

Visual Characterization and Diversity Quantification of Chemical Libraries: 1. Creation of Delimited Reference Chemical Subspaces

Vincent Le Guilloux,[†] Lionel Colliandre,[†] Stéphane Bourg,[‡] Guillaume Guénégou,[†] Julie Dubois-Chevalier,^{†,§} and Luc Morin-Allory^{*†}

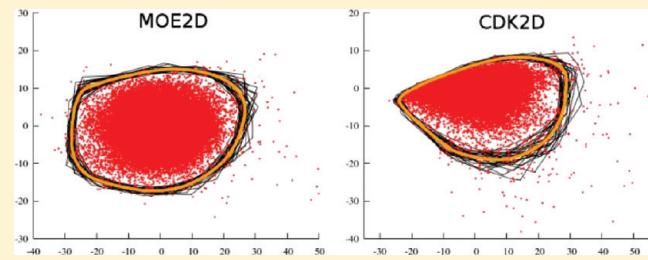
[†]Institut de Chimie Organique et Analytique (ICOA), Université d'Orléans, UMR CNRS 6005 B.P. 6759, rue de Chartres, 45067 Orléans Cedex 2, France

[‡]Fédération de Recherche, Physique et Chimie du Vivant, Université d'Orléans-CNRS; FR 2708, Avenue Charles Sadron, 45071 Orléans Cedex 2, France

[§]Laboratoire d'Informatique Fondamentale d'Orléans (LIFO), Université d'Orléans, Rue de Chartres, 45067 Orléans Cedex 2, France

 Supporting Information

ABSTRACT: High-throughput screening (HTS) is a well-established technology which can test up to several million compounds in a few weeks. Despite these appealing capabilities, available resources and high costs may limit the number of molecules screened, making diversity analysis a method of choice to design and prioritize screening libraries. With a constantly increasing number of molecules available for screening, chemical space has become a key concept for visualizing, analyzing, and comparing chemical libraries. In this first article, we present a new method to build delimited reference chemical subspaces (DRCS). A set of 16 million screening compounds from 73 chemical providers has been gathered, resulting in a database of 6.63 million standardized and unique molecules. These molecules have been used to create three DRCS using three different sets of chemical descriptors. A robust principal component analysis model for each space has been obtained, whereby molecules are projected in a reduced two-dimensional viewable space. The specificity of our approach is that each reduced space has been delimited by a representative contour encompassing a very large proportion of molecules and reflecting its overall shape. The methodology is illustrated by mapping and comparing various chemical libraries. Several tools used in these studies are made freely available, thus enabling any user to compute DRCS matching specific requirements.



INTRODUCTION

“Chemical space—which encompasses all possible small organic molecules, including those present in biological systems—is vast. So vast, in fact, that so far only a tiny fraction of it has been explored” observed Dobson in a reference paper in 2004.¹ He provided the following definition of chemical space: “Chemicals can be characterized by a wide range of ‘descriptors’, such as their molecular mass, lipophilicity, and topological features. Chemical space is a term often used in place of multidimensional descriptor space; it is a region defined by a particular choice of descriptors and the limits placed on them. Chemical space is defined as the total descriptor space that encompasses all the small carbon-based molecules that could in principle be created”.

The notion of chemical space has been used in drug discovery for over 10 years now^{2–8} and is still an active field of research.^{9–15} In a rich review, Medina-Franco et al.¹⁶ quoted another definition which is “the set of all possible molecular structures”.¹⁷ This definition is however too trivial and lacks essential characteristics, namely descriptions and rules. Without descriptions and rules, objects cannot be compared to each other. In fact, both in the mathematical meaning and in the common-sense usage [e.g., our

three-dimensional (3D) universe], spaces are defined by objects having several properties and related by mathematical rules. Thus, a chemical space must be defined by a set (finite or infinite) of compounds and by rules defining their relations (e.g., positions in a multidimensional descriptor space, similarity or dissimilarity metrics, a graphical representation with specific rules, etc.).

While Dobson’s definition fulfills these requirements and gives a complete and appropriate definition of chemical space, it has a limited practical use since the potential number of compounds and descriptors to be calculated is way too large to be mined. Therefore, chemoinformaticians have in practice used a “restricted” number of compounds and, for each one, a restricted number of descriptors to represent chemical spaces.

This number of compounds varies from thousands to hundreds of millions.^{3,18–26} Virtual and Tangible chemical spaces as defined by Hann and Oprea²⁷ are a good illustration of this. Virtual space is defined using all the imaginable molecules. In such a space, the number of molecules is almost infinite; even

Received: February 3, 2011

Published: July 17, 2011

when restricted to small drug-like molecules, it is estimated at over 10^{60} .^{4,28} Despite the increasing power of recent computers, navigating through such a large set of molecules is impossible with tools available nowadays. Tangible space corresponds to compounds that can reliably be made (many approaches have been tested for the creation of such real or virtual chemical libraries, e.g., ref 29); it shows similar limitations. Besides the difficulty of defining what is (and will be) a compound that "can reliably be made", the potential number of molecules involved is still out of reach. Restricting the space to smaller molecules (up to 11 or 13 heavy atoms)^{18,19} is possible and yields more usable sets of compounds, but obviously, these spaces represent only part of the real space.

Similarly, obtaining the "total descriptor space" would require the calculation of all the molecular descriptors available. In their reference book "Molecular Descriptors for Chemoinformatics"³⁰ Todeschini and Consonni reference more than 3000 descriptors, some of which are fingerprints, graphs, or arrays of numerical values. It is therefore possible to characterize each compound by much more than 10 000 values. Working in a space of such a huge dimensionality is practically impossible. One has to reduce the number of descriptors by selecting those which are the most appropriate for a given problem and/or to use data analysis approaches to reduce the apparent dimensionality.^{2,31–43}

To this end, various multivariate methods have already been used to reduce this dimensionality and facilitate the interpretation and navigation through chemical spaces. Nonlinear methods, such as self-organizing map (SOM),^{16,19,44} generative topographic map,⁴⁵ or multifusion similarity approaches^{16,46} have been successfully used to mine and visualize chemical spaces. On the other hand, principal component analysis (PCA) is a gold standard linear method used for several decades now for dimensionality reduction problems. It was used by Oprea et al. to create the ChemGPS chemical space navigation system;³ they suggested the term of "chemography" as the "art of navigating chemical spaces". This original approach was subsequently extended toward pharmacokinetic properties⁴⁷ and natural products,^{48,49} allowing to position chemical entities in different reference spaces. PCA has also been used to analyze active cancer medicinal chemistry compounds, showing that active molecules cover a different chemical space compared to nonactive and hit-like molecules.⁵⁰ Various "high-throughput screening (HTS)-like" screening libraries (Lipinski molecules, natural products, fragments, etc.) were also described and compared by Shelat and Guy⁵¹ using PCA. Reymond et al. recently used both PCA and SOM to represent a virtual chemical space built using more than 26.4 million virtual molecules described by standard physicochemical and autocorrelation descriptors.¹⁹ The same group subsequently used PCA to determine the chemical space repartition of PubChem compounds which were described using molecular quantum numbers.¹⁵ The visualization of binding site-centric chemical space using PCA was also reported by Macchiarulo et al.⁵² Recently, Singh et al.⁵³ compared combinatorial libraries, natural products, and molecules from Molecular Libraries Small Molecule Repository using various molecular descriptions, including multifusion similarity maps and PCA analysis to compare their chemical space coverage. They further highlighted the importance of comparing chemical libraries from different points of view. The wide usage of PCA can be explained by the simplicity of the underlying methodology, the absence of parameters, and the possibility to apply it on a very large data sets without much effort, which is generally more difficult to

perform using more complex methods such as those cited previously.

In the framework of our research dealing with chemical libraries,^{22,54} we intend to render the use of chemical space and the notion itself more accessible. The notion of chemical space has two main applications in drug discovery: the comparison of chemical libraries and the quantification of the overall diversity of a given library.^{16,22,55–57} The role of chemical diversity in drug discovery and more precisely in the HTS of a chemical library has been the focus of debate for many years.^{12,13,58–65} In 2009, the positive influence of diversity on the ratio of hits has been clearly proven by Sukuru et al.⁶⁶ in a retrospective study analyzing 35 recent Novartis HTS campaigns.

In this series of two papers, we propose a new way to create and use viewable chemical subspaces, referred to hereafter as delimited reference chemical subspaces (DRCS). Each DRCS is defined by a set of normalized molecules used to build the space and a set of descriptors encoding each compound. Molecules are then represented in a 2D reduced space using PCA. Furthermore, each DRCS is delimited by a visual contour encompassing a given proportion of molecules and representing the overall shape of the chemical space.

The present paper describes the basic methodology used to build a DRCS, illustrated by its creation for HTS molecules based on a set of more than 16 million available compounds. The second paper will describe the application of DRCS to compare the relative molecular diversity of chemical libraries using DRCS based diversity indices which are independent of the size of the library. The delimitation of chemical spaces makes it possible to use numerous mathematical methods to compute this diversity. Various chemical libraries will be profiled and characterized in terms of diversity.

The practical use of the method, along with DRCS defined in this paper, is made possible for everyone through the release under open-source licenses of several in-house tools used in these studies.⁶⁷

■ DATA PREPARATION

To illustrate the methodology described herein, three DRCS have been defined for a large set of available HTS compounds. This section describes the gathering and preparation of the data used for this purpose.

Collections of Compounds. Collections of compounds from 73 chemical providers have been retrieved, resulting in a set of around 16 million nonunique compounds. The detailed list of providers and the number of molecules for each of them is given in Tables S1 and S2 in Supporting Information. Each provider usually proposes a unique data set of compounds. Sometimes, however, various types of collections are made available: screening compounds, building blocks, fragments, target-focused compounds, etc. In such a case, explicit building blocks libraries were skipped for reactivity reasons. We have chosen to create our own internal database rather than to use publicly available databases because specific standardization and filtering procedures may have been applied in their development (e.g., the ZINC database).²⁶

Data Standardization. Data preparation is a crucial step in chemoinformatics,⁶⁸ especially when chemical descriptors have to be computed. Two identical molecules represented with different ionization states will have different descriptor values, e.g., the number of hydrogen-bond acceptors and donors, the formal charge, etc. Moreover, *in silico* representation of molecules is

Table 1. Number of Compounds Rejected at Each Step of the Overall Protocol

| steps of the standardization protocol | | removed compounds ^a no. | remaining compounds no. |
|---|-----------|------------------------------------|-------------------------|
| initial collection from 73 providers | | — | 16 068 877 |
| file reading | 10 | (0.6 ppm) | 16 068 867 |
| SDF entry without coordinates | 7105 | (442 ppm) | 16 061 762 |
| compounds with isotopic atoms | 4855 | (302 ppm) | 16 056 907 |
| compounds containing “alias” atoms | 1554 | (97 ppm) | 16 055 353 |
| kekulization | 21 | (1.3 ppm) | 16 055 332 |
| compounds with exotic atoms | 3496 | (218 ppm) | 16 051 836 |
| compounds containing atoms with bad valence | 621 | (39 ppm) | 16 051 215 |
| Corina 3D conformation calculation | 7158 | (446 ppm) | 16 044 057 |
| InChI calculation | 5 | (0.3 ppm) | 16 044 052 |
| internal duplicates removal | 945 622 | (58 939 ppm) | 15 098 430 |
| global duplicates removal | 7 954 015 | (473 189 ppm) | 7 144 415 |
| reactive compounds filtering | 515 476 | (72 151 ppm) | 6 628 939 |
| final collection from 73 providers | | — | 6 628 939 |

^a Proportions in ppm are calculated with reference to the compounds remaining at the previous step.

prone to various types of errors, often leading to inconsistent data. A recent publication by Fourches et al.⁶⁹ shed light on various practical issues related to data consistency in chemical libraries and provided several guidelines to minimize the risk of errors and obtain clean data sets. To this end, we have developed and applied an 11 step Pipeline Pilot⁷⁰ protocol to obtain the final standardized 3D structures. The main steps of this protocol are summarized in Table 1, and a complete description is given in Supporting Information. All molecules considered in this paper were standardized using this protocol. We insist on the fact that using the DRCS defined in this article raises the need to strictly apply this protocol (or an equivalent implementation) to any chemical library subsequently mapped in order to obtain interpretable results.

Duplicates Removal. Chemical libraries often contain duplicated molecules, and the standardization process may also create some additional redundancy; some compounds initially present with different counterions will be identical after the standardization procedure. To obtain a unique and nonredundant chemical library, an in-house InChI⁷¹ based script was applied to remove duplicate compounds for each provider’s library. The high value of the mean percentage of duplicates for each provider (5.89%) can be explained by the fact that the libraries of several providers contain up to 50% of duplicates (Table S2 in Supporting Information). This is due to the redundancy existing between multiple SD files for the same provider, e.g., different SD files corresponding to different quantities of the same list of products.

All the standardized libraries were finally merged to obtain a final screening database of 15.10 million unique compounds. As duplicates may also exist between providers, the InChI-based script was again applied on the whole data set. Finally, an SD file containing 7.14 million unique and normalized compounds was obtained. This step allows us to compute the “originality” of each provider among this set of providers, defined as the proportion of molecules that are present only in the provider’s library⁷² (Table S2 in Supporting Information).

Reactive Compounds Removal. An additional filtering step was applied to remove reactive compounds (7.22%), obviously not meant to be included in an HTS chemical space. The sdfilter batch program from MOE⁷³ was used for this purpose. The reactive patterns used are based on those defined by Oprea.⁷⁴

A database of 6.63 million unique, normalized, and nonreactive molecules was finally obtained.

An overview of all preparation steps is given in Table 1.

CALCULATION OF DESCRIPTORS

Molecular description is the central part of any methodology seeking to represent and compare molecules. In this article, we aim to illustrate a new methodology intended to be coupled with the PCA dimensionality reduction method in order to describe and visualize chemical spaces. Therefore, we do not address the issue of selecting the appropriate descriptors for a specific problem. Since the PCA method is suited to continuous variables, we do not consider fingerprint-based descriptors. Although structural fingerprints have shown their usefulness in many applications of chemoinformatics (e.g., retrieving bioactive compounds),⁷⁵ the nature of these descriptors would require the use of an appropriate multivariate method.^{76–78} A similar work, as presented in this article, could be performed in this direction, but we will focus here on the widely accepted PCA method. Three sets of descriptors were subsequently used: two sets of 2D and one set of 3D descriptors (Table 2).

The choice of MOE 2D descriptors is justified by its good coverage of various standard types of 2D descriptors.⁷⁹ The free and open source CDK 2D molecular descriptors^{80,81} were selected so as to allow anyone to use and navigate through the DRCS defined in this study and for comparison purposes. Only a few descriptors were removed for statistical reasons (see Table S3 in Supporting Information for details).

Using a simple classification scheme, Table 2 shows that CDK descriptors are highly biased toward constitutional and topological descriptors, while MOE descriptors have a much more balanced distribution. Although the assignment of a descriptor to a given category may be subject to discussion, the differences are notable enough to support this observation.

In addition to these two pools of descriptors, MOE 3D descriptors were also used to assess eventual differences in chemical libraries mapping between 2D and 3D descriptors. Only a few descriptors were removed for various reasons (see Table S3 in Supporting Information for details).

The sddesc program from MOE was used to calculate 2D and 3D descriptors. Prior to descriptor calculation, partial charges

Table 2. the Three Sets of Descriptors Used to Compute Each Corresponding DRCS^a

| descriptors | initial count | manually removed | null variance | final count | composition for 2D descriptors | | |
|-------------|---------------|------------------|---------------|-------------|--------------------------------|----------------|-------------|
| | | | | | physicochemical | constitutional | topological |
| MOE 2D | 185 | 9 | 2 | 174 | 47.0% | 26.0% | 28.0% |
| CDK 2D | 215 | 1 | 35 | 179 | 3.0% | 50.0% | 47.0% |
| MOE 3D | 148 | 31 | 0 | 117 | — | — | — |

^a Manually removed descriptors correspond to various meta-descriptors, such as Lipinski's "drug-like" flag.

were computed using the MMFF94 force field.⁸² CDK descriptors were calculated using an in-house JAVA program made freely available with the other tools used in this study.

■ DRCS DEFINITION

Based on the previous data, the three chemical spaces are basically represented by all the descriptors, one space for each descriptor set. Their large dimensionality (one dimension per chemical descriptor: 174, 179, and 117 dimensions) makes visual analysis impossible. A reduction to a lower dimensionality (e.g., 2 or 3) is thus required to obtain visual and intuitive representations. This can be easily performed using PCA.

Briefly, PCA is a popular linear projection method used to transform an N -dimensional space into an M -dimensional one ($M < N$) created by M uncorrelated vectors called principal components (PCs). PCs are actually defined by the eigenvectors of the variance–covariance matrix of the input matrix. In our case, the input matrix is centered and scaled to unit variance prior to any calculation. PCs are then calculated from the variance–covariance matrix of the scaled original matrix. When ordered by decreasing eigenvalues, the first components correspond to the largest eigenvalues and explain the largest amount of the total original variance. Finally, each PCA model is defined by a set of PCs and a set of means and standard deviations used to scale and center each descriptor.

Graphical representation of the resulting chemical space can then be obtained by projecting the original molecules (N -dimensions) onto a reduced PCA space using the first two or three components. We are aware that the transformation of such a high dimensional space into 2D or 3D space leads to a substantial loss of information. But priority has been given to intuitive and easy visualization and interpretation of each DRCS, and despite the dimension reduction required by the visual analysis, working on the N -dimensional DRCS remains possible. Practically, two proximal points on this representation can be distant in real space if they differ only by descriptors not represented in the first two components (i.e., having low weights). But, as a large part of the information is available in the first components, two distant points in the representation are expected to be rather distant in the real space. Yet, as some descriptors are not represented in the first components, these two distant points can still be rather closed when considering only other components. But their overall distance in the real space is, at least, their distance in the plane created by the two first components.

To compare the relative positions of chemical libraries within a reference chemical space, both molecules of interest and molecules of reference have to be plotted in the reduced PCA space if no boundaries have been defined. This raises two potential problems: (1) descriptor values of the reference molecules must be available and (2) in our case, plotting several million

molecules for each chemical library of interest cannot be done rapidly and may generate graphs which are difficult to interpret. To address these issues, each chemical space will be delimited by a representative contour that can be used independently of the original data set. This contour defines a zone of the space containing a given proportion of molecules and represents a visual aid enabling chemical libraries to be compared and positioned in an intuitive way. A DRCS is finally defined by the following two components:

- (1) A PCA model. For visual mapping, only the first three components have been retained for each chemical space. Based on these eigenvectors, 2D and 3D space coordinates can easily be assigned to each compound. In the remainder of this paper, only 2D spaces will be considered, as visual interpretation is easier on a 2D medium, but the extension to three dimensions is of course possible.
- (2) A contour defining the chemical subspace boundaries. Each contour is computed using coordinates of molecules in the 2D space. This contour is based on a convex hull calculation.⁸³ The convex hull of a set of points S is the minimal subset of points creating a convex polygon that completely encloses S . It is finally defined as a set of ordered points creating a polygon. An intuitive illustration of the concept may be to imagine an elastic band stretched open to encompass the given objects; when released, the elastic will enclose all the objects, reflecting its convex boundaries. If a 3D space is used, a 3D convex hull is created, generating a potato-shaped volume.

■ DRCS CONSTRUCTION

The most intuitive idea to build the DRCS based on a set of compounds would be to compute both PCA and contour on the whole set. This simple idea however is neither appropriate nor efficient for two main reasons:

- (1) Computing the PCA on the whole data set may not be possible in a reasonable time scale for a large number of molecules. The cost of descriptor calculation (computation time, cost of the licenses) can significantly decrease our capacity to calculate a great number of descriptors. This is typically the case for the MOE 3D descriptors used in this study.
- (2) The convex hull is by nature very sensitive to the data set composition, because it is defined by the objects located in the extreme zones of the space. Exotic molecules are typically mapped in these extreme zones of the reduced PCA space and decrease the representativeness of the boundaries. Moreover, boundaries derived from a single convex hull are rather rough and located at the exact position of molecules defining the convex hull. A single

exotic compound can substantially modify the final shape of the contour.

The procedure used to build each DRCS has therefore been broken down into two steps: PCA model definition and contour calculation. The following sections will describe these two steps and how their respective issues have been tackled. The final DRCS creation procedure will then be outlined.

PCA Model Definition. The limited number of licenses for commercial software decreases our capacity to calculate descriptors which require an extensive calculation time, e.g., MOE 3D descriptors. To overcome this problem, we tested the hypothesis that an equivalent PCA model could be obtained by averaging a limited number of PCA models computed on random subsets containing a limited number of molecules. The number of subsets and the number of molecules per subset were determined in order to obtain PCA models that could be presumed to be very similar to those computed on the entire data set. This approach was validated using the MOE 2D data set, as the calculation of all descriptors for the whole data set took only slightly over three days. The Supporting Information details this procedure and the algorithm used to compute each averaged PCA model. Following this analysis, 30 random subsets of 20 000 molecules were extracted from the whole data set. These subsets were subsequently used to compute average DRCS PCA models for each descriptors pool.

It should be emphasized that no molecules were filtered out during the model computation. The model is actually defined to be representative of the whole data set, and no particular bias is introduced. We will see in the next section that this is not the case for contour calculation.

Contour Calculation. Computing a single convex hull on the whole data set is clearly not a satisfactory way to get a representative contour. One would obtain a rough contour enclosing all molecules but no information describing their actual distribution and density. Our aim, in contrast, is to create a contour representing the part of the chemical space occupied by a very large proportion of molecules, excluding the outliers located in extreme or isolated zones of the space. This is obviously a simplified representation of a chemical space, highlighting the most important characteristics of the compounds' distribution. The obvious advantage is its ease of interpretation. It is conceptually similar to a caricature which exaggerates some characteristics and oversimplifies others, allowing faster recognition of a face than using the original photograph.³⁴

To build such a contour, the selected strategy is to compute several convex hulls on several subsets of compounds, with an additional step removing obvious outliers for each subset prior to each individual contour calculation. It should provide a way to: (1) make the boundaries representative of a given proportion of the original set by excluding isolated compounds and (2) smooth the boundaries by averaging several contours, thus increasing the representativeness. The outlier removal procedure was defined so as to obtain a good trade-off between stability of the final shape and a small proportion of molecules located outside the final consensus envelope. For the sake of simplicity, we used the same 30 subsets as those used to compute each consensus model. For each subset, a given proportion of outliers was first removed; a convex hull was then computed using the remaining molecules.

Removing outliers raises the need to provide a definition of what constitutes an outlier. In the Oprea et al. study³ outliers (referred to as 'satellites') are defined both explicitly (e.g., benzene, cubane) or based on extreme property values. In another context,

Baskin et al.⁸⁵ used one-class support vector machine to surround chemical space regions having a high density of data points, making it possible to identify isolated compounds. Actually, there is no absolute definition of what an outlier is, and there could be many different and equally valid ones, depending on the context. In this study, the removal of outliers is needed to obtain a stable and representative envelope encompassing a given, preferably large proportion of molecules in the reduced space. Typically, outliers in the reduced space are located either in extreme zones of the space or isolated in a sparse region where the neighborhood density is very low. Thus, outliers have been defined based on their distance from the barycenter of the cloud (hereafter named barycenter method) and on their neighborhood density (hereafter named density method).

The following procedure was applied prior to the calculation of each contour for each subset:

- (1) Compute the reduced space coordinates for each molecule of the current subset.
- (2) Remove $b\%$ of molecules having the highest distance from the barycenter of the resulting cloud of points.
- (3) Remove $d\%$ of molecules having the lowest neighborhood density. For a given molecule, the neighborhood density is defined as the number of molecules located at a Euclidean distance D in the 2D reduced space.

The combination of these two different ways of removing outliers has been found to give the best results in terms of stability, representativeness and number of molecules excluded, a number which we try to minimize here. Both methods (and other in-house methods, not discussed here) were first tested separately with different values of parameters and numbers of molecules in the subset.

The consensus contour is finally obtained as the average of the contours of the 30 subsets. From the origin of the referential, which is located close to the center of the cloud of points due to the data centering procedure, 360 successive vectors are drawn making angles from 0° to 359° with the X axis. Each vector intersects the 30 convex hulls in 30 points. The average coordinates of these points are computed and represent the consensus point corresponding to each angle. These 360 average points finally form the consensus hull.

To investigate the differences between the two methods and the advantage of their combined use, both parameters b and d were successively set to 0, while the other varied from 0 to 5% with a step of 0.025%. Each consensus hull was subsequently calculated using the three DRCS models used in this study.

Figure 1 shows the resulting consensus hulls obtained in all cases. The barycenter method used alone leads to consensus hulls showing a clear tendency to shift toward a circle shape, which obviously does not necessarily reflect the actual shape of a given chemical space. The density-based method alone gives much better results, reflecting the shape of each space based on the neighborhood density, but the number of molecules that need to be removed in order to obtain stable envelopes appears to be slightly higher. Figure 1 also demonstrates the importance of removing outliers. Chemical subspace boundaries reach stability (i.e., a stable shape which only becomes smaller when increasing the percentage of outliers removed) only when a given proportion of outliers has been removed. Based on empirical observations and determination of stability (see Supporting Information), b and d were set to 0.05 and 0.20%, respectively. We have verified that this procedure preserves the ratio of points

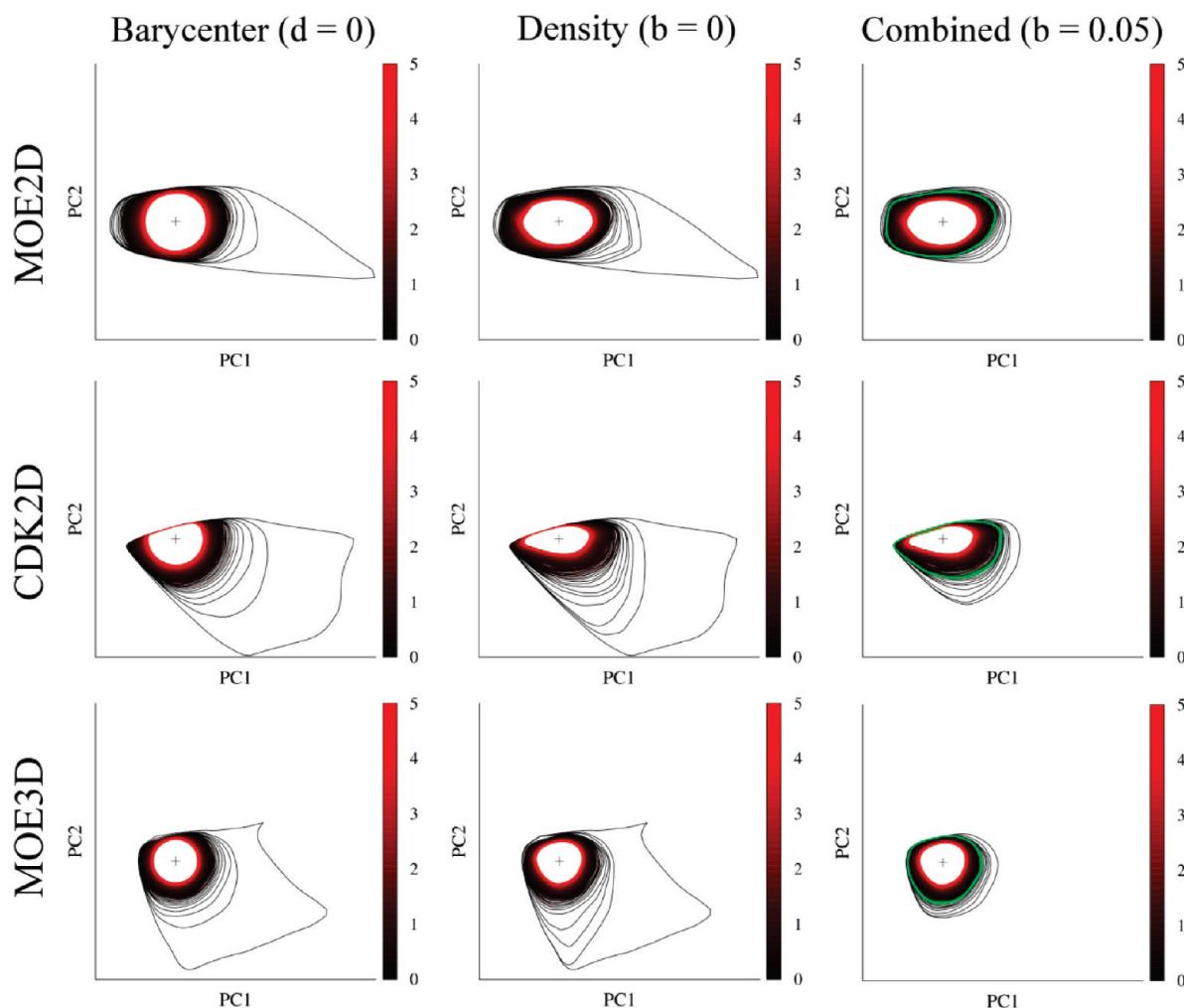


Figure 1. Comparison of three outlier removal methods for the three DRCS. Consensus hulls are computed on the 30 subsets of 20 000 molecules using PC-1 and PC-2 coordinates with $D = 3$ and with a percentage of outliers removed ranging from 0 (black) to 5 (red) with a step of 0.025. Left: only the barycenter-based method is used ($d = 0$, b ranges from 0 to 5). Middle: only the density-based method is used ($b = 0$, d ranges from 0 to 5). Right: both methods are combined, with b fixed at 0.05 and d ranging from 0 to 5. The green consensus contours correspond to the final d value used to compute each DRCS contour. Axis boundaries have been scaled to better highlight the barycenter method.

inside and outside the final contour as defined by the two parameters b and d . Plotting the 600 000 molecules used to compute each DRCS yields 0.33, 0.30, and 0.32% of molecules outside the hull for MOE 2D, CDK 2D, and MOE 3D subspaces, respectively. These ratios are very close to that defined for the creation of each single convex hull (0.25%), the difference being explained by the averaging procedure.

Final Procedure. The following final procedure was used to create each DRCS:

- (1) Pick 30 random subsets of 20 000 molecules in the HTS database.
- (2) Compute each subset PCA.
- (3) Compute the DRCS model by averaging PCA computed on each subset.
- (4) Find and filter out outliers for each subset.
- (5) Compute the convex hulls on the filtered subsets.
- (6) Compute the consensus contour by averaging these convex hulls.
- (7) Define the DRCS as being the PCA model and the consensus contour.

Of course, steps 1–3 could be replaced by a single PCA model computation procedure on the total library. The 30 reference subsets were finally projected in the three DRCS, as shown in Figure 2. Although exotic compounds had been removed prior to calculating each convex hull, this figure also shows that averaging several contours is needed to obtain stability and representativeness. These contours encompass more than 99.6% of the reference compounds and provide useful visual boundaries to describe and compare chemical libraries in this HTS space, as illustrated in the following section.

■ USE OF DRCS FOR VISUALIZATION

Visualization is a rapid and intuitive way to explore and describe the content of a library. Selection of a screening library among all those proposed by the different providers raises the need to rapidly determine the chemical space coverage of each library and to compare these libraries in order to select the one(s) having the required chemical space coverage. Rapid identification of unexplored chemical space zones also greatly facilitates library

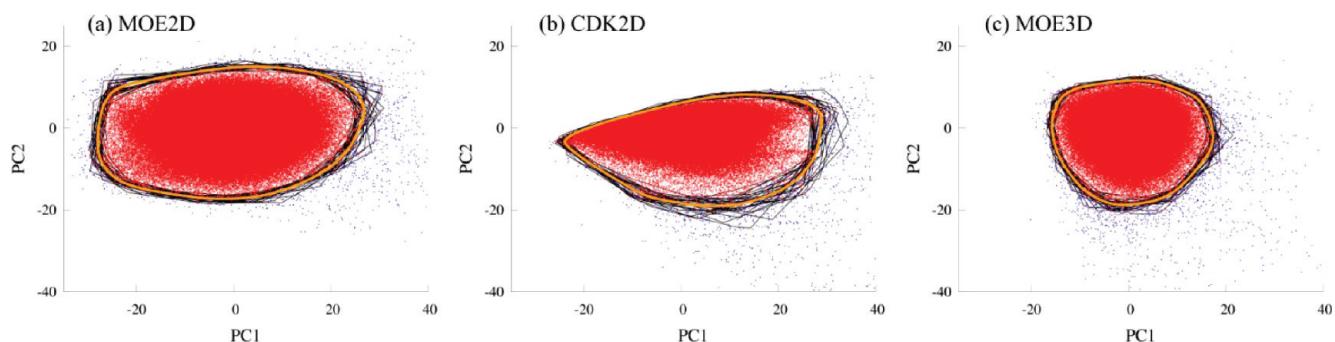


Figure 2. Projection of the 30 subsets (for a total of 600 000 molecules) used to compute each DRCS: (a) MOE 2D, (b) CDK 2D, and (c) MOE 3D spaces. Blue dots represent molecules that have been excluded for at least one contour calculation. Black polygons represent convex hulls of each subset. The orange shape represents the final consensus contour.

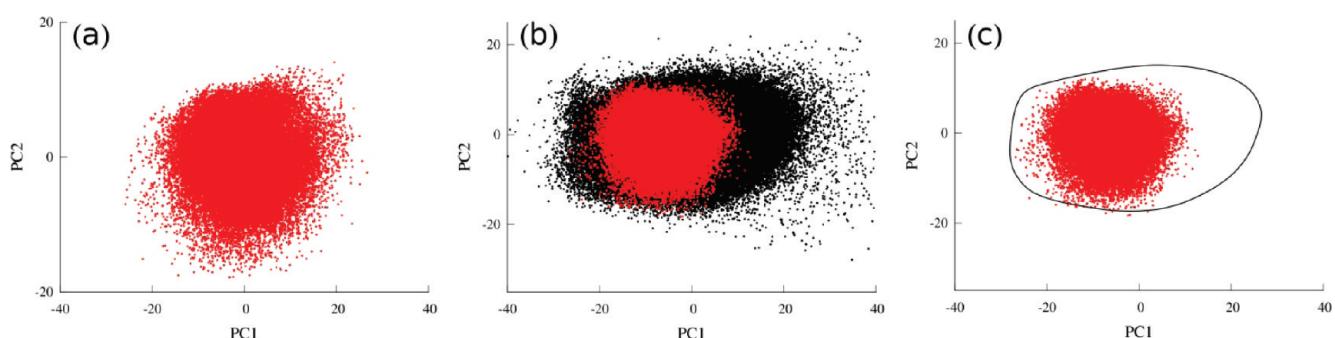


Figure 3. Mapping of a chemical library: (a) using its own MOE 2D PCA model, (b) in the MOE 2D space (red) with the 30 subsets (black) used to compute the MOE 2D model, and (c) in the MOE 2D DRCS plotted with the contour.

enrichment with compounds having complementary properties regarding the chemical space under consideration.

By projecting chemical libraries using the described DRCS methodology, one is able to tackle these issues. Figure 3 illustrates the advantage of using the consensus contour. A chemical library [the Cambridge Diverset Library⁸⁶ (CDL)] has been projected in various ways. In Figure 3a the library has been plotted in its own reduced space. It can be seen that it is very difficult to derive any useful information in terms of chemical space coverage using this figure, as there is no reference to compare it with. When it is plotted in the MOE 2D reduced space against the 30 random subsets used to build the space (Figure 3b), it becomes much easier to compare both libraries and to assess the HTS chemical space coverage of the CDL; the figure clearly shows that the library has a limited coverage on this space. In Figure 3c these 30 subsets have been replaced by the consensus contour. One can see that a very similar interpretation can be derived, without the need to use the descriptor values of the original molecules, which may be either unavailable or difficult to plot (especially when it contains millions of molecules). Moreover, with such a huge data set, it is very difficult in (b) to view the isolated compounds of the library which are located in dense regions of the reference space.

■ DRCS INTERPRETATION AND COMPARISON

Three DRCS have been obtained using the methodology described in the previous sections, one per descriptor set: DRCS-MOE 2D, DRCS-MOE 3D, and DRCS-CDK 2D. The percentages

of explained variances are shown in Table 3 for components 1–5 of each DRCS model.

The percentage of explained variance varies depending on the descriptor set under consideration. It is important to realize that these DRCS cannot be compared based on their respective explained variance as the number and the nature of descriptors differ for each set. In fact, a large value for the total explained variance might rather suggest the presence of very strong inter-correlations between descriptors. Intuitively, one cannot expect to obtain a large proportion of explained variance by reducing the dimensionality, e.g., from 174 to 2, unless strong redundancy exists within the original descriptor space. In other words, large explained variances are not necessarily synonymous with better chemical space representativeness.

For visualization purposes, only the first two PCs were used to obtain viewable 2D DRCS, leading to a cumulated explained variance of 43.51, 31.93, and 35.49% for DRCS-MOE 2D, DRCS-CDK 2D, and DRCS-MOE 3D, respectively. For MOE 2D and MOE 3D DRCS, PC-3 explains around 10% of the variance, resulting in a 3D chemical space accounting for around 50% of the total variance. Although the third component of the CDK 2D DRCS does not explain as much amount of variance, taking into account this component may lead to a substantial gain in information using a 3D visualization device.

The chemical interpretation for each axis can be derived by comparing the relative weights of descriptors for each PC (see Supporting Information for detailed figures showing PC weights) or using simple visual analysis of the compounds' distribution. PC-1 of DRCS-MOE 2D and DRCS-MOE 3D globally describes

Table 3. Mean Explained Variance of the First Five Components of the PCA Consensus Model for Each Set of Descriptors^a

| | MOE 2D (%) | | CDK 2D (%) | | MOE 3D (%) | |
|------|------------|--------------|------------|--------------|------------|--------------|
| | component | cumulative | component | cumulative | component | cumulative |
| PC-1 | 31.31 | 31.31 | 25.84 | 25.84 | 19.01 | 19.01 |
| PC-2 | 12.20 | 43.51 | 6.09 | 31.93 | 16.48 | 35.49 |
| PC-3 | 10.99 | 54.50 | 5.07 | 37.01 | 10.33 | 45.82 |
| PC-4 | 4.74 | 59.24 | 3.89 | 40.89 | 8.16 | 53.98 |
| PC-5 | 3.84 | 63.08 | 3.49 | 44.39 | 5.98 | 59.97 |

^a Corresponding cumulative variances are presented in bold. The total explained variances for each 2D reduced space are in italic.

molecular size through large weights for descriptors related to volume, VDW surface area, or molecular weight/atom counts. PC-1 of CDK 2D is more related to molecular complexity, probably explained by the larger number of topological descriptors present in the CDK library. PC-2 of both DRCS-MOE 2D and DRCS-MOE 3D spaces order compounds by increasing hydrophobicity. In PC-2 of DRCS-CDK 2D, the compounds are classified in a slightly different manner, ranging from aliphatic to aromatic compounds. Furthermore, the shape of the DRCS-CDK 2D (roughly triangular) is very similar to the shape of the space described by van Deursen et al.,¹⁵ where the PCA was computed based on PubChem⁸⁷ compounds. The MQN descriptors used in their study are primarily constitutional and topological descriptors, which explains the similarity observed with our CDK 2D space. Finally, these tendencies show that different chemical descriptors lead to spaces closely related to weight = $f(\text{hydrophobicity})$, as it was previously found by Egan et al.⁸⁸ However, each DRCS has its own specificities that will be highlighted in the next section.

■ LIBRARY MAPPING AND COMPARISON

In this section, the use of DRCS to view and compare libraries will be illustrated, focusing in particular on four different sets of molecules:

- (1) The Comprehensive Medicinal Chemistry (CMC), a database of pharmaceutical compounds.⁸⁹
- (2) The Prestwick Chemical Library containing marketed drugs.⁹⁰
- (3) The Pyxis Discovery Smart Fragment Library,⁹¹ a fragment library based on scaffolds found in existing drugs.
- (4) Estrogen receptor (ER) agonists and antagonists provided by the Directory of Useful Decoys (DUD) benchmark data set.^{92,93}

Prior to any library mapping to a DRCS, each library was standardized using the protocol described in a previous section, and the three sets of descriptors calculated. The results are given in Figure 4.

The 8773 compounds of the CMC database cover a very large proportion of the three DRCS contours. Only one small region is not covered for each of these contours. It corresponds to large and very lipophilic compounds and could be associated to compounds having poor solubility and/or bioavailability. Around 9.5, 8.1, and 10.4% of the compounds are located outside the contour in the MOE 2D, CDK 2D, and MOE 3D DRCS, respectively. On the negative extreme part of PC-1, outside the

hull, very small and volatile compounds can be found (e.g., NO or cyclopropane which are gases used for anesthesia). On the extreme positive part, complex compounds, such as cyclosporin A, a natural macrocycle, are typically present. These classes of compounds are unlikely to be exploited by the providers forming our initial HTS database, and it thus makes sense to find them outside the contour. The coverage of the DRCS contours and the large proportion of compounds outside these contours show that the chemical space of these pharmaceutical compounds has a significantly different distribution from that of commercial HTS compounds.

The Prestwick library, which is commercially available for screening purposes, shows similar space coverage to that of CMC. It contains 1200 small molecules, 100% being marketed drugs. Around 5.8, 7.3, and 8.1% of the compounds are located outside the contour in the MOE 2D, CDK 2D, and MOE 3D DRCS, respectively. As the compounds are similar to those of the CMC (both are pharmaceutical compounds), it seems natural to observe the same distribution. With so few compounds, this library is clearly a good starting point for the creation of an HTS diverse library covering the current marketed drug space.

The projection of these two collections suggests some similarities between the HTS chemical space and the so-called “drug-like space” as defined using the Lipinski⁹⁴ rules for oral bioavailability. The “Lipinski” filter implemented by the MOE software was applied to the 6.6 million data set, yielding around 6.2 million drug-like (93.8%) and 0.4 million nondrug-like (6.2%) molecules. Figure 5 shows the CMC plus Prestwick (CMC-P) set plotted in the MOE 2D subspace against the resulting drug-like molecules (Figure 5a) and nondrug-like molecules (Figure 5b). A rather good correspondence can be observed between the CMC-P set and the Lipinski drug-like molecules space coverage, as shown in Figure 5a. An empty zone can be found in both sets on the top right zone of the DRCS, which corresponds to large lipophilic compounds, showing that the Lipinski filter clearly makes sense in this space area regarding pharmaceuticals and marketed drug compounds. Conversely, Figure 5b shows that this same empty zone is filled by molecules not matching the Lipinski filter. On the other hand, a higher compounds density can be observed, outside the hull, on the bottom right part of the space for the CMC-P set, corresponding roughly to large hydrophilic molecules. The Lipinski drug-like molecules are almost absent in this area (Figure 5a), while nondrug-like molecules are much abundant (Figure 5b). It can also be seen that inside the contour, a significant overlap appears between zones covered by drug- and nondrug-like molecules, showing that pharmaceutical and marketed compounds can also be found in mixed (drug- + nondrug-like) as well as in exclusively nondrug-like zones of the space. All these trends highlight the advantages and limitations of the drug-like filters.

The Pyxis library is entirely located within the DRCS boundaries, in the left part of the subspace. The rule of three⁹⁵ used to define what a fragment is implicitly explains this observation. Both the lower and upper limits applied for the molecular mass and consensus properties defined by this rule explain the tight location of the library. It gives us an idea of where the fragment subspace of the HTS library is. The same simple idea could also be applied to other rules (e.g: lead-like) or subsets of molecules (target-specific), and focused subcontours could also be drawn within the original space using the same consensus methodology described in this paper.

The ER agonists and antagonists are located in two distinct zones of each subspace. Global physicochemical trends separating

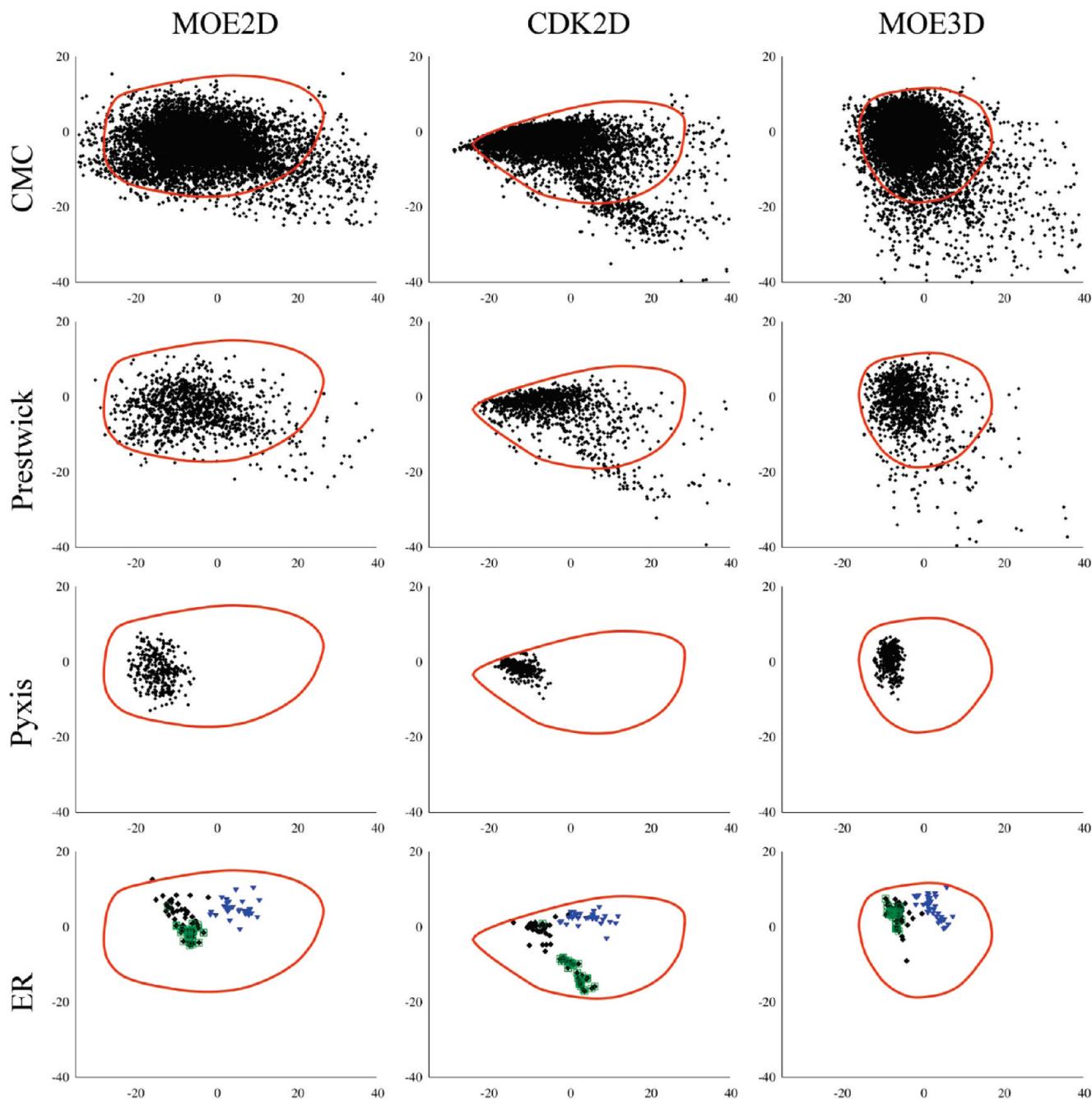


Figure 4. Various libraries mapped in the DRCS described in this article, MOE 2D (left), CDK 2D (middle), and MOE 3D (right). The CMC, Prestwick, Pyxis, and estrogen receptor agonist (black) and antagonists (blue triangles) are shown from top to bottom, respectively. Steroid-like agonists are highlighted in green.

the two classes of compounds can thus be easily captured using the three spaces axis interpretations. For the agonists however, each DRCS yields a slightly different compounds distribution when looking at the compounds' structures. Indeed, the CDK 2D contains more tight and isolated clusters, corresponding to different classes of molecules. In Figure 4, steroid-like agonists have been manually highlighted in green. These compounds are clearly clustered in the CDK 2D subspace and to a lesser extent in the MOE 2D subspace. In the CDK 2D space, they are found clearly isolated from other agonists. In contrast, no obvious specific cluster can be found in the MOE 3D space, and the

steroids are found among the other agonists, suggesting that 3D descriptors for both PC under consideration provide different information compared to CDK 2D and MOE 2D spaces.

Interestingly, the 3D information does not seem to provide any additional information in terms of space coverage compared to 2D descriptors. Although the CDK 2D space seems to give some topological information, the overall conclusion in terms of chemical space coverage remains quite similar despite the different nature of the descriptors used in each DRCS. This shows that, for a rapid overview of the diversity of a chemical library, any of these three DRCS could be used. The simplicity and speed of this

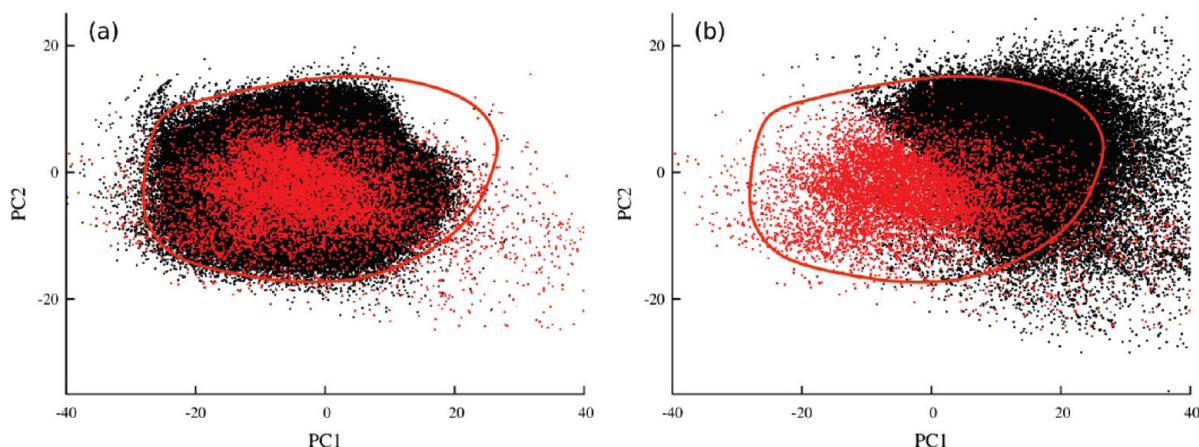


Figure 5. Prestwick and CMC (red dots) plotted in the DRCS-MOE 2D against: (a) a random subset of 2 million “drug-like” molecules filtered using the MOE Lipinski filter and (b) all nondrug-like molecules (same filter) extracted from the entire data set.

methodology obviously ease its systematic application to the analysis of new chemical libraries.

■ CONCLUSION AND PERSPECTIVES

A new methodology has been described to create delimited reference chemical subspaces (DRCS). These so-called DRCS introduced a new way of associating each chemical space with boundaries that reflect their overall shape, encompassing a large majority of compounds in a reduced and viewable 2D subspace. Such a visual aid is of great interest for exploring and comparing chemical libraries by assessing their relative chemical space coverage, leading to more rational library design and selection. Furthermore, we showed that PCA models, very similar to that computed using the entire data set, can be obtained by averaging PCA models computed on several random subsets of molecules. This observation could be useful if the number of molecules to be processed is very large and if computational resources and/or software licenses are limiting factors, though care must be taken with respect to several issues highlighted in Supporting Information of this article. Besides the original purpose, data set cleaning and standardization has emerged as a critical issue. We believe that this step should never be neglected and must be systematically detailed in a modern chemoinformatics study, as results might be dramatically affected by inconsistencies and undetectable errors.

The crucial descriptor selection step has not been addressed here since it is usually context dependent. Obviously, depending on the nature of the target and the amount of knowledge available, focused library design may require careful descriptors selection to create the appropriate DRCS and assess “target subspace” coverage. In contrast, the general purpose HTS chemical spaces described in this study use a wide and diverse range of descriptors, without any particular assumption. They could be useful to assess the global diversity of a chemical library, typically for a blind HTS campaign where only little information is available on the targeted entity. Yet, even with no particular descriptor selection, we showed that different information could still be captured depending on the set of descriptors used. Interestingly, the overall conclusion seems to remain quite similar in terms of chemical space coverage, regardless of which DRCS is under consideration. Rational selection of descriptors

would probably provide more contrasted results for specific purposes.

The conclusions that can be derived using a chemical space based on PCs are obviously limited to the information available in the components under consideration. In the spaces described here, the first two components only explain around 40% of the information available in the original descriptor space and are related to global properties, like size, complexity, and lipophilicity. These properties are useful to determine zones of the space spanned by favorable physicochemical properties and allow one to easily identify bias in chemical libraries. Further insight into chemical space coverage could nevertheless be obtained by computing DRCS using other components. This would provide additional information and complementary molecular representations, and one could obtain a more complete picture of the chemical space, although at the expense of simplicity and ease of interpretation.

The contour associated with each DRCS provides a useful visual aid to compare chemical libraries. By encompassing the densest region and reflecting its overall shape, the most important characteristics of the reduced subspace can be highlighted. A further extension of this approach would be to define subspaces encompassing specific types of molecules, making it possible to locate more focused subspaces.

Besides a simple visual aid, the contour makes it much easier to obtain a more representative grid-based partition of each subspace. This way, grid-based diversity indices can be easily applied and would reflect more accurately the coverage of the most explored regions, thereby making easier to determine an appropriate grid resolution.

From a more practical point of view (e.g., diversity indices calculation or diverse subset extraction), we suggest that the densest and more sparse regions of a chemical space should be treated separately. By creating a clear delimitation for the densest region, the focus is set on the widely explored part of a given chemical space, which a good general purpose/prospective screening library should cover anyway. More exotic compounds that usually represent unexplored chemotypes would require more attention and should thus be treated and sampled differently. The convex delimitation of each space provides a way to perform such analysis; a second paper will explore the creation of specific subspaces as well as various diversity indices adapted to the DRCS methodology.

Finally, the open source availability of all the tools used in this study, both as standalone tools and as screening assistant⁹⁶ platform features, opens the way to easy and interactive DRCS navigation. Free, open, and validated tools/data are still a missing piece in chemoinformatics, especially compared to the bioinformatics field. Fortunately, more and more valuable and very promising initiatives have been recently reported, such as CDK,⁹⁷ Bioclipse,⁹⁸ RDKit,⁹⁹ CDK-Taverna,¹⁰⁰ KNIME¹⁰¹, to cite but a few, and the work described herein is clearly striving to move in this promising direction.

■ ASSOCIATED CONTENT

S Supporting Information. Tables S1 and S2: List of providers and quantitative data about the products; Table S3: List of removed descriptors; Detailed preparation of chemical structures; Description of tools and computational performances; Methodology for the computation of the consensus PCA; Parameters of the contour calculation; The file Models.zip contains all the figures describing each PCA model in detail and all the figures comparing the global and the average MOE 2D model. 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.

■ ACKNOWLEDGMENT

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 the referees and Peter Schmidtke for helpful comments on the manuscript. V.L.G. thanks the “Conseil Général du Loiret” and J.D.C. the “Conseil Régional du Centre” for funding their respective PhDs.

■ REFERENCES

- (1) Dobson, C. M. Chemical space and biology. *Nature* **2004**, *432* (7019), 824–828.
- (2) Willett, P. Computational tools for the analysis of molecular diversity. *Perspect. Drug Discovery Des.* **1997**, *7/8*, 1–11.
- (3) Oprea, T. I.; Gottfries, J. Chemography: the art of navigating in chemical space. *J. Comb. Chem.* **2001**, *3* (2), 157–166.
- (4) Bohacek, R. S.; McMurtin, C.; Guida, W. C. The art and practice of structure-based drug design: a molecular modeling perspective. *Med. Res. Rev.* **1996**, *16* (1), 3–50.
- (5) Lahana, R. How many leads from HTS?. *Drug Discovery Today* **1999**, *4* (10), 447–448.
- (6) Gorse, D.; Rees, A.; Kaczkorek, M.; Lahana, R. Molecular diversity and its analysis. *Drug Discovery Today* **1999**, *4* (6), 257–264.
- (7) Teague, S. J.; Davis, A. M.; Leeson, P. D.; Oprea, T. The Design of Leadlike Combinatorial Libraries. *Angew. Chem., Int. Ed. Engl.* **1999**, *38* (24), 3743–3748.
- (8) Brown, R. D. Descriptors for diversity analysis. *Perspect. Drug Discovery Des.* **1997**, *7/8*, 31–49.
- (9) Chuprina, A.; Lukin, O.; Demoiseaux, R.; Buzko, A.; Shvanyuk, A. Drug- and lead-likeness, target class, and molecular diversity analysis of 7.9 million commercially available organic compounds provided by 29 suppliers. *J. Chem. Inf. Model.* **2010**, *50* (4), 470–479.
- (10) Gu, Q.; Xu, J.; Gu, L. Selecting Diversified Compounds to Build a Tangible Library for Biological and Biochemical Assays. *Molecules* **2010**, *15* (7), 5031–5044.
- (11) Reymond, J.-L.; van Deursen, R.; Blum, L. C.; Ruddigkeit, L. Chemical space as a source for new drugs. *Med. Chem. Commun.* **2010**, *1*, 30–38.
- (12) Schneider, G.; Hartenfeller, M.; Reutlinger, M.; Tanrikulu, Y.; Proschak, E.; Schneider, P. Voyages to the (un)known: adaptive design of bioactive compounds. *Trends Biotechnol.* **2009**, *27* (1), 18–26.
- (13) Schneider, P.; Tanrikulu, Y.; Schneider, G. Self-organizing maps in drug discovery: compound library design, scaffold-hopping, repurposing. *Curr. Med. Chem.* **2009**, *16* (3), 258–266.
- (14) Varnek, A.; Baskin, I. I. Chemoinformatics as a Theoretical Chemistry Discipline. *Mol. Inf.* **2011**, *30*, 20–32.
- (15) van Deursen, R.; Blum, L. C.; Reymond, J. L. A searchable map of PubChem. *J. Chem. Inf. Model.* **2010**, *50* (11), 1924–1934.
- (16) Medina-Franco, J. L.; Karina, M.-M.; Giulianotti, Marc A.; Houghten, Richard A.; Pinilla, Clemencia Visualization of the Chemical Space in Drug Discovery. *Curr. Comput.-Aided Drug Des.* **2008**, *4* (4), 322–333.
- (17) This definition is similar to the one given by Wikipedia, English version; http://en.wikipedia.org/wiki/Chemical_space. Accessed January 15, 2011.
- (18) Blum, L. C.; Reymond, J. L. 970 million drug-like small molecules for virtual screening in the chemical universe database GDB-13. *J. Am. Chem. Soc.* **2009**, *131* (25), 8732–8733.
- (19) Fink, T.; Reymond, J. L. Virtual exploration of the chemical universe up to 11 atoms of C, N, O, F: assembly of 26.4 million structures (110.9 million stereoisomers) and analysis for new ring systems, stereochemistry, physicochemical properties, compound classes, and drug discovery. *J. Chem. Inf. Model.* **2007**, *47* (2), 342–353.
- (20) Oprea, T. I. Chemical space navigation in lead discovery. *Curr. Opin. Chem. Biol.* **2002**, *6* (3), 384–389.
- (21) Oprea, T. I.; Allu, T. K.; Fara, D. C.; Rad, R. F.; Ostropovici, L.; Bologa, C. G. Lead-like, drug-like or “Pub-like”: how different are they?. *J. Comput.-Aided Mol. Des.* **2007**, *21* (1–3), 113–119.
- (22) 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.
- (23) Koch, M. A.; Schuffenhauer, A.; Scheck, M.; Wetzel, S.; Casaulta, M.; Odermatt, A.; Ertl, P.; Waldmann, H. Charting biologically relevant chemical space: a structural classification of natural products (SCONP). *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102* (48), 17272–17277.
- (24) Schuffenhauer, A.; Ertl, P.; Roggo, S.; Wetzel, S.; Koch, M. A.; Waldmann, H. The scaffold tree--visualization of the scaffold universe by hierarchical scaffold classification. *J. Chem. Inf. Model.* **2007**, *47* (1), 47–58.
- (25) Wetzel, S.; Klein, K.; Renner, S.; Rauh, D.; Oprea, T. I.; Mutzel, P.; Waldmann, H. Interactive exploration of chemical space with Scaffold Hunter. *Nat. Chem. Biol.* **2009**, *5* (8), 581–583.
- (26) Irwin, J. J.; Shoichet, B. K. ZINC—a free database of commercially available compounds for virtual screening. *J. Chem. Inf. Model.* **2005**, *45* (1), 177–182.
- (27) Hann, M. M.; Oprea, T. I. Pursuing the leadlikeness concept in pharmaceutical research. *Curr. Opin. Chem. Biol.* **2004**, *8* (3), 255–263.
- (28) Gorse, A. D. Diversity in medicinal chemistry space. *Curr. Top. Med. Chem.* **2006**, *6* (1), 3–18.
- (29) Andrews, K.; Cramer, R. D. Toward General Methods of Targeted Library Design: Topomer Shape Similarity Searching with Diverse Structures as Queries. *J. Med. Chem.* **2000**, *43* (9), 1723–1740.
- (30) Todeschini, R.; Consonni, V., Molecular Descriptors for Chemoinformatics. In *Methods and Principles in Medicinal Chemistry*; Mannhold, R., Kubinyi, H., Folkers, G., Eds. Wiley-VCH: Weinheim, Germany, 2009; Vol. 41.
- (31) Dunbar, J. B., Jr Cluster-based selection. *Perspect. Drug Discovery Des.* **1997**, *7/8*, 51–63.
- (32) Godden, J. W.; Bajorath, J. Chemical descriptors with distinct levels of information content and varying sensitivity to differences

- between selected compound databases identified by SE-DSE analysis. *J. Chem. Inf. Comput. Sci.* **2002**, *42* (1), 87–93.
- (33) Godden, J. W.; Bajorath, J. Differential Shannon Entropy as a sensitive measure of differences in database variability of molecular descriptors. *J. Chem. Inf. Comput. Sci.* **2001**, *41* (4), 1060–1066.
- (34) 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.
- (35) Lajiness, M. S. Dissimilarity-based compound selection techniques. *Perspect. Drug Discovery Des.* **1997**, *7/8*, 65–84.
- (36) Landon, M. R.; Schaus, S. E. JEDA: Joint entropy diversity analysis. An information-theoretic method for choosing diverse and representative subsets from combinatorial libraries. *Mol. Diversity* **2006**, *10* (3), 333–339.
- (37) Mason, J. S.; Pickett, S. D. Partition-based selection. *Perspect. Drug Discovery Des.* **1997**, *7/8*, 85–114.
- (38) Stahura, F. L.; Godden, J. W.; Bajorath, J. Differential Shannon entropy analysis identifies molecular property descriptors that predict aqueous solubility of synthetic compounds with high accuracy in binary QSAR calculations. *J. Chem. Inf. Comput. Sci.* **2002**, *42* (3), 550–558.
- (39) Vogt, I.; Bajorath, J. Design and exploration of target-selective chemical space representations. *J. Chem. Inf. Model.* **2008**, *48* (7), 1389–1395.
- (40) Willett, P.; Barnard, J. M.; Downs, G. M. Chemical Similarity Searching. *J. Chem. Inf. Comput. Sci.* **1998**, *38* (6), 983–996.
- (41) Xue, L.; Bajorath, J. Molecular descriptors in chemoinformatics, computational combinatorial chemistry, and virtual screening. *Comb. Chem. High Throughput Screen.* **2000**, *3* (5), 363–372.
- (42) Xue, L.; Bajorath, J. Molecular descriptors for effective classification of biologically active compounds based on principal component analysis identified by a genetic algorithm. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (3), 801–809.
- (43) Xue, L.; Godden, J. W.; Bajorath, J. Evaluation of descriptors and mini-fingerprints for the identification of molecules with similar activity. *J. Chem. Inf. Comput. Sci.* **2000**, *40* (5), 1227–1234.
- (44) Sadowski, J.; Wagener, M.; Gasteiger, J. Assessing Similarity and Diversity of Combinatorial Libraries by Spatial Autocorrelation Functions and Neural Networks. *Angew. Chem., Int. Ed. Engl.* **1996**, *34* (23–24), 2674–2677.
- (45) Maniyar, D. M.; Nabney, I. T.; Williams, B. S.; Sewing, A. Data Visualization during the Early Stages of Drug Discovery. *J. Chem. Inf. Model.* **2006**, *46* (4), 1806–1818.
- (46) Medina-Franco, J. L.; Maggiore, G. M.; Giulianotti, M. A.; Pinilla, C.; Houghten, R. A. A Similarity-based Data-fusion Approach to the Visual Characterization and Comparison of Compound Databases. *Chem. Biol. Drug Des.* **2007**, *70* (5), 393–412.
- (47) Oprea, T. I.; Zamora, I.; Ungell, A.-L. Pharmacokinetically Based Mapping Device for Chemical Space Navigation. *J. Comb. Chem.* **2002**, *4* (4), 258–266.
- (48) Larsson, J.; Gottfries, J.; Muresan, S.; Backlund, A. ChemGPS-NP: Tuned for Navigation in Biologically Relevant Chemical Space. *J. Nat. Prod.* **2007**, *70* (5), 789–794.
- (49) Rosén, J.; Lövgren, A.; Kogej, T.; Muresan, S.; Gottfries, J.; Backlund, A. ChemGPS-NPWeb: chemical space navigation online. *J. Comput.-Aided Mol. Des.* **2009**, *23* (4), 253–259.
- (50) Lloyd, D. G.; Golfs, G.; Knox, A. J. S.; Fayne, D.; J., M. M.; Oprea, T. I. Oncology exploration: charting cancer medicinal chemistry space. *Drug Discovery Today* **2006**, *11* (3/4), 149–159.
- (51) Shelat, A. A.; Guy, R. K. The interdependence between screening methods and screening libraries. *Curr. Opin. Chem. Biol.* **2007**, *11* (3), 244–251.
- (52) Macchiarulo, A.; Pellicciari, R. Exploring the other side of biologically relevant chemical space: Insights into carboxylic, sulfonic and phosphonic acid bioisosteric relationships. *J. Mol. Graphics Modell.* **2007**, *26* (4), 728–739.
- (53) Singh, N.; Guha, R.; Giulianotti, M. A.; Pinilla, C.; Houghten, R. A.; Medina-Franco, J. L. Chemoinformatic Analysis of Combinatorial Libraries, Drugs, Natural Products, and Molecular Libraries Small Molecule Repository. *J. Chem. Inf. Model.* **2009**, *49* (4), 1010–1024.
- (54) 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.
- (55) Baurin, N.; Baker, R.; Richardson, C.; Chen, I.; Foloppe, N.; Potter, A.; Jordan, A.; Roughley, S.; Parratt, M.; Greaney, P.; Morley, D.; Hubbard, R. E. Drug-like annotation and duplicate analysis of a 23-supplier chemical database totalling 2.7 million compounds. *J. Chem. Inf. Comput. Sci.* **2004**, *44* (2), 643–651.
- (56) 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.
- (57) Verheij, H. J. Leadlikeness and structural diversity of synthetic screening libraries. *Mol. Diversity* **2006**, *10* (3), 377–388.
- (58) Potter, T.; Matter, H. Random or rational design? Evaluation of diverse compound subsets from chemical structure databases. *J. Med. Chem.* **1998**, *41* (4), 478–488.
- (59) Schuffenhauer, A.; Popov, M.; Schopfer, U.; Acklin, P.; Stanek, J.; Jacoby, E. Molecular diversity management strategies for building and enhancement of diverse and focused lead discovery compound screening collections. *Comb. Chem. High Throughput Screen.* **2004**, *7* (8), 771–781.
- (60) Jacoby, E.; Schuffenhauer, A.; Popov, M.; Azaoui, K.; Havill, B.; Schopfer, U.; Engeloch, C.; Stanek, J.; Acklin, P.; Rigollier, P.; Stoll, F.; Koch, G.; Meier, P.; Orain, D.; Giger, R.; Hinrichs, J.; Malagu, K.; Zimmermann, J.; Roth, H. J. Key aspects of the Novartis compound collection enhancement project for the compilation of a comprehensive chemogenomics drug discovery screening collection. *Curr. Top. Med. Chem.* **2005**, *5* (4), 397–411.
- (61) Crisman, T. J.; Jenkins, J. L.; Parker, C. N.; Hill, W. A.; Bender, A.; Deng, Z.; Nettles, J. H.; Davies, J. W.; Glick, M. "Plate cherry picking": a novel semi-sequential screening paradigm for cheaper, faster, information-rich compound selection. *J. Biomol. Screen.* **2007**, *12* (3), 320–327.
- (62) Valler, M. J.; Green, D. Diversity screening versus focussed screening in drug discovery. *Drug Discovery Today* **2000**, *5* (7), 286–293.
- (63) Shoichet, B. K. Virtual screening of chemical libraries. *Nature* **2004**, *432* (7019), 862–865.
- (64) Hamprecht, F. A.; Thiel, W.; van Gunsteren, W. F. Chemical library subset selection algorithms: a unified derivation using spatial statistics. *J. Chem. Inf. Comput. Sci.* **2002**, *42* (2), 414–428.
- (65) Willett, P. Chemoinformatics - similarity and diversity in chemical libraries. *Curr. Opin. Biotechnol.* **2000**, *11* (1), 85–88.
- (66) 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.
- (67) DRCS Tools; ICOA-CNRS: Orleans, France; <http://www.univ-orleans.fr/icoa/DRCS/>. Accessed January 15, 2011.
- (68) Bologa, C. G.; Olah, M. M.; Oprea, T. I. Chemical database preparation for compound acquisition or virtual screening. *Methods Mol. Biol.* **2006**, *316*, 375–388.
- (69) Fourches, D.; Muratov, E.; Tropsha, A. Trust, but verify: on the importance of chemical structure curation in chemoinformatics and QSAR modeling research. *J. Chem. Inf. Model.* **2010**, *50* (7), 1189–1204.
- (70) Pipeline Pilot, student ed.; Accelrys: San Diego, CA, 2010.
- (71) InChI, 1.03; IUPAC: Research Triangle Park, NC, 2010; <http://www.iupac.org/inchi/>. Accessed January 15, 2011.
- (72) Originality is a relative concept which not only depends both on the list of products of the provider but also on the list of compounds of the other providers. If the compounds of a provider are remarketed by another one of the list, then the originality connected to these compounds is null for both providers. It is why the comparison with the results obtained in previous work is without object. It is also why this value cannot be used as a commercial argument (either positive or negative).
- (73) MOE, version 2009–10; Chemical Computing Group: Montreal, Quebec, Canada, 2009.
- (74) Oprea, T. I. Property distribution of drug-related chemical databases*. *J. Comput.-Aided Mol. Des.* **2000**, *14* (3), 251–264.

- (75) 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.
- (76) Lee, S.; Huang, J. Z.; Hu, J. Sparse logistic principal components analysis for binary data. *Ann. Appl. Stat.* **2010**, *4* (3), 1579–1601.
- (77) Nikolaj, T. In *What is the Dimension of Your Binary Data?*: 6th IEEE International Conference on Data Mining (ICDM'06), Hong Kong, China, December 18–22, 2006; Taneli, M., SAristides, G., Heikki, M., Eds.; IEEE Computer Society: Los Alamitos, CA, 2006; pp 603–612.
- (78) 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.
- (79) Labute, P. A widely applicable set of descriptors. *J. Mol. Graphics Modell.* **2000**, *18* (4–5), 464–477.
- (80) Steinbeck, C.; Han, Y.; Kuhn, S.; Horlacher, O.; Luttmann, E.; Willighagen, E. The Chemistry Development Kit (CDK): an open-source Java library for Chemo- and Bioinformatics. *J. Chem. Inf. Comput. Sci.* **2003**, *43* (2), 493–500.
- (81) Steinbeck, C.; Hoppe, C.; Kuhn, S.; Floris, M.; Guha, R.; Willighagen, E. L. Recent developments of the chemistry development kit (CDK) - an open-source java library for chemo- and bioinformatics. *Curr. Pharm. Des.* **2006**, *12* (17), 2111–2120.
- (82) Halgren, T. A. Merck molecular force field. I. Basis, form, scope, parameterization, and performance of MMFF94. *J. Comput. Chem.* **1996**, *17* (5–6), 490–519.
- (83) Graham, R. L. An Efficient Algorithm for Determining the Convex Hull of a Finite Planar Set. *Inf. Process. Lett.* **1972**, *1*, 132–133.
- (84) Leopold, D. A.; Bondar, I. V.; Giese, M. A. Norm-based face encoding by single neurons in the monkey inferotemporal cortex. *Nature* **2006**, *442* (7102), 572–575.
- (85) Baskin, I. I.; Kireeva, N.; Varnek, A. The One-Class Classification Approach to Data Description and to Models Applicability Domain. *Mol. Inf.* **2010**, *29* (8–9), 581–587.
- (86) Chembridge; ChemBridge: San Diego, CA; <http://www.chembridge.com>. Accessed January 15, 2011.
- (87) Wang, Y.; Xiao, J.; Suzek, T. O.; Zhang, J.; Wang, J.; Bryant, S. H. PubChem: a public information system for analyzing bioactivities of small molecules. *Nucleic Acids Res.* **2009**, *37* (Web issue), W623–633.
- (88) Egan, W. J.; Merz, K. M.; Baldwin, J. J. Prediction of Drug Absorption Using Multivariate Statistics. *J. Med. Chem.* **2000**, *43* (21), 3867–3877.
- (89) CMC; AKos Consulting and Solutions GmbH: Steinen, Duetschland; <http://www.akosgmbh.de/Symyx/software/databases/cmc-3d.htm>. Accessed January 15, 2011.
- (90) Prestwick; Prestwick Chemical: Illkirch, France; <http://www.prestwickchemical.com/>. Accessed January 15, 2011.
- (91) Pyxis; Chemonaut: Delft, The Netherlands; <https://www.chemonaut.com>. Accessed January 15, 2011.
- (92) Huang, N.; Shoichet, B. K.; Irwin, J. J. Benchmarking sets for molecular docking. *J. Med. Chem.* **2006**, *49* (23), 6789–6801.
- (93) DUD, A Directory of Useful Decoys; University of California, San Francisco: San Francisco, CA; <http://dud.docking.org/>. Accessed January 15, 2011.
- (94) 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.
- (95) 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.
- (96) ScreeningAssistant 2; ICOA-CNRS: Orleans, France; <http://www.univ-orleans.fr/icoa/modelisation/index.php?h=2>. Accessed June 15, 2011.
- (97) CDK; Geeknet, Inc.: Fairfax, VA; <http://sourceforge.net/projects/cdk/>. Accessed January 15, 2011.
- (98) Spjuth, O.; Helmus, T.; Willighagen, E. L.; Kuhn, S.; Eklund, M.; Wagener, J.; Murray-Rust, P.; Steinbeck, C.; Wikberg, J. E. Bioclipse: an open source workbench for chemo- and bioinformatics. *BMC Bioinformatics* **2007**, *8*, 59.
- (99) RDKit; Geeknet, Inc.: Fairfax, VA; <http://rdkit.org/>. Accessed January 15, 2011.
- (100) Kuhn, T.; Willighagen, E.; Zielesny, A.; Steinbeck, C. CDK-Taverna: an open workflow environment for cheminformatics. *BMC Bioinformatics* **2010**, *11*, 159.
- (101) KNIME; KNIME.com GmbH: Zurich, Switzerland; <http://www.knime.org/>. Accessed Jasuary 15, 2011.