

# Visual Characterization and Diversity Quantification of Chemical Libraries: 2. Analysis and Selection of Size-Independent, Subspace-Specific Diversity Indices

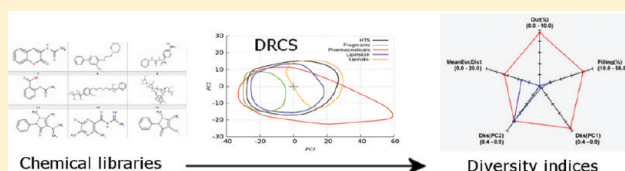
Lionel Colliandre,<sup>†</sup> Vincent Le Guilloux,<sup>†</sup> Stephane Bourg,<sup>‡</sup> and Luc Morin-Allory<sup>†,\*</sup>

<sup>†</sup>Institut de Chimie Organique et Analytique (ICOA), Université d'Orléans-CNRS, UMR 7311 B.P. 6759, rue de Chartres, 45067 Orléans Cedex 2, France

<sup>‡</sup>Fédération de Recherche, Physique et Chimie du Vivant, Université d'Orléans-CNRS; FR 2708, avenue Charles Sadron, 45071 Orléans Cedex 2, France

## S Supporting Information

**ABSTRACT:** High Throughput Screening (HTS) is a standard technique widely used to find hit compounds in drug discovery projects. The high costs associated with such experiments have highlighted the need to carefully design screening libraries in order to avoid wasting resources. Molecular diversity is an established concept that has been used to this end for many years. In this article, a new approach to quantify the molecular diversity of screening libraries is presented. The approach is based on the Delimited Reference Chemical Subspace (DRCS) methodology, a new method that can be used to delimit the densest subspace spanned by a reference library in a reduced 2D continuous space. A total of 22 diversity indices were implemented or adapted to this methodology, which is used here to remove outliers and obtain a relevant cell-based partition of the subspace. The behavior of these indices was assessed and compared in various extreme situations and with respect to a set of theoretical rules that a diversity function should satisfy when libraries of different sizes have to be compared. Some gold standard indices are found inappropriate in such a context, while none of the tested indices behave perfectly in all cases. Five DRCS-based indices accounting for different aspects of diversity were finally selected, and a simple framework is proposed to use them effectively. Various libraries have been profiled with respect to more specific subspaces, which further illustrate the interest of the method.



## INTRODUCTION

Molecular diversity is an established concept that is now routinely used to improve the results of a High Throughput Screening (HTS) campaign. The underlying assumption is that a diverse set of compounds will increase the ratio of hits by maximizing the number of chemically distinct molecules to be tested, thereby avoiding redundant tests.<sup>1</sup> It is thus necessary to have tools that allow the description and comparison of different chemical libraries to be performed and, hence, facilitate the selection of the most promising one(s).<sup>2–4</sup>

The concept of chemical space provides a way to represent chemical libraries in a fixed referential, making it possible to assess their relative coverage and absolute diversity.<sup>5</sup> In a previous article,<sup>6</sup> a new graphical representation that can be used to represent a particular library in a reduced chemical space was introduced. Delimited Reference Chemical Subspaces are defined by the combination of a Principal Component Analysis (PCA) model (DRCS model) and a subspace delimitation (DRCS contour) intended to encompass a large proportion of compounds. The delimitation is computed on a reduced (2D) space obtained using the PCA model and is based on the convex-hull calculation. Isolated compounds (outliers) are excluded prior to the creation of this delimitation, which finally represents the most populated (dense) subspace spanned by the reference library.

Although an intuitive visual inspection is a mandatory step to determine the space coverage of a particular library, a quantitative assessment of molecular diversity can provide complementary information. Automated library design, in particular, makes extensive use of diversity indices as target values that need to be optimized. Diversity indices can also be used as a complementary decision support when visual inspection does not permit to choose between two similar libraries.

Many different types of diversity indices have been developed in past decades.<sup>7</sup> Although intuitive in appearance, molecular diversity is actually difficult to quantify precisely because what we call “diversity” is obviously strongly dependent on the method used to describe molecules, the metric used to compare these molecules, and the function used to quantify the diversity itself. In terms of applicability domain, molecular diversity indices are typically used to compare and optimize collections of fixed size extracted from a reference library.<sup>8</sup> However, when one needs to analyze and compare two libraries of different sizes, e.g., after an enrichment process, some of the indices are found inappropriate. Previous research has shown that the size

**Special Issue:** 2011 Noordwijkerhout Cheminformatics

**Received:** November 9, 2011

**Published:** December 19, 2011

of the library usually has a notable influence on the final value of most diversity indices.<sup>9–11</sup>

Diversity functions can be classified in three broad categories: (1) distance-based methods, which compute the diversity using compounds' similarity/dissimilarity measures, (2) cell-based methods, which partition a multidimensional descriptor space into a finite number of cells used to assess the space coverage, and (3) structure-based methods, which seek to maximize the number of different substructures in the target library (e.g., scaffolds). A decade ago, Waldman et al.<sup>9,12</sup> proposed in two excellent papers a framework to better describe the expected behavior of a good diversity function. In particular, they defined a diversity function as being any protocol that quantitatively assesses the coverage of a particular descriptor space.<sup>9</sup> Such a definition implies that a working descriptor space must be defined. Furthermore, as working in an open space is inappropriate in most cases, a way to delimit this space is usually required to focus attention on the populated regions. To this end, cell-based methods are probably the simplest, fastest, and most appropriate ways of partitioning a particular descriptor space and have been extensively used for this purpose<sup>13,14</sup> to characterize chemical libraries. It is however well recognized that they usually suffer from the presence of sparse regions, which usually leads to an overconsideration of outlier compounds.<sup>15</sup> The choice of the binning scheme used to partition the chemical space usually remains empirical, although valuable improvements have been proposed.<sup>11</sup>

In this work, we will show that the DRCS methodology can be used to tackle the problem of delimiting a reduced descriptor space and to obtain Relative Diversity Indices (RDI) that are independent of the cardinality of the library. To do so, we suggest that the densest region of a particular subspace, as defined by the DRCS contour, should be focused on to apply cell-based diversity functions, leaving the outliers to be analyzed separately. As the DRCS contour reflects the overall shape of the subspaces spanned by a chosen library, a more relevant and flexible partition of it can be obtained compared to the classical hypercube partitioning. A total of 22 diversity indices have been implemented and analyzed in combination with the DRCS methodology. As no index behaves perfectly, five indices that account for different aspects of diversity were finally selected. Various libraries were profiled using different DRCS subspaces obtained using more specific reference libraries (e.g., Lipinski compliant, fragment, or pharmaceutical compounds), which further illustrate the interest of the DRCS methodology.

## METHOD

Various real and fictive chemical libraries were used to study the behavior of 22 diversity indices. All these libraries and the indices will be further described.

**Publicly Available Data Sets.** A database of 6.63 million unique, standardized, and nonreactive compounds<sup>6</sup> that were gathered from 73 vendor collections was used to create the following subsets:

- 0.12 M compounds satisfying the rule of three as defined by Congreve et al. (molecular weight <300, number of hydrogen bond donors  $\leq 3$ , number of hydrogen bond acceptors  $\leq 3$ , and ClogP  $\leq 3$ ; no violation of these four rules allowed)<sup>16</sup>
- 6.22 M compounds satisfying the rule of five as defined by Lipinski et al.<sup>17</sup>

- 0.41 M compounds not satisfying the rule of five

MOE<sup>18</sup> descriptors were used to create these filtered sets. Several commercial or publicly available compound collections were also used in this study:

1. The Prestwick Chemical Library<sup>19</sup> containing only marketed drugs
2. The Comprehensive Medicinal Chemistry (CMC),<sup>20</sup> a database of pharmaceutical compounds
3. The Chembridge Kinaset,<sup>21</sup> a target-based library containing compounds similar to kinase inhibitors
4. The Pyxis Discovery Smart Fragment Library,<sup>22</sup> a fragment library based on scaffolds found in existing drugs
5. The EPA Fathead Minnow Acute Toxicity (EPAFHM),<sup>23,24</sup> a database containing compounds with known toxicity, often used for the development of predictive quantitative structure–activity relationship (QSAR) models of toxicity
6. A set of 1420 accepted marketed compounds retrieved from the DrugBank database<sup>25</sup>
7. Two combinatorial libraries used in the article of Owen et al.:<sup>26</sup> A3 and B5 libraries. The original denomination of these libraries was kept.

A combined library was created by merging the DrugBank, Prestwick, and CMC libraries. This combined library is intended to represent “pharmaceutical compounds” in a broad sense. Moreover, random selections were made to create 60 subsets from the CMC database (8773 compounds). The subsets vary from 100 to 8000 compounds (12 sizes of subsets; 5 subsets per size).

Each library was standardized using the protocol described in the first article, and the MOE2D descriptors were subsequently computed.

**Fictive Data Sets.** A set of 33 fictive libraries was created to study the behavior of diversity indices in various extreme situations. On the basis of the DRCS, a grid was applied to the 2D subspace, strictly encompassing the DRCS contour. The fictive libraries were created by adding points at the center of selected cells of the grid. Each point represents one compound. Only the cells inside the DRCS can be completed, i.e., cells with at least one corner inside the DRCS. Except for libraries H, I, J, and K, a grid containing around 2500 cells inside the DRCS was used. The exact number of cells depends on the subspace used, which might lead to an optimal grid that contains a slightly different number of cells (see the Group 3 section). Schematic views of each library can be found in the Results section and in the Supporting Information.

- Library A was created to perfectly occupy the subspace, i.e., one point was added in each cell.
- Libraries B, C, and D were created by duplicating library A two, four, or eight times, respectively.
- Libraries E, F, and G: Two points per cell were created for half of the cells. The occupied cells are regularly distributed in the delimited space in library E (two points in the first cell, then zero points in the next cell, etc.). In F, only the cells in the left part of the subspace are occupied, and in G, only the cells in the bottom part of the subspace are occupied.
- Libraries H, I, J, and K were created to perfectly occupy the subspace, i.e., one point was added in each cell. The four libraries differ in the grids used to create the sets of points (from around 5000 to 50 000 cells inside the subspace) and, thus, in the number of compounds.

- Libraries L to Q: Points were added in cells of two independent regions of the delimited space leading to libraries of around 2500 compounds. From L to Q, only the distance between the barycenters of the two groups of points was increased.
- Library R: Six points were added on a set of 416 contiguous cells.
- Combined libraries: A series of 15 libraries was designed corresponding to the association of  $x\%$  of randomly selected points taken from the fictive library R and  $(100 - x)\%$  of the fictive library A, with  $x$  varying from 1 to 99.

**Construction of DRCS and Specific Subspaces.** The MOE2D DRCS, as defined in the previous article,<sup>6</sup> was used for all experiments. Briefly, 174 MOE2D descriptors were computed on 30 subsets of 20 000 molecules randomly selected in the database of 6.63 million molecules. These subsets were subsequently used to build a consensus MOE2D PCA model (43.5% of the variance is explained by the two first principal components). The HTS subspace (as defined by the DRCS contour) was subsequently built upon the same subsets using their projection coordinates in the first two Principal Components. Briefly, a subspace is defined by a 2D convex delimitation that is obtained by averaging a set of convex hulls computed on each subset of molecules. Prior to the calculation of each individual convex hull, a certain proportion of molecules (referred to as outliers) is removed. This proportion is determined by identifying the smallest percentage that can be used to obtain a stable shape (see ref 6 for details). The HTS subspace finally encompassed 99.67% of the compounds.

Using the specific sets of compounds (see above) and the MOE2D PCA model, four new subspaces were derived:

1. The “Lipinski+” and “Lipinski–” subspaces, based on molecules that satisfy (respectively do not satisfy) the Lipinski rule of five
2. The “Fragment” subspace, based on molecules that satisfy the rule of three
3. The “Pharmaceutical” subspace, based on the combined library of pharmaceutical compounds

The detailed methodology used to obtain all these contours can be found in the Supporting Information.

**Generalities on Chemical Diversity Indices: What Is a Usable Diversity Index?** On the basis of their definition of a diversity function, Waldman et al.<sup>9</sup> proposed a set of theoretical requirements that a “perfect” diversity function should satisfy:

1. Adding redundant molecules to a system does not change its diversity.
2. Adding nonredundant molecules always increases the diversity of the system.
3. Space-filling behavior of diversity space should be preferred.
4. Perfect (i.e., infinite) filling of a finite descriptor space should result in a finite value for the diversity function.
5. If the dissimilarity or distance of one molecule to all others is increased, the diversity of the system should increase.

We completely agree with these requirements for the definition of an *absolute diversity index* characterizing the overall diversity of a library. But taking only these criteria into account, the most natural way to increase the diversity is to increase the number of compounds. It must be remembered that the general

goal of designing diverse screening libraries is to increase the ratio and the interest of the hits. Hence, from the practical point of view of an experimentalist doing screening, the interest of an index observing rules 1 and 2 is quite debatable. For such an application, a Relative Diversity Index (RDI) describing the mean diversity *per molecule* is certainly more valuable and better adapted. The most natural way to get such a value is to divide the value of the diversity of a library (the absolute diversity index) by the cardinality of this library. It will discriminate two libraries with the same overall diversity but different cardinalities. Moreover, this index will be comparable between two libraries containing a different number of molecules, which is clearly not possible using an absolute diversity index that is expected to increase with the number of compounds, regardless how similar are these compounds. In their second paper, Waldman et al.<sup>12</sup> went into this problem but did not explicitly adapt the previous rules.

For a RDI, rule one must be modified. The addition of redundant molecules, which does not increase the overall diversity but increases the number of compounds, decreases the index (from the experimentalist point of view, adding redundant molecules decreases the overall interest of a chemical library: the diversity will not increase, but the costs will).

Dealing with the overall diversity of a library, rule two is obvious. Any new (i.e., not previously present) compound will increase this diversity, no matter how different it is compared to the existing library. But regarding the relative diversity (i.e., the mean contribution to the diversity per product), the addition of a nonredundant molecule has no obvious effect. The RDI can either increase or decrease depending on the diversity induced by the new product versus the mean diversity of the library. A new compound very different from the previous ones is likely to increase the RDI, but another one very similar to some products is likely to decrease this RDI. Thus, the second rule has no interest for an index such as RDI.

Rules 3, 4, and 5 concern the space coverage and have the same interest for an absolute diversity index as for an RDI. Their application is sometimes rather difficult. One can note that rule 4, “perfect filling...” implies the use of a finite descriptor space, which can be achieved using the DRCS contour. Furthermore, rule 5 also implies that the RDI value must regularly increase when diverse compounds are added to a library.

The use of an RDI implies one other rule. A representative subset of a library should have the same RDI (or a very similar value) as its parent library. There is no unique definition of representativity in chemoinformatics.<sup>27,28</sup> In this work, the statistical definition of representativity is used: one set is representative of another if the two have the same statistical distribution of their properties. The easiest way to obtain such a subset is to perform a simple random selection. Differences can appear between a random sample and the whole library due to the random process, but they decrease when the cardinality of the sample increases. This criterion, which is quite obvious when using basic descriptors, is not always satisfied with the currently used indices, as shown in the results.

Finally a perfect *relative diversity index* should satisfy the following rules:

1. The RDI of a representative subset of a library is equal (or very similar) to that of the whole library.
2. Adding redundant molecules to a system decreases its RDI.



3. Space-filling behavior of diversity space should be favored by the RDI.
4. Perfect (i.e., infinite) filling of a finite descriptor space should result in a finite value of the RDI.
5. If the dissimilarity or distance of one molecule to all others is increased, the RDI of the system should increase.

In the following part of this paper we will analyze the various indices using these rules.

**Indices Implemented.** Various indices have been previously developed to describe a chemical library, especially to describe and quantify their chemical diversity.<sup>9–11,13,14,29–34</sup> Twenty-two of these indices were implemented in our study. Our goal is to select the index(ices) satisfying the rules of an RDI and allowing the best description and characterization of the chemical diversity of libraries through the DRCS methodology.

Six DRCS-free indices were implemented for comparison purposes only. The other 16 indices were applied in combination with the DRCS methodology (in the 2D representation). The indices were classified in four groups:

- Reference indices: Widely used indices that do not depend on the DRCS contour
- Group 1: Index characterizing the number of compounds projected outside the DRCS contour
- Group 2: Noncell based indices applied on DRCS
- Group 3: Cell-based indices applied on DRCS

For all the four groups, each measurement methodology will be further described and compared. It is important to note that, for comparison purposes, all the indices were computed using the compounds projected inside the DRCS contour.

**Reference Indices.** For these indices, the DRCS contours only serve as a way to remove outliers prior to calculation.

**Molecular Scaffolds and Frameworks.** Molecular scaffolds and frameworks are simplified representations of chemical structures.<sup>35</sup> For each chemical library, the ratio between the number of different scaffolds/frameworks and the number of compounds was computed:

- *Scaffolds*: The Bemis and Murcko<sup>35</sup> definition was used, where only rings and linkers between the rings are kept.
- *Frameworks* are based on the scaffolds, but atom types are removed, and bond orders are all set to one. However, unlike the original implementation, we differentiate aromatic from nonaromatic bonds for six-member rings.<sup>3</sup>

For the two indices, an in-house InChI-based<sup>36,37</sup> script was applied to remove duplicate scaffolds and frameworks.

**Fingerprint-Based Indices.** Molecular fingerprints are binary strings that encode the presence of a set of chemical features in a compound. They were previously used<sup>38</sup> for the characterization of chemical libraries. In this study, two fingerprints implemented in Pipeline Pilot<sup>39</sup> called ECFP\_4#S and EPFP\_4#S that take into account the stereochemistry of the compounds were computed. Briefly, ECFP\_4#S generates extended-connectivity fingerprints<sup>38,40</sup> whereas EPFP\_4#S generates Daylight-style path-based fingerprints<sup>41</sup>. In these methodologies, for a given chemical library, each bit of the fingerprint represents a chemical feature found in a chemical structure. Finally, the length of the fingerprints corresponds to the number of different chemical features found in all the compounds of the library.

The following indices were derived for both fingerprints:

- NumFPFeatures: Number of chemical features present in a database scaled by the number of compounds
- AvgDistance: The average Tanimoto distance for all pairs of molecules

**Group 1: External Compounds.** The DRCS contour encompasses 99.67% of the HTS compounds used to compute the DRCS model.<sup>6</sup> However, using the same DRCS (model + contour), the percentage of compounds projected outside the contour can vary significantly depending on the library under consideration. These compounds represent “exotic” molecules with respect to the library used to delimit the subspace, i.e., molecules having a combination of molecular properties not found in the reference library. Thus, for each chemical library, the percentage of compounds projected outside the DRCS contour will be computed. It will be further referred to as the Out(%) index. These external compounds will not be considered for the calculation of all the other indices.

**Group 2: Noncell-Based Indices Applied on DRCS.** For these indices compounds will be characterized by their coordinates in the DRCS model.

**Euclidean Distance-Based Functions.** The Euclidian distances in the reduced 2D space were used to derive the following indices:

1. MeanMinEucDist: The average Euclidean distance between each molecule and its closest neighbor.
2. MeanEucDist: The average Euclidean distance for all pairs of molecules in the library.

These indices require the calculation of all intermolecular distances, which is computationally very expensive (complexity in  $O(N^2)$ ).

**Diversity Integral-Based Indices.** Other distance-based functions exist that are less time consuming. In particular, the Diversity Integral methodology of Cerius2 C<sup>2</sup>-Lib,<sup>42</sup> used by Pascual et al. and others<sup>11,31</sup> is comparable to the calculation of the MeanMinEucDist index. It is based on the average Euclidean distance between random points and their closest compound in the library under consideration.

Similarly, in this work, a set of fixed reference points was used instead of random points to define the Diversity Integral index (DivInt). This implies that the complexity of the methodology drops to  $O(N)$ . For each subspace, a set of 1000 reference points was equally dispersed to cover the entire subspace (see the Supporting Information) and, for each reference point, the minimum distance to a product of the considered library is computed. The average of these distances is calculated to obtain the DivInt index. Hence, the lower the DivInt, the better the distribution of the projected compounds inside the DRCS.

**Group 3: Cell-Based Indices Applied on DRCS.** As mentioned previously, cell-based partitioning is an intuitive way of assessing the space coverage of a particular library. Because chemical space is almost infinite, the space is usually partitioned into a hypercube binned on each descriptor, and the diversity is usually expressed using the proportion of occupied cells. Pascual et al.<sup>11,31</sup> divided the ranges of principal components in such a way that the number of occupied cells is always less than or equal to the number of molecules to select. This leads to a large number of empty cells at the extreme part of the space, which are irrelevant to consider. To avoid the use of parts of the chemical space that are not informative, it is necessary to focus on what Agrafiotis<sup>15</sup> called the “Accessible space” determined by focusing on a representative subspace that minimizes the influence of

outer regions. Cummins et al.<sup>29</sup> proposed to remove the compounds present in the low density space (e.g., cells that are occupied by few compounds) in order to optimize the size of the considered space and of the grid. This allows one to focus on the most representative subspace.

The DRCS contour provides an accurate delimitation of the densest part of a particular delimited subspace. Consequently, a partitioning of this subspace is expected to be more relevant than using a traditional hyper cubic or hyper spherical delimitation. A specific partitioning was therefore applied to each subspace delimited by its corresponding DRCS contour, as described in the following section.

**DRCS-Based Partitioning.** To assess the coverage of a library for a particular partitioned subspace, the following assumption was made: "If the chemical library covers the subspace defined by the DRCS contour in an optimal manner, the proportion of occupied cells must be 100 %". In other words, the partitioning should be defined in such a way that each cell will be filled by one and only one compound in the case of an ideal library. Using a unique grid to compare libraries of different sizes is thus impossible. This leads to one important practical consequence in the present case: The number of cells falling inside the DRCS must be equal to the number of compounds, leading to a different grid for each library of different cardinality.

The first step is therefore to create a 2D grid that strictly encompasses the entire subspace. In the case of a typical squared-shaped subspace, the number of bins  $N_{\text{bins}}$  for each dimension of the grid (which is assumed to be the same) would be determined as

$$N_{\text{bins}} = \sqrt{N_{\text{mol}}}$$

where  $N_{\text{mol}}$  is the number of molecules in the library. In our case however, the shape of each contour is not straight and regular. Thus, encompassing the entire subspace using such a  $N_{\text{bins}} \times N_{\text{bins}}$  grid will inevitably lead to cells falling outside it (Figure 1), and the final number of cells located inside the



**Figure 1.** Simulation of the projection in a DRCS of a library having four clusters of occupied cells. The cells considered to be outside the DRCS are in gray and the DRCS contour is in black. The four clusters of occupied cells are colored differently.

contour will be clearly different from the number of molecules. This suggests that a relationship might exist between the number of bins  $N_{\text{bins}}$  and the final number of cells located inside a particular DRCS contour  $N_{\text{cells-in}}$ . For the HTS contour, a set of grids was computed with  $N_{\text{bins}}$  varying from 10 to 2500. For each grid, the actual number of cells falling inside the contour was computed (a cell is considered as being inside the contour if at least one of its corners is inside the contour). Various regression-based relationships between  $N_{\text{bins}}$  and  $N_{\text{cells-in}}$  were tested. Interestingly, we found that a power regression is able to establish a strong correlation for all the cases (i.e., for all the various subspaces considered in the previous and in the present work). Following this analysis, a power relationship was established and validated on all subspaces to determine the

optimal number of bins based on the number of molecules. Finally, given a subspace  $S$  and a library containing  $N_{\text{mol}}$  molecules, the number of bins in each dimension of the grid is determined as

$$N_{\text{bins}} = \alpha_S \times N_{\text{mol}}^{\beta_S}$$

where  $\alpha_S$  and  $\beta_S$  are the two parameters determined by the power regression obtained on the subspace  $S$ . The coefficient values for each subspace presented herein can be found in the Supporting Information. Consequently, each grid is specific both to the subspace under consideration and to the library to be analyzed. Once  $N_{\text{bins}}$  has been determined, the grid is positioned to strictly encompass subspace  $S$ , and cells located outside the subspace (i.e., cells that have all their corners outside the contour) are subsequently removed. All the indices described below were applied using this optimized grid.

**Filling(%).** The simplest index that can be defined based on a cell-based partitioning is the percentage of occupied cells. Given  $N_{\text{occ}}$  occupied cells among  $N_{\text{total}}$  cells in an optimized grid, the filling percentage (Filling(%)) is defined as

$$\text{Filling(\%)} = \frac{N_{\text{occ}}}{N_{\text{total}}} \times 100$$

It quantifies the overall coverage of a particular subspace, but does not take into account the evenness of this distribution.

**Shannon Entropy.** Shannon entropy (SE) was originally developed for application in digital communication theory.<sup>43</sup> It was transferred to the chemical domain and applied to measure the information of molecular descriptors distribution in compound database.<sup>44–46</sup>

Shannon entropy was also successfully applied to quantify the distribution uniformity of compounds in cells when chemical spaces are partitioned by a grid.<sup>13,31,33</sup> It is defined as

$$\text{SE} = - \sum_{i=1}^{N_{\text{full cells}}} p_i \log_2 p_i$$

where  $p_i = N_i/N_{\text{cpds}}$  with  $N_i$  being the number of compounds in cell  $i$  and  $N_{\text{cpds}}$  the number of compounds in the DRCS.

To compare this SE value between libraries of different sizes (hence different numbers of cells inside the DRCS contour), it has to be scaled. The meaning of the scaled value depends on the reference used for the scaling. If one uses, as Godden and Bajorath<sup>46</sup> did, the total number of cells as reference (total number of bins), one obtains

$$\text{sSE}_{\text{all cells}} = \frac{\text{SE}}{\log_2(N_{\text{cells}})}$$

with  $N_{\text{cells}}$  being the number of cells inside the DRCS contour. Thus,  $\text{sSE}_{\text{all cells}}$  captures the uniformity of the distribution of the compounds through the entire grid (i.e., all the chemical space defined by the DRCS).

The latter index depends on the percentage of occupied cells (i.e., the Filling(%)), thus Shannon entropy can be scaled using the number of occupied cells instead of the total number of cells:

$$\text{sSE}_{\text{occ cells}} = \frac{\text{SE}}{\log_2(N_{\text{occ cells}})}$$

with  $N_{\text{occ cells}}$  being the number of occupied cells inside the DRCS contour. The  $\text{sSE}_{\text{occ cells}}$  captures the uniformity of the

distribution of the compounds only through the occupied cells, i.e., the occupied chemical space defined by the DRCS. These two sSE values vary from zero to one (maximum uniformity of the distribution).

Shanmugasundaram and Maggiora<sup>33</sup> pointed out that the cell-based Shannon entropy treats cells as “positionally independent”. To avoid the loss of this information they thus introduced a new Shannon-like index to measure the uniformity of the distribution of occupied cells along each dimension of the grid. The application of this index to our DRCS gives

$$SE_{PCx} = - \sum_{i=0}^{N_{bins}} p_i \log_2 p_i$$

with  $PCx$  being the principal component  $X$  of the DRCS space and  $p_i = N_{occ\ cells\ i} / N_{occ\ cells}$  where  $N_{occ\ cells}$  is the number of occupied cells inside the DRCS and  $N_{occ\ cells\ i}$  is the number of occupied cells in the  $i^{th}$  bin. Section 6 of Supporting Information illustrates the overall process.

Shanmugasundaram and Maggiora do not scale this entropy value using the number of intervals in each dimension but average the values obtained for all the dimensions. For the application of this index to our DRCS, scaling is mandatory because of the use of grids of variable sizes. But our grid scaling is not trivial because the bins are not all equivalent, i.e., the bins do not all contain the same number of cells. To avoid this problem, the proportion of occupied cells for each bin was used instead of the number of occupied cells for each bin. We thus obtain a new  $SE_{PCx}$  value where  $p_i = p_{cell\ i} / \sum p_{cell\ i}$  and  $p_{cell\ i} = N_{occ\ cells\ i} / N_{cells\ i}$ . This value can be scaled

$$sSE_{PCx} = \frac{SE_{PCx}}{\log_2(\text{total number of intervals of the grid})}$$

The  $sSE_{PCx}$  values were kept for the first two principal components of the DRCS. An additional index  $sSE_{PCmean}$  was also defined as the average value of  $sSE_{PC1}$  and  $sSE_{PC2}$ .

**Cluster Diversity Index.** Data clustering is broadly used to assess the quality/diversity of data sets.<sup>47</sup> The CLIQUE<sup>48</sup> (CLustering In QUEst) algorithm was used to find subspaces of a partitioned space with high density clusters. A simplified version of this algorithm was used to evaluate the number of clusters in our DRCS, where the clusters are based on all the occupied cells and not only on the most occupied ones (the “densest cells”). A cluster is defined as a set of contiguous occupied cells (two cells are contiguous if they have a common face, Figure 1). Only the number of clusters is kept.

The homogeneity of the projection can be characterized by the number of clusters scaled by the number of occupied cells, leading to the Cluster Diversity Index (ClusterDiv):

$$\text{ClusterDiv} = \frac{\text{number of clusters}}{\text{number of occupied cells}}$$

**Kolmogorov–Smirnov Index.** The Kolmogorov–Smirnov (KS) criterion was first applied by Rassokhin and Agrafiotis<sup>10,49</sup> to quantify molecular diversity. The KS criterion measures how well an experimental distribution is approximated by a particular distribution function.<sup>50</sup> It is defined as the maximum value of the absolute difference between two cumulative functions

$$KS = \max_{-\infty < x < +\infty} |P(x) - P^*(x)|$$

where  $P(x)$  is a known cumulative distribution of a uniform sample and  $P^*(x)$  is the experimental cumulative distribution. Thus, the KS criterion measures the extent to which the real distribution deviates from the theoretical one. In our study, the optimal (known) distribution is a uniform distribution, where all molecules are evenly distributed over all the cells of the grid. Because the KS criterion is a measure of dissimilarity,  $KS = 0$  is obtained when the experimental distribution is optimal. As we want to obtain a maximum value when the distribution is optimal, a KS-based diversity index called Dks was defined as

$$Dks = 1 - \max_{-\infty < x < +\infty} |P(x) - P^*(x)|$$

The Manhattan distances between cells of the optimal grid were used to compute the Dks(ManhDist) index. For each possible value of Manhattan distance (between zero and twice the number of bins), the ratio between the number of distances and the total number of distances is considered. Then, the theoretical cumulative distribution is based on the Manhattan distances for all the cells ( $P^*(R_{i(\text{ManhDist})})$ ) and the experimental cumulative distribution on the Manhattan distances for the occupied cells ( $P(R_{i(\text{ManhDist})})$ ). The comparison of the two cumulative distributions leads to the Dks(ManhDist) index

$$Dks(\text{ManhDist}) = 1 - \max_{0 < i < 2N_{bins}} |P(R_{i(\text{ManhDist})}) - P^*(R_{i(\text{ManhDist})})|$$

Since the Dks(ManhDist) implies the calculation of the Manhattan distances for all the occupied cells inside the DRCS contours, it is computationally expensive. Two new Dks indices were created for each principal component of the DRCS model (Dks(PC1) and Dks(PC2)). On the basis of the optimized grid, the ratio between the number of occupied cells and the total number of occupied cells is considered from the first to the last bin of the grid for each PC. Then the experimental ( $P(R_{occ\ cells})$ ) and theoretical ( $P^*(R_{occ\ cells})$ ) cumulative distributions are based on the real and expected distributions of the occupied cells through all the bins (the shape of the DRCS contour is taken into account). The comparison of these two cumulative distributions thus leads to the Dks(PCx) indices

$$Dks(PCx) = 1 - \max_{0 < i < N_{bins}} |P(R_{i(occ\ cells)}) - P^*(R_{i(occ\ cells)})|$$

These Dks(PCx) indices computed on one dimension are faster to compute than the Dks(ManhDist). Section 7 of Supporting Information illustrates the overall process.

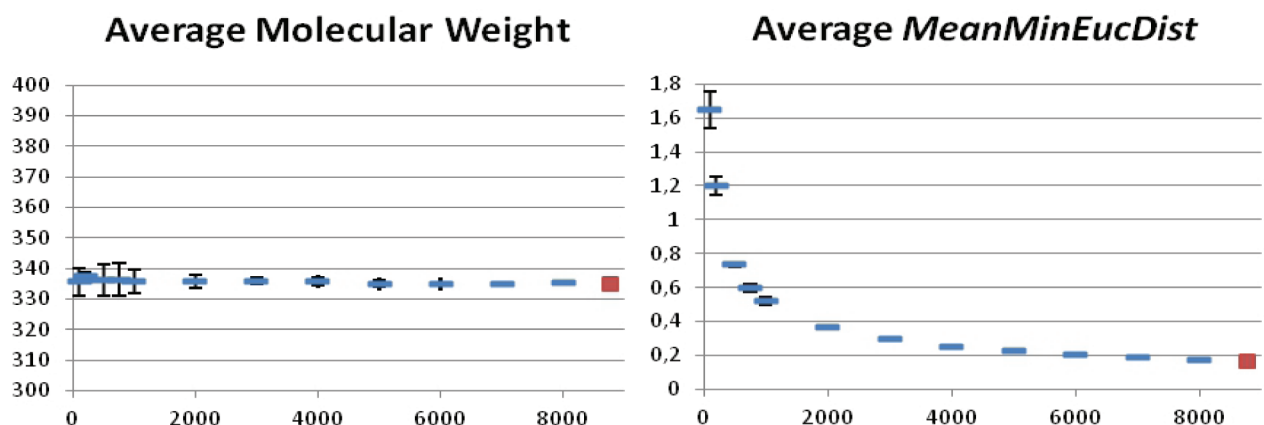
**Manhattan Distance-Based Indices.** Similarly to noncell-based indices, distance-based functions can be applied on the occupied cells. The following indices were thus computed

1. MeanMinManhDist: The average Manhattan distance between each occupied cell and its closest neighbor
2. MeanManhDist: The average Manhattan distance for all pairs of occupied cells

**Final Procedure To Compute DRCS-Based Indices.** From a precomputed DRCS (model + contour) based on a reference library, the following final procedure was used to compute DRCS-based indices of a new library:

1. Standardize the molecules of the studied library





**Figure 2.** Average molecular weight and average MeanMinEucDist for five subsets randomly selected from the CMC database. The size of the subsets varies from 100 to 8000 compounds. The error bars correspond to one standard deviation. The red point indicates the indices values for the whole CMC database.

2. Compute the descriptors corresponding to the DRCS
3. Compute the compounds' coordinates in the DRCS model
4. Test whether the molecules are inside or outside the DRCS contour
5. Compute the optimal number of grid bins
6. Test whether the grid cells are inside or outside the DRCS contour
7. For compounds and grid cells inside the DRCS contour, attribute each compound to a grid cell
8. Compute the DRCS-based diversity indices

## RESULTS AND DISCUSSIONS

Twenty-two indices were implemented to study the diversity of chemical libraries (see Method section). The behavior of these indices was further explored with respect to the rules that a relative diversity index (RDI) should satisfy (see Method section). The objective of this work is to select a few indices based on the DRCS methodology that respect the maximum number of these rules and provide nonredundant information on diversity.

**Representativity of Diversity Indices.** Rule 1 implies that the RDI values of representative subsets of a library must be equal to the RDI value of the entire library. This property was verified by computing and studying the variation of the indices for the 60 random subsets extracted from the CMC library. The CMC database was used in order to study the implemented indices on a real library. Furthermore, the CMC database has one of the highest coverages of chemical space that we have found.<sup>6</sup>

In order to obtain references, two simple indices were studied: the average molecular weight and the MeanMinEucDist. According to the statistical definition of representativity, the average molecular weight must be constant regardless of the cardinality of the subset. On the contrary, the average of the distances of each compound from its closest neighbor (i.e., the MeanMinEucDist index) should obviously decrease when the number of compounds increases. The observations derived from Figure 2 are consistent with these predictions. For each size of subset, the average value of the five subsets was plotted with the error bars limited to one standard deviation. Obviously the error bars for the small sizes of subset are the largest. The average molecular weight remains quite constant (within the error bars) whereas the average MeanMinEucDist shows strong variations. The first index thus satisfies the first rule of an RDI but the second one is unusable.

Globally, all the indices show either of the two variations observed previously. For some of them it is difficult to analyze these variations for the small subsets (under 1000 compounds) because the random variations are too great, but for the subsets between 1000 and 8000 the results are quite clear. All the graphs are presented in Table S2 of the Supporting Information, and the results are summarized in the column "Stability" in Table 1.

**Representativity of NonDRCS-Based Indices.** Among the indices of the reference group, Scaffolds and Frameworks depend on the cardinality of the subsets. For the fingerprint-based indices, NumFPFeature also depends on the size of the subsets, but not AvgDistance. The same holds for both types of fingerprint. AvgDistance is hence the only representative index in the reference group.

**Representativity of DRCS-Based Indices.** The percentage of compounds projected outside the contour (Out(%)) is obviously stable for all the subsets. For the noncell-based indices applied to the DRCS (group 2), MeanMinEucDist, and DivInt vary with the cardinality of the subsets, whereas MeanEucDist remains constant. For cell-based indices (group 3), the Shannon entropy based indices and MeanManhDist are found to be unstable, while all the other cell-based indices are stable (see Table S2 of the Supporting Information).

Finally, among the 22 indices implemented, only 11 (50%) are independent of the cardinality of the subsets and hence satisfy the first rule. The same conclusions were obtained for larger data sets ( $N$  between 20 000 and 100 000) extracted from our in-house database (data not shown).

These results show that many diversity indices, including some DRCS-based indices, are not appropriate for comparing libraries of different cardinalities. In the remainder of this section, only the DRCS-based indices that satisfy rule 1 of an RDI will be considered.

**Behavior of DRCS-Based Indices.** To examine rules 2 to 5,  $N$  fictive library projections were used to represent various extreme cases (see Method section). As these libraries were created inside the DRCS contour, the Out(%) index is systematically 0, and only the seven remaining stable indices were analyzed. The results are summarized in Table 2.

**Redundancy.** Fictive library A has an optimum coverage of the DRCS. On the basis of it, libraries B, C, and D were assembled by duplicating library A two, four, or eight times. Thus, in libraries B, C, and D, each product is present in two, four, or eight copies.

Table 1. Summary of Results for Rule 1 of RDI<sup>a</sup>

group	description of the indices	name of the indices	rule 1 of an RDI: stability
REF	percentage of scaffolds	Scaffolds	—
	percentage of Frameworks	Frameworks	—
	percentage of chemical features in molecular fingerprints ECFP_4#S	NumFPFeatures (ECFP_4#S)	—
	percentage of chemical features in molecular fingerprints EPFP_4#S	NumFPFeatures (EPFP_4#S)	—
	average Tanimoto distance for ECFP_4#S molecular fingerprints	AvgDistance (ECFP_4#S)	+
	average Tanimoto distance for EPFP_4#S molecular fingerprints	AvgDistance (EPFP_4#S)	+
1	percentage of compounds outside the DRCS contour	Out(%)	+
2	mean of minimal Euclidean distances	MeanMinEucDist	—
	mean of Euclidean distances	MeanEucDist	+
	diversity integral	DivInt	—
3	Filling percentage for optimized grid	Filling(%)	+
	scaled Shannon entropy with reference to the total of cells	sSE <sub>all cells</sub>	—
	scaled Shannon entropy with reference to the occupied cells	sSE <sub>occ cells</sub>	—
	scaled Shannon entropy for PC1	sSE <sub>PC1</sub>	—
	scaled Shannon entropy for PC2	sSE <sub>PC2</sub>	—
	mean of Shannon entropy for all PC	sSE <sub>PCmean</sub>	—
	scaled number of clusters for the occupied cells	ClusterDiv	+
	Kolmogorov–Smirnov diversity based on Manhattan interdistances of occupied cells	Dks(ManhDist)	+
	Kolmogorov–Smirnov diversity based on the occupied cells distribution on the PC1	Dks(PC1)	+
	Kolmogorov–Smirnov diversity based on the occupied cells distribution on the PC2	Dks(PC2)	+
	mean of minimal Manhattan distances	MeanMinManhDist	+
	mean of Manhattan distances	MeanManhDist	—

<sup>a</sup>+: Stable index. —: Size-dependent index.

A good index should discriminate between these four fictive libraries, ranking them in the ABCD order. As shown in Table 3, the Filling(%) value decreases with the redundancy (the small differences obtained compared to the expected values

[100, 50, 25, and 12.5 for A, B, C, and D, respectively] stem from the rounding of  $N_{\text{bins}}$  induced by the power regression). The MeanEucDist and Dks indices do not seem to change. ClusterDiv and MeanMinManhDist values increase, while ClusterDiv seems to reach a plateau for libraries C and D. This shows that MeanEucDist, ClusterDiv, and the Dks indices do not satisfy rule 2.

**Space-Filling.** In fictive libraries E, F, and G, 50% of the cells are occupied. In library E, the occupied cells are regularly distributed on the grid. In library F and G, the cells are filled in the left and bottom part of the subspace, respectively (Figure 3).

With the same number of occupied cells, fictive library E is obviously better than F or G. The indices that differentiate these situations will thus satisfy rule 3 of an RDI. Results are given in Table 4. For all three fictive libraries, the Filling(%) is constant (the small difference for library G comes from an artifact of the construction of the fictive library as explained above). ClusterDiv, Dks(ManhDist), and MeanMinManhDist have equal values for F and G but a different one for E. The combination of the two Dks(PC) indices makes it possible to differentiate the three libraries. Indeed, for these indices we have the maximum value ( $Dks(PC) = 1$ ) when the points are equally projected along each axis and the middle value ( $Dks(PC) = 0.5$ ) when the points are projected only on half of the axis. Finally, only MeanEucDist gives three different values for the three libraries.

Rule 3 of an RDI implies that the size of the empty regions of a chemical space must be minimized. With the same coverage, libraries F and G have one empty region representing 50% of the space, whereas library E has the same proportion of empty regions but spread out through all the delimited space. Thus, because they differentiate library E from the other two, all the indices except Filling(%) satisfy rule 3.

**Perfect-Filling.** On the basis of a given number of compounds, fictive libraries A, H, and I were created to simulate perfect filling of the subspace. This perfect filling is defined as the completion of all the cells of the optimal grid in the finite and delimited space. All the indices have finite values for these libraries (Table 5) and, hence, satisfy rule 4 of an RDI. It can be noted that some of the indices are constant (Filling(%), Dks, and MeanMinManhDist). Among the others, ClusterDiv tends to zero because in all the cases there is only one cluster but each time the number of compounds increases. MeanEucDist also decreases, but this does not question the stability (rule 1) of this index.

**Dissimilarity.** Rule 5 of an RDI specifies that when the distance between the compounds (their dissimilarity) increases,

Table 2. Summary of Results for Stable DRCS-Based Indices According to Rules 2 to 5 of RDI<sup>a</sup>

group	name of the indices	rules of an RDI				
		2	3	4	5	
		redundancy	space-filling	perfect filling	dissimilarity	monotony
1	Out(%)	not appropriate for describing the compounds inside the DRCS				
2	MeanEucDist	—	+	+	+	+
3	Filling(%)	+	—	+	—	+
	ClusterDiv	—	+	+	—	—
	Dks(ManhDist)	—	+	+	—	±
	Dks(PC1)	—	+	+	—	±
	Dks(PC2)	—	+	+	—	±
	MeanMinManhDist	+	+	+	—	—

<sup>a</sup>+: Index satisfying the rule. —: Index not satisfying the rule. ±: Index satisfying the rule with some applicability limits.



Table 3. Values of Eight Stable Indices for the Fictive Libraries A, B, and C

fictive library	Mean EucDist	Filling(%)	ClusterDiv	Dks (ManhDist)	Dks (PC1)	Dks (PC2)	MeanMin ManhDist
A	20.84	100.00	0.04	1.00	1.00	1.00	1.00
B	20.84	49.73	18.03	1.00	0.99	0.99	1.01
C	20.84	24.89	100.00	0.99	0.99	0.99	2.00
D	20.83	12.51	100.00	0.99	0.99	0.99	2.49

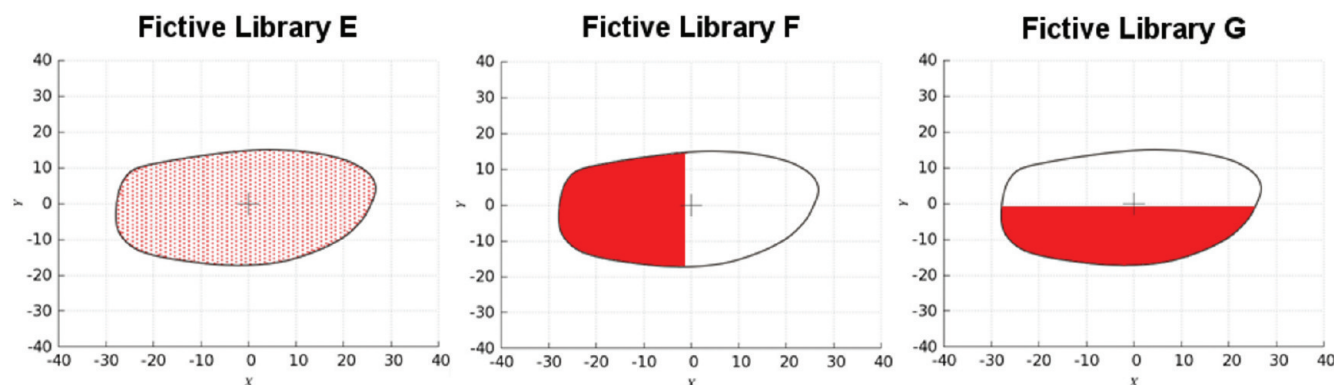


Figure 3. Projections of the fictive libraries E, F, and G inside the DRCS.

Table 4. Values of Seven Stable Indices for the Fictive Libraries E, F, and G

fictive library	Mean EucDist	Filling(%)	ClusterDiv	Dks (ManhDist)	Dks (PC1)	Dks (PC2)	MeanMin ManhDist
E	20.86	50.16	100.00	0.99	1.00	1.00	2.00
F	14.22	50.16	0.08	0.79	0.50	0.97	1.00
G	17.52	49.09	0.08	0.78	0.96	0.49	1.00

Table 5. Values of Seven Stable Indices for the Fictive Libraries A, H, I, J, and K<sup>a</sup>

fictive library	number of fictive compounds	Mean EucDist	Filling (%)	ClusterDiv	Dks (ManhDist)	Dks (PC1)	Dks (PC2)	MeanMin ManhDist
A	2540	20.84	100.00	0.04	1.00	1.00	1.00	1.00
H	4930	20.75	100.00	0.02	1.00	1.00	1.00	1.00
I	9982	20.68	100.00	0.01	1.00	1.00	1.00	1.00
J	19956	20.62	100.00	0.00	1.00	1.00	1.00	1.00
K	49779	20.57	100.00	0.00	1.00	1.00	1.00	1.00

<sup>a</sup>The number of fictive compounds in the library is indicated.

the diversity must increase. Six fictive libraries with the same coverage of the chemical space were created to study this rule (Figure 4). For each library, the compounds were separated into two groups covering two small regions of the delimited space. From L to Q, the distance between the barycenters of the two groups increases, thus according to rule 5 the diversity of the libraries increases.

Filling(%), ClusterDiv, Dks(PC2), and MeanMinManhDist are constant for the six libraries (Table 6). This means that they are not able to characterize alone the difference of diversity. Note that Dks(PC2) is constant because the projection of the compounds on the *y*-axis of the DRCS does not change from L to Q. Dks(PC1) and Dks(ManhDist) increase and rapidly reach a plateau. This limit is clearly due to the very specific examples chosen, but this demonstrates that the ability of the Dks indices to measure the homogeneity of the compounds projection inside the DRCS is limited. Finally, MeanEucDist increases from L to Q. It is the only index that strictly satisfies rule 5 of an RDI.

The previous fictive libraries provide some indication about the behavior of the computed indices for such extreme cases, but no indication is available on their evolution between these cases. One consequence of rule 5 of an RDI is that the

evolution of the indices must be monotone, i.e., the indices values must increase regularly from the least to the most diverse case.

A series of 17 fictive libraries was created. It corresponds to the evolution from the least diverse library R to the most diverse case A (Figure 5). Results of the indices are given in Table 7. Figure 6 shows the evolution of the seven stable indices values versus the percentage of compounds of the fictive library A in the combined libraries.

Clearly, Filling(%) has a linear evolution from the fictive libraries R to A. It well characterizes the linear increase in the chemical space coverage by the fictive combined libraries.

MeanEucDist and all the Dks indices have a non linear but monotone evolution. The combined libraries containing more than 40–60% of the library A (Filling(%) > 50–60%) have similar values of the Dks indices. These indices are thus difficult to use for comparing chemical libraries having Filling(%) value higher than 50%. This is not a problem, however, because this is not the case for the majority of real chemical libraries.

Finally, ClusterDiv and MeanMinManhDist indices form bell curves. Their values increase up to a maximum and then decrease with the diversity. These indices thus do not have a monotone behavior. Furthermore, the possibility of having the

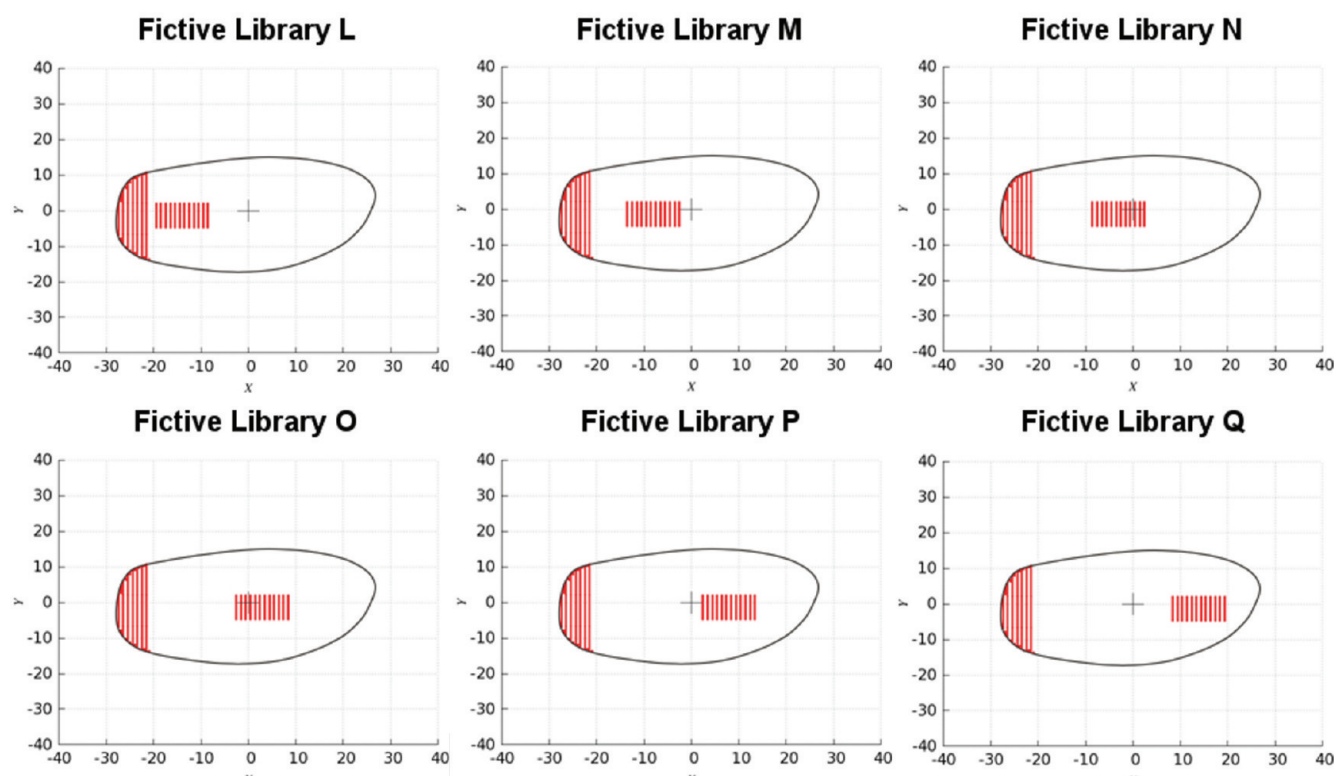


Figure 4. Projections of the fictive libraries L to Q inside the DRCS.

Table 6. Values of Seven Stable Indices for the Fictive Libraries L to Q

fictive library	Mean EucDist	Filling(%)	ClusterDiv	Dks (ManhDist)	Dks (PC1)	Dks (PC2)	MeanMin ManhDist
L	9.37	15.71	0.50	0.51	0.35	0.81	1.00
M	11.99	15.71	0.50	0.61	0.46	0.81	1.00
N	14.32	15.71	0.50	0.71	0.46	0.81	1.00
O	17.19	15.71	0.50	0.79	0.46	0.81	1.00
P	19.62	15.71	0.50	0.79	0.46	0.81	1.00
Q	22.55	15.71	0.50	0.79	0.46	0.81	1.00

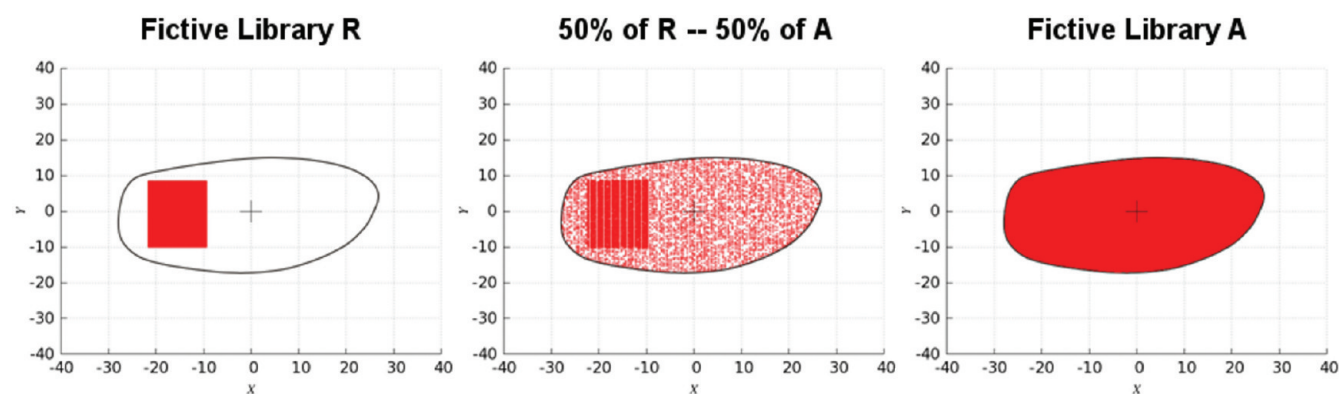


Figure 5. Projections of the fictive libraries R and A and of one combined library comprising 50% of compounds from the fictive libraries R and A.

same value for two differently diverse projections is clearly not acceptable.

So, Filling(%), MeanEucDist, and Dks indices satisfy the consequence of rule 5 concerning their evolution (see column "Monotony" Table 2). This result shows the importance of testing the monotony of the evolution of the indices values, and whether this evolution is similar to that of the diversity.

**Selection of the DRCS-Based Diversity Indices.** Sixteen DRCS-based indices were implemented and tested against the rules that an RDI should satisfy. Tables 1 and 2 summarize the results. None of the implemented indices satisfy all the rules. Some indices do not satisfy some rules in a critical way, e.g., the nonstable and nonmonotone indices. Other indices are characteristic of some rules and do not satisfy other rules

Table 7. Values of Seven Stable Indices for the Fictive Libraries R, A and the Combined Libraries from R to A

fictive library	Mean EucDist	Filling (%)	ClusterDiv	Dks (ManhDist)	Dks (PC1)	Dks (PC2)	MeanMin ManhDist
R (100–0)	8.34	16.38	0.24	0.48	0.33	0.83	1.00
99–1	8.56	17.50	4.58	0.53	0.36	0.84	1.22
95–5	9.53	20.87	17.67	0.70	0.48	0.88	1.37
90–10	10.55	24.95	26.67	0.80	0.59	0.91	1.36
80–20	12.70	33.31	30.20	0.91	0.73	0.93	1.25
70–30	14.37	41.40	26.00	0.95	0.81	0.95	1.16
60–40	16.10	50.15	17.44	0.97	0.87	0.96	1.09
55–45	16.73	54.13	13.62	0.97	0.89	0.97	1.06
50–50	17.32	58.01	10.91	0.98	0.90	0.98	1.05
45–55	17.98	62.38	7.60	0.99	0.92	0.98	1.04
40–60	18.55	66.21	4.37	0.99	0.93	0.99	1.02
30–70	19.39	74.50	1.29	0.99	0.96	0.99	1.01
20–80	20.06	82.36	0.23	1.00	0.98	1.00	1.00
10–90	20.44	90.86	0.04	1.00	0.99	1.00	1.00
5–95	20.63	95.21	0.01	1.00	1.00	1.00	1.00
1–99	20.67	99.03	0.01	1.00	1.00	1.00	1.00
A (0–100)	20.84	100.00	0.04	1.00	1.00	1.00	1.00

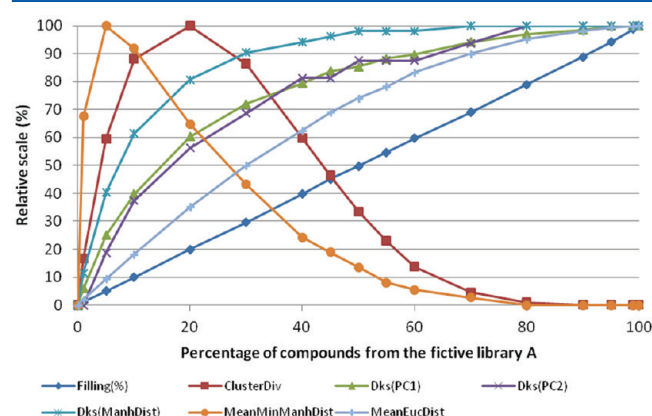


Figure 6. Evolution of seven stable indices computed on combined fictive libraries between R and A (the indices values are indicated on a relative scale between 0 and 100%).

without altering their interpretability, e.g., the Filling(%) characterizes the coverage of the delimited subspace. Thus, depending on their characteristics, a selection of indices can be made to produce a simple description of the chemical diversity of a library with regard to the DRCS used. In our opinion, these indices must be used in the following order if they are to be interpreted clearly:

1. Out(%): Calculation of the proportion of external compounds with regard to the DRCS used. A chemical library has to be diverse in a given subspace, and it should also contain external compounds to explore other subspaces. The calculation of this proportion of external compounds is mandatory as it enables one to know the proportion of compounds on which the other indices are computed. That is why the more compounds projected outside the DRCS contour there are, the less usable the computed indices are. No strict rule can be implemented, but we consider that the indices are representative and interpretable if the library contains less than 10% of compounds outside the DRCS contour used. If Out(%) is higher than 10%, the DRCS used is maybe ill-adapted. Further, the other selected indices will not be computed in this situation.

2. Filling(%): Coverage of the studied delimited subspace. This second index gives the proportion of the subspace that the library explores. The higher the Filling(%), the higher the diversity.
3. Dks: Homogeneity of the distribution of the occupied cells inside the delimited subspace. For similar Filling(%) values, two libraries can have very different distributions of the occupied cells. The Dks indices enable one to know whether the occupied cells are grouped together or not. Dks(ManhDist) index gives information about the homogeneity of the distribution of the occupied cells in all the DRCS. Dks(PC1 and 2) indices are faster to compute and provide the same information but through each axis of the DRCS model. They thus offer a greater description of the diversity. That is why they will be further used to characterize libraries. Nevertheless, Dks(ManhDist) could be preferred for an automatic application to optimize the chemical diversity of libraries.
4. MeanEucDist: Homogeneity of the distribution of the compounds inside the delimited subspace. Similar Dks values can mask a nonhomogeneous distribution of the compounds inside the DRCS. This information can be obtained with the MeanEucDist index. However, it depends on the calculation of the Euclidean distances for all the pairs of compounds, which implies a high computational cost. Thus, this index will be used if the previous indices are unable to differentiate the compared libraries.

The previous selected diversity indices were implemented in radar graphs, with which various chemical libraries can be easily compared with regard to the computed indices.

**Specific DRCS Contours.** DRCS-based diversity indices characterize the coverage and the homogeneity of the projection of a chemical library inside a DRCS. They are representative of the diversity with respect to the contour used and thus to the chemical subset used to construct the contour. It is therefore important to know which DRCS contour has to be used, e.g., using the HTS contour for characterizing fragment subsets is totally inappropriate.

On the basis of the DRCS model and on the databases prepared (see Method section), four specific contours were computed (see Supporting Information). The initial DRCS contour



(representative of the HTS compounds) and the other four specific contours are shown in Figure 7.

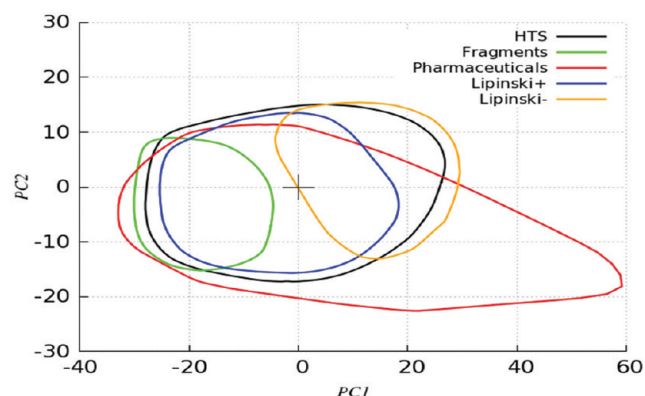


Figure 7. Multiple contours representation in the DRCS model.

In the same way as for the initial DRCS contour, the use of representative contours instead of compound projections facilitates the visualization and the characterization of chemical libraries. The projection of a library means that the coverage of different chemical subspaces represented by specific contours can be rapidly estimated, and unexplored zone(s) in these subspaces identified.

Moreover, because the contours represent specific chemical subspaces, the spaces can be compared without the need to project all the compounds. Figure 7 shows that the subspace of pharmaceutical compounds (red contour) is much broader than the subspace of HTS commercial compounds (black contour). Likewise, the subspace of fragment compounds (green contour) clearly occupies only a small part of the pharmaceutical or HTS commercial compounds spaces. The blue and orange contours represent the subspaces of the compounds that respectively satisfy or do not satisfy the “rule of five”. One region of the space is common to these two contours. This fuzzy limit between the two subspaces is the consequence of the exception accepted in the “rule of five”.

The distribution of a chemical library in different contours reflects different occupations and chemical diversities with reference to each contour. Consequently, the DRCS-based indices will not be equivalent but depend on the DRCS contour used. The contour will modify the interpretation of the diversity indices.

**Application to the Characterization of Libraries.** The DRCS-based methodology (graphical tool and diversity indices) was applied to six publicly available chemical libraries: Prestwick Chemical Library, Chembridge Kinaset, two combinatorial libraries (CL-A3 and CL-B5), Pyxis Discovery Smart Fragment Library, and EPAFHM (see Method section for description). This application will illustrate how this methodology should be used and what kind of information it provides.

**Compounds Projection.** The six libraries were projected in the DRCS. Figure 8 shows the compounds projections in the five precomputed contours and allows their visual comparison. The Prestwick collection has a substantial coverage of the HTS contour. Moreover, the majority of the compounds projected outside the DRCS contour appear to be in the pharmaceuticals subspace. The compounds of the Chembridge Kinaset collection are concentrated in a reduced part of the HTS or pharmaceuticals contours. This is what is expected for target-specific libraries. The two combinatorial libraries have the same

behavior as the Chembridge Kinaset library. They focus on a small part of the HTS contour. The two libraries contain the same number of compounds (1000), but visually the A3 combinatorial library seems to occupy a higher subspace than the B5 library. The Pyxis collection of fragment compounds homogeneously occupies a small part of the HTS or pharmaceuticals subspaces but seems to cover a high part of the fragments contour. The EPAFHM collection also seems cover a high part of the fragments subspace with a high proportion of compounds outside this subspace.

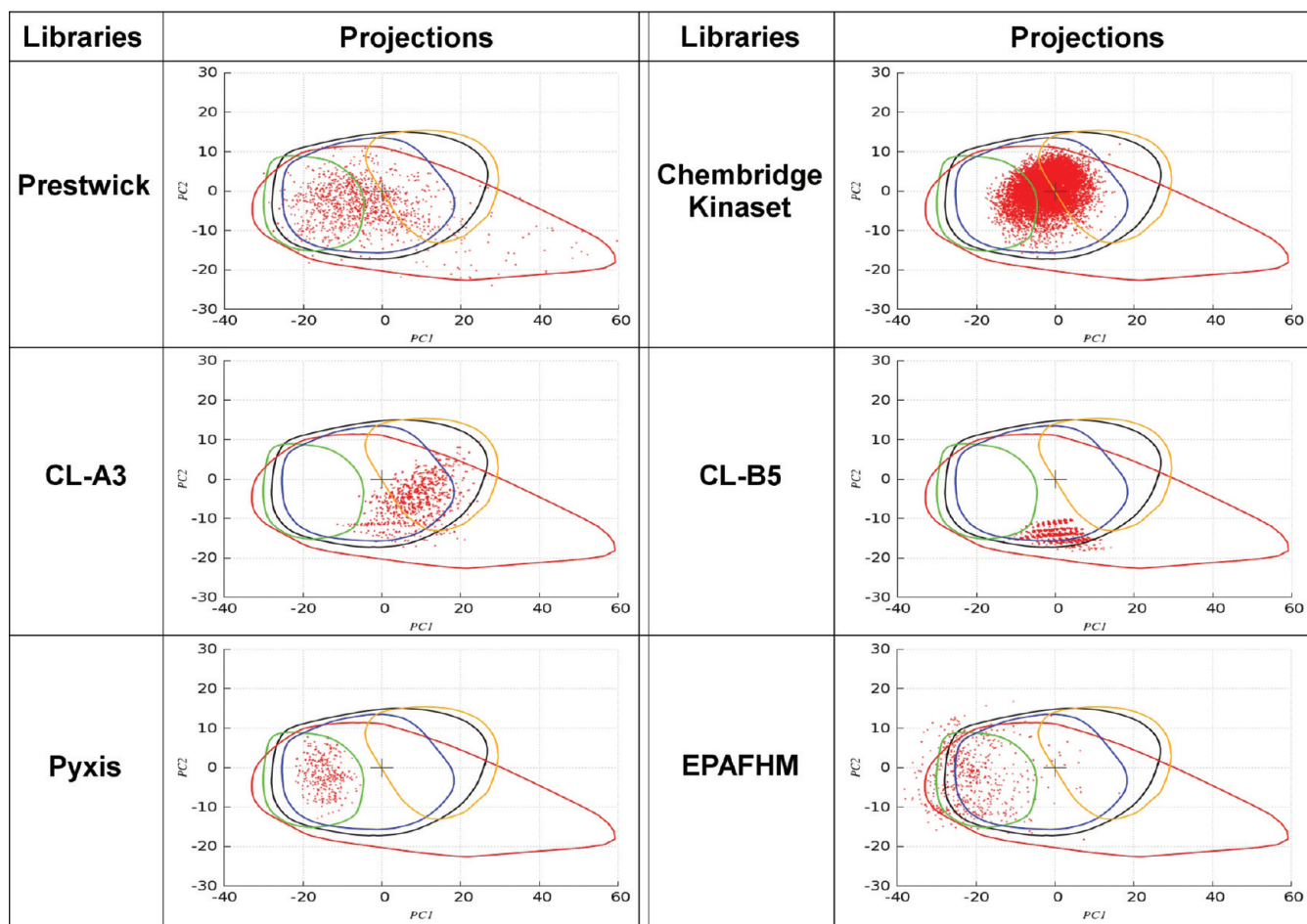
**DRCS-Based Diversity Indices Calculation.** The five selected DRCS-based diversity indices were computed for the same six chemical libraries. The DRCS model and four representative contours (HTS, Pharmaceuticals, Lipinski+ and Fragments; see Methods section for description) were used. Figures 9–12 show the four resulting radar graphs (one for each contour).

It is important to note that the EPAFHM library is not represented in the radar graphs. This is because of the large proportion of compounds projected outside the contour (>10%) for all subspaces. In the light of its projection shown in Figure 8, the applicability domains of QSPR models built using this reference library could obviously be questioned. This library indeed contains a high proportion of very small compounds compared to the other projected libraries, or to the libraries used to build the specific contours. As this seems to demonstrate that none of the subspaces are appropriate to quantify the diversity of this library, it will not be further studied.

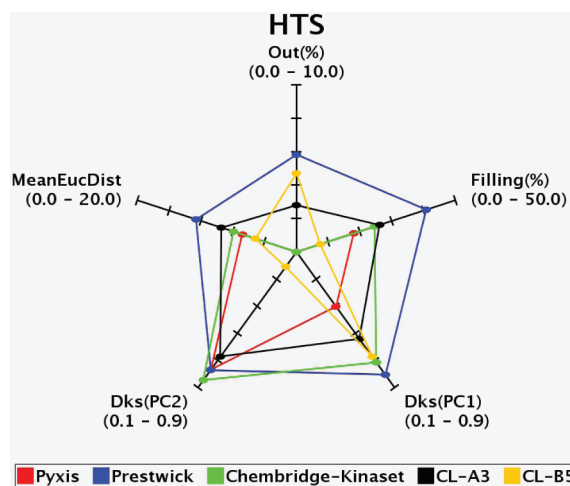
**HTS Contour.** The five libraries have less than 10% of compounds projected outside the HTS contour with the Pyxis and Chembridge Kinaset collections having 100% of the compounds projected inside the contour. Thus, these two libraries do not have distinct compounds with respect to the HTS commercial compounds.

The Filling(%) shows that the Prestwick library has the highest coverage of the HTS subspace. With 5.81% of external compounds and 40.85% of coverage, this library is the most diverse and covers the subspace of HTS commercial compounds quite well. The Pyxis and CL-B5 libraries, on the contrary, have a low coverage of the subspace (18.13% and 7.57% respectively). They cover the two axes differently, as shown by the significant differences found for the Dks indices. The Pyxis library has a better coverage of the PC2 axis ( $Dks(PC1) = 0.79$ ) compared to the PC1 axis ( $Dks(PC2) = 0.42$ ), while the opposite is true for CL B5, which has a better coverage on the PC1 axis ( $Dks(PC1) = 0.72$ ) than on the PC2 axis ( $Dks(PC2) = 0.19$ ). These differences are a good illustration of the complementarity of the selected indices, which account for different information that would not be captured by a single index. Chembridge Kinaset and CL-A3 libraries have similar coverage of the subspace but the  $Dks(PC1)$  and  $PC2$  indices indicate that Chembridge Kinaset library has a better distribution on each axis. Nevertheless, CL-A3 has a higher value for the MeanEuDist index compared to the Chembridge Kinaset library, indicating a better homogeneity of the coverage in the occupied regions. This confirms the visual interpretation of their projection in the DRCS (Figure 8).

**Pharmaceuticals Contour.** The libraries have a similar profile to that in the HTS contour (Figure 10). However, differences can be observed in the  $Dks(PC1)$  and  $PC2$  indices. The distribution of the occupied cells in the Prestwick and CL-A3 libraries is the same through each axis of the DRCS model, whereas the

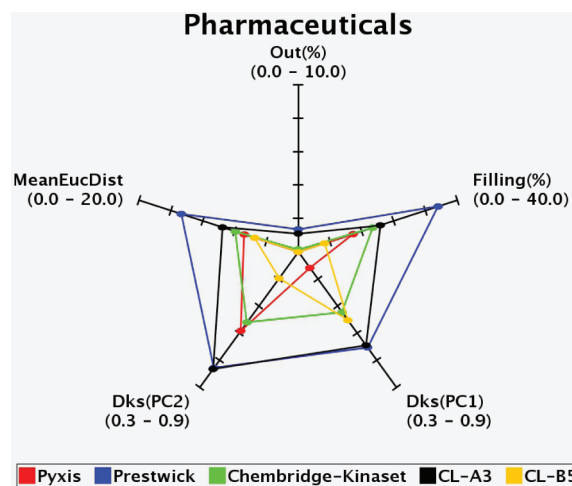


**Figure 8.** Projections of chemical libraries on the DRCS. Compounds are shown by red dots. Contours are presented in the same way as in Figure 7.



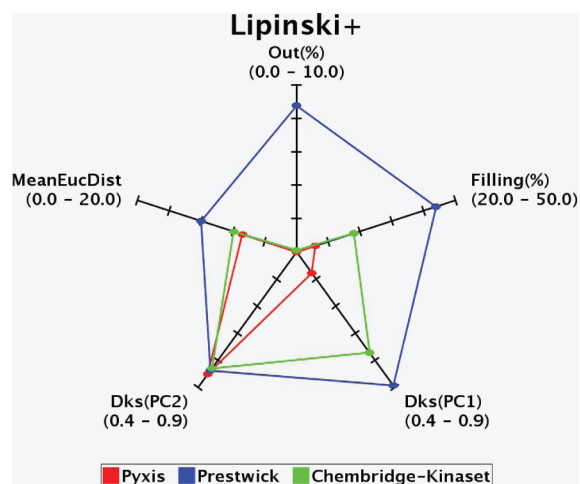
**Figure 9.** Radar graph of the DRCS-based diversity indices computed for five chemical libraries. The DRCS model and the HTS commercial compounds contour were used for all the calculations.

homogeneity of Chembridge Kinaset is lower. It can also be noted that in comparison with the HTS contour, the values of the indices are lower in the pharmaceuticals contour for all the libraries. All these differences reflect the incapacity of the libraries to cover the region of pharmaceutical compounds with high values on the PC1 axis (Figure 8).



**Figure 10.** Radar graph of the DRCS-based diversity indices computed for five chemical libraries. The DRCS model and the pharmaceuticals contour were used for all the calculations.

*Lipinski+ Contour.* It does not appear to be suitable for studying libraries CL-A3 and CL-B5 (14.00% and 24.80% of external compounds, respectively (Figure 11)). With a high proportion of external compounds (8.76%), the Prestwick collection has the best coverage of this subspace (46.33%), while Chembridge Kinaset and Pyxis libraries have a lower diversity. It can be noted that differences in the coverage of the subspace

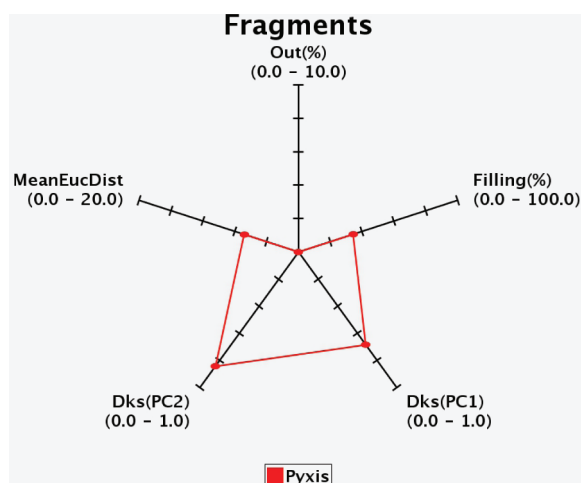


**Figure 11.** Radar graph of the DRCS-based diversity indices computed for three chemical libraries. The DRCS model and the Lipinski+ contour were used for all the calculations.

by the three libraries impacts only on their distribution through PC1 (similar values for Dks(PC2)).

In this case, MeanEucDist calculation is not mandatory, as the previous indices clearly show the differences between the libraries.

**Fragments Contour.** Only the Pyxis library has less than 10% of external compounds with regard to the fragments contour (Figure 12). This reflects the high specificity of this subspace for



**Figure 12.** Radar graph of the DRCS-based diversity indices computed for the Pyxis library. The DRCS model and the fragments contour were used for all the calculations.

the fragments compounds. Here, the Pyxis library has a quite high coverage of this subspace (34.50%), comparable to the coverage of the HTS subspace by the Prestwick collection (Figure 9). Finally, in contrast to the other subspaces, the Pyxis library has a similar distribution through the two axes of the DRCS model. This shows that this library is clearly more diverse and better suited to this subspace.

It should be noted that the Lipinski- contour was not used. All the chemical libraries studied have more than 10% of external compounds with respect to this contour (in the best case CL-A3 has 26.30% of compounds outside this contour). This is a typical example of a subspace that, in most cases, will

not be useful to estimate the diversity of chemical libraries, as most screening libraries are designed to contain drug-like molecules. Such a subspace may rather be used to estimate how a particular library covers an undesired subspace, and if this coverage is acceptable for the problem at hand.

## CONCLUSION AND PERSPECTIVES

In a previous paper, the DRCS methodology was introduced as the combination of a PCA model and a subspace delimitation. A subspace, as defined by the DRCS contour, provides a relevant delimitation of the densest zone of the chemical space spanned by the reference library. In this work, the DRCS methodology was used to create a new set of subspaces representing specific types of compounds and to derive new diversity indices that are independent of the cardinality of the library. As it obviously does not make much sense to, for example, assess the space coverage of a fragment library in general-purpose HTS subspace, we showed that the DRCS methodology can be used to focus on the relevant part of the chemical space that corresponds to the problem at hand. The general-purpose Lipinski, Pharmaceutical, and Fragment subspaces were presented, and we showed that despite some significant overlaps, they all span a different zone of the chemical space.

The DRCS contour was also used to obtain a relevant and flexible partition of the corresponding subspace, making it possible to derive diversity indices. The behavior of reference as well as DRCS-based diversity indices was assessed in the light of a new set of rules that a diversity index should satisfy when comparing libraries of different cardinalities. Following this analysis, it was found that, even in simplified fictive situations, none of these indices can account for all the important aspects that characterize the diversity of screening libraries (space coverage, redundancy, uniformity of the distribution...). Hence, five complementary indices were finally selected that account for these different aspects of diversity, and a simple framework is proposed to use them effectively. These indices can be applied to compare libraries containing different numbers of molecules (e.g., after an enrichment operation or to decide between two libraries proposed by external vendors), as well as to other classical diversity problems (e.g., automatic library design).

In conclusion, the methodology proposed in these two papers provides complementary visual and numerical tools that can be used to analyze the diversity of chemical libraries. The overall idea can be summarized as the following: intuitively, the real chemical space is defined by both very dense and very sparse regions. Dense regions typically contain the large majority of available compounds, corresponding to widely studied chemotypes. In contrast, sparse regions generally contain specific and somewhat exotic molecules, which might nevertheless be included in a screening campaign as well, depending on the context. It becomes evident that the sampling of the two regions should not be performed the same way, and the delimitation of the densest subspace represents a simple way of tackling this issue. The diversity indices as well as the visual characterization provides a way to assess the coverage of a library with respect to the most explored region of a reference library, leaving the sparse regions to a separate analysis. The use of a 2D representation obtained by PCA projection is certainly the most limiting factor of the methodology. The easiest way to overcome this is either to extend the analysis to a third dimension or to perform the same analysis using different Principal Components, although at the expense of simplicity and ease of interpretation. Another possibility would be the use of nonlinear



multidimensional scaling techniques, such as Generative Topographic Maps.<sup>26</sup> These nonlinear methods would be especially useful to define biologically relevant subspaces, e.g., target-specific subspaces, without losing the advantage of visual analysis.

## ■ ASSOCIATED CONTENT

### ■ Supporting Information

S1: Schematic view of the creation of the fictive libraries A to K. S2: Creation of specific subspaces. S3: Determination of the mathematical relation between  $N_{bins}$  of the optimal grid and  $N_{mol}$  of a library for a given DRCS. S4: *DivInt* index: Definition of the reference points. S5: Behavior of DRCS-based indices: rule 1, influence of the size of the library. S6: Illustration of the overall process of Shannon entropy calculation on PC1. S7: Illustration of the overall process of the Kolmogorov–Smirnov index on PC1. Description of tools and computational performances. This material is available free of charge via the Internet at <http://pubs.acs.org>.

## ■ AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [luc.morin-allory@univ-orleans.fr](mailto:luc.morin-allory@univ-orleans.fr).

## ■ ACKNOWLEDGMENTS

The authors thank Accelrys for providing, free of charge, the software “Pipeline Pilot Student edition”, and the “Conseil Régional du Centre” for supporting this research. The authors also thank Peter Schmidtko for helpful comments on the manuscript. V.L.G. thanks the “Conseil Général du Loiret” for funding his Ph.D.

## ■ REFERENCES

- (1) Sukuru, S. C.; Jenkins, J. L.; Beckwith, R. E.; Scheiber, J.; Bender, A.; Mikhailov, D.; Davies, J. W.; Glick, M. Plate-based diversity selection based on empirical HTS data to enhance the number of hits and their chemical diversity. *J. Biomol. Screen.* **2009**, *14* (6), 690–699.
- (2) Sirois, S.; Hatzakis, G.; Wei, D.; Du, Q.; Chou, K. C. Assessment of chemical libraries for their druggability. *Comput. Biol. Chem.* **2005**, *29* (1), 55–67.
- (3) Monge, A.; Arrault, A.; Marot, C.; Morin-Allory, L. Managing, profiling and analyzing a library of 2.6 million compounds gathered from 32 chemical providers. *Mol. Diversity* **2006**, *10* (3), 389–403.
- (4) Dubois, J.; Bourg, S.; Vrain, C.; Morin-Allory, L. Collections of compounds: How to deal with them? *Curr. Comput.-Aided Drug Des.* **2008**, *4* (3), 156–168.
- (5) Medina-Franco, J. L.; Martinez-Mayorga, K.; Giulianotti, M. A.; Houghten, R. A.; Pinilla, C. Visualization of the chemical space in drug discovery. *Curr. Comput.-Aided Drug Des.* **2008**, *4* (4), 322–333.
- (6) Le Guilloux, V.; Colliandre, L.; Bourg, S.; Guenegou, G.; Dubois-Chevalier, J.; Morin-Allory, L. Visual characterization and diversity quantification of chemical libraries: 1. Creation of delimited reference chemical subspaces. *J. Chem. Inf. Model.* **2011**, *51* (8), 1762–1774.
- (7) Gillet, V.; Dean, P.; Lewis, R. Background Theory of Molecular Diversity. In *Molecular Diversity in Drug Design*; Springer: Netherlands: 2002; pp 43–66.
- (8) 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* (6), 501–506.
- (9) Waldman, M.; Li, H.; Hassan, M. Novel algorithms for the optimization of molecular diversity of combinatorial libraries. *J. Mol. Graphics Modell.* **2000**, *18* (4–5), 412–426, 533–536.
- (10) Agrafiotis, D. K.; Rassokhin, D. N.; Lobanov, V. S. Multi-dimensional scaling and visualization of large molecular similarity tables. *J. Comput. Chem.* **2001**, *22* (5), 488–500.
- (11) Rabal, O.; Pascual, R.; Borrell, J.; Teixido, J. Cell-integral-diversity criterion: A proposal for minimizing cluster artifact in cell-based selections. *J. Chem. Inf. Model.* **2007**, *47* (5), 1886–1896.
- (12) Brown, R. D.; Hassan, M.; Waldman, M. Combinatorial library design for diversity, cost efficiency, and drug-like character. *J. Mol. Graphics Modell.* **2000**, *18* (4–5), 427–437.
- (13) Jamois, E. A.; Hassan, M.; Waldman, M. Evaluation of reagent-based and product-based strategies in the design of combinatorial library subsets. *J. Chem. Inf. Comput. Sci.* **1999**, *40* (1), 63–70.
- (14) Bayley, M. J.; Willett, P. Binning schemes for partition-based compound selection. *J. Mol. Graphics Modell.* **1999**, *17* (1), 10–18.
- (15) Agrafiotis, D. K. A constant time algorithm for estimating the diversity of large chemical libraries. *J. Chem. Inf. Comput. Sci.* **2001**, *41* (1), 159–167.
- (16) Congreve, M.; Carr, R.; Murray, C.; Jhoti, H. A “rule of three” for fragment-based lead discovery? *Drug Discovery Today* **2003**, *8* (19), 876–877.
- (17) Lipinski, C.; Lombardo, F.; Dominy, B.; Feeney, P. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Adv. Drug. Del. Rev.* **1997**, *23*, 3–25.
- (18) MOE, version 2009-10; Chemical Computing Group; Montreal, Quebec, Canada.
- (19) Prestwick. <http://www.prestwickchemical.com/> (accessed January 15, 2011).
- (20) CMC. <http://www.akosgmbh.de/Symyx/software/databases/cmc-3d.htm> (accessed January 15, 2011).
- (21) Chembridge. <http://www.chembridge.com> (accessed January 15, 2011).
- (22) Pyxis. <https://www.chemonaut.com> (accessed January 15, 2011).
- (23) EPAFHM. U.S. EPA Computational Toxicology Program. [http://www.epa.gov/ncct/dssto/sdf\\_epafhm.html](http://www.epa.gov/ncct/dssto/sdf_epafhm.html). (accessed November 2, 2010).
- (24) Russom, C. L.; Bradbury, S. P.; Broderius, S. J.; Hammermeister, D. E.; Drummond, R. A. Predicting modes of toxic action from chemical structure: Acute toxicity in the fathead minnow (*Pimephales promelas*). *Environ. Toxicol. Chem.* **1997**, *16* (5), 948–967.
- (25) DrugBank. <http://www.drugbank.ca/> (accessed January 15, 2011).
- (26) Owen, J. R.; Nabney, I. T.; Medina-Franco, J. L.; Lopez-Vallejo, F. Visualization of molecular fingerprints. *J. Chem. Inf. Model.* **2011**, *51* (7), 1552–1563.
- (27) Clark, R. D.; Langton, W. J. Balancing representativeness against diversity using optimizable k-dissimilarity and hierarchical clustering. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (6), 1079–1086.
- (28) Bayada, D. M.; Hamersma, H.; van Geerestein, V. J. Molecular diversity and representativity in chemical databases. *J. Chem. Inf. Comput. Sci.* **1999**, *39* (1), 1–10.
- (29) 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* (4), 750–763.
- (30) Reynolds, C. H.; Druker, R.; Pfahler, L. B. Lead discovery using stochastic cluster analysis (SCA): A new method for clustering structurally similar compounds. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (2), 305–312.
- (31) Pascual, R.; Mateu, M.; Gasteiger, J.; Borrell, J. I.; Teixido, J. Design and analysis of a combinatorial library of HEPT analogues: Comparison of selection methodologies and inspection of the actually covered chemical space. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (1), 199–207.
- (32) Akella, L. B.; DeCaprio, D. Cheminformatics approaches to analyze diversity in compound screening libraries. *Curr. Opin. Chem. Biol.* **2011**, *14* (3), 325–330.
- (33) Shanmugasundaram, V.; Maggiora, G. Application of Shannon-like diversity measures to cell-based chemistry spaces. *J. Math. Chem.* **2011**, *49* (2), 342–355.
- (34) Viswanadhan, V. N.; Rajesh, H.; Balaji, V. N. Atom type preferences, structural diversity, and property profiles of known drugs,

leads, and nondrugs: A comparative assessment. *ACS Comb. Sci.* **2011**, *13* (3), 327–336.

(35) Bemis, G. W.; Murcko, M. A. The properties of known drugs. 1. Molecular frameworks. *J. Med. Chem.* **1996**, *39* (15), 2887–2893.

(36) Stein, S. E.; Heller, S. R.; Tchekhovskoi, D. An Open Standard for Chemical Structure Representation: The IUPAC Chemical Identifier. In *2003 International Chemical Information Conference*, Collier, H., Ed.; Infonortics Ltd.: Tetbury, U.K., 2003; pp 131–143.

(37) *InChI*, 1.03, IUPAC, 2010. <http://www.iupac.org/inchi/> (accessed January 15, 2011).

(38) Rogers, D.; Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **2010**, *50* (5), 742–754.

(39) *Pipeline Pilot*, student ed.; Accelrys: San Diego, CA, 2010.

(40) Hert, J.; Willett, P.; Wilton, D. J.; Acklin, P.; Azzaoui, K.; Jacoby, E.; Schuffenhauer, A. Comparison of fingerprint-based methods for virtual screening using multiple bioactive reference structures. *J. Chem. Inf. Comput. Sci.* **2004**, *44* (3), 1177–1185.

(41) Daylight Chemical Information Systems, Inc., PO Box 7737, Laguna Niguel, CA 92677, U.S.A.

(42) Accelrys, Inc., 10188 Telesis Court, Suite 100, San Diego, CA 92121, U.S.A.

(43) Shannon, C. E. A mathematical theory of communication. *Bell Syst. Tech. J.* **1948**, *27*, 379–423.

(44) Dash, M.; Liu, H.; Terano, T.; Chen, A. *Feature Selection for Clustering: Knowledge Discovery and Data Mining. Current Issues and New Applications*; Springer: Berlin/Heidelberg: 2000; Vol. 1805, pp 110–121.

(45) Godden, J. W.; Stahura, F. L.; Bajorath, J. Variability of molecular descriptors in compound databases revealed by Shannon entropy calculations. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (3), 796–800.

(46) Godden, J. W.; Bajorath, J. An information-theoretic approach to descriptor selection for database profiling and QSAR modeling. *QSAR Comb. Sci.* **2003**, *22* (5), 487–497.

(47) MacQueen, J. B. Some Methods for Classification and Analysis of Multivariate Observations. In *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*, University of California Press: Berkeley, 1967; Vol. 1, pp 281–297.

(48) Agrawal, R.; Gehrke, J.; Gunopulos, D.; Raghavan, P. Automatic Subspace clustering of high dimensional data. *Data Min. Knowl. Discovery* **2005**, *11* (1), 5–33.

(49) Rassokhin, D. N.; Agrafiotis, D. K. Kolmogorov–Smirnov statistic and its application in library design. *J. Mol. Graphics Modell.* **2000**, *18* (4–5), 368–382.

(50) von Mises, R. *Mathematical Theory of Probability and Statistics*; Academic Press: New York, 1997.