: +36 1 214 2306; e-mail: df.computdrug@worldnet.att.net.

[†] AMRI Hungary.

[‡] ComGrid Ltd.

[§] CompuDrug International.

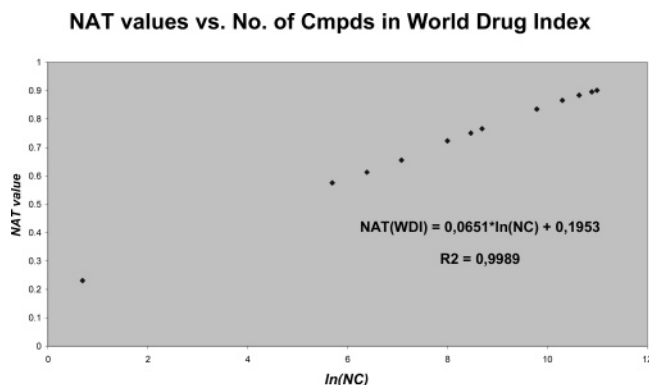


Figure 1. Relationship between NAT and the number of compounds in the World Drug Index.

adapted the procedure of Patterson et al.²⁸ First, the 2D Unity Fingerprints of the Sybyl program package²⁹ are computed, followed by the calculation of the Tanimoto similarity coefficients. The nearest neighbor of each compound is then selected, and the “Nearest-Neighbor Average Tanimoto” (NAT) is produced as the average of the Tanimoto coefficient of the selected duplets. The use of the nearest-neighbor Tanimoto value is more advantageous to describe structural dissimilarity than the “mean pairwise intermolecular dissimilarity”⁷ because it expresses the dissimilarity from the most similar compounds rather than giving a simple average dissimilarity. In the case of compound sets that have been investigated in this paper, the distribution of the nearest-neighbor (NN) Tanimoto values is usually not normal because a very high percentage (ca. 80%) of the values are between 0.9 and 1, while usually, there is a portion of compounds having smaller NN Tanimoto values giving a nonsymmetrical distribution curve. The reason for the surprisingly high similarity values is that the examples and validation studies, in accordance with the present-day discovery practice, represent combinatorial libraries, and these types of compound sets usually contain very similar structures.

The diversity of the compounds are compared to the diversity of the compounds listed by the World Drug Index (WDI);³⁰ WDI serves as a generally recognized and acknowledged compendium of chemical structures used in clinical practice, giving a good representation of the druglike chemical space.^{23,31,32}

Another important part of the EDI concept is to incorporate the size of the library. In the case of sets having high numbers of compounds, the probability of finding similar structures is always higher than that in a smaller set. As a consequence of this, the diversity values of libraries with different sizes are not comparable. To eliminate the effect of the library size, the NAT value of the studied library is compared to the NAT value calculated for a WDI subset of the same size. We created an adequate number of independent, differently sized subsets of the WDI by random selection (repeated three times), for which NAT values were calculated and averaged (see the Supporting Information, Tables 4 and 5). The relationship between the number of compounds and the NAT values for 11 subsets are plotted in Figure 1.

By approaching the relationship between NAT and the number of compounds with a power function, the NAT value for a WDI subset of deliberate size can be calculated by extrapolation.

For the final step, structural dissimilarity (SDI) is calculated with eq 1:

$$SDI = \frac{1 - NAT_{\text{current}}}{1 - NAT_{\text{WDI}}} = \frac{1 - NAT_{\text{current}}}{0.8047 - 0.065 \ln(NC)} \quad (1)$$

where NAT_{current} is the NAT value of the compound set in question and NAT_{WDI} is the return value of the NAT function determined for WDI at the number of compounds in the compound set in question (and \ln is the natural logarithm). Instead of directly comparing these NAT values, this formula amplifies the difference in the increment of diversity over a compound set including only identical compounds (for which $NAT = 1$). In practical situations, SDI is rarely higher than 1 (because it means a better structural diversity than that of WDI), and usually, it is between 0.8 and 1.

Shortened Calculation Method for NAT Value. The time necessary for calculating NAT values has a quadratic dependence on the number of compounds, which makes the calculation of NAT for large compound sets cumbersome. The most time-consuming step in this procedure is the Tanimoto calculation for each compound pair. There are several ways of finding the nearest neighbor of a compound in a data set.^{33–36} We have also developed a new algorithm, the Shortened Nearest-Neighbor Average Tanimoto (SNAT) method. It saves computational time by identifying a smaller subset of the compounds, which contains the nearest neighbor of a compound with high probability, and therefore, the Tanimoto coefficients have to be calculated only for a reduced number of compound pairs. The algorithm starts with selecting a compound randomly (compound “A”) and calculating the Tanimoto coefficients with all other compounds in the entire set. The compounds are then sorted by Tanimoto values; the nearest neighbor of compound “A” is the first compound in this list (and vice versa), while the second one will be the next investigated compound (compound “B”). The basic assumption of the SNAT algorithm is that the nearest neighbor of compound “B” will be included in the set of closest neighbors of compound “A”; therefore, a predefined number of compounds from the top of the sorted neighbor list of compound “A” is selected as a representative set for compound “B” and called the set of “closest neighbors”. When searching for the nearest neighbor of compound “B”, the Tanimoto coefficients are calculated only for this representative set. Because the list is sorted after each cycle, this procedure ensures that the particular list will be representative for the compound that will be selected in the next cycle. After a specified number of cycles, the content of the representative set is refreshed by calculating the Tanimoto coefficients for the entire set again. We compared the NAT values calculated with the original and those with the shortened algorithm by calculating the “correctness (%)” of the calculation. This is the percentage of compounds that have the same nearest neighbor with the two different calculation algorithms (Table 1). In the case of a usual-sized compound set (some tens of thousands of compounds), the correctness is very high, even in the case of using the smallest number of closest neighbors as an option. Furthermore, which is more important, the difference between SNAT and NAT is negligible (less than 0.5%). On the other hand, the contribution of the difference of SNAT and NAT in the $1 - NAT$ value is a little bit higher, but this effect is considerable

Table 1. Comparison of the NAT Values Calculated with the Original and the Shortened Algorithm for Two Different Sets

set size	54 664				166 110
number of closest neighbors	1086	2717	5435	6017	1000
correctness (%)	93.7	95.2	96.7	96.8	98.1
SNAT	0.893	0.894	0.896	0.896	0.977
NAT			0.898		0.977
(SNAT - NAT)/ (1 - NAT) ^a (%)	4.7	3.7	1.9	1.9	0.2

^a (SNAT - NAT)/(1 - NAT) shows the inaccuracy caused by the use of the shortened algorithm for the calculation of NAT.

only in the case of smaller compound sets and can be eliminated by increasing the parameter used for the number of closest neighbors. On the whole, the SNAT algorithm can be used successfully instead of the longer NAT calculation algorithm.

Calculation of Core Representativeness. In accordance with the standard definition in organic chemistry, “core” (scaffold, skeleton, or chemotype) is the common substructure of a compound library synthesized by the same or analogous synthetic routes. The term “common substructure” refers to the stable and repetitive core structural element of the library, which can typically be represented by a Markush structure. It can be either an open-chain or a ring (ring system) together with its substituent pattern.

“Cores” are represented mostly in different degrees in a set of similar compounds. In our concept, “core representativeness” gives a general description of the degree of the representation of the cores within the library members.

We constructed the core representative component of the diversity in the following three steps:

(a) Identify the parameters that may influence the numerical value of “core representativeness” (CR). Evidently, CR depends on the number of compounds in the whole data set (NC) and on the number of core structures (CS):

$$CR = f(NC, CS)$$

(b) Select direct or inverse dependencies of the parameters on CR. As intuitively evident, CR decreases with the average number of compounds per core structure and increases with the number of core structures within the library:

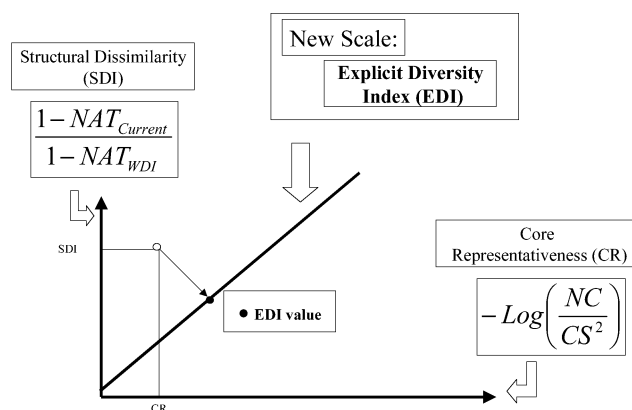
$$CR = f(NC/CS, 1/CS)$$

(c) Define the exact form of the particular function. For that purpose, CR is expressed as the negative logarithm (base 10 logarithm, hereinafter referred to as “log”) of the ratio of NC and the square of CS (eq 2):

$$CR = -\log\left[\left(\frac{NC}{CS}\right)\frac{1}{CS}\right] = -\log\left(\frac{NC}{CS^2}\right) \quad (2)$$

In the case of compounds produced by combinatorial or parallel synthesis, CS is equal to the number of synthetic libraries. In other cases, CS can be obtained by clustering the compound set on the basis of maximum common substructure searches.^{37,38}

Calculation of the EDI Value. The combination of structural dissimilarity and “core representativeness” provides a novel, practical measure of library attractiveness and helps the screening practitioners in their selection efforts.

**Figure 2.** Overview of the EDI algorithm.

In our approach, structural dissimilarity and core representativeness are combined to a single number, EDI, by using the Spectral Mapping technique.³⁹ This is a simple linear mathematical method for converting the parameter data of the investigated databases in a short and rapid way to one information-rich number. Because of its calculation algorithm, it is easy to apply as the input of a principal component analysis or clustering. The calculation of EDI is demonstrated in Figure 2.

Because in the expression of EDI there is no reason to intensify the role of either SDI or CR, we applied the even weighting of these components, and the spectral mapping is simplified to a linear combination of SDI and CR, as it is expressed by eq 3:

$$EDI_{\text{unscaled}} = \frac{\sqrt{2}}{2}(SDI + CR) \quad (3)$$

Scaling the Value of EDI. To facilitate the interpretation of EDI and to establish theoretical upper and lower endpoints of the scale, we created two virtual compound sets, one with maximum and another with minimum EDI values. Theoretically, the upper endpoint of the EDI scale is a very large, heterogeneous compound set, with 2⁹⁸⁸ compounds (see the explanation in the Supporting Information), each belonging to a separate core, while the lower endpoint is a structurally homogeneous set containing 2⁹⁸⁸ identical compounds. For the upper endpoint, the EDI value is infinite, while for the lower endpoint, it is -210.3 (see the explanation in the Supporting Information).

To create a normalized scale between 0 and 100 for EDI (Figure 3), the raw EDI values are transformed by a tangent hyperbolic function (eq 3), which gives zero for the structurally homogeneous set (formerly EDI = -210.3) and gives 100 for the theoretically heterogeneous set (formerly EDI = ∞). To achieve the optimal sensitivity of the EDI in the domain of practical diversity values, the value of the tangent hyperbolic function was divided by 3. The advantage of this transformation is that the generally applicable EDI values are in essence linearly transformed and can be defined to any extreme value, giving a number between 0 and 100 (eq 4).

$$EDI_{\text{scaled}} = \frac{\tanh\left(\frac{EDI_{\text{unscaled}}}{3}\right) + 1}{2} \times 100 \quad (4)$$

$$EDI_{scaled} = \frac{\tanh\left(\frac{EDI_{unscaled}}{3}\right) + 1}{2} \cdot 100 \quad (4)$$

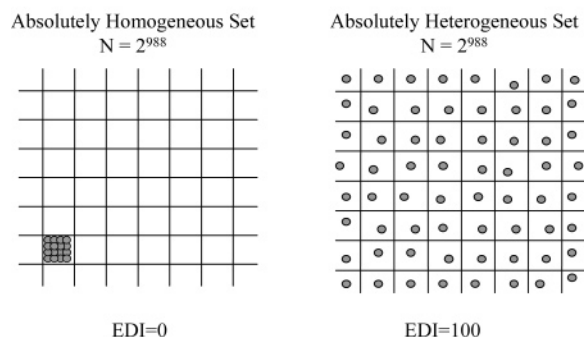


Figure 3. Theoretical endpoints of the EDI scale.

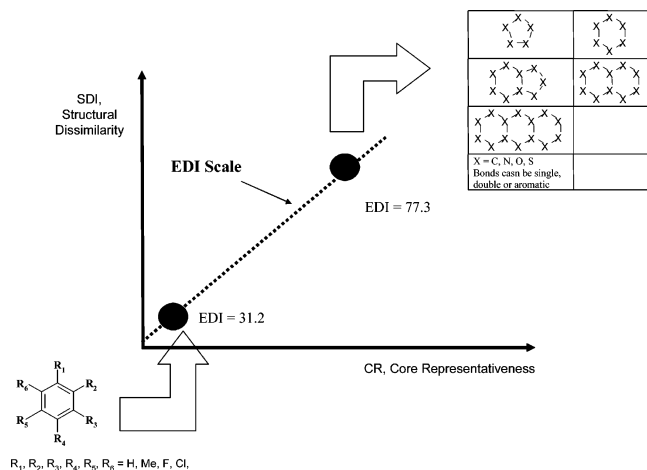


Figure 4. Defining an experimental scale for the EDI approach.

VALIDATION

The purpose of validation is to prove the effectiveness of the EDI approach. As a first validation step, we investigated the compatibility of the EDI approach with the empirical assessment of diversity. We constructed two libraries, one was a typical heterogeneous library and the other a typical homogeneous one, displaying a significant difference in their diversity by simple visual observation. The homogeneous set contains very similar structures: 67 substituted benzene analogues (set A). The second is constructed from 67 simple heterocyclic rings (set B; Figure 4) considered as singletons with different cores (see the Supporting Information, Tables 6 and 7). EDI was calculated as 31.2 for set A and 77.3 for set B. This difference shows that EDI properly reflects the diversity differences. At the same time, the calculated EDI values for these model sets define somehow the rough boundaries of the diversity scale in a daily combinatorial synthesis.

In the second validation step, the compatibility of the EDI approach was studied on several libraries produced by the commonly used OptiSim⁴⁰ diversity selection algorithm. The selection was carried out in five consecutive steps, reducing the number of structures from 10 000 to 500, resulting in five different compound sets, with continuously increasing diversity. The calculated EDI values again properly describe the tendency: the EDI is 56.9 for the largest, least diverse library and 73.3 for the smallest, most diverse set (Figure 5).

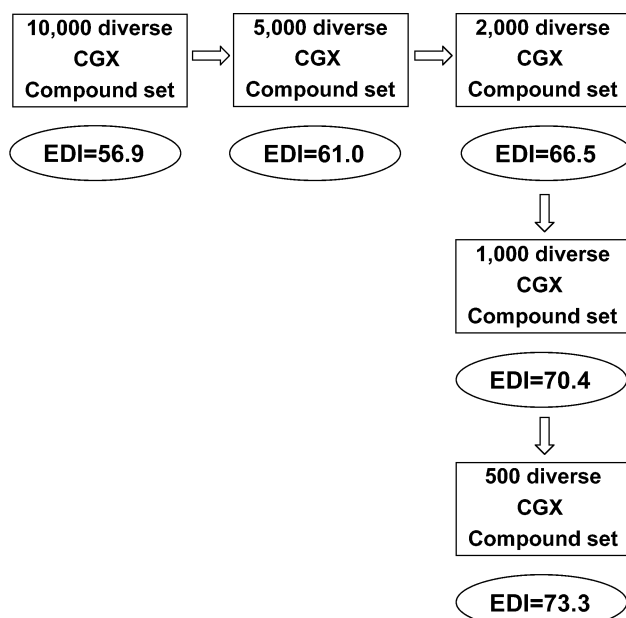


Figure 5. Comparison of the EDI values with the commonly used OptiSim⁴⁰ diversity selection algorithm.

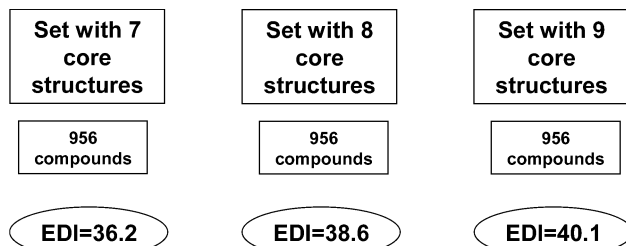
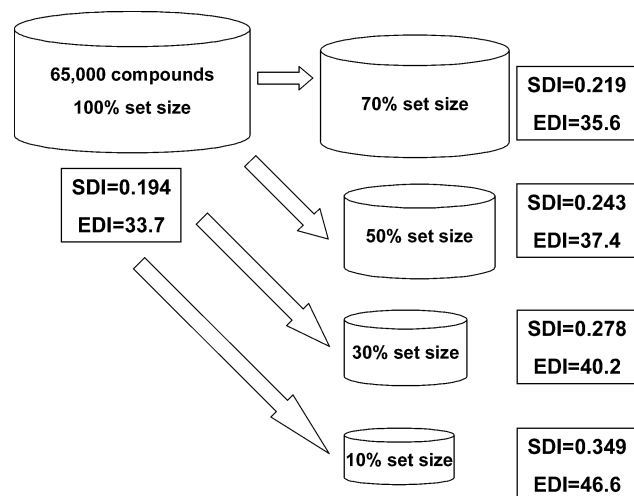


Figure 6. Investigating the effect of changing core representativeness systematically.

The purpose of the next validation step was to investigate the effect of increasing and decreasing core representativeness by systematically changing the number of cores in the database. Three different sets were constructed, all of them containing 956 compounds belonging to seven, eight, and nine core structures (Figure 6).

The starting set contained seven core structures. The second set was created by adding compounds belonging to a new core structure and, at the same time, deleting an equal number of redundant compounds belonging to one of the original core structures. This way, the number of core structures was raised from seven to eight and, in case of the third set, from eight to nine. During the selection procedure, in addition to keeping the number of compounds steady, we carefully strived to keep the structural complexity of the compounds on the same level by inspecting the complexity of the removed and replaced compounds. As it is shown in Figure 6, by reducing the number of core structures in the database, the diversity is decreased together with the EDI value of the compound sets. When increasing the number of core structures, and thus the diversity of the compound set, the EDI value increases.

The effect of the library size on the value of EDI has also been studied. As we mentioned in the Method Description-section, the EDI approach compensates for the effect of set size (Figure 7). We selected 10%, 30%, 50%, and 70% subsets of a 166 110-member library by random selection

**Figure 7.** Investigating the effect of set size.**Table 2.** Comparing Different Libraries with the EDI Approach

compound set	NC	CS	CR	NAT _{WDI}	NAT _{current}	SDI	EDI
library 1	119 702	213	-0.421	0.955	0.957	0.967	56.49
library 2	6122	67	-0.135	0.763	0.893	0.448	53.68
library 3	166 110	89	-1.322	0.976	0.977	0.956	45.69
library 4	29 183	23	-1.756	0.867	0.979	0.152	31.96

^a NC: number of compounds; CS: number of core structures; CR: core representativeness; NAT_{WDI}: NAT value for the whole World Drug Index; NAT_{current}: NAT value for the current database; SDI: structural dissimilarity; EDI: scaled explicit diversity index.

and investigated how size changes are reflected in the EDI values.

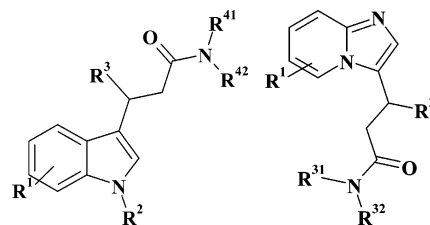
The results presented in Figure 7 show that SDI and EDI remain approximately constant for all sets. The single exception is the smallest set, where the effect of set size on the values of structural dissimilarity and EDI are not compensated by the present algorithm. (The detailed results of the validation calculations are summarized in the Supporting Information, Table 8.)

APPLICATION EXAMPLES

Example 1. Comparing Inventory Libraries. It is always difficult to decide which compound set should be purchased or synthesized for screening purposes. The comparison of different inventory libraries is very convenient by calculating EDI values for all optional libraries and selecting the one with the largest EDI value. Data and EDI values of different commercial libraries are shown in Table 2, two of them containing more than 100 000 compounds, the others having only a few thousand structures. There are large differences in the structural dissimilarity values; libraries 1 and 3 display higher values, which shows that their structural properties are very similar to the World Drug Index of the same size. The core representativeness value is the highest for library 2, because this one incorporates the lowest number of compounds with a high number of core structures. The opposite is true for library 4, which has a relatively small number of compounds and just a few core structures. By combining all of these properties, library 1 has the best EDI value, because it has fairly high core representativeness and the highest structural dissimilarity.

Table 3. Effect of the Diversity in the Second Derivatization Step on the Explicit Diversity Index (EDI)

	library 1	library 2
multiplication in the first derivatization step	6	5
multiplication in the second derivatization step	10	4
multiplication in the third derivatization step	82	90
number of virtual compounds	4400	1800
total number of reactions	98	99
EDI values of the libraries	50.1	44.1



Example 2. Design of Combinatorial Libraries. EDI can be applied to select the theoretically most diverse combinatorial library for synthesis (Table 3). For example, we devised two relatively similar library cores: substituted indole-3-propionic acids (library 1) and their “azalog” 3-imidazo-[1,2-a]-pyridin-3-yl-propionic acids (library 2) using a three-step reaction sequence.⁴¹

On the basis of the available building blocks and their synthetic feasibility, we constructed virtual libraries with different combinatorial matrix arrangements. The question was: which library should be synthesized covering more chemical space with enhanced diversity? In the first derivatization step, the multiplying factors are almost the same, 6 and 5, and in the last step, this factor is around 90 in both of the libraries. The main difference is in the second step, where in library 1 10 reagents are used and in library 2 only four reagents. On the basis of the calculated EDI values, the chemists involved in the project selected library 1 for synthesis.

When comparing two or more individual libraries that were synthesized starting from a single core structure, the “core” is defined as the number of multiplying reagents at the different reaction steps.

DISCUSSION AND SUMMARY

The EDI method has been used successfully over the past three years at ComGenex,^{42,43} accelerating the routine compound design decisions significantly. The run time of the software calculating EDI was dependent on the size of the libraries; for focused libraries (100–1000 members), it takes a few seconds; for discovery libraries (1000–4000), it takes several minutes, while for our molecular banks comprising normally a few hundred thousand molecules, the running time was on the order of 1–2 h, still allowing a comfortable decision.

The assessment of the diversity of different compound sets is a crucial factor in today’s high-throughput drug discovery

problems. We developed an easily applicable diversity measure, EDI, which expresses the diversity of a compound set in the form of a single number. EDI is intended to be used in the field of drug discovery. It combines structural and synthetic dissimilarity elements, giving a better approximation of biological diversity over the traditional measures on the basis of purely structural features. For the rapid calculation of EDI, we also developed a novel algorithm for obtaining the average nearest-neighbor Tanimoto values (NAT). The SNAT algorithm can determine the NAT value with over 99% correctness and significantly reduces the time needed for the calculation.

The presented EDI approach enables primarily the comparison of molecular collections of different sizes of up to approximately 80 000 structures. This limitation comes from the extrapolation limitations of the function determined for the reference data set (i.e., the copy of the World Drug Index we used contained only 60 000 compounds). At the moment, this does not represent serious limitations in practice because, in most of the comparisons, diverse subsets of compound collections were used, which fall into this range. The system allows the upper limit to be custom modified (increased) if another suitable database is used as a reference collection containing a significantly higher number than the WDI (even the most up-to-date version of WDI consists of only 80 000 compounds). Companies and laboratories would be able to construct a proprietary reference library combining various databases. With an approximately 150 000-reference library, the upper limit of the comparison could reach the 200 000-compound library size.

On the basis of the presented validation studies, we believe that EDI is a proper measure of diversity and is appropriate for solving a number of practical diversity problems emerging in daily drug discovery practice. This includes comparing the diversity of in-house or commercially available libraries, designing combinatorial libraries with enhanced diversity, continuously monitoring the diversity of in-house inventory libraries, and controlling the efficiency of diverse selection methods.

Supporting Information Available: The description of the scaling of EDI, the sensitivity of the index, and applicability studies. This material is available free of charge via the Internet at <http://pubs.acs.org>.

REFERENCES AND NOTES

- Martin, Y. C. Molecular Diversity: How We Measure It? Has It Lived Up to Its Promise? *Farmaco* **2001**, *56*, 137–139.
- Bradley, M. P. An Overview of the Diversity Represented in Commercially Available Databases. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 301–309.
- Cummins, D. J.; Andrews, C. W.; Bentley, J. A.; Cory, M. Molecular Diversity in Chemical Databases: Comparison of Medicinal Chemistry Knowledge Bases and Databases of Commercially Available Compounds. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 750–763.
- Voigt, J. H.; Bienfait, B.; Wang, S.; Nicklaus, M. C. Comparison of the NCI Open Database with Seven Large Chemical Structural Databases. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 702–712.
- Agrafiotis, D. K.; Lobanov, V. S.; Rassokhin, D. N.; Izrailev, S. The Measurement of Molecular Diversity. In *Virtual Screening for Bioactive Molecules*; Wiley: Weinheim, Germany, 2000; pp 265–300.
- Waldman, M.; Li, H.; Hassan, M. Novel Algorithms for the Optimization of Molecular Diversity of Combinatorial Libraries. *J. Mol. Graphics Modell.* **2000**, *18*, 412–426.
- Turner, D. B.; Tyrrell, S. M.; Willett, P. Rapid Quantification of Molecular Diversity for Selective Database Acquisition. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 18–22.
- Golbraikh, A. Molecular Dataset Diversity Indices and Their Applications to Comparison of Chemical Databases and QSAR Analysis. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 414–425.
- Zheng, W.; Cho, S. J.; Waller, C. L.; Tropsha, A. Rational Combinatorial Library Design. 3. Simulated Annealing Guided Evaluation (SAGE) of Molecular Diversity: A Novel Computational Tool for Universal Library Design and Database Mining. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 738–746.
- Martin, E. J.; Critchlow, R. E. Beyond Mere Diversity: Tailoring Combinatorial Libraries for Drug Discovery. *J. Comb. Chem.* **1999**, *1*, 32–45.
- Reynolds, C. H.; Tropsha, A.; Pfahler, L. B.; Drujer, R.; Chakravorty, S.; Ethiraj, G.; Zheng, W. Diversity and Coverage of Structural Sublibraries Selected Using the SAGE and SCA Algorithms. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 1470–1477.
- Lobanov, V. S.; Agrafiotis, D. K. Stochastic Similarity Selections from Large Combinatorial Libraries. *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 460–470.
- Lin, S. K. The Similarity Principle and Its Application. *Molecules* **1996**, *1*, 57–67.
- Agrafiotis, D. K. On the Use of Information Theory for Assessing Molecular Diversity. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 576–580.
- Jorgensen, M. A. M.; Pedersen, J. T. Structural Diversity of Small Molecule Libraries. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 338–345.
- Trepalin, S. V.; Gerasimenko, V. A.; Kozyukov, A. V.; Savchuk, N. P.; Ivaschenko, A. A. New Diversity Calculations Algorithms Used for Compound Selection. *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 249–258.
- Mount, J.; Ruppert, J.; Welch, W.; Jain, A. N. Structural Diversity of Small Molecule Libraries. *J. Med. Chem.* **1999**, *42*, 60–66.
- Todorov, N. P.; Dean, P. M. Evaluation of a Method for Controlling Molecular Scaffold Diversity in de Novo Ligand Design. *J. Comput.-Aided Mol. Des.* **1997**, *11*, 175–192.
- Clark, R. D.; Kar, J.; Akella, L.; Soltanshahi, F. OptDesign: Extending Optimizable k-Dissimilarity Selection to Combinatorial Library Design. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 829–836.
- Gillet, V. J.; Willett, P.; Bradshaw, J. The Effectiveness of Reactant Pools for Generating Structurally-Diverse Combinatorial Libraries. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 731–740.
- Koehler, R. T.; Dixon, S. L.; Villar, H. O. LASSOO: A Generalized Directed Diversity Approach to the Design and Enrichment of Chemical Libraries. *J. Med. Chem.* **1999**, *42*, 4695–4704.
- Martin, Y. C.; Kofron, J. L.; Traphagen, L. M. Do Structurally Similar Molecules Have Similar Biological Activity? *J. Med. Chem.* **2002**, *45*, 4350–4358.
- Nilakantan, R.; Immermann, F.; Haraki, F. A Novel Approach to Combinatorial Library Design. *Comb. Chem. High Throughput Screening* **2002**, *5*, 105–110.
- Sauer, W. H.; Schwarz, M. K. Molecular Shape Diversity of Combinatorial Libraries: A Prerequisite for Broad Bioactivity. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 987–1003.
- Burke, M. D.; Schreiber, S. L. A Planning Strategy for Diversity-Oriented Synthesis. *Angew. Chem., Int. Ed.* **2004**, *43*, 46–58.
- Hogan, J. C. Directed Combinatorial Chemistry. *Nature* **1996**, *384*, 17–19.
- Walters, W. P.; Stahl, M. T.; Murcko, M. A. Virtual Screening — An Overview. *Drug Discovery Today* **1998**, *3* (4), 160–178.
- Patterson, D. E.; Cramer, R. D.; Ferguson, A. M.; Clark, R. D.; Weinberger, L. E. Neighborhood Behavior: A Useful Concept for Validation of “Molecular Diversity” Descriptors. *J. Med. Chem.* **1996**, *39*, 3049–3059.
- SYBYL, version 6.8; Tripos Associates, Inc.: St. Louis, MO, 2001.
- World Drug Index, sample SD edition; Derwent: Philadelphia, PA, 2002.
- Nilakantan, R.; Bauman, N.; Haraki, K. S. Database Diversity Assessment: New Ideas, Concepts, and Tools. *J. Comput.-Aided Mol. Des.* **1997**, *11*, 447–452.
- Schneider, G.; Lee, M.-L.; Stahl, M.; Schneider, P. De Novo Design of Molecular Architectures by Evolutionary Assembly of Drug-Derived Building Blocks. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 487–494.
- Agrafiotis, D. K.; Lobanov, V. S. An Efficient Implementation of Distance-Based Diversity Measures Based on k-d Trees. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 51–58.
- Hudson, B. D.; Hyde, R. M.; Rahr, E.; Wood, J. Parameter Based Methods for Compound Selection from Chemical Databases. *Quant. Struct.-Act. Relat.* **1996**, *15*, 285–289.
- Holliday, J. D.; Ranade, S. S.; Willett, P. A Fast Algorithm for Selecting Sets of Dissimilar Molecules from Large Chemical Databases. *Quant. Struct.-Act. Relat.* **1995**, *14*, 501–506.

- (36) Allen, B. C. P.; Grant, G. H.; Richards, W. G. Similarity Calculations Using Two-Dimensional Molecular Representations. *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 330–337.
- (37) Raymond, J. W.; Willett, P. Maximum Common Subgraph Isomorphism Algorithms for the Matching of Chemical Structures. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 521–533.
- (38) Stahl, M.; Mause, H. Database Clustering with a Combination of Fingerprint and Maximum Common Substructure Methods. *J. Chem. Inf. Model.* **2005**, *45*, 542–548.
- (39) Lewi, P. J. Spectral Mapping, a Technique for Classifying Biological Activity Profiles of Chemical Compounds. *Arzneim. Forsch.* **1976**, *26*, 1295–1299.
- (40) Clark, R. D. OptiSim: An Extended Dissimilarity Selection Method for Finding Diverse Representative Subsets. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 1181–1188.
- (41) Gerencser, J.; Nagy, T.; Panka, G.; Egyed, O.; Urge, L.; Darvas, F. A Versatile Procedure for the Preparation of 3-Imidazo[1,2-a]pyridin-3-yl-propionic Acids. A Novel Application of Meldrum's Acid. ACS Meeting, New Orleans, LA, March 23, 2003.
- (42) Gulyas-Forro, A.; Urge, L.; Papp, A.; Darvas, F. Novel Method for Calculating the Absolute Diversity Index (ADI) for Drug Candidate Libraries. EuroCombi-2, Copenhagen, Denmark, June 30–July 3, 2003.
- (43) Gulyas-Forro, A.; Urge, L.; Papp, A.; Dorman, G.; Darvas, F. Novel Method for Calculating the Absolute Diversity Index (ADI) for Drug Candidate Libraries. Advances in Chromatography/Electrophoresis and Computational Chemistry, Budapest, Hungary, October 27–29, 2003.

CI060074F