

# Element Interaction Manifolds for Systematic Analysis of Geometric Representations in Protein–Ligand Recognition

Alireza Shahi<sup>1</sup>, Md Masud Rana<sup>2</sup>, Duc Duy Nguyen<sup>1\*</sup>

<sup>1</sup>Department of Mathematics, University of Tennessee, Knoxville, TN 37996, USA

<sup>2</sup>Department of Mathematics, Kennesaw State University, Kennesaw, GA 30144, USA

## Abstract

Accurate modeling of protein–ligand interactions requires representations that capture both molecular geometry and chemically specific interactions. Existing approaches encode molecular recognition through several geometric formalisms, including differential geometry and spherical harmonic surface descriptors such as three-dimensional Zernike descriptors (3DZD). However, the relative roles of local curvature, surface area, volume, global shape, and chemical partitioning in protein–ligand recognition remain insufficiently understood. In this work, we present a unified geometric study centered on the Element Interaction Manifold (EIM) framework. We integrate element-interactive curvature, surface-area, and volume descriptors into a common representation and systematically evaluate their roles in binding affinity prediction and ligand-aware binding-site similarity. To place differential geometry in context, we compare EIM against global 3DZD and a newly developed chemically partitioned extension termed element-pair 3DZD (EP-3DZD). Across PDBbind/CASF-2016 affinity prediction and unsupervised ligand-identity recognition tasks, chemically resolved local differential geometry consistently provides the strongest predictive performance and the best accuracy–complexity trade-off. Combining EIM with spherical harmonic descriptors yields only marginal improvements, suggesting limited complementarity between local interaction geometry and global shape representations. Overall, these results suggest that chemically resolved local differential geometry provides a compact, interpretable, and transferable foundation for protein–ligand recognition.

**Keywords:** Geometric Data Analysis, Protein–Ligand Interaction, Differential Geometry, Molecular Surface, Binding Affinity Prediction, Binding-Site Similarity, Geometric Machine Learning, Drug Discovery

---

\*Address correspondence to Duc Duy Nguyen. E-mail: [ducnguyen@utk.edu](mailto:ducnguyen@utk.edu)

# 1 Introduction

Accurately modeling protein–ligand interactions is a central goal in structure-based drug discovery and computational structural biology. Reliable molecular representations are needed for binding affinity prediction, virtual screening, off-target prediction, polypharmacology analysis, and rational lead optimization [1–8]. Classical empirical and machine-learning scoring functions, including X-Score [9], AutoDock Vina [10], RF-Score [11], and ID-Score [12], describe protein–ligand binding using combinations of atom-type contacts, van der Waals interactions, hydrogen bonding, hydrophobic effects, desolvation, and shape complementarity [1, 4, 13]. These methods have demonstrated that chemically meaningful interaction categories are informative for affinity prediction, but they often rely on predefined descriptor classes and do not directly resolve the continuous geometry of molecular recognition.

A fundamental challenge is how to represent molecular geometry in a way that is both mathematically expressive and chemically interpretable. Protein–ligand recognition is influenced by shape complementarity, surface packing, electrostatics, hydrogen bonding, hydrophobic contacts, and local chemical environments [2, 14–19]. Foundational work emphasized solvent-accessible and solvent-excluded surfaces, surface area, packing, reduced surfaces, and space-filling representations [20–22]. Computational surface generation has remained an active area of research, with methods based on Marching Tetrahedra, Alpha Shapes, Euclidean distance transforms, and level-set formulations [16, 19, 23, 24]. These developments reflect a long-standing view that molecular recognition is not only a chemical problem but also a geometric one.

Geometry enters protein–ligand modeling through several complementary formalisms. One major class is differential geometry, where molecular surfaces are described by local curvature quantities such as principal curvatures, mean curvature, and Gaussian curvature [25–27]. Differential-geometric molecular surface models have been used in solvation modeling, electrostatics, molecular surface construction, and binding affinity prediction [2, 25, 28–30]. In our previous DG-GL framework, element-interactive curvatures were introduced to encode chemically resolved protein–ligand interactions through differentiable density manifolds [31]. Separately, the EISA-score framework developed element-interactive surface-area and volume descriptors for binding affinity prediction [32]. These studies demonstrated the predictive potential of individual geometric descriptor families, but they did not systematically examine how curvature, surface area, volume, global shape, and chemical partitioning complement or overlap with one another in protein–ligand recognition.

A second class of methods represents molecular geometry through global shape expansions. Three-dimensional Zernike descriptors (3DZD) project molecular surfaces or volumetric grids onto an orthonormal basis on the unit ball, producing compact rotation-invariant shape descriptors [18, 24, 33, 34]. These descriptors have been widely used in pocket matching, protein shape comparison, and virtual screening. More recent surface-comparison approaches, such as PL-PatchSurfer3, incorporate local patches and visibility information to improve binding-site comparison [35]. Community benchmarks such as SHREC 2025 further highlight the growing interest in integrating surface shape, electrostatics, and local pocket properties for protein binding-site analysis [36]. However, global spherical harmonic descriptors may compress away fine-scale interaction geometry, especially when chemically distinct protein–ligand contacts produce similar overall shapes.

A third family of approaches encodes protein–ligand complexes through contact patterns, interaction fingerprints, or graph-based representations. RF-Score and related contact-count descriptors represent complexes by atom-type pair counts within distance bins [1]. Pharmacophore and interaction-fingerprint methods, including IChem GRIM, encode hydrogen bonds, hydrophobic contacts, ionic interactions, and other interaction motifs as graph-based or pseudoatom-based patterns [37]. These methods are directly ligand-aware and interpretable, but their discrete nature may limit their ability to represent continuous surface geometry, curvature, and shape complementarity. They therefore provide an important contrast to differential-geometric descriptors: contact and graph methods encode explicit interaction events, whereas surface-based descriptors encode the geometric environment in which those events occur.

Modern learning architectures provide an important broader context for geometry in protein–ligand modeling. Three-dimensional convolutional neural networks, graph neural networks, equivariant architectures, protein language models, and diffusion-based docking methods have demonstrated the value of geometric inductive bias in affinity prediction, pose prediction, and molecular design [38–44]. However, these methods differ substantially in training data, architecture, supervision, and target task, making them difficult to compare directly in a controlled descriptor-level study. The present work therefore focuses on explicit and physically interpretable representations, including differential geometry, spherical harmonics, contact histograms, ligand-localized pocket geometry, and interaction graphs. This allows us to ask which geometric signals themselves are most informative for protein–ligand recognition.

Taken together, existing methods represent protein–ligand geometry at different levels: local curvature, surface area and volume, global spherical harmonic shape, element-specific shape, atom-type contact histograms, interaction graphs, learned neural representations, and generative docking models. However, the relative roles of these geometric signals remain insufficiently understood. In particular, it is unclear whether binding affinity and ligand-aware binding-site similarity are driven primarily by global shape, local surface morphology, chemically partitioned shape, explicit contact patterns, or chemically resolved differential geometry. This motivates a unified descriptor-level study that disentangles the contributions of curvature, surface area, volume, global shape, chemical partitioning, ligand-localized pocket morphology, and interaction graphs under a common evaluation framework.

In this work, we present a unified Element Interaction Manifold (EIM) framework for protein–ligand recognition. Rather than introducing curvature, surface area, or volume descriptors in isolation, the present study consolidates previously developed element-interactive geometric components into a common representation and systematically analyzes their roles. EIM constructs element-pair-specific density manifolds from protein–ligand atomic coordinates and extracts surface area, volume, mean curvature, Gaussian curvature, and principal curvature descriptors at both global and local spatial scales. This formulation integrates three principles: element-level chemical specificity, continuous differential geometry, and spatial localization around the protein–ligand interface.

The novelty of this work lies in the unified comparison and role analysis of explicit geometric representations for protein–ligand modeling. First, we integrate curvature, surface-area, and volume descriptors into a common EIM representation and evaluate their combined utility for binding affinity prediction. Second, we introduce element-pair 3D Zernike de-

scriptors (EP-3DZD), a chemically partitioned extension of classical 3DZD, to test whether global spherical harmonic descriptors recover interaction-level information when element specificity is added. Third, we construct a ligand-aware ESES pocket geometry baseline based on localized solvent-excluded surface statistics, enabling a controlled comparison between pocket morphology and interaction geometry. Fourth, we evaluate all descriptors on two complementary tasks: supervised binding affinity prediction on PDBbind/CASF-2016 and unsupervised ligand-identity recognition, where descriptors are tested by their ability to distinguish complexes sharing the same ligand from complexes binding different ligands.

Through this design, the paper surveys the role of geometry in protein–ligand interactions rather than simply reporting another scoring function. Our results suggest that chemically resolved local differential geometry provides highly informative signals for both binding affinity prediction and ligand-aware binding-site similarity. In contrast, global spherical harmonic descriptors, including chemically partitioned variants, provide limited additional complementarity when combined with interaction-aware geometric features. These findings support the view that local, element-specific surface geometry is a useful and interpretable source of transferable information in protein–ligand recognition.

The rest of the paper is organized as follows. Section 2 introduces the mathematical background for element interaction manifolds, differential geometry, and 3D Zernike descriptors. Section 3 describes the EIM feature extraction pipeline, EP-3DZD, ligand-aware ESES pocket descriptors, interaction-based baselines, and evaluation protocols. Section 4 presents results for binding affinity prediction and ligand-aware binding-site similarity. Section 5 concludes the paper.

## 2 Mathematical Background

### 2.1 Element Interaction Manifolds via Discrete-to-Continuum Mapping

Following the approach of Rana and Nguyen [32], we construct element-pair-specific interaction surfaces through a discrete-to-continuum mapping. This framework transforms discrete atomic coordinates into continuous density fields, allowing for the application of differential geometry to fragmented molecular data.

Given a protein–ligand complex with  $N$  atoms, let

$$\mathcal{X} = \{(\mathbf{r}_i, \alpha_i) \mid \mathbf{r}_i \in \mathbb{R}^3, \alpha_i \in \mathcal{T}, i = 1, 2, \dots, N\} \quad (1)$$

denote the collection of atoms annotated by their coordinates  $\mathbf{r}_i$  and the element types  $\alpha_i$ . Here

$$\mathcal{T} = \{\text{H, C, N, O, S, P, F, Cl, Br, I}\} \quad (2)$$

represents the set of element types considered in this work.

For two type of elements  $\mathcal{T}_k$  and  $\mathcal{T}_{k'}$ , the element interactive collection is defined as

$$\mathcal{X}_{kk'} = \{\mathbf{r}_j \mid \alpha_j \in \mathcal{T}_{k'}, \exists i \text{ with } \alpha_i \in \mathcal{T}_k \text{ such that } \|\mathbf{r}_i - \mathbf{r}_j\| \leq d_c\} \quad (3)$$

where  $d_c$  denotes a cutoff distance that defines the interaction region. The corresponding element interactive domain  $D_{kk'}$  is defined as the union of balls centered at atoms in  $\mathcal{X}_{kk'}$ ,

$$D_{kk'} = \bigcup_{\mathbf{r}_j \in \mathcal{X}_{kk'}} B(\mathbf{r}_j, d_c), \quad (4)$$

where  $B(\mathbf{r}_j, d_c)$  denotes a ball with center  $\mathbf{r}_j$  and radius  $d_c$ . The closure of this domain is denoted by  $\bar{D}_{kk'}$ .

To construct a smooth geometric representation, we generate a molecular density for each element-pair interaction through a discrete-to-continuum mapping using a  $C^2$  correlation kernel  $\Phi$ . The global density associated with the element types  $\mathcal{T}_k$  and  $\mathcal{T}_{k'}$  is defined as

$$\rho_{kk'}(\mathbf{r}, \Phi) = \sum_{\mathbf{r}_j \in \mathcal{X}_{kk'}} \Phi(\|\mathbf{r} - \mathbf{r}_j\|; \eta_{kk'}), \quad \mathbf{r} \in \bar{D}_{kk'}, \quad (5)$$

where the kernel  $\Phi$  satisfies the admissibility conditions

$$\Phi(d) \rightarrow 1 \quad \text{as } d \rightarrow 0, \quad \Phi(d) \rightarrow 0 \quad \text{as } d \rightarrow \infty. \quad (6)$$

In this work, we employ the generalized exponential kernel

$$\Phi(\|\mathbf{r} - \mathbf{r}_j\|; \eta_{kk'}) = \exp\left(-\frac{\|\mathbf{r} - \mathbf{r}_j\|^\kappa}{\eta_{kk'}^\kappa}\right), \quad \kappa > 0. \quad (7)$$

The kernel scale parameter is defined as

$$\eta_{kk'} = \tau(\bar{r}_k + \bar{r}_{k'}), \quad (8)$$

where  $\bar{r}_k$  and  $\bar{r}_{k'}$  denote the van der Waals radii of the corresponding element types. Parameters  $\tau$  and  $\kappa$  are determined through cross-validation.

To bound the density field, we define the normalized global density

$$\hat{\rho}_{kk'}(\mathbf{r}, \Phi) = \frac{\rho_{kk'}(\mathbf{r}, \Phi)}{\max_{\mathbf{r} \in \bar{D}_{kk'}} \rho_{kk'}(\mathbf{r}, \Phi)}. \quad (9)$$

In addition to the global representation, we further examine a local density centered at a specific atom  $\mathbf{r}_i$  with element type  $\alpha_i = \mathcal{T}_k$ . This local density describes the interaction between the atom  $i$  and all atoms of type  $\mathcal{T}_{k'}$ :

$$\rho_{kk'}^i(\mathbf{r}, \Phi) = \Phi(\|\mathbf{r} - \mathbf{r}_i\|; \eta_{kk'}) + \sum_{j \neq i, \alpha_j = \mathcal{T}_{k'}} \Phi(\|\mathbf{r} - \mathbf{r}_j\|; \eta_{kk'}), \quad \mathbf{r} \in \bar{D}_{kk'}^i, \quad (10)$$

here  $\bar{D}_{kk'}^i$  denotes a local element interactive domain enclosing

$$D_{kk'}^i = B(\mathbf{r}_i, d_c) \cup \bigcup_{j \neq i, \alpha_j = \mathcal{T}_{k'}} B(\mathbf{r}_j, d_c). \quad (11)$$

It follows that the global interaction domain can be expressed as the union of these local components,

$$D_{kk'} = \bigcup_{i: \alpha_i = \mathcal{T}_k} D_{kk'}^i. \quad (12)$$

The assembly of local densities  $\rho_{kk'}^i$  enables the model to capture fine-grained atomic interactions, allowing the geometric representation to encode essential physical and chemical characteristics at the atom-pair level.

## 2.2 Differential Geometry of Element Interactive Surfaces

For each element-pair interaction density  $\rho_{kk'}(\mathbf{r})$ , the resulting isosurface is treated as a smooth manifold embedded in  $\mathbb{R}^3$ . The intrinsic and extrinsic geometric properties are extracted by analytically evaluating curvature quantities derived from the density field.

Let the isosurface be locally parameterized by  $X(u_1, u_2)$ . The second fundamental form is determined from the shape operator

$$H(X_i, X_j) = (h_{ij})_{i,j=1,2} = \left( \left\langle -\frac{\partial \mathbf{N}}{\partial u_i}, X_j \right\rangle \right), \quad (13)$$

where  $X_i = \frac{\partial X}{\partial u_i}$ ,  $i = 1, 2$  and the unit normal vector is

$$\mathbf{N}(u_1, u_2) = \frac{X_1 \times X_2}{\|X_1 \times X_2\|}. \quad (14)$$

The mean curvature is given by

$$H = \frac{1}{2} g^{ij} h_{ij}, \quad (15)$$

where  $(g^{ij}) = (g_{ij})^{-1}$  and the Einstein summation is assumed. The Gaussian curvature is defined as

$$K = \frac{\det(h_{ij})}{\det(g_{ij})}. \quad (16)$$

### 2.2.1 Element Interactive Curvatures

Based on the discrete-to-continuum mapping, the Gaussian curvature ( $K$ ) and mean curvature ( $H$ ) of the element-interactive density  $\rho_{kk'}(\mathbf{r})$  can be analytically evaluated [31]. Let

$$\rho_x = \frac{\partial \rho}{\partial x}, \quad \rho_{xy} = \frac{\partial^2 \rho}{\partial x \partial y},$$

and similarly for the other derivatives. Define

$$g = \rho_x^2 + \rho_y^2 + \rho_z^2. \quad (17)$$

The Gaussian curvature is

$$\begin{aligned} K = \frac{1}{g^2} [ & 2\rho_x \rho_y \rho_{xz} \rho_{yz} + 2\rho_x \rho_z \rho_{xy} \rho_{yz} + 2\rho_y \rho_z \rho_{xy} \rho_{xz} \\ & - 2\rho_x \rho_z \rho_{xz} \rho_{yy} - 2\rho_y \rho_z \rho_{xx} \rho_{yz} - 2\rho_x \rho_y \rho_{xy} \rho_{zz} \\ & + \rho_z^2 \rho_{xx} \rho_{yy} + \rho_x^2 \rho_{yy} \rho_{zz} + \rho_y^2 \rho_{xx} \rho_{zz} \\ & - \rho_x^2 \rho_{yz}^2 - \rho_y^2 \rho_{xz}^2 - \rho_z^2 \rho_{xy}^2 ], \end{aligned} \quad (18)$$

and the mean curvature is

$$H = \frac{1}{2g^{3/2}} [2\rho_x \rho_y \rho_{xy} + 2\rho_x \rho_z \rho_{xz} + 2\rho_y \rho_z \rho_{yz} - (\rho_y^2 + \rho_z^2) \rho_{xx} - (\rho_x^2 + \rho_z^2) \rho_{yy} - (\rho_x^2 + \rho_y^2) \rho_{zz}]. \quad (19)$$

The principal curvatures are obtained as

$$\kappa_{\min} = H - \sqrt{H^2 - K}, \quad \kappa_{\max} = H + \sqrt{H^2 - K}. \quad (20)$$

These element-interactive curvatures (EICs) are continuous functions of the density field. In this work, we evaluate EICs at atomic centers to define the element-interactive Gaussian curvature (EIGC) for  $k \neq k'$ :

$$K_{kk'}^{\text{EI}}(\eta_{kk'}) = \sum_{\alpha_i \in \mathcal{T}_k} K_{kk'}(\mathbf{r}_i, \eta_{kk'}). \quad (21)$$

Similar summations define the element-interactive mean curvature  $H_{kk'}^{\text{EI}}$  and the principal curvature descriptors  $\kappa_{kk',\min}^{\text{EI}}$  and  $\kappa_{kk',\max}^{\text{EI}}$ . Because the density  $\rho$  is constructed from  $C^2$  correlation kernels, all derivatives and curvature quantities are analytically evaluated, avoiding numerical differentiation errors.

### 2.2.2 Surface Area and Volume

Let  $\Gamma$  denote the isosurface generated from the density field and  $\Omega$  the enclosed region. The surface area  $A$  and the enclosed volume  $V$  are defined as

$$A = \int_{\Gamma} dS, \quad V = \int_{\Omega} dV. \quad (22)$$

On a Cartesian grid with uniform mesh size  $h$ , these integrals are approximated following the method of [32]:

$$\int_{\Gamma} f(x, y, z) dS \approx \sum_{(i,j,k) \in I_o} \left( f(x_o, y_j, z_k) \frac{|n_{o,x}|}{h} + f(x_i, y_o, z_k) \frac{|n_{o,y}|}{h} + f(x_i, y_j, z_o) \frac{|n_{o,z}|}{h} \right) h^3, \quad (23)$$

where  $(x_o, y_j, z_k)$  denotes the intersection point between the isosurface  $\Gamma$  and the  $x$ -mesh line,  $n_{o,x}$  is the  $x$ -component of the unit normal vector, and  $I_o$  represents the set of irregular grid points where the interface intersects the grid.

The enclosed volume is computed similarly:

$$\int_{\Omega} f(x, y, z) dV = \frac{1}{2} \left( \sum_{(i,j,k) \in I_1} f(x_i, y_j, z_k) h^3 + \sum_{(i,j,k) \in I_1 \cup I_o} f(x_i, y_j, z_k) h^3 \right), \quad (24)$$

where  $I_1$  denotes the set of grid points inside the enclosed region  $\Omega$ .

The surface and volume integrals are numerically evaluated using the discrete differential-geometry framework implemented in the DG-EIM and EISA pipelines [31, 32].

## 2.3 3D Zernike Descriptors

Following [33, 34, 45–47], we represent molecular surfaces using 3D Zernike descriptors (3DZD), which provide a rotation-invariant characterization of molecular shape.

The 3D Zernike–Canterakis basis functions are defined by

$$Z_{nl}^m(r, \theta, \varphi) = R_{nl}(r)Y_l^m(\theta, \varphi), \quad (25)$$

where  $-l \leq m \leq l$ ,  $0 \leq l \leq n$ , and  $n - l$  is even. The functions  $Y_l^m(\theta, \varphi)$  are the spherical harmonics:

$$Y_l^m(\theta, \varphi) = \sqrt{\frac{(2l+1)(l-m)!}{4\pi(l+m)!}} P_l^m(\cos \theta) e^{im\varphi}, \quad (26)$$

where  $P_l^m(\cos \theta)$  are the associated Legendre polynomials,  $\theta$  and  $\varphi$  denote the spherical coordinates, and  $R_{nl}(r)$  is the radial function.

The Zernike–Canterakis basis functions  $Z_{nl}^m(\mathbf{r})$  can be expressed as polynomials in the Cartesian coordinates  $(x, y, z)$ .

For a 3D function  $f(\mathbf{x})$  defined on the unit ball  $\|\mathbf{x}\| \leq 1$ , the 3D Zernike moments are computed as

$$\Omega_{nl}^m = \frac{3}{4\pi} \int_{\|\mathbf{x}\| \leq 1} f(\mathbf{x}) \overline{Z_{nl}^m(\mathbf{x})} dV, \quad (27)$$

where  $\overline{Z_{nl}^m}$  denotes the complex conjugate.

The solvent–excluded surface (SES) of each protein or binding pocket is constructed using a triangulated surface representation. The surface is then voxelized onto a cubic grid with spacing 1 Å, and a binary volumetric function  $f(\mathbf{x})$  is defined as

$$f(\mathbf{x}) = \begin{cases} 1, & \text{if the voxel intersects the SES surface,} \\ 0, & \text{otherwise.} \end{cases} \quad (28)$$

### 3 Methods

Although several geometric descriptor components used in this work were introduced in prior studies [31–33], the present work focuses on their unified integration and systematic comparative evaluation under a common geometric framework. Rather than proposing a single isolated scoring function, this study investigates how different geometric representations, including local curvature, surface area, volume, global shape, chemical partitioning, and ligand-localized pocket morphology, contribute to protein–ligand recognition across complementary supervised and unsupervised tasks.

We first introduce the unified Element Interaction Manifold (EIM) framework, which serves as the central differential-geometric representation investigated throughout this study. EIM features quantify geometric and topological properties of the protein–ligand interface through chemically resolved interaction manifolds. Two complementary extraction regimes are considered, namely global EIM and local EIM. Although both use the same underlying density kernel and curvature formulations, their aggregation strategies and geometric interpretations differ substantially.

### 3.1 Global EIM Features

For each ligand atom type  $\ell$  and protein atom type  $p$ , all ligand atoms of type  $\ell$  and all protein atoms of type  $p$  within a global cutoff distance are collected

$$L_\ell = \{\mathbf{x}_i : \text{ligand atom of type } \ell\}, \quad P_p = \{\mathbf{y}_j : \text{protein atom of type } p\}.$$

The union of these atoms defines a 3D region:

$$\Omega_{\ell,p} = L_\ell \cup \{\mathbf{y}_j \in P_p : \exists \mathbf{x}_i \in L_\ell \text{ such that } \|\mathbf{x}_i - \mathbf{y}_j\| < d_c\}, \quad (29)$$

where  $d_c$  is the predefined cut-off.

**Surface area and volume over multiscale isovalues.** For a sequence of isovalues  $\mathcal{I} = \{0.05, 0.10, \dots, 0.75\}$ , a uniform grid  $G$  with mesh size  $h = 0.5\text{\AA}$  is constructed over the bounding box of  $\Omega_{\ell,p}$ . The normalized density  $\hat{\rho}(\mathbf{x}) = \rho(\mathbf{x})/\rho_{\max}$  is evaluated at each grid point with kernel parameters  $\tau = 1, \kappa = 2$ . Then we compute:

$$\text{SA}_\alpha, \quad \text{Vol}_\alpha, \quad \alpha \in \mathcal{I}.$$

Global features summarize these values using six statistics (sum, mean, median, standard deviation, minimum, and maximum).

**Curvature evaluation.** For every ligand atom position  $\mathbf{x}_i \in L_\ell$ , curvature is evaluated using kernel-based differential operators

$$H(\mathbf{x}_i), \quad K(\mathbf{x}_i), \quad \kappa_{\min}(\mathbf{x}_i), \quad \kappa_{\max}(\mathbf{x}_i).$$

Each curvature vector over all ligand atoms is summarized by the same six statistics. Thus, global features represent pairwise element-type interaction geometry across the entire interface region.

### 3.2 Local EIM Features

In contrast to the global formulation, local EIM constructs features independently around each ligand atom position  $\mathbf{x}_i$ .

**Local atom-centered neighborhood.** For each ligand atom of type  $\ell$ , we define

$$B_i = \{\mathbf{y}_j \in P_p : \|\mathbf{x}_i - \mathbf{y}_j\| < d_c\}, \quad (30)$$

a spherical neighborhood around  $\mathbf{x}_i$  is formed

$$\Omega_i = B_i \cup \{\mathbf{x}_i\}. \quad (31)$$

**Local density field and single isovalue.** A cubic grid centered at  $\mathbf{x}_i$  with  $d_c + 2 \text{ \AA}$  is used on each side to compute the density field  $\hat{\rho}_i$ . Unlike global EIM, local EIM uses a fixed single isovalue:  $\alpha = 0.25$ , and parameters  $\tau = 1, \kappa = 2$ . For this configuration,

$$SA_i, \quad \text{Vol}_i$$

are computed for each ligand atom and aggregated using the same statistics across all ligand atoms of type  $\ell$  for the pair  $(\ell, p)$ .

**Local curvature.** Curvatures are evaluated at the positions of the ligand atoms  $\mathbf{x}_i$ , using only atoms within the local ball

$$H(\mathbf{x}_i), K(\mathbf{x}_i), \kappa_{\min}(\mathbf{x}_i), \kappa_{\max}(\mathbf{x}_i),$$

and concatenated by the same statistics across all ligand atoms of type  $\ell$ .

Thus, local EIM captures fine-grained, atom-centered interaction geometry that reflects the immediate chemical neighborhood surrounding each ligand atom. In contrast, global EIM builds multiscale surface and curvature descriptors over large interaction regions using multiple isovalues, providing a macroscopic view of the binding interface. Together, the two regimes offer complementary structural perspectives and are often combined into hybrid descriptors for improved binding affinity prediction.

### 3.3 Hybrid EIM Features

We evaluate four feature configurations that span different scales of molecular recognition. The Global (1G) configuration employs a single kernel with bandwidth  $\tau = 1.0$ , cutoff of  $12 \text{ \AA}$ , and power of 2.0, generating 1,440 interface-level descriptors that capture overall binding-site geometry and composition. The Local (1L) configuration uses a kernel with bandwidth  $\tau = 1.0$ , cutoff of  $7 \text{ \AA}$ , power of 2.0, and isovalue threshold of 0.25, producing 1,440 contact-level descriptors that encode short-range atom–atom interactions within the binding interface. The Hybrid (1G+1L) configuration concatenates global and local features to form a balanced 2,880-dimensional multi-scale representation. The Two-Kernel Hybrid (2G+2L) extends this approach by incorporating additional kernels at different spatial scales, Global Kernel 2 (GK2,  $\tau = 0.5$ , cutoff  $18 \text{ \AA}$ ) and Local Kernel 2 (LK2,  $\tau = 0.5$ , cutoff  $12 \text{ \AA}$ ), yielding a comprehensive 5,760-dimensional feature space. For curvature, these features capture complementarity between protein and ligand shapes. Element-specific interactions naturally encode chemical information, and curvature evaluated at ligand positions is mechanistically relevant, reflecting how well the ligand fits within the binding pocket. A summary of the kernel settings and resulting feature dimensions is provided in Table 1.

Table 1: Parameters for multi-scale EIM feature configurations

Kernel	Scale	$\tau$	Kernel Power	Cutoff (Å)	Isovalue	Feature Dimension
GK1	Global	1.0	2.0	12	N/A	1,440
GK2	Global	0.5	2.0	18	N/A	1,440
LK1	Local	1.0	2.0	7	0.25	1,440
LK2	Local	0.5	2.0	12	0.25	1,440
Hybrid (1G+1L): GK1 + LK1						2,880
Two-Kernel (2G+2L): GK1 + GK2 + LK1 + LK2						5,760

The complete EIM pipeline, including element-wise decomposition, manifold generation, and feature aggregation, is depicted in Figure 1.

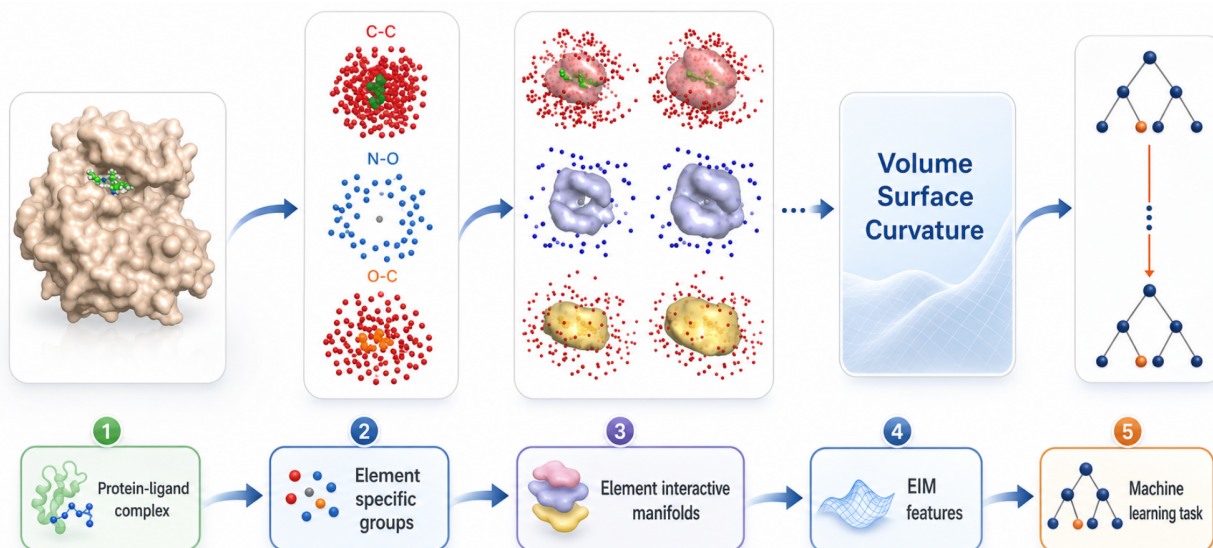


Figure 1: EIM Learning Strategy – Complex 5dwr

### 3.4 Standard 3D Zernike Descriptors Feature Extraction

To establish a baseline for global shape-based molecular comparison, we employ 3D Zernike Descriptors (3DZD). This method provides a compact, rotation-invariant representation of a three-dimensional object by decomposing its geometry into a series of orthogonal volumetric moments. Protein structures are processed through a standardized surface-based computational pipeline as follows:

1. **Surface Generation (EDTSurf):** A molecular surface is generated from each protein structure using the EDTSurf utility, which constructs a triangulated surface representation based on Euclidean distance transforms. Structures that produce missing or degenerate surfaces are excluded from further analysis.

2. **Mesh Format Conversion:** The resulting surface mesh is converted from PLY to OBJ format to ensure compatibility with subsequent voxelization and Zernike moment computation tools.
3. **Voxelization (obj2grid):** The triangulated surface is discretized onto a  $64^3$  cubic occupancy grid using `obj2grid`. During this step, the surface is centered and uniformly scaled to fit within a unit sphere, ensuring translation invariance and scale normalization.
4. **3D Zernike Moment Computation (map2zernike):** The volumetric grid is projected onto the 3D Zernike polynomial basis up to expansion order  $n = 20$  using `map2zernike`. This produces a set of 121 rotation-invariant Zernike descriptors stored in a `.inv` file.

**Mathematical Representation.** Rotation-invariant 3D Zernike descriptors are obtained by computing the  $L_2$ -norm of the complex Zernike moments across all magnetic quantum numbers  $m$  for each pair of expansion order  $n$  and angular momentum  $l$ :

$$F_{nl} = \sqrt{\sum_{m=-l}^l |\Omega_{nl}^m|^2}. \quad (32)$$

The final representation of each protein surface is a fixed-length, 121-dimensional descriptor vector ordered by increasing  $(n, l)$

$$\mathbf{F}_{3DZD} \in \mathbb{R}^{121}, \quad (33)$$

which encodes global shape information in a rotation-invariant manner.

A schematic overview of the 3DZD feature extraction process is shown in Figure 2. The pipeline includes surface generation, voxelization, and projection onto the Zernike polynomial basis.

### 3.5 Element-pair 3D Zernike Descriptors

To investigate whether element-specific variants of spherical harmonic descriptors can recover complementary information beyond standard global 3D Zernike Descriptors (3DZD), we constructed an element-pair 3DZD (EP-3DZD) representation as a chemically partitioned extension of classical 3DZD for controlled geometric comparison.

EP-3DZD extends the standard 3DZD framework by computing Zernike moments independently for element-pair-restricted molecular surfaces, analogous to the element-specific construction used in EIM. Rather than extracting a single global shape descriptor from the entire protein surface, EP-3DZD decomposes each protein–ligand complex into chemically defined subsets based on element identity. All protein–ligand atom pairs within a 12 Å cutoff are identified for the 40 possible element combinations (10 ligand atom types  $\times$  4 protein atom types). For each element pair (e.g., C–C, C–O, N–N), a localized molecular surface is generated and encoded independently using 3D Zernike moments. This design tests whether

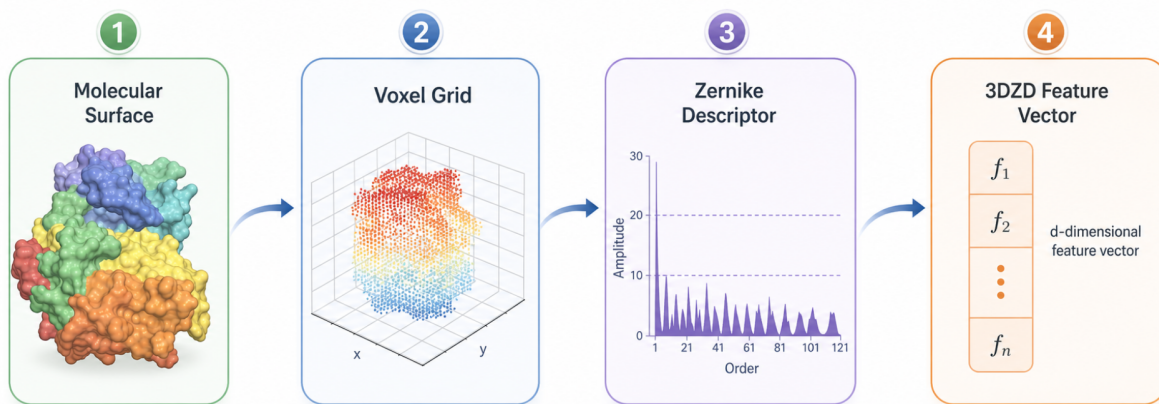


Figure 2: 3DZD Learning Strategy–Complex 5dwr

the incorporation of chemical specificity into Zernike-based representations can improve binding site similarity performance. It combines EIM’s element-aware philosophy with 3DZD’s global geometric formalism.

EP-3DZD extraction proceeds through a three-stage pipeline optimized for high-throughput computation. **Stage 1** performs element-pair decomposition: atoms are filtered by element type, and interacting protein–ligand atom groups within 12 Å are identified for each element combination. Each group is written to a separate, coordinate-centered PDB file to improve numerical stability during surface generation. This procedure yields up to 40 localized structural regions per complex, each representing a chemically homogeneous surface patch with respect to element identity.

**Stage 2** computes 3D Zernike descriptors using the standardized workflow used in the 3D-AF-Surfer toolkit. For each element-pair PDB, (1) **EDTSurf** generates a triangulated molecular surface in PLY format; (2) surfaces are converted from PLY to OBJ format; (3) **obj2grid** discretizes each surface into a  $64^3$  voxel grid using `obj2grid -g 64`, automatically normalizing the structure to fit within a unit sphere; and (4) **map2zernike** computes 3D Zernike moments up to expansion order  $n = 20$ , producing 121 rotation-invariant coefficients per element pair. Concatenation across all 40 element pairs yields a 4,840-dimensional EP-3DZD feature vector. In practice, many element pairs exhibit sparse or null interaction density, resulting in highly sparse feature representations.

**Stage 3** implements an automated pipeline with fault tolerance, checkpoints, and detailed logs to ensure reproducibility across large-scale datasets. Each extraction task is isolated with timeout protection to prevent failure on degenerate surface geometries, and intermediate results are cached to enable efficient restart. The final EP-3DZD descriptors are stored as compressed feature matrices (`ep3dzd_features.npz`), providing chemically interpretable, element-resolved shape fingerprints suitable for direct comparison with EIM-based representations. An overview of the EP-3DZD pipeline is shown in Figure 3. It outlines the decomposition into element-pair-specific regions followed by independent surface encoding and descriptor aggregation.

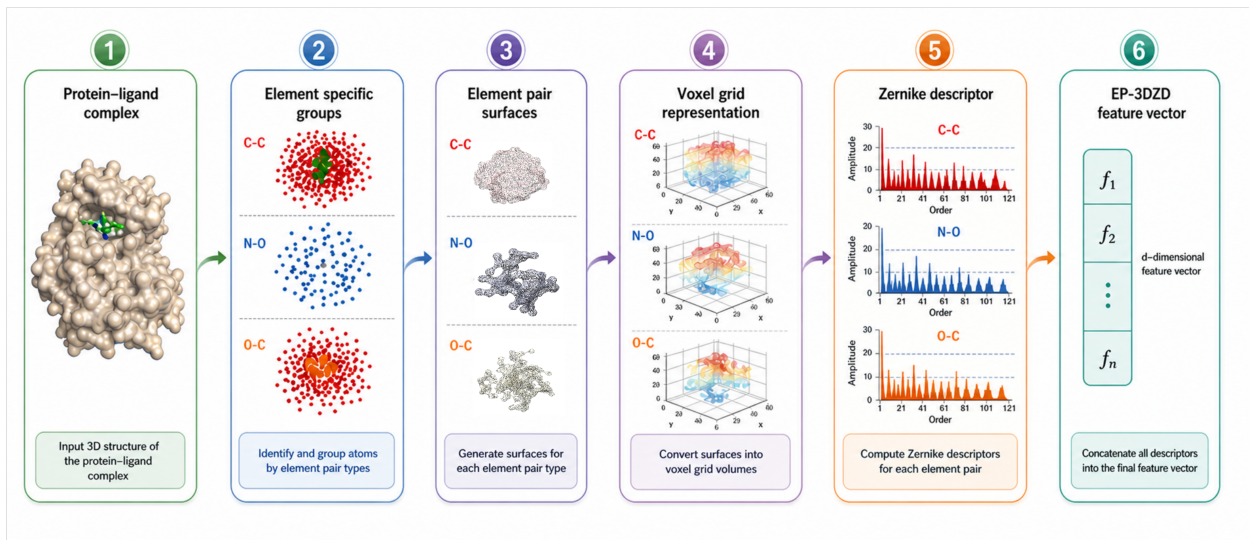


Figure 3: EP-3DZD Learning Strategy-Complex 5dwr

### 3.6 Ligand-Aware ESES Pocket Geometry Baseline

To construct a ligand-aware geometric baseline under the same ligand-identity recognition protocol used for EIM, we adapted the Eulerian solvent excluded surface (ESES) pipeline to compute pocket-level geometric descriptors defined relative to the bound ligand [20]. Our implementation builds upon the publicly available ESES software <https://github.com/rdzhao/ESES>.

Unlike 3DZD, which encodes global protein surface geometry, and unlike RF-Score v2 or EIM, which explicitly model protein-ligand contacts, this baseline relies solely on ligand-conditioned surface geometry. The ligand pose determines which region of the protein surface is summarized, while the extracted features remain purely geometric. From each complex, a 7-dimensional feature vector was constructed from the pre-computed pocket mesh statistics at  $R = 6 \text{ \AA}$ : log-transformed pocket surface area, total mesh volume, and pocket concave area, together with raw pocket concavity fraction, vertex count, face count, and pocket face count.

## 4 Results

The two evaluation tasks considered in this study probe complementary aspects of molecular recognition. Binding affinity prediction evaluates whether geometric descriptors encode quantitative interaction strength, whereas ligand-aware binding-site similarity evaluates whether descriptors preserve transferable interaction signatures across structurally diverse protein targets. Together, these tasks provide a broader assessment of the representational quality of geometric protein-ligand descriptors.

## 4.1 Binding Affinity Prediction Performance

### 4.1.1 Dataset and Evaluation Protocol

All experiments were conducted on the PDBbind v2016 dataset, which comprises 4,057 protein-ligand complexes in the Refined Set and 285 complexes in the CASF-2016 Core Set. Following standard PDBbind benchmarking practice, the Core Set serves as an independent test set, while the Refined Set (excluding Core Set complexes) constitutes the training set. After filtering for data availability and feature completeness across all methods, the final dataset consisted of 3,772 training complexes and 285 test complexes. This train-test split was applied identically across all feature configurations and baseline methods to ensure unbiased performance comparison.

Model performance was quantified using multiple complementary metrics. Binding affinity predictions in pK units ( $pK = -\log_{10} K_d$ , where  $K_d$  is in molar units) were evaluated using root mean squared error (RMSE), mean absolute error (MAE), Pearson correlation coefficient ( $R$ ), and coefficient of determination ( $R^2$ ). For comparison with methods reporting results in kcal/mol, we applied the standard conversion factor of 1.363 kcal/mol per pK unit. Cross-validation was performed using 5-fold splitting on the training set to assess model stability and prevent overfitting. All reported test metrics correspond to predictions on the held-out CASF-2016 Core Set, providing an estimate of generalization performance on unseen complexes.

### 4.1.2 Machine Learning Model

Binding affinity prediction was performed using Gradient Boosting Tree (GBT) with hyperparameters held constant across all feature configurations to ensure fair comparison. The model consisted of 10,000 estimators with a learning rate of 0.01, maximum tree depth of 7, minimum samples per split of 3, and subsample ratio of 0.3. At each split, a random subset of  $\sqrt{n}$  features (where  $n$  is the total feature dimensionality) was considered. The loss function was squared error, and the random seed was fixed at 42 to ensure reproducibility. These hyperparameters were selected to balance model expressiveness with generalization capacity, preventing overfitting while maintaining sufficient flexibility to capture complex interaction patterns.

### 4.1.3 Performance of EIM Features

**Performance of Global EIM Features** To establish a baseline for geometric representation, we first evaluated the Global Element Interaction Manifold (Global EIM) approach. In this configuration, manifold extraction is performed across the entire protein surface using a broader 12.0 Å cutoff for pairwise interactions. The global representation captures the total topological footprint of the protein-ligand complex, providing a macroscopic view of the surface curvature and volumetric density. This allows us to assess whether a holistic geometric description can effectively characterize binding affinity or if such global features are potentially susceptible to structural noise from surface regions distant to the active site. The robustness of the EIM global features was assessed using 5-fold cross-validation on the training dataset. The model exhibited consistent performance across all folds, achieving an

average Pearson correlation coefficient ( $R$ ) of 0.7348 and a Root Mean Square Error (RMSE) of 1.3557. The low standard deviation across folds ( $\pm 0.0119$  for Pearson  $R$ ) indicates that the high-dimensional EIM descriptor (1,440 features) provides a stable representation of the protein-ligand interface without significant overfitting. The final model, trained on the full training set, was evaluated on the CASF-2016 benchmark. The global EIM model achieved a  $R$  value of 0.7836 and an RMSE of 1.4449. These results demonstrate that the multiscale geometric features—integrating surface area, volume, and curvature across element-specific manifolds—capture essential physical characteristics of binding affinity. Analysis of the GBT feature importance reveals the significance of specific chemical interactions. The most influential features were dominated by Carbon-Carbon ( $C-C$ ) and Carbon-Nitrogen ( $C-N$ ) interactions, specifically those involving mean curvature ( $H$ ), minimum curvature ( $\kappa_{\min}$ ), and volumetric descriptors. Notably, the top-ranked feature was the sum of mean curvatures for  $C-C$  pairs, suggesting that the local geometry of hydrophobic interfaces is strongly associated with binding affinity prediction performance.

**Performance of Local EIM Features** To investigate the impact of spatial localization on geometric representations, we evaluated the Local Element Interaction Manifold (Local EIM) approach. In this configuration, manifold extraction is restricted to a local region defined by a 7.0 Å cutoff centered at the binding site, specifically targeting the immediate protein-ligand interface while filtering out distant surface noise. The Local EIM model was trained using Gradient Boosting Trees (GBT) with parameters consistent with our global evaluation. The stability of these localized features was confirmed via 5-fold cross-validation on the training set, yielding an average Pearson  $R$  of 0.7462. This indicates that even within a constrained spatial volume, the element-specific manifolds provide a dense and reliable signal for binding affinity. When applied to the independent CASF-2016 Core Set, the Local EIM model achieved a Pearson  $R$  of 0.8022 and an RMSE of 1.3841. The localizing the manifold extraction led to a significant performance boost over the global approach (Pearson  $R$  0.8022 vs. 0.7836). The coefficient of determination ( $R^2 = 0.5933$ ) suggests that approximately 59% of the variance in experimental binding affinity can be explained by localized differential geometry alone. Analysis of the GBT feature importance reveals a shift in the descriptors prioritized at the interface. Unlike the global model, which favored mean curvature, the local model relies heavily on the surface area sums for ( $C-C$ ) and ( $C-S$ ) pairs, respectively.

**Performance of Hybrid EIM Features** To combine global structural context with fine-grained binding pocket information, we evaluated a Hybrid EIM approach. This model utilizes a concatenated feature vector of 2,880 dimensions, merging 1,440 global and 1,440 local EIM descriptors into a 2,880-dimensional feature vector. The 5-fold cross-validation on the training set yielded an average RMSE of **1.3100**, demonstrating lower prediction error than either standalone representation. The final evaluation on the independent CASF-2016 Core Set reveals that the Hybrid model achieves a Pearson correlation of **0.8079** and an RMSE of **1.3575**. This performance surpasses both the Global EIM ( $R = 0.7836$ ) and the Local EIM ( $R = 0.8022$ ). The improvement indicates a synergistic effect: local manifolds capture high-resolution geometric signals critical for affinity estimation, while global mani-

fold provide structural context related to overall shape and solvent exposure. The feature importance ranking for the Hybrid model reveals a synergistic dependency between local geometric specificity and global volumetric context. The primary predictor is the local  $C-C$  minimum curvature ( $\kappa_{\min}$ ) sum, confirming that fine-grained carbon-carbon packing is the strongest driver of affinity prediction.

**Performance of Hybrid Two-Kernel Features** To determine the optimal spatial and multi-scale configuration for binding affinity prediction, we compared four distinct EIM representations. These range from single-scale global and local descriptors to high-dimensional multiscale fusions. The results on the CASF-2016 Core Set (Table 2) reveal a clear trend: increasing the resolution and multi-scale diversity of the geometric features leads to superior predictive accuracy.

The Hybrid EIM (2,880 features) and the Two-Kernel Fusion (5,760 features) represent our highest-performing configurations. The Two-Kernel Fusion, which integrates four distinct Gaussian kernel bandwidths across both global and local contexts, achieved a Pearson correlation of 0.8074, slightly below the Hybrid model.

Table 2: Performance comparison of EIM representations on the CASF-2016 Core Set.

Model Configuration	Dims	RMSE	MAE	$R^2$	Pearson $R$
Global EIM (Single Scale)	1,440	1.4449	1.1782	0.5569	0.7836
Local EIM (Single Scale)	1,440	1.3841	1.1089	0.5933	0.8022
<b>Hybrid EIM (Global + Local)</b>	<b>2,880</b>	<b>1.3575</b>	<b>1.0837</b>	<b>0.6088</b>	<b>0.8079</b>
Two-Kernel Fusion	5,760	1.3734	1.1148	0.5997	0.8074

The feature importance analysis for the Two-Kernel model reveals a sophisticated multi-scale hierarchy. The top predictor is the LK2  $C-C$  Surface Area sum, followed by the GK1  $C-C$  Volumetric Mean. Notably, the mid-range local kernel (LK2) appears more frequently among the top 10 features than the tight local kernel (LK1), suggesting that geometric information within the 12Å shell provides a richer geometric signal for binding affinity than the immediate 7Å first-shell alone. This multi-scale dependency supports the use of combined differential geometry descriptors to capture both the fine-grained atom-pair contacts and the broader macromolecular context. Table 3 presents top 10 important features for EIM representations.

Table 3: Top 10 most important features for each EIM configuration.

	<b>Global (1G)</b>	<b>Local (1L)</b>	<b>Hybrid (1G+1L)</b>	<b>Two-Kernel (2G+2L)</b>
<b>Rank</b>	Feature	Feature	Feature	Feature
1	C_C_H_sum	local_C_C_surface_area_sum	local_C_C_kappa_min_sum	lk2_C_C_surface_area_sum
2	C_C_volume_max	local_C_S_surface_area_sum	global_C_C_H_sum	gk1_C_C_volume_mean
3	C_C_kappa_min_sum	local_C_C_kappa_min_sum	global_C_N_volume_mean	gk1_C_C_kappa_min_sum
4	C_N_kappa_min_sum	local_C_C_H_sum	global_C_C_volume_median	lk2_C_N_volume_sum
5	C_C_volume_median	local_C_C_volume_sum	local_C_O_kappa_min_sum	lk2_C_C_H_sum
6	C_O_volume_mean	local_C_N_surface_area_sum	local_C_N_volume_sum	lk2_C_N_surface_area_sum
7	H_O_kappa_max_sum	local_H_O_surface_area_sum	local_H_S_volume_sum	lk1_C_O_volume_sum
8	H_O_volume_median	local_C_S_volume_sum	global_C_O_kappa_min_sum	lk1_C_C_surface_area_sum
9	C_N_surface_area_median	local_C_O_surface_area_sum	local_C_C_surface_area_sum	gk1_C_N_volume_max
10	C_O_volume_min	local_H_C_surface_area_sum	global_C_C_kappa_min_sum	lk2_C_C_volume_sum

#### 4.1.4 Performance of 3DZD and EP-3DZD Features

To evaluate the predictive power of EIM, we performed two distinct comparative studies. The first evaluates EIM against traditional global shape descriptors (Surface 3DZD), and the second evaluates it against chemically-stratified shape descriptors (EP-3DZD).

**Case 1: Surface 3DZD** In this evaluation, we benchmarked our method against the standard Surface 3D Zernike Descriptors, which capture the global geometric envelope of the protein surface without chemical distinction.

As shown in Table 4, the traditional Surface 3DZD model exhibits significantly higher error and lower correlation compared to the hybrid models. Notably, the addition of global surface information to the EIM framework (Hybrid EIM + Surface 3DZD) resulted in a negligible  $\Delta$ RMSE of only  $-0.005$  pK ( $-0.4\%$ ) and no improvement in Pearson  $R$  for EIM (it improved 3DZD significantly). This suggests that the local interaction patterns captured by EIM already provide a highly informative representation of binding affinity.

Table 4: Comparison of traditional surface-based 3DZD and Hybrid EIM representations.

<b>Configuration</b>	<b>Dims</b>	<b>CV RMSE</b>	<b>Test RMSE</b>	<b>Test <math>R</math></b>	<b>Test <math>R^2</math></b>
Surface 3DZD	4,840	$1.594 \pm 0.052$	1.594	0.705	0.461
Hybrid EIM + Surface 3DZD	7,720	$1.305 \pm 0.039$	1.361	0.808	0.607

**Case 2: EIM vs. Element-Pair 3DZD (EP-3DZD)** The second evaluation compared EIM against the chemically resolved Element-Pair 3DZD, which utilizes 40 distinct chemical manifolds. This comparison tests whether the limited improvement from traditional 3DZD arises from a lack of chemical specificity.

Table 5 summarizes the results. While EP-3DZD performs better than its element-agnostic counterpart, it still underperforms relative to the combined model. The "Combined" model (EIM + EP-3DZD) achieved a Pearson  $R$  of 0.811, representing a marginal gain of only  $+1.1\%$  over the baseline EIM performance. This weak complementarity suggests that important geometric signals relevant to binding affinity appear to already be captured

by local manifold geometry, rather than by global spherical harmonic representations of chemically partitioned surfaces.

Table 5: Comparative analysis of EP-3DZD and Combined (EIM + EP-3DZD) features.

Method	Dims	CV RMSE	Test RMSE	Test $R$	Test $R^2$
EP-3DZD	4,840	$1.360 \pm 0.036$	1.456	0.779	0.550
Combined (EIM + EP-3DZD)	7,720	$1.306 \pm 0.036$	1.353	0.811	0.611

#### 4.1.5 Performance of ESES Pocket Geometry Features

Five-fold cross-validation on the training set yielded a mean Pearson correlation  $R$  of 0.4950 ( $\pm 0.0152$ ) and a mean RMSE of 1.7849 pK ( $\pm 0.0549$ ), indicating a modest but stable predictive signal across folds. On the Core Set, the ligand-aware ESES pocket model achieved a Pearson correlation  $R$  of 0.5603 ( $R^2 = 0.3140$ ) and an RMSE of 1.8041 pK (2.4590 kcal/mol), with a MAE of 1.4562 pK (1.9848 kcal/mol).

Although the ligand pose localizes the binding pocket and removes uninformative distal surface geometry, the resulting descriptor remains substantially below the Global EIM model ( $R = 0.7836$ , RMSE = 1.4449 pK). This gap suggests that geometric pocket statistics, even when conditioned on the bound ligand pose, do not fully capture the interaction-level specificity encoded by element-specific manifolds. These results suggest that explicit modeling of protein–ligand interaction geometry contributes substantially to competitive binding affinity prediction performance.

As shown in Table 6, the Hybrid EIM achieves near state-of-the-art predictive performance (Pearson  $R = 0.808$ ), while requiring significantly fewer features than higher-dimensional combined models such as EIM + EP-3DZD ( $R = 0.811$ ). This highlights its effectiveness as a compact and efficient representation of protein–ligand interactions.

## 4.2 Binding Site Similarity Evaluation

### 4.2.1 Dataset and Evaluation Protocol

Beyond binding affinity prediction, a key application of molecular interaction descriptors is the identification of similar binding sites across different protein targets. Such similarity metrics enable a range of applications, including polypharmacology analysis, off-target prediction, drug repurposing, and virtual screening against novel targets. Recent advances in molecular recognition have highlighted the importance of explicitly modeling ligand properties alongside protein structure [48]. Motivated by this paradigm, we evaluate binding-site similarity from a ligand-aware perspective. Unlike sequence-based or structure-based methods that ignore ligand context, this framework defines similarity between binding events (protein–ligand complexes) rather than proteins alone. If two binding sites bind the same ligand, they are likely to share compatible geometry and chemistry, even if their sequences or global folds differ substantially.

Table 6: Comparison of geometric, interaction-based, and learning-based descriptors on the CASF-2016 Core Set. Methods are ordered by Pearson correlation ( $R$ ). CP denotes chemically partitioned descriptors. Our proposed methods are shown in bold.

Method	Descriptor Type	Dims	$R$
<b>Hybrid EIM + EP-3DZD</b>	Interaction Geometry + Shape (CP)	7,720	0.811
<b>Hybrid EIM (Global + Local)</b>	Interaction Geometry	2,880	0.808
<b>Hybrid EIM + Surface 3DZD</b>	Interaction Geometry + Shape	7,720	0.808
<b>Two-Kernel Fusion (EIM)</b>	Interaction Geometry	5,760	0.807
<b>Local EIM</b>	Interaction Geometry	1,440	0.802
RF-Score [11]	Distance-based Descriptor	–	0.800
<b>Global EIM</b>	Interaction Geometry	1,440	0.784
Pafnucy [32]	Deep Learning (3D CNN)	–	0.780
<b>EP-3DZD</b>	Shape (CP)	4,840	0.779
Surface 3DZD (global) [33]	Shape	4,840	0.705
ICChem GRIM [37]	Interaction (Pharmacophore)	37	0.671
X-Score [32]	Empirical Scoring	–	0.631
$\Delta$ SAS [32]	Surface Area-based	–	0.625
Autodock Vina [32]	Force Field-based	–	0.604
<b>ESES Pocket Geometry</b>	Surface Geometry	7	0.560

**Dataset Construction** We define a notion of binding-site similarity based on shared ligand identity. Intuitively, two protein binding sites are considered similar if they exhibit the capacity to bind the same ligand molecule. This indicates the presence of conserved interaction patterns and geometric arrangements capable of accommodating a common molecular scaffold, regardless of whether the proteins share high sequence identity.

Formally, let  $\mathcal{P}$  denote the set of protein binding sites and  $\mathcal{L}$  the set of ligands. For sites  $P_i, P_j \in \mathcal{P}$ , we define the binary relation [49]

$$P_i \sim_{\text{lig}} P_j \iff \exists L \in \mathcal{L} \text{ such that } (P_i, L) \text{ and } (P_j, L) \text{ are observed complexes.} \quad (34)$$

This definition does not assume identical binding poses or interaction modes, but rather captures the existence of at least one experimentally observed compatible interaction geometry. Under this definition, similarity is induced by ligand co-binding. This reflects a ligand-aware perspective in which ligand identity serves as a proxy for the structural and chemical environment of the pocket. It is important to note that  $\sim_{\text{lig}}$  is not an equivalence relation, since transitivity does not generally hold. This non-transitive behavior is biologically meaningful because it represents a local proximity relation within the broader landscape of protein–ligand interactions rather than a global categorization.

We constructed balanced evaluation sets from the 4,057 PDBbind v2016 complexes as follows. For each ligand with  $\geq 2$  protein binders, all within-ligand pairwise combinations were enumerated as positive examples. An equal number of negative pairs was randomly sampled from complexes binding different ligands, ensuring 50/50 class balance and preventing trivial classification via label frequency. For computational efficiency, performance was evaluated

on 15,518 randomly sampled test pairs balanced between positive (same ligand) and negative (different ligand) examples, rather than exhaustively evaluating all  $\binom{4057}{2} \approx 8.2 \times 10^6$  possible pairs. Although positive pairs share ligand identity and therefore partially share chemical information, the task remains nontrivial because the corresponding proteins often span diverse folds, binding environments, and functional classes.

**Evaluation Protocol** Given two complexes  $i$  and  $j$  with feature vectors  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$  (e.g., EIM or other descriptor representations), pairwise similarity was computed using cosine similarity:

$$S_{\cos}(i, j) = \frac{\mathbf{x}_i \cdot \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}. \quad (35)$$

This similarity score ranges from  $-1$  to  $1$ , with higher values indicating greater similarity, and serves directly as a classification score without supervised learning. Performance was quantified using the area under the receiver operating characteristic curve (AUC), computed from cosine similarity scores between feature vectors. This evaluation measures the intrinsic ability of the descriptors to distinguish ligand-compatible binding environments without supervised training or parameter optimization.

#### 4.2.2 Performance of Hybrid EIM Features

EIM achieves strong unsupervised ligand identity prediction performance without supervised learning (AUC = 0.8476). Complexes binding the same ligand exhibit systematically higher feature-space similarity than complexes binding different ligands, even across diverse protein folds and binding pocket architectures. These results indicate that EIM captures interaction signatures that generalize across distinct structural contexts while preserving ligand-specific geometric and chemical information.

To assess whether performance varies systematically with ligand chemistry, we analyzed AUC scores stratified by ligand identity for all ligand types with  $\geq 10$  protein binders in the filtered configuration. Performance is highly heterogeneous across ligand categories, ranging from near-perfect discrimination (AUC = 1.0000 for BTN, AZM, and PGA) to moderate performance for structurally diverse ligand families. The top-performing ligands (Table 7; top 10 of 21 categories shown) share two key characteristics: structural rigidity (e.g., biotin [BTN] and benzamidine [BEN]) and chemically distinctive functional group patterns that are uncommon among other ligands.

Notably, nucleotide cofactors (ATP, ADP, and GDP) achieve near-perfect discrimination (AUC > 0.96) despite binding to diverse protein families spanning kinases, GTPases, and metabolic enzymes. This suggests that EIM captures conserved interaction geometries associated with the adenine/guanine scaffold and phosphate groups, independent of the broader protein structural context.

The mean AUC across all 21 ligand categories is  $0.9108 \pm 0.1132$ , illustrating substantial variability in ligand discriminability. Lower-performing ligands tend to be small and flexible molecules containing common pharmacophores (e.g., carboxylates, primary amines, and simple aromatic rings) that recur across many chemically distinct ligands, making them more difficult to distinguish using interaction geometry alone. This performance stratification further supports that EIM captures chemically meaningful interaction signatures: ligands

with unique three-dimensional interaction geometries and distinctive scaffolds are readily distinguished, whereas ligands sharing common interaction motifs require more subtle discrimination that approaches the limits of geometric descriptors relying only on element-level chemical information.

Table 7: Per-ligand AUC for ligand identity prediction (filtered configuration, ligands with  $\geq 10$  binders). Top 10 of 21 categories shown. Ligand names obtained from RCSB PDB Chemical Component Dictionary.

Rank	Code	Ligand Name	Binders	AUC
1	BTN	Biotin	10	1.0000
2	AZM	Acetazolamide	12	1.0000
3	PGA	Pteridine-6-carboxylic acid	12	1.0000
4	478	N-Cyclohexyl-N-methylglycine	14	0.9983
5	017	(S)-4-Methyl-2-oxopentanoic acid	23	0.9940
6	BEN	Benzamidine	11	0.9848
7	GDP	Guanosine-5'-diphosphate	21	0.9824
8	ADP	Adenosine-5'-diphosphate	27	0.9780
9	GLU	Glutamic acid	15	0.9687
10	ATP	Adenosine-5'-triphosphate	20	0.9686
<i>Mean <math>\pm</math> SD across 21 ligand categories:</i>				0.9108 $\pm$ 0.1132

### 4.2.3 Performance of 3D Zernike Features

Ligand identity prediction using 3D Zernike Descriptors (3DZD) followed the same protocol established for EIM. Pairwise similarity was computed using cosine similarity between feature vectors, and ligand identity classification was evaluated directly from the resulting similarity scores without supervised learning.

EIM (Hybrid 1G+1L, 2,880 dimensions) achieved an AUC of 0.8476, whereas 3DZD (order 20, 121 dimensions) achieved an AUC of 0.5668. This corresponds to an absolute improvement of +0.2808 AUC (49.5% relative increase over 3DZD).

Relative to random classification (AUC = 0.5), 3DZD achieves a 13.4% improvement, while EIM achieves a 69.5% improvement, highlighting the substantially stronger discriminative capability of element-specific differential geometry compared to global spherical harmonic-based shape descriptors.

### 4.2.4 Performance of Element-Pair 3D Zernike (EP-3DZD) Features

We further compared EIM against Element-Pair 3D Zernike Descriptors (EP-3DZD), a shape-based representation that extends classical 3D Zernike descriptors by incorporating element-pair specificity. EIM achieves an AUC of 0.8476, outperforming EP-3DZD (AUC = 0.8320) by an absolute margin of 0.0156 AUC, corresponding to a relative improvement of 1.88% over EP-3DZD.

While both methods substantially exceed random classification ( $AUC = 0.50$ ), the consistent advantage of EIM indicates that element-specific differential geometry descriptors capture ligand-discriminative binding-site characteristics more effectively than purely shape-based representations.

Combining EP-3DZD with EIM features does not yield further performance gains, with combined models achieving AUC values comparable to EIM alone (e.g., 0.8478 versus 0.8476). This suggests that global shape descriptors provide limited complementary information beyond interaction-aware geometric features.

#### 4.2.5 Performance of Ligand-Aware ESES Pocket Geometry Features

The ligand-aware ESES pocket geometry model, using a pocket radius of  $R = 6 \text{ \AA}$  and cosine similarity on z-scored feature vectors, achieves an AUC of 0.7125. Although substantially above random classification ( $AUC = 0.50$ ), its performance remains considerably lower than both RF-Score (0.7726) and EIM (0.8476).

This performance gap suggests that geometric pocket statistics alone, including surface area, concavity, and mesh-scale descriptors, do not sufficiently encode the interaction-level specificity required for robust ligand discrimination, even when localized by the bound ligand pose. The results further indicate that explicit modeling of protein–ligand interaction structure, rather than pocket morphology alone, appears necessary to achieve high ligand-identity discrimination.

Table 8 gives complete binding-site similarity performance comparison.

Table 8: Comprehensive comparison of binding-site similarity performance. Our developed methods are in bold.

Method	Representation	Dim.	Pairs	AUC
<b>Combined (EIM + EP-3DZD)</b>	Hybrid geometry + spherical harmonics	7,720	15,518	0.8478
<b>EIM (Hybrid 1G+1L)</b>	Element-specific differential geometry	2,880	15,518	0.8476
<b>EP-3DZD</b>	Element-pair spherical harmonics	4,840	15,518	0.8320
RF-Score v2 [11]	Binned atom-type contact histogram	216	15,518	0.7726
<b>Ligand-Aware ESES</b>	Ligand-localized pocket geometry	7	15,518	0.7125
IChem GRIM [37]	Protein–ligand interaction graph	37	15,518	0.6274
3DZD (order 20) [33]	Global spherical harmonics	121	15,518	0.5668

## 5 Conclusion

This work presented a systematic geometric analysis of protein–ligand representations across both binding affinity prediction and ligand-aware binding-site similarity tasks. Rather than evaluating a single descriptor family in isolation, this work examined how local curvature, surface area, volume, global shape, chemical partitioning, and interaction-level representations contribute to molecular recognition.

By integrating previously developed element-interactive curvature, surface-area, and volume descriptors into a common Element Interaction Manifold (EIM) framework, we showed that chemically resolved local differential geometry consistently provides highly informative signals for both supervised binding affinity prediction and unsupervised ligand-identity recognition. In particular, local interaction geometry substantially outperformed global spherical harmonic shape descriptors, ligand-localized pocket morphology, and explicit interaction graph matching in cross-target ligand recognition tasks.

To further investigate the role of global shape representations, we introduced element-pair 3D Zernike descriptors (EP-3DZD), a chemically partitioned extension of classical 3DZD. Although EP-3DZD improves substantially over standard global 3DZD, its combination with EIM provides only limited additional gains, suggesting that much of the relevant interaction information is already captured by local differential geometry. Similarly, ligand-aware ESES pocket geometry and contact histogram descriptors achieve moderate performance but remain substantially below interaction-aware geometric representations.

Taken together, these results suggest that local, chemically resolved interaction geometry provides a compact and transferable representation of protein–ligand recognition. More broadly, this work provides a systematic benchmark of geometric representations for protein–ligand modeling and clarifies the relative roles of curvature, surface area, volume, global shape, and interaction topology in molecular recognition tasks.

Several limitations remain. The present study focuses on static protein–ligand complexes and does not explicitly model conformational dynamics, electrostatics, or solvent effects beyond geometry-induced manifolds. In addition, although EIM descriptors are interpretable and computationally efficient relative to high-dimensional voxel representations, further optimization may be beneficial for large-scale virtual screening. Future work will explore integration with equivariant neural architectures, protein language models, and diffusion-based docking frameworks, as well as applications to protein–protein and protein–nucleic acid interactions.

Overall, the present work suggests that chemically resolved differential geometry provides a promising and interpretable foundation for next-generation geometric representations in structure-based molecular modeling.

## Acknowledgements

This work is supported in part by funds from the National Science Foundation (NSF: # 2516126, # 2151802, and # 2534947).

## Author contributions

A.S., M.M.R., and D.D.N. prepared the manuscript. A.S. and D.D.N. finalized the manuscript. A.S. collected the data and performed the analysis. D.D.N. supervised the project.

## Availability of data and materials

Codes for this study, including implementations of EIM-Score and EP-3DZD for protein–ligand binding affinity prediction and binding site similarity analysis, is available at: <https://github.com/MathIntelligence/Element-Interaction-Manifolds>.

## Declarations

**Competing interests** The authors declare no competing interests.

## References

- [1] P. J. Ballester and J. B. O. Mitchell. A machine learning approach to predicting protein–ligand binding affinity with applications to molecular docking. *Bioinformatics*, 26(9):1169–1175, 2010.
- [2] Y. Cao and L. Li. Improved protein–ligand binding affinity prediction using a curvature-dependent surface model. *Bioinformatics*, 30(12):1674–1680, 2014.
- [3] L. Dong, X. Qu, Y. Zhao, and B. Wang. Prediction of binding free energy of protein–ligand complexes with a hybrid molecular mechanics/generalized born surface area and machine learning method. *ACS Omega*, 6(48):32938–32947, 2021.
- [4] M. D. Eldridge, C. W. Murray, T. R. Auton, G. V. Paolini, and R. P. Mee. Empirical scoring functions: I. the development of a fast empirical scoring function to estimate the binding affinity of ligands in receptor complexes. *Journal of Computer-Aided Molecular Design*, 11(5):425–445, 1997.
- [5] H. Li, K.-S. Leung, M.-H. Wong, and P. J. Ballester. Improving autodock vina’s protein–ligand binding affinity prediction with machine learning. *Molecular Informatics*, 34(2–3):115–126, 2015.
- [6] X. Liu, H. Zhao, R. Zhao, X. Feng, and K. Xia. Dcml: Deep cross-molecular learning for protein–ligand binding affinity prediction. *PLoS Computational Biology*, 18(4):e1009943, 2022.
- [7] M. M. Stepniewska-Dziubinska, P. Zielenkiewicz, and P. Siedlecki. Development and evaluation of a deep learning model for protein–ligand binding affinity prediction. *Bioinformatics*, 34(21):3666–3674, 2018.

- [8] J. Wee and K. Xia. Ollivier persistent ricci curvature models for the protein–ligand binding affinity prediction. *Journal of Chemical Information and Modeling*, 61(4):1617–1626, 2021.
- [9] Renxiao Wang, Luhua Lai, and Shaomeng Wang. Further development and validation of empirical scoring functions for structure-based binding affinity prediction. *Journal of Computer-Aided Molecular Design*, 16(1):11–26, 2002.
- [10] Oleg Trott and Arthur J. Olson. Autodock vina: Improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading. *Journal of Computational Chemistry*, 31(2):455–461, 2010.
- [11] Maciej Wójcikowski, Pedro J. Ballester, and Paweł Siedlecki. Performance of machine-learning scoring functions in structure-based virtual screening. *Scientific Reports*, 7:46710, 2017.
- [12] Guo-Bo Li, Ling-Ling Yang, Wen-Jing Wang, Lin-Li Li, and Sheng-Yong Yang. Id-score: A new empirical scoring function based on a comprehensive set of descriptors related to protein–ligand interactions. *Journal of Chemical Information and Modeling*, 53(3):592–600, 2013.
- [13] G.-B. Li, L.-L. Yang, W.-W. Fu, L.-H. Liang, L. Xing, and S.-Y. Yang. Id-score: A new empirical scoring function based on a comprehensive set of descriptors related to protein–ligand interactions. *Journal of Chemical Information and Modeling*, 53(3):592–600, 2013.
- [14] P. W. Bates, Z. Chen, Y. Sun, G.-W. Wei, and S. Zhao. Geometric and potential driving formation and evolution of biomolecular surfaces. *Journal of Mathematical Biology*, 59(2):193–231, 2009.
- [15] C. A. S. Bergström, M. L. Strafford, L. Lazorova, A. Avdeef, K. Luthman, and P. Artursson. Absorption classification of oral drugs based on molecular surface properties. *Journal of Medicinal Chemistry*, 46(4):558–570, 2003.
- [16] S. L. Chan and E. O. Purisima. Molecular surface generation using marching tetrahedra. *Journal of Computational Chemistry*, 19(11):1268–1277, 1998.
- [17] W. Chen, J. Zheng, and Y. Cai. Kernel modeling for molecular surfaces using a uniform solution. *Computer Aided Design*, 42(3):267–278, 2010.
- [18] S. Daberdaku and C. Ferrari. Voxelised representations of macromolecular surfaces and application to protein–protein interaction site prediction. *International Journal of High Performance Computing Applications*, 32(3):407–432, 2018.
- [19] R. Egan and F. Gibou. Level-set formulation for the construction of solvent-excluded surfaces for large biomolecules. *Journal of Computational Physics*, 374:91–120, 2018.
- [20] F. M. Richards. Areas, volumes, packing and protein structure. *Annual Review of Biophysics and Bioengineering*, 6(1):151–176, 1977.

- [21] W. L. Koltun. Precision space-filling atomic models. *Biopolymers*, 3(6):665–679, 1965.
- [22] M. F. Sanner, A. J. Olson, and J. C. Spehner. Reduced surface: An efficient way to compute solvent-excluded surfaces. *Biopolymers*, 38(3):305–320, 1996.
- [23] J. Liang, H. Edelsbrunner, P. Fu, P. V. Sudhakar, and S. Subramaniam. Analytical shape computation of macromolecules: I. molecular area and volume through alpha shape. *Proteins*, 33(1):1–17, 1998.
- [24] D. Xu and Y. Zhang. Generating triangulated macromolecular surfaces by euclidean distance transform. *PLoS One*, 4(12):e8140, 2009.
- [25] Z. Chen, N. A. Baker, and G.-W. Wei. Differential geometry based solvation model i: Eulerian formulation. *Journal of Computational Physics*, 229(22):8231–8258, 2010.
- [26] Z. Chen, N. A. Baker, and G.-W. Wei. Differential geometry based solvation model ii: Lagrangian formulation. *Journal of Mathematical Biology*, 63(6):1139–1200, 2011.
- [27] Z. Chen and G.-W. Wei. Differential geometry based solvation model iii: Quantum formulation. *Journal of Chemical Physics*, 135(19):194108, 2011.
- [28] P. W. Bates, G.-W. Wei, and S. Zhao. The minimal molecular surface. *arXiv*, 2006. q-bio/0610038.
- [29] P. W. Bates, G.-W. Wei, and S. Zhao. Minimal molecular surfaces and their applications. *Journal of Computational Chemistry*, 29(3):380–391, 2008.
- [30] D. D. Nguyen, T. Xiao, M. Wang, and G.-W. Wei. Rigidity strengthening: A new mechanism for protein–ligand binding affinity prediction. *Journal of Chemical Information and Modeling*, 57(7):1715–1721, 2017.
- [31] D. D. Nguyen and G.-W. Wei. Dg-gl: Differential geometry-based geometric learning of molecular forms and interactions. *International Journal for Numerical Methods in Biomedical Engineering*, 35(3):e3179, 2019.
- [32] M. M. Rana and D. D. Nguyen. Eisa-score: Element interactive surface area score for protein–ligand binding affinity prediction. *Journal of Chemical Information and Modeling*, 62(18):4329–4341, 2022.
- [33] D. Kihara, L. Sael, R. Chikhi, and J. Esquivel-Rodriguez. Molecular surface representation using 3d zernike descriptors for protein shape comparison and docking. *Current Protein and Peptide Science*, 12(6):520–530, 2011.
- [34] V. Venkatraman, L. Sael, and D. Kihara. Potential for protein surface shape analysis using zernike descriptors and associated feature extraction techniques. *Proteins*, 76(2):384–398, 2009.
- [35] W.-H. Shin and D. Kihara. Pl-patchsurfer3: Improved structure-based virtual screening for structure variation using 3d zernike descriptors. *bioRxiv*, 2024.

- [36] T. Yacoub, C. Depenveiller, A. Tatsuma, T. Barisić, E. Rusakov, U. Göbel, et al. Shrec 2025: Protein surface shape retrieval including electrostatic potential. *Computers & Graphics*, 132:104394, 2025.
- [37] F. Da Silva, J. Desaphy, and D. Rognan. Ichem: A versatile toolkit for detecting, comparing, and predicting protein–ligand interactions. *ChemMedChem*, 13(6):507–510, 2018.
- [38] Marta M. Stepniewska-Dziubinska, Piotr Zielenkiewicz, and Pawel Siedlecki. Development and evaluation of a deep learning model for protein–ligand binding affinity prediction. *Bioinformatics*, 34(21):3666–3674, 2018.
- [39] Kenneth Atz, Francesca Grisoni, and Gisbert Schneider. Geometric deep learning on molecular representations. *Nature Machine Intelligence*, 3(12):1023–1032, 2021.
- [40] Fabian Fuchs, Daniel Worrall, Volker Fischer, and Max Welling. Se(3)-transformers: 3d roto-translation equivariant attention networks. In H. Larochelle, M. Ranzato, R. Hassel, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1970–1981. Curran Associates, Inc., 2020.
- [41] Víctor Garcia Satorras, Emiel Hoogeboom, and Max Welling. E(n) equivariant graph neural networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 9323–9332. PMLR, July 2021.
- [42] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.
- [43] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- [44] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking, 2023.
- [45] B. Hu, X. Zhu, L. Monroe, M. G. Bures, and D. Kihara. Pl-patchsurfer: A novel molecular local surface-based method for exploring protein–ligand interactions. *International Journal of Molecular Sciences*, 15(9):15122–15145, 2014.
- [46] H. Mirzaei and D. Kihara. Application of deep learning to protein surface analysis using 3d zernike-based representations. *Bioinformatics*, 36, 2020. Supplement 1, i30–i38.
- [47] L. Sael, D. La, B. Li, R. Rustamov, and D. Kihara. Rapid comparison of protein surfaces using geometric 3d zernike descriptors. *Proteins*, 388(2):457–466, 2009.

- [48] Z. Zhang, L. Quan, J. Wang, L. Peng, Q. Chen, B. Zhang, L. Cao, Y. Jiang, G. Li, L. Nie, T. Wu, and Q. Lyu. Labind: identifying protein binding ligand-aware sites via learning interactions between ligand and protein. *Nature Communications*, 16:7712, 2025.
- [49] Y. Chen, R. Tolbert, A. Aronov, G. McGaughey, W. Walters, and L. Meireles. Prediction of protein pairs sharing common active ligands using protein sequence, structure, and ligand similarity. *Journal of Chemical Information and Modeling*, 56:1734–1745, 2016.