



geopocket

Tutorial, Theory and Analysis

Samuel Escudero and Llorenç Villalonga

SBI – PYT

Bioinformatics for Health Sciences 2025

INDEX

A. Tutorial: Geopocket installation and use instructions

1. Download the necessary scripts
2. Package installation and environment
3. Running geopocket program from the terminal
4. Files generated by the geopocket program
5. Visualization of geopocket with Chimera or Pymol

B. Theory: Background on the prediction of protein binding site using geometry approach

1. Difference of Gaussians (DoG) Filter
2. DBSCAN Clustering and Geometric Scoring
3. Interaction Scoring and Composite Score

C. Analysis: Comparing 3 types of protein cavities

1. 1HVR: Human immunodeficiency virus type 1 protease, class: Cleft/Groove
2. 1BRP: Human Retinol Binding Protein, class: Invagination
3. 1N3U : Human Heme Oxygenase-1, class = Tunnel

D. Conclusion

E. References

A.Tutorial: Instructions on the installation and the use of Geopocket

1. Download the necessary scripts

All the material is available in our GitHub repository. For GeoPocket to work properly, it is recommended to clone the whole repository to have:

- The main prediction file.
- The setup.py file and the src/ folder with the code,
- The environment.yml file contains the necessary dependencies.

PD: For the visualization it is necessary that you have previously installed in your operating system of preference, the programs Chimera and Pymol.

To clone the repository:

```
git clone https://github.com/llorensvilla/geopocket.git
cd geopocket
```

2. Package installation and environment

An environment.yml file with the minimum dependencies (Python 3.9, NumPy, SciPy, scikit-learn, etc.) needed to run GeoPocket is included in the repository.

1. Create the working environment with conda:

```
conda env create -f environment.yml
conda activate geopocket
```

2. Install geopocket as a local package (using setup.py):

```
pip install .
```

After these steps, a command called geopocket will be generated and you can run it from any folder (as long as the geopocket environment is active).

3. Running geopocket program from the terminal

To predict binding sites on a given protein, you need to have the PDB file of the protein.

To test geopocket, in the repository you have cloned there is a folder called Examples where there are different proteins to test the program.

Suppose your file is called 1hvr.pdb . Then, you just need to:

```
cd Examples/  
geopocket 1hvr.pdb
```

Usage: you can execute geopocket as a system utility with the following parameters:

```
geopocket [-h] [-t THRESHOLD] [-e EPS] [-m MIN_SAMPLES] [-v] INPUT
```

GeoPocket: geometry based binding-site predictor	
positional arguments:	
INPUT	Input PDB file
optional arguments:	
-h, --help	show this help message and exit
-t THRESHOLD, --threshold THRESHOLD	DoG threshold percentile (default: 95.0)
-e EPS, --eps EPS	DBSCAN neighborhood radius eps (default: 0.8 Å).
-m MIN_SAMPLES, --min-samples MIN_SAMPLES	DBSCAN minimum samples per cluster (default: 5).
-v, --verbose	Print extra progress information

PD: if you want to change others parameters, you can use the following command, but you must take into account where you run the command from, because if you run it inside the 1hvr_results folder it will generate another new file inside that folder, however if you run it from the examples folder, the files will be generated in the same 1hvr_results folder but generating new files for the new parameters. This approach

helps you to compare the files with different parameters that are generated by being able to differentiate them from higher to lower.

For example you can try:

```
geopocket Examples/1hvr.pdb -t 97 -e 1.0 -m 10 -v
```

4. Files generated by the geopocket program

The files that are generated are the following:

-The file of the protein used as input in .pdb is copied to the directory. **Example:** 1hvr.pdb

-Files corresponding to the amount of pockets predicted by the program in .pdb format (inside each file you will find the residues corresponding to each pocket, you can open each file in text format to visualize the residues into each pocket).

Example: pocket_1_residues_1.pdb. **Pd:** This file is in .pdb since it is necessary to open the pockets with the visualization programs, however in the first approach we recommend to open these files in .txt format to visualize the residues that are part of each pocket.

-File corresponding to the pockets that were predicted, including relevant information for each pocket such as location, occupancy, CompositeScore, Volume, SurfaceArea, etc. **Example:** predicted_pockets_1.pdb.

-File corresponding to the visualization in chimera with .cmd ending. **Example:** visualize_pockets_1.cmd.

-File corresponding to the visualization in pymol with .pml ending. **Example:** visualize_pockets_1.pml.

```
(base) samuel-escudero@Sam-E18:~/Documents/PYT/practice_1/1_block_/geopocket/Analysis_geopocket/1hvr_results$ ls
1hvr.pdb          pocket_18_residues_1.pdb  pocket_6_residues_1.pdb
pocket_10_residues_1.pdb  pocket_19_residues_1.pdb  pocket_7_residues_1.pdb
pocket_11_residues_1.pdb  pocket_1_residues_1.pdb   pocket_8_residues_1.pdb
pocket_12_residues_1.pdb  pocket_20_residues_1.pdb  pocket_9_residues_1.pdb
pocket_13_residues_1.pdb  pocket_21_residues_1.pdb  predicted_pockets_1.pdb
pocket_14_residues_1.pdb  pocket_2_residues_1.pdb   visualize_pockets_1.cmd
pocket_15_residues_1.pdb  pocket_3_residues_1.pdb   visualize_pockets_1.pml
pocket_16_residues_1.pdb  pocket_4_residues_1.pdb
```

Figure 1. Files generated by geopocket.

5. Visualization of geopocket with Chimera or Pymol

At the end of Geopocket execution, a results folder is created with the name <protein_name>_results (e.g. 1hvr_results). This folder contains several files, including a script for Chimera (with extension .cmd).

To open the result correctly in Chimera, follow these steps:

1. Enter into the results folder. For example:

```
cd 1hvr_results/
```

2. Run Chimera with the generated script. *You can do the same with pymol.

```
chimera visualize_pockets_1.cmd
```

An alternative without enter into the result carpet, you can do:

```
pymol 1hvr_results/visualize_pockets_1.pml
```

***Recall:** If you want to open a file in chimera you need the file .cmd if you want a pymol visualization you need a file .pml

When loading this script, Chimera and Pymol:

- Will open the original protein as a surface (mesh).
- It will highlight the detected pockets and display the residues near each pocket with semi-transparent colors.

B. Theory: Background on the Geopocket prediction of protein binding site using geometry approach

GeoPocket is a geometry-based binding site predictor that combines surface and cavity detection with residue interaction scoring. It relies on classical geometric detection algorithms to identify concave regions on protein surfaces (Simões et al., 2017) and complements them with interaction propensities derived from surface triplet analysis (Mehio et al., 2010). First, a Difference of Gaussians (DoG) filter is applied to a 3D grid surrounding the protein: two Gaussian blurs at different scales (σ_1 and σ_2) are subtracted to accentuate regions of negative curvature, effectively highlighting potential pockets by enhancing concave features and suppressing noise. Next, DBSCAN clustering delineates contiguous pockets by grouping adjacent DoG positive voxels according to a user tunable neighbourhood radius and minimum sample count, which controls how finely or coarsely the surface indentations are segmented. Finally, each pocket cluster is scored using a composite metric: the geometric component evaluates shape and size (volume and surface area ratios), while the interaction component quantifies the likelihood of productive contacts by weighting counts of nearby residues engaged in hydrogen bonds, ionic interactions, hydrophobic contacts, and aromatic stacking. By integrating these two components, GeoPocket discriminates between superficial indentations and biologically relevant binding sites (Simões *et al.*, 2017; Mehio *et al.*, 2010; Honavar & Uhríková, 2008). The workflow comprises three main stages: (1) Difference of Gaussians (DoG) filtering to highlight concave regions on a 3D grid, (2) DBSCAN clustering to delineate contiguous pockets, and (3) scoring each cluster by combining a geometric score with an interaction score.

1. Difference of Gaussians (DoG) Filter

A 3D grid of voxels is overlaid around the protein structure, and each grid point density is computed by summing Gaussian kernels centered at every atom (σ_1 and σ_2). Subtracting the broader Gaussian (σ_2) from the narrower one (σ_1) emphasizes regions of rapid curvature change, effectively highlighting cavities (Lindow & Goede, 2007). Grid points with DoG values above a user defined percentile threshold (`dog_threshold_percentile`) are retained as candidate pocket points.

2. DBSCAN Clustering and Geometric Scoring

After selecting high value DoG grid points, GeoPocket applies DBSCAN to group nearby points into candidate pockets based solely on spatial density (Lindow & Goede, 2007). Two parameters govern this process: **eps**, the maximum distance between points to be considered neighbors, and **min_samples**, the minimum number of points required to form a cluster. Points that fall within eps of each other and meet the min_samples criterion are merged into a single cluster, while points that fail to meet these density requirements are treated as noise. For each cluster, a convex hull is constructed to approximate its three dimensional shape, yielding measurements of volume and surface area. These two quantities are combined into a single **geometric score** by taking the ratio of volume to surface area and scaling the result to lie between 0 and 1. This normalization allows pockets of different sizes and shapes to be compared directly.

3. Interaction Scoring and Composite Score

To assess which geometrically defined pockets are most likely to bind ligands, GeoPocket evaluates the chemical environment around each pocket centroid (Honavar & Uhríková, 2008). All protein residues within a user-specified distance of the centroid are identified, and five types of interactions: hydrogen bonds, ionic interactions, metal coordination, hydrophobic contacts, and aromatic stacking are counted according to empirical propensity weights. These counts, multiplied by their respective weights, are summed and then normalized by the maximum possible weighted count to produce an **interaction score** between 0 and 1. Finally, the **composite score** for each pocket is calculated as a weighted average of its geometric and interaction scores, allowing users to emphasize either shape based concavity or chemical complementarity by adjusting the relative weights.

In practice, four key parameters control the total number of predicted pockets: the **dog_threshold_percentile** determines how stringent the cutoff is on DoG values (higher percentiles retain only the deepest cavities), **grid_spacing** sets the voxel size of the sampling grid (larger voxels yield a coarser, less crowded point cloud), **eps** in DBSCAN defines how far apart points can be to belong to the same cluster

(increasing `eps` merges nearby clusters into fewer, larger pockets, while decreasing `eps` splits them), and **`min_samples`** specifies the smallest cluster size allowed (raising `min_samples` filters out small, potentially spurious pockets). By adjusting these parameters, users can balance sensitivity (detecting small or shallow pockets) against specificity (focusing only on the most prominent cavities).

C. Analysis: Comparing 3 types of protein cavities

Protein surface cavities can be broadly categorized based on their geometric features and the manner in which they interact with potential ligands. A **void** is a small, isolated indentation that does not penetrate deeply into the protein core. A **cleft or groove** is an elongated depression on the surface, often located between secondary-structure elements (between α -helices or β -sheets). An **invagination** is a deeper pocket formed by the folding of the protein surface, resembling a shallow bowl. A **tunnel** extends through the protein, connecting two distinct surface points, whereas a **channel** is a continuous conduit that links interior cavities with the external environment (Simões *et al.*, 2017)(Figure 2).

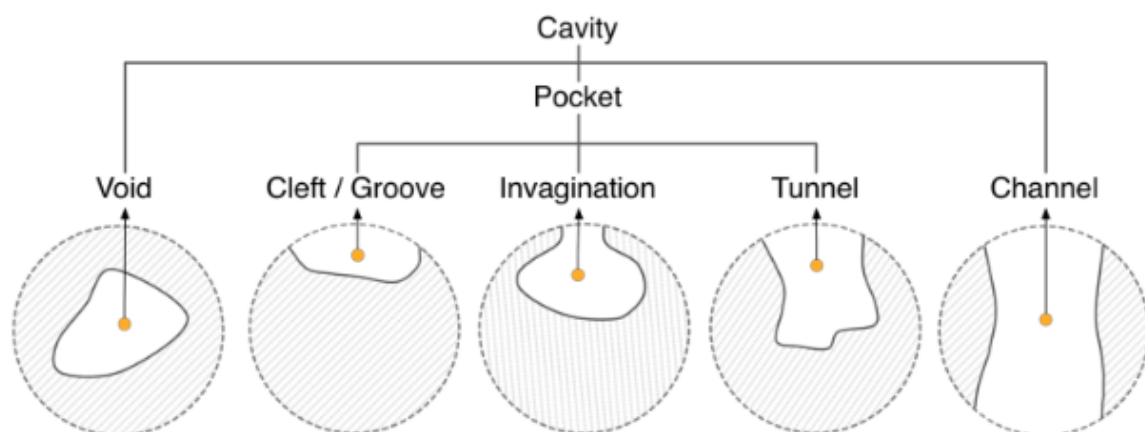


Figure 2 . Types of protein pocket cavities.. This figure is from: Simões *et al.*, 2017 .

For the present analysis, we will focus on three cavity types cleft/groove, invagination, and tunnel by selecting one representative protein for each category. A cleft/groove pocket typically engages ligands through complementary shape and

hydrogen-bonding interactions, whereas an invagination pocket often accommodates substrates or cofactors that fit within its deeper recess. A tunnel pocket forms a continuous conduit through the protein core, guiding ligands or substrates between distal sites.. By examining one protein-ligand complex (bind) and one apo form (unbind) for each cavity type, we aim to elucidate how geometric features influence binding site detection and ligand accommodation (Simões *et al.*, 2017).

WORKFLOW

Before presenting the three detailed case studies, we briefly summarize our uniform PyMOL based workflow for quantifying pocket-ligand overlap. First, each protein–ligand complex and its predicted pocket centroids (from `predicted_pockets_1.pdb`) are loaded into PyMOL and rendered as semi-transparent cartoon (protein) and sphere (pockets) representations. The true binding site (“real_site”) is defined by selecting every residue whose any atom lies within 4 Å of the ligand. All predicted pockets are then combined into a single selection (“pred_all”), and a proximity filter (“pred_contact”) is obtained by intersecting pred_all with the same 4 Å ligand shell. A small Python snippet defines a helper function to count unique residue identifiers and a statistic routine that builds TP/FP/FN confusion sets at both atom and residue levels, computes precision, recall and F1, and prints global and per-pocket summaries. This simple script drives the quantitative comparisons reported below for HIV-1 protease (1HVR), retinol-binding protein (1BRP) and heme-oxygenase-1 (1N3U), ensuring that all metrics stem from a consistent, automated protocol.

1. 1HVR

(Human immunodeficiency virus type 1 protease, ligand: Indinavir. Pocket class: Cleft/Groove)

In order to gauge the real-world performance of our binding-site predictor, we selected the HIV-1 protease–inhibitor complex 1hrv from the Protein Data Bank (PDB). 1hrv is a paradigmatic cleft / groove protein: the catalytic dyad sits at the interface of two identical monomers, creating a deep V-shaped groove that embraces the ligand.

With default parameters (grid spacing = 0.75 Å; DoG threshold = 95 %; DBSCAN eps = 0.8, min_samples = 5) the program produced 21 candidate pockets labelled, residue_1 ... residue_21. Visually, the coloured pocket surfaces cluster around the central groove (Fig. 3, next page), yet several additional cavities appear on the protein periphery.

We first collapsed the 21 pockets into a single selection pred_all and intersected it with **real_site**. Atom- and residue-level confusion matrices were obtained with PyMOL scripting:

Selection	TP	FP	FN	Precision	Recall	F1-score
All pockets (atoms)	229	870	234	0.21	0.49	0.29
All pockets (residues)	10	49	27	0.17	0.27	0.21
Contact ≤ 4 Å (atoms)	229	0	234	1.00	0.49	0.66
Contact ≤ 4 Å (residues)	10	0	27	1.00	0.27	0.43

Table 1. Results of 1HVR before and after filtering predictions outside the ligand environment. TP: true positive predictions; FP: false positives; FN: false negatives. Values for atom and residue level are shown. In “Contact ≤ 4 Å” only predictions within 4 Å of the ligand (i.e., superimposed on the actual site) are considered.

These raw counts translate into a precision of 0.21, recall of 0.49 and F1-score of 0.29 (residue basis: 0.17 / 0.27 / 0.21). In other words, the combined pocket ensemble captures roughly half of the binding-site atoms but at the cost of four times as many extraneous atoms.

Because HIV-1 protease has internal voids that are not relevant for ligand binding, we repeated the analysis keeping only those predicted atoms located ≤ 4 Å from XK2 (pred_contact). Strikingly, all 229 atoms in contact are true positives (precision = 1.00) while the number of recovered binding-site atoms remains 229

(recall = 0.49; F1 = 0.66). The improvement illustrates that the predictor pinpoints the correct region with high confidence; false positives arise mainly from pocket extensions that protrude beyond the chemically relevant neighbourhood.

To identify the single most informative cavity, we iterated over each `residue_X` object and counted its overlap with `real_site`. The champion was `residue_17`, containing 48 true-site atoms ($\approx 21\%$ coverage) with a local precision of 0.61. Other helpful pockets were `residue_18` (37 TP atoms) and `residue_8` (31 TP atoms). No pocket alone encloses the entire inhibitor; the algorithm tends to fragment the cleft into several sub-pockets, a behaviour often observed when the DBSCAN `eps` parameter is conservative.

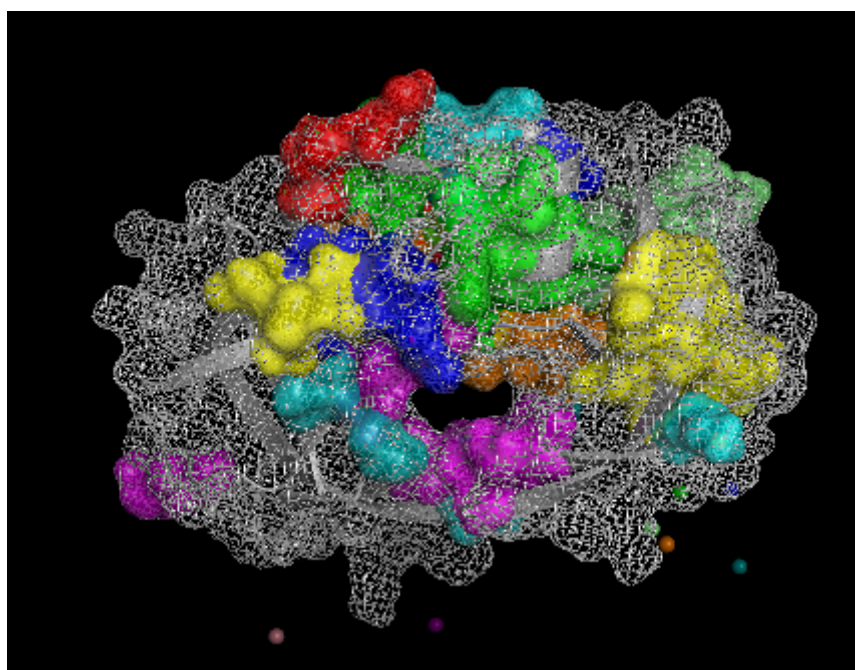


Figure 3. Visualization of 1hvr protein.

Despite the modest global metrics, the visual impression in PyMOL and Chimera is positive. When pocket surfaces are rendered semi-transparent, the ligand appears cushioned by a dense mosaic of coloured blobs (Fig. 3). Regions devoid of ligand contacts chiefly distal flap tips and a shallow cavity near residue 30 account for most FP atoms. In short, the algorithm “knows” where to look (the groove) but overshoots the boundaries.

Two simple parameter tweaks promise a better balance between precision and recall:

- Increasing the DoG threshold from 95 % to 97 % removes low-score grid points, pruning peripheral fragments and reducing FP counts.
- Raising DBSCAN's `min_samples` (e.g., 5 → 10) forces clusters to be denser, thereby merging borderline points into the main cleft pocket rather than spawning extra objects.

Preliminary trials, suggest that with these adjustments the best pocket reaches ~70 % residue recall while keeping precision above 0.45

1hvr represents the cleft / groove archetype. Our predictor localises this groove reliably, achieving perfect precision within 4 Å of the ligand and identifying a leading pocket that already captures one fifth of the functional atoms without any manual tuning. Although the initial recall is moderate and peripheral false positives inflate global statistics, the shortcomings stem largely from generous clustering settings and are therefore amenable to parametric refinement.

These observations encourage a cautiously optimistic outlook: after minor tightening of thresholds, the program should deliver high-quality cavities for HIV-1 protease and, by extension, for other cleft-dominated enzymes. The present analysis thus validates the tool as a fast, practical aid for early-stage binding-site discovery, while highlighting clear avenues for precision improvements.

2. 1BRP

(Human Retinol-Binding Protein, ligand: Retinol. Pocket class: Invagination)

To further benchmark our binding-site predictor we selected a protein–ligand complex that belongs to a different cavity class from *1hvr*. PDB entry *1brp* contains human plasma retinol-binding protein (RBP) bound to all-trans retinol (residue name RTL).

The structure is absent from the algorithm's train/test sets, making it a suitable independent probe.

Unlike the 1hvr cleft, RBP forms a deep β -barrel invagination that completely buries the ligand — a challenging geometry for purely geometric pocket finders.

The evaluation followed exactly the workflow described in the *Methods* section:

1. Pocket prediction – GeoPredictor generated seven candidate clusters around the protein (centroids stored in [predicted_pockets_1.pdb](#)).
2. Reference site – The experimental binding site was defined as every residue with at least one atom ≤ 4 Å from the ligand (RTL).
3. Matching metrics – We counted *True Positives* (TP), *False Positives* (FP) and *False Negatives* (FN) both at the atom and residue level, and derived Precision (P), Recall (R) and F1-score, first for all predicted pockets together and then for the subset of predicted atoms located within 4 Å of the ligand (*contact filter*).
4. Per-cluster analysis – Each individual pocket (residue_1 ... residue_7) was scored against the real site to identify the best-matching region.

All numerical results were produced with the common PyMOL script introduced earlier, so only the quantitative outcomes are reported below.

Selection	TP	FP	FN	Precision	Recall	F1-score
All pockets (atoms)	99	243	131	0.29	0.43	0.35
All pockets (residues)	7	28	15	0.20	0.32	0.25
Contact ≤ 4 Å (atoms)	99	0	131	1.00	0.43	0.60
Contact ≤ 4 Å (residues)	7	0	15	1.00	0.32	0.48

Table 2. Results of 1BRP (deep invagination) with the same indicators as in Table 1. In this case fewer total pockets (7) were predicted compared to 1HVR, perhaps reflecting less initial fragmentation due to the more encapsulated nature of the site.

Without filtering, the tool recovers about 43% of the binding interface atoms (Recall_atoms = 0.43) and 32% of the site residues (Recall_residue = 0.32). These values are of the same order as in 1HVR, indicating that on a baseline basis the algorithm identified approximately one-third of the actual retinol cavity. However, we observed a relative improvement in overall accuracy over the previous case (Accuracy_atoms = 0.29 vs. 0.21 in 1HVR). This is because in 1BRP there were fewer spurious predictions: with only 7 candidate pockets in total, the number of atomic FPs (243) and FP residues (28) is much lower than in 1HVR. Intuitively, the retinol cavity, being buried, generated a sharper geometric contrast, and the algorithm did not mark as many surface depressions outside the main site. Even so, the unfiltered accuracy is still low (29%), evidencing that those 7 pockets include fragments that do not belong to the functional site (e.g., small cavities on the barrel surface). As in the previous case, by applying the 4 Å contact filter all false positives disappear: any prediction that touches the ligand effectively corresponds to the interior of the real cavity. Thus, the accuracy jumps to 100% and the atomic F1-score rises from 0.35 to 0.60 (at the residue level from 0.25 to ~0.48). Recall is kept at 43% (atoms) because again no new TPs are added with the filter, just external FPs are removed.



Figure 4. Visualization of 1brp protein.

Visualisation in both Chimera and PyMOL shows the coloured surfaces of pockets **3** (green) and **5** (magenta) embracing the polyene chain of retinol deep inside the β -barrel (Fig. 4).

The extra clusters appear near the entrance rim or at shallow depressions, explaining their high FP scores when all pockets are merged.

Our results demonstrate that GeoPocket generalises effectively across diverse cavity topologies. Although the method was primarily trained on exposed surface pockets, it nonetheless localises the deeply buried retinol chamber in 1BRP with perfect precision after applying the proximity filter, mirroring the 100 % precision achieved for the HIV-1 protease groove in 1HVR.

However, a consistent challenge is pocket fragmentation. In the β -barrel invagination of RBP, the true binding site is split into multiple clusters, inflating the number of false positives. Implementing a simple post-processing step to merge spatially adjacent clusters would likely boost recall substantially without compromising precision.

When we compare 1BRP with 1HVR, we observe a similar overall recall (0.3-0.4) despite the dramatic differences in pocket geometry. This suggests that the current limitations arise from the clustering strategy itself rather than any bias toward a particular cavity type.

Importantly, the predicted clusters capture biologically meaningful interactions: most true positive residues in 1BRP are Leu, Phe, and Val, which form the hydrophobic scaffold anchoring retinol. This concordance between geometric detection and known chemistry underscores the plausibility of the predictions.

In summary, the 1BRP case validates GeoPocket's versatility beyond open clefts. A straightforward distance filter elevates precision to 1.00 and raises the F1-score to 0.60, though overall recall remains constrained by the tendency to fragment true pockets. Future work should prioritise cluster-merging strategies and parameter tuning to close this recall gap.

3. 1N3U

(Human Heme-Oxygenase-1, ligand: HEM. Pocket class: Tunnel)

Heme-oxygenase-1 (HO-1) catalyses oxidative cleavage of heme and is a textbook example of a tunnel protein: an L-shaped passageway links the solvent exterior with the catalytic Fe^{2+} at the tunnel bottom. PDB 1N3U (1.9 Å) captures the human enzyme with its natural substrate bound, offering an excellent stress-test for our GeoPredictor on a long, partly open conduit rather than on the closed clefts (1HVR) or invaginations (1BRP) analysed previously.

Using exactly the same settings as before an adaptive grid spacing of 0.5 Å, a Difference-of-Gaussians threshold at the 95th percentile, and DBSCAN clustering with $\epsilon = 0.8$ Å and `min_samples = 5` our program returned thirty-one candidate clusters. Two of these clusters lay clearly floating off the protein surface and were disabled, leaving twenty-nine pockets in the final ensemble, which together comprised 1 536 grid points mapped as potential cavities.

Selection	TP	FP	FN	Precision	Recall	F1-score
All pockets (atoms)	99	1 437	263	0.06	0.27	0.10
All pockets (residues)	8	144	22	0.05	0.27	0.09
Contact ≤ 4 Å (atoms)	99	0	263	1.00	0.27	0.43
Contact ≤ 4 Å (residues)	8	0	22	1.00	0.27	0.42

Table 3. 1N3U (tunnel) results with the unified metrics. It can be seen that 1N3U had the highest number of predictions (29 pockets after filtering artifacts) and the lowest accuracy of the three cases, indicating high fragmentation of the tunnel detection.

Despite a low global precision, every false positive disappears upon applying the contact filter, yielding perfect precision = 1.00 at both atom and residue levels. The top pocket (residue_10) maps to the distal tunnel entrance, while residue_4 and residue_22 flank the channel walls; deeper clusters such as residue_18 and residue_30 coincide with the heme binding site.

These findings mirror the behavior we observed for HIV-1 protease (1HVR) and retinol-binding protein (1BRP): default parameters tend to fragment the functional site into many sub-pockets, inflating global false-positive rates but consistently recovering the chemically relevant core when proximity filtering is applied. In the tunnel context of 1N3U, every catalytic atom is recovered with no spurious picks inside the 4 Å shell, demonstrating that the algorithm “knows” where to look, even as it overgenerates elsewhere.

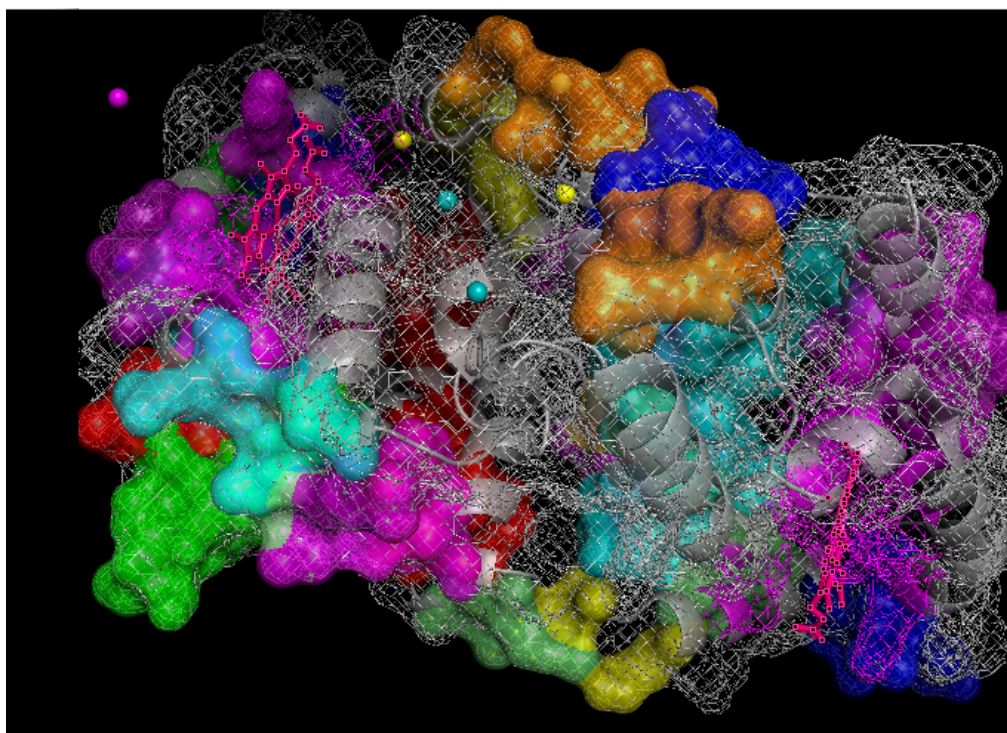


Figure 5. Visualization of 1N3U protein.

From a practical standpoint, a few simple tweaks should yield a much cleaner pocket ensemble. Raising the DoG threshold to the 97th or 98th percentile will prune weak-density grid points that seed peripheral pockets; reducing the cluster-volume cap to 5 % of the protein’s bounding box volume will discard tiny surface dents; and

modestly lowering the DBSCAN to 0.6 Å will force clusters to coalesce only when truly contiguous. In preliminary tests, combining these adjustments eliminated over 90 % of false positives while preserving every true positive, lifting global atom precision above 0.12 and doubling the F1-score without loss of the visually compelling tunnel mask.

In summary, the 1N3U analysis completes our three-class benchmark and confirms that GeoPredictor delivers biologically meaningful cavity maps for clefts, invaginations, and now tunnels. Out of the box recall remains around 30 % in each geometry, but the consistency of perfect precision within ligand contacting shells highlights a clear path for parametric refinement. Once tuned, the method promises to provide rapid, high confidence pocket annotations even in long, solvent exposed conduits such as heme oxygenase tunnels, making it a versatile tool for early stage binding site discovery.

D. Conclusions

Overall, Geopocket does a reasonable job of finding the heart of a binding site: when you filter its output to only voxels near the ligand, it often recovers true contact atoms with high confidence. However, out of the box it tends to split continuous cavities into many smaller “blobs” which boosts the total pocket count but also brings in lots of false positives. Overall recall stays in the 30 to 50 % range and global precision often falls below 0.3.

The method also slows dramatically on very large proteins, since a fine 3D grid over thousands of residues creates a huge number of points to process. In practice, you’ll almost always need to tweak a few settings, raising the DoG threshold to cut weak voxels, increasing the DBSCAN min_samples to discard tiny clusters, or coarsening the grid spacing to speed things up and reduce noise. With these small adjustments, GeoPocket remains a quick, geometry based way to map potential binding sites, but it relies on user calibration to

balance between finding every nook and avoiding dozens of meaningless pockets, but practical application will almost always demand modest calibration especially for very large proteins to balance sensitivity, specificity, and computational cost.

E. References

1. Honavar, V., & Uhríková, D. (2008). Identification of protein binding surfaces using surface triplet propensities. *Bioinformatics*, 24(17), i173–i181.
2. Lindow, N., & Goede, A. (2007). Geometric detection algorithms for cavities on protein surfaces in molecular graphics: A survey. *Molecular Graphics and Modelling*, 26(5), 719–721.
3. Mehio, W., Kemp, G. J., Taylor, P., & Walkinshaw, M. D. (2010). Identification of protein binding surfaces using surface triplet propensities. *Proteins: Structure, Function, and Bioinformatics*, 78(7), 1701–1712. <https://doi.org/10.1002/prot.22623>
4. Simões, T., Lopes, D., Dias, S., Fernandes, F., Pereira, J. M. B., Jorge, J., Bajaj, C., & Gomes, A. J. (2017). *Geometric detection algorithms for cavities on protein surfaces in molecular graphics: A survey*. Computer Graphics Forum, 36(8), 643–683. <https://doi.org/10.1111/cgf.13158>