# SP-Dock: Protein-Protein Docking using Shape and Physicochemical Complementarity

Apostolos Axenopoulos, Petros Daras, *Member*, IEEE, Georgios E. Papadopoulos, and Elias Houstis

**Abstract**— In this paper, a framework for protein-protein docking is proposed, which exploits both shape and physicochemical complementarity to generate improved docking predictions. Shape complementarity is achieved by matching local surface patches. However, unlike existing approaches, which are based on single-patch or two-patch matching, we developed a new algorithm that compares simultaneously, groups of neighboring patches from the receptor with groups of neighboring patches from the ligand. Taking into account the fact that shape complementarity in protein surfaces is mostly approximate rather than exact, the proposed group-based matching algorithm fits perfectly to the nature of protein surfaces. This is demonstrated by the high performance that our method achieves especially in the case where the unbound structures of the proteins are considered. Additionally, several physicochemical factors, such as desolvation energy, electrostatic complementarity, hydrophobicity, Coulomb potential and Lennard-Jones potential are integrated using an optimized scoring function, improving geometric ranking in more than 60% of the complexes of Docking Benchmark 2.4.

**Index Terms** — protein docking, local descriptors, shape complementarity, physicochemical complementarity.

— — — — — — — — — Φ — — — — — — — — —

## 1 INTRODUCTION

PROTEIN-PROTEIN DOCKING has attracted increasing interest during the last years and still remains a hot research topic in Bioinformatics. It deals with the prediction of the conformation and orientation of one protein (ligand) within the binding site of another protein (receptor). Despite the extensive research in protein-protein docking, a complete solution has yet to be achieved due to the large complexity of the problem. It has been proposed that shape complementarity alone cannot achieve highly accurate docking predictions [1]. Since geometric docking is based on approximate surface complementarity, a large number of false-positive predictions may occur. On the other hand, multiple physicochemical factors, such as Coulomb potentials, van der Waals forces, hydrophobicity, etc., can affect the docking predictions, but they need to be appropriately merged with a geometric docking approach. Last but not least, protein interactions can involve significant conformational changes, thus docking techniques should take into account the side-chain and the backbone flexibility. A computational method that will accurately predict protein-protein interactions and within a short time frame would become a valuable tool for biologists and biochemists. Successful docking will predict binding site amino acids crucial for the complex stability, which will assist biochemists perform concrete mutations in order to test their impact for the protein function. In industrial drug design, the economic impact of protein-protein docking is very high, since an accurate and fast docking algorithm will enable rapid

scanning of structural data bases for matches with specific targets, which will speed-up the design process of new drugs and increase productivity. Thus, it is not surprising why protein-protein docking is still a very hot research topic and a lot of effort is put towards investigation of a more accurate computational docking solution, which is expected to provide additional insight into the nature of macromolecular recognition.

## 1.1 Related Work

Shape-based docking approaches can be classified into two main categories: brute-force scanning and local shape feature matching. The former consists of methods based on exhaustive scanning of the transformation space [2], [3]. These begin with a simplified rigid body representation by projecting the protein onto a 3D Cartesian grid; then, they distinguish grid cells according to whether they are near or intersect the protein surface, or are deeply buried within the core of the protein. Complementarity is computed by scoring the degree of overlap between pairs of grids in different relative orientations. To speedup this procedure, FFT-based docking approaches have been introduced [4], [5], [6]. In [7], a method based on Spherical Polar Fourier (SPF) is presented, which calculates rotational correlations using 1D FFTs. ZDOCK introduces a shape complementarity scoring function called Pairwise Shape Complementarity (PSC) [10], which computes the total number of receptor-ligand atom pairs within a distance cutoff. PSC does not rely on excaustive scanning of the entire rotational space resulting in low computation times. One of the most recent approaches of this category is presented in [34]. The so-called F²Dock is an extension of a NFFT-based docking algorithm, where an adaptive search phase (rotational and translational) has been incorporated to achieve faster running times.

Since they are based on exhaustive scanning of translational and rotational space, brute-force methods are able to detect at least one near-native pose in almost every complex. On the other hand, this may lead to an extraordinary big number of candidate docking poses, where, due to the existence of false positives, the near-native poses may not be ranked at the first positions. Such phenomena could be avoided with the use of local shape feature matching methods, which detect points of interest on the protein surfaces. These methods require a representation of the molecular surface, attempting to find critical patches on the surface. Then, pairwise complementarity matching is applied on these patches. In [8], a method based on geometric hashing is presented. Each protein surface is pre-processed to give a list of critical points ("pits", "caps", and "belts"), which are compared, using geometric hashing, to generate a relatively small number of candidate docking poses. The method requires low computation times, however, it does not produce very accurate predictions, since pits, caps and belts do not encode significant shape information. Context Shapes [9] extract local features from the protein surface, which are boolean data structures and correspond to significantly large parts

of the protein surface. Complementarity matching is achieved using boolean operations. The method demonstrates superior performance comparing with the previous one, however, the exhaustive search of relative orientations for each local feature increases the computational time and the memory requirements. In an attempt to deal with the above limitations, in [35], a rotation-invariant shape descriptor is utilised, namely the Shape Impact Descriptor [39], to produce more accurate docking poses with lower computational cost.

In [32], the method LZerD is introduced, which is based on 3D Zernike Descrptors (3DZD). These are a series expansion of a 3D function (i.e. protein surface) allowing for a compact representation of the 3D function. 3DZD are extracted on local patches that are derived on uniformly distributed points of the protein surface. Partial matches are computed using geometric hashing. Surface Histograms (shDock) [33] is a local shape descriptor, which captures the local geometry around a set of two points with given normals on the surface of a protein. The docking pose is obtained automatically by matching two surface histograms. shDock has achieved the best performance among existing methods in Docking Benchmark 2.4, in the bound docking case, i.e. where the candidate proteins are taken directly from the crystallized complex. However, when dealing with the unbound case, the performance of shDock decreases significantly.

Since docking based only on shape complementarity does not provide the best possible results, other non-geometric factors such as desolvation, hydrophobicity, and electrostatics have been also investigated [10], [19]. Recent attempts focus on combining geometric and physicochemical properties in order to produce more accurate predictions. In [20], shape complementarity matching along with knowledge-based potentials, electrostatics, atom desolvation energy, residue contact preferences and Van-derWaals potential are combined, demonstrating remarkable results on a test set of 68 bound and 30 unbound test cases. Although the contribution of each individual non-geometric factor was not assessed in [20], an important conclusion can be drawn: shape complementarity should be combined with physicochemical complementarity to increase the accuracy of docking predictions. F²Dock [34] computes separately shape complementarity scores and electrostatics scores and combines them. This leads to an improvement of shape-only docking in 54% of the complexes of Docking Benchmark 2.0. The most straightforward way to incorporate geometric and non-geometric properties is to represent the final scoring function as a weighted sum of those factors and determine the optimal weights that each factor contributes to the overall scoring. Such an approach is presented in this paper.

Towards the direction of improving existing docking approaches and investigating new approaches, the CAPRI experiment [40] (Critical Assessment of Predicted Interactions) has become an ideal arena for testing docking algorithms. More specifically, in CAPRI, new protein-protein complexes are subjected to structure

prediction before they are published. The complexes are submitted by several predictor groups and they are assessed by comparing their geometry to the original structure. Some of the most well-known docking algorithms, such as PatchDock [8] and ZDock [10], have participated in CAPRI experiment producing acceptable solutions for several CAPRI targets [41][42].

### 1.2 Method Overview and Contributions

In Fig. 1, the block diagram of the proposed method is depicted. The PDB files [11] of the receptor and ligand proteins are given as input and their Solvent Excluded Surfaces (SESs) are extracted. Then, by computing the curvature of the SES, a set of critical points is extracted, which correspond to the centers of small elementary patches (either convex or concave). Each elementary patch is expanded in size in order to cover a wider area producing a Geodesic Surface Patch (GSP). For each GSP an appropriate local shape descriptor is extracted, which uniquely characterizes its shape. During complementarity matching, each GSP that corresponds to a convex (or concave) elementary patch of the receptor protein is matched with all GSPs that correspond to concave (or convex) elementary patches of the ligand protein. As a next step, several neighboring GSPs are grouped together to generate candidate binding regions on the surfaces of the receptor and the ligand. For aligning the two proteins, one candidate region of the ligand is aligned with respect to a complementary candidate region of the receptor using the Iterative Closest Point (ICP) algorithm. At the final step of the algorithm, the aligned poses are scored using both geometric and physicochemical properties. The weights of the scoring function are optimized via training to achieve improved docking results.
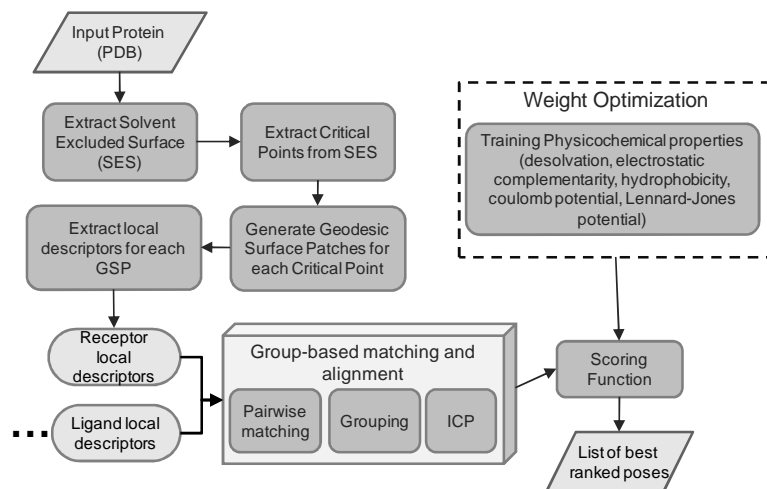


**Fig. 1.** Block diagram of the proposed method

Comparing the proposed SP-Dock (Shape-Physicochemical Docking) method with the approach presented in [35], both are based on local shape feature matching of surface patches corresponding to convex or concave

elementary shape patches. However, numerous novel features are introduced in SP-Dock, which are explained below.

First of all, *a more discriminative local surface descriptor* has been adopted in SP-Dock for patch complementarity matching, instead of the Shape Impact Descriptor (SID) that was used in [35]. Some of the most well-known local descriptors [22] were compared to SID in a dataset of known complexes to select the most appropriate descriptor. The Local Spectral Descriptor [23] has been proven to be the most discriminative among others.

Another notable innovation of SP-Dock is *the group-based matching algorithm*. It introduces a new approach for shape complementarity matching beyond traditional local shape feature matching techniques. It has been inspired by the fact that the shape complementarity between a pair of local surface patches (one from the receptor and one from the ligand), which correspond to a near-native pose, is mostly approximate rather than exact, while at the same time there are plenty of pairs of patches corresponding to non-native poses that have similar or even better shape complementarity than the near-native ones. Thus, existing local shape matching approaches, which rely on single-patch-to-single-patch or two-patch-to-two-patch complementarity matching, may predict a large number of false-positive docking poses and fail to detect near-native poses. The approach presented in this paper intuitively groups neighboring patches from both the receptor and the ligand so as to create larger candidate binding regions. This increases the confidence of a receptor patch to be complementary to a ligand patch, since, according to the grouping criterion, the neighbours of the receptor patch should be complementary to the neighbors of the ligand patch as well. The effectiveness of the proposed approximate complementarity matching is convincingly reflected in the unbound docking case where SP-Dock clearly outperforms similar docking approaches.

Additionally, the paper *proposes the adoption of the Iterative Closest Point (ICP) algorithm for fast alignment of the complementary candidate regions*. ICP has been extensively used for surface registration in 3D reconstruction problems. Although 3D reconstruction involves alignment of surfaces with near-exact similarity, we prove that ICP is also appropriate for aligning surfaces with approximate similarity, as is the case of geometric docking. It is the first time, to the best of our knowledge, that ICP has been used for alignment of protein surfaces. It is also worth mentioning that surface similarity is equivalent to surface complementarity, if the surface of the ligand is turned upside-down [35].

Finally, the paper *assesses the contribution of physicochemical factors to achieve more accurate docking predictions*. Several non-geometric factors, namely the Atom Desolvation Energy, Interface Residue Contact Preferences,

Generic Residue Contact Preferences, Electrostatic Complementarity, Coulomb Potential, Hydrophobicity and Van der Waaks Potential, were computed and combined with the geometric properties into a unified scoring function. These factors have been already discussed in previous works and are summarized in [20]. In this paper, the contribution of each factor is assessed and the optimal weight, with which each factor participates in the scoring function, is estimated using an appropriately selected optimization method. The improvement of docking predictions by combining the geometric with the physicochemical factors, in Docking benchmark 2.4, is impressive.

As it is a local feature matching method, the proposed algorithm shares similarities with the well-known PatchDock method [8]. More specifically, the step of critical points extraction produces similar sparse surface representations for both PatchDock and SP-Dock (although a different algorithm is used in each case to extract the critical points). The geometric scoring step is also similar in both methods, since they generate a 3D distance grid around the receptor, which is accessed by the surface points of the ligand. On the other hand, their surface complementarity matching stages, which constitute core parts of the docking process, are completely different. First of all, SP-Dock does not rely on shape matching of the small convex and concave patches of the sparse surface, but it generates bigger surface patches (the GSPs), which cover a wider surface area around a critical point. These GSPs enclose more significant shape information than the local patches of PatchDock, which is important especially in filtering out a lot of false positive matches. Additionally, instead of the rather simple geometric features that describe the shape of a patch (or a pair of patches) in PatchDock, SP-Dock utilizes state-of-the art local shape descriptors, which makes the method more discriminative in terms of local complementarity matching. Then, in order to match multiple complementary pairs simultaneously and enhance the certainty of pairwise matches, the proposed SP-Dock method does not use geometric hashing (as in PatchDock) but it introduces a new grouping algorithm. This algorithm groups intuitively pairs of complementary GSPs and allows for slight flexibility in the relative positions of the corresponding GSPs within a group. The latter increases the robustness of the method and makes it more appropriate for unbound docking cases, where slight side-chain flexibilities are allowed. The superiority of the proposed method over PatchDock is demonstrated in the experiments section, where SP-Dock outperforms PatchDock especially in the unbound case.

The rest of the paper is organized as follows: in Section 2, the preprocessing phase is described, which includes the surface representation and extraction of local patches, as well as the local shape descriptor extraction for each patch. Section 3 analyzes the new group-based matching and alignment algorithm, while in Sec-

tion 4 the geometric and physicochemical scoring procedure of the candidate docking poses is given. Concerning the physicochemical scoring, the optimization process that assigns a set of weights for each physicochemical factor is provided. Then, in Section 5, the experimental results are presented, where the proposed method is compared to other existing docking approaches. Finally, conclusions are drawn in Section 6.

## 2   PREPROCESSING

This section describes the preprocessing procedure, which involves two phases: during the first phase, an appropriate representation of the molecular surface is generated from the input PDB file, a set of critical points is extracted and a GSP is created for each critical point. The second phase involves the extraction of low-level geometric descriptors for each GSP, which uniquely characterize its shape.

### 2.1 Surface Representation and Extraction of Local Patches

Extraction of 3D shape descriptors from a protein initially requires an appropriate representation of its 3D structure. Several representations have been proposed so far, namely the volumetric representation [7], the Solvent Excluded Surface (SES) [12], Sparse Surface [14] and Alpha Shapes [21]. In this paper, the SES method has been selected, which produces a 3D triangulated surface of the protein. In order to generate a SES, the Maximal Speed Molecular Surface (MSMS) [13] algorithm has been utilized.



**Fig. 2.** A Geodesic Surface Patch (GSP) is centered at the critical point **p**.

Computation of critical points on the SES offers a sufficient approximation of the protein surface and constitutes a preliminary step that is followed by almost all the local shape feature matching approaches. In this paper, a method for generating critical points based on the local curvature of the surface has been followed. This approach has been introduced in [35], it is applied directly to the 3D triangulated mesh and it is applicable to all types of triangulated meshes. The extracted critical points are the centers of concave and convex

regions of the molecular surface. A detailed description of the algorithm is available in [35].

For each critical point, a GSP is created (Fig. 2), which spreads over a wider surface area around that point. More specifically, a GSP consists of all points of SES whose geodesic distance from the critical point is less than a predefined threshold ($G_{max}$). GSP differs from the Extended Surface Patch (ESP) that was defined in [35] in the sense that the latter uses the Euclidean distance as an initial threshold, while the geodesic distance is used only as a post-filter to remove unconnected surface parts. However, it was proven experimentally that the GSP-based approach achieves better accuracy than the ESP-based approach. It was also experimentally found that an optimal value for $G_{max}$ is 16Å.

## 2.1 Local Descriptor Extraction

In protein-protein docking problems, local-shape-feature-based methods rely on pairwise matching of local surface regions between the receptor and the ligand. The most complementary surface regions are, then, selected as candidate poses. The approach presented in this paper uses shape similarity descriptors to measure surface complementarity. It has been proven in [35] that complementarity matching of surface patches can be reduced to a similarity matching problem, if the inner surface part of the ligand patches is treated as outer and vice versa. This concept is illustrated in Fig. 3, where a pair of complementary surface patches of the 1CGI complex is depicted. In Fig. 3a and 4b, the outer parts of the surface patches are shown. In Fig. 3c, the inner part of the ligand patch is depicted. It is obvious that the latter patch has similar shape with the patch in Fig. 3a.

In the approach presented in this paper, the GSPs of the receptor that correspond to convex (or concave) critical points are matched with the GSPs of the ligand that correspond to concave (or convex) critical points. The matching relies on the shape complementarity between the GSPs. Unlike the method in [35], where only the SID was used for shape similarity, in this paper, three local shape descriptors have been tested in order to find the most appropriate one for our problem. The most well-known local shape descriptors for 3D meshes have been presented in SHREC 2011 (Shape Retrieval Contest on Non-rigid 3D Watertight Meshes) [22]. Two local shape descriptors that achieved high accuracy in SHREC 2011 have been tested in our docking framework and compared with the Shape Impact Descriptor. The selection of descriptors has been performed according to the following criteria:

*Rotation Invariance*: the local patches of the two protein surfaces have arbitrary orientations. In order to be matched, they should be either aligned or a rotation-invariant descriptor can be used. Alignment is usually based on the directions of the patch normals; however, the latter do not provide a robust measure, which

leads in inaccurate alignment. Rotation-invariant descriptors are able to match two surface patches irrespective of their pose.

*Compactness and fast extraction*: local descriptors are applied to a relatively big number of surface patches. This implies that descriptor extraction and pairwise matching of single patches should be extremely fast. Fast matching is achieved by using very compact descriptor vectors (usually up to 100 values). Thus, shape descriptors with high computational complexity are not appropriate in our case.

Finally, the candidate shape descriptor *should be applied to the surface of the protein*, which automatically excludes descriptors based on the volume of a 3D object. The descriptors that will be described in the sequel, for the sake of completeness, fulfil all the above requirements. A more detailed description is provided for the first descriptor (Local Spectral Descriptor) since it is the one that has been eventually chosen.
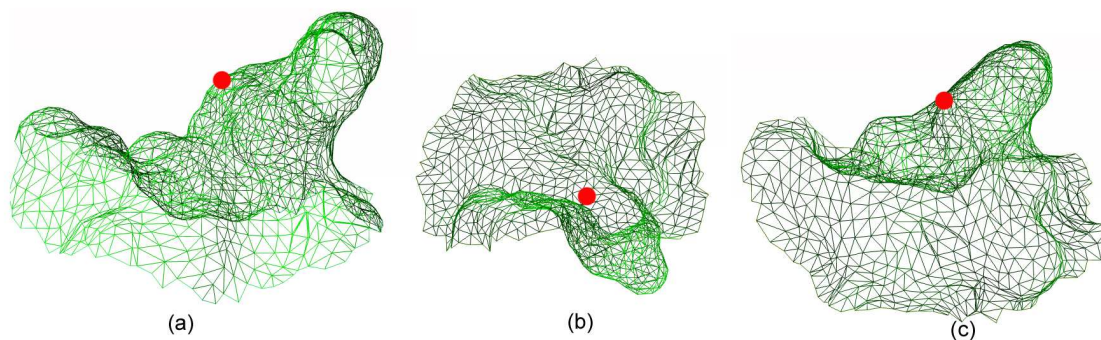


**Fig. 3.** a) a surface patch of the receptor of the 1CGI complex (large protrusion); b) a surface patch of the ligand (deep cavity); c) the patch of b) turned upside down so that the inner surface is visible. The patches in a) and c) have approximately similar shapes, thus, patches a) and b) are complementary.

*Local Spectral Descriptor*

This local descriptor has been proposed in [23] for retrieval of non-rigid 3D meshes. It is based on the extraction of geometric descriptors from a surface patch $P_i$ centered around a sample point $\mathbf{p}_i$ on the mesh. The method computes the Fourier spectra of the patch by projecting the geometry on the eigenvectors of the Laplace-Beltrami operator (LBO). LBO is defined as the divergence of the gradient for functions that are defined over manifolds. The eigenvalues and eigenvectors of this operator satisfy the following equation:

$$-\mathbf{Q}\mathbf{h}^k = \lambda_k \mathbf{D}\mathbf{h}^k \qquad (1)$$

where $\lambda_k$ is the $k$th eigenvalue, $\mathbf{h}^k$ is the $k$th eigenvector $\mathbf{h}^k = [H_1^k, \ldots H_m^k]$ and $m$ is the total number of vertices of the surface patch. $\mathbf{D}$ is the Lumped Mass matrix and $\mathbf{Q}$ is the Stiffness matrix that are described in [24]. In order to compute the $k$th spectral coefficient, the inner product between the patch surface and the $k$th eigenvector is calculated:

$$\widetilde{x}_k = <\mathbf{x}, \mathbf{h}^k> = \sum_{i=1}^{m} x_i D_{i,i} H_i^k \qquad (2)$$

where $x_i$ is the x-coordinate of the i$^{th}$ vertex of the surface patch. Similar equations hold for $\widetilde{y}_k$ and $\widetilde{z}_k$, which correspond to the y and z coordinates, respectively. Finally, the $k^{th}$ spectral coefficient is given by:

$$c_k = \sqrt{(\widetilde{x}_k)^2 + (\widetilde{y}_k)^2 + (\widetilde{z}_k)^2} \qquad (3)$$

The Local Spectral Descriptor for patch $P_i$ around point $\mathbf{p}_i$ is the vector $\mathbf{c}^i = [c_1^i, \ldots c_n^i]$, where $k=1,\ldots n$ the first spectral coefficients. The dimensionality of the descriptor has been experimentally found to be $n=50$.

*ShapeDNA*

The ShapeDNA descriptor has been proposed in [25] for non-rigid shape analysis. It presents similarities with the Local Spectral Descriptor in the sense that they are both based on solving the eigenvalue problem of the Laplace-Beltrami operator. However, in ShapeDNA the descriptors are the first smallest $N$ eigenvalues, which are the solutions of the Laplacian eigenvalue problem (1), while in Local Spectral Descriptors, the descriptors are extracted by projecting the geometry of the surface on the eigenvectors of the Laplace-Beltrami operator. In general, a small number of egenvalues (10 to 15) provide a sufficient number of descriptors. In our experiments, $N=14$ was experimentally found to give the optimal results. A detailed description of the ShapeDNA descriptor is available in [25].

*Shape Impact Descriptor (SID)*

SID was firstly introduced in [15] and extended in [39] as a shape similarity measure for 3D objects. The key idea of SID is the description of the resulting phenomena that occur by the insertion of the 3D object in the space. It is expected that similar objects will result in similar physical phenomena. Some obvious selections of surrounding fields are the traditional electrostatic force field and the Newtonian force field. Any 3D object can be considered as a distributed mass (or a distributed charge) with a specific distribution, resulting in a static field around it. SID is composed of three major histograms created by a) the field potential values, b) the field density Euclidean norms and c) the radial component of the field density, computed in points that are equidistant from the object surface. The computation of histograms involves only relative distances, thus the descriptor is rotation-invariant. A more detailed description of SID is available in [15], [39].

## 3 GROUP-BASED MATCHING AND ALIGNMENT

Most of the existing local shape feature docking approaches are based on either one-patch-to-one-patch or two-patch-to-two-patch complementarity matching between the local patches of the receptor and the ligand.

Then, the most complementary pairs are aligned in order to produce the final poses as follows: a) for methods based on single-patch matching [9], [35], the ligand is translated so that the patch center of the ligand patch coincides with the patch center of the receptor patch; b) for methods based on two-patch matching [8], [36], the ligand is translated so that its first critical point coincides with one of the receptor and then it is rotated so that its second critical point coincides with the second critical point of the receptor. These approaches suffer from the following limitations:

*Alignment is not always accurate*: the patch centers (critical points) of the receptor and the ligand do not always coincide with their real contact points, producing docking poses that may be far from the near-native poses. An approach that is usually followed is to increase the number of samples on the protein surfaces, which in turn dramatically increases the computation time.

*Low shape complementarity between surface patches:* this is due to the fact that shape complementarity in protein surfaces is mostly approximate rather than exact (Fig. 4). This results in relatively low complementarity scores of patches that correspond to near-native poses comparing to scores of patches that correspond to non-native poses, causing a high number of false positive predictions.

Instead of applying single-patch or two-patch matching, in this paper, a novel approach is presented, where several neighboring GSPs are grouped together to generate candidate binding regions on the surfaces of the receptor and the ligand. This increases the confidence of a receptor patch to be complementary to a ligand patch, since, according to the grouping criterion, the neighbours of the receptor patch should be complementary to the neighbors of the ligand patch as well.

### 3.1 Creating groups of neighboring complementary GSPs

The steps of the group-based matching algorithm are summarized in Fig. 4. Let $N_R$, $N_L$ be the GSPs of receptor and ligand, respectively, and $D_R^i$, $D_L^i$ their corresponding local shape descriptors, where $i=1,\ldots N_R$ (or $N_L$). Let, also, the function that represents the convexity or concavity of a GSP be:

$$Cur(i) = \begin{cases} 1, \text{if } i \text{ is a convex GSP} \\ -1, \text{if } i \text{ is a concave GSP} \end{cases} \qquad (4)$$

Each receptor GSP $i$ is matched with all ligand GSPs $j$ of different type ($Cur(i) \cdot Cur(j) = -1$). A dissimilarity metric is calculated for each pair $(i,j)$ as:

$$Dissimilarity(i, j) = dis(D_R^i, D_L^j) \qquad (5)$$

where *dis()* is an appropriate distance metric applied on the descriptor vectors $D_R^i$ and $D_L^i$. The distance

metric depends on the selected descriptor. In our experiments, the Manhattan distance ($L_1$), the Euclidean

distance ($L_2$) and the diffusion distance [16] have been selected for matching of the Local Spectral Descriptors,

the ShapeDNA descriptors and the SID descriptors, respectively. After computation of the dissimilarities, the

GSPs of ligand are sorted with respect to similarity to the receptor GSP $i$ and the $k$-first are selected to form a

ranked list $RL_R^i$. It is worth mentioning that similarity of the local descriptors is equivalent to complemen-

tarity as it was explained in Section 2.1.

```
INPUT:
  N_R, N_L the GSPs of receptor and ligand, respectively
  D_R^i, D_L^i  their local shape descriptors, i=1,…N_R (or N_L)
OUTPUT:
  G = {G_1, G_2, …, G_M} the set of patch groups
ALGORITHM:
  Set G ← {}
  For each receptor GSP i
    For each ligand GSP j
      If Cur(i)·Cur(j) = -1
        Calculate dis(D_R^i, D_L^j)
    Sort GSPs of ligand
    Keep k-first ligand GSPs and create ranked list RL_R^i
    For each ligand GSP j of RL_R^i
      For each group G_k ∈ G
        If pair (i,j) fulfils Grouping Criterion for G_k
          Then add pair (i,j) to G_k
      If (i,j) not added to any group
        Then create new group and add to G
```

**Fig. 4.** The Group-based Matching algorithm.

The output of the algorithm is a set of groups $G$, which is defined as follows:

$$G = \{G_1, G_2, \ldots, G_M\} \qquad (6)$$

where $G_k = \{(I_R^1, I_L^1), (I_R^2, I_L^2), \ldots, (I_R^g, I_L^g)\}$ is a group that consists of the pairs $(I_R^i, I_L^i)$, $I_R^i$ is the index of

a receptor GSP ($I_R^i = 1, \ldots N_R$) and $I_L^i$ is the index of a ligand GSP ($I_L^i = 1, \ldots N_L$). In order for the above

pairs to form a group, the following grouping criterion must hold:

$$d_{Geod}(I_R^i, I_R^j) < gThres \qquad i, j \in [1, \ldots g] \text{ and}$$

$$d_{Geod}(I_L^i, I_L^j) < gThres \qquad i, j \in [1, \ldots g] \qquad (7)$$

where $d_{Geod}$ is the geodesic distance between two GSPs of either the receptor or the ligand and *gThres* an

appropriately selected geodesic threshold.



**Fig. 5.** (a) and (b): a pair of complementary groups, one from receptor and one from ligand, respectively, for the 1AVX complex. The first patch of the ligand is complementary with one patch of the receptor and the second patch of the ligand is complementary with three patches of the receptor. (c) and (d) the corresponding point clouds that are given as input to ICP for the alignment step.

The candidate pairs of a group $G_k$ are created by combining each receptor GSP $i$ with the $k$ most similar ligand GSPs, i.e. those included in the ranked list $RL_R^i$. Consequently, a group $G_k$ consists of neighboring receptor and ligand GSPs and each receptor GSP of the group is complementary with at least one ligand GSP of the group. This is illustrated in Fig. 5 (a) and (b). In Fig. 5 (a), the group of the receptor consists of four patches, whose centers are represented by blue spheres. Three of these patches are complementary with one patch of the ligand group, while the second patch of the ligand group is complementary with the fourth patch of the receptor. In general, the proposed grouping algorithm allows many–to-many correspondences of local patches increasing the confidence of complementarity between pairs of groups.

The ranges of *gThres* and *k* values have been experimentally determined to be at the ranges of 4-6 Å and 8-11, respectively. Higher values of *gThres* result in a smaller number of larger groups, while lower values of *gThres* result in a larger number of smaller groups. In the former case, the algorithm may fail to predict some near-native poses, while in the latter case, more false positive results may occur. Similar observations are made for *k*, if we decrease it (*k*<8) or increase it (*k*>11), respectively. For the experiments presented in Section

5, $gThres$=5Å and $k$=10 lead to the best results.

The above process produces $M$ groups, which will be given as input to alignment and final scoring function and will result in $M$ predicted docking poses. There is no need to define an additional cutoff threshold, as required in [35] to select the first most complementary pairs of patches, since the average number of $M$ is 2500-3000, while the average number of the patch pairs in [35] is ~500000. Another advantage of the proposed grouping algorithm, comparing with the method in [35], is that patch pairs that lead to almost the same docking poses are not taken as separate cases but are grouped together (due to the neighborhood criterion). This results in a significantly smaller number of false positive predictions, which improves the final rank of the near-native predictions.

### 3.3 Alignment of groups

During the alignment phase, a rigid transformation of the ligand is computed for each of the $M$ groups created using the group-based matching algorithm. Let $C_R^i$ (or $C_L^i$) be the point cloud that consists of all points of the $i$th receptor GSP (or ligand GSP). The receptor point cloud $GC_R^k$ of group $G_k$ is given by:

$$GC_R^k = C_R^{I_R^1} \bigcup \ldots \bigcup C_R^{I_R^g} \tag{8}$$

i.e. it is the union of the receptor point clouds $C_R^i$ of the GSPs within group $G_k$. The ligand point cloud $GC_L^k$ of group $G_k$ is computed in a similar manner. The required rigid transformation translates and rotates $GC_L^k$ so as to optimally fit to $GC_R^k$. Then, the same rigid transformation is applied to the entire ligand molecule in order to compute the final score of the predicted pose.

The optimal alignment of two point clouds is a surface registration problem. One of the most well-known techniques for surface registration is the Iterative Closest Point (ICP) algorithm [37]. Let $GC_R = \left\{ \mathbf{c}_R^1, \mathbf{c}_R^2, \ldots, \mathbf{c}_R^{n_R} \right\}$ and $GC_L = \left\{ \mathbf{c}_L^1, \mathbf{c}_L^2, \ldots, \mathbf{c}_L^{n_L} \right\}$ be the two point clouds to be aligned, and $\left\| \mathbf{c}_L^j - \mathbf{c}_R^i \right\|$ be the Euclidean distance between point $\mathbf{c}_R^i \in GC_R$ and $\mathbf{c}_L^j \in GC_L$. Let also $CP(\mathbf{c}_L^j, GC_R)$ the closest point of $GC_R$ to the point $\mathbf{c}_L^j$. It is useful to launch ICP with an initial estimate $T^0$ of the rigid transformation. This is usually computed by translating the median point of $GC_L$ to coincide with the median point of $GC_R$ and rotating $GC_L$ so that its average normal (the average of the normals of all points $\mathbf{c}_L^j$) is aligned with the average normal of $GC_R$. Then, an iterative process is repeated ($t$=1,...,$t_{max}$ iterations) until convergence. For the $t$th iteration, the set of correspondences is computed by:

$$Corr^t = \bigcup_{i=1}^{n_L} \left\{ \left( \mathbf{c}_L^i, CP\left(T^{t-1}(\mathbf{c}_L^i), GC_R\right) \right) \right\} \tag{9}$$

Then, the new transformation $T^t$ that minimizes the mean square error between point pairs in $Corr^t$ is computed. In Fig. 5 (c) and (d), the point clouds that correspond to the pair of complementary groups (a) and (b) are depicted. These are given as input to ICP at the alignment phase. In Fig. 6, two results of alignment using ICP are provided for the 1AVX and 1HIA complexes. It is obvious that a highly accurate alignment is achieved.



(a)                                                                    (b)

**Fig. 6.** Aligment results using ICP for the (a) 1AVX and (b) 1HIA complexes. A surface representation is used for the receptor and a backbone representation for the ligand. The blue line corresponds to the original position of the ligand and the magenta line corresponds to the pose predicted using ICP.

## 4 SCORING OF CANDIDATE POSES

In this section, the final stage of the proposed SP-Dock method is described, which involves scoring of the candidate poses that were produced during the group-based matching and alignment phase. Apart from the geometric complementarity, the effect of several (non-geometric) physicochemical factors on the accuracy of docking predictions is also investigated. The final scoring function is a weighted sum of the geometric score and the scores obtained from each separate physicochemical factor. The predicted docking poses are sorted in descending order, with the poses of the highest overall score to appear first in the ranked list.

### 4.1 Geometric Scoring

For the geometric scoring of each candidate pose, the 3D distance grid, which was presented in [35], is used. The receptor protein and its surrounding space is represented by a 3D function $DT(i, j, k)$:

$$DT(i, j, k) = \begin{cases} 0, & \text{if at least one surface point lies inside the voxel} \\ < 0, & \text{if the voxel lies inside the molecule} \\ > 0, & \text{if the voxel lies outside the molecule} \end{cases} \tag{10}$$

The absolute value of each voxel corresponds to the Euclidean distance from the closest surface point.

Then, the distance grid is divided into 6 shells (Table I) according to the distance from the molecular surface.

The shell ranges have been experimentally determined. It is worth mentioning that the shell ranges are similar to the ones obtained by the PatchDock method [8] (with a 0.2 Å shift), which was expected since the geometric scoring step is quite similar in both methods.

**Table I:** The shells in which the distance grid is divided.

| Shell 1 | [1.2, ∞) | The range (in Å) of the first shell of the distance grid |
|---|---|---|
| Shell 2 | [-1.2, 1.2) | The range of the second shell of the distance grid |
| Shell 3 | [-2.4, -1.2) | The range of the third shell of the distance grid |
| Shell 4 | [-3.8, -2.4) | The range of the fourth shell of the distance grid |
| Shell 5 | [-5.2, -3.8) | The range of the fifth shell of the distance grid |
| Shell 6 | [−∞, -5.2) | The range of the sixth shell of the distance grid |
| $a_{1-6}$ | 0, 1, -1, -18, -190, -10000 | The values of the weights in the scoring function (13) |

For each of the docking poses predicted in the group-based matching and alignment phase, the translated and rotated ligand $L$ enters the 3D distance grid of the receptor $R$. $L$'s surface points access the voxels of the 3D grid and are assigned a value according to the distance from $R$'s molecular surface. The score $E_S$ of the transformation is given by:

$$E_S = \sum_{i=1}^{6} a_i N_i \qquad\qquad (11)$$

where $N_i$ is the number of $L$ points in shell $i$ of the distance grid and $a_i$ is the weight of the $i$-th shell (Table I).

After the ICP-based alignment step (Section 3.3), $M$ different poses of the ligand are taken (equal to the $M$ generated groups $G$). For the pose that corresponds to a group $G_k$, an additional refinement step is applied, which involves +/-2 Å translation of the ligand towards the direction of $GC_L^k$'s average normal (Section 3.3) and +/-25$^O$ rotation of the ligand around $GC_L^k$'s average normal. This results in a total of 9 poses for each group $G_k$, which is significantly faster than the method in [35] that requires 1872 different poses for each pair of complementary ESPs. The reason for taking only 9 poses is that the final transformation has been already approximated using ICP, thus, only a slight refinement is required. Taking also into account that ICP is significantly faster than distance-grid-based scoring, the significance of ICP in our approach is obvious.

The computation time required for the distance-grid-based scoring is proportional to the size and the resolution of the ligand's surface. In order to achieve low computation times, two different resolutions of the ligand SES are used: a) the low-resolution surface with point density of 1 point per Å$^2$ and b) the high-resolution surface with point density of 4 points per Å$^2$. The low-resolution surface is used to score all 9 poses for each group $G_k$, and the high-resolution surface is used for the best among the 9 poses.

## 4.2 Physicochemical Factors Assessment

Among the several non-geometric physicochemical factors that may affect the accuracy of protein-protein docking, the following have been assessed in this paper:

*Atom Desolvation Energy (ADE)*: the atomic contact potential, which is used to estimate the desolvation energy for the replacement of protein-water contacts with protein-protein contacts, is given by [26]:

$$E_{ADE} = \sum_{i=1}^{N} \sum_{j=1}^{M} e_{ij} \tag{12}$$

where $e_{ij}$ is the non-scaled contact value of a contact between atom $i$ from receptor and atom $j$ from ligand. The contact values are summed over all atoms of receptor that are within 6 Å distance to at least one atom of ligand and vice-versa.

*Interface Residue Contact Preferences (RCP)*: these are volume-normalized pair probabilities that represent the pairing preferences of aminoacids at the protein-protein interface [27]:

$$E_{RCP} = \sum_{i=1}^{N} \sum_{j=1}^{M} \frac{e_{ij}}{r_{ij} + 1.5} \tag{13}$$

where $e_{ij}$ is the volume-normalised pairing preference between aminoacid $i$ from the receptor and aminoacid $j$ from ligand and $r_{ij}$ is the distance between their corresponding $C_{\beta}$ atoms. The value 1.5 has been added to avoid unrealistic close contacts.

*Generic Residue Contact Preferences (GCP)*: it is calculated in a similar manner as in the case of Interface Residue Contact Preferences. In this case, $e_{ij}$ is the pairing probability of aminoacids in protein structures [28].

*Electrostatic Complementarity (EC)*: the electrostatic complementarity at the interface is calculated by [29]:

$$E_{EC} = \sum_{i=1}^{N} \sum_{j=1}^{M} e_{ij} \tag{14}$$

$$e_{ij} = \begin{cases} 0, & \text{if } r_{ij} > r_{\max} \\ A(p_i, p_j), & \text{if } r_{ij} \le r_{\max} \end{cases} \tag{15}$$

where $A(p_i, p_j)$ are statistical interaction energies, $r_{ij}$ is the distance between atoms $i,j$, and $r_{\max}$=4Å if both atoms are apolar and 3.4 Å otherwise.

*Coulomb Potential (CP)*: Coulomb potential is given by the following equation:

$$E_{CP} = \frac{q_i q_j}{(r_{ij} + c)^2} \tag{16}$$

where $q_i$, $q_j$ are the partial charges of each atom. The constant $c$ is equal to 1.5 Å to avoid strong influence of very close atoms [30].

*Hydrophobicity (HP)*: it is calculated using the following equation [31]:

$$E_{HP} = \frac{hh}{hh + pp + hp} \tag{17}$$

where *hh* is the number of contacts between hydrophobic atoms, *pp* is the number of contacts between two polar atoms and *hp* is the number of contacts between polar and hydrophobic atoms.

*Van-der-Waals Potential (vdW)*: here, the modified 6-12 Lennard-Jones Potential is calculated by :

$$E_{vdW} = \begin{cases} \varepsilon_{ij}\left( \dfrac{\sigma_{ij}^{12}}{r_{ij}^{12}} - 2\dfrac{\sigma_{ij}^{6}}{r_{ij}^{6}} \right), & \text{if } r_{ij} > 0.6\sigma_{ij} \\[2ex] \varepsilon_{ij}(A + (r_{ij} - 0.6\sigma_{ij})B), & \text{otherwise} \end{cases} \tag{18}$$

$$A = \frac{\sigma_{ij}^{12}}{0.6\sigma_{ij}^{12}} - 2\frac{\sigma_{ij}^{6}}{0.6\sigma_{ij}^{6}}, B = -12\frac{\sigma_{ij}^{12}}{0.6\sigma_{ij}^{13}} + 12\frac{\sigma_{ij}^{6}}{0.6\sigma_{ij}^{7}} \tag{19}$$

where $\sigma_{ij}$ is the sum of van-der-Waals radii and $r_{ij}$ is the distance between atoms *i* and *j*. The potential is calculated for atoms with interatomic distances of less than 6 Å.

It should be stressed that the above factors are not the only ones that affect the protein interactions. A variety of additional physicochemical properties could be also found and integrated into a compound scoring function. An extensive survey on all possible factors is not within the scope of this paper, but it constitutes a significant challenge for future work. The factors presented above are also summarized in [20], where it is stated that they are able to improve docking predictions when merged with geometric docking. However, no information about the contribution of each separate factor is given in [20]. In this paper, an assessment of each factor is provided through an appropriate optimization method. More specifically, the overall score of each docking pose is given as the weighted sum of the geometric score (Section 4.1) and the scores of the factors described above. The weights are optimized on a training dataset (59 test cases of Docking Benchmark v1.0 [18]) using Particle Swarm Optimization (PSO) [38]. The overall scoring function is given by:

$$Score_{Total} = w_s E_s + w_{ADE} E_{ADE} + w_{RCP} E_{RCP} + w_{GCP} E_{GCP} + w_{EC} E_{EC} + w_{HP} E_{HP} + w_{CP} E_{CP} + w_{vdW} E_{vdW} \tag{20}$$

PSO is a global optimization algorithm, similar to a genetic algorithm, motivated by social behavior of organisms such as bird flocking and fish schooling. PSO iteratively tries to improve a candidate solution with respect to a given measure of quality (fitness function). PSO establishes a population (swarm) of candidate solutions, known as particles that move around in the search space, and are guided by the best found positions, updated when better positions are found by the particles.

In our approach, the population of candidate solutions is the 8 weights *w* of (22), which can take arbitrary

real values within the range [0,1]. The values of scores $E$ (22) have been normalized so that their value range is within [0,100]. The key of success of the PSO method is the selection of an appropriate fitness function. In our experiments, two fitness functions are determined. The first one is the *Average Precision* of the first-ranked *hit* for all the complexes of the training dataset. The *Precision* of the first-ranked hit for one complex is given by:

$$F_1 = \frac{n_{hit}}{n_{retrieved}} = \frac{1}{n_{retrieved}} \tag{21}$$

where $n_{hit}$ is the total number of hits (i.e. near-native poses) that are retrieved and $n_{retrieved}$ is the total number of predicted docking poses that are retrieved. If we select $n_{retrieved}$ to be equal to the number of retrieved poses until the first hit is retrieved, then the numerator of (23) is equal to 1. As an example, if the first hit is retrieved in the fourth position, then the Precision for this complex is $F_1$=0.25, or 25%. The Average Precision of the first-ranked hit provides an acceptable metric to be used as a fitness function, however, it suffers from the following limitation: it favors those complexes in which the first ranked hit is retrieved at the first positions (1 - 10), while the complexes, in which the first hit has rank >100, have insignificant contribution to the calculation of the Average Precision. In other words, an improvement of the hit's rank from 200 to 100 contributes with 0.01 to the average precision, while an improvement from 2 to 1 contributes with 1 to average precision. This is not desired since in the former case the improvement is much more significant and should contribute more to the average precision.

To overcome the above limitation a new fitness function was determined, which is given by:

$$F_2 = \frac{1}{N_C} \sum_{i=1}^{N_C} \frac{rank_S^i - rank_{SP}^i}{N_{Poses}^i} \tag{22}$$

where $N_C$ is the number of complexes of the training dataset, $rank_S^i$ is the rank of the first hit (of complex $i$) that is retrieved using only shape complementarity, $rank_{SP}^i$ is the rank of the first hit using the weighted score (20) and $N_{Poses}^i$ is the number of predicted poses of complex $i$.

## 5 RESULTS AND DISCUSSION

The proposed (Shape-Physicochemical) SP-Dock method was experimentally evaluated using the protein-protein docking benchmark v2.4 [17], which consists of 84 known complexes (63 rigid-body cases, 13 cases of medium difficulty, and 8 difficult cases). To evaluate the performance of the method, for each complex, the receptor and ligand are separated from each other and the ligand is translated and rotated arbitrarily. In or-

der to increase the confidence of the results, the docking algorithm has been repeated for three different initial rotations of the ligand. Eventually, we observed that these three arbitrary rotations produced only very slight modifications on the final poses (mainly due to the outcome of the ICP algorithm), which did not affect the final rankings. Thus, the result of only one of the three iterations (the first one) is presented in the following subsections. The docking algorithm described in the previous sections is applied to generate a set of candidate poses of the ligand. The predicted pose of the ligand is compared to its original pose in the complex in terms of interface Root Mean Square Deviation (*iRMSD*):

$$iRMSD = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left\| \mathbf{a}_i^p - \mathbf{a}_i^o \right\|^2} \tag{23}$$

where $\mathbf{a}_i^p$ is the $i^{th}$ interface $C_a$ atom (in x, y, z coordinates) of the ligand in the predicted pose and $\mathbf{a}_i^o$ is the corresponding $C_a$ atom of the ligand in the original pose (crystallized complex). Interface $C_a$ atoms of the ligand are those that are within the distance of 10 Å from the receptor. A predicted pose is called a hit if the iRMSD between the ligand in that pose and the ligand in the original complex is less than 2.5 Å.

## 5.1 Evaluation of Local Descriptors and Physicochemical Factors

The choice of the appropriate local shape descriptor is crucial for the accuracy of the docking predictions. In Fig. 7(a), a comparison of the three shape descriptors of Section 2.1 is given. The diagram depicts the distribution of the ranks of the first prediction within 2.5 Å of the native complex structure. As an example, in the case of the Local Spectral Descriptor (blue column), the value of the first bin is 17, which means that in 17 out of 84 complexes of Docking Benchmark 2.4 the algorithm returned a hit at the first position. Similarly, the value of second bin is the number of complexes where a hit is predictred within the first five positions and so on. The value of the last bin is 76, i.e. in 76 complexes the algorithm returned a hit within the first 3600 positions, thus failed only in 8 cases. It is also clear from Fig. 7(a) that the Local Spectral Descriptor produces better results than ShapeDNA and SID, thus it was eventually selected for our SP-Dock method.

In Fig. 7(b), the effect of using physicochemical properties along with shape complementarity is demostrated. The red and green columns depict the ranks distribution in Docking Benchmark 2.4 using the unified scoring function of (20) optimized with the fitness functions $F_1$ (21) and $F_2$ (22), respectively. In both cases, the use of physicochemical properties improves the docking predictions of the shape-only approach. However, the weighted function optimized with $F_1$ demonstrates better improvement at the first ranks (1-10), while the weighted function optimized with $F_2$ is better at the higher ranks (100-2000). This makes sense taking into

account the fact that $F_1$ favors those complexes in which the first ranked hit is retrieved at the first positions, as explained in Section 4.2. Eventually, the results obtained with the fitness function $F_2$ were selected since they provide better overall improvement over the shape-only approach (65.3% improvement comparing with 57.3% obtained with $F_1$).



(a)                                                                                      (b)

**Fig. 7.** Distribution of the ranks of the first prediction within 2.5 Å of the native complex structure for all test cases in Docking Benchmark 2.4, a) for different local shape descriptors and b) comparison of our method using only shape complementarity with our method using shape and physicochemical complementarity and different fitness functions for weight optimization of the scoring function (F-1 is the precision of the first-ranked hit and F-2 is the function described in (22)).

The weights of (20) that produced the results presented in Fig. 7(b) have been optimized by training on a dataset (59 complexes) of the docking benchmark v1.0 [18]. These weights for both $F_1$ and $F_2$ fitness functions are depicted in Table II.

**Table II:** The optimized weights for each factor in (20) obtained by the two fitness functions and Particle Swarm Optimization.

| Fitness Function | $w_S$ | $w_{ADE}$ | $w_{RCP}$ | $w_{GCP}$ | $w_{EC}$ | $w_{HP}$ | $w_{CP}$ | $w_{vdW}$ |
|---|---|---|---|---|---|---|---|---|
| $F_1$ | 0.114 | 0.158 | 0.008 | 0.006 | 0.143 | 0.01 | 0.05 | 0.51 |
| $F_2$ | 0.231 | 0.013 | 0.007 | 0.005 | 0.104 | 0.079 | 0.047 | 0.512 |

## 5.2 Comparison with PatchDock, ZDock, LZerD, shDock and F²Dock

The results of the proposed method were compared to those of the following five methods: a) Local 3D Zernike descriptor-based Docking (LZerD) [32], b) Surface Histograms (shDock) [33], c) Fast Fourier Protein-Protein Docking (F²Dock) [34], d) PatchDock [8] and e) ZDock [10]. These are the most recent works related to geometric protein-protein docking and they have achieved the best docking accuracy reported so far. In our experiments, both *R-bound/L-bound* and *R-unbound/L-unbound* cases were evaluated. It is worth mentioning that the last two methods, PatchDock and ZDock, have participated in the CAPRI experiment, a well-established arena for testing docking algorithms.

Two variations of the proposed method have been tested: the first (S-Dock) is based only on geometric properties, while in the second (SP-Dock), both shape and physicochemical properties are integrated as described in Section 4.2. An analytic comparison of the proposed method with the other five methods, in the *R-bound/L-bound* case, is available in supplemental material (Appendix 1, Table IX). Summing up the results of Table IX**Σφάλμα! Το αρχείο προέλευσης της αναφοράς δεν βρέθηκε.**, the proposed approach failed to return a hit in 8 out of 84 cases, while shDock failed in 10 cases, LZerD in 25, F²Dock (S) in 22, F²Dock (S-E) in 25, PatchDock in 39 and ZDock in 23 cases. SP-Dock was ranked first in 40 out of 84 cases, far beyond the shDock method that was ranked first in 32 cases. In Table IV, the number of cases, where at least one hit is found at different docking thresholds, is presented for all methods. In the R-bound/L-bound results, the proposed SP-Dock method outperforms all other five methods for thresholds 1, 1000, 2000 and 3600, while for thresholds 5, 10, 100 only shDock outperforms the proposed method. In Table III, the win-tie-loss-failure records for the proposed method versus shDock, LZerD, F²Dock, PatchDock and ZDock is presented. Comparing our shape-only approach (S-Dock) with shDock, S-Dock returns a better ranked hit in 31 cases, whereas shDock returns a better hit in 29 cases. The methods tie in 20 cases, and both fail in 4 cases. Comparing against LZerD, F²Dock (S), PatchDock and ZDock, the proposed method clearly outperforms them; it has 50-21 win-loss record against LZerD, 66-8 win-loss record against F²Dock (S), 54-17 win-loss record against PatchDock and 56-21 win-loss record against ZDock (S). The accuracy of our method is further improved when shape complementarity is merged with physicochemical complementarity (Table III). Note that S-Dock is compared to the shape-only version of F²Dock, while SP-Dock is compared to the F²Dock (S-E), where electrostatics are merged with the geometric properties.

**Table III:** R-bound/L-bound: the win-tie-loss-failure records for the proposed method versus shDock, LZerD and F²Dock.

| S-Dock vs | Win | Tie | Loss | Both fail |
|---|---|---|---|---|
| shDock | 31 | 20 | 29 | 4 |
| LZerD | 50 | 9 | 21 | 4 |
| F²Dock (S) | 66 | 4 | 8 | 6 |
| PatchDock | 54 | 9 | 17 | 4 |
| ZDock | 56 | 4 | 21 | 3 |
| SP-Dock vs | Win | Tie | Loss | Both fail |
| shDock | 39 | 18 | 23 | 4 |
| LZerD | 51 | 11 | 18 | 4 |
| F²Dock (S-E) | 59 | 6 | 13 | 6 |
| PatchDock | 58 | 11 | 11 | 4 |
| ZDock | 60 | 6 | 15 | 3 |

The above experiments have been performed using the bound molecules of both the receptor and the li-

gand. In Table IV, the number of cases, where at least one hit is found at different docking thresholds, is presented also for the unbound case. In the R-unbound/L-unbound results, the proposed SP-Dock method is ranked first for all docking thresholds.

In Table IV, the average iRMSD is also presented for both the bound and the unbound cases. The average iRMSD is calculated as follows: for each case that succeeds in finding a hit in the top 3600 predictions, the iRMSD of the best ranked hit is taken. For all these cases, the average iRMSD is computed. It is worth mentioning that the average iRMSD of the proposed method is greater (i.e. less accurate) than the iRMSD of the other methods in the bound case, while it is comparable in the unbound case. This can be explained by the fact that the proposed method provides an approximate estimation of the docking pose, while other methods provide more exact estimations. However, the approximate complementarity matching of SP-Dock allows identification of complementary pairs of patches even after small conformational changes (unbound docking). This is the reason why the proposed method performs better in unbound docking than other methods, while at the same time its iRMSD is not significantly affected (as it happens with the other methods).

**Table IV:** Number of test cases where at least one hit is found for different thresholds (1, 5, 10, 100, 1000, 2000 and 3600) and the average iRMSD, for both R-bound/L-bound and R-unbound/-unbound cases.

| | PatchDock | ZDock | shDock | LZerD | F²Dock (S) | F²Dock (S-E) | S-Dock | SP-Dock |
|---|---|---|---|---|---|---|---|---|
| R-bound/L-bound | | | | | | | | |
| Rank = 1 | 13 | 6 | 23 | 16 | 8 | 8 | 17 | **26** |
| Rank ≤ 5 | 15 | 11 | **37** | 20 | 10 | 13 | 25 | 31 |
| Rank ≤ 10 | 17 | 18 | **41** | 20 | 12 | 17 | 30 | 34 |
| Rank ≤ 100 | 28 | 37 | **57** | 39 | 23 | 33 | 51 | 56 |
| Rank ≤ 1000 | 39 | 48 | 67 | 56 | 46 | 52 | 73 | **75** |
| Rank ≤ 2000 | 41 | 54 | 69 | 58 | 55 | 57 | 75 | **76** |
| Rank ≤ 3600 | 42 | 55 | 74 | 60 | 62 | 59 | 76 | **76** |
| Avg, iRMSD (Å) | 1.53 | 1.73 | 0.69 | 1.17 | 1.01 | 0.95 | 1.73 | 1.68 |
| R-unbound/L-unbound | | | | | | | | |
| Rank = 1 | 0 | **2** | 0 | 1 | 1 | 1 | **2** | **2** |
| Rank ≤ 5 | 1 | 4 | 1 | 2 | 2 | 2 | 3 | **6** |
| Rank ≤ 10 | 1 | 5 | 2 | 2 | 2 | 2 | 9 | **11** |
| Rank ≤ 100 | 9 | 11 | 6 | 14 | 9 | 11 | 23 | **30** |
| Rank ≤ 1000 | 23 | 30 | 22 | 29 | 24 | 27 | **53** | **53** |
| Rank ≤ 2000 | 31 | 35 | 33 | 36 | 31 | 33 | 55 | **56** |
| Rank ≤ 3600 | 37 | 42 | 41 | 38 | 33 | 37 | **56** | **56** |
| Avg, iRMSD (Å) | 1.76 | 1.84 | 1.89 | 1.87 | 1.57 | 1.59 | 1.88 | 1.84 |

Similar conclusions can be drawn in the win-tie-loss-failure records (Table V). The proposed approach clearly outperforms all five methods, even in the case when only shape complementarity is used (S-Dock). If, instead of the geometric-only scoring, the shape-physicochemical scoring of (20) is used, the hit ranks are improved in 60.4% of the cases of Benchmark 2.4. The performance of all five methods for the unbound case

in Benchmark 2.4 is shown in Table VI. It should be stressed that the proposed SP-Dock method does not return a fixed number of docking poses. The number of docking poses corresponds to the number $M$ of patch groups that was presented in Section 3.2 and it varies depending on the size of the interacting proteins. In case $M>3600$, only the first 3600 ranked poses are kept and presented in Table VI. In complexes where $M<3600$, all poses fulfill the constraint of "first 3600 predictions", thus, all hits can be included in Table VI.

**Table V:** R-unbound/L-unbound: the win-tie-loss-failure records for the proposed method versus shDock, LZerD and F²Dock.

| S-Dock vs | Win | Tie | Loss | Both fail |
|---|---|---|---|---|
| shDock | 47 | 0 | 18 | 19 |
| LZerD | 46 | 1 | 19 | 18 |
| F²Dock (S) | 51 | 0 | 13 | 20 |
| PatchDock | 47 | 0 | 15 | 7 |
| ZDock | 41 | 1 | 22 | 5 |
| | | | | |
| SP-Dock vs | Win | Tie | Loss | Both fail |
| shDock | 49 | 0 | 16 | 19 |
| LZerD | 50 | 1 | 15 | 18 |
| F²Dock (S-E) | 47 | 0 | 17 | 20 |
| PatchDock | 52 | 0 | 11 | 6 |
| ZDock | 45 | 1 | 19 | 4 |

**Table VI:** R-unbound/L-unbound: Comparisons between S-Dock, SP-Dock, LZerD, shDock and F²Dock on 84 test cases from Benchmark v2.4. PDB gives the PDB id for the protein complex. RMSD and Rank give the iRMSD and rank of the best ranked hit (2.5 Å cut-off). In 15 cases none of the four methods returned a hit in the first 3600 predictions.

| | PatchDock | | ZDock | | LzerD | | shDock | | F²Dock (S) | | F²Dock (S-E) | | S-Dock | | SP-Dock | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PDB | Rank | RMSD | Rank | RMSD | Rank | RMSD | Rank | RMSD | Rank | RMSD | Rank | RMSD | Rank | RMSD | Rank | RMSD |
| Enzyme–Inhibitor or Enzyme–Substrate | | | | | | | | | | | | | | | | |
| 1ACB | – | – | – | – | – | – | – | – | – | – | – | – | 26 | 2.24 | **10** | **2.17** |
| 1AVX | 2053 | 2.22 | 2863 | 2.23 | 786 | 2.41 | 1199 | 2.5 | 1769 | 1.75 | 1909 | 1.75 | 122 | 2.21 | **97** | **2.21** |
| 1AY7 | 679 | 1.14 | – | – | 1884 | 1.98 | 733 | 1.56 | 94 | 0.87 | **32** | **0.98** | 37 | 2.15 | 83 | 2.15 |
| 1BVN | 110 | 1.68 | 502 | 1.97 | 27 | 2.32 | 82 | 2.15 | 72 | 1.58 | 54 | 1.58 | 6 | 1.65 | **5** | **1.65** |
| 1CGI | – | – | 145 | 2.44 | – | – | – | – | 39 | 2.5 | 45 | 2.5 | **7** | **2.12** | 21 | 2.12 |
| 1D6R | – | – | 2951 | 2.03 | 2619 | 2.24 | – | – | 177 | 1.45 | **170** | **1.45** | 634 | 2.19 | 449 | 2.19 |
| 1DFJ | – | – | **9** | **2.27** | – | – | – | – | 243 | 1.15 | 22 | 1.14 | – | – | – | – |
| 1E6E | 38 | 2.12 | – | – | 52 | 2.13 | 1014 | 1.52 | – | – | 3526 | 2.41 | 728 | 1.72 | **7** | **1.34** |
| 1EAW | 59 | 2.1 | **3** | **1.54** | 20 | 2.42 | 324 | 2.07 | 517 | 1.7 | 454 | 1.52 | 9 | 1.14 | 18 | 1.14 |
| 1EWY | 88 | 2.46 | 259 | 2.32 | 349 | 2.36 | 175 | 2.15 | **4** | **1.21** | **4** | **1.17** | 76 | 2.12 | 15 | 2.12 |
| 1EZU | – | – | 1100 | 1.94 | 824 | 1.21 | **784** | **2.24** | – | – | – | – | – | – | – | – |
| 1F34 | 30 | 1.57 | 5 | 2.2 | – | – | 1528 | 2.14 | 98 | 1.34 | 60 | 1.34 | **1** | **0.72** | **1** | **0.72** |
| 1HIA | – | – | – | – | – | – | – | – | – | – | – | – | **49** | **1.93** | 336 | 1.93 |
| 1KKL | – | – | – | – | – | – | – | – | – | – | – | – | **4** | **2.13** | 8 | 2.13 |
| 1MAH | 1184 | 0.83 | **92** | **1.31** | 92 | 0.87 | 2252 | 2.13 | – | – | 3327 | 2.07 | 1614 | 1.94 | 155 | 1.34 |
| 1PPE | 12 | 1.51 | **1** | **0.57** | **1** | **0.83** | 8 | 1.86 | 355 | 1.12 | 392 | 1.12 | **1** | **0.62** | **1** | **0.62** |
| 1TMQ | **3** | **1.16** | 314 | 1.88 | 50 | 1.45 | 186 | 1.18 | 247 | 1.63 | 241 | 1.63 | 225 | 1.56 | 5 | 1.56 |
| 1UDI | 261 | 1.55 | 258 | 2.17 | 59 | 2.36 | – | – | – | – | 3043 | 1.74 | **25** | **1.56** | 25 | 1.56 |
| 2MTA | 1086 | 0.83 | – | – | 606 | 1.64 | 2423 | 2.11 | 1378 | 1.58 | 1124 | 1.58 | 208 | 1.42 | **24** | **1.19** |
| 2PCC | – | – | – | – | – | – | – | – | – | – | 843 | 0.66 | **14** | **2.21** | 70 | 2.21 |
| 2SIC | 113 | 1.24 | 173 | 1.86 | **12** | **2.04** | 35 | 1.94 | 1072 | 1.79 | 1429 | 2.35 | – | – | – | – |
| 2SNI | – | – | – | – | – | – | – | – | 362 | 1.92 | 377 | 1.92 | 72 | 2.15 | **51** | **1.98** |
| 7CEI | 241 | 2.49 | 106 | 1.97 | – | – | 1515 | 2.05 | 1188 | 1.04 | 598 | 0.85 | 197 | 1.46 | **39** | **1.46** |
| Antibody – Antigen | | | | | | | | | | | | | | | | |
| 1AHW | 168 | 1.3 | 268 | 2.28 | **5** | **1.34** | 1419 | 2.05 | – | – | – | – | 319 | 2.29 | 269 | 2.29 |

| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1BVK | – | – | – | – | – | – | – | – | 801 | 2.21 | **560** | **2.21** | 935 | 2.14 | 576 | 2.14 |
| 1DQJ | – | – | 2287 | 2.48 | – | – | – | – | – | – | – | – | **933** | **1.32** | 1009 | 1.32 |
| 1E6J | 3483 | 2.29 | **15** | **1.56** | 439 | 2.18 | 3065 | 2.49 | – | – | – | – | 631 | 2.21 | 83 | 2.12 |
| 1JPS | 1185 | 1.89 | **171** | **1.81** | 292 | 0.9 | 469 | 1.56 | 484 | 1.24 | 702 | 1.17 | – | – | – | – |
| 1MLC | 847 | 0.98 | 110 | 1.19 | 1834 | 1.16 | 1027 | 0.96 | – | – | – | – | 524 | 2.18 | **54** | **2.18** |
| 1VFB | 1541 | 2.48 | 2734 | 1.79 | 1303 | 1.69 | **207** | **2.23** | 310 | 0.75 | 213 | 0.75 | 611 | 2.09 | 526 | 2.09 |
| 1WEJ | 2152 | 1.25 | 465 | 2.37 | – | – | – | – | – | – | – | – | 386 | 2.26 | **269** | **2.34** |
| 2VIS | – | – | **2747** | **2.49** | – | – | 3027 | 1.45 | – | – | – | – | – | – | – | – |
| Antigen–Bound Antibody | | | | | | | | | | | | | | | | |
| 1BJ1 | – | – | 129 | 0.86 | 298 | 1.86 | 2052 | 1.58 | – | – | – | – | 143 | 2.32 | **14** | **2.17** |
| 1FSK | 420 | 2.08 | **1** | **1.63** | 15 | 2.4 | 47 | 0.62 | – | – | – | – | 46 | 0.91 | 64 | 0.91 |
| 1I9R | – | – | **50** | **2.45** | 95 | 2.39 | 302 | 2.48 | 2739 | 1.51 | 2090 | 1.51 | 71 | 2.28 | 96 | 2.28 |
| 1IQD | 3228 | 2.12 | 612 | 2.27 | **41** | **1.2** | – | – | – | – | – | – | 374 | 2.13 | 75 | 2.13 |
| 1K4C | – | – | – | – | **1188** | **1.43** | – | – | – | – | – | – | – | – | – | – |
| 1KXQ | 11 | 1.5 | 212 | 1.91 | 73 | 1.68 | **30** | **1.41** | 646 | 1.36 | 528 | 1.39 | 308 | 2.2 | 387 | 2.18 |
| 1NSN | 1254 | 1.76 | 185 | 1.96 | 945 | 2.29 | 1364 | 2.03 | – | – | – | – | 179 | 2.01 | **115** | **2.01** |
| 1NCA | 575 | 1.45 | **14** | **1.93** | – | – | 600 | 0.85 | – | – | – | – | 232 | 1.21 | 257 | 1.21 |
| 1QFW | 1457 | 1.85 | 257 | 1.14 | **108** | **1.24** | 759 | 1.08 | 1372 | 1.34 | 1212 | 1.34 | – | – | – | – |
| 2JEL | 1142 | 0.95 | 45 | 1.79 | 133 | 2.49 | – | – | – | – | – | – | 83 | 2.16 | **9** | **1.89** |
| 2HMI | – | – | **237** | **2.5** | – | – | – | – | – | – | – | – | – | – | – | – |
| Others | | | | | | | | | | | | | | | | |
| 1A2K | – | – | – | – | – | – | 237 | 2.45 | – | – | – | – | **13** | **2.21** | 27 | 2.21 |
| 1AKJ | – | – | – | – | – | – | 292 | 2.23 | 102 | 1.45 | **46** | **1.45** | 47 | 2.04 | 484 | 2.04 |
| 1B6C | 201 | 2.14 | 1717 | 2.43 | 1001 | 2.41 | – | – | 1862 | 1.96 | 1687 | 1.96 | 172 | 2.06 | **159** | **2.06** |
| 1BUH | 625 | 2.37 | – | – | – | – | 391 | 1.78 | 65 | 0.75 | **64** | **0.75** | 357 | 1.62 | 735 | 1.62 |
| 1E96 | – | – | 3094 | 2.26 | 216 | 2.14 | 3526 | 2.5 | 300 | 1.79 | **193** | **1.79** | – | – | – | – |
| 1F51 | 650 | 2.03 | 230 | 2.18 | 3545 | 1.58 | 3561 | 2.12 | – | – | – | – | **79** | **2.21** | 154 | 2.21 |
| 1FAK | – | – | – | – | – | – | – | – | – | – | – | – | 993 | 2.11 | **146** | **1.87** |
| 1FQJ | 3004 | 2.46 | – | – | – | – | – | – | 27 | 2.12 | 30 | 2.1 | **7** | **2.17** | 19 | 2.17 |
| 1GCQ | – | – | – | – | – | – | 1787 | 2.21 | – | – | – | – | 681 | 1.93 | **231** | **1.93** |
| 1GP2 | – | – | – | – | – | – | – | – | – | – | – | – | 764 | 2.21 | **288** | **2.21** |
| 1GRN | 831 | 1.54 | 1704 | 2.34 | 1407 | 2.18 | 1724 | 1.61 | 1264 | 2.23 | 674 | 2.23 | **501** | **2.17** | 692 | 2.17 |
| 1HE1 | 33 | 2.16 | – | – | 267 | 1.98 | 3107 | 1.41 | **1** | **1.12** | **1** | **1.12** | – | – | – | – |
| 1HE8 | – | – | – | – | – | – | **646** | **2.27** | – | – | – | – | 1589 | 2.32 | 898 | 2.32 |
| 1I2M | – | – | – | – | – | – | – | – | – | – | – | – | **210** | **1.57** | 482 | 1.57 |
| 1I4D | – | – | – | – | – | – | – | – | – | – | – | – | **647** | **2.14** | 1186 | 2.24 |
| 1IB1 | – | – | – | – | – | – | – | – | – | – | – | – | 16 | 2.27 | **4** | **2.27** |
| 1IJK | – | – | – | – | – | – | 1639 | 2.42 | 2221 | 2.5 | **1426** | **2.43** | – | – | – | – |
| 1KAC | – | – | 2896 | 2.33 | 655 | 2.18 | 138 | 2.38 | 747 | 1.67 | 672 | 1.67 | 8 | 2.12 | **5** | **2.12** |
| 1KLU | – | – | – | – | – | – | – | – | – | – | – | – | 2450 | 1.89 | **1861** | **1.67** |
| 1KTZ | – | – | – | – | – | – | – | – | – | – | – | – | 286 | 1.23 | **17** | **1.23** |
| 1KXP | 37 | 2.49 | 1734 | 2.36 | – | – | **3** | **1.7** | 306 | 2.01 | 157 | 2.01 | 100 | 0.87 | 8 | 0.87 |
| 1ML0 | 450 | 1.59 | **36** | **1.56** | 559 | 2.38 | 303 | 1.87 | – | – | – | – | 714 | 1.92 | 522 | 1.92 |
| 1QA9 | 3039 | 1.86 | – | – | 1381 | 2.19 | **1264** | **2.16** | – | – | – | – | – | – | – | – |
| 1WQ1 | – | – | 1101 | 2.49 | 141 | 1.87 | – | – | 96 | 1.95 | 62 | 1.95 | **7** | **1.76** | 102 | 1.76 |
| 2BTF | – | – | – | – | – | – | – | – | – | – | – | – | **236** | **1.72** | 363 | 1.56 |
| 2QFW | 1018 | 1.68 | 832 | 2.29 | **68** | **1.55** | – | – | 525 | 1.18 | 427 | 1.18 | – | – | – | – |

## 5.2 Computational Issues

In Table VII, the average computation times for various tasks of the proposed approach are presented. The average time required for extraction of the Local Spectral Descriptor for a GSP is 0.1s. The time required for matching between a pair of GSPs (using the Local Spectral Descriptor) is ~ 0.001ms. It is obvious that descriptor matching is $8 \cdot 10^4$ times faster than the geometric scoring based on distance grid, which demonstrates the importance of the Local Spectral Descriptor as a fast filtering stage.

**Table VII:** Average computation times for various tasks of the proposed approach

| Activity | Average Computation Time |
|---|---|
| Local Spectral Descriptor Extraction / GSP | 100ms |
| Complementarity matching of a pair of GSPs | 0.001ms |
| Scoring (distance grid) of a pose | 86ms |

The average running time of SP-Dock is the sum of a) the time required for preprocessing and descriptor extraction; b) descriptor matching, grouping and alignment; c) distance-grid-based geometric scoring and d) physicochemical scoring. The most time-consuming parts are the geometric and physicochemical scoring, while the fastest part is the descriptor matching, grouping and alignment. The average running time for small-to-medium-sized complexes is approximately one hour (Table VIII), while it takes a few hours for large complexes. The running times for all complexes are given in Table X, in Appendix 2 of the supplementary material. The times reported in this paper were obtained using a PC with a dual-core 2.4 GHz processor and 8GB RAM.

Although we did not run the other three methods, we compare our algorithm with the times reported in the related articles. As stated in [32], LZerD requires 1-2 hours for small proteins and it may take longer for larger proteins. These numbers were obtained using a computer with dual-core 2.1 GHz processor with 8 GB RAM, i.e. similar to the PC that we conducted our experiments. Thus, our approach is slightly faster than LZerD, while at the same time it clearly outperforms LZerD. In [33], authors use a computer with i7 quad-core processor at 3.2GHz and 12GB RAM. The average running time for shDock is reported to be 2758s, i.e. a bit less than SP-Dock and LZerD. Taking into account the fact that in shDock a higher performance computer is used, it can be inferred that the average running time is comparable to SP-Dock and LZerD. Finally, in [34], no specific running time is reported for F$^2$Dock. It should be stressed that the running time for SP-Dock includes also the time for physicochemical scoring, while in the cases of LZerD and shDock only geometric docking is considered. If we keep only the geometric part, our method becomes much faster than LZerD and shDock. On the other hand, if we use both shape and physicochemical properties, we produce much better docking results within approximately the same running time.

**Table VIII:** Average running time of the proposed method

| Average Running Time | | | | |
|---|---|---|---|---|
| Preprocessing/ Descriptor extraction | Descriptor Matching, Grouping, Alignment | Geometric Scoring | Physicochemical Scoring | Average Running Time |
| 300s | 85s | 1935s | 1340s | 3660s |

## 6  CONCLUSIONS

We have presented a unified framework for protein-protein docking based on both shape and physicochemi-

cal complementarity. For shape complementarity, a new approach has been implemented, which utilises an effective local descriptor. The so-called Local Spectral Descriptror is compact, fast to extract and capable to capture similatities of local surface patches. As a next step, multiple pairs of complementary local patches from the receptor and the ligand are grouped together using a new grouping algorithm. The above grouping algorithm was inspired by the observation that shape complementarity in protein surfaces is mostly approximate rather than exact, thus single-patch or two-patch complementary matching generates numerous false-positive predictions. Additionally, shape complementarity is enhanced by physicochemical complementarity. Several non-geometric factors were tested and their contribution to the improvement of the shape-only docking predictions was assessed. Particle Swarm Optimization was applied to train the weights that each factor contributes to the overall scoring function. The most significant improvement is achieved when Atom Desolvation Energy, Electrostatic Complementarity, Hydrophobicity, Coulomb Potential and van der Waals Potential are introduced along with the shape complementarity, while Residue Contact Preferences and Generic Contact Preference seem to have insignificant contribution. This was an initial selection of the most well-known non-geometric factors. More factors that are available in the literature can be tested and assessed in a similar manner, which is planned for future work.

The proposed method advances the state of the art mainly in the parts of local surface complementarity matching and alignment. The Local Spectral Descriptor provides a more robust measure for shape complementarity of local patches, while the new grouping algorithm enhances the certainty of a wider surface region of receptor to be complementary to a wider surface region of the ligand. Additionally, instead of superimposing the sparse points of the ligand on the matching points of the receptor, as it is the case with most of the existing local-patch-based docking approaches, the ICP algorithm used by SP-Dock achieves alignment of the two proteins by taking into account the overall shape of the complementary regions. While this feature provides less accurate poses (i.e. with higher iRMSD) in the bound case, it significantly improves the unbound case. The reason is that the surfaces of the two proteins at their binding interfaces have approximate complementarity in the unbound case. Thus, a method based on exact matching and alignment would probably fail to retrieve a near-native pose within the list of predicted poses, while a more approximate method, such as SP-Dock, is more likely to achieve a correct prediction. This is an interesting conclusion and could assist in further research in protein-protein docking by proposing ideas on how to deal with unbound docking and slight side-chain flexibility. Another advancement of SP-Dock is the new scoring process based on geometric and physicochemical factors. Several works have been presented so far dealing with the assess-

ment of physicochemical factors but only few of them address both geometric and physicochemical complementarity. The proposed scoring of SP-Dock can be used as a starting point for further research, where the effect of additional factors, apart from Atom Desolvation Energy, Electrostatic Complementarity, Hydrophobicity, Coulomb Potential and van der Waals Potential, could be investigated.

Results performed on the 84 complexes of the Docking Benchmark 2.4 demonstrate the superiority of the proposed SP-Dock method over five similar docking approaches. While in the case of bound complexes our method performs slightly better than the best docking methods reported so far, in the unbound case our approach clearly outperforms them. This confirms the assumption that shape complementarity should be approximate (not exact) in order to take into account small side-chain conformations on the protein surface. Additionally, when several physicochemical factors are introduced (SP-Dock), the shape-only docking predictions are improved in both bound and unbound cases. Despite the improvements of the proposed SP-Dock method presented above and the interesting conclusions regarding the protein-protein docking problem, there is still a lot of work to be done in this direction. In terms of accuracy, research should focus on the following two goals: i) to appropriately model the contribution of each factor (geometric or non-geometric) to protein interactions; ii) to appropriately model the flexibility (both side-chain and backbone) of the interacting proteins. Existing methods have reached an acceptable level in terms of computation time, though not adequately modeling the flexibility. If a deeper analysis of the flexibility takes place, then the computational time increases prohibitively. This tradeoff between accuracy and computation time should be considered, until a method that will address both problems is proposed.

## REFERENCES

[1] D. W. Richie, "Recent Progress and Future Directions in Protein-Protein Docking", *Current Protein and Peptide Science*, 2008, 9, 1-15.

[2] J.C. Camacho, D.W. Gatchell, S.R. Kimura, and S. Vajda. "Scoring docked conformations generated by rigid body protein–protein docking". *PROTEINS: Structure, Function and Genetics*, 40:525–537, 2000.

[3] R. Chen and Z Weng. "Docking unbound proteins using shape complementarity, desolvation, and electrostatics". *PROTEINS: Structure, Function and Genetics*, 47:281–294, 2002.

[4] Carter, P., Lesk, V.A., Islam, S.A. and Sternberg, M.J.E. (2005) *Proteins: Struct. Func. Bioinf.*, 60, 281-288.

[5] Kozakov, D., Brenke, R., Comeau, S.R. and Vajda, S. (2006) *Proteins: Struct. Func. Bioinf.*, 65, 392-406.

[6] Eisenstein, M. and Katchalski-Katzir, E. (2004) *Comptes Rendus Biologies*, 327, 409-420.

[7] Ritchie, D.W. and Kemp, G.J.L. "Protein Docking Using Spherical Polar Fourier Correlations", (2000) *Proteins: Struct. Func. Genet.*, 39(2) 178-194.

[8] D. Duhovny, R. Nussinov, and H. J. Wolfson. "Efficient unbound docking of rigid molecules". *In 2'nd Workshop on Algorithms in Bioinformatics*, pages 185–200, 2002.

[9] Zujun Shentu, Mohammad Al Hasan, Chris Bystroff and Mohammad J. Zaki, "Context Shapes: Efficient Complementary Shape Matching for Protein-Protein Docking". *Proteins: Structure, Function and Bioinformatics*, 70(3):1056-1073. February 2008.

[10] R. Chen and Z. Weng. "ZDOCK: An initial-stage protein-docking algorithm". *Proteins: Structure, Function and Genetics*, 52:80–87, 2003.

[11] H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N. Shindyalov, and P.E. Bourne, "The Protein Data Bank," *Nucleic Acids Research*, vol. 28, pp. 235-242, 2000.

[12] M.L. Connolly. "Solvent-accessible surfaces of proteins and nucleic acids". *Science*, 221:709–713, 1983.

[13] M.F. Sanner, A.J. Olson, and J.-C. Spehner. "Fast and robust computation of molecular surfaces". *In 11th ACM Symposium on Computational Geometry*, 1995.

[14] D. Fischer, S.L. Lin, H.L. Wolfson, and R. Nussinov. "A geometry-based suite of molecular docking processes". *J. Mol. Bio.*, 248:459–477, 1995.

[15] A.Mademlis, P.Daras, D.Tzovaras and M.G.Strintzis, "3D Object Retrieval based on Resulting Fields" *29th International conference on EUROGRAPHICS 2008, workshop on 3D object retrieval*, Crete, Greece, Apr 2008

[16] LING H., OKADA K.: "Diffusion distance for histogram comparison". *In CVPR '06: Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (Washington, DC, USA, 2006), IEEE Computer Society, pp. 246–253.

[17] J. Mintseris, K. Wiehe, B. Pierce, R. Anderson, R. Chen, J. Janin, and Z. Weng. "Protein-protein docking benchmark 2.0: An update". *Proteins: Structure, Function and Genetics*, 60(2):214–216, 2005.

[18] R. Chen, J. Mintseris, J. Janin, and Z. Weng. "A protein-protein docking benchmark". *Proteins: Structure, Function and Genetics*, 52(1):88–91, 2003.

[19] Brian Pierce and ZhipingWeng. "ZRANK: Reranking Protein Docking Predictions With an Optimized Energy Function", *PROTEINS: Structure, Function, and Bioinformatics*, 67:1078–1086 (2007)

[20] Tim Geppert, Ewgenij Proschak, Gisbert Schneider, "Protein-protein docking by shape-complementarity and property matching", *Journal of Computational Chemistry*, Volume 31 Issue 9, Pages 1919 – 1928, 2010.

[21] H. Edelsbrunner, E. P. Miicke, "Three-dimensional alpha shapes", *Proceedings of the 1992 workshop on Volume visualization (VVS '92)*, pp 75-82, NY, USA, 1992.

[22] Z. Lian, A. Godil, B. Bustos, M. Daoudi, J. Hermans, S. Kawamura, Y. Kurita, G. Lavoué, H.V. Nguyen, R. Ohbuchi, Y. Ohkita, Y. Ohishi, F. Porikli, M. Reuter, I. Sipiran, D. Smeets, P. Suetens, H. Tabia, D. Vandermeulen, "SHREC'11 Track: Shape Retrieval on Non-rigid 3D Watertight Meshes", *Eurographics Workshop on 3D Object Retrieval (2011)*, April 10, 2011, Llandudno (UK).

[23] G. Lavoué, "Bag of Words and Local Spectral Descriptor for 3D Partial Shape Retrieval", *Eurographics Workshop on 3D Object Retrieval (2011)*, April 10, 2011, Llandudno (UK).

[24] B. Vallet, B. Levy, "Spectral geometry processing with manifold harmonics", *Computer Graphics Forum* 27, 2 (2008), 251–260.

[25] M. Reuter, F.-E. Wolter and N. Peinecke, "Laplace-Beltrami spectra as "Shape-DNA" of surfaces and solids", *Computer-Aided Design* 38 (4), pp.342-366, 2006.

[26] C. Zhang, G. Vasmatzis, J. L. Cornette, C. DeLisi, "Determination of atomic desolvation energies from the structures of crystallized proteins", *J. Mol Biol*, 1997, 267, 707.

[27] F. Glaser, D.M. Steinberg, I.A. Vakser, N. Ben-Tal, "Residue frequencies and pairing preferences at protein-protein interfaces", *Proteins*, 2001 May 1;43(2):89-102.

[28] M. Berrera, H. Molinari, F. Fogolari, "Amino acid empirical contact energy definitions for fold recognition in the space of contact maps", *BMC Bioinformatics 2003*, 4:8.

[29] G. Ausiello, G. Cesareni, M. Helmer-Citterich, "ESCHER: a new docking procedure applied to the reconstruction of protein tertiary structure", *Proteins* 1997, 28(4), 556-567.

[30] N. P. Palma, L. Krippahl, J. E. Wampler, J. J. G. Moura, "BiGGER: A new (soft) docking algorithm for predicting protein interactions", *Proteins: Structure, Function, and Bioinformatics*, 2000, 39, 372.

[31] R. Norel, D. Petrey, H.J. Wolfson, R. Nussinov, "Examination of shape complementarity in docking of unbound proteins", *Proteins* 1999, 36, 307.

[32] V. Venkatraman, Y. Yang, L. Sael, D. Kihara, "Protein-protein docking using region-based 3D Zernike descriptors", *BMC Bioinformatics*, 2009;10:407.

[33] S. Gu, P. Koehl, J. Hass, N. Amenta, "Surface-histogram: A new shape descriptor for protein-protein docking", *Proteins 2012*, 80:221–238.

[34] C. Bajaj, R. Chowdhury, V. Siddavanahalli, "F$^2$Dock: Fast Fourier Protein-Protein Docking", *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, Vol. 8, No. 1, Jan/Feb 2011.

[35] A. Axenopoulos, P. Daras, G. Papadopoulos, E. Houstis, "A Shape Descriptor for Fast Complementarity Matching in Molecular Docking", *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, Vol. 8, Issue: 6, Pages: 1441-1457, 2011.

[36] A. Axenopoulos, P. Daras, G. Papadopoulos, E. Houstis, "3D Protein-Protein Docking using Shape Complementarity and Fast Alignment", *IEEE int Conference on Image Processing*, ICIP 2011, Sep 11-14, Brussels.

[37] P.J. Besl and N.D. McKay, Reconstruction of Real-World Objects via Simultaneous Registration and Robust Combination of Multiple Range Images, *IEEE Trans. Pattern Analysis and Machine Intelligence*, 14:2(239-256), 1992.

[38] Particle Swarm Optimization, Available Online: http://www.swarmintelligence.org/tutorials.php.

[39] A.Mademlis, P.Daras, D.Tzovaras and M.G.Strintzis, "3D Object Retrieval using the 3D Shape Impact Descriptor" *ELSEVIER, Pattern Recognition*, Volume 42 , Issue 11, pp. 2447-2459, Nov 2009.

[40] J. Janin, K. Henrick, J. Moult, LT. Eyck, MJ. Sternberg, S. Vajda, I. Vakser, SJ. Wodak, "CAPRI: a Critical Assessment of PRedicted Interactions", *Proteins* 2003;52:2–9.

[41] E. Mashiach, D. Schneidman-Duhovny, A. Peri, Y. Shavit, R. Nussinov and H. J. Wolfson, "An Integrated Suite of Fast Docking Algorithms", *Proteins*. 2010 November 15; 78(15): 3197–3204.

[42] H. Hwang, T. Vreven, B. G. Pierce, J. Hung, and Z. Weng, "Performance of ZDOCK and ZRANK in CAPRI rounds 13–19", *Proteins, Structure Function Bioinformatics*, Vol. 78, Issue 15, pages 3104–3110, 15 November 2010.

**Apostolos Axenopoulos** was born in Thessaloniki, Greece, in 1980. He received the Diploma degree in electrical and computer engineering and the M.S. degree in advanced computing systems from Aristotle University of Thessaloniki, Thessaloniki, Greece, in 2003 and 2006, respectively. He is an Associate Researcher at the Information Technologies Institute, Thessaloniki. His main research interestsinclude 3D object indexing, content-based search and retrieval and bioinformatics. Currently, he is a PhD candidate in Computer & Communication Engineering Department, University of Thessaly.

**Petros Daras** (M'07) was born in Athens, Greece, in 1974. He received the Diploma degree in electrical and computer engineering, the M.Sc. degree in medical informatics, and the Ph.D. degree in electrical and computer engineering, all from the Aristotle University of Thessaloniki, Thessaloniki, Greece, in 1999, 2002, and 2005, respectively. He is a Researcher Grade C, at the Information Technologies Institute (ITI) of the Centre for Research and Technology Hellas (CERTH). His main research interests include search, retrieval and recognition of 3D objects, 3D object processing, medical informatics applications, medical image processing, 3D object watermarking and bioinformatics. He regularly serves as a reviewer/evaluator of European projects and he is a member of IEEE, a key member of the IEEE MMTC 3DRPC IG and chair of the IEEE Image, Video and Mesh Coding IG.

**Georgios Papadopoulos** received his diploma in Physics from the Aristotle University of Thessaloniki/Greece in 1977. He studied further theoretical Biophysics in the department of Physics of the Freie Universitat Berlin/Germany and received his Ph.D degree in Biophysics from

the same department in 1989. He worked with short term contracts as guest scientist in the Hahn-Meitner Institut/Berlin/Germany and in the Fosrchungszentrum Julich/Germany. Since 1994 he has been teaching Physics, Biostatistics, Bioinformatics, Physical Chemistry and Biophysics in the University of Thessaly/Greece, Democritus University of Thrace/Greece and the Aristotle University of Thessaloniki/Greece. Since 2009 he is lecturer of Biophysics in the department of Biochemistry & Biotechnology/UTh. His Research interests are focused on the study of the structure of biological macromolecules and of their interactions using theoretical and computational methods.

**Elias. N. Houstis** is currently a full Professor of Computer Engineering and Communications department at University of Thessaly, Greece, Director of Research Center of Thessaly (CE.RE.TE.TH.), and Emeritus Professor of Purdue University. USA. Most of his academic career is associated with Purdue University. He has been a Professor of Computer Science and Director of the Computational Science & Engineering Program of Purdue University. He is a member of working groups WG2.5 IFIP on mathematical software and European ICT Directors. Houstis' current research interests are in the areas of problem solving environments, networking and parallel computing, enterprise systems, computational intelligence and finance, and e-services.