

Three-Dimensional Pharmacophore Methods in Drug Discovery

Andrew R. Leach,^{*,||} Valerie J. Gillet,[§] Richard A. Lewis,[‡] and Robin Taylor[†]

^{||}Computational and Structural Chemistry, GlaxoSmithKline Research & Development, Gunnels Wood Road, Stevenage, Hertfordshire SG1 2NY, U.K., [§]Department of Information Studies, University of Sheffield, Regent Court, 211 Portobello Street, Sheffield S1 4DP, U.K., [‡]Novartis Institutes for BioMedical Research, CH-4002 Basel, Switzerland, and [†]Taylor Cheminformatics Software, 54 Sheffield Avenue, Rickmansworth, Herts WD3 1NL, U.K.

Received June 5, 2009

Introduction

The term pharmacophore has been used in medicinal chemistry for many years. The official 1998 IUPAC definition¹ is as follows:

“A pharmacophore is the ensemble of steric and electronic features that is necessary to ensure the optimal supramolecular interactions with a specific biological target structure and to trigger (or to block) its biological response.”

The term is unfortunately also used incorrectly in medicinal chemistry, as implied in the notes that accompany the above definition:¹

“A pharmacophore does not represent a real molecule or a real association of functional groups, but a purely abstract concept that accounts for the common molecular interaction capacities of a group of compounds towards their target structure. The pharmacophore can be considered as the largest common denominator shared by a set of active molecules. This definition discards a misuse often found in the medicinal chemistry literature which consists of naming as pharmacophores simple chemical functionalities such as guanidines, sulfonamides or dihydroimidazoles (formerly imidazolines), or typical structural skeletons such as flavones, phenothiazines, prostaglandins or steroids.”

Central to the pharmacophore concept is the notion that the molecular recognition of a biological target shared by a group of compounds can be ascribed to a (small) set of common features that interact with a set of complementary sites on the biological target. In pharmacophore research quite general features such as hydrogen-bond donors, hydrogen-bond acceptors, positively and negatively charged groups, and hydrophobic regions are typically used. As such, there is a close link between the pharmacophore concept and the widely used principles of bioisosterism.

The other key component of contemporary pharmacophore research is the incorporation of information about the three-dimensional nature of molecular interactions. The focus of this perspective is on 3D pharmacophore methods in which the spatial relationship between the pharmacophore features is also specified. Appreciation of the importance of molecular conformation grew during the 1970s and 1980s, spurred by the increased availability of relevant experimental data (e.g., via the Cambridge Structural Database²) and the development of

new computational methods for the calculation and visualization of stable molecular conformations. Marshall and colleagues in particular recognized the challenge of identifying the binding conformation from among the potentially very large number of accessible structures through their development of the “active analogue” approach.^{3,4}

The early pharmacophore studies (note that the origins of the term pharmacophore have been investigated by Van Drie⁵ who credits Keir with the introduction of the concept in a series of publications in the late 1960s and early 1970s^{6–8}) were typically performed using small numbers of compounds in which the key features required for activity could be readily identified by visual inspection of the molecular structure(s). Such molecules also tended to possess limited conformational flexibility. It was thus generally possible to deduce the pharmacophore manually, possibly assisted through the use of interactive molecular graphics visualization programs. The diversity and complexity of molecular structures that characterize drug discovery today have led to the development of sophisticated computer algorithms for the elucidation, manipulation, and use of pharmacophore models. Nevertheless, the basic concept of a pharmacophore as a simple geometric representation of the key molecular interactions remains unchanged. Pharmacophores have found widespread use in medicinal chemistry for hit and lead identification and during the subsequent lead optimization. The simplicity of the pharmacophore representation does also inevitably mean that it cannot explain everything; understanding the limitations of the concept is essential to a successful application.

In the first part of this Perspective we provide an overview of the computational 3D pharmacophore methods most commonly used in drug discovery. We focus in particular on the key problem of pharmacophore elucidation: the identification from a set of active molecules and their biological activities of the key common features and their relative orientations (also called pharmacophore mapping). From this core concept has arisen a number of other widely used techniques that we also review in outline. These include the use of 3D pharmacophore databases in focused screening, pharmacophore fingerprints, the derivation and application of pharmacophore information from protein structures, the use of molecular fields as an alternative form of pharmacophore representation, and 3D QSAR methods. In the second part we reflect on the challenges in this field and provide some thoughts on future directions. It is not our aim to provide a detailed survey of recent developments or a comprehensive historical review, which can be found

*To whom correspondence should be addressed. Phone: +44 1438 763383. Fax: +44 1438 763352. E-mail andrew.r.leach@gsk.com.

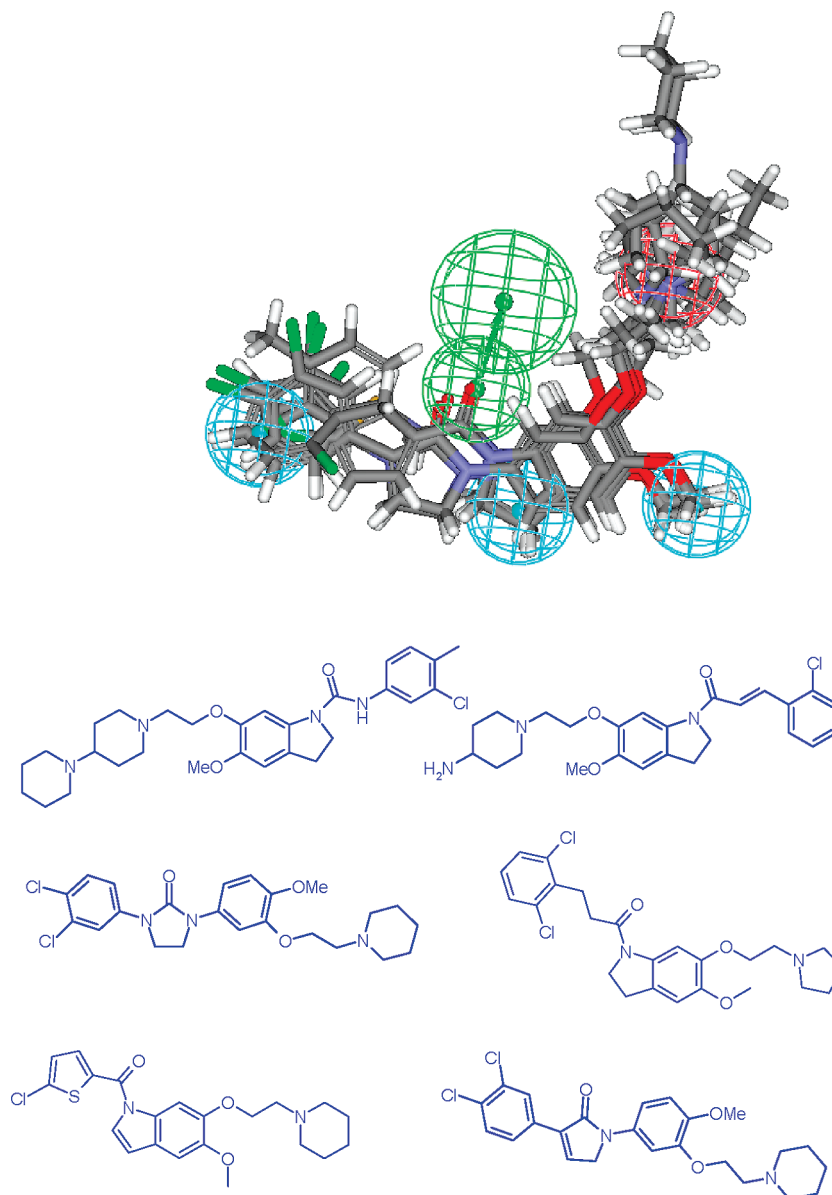


Figure 1. An illustration of the basic pharmacophore concept. The 5HT_{2c} antagonists illustrated can be aligned to generate the 3D pharmacophore model shown with the following color coding: H-bond acceptor (green), positive ionizable (red), hydrophobic (cyan).

elsewhere.^{9–14} Nor is it our intention to provide more than outline information on the various algorithms and software approaches. Our primary goal is to provide the interested reader with a fair assessment of the current state-of-the-art in this field with an emphasis on the methods and software that in our experience are most widely used in real drug discovery projects. We will cover both the practical aspects and some of the inherent limitations of such methods.

Pharmacophore Elucidation

Pharmacophore elucidation is a molecular alignment problem, the aim being to superimpose a set of active ligands, all of which bind to the same protein of unknown 3D structure, so that the features they have in common become evident. Variants of the problem in which the set of ligands includes inactives as well as actives, and/or the 3D structure of the binding site is known, are considered later. To illustrate the basic concept, we show in Figure 1 a series of structures together with the corresponding 3D pharmacophore.¹⁵

A number of programs for pharmacophore elucidation are widely used largely because of their availability in commercial software packages. These include CATALYST,¹⁶ GALAHAD,¹⁷ GASP,¹⁸ the pharmacophore module of MOE,¹⁹ and PHASE.²⁰ However, many other pharmacophore-elucidation algorithms have been described, several of which have been published in the past 3 years or so.^{21–29} This demonstrates that the problem is not considered solved. Its difficulty stems from two sources. First, the calculation is compute-intensive because it involves searching the conformational space of flexible molecules and finding matching subsets of feature points from the different molecules. Each of these on its own is challenging; the combination is formidable. Second, identifying the correct solution can be extremely difficult (we return to this point later). In addition, for an algorithm to be useful, it must be able to align nontrivial sets of ligands. This might go without saying were it not that pharmacophore elucidation programs are occasionally “validated” by their authors on series of ligands that any competent modeler could mentally overlay in a few seconds.

All pharmacophore elucidation algorithms must include methods for (a) representing the ligands (i.e., placing points on or around the molecules to represent the various pharmacophoric features they contain), (b) searching for candidate alignments, (c) scoring those alignments. These aspects are considered separately.

Feature Representation. This involves two steps. First, the molecule must be partitioned into a set of features, each capable of a particular type of intermolecular interaction with the protein. Second, each feature must be represented by one or more points that can be used for least-squares fitting of one molecule onto another. Greene et al. (henceforth “GKSST”^a) wrote an influential paper describing how these steps were accomplished in the CATALYST program.³⁰ Many of the popular pharmacophore elucidation methods in use today employ a similar approach to feature representation,³¹ and so we use their recommendations as a basis for discussion.

Hydrogen-Bond Donor Features. GKSST define hydroxyl groups, nitrogen-bound hydrogens, thiols, and acetylenic CH groups as donors. The last two are controversial, and many workers exclude them as being too weak. Conversely, some argue that other types of CH, in addition to acetylenes, should be considered donors; e.g., CH groups in nitrogen heterocycles are implicated in the binding of some kinase ligands.³² Ionization complicates the picture. GKSST exclude the hydroxyl of $-\text{CO}_2\text{H}$ as a donor because the acid is likely to be ionized at physiological pH, whereas basic amines (e.g., RCH_2NMe_2) are included because it is reasoned that they will be protonated. Uncertainties in ionization and tautomeric states bedevil pharmacophore analysis^{33,34} because they can reverse the nature of a feature (i.e., change an acceptor to a donor or vice versa). One approach is to define explicitly all the variants in which hydrogen-bonding groups might be presented to the pharmacophore elucidation program (e.g., the use of a Daylight SMARTS string³⁵ that matches both RCH_2NMe_2 and $\text{RCH}_2\text{NHMe}_2^+$ could be used to ensure that both will be counted as donors). Many programs provide algorithms for “cleansing” ligand structures in order to achieve compatibility with the subsequent stages of pharmacophore generation and application. Alternatively, the assumption can be made that molecules will be presented to the pharmacophore elucidation program in their appropriate protonation states, thus delegating the entire problem to specialist ionization-state- and tautomer-prediction programs. The substructural definitions of hydrogen-bonding features then become much simpler.

Hydrogen-Bond Acceptor Features. GKSST define any N, O, or S with at least one available (i.e., nondelocalized) lone pair as an acceptor. There is a strong argument for excluding some types of oxygen atoms, e.g., those in rings such as furan and oxazole, which theoretical and crystallographic evidence indicates are very weak acceptors.³⁶

Donor/Acceptor Feature-Point Positioning. It is not sufficient merely to identify which features are present in each

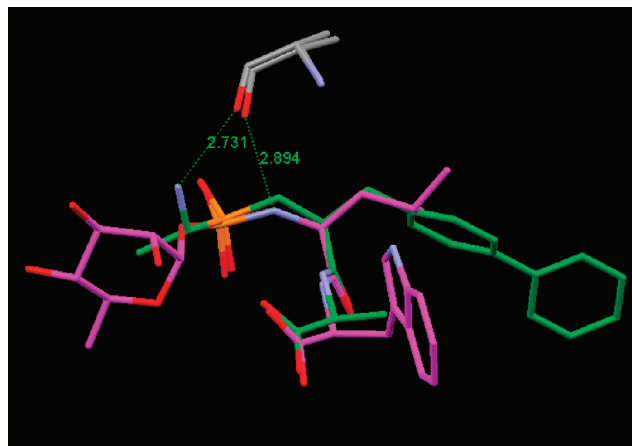


Figure 2. Two neprilysin ligands from PDB structures 1r1h (ligand carbon atoms colored green) and 1dmt (purple carbons). The $-\text{NH}_3^+$ nitrogen of the former and the amidophosphate nitrogen of the latter both donate hydrogen bonds to the backbone carbonyl oxygen of Ala543, but they interact with different lone pairs and so occupy very different positions in the binding site.

ligand molecule; it is also necessary to determine the locations of the associated feature points that will be overlaid in the resulting pharmacophore. Most programs in common use today associate donor and acceptor features not only with the corresponding ligand atoms but also with the presumed location(s) of the complementary protein atom involved in the interaction. For example, GKSST place feature points 3 Å from the heavy atom in the direction(s) of the donor hydrogen or acceptor lone pair(s), provided they are surface exposed. This strategy has the advantage that it permits the overlay of two ligands that form hydrogen bonds to the same protein atom but from different locations and directions (e.g., two ligand donors donating to different sp^2 lone pairs of a protein carbonyl oxygen; see Figure 2). However, it implicitly makes assumptions about the directionality of hydrogen bonding, and if these assumptions are untrue, the points may be placed at inappropriate positions and the probability of generating a correct overlay reduced. One way of relaxing the directionality assumptions is to place points not only along lone-pair directions but also at intermediate positions (e.g., along the extension of the $\text{C}=\text{O}$ bond of a carbonyl group). Alternatively, points can be placed on donor-hydrogen and “lone-pair positions” (typically chosen to be 1 Å from the acceptor atom), therefore reducing the leverage the points have during least-squares fitting, or on the donor or acceptor heavy-atom positions. One advantage of the latter is that it reduces the total number of feature points, which simplifies the search space.

GKSST account for the possibility of bond rotation in, for example, an $\text{R}-\text{OH}$ group by placing points at intervals on the circle swept out as the $\text{R}-\text{O}$ bond rotates. This stratagem is unnecessary if hydroxyl groups are explicitly rotated during pharmacophore elucidation. Similar considerations apply to the lone-pair or virtual points associated with, for example, the oxygen atoms of an ionized phosphate.

Positive and Negative Features. GKSST define atoms bearing formal charges (taking into account probable ionization states *in vivo*) as positive or negative features unless they are bonded to an atom with the opposite charge. Delocalized groups bearing a net formal charge are also counted as positive or negative features. Feature points are placed on the centroid of the heteroatoms of the group

^a Abbreviations: GKSST, Greene, Kahn, Savoj, Sprague, Teig; rmsd, root mean-square deviation; HTS, high-throughput screening; D, donor; H, hydrophobe; GA, genetic algorithm; MIF, molecular interaction field; SIFt, structural interaction fingerprint; MEP, molecular electrostatic potential; XED, extended electron distribution; DHFR, dihydrofolate reductase; CDK2, cyclin-dependent kinase 2; HIV, human immunodeficiency virus; PDB, Protein Data Bank; GPCR, G-protein-coupled receptor.

(e.g., the centroid of the two oxygen atoms of $-\text{CO}_2^-$). Some authors are careful to refer specifically to *ionizable* positive and negative features, emphasizing, for example, that they count RNH_3^+ as a positive feature but not RNMe_3^+ . This makes sense because the intermolecular interactions of these two groups are very different; the former can donate very strong hydrogen bonds, the latter at best only very weak hydrogen bonds.

Although most pharmacophore programs offer positive and negative features, it is arguable whether they are necessary. The underlying assumption is that a particular part of a protein binding site will accommodate either an ionized hydrogen-bonding group or an un-ionized group but not both. In practice, there are certainly exceptions to this rule. An alternative approach is to dispense with positive and negative features and use instead the concept of hydrogen-bond similarity. Each donor or acceptor group is assigned a strength (ionized groups having higher strengths than un-ionized); when an overlay is subsequently scored, greater credit is given if groups of similar strength are overlaid on one another.^{37–39} Of course, one should always be aware that the specific system under study may not necessarily conform to the general rule because of the complex interplay between protein–ligand interactions and solvation/desolvation effects.

Hydrophobic Features. Deciding which atoms should be considered hydrophobic is not straightforward. For example, a peptide linkage, though not at first sight a hydrophobe, often forms contacts to hydrophobic residues above and below the plane of the peptide atoms. The GKSST algorithm for placing hydrophobic feature points first assigns a hydrophobicity value to each atom, using a set of empirical rules based on the judgements of medicinal chemists. For example, a carbon atom three bonds from a double-bonded oxygen is assigned a hydrophobicity of 0.6. Each atom's hydrophobicity is modified to take account of its solvent accessibility (hence, the algorithm is conformation-dependent; many other programs neglect the solvent-accessibility check to avoid this complication). Finally, neighboring atoms with sufficiently large hydrophobicity values are clustered together into groups and a feature point is placed at the centroid of each. The algorithm looks first for rings (weighting the centroid by the hydrophobicities of the ring atoms), then for groups such as $-\text{CF}_3$, then for chains. Chains are divided into contiguous, small groups of atoms. Many other programs in common use (e.g., PHASE and MOE) implement similar procedures, though they vary in detail.³¹ Some programs, e.g., QUASI,²⁷ place feature points on all atoms that are not donors or acceptors, which has the merit of being simple and producing a fine-grained representation of molecular shape. These feature points are classified in QUASI as “steric” rather than “hydrophobic”, which emphasizes that there are two distinct, albeit closely related, aims: to get a good match of hydrophobic regions of molecules and to achieve a good steric match. This distinction is explicit in the software of Cheeseright et al., which has both van der Waals points and hydrophobic points.⁴⁰ Ultimately, there is no “right” way of placing the points, since the hydrophobic regions of different ligands may not overlay in a tidy group-on-group fashion. For example, a common strategy is to place points at the centers of aromatic rings, but experimentally observed overlays (e.g., Figure 3) do not necessarily have such rings neatly positioned on top of one another. Many programs represent certain types of

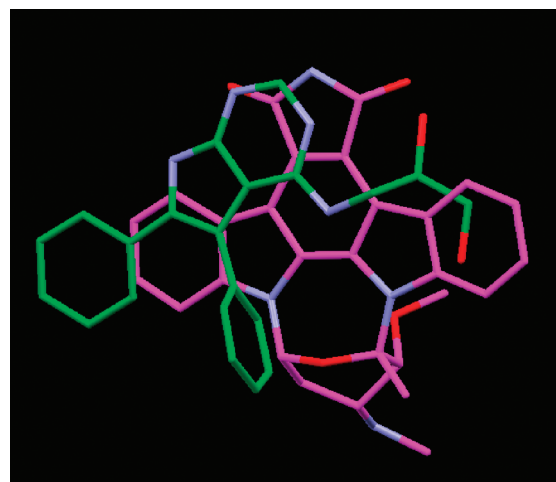


Figure 3. Experimentally observed overlay of two Chk1 kinase ligands, from PDB structures 2bro (carbon atoms colored green) and 1nvq (carbons purple), showing a poor match of ring centers.

hydrophobes, notably aromatic rings, by vectors rather than points so that directional interactions like π stacking can be taken into account when overlaying. In general, however, hydrophobic interactions are often not as well-defined geometrically as hydrogen bonds and so are hard to localize to a small region of a molecule.

The most important purpose of features is to enable molecules to be overlaid in a sensible way, by least-squares fitting of feature points. In addition, feature points may be used in scoring; i.e., the quality of an overlay might be judged by computing the root-mean-square deviation (rmsd) between matched feature points. However, scoring can also be done without reference to the feature points or with reference only to points relating to certain types of features. The placement of feature points is more critical if they are used for both generating and scoring an overlay rather than just the former. A balance needs to be achieved between including all possible features and having sufficient information to give a meaningful and parsimonious model. In particular, if there are many more hydrophobes than other features, then these may dominate the alignment. While this may give a better model according to scoring metrics, the model itself might be less useful for further work, as it will be less discriminating.

Customized Features. It is unlikely that there is a universally “best” set of feature types, though most would agree on an irreducible minimum of hydrogen-bond donor, hydrogen-bond acceptor, and hydrophobe. Depending on the problem, other desirable types might include positive and negative centers, metal coordinators, and different types of hydrophobes (e.g., aromatic and aliphatic). Interactions involving halogen atoms have been the subject of recent studies and may in some cases be usefully incorporated into a pharmacophore model.^{41,42} Further, the user may wish to define a very specific type of feature to exploit a known structure–activity relationship. For example, a chemist working on carbonic anhydrase II inhibition, aware that most known inhibitors have a terminal sulfonamide or sulfamate group, might wish to define a feature of precisely this type. As a consequence, many programs allow user-defined feature types.

Another adjustable parameter is feature separation (the minimum distance between two features for both to be

counted). Does a carboxamide count as both donor and acceptor, and does phenol count as both hydrophobe and donor/acceptor? While these two cases argue for a small separation, does a carboxylate count for two acceptor atoms or, worse, four virtual points? This can lead to overweighting of certain functional groups in the model, which is particularly harmful if the main function of the group is to improve the developability properties of the molecule (e.g., solubility) rather than to interact with the protein.

Searching for Candidate Alignments. This is the most difficult part of the methodology to review because of the great diversity of published approaches. To reiterate: alignment can be conceptually divided into a feature matching problem and a conformational search problem. Perhaps the most telling difference between algorithms is how they deal with molecular flexibility: in particular, whether they operate on a pregenerated set of conformations for each ligand or, alternatively, alter ligand conformations on the fly during the alignment process. (Strangely, few programs do both; i.e., start with a set of pregenerated conformers but then tweak them to optimise the alignment.) Wolber et al. argue in favor of the former:⁴³ “Methods that tweak the molecule while fitting must dramatically reduce the search space while aligning in order to stay efficient and therefore bear the danger of falling into a local minimum.” However, using precomputed conformations is itself a reduction in search space, especially for highly flexible molecules when the granularity of the conformational sampling is likely to be coarse. Thus, the possibility exists that the true answer will not be found because it lies between the conformations available to the algorithm. The ability of widely used programs to search conformational space has been extensively examined in recent years.^{44–48} Some of these studies have also investigated the conformational energy landscape and the strain energy of bound ligand conformations.^{44,46} These studies provide useful guidelines for use not only in pharmacophore elucidation but also in applications such as database searching (vide infra). An important practical point is that while most pharmacophore elucidation programs contain their own conformer search procedures, they do also permit the import of conformations generated with an external specialized conformer search procedure. Such flexibility provides the user with the ability to choose the best method for each distinct task in the pharmacophore elucidation process.

Algorithms Using Pregenerated Conformers. To avoid confusion, we use consistently the terms defined by Barnum et al.¹⁶: a configuration is a set of points in 3D space, each associated with a type of feature; a partition is an object associated with all configurations of a particular type; e.g., the DHH partition is associated with all configurations consisting of one hydrogen-bond donor (D) feature and two hydrophobic (H) features. If two molecules both contain a configuration that belongs to the same partition and that can be superimposed, one onto the other, within a specified threshold rmsd, this configuration is common to those molecules.

The classic way to identify common pharmacophore alignments from precomputed conformations is to use clique detection, first employed in the DISCO program^{49,50} and a mainstay of more recently published algorithms. However, the CATALYST/HipHop algorithm has probably been more widely used than any other. It operates on precomputed conformers and works by finding small common

configurations and then attempting to enlarge them. The authors describe it as a “pruned exhaustive search”. To begin, one or more molecules in the set to be aligned are chosen as reference molecules. All two-point configurations in all conformations of the reference molecules are found. Configurations are rejected if either of the features is not surface accessible or if the interfeature distance is below a user-specified limit (this is to encourage diversity among the final answers). For each of the survivors, a search of the conformations of all the other molecules is performed to see whether the configuration is common to all ligands. If so, it is kept. Eventually, a list of all common two-point configurations is compiled.

The algorithm then enumerates all three-point configurations in the reference molecule(s). Any such configuration contains within it three two-point configurations (e.g., a configuration belonging to the partition DHH will contain two DH configurations and one HH configuration). It can be rapidly ascertained (by checking the list of common two-point configurations) whether these three two-point configurations are common to all ligands. If this is not the case, the three-point configuration is rejected (this is the pruning step). Otherwise, the remaining molecules are checked to see whether the three-point configuration is common to all, and a list of common three-point configurations is built up.

Repetition of the procedure allows increasingly large common configurations to be found until no further increase in size is possible. As well as pruning, computer time is also saved by using a two-step procedure to ascertain whether a given configuration is present in a molecule. First, a precomputed list of all the interfeature distances in the molecule is checked to see whether the configuration might be present. Only if this filter is passed will a final least-squares fit of feature points be performed to prove the matter one way or the other. Least-squares fitting is used in almost all pharmacophore programs as the final arbiter of whether a configuration is common to two molecules, since it distinguishes mirror images.

The claim that the algorithm is exhaustive is debatable. First, the number of solutions found depends on the granularity of the conformations presented to the algorithm. Second, the decision on whether a configuration is common or not depends on the threshold rmsd. Third, the algorithm will, if the user requests, allow certain molecules to miss a feature in a “common” configuration, provided the total number of such molecules remains below a specified limit. This is a highly desirable option, since such situations occur frequently in reality, but it again means that the number of solutions produced by the algorithm depends on a user-chosen parameter. In truth, no pharmacophore elucidation algorithm is exhaustive, despite occasional claims to the contrary.

The PHASE algorithm employs a tree-based partitioning method that rapidly groups together similar configurations according to their interpoint distances. For each conformer of each molecule, the program determines all k -point configurations that are present. Each configuration is represented by its vector of $n = k(k - 1)/2$ interpoint distances. Each interpoint distance is filtered through a binary tree to assign it to a particular distance range. Once all n distances have been filtered, the configuration has effectively been assigned to an n -dimensional box whose sides are equal in length to the range corresponding to the bottom of the

decision tree. Now, if two molecules each have a four-point configuration belonging to the same partition (e.g., DDHH) and falling in the same box, that configuration might be common to those molecules. (In addition, it is necessary to consider neighboring boxes, since similar distances may fall on either side of a box boundary; Wolber et al.⁴³ and Zhu and Agrafiotis²⁸ describe methods for circumventing this troublesome problem.) Potential common configurations of various sizes are therefore identified rapidly. A least-squares fit is then performed to determine which of these are truly common. A given k -point configuration is only accepted if it occurs in a user-specified minimum number of molecules.

Algorithms That Alter Conformations on the Fly. An early but instructive example of this type of algorithm is GASP,¹⁸ which uses a genetic algorithm (GA) to optimize an initial population of random overlays. The ligand with the smallest number of features is chosen as a reference molecule. Each population member is represented by a chromosome containing (a) columns specifying torsion-angle values for all the rotatable bonds in all of the molecules and (b) "mapping columns", which define pairings between features of each nonreference molecule and features of the reference molecule. This is enough information to generate a molecular overlay, by setting each molecule to the conformation defined by the appropriate torsion-angle columns and then least-squares fitting using the feature-point pairings in the mapping columns. A scoring function is used to estimate the fitness of each overlay, taking into account the quality of geometric matching between paired feature points, the extent of volume overlap, and the ligand strain energies. In each step of the GA, mutation of one parent chromosome (changing either a torsion angle or a mapping column) or crossover of two is used to generate a child (two children in the case of crossover), which, if sufficiently fit, can be used to replace the least-fit member(s) of the population.

The successor to GASP, viz. GALAHAD,¹⁷ divides the problem into two steps. In stage 1, which operates solely in ligand torsional space (i.e., the ligands are never actually overlaid), the GA is used to find a set of ligand conformations that tend to have large common volume and low steric energy and are similar in the feature configurations they contain. The latter is determined by generating fingerprints for the ligands that capture the various triplets, quadruplets, etc. of features and their interfeature distances. Stage 2 employs an algorithm from image recognition to produce the final alignment of the ligands, held rigid in their conformations from the first step.

In summary, while the diversity of approaches defies easy classification, some common characteristics can be identified. Most algorithms have strategies for reducing compute time, most commonly the use of pruning or fingerprints. The use of clique detection and search techniques such as genetic algorithms is common. Least-squares fitting of matched feature points is almost always used to produce the final alignments.

Scoring. Most algorithms rank candidate overlays by means of a scoring function containing some or all of the following terms: feature matching; volume overlap; strain energy; selectivity. Some of these terms have an underlying physical rationale; others are more subjective, being included to increase the score of the more "relevant" pharmacophores. Often, a particular molecule in the overlay, the reference, has a special status in the scoring procedure.

The reference may be chosen arbitrarily, or its selection may be biased toward particularly active or less flexible molecules.

The function used in PHASE, which may be written as

$$\text{score} = F + w_v V - w_e E + w_m M^{-1} + w_s S$$

serves as a typical example. w_v , w_e , w_m , and w_s are user-defined weights. F measures the quality of alignment of features in the reference molecule with the corresponding matched features in each molecule. Two criteria, appropriately weighted, are used: the rmsd of the matched feature points and the average cosine of the angles formed by matched pairs of vector features (e.g., aromatic rings). V is the average over the nonreference molecules of their volume overlap with the reference, measured as (intersection volume)/(union volume). E is the strain energy of the reference ligand conformer from which the pharmacophore is derived. PHASE allows configurations to be found that are common to some but not all of the molecules. M in the term $w_m M^{-1}$ is the number of molecules that contain the configuration; so, depending on the user-chosen value of w_m , this term can be used to reward overlays in which M is large. S is a measure of selectivity, i.e., an estimate of the fraction of molecules in a random database likely to match the common configuration of the overlaid molecules. The smaller this fraction is (i.e., the more unusual the arrangement of features in the common configuration), the larger S will be. The logic behind this term is that if the common configuration is unusual, it is reasonable to conclude that it occurs in the overlaid molecules because it is a requirement of their activity. It is too time-consuming to estimate S by searching a large database for the common configuration; Dixon et al.²⁰ and Barnum et al.¹⁶ describe how it may be approximated.

Labute et al. describe a scoring mechanism based on the use of Gaussian functions.⁵¹ In their formalism, all properties to be considered in the scoring (molecular volume, distribution within the molecule of hydrogen-bond donor and acceptor groups, etc.) are described in a consistent manner. Thus, for a given conformation of a molecule, the density of a property P at a point in space x is described by the equation

$$f_P(x, x_1, \dots, x_n) = \sum \{ (w_i/n) [a^2 / (2\pi r_i^2)]^{3/2} \exp[-a^2(x - x_i)^2 / (2r_i^2)] \}$$

where (x_1, \dots, x_n) are the positions of the n atoms of the ligand, r_i is the van der Waals radius of atom i , a is an empirical parameter, w_i is the weight of property P at x_i , and the summation is over the n atoms. For example, if property P is hydrogen-bond donor ability, then w_i could be set to 1 if atom i is a donor and zero if it is not. The degree of overlap of property P for a pair of overlaid molecules can thus be calculated rapidly as a sum-of-Gaussians density in the interatomic distances. The degree of overlap of several properties is a weighted summation of the individual property-overlap equations, the weights reflecting the relative contribution of each property to the overall score. The consistent manner in which properties are treated lends itself to programming simplicity and makes it easy to add new properties to the scoring function.

A problem with the scoring methods described so far is that they require weights to be assigned to the different terms, and these weights are inevitably arbitrary. Some programs avoid the problem by Pareto ranking,⁵² which

produces a population of possible overlays, each representing a different compromise between conflicting criteria (e.g., ligand strain energy, volume overlap, feature-matching). Specifically, the aim is to converge on a set of nondominated overlays such that for any given overlay there is no other overlay in the set that scores better against all criteria. The approach therefore recognizes that it is impossible to predict, for any given set of ligands, which criterion will be most important.

Practical Aspects and Applications. It is to the credit of software developers that pharmacophore elucidation methods are no longer limited to a small number of expert computational chemistry practitioners but are much more widely accessible. However, it would be unfortunate were the reader to assume that such methods have reached a degree of accuracy that they can be used “blind”. Software programs such as these are no different to a sensitive instrument that requires knowledge and expertise to obtain the maximum benefit. As we have already indicated, many programs contain a number of user-selectable parameters that govern (sometimes dramatically so) the amount and quality of output produced.

Three main stages can be identified in the elucidation of a pharmacophore. First, prepare the data set. Second, generate possible pharmacophores. Third, validate the pharmacophore(s). The compounds used to construct a 3D pharmacophore should all have the same mechanism of action (e.g., agonist, antagonist, inhibitor, binder) and should preferably have a high affinity as measured in an appropriate biological assay (e.g., full-curve dose response). If a competitive binding assay is available, then this may provide additional confidence that the ligands bind in the same region, though it cannot guarantee that they share a common binding mode. The ideal data set contains ligands from a number of different chemical series having limited conformational flexibility and without too many heteroatoms. Each molecule should be inspected to ensure that it is represented as the appropriate tautomer and that the appropriate ionization state is used; depending on the software used, these aspects may have a significant impact on the quality of conformations generated. It is also important to check the output from any cleansing procedures that may be applied to the input ligand structures. As the typical pharmacophore elucidation data set contains a relatively small number of molecules, it should not be necessary to sacrifice speed for quality; it is also advisable to visually inspect the conformations for potential issues such as *cis* amides, high-energy axial ring substituents, and close contacts. Most pharmacophore generation programs will produce a number of results and will attempt to score them as outlined above. These scores can be a useful guide, but it is important to check for the quality of fit of each molecule to the pharmacophore, for a sensible common volume overlap and to ensure that the pharmacophore is consistent with the known SAR. Stereoisomers that show different activities can be very informative in helping to distinguish between possible pharmacophores. Inactive molecules can also be very useful. For example, if an inactive molecule gives a good match to the pharmacophore, then this may help to identify excluded regions. Sometimes it proves impossible to derive a satisfactory pharmacophore model. Clearly this may happen if the fundamental principle of a common binding mode does not apply. It may also be helpful to check the structures for functional groups that may not be involved in molecular

recognition with the biological target (e.g., are solubilizing groups). Other strategies include relaxing various criteria (permitting partial matches). Ignoring the inactive compounds on the basis that a pharmacophore is a necessary but not sufficient requirement for activity may also help.

As will be discussed in the next section, 3D pharmacophores are widely used for virtual screening of databases in order to identify new lead series. They also have significant utility in other areas, particularly when no detailed structural information on the biological target is available. Thus, a pharmacophore can help guide the design of focused arrays, can help to rationalize the SAR of a chemical series, and can enable different chemical series to be combined and information to be transferred from one series to another. In these and other applications the ability to visualize a 3D pharmacophore is an especially useful aspect. In all such uses, however, one should always remember that a pharmacophore is but a hypothesis that should be continually tested, refined, and possibly rejected.

3D Database Searching

One of the common purposes of making pharmacophore models is to search for novel chemical matter. At the time that the first pharmacophore searching algorithms were implemented in the early 1990s the discovery landscape was quite different. X-ray structures of proteins were relatively rare, and it was much more common to work with leads derived from endogenous ligands or from the literature. At the same time, the first generation of screening and compound handling robots were being developed, together with electronic databases of corporate archives. Most chemists are familiar with 2D substructure-based searches, but these typically find only compounds with similar scaffolds or members of the same structural family. The pharmacophore represents an abstraction that can be used to find alternative chemotypes (i.e., chemical series with a different underlying framework, scaffold, or common moiety). All that matters is the disposition of molecular recognition features, not the underlying pattern of atoms and bonds. The stage was therefore set to search large databases for compounds (both “real” and “virtual”) that could exhibit the desired pharmacophore and thus (hopefully) have the same biological activity.^{53–56} Depending on the precision of the query, one can find numbers of hits from 10s to 1000s, which was in line with the screening capacities available at the time. Many will be false positives and show no activity in the screen, but generally, the hit rates from pharmacophore searches are much higher than from random screening. The hits can also sample very novel and diverse chemotypes, allowing the medicinal chemist the luxury of pursuing the series with the best overall profile. Today, it can be equally cost-effective to screen the whole corporate collection as to screen a significant subset. Thus, the main role of pharmacophore searching in industry has changed, with a principal emphasis being on the creation of small focused sets for low-throughput, higher quality assays to enhance the lead identification process in parallel with high-throughput screening (HTS). The sources of the compounds in such focused sets can be both internal and external (i.e., compound vendor). In academia, pharmacophore searching still plays a useful role in reducing screening set size to manageable proportions due to assay throughput; moreover, the main source of compounds is from commercial catalogues.

Database Generation. All the issues highlighted previously now return in starker form, when a 3D pharmacophore is used as a search query for a large corporate database. It is desirable to build the database so that it can be used generally, without further modification. The more information that can be precomputed, the faster the search will be but the larger the database will also become. However, if the hard-coded parameters used to build the database do not match those used ad hoc to build the query, the retrieval rates will be greatly reduced. The key parameters to be considered are tautomerism, stereochemistry, conformational sampling, and feature definition. As the usual approach is to capture as many actives as possible, even if that means more false positives, most databases are built to include all possibilities; that is, all tautomer and stereoisomer forms are added as duplicates. Conformational sampling is a particular headache because of the trade-off between coverage and size of the database/search time. Broadly, there are two strategies: (1) compute the conformations on the fly; (2) precompute and store. One can also use single-conformer databases; if a good quality structure generation program is used, the single conformer ought to be close to, or even at, the global energy minimum. If this conformation can match the pharmacophore without further adjustment, the physical interpretation is that no strain energy is present, thus improving the energetics of binding. This low-energy match should have a better chance of translating into an experimentally active compound. However, the hit rate using flexible searching has been found to be 2–10 times higher.⁵⁶ The business rules used to construct the database need to be clearly documented for subsequent users so that the pharmacophore queries can be adapted accordingly. The basic workflow involves two key steps: for each structure in the database, determine if the structure contains a predefined number of features from the query, and if it does, can the structure be made to display those features in the right geometry, within a defined tolerance?

Feature Matching. Feature definition is generally hard-coded, leading to very fast search times; one can prune any molecule from the search that does not contain the partition, before searching for the configuration. The feature definitions are encoded in a dictionary; when the database is built, the type and number of features in each molecule can be determined and stored. The same is done for the query, and only molecules with the same number (or more) and type of features are retained for the time-consuming next step. However, this can lead to a loss of applicability if the pharmacophore query uses even slightly nonstandard definitions. The alternative, that of making the features very wide-ranging, will reduce the power of a query to effectively discriminate between true hits and false positives. In the event that no matches are found, there is often the option of allowing partial matches (e.g., to match four out of five features). This assumes that the missing features can be designed in later (or perhaps sufficient activity can be obtained even with a partial match). Some programs allow relative weighting of features. This facility permits one to designate which features contribute most to the SAR, for example, the basic nitrogen for serotonin mimetics. This can also be done after the search of the final hit list. Feature matching can give high screen-out rates (80–100%), depending on the complexity of the query.

Geometry Matching. Once the structures that might match the query are found, the next step is to examine the

conformations of those structures to see if the geometric disposition of the features can be matched. Note also that a pure distance comparison cannot distinguish enantiomers; for this, an actual fit (in 3D space) to the query is required. The goodness of fit to the query is defined by the tolerance allowed in the match between features. One might be looking for a donor–acceptor pair separated by 5 Å, but be prepared to accept a value between 4 and 6 Å. Not all interfeature tolerances need to be the same; they should be set using knowledge of the SAR. Matching of the features should be possible without introducing undue strain into the structure. This is another parameter that can be defined. It is not possible to give specific rules, as the acceptable limit is force-field-dependent, but a strain energy of more than a few kcal/mol will lead to many false positives. Features can also be negative, for example, excluded volumes. Here, the aim is to avoid rather than match. Inclusion volumes, for example, derived from the ensemble of aligned active ligands, can also be incorporated. Conformational analysis of all but the simplest structures quickly becomes very computationally expensive. Three strategies can be used, each with their good and bad points. The first strategy is used in UNITY,⁵⁷ which stores single conformers and uses the directed tweak algorithm⁵⁸ to perform flexible searches. The query is mapped onto the molecule, feature by feature, and a minimization in torsion space is conducted to satisfy each interfeature distance. The results are then checked for steric clashes. This is important to avoid unreasonable conformations with high internal energy. Multiple mappings of the query are possible, but once one good match is found, the search can be terminated. Essentially, the molecule is pulled by the query into the right shape, so the strain energy is the key parameter. Screen-out rates are lower than for the other methods, but with any rise in true hits found, there will also be a rise in false positives. This is a strategic decision that should be made in the wider context of the project. The second strategy is adopted by MOE, CATALYST, and PHASE, which all use a database of conformations. Each conformation can be quickly matched as a rigid body to the query, and if it does not fit, it is discarded. The strain energy can be preset so that all conformations are reasonable. The disadvantage is that normally only a relatively small number of conformations per structure can be stored (typically 250 or fewer). The selection of conformers may be performed using a poling approach to ensure “conformational diversity”^{59,60} and predefined preferences, for example, cis/trans ratios for carboxamides or axial/equatorial ratios for saturated rings. The key issue here is whether the conformational sampling is of sufficient resolution to be useful, especially for more flexible molecules. In the third approach, compounds that pass the feature mapping stage are subjected to conformational analysis on the fly according to a preset recipe (typically a rule-based enumeration that uses a predefined number of torsional increments for each type of rotatable bond). For each conformation, a rigid body match is performed as before and the analysis aborted once a fit is found. Only the torsions that can affect the interfeature distances are analyzed. If the query is only loosely defined (few features, wide tolerance allowed in the fit), this approach can be very time-consuming. A sensible upper limit for the number of rotatable bonds is usually defined to prevent the search time being dominated by the matching of just a few very flexible

structures. The directed tweak can produce strained conformations (leading to a higher false positive rates), but the other approaches can fail to generate the matching conformation (higher false negative rate). This is the reason why it is not possible to say that one strategy is superior to the other; so much will depend on the query and the database and the nature of the decoys. A final variant used by CATALYST is its BEST option, which is to take structures that nearly match the query, then tweak them.

A very important test is that the search procedure should find the molecules used to build the original pharmacophore; failure to do so is an excellent way of losing any credibility! As with pharmacophore elucidation, there are a number of locally adjustable parameters. These include the tolerance allowed in the fit and the number of features which can be skipped; naturally, loosening the search will increase the number of hits returned, so the capacity of the biological screen should be carefully kept in mind. Validation of the search parameters is a very valuable exercise that can be performed on a small test database, seeded with known actives and decoys, to see how changes to the query can change the enrichment statistics. It is strongly recommended to tune a query this way before running the main search.

Postprocessing. A hit list of molecules that match the pharmacophore query should be postprocessed. Highly flexible, feature-rich molecules such as peptides can match most queries. Parsimonious matches might be more promising. Where the information exists, exclusion or inclusion volumes are often used; this will be discussed further in the section on using protein structures. Alternatively, simple measures of how much of the molecule falls within the pharmacophore, and how much outside, can be used to prioritize the virtual hits for inclusion in the final set for screening. Using the same scoring function as that used to create the query is possible, but given the inaccuracy of the scoring function, visual inspection is often a better guide. The simplest score is the rmsd fit to the query. One can also filter using physicochemical criteria (i.e., only returning hits that are druglike or leadlike in nature). However, a drawback of pruning too severely at this early stage is that novel chemical matter that could subsequently be modified to give a viable lead series would be lost.

Pharmacophore Keys

A pharmacophore key is a binary descriptor of the partitions and configurations accessible to the molecule. Such a key (also called a fingerprint) is generated by performing a conformational analysis for each molecule and abstracting the interfeature geometries for each conformation sampled. Configurations are handled by binning the distances; for example, one specific bit in the fingerprint is set if a donor and acceptor can be positioned 4–6 Å apart. By analogy with the use of the 2D fingerprints used in substructure and similarity searching, the pharmacophore key can provide a way to improve search speeds for pharmacophore searches of 3D databases: the first phase of searching becomes a simple AND operation such that only molecules that pass this filter are matched further. Two-center pharmacophore keys, in which bits are assigned to pairs of features, are probably not specific enough. If we consider 7 features (i.e., Donor, Acceptor, Hydrophobe, Aromatic, Positive, Negative, and Donor/Acceptor) and 17 distance bins, the key would contain $7 \times 6 \times 17 = 714$ bits, and many bits would be turned on. The ability to discriminate

between the known actives and inactive decoys is diminished. A three-center key, in which each bit represents a triplet of features, needs to be 184 884 bits long. This type of key is sparse and so much more selective. In addition to binary keys, one can also count the frequency that a triplet is found during the analysis of a structure. Four-center keys are even larger.⁶¹ Using 6 feature types and 10 distance bins leads to 24 million distinct quadruplets, or tetrahedra. The advantage is that such keys are now able to distinguish chirality.

Pharmacophore keys capture information about shape and intermolecular interaction propensity. In contrast to the “traditional” pharmacophores described thus far, they do not directly identify the particular arrangement(s) of features required for a specific biological activity. Rather, they are whole-molecule descriptors that have been used for a variety of similarity and diversity applications, driven in particular by the needs of technologies such as high-throughput screening and combinatorial chemistry where the number of chemicals to be processed is very large, for example, screening decks or virtual combinatorial libraries. In library design one might like to answer the question, “If I add this reagent or remove that reagent, will the properties of my library improve?” For example, in designing a library of type I kinase inhibitors, one will want to improve the overlap to the type I pharmacophore whereas if the library is being designed for more general enrichment purposes, one may want to minimize the overlap with the current screening deck but maximize the overlap with the space of known drugs. Pairwise comparison of even a small library of 1000 against a deck of 1 000 000 is not feasible. Pharmacophore keys offer an alternative approach.

Ranking a set of compounds in a database according to their similarity to a query molecule is perhaps the most straightforward application of pharmacophore keys. New and modified similarity metrics have been developed to deal with the particular requirements of the pharmacophore key, which is in general larger and more sparse than a standard 2D fingerprint.⁵⁵ The union of the pharmacophore keys for a collection of molecules provides an overall profile that can be considered a diversity metric and can also be used to assess the extent to which the set adds new pharmacophores to an existing collection when purchasing compounds or in library design.⁶² A pharmacophore key can also be used as a descriptor vector in QSAR modeling.^{63,64} Here, the bits provide the basis of discriminating between actives and inactives. The nature of the pharmacophore key means that such models can in principle deal with multiple binding modes. Depending on the model-building technique used, it may also be possible to interrogate the resulting model to determine which pharmacophore(s) are responsible for the activity. As always, it is important to ensure that such models are robust and do not suffer from overfitting.

3D Pharmacophores and Protein Structure

Pharmacophores can be derived from the binding site of a protein. This is particularly relevant, as the number of protein structures has greatly increased not only in simple numerical terms but also with regard to the breadth of coverage of gene families of interest to drug discovery. Moreover, the distinction between “structure-based” and “ligand-based” drug design methods is now much more diffuse and the judicious use of traditional ligand-based methods such as 3D pharmacophores can greatly enhance the efficiency and effectiveness of structure-based design.

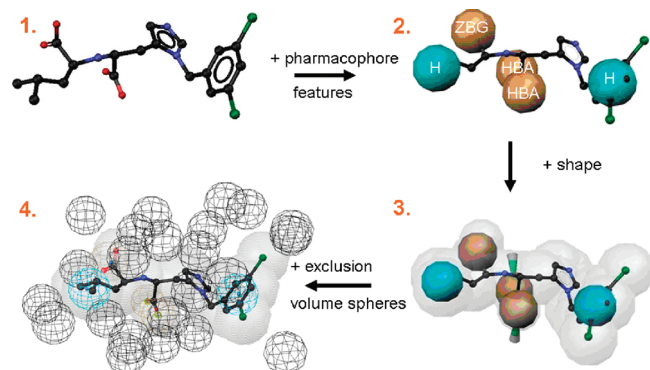


Figure 4. Illustration of how a 3D pharmacophore can be enhanced by the addition of shape and excluded volume information, in this case to identify inhibitors of angiotensin converting enzyme. Figure is reproduced with permission from Rella et al. *J. Chem. Inf. Model.* 2006, 46 (2), 708–716.⁶⁷

The most straightforward way to derive a 3D pharmacophore from a protein structure is through the direct observation of specific interactions between protein and ligand(s). Such a pharmacophore can then be used in the usual way, for example, to search a 3D database in order to identify compounds for focused screening. Database searching methods based on 3D pharmacophores are in general much faster than structure-based methods such as docking, and so this can be a more effective way to screen very large databases. Alternatively, the 3D pharmacophore search can act as the first stage in a docking workflow. A large database can be prescreened, using a pharmacophore query, to create a much smaller subset for docking. A pharmacophore, if properly constructed, will give very few false negatives (false negatives might arise from different binding modes or poor feature definition). The scores of docked hits may then become more useful.

Knowledge of the protein structure enables the pharmacophore to incorporate more detailed information about regions that are not accessible to the ligand. This is most commonly achieved through the use of exclusion volumes and/or inclusion regions.^{65,66} An illustration of how a 3D pharmacophore can be enhanced by the addition of such information is given in Figure 4 for the development of a pharmacophore to identify inhibitors of angiotensin converting enzyme 2.⁶⁷ When used for database searching, such pharmacophore queries have the distinct advantage of finding hits that not only contain the key binding elements but are also able to fit into the active site, thereby reducing the false positive rate. The extent to which an ensemble of excluded volumes can represent the atomic detail of a protein structure is obviously limited (though one should remember that proteins can be flexible, so the apparent precision of an atomic-resolution protein structure may be illusory). In addition, the widespread use of a two-stage algorithm further limits their utility. In this approach matches to the underlying pharmacophore features are first identified. The matching conformation is then tested for violation of the exclusion or inclusion criteria, often without any attempt to modify the conformation to fit these constraints. This can lead to potential hits being missed even though a relatively simple bond rotation could alleviate a clash. Conversely, the use of a limited number of exclusion volumes to represent an entire protein active site can result in “holes” through which a ligand may protrude. We also note that widely used shape searching methods^{68–70} provide an alternative way to tackle this problem. Such methods will be

the subject of a separate Perspective and so will not be considered further in any detail here except to note that those methods based purely on shape often incorporate elements of pharmacophore recognition in order to improve the efficiency and sensitivity of their algorithms.

Structure-based 3D pharmacophores derived solely on the interactions observed in known protein–ligand complexes may be unnecessarily restrictive. An alternative is to define pharmacophores based on an analysis of the “hot spots” in the active site. A number of methods can in principle be used to identify such hot spots (or site points). These include programs such as GRID⁷¹ that probe the site with small molecules or functional groups and calculate the enthalpy of the interaction between the probe and the protein atoms at points on a grid lattice to generate a molecular interaction field (MIF). These fields can then be contoured by energy to find the most favorable regions for an acceptor or a donor (or any other type of feature) to interact with the protein. Programs such as LUDI⁷² and SUPERSTAR⁷³ use a knowledge-based approach in which rules are used to generate a set of interaction sites for each atom or functional group of the protein that is capable of participating in a nonbonded contact. The rules are largely based on statistical analysis of experimental structures from the Protein Data Bank or small molecule crystal structures and take into account the chemical nature of the atoms as well as the orientational preferences of features such as hydrogen bond donors/acceptors. From the locations of the site points it is then necessary to construct one or more 3D pharmacophores. The simplistic approach would be to combine all such points into a single pharmacophore, but the resulting query would typically be matched (if at all) only by a molecule filling the entire active site. Rather, all possible 3D pharmacophores (containing three, four, or more features) are enumerated from the set of site points. The most important step in such a procedure is to triage the set of pharmacophore queries. Commercially available programs that can perform the entire process from site to pharmacophores include Structure-Based Focusing⁷⁴ and LigandScout.⁷⁵

The combination of pharmacophore keys discussed previously and techniques for site analysis allows one to avoid making (possibly subjective) choices about which features are important or even which model of the active site is the bioactive form. The three-center or four-center key from a binding site is less sparse than the key of a ligand but is still discriminating enough to be useful; it captures the overall character of the pocket. A comparison between the key of a protein and a ligand can quickly highlight common triples or tetrahedra, allowing fast alignment between the two objects. In the FLAP program⁷⁶ four-point pharmacophore fingerprints are derived from a GRID analysis of a protein active site and combined with a shape-based description of the binding site. When combined with a set of ligand pharmacophore fingerprints, this provides an efficient way to perform structure-based virtual screening. Another well-known application of this technique is embodied in the Metasite program.⁷⁷ Here, binding sites of cytochrome P450s have been characterized using the GRID approach and the key fields selected on the basis of experimentally determined sites of oxidation in substrates. New molecules are aligned to the pharmacophore keys to see which C–H bonds can come close to the oxidation center, and the results can be further ranked according to the reactivity of the C–H bond. In tests, the actual site of metabolism was in the top three predictions over 80% of the time.

An obvious extension to searching for ligands is to compare site-based pharmacophore keys across proteins. If there is a high degree of similarity, the pockets are likely to bind the same type of chemical moiety. In addition to the CavBase program,⁷⁸ Nussinov et al.⁷⁹ have developed a much faster method based in essence on pharmacophore triplets. As always, there is a balance between speed and resolution, and the pharmacophore features are crude and not influenced by any local electrostatic perturbations.

Another recent structure-based fingerprint is the structural interaction fingerprint (SIFt),^{80,81} in which each residue in the binding site is represented by a seven-bit descriptor that characterizes whether (1) it is in contact with the ligand, (2) any main chain atom is involved in the contact, (3) any side chain atom is involved in the binding, (4) a polar interaction is involved, (5) a nonpolar interaction is involved, (6) the residue provides hydrogen bond acceptor(s), and (7) the residue provides hydrogen bond donors. The whole interaction fingerprint is constructed by concatenating these residue bit strings together according to their sequence order. The SIFts have been used for a number of purposes. One straightforward application is to compare sets of docking results by calculating pairwise similarities between the fingerprints followed by a cluster analysis. This enables distinct groups of docking poses to be easily identified. By comparison of the interaction fingerprints from docking results with the fingerprints for known protein–ligand complexes, it was possible to achieve enrichments in virtual screening experiments that were superior to those obtained with commonly used scoring functions. Other variations and extensions to the SIFt concept include the use of weighted profile-like fingerprints that represent the interactions of multiple complexes⁸² and encoding more detailed information about hydrogen bonding strengths and geometries.⁸³ Finally, SIFts and other protein structural fingerprints can be used to cluster and compare the structures of different proteins across a gene family, most notably within the kinases. As with the ligand-based pharmacophore fingerprints, the structural fingerprints enable complex data to be stored and manipulated efficiently using well-understood algorithms, albeit with some information loss.

Field-Based Pharmacophores

The pharmacophore generation methods described so far are based on the assignment of features to functional groups within the molecules and are therefore essentially atom-based. An alternative approach is to base the alignment on the molecular fields exhibited by the molecules, recognizing that the interactions between molecules are governed by their overall electrostatic and van der Waals properties. Field-based approaches⁸⁴ thus aim to describe what the receptor “sees” in terms of charge distribution and shape rather than focusing on the underlying structural skeleton.

The molecular electrostatic potential (MEP) around a molecule depends on the distribution of atomic charge and is often modeled on a 3D grid or surface. The MEP at each grid point surrounding a molecule is calculated using Coulomb’s law as the interaction between a probe atom of unit positive charge (a proton) and partial charges that are centered on the atoms. The resulting grids typically consist of a large number of data points so that direct alignment of the grids themselves is generally considered too computationally intensive to be feasible. However, Good et al.⁸⁵ showed that

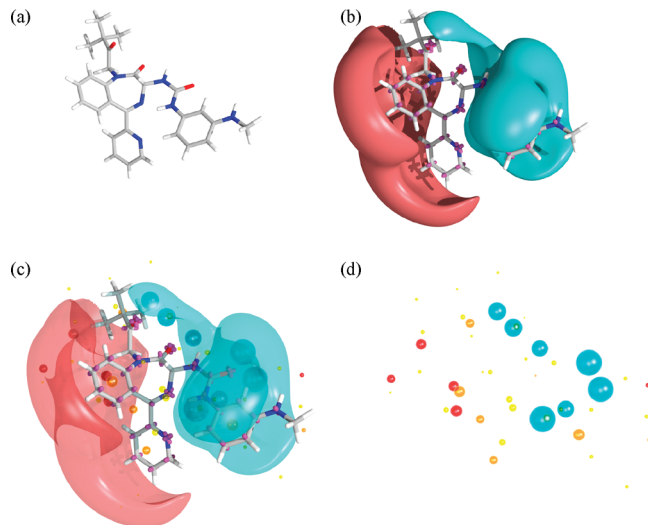


Figure 5. Steps involved in generating field points around a molecule using the method of Cheeseright et al.⁹² (a) a molecular conformation; (b) the electrostatic potential map calculated using XED; (c) field points superimposed onto (b); (d) the final set of field points including electrophilic (red), nucleophilic (blue), van der Waals attractive (yellow), and hydrophobic (orange) points. Figure is reproduced with permission from Cheeseright et al. *J. Chem. Inf. Model.* **2006**, *46*, 665–676.⁹²

the MEP can be modeled using atom-centered Gaussian functions; these provide an elegant way of overcoming the problem of handling the large number of data points in a grid. The Gaussian approximation allows the rapid comparison of the similarity of two fields using the cosine similarity coefficient. Such a similarity can then be used as the function to be maximized in an optimization of the alignment using, for example, a gradient-based method⁸⁶ or a genetic algorithm.⁸⁷

The MIFs produced by the GRID program⁷¹ and mentioned earlier in the context of protein binding sites can also be used to represent small molecules. However, MIF fields are not straightforward to model using Gaussians⁸⁸ so that identifying the optimum alignment of such fields remains a challenging task. One approach to using GRID fields that obviates the need for alignment is the GRIND method⁸⁹ in which extrema in the field are identified and mapped to a fingerprint to give a vector representation that can then be used in similarity searching as for the pharmacophore fingerprints described previously. It is not trivial, however, to determine which of the grid points should be included in the descriptor, and the original GRIND approach requires that the user define the number of points that are extracted. More recently, an automated method for extracting hot spots from a field has been developed that avoids the need for parametrization.⁹⁰

The use of atom-centered charges can lead to an inadequate representation of the MEP. For example, a carbonyl oxygen will give rise to a single field maximum extended outward from the C=O double bond. The extended electron distribution (XED) force field attempts to model a quantum orbital description of the distribution of charge.⁹¹ In XED, the charge on electronegative and π atoms is extended away from the nucleus of the atom to give rise to multipoles. For example, the carbonyl oxygen is modeled by two field maxima that correspond to the positions of the lone pairs. The XED force field is the basis of a method to generate a pairwise alignment of two molecules based on the extrema (called field points) in their

respective fields⁹² (Figure 5). Four types of field points can be calculated: positive, negative, hydrophobic, and van der Waals, depending on the potential used. Clique detection based on the field points is then used to find starting points for the alignment which are optimized using a simplex algorithm. Similarity is measured by using the field points in each molecule to sample the whole field in the other, with the size of a “point” determined by the depth of the energy well, to give a more accurate reflection of the field than simply using the points themselves. Using such a multistage algorithm enables large databases to be searched in a reasonable time period.

As with the atom-based pharmacophore elucidation methods, conformational flexibility is most often handled using the ensemble approach in which conformers are pre-computed and handled one at a time. It is also possible to vary conformation on the fly, usually under the assumption that the partial charges on the atoms are invariant to conformation so that it is only necessary to recalculate the fields for different conformers.⁹³

The most common application of field-based alignment has been to virtual screening where the molecules in a database are aligned, one at a time, to a query compound and scored. Usually, a single conformer of the query (the presumed bioactive conformation) is considered with the molecules in the database treated as flexible. Field-based methods have also been used to align multiple molecules. The simplest way to achieve this is to use a single molecule as the template or reference molecule,⁹⁴ but this approach is unlikely to lead to an optimal alignment. An improved approach is to first identify an overlay of two molecules by considering all pairwise alignments of the respective sets of conformers.⁹⁵ High scoring conformer pairs that share common fields can then be extended to a third, fourth, etc. molecule to give a multi-conformer alignment. Such an alignment can then be used as a field-template for database screening.

Pharmacophores and 3D QSAR Methods

Field-based approaches have been widely used to generate 3D QSAR models from sets of aligned molecules, for example, in the widely used CoMFA program. In this case, a variety of different methods can be used to generate the alignment prior to model building, including overlays based on a 3D pharmacophore. In addition, some of the more commonly used pharmacophore methods also use 3D QSAR methods to refine their models or to closely integrate the two techniques. In the HypoGen algorithm in CATALYST the degree to which each molecule matches the query (in a geometrical sense) is assumed to correlate with the observed activity. The PHASE program also provides a facility for constructing 3D QSAR models. These models can be atom-based (i.e., taking all of the atoms in each ligand into account) or pharmacophore-based (i.e., only the features involved in matching the common pharmacophore are used). A rectangular grid is defined to surround the molecules aligned to match the pharmacophore. Each cube in this grid is then characterized according to the atoms or pharmacophore features that fall within it, resulting in a set of binary strings (one per molecule). Partial least squares regression applied to these bit strings helps to identify which features at which locations lead to an increase or decrease in activity.

3D QSAR presents challenges even beyond traditional pharmacophore elucidation methods, though in some cases they can clearly provide additional (and useful) information.

In particular, the ability to incorporate inactive molecules is very appealing, though this should be done with great care; molecules can be inactive for a variety of reasons too subtle to be captured in the model, and they may only serve to add noise rather than signal. Any 3D QSAR protocol can be used to examine alignments, but bear in mind that pharmacophore models are underdetermined, and by use of powerful statistical tools, it is easy to find patterns in the data that are not physically real and therefore have no utility for finding novel active ligands.⁹⁶ As with a pharmacophore model, the output of a 3D QSAR protocol can be interpreted as a pseudo-receptor or as a more refined pharmacophore query. In either case, one should complete a cycle of experimental validation to see if the models find new hits or explain existing unseen experimental data rather than relying only on statistical validation criteria. There is also a paradox between the pharmacophore alignment and the way in which 3D QSAR models are set up. A pharmacophore will emphasize the points of commonality; therefore, the variation in the fields around the common features will be low so that these descriptor points can be removed from the QSAR model.

Pseudo-Receptors. A pseudo-receptor is a surrogate for the true binding site.⁹⁷ The pharmacophore alignment is the starting point, around which the pseudo-receptor is assembled, and optimized against the observed binding energies of the pharmacophore training set. As with all QSAR methods, a good spread of activity (3–4 orders of magnitude) is required. The pseudo-receptor can be built from grids, isosurfaces, Voronoi polyhedra, atomic shrink-wraps, or fragment packing. In the last two, atoms or fragments are placed to make the key putative interactions, followed by in-filling with mainly hydrophobic atoms or fragments. As the problem is greatly underdetermined, most researchers recommend that a family of receptors is generated. The shrink-wrapping algorithms tend to reduce the cavities that are usually observed experimentally in receptor–ligand complexes. The advantage of these methods is the correlation of fit to observed activity. This allows the results of a database search or a newly designed molecule to be more rigorously assessed. Pseudo-receptors have even been used to drive de novo design. However, such models are very speculative and hide two levels of ambiguity behind a seeming tangible receptor binding site.

Evaluation of Pharmacophore Methods

A lack of good data sets where the “true” pharmacophore is known has meant that the evaluation of new pharmacophore generation methods has tended to be less rigorous than programs in other areas such as protein–ligand docking. Very often, programs have been evaluated on a small number of test cases (sometimes only one), and few comparative studies of different programs have been carried out. While such a situation is natural as a field develops, pharmacophore generation has now reached a state of maturity where more rigorous evaluations are appropriate.

The earliest automated pharmacophore generation methods tended to be evaluated by comparing the generated pharmacophores to those previously published in the literature that had been hand-crafted and/or were amply supported by experimental data. An alternative approach has been to fit compounds that are known to be active but that were not included in model generation to the pharmacophore. Ideally the model should also explain inactive molecules as

well; for example, to determine whether the inactive compounds occupy regions outside the common volume of the pharmacophore.

The development of 3D database searching techniques made it possible to evaluate pharmacophore hypotheses in retrospective virtual screening experiments. A database of compounds can be spiked with compounds that share the same activity as those used to construct the pharmacophore. The database can then be either partitioned into those that match the pharmacophore and those that do not, or alternatively the compounds may be scored and ranked on their fit to the pharmacophore. Ideally, all active compounds will be retrieved or appear high in the ranked list. Database searching can also be used to choose between alternative hypotheses that may have been generated. Pharmacophore queries do, however, tend to produce large numbers of false positive hits, especially if used without additional features such as excluded volumes. There has been much discussion in the literature concerning the “proper” ways to perform such retrospective virtual screening evaluations.^{98–103} From a practical drug discovery perspective the most important consideration is whether the search provides novel chemotypes that provide a new direction for a project to explore.

When the pharmacophore generation method is linked to the building of a 3D QSAR, as is possible in the CATALYST and PHASE programs, the pharmacophores can be evaluated on the statistical quality of the resulting models. In many cases, models have been evaluated on their internal predictivities, but as with all QSAR methods, a more rigorous evaluation involves the assessment of an external test set. In a recent comparison of PHASE and CATALYST, eight series of compounds with known activities against different targets were divided into training and test sets. The compounds in the training sets were used to generate pharmacophore hypotheses that were then tested on their abilities to predict the activities of compounds in the test sets.¹⁰⁴ Acceptable models were found for only four of the data sets. Moreover, in some cases the programs were found to be highly sensitive to the chosen parameters, as has been found in other comparative studies.¹⁰⁵ Furthermore, the models with the highest predictivity did not correspond to the hypotheses that were scored highest by the programs, suggesting that such scores should be treated with caution.

The increased numbers and diversity of X-ray structures of protein–ligand complexes over the past decade have led to the development of test sets for the evaluation of protein–ligand docking programs. Performance is measured by the ability to reproduce the ligand binding modes. The initial test sets were of limited quality and diversity, but these are now of a much higher quality with the data carefully examined for accuracy.¹⁰⁶ The availability of a wide diversity of protein–ligand complexes also provides opportunities for the development of test sets for pharmacophore methods, although this is a more complex task than for docking, since each test case requires the identification of diverse ligands bound to the same protein target. After identification of such a set of ligands, it is then necessary to superimpose the complexes based on the active site residues in order to generate an alignment of the ligands. Finally, the common features that comprise the pharmacophore can be identified, typically through visual inspection. This was the approach taken by Patel et al. in their evaluation of the programs CATALYST, DISCO, and GASP,¹⁰⁷ which involved five protein targets (DHFR, thrombin, CDK2, thermolysin, and HIV reverse transcriptase) with up

to 10 ligands for each target. The complexes were visually analyzed to identify the “true” pharmacophore, and the three programs were evaluated on their ability to reproduce it. One outcome of this study was that the assessment of the generated pharmacophores is nontrivial. The ultimate goal is clear: the predicted 3D pharmacophore should consist of the X-ray conformation of each ligand superimposed exactly as seen in the crystal structure alignment. However, the definition of an acceptable solution is less clear. How should a pharmacophore be rated that contains the correct feature points (configuration) but the “wrong” conformations? What is an acceptable rmsd on the pharmacophoric distances (and conformations)? What if the variation in conformation is in a region of the ligand that is exposed to the solvent and not directly involved in the pharmacophore? What if most of the ligands are aligned correctly but one does not fit the pharmacophore? Multiple criteria were thus devised for the evaluation of the programs. The results, however, were rather disappointing, demonstrating several inadequacies of the existing programs and also the need for more rigorous evaluations.

Patel's data set was perhaps the first attempt to develop a standard test set for pharmacophore elucidation, and it has subsequently been used by other groups to evaluate new methods.^{20,17} As with the early attempts to develop test sets for docking, careful examination of these initial pharmacophore data sets by other authors revealed some errors (for example, in protonation states of the ligands) that were subsequently corrected. Thus, as for the docking test sets, it is important that careful consideration is given to issues such as tautomerism, protonation state, and fit to electron density when identifying the true pharmacophore. Furthermore, the small size of the Patel data set leads to a real risk of programs being overtrained on a few examples. There is a clear need for the further development of such data sets that are accepted within the community as the minimum required for the evaluation of any new pharmacophore elucidation program.

Exemplar Applications of 3D Pharmacophore Methods

Many applications of 3D pharmacophore methods can be found in the literature. In addition, several reviews and summaries of literature applications are available, often providing a useful practical perspective.^{14,108–111} Here, we highlight a very limited number of applications from the recent literature to illustrate some of the main ways in which 3D pharmacophore models are used in contemporary drug discovery.

A major use of 3D pharmacophore methods is for the identification of new hits and leads via 3D database searching. Such applications fall into two broad categories depending on whether the pharmacophore incorporates protein structural information or not. Wang and colleagues recently described the development of a cannabinoid CB1 receptor pharmacophore model from eight known active compounds.¹¹² The resulting pharmacophore contained five features. When used to search the Schering Plough compound collection of approximately 500 000 structures, almost 23 000 hits were obtained (significantly more than the available screening capacity). This is a common occurrence. The authors also describe in some detail how they subsequently reduced this large hit list down to a more manageable size (420 compounds) using a combination of filters, a Bayesian model of CB1

activity, and cluster analysis. A number of active compounds were identified, the most potent of which had an activity (K_i) of 53 nM. An example of hit identification using a structure-based pharmacophore comes from the work of Rella and colleagues using the crystal structure of the inhibitor MLN-4760 bound to angiotension converting enzyme 2 (ACE2).⁶⁷ This provided the 3D pharmacophore illustrated in Figure 4. The pharmacophore was validated and refined to ensure not only activity against ACE2 but also selectivity against ACE. A small number of compounds were purchased for biological testing in an ACE2 assay, resulting in some weak inhibitors from different chemical series. Our third example illustrates the use of 3D pharmacophore models in lead optimization and comes from work to identify nonsteroidal glucocorticoid (GR) agonists.¹¹³ Specifically, a 3D pharmacophore was derived from a docking model of an established series of GR agonists and used to design an array to identify replacements for a metabolically labile benzoxazinone moiety. A three-step procedure was used for the array design. First, a large set of potential building blocks was identified and combined with the core template to generate an enumerated library. The virtual products were filtered using modified Lipinski criteria. The remaining structures were then converted to a 3D database and searched using the pharmacophore query. This process resulted in 200 compounds for synthesis, some of which did indeed show activity in a GR binding assay.

Discussion and Future Directions

In this final section, we focus on a small number of issues that seem to us particularly important: the difficulty of aligning molecules and its consequences; the need to improve validation standards; limitations in the way features are typically represented; the growing emphasis on the development of algorithms that take account of the possibility of multiple binding modes; and our views on how pharmacophore methods are most likely to be of value in the next few years.

Challenge of Molecular Alignment. Sadly, it is almost impossible to predict with confidence the correct way of overlaying a set of ligands that bind to a protein of unknown 3D structure (except perhaps in trivial cases involving small numbers of ligands that share a large common substructure). A typical protein binding site contains many surface-exposed functional groups, and different ligands frequently interact with different selections of those groups. In addition, some of the interactions may be mediated by water molecules and the protein may be flexible. Overlaying ligands is especially difficult when the binding is largely driven by hydrophobic interactions, due to the weaker directional preferences than for hydrogen bonds. We discuss below how these problems are manifested in three different situations.

Alignment of Ligands for 3D QSAR. In principle, this is the easiest case, as the molecules are likely to be fairly close analogues. A study of protein–ligand complexes found that structurally similar ligands generally exhibit a high degree of structural conservation; changes are more likely to be found in water molecule architecture and/or side chain movements.¹¹⁴ Indeed, one might reasonably question whether it is necessary to use a computer program in such cases. Nevertheless, difficulties can arise because of local symmetry or approximate symmetry in the ligand structures. Consider the neuraminidase ligands in Figure 6 (from PDB structures 1a4g, 1a4q, 1b9t, 1inf). The obvious (and indeed correct) way

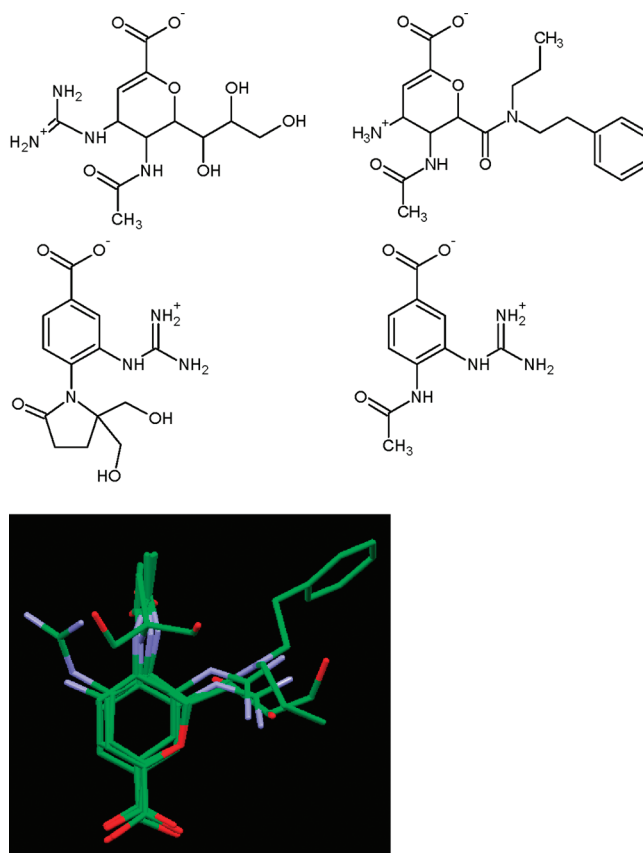


Figure 6. Four neuraminidase ligands superimposed according to their binding mode. The ring systems and the carboxylate group overlay; the guanidinium, and ammonium groups do not.

of aligning these is to superimpose the carboxylate groups and the ring atoms. However, each of the rings has two possible orientations (i.e., can be flipped to place the guanidinium or ammonium group on the “left” or “right” side). Consequently, there is a significant number of permutations from which to choose. Moreover, as is often the case, the correct answer is not obvious: two of the guanidinium groups lie on one side, the remaining guanidinium and the ammonium group on the other. The problem is exacerbated by the fact that substituents on synthetic ligands may have been introduced to moderate physical properties and may not form significant binding interactions with the protein.

Rationalization of the Shared Activity of Two Different Types of Ligands. A frequent requirement is to overlay molecules that fall into distinct structural types in order to understand their common activity against the same target protein. An example is the pair of dihydrofolate reductase ligands shown in Figure 7 (from PDB structures 1drf, 1mvt), the heterocyclic portions of which famously do not overlay in the way that maximizes steric complementarity but instead align to optimize the overlap of similar hydrogen-bonding functions. If the binding-site structure were unknown, there would be no reason to favor the correct heterocycle overlay over the incorrect; both would be perfectly feasible given that some of the hydrogen-bonding groups might be solvent exposed.

Analysis of Hits from a Random Screen. Overlaying a set of diverse molecules (e.g., hits from a random screen) is highly challenging. Consider the adenosine deaminase ligands in Figure 8 (from PDB structures 1krm, 1ndw, 1wxy, 1ndv).

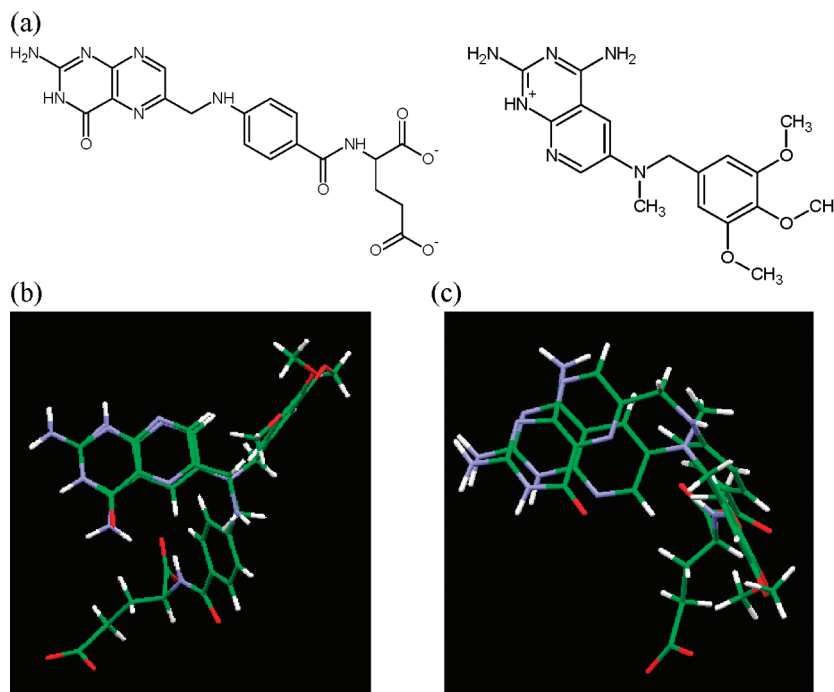


Figure 7. Heterocyclic moieties of the pair of dihydrofolate reductase ligands (a) do not overlay to maximize steric complementarity (b) but rather align to optimize the overlap of similar hydrogen-bonding functions (c).

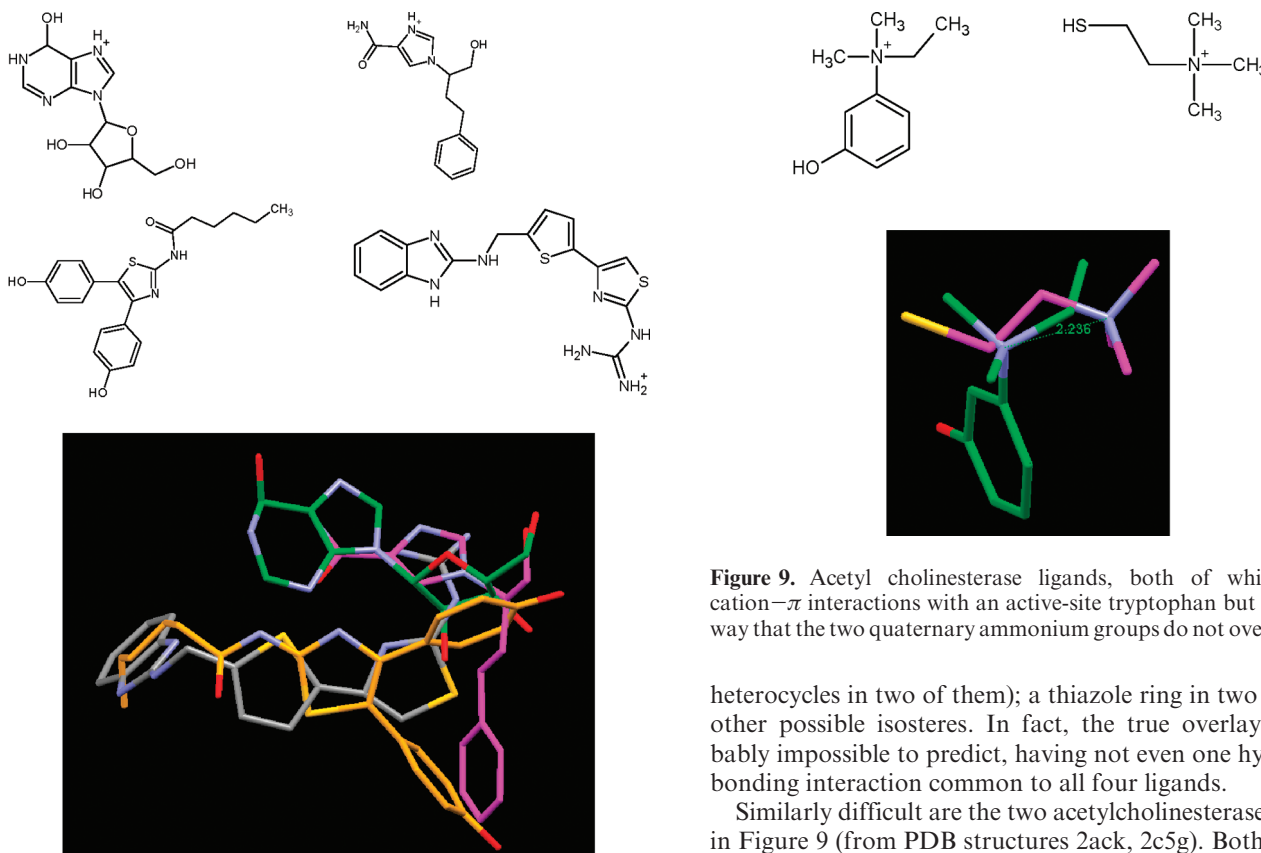


Figure 8. Overlay of four adenosine deaminase in their binding modes. Not even one hydrogen-bonding interaction is common to all four ligands.

There appear to be several clues that might point to the correct overlay: the presence of a positively charged donor in three of the molecules (and occurring in rather similar

Figure 9. Acetyl cholinesterase ligands, both of which form cation- π interactions with an active-site tryptophan but in such a way that the two quaternary ammonium groups do not overlay well.

heterocycles in two of them); a thiazole ring in two ligands; other possible isosteres. In fact, the true overlay is probably impossible to predict, having not even one hydrogen-bonding interaction common to all four ligands.

Similarly difficult are the two acetylcholinesterase ligands in Figure 9 (from PDB structures 2ack, 2c5g). Both ligands form cation- π interactions with an active-site tryptophan, but one interacts with the six-membered ring of the side chain and the other with the five-membered ring, so the two quaternary ammonium groups do not overlay well. Other ligands of this protein contain a huge variety of chemically diverse, electron-deficient hydrophobic groups that bind to the same tryptophan. Overlaying these groups is near impossible. (We are aware of the argument that the correct

pharmacophore can sometimes be predicted even if the correct overlay is not, but getting the right answer for the wrong reason makes us uneasy.)

In addition to the breakdown of the “single binding mode” assumption, it can also be shown that small changes to the parameters chosen for conformational analysis, feature definition, or even, sadly, the order in which ligands are considered and the version of software used, can lead to alternative, equally convincing pharmacophore models. In one study¹⁰⁵ several published data sets were taken, and attempts were made to reproduce the pharmacophore models, with only limited success. This highlights the dangers of using models without follow-up experimental validation.

Despite these challenges, pharmacophore methods can (and do) make significant contributions to the advancement of drug discovery projects. Our point is that any attempt to rank a set of overlays with the expectation that the top-ranked one will be correct is unduly optimistic, not because of deficiencies in the algorithm but because the problem is underdetermined. We are better off acknowledging this unfortunate reality and accepting that multiple hypotheses must be generated, regarded as equally likely, and challenged experimentally.

Validation Standards. Quite simply, validation standards are inadequate. In stark contrast to protein–ligand docking, where construction of ever-more exacting test sets is almost a cottage industry, pharmacophore elucidation programs are typically validated on pitifully small numbers of ligands, often with large common substructures. Sufficient protein–ligand series now exist in the PDB for this to be no longer excusable. Admittedly, test-set construction is a harder task than for docking, since the detailed protein–ligand interactions in each test complex must be identified and properly understood. Nevertheless, test sets of reasonable size can and should be achieved. Two of us (R.T., V.J.G.) are currently using such a set to test a pharmacophore elucidation program and anticipate that the set will be made freely available to the community in due course. Of course, the obvious limitation of a validation protocol based on protein–ligand structures is that it means that certain gene families of significant pharmaceutical interest (notably GPCRs, ion channels, and transporters) will not be represented, at least at present. The alternative approach, that of validation based on retrospective virtual screening experiments, is therefore very popular. As evidenced from the surfeit of papers on this topic, there continues to be much debate concerning the “best” way to perform such experiments, from construction of the data set to what performance metric to employ. One major drawback of such studies compared to validation based on 3D structure is that they are indirect; the right answer may be found for the wrong reasons. The challenge is that “real” drug discovery is not so easily quantified and that true success (i.e., that a compound identified by virtual screening leads to a development candidate, let alone a marketed drug) is extremely rare.

Pharmacophore Feature Representation. Feature representation has probably had less attention recently than it deserves. With the exception of field-based methods, most programs use features that differ little from those described over a decade ago by Greene et al.³⁰ While the seminal nature of that work is acknowledged, some questions can be posed. Some programs perform a surface-accessibility check before placing hydrophobic points, making the algorithm

conformation dependent and hence slower. Is the extra computation justified? Most programs offer positive and negative feature types. Is this necessary, or is the hydrogen-bond similarity method of Jones et al.¹⁸ a better alternative? How important is it to align aromatic ring normals, and if it is important, is it equally vital to achieve coplanarity of peptide and other groups that are hydrophobic in some directions but not in others (or is that achieved as a byproduct of matching the hydrogen-bonding features of such groups)? Is it sensible to generate alignments by least-squares superposition of virtual points, placed along lone-pair or donor-hydrogen directions at the “expected” positions of the complementary protein atoms? It may assume greater hydrogen-bond directionality than exists in practice, especially given that the protein may flex a little from one ligand complex to another. Also, the situation that virtual points aim to deal with (Figure 2) may be quite rare.

One important observation in protein–ligand docking is that greatly improved results can often be obtained if the user’s knowledge of the problem in hand can be exploited (for example, by allowing constraints to be set). The same tactic can be similarly advantageous in pharmacophore analysis. For example, most programs do not have a separate “metal-coordinator” feature type, effectively assuming that the hydrogen-bond acceptor feature will be a suitable proxy. This assumption can break down, as in the case of the thiol-containing ligands. Thiol is a very weak hydrogen-bond acceptor but a good coordinator of zinc, which normally induces it to deprotonate. A user, aware that the target protein contains zinc, will obtain better results by overriding the usual program behavior and ensuring that the importance of this group is taken properly into account.

One problem with ad hoc feature definition is that it causes difficulty at the searching stage, since databases are normally constructed with precomputed features and screens are based on those features. These are likely to be rendered useless if the query contains features assigned according to different rules. Nevertheless, on-the-fly feature assignment on large databases and searches of unscreened or poorly screened databases are not impossible. It is perhaps better to get a good answer slowly than a bad answer quickly.

Variable Binding Modes. One limitation of the pharmacophore elucidation programs in common use is the assumption of a common binding mode for the ligands in the data set. Some relaxation of this requirement may be permitted (for example, some compounds can match a subset of the common features), but data sets where the ligands correspond to two or more nonoverlapping pharmacophores will prove difficult or impossible. In addition, the nature of the algorithms can restrict the application of current methods to rather small data sets. Insofar as this may therefore result in a careful examination of the ligand structures, their activities, and mechanism of action, such limitations may be considered valuable. However, in an era when large volumes of data can be generated using automated screening technologies it is useful to have tools that not only can deal with larger sets but also can cope with the possibility of multiple binding modes. Some of the QSAR methods, particularly those based on pharmacophore keys, are able in part to accommodate such situations; more recent pharmacophore elucidation methods are also designed to deal with multiple binding modes while also retaining the link to the underlying three-dimensional models.¹¹⁵

The Future: Maximizing the Value of Pharmacophore Methods. We have alluded to a number of challenges in the course of this Perspective and have highlighted some areas for improvement. Nevertheless, it is undoubtedly true that the pharmacophore concept has proved to be extremely useful over the past 30 or so years. In part this is due to its simplicity, enabling the complexities of the molecular interaction between ligand and receptor (with attendant solvent, ions, cell membranes, etc.) to be reduced to a handful of features in a geometrical relationship. As such, it is perhaps surprising that it works at all. Pharmacophore methods are now commonly used as part of more complex workflows, for example, as a prefilter prior to protein–ligand docking. Such workflows will continue to be further refined, driven by the continued growth in the number of protein structures of pharmaceutical relevance. It should also not be forgotten that integral membrane proteins (GPCRs, ion channels, transporters) constitute many of the targets where pharmacophore analysis has been successfully applied for many years. Recent advances in X-ray structure determination for these families of proteins are providing new insights, but it is likely to take some years before we achieve sufficient breadth of coverage and depth of understanding (especially of the more complex mechanisms of action) that will truly enable structure-based design techniques. Until then, simpler models as embodied in the pharmacophore concept will continue to play a key role.

So far as lead discovery is concerned, why not just use HTS? There can be a naive expectation that HTS will return all known actives. However, HTS will miss some compounds for a number of good experimental reasons. It may also return too many false positives. Pharmacophore analysis is more hypothesis driven and is intended to target small regions of chemical space for novel chemical matter; this can be of particular value when trying to accelerate lead discovery (e.g., by screening a small focused set prior to running a HTS) or when the available assay has low throughput. More importantly, pharmacophore methods can be used to search virtual chemical space, both in the form of known but not yet executed chemistry and in the form of compounds unavailable in-house but (for example) accessible from a compound vendor.

It is now 4 decades since the modern-day concept of 3D pharmacophores was first introduced. It did not take long for scientists to appreciate the potential applications in drug discovery,¹¹⁶ since when there has been a continuous stream of publications describing its use in the selection and design of new chemical entities with interesting pharmacological properties, a trend that continues until the present day. It is truly a concept that has stood the test of time and one that is likely to play an important role in drug discovery for many more years to come.

Acknowledgment. We thank the following colleagues who have read and commented on earlier versions of this manuscript: Anna-Maria Capelli, Giovanna Tedesco, Colin Groom, and Rajeshri Karki.

Biographies

Andrew R. Leach received his B.A. (1986) and D.Phil. (1989) from Oxford University, U.K., under the direction of Keith Prout in the field of computational approaches to conformational analysis. Following postdoctoral studies as a NATO fellow in protein–ligand interactions and macromolecular

simulations working with Bob Langridge, Tack Kuntz, and Peter Kollman at the University of California, San Francisco, he joined Southampton University in 1991 as an EPSRC Advanced Fellow. In 1994 he moved to Glaxo Group Research and is currently a Director in Computational and Structural Chemistry at GlaxoSmithKline. Dr. Leach is an Editor-in-Chief for the *Journal of Computer-Aided Molecular Design* and has published over 60 papers and 3 books. He has long-standing interests in computational chemistry, cheminformatics, and scientific education in these fields.

Valerie J. Gillet is Professor of Chemoinformatics at The University of Sheffield, U.K. Her research interests include applications of evolutionary algorithms to problems in chemoinformatics, including de novo design, library design, structure–activity relationships, and pharmacophore elucidation. She has authored over 90 research papers and has collaborated extensively with industry. She serves on the Editorial Advisory Board of *Journal of Chemical Information and Modeling* and regularly reviews manuscripts for this and related journals including *Journal of Medicinal Chemistry*. She is program coordinator of the first Masters program in chemoinformatics worldwide, coauthor of the textbook *An Introduction to Chemoinformatics*, and organizer of the Sheffield triennial conference in chemoinformatics and the annual short course to industry, *A Practical Introduction to Chemoinformatics*, both of which attract delegates from around the globe.

Richard A. Lewis earned his B.A. (1984) and Ph.D. (1988) from Cambridge University, U.K., under the supervision of Dr. Philip Dean. He performed postdoctoral studies as a Fulbright Scholar in the laboratory of Prof. Tack Kuntz and as a Royal Commission Fellow in the laboratory of Prof. Mike Sternberg before joining Rhone-Poulenc Rorer in 1991. In 1998 he joined Eli Lilly, and then in 2004 he moved to Novartis as Head of Computer-Aided Drug Discovery. He has authored over 50 papers and patents and is a fellow of the Royal Society of Chemistry.

Robin Taylor obtained B.A. (1973) and Ph.D. (1976) degrees from Oxford and Cambridge Universities, U.K., respectively. After spells at Westminster Medical School and York and Pittsburgh Universities, he joined the Cambridge Crystallographic Data Centre (CCDC) in 1980, where he did research into hydrogen bonding and the use of crystallographic databases. He moved to ICI Agrochemicals in 1985, becoming head of their molecular design team in 1989. In 1994, he rejoined CCDC to head the company's software development team, working on programs such as GOLD, ConQuest, and Relibase. In 2008, he left CCDC to become a self-employed software developer trading under the name Taylor Cheminformatics Software.

References

- (1) Wermuth, C. G.; Ganellin, C. R.; Lindberg, P.; Mitscher, L. A. Glossary of terms used in medicinal chemistry (IUPAC Recommendations 1998). *Pure Appl. Chem.* **1998**, *70*, 1129–1143.
- (2) Allen, F. H. The Cambridge Structural Database: a quarter of a million crystal structures and rising. *Acta Crystallogr.* **2002**, *B58*, 380–388.
- (3) Marshall, G. R.; Barry, C. D.; Bosshard, H. E.; Dammkoehler, R. A.; Dunn, D. A. The Conformational Parameter in Drug Design: The Active Analog Approach. In *Computer-Assisted Drug Design*; Olson, E. C., Christoffersen, R. E., Eds.; American Chemical Society: Columbus, OH, 1979; pp 205–226.
- (4) Mayer, D.; Naylor, C. B.; Motoc, I.; Marshall, G. R. A unique geometry of the active site of angiotensin-converting enzyme. *J. Comput.-Aided Mol. Des.* **1987**, *1*, 3–16.
- (5) Van Drie, J. H. Monty Kier and the origin of the pharmacophore concept. *Internet Electron. J. Mol. Des.* **2007**, *6*, 271–279.
- (6) Kier, L. B. Molecular orbital calculation of preferred conformations of acetylcholine, muscarine, and muscarone. *Mol. Pharmacol.* **1967**, *3*, 487–494.
- (7) Kier, L. B. Receptor Mapping Using MO Theory. In *Fundamental Concepts in Drug–Receptor Interactions*; Danielli, J. F., Moran, J. F., Triggler, D. J., Eds.; Academic Press: New York, 1970.

- (8) Kier, L. B. *MO Theory in Drug Research*; Academic Press: New York, 1971; pp 164–169.
- (9) Guner, O. F., Ed. *Pharmacophore Perception, Development and Use in Drug Design*; International University Line: La Jolla, CA, 2000.
- (10) Mason, J. S.; Good, A. C.; Martin, E. J. 3-D pharmacophores in drug discovery. *Curr. Pharm. Des.* **2001**, *7*, 567–597.
- (11) Van Drie, J. H. Pharmacophore discovery. Lessons learned. *Curr. Pharm. Des.* **2003**, *9*, 1649–1664.
- (12) Langer, T.; Hoffmann, R. D., Eds. *Pharmacophores and Pharmacophore Searches*; Wiley-VCH: Weinheim, Germany, 2006.
- (13) Martin, Y. C. Pharmacophore Modeling: 1. Methods. In *Comprehensive Medicinal Chemistry II*; Taylor, J. B., Triggle, D. J., Mason, J. S., Eds.; Computer-Assisted Drug Design, Vol. 4; Elsevier: Amsterdam, 2007; pp 119–147.
- (14) Martin, Y. C. Pharmacophore Modeling: 2. Applications. In *Comprehensive Medicinal Chemistry II*; Taylor, J. B., Triggle, D. J., Mason, J. S., Eds.; Computer-Assisted Drug Design, Vol. 4; Elsevier: Amsterdam, 2007; pp 515–536.
- (15) Hamprecht, D.; Micheli, F.; Tedesco, G.; Checchia, A.; Donati, D.; Petrone, M.; Terreni, S.; Wood, M. Isoindolone derivatives, a new class of 5-HT_{2C} antagonists: synthesis and biological evaluation. *Bioorg. Med. Chem. Lett.* **2007**, *17*, 428–433.
- (16) Barnum, D.; Greene, J.; Smellie, A.; Sprague, P. Identification of common functional configurations among molecules. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 563–571.
- (17) Richmond, N. J.; Abrams, C. A.; Wolohan, P. R. N.; Abrahamian, E.; Willett, P.; Clark, R. D. GALAHAD: 1. Pharmacophore identification by hypermolecular alignment of ligands in 3D. *J. Comput.-Aided Mol. Des.* **2006**, *20*, 567–587.
- (18) Jones, G.; Willett, P.; Glen, R. C. A genetic algorithm for flexible molecular overlay and pharmacophore elucidation. *J. Comput.-Aided Mol. Des.* **1995**, *9*, 532–549.
- (19) *Molecular Operating Environment*; Chemical Computing Group: Montreal, Canada; <http://www.chemcomp.com/>.
- (20) Dixon, S. L.; Smondyrev, A. M.; Knoll, E. H.; Rao, S. N.; Shaw, D. E.; Friesner, R. A. PHASE: a new engine for pharmacophore perception, 3D QSAR model development, and 3D database screening: 1. Methodology and preliminary results. *J. Comput.-Aided Mol. Des.* **2006**, *20*, 647–671.
- (21) Anghelescu, A. V.; DeLisle, R. K.; Lowrie, J. F.; Klon, A. E.; Xie, X.; Diller, D. J. Technique for generating three-dimensional alignments of multiple ligands from one-dimensional alignments. *J. Chem. Inf. Model.* **2008**, *48*, 1041–1054.
- (22) Cho, S. J.; Sun, Y. FLAME: a program to flexibly align molecules. *J. Chem. Inf. Model.* **2006**, *46*, 298–306.
- (23) Feng, J.; Sanil, A.; Young, S. S. PharmID: pharmacophore identification using Gibbs sampling. *J. Chem. Inf. Model.* **2006**, *46*, 1352–1359.
- (24) Marialke, J.; Koerner, R.; Tietze, S.; Apostolakis, J. Graph-based molecular alignment (GMA). *J. Chem. Inf. Model.* **2007**, *47*, 591–601.
- (25) Podolyan, Y.; Karypis, G. Common pharmacophore identification using frequent clique detection algorithm. *J. Chem. Inf. Model.* **2009**, *49*, 13–21.
- (26) Schneidman-Duhovny, D.; Dror, O.; Inbar, Y.; Nussinov, R.; Wolfson, H. J. Deterministic pharmacophore detection via multiple flexible alignment of drug-like molecules. *J. Comput. Biol.* **2008**, *15*, 737–754.
- (27) Todorov, N. P.; Alberts, I. L.; de Esch, I. J. P.; Dean, P. M. QUASI: a novel method for simultaneous superposition of multiple flexible ligands and virtual screening using partial similarity. *J. Chem. Inf. Model.* **2007**, *47*, 1007–1020.
- (28) Zhu, F.; Agrafiotis, D. K. Recursive distance partitioning algorithm for common pharmacophore identification. *J. Chem. Inf. Model.* **2007**, *47* (4), 1619–1625.
- (29) Taminau, J.; Thijss, G.; De Winter, H. Pharao: pharmacophore alignment and optimization. *J. Mol. Graphics Modell.* **2008**, *27*, 161–169.
- (30) Greene, J.; Kahn, S.; Savoj, H.; Sprague, P.; Teig, S. Chemical function queries for 3D database search. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 1297–1308.
- (31) Wolber, G.; Seidel, T.; Bendix, F.; Langer, T. Molecule-pharmacophore superpositioning and pattern matching in computational drug design. *Drug Discovery Today* **2008**, *13*, 23–29.
- (32) Pierce, A. C.; Sandretto, K. L.; Bemis, G. W. Kinase inhibitors and the case for CH...O hydrogen bonds in protein–ligand binding. *Proteins* **2002**, *49*, 567–576.
- (33) Oellien, F.; Cramer, J.; Beyer, C.; Ihlenfeldt, W.-D.; Selzer, P. M. The impact of tautomer forms on pharmacophore-based virtual screening. *J. Chem. Inf. Comput. Sci.* **2006**, *46*, 2342–2354.
- (34) Milletti, F.; Storchi, L.; Sforza, G.; Cross, S.; Cruciani, G. Tautomer enumeration and stability prediction for virtual screening on large chemical databases. *J. Chem. Inf. Comput. Sci.* **2009**, *49*, 68–75.
- (35) *SMARTS. Language for Describing Molecular Patterns*; Daylight Chemical Information Systems, Inc.: Aliso Viejo, CA; <http://www.daylight.com/>.
- (36) Boehm, H.-J.; Brode, S.; Hesse, U.; Klebe, G. Oxygen and nitrogen in competitive situations: Which is the hydrogen-bond acceptor? *Chem.—Eur. J.* **1996**, *2*, 1509–1513.
- (37) Chau, P.; Dean, P. Electrostatic complementarity between proteins and ligands. 2. Ligand moieties. *J. Comput.-Aided Mol. Des.* **1994**, *8*, 527–544.
- (38) Abraham, M. H.; Ibrahim, A.; Zissimos, A. M.; Zhao, Y. H.; Comer, J.; Reynolds, D. P. Application of hydrogen bonding calculations in property based drug design. *Drug Discovery Today* **2002**, *7*, 1056–1063.
- (39) Laurence, C.; Brameld, K. A.; Graton, J.; Le Questel, J.-Y.; Renault, E. The pK_{BHX} database: towards a better understanding of hydrogen-bond basicity for medicinal chemists. *J. Med. Chem.* **2009**, *52*, 4073–4086.
- (40) Cheeseright, T. J.; Mackey, M. D.; Melville, J. L.; Vinter, J. G. FieldScreen: virtual screening using molecular fields. Application to the DUD data set. *J. Chem. Inf. Model.* **2008**, *48*, 2108–2117.
- (41) Voth, A. R.; Khoo, P.; Oishi, K.; Ho, P. S. Halogen bonds as orthogonal molecular interactions to hydrogen bonds. *Nat. Chem.* **2009**, *1*, 74–79.
- (42) Lu, Y.; Shi, T.; Wang, Y.; Yang, H.; Yan, X.; Luo, X.; Jiang, H.; Zhu, W. Halogen bonding. A novel interaction for rational drug design? *J. Med. Chem.* **2009**, *52*, 2854–2862.
- (43) Wolber, G.; Dornhofer, A. A.; Langer, T. Efficient overlay of small organic molecules using 3D pharmacophores. *J. Comput.-Aided Mol. Des.* **2006**, *20*, 773–788.
- (44) Bostrom, J.; Norrby, P.-O.; Liljefors, T. Conformational energy penalties of protein-bound ligands. *J. Comput.-Aided Mol. Des.* **1998**, *12*, 383–396.
- (45) Bostrom, J. Reproducing the conformations of protein-bound ligands: a critical evaluation of several popular conformational searching tools. *J. Comput.-Aided Mol. Des.* **2001**, *15*, 1137–1152.
- (46) Perola, E.; Charifson, P. S. Conformational analysis of drug-like molecules bound to proteins: an extensive study of ligand reorganization upon binding. *J. Med. Chem.* **2004**, *47*, 2499–2510.
- (47) Agrafiotis, D. K.; Gibbs, A. C.; Zhu, F.; Izrailev, S.; Martin, E. Conformational sampling of bioactive molecules: a comparative study. *J. Chem. Inf. Model.* **2007**, *47*, 1067–1086.
- (48) Chen, I.-J.; Foloppe, N. Conformational sampling of druglike molecules with MOE and Catalyst: implications for pharmacophore modeling and virtual screening. *J. Chem. Inf. Model.* **2008**, *48*, 1773–91.
- (49) Martin, Y. C.; Bures, M. G.; Danaher, A. A.; DeLazzer, J.; Lico, I.; Pavlik, P. A. A fast new approach to pharmacophore mapping and its application to dopaminergic and benzodiazepine agonists. *J. Comput.-Aided Mol. Des.* **1993**, *7*, 83–102.
- (50) Martin, Y. C. DISCO: What We Did Right and What We Missed. In *Pharmacophore Perception, Development and Use in Drug Design*; Guner, O. F., Ed.; International University Line: La Jolla, CA, 2000; pp 51–66.
- (51) Labute, P.; William, C.; Feher, M.; Sourial, E.; Schmidt, J. M. Flexible alignment of small molecules. *J. Med. Chem.* **2001**, *44*, 1483–1490.
- (52) Cottrell, S. J.; Gillet, V. J.; Taylor, R.; Wilton, D. J. Generation of multiple pharmacophore hypotheses using multiobjective optimisation techniques. *J. Comput.-Aided Mol. Des.* **2004**, *18*, 665–682.
- (53) Martin, Y. C. 3D database searching in drug design. *J. Med. Chem.* **1992**, *35*, 2145–2154.
- (54) Manallack, D. T. Getting that hit: 3D database searching in drug discovery. *Drug Discovery Today* **1996**, *1*, 231–238.
- (55) Clark, D. E.; Westhead, D. R.; Sykes, R. A.; Murray, C. W. Active-site-directed 3D database searching: pharmacophore extraction and validation of hits. *J. Comput.-Aided Mol. Des.* **1996**, *10*, 397–416.
- (56) Good, A. C.; Mason, J. S. Three-Dimensional Structure Database Searches. In *Reviews in Computational Chemistry*; Lipkowitz, K. B., Boyd, D. B., Eds.; VCH Publishers: New York, 1996; Vol. 7, pp 67–117.
- (57) UNITY, Tripos. <http://www.tripos.com/>
- (58) Hurst, T. Flexible 3D Searching: The Directed Tweak Technique. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 190–196.
- (59) Smellie, A. S.; Kahn, S. D.; Teig, S. L. Analysis of conformational coverage. 1. Validation and estimation of coverage. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 285–294.

- (60) Smellie, A. S.; Teig, S. L.; Towbin, P. Poling: promoting conformational variation. *J. Comput. Chem.* **1995**, *16*, 171–187.
- (61) Mason, J. S.; Morize, I.; Menard, P. R.; Cheney, D. L.; Hulme, C.; Labaudiniere, R. F. New 4-point pharmacophore method for molecular similarity and diversity applications: overview of the method and applications, including a novel approach to the design of combinatorial libraries containing privileged substructures. *J. Med. Chem.* **1999**, *42*, 3251–3264.
- (62) Good, A. C.; Lewis, R. A. New methodology for profiling combinatorial libraries and screening sets: cleaning up the design process with HARPick. *J. Med. Chem.* **1997**, *40*, 3926–3936.
- (63) McGregor, M. J.; Muskal, S. M. Pharmacophore fingerprinting. 1. Application to QSAR and focused library design. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 569–574.
- (64) DeAnda, F.; Stewart, E. L. Application of the PharmPrint methodology to two protein kinases. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1803–1809.
- (65) Greenidge, P. A.; Carlsson, B.; Bladh, L.-G.; Gillner, M. Pharmacophores Incorporating numerous excluded volumes defined by X-ray crystallographic structure in three-dimensional database searching: application to the thyroid hormone receptor. *J. Med. Chem.* **1998**, *41*, 2503–2512.
- (66) Tintori, C.; Corradi, V.; Magnani, M.; Manetti, F.; Cotta, M. Targets looking for drugs: a multistep computational protocol for the development of structure-based pharmacophores and their applications for hit discovery. *J. Chem. Inf. Model.* **2008**, *48*, 2166–2179.
- (67) Rella, M.; Rushworth, C. A.; Guy, J. L.; Turner, A. J.; Langer, T.; Jackson, R. M. Structure-based pharmacophore design and virtual screening for novel angiotensin converting enzyme 2 inhibitors. *J. Chem. Inf. Model.* **2006**, *46* (2), 708–716.
- (68) Grant, J. A.; Gallardo, M. A.; Pickup, B. T. A fast method of molecular shape comparison. A simple application of a Gaussian description of molecular shape. *J. Comput. Chem.* **1996**, *17*, 1653–1666.
- (69) Hahn, M. Three-dimensional shape-based searching of conformationally flexible compounds. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 80–86.
- (70) Lemmen, C.; Lengauer, T.; Klebe, G. FlexS: a method for fast flexible ligand superposition. *J. Med. Chem.* **1998**, *41*, 4502–4520.
- (71) Goodford, P. J. A computational-procedure for determining energetically favorable binding-sites on biologically important macromolecules. *J. Med. Chem.* **1985**, *28* (7), 849–857.
- (72) Bohm, H. J. The computer program LUDI: a new method for the de novo design of enzyme inhibitors. *J. Comput.-Aided Mol. Des.* **1992**, *6*, 61–78.
- (73) Verdonk, M. L.; Cole, J. C.; Taylor, R. SuperStar: a knowledge-based approach for identifying interaction sites in proteins. *J. Mol. Biol.* **1999**, *289*, 1093–1108.
- (74) Venkatachalam, C. M.; Kirchoff, P.; Waldman, M. Receptor-Based Pharmacophore Perception and Modeling. In *Pharmacophore Perception, Development and Use in Drug Design*; Guner, O. F., Ed.; IUL Biotechnology Series: La Jolla, CA, 2000; pp 339–350.
- (75) Wolber, G.; Langer, T. LigandScout: 3-D pharmacophores derived from protein-bound ligands and their use as virtual screening filters. *J. Chem. Inf. Model.* **2005**, *45*, 160–169.
- (76) Baroni, M.; Cruciani, G.; Sciabola, S.; Perruccio, F.; Mason, J. S. A common reference framework for analyzing/comparing proteins and ligands. fingerprints for ligands and proteins (FLAP): theory and application. *J. Chem. Inf. Model.* **2007**, *47*, 279–294.
- (77) Cruciani, G.; Carosati, E.; De Boeck, B.; Ethirajulu, K.; Mackie, C.; Howe, T.; Vianello, R. MetaSite: understanding metabolism in human cytochromes from the perspective of the chemist. *J. Med. Chem.* **2005**, *48*, 6970–6979.
- (78) Kuhn, D.; Weskamp, N.; Schmitt, S.; Hullermeier, E.; Klebe, G. From the similarity analysis of protein cavities to the functional classification of protein families using Cavbase. *J. Mol. Biol.* **2006**, *359* (4), 1023–1044.
- (79) Shatsky, M.; Shulman-Peleg, A.; Nussinov, R.; Wolfson, H. J. The multiple common point set problem and its application to molecule binding pattern detection. *J. Comput. Biol.* **2006**, *13*, 407–428.
- (80) Deng, Z. D.; Chuaqui, C.; Singh, J. Structural interaction fingerprint (SIFt): a novel method for analyzing three-dimensional protein–ligand binding interactions. *J. Med. Chem.* **2004**, *47*, 337–344.
- (81) Brewerton, S. C. The use of protein–ligand interaction fingerprints in docking. *Curr. Opin. Drug Discovery Dev.* **2008**, *11*, 356–364.
- (82) Chuaqui, C.; Deng, Z.; Singh, J. Interaction profiles of protein kinase–inhibitor complexes and their application to virtual screening. *J. Med. Chem.* **2005**, *48*, 121–133.
- (83) Kelly, M. D.; Mancera, R. L. Expanded interaction fingerprint method for analyzing ligand binding modes in docking and structure-based drug design. *J. Chem. Inf. Comput. Sci.* **2004**, *44*, 1942–1951.
- (84) Cross, S.; Cruciani, G. Molecular fields in drug discovery: getting old or reaching maturity? *Drug Discovery Today* **2009**, 1–10.
- (85) Good, A. C.; Hodgkin, E. E.; Richards, W. G. Utilization of Gaussian functions for the rapid evaluation of molecular similarity. *J. Chem. Inf. Comput. Sci.* **1992**, *32*, 188–191.
- (86) McMahon, A. J.; King, P. M. Optimization of carbo molecular similarity index using gradient methods. *J. Comput. Chem.* **1997**, *18*, 151–158.
- (87) Wild, D. J.; Willett, P. Similarity searching in files of three-dimensional chemical structures. Alignment of molecular electrostatic potential fields with a genetic algorithm. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 159–167.
- (88) Nissink, J. W. M.; Verdonk, M. L.; Klebe, G. Simple knowledge-based descriptors to predict protein–ligand interactions. Methodology and validation. *J. Comput.-Aided Mol. Des.* **2000**, *14*, 787–803.
- (89) Pastor, M.; Cruciani, G.; McLay, I.; Pickett, S.; Clementi, S. GRIND-INdependent descriptors (GRIND): a novel class of alignment-independent three-dimensional molecular descriptors. *J. Med. Chem.* **2000**, *43*, 3233–3243.
- (90) Durán, A.; Martínez, G. C.; Pastor, M. Development and validation of AMANDA, a new algorithm for selecting highly relevant regions in molecular interaction fields. *J. Chem. Inf. Model.* **2008**, *48*, 1813–1823.
- (91) Vinter, J. G. Extended electron distributions applied to the molecular mechanics of some intermolecular interactions. *J. Comput.-Aided Mol. Des.* **1994**, *8*, 653–668.
- (92) Cheeseright, T.; Mackey, M.; Rose, S.; Vinter, A. Molecular field extrema as descriptors of biological activity: definition and validation. *J. Chem. Inf. Model.* **2006**, *46*, 665–676.
- (93) Thorner, D. A.; Wild, D. J.; Willett, P.; Wright, P. M. Similarity searching in files of three-dimensional chemical structures: flexible field-based searching of molecular electrostatic potentials. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 900–908.
- (94) Jewell, N. E.; Turner, D. B.; Willett, P.; Sexton, G. J. Automatic generation of alignments for 3D QSAR analyses. *J. Mol. Graphics Modell.* **2001**, *20*, 111–121.
- (95) Cheeseright, T.; Mackey, M.; Rose, S.; Vinter, A. Molecular field technology applied to virtual screening and finding the bioactive conformation. *Expert Opin. Drug Discovery* **2007**, *2*, 131–144.
- (96) Doweiko, A. M. QSAR: dead or alive? *J. Comput.-Aided Mol. Des.* **2008**, *22*, 81–89.
- (97) Tanrikulu, Y.; Schneider, G. Pseudoreceptor models in drug design: bridging ligand- and receptor-based virtual screening. *Nat. Rev. Drug Discovery* **2008**, *7*, 667–677.
- (98) Jain, A. N.; Nicholls, A. Recommendations for evaluation of computational methods. *J. Comput.-Aided Mol. Des.* **2008**, *22*, 133–139.
- (99) Hawkins, P. C. D.; Warren, G. L.; Skillman, A. G.; Nicholls, A. How to do an evaluation: pitfalls and traps. *J. Comput.-Aided Mol. Des.* **2008**, *22*, 179–190.
- (100) Good, A. C.; Oprea, T. O. Optimization of CAMD techniques 3. Virtual screening enrichment studies: a help or hindrance in tool selection? *J. Comput.-Aided Mol. Des.* **2008**, *22*, 169–178.
- (101) Truchon, J.-F.; Bayly, C. I. Evaluating virtual screening methods: good and bad metrics for the early recognition problem. *J. Chem. Inf. Model.* **2007**, *47*, 488–508.
- (102) Good, A. C.; Hermsmeier, M. A.; Hindle, S. A. Measuring CAMD technique performance: a virtual screening case study in the design of validation experiments. *J. Comput.-Aided Mol. Des.* **2004**, *18*, 529–536.
- (103) Mackey, M. D.; Melville, J. L. Better than random? The chemotype enrichment problem. *J. Chem. Inf. Model.* **2009**, *49*, 1154–1162.
- (104) Evans, D. A.; Doman, T. N.; Thorner, D. A.; Bodkin, M. J. 3D QSAR methods: phase and catalyst compared. *J. Chem. Inf. Model.* **2007**, *47*, 1248–1257.
- (105) Kristam, R.; Gillet, V. J.; Lewis, R. A.; Thorner, D. Comparison of conformational analysis techniques to generate pharmacophore hypotheses using Catalyst. *J. Chem. Inf. Model.* **2005**, *45*, 461–476.
- (106) Hartshorn, M. J.; Verdonk, M. L.; Chessari, G.; Brewerton, S. C.; Mooij, W. T.; Mortenson, P. N.; Murray, C. W. Diverse, high-quality test set for the validation of protein–ligand docking performance. *J. Med. Chem.* **2007**, *50*, 726–741.
- (107) Patel, Y.; Gillet, V. J.; Bravi, G.; Leach, A. R. A comparison of the pharmacophore identification programs: Catalyst, DISCO and GASP. *J. Comput.-Aided Mol. Des.* **2002**, *16*, 653–681.

- (108) Schneider, G.; Schneider, P.; Renner, S. Scaffold-hopping: how far can you jump? *QSAR Comb. Sci.* **2006**, *12*, 1162–1171.
- (109) Clark, D. E. What has virtual screening ever done for drug discovery? *Expert Opin. Drug Discovery* **2008**, *3*, 841–851.
- (110) Green, D. V. S. Virtual screening of chemical libraries for drug discovery. *Expert Opin. Drug Discovery* **2008**, *3*, 1011–1026.
- (111) Sun, H. Pharmacophore-based virtual screening. *Curr. Med. Chem.* **2008**, *15*, 1018–1024.
- (112) Wang, H.; Duffy, R. A.; Boykow, G. C.; Chackalamannil, S.; Madison, V. Identification of novel cannabinoid CB1 receptor antagonists by using virtual screening with a pharmacophore model. *J. Med. Chem.* **2008**, *51*, 2439–2446.
- (113) Clackers, M.; Coe, D.; Demaine, D. A.; Hardy, G. W.; Humphreys, D.; Inglis, G. G. A.; Johnston, M. J.; Jones, H. T.; House, D.; Loiseau, R.; Minick, D. J.; Skone, P. A.; Uings, I.; McLay, I. M.; Macdonald, S. J. F. Non-steroidal glucocorticoid agonists. The discovery of aryl pyrazoles as A-ring mimetics. *Bioorg. Med. Chem. Lett.* **2007**, *17*, 4737–4745.
- (114) Bostrom, J.; Hogner, A.; Schmitt, S. Do structurally similar ligands bind in a similar fashion? *J. Med. Chem.* **2006**, *49*, 6716–6725.
- (115) Feng, J.; Sanil, A.; Young, S. S. PharmID: pharmacophore identification using Gibbs sampling. *J. Chem. Inf. Model.* **2006**, *46*, 1352–1359.
- (116) Martin, Y. C.; Jarboe, C. H.; Krause, R. A.; Lynn, K. R.; Dunnigan, D.; Holland, J. B. Potential anti-Parkinson drugs designed by receptor mapping. *J. Med. Chem.* **1973**, *16*, 147–150.