



steps with respect to each other,  $n_{\max}$  – the maximum circle rank (the resolution). The symbol  $\oplus$  designates the sum modulo 7.

The resemblance varies when the two neighborhoods are rotated around the common normal. For each pair of surface atoms  $A$  and  $B$ , the maximum and the minimum resemblance are defined by:

$$R_{\max}(A, B) = \max_{h \in \{0, \dots, a-1\}} R(A, B, h) \quad (5)$$

$$R_{\min}(A, B) = \min_{h \in \{0, \dots, a-1\}} R(A, B, h) \quad .$$

Because two interacting molecules can usually move with respect to each other, these magnitudes are more significant when comparing or when estimating the interaction of two surface atom neighborhoods.

The *similitude* and the *interaction* of a pair of atom neighborhoods are defined as their *resemblance* for parallel and, respectively, anti-parallel orientations of the vectors normal on the molecular surfaces in the superposed contact points  $A=B$ , as shown in Fig. 2.

For illustration, Fig. 3 shows the histogram of the maximum interactions for all pairs of surface atoms of the proteins with the labels 135L and 1HZH in PDB. All mutual orientations of the two surfaces are considered for each pair of surface atoms and the largest value of the interaction is retained in each case.

Similarly, Fig. 4 shows the histogram of the minimum interactions for all pairs of surface atoms of the same proteins 135L and 1HZH. The regularity of these histograms indicates a limited variety of the surface atom neighborhoods.

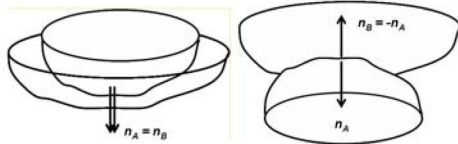


Figure 2. Parallel and anti-parallel orientation of vectors normal on the molecular surfaces in the superposed contact points  $A=B$ .

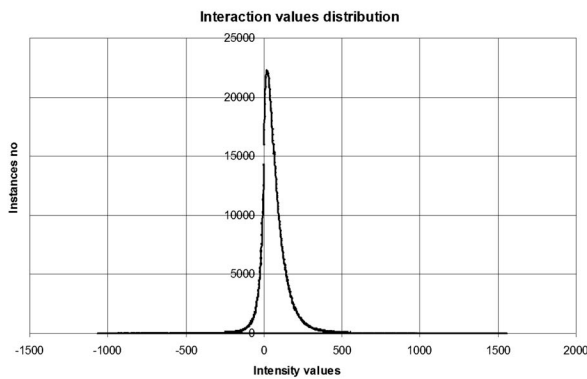


Figure 3. Maximum interaction histograms for all pairs of surface atoms of the proteins 135L and 1HZH.

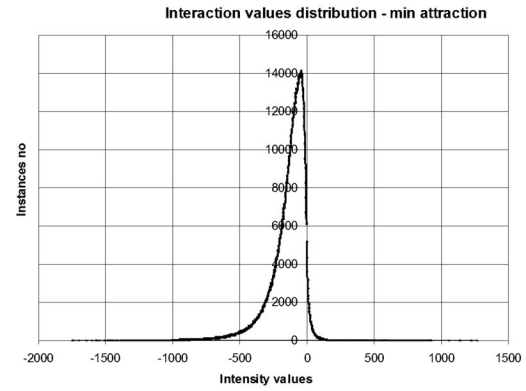


Figure 4. Minimum interaction histograms for all pairs of surface atoms of the proteins 135L and 1HZH.

#### IV. VECTOR CLASSIFICATION OF SURFACE ATOM NEIGHBORHOODS

A functional image oriented representation of protein surfaces and the corresponding software tool have been developed to interactively explore the hydrophobicity distribution on the molecular surface. The surface atom neighborhoods have been described by vectors of 65 components which specify the hydrophobicity densities (2) in each patch of the standardized octagonal frame defined in Section II. The comparison of the neighborhoods has been performed both globally, in terms of *similitude* and *interaction*, and by their clustering using the vector description.

Because the relative angular position of two surface atom neighborhoods can vary arbitrarily, it is necessary to pre-process the vectors describing the neighborhoods, to bring them in similar positions before applying any classification algorithm. The calibration pattern shown in Fig. 5 has been used for this purpose. The pattern has patches of clockwise decreasing hydrophobicities, starting from the highest hydrophobicity density (3.89) in the sector 1, passing through zero hydrophobicity density in sector 6, and reaching the lowest value (-1.007) in sector 8 (7 and 8 are thus hydrophilic). Each surface atom neighborhood is rotated to a position in which its resemblance to the calibration pattern is maximized.

The WEKA (Waikato Environment for Knowledge Analysis) API was used as clustering engine. Fig. 6. shows the results of the clustering in terms of the *Average Hydrophobicity* vs. *Average Resemblance* dependence.



Figure 5. Surface atom neighborhoods calibration pattern with clockwise decreasing hydrophobicities.

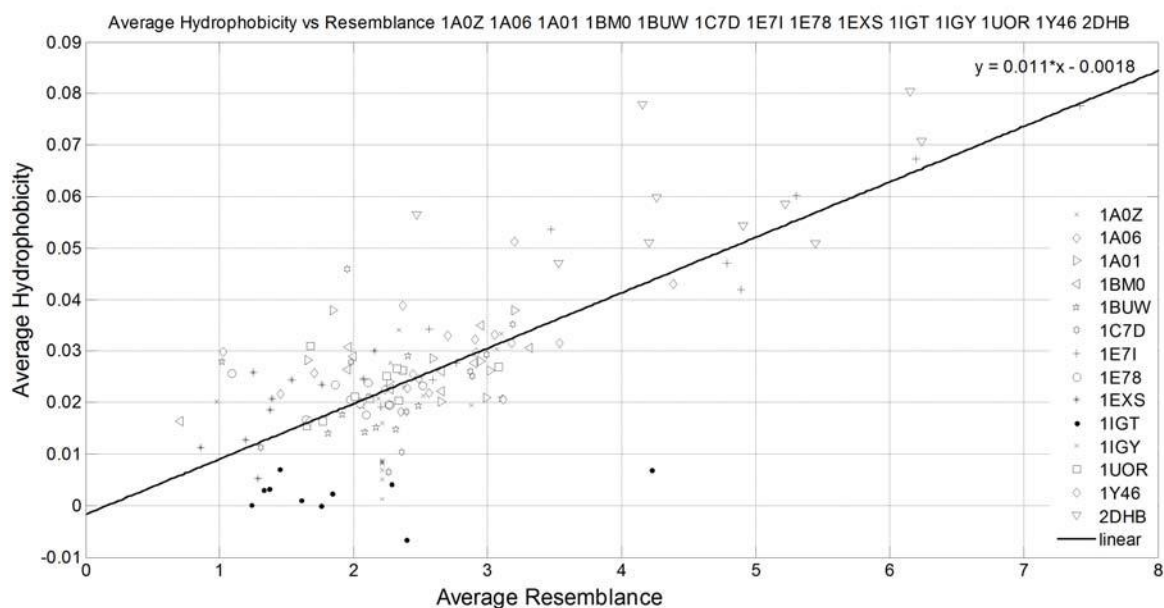


Figure 6. Clustering of surface atom neighborhoods for the 1A0Z, 1A06, 1A01, 1BM0, 1BUW, 1C7D, 1E71, 1E78, 1EXS, 1IGT, 1IGY, 1UOR, 1Y46, 2DHB proteins [], expressed by the dependence of the *Average Hydrophobicity* on the *Average Resemblance* of clusters to a reference pattern, after re-orienting individual surface atom neighborhoods to maximize their resemblance to the reference pattern.

The *Average Hydrophobicity* of the surface atom neighborhood clusters generated by WEKA is given as a function of their *Average* (maximum) *Resemblance* to the reference pattern in Fig. 5, after the re-orientation of the individual pattern in each cluster. Data for the larger molecules (1A0Z, 1A06, 1A01, 1BM0, 1BUW, 1C7D, 1E71, 1E78, 1EXS, 1IGT, 1IGY, 1UOR, 1Y46, 2DHB) in the set of 36 studied proteins have been used in Fig. 6. An approximately linear dependence has been found for the 140 clusters, with a larger spreading of data for the small values of average resemblance. Notice that the resemblance to the reference pattern is maximized for each individual neighborhood in a cluster, whereas the average resemblance of for all the patterns is given for every cluster.

Results of clustering at the vector level are given in Figs. 7-10, for the 1E71, 1Y46, 1UOR and 2DHB protein surface atom neighborhoods.

Each neighborhood is described by a hydrophobicity density vector with 65 components, corresponding to the central circle and to the 8x8 annular sectors. The goal of the analysis is to find patterns on the protein surface neighborhoods in terms of hydrophobicity. Because there are no priorly defined classes of hydrophobicity densities, an unsupervised learning scenario has to be used. Clustering is used to group items that seem to fall naturally together. The output takes the form of a list specifying how the instances fall into clusters. The success of clustering is often measured subjectively in terms of how useful the result appears to a human user.

The basic clustering technique is the *k*-means. The user specifies in advance how many clusters are being sought, the *k* parameter. Then, *k* points are chosen at random as cluster centers. All instances are assigned to their closest cluster center according to the ordinary Euclidean distance metrics.

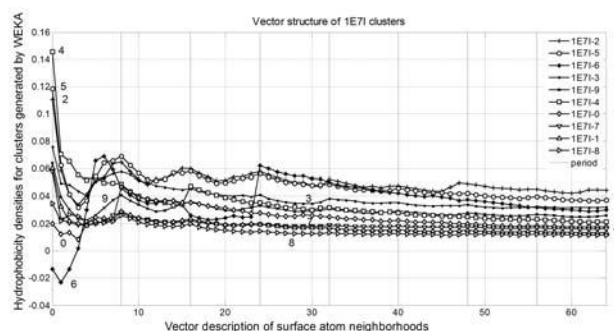


Figure 7. Vector clustering of 1E71 protein surface atom neighborhoods.

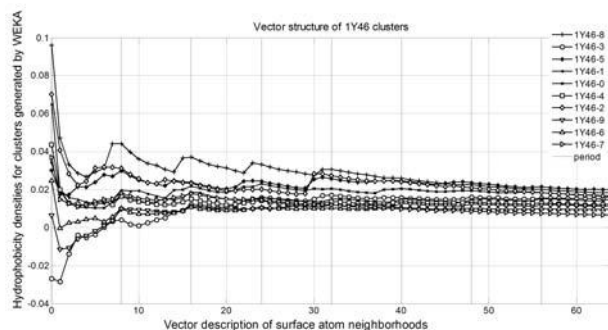


Figure 8. Vector clustering of 1Y46 protein surface atom neighborhoods.

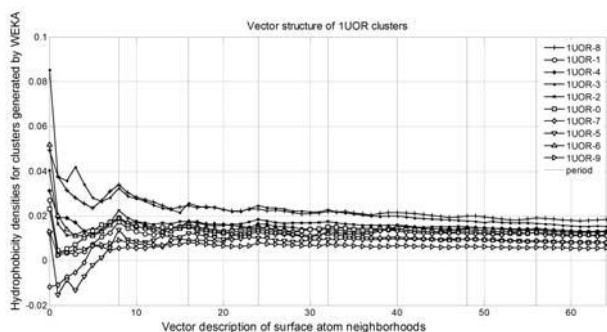


Figure 9. Vector clustering of 1UOR protein surface atom neighborhoods.

We have used the expectation-maximization algorithm for clustering the hydrophobicity densities of the atoms vicinities, one of the most powerful methods for finding a maximum likelihood solution for models with hidden variables [11,12].

The WEKA EM implementation requires the input parameters:

- Number of clusters to generate, chosen 10 based on the results obtained from tests with 5 to 25 clusters;
- Maximum number of iterations, chosen to 100 iterations;

Maximum allowable standard deviation for the density calculation, chosen  $10^{-6}$ .

The figures give the average hydrophobicity densities for the 65 components describing each of the ten clusters in which are classified similar patterns. Prototype images have been constructed for each of the clusters by using their average hydrophobicity vectors. These images, not shown here to avoid the need of colors, represent a set of “physiognomies” corresponding to the pattern at the surface of each studied protein.

## V. CONCLUSIONS

We continued the study of protein surfaces has with the classification of local molecule properties using resemblance and vector descriptions [10]. The global values of these magnitudes, the histograms of their distribution for all the surface atoms, and the actual structure of the surface atom neighborhoods are taken into account.

The classification allows predicting the behavior of protein molecules when interacting with each other, or with a nanostructured surface. Further work will include not only the study of the interactions determined by the hydrophobicities, de-convoluted at the level of atoms [6], but also the effects of the electrical interactions and the way these interactions are influenced by the pH. The cumulated effect of the two types of interactions will be expressed in a coherent way, allowing to compute the resulting similitude and interaction of two surface atom neighborhoods. This approach will result in a better description of the complex phenomena involved in the protein multi-parameter interactions. A web approach is also considered, to facilitate the access of academia and industry to the new results.

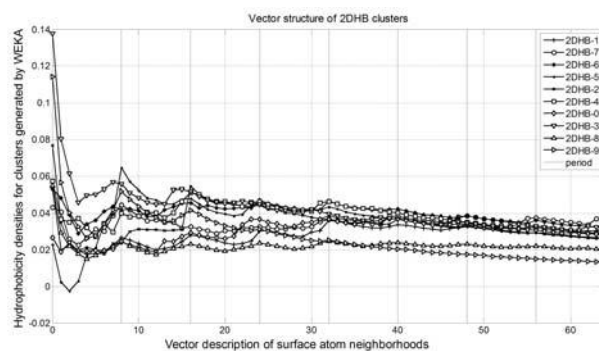


Figure 10. Vector clustering of 2DHB protein surface atom neighborhoods.

## ACKNOWLEDGMENT

The work was partially supported by the project 214538 – 2008 - BISNES – “Bio-Inspired Self-assembled Nano-Enabled Surfaces,” in the framework of the NMP-2007-1.1-2 Self-assembling and self-organization, NMP-2007-1.1-1 Nano-scale mechanisms of bio/non-bio interactions.

## REFERENCES

- [1] Protein Data Bank [Online], <http://www.rcsb.org/pdb/home/home.doc>.
- [2] Project 214538 – 2008 - BISNES – “Bio-Inspired Self-assembled Nano-Enabled Surfaces,” <http://www.bisnes4eu.com/>
- [3] M. L. Connolly, “MS: Molecular Surface Program,” *QCPE Program 429*, Quantum Chemistry Program Exchange, Univ. of Indiana, Bloomington, 1983.
- [4] M. L. Connolly, “Molecular Surfaces: A Review,” *Network Science (Online)*, vol.2(4), 1996, <http://www.awod.com/netsci/Science/Compchem/feature14.html>.
- [5] P. A. Karplus, “Hydrophobicity regained,” *Protein Science*, vol. 6, pp. 1302-1307, 1997.
- [6] M. Held and D. V. Nicolau, (2007), “Estimation of atomic hydrophobicities using molecular dynamics simulation of peptides,” *Proc. of SPIE*, 6799, Modelling and THZ Technology, 6799-16, 1-7.
- [7] D. V. Nicolau, F. Fulga, and D. V. Nicolau, (2003), “A new program to compute the surface properties of biomolecules,” *Asia-Pacific Biotech*, 7(3), 29-34.
- [8] P. D. Cristea, Rodica Tuduce, O. Arsene, D. V. Nicolau, F. Fulga, “Multi-threading Protein Surface Functional Description,” *NEUREL 2010*, Belgrade, Serbia, September 23-25, 2010, *Proc. of NEUREL 2010*, pp. 1075-1078.
- [9] P. D. Cristea, Rodica Tuduce, O. Arsene, Alina Dinca, D. V. Nicolau and F. Fulga, “Modeling of Biological Nanostructured Surfaces,” *Proc. of SPIE*, 7574, Nanoscale Imaging, Sensing, and Actuation for Biomedical Applications VII, 2010.
- [10] P. D. Cristea, Rodica Tuduce, O. Arsene, D. V. Nicolau, “Functional Nanoscale Imaging of Protein Surfaces,” *BiOS SPIE Photonics West – Nanoscale Imaging, Sensing, and Actuation for Biomedical Applications Conference*, San Francisco, California, USA, 2011.
- [11] Russell, S., Norvig, P., “Artificial Intelligence. A Modern Approach,” Prentice Hall, 3rd edition, 2010.
- [12] Bishop, C., “Pattern Recognition and Machine Learning,” Springer, 2006.