

A new method for protein characterization and classification using geometrical features for 3D face analysis: an example of tubulin structures

Luca Di Grazia ¹, Maral Aminpour ^{2,3} Enrico Vezzetti ¹, Vahid Rezaia ⁴, Federica Marcolin ¹, and Jack Adam Tuszynski ^{1,2,3*}

¹ Politecnico di Torino, corso Duca degli Abruzzi 24, 10129 Torino, Italy;

² Department of Physics, University of Alberta, Edmonton, Alberta, Canada

³ Department of Oncology, University of Alberta, Edmonton, Canada

⁴ Department of Physical Sciences, MacEwan University, Edmonton, Alberta, Canada

* Correspondence: jacek.tuszynski@polito.it;

Short title/running title:

Protein characterization and classification using geometrical features for face analysis

Abstract: This paper reports on the results of research aimed to translate biometric 3D face recognition concepts and algorithms into the field of protein biophysics in order to precisely and rapidly classify morphological features of protein surfaces. Both human faces and protein surfaces are free-forms and some descriptors used in differential geometry can be used to describe them applying the principles of feature extraction developed for computer vision and pattern recognition. The first part of this study focused on building the protein dataset using a simulation tool and performing feature extraction using novel geometrical descriptors. The second part tested the method on two examples, first involved a classification of tubulin isotypes and the second compared tubulin with the FtSZ protein, which is its bacterial analogue. An additional test involved several unrelated proteins. Different classification methodologies have been used: a classic approach with a Support Vector Machine (SVM) classifier and an unsupervised learning with a k-means approach. The best result was obtained with SVM and the radial basis function (RBF) kernel. The results are significant and competitive with the state-of-the-art protein classification methods. This opens a new area for protein structure analysis.

Keywords: 3D Face Analysis; Protein Classification; Tubulin; SVM; Geometrical Descriptors; Differential Geometry; Machine Learning

1. Introduction

The structure of a protein is an important indicator of its potential biological functions, especially its surface, which is exposed to the solvent and participates in interactions with other proteins and ligands. In a recently published work [1] it was shown how to capture fingerprints of a protein using deep learning methodology and a strong correlation was demonstrated between the structure of a protein and its biological behavior. Another work [2] showed the relevant role of protein-protein interactions using local structural features. In this latter paper geometrical features were found to be interesting in this context.

The first step in the process of classifying proteins is to acquire a realistic (usually experimental) 3D dataset regarding a protein's structure. X-ray crystallography has made the largest and most important contribution to our understanding of protein structure. Nuclear Magnetic Resonance (NMR) and cryogenic electron microscopy (cryo-EM) are other methods by which to determine the protein structure [3] but they have various limitations. As an alternative to crystallographic structure determination, a computational method can be used to generate its prediction using a three-dimensional model [4]. However, proteins are non-static molecular structures, thus a crystallography-generated image is only a snapshot in time of a protein structure and not a fully realistic representation of all protein states, which can be quite dynamic. Therefore, molecular dynamics (MD) is a useful computational tool that can be used to produce atomic coordinate trajectories in order to provide a sampling of structural representations of a given protein. The method we propose in this paper is agnostic to the origin of the data, which in the case of proteins can either be obtained from experiments such as cryo-EM or

91 synthetically generated from computational approaches such as MD. The key aspect is to
92 have an atomistic model of the objects studied [3], which serves as the starting point for
93 feature extraction based on the protein surface. Such a model provides a high-resolution
94 representation of the object of interest, which is later on processed and characterized by a
95 manageable number of parameters.

96 A protein can have different equilibrium conformational states that depend on
97 ambient conditions. Moreover, some proteins are expressed by several genes leading to
98 different isotypes with a high degree of structural similarity making accurate comparison
99 important, so a good dataset with different frames is important in order to have a
100 statistically significant and valid test set. The most difficult task would be to distinguish
101 between very closely related proteins or indeed the same protein in its wild type form and
102 a mutated protein structure. For clearly distinct protein structures, standard approaches
103 for their comparisons such as the use of the RMSD (root mean squared deviation) may
104 work reasonably well but providing a single parameter only for structure comparisons
105 may not always be useful or sensitive enough to distinguish subtle structural changes
106 involving, for example, single point mutations or a small number of amino acid
107 substitutions. It should also be mentioned that while sequence comparison methods are
108 rapid and reliable, since there is no general solution to the protein folding problem,
109 sequence comparisons are insufficient by themselves to inform us about subtle structural
110 changes that can distinguish between highly similar protein structures.

111 Some experimentation has already been undertaken to classify proteins according to
112 their states. Tsuda et al. adopted a Support Vector Machine (SVM) classifier for fast
113 protein classification [5]. They obtained 13 classes and reached an accuracy of about

90%. Weston et al. [6] used a semi-supervised classification with a kernel cluster and reached a result of 94.3%. Another interesting result has been obtained using a random forest approach and fifteen different supervised methods with about 11,000 pairs of protein domains leading to an accuracy of 97.0% [7]. Our focus in this paper is on accurate differentiation between structurally-similar proteins, which is a much harder problem to solve than comparing vastly different protein structures. Many cases of protein families can be found and it is important to be able to find characteristic features distinguishing proteins belonging to the same family between each other. This could be valuable with respect to their functional roles in cell biology as well as potential applications in rational drug design. To the best of our knowledge, no methodology exists in the literature that deals with this particular situation.

One of the most important proteins abundantly expressed in all eukaryotic cells is the family of tubulin proteins, which will be studied in this paper as a challenging test case for this methodology. It is also highly homologous with its bacterial ancestor, FtSZ, which will also be used here for comparison. We should stress again that comparing protein sequences is a trivial problem in bioinformatics while 3D structural features of folded proteins pose a much greater challenge, which is addressed here. These structures are obtained from various experiments such as X-ray, NMR or cryo-electron crystallography or from computational simulations such as MD, as mentioned above.

In the computational experiment reported below SVM was used because the quantity of data tested was relatively low, and a deep learning approach requires large data sets to achieve a high level of confidence. The novelty of our approach rests with the feature extraction using geometrical descriptors and its general applicability to 3D structure

137 characterization, because geometric feature surfaces were used with significant results in
138 many other applications before, e.g. [8, 9]. We believe that the classification provided
139 here can be further improved with more data, more classes and a complex neural
140 network, all of which is planned for future work, especially within the context of
141 geometric deep learning [10], which nowadays is the state-of-the art of classification.

142 Tubulin is a key cytoskeletal protein, which has been exhaustively studied for its
143 applications in several fields including being the target for various anti-cancer drugs [11]
144 and the discrimination of the *Saccharomyces* complex [12]. It is a globular protein with a
145 molecular weight of 55 kDa per monomer and its numerous isotypes expressed by
146 separate genes have a broad distribution in animal and plant cells [13]. Tubulin is a
147 building block of microtubules (MTs) and its stable form is an $\alpha\beta$ -heterodimer. MTs play
148 various important roles in all eukaryotic cells including cell motility, material transport
149 and most importantly cell division where MTs form mitotic spindles [12-13].

150 The novelty of the present work rests with the application of geometrical descriptors
151 coming from the field of face analysis to the classification of surfaces of proteins, with
152 the aim of adopting this geometrical information as descriptive features and
153 discriminating elements to classify proteins. Here, we test the method on the examples of
154 tubulin isotypes and related proteins (e.g. FtsZ). The method can, of course, be applied to
155 an arbitrary protein or indeed a protein complex but being able to discriminate between
156 highly homologous proteins based on the geometrical shapes of their surfaces opens the
157 door to numerous applications across the field of protein science. The idea comes from
158 the realization that geometrical properties can well describe the surface of a 3D object
159 such as a protein and could identify characteristic features when comparing two or more

160 similar structures. Proteins surfaces can be split into two outer surfaces by cutting a plane
161 through the data set including the main axis of rotational symmetry. These two halves of
162 the outer surface, similarly to human faces, differ from one another depending on the
163 protein type, and also can change their conformational states dynamically, similarly to
164 human facial expressions. Thus, what in the field of pattern recognition is called face
165 recognition could be transferred to the context of protein classification according to the
166 typology. These common points have fostered the interest of uncovering the potentiality
167 of cross-fertilization between these two fields with the aim of better categorization.

168 All eukaryotic organisms carry multiple genes coding for α and β tubulin (and other
169 variants, e.g. γ), which are referred to as isoforms when comparing tubulin expressed by
170 different organisms. When a single organism is discussed, various tubulin genes code for
171 what are called tubulin isotypes. Isotypes have highly homologous amino acid sequences
172 that appear to have diverged as a result of accumulated mutations since their separation
173 by distinct speciation events [14]. Amino acid sequence similarity is very high for all
174 tubulin proteins both within and between diverse species making structural comparisons
175 difficult. At the cellular level, the roles of the α and β tubulin isotypes are essential, a
176 result of subtle structural variations within their sequences [15] Several isotypes of the α
177 and β tubulins have been identified in human cells, their existence and distribution
178 providing a link to their specific roles in the polymerization and stability of MTs, among
179 other roles [8] making structural differences correlate with functional roles in cells,
180 importantly including cancer cells. For example, β II tubulin has been a common target
181 for chemotherapy drug action and is involved in protein-protein interactions [2]. Hence
182 again, the structural differences between tubulin isotypes significantly assist in drug

design targeting specific isotypes such as β III, which is overexpressed in all cancer cells. Through a search of available protein sequence databases, a total of ten unique β tubulin isotypes can be found, all of which have highly similar amino acid sequences and are generally well conserved. Sequence alignment, similarity and identity values of the studied isotype proteins (see below for details) range between 78% and 98%, indicating a major level of similarity between these structures. The question that remains is how do these sequence variations translate into structural differences.

As stated above, MTs are dynamic cytoskeleton polymers present in all eukaryotic cells made up of the protein tubulin. FtsZ is a close structural homologue of tubulin within prokaryotic cells, and plays an important functional role during bacterial cell division. A close relationship between FtsZ and tubulin can be seen from their very similar protein structures (Figure 1a). Both α and β tubulin share an approximate 35% sequence identity with FtsZ [16]. Both FtsZ and tubulin can assemble to form straight filaments. This association is regulated by guanosine triphosphate (GTP), which is bound in the junction between adjacent monomers (Figure 1b). FtsZ forms long protofilaments consisting of a single string of FtsZ proteins in contrast to tubulin, which makes cylindrical MTs. Unlike tubulin, FtsZ does not appear to provide a structural role throughout the cell cycle, but instead just plays a structural role during bacterial cell division, when it forms a band, known as the Z-ring, around the inner cell wall at the location where the cell will divide.

Figure 1

The main goal of the research reported here has been to investigate the following issues:

- whether it is possible to rely on features coming from the field of pattern recognition and face analysis to geometrically describe (and classify) the geometrical properties of the protein surface;
- whether it is possible to recognize different isotypes of the same protein from a different set of molecular dynamics snapshots;
- whether it is possible distinguish between two highly structurally similar but not identical proteins such as tubulin and FtsZ, and whether it is possible to distinguish arbitrary proteins with no relation to each other.

It is worth stating in this context that in general the main goal of a classifier is to separate objects belonging to different classes using a number of possible linear separators as shown in the examples presented in Figure 2.

Figure 2

It is reasonable to expect that using one of these separators one can get a datum that is on the other side of the hyperplane, which would then be misclassified because the hyperplane is really near the ham data [17]. SVM is able to find a solution with a larger margin for the two-separator classifier as shown in Figure 2(a). This hyperplane works better than others as it is expected to reduce the number of misclassifications, because it is the one with the highest margins from the two sets of data.

The first part of this paper describes the development of the dataset using tubulin isotypes and FtsZ protein as test cases. Then, geometrical descriptors are computed on the 3D surface of these proteins. They are then converted into histograms and saved in a file. This file is the input of the classifiers. The code is provided in a github repository [18]. The entire process is summarized in Figure 3.

Figure 3

This paper is organized as follows. In Section 2 geometrical descriptors used for implementing the feature extraction are described. Section 3 is the core of the paper and it outlines feature extraction and classification methods with a detailed description of the strategies and techniques performed. Section 4 summarizes and discusses the results comparing them with the-state-of-art results. Finally, Section 5 summarizes the work and discusses future developments.

2. Material and methods

2.1 Geometrical descriptors

The surfaces representing both human faces and proteins are geometrically considered as a free form. Thus, features coming from the field of differential geometry can be applied in order to understand their local and global properties. Geometrical descriptors are widely used in the area of 3D face recognition with significant results reported elsewhere in the literature [19, 20]. They underline different characteristics of a free-form and are an important tool for feature extraction [21] within the context of face analysis [22]. In this work, for the first time we apply these descriptors to proteins and use them for structural classification purposes [19].

The geometrical descriptors used in this research are the following novel geometrical descriptors [22, 23]: mean curvature (H_{mean}), Gaussian curvature (g_{mean}), principal curvatures ($k_{1_{mean}}$ and $k_{2_{median}}$), the shapes type of a surface (S_{mean}), and the symmetry property (F_{den2}). Considering that these descriptors rely on the derivatives of the surface (h_x , h_y), they well describe the changes in surface curvature ($k_{1_{mean}}$, g_{mean} , $k_{2_{median}}$, H_{mean}), depressions and peaks (local minima and maxima) of the

251 surface $(k_{1\text{mean}}, \text{sing}, k_{2\text{median}}, g_{\text{mean}}, H_{\text{mean}})$, the shapes in terms of the types of
 252 surfaces (S_{mean}) , and the surface's symmetry property $(F_{\text{den2}},)$. These parameters are
 253 highly informative of the investigated surface's geometrical properties. Each descriptor
 254 can underline a specific characteristic of a certain surface. These descriptors are briefly
 255 described below in regard to their conceptual order. The first and second fundamental
 256 forms provide the first six descriptors of the set. They are used to measure distance on
 257 surfaces and are defined by the formula

$$258 \quad ds^2 = Edu^2 + 2Fdudv + Gdv^2 \quad (1)$$

259 where E, F, G, e, f and g are their coefficients given by:

260

$$261 \quad E = 1 + h_x^2, \quad (2)$$

$$262 \quad F = h_x h_y, \quad (3)$$

$$263 \quad G = 1 + h_y^2, \quad (4)$$

$$264 \quad e = \frac{h_{xx}}{\sqrt{1+h_x^2+h_y^2}}, \quad (5)$$

$$265 \quad f = \frac{h_{xy}}{\sqrt{1+h_x^2+h_y^2}}, \quad (6)$$

$$266 \quad g = \frac{h_{yy}}{\sqrt{1+h_x^2+h_y^2}}, \quad (7)$$

267 Curvatures are used to measure how a regular surface x bends in. If D is the
 268 differential and N is the normal plane to a surface, then the determinant of DN is the
 269 product of the principal curvatures, and the trace of DN is the negative of the sum of
 270 principal curvatures. At point P , the determinant is the Gaussian curvature K of x at P .
 271 The negative of half of the trace of DN is called the mean curvature H of x at P .

272 The principal curvatures k_1, k_2 are the roots of the quadratic equation given below:

273 $x^2 - 2Hx + K = 0$ (8)

274 Thus, we can choose k_1 and k_2 so that:

275 $k_1: H + \sqrt{H^2 - K}$ and $k_2: H - \sqrt{H^2 - K}$ (9)

276 where

277 $K = \frac{eg-f^2}{EG-F^2}$ (10)

$$H = \frac{eG - 2fF + gE}{2(EG - F^2)} \quad (11)$$

278 In terms of the principal curvatures, Gaussian (K) and mean curvatures (H) can be written
279 as

280 $K = k_1 k_2$, (12)

281 $H = \frac{k_1 + k_2}{2}$ (13)

282 where h is a differentiable function representing the three-dimensional surface.

283 The curvedness index S , which describes the shape of the surface, is defined as:

284

285 $S = -\frac{2}{\pi} \arctan \frac{k_1 + k_2}{k_1 - k_2}, S \in [-1, 1], k_1 \leq k_2$. (14)

286

287 Some descriptors highlight particular facial lines, such as F_{den2} , which shows visible
288 facial part contours. It can be computed using the formula:

$$\frac{F}{1 + (h_x)^2 + (h_y)^2}, \quad (15)$$

289 where:

290 • h is the differentiable function $z = h(x, y)$ representing the face/protein surface;

- h_x and h_y are the first derivatives of h with respect to x and y [24].

In a protein, F_{den2} can underline different trends of the free form analyzed. There is a representation focus in these parameters on the local maxima and local minima of the protein surface.

The surfaces of human faces are given by depth maps, which are manageable as matrixes ($X\ Y\ Z$). For each coordinate pair X, Y , there is a unique value of Z . Since proteins do not have a default form, their surfaces are split up in two parts divided into two opposite faces: surfaces with a positive Z -axis and those with a negative Z -axis in order to yield two shells that complete the protein surface.

The descriptors were mapped onto these maps as follows. Considering that the Z coordinate of the depth map is the one describing the “surface”, and is represented in these formulas as h , the derivatives h_x, h_y, \dots were evaluated so that other same-sized matrices representing the first derivative with respect to x , the first derivative with respect to y , etc., were generated and stored. The derivatives are computed with the Matlab function "gradient". Then, the implementation formulas for the descriptors were calculated on the matrices previously computed and new same-sized matrices were obtained representing every geometrical descriptor.

The descriptors used are mapped onto the surfaces as described in Section 3.3. These descriptors are calculated for all protein faces considered in the following. An example of F_{den2} applied to both a human face and a protein is shown in Figure 4a. The descriptor *sing* is built from the application of the sine standard function applied to the third coefficient of the second fundamental form (e) (see Figure 4b) [23]. Mean and median filters have been applied to the primary descriptors $S, k1, k2,$

g, and H. Mean and median values are computed in squared neighborhoods of side 5 around each point of the facial depth maps [23]. These descriptors are obtained as follows: S_{mean} , (see Figure 4c), k_{1mean} (see Figure 4d), $k_{2median}$ (see Figure 4e), g_{mean} (see Figure 4f) and H_{mean} (see Figure 4g).

Figure 4

The process we follow in this paper starts with the collection of protein data. In the present example we focus on tubulin whose bovine structure has been crystallized and can be found in the Protein Data Bank (PDB). However, its various isotypes have not been crystallized and hence these structures need to be generated by homology modeling using the bovine (not human) variant of this protein as a template. To obtain frames of the protein structure, it is necessary to run MD simulations for some time, typically 10-100 nanoseconds and take snapshots, approximately every nanosecond, at the very moment when the structure relaxes to an equilibrium conformation. Only the atoms comprising the protein are kept in the file used for these MD simulations with the ligand atoms removed in order to avoid false representations of the protein since ligands are not part of the protein and can form an occlusion during the process of protein recognition. The next step in this computational experiment is to analyze similar but not identical proteins and their states, for example tubulin isotypes with each other or a tubulin isotype and FtsZ and to compare the two for similarities and differences.

The result of these MD simulations is in each case a PDB-formatted file that is a 3D representation of a protein, which is converted into a MAT file using a MATLAB script. In the current work several software packages are used: Matlab 9.5 (R2018b) [25] for the feature extraction using geometrical descriptors, Anaconda 1.9.6 [26] with Python 3.7

[27] and the library sklearn 0.22 [28] for the implementation of classification methods and R-3.5.3 for the K-means algorithm [29].

2.2 Molecular dynamics simulations

The tubulin crystal structures available in the PDB are those for bovine protein. The bovine tubulin structure of tubulin (PDB ID: 1JFF) [30] was used as a template to construct the homology model for human $\alpha\beta$ tubulin isotypes (β I (UniProtKb: P07437), β IIa (UniProtKb: UniProtKb: Q13885), β IIb (UniProtKb: Q9BVA1), β III (UniProtKb: Q13509), β IVa (UniProtKb: P04350), $\alpha\beta$ IVb (UniProtKb: P68371), $\alpha\beta$ V (UniProtKb: Q9BUF5), $\alpha\beta$ VI (UniProtKb: Q9H4B7) and β VIII(UniProtKb: Q3ZCM7)) using the Molecular Operating Environment (MOE) software package [31]. Multiple sequence alignment results contained in Figure 5 show that human β -tubulin isotypes exhibit residue composition variations at different locations.

Figure 5

Sequence similarity matrix and sequence identity matrix of the tubulin isotypes are shown in Figure 6(a) and (b), respectively. The matrix values (i, j) for the percentage identity and similarity metrics are equal to the number of sequence matches between chains i and j, divided by the number of residues in chain i. Residues are considered identical if their single-letter code is the same (note that MSE-Selenomethionine and MET-Methionine are considered "identical"). Residues are "similar" if their BLOSUM62 substitution score is greater than zero.

Figure 6

The atomic coordinates of similar but not identical FtsZ dimer were obtained from the Protein Data Bank as (PDB ID: 1W5B) [32]. The coordinates for the missing residues

of the proteins were obtained by modeling using the MOE package [31]. Since the C-terminus has not been included in the electron crystallography data for the tubulin structure, we did not consider it in our calculations. The missing hydrogens for heavy atoms were added using the tLEAP module of AMBER [33] with the AMBER14SB force field. The protonation states of all ionizable residues were determined at pH = 7 using the MOE program. Each protein model was solvated in a 12 Å box of TIP3P water. Na⁺ and Cl⁻ ions were added in order to bring the salt concentration to the physiological value of 0.15 M. After minimization, the MD simulations were carried out in three steps: heating, density equilibration, and production. First, each solvated system was heated to 300 K for 50 ps, with weak restraints on all backbone atoms. Next, density equilibration was carried out for 50 ps of constant pressure equilibration at 300 K, with weak restraints. Finally, MD production runs were performed on all systems for 100 ns. Ligands and ions were all removed from the complex after equilibration in order to avoid false representations of the protein since ligands can form an occlusion during the process of protein recognition. After equilibration, density-based clustering algorithm from the AMBER software was used for cluster analysis of MD trajectories (20). Several snapshots from top clusters were selected for all further calculations in the study. The result of our simulation is a PDB-formatted file (a 3D representation of all atoms comprising the protein), which is converted into a MAT file using a MATLAB script.

2.3 Data augmentation

To expand the dataset for FtsZ, a data augmentation technique is used where each structure is rotated around the Z-axis in 40° steps. Subsequently, the 3D protein representation is ready to be used for feature extraction. It was not necessary to follow the

same procedure for tubulin since we have many examples available. The purpose of reorienting the z-axis is not only to obtain additional examples, but also in order to not have a bias inside the classifier, in fact most of the rotated proteins were used during the test phase. Both hemispheres of the protein were used to have a complete dataset. Then, to avoid the over-fitting problem a k-fold cross validation is implemented with $k = 5$. One such example is shown in Figure 7 (<https://probis.nih.gov/>) [34].

Figure 7

At this point the 3D protein representation is ready and the feature extraction can be performed.

2.4 Protein samples

In this computational experiment, we used 889 examples of tubulin structure files, divided into 9 isotypes, as shown in Table 1.

Table 1

Using data augmentation, the 13 FtSZ protein samples were rotated in order to create 65 samples, most of them used only during the test phase. The binary classification between tubulin and FtSZ was performed using the samples shown in Table 2.

Table 2

2.5 Data processing

The x-, y- and z-coordinates were extracted from the PDB file. First, the data were shifted in order to be symmetric with respect to x-, y- and z- axes, i.e. the center of the coordinate systems is the center of the symmetry of the dataset. Then, the data were divided into two groups of positive and negative z-values. Finally, for each group, the exterior surface with a desired resolution was calculated using "meshgrid" and "griddata"

406 commands in Matlab with the cubic interpolation method. The descriptors were mapped
407 onto the surfaces as follows. The surfaces were given by point clouds where points are
408 non-connected (not a mesh) and arranged in a square grid. This type of data is called
409 depth map and can be described by matrices: X, Y, Z, where Z is the one describing the
410 “surface” and is represented in these formulas as h. Through Matlab “gradient” function,
411 the derivatives h_x , h_y ... were evaluated so that other matrices representing the first
412 derivative with respect to x, the first derivative with respect to y, etc., were generated and
413 stored. Then, the implementation formulas for the descriptors were calculated on the
414 matrices previously computed and new matrices were obtained representing every
415 geometrical descriptor.

416 For each protein the Z axis was divided in two files: one for the positive part and
417 the second for the negative part using the formula: $z - \max(z) + (|\max(z) - \min(z)|)/2$.
418 Each part represents a “face” of the protein and the geometrical features were computed
419 for both the faces. Then, for every geometrical descriptor a 9-bin histogram was created
420 with the same equidistance for the X-axis.

421 The MATLAB code loaded all data and the following processing steps were performed
422 for all the datasets:

- 423 • the class of the protein was extracted from the filename and the class was
424 recorded in the first column of the dataset matrix;
- 425 • geometrical descriptors were computed from matrix Z (positive and negative);
- 426 • histograms were created and each bin was written in the right column of the
427 dataset matrix;
- 428 • at the end of each loop the dataset matrix became the input for the classifier.

429 The entire process is summarized in Figure 8.

430 **Figure 8**

431 In this computational experiment, 9 isotypes were used (indeed, the classifier will work
432 with 9 classes). The classes were chosen 1 to 9 in an ascendant order as shown in Table 3.

433 **Table 3**

434 This task was performed using a switch case construct. The right class was written in the
435 first column of the Features Matrix.

436 **2.6 Feature extraction**

437 For every geometrical descriptor, a 9-bin histogram was created. Since it is
438 possible that some descriptors have values $\in \mathbb{C}$ (complex), a check was performed first.
439 The geometrical descriptors were calculated using 9 bins and the X-axis values were
440 compressed between -0.2 and 0.2, then the Y-axis values were saved and used as features.
441 Some examples of histograms are shown in Figure 9.

442 **Figure 9**

443 Finally, when all descriptors for all protein data were computed, the resultant
444 matrix was copied into a file. For tubulin and other proteins these descriptors can
445 underline specific characteristic of a certain surface. They can indicate different trends of
446 the free form analyzed and they can describe the shape of the surface. The features are
447 extracted with multiple geometrical descriptors to extract more details; using this
448 approach, also small differences in convexity and concavity can be recognized during the
449 classification. Analyzing the features extracted, the most important features were found
450 from parameter values of Fden2 and sing, because analyzing the data these values were
451 sufficiently different to help the classifier select the right class.

2.7 Classification

The adopted classifiers were k-means and SVM. First, an unsupervised method was tested (K-means) using 9 clusters and a limited number of iterations, then a supervised method (SVM) using linear and non-linear kernels was used. In these cases, it is not a simply binary classification, but there are many classes (9) and many features (more than 100), so some distributions cannot separate the dataset in a linear way or with a linear separator as a high misclassification rate is reached. An interesting improvement is to use a non-linear separator or a kernel trick. An example of a non-linear kernel is the RBF kernel, which in this test led to positive results.

A linear and a nonlinear kernel (RBF in our case) were chosen in order to see whether a non-linear kernel can reach better results. The difference between linear and non-linear kernel is on the way they divided dataset into classes. A linear kernel uses a linear function to divide it and it is less time consuming but also less precise. A non-linear kernel uses a non-linear function, so it can divide better the dataset. The cross validation is not performed because the results were positive, and the validation part was performed using a large number of parameters and the best ones were selected for the testing part.

2.8 K-means

An unsupervised approach was performed using a k-means classifier implemented in R. The matrix file was loaded and the column with the label was deleted. Then, the classifier was tasked with finding 9 clusters in the input data and at the end there was a comparison made between the clustering and the right label.

K-means works in an iterative way and it performs three steps. In the first step, the dataset is loaded, and the number of clusters is chosen. The centroids are created in a

475 random position. In the second step, each data point is assigned to a nearest cluster. The
476 range for the initialization of the centroids of K-means is set from 2 to 10. The Euclidean
477 distance is computed between a point and every centroid. The minimum distance centroid
478 is chosen as the following cluster:

$$\operatorname{argmin} \operatorname{dist}(c_i, x)^2,$$

479 where c is the centroid and x the data points. In this last phase the centroids are computed
480 again as the mean of all the data points of the cluster:

$$c_i = \frac{1}{|S_i|} \sum x_i,$$

481 where S_i is the sum of a single cluster. Therefore, new centroid positions are computed,
482 and this loop continues until the centroid positions do not change significantly.

483 The stop condition is given by the following criteria:

- 484 • no data points change the cluster;
- 485 • the sum of distances is at the minimum;
- 486 • the maximum number of iterations is reached.

487 Therefore, when the convergence is obtained the algorithms stops.

488 The final result achieved in this example was 76.6%, which is an acceptable
489 result, considering that it is an unsupervised method. Nonetheless, in order to improve the
490 method's accuracy, other types of classifications were tested by us and we discuss them
491 below.

492 **2.9 Support Vector Machine**

493 The first test was performed using a linear kernel where λ is a key parameter of SVM. In
494 fact, the main factors in SVM are setting a large margin and reducing the
495 misclassification rate. These two properties are inversely proportional, and the λ

parameter helps to find a trade-off. A large value of λ is for a small margin, whereas a small value of λ is for a large margin. The right λ parameter depends on the test data. The steps used are as follows:

- the dataset is loaded and features and labels are divided;
 - the dataset is randomly split into 60% training set, 10 % validation set and 30% test set;
 - the training is performed using a linear kernel. We then use different values of λ in the range 10^{-5} to 10^5 and it is evaluated on the validation set. The best parameter found on the validation set is $\lambda = 10^{-5}$ with a score of 95.1%;
- the model is tested and scored on the validation set with the best parameters.

The accuracy obtained changes using different λ values. As a matter of fact, by increasing the λ value, the optimization will choose a smaller margin hyperplane, but the best parameters depend on the dataset and in this case the best value is obtained as $\lambda = 10^{-5}$. The final evaluation on the test set with the best parameter $\lambda = 10^{-5}$ was found to be 92.4%.

The dataset was built using 9 different Tubulin isotypes. Hence, the number of classes used for the SVM classifier was 9; the same number was used in the k-Means test, in order to have comparable results. The confusion matrix is an important tool to evaluate the results, since it gives precise information about misclassification. A confusion matrix without normalization and a normalized confusion matrix are represented in Figure 10. In this case, the accuracy is very high, since there is misclassification found only in one class.

Figure 10

The second test was performed using an RBF kernel. The number of features used was 112 and the dataset was not large, so an approximation of the RBF kernel was not taken into consideration (22). The steps used are as follows:

- the dataset is loaded and features and labels are divided;
- the dataset is randomly split into 60% training set, 10 % validation set and 30% test set;
- the training is done using an RBF kernel. We then use different λ and gamma parameters in the range between 10^{-5} to 10^{15} and it is evaluated on the validation set. The best parameters on the validation set are found to be: $\lambda = 100$ and gamma = 10^{-9} with a score of 98.0%;
- the model is tested and scored on the validation set with the best parameters.

Note that the achieved accuracy changes significantly using different λ and gamma values. Indeed, by increasing the λ value, the optimization will choose a smaller margin hyperplane, but the best parameter depends on the dataset selected and, in this case, the best is 100. The final evaluation on the test set with the best parameter $\lambda = 100$, gamma = 10^{-9} and the accuracy obtained was 96.5%.

The same methodology was applied to tubulin and FtsZ classifications.

3. Results and discussion

In the case of tubulin isotype comparison, the best result was given by the SVM classifier with an RBF kernel. All results are summarized in Table 4.

Table 4

In the case of tubulin and FtsZ comparison, the best result is also given by the SVM classifier with an RBF kernel. All results are summarized in Table 5.

Table 5

These results are competitive with the state-of-the-art results found in the literature. A fast protein classification method [5] based on an SVM classifier reached an accuracy of about 90% with 13 classes. Another study [7] used a semi-supervised classification with a kernel cluster and achieved a 94.3 % accuracy. Consequently, the results of the present study appear to be significant. This work is a starting point toward protein classification based on geometrical features and we expect that even better results can be reached in the future. A natural continuation of this work can be to study important features of a protein, for example characterization of a binding pocket [35] for a ligand, a catalytic domain recognition or a protein-protein interaction interface.

A larger experiment was performed using several additional proteins in order to provide an increased validation for the method proposed in this paper. This test involved four arbitrarily chosen FtsZ protein structures, namely: 2R6R, 2VAW, 2VAP and 2VAM. These structures correspond, respectively, to the following biological species: *B. subtilis*, *Pseudomonas aeruginosa*, *M. jannaschii* and *Aquifex aeolicus*. In this test 683 samples were used as listed in Table 6.

Table 6

The results of this test are very encouraging as shown in Table 7, which summarizes the use of various classifiers for different tests performed and their accuracy levels achieved.

Table 7

To avoid over-fitting and to generalize the method in a better way a 5-fold cross validation is performed. In this way the classifier is not biased by the test set and it also works well with other proteins. The last experiment showed that it also works well with

four very different proteins. In this test a k-cross validation method was applied using $k=5$.

4. Conclusions

A novel method for protein characterization and classification has been proposed in this paper, which is inspired by and uses the algorithms from the facial recognition field. The first application of this method involves a challenging case of classification of highly homologous tubulin isotypes using as features some geometrical descriptors typically found within the context of face recognition analysis. While human faces and proteins represent very different biological structures, they are both free-form surfaces and the same types of geometrical features are adopted for their classification and recognition.

The aim of this study has been to implement different classifiers to be tested on the dataset previously built. In this work we used the following approaches: SVM with a linear RBF kernel, and a K-means algorithm. This methodology and the geometrical descriptors have been used for protein classification. The first classification was performed using the tubulin protein and 9 of its isotypes. The second application performed used two structurally similar proteins: tubulin and FtsZ and third application involved four unrelated proteins. In all cases very encouraging results were obtained.

It should be stressed that until now the use of RMSD as a measure of similarity has been prevalent in protein biophysics, especially regarding structural comparisons. However, this approach relies on a single number, which does not allow for feature extraction or more detailed shape comparisons, which the present methodology provides. A single parameter such as an RMSD value can answer the question if two proteins are

587 structurally similar or not but does not address the issue regarding which features differ
588 between them. For this reason, our method can assist in identifying structure-function
589 dependence when comparing various proteins, even highly similar ones. Since we only
590 investigate geometrical features, both physical and chemical properties are not directly
591 involved in our method but can eventually be extracted by mapping geometrical features
592 back onto to amino acid distributions underlying them.

593 In this study, MD has been used to generate additional models of each protein for
594 the training purpose where each of the models is extracted from equilibrated MD
595 trajectories after clustering. Clustering of the trajectory provides us with different
596 conformations of the same protein from MD trajectories. We used several snapshots from
597 each structural cluster, which makes it possible to probe diverse sampling of the
598 trajectory. In future work, a larger set of protein structures will be used to address the
599 issue of structural diversity across the entire PDB dataset consisting of over 150,000
600 entries.

601 The results obtained and reported here are significant: a 96.5 % accuracy for
602 tubulin isotype classification, a 98.2 % accuracy for tubulin and FtsZ classification and a
603 98% accuracy for a set of four arbitrarily chosen protein structures. SVM is a classifier
604 with competitive performance using a small dataset (< 3000 samples) and in this case the
605 results are significant. The application of a neural network can be a future development
606 using a convolutional type on a larger dataset (> 10,000 samples). The conclusion is that
607 these geometrical descriptors work properly with the description of protein surfaces and
608 they are accurate enough to properly describe protein surfaces.
609 Several future developments can be taken in consideration, namely:

- building a database adding more samples and more proteins;
- computing more features and testing classifiers, using more geometrical descriptors and filters;
- developing more data augmentation techniques to enlarge the dataset;
- identifying specific important features on a protein, for example a binding pocket for a ligand or a protein-protein interaction interface.

Other important improvements will be performed in future tests. First, we will employ neural networks that were applied here with significant results with 3D geometrical descriptors [19]. Second, using a large dataset with unnecessarily numerous features the classifier could be slow, so some feature optimization techniques will be implemented in order to [36] accelerate the training of the kernel machine.

Author Contributions: For research articles with several authors, a short paragraph specifying their individual contributions must be provided. The following statements should be used “conceptualization, F.M. and JAT; methodology, FM, JAT and LDG; software, LDG; validation, LDG; formal analysis, F.M. and JAT ; investigation, FM, JAT and LDG ; Matlab scripts, VR; resources, MA and EV; data curation, MA; writing—original draft preparation, LDG and MA; writing—review and editing, MA, F.M. and JAT; visualization, LDG and MA; supervision, F.M. and JAT; project administration, JAT; funding acquisition, JAT and EV”, please turn to the [CRediT taxonomy](#) for the term explanation.

Funding: This research received no external funding.

Acknowledgments: Computational infrastructure of the Pharma-matrix cluster at the Cross Cancer Institute is gratefully acknowledged.

Conflicts of Interest: The authors declare no conflict of interest.

References:

1. Gainza, P., et al., *Deciphering interaction fingerprints from protein molecular surfaces*. bioRxiv, 2019: p. 606202.
2. Planas-Iglesias, J., et al., *Understanding Protein–Protein Interactions Using Local Structural Features*. Journal of Molecular Biology, 2013. **425**(7): p. 1210-1224.
3. Rupp, B. and J. Wang, *Predictive models for protein crystallization*. Methods, 2004. **34**(3): p. 390-407.
4. Saberi Fathi, S.M., D.T. White, and J.A. Tuszynski, *Geometrical comparison of two protein structures using Wigner-D functions: Geometrical Comparison of Protein Structures*. Proteins: Structure, Function, and Bioinformatics, 2014. **82**(10): p. 2756-2769.
5. Tsuda, K., H. Shin, and B. Schölkopf, *Fast protein classification with multiple networks*. Bioinformatics, 2005. **21**(suppl_2): p. ii59-ii65.
6. Weston, J., et al., *Semi-supervised protein classification using cluster kernels*. Bioinformatics, 2005. **21**(15): p. 3241-3247.
7. Jain, P., J.M. Garibaldi, and J.D. Hirst, *Supervised machine learning algorithms for protein structure classification*. Computational Biology and Chemistry, 2009. **33**(3): p. 216-223.
8. Masci, J., et al. *Geodesic Convolutional Neural Networks on Riemannian Manifolds*. 2015. IEEE Computer Society.

- 654 9. Monti, F., et al. *Geometric Deep Learning on Graphs and Manifolds Using*
655 *Mixture Model CNNs*. in *Proceedings of the IEEE Conference on Computer Vision and*
656 *Pattern Recognition*. 2017.
- 657 10. Bronstein, M.M., et al., *Geometric Deep Learning: Going beyond Euclidean data*.
658 *IEEE Signal Processing Magazine*, 2017. **34**(4): p. 18-42.
- 659 11. Espinosa, E., et al., *Classification of anticancer drugs—a new system based on*
660 *therapeutic targets*. *Cancer Treatment Reviews*, 2003. **29**(6): p. 515-523.
- 661 12. Huang, C.-H., F.-L. Lee, and C.-J. Tai, *The β -tubulin gene as a molecular*
662 *phylogenetic marker for classification and discrimination of the *Saccharomyces sensu**
663 *stricto complex*. *Antonie van Leeuwenhoek*, 2009. **95**(2): p. 135-142.
- 664 13. Ludueña, R.F., *Are tubulin isotypes functionally significant*. *Molecular Biology of*
665 *the Cell*, 1993. **4**(5): p. 445-457.
- 666 14. Fitch, W.M., *Homology: a personal view on some of the problems*. *Trends in*
667 *Genetics*, 2000. **16**(5): p. 227-231.
- 668 15. Richards, K.L., et al., *Structure–Function Relationships in Yeast Tubulins*.
669 *Molecular Biology of the Cell*, 2000. **11**(5): p. 1887-1903.
- 670 16. Schlieper, D., et al., *Structure of bacterial tubulin BtubA/B: Evidence for*
671 *horizontal gene transfer*. *Proceedings of the National Academy of Sciences of the United*
672 *States of America*, 2005. **102**(26): p. 9170-9175.
- 673 17. Gunn, S.R., *Support Vector Machines for Classification and Regression*. p. 66.
- 674 18. [https://github.com/lucaresearch/A-new-protein-characterization-and-](https://github.com/lucaresearch/A-new-protein-characterization-and-classification-method-using-3D-face-recognition-algorithms)
675 [classification-method-using-3D-face-recognition-algorithms](https://github.com/lucaresearch/A-new-protein-characterization-and-classification-method-using-3D-face-recognition-algorithms).

- 676 19. Ciravegna, G., et al., *Assessing Discriminating Capability of Geometrical*
677 *Descriptors for 3D Face Recognition by Using the GH-EXIN Neural Network*, in *Neural*
678 *Approaches to Dynamics of Signal Exchanges*, A. Esposito, et al., Editors. 2020,
679 Springer: Singapore. p. 223-233.
- 680 20. Cirrincione, G., et al., *Intelligent Quality Assessment of Geometrical Features for*
681 *3D Face Recognition*, in *Neural Advances in Processing Nonlinear Dynamic Signals*, A.
682 Esposito, et al., Editors. 2019, Springer International Publishing: Cham. p. 153-164.
- 683 21. Li, S.Z. and A.K. Jain, *Handbook of Face Recognition*. 2 ed. 2011, London:
684 Springer-Verlag.
- 685 22. Marcolin, F., et al., *Three-dimensional face analysis via new geometrical*
686 *descriptors*, in *Advances on Mechanics, Design Engineering and Manufacturing :*
687 *Proceedings of the International Joint Conference on Mechanics, Design Engineering &*
688 *Advanced Manufacturing (JCM 2016), 14-16 September, 2016, Catania, Italy*, B. Eynard,
689 et al., Editors. 2017, Springer International Publishing: Cham. p. 747-756.
- 690 23. Marcolin, F. and E. Vezzetti, *Novel descriptors for geometrical 3D face analysis*.
691 *Multimedia Tools and Applications*, 2017. **76**(12): p. 13805-13834.
- 692 24. Vezzetti, E. and F. Marcolin, *Geometrical descriptors for human face*
693 *morphological analysis and recognition*. *Robotics and Autonomous Systems*, 2012.
694 **60**(6): p. 928-939.
- 695 25. *MATLAB*. 2018, The MathWorks Inc.: Natick, Massachusetts.
- 696 26. *Anaconda Software Distribution. Computer software*. 2019, Anaconda,.
- 697 27. Van Rossum, G. and F.L. Drake Jr, *Python*. 2019, Centrum voor Wiskunde en
698 Informatica Amsterdam, The Netherlands.

- 699 28. Version 0.22.0 — scikit-learn 0.22 documentation - [https://scikit-](https://scikit-learn.org/stable/whats_new/v0.22.html)
- 700 [learn.org/stable/whats_new/v0.22.html](https://scikit-learn.org/stable/whats_new/v0.22.html).
- 701 29. Download R-3.5.3 for Windows. The R-project for statistical computing -
- 702 <https://cran.r-project.org/bin/windows/base/old/3.5.3/>.
- 703 30. Löwe, J., et al., *Refined structure of $\alpha\beta$ -tubulin at 3.5 Å resolution* | Edited by I.
- 704 A. Wilson. Journal of Molecular Biology, 2001. **313**(5): p. 1045-1057.
- 705 31. Molecular Operating Environment (MOE). Group, Chemical Computing. 2012:
- 706 Montreal, QC, Canada.
- 707 32. Oliva, M.A., S.C. Cordell, and J. Löwe, *Structural insights into FtsZ*
- 708 *protofilament formation*. Nature Structural & Molecular Biology, 2004. **11**(12): p. 1243-
- 709 1250.
- 710 33. D.A. Case, et al., *AMBER 2014*. 2014: University of California, San Francisco.
- 711 34. Konc, J., et al., *ProBiS-CHARMMing: Web Interface for Prediction and*
- 712 *Optimization of Ligands in Protein Binding Sites*. Journal of Chemical Information and
- 713 Modeling, 2015. **55**(11): p. 2308-2314.
- 714 35. Saberi Fathi, S.M. and J.A. Tuszynski, *A simple method for finding a protein's*
- 715 *ligand-binding pockets*. BMC Structural Biology, 2014. **14**: p. 18.
- 716 36. Rahimi, A. and B. Recht, *Random Features for Large-Scale Kernel Machines*. p.
- 717 10.

718 **Figures' captions:**

719 **Figure 1:** Structural similarities between tubulin and FtsZ proteins. The tubulin dimer
720 consists of an α -tubulin and a closely related β -tubulin monomer. $\alpha\beta$ -tubulin
721 heterodimers associate head to tail to form protofilaments and laterally to form the

722 cylindrical MT wall. GTP and GDP nucleotides (ball and stick models) are bound to a
723 and β tubulin, respectively. (b) The FtsZ dimer consists of two identical monomers with
724 GTP bound to N-terminals (blue). In both (a) and (b) N-terminals (blue) and C-terminals
725 (red) are separated by H7 helices (green). N-terminal regions show the typical nucleotide-
726 binding motif with parallel β sheets connected by a helices known as the Rossmann fold.
727 By comparing the two protein structures, the differences in C-terminal regions are
728 obvious. GDP and GTP are shown in ball and stick models. The figures were rendered
729 using the MOE (Molecular Operating Environment) software. PDB ID for tubulin: 1JFF.
730 PDB ID for FtsZ: 1W5B.

731 **Figure 2:** Valid solutions can be found with perceptron in a binary case(a) and the best
732 theoretical solution that a SVM classifier can find (b).

733 **Figure 3:** Flow chart of the entire protein characterization and classification process.

734 **Figure 4:** Effects of applying different descriptors (a) F_den2 ,(b) $sing$ (c) , S_{mean} , (d)
735 k_{1mean} ,(e) $k_{2median}$, (f), g_{mean} ,and (g) H_{mean} to a human face (left column) and to
736 the tubulin protein (right column)

737 **Figure 5:** Sequence alignment of β tubulin isotypes. Each of the human β tubulin
738 isotypes that were identified in our screen of the UniProt databases were aligned using
739 the MOE package. Prior to performing the alignment, the highly variable carboxy-
740 terminal residues were removed from each sequence. This was done as the template
741 structure, 1JFF, does not contain any of these residues. At each position within the
742 alignment, dark blue boxes indicate identical residues; light blue boxes indicate residues
743 that are conserved, while red boxes indicate residues that are divergent (poorly aligned).

Figure 6: a) Sequence similarity matrix and (b) sequence identity matrix of the studied tubulin isotypes

Figure 7: Tubulin protein image for two different rotations with respect to the Z-axis.

Figure 8: Protein data processing overview. The input consists of a 3D structure of a protein from either the PDB database or from homology modeling combined with MD simulations. The color selection in the input structure is arbitrarily chosen for better visualization. The output consists of geometrical descriptor values obtained from a facial recognition algorithm.

Figure 9: 9 bin histograms calculated using (a) F_{den2} , (b) g_{mean} and (c) H_{mean} geometrical descriptor

Figure 10: Confusion matrix of SVM classifier using the RBF kernel.

Figures:

Figure 1

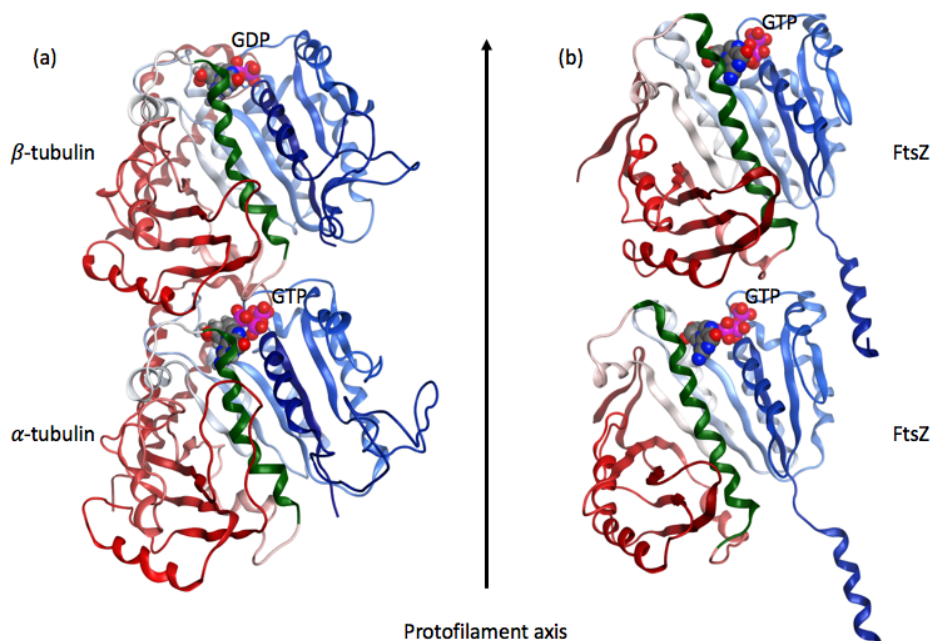


Figure 2

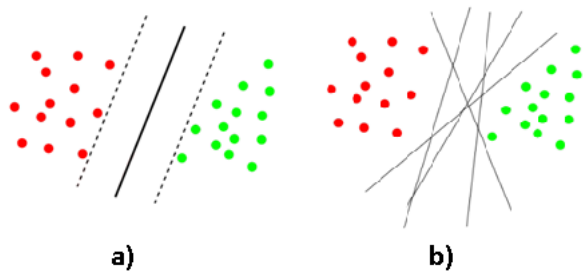
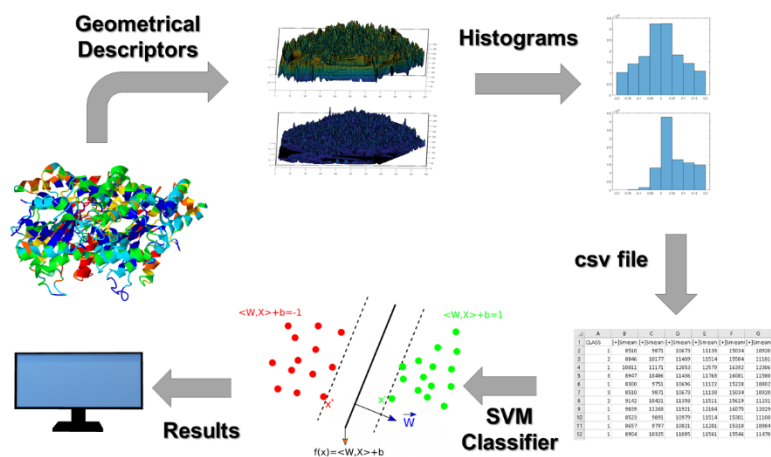
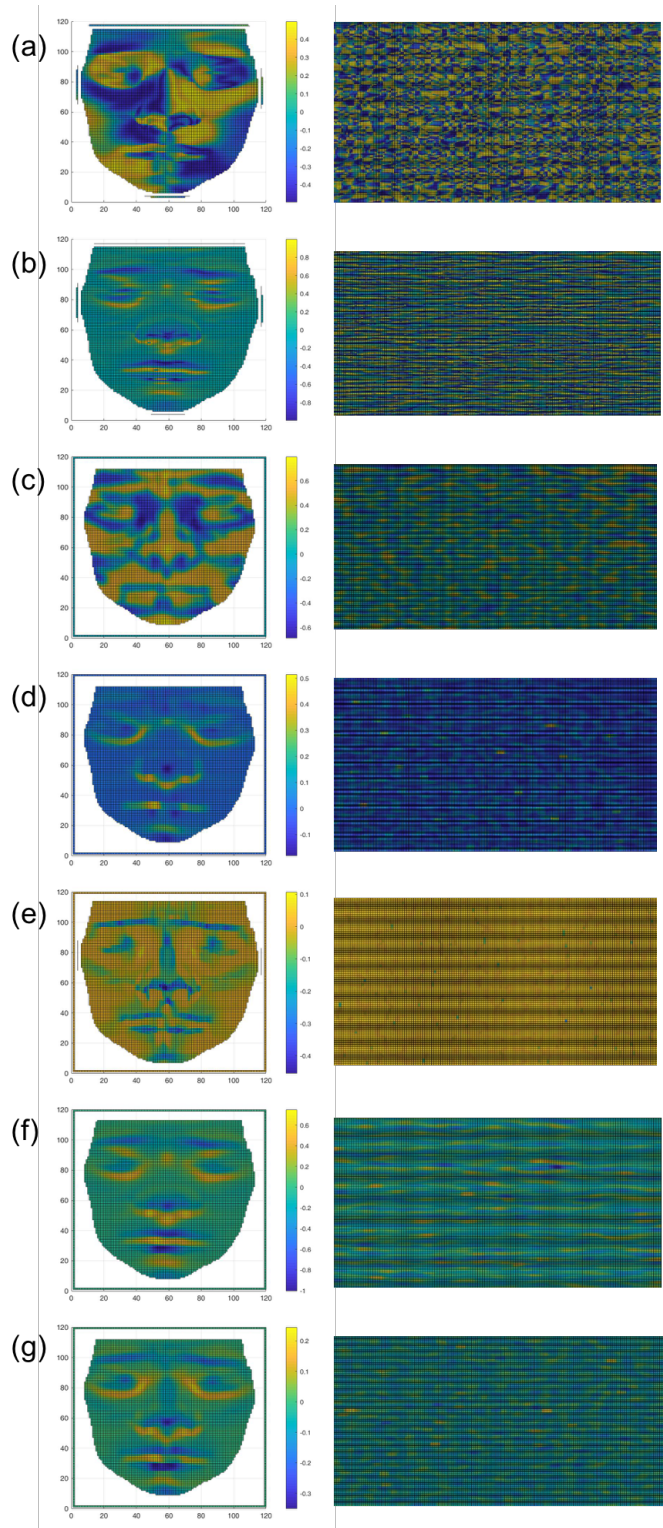


Figure 3



773 **Figure 4**



774

775

776 **Figure 5**

Chain	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75																																																										
1: beta1	M	R	E	I	V	H	I	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	T	G	T	Y	H	G	S	D	L	Q	L	R	I	S	V	Y	N	E	A	T	G	G	K	Y	V	P	R	A	I	L	V	D	L	E	P	G	T	M	D	S				
2: betaIIA	M	R	E	I	V	H	I	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	T	G	S	Y	H	G	S	D	L	Q	L	E	R	I	N	V	Y	N	E	A	A	G	N	K	Y	V	P	R	A	I	L	V	D	L	E	P	G	T	M	D	S			
3: betaIIB	M	R	E	I	V	H	I	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	T	G	S	Y	H	G	S	D	L	Q	L	E	R	I	N	V	Y	N	E	A	T	G	N	K	Y	V	P	R	A	I	L	V	D	L	E	P	G	T	M	D	S			
4: betaIII	M	R	E	I	V	H	I	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	S	G	N	Y	V	G	S	D	L	Q	L	E	R	I	S	V	Y	N	E	A	S	H	K	Y	V	P	R	A	I	L	V	D	L	E	P	G	T	M	D	S				
5: betaIVA	M	R	E	I	V	H	L	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	T	G	T	Y	H	G	S	D	L	Q	L	E	R	I	N	V	Y	N	E	A	T	G	G	N	Y	V	P	R	A	V	L	V	D	L	E	P	G	T	M	D	S			
6: betaIVB	M	R	E	I	V	H	L	Q	A	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	G	I	D	P	T	G	T	Y	H	G	S	D	L	Q	L	E	R	I	N	V	Y	N	E	A	T	G	G	K	Y	V	P	R	A	V	L	V	D	L	E	P	G	T	M	D	S			
7: betaV	M	R	E	I	V	H	I	Q	A	G	C	G	N	I	G	T	K	F	W	E	V	I	S	D	E	H	G	I	D	P	A	G	G	Y	V	G	S	A	L	Q	L	E	R	I	N	V	Y	N	E	S	S	Q	K	Y	V	P	R	A	L	V	D	L	E	P	G	T	M	D	S					
8: betaVI	M	R	E	I	V	H	I	Q	G	C	G	N	I	G	A	K	F	W	E	M	I	G	E	H	G	I	D	L	A	G	S	D	R	G	A	S	A	L	Q	L	E	R	I	S	V	Y	N	E	A	Y	G	R	K	Y	V	P	R	A	V	L	V	D	L	E	P	G	T	M	D	S				
9: betaVIII	M	R	E	I	V	L	T	I	G	C	G	N	I	G	A	K	F	W	E	V	I	S	D	E	H	A	I	D	S	A	G	T	Y	H	G	S	H	L	Q	L	E	R	I	N	V	Y	N	E	A	S	G	G	R	Y	V	P	R	A	V	L	V	D	L	E	P	G	T	M	D	S				
Chain	76	80	85	90	95	100	105	110	115	120	125	130	135	140	145	150																																																										
1: beta1	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	V	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	S	V	L	D	V	V	R	K	E	A	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
2: betaIIA	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	V	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	S	V	L	D	V	V	R	K	E	S	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
3: betaIIB	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	V	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	S	V	L	D	V	V	R	K	E	S	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
4: betaIII	V	R	S	G	A	F	G	H	L	F	R	P	D	N	F	I	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	S	V	L	D	V	V	R	K	E	C	E	N	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
5: betaIVA	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	V	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	A	V	L	D	V	V	R	K	E	A	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
6: betaIVB	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	V	F	G	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	S	V	L	D	V	V	R	K	E	A	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L
7: betaV	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	I	F	G	Q	T	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	V	D	A	V	L	D	V	V	R	K	E	C	H	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L	
8: betaVI	I	R	S	S	K	L	G	A	L	F	Q	P	D	S	F	V	H	G	N	S	G	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	I	N	V	L	E	V	V	R	H	E	S	E	S	C	D	C	L	Q	G	F	I	V	H	S	L	G	G	G	T	G	S	G	M	G	T	L	
9: betaVIII	V	R	S	G	P	F	G	Q	I	F	R	P	D	N	F	I	F	G	C	A	G	N	N	W	A	K	G	H	Y	T	E	G	A	E	L	M	E	S	V	M	D	V	V	R	K	E	A	E	S	C	D	C	L	Q	G	F	Q	L	T	H	S	L	G	G	G	T	G	S	G	M	G	T	L	
Chain	151	155	160	165	170	175	180	185	190	195	200	205	210	215	220	225																																																										
1: beta1	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
2: betaIIA	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
3: betaIIB	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
4: betaIII	L	I	S	K	V	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	I	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	A	P	T	Y	G	D	L							
5: betaIVA	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
6: betaIVB	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
7: betaV	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	T	D	E	T	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L						
8: betaVI	L	M	N	K	I	R	E	E	P	D	R	I	N	T	F	S	V	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	I	H	Q	L	E	N	A	D	A	C	F	C	I	D	N	E	A	L	D	I	C	F	R	T	L	K	L	T	T	P	T	Y	G	D	L					
9: betaVIII	L	I	S	K	I	R	E	E	P	D	R	I	N	T	F	S	I	L	P	S	P	K	V	S	D	T	V	V	E	P	N	A	T	L	S	V	H	Q	L	E	N	A	D	E	T	F	C	I	D	N	E	A	L	D	I	C	S	K	T	L	K	L	P	T	P	T	Y	G	D	L				
Chain	226	230	235	240	245	250	255	260	265	270	275	280	285	290	295	300																																																										
1: beta1	N	H	L	V	S	A	T	M	S	G	V	T	T	C	L	R	F	P	G	Q	L	N	A	D	L	R	K	L	A	V	N	M	V	P	F	P	R	L	H	F	M	P	G	F	A	P	L	T	S	R	G	S	Q	Q	Y	R	A	L	T	V	P	E	L	T	Q	Q	V	F	D	A	K	N	M	M
2: betaIIA	N	H	L	V	S	A	T	M	S	G	V	T	T	C	L	R	F	P	G	Q	L	N	A	D	L	R	K	L	A	V	N	M	V	P	F	P	R	L	H	F	M	P	G	F	A	P	L	T	S	R	G	S	Q	Q	Y	R	A	L	T	V	P	E	L	T	Q	Q	M	F	D	A	K	N	M	M
3: betaIIB	N	H	L	V	S	A	T	M	S	G	V	T	T	C	L	R	F	P	G	Q	L	N	A	D	L	R	K	L	A	V	N	M	V	P	F	P	R	L	H	F	M	P	G	F	A	P	L	T	S	R	G	S	Q	Q	Y	R	A	L	T	V	P	E	L	T	Q	Q	M	F	D	A	K	N	M	M
4: betaIII	N	H	L	V	S	A	T	M	S	G	V	T	T	S	L	R	F	P	G	Q	L	N	A	D	L	R	K	L	A	V	N	M	V	P	F	P	R	L	H	F	M	P	G	F	A	P	L	T	A	R	G	S	Q	Q	Y	R	A	L	T	V	P	E	L	T	Q	Q	M	F	D	A	K	N	M	M
5: betaIVA	N	H	L	V	S	A	T	M	S	G	V	T	T	C	L	R	F	P	G	Q	L	N																																																				

Figure 6

(a)	1	2	3	4	5	6	7	8	9
1:betaI		99.1	99.5	97.0	99.5	99.8	96.3	90.4	96.0
2:betaIIA	99.1		99.5	97.4	99.1	99.3	96.3	90.4	95.8
3:betaIIB	99.5	99.5		97.7	99.5	99.8	96.5	90.6	96.0
4:betaIII	97.0	97.4	97.7		97.2	97.2	97.9	90.9	93.7
5:betaIVA	99.5	99.1	99.5	97.2		99.8	96.5	90.4	95.8
6:betaIVB	99.8	99.3	99.8	97.2	99.8		96.5	90.6	96.3
7:betaV	96.3	96.3	96.5	97.9	96.5	96.5		90.2	93.2
8:betaVI	90.4	90.4	90.6	90.9	90.4	90.6	90.2		89.2
9:betaVIII	96.0	95.8	96.0	93.7	95.8	96.3	93.2	89.2	

(b)	1	2	3	4	5	6	7	8	9
1:betaI		97.0	97.4	93.9	97.4	98.4	92.5	80.1	89.7
2:betaIIA	97.0		99.5	93.4	96.3	97.7	92.5	80.8	89.9
3:betaIIB	97.4	99.5		93.7	96.7	98.1	92.7	81.0	90.2
4:betaIII	93.9	93.4	93.7		93.2	93.9	94.4	80.1	87.6
5:betaIVA	97.4	96.3	96.7	93.2		98.6	93.4	80.3	90.2
6:betaIVB	98.4	97.7	98.1	93.9	98.6		92.7	80.3	91.1
7:betaV	92.5	92.5	92.7	94.4	93.4	92.7		80.1	87.1
8:betaVI	80.1	80.8	81.0	80.1	80.3	80.3	80.1		78.2
9:betaVIII	89.7	89.9	90.2	87.6	90.2	91.1	87.1	78.2	

Figure 7

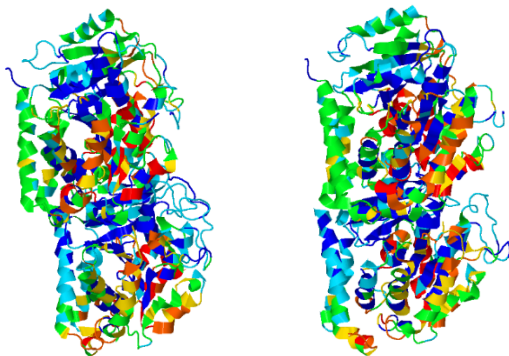


Figure 8

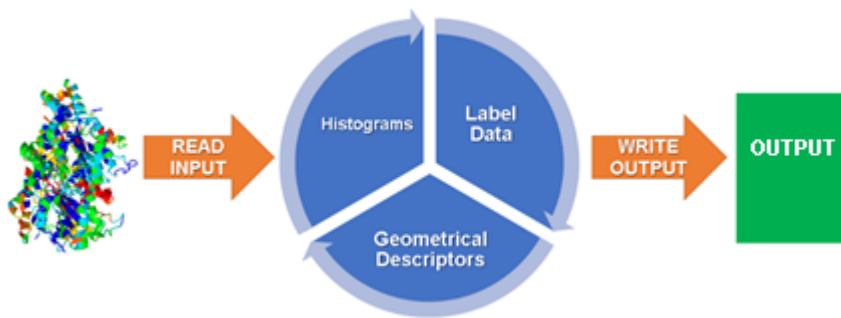


Figure 9

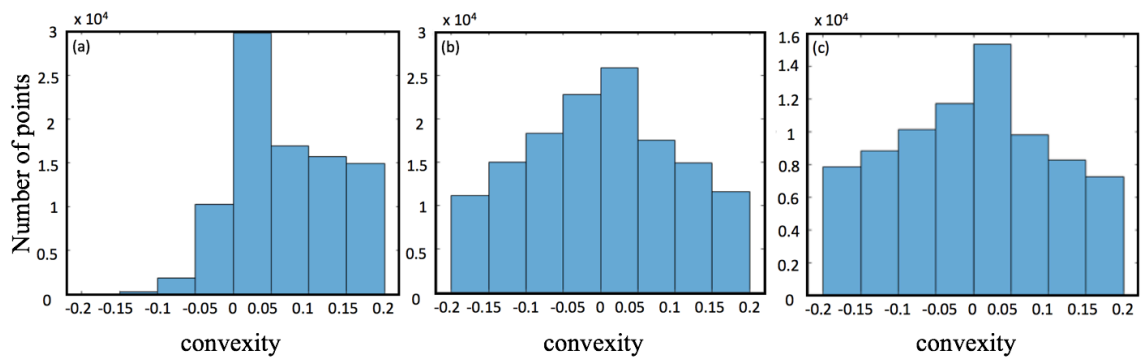
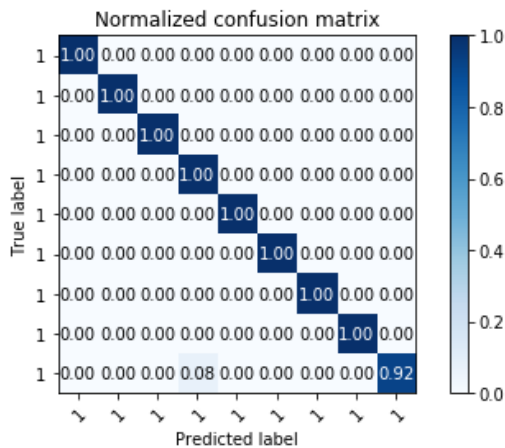


Figure 10.



Tables' captions:

Table 1: Numbers of tubulin isotype structures used.

Table 2: Sample numbers in the binary classification between tubulin and FtsZ.

Table 3: Number of Tubulin isotypes used.

Table 4: Tubulin isotypes accuracy results.

Table 5: Accuracy results for the tubulin and FtsZ binary classification.

Table 6: 2R6R, 2VAM, 2VAP and 2VAM samples.

Table 7: 2R6R, 2VAM, 2VAP and 2VAM experiment.

Tables:

803 **Table 1**

Isotypes	Beta I	Beta IIa	Beta IIb	Beta III	Beta IVa	Beta IVb	Beta V	Beta VI	Beta VIII
Samples	123	128	94	57	128	68	107	62	125

804

805 **Table 2**

Protein	Samples
Tubulin	112
FtsZ	65

806

807 **Table 3**

Isotypes	Beta I	Beta IIa	Beta IIb	Beta III	Beta IVa	Beta IVb	Beta V	Beta VI	Beta VIII
Samples	1	2	3	4	5	6	7	8	9

808

809 **Table 4**

Classifier	Accuracy
SVM with RBF kernel	96.5 %
SVM with linear kernel	92.4 %
K-means	76.6 %

810

811 **Table 5**

Classifier	Accuracy
SVM with RBF kernel	98.2 %
SVM with linear kernel	97.0 %
k-means	72.3 %

812

813 **Table 6**

Proteins Samples	2R6R 175	2VAW 170	2VAP 168	2VAM 170
------------------	-------------	-------------	-------------	-------------

814

815 **Table 7**

Classifier	Accuracy
SVM with RBF kernel	97.1 %
SVM with linear kernel	98.0 %
K-means	62.3 %