

This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at SciVerse ScienceDirect

Journal of Biotechnology

journal homepage: [www.elsevier.com/locate/jbiotec](http://www.elsevier.com/locate/jbiotec)

# Structure-based computational analysis of protein binding sites for function and druggability prediction

Britta Nisius, Fan Sha, Holger Gohlke\*

Department of Mathematics and Natural Sciences, Institute of Pharmaceutical and Medicinal Chemistry, Heinrich-Heine University Düsseldorf, Germany

## ARTICLE INFO

### Article history:

Received 2 October 2011  
Received in revised form 2 December 2011  
Accepted 6 December 2011  
Available online 14 December 2011

### Keywords:

Computational analysis  
Druggability  
Protein binding sites  
Protein function prediction  
Protein structure

## ABSTRACT

Protein binding sites are the places where molecular interactions occur. Thus, the analysis of protein binding sites is of crucial importance to understand the biological processes proteins are involved in. Herein, we focus on the computational analysis of protein binding sites and present structure-based methods that enable function prediction for orphan proteins and prediction of target druggability. We present the general ideas behind these methods, with a special emphasis on the scopes and limitations of these methods and their validation. Additionally, we present some successful applications of computational binding site analysis to emphasize the practical importance of these methods for biotechnology/bioeconomy and drug discovery.

© 2011 Elsevier B.V. All rights reserved.

## 1. Introduction

The interactions of a protein with other molecules, e.g., ligands, nucleic acids, or other proteins, are critical to its biochemical function. Usually, not all residues on a protein's surface participate in these interactions; rather the interactions occur at defined locations, the protein binding sites. Thus, the identification and characterization of these protein binding sites is crucial to understand molecular interactions and recognition. The binding of a molecule to a protein's binding site depends on their physico-chemical and shape complementarity. Size, buriedness, and flexibility of the binding site are additional key factors for molecular recognition. The role of these key factors will be described in more detail below in this section.

Due to the importance of protein binding sites in molecular recognition and interactions, various approaches aiming at the structure-based computational *binding site analysis* (BSA) have been developed in recent years. In this review, we initially focus on BSA methods to perform *binding site comparison* (BSC). By applying such methods, one can detect binding sites in a set of protein structures that are similar to a given binding site. From a biotechnological/bioeconomical point of view, this characterization of binding pockets allows de-orphanization of (biochemical)

protein function by comparing binding pockets of multiple proteins and inferring the function of the orphan protein from the most similar protein(s). This is valuable for the identification of novel enzymes, which can subsequently be used for biocatalytic compound transformation, aiming at more sustainable production pathways (see Section 2). From a pharmaceutical–medicinal chemistry point of view, these methods also help to rationalize and predict cross-target drug interactions and toxicity (see Section 2). Another important question in drug discovery relates to whether a protein binding site is amenable to binding drug-like molecules, i.e., “*druggability*” prediction (DP). BSA methods that aim at answering this question are reviewed in Section 3. Rather than presenting algorithmic details of all available methods, we aim at presenting the general ideas of selected methods together with practical applications. Finally, we discuss the scopes and limitations of the presented methods in Section 4.

### 1.1. Diversity of protein binding sites

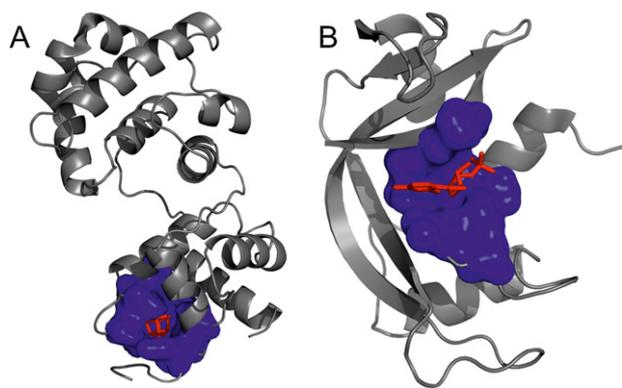
Since proteins are able to interact with a wide range of molecules, the binding sites involved in these interactions are diverse: The active site of an enzyme is often characterized by a particularly large and deep cleft, whereas protein–protein interfaces are usually flat and unstructured (Laskowski et al., 1996). Furthermore, even binding sites of enzymes can vary significantly. For instance, the binding site of endonuclease is a spherical cavity containing a deeply buried ligand, whereas the binding site of ribonuclease is an elongated groove containing a rather exposed ligand (Fig. 1).

Abbreviations: BSA, binding site analysis; BSC, binding site comparison; DP, druggability prediction; PPI, protein–protein interaction.

\* Corresponding author at: Universitätstraße 1, 40225 Düsseldorf, Germany.

Tel.: +49 211 81 13662; fax: +49 211 81 13847.

E-mail address: [gohlke@uni-duesseldorf.de](mailto:gohlke@uni-duesseldorf.de) (H. Gohlke).



**Fig. 1.** Diverse enzyme binding sites. The 3D structures of endonuclease (shown in A, PDB code 2ABK) and ribonuclease (shown in B, PDB 1ROB) and their ligand binding sites (highlighted in blue) are shown. Endonuclease has a spherical cavity with a deeply buried ligand, whereas ribonuclease has an elongated binding site and a rather exposed ligand.

For enzyme binding sites, the largest cavity is most often the active site (Laskowski et al., 1996; Liang et al., 1998). Furthermore, while larger proteins tend to have *more* binding sites, they do not necessarily have *larger* binding sites (Liang et al., 1998). Binding sites involved in protein–ligand interactions are often characterized by the presence of regions with very low and other regions with very high structural stability. Thus, these binding sites exhibit a dual character (Luque and Freire, 2000). The stable part of the binding site usually contains residues involved in interactions requiring a well-defined stereochemical arrangement, e.g., hydrogen bonds. By contrast, the flexible part of the binding site enables an adjustment of the binding pocket's shape to ensure an optimal fit and buriedness of the ligand and/or the accommodation of structurally different ligands. Furthermore, the low stability regions are proposed to play a crucial role in the transmission of information from an allosteric binding site to a catalytic binding site. Thus, shape and size of a ligand binding site are crucial determinants of its recognition power. However, geometrical complementarity alone is not sufficient to fully account for molecular recognition (Kahraman et al., 2007). Additionally, a physico-chemical complementarity is important. For instance, it was shown that some specific amino acids (Arg, His, Trp, and Tyr) occur substantially more frequently in protein binding sites than in the entire protein (Villar and Kauvar, 1994). Furthermore, the amino acid composition among protein binding sites can vary significantly, e.g., neuraminidase has a highly charged binding site whereas the binding site of avidin contains no charged residue (Hou et al., 2011).

Interfaces involved in protein–protein interactions (PPIs) are typically flat and significantly larger than protein–ligand binding sites. Additionally, these larger protein–protein interfaces are most often composed of multiple epitopes that are not sequentially connected. The epitopes can be divided into “functional epitopes”, which actually contribute to binding, and additional “structural epitopes” (Grimme et al., *in press*). In fact, mutagenesis studies revealed that only a small subset of all amino acids flanking the protein–protein interface significantly contribute to binding affinity. These residues are called “hotspots” (Bogan and Thorn, 1998). Furthermore, protein–protein interfaces show a significantly higher degree of inherent flexibility and plasticity than protein–ligand binding sites (Grimme et al., *in press*).

Thus, since proteins are involved in complex and diverse molecular interactions, a full characterization of protein binding sites requires a detailed analysis of the various factors contributing to molecular recognition. Up to now there is no standard definition of what constitutes a binding site, which represents a major complication in BSA (Perot et al., 2010). However, since the relative

importance of the key factors significantly varies for different binding sites, it is indeed very difficult to develop procedures that are generally applicable across diverse sites (Henrich et al., 2009).

### 1.2. Protein flexibility

Another complication in BSA arises from the flexibility of proteins, which enables a range of possible movements, from single side-chain rotations to drastic structural rearrangements (Ahmed et al., 2007; Cozzini et al., 2008). Thus, it is not always sufficient to use just one static structure for BSA: First, protein flexibility and plasticity can allow for the opening of novel binding sites (so-called “transient” or “cryptic” pockets (Eyrisch and Helms, 2007; Metz et al., *in press*)) that may not be detectable in the one single structure selected for analysis. Second, binding sites can also change their sizes and shapes upon binding. This is in line with the “conformational selection model” (Tsai et al., 1999), which proposes that, from various rapidly interconverting conformations of the unbound protein, that conformation is picked by a binding partner that has a binding site most complementary with the characteristics of the partner. Thus, based on this model, it has to be assumed that a binding site's shape and size strongly depends on the interacting ligand and, therefore, cannot be analyzed independently of the ligand (Ma et al., 2002). Hence, whenever knowledge about moving protein parts is available, it should be included in the analysis of protein binding sites. This knowledge can be gained from experimental information, e.g., multiple structures solved by crystallography or an ensemble of structures determined by NMR, as well as from computational approaches such as molecular dynamics simulations, graph-theoretical approaches, or normal mode analysis (Ahmed et al., 2007; Cozzini et al., 2008).

### 1.3. Computational approaches for binding pocket identification

Before protein binding sites can be analyzed by computational means, their location on the surface of a protein has to be identified. Unless a co-crystallized ligand readily provides this information, the detection of potential binding pockets is the first step in computational BSA. Many computational methods have been developed with that aim, which, given a 3D structure of a protein, scan the surface for cavities or pockets that most likely represent binding sites. Since several recent reviews provide detailed insights into these methods (Henrich et al., 2009; Laurie and Jackson, 2006; Perot et al., 2010), we only summarize the main findings as to their detection performance here.

In general, pocket prediction methods can be divided into two categories: energy- and geometry-based algorithms. Energy-based methods aim at finding pockets by computing the interaction energy between protein atoms and a small-molecule probe. By contrast, geometry-based methods try to detect solvent accessible regions that are embedded in the protein surface solely using geometric criteria. A recent comparison of energy- and geometry-based algorithms revealed that, in general, both types of binding pocket detection algorithms exhibit a very good performance (Schmidtke et al., 2010). Especially for *holo* structures, the performance of the compared methods is very similar because all methods correctly predict around 95% of the known binding sites, even though the underlying methods are rather diverse. Still, in a large-scale prediction of potential binding pockets, geometry-based algorithms were found to have some inherent advantages over energy-based algorithms because the former are faster and more robust against structural variations or missing atoms/residues in the input structures (Schmidtke et al., 2010). Notably, many pocket detection algorithms are freely available via web-servers, or are accessible via commercial software packages (see Table 1 in Perot et al., 2010 for a detailed overview).

After successfully identifying the potential binding pockets of a protein, the pockets can then be analyzed and characterized by computational means for BSC and DP.

## 2. Detection of similarity among protein binding sites

The number of protein structures deposited in the protein data bank (PDB) (Berman et al., 2000) is steadily growing. For some proteins, their 3D structures are obtained even before their functions are revealed, e.g., as a result of structural genomics projects. However, the knowledge of a protein's function is of tremendous importance in biotechnology and drug discovery. Most frequently, experimental methods are used in combination with computational approaches to determine this function (Rentzsch and Orengo, 2009). Importantly, when predicting a protein's function, it has to be considered that the definition of biological function is highly contextual (Godzik et al., 2007). Below, we concentrate on methods that predict *biochemical* protein functions.

Following the hypothesis that homologous proteins of similar function have conserved sequences and structures, many computational methods have been developed that aim at detecting functional similarities based on overall protein sequence or structure (Devos and Valencia, 2000). However, proteins of dissimilar sequences or structures can have a similar function, too (Gold and Jackson, 2006a). A typical example for this so-called convergent evolution is given by the two serine proteases trypsin and subtilisin. Although these proteins are dissimilar in overall sequence (sequence identity: ~17%) or fold (CATH class 2.40.10.10 for trypsin versus 3.40.50.200 for subtilisin; Fig. 2), they both hydrolyze peptides. This can be explained by a high local similarity of the active sites of these proteins, which is markedly shown by both enzymes containing the same catalytic triad (His, Asp, Ser). Thus, by identifying local similarities within binding sites, functional relationships among proteins shall be detectable, too.

Following the hypothesis that proteins with similar biochemical functions also bind similar ligands, local similarities within binding sites can be identified by comparing molecular recognition features in these sites, e.g., size, shape, and physico-chemical properties. In fact, many such methods have been developed recently. Besides being applied for protein function de-orphanization, these methods can also be used to understand drug selectivity, or to predict cross-reactivities of drugs.

In general, all BSC approaches consist of three main steps, which include an encoding of the molecular recognition features of the site in a computer-accessible manner (the "binding site representation"), a search for similar sites, and a similarity scoring. Since the quality of the representation greatly influences the subsequent similarity searching and scoring steps, it is the first step that

determines the predictive power of the BSC approach to a large extent. In the subsequent similarity searching step, the optimum superposition of two binding site representations is determined. To this end, established similarity assessment approaches such as sequence alignment (Binkowski et al., 2003), geometric hashing (Shulman-Peleg et al., 2004), or graph-based clique detection methods (Kinoshita et al., 2001) can be adopted. A detailed and comprehensible description of similarity searching approaches for protein binding sites is given in a recent review by Kellenberger et al. (2008). In the final similarity scoring step, the degree of similarity between the superimposed representations is quantified. For this purpose, geometric criteria like root-mean-square deviation (RMSD) (Jambon et al., 2003), residue conservation (Stark, 2003), physico-chemical property similarity (Kinoshita et al., 2001), or fingerprint overlap (Schalon et al., 2008) are frequently utilized.

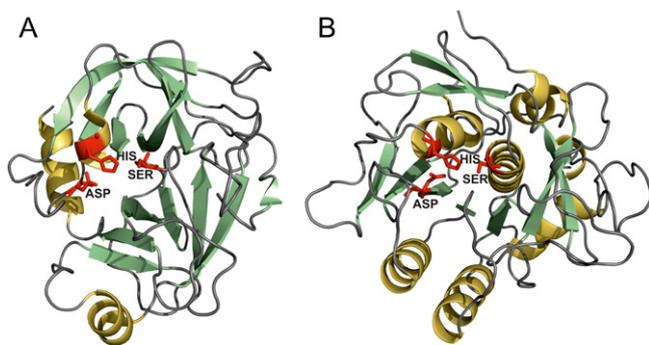
Due to the importance of the binding site representation step, we decided to categorize the BSC approaches according to the representation method. In the following, we will present a general description of binding site representation approaches and discuss their scopes and limitations. Since it is our aim to present the principal ideas of how binding sites can be represented and compared, we describe only selected methods within each category rather than giving a comprehensive overview of all available methods.

### 2.1. How to best characterize binding sites for a pair-wise comparison? Binding site representations with different levels of abstraction

Depending on the binding site representation, all available methods can be divided into two principal categories: geometry-based and signature-based approaches. The geometry-based algorithms make use of the exact 3D coordinates of a set of amino acids, atoms, pseudo-atoms, or surface points flanking the binding site to represent and compare binding sites. By contrast, signature-based approaches use numerical descriptors to encode characteristic features of the entire binding site, e.g., its shape or physico-chemical properties.

#### 2.1.1. Geometry-based binding site representations

The 3D coordinates of atoms flanking the binding sites can be incorporated using different levels of abstraction. Additionally, characteristic features such as atom type, residue type, or physico-chemical properties can be included into the similarity assessment. A simple and rather abstract representation of the binding site uses all amino acids flanking the binding site. The rationale behind this approach is given by the observation that, when comparing two proteins, binding sites have a ~20% higher sequence identity than the overall proteins (Tseng et al., 2009a). Consequently, local amino acid sequences extracted from binding sites can still uncover similarity relationships among proteins that could not be detected using overall sequence identity. Often, amino acids flanking a binding site are widely scattered in the primary sequence of the protein. Hence, the binding site is actually represented by a set of disconnected fragments, each containing a few amino acids only. Therefore, the local amino acid sequences of two binding sites are best compared using standard sequence alignment algorithms in combination with customized scoring matrices that allow a reasonable alignment of unconnected sequences (Binkowski et al., 2003; Powers et al., 2006; Tseng et al., 2009a). However, some of these methods depend on the order of the amino acid within the aligned sequences. Hence, they cannot detect similarities among local sequences in which the critical amino acids have different orderings. Recently, a BSC approach called SOIPPA (Xie and Bourne, 2008) was developed that represents a protein structure as a 3D graph with position-specific geometric and evolutionary information assigned to each C $\alpha$  atom. Then, a sequence order-independent



**Fig. 2.** 3D structure and catalytic triad of trypsin and subtilisin. The 3D structures of trypsin (shown in A, PDB code 1PPH) and subtilisin (shown in B, PDB 1BH6), which are two evolutionarily unrelated serine proteases, are shown. Even though the folds of these enzymes are distinct, each active site contains a conserved catalytic triad consisting of histidine (HIS), serine (SER), and aspartic acid (ASP).

alignment approach is applied to these graph representations, rendering this approach applicable to local amino acid sequences with different orders in different proteins.

Instead of using entire amino acids, binding sites can be represented in a more detailed way by using the 3D coordinates of individual flanking atoms (Angaran et al., 2009; Brakoulias and Jackson, 2004; Najmanovich et al., 2008). Furthermore, the program SitesBase (Gold and Jackson, 2006b) combines these two binding site representations by calculating a binding site similarity based on geometric matching of exact atom coordinates in combination with an alignment of amino acid sequences. As a principal limitation of the detailed binding site representation, small structural changes might significantly alter the resulting similarity score. Hence, these approaches cannot be applied to low resolution crystal structures or proteins showing even moderate flexibility within the binding site. Furthermore, atoms or amino acids may have different positions while playing equivalent roles in ligand binding. This fact is considered in the recently developed sup-CK approach (Hoffmann et al., 2010), which represents each binding site by a cloud of atoms. That way, binding site similarity can be computed based on similarities in atom densities in 3D space, which is considered to render the method more robust against small structural variations within the binding site.

Another level of abstraction can be incorporated into the binding site representation step by focusing on the physico-chemical properties of the atoms or amino acids flanking the binding site. This way of representing binding sites was initially realized in the Cavbase approach (Schmitt et al., 2002), which encodes amino acids as pseudo centers associated with one of five generic properties essential for molecular recognition (hydrogen bond donors and acceptors, mixed donors/acceptors, and hydrophobic or aromatic contacts). Subsequently, this approach was broadly adopted by other methods including ProBis (Konc and Janezic, 2010), SiteEngine (Shulman-Peleg et al., 2004), and MultiBind (Shulman-Peleg et al., 2008).

### 2.1.2. Signature-based binding site representations

All methods presented so far make use of the exact 3D coordinates of atoms, pseudo-atoms, or amino acids to superimpose two binding site representations and subsequently quantify the similarity between them. By contrast, signature-based approaches allow a similarity assessment among pairs of binding sites irrespective of exact 3D coordinates making these methods in general more robust against small structural changes within the binding site.

The SiteAlign method developed by Schalon et al. (2008) accomplishes a similarity assessment by mapping binding site properties onto a discretized sphere that is placed at the center of the ligand binding site. The properties of each amino acid flanking the binding site are encoded using three topological and five physico-chemical descriptors, which are then projected onto the sphere resulting in a fixed-length cavity fingerprint. Aligning two binding sites is performed by finding the highest possible similarity between the respective fingerprints. The PocketMatch approach (Yeturu and Chandra, 2008) describes binding sites as lists of distances between selected atoms of each amino acid flanking the binding site. These lists are then aligned using an incremental alignment method to obtain the final similarity score.

Both of the above binding site representations still require the computation of an optimal alignment, which is frequently time-consuming, and erroneous alignments lead to underestimated similarity scores. Thus, most recently, it was attempted to develop alignment-independent binding site representations. The FuzCav method (Weill and Rognan, 2010) represents a binding site using pharmacophoric fingerprints. Here, the  $C_{\alpha}$  atoms of all amino acids lining the binding site are annotated using six pharmacophoric properties (hydrogen bond donor or acceptor, positive

ionizable, negative ionizable, aromatic, aliphatic). Then, fingerprints are generated by considering all pharmacophoric triplets (three properties and three related distances) occurring at binned inter-feature distances. Finally, two fingerprint representations are compared by counting the number of shared features, which provides a simple and robust measure. Surface shapes of binding sites have also been encoded in an alignment-independent manner using 3D Zernike descriptors (Chikhi et al., 2010). The 3D Zernike function expansion allows a translation- and rotation-invariant series expansion of any 3D function. This represents the binding site shape in a compact manner as a vector of weights assigned from the terms in the function expansion, which allows an efficient comparison of binding sites afterwards.

### 2.2. How to evaluate the performance of binding site comparison approaches? Datasets and context-dependent performance criteria

For thoroughly comparing the algorithms with respect to their performances, high-quality benchmarking datasets and consistent performance criteria are required. An initial dataset that has been frequently used to compare the performance of BSC algorithms was assembled by Kahraman et al. (2007). This dataset contains 100 protein binding sites bound to one of nine ligand types of varying size and flexibility. However, a comparison of methods on this dataset revealed that, strikingly, no method performed significantly better in clustering proteins binding the same ligand than the descriptor “binding site volume” alone (Hoffmann et al., 2010). This clearly shows that even though the Kahraman dataset was frequently utilized to assess BSC methods (Chikhi et al., 2010; Kahraman et al., 2007; Najmanovich et al., 2008; Sael and Kihara, 2010; Spitzer et al., 2011), it is not well suited for this purpose because its protein binding sites can be discriminated already when using a simple descriptor. Consequently, a dataset was assembled by Hoffmann et al. (2010) that aims at mimicking the Kahraman dataset but, at the same time, limiting the effects of diverse ligand sizes and pocket volumes on the binding pocket comparison. This dataset comprises 100 protein structures in complex with ten ligands of a similar size. It was shown that this novel dataset is not biased against binding site volume, because volume alone does not perform better than a random classifier on this dataset (Hoffmann et al., 2010). Thus, we recommend using the Hoffmann dataset as a benchmarking dataset in the future. Both the Hoffmann and Kahraman datasets have been assembled following the assumption that proteins binding similar ligands have binding sites of similar physico-chemical or biochemical properties. Yet, there is the possibility that a ligand can bind to more than one binding site in different orientations, at which case this assumption is not valid. Additionally, it has to be considered that proteins binding identical ligands can still have diverse (biochemical) functions. As an example, the Hoffmann dataset contains ten proteins binding to *S*-adenosyl methionine; most of these are enzymes with a transferase activity whereas some are involved in DNA- or RNA-binding. Hence, binding the same ligand does not always imply identical protein function.

Since protein function prediction is the ultimate goal of many BSC approaches, the performance of methods can also be evaluated based on their ability to correctly detect similarities among proteins having the same function. One of the most extensive validations regarding protein function prediction was performed by Tseng et al. (2009a) using Enzyme Commission (EC) numbers (Bairoch, 2000) and Gene Ontology (GO) annotations (Lomax, 2005). Using an acetylcholinesterase binding site as a query cavity, 70 proteins were predicted to have similar binding sites, and indeed all of them have the same EC numbers as the query. Furthermore, the binding site of deformylase was used as a template, and a total of 94 proteins were found to have similar binding sites, of which 50

proteins share all three GO terms and 40 proteins have no GO annotation. Only four proteins were found to have different GO terms and, therefore, were considered as incorrect predictions. Finally, EC numbers were used for a large-scale enzyme function prediction on a dataset of 100 enzyme families, each containing at least 10 structures and each being represented by a different EC number. Thus, in principle, datasets are available for the evaluation of a method's ability to correctly detect similarity among binding sites of proteins that share the same protein function annotation, even though none of these datasets was intentionally designed as a benchmark. However, when comparing methods on this account, it has to be considered that the definition of biological function is highly contextual (Godzik et al., 2007) and that protein function annotation from different annotation systems might be inconsistent or erroneous (Devos and Valencia, 2000). However, this complication is common to all protein function prediction approaches, including sequence or fold-based approaches, as comprehensively outlined in a review by Godzik et al. (2007).

Finally, the ability to correctly identify similarity relationships within target families can be used as a third validation criterion. To this end, various methods were validated based on their ability to quantify binding site similarities across serine proteases (Jambon et al., 2003; Milik et al., 2003; Schalon et al., 2008; Schmitt et al., 2002; Shulman-Peleg et al., 2004; Weill and Rognan, 2010), which is a protein family that exhibits different SCOP folds and substrate cleavage specificities. For instance, the Cavbase approach was validated using a trypsin cavity query to rank a large data set of more than 5000 cavities. On the top ranks, only members of the trypsin family were detected; these were followed by other serine proteases adopting a similar fold but having a decreased sequence similarity. Furthermore, structurally unrelated serine proteases from the subtilisin family were found among the top 3% of the data set, highlighting the ability of binding site-based approaches to detect similarities among functionally related proteins that exhibit diverse sequences and folds. Furthermore, kinases were frequently utilized to assess similarity relationships within large protein families (Kinnings and Jackson, 2009; Powers et al., 2011; Spitzer et al., 2011; Weill and Rognan, 2010). A comprehensive polypharmacological analysis based on binding site similarity among kinases was published by Milletti and Vulpetti (2010). Here, the selectivity profile of 17 kinase inhibitors against 189 kinases was rationalized using predicted similarities and dissimilarities of the ATP binding sites. The selectivity profiles had been taken from the kinase interaction map published by Karaman et al. (2008), which is currently one of the most comprehensive and publicly available studies of kinase inhibitor selectivity.

Thus, in total three validation criteria to assess the performance of binding site similarity detection approaches are available, and it strongly depends on the general aim of the method which one of these criteria is most useful for the validation. Fortunately, for each of the three categories at least one publicly available dataset was published that can be used as a benchmarking dataset. Up to now, only a few newly published approaches were compared to already existing methods (Brylinski and Skolnick, 2008; Chikhi et al., 2010; Spitzer et al., 2011; Weill and Rognan, 2010; Xie and Bourne, 2008; Yeturu and Chandra, 2008). However, things are getting better because more and more recently published methods contain a comparison to at least one other method, whereas methods published before 2008 almost always completely lack this type of validation. One of the most comprehensive performance comparisons was carried out by Yeturu and Chandra (2008), who compared their PocketMatch approach to other methods including SitesBase (Gold and Jackson, 2006a,b), SuMo (Jambon et al., 2003), and ProFunc (Laskowski et al., 2005). Hence, algorithms with diverse binding site representations and diverse similarity matching approaches were compared. Very consistent results were observed for obviously

and highly similar proteins, whereas the detection of protein pairs with medium similarity was not that consistent. This observation was justified by the different site representations included in the methods, which all capture different aspects of the binding site. Furthermore, it was elucidated that some similarities among related binding sites were not detected due to wrong alignments by some approaches. This observation is in line with the results of Weill and Rognan (2010) who compared their alignment-independent FuzCav approach to three other alignment-dependent methods, SiteEngine (Shulman-Peleg et al., 2004), SiteAlign (Schalon et al., 2008), and BSAAlign (Aung and Tong, 2008). This study revealed a higher accuracy of FuzCav due to misalignments of unrelated proteins by the alignment-dependent methods.

In summary, only a few of the many available methods for the detection of local similarities among protein binding sites have been directly compared up to now. However, an extensive comparison of methods on a high-quality benchmarking dataset is required to unmask strengths and weaknesses of the binding site representation categories, to judge the general performance of the similarity matching algorithms, and to conclude whether certain methods are in general superior to others.

As outlined in Section 1, protein binding sites can exhibit substantial conformational variations. This is important in view of the fact that protein function predictions are most often desired for orphan targets in the *apo* state. Yet, methods for such predictions are frequently validated against *holo* complexes. Thus, the comparison (and the check of transferability) of similarity scores for *apo* and *holo* structures is of crucial importance. Consequently, a careful validation of a novel BSC approach should assess the influence of structural changes on its outcome by all means. However, this kind of validation is lacking in almost all currently published studies. Two notable exceptions are (I) the validation of FuzCav (Weill and Rognan, 2010) and (II) the 3D Zernike surface-based approach by Chikhi et al. (2010). Both methods were found to be quite robust in detecting true similarities also among *apo* structures. Hence, while these initial investigations are encouraging, a more comprehensive analysis of the influence of structural variations on BSC is highly desirable.

### 2.3. What is the benefit of binding site comparison approaches? Selected examples and prospective applications

The advantage of local BSC methods over global sequence, structure, or fold-based similarity detection methods originates from their ability to detect similar binding sites among proteins with even remote evolutionary relation. To emphasize this benefit of binding site-based approaches, we will present selected examples of detecting local binding site similarities among remotely related proteins. Binkowski et al. (2003) detected an unexpected similarity among binding sites of HIV-1 protease and heat shock protein-90 (HSP90). An alignment of 15 amino acids flanking the binding sites yielded a match of ten residues, showing that the key pocket residues from HIV-1 protease are conserved in HSP90. Furthermore, Minai et al. (2008) reported an unexpected similarity between a porcine pancreatic elastase binding site and the ibuprofen-binding region of phospholipase A2. Subsequently, the authors were able to experimentally confirm the previously unknown binding between ibuprofen and porcine pancreatic elastase using NMR spectroscopy.

Other recent prospective applications of BSC approaches were able to detect cross-reactivities of known drugs with other proteins than their primary target. The SiteAlign (Schalon et al., 2008) algorithm was used to find binding sites that are similar to the ATP binding site of PIM-1 kinase. Besides obviously related ATP binding sites of other kinases, an ATP binding pocket from synapsin I was predicted to be similar to the query pocket. Thus, a cross-reactivity of protein kinase inhibitors with synapsin I was suggested

and, subsequently, confirmed by experiment for four selected protein kinase inhibitors (De Franchi et al., 2010). Additionally, the SOIPPA (Xie and Bourne, 2008) approach enabled the detection of similarities among the binding sites of human catechol-O-methyltransferase (COMT), a target in Parkinson's disease, and enoyl-acyl carrier protein reductase (InhA), a tuberculosis target (Kinnings et al., 2009). Accordingly, it was proposed that two marketed drugs inhibiting COMT might also inhibit InhA and so can also be used to treat tuberculosis. The inhibition was subsequently confirmed by kinetic assays. Thus, these selected applications demonstrate the practical importance of computational BSC methods in detecting cross-reactivities of known drugs or predicting novel targets for known drugs.

Furthermore, Tseng et al. (2009a) showed the practical usability of their binding pocket comparison approach in predicting protein function for orphan targets. BioH is a protein from *Escherichia coli* with unknown function, which was conjectured to be involved in biotin biosynthesis. All known structural homologs have a sequence identity below 20%. Thus, sequence-based approaches cannot be applied to predict protein function for this orphan protein. However, the pevoSOAR (Tseng et al., 2009a) approach allows the prediction of a binding profile by BSA, which suggests that BioH is most likely a carboxylic ester hydrolase. This prediction is in line with an experimental analysis of BioH that also revealed a carboxylesterase activity of BioH (Sanishvili et al., 2003). Hence, this exemplary application to protein function prediction shows that computational BSC approaches are of practical importance in de-orphanizing protein functions: Either by directly predicting the protein function or by suggesting a limited number of biochemical assays, which can then be carried out to finally reveal a protein's function.

### 3. Computational prediction of druggability

Besides detecting similarities among protein binding sites, computational BSA approaches can also be utilized for DP. Due to the overall high costs and high failure rates even in late stages of drug discovery (Bains, 2004), a key step in early stages is the selection, prioritization, and validation of suitable drug targets. Obviously, to be considered as a potential target, the respective protein has to be involved in a particular metabolic or signaling pathway that is linked to the disease of interest (Lindsay, 2003). Additionally, the target protein's biological function should be tunable by binding small molecules (the drug candidates). Both requirements are combined in the term "druggability", which was initially defined by Hopkins and Groom (2002) as the ability of a protein linked to a specific disease to bind small, drug-like molecules with high affinity.

Experimental approaches to assess a target's druggability include high-throughput screening (Gupta et al., 2009) or NMR fragment screening (Hajduk et al., 2005). Here, binding of a large amount of diverse compounds to the target of interest is tested, and the observed hit rates are used as a measure of druggability. However, besides being time-consuming and expensive, these approaches heavily depend on the quality and diversity of the screening database due to the vastness of the drug-like chemical space (Oprea, 2002).

From the view point of binding pockets, a druggable target needs to contain a binding site that possesses characteristic structural and physico-chemical properties that favor small-molecule binding with high affinity and specificity. Accordingly, some computational approaches were recently developed that aim at characterizing the binding site(s) of a protein and, then, predicting the druggability based on these characteristics. The ultimate goal of these approaches is a fast and cost-effective prioritization of potential drug targets using only 3D structures of the proteins as input.

In general, these approaches proceed in three major steps. Initially, the potential binding pockets of the protein of interest are detected (*vide supra*). Subsequently, the characteristic features of the protein binding site(s) are determined and described by computational means. Finally, a model is developed that transforms these characteristic features into a quantitative estimate of the protein's druggability.

#### 3.1. What makes a binding pocket druggable? Useful descriptors for computational druggability prediction

The characteristic features of a binding site that are crucial for binding a drug-like molecule have to be described as accurate as possible. The first attempt of computational DP based on binding pocket descriptors was made by Nayal and Honig (2006). They used in total 408 pocket descriptors including size, shape, electrostatics, hydrogen bonding, hydrophobicity, polarity, amino acid composition, rigidity, and secondary structure. This large descriptor pool was utilized to develop random forest-based classification models. Interestingly, only 18 of these pocket descriptors, mostly related to the size, shape, and rigidity of the binding pockets, were found to have a significant influence in their final model.

Similar observations were made in several following publications that all used a large set of diverse pocket descriptors in combination with machine learning or regression models for DP (Hajduk et al., 2005; Halgren, 2009; Krasowski et al., 2011; Schmidtke and Barril, 2010; Weisel et al., 2009). All approaches consistently found that there is no single pocket parameter that is solely able to explain a pocket's druggability. Rather, druggability depends on the interplay of various pocket parameters: Usually pocket size, buriedness, and hydrophobicity positively affect druggability, whereas at least three publications describe a negative contribution of polar contact area or hydrophilicity to druggability (Hajduk et al., 2005; Halgren, 2009; Krasowski et al., 2011). The strong influence of hydrophobicity is in line with a notable model published by Cheng et al. (2007). This model allows performing DP exclusively using biophysical modeling of non-polar desolvation based on hydrophobic surface area and curvature of the binding pocket. However, it is well-known that, besides non-polar interactions, also polar and ionic interactions can significantly contribute to drug-target interactions. For instance, the specificity of drug-target interactions is usually associated with polar interactions, as these interactions are involved in maintaining a stable 3D arrangement of the ligand in the binding pocket. Thus, the influence of polar groups to druggability should not be neglected.

A first model including a positive contribution of polarity as well as hydrophobicity and pocket size to druggability was published by Schmidtke and Barril (2010). Another DP approach that incorporates hydrophobic contacts as well as polar interactions uses molecular dynamics (MD) simulations (Seco et al., 2009). This approach is rooted in the finding that binding sites have a tendency to become desolvated upon interacting with small organic molecules. Thus, potential binding sites of the protein are initially identified by performing MD simulations with organic solvent to detect regions in the protein that have a high preference to get desolvated. To quantify the druggability of the identified binding sites, a 3D grid is placed into each binding site to count the number of times that a solvent feature falls within each grid element. By comparing the observed frequencies with the expected values, associated free energies are computed for each grid point, which are then used for DP of the binding pocket. Since the utilized molecular probe (isopropyl alcohol) consists of a polar and a lipophilic part, this approach allows a reasonable estimate of the maximal affinity for molecules containing polar and lipophilic features, whereas ionic interactions are not taken into account.

While the previously described machine learning or regression models allow a fast and efficient DP for many targets, the aim of the MD simulation-based approach is different. Due to the high computational costs of MD simulations, it is obvious that this approach cannot be applied on a large scale. Rather, selected targets can be analyzed in detail, and druggability can be rationalized on an atomic level. Additionally, the MD simulation approach does not depend on any parameters that have to be optimized using a training set. This is in contrast to the pocket descriptor-based approaches and leads to the MD simulation approach being not biased toward those target classes dominating the training set. Thus, the latter approach can also be applied to atypical binding sites, e.g., to protein–protein interaction sites (*vide infra*).

### 3.2. How to quantify druggability? Measures based on different definitions of druggability

After describing the characteristic features of a pocket by computational means, these descriptors need to be transformed into a quantitative druggability score by a computational model. As mentioned before, to this end machine learning or regression models are frequently utilized. However, finding the optimal parameter set for these models or evaluating the performance of any type of DP model strongly depends on the criterion to quantify druggability. Unfortunately, as there is no unique definition of druggability, there is no unique measure of druggability either.

The first described measure of druggability goes in line with the experimental assessment of druggability using screening hit rates. Hajduk et al. (2005) developed a simple regression model using a few pocket descriptors, aiming at reproducing the hit rates observed in NMR fragment screening. The same druggability measure was used in another study aiming at reproducing hit rates from NMR fragment screening by performing docking-based virtual fragment screening (Huang and Jacobsen, 2010). However, irrespective of whether experimental or virtual fragment screening is performed, the quality of the resulting DP strongly depends on the size and structural and functional diversity of the screened fragment collection.

Another measure of druggability is based on Hopkins' and Groom's initial definition of druggability by quantifying a target's druggability according to the maximal affinity of the most potent known drug-like ligand (Cheng et al., 2007; Halgren, 2009; Schmidtke and Barril, 2010; Seco et al., 2009). The quality of this measure strongly depends on the quality of the utilized binding affinity values because it is well-known that assays performed in different laboratories might produce significantly varying binding affinities for identical molecules. Finally, Sheridan et al. (2010) correctly state that the term “druggability” implies too much because, for a drug-like molecule to become a drug, much more is required than just high-affinity binding. Accordingly, this study only aims at predicting the “chemical tractability” or “ligandability” of a target. To develop the ligandability prediction model, a measure called drug-like density is introduced that reflects the number of neighboring pockets containing a drug-like ligand in comparison to all neighboring pockets in pocket space. For this, the pocket space was spanned by the simple pocket descriptors volume, buriedness, and hydrophobicity.

These three measures clearly show that, up to now, there is no standard definition of how to quantify a protein's druggability. Thus, a druggability score given to a specific target might vary depending on the utilized measure. However, even though all principal druggability measures have some weaknesses, they allow a qualitative classification of targets into “druggable”, “non-druggable”, and “medium druggable/difficult” targets, which is already extremely helpful for target prioritization.

### 3.3. How to develop high-quality druggability models? Datasets and validation strategies

The predictivity and applicability of a druggability model strongly depends on the quality, diversity, and size of the dataset that is used to develop and validate the model, including the quality of the experimental measurements that are used to quantify the druggability. Thus, the dataset should be carefully selected to ensure that the druggability assignments in the dataset are as accurate as possible. Besides finding targets with experimentally validated druggability of acceptable confidence, it is often even more difficult to collect sufficiently many “confirmed” negative examples (the undruggable targets). Furthermore, the dataset should contain as many targets from different protein families as possible. Otherwise, the transferability to novel targets may be limited. Finally, also the size of the dataset is of crucial importance because only a dataset of reasonable size can be split up into an adequate training and test set. Needless to say, the training and test set should be selected in an appropriate way to ensure a careful validation of the final model and to prevent over-fitting.

One of the first DP models was developed on a dataset containing 28 binding sites on 23 proteins, on which 10,000 fragments were tested via NMR screening (Hajduk et al., 2005). Here, a major limitation is the proprietary screening data. Additionally, an extension of this dataset is expensive and requires the same screening library and technology to keep the dataset consistent. Thus, even though in total two predictive models were based on this dataset, the re-usability of this dataset is limited.

A dataset of higher re-usability containing 63 protein–ligand complexes for 27 targets was presented by Cheng et al. (2007). Since all the structures were taken from the PDB, and all binding affinities were taken from the literature, this dataset is publicly available. Additionally, it contains multiple structures per target and, thus, allows assessing the consistency of DP across multiple structures of one target. However, this dataset is highly dominated by enzyme targets, which probably restricts the applicability of models developed with this dataset to this class of targets.

In order to facilitate further studies and to establish a benchmarking dataset for DP, the Cheng dataset was enlarged by Schmidtke and Barril (2010). The extended dataset, termed the “Druggable Cavity Directory”, contains 1070 structures for 70 targets and is publicly available via <http://fpocket.sourceforge.net/dcd>. Using a sequence identity cut-off of 70%, this dataset was also utilized to compile a non-redundant sub-dataset containing 45 druggable, 20 non-druggable, and 5 prodrug binding sites. Besides containing multiple protein–ligand complexes per target, the dataset assembled by Schmidtke and Barril also contains apo structures for some targets. Furthermore, since it also contains multiple structures for non-enzymes like G protein-coupled receptors or nuclear hormone receptors, the structural and functional diversity of this dataset is significantly enlarged in comparison to Cheng's dataset. As outlined before, the druggability score of a specific target depends on the definition and measure of druggability and might also change over time. Thus, the inventors of the Druggable Cavity Directory leave the final druggability classification of the contained targets open for discussion and explicitly request that scientists from the field participate in adjusting and extending the current dataset.

Only recently, an even larger, non-redundant dataset containing crystal structures of 71 druggable and 44 less druggable protein binding sites was published (Krasowski et al., 2011). For the assembly of this dataset, a binding site is defined to be druggable if it can non-covalently bind small drug-like ligands that are orally available and do not require administration as prodrugs. By contrast, a protein is defined to be less druggable if either no orally available (according to Lipinski's rule of five) ligand exists or if such ligands

are too hydrophilic ( $c \log P < -2$ ) and/or lack a sufficient ligand efficiency ( $<0.3$  kcal/mol per heavy atom). This dataset is currently the largest publicly available dataset, which also contains most of the data from Cheng's as well as Schmidtke and Barril's datasets. However, due to varying definitions of druggability, the classifications for some targets differ, e.g.  $\beta$ -lactamase is classified as druggable by Schmidtke and Barril, whereas Krasowski et al. classified this protein as less-druggable, because all marketed drugs are covalent inhibitors.

Besides using a dataset of appropriate quality, diversity, and size, a proper validation of computational DP approaches also includes the comparison of the novel methods to other previously published methods. The most comprehensive comparison of DP models was performed by Sheridan et al. (2010). Using Cheng's dataset, the performances of four methods (Cheng et al., 2007; Halgren, 2009; Schmidtke and Barril, 2010; Sheridan et al., 2010) were compared yielding a significant correlation of all methods, even though the predictions for individual targets might vary in certain cases. Thus, since all druggability scores gave roughly similar results on this dataset, none of the compared methods was found to be superior to any other.

#### 3.4. Reproducible druggability scores for different conformations? Investigating the influence of conformational variability on druggability prediction

Since computational druggability models solely make use of the 3D structure of the targets, it is of crucial importance to investigate the influence of structural variations on the predicted druggability score and, hence, the reproducibility of the druggability scores for conformations of a protein bound to different ligands. The influence of structural variations among multiple structures of the same protein was initially investigated by Cheng et al. who analyzed multiple crystal structures for 14 proteins and concluded that their DP are reasonably robust against small structural variations such as side chain conformational changes. Similar observations were made in two following publications also reporting reasonably consistent results for proteins showing limited structural changes (maximal RMSD  $< 3$  Å) (Huang and Jacobsen, 2010; Schmidtke and Barril, 2010).

Furthermore, Schmidtke and Barril (2010) investigated three targets known to exhibit significant conformational variability resulting in significantly different predictions for two proteins (phosphodiesterase 4D, carbonic anhydrase) and comparable predictions for one protein (renin). Similar results were obtained by Huang and Jacobsen (2010) who found that a larger structural variation of the protein usually yields a larger variation in the druggability scores. This publication also compares druggability scores for *apo* and *holo* structures; strikingly, only little variations in the DP were observed.

Additionally, it was stated by Sheridan et al. (2010) that structural artifacts, e.g., missing residues or chains, might lead to significant variations in the DP for multiple structures of the same target. To overcome this limitation, the authors propose to use as many structures as possible per target and, then, use the median druggability score over all structures as the final score.

Besides using multiple experimental structures per target, an alternative approach is to use MD simulations to create conformations of a protein. Brown and Hajduk (2006) utilized this approach for an in-depth study of the conformation dependence of DP for three exemplary targets. For two targets (FK506 binding protein and the PH domain of Akt) the DP based on the static structures agreed with the predictions obtained by averaging over the MD trajectories. In contrast, for the protein–protein complex Bcl-xL, the static *apo* form is predicted to have a low druggability although Bcl-xL is known to be druggable. However, when performing MD

simulations, conformations that are predicted to be druggable dominate the trajectories. This result is in line with observations by Huang and Jacobsen (2010) who also observed a significantly stronger conformational dependence on DP for protein–protein complexes than for protein–ligand complexes.

Finally, it should be noted that the conformation dependence of any DP model is strongly influenced by the dependence of the utilized pocket prediction approach on structural variations.

#### 3.5. Selected applications of druggability: protein–protein interactions and prospective druggability predictions

PPIs are linked with many biological processes. Therefore, a computational DP for protein–protein interfaces is of great practical importance. However, as already mentioned above, the applicability of machine learning or regression models is mostly limited to protein–ligand binding sites because most of the available training datasets are dominated by protein–ligand complexes. However, two DP models were published that do not require a parameter optimization on training data: the MD simulation-based approach by Seco et al. (2009) and the virtual fragment screening approach by Huang and Jacobsen (2010). Thus, it can be expected that these approaches are also applicable to, as of now, uncommon protein binding sites such as in protein–protein interfaces, which was also shown for both approaches. The MD simulation-based approach was applied to two exemplary protein–protein interaction sites (MDM2-p53 and LFA1-ICAM1). Both protein–protein interaction sites are predicted to be druggable, with MDM2-p53 being predicted to be even more druggable than LFA1-ICAM1. These predictions are in line with the potencies of the most-potent known inhibitors ( $K_d = 3$  nM for MDM2-p53 and  $K_d = 18.3$  nM for LFA1-ICAM1). The virtual fragment screening approach was applied to six targets involved in PPIs using at least two structures per protein. For the complexes MDM2-p53 and Bcl-xL-BAK high-affinity ligands are known, and both complexes are correctly predicted to be druggable. The complexes IL2-ILR and HPV-E2 also have high-affinity ligands, but correct DP are only possible if crystal structures with a small-molecule inhibitor are used. By contrast, these PPIs are predicted to be undruggable when using crystal structures with a peptide or a protein bound. In general, the structural variations among conformations of a PPI target are larger than the variations of protein–ligand complexes. Thus, there is also a stronger dependency on the utilized structures for the DP of PPI.

The practical usefulness of computational DP was additionally shown by a prospective application of the druggability model developed by Cheng et al. (2007) to two novel targets (fungal homoserine dehydrogenase: HDS and hematopoietic prostaglandin D synthase: H-PGDS). HDS is predicted to be undruggable, and H-PGDS is predicted to be druggable. These predictions are validated using HTS hit rates as an experimental measure for druggability. The HTS results confirm the computational DP (16 hits for HDS and 200 hits for H-PGDS). Considerable follow-up efforts resulted in 11 drug-like high-affinity ligands for H-PGDS, but no drug-like ligand for HDS. This is clearly an encouraging result in that it shows the reliability of such predictions and highlights the usefulness of computational DP approaches for practical target prioritization.

#### 4. Scopes and limitations of structure-based computational analysis of protein binding sites

The two previous sections have extensively outlined the available methods aiming at BSA. Furthermore, diverse applications of BSA were presented ranging from protein function de-orphanization to the rationalization of drug selectivity and cross-reactivity to the prediction of target druggability. Irrespective of the specific aim, each method consists of three major steps:

(I) binding site identification, (II) binding site encoding, and (III) BSA. The first two steps can be performed irrespective of whether the prediction of protein function, off-target effects, or target drug-gability is the ultimate aim of the method.

#### 4.1. General dependence of binding site analysis approaches on binding site definitions

The initial binding site identification step can either be performed by focusing only on the binding site that is occupied by the ligand in the given 3D structure or by using a computational binding pocket detection approach (see Section 1). The majority of all BSC approaches were validated using binding sites that are defined by a given ligand in the 3D structure (e.g., see Gold and Jackson, 2006a or Spitzer et al., 2011). Still, a few approaches are available that allow the comparison of protein binding sites predicted by pocket detection methods, e.g., Cavbase (Schmitt et al., 2002) uses Ligsite (Hendlich et al., 1997) and ProFunc (Laskowski et al., 2005) uses SURFNET (Laskowski, 1995). By contrast, DP approaches mostly make use of computational methods to identify potential binding pockets. Subsequently, some of the machine learning-based DP approaches use given ligands in the crystal structures to discriminate druggable and non-druggable pockets (Nayal and Honig, 2006; Sheridan et al., 2010; Weisel et al., 2009). Furthermore, a few methods were developed recently that are completely independent from the pre-definition of binding sites because they scan the entire protein for similarities against a given set of binding sites (e.g., MED-Sumo (Doppelt-Azeroual et al., 2009) and SOIPPA (Xie and Bourne, 2008)) or, more specifically, druggable binding sites (Seco et al., 2009).

Identifying binding sites based on ligands in *holo* structures allows using experimentally confirmed binding sites. However, it also restricts the applicability to targets for which at least one 3D structure of a protein–ligand complex is available. Nevertheless, it has to be considered that protein de-orphanization and DP are most often applied in early research phases, and it is therefore questionable whether *holo* structures are available at this early stage. By contrast, computational pocket prediction algorithms are also applicable to *apo* structures. Furthermore, these methods allow analyzing a protein with respect to multiple binding sites per protein. This is predominantly important for DP, when multiple pockets within one protein shall be ranked according to their individual druggabilities.

Current BSA methods most often focus on analyzing whole binding sites. However, binding sites frequently consist of multiple subpockets binding specific substructural motifs. Particularly in BSC, the detection of local similarities among binding sites is an important but difficult problem. For instance, while the overall binding sites of two ATP- and NAD-binding proteins look different, the proteins nevertheless share a similar subpocket for binding adenine (Xie and Bourne, 2008). Hence, a pair of proteins whose binding sites differ significantly when compared in their entirety may still share similarity at the subpocket level (Wallach and Lilien, 2009). This complication of BSC was taken into account in the validation of some recent BSC methods (McGready et al., 2009; Weill and Rognan, 2010; Xie and Bourne, 2008) by focusing on finding local similarity among binding sites for the adenine substructural motif. Furthermore, Wallach and Lilien (2009) developed a novel approach explicitly aiming at the detection of subpocket similarity.

Since the assessment of subpocket similarity strongly depends on the way subpockets are defined, the recent development of DoGSite (Volkamer et al., 2010), a subpocket prediction approach, is of particular importance. DoGSite splits predicted pockets into subpockets, revealing a refined description of the topology of binding sites and, thereby, an improved ligand coverage (Volkamer et al., 2010). Since pocket size is an important parameter in most DP models, the correct definition of the ligand-binding subpocket is also of

importance in DP. Hence, irrespective whether whole binding sites or motif-binding subpockets are analyzed, the correct identification of the binding pocket is an important step in BSA.

#### 4.2. Scopes of binding site representations and subsequent binding site analysis methods

After identifying the potential binding site(s) of a protein structure, the characteristic features of these sites have to be determined and quantified to allow a computational assessment of binding site similarity or druggability. As outlined in Sections 2 and 3, many binding site representation methods have been developed, especially with the aim of determining similarities among binding sites. Despite this fact, the overlap of methods for BSC and DP is still small. For instance, geometric binding site representations currently dominate the available BSC methods, whereas DP approaches almost always make use of rather abstract encodings of binding site characteristics that do not rely on the exact 3D coordinates of atoms or amino acids flanking the binding site.

From a computational point of view, BSA requires a simplified computer-accessible representation of the binding site. From a biological perspective, this binding site representation must contain information about shape and physico-chemical properties of the binding sites that are crucial for molecular recognition (Kellenberger et al., 2008). Intuitively, one may expect that the highest level of information leads to the best prediction of similarity or druggability. However, it has to be considered that different amino acids can play equivalent roles in molecular interactions and that structural flexibility might lead to varying coordinates of amino acids flanking the binding site. Therefore, all binding site representations require a certain level of abstraction or, phrased differently, the most detailed description, e.g., in terms of atom coordinates, is not necessarily the best approach to represent binding sites.

In general, two principal aims of BSA approaches can be differentiated. On the one hand, BSA approaches can aim at the prediction of protein function or druggability on a large scale. On the other hand, for a limited set of proteins, methods can aim at rationalizing why two binding sites are similar or why they are druggable. The majority of all available approaches are aimed at being fast and applicable to large data sets. Among the BSC approaches, in particular the methods FuzCav (Weill and Rognan, 2010) and PocketMatch (Yeturu and Chandra, 2008) were shown to be efficient such that the screening of large databases like sc-PDB (Kellenberger et al., 2006) and PDBbind (Wang et al., 2004) at ultra-high speed (a few ms per average pair of binding sites) is possible. Among the DP approaches, in particular the fpocket approach developed by Schmidtke and Barril (2010) was designed as a high-throughput method. It was shown that, on average, this method takes two to four seconds per structure. Thus, the method can process the entire PDB within a few days on a normal computer.

A clear drawback of most of these high-throughput methods is the lack of interpretability of the results. For instance, even though the FuzCav approach is fast and alignment-independent, it is currently not possible to explain why two cavity fingerprints are similar and which residues are responsible for the observed similarity among the binding sites (Weill and Rognan, 2010). Hence, the authors propose a two-step approach: Initially, the FuzCav fingerprint should be used as a filter to check a putative similarity; subsequently, a 3D alignment tool should be used to visualize and rationalize the predicted similarity between two binding sites. A similar multi-step approach has already been realized in PocketAlign (Yeturu and Chandra, 2011) in which an initial fast seed alignment is subsequently extended by a more accurate atomic alignment.

However, for both BSC and DP also methods are available that allow a rationalization of the prediction. For instance, the MD

simulation-based approach developed by Seco et al. (2009) allows the rationalization of druggability on an atomic level. Furthermore, the MolLoc approach (Angaran et al., 2009) enables a detailed comparison of up to 20 protein binding sites using geometry and physico-chemical properties. Also, the PocketAlign approach (Yeturu and Chandra, 2011) and Cavbase (Schmitt et al., 2002) include a detailed alignment of two binding sites in combination with a visualization of the alignment to provide insights as to why two binding sites are similar.

#### 4.3. Assessing the quality of computational binding site analysis approaches

To compare the BSA methods and to determine whether some methods are superior to others, high-quality benchmarking datasets and consistent performance measures are required. As outlined in Sections 2 and 3, both DP and BSC approaches share the problem of inconsistent performance metrics. For instance, different DP approaches were developed using different definitions and measures of druggability. A similar situation is given for BSC methods because a unique criterion that defines the similarity between pairs of binding sites is also missing. Hence, this lack of consistent performance criteria makes it rather difficult to compare multiple BSA approaches.

Furthermore, high-quality benchmarking datasets are required to compare the performance of methods. As one such dataset, the Druggable Cavity Directory (Schmidtke and Barril, 2010) was developed for the comparison of DP methods. To compare BSC methods, only the Hoffmann dataset is presently available for future works (Hoffmann et al., 2010). This dataset allows the comparison of methods based on their ability to correctly detect similarities among proteins that bind the same ligand. However, as outlined in Section 2, there are also other criteria that can be used for a comparison. Yet, no dataset has so far been put forward as a benchmark for evaluating methods based on their ability to correctly detect similarities among proteins with the same protein function annotations

or among proteins belonging to the same protein family. However, it is of crucial importance to compare as many methods as possible on a high-quality benchmarking dataset to find out whether certain binding site representation approaches are better, or whether certain similarity searching methods are faster than others. Only by doing so, the scopes and limitations of these methods can be uncovered.

Besides comparing the performance of methods, an appropriate validation should also investigate the influence of conformational variability on the BSA. Whereas a sufficiently large amount of publications dealing with DP include this validation step, it is almost always missing in publications presenting novel approaches for BSC. At least, some of these methods claim to be more robust against structural variants than others. However, this is rarely effectively shown. Due to sometimes significant structural variations among conformations of the same protein (see Section 1), we recommend performing a careful validation of the robustness against structural variations for every newly developed BSA approach.

#### 4.4. Practical importance of computational binding site analysis

Both BSC and DP approaches are of practical importance in biotechnology and drug discovery. The major aim of most BSC approaches is the prediction of protein function. For the orphan protein BioH it was shown that computational BSC methods can indeed be used to perform protein de-orphanization. Furthermore, it was shown that local BSC methods are able to uncover functional similarities among proteins, even in the absence of sequence or fold similarity. Besides predicting protein function, these methods can also be used to rationalize drug selectivity or to predict novel targets for known drugs. Two prospective studies were presented that predict and, subsequently, validate by experiment the interactions of known drugs with targets other than the drug's primary ones, which, again, highlight the benefit of these methods. Finally, many BSC approaches are also publicly available via web-servers or as programs that can be downloaded free of charge (Table 1).

**Table 1**  
Publicly available methods for protein binding site comparison.

Name	Reference	Web address
Cavbase	Schmitt et al. (2002)	<a href="http://relibase.ccdc.cam.ac.uk">http://relibase.ccdc.cam.ac.uk</a> <a href="http://relibase.rutgers.edu">http://relibase.rutgers.edu</a>
CPASS	Powers et al. (2011)	<a href="http://bionmr-c1.unl.edu/CPASS.OV/CPASS.htm">http://bionmr-c1.unl.edu/CPASS.OV/CPASS.htm</a>
eF-seek	Kinoshita et al. (2007)	<a href="http://ef-site.hgc.jp/eF-seek/top.do">http://ef-site.hgc.jp/eF-seek/top.do</a>
FINDSITE	Brylinski and Skolnick (2008)	<a href="http://cssb.biology.gatech.edu/skolnick/files/FINDSITE/">http://cssb.biology.gatech.edu/skolnick/files/FINDSITE/</a>
fPOP	Tseng et al. (2009b)	<a href="http://pocket.uchicago.edu/fpop/">http://pocket.uchicago.edu/fpop/</a>
FuzCav	Weill and Rognan (2010)	<a href="http://bioinfo-pharma.u-strasbg.fr/labwebsite/download.html">http://bioinfo-pharma.u-strasbg.fr/labwebsite/download.html</a>
Isocleft	Najmanovich et al. (2008)	Can be obtained from the authors
MolLoc	Angaran et al. (2009)	<a href="http://bcb.dei.unipd.it/MolLoc/">http://bcb.dei.unipd.it/MolLoc/</a>
MultiBind	Shulman-Peleg et al. (2008)	<a href="http://bioinfo3d.cs.tau.ac.il/MultiBind/">http://bioinfo3d.cs.tau.ac.il/MultiBind/</a>
PDBSiteScan	Ivanisenko et al. (2004)	<a href="http://www.mgs.bionet.nsc.ru/mgs/gnw/pdbsitescan/">http://www.mgs.bionet.nsc.ru/mgs/gnw/pdbsitescan/</a>
PESD	Das et al. (2009)	<a href="http://reccr.chem.rpi.edu/Software/pesdserve/">http://reccr.chem.rpi.edu/Software/pesdserve/</a>
PevoSOAR	Tseng et al. (2009a)	<a href="http://sts.bioengr.uic.edu/pevosoar/">http://sts.bioengr.uic.edu/pevosoar/</a>
PINTS	Stark (2003)	<a href="http://www.russelllab.org/cgi-bin/tools/pints.pl">http://www.russelllab.org/cgi-bin/tools/pints.pl</a>
PocketAlign	Yeturu and Chandra (2011)	<a href="http://proline.physics.iisc.ernet.in/pocketalign/">http://proline.physics.iisc.ernet.in/pocketalign/</a>
PocketMatch	Yeturu and Chandra (2008)	<a href="http://proline.physics.iisc.ernet.in/pocketmatch/">http://proline.physics.iisc.ernet.in/pocketmatch/</a>
ProBiS	Konc and Janezic (2010)	<a href="http://probis.cmm.ki.si/">http://probis.cmm.ki.si/</a>
ProFunc	Laskowski et al. (2005)	<a href="http://www.ebi.ac.uk/thornton-srv/databases/ProFunc/">http://www.ebi.ac.uk/thornton-srv/databases/ProFunc/</a>
Query3d	Ausiello et al. (2005)	<a href="http://pdbsfun.uniroma2.it/">http://pdbsfun.uniroma2.it/</a>
SiteAlign	Schalon et al. (2008)	<a href="http://bioinfo-pharma.u-strasbg.fr/labwebsite/download.html">http://bioinfo-pharma.u-strasbg.fr/labwebsite/download.html</a>
SitesBase	Gold and Jackson (2006b)	<a href="http://www.modelling.leeds.ac.uk/sb/">http://www.modelling.leeds.ac.uk/sb/</a>
SiteEngine	Shulman-Peleg et al. (2004)	<a href="http://bioinfo3d.cs.tau.ac.il/SiteEngine/">http://bioinfo3d.cs.tau.ac.il/SiteEngine/</a>
SMAP-WS (SOIPPA)	Xie and Bourne (2008)	<a href="http://kryptonite.nbc.net/smap-ws/">http://kryptonite.nbc.net/smap-ws/</a>
SuMo	Jambon et al. (2003)	<a href="http://sumo-pbil.ibcp.fr">http://sumo-pbil.ibcp.fr</a>
SURFCOMP	Hofbauer et al. (2004)	Can be obtained from the authors

The second category of computational BSA approaches aims at DP. Rather few of these methods have been developed, in contrast to many methods for BSC. Notably, most of the DP methods are very carefully validated, and the practical applicability of these methods in prioritizing target proteins was proven in a prospective study. Furthermore, it was shown that selected methods are also applicable to PPIs, which are linked to many important biological processes. Finally, at least one DP method is currently publicly available via a web-server (the fpocket approach can be accessed via <http://fpocket.sourceforge.net>).

Clearly, the recent progress in the area of structure-based approaches for the computational analysis of protein binding sites is impressive, and the interested reader is encouraged to apply these methods in her/his own work.

## Acknowledgments

We are grateful to the Ministry of Innovation, Science, and Research of North Rhine-Westphalia and Heinrich-Heine-University Düsseldorf for a scholarship to FS within the CLIB-Graduate Cluster Industrial Biotechnology.

## References

- Ahmed, A., Kazemi, S., Gohlke, H., 2007. Protein flexibility and mobility in structure-based drug design. *Front. Drug Des. Discov.* 3, 455–476.
- Angaran, S., Bock, M.E., Garutti, C., Guerra, C., 2009. MolLoc: a web tool for the local structural alignment of molecular surfaces. *Nucleic Acids Res.* 37, W565–W570.
- Aung, Z., Tong, J.C., 2008. BSAalign: a rapid graph-based algorithm for detecting ligand-binding sites in protein structures. *Genome Inf.* 21, 65–76.
- Ausiello, G., Via, A., Helmer-Citterich, M., 2005. Query3d: a new method for high-throughput analysis of functional residues in protein structures. *BMC Bioinformatics* 6 (Suppl. 4), S5.
- Bains, W., 2004. Failure rates in drug discovery and development: will we ever get any better? *Drug Discov. World* (Fall), 9–18.
- Bairoch, A., 2000. The ENZYME database in 2000. *Nucleic Acids Res.* 28, 304–305.
- Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N., Bourne, P.E., 2000. The Protein Data Bank. *Nucleic Acids Res.* 28, 235–242.
- Binkowski, T.A., Adamian, L., Liang, J., 2003. Inferring functional relationships of proteins from local sequence and spatial surface patterns. *J. Mol. Biol.* 332, 505–526.
- Bogan, A.A., Thorn, K.S., 1998. Anatomy of hot spots in protein interfaces. *J. Mol. Biol.* 280, 1–9.
- Brakoulias, A., Jackson, R.M., 2004. Towards a structural classification of phosphate binding sites in protein–nucleotide complexes: an automated all-against-all structural comparison using geometric matching. *Proteins* 56, 250–260.
- Brown, S.P., Hajduk, P.J., 2006. Effects of conformational dynamics on predicted protein druggability. *Nat. Biotechnol.* 25, 71–75.
- Brylinski, M., Skolnick, J., 2008. A threading-based method (FINDSITE) for ligand-binding site prediction and functional annotation. *Proc. Natl. Acad. Sci. U.S.A.* 105, 129–134.
- Cheng, A.C., Coleman, R.G., Smyth, K.T., Cao, Q., Soulard, P., Caffrey, D.R., Salzberg, A.C., Huang, E.S., 2007. Structure-based maximal affinity model predicts small-molecule druggability. *Nat. Biotechnol.* 25, 71–75.
- Chikhi, R., Sael, L., Kihara, D., 2010. Real-time ligand binding pocket database search using local surface descriptors. *Proteins* 78, 2007–2028.
- Cozzini, P., Kellogg, G.E., Spyralis, F., Abraham, D.J., Costantino, G., Emerson, A., Fanelli, F., Gohlke, H., Kuhn, L.A., Morris, G.M., Orozco, M., Pertinhez, T.A., Rizzi, M., Sotriffer, C., 2008. Target flexibility: an emerging consideration in drug discovery and design. *J. Med. Chem.* 51, 6237–6255.
- Das, S., Kokardekar, A., Breneman, C.M., 2009. Rapid comparison of protein binding site surfaces with property encoded shape distributions. *J. Chem. Inf. Model.* 49, 2863–2872.
- De Franchi, E., Schalon, C., Messa, M., Onofri, F., Benfenati, F., Rognan, D., 2010. Binding of protein kinase inhibitors to synapsin I inferred from pair-wise binding site similarity measurements. *PLoS One* 5, e12214.
- Devos, D., Valencia, A., 2000. Practical limits of function prediction. *Proteins* 41, 98–107.
- Doppelt-Azeroual, O., Moriaud, F., Adcock, S.A., Delfaud, F., 2009. A review of MED-SuMo applications. *Infect. Disord.: Drug Targets* 9, 344–357.
- Eyrisch, S., Helms, V., 2007. Transient pockets on protein surfaces involved in protein–protein interaction. *J. Med. Chem.* 50, 3457–3464.
- Godzik, A., Jambon, M., Friedberg, I., 2007. Computational protein function prediction: are we making progress? *Cell. Mol. Life. Sci.* 64, 2505–2511.
- Gold, N.D., Jackson, R.M., 2006a. Fold independent structural comparisons of protein–ligand binding sites for exploring functional relationships. *J. Mol. Biol.* 355, 1112–1124.
- Gold, N.D., Jackson, R.M., 2006b. SitesBase: a database for structure-based protein–ligand binding site comparisons. *Nucleic Acids Res.* 34, D231–D234.
- Grimme, D., Gonzalez-Ruiz, D., Gohlke, H. Computational strategies and challenges for targeting protein–protein interactions with small molecules. In: Luque, F.J., Barril, X. (Eds.), *Physico-chemical and computational approaches to drug discovery*, in press.
- Gupta, A., Gupta, A.K., Seshari, K., 2009. Structural models in the assessment of protein druggability based on HTS data. *J. Comput. Aided Mol. Des.* 23, 583–592.
- Hajduk, P.J., Huth, J.R., Fesik, S.W., 2005. Druggability indices for protein targets derived from NMR-based screening data. *J. Med. Chem.* 48, 2518–2525.
- Halgren, T.A., 2009. Identifying and characterizing binding sites and assessing drug-gability. *J. Chem. Inf. Model.* 49, 377–389.
- Hendlich, M., Rippmann, F., Barnickel, G., 1997. LIGSITE: automatic and efficient detection of potential small molecule-binding sites in proteins. *J. Mol. Graph. Model.* 15, 359–363.
- Henrich, S., Salo-Ahen, O.M.H., Huang, B., Rippmann, F., Cruciani, G., Wade, R.C., 2009. Computational approaches to identifying and characterizing protein binding sites for ligand design. *J. Mol. Recognit.* 23, 209–219.
- Hofbauer, C., Lohninger, H., Aszódi, A., 2004. SURFCOMP: a novel graph-based approach to molecular surface comparison. *J. Chem. Inf. Comput. Sci.* 44, 837–847.
- Hoffmann, B., Zaslavskiy, M., Vert, J.P., Stoven, V., 2010. A new protein binding pocket similarity measure based on comparison of clouds of atoms in 3D: application to ligand prediction. *BMC Bioinformatics* 11, 99.
- Hopkins, A.L., Groom, C.R., 2002. The druggable genome. *Nat. Rev. Drug Discov.* 1, 727–730.
- Hou, T., Wang, J., Li, Y., Wang, W., 2011. Assessing the performance of the MM/PBSA and MM/GBSA methods. 1. Accuracy of binding free energy calculations based on molecular dynamics simulations. *J. Chem. Inf. Model.* 51, 69–82.
- Huang, N., Jacobsen, M.P., 2010. Binding-site assessment by virtual fragment screening. *PLoS One* 5, e10109.
- Ivanisenko, V.A., Pintus, S.S., Grigorovich, D.A., Kolchanov, N.A., 2004. PDBSiteScan: a program for searching for active, binding and posttranslational modification sites in the 3D structures of proteins. *Nucleic Acids Res.* 32, W549–W554.
- Jambon, M., Imberty, A., Deleage, G., Geourjon, C., 2003. A new bioinformatic approach to detect common 3D sites in protein structures. *Proteins* 52, 137–145.
- Kahraman, A., Morris, R.J., Laskowski, R.A., Thornton, J.M., 2007. Shape variation in protein binding pockets and their ligands. *J. Mol. Biol.* 386, 283–301.
- Karaman, M.W., Herrgard, S., Treiber, D.K., Gallant, P., Atteridge, C.E., Campell, B.T., Chan, K.W., Ciceri, P., Davis, M.I., Edeen, P.T., Faraoni, R., Floyd, M., Hunt, J.P., Lockhart, D.J., Milanov, Z.V., Morrison, M.J., Pallares, G., Patel, H.K., Pritchard, S., Wodicka, L.M., Zarrinkar, P.P., 2008. A quantitative analysis of kinase inhibitor selectivity. *Nat. Biotechnol.* 26, 127–132.
- Kellenberger, E., Muller, P., Schalon, C., Bret, G., Foata, N., Rognan, D., 2006. sc-PDB: an annotated database of druggable binding sites from the Protein Data Bank. *J. Chem. Inf. Model.* 46, 717–727.
- Kellenberger, E., Schalon, C., Rognan, D., 2008. How to measure the similarity between protein ligand-binding sites? *Curr. Comput.-Aided Drug Des.* 4, 209–220.
- Kinnings, S.L., Jackson, R.M., 2009. Binding site similarity analysis for the functional classification of the protein kinase family. *J. Chem. Inf. Model.* 49, 318–329.
- Kinnings, S.L., Liu, N., Buchmeier, N., Tonge, P.J., Xie, L., Bourne, P.E., 2009. Drug discovery using chemical systems biology: repositioning the safe medicine Comtan to treat multi-drug and extensively drug resistant tuberculosis. *PLoS Comput. Biol.* 5, e1000423.
- Kinoshita, K., Furui, J., Nakamura, H., 2001. Identification of protein functions from a molecular surface database, eF-site. *J. Struct. Funct. Genomics* 2, 9–22.
- Kinoshita, K., Murakami, Y., Nakamura, H., 2007. eF-seek: prediction of the functional sites of proteins by searching for similar electrostatic potential and molecular surface shape. *Nucleic Acids Res.* 35, W398–W402.
- Konc, J., Janezic, D., 2010. ProBiS algorithm for detection of structurally similar protein binding sites by local structural alignment. *Bioinformatics* 26, 1160–1168.
- Krasowski, A., Muthas, D., Sarkar, A., Schmitt, S., Brenk, R., 2011. DrugPred: a structure-based approach to predict protein druggability developed using an extensive nonredundant data set. *J. Chem. Inf. Model.* 51, 2829–2842.
- Laskowski, R.A., 1995. SURFNET: a program for visualizing molecular surfaces, cavities, and intermolecular interactions. *J. Mol. Graph.* 13, 323–330, 307–308.
- Laskowski, R.A., Luscombe, N.M., Swindells, M.B., Thornton, J.M., 1996. Protein clefts in molecular recognition and function. *Protein Sci.* 5, 2438–2452.
- Laskowski, R.A., Watson, J.D., Thornton, J.M., 2005. ProFunc: a server for predicting protein function from 3D structure. *Nucleic Acids Res.* 33, W89–W93.
- Laurie, A.T.R., Jackson, R.M., 2006. Methods for the prediction of protein–ligand binding sites for structure-based drug design and virtual ligand screening. *Curr. Protein Pept. Sci.* 7, 395–406.
- Liang, J., Edelsbrunner, H., Woodward, H., 1998. Anatomy of protein pockets and cavities: measurements of binding site geometry and implications for ligand design. *Protein Sci.* 7, 1884–1897.
- Lindsay, A.M., 2003. Target discovery. *Nat. Rev. Drug Discov.* 2, 831–838.
- Lomax, J., 2005. Get ready to GO! A biologist's guide to the Gene Ontology. *Brief Bioinform.* 6, 298–304.
- Luque, I., Freire, E., 2000. Structural stability of binding sites: consequences for binding affinity and allosteric effects. *Proteins Suppl.* 4, 63–71.
- Ma, B., Shatsky, M., Wolfson, H.J.M., Nussinov, R., 2002. Multiple diverse ligands binding at a single protein site: a matter of pre-existing populations. *Protein Sci.* 11, 184–197.
- McGready, A., Stevens, A., Lipkin, M., Hudson, B.D., Whitley, D.C., Ford, M.G., 2009. Vicinity analysis: a methodology for the identification of similar protein active sites. *J. Mol. Model.* 15, 489–498.

- Metz, A., Pfeleger, C., Kopitz, H., Pfeiffer-Marek, S., Baringhaus, K.-H., Gohlke, H., 2011. Hot spots and transient pockets: predicting the determinants of small-molecule binding to a protein–protein interface. *J. Chem. Inf. Model.*, doi:10.1021/ci200322s.
- Milik, M., Szalma, S., Olszewski, K.A., 2003. Common Structural Cliques: a tool for protein structure and function. *Protein Eng.* 16, 543–552.
- Milletti, F., Vulpetti, A., 2010. Predicting polypharmacology by binding site similarity: from kinases to the protein universe. *J. Chem. Inf. Model.* 50, 1418–1431.
- Minai, R., Matsuo, Y., Onuki, H., Hirota, H., 2008. Method for comparing the structures of protein ligand-binding sites and application for predicting protein–drug interactions. *Proteins* 72, 367–381.
- Najmanovich, R., Kurbatova, N., Thornton, J., 2008. Detection of 3D atomic similarities and their use in the discrimination of small molecule protein-binding sites. *Bioinformatics* 24, i105–i111.
- Nayal, M., Honig, B., 2006. On the nature of cavities on protein surfaces: application to the identification of drug-binding sites. *Proteins* 63, 892–906.
- Oprea, T., 2002. Chemical space navigation in lead discovery. *Curr. Opin. Chem. Biol.* 6, 384–389.
- Perot, S., Sperandio, O., Miteva, M.A., Camproux, A.-C., Villoutreix, B.O., 2010. Druggable pockets and binding site centric chemical space: a paradigm shift in drug discovery. *Drug Discov. Today* 15, 656–667.
- Powers, R., Copeland, J.C., Germer, K., Mercier, K.A., Ramanathan, V., Revesz, P., 2006. Comparison of protein active site structures for functional annotation of proteins and drug design. *Proteins* 65, 124–135.
- Powers, R., Copeland, J.C., Stark, J.L., Caprez, A., Guru, A., Swanson, D., 2011. Searching the protein structure database for ligand-binding site similarities using CPASS v.2. *BMC Res. Notes* 4, 17.
- Rentzsch, R., Orengo, C.A., 2009. Protein function prediction – the power of multiplicity. *Trends Biotechnol.* 27, 210–219.
- Sael, L., Kihara, D., 2010. Binding ligand prediction for proteins using partial matching of local surface patches. *Int. J. Mol. Sci.* 11, 5009–5026.
- Sanishvili, R., Yahunin, A.F., Laskowski, R.A., Skarina, T., Evdokimova, E., Doherty-Kirby, A., Lajoie, G.A., Thornton, J.M., Arrowsmith, C.H., Savchenko, A., Joachimiak, A., Edwards, A.M., 2003. Integrating structure, bioinformatics, and enzymology to discover function: BioH, a new carboxylesterase from *Escherichia coli*. *J. Biol. Chem.* 278, 26039–26045.
- Schalon, C., Surgand, J.-S., Kellenberger, E., Rognan, D., 2008. A simple and fuzzy method to align and compare druggable ligand-binding sites. *Proteins* 71, 1755–1778.
- Schmidtke, P., Barril, X., 2010. Understanding and predicting druggability. A high-throughput method for detection of drug binding sites. *J. Med. Chem.* 53, 5858–5867.
- Schmidtke, P., Souaille, C., Estienne, F., Baurin, N., Kroemer, R.T., 2010. Large-scale comparison of four binding site detection algorithms. *J. Chem. Inf. Model.* 50, 2191–2200.
- Schmitt, S., Kuhn, D., Klebe, G., 2002. A new method to detect related function among proteins independent of sequence and fold homology. *J. Mol. Biol.* 323, 387–406.
- Seco, J., Luque, J., Barril, X., 2009. Binding site detection and druggability index from first principles. *J. Med. Chem.* 52, 2363–2371.
- Sheridan, R.P., Maiorov, V.N., Hooloway, M.K., Cornell, W.D., Gao, Y.-D., 2010. Drug-like density: a method of quantifying the bindability of a protein target based on a very large set of pockets and drug-like ligands from the protein data bank. *J. Chem. Inf. Model.* 50, 2029–2040.
- Shulman-Peleg, A., Nussinov, R., Wolfson, H.J., 2004. Recognition of functional sites in protein structures. *J. Mol. Biol.* 339, 607–633.
- Shulman-Peleg, A., Shatsky, M., Nussinov, R., Wolfson, H.J., 2008. MultiBind and MAPPIS: web servers for multiple alignment of protein 3D-binding sites and their interactions. *Nucleic Acids Res.* 36, W260–W264.
- Spitzer, R., Cleves, A.E., Jain, A.N., 2011. Surface-based protein binding pocket similarity. *Proteins* 79, 2746–2763.
- Stark, A., 2003. Annotation in three dimensions. PINTS: Patterns in Non-homologous Tertiary Structures. *Nucleic Acids Res.* 31, 3341–3344.
- Tsai, C.-J., Ma, B., Nussinov, R., 1999. Folding and binding cascades: shifts in energy landscapes. *Proc. Natl. Acad. Sci. U.S.A.* 96, 9970–9972.
- Tseng, Y.Y., Dundas, J., Liang, J., 2009a. Predicting protein function and binding profile via matching of local evolutionary and geometric surface patterns. *J. Mol. Biol.* 387, 451–464.
- Tseng, Y.Y., Chen, Z.J., Li, W.H., 2009b. fPOP: footprinting functional pockets of proteins by comparative spatial patterns. *Nucleic Acids Res.* 38, D288–D295.
- Villar, H.O., Kauvar, L.M., 1994. Amino acid preferences at protein binding sites. *FEBS Lett.* 349, 125–130.
- Volkamer, A., Griewel, A., Grombacher, T., Rarey, M., 2010. Analyzing the topology of active sites: on the prediction of pockets and subpockets. *J. Chem. Inf. Model.* 50, 2041–2052.
- Wallach, I., Lilien, R.H., 2009. Prediction of sub-cavity binding preferences using an adaptive physicochemical structure representation. *Bioinformatics* 25, i296–i304.
- Wang, R., Fang, X., Lu, Y., Wang, S., 2004. The PDBbind database: collection of binding affinities for protein–ligand complexes with known three-dimensional structures. *J. Med. Chem.* 47, 2977–2980.
- Weill, N., Rognan, D., 2010. Alignment-free ultra-high-throughput comparison of druggable protein–ligand binding sites. *J. Chem. Inf. Model.* 50, 123–135.
- Weisel, M., Proschak, E., Kriegl, J.M., Schneider, G., 2009. Form follows function: shape analysis of protein cavities for receptor-based drug design. *Proteomics* 9, 451–459.
- Xie, L., Bourne, P.E., 2008. Detecting evolutionary relationships across existing fold space, using sequence order-independent profile–profile alignments. *Proc. Natl. Acad. Sci. U.S.A.* 105, 5441–5446.
- Yeturu, K., Chandra, N., 2008. PocketMatch: a new algorithm to compare binding sites in protein structures. *BMC Bioinformatics* 9, 543.
- Yeturu, K., Chandra, N.R., 2011. PocketAlign: a novel algorithm for aligning binding sites in protein structures. *J. Chem. Inf. Model.* 51, 1725–1736.