

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

Luca Di Grazia<sup>1</sup> | Maral Aminpour<sup>2,3</sup> | Enrico Vezzetti<sup>1</sup> | Vahid Rezania<sup>4</sup> | Federica Marcolin<sup>1</sup> | Jack Adam Tuszyński<sup>1,2,3</sup> 

<sup>1</sup>DIGEP, Politecnico di Torino, Torino, Italy

<sup>2</sup>Department of Physics, University of Alberta, Edmonton, Alberta, Canada

<sup>3</sup>Department of Oncology, University of Alberta, Edmonton, Alberta, Canada

<sup>4</sup>Department of Physical Sciences, MacEwan University, Edmonton, Alberta, Canada

## Correspondence

Jack Adam Tuszyński, Department of Oncology, University of Alberta, Edmonton, AB, Canada.

Email: jacek.tuszyński@polito.it

## Abstract

This article 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 analog. 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 kernel. The results are significant and competitive with the state-of-the-art protein classification methods. This leads to a new methodological direction in protein structure analysis.

## KEY WORDS

3D face analysis, differential geometry, geometrical descriptors, machine learning, protein classification, support vector machine, tubulin

## 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<sup>1</sup> 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<sup>2</sup> showed the relevant role of protein-protein interactions using local structural features. In this latter article, 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 and cryogenic electron microscopy (cryo-EM) are other methods by which to determine the protein structure<sup>3</sup> 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.<sup>4</sup> However, proteins are nonstatic 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 article 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 synthetically generated from computational approaches such as MD. The key aspect is to have an atomistic model of the objects studied,<sup>3</sup> which serves as the starting point for feature extraction based on the protein surface. Such a model provides a high-resolution representation of the object of interest, which is later on processed and characterized by a manageable number of parameters.

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

Some experimentation has already been undertaken to classify proteins according to their states. Tsuda et al adopted a support vector machine (SVM) classifier for fast protein classification.<sup>5</sup> They obtained 13 classes and reached an accuracy of about 90%. Weston et al<sup>6</sup> used a semisupervised classification with a kernel cluster and reached a result of 94.3%. Another interesting result has been obtained using a random forest approach and 15 different supervised methods with about 11 000 pairs of protein domains leading to an accuracy of 97.0%.<sup>7</sup> Our focus in this article 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. This could be valuable with respect to their functional roles in cell biology as well as potential applications in rational drug design.

One of the most important proteins abundantly expressed in all eukaryotic cells is the family of tubulin proteins, which is studied in this article 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.

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 characterization, because geometric feature surfaces were used with significant results in many other applications before, for example, References<sup>8,9</sup> We believe that the classification provided here can be further improved with more data, more classes, and a complex neural network. A complex neural network is one of the applications we are planning to implement in the near future. We intend to use a convolutional neural network to minimize the cost function to cluster the inputs correctly, because this could be an efficient way to find a pattern in the input data and it can be a significant improvement for our objectives. All of which is planned for future work, especially within the context of geometric deep learning,<sup>10</sup> which nowadays is the state-of-the art of classification.

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

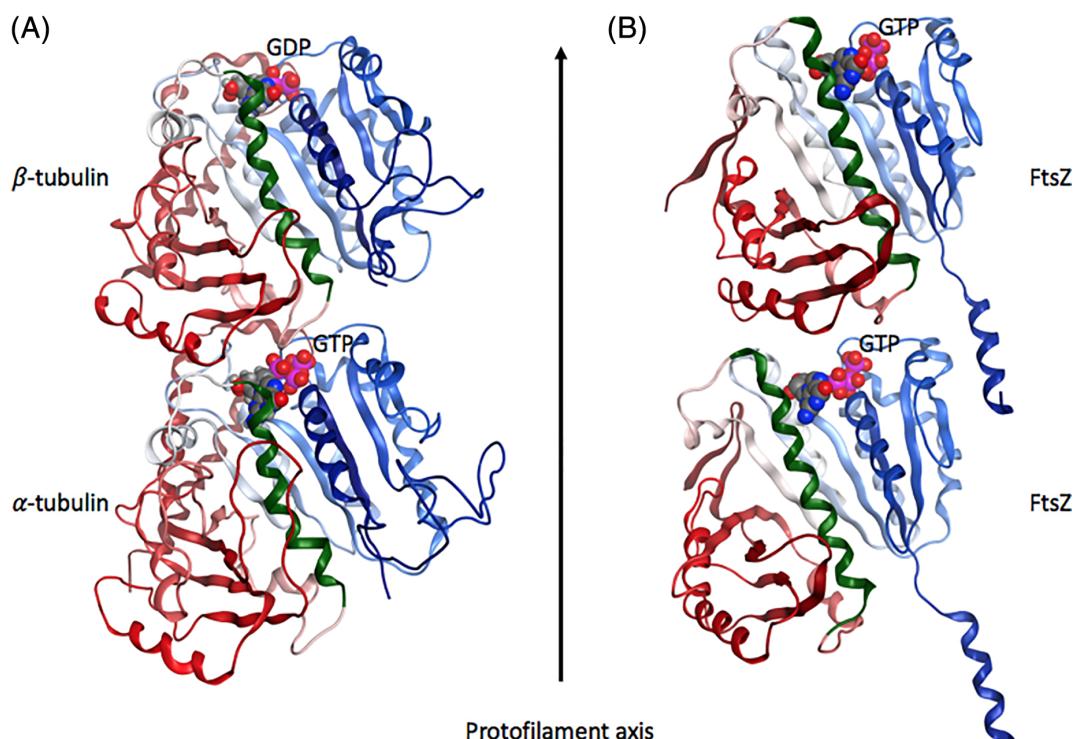
The novelty of the present work rests with the application of geometrical descriptors coming from the field of face analysis to the classification of surfaces of proteins, with the aim of adopting this geometrical information as descriptive features and discriminating elements to classify proteins. Here, we test the method on the examples of tubulin isotypes and related proteins (eg, FtsZ). The method can, of course, be applied to an arbitrary protein or indeed a protein complex but being able to discriminate between highly homologous proteins based on the geometrical shapes of their surfaces opens the door to numerous applications across the field of protein science. The idea comes from the realization that geometrical properties can well describe the surface of a 3D object such as a protein and could identify characteristic features when comparing two or more similar structures. Proteins surfaces can be split into two outer surfaces by cutting a plane through the data set including the main axis of rotational symmetry. These two halves of the outer surface, similarly to human faces, differ from one another depending on the protein type, and also can change their conformational states dynamically, similarly to human facial expressions. Thus, what in the field of pattern recognition is called face recognition could be transferred to the context of

protein classification according to the typology. These common points have fostered the interest of uncovering the potentiality of cross-fertilization between these two fields with the aim of better categorization.

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

in drug design targeting specific isotypes such as  $\beta$ III, which is over-expressed in all cancer cells. Through a search of available protein sequence databases, a total of 10 unique  $\beta$  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 homolog 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  $\alpha$  and  $\beta$  tubulin share an approximate 35% sequence identity with FtsZ.<sup>16</sup> 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 bacterial cell cycle, but instead just plays a structural role during bacterial cell division, when it forms a



**FIGURE 1** Structural similarities between tubulin and FtsZ proteins. The tubulin dimer consists of an  $\alpha$ -tubulin and a closely related  $\beta$ -tubulin monomer.  $\alpha\beta$ -tubulin heterodimers associate head to tail to form protofilaments and laterally to form the cylindrical microtubule wall. GTP and GDP nucleotides (ball and stick models) are bound to  $\alpha$  and  $\beta$  tubulin, respectively. (B) The FtsZ dimer consists of two identical monomers with GTP bound to N-terminals (blue). In both (A) and (B) N-terminals (blue) and C-terminals (red) are separated by H7 helices (green). N-terminal regions show the typical nucleotide-binding motif with parallel  $\beta$  sheets connected by  $\alpha$  helices known as the Rossmann fold. By comparing the two protein structures, the differences in C-terminal regions are obvious. GDP and GTP are shown in ball and stick models. The figures were rendered using the MOE (Molecular Operating Environment) software. PDB ID for tubulin: 1JFF. PDB ID for FtsZ: 1W5B

band, known as the Z-ring, around the inner cell wall at the location where the cell will divide.

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 MD 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.

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.<sup>17</sup> SVM is able to find a solution with a larger margin for the two-separator classifier as shown in Figure 2A. 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 article 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 pCloud repository.<sup>18</sup> The entire process is shown in Figure 3.

This article is organized as follows. In Section 2 geometrical descriptors used for implementing the feature extraction are described. Section 3 is the core of the article 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 | GEOMETRICAL DESCRIPTORS

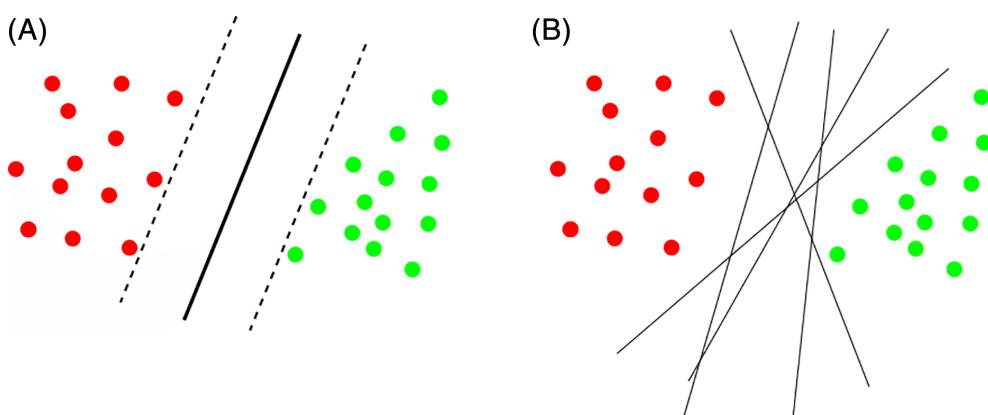
The surfaces representing both human faces and proteins are geometrically considered as a free form. Thus, features coming from the field of differential geometry can be applied in order to understand their local and global properties. Geometrical descriptors are widely used in the area of 3D face recognition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

The geometrical descriptors used in this research are the following geometrical descriptors<sup>22,23</sup>: mean curvature ( $H_{\text{mean}}$ ), principal curvatures ( $k_{1\text{mean}}$  and  $k_{2\text{median}}$ ), the shape index ( $S_{\text{mean}}$ ), the third coefficient of the second fundamental form ( $g_{\text{mean}}$  and  $\sin g$ ), and a descriptor enlightening the symmetry property ( $F_{\text{den}2}$ ). Considering that these descriptors rely on the derivatives of the surface ( $h_x$ ,  $h_y$ ,  $k_{1\text{mean}}$ ,  $\sin g$ ,  $k_{2\text{median}}$ ,  $g_{\text{mean}}$ ,  $H_{\text{mean}}$ ), depressions and peaks (local minima and maxima) of the surface ( $k_{1\text{mean}}$ ,  $\sin g$ ,  $k_{2\text{median}}$ ,  $g_{\text{mean}}$ ,  $H_{\text{mean}}$ ), the shapes in terms of the types of surfaces ( $S_{\text{mean}}$ ), and the surface's symmetry property ( $F_{\text{den}2}$ ). These parameters are highly informative of the investigated surface's geometrical properties. Each descriptor can underline a specific characteristic of a certain surface. These descriptors are briefly described below in regard to their conceptual order. The first and second fundamental forms provide the first six descriptors of the set. They are used to measure distance on surfaces and are defined by the formula:

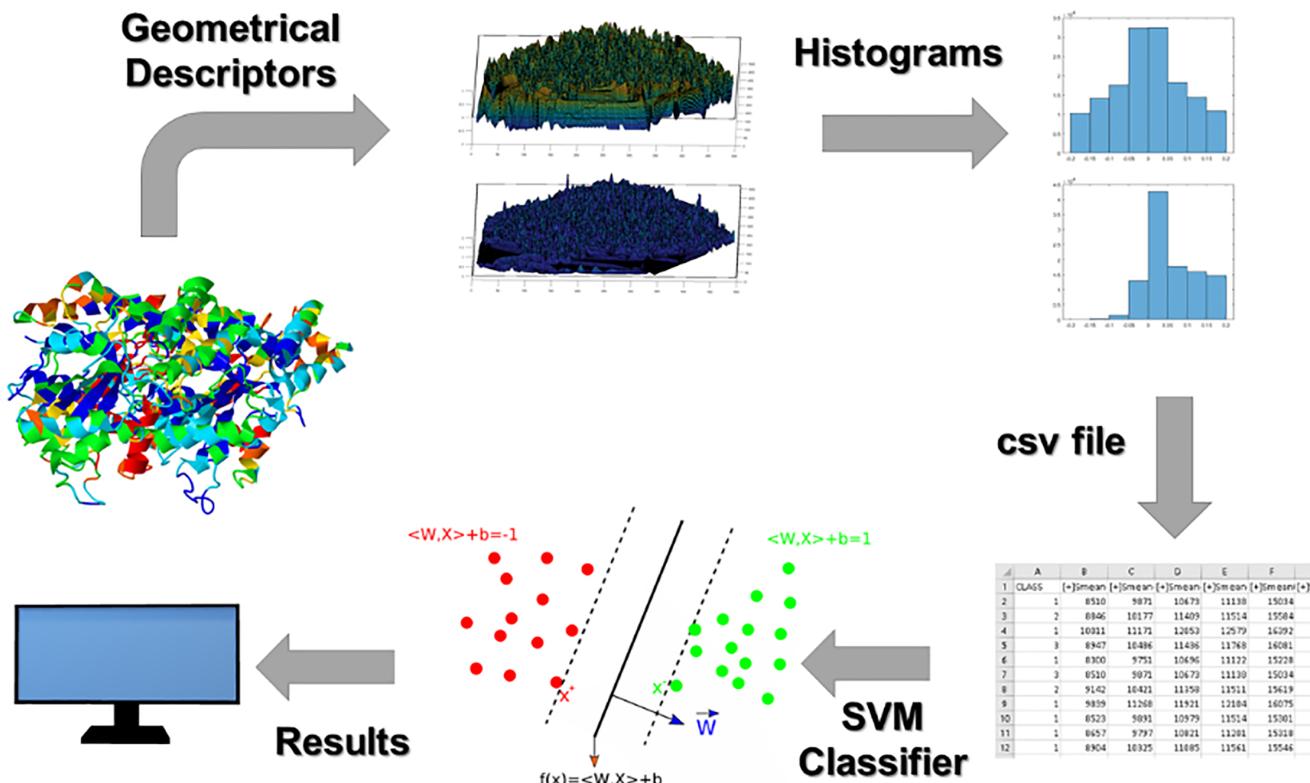
$$ds^2 = Edu^2 + 2Fdudv + Gdv^2, \quad (1)$$

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

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



**FIGURE 2** Valid solutions can be found with perceptron in a binary case (A) and the best theoretical solution that a support vector machine classifier can find (B) [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]



**FIGURE 3** Flow chart of the entire protein characterization and classification process. SVM, support vector machine [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

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

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

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

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

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

where  $h$  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  $y$ ,  $h_{xx}$ ,  $h_{yy}$ , and  $h_{xy}$  are respectively the second and mixed derivatives.

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

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

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

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

$$k_1 = H + \sqrt{H^2 - K} \text{ and } k_2 = H - \sqrt{H^2 - K}, \quad (9)$$

where,

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

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

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

$$K = k_1 k_2, \quad (12)$$

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

The shape index  $S$ , which describes the shape of the surface, is defined as<sup>24,25</sup>:

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

Some descriptors highlight particular facial lines, such as  $F_{\text{den}2}$ , which shows visible facial part contours. It can be computed using the formula:

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

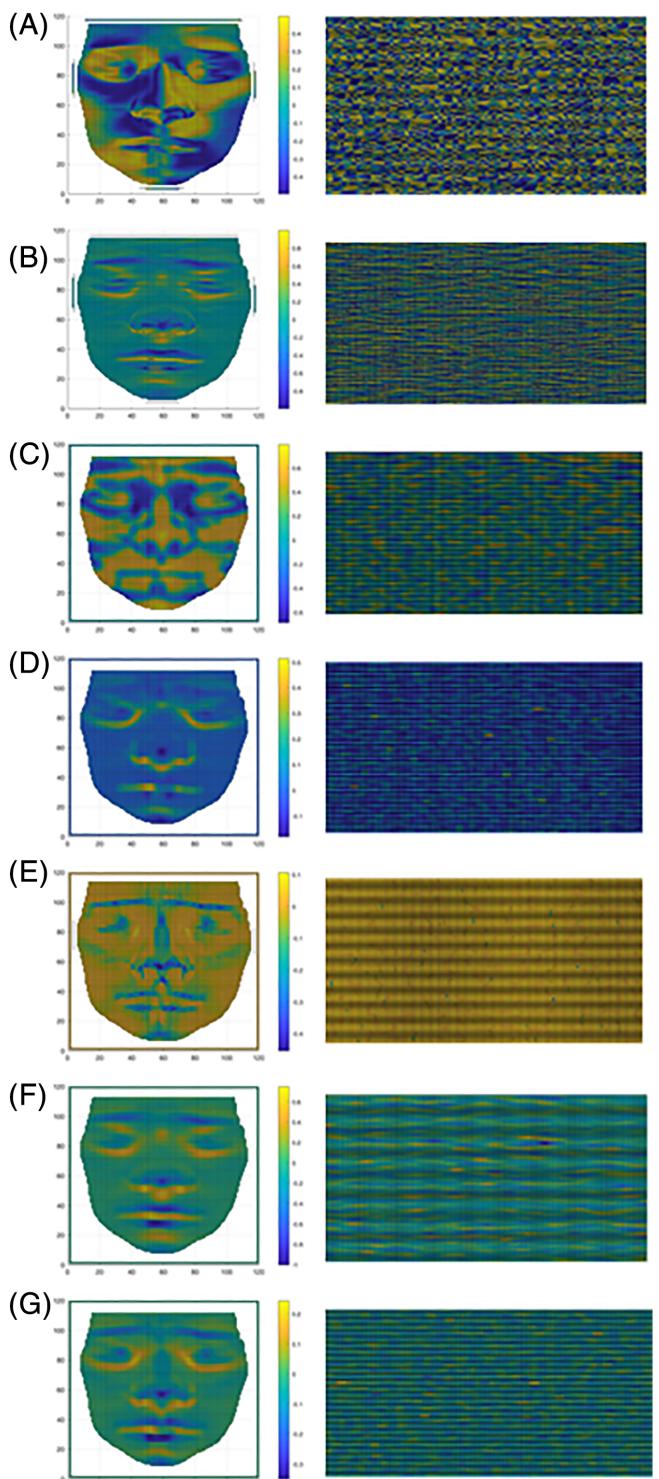
In a protein,  $F_{\text{den}2}$  can underline different trends of the free form analyzed. In particular, this descriptor has high and low values in correspondence to concavities and convexities, and values approximately equal to zero on critical points.

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 used are mapped onto the surfaces as described in Section 3.4. These descriptors are calculated for all protein faces considered in the following. An example of  $F_{\text{den}2}$  applied to both a human face and a protein is shown in Figure 4A. The descriptor  $\sin g$  is built from the application of the  $\sin e$  standard function applied to the third coefficient of the second fundamental form ( $g$ ) (see Figure 4B).<sup>23</sup> Mean and median filters have been applied to the primary descriptors  $S$ ,  $k_1$ ,  $k_2$ ,  $g$ , and  $H$ . Mean and median values are computed in squared neighborhoods of side 5 around each point of the facial depth maps.<sup>23</sup> These descriptors are labeled as follows:  $S_{\text{mean}}$ , (see Figure 4C),  $k_{1\text{mean}}$  (see Figure 4D),  $k_{2\text{median}}$  (see Figure 4E),  $g_{\text{mean}}$  (see Figure 4F), and  $H_{\text{mean}}$  (see Figure 4G).

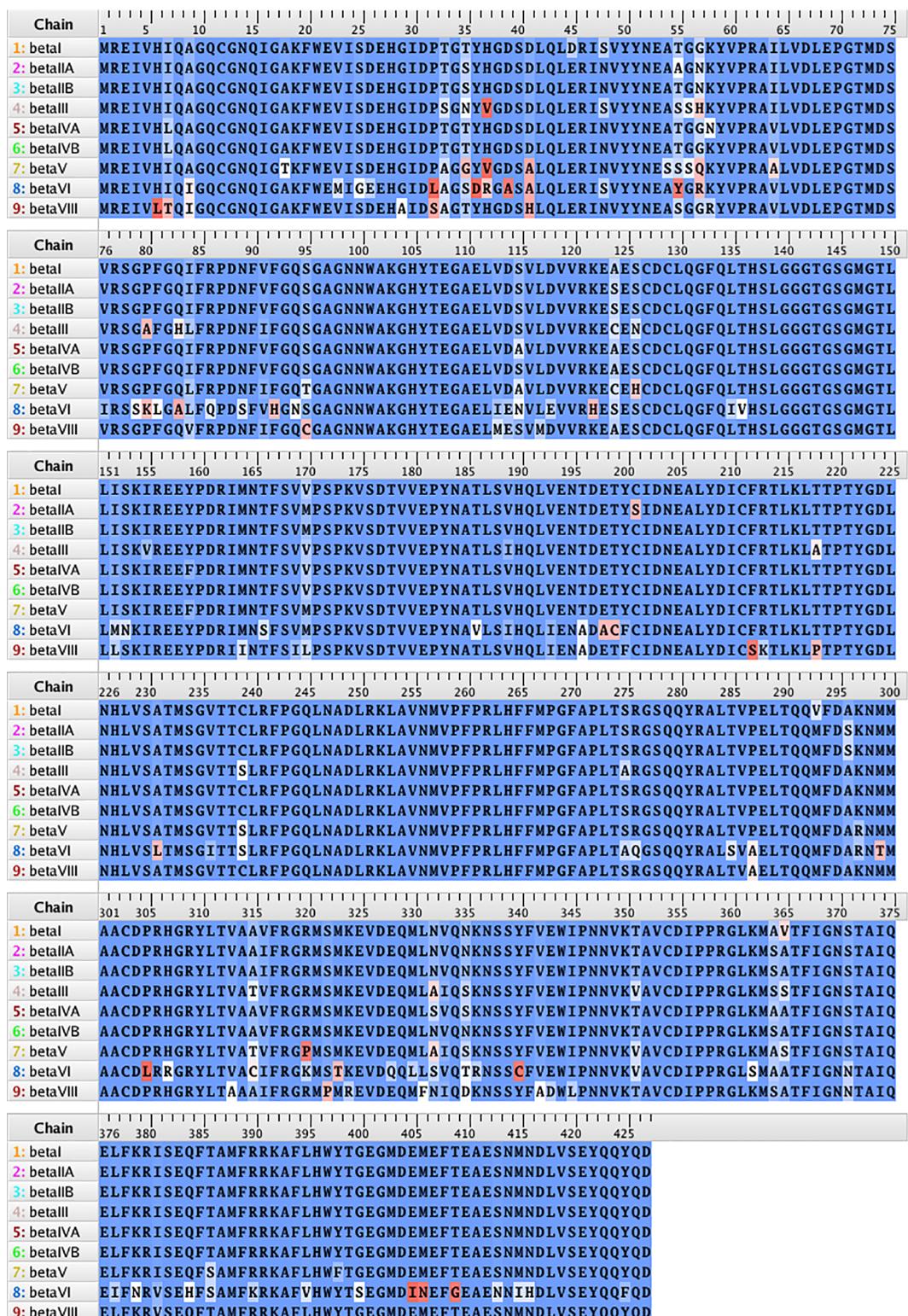
### 3 | MATERIAL AND METHODS

At the beginning of this section, we give a brief introduction to some basic concepts related to machine learning (ML), which can be useful for understanding the methods used in this article. ML is a subset of artificial intelligence (AI) tools that include mathematical and statistical models, which complete tasks with experience gained through training. The quality and amount of the training data have an important role in this process. ML classifiers can be divided into two types based on their training methods: supervised and unsupervised learning. Supervised learning needs a training phase with labeled training data (ie, sample data containing input-output pairs) in order to learn the relationship between the input and output data. On the other hand, unsupervised learning algorithms do not employ labeled training data and they aim to divide the dataset into clusters without the training phase. In this work, we use a discriminative model (a supervised model) that employs SVM. The aim of this model is to determine the division of different clusters without considering how data are generated, unlike generative models, which do consider how the data are



**FIGURE 4** Effects of applying different descriptors (A)  $F_{\text{den}2}$ , (B)  $\sin g$  (C),  $S_{\text{mean}}$ , (D)  $k_{1\text{mean}}$ , (E)  $k_{2\text{median}}$ , (F),  $g_{\text{mean}}$ , and (G)  $H_{\text{mean}}$  to a human face (left column) and to the tubulin protein (right column) [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

generated during the process. In our model, dot-product kernels are used to compute the similarity between two vectors in a higher dimensional feature in a more efficient manner. For the SVM, we tried both linear and nonlinear kernels. As the linear kernel essentially



performs the normal dot-product, the similarity score is calculated as the length of the projection of one vector onto another. The nonlinear kernel can perform the dot product in a higher dimensional feature space. Even though nonlinear kernels may be slower to use due to the computational complexity, they usually yield more favorable results. Geometric deep learning is a new field in deep learning that aims to build neural networks that can learn from non-Euclidean data, for example from graphs or complex surfaces.

The process we follow in this article 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 to 100 ns 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)<sup>26</sup> for the feature extraction using geometrical descriptors, Anaconda 1.9.6<sup>27</sup> with Python 3.7<sup>28</sup> and the library sklearn 0.22<sup>29</sup> for the implementation of classification methods and R-3.5.3 for the k-means algorithm.<sup>30</sup>

| (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 |      |

### 3.1 | 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)<sup>31</sup> was used as a template to construct the homology model for human  $\alpha\beta$  tubulin isotypes ( $\beta I$  (UniProtKB: P07437),  $\beta IIa$  (UniProtKB: UniProtKB: Q13885),  $\beta IIb$  (UniProtKB: Q9BVA1),  $\beta III$  (UniProtKB: Q13509),  $\beta IVa$  (UniProtKB: P04350),  $\alpha\beta IVb$  (UniProtKB: P68371),  $\alpha\beta V$  (UniProtKB: Q9BUF5),  $\alpha\beta VI$  (UniProtKB: Q9H4B7), and  $\beta VIII$  (UniProtKB: Q3ZCM7)) using the Molecular Operating Environment (MOE) software package.<sup>32</sup> Multiple sequence alignment results contained in Figure 5 show that human  $\beta$ -tubulin isotypes exhibit residue composition variations at different locations.

Sequence similarity matrix and sequence identity matrix of the tubulin isotypes are shown in Figure 6A,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.

The atomic coordinates of similar but not identical FtsZ dimer were obtained from the Protein Data Bank as (PDB ID: 1W5B).<sup>33</sup> The coordinates for the missing residues of the proteins were obtained by modeling using the MOE package.<sup>32</sup> 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<sup>34</sup> 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<sup>+</sup> and Cl<sup>-</sup> ions were added in order to bring the salt concentration

| (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 |      |

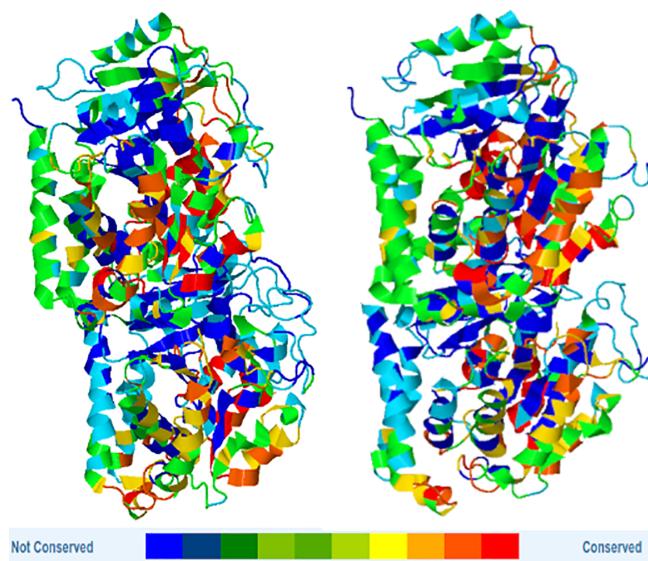
**FIGURE 6** (A) Sequence similarity matrix and (B) sequence identity matrix of the studied tubulin isotypes. The matrices are heatmap color-coded (the darker the shade, the more similar the values are) [Color figure can be viewed at wileyonlinelibrary.com]

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.

### 3.2 | 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



**FIGURE 7** Tubulin protein image for two different rotations with respect to the Z-axis. The blue color-code represents not conserved and red color represents the more conserved as it shown in the scale bar. The images were taken from <https://probis.nih.gov/> [Color figure can be viewed at wileyonlinelibrary.com]

**TABLE 1** Numbers of tubulin isotype structures used

| 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       |

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  $k = 5$ .

Cross validation is a powerful technique used to avoid overfitting. When the model is trained and tested on the same dataset, high scores can be easily obtained since the model becomes biased. In this case, low score results are obtained when the model is tested on an unseen dataset. Using cross validation, the dataset is divided into  $k$  sub parts, called folds. Then, the training is performed iteratively on the  $k-1$  folds and the remaining fold is used for the testing phase. In this way, the test set will be a truly unseen dataset for the model. One such example is shown in Figure 7 (<https://probis.nih.gov/>).<sup>35</sup>

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

### 3.3 | Protein samples

In this computational experiment, we used a total of 889 examples of tubulin structure files for 9 isotopes, as shown in 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.

### 3.4 | Data processing

The x-, y-, and z-coordinates were extracted from the PDB file. First, the data were shifted in order to be geometrically symmetric with respect to x-, y-, and z- axes, that is, the center of the coordinate systems is the geometric center of the dataset:  $(x, y, z) \rightarrow (x - \Delta x, y - \Delta y, z - \Delta z)$ , where  $\Delta x = (x_{\max} - x_{\min})/2$ ,  $\Delta y = (y_{\max} - y_{\min})/2$ , and  $\Delta z = (z_{\max} - z_{\min})/2$ . 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

**TABLE 2** Sample numbers in the binary classification between tubulin and FtsZ

| Protein | Samples |
|---------|---------|
| Tubulin | 112     |
| FtsZ    | 65      |

"griddata" commands in the Matlab with the cubic interpolation method.

The descriptors were mapped onto the surfaces as follows. The surfaces were given by point clouds where points are nonconnected (not a mesh) and arranged in a square grid. This type of data is called depth map and can be described by matrices: X, Y, Z, where Z is the one describing the "surface" and is represented in these formulas as  $h$ . Through Matlab "gradient" function, the derivatives  $h_x, h_y \dots$  were evaluated so that other matrices representing the first derivative with respect to x, the first derivative with respect to y, and so forth, were generated and stored. Then, the implementation formulas for the descriptors were calculated on the matrices previously computed and new matrices were obtained representing every geometrical descriptor.

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

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

- The class of the protein was extracted from the filename and the class was recorded in the first column of the dataset matrix;
- Geometrical descriptors were computed from matrix Z (positive and negative);
- Histograms were created and each bin was written in the right column of the dataset matrix;
- At the end of each loop the dataset matrix became the input for the classifier.

The entire process is shown in Figure 8.

In this computational experiment, nine isotopes were used (indeed, the classifier will work with nine classes). The classes were chosen one to nine in an ascendant order as shown in Table 3.

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

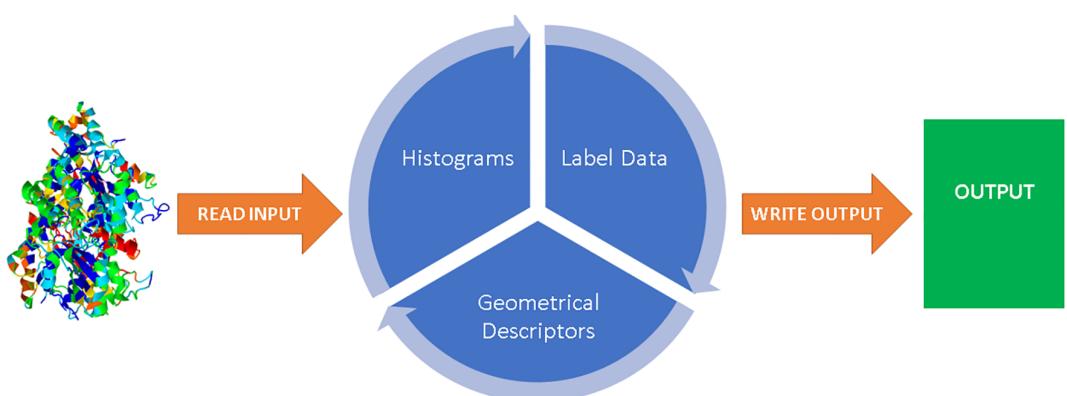
### 3.5 | Feature extraction

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

Finally, when all descriptors for all protein data were computed, the resultant matrix was copied into a file. For tubulin and other proteins, these descriptors can underline specific characteristic of a certain surface. They can indicate different trends of the free form analyzed and they can describe the shape of the surface. The features are extracted with multiple geometrical descriptors to extract more details; using this approach, also small differences in convexity and concavity can be recognized during the classification. Analyzing the features extracted, the most important features were found from parameter values of  $F_{\text{den}2}$  and  $\sin g$ , because analyzing the data these values were sufficiently different to help the classifier select the right class. In particular,  $F_{\text{den}2}$  is meant to be descriptive for its behavior in the loci of critical points, and  $\sin g$  for curvature changes, local minimums in convexities and local maximums in concavities, respectively.

### 3.6 | Classification

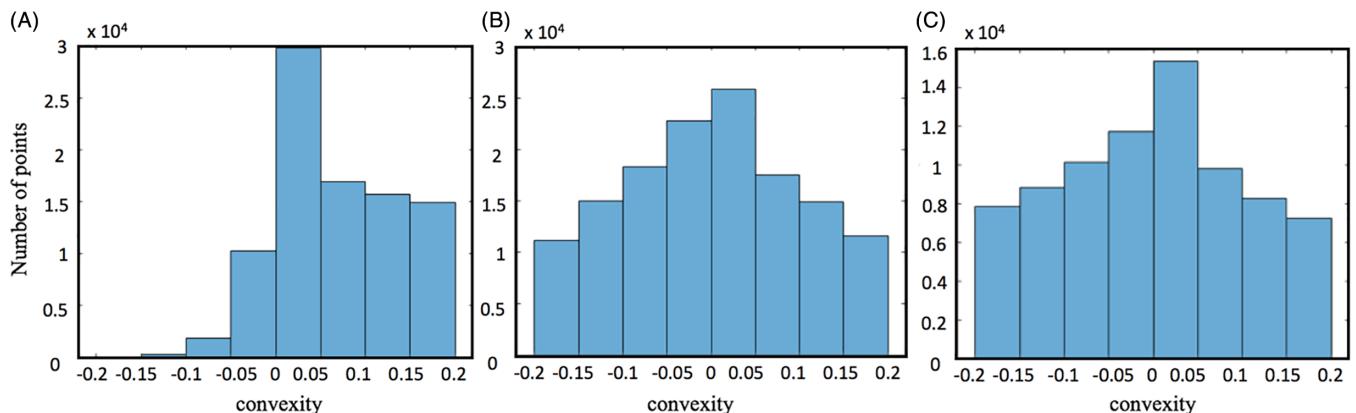
The adopted classifiers were k-means and SVM. First, an unsupervised method was tested (k-means) using 9 clusters and a limited number of iterations, then a supervised method (SVM) using linear and nonlinear kernels was used. In these cases, it is not a simply binary classification, but there are many classes (9) and many features (more than 100), so some distributions cannot separate the dataset in a linear way or with



**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 molecular dynamics 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 [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

**TABLE 3** Number of tubulin isotypes used

| 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         |

**FIGURE 9** 9 bin histograms calculated using (A)  $F_{\text{den}2}$ , (B)  $g_{\text{mean}}$ , and (C)  $H_{\text{mean}}$  geometrical descriptor [Color figure can be viewed at wileyonlinelibrary.com]

a linear separator as a high misclassification rate is reached. An interesting improvement is to use a nonlinear separator or a kernel trick. An example of a nonlinear kernel is the radial basis function (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 nonlinear kernel can reach better results. The difference between linear and nonlinear 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 nonlinear kernel uses a nonlinear function, so it can divide the dataset better. The cross validation has not been performed here because the results were positive, and hence we have already avoided the over-fitting problem. The validation part was performed using a large number of parameters and the best ones were selected for the testing part.

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

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

where  $c$  is the centroid and  $x$  the data points.

In this last phase, the centroids are computed again as the mean of all the data points of the cluster:

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

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

The stop condition is given by the following criteria:

- No data points change the cluster;
- The sum of distances is at the minimum;
- The maximum number of iterations is reached.

Therefore, when the convergence is obtained the algorithms stops.

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

### 3.6.2 | Support vector machine

The first test was performed using a linear kernel where  $\lambda$  is a key parameter of SVM. In fact, the main factors in SVM are setting a large margin and reducing the misclassification rate. These two properties are inversely proportional, and the  $\lambda$  parameter helps to find a trade-

off. A large value of  $\lambda$  is for a small margin, whereas a small value of  $\lambda$  is for a large margin. The right  $\lambda$  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  $\lambda$  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  $\lambda$  values. As a matter of fact, by increasing the  $\lambda$  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%.

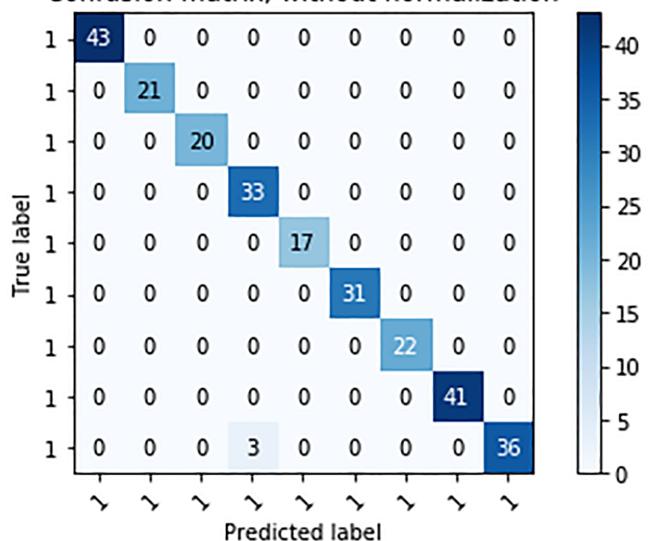
The dataset was built using nine different Tubulin isotypes. Hence, the number of classes used for the SVM classifier was nine; the same number was used in the k-means test, in order to have comparable results. The confusion matrix is an important tool to evaluate the results, since it gives precise 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.

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  $\lambda$  and gamma parameters in the range between  $10^{-5}$  and  $10^{15}$  and it is evaluated on the validation set. The best parameters on the validation set are found to be:  $\lambda = 100$  and  $\text{gamma} = 10^{-9}$  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  $\lambda$  and gamma values. The gamma parameter that is used in the RBF kernel function is the inverse of the SD of the RBF kernel, which is used as a similarity function. A small value of gamma indicates a large variance where two points can be matched as similar. This results in a smoother decision-making by the model. A higher gamma value has the opposite effect on the process. The challenge will be to find an optimum value of gamma for the given data set. Indeed, by

Confusion matrix, without normalization



Normalized confusion matrix

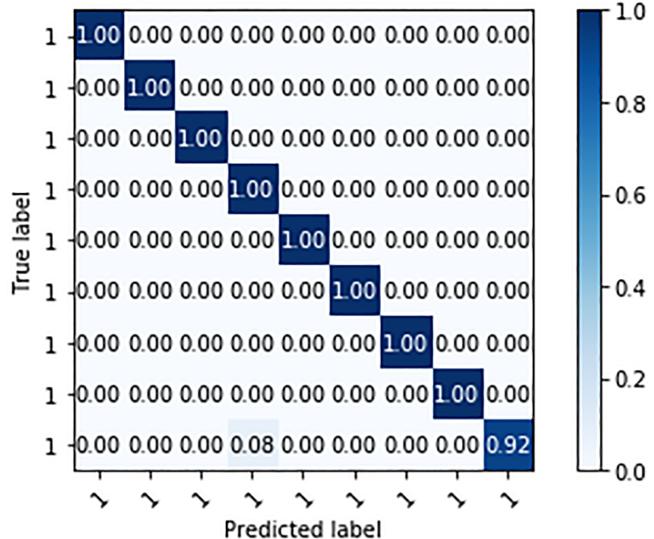


FIGURE 10 Confusion matrix of support vector machine (SVM) classifier using the radial basis function kernel [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

increasing the  $\lambda$  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  $\lambda = 100$ ,  $\text{gamma} = 10^{-9}$  and the accuracy obtained was 96.5%.

The same methodology was applied to tubulin and FtsZ classifications.

## 4 | 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** Tubulin isotypes accuracy results

| Classifier             | Accuracy (%) |
|------------------------|--------------|
| SVM with RBF kernel    | 96.5         |
| SVM with linear kernel | 92.4         |
| k-means                | 76.6         |

Abbreviations: RBF, radial basis function; SVM, support vector machine.

**TABLE 5** Accuracy results for the tubulin and FtsZ binary classification

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

Abbreviations: RBF, radial basis function; SVM, support vector machine.

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.

These results are competitive with the state-of-the-art results found in the literature. A fast protein classification method<sup>5</sup> based on an SVM classifier reached an accuracy of about 90% with 13 classes. Another study<sup>7</sup> used a semisupervised 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<sup>36</sup> 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 article. This test involved four arbitrarily chosen FtsZ protein structures, namely: 2R6R, 2VAW, 2VAP, and 2VAM. These structures correspond, respectively, to the following biological species: *Bacillus subtilis*, *Pseudomonas aeruginosa*, *Mathanococcus jannaschii*, and *Aquifex aeolicus*. In this test 683 samples were used as given in 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.

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  $k = 5$ .

**TABLE 6** 2R6R, 2VAM, 2VAP, and 2VAM samples

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

**TABLE 7** 2R6R, 2VAM, 2VAP, and 2VAM experiment

| Classifier             | Accuracy (%) |
|------------------------|--------------|
| SVM with RBF kernel    | 97.1         |
| SVM with linear kernel | 98.0         |
| k-means                | 62.3         |

Abbreviations: RBF, radial basis function; SVM, support vector machine.

## 5 | CONCLUSIONS

A novel method for protein characterization and classification has been proposed in this article, 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 nine of its isotypes. The second application performed used two structurally similar proteins: bovine 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 structurally similar or not but does not address the issue regarding which features differ between them. For this reason, our method can assist in identifying structure-function dependence when comparing various proteins, even highly similar ones. Since we only investigate geometrical features, both physical and chemical properties are not directly involved in our method but can eventually be extracted by mapping geometrical features back onto to amino acid distributions underlying them. Also, the number of potential mutations of any protein, in particular tubulin, is astronomical. Consequently, brute force methods are not viable in classifying the role of specific mutations regarding the root causes of the conformational

changes resulting in dysfunction of a given protein. However, our methodology based on ML approaches may offer a viable alternative with numerous potential applications in protein biophysics and beyond.

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

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

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;
- Applying our method to different data set for the purpose of protein classification such as hemoglobin classification.<sup>37</sup> Additional proteins of interest that could be investigated using our methods involve those with significant roles in neurodegenerative diseases that have been previously investigated using MD simulations, for example: Josephine domain protein involved in spinocerebellar ataxia<sup>38</sup> as well as Ataxin-1,<sup>39</sup> amyloid beta involved in Alzheimer's disease<sup>40</sup> and a host of MT-associated proteins such as MAP-tau.
- 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.<sup>19</sup> 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<sup>41</sup> accelerate the training of the kernel machine.

## ACKNOWLEDGMENT

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

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## 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, Federica Marcolin and Jack Adam Tuszynski; methodology, Federica Marcolin, Jack Adam Tuszynski and Luca Di Grazia; software, Luca Di Grazia; validation, Luca Di Grazia; formal analysis, Federica Marcolin and Jack Adam Tuszynski; investigation, Federica Marcolin, Jack Adam Tuszynski and Luca Di Grazia; Matlab scripts, Vahid Rezania; resources, Maral Aminpour and Enrico Vezzetti; data curation, Maral Aminpour; writing—original draft preparation, Luca Di Grazia and Maral Aminpour; writing—review and editing, Maral Aminpour, Federica Marcolin and Jack Adam Tuszynski; visualization, Luca Di Grazia and Maral Aminpour; supervision, Federica Marcolin and Jack Adam Tuszynski; project administration, Jack Adam Tuszynski; funding acquisition, Jack Adam Tuszynski and Enrico Vezzetti", please turn to the CRediT taxonomy for the term explanation.

## PEER REVIEW

The peer review history for this article is available at <https://publons.com/publon/10.1002/prot.25993>.

## ORCID

Jack Adam Tuszynski  <https://orcid.org/0000-0001-9976-0429>

## REFERENCES

1. Gainza P, Sverrisson F, Monti F, Rodola E, Bronstein MM, Correia BE. Deciphering interaction fingerprints from protein molecular surfaces. *bioRxiv*. 2019;606202.
2. Planas-Iglesias J, Bonet J, García-García J, Marín-López MA, Feliu E, Oliva B. Understanding protein–protein interactions using local structural features. *J Mol Biol*. 2013;425(7):1210–1224.
3. Rupp B, Wang J. Predictive models for protein crystallization. *Methods*. 2004;34(3):390–407.
4. Saberi Fathi SM, White DT, Tuszynski JA. Geometrical comparison of two protein structures using Wigner-D functions: geometrical comparison of protein structures. *Proteins*. 2014;82(10):2756–2769.
5. Tsuda K, Shin H, Schölkopf B. Fast protein classification with multiple networks. *Bioinformatics*. 2005;21(2):59–65.
6. Weston J, Leslie C, le E, Zhou D, Elisseeff A, Noble WS. Semi-supervised protein classification using cluster kernels. *Bioinformatics*. 2005;21(15):3241–3247.
7. Jain P, Garibaldi JM, Hirst JD. Supervised machine learning algorithms for protein structure classification. *Comput Biol Chem*. 2009;33(3):216–223.
8. Masci J., Boscaini D, Bronstein M, Vandergheynst P. Geodesic convolutional neural networks on riemannian manifolds. Paper presented at: Proceedings of the IEEE International Conference on Computer Vision Workshops 2015. IEEE Computer Society.
9. Monti, F., Boscaini D, Masci J, Rodola E, Svoboda J, Bronstein MM. Geometric deep learning on graphs and manifolds using mixture model CNNs. Paper presented at: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition 2017.

10. Bronstein MM, Bruna J, LeCun Y, Szlam A, Vandergheynst P. Geometric deep learning: going beyond Euclidean data. *IEEE Signal Process Mag.* 2017;34(4):18-42.
11. Espinosa E, Zamora P, Feliu J, González Barón M. Classification of anticancer drugs—a new system based on therapeutic targets. *Cancer Treat Rev.* 2003;29(6):515-523.
12. Huang C-H, Lee F-L, Tai C-J. The  $\beta$ -tubulin gene as a molecular phylogenetic marker for classification and discrimination of the Saccharomyces sensu stricto complex. *Antonie Van Leeuwenhoek.* 2009;95(2):135-142.
13. Ludueña RF. Are tubulin isotypes functionally significant. *Mol Biol Cell.* 1993;4(5):445-457.
14. Fitch WM. Homology: a personal view on some of the problems. *Trends Genet.* 2000;16(5):227-231.
15. Richards KL, Anders KR, Nogales E, Schwartz K, Downing KH, Botstein D. Structure-function relationships in yeast tubulins. *Mol Biol Cell.* 2000;11(5):1887-1903.
16. Schlieper D, Oliva MA, Andreu JM, Lowe J. Structure of bacterial tubulin BtubA/B: evidence for horizontal gene transfer. *Proc Natl Acad Sci USA.* 2005;102(26):9170-9175.
17. Gunn, S.R.. Support vector machines for classification and regression. ISIS Technical Report. 1998:52.
18. pCloud. A new protein characterization and classification method using 3D face recognition algorithms. Available from: <https://u.pcloud.link/publink/show?code=XZwyRNkZdgxbscKvDcz9RcNn832cPYuD3pRV>.
19. Ciravegna G, Cirrincione G, Marcolin F, Barbiero P, Dagnes N, Piccolo E. Assessing discriminating capability of geometrical descriptors for 3D face recognition by using the GH-EXIN neural network. In: Esposito A, ed. *Neural Approaches to Dynamics of Signal Exchanges.* Singapore: Springer; 2020:223-233.
20. Cirrincione G, Marcolin F, Spada S, Vezzetti E. Intelligent quality assessment of geometrical features for 3D face recognition. In: Esposito A, ed. *Neural Advances in Processing Nonlinear Dynamic Signals.* Cham: Springer International Publishing; 2019:153-164.
21. Li SZ, Jain AK. *Handbook of Face Recognition.* 2nd ed. London: Springer-Verlag; 2011.
22. Marcolin F, Violante MG, Sandro MO, et al. Three-dimensional face analysis via new geometrical descriptors. In: Eynard B, ed. *Advances on Mechanics, Design Engineering and Manufacturing.* Cham: Springer International Publishing; 2017:747-756.
23. Marcolin F, Vezzetti E. Novel descriptors for geometrical 3D face analysis. *Multimed Tools Appl.* 2017;76(12):805-834.
24. Koenderink JJ, Van Doorn AJ. Surface shape and curvature scales. *Image Vision Comput.* 1992;10(8):557-564.
25. Vezzetti E, Marcolin F. Geometrical descriptors for human face morphological analysis and recognition. *Robot Auton Syst.* 2012;60(6):928-939.
26. MATLAB. Statistics and machine learning toolbox 2018, The MathWorks Inc.: Natick, MA.
27. Anaconda. Anaconda Software Distribution. Computer software. 2019.
28. Van Rossum G, Drake FL Jr. *Python.* Amsterdam, The Netherlands: Centrum voor Wiskunde en Informatica; 2019.
29. Version 0.22.0—scikit-learn 0.22 documentation. Available from: [https://scikit-learn.org/stable/whats\\_new/v0.22.html](https://scikit-learn.org/stable/whats_new/v0.22.html).
30. Download R-3.5.3 for Windows. The R-project for statistical computing. Available from: <https://cran.r-project.org/bin/windows/base.old/3.5.3/>.
31. Löwe J, Li H, Downing KH, Nogales E. Refined structure of  $\alpha\beta$ -tubulin at 3.5 Å resolution. *J Mol Biol.* 2001;313(5):1045-1057.
32. Molecular Operating Environment (MOE). Group, chemical computing. 2012: Montreal, QC, Canada.
33. Oliva MA, Cordell SC, Löwe J. Structural insights into FtsZ protofilament formation. *Nat Struct Mol Biol.* 2004;11(12):1243-1250.
34. D.A. Case, Darden TA, Cheatham III TE, Simmerling CL, Wang J, Duke RE, Luo R, Merz KM, Pearlman DA, Crowley M, Walker RC. AMBER 2014. 2014: University of California, San Francisco.
35. Konc J, Miller BT, Štular T, et al. ProBiS-CHARMMing: web interface for prediction and optimization of ligands in protein binding sites. *J Chem Inf Model.* 2015;55(11):2308-2314.
36. Saberi Fathi SM, Tuszyński JA. A simple method for finding a protein's ligand-binding pockets. *BMC Struct Biol.* 2014;14:18.
37. Cang Z, Mu L, Wu K, Opron K, Xia K, Wei GW. A topological approach for protein classification. *Mol Based Math Biol.* 2015;3:140-162.
38. Deriu MA, Grasso G, Licandro G, et al. Investigation of the Josephin domain protein-protein interaction by molecular dynamics. *PLoS One.* 2014;9(9):108677.
39. Gianvito G, Deriu MA, Tuszyński JA, Gallo D, Morbiducci U, Danani A. Conformational fluctuations of the AXH monomer of Ataxin-1. *Proteins.* 2015;84(1):52.
40. Gianvito G, Rebella M, Muscat S, et al. Conformational dynamics and stability of U-shaped and S-shaped amyloid  $\beta$  assemblies. *Int J Mol Sci.* 2018;19(2):571.
41. Rahimi A, Recht B. Random features for large-scale kernel machines. *Advances in Neural Information Processing Systems.* Cambridge, MA: MIT Press; 2008:10.

**How to cite this article:** Di Grazia L, Aminpour M, Vezzetti E, Rezania V, Marcolin F, Tuszyński JA. A new method for protein characterization and classification using geometrical features for 3D face analysis: An example of tubulin structures. *Proteins.* 2021;89:53-67. <https://doi.org/10.1002/prot.25993>